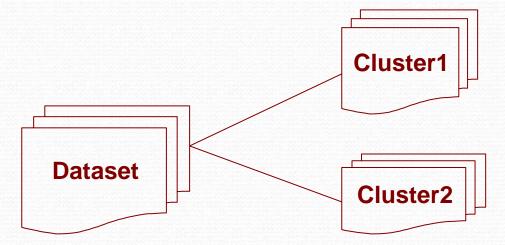
Topic-Wise, Sentiment-Wise, or Otherwise? Identifying the Hidden Dimension for Unsupervised Text Classification

Sajib Dasgupta and Vincent Ng Human Language Technology Research Institute University of Texas at Dallas

Goal

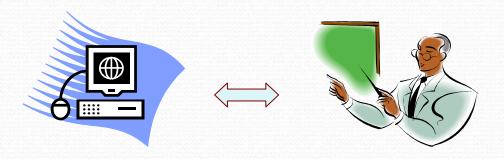


- Unsupervised Text Classification
 - Why unsupervised? Labeled data is expensive

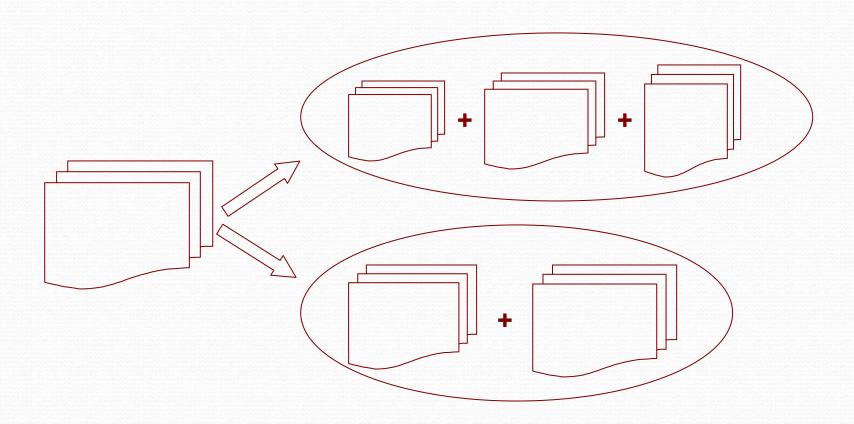
Goal

- Investigate"Why text clustering fails more often?"
- Possible Reasons:
 - Complexity of the task: Ambiguity
 - Large and complex feature space
 - No labeled data
- One more reason!

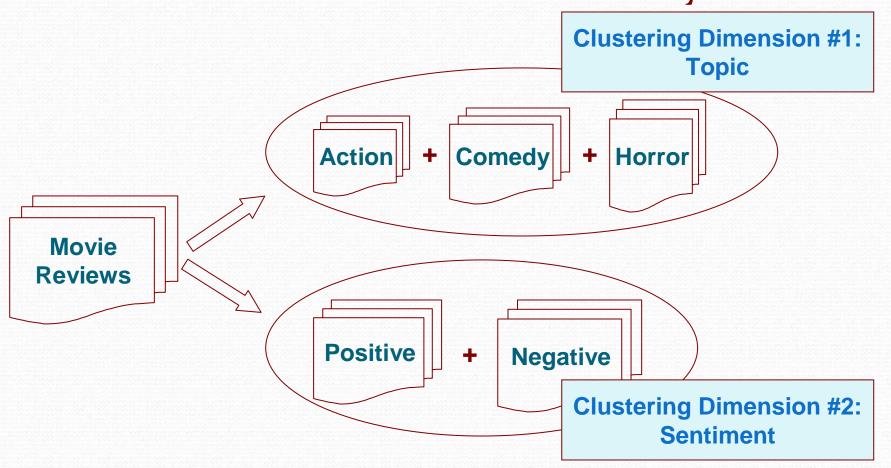
Conflict of interest between machine and human!



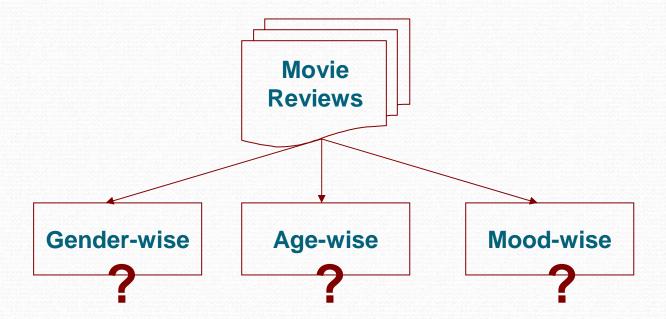
• Same data can be clustered different way



• Same data can be clustered different way



Same data can be clustered different way



- What clustering dimension user wants?
 - Topic-wise?
 - Sentiment-wise?
 - Otherwise?



- What clustering-dimension user wants?
 - Topic-wise?
 - Sentiment-wise?
 - Otherwise?



- What clustering-dimension clustering algorithm produces?
 - Topic-wise?
 - Sentiment-wise?
 - Otherwise?



Clustering algorithm fails if it fails to meet user demand

Goal

• Cluster a set of documents along the userspecified dimension

Dataset

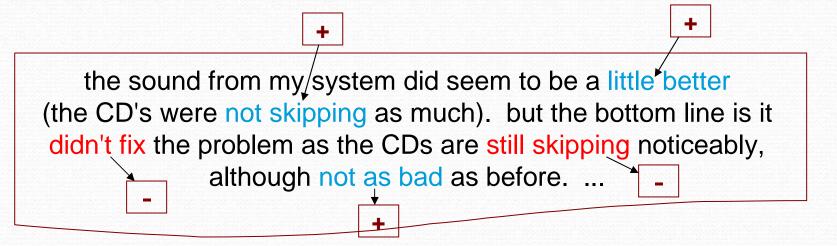
- Sentiment Domain:
 - 2-way clustering
 - Classify a review as "thumbs up" or "thumbs down"

Why sentiment domain?

Why Sentiment Domain?

- Clustering of sentiment is tough
- A lot tougher than topic-based clustering

- Sentimental ambiguity
 - A review may contain both positive and negative words!



