

# Multivariate probability

Adam Boulton ([www.bou.lt](http://www.bou.lt))

February 10, 2023

# Contents

Preface	2
<b>I Exploratory data analysis</b>	<b>3</b>
1 Association rules	4
2 Summary statistics and visualisation for multiple variables	7
3 Privacy	10
4 Distance metrics and outliers	11
<b>II Estimating discriminative probability distributions</b>	<b>13</b>
5 Bayesian parameter estimation of discriminative models	14
6 Point variable estimates for discriminative models	17
7 Using F-tests to compare regression models	22
8 Test sets and validation sets	23
9 Choosing parametric discriminative probability distributions	25
<b>III Unsupervised machine learning</b>	<b>28</b>
10 Dimensionality reduction with Principal Component Analysis (PCA)	29
11 K-means and k-medoids clustering	31

<i>CONTENTS</i>	2
<b>IV Estimating generative probability distributions</b>	<b>34</b>
12 M-estimators	35
13 The Generalised Method of Moments (GMM)	37
<b>V Supervised linear regression</b>	<b>40</b>
14 Ordinary Least Squares for prediction	41
15 Regularising linear regression for prediction	47
16 Choosing linear models for prediction	51
17 Generalised linear models, the delta rule and binary classification	52
18 Generalised linear models and multiclass classification	57
<b>VI Supervised machine learning</b>	<b>62</b>
19 Classification trees	63
20 Regression trees	66
21 Bayesian trees	68
22 Support Vector Machines (SVMs)	70
23 Variational Bayes	72
24 The Naive Bayes classifier	73
25 The K-Nearest Neighbours (KNN) classifier	75
26 Discriminant analysis	76
27 Non-parametric regression	77
28 Ensemble methods	80
29 Ensemble methods for trees	84
30 Regularising black box models	85
31 Confidence intervals of black box models	86

<i>CONTENTS</i>	3
<b>32 Interpreting black box models</b>	<b>87</b>
<b>33 Semi-supervised learning</b>	<b>89</b>
<b>34 Imputing missing data for prediction</b>	<b>90</b>
<b>35 Recommenders</b>	<b>91</b>
<b>VII Linear regression for inference</b>	<b>93</b>
<b>36 Ordinary Least Squares for inference</b>	<b>94</b>
<b>37 Testing regression parameter estimates with Z-tests and T-tests</b>	<b>99</b>
<b>38 Multiple hypothesis testing</b>	<b>100</b>
<b>39 Generalised Least Squares</b>	<b>102</b>
<b>40 General Linear Models</b>	<b>104</b>
<b>VIII Advanced inference</b>	<b>108</b>
<b>41 Analysis of variance (ANOVA)</b>	<b>109</b>
<b>42 Instrumental Variables</b>	<b>110</b>
<b>43 Imputing missing data for inference</b>	<b>118</b>
<b>44 Measurement error and inference</b>	<b>119</b>
<b>45 Semi-parametric regression</b>	<b>120</b>

# Preface

This is a live document, and is full of gaps, mistakes, typos etc.

## Part I

# Exploratory data analysis

# Chapter 1

## Association rules

### 1.1 Association rules

#### 1.1.1 Association rules

##### The data

We have a transaction dataset,  $D$ .

This includes transactions of items in  $I$ .

Any subset of  $I$  is an itemset.

A subset of size  $k$  is a  $k$ -itemset.

Transactions are a  $k$ -itemset with a unique id, tid.

The set of all transactions is  $T$ .

A tidset is a subset of  $T$ .

##### Forming a lattice

We have a total order on the items.

An itemset  $ab$  is greater than  $a$ , for example.

The two points of the lattice are the nullset, and  $I$ .

##### Mappings

We have a mapping from  $I$  to  $T$  called  $t$ .

We have another mapping from  $T$  to  $I$  called  $i$ .

**Frequency**

We define the frequency of an itemset as the number of transactions it appears in.

We can write the frequency of  $A$  as  $\sum A$ .

**1.1.2 Strong rules****Strong rules**

We use frequent patterns to generate strong rules,  $R$ .

An example of a strong rule is  $a \rightarrow b$ .

We can look at this by comparing the support of  $a$  to  $a \wedge b$ .

$$\text{supp}(A \rightarrow B) = P(A \wedge B)$$

$$\text{conf}(A \rightarrow B) = P(B|A)$$

$$\text{conf}(A \rightarrow B) = \frac{P(A \wedge B)}{P(A)}$$

**1.1.3 Support****Support**

We define the support of an itemset as the proportion of transactions which contain the itemset.

$$\text{supp}(A) = \frac{\sum A}{n}.$$

We can also consider this as:

$$\text{supp}(A) = P(A)$$

**1.1.4 Frequent patterns****Frequent patterns**

A frequent itemset is one where the support is above a minimum.

We know that if an itemset is frequent, then all its subsets are also frequent.

We look for frequent patterns,  $F$ , between the items  $I$ .

An example of a frequent pattern is  $\{a, b\}$ .

**1.1.5 Confidence****Confidence**

The confidence of a frequent pattern is defined as:



$$\text{conf}(A \rightarrow B) = \frac{\text{supp}(A \wedge B)}{\text{supp}(A)}$$

### 1.1.6 Finding strong rules using the Apriori algorithm

#### Finding frequent patterns

We can use search algorithms to find frequent patterns. Starting at the empty set.

#### Apriori algorithm

Breadth first search to generate candidate set of itemsets with support above some value.

Start with a 1-itemset, and increase  $k$  once done.

Once we have found a frequent pattern, we can immediately identify other frequent patterns associated with it.

We can do this by looking at confidence, not support.

### 1.1.7 Interest

#### Interest

An alternative measure for finding rules is to use interest.

$$\text{Interest}(A \rightarrow B) = \frac{P(A \wedge B)}{P(A)P(B)}$$

$$\text{Interest}(A \rightarrow B) = \frac{\text{supp}(A \wedge B)}{P(\text{supp}(A))P(\text{supp}(B))}$$

If this is 1, then they are independent.

If this is greater than 1, they are positively dependent.

If this is less than 1, they are negatively dependent.

### 1.1.8 Quantitative association rules

#### Quantitative association rules

The search space is infinite in size. For example continuous age.

We choose intervals instead.

## Chapter 2

# Summary statistics and visualisation for multiple variables

### 2.1 Statistics for two variables

#### 2.1.1 Sample covariance

We previously defined the population covariance as  $\sigma_{XY} = E[(X - \mu_X)^T(Y - \mu_Y)]$ .

We define the sample covariance as  $\sigma_{XY} = \frac{1}{n} \sum_i (x_i - \bar{x})(y_i - \bar{y})$ .

We can calculate this using matrices:

$$M = X - \bar{x}$$

$$N = Y - \bar{y}$$

$$\sigma_{XY} = \frac{1}{n} M^T N.$$

#### 2.1.2 Sample correlation

$$\rho_{XY} = \frac{\sigma_{XY}}{\sigma_X \sigma_Y}$$

#### 2.1.3 Covariance matrix

If we have  $n$  variables we can have a  $n \times n$  matrix  $\Sigma$  where:

$$\Sigma_{ij} = \sigma_{ij} = \frac{1}{n} (X_i - \bar{x}_i)^T (X_j - \bar{x}_j)$$

### 2.1.4 Centred covariance

If  $\bar{x} = \bar{y} = 0$  then:

$$\sigma_{XY} = \frac{1}{n} X^T Y$$

### 2.1.5 Correlation matrix

Here each entry is the correlation rather than the covariance.

## 2.2 Correlation coefficients

### 2.2.1 Pearson correlation coefficient

The Pearson correlation coefficient is defined as the covariance normalised by the individual variances.

It is between  $-1$  (total negative linear correlation),  $0$  (no linear correlation) and  $1$  (total positive linear correlation).

$$\rho_{X,Y} = \frac{\text{cov}(X, Y)}{\sigma_X \sigma_Y}$$

### 2.2.2 Spearman rank correlation

For each of 2 variables we create a ranking of them.

From  $X$  and  $Y$  we then have  $R_X$  and  $R_Y$ .

We then calculate the Pearson correlation coefficient between the rankings.

$$r_S = \frac{\text{cov}(R_X, R_Y)}{\sigma_{R_X} \sigma_{R_Y}}$$

### 2.2.3 Kendall rank correlation

Concordant and discordant pairs

$$\tau = \frac{n_{\text{concordant}} - n_{\text{discordant}}}{\binom{n}{2}}$$

### 2.2.4 General correlation coefficient

## 2.3 Updating statistics

### 2.3.1 Updating the covariance

If it is centred:

$$\sigma_{XY}^n = \frac{1}{n} X_n^T Y_n$$

So:

$$\sigma_{XY}^{n+1} = \frac{n\sigma_{XY}^n + x_{n+1}^t y_{n+1}}{n+1}$$

## 2.4 Visualising multiple continuous variables

### 2.4.1 Time series

### 2.4.2 Scatter plots (with size as variable)

### 2.4.3 Q-Q plots

Plot quartiles of variables against each other.

## 2.5 Visualising a single class variable

### 2.5.1 Bar and column charts

### 2.5.2 Pie charts

## 2.6 Visualising multiple class variables

### 2.6.1 Stacked bar and column charts

## 2.7 Visualising class and continuous variables

### 2.7.1 Multiple box and whiskers

### 2.7.2 Scatter plots with colour

## 2.8 Visualising geographic data

## 2.9 Visualising time series

### 2.9.1 Heat maps

### 2.9.2 Sparklines

## Chapter 3

# Privacy

### 3.1 Differential privacy

## Chapter 4

# Distance metrics and outliers

### 4.1 Measuring distance between vectors

#### 4.1.1 $L_p$ norms

$L_p$  norms can be used to measure the distance between two metrics.

If we have data points  $v$  and  $w$  the distance is:

$$\|v - w\| = (\sum_i |v_i - w_i|^p)^{\frac{1}{p}}$$

If  $p = 2$  we have the Euclidian norm. If  $p = 1$  we have the Manhattan norm.

#### 4.1.2 Dot product

Given two vectors we can calculate:

$$\frac{a \cdot b}{\|a\| \|b\|}$$

If the two vectors are identical, this is 1. If they are orthogonal this is 0. If they are opposite, this is  $-1$ .

#### 4.1.3 Kernels

This is a generalisation of the dot product function, where we want to find similarity between two vectors.

If we have data points  $v$  and  $w$  the distance is:

$$K(v, w)$$

#### 4.1.4 Mahalanobis distance

We have a point. How far away is this from the mean.

For a single dimension: number of standard deviations.

What about multidimensional data?

Could do sd for all distances, but correlations between variables. If two variables are highly correlated, it's not really twice as far.

We use this:

$$D_M(\mathbf{x}) = \sqrt{(\mathbf{x} - \boldsymbol{\mu})^T S^{-1} (\mathbf{x} - \boldsymbol{\mu})}$$

## 4.2 Measuring distance between matrices

### 4.2.1 Frobenius norm

If we have matrices  $A$  and  $B$  the distance is:

$$\|A - B\| = \sqrt{\sum_i \sum_j |a_{ij} - b_{ij}|^2}$$

This is a Euclidian norm.

## 4.3 Measuring distance between time series

### 4.3.1 Dynamic time warping

We may want to examine the similarity between two sequences.

We want to match a sample from one sequence to a sample from the other sequence.

Simply matching at the same time point is naive, as samples may move at different speeds, or have offsets.

## 4.4 Finding neighbours

### 4.4.1 Nearest Neighbour Search (NSS)

#### 4.4.2 Finding neighbours

Say we have a distance function and a sample. How can we identify the  $k$ -nearest neighbours?

We can find the distance for all points, sort this and take the top  $k$  observations.

### 4.4.3 k-Nearest Neighbour Search

## Part II

# Estimating discriminative probability distributions



## Chapter 5

# Bayesian parameter estimation of discriminative models

### 5.1 Introduction

#### 5.1.1 Generative and discriminative models

##### Recap

For parametric models without dependent variables we have a form:

$$P(y|\theta)$$

And we have various ways of estimating  $\theta$ .

