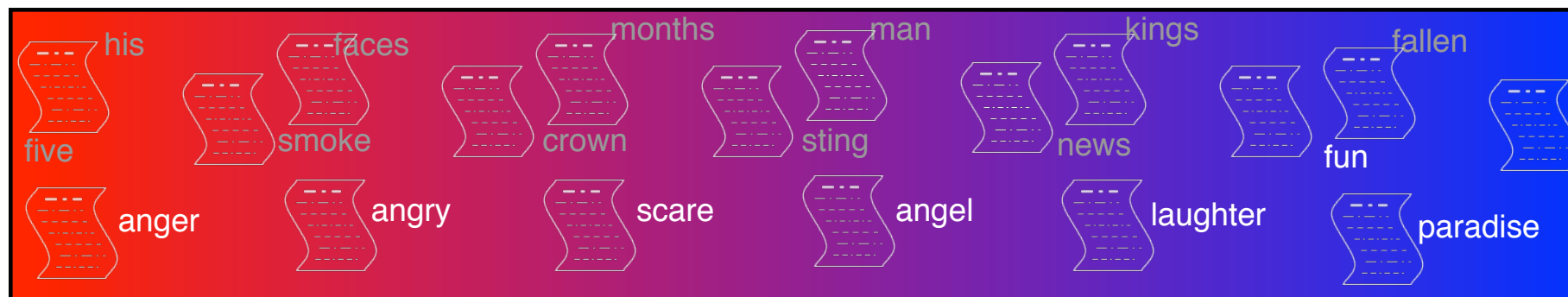


# Multilingual Sentiment Analysis

## Using Latent Semantic Indexing and Machine Learning



Brett Bader, Digital Globe, [bbader@digitalglobe.com](mailto:bbader@digitalglobe.com)

Philip Kegelmeyer, Sandia National Laboratories, [wpk@sandia.gov](mailto:wpk@sandia.gov)

Peter Chew, Galisteo, [PeterAChew@aol.com](mailto:PeterAChew@aol.com)



*Sandia is a multiprogram laboratory operated by Sandia Corporation, a Lockheed Martin Company, for the United States Department of Energys National Nuclear Security Administration under contract DE-AC04-94AL85000.*



SENTIRE, December 11, 2011

# Overview

We treat **multilingual document sentiment classification**

as a supervised machine learning problem, which requires:

- multilingual document attributes
- monolingual ground truth documents

We assess performance,

raise objections,

and address one of them.

## Document Sentiment Valence

### Psalm 126:2–4

Then was our mouth filled with laughter, and our tongue with singing: then said they among the heathen, The LORD hath done great things for them. The LORD hath done great things for us; whereof we are glad. Turn again our captivity, O LORD, as the streams in the south.

### Revelation 9:18–19

By these three was the third part of men killed, by the fire, and by the smoke, and by the brimstone, which issued out of their mouths. For their power is in their mouth, and in their tails: for their tails were like unto serpents, and had heads, and with them they do hurt.

# Document Sentiment Valence, Multilingually

## Salmos 126:2–4

Entonces nuestra boca se henchirá de risa, Y nuestra lengua de alabanza; Entonces dirán entre las gentes: Grandes cosas ha hecho Jehová con éstos. Grandes cosas ha hecho Jehová con nosotros; Estaremos alegres. Haz volver nuestra cautividad oh Jehová, Como los arroyos en el austro.

## Apocalipsis 9:18–19

De estas tres plagas fué muerta la tercera parte de los hombres: del fuego, y del humo, y del azufre, que salan de la boca de ellos. Porque su poder está en su boca y en sus colas: porque sus colas eran semejantes serpientes, y tenían cabezas, y con ellas dañan.

