

Elements of Statistical Learning for Regression Analysis

Predictive Modeling & Statistical Learning

Gaston Sanchez

CC BY-SA 4.0

So Far ...

What's coming next

We've talked about Linear Regression—simple and multiple via OLS—from a traditional/classic perspective.

In order to continue introducing more contemporary (modernish) approaches, we need to discuss ideas like:

- ▶ modeling purposes
- ▶ measuring predictive accuracy
- ▶ bias-variance trade-off
- ▶ over-fitting
- ▶ learning and test sets
- ▶ resampling methods (cross-validation and bootstrapping)

Modeling ... What for?

Statistical Modeling: Two Cultures

Statistical Science
2001, Vol. 16, No. 3, 199–231

Statistical Modeling: The Two Cultures

Leo Breiman

Abstract. There are two cultures in the use of statistical modeling to reach conclusions from data. One assumes that the data are generated by a given stochastic data model. The other uses algorithmic models and treats the data mechanism as unknown. The statistical community has been committed to the almost exclusive use of data models. This commitment has led to irrelevant theory, questionable conclusions, and has kept statisticians from working on a large range of interesting current problems. Algorithmic modeling, both in theory and practice, has developed rapidly in fields outside statistics. It can be used both on large complex data sets and as a more accurate and informative alternative to data modeling on smaller data sets. If our goal as a field is to use data to solve problems, then we need to move away from exclusive dependence on data models and adopt a more diverse set of tools.

Understanding vs Prediction

Models for Understanding versus Models for Prediction

Gilbert Saporta

Chaire de statistique appliquée & CEDRIC, CNAM
292 rue Saint Martin, Paris, France *saporta@cnam.fr*

Abstract. According to a standard point of view, statistical modelling consists in establishing a parsimonious representation of a random phenomenon, generally based upon the knowledge of an expert of the application field: the aim of a model is to provide a better understanding of data and of the underlying mechanism which have produced it. On the other hand, Data Mining and KDD deal with predictive modelling: models are merely algorithms and the quality of a model is assessed by its performance for predicting new observations. In this communication, we develop some general considerations about both aspects of modelling.

To Explain or to Predict?

Statistical Science

2010, Vol. 25, No. 3, 289–310

DOI: 10.1214/10-STS330

© Institute of Mathematical Statistics, 2010

To Explain or to Predict?

Galit Shmueli

Abstract. Statistical modeling is a powerful tool for developing and testing theories by way of causal explanation, prediction, and description. In many disciplines there is near-exclusive use of statistical modeling for causal explanation and the assumption that models with high explanatory power are inherently of high predictive power. Conflation between explanation and prediction is common, yet the distinction must be understood for progressing scientific knowledge. While this distinction has been recognized in the philosophy of science, the statistical literature lacks a thorough discussion of the many differences that arise in the process of modeling for an explanatory versus a predictive goal. The purpose of this article is to clarify the distinction between explanatory and predictive modeling, to discuss its sources, and to reveal the practical implications of the distinction to each step in the modeling process.

Modeling for what?

Goals

Understanding -vs- Prediction

Increasing acknowledge about the Prediction-vs-Understanding spectrum

Paradox 1

Typical conception of a model:

$$Y = f(X; \theta) + \varepsilon$$

Models for understanding

A “good” statistical model is one in which $f()$ is a parsimonious function that helps us explain how Y is related to X ;

- ▶ we assume there is a “true” function $f()$ that we want to approximate
- ▶ we assume there is some stochastic mechanism that generates the data
- ▶ a model should be simple with interpretable parameters e.g. rate of change, elasticity, odds-ratio, etc.

Paradox 1

Paradox 1

A “good” statistical model does not necessarily give accurate predictions.

For instance, in epidemiology it is more important to find risk factors than having an accurate individual prediction of getting some disease.

Models for Prediction

Prediction?

A “good” statistical model is one that provides accurate predictions: predict new observations with “good” accuracy without necessarily providing an explanation about the data generation/variation mechanism.

- ▶ no need for a theory of consumer to predict marketing target
- ▶ a model may be just simply an algorithm

Paradox 2

Paradox 2

We can predict without understanding (i.e. model is a black box).

- ▶ The aim is not to approximate the true function $f()$
- ▶ The aim is to find $f()$ that performs well by predicting new observations.