We can write this as a likelihood function:

$$L(\theta; y) = P(y|\theta)$$

##### Discriminative models

In discriminative models we learn:

$$P(y|X, \theta)$$

Which we can write as a likelihood function:

$$L(\theta; y, X) = P(y|X, \theta)$$

##### Generative models

In generative models we learn:

$$P(y, X|\theta)$$

Which we can write as a likelihood function:

$$L(\theta; y, X) = P(y, X|\theta)$$

We can use the generative model to calculate dependent probabilities.

$$P(y|X, \theta) = \frac{P(y, X|\theta)P(\theta)}{P(X, \theta)}$$

$$P(y|X, \theta) = \frac{P(y, X|\theta)}{P(X|\theta)}$$

## 5.2 Bayesian parameter estimation

### 5.2.1 Bayesian parameter estimation for dependent models

#### Recap

For non-dependent models we had:

$$P(\theta|y) = \frac{P(y, \theta)}{P(y)}$$

$$P(\theta|y) = \frac{P(y|\theta)P(\theta)}{P(y)}$$

The bottom bit is a normalisation factor, and so we can use:

$$P(\theta|y) \propto P(y|\theta)P(\theta)$$

We have here:

- Our prior -  $P(\theta)$
- Our posterior -  $P(\theta|y)$
- Our likelihood function -  $P(y|\theta)$

#### Bayesian regression for generative models

We know:

$$P(\theta|y, X) = \frac{P(y, \theta, X)}{P(y, X)}$$

$$P(\theta|y, X) = \frac{P(y, X|\theta)P(\theta)}{P(y, X)}$$

The bottom bit is a normalisation factor, and so we can use:

$$P(\theta|y, X) \propto P(y, X|\theta)P(\theta)$$

We have here:

- Our prior -  $P(\theta)$
- Our posterior -  $P(\theta|y, X)$
- Our likelihood function -  $P(y, X|\theta)$

### Bayesian regression for discriminative models

We know:

$$P(\theta|y, X) = \frac{P(y, \theta, X)}{P(y, X)}$$

$$P(\theta|y, X) = \frac{P(y|\theta, X)P(\theta, X)}{P(y, X)}$$

$$P(\theta|y, X) = \frac{P(y|\theta, X)P(\theta)P(X|\theta)}{P(y, X)}$$

We assume  $P(X|\theta) = X$  and so:

$$P(\theta|y, X) = \frac{P(y|\theta, X)P(\theta)P(X)}{P(y, X)}$$

The bottom bit is a normalisation factor, and so we can use:

$$P(\theta|y, X) \propto P(y|X, \theta)P(\theta)$$

We have here:

- Our prior -  $P(\theta)$
- Our posterior -  $P(\theta|y, X)$
- Our likelihood function -  $P(y|X, \theta)$

### 5.2.2 Prior and posterior predictive distributions for dependent variables

#### Prior predictive distribution

Our prior predictive distribution for  $P(y|X)$  depends on our prior for  $\theta$ .

$$P(y|X) = \int_{\Theta} P(y|X, \theta)P(\theta)d\theta$$

#### Posterior predictive distribution

Once we have calculated  $P(\theta|\mathbf{y}, \mathbf{X})$ , we can calculate a posterior probability distribution for  $P(y|X)$ .

$$P(y|\mathbf{x}, \mathbf{y}, \mathbf{X}) = \int_{\Theta} P(y|\mathbf{x}, \theta)P(\theta|\mathbf{y}, \mathbf{X})d\theta$$

## Chapter 6

# Point variable estimates for discriminative models

### 6.1 Predictions and residuals

#### 6.1.1 Predictions

Our data  $(\mathbf{y}, \mathbf{X})$  is divided into  $(y_i, \mathbf{x}_i)$ .

We create a function  $\hat{y}_i = f(\mathbf{x}_i)$ .

The best predictor of  $y$  given  $x$  is:

$$g(X) = E[Y|X]$$

The goal of regression is to find an approximation of this function.

#### 6.1.2 Residuals

$$\epsilon_i = y_i - \hat{y}_i$$

#### 6.1.3 Residual sum of squares (RSS)

$$RSS = \sum_i \epsilon_i^2$$

$$RSS = \sum_i (y_i - \hat{y}_i)^2$$

#### 6.1.4 Explained sum of squares (ESS)

$$ESS = \sum_i (\bar{y} - \hat{y}_i)^2$$

**6.1.5 Total sum of squares (TSS)**

$$TSS = \sum_i (y_i - \bar{y})^2$$

**6.1.6 Relationship between prediction and probability distribution**

$$P(y|X, \theta)$$

$$\hat{y} = f(\mathbf{x})$$

Through integration?

$$E[y] = \int P(y|X, \theta) dy$$

**6.1.7 Coefficient of determination ( $R^2$ )**

$$R^2 = 1 - \frac{RSS}{TSS}$$

**6.2 Classification****6.2.1 Binary classification**

Classification models are a type of regression model, where  $y$  is discrete rather than continuous.

So we want to find a mapping from a vector  $X$  to probabilities across discrete  $y$  values.

A classifier takes  $X$  and returns a vector.

For a classifier we have  $K$  classes.

**6.2.2 Classification**

Confusion matrix. true positive, false positive, false negative, true negative

Can use this to get

Accuracy: percentage correct

Precision: percentage of positive predictions which are correct

Recall (sensitivity): percentage of positive cases that were predicted as positive

Specificity: percentage of negative cases predicted as negative

**6.2.3 Multiclass classification**

Multiclass classification

What if can be email for work, friends, family, hobby?

### 6.2.4 Confusion matrix

Include error types here

### 6.2.5 Hard and soft classifiers

A hard classifier can return a sparse vector with 1 in the relevant classification.

A soft classifier returns probabilities for each entry in the vector.

The vector represents  $P(Y = k|X = x)$

### 6.2.6 Transforming soft classifiers into hard classifiers

We can use a cutoff.

If there are more than two classes we can choose the one with the highest score.

## 6.3 Loss functions for point predictions

### 6.3.1 Minimum Mean Square Error (MMSE)

Mean estimate.

Can do for a parameter, or for a predicted estimate for  $y$ .

Linear models

MLE is same as  $y^2$  loss

MAP is same as  $y^2$  loss with regularisation

### 6.3.2 Loss functions for soft classifiers

**Hinge loss**

**Brier score**

### 6.3.3 Loss functions for hard classifiers

Don't want answers outside 0 and 1.

**F score**

**F1 score**

$$F_1 \text{ score: } \frac{2PR}{(P + R)}$$

may not just care about accuracy, eg breast cancer screening

high accuracy can result from v basic model (ie all died on titanic)

Receiver Operating Characteristic (ROC) Area Under Curve (AUC)

## 6.4 Other

### 6.4.1 Estimating other priors

Estimating  $P(k|T)$  - Which variables we split by, given the tree size

Estimating  $P(r|T, k)$  - The cutoff, given the tree size and the variables we are splitting by

Estimating  $P(\theta|T, k, r)$

### 6.4.2 Maximum Likelihood Estimation (MLE) for generative and discriminative models

### 6.4.3 Maximum A-Priori estimation (MAP) for generative models

Bayesian regression for generative models

We know:

$$P(\theta|y, X) = \frac{P(y, \theta, X)}{P(y, X)}$$

$$P(\theta|y, X) = \frac{P(y, X|\theta)P(\theta)}{P(y, X)}$$

The bottom bit is a normalisation factor, and so we can use:

$$P(\theta|y, X) \propto P(y, X|\theta)P(\theta)$$

We have here:

- Our prior -  $P(\theta)$
- Our posterior -  $P(\theta|y, X)$
- Our likelihood function -  $P(y, X|\theta)$

### 6.4.4 Bayesian classifier

**Classification risk**

We can measure the risk of a classifier. This is the chance of misclassification.

$$R(C) = P(C(X) \neq Y)$$

**The Bayesian classifier**

This is the classifier  $C(X)$  which minimises the chance of misclassification.

*CHAPTER 6. POINT VARIABLE ESTIMATES FOR DISCRIMINATIVE MODELS*23

It takes the output of the soft classifier and chooses the one with the highest chance.



## Chapter 7

# Using F-tests to compare regression models

### 7.1 Hypothesis testing

#### 7.1.1 Power of tests

#### 7.1.2 Type I and type II errors

#### 7.1.3 Sensitivity tests

### 7.2 F-test

#### 7.2.1 F test for equal population means

#### 7.2.2 F test for additional variables

### 7.3 Other

#### 7.3.1 Cohen's d

# Chapter 8

## Test sets and validation sets

### 8.1 Test sets

#### 8.1.1 Overfitting

Overfitting is a risk. Instead we split to test, train. risk of using too many features. more features always improve training score, not necessarily test score.

As model gets more complex, both test and train do better. however at some point, test stops doing better, overfitting

Structural risk minimisation can address this trade off. use test and training sets. train model on train, rate it on test

Structural minimisation curve has accuracy of boths sets over complexity

To avoid overfitting:

+ reduce number of features + do a model selection + use regularisation + do cross validation

Can choose other model parameters

How to evalute model?

#### 8.1.2 K-fold cross-validation

Can do k-fold cross validation. given algo A and dataset D, divide D into k equal sized subsets

For each subset, train the model on all other subsets and test on the other subset. average error between folds

## 8.2 Splitting

### 8.2.1 K-folds

The problem of different sample sizes (sample size for validation sets is lower, different hyper parameters could be more appropriate)

### 8.2.2 Leave-One-Out

## 8.3 Hyperparameters

### 8.3.1 Learning rate, batching and momentum

## 8.4 Search methods

### 8.4.1 Grid-search

## 8.5 Validation sets

### 8.5.1 Validation sets

which features? remove, add?

change lambda, regularisation

change polynomial features

## Chapter 9

# Choosing parametric discriminative probability distributions

### 9.1 Sample sizes

#### 9.1.1 Sample sizes

If you're modelling house prices using just size, getting a large sample size won't help too much

Can improve low bias models

Is data size an issue? can artificially restrict training data size and then evaluate error

Training:

+ zero error for low  $m$  + increases error as  $m$  increases, as degrees of freedom/ $m$  falls

cv:

+ error decreases as data set increases, more accurate  $\theta$

The two curves converge towards each other for  $v$  large  $m$

When are large datasets useful?

When all features available:

Predicting house price using just size, won't benefit from more data...

If choosing correct word in sentence (to, too, two), more helpful

If human expert can do it, then more data probably helpful

Expert realtor probably couldn't do much with just size, but speaker could answer other q

Could expert do it?

Low bias algorithms do well with more data

More data good if large number of parameters, or lot of hidden units.

## 9.2 Overfitting

### 9.2.1 Overfitting

Role of lambda: high makes impact of more variables lower  $\Rightarrow$  high bias

Low makes impacts of more variables strong  $\Rightarrow$  high variance

Can trade off using cut off. only make positive if above 0.7

How to use? difficult, as lambda within cost!

Can do similarly to d:

Run for a range of lambda (eg 0, 0.01, 0.02, 0.04, 0.08: 10), then pick from cross validation set

Low lambda always has low cost for training set, but not for cv set..

Regularisation: add to error term the size of the term. penalised large parameters

May not fit outside sample

High bias: eg house prices and size. linear would have high bias for out of scope sample (underfitting)

High variance: making polynomial passing through all data (overfitting)

Can reduce overfitting by reducing features either manually or using models

OR regularisation: keep all features, but reduce magnitude of theta

### 9.2.2 Regularisation

Make cost function include size of  $\theta^2$  values

$$\min \frac{1}{2m} [\sum (h(x) - y)^2 + 1000\theta_3^2 + 1000\theta_4^2]$$

or more broadly:

$$\min \frac{1}{2m} [\sum \dots + \lambda \sum \theta_j^2]$$

Tend to not include  $\theta_0$  as convention, no regularisation

Update for linear regression is

$$\theta_j = \theta_j - \alpha \left( \frac{1}{m} \right) * \text{sum}(h(x) - y)x_j + (\lambda/m)\theta_j$$

$$\theta_j = \theta_j(1 - \alpha\lambda/m) - \text{alpha}(1/m) * \sum(h(x) - y)x_j$$

This is the same as before, but  $\theta_j$  updates from a smaller  $\theta_j$  each time.

Normal equation needs a change

$$(X'X)^{-1}X'y = \theta''$$

Now is

$$(X'X + \lambda I)^{-1}X'y'$$

although for  $\theta_0$ ,  $\lambda$  zero, so identity matrix, but first element 0

REGULARISATION FOR REGULARISATION

add to end of  $J(\theta)$ :

$$\frac{\lambda}{2m} \sum \theta_j^2$$

update for  $\theta_j$   $j > 0$ : is as linear regression, but  $h(x)$  is a different function

## Part III

# Unsupervised machine learning

## Chapter 10

# Dimensionality reduction with Principal Component Analysis (PCA)

### 10.1 Dimensionality reduction

#### 10.1.1 Classical principal component analysis

##### Introduction

Principal component analysis takes a dataset  $X$  with  $m$  variables and returns a principal component matrix  $A$  with size  $m \times k$ .

Each new dimension is a linear function of the existing data.  $Z = XA$ .

Each dimension is uncorrelated, and ordered, in order of descending explanation of variability.

The problem of principal component analysis is to find these weightings  $A$ .

##### Classical PCA

We take the first  $k$  eigenvectors of the covariance matrix, ordered by eigenvalue.

##### Getting the eigenvectors using SVD

We can decompose  $X = U\Sigma A^T$ .

We can take the eigenvectors from  $A$ .



### Choosing the number of dimension

We can choose  $k$  such that a certain percentage of the variance is retained.

### 10.1.2 Robust principal component analysis

#### Robust PCA

Robust PCA can be used to deal with corrupted data, such as corrupted image data.

Rather than data  $X$  we have  $M = L_0 + S_0$  where  $L_0$  is what we want to recover (and is low rank), and  $S_0$  is noise (and sparse).

In video footage,  $L_0$  can correspond to the background, while  $S_0$  corresponds to movement.

# Chapter 11

## K-means and k-medoids clustering

### 11.1 Clustering

#### 11.1.1 Evaluating clusterings

##### Davies-Bouldin index

The Davies-Bouldin index is a method for evaluating clustering algorithms, such as k-means.

It examines the distance between centroids, and the tightness of centroids.

#### 11.1.2 k-means clustering

##### Introduction

K-means clustering is the most widely used unsupervised model.

In k-means clustering we identify  $k$  centroids in the feature space. We then calculate the distance from each data point to each of the centroids, and allocate the data point to the nearest centroid.

This requires a method for calculating the location of the centroids.

##### Identifying the centroids

We apply an iterative approach to identifying the centroids.

We first initialise by assigning centroids randomly to existing data points.

We then iteratively perform the following:

- Calculate the distances between each data point and each centroid.
- Assign each data point to the closed centroid.
- Update each centroid location to the mean of the data points allocated to it.

### Calculating distances

This method requires us to calculate the distance between two points in the feature space.

For k-means we use the Euclidian distance.

### Potential issues

It is possible for a centroid to have no data assigned to it. If this happens we can eliminate the cluster, or reassign some data points.

The algorithm may only arrive at a local minima. In order to maximise the chance of an effective clustering, we can do k-means under different initialisations of the centroids in order to minimise risk of bad local optima.

