

Introduction to Machine Learning for Social Scientists

Class 8: Re-sampling and intro to regularization

Edgar Franco Vivanco

Stanford University
Department of Political Science

edgarf1@stanford.edu

Summer 2018

Group project

Lottery: Monday July 30th

15% of your grade

Submit your team by tomorrow

On Group Project

- ▶ Detailed instructions on Monday.
- ▶ Final presentation and a final memo (due Aug 15th).
- ▶ Topics: **Random Forest, Neural Networks, Support Vector Machines, Bagging and boosting, Topic Modeling**

So far:

- ▶ Supervised learning.
- ▶ Regression and classification.
- ▶ Performance.
- ▶ In sample and out of sample.

Next weeks:

- ▶ Today: Re-sampling, text analysis ($p \geq n$).
- ▶ Next week: Regularization.
- ▶ ggplot and data manipulation workshop (Friday August 3rd)
- ▶ Supervised learning.
- ▶ Cluster analysis.
- ▶ Guest speakers.

Today:

- ▶ Key concepts:
 - ▶ Validation approach
 - ▶ LOOCV
 - ▶ K-fold
 - ▶ Document-term matrix and the problem of high dimensional data
- ▶ Key R methods:
 - ▶ `cv.glm()`

Evaluating fit:

- ▶ In sample: dependent variable in “training” data

Evaluating fit:

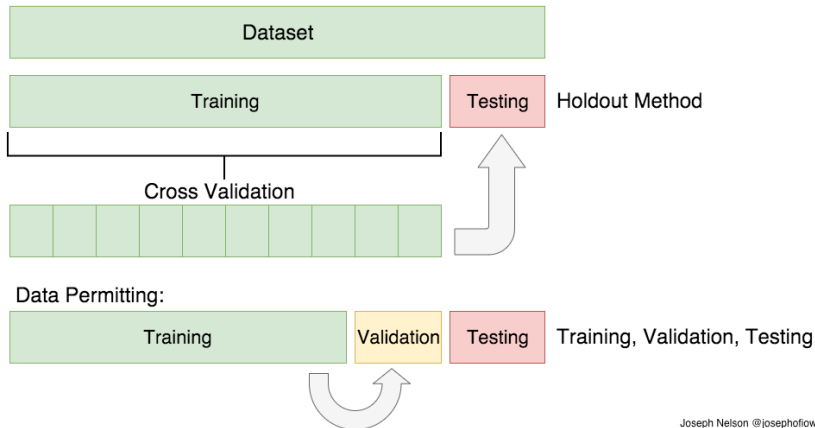
- ▶ In sample: dependent variable in “training” data
- ▶ But in general, we do not really care how well the method works on the training data. Rather, we are interested in the accuracy of the predictions that we obtain when we apply our method to previously unseen test data.

Evaluating fit:

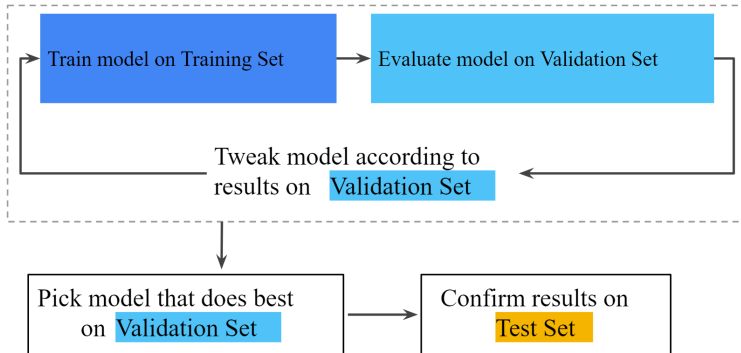
- ▶ In sample: dependent variable in “training” data
- ▶ But in general, we do not really care how well the method works on the training data. Rather, we are interested in the accuracy of the predictions that we obtain when we apply our method to previously unseen test data.
- ▶ Out of sample: **held out** data, test data

Evaluating fit:

- ▶ In sample: dependent variable in “training” data
- ▶ But in general, we do not really care how well the method works on the training data. Rather, we are interested in the accuracy of the predictions that we obtain when we apply our method to previously unseen test data.
- ▶ Out of sample: **held out** data, test data
- ▶ Best practice: evaluate t with **gold standard** data



Joseph Nelson @josephoflowa



Training error vs Test error

- Recall the distinction between the **test error** and the **training error**.

