



THE UNIVERSITY OF
SYDNEY

Lecture 6: Model evaluation and bootstrap sampling

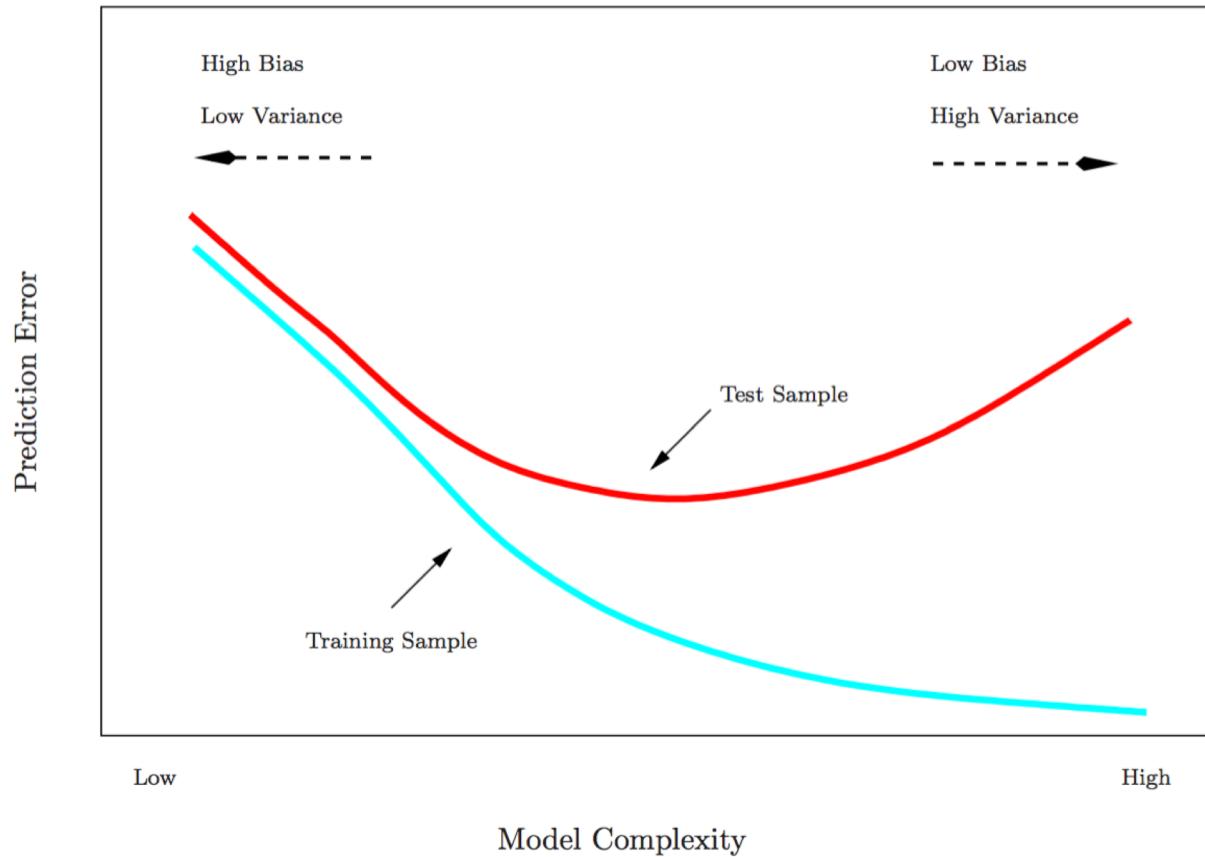
STAT5003

Pengyi Yang

Training error versus test error

- Recall the distinction between the *test error* and the *training error*
- The *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.
- In contrast, the *training error* can be easily calculated by applying the statistical learning method to the observations used in its training.
- The training error rate often is quite different from the test error rate, and in particular the former can dramatically *underestimate* the latter.

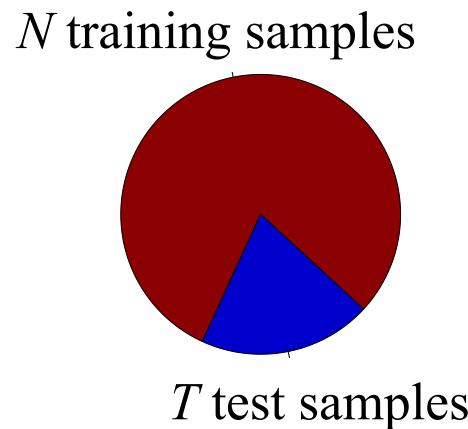
Training set versus test set error



Demonstration

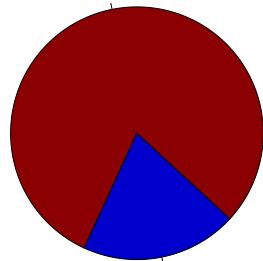
More on prediction error estimates

- Best solution: a large designated test set. Often not available
- Some methods (e.g. *BIC*) make a mathematical adjustment to the training error rate in order to estimate the test error rate. (discuss in future)
- Here we instead consider a class of methods that estimate the test error by *holding out* a subset of the training observations from the fitted process, and then applying the statistical learning method to those held out observations.



