many document have associated metadate in the form of classes or labels (categorical variables)

approximate a function that can map documents features onto class information

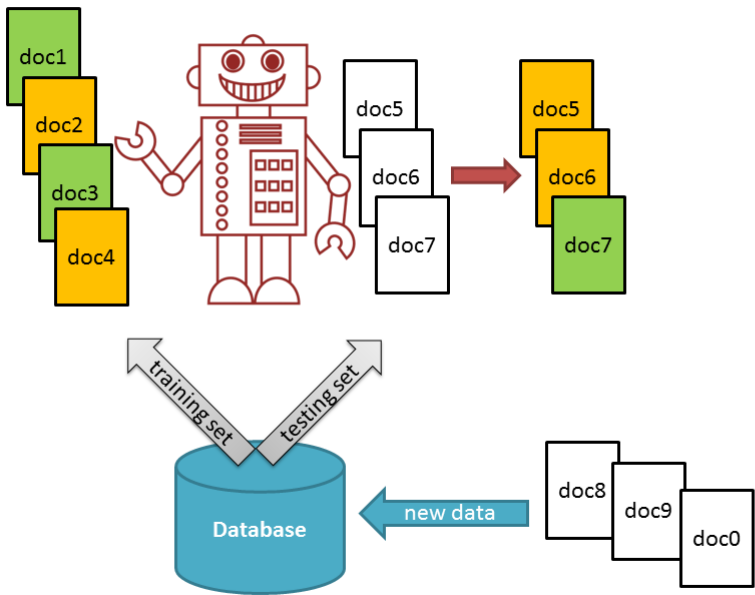preferably, our model should be the best performing classifier in the set of possible classifiers

given labeled data (supervised learning), a classification algorithm will output a solution that categorizes new examples $\rightarrow$ associate labels with subsets of the data

while clustering (unsupervised learning) searches for groups within the corpus, classification learns to map a collection of documents onto a categorical class values or labels $\rightarrow$ find mapping function
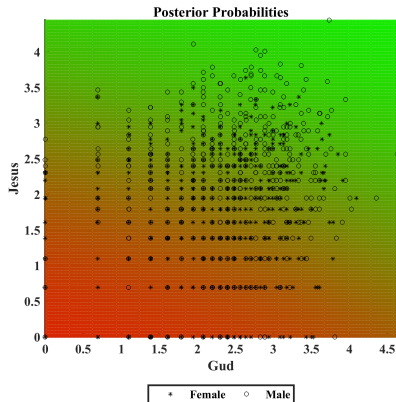
data (features) with class values ($\sim$ labeled data), excellent opportunity to make use of metadata

vast majority of models are **black box models**

workflow: separate data set in training and test subsets (training, test, and validation) $\rightarrow$ train model $\rightarrow$ test model $\rightarrow$ apply model to new data

doc1 doc2 doc3 doc4

doc5 doc6 doc7

doc5 doc6 doc7

training set

testing set

Database

new data

doc8 doc9 doc0

# classification in the humanities

binary and multiclass classification problems [1]

**naive bayes** probabilistic classifier that is fast and popular for in text categorization, but assumes independence between features (naive)

**neural network** broad framework for machine learning, which is very extremely flexible. Training can be very slow, but classification fast. Prone to overfitting

**decision tree** versatile and creates sets of rules (binary decisions) that are simple and can be understood (leaves are classes and branches features) $\rightarrow$ white box method

**support vector machines** works on small datasets (typically binary) with high dimensional data (features $>$ objects) and very memory efficient (only uses the support vectors). Bad performance on noisy data (overlapping classes)

---

[1]Can be advantageous to reformulate multiclass problems as binary

**labeled data** the correct class information is available

- metadata is readily available, e.g. author, genre, year of publication
- labels from an external source/databases, e.g. reviews, ratings, reads
- annotate data (expert or raters)

**evaluate** performance (error rate) of a classifier and compare to other classifiers

most metrics are developed for binary classification problems

confusion matrix: table for describing performance of classifier on training and/or testing data

|  |  | true | |
|---|---|---|---|
|  |  | positive | negative |
| predicted | positive | TP | FP |
|  | negative | FN | TN |

**True Positive** correctly assigns positive class membership
**True Negative** correctly rejects class membership
**False Positive** fail to rejects class membership
**False Negative** reject class membership incorrectly

we train a naive bayes classifier on 1500 verses of the kjv bible labeled with collection data (nt: new testament ot: old testament)

Confusion matrix for binary classification problem:

|     | nt  | ot  |
| --- | --- | --- |
| nt  | 644 | 89  |
| ot  | 106 | 661 |

, verses: $644 + 106 + 89 + 661 = 1500$

**accuracy** measures in how many cases the predicted class conformed with the correct class: $\dfrac{TP + NP}{TP + TN + FP + FN}$

|    | nt  | ot  |
|----|-----|-----|
| nt | 644 | 89  |
| ot | 106 | 661 |

, accuracy: $\dfrac{(644 + 661)}{1500} = 0.87\ (87\%)$

**precision** measures the number of selected verses that are relevant, i.e., how certain are we that a classified verse is correctly classified ($\sim$ how many time did the model positively predict a class): $\dfrac{TP}{TP + FP}$

| | nt | ot |
|---|---|---|
| nt | 644 | 89 |
| ot | 106 | 661 |

, $precision_{NT}$: $\dfrac{(644)}{644 + 89} = 0.88$

for each class label: How many of the items that got the label should have gotten it? How many should have gotten other labels?