Training error vs Test error

- ▶ Recall the distinction between the **test error** and the **training error**.
- ▶ **Test error**: is the average error that results from using a statistical learning method to predict the response on a new observation, one that was not used in training the method.
- ▶ **Training error**: can be easily calculated by applying the statistical learning method to the observations used in its training.

Training error vs Test error

- ▶ Recall the distinction between the **test error** and the **training error**.
- ▶ **Test error**: is the average error that results from using a statistical learning method to predict the response on a new observation, one that was not used in training the method.
- ▶ **Training error**: can be easily calculated by applying the statistical learning method to the observations used in its training.
- ▶ But the training error rate often is quite different from the test error rate, and in particular the former can dramatically underestimate the latter.

Re-sampling: Steps

1. Randomly divide the data into training and validation

Re-sampling: Steps

1. Randomly divide the data into training and validation
2. Fit the model on the training set

Re-sampling: Steps

1. Randomly divide the data into training and validation
2. Fit the model on the training set
3. The fitted model is used to predict the observations in the validation set

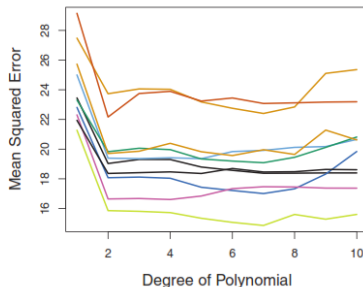
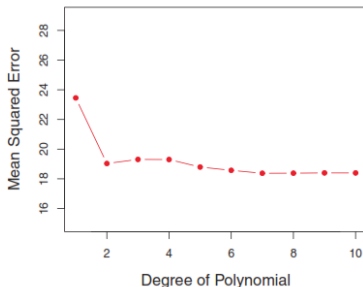
Re-sampling: Steps

1. Randomly divide the data into training and validation
2. Fit the model on the training set
3. The fitted model is used to predict the observations in the validation set
4. The resulting validation set error rate typically assessed using MSE in the case of a quantitative response provides an estimate of the test error rate.



Validation set approach:

- ▶ Linear vs Polynomial
- ▶ Split observations into two sets of same size.

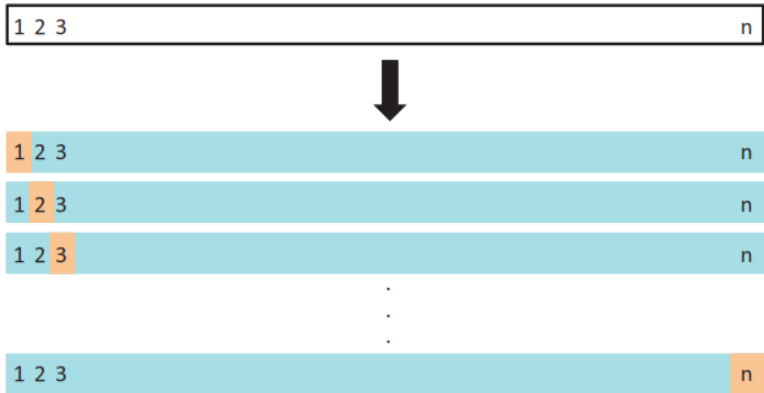


Two potential drawbacks:

- ▶ The validation estimate of the test error can be highly variable.
- ▶ *Overestimates* the test error for the model fit on the entire dataset (because we use few observations in training set).

Leave-One-Out Cross-Validation

- ▶ Like the validation set approach, LOOCV involves splitting the set of observations into two parts.
- ▶ However, instead of creating two subsets of comparable size, a single observation (x_1, y_1) is used for the validation set, and the remaining observations $(x_2, y_2), \dots, (x_n, y_n)$ make up the training set.
- ▶ The statistical learning method is trained on the $n - 1$ training observations, and a prediction \hat{y}_1 is made for the excluded observation, calculating the error $(y_1 - \hat{y}_1)^2$
- ▶ We repeat the process n times
- ▶ We estimate a mean-cross validation error.



