

Ruprecht-Karls-Universität Heidelberg
Institut für Informatik
Sommersemester 2010
Seminar: Text Mining
Dozenten: Jannik Strötgen
Prof. Dr. Michael Gertz

Seminararbeit
Coreference of Named Entities

Name:	Philipp Schäfer, Daniel Kruck
Matrikelnummer:	2612579 (Philipp), 2440234 (Daniel)
Studiengang:	Angewandte Informatik (6. Semester)
Email:	trashzopf@googlemail.com (Philipp), daniel.kruck@gmx.net (Daniel)
Datum der Abgabe:	28. Juli 2010

Hiermit versichere ich **Philipp Schäfer, Daniel Kruck**, dass ich die Hausarbeit mit dem Titel **Coreference of Named Entities** im Seminar **Text Mining** im **Sommersemester 2010** bei **Jannik Strötgen** und **Prof. Dr. Michael Gertz** selbstständig und nur mit den in der Arbeit angegebenen Hilfsmitteln verfasst habe. Zitate sowie der Gebrauch fremder Quellen, Texte und Hilfsmittel habe ich nach den Regeln wissenschaftlicher Praxis eindeutig gekennzeichnet. Mir ist bewusst, dass ich fremde Texte und Textpassagen nicht als meine eigenen ausgeben darf und dass ein Verstoß gegen diese Grundregel des wissenschaftlichen Arbeitens als Täuschungs- und Betrugsversuch gilt, der entsprechende Konsequenzen nach sich zieht. Diese bestehen in der Bewertung der Prüfungsleistung mit "nicht ausreichend"(5,0) sowie ggf. weiteren Maßnahmen.

Außerdem bestätige ich, dass diese Arbeit in gleicher oder ähnlicher Form noch in keinem anderen Seminar vorgelegt wurde.

Heidelberg, den 28. Juli 2010

Inhaltsverzeichnis

1	Einleitung	1
2	CEEF	2
2.1	Plausible Annahme über die Verteilung von Koreferenzen . .	2
2.2	Formalisierung	3
2.3	Eliteness	4
2.4	Poisson	5
3	Zusätzliche Bearbeitungsschritte	9
3.1	Automatische Klassifizierung der named entities	9
3.2	Automatische Erkennung von Koreferenzen	11
3.3	Retrieval Method	11

1 Einleitung

In letzter Zeit hat sich das Sortier- und Suchverhalten der Menschheit geändert. Sortierte man früher noch seine Dokumente in Ordner, ist man heute glücklich, wenn man mit leistungsstarken Suchalgorithmen schnell und präzise das gewünschte Dokument findet.

Dabei beschränken sich Suchanfragen nicht nur auf lokale Daten, sondern werden sogar großteils ans Web gestellt. Eine häufige Anfrageform an Suchmaschinen sind hierbei die *named entity queries*[1].



Tabelle 1: Häufige Suchanfrage an Suchmaschinen

Ein Problem bei solchen Suchanfragen ist es, die Häufigkeit der **named entity** in Dokumenten richtig zu erfassen[4]. Denn die Referenzierung des Objektes, das sich hinter einem Eigennamen verbirgt, wird sowohl mit dem Eigennamen selbst, als auch mit kompakteren Ausdrücken vorgenommen. So wird beispielsweise in einem Text, der von Peter Chen handelt, die Person *Peter Chen* auch mit „er“ referenziert.

„**Chen** has recieved serveral awards in the fields of Information Technology. **He** received the Data Resource Management Technology Award [...]“ [2]

In der folgenden Hausarbeit wird ein statistisches Verfahren erklärt, welches die Häufigkeit von Eigennamen in Dokumenten schätzen soll. Grundlage für diese Ausarbeitung ist ein Paper der Herren Na und Ng über „A 2-Poisson Model for Probabilistic Coreference of Named Entities for Improved Text Retrieval“[4].

