

Input data
 - subject
 - body

→ ML classifier

- Non-architectural
- Existence
- Technology
- Process
- Property

Input data

↓ raw data

Text preprocessing

↓ preprocessed data = clean string

- concatenate subject + body
- lowercase → NLTK, SKLEARN
- remove replies ← Handmatig
- remove HTML?

BOW

Sparse

1 Ik gooi een bal
 2 Ik gooi een steen] dataset → [11110] → [0123]
 0 Ik
 1 gooi
 2 een
 3 bal
 4 steen

Sparse vector

Feature generation

↓ vectors

ML classifier

- Bag of Words - CountVectorizer sklearn
- TFIDF - TfidfVectorizer sklearn
- Decision tree classifier
- Random Forest classifier
- SVM - linear SVC
- Naive Bayes

Gridsearch
 CV
 sklearn

- Training ↗ k-fold cross validation
- Testing ↖
- Optimiseren ML classifier

↓

Resultaten

1 Ik gooi een een bal
 2 Ik gooi een steen

TF = 2
 DF = 2

TF = 1
 DF = 2

1
 $\frac{1}{2}$