# Overview

We treat multilingual document sentiment classification

as a **supervised machine learning** problem, which requires:

- multilingual document attributes
- monolingual ground truth documents

We assess performance,

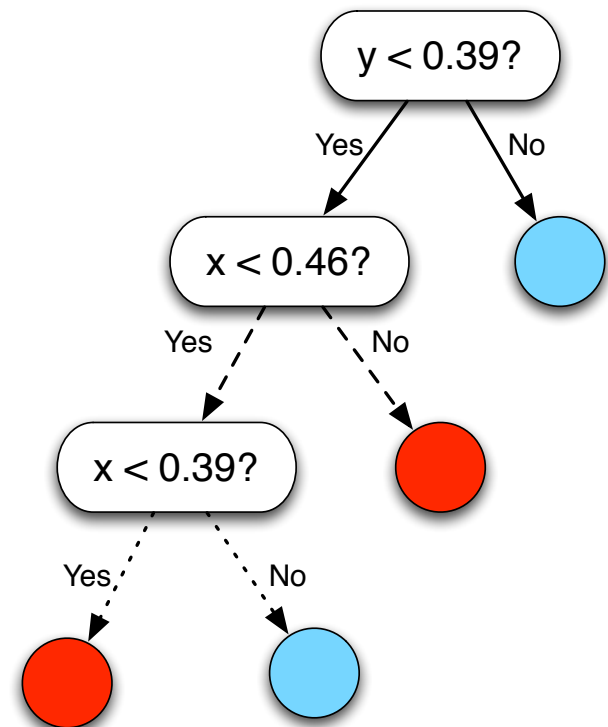
raise objections,

and address one of them.

# Supervised Machine Learning Overview

Also known as: pattern recognition, statistical inference, data mining.

- Input: “ground truth” data.
  - Samples, with attributes and labels.
  - For document sentiment analysis:
    - \* Samples: documents
    - \* Attributes: concept weights (to be described)
    - \* Labels: **positive**, **negative**
- Apply suitable method:  
decision trees, neural nets, SVMs.
- Output:  
rules for labeling new, *unlabeled* documents.



Decision tree representation.

