Assignment Project Exam Help

https://powcoder.com

Add WeChat powcoder

# Introduction

## MAST90083 Computational Statistics and Data Mining

Dr Karim Seghouane

School of Mathematics & Statistics

The University of Melbourne

**Admin**
ooooooo

**Overview**
ooooooooooooo

**Basic concepts**
ooooo

## Outline

Assignment Project Exam Help

§i. Admin

https://powcoder.com

§ii. Introduction & overview

Add WeChat powcoder

§iii. Basic concepts

**Admin**
●○○○○○

Overview
○○○○○○○○○○○○

**Basic concepts**
○○○○○

## Admin

- Lectures - Dr Karim Seghouane
  - Mon. 14:15-16:15
  - Baldwin Spencer Theatre or via zoom
- Practical lab - Jiadong Mao
  - 8 Groups
  - f2f: G2/Tues. 9:00-10:00, G3/Tues. 14:15-15:15 in PAR-Peter Hall-G70 (Wilson Laboratory) G4/Fri. 10:00-11:00 in PAR-Peter Hall-G69 (Thompson Lab)
  - Online: G1/Wedn. 15:15-16:15, G5/Tues. 16:15-17:15, G6/Thur. 16:15-17:15, G7/Wedn. 16:15-17:15 and G8/Thur. 12:00-13:00
- Consultation time
  - Wedn. 08:00-09:00
  - Frid. 08:00-09:00

**Admin**
○ ● ○ ○ ○ ○

**Overview**
○ ○ ○ ○ ○ ○ ○ ○ ○ ○ ○ ○

**Basic concepts**
○ ○ ○ ○ ○

## Admin - Assessment

Assignment Project Exam Help

- Problem solving assignments - 45%

https://powcoder.com
  - Three assignments due early, mid and late semester
  - Each written assignment is worth 15%

- Exam. - 55%

Add WeChat powcoder

## Admin - LMS

- All relevant teaching material will be posted on LMS (including the supplementary and the additional material)
- Due to the time limits during the lectures, you will need to go over some mathematical details & deepen your knowledge outside of the lectures time.
- Your assignment must be submitted via LMS
- Discussion board?

**Admin**
○○○●○○

**Overview**
○○○○○○○○○○○○

**Basic concepts**
○○○○○

## Admin - References

- **Elements of Statistical Learning** by Hastie Trevor, Tibshirani Robert & Friedman Jerome (2009).

- **An Introduction to Statistical Learning** by James Gareth, Witten Daniela, Hastie Trevor & Tibshirani Robert (2013).

- **Introducing Monte Carlo Methods** with R by R.P. Christian & G. Casella (2010).

- **Computational Statistics** by G.H. Givens & J.A. Hoeting (2005).

- Academic **articles**, links to **blogs** & **videos** on LMS.

**Admin**
○○○○●○

**Overview**
○○○○○○○○○○○○

**Basic concepts**
○○○○○

## Admin - Communication

- Note that email regarding the course material, laboratories and assignments will be addressed during office hours and if time permits.

- It is expected that questions regarding these matters will be asked during consultation hours or during the laboratories.

- There will be no consultation hours during non-teaching periods.

- I will be out of office during SWOTVac week and the first two weeks of the examination period.

- I will provide an extra consultation time on week 12.

- Students should plan ahead with any queries regarding assignments or course material.

**Admin**
○○○○○●

**Overview**
○○○○○○○○○○○○

**Basic concepts**
○○○○○

## Lecture schedule (provisional)

- **Data mining** ( 7 weeks ):
  - linear model selection and regularisation; kernel and local regression; basis expansion and spline regression; general additive models (GAM); classification and regression trees; bagging, random forests and boosting; support vector machines (SVM); component analysis and deep learning.
- **Computational statistics** ( 4 weeks ):
  - EM algorithm; Bayes computing; Monte Carlo methods; and bootstrap methods.

**Admin**
○○○○○○

**Overview**
●○○○○○○○○○○○

**Basic concepts**
○○○○○

## Question to the class...

How would you explain in a few sentences to a general audience

► What a statistician / data scientist does?

► What is the main purpose of that?