### Choosing $k$

If the points in each cluster follow a normal distribution, that's a good sign. This can be tested with Anderson-Darling.

If it's not normal, we can split the cluster into 2.

### Using clustering as part of data analysis

We can choose  $k$  if output is being used in later data analysis (eg type assignment, complaint level or something)

## 11.1.3 k-medoids

### Introduction

k-medoids is similar to k-means clustering, with two key differences:

- Centroids are now always located on data points, rather than floating freely.
- We minimise  $l_1$  distance, rather than  $l_2$ .

### Partitioning Around Medoids (PAM) algorithm

This is the most common approach for k-medoids.

We initialise randomly, as we do for k-means.

We then iterate the following:

- Calculate the loss for the current allocation
- For each medoid, see if swapping allocation with another (non-medoid) data point decreases the cost.
- If it does, make the swap.

## Part IV

# Estimating generative probability distributions

# Chapter 12

## M-estimators

### 12.1 M-estimators

#### 12.1.1 Introduction

page setting out linear stuff to come

OLS, generalised linear models etc are m-estimators, as are gmm

h3 on parametric

With maximum likelihood estimation we maximise a function.

We could choose other functions to maximise or minimise.

$$\sum_i f(x_i, \theta)$$

If  $f(x_i, \theta)$  is differentiable wrt to  $\theta$  this can be solved by finding the stationary point.

This is a  $\phi$  type.

Otherwise it is a  $\rho$  type.

page on influence functions there

Generalisation of MLE.

$$m_\theta = m_\theta(x, \theta)$$

Z-estimator is where this is met, through diff

$$\frac{\delta}{\delta \theta} m_\theta = z_\theta(\theta, x) = 0$$

M-estimator for mean

$$m_\theta(\theta) = -(x - \theta)^2$$

$$z_{\theta}(\theta) = x - \theta$$

## Chapter 13

# The Generalised Method of Moments (GMM)

### 13.1 Generalised Method of Moments (GMM)

#### 13.1.1 Difference from Method Of Moments (MOM)

More conditions than data.

#### 13.1.2 Generalised Method of Moments (GMM)

We have a function on the output and a parameter:

$$g(y, \theta)$$

A moment condition is that the expectation of such a function is 0.

$$m(\theta) = E[g(y, \theta)] = 0$$

To do GMM, we estimate this using:

$$\hat{m}(\theta) = \frac{1}{n} \sum_i g(y_i, \theta)$$

We define:

$$\Omega = E[g(y, \theta)g(y, \theta)^T]$$

$$G = E[\Delta_\theta g(y, \theta)]$$

And then minimise the norm:

$$\|\hat{m}(\theta)\|_W^2 = \hat{m}(\theta)^T W \hat{m}(\theta)$$

Where  $W$  is a positive definite matrix for the norm.



$\Omega^{-1}$  is most efficient. But we don't know this. It depends on  $\theta$ .

We can estimate it if IID:

$$\hat{W}(\hat{\theta}) = \left( \frac{1}{n} \sum_i g(y, \hat{\theta}) g(y, \hat{\theta})^T \right)^{-1}$$

### 13.1.3 Two-step feasible GMM

Estimate using  $\mathbf{W} = \mathbf{I}$

Consistent, but not efficient.

### 13.1.4 Moment conditions

OLS:

$$E[x(y - x\theta)] = 0$$

WLS

$$E[x(y - x\theta)/\sigma^2(x)] = 0$$

IV

$$E[z(y - x\theta)] = 0$$

MLE

$$E[\Delta_\theta \ln f(x, \theta)] = 0$$

### 13.1.5 New GMM

$$m(\theta_0) = E[g(\mathbf{x}_i, \theta_0)]$$

We replace this with sample moment

$$\hat{m}(\theta) = \frac{1}{n} \sum_i g(\mathbf{x}_i, \theta)$$

We have the "score"

$$\nabla_\theta g(\mathbf{x}_i, \theta_0)$$

Information

$$G = E[\nabla_\theta g(\mathbf{x}_i, \theta_0)]$$

Variance-covariance loss matrix

$$\Omega = E[g(\mathbf{x}_i, \theta_0) g(\mathbf{x}_i, \theta_0)^T]$$

We want to minimise moment loss

$$\|\hat{m}(\theta)\|_W^2 = \hat{m}(\theta)^T W \hat{m}(\theta)$$

$$\hat{\theta} = \operatorname{argmin}_\theta \left( \frac{1}{n} \sum_i g(\mathbf{x}_i, \theta) \right)^T \hat{W} \left( \frac{1}{n} \sum_i g(\mathbf{x}_i, \theta) \right)$$

### 13.1.6 Asymptotic

CLT means normal.

They are consistent IF moment condition is true.

There is an explicit formula for variance.

$$\sqrt{n}(\hat{\theta} - \theta_0) \rightarrow^d N[0, (G^T W G)^{-1} G^T W \Omega W^T G (G^T W^T G)^{-1}]$$

If we choose  $W \propto \Omega^{-1}$  then:

$$\sqrt{n}(\hat{\theta} - \theta_0) \rightarrow^d N[0, (G^T \Omega^{-1} G)^{-1}]$$

Problem: we need to estimate  $\Omega$  and  $G$ .

$\Omega$ : estimate from sample. allows us to choose estimator, but still leaves variance unidentified.

Do the above from OLS? This is where robust etc stuff comes from

If it is specified. Moment conditions are equal to the number of moments, then  $W$  doesn't matter. This is normal Method of Moments.

Estimating the weighting matrix

### 13.1.7 Iterated GMM

### 13.1.8 Moment-covariance matrix

### 13.1.9 Bias and variance of the GMM estimator

page on Bias and variance of the GMM estimator (cluster assumption should be part of moment condition?) part of later calculation of weighting?

Can do robust, hac, clustering as part of GMM too.

## Part V

# Supervised linear regression

## Chapter 14

# Ordinary Least Squares for prediction

### 14.1 Constructing a linear model

#### 14.1.1 Defining linear models

##### Defining

One option for  $f(X)$  is a linear model.

$$f(X_i) = \hat{Y}_i = \beta_0 + \sum_{j=1}^p \beta_j X_{ij}$$

The values for  $\beta$  are the regression coefficients.

So we have:

$$Y_i = \beta_0 + \sum_{j=1}^p \beta_j X_{ij} + e(X_i) + e_i$$

We define the error of the estimate as:

$$\epsilon_i = Y_i - \hat{Y}_i$$

$$\epsilon_i = e(X_i) + e_i$$

So:

$$Y_i = \beta_0 + \sum_{j=1}^p \beta_j X_{ij} + \epsilon_i$$

The linear model could be wrong for two reasons. No linear model could be appropriate, or the wrong coefficients could be provided for a linear model.

Linear regression if  $f$  is a linear function on  $w$ . NB: not linear in  $x$  necessarily. could have  $x^2$  etc, but still linear in  $w$ .

### 14.1.2 Intercept

### 14.1.3 Modelling non-linear functions as linear

#### Polynomials

The function  $y = x^2$  is not linear, however we can model it as linear, by including  $x^2$  as a variable.

We can expand this, and using linear models to estimate parameters for functions such as:

$$y = ax^3 + bx^2 + cx$$

#### Logarithms and exponentials

We can also transform data using logarithms and exponents.

For example we can model

$$\ln y = \theta \ln x$$

### 14.1.4 Geometric interpretation of OLS

#### Best Approximation Theorem

## 14.2 Calculating Ordinary Least Squares (OLS) estimators

### 14.2.1 Normal equation

#### Least squares

The square error is  $\sum_i (\hat{y}_i - y_i)^2$ .

The differential of this with respect to  $\hat{\theta}_j$  is:

$$2 \sum_i \frac{\delta \hat{y}_i}{\delta \hat{\theta}_j} (\hat{y}_i - y_i)$$

The stationary point is where this is zero:

$$\sum_i \frac{\delta \hat{y}_i}{\delta \hat{\theta}_j} (\hat{y}_i - y_i) = 0$$

#### Linear least squares

Here,  $\hat{y}_i = \sum_j x_{ij} \hat{\theta}_j$

Therefore:  $\frac{\delta \hat{y}_i}{\delta \hat{\theta}_j} = x_{ij}$

And so the stationary point is where

$$\sum_i x_{ij}(\sum_j x_{ij}\hat{\theta}_j - y_i) = 0$$

$$\sum_i x_{ij}(\sum_j x_{ij}\hat{\theta}_j) = \sum_i x_{ij}y_i$$

**Normal equation**

We can write this in matrix form.

$$X^T X \hat{\theta} = X^T y$$

We can solve this as:

$$\hat{\theta} = (X^T X)^{-1} X^T y$$

**Perfectly correlated variables**

If variables are perfectly correlated then we cannot solve the normal equation.

Intuitively, this is because for perfectly correlated variables there is no single best parameter, as changes to one parameter can be counteracted by changes to another.

**14.2.2 Mean and variance of predictions**

**Bias**

$$\hat{y} = \theta x$$

$$E[\hat{y} - y] = E[\theta x - y]$$

$$y = \hat{y} + \epsilon$$

$$E[y - \hat{y}|X]$$

$$E[\epsilon|X]$$

Unbiased so long as independent of error term.

**Variance**

$$Var[\hat{y} - y] = Var[\theta x - y]$$

$$Var[y - \hat{y}|X]$$

$$Var[\epsilon|X]$$

**14.2.3 The Moore-Penrose pseudoinverse**

For a matrix  $X$ , the pseudoinverse is  $(X^* X)^{-1} X^*$ .

For real matrices, this is:  $(X^T X)^{-1} X^T$

The pseudoinverse can be written as  $X^+$

Therefore  $\theta$  is the pseudoinverse of the inputs, multiplied by the outputs. Or:

$$\theta = X^+y$$

The pseudoinverse satisfies:

$$XX^+X = X$$

$$X^+XX^+ = X^+$$

### 14.2.4 Leverage

#### Introduction

Leverage measures how much the predicted value of  $y_i$ ,  $\hat{y}_i$ , changes as  $y_i$  changes.

We have:

$$\mathbf{y} = \mathbf{X}\theta + \mathbf{u}$$

$$\hat{\theta} = X(X^T X)^{-1}X^T y$$

$$\hat{\theta} = P_X y$$

The leverage score is defined as:

$$h_i = P_{ii}$$

## 14.3 Making forecasts with OLS

### 14.3.1 The projection and annihilation matrices

#### The projection matrix

We have  $X$ .

The projection matrix is  $X(X^T X)^{-1}X^T$

The projection matrix maps from actual  $y$  to predicted  $\hat{y}$

$$\hat{y} = P y$$

Each entry refers to the covariance between actual and fitted

$$p_{ij} = \frac{\text{Cov}(\hat{y}_i, y_j)}{\text{Var}(y_j)}$$

#### The annihilation matrix

We can get residuals too:

$$u = y - \hat{y} = y - p y = (1 - P)y$$

$1 - P$  is called the annihilator matrix

We can now use the propagation of uncertainty

$$\Sigma^f = A\Sigma^x A^T$$

To get:

$$\Sigma^u = (I - P)\Sigma^y(I - P)$$

Annihilator matrix is:

$$M_X = I - X(X^T X)^{-1} X^T$$

Called this because:

$$M_X X = X - X(X^T X)^{-1} X^T X$$

$$M_X X = 0$$

Is called residual maker

## 14.4 Frisch-Waugh-Lovell theorem

### 14.4.1 Introduction

If we have a partitioned linear regression model:

$$\mathbf{y} = \mathbf{X}\theta + \mathbf{Z}\beta + \mu$$

Use the annihilator matrix:

$$M_X \mathbf{y} = M_X \mathbf{X}\theta + M_X \mathbf{Z}\beta + M_X \mu$$

$$M_X \mathbf{y} = M_X \mathbf{Z}\beta + M_X \mu$$

We can then estimate  $\beta$ .

Frisch-Waugh-Lovell theorem says that this is the same estimate as the original regression.

## 14.5 Trimming

### 14.5.1 Introduction

OLS:

$$\hat{\theta} = \frac{\sum_i (X_i - \mu_X)(y_i - \mu_y)}{\sum_i (x_i - \mu_X)^2}$$

Trimming

$$\hat{\theta} = \frac{n^{-1} \sum_i (X_i - \mu_X)(y_i - \mu_y) \mathbf{1}_i}{n^{-1} \sum_i (x_i - \mu_X)^2 \mathbf{1}_i}$$

Where:

$$\mathbf{1}_i = \mathbf{1}(\hat{f}(z_i) \geq b)$$

Where  $b = b(n)$  is a trimming parameter, where:



$b \rightarrow 0$  as  $n \rightarrow \infty$

## 14.6 Best linear predictor

### 14.6.1 Introduction

The best linear predictor is the one which minimises:

$$E[Y - X\theta]$$

Under what circumstances is this the same as OLS? When  $n_{i,j}$ . When  $n$  is not, then other linear estimators (like LASSO) can be better.

## 14.7 Other

### 14.7.1 Cook's distance

Cook's distance measures the effect of deleting outliers. work out predictions if outlier was removed, sum all differences in  $\hat{y}$

Outliers have a high Cook's distance.

### 14.7.2 Bayesian linear regression

In linear regression we have

$$P(y|X, \theta, \sigma_\epsilon^2)$$

For Bayesian linear regression we want:

$$P(\theta, \sigma_\epsilon^2|y, X)$$

We can use Bayes rule:

$$P(\theta, \sigma_\epsilon^2|y, X) \propto P(y|X, \theta, \sigma_\epsilon^2)P(\theta|\sigma_\epsilon^2)P(\sigma_\epsilon^2)$$

# Chapter 15

## Regularising linear regression for prediction

### 15.1 OLS predictions with many parameters

#### 15.1.1 Too many variables

If there are more independent variables than samples then OLS will not work. There will be an infinite number of perfect fits.