Validation set approach

N training samples



T test samples

- Here we randomly divide the available set of samples into two parts: a *training set* and a validation or *hold-out set*
- The model is fit on the training set, and the fitted model is used to predict the responses for the observations in the validation set.
- The resulting validation-set error provides an estimate of the test error. This is typically assessed using MSE in the case of a quantitative response and misclassification rate in the case of a qualitative (discrete) response.

The validation partition



A random splitting into two halves: left part is training set, right part is validation set.

Drawbacks of validation set approach

- The validation estimate of the test error can be highly variable, depending on precisely which observations are included in the training set and which observations are included in the validation set.
- In the validation approach, only a subset of the observations – those that are included in the training set – are used to fit the model. This suggests that the validation set error may tend to *overestimate* the test error for the model fit on the entire data set.

Why?

K -fold cross-validation

- *Widely used approach* for estimating test error.
- Estimates can be used to select best model, and to give an idea of the test error of the final chosen model.
- Idea is to randomly divide the data into K equal-sized parts. We leave out part k , fit the model to the other $K-1$ parts (combined), and then obtain predictions for the left-out k th part.
- This is done in term for each part $k=1, 2, \dots, K$, and then the results are combined.

K -fold cross-validation in detail

Divide data into K roughly equal-sized parts ($K=5$ here)

1	2	3	4	5
Validation	Train	Train	Train	Train

Cross-validation formulation

- Let the K parts be C_1, C_2, \dots, C_K , where C_k denotes the indices of the observations in part k . There are n_k observations in part k : if n is a multiple of K , then $n_k = n/K$.
- Compute

$$\text{CV}_{(K)} = \sum_{k=1}^K \frac{n_k}{n} \text{MSE}_k$$

where $\text{MSE}_k = \sum_{i \in C_k} (y_i - \hat{y}_i)^2 / n_k$, and \hat{y}_i is the fit for observation i , obtained from the data with part k removed.

- Setting $K = n$ yields n -fold or *leave-one out cross-validation* (LOOCV).

Cross-validation for classification problems

- We divide the data into K roughly equal-sized parts C_1, C_2, \dots, C_K . C_k denotes the indices of the observations in part k . There are n_k observations in part k : if n is a multiple of K , then $n_k = n/K$.
- Compute

$$\text{CV}_K = \sum_{k=1}^K \frac{n_k}{n} \text{Err}_k$$

where $\text{Err}_k = \sum_{i \in C_k} I(y_i \neq \hat{y}_i)/n_k$.

- The estimated standard deviation of CV_K is

$$\widehat{\text{SE}}(\text{CV}_K) = \sqrt{\sum_{k=1}^K (\text{Err}_k - \overline{\text{Err}_k})^2 / (K - 1)}$$

Overall classification accuracy rate

Overall classification accuracy:

$$ACC = 1 - \frac{1}{N} \sum_{i=1}^N I(y_i \neq \hat{y}_i)$$

Disadvantages:

- Makes no distinction about the type of errors being made. In spam filtering, the cost of erroneous deleting an important email is likely to be higher than incorrectly allowing a spam email past a filter.
- Does not consider the natural frequencies of each class.

The confusion matrix (two-class classification)

		Actual	
		True	False
Classifier	True	True Positive	False Positive
	False	False Negative	True Negative