A machine is more likely to label it positive.

- Sentimental ambiguity
 - A review may contain both positive and negative words!

the sound from my system did seem to be a little better (the CD's were not skipping as much). but the bottom line is it didn't fix the problem as the CDs are still skipping noticeably, although not as bad as before. ...

It's the bottom line that determines sentiment!

• But ... it's actually a negative review!

Subjective vs. Objective material

John Lynch wrote a classic in Spanish-American Revolutions 1808-1826.

He describes all the events that led to the independence of Latin America from Spain.

The book starts in Rio de La Plata and ends in Mexico and Central America.

Curiously one can note a common pattern of highly stratified societies lead by Spanish

The reluctance of Spanish Monarchy (and later even of liberals) led to independence

In all of those who are interested in a better understanding of Latin this great book is a must.

In all cleverly combines historical and economic facts about the Hispanic American societies

- Too much objective material might mislead a machine
 - if not human!

- Objective Features --> Topic-Wise Clusters
- Subjective Features --> Sentiment-Wise Clusters
- Too much of objective content might influence clustering
- Clustering algorithm is more likely to produce topic-wise clusters

Goal

Given all these complexities,

Is it possible to cluster a set of reviews along the sentiment dimension?

Goal

- Is it possible to cluster a set of reviews along the sentiment dimension?
- Yes, with the help of
 - Spectral learning techniques
 - A simple user feedback scheme

Lets start with spectral clustering,

- Algorithm for 2-way clustering (Ng et al. 2002)
- Given a data matrix D (dimension: $n \times f$),

- Algorithm for 2-way clustering (Ng et al. 2002)
- Given a data matrix *D* (dimension: *n* x *f*),
 - Form similarity matrix $S=\emptyset(D)$ (dimension: $n \times n$)

We used dot product

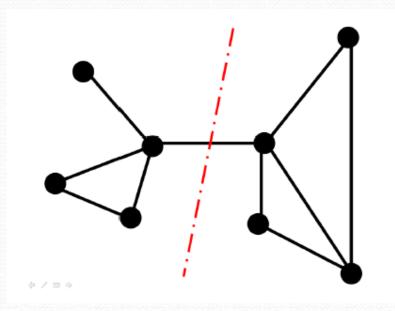
- Algorithm for 2-way clustering (Ng et al. 2002)
- Given a data matrix D (dimension: $n \times f$),
 - Form similarity matrix $S=\emptyset(D)$ (dimension: $n \times n$)
 - Form diagonal matrix G (dimension: $n \times n$), where G(i,i)=sum of the i-th row of S
 - Form Laplacian matrix $L=G^{-1/2} S G^{1/2}$

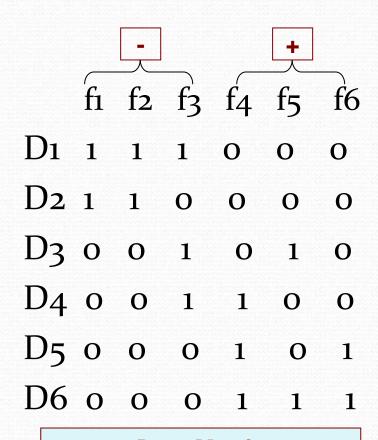
- Algorithm for 2-way clustering (Ng et al. 2002)
- Given a data matrix D (dimension: $n \times f$),
 - Form similarity matrix $S=\emptyset(D)$ (dimension: $n \times n$)
 - Form diagonal matrix G (dimension: $n \times n$), where G(i,i)=sum of the i-th row of S
 - Form Laplacian matrix $L=G^{-1/2} S G^{1/2}$
 - Find the eigenvector e corresponding to 2^{nd} largest eigenvalue of L

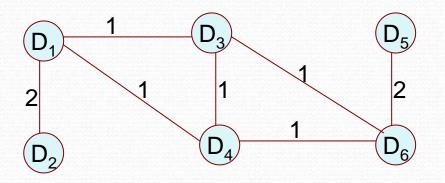
- Algorithm for 2-way clustering (Ng et al. 2002)
- Given a data matrix D (dimension: $n \times f$),
 - Form similarity matrix $S=\emptyset(D)$ (dimension: $n \times n$)
 - Form diagonal matrix G (dimension: $n \times n$), where G(i,i)=sum of the i-th row of S
 - Form Laplacian matrix $L=G^{-1/2} S G^{1/2}$
 - Find the eigenvector e corresponding to 2^{nd} largest eigenvalue of L
 - Use 2-means to cluster n data points using e

Why 2nd Eigenvector?

- Shi and Malik (2000): Normalized Cut and Image Segmentation
- 2nd eigenvector of the Laplacian induces the **normalized mincut** of a graph formed from the similarity matrix S

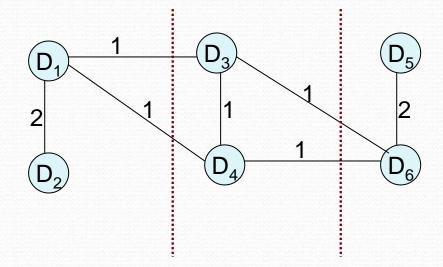






Similarity Graph

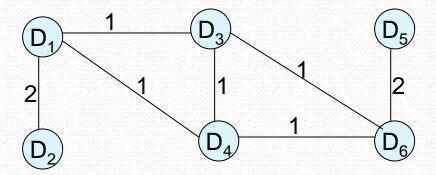
• Similarity Graph



Two possible normalized mincut partitions

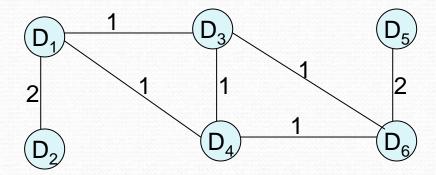
Let's see what partition 2nd eigenvector produces

• Similarity Graph

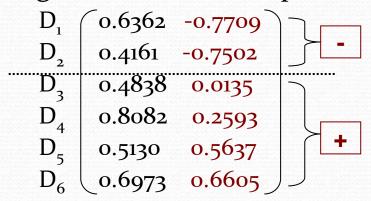


• Top 2 eigenvectors of its Laplacian

• Similarity Graph



• Top 2 eigenvectors of its Laplacian



2nd Eigenvector

- 2nd eigenvector captures the **most prominent** clustering dimension
- But does it capture sentiment?
 - Sentiment might be hidden

How to cluster along the desired (possible hidden) dimension?