Model Performance

Model Performance

How do we define what a “good” model is?

- ▶ A model that fits the data well?
(e.g. minimize resubstitution error)
- ▶ A model with optimal parameters?
(e.g. most likely coefficients)
- ▶ A model that adequately predicts new (unseen) observations?
(e.g. minimize generalization error)

In the Predictive Modeling arena ...

How do we measure prediction accuracy?

What do we mean by “prediction”?

Predictive accuracy

Assessing predictions

Observed -vs- Predicted
 y_i \hat{y}_i

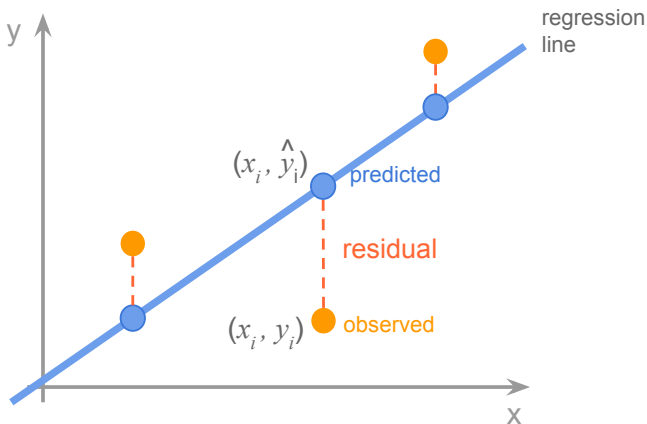
What measure of accuracy can we use?

Starting Point: Residuals

- ▶ The main idea consists of comparing observed inputs y_i against predicted inputs \hat{y}_i .
- ▶ We need to measure the discrepancy between observed inputs and predicted inputs.
- ▶ This means our starting point involves considering residuals $e_i = y_i - \hat{y}_i$
- ▶ The summary considered up to this point is the RSS:

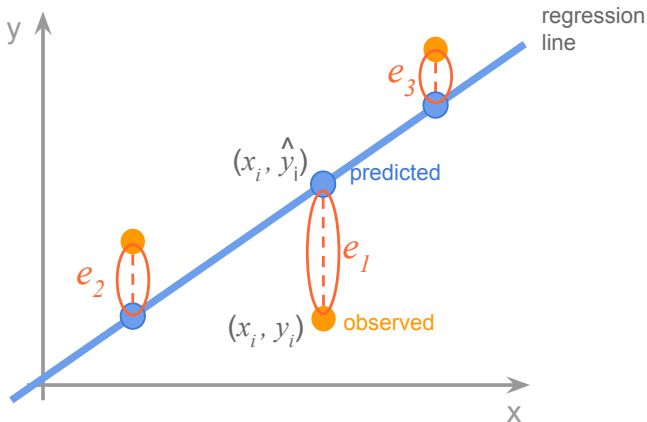
$$\text{RSS} = \sum_{i=1}^n e_i^2 = \sum_{i=1}^n (y_i - \hat{y}_i)^2$$

Looking at the amount of residuals



Measuring the overall size of the differences

Looking at the amount of residuals



How big are these differences?

Is RSS a good measure of model performance?

kind of ...

Residual Sum of Squares

One issue with the RSS is that it is a *total* value that depends on the number of residuals:

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

Preferably we would like to have a summary that is a “representative” measure of the **typical** error

Mean Squared Error (MSE)

A more “representative” measure of the typical error can be achieved by averaging RSS, getting what is called the **Mean Squared Error (MSE)**

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

Think of MSE as the “typical” squared error

Mean Squared Error (MSE)

- ▶ The MSE says how far typical points are above or below the regression line.
- ▶ Think of MSE as the typical size of prediction errors.
- ▶ In a general sense, MSE is a measure of scatter for residuals.
- ▶ MSE is a measure of how accurate our predictions are.

Root Mean Squared Error (RMSE)

A side effect of using squared residuals is that of having to deal with square units.

Some authors and practitioners prefer to take the square root of MSE in order to **recover the original units**. This produces the **Root Mean Squared Error (RMSE)**:

$$\text{RMSE} = \sqrt{\text{MSE}} = \sqrt{\frac{1}{n} \sum_{i=1}^n (y_i - \hat{y}_i)^2}$$

Mean Absolute Error (MAE)

MSE (and RMSE) are not the only possible options.

