

Klassifikation von Textabschnitten

Am Beispiel von **Stellenanzeigen**
(**JASC - Job Ads Section Classifier**)



Gliederung

1. Einführung: Zu welchem Zweck machen wir das?
2. Klassifikation – ein kurzer Überblick
3. Umsetzung: Der Job Ads Section Classifier (JASC)
4. Systemdemonstration
5. Probleme und Ausblick



Informationsextraktion aus Jobangeboten

- Auftraggeber: Bundesinstitut für Berufsbildung (BIBB).
- Zuständig u.a. für den Zuschnitt von Lehr- und Ausbildungsberufen
- Datenbank der Bundesanstalt für Arbeit (BfA) über 1.000.000 Stellenzeigen, wachsend
- Strenge Datenschutzauflagen verlangen „Vercodung“
- Empirischer Ansatz: Unterstützung von Expertengruppe durch realwirtschaftliche Daten



Informationsextraktion: Templates

Für unser Kaufhaus in ##### haben wir zum ##.##.#### eine Ausbildungsstelle zum/zur Kaufmann/-frau im Einzelhandel zu besetzen.

Wir erwarten:

- Guter Hauptschulabschluss oder Mittlere Reife
- Gute Deutschkenntnisse in Wort u. Schrift
- Gute bis sehr gute Ausdrucksmöglichkeiten
- Französische Sprachkenntnisse wären von Vorteil

Ihre Aufgaben werden sein:

- Verkauf u. Beratung von Kunden
- Sortimentsgestaltung
- Marketing
- Einkauf und Lagerwesen

Interesse?

Dann freuen wir uns auf Ihre Bewerbung!

Wert	Attribut
	Datum
	Branche
	Beruf
	Gefordert
	Gewünscht
	Aufgaben

Informationsextraktion: Templates

Für unser Kaufhaus in ##### haben wir zum ##.##.####
eine Ausbildungsstelle zum/zur Kaufmann/-frau im
Einzelhandel zu besetzen.

Wir erwarten:

- Guter Hauptschulabschluss oder Mittlere Reife
- Gute Deutschkenntnisse in Wort u. Schrift
- Gute bis sehr gute Ausdrucksmöglichkeiten
- Französische Sprachkenntnisse wären von Vorteil

Ihre Aufgaben werden sein:

- Verkauf u. Beratung von Kunden
- Sortimentsgestaltung
- Marketing
- Einkauf und Lagerwesen

Interesse?
Dann freuen wir uns auf Ihre Bewerbung!

Wert	Attribut
	Datum
	Branche
	Beruf
	Gefordert
	Gewünscht
	Aufgaben

Informationsextraktion: Templates

Für unser Kaufhaus in ##### haben wir zum ##.##.####
eine Ausbildungsstelle zum/zur Kaufmann/-frau im
Einzelhandel zu besetzen.

Wir erwarten:

- Guter Hauptschulabschluss oder Mittlere Reife
- Gute Deutschkenntnisse in Wort u. Schrift
- Gute bis sehr gute Ausdrucksmöglichkeiten
- Französische Sprachkenntnisse wären von Vorteil

Ihre Aufgaben werden sein:

- Verkauf u. Beratung von Kunden
- Sortimentsgestaltung
- Marketing
- Einkauf und Lagerwesen

Interesse?

Dann freuen wir uns auf Ihre Bewerbung!

Wert	Attribut
	Datum
	Branche
Kaufmann/-frau	Beruf
	Gefordert
	Gewünscht
	Aufgaben

Informationsextraktion: Templates

Für unser Kaufhaus in ##### haben wir zum ##.##.####
eine Ausbildungsstelle zum/zur Kaufmann/-frau im
Einzelhandel zu besetzen.

Wir erwarten:

- Guter Hauptschulabschluss oder Mittlere Reife
- Gute Deutschkenntnisse in Wort u. Schrift
- Gute bis sehr gute Ausdrucksmöglichkeiten
- Französische Sprachkenntnisse wären von Vorteil

Ihre Aufgaben werden sein:

- Verkauf u. Beratung von Kunden
- Sortimentsgestaltung
- Marketing
- Einkauf und Lagerwesen

Interesse?
Dann freuen wir uns auf Ihre Bewerbung!

Wert	Attribut
	Datum
	Branche
Kaufmann/-frau	Beruf
	Gefordert
	Gewünscht
	Aufgaben

Informationsextraktion: Templates

Für unser Kaufhaus in ##### haben wir zum ##.##.####
eine Ausbildungsstelle zum/zur Kaufmann/-frau im
Einzelhandel zu besetzen.