Our Algorithm

- We propose a novel feedback-oriented spectral clustering algorithm
- It allows us to achieve the desired clustering

• It has 4 steps,

Our Algorithm

Steps:

- Identify important dimensions
- Identify unambiguous reviews
- Identify relevant features and get user-feedback
- Cluster along the selected dimension

Our Algorithm

Steps:

- Identify important dimensions
- Identify unambiguous reviews
- Identify relevant features and get user-feedback
- Cluster along the selected dimension

Identify Important Dimensions

- Spectral clustering setup:
 - Features: Bag of words
 - Similarity metric: dot product

Identify Important Dimensions

Given a data matrix *D*,

- We form similarity matrix S, diagonal matrix G, and Laplacian matrix $L=G^{-1/2}SG^{-1/2}$
- Compute the **5 eigenvectors** (with largest eigenvalues) from the Laplacian matrix

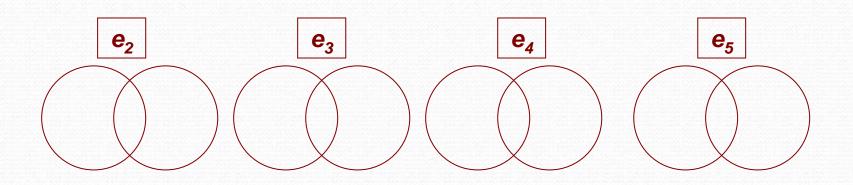
Identify Important Dimensions

Given a data matrix *D*,

- We form similarity matrix S, diagonal matrix G, and Laplacian matrix $L=G^{-1/2}SG^{-1/2}$
- Compute the **5 eigenvectors** (with largest eigenvalues) from the Laplacian matrix
 - Each eigenvector (starting from 2nd) captures an important clustering dimension
 - They capture the largest variance in the data
 - First eigenvector is uninformative

Identify Important Dimensions

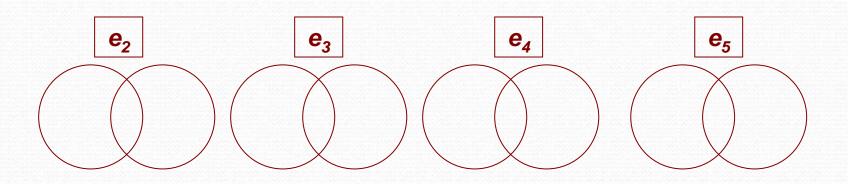
• Each of the 2nd to 5th eigenvector leads to a certain clustering along an important dimension



Which one is sentiment-oriented?

Identify Important Dimensions

• Each of the 2nd to 5th eigenvector leads to a certain clustering along an important dimension



Which one is sentiment-oriented?

Ask human!

Human Feedback

- We don't give *clusters* to humans to judge!
 - It's time consuming
 - It will make it a supervised method
- We give them *features* representative of each cluster

Human Feedback

- We don't give *clusters* to humans to judge!
 - It's time consuming
 - It will make it a supervised method
- We give them *features* representative of each cluster
 - It's simple, as simple as a cursory look!
 - At least, it's simple to 5 graduate students

• Before we proceed to human-feedback step,

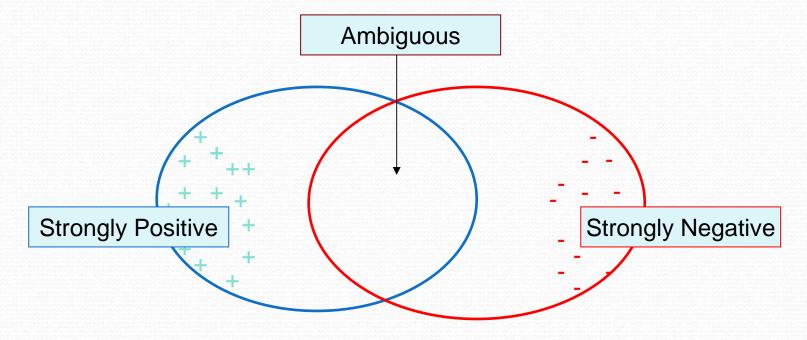
We have one other important step

Our Algorithm

Steps:

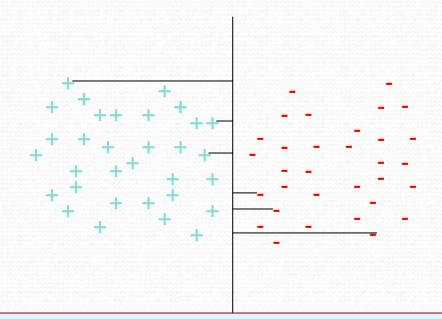
- Identify important dimensions
- Identify unambiguous reviews
- Identify relevant features and get user-feedback
- Cluster along the selected dimension

- For each partition corresponding to e₂ to e₅
- We separate out unambiguous reviews
 - in an unsupervised manner



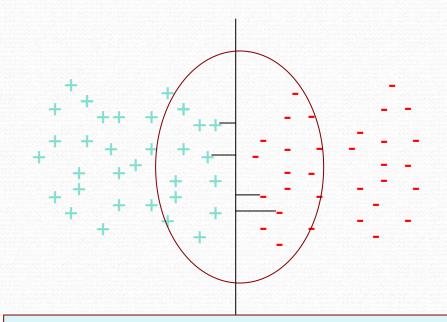
- Why separate unambiguous reviews?
 - Feature distribution learned from only unambiguous portion is more reliable
 - Ambiguous reviews can only mislead!

• How to separate unambiguous reviews?



Orthogonal projections on a learned-dimension

• How to separate unambiguous reviews?



