

Intro To Statistical Learning: Notes

Mo D Jabeen

December 13, 2022

1 General

Statistics learning is based on making predictions or inferences on data inputs. Via approximating $f(x)$, where $y = f(x) + \text{error}$.

1.1 What are the types of statistical problems?

- Regression: Determine outcome variable based on predictors, continuous problem ie Range of values
- Classification: Discrete choice of answers (normally qualitative)
- Clustering: Determine similar groups of data (no natural output variable)

1.2 Notation

$$x_{ij}, i : 1, 2, \dots, n, j : 1, 2, \dots, p \quad (1.1)$$

$$x_i = (x_{i1}, x_{i2}, x_{i3} \dots, x_{ip}) \quad (1.2)$$

i is the observation and j is the predictor.

1.3 What is parametric and non parametric?

Parametric: Assume form of the desired function and calculate parameters based on the assumption.

Non Parametric: Do not make any explicit assumptions, instead fit function best to the data given.

Choosing a non parametric avoids the problem of having a function form very different to reality however opens up the possibility of overfitting (following noise too closely) and needs more data for an accurate form.

Furthermore, if the main goal is inference, restrictive (parametric) are much easier to interpret.

1.4 What is unsupervised learning?

There is no response/outcome variable, ie clustering.

Semi supervised learning: some response variables.

1.5 What is RSS?

Residual sum of squares is minimised in regression to calculate the parameters.

$$RSS = (y_1 - \beta_0 - \beta_1 x_1)^2 + \dots + (y_n - \beta_0 - \beta_1 x_n)^2 \quad (1.3)$$

2 Quality of fit

Measure how well the model fits the true model.

2.1 How is Mean Squared Error used?

Regression compares the predicted outcome to the true value and measures the MSE. However many statistical methods minimise the **training** MSE, therefore a test MSE should be used!

If there is a small training MSE and a large test MSE this can indicate overfitting.

$$MSE = \frac{1}{n} \sum_{i=1}^n (y_i - \hat{f}(x_i))^2 \quad (2.1)$$

2.2 What is Bias Variance Trade off?

MSE can be deconstructed to give variance of the estimate function, squared bias of the estimate function and variance of error.

$$E(y_o - \hat{f}(x_o))^2 = Var(\hat{f}(x_o)) + (Bias(\hat{f}(x_o)))^2 + Var(\epsilon) \quad (2.2)$$

Aim is to minimise variance and bias as error can not be removed.

The rate of change between var and bias determines optimal flexibility of a model.

2.2.1 What is variance?

The change in outcome/model if the data set is changed. Higher flexibility of model often increases this.

2.2.2 What is bias?

Error from approximations.

2.3 Error Rate

Classification accuracy is measured by error rate, the frequency of values that predicted the wrong class. In this case test error rate is also preferred.

$$Error\ rate = \frac{1}{n} \sum I(y_i \neq \hat{y}_i), I: y_i \neq \hat{y}_i? 1:0 \quad (2.3)$$

2.4 What is Bayes Classification?

Choosing the max probability an observation is a class will minimise the error rate.

If there are only two classes, choose the class that fits:

$$P(Y = 1|X = x_o) > 0.5 \quad (2.4)$$

2.5 What is Bayes error rate?

$$\text{Bayes Error Rate} = 1 - E(\max P(Y = j|X = x_o)) \quad (2.5)$$

E is the average for all values of X. ie if the probability of a 2 class, setup where one class is 0.7 the error rate will be 0.3. Not possible to actually use Bayes as the probability is unknown, however, this is the gold standard.

3 Classification

3.1 What is K nearest neighbours?

KNN identifies the K points closest to training point x, shown as η_o . The conditional probability is then the fraction of the points in η_o that equal class j.

$$P(Y = j|X = X_o) = \frac{1}{K} \sum_{i \in \eta_o} I(y_i = j) \quad (3.1)$$

Then classify point x as the class with the highest probability. K=1 has low bias but high variance.

3.2 Why logistic regression instead of linear regression?

To create a regression problem from a classification, you would be forced to create some type of ordering of the qualitative variables if >2. This ordering may not be true for the data set. However with a binary classification least squares is possible, the issue is that least squares does not have the required boundaries for a binary decision ie 0-1 so you would create a useless areas.

3.3 What is logistic regression?

To bound p(X) between 0 and 1, exponential equation is used:

$$p(X) = P(Y = 1|X) \quad (3.2)$$

$$p(X) = \frac{e^{\beta_0 + \beta_1 X}}{1 + e^{\beta_0 + \beta_1 X}} \quad (3.3)$$

This will produce an S curve, and if the log is taken shows that 1 unit increase in X gives a change of log odds by β_1 .

$$\log \text{ odds} : \ln\left(\frac{p(X)}{1 - p(X)}\right) = \beta_0 + \beta_1 X \quad (3.4)$$

This also shows that if β_1 is positive, increasing X will increase p(X), if negative the increasing X will decrease p(X).