2 CEEF

2.1 Plausible Annahme über die Verteilung von Koreferenzen

Die Häufigkeit von Eigennamen in Dokumenten ist schwer zu erfassen. Zählt man nur die *named entites* selbst, erhält man die raw entity frequency:

$$tf(e; d) = \#(\text{named entity})$$

Die raw entity frequency beachtet noch keine zusätzlichen Koreferenzen auf das Objekt hinter dem gesuchten Eigennamen. Um später ein besseres Ranking ausführen zu können, möchten wir möglichst alle Referenzen erfassen. Also addieren wir die Häufigkeit der Koreferenzen zu unserer raw entity frequency:

$$tf(e; d) \leq tf_{true}(e; d) = tf(e; d) + \underbrace{atf(e_Q; A, d)}_{\text{Koreferenzen}}$$

Theoretisch sind wir nun am Ziel. Praktisch ist es aber nicht so einfach, die Häufigkeit der Koreferenzen zu bestimmen. Mit Programmen, die Koreferenzen vollständig auflösen, geht ein überdimensionaler Berechnungsaufwand einher. Zudem sind diese Programme noch nicht besonders Präzise und ordnen weniger als 70% der Koreferenzen richtig zu.[4]

Seung-Hoon Na und Hwee Tou Ng untersuchten deswegen einen statistischen Ansatz zur Schätzung der Koreferenzen. Dafür nehmen sie folgendes an:

“Our key assumption is that the frequency of anaphoric expressions is distributed over named entities in a document according to the probabilities of whether the document is elite for the named entities.”[4]

Man geht also davon aus, dass sich mögliche Koreferenzen in einem Dokument auf die verschiedenen Eigennamen verteilen. Die Häufigkeit der Korelation zwischen Koreferenz und Eigennamen hängt dabei von der Wahrscheinlichkeit ab, dass sich ein Dokument um eben diesen Eigennamen dreht.

Ausgehend von dieser Annahme wird nun ein statistisches Modell erstellt, das in der Lage ist, die Häufigkeit von Eigennamen in Dokumenten zu schätzen. Dazu wird der Kontext vernachlässigt und die Anzahl der anaphorischen Ausdrücke mit einem Faktor multipliziert, der die Wahrscheinlichkeit eines

Bezuges auf die gesuchte Named Entity abschätzt.

$$tf_{true}(e_Q; d) \approx tf(e_Q; d) + P(e_Q|A, d) tf(A; d)$$

Im folgenden gilt es diesen Faktor $P(e_Q|A, d)$ zu bestimmen.

2.2 Formalisierung

Bevor wir zur eigentlichen Statistik kommen, legen wir noch die mathematische Schreibweise für diese Aufgabe fest. Zunächst sei Q unsere Query. Die Query besteht hier zur Vereinfachung aus einem einzigen Eigennamen.

Weiter sei e ein beliebiger Eigennamen. Entspricht e der Suchanfrage der Query Q , so nennen wir den Eigennamen gesuchte Entität e_Q . Ist e ein beliebiger nicht gesuchter Eigenname, so nennen wir die Entität nicht gesuchter Eigenname e_N .

Die Koreferenzen r werden *mögliche Koreferenzen* zu einem Eigennamen genannt, wenn sie semantisch zu dem Eigennamen passen. Unterteilen wir beispielsweise die Eigennamen in Objekte und Personen, so sind $\{“he”, “she”\}$ mögliche Koreferenzen zu Personen. Eine mögliche Koreferenz auf Objekte ist $\{“it”\}$. Die Menge aller möglichen Koreferenzen wird A genannt.

Zählt man alle möglichen Koreferenzen A in einem Dokument d , so schreibt man $tf(A; d)$. Die Menge aller möglichen Eigennamen zu einer Menge von Koreferenzen in einem Dokument d wird mit $\epsilon(A; d)$ beschrieben.