Ambiguous points have small projections

Our Algorithm

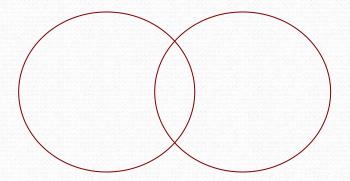
- For each of the 2nd to 5th eigenvector e,
 - we sort the data points according to *e* components
 - keep only top $\alpha/2$ and bottom $\alpha/2$ points
 - create two clusters U_1 and U_2 with top $\alpha/2$ and bottom $\alpha/2$ points respectively

Our Algorithm

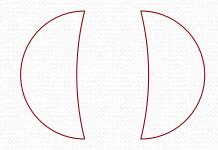
- For each of the 2nd to 5th eigenvector e,
 - we sort the data points according to *e* components
 - keep only top $\alpha/2$ and bottom $\alpha/2$ points
 - create two clusters U_1 and U_2 with top $\alpha/2$ and bottom $\alpha/2$ points respectively
 - We use $\alpha = n/4$
 - Assuming that at least 25% reviews are unambiguous

Clustering Accuracy

All data points



Clustering Accuracy: Low **Unambiguous data points**



Clustering Accuracy: High

Clustering accuracy of unambiguous reviews

• Accuracy of 500 unambiguous reviews out of 2000.

	Clustering
	Accuracy (%)
Movie	87.0
Kitchen	87.6
Electronics	77.6
Books	86.2
DVD	87.4

Clustering accuracy of unambiguous reviews

Accuracy of 500 unambiguous reviews out of 2000.

	Clustering
	Accuracy (%)
Movie	87.0
Kitchen	87.6
Electronics	77.6
Books	86.2
DVD	87.4

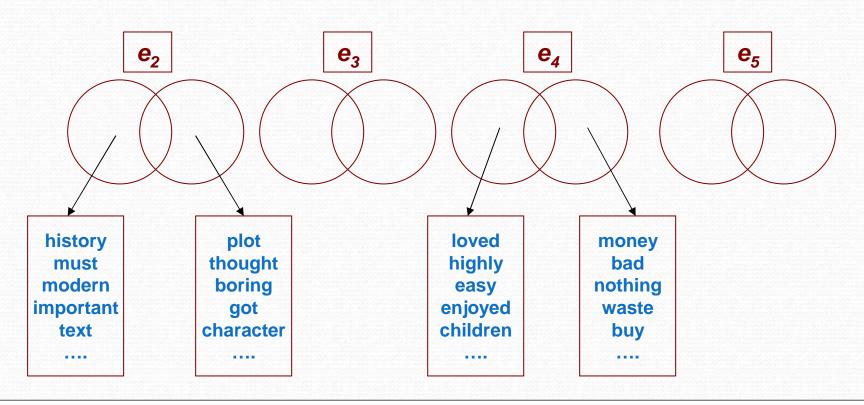
High-accuracy is essential to learn reliable feature distribution

Our Algorithm

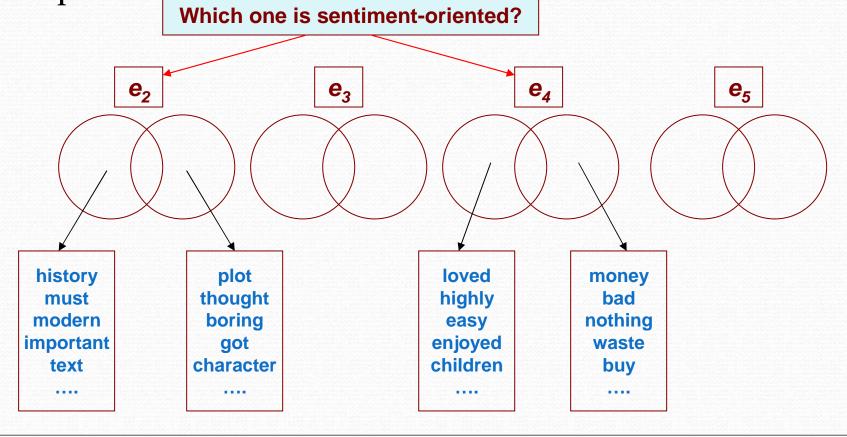
Steps:

- Identify important dimensions
- Identify unambiguous reviews
- Identify relevant features and get user-feedback
- Cluster along the selected dimension

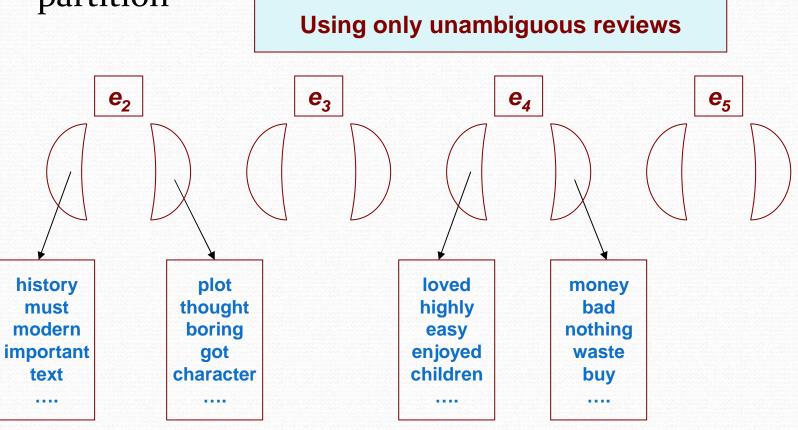
• We extract informative features characterizing each partition



• We extract informative features characterizing each partition



• We extract informative features characterizing each partition



 How to extract informative features characterizing a partition?

- How to extract informative features characterizing a partition?
 - Use Support Vector Machine
 - Why?

