

# INF519 Machine Learning 2: Homework 2

## Filtering spam messages

Kyota Lannelongue – Hugo Cambier – Thomas Charuel

The goal of this homework is to explore the Naïve Bayes method compared to the SVM method from the previous homework.

Whereas SVM is a discriminative model, Naïve Bayes is a Generative model. In discriminative models, we model the dependence of  $y$  on  $x$ , and by finding the conditional probability distribution  $P(y|x)$ , we ultimately train our dataset to predict  $y$  from  $x$ . It simply categorizes a given signal. Generative models analyze both the dependence of  $y$  on  $x$  and  $x$  on  $y$  to generate both inputs and outputs. It models how the data was generated to categorize a signal. One of the advantages of generative algorithms is that you can use  $P(x, y)$  to generate new data similar to existing data. On the other hand, discriminative algorithms generally give better performance in classification tasks.

For this exercise, we will take emails text and categorize it as spam or ham by comparing the content with a dictionary constructed with all the words from the emails. We then analyze each email content and generate a feature vector indicating which words are contained.

We manually divided the data in two groups to create the testing and the training sets.

```
[('to', 1804), ('you', 1381), ('the', 994), ('and', 717), ('in', 671), ('is', 665), ('for', 555), ('my', 529),
 ('of', 499), ('me', 487), ('your', 485), ('on', 402), ('have', 381), ('that', 366), ('are', 349), ('it', 327),
 ('or', 318), ('call', 314), ('be', 304), ('at', 303), ('with', 302), ('not', 296), ('will', 283), ('get', 279),
 ('so', 262), ('ur', 260), ('can', 259), ('but', 236), ('You', 222), ('from', 216), ('do', 212), ('if', 211),
 ('up', 205), ('go', 203), ('just', 201), ('we', 191), ('when', 190), ('all', 189), ('this', 186), ('like', 185),
 ('know', 183), ('out', 181), ('got', 176), ('now', 172), ('was', 167), ('come', 167), ('am', 164), ('Call', 133),
 ('want', 132), ('time', 132), ('by', 130), ('only', 129), ('about', 127), ('need', 125), ('send', 124), ('then', 121),
 ('still', 120), ('going', 118), ('what', 118), ('But', 117), ('How', 117), ('as', 113), ('its', 111), ('If', 111),
 ('one', 110), ('he', 110), ('text', 107), ('So', 107), ('Just', 106), ('been', 106), ('our', 105), ('No', 102),
 ('has', 102), ('We', 100), ('some', 99), ('good', 99), ('no', 98), ('see', 98), ('any', 96), ('Do', 95), ('think', 95),
 ('love', 95), ('there', 94), ('an', 93), ('What', 91), ('how', 90), ('tell', 90), ('home', 87), ('back', 86),
 ('free', 86), ('take', 85), ('Ok', 85), ('her', 85), ('Your', 85), ('day', 84), ('dont', 83), ('My', 82), ('who', 81),
 ('mobile', 81), ('And', 80), ('give', 79), ('The', 77), ('phone', 77), ('new', 76), ('Have', 76), ('they', 76),
 ('FREE', 75), ('much', 74), ('him', 74), ('more', 74), ('make', 73), ('Are', 69), ('To', 69), ('reply', 68),
 ('later', 68), ('ask', 67), ('great', 66), ('txt', 65), ('Good', 65), ('had', 64), ('meet', 64), ('she', 64),
 ('say', 64), ('should', 63), ('Hey', 62), ('here', 62), ('This', 61), ('number', 60), ('claim', 60), ('da', 59),
 ('them', 59), ('after', 59), ('find', 59), ('too', 59), ('really', 59), ('contact', 59), ('way', 57), ('Can', 57),
 ('said', 57), ('would', 57), ('Its', 56), ('sent', 55), ('night', 55), ('Pls', 55), ('pick', 54), ('doing', 54),
 ('work', 54), ('miss', 54), ('Ur', 54), ('very', 54), ('week', 53), ('Please', 53), ('It', 53), ('ok', 53), ('right', 53),
 ('last', 53), ('Hi', 52), ('every', 52), ('did', 52), ('gonna', 51), ('stop', 50), ('tomorrow', 50), ('Txt', 50),
 ('feel', 49), ('told', 49), ('per', 49), ('Sorry', 48), ('where', 48), ('Did', 48), ('Oh', 47), ('could', 47),
 ('cos', 47), ('were', 47), ('sure', 46), ('Am', 46), ('Not', 46), ('his', 45), ('before', 45), ('Then', 45), ('keep', 45),
 ('let', 45), ('around', 45), ('hope', 44), ('buy', 44), ('When', 44), ('next', 44), ('He', 44), ('Reply', 44),
 ('cant', 44), ('already', 43), ('Many', 43), ('always', 43), ('Happy', 43), ('Is', 43), ('cash', 42), ('msg', 42),
 ('Yeah', 42), ('message', 42), ('being', 41), ('down', 41), ('us', 41), ('Gud', 41), ('wan', 41), ('Hope', 40),
 ('why', 40), ('dun', 40), ('Nokia', 40), ('wait', 40), ('which', 40), ('anything', 39), ('place', 39), ('other', 39),
 ('even', 39), ('service', 39), ('Dear', 39), ('please', 39), ('care', 38), ('STOP', 38), ('today', 38), ('someone', 38),
 ('something', 38), ('Lol', 38), ('off', 38), ('went', 37), ('Text', 37), ('won', 37), ('leave', 37), ('coming', 37),
 ('Get', 36), ('trying', 36), ('Will', 36), ('getting', 36), ('first', 36), ('awarded', 36), ('prize', 36), ('dear', 35),
 ('thing', 35), ('also', 35), ('Where', 35), ('She', 35), ('waiting', 35), ('try', 35), ('As', 34), ('Thanks', 34),
 ('finish', 34), ('over', 34), ('things', 34), ('people', 33), ('win', 33), ('few', 33), ('life', 33), ('im', 33),
 ('Or', 33), ('Free', 33), ('sleep', 32), ('For', 32), ('In', 32), ('bit', 31), ('That', 31), ('happy', 31), ('dat', 31),
 ('draw', 31), ('having', 30), ('Send', 30), ('car', 30), ('Yes', 30), ('another', 30), ('stuff', 30), ('use', 30),
 ('shows', 30), ('receive', 30), ('wat', 29), ('Now', 29), ('morning', 29), ('customer', 29), ('thk', 29), ('late', 29),
 ('well', 29), ('lunch', 29), ('pls', 29), ('thought', 29), ('money', 29), ('talk', 28), ('YOU', 28), ('Well', 28), ('long', 28), ('than', 28),
```

Here is the dictionary (part of it) that we created from the training data. And below is the beginning of the feature extracted from the training data. This is the feature vector that indicate the presence of the words.

```
[ 0.]
[ 1.]
[ 0.]
[ 0.]
[ 0.]
[ 0.]
[ 0.]
[ 0.]
[ 0.]
```

After extracting the features from both the training and the testing sets, we use both Naïve Bayes and SVM to fit the models to the training set. The first confusion matrix corresponds to the SVM and the second to the multinomial Naïve Bayes.

```
[[245  2]
 [ 7 31]]
[[247  0]
 [ 6 32]]
```

We can see that the Naïve Bayes model gives slightly better results compared to the SVM. For tasks such as classification and regression that do not require the joint distribution, discriminative models

can yield superior performance. However, in this case the two models have pretty much the same results. As we know that Generative models are harder to implement, Discriminative models work very well in classification and regression task. This example shows that the two techniques are two different approach of the problem, and that each one has its perks.