

Favourite Counts of Tweets of Donald Trump

Felix Herron

TU Berlin

15.5.2020



Table of Contents

Introduction

Erste Ergebnisse

- Preprocessing

- Erste Ergebnisse

Erster Ansatz

- Warum Classification Eigentlich?

- N-Gram Classifier

- Weniger Interessante Versuche

- Boosting

Der Fluch der Zeitverschiebung

- Die Hürde

- Das eventuelle Elixier

Der Fluch des Datenmangels sowie Qualität

- Die Hürde

- Das eventuelle Elixier

Nächste Schritte

Introduction

- ▶ Donald Trump ist US-Amerikanischer (bzw. mein) Präsident
#notMyPresident
- ▶ Seit Kontoeröffnung beinah 50.000 mal getweeted
- ▶ Einer der aller größten Twitterpersönlichkeiten ¹
- ▶ Ein geiles Thema für die Bachelorarbeit

¹<https://www.brandwatch.com/blog/most-twitter-followers/>

Introduction

- ▶ Donald Trump ist US-Amerikanischer (bzw. mein) Präsident
#notMyPresident
- ▶ Seit Kontoeröffnung beinah 50.000 mal getweeted
- ▶ Einer der aller größten Twitterpersönlichkeiten ¹
- ▶ Ein geiles Thema für die Bachelorarbeit

Ziele

1. Fav-count untersuchen/vorhersagen (heutiges Thema)
2. Mit topics verknüpfen (noch in Bearbeitung)
3. Feedback von euch sammeln :)

¹<https://www.brandwatch.com/blog/most-twitter-followers/>

Table of Contents

Introduction

Erste Ergebnisse

- Preprocessing

- Erste Ergebnisse

Erster Ansatz

- Warum Classification Eigentlich?

- N-Gram Classifier

- Weniger Interessante Versuche

- Boosting

Der Fluch der Zeitverschiebung

- Die Hürde

- Das eventuelle Elixier

Der Fluch des Datenmangels sowie Qualität

- Die Hürde

- Das eventuelle Elixier

Nächste Schritte

Tables and Figures

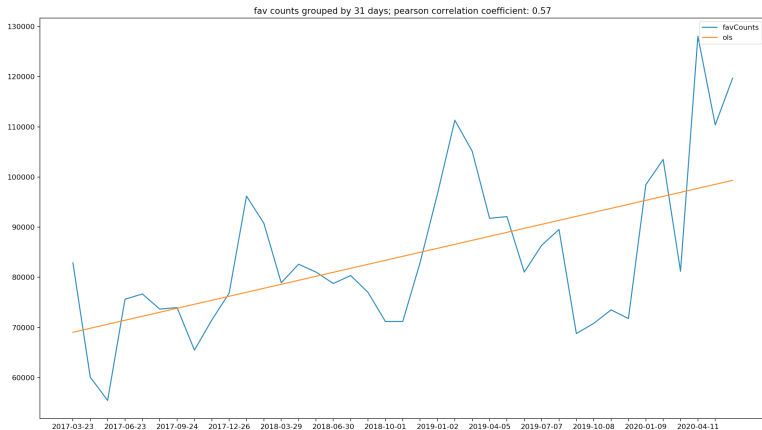
1. Tweets heruntergeladen, in eine DB geladen
2. Tweets gecleanet:
 - ▶ “RT”, links, “#”, u.a. entfernen
 - ▶ lemmas behalten

Tables and Figures

1. Tweets heruntergeladen, in eine DB geladen
2. Tweets gecleanet:
 - ▶ “RT”, links, “#”, u.a. entfernen
 - ▶ lemmas behalten
3. Mit verschiedenen Paramtern beschriftet:
 - ▶ isRT, isReply, isDeleted, mediaType, publishTime ...
 - ▶ n-Grams (1-4)

