

DDC Document Classification with SVMs

Bielefeld Conference 2012
Bielefeld, Germany

Tim vor der Brück, Alexander Mehler
Goethe-Universität Frankfurt

April 2012



- 1 DDC Document classification
- 2 Integration in e-Humanities-Desktop
- 3 Conclusion

DDC Document Classification with Support Vector Machines

Task

- Assign topics to documents (n:m)
- Classification of abstracts
- Example applications:
 - Looking for literatur about certain topics
 - Customer Emails
 - Spam/non-spam

Document Classification Scheme: DDC

- DDC: Dewey Decimal Classification
- Hierarchical Classification
- 3 digits and one optional decimal

DDC (contd.)

DDC-Category: 0 6 6 (Organizations in Italy & adjacent islands)

Level 1: 0 (Computer science, information & general works)

Level 2: 06 (Associations, organizations & museums)

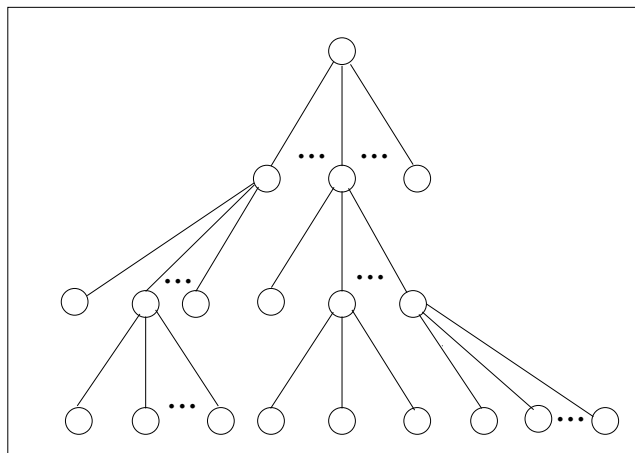
Method

- Extract term/document matrix from text corpus
- Stop words, pruning
- Use weighted term occurrences as features
- Weights: $tf \times idf$, GSS, Leipzig Web Service
- Class/topic is estimated by an SVM
- Feature expansion with hypernyms and synonyms

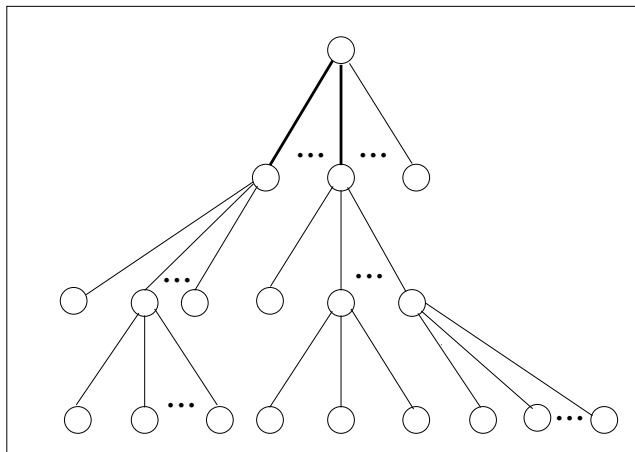
Selection procedure I

- Select most probable options at each level
- recursively repeat the selection process

