

Use of latent semantic analysis at Otec
presentation at L3S

Hannover

15 May 2006

Jan van Bruggen

Outline

- Background
- LSA in a nutshell
- Applications
 - Positioning and APL
 - Question answering and community formation

Background

- Otec development program
 - Learning networks
 - Bottom-up / data-driven
- Methodology
 - IT tools and technique
 - Psychological research methods
- LSA
 - Data driven
 - Latent variables explain data

Latent semantic analysis is like ...

- Principal component analysis / factor analysis
 - Symmetric matrix **M** - correlations
 - Factor analysis: insert communalities
 - Eigenvalue en eigenvectors
 - **$M = U \Lambda U'$**
 - **Λ** is diagonal matrix with sorted eigenvalues
 - Reproduction: remove smallest eigenvalues in **Λ** and columns and rows in **U** and **U'**
 - Rotate and interpret factor solution

Latent semantic analysis

- Asymmetric matrix (data-matrix)
 - Terms by documents
 - Word frequencies
- Singular value decomposition: $\mathbf{D} = \mathbf{L} \mathbf{S} \mathbf{R}'$
- LSA: reproduction based on a model with less dimensions
- If $\mathbf{M} = \mathbf{D}\mathbf{D}'$ then no difference with PCA

Application areas

- Document retrieval (LSI)
- Cognitive science
 - Semantics of text (Kintsch)
 - Concept learning
- Education
 - Essay rating
 - Selection and sequencing of material
 - Writing

Application domains

Document retrieval

- Large corpora
- Heterogeneous
- High-dimensional

Educational uses

- Small corpora
- Homogeneous
- Low dimensional

Methods:

Minimum size of corpus

Number of dimensions to retain

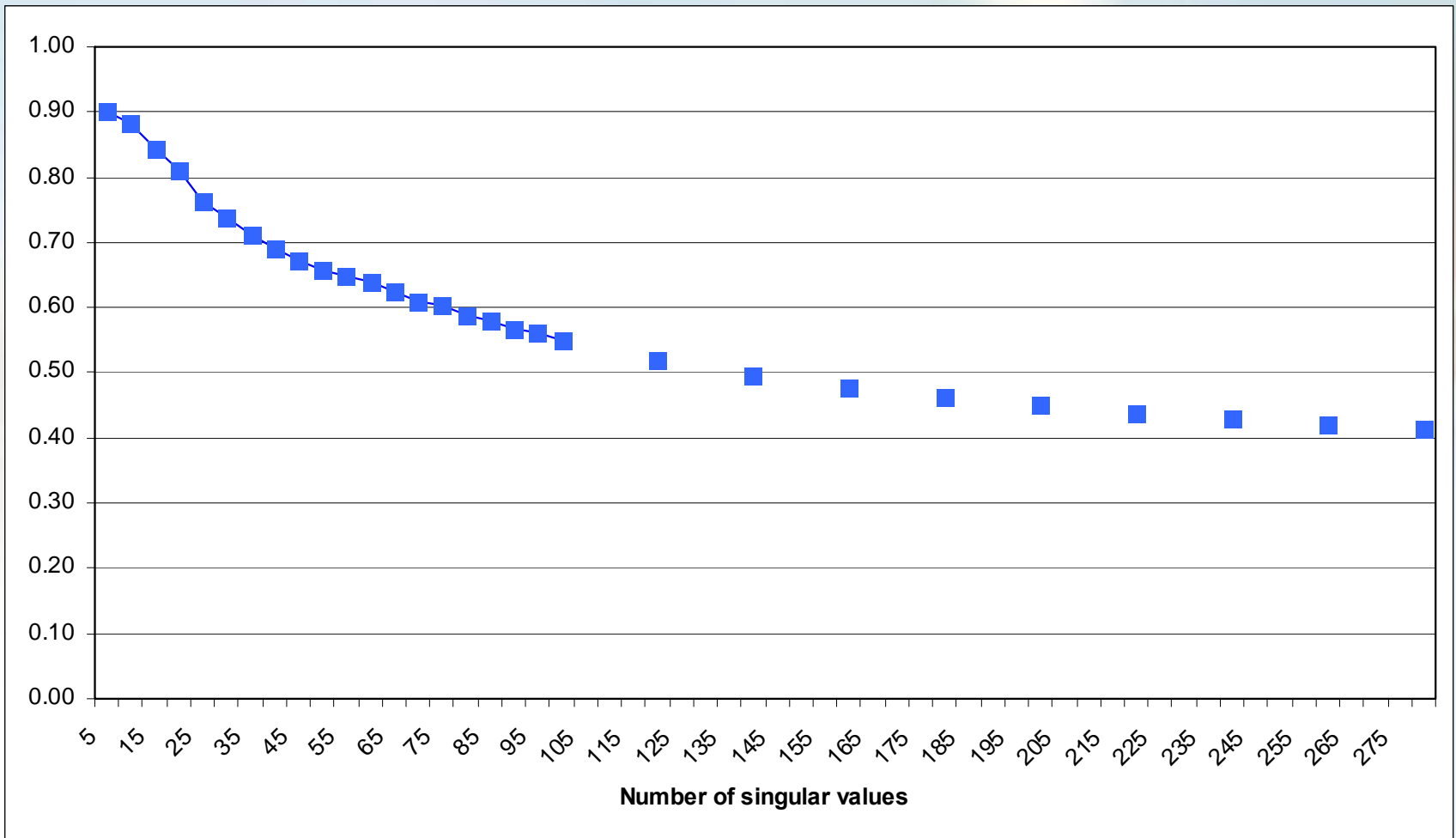
Term frequency measure

Filtering noise words

Methodological considerations

- LSA is not a statistical method
- SVD is a least squares technique
- Error rulez!
 - $x = \text{True}(x) + \text{error}$
 - $\text{covariance}(xy) = \text{covar}(\text{True}(x)\text{True}(y))$
 - $\text{correlation} = \text{cov}(xy) / \text{sqrt}(\text{varx.vary})$
 - Thus: the better you reproduce the data, the worse your correlations get

Example



Explained variance and # SVs

Reliability limits explained variance (.80 is good!)

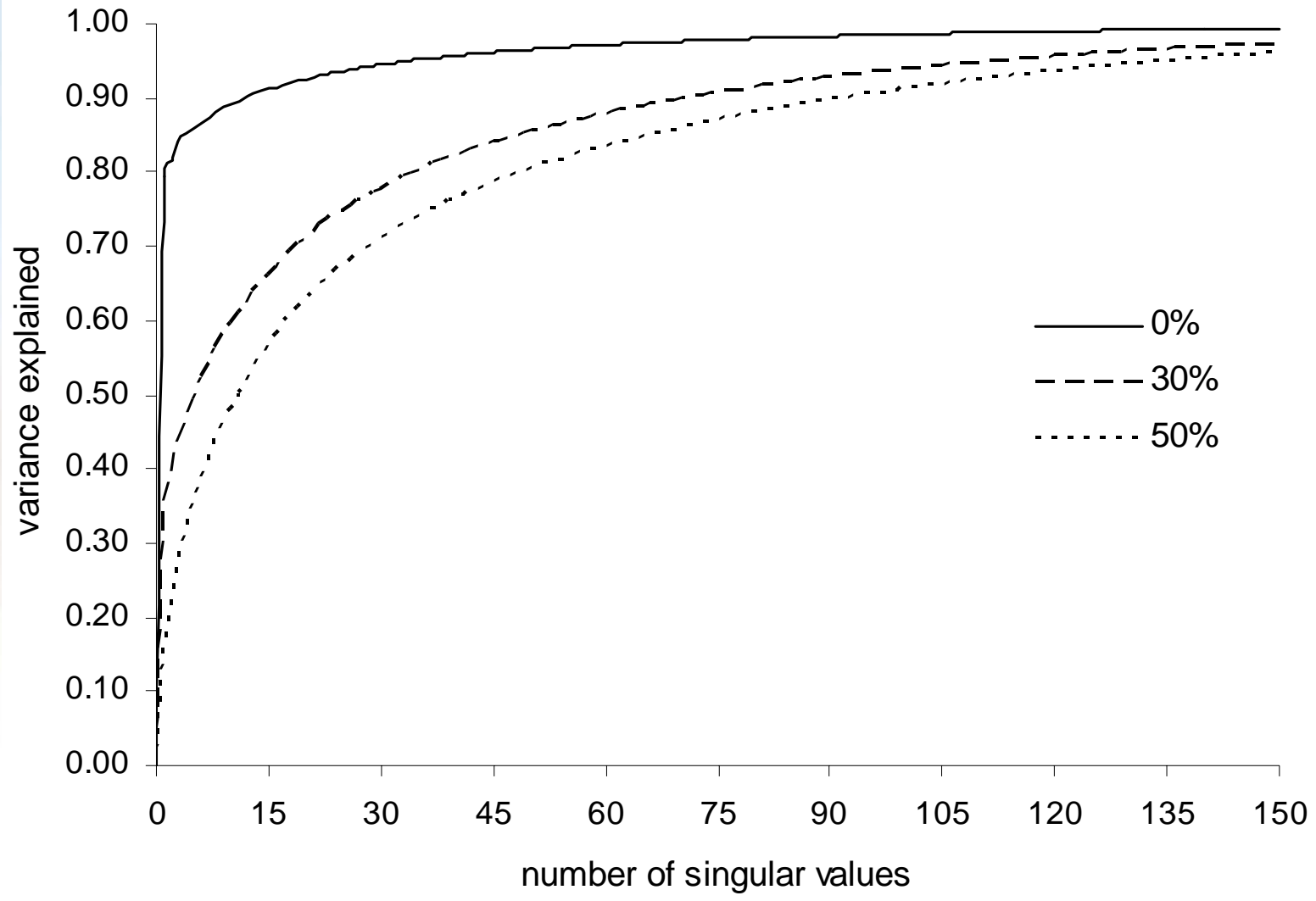
$$\frac{\sum x^2}{n} - M^2 = \sum S^2 = \sum x^2$$