FavCount x Zeit

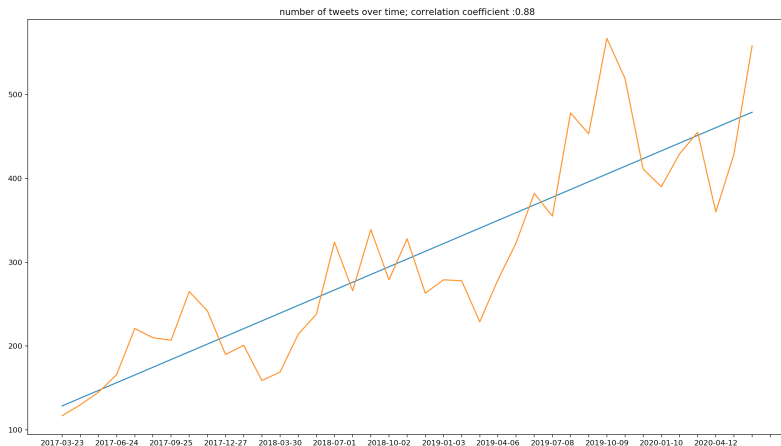
Die FavCounts von Trump steigen mit der Zeit



=> Bei der Vorhersage muss das in Betracht kommen... mehr dazu später

TweetCount x Zeit

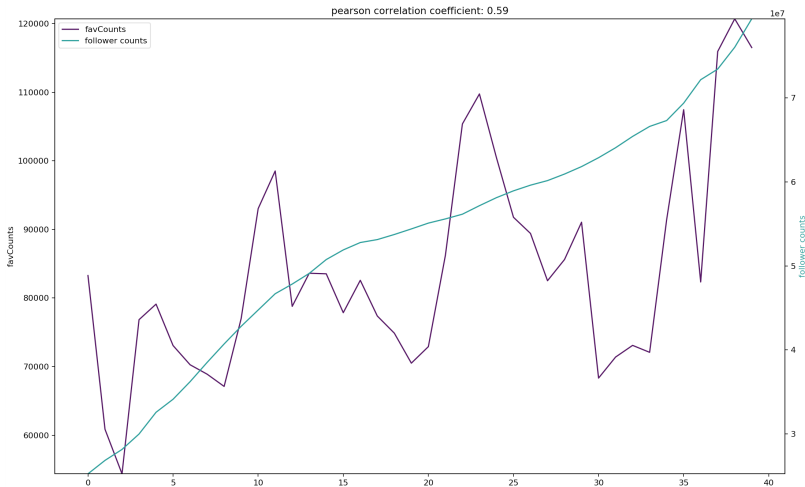
Die Anzahl an Tweets von Trump steigen mit der Zeit



=> warum?

FavCount x FollowerCount (bzw. Zeit)

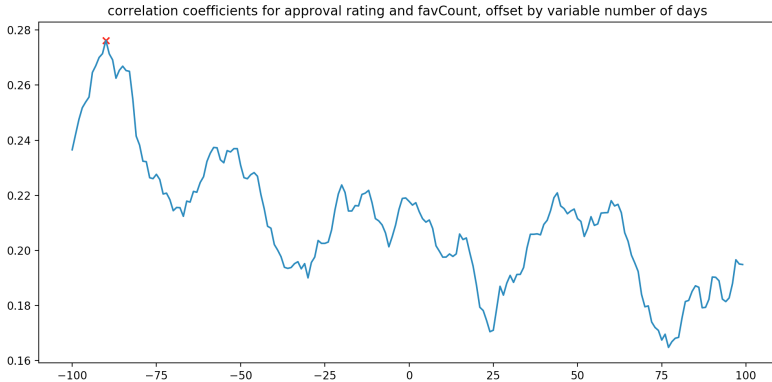
Die Anzahl seiner Follower steigt stetig mit der Zeit



=> Verknüpfung FollowerCount/FavCount? Unklar...

