

Interface 2008

Disparate Information Fusion

Carey E. Priebe

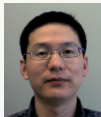
Department of Applied Mathematics & Statistics
Johns Hopkins University

May 21 – 24, 2008
Durham, NC

Interface 2008

Disparate Information Fusion

Carey E. Priebe

Department of Applied Mathematics & Statistics
Johns Hopkins UniversityAdam
Cardinal-
StakenasZhiliang
MaYoungser
ParkMay 21 – 24, 2008
Durham, NC

Introduction

Disparate Data Types

Dissimilarities

Approaches

Motivation

Disparate Information Fusion Framework

Examples

Image/Caption Fusion

Brain Shape Comparison

Streaming Content in Context

Disparate Data Types

- Text: emails, ships' logs, computer logs, ... ,
- Audio: speech, cell phone, other communications, ... ,
- Image: UAV, satellite, mutlispectral, hyperspectral, ... ,
- Fingerprint data,
- Facial/Character recognition data,
- Other abstract data.

Disparate Data Types

- Text: emails, ships' logs, computer logs, ... ,
- Audio: speech, cell phone, other communications, ... ,
- Image: UAV, satellite, mutlispectral, hyperspectral, ... ,
- Fingerprint data,
- Facial/Character recognition data,
- Other abstract data.

Goal:

To obtain superior performance in the exploitation tasks (inference, prediction, synthesis, and general mining) through the fusion of available disparate data types.

Dissimilarities

Setup:

$$(X_i, Y_i) \stackrel{i.i.d}{\sim} F_{XY} \text{ with } X_i : \Omega \rightarrow \Xi_1 \times \cdots \times \Xi_K, \text{ and} \\ Y_i : \Omega \rightarrow \{1, \dots, J\}.$$

Definition

A *dissimilarity measure* is a function $\delta : \Xi \times \Xi \rightarrow \mathbb{R}^+ \cup \{0\}$ with:

1. Positivity: $\delta(x_1, x_2) \geq 0$,
2. Symmetry: $\delta(x_1, x_2) = \delta(x_2, x_1)$,
3. Reflexivity: $\delta(x, x) = 0$,
4. Identifiability: $\delta(x_1, x_2) = 0 \Rightarrow x_1 = x_2$.

Dissimilarities

Setup:

$(X_i, Y_i) \stackrel{i.i.d.}{\sim} F_{XY}$ with $X_i : \Omega \rightarrow \Xi_1 \times \cdots \times \Xi_K$, and
 $Y_i : \Omega \rightarrow \{1, \dots, J\}$.

Definition

A *dissimilarity measure* is a function $\delta : \Xi \times \Xi \rightarrow \mathbb{R}^+ \cup \{0\}$ with:

1. Positivity: $\delta(x_1, x_2) \geq 0$,
2. Symmetry: $\delta(x_1, x_2) = \delta(x_2, x_1)$,
3. Reflexivity: $\delta(x, x) = 0$,
4. Identifiability: $\delta(x_1, x_2) = 0 \Rightarrow x_1 = x_2$.

A *dissimilarity representation* for a set of n objects is expressed as a symmetric, nonnegative and hollow matrix Δ .

Introduction

Disparate Data Types

Dissimilarities

Approaches

Motivation

Disparate Information Fusion Framework

Examples

Image/Caption Fusion

Brain Shape Comparison

Streaming Content in Context

Motivation

on the joint vs the product of marginals in disparate dissimilarity fusion

1. “The product of the marginals has less information than the joint.”

- We all understand that we can derive the marginals from the joint,
- but that we cannot (necessarily) recover the joint from the marginals.
- *So . . . how shall we fuse disparate dissimilarities?*

Motivation

on the joint vs the product of marginals in disparate dissimilarity fusion

2. A Toy Example:

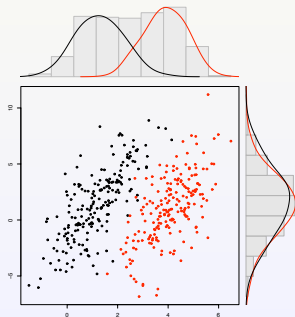
- Consider four objects: ①, ②, ③, ④.
- Consider three “marginal dissimilarities”: `size`, `shape`, `color`.
- Each **marginal** dissimilarities alone produces its own clusters:
 - `size` $\rightarrow \{1\}, \{2,3,4\}$;
 - `shape` $\rightarrow \{1,2\}, \{3,4\}$;
 - `color` $\rightarrow \{1,2,3\}, \{4\}$.
- With the **joint** information, we can distinguish all four objects.
- **Combining marginal information**, in this case, also provides full information.

