# Text Corpus

Salar Mohtaj | DFKI

# Text corpus

- What is a text corpus?

- Sample text corpora

- Corpus annotation

- Underfitting and overfitting

- Data splitting

- Text corpora in Python

# Text corpus

- What is a text corpus?

- Sample text corpora

- Corpus annotation

- Underfitting and overfitting

- Data splitting

- Text corpora in Python

# What is a text corpus?

- Text corpus is a collection of text, usually contains several documents
  - Wikipedia articles
  - Collection of movies reviews
  - Internet comments
  - Collection of tweets

- A corpus may be quite small, for example, containing only **thousands words** of text, or very large, containing **millions of words**

- NLP corpora are in different standards based on the target task

# What is a text corpus?

- Text corpora can be compared based on different factors
  - **Size**: larger corpora are better for training deep learning models
  - **Domain**: unless corpus has been collected for specific tasks, it should include different genres such as newspapers, magazines, blogs, academic journals, etc
  - **Metadata**: metadata should indicate the sources, assumptions, limitations and what's included in the corpus
  - **Clean**: a wordlist giving word forms of the same word can be messy to process

# Text corpus

- What is a text corpus?

- **Sample text corpora**

- Corpus annotation

- Underfitting and overfitting

- Data splitting

- Text corpora in Python

# Sample text corpora

- Common file formats
- TXT
- CSV
- JSON
- XML

# Sample text corpora

- Where to find text corpora?
- NLP competitions
  - NLP shared task
- Kaggle
- Active NLP groups websites

# Text corpus

- What is a text corpus?

- Sample text corpora

- **Corpus annotation**

- Underfitting and overfitting

- Data splitting

- Text corpora in Python

# Corpus annotation

- It would happen that you have to compile your own corpus for the task and domain of interests
  - e.g., generating a text summarization corpus for medical texts in German

- Apart from the pure text, a corpus can also be provided with additional linguistic information, called **annotation**

- Corpus annotation is the practice of adding interpretative linguistic information to a piece of text

# Corpus annotation

# Corpus annotation

# Corpus annotation

- The annotation process could meet lots of questions and ambiguities

- Annotation guideline
  - Describe the annotation procedure as generic as possible but as precise as necessary
  - So that human annotators can annotate the concept or phenomenon in any text without running into problems or ambiguity issues

PROPOSED SEMANTIC-ROLE BASED SENTIMENT QUESTIONNAIRE

Q1. From reading the text, the speaker's emotional state can best be described as:

- *positive state*: there is an explicit or implicit clue in the text suggesting that the speaker is in a positive state, i.e., happy, admiring, relaxed, forgiving, etc.

- *negative state*: there is an explicit or implicit clue in the text suggesting that the speaker is in a negative state, i.e., sad, angry, anxious, violent, etc.

- *both positive and negative, or mixed, feelings*: there is an explicit or implicit clue in the text suggesting that the speaker is experiencing both positive and negative feelings

- *unknown state*: there is no explicit or implicit indicator of the speaker's emotional state

Q2. From reading the text, identify the entity towards which opinion is being expressed or the entity towards which the speaker's attitude can be determined.

# Corpus annotation

- Laboratory vs. crowdsourcing
  - Laboratory
    - asking experts (e.g., linguists) for data annotation
  - Crowdsourcing
    - Is a participatory method of building a dataset with the help of a large group of people

# Corpus annotation

- Crowdsourcing

# Corpus annotation

## Crowdsourcing

- Cheaper
- Faster
- Scalable
- Diverse group of participants
- No moderator
- Unreliable data

## Laboratory

- High reliability
- Controlled unwanted factors
- Environment constant, no noise
- Participants: one-to-one contact to moderator

# Text corpus

- What is a text corpus?

- Sample text corpora

- Corpus annotation

- **Underfitting and overfitting**

- Data splitting

- Text corpora in Python

# Underfitting and overfitting

- Overfitting
  - Overfitting refers to a model that models the training data too well
  - It happens when a model learns the detail and noise in the training data to the extent that it negatively impacts the performance of the model on new data
  - This means that the noise or random fluctuations in the training data is picked up and learned as concepts by the model
  - The problem is that these concepts do not apply to new data and negatively impact the models ability to generalize