Approval Rating x Zeit

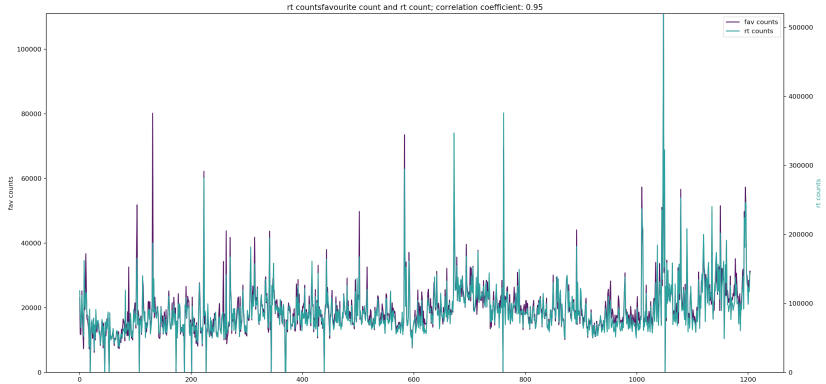
Die Approval Rating von Trump steigt schwach mit seiner FavCount



=> Wahrscheinlich Zufall

FavCount x RTCount

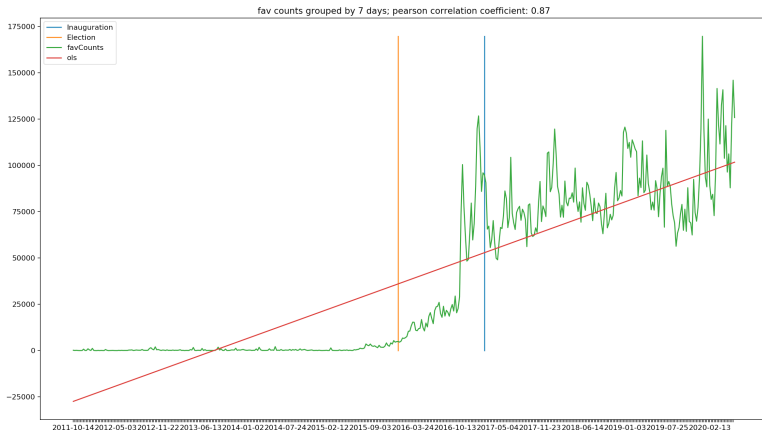
Fav Count \propto RT Count



=> FavCount und RT Count sind quasi ein Feature

FavCount x Zeit (Präsident)

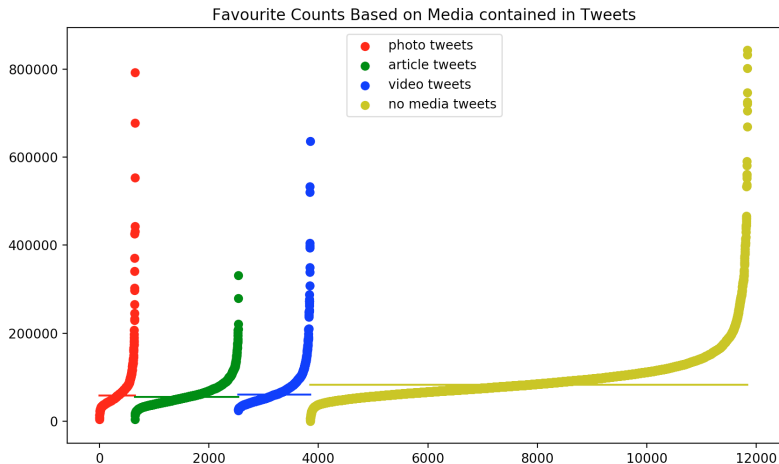
Die Reaktion auf seine Tweets hat sich seit seinem Amtsantritt wesentlich verändert