Motivation

on the joint vs the product of marginals in disparate dissimilarity fusion

3. Tilted Parallel Cigars:

- Neither marginal yields much discriminatory power, but the **joint** does.
- If the two canonical Euclidean marginal dissimilarities are considered, the joint information can be recovered.
- This is not, in general, the case.

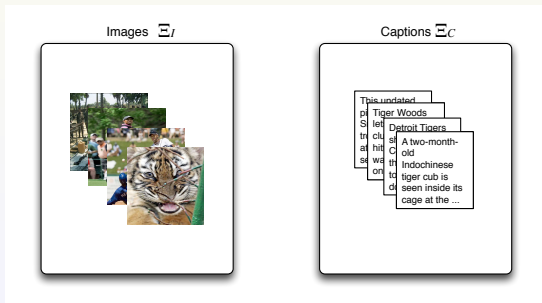


Motivation

on the joint vs the product of marginals in disparate dissimilarity fusion

4. Image/Caption Fusion:

- Images and text are most naturally considered **marginally**,
- Our fusion methodologies allow for subsequent **“joint”** analyses.

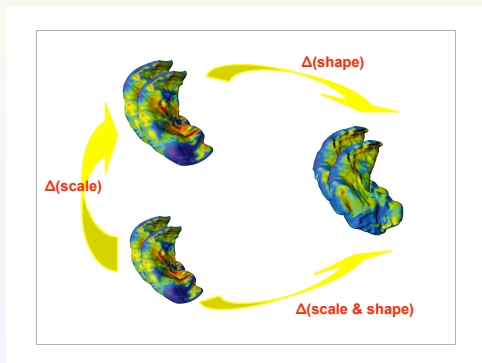


Motivation

on the joint vs the product of marginals in disparate dissimilarity fusion

5. Brain Shape Comparison:

- left & right hippocampi
- scale & shape



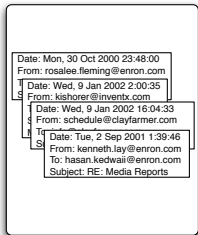
Motivation

on the joint vs the product of marginals in disparate dissimilarity fusion

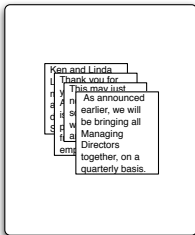
6. Streaming Content in Context:

- Here we have T_E and T_C , and
- $T = g(T_E, T_C)$ provides superior inference than either T_E alone or T_C alone.

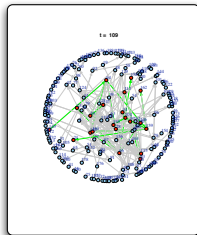
Externals Ξ_E



Content Ξ_C



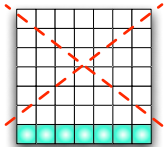
Graph Ξ_G



Framework

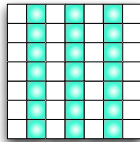
Three dissimilarity classification approaches

Dissimilarity Matrix ($n \times n$)



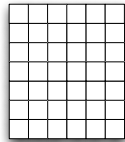
Neighborhood-based

Dissimilarity Matrix ($n \times n$)



Dissimilarity Space-based

Feature Matrix ($n \times d$)

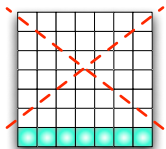


Embedding

Framework

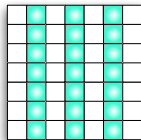
Three dissimilarity classification approaches

Dissimilarity Matrix ($n \times n$)



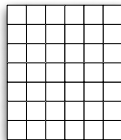
Neighborhood-based

Dissimilarity Matrix ($n \times n$)



