



MONASH University

ETC3250

Business Analytics

Week 2.
Statistical learning

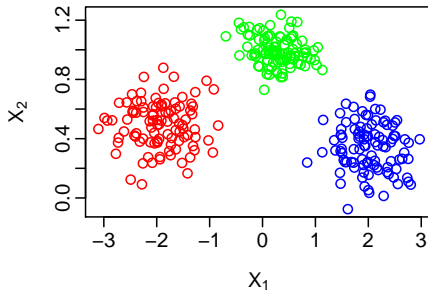
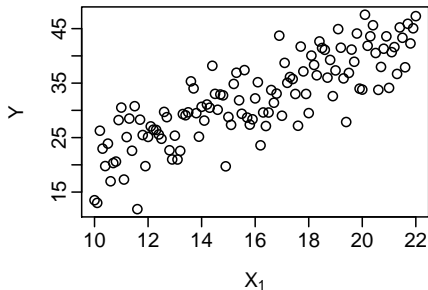
3 August 2015

Learning from data

- **Better understand** or **make predictions** about a certain phenomenon under study
- **Construct a model** of that phenomenon by finding relations between **several variables**
- If phenomenon is complex or depends on a large number of variables, an **analytical solution** might not be available
- However, we can **collect data** and learn a model that **approximates** the true underlying phenomenon

Learning from a dataset

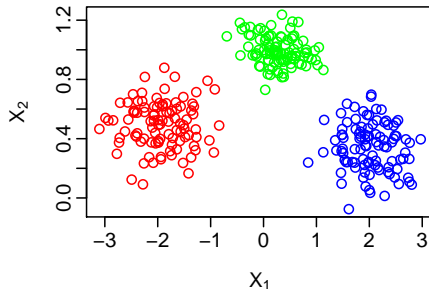
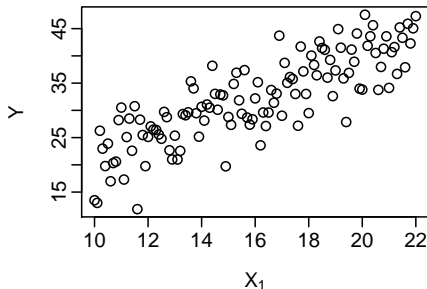
$$\mathcal{D} = \{(x_i, y_i)\}_{i=1}^N \text{ with } x_i = (x_{i1}, \dots, x_{ip})^T$$



Statistical learning provides a framework for constructing models from \mathcal{D} .

Learning from a dataset

$$\mathcal{D} = \{(x_i, y_i)\}_{i=1}^N \text{ with } x_i = (x_{i1}, \dots, x_{ip})^T$$



Statistical learning provides a framework for constructing models from \mathcal{D} .

Different learning problems

- Supervised learning
 - Regression (or prediction)
 - Classification

→ y_i **available for all** x_i
- Unsupervised learning

→ y_i **unavailable for all** x_i
- Semi-supervised learning

→ y_i **available only for few** x_i
- Other types of learning: reinforcement learning, online learning, active learning, etc.

Identification of the best learning problem is important in practice

What is Statistical Learning?

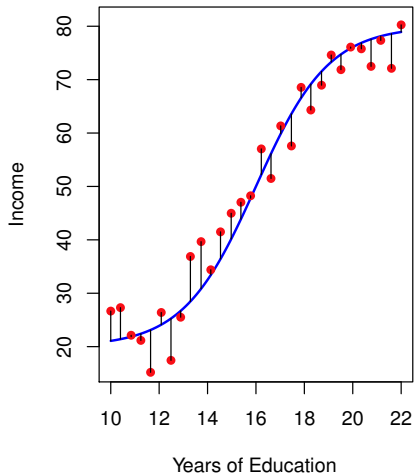
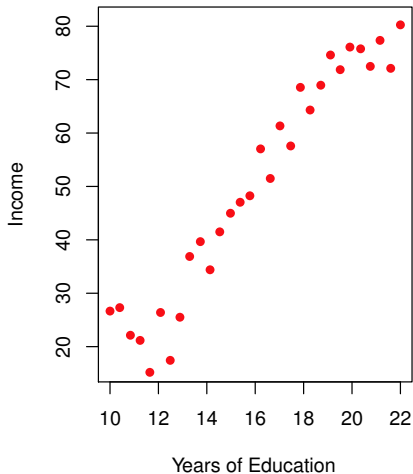
$$\mathcal{D} = \{(x_i, y_i)\}_{i=1}^N$$

$$Y = f(\overbrace{X_1, \dots, X_p}^X) + \varepsilon$$

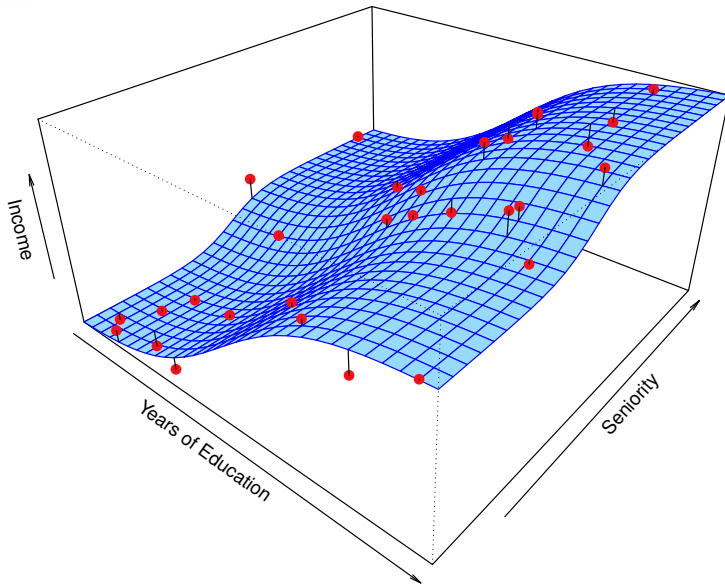
- Y : response (output)
- f : unknown function
- X : set of p predictors (inputs)
- ε : error term

Learn (or estimate) the function f using \mathcal{D}

What is Statistical Learning?