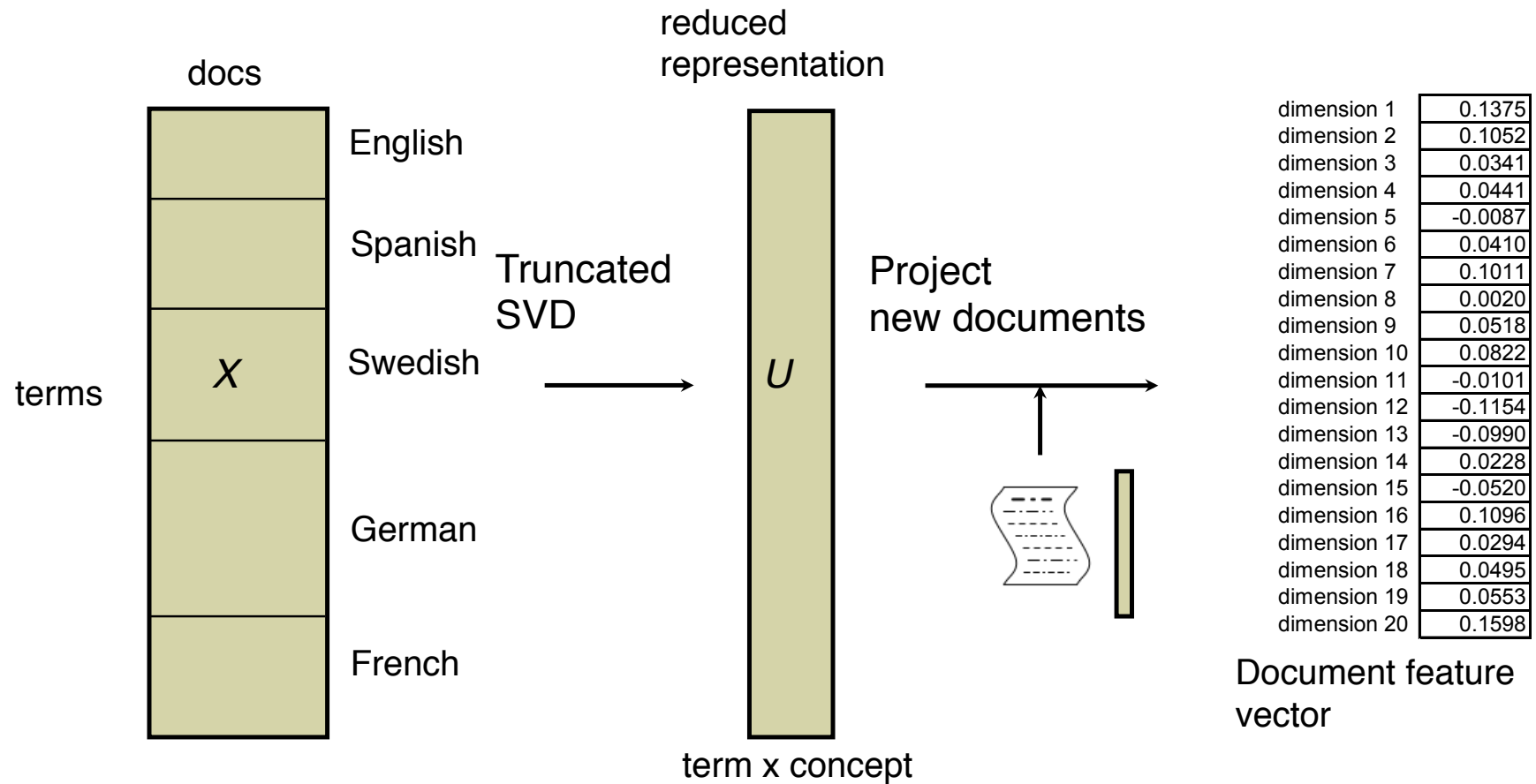
# Overview

We treat multilingual document sentiment classification  
as a supervised machine learning problem, which requires:

- **multilingual document attributes**
- monolingual ground truth documents

We assess performance,  
raise objections,  
and address one of them.

# Multilingual Latent Semantic Analysis

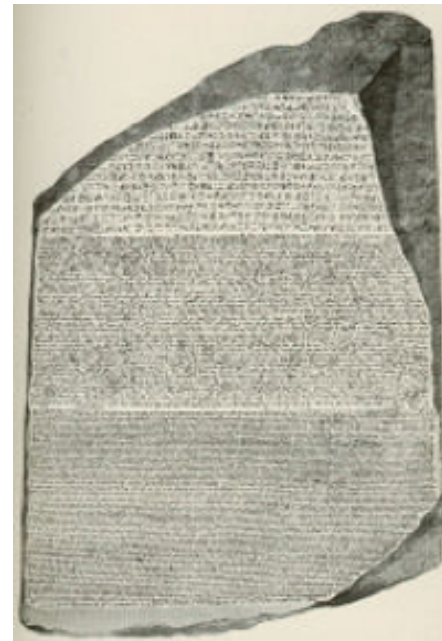


Result: documents represented by *language-independent* features.



# Europarl Corpus as Our “Rosetta Stone”

- Translations of the proceedings of the European Parliament.
- Sentence aligned text  
16M sentences across 11 languages.
- 1,247,832 speeches,
- 1,249,253 terms across 11 languages.



Caption

# Overview

We treat multilingual document sentiment classification  
as a supervised machine learning problem, which requires:

- multilingual document attributes
- **monolingual ground truth documents**

We assess performance,  
raise objections,  
and address one of them.

# Groundtruthing the Bible for Sentiment

- Could use exhaustive human reading and judgment.  
(Requires exhaustion. And judgment).
- We used a sentiment lexicon to bootstrap the process.  
(Sentiment lexicon not strictly required;  
any accurate labeling mechanism suffices.)
- A sentiment lexicon[1] maps terms to “valence”

Term	Valence, 0 to 1
ace	6.88
ache	2.26
...	...
fun	8.27
funeral	1.39

...

Chronicles 4

Chronicles 5

Chronicles 6

...

Psalms 125

Psalms 126

Psalms 127

...

Revelation 8

Revelation 9

Revelation 10

Revelation 11

...

## Initial Scoring for Each Bible Chapter

- For each chapter
  - Add up (centered) valences
  - Normalize by the number of terms
- Find the 100 most positive, 100 most negative.
- Inspect only those, to confirm.

	Term	Valence
In the beginning	God	8.15
created the	heavens	7.30
and the	earth .	7.15
And the	earth	7.15
was	waste	2.93
and void...		
(Genesis, Chapter 1, ranks 227 out of 1188 chapters.)		