Max-margin systems learn feature distribution well

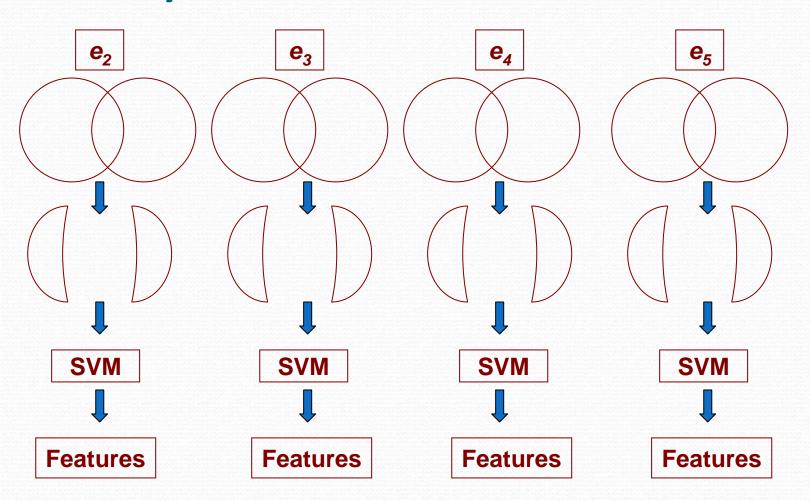
• Use SVM to learn the max-margin hyperplane (i.e., $w \cdot x - b = 0$)

- Use w to extract informative features
 - Features with large positive weight indicative of positive class
 - Features with large negative weight indicative of negative class

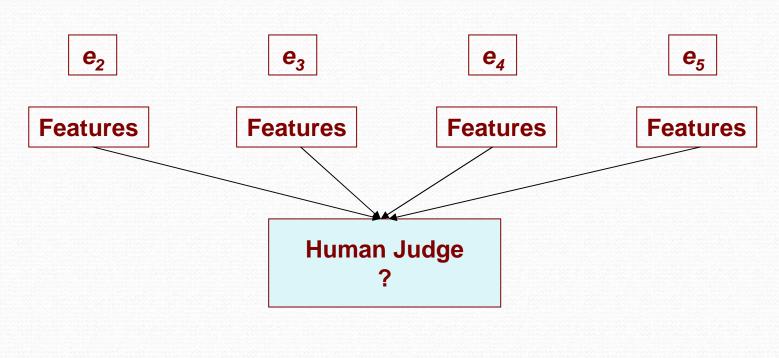
- For each of the 2nd to 5th eigenvector e,
 - Train SVM classifier on unambiguous clusters U_1 and U_2
 - Sort the features according to learned **w**
 - List₁ -> top 100
 - List₂ -> bottom 100

List₁: Informative Features for + Class

List₂: Informative Features for - Class



Get Human Feedback



Which one is sentiment-oriented?

Our Algorithm

Steps:

- Identify important dimensions
- Identify unambiguous reviews
- Identify relevant features and get user-feedback
- Cluster along the selected dimension

Cluster along selected dimension

- Get the dimension (eigenvector) selected by human
- Use 2-means to cluster all data points using the selected eigenvector

End of algorithm

• The algorithm for producing the feature list is unsupervised

 Except for human feedback, the whole clustering process is unsupervised • Except for human feedback, the whole clustering process is unsupervised

But, does it actually work?

• Let's do the evaluation,

- Datasets
 - Movie (Pang et al., 2002)
 - 4 datasets from Blitzer et al. (2007)
 - Kitchen and Housewares
 - Electronics
 - Books
 - DVD
 - each has 2000 points (1000 positives & 1000 negatives)

- One more dataset
 - 2-Newsgroup: sci.crypt and talks.politics
 - Goal: Science vs. Politics
 - It has 2000 points (1000 Science & 1000 Politics)

Show the difference between topic-based clustering and sentiment-based clustering

- Show the SVM-learned feature list used for user feedback
 - Is our feedback system easy and feasible?
- 2. Human Agreement Rate
 - Did humans actually find it easy?
- **3.** Clustering Accuracy
 - Did we achieve sentiment-oriented clustering?

- 1. Show the SVM-learned-feature-list used for user-feedback
 - Is our feedback system easy and feasible?
- 2. Human Agreement Rate
 - Did humans actually find it easy?
- 3. Clustering Accuracy
 - Did we achieve sentiment-oriented clustering?

Top Features

• DVD

2nd eigenvector

C₁

worth
bought
series
money
got
season
fan
collection
music
tv
thought
quality

C_2

young
between
actors
men
cast
seems
job
beautiful
around
us
our
director

3rd eigenvector

C₁

music
collection
excellent
wonderful
must
loved
perfect
highly
makes
special
performance
picture

....

 C_2

worst

money
thought
boring
nothing
minutes
waste
saw
pretty
reviews
interesting
maybe

.

2nd eigenvector

C₁

worth

bought series money got season fan collection

> music tv

thought quality

.

C_2

young between actors men

cast

job

beautiful

around

us

our

director

....

3rd eigenvector

C

music collection

excellent

wonderful

must

loved

perfect

highly

makes

special performance

picture

....

C_2

worst

money

thought

boring nothing

minutes

waste

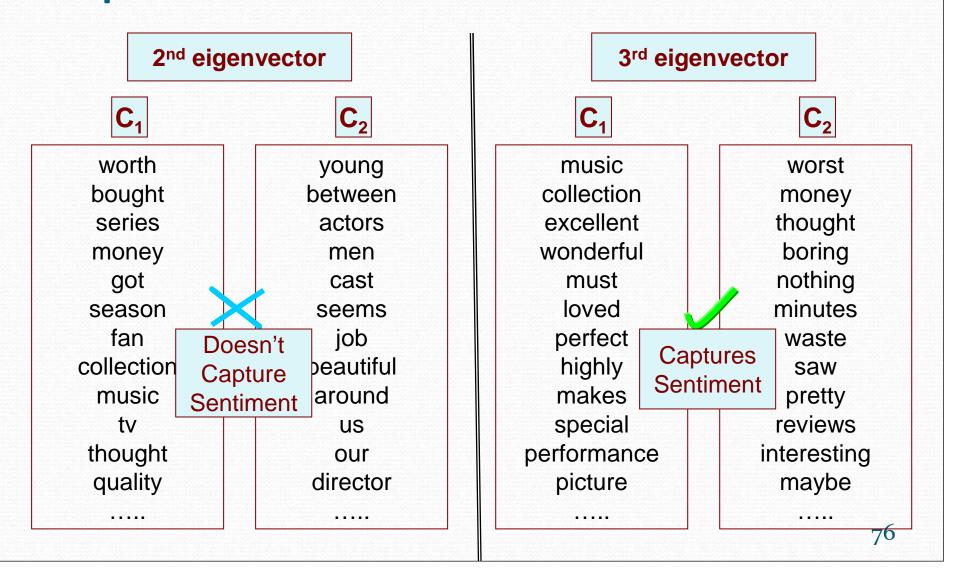
saw

pretty

reviews

interesting maybe

.



