

Enron emails

(Text analytics)

Tool: Bag of words model

~~Other~~ ~~Boys~~, CRR, Random forest
Naive Bayes

The Analytics Edge:

Predictive coding provides an innovative way for lawyers to review large sets of documents. lawsuits often mean the review of letters, memoranda, files, to determine which ones are relevant. With more email communications, this data has exploded. It is no longer economical to review these individually to find the typically smaller set of sensitive documents. Predictive coding uses sample set reviewed by senior attorney & analytics help predict from the many documents, the more & less sensitive one

Overview

Enron, an American energy company was founded in 1985. It was originally involved in transmitting and distributing electricity and natural gas throughout the US. Before 2001, Enron employed nearly 20000 staff, with claimed revenues of nearly \$111 billion during 2000.

By the end of 2001, it was revealed that the reported financial condition came about due to systematic accounting fraud (the Enron scandal).

Enron Corpus

The Enron Corpus is a publicly available dataset of the email messages sent or received by about 150 senior managers of Enron. This data was obtained by the Federal Energy Regulatory Commission during its investigation. This corpus is one of the few publicly available mass collection of real emails available for study as such collections often are bound by privacy laws.

Predictive coding

Predictive coding allows software to take information entered by people and generalize it to a larger group of documents.

Traditionally, one searches a set of documents for a term to see and identify documents relevant to a case. Rather predictive coding considers the context. Instead of having legal experts going through all the documents, one uses a training set from a broader set which is categorized by legal experts. Then use of predictive analytics tools help predict the categorization of future documents.

California energy crisis

The California energy crisis in 2000-2001 was a situation in which California had a shortage of electric supply caused by market manipulations, illegal shutdown of pipelines by the Texas energy consortium Enron. The state suffered from multiple blackouts, while Enron gained the market which was deregulated.

The Federal Energy Regulatory Commission (FERC) investigated Enron's involvement in the crisis and this ultimately led to a \$1.52 billion settlement.

As part of the investigation FERC released the emails from some of the executives at Enron.

One approach is to search for keywords like "electricity bid", "energy schedule" in emails and then carefully review which ones are responsive.

Predictive coding using text analytics is an alternate & newer way to do this. The data comes from a 2010 Text Retrieval Conference - legal track where attorneys labeled emails responsive to energy schedules or bids.

Analytics on Enron dataset

```
energy <- read.csv("energy-bid.csv", stringsAsFactors = FALSE)
```

Str(email)

The data frame contains a total of 851 emails where

read(email)

\$email = Content of email

\$responsive = $\begin{cases} 1 & \text{if email is responsive to query about energy bids \& schedules} \\ 0 & \text{otherwise} \end{cases}$

energy\$email[1]

lists out the entire email but this is hard to read

strwrap(email[1])

This command wraps character strings to format paragraphs to be easy to read

energy\$responsive[1]

This takes a value 0 since the email is not responsive to energy bid & schedule

strwrap (energy & email [4])

energy & response [4]

The fourth email deals with the search on the energy bids & schedules. Hence the response variable takes a value of 1 for this email

table (energy & response)

	0	1
	714	137

} Many of the emails do not relate to the energy bids & schedules.

Preprocessing

library (tm)

corpus <- Corpus (VectorSource (energy & email))

Creates the corpus

strwrap (as.character (corpus [[4]]))

Read the email in the corpus

Corpus \leftarrow tm-map (Corpus, ~~tm-map (Corpus, tolower)~~)

Corpus \leftarrow tm-map (Corpus, remove Words,
Stopwords ("english"))

Corpus \leftarrow tm-map (Corpus, remove Punctuation)

Corpus \leftarrow tm-map (Corpus, remove Numbers)

Corpus \leftarrow tm-map (Corpus, Stem Document)

strwrap (as.character (Corpus[[4]]))

Note the email is much harder for us
as humans to make sense of after the
preprocessing but it is much more
suitable to apply analytics

freq \leftarrow DocumentTermMatrix (Corpus)

The document-term matrix has 851
documents, ~~13770~~¹³⁷⁷⁰ terms. The matrix is
again very sparse (99 %)

freq ← remove Sparse Terms (freq, 0.97)

Removes all terms that do not occur in
at least 3% of documents.

