

# Checkliste zu NLP-Klassifikations-Experimenten

Stand: 14. November 2024

## Hinweise

- Wenn Ihr Task kein Klassifikationstask ist, dann ist dieser Fragebogen nicht für Sie.
  - Keine Klassifikationstasks sind z.B. Generierung oder Übersetzung.
- Markieren Sie alle Punkte die Sie planen umzusetzen.
- Machen Sie einen Termin in der Sprechstunde, wenn Ihnen Punkte unklar sind.
- Der Fragebogen ist keine Prüfung, sondern dient als Hilfestellung, bei der Experimentplanung an alles zu denken, auf Ideen zu kommen und ggf. die richtigen Fragen zu stellen. Diese klären wir dann idealerweise im Gespräch.
- Der Fragebogen repräsentiert eine Planung. Gewisse Abweichungen von der Planung sind normal und zu erwarten.
- Bei Fragen schreiben Sie gerne eine E-Mail an [nils.reiter@uni-koeln.de](mailto:nils.reiter@uni-koeln.de) oder melden Sie sich gerne zu einer Sprechstunde an. Wenn Sie rein technische Fragen zur Infrastruktur haben, schreiben Sie bitte an [spinfo-admin@uni-koeln.de](mailto:spinfo-admin@uni-koeln.de).

## Der Task

1. Die Aufgabe heißt: Erkennung des Texttyps "Roman"
2. Es handelt sich um  Textklassifikation,  Sequence Labeling, oder  Sonstiges: \_\_\_\_\_
3. Die zu klassifizierenden Instanzen sind: Texte (400-500 Wörter Länge)
4. Es gibt 2 Kategorien/Klassen.
5. Einer Instanz kann  genau eine oder  mehrere Klassen zugewiesen werden.

## Die Daten

1. Annotierte Daten  liegen bereits vor oder  müssen noch erstellt werden.
2. In den Daten sind 10000 Instanzen (von o.g. Typ) annotiert.
3. Die Klassen sind
  - gleichverteilt (d.h. jede Klasse ist ungefähr gleich häufig) (Trainingsdaten)
  - unterschiedlich verteilt, und zwar: 3:1 (Testdaten)

## Die Annotationen

*Nur relevant, wenn neue Daten annotiert werden sollen.*

1. Annotationsrichtlinien
  - Ich verwende die folgenden, bereits existierenden Annotationsrichtlinien: \_\_\_\_\_
    - Mit diesen wurde ein Inter-Annotator-Agreement von \_\_\_\_\_ erzielt (Metrik: \_\_\_\_\_).
  - Ich schreibe neue Annotationsrichtlinien.

## 2. Annotator:innen

- Ich annotiere selbst.
- Ich rekrutiere Annotator:innen aus meinem Freundes-/Bekannteskreis.
- Ich sammle Annotationen über eine Umfrage, z.B. mittels LimeSurvey.
  - Ich verfüge über einen Verteiler, über den ich die Umfrage bewerbe.
- Ich sammle Annotationen über crowd sourcing.

## 3. Annotationsworkflow

- Annotator:innen treffen eine Annotationsentscheidung auf der Basis eines Kontextes von \_\_\_  Wörtern,  Sätzen,  Zeilen,  Absätzen,  \_\_\_\_\_; oder  sie verwenden den gesamten Text als Kontext.
- Sie können dabei außerdem die folgenden Wissensquellen verwenden:  Wikipedia,  Lexika,  Wörterbücher,  Sonstige: \_\_\_\_\_

## 4. Anforderungen an Annotationssoftware

- Annotator:innen müssen Spannen selbst markieren können.
- Annotator:innen müssen neue Kategorien oder Labels ergänzen können.

## Die Baseline

- Weil die Klassen ungleich verteilt sind, bietet sich eine majority baseline an. Diese erzielt eine Accuracy von 75 %.
- Weil die Klassen gleich verteilt sind, bietet sich eine random baseline an. Diese erzielt eine Accuracy von \_\_\_\_\_ %.
- Eine weitere mögliche Baseline ist: Wir schauen nach Anführungszeichen. Diese erzielt eine Accuracy von ... %.
- Eine weitere mögliche Baseline ist: \_\_\_\_\_ . Diese erzielt eine Accuracy von \_\_\_\_\_ %.

## Das Experiment

Ich möchte das folgende oder die folgenden Verfahren verwenden:

- Klassisches Machine Learning
  - Entscheidungsbaum / Decision Tree (DT)
  - Naive Bayes
  - Support Vector Machines (SVM)
  - Logistic Regression
  - Neuronale Netze (NN)
    - Feed-Forward Neural Networks
    - Convolutional Neural Networks
    - Recurrent Neural Networks
- Große Sprachmodelle / Transformer (inkl. BERT & co.)
  - Welche(s) Modell(e) genau? \_\_\_\_\_
  - Fine-Tuning
  - Prompting
- Sonstige: \_\_\_\_\_

Klassisches Machine Learning

Nur relevant, wenn klassisches Machine Learning verwendet werden soll.