3.4 How is logistic regression used?

For binary class scenarios you are trying to maximise the likelihood function to estimate the beta parameters.

Likelihood function:

$$(\beta_0, \beta_1) = \prod_{i: Y_i=1} p(x_i) + \prod_{i: Y_i=0} (1 - p(x_i)) \quad (3.5)$$

3.5 What is standard error?

$$SE(\alpha^2) = x \quad (3.6)$$

3.6 shows for each sample expect α to vary by x on average.
The below is the standard error of the mean ie std/n :

$$SE(\mu)^2 = \frac{\sigma^2}{n} \quad (3.7)$$

The below is based on how the Betas are calculated (this is for linear regression):

$$SE(\hat{\beta}_0)^2 = \sigma^2 \left(\frac{1}{n} + \frac{\bar{x}^2}{\sum (x_i - \bar{x})^2} \right) \quad (3.8)$$

$$SE(\beta_1)^2 = \frac{\sigma^2}{\sum (x_i - \bar{x})^2} \quad (3.9)$$

3.6 What ways can you validate the parameters?

Use confidence intervals, z/t tests and hypothesis tests.

The null hypothesis is that there is no relationship between predictor and outcome.

3.7 What are dummy variables?

If there is qualitative category for the observational data, this can be used a dummy value which is used to represent the observations. This then requires new parameters to be calculated.

I.e. set x_i as 1 if male, 0 if female. If more than one category use multiple pairs as predictors, ie $x_{i1} = X == \text{Asian?} : 1 : 0, x_{i2} = X == \text{Jamican?} : 1 : 0$.

These parameters from this can be validated in the same way.

3.8 How do you include multiple predictors in logistic regression?

Match the number of predictors to the number of parameters:

$$\ln \frac{p(x)}{1 - p(x)} = \beta_0 + \beta_1 x_{i1} + \beta_2 x_{i2} + \cdots \beta_p x_{ip} \quad (3.10)$$

Multiple predictors can show how the interweaving between each other and the outcome. Not shown when not included if not included in the equation.

Multi-class logistic regression exists but is not really used.

3.9 What is linear discriminant analysis?

Model X for each class Y as a distribution, then use Bayes theorem and compare all dist to choose the max $P(Y=k|X=x)$. **Assumes normal distribution for each class.**

3.9.1 LDA Equation

Π_k : Probability of the randomly chosen observation X being in class K . (Fraction of observations that belong to the k th class)

$f_k(X)$: Density function for X belonging to class K .

$$P(Y = k|X = x) = \frac{\Pi_k f_k(x)}{\sum_{l=1}^K \Pi_l f_l(x)} \quad (3.11)$$

$$P(Y = k|X = x) = \frac{\text{Overall Prob * density of } x \text{ for } k}{\sum_{\text{all classes}} (\text{Overall Prob * density of } x \text{ for } k)} \quad (3.12)$$

The focus is estimating $f(x)$ to fit Bayes Classifier. A point on the density function should be maximised as shown by: **Assuming the var is common for all classes and a normal dist**

$$\delta_k(x) = x \frac{\mu_k}{\sigma^2} - \frac{\mu_k^2}{2\sigma^2} \quad (3.13)$$

Calculating the above may not be possible as the population is not available so estimates of mean and std are used:

$$\hat{\mu}_k = \frac{1}{n_k} \sum_{i:y_i=k} x_i \quad (3.14)$$

$$\hat{\sigma}^2 = \frac{1}{n - K} \sum_{k=1}^K \sum_{i:y_i=k} (x_i - \hat{\mu}_k)^2 \quad (3.15)$$

This being the average sample var for each class (on the premise they are common).

These are then plugged in to give the boundary points as:

$$\delta_k(\hat{x}) = \hat{x} \frac{\hat{\mu}_k}{\hat{\sigma}^2} - \frac{\hat{\mu}_k^2}{2\hat{\sigma}^2} \quad (3.16)$$

If possible compare to Bayes error rate to determine classifier performance.

3.9.2 What if the number of predictors > 1?

Multi-variate normal dist is then used. From which a class specific mean vector is used and a common covariance matrix.

Each predictor follows a one dimensional normal dist with some correlation between predictors.