For example if we regression genetic information on height with 1000 people, there will be too little data to fit using OLS.

This is due to colinearity.

We could also have too many variables through the use of derived variables. For example if we choose to use  $x$ ,  $x^2$ ,  $x^3$  etc.

#### **Optimal sparse regression**

Optimal is  $\lambda = \sigma 2\sqrt{2\log(pn)/n}$

Relies on knowing  $\sigma$ , which we may not.

Instead we can use root LASSO.

Minimise the squareroot of the sum of squares loss (over  $n$ ), and use  $\lambda = \sqrt{2\log(pn)/n}$

Doesn't have  $\sigma$

Lasso biased, estimators 0 for many.

**Post-LASSO**

We can use LASSO for model selection, then use OLS on only those estimators.

## 15.2 Least Absolute Shrinkage and Selection Operator (LASSO)

### 15.2.1 Least Absolute Shrinkage and Selection Operator (LASSO)

**Introduction**

With LASSO we add a constraint to  $\hat{\theta}$ .

$$\sum_i \hat{\theta}_i \leq t$$

Regularisation of LLS. Sum of thetas are constrained to be below hyperparameter  $t$

L1 regularisation

This is also known as sparse regression, because many weights are set to 0.

This now looks like:

$$w_{lasso} = \arg \min \|y - Xw\|_2^2 + \lambda \|w\|_1$$

**Hyperparameter**

$t$  is a hyperparameter.

### 15.2.2 Optimal hyperparameter for LASSO

### 15.2.3 Feature scaling

## 15.3 Ridge regression

### 15.3.1 Ridge regression

Regularisation of LLS. The cost function now includes a norm on  $M\theta$ .

$L_2$  regularisation

This allows us to solve problems where there are too many features.  $L_1$  also allows us to do this.

**Overspecified**

If  $n > d$  we can minimise weights subject to  $Xw=y$ . This is the same as the least norm.

**Maximum A-Priori (MAP) estimator for linear regression**

Maximum a priori estimation. equiv to ridge regression with a priori estimate of 0

$$W_{RR} = (\lambda I + X^T X)^{-1} X^T y$$

$$E[w_{RR}] = (\lambda I + X^T X)^{-1} X^T X w$$

$$Var[W_{RR}] = a$$

**15.3.2 Elasticnet**

Regularisation of LLS. Combines lasso and ridge regression.

$L_1$  and  $L_2$  regularisation

**15.3.3  $L_p$  regularisation****Introduction**

We can generalise this to:

$$w_{l_p} = \arg \min \|y - Xw\|_2^2 + \lambda \|w\|_p^p$$

For ridge regression there is always a solution.

For least squares there is a solution if  $X^T X$  is invertible

For Lasso we must use numerical optimisation.

lasso and  $L_1$  induces sparsity

Goal is  $\min \|y - f(x)\| + \lambda g(w)$

Ridge regression:  $g(w) = \|w\|^2$

If  $\lambda = 0$ , OLS, if infinite,  $w$  goes to 0.

Normal equation changes to:  $(\lambda I + X^T X)^{-1} X^T y$

We can preprocess to avoid processing of 1s. shift mean of  $y$  to 0. normalise  $x$  mean 0 var 1.

**15.3.4 Lava****Introduction**

Alternative to ElasticNet

Each parameter is split into

$$\theta_i = \rho_i + \phi_i$$

There is  $L_2$  loss on  $\rho$  and  $L_1$  loss on  $\phi$ .

This means that large coefficients can be penalised like  $L_1$  and small coefficients like  $L_2$ .

## 15.4 Tests

### 15.4.1 The Ramsey RESET test

The Ramsey Regression Equation Specification Error Test (RESET)

Once we have done our OLS we have  $\hat{y}$ .

The Ramsey RESET test is an additional stage, which takes these predictions and estimates:

$$y = \theta x + \sum_{i=1}^3 \alpha_i \hat{y}^i$$

We then run an F-test on  $\alpha$ , with the null that  $\alpha = 0$ .

### 15.4.2 The Link test

#### Introduction

Alternative to RESET

We have  $\hat{y}$ .

We regress  $y = \alpha + \beta \hat{y} + \gamma \hat{y}^2$ .

We test that  $\gamma = 0$ .

If it is not, then this suggests the model is misspecified.

## 15.5 Bias trade-off

### 15.5.1 Introduction

Trade-off between parameter accuracy and prediction accuracy.

## Chapter 16

# Choosing linear models for prediction

### 16.1 Other

#### 16.1.1 Selection

How to restrict variables? best subset. iterate through all and test.

not feasible for large numbers of variables. forward selection, backward selection, L1 alternatives.

removing variable does: increase bias. may reduce variance of prediction

## Chapter 17

# Generalised linear models, the delta rule and binary classification

### 17.1 Introduction

#### 17.1.1 Link/activation functions

### 17.2 Estimating parameters

#### 17.2.1 Delta rule

##### Introduction

We want to train the parameters  $\theta$ .

We can do this with gradient descent, by working out how much the loss function falls as we change each parameter.

The delta rule tells us how to do this.

##### The loss function

The error of the network is:

$$E = \sum_j \frac{1}{2} (y_j - a_j)^2$$

We know that  $a_j = a(\theta x_j)$  and so:

$$E = \sum_j \frac{1}{2} (y_j - a(\theta x_j))^2$$

**Minimising loss**

We can see the change in error as we change the parameter:

$$\frac{\delta E}{\delta \theta_i} = \sum_j \frac{\delta E}{\delta a_j} \frac{\delta a_j}{\delta z_j} \frac{\delta z_j}{\delta \theta_i}$$

$$\frac{\delta E}{\delta \theta_i} = - \sum_j (y_j - a_j) a'(z_j) x_{ij}$$

**Delta**

We define delta as:

$$\delta_i = - \frac{\delta E}{\delta z_j} = \sum_j (y_j - a_j) a'(z_j)$$

So:

$$\frac{\delta E}{\delta \theta_i} = \delta_i x_{ij}$$

**The delta rule**

We update the parameters using gradient descent:

$$\Delta \theta_i = \alpha \delta_i x_{ij}$$

**17.2.2 Maximum likelihood****17.2.3 Identity link function****The function**

$$a(z) = z$$

**The derivative**

$$a'(z) = 1$$

**Notes**

This is the same as ordinary linear regression.

**17.3 Link/activation functions: Classification****17.3.1 The binomial data generating process****Introduction**

For linear regression our data generating process is:

$$y = \alpha + \beta x + \epsilon$$



For linear classification our data generating process is:

$$z = \alpha + \beta x + \epsilon$$

And set  $y$  to 1 if  $z > 0$

Or:

$$y = \mathbf{I}[\alpha + \beta x + \epsilon > 0]$$

### Probability of each class

The probability that an individual with characteristics  $x$  is classified as 1 is:

$$P_1 = P(y = 1|x)$$

$$P_1 = P(\alpha + \beta x + \epsilon > 0)$$

$$P_1 = \int \mathbf{I}[\alpha + \beta x + \epsilon > 0]f(\epsilon)d\epsilon$$

$$P_1 = \int \mathbf{I}[\epsilon > -\alpha - \beta x]f(\epsilon)d\epsilon$$

$$P_1 = \int_{\epsilon=-\alpha-\beta x}^{\infty} f(\epsilon)d\epsilon$$

$$P_1 = 1 - F(-\alpha - \beta x)$$

### Example: The logistic function

Depending on the probability distribution of *epsilon* we have different classifiers.

For the logistic case we then have

$$P(y = 1|x) = \frac{e^{\alpha+\beta x}}{1 + e^{\alpha+\beta x}}$$

## 17.3.2 Perceptron (step function)

### The function

If the sum is above 0,  $a(z) = 1$ . Otherwise,  $a(z) = 0$ .

### The derivative

This has a differential of 0 at all point except 0, where it is undefined.

### Notes

This function is not smooth.

These is the activation function used in the perceptron.

Perceptron data needs to be linearly separable to train.

Even if linearly separable, doesn't necessarily get the best outcome?

### 17.3.3 Perceptron

Perceptron: one node neural network. is one or zero depending if weighted inputs enough. therefore is classification

If error, update weights

Only works if linearly separable. ie can draw linear line completely separating all inputs

Neural network has more layers

Works if data is linear

How to treat node inputs: raw, sigmoid, 0,1

For all of these want the cost function have only one solution, like least squares does. not guaranteed for all

For logistic, want to make it convex. loss =  $-\log(f(x))$  or  $-\log(1-f(x))$  depending on correct y. this is convex

How to create node inputs: sigmoid, binary cutoff

### 17.3.4 Logistic function (AKA sigmoid, logit)

**The function**

$$\sigma(z) = \frac{1}{1 + e^{-z}}$$

The range of this activation is between 0 and 1.

**The derivative**

$$\sigma'(z) = \frac{e^{-z}}{(1 + e^{-z})^2}$$

$$\sigma'(z) = \sigma(z) \frac{1 + e^{-z} - 1}{1 + e^{-z}}$$

$$\sigma'(z) = \sigma(z)[1 - \sigma(z)]$$

**Notes**

### 17.3.5 Probability unit (probit)

**The function**

The cumulative distribution function of the normal distribution.

$$\Phi(z)$$

**The derivative**

The normal distribution:

$$\Phi'(z) = \phi(z)$$

**17.3.6 tanH**

**17.3.7 ArcTan**

**17.3.8 Radial Basis Function (RBF) activation function**

$$a(x) = \sum_i a_i f(\|x - c_i\|)$$

$$a(x) = \sum_i a_i f e^{\|x - c_i\|^2}$$

**17.3.9 Linear probability model**

$p = xB$ . can be outside  $[0, 1]$ .

**17.4 Generalised Additive Models (GAMs)**

**17.4.1 Introduction**

# Chapter 18

## Generalised linear models and multiclass classification

### 18.1 Multinomial classification

#### 18.1.1 The multinomial data generating process

##### Introduction

In the binomial case we had:

$$z_i = \alpha + \beta x_i + \epsilon_i$$

And set  $y_i$  to 1 if  $z_i > 0$

In the multinomial case we have  $m$  alternatives

$$z_{ij} = \alpha + \beta x_{ij} + \epsilon_{ij}$$

And set  $y_{ij} = 1$  if  $z_{ij} > z_{ik} \forall k \neq j$

##### Generalised version

We can rewrite this as:

$$z_{ij} = v_{ij} + \epsilon_{ij}$$

Where:

$$v_{ij} = \alpha + \beta x_{ij}$$

In this case  $v$  does not depend on  $j$ , but in other formulations it could.

**Probabilities**

$$P_{ij} = P(y_{ij} = 1 | x_{ij})$$

$$P_{ij} = P(z_{ij} > z_{ik} \forall k \neq j)$$

$$P_{ij} = P(\epsilon_{ik} < v_{ij} - v_{ik} + \epsilon_{ij} \forall k \neq j)$$

**The form of the multinomial model: Intercepts**

Previously we described the multinomial model

$$z_{ij} = v_{ij} + \epsilon_{ij}$$

Where:

$$v_{ij} = \alpha + \beta x_{ij}$$

The probability of  $j$  being chosen is.

$$P_{ij} = P(\epsilon_{ik} < v_{ij} - v_{ik} + \epsilon_{ij} \forall k \neq j)$$

Intercepts in  $v$  cancel out. Therefore in the basic model there is no need to use

$$v_{ij} = \alpha + \beta x_{ij}$$

We can instead use:

$$v_{ij} = \beta x_{ij}$$

**The form of the multinomial model: Conditional model**

We have :

$$v_{ij} = \beta x_{ij}$$

What do we include in  $x_{ij}$ ?

We can include observable characteristics for each product:

$$v_{ij} = \alpha_j + \beta x_j$$

One of the  $\alpha_j$  must be normalised to 0, as only differences matter. We cannot tell the difference if all  $\alpha$  are raised by the same amount.

For consistency with other models we can write this as:

$$v_{ij} = \beta x_{ij}$$

Even though this does not vary from individual to individual.

Here  $\beta$  represents average preferences for each product characteristic.

**The form of the multinomial model: The multinomial model**

We have differing characteristics for each individual:

$$v_{ij} = \beta x_i$$

However this adds a constant for each product. For this to discriminate we need varying coefficients.

$$v_{ij} = \beta_j x_i$$

As we only observe differences, one of the  $\beta_j$  must be normalised to 0.

We can rewrite this.

$$v_{ij} = \sum_k \beta_k \delta_{kj} x_i$$

$$v_{ij} = \beta z_{ij}$$

The original  $x_i$  is dense and contains data about the individual.

$z_{ij}$  is sparse and only has entries in the  $\{j\}$  section.