What is Statistical Learning?



Why estimate f ?

- Prediction: $\hat{Y} = \hat{f}(X)$

$$\begin{aligned} E[(Y - \hat{Y})^2] &= E[(f(X) + \varepsilon - \hat{Y})^2] \\ &= \underbrace{E[(f(X) - \hat{f}(X))^2]}_{\text{Reducible}} + \underbrace{\text{Var}(\varepsilon)}_{\text{Irreducible}} \end{aligned}$$

- Inference (or explanation):
 - Which predictors are associated with the response?
 - What is the relationship between the response and each predictor?

How do we estimate f ?

■ Parametric methods

- Assumption about the form of f , e.g. linear:
 $f(X) = \beta_0 + \beta_1 X_1 + \dots + \beta_p X_p$ and $\hat{Y}(x) = \hat{f}(x)$
- 😊 The problem of estimating f reduces to estimating a set of parameters
- 😊 Usually a good starting point for many learning problems
- 😞 Poor performance if linearity assumption is wrong

■ Non-parametric methods

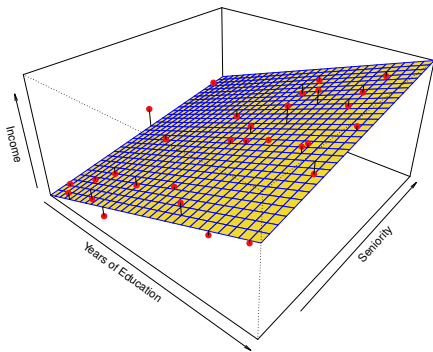
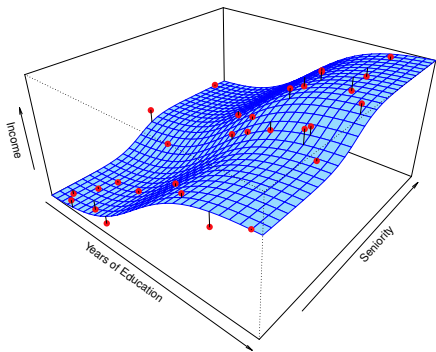
- No *explicit* assumptions about the form of f , e.g. nearest neighbours: $\hat{Y}(x) = \frac{1}{k} \sum_{x_i \in N_k(x)} y_i$
- 😊 High flexibility: it can potentially fit a wider range of shapes for f
- 😞 A large number of observations is required to estimate f with good accuracy

How do we estimate f ?

For each of parts (a) through (d), indicate whether we would generally expect the performance of a **flexible** statistical learning method to be better or worse than an **inflexible** method.

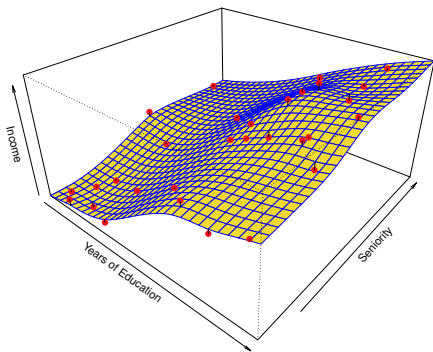
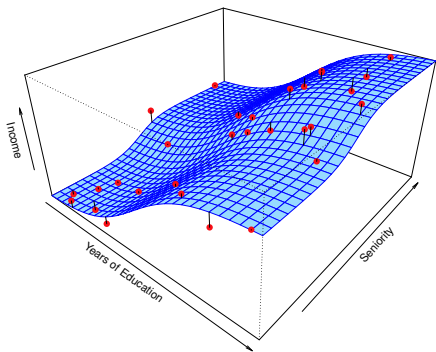
- The sample size n is extremely large, and the number of predictors p is small.
- The number of predictors p is extremely large, and the number of observations n is small.
- The relationship between the predictors and response is highly non-linear.
- The variance of the error terms, i.e. $\sigma^2 = \text{Var}(\varepsilon)$ is extremely high.

How do we estimate f ?

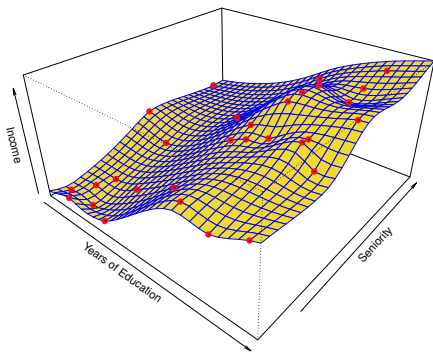
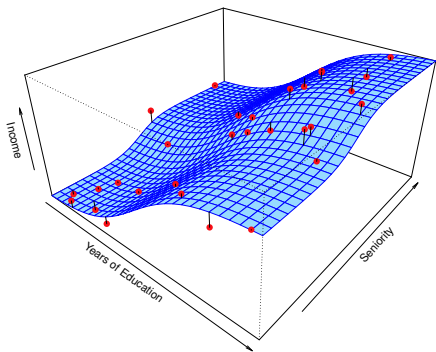


$$\hat{f}(\text{education}, \text{seniority}) = \hat{\beta}_0 + \hat{\beta}_1 \times \text{education} + \hat{\beta}_2 \times \text{seniority}$$

How do we estimate f ?



How do we estimate f ?



Prediction Accuracy vs Model Interpretability