- Q = query (Suche)
- e_Q = query entity (gesuchte Entität)
- e_N = non-query entity
- A = Menge plausibler anaphorischer Ausdrücke
- $tf(A; d)$ = Anzahl von A in Dokument d
- $\epsilon(A; d)$ = Menge plausibler Entitäten in Dokument d
- $tf(e; d)$ = raw entity frequency
- $atf(e; A, d)$ = anaphoric entity frequency

Ausgerüstet mit einer einheitlichen Schreibweise wird zunächst Eliteness diskutiert und dann der Kern des Ganzen - das Poisson-Modell für statistische Häufigkeit von Koreferenzen - hergeleitet.

2.3 Eliteness

Definition:

A document is elite for a term if the document is “about” the concept represented by the term [3].

Bei genauerer Begutachtung der Definition liegt die Annahme nahe, dass die Bezüge der anaphorischen Ausdrücke auf die Entitäten im Zusammenhang stehen, ob ein Dokument bezüglich dieser Entität „elite“ ist, oder nicht. Wir definieren uns also einen weiteren Faktor:

Definition:

$P(\mathbf{E}(e) = 1|d)$ ist die Wahrscheinlichkeit, dass ein Dokument für eine Entität e „elite“ ist.

Jetzt können wir folgendes setzen:

$$P(e_Q|A, d) = \frac{P(\mathbf{E}(e) = 1|d)}{\sum_{e \in \mathcal{E}(A; d)} P(\mathbf{E}(e) = 1|d)}$$

Wir machen also nichts anderes, als die Wahrscheinlichkeit der Eliteness des Dokumentes bezogen auf unsere gesuchte Entität in Relation mit den Wahrscheinlichkeiten der anderen Named Entities zu stellen.

Jetzt vereinfachen wir die Formel noch ein wenig, indem wir zuerst alle non-query Entitäten nach der Wahrscheinlichkeit der Eliteness sortieren. Danach werden die K non-query Entitäten mit der höchsten Wahrscheinlichkeit selektiert, sodass sich unsere Formel auf Folgendes beschränkt:

$$P(e_Q|A, d) \approx \frac{P(\mathbf{E}(e) = 1|d)}{P(\mathbf{E}(e) = 1|d) + \sum_{i=1}^K P(\mathbf{E}(e_n^{(i)}) = 1|d)}$$

K ist in dem Fall frei gewählt und kann entsprechend angepasst werden. So dass tatsächlich nur Wahrscheinlichkeiten vernachlässigt werden, die nicht ins Gewicht fallen.

Nun suchen wir uns eine repräsentative non-query Entität heraus und setzen die Wahrscheinlichkeit der Eliteness aller anderen K Entitäten mit der

Wahrscheinlichkeit der Herausgesuchten gleich. Somit erhalten wir eine stark vereinfachte, aber trotzdem recht genaue Alternative der obigen Formel.

$$P(e_Q|A, d) \approx \frac{P(\mathbf{E}(e) = 1|d)}{P(\mathbf{E}(e) = 1|d) + K \cdot P(\mathbf{E}(e_N) = 1|d)}$$

Im Folgenden werden wir eine performante und einfache Methode beschreiben, mit der es möglich ist, den Faktor $P(\mathbf{E}(e) = 1|d)$ zu berechnen.

2.4 Poisson

Nochmal: die grundlegende Annahme von Seung-Hoon Na und Hwee Tou Ng lautet:

“Our key assumption is that the frequency of anaphoric expressions is distributed over named entities in a document according to the probabilities of whether the document is elite for the named entities.”[4]

Um diese Annahme in ein mathematisches Modell zu übersetzen, muss man die Wahrscheinlichkeit abschätzen, dass ein Dokument “elite” ist. Wir suchen also eine Formel für

$$P(E(e) = 1|d) = ?$$

Zunächst interessiert uns hierfür die Wahrscheinlichkeit, dass eine beliebige Entität tf -fach in einem Dokument vorkommt. Dafür wird ein 2-Poisson Mixture Modell herangezogen.