We can also consider the **absolute value** of residuals $|y_i - \hat{y}_i|$.
and the corresponding **Mean Absolute Error** (MAE)

$$\text{MAE} = \frac{1}{n} \sum_{i=1}^n |y_i - \hat{y}_i|$$

Absolute value is more robust than quadratic value because squaring may exacerbate atypical values (e.g. very large residuals).

Course Assumption

*For this course, unless stated otherwise, we are going to use **MSE** as the standard measure of model performance.*

In the Predictive Modeling arena ...

Used (seen) observations

- ▶ Prediction \hat{y}_i for y_i that was used to build the model
- ▶ *resubstitution* error: $e_i = y_i - \hat{y}_i$
- ▶ “already seen” MSE (less honest measure of accuracy)

Unseen observations

- ▶ Prediction \hat{y}_0 for y_0 that was NOT used to build the model
- ▶ *generalization* error: $e_0 = y_0 - \hat{y}_0$
- ▶ “unseen” MSE (more honest measure of accuracy)

In the Predictive Modeling arena ...

- ▶ A “good” model is one which gives accurate predictions.
- ▶ By *predictions* we mean predictions of new data.
- ▶ Therefore we focus on the generalization ability of the model to predict unobserved data
- ▶ This involves finding a measure of accuracy for predictions.

Learning

Idea of Learning

If what we want is a model $f()$ that gives accurate predictions (i.e. predictions of unseen data), then **how** should we build $f()$?

By “how” I’m referring to the overall mechanics/strategy of the model building process.

This is where the concept of **learning** comes in.

Idea of Learning

- ▶ We are interested in getting a model with good generalization power
- ▶ We want our model to *learn* as much as possible
- ▶ It sounds reasonable to use as much data as possible so the obtained model is able to generalize adequately

Learning Issues

Building a model by using as many observations as possible looks like a good strategy.

And this strategy should work (in theory) as long as the used data is *representative* of the phenomenon under study.

The problem is: most of the time we don't know for sure if the available data is really representative.

Learning Issues

We run the risk that we use data that is not representative.

Perhaps the built model fits the used data extremely well, but it may lack generalization ability.

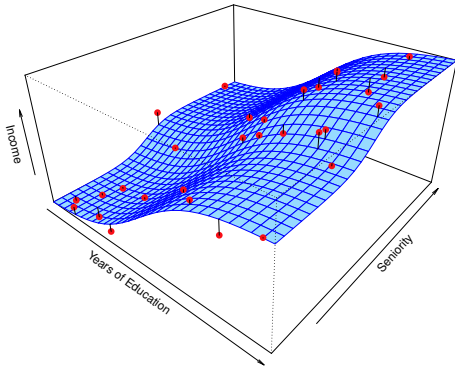
Moreover, by using all the available data we don't really have an honest measurement of the model accuracy (we only know the resubstitution MSE).

