

ZUSAMMENFASSUNGEN COLLOQUIUM, REPETITORIUM

Pascal Guldener

19.6.2017

1 Colloquium

1.1 Musik und Wittgenstein - Semantische Suche Im Wittfind

Die Bachelorarbeit von Ines Röhrer betreut von Dr. Hadersbeck, hatte zum Ziel das vom CIS entwickelten und betreuten Tool, Wittfind, eine Semantische Kategorie zur Suche im Nachlass von Ludwig Wittgenstein hinzuzufügen. Die Motivation ist die Annahme, dass Musik im Leben und Werk des Sprachphilosophen eine einschlägige Rolle gespielt habe und auch das in seinem geistigen Nachlass widerspiegele. Dies motivierte die Autorin zur Zusammenarbeit mit einem Musikwissenschaftler zur Identifizierung der relevantesten Begriffe aus der Domäne Musik im Nachlass von Wittgenstein. Zunächst erläuterte Sie die Verwendung semantischer Kategorien im Wittfind-tool anhand von der Kategorie "Farbe" welche eine Klasse von farbbeischreibenden Begriffen bildet. Als Ihre Aufgabe beschrieb Sie nun, diese Klasse analog für musikalische Begriffe zu definieren und eventuell auch eine Ontologie zu entwickeln die es erlaubt Relationen zwischen diesen Begriffen zu finden.

Die Basis ihrer Arbeit bildete alle Vollformenextraktion aller relevanten Begriffe. In Ermangelung geeigneter Tools hatte sie sich dafür ein eigenes Vollformenlexikon geschrieben welches aus Endungslisten und denzugehörigen Lemmata der Kategorie bestand. Während für Sie die Nachteile und Grenzen dieses Verfahrens auf der Hand liegen, stellte dies eine vorläufig akzeptable Grundlage für Ihre Arbeit dar. Als weitere große Herausforderung stellte sich die Disambiguierung heraus. Frau Röhrer erläuterte dies anhand der verschiedenen Verwendungen des Buchstabens "C". Dieser kann als Variable in einem mathematischen Kontext genauso auftauchen wie bei der Benennung des Tons in einem musikalischen Kontext. Um die Kontextsensitivität bei der Disambiguierung solcher Begriffe zu auszunutzen implementierte sie einen Ringpuffer in welchem auf semantisch verwandte Begriffe in einen gewissen Bestand geprüft werden kann.

Der grösste Teil des praktischen Teils bezog sich auf die technischen Herausforderungen, insbesondere die GUI des Wittfind um die notwendigen JavaScript-Komponenten zu erweitern welche für die praktische Verwendung nötig sind. Dabei kamen auch Unzulänglichkeiten in der bevorzugten Darstellung (linear in Abhängigkeit zu Ihrer Häufigkeit) dieser "Wortfelder", bzw. Klassen, als Wortcloud zur Sprache. Viele Wörter kamen so selten vor, dass sie kaum mehr wahrnehmbar dargestellt wurden während andere den Großteil des verfügbaren Platzes einnahmen.

Aus Zeitmangel kam es allerdings nicht mehr zur Ausarbeitung einer Ontologie und Modellierung der Relationen von Begrifflichkeiten zueinander.

1.2 Machine-Learning basierte automatische OCR-Korrektur

Die Bachelorarbeit von Michael Strohmeier, betreut von Prof. Schulz, behandelt die Frage ob sich das Ergebnis einer OCR Korrektur durch mit Hilfe einer Klassifizierung durch ein Machine Learning System verbessern lässt. Als features verwendete er dabei die Ausgabe des Profilers des von ihm verwendeten OCR-Systems. OCR-Systeme werden vor mannigfaltige Probleme gestellt, dazu gehören nicht nur die Erkennung von Handschriften oder Robustheit gegenüber verschiedenen

Fonts sondern auch die Veränderung der korrekten Schreibweise eines Worts über die Zeit. Diese Unwägbarkeiten machen es sehr schwierig ein System zu bauen welches keiner Korrektur durch Menschen bedarf. Die Ausgabe des verwendeten OCR-Systems bestand aus einer Liste von Korrekturvorschlägen mit Konfidenzwert und weiteren Informationen, wie dem Levenstheinabstand. Während diese vom OCR-System als interaktive Liste zur manuellen Auswahl zur Verfügung gestellt werden, wurden nun ein Teil dieser Informationen als Features für den Klassifizierer verwendet. Weitere hinzugefügte Features waren die Längendifferenz, der Konfidenzwert des zweitgerankten Korrekturvorschlags und die Häufigkeit des betreffenden Wortes. Als Klassifikatoren wurden ein Naive-Bayes Klassifikator aus der Scit-Learn Library sowie eine durch die Libsvm bereitgestellte Support Vector Machine verwendet. Das Goldstandard-Datenset bestand aus den Korpora "Pardiesgärtlein", "Curioser Botanicus" und dem aus 33 Kräuterkundetexten bestehenden RIDGES Korpus, welche vom CIS in Kooperation mit der Humboldt Universität Berlin erstellt wurden. Herr Strohmeier berichtete von anfangs eher schlechten Resultaten, die er allerdings auf eine Verkürzung des Konfidenzwerts in der Ausgabe des OCR-Systems zurückführen konnte. Einerseits konnte dies durch Kreuzevaluierung, also der Einführung eines Developmentdatensets, ausgleichen, andererseits hat sich dadurch auch die starke Abhängigkeit vom Feature des vorberechneten Konfidenzwerts gezeigt. Auf Nachfrage an Herrn Schulz ob man daher nicht das OCR-System selbst bereits als (Multiclass) Klassifikator anzusehen sei, äusserte dieser sich skeptisch und verwies auf die geringe Zuverlässigkeit, welche ein Beibehalten des zu korrigierenden Wortes oft als die bessere Wahl erscheinen lässt.

Beim Vergleich der beiden Klassifikatoren zeigten sich sehr grosse Unterschiede in der Performance was die Laufzeiten bei Training und Anwendung angehen. Während die SVM 12s für das Training und 3.26s für eine Binär-Klassifizierung benötigte brauchte der Naive-Bayes Klassifikator dafür 0.32s und 0.009s. Da die SVM bessere Ergebnisse lieferte wurde Sie für die Evaluation herangezogen. Vergleiche über Basisfeatures und extrahierten Features zeigten einen deutlichen Performancegewinn mit einem finalen F-Score von über 80%. Herr Strohmeiers Fehleranalyse ergab, dass die Abdeckung der zu klassifizierenden Wörter durch den Goldstandard sowie die Qualität des Konfidenzwerts entscheidende Kriterien für die Performance der Klassifikatoren waren.

2 Repetitorium

2.1 Versionierung mit Git

Herr Roth ging auf die Entstehungsgeschichte von Git, insbesondere die Rolle von Git bei der Entwicklung des Linux-Kernels, Unterschiede zu älteren Versionierungssystemen und praktische Verwendung ein. Insbesondere illustrierte er die verschiedenen Zustände die eine Datei in Git (add,stage,commit) annehmen kann und erklärte die fundamentalen Konzepte von Branching und Merging. Ausserdem empfahl er publickey authentication und erläuterte grob deren Einrichtung. Zuletzt führte er in grundlegende git Befehle wie git rm , git diff und die Funktionalität von .gitignore ein. Auch die Verwendung von diffs zum Vergleich alter mit neueren Versionen und die Commit-History wurden knapp erklärt.