$$P(tf) = \pi_e \frac{e^{-\lambda_e} \lambda_e^{tf}}{tf!} + (1 - \pi_e) \frac{e^{-\mu_e} \mu_e^{tf}}{tf!} \quad (1)$$

Der erste Term repräsentiert hierbei die Wahrscheinlichkeit des tf -fachen Auftretens der Entität in der Funktion als “elite”-Eigennamen, der zweite Term steht für die Wahrscheinlichkeit des tf -fachen Auftretens der Entität als “non-elite”-Eigennamen.

$$P(tf) = \underbrace{\pi_e \frac{e^{-\lambda_e} \lambda_e^{tf}}{tf!}}_{\text{Wahrscheinlichkeitsanteil—elite}} + \underbrace{(1 - \pi_e) \frac{e^{-\mu_e} \mu_e^{tf}}{tf!}}_{\text{Wahrscheinlichkeitsanteil—nonelite}}$$

Dabei deutet π_e auf den Erwartungswert hin, ob ein Dokument “elite” ist oder nicht.

π_e , λ_e und μ_e werden festgelegt als:

$$\pi_e = \frac{df(e)}{df(e) + N}, \quad \lambda_e = \frac{cf(e)}{df(e)}, \quad \mu_e = \frac{cf(e)}{N}$$

wobei

- N : Anzahl der Dokumente
- $df(e)$: Anzahl der Dokumente, in denen e vorkommt
- $cf(e)$: Anzahl von e in allen Dokumenten

Damit haben wir ein Modell, dass die Wahrscheinlichkeit der Häufigkeit von Entitäten in Dokumenten schätzt. Wichtig ist hierbei, dass sowohl das Auftreten von “elite”-Entitäten in Dokumenten mit einer Poissonverteilung modelliert wird, als auch das “zufällige” Auftreten der “nonelite”-Entitäten. Eine einfache Poissonverteilung sieht in Formel und im Diagramm wie folgt aus:

$$P_\lambda(k) = \frac{\lambda^k}{k!} e^{-\lambda}$$

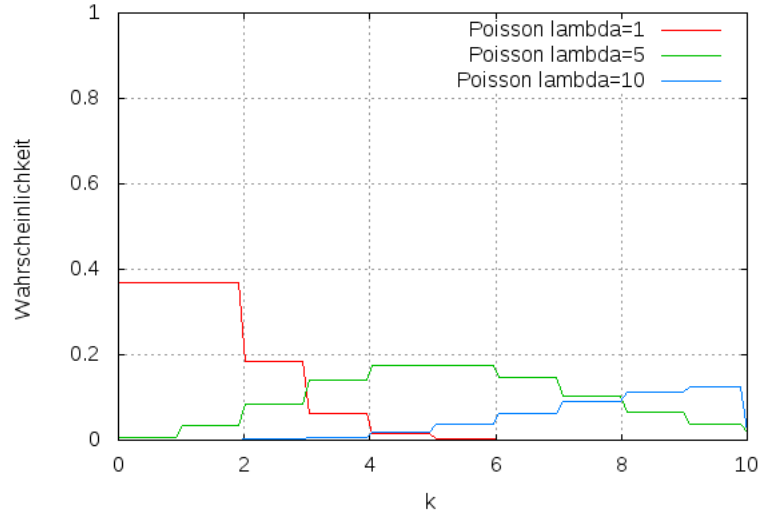


Tabelle 2: Einfache Poissonverteilung

Mit Hilfe von (1) können wir nun $P(\mathbf{E}(e) = 1|d)$ bestimmen. Wir setzen

$$P(\mathbf{E}(e) = 1|d) = \frac{\pi_e P(tf(e; d) | \mathbf{E}(e) = 1)}{\pi_e P(tf(e; d) | \mathbf{E}(e) = 1) + (1 - \pi_e) P(tf(e; d) | \mathbf{E}(e) = 0)}$$

