



Università degli Studi di Padova
Facoltà di Ingegneria

Corso di Laurea Magistrale in Ingegneria
dell'automazione

Tesi di laurea magistrale

Analisi di forma e movimento di cellule in strutture reticolari

Candidato:
Camilla Corfini
Matricola 1014127

Relatore:
Prof. Angelo Cenedese

Anno Accademico 2012–2013

Sometimes the winner is simply a dreamer who never gave up.

— Jim Morrison.

Dedicato alla mia famiglia.
Ad Alessandro e la sua famiglia.
Ai miei amici.

*C'è qualcos'altro che devi sempre ricordare.
Ci sono persone magnifiche su questa terra,
che se ne vanno in giro travestite da normali esseri umani.
Non scordarlo mai.*

— Fannie Flag, Pomodori verdi fritti.

RINGRAZIAMENTI

Ringrazio i miei genitori e mio fratello per avermi sostenuta in ogni momento della mia vita e per aver arricchito di passione e idee ogni mio progetto. Grazie per l'esempio meraviglioso che siete.
Ringrazio Alessandro per aver sempre tenuto alto l'umore della squadra, per avermi incoraggiata quando ne avevo più bisogno e per aver sempre creduto in me. Grazie per l'uomo stupendo che sei.
Ringrazio la sua famiglia perché mi ha accolta a braccia aperte.
Ringrazio i miei amici per essere sempre presenti.
Infine ringrazio il mio relatore, il Professor Angelo Cenedese, per la sua disponibilità e il suo supporto durante lo svolgimento di questo lavoro.

Padova, marzo 2013

C. C.

INDICE

1	CONTESTO BIOLOGICO	1
1.1	Richiami di biologia dello sviluppo	1
1.2	Caso di studio: <i>Drosophila Melanogaster</i>	2
1.2.1	Modello sperimentale	2
1.2.2	Caratteristiche del <i>Drosophila melanogaster</i>	3
1.2.3	Il moscerino della frutta modello della ricerca	4
1.3	Visione computazionale e biologia	6
1.3.1	Cenni su morfogenesi dei tessuti e migrazione cellulare	7
2	RICHIAMI DI VISIONE COMPUTAZIONALE	11
2.1	Riconoscimento e analisi di forme	12
2.2	Modelli di rappresentazione di forme	14
2.3	<i>Single shape detection</i>	16
2.3.1	<i>Standard Active Contours (snake)</i>	17
3	TRACKING E ANALISI DI UNA CELLULA	21
3.1	<i>Dynamic shape detection</i>	21
3.1.1	<i>Optical Flow</i> e J-maps	23
3.1.2	Algoritmo di <i>tracking</i> della struttura reticolare	25
3.2	<i>Shape Deformation Analysis</i>	25
3.2.1	Moto e deformazione	27
4	ANALISI DI STRUTTURE RETICOLARI	39
4.1	<i>Shape detection</i> in strutture reticolari	39
4.1.1	<i>The Random Walk Agents Algorithm</i>	40
4.2	Implementazione: video di <i>Drosophila wing</i>	43
4.3	<i>Tracking</i> della struttura reticolare	45
4.3.1	Analisi di deformazione di gruppi di cellule limitrofe	45
5	FILTRAGGIO STATISTICO	51
5.1	Il problema della stima	51
5.2	Le teorie del filtraggio statistico	53
5.3	Filtro di Kalman	54
5.3.1	Le equazioni del filtro di Kalman	55
5.4	Applicazione del filtraggio al modello di riferimento	57
6	CONCLUSIONI E SVILUPPI FUTURI	63
A	RICHIAMI DI REGRESSIONE STATISTICA	65
A.0.1	Scatterplots	65

A.o.2	Definizioni	65	
A.o.3	Proprietà della covarianza	67	
A.o.4	Proprietà della correlazione	70	
A.o.5	Regressione lineare	70	
A.o.6	Residui	72	
B	PRINCIPAL COMPONENT ANALYSIS	75	
B.1	Formalizzazione della PCA	75	
	BIBLIOGRAFIA	81	

ELENCO DELLE FIGURE

Figura 1	Esemplare di drosophila	4
Figura 2	Epitelio di drosophila	6
Figura 3	Microscopio a fluorescenza ottica	7
Figura 4	Meccanismi di locomozione cellulare	9
Figura 5	Trasformazioni Euclidee della stessa forma	13
Figura 6	Landmarks corrispondenti a diverse forme	13
Figura 7	Funzioni implicite	16
Figura 8	Primo template dell'occhio umano	17
Figura 9	Single Shape Detection	17
Figura 10	Tracking reticolare su <i>Drosophila epithelium</i>	26
Figura 11	Medusa in moto-deformazione	27
Figura 12	Traslazione di una cellula in 20 frames	29
Figura 13	Traslazione di Xc e Yc	29
Figura 14	Rotazione di una cellula	31
Figura 15	Radius Vector Function	33
Figura 16	Distanza di Hausdorff	34
Figura 17	Metriche di deformazione	35
Figura 18	Contour plot dell'andamento del raggio	36
Figura 19	Reticolo cellule.	37
Figura 20	Contour plot dell'andamento del raggio	37
Figura 21	Strutture reticolari	39
Figura 22	Schema logico Random Walk Algorithm	43
Figura 23	Gruppi cellule adiacenti	46
Figura 24	Boundary Activity per il Gruppo 1	46
Figura 25	Elongazione per il Gruppo 2	47
Figura 26	Scatterplots DIFF Vs DELTA	49
Figura 27	Scatterplots DIFF Vs DELTA	50
Figura 28	Kalman Algorithm	57
Figura 29	Random walk e filtro di kalman	59
Figura 30	Kalman su random walk	60
Figura 31	Kalman su random walk - stime	61

ABSTRACT

In questa tesi viene presentata un'implementazione di *dynamic shape detection*, applicata ad una sequenza di immagini di epitelio di *Drosophila Melanogaster*, sulla base della quale viene infine condotta un'analisi su moto e deformazione della struttura cellulare osservata. Lo scopo dell'elaborato consiste in primo luogo nell'estrarre una buona rappresentazione della struttura reticolare dell'epitelio cellulare, grazie alla tecnica di *Active Contours*. Quindi, estraendo da questa una serie di features per ogni cellula (punti di contorno, perimetro area, posizione centroide, etc.) ed estendendo il metodo di *shape detection* ad una sequenza di frames, riuscire a tracciare le cellule durante il filmato e infine condurre un'analisi sulla loro evoluzione di forma e sul moto.

SOMMARIO

Nel Capitolo 1 vengono introdotte alcune nozioni generali di biologia dello sviluppo, riguardo in particolare lo studio sull'evoluzione di cellule e tessuti. Il Capitolo viene concluso con la descrizione del caso di studio analizzato in questo lavoro di tesi. Nel Capitolo 2 invece viene introdotto il concetto di *forma (shape)*; poi viene brevemente trattata la tecnica di riconoscimento statico di forme (in inglese *shape detection*) e l'algoritmo che viene usato nel lavoro per riconoscere una cellula all'interno di un'immagine digitale. Il Capitolo 3 estende l'approccio di *single shape detection* ad una sequenza di immagini, e quindi esplora la possibilità di effettuare un *tracking* della cellula lungo un filmato (sequenza di frames consecutivi). Viene concluso il capitolo con l'analisi di moto e deformazione di una singola cellula nel tempo e per fare questo si definiscono alcune metriche utili. Il Capitolo 4 estende dapprincipio la tecnica di *static single shape detection*, trattata nel Capitolo 2, a forme di struttura reticolare (reticolo di cellule); in seguito estende la *reticular detection* ad una sequenza di immagini, analogamente a ciò che si fa nel Capitolo 3. L'ultima parte del lavoro tratta, nel Capitolo 5, di filtraggio statistico e in particolare mostra un approccio di filtraggio alla Kalman al modello di moto del reticolo cellulare.

INTRODUZIONE

L'elaborazione di immagini biologiche e mediche, negli anni, ha acquisito un'importanza sempre crescente ed è, ad oggi, in grado di fornire strumenti di alto livello tecnologico per affiancare il lavoro di medici e biologi. Lo scopo dell'elaborazione automatica di questo tipo di dati visivi, che possono essere sia di tipo istantaneo (un singolo frame) che di tipo dinamico (video), è in primo luogo legato all'analisi di sistemi biologici molto complessi, per poterne osservare e fornire in maniera automatica una rappresentazione schematica che li rappresenti nel modo più accurato e semplice possibile; in secondo luogo quindi alla possibilità fornire degli strumenti utili nella diagnosi e nella ricerca per acquisire una certa quantità di informazioni dalle suddette rappresentazioni.

Per raggiungere questo secondo scopo è necessario fare un passo indietro e capire come "estrarre" informazioni in maniera automatica a partire da un file immagine. L'obiettivo è cercare di sviluppare un algoritmo il più possibile robusto e adatto alla specifica applicazione che si vuole trattare.

Dal punto di vista tecnico, il primo passo da fare consiste nella "pulizia" delle immagini da fonti di rumore o più in generale di distorsione, in modo da migliorare i dati immagine o intensificare qualche caratteristica dell'immagine utile per le successive fasi di processamento; questa fase è detta appunto preprocessamento e varia da caso a caso a seconda del tipo di immagini che si stanno trattando. Le trasformazioni che comprende sono vari tipi di filtraggio ma anche trasformazioni geometriche (come rotazioni, scalamenti o traslazioni) sono classificate come metodi di *preprocessing*.

Il passo successivo consiste nella *segmentazione* dell'immagine e quindi nell'estrazione di un modello astratto di essa che riduca la quantità di informazioni da analizzare nei passaggi successivi. A questo punto si colloca il riconoscimento di forme (*shape detection*) che mediante diverse tecniche, come si intuisce, conduce al tracciamento di un particolare oggetto visibile nell'immagine digitale. Nel caso specifico che si è analizzato, l'obiettivo è stato tracciare un reticolo che seguisse la struttura dell'epitelio di *Drosophila*, sfruttando l'algoritmo *Random Walkers* recentemente presentato in [Silletti, 2007-2009](#).

Una volta riconosciuta la struttura di interesse è stato possibile generare un modello che la rappresentasse e dal quale fossero estraibili *features* interessanti per condurre un *analisi* sulla deformazione e sul moto delle cellule osservate. L'analisi dell'immagine è quasi sempre l'obiettivo finale del riconoscimento di forme, infatti avere a disposizione una forma senza poter estrarre informazioni significative per

l'applicazione è inutile. In pratica ciò che si vuole fornire, con l'analisi di metriche rilevanti su moto e deformazione delle cellule, sono una serie di dati quantitativi che un biologo possa sfruttare per tracciare e studiare dei comportamenti tipo o magari registrare il comportamento delle cellule per esempio a seguito di determinati stimoli. Questo può aiutare i biologi soprattutto nel quantificare in maniera dettagliata e automatica il comportamento dinamico delle cellule, durante alcune fasi importanti di evoluzione di tessuti.

Una delle motivazioni per cui è interessante un progetto come questo è la sfida, che sta alla base di esso, di mimare l'effetto della visione umana nel percepire e comprendere un'immagine attraverso la visione computerizzata. L'intenzione di questo lavoro è quella di fornire questa possibilità, calandosi in un contesto pratico, quello dello studio dei tessuti nel modo più apprezzabile che si possa. Si cerca di fornire, all'interno di un contesto pratico, l'utilità di tecniche come l'*Active Contour* e i *Random Walk Agents* per condurre un'analisi su immagini aventi strutture reticolari, per altro molto frequenti in biologia.

1

CONTESTO BIOLOGICO

Indice

1.1	Richiami di biologia dello sviluppo	1
1.2	Caso di studio: <i>Drosophila Melanogaster</i>	2
1.2.1	Modello sperimentale	2
1.2.2	Caratteristiche del <i>Drosophila melanogaster</i>	3
1.2.3	Il moscerino della frutta modello della ricerca	4
1.3	Visione computazionale e biologia	6
1.3.1	Cenni su morfogenesi dei tessuti e migrazione cellulare	7

Questa sezione fornisce una breve panoramica del background scientifico-biologico del progetto. Questo progetto, come si intuisce, si colloca all'interno di quella disciplina nota col nome di bioingegneria che ha come caratteristica quella di mettere tecniche e strutture tipiche dell'ingegneria al servizio di studi di natura medica o biologica, come nel nostro caso.

1.1 RICHIAMI DI BIOLOGIA DELLO SVILUPPO

La biologia dello sviluppo è la disciplina che studia i meccanismi molecolari e fisiologici che controllano le varie fasi embrionali e la formazione di cellule, organi, tessuti costituente un lento processo di cambiamenti progressivi, chiamato sviluppo. Con il termine sviluppo si intendono: l'arco temporale che, dalla fecondazione, porta alla formazione di un organismo vivente adulto (embriologia) in grado a sua volta di riprodursi; e i cambiamenti che intercorrono nella crescita ed organizzazione dell'organismo. Per esempio, nell'essere umano ogni giorno circa un grammo di cellule della cute vengono sostituite ed ogni minuto il midollo osseo produce milioni di nuovi globuli rossi. In molti organismi, infatti, lo sviluppo prosegue per tutta la durata della loro esistenza. La biologia dello sviluppo usa metodi di biologia cellulare (citologia), genetica, biologia molecolare, biochimica e microscopia, e studia soprattutto alcuni organismi chiamati organismi modello.

Gli eventi che caratterizzano lo sviluppo sono:

- **Gametogenesi:** per gametogenesi si intende quel processo che ha luogo nelle gonadi e porta alla formazione dei gameti, ossia delle cellule sessuali mature, capaci quindi di fecondare o di essere fecondate.
- **Fecondazione:** la fecondazione avviene solo nella riproduzione sessuata anfigonica, ossia con la fusione di gameti. Il risultato della fecondazione è una nuova cellula, diversa dai gameti e unica nella sua specie, chiamata zigote.
- **Segmentazione:** in biologia, con il termine segmentazione si intende un intenso processo in cui l'ovulo fecondato subisce una serie di divisioni mitotiche che portano alla divisione in cellule chiamate blastomeri.
- **Gastrulazione:** la gastrulazione è un tipico processo embrionale che consiste in movimenti morfogenetici e di differenziazione utili alla sistemazione dei foglietti embrionali primari (ectoderma, endoderma) e di quello secondario (mesoderma).
- **Morfogenesi:** in biologia la morfogenesi è lo sviluppo della forma e della struttura di un organismo, sia da un punto di vista evolutivo, sia dal punto di vista dello sviluppo ontogenetico del singolo organismo a partire dalla cellula fecondata (sviluppo embrionale)
- **Organogenesi:** l'organogenesi è il meccanismo di costruzione e crescita delle varie parti dell'embrione che rispetta parametri quantitativi e qualitativi tali da far riconoscere un individuo come appartenente ad una determinata specie.

1.2 CASO DI STUDIO: DROSOPHILA MELANOGASTER

1.2.1 Modello sperimentale

In generale, un modello sperimentale è un sistema semplice ed idealizzato che è accessibile e facile da manipolare. Nella scelta di organismi viventi da utilizzare come modelli sperimentali, vengono adottati alcuni criteri che dipendono in parte dalle finalità sperimentali. Tuttavia, vi sono alcune caratteristiche che sono comuni alla maggior parte degli organismi modello, come ad esempio:

1. Lo sviluppo rapido ed un ciclo vitale breve;
2. Le dimensioni ridotte;
3. La facile reperibilità;

4. La facile manipolazione.

Tali caratteristiche mirano principalmente a contenere le spese di allevamento, dati i limiti di spazio fisico e di budget normalmente a disposizione dei laboratori. Generalmente parlando, gli organismi modello fungono da surrogati che permettono la conduzione di esperimenti che non potrebbero altrimenti essere condotti sull'organismo di reale interesse (come ad esempio l'uomo). Le scoperte compiute su organismi modello hanno spesso contribuito in modo significativo alla cura di malattie umane ed alla comprensione di processi di fondamentale importanza in biologia.

1.2.2 Caratteristiche del *Drosophila melanogaster*

Gli organismi modello vengono impiegati da tempo per comprendere meglio concetti relativi a diverse discipline. Per esempio, il dittero *Drosophila melanogaster* (noto anche come: moscerino della frutta), è stato spesso impiegato come organismo modello nell'ambito della genetica e nella divulgazione biologica. Il moscerino della frutta, *Drosophila melanogaster* è da tempo uno degli organismi eucarioti più popolari per la ricerca: è stato oggetto di importanti indagini nel campo dell'ereditarietà dei caratteri e nelle ricerca biomedica. *Drosophila* costituisce un modello eccellente per la comprensione dei principi della genetica mendeliana. Il successo di quest'organismo è legato in parte alla brevità del suo ciclo vitale (2 settimane da uovo ad adulto), che consente quindi di svolgere studi anche di diverse generazioni nell'arco di un singolo anno. Esso è un organismo facile da allevare e si può mantenere in grandi numeri con costi contenuti. Le larve mature mostrano cromosomi politenici nelle ghiandole salivari; hanno solo 4 paia di cromosomi: 3 autosomi e 1 sessuale; inoltre i maschi non mostrano ricombinazioni genetiche, facilitando gli studi genetici. Le sue dimensioni sono tali da consentirne l'osservazione di molte caratteristiche fenotipiche ad occhio nudo o a basso ingrandimento. I moscerini della frutta di tipo selvatico hanno gli occhi color rosso mattone, sono di colore giallo-marrone, e presentano anelli trasversali neri in tutto l'addome. I maschi si distinguono facilmente dalle femmine sulla base di differenze di colore, hanno una macchia nera distinta sull'addome, meno evidente nelle mosche recentemente osservate, e le *sexcombs* (una fila di setole scure sul tarso della prima zampa). Inoltre, i maschi hanno un gruppo di peli appuntiti (appendici) che circondano le parti riproduttive che sono utilizzati per attaccarsi alla femmina durante l'accoppiamento. Il periodo di sviluppo di *Drosophila melanogaster* varia con la temperatura, come per molte specie ectotermici. Il minor tempo di sviluppo (da uovo ad adulto), 7 giorni, si raggiunge a 28°C. Le femmine depongono circa 400 uova (embrioni), circa cinque alla volta. Le uova, di circa 0,5 millimetri di lunghezza, si schiudono dopo 1215 ore. Le larve risultanti crescono



Figura 1.: Esemplare di adulto di *Drosophila melanogaster*. (2.5 x 0.8 mm)

per circa 4 giorni, mentre la muta due volte (in larve di 2° e 3° instar), a circa 24 e 48 ore dopo la schiusa. Durante questo tempo, si nutrono dei microrganismi che decompongono la frutta, così come per gli zuccheri del frutto stesso. Quindi le larve si incapsulano nel puparium e subiscono una metamorfosi lunga quattro giorni, dopo di che emerge un adulto.

1.2.3 Il moscerino della frutta modello della ricerca

Il fatto che questo insetto sia stato fra i modelli più gettonati nella ricerca biologica rientra fra le tante caratteristiche che lo rendono un interessante modello di analisi, ricco di letteratura da sfruttare come base da cui partire per effettuare nuovi test. La *Drosophila* ha una lunga storia nell'ambito della ricerca biologica (sin dai primi '900); di conseguenza sono disponibili informazioni dettagliate su moltissimi aspetti della biologia di quest'organismo, oltre a raccolte di mutanti e di ceppi con caratteristiche particolari. Da un punto di vista genetico l'uomo e il moscerino della frutta sono abbastanza simili. Circa il 60% delle malattie genetiche conosciute si possono verificare nel patrimonio genetico del moscerino. E circa il 50% delle proteine della *Drosophila* hanno un analogo nei mammiferi. La *Drosophila* viene usata come modello genetico per varie malattie umane, inclusi i disturbi neurodegenerativi come il morbo di Parkinson, la corea di Huntington e il morbo di Alzheimer. La mosca viene utilizzata anche per studiare il meccanismo biologico del sistema immunitario, del diabete, del cancro, dell'intelligenza, dell'invecchiamento e persino dell'abuso di sostanze stupefacenti.

Parlando della *Drosophila*, è quindi impossibile non citare Thomas H. Morgan (USA, 1866 -1945), che con i suoi allievi formò il celebre gruppo della stanza dei moscerini, e che portò quest'insetto senza titoli a diventare una vera e propria star della genetica. Morgan non era convinto della teoria cromosomica dell'eredità, che in quegli anni

si andava diffondendo tra i genetisti e che si basava sull'assunto che i cromosomi, i filamenti presenti in tutti i nuclei cellulari, fossero la sede dei "fattori ereditari mendeliani", più tardi battezzati geni. Per dimostrare che non c'era nessuna prova incontrovertibile circa la validità di questa teoria, Morgan finirà, paradossalmente, per fornirne la prova cruciale. La *Drosophila* è caratterizzata dal fatto che delle sue quattro coppie di cromosomi solo una trasmette al nascituro il sesso e quei caratteri sessuali che permettono di distinguere un maschio da una femmina. Questa coppia è costituita dal cromosoma X e ricco di geni e dal cromosoma Y più piccolo e con pochi geni: se nelle cellule è presente la coppia XX nasce una femmina, altrimenti (XY) un maschio. Anche nella specie umana la determinazione del sesso avviene in modo simile a quello della *Drosophila*, ma non è così in tutti i viventi. Morgan studiando l'eredità di un mutante "occhi bianchi" trovò che, contrariamente ad altri caratteri mendeliani, i risultati degli incroci dipendevano dal sesso del genitore che portava tale carattere. Trovò che altri caratteri si comportavano allo stesso modo come, per esempio, la D. "corpo giallo" o la D. "ali miniatura", che vennero perciò chiamati caratteri (ovvero geni) legati al sesso, in quanto localizzati sul cromosoma X. Quindi è sui cromosomi che sono distribuiti i geni dei caratteri dell'individuo e in numero ben superiore a quello dei cromosomi. Morgan e i suoi allievi, ipotizzando per la prima volta l'associazione di più geni su uno stesso cromosoma (linkage), si accorsero che durante la meiosi i geni localizzati sui cromosomi spesso si separavano (linkage incompleto). La ricombinazione di geni associati, spiegò Morgan, è il risultato del crossing-over (scambio di parti tra cromosomi omologhi) e la quantità di linkage incompleto rappresenta la distanza lineare tra i geni lungo i cromosomi.

In origine, come si è detto sopra, questo organismo è stato utilizzato soprattutto nel campo della genetica, per esempio, per scoprire che i geni erano legati alle proteine e di studiare le regole di eredità genetica. Più di recente, è stato usato per lo più in biologia dello sviluppo, cercando di osservare come un organismo complesso nasca da una struttura relativamente semplice di uovo fecondato. L'attenzione è posta maggiormente sullo sviluppo embrionale, ma c'è anche un grande interesse sul come varie strutture adulte si sviluppino durante lo stadio pupa. Grande interesse è focalizzato principalmente sullo sviluppo della composizione dell'occhio, ma anche sulle ali, gambe e altri organi.

I microscopi a fluorescenza ottica sono diventati strumenti essenziali nella tecnica di imaging per la ricerca in biologia, in particolare nel campo della biologia dello sviluppo. Recenti progressi tecnologici hanno permesso lo sviluppo di un tipo di microscopi più veloci che offre una maggiore risoluzione per raccogliere immagini di campioni

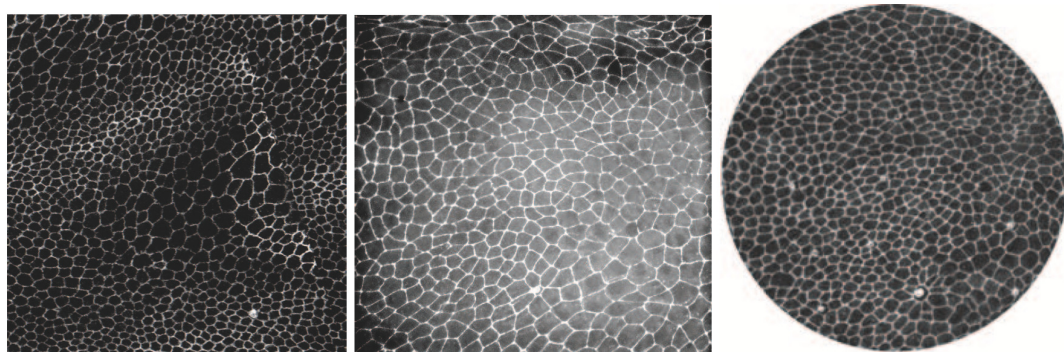


Figura 2.: Uno scorcio di epitelio di *Drosophila*. Epitelio di *Drosophila* osservato in un arco temporale da un microscopio confocale a scansione laser, durante la fase di sviluppo embrionale. In bianco: si può tracciare la proteina E-caderina marcata con una proteina fluorescente bianca, che svolge un ruolo importante nella adesione cellulare

biologici viventi. In parallelo, importanti progressi sono stati fatti nello sviluppo di sonde fluorescenti stabili e luminose. Un microscopio confocale ci ha fornito una sequenza video di quaranta frame che mostra la fase di sviluppo embrionale (vedi figura 2). Le cellule vengono evidenziate in bianco mediante una proteina marker.

1.3 VISIONE COMPUTAZIONALE E BIOLOGIA

Negli ultimi decenni i grandi progressi ottenuti nel campo dell'informatica hanno rivoluzionato la nostra capacità di ottenere ed analizzare i dati medici e biologici. Ciò ha avuto impatto su una vastissima gamma di applicazioni nel campo della ricerca in campo biologico e medico. Si può affermare che la visione computazionale ad oggi è diventata un grande alleato per i ricercatori biologi e gli operatori medici, nonché una tecnologia di supporto per il personale tecnico che lavora al fianco dei medici. In particolare i ricercatori che lavorano nel campo della ricerca sulle cellule staminali esprimono la necessità di poter usufruire di metodologie e procedure per quantificare i cambiamenti che avvengono a cellule e tessuti in seguito a stimoli specifici o lungo la loro naturale evoluzione.