Think of one dimensional as collapsing into the left side or down. If the var is equal and correlation=0 then the image on the left in 3.1 with a circular bottom is produced, if otherwise skewed, as shown on the right.

$$X \sim N(\mu, \Sigma); \mu: \text{mean vector}, \Sigma: \text{covariance matrix} \quad (3.17)$$

3.9.3 How is multi predictor LDA used?

Using the multi variate normal dist values can create the vector/matrix version of the Bayes classifier boundary. If there are >2 classes the boundary will be shown as pair of classes.

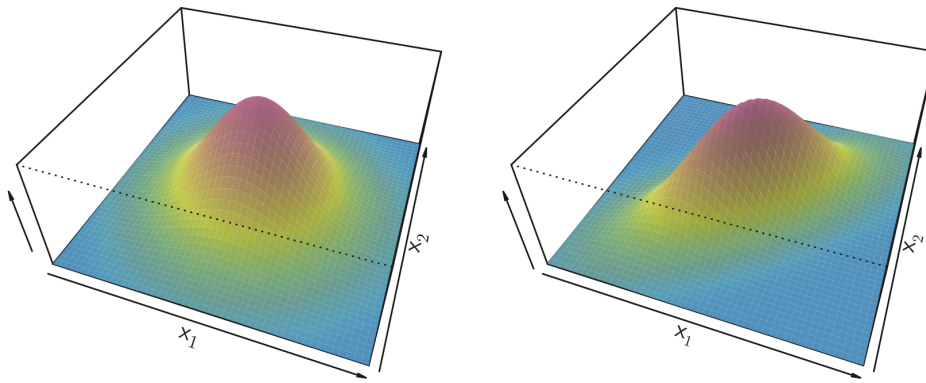


FIGURE 4.5. Two multivariate Gaussian density functions are shown, with $p = 2$. Left: The two predictors are uncorrelated. Right: The two variables have a correlation of 0.7.

Figure 3.1: Multivariate Normal Dist

3.9.4 Issues to be aware of from multi LDA

Small distribution between classes in the data set can result in good training error rates, but poor test error rates.

3.10 What is a confusion matrix?

Determine which type of error is being made in terms of the classes, comparing the error rate for each class.

Sensitivity: Sensitivity is the percentage of true positives (e.g. 90% sensitivity = 90% of people who have the target disease will test positive)

Specificity: Specificity is the percentage of true negatives (e.g. 90% specificity = 90% of people who do not have the target disease will test negative).

3.11 How to accomdate non equal classes?

If the bias between classes is not equal the boundary can be altered to instead of being the max, match the criteria of the problem.

3.12 What is ROC curve?

Shows the error rate using different thresholds, the overall peformance of the model is the area under the curve.

If area under the curve is ≤ 0.5 its assumed the performance is no better than chance.

3.13 Different classification measurments

3.14 Quadratic Discriminant Analysis

The main differnce to LDA is that there is not an assumption of common variance between distribution, this leads to a quadratic boundary function.

Name	Definition	Synonyms
False Pos. rate	FP/N	Type I error, 1–Specificity
True Pos. rate	TP/P	1–Type II error, power, sensitivity, recall
Pos. Pred. value	TP/P*	Precision, 1–false discovery proportion
Neg. Pred. value	TN/N*	

Figure 3.2: Types of classification measurements

3.14.1 What are the pros/cons of QDA?

More flexible, accommodates differing class distributions and therefore has lower bias.

It is more compute intensive and has higher variance.

3.14.2 How do you choose between LDA and QDA?

Compare the correlation between class and observations for each class, if they are all similar shows the variance is likely similar. And therefore LDA is a good choice otherwise QDA is.

3.15 How do you choose between LDA and Logistic for binary scenarios?

Main difference is that LDA assumes normal distribution of class, logistic does not.

3.16 Differences with KNN?

Decision boundaries can be weighted towards classes, however it is the most flexible choice.

3.17 Methods to improve non linear performance?

Can improve non linear performance by adding transformed predictors, ie X_1^2, X_2^2 etc.

4 Resampling methods

Completely separate data as the test set may not make sense to use as that could be used as a valuable method to train the original data set. So instead you can repeatedly draw samples from the data set you have and refit the model and compare to other samples as test data to gain info about the model's setup.

Model assessment: Evaluating a model's performance.

Model selection: Determining a model's flexibility.

Side note: Can determine if a transformation is valid using hypothesis tests.

4.1 Cross Validation

4.1.1 What is the validation set approach?

Randomly divide the data into a training and test set. The issue is there will be less data used for training the model with each split, will most likely overestimate the test error.

4.1.2 What is leave one out cross validation (LOOCV)?

Use a single observation as the validation set, the rest as training and loop through all data points.

