

Statistical Analysis of Text - a mini-course

Text classification

Mattias Villani

Statistiska institutionen
Stockholms universitet

Institutionen för datavetenskap
Linköpings universitet



mattiasvillani.com



[@matvil](https://twitter.com/matvil)



[mattiasvillani](https://github.com/mattiasvillani)

- Text classification
- Regularization
- R's `tm` package (demo) [TMPackageDemo.R]

Supervised classification

- Predict the **class label** $s \in S$ using a set of **features**.
- Feature = Explanatory variable = Predictor = Covariate
- Binary classification: $s \in \{0, 1\}$
 - ▶ Movie reviews: $S = \{\text{pos}, \text{neg}\}$
 - ▶ E-mail spam: $S = \{\text{Spam}, \text{Ham}\}$
 - ▶ Bankruptcy: $S = \{\text{Not bankrupt}, \text{Bankrupt}\}$
- Multi-class classification: $s \in \{1, 2, \dots, K\}$
 - ▶ Topic categorization of web pages:
 $S = \{\text{'News'}, \text{'Sports'}, \text{'Entertainment'}\}$
 - ▶ POS-tagging: $S = \{\text{VB}, \text{JJ}, \text{NN}, \dots, \text{DT}\}$

Supervised classification, cont.

■ Example data:

- ▶ Larry Wall, born in British Columbia, Canada, is the original creator of the programming language Perl. Born in 1956, Larry went to ...
- ▶ Bjarne Stroustrup is a 62-years old computer scientist ...

Person	Income	Age	Single	Payment remarks	Bankrupt
Larry	10	58	Yes	Yes	Yes
Bjarne	15	62	No	Yes	No
⋮	⋮	⋮	⋮	⋮	⋮
Guido	27	56	No	No	No

■ Classification: construct prediction machine

Features \rightarrow Class label

■ More generally:

Features $\rightarrow \text{Pr}(\text{Class label}|\text{Features})$

Features from text documents

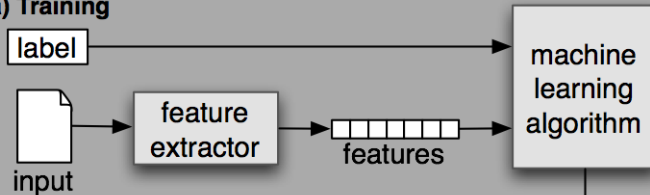
- Any quantity computed from a document can be used as a **feature**:
 - ▶ Presence/absence of individual words
 - ▶ Number of times an individual word is used
 - ▶ Presence/absence of pairs of words
 - ▶ Presence/absence of individual bigrams
 - ▶ Lexical diversity
 - ▶ Word counts
 - ▶ Number of web links from document, possibly weighted by Page Rank.
 - ▶ etc etc

Document	has('ball')	has('EU')	has('political_arena')	wordlen	Lex. Div.	Topic
Article1	Yes	No	No	4.1	5.4	Sports
Article2	No	No	No	6.5	13.4	Sports
⋮	⋮	⋮		⋮	⋮	⋮
ArticleN	No	No	Yes	7.4	11.1	News

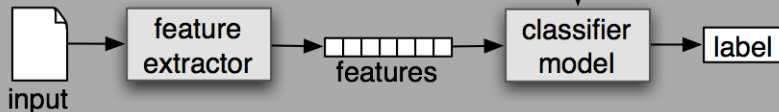
- Constructing clever **discriminating features** is the name of the game

Supervised learning for classification

(a) Training



(b) Prediction



The Bayesian classifier

■ Bayesian classification

$$\operatorname{argmax}_{s \in S} p(s|x)$$

where $x = (x_1, \dots, x_n)$ is a feature vector.

■ By Bayes' theorem

$$p(s|x) = \frac{p(x|s)p(s)}{p(x)} \propto p(x|s)p(s)$$

■ Bayesian classification

$$\operatorname{argmax}_{s \in S} p(x|s)p(s)$$

- $p(s)$ can be easily estimated from training data by relative frequencies.
- **Main problem:** Even with binary features [*has(word)*] the outcome space of $p(x|s)$ is huge (=data are sparse).