=> ich betrachten nur die Tweets seit Amtsantritt

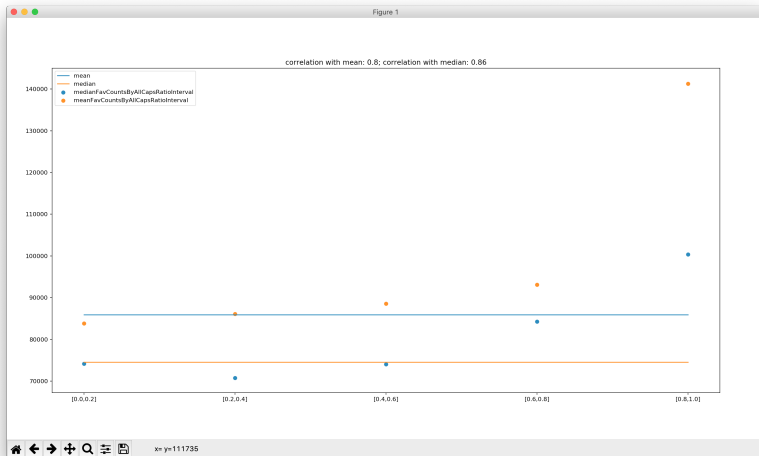
FavCount x MediaType

Die Art der beigefügten Medien variiert leicht mit der FavCount



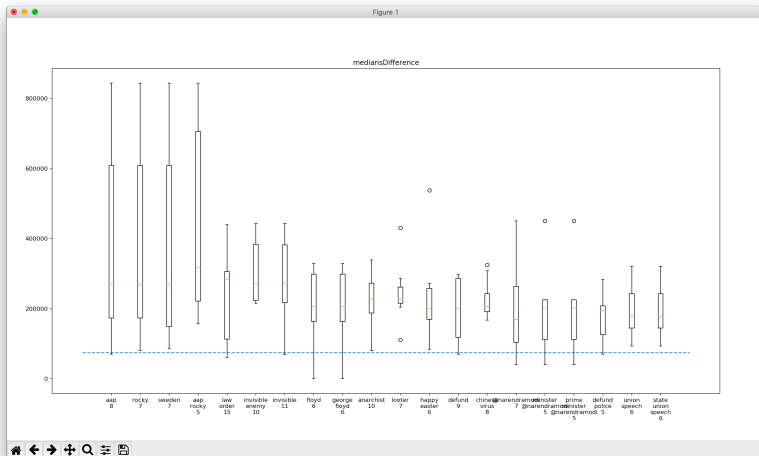
FavCount x allCaps Percentage

Die Anzahl an groß geschriebene Wörter hat einen starken Einfluß auf die Populärität



FavCount x nGrams Used

Die Wortwahl der Tweets hat ein Effekt auf deren FavCount



USW.

und so weiter und so weiter...²

²LADY HARRIET DURHAM

Table of Contents

Introduction

Erste Ergebnisse

Preprocessing

Erste Ergebnisse

Erster Ansatz

Warum Classification Eigentlich?

