

Projektplan

Bereich	Lösung	Erläuterung	Material, Ideen, Bibliotheken
Thema	Tippfehlerkorrektur/Autokorrektur	Tippfehler in freien Texten sollen erkannt und korrigiert werden, die korrekten Wörter müssen dazu im Vorhinein bereits vorliegen/ eingelesen werden damit im Nachhinein festgestellt werden kann, wie erfolgreich die Korrektur war. Es soll ein Vergleich zwischen verschiedenen Algorithmen stattfinden.	
Benötigte Daten	Lexikon	Eigenes kleines Lexikon, dieses ist dann erweiterbar <u>Alternative:</u> Es muss ein Lexikon vorhanden sein, auf das zugegriffen werden kann, um die falschen Wörter zu erkennen und sie durch das (am wahrscheinlichsten) richtige Wort zu ersetzen	<u>Methode:</u> Lexikon in das Programm einlesen Lexikon liegt in Form einer .txt-Datei vor https://wordnet.princeton.edu/ http://elexicon.wustl.edu/ oder mittels PyEnchant Wörter als falsch erkennen lassen
	Korrekte Texte (zu Anfang sollen die Trainingstexte eher Kurz sein, also ca.100 – 200 Wörter lang)	Es muss ein korrekter Ausgangstext ohne Fehler vorliegen. Dieser wird benötigt um den Erfolg der Algorithmen zu berechnen → Wie viele Fehler wurden erkannt → Wie viele der fehlerhaften Wörter wurden richtig korrigiert Hierbei kann auf die Tabellen (Gold Labels) von Naive Bayes zurückgegriffen werden	<u>Methode:</u> Wortabgleich <u>Textbeispiel:</u> http://www.writerswrite.com/books/excerpts/inkheart-excerpt-400014
	Fehlerhafte Texte	Der richtig vorliegende Text muss intellektuell so verändert werden, dass Fehler enthalten sind. Hierbei	<u>Methode:</u> Fehlerhafte Wörter erkennen Frage: Liegt mein Wort in Lexikon vor → Ja: Wort ist richtig

		<p>muss beachtet werden, dass es verschiedene Arten von Fehlern gibt (die auch besonders Häufig sind)</p> <ul style="list-style-type: none"> - Buchstabe fehlt - Buchstabe zu viel - Buchstaben verdreht - Leerzeichen falsch (zu viel oder zu wenig) - Andere Fehler 	<p>➔ Nein: Berechnung des „richtigen“ Wortes</p> <p><u>Methode:</u> Fehlerhafte Wörter ersetzen</p> <p>➔ Gebe dem Nutzer das „richtige“ Wort</p>
Methoden ➔ Diese können im Programm umgesetzt werden	N-Grams (z.B. Henrichs's Pentagram Index)	<ul style="list-style-type: none"> - Es wird kontrolliert, wie viele n-Gramme (z.B. Trigramme) von Wörtern miteinander übereinstimmen 	<u>Methode:</u> checkNGrams
	Acquaintance: N-Grams in the Vector Space, Euklidische Distanz oder Cosinusähnlichkeit	<ul style="list-style-type: none"> - Ähnlichkeit von Vektoren (Ein Wort wird in Form von einem Vektor repräsentiert) - Distanz zwischen zwei Vektoren 	<u>Methode:</u> checkVectorSpace checkVectorDistance
	Edit distance	<ul style="list-style-type: none"> - Misst wie unähnlich sich zwei Strings sind ➔ Es wird die Anzahl an Operationen gemessen, die benötigt wird, um den falschen String in den richtigen zu verwandeln 	<p><u>Methode:</u> checkEditDistance</p> <p><u>Beispiel:</u> http://stackoverflow.com/questions/13928155/spell-checker-for-python</p> <p>http://norvig.com/spell-correct.html</p>
	Damerau-Levenshtein distance	<ul style="list-style-type: none"> - Misst die Anzahl der Veränderungen um einen String in einen anderen (den richtigen) zu verwandeln 	<p><u>Methode:</u> checkDLDistance</p> <p><u>Library:</u> https://pypi.python.org/pypi/pyxDamerauLevenshtein/1.4.1</p>
	Bayes Theorem	<ul style="list-style-type: none"> - Wir versuchen die Ersetzung mit der höchsten Wahrscheinlichkeit (gemeint zu sein) zu finden 	<p><u>Methode:</u> checkBayes</p> $\operatorname{argmax}_{c \in \text{candidates}} P(c w)$ $\operatorname{argmax}_{c \in \text{candidates}} P(c) P(w c) / P(w)$
Aufgaben	Textbearbeitung	<p>Text:</p> <ul style="list-style-type: none"> - Einlesen 	<u>Beispiel:</u>

		<ul style="list-style-type: none"> - Säubern - In Wörter zertrennen - In Teile Zertrennen, mit denen später gearbeitet werden kann (N-Grams, Buchstaben o.ä.) 	<ul style="list-style-type: none"> > import string > clean = [] > with open('book.txt') as f: // .txt einlesen > data = f.readlines() > for element in data: > s = element.translate(None, string.punctuation) // Punktation entfernen > clean.append(s)
	Linguistische Methoden anwenden	<p>Die implementierten Methoden sollen auf einen fehlerhaften Text angewendet werden.</p> <ul style="list-style-type: none"> - Fehlererkennung - Fehlerkorrektur - Darstellung der Ergebnisse (Gesamt oder einzelne Wörter mit ihrer Korrektur) 	Verschiedene Algorithmen, die einzeln für einen Fehler das „richtige“ Wort berechnen
	Ausgabe erstellen	<ul style="list-style-type: none"> - Wie viele Worte sind Falsch - Wie sehen sie (wahrscheinlich) richtig aus - Dies soll in einer kleinen GUI dargestellt werden 	<p>Darstellung jeweils gesplittet zwischen den einzelnen Methoden</p> <p>Vergleich zwischen den Algorithmen:</p> <ul style="list-style-type: none"> - War das Wort wirklich richtig - Darstellung aller Wortvorschläge in GUI Ausgabe
	Eigene Wörter dem Lexikon hinzufügen	<ul style="list-style-type: none"> - Button in der Haupt-GUI, die in zweites Fenster leitet → hier ist das Hinzufügen eigener Wörter möglich 	<p><u>Textfeld</u>: Hier kann ein neuer Lexikoneintrag hineingeschrieben werden</p> <p><u>Hinzufügen – Button</u>: Hiermit wird das Wort im Textfeld in das Textdokument hinzugefügt, indem sich das bisherige Lexikon befindet</p>
	Methoden Vergleichen/ Analysieren	<p>Welche der Methoden ist wie erfolgreich, wie viele Wörter werden richtig korrigiert. Dies kann in Form von Graphen dargestellt werden Analyse: Werden die Methoden heute noch genutzt, ist es sinnvoll sie anzuwenden?</p>	z.B.: allgemeiner Vergleich

	Visualisierung	z.B. mit Matplotlib direkt mit ausgeben <ul style="list-style-type: none"> - als falsch erkannte Wörter - Wörter die richtig als falsch erkannt wurden → Gold Labels - Vergleich der Methoden z.B. als Balkendiagramme 	<u>Library:</u> Matplotlib
Technische Details	Programmiersprache	Python	Anaconda
	Visualisierung	Matplotlib für Graphen etc.	<u>Library:</u> Matplotlib
	Präsentation der Ergebnisse/ Arbeitsoberfläche	Konsole und GUI	<u>Library:</u> Tkinter