Dissimilarity Space-based

Feature Matrix ($n \times d$)



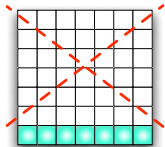
Embedding

1. The *neighborhood-based approach* interprets dissimilarities as neighborhood relations.

Framework

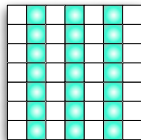
Three dissimilarity classification approaches

Dissimilarity Matrix ($n \times n$)



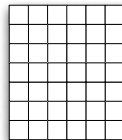
Neighborhood-based

Dissimilarity Matrix ($n \times n$)



Dissimilarity Space-based

Feature Matrix ($n \times d$)



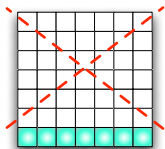
Embedding

1. The *neighborhood-based approach* interprets dissimilarities as neighborhood relations.
2. The *dissimilarity space approach* defines a representation set $R = \{p_1, \dots, p_r\}$, and interprets dissimilarities from a point to each element of the representation set as features of this point.

Framework

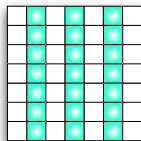
Three dissimilarity classification approaches

Dissimilarity Matrix ($n \times n$)



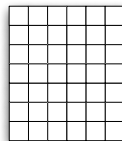
Neighborhood-based

Dissimilarity Matrix ($n \times n$)



Dissimilarity Space-based

Feature Matrix ($n \times d$)

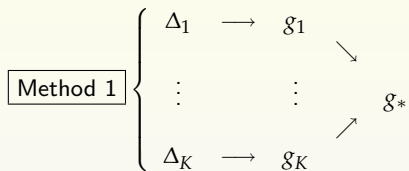


Embedding

1. The *neighborhood-based approach* interprets dissimilarities as neighborhood relations.
2. The *dissimilarity space approach* defines a representation set $R = \{p_1, \dots, p_r\}$, and interprets dissimilarities from a point to each element of the representation set as features of this point.
3. The *embedding approach* embeds dissimilarities into \mathbb{R}^d in such a way that the configuration's interpoint distances approximate the original dissimilarities.

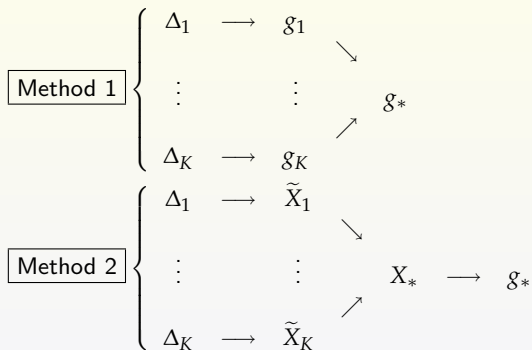
Framework

Three DIF approaches



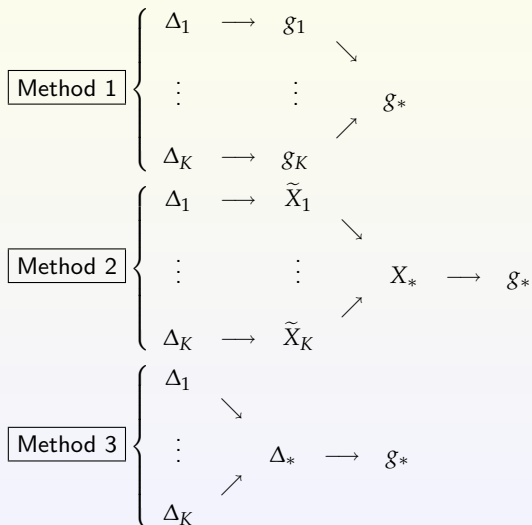
Framework

Three DIF approaches



Framework

Three DIF approaches



Introduction

Disparate Data Types

Dissimilarities

Approaches

Motivation

Disparate Information Fusion Framework

Examples

Image/Caption Fusion

Brain Shape Comparison

Streaming Content in Context

Image/Caption Fusion

- 140,577 *images & captions* were collected from Yahoo! Photos website [Jeff Solka and his team].