N-Gram Classifier

Weniger Interessante Versuche

Boosting

Der Fluch der Zeitverschiebung

Die Hürde

Das eventuelle Elixier

Der Fluch des Datenmangels sowie Qualität

Die Hürde

Das eventuelle Elixier

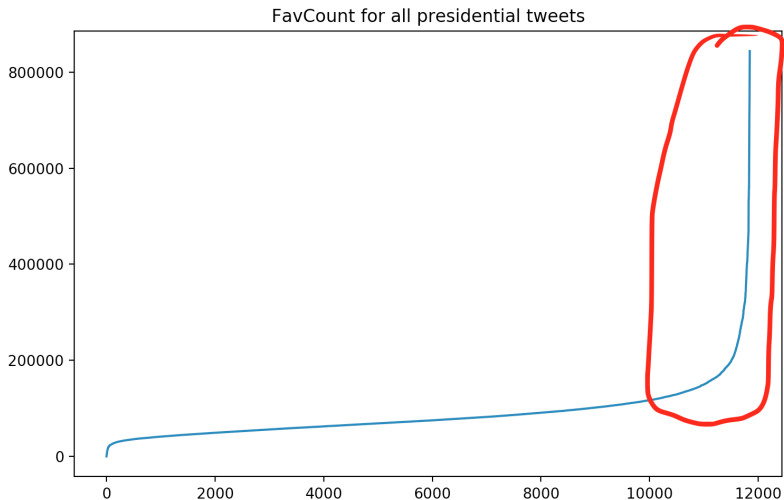
Nächste Schritte

Warum Classification Eigentlich?

- ▶ Eigentlich wäre Regression richtig
- ▶ Zu wenige Daten \Rightarrow Vorhersagen wären unpräzis
- ▶ Unpräzise Regression \approx Klassifikation

Viral Classification Motivation

- ▶ Weitere Trennungen kommen später
- ▶ Z.B. Viral/Dud
- ▶ Hier: Flache Mitte, Viral Tail



N-Gram Classifier

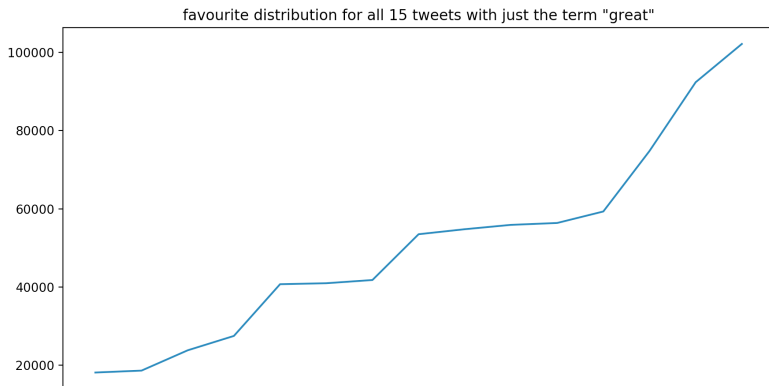
Erstmal aber Popular/Unpopular

- ▶ Idee: für jedes (relevante) N-Gram in einem Tweet, bestimme erwartete Skew
- ▶ Summiere (gewichtet) über alle Skews)

$$prediction(t) = c \sum_{nGram \in t} skew_{nGram} \cdot \sigma(favs_{nGram})$$

N-Gram Classifier Nachteile

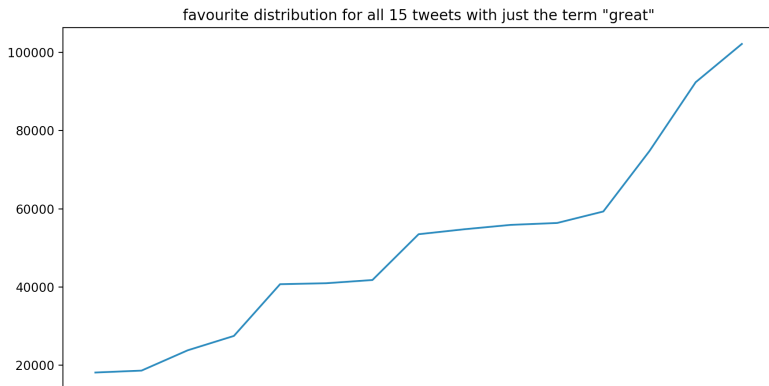
Viele Tweets haben sehr wenige Wörter



Wie unterscheidet man zwischen denen?

N-Gram Classifier Nachteile

Viele Tweets haben sehr wenige Wörter



Wie unterscheidet man zwischen denen?

Hyper- bzw. weitere parameter, z.B.