Si definisce *Biologia Computazionale* la branca scientifica che si occupa di utilizzare l'insieme degli strumenti informatici, delle relative applicazioni e delle ricerche di base e applicate che possono essere sviluppate tramite di essi e alla definizione di algoritmi di analisi, per comprendere più velocemente e facilmente i complessi problemi che la ricerca di laboratorio affronta ogni giorno. In questo lavoro in particolare si è analizzato un video di morfogenesi dell'ala di una *Drosophila melanogaster*. Recentemente, sono stati studiati da mo-

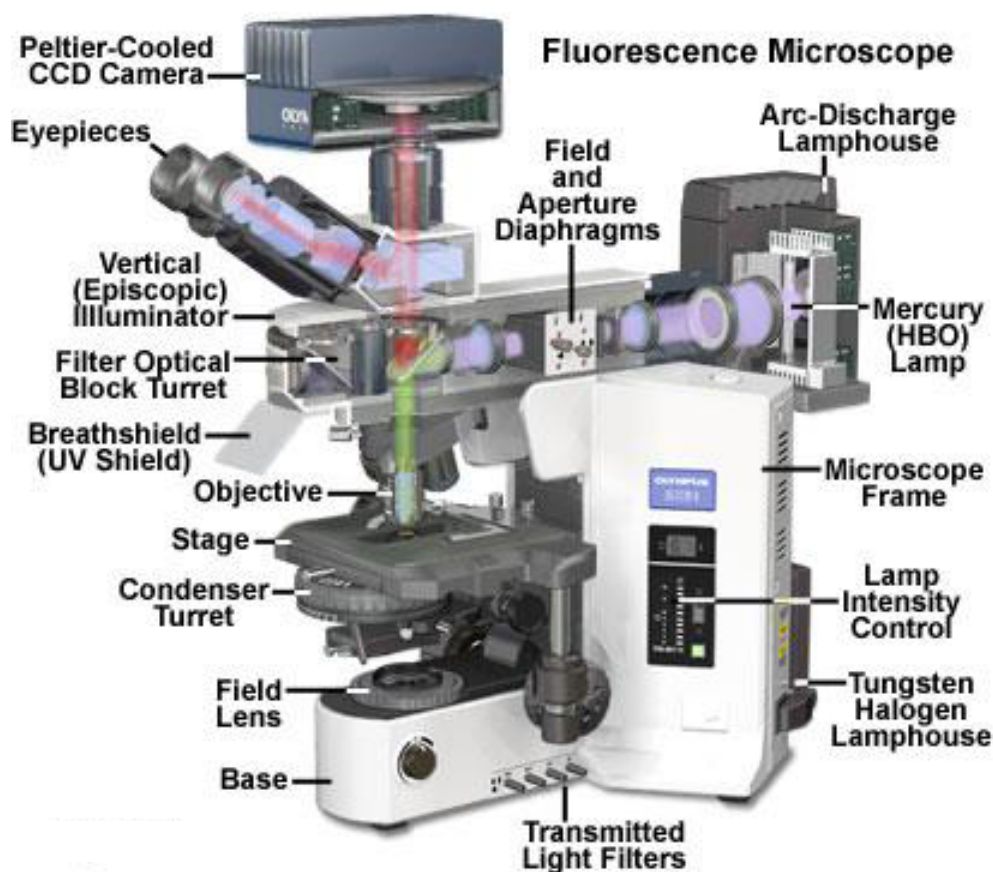


Figura 3.: Esempio di microscopio a fluorescenza ottica

delli matematici alcuni pattern di evoluzione della *Drosophila*, con l'obiettivo di ottenere una conoscenza quantitativa e più approfondita riguardo la morfogenesi di questo insetto. Vi è una necessità di conoscere in maniera accurata la dinamica della struttura cellulare epiteliale e l'organizzazione all'interno dell'ala mosca, per favorire la comprensione di un fenomeno noto come polarità cellulare planare.

1.3.1 Cenni su morfogenesi dei tessuti e migrazione cellulare

Le cellule contengono dei meccanismi complessi che ne permettono la motilità e questa caratteristica è uno dei passaggi cruciali dell'evoluzione. Con molta probabilità, infatti, le cellule primordiali avevano scarse caratteristiche motorie e venivano per lo più trasportate dal mezzo in cui vivevano; in seguito, con la specializzazione delle funzioni, le cellule hanno acquisito un grande capacità di muoversi autonomamente, per dirigersi in zone più favorevoli alla loro crescita e un raffinato sistema di trasporto e movimento delle sostanze al loro interno. L'evoluzione in organismi pluricellulari e complessi ha reso necessario che gruppi di cellule diverse raggiungano zone diverse

durante l'*embriogenesi*¹ e che le cellule deputate alla difesa contro le infezioni raggiungano e distruggano gli agenti patogeni. I movimenti non sono casuali ma sono frutto di raffinati e precisi sistemi che sono strettamente controllati dalla cellula. A livello microscopico i meccanismi del movimento hanno in generale necessità di energia che nella cellula è fornito dall'ATP² ed esistono una serie di proteine che convertono l'energia di questa molecola in movimento. La struttura del citoscheletro è coinvolta in molti processi cellulari e rappresenta l'impalcatura della cellula; essa è tuttavia una struttura dinamica che va incontro a continui riarrangiamenti che producono movimento.

Come si è detto la morfogenesi (generazione di nuova forma) è lo sviluppo della conformazione e della struttura di un organismo, ed assieme all'accrescimento (aumento dimensioni) costituisce l'insieme dei processi tramite i quali le cellule si organizzano a formare tessuti. La dinamica di rimodellamento dei tessuti durante la morfogenesi e la migrazione cellulare giocano un ruolo centrale in una varietà di fenomeni biologici. In embriogenesi, le migrazioni cellulari sono un tema ricorrente nei processi morfogenetici importanti che vanno dalla gastrulazione allo sviluppo del sistema nervoso. La migrazione rimane un fenomeno importante nell'organismo adulto, nel suo normale contesto fisiologico e patologico. Nella risposta infiammatoria, ad esempio, i leucociti migrano in aree di infezione, dove mediano funzioni fagocitarie ed immunitarie. Ancora, la migrazione di fibroblasti e di cellule endoteliali vascolari è essenziale per la guarigione della ferita. In metastasi, cellule tumorali migrano dalla massa tumorale iniziale attraverso il sistema circolatorio. Infine, la migrazione cellulare è fondamentale per applicazioni tecnologiche come l'ingegneria dei tessuti, giocando un ruolo essenziale nella colonizzazione dei ponteggi di biomateriali. Come per molti altri processi cellulari, i componenti molecolari coinvolti nella migrazione cellulare sono stati individuati in tempi piuttosto rapidi, mentre la determinazione di come essi partecipano effettivamente alla migrazione è fatto un po' più complesso e gli studi circa questi comportamenti sono avvenuti in tempi più lunghi. Il modo in cui questi componenti evolvono insieme come un processo dinamico, un sistema integrato per dare origine alla migrazione lascia ancora molto spazio alla ricerca. Comprendere la migrazione cellulare come un processo integrato infatti richiede una valutazione delle proprietà chimiche e fisiche delle strutture multicomponenti e gruppi, compresa la loro termodinamica, cinetica, e le

1 L'embriogenesi, detta anche ontogenesi o sviluppo embrionale, è lo sviluppo dell'embrione dall'uovo fecondato e consiste nell'ordinata sequenza dei fenomeni di accrescimento, differenziamento e di organogenesi che conducono alla formazione di un individuo (segmentazione, gastrulazione e organogenesi).

2 Molecola chiamata *adenosin trifosfato* è costituita da una base azotata (adenina), da uno zucchero a cinque atomi di carbonio (ribosio) e da tre gruppi fosfato. Quando un gruppo fosfato si stacca per idrolisi si liberano circa 7 kilocalorie di energia per mole.

caratteristiche meccaniche, perché la migrazione è un processo che è fisicamente coordinato nel tempo e nello spazio.

La locomozione cellulare è il risultato di movimenti generati in diverse parti della cellula e, per esempio, in un leucocita (la cellula che deve muoversi all'interno del corpo per cercare batteri che invadono i tessuti), i filamenti di actina producono una grande protrusione della membrana chiamata pseudopodio. Questa struttura si fissa al substrato (figura 4) e la parte della membrana che lo circonda viene riempita di citoplasma; la cellula avanza ripetendo questa operazione più volte.

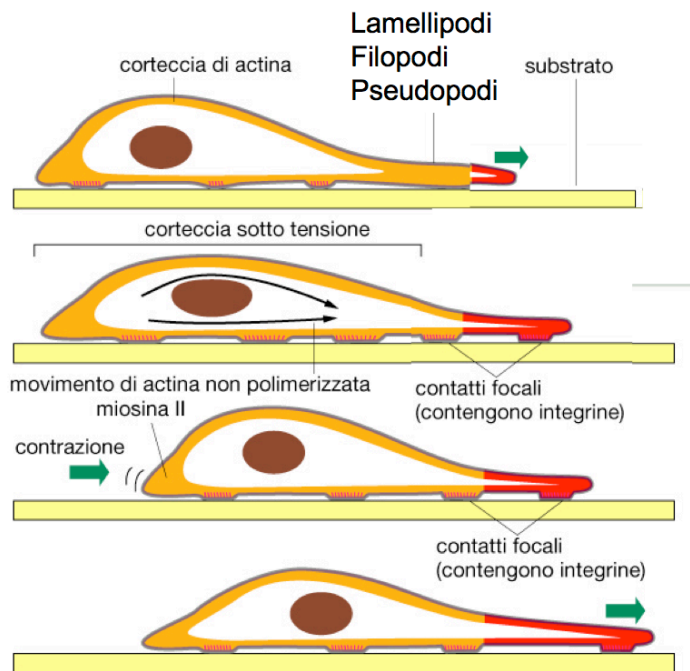


Figura 4.: Meccanismi di locomozione cellulare: protrusione ossia formazione di protrusioni che esplorano l'ambiente; ancoraggio cioè adesione al substrato e stiramento; trazione e traslazione.

Una delle sfide più importanti per la biologia dello sviluppo è capire come l'informazione molecolare conduce al movimento individuale e collettivo delle cellule che formano i tessuti sia attraverso lo stress intrinseco della cellula che alle risposte cellulari allo stress applicato. Con gli enormi progressi nel campo della biologia molecolare, della genetica, delle tecniche di imaging e il monitoraggio automatizzato di strutture pluricellulari, è ora possibile tracciare l'evoluzione morfo-genetica di alcuni fenotipi durante lo sviluppo, in funzione di perturbazioni molecolari e manipolazioni fisiche. Questo pone le basi per l'identificazione e la quantificazione della geometria dei cambiamenti di forma in termini tassi di deformazione per unità di tempo. Gli attuali approcci per la caratterizzazione statistica dei parametri morfo-genetici si basano sulla somiglianza dei tessuti a schiume e

materiali granulari, per cui vengono utilizzati metodi topologici che si basano sulla connettività.

Con un tool di visione di bio-immagini, si cerca quindi di estrarre dei dati quantitativi sul movimento e la deformazione cellulare; questi dati devono essere sintetizzabili in statistiche osservabili e analizzabili da biologi e ricercatori. Il contesto scientifico in cui si collocano queste tecniche di *bio-imaging* è variegato. Una importantissima applicazione che si sta sviluppando negli ultimi anni è quello dell'ingegneria tessutale o più in generale quello della medicina rigenerativa. La medicina rigenerativa si pone come obiettivo principale la riparazione di organi e tessuti danneggiati da eventi patologici, invecchiamento o traumi in maniera da ripristinare o migliorare il loro funzionamento biologico. Il termine "medicina rigenerativa" viene comunemente utilizzato per indicare quelle strategie mediche in ambito terapeutico o di ricerca che fanno uso dello straordinario potenziale di un particolare tipo di cellule, le cellule staminali, progenitori immaturi dotati del potenziale di differenziarsi nei diversi tipi cellulari.

La medicina rigenerativa (o ingegneria tessutale) è un campo multidisciplinare in rapida crescita che coinvolge le scienze mediche, umane ed ingegneristiche e che cerca di sviluppare cellule funzionali, tessuti o sostituti di organi, allo scopo di riparare, rimpiazzare o migliorare le funzioni biologiche che sono state perse a causa di anomalie congenite, traumi, malattie o come conseguenza dell'invecchiamento.

Lo studio delle cellule staminali permette ai ricercatori di indagare anche quali sistemi si attivano in tali cellule quando riparano un danno all'interno di un organismo. Le conoscenze attuali sulla distrofia muscolare e sull'atrofia muscolare sono un esempio dei risultati di questo tipo di studi.

In questo campo è chiaro che il ruolo delle tecniche ingegneristiche di visione computazionale dà un contributo fondamentale nello studio di sistemi biologici in evoluzione. La tecnologia supporta una svariata quantità di indagini visive in questo campo, perciò si rende necessaria la progettazione di algoritmi di visione robusti e il più possibile versatili. Esistono aziende che in questo settore già commercializzano sistemi di visione automatica per colonie cellulari. Sono sistemi che probabilmente risentono ancora molto della variabilità di dati cui vengono sottoposti (infatti esistono alcuni tipi di colonie cellulari che, anche per un tecnico esperto, sono visivamente ancora difficili da catalogare). La sfida futura sarà quella di superare queste barriere e implementare sistemi di visione automatica effettivamente robusti ed efficienti.

2

RICHIAMI DI VISIONE COMPUTAZIONALE

Indice

2.1	Riconoscimento e analisi di forme	12
2.2	Modelli di rappresentazione di forme	14
2.3	<i>Single shape detection</i>	16
2.3.1	<i>Standard Active Contours (snake)</i>	17

La visione computazionale (o visione artificiale; in inglese "*computer vision*") si occupa dello sviluppo di algoritmi e sistemi per l'analisi di immagini. Essa è complementare alla "*computer graphics*", che invece si occupa della sintesi di immagini. Semplificando, in *computer vision* telecamere e computer svolgono il ruolo che nella visione biologica hanno occhi e cervello. Le applicazioni della visione computazionale sono pressoché inesauribili. Alcuni temi sono la sintesi visiva di scene a partire da collezioni di immagini, l'acquisizione della forma (anche tridimensionale) e tessitura di oggetti da sequenze video, l'elaborazione di tecniche super-risoluzione, la cattura visiva di movimenti e gesti del corpo umano per lo sviluppo di interfacce naturali uomo-macchina, il controllo di robot basato su visione . . . Il crescente sviluppo e l'espansione delle tecnologie di imaging medicale sta rivoluzionando la medicina. Questo permette a scienziati e fisici di raccogliere informazioni potenzialmente in grado di salvare vite attraverso un approccio non invasivo di osservazione del corpo umano. Il ruolo dell'imaging medicale si è espanso oltre alla semplice visualizzazione e analisi delle forme anatomiche. Esso è diventato uno strumento per pianificazioni chirurgiche e simulazioni, navigazione intra-operativa, pianificazione radioterapica, e per tracciare lo sviluppo di una certa malattia. Per esempio accertare la forma dettagliata e l'organizzazione delle strutture anatomiche prima dell'intervento consente ad un chirurgo di pianificare l'approccio ottimale per operare. In radioterapia l'imaging medicale permette la somministrazione di una dose di radiazioni minimizzando i danni collaterali ai tessuti sani. Con il prominente ruolo giocato da questa tecnologia sulla diagnostica e trattamento delle malattie, la comunità di analisi di immagini biomedicali si è interessata al problema dell'estrazione, senza l'ausilio dei computers, di informazioni cliniche sulle strutture anatomiche attraverso modalità come TAC, RM, PET. Sebbene i moderni dispositivi di imaging garantiscano visioni eccezionali delle anatomie interne, l'uso dei computers per analizzare e quantificare tali strutture con accuratezza ed efficienza è limitato. Ciò che conta è riuscire ad estrarre dati

che siano accurati, ripetibili in modo da assistere l'ampia gamma di ricerche biomediche e attività cliniche per la diagnosi, la radioterapia e la chirurgia.

2.1 RICONOSCIMENTO E ANALISI DI FORME

Chiunque crei o utilizzi algoritmi o dispositivi per l'elaborazione di immagini digitali dovrebbe tenere conto dei principi della percezione dell'immagine umana. Se l'immagine deve essere analizzata da un essere umano le informazioni devono essere espresse utilizzando variabili che siano facili da percepire, che siano parametri psico-fisici quali contrasto, bordo, forma, consistenza, colore, ecc. Gli esseri umani sono in grado di trovare oggetti all'interno di immagini solo se ne possono distinguere facilmente dallo sfondo.

Il punto di partenza è specificare che cosa si intende per "forma" (in inglese *shape*). Non esiste un unico modo per definire questa entità, tipicamente per forma si intende la "parte di spazio occupata da un oggetto", determinata dal suo bordo esterno, astraendo questo concetto da proprietà come colore, contenuto e composizione fisica e astraendolo allo stesso tempo da proprietà spaziali proprie dell'oggetto come posizione, grandezza e orientazione nello spazio. Si può dare una definizione di *shape* nel seguente modo:

Definizione 1 (Shape). *Shape* è tutta l'informazione geometrica che rimane quando posizione, scala ed effetti rotazionali vengono filtrati dall'oggetto.

In altre parole, in base a questa definizione la forma è invariante rispetto alle trasformazioni di similarità Euclidee e ciò si intuisce più facilmente se si guarda la figura 5.

Denoteremo con X l'oggetto di interesse e con $\Phi(X)$ la sua forma e intenderemo con "intuizione" il processo di estrazione di $\Phi(X)$ da X . Come si è detto sopra questo concetto di estrazione della forma è fortemente radicato nella capacità umana di riconoscimento di pattern. In realtà il processo di intuizione della forma non avviene direttamente da X bensì dall'immagine di X impressa nella nostra retina attraverso gli occhi, che possiamo immaginare come il dispositivo o sensore di acquisizione delle immagini. Formalmente identifichiamo con I l'uscita dei sensori e quindi il processo di intuizione della forma sarà del tipo

$$X \xrightarrow{\text{sensori}} I \xrightarrow{\text{intuizione}} \Phi(X) \quad (1)$$

La seconda domanda che sorge spontanea a questo punto riguarda il modo in cui è utile, da un punto di vista ingegneristico, descrivere una forma. A tal proposito si può affermare che un modo per descri-

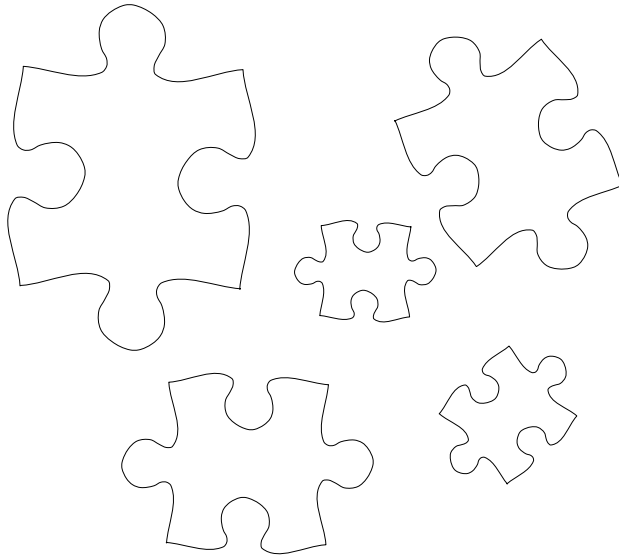


Figura 5.: Cinque copie della stessa forma, ma dopo differenti trasformazioni Euclidee.

vere una forma è quello di fissare un numero finito di punti *landmarks* sul contorno.

Definizione 2 (Landmark). *Landmark* è un punto di corrispondenza identificabile su ciascun oggetto della stessa famiglia di forme.

La fase di rappresentazione è assai delicata perché deve essere in grado di fornire un modello analitico \mathcal{R} per $\Phi(X)$. La scelta è cruciale poiché influenza pesantemente le fasi successive di analisi dell'immagine. La rappresentazione rende possibile mettere in pratica algoritmi di calcolo per l'analisi di immagini; si ha quindi la seguente successione di passi

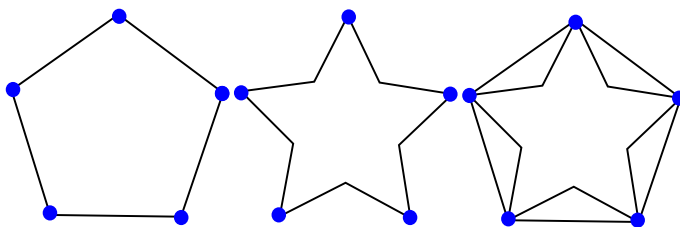


Figura 6.: La figura mostra due forme definite mediante la stessa scelta di landmarks. Risulta evidente che se la scelta dei landmarks è troppo grossolana si ha una perdita di informazione tale per cui non è più possibile riconoscere la differenza che in origine era visibile.

$$X \xrightarrow{\text{sensori}} I \xrightarrow{\text{intuizione}} \Phi(X) \xrightarrow{\text{rappresentazione}} \mathcal{R} \quad (2)$$

Come si evince dalla figura 6 la rappresentazione del modello, che conduce inevitabilmente a perdita di informazione, deve essere tale da definire uno spazio metrico (anche se non è richiesto esplicitamente che definisca un prodotto estero né una metrica di distanza). Una rappresentazione ideale dovrebbe garantire che

$$X_a = X_b \iff \Phi(X_a) = \Phi(X_b) \quad (3)$$

In letteratura sono stati proposti molti modelli \mathcal{R} di rappresentazione di forme: basati sulla scelta di specifici landmarks, su funzioni implicite, su curve chiuse \mathcal{C}^2 , su grafi ... In generale la scelta di una rappresentazione piuttosto che un'altra può condurre a risultati consistentemente differenti. Nel prossimo paragrafo si dà una breve panoramica degli approcci appena nominati.

RICONOSCIMENTO DI FORME Il riconoscimento (in inglese *shape detection*) è il processo complessivo che a partire dall'oggetto X fornisce la rappresentazione della sua forma $\Phi(X)$

$$X \xrightarrow[\Phi(X)]{\text{detection}} \mathcal{R} \quad (4)$$

ANALISI DI FORME (in inglese *shape analysis*) prende come punto di partenza la rappresentazione \mathcal{R} e restituisce una serie di output che possono essere di varia natura dipendentemente dal contesto applicativo. Il tipo di analisi e l'interpretazione di queste misure dipendono dalla specifica funzione e possono includere metriche come lunghezze, classificazioni, descrittori, particolari quantità riconducibili ad una diagnostica specifica etc. In questo lavoro ad esempio si analizza la deformazione delle cellule di *Drosophila melanogaster* per mezzo dell'estrazione di metriche rilevanti.

$$X \xrightarrow[\Phi(X)]{\text{detection}} \mathcal{R} \xrightarrow{\text{analysis}} \mathbb{R}^k \quad (5)$$

dove \mathbb{R}^k è il set di numeri che codificano l'output dell'analisi in questione e quindi corrispondono alle metriche scelte per computarla.

2.2 MODELLI DI RAPPRESENTAZIONE DI FORME

Come si è detto una "rappresentazione" di una forma è un modello che rappresenti il concetto astratto di forma $\Phi(X)$, tale cioè da

descrivere la categoria di forma che appartiene ad una certa configurazione prestabilita, come ad esempio conformazioni di tipo reticolare o superfici continue. Si presentano in seguito alcuni dei modelli di rappresentazione di forme proposti in letteratura.

SUPERFICI \mathcal{C}^d Le forme sono incluse in uno spazio d-dimensionale come la superficie continua $\mathcal{C}^d \in \mathcal{R}^d$ del tipo $\mathcal{C}^d(s_1, s_2, \dots, s_d)$ dove $s_i \in [0, 1] \mid \text{for } i = 1, \dots, d$. Le curve e le superfici *spline* sono l'implementazione più comune.

FUNZIONI $\mathcal{R}^{n \times m}$ Le funzioni continue modellizzano $\Phi(X)$ come un set di n funzioni m-dimensionali. Ogni funzione è legata puntualmente alla *shape* e ne descrive un aspetto del contorno in quel punto. Un esempio di funzione è quello della distanza dal centroide che mette in relazione la distanza del centroide dal contorno (raggio medio o magnitudine) con l'angolo formato da esso $r_x(\phi)$.

GRAFI $G = (V, E)$ Questo è un modello che utilizza i vertici (V) come punti rilevanti insieme con le loro connessioni (edges E). Modelli come questo basati su grafi sono particolarmente adatti a rappresentare strutture reticolari come le cellule di epitelio (*Drosophila wing*) che si analizzano in questa tesi.

LANDMARK Si assume che una shape sia rappresentabile mediante le coordinate di un set finito di landmarks $L = \{l_1, \dots, l_n\}$ che sono punti particolari che evidenziano determinate particolarità dell'oggetto osservato. Esistono landmark di tipo *anatomico* che sono punti corrispondenti fra diversi organismi e significativi in qualche modo biologico (come gli angoli della bocca, le narici etc.); landmark di tipo *matematico* che corrispondono a punti che obbediscono a determinate proprietà geometriche o matematiche (come curvature, punti estremi, etc.).

FUNZIONI IMPLICITE ϕ Sono funzioni del tipo $\phi(x, y) = x^2 + y^2 - 1$ che definiscono la forma per mezzo dei punti $(x, y) \in \mathcal{R}^2$ tali che $\phi(x, y) = 0$. Formalmente supponendo di avere una funzione $\phi(x)$, una shape è definita come il set di tutti i punti x appartenenti all'isocontorno $\phi(x) = 0$. Un esempio di questo tipo di modello è mostrato in figura 7.

MODELLI DEFORMABILI Sono detti anche *template* e sono modelli standard di forme associati a diverse classi di oggetti. Questi modelli contengono un certo numero di parametri e una funzione di costo, dipendente da questi parametri, che rappresenta la bontà della rappresentazione per il dato oggetto. Sono rappresentazioni adeguate a situazioni applicative in cui è abbondante la conoscenza a priori sul sistema e le forme da riconoscere sono

definite in maniera chiara e definita (es. problemi di "face detection" o più in generale di "face's features detection"). Un esempio di questo tipo di modello è mostrato in figura 8.

2.3 *single shape detection*

In questa sezione si espone il problema di cercare una (singola) *shape* Φ all'interno di un'immagine digitale I . Questa fase ci riporta alla scelta di una $\mathcal{R}(\Phi_c) \in S$ ottima tale da minimizzare un certo funzionale di energia $E(I, \Phi_c)$ sul set S di tutte le possibili rappresentazioni

$$\mathcal{R}(\Phi) = \arg \min_{\mathcal{R}(\Phi_c) \in S} E(I, \Phi_c) \quad (6)$$

E è una funzione scalare dell'immagine I e della *shape* Φ_c che si cerca di tracciare. Nel seguito si farà riferimento a Φ intendendo $\mathcal{R}(\Phi)$ in modo da rendere più leggibili le notazioni. Quindi l'equazione 6 diventa

$$\Phi = \arg \min_{\Phi_c \in S} E(I, \Phi_c) \quad (7)$$

Anche se la dichiarazione appare semplice nella realtà ci sono due fattori principali che giocano un ruolo essenziale nella costruzione di una soluzione a questo problema. In primo luogo il calcolo di un minimo globale è quasi sempre impossibile. In altre parole, anche per problemi banali ma con una struttura ben definita (figura 9), è impossibile trovare una forma chiusa per la minimizzazione o la conoscenza esaustiva del set S . A seconda della natura del problema, la cardinalità di S potrebbe essere addirittura infinita. Quindi nella maggior parte dei casi applicativi gli algoritmi condurranno a soluzioni non ottimali in senso stretto. In secondo luogo è importante osservare che la definizione di E non è unica, inoltre non esistono in linea di principio modalità generali su come costruirlo. Si può dire che una buona scelta di E è un funzionale che sia concavo avente un

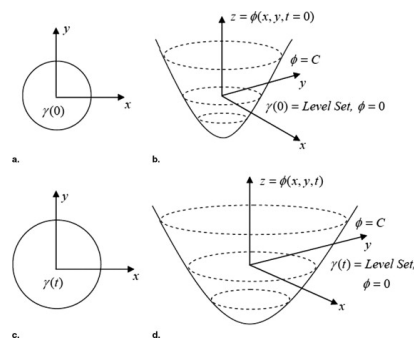


Figura 7.

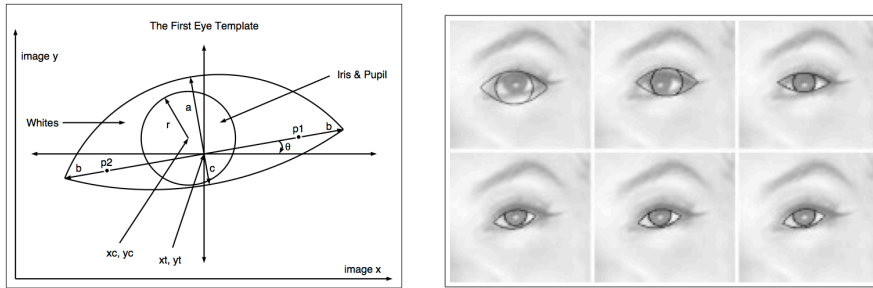


Figura 8.: Primo template per l'occhio umano definito attraverso variabili che rappresentano parametri geometrici. Si cerca una configurazione di equilibrio che "avvolga" l'occhio in maniera più precisa possibile. Immagine presa da *Active Contour*.