Have cancer and predicted to have cancer

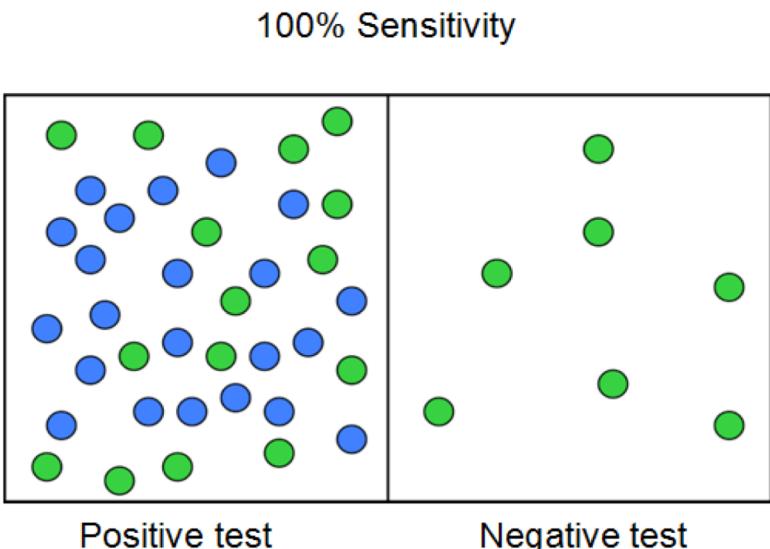
No cancer but predicted to have cancer

Have cancer but predicted to have no cancer

No cancer and predicted to have no cancer

Sensitivity and specificity

- Accuracy $ACC = \frac{(TP + TN)}{(TP + FP + FN + TN)}$
- Sensitivity $Sen = TP / (TP + FN)$
- Specificity $Spe = TN / (TN + FP)$

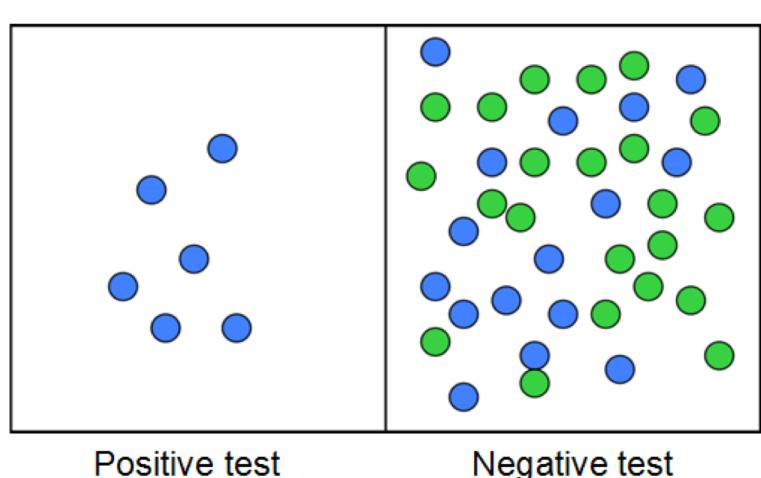


$$F_1 = \frac{2 * TP}{2 * TP + FP + FN}$$

(harmonic mean)

$$GM = \sqrt{\frac{TP}{TP + FN} * \frac{TP}{TP + FP}}$$

(geometric mean)



Blue = Individual that has cancer
Green = Individual that are normal

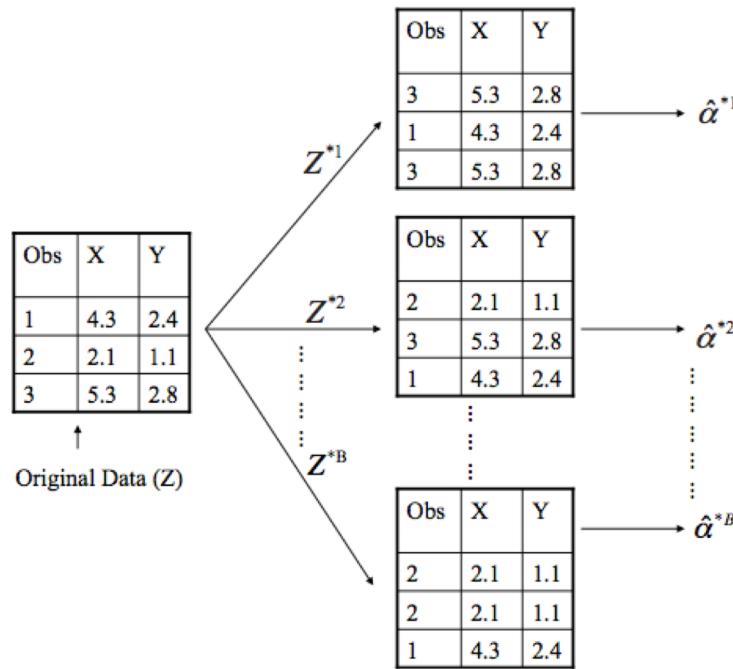
The Bootstrap

The Bootstrap

- The *bootstrap* is a flexible and powerful statistical tool that can be used to quantify the uncertainty associated with a given estimator or statistical learning method.
- For example, it can provide an estimate of the standard error of a coefficient, or a confidence interval for that coefficient.

$$y_i = \beta_0 + \beta_1 x_i + e_i$$

A simple example with 3 observations



A graphical illustration of the bootstrap approach on a small sample containing $n = 3$ observations. Each bootstrap data set contains n observations, sampled with replacement from the original data set. Each bootstrap data set is used to obtain an estimate of α .