The LOOCV estimate for the test MSE is the average of these n test error estimates.

$$CV_{(n)} = \frac{1}{n} \sum_{i=1}^n (y_i - \hat{y}_i)^2$$

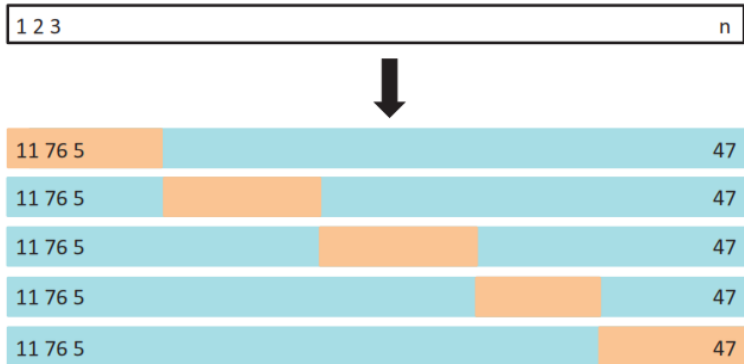
LOOCV

Advantages:

- ▶ Less bias: Because we use $n-1$ observations, the LOOCV tends not to overestimated the test error rate.
- ▶ Since there is no randomness in the training/validation set splits, LOOCV will always yield the same results.

K-fold Cross Validation

- ▶ This approach involves randomly dividing the set of observations into k groups, or folds of approximately equal size.
- ▶ The first fold is treated as the validation set and the method is fit on the remaining $k-1$ folds.
- ▶ The mean squared error is then computed on the observations in the held-out fold (MSE_1)
- ▶ The procedure is repeated k times; each time, a different group of observations is treated as a validation set.
- ▶ This process results in k estimates of the test error $MSE_1, MSE_2, \dots, MSE_k$



K-fold Cross Validation

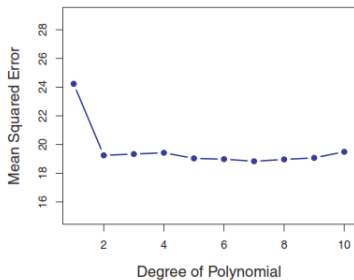
The k-fold CV estimate is computed by averaging these values:

$$CV_{(k)} = \frac{1}{k} \sum_{i=1}^k MSE_i$$

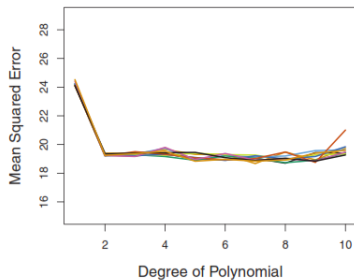
K-fold Cross Validation

- ▶ LOOCV is a special case of k-fold in which k is set equal n .
- ▶ Advantages of K-fold:
 - ▶ Computational speed
 - ▶ Intermediate level of bias and variance.
 - ▶ More bias than LOOCV but less than validation approach
 - ▶ Lower variance than LOOCV.

LOOCV



10-fold CV



Cross Validation on classification problems

- ▶ We can use cross-validation when Y is qualitative/
- ▶ We use the number of misclassified observations:

$$CV_{(n)} = \frac{1}{n} \sum_{i=1}^n I(y_i \neq \hat{y}_i)$$

R!



Text analysis

- ▶ A large amount of social interaction occurs in texts:
 - ▶ Congressional speeches, press releases, newsletters,
 - ▶ Facebook posts, tweets. emails, cell phone records.
 - ▶ Newspapers, magazines, news broadcasts,
 - ▶ Treaties, sermons, fatwas
- ▶ Pre 2000's social scientists avoided using texts/speech.
- ▶ Why?