Naive Bayes

- **Naive Bayes (NB):** features are assumed independent

$$p(x|s) = \prod_{j=1}^n p(x_j|s)$$

- Naive Bayes solution

$$\operatorname{argmax}_{s \in S} \left[\prod_{j=1}^n p(x_j|s) \right] p(s)$$

- With binary features, $p(x_j|s)$ can be easily estimated by

$$\hat{p}(x_j|s) = \frac{C(x_j, s)}{C(s)}$$

- Example: $s = \text{news}$, $x_j = \text{has('ball')}$

$$\hat{p}(\text{has(ball)}|\text{news}) = \frac{\text{Number of news articles containing the word 'ball'}}{\text{Number of news articles}}$$

Naive Bayes

- **Continuous features** (e.g. lexical diversity) can be handled by:
 - ▶ Replacing continuous feature with several binary features ($1 \leq \text{lexDiv} < 2$, $2 \leq \text{lexDiv} \leq 10$ and $\text{lexDiv} > 10$)
 - ▶ Estimating $p(x_j|s)$ by a density estimator (e.g. kernel estimator)
- Finding the **most discriminatory features**. Sort from largest to smallest

$$\frac{p(x_j|s = \text{pos})}{p(x_j|s = \text{neg})} \text{ for } j = 1, \dots, n.$$

- **Problem with NB**: features are seldom independent \Rightarrow double-counting the evidence of individual features.
- **Advantages of NB**: simple and fast, yet often surprising accurate classifications.

Multinomial regression

- Logistic regression (Maximum Entropy/**MaxEnt**):

$$p(s = 1|x) = \frac{\exp(x'\beta)}{1 + \exp(x'\beta)}$$

- Classification rule: Choose $s = 0$ if $p(s|x) < 0.5$ otherwise choose $s = 1$.
- ... at least when consequences of different choices of s are the same. Loss/Utility function.
- Multinomial regression for multi-class data with K classes

$$p(s = s_j|x) = \frac{\exp(x'\beta_j)}{\sum_{k=1}^K \exp(x'\beta_k)}$$

- Classification

$$\operatorname{argmax}_{s \in \{s_1, \dots, s_K\}} p(s|x)$$

- $P \times (S - 1)$ number of coefficients
- Classification with text data is like any multi-class regression

Regularization - Variable selection

- Select a subset of the covariates.
- Old school: **Forward** and **backward selection**.
- New school: **Bayesian variable selection**.
- For each β_i introduce binary indicator I_i such that

$I_i = 1$ if covariate is in the model, that is $\beta_i \neq 0$

$I_i = 0$ if covariate is in the model, that is $\beta_i = 0$

- Use Markov Chain Monte Carlo (MCMC) simulation to approximate $\Pr(I_i | \text{Data})$ for each i .
- Example $S = \{\text{News, Sports}\}$. $\Pr(\text{News} | x)$.

	has('ball')	has('EU')	has('political_arena')	wordlen	Lex. Div.
$\Pr(I_i \text{Data})$	0.2	0.90	0.99	0.01	0.85

Regularization - Shrinkage

- Keep all covariates, but **shrink** their β -coefficient to zero.
- **Penalized likelihood**

$$L_{Ridge}(\beta) = LogLik(\beta) - \lambda \beta' \beta$$

where λ is the **penalty parameter**.

- Maximize $L_{Ridge}(\beta)$ with respect to β . Trade-off of fit ($LogLik(\beta)$) against complexity penalty $\beta' \beta$.
- **Ridge regression** if regression is linear.
- The penalty can be motivated as a **Bayesian prior**
 $\beta_i \stackrel{iid}{\sim} N(0, \lambda^{-1})$.
- λ can be estimated by cross-validation or Bayesian methods.

Lasso - Shrinkage and variable selection

- Replace Ridge penalty

$$L_{Ridge}(\beta) = \text{LogLik}(\beta) - \lambda \sum_{j=1}^n \beta_j^2$$

by

$$L_{Lasso}(\beta) = \text{LogLik}(\beta) - \lambda \sum_{j=1}^n |\beta_j|$$

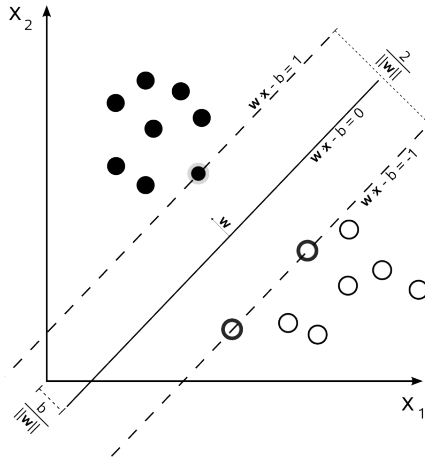
- The β that maximizes $L_{Lasso}(\beta)$ is called the **Lasso estimator**.
- Some parameters are shrunk exactly to zero \Rightarrow Lasso does **both shrinkage and variable selection**.
- Lasso penalty is equivalent to a double exponential prior