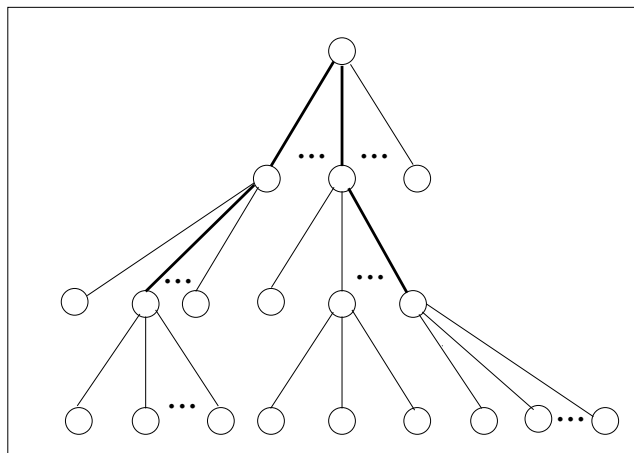
Selection procedure II



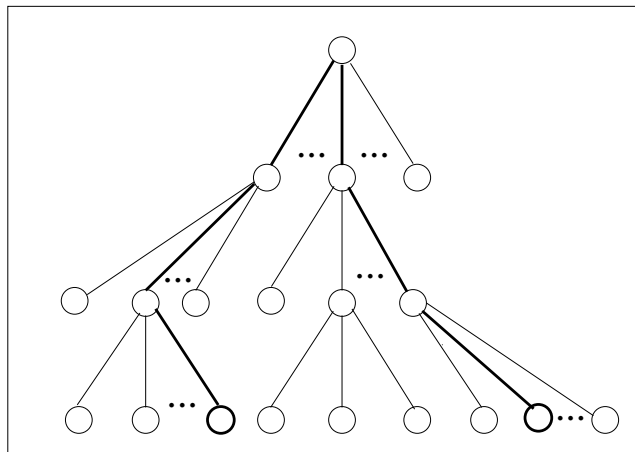
Selection procedure III



Selection procedure IV



Selection procedure V



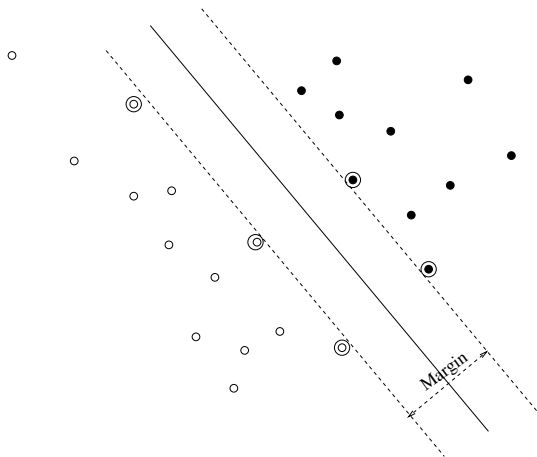
Selection procedure

Predict the category for a new text

- Test one alternative against all others employing the SVM
- Redo the test with all alternatives
- Select the categories with the highest scores

Classification with Support Vector Machine

- Data set is divided in to areas by a separating hyperplane
- Each side is classified identical (document is contained in a category or not)
- Nonlinear classification can be done by applying a kernel function which can transfer the data in higher dimensional space



e-Humanities Desktop

- e-Humanities Desktop is a desktop for people dealing with linguistic corpora
- Words selected by the classifier are marked
- Lemma weights are visualized in a Word Cloud
 - Positive: Lemmas with higher weights are printed larger
 - Inverse: Lemmas with lower weights are printed larger

Selected words

Text View for 367818 WordCloud + WordCloud -

Verflechtung,quantitativen Gesichtspunkt,tatsächliche Lehrbelastung
Name: Corpus; Source: CorpusRes
Statistik über die Studierenden der Universität Konstanz.Studienjahr 2008/2009, 2. Studienabschnitt (Sommersemester).
Statistik über die Studierenden der Universität Konstanz.Studienjahr 2008/2009, 2. Studienabschnitt (Sommersemester).
Alle Statistiken der Studierenden - mit Ausnahme der Statistik Studienfachbelegung nach Abschlusszielen - sind als Kopfstatistik erstellt worden.Damit ist der Aussagewert dieser Statistiken unter quantitativen Gesichtspunkten eingeschränkt.Der Vorteil der Kopfstatistiken liegt in der einfachen Vergleichbarkeit mit Statistiken anderer Hochschulen und staatlichen Statistiken, die überwiegend als Kopfstatistik erstellt werden.Alle

Select

Inverted Word Cloud



Ongoing work

- GSS score
- Feature expansion
- Metaclassifier

Conclusion

Introduced:

- Document classification for DDC
- Classification done my machine learning method (support vector machines)
- Classifier results are visualized in e-Humanitites Desktop

Thank you for your attention.

Any questions?