Use the average MSE from all models and validations point:

$$CV_n = \frac{1}{n} \sum MSE \quad (4.1)$$

The benefits to this are:

1. less bias in model as the training set is much bigger
2. No randomness in method so consistent results produced

This is a very flexible method, however, does require large computation.

4.1.3 What is K fold cross validation?

Separate into K sets, 1 set is the validation and the rest are used for training. Loop through switching the validation set to each Kth set, taking the average across n sets as the error. The main benefit is the computational costs.

Generally the result produced for using K as a reasonable number (>5), produces a result similar to LOOCV (K=n).

The mean of many highly correlated samples has a higher variance than if uncorrelated. And therefore kCV will produce a low variance than LOOCV but higher bias (less training data used).

A good var-bias balance is found around k=5 and k=10, and cleaner to use a multiple of the n as k. K fold is the main method to compare different models looking for the the bottom of the error over the different models.

4.2 Bootstrap

Bootstrapping involves altering the original data set to make estimates on the accuracy of the data. An example of bootstrapping is randomly duplicating and replacing observations throughout the data set to get new data to train and validate on. Other types include "Smooth sampling" zero centered (mean) adding noise to normal dist data.

This can be applied to a wide range of statistical methods, especially useful for models whose variance is difficult to measure.

5 Regularisation

If the number of predictors is greater than the number of observations, least square method cannot be used, it will produce incorrect results.

Irrelevant predictors will lead to unwanted complexity in the model and wasteful computation, the following are methods to reduce the number of predictors.

Methods below normally assume linear regression is being used and therefore RSS is being minimised.

1. Subset: Choose a subset of predictors.
2. Shrinkage (Normalisation): Set some parameters to 0.
3. Dimension reduction: Project p onto M space.

5.1 How do you use subsets?

Compare all possible combination of predictors (2^p). Choose the best in terms of RSS or R^2 for each set number of predictors. Choose the best overall with AIC, BIC or adjusted R^2 .

5.1.1 How does it work with logistic regression?

Use deviance instead of RSS, which is the $-2 \times$ maximum log likelihood.

$$D_M = -2 \log(L_M) \quad (5.1)$$

L_M : Max achievable likelihood.

5.2 How do you perform step wise subset ?

Add one predictor on every iteration, each time trying all possible predictor combination for that set number of predictors. Then based on the previous best add another predictor.

When comparing two models with the same number of predictors can use R^2 or RSS, however differing predictor number AIC, BIC or adjusted R^2 should be used.

5.2.1 What is the pros/cons to this method?

Reduce the iterations to $p(1+p)/2$. The con is possibly missing different combinations that together produce a better model than working in the best approach method.

To combat this backward stepwise can be used as well as some hybrid methods that allow better scoping of possible combinations.

5.3 How do you compare test error for different models?

Two methods: Adjust training data to remove bias or use cross validation.

How to determine the error parameter to compare:

- C_p
- Akaike Information Criterion (AIC)
- Bayesian Information Criterion (BIC)
- R^2

$$C_p = \frac{1}{n} (RSS + 2d\sigma^2); d : \text{predictors} \quad (5.2)$$

This adds a penalty of $2d$ var to the RSS error. For least squares AIC is proportional to C_p .

$$AIC = \frac{1}{n\sigma^2} (RSS + 2d\sigma^2) \quad (5.3)$$

$$BIC = \frac{1}{n} (RSS + \log(n)d\sigma^2) \quad (5.4)$$

BIC provides a larger penalty for more predictors as when $n > 7 \log(n) > 2$.

$$\text{Adjusted } R^2 = 1 - \frac{RSS/(n-d-1)}{TSS/(n-1)} \quad (5.5)$$

$$TSS = \sum (y_i - \bar{y})^2 \quad (5.6)$$

Balances RSS by increasing the denominator in terms of the number of predictors, ensuring the number of predictors provides enough benefit in RSS to make it worth it.

5.4 How does cross validation compare?

Cross validation does not assume any model form and so is preferred in comparing models.

5.5 What is the one standard error rule?

Choose the lowest predictor error score in a cross validation which is within one standard error (of the MSE for each model) of the lowest CV error.

5.6 What is Shrinkage?

Reduce the parameters of the model, two methods: Ridge regression and LASSO.

5.6.1 What is ridge regression

Adds a penalty to the standard RSS, to push during the minimisation the parameters towards 0.

$$RSS + \lambda \sum_{j=1}^p \beta_j^2 \quad (5.7)$$

$\lambda = 0$: no effect

$\lambda \rightarrow \infty$: Impact grows

$\beta_0 = \bar{y}, x_i = 0$

5.6.2 How is ridge regression effected by scaling?