4th eigenvector

C₁

video
music
found
feel
bought
workout
daughter
recommend

disappointed

our

right moves

....

 C_2

series
cast
fan
stars
original
comedy

comedy actors worth classic action

season big 5th eigenvector

C₁

saw watched

loved enjoy

whole

got

family

series

season liked

entertaining

lot

....

 C_2

money quality

video

worth

found

version

picture

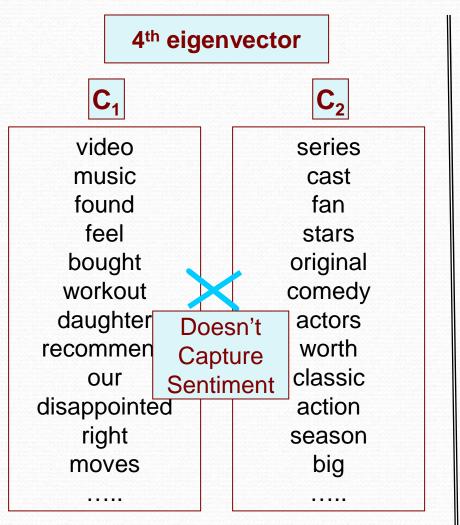
waste

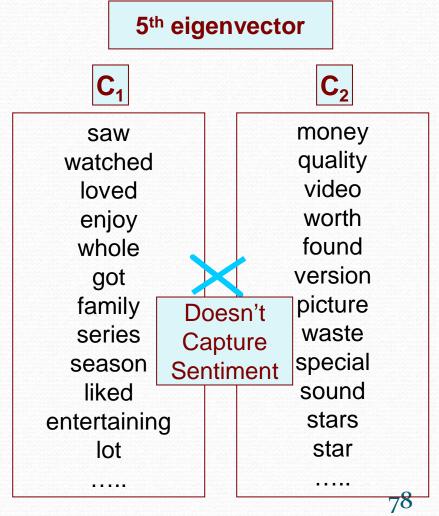
special

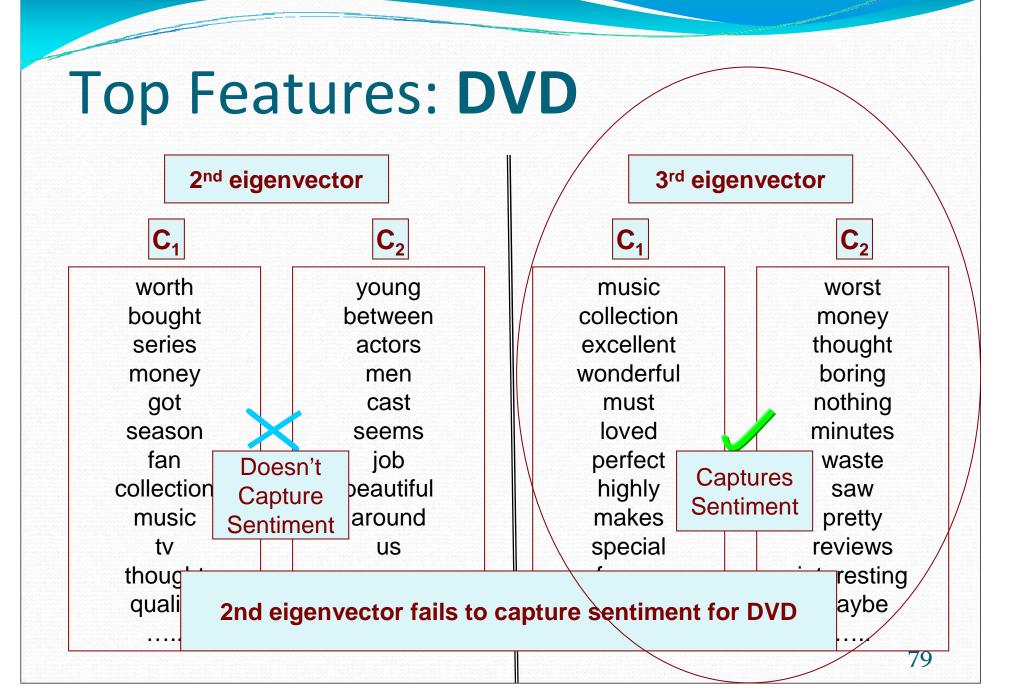
sound

stars star

.







Top Features

Books

2nd eigenvector

history must modern important text reference

excellent

provides business both understanding given

plot didn thought boring got character couldn ending fan pages funny

3rd eigenvector

series man history character death between war seems political american during plot