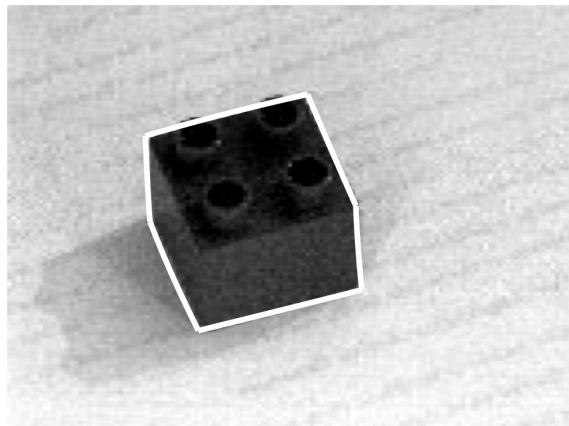


Figura 9.: Esempio di riconoscimento di un oggetto.

solo minimo e che sia in qualche maniera correlato alla percezione (ad alto livello) che si ha di Φ .

Anche ammesso che esista una funzione ideale di rappresentazione da ricercare all'interno del set S non è garanzia del successo della procedura globale di detection, dal momento che anche la definizione di E contribuisce alla definizione dei limiti qualitativi di scelta. In maniera analoga si pensi di avere a disposizione una pessima conoscenza di S e quindi di avvalersi di una inadeguata procedura di ricerca su S , anche avendo a disposizione un funzionale E ottimo non sarebbero garantiti buoni risultati. Ciò implica che un buon algoritmo di shape detection debba basarsi su un'implementazione il più accurata possibile di entrambi questi aspetti.

2.3.1 Standard Active Contours (snake)

Uno dei compiti generali nell'elaborazione delle immagini è quello di riuscire a dividere l'immagine in più parti. Si tratta di separare gli "oggetti", riferiti anche con l'espressione "primo piano", dall'insieme dei pixels dello sfondo. Questo processo è detto *segmentazione*.

Esistono due modi per definire i segmenti di un'immagine: il primo è quello di individuare i pixels appartenenti all'area dell'oggetto; il secondo modo è di definire i confini della zona di interesse. Esiste un ampio spettro di tecniche di segmentazione ben note come quella della *simple thresholding* fino a diversi tipi di tecniche di machine learning, come i metodi di clustering o segmentazione basata su reti neurali. I modelli di contorno attivi appartengono alla classe dei metodi di rilevamento del bordo. Il primo modello di snake è stato introdotto da Michael Kass, Andrew Witkin e Demetri Terzopoulos nel 1988 (Kass *et al.*, 1988). Anche se originariamente vennero sviluppati per applicazioni a problemi di computer vision e computer graphics, la potenzialità dei modelli deformabili per l'uso in analisi di immagini mediche è stato subito capito. Essi sono stati applicati alle immagini generate mediante tecniche di imaging diverse, come i raggi X, la tomografia computerizzata (CT), l'angiografia, la risonanza magnetica (RM), e gli ultrasuoni. Sono stati utilizzati modelli deformabili a due o tre dimensioni per segmentare, visualizzare, monitorare e quantificare una serie di strutture anatomiche che vanno in scala dal macroscopico al microscopico. Questi includono il cervello, il cuore, la faccia, le arterie coronarie e della retina, i reni, i polmoni, lo stomaco, il fegato, il cranio, le vertebre, oggetti come tumori cerebrali, un feto, e anche le strutture cellulari come i neuroni e i cromosomi. Gli *active contours* sono stati utilizzati per monitorare il movimento non rigido del cuore o il moto di eritrociti. Ancora sono stati utilizzati per localizzare strutture nel cervello, e registrare immagini di tessuto della retina, della vertebra e tessuto neuronale. Il contesto applicativo in cui si colloca questo lavoro, come più volte detto, sarà quello dell'applicazione di modelli deformabili per l'interpretazione di immagini tratte da un video di embriogenesi di *Drosophila Melanogaster*; in particolare di cercherà di ottenere una buona segmentazione di immagini per arrivare a fare un'analisi del movimento cellulare.

Gli *Active Contours* sono una tecnica di segmentazione che pone le sue basi sulla minimizzazione di un certo funzionale di energia che soddisfi a determinati vincoli esterni e sia influenzata dalle "forze" dell'immagine. L'algoritmo agisce deformando un contorno con l'obiettivo di agganciare certe caratteristiche di interesse su un'immagine. Di solito le caratteristiche di interesse sono i bordi, le linee o i confini di una forma. La procedura funziona in 2D, 3D o in dimensioni ancora più elevate. La versione 3D è spesso chiamato modelli deformabili o superfici attive.

Formalmente gli *active contours*, chiamati anche *snakes*, sono curve che si muovono all'interno delle immagini per trovare i contorni degli oggetti:

$$C(s) \in \mathbb{R}^d, \quad s \in [0, 1]^d \quad (8)$$

La forma del contorno $C(s)$ di un'immagine evolve (in pseudo-tempo) secondo il funzionale di energia $E(C)$:

$$E(C) = S(C) + P(C) \quad (9)$$

dove

$$S(C) = \int_{\partial C} \alpha(s) \left| \frac{\partial C}{\partial s} \right|^2 + \beta(s) \left| \frac{\partial^2 C}{\partial s^2} \right|^2 \quad (10)$$

$$P(C) = \int_{\partial C} \mathcal{D}[I(C(s))] ds \quad (11)$$

$S(C)$ è detta energia interna ed impone una accentuata curvatura (*smoothness*) alla curva C . Ciò si riferisce all'ipotesi di alto livello che indica che superfici lisce appartenenti al mondo reale vengano mappate in "dati smussati" del sensore. Il termine $\left| \frac{\partial C}{\partial s} \right|^2$ obbliga il risultante snake ad essere una curva corta, mentre il termine $\left| \frac{\partial^2 C}{\partial s^2} \right|^2$ impone valori di curvatura bassi. In altre parole il primo contributo attribuisce energia alta a contorni elongati (forza elastica) mentre quest'ultimo attribuisce energia alta a contorni ad elevata curvatura (forza rigida). I termini $\alpha(s)$ e $\beta(s)$ sono scalari che si occupano di bilanciare la grandezza del termine energia in ogni segmento di C .

L'energia esterna $P(C)$ è minima quando C si trova sul contorno dell'oggetto da riconoscere. L'approccio più semplice consiste nella scelta di $D(I) = -|\nabla I|$. Di nuovo, questa è una assunzione sulla natura dei dati, che impone che il contorno dell'oggetto reale venga mappato in punti dell'immagine di derivata elevata. A regime questo vincolo spinge la forma finale ottenuta verso i bordi dell'immagine. Secondo il principio variazionale le forze agenti su $C(s)$ possono essere derivate come:

$$F(s) = \frac{\partial E}{\partial s} \quad (12)$$

$$F(s) = \frac{\partial S(C)}{\partial s} + \frac{\partial P(C)}{\partial s} \quad (13)$$

$$F(s) = -\nabla D [I] + 2 \frac{\partial}{\partial s} \left(\alpha(s) \frac{\partial C}{\partial s} \right) + 2 \frac{\partial^2}{\partial s^2} \left(\beta(s) \frac{\partial^2 C}{\partial s^2} \right) \quad (14)$$

e la forma a riposo minima viene raggiunta per $F(s) = 0$ e cioè per:

$$2 \frac{\partial}{\partial s} \left(\alpha(s) \frac{\partial C}{\partial s} \right) + 2 \frac{\partial^2}{\partial s^2} \left(\beta(s) \frac{\partial^2 C}{\partial s^2} \right) - \nabla D [I] = 0 \quad (15)$$

Questa equazione differenziale parziale esprime l'equilibrio delle forze interne ed esterne, quando il contorno si ferma all'equilibrio. Richiamando ciò che si è detto sopra diciamo che i primi due termini

rappresentano le forze interne di allungamento e curvatura, rispettivamente, mentre il terzo termine rappresenta le forze esterne che agganciano lo snake ai dati di immagine.

Gli *active contours* si comportano come nel quadro fornito alla sezione precedente dove $E = E(C) = S(C) + P(C)$ ed $\mathcal{R}(\Phi) = C(s)$. In questo caso il funzionale di energia è ben definito ed è in grado di essere adattato ad una serie di applicazioni, semplicemente accordando i parametri $\alpha(s)$ e $\beta(s)$.

3

TRACKING E ANALISI DI UNA CELLULA

Indice

3.1	<i>Dynamic shape detection</i>	21
3.1.1	<i>Optical Flow</i> e J-maps	23
3.1.2	Algoritmo di <i>tracking</i> della struttura reticolare	25
3.2	<i>Shape Deformation Analysis</i>	25
3.2.1	Moto e deformazione	27

In questo capitolo si illustrerà come estendere il quadro esposto al capitolo precedente di *single shape detection* ad un contesto dinamico, ossia ad una sequenza di frames piuttosto che ad un'immagine statica grazie all'approccio descritto in [Silletti, 2007-2009](#). Successivamente si illustreranno le metriche e i parametri di forma che si sono utilizzati per *analizzare* la dinamica di una specifica cellula, osservata in un video di morfogenesi epiteliale di *Drosophila*. Le metriche e le statistiche scelte sono state poi messe a confronto per cercare di estrarne delle statistiche evolutive.

L'idea di tracciare oggetti in immagini tempo-varianti usando modelli deformabili venne inizialmente proposta nel contesto generico del Computer Vision ([Kass et al., 1988](#); [Terzopoulos e Kurt, 1988](#)). I modelli deformabili sono stati utilizzati per tracciare strutture microscopiche e macroscopiche non rigide in movimento, come le cellule sangue, i coni di crescita dei neuroni in cine-microscopia o le arterie coronarie in cine-angiografia. Ad ogni modo il primo campo in cui vennero applicati i modelli deformabili per il tracking di immagini mediche fu in cardiologia e in particolare per misurare il comportamento dinamico del cuore umano, in particolare del ventricolo sinistro. La descrizione del movimento della parete cardiaca è necessario per escludere la gravità e la portata di malattie come ischemia e quindi si capisce l'importanza che ha lo studio di queste tecnologie di imaging.

3.1 *dynamic shape detection*

Consideriamo ora una sequenza di immagini $\{I_1, \dots, I_t, \dots, I_N\}$. Chiamiamo il problema *dynamic shape detection* in quanto ci riferiamo all'evoluzione (dinamica) dell'oggetto nel tempo. L'obiettivo è ora quello di riconoscere l'oggetto X in ogni immagine I_t ottenendo

una sequenza di shapes $\{\Phi_1(X), \dots, \Phi_t(X), \dots, \Phi_N(X)\}$. Quello che si fa è applicare l'equazione 7 ad ogni passo temporale:

$$\begin{aligned}\Phi_1(X) &= \arg \min_{\Phi_c \in S} E_1(\{I_1, \dots, I_t, \dots, I_N\}, \Phi_c) \\ &\dots \\ \Phi_t(X) &= \arg \min_{\Phi_c \in S} E_t(\{I_1, \dots, I_t, \dots, I_N\}, \Phi_c) \\ &\dots \\ \Phi_N(X) &= \arg \min_{\Phi_c \in S} E_N(\{I_1, \dots, I_t, \dots, I_N\}, \Phi_c)\end{aligned}\quad (16)$$

Assumiamo in prima analisi che E_t dipenda solo dall'istante di tempo corrente e cioè si abbia $E_t(\{I_1, \dots, I_t, \dots, I_N\}) = E_t(I_t)$ e quindi $E_1 = \dots = E_N$. Si può riscrivere la 16 alleggerendo la notazione come

$$\begin{aligned}\Phi_1(X) &= \arg \min_{\Phi_c \in S} E_1(I_1, \Phi_c) \\ &\dots \\ \Phi_t(X) &= \arg \min_{\Phi_c \in S} E_t(I_t, \Phi_c) \\ &\dots \\ \Phi_N(X) &= \arg \min_{\Phi_c \in S} E_N(I_N, \Phi_c)\end{aligned}\quad (17)$$

Arricchiamo il modello tenendo conto del fatto che l'immagine digitale I_t deve essere simile all'immagine del frame successivo I_{t+1} , cioè

$$I_{t+1} = I_t + \Delta_t^{t+1} \quad (18)$$

dove abbiamo introdotto questa "coerenza temporale" mediante il termine Δ_t^{t+1} il quale misura la differenza fra due frames consecutivi. Chiaramente si ha coerenza temporale alta per $|\Delta| \approx 0$, bassa per $|\Delta| \gg 0$. Analogamente si può fare per Φ_{t+1}

$$\Phi_{t+1} = \Phi_t + \delta_t^{t+1} \quad (19)$$

In base all'equazione 7, il valore ottimo di δ_t^{t+1} soddisfa a

$$\Phi_{t+1} = \arg \min_{\Phi_c \in S} E(I_{t+1}, \Phi_c) \quad (20)$$

di conseguenza la soluzione per 21 è

$$\begin{aligned}
\Phi_1 &= \arg \min_{\Phi_c \in S} E_1(I_1, \Phi_c) \\
&\dots \\
\Phi_2 &= \Phi_1 + \delta_1^2 \\
&\dots \\
\Phi_t &= \Phi_{t-1} + \delta_{t-1}^t \\
&\dots \\
\Phi_N &= \Phi_{N-1} + \delta_{N-1}^N
\end{aligned} \tag{21}$$

3.1.1 Optical Flow e J-maps

L'Optical Flow \mathcal{F} è un concetto applicato all'analisi del moto di un oggetto all'interno di una rappresentazione visuale digitale. Date due immagini I_t e I_{t+1} lo scopo del flusso ottico è quello di assegnare a ciascun pixel appartenente al frame corrente I_t un *motion vector* che punta verso la posizione dello stesso pixel in un frame di riferimento successivo I_{t+1} . Un punto di partenza comune per la stima dell'optical flow è di assumere che le intensità dei pixel siano trasmesse da un fotogramma al successivo. L'ipotesi è che la luminosità della superficie rimanga fissa da un fotogramma al successivo. Comunemente l'algoritmo Optical Flow ha come output una funzione $\mathcal{F} : \mathbb{R}^2 \rightarrow \mathbb{R}^2$ o meglio un campo vettoriale i cui elementi $\mathcal{F} = (\mathbf{V}_x, \mathbf{V}_y)$ descrivono lo spostamento del pixel (x, y) dell'immagine I_t nell'immagine I_{t+1} ; in simboli si ha:

$$I_t(x, y) = I_{t+1}(x + v_x, y + v_y) \tag{22}$$

Indicando con $I_t(x, y) = I(x, y, t)$ e per piccoli spostamenti riscriviamo 22 e otteniamo l'equazione di vincolo dell'immagine:

$$I(x, y, t) = I(x + \delta x, y + \delta y, t + \delta t) \tag{23}$$

Espandendo $I(x, y, t)$ come serie di Taylor otteniamo

$$I(x + \delta x, y + \delta y, t + \delta t) = I(x, y, t) + \frac{\delta I}{\delta x} \delta x + \frac{\delta I}{\delta y} \delta y + \frac{\delta I}{\delta t} \delta t + \text{t.o.s.} \tag{24}$$

ove t.o.s. sono termini di ordine superiore trascurabili; ponendo $\frac{\delta I}{\delta x} \delta x + \frac{\delta I}{\delta y} \delta y + \frac{\delta I}{\delta t} \delta t = 0$ si ottiene

$$\frac{\delta I}{\delta x} v_x + \frac{\delta I}{\delta y} v_y + \frac{\delta I}{\delta t} v_t = 0 \tag{25}$$

Per risolvere l'Optical Flow sono necessarie altre equazioni, derivanti da vincoli addizionali, molto spesso assumendo un vincolo di

smoothing $|\Delta\mathcal{F}| \approx 0$, il che equivale ad assumere che il campo di moto non vari drasticamente fra regioni "vicine". Conseguentemente il grafo \mathcal{G}_t deformerà in \mathcal{G}_{t+1} e il suo aggiornamento (*morphing*) corrisponderà al optimal flow ristretto ai nodi N_t . Il calcolo di tale deformazione è ottenuto come soluzione al seguente problema:

Problema 1. Per una specifica locazione $C(s)$, dato s , si considera una porzione circolare $N_t(C_t)$ centrata in $C_t(s)$ di raggio I_t . $C_t(s)$ sia un punto di \mathbb{R}^2 . Prendiamo ora una porzione circolare analoga in I_{t+1} ma centrata in $C_t(s) + [\hat{v}_x, \hat{v}_y]$ e la chiamiamo $N_t(C_t + [\hat{v}_x, \hat{v}_y])$. Per calcolare l'optical flow $\mathcal{F}(C(s)) = \mathcal{F}(x, y) = (v_x, v_y)$ si minimizza l'integrale della distanza fra i due settori su un generico dominio A bidimensionale:

$$\mathcal{F}(C(s)) = \arg \min_{\hat{v}_x, \hat{v}_y} \underbrace{\int_A |N_t(C_t) - N_{t+1}(C_t + [\hat{v}_x, \hat{v}_y])| da}_{J_s} \quad (26)$$

J-map è un metodo per computare l'algoritmo di Optical Flow e altro non è che una rappresentazione di questo funzionale da minimizzare $J_s(\hat{v}_x, \hat{v}_y)$ relativo alla posizione $C(s)$. L'equazione 26 nella pratica porta ad un risultato sub-ottimo e ciò è dovuto per esempio alla compresenza di molti minimi nella stessa magnitudine (individuati quindi nello stesso settore). Per ottenere il vero minimo si utilizza un set di J-Maps "vicine", al posto di considerarne una sola, in modo da poterle correggere mutuamente. Infatti, considerando di applicare uno *smoothing* all'optical flow, e cioè

$$\mathcal{F}(x, y) \text{ è tale che } |\Delta\mathcal{F}| \approx 0$$

di conseguenza le J-Map "vicine" avranno valori di minimo simili. Preso un set di punti $\{C(s_1), \dots, C(s_n)\}$ dopo aver applicato lo *smoothing* si otterrà

$$\mathcal{F}(C(s_1)) \approx \dots \approx \mathcal{F}(C(s_n)) \text{ e quindi } J_{s_1} \approx \dots \approx J_{s_n}$$

Il processo di correzione è iterativo. Per ogni immagine (J-Map corrispondente al punto $C(s)$) J_s si considera un set di immagini $\{J_{s_1}, \dots, J_{s_n}\}$ "vicine" (J-Maps corrispondenti al set di punti $\{C(s_1), \dots, C(s_n)\}$) e si calcola il minimo globale su ognuna di esse. Indichiamo con $\{(v_{x_1}, v_{y_1}), \dots, (v_{x_n}, v_{y_n})\}$ il cui valore medio e matrice di covarianza indicano la direzione media dell'optical flow in $C(s)$ e forniscono una confidenza di tale misura.

$$v_x = \sum_{i=1}^n v_{x_i} \quad v_y = \sum_{i=1}^n v_{y_i} \quad (27)$$

$$\Sigma_{xy} = \begin{bmatrix} \sigma_{xx} & \sigma_{xy} \\ \sigma_{yx} & \sigma_{yy} \end{bmatrix} \quad (28)$$

Si è scelto di correggere J_s con

$$J_s = J_s \cdot (1 - \mathcal{N}([v_x, v_y], \Sigma_{xy})) \quad (29)$$

dove con \mathcal{N} si intende *distribuzione Gaussiana*. Dopo il passo di correzione la J_s corretta viene sovrascritta sulla "vecchia" J_s .

3.1.2 Algoritmo di *tracking* della struttura reticolare

L'implementazione¹ pratica di questo algoritmo è stata testata su una sequenza video di 50 frames di epitelio di *Drosophila*. Nel video si osserva un reticolo di circa 400 cellule che si dividono, si uniscono, si deformano e migrano; in cui alcune cellule escono dal frame mentre altre nuove entrano. Per inizializzare l'algoritmo si opera uno *static detection* al primo frame (come si è visto al capitolo precedente) utilizzando l'algoritmo Random Walk Agents il quale fornisce una rappresentazione grafica del reticolo. Utilizzando le J-Maps si effettua una trasformazione del grafo secondo l'equazione 21 ottenendo una rappresentazione del reticolo per il frame successivo. Successivamente si converte il grafo in una struttura di tipo Active-Contour (descritta al capitolo precedente). Infine si genera un set di random walk agent per ciascun nodo del grafo in modo da riempire possibili buchi e da espandere, ove possibile, il reticolo alla porzione di frame visibile. Questo processo permette di tracciare, frame per frame, il reticolo cellulare di *Drosophila*. Una volta che la struttura viene tracciata, a seguito del processo di *dinamic shape detection*, è possibile effettuare un'analisi sul comportamento dinamico del reticolo. Il dettaglio di come viene condotto questo tipo di analisi viene rimandato al capitolo seguente. Nella sezione che segue si illustrerà la tecnica per l'analisi della dinamica di una singola cellula in modo da poterla poi estendere a quella di tutto il reticolo.

3.2 shape deformation analysis

L'algoritmo appena descritto è stato il punto di partenza del lavoro di *analisi* che si è fatto in questo lavoro. In questa sezione, come si intuisce dal titolo, si tratterà del punto focale del lavoro, ossia *analisi di forma e movimento*. Per cominciare si darà una spiegazione formale e teorica su ciò che si intende per *shape analysis*. Poi si cercherà di caratterizzare in maniera più dettagliata ciò che si intende per *moto* e ciò che invece si intende per *deformazione* in modo da dare un senso alla frase: "*analisi del moto di un oggetto deformabile*".

¹ scritta in linguaggio MATLAB.

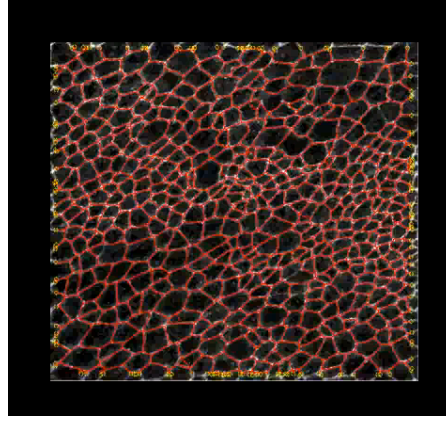


Figura 10.: Un frame risultato del tracking sul reticolo di cellule epiteliali di Drosophila.

Si può definire *shape analysis* come il l'insieme dei processi per estrarre metriche e descrittori a partire da un modello di rappresentazione di una singola forma o un insieme di esse. In altri termini, dato un set di forme $\{\Phi_0, \dots, \Phi_t, \dots, \Phi_N\}$ e il set delle loro rappresentazioni $\{\mathcal{R}(\Phi_0), \dots, \mathcal{R}(\Phi_t), \dots, \mathcal{R}(\Phi_N)\}$. L'obiettivo del processo di analisi, come anticipato nel cap 2, è quello di estrarre "numeri" o "figure" rilevanti; generalizzando l'equazione 5 si ha

$$\{\mathcal{R}(\Phi_0), \dots, \mathcal{R}(\Phi_t), \dots, \mathcal{R}(\Phi_N)\} \xrightarrow{\text{analysis}} \mathbb{R}^k \quad (30)$$

Nell'equazione 5 si considerava semplicemente $N = 1$. In questo caso è esplicitato il fatto che si può effettuare l'analisi della singola forma usando le informazioni relative al singolo istante temporale t oppure dell'intero set di forme sfruttando tutte le informazioni temporali.

data una forma piana $A \subset \mathbb{R}^2$ un modo semplice per estrarre la misura $f(A) \in \mathbb{R}^n$ è fare riferimento ai filtri lineari del tipo

$$f(A) = \int \psi(A) da$$

Ad esempio un set di funzioni polinomiali del tipo $\psi(x, y) = x^p y^q$ generano i momenti $m_{p,q}(C)$ di ordine $(p + q)$; dal momento che un'analisi effettuata solo mediante i momenti non da una percezione della deformazione principale, essi non sono metriche adeguate ad esplicitare una descrizione sulle misure effettive.

$$f_i(S) = \int_S \psi_i(x, y) dx dy \quad (31)$$

$$m_{p,q}(C) = \oint_C x^p(s) y^q(s) \delta(s) ds \quad (32)$$

dove δ è la densità lineare del contorno e, conoscendo il centro di massa della curva (x_C, y_C) , si possono scrivere i momenti centrali μ che sono *invarianti per traslazione*

$$\mu_{p,q}(C) = \oint_C (x(s) - x_C)^p (y(s) - y_C)^q \delta(s) ds \quad (33)$$

L'analisi della forma ha come obiettivo la qualificazione di due concetti legati: la deformazione e il moto. In particolare si cerca di dare una vera e propria quantificazione, mediante l'esplicitazione di determinate metriche di interesse, del fenomeno di "deformation" osservato per cercare di studiarne i comportamenti statistici e riconoscere eventuali correlazioni fra parametri.

3.2.1 Moto e deformazione

Cosa significa "moto di un oggetto deformabile"? Dal punto di vista formale (algebrico) è possibile separare il moto intendendolo come una funzione di gruppo (finito-dimensionale) dalla deformazione intesa come diffeomorfismo.

Quando si osserva un generico oggetto deformabile in moto si cerca da un lato di conservare una *nozione globale di moto* e dall'altro di descrivere "l'allontanamento dalla rigidità" mediante una *deformazione*. Un esempio eclatante di ciò che si cerca di descrivere è il moto di una medusa, infatti essa si sposta nello spazio deformandosi.

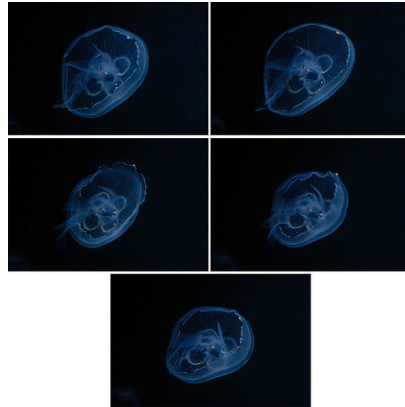


Figura 11.: Una medusa che si deforma muovendosi nello spazio.

Intuitivamente se si osserva invece un oggetto rigido che, ad esempio, cade, è possibile definirne il *moto* semplicemente fornendo le coordinate di una sua particella e l'orientamento di un sistema di riferimento ortogonale fisso con essa. Se invece si osserva un oggetto non rigido in movimento per descriverlo in ogni istante di tempo bisognerebbe specificare la traiettoria seguita da ciascuna sua particella. Cioè, detta γ_0 l'iniziale collezione di particelle, bisogna esprimere la funzione f che descrive come l'intero set di particelle evolve nel tempo:

$$\gamma_t = f(\gamma_0, t)$$

Se poi si ipotizza che ogni particella possa muoversi "indipendentemente" potrebbe non avere più senso cercare di quantificare un con-

cetto globale di moto quanto piuttosto sembrerebbe più appropriato la descrizione di f come una *deformazione* dell'oggetto di partenza.

Il *moto rigido* è un caso specifico ed è descritto globalmente mediante una matrice di rotazione $\mathbf{R}(t) \in SO(3)^2$ e da un vettore di traslazione $\mathbf{T}(t) \in \mathbb{R}^3$ tali che

$$\gamma_t = f(\gamma_0, t) = \mathbf{R}(t)\gamma_0 + \mathbf{T}(t)$$

Matematicamente possiamo esprimere il moto di un oggetto deformabile come composizione di un *moto rigido* $(\mathbf{R}(t), \mathbf{T}(t))$ e di una funzione globale di *deformazione* $h(\cdot, t)$ tale che

$$\gamma_t = h(\mathbf{R}(t)\gamma_0 + \mathbf{T}(t), t)$$

Si osservi che non è sempre possibile o sensato separare il concetto di moto globale dalla generica deformazione, dipende chiaramente dall'oggetto della propria analisi. Nel caso specifico studiato in questo progetto è utile cercare di definire un concetto di moto "globale", intendendolo come moto dell'intero reticolo, in relazione alla "mutua" deformazione fra le cellule. I concetti descritti sopra vengono in qualche modo estesi e complicati dal fatto di lavorare con un oggetto dalla struttura reticolare. Prima di porsi il problema più complesso dell'analisi del moto dell'intera struttura reticolare (che si rimanda al capitolo 4), ci concentriamo sul moto e la deformazione di una singola cellula. In questo modo cerchiamo di formalizzare la scelta e lo studio di alcune metriche di deformazione che possano risultare interessanti per lo studio della forma che assume mediamente una cellula epiteliale di *Drosophila*.