mit

$$\begin{aligned} P(tf(e; d) | \mathbf{E}(e) = 1) &= \frac{e^{-\lambda_e} \lambda_e^{tf}}{tf!} \\ P(tf(e; d) | \mathbf{E}(e) = 0) &= \frac{e^{-\mu_e} \mu_e^{tf}}{tf!} \end{aligned}$$

und erhalten so

$$\begin{aligned} P(\mathbf{E}(e) = 1 | d) &= \frac{\pi_e P(tf(e; d) | \mathbf{E}(e) = 1)}{\pi_e P(tf(e; d) | \mathbf{E}(e) = 1) + (1 - \pi_e) P(tf(e; d) | \mathbf{E}(e) = 0)} \\ &= \frac{\pi_e \frac{e^{-\lambda_e} \lambda_e^{tf}}{tf!}}{\pi_e \frac{e^{-\lambda_e} \lambda_e^{tf}}{tf!} + (1 - \pi_e) \frac{e^{-\mu_e} \mu_e^{tf}}{tf!}} \\ &= \frac{\pi_e e^{-\lambda_e} \lambda_e^{tf}}{\pi_e e^{-\lambda_e} \lambda_e^{tf} + (1 - \pi_e) e^{-\mu_e} \mu_e^{tf}} \\ &= \frac{\pi_e}{\pi_e + (1 - \pi_e) \frac{e^{-\mu_e} \mu_e^{tf}}{e^{-\lambda_e} \lambda_e^{tf}}} \\ &= \frac{\pi_e}{\pi_e + (1 - \pi_e) e^{\lambda_e - \mu_e} \left(\frac{\mu_e}{\lambda_e} \right)^{tf(e; d)}} \end{aligned}$$

Jetzt werden 2 Fälle unterschieden:

1. $e = e_Q$

Wir machen die plausible Annahme, dass $tf(e; d) \geq 1$, wenn das Dokument d elite ist für die gesuchte Entität e . Alle anderen Dokumente im Test seien non-elite bezüglich e .

Dann können wir die π_e, λ_e und μ_e wie oben beschrieben ersetzen und erhalten

$$\begin{aligned} P(\mathbf{E}(e_Q) = 1 | d) &= \frac{\pi_e}{\pi_e + (1 - \pi_e) e^{\lambda_e - \mu_e} \left(\frac{\mu_e}{\lambda_e} \right)^{tf(e; d)}} \\ &= \frac{\frac{df(e)}{df(e) + N}}{\frac{df(e)}{df(e) + N} + \left(1 - \frac{df(e)}{df(e) + N} \right) e^{\frac{cf(e)}{df(e)} - \frac{cf(e)}{N}} \left(\frac{\frac{cf(e)}{N}}{\frac{cf(e)}{df(e)}} \right)^{tf(e; d)}} \\ &= \frac{1}{1 + \left(\frac{df(e)}{N} \right)^{tf(e; d) - 1} e^{\frac{cf(e)}{df(e)} - \frac{cf(e)}{N}}} \end{aligned}$$

2. $e \neq e_Q$

Auch hier machen wir zu Beginn zwei Annahmen. Zum einen seien alle non-query Entitäten vom selben Typ, wie die query Entität und zum anderen sei die Termfrequency ($tf(e; d)$) einer Top- K non-query Entität gleich dem Durchschnittsvorkommen einer query Entität in einem elite Dokument.