Approximate variance by SSQ(sv)

Determine bandwidth for number of SVs

Still filtering (stopping) is often needed

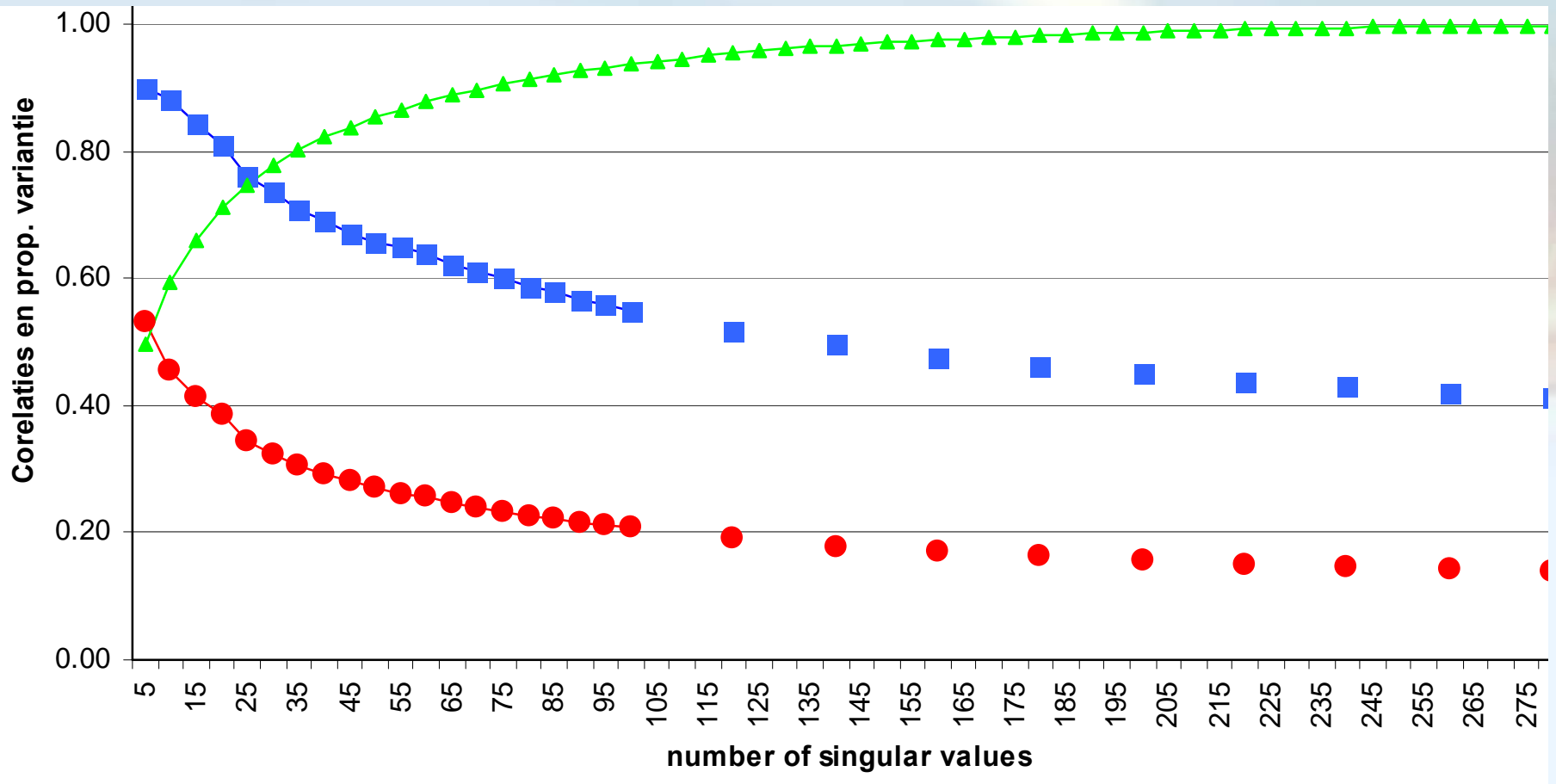
Effect of stopping



Positioning

- Positioning and APL
 - The APL problem
 - Equivalence of outcomes
- Content as a proxy for outcomes
- Performance:
 - Discriminate like an expert
 - Explain sufficient variance
- Trial