We now have 612 terms left.

energy_sparse ← as.data.frame(as.matrix(freq))

energy_sparse\$response ← energy\$response

colnames(energy_sparse) ← make.names(colnames(
energy_sparse))

library(caTools)

set.seed(1978)

spl ← sample.split(energy_sparse\$response,
SplitRatio = 0.7)

train ← subset(energy_sparse, spl == TRUE)

test ← subset(energy_sparse, spl == FALSE)

Base model

table (train \$ response)

0	1
500	96

} Majority of emails are non-responsive
Note false negatives are more critical here than false positives.

table (test \$ response)

0	1
214	41

} Baseline model will give accuracy = $\frac{214}{214+41} = 0.83$

CART

library (rpart)

library (rpart. plot)

set.seed (1)

model 1 ← rpart (as.factor (response) ~ . ,
data = train)

prp (model 1, type = 4, extra = 2)

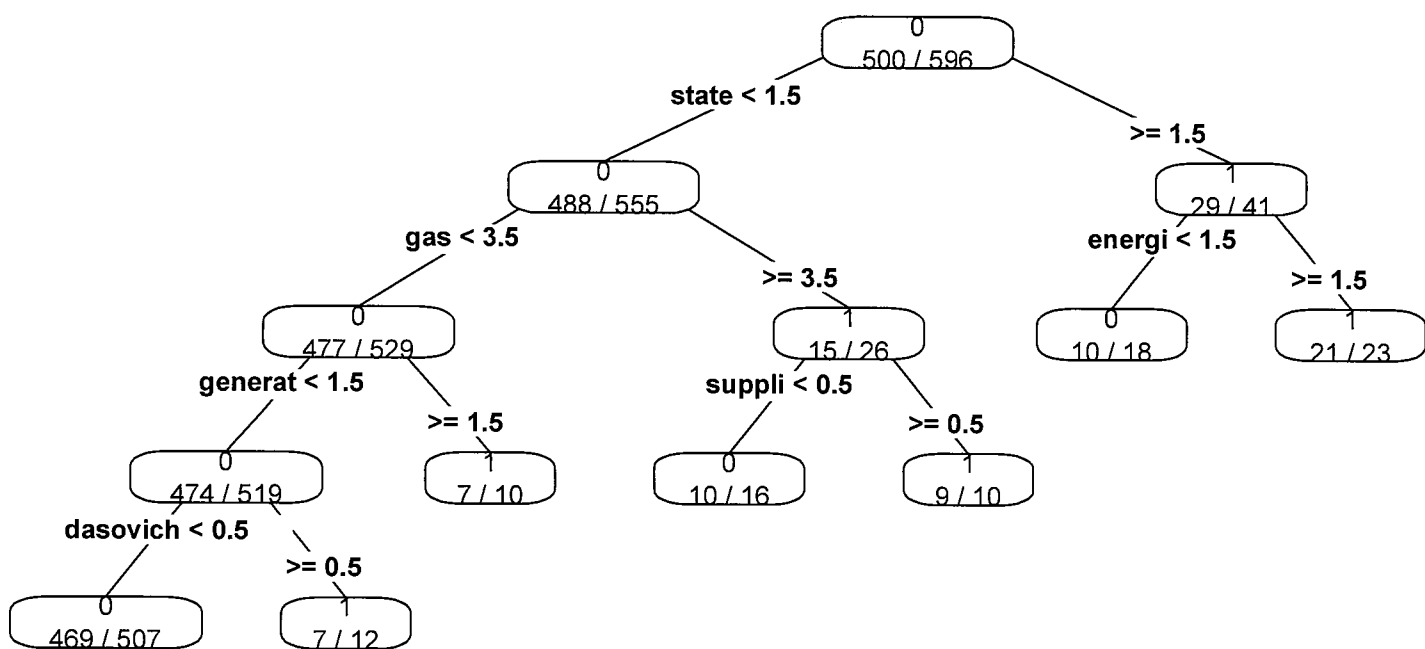
summary (model 1)

predict 1 ← predict (model 1, newdata = test, type = "class")

table (predict 1, test \$ response)

	0	1
0	199	26
1	15	15

Same accuracy as baseline but true positives have increased.



Randon forest

library (random Forest)

Set. seed (1979)

model 2 \leftarrow Random Forest (as. factor (response) ~ .,
data = train)

```
predict 2 <- predict (model 2, newdata = test,
                      type = "class")
```

Table (predict 2, test & responsive)

		Actual	
		0	1
Predict	0	204	26
	1	10	15

$$\text{Accuracy} = 0.858$$

has improved.