Image/Caption Fusion

- 140,577 *images & captions* were collected from Yahoo! Photos website [Jeff Solka and his team].
- 1,600 pairs were selected using query word “tiger” on captions.

Image/Caption Fusion

- 140,577 *images & captions* were collected from Yahoo! Photos website [Jeff Solka and his team].
- 1,600 pairs were selected using query word “tiger” on captions.
- They were labeled manually based only on captions:

label	#
animal tiger	148
Detroit Tigers baseball team	145
Tiger Woods the golfer	897
Tamil Tigers soldiers of Sri Lanka	330
Leicester Tigers rugby team	48
others	32

Image/Caption Fusion

- 140,577 *images & captions* were collected from Yahoo! Photos website [Jeff Solka and his team].
- 1,600 pairs were selected using query word “tiger” on captions.
- They were labeled manually based only on captions:

label	#
animal tiger	148
Detroit Tigers baseball team	145
Tiger Woods the golfer	897
Tamil Tigers soldiers of Sri Lanka	330
Leicester Tigers rugby team	48
others	32

- Two class problem: “Tiger Woods” and “Tamil Tigers”.

“Tiger” Dissimilarity Matrices

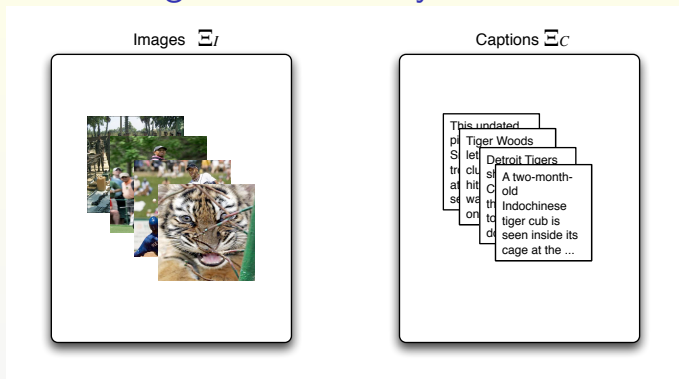
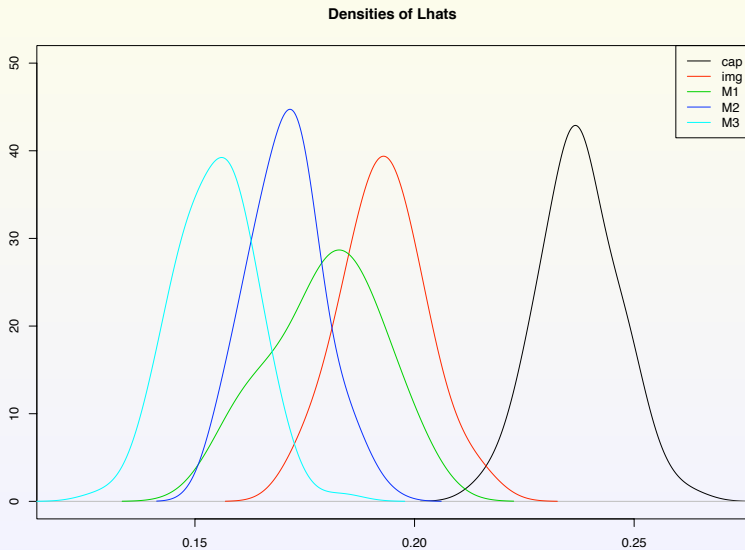


Figure: Conceptual depiction of the “tiger data set”

- *image* features: first & second order pixel derivatives [Jain].
- *caption* features: Mutual Information [Lin & Pantel].
- **dissimilarities**: $1 - (\text{Random Forest proximity})$ [Breiman].