A simple example

- Suppose that we wish to invest a fixed sum of money in two financial assets that yield returns of X and Y , respectively, where X and Y are random quantities.
- We will invest a fraction α of our money in X , and will invest the remaining $1 - \alpha$ in Y .
- We wish to choose α to minimize the total risk, or variance, of our investment. In other words, we want to minimize $\text{Var}(\alpha X + (1 - \alpha)Y)$.
- One can show that the value that minimizes the risk is given by

$$\alpha = \frac{\sigma_Y^2 - \sigma_{XY}}{\sigma_X^2 + \sigma_Y^2 - 2\sigma_{XY}},$$

where $\sigma_X^2 = \text{Var}(X)$, $\sigma_Y^2 = \text{Var}(Y)$, and $\sigma_{XY} = \text{Cov}(X, Y)$.

Example continued

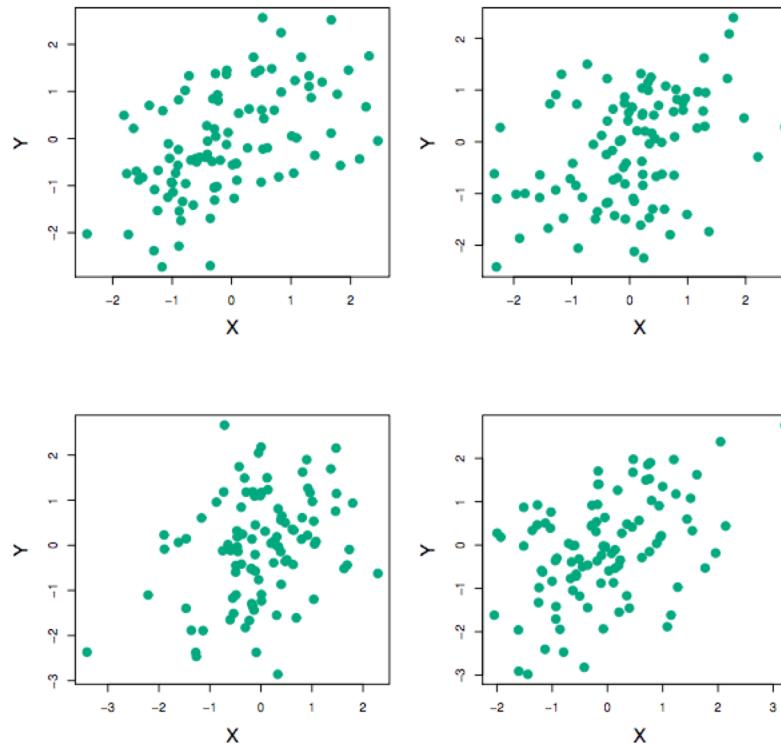
- But the values of σ_X^2 , σ_Y^2 , and σ_{XY} are unknown.
- We can compute estimates for these quantities, $\hat{\sigma}_X^2$, $\hat{\sigma}_Y^2$, and $\hat{\sigma}_{XY}$, using a data set that contains measurements for X and Y .
- We can then estimate the value of α that minimizes the variance of our investment using

$$\hat{\alpha} = \frac{\hat{\sigma}_Y^2 - \hat{\sigma}_{XY}}{\hat{\sigma}_X^2 + \hat{\sigma}_Y^2 - 2\hat{\sigma}_{XY}}.$$

- To estimate the standard deviation of $\hat{\alpha}$, we repeated the process of simulating 100 paired observations of X and Y , and estimating α 1,000 times.
- We thereby obtained 1,000 estimates for α , which we can call $\hat{\alpha}_1, \hat{\alpha}_2, \dots, \hat{\alpha}_{1000}$.

Simulations

- For these simulations the parameters were set to $\sigma_X^2 = 1$, $\sigma_Y^2 = 1.25$, and $\sigma_{XY} = 0.5$, and so we know that the true value of α is 0.6



Each panel displays 100 simulated returns for investments X and Y . From left to right and top to bottom, the resulting estimates for α are 0.576, 0.532, 0.657, and 0.651.