Text analysis

- ▶ A large amount of social interaction occurs in texts:
 - ▶ Congressional speeches, press releases, newsletters,
 - ▶ Facebook posts, tweets. emails, cell phone records.
 - ▶ Newspapers, magazines, news broadcasts,
 - ▶ Treaties, sermons, fatwas
- ▶ Pre 2000's social scientists avoided using texts/speech.
- ▶ Why?
 - ▶ Hard to find
 - ▶ Time consuming
 - ▶ Not generalizable (each new data set ... new coding scheme)
 - ▶ Difficult to store/search
 - ▶ Idiosyncratic to coders/researchers
 - ▶ Statistical methods/algorithms, computationally intensive

Text analysis

Today:

- ▶ Massive increase in availability of unstructured text: In 2017, the number of emails sent and received per day total over 260 billion.
- ▶ Cheap storage: 1981: \$ 500,000 per GB. 2017: \$ 0.03 per GB
- ▶ Explosion in methods and programs to analyze texts
 - ▶ Generalizable
 - ▶ Systematic
 - ▶ Cheap

Preprocessing Texts

One (of many) recipe for preprocessing: retain **useful** information

Preprocessing Texts

One (of many) recipe for preprocessing: retain **useful** information

- 1) Remove capitalization, punctuation

Preprocessing Texts

One (of many) recipe for preprocessing: retain **useful** information

- 1) Remove capitalization, punctuation
- 2) Discard Word Order: (**Bag of Words** Assumption)

Preprocessing Texts

One (of many) recipe for preprocessing: retain **useful** information

- 1) Remove capitalization, punctuation
- 2) Discard Word Order: (**Bag of Words** Assumption)
- 3) Discard stop words

Preprocessing Texts

One (of many) recipe for preprocessing: retain **useful** information

- 1) Remove capitalization, punctuation
- 2) Discard Word Order: (**Bag of Words** Assumption)
- 3) Discard stop words
- 4) Combine similar terms: Stem, Lemmatize

Preprocessing Texts

One (of many) recipe for preprocessing: retain **useful** information

- 1) Remove capitalization, punctuation
- 2) Discard Word Order: (**Bag of Words** Assumption)
- 3) Discard stop words
- 4) Combine similar terms: Stem, Lemmatize
- 5) Discard less useful features \rightsquigarrow depends on application

Preprocessing Texts

One (of many) recipe for preprocessing: retain **useful** information

- 1) Remove capitalization, punctuation
- 2) Discard Word Order: (**Bag of Words** Assumption)
- 3) Discard stop words
- 4) Combine similar terms: Stem, Lemmatize
- 5) Discard less useful features \rightsquigarrow depends on application
- 6) Other reduction, weighting

Preprocessing Texts

One (of many) recipe for preprocessing: retain **useful** information

- 1) Remove capitalization, punctuation
- 2) Discard Word Order: (**Bag of Words** Assumption)
- 3) Discard stop words
- 4) Combine similar terms: Stem, Lemmatize
- 5) Discard less useful features \rightsquigarrow depends on application
- 6) Other reduction, weighting
- 7) **Output**: Count vector, each element counts occurrence of terms

Preprocessing Texts

One (of many) recipe for preprocessing: retain **useful** information

- 1) Remove capitalization, punctuation
- 2) Discard Word Order: (**Bag of Words** Assumption)
- 3) Discard stop words
- 4) Combine similar terms: Stem, Lemmatize
- 5) Discard less useful features \rightsquigarrow depends on application
- 6) Other reduction, weighting
- 7) **Output**: Count vector, each element counts occurrence of terms

1. Remove capitalization, punctuation, numbers

Assumption: capitalization, punctuation does not provide useful information.

1. Remove capitalization, punctuation, numbers

Assumption: capitalization, punctuation does not provide useful information.

Now we are engaged in a great civil war, testing
whether that nation, or any nation

1. Remove capitalization, punctuation, numbers

Assumption: capitalization, punctuation does not provide useful information.

Now we are engaged in a great civil war, testing
whether that nation, or any nation

now we are engaged in a great civil war testing
whether that nation or any nation

1. Remove capitalization, punctuation, numbers

Assumption: capitalization, punctuation does not provide useful information.