Here  $\beta$  represents how the

**The form of the multinomial model: Combined multinomial and conditional model**

If we have observations of the characteristics of both individuals and alternatives we can write:

$$v_{ij} = \beta_m m_{ij} + \beta_c c_{ij}$$

$$v_{ij} = \beta x_{ij}$$

Here  $\beta$  represents both:

- Average preferences for customer characteristics (conditional)
- How preferences change as individual characteristics change (multinomial)

**18.1.2 Extreme IID multinomial****IID**

The probability of  $j$  being chosen is:

$$P_{ij} = P(\epsilon_{ik} < v_{ij} - v_{ik} + \epsilon_{ij} \forall k \neq j)$$

If these are independent then we have:

$$P_{ij} = \prod_{k \neq j} P(\epsilon_{ik} < v_{ij} - v_{ik} + \epsilon_{ij})$$

$$P_{ij} = \prod_{k \neq j} F_\epsilon(v_{ij} - v_{ik} + \epsilon_{ij})$$

We do not know  $\epsilon_{ij}$  so we have to integrate over possibilities.

$$P_{ij} = \int [\prod_{k \neq j} F_{\epsilon}(v_{ij} - v_{ik} + \epsilon_{ij})] f_{\epsilon}(\epsilon_{ij}) d\epsilon_{ij}$$

**Extreme values**

We have:

$$P_{ij} = \int [\prod_{k \neq j} F_{\epsilon}(v_{ij} - v_{ik} + \epsilon_{ij})] f_{\epsilon}(\epsilon_{ij}) d\epsilon_{ij}$$

If  $\epsilon$  is extreme value type-I this gives us:

$$P_{ij} = \frac{e^{v_{ij}}}{\sum_k e^{v_{ik}}}$$

**Independence of irrelevant alternatives**

Consider the ratio two probabilities:

$$\frac{P_{ij}}{P_{im}} = \frac{e^{v_{ij}}}{e^{v_{im}}}$$

This means that changes to any other products do not affect relative odds.

This can be undesirable. For example removing one option may cause unbalanced substitution.

For example raising the price of buses may cause more substitution to trains than helicopter, for a commute.

**18.1.3 Estimating multinomial logit models**

Estimating with individual level data.

Estimating with market share level data.

**18.1.4 Nested logit**

The probability of  $j$  being chosen is:

$$P_{ij} = P(\epsilon_{ik} < v_{ij} - v_{ik} + \epsilon_{ij} \forall k \neq j)$$

If the error terms are not IID this is more difficult to calculate.

We divide the  $J$  alternatives into nests. Within each of these we assume IID error terms, but allow variation between them.

For example we could have a nest of public/private transport. We could have a nest of types of product, and within that the firms offering the product.

The nested logit model does 2 or more sequential IID logit models. One to select the nest, and the other to select the alternative within the nest.

### 18.1.5 Mixed logit (random coefficients)

#### Introduction

In our standard model we have:

$$z_{ij} = \beta x_{ij} + \epsilon_{ij}$$

If we allow the parameters to vary for each individual we have:

$$z_{ij} = \beta_i x_{ij} + \epsilon_{ij}$$

The probability of choosing  $j$  now depends on the distribution of  $\beta$ .

In the IID case we had:

$$P_{ij} = \frac{e^{\beta x_{ij}}}{\sum_k e^{\beta x_{ik}}}$$

Rather than evaluate this at a single point  $\beta$  we integrate.

$$P_{ij} = \int \frac{e^{\beta x_{ij}}}{\sum_k e^{\beta x_{ik}}} f(\beta) d\beta$$

If  $\beta$  is degenerate this reduces to the standard logit model.

### 18.1.6 Multinomial probit

This relaxes the IID and extreme value assumption.

Errors have a normal variance-covariance matrix.

### 18.1.7 Softmax

The softmax function is often used in the last layer of a classification network.

It takes a vector of dimension  $k$  and returns another vector of the same size. Only, this time all numbers are between 0 and 1 and the values sum to 1.

The softmax function is based on the sigmoid function.

$$a_j(z) = \frac{e^{z_j}}{\sum_i e^{z_i}}$$

### 18.1.8 Temperature for Softmax



## Part VI

# Supervised machine learning

# Chapter 19

## Classification trees

### 19.1 Simple decision trees

#### 19.1.1 Tree traversal

#### 19.1.2 Training decision trees with information gain

We can train a decision tree by starting out with the most simple tree - all outcomes in same node.

We can then do a greedy search to identify which split on the node is best.

We can then iterate this process on future nodes.

#### **Training with information gain**

We split nodes to increase maximum entropy.

Entropy is:

$$E = - \sum_i^n p_{i=1} \log_2 p_i$$

Where we are summing across all nodes.

#### **Information gain**

The gain in entropy is the original entropy - weighted by size entropy of each branch

**Information gain ratio****19.1.3 Training decision trees with Gini impurity****19.1.4 Pruning decision trees**

Training a decision tree until there is only one entry from the training set will result in overfitting.

We can use pruning to regularise trees.

**Pruning****Reduced error pruning**

From bottom, replace each node with a leaf of the most popular class. Accept if no reduction in accuracy.

**Cost complexity pruning**

Take full tree  $T_0$

Iteratively find a subtree to replace with a leaf. Cost function is accuracy and number of leaves.

Remove this generating  $T_{i+1}$

When we have just the root, choose a single tree using CV.

**Growing and pruning**

Generally we would split the data up. Grow the tree with one set and then prune with the other.

We can split our data up and iterate between growing and pruning.

When pruning, for each pair of leaves we test to see if they should be merged.

If our two sets are  $A$  and  $B$  we can do:

- $A$ : Grow
- $B$ : Prune
- $B$ : Grow
- $A$ : Prune

And repeat this process.

**Partial regression trees**

Once we have built a tree, we keep a single leaf and discard the rest.

### **19.1.5 Decision trees with many classes**

Training decision trees with many classes

## **19.2 Other**

### **19.2.1 Training with unbalanced data**

Unbalanced dataset: more of one class than others.

Can reduce sample of majority, or synthetically generate minority.

# Chapter 20

## Regression trees

### 20.1 Regression trees

#### 20.1.1 Classic regression trees

In a classical regression tree, we follow a decision process as before, but the outcome is real number.

Within each leaf, all inputs are assigned that same number.

#### Training

With a regression problem we cannot split nodes the same way as we did for classification.

Instead by split by the residual sum of squares.

#### 20.1.2 Training decision trees with Mean Squared Error (MSE)

#### 20.1.3 Mixed regression trees

In classical trees all items in a leaf are assigned the same values. In this model, all are given  $\theta$  for a parametric model.

This makes the resulting trees smoother.

We have some  $\hat{y}_i = f(\mathbf{x}_i, \theta) + \epsilon$

The approach generalises classic regression trees. There the estimate was  $\bar{y}$ . Here it's a regression.

**Training**

At each node we do OLS. If the  $R^2$  of the model is less than some constant, we find a split which maximises the minimum of the two new  $R^2$ .

**20.1.4 Classifying with probabilistic decision trees**

Previously our decision tree classifier was binary.

We can instead adapt the mixed tree model and using a probit model at each leaf.

# Chapter 21

## Bayesian trees

### 21.1 Bayesian trees

#### 21.1.1 Priors of trees

##### Priors for simple trees

We can define a tree as a set of nodes:  $T$ .

For each node we define a splitting variable  $k$  and a splitting threshold  $r$ .

Our prior is  $P(T, k, r)$ .

We split this up to:

$$P(T, k, r) = P(T)P(k, r|T)$$

$$P(T, k, r) = P(T)P(k|T)P(r|T, k)$$

So we want to estimate:

- $P(T)$  - The number of nodes.
- $P(k|T)$  - Which variables we split by, given the tree size.
- $P(r|T, k)$  - The cutoff, given the tree size and the variables we are splitting by.

##### Priors for mixed trees

If at the leaf we have a parametric model, our prior is instead:

$$P(T, k, r, \theta) = P(T)P(k|T)P(r|T, k)P(\theta|T, k, r)$$

We then need to additionally estimate  $P(\theta|T, k, r)$ .

### 21.1.2 The pinball prior

We can generate a tree with a fixed number of leaves, according to our prior.

As we start the tree we associate the root node with a count of all leaves.

As we split a node, the remaining leaf counts are divided between the directions. If there is only one leaf left, we do no further splitting.

### 21.1.3 Estimating other priors

#### 21.1.4 Bayesian CART

##### Our prior

Call the collective parameters of the tree  $\Theta = (T, k, r)$  and  $\theta$ .

Collectively our prior is defined by  $P(\Theta)$  and  $P(\theta)$

##### Bayes' theorem

We want to know the posterior given our data  $X$ .

$$P(\Theta|X) = \frac{P(X|\Theta)P(\Theta)}{P(X)}$$

$$P(\Theta|X) \propto P(X|\Theta)P(\Theta)$$

##### Expanded posterior

We now explore  $P(X|\Theta)$

$$P(X|\Theta) = \int P(X|\theta, \Theta)P(\theta)d\theta$$

This means our posterior is:

$$P(\Theta|X) \propto P(\Theta) \int P(X|\theta, \Theta)P(\theta)d\theta$$

##### Estimation

This can be estimated with MCMC.



## Chapter 22

# Support Vector Machines (SVMs)

### 22.1 Linear Support Vector Classifiers (SVCs)

#### 22.1.1 Hard-margin SVC

##### Linear separators

We want to create a hyperplane to separate classes.

For classification problem  $(x, y)$

Hyperplane is  $wx-b=0$

##### Hard margin

If data is linearly separable then a hyperplane exists such that all data can be correctly classified

There are an infinite number that could work.

We select two parallel with distance between as large as possible. the region between these two is the margin

The maximum margin hyperplane is the one between the two margin planes

We can rescale the two hyperplanes to:

$$wx-b=1$$

$$wx-b=-1$$

The distance between the two parallel hyperplanes is  $\frac{2}{\|w\|}$

So we minimise  $\|w\|$  conditional on all points being correctly classified

$$y_i(wx_i - b) \geq 1$$

We select  $w$  and  $b$  to solve this.

### 22.1.2 Support vectors

Support vectors are those that make up the classifier boundary.

### 22.1.3 Estimating the SVC using quadratic equations

### 22.1.4 Soft-margin SVC

#### Soft margin

Soft margin

Data may not be linearly separable, so we introduce a hinge loss function

$$\text{Max}(0, 1 - y_i(wx - b))$$

We then minimise

$$\lambda\|w\|^2 + \left[\frac{1}{n} \sum_{i=1}^n \max(0, 1 - y_i(wx_i - b))\right]$$

This introduces  $\lambda$  as a parameter.

### 22.1.5 Regularising the SVC

## 22.2 Multiple classes

### 22.2.1 Support vector classifiers for multiple classes

## 22.3 Non-linear support vector classifiers

### 22.3.1 The dot product SVC

### 22.3.2 The kernel trick

We can use kernels as an alternative to the dot product.

### 22.3.3 The radial basis function (RBF) SVC

## Chapter 23

# Variational Bayes

### 23.1 Variational Bayes

#### 23.1.1 Introduction

# Chapter 24

## The Naive Bayes classifier

### 24.1 Naive Bayes

#### 24.1.1 The Naive Bayes posterior

##### Bayes theorem

Consider Bayes' theorem

$$P(y|x_1, x_2, \dots, x_n) = \frac{P(x_1, x_2, \dots, x_n|y)P(y)}{P(x_1, x_2, \dots, x_n)}$$

Here,  $y$  is the label, and  $x_1, x_2, \dots, x_n$  is the evidence. We want to know the probability of each label given evidence.

The denominator,  $P(x_1, x_2, \dots, x_n)$ , is the same for all, so we only need to identify:

$$P(y|x_1, x_2, \dots, x_n) \propto P(x_1, x_2, \dots, x_n|y)P(y)$$

##### The assumption of Naive Bayes

We assume each  $x$  is independent. Therefore:

$$P(x_1, x_2, \dots, x_n|y) = P(x_1|y)P(x_2|y)\dots P(x_n|y)$$

$$P(y|x_1, x_2, \dots, x_n) \propto P(x_1|y)P(x_2|y)\dots P(x_n|y)P(y)$$

#### 24.1.2 The Naive Bayes classifier

##### Calculating the Naive Bayes estimator

With the Naive Bayes assumption we have:

$$P(y|x_1, x_2, \dots, x_n) \propto P(x_1|y)P(x_2|y)\dots P(x_n|y)P(y)$$

We now choose  $y$  which maximises this.

This is easier to calculate, as there is less of a sample restriction.

This is used when evidence is also in classes, as the chance of any individual outcome on a continuous probability is 0.

### Estimating $P(y)$

We can easily calculate  $P(y)$ , by looking at the frequency across the sample.

### Estimating $P(x_1|y)$

Normally,  $P(x_1|y) = \frac{n_c}{n_y}$ , where:

- $n_c$  is the number of instances where the evidence is  $c$  and the label is  $y$ .
- $n_y$  is the number of instances where the label is  $y$ .

### Regularising the Naive Bayes estimator

To reduce the risk of specific probabilities being zero, we can adjust them, so that:

$P(x_1|y) = \frac{n_c + mp}{n_y + m}$ , where:

- $p$  is the prior probability. If this is unknown, use  $\frac{1}{k}$ , where  $k$  is the number of classes.
- $m$  is a parameter called the equivilant sample size.

## 24.1.3 Text classification using Naive Bayes

### Naive Bayes and text classification

Naive Bayes can be used to classify text documents. The  $x$  variables can be appearances of each word, and  $y$  can be the document classification.

## Chapter 25

# The K-Nearest Neighbours (KNN) classifier

### 25.1 K-Nearest Neighbours (KNN)

#### 25.1.1 K-nearest neighbours

K-nearest neighbours is a non-parametric classifier. For a point with an unknown class we identify the classes of the  $K$ -nearest neighbours and assign the most common class.

For this we need to find the distance between observations. We can do this using norms.

Weightings of the dependent variables is important here.

#### **Practicality**

Requires space to store

$N$  samples,  $d$  features

times:  $O(n \cdot d)$

#### 25.1.2 Choosing $K$ for K-Nearest Neighbours

## Chapter 26

# Discriminant analysis

### 26.1 Discriminant analysis

#### 26.1.1 Linear Discriminant Analysis (LDA)

Assume data is mixed gaussian. Use this to estimate classes.

#### 26.1.2 Kernel Fisher discriminant analysis

Use LDA with kernel feature spaces.