# Underfitting and overfitting

- Overfitting



Appropriate-fitting

Overfitting

# Underfitting and overfitting

- Underfitting
  - Underfitting refers to a model that can neither model the training data nor generalize to new data
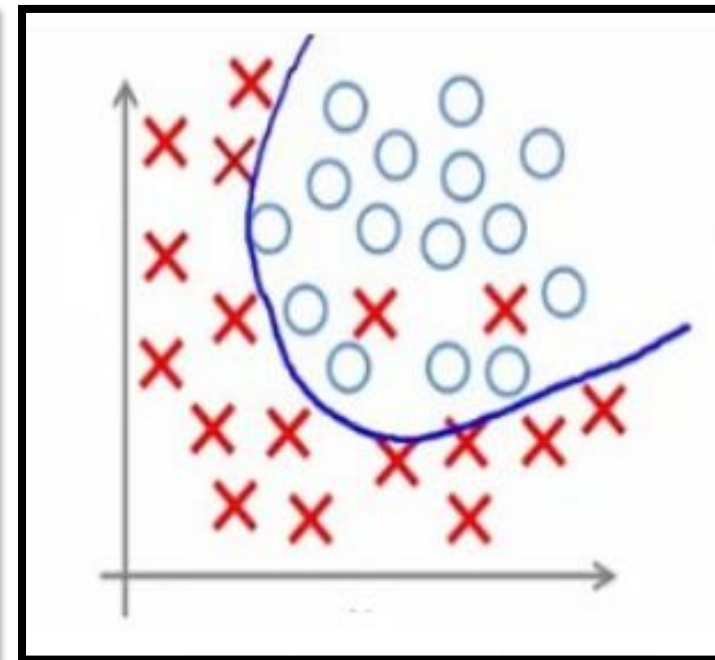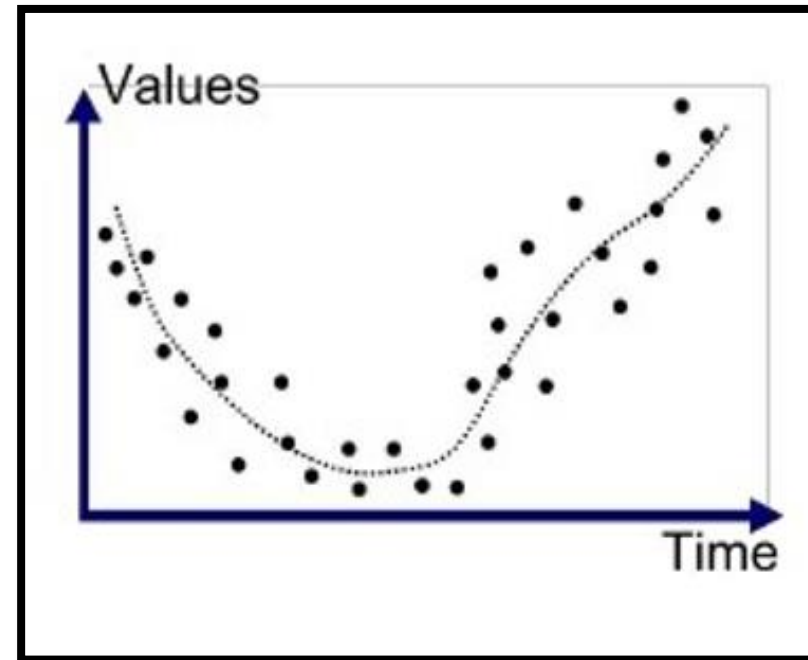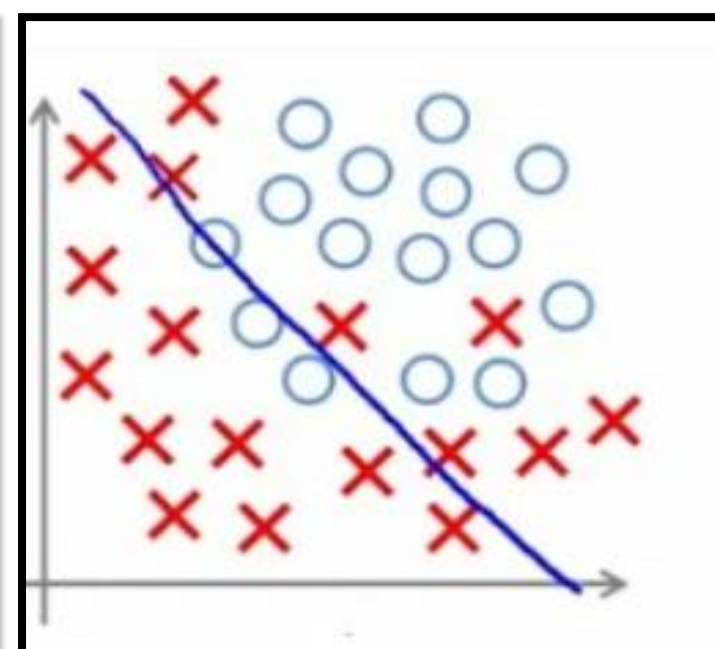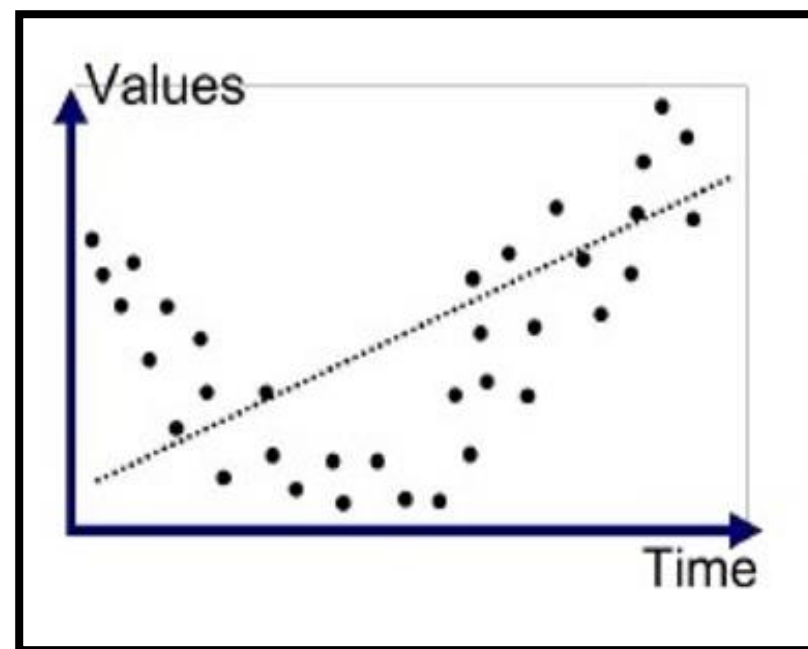  - An underfit machine learning model is not a suitable model and will be obvious as it will have poor performance on the training data