Analisi di movimento

Prima di occuparsi di analizzare come e in che direzione avvengono cambiamenti nella forma di una cellula, si è raccolta tutta l'informazione saliente riguardo al suo moto, ossia: traslazione e rotazione rispetto alla posizione (frame) iniziale. Come abbiamo detto è fondamentale separare ciò che caratterizza il moto da ciò che più precisamente si intende come perturbazione di forma.

Per ogni frame estratto dal video, conoscendo gli n punti

$\mathbf{P} = \begin{bmatrix} \mathbf{X} \\ \mathbf{Y} \end{bmatrix}$ che definiscono il contorno, calcoliamo la posizione del centroide (x_c, y_c) di ogni cellula come:

$$x_c = \frac{\sum_{i=1}^n \mathbf{X}(i) \cdot dl_i}{\sum_{i=1}^n dl_i}$$

$$y_c = \frac{\sum_{i=1}^n \mathbf{Y}(i) \cdot dl_i}{\sum_{i=1}^n dl_i}$$

² $SO(3)$ è il gruppo delle matrici ortogonali di dimensione 3.

dove

$$dl_i = \|\mathbf{P}(j) - \mathbf{P}(i)\| = \left\| \begin{bmatrix} X(j) \\ Y(j) \end{bmatrix} - \begin{bmatrix} X(i) \\ Y(i) \end{bmatrix} \right\|$$

con $j = \text{mod}(i - 2, n) + 1$

A questo punto è immediato associare lo spostamento nello spazio cartesiano della cellula allo spostamento del suo centroide. Nella figura seguente si può osservare lo spostamento di una cellula in 20 frames consecutivi; come si può notare la cellula mentre trasla può ruotare e deformarsi

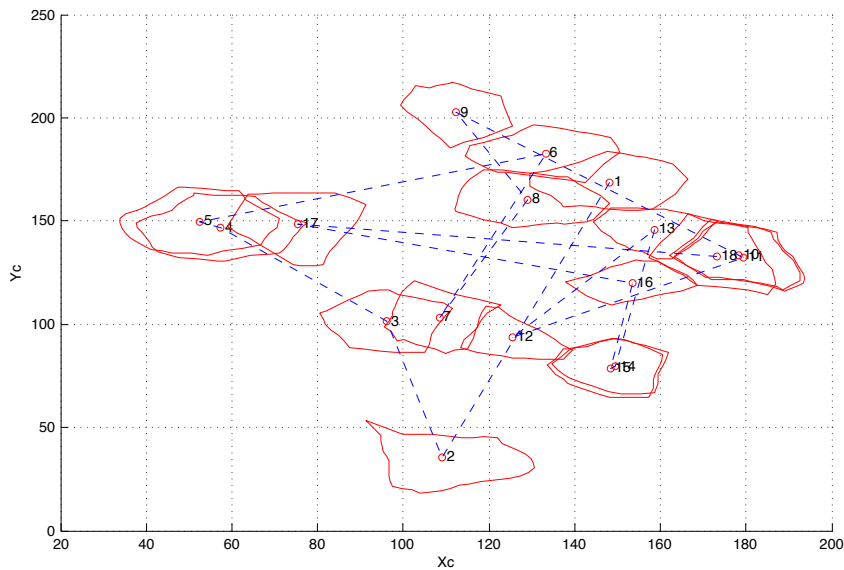


Figura 12.: Traslazione di una specifica cellula (nel codice identificata dalla variabile cellID) in 20 frames consecutivi. I numeri all'interno della cellula rappresentano la variabile temporale.

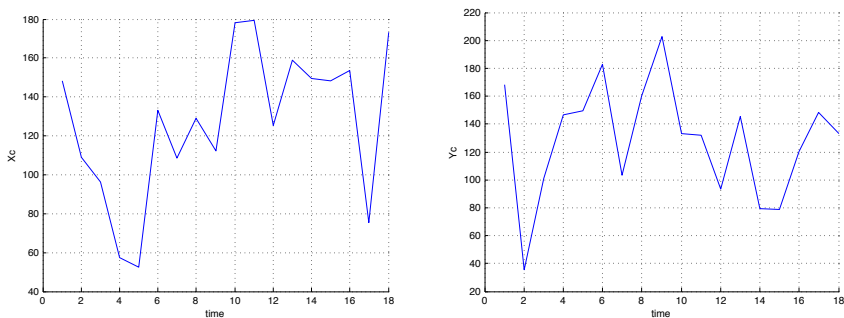


Figura 13.: Traslazione di ascissa e ordinata di una cellula in 20 frames consecutivi.

Per calcolare la rotazione è necessario riferirsi ad un asse rispetto al quale calcolarla. A tal proposito si sfrutta la ben nota tecnica di Ana-

lisi delle Componenti Principali³ (PCA *Principal Component Analysis*). PCA consiste in una trasformazione lineare dalle variabili originali \mathbf{X} ad altre \mathbf{Z} che esprimono la stessa informazione ma sono fra loro incorrelate (Componenti Principali). La trasformazione cercata è la similitudine \mathbf{W} fra la matrice di correlazione e la matrice diagonale degli autovalori, tale che

$$\mathbf{L} = \text{diag}(\sigma_1^2, \sigma_2^2, \dots, \sigma_p^2) = \mathbf{W}^T \cdot \mathbf{C} \cdot \mathbf{W}$$

$$\mathbf{Z} = \mathbf{X} \cdot \mathbf{W}$$

In poche parole la PCA può essere vista come la ricerca delle "direzioni privilegiate" che massimizzano la variazioni dei dati ed eliminano le correlazioni:

$$\mathbf{Z}^T = \mathbf{W}^T \cdot \mathbf{X}^T$$

$$\begin{bmatrix} z_1 \\ z_2 \end{bmatrix} = \begin{bmatrix} w_{11} & w_{21} \\ w_{12} & w_{22} \end{bmatrix} \cdot \begin{bmatrix} x_1 \\ x_2 \end{bmatrix}$$

ove

$$z_1 = x_1 \cdot w_{11} + x_2 \cdot w_{21}$$

rappresenta la direzione di massima varianza di \mathbf{X} mentre

$$z_2 = x_1 \cdot w_{12} + x_2 \cdot w_{22}$$

rappresenta la direzione di massima varianza di \mathbf{X} , escluso x_1 . La matrice \mathbf{W} formata dagli autovettori ordinati per autovalori decrescenti indicano le direzioni di massima varianza. La matrice \mathbf{L} (diagonale) degli riporta i valori delle varianze nel nuovo riferimento PCA

In MATLAB PCA è implementata mediante la funzione `[COEFF, SCORE, latent]=princomp(P)`, dove `COEFF (W)` è la matrice dei coefficienti (pesi) delle componenti principali di \mathbf{P} ; ogni riga di \mathbf{P} corrisponde alle osservazioni (in tal caso i punti del contorno della cellula) ed ogni colonna alle variabili (\mathbf{X} e \mathbf{Y} di tutti i punti del contorno). `COEFF` è una matrice $p \times p$ (dove $p = \#$ osservazioni) le cui colonne contengono i coefficienti per una PC in ordine decrescente di varianza. `SCORE` rappresenta \mathbf{P} nello spazio PC; le sue righe corrispondono alle osservazioni mentre le colonne alle componenti. `latent (L)` infine è il vettore contenente gli autovalori della matrice di covarianza di \mathbf{P} . Di conseguenza si possono ottenere gli assi principali (maggiore e minore) della cellula come:

```
paxis = COEFF(:,1)*sqrt(latent(1)*2);
maxis = COEFF(:,2)*sqrt(latent(2)*2);
```

³ si veda appendice A2.

dove l'asse principale *paxis* deriva infatti dalla componente principale a varianza maggiore invertendo la relazione che le lega: $\text{latent}(1)$ è il rappresentante dell'emipioiezione dell'asse principale cioè $\alpha^2/2$, da cui si evince l'espressione per l'asse principale.

L'angolo che si considera è quello fra l'asse principale e la direzione verticale (o equivalentemente quello fra l'asse minore e la direzione orizzontale); esso è fissato 0° nel verso verticale e cresce in senso antiorario.

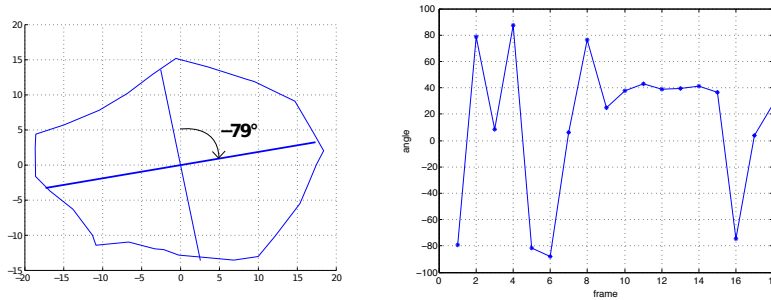


Figura 14.: A sinistra è raffigurata una cellula con in evidenza i suoi assi principali. L'errore che si ha sulla lunghezza degli assi corrisponde ad un errore di misura dovuto al campionamento finito ($N = 1000$ punti descrivono una cellula). A destra si ha l'andamento dell'angolo (angle) nel tempo.

Metriche di deformazione

In questa sezione si entrerà più nel dettaglio e ci si occuperà di tradurre le informazioni "di forma" in parametri utili a descrivere il comportamento dinamico delle cellule. Una delle principali attività nell'analisi di immagini è la discriminazione di oggetti in base alla loro apparenza. Possono essere misurati vari aspetti riguardo a come un oggetto appare: la struttura, il colore, la forma... La forma è lo strumento più potente forse per descrivere e discriminare diversi oggetti ed è perciò che viene applicata in moltissime aree della visione computazionale e non solo. Nonostante siano state introdotte molte nuove tecniche di modellizzazione matematica, il problema dell'analisi della forma rimane un argomento che non lascia spazio ad un'unica interpretazione. Una delle difficoltà risiede proprio nel capire quali siano delle buone metriche per condurre un'analisi sulla forma. La scelta dipende ovviamente dal caso specifico che si vuole analizzare; infatti è possibile che non tutte le variazioni di forma siano significative per la data situazione, o ad esempio è bene non dare molto peso a piccole variazioni di forma causate dal rumore.

L'approccio che si è scelto di seguire in questa sezione è considerare la deformazione di una singola cellula rispetto ad una forma di riferimento: un'ellisse avente stessi assi della cellula. L'ellisse di riferimento viene costruita al primo frame. Nei frames successivi si

valutano quattro parametri di deformazione rispetto all'ellisse di riferimento: *boundary activity*⁴ v_C , differenza di area percentuale (Δ_A), differenza media di area⁵ (δ), distanza di Hausdorff (d_H) ed elongazione (κ).

Il parametro attività del contorno è definito mediante la seguente espressione

$$v_C(t) = \left(\Lambda_C(t) - \frac{1}{n} \sum_{j=t-n}^{t-1} \Lambda_C(j) \right)^2$$

dove $\Lambda_C(t)$ si riferisce al perimetro della cellula C (al generico istante t , fissato). Questo parametro viene calcolato sfruttando l'informazione dell'intervallo temporale degli n istanti che precedono t .

La differenza di area percentuale è definita come

$$\Delta_A = \frac{|A_E - A_C|}{A_C \cdot 100}$$

con A_E area dell'ellisse e A_C area della cellula, calcolate con l'approssimazione trapezoidale a partire dai punti del contorno in coordinate polari; essa viene normalizzata rispetto all'area della cellula.

La differenza media di area invece è definita come

$$\delta = \frac{1}{N} \sum_{i=1}^N \frac{|\rho - \bar{\rho}|^2}{N \cdot \bar{\rho}^2}$$

con ρ e $\bar{\rho}$ raggi medi di cellula ed ellisse, cioè distanza tra centroide e contorno; e viene normalizzata rispetto ad $N \cdot \bar{\rho}$ che è in un certo senso pensabile come "area media" della cellula. Come si è detto nella nota in questa metrica si va a stimare la differenza di "area sottesa" fra le curve in coordinate polari di ellisse di riferimento e cellula. Non è un parametro che stima l'area vera e propria. Questa metrica riprende il concetto di *radius vector function* definito in [Kindratenko, 2003](#). Richiamiamo brevemente la definizione di *radius vector function* definita nell'articolo di Kindratenko:

Definizione 3 (Radius Vector Function). Si scelga

- un punto di riferimento O interno alle figura X, solitamente il centro di gravità;

⁴ definita in [Silletti, 2007-2009](#).

⁵ con "area" in questo caso si intende "area sottesa alla curva che rappresenta ellisse e cellula in coordinate polari".

- un asse di riferimento l che intersechi il punto O , solitamente viene scelto parallelo all'asse x o y ;

allora la *radius vector function* $r_X(\phi)$ è la distanza del punto di riferimento O dal contorno della forma X nella direzione dell'angolo ϕ , con $0 \leq \phi \leq 2\pi$.

Una condizione necessaria affinché la funzione *radius vector* definisca completamente la figura (cioè tale che, dato $r_X(\phi)$, allora la figura può essere ricostruita univocamente) è che la figura X sia "a forma di stella" rispetto a O . Vale a dire che ogni segmento da O al contorno p sia interamente contenuto nella forma X .

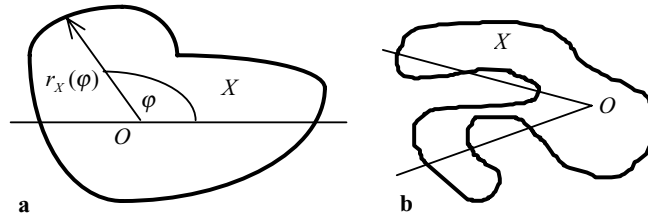


Figura 15.: (a) Definizione di *radius vector function*; (b) definizione mal posta nel caso di forme non convesse.

Nel nostro caso il punto O corrisponde al centroide e la linea di riferimento l passante per il punto O corrisponde all'asse verticale della cellula; l'angolo θ che abbiamo considerato è definito esattamente come ϕ . Infine la definizione è consistente in quanto la cellula ha sempre una forma *star-shaped* rispetto al centroide.

La distanza di Hausdorff è definita dalle seguenti relazioni

$$\begin{aligned} d_H(C, E) &= \max\{d_{H^+}(C, E), d_{H^-}(C, E)\} \quad \text{con} \\ d_{H^+}(C, E) &= \max\{d_{\text{Euclid}}(c, E) : c \in C\} \\ d_{H^-}(C, E) &= \max\{d_{\text{Euclid}}(e, C) : e \in E\} \end{aligned}$$

e viene normalizzata rispetto alla numerosità degli insiemi C ed E dei punti del contorno di cellula ed ellisse.

$$d_H = \frac{1}{N} \cdot d_H(C, E)$$

Infine il parametro di elongazione è definito come rapporto fra asse maggiore ed asse minore

$$\kappa = \frac{\text{asse}_p}{\text{asse}_m}$$

Queste metriche vengono calcolate per una cellula in ogni frame e danno come risultati degli andamenti riportati in figura 17

Ad esclusione della distanza di Hausdorff, confrontando gli andamenti delle altre tre metriche si possono osservare dei picchi negli stessi istanti (ad esempio si veda l'istante 12) il che è prevedibile visto che le metriche corrispondono a descrizioni molto simili della deformazione che la cellula sta maturando.

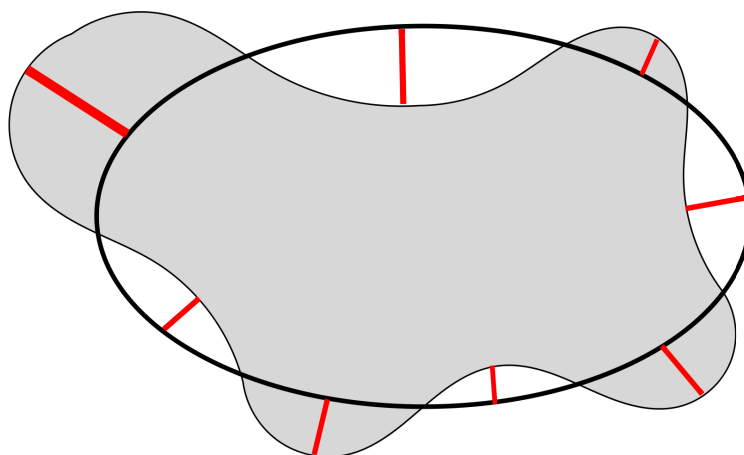


Figura 16.: Distanza di Hausdorff fra cellula ed ellisse.

La distanza di Hausdorff è la metrica, tra quelle scelte, che descrive la deformazione in maniera più "puntuale". In effetti ciò che essa valuta è la massima distanza fra i punti delle due curve. A parità di area le curve potrebbero essere molto diverse, ad esempio la cellula potrebbe avere degli spigoli molto marcati. Nonostante la sua diversità e la difficoltà nel paragonarla ad altre metriche essa da un'informazione importante dal punto di vista della "forma" vera e propria della cellula. Quantifica in maniera dettagliata quanto il contorno della cellula sia irregolare rispetto ad una geometria ellittica presa come modello di riferimento.

Ogni cellula durante tutto il video si deforma e trasla nello spazio immagine. Una metrica che mostra in maniera esplicita la variazione che avviene nel contorno della cellula è il raggio, ossia la distanza fra il suo centroide e i punti del bordo. In figura 20 viene rappresentato il raggio della cellula in funzione dell'angolo (fra l'asse principale e la direzione verticale) e del tempo.

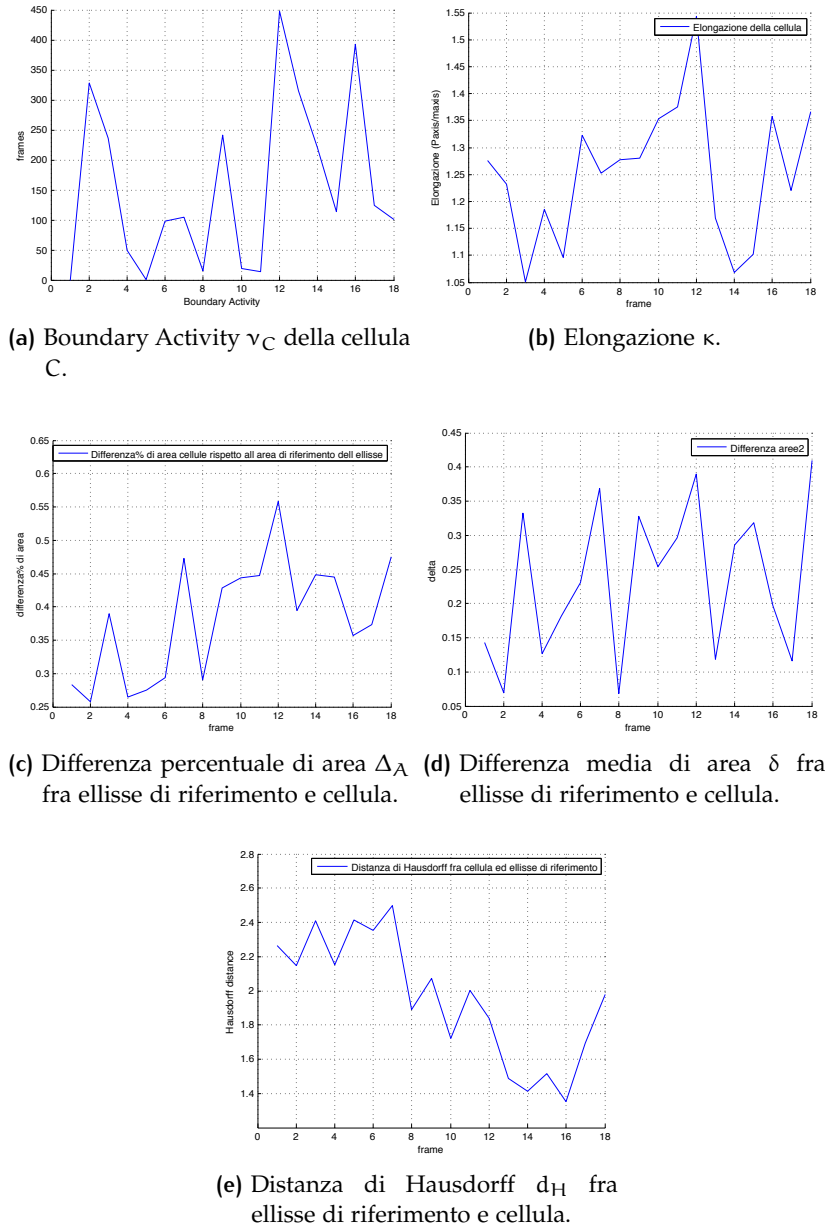
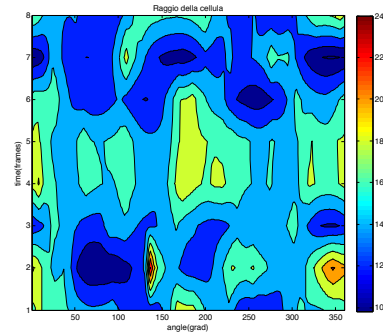
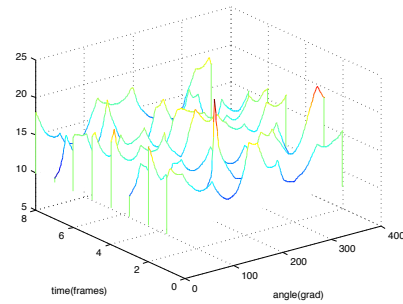


Figura 17.: Andamento delle metriche di deformazione in 20 frames consecutivi per una singola cellula.

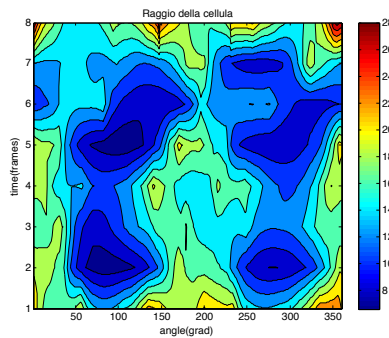
Queste immagini rappresentano la variazione del raggio della cellula durante un intervallo temporale di 10 frames. In particolare quello che la figura 20(c) mostra è che la cellula 43 in linea di massima, per



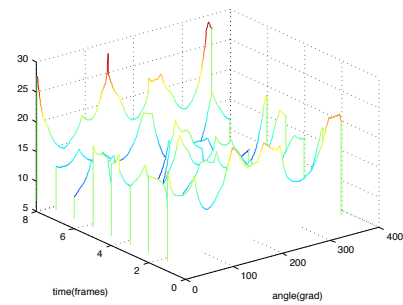
(a) Raggio della cellula 10.



(b) Altra rappresentazione del raggio della cellula 10.



(c) Raggio della cellula 43.



(d) Altra rappresentazione del raggio della cellula 43.

Figura 18.: Raggio delle cellule CellID = 10 e CellID = 43 in funzione del tempo e dell'angolo.

tutti i frames, mantiene come direzione principale quella che presenta un angolo fra i 150 e 200 gradi rispetto alla verticale. Un picco del raggio infatti corrisponde alla direzione di "elongazione" della cellula. La figura 20(a) invece può far intuire che inizialmente (al frame 2) la cellula 10 presenti un picco di raggio, e quindi la direzione principale, ad un angolo di 140° ; dopo di che il picco si abbassa e si "stabilizza" attorno ai 200° fino al frame 6, per poi rimpicciolirsi negli ultimi frame. Valutazioni di questo tipo sono utili per cercare di riconoscere le direzioni in cui le cellule si allungano e ciò interessa poiché la condizione di grande allungamento di una cellula corrisponde ad una condizione di NON equilibrio. Di solito questa situazione di allungamento precede uno *split* ossia una divisione della cellula originaria in due nuove cellule.

Si osservi ad esempio la posizione "finale" (cioè all'ultimo frame) del reticolo cellulare in figura 19 e si noti a tal proposito come le cellule adiacenti 30 e 40 siano pressapoco disposte lungo la stessa direzione principale (circa 190°).

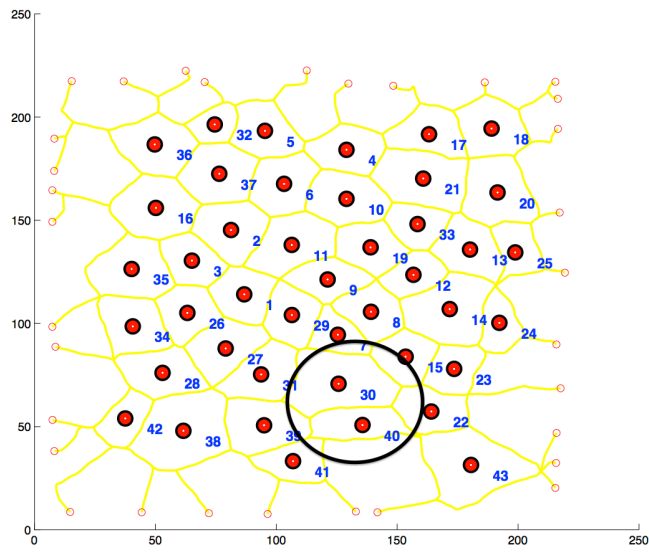
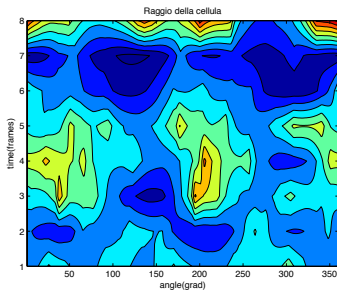
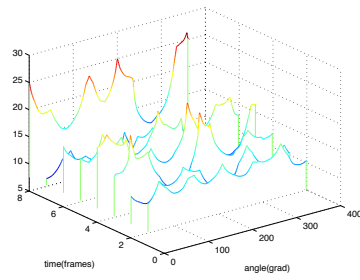


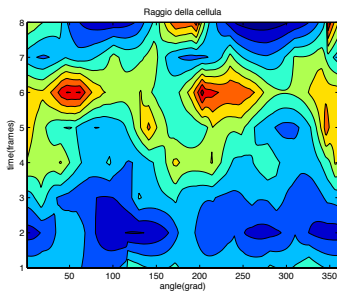
Figura 19.



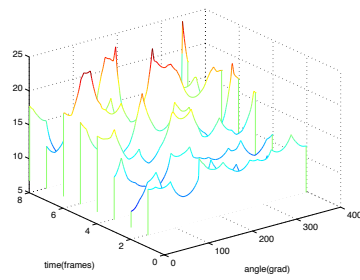
(a) Raggio della cellula 30.



(b) Altra rappresentazione del raggio della cellula 30.



(c) Raggio della cellula 40.



(d) Altra rappresentazione del raggio della cellula 40.

Figura 20.: Raggio delle cellule CellUID = 30 e CellUID = 40 in funzione del tempo e dell'angolo.

4

ANALISI DI STRUTTURE RETICOLARI

Indice

4.1	<i>Shape detection</i> in strutture reticolari	39
4.1.1	<i>The Random Walk Agents Algorithm</i>	40
4.2	Implementazione: video di <i>Drosophila wing</i>	43
4.3	<i>Tracking</i> della struttura reticolare	45
4.3.1	Analisi di deformazione di gruppi di cellule limitrofe	45

4.1 *shape detection* IN STRUTTURE RETICOLARI

Come si può facilmente immaginare esistono in natura e nella realtà umana un gran numero di esempi di strutture reticolari.