## Chapter 27

# Non-parametric regression

### 27.1 Kernel regression

#### 27.1.1 Kernel regression

##### Introduction

For parametric regression we have:

$$y = f(X)$$

Where the form of  $f(X)$  is fixed, such as for linear regression.

For non-parametric regression we have:

$$y = m(X)$$

Where  $m(X)$  is not fixed.

We can estimate  $m(X)$  using kernel regression.

$$m(X) = \frac{\sum_{i=1}^n K_h(x - x_i) y_i}{\sum_{i=1}^n K_h(x - x_i)}$$

We know this because we have:

$$E(y|X) = \int y f(y|x) dy = \int y \frac{f(x, y)}{f(x)} dy$$

We then use kernel density estimation for both.

### 27.2 Splines

#### 27.2.1 Multivariate Adaptive Regression Splines (MARS)

A linear model looks like:



$$\hat{y} = c + \sum_i x_i \theta_i$$

MARS instead produces a linear model for subsets of X.

$$\hat{y} = c + \sum_i B_j(x_i, a_j) \theta_i$$

Where:

- $B_j = \max(0, x_i - a_j)$ ; or
- $B_j = -\max(0, a_j - x_i)$

This is trained using a forward pass and a backward pass.

**Forward pass**

**Backward pass**

## 27.2.2 Bayesian splines

## 27.3 Other

### 27.3.1 Local regression

### 27.3.2 LOWESS

### 27.3.3 LOESS

### 27.3.4 Kernel regression

**Quantile regression**

In other supervised?

Normally we return a central estimate, commonly the mean.

Quantile regression returns an estimate of the  $i$ th quartile instead.

Goal is to find  $x$ th quartile of variance.

**Linear quantile regression**

**Tree quantile regression**

## 27.3.5 Principal component regression

Do PCA on X.

Do OLS with this.

Transform parameters by reversing PCA procedure on parameters.

### **27.3.6 Partial least squares regression**

This expands on principal component regression.

Both  $X$  and  $Y$  are mapped to new spaces.

# Chapter 28

## Ensemble methods

### 28.1 Combining learners

#### 28.1.1 Majority voting

##### Combining classifiers

If we generate different classifiers then each can give different predictions for the same input.

If we have different predictions how should we proceed?

##### Using one model or multiple models

It is not obvious that we should want to use more than one model. If one model was superior to another then there may be no benefit to using the information from the additional model.

However if the errors in different models are varied, then combining multiple models can lead to better performance as each individual model can have unique information.

##### Majority voting

One approach to using different predictions is to use majority voting.

If we have a collection of hard classifiers then we choose the classification with the most votes.

##### Condorcet's Jury Theorem

Consider a collection of classifiers. For each classifier there is a chance  $p_i$  of the classification being correct.

If the voters are independent

If voters are independent, and the chance of one vote being right is greater than 0.5, then the more voters, the better.

A weak model can still be useful, if it is independent.

### 28.1.2 Averaging regression predictions

If we have multiple predictions, we can take an average of these, possibly weighted.

We have  $m$  regressors  $g_j(\mathbf{x}_i)$ .

Our output is:

$$h(\mathbf{x}_i) = \sum_j w_j g_j(\mathbf{x}_i)$$

### 28.1.3 Stacking and SuperLearner

#### Introduction

With stacking we take the predictions from each of our classifiers, and then train a new model using these predictions as inputs.

#### Hard and soft inputs

Hard classification (0 or 1) and soft classification (between 0 and 1) can be used as inputs.

#### Cross validation

We can select hyper-parameters using cross validation, however there is an issue of using cross validation twice on the same data. Once for the underlying classifiers, and again for the stacked model.

#### SuperLearner

We have  $m$  models.

Part 1: Train each of them on all data.

Part 2: Split the data into  $k$  sets

For each set, associate all other data as training

For each fold:

- Fit each model on the training data
- Predict on other
- Create weighted predictor. choose weightings to minimise error

Part 3: Use these weightings on the original (unrestricted data) model

## 28.2 Generating learners with boosting

### 28.2.1 $L_2$ boosting

### 28.2.2 Adaboost

#### Introduction

Boosting is a way to create multiple learners for use in an ensemble predictor.

The goal is to create many predictors, which may not be themselves very accurate, but have a high degree of independence.

#### AdaBoost

AdaBoost is a popular algorithm for boosting.

It works by:

- Creating a set of weak learners using different restrictions on features in the training data.
- Choosing the weak learner that most reduces the error of the combined learners, and give it a weighting which most reduces the error of the combined learners.
- Creating a new weighting for the dataset, where ones poorly predicted (by the combination of learners) are given high weights.
- Repeating the process a fixed number of times.

### 28.2.3 Gradient boosting

Gradient boosting does not iterative change the weights for the learners. Instead, it trains on different errors.

While AdaBoost trains to reduce the absolute error for each weak classifier, gradient boosting trains on the difference between the actual classification and the current classification.

## 28.3 Generating learners with Bootstrapped AGGregation (bagging)

### 28.3.1 Bagging

#### Introduction

Bagging, or Bootstrap AGGregation, is a way of generating weak learners.

### **Bootstrapping**

Bootstrapping refers to taking samples with replacement from the training set. This is how the datasets for each of the weak learners are formed.

### **How to do bagging**

We take samples from the training set, with replacement, and train each of these separately. This gives us our weak learners.

## Chapter 29

# Ensemble methods for trees

### 29.1 Ensemble methods for trees

#### 29.1.1 Gradient tree boosting

This applies gradient boosting to tree.

#### 29.1.2 Multiple Additive Regression Trees (MART)

#### 29.1.3 Bayesian Additive Regression Trees (BART)

#### 29.1.4 Extra trees

#### 29.1.5 Random forests

These use bagging techniques with random trees.

At each node, rather than sample the whole data we sample a random selection.

Get  $d$  dimensions, and sample  $m$  of them at each node.

Choose  $m \leq \sqrt{d}$

#### 29.1.6 Regression forests

## Chapter 30

# Regularising black box models

### 30.1 Regularising classifiers

#### 30.1.1 Label smoothing

Regularising of classifiers assume some incorrect (similar for regression?)



## Chapter 31

# Confidence intervals of black box models

### 31.1 Confidence intervals of black box models

#### 31.1.1 Introduction

ensemble confidence intervals non-parametric confidence intervals semi-parametric confidence intervals

one-sided, 2 sided confidence intervals

jackknife for bagging confidence interval

### 31.2 Bootstrapping moments of ensemble statistics

#### 31.2.1 Bootstrapping mean and variance

Sample size selection

#### 31.2.2 Bootstrapping confidence intervals

Need to know estimate is unbiased for this.

# Chapter 32

## Interpreting black box models

### 32.1 Interpreting black box models

#### 32.1.1 Introduction

partial dependence plots can be used on black box models

Interpretation: if its interpretable then you can adjust an interpretable part manually part way?

Transparency of models: Sparse linear models are more transparent. Decomposition: Can each part of the model have input, output, parameters which can be interpreted? Complex feature selection means loss of this. Complex models. Boosting. Loses Can we say formal things about performance? We can for linear models (?), but not for others For agents, we can't validate their behaviour. We can for manually defined rules. We can for interpretive models. Post-hoc interpretability

LIME Local Interpretable Model-Agnostic Explanations. Page on explainable models, h3 in that on locally explainable models

Explaining models: We may want to understand how it works. Black box algorithms are hard to understand.

This is important if the algorithm is used in high stakes cases, or where data is different to the static case used for training.

We can create explainable models (sparse linear models)

Take black box models and make them explainable.

SLAM algorithms?

Local explanation? Saliency maps?

Why do we care about transparency?

## Chapter 33

# Semi-supervised learning

### 33.1 Introduction

#### 33.1.1 Introduction

We seek to make a discriminative model using both labelled and unlabelled data.

#### 33.1.2 Generative models

#### 33.1.3 Low density separation

## Chapter 34

# Imputing missing data for prediction

### 34.1 Introduction

#### 34.1.1 Deleting observations

Deleting whole row if missing data (bias if not random)

#### 34.1.2 Linear interpolation

#### 34.1.3 Missing Completely At Random (MCAR)

Different to MAR

#### 34.1.4 Missing at Random (MAR)

#### 34.1.5 Other

Interpolation: mean, conditional (in IID, different for time series)

## Chapter 35

# Recommenders

### 35.1 Non-negative matrix factorisation

#### 35.1.1 Non-negative matrix factorisation

### 35.2 Recommenders

#### 35.2.1 The recommendation problem

We have a collection of users and a collection of products. We want to recommend products to customers.

Once a customer "consumes" something, we get feedback. however for vast majority of items, not consumed. goal: predict feedback score beforehand to make good recommendations.

We have a matrix of how much each customer "likes" things. however this is sparse.

### 35.3 Approaches

#### 35.3.1 Content filtering

We have metadata on customers and products. Eg stated preferences, genres, actors etc.

We use this to create recommendations.

#### 35.3.2 Collaborative filtering

We look at the actions of customers. We then recommend things based on customers who are similar.

Doesn't need stated preferences, or metadata on content.

Needs data on customers (behaviours, etc)

2 steps:

- Look for similar users
- Use ratings from those users to make predictions for missing items

### **35.3.3 Cold start**

When you start you have no data on views. This is the cold start problem.

## Part VII

# Linear regression for inference



## Chapter 36

# Ordinary Least Squares for inference

### 36.1 Bias of OLS estimators

#### 36.1.1 Expectation of OLS estimators

##### Expectation in terms of observables

We have:  $\hat{\theta} = (X^T X)^{-1} X^T y$

Let's take the expectation.

$$E[\hat{\theta}] = E[(X^T X)^{-1} X^T y]$$

##### Expectation in terms of errors

Let's model  $y$  as a function of  $X$ . As we place no restrictions on the error terms, this is not an assumption.

$$y = X\theta + \epsilon.$$

$$E[\hat{\theta}] = E[(X^T X)^{-1} X^T (X\theta + \epsilon)]$$

$$E[\hat{\theta}] = E[(X^T X)^{-1} X^T X\theta] + E[(X^T X)^{-1} X^T \epsilon]$$

$$E[\hat{\theta}] = \theta + E[(X^T X)^{-1} X^T \epsilon]$$

$$E[\hat{\theta}] = \theta + E[(X^T X)^{-1} X^T] E[\epsilon] + cov[(X^T X)^{-1} X^T, \epsilon]$$

##### The Gauss-Markov: Expected error is 0

$$E[\epsilon] = 0$$

This means that:

$$E[\hat{\theta}] = \theta + cov[(X^T X)^{-1} X^T, \epsilon]$$

**The Gauss-Markov: Errors and independent variables are uncorrelated**

If the error terms and  $X$  are uncorrelated then  $E[\epsilon|X] = 0$  and therefore:

$$E[\hat{\theta}] = \theta$$

So this is an unbiased estimator, so long as the condition holds.

## 36.2 Variance of OLS estimators

### 36.2.1 Variance of OLS estimators

**Variance-covariance matrix**

We know:

$$\hat{\theta} = (X^T X)^{-1} X^T y$$

$$y = X\theta + \epsilon$$

Therefore:

$$\hat{\theta} = (X^T X)^{-1} X^T (X\theta + \epsilon)$$

$$\hat{\theta} = \theta + (X^T X)^{-1} X^T \epsilon$$

$$\hat{\theta} - \theta = (X^T X)^{-1} X^T \epsilon$$

$$Var[\hat{\theta}] = E[(\hat{\theta} - \theta)(\hat{\theta} - \theta)^T]$$

$$Var[\hat{\theta}] = E[(X^T X)^{-1} X^T \epsilon (X^T X^{-1} X^T \epsilon)^T]$$

$$Var[\hat{\theta}] = E[(X^T X)^{-1} X^T \epsilon \epsilon^T X (X^T X)^{-1}]$$

$$Var[\hat{\theta}] = (X^T X)^{-1} X^T E[\epsilon \epsilon^T] X (X^T X)^{-1}$$

We write:

$$\Omega = E[\epsilon \epsilon^T]$$

$$Var[\hat{\theta}] = (X^T X)^{-1} X^T \Omega X (X^T X)^{-1}$$

Depending on how we estimate  $\Omega$ , we get different variance terms.

**Variance under IID**

If IID:

$$\Omega = I\sigma_\epsilon^2$$

$$Var[\hat{\theta}] = (X^T X)^{-1} X^T I\sigma_\epsilon^2 X (X^T X)^{-1}$$

$$Var[\hat{\theta}] = \sigma_\epsilon^2 (X^T X)^{-1}$$

### 36.2.2 Heteroskedasticity-Consistent (HC) standard errors

#### Variance of OLS estimators

$$\text{Var}[\hat{\theta}] = (X^T X)^{-1} X^T \Omega X (X^T X)^{-1}$$

#### Robust standard errors for heteroskedasticity

$$\Omega_{ij} = \delta_{ij} \epsilon_i \epsilon_j$$

These are also known as the Eicker-Huber-White standard errors, or the White correction.

These are also referred to as robust standard errors.