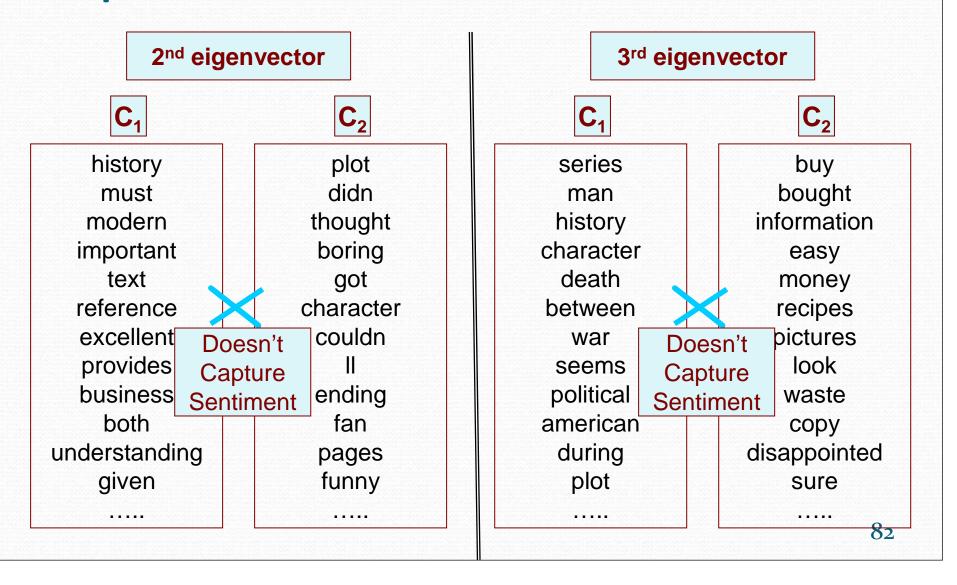
buy bought information

easy

money recipes pictures look waste copy disappointed

sure

81



4th eigenvector

C₁

loved
highly
easy
enjoyed
children
again
although
excellent

understand three both looking

.

 C_2

money
bad
nothing
waste
buy
anything
doesn
already
instead
seems
believe
thought

5th eigenvector

C₁

must wonderful

old feel away children

someone man made

year

thing buy

....

O₂

boring series

history

pages

information

between

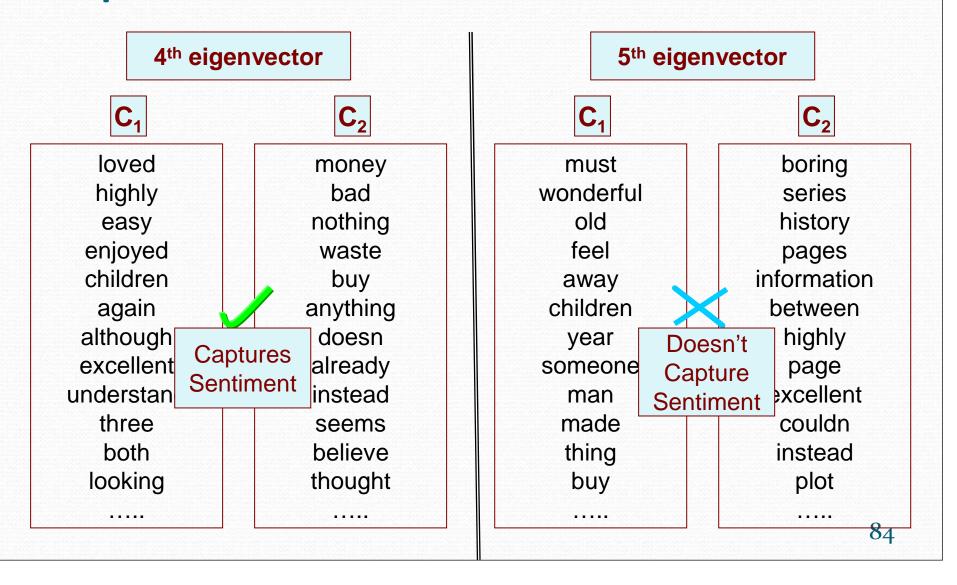
highly

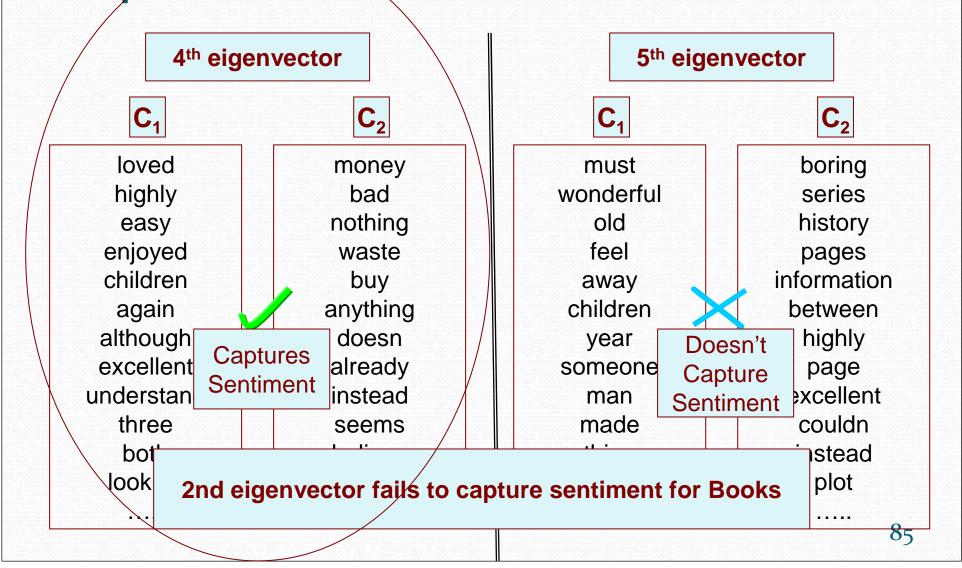
page

excellent

couldn instead plot

. . . .