Combining Dissimilarities

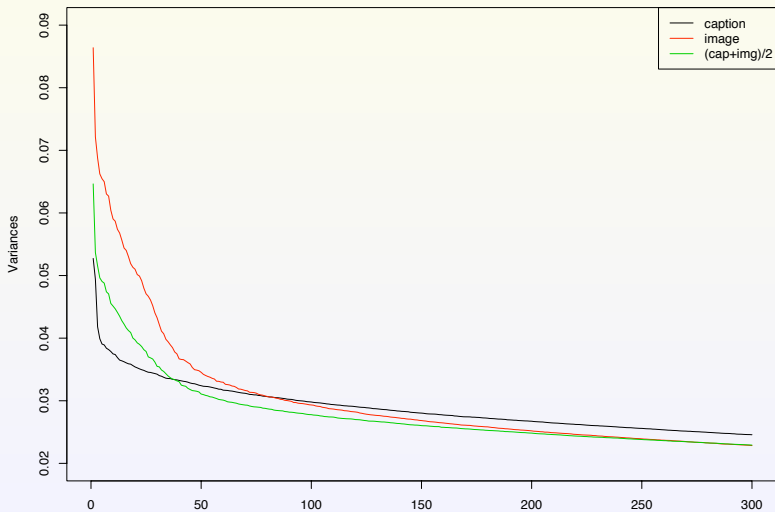
Methods 1, 2 & 3 compared with “image” only and “caption” only



Combining Dissimilarities

Determining the dimensions for exploitation task

Scree plot



Semisupervised Learning

Self-Training [Zhu2005]:

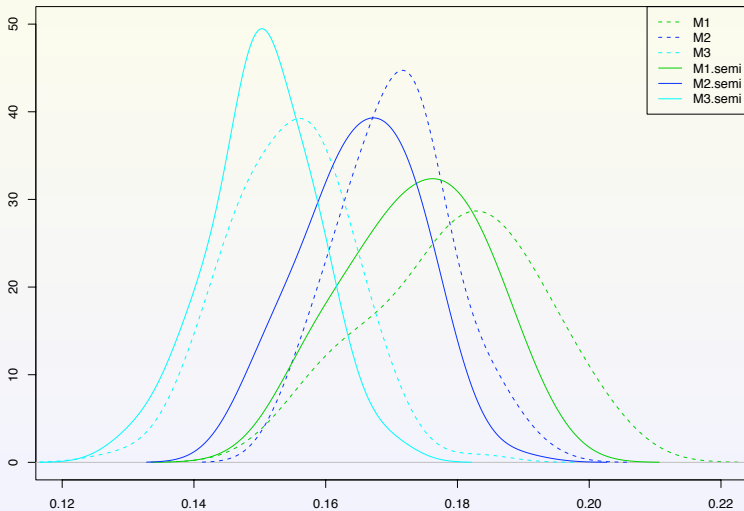
“Self-training is a commonly used technique for semi-supervised learning. In self-training a classifier is first trained with the small amount of labeled data. The classifier is then used to classify the unlabeled data. Typically the most confident unlabeled points, together with their predicted labels, are added to the training set.”

- X. Zhu. *Semi-supervised learning literature survey*. Technical Report 1530, Computer Sciences, University of Wisconsin-Madison, 2005.