Weniger Interessante Versuche

Naive Bayes Classifier:

- ▶ Gleiche Daten, andere Perspektive als bei Boosting
- ▶ 65 % Genauigkeit, viel Overfitting

$$p(C_k \mid \mathbf{x}) = \frac{p(C_k) p(\mathbf{x} \mid C_k)}{p(\mathbf{x})}$$

Weniger Interessante Versuche

Naive Bayes Classifier:

- ▶ Gleiche Daten, andere Perspektive als bei Boosting
- ▶ 65 % Genauigkeit, viel Overfitting

$$p(C_k | \mathbf{x}) = \frac{p(C_k) p(\mathbf{x} | C_k)}{p(\mathbf{x})}$$

OLS:

- ▶ Jeder Tweet bekommt sparsen Eingabevektor, -1 oder 1 als Target
- ▶ 68 % Genauigkeit, viel Overfitting, Hyperparameter müssen noch untersucht werden

Immerhin Fehler auf nicht den selben Tweets

Boosting/Classifier Committee

- ▶ Es gibt viele Merkmale, nicht nur auf n-Gram-skew basiert, die Tweet Popularität vorhersagen
- ▶ Allerdings nicht so stark
- ▶ Sie machen Fehler bei unterschiedlichen Tweets

³<https://web.stanford.edu/hastie/Papers/buehlmann.pdf>

Boosting/Classifier Committee

- ▶ Es gibt viele Merkmale, nicht nur auf n-Gram-skew basiert, die Tweet Populärkeit vorhersagen
- ▶ Allerdings nicht so stark
- ▶ Sie machen Fehler bei unterschiedlichen Tweets

=> Classifier, der Ergebnisse aus vielen schwachen Classifiern zusammenträgt³

³<https://web.stanford.edu/hastie/Papers/buehlmann.pdf>

Boosting/Classifier Committee p. 2

Z.B.:

- ▶ n-Gram median discrepancy classifier sagt -1 mit geringer Sicherheit
- ▶ all-Caps Classifier sagt 1 mit hoher Sicherheit
- ▶ ...
- ▶ Prediction: 1

Ergebnisse aus diesem Ansatz

- ▶ 72% test Genauigkeit
- ▶ Noch nicht viele Classifier involviert
- ▶ Verbesserungsbedürftigkeiten bei umfangreicherem Zusammentun
- ▶ Gewichtungen lernen!

Ergebnisse aus diesem Ansatz

- ▶ 72% test Genauigkeit
- ▶ Noch nicht viele Classifier involviert
- ▶ Verbesserungsbedürftigkeiten bei umfangreicherem Zusammentun
- ▶ Gewichtungen lernen!

Wichtig: Abspaltung von Stärken/Schwächen

Table of Contents

Introduction

Erste Ergebnisse

- Preprocessing

- Erste Ergebnisse

Erster Ansatz

- Warum Classification Eigentlich?

- N-Gram Classifier

- Weniger Interessante Versuche

- Boosting

Der Fluch der Zeitverschiebung

- Die Hürde

- Das eventuelle Elixier

Der Fluch des Datenmangels sowie Qualität

- Die Hürde

- Das eventuelle Elixier

Nächste Schritte

Die Hürde

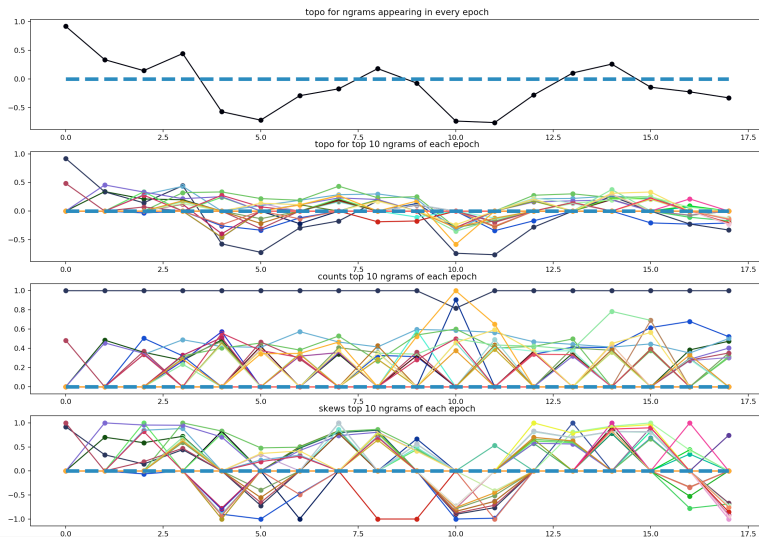
“It’s no use going back to yesterday, because I was a different person then”⁴