# Hand Inspection Was Necessary

Revelation 9:1–12 (a demonic plague of locusts) scored *positive*.

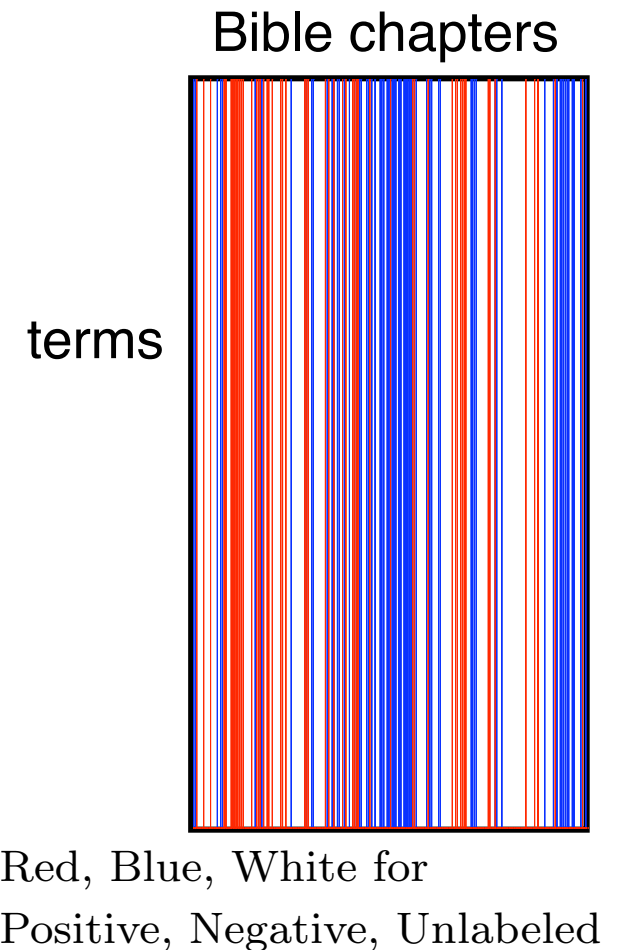
1 The fifth angel sounded his trumpet, and I saw a star that had fallen from the sky to the earth. The star was given the key to the shaft of the Abyss. 2 When he opened the Abyss, smoke rose from it like the smoke from a gigantic furnace. The sun and sky were darkened by the smoke from the Abyss. 3 And out of the smoke locusts came down upon the earth and were given power like that of scorpions of the earth. 4 They were told not to harm the grass of the earth or any plant or tree, but only those people who did not have the seal of God on their foreheads. 5 They were not given power to kill them, but only to torture them for five months. And the agony they suffered was like that of the sting of a scorpion when it strikes a man. 6 During those days men will seek death, but will not find it; they will long to die, but death will elude them.

7 The locusts looked like horses prepared for battle. On their heads they wore something like crowns of gold, and their faces resembled human faces. 8 Their hair was like women's hair, and their teeth were like lions' teeth. 9 They had breastplates like breastplates of iron, and the sound of their wings was like the thundering of many horses and chariots rushing into battle. 10 They had tails and stings like scorpions, and in their tails they had power to torment people for five months. 11 They had as king over them the angel of the Abyss, whose name in Hebrew is Abaddon, and in Greek, Apollyon.

12 The first woe is past; two other woes are yet to come. ...

# Final Sentiment Groundtruth Dataset

- Manual inspection turned up a few “Revelations 9” problems.
- Weeded by hand, and re-seeded.
- Final result (out of 1188 chapters)
  - 115 positive chapters
  - 78 negative chapters
  - Positive/negative ratio is 59.6%.



# Overview

We treat multilingual document sentiment classification  
as a supervised machine learning problem, which requires:

- multilingual document attributes
- monolingual ground truth documents

We **assess performance**,

raise objections,

and address one of them.

## Pause: What Have Other People Done?

How have others attempted multilingual sentiment assignment?