**Admin**
○○○○○○

**Overview**
○●○○○○○○○○○○○
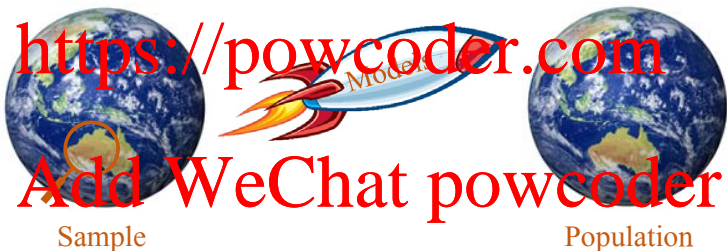
**Basic concepts**
○○○○○

## Some answers

Assignment Project Exam Help

- ► What a statistician / data scientist does?
- ► Creat a model (box) to understand the releathinseep between several variable

https://powcoder.com

- ► make sense of data, extract important patterns and trends, we can call this learning from data

- ► What is the main purpose of that?

Add WeChat powcoder

- ► Predict/decide or describe/understand

## The power of (statistical) model



Sample

Population

Admin
○○○○○○

Overview
○○○●○○○○○○○○

Basic concepts
○○○○○

## Once upon a time …



Image source: http://freshlearners.blogspot.com.au/2015/07/most-recent-communications-technology.html

**Admin**
○○○○○○

**Overview**
○○○○●○○○○○○○

**Basic concepts**
○○○○○

## Today

The real power is in KNOWING WHAT TO DO with the data



Image source: http://psmit.com/about.html

Admin
○○○○○○

Overview
○○○○○○●○○○○○○

Basic concepts
○○○○○

# Taxonomy of covered methods

**Admin**
○○○○○○

**Overview**
○○○○○○○●○○○○○○

**Basic concepts**
○○○○○

## Regression vs. Clasification

Assignment Project Exam Help

https://powcoder.com

- ▶ **Regression** - The aim is to predict the numerical outcome of a subject based on its features.
- ▶ **Clasification** - The aim is to predict the class belonging of a subject based on its features.

Add WeChat powcoder

**Admin**
○○○○○○

**Overview**
○○○○○○○○●○○○○

**Basic concepts**
○○○○○

## Parametric vs. Nonparametric

- **Parametric** approach - Makes an explicit assumption about the functional form. (Restriction on the shape)
- **Nonparametric** approach - does not assume any functional form for the underlying model structure.

Admin
○○○○○○

Overview
○○○○○○○○○●○○○○

Basic concepts
○○○○○

## Parametric models

### Advantages
- Simple to understand and interpret.
- Fast to fit (CPU time)
- Requires to fit a small number of parameters, therefore requires relatively small data sets.

### Disadvantages
- Highly constrained to the specified function form.
- It is hard to match the mathematical form of the model to the unknown 'DGP'.
- Can fit well relatively simple data structures (complex models -> danger of over-fitting).

## Nonparametric models

### Advantages
- Very flexible
- Weak(er) assumptions about the underlying function.
- Outperforms in prediction context, mostly.

### Disadvantages
- Slow(er) to fit (CPU time).
- Requires large(r) data sets
- Harder to interpret
- Higher danger of over-fitting

Admin
○○○○○○

Overview
○○○○○○○○○○●○○

Basic concepts
○○○○○

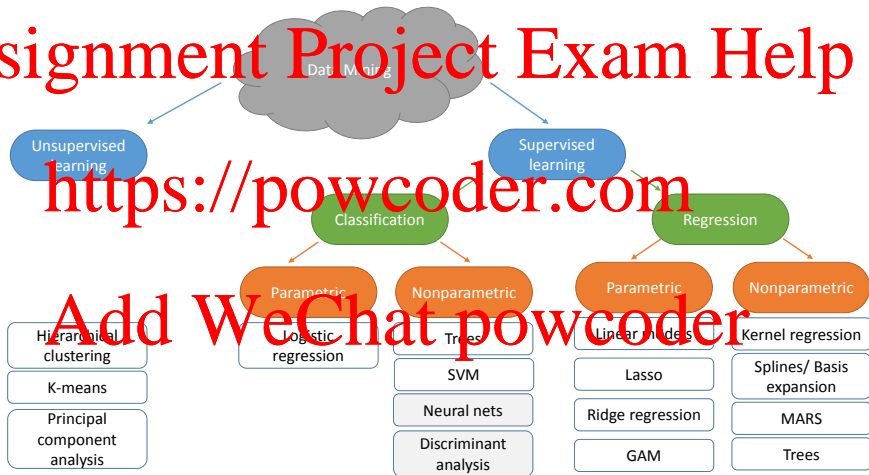## Supervised vs. Unsupervised

Assignment Project Exam Help

- ▶ **Supervised** learning - In case our data set contains the response (outcome) measurements, the fitted model relates the different features to the response.

https://powcoder.com

- ▶ **Unsupervised** learning - In case our data set contains only information about the different features of the subjects, the fitted model aims to segment the subjects into groups or just learn about the relationship of the features.