**recall** measures the number of relevant verses that are selected, i.e., how good is the classifier at detecting verses within a given class: $\dfrac{TP}{TP + FN}$

|     | NT  | OT  |
| --- | --- | --- |
| NT  | 644 | 89  |
| OT  | 106 | 661 |

, $recall_{NT}$: $\dfrac{(644)}{644 + 106} = 0.86$

For each class label: How many items that should have gotten the label did get it? How many were missed?

**F-score** composite (general) measure of a classifier's accuracy

$$F_1 = 2 \times \frac{percision \times recall}{precision + recall}$$

$$F_1 : 2 \times \frac{.88 \times .86}{.88 + .86} = 0.87$$

$F$ is the harmonic mean of precision and recall.

if a model is sufficiently complex and gets enough data, it can basically memorize the data set (overfitting) $\rightarrow$ need to test the model on held-out data

**validation** when building a predictive model, we need a way to evaluate the capability of the model on unseen data

- data Split (conventional validation)
- cross validation
- bootstrap

## classification with scikit-learn

```
1   datapath = '/home/kln/corpora/kjv_books'
2   docs = vanilla_folder(datapath)
3
4   import pandas as pd
5   import numpy as np
6   metadata = pd.read_csv('/home/kln/corpora/kjv_metadata.csv')
7   class_id = metadata['class'].tolist()
8   class_u, class_int = np.unique(class_id, return_inverse = True)
9
10  from sklearn.feature_extraction.text import CountVectorizer
11  countvect = CountVectorizer()
12  vectspc = countvect.fit_transform(docs)
13  vectspc.shape
14
15  # index value of a word in the vocabulary
16  countvect.vocabulary_.get(u'god')
17  countvect.vocabulary_.get(u'woman')
18
19  # build vector space model
20  from sklearn.feature_extraction.text import TfidfTransformer
21  tfidf_transformer = TfidfTransformer()
22  vectspc_tfidf = tfidf_transformer.fit_transform(vectspc)
23  vectspc_tfidf.shape
24
25  # train naive bayes classfier
26  from sklearn.naive_bayes import MultinomialNB
27  nb_class = MultinomialNB().fit(vectspc_tfidf, class_id)
28  # classifier training performance
29  predicted = nb_class.predict(vectspc_tfidf)
30  np.mean(predicted == class_id)
31
32  # svm for comparison
33  from sklearn.linear_model import SGDClassifier
34  svm_class = SGDClassifier(loss='hinge', penalty='l2',alpha=1e-3, n_iter=5,
35          random_state=42).fit(vectspc_tfidf, class_id)
36  predicted = svm_class.predict(vectspc_tfidf)
37  np.mean(predicted == class_id)
```

## classification with RTextTools

```
1   library(RTextTools)
2   ## separate training and testing set and create a container
3   # random sample for testing data from data set
4   trainidx.v <- 1:nrow(text.dtm)
5   testidx.v <- sort(sample(trainidx.v, nrow(text.dtm)*.1, replace = FALSE, prob = NULL))
6   trainidx.v <- sort(trainidx.v[! trainidx.v%in%testidx.v])
7
8   # change object type, create_analytics() only handles numeric
9   classnum.v <- as.numeric(as.factor(class.v))
10    # to transform back to original
11    factor(classnum.v, labels = unique(class.v))
12  # create container
13  container <- create_container(text.dtm, classnum.v, trainSize=trainidx.v,
14                                testSize=testidx.v, virgin=FALSE)
15  # training models
16  mdl1.l <- train_models(container, algorithms='SVM')
17  mdl2.l <- train_models(container, algorithms = c('SVM','NNET','TREE') )
18
19  # Classifying data
20  res.df <- classify_models(container, mdl2.l)
21  head(res.df)
22  confusion.mat <- as.matrix(table(res.df$SVM_LABEL, container@testing_codes))
23  rownames(confusion.mat) <- colnames(confusion.mat) <- unique(class.v)
24  print(confusion.mat)
25  accuracy <- sum(diag(confusion.mat))/sum(confusion.mat)
26
27  # performance metrics
28  analytics <- create_analytics(container, res.df)
29  class(analytics)
30  summary(analytics)
```