$$p(\beta_i) = \frac{\lambda}{2} \exp(\lambda |\beta_i - 0|)$$

Support vector machines

- One of the best off-the-shelf classifiers around.
- Finds the line in covariate space that maximally separates the two classes.
- When the points are not linearly separable: add a slack-variable $\xi_i > 0$ for each observation. Allow misclassification, but make it costly.
- Non-linear separating curves can be obtained by basis expansion (think about adding x^2 , x^3 and so on)
- The kernel trick makes it possible to handle many covariates.
- **Drawback:** not so easily extended to multi-class.
- `svm` function in R-package `e1071` [or `nltk.classify.svm`]

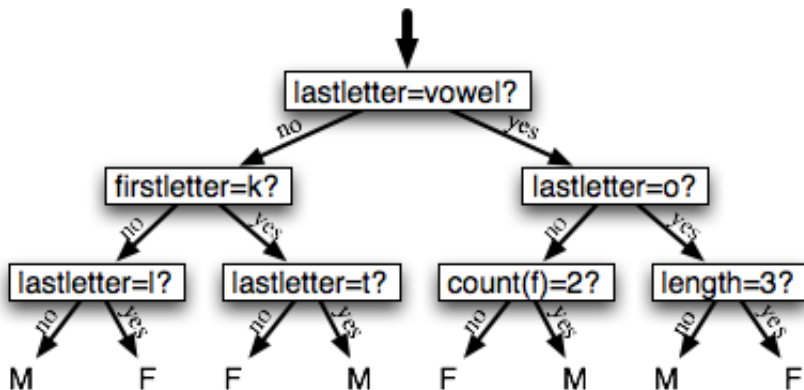
Linear SVMs



Regression trees and random forest

- Binary partitioning tree.
- At each internal node decide:
 - ▶ Which covariate to split on
 - ▶ Where to split the covariate ($X_j < c$. Trivial for binary covariates)
- The optimal splitting variables and split-points are chosen to minimize the mis-classification rate (or other similar measures).
- **Random forest (RF)** predicts using an average of many small trees.
- Each tree in RF is grown on a random subset of variables. Makes it possible to handle **many covariates**. Parallel.
- Advantage of RF: better predictions than trees.
- RF harder to interpret, but provide variable importance scores.
- **R packages:** tree and rpart (trees), randomForest (RF).

Regression trees



Evaluating a classifier: Accuracy and Error

■ Confusion matrix:

		Truth	
		Spam	Not Spam
Decision	Spam	tp	fp
	Not Spam	fn	tn

- tp = true positive, fp = false positive
- fn = false negative, tn = true negative
- **Accuracy** is the proportion of correctly classified items

$$\text{Accuracy} = \frac{tp + tn}{tp + tn + fn + fp}$$

- **Error** is the proportion of wrongly classified items

$$\text{Error} = 1 - \text{Accuracy}$$

Accuracy can be misleading

- **Accuracy is problematic when tn is large.** High accuracy can then be obtained by not acting at all!

		Truth	
		Spam	Not Spam
Choice	Spam	0	0
	Not Spam	100	900

Evaluating a classifier: the F-measure

■ Confusion matrix:

		Truth	
		Spam	Good
Choice	Spam	tp	fp
	Good	fn	tn

- **Precision** = proportion of selected items that the system got right

$$\text{Precision} = \frac{tp}{tp+fp}$$

- **Recall** = proportion of spam that the system classified as spam

$$\text{Recall} = \frac{tp}{tp+fn}$$

- **F-measure** is a harmonic mean between Precision and Recall

$$F = \frac{1}{\alpha \frac{1}{\text{Precision}} + (1 - \alpha) \frac{1}{\text{Recall}}}$$