Figura 21.: Esempi di strutture reticolari: cellule epiteliali di *Drosophila Melanogaster*, un alveare, Venezia dall'alto.

Per questo tipo di strutture, dato l'elevato numero di gradi di libertà, approcci comuni alla minimizzazione come quelli presentati nella sezione precedente non portano a risultati soddisfacenti. Lavorare su questo tipo di immagini in primo luogo presenta una complessità di calcolo maggiore rispetto al caso semplice descritto alla sezione 2.3, e in secondo luogo non è ben chiaro come sia da progettare il termine energetico E . Si proporrà in seguito l'approccio (presentato in [Silletti, 2007-2009](#)) di tipo *random walk* che modella, emulandola, la capacità di visione umana nel catturare immagini reticolari. Si immagina far processare l'immagine I a delle entità intelligenti, che chiameremo *random walk agents*, ciascuna delle quali dovrà trovare un cammino all'interno della tessitura digitale.

L'immagine tipo su cui lavora l'algoritmo ha uno sfondo nero in cui si dirama un reticolo bianco. L'immagine digitale è quindi considerata come uno sfondo da cui estrarre i percorsi chiari, individuati grazie al contrasto con i pixel scuri. Ogni frame è infatti visto come una zona che deve essere esplorata e i percorsi definiti da sequenze di pixel chiari, appartenenti al perimetro della cellula, sono considerati come strade percorribili, mentre i pixel scuri, la parte interna alle cellule, devono rimanere aree inesplorate. L'agente avanza attraverso la scansione delle zone limitrofe alla posizione in cui si trova alla ricerca di percorsi esplorabili. Quando si verifica una biforcazione (più direzioni di avanzamento possibili), l'agente genera uno o più "fratelli" che iniziano a muoversi indipendentemente finché l'intera rete è stata esplorata. I cammini vengono creati minimizzando localmente un termine di energia \hat{E} il quale funge da approssimazione analitica di E . In pratica questo approccio sposta il problema della determinazione di E al sottoproblema più semplice di progettare il termine \hat{E} .

4.1.1 The Random Walk Agents Algorithm

Un *random walk agent* è un'entità \mathcal{A} che si muove sull'immagine digitale tracciando il percorso del suo passaggio. Si può immaginare un agente come un sistema a tempo discreto che descriva il modello di cammino e abbia come ingresso un termine di energia esterna \hat{E} . Formalmente si ipotizza che ogni agente \mathcal{A}_i caratterizzato all'istante t dalla coppia $\{p_i, \theta_i\}(t)$ della posizione nell'immagine digitale e della direzione di avanzamento implichi una legge di moto del tipo:

$$\mathbf{p}_i(t+1) = \mathbf{p}_i(t) + ke^{j\theta_i(t)} \quad (34)$$

dove la costante k è un valore scalare che rappresenta la taglia del passo di avanzamento, j è l'unità immaginaria e la direzione di avanzamento $\theta_i(t) \in \Theta = [0, 2\pi)$ è l'angolo polare visto dall'agente. In generale quindi il sistema associato al random walk agent sarà del tipo

$$\mathbf{p}(t+1) = \mathbf{p}(t) + k \cdot g(\hat{E}(t)) \quad (35)$$

dove

$$\mathbf{p}(t) = \begin{bmatrix} x(t) \\ y(t) \end{bmatrix} \in \mathbb{L} \quad (36)$$

è la posizione corrente nel dominio $\mathbb{L} \subset \mathbb{R}^2$ dell'immagine e $g(\cdot)$ è la funzione di direzione che dipende da una certa funzione di energia \hat{E} , dalla posizione, dal verso e dalla velocità del moto.

L'equazione di osservazione è la seguente

$$\mathbf{y}(t) = \begin{bmatrix} x(t) \\ y(t) \\ \theta(t) \end{bmatrix} = \begin{bmatrix} 1 & 0 \\ 0 & 1 \\ 0 & 0 \end{bmatrix} \mathbf{p}(t) + \begin{bmatrix} 0 \\ 0 \\ 1 \end{bmatrix} \theta(t) \quad (37)$$

La variabile temporale t si riferisce allo pseudo-tempo relativo dell'agente, non è quindi una sequenza temporale che fa riferimento ad una collezione di reali immagini digitali consecutive. Sono attive contemporaneamente varie istanze di agente $\{\mathcal{A}_i, i = 1, \dots, N_{\mathcal{A}}\}$ ognuna delle quali ad ogni istante genera più di una direzione candidata per l'avanzamento dell'esplorazione. Più precisamente, per ogni coppia $\{x, y\}_i$ viene salvato nel vettore $\Theta_i \in \mathbb{R}^{m_i}$ il set delle m_i possibili direzioni $\{\theta_{i,j}, j = 1, \dots, m_i\}$ che generano $(m_i - 1)$ nuovi agenti.

Il ruolo giocato da \hat{E} nell'equazione 35 corrisponde a condurre l'agente \mathcal{A} attraverso le locazioni dell'immagine caratterizzate dalla presenza della struttura reticolare da tracciare. Consideriamo che l'agente \mathcal{A} sia nella locazione $[x(t), y(t)]$ del frame. Lo scopo di $g(\hat{E})$ è di fornire un set di direzioni per l'avanzamento $\Theta = \{\theta_1, \dots, \theta_{N_{\mathcal{A}}}\}$ tali che $[x(t+1), y(t+1)]$ sia ancora una locazione appartenente alla struttura reticolare. L'idea è quella di esplorare le zone limitrofe alla posizione corrente e intuitivamente muoversi da una buona locazione ad un'altra buona locazione.

$\hat{E} : [0, 2\pi) \rightarrow \mathbb{R}$ raccoglie informazioni locali. Per ogni settore (ad esempio circolare) Ω_i centrato nella posizione corrente \mathbf{p} , la struttura dell'immagine viene confrontata con una di riferimento:

$$\hat{E}(\theta_i) = \frac{\int_{\Omega_i} \sqrt{(I(\omega) - I_{ref}(\omega))^2} d\omega}{\int_{\Omega_i} d\omega} \quad (38)$$

Il set delle direzioni di avanzamento è fornito dall'espressione

$$\Theta = \{\theta_1, \dots, \theta_{N_{\mathcal{A}}}\} = \arg \frac{\partial \hat{E}}{\partial \theta} = 0 \text{ and } \hat{E} < \text{threshold} \quad (39)$$

da cui si può dedurre che il valore threshold elimina tutte le direzioni corrispondenti ad energie elevate e $\frac{\partial \hat{E}}{\partial \theta} = 0$ impone che il punto sia un minimo (locale della funzione di energia). Le direzioni di avanzamento sono ingressi per l'equazione 35 di conseguenza generano corrispondentemente un numero di agenti $\{\mathcal{A}_i, i = 1, \dots, N_{\mathcal{A}}\}$.

L'algoritmo si occupa anche di rimuovere piccoli loop che spesso si tracciano in zone di colore uniforme nell'immagine, poiché in queste il funzionale di energia non è incline a distinguere tra indicazioni ridondanti. In tali situazioni ogni direzione si equivale in termini di energia e quindi la scelta del passo successivo viene fatta in maniera random. In simboli ciò accade poiché $|\Theta| \gg 1$ e $\theta_i \approx \theta_{i+1} \forall \theta_i \in \Theta$. In altre parole tutti i valori $\theta_i \in \Theta$ sono simili e queste direzioni

si estende uniformemente sull'ambiente circostante e quindi non c'è nessuna particolare ragione per sceglierne una piuttosto che un'altra.

Per rimuovere i suddetti anelli piccoli, dopo l'estrazione di un punto dalla coda, viene testato sulla creazione di loop nel grafico. Loop di grandi dimensioni sono accettati in quanto corrispondono a reali tracciati chiusi (corrispondenti a cellule) mentre i cicli piccoli vengono ignorati. Una procedura di unione è inoltre necessaria per colmare le lacune di piccole dimensioni e chiudere grandi loop: se due agenti \mathcal{A}_i e \mathcal{A}_j finiscono vicini in modo tale che $|\{x, y\}_i - \{x, y\}_j| < \epsilon$, dove ϵ è una soglia, essi vengono uniti.

L'algoritmo segna anche come nodi di bordo quei nodi vicino al confine dell'immagine digitale: il processo di espansione finisce lì.

Gli agenti random walk generano una struttura che si espande sulla struttura reticolare originale, ma il processo di espansione potrebbe creare dei rami morti. Questi rami sono caratterizzati da nodi estremali non di bordo di grado 1. Si richiede una pulizia post-processo che prenda il grafico come input, e rimuova iterativamente tutti questi punti. La procedura è iterativa e si ferma quando non si verifica più alcuna rimozione. Una procedura di aggancio unisce quindi i seguenti tipi di punti:

- una coppia di nodi vicini di grado $d > 2$
- un nodo di grado $d > 2$ e il suo vicino nodo di bordo
- una coppia di nodi di bordo vicini

Anche questa procedura è iterativa e si conclude quando non ci sono più fusioni da fare.

Globalmente la descrizione che viene fornita dai random walkers, ottenuta dall'unione di tutte le osservazioni $\mathbf{y}_i(t)$ associate a ciascun agente \mathcal{A}_i , produce un modello di grafo del tipo $\mathcal{G} = (\mathcal{N}, \mathcal{E})$, dove con \mathcal{N} si denotano i nodi e con \mathcal{E} le connessioni (edges). Nel grafo ogni nodo $n_i \in \mathcal{N}$ rappresenta lo stato $\{x, y, \theta\}$ della posizione visitata da ciascuno degli agenti \mathcal{A}_i e ogni connessione tiene traccia del cammino percorso da ciascun agente. Questo grafo fornisce una buona rappresentazione della struttura recuperata. Per avere un modello che sia più compatto e adatto ad essere usato sia per simulare che per predire il comportamento della struttura, una procedura si occupa di astrarre le informazioni essenziali e generare un nuovo modello. Si considerano interessanti solo un sottoinsieme $\bar{\mathcal{N}} \subset \mathcal{N}$ (ed i corrispondenti edges $\bar{\mathcal{E}}$) associati alle biforcazioni. Il nuovo grafo $\bar{\mathcal{G}} = (\bar{\mathcal{N}}, \bar{\mathcal{E}})$ è ottenuto selezionando i nodi che abbiano grado maggiore di 2, cioè con $\delta(\bar{\mathcal{G}}) = 3$. Questa procedura restituisce una struttura dati compatta che permette di definire delle metriche per poi poter effettuare delle valutazioni quantitative sulla struttura di interesse, nonché di applicare metodi e modelli tipici della teoria dei grafi.

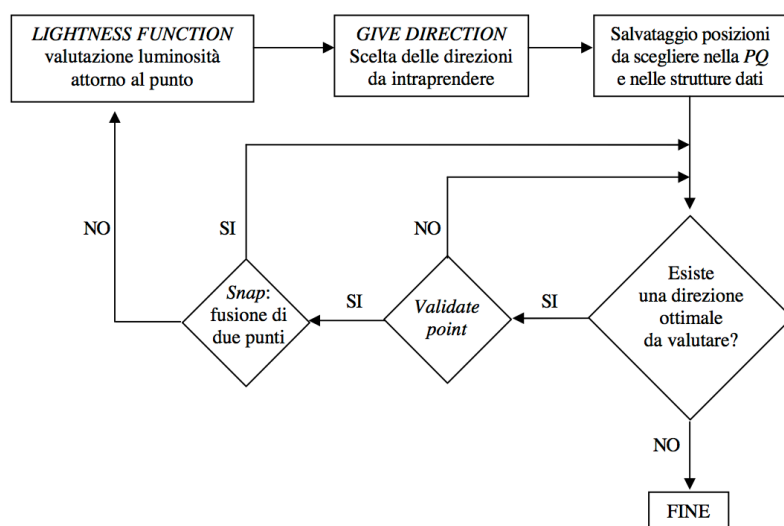


Figura 22.: Schema logico di implementazione del Random Walk Algorithm. *LIGHTNESS FUNCTION* scansiona le zone limitrofe mediante dei settori rettangolari centrati nella posizione corrente e calcola per ognuno la corrispettiva luminosità. *GIVE DIRECTION* valuta le direzioni ottimali di avanzamento (massima lightness). *Priority Queue* gestisce i nodi che devono ancora essere convalidati. Le *strutture dati* servono ad immagazzinare i risultati dei cammini già esplorati. *VALIDATE POINT* funzione che determina se il punto in esame è valido (o sul bordo). *SNAP* fonde punti vicini.

4.2 IMPLEMENTAZIONE: VIDEO DI DROSOPHILA WING

Recentemente sono stati studiati alcuni pattern evolutivi della *Drosophila* con l'obiettivo di acquisire maggiore conoscenza circa la morfogenesi di questo insetto. Al fine di raggiungere questo obiettivo i ricercatori biologi hanno espresso la necessità di poter comprendere meglio la dinamica della struttura cellulare dell'epitelio delle ali di *Drosophila*, in modo da favorire oltretutto la comprensione di un fenomeno noto come polarità cellulare planare¹. Un approccio meccanicistico per studiare questo fenomeno è stato in grado di fornire maggior materiale per comprendere il comportamento delle strutture intracellulari e delle proteine di regolamentazione nelle cellule epiteliali dell'ala mosca (Amonlirdviman *et al.*, 2005). Nel modello dinamico si assume che le cellule siano distribuite regolarmente, in una struttura a nido d'ape, su tutta la superficie alare. Per capire come estendere questi modelli, ha senso cercare di dedurre la struttura e il movimento della rete di cellule epiteliali, osservando il fenotipo dell'ala mosca con filmati ripresi in laboratorio.

¹ Lo stabilirsi e il mantenimento della polarità cellulare si basa su una distribuzione asimmetrica del citoscheletro e dell'attività di altre proteine all'interno della cellula.

Il lavoro proposto affronta, mediante l'uso dell'algoritmo presentato in Silletti, 2007-2009 e , il problema del recupero di una tale struttura utilizzando gli strumenti classici della teoria dei sistemi dinamici. L'obiettivo è quello di riuscire ad ottenere da una parte un rilevamento visivo accurato e dall'altra fornire una rappresentazione del fenomeno osservato in modo preciso e coerente con gli interessi biologici per il quale è stato progettato. Particolare cura è stata messa a tal proposito nella fase iniziale di progetto per ottenere un modello (rappresentazione) di questa struttura, le cui caratteristiche principali siano, come abbiamo anticipato, la compattezza e la semplicità.

Il primo passo consiste, a partire da un'immagine di epitelio di mosca, nel costruire una rete che rappresenti la struttura cellulare reticolare. Nella fase successiva invece, dato un filmato (sequenza di frames) dell'epitelio, correlare le reti generate per ogni singolo frame all'interno di un modello dinamico temporizzato (il tempo corrisponde alla sequenza di frames). Questo approccio ha come duplice obiettivo quello di estrarre la struttura reticolare in modo che sia una buona rappresentazione del reticolo epiteliale dell'ala di *Drosophila*, e che tale rete permetta di astrarre le informazioni relative alla struttura in un modello gestibile. Tutto ciò è ottenuto mediante l'analisi dei dati visivi attraverso l'algoritmo presentato alla sezione precedente.

Le immagini acquisite da esperimenti biologici, come quelle utilizzate, sono particolarmente rumorose e mostrano scarso contrasto, inoltre spesso lo sfondo ha un'illuminazione non uniforme. Conseguenza inevitabile di ciò sono i contorni poco marcati della struttura tali da non poter essere adeguatamente segmentate e riconosciute. Come spesso accade nelle applicazioni di Visual Computing è utile operare un pre-processamento ai dati visivi in modo da "pulirli" il più possibile e cercare di aumentare il contrasto fra ciò che si considera sfondo e la *feature* di interesse. Si effettuano sequenzialmente le seguenti operazioni sull'immagine:

- viene applicato un filtro gaussiano passa-basso per attenuare il rumore in alta frequenza;
- viene applicato un filtro di erosione che sopprime i pixel bianchi isolati e diminuisca l'intensità dei bordi delle cellule, mantenendo tutte le informazioni significative
- si effettua un *histogram stretch* in modo da recuperare parzialmente il range di colore dinamico iniziale;
- si aumenta l'intensità dell'immagine per ottenere maggior contrasto.

Le immagini di *Drosophila wing epithelium* che si sono analizzate, in particolare, presentano una struttura reticolare caratterizzata da tratti

luminosi sovrainposti ad uno sfondo scuro. La luminosità che si può osservare è dovuta alla presenza di proteine marker che aderiscono sul bordo della cellula e appaiono bianche quando vengono visualizzate per mezzo di un microscopio. Questa luminosità mette in risalto una struttura globale reticolare che appare "bianca" su fondo nero. Nel progettare il termine energetico bisogna tenere conto di questo e far sì che \hat{E} rifletta localmente le proprietà globali della struttura da rilevare. In particolare si deve scegliere I_{ref} a valore costante 1 nell'equazione 38. Porre $I_{ref} = 1$ corrisponde col prendere come riferimento il colore *bianco*, infatti 1 è il codice di colore bianco nelle immagini RGB:

$$\hat{E}(\theta_i) = \frac{\int_{\Omega_i} \sqrt{(I(\omega) - 1)^2} d\omega}{\int_{\Omega_i} d\omega} \quad (40)$$

4.3 *tracking* DELLA STRUTTURA RETICOLARE

La dinamica di rimodellamento dei tessuti durante la morfogenesi è il risultato della combinazione dei comportamenti delle singole cellule e dei loro riarrangiamenti collettivi. Una delle sfide più importanti per la biologia dello sviluppo è quella di capire come le informazioni molecolari dipendano dal movimento individuale e collettivo delle cellule che formano i tessuti sia attraverso gli stimoli fra le cellule dei tessuti che attraverso la risposta a stimoli esterni applicati. Grazie agli enormi progressi nella biologia molecolare, nella genetica, nelle tecniche di imaging e nell'inseguimento automatico di molte cellule in parallelo, diventa possibile tracciare l'evoluzione dei fenotipi morfo-genetici, durante la fase dello sviluppo, come funzione di perturbazioni molecolari e manipolazione fisiche. Questo pone le basi per l'identificazione e la quantificazione del cambiamento nella geometria delle forme in termini di tassi di deformazione di uno specifico ceppo.

4.3.1 Analisi di deformazione di gruppi di cellule limitrofe

Durante la morfogenesi movimenti di convergenza ed estensione di rimodellano i tessuti. Quando si analizza l'evoluzione di un reticolo cellulare, per quantificare la deformazione locale del tessuto, può essere utile concentrarsi su domini definiti o da una cella centrale circondata da una corona di cellule vicine o da un insieme di cellule che condividono un vertice. Si vuole cercare di capire se c'è una sorta di trend di rimodernamento d'insieme fra esse. Il software di monitoraggio individua le cellule e le collega in un processo iterativo mediante l'algoritmo dei Random Walkers descritto al paragrafo 4.1.1. In ogni istante di tempo il programma memorizza, per ogni

cellula, le coordinate del suo centroide e anche quelle dei punti che costituiscono il poligono che ne descrive il contorno. Si osservi a tal proposito la figura 23

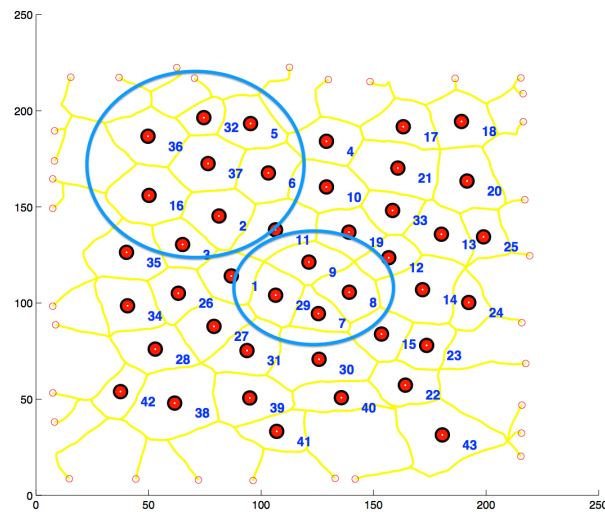


Figura 23.: Reticolo ricavato dall'epitelio di Drosophila. In evidenza due gruppi di cellule adiacenti.

per questi due gruppi di cellule:

$$\text{Gruppo1} = \{7, 8, 9, 29\} \quad \text{Gruppo2} = \{2, 3, 5, 6, 16, 32, 36, 37\}$$

in effetti, confrontando gli andamenti delle metriche descritte al capitolo precedente, hanno mostrato dei comportamenti consistenti. Si veda in figura 24 come per il Gruppo1 si abbia un andamento del Boundary Activity coerente fra le cellule

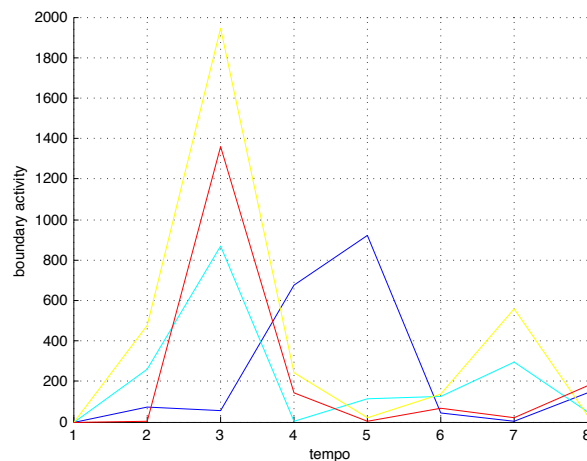


Figura 24.: Boundary Activity per il Gruppo 1.

Si veda in figura 25 invece come per il Gruppo2 si abbia un andamento dell'elongazione coerente fra le cellule

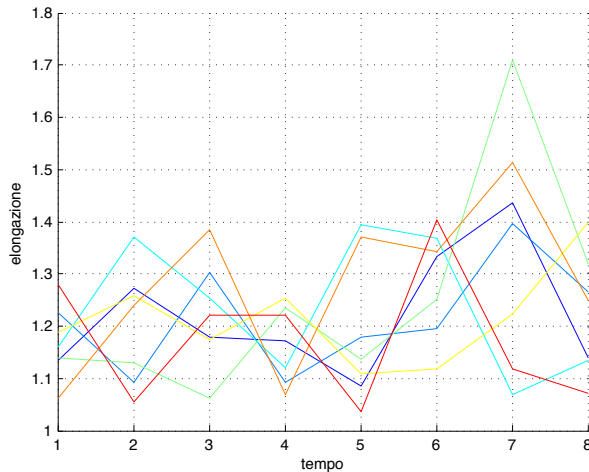


Figura 25.: Elongazione per il Gruppo 2.

A questo punto è interessante cercare di capire se esiste una legge di correlazione fra i i valori di tutte le metriche, calcolate per ogni cellula. Nel paragrafo seguente si mostrerà come è stata condotta l'analisi statistica dei dati di deformazione.

Analisi di correlazione dei parametri di deformazione

Una volta calcolate tutte le metriche descritte al capitolo 3 per tutte le n cellule si è cercato di dare una descrizione statistica dei risultati. Per farlo si sono calcolate le statistiche delle metriche messe a confronto.

$$m(\mathbf{x}) = \frac{1}{n} \sum_{i=1}^n x_i$$

e le varianze campionarie

$$s^2(\mathbf{x}) = \frac{1}{n-1} \sum_{i=1}^n [x_i - m(\mathbf{x})]^2$$

Il primo passo da fare nell'analisi esplorativa dei dati è spesso quello di dar un senso visuale della relazione statistica tra le variabili attraverso i cosiddetti *scatterplots*. In particolare è di interesse capire se la nuvola di punti assume un trend lineare o curvilineo. L'informazione che si vuole estrarre è come una variabile x possa essere utilizzata per predire la variabile y .

Il passo successivo è quello di definire delle grandezze in grado di misurare l'associazione tra i dati x ed y .

Si definisce *covarianza campionaria*

$$s(\mathbf{x}_1, \mathbf{x}_2) = \frac{1}{n-1} \sum_{i=1}^n [x_{1i} - m(\mathbf{x}_1)][x_{2i} - m(\mathbf{x}_2)]$$

e le correlazioni campionarie

$$c(\mathbf{x}_1, \mathbf{x}_2) = \frac{s(\mathbf{x}_1, \mathbf{x}_2)}{s^2(\mathbf{x}_1) \cdot s^2(\mathbf{x}_2)}$$

infine, avendo questi dati, si può calcolare la retta di regressione lineare² per ogni scatterplot

$$\mathbf{y} = m(\mathbf{y}) + \frac{s(\mathbf{x}, \mathbf{y})}{s^2(\mathbf{x})} [\mathbf{x} - m(\mathbf{x})]$$

con \mathbf{y} variabile da stimare in base a \mathbf{x} .

Si rimanda all'appendice A per maggiori dettagli.

Prima di mostrare i risultati ottenuti è bene specificare quali siano le metriche che effettivamente ha senso confrontare. Una metrica fra tutte si differenzia per il fatto di essere *dinamica*, ossia essa sfrutta in maniera diretta la dipendenza da uno slot temporale fissato: il Boundary Activity. Quello che ci si aspetta però è che essa sia in qualche modo legata alla variazione nell'elongazione della cellula. Infatti ha senso se si immagina che quando una cellula si stira e assume una forma più allungata, il suo perimetro cambi in maniera proporzionale. Dal punto di vista biologico, come si è anticipato poco fa, la condizione di elevata elongazione è indice di instabilità per la forma della cellula; quindi cercare di quantificare questo parametro può dare un'idea di quale sia il valore di soglia che una cellula può raggiungere prima di dividersi e in che maniera questo sia influenzato dal comportamento delle cellule adiacenti. Le due metriche che per ovvie ragioni sono consistentemente confrontabili sono le due differenze di aree, percentuale e media si veda la figura 26. Le figure che si riportano di seguito mostrano il risultato di una regressione lineare statistica sui dati di deformazione: la 26 mostra la relazione lineare che lega le due metriche di variazione di area, mentre la 27 quella che lega (inaspettatamente) le metriche di variazione percentuale di area e distanza di Hausdorff. In effetti va precisato che però tali scatterplots permettono sì di calcolare una retta di regressione lineare fra i punti, ma la nuvola di dati presenta valori troppo sparsi per poter effettivamente concludere che le metriche sono in relazione lineare.