Also we now do a better job of identifying true negatives.

var Used (model 2)

This helps identify by frequency which predictors are used in the random forest.

Order 2 \leftarrow Sort (varUsed(model2), index.return = TRUE)
Returns sorted frequency & indices

```
names(test[, tail(order2$ixc)])
```

~~Jeff Richard Pleas attach forward Subject~~

library (ROCR)

```
predict prob1 <- predict (model1, newdata = test,  
                           type = "prob")
```

```
predict prob2 <- predict (model2, newdata = test,  
                           type = "prob" )
```

```
pred1 <- prediction (predict prob1[,2], test$response)
```

```
pred2 <- prediction (predict prob2[,2], test$response)
```

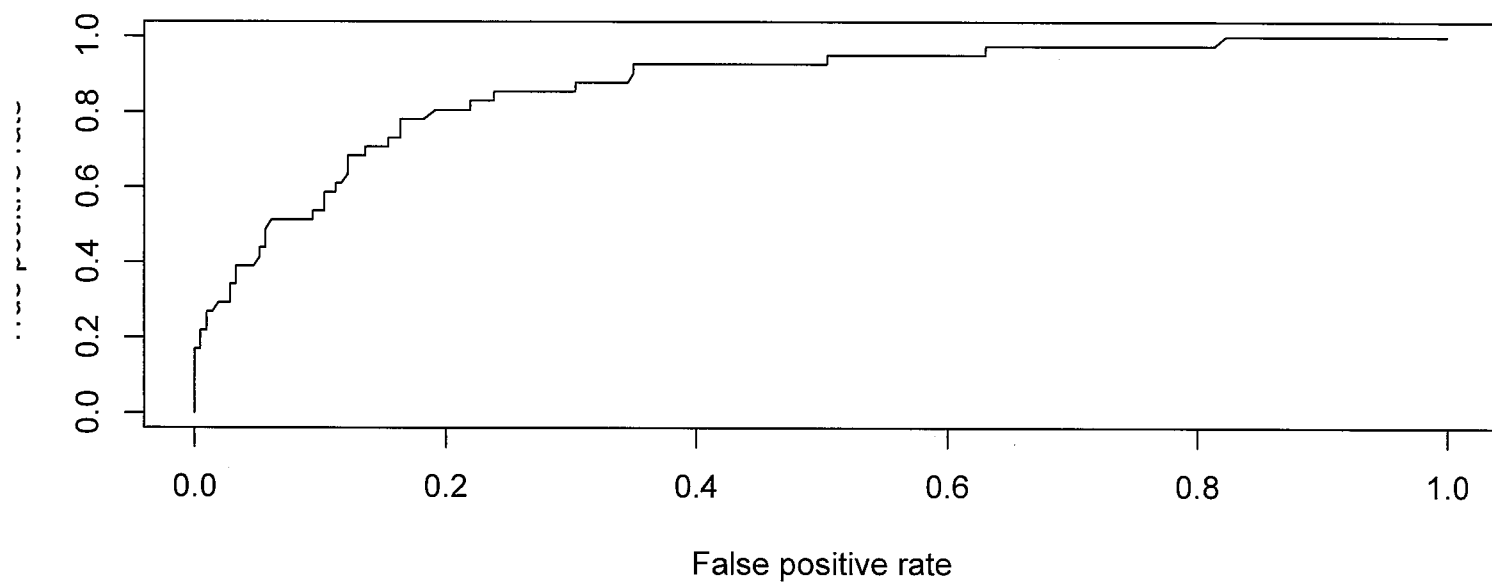
```
performance (pred1, measure = "auc")
```

```
performance (pred2, measure = "auc")
```

```
plot (performance (pred2, measure = "tpr", x.measure = "fpr"),  
      AUC of random forest is very good)
```

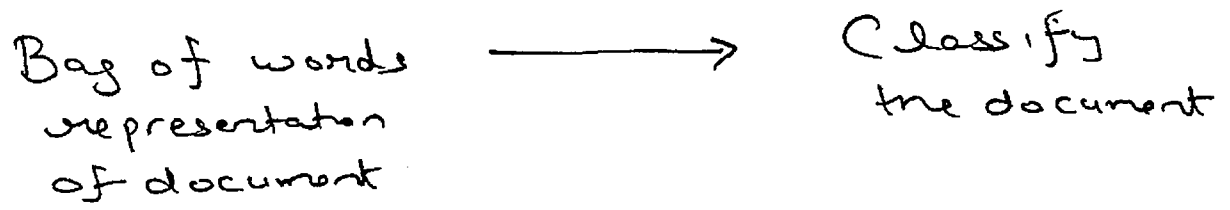
Here AUC is around 0.86.

