A Specialized Inter-Curve Similarity Measure for Agglomerative Diffusion MRI Streamline Clustering

Jadrian Miles (jadrian@cs.brown.edu) Visualization Research Lab Computer Science Department Brown University, Providence, RI Talk at NIH, 2009-05-19

Overview

Motivation

- Agglomerative clustering
- Inter-curve similarity measures
- Efficient clustering implementation

Overview

Motivation

- Agglomerative clustering
- Inter-curve similarity measures
- Efficient clustering implementation

What are we dealing with?



What are we dealing with?

- DTI dataset (Imm isotropic voxels)
- Principle eigenvector streamlines
- 323k space curves
- Lots of short / broken segments

What are we dealing with?



Why clustering?

- Approximates tracts for:
 - Visualization and selection
 - Statistics
 - Non-local modeling



The Goal

- Automatic clustering
- Dense whole-brain tractogram
- Approximate anatomical tracts

The Hypothesis

Short segments can aid clustering

Overview

Motivation

- Agglomerative clustering
- Inter-curve similarity measures
- Efficient clustering implementation

Overview

Motivation

- Agglomerative clustering
- Inter-curve similarity measures
- Efficient clustering implementation

How does it work?

- I. Every curve begins as its own cluster
- 2. Join clusters of two closest curves
- 3. Repeat step 2 until termination

How does it work?

- Single-link distance
- $A \sim_{\varepsilon} B \Leftrightarrow \exists \{X_1, X_2, ..., X_k\} \text{ s.t. } d(A, X_1) < \varepsilon,$ $d(X_k, B) < \varepsilon, d(X_i, X_{i+1}) < \varepsilon \forall I \le i < k$

Why agglomerative?

- Provides a cluster hierarchy
- Easy to understand
- Parallelizable
- Flexible termination conditions

Overview

Motivation

- Agglomerative clustering
- Inter-curve similarity measures
- Efficient clustering implementation

Overview

Motivation

- Agglomerative clustering
- Inter-curve similarity measures
- Efficient clustering implementation

Why a new measure?

- Previous work:
 - Corouge
 - Brun
 - Zhang
 - Ding

• Corouge, et al., ISBI 2004

- Minimum-distance vertex pairs
- Minimum, mean, max of pair distances

• Brun, et al., EUROCAST 2003

• Endpoint distance

Zhang & Laidlaw, TVCG 2008

Mean distance above threshold

Zhang & Laidlaw, TVCG 2008

 Mean distance above threshold ... from shorter curve

• Zhang & Laidlaw, TVCG 2008

- Mean distance above threshold ... from shorter curve
- But!

- Ding, et al., Vis 2001
 - Corresponding segment
 - Mean vertex distance in segment

- Ding, et al., Vis 2001
 - Corresponding segment
 - Mean vertex distance in segment

Detects overlapping segment
Penalizes skewed curves

• Trimming

- Induced relative orientation
- Matching order
- Match length ratio

Trimming



Trimming



What we want



What we (can) get

Trimming Excess Matches

Trimming Excess Matches

d(A,B) = mean distance in trimmed segment

• Trimming

- Induced relative orientation
- Matching order
- Match length ratio

Relative Orientation

Relative Orientation

A Good Match

Not So Good

Forcing Relative Orientation

Forcing Relative Orientation

d(A,B) = mean distance in trimmed segment, with forced relative orientation

• Trimming

- Induced relative orientation
- Matching order
- Match length ratio

Without Orientation Forcing or Trimming

With Forced Orientation: Whoa!

Matching Order

Matching Order

Matching Order

All Matching Orders

All Matching Orders

d(A,B) = min mean distance in trimmed segment with forced orientation, over all four matching orders

• Trimming

- Induced relative orientation
- Matching order
- Match length ratio

Success!

Success!

Still Not Quite Right

An Extreme Case

Ъŋ Г

An Extreme Case

An Extreme Case

Match Length Ratio

Infinite for perpendicular segments
Infinite for end-to-end segments

- d(A,B) = minimum over all four matching orders of:
 - Mean distance between points
 - In trimmed segment
 - With forced relative orientation
 - Multiplied by the match length ratio

Detects overlapping segment
Penalizes skewed curves

• Exploits agglomerative clustering with singlelink cluster distance

A

B

- Exploits agglomerative clustering with singlelink cluster distance
- $d(A,C) \le \epsilon, d(B,C) \le \epsilon$
- $d(A,B) = \infty$ and yet $A \sim_{\epsilon} B$

A Failure Case

- A single short segment can join two otherwise strongly separated clusters
- Post-processing the clustering is likely necessary

Overview

Motivation

- Agglomerative clustering
- Inter-curve similarity measures
- Efficient clustering implementation

Overview

Motivation

- Agglomerative clustering
- Inter-curve similarity measures
- Efficient clustering implementation

What We'd Like

- I. Compute all pairwise curve similarities
- 2. Find highest similarity
- 3. Merge those clusters
- 4. GOTO 2

What We'd Like

- Can compute 1600 similarities / s
- But we have 300k curves!
 - 45 billion similarities, 400 GB, 11.5 months
- Even decimated, 30k curves
 - 450 million similarities, 5.1 GB, 83 hours

What we need

- Parallelism
- Only the highest elements of the similarity matrix in memory

Efficient Clustering

Efficient Clustering

- 30k curves with 13 hosts
- 6.5 hours, 5.1 GB total
- I0 priority queues of increasing size: 8MB, 22MB, 56MB, I30MB... I.IGB
- Only one priority queue must be in memory at a time to do clustering

Thanks!

David Laidlaw, my advisor
Evren Ozarslan for inviting me

Questions?

Motivation

- Agglomerative clustering
- Inter-curve similarity measures
- Efficient clustering implementation