Over-Fitting Idea

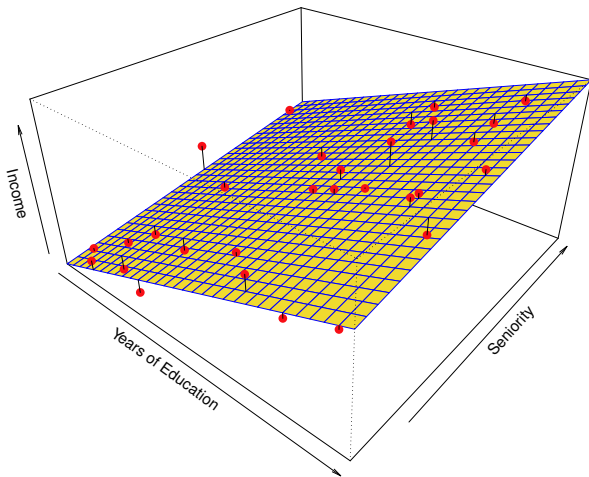
Summary



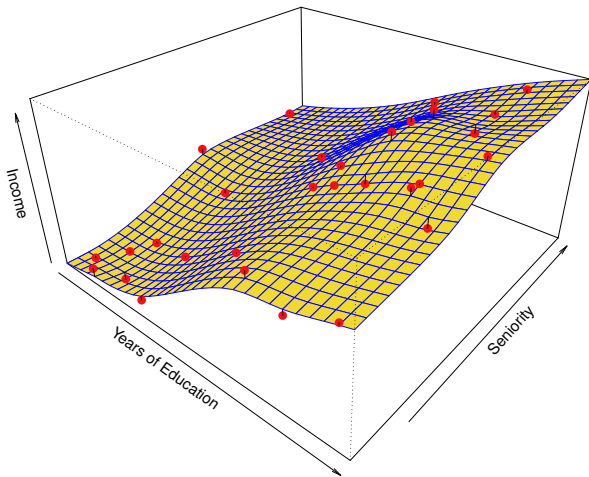
- ▶ A “good” model is one which gives accurate predictions.
- ▶ Predictions of observed data? Predictions of unobserved data?
- ▶ Prediction of unobserved data is different from fitting observed data.
- ▶ We refer to predictive performance over new observations (generalization).



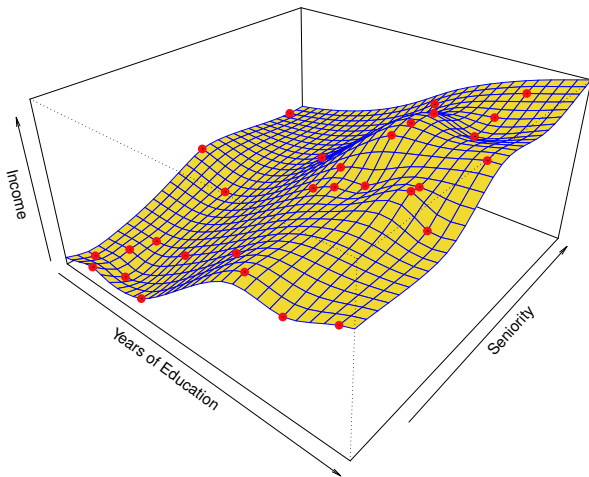
ISL Example (Chap 2): red points are simulated values for $\text{income} = f(\text{education}, \text{seniority}) + \epsilon$; where $f()$ is the blue surface



Underfitting: linear regression model fit to the simulated data.



OK Fitting: more flexible regression model fit with a thin-plate-spline.



Overfitting: even more flexible regression model fit with no errors.

Motivation

- ▶ Suppose you have to fit a regression model.
- ▶ You decide to use all your available data to get predicted values \hat{y}
- ▶ You use MSE to measure the predicting performance of the model.
- ▶ How reliable is the MSE value that you would obtain?
(How “honest” would be the MSE?)

Over-fitting

When you use all the available data to train (fit) a model, you run the risk to obtain a model that has learned not only the systematic part of the model, but also the unique noise of the data.

This situation is called **over-fitting**. And any measure of model performance will tend to be too optimistic.

Over-fitting

- ▶ A highly accurate model may suffer from overfitting
- ▶ A very robust model (rigid) won't be able to adequately fit the data

Motivation

Evaluating the model by using the data employed to fit the model produces a resubstitution or apparent measure of error.

We are not really evaluating the model with new/unseen data.

In other words, we are not really evaluating the generalization ability of the model.

Learning Dilemma

On one hand ...

- ▶ You want to use as much data as possible to train a model.
- ▶ You want to feed your model with as many examples as possible.

Learning Dilemma

On one hand ...

- ▶ You want to use as much data as possible to train a model.
- ▶ You want to feed your model with as many examples as possible.

On the other hand ...

- ▶ You also want to know how will your model behave when new input data is available?
- ▶ How will your model generalize (when predicting unseen data)?

Dilemma

We face a BIG dilemma

- ▶ On one hand, you want to use as much data as possible to train a model.
- ▶ On the other hand, you want to have new/unseen data to evaluate the performance of your model and see how well it generalizes.