Wir erwarten:

- Guter Hauptschulabschluss oder Mittlere Reife
- Gute Deutschkenntnisse in Wort u. Schrift
- Gute bis sehr gute Ausdrucksmöglichkeiten
- Französische Sprachkenntnisse wären von Vorteil

Ihre Aufgaben werden sein:

- Verkauf u. Beratung von Kunden
- Sortimentsgestaltung
- Marketing
- Einkauf und Lagerwesen

Interesse?

Dann freuen wir uns auf Ihre Bewerbung!

Wert	Attribut
##.##.####	Datum
Einzelhandel	Branche
Kaufmann/-frau	Beruf
Hauptschule / ...	Gefordert
Französisch	Gewünscht
Verkauf / ...	Aufgaben

Vorstufe: Klassifikation von Abschnitten

Für unser Kaufhaus in ##### haben wir zum ##.##.#### eine Ausbildungsstelle zum/zur Kaufmann/-frau im Einzelhandel zu besetzen.

Wir erwarten:

- Guter Hauptschulabschluss oder Mittlere Reife
- Gute Deutschkenntnisse in Wort u. Schrift
- Gute bis sehr gute Ausdrucksmöglichkeiten
- Französische Sprachkenntnisse wären von Vorteil

Ihre Aufgaben werden sein:

- Verkauf u. Beratung von Kunden
- Sortimentsgestaltung
- Marketing
- Einkauf und Lagerwesen

Interesse?

Dann freuen wir uns auf Ihre Bewerbung!

Klassen

1: Firma

2: Job

3: Anforderungen

4: Sonstiges



Vorstufe: Klassifikation von Abschnitten

Für unser Kaufhaus in ##### haben wir zum ##.##.#### eine Ausbildungsstelle zum/zur Kaufmann/-frau im Einzelhandel zu besetzen.

Wir erwarten:

- Guter Hauptschulabschluss oder Mittlere Reife
- Gute Deutschkenntnisse in Wort u. Schrift
- Gute bis sehr gute Ausdrucksmöglichkeiten
- Französische Sprachkenntnisse wären von Vorteil

Ihre Aufgaben werden sein:

- Verkauf u. Beratung von Kunden
- Sortimentsgestaltung
- Marketing
- Einkauf und Lagerwesen

Interesse?

Dann freuen wir uns auf Ihre Bewerbung!

Klassen

1: Firma

2: Job

3: Anforderungen

4: Sonstiges



Vorstufe: Klassifikation von Abschnitten

Für unser Kaufhaus in ##### haben wir zum ##.##.#### eine Ausbildungsstelle zum/zur Kaufmann/-frau im Einzelhandel zu besetzen.

Wir erwarten:

- Guter Hauptschulabschluss oder Mittlere Reife
- Gute Deutschkenntnisse in Wort u. Schrift
- Gute bis sehr gute Ausdrucksmöglichkeiten
- Französische Sprachkenntnisse wären von Vorteil

Ihre Aufgaben werden sein:

- Verkauf u. Beratung von Kunden
- Sortimentsgestaltung
- Marketing
- Einkauf und Lagerwesen

Interesse?

Dann freuen wir uns auf Ihre Bewerbung!

Klassen

1: Firma

2: Job

3: Anforderungen

4: Sonstiges



Voraussetzungen und Ziele

- **Was ist vorhanden?**
 - Datenbank mit >1 Mio datengeschützten, nicht ausgezeichneten Stellenanzeigen
- **Was wird gebraucht?**
 - Eine Anbindung eines Klassifikationswerkzeugs an die Datenbank für Stellenanzeigen, dafür benötigt man ...
 - Einen Klassifizierer, der getestet und für brauchbar befunden wurde, dieser benötigt ...
 - Modelle für die Klassen, in die Abschnitte eingeordnet werden sollen, was wiederum voraussetzt ...
 - Ein Textkorpus mit > 1000 anonymisierten, vorklassifizierten Abschnitten.



Automatische Klassifikation - Ansätze

Einordnung von Objekten in vorgegebene Klassen

2 Ansätze:

1. Maschinelles Lernen

- Lernen aus Erfahrung/Beispielen
- Voraussetzung: Trainingskorpus

2. Regelbasierte Klassifikation

- Formulierung fester Regeln



Klassifikation - Grundprinzip

Voraussetzungen für den Vergleich von Objekten

1. Definition von Merkmalen
2. numerische Repräsentation der Merkmale (Vektor)

Beispiel: Bauklötze

M1 = Form → Anzahl der Ecken
M2 = Größe → Gewicht in Gramm



:

Klassifikation - Grundprinzip

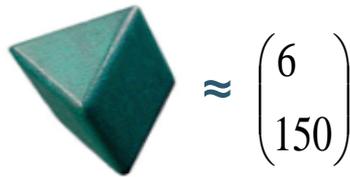
Beispiel: Bauklötze

M1 = Form

→ Anzahl der Ecken

M2 = Größe

→ Gewicht in Gramm



Ähnlichkeit = Vektorähnlichkeit
(z.B. euklidische Distanz oder Cosinus-Ähnlichkeit)