Now we are engaged in a great civil war, testing
whether that nation, or any nation

now we are engaged in a great civil war testing
whether that nation or any nation

Caution

‘‘Turkey’’ = ‘‘turkey’’

2. Discard Word Order (Bag of Words) \rightsquigarrow Tokenize

Assumption: Word Order Doesn't Matter.

2. Discard Word Order (Bag of Words) \rightsquigarrow Tokenize

Assumption: Word Order Doesn't Matter.

now we are engaged in a great civil war testing
whether that nation or any nation

2. Discard Word Order (Bag of Words) \rightsquigarrow Tokenize

Assumption: Word Order Doesn't Matter.

now we are engaged in a great civil war testing
whether that nation or any nation

[now, we, are, engaged, in, a, great, civil, war,
testing, whether, that, nation, or, any, nation]

2. Discard Word Order (Bag of Words) \rightsquigarrow Tokenize

Assumption: Word Order Doesn't Matter.

now we are engaged in a great civil war testing
whether that nation or any nation

[now, we, are, engaged, in, a, great, civil, war,
testing, whether, that, nation, or, any, nation]

[a, any, are, civil, engaged, great, in, nation, now,
or, testing, that, war, we, whether]

2. Discard Word Order (Bag of Words) \rightsquigarrow Tokenize

Assumption: Word Order Doesn't Matter.

now we are engaged in a great civil war testing
whether that nation or any nation

[now, we, are, engaged, in, a, great, civil, war,
testing, whether, that, nation, or, any, nation]

[a, any, are, civil, engaged, great, in, nation, now,
or, testing, that, war, we, whether]

Tokenization

Tokenization

Unigrams now we are engaged in a great civil war
testing whether that nation or any nation

Tokenization

Unigrams now we are engaged in a great civil war
testing whether that nation or any nation

Bigrams [now we, we are, are engaged, engaged in, in
a, a great, great civil, civil war, war testing,
testing whether, whether that, that nation, nation
or, or any, any nation]

Tokenization

Unigrams now we are engaged in a great civil war
testing whether that nation or any nation

Bigrams [now we, we are, are engaged, engaged in, in
a, a great, great civil, civil war, war testing,
testing whether, whether that, that nation, nation
or, or any, any nation]

Trigrams [now we are, we are engaged, are engaged in,
engaged in a, in a great, a great civil, great civil
war, civil war testing, war testing whether, testing
whether that, whether that nation, that nation or,
nation or any, or any nation]

Document Term Matrix:

	T1	T2	T3	T4	T5	T6	T7	T8
Doc1	2	0	4	3	0	1	0	2
Doc2	0	2	4	0	2	3	0	0
Doc3	4	0	1	3	0	1	0	1
Doc4	0	1	0	2	0	0	1	0
Doc5	0	0	2	0	0	4	0	0
Doc6	1	1	0	2	0	1	1	3
Doc7	2	1	3	4	0	2	0	2

The problem of high dimensionality :

- ▶ Most traditional statistics for regression and classification are intended for the *low dimension* setting ($n \geq p$).

The problem of high dimensionality :

- ▶ Most traditional statistics for regression and classification are intended for the *low dimension* setting ($n \geq p$).
- ▶ New technologies allow to collect an almost unlimited number of feature measurements.

The problem of high dimensionality :

- ▶ Most traditional statistics for regression and classification are intended for the *low dimension* setting ($n \geq p$).
- ▶ New technologies allow to collect an almost unlimited number of feature measurements.
- ▶ Text analysis is an example of $n \leq p$.

The problem of high dimensionality :

- ▶ Most traditional statistics for regression and classification are intended for the *low dimension* setting ($n \geq p$).
- ▶ New technologies allow to collect an almost unlimited number of feature measurements.
- ▶ Text analysis is an example of $n \leq p$.
- ▶ What goes wrong in high dimensions?