Example continued

- The mean over all 1,000 estimates for α is

$$\bar{\alpha} = \frac{1}{1000} \sum_{r=1}^{1000} \hat{\alpha}_r = 0.5996,$$

very close to $\alpha = 0.6$, and the standard deviation of the estimates is

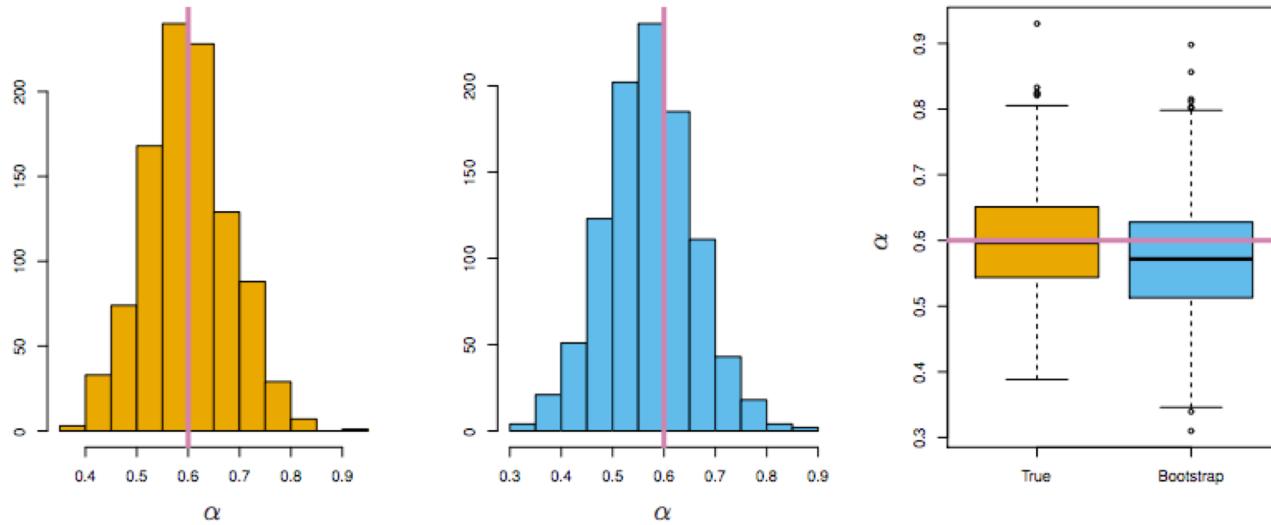
$$\sqrt{\frac{1}{1000 - 1} \sum_{r=1}^{1000} (\hat{\alpha}_r - \bar{\alpha})^2} = 0.083.$$

- This gives us a very good idea of the accuracy of $\hat{\alpha}$: $SE(\hat{\alpha}) \approx 0.083$.
- So roughly speaking, for a random sample from the population, we would expect $\hat{\alpha}$ to differ from α by approximately 0.08, on average.

The real world situation

- The procedure outlined above cannot be applied, because for real data we cannot generate new samples from the original population.
- However, the bootstrap approach allows us to use a computer to mimic the process of obtaining new data sets, so that we can estimate the variability of our estimate without generating additional samples.
- Rather than repeatedly obtaining independent data sets from the population, we instead obtain distinct data sets by repeatedly sampling observations from the original data set *with replacement*.
- Each of these “bootstrap data sets” is created by sampling *with replacement*, and is the *same size* as our original dataset. As a result some observations may appear more than once in a given bootstrap data set and some not at all.

Results



Left: A histogram of the estimates of α obtained by generating 1,000 simulated data sets from the true population. *Center:* A histogram of the estimates of α obtained from 1,000 bootstrap samples from a single data set. *Right:* The estimates of α displayed in the left and center panels are shown as boxplots. In each panel, the pink line indicates the true value of α .

Generic bootstrap procedure

- Denoting the first bootstrap data set by Z^{*1} , we use Z^{*1} to produce a new bootstrap estimate for α , which we call $\hat{\alpha}^{*1}$
- This procedure is repeated B times for some large value of B (say 100 or 1000), in order to produce B different bootstrap data sets, $Z^{*1}, Z^{*2}, \dots, Z^{*B}$, and B corresponding α estimates, $\hat{\alpha}^{*1}, \hat{\alpha}^{*2}, \dots, \hat{\alpha}^{*B}$.
- We estimate the standard error of these bootstrap estimates using the formula

$$\text{SE}_B(\hat{\alpha}) = \sqrt{\frac{1}{B-1} \sum_{r=1}^B (\hat{\alpha}^{*r} - \bar{\hat{\alpha}}^*)^2}.$$

- This serves as an estimate of the standard error of $\hat{\alpha}$ estimated from the original data set.

See figures on the previous slide. Bootstrap results are in blue.
For this example $\text{SE}_B(\hat{\alpha}) = 0.087$.

Schematic illustration of the bootstrap procedure