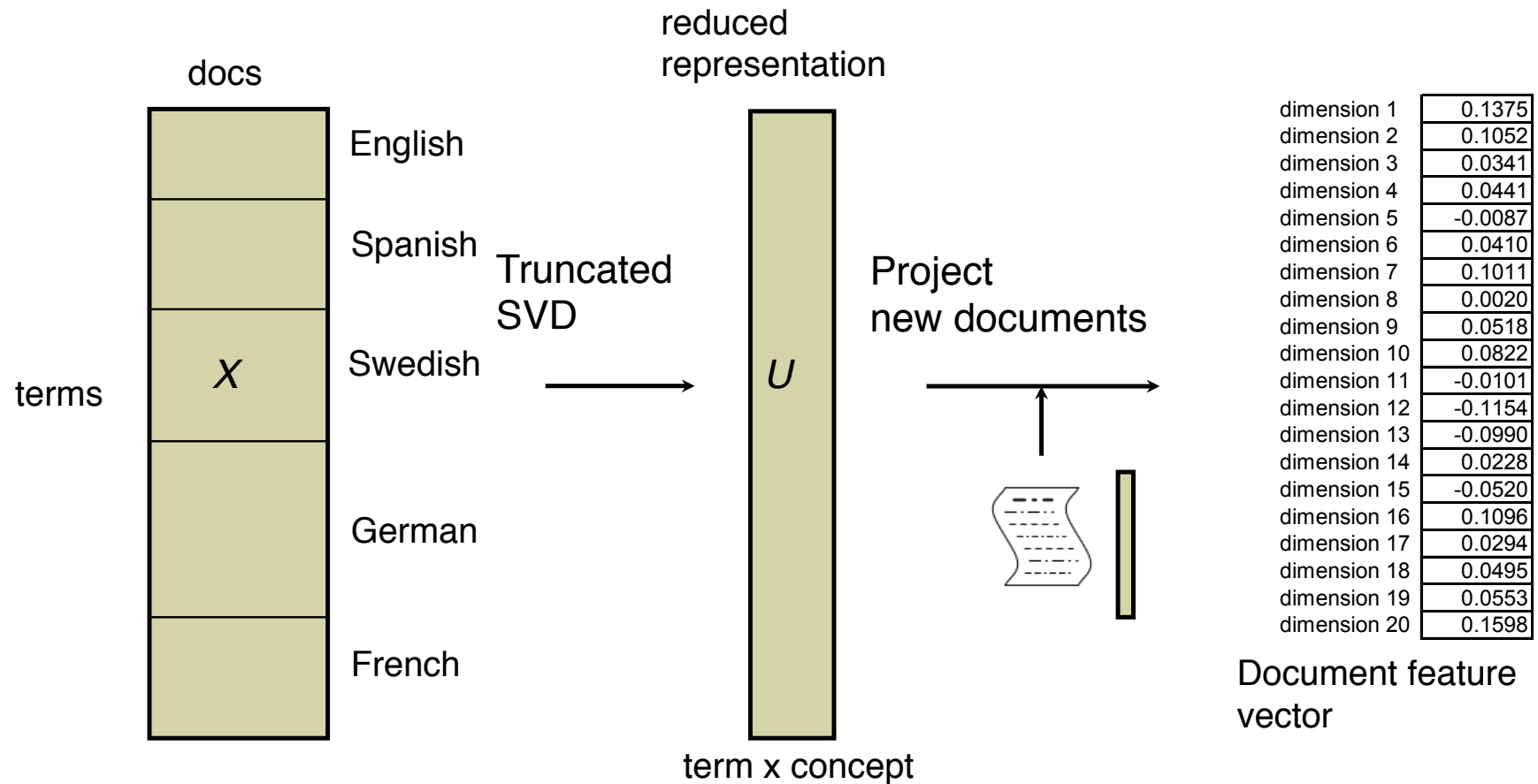
- Translate documents into English (for which sentiment-tagged data is more readily available), then apply sentiment analysis[2].
- Or the reverse: translate sentiment lexicons from English into the target language, classify documents directly in the other language[3].
- Or simply invest the brute force labor needed to create a sentiment lexicon in every language of interest[4].

(Our method requires no translation.

Further, it needs a sentiment lexicon only as one way to boot strap.)