The problem of high dimensionality :

- ▶ Most traditional statistics for regression and classification are intended for the *low dimension* setting ($n \geq p$).
- ▶ New technologies allow to collect an almost unlimited number of feature measurements.
- ▶ Text analysis is an example of $n \leq p$.
- ▶ What goes wrong in high dimensions?
- ▶ LASSO

Document Term Matrix:

American Political Science Review (2018) 112, 2, 358–375
doi:10.1017/S0003005417000570

© American Political Science Association 2017

Reading Between the Lines: Prediction of Political Violence Using Newspaper Text

HANNES MUELLER *Institut d'Anàlisi Econòmica*
CHRISTOPHER RAUH *University of Montreal*

This article provides a new methodology to predict armed conflict by using newspaper text. Through machine learning, vast quantities of newspaper text are reduced to interpretable topics. These topics are then used in panel regressions to predict the onset of conflict. We propose the use of the within-country variation of these topics to predict the timing of conflict. This allows us to avoid the tendency of predicting conflict only in countries where it occurred before. We show that the within-country variation of topics is a good predictor of conflict and becomes particularly useful when risk in previously peaceful countries arises. Two aspects seem to be responsible for these features. Topics provide depth because they consist of changing, long lists of terms that make them able to capture the changing context of conflict. At the same time, topics provide width because they are summaries of the full text, including stabilizing factors.

The conflict literature has made significant progress in understanding which countries are more at risk of suffering an armed conflict.¹ However, many factors that have been identified as leading to increased risk, like mountainous terrain or ethnic polarization, are time invariant or very slow moving, and therefore not useful in predicting the timing of conflict. Other factors, like GDP levels or political institutions, still vary more between countries than within countries over time. This means it is easier to predict whether a country is at risk in general rather than when a country is particularly at risk. Yet, understanding the timing of conflict is critical for policy.

An additional problem of forecasting the timing of armed conflict is that it is rare and at the same time relatively concentrated in some countries. This is problematic because it implies that the variation between countries can dominate the analysis unless the between- and

within-country variations are separated explicitly. Empirical models that are overall quite accurate can therefore be of little use on the time dimension. We show, using a simple panel regression model, that many variables commonly used in the literature indeed face this problem. This means they predict conflict where it occurred before, and therefore fail to predict conflicts in previously peaceful countries.

As a solution to this problem, we propose data generated from news sources. To this end, we implement an automated method to quantify the content of news using the latent Dirichlet allocation (LDA) model (Blei, Ng, and Jordan 2003), which we apply to over 700,000 newspaper articles from English-speaking newspapers. There are two advantages that topics have over existing methods of analyzing text. First, topics provide depth because, by design, they put words into context. The context can be useful for forecasting. Second, topics provide width because they allow us to use the whole text, including stabilizing factors, when forecasting conflict. This means we can let the data speak without losing interpretability of the results.

At the prediction stage, we rely on a simple panel regression model, which uses all generated topics as explanatory variables. The result is a model able to forecast out of sample the onset of civil war, armed conflict, and even movements of refugees a year before they occur. It relies entirely on news text and can therefore provide forecasts without the need to extrapolate or wait for other data sources. Furthermore, the procedure can be implemented with only minimal personal judgment and appears to generate consistent summaries of

Hannes Mueller is a tenured scientist at IAE (CSIC), Barcelona GSE Institut d'Anàlisi Econòmica, CSIC Campus UAB, 08193 Bellaterra, Spain (h.mueller.unf@gmail.com).

Christopher Rauh is an Assistant Professor at University of Montreal, Département de Sciences Économiques, Université de Montréal, CP6128 succ. Centre-Ville, Montréal H3C 3J7, Canada (christopher.rauh@umontreal.ca).

We thank Tim Besley, Melissa Dell, Vincenzo Galasso, Hector Galindo, Matt Gentzkow, Stephen Hansen, Ethan Kaptein, Daniel Ohayon, Akash Raja, Bernhard Reinsberg, Anand Shrivastava, Ron Smith, Jack Willis, Stéphane Wolton, and the participants of the workshops and conferences ENCoRe, Barcelona, Political Econ-

NEXT

- ▶ Lottery
- ▶ LASSO
- ▶ Reading on text analysis
- ▶ ggplot and data manipulation workshop (Friday August 3rd)