1. Ich möchte die folgenden Features verwenden

Metadaten:   **Autor**    
\_\_\_\_\_  
\_\_\_\_\_

Inhaltsdaten, z.B. aus Texten:

Worthäufigkeiten (von allen Wörtern), auch bekannt als bag of words

Häufigkeiten von Wörtern aus folgenden Wortlisten:   **Redewiedergabe**  

Embeddings (z.B. Word Embeddings)

Sequenzielle Information (d.h. Klassifikationsergebnisse für Elemente davor oder danach)

N-Gram-Häufigkeiten, mit  $N \leq$  \_\_\_\_\_

Thematische Informationen aus einem Topic Model (z.B. Latent Dirichlet Allocation, LDA)

**Anzahl der Zeilenumbrüche**

2. Meine Features haben die folgenden Datentypen:

**Klasse und Autor**

Numerisch:   **alle - 2**   (Anzahl Features)

Kategorial:   **2**   (Anzahl Features)

Große Sprachmodelle / Transformer

Nur relevant, wenn große Sprachmodelle verwendet werden sollen.

Ich möchte Fine-Tuning verwenden

Ich möchte Prompting verwenden, und verwende dabei die folgenden Strategien:

Hard prompting (d.h., der Prompt besteht aus Wörtern) mit den Komponenten:

Beschreibung der Aufgabe

Versprechen einer Belohnung

Aufforderung, Schritt für Schritt vorzugehen

Definition einer Rolle, die das Modell einnehmen soll

Beispiele –  Positive oder  Negative

Sonstiges: \_\_\_\_\_

Soft prompting (d.h., der Prompt besteht aus Vektoren) mit dieser/diesen Technik/en: \_\_\_\_\_  
\_\_\_\_\_  
\_\_\_\_\_

Das Modell soll folgenden Output produzieren

Numerischen

Sprachlichen. Dieser soll wie folgt ausgewertet werden: \_\_\_\_\_  
\_\_\_\_\_  
\_\_\_\_\_

JSON oder YAML.

Sonstigen: \_\_\_\_\_

Mit nicht-wohlgeformtem Output möchte ich wie folgt umgehen: \_\_\_\_\_  
\_\_\_\_\_  
\_\_\_\_\_

Beide Varianten

1. Testdaten

- Ich teile meinen o.g. Datensatz selbst in Trainings- und Testdaten auf. \_\_\_\_\_% der Instanzen werden als Trainingsdaten verwendet.
- Ich verwende  $N$ -fold cross validation, mit  $N =$  \_\_\_\_\_.
- Trainings- und Testdaten sind bereits aufgeteilt, z.B. weil es Daten aus einem shared task sind.

2. Ich variiere und vergleiche

- die Größe des Trainingsdatensatzes (z.B. 100, 1000, 10 000 Instanzen für den Trainingsdatensatz)
- die Menge an oder Art von Features die verwendet werden (z.B. inhaltliche vs. sprachliche Features)
- das Verfahren als solches oder Parameter davon (z.B. NN vs. SVM)
- die Vorverarbeitung (z.B. Groß- und Kleinschreibung)
- die Prompt-Komponenten (z.B. mit unterschiedlichen Belohnungen)

3. Meine Hypothese ist: Mehr Trainingsdaten → Bessere Performance

---

---

---

## Die Auswertung und Evaluation

1. Ich verwende die Evaluationsmetrik(en)

- Accuracy    Precision    Recall    F-Measure    Area under curve (AUC)
- Sonstige: \_\_\_\_\_

- Meine Testdaten sind stark unbalanciert, daher verwende ich die Metriken in der Mikro- und Makro-Average-Variante.

2. Für meine Fehleranalyse inspiziere ich 20 Instanzen manuell.

## Die praktische Umsetzung

1. Ich verwende die Programmiersprache

- Python
- Java
- R
- \_\_\_\_\_

2. Hardware-Ausstattung und Vorkenntnisse

- Ich verfüge über einen Computer
  - der auch mal über Nacht durchlaufen kann, wenn eine Berechnung etwas länger dauert.
  - der eine GPU mit CUDA-Unterstützung hat oder ein Mac mit M1/M2-Prozessor ist.
  - der ausreichend freien Plattenspeicher hat.
- Ich möchte Berechnungen auf einem Server der Universität laufen lassen.
  - Ich kann mich per SSH auf einem Server einloggen.
  - Ich weiß wie ich auf einer Kommandozeile ein Programm laufen lasse.