² in MATLAB implementata dai metodi polyfit e polyval

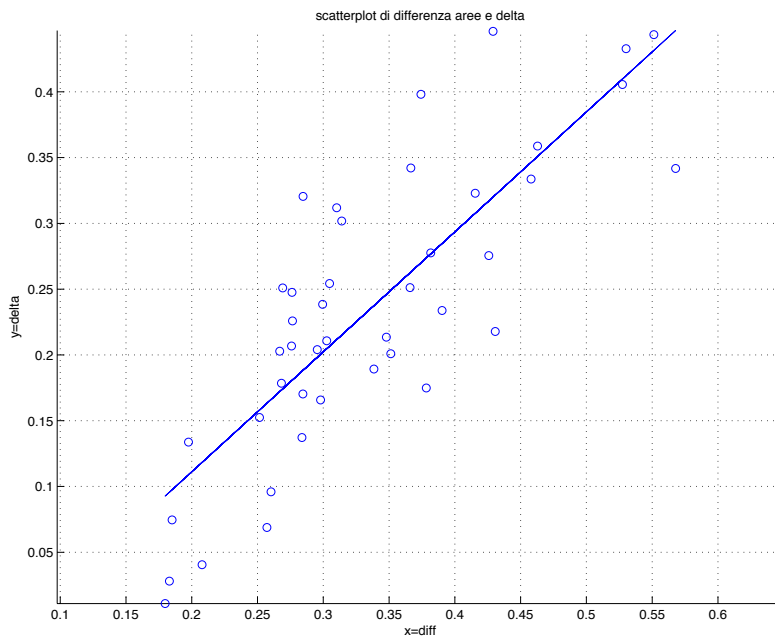
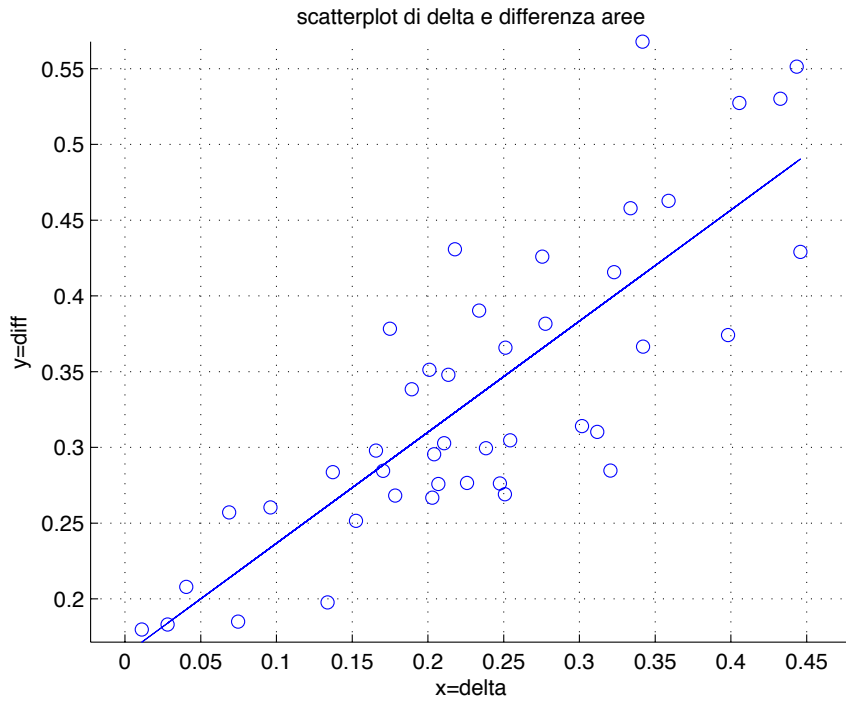


Figura 26.: Scatterplots DIFF Vs DELTA.

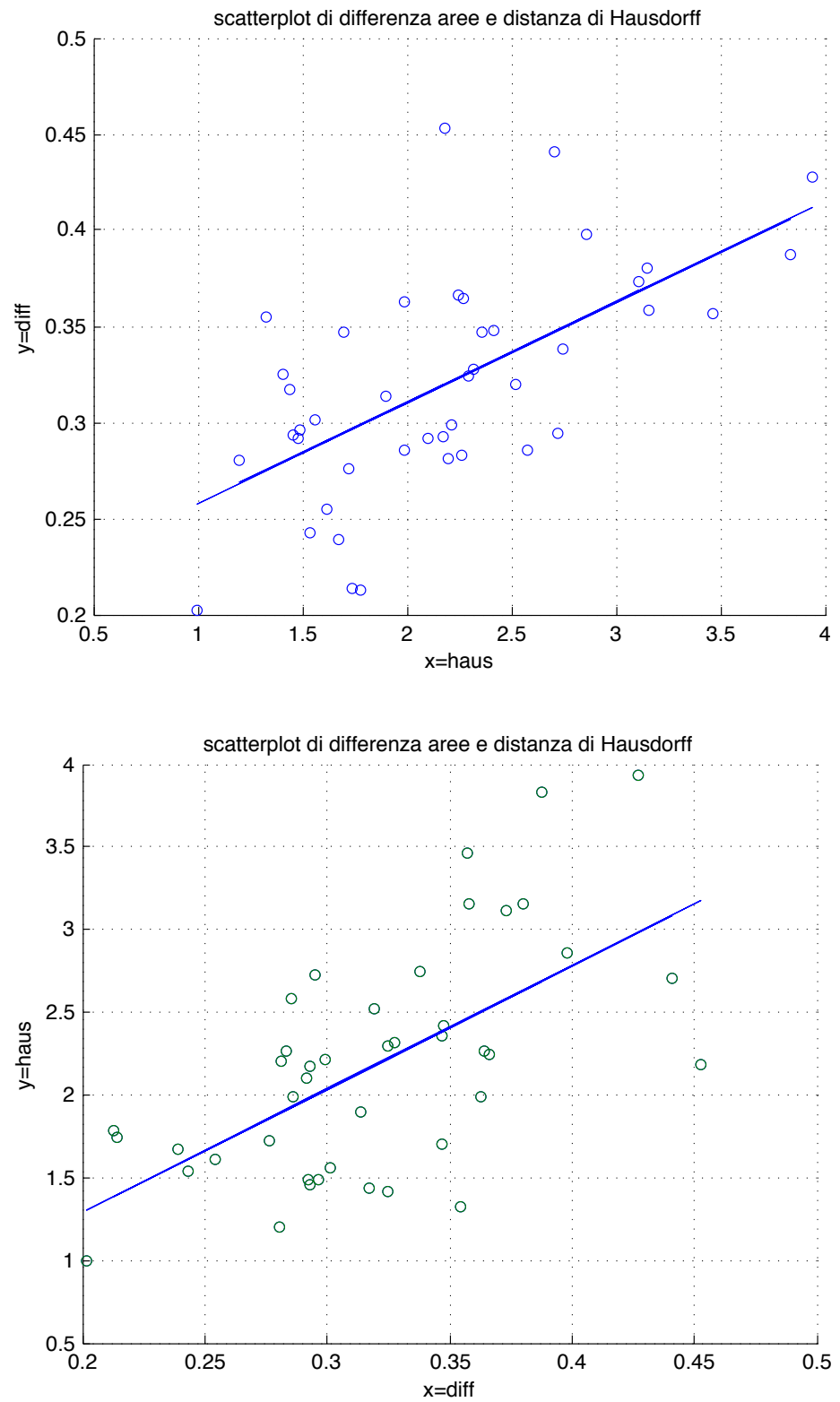


Figura 27.: Scatterplots DIFF Vs HAUS.

5

FILTRAGGIO STATISTICO

Indice

5.1	Il problema della stima	51
5.2	Le teorie del filtraggio statistico	53
5.3	Filtro di Kalman	54
5.3.1	Le equazioni del filtro di Kalman	55
5.4	Applicazione del filtraggio al modello di riferimento	57

Questo capitolo è dedicato alla descrizione della fase del progetto in cui si è utilizzata la teoria del filtraggio statistico per stimare il moto delle cellule in analisi.

Il problema del filtraggio (o della misura) corrisponde alla progettazione di algoritmi per la ricostruzione dei segnali non direttamente accessibili, come segnali che sono stati trasmessi da sorgenti remote o descrivono variabili "interne" essenziali del sistema, non "visibili" direttamente. La ricostruzione opera a partire da altre variabili osservabili che sono accessibili direttamente e che, tipicamente, rappresentano versioni distorte e rumorose (affette da errori di misura) delle prime. La "misura" in questa accezione generale comprende classi di problemi molto generali chiamati *problemi inversi* in matematica applicata; la sottoclasse cui si fa riferimento sono i problemi di *stima e filtraggio statistico*.

Il primo passo da compiere nell'affrontare un problema di filtraggio statistico è quello dell'*identificare* un modello matematico adeguato per rappresentare il sistema cui si vuole analizzarne il comportamento. Quello che bisogna fare è trovare una descrizione matematica, il più possibile semplice, del sistema e dei segnali che ne influenzano il comportamento.

5.1 IL PROBLEMA DELLA STIMA

In questo capitolo consideriamo sia i dati che la grandezza da stimare come *processi stocastici* $\{y(t)\}$ e $\{x(t)\}$ e lo stimatore della classe degli *stimatori lineari*, in modo tale che sarà sufficiente disporre di descrizioni che specificano *momenti congiunti* del primo e del secondo ordine di $\{y(t)\}$ e $\{x(t)\}$. La soluzione dei problemi di stima dinamica sui processi stocastici è basata fortemente sulle proprietà che caratte-

rizzano il sistema dinamico modellizzato.

Di base ci sono tre tipi di problemi legati alla stima:

1. Predizione

Si vuole stimare il processo

$$\mathbf{x}(t) = \mathbf{y}(t+h) \quad h \in \mathbb{Z}^+$$

Si cerca quindi, per ogni t , la migliore predizione "h passi in avanti" di $\{\mathbf{y}(t)\}$ basata su tutte le misure disponibili all'istante t , $\{\mathbf{y}(s); t_0 \leq s \leq t\}$ e si denota con $\hat{\mathbf{y}}(t+h | t)$.

2. Filtraggio

Le misure $\{\mathbf{y}(t)\}$ sono elaborate per ricostruire un segnale $\{\mathbf{x}(t)\}$ in generale distorto o affetto da rumore:

$$\mathbf{y}(t) = H \circ \mathbf{x}(t) + \mathbf{n}(t)$$

dove $H \circ$ denota una trasformazione nota e solitamente $\mathbf{x}(t)$ e $\mathbf{n}(t)$ sono ipotizzati scorrelati. La ricostruzione (stima) di $\mathbf{x}(t)$ deve essere calcolata ad ogni istante sfruttando tutta l'informazione disponibile fino a quell'istante ($I = [t_0, t]$); e si indica col simbolo $\hat{\mathbf{x}}(t | t)$.

3. Interpolazione

Nei problemi di interpolazione si considera possibile utilizzare anche i dati "futuri" per ottenere le stime di $\mathbf{x}(t)$. In genere si usa quando non serve la stima del segnale per prendere delle decisioni in tempo reale o comunque non è importante se la sua ricostruzione viene fatta con un certo ritardo, ma si richiede piuttosto la bontà della ricostruzione del messaggio $\mathbf{x}(t)$. Solitamente la base temporale è ipotizzata fissa del tipo $I = [t_0, t_1]$ e si chiede di ricostruire $\{\mathbf{x}(t)\}$ per ogni $t \in I$. L'interpolatore si indica con il simbolo $\hat{\mathbf{x}}(t | I)$.

In generale la ricerca della soluzione di un problema di stima si riconduce all'implementazione pratica delle formule risolutive. Prendiamo un esempio di stima in tempo reale della variabile corrente del processo $\{\mathbf{x}(t)\}$ basata su osservazioni $\{\mathbf{y}(t)\}$ che supponiamo vengano acquisite per un tempo indefinito. Prendiamo istante iniziale definito t_0 ; uno *stimatore lineare* di $\mathbf{x}(t)$ in base ai dati osservati $\{\mathbf{y}(s); t_0 \leq s \leq t\}$ avrà la forma

$$\hat{\mathbf{x}}(t | t) = \sum_{t_0}^t H(t, k) \mathbf{y}(k)$$

$H(t, \cdot)$ può essere immaginata come risposta impulsiva di un sistema lineare con ingresso i dati di misura e in uscita le stime $\hat{\mathbf{x}}(t | t)$ di $\mathbf{x}(t)$.

5.2 LE TEORIE DEL FILTRAGGIO STATISTICO

Ogni problema pratico che si intende risolvere con algoritmi di stima su processi stocastici porta con sé il fatto di dover introdurre delle ipotesi riguardanti i segnali aleatori che si utilizzano nel modello dinamico. Ciò è legato alla possibilità di ottenere soluzioni praticamente computabili mediante un procedimento teorico di stima. Esistono tre principali "teorie del filtraggio statistico" le quali si differenziano principalmente sulle assunzioni di partenza fatte sui segnali in gioco. Le ipotesi di partenza risulteranno più o meno limitative a seconda dello specifico contesto applicativo, ed è chiaro che a rigore nessuna ipotesi è mai verificata "esattamente" nei problemi reali. Va da sé che si cercherà di utilizzare la teorizzazione più robusta rispetto alle variazioni sui segnali in modo tale da poterla applicare a situazioni a stretto rigore non contemplate nelle ipotesi fatte, ma ottenendo risultati comunque accettabili.

WIENER-KOLMOGOROV Il primo approccio storicamente formalizzato è quello dovuto a N. Wiener e A. N. Kolmogorov (anni '40). Tale approccio si basa su due ipotesi:

- L'intervallo di osservazione è illimitato inferiormente $I = (-\infty, t]$
- i processi $\{x(t)\}$ e $\{y(t)\}$ sono congiuntamente stazionari (in senso debole).

Si ha che in queste ipotesi la risposta impulsiva dello stimatore $H(t, k)$ dipende solo dalla differenza dei due argomenti temporali t e k . Perciò lo stimatore si può vedere come un sistema lineare invariante nel tempo e, note le statistiche dei segnali, può essere calcolato una volta per tutte. Ciò che si cerca in un approccio di W.-K. è quello di calcolare esplicitamente la risposta impulsiva $\{H(k \mid k \geq 0)\}$ oppure la corrispondente funzione di trasferimento $\hat{H}(z)$ in forma chiusa, supponendo le covarianze $\Sigma_{xy}(k)$ e $\Sigma_y(k)$ o gli spettri $S_{xy}(z)$ e $S_y(z)$ noti analiticamente.

LEVINSON Il secondo approccio viene fornito da Levinson il quale rinuncia all'ipotesi di intervallo di osservazione infinito ed introduce il fatto che la matrice di varianza dei processi abbia la struttura di una "matrice di Toeplitz"; ciò pone enfasi sul ruolo essenziale che gioca la stazionarietà. In tale contesto comunque la risposta impulsiva dello stimatore varia al variare di t ed è quindi comunque non stazionaria. Infine la soluzione di Levinson propone uno schema ricorsivo per il calcolo di $H(t, \cdot)$, cioè un algoritmo che aggiorna $H(t, \cdot)$ (e la stima) in modo veloce ed efficiente.

KALMAN Il terzo approccio (che viene sfruttato in questo progetto) è particolarmente adatto ad essere implementato su sistemi digitali e si caratterizza per la sua struttura ampiamente ricorsiva. Le soluzioni alla Kalman hanno la forma di equazioni di evoluzione di stato tipiche di un sistema dinamico, come

$$\hat{\mathbf{x}}(t+1) = \mathbf{A}\hat{\mathbf{x}}(t) + \mathbf{K}\mathbf{y}(t+1)$$

che è un algoritmo a memoria fissa che ad ogni passo aggiorna la "vecchia" stima $\hat{\mathbf{x}}(t)$ conseguentemente all'acquisizione della "nuova" misura $\mathbf{y}(t+1)$. Questo approccio si basa fondamentalmente quindi sull'ipotesi che i processi in gioco siano descrivibili mediante modelli dinamici piuttosto che mediante dati probabilistici come le covarianze.

5.3 FILTRO DI KALMAN

Il filtro di Kalman è uno strumento statistico che permette di stimare variabili, anche non misurabili, a partire da valori misurati di alcune grandezze ad esse legate. Si basa su alcune ipotesi fondamentali, quali la dipendenza dello stato attuale del sistema solo dallo stato al tempo precedente (processo markoviano del primo ordine) e la linearità di tale dipendenza. Esso affronta il problema generale di cercare di stimare dei parametri di interesse a partire da misure indirette, imprecise (affette cioè da errore di misura) e incerte. Grazie alla sua struttura ricorsiva, nuovi dati di misurazione possono essere reimmessi nel sistema, e ciò permette in linea generale che esso sia utilizzabile in un sistema di elaborazione real-time. Il filtro di Kalman genera infatti la stima di un processo utilizzando una forma di controllo in retroazione: stima lo stato di un processo ad un certo istante di tempo e successivamente ottiene feedback nella forma di misure affette da rumore. Per applicare il filtro di Kalman è necessario determinare due relazioni: l'equazione di stato e quella di misura. L'equazione di stato esprime la relazione tra i parametri di stato di due epoche successive e rappresenta il modello matematico scelto per descrivere il fenomeno che si sta analizzando. L'equazione di misura invece lega il valore dei parametri alle misure di grandezze ad essi legate. Entrambe le equazioni devono essere lineari. Le equazioni di aggiornamento sono responsabili della proiezione in avanti dello stato attuale e le stime di covarianza di errore servono ad ottenere le stime a priori per la fase successiva. Le equazioni di misura di aggiornamento rappresentano il feedback, cioè l'inserimento di una nuova misura nella stima a priori in modo da ottenere una stima migliore a posteriori. Il fatto chiave è l'ipotesi che i segnali in gioco siano descrivibili per mezzo di *sistemi dinamici (lineari) di dimensione finita*. La fondamentale novità dell'approccio di Kalman sta nell'introdurre

questo tipo di modelli fin dall'inizio, dall'atto stesso di formulazione del problema.

Il filtro di Kalman ha una struttura algoritmica fissa e nota a priori a meno di un parametro. Si tratta di un sistema dinamico lineare a tempo discreto e a memoria finita, pilotato dalle osservazioni, che costituisce quindi uno schema numerico ricorsivo adatto all'elaborazione di segnali in linea. I suoi parametri sono noti a priori a partire dalle specifiche del modello, a meno di una matrice di guadagno che può essere calcolata analiticamente fuori linea risolvendo un'equazione alle differenze non lineare di Riccati.

Il filtro di Kalman si divide in due parti: la prima, il *filtering*, permette di determinare la miglior stima dei parametri al tempo corrente sfruttando le misure fino al tempo attuale; la seconda, detta *smoothing*, permette di determinare la miglior stima dei parametri dei tempi precedenti tenendo conto di tutte le misure fino all'ultimo istante di tempo. Gli stimatori sono ottimi nel senso di Wiener-Kolmogorov, hanno minima varianza e sono distribuiti normalmente. L'ottimalità tuttavia è garantita solo finché le assunzioni fatte nel modello sono valide. Per garantire la stabilità del filtro di Kalman è richiesto unicamente che siano soddisfatte le condizioni di osservabilità e controllabilità. Uno dei principali vantaggi di tale metodo consiste nel fatto che le dimensioni delle matrici coinvolte nel calcolo della stima dipendono solamente dalle misure e dai parametri del tempo attuale e al più di quello successivo.

5.3.1 Le equazioni del filtro di Kalman

Ciò che si vuole fare è cercare un modello di stato che descriva il segnale di osservazione $\{\mathbf{y}(t)\}$. Poiché i modelli di stato possono essere combinati linearmente fra loro mantenendo la medesima struttura, si ha che tutti i dati a priori sul problema di stima possono essere rappresentati mediante un unico modello lineare di dimensione finita del tipo

$$\begin{cases} \mathbf{x}(t+1) = \mathbf{A}\mathbf{x}(t) + \mathbf{v}(t) \\ \mathbf{y}(t) = \mathbf{C}\mathbf{x}(t) + \mathbf{w}(t) \end{cases} \quad t \geq t_0 \quad (41)$$

dove $\{\mathbf{y}(t)\}$ è il processo m -dimensionale delle osservazioni, $\{\mathbf{v}(t)\}$ e $\{\mathbf{w}(t)\}$ sono *rumori bianchi* (di media zero) aventi matrice di covarianza

$$\mathbb{E} \left\{ \begin{bmatrix} \mathbf{v}(t) \\ \mathbf{w}(t) \end{bmatrix} \begin{bmatrix} \mathbf{v}(s)' & \mathbf{w}(s)' \end{bmatrix} \right\} = \begin{bmatrix} \mathbf{Q} & \mathbf{S} \\ \mathbf{S}' & \mathbf{R} \end{bmatrix} \delta(t-s) \quad (42)$$

lo stato $\{\mathbf{x}(t)\}$ del modello 41 è un processo n -dimensionale non direttamente osservabile, il cui valore iniziale $\mathbf{x}(t_0) = \mathbf{x}_0$ è un vettore aleatorio scorrelato dall'andamento presente e futuro dei rumori di modello $\{\mathbf{v}(t)\}$ e di misura $\{\mathbf{w}(t)\}$

$$E \{ \mathbf{x}_0 [\mathbf{v}(s)' \quad \mathbf{w}(s)'] \} = 0 \quad \forall t \geq t_0 \quad (43)$$

Si assumerà anche che media e varianza di \mathbf{x}_0 siano note

$$E \mathbf{x}_0 = \mu_0 \quad \text{Var}\{\mathbf{x}_0\} = P_0 \quad (44)$$

Si osservi che $\mathbf{v}(t)$ è correlata solo con $\mathbf{w}(t)$ essendo per ipotesi $\mathbf{v}(t) \perp \mathbf{w}(s)$ per $s \neq t$. Si può quindi definire l'errore di stima come

$$\tilde{\mathbf{v}}(t) := \mathbf{v}(t) - \hat{E}[\mathbf{v}(t) | \mathbf{H}\{\mathbf{w}\}] = \mathbf{v}(t) - \hat{E}[\mathbf{v}(t) | \mathbf{w}(t)] \quad (45)$$

che è ortogonale all'intero processo rumore di misura $\{\mathbf{w}(t)\}$ e si può riscrivere come

$$\tilde{\mathbf{v}}(t) = \mathbf{v}(t) - \mathbf{S}\mathbf{R}^{-1}\mathbf{w}(t) \quad (46)$$

quindi $\{\tilde{\mathbf{v}}(t)\}$ è rumore bianco di varianza

$$\tilde{\mathbf{Q}} = \mathbf{Q} - \mathbf{S}\mathbf{R}^{-1}\mathbf{S}' \quad (47)$$

sostituendo la formula 46 nel modello 41 e considerando che $\mathbf{w}(t) = \mathbf{y}(t) - \mathbf{C}\mathbf{x}(t)$ si può riscrivere il modello 41 nella forma

$$\begin{cases} \mathbf{x}(t+1) = \mathbf{F}\mathbf{x}(t) + \mathbf{S}\mathbf{R}^{-1}\mathbf{y}(t) + \tilde{\mathbf{v}}(t) \\ \mathbf{y}(t) = \mathbf{C}\mathbf{x}(t) + \mathbf{w}(t) \end{cases} \quad (48)$$

dove i rumori bianchi $\{\tilde{\mathbf{v}}(t)\}$ e $\{\mathbf{w}(t)\}$ sono tra loro scorrelati e hanno variante $\tilde{\mathbf{Q}}$ ed \mathbf{R} ; la matrice \mathbf{F} è data da

$$\mathbf{F} = \mathbf{A} - \mathbf{S}\mathbf{R}^{-1}\mathbf{C} \quad (49)$$

Il termine forzante $\mathbf{S}\mathbf{R}^{-1}\mathbf{y}(t)$ ha proprio il significato di retroazione dall'uscita sullo stato.

Algoritmo 1 (Filtro di Kalman). *Gli stimatori lineari a minima varianza $\hat{\mathbf{x}}(t+1 | t)$ e $\hat{\mathbf{x}}(t | t)$ dello stato del modello lineare 41 all'istante $t+1$ e t , in base alle osservazioni $\{\mathbf{y}(s); t_0 \leq s \leq t\}$, sono calcolabili mediante il seguente algoritmo ricorsivo*

1. *Stime a priori* (Aggiornamento temporale)

$$\hat{\mathbf{x}}(t+1 | t) = \mathbf{F}\hat{\mathbf{x}}(t | t) + \mathbf{S}\mathbf{R}^{-1}\mathbf{y}(t) \quad (50)$$

$$\mathbf{P}(t+1 | t) = \mathbf{F}\mathbf{P}(t | t)\mathbf{F}' + \tilde{\mathbf{Q}} \quad t \geq t_0 \quad (51)$$

2. *Stime a posteriori* (Aggiornamento rispetto alle misure)

$$\hat{\mathbf{x}}(t+1 | t+1) = \hat{\mathbf{x}}(t+1 | t) + \mathbf{L}(t+1)[\mathbf{y}(t+1) - \mathbf{C}\hat{\mathbf{x}}(t+1 | t)]$$

$$\begin{aligned} \mathbf{P}(t+1 | t+1) &= \mathbf{P}(t+1 | t) - \\ &- \mathbf{P}(t+1 | t)\mathbf{C}'\mathbf{\Lambda}(t+1)^{-1}\mathbf{C}\mathbf{P}(t+1 | t) \end{aligned} \quad (52)$$

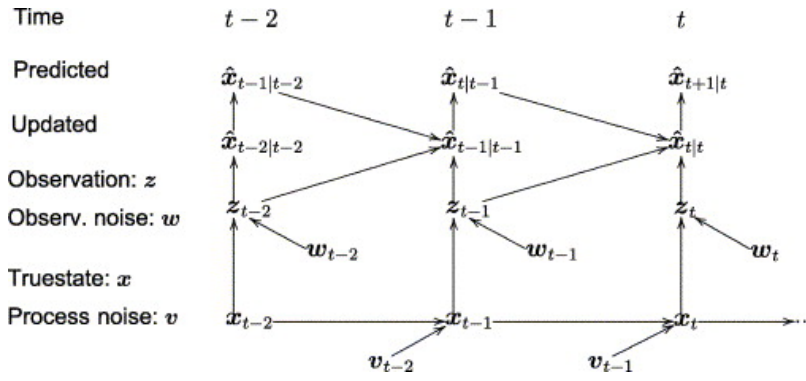


Figura 28.: Algoritmo per il calcolo della stima mediante filtro di Kalman

3. Condizioni iniziali

$$\hat{x}(t_0 | t_0 - 1) = \mu_0 \quad P(t_0 | t_0 - 1) = P_0 \quad (53)$$

Dove le matrici F e \tilde{Q} sono definite dalle 49 e 47; $P(t+1 | t)$ e $P(t | t)$ sono le varianze di errore di predizione e filtraggio

$$P(t+1 | t) = E \{ \tilde{x}(t+1 | t) \tilde{x}(t+1 | t)' \} \quad (54)$$

$$P(t | t) = E \{ \tilde{x}(t | t) \tilde{x}(t | t)' \} \quad (55)$$

$\Lambda(t)$ è la varianza del processo di innovazione $e(t) = y(t) - C\hat{x}(t | t-1)$

$$\Lambda(t) = CP(t | t-1)C' + R \quad (56)$$

e il guadagno del filtro è definito dalla

$$L(t) = P(t | t-1)C'\Lambda(t)^{-1} . \quad (57)$$

5.4 APPLICAZIONE DEL FILTRAGGIO AL MODELLO DI RIFERIMENTO

Proponiamo un tracking basato sul filtro di Kalman applicato ad un modello di *random walk* per predire la posizione delle cellule (rappresentata schematicamente dalla posizione dei loro centroidi) in un range di frames, quindi aggiornare tale predizione usando come osservazioni i valori misurati delle posizioni dei centroidi.

Consideriamo come sistema di random walk associato al moto della cellula un sistema di aggiornamento dello stato del tipo 41

$$\begin{cases} x(t+1) = Ax(t) + Bv(t) \\ y(t) = Cx(t) + Dw(t) \end{cases} \quad t \geq t_0 \quad (58)$$

con matrici definite come

$$\begin{aligned}\mathbf{A} &= \begin{bmatrix} 1 & 0 \\ 0 & 1 \end{bmatrix} \\ \mathbf{B} &= dt \cdot \begin{bmatrix} 1 & 0 \\ 0 & 1 \end{bmatrix} \\ \mathbf{C} &= \begin{bmatrix} 1 & 0 \\ 0 & 1 \end{bmatrix} \\ \mathbf{D} &= \begin{bmatrix} 1 & 0 \\ 0 & 1 \end{bmatrix}\end{aligned}$$

dove dt è il passo di campionamento che rappresenta la velocità di spostamento della cellula. Dato che non è esplicito mediante la legge di moto delle cellule come modellizzare questo parametro, è stato possibile sceglierne un valore consistente grazie ad un tuning manuale fatto osservando l'andamento del modello confrontato col moto vero di una cellula. In particolare esso è stato fissato pari a $dt = 20$.