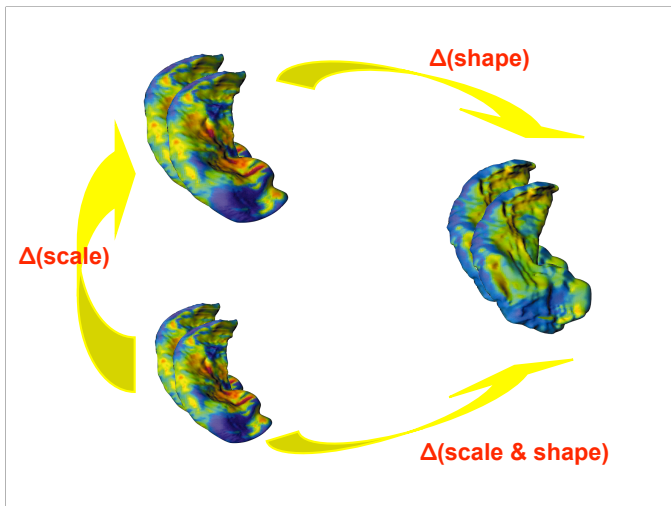
Combining Dissimilarities

Methods 1, 2 & 3 compared with Self-training

Densities of Lhats

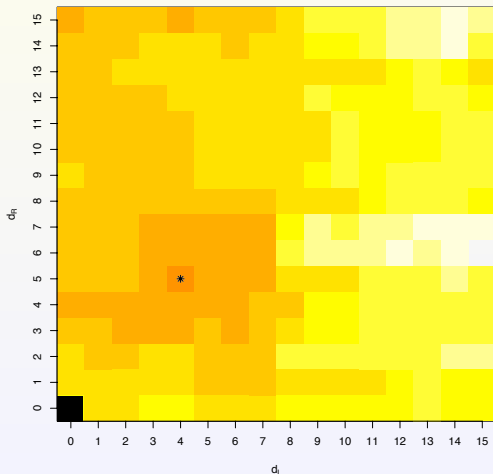


Brain Shape Comparison

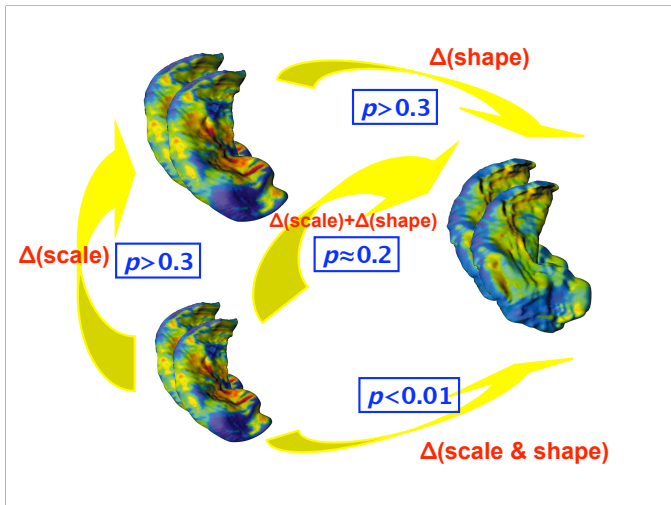


Brain Shape Comparison

- $\hat{L}_{d_L, d_R} = \frac{1}{m} \sum_{j=n+1}^{n+m} I\{\hat{Y}_j \neq Y_j | (X_1, Y_1), \dots, (X_n, Y_n)\}$
- $(d_L^*, d_R^*) \in \arg \min_{d_L, d_R} \hat{L}_{d_L, d_R}$



Brain Shape Comparison



Streaming Content in Context

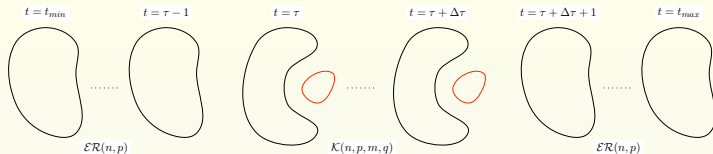


Figure: Anomaly detection in a time series of graphs

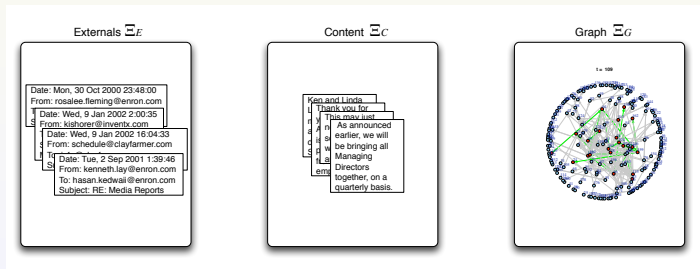


Figure: Conceptual depiction of the "Enron email data set"

Streaming Content in Context

Theorem

A test based on *fusion of externals and content* can be more powerful than a test based on either *externals alone* or *content alone*.

Proof.