# Underfitting and overfitting

- Underfitting



Appropriate-fitting

Underfitting

# Underfitting and overfitting

- Overfitting
  - Increase training data
  - Reduce model complexity
  - Cross-validation
  - Early stopping during the training phase
  - Regularization

- Underfitting
  - Increase model complexity
  - Increase number of features
  - Remove noise from the data

# Text corpus

- What is a text corpus?

- Sample text corpora

- Corpus annotation

- Underfitting and overfitting
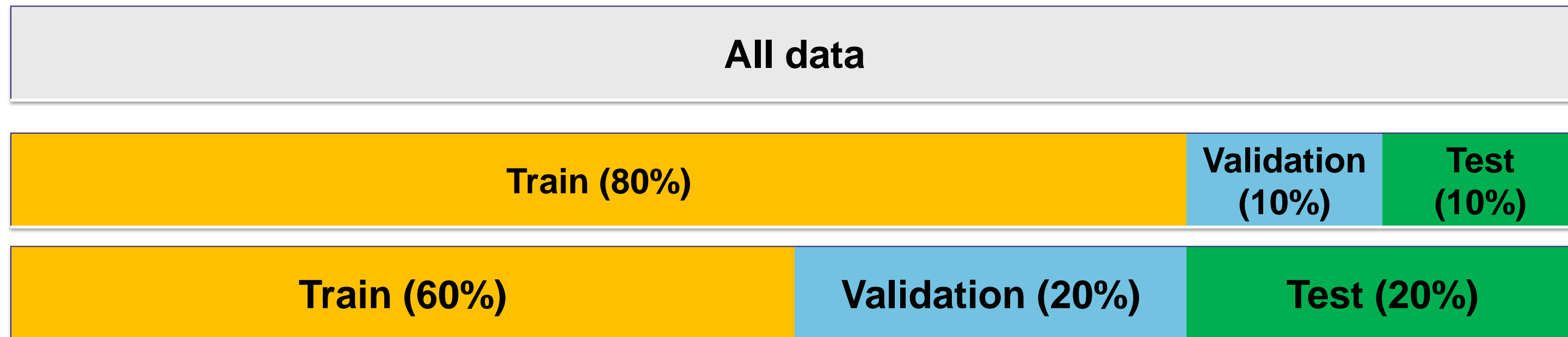
- **Data splitting**