Add WeChat powcoder

Admin
○○○○○○

Overview
○○○○○○○○○○○●○

Basic concepts
○○○○○

# Taxonomy of covered methods

**Admin**
ooooooo

**Overview**
oooooooooooooo

**Basic concepts**
●oooo

## How models are fitted

Assignment Project Exam Help

https://powcoder.com

- ▶ By optimization
- ▶ We minimize some criterion (e.g. squared error with some additional penalty)

Add WeChat powcoder

Admin
oooooo

Overview
oooooooooooo

Basic concepts
o●oooo

## Trade-offs

Assignment Project Exam Help

▶ Flexibility (predictability) - Interpretability trade-off

https://powcoder.com

▶ Bias - Variance trade-off.
  ▶ variance - how much $\hat{f}$ changes with different data set.
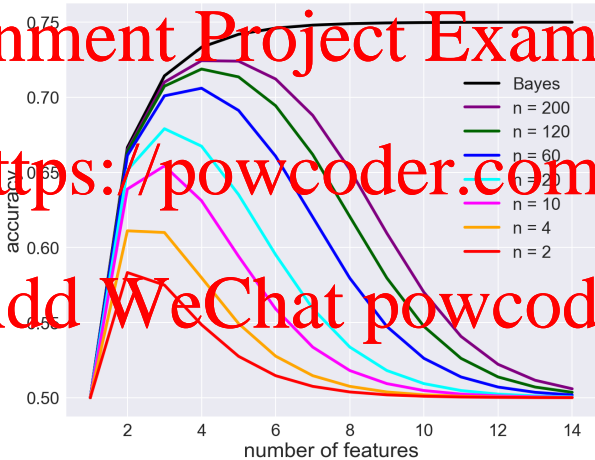  ▶ bias - the error from the approximation.

Add WeChat powcoder

# Complexity

► Complexity trade-offs involving sample size, dimensionality and empirical performance. It is a characteristic feature of supervised learning methods.

► Notion of *curse of dimensionality*: for a fixed sample size, the expected classification error will improve by increasing the number of features, but eventually will decrease. This is a consequence of the large size of high-dimensional spaces, which require correspondingly large training sample sizes.

► *Scissors effect*: the expected error typically decreases as sample size increases, and more complex classification rules achieve smaller error for large sample sizes; however, simpler classification rules can perform better under small sample sizes, by virtue of needing less data.

## Complexity



Expected accuracy in a discrete classification problem for various training sample sizes as a function of the number of predictors.
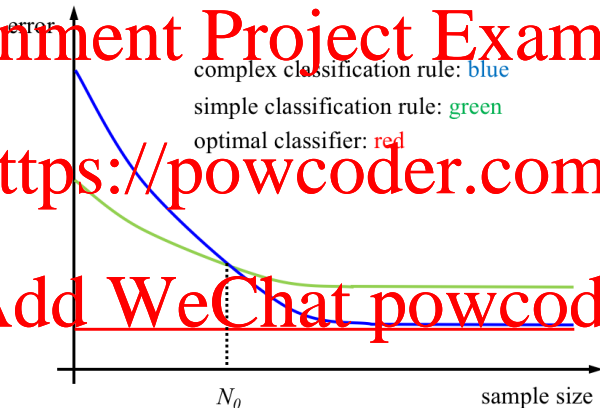
Admin
ooooooo

Overview
oooooooooooooo

Basic concepts
ooooo●

## Complexity



error

complex classification rule: blue

simple classification rule: green

optimal classifier: red

$N_0$

sample size

Expected error as a function of sample size for two classification rules. There is a problem-dependent critical sample size N0, under which one should use the simpler classification rule.