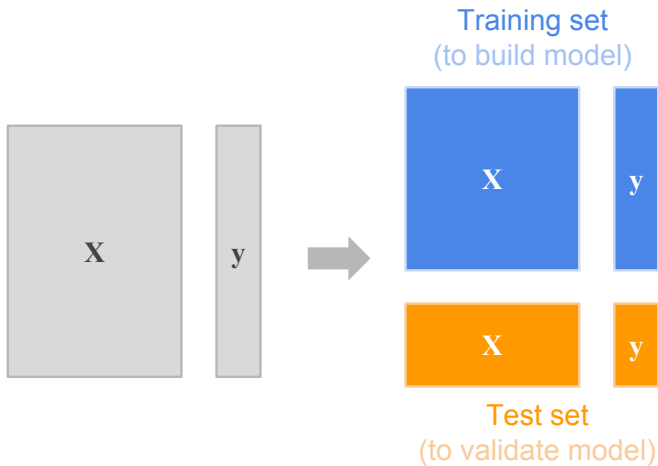
Training and Test Sets

The solution to this conundrum is to dispose of two data sets:

- ▶ **Training Dataset**: used to train the model
- ▶ **Test Dataset**: used to test the model and evaluate predicting performance

We assume that both the training and the test data sets are representative of the studied phenomenon. And that the test dataset is NOT used in the training stage.

Training and Test Sets

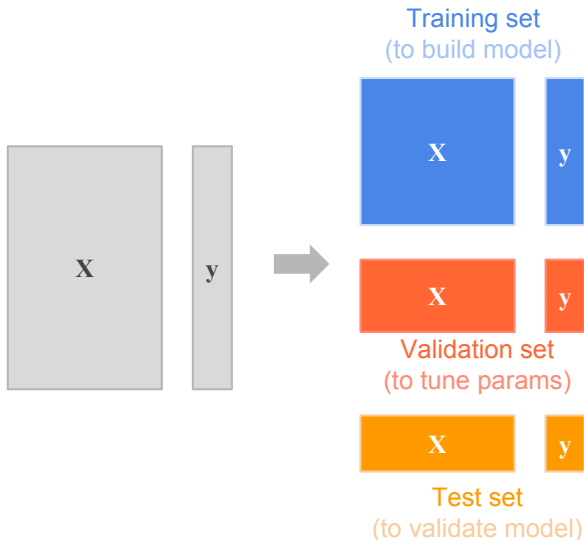


Training and Test Sets

Some authors go further and propose a three dataset scheme:

- ▶ Training set to fit the model
- ▶ Validation set to tune parameters
- ▶ Test set to assess predicting performance and select the best model

Training-Validation-Test Sets



Training and Test Sets

How do you decide what samples go into the training set, and what samples go into the test set?

Training and Test Sets

Selecting an adequate strategy to form the training and test sets obviously mainly depends on the amount of data.

With a vast amount of data, one could split the data set in halves, one half as the training set, the other half as the test set.

Training and Test Sets

The problem is: we don't always have vast amounts of data.

Either because the size of the data is small (“few observations”), or because the quality of the data is not good, and the best samples form a small subset.

Test-MSE

In order to have an “honest” measure of performance, we calculate **MSE on test data** (the so-called *test-MSE*).

Evaluating the predicting power on unseen data will give us a better idea of the model performance than when we just use the train-MSE.

Idea: Various Test Sets

Depending on how you form your train and test sets, you may end up with sets that are not truly representative of the studied phenomenon.

So instead of using just one test set, some authors propose to use several training-test sets.

Of course, with this suggestion we go back again to the issue of the size of the data. The solution to this limitation comes from using [re-sampling](#) methods.

Cross-Validation

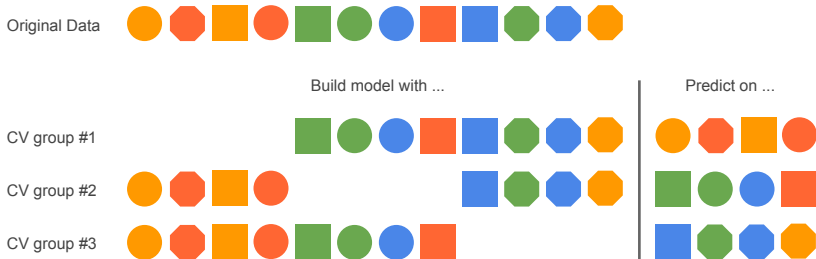
Cross-Validation (CV)

In Cross-Validation the main idea is to hold-out part of the data.

CV uses a systematic way of sampling the data.