JASC: Workflow

1. Preprocessing & Trainingskorpus
2. Feature Engineering
3. Feature Quantifying
4. Classifying
5. Evaluation
6. Ranking der Verfahren



1. Preprocessing

- **Import der Stellenanzeigen aus der BIBB-Datenbank bzw. CSV-Datei**
 - DatabaseConnector bzw. BIBBReader
- **Zerlegung in Paragraphen (*ClassifyUnits*)**
 - ClassifyUnitSplitter
- **Trainingskorpus erstellen**
(manuelle Klassifikation von ca. 280 Jobangeboten, mit insgesamt etwa 1500 Abschnitten)
 - TrainingdataGenerator



Preprocessing



1	_____
1	_____
1	_____
2	_____
2	_____
2	_____
2	_____
2	_____
3	_____
3	_____
3	_____

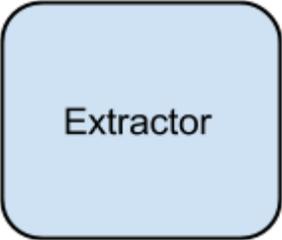
Extracted JobAds

1	_____	
1	_____	
1	_____	
2	_____	
2	_____	
2	_____	
2	_____	
2	_____	
3	_____	
3	_____	
3	_____	

Splitted ClassifyUnits

1	_____	2
1	_____	3
1	_____	4
2	_____	1
2	_____	1
2	_____	2
2	_____	3
2	_____	4
3	_____	1
3	_____	2
3	_____	2

Labeled ClassifyUnits



ClassifyUnit nach Preprocessing

ClassifyUnit

- `int` ID
- `String` content
- `boolean[]` classes



2. Merkmalsextraktion

Schritt 1: Tokenisierung

→ Merkmale = Strings

Schritt 2: Filtern

→ Stoppwörter filtern?

→ ‚unwichtige‘ Wörter filtern?

Schritt 3: Modifizieren

→ Stemming? Normalisieren?

→ NGramme?



ClassifyUnit nach Merkmalsextraktion

ClassifyUnit

- `int` ID
- `String` content
- `boolean[]`
classes
- `List<String>`
FeatureUnits



3. Merkmalsgewichtung

Ziel: numerische Repräsentation der FeatureUnits

4 Möglichkeiten:

1. absolute Häufigkeit
2. relative Häufigkeit (rel. zur Textlänge)
3. Tf-Idf-Wert
4. LogLikelihood-Wert



ClassifyUnit nach Merkmalsgewichtung

ClassifyUnit

- `int` ID
- `String` content
- `boolean[]` classes
- `List<String>`
FeatureUnits
- `double[]`
FeatureVector



4. Classifying

Implementation von vier verschiedener Klassifikatoren

1. K-nearest-neighbor-Klassifikator
2. Rocchio-Klassifikator
3. NaiveBayes-Klassifikator
4. (RegEx-Klassifikator)



Rocchio-Klassifikation

- **Objekte = Feature-Vektoren**
- **Trainingsphase:**
 - Bilde für jede Klasse einen Zentroid-Vektor
(= Mittelwert aller Klassenobjekte)
- **Klassifikation eines neuen Objekts X:**
 - Ordne X der Klasse mit dem nächsten Zentroid-Vektor zu
- **Konfiguration**
 - Distanzmaß



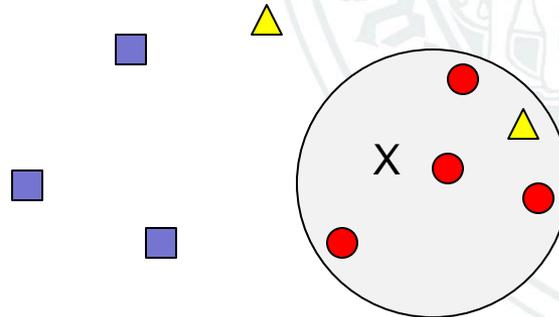
K-nearest-Neighbor-Klassifikation

- **Objekte = Feature-Vektoren**
- **Klassifikation eines neuen Objekts X:**
 1. Ermittle die zu X k nächsten/ähnlichsten Objekte im Trainingskorpus
 2. Ordne X der Klasse zu, die unter den k nächsten Nachbarn am häufigsten vertreten ist
- **Konfiguration:**
 - Distanzmaß
 - Anzahl der Nachbarn (k)



K-nearest-neighbor-Klassifikation

$K = 5$



$X = \bullet$



Naive Bayes Klassifikation

- Arbeitet mit *FeatureUnits* (Objekt = String-List)
- Beruht auf dem Bayes'schen Theorem

$$P(A/B) = \frac{P(B/A) * P(A)}{P(B)}$$

- Berechnet Klassenwahrscheinlichkeiten auf der Grundlage von Wort-Häufigkeiten



5 Evaluation

Zahlreiche Konfigurationsmöglichkeiten:

- Tokens vs. NGramme
- Stoppwortfilter vs. MI-Filter vs. kein Filter
- Stemming?
- Normalisieren?
- Rocchio vs. NaiveBayes vs. KNN-Classifyer
- Distanzmaß?