# Project English Bible Chapters through Europarl ...



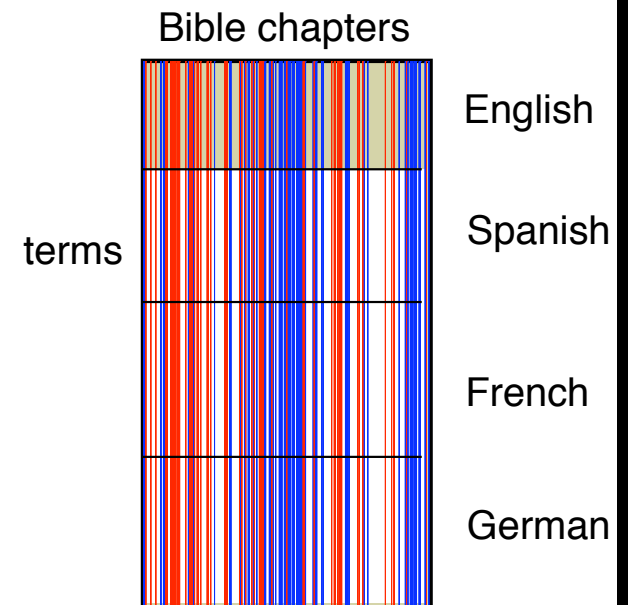
## Use only the Labeled Chapters ...

...to generate training data that looks like:

Chapter	Valence	$f_1$	$f_2$	$f_3$	...	$f_{300}$
$c_1$ , Psalm 126	Positive	0.12	0.03	0.97	...	0.12
$c_2$ , Psalm 127	Positive	0.99	0.02	0.33	...	0.03
$c_3$ , Chronicles 5	Negative	0.30	0.27	0.12	...	0.13
$c_4$ , Revelation 10	Positive	0.16	0.83	0.08	...	0.58
$c_5$ , Chronicles 5	Negative	0.17	0.65	0.36	...	0.64
$c_6$ , Ezra 10	Negative	0.44	0.12	0.29	...	0.42
$c_7$ , Ezekiel 5	Negative	0.42	0.24	0.33	...	0.88
$c_8$ , James3	Positive	0.78	0.42	0.44	...	0.52
⋮	⋮	⋮	⋮	⋮		⋮
$c_{193}$ , Revelation 9	Negative	0.12	0.41	0.92	...	0.17

## Test on the Foreign, Labeled Chapters

- Build an ensemble of bagged decision trees.
- Project foreign language chapters through Europarl.
- Assume that sentiment is preserved across languages.
- Use the ensemble to classify the 3x193 chapters in Spanish, French, German.
- Result:
  - Accuracy of 74.9%
  - Statistically significantly (one-sample  $t$ -test,  $\alpha = 0.01$ ) better than the ...
  - Baseline random accuracy of 56.9%



# Overview

We treat multilingual document sentiment classification  
as a supervised machine learning problem, which requires:

- multilingual document attributes
- monolingual ground truth documents

We assess performance,

**raise objections,**

and address one of them.

# Have We Simply Learned Topic, Not Sentiment?

Maybe. It does seem plausible that topic and sentiment are intertwined.

So we changed the training data by shuffling versions across chapters.

Chapter 10, Verse 1	Chapter 88, Verse 2
Chapter 10, Verse 2	Chapter 10, Verse 4
Chapter 10, Verse 3	Chapter 92, Verse 1
Chapter 10, Verse 4	Chapter 10, Verse 3

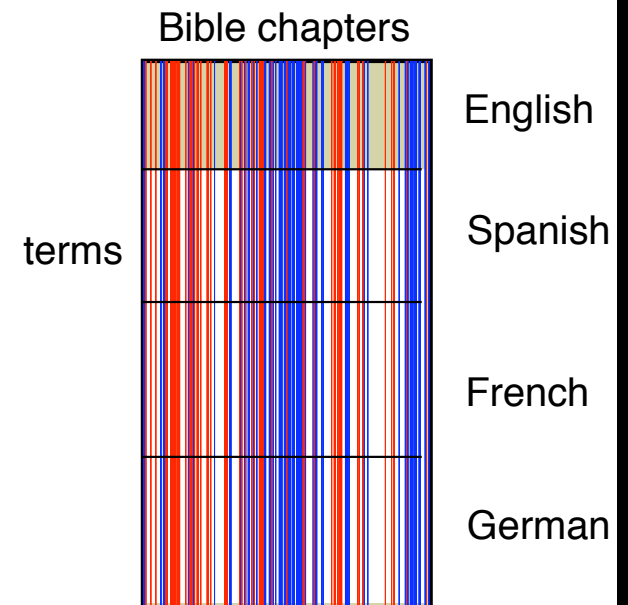
Chapter 88, Verse 1	→	Chapter 88, Verse 3
Chapter 88, Verse 2		Chapter 10, Verse 1
Chapter 88, Verse 3		Chapter 92, Verse 2