We construct an example wherein the result holds.

$$H_0 : \mathcal{E}\mathcal{R}(n, p, \theta)$$

$$H_A : \mathcal{K}(n, p, \theta, m, q, \theta')$$

$T_1 = \text{size}(E)$: test statistic based on **externals** only

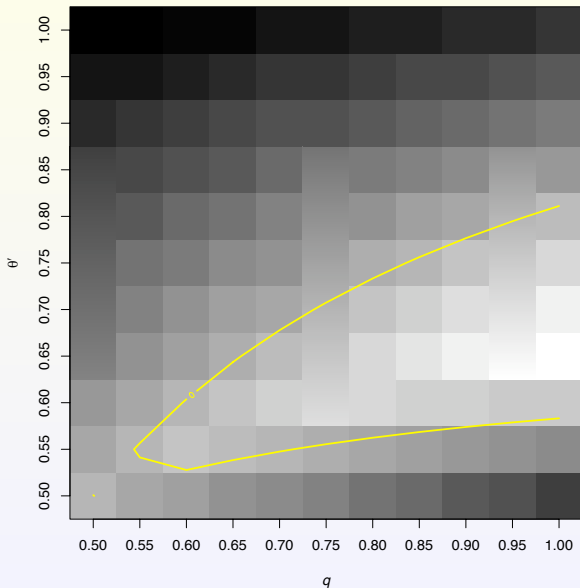
$T_2 = \text{size}(C)$: test statistic based on **content** only

$T = g(T_1, T_2)$: **disparate information fusion!**

We demonstrate that $\beta > \max\{\beta_1, \beta_2\}$ for appropriate choices of $n, p, \theta, \alpha, m, q, \theta'$. □

Streaming Content in Context

power difference: $\beta - \max(\beta_1, \beta_2)$

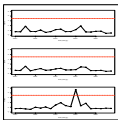


Streaming Content in Context

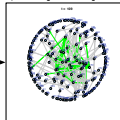
enron content in context

context analysis

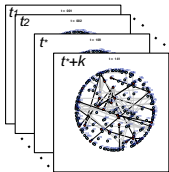
Detection at k^*, t^*, v^*



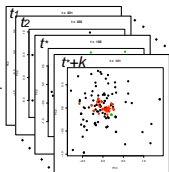
$\Omega(N_{k^*}[v^*; t^*]; t^*)$



Finally, tracking vertices incident to msgs C_t^* , denoted V_t^* , forwards and backwards in time from t^* provides genesis and outcome of topic of interest C_j .



Cluster msgs $M_t^* = M_t \cup \{S_j\}$. That cluster C_t^* s.t. $S_j \in C_t^*$ is the collection of msgs at time t associated with S_j .

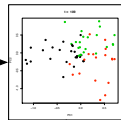


content analysis

Cluster msgs M_{t^*}

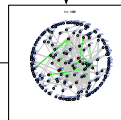
$$E_{t^*} = E(\Omega(N_{k^*}[v^*; t^*]; t^*))$$

$$M_{t^*} = \{\text{msgs } m_{vw} \text{ twixt } v \& w \text{ at time } t^* \text{ s.t. } vw \in E_{t^*}\}$$



Choose cluster C_j , a collection of msgs of interest which make up (part of) our detection - aka, a "topic".

$$S_j = \text{summary}(C_j)$$

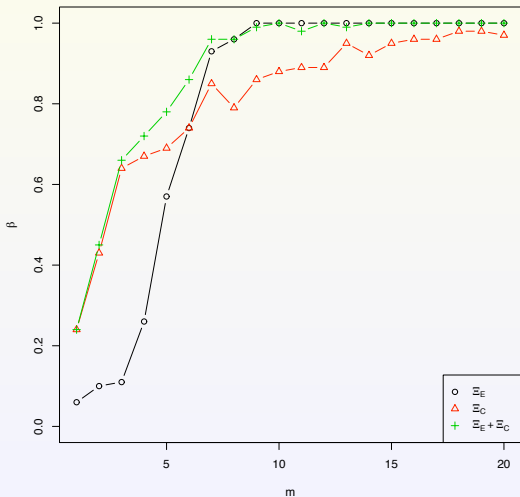


$\Omega(\{v \text{ s.t. } \text{msg } m_{vw} \in C_j\}; t^*)$

Streaming Content in Context

1000 Monte Carlo replicates at $m = 6$ yields statistical significance:

$$\hat{\beta} \approx 0.88 > \max\{\hat{\beta}_1 \approx 0.78, \hat{\beta}_2 \approx 0.76\}$$



Leopold Kronecker to Hermann von Helmholtz:

“The wealth of your practical experience with sane and interesting problems will give to mathematics a new direction and a new impetus.”



Leopold Kronecker



Hermann von Helmholtz