- ▶ k-fold CV
- ▶ Leave-One-Out CV (loocv)
- ▶ 10-fold CV

Cross-validation three-fold



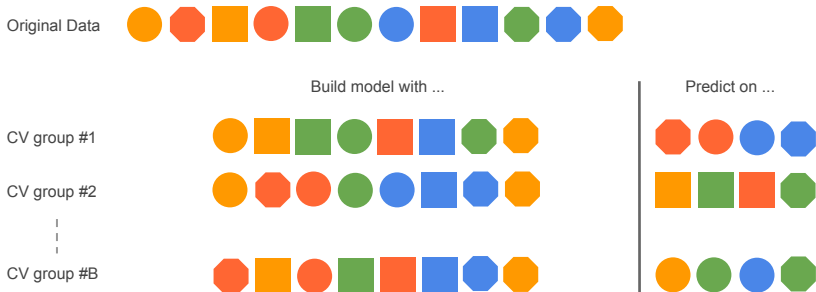
k -fold Cross-Validation

- ▶ The samples are randomly partitioned into k sets of roughly equal size.
- ▶ A model is fit using all the samples except the first subset.
- ▶ The hold-out samples are predicted by this model and used to estimate performance measures.
- ▶ The first subset is returned to the training set, and the procedure repeats with the second subset held out.
- ▶ The k resampled estimates of performance are summarized (usually with the mean and standard error).
- ▶ The choice of k is usually 5 or 10, but there is no formal rule.

Leave-One-Out Cross-Validation (loocv)

- ▶ A special version is the leave-one-out cross-validation.
- ▶ This is a special case for $k = 1$.
- ▶ Only one sample is held out at a time.
- ▶ The overall performance is calculated from the k individual held out predictions.

B Repeated cross-validation



B Repeated cross-validation

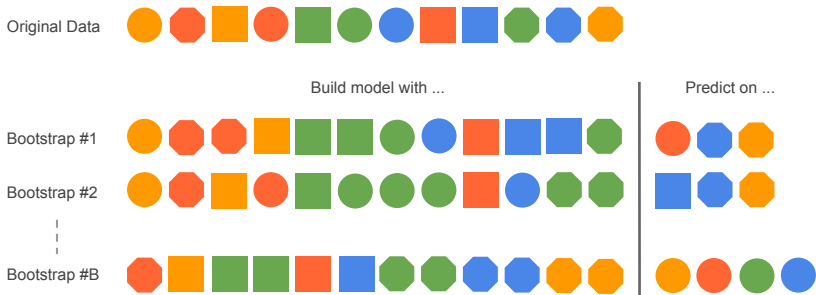
- ▶ Repeated training/set splits is also known as “leave-group-out” cross-validation.
- ▶ You simply create multiple splits of the data into training and test sets.
- ▶ A good rule of thumb is about 75-80% training, 25-20% test.
- ▶ Increasing the number of subsets has the effect of decreasing the uncertainty of the performance estimates.
- ▶ It is suggested to choose a larger number of repetitions (say 50-200).

The Bootstrap

Bootstrap Resampling

- ▶ A bootstrap sample is a random sample of the data taken *with replacement*.
- ▶ The bootstrap sample is the same size as the original data set
- ▶ As a result, some samples will be represented multiple times in the bootstrap sample while other will not be selected at all.
- ▶ The samples not selected are usually referred to as the out-of-bag samples.

Bootstrap Resampling



For a given iteration, a model is built on the selected samples and is used to predict the out-of-bag samples.

Bootstrap Resampling

- ▶ On average, 63.2% of the data points in the bootstrap sample are represented at least once.
- ▶ Bootstrap resampling has bias similar to k -fold cross-validation when $k = 2$.
- ▶ If the training set size is small, this bias may be problematic, but will decrease as the training set sample size becomes larger.

References

- ▶ **Statistical Modeling: The Two Cultures** by Leo Breiman (2001). *Statistical Science*, Vol 16 (3), 199-231.
- ▶ **Models for Understanding versus Models for Prediction** by Gilbert Saporta (2008). COMPSTAT 2008. Physica-Verlag.
- ▶ **To Explain or to Predict?** by Galit Shmueli (2010). *Statistical Science*, Vol 25 (3), 289-310.