Chapter 92, Verse 1	Chapter 10, Verse 2
Chapter 92, Verse 2	Chapter 88, Verse 1

Break up topics by re-distributing their sentences.

## Train on Shuffled English, Test on Foreign

- Project shuffled chapters through SVD.
- Generates new, “shuffled” training data.
- Use the new ensemble to classify the 3x193 chapters in Spanish, French, German.
- Result:
  - Accuracy of 72.0%
  - Still significantly better than the 56.9% baseline.
  - But lower than 74.9%.
- Indicates that some, but not all, of sentiment is bound up in topic.



## Conclusion

- We have demonstrated a supervised machine learning approach to determine sentiment in multilingual documents.
  - Does not require translation
  - Uses a sentiment lexicon only for bootstrapping sentiment labels
  - Uses LSA to project new documents into a language-independent concept space.
  - Uses machine learning on these features to build a predictive model

### Extensions:

- Could easily be used with other topic models, such as LDA or NMF.
- Could be applied to other emotional dimensions or meta-properties, such as “framing language”; prior similar application has been seen in characterizing ideology[8] in multilingual text.

## References

- [1] M. M. Bradley and P. J. Lang, “Affective norms for English words (ANEW): Instruction manual and affective ratings,” *Technical Report C-1, The Center for Research in Psychophysiology University*, 1999.
- [2] K. Denecke, “Using SentiWordNet for multilingual sentiment analysis,” in *ICDE Workshops*. IEEE Computer Society, 2008, pp. 507–512.
- [3] R. Mihalcea, C. Banea, and J. Wiebe, “Learning multilingual subjective language via cross-lingual projections,” in *Proceedings of the 45th Annual Meeting of the Association of Computational Linguistics*, 2007, pp. 976–983.
- [4] A.-L. Ginsca, E. Boros, A. Iftene, D. Trandabat, M. Toader, M. Corici, C.-A. Perez, and D. Cristea, “Sentimatrix – multilingual sentiment analysis service,” in *Proceedings of the 2nd Workshop on Computational Approaches to Subjectivity and Sentiment Analysis (WASSA 2.011)*. Portland, Oregon: Association for Computational Linguistics, June 2011, pp. 189–195. [Online]. Available: <http://www.aclweb.org/anthology/W11-1725>
- [5] P. A. Chew, B. W. Bader, S. Helmreich, A. Abdelali, and S. J. Verzi, “An information-theoretic, vector-space model approach to cross-language information retrieval,” *Journal of Natural Language Engineering*, 2010.
- [6] B. Bader and P. Chew, *Text Mining: Applications and Theory*. Wiley, 2010, ch. Algebraic Techniques for Multilingual Document Clustering.
- [7] —, “Enhancing multilingual latent semantic analysis with term alignment information,” in *COLING 2008*, 2008.
- [8] P. Chew, P. Kegelmeyer, B. Bader, and A. Abdelali, “The knowledge of good and evil: Multilingual ideology classification with PARAFAC2 and machine learning,” *Language Forum*, vol. 34, no. 1, pp. 37–52, 2008.
- [9] P. A. Chew, B. W. Bader, T. G. Kolda, and A. Abdelali, “Cross-language information retrieval using PARAFAC2,” in *KDD '07: Proceedings of the 13th ACM SIGKDD International Conference on Knowledge Discovery and Data Mining*. ACM Press, 2007, pp. 143–152.
- [10] P. Chew and A. Abdelali, “Benefits of the massively parallel rosetta stone: Cross-language information retrieval with over 30 languages,” in *Proceedings of the Association for Computational Linguistics*, 2007, pp. 872–879.