Il vettore delle misure sia

$$\mathbf{z}(t, :) = [x_c(t) \quad y_c(t)]$$

in cui $x_c(t)$ e $y_c(t)$ sono le reali posizioni dei centroidi di una cellula all'istante t nello spazio Cartesiano. Si è scelto di fornire come osservazioni i valori di \mathbf{z} mediati su tutti i tempi. Vale a dire, presa $m = \text{mean}(\mathbf{z})$ si prendono come osservazioni i valori

$$\mathbf{z}_2(t, :) = [x_c(t) - m(1) \quad y_c(t) - m(2)]$$

quindi i valori delle osservazioni sottratte della media.

Si scelgono come condizioni iniziali dell'algoritmo

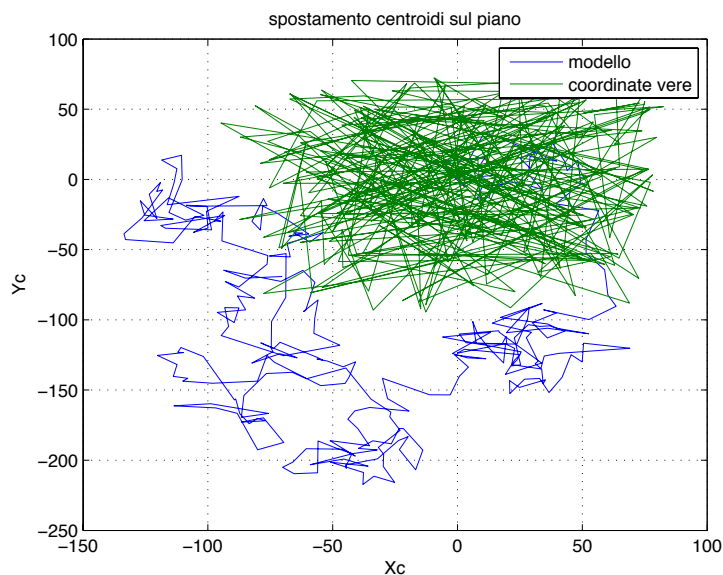
$$\hat{\mathbf{x}}(t_0 | t_0 - 1) = \begin{bmatrix} 0 \\ 0 \end{bmatrix} \quad (59)$$

$$P(t_0 | t_0 - 1) = \alpha \cdot \begin{bmatrix} 0 & 0 \\ 0 & 0 \end{bmatrix} \quad \text{con } \alpha = 100 \quad (60)$$

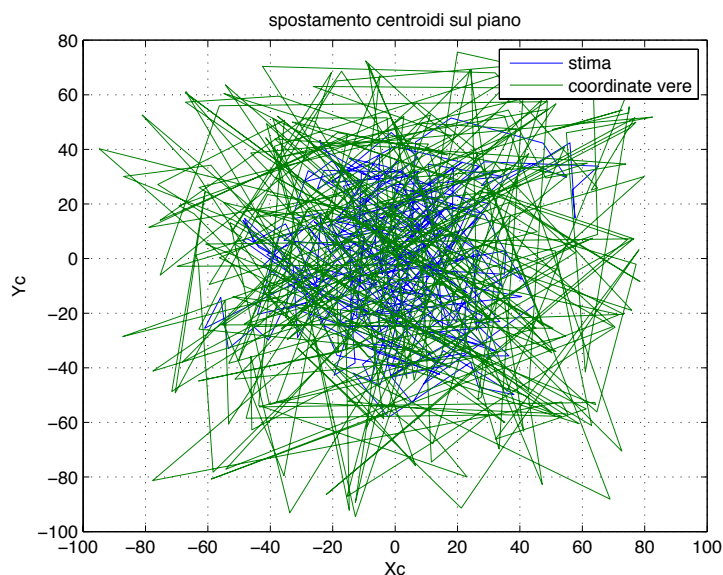
dopodiché si sono calcolate le matrici del filtro

$$\begin{aligned}\mathbf{Q} &= \mathbf{B} \cdot \sigma_v^2 \mathbf{B}^T \\ \tilde{\mathbf{Q}} &= \mathbf{Q} - \mathbf{S} \cdot (\mathbf{R}^{-1}) \cdot \mathbf{S}^T \\ \mathbf{S} &= \begin{bmatrix} 0 & 0 \\ 0 & 0 \end{bmatrix} \\ \mathbf{R} &= \mathbf{D} \cdot \sigma_w^2 \mathbf{D}^T \\ \mathbf{F} &= \mathbf{A} - \mathbf{S} \cdot (\mathbf{R}^{-1}) \cdot \mathbf{C}\end{aligned}$$

I vettori $\mathbf{v}(t)$ e $\mathbf{w}(t)$ sono stati inizialmente scelti come rumori bianchi ($\mu_v = \mu_w = 0$ e $\sigma_v = 0.002$ e $\sigma_w = 0.1$) e rappresentano, il primo il rumore di modello e il secondo il rumore di misura. Scelti in questo modo impongono al filtro di basarsi maggiormente sui valori delle misure fornite piuttosto che sui valori dello stato calcolato mediante il modello di random walk.



(a) Modello di random walk Vs moto cellula di riferimento.



(b) Stima calcolata col filtro di Kalman Vs moto cellula di riferimento.

Figura 29.

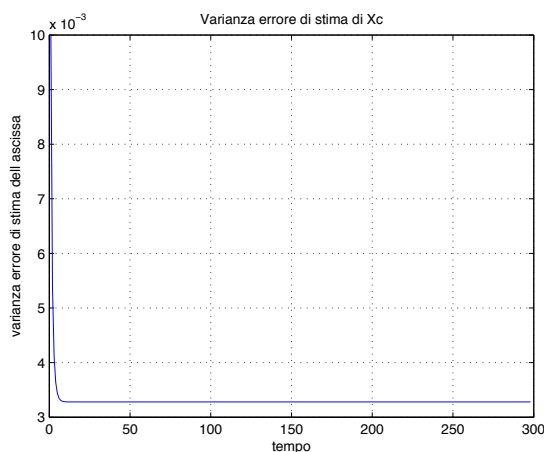
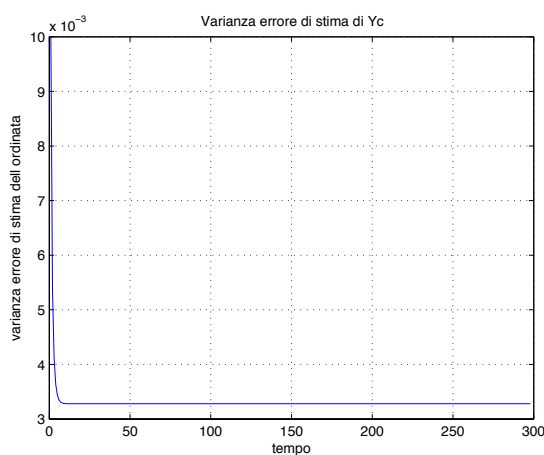
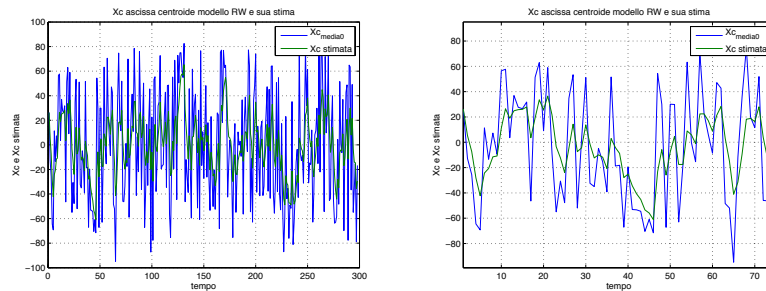
(a) Varianza errore di stima dell'ascissa X_C del centroide.(b) Varianza errore di stima dell'ordinata Y_C del centroide.

Figura 30.

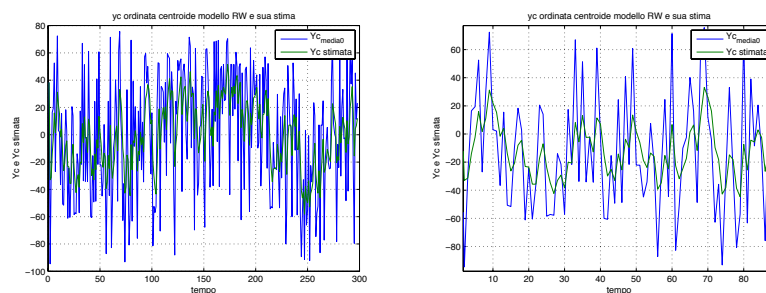
Dalle figure 31 si può osservare l'andamento della stima dell'ascissa e dell'ordinata della cellula. La stima insegue in maniera grossolana l'andamento vero di entrambe, presentando poi alla fine un errore di stima che converge a zero in pochi passi: si osservino le figure 30 (a) e (b). Si può quindi dire che la stima fatta è soddisfacente.

L'intento di questo modello è quello di fornire una predizione sullo spostamento delle cellule all'interno dei frames.

Per capire in che misura la modellizzazione fatta rispecchia il reale moto delle cellule osserviamo la figura 29 (a). Nel grafico seguente il moto vero dei centroidi delle cellule è messo in comparazione con il modello di random walk scelto per rappresentarlo. In effetti modellare il moto delle cellule come semplice random walk (rappresentato dalla spezzata di colore blu) non è del tutto ingenuo e anzi sembra proprio avere senso. L'unica accortezza che si può avere è considerare che una cellula si muova in maniera random ma all'interno di un'area abbastanza ben delimitata che la circonda. Risulta



(a) Stima dell'ascissa X_C del centri- (b) Zoom stima dell'ordinata Y_C del
de. centroide.



(c) Stima dell'ascissa X_C del centri- (d) Zoom stima dell'ordinata Y_C del
de. centroide.

Figura 31.

evidente a occhio, nonché suggerito dal buon senso, che una cellula all'interno del video "spazzoli" un'area di spazio non eccessivamente estesa. Infatti si osservi la spezzata di colore verde e si noti come il moto non si estende al di fuori di una nuvola quadrata compresa fra i punti $[-100, 70]$ in ascissa e $[-100, 60]$ in ordinata. Considerato che la lunghezza media dell'asse maggiore di una cellula, ad esempio $CellID = 10$, è 33.9 si può intuire che la cellula copra, muovendosi in maniera casuale, uno spazio di cinque volte la sua grandezza, sia in orizzontale che in verticale.

Quello che si potrebbe fare quindi è di forzare il modello di random walk all'interno di quest'area anche se dal punto di vista della stima non cambia il risultato visto che comunque il filtro ha come input i valori delle misure vere delle posizioni del centroide di riferimento.

6

CONCLUSIONI E SVILUPPI FUTURI

In questa tesi si è studiata una tecnica per analizzare strutture reticolari presenti in immagini biologiche. L'obiettivo dello studio è stato quello di implementare un sistema di riconoscimento automatico per una struttura reticolare di cellule, al fine di presentare delle linee guida per l'analisi di moto e deformazione delle cellule appartenenti al reticolo. Per giungere agli obiettivi suddetti si sono inizialmente studiate alcune tecniche di *static detection* di una singola forma (cellula) all'interno di un'immagine; successivamente, estendendo queste tecniche al caso dinamico, ci si è concentrati su metodi di *dynamic detection* (o *tracking*) di una cellula all'interno di una sequenza temporizzata di immagini (video). In particolare, a tal proposito, si è studiato l'algoritmo di *shape detection* basato sulla minimizzazione di un funzionale di energia: l'algoritmo *Random Walk Agents*. Grazie all'implementazione di questo algoritmo (in linguaggio MATLAB) si è costruita una rappresentazione schematica del reticolo cellulare e si sono memorizzate per ogni cellula, in opportune variabili, le informazioni che sono poi servite a studiare la dinamica del moto e della deformazione. Sono state introdotte e testate in questo lavoro alcune metriche per condurre uno studio sull'evoluzione della forma di singole cellule e di gruppi di esse all'interno del reticolo. Successivamente è stata condotta un'analisi di regressione statistica sulla correlazione fra le metriche proposte, tuttavia i risultati che si sono ottenuti non forniscono sufficiente evidenza di una relazione statistica fra di esse. Nella parte finale della tesi l'obiettivo è stato quello di proporre un predittore del moto delle cellule del reticolo. A tale scopo si è proposto un filtraggio alla Kalman applicato ad un modello di moto delle cellule di tipo random walk. Questo modello è isotropico e non polarizzato e ciò significa che ogni cellula ha la stessa probabilità di muoversi in qualunque direzione indipendentemente dalle direzioni dei passi precedenti. Si è osservato che in effetti tale modello è adeguato a rappresentare il moto reale dei centroidi delle cellule. Il filtro di Kalman applicato al modello di random walk è riuscito a produrre la stima della posizione di una cellula con una varianza d'errore molto bassa.

Sviluppi futuri possono essere condotti a ciascun livello del lavoro proposto. Dal un punto di vista dell'analisi di moto e deformazione, si potrebbe cercare di andare più a fondo nelle relazioni che legano la deformazione locale di una singola cellula alla deformazione del dominio di cellule che la circondano. Ancora, sarebbe possibi-

le investigare sulla modellizzazione di alcune metriche dinamiche di deformazione dipendenti da un insieme di parametri di *deformation* di ogni singola cellula. Un'altra analisi che si potrebbe condurre per studiare il moto del reticolo cellulare, potrebbe avere l'obiettivo di capire se esiste un legame fra la direzione del moto di una cellula e la sua direzione di massima estensione. Da un punto di vista biologico inoltre può essere utile cercare di riconoscere ed inseguire i fenomeni di "*cell-intercalation*" in cui una cellula si inserisce in un nuovo dominio di cellule. Questa analisi potrebbe risultare rilevante nel cercare le relazioni che legano la deformazione locale a quella globale del reticolo. La ricerca di pattern di geometrie particolari, come proposto anche in [Gibson *et al.*, 2006](#) e in [Gibson *et al.*, 2007](#), si potrebbe inoltre applicare al reticolo di cellule epiteliali di *Drosophila*.

Dal punto di vista del filtraggio statistico infine, si potrebbe proporre un modello esteso per il moto delle cellule che consideri anche al suo interno parametri, come l'angolo di inclinazione della cellula o la lunghezza dei suoi assi.

A

RICHIAMI DI REGRESSIONE STATISTICA

Il modello base della statistica prevede che si consideri una popolazione di oggetti di interesse e diverse misure, trattate come variabili, effettuate sugli oggetti stessi. Si selezionino quindi degli oggetti dalla popolazione costituenti un *campione* e si registrino le variabili associate che avranno il ruolo di dati. Al fine di una descrizione qualitativa si assuma che i dati siano ottenuti in maniera *empirica* ovvero non originati da una sottostante distribuzione di probabilità.

Si supponga che x ed y siano variabili reali per una popolazione, e che $((x_1, y_1), (x_2, y_2), \dots, (x_n, y_n))$ è un campione osservato da (x, y) di dimensione n . Sia $\mathbf{x} = (x_1, x_2, \dots, x_n)$ una campione tratto da x e $\mathbf{y} = (y_1, y_2, \dots, y_n)$ il campione di y . In questo paragrafo è interessante valutare la *misura di associazione* tra \mathbf{x} ed \mathbf{y} e l'identificazione della linea (o curva) che meglio rappresenti i dati di misura in termini di fitting.

Si richiamano le medie campionarie

$$m(x) = \frac{1}{n} \sum_{i=1}^n x_i, \quad m(y) = \frac{1}{n} \sum_{i=1}^n y_i \quad (61)$$

e le varianze campionarie

$$s^2(x) = \frac{1}{n-1} \sum_{i=1}^n [x_i - m(x)]^2, \quad s^2(y) = \frac{1}{n-1} \sum_{i=1}^n [y_i - m(y)]^2 \quad (62)$$

A.0.1 Scatterplots

Il primo passo da fare nell'analisi esplorativa dei dati è spesso quello di dar un senso visuale della relazione statistica tra le variabili attraverso i cosiddetti *scatterplots*

In particolare è di interesse capire se la nuvola di punti assume un trend lineare o curvilineo, l'informazione che sarebbe gradita estrarre è come una variabile x possa essere utilizzata per predire la variabile y .

A.0.2 Definizioni

Il passo successivo è quello di definire delle grandezze in grado di misurare l'associazione tra i dati x ed y . Si definisce *covarianza campionaria*

$$s(\mathbf{x}, \mathbf{y}) = \frac{1}{n-1} \sum_{i=1}^n [x_i - m(\mathbf{x})][y_i - m(\mathbf{y})] \quad (63)$$

Si noti che la covarianza campionaria è una media del prodotto delle deviazioni dei dati x ed y dalle rispettive medie. La *correlazione campionaria* è definita come

$$r(\mathbf{x}, \mathbf{y}) = \frac{s(\mathbf{x}, \mathbf{y})}{s(\mathbf{x})s(\mathbf{y})} \quad (64)$$

posto che i vettori dei dati non siano costanti, così che le deviazioni standard siano positive. La correlazione è una versione standardizzata della covarianza. In particolare la correlazione è una grandezza adimensionale dato che le unità di misura al numeratore e denominatore coincidono. Si noti inoltre che correlazione e covarianza hanno lo stesso segno: positivo, negativo o nullo. Nel primo caso si dice che x ed y sono *positivamente correlati*, nel secondo caso *negativamente correlati* e nell'ultimo caso x ed y sono detti *incorrelati*.

Per mostrare che la covarianza è una semplice misura di associazione si ricordi che il punto $(m(\mathbf{x}), m(\mathbf{y}))$ è una misura del centro dei dati bivariati. Infatti, se ciascun punto fosse la posizione di un'unità di massa, allora $(m(\mathbf{x}), m(\mathbf{y}))$ sarebbe il *centro di massa* definito in fisica. Le linee orizzontale e verticale attraverso questo centro dividono il piano in quattro quadranti. Il prodotto $[x_i - m(\mathbf{x})][y_i - m(\mathbf{y})]$ è positivo nel primo e nel terzo quadrante e negativo nei restanti due. Con lo studio della *regressione lineare* che ci si accinge a trattare, si comprenderà in senso più profondo che cosa la covarianza misura.

In un primo momento è possibile rimanere perplessi sul fatto che le deviazioni sono mediate dal fattore $1/(n-1)$ piuttosto che da $1/n$. La migliore spiegazione è che nella statistica bivariata la covarianza campionaria è uno stimatore non polarizzato (altrimenti detto non corretto) della covarianza. Ad ogni una spiegazione del modo in cui si è mediato si può derivare in termini di *gradi di libertà* come è stato fatto per la varianza campionaria. Inizialmente si hanno $2n$ gradi di libertà nei dati bivariati. Ne vengono persi 2 realizzando le medie campionarie $m(\mathbf{x})$ ed $m(\mathbf{y})$. Dei rimanenti $2n-2$ gradi di libertà ne vengono persi $n-1$ calcolando i prodotti delle deviazioni. In questo modo ne rimangono un totale di $n-1$. Tipico della statistica è mediare dividendo non per il numero di termini della sommatoria ma piuttosto per il numero di gradi di libertà di tali termini. Ad ogni modo da un punto di vista puramente qualitativo sarebbe stato ragionevole dividere anche per n .

Come noto, associata ai dati vi è sempre una distribuzione di probabilità naturale detta distribuzione empirica, che dà probabilità $1/n$ a ciascun punto (x_i, y_i) . (Perciò se tali punti sono distinti questa coincide con la distribuzione uniforme discreta dei dati.) Le medie campionarie risultano semplicemente i valori attesi di questa distribuzione

bivariata e le variante campionarie sono semplicemente le varianze delle distribuzione bivariata eccezion fatta per la costante moltiplicativa (divisione per $1/(n-1)$ anziché $1/n$). Similmente, ad eccezione della stessa costante moltiplicativa, anche la covarianza campionaria è la covarianza della distribuzione bivariata e la correlazione campionaria è la correlazione della distribuzione bivariata. Tutti i risultati presentati di qui in avanti circa la statistica descrittiva sono casi particolari di risultati più generali sulle distribuzioni di probabilità.

A.o.3 Proprietà della covarianza

Si illustrano di seguito proprietà fondamentali della covarianza. I simboli in grassetto denotano campioni di dimensione n , ossia vettori costituiti da n variabili reali.

Proposizione 1. *Detto $\mathbf{xy} = (x_1y_1, x_2y_2, \dots, x_ny_n)$ la covarianza campionaria può essere ottenuta come*

$$s(\mathbf{x}, \mathbf{y}) = \frac{1}{n-1} \sum_{i=1}^n x_i y_i - \frac{n}{n-1} m(\mathbf{x}) m(\mathbf{y}) = \frac{n}{n-1} [m(\mathbf{xy}) - m(\mathbf{x}) m(\mathbf{y})] \quad (65)$$

Dimostrazione.

$$\begin{aligned} \sum_{i=1}^n [x_i - m(\mathbf{x})][y_i - m(\mathbf{y})] &= \sum_{i=1}^n [x_i y_i - x_i m(\mathbf{y}) - y_i m(\mathbf{x}) + m(\mathbf{x}) m(\mathbf{y})] \\ &= \sum_{i=1}^n x_i y_i - m(\mathbf{y}) \sum_{i=1}^n x_i - m(\mathbf{x}) \sum_{i=1}^n y_i + n m(\mathbf{x}) m(\mathbf{y}) \\ &= \sum_{i=1}^n x_i y_i - n m(\mathbf{y}) m(\mathbf{x}) - n m(\mathbf{x}) m(\mathbf{y}) + n m(\mathbf{x}) m(\mathbf{y}) \\ &= \sum_{i=1}^n x_i y_i + n m(\mathbf{x}) m(\mathbf{y}) \end{aligned}$$

□

Un'altra espressione utile per la covarianza campionaria è espresso di seguito

Proposizione 2. *La covarianza campionaria può essere espressa come*

$$s(\mathbf{x}, \mathbf{y}) = \frac{1}{2n(n-1)} \sum_{i=1}^n \sum_{j=1}^n (x_i - x_j)(y_i - y_j) \quad (66)$$

Dimostrazione.

$$\begin{aligned} \sum_{i=1}^n \sum_{j=1}^n (x_i - x_j)(y_i - y_j) &= \frac{1}{2n} \sum_{i=1}^n \sum_{j=1}^n [x_i - m(\mathbf{x}) + m(\mathbf{x}) - x_j][y_i - m(\mathbf{y}) + m(\mathbf{y}) - y_j] \\ &= \sum_{i=1}^n \sum_{j=1}^n \left([x_i - m(\mathbf{x})][y_i - m(\mathbf{y})] + [x_i - m(\mathbf{x})][m(\mathbf{y}) - y_j] + \right. \\ &\quad \left. + [m(\mathbf{x}) - x_j][y_i - m(\mathbf{y})] + [m(\mathbf{x}) - x_j][m(\mathbf{y}) - y_j] \right) \end{aligned}$$

Il primo termine risulta

$$n \sum_{i=1}^n [x_i - m(\mathbf{x})][y_i - m(\mathbf{y})]$$

i secondi due termini sommano a zero. L'ultimo termine invece

$$n \sum_{i=1}^n [m(\mathbf{x}) - x_j][m(\mathbf{y}) - y_j] = n \sum_{i=1}^n [x_i - m(\mathbf{x})][y_i - m(\mathbf{y})]$$

dividendo l'intera somma per $2n(n-1)$ porta alla tesi. \square

Proposizione 3. *Varianza campionaria*

$$s(\mathbf{x}, \mathbf{x}) := s^2(\mathbf{x}) \quad (67)$$

Proposizione 4. *Simmetria della varianza*

$$s(\mathbf{x}, \mathbf{y}) = s(\mathbf{y}, \mathbf{x}) \quad (68)$$

Proposizione 5. *Linearità della covarianza rispetto al primo argomento* Se \mathbf{x}, \mathbf{y} e \mathbf{z} sono vettori di dati tratti da una popolazione di variabili x, y e z rispettivamente, sia c una costante, allora

1.

$$s(\mathbf{x} + \mathbf{y}, \mathbf{z}) = s(\mathbf{x}, \mathbf{z}) + s(\mathbf{y}, \mathbf{z}) \quad (69)$$

2.

$$s(c\mathbf{x}, \mathbf{y}) = cs(\mathbf{x}, \mathbf{y}) \quad (70)$$

Dimostrazione. 1. Ricordando che $m(\mathbf{x} + \mathbf{y}) = m(\mathbf{x}) + m(\mathbf{y})$

$$\begin{aligned} s(\mathbf{x} + \mathbf{y}, \mathbf{z}) &= \frac{1}{n-1} \sum_{i=1}^n [x_i + y_i - m(\mathbf{x} + \mathbf{y})][z_i - m(\mathbf{z})] \\ &= \frac{1}{n-1} \sum_{i=1}^n \left([x_i - m(\mathbf{x})] + [y_i - m(\mathbf{y})] \right) [z_i - m(\mathbf{z})] \\ &= \frac{1}{n-1} \sum_{i=1}^n [x_i - m(\mathbf{x})][z_i - m(\mathbf{z})] + \frac{1}{n-1} \sum_{i=1}^n [y_i - m(\mathbf{y})][z_i - m(\mathbf{z})] \\ &= s(\mathbf{x}, \mathbf{z}) + s(\mathbf{y}, \mathbf{z}) \end{aligned}$$

2.

$$\begin{aligned} s(\mathbf{cx}, \mathbf{y}) &= \frac{1}{n-1} \sum_{i=1}^n [cx_i - m(\mathbf{cx})][y_i - m(\mathbf{y})] \\ &= \frac{1}{n-1} \sum_{i=1}^n [cx_i - cm(\mathbf{x})][y_i - m(\mathbf{y})] = cs(\mathbf{x}, \mathbf{y}) \end{aligned}$$

□

Naturalmente si arriverebbe alle stesse conclusioni ragionando sul secondo argomento. Questo porta al risultato più generale:

Proposizione 6. *Bilinearità della covarianza*

Sia \mathbf{x}_i un vettore di dati ottenuti da una popolazione di variabili x_i , con $i \in \{1, 2, \dots, k\}$ e sia \mathbf{y}_j un vettore di dati ottenuti da una popolazione di variabili y_j , con $j \in \{1, 2, \dots, l\}$. Si supponga inoltre a_1, a_2, \dots, a_k e b_1, b_2, \dots, b_l siano costanti. Allora

$$s\left(\sum_{i=1}^n a_i \mathbf{x}_i, \sum_{j=1}^n b_j \mathbf{y}_j\right) = \sum_{i=1}^n \sum_{j=1}^n a_i b_j s(\mathbf{x}_i, \mathbf{y}_j) \quad (71)$$

Proposizione 7. *La proprietà di bilinearità offre un modo interessante per calcolare la varianza della somma:*

$$s^2(\mathbf{x} + \mathbf{y}) = s^2(\mathbf{x}) + 2s(\mathbf{x}, \mathbf{y}) + s^2(\mathbf{y}) \quad (72)$$

Dimostrazione. Utilizzano le proposizioni 3, 4 e 5

$$\begin{aligned} s^2(\mathbf{x} + \mathbf{y}) &= s(\mathbf{x} + \mathbf{y}, \mathbf{x} + \mathbf{y}) = s(\mathbf{x}, \mathbf{x}) + s(\mathbf{x}, \mathbf{y}) + s(\mathbf{y}, \mathbf{x}) + s(\mathbf{y}, \mathbf{y}) \\ &= s^2(\mathbf{x}) + 2s(\mathbf{x}, \mathbf{y}) + s^2(\mathbf{y}) \end{aligned} \quad (73)$$

□

Il risultato è generalizzabile dicendo che la varianza campionaria della somma è la somma di tutte le coppie di covarianze campionarie. Si noti che la varianza campionaria della somma può essere a priori maggiore, minore o pari alla somma di tutte le varianze campionarie. In particolare, se i vettori sono a due a due incorrelati allora la varianza della somma coincide con la somma delle varianze.