AUC is a good measure across thresholds.



Naive Bayes Classifier

The naive Bayes classifier is a simple classification rule based on Bayes rule and it can be used with the simple bag of words representation of text.



$$P(\text{Document } d \in \underbrace{\text{Class } k}_{C_k} \mid \underbrace{x_{i1}, \dots, x_{ip}}_{\text{Features of document}})$$

$$P(C_k \mid x_{i1}, \dots, x_{ip}) = \frac{P(C_k \cap (x_{i1}, \dots, x_{ip}))}{P(x_{i1}, \dots, x_{ip})}$$

Naive Bayes classifier assumes independence of the features given the category

$$P(C_k \mid x_{i1}, \dots, x_{ip}) \propto P(C_k) P(x_{i1} \mid C_k) \dots P(x_{ip} \mid C_k)$$

For example: A fruit is an apple (category) if it is red, round and around 3 cm in diameter.

A Naive Bayes classifier considers each of these features to contribute to the probability that the fruit is an apple

$$\therefore P(C_k | x_{i1}, \dots, x_{ip}) = \frac{1}{Z} P(C_k) \prod_{i=1}^p P(x_{ki} | C_k)$$

where $Z = P(x_{i1}, \dots, x_{ip})$ depends only on factors

The prediction rule is to assign a class k such that

$$\arg \max_{k=1, \dots, K} \underbrace{P(C_k)}_{\downarrow} \underbrace{\prod_{i=1}^p P(x_{ki} | C_k)}_{\downarrow}$$

Estimated from training set
and then applied on test set

For example, assuming Gaussian prediction variables:

$$P(x_{ki} | C_k) = \frac{1}{\sqrt{2\pi} \sigma_{ki}} e^{-\frac{(x_{ki} - \mu_{ki})^2}{2\sigma_{ki}^2}}$$

where μ_{ki} and σ_{ki}^2 are mean and variances of values in category k for variable \tilde{x}_{ki} .

Naïve Bayes

$$P(C_k | x_1, \dots, x_p) = \frac{P(C_k) P(x_1 | C_k) \dots P(x_p | C_k)}{\sum_c P(C_c) P(x_1 | C_c) \dots P(x_p | C_c)}$$

In general,

$$\begin{aligned} P(C_k | x_1, \dots, x_p) &= \frac{P(C_k) P(x_1, \dots, x_p | C_k)}{P(x_1, \dots, x_p)} \\ &= \frac{P(C_k) P(x_1 | x_2, \dots, x_p, C_k) P(x_2 | x_3, \dots, x_p, C_k) \dots P(x_p | C_k)}{P(x_1, \dots, x_p)} \end{aligned}$$

To estimate $P(x_1, \dots, x_p | C_k)$ where each x_i takes 2 values and C_k takes 2 values, we need to estimate approximately $2(2^p)$ parameters (probabilities).

In Naive Bayes, we use & assume:

$$P(x_1, \dots, x_p | C_k) = P(x_1 | C_k) \dots P(x_p | C_k),$$

Now we need to estimate approximately

$2p(2)$ parameters (probabilities)

Naive Bayes makes the problem much smaller to solve & estimate. It works with less data too.

Naive Bayes classifier

install.packages("e1071")
library(e1071)

This provides a package to use a naive Bayes classifier.

model3 <- naiveBayes(as.factor(response) ~ ,
data = train)

naiveBayes(.) in the e1071 package computes the conditional posterior probabilities of a categorical variable given independent predictor variables using Bayes rule

Summary(model3)

model3 \$ apriori

0	1
500	96

 } Distribution of dependent variable

model3 \$ tables

Justs tables one for each predictor. For each numeric variable, it gives target class, mean & standard deviation of variable

predict3 <- predict(model3, newdata = test, type = "class")

table(predict3, test\$response)

	0	1
0	184	15
1	30	26

Accuracy is lower than logistic but true positives has increased.