- ▶ $reaktion(nGram_1, time_1) \neq reaktion(nGram_1, time_2)$
- ▶ Wörter bekommen veränderte Signifikanz/Geschmack
- ▶ Z.B.: Great beginnt als positives Wort, am Ende ist es negativ
- ▶ Verbleibt häufig genutzt

⁴Lewis Carroll

Veranschaulichung des Fluches

Schützen Sie Sich die Augen!



Das eventuelle Elixier

- ▶ Zeit-lokalisiert trainieren
- ▶ Wird er fähiger?
- ▶ Bzw. veränderte Wortwahl \implies favCount \uparrow ?

Das eventuelle Elixier

- ▶ Zeit-lokalisiert trainieren
- ▶ Wird er fähiger?
- ▶ Bzw. veränderte Wortwahl \implies favCount \uparrow ?

\implies Das führt zu...

Table of Contents

Introduction

Erste Ergebnisse

- Preprocessing

- Erste Ergebnisse

Erster Ansatz

- Warum Classification Eigentlich?

- N-Gram Classifier

- Weniger Interessante Versuche

- Boosting

Der Fluch der Zeitverschiebung

- Die Hürde

- Das eventuelle Elixier

Der Fluch des Datenmangels sowie Qualität

- Die Hürde

- Das eventuelle Elixier

Nächste Schritte

Die Hürde

- ▶ Er hat seit seinem Amtsantritt nur etwa 12.000 mal getweetet
 - ▶ NLP Anwendungen aus der Forschung haben viel größere Corpora
- ▶ Tweets sind kurz
 - ▶ “bis auf die Knochen abgemagert” Dokumente
 - ▶ Viele sind Unschlüssig
 - ▶ Je raffinierter bei der Dokumentenwahl, desto weniger Daten!

Das eventuelle Elixier

Daten aus den Replies mitnutzen.

- ▶ Seeehr viele Leute antworten auf seine Tweets
- ▶ Die Inhalte dieser Antworten könnten viel über den Tweet enthüllen

Table of Contents

Introduction

Erste Ergebnisse

- Preprocessing

- Erste Ergebnisse

Erster Ansatz

- Warum Classification Eigentlich?

- N-Gram Classifier

- Weniger Interessante Versuche

- Boosting

Der Fluch der Zeitverschiebung

- Die Hürde

- Das eventuelle Elixier

Der Fluch des Datenmangels sowie Qualität

- Die Hürde

- Das eventuelle Elixier

Nächste Schritte

Analyse der Antworten

- ▶ Großer Datensatz
- ▶ Algorithmen anwenden, die zuvor ausgeschlossen waren
- ▶ Z.B. Word-embeddings weitertrainieren
- ▶ Z.B. Topic Modelling
- ▶ => Aktuelle Forschung brauchbarer

Topic/Sentiment Analysis

- ▶ Tweets nach Sentiment/Topic deren Antworten filtern
- ▶ Daraus FavCount vorhersagen?
- ▶ Hähnchen und Ei...

das wars

Vielen Dank, dass Sie zugehört haben

Und Danke an Stephi für die bisherige Betreuung!

das wars

Vielen Dank, dass Sie zugehört haben

Und Danke an Stephi für die bisherige Betreuung!

Ich freue mich über jeden Vorschlag :)