Results



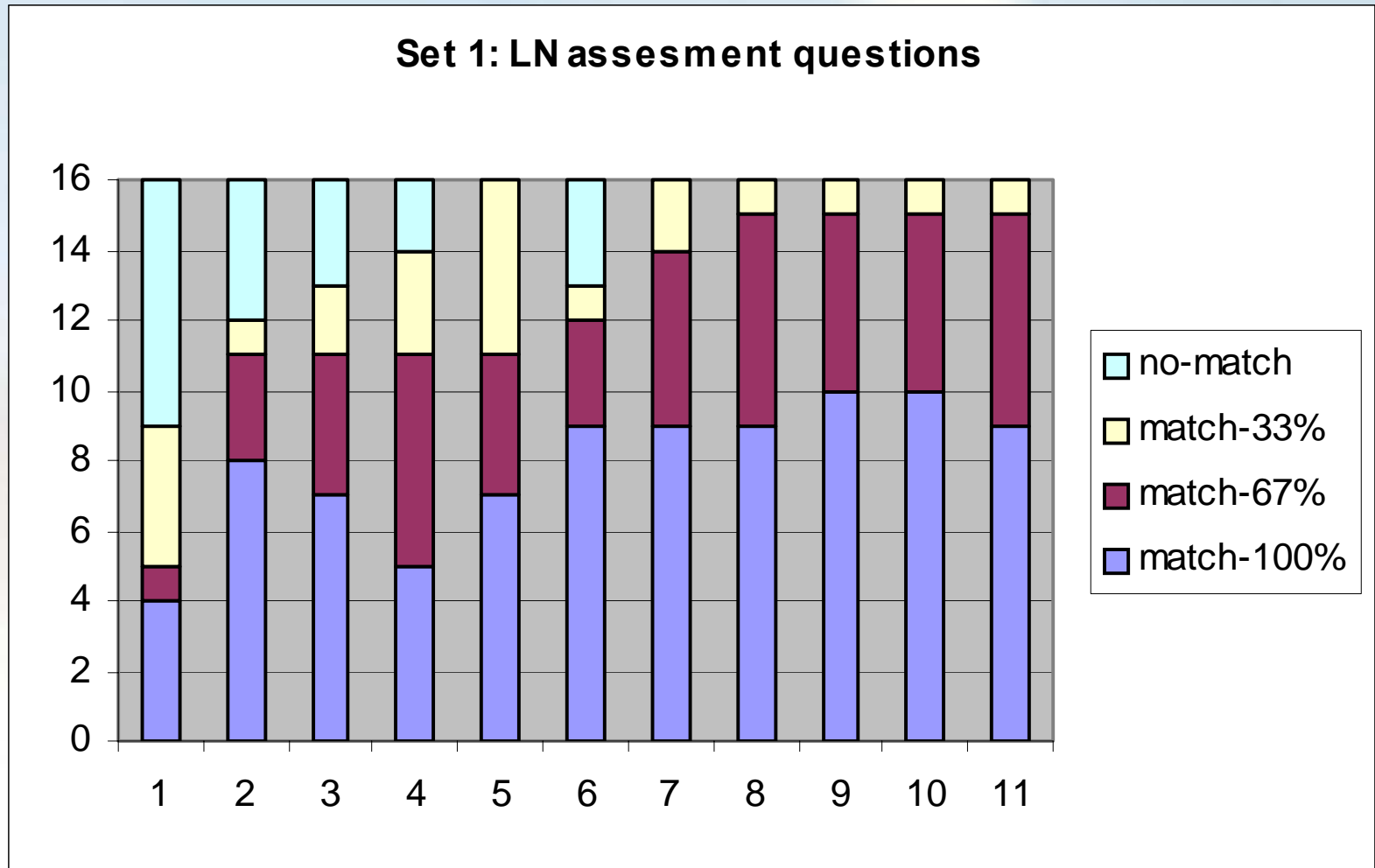
Ad-hoc community formation

- Reduce tutor load
- Question answering by peers
- Find relevant documents
- Find relevant peers
- Create ad-hoc community to answer the question

Methods

- Internet course
- Match assessment questions to three best fitting course elements
- Compare to expert match
- Stepwise optimization

Results



Conclusions

- Preliminary results encouraging
- Optimization itself can be automated
- Extend / Combine techniques:
 - Positioning project Marco Kalz
 - Recommender system Cooper