Unter diesen Vorraussetzungen dürfen die λ_e, μ_e und π_e wie oben gewählt werden. Dies führt uns erneut zu der Formel

$$P(\mathbf{E}(e) = 1|d) = \frac{1}{1 + \left(\frac{df(e)}{N}\right)^{tf(e;d)-1} e^{\frac{cf(e)}{df(e)} - \frac{cf(e)}{N}}}$$

Aufgrund unserer 2. Annahme gilt jetzt aber, dass $tf(e; d) = \frac{cf(e)}{df(e)}$, weshalb wir in diesem Fall zu dem Ergebnis kommen, dass

$$P(\mathbf{E}(e_N) = 1|d) = \frac{1}{1 + \left(\frac{df(e)}{N}\right)^{\frac{cf(e)}{df(e)}-1} e^{\frac{cf(e)}{df(e)} - \frac{cf(e)}{N}}}$$

Mit diesen Ergebnissen können wir jetzt unseren Wahrscheinlichkeitsfaktor

$$P(e_Q|A, d) \approx \frac{P(\mathbf{E}(e) = 1|d)}{P(\mathbf{E}(e) = 1|d) + K \cdot P(\mathbf{E}(e_N) = 1|d)}$$

abschätzen und damit die gesuchte wahre Termfrequency

$$tf_{true}(e_Q; d) \approx tf(e_Q; d) + P(e_Q|a, d)tf(A; d)$$

Wir wollen zum Schluss noch darauf aufmerksam machen, dass es sinnvoller ist, statt der normalen Termfrequency, die so genannte „normalized Termfrequency“ zu verwenden.

Definition: normalized Termfrequency

$$ntf(e; d) = \frac{tf(e; d) \cdot avglen}{len(d)}$$

Wobei *avglen* die Durchschnittslänge eines Dokuments im Test ist und *len(d)* die Länge des aktuellen Dokuments.

Die normalized Termfrequency setzt also das Vorkommen der Entitäten in Relation mit der Länge der Dokumente. Das ist wichtig, da die Termfrequency in langen Dokumenten ja meist höher ist als in kurzen, dies aber das Dokument nicht unbedingt relevanter macht. Somit setzen wir also

$$ntf_{true}(e_Q; d) \approx ntf(e_Q; d) + P(e_Q|a, d)ntf(A; d)$$

und sind damit fertig.

3 Zusätzliche Bearbeitungsschritte

3.1 Automatische Klassifizierung der named entities

Nachdem man nun die Wahrscheinlichkeit schätzen kann, mit der ein Dokument “elite” ist, benötigt man nun noch die Häufigkeit der möglichen Koreferenzen.

Dafür müssen wir zunächst herausfinden, von welchem Typ die Entity Query ist. Für Personen wird

$$A_P = \{ 'he', 'she', 'his', 'her', 'himself', 'herself' \}$$

als sinnvolle Menge von Koreferenzen angesehen. Für Objekte wird eine erweiterbare Menge an Koreferenzen um

$$A_0 = \{ 'it', 'its', \dots \}$$

bei Bedarf verwendet. Das hat den Grund, dass Objekte oft mit Nomen referenziert werden. So wird beispielsweise “Google” eventuell mit “Suchmaschine”, “Internetriese” oder ähnlichem bezeichnet.

Für die Unterscheidung zwischen Person und Objekt wird ein Support Vector Machine Klassifikator eingesetzt.

Dafür werden zunächst Entitäten aus einer Trainingsmenge betrachtet. Bei diesen Eigennamen wird dem Computer als zusätzliche Information mitgeteilt, ob es sich um eine Person oder um ein Objekt handelt. Der Computer sucht dann die besten M Treffer zu den jeweiligen Eigennamen herausfinden und errechnet charakteristische Werte für diesen Eigennamen. Da wir hier herausbekommen wollen, ob es sich um ein Objekt oder um eine Per-

son handelt, werden bei den Top M Einträgen die Koreferenzen betrachtet. Kommt häufig $A_P = \{\text{'he'}, \text{'she'}, \text{'his'}, \text{'her'}, \text{'himself'}, \text{'herself'}\}$ vor, so wird es wahrscheinlicher, dass es sich die betrachtete Entität eine Person ist. Kommt $A_O = \{\text{'it'}, \text{'its'}\}$ häufig vor, so handelt es sich wahrscheinlich um ein Objekt. Die Quantifizierung erfolgt normiert nach:

$$f_P(Q) = \frac{1}{|F(Q)|} \sum_{d \in F(Q)} \frac{tf(A_P; d)}{len(d)}$$

$$f_O(Q) = \frac{1}{|F(Q)|} \sum_{d \in F(Q)} \frac{tf(A_O; d)}{len(d)}$$

Dabei stellt $F(Q)$ die Menge der besten M Treffer dar, $|F(Q)|$ die Kardinalität dieser Menge, $tf(A_P; d)$ die Häufigkeit der Koreferenzen auf Personen in dem jeweiligen Dokument und $len(d)$ gibt die Anzahl der Worte eines Dokumentes aus.

Der Computer errechnet $f_O(Q)$ und $f_P(Q)$ für jeden Eigennamen der Trainingsmenge. Visualisiert könnte ein Diagramm dazu wie folgt aussehen:

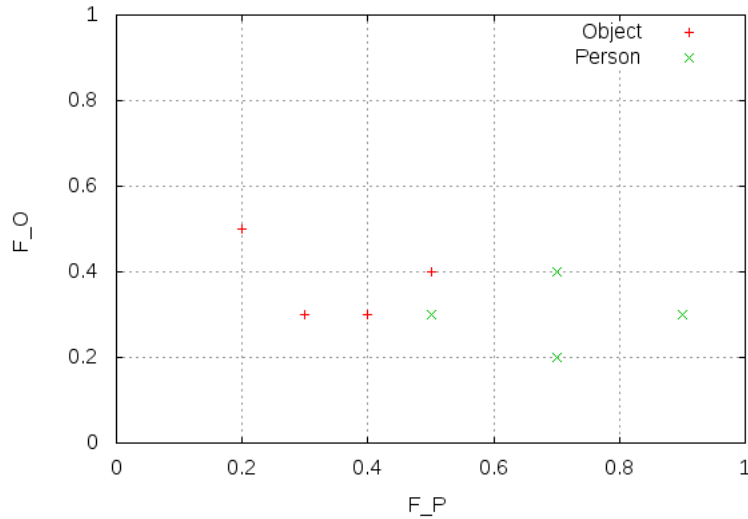


Tabelle 3: Mögliche Visualisierung der SVM nach dem Trainingsschritt ohne Hyperebene

Nachdem sich die Eigenschaften von Personen und Objekten hoffentlich etwas unterscheiden und die verschiedenen Typen im Diagramm getrennt er-

scheinen, muss man nun noch eine Hyperebene mit möglichst großem Abstand zum Rand zwischen den aufgetrennten Typen einfügen. Dafür wurde im Paper über “A 2-Poisson Model for Probabilistic Coreference of Named Entities for Improved Text Retrieval” [4] eine Radial Basis Function herangezogen.

3.2 Automatische Erkennung von Koreferenzen

Da die Menge an Koreferenzen

$$A_O = \{\text{'it'}, \text{'its'}\}$$

für Objekte nicht ausreichend ist, haben die Autoren des Papers “A 2-Poisson Model for Probabilistic Coreference of Named Entities for Improved Text Retrieval”[4]

3.3 Retrieval Method

Literatur

- [1] Guha, R. und A. Garg. Disambiguating people in search. *In TAP: Building the Semantic Web*, 2003.
- [2] Modified: 21:02, 21 May 2010 Nineball. Peter chen, 2010.
- [3] Robertson, S. E. und S. Walker. Some simple effective approximations to the 2-poisson model for probabilistic weighted retrieval. *In SIGIR'94*, Seite 232–241, 1994.
- [4] Seung-Hoon Na, Hwee Tou Ng. A 2-poisson model for probabilistic co-reference of named entities for improved text retrieval. *Annual ACM Conference on Research and Development in Information Retrieval*, 2009.