Proposizione 8. *Sia \mathbf{c} un set di dati costanti allora*

$$s(\mathbf{x}, \mathbf{c}) = 0 \quad (74)$$

Dimostrazione. Segue direttamente dalla definizione. Se $c_i = c \forall i$, allora $m(\mathbf{c}) = c \Rightarrow c_i - m(\mathbf{c}) = 0 \forall i$

□

Combinando questi risultati si ottiene che la covarianza non cambia se si sommano costanti ai dati:

$$s(\mathbf{x} + \mathbf{c}, \mathbf{y} + \mathbf{d}) = s(\mathbf{x}, \mathbf{y}) \quad (75)$$

A.o.4 Proprietà della correlazione

Molte proprietà della correlazione sono corrispondenti a quelle della covarianza. Si richiamano le *standard scores*, grandezze adimensionali a media nulla e varianza unitaria:

$$u_i = \frac{x_i - m(\mathbf{x})}{s(\mathbf{x})}, \quad v_i = \frac{y_i - m(\mathbf{y})}{s(\mathbf{y})} \quad (76)$$

Proposizione 9. *La correlazione tra \mathbf{x} ed \mathbf{y} è la covarianza delle loro standard score \mathbf{u} e \mathbf{v} .*

Dimostrazione. Si noti che

$$\mathbf{u} = \frac{1}{s(\mathbf{x})}[\mathbf{x} - m(\mathbf{x})], \quad \mathbf{v} = \frac{1}{s(\mathbf{y})}[\mathbf{y} - m(\mathbf{y})] \quad (77)$$

Perciò il risultato discende direttamente dalle proposizioni 5 e 8:

$$s(\mathbf{u}, \mathbf{v}) = \frac{1}{s(\mathbf{x})s(\mathbf{y})}s(\mathbf{x}, \mathbf{y}) =: r(\mathbf{x}, \mathbf{y}) \quad (78)$$

□

A.o.5 Regressione lineare

L'interesse di questo paragrafo è quello di cercare la retta $y = a + bx$ che meglio si adatti, dall'inglese *fitting*, a rappresentare il campione di punti $((x_1, y_1), (x_2, y_2), \dots, (x_n, y_n))$. Questo è un problema fondamentale di molte branche della matematica, non solo della statistica. Si pensi ad \mathbf{x} come variabile predittrice ed a \mathbf{y} come variabile di risposta. In questo modo si vogliono trovare i coefficienti a e b , quindi una linea, che minimizzi la media degli errori quadratici tra i valori y nei dati e i valori y di predizione:

$$\text{mse}(a, b) = \frac{1}{n-1} \sum_{i=1}^n [y_i - (a + bx_i)]^2 \quad (79)$$

Si noti che nell'operare la minimizzazione della funzione convessa sopra definita si otterrebbe come punto (a, b) di minimo globale esattamente lo stesso punto risultante dalla minimizzazione di una funzione simile, ma anch'essa convessa. Utilizzando una terminologia propria della teoria dell'ottimizzazione il valore (a, b) ottimo è invariante alla valutazione di problemi equivalenti. L'espressione scelta per l'mse rappresenta tuttavia la più conveniente da un punto di vista statistico.

Proposizione 10. *Il grafico dell'mse è un paraboloide a concavità rivolta verso l'alto. La funzione raggiunge il punto di minimo globale quando*

$$b(\mathbf{x}, \mathbf{y}) = \frac{s(\mathbf{x}, \mathbf{y})}{s^2(\mathbf{x})}$$

$$a(\mathbf{x}, \mathbf{y}) = m(\mathbf{y}) - b(\mathbf{x}, \mathbf{y})m(\mathbf{x}) = m(\mathbf{y}) - \frac{s(\mathbf{x}, \mathbf{y})}{s^2(\mathbf{x})}m(\mathbf{x})$$

Dimostrazione. L'espressione algebrica dell'mse rivela da sola la natura di paraboloide con concavità rivolta verso l'alto. Per trovare l'unico punto che minimizzi l'mse si noti che

$$\begin{aligned}\frac{\partial}{\partial a} \text{mse}(a, b) &= \frac{1}{n-1} \sum 2[y_i - (a + bx_i)](-1) \\ &= \frac{2}{n-1} \left[-\sum_{i=1}^n y_i + na + b \sum_{i=1}^n x_i \right]\end{aligned}\quad (80)$$

$$\begin{aligned}\frac{\partial}{\partial b} \text{mse}(a, b) &= \frac{1}{n-1} \sum 2[y_i - (a + bx_i)](-x_i) \\ &= \frac{2}{n-1} \left[-\sum_{i=1}^n x_i y_i + a \sum_{i=1}^n x_i + b \sum_{i=1}^n x_i^2 \right]\end{aligned}\quad (81)$$

Risolvere $\frac{\partial}{\partial a} \text{mse}(a, b) = 0$, porta a $a = m(\mathbf{y}) - bm(\mathbf{x})$. Sostituendo il valore in $\frac{\partial}{\partial b} \text{mse}(a, b) = 0$ e risolvendo rispetto a b otteniamo

$$b = \frac{n[m(4\mathbf{xy}) - m(\mathbf{x})m(\mathbf{y})]}{n[m(\mathbf{x}^2) - m^2(\mathbf{x})]}$$

Dividendo numeratore e denominatore per $(n-1)$ e utilizzando la Proposizione 1 otteniamo $b = s(\mathbf{x}, \mathbf{y})/s^2(\mathbf{x})$.

□

Risulta chiaro che i valori di a e b sono statistici, ovvero funzioni dei dati. La linea di regressione campionaria risulta essere

$$\mathbf{y} = m(\mathbf{y}) + \frac{s(\mathbf{x}, \mathbf{y})}{s^2(\mathbf{x})}[\mathbf{x} - m(\mathbf{x})]$$

Si noti che la retta di regressione passa per il punto $(m(\mathbf{x}), m(\mathbf{y}))$, centro del campione di punti.

Proposizione 11. *L'errore quadrato medio risulta*

$$\text{mse}[a(\mathbf{x}, \mathbf{y}), b(\mathbf{x}, \mathbf{y})] = s^2(\mathbf{y})[1 - r^2(\mathbf{x}, \mathbf{y})]\quad (82)$$

Dimostrazione. Sostituendo $a(\mathbf{x}, \mathbf{y})$ e $b(\mathbf{x}, \mathbf{y})$ nell'espressione dell'mse e emplificando, si ottiene la tesi. □

Proposizione 12. *La correlazione campionaria soddisfa le seguenti proprietà:*

1. $-1 \leq r(\mathbf{x}, \mathbf{y}) \leq 1$
2. $-s(\mathbf{x})s(\mathbf{y}) \leq s(\mathbf{x}, \mathbf{y}) \leq s(\mathbf{x})s(\mathbf{y})$
3. $r(\mathbf{x}, \mathbf{y}) = -1$ se e solo se i punti del campione giacciono su una curva con concavità rivolta verso il basso

4. $r(\mathbf{x}, \mathbf{y}) = 1$ se e solo se i punti del campione giacciono su una curva con concavità rivolta verso l'alto

Dimostrazione. Si noti che $mse \geq 0$ perciò dalla Proposizione 11 deve essere $r^2(\mathbf{x}, \mathbf{y}) \leq 1$ e questo dimostra il punto 1. Dalla definizione di correlazione campionaria allora $1 \Rightarrow 2$. Per dimostrare le equivalenze 3 e 4 si noti infine che $mse(a, b) = 0$ se e solo se $y_i = a + bx_i \forall i$, e inoltre $b(\mathbf{x}, \mathbf{y})$ ha lo stesso segno di $r(\mathbf{x}, \mathbf{y})$ \square

Con questi strumenti siamo in grado di capire come la covarianza e la correlazione vadano a misurare il grado di linearità dei punti del campione.

Risulta evidente che la costante a che minimizza l'espressione

$$mse(a) = \frac{1}{n-1} \sum_{i=1}^n [y_i - a]^2 \quad (83)$$

coincide con la media campionaria $m(\mathbf{y})$ e il valore minimo dell' mse coincide con la varianza campionaria $s^2(\mathbf{y})$. La differenza tra questo valore e quello in (82), ovvero $s^2(\mathbf{y})r^2(\mathbf{x}, \mathbf{y})$, è la riduzione in variabilità nei dati \mathbf{y} quando il termine lineare in x è aggiunto al predittore. La riduzione frazionaria è $r^2(\mathbf{x}, \mathbf{y})$ e quindi questa statistica è chiamata *coefficiente di determinazione*. Si noti che se i vettori dei dati \mathbf{x} ed \mathbf{y} sono incorrelati, allora x non ha alcun valore come predatore di y : la retta di regressione in questo caso risulta essere la retta orizzontale $y = m(\mathbf{y})$ e l'errore quadrato medio è $s^2(\mathbf{y})$.

Le scelte del predatore e delle variabili di risposta sono importanti.

Proposizione 13. *La retta di regressione campionaria con variabile di predizione x e variabile di risposta y non è la stessa retta di regressione ottenuta con variabile di predizione y e variabile di risposta x , eccezion fatta per il caso limite $r(\mathbf{x}, \mathbf{y}) = \pm 1$ in cui i punti del campione giacciono tutti su una linea.*

A.o.6 Residui

La differenza tra il valore vero y di un dato e il valore predetto dalla retta di regressione è chiamato *residuo* del dato. Da questa definizione segue che il residuo corrispondente a (x_i, y_i) è $d_i = y_i - \hat{y}_i$ dove \hat{y}_i è il valore assunto dalla retta di regressione in x_i :

$$\hat{y}_i = m(\mathbf{y}) + \frac{s(\mathbf{x}, \mathbf{y})}{s^2(\mathbf{x})} [x_i - m(\mathbf{x})] \quad (84)$$

Si noti che il valore predetto \hat{y}_i e il residuo d_i sono *statistiche* ovvero funzioni dei dati (\mathbf{x}, \mathbf{y}) ; ciò non è stato evidenziato dalla notazione adottata per ovvie ragioni di chiarezza espositiva.

Proposizione 14. *La somma dei residui è zero:*

$$\sum_{i=1}^n d_i = 0 \quad (85)$$

Dimostrazione. Ciò segue immediatamente dalla definizione ed è una riformulazione del fatto che la retta di regressione passa per il centro del set dei dati $(m(\mathbf{x}), m(\mathbf{y}))$. \square

B

PRINCIPAL COMPONENT ANALYSIS

Con *Principal Component Analysis* (PCA) si intende una tecnica basata su sofisticati strumenti matematici utile alla trasformazione di un insieme di variabili correlate in un insieme più piccolo di variabili chiamate componenti principali. La tecnica trae origine dall'analisi di dati multivariati anche se è stata successivamente applicata a diversi altri campi.

In termini generali, la PCA realizza una trasformazione di spazi vettoriali per ridurre la dimensione di grandi insiemi di dati. Utilizzando proiezioni matematiche l'insieme di dati originale, interessato da diverse variabili, può essere interpretato attraverso poche variabili dette appunto componenti principali. Risulta facilmente intuibile come questa tecnica renda maggiormente accessibile l'analisi di dati multivariati.

B.1 FORMALIZZAZIONE DELLA PCA

Lo scopo principale della principal component analysis, in prima battuta, è quello di rispondere alla seguente domanda: *a partire da una base con la quale poter esprimere i dati è possibile ottenerne un'altra che sia combinazione lineare della base originale e che esprima i dati in modo altrettanto ottimale?* Posto in questo modo il problema potrebbe risultare ambiguo, si illustrerà di seguito una formulazione più accurata.

Si assuma di partire da un insieme di dati rappresentati da una matrice $\mathbf{X}_{m \times n}$, dove le n colonne rappresentano i campioni (cioè le osservazioni) e le m righe rappresentano le variabili. Si vuole ora compiere una trasformazione lineare su \mathbf{X} attraverso la matrice $\mathbf{P}_{m \times m}$:

$$\mathbf{Y}_{m \times n} = \mathbf{P}_{m \times m} \mathbf{X}_{m \times n} \quad (86)$$

Questa equazione rappresenta il cambio di base. Se si considera ciascuna riga di \mathbf{P} essere i vettori riga $\mathbf{p}_1, \mathbf{p}_2, \dots, \mathbf{p}_m$ e le colonne di \mathbf{X} i vettori colonna $\mathbf{x}_1, \mathbf{x}_2, \dots, \mathbf{x}_n$ allora la 86 può così essere interpretata:

$$\mathbf{PX} = (\mathbf{Px}_1 \quad \mathbf{Px}_2 \quad \dots \quad \mathbf{Px}_n) = \begin{pmatrix} \mathbf{p}_1\mathbf{x}_1 & \mathbf{p}_1\mathbf{x}_2 & \dots & \mathbf{p}_1\mathbf{x}_n \\ \mathbf{p}_2\mathbf{x}_1 & \mathbf{p}_2\mathbf{x}_2 & \dots & \mathbf{p}_2\mathbf{x}_n \\ \vdots & \vdots & \ddots & \vdots \\ \mathbf{p}_m\mathbf{x}_1 & \mathbf{p}_m\mathbf{x}_2 & \dots & \mathbf{p}_m\mathbf{x}_n \end{pmatrix} = \mathbf{Y} \quad (87)$$

Si noti che $\mathbf{p}_j, \mathbf{x}_j \in \mathbb{R}^m$ perciò $\mathbf{p}_j\mathbf{x}_j$ è il prodotto interno euclideo tra vettori. Questo conferma che i dati originali \mathbf{X} vengono proiettati sulle colonne di \mathbf{P} . In questo modo le righe $\{\mathbf{p}_1, \mathbf{p}_2, \dots, \mathbf{p}_m\}$ sono una nuova base per rappresentare le colonne di \mathbf{X} . Le righe di \mathbf{P} diventeranno le direzioni delle componenti principali.

Bisogna ora affrontare il problema di come deve essere la nuova base in termini di indipendenza tra le componenti principali della stessa.

La PCA definisce questa indipendenza considerando la varianza dei dati nella base originale. L'intento è quello di decorrelare i dati originali trovando le direzioni in cui la varianza è massimizzata e quindi usare queste direzioni per definire la nuova base. Si richiama la definizione di varianza per una variabile aleatoria Z di media μ :

$$\sigma_Z^2 = E[(Z - \mu)^2] \quad (88)$$

Si supponga di avere un vettore di n misurazioni discrete, $\tilde{\mathbf{r}} = (\tilde{r}_1, \tilde{r}_2, \dots, \tilde{r}_n)$ di media μ_r . Sottraendo la media da ciascuna misura si ottiene il set traslato $\mathbf{r} = (r_1, r_2, \dots, r_n)$ a media nulla. In questo modo la varianza di queste misure si può esprimere come:

$$\sigma_r^2 = \frac{1}{n} \mathbf{r}\mathbf{r}^T \quad (89)$$

Disponendo di un secondo vettore di n misure, $\mathbf{s} = (s_1, s_2, \dots, s_n)$ a media nulla, allora possiamo generalizzare il procedimento ed ottenere la *covarianza* tra \mathbf{r} ed \mathbf{s} :

$$\sigma_{rs}^2 = \frac{1}{n-1} \mathbf{r}\mathbf{s}^T \quad (90)$$

Estendiamo ora questa idea al caso della matrice di dati $\mathbf{X}_{m \times n}$. Quest'ultima si esprime come

$$\mathbf{X} = \begin{pmatrix} x_{1,1} & x_{1,2} & \dots & x_{1,n} \\ x_{2,1} & x_{2,2} & \dots & x_{2,n} \\ \vdots & \vdots & \ddots & \vdots \\ x_{m,1} & x_{m,2} & \dots & x_{m,n} \end{pmatrix} = \begin{pmatrix} \mathbf{x}_1 \\ \mathbf{x}_2 \\ \vdots \\ \mathbf{x}_m \end{pmatrix} \in \mathbb{R}^{m \times n}, \quad \mathbf{x}_i^T \in \mathbb{R}^n \quad (91)$$

Dato che vi è un vettore riga per ogni variabile ciascun vettore riga contiene tutti i campioni per quella data variabile. Ad esempio \mathbf{x}_i è il

vettore degli n campioni per la i -esima variabile. Per questo motivo è lecito considerare il seguente prodotto:

$$\mathbf{C}_X = \frac{1}{n-1} \mathbf{X}\mathbf{X}^T = \frac{1}{n-1} \begin{pmatrix} \mathbf{x}_1\mathbf{x}_1^T & \mathbf{x}_1\mathbf{x}_2^T & \dots & \mathbf{x}_1\mathbf{x}_m^T \\ \mathbf{x}_2\mathbf{x}_1^T & \mathbf{x}_2\mathbf{x}_2^T & \dots & \mathbf{x}_2\mathbf{x}_m^T \\ \vdots & \vdots & \ddots & \vdots \\ \mathbf{x}_m\mathbf{x}_1^T & \mathbf{x}_m\mathbf{x}_2^T & \dots & \mathbf{x}_m\mathbf{x}_m^T \end{pmatrix} \in \mathbb{R}^{m \times m} \quad (92)$$

In questa matrice appaiono tutte le possibili coppie di covarianza tra le m variabili, in particolare le varianze giacciono sulla diagonale. Questa matrice è detta *matrice di covarianza*.

Tornando al problema originale della trasformazione lineare $\mathbf{Y} = \mathbf{P}\mathbf{X}$ per qualche \mathbf{P} , è necessario capire quali caratteristiche vogliamo che \mathbf{Y} esibisca e in che modo queste si riflettano sulla matrice di covarianza \mathbf{C}_Y . Dato che la PCA parte dal presupposto che le variabili trasformate debbano essere il più incorrelate possibile è necessario che nella matrice di covarianza \mathbf{C}_Y gli elementi afferenti a coppie di variabili diverse siano quanto più vicini al valore nullo, fermo restando che le matrici di covarianza sono oggetti definiti positivi o al più semidefiniti positivi. Dualmente siamo interessati ad elevati valori delle varianze sinonimo di una dinamica apprezzabile del sistema. La matrice di covarianza \mathbf{C}_Y va costruita perciò secondo due direttive principali:

1. Massimizzare il segnale misurato dalla varianza (massimizzare i valori sulla diagonale)
2. Minimizzare la covarianza tra variabili diverse (minimizzare i restanti valori)

Cerchiamo perciò una matrice \mathbf{P} tale che la matrice di covarianza \mathbf{C}_Y sia *diagonale*.

Non comporta alcuna restrizione assumere ora che i vettori della nuova base $\mathbf{p}_1, \mathbf{p}_2, \dots, \mathbf{p}_m$ siano ortogonali. Questo ci permette di trovare una soluzione al problema utilizzando gli strumenti dell'algebra lineare. Si consideri la formula per la covarianza \mathbf{C}_Y e la nostra interpretazione di \mathbf{Y} in termini di \mathbf{X} e \mathbf{P} .

$$\begin{aligned}
\mathbf{C}_Y &= \frac{1}{n-1} \mathbf{Y}\mathbf{Y}^T \\
&= \frac{1}{n-1} (\mathbf{P}\mathbf{X})(\mathbf{P}\mathbf{X})^T \\
&= \frac{1}{n-1} (\mathbf{P}\mathbf{X})(\mathbf{X}^T \mathbf{P}^T) \\
&= \frac{1}{n-1} \mathbf{P}(\mathbf{X}\mathbf{X}^T) \mathbf{P}^T \\
&= \frac{1}{n-1} \mathbf{P}\mathbf{S}\mathbf{P}^T \quad \text{dove } \mathbf{S} = \mathbf{X}\mathbf{X}^T \tag{93}
\end{aligned}$$

Con \mathbf{S} è indicata la matrice simmetrica t.c. $(\mathbf{X}\mathbf{X}^T)^T = (\mathbf{X}^T)^T \mathbf{X}^T = \mathbf{X}\mathbf{X}^T$. Un teorema dell'algebra lineare afferma che ciascuna matrice quadrata simmetrica è diagonalizzabile ortogonalmente:

$$\mathbf{S} = \mathbf{E}\mathbf{D}\mathbf{E}^T \tag{94}$$

dove $\mathbf{E}_{m \times m}$ matrice ortonormale le cui colonne sono gli autovettori ortonormali ad \mathbf{S} , e \mathbf{D} è una matrice diagonale i cui elementi sono gli autovalori di \mathbf{S} . Il rango r di \mathbf{S} è il numero degli autovettori ortonormali in essa presenti. Se \mathbf{B} fosse t.c. $r < m$ sarà semplicemente necessario generare $m - r$ vettori ortonormali per riempire le rimanenti colonne di \mathbf{S} .

A questo punto si fa la scelta della matrice di trasformazione \mathbf{P} in modo tale che le sue righe siano gli autovettori di \mathbf{S} così da assicurare che $\mathbf{P} = \mathbf{E}^T$ e viceversa. Fatto ciò, sostituendo nell'espressione derivata in (93) otteniamo

$$\begin{aligned}
\mathbf{C}_Y &= \frac{1}{n-1} \mathbf{P}\mathbf{S}\mathbf{P}^T \\
&= \frac{1}{n-1} \mathbf{E}^T (\mathbf{E}\mathbf{D}\mathbf{E}^T) \mathbf{E} \tag{95}
\end{aligned}$$

Ora, dato che \mathbf{E} è una matrice ortonormale si ha $\mathbf{E}^T \mathbf{E} = \mathbf{I}_{m \times m}$ dove \mathbf{I} è la matrice identità. Dunque per questa particolare scelta di \mathbf{P} si avrà:

$$\mathbf{C}_Y = \frac{1}{n-1} \mathbf{D} \tag{96}$$

Un'ultima nota è spesa per sottolineare che questo metodo estrae automaticamente informazione dalle varianze sull'importanza relativa di ciascuna componente principale. La varianza più grande corrisponde alla prima componente principale, la seconda varianza più grande corrisponde alla seconda componente principale e così via. In questo modo si è stabilito un metodo per organizzare i dati nello stadio di diagonalizzazione. Una volta ottenuti gli autovettori di

$\mathbf{S} = \mathbf{X}\mathbf{X}^T$ si organizzano gli autovalori in ordine decrescente e si posizionano in questo ordine sulla diagonale di \mathbf{D} . Quindi si costruisce la matrice ortonormale \mathbf{E} posizionando nello stesso ordine gli autovettori associati e formando così le colonne di \mathbf{E} (questo vuol dire che si posiziona l'autovettore che corrisponde all'autovalore più grande nella prima colonna, l'autovettore che corrisponde al secondo autovalore più grande nella seconda colonna e così via).

In questo modo è stato raggiunto l'obiettivo di diagonalizzare la matrice di covarianza dei dati avendone fatto una trasformazione preventiva. Le componenti principali, le righe di \mathbf{P} , sono gli autovettori della matrice di covarianza $\mathbf{X}\mathbf{X}^T$ e le righe appaiono in ordine di 'importanza', rivelando quanto 'principale' una componente effettivamente sia.

BIBLIOGRAFIA

- Amonlirdviman, Keith, Narmada A. Khare, David R.P. Tree, Wei-Shen Chen, Jeffrey D. Axelrod e Claire J. Tomlin
2005 *Mathematical modeling of planar cell polarity to understand dominating non-autonomys*, Stanford University. (Citato a p. 43.)
- Blake, Andrew e Michael Isard
Active Contour, Springer. (Citato a p. 17.)
- Blanchard, Guy B., Alexandre J. Kabla, Nora L. Schultz, Lucy C. Butler, Benedicte Sanson, Nicole Gorfinkel, L. Mahavedan e Richard J. Adams
2009 *Tissue tectonics: morphogenetic strain rates, cell shape change and intercalation*, University of Cambridge, Harvard University, Harvard Medical School.
- Butler, Lucy C., Guy B. Blanchard, Alexandre J. Kabla, Nicola J. Lawrence, David P. Welchman, L. Mahavedan, Richard J. Adams e Benedicte Sanson
2009 *Cell shape changes indicate a role for extrinsic tensile forces in Drosophila germ-band extension*, University of Cambridge, Ecole Polytechnique Federale de Lausanne, Harvard University, Harvard Medical School.
- Cenedese, Angelo e Alessandro Beghi
2006 *Optimal Approach to Shape Parameter Control*, University of Padova, Bali, Indonesia.
- Cenedese, Angelo e Alberto Silletti
2009 *A robust active contour approach for studying cell deformation from noisy images*, Swansea, UK.
- Gibson, Matthew C., Ankit B. Patel, Radhika Nagpal e Norbert Perrimon
2006 «The Emergence of Geometric Order in Proliferating Metazoon Epithelia», *Nature*. (Citato a p. 64.)
2007 *Supplementary Information for "The Emergence of Geometric Order in Proliferating Metazon Epithelia"*. (Citato a p. 64.)
- Kass, M, A. Witkin e Demetri Terzopoulos
1988 *Snakes: Active contour models*, 8(2):321–331. (Citato alle p. 18, 21.)

- Kindratenko, Volodymyr V.
2003 *On Using Functions to Describe the Shape*, 8: 225–245. (Citato a p. 32.)
- Lauffenburger, Douglas A. e Alan F. Horwitz
1996 *Cell Migration: A Physically Integrated Molecular Process*, Massachusetts Institute of Technology, University of Illinois at Urbana–Champaign, cap. Cell, Vol. 84, pages.
- Maeda, Yusuke T., Junya Inose, Miki Y. Matsuo, Suguru Iwaya e Masaki Sano
2008 *Ordered Patterns of Cell Shape and Orientational Correlation during Spontaneous Cell Migration*, University of Tokyo.
- McInerney, Tim e Demetri Terzopoulos
1996 *Deformable Models in Medical Image Analysis: A Survey*, 1(2):91–108.
- Picci, Giorgio
2007 *Filtraggio statistic (Wiener, Levinson, Kalman) e applicazioni*, Edizioni Libreria Progetto Padova.
- Rosin, Paul L.
2003 *Measuring shape: ellipticity, rectangularity, and triangularity*, Georgia Institute of Technology Atlanta, University of California Los Angeles, cap. Computer Vision 53(2), pages.
- Silletti, Alberto
2007-2009 *Dynamic shape detection and analysis of deformable structures in biomedical imaging*, doctoral thesis, Università degli studi di Padova. (Citato alle p. xiii, 21, 32, 39, 44.)
- Silletti, Alberto, Alessandro Abate, Jeffrey D. Axelrod e Claire J. Tomlin
2011 *Versatile spectral methods for point set matching*, University of Padova, Delft University of Technology, Stanford University School of Medicine, University of California at Berkeley, cap. Pattern Recognition Letters, Volume 32, Issue 5, Pages 731–739, journal%20homepage:%20www.elsevier.com/locate/patrec.
- Silletti, Alberto, Angelo Cenedese e Alessandro Abate
2009 *The emergent structure of the Drosophila wing. A dynamical Model Generator*, University of Padova, Stanford University.
- Sonka Milan, Hlavac Vaclav e Boyle Roger
2008 *Image Processing, Analysis, and Machine Vision*, vol. 3, Tomson.
- Stegmann, Mikkel B. e David Delgado Gomez
2002 *A Brief Introduction to Statistical Shape Analysis*, report, Informatics e Mathematical Modelling, Technical University of Denmark.

Terzopoulos, Demetri e Fleischer Kurt

1988 *Deformable Models*, *The Visual Computers* 4:306-331. (Citato a p. 21.)

Yezzi, Anthony J. e Stefano Soatto

2003 *Deformation: Deforming Motion Shape Average and the Joint Registration and Approximation of Structure in Images*, Cardiff University, cap. *Machine Vision and Applications*, 14, pages.