- Text corpora in Python

# Data splitting
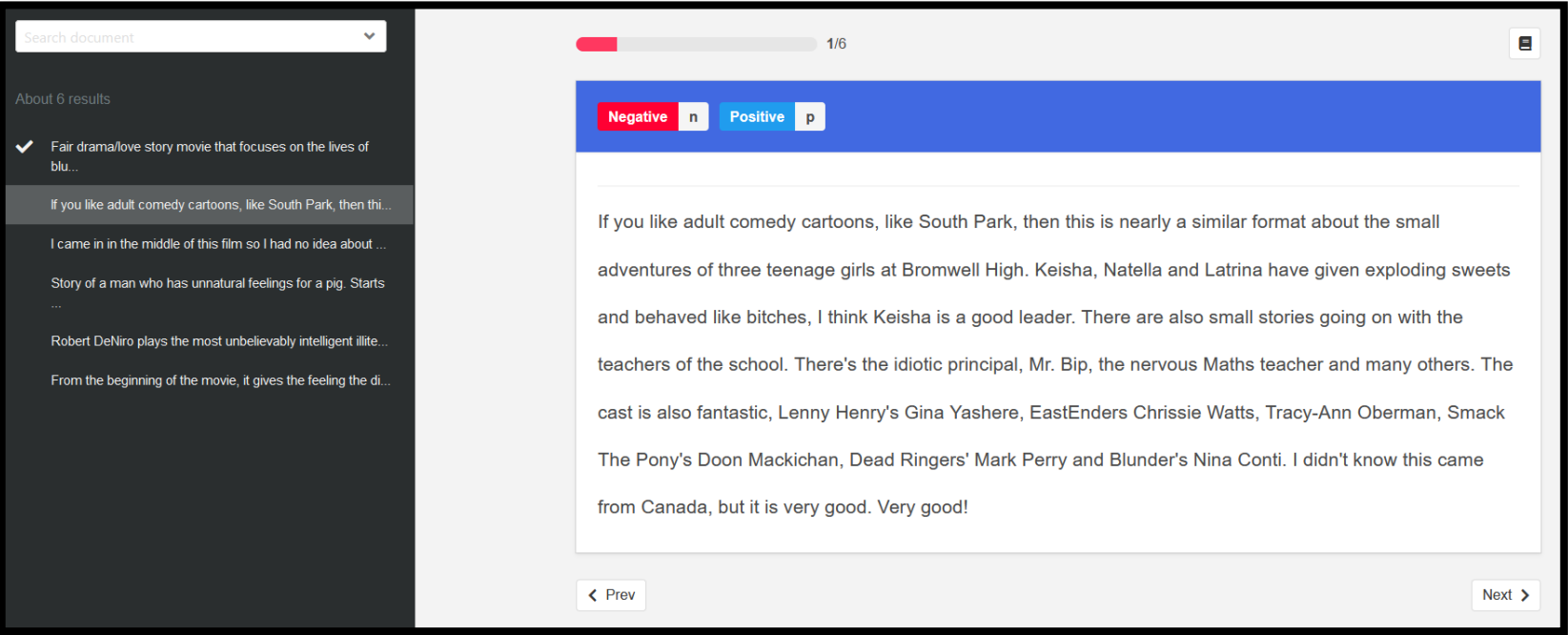
- As a common practice in machine learning (NLP), we have to split the dataset (corpus) into different part in order to train a model and test it

- Two main approaches for splitting data
  - Train / Validation / Test
  - Cross validation

# Data splitting

- Train / Validation / Test
  - **Training Dataset**: The sample of data used to fit the model
  - **Validation Dataset**: The sample of data used to provide an unbiased evaluation of a model fit on the training dataset while tuning model hyperparameters
  - **Test Dataset**: The sample of data used to provide an unbiased evaluation of a final model fit on the training dataset

| All data |
|:---:|

| Train (80%) | Validation (10%) | Test (10%) |
|:---:|:---:|:---:|

| Train (60%) | Validation (20%) | Test (20%) |
|:---:|:---:|:---:|

# Data splitting

- Cross validation
  - Cross-validation is a resampling procedure used to evaluate machine learning models on a limited data sample

- The procedure has a single parameter called k that refers to the number of groups that a given data sample is to be split into
  - the procedure is often called k-fold cross-validation

# Data splitting

# Text corpus

- What is a text corpus?

- Sample text corpora

- Corpus annotation

- Underfitting and overfitting

- Data splitting

- **Text corpora in Python**

# Text corpora in Python

- Reuters Corpus
- Brown Corpus
- Web and Chat Text
- Gutenberg Corpus
- …

- And Corpora in Other Languages

NLTK
Natural Language Toolkit

```
>>> from nltk.corpus import reuters
>>> reuters.categories()
['acq', 'alum', 'barley', 'bop',
'carcass', 'castor-oil', 'cocoa',
'coconut', 'coconut-oil', 'coffee', ...]
>>> from nltk.corpus import brown
>>> brown.words()
['The', 'Fulton', 'County', 'Grand',
'Jury', 'said', ...]
```

# Summary

- A corpus is a collection of texts, **written** or spoken

# Summary

# Summary