## 36.3 Properties of the OLS estimator

### 36.3.1 Maximum Likelihood Estimator (MLE) and OLS equivalence

#### The OLS estimator

$$\hat{\theta}_{OLS} = (X^T X)^{-1} X^T y$$

$$E[\hat{\theta}_{OLS}] = w$$

$$\text{Var}[\hat{\theta}_{OLS}] = \sigma^2 (X^T X)^{-1}$$

#### The MLE estimator

$$y_i = \mathbf{x}_i \theta + \epsilon_i$$

$$P(y = y_i | x = x_i) = P(\epsilon_i = y_i - \mathbf{x}_i \theta)$$

If we assume  $\epsilon_i \sim N(0, \sigma_\epsilon^2)$  we have:

$$P(y = y_i | x = x_i) = \frac{1}{\sqrt{2\pi\sigma_\epsilon^2}} e^{-\frac{(y_i - \mathbf{x}_i \theta)^2}{2\sigma_\epsilon^2}}$$

$$L(X, \theta) = \prod_{i=1}^n \frac{1}{\sqrt{2\pi\sigma_\epsilon^2}} e^{-\frac{(y_i - \mathbf{x}_i \theta)^2}{2\sigma_\epsilon^2}}$$

$$l(X, \theta) = \sum_{i=1}^n -\frac{1}{2} \ln(2\pi\sigma_\epsilon^2) - \frac{(y_i - \mathbf{x}_i \theta)^2}{2\sigma_\epsilon^2}$$

$$\frac{\delta l}{\delta \theta_j} = \sum_{i=1}^n 2x_{ij} \frac{y_i - \mathbf{x}_i \theta}{2\sigma_\epsilon^2}$$

$$\sum_{i=1}^n x_{ij} (y_i - \hat{\theta}_{MLE} \mathbf{x}_i) = 0$$

$$X^T(y - X\hat{\theta}_{MLE}) = 0$$

$$X^T y = X^T X \hat{\theta}_{MLE}$$

$$\hat{\theta}_{MLE} = (X^T X)^{-1} X^T y$$

### Equivalence

If errors are normally IID then:

$$\hat{\theta}_{OLS} = \hat{\theta}_{MLE}$$

### 36.3.2 Gauss-Markov theorem

Mean of errors zero + If the model should only have errors on upside or downside for some reason, OLS will not provide this.

Homoscedastic (all have the same variance) + The results aren't biased, but variances etc are

Errors are uncorrelated + (this would mean you should add lagged variables etc)

show bias from each GM violation

OLS is BUE under normally distributed errors

OLS is BLUE for non-normally distributed errors

## 36.4 Selection

### 36.4.1 T-test selection

### 36.4.2 Post-LASSO

## 36.5 Heteroskedasticity

### 36.5.1 Checking for heteroskedasticity using the White test

### 36.5.2 Robust standard errors

### 36.5.3 Noise

### 36.5.4 Regression dilution

Noise in  $y$  doesn't cause bias.

Noise in  $x$  does cause bias.

Need to correct.

### **36.5.5 Causality**

### **36.5.6 Introduction**

Causality v correlation. If just getting correlation, could have bad out of sample performance

Section on causality. Difference between disease causes symptom and symptom causes disease

Linear models can be manipulated to have any variable on the left.

## Chapter 37

# Testing regression parameter estimates with Z-tests and T-tests

## Chapter 38

# Multiple hypothesis testing

### 38.1 Multiple hypothesis testing

#### 38.1.1 P-hacking

Likely to see some significant results from random chance.

#### 38.1.2 Family-Wise Error Rate (FWER)

What is the chance of making at least one false positive result?

Number of tests:  $m$

Number of false positive results:  $V$

$$FWER = P(V > 0)$$

#### 38.1.3 False Discovery Rate (FDR)

The proportion of false discoveries is:

$$Q = \frac{V}{V+S}$$

Where:  $V$  is the number of false positives

$S$  is the number of true positives

The FRD is  $E[Q]$ .

#### 38.1.4 The Bonferroni correction

We change the significance level.

reject if  $p \leq \frac{\alpha}{m}$

If  $m = 1$  this is the standard test.



## Chapter 39

# Generalised Least Squares

### 39.1 Generalised Least Squares (GLS)

#### 39.1.1 The Generalised Least Squares (GLS) estimator

##### Introduction

We make the same assumptions as OLS.

$$\mathbf{y} = \mathbf{X}\theta + \epsilon$$

We assume:

- $E[\epsilon|\mathbf{X}] = 0$
- $Cov[\epsilon|\mathbf{X}] = \Omega$

##### The GLS estimator

GLS estimator is:

$$\hat{\theta}_{GLS} = \operatorname{argmin}_b (y - Xb)^T \Omega^{-1} (y - Xb)$$

$$\hat{\theta}_{GLS} = (X^T \Omega^{-1} X)^{-1} X^T \Omega^{-1} y$$

This is the vector that minimises the Mahalanobis distance.

This is equivalent to doing OLS on a linearly transformed version of the data.

##### Identifying $\Omega$

If  $\Omega$  is known, we can proceed. Generally, however,  $\Omega$  is not known, and so the GLS estimate is infeasible.

## **39.2 Feasible Generalised Least Squares (FGLS)**

### **39.2.1 The Feasible Generalised Least Squares (FGLS) estimator**

#### **Introduction**

We do OLS to get a consistent estimate of  $\Omega$ ,  $\hat{\Omega}$ .

We then plug this into the GLS estimator.

## **39.3 Heteroskedasticity**

### **39.3.1 Weighted least squares**

## **39.4 Bias and variance of the GLS estimator**

### **39.4.1 Introduction**

you have the same sandwich term as before, so same process, right?

## **39.5 Linear discriminant analysis**

# Chapter 40

## General Linear Models

### 40.1 Cross-sectional regression

#### 40.1.1 The cross-sectional model

##### Hierarchical data

Our standard linear model is:

$$y_i = \alpha + X_i\theta + \epsilon_i$$

If we had two sets of data we could view these as:

$$y_{i,0} = \alpha_0 + X_{i,0}\theta_0 + \epsilon_{i,0}$$

$$y_{i,1} = \alpha_1 + X_{i,1}\theta_1 + \epsilon_{i,1}$$

Here, the data from 1 does not affect the parameters in 2.

##### Pooled data

If we think the data generating process is similar between models, then by restricting the freedom of parameters between models we can get more data for each estimate.

For example if we think that all parameters are the same between the models we can estimate:

$$y_{i,0} = \alpha + X_{i,0}\theta + \epsilon_{i,0}$$

$$y_{i,1} = \alpha + X_{i,1}\theta + \epsilon_{i,1}$$

Or:

$$y_{ij} = \alpha + X_{ij}\theta + \epsilon_{ij}$$

**Fixed slopes**

Intercepts may be different between the groups. In this case we can instead use the model:

$$y_{ij} = \alpha + X_{ij}\theta + \xi_j + \epsilon_{ij}$$

There are different ways of estimating this model:

- Pooled OLS
- Fixed effects
- Random effects

**40.1.2 Unbalanced data****40.2 The pooled OLS estimator****40.2.1 Pooled OLS****Introduction**

Our model is:

$$y_{ij} = \alpha + X_{ij}\theta + \xi_j + \epsilon_{ij}$$

**The pooled OLS estimator****40.3 The fixed effects estimator****40.3.1 Within and between transformation****Introduction**

We can group the data in two ways, one gets between differences and the other within differences.

In the above example, we could find the effects of schools, or of departments.

$$y_{ij} = \alpha + X_{ij}\theta + \epsilon_{ij}$$

$$(y_{ij} - \bar{y}_j) = (\alpha - \bar{\alpha}) + (X_{ij} - \bar{X}_j)\theta + (\epsilon_{ij} - \bar{\epsilon}_j)$$

$$(y_{ij} - \bar{y}_j) = (X_{ij} - \bar{X}_j)\theta + (\epsilon_{ij} - \bar{\epsilon}_j)$$

Or alternatively:

$$(y_{ij} - \bar{y}_i) = (X_{ij} - \bar{X}_i)\theta + (\epsilon_{ij} - \bar{\epsilon}_i)$$

Regardless of the form we choose, we can write this as:

$$\ddot{y}_{ij} = \ddot{X}_{ij}\theta + \ddot{\epsilon}_{ij}$$

### 40.3.2 The fixed effects estimator

#### Recap on the model

Our model is:

$$y_{ij} = \alpha + X_{ij}\theta + \xi_j + \epsilon_{ij}$$

#### The fixed effects estimator

With fixed effects we assume that  $U_{ij}$  is a constant for each group. That is:

$$U_{ij} = \delta_{ij}U_j$$

$$y_{ij} = \alpha + X_{ij}\theta + \epsilon_{ij} + \delta_{ij}U_j$$

We can use this in a regression if the standard assumptions of OLS are met. In particular, that group membership is uncorrelated with the error term.

We add these dummies to  $X_{ij}$  and regress:

$$y_{ij} = \alpha + X_{ij}\theta + \epsilon_{ij}$$

The parameter for the dummy is the fixed effect of group membership.

As we are including membership in the dependent variables, there is no problem if group membership correlates with other independent variables.

#### Using the within and between transformations

$$(y_{ij} - \bar{y}_i) = (X_{ij} - \bar{X}_i)\theta + (U_{ij} - \bar{U}_i) + (\epsilon_{ij} - \bar{\epsilon}_i)$$

Or:

$$\ddot{y}_{ij} = \ddot{X}_{ij}\theta + \ddot{\epsilon}_{ij}$$

This this get the same outcome, but is a different computational process.

## 40.4 The random effects estimator

### 40.4.1 The random effects estimator

#### Introduction

Our model is:

$$y_{ij} = \alpha + X_{ij}\theta + \xi_j + \epsilon_{ij}$$

#### FGLS recap

#### The random effects estimator

For fixed effects, we had the requirement that group membership be uncorrelated with the error term, but that it could be correlated with other independent

variables.

For random effects models, group membership cannot be correlated with other variables.

We have:

$$y_{ij} = \alpha + X_{ij}\theta + \epsilon_{ij} + U_{ij}$$

We now model  $U_{ij} = \bar{U}_j + \rho_j$ .

$$y_{ij} = \alpha + X_{ij}\theta + \epsilon_{ij} + \bar{U}_j + \rho_j$$

This randomness of the effect implies, for example, that if we ran the survey again we would expect a different effect

### **Clustering standard error**

#### **Estimation**

We use GLS.

## **40.5 Choosing the model form**

### **40.5.1 The Hausman specification test**

#### **Introduction**

The Hausman specification test allows you to choose between a fixed effects model and a random effects model.

#### **Efficiency**

Random effects models are more efficient.

## **40.6 The mixed effects estimator**

### **40.6.1 The mixed effects estimator**

#### **Introduction**

## **40.7 Manipulating data**

### **40.7.1 Disaggregation**

Used in polls

### **40.7.2 Multilevel Regression with Poststratification (Mr P)**

## Part VIII

# Advanced inference

## Chapter 41

# Analysis of variance (ANOVA)

### 41.1 Cross-sectional data

#### 41.1.1 Cross-sectional data

#### 41.1.2 Group means and the grand mean

Introduction

#### 41.1.3 Within-group variance and between-group variance

Introduction

### 41.2 Analysis of variance (ANOVA)

#### 41.2.1 Analysis of variance (ANOVA) table



## Chapter 42

# Instrumental Variables

### 42.1 Motivation

#### 42.1.1 Bias of OLS estimator from omitted variables

#### 42.1.2 Bias of OLS estimator from measurement error

### 42.2 Parameter estimation for simultaneous equations

#### 42.2.1 Structural and reduced forms

#### 42.2.2 Parameter identification problem with simultaneous equations

#### Identification terminology

A system is under-identified if there are not enough estimators for all structural parameters.

A system is exactly identified if there are the same number of estimators as structural parameters.

A system is over-identified if there are more estimators than structural parameters.

In general we have in our structural form:

$$\sum_i^n \beta_{ij} y_i = \sum_i^m \gamma_{ij} x_i + \epsilon_j$$

This is a system with  $n$  endogenous variables and  $m$  exogenous variables.

We can write this in matrix form.

$$B\mathbf{y} = \Gamma\mathbf{x} + \epsilon$$

We can use this to get:

$$\mathbf{y} = B^{-1}\Gamma\mathbf{x} + B^{-1}\epsilon$$

We estimate by placing restrictions on  $\Gamma$ .

### Structural models

If our data generating process is:

$$Q = \alpha + \beta P + \epsilon$$

We can estimate  $\alpha$  and  $\beta$  through measuring  $P$  and  $Q$ .

If, however the data generating process involves simultaneous equations, we can have:

$$Q = \alpha_1 + \beta_1 P + \epsilon_1$$

$$Q = \alpha_2 + \beta_2 P + \epsilon_2$$

### Reduced form

We can reduce this:

$$\alpha_1 + \beta_1 P + \epsilon_1 = \alpha_2 + \beta_2 P + \epsilon_2$$

$$(\alpha_1 - \alpha_2) + (\beta_1 - \beta_2)P + (\epsilon_1 - \epsilon_2) = 0$$

$$P = \frac{\alpha_2 - \alpha_1}{\beta_1 - \beta_2} + \frac{\epsilon_2 - \epsilon_1}{\beta_1 - \beta_2}$$

We can rewrite this as:

$$P = \pi_1 + \tau_1$$

Similarly we can reduce for  $Q$ :

$$Q = \frac{\alpha_2\beta_1 - \alpha_1\beta_2}{\beta_1 - \beta_2} + \frac{\beta_1\epsilon_2 - \beta_2\epsilon_1}{\beta_1 - \beta_2}$$

$$Q = \pi_2 + \tau_2$$

### We can't directly estimate structural models

If  $P$  is correlated with  $\epsilon_1$  or  $\epsilon_2$  then our estimates for  $\beta_1$  and  $\beta_2$  will be biased.

This also affects  $Q$ .

From the reduced forms we can see that  $P$  will be correlated, due to simultaneity.

**The identification problem**

We can estimate  $\pi_1$  and  $\pi_2$ , but this does not allow us to identify any of the structural parameters.

We have 2 estimators, but 4 parameters.

This is the identification problem.

**42.3 2 Stage OLS****42.3.1 2 Stage OLS (2SOLS) estimator****Motivation**

If  $x$  is correlated with the error term the OLS estimate will be biased.