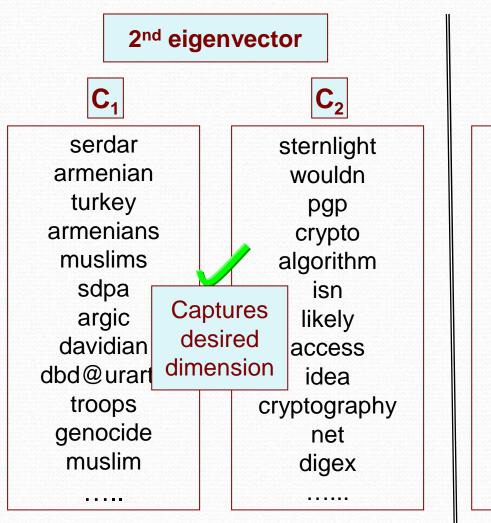


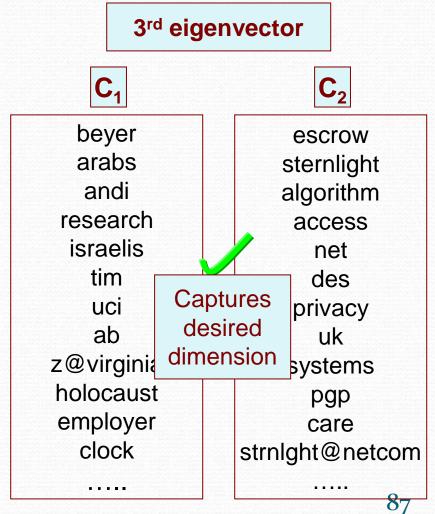


Human Feedback

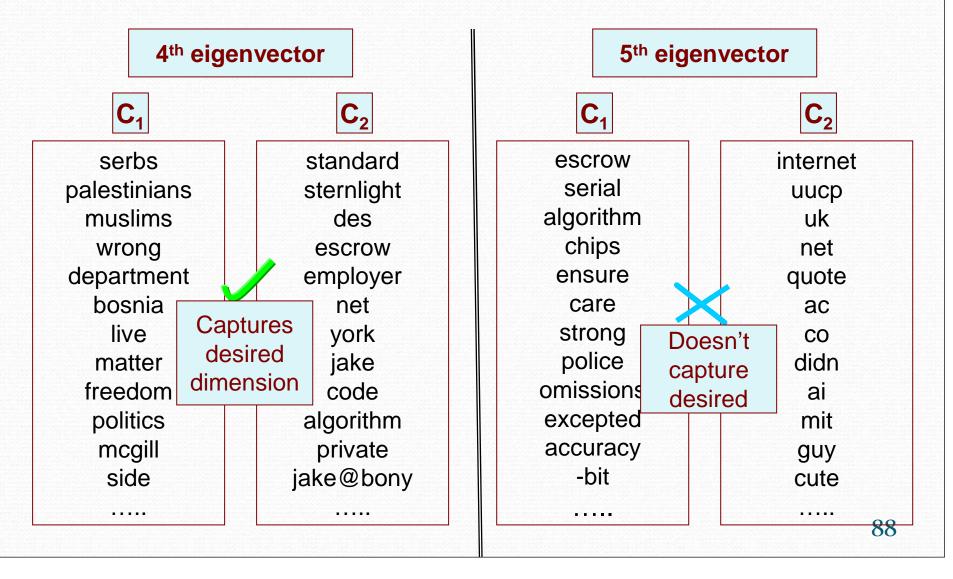
- Lets see what happens to 2-Newsgroup
- It will show the difference between topic clustering and sentiment clustering

Top Features: Politics vs. Science





Top Features: Politics vs. Science



Top Features: Politics vs. Science

2nd, 3rd, 4th eigenvectors all capture right dimension

Topic based clustering is different from sentiment based clustering

Evaluation

- Show the SVM-learned-feature-list used for userfeedback
 - Is our feedback system easy and feasible?

2. Human Agreement Rate

- Did humans actually find it easy?
- 3. Clustering Accuracy
 - Did we achieve sentiment-oriented clustering?

Human Agreement Rate

- Human Judges:
 - 5 computer science graduate students

• Guidelines:

- We show top 100 features for each cluster
- We inform them of desired clustering dimension beforehand

Human Agreement Rate

• Did they find sentiment dimension correctly?

	Correctly found
Movie	5/5
Books	5/5
DVD	5/5
Electronics	5/5
Kitchen	4/5

Human Agreement Rate

- Almost perfect human agreement rate
- Feedback approach is feasible

Human Feedback

• Eigenvector selected by all human-judges

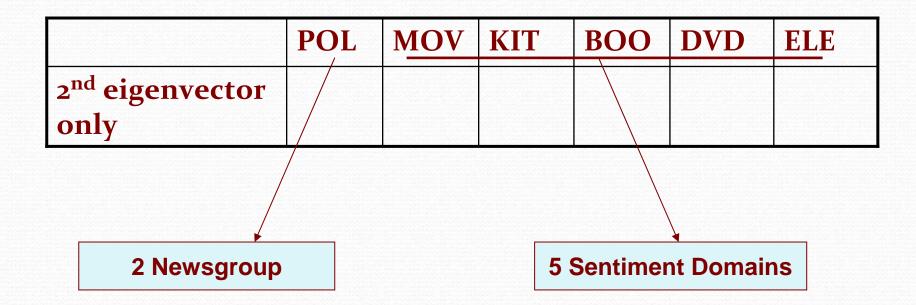
	Eigenvector
Movie	2
Books	4
DVD	3
Electronics	3
Kitchen	2

For 3 out of 5 domains, sentiment is hidden

Evaluation

- Show the SVM-learned-feature-list used for userfeedback
 - Is our feedback system easy and feasible?
- 2. Human Agreement Rate
 - Did humans actually find it easy?
- **3.** Clustering Accuracy
 - Did we achieve sentiment-oriented clustering?

- Baselines:
- 2nd eigenvector only
- Evaluation Metric:
 - Accuracy



	POL	MOV	KIT	ВОО	DVD	ELE
2 nd eigenvector	93.7%	70.9%	69.7%	58.9%	55.3%	50.8%
only						

Clustering accuracy high for POL

Clustering of sentiment is tough

	POL	MOV	KIT	ВОО	DVD	ELE
2 nd eigenvector	93.7%	70.9%	69.7%	58.9%	55.3%	50.8%
only						

Clustering accuracy low 2nd Eigenvectors fails to capture sentiment

	POL	MOV	KIT	ВОО	DVD	ELE
2 nd eigenvector only	93.7%	70.9%	69.7%	58.9%	55.3%	50.8%
Our approach	93.7%	70.9%	69.7%	69.5%	70.8%	65.8%

We achieve best performance for all sentiment domain

	POL	MOV	KIT	ВОО	DVD	ELE
2 nd eigenvector only	93.7%	70.9%	69.7%	58.9%	55.3%	50.8%
Our approach	93.7%	70.9%	69.7%	69.5%	70.8%	65.8%

Human helps to identify the hidden sentiment dimension

Conclusion

- Proposed a new feedback-oriented clustering algorithm:
 - It produces the desired clustering
 - It requires only simple human feedback

Conclusion

• Finally, we have a system where

Machine meets human



• Thank you