→ > 5000 Konfigurationsmöglichkeiten

Welche liefert das beste Ergebnis?



5 Evaluation

Kreuzvalidierung

- Aufteilung des Trainingskorpus in Trainings- und Testset
- Vergleich der manuell zugeordneten Klassen mit den maschinell zugeordneten

Evaluationsmaße

- Precision
- Recall
- Accuracy
- F-Measure



Classify and Evaluate

1	100450001001053112000102	2
1	0111323451000000010001230	3
1	0200128881010201234130000	4
2	0000123413000000001234130	1
2	1110000123413000011250003	1
2	1112400000000001234130000	2
2	6101020123461010201234112	3
2	1120008101020123411300222	4
3	0000000012341300234130023	1
3	0123413000001234130000112	2
3	2341300124400023422320000	2

Units with FeatureVectors

1	100450001001053112000102	2	2
1	0111323451000000010001230	3	3
1	0200128881010201234130000	4	1
2	0000123413000000001234130	1	1
2	1110000123413000011250003	1	1
2	1112400000000001234130000	2	2
2	6101020123461010201234112	3	3
2	1120008101020123411300222	4	4
3	0000000012341300234130023	1	2
3	0123413000001234130000112	2	2
3	2341300124400023422320000	2	2

Classified Units

Accuracy	0.94
Precision	0.90
Recall	0.93
F-Score	0.91

Measures



6. Ranking

- mehr als 5000 Klassifikationsergebnisse
- Jeweils 4 Evaluationswerte

→ **Implementation einer Ranking-Komponente**
Erstellt für jeden Evaluationswert eine Ranking-Tabelle



6. Ranking: F-Measure

fscore	prec	rec	accu	Classifier	Dist	Quant	NGrams	MI	SW	Norm
0,9882	0,9913	0,9852	0,9935	KNN(k=1)	COS	LogLike	-2-3	0	false	true
0,9882	0,9913	0,9852	0,9935	KNN(k=1)	COS	LogLike	-2-4	0	false	true
0,9882	0,9925	0,9840	0,9935	KNN(k=1)	COS	LogLike	-2-3-4	0	false	true
0,9615	0,9845	0,9395	0,9790	KNN(k=4)	COS	LogLike	-2-4	0	false	true
0,9190	0,9102	0,9278	0,9530	NaiveBayes	null	null	null	60	true	false
0,9190	0,9102	0,9278	0,9530	NaiveBayes	null	null	null	60	true	false



7. Ergebnisse

- **Das KNN-Verfahren klassifiziert am besten, braucht aber am längsten**
 - F-Score = 0,99 bei $k = 1$ (0,96 bei $k = 4$)
 - Cosinus-Distanz
 - NGramme
 - LogLikelihood
- **Naive Bayes wesentlich performanter und gering schlechter**
 - F-Score = 0,92



Systemvorführung

Applications:

- `DefaultExperimentsGenerator`: Ausführung von 2500 verschiedenen Experimenten mit unterschiedlicher Konfiguration
- `SingleExperimentExecutor`: Ausführung eines zu konfigurierenden Experiments
- `RankResultsApplication`: Ranken der durchgeführten Experimente anhand ihrer Evaluationswerte.
- `ConfigurableDatabaseClassifier`: Klassifizieren von Datensätzen aus einer Datenbank anhand ausgewählter Klassifikatoren.

Code (Repo wird veröffentlicht nach Klärung zum Datenschutz):

<https://github.com/spinfo/bibb-jasc>



Ausblick

- Probleme bei der Übertragung auf die BIBB-Datenbank:
 - Anonymisierte Modelle vs. nicht anonymisierte Klassifikationsdaten
 - Gemischtes Encoding in der Datenbank
 - Zeitbedarf des KNN-Klassifikators
- Möglichkeiten der weiteren Kooperation:
 - Behandlung Encoding, Modelle über nicht-anonymisierte Trainingsdaten
 - Erprobung schnellerer Klassifikatoren
 - Einbeziehung von Metainformationen aus der Datenbank
 - Studie Informationsextraktion
 - Auswertung gesuchter Kompetenzen und Kompetenz-Kombinationen