**2 Stage OLS - first stage**

We have

$$y_i = x_i\theta + \epsilon_i$$

$$x_i = z_i\rho + \mu_i$$

We do OLS on the second to get  $\hat{\rho}$ .

$$\hat{\rho} = (Z^T Z)^{-1} Z^T X$$

We use this to get predicted values of  $X$ .

$$\hat{X} = Z\rho = Z(Z^T Z)^{-1} Z^T X = P_Z X$$

**2 Stage OLS - second stage**

We then regress  $y$  on the estimated  $X$ :

$$y_i = \hat{x}_i\theta + \epsilon_i$$

Our prediction is then:

$$\theta_{2SOLS} = (\hat{X}^T \hat{X})^{-1} \hat{X}^T y$$

$$\theta_{2SOLS} = ((P_Z X)^T P_Z X)^{-1} (P_Z X)^T y$$

$$\theta_{2SOLS} = (X^T P_Z X)^{-1} X^T P_Z y$$

If the dimension of  $Z$  is the same as  $X$  this collapses to:

$$\theta_{2SOLS} = (Z^T X)^{-1} Z^T y$$

**42.3.2 Bias of the 2SOLS estimator****42.3.3 Variance of the 2SOLS estimator****42.4 More****42.4.1 Identification through exogeneous variables**

Previously our structural model was:

$$Q = \alpha_1 + \beta_1 P + \epsilon_1$$

$$Q = \alpha_2 + \beta_2 P + \epsilon_2$$

And our reduced form:

$$P = \frac{\alpha_2 - \alpha_1}{\beta_1 - \beta_2} + \frac{\epsilon_2 - \epsilon_1}{\beta_1 - \beta_2}$$

$$Q = \frac{\alpha_2\beta_1 - \alpha_1\beta_2}{\beta_1 - \beta_2} + \frac{\beta_1\epsilon_2 - \beta_2\epsilon_1}{\beta_1 - \beta_2}$$

Or:

$$P = \pi_1 + \tau_1$$

$$Q = \pi_2 + \tau_2$$

**Adding another variable**

This time we add another measured variable,  $I$ .

$$Q = \alpha_1 + \beta_1 P + \theta_1 I + \epsilon_1$$

$$Q = \alpha_2 + \beta_2 P + \theta_2 I + \epsilon_2$$

The reduced form is now:

$$P = \frac{\alpha_2 - \alpha_1}{\beta_1 - \beta_2} + \frac{\theta_2 - \theta_1}{\beta_1 - \beta_2} I + \frac{\epsilon_2 - \epsilon_1}{\beta_1 - \beta_2}$$

$$Q = \frac{\alpha_2\beta_1 - \alpha_1\beta_2}{\beta_1 - \beta_2} + \frac{\theta_2\beta_1 - \theta_1\beta_2}{\beta_1 - \beta_2} I + \frac{\beta_1\epsilon_2 - \beta_2\epsilon_1}{\beta_1 - \beta_2}$$

Or:

$$P = \pi_{11} + \pi_{12} I + \tau_1$$

$$Q = \pi_{21} + \pi_{22} I + \tau_2$$

We can estimate  $\pi_1$  and  $\pi_2$  as  $\hat{\pi}_1$  and  $\hat{\pi}_2$  respectively.

We can now create estimators  $\hat{\pi}_{11}$ ,  $\hat{\pi}_{12}$ ,  $\hat{\pi}_{21}$  and  $\hat{\pi}_{22}$ .

**Identification with an exogeneous variable**

We now have 4 estimators and 6 parameters, meaning that we still cannot identify the model.

**Partial identification**

Can we use  $\hat{\pi}$  to identify any of the structural parameters?

We know that:

- $\pi_{11} = \frac{\alpha_2 - \alpha_1}{\beta_1 - \beta_2}$
- $\pi_{12} = \frac{\theta_2 - \theta_1}{\beta_1 - \beta_2}$
- $\pi_{21} = \frac{\alpha_2\beta_1 - \alpha_1\beta_2}{\beta_1 - \beta_2}$
- $\pi_{22} = \frac{\theta_2\beta_1 - \theta_1\beta_2}{\beta_1 - \beta_2}$

If the exogenous variable only affects one side of the equation, so  $\theta_1 = 0$ , we have:

- $\pi_{11} = \frac{\alpha_2 - \alpha_1}{\beta_1 - \beta_2}$
- $\pi_{12} = \frac{\theta_2}{\beta_1 - \beta_2}$
- $\pi_{21} = \frac{\alpha_2\beta_1 - \alpha_1\beta_2}{\beta_1 - \beta_2}$
- $\pi_{22} = \frac{\theta_2\beta_1}{\beta_1 - \beta_2}$

So we can see that:

$$\hat{\beta}_1 = \frac{\hat{\pi}_{22}}{\hat{\pi}_{12}}$$

This means we now have:

- $\pi_{11} = \frac{\pi_{12}(\alpha_2 - \alpha_1)}{\pi_{22} - \pi_{12}\beta_2}$
- $\pi_{12} = \frac{\pi_{12}\theta_2}{\pi_{22} - \pi_{12}\beta_2}$
- $\pi_{21} = \frac{\pi_{12}(\alpha_2\beta_1 - \alpha_1\beta_2)}{\pi_{22} - \pi_{12}\beta_2}$
- $\pi_{22} = \frac{\pi_{12}\theta_2\beta_1}{\pi_{22} - \pi_{12}\beta_2}$

We can use this to also identify  $\alpha_1$ .

### Complete identification

If we have independent variables for each of the two equations, we can fully identify the model.

We will have 6 estimators and 6 parameters.

We are estimating:

$$Q = \alpha_1 + \beta_1 P + \theta_1 I + \epsilon_1$$

$$Q = \alpha_2 + \beta_2 P + \theta_2 J + \epsilon_2$$

$I$  and  $J$  are essentially instrumental variables for the model.

$I$  is an instrumental variable for demand shocks, and  $J$  is an instrumental variable for supply shocks.

#### 42.4.2 Power of instruments

### 42.5 The Instrumental Variable (IV) estimator

#### 42.5.1 Instrumental Variable (IV) estimator

$$\theta_{IV} = (Z^T X)^{-1} Z^T y$$

2SOLS collapses to IV in some circumstances.

#### 42.5.2 Bias of the IV estimator

Equal to actual parameter so long as  $\epsilon$  uncorrelated with  $Z$ .

#### 42.5.3 Variance of the IV estimator

In OLS we had:

$$\theta_{OLS} = (X^T X)^{-1} X^T y$$

$$Var[\theta_{OLS}] = (X^T X)^{-1} X^T \Omega X (X^T X)^{-1}$$

With IV we have

$$\theta_{IV} = (Z^T X)^{-1} Z^T y$$

$$Var[\theta_{IV}] = (Z^T X)^{-1} Z^T \Omega Z (Z^T X)^{-1}$$

We can use weighted least squares for  $\Omega$ .

## 42.6 Choosing instrumental variables

### 42.6.1 Double selection

## 42.7 Other

### 42.7.1 Natural experiments

### 42.7.2 Non-linear models in the first stage

### 42.7.3 Random Effects Instrumental Variables (REIV)

### 42.7.4 Fixed Effects Instrumental Variables (FEIV)

### 42.7.5 SORT

synthetic IV indep on nuisance as alternative to matching.

IV: h3 on non-linear models for first stage

discontinuity

controlled experiments

two sources: missing data and simultaneous

variations in government rollouts, lotteries

IV may only affect subset of individuals

For example IV of draft number for military service. This only is an instrument for conscripts, not volunteers

generally, need to rationalise this and time series. There's stuff there on natural experiments etc

define confounding in IV? or in dependent variables? is different issue to the one of correlation with error?

h3 on Limited Information Maximum Likelihood

h3 on K-class estimation

Contrast loss and Siamese h3? One shot classification

IV: frame around parameter estimation when don't observe some variables. This can mean the direct variable can't be measured, or that some controls can't be measured

which factors to include? All?

page on structural and reduced forms

h3 on simultaneous equations there? Eg  $y = c_1 + \theta_1 X + \epsilon_1$   $y = c_2 + \theta_2 X + \rho Z \epsilon_2$

We can turn this into the reduced form:  $y = c_3 + \theta_3 Z + \epsilon_3$   $y = c_4 + \theta_4 Z + \epsilon_4$

difference between confounding and correlation with error?



## Chapter 43

# Imputing missing data for inference

### 43.1 Introduction

#### 43.1.1 Techniques for inference

Deleting whole row if missing data.

As with techniques for prediction, there is bias if not random.

## Chapter 44

# Measurement error and inference

### 44.1 Other

#### 44.1.1 Omitted variable bias

### 44.2 Measurement error

## Chapter 45

# Semi-parametric regression

### 45.1 The Robinson estimator

#### 45.1.1 Partially linear models

#### 45.1.2 The Robinson estimator

Partially out

$$y_i = \mathbf{x}_i\theta + g(\mathbf{z}_i) + \epsilon_i$$

Consider:

$$E(y_i|\mathbf{z}_i) = E(\mathbf{x}_i\theta + g(\mathbf{z}_i) + \epsilon_i|\mathbf{z}_i)$$

$$E(y_i|\mathbf{z}_i) = E(\mathbf{x}_i\theta|\mathbf{z}_i) + E(g(\mathbf{z}_i)|\mathbf{z}_i) + E(\epsilon_i|\mathbf{z}_i)$$

$$E(y_i|\mathbf{z}_i) = E(\mathbf{x}_i|\mathbf{z}_i)\theta + g(\mathbf{z}_i)$$

We can now remove the parametric part:

$$y_i - E(y_i|\mathbf{z}_i) = \mathbf{x}_i\theta + g(\mathbf{z}_i) + \epsilon_i - E(\mathbf{x}_i|\mathbf{z}_i)\theta - g(\mathbf{z}_i)$$

$$y_i - E(y_i|\mathbf{z}_i) = (\mathbf{x}_i - E(\mathbf{x}_i|\mathbf{z}_i))\theta + \epsilon_i$$

We define:

- $\bar{y}_i = y_i - E(y_i|\mathbf{z}_i)$
- $\bar{\mathbf{x}}_i = \mathbf{x}_i - E(\mathbf{x}_i|\mathbf{z}_i)$

$$\bar{y}_i = \bar{\mathbf{x}}_i\theta + \epsilon_i$$

**Estimating  $\bar{y}_i$  and  $\bar{\mathbf{x}}_i$**

So we can use OLS if we can estimate.

- $E(y_i | \mathbf{z}_i)$
- $E(\mathbf{x}_i | \mathbf{z}_i)$

We can do this with non-parametric methods.

### 45.1.3 Bias and variance of the Robinson estimator

robinson: can't have confounded in dummy. but can in real. general result of propensity stuff?

Framing: Partialling out is an alternative to OLS where  $n \ll p$  doesn't hold. alternative to LASSO etc

$$\hat{\theta} \approx N(\theta, V/n)$$

$$V = (E[\hat{D}^2])^{-1} E[\hat{D}^2 \epsilon^2] (E[\hat{D}^2])^{-1}$$

These are robust standard errors.

#### Moments of the Robinson estimator

If IID then

$$Var(\hat{\theta}) = \frac{\sigma_\epsilon^2}{\sum_i (x_i - \hat{X}_i)^2}$$

Otherwise, can use GLM

What are the properties of the estimator?

$$E[\hat{\theta}] = E\left[\frac{\sum_i (X_i - \hat{X}_i)(y_i - \hat{y}_i)}{\sum_i (x_i - \hat{X}_i)^2}\right]$$

### 45.1.4 Non-linear treatment effects in the Robinson estimator

Page on reformulating as non-linear. can do it. show can be estimated using arg min <https://arxiv.org/pdf/1712.04912.pdf>

### 45.1.5 DML

in DML. page on orthogonality scores, page on constructing them; page on using them to estimate parameters (GMM)

$$\text{We have } P(X) = f(\theta, \rho) \quad \hat{\theta} = f(X, n) \quad \theta = g(\rho, X)$$

$$\text{So error is: } \hat{\theta} - \theta = f(X, n) - g(\rho, X)$$

$$\text{Bias is defined as: } Bias(\hat{\theta}, \theta) = E[\hat{\theta} - \theta] = E[\hat{\theta}] - \theta \quad Bias = E[\hat{\theta} - \theta] = E[f(X, n) - g(\rho, X)] \quad Bias = E[\hat{\theta} - \theta] = E[f(X, n)] - g(\rho, X)$$

double ML: regression each parametric parameter on ML of other variables. eg: get  $e(x|z)$   $e(d|x)$   $d = m(x) + v$   $d$  is correlated with  $x$  so bias.  $v$  is correlated with  $d$  but not  $x$ . use as "iv". Still need estimate for  $g(x)$ .

for iterative, process is: + estimate  $g(x)$  + plug into other and estimate theta + this section should be in sample splitting. rename iterative estimation. separate pages for bias, variance + how does this work?? paper says random forest regression and OLS. initialise  $\theta$  randomly? + page on bias, variance, efficiency? + page on sample splitting, why?

+ page on goal:  $x$  and  $z$  orthogonal for split sampling + page on  $X = m_0(Z) + \mu$ , first stage machine learning, synthetic instrumental variables? h3 on that for multiple variables on interest. regression for each

### 45.1.6 DML1

Divide into  $k$ .

For each do ML on nuisance (how???) use all instances outside of sample

Then do GMM using orthogonality condition to calculate  $\theta$ . (how??) use instances in sample

Average  $\theta$  from each class

### 45.1.7 Last stage Robinson

Separate page for last stage: note we can do OLS, GLS etc with choice of  $\Omega$ .

## 45.2 Causal trees