#### **Lecture 1: Introduction**

Felix Held, Mathematical Sciences

MSA220/MVE441 Statistical Learning for Big Data

22<sup>nd</sup> March 2021



# What is Big Data?

.

#### Just a buzz word?



#### Cancer treatment is on the brink of a data revolution

Lydia Ramsey Sep. 22, 2015, 4:29 PM



Business Insider<sup>1</sup>

https://www.businessinsider.com/big-data-and-cancer-2015-9?r=US&IR=T&IR=T

**Big Data - Big Problems?** 

### **BIG DATA**

# The Parable of Google Flu: Traps in Big Data Analysis

David Lazer, 1,2\* Ryan Kennedy, 1,3,4 Gary King, 3 Alessandro Vespignani 5,6,3

Scientific discussion article<sup>1</sup>

<sup>&</sup>lt;sup>1</sup> Lazer et al. (2014) The Parable of Google Flu: Traps in Big Data Analysis. Science 343 (6176):1203–1205. DOI 10.1126/science.1248506

#### **Big Data - Big Problems?**



#### Financial Times<sup>1</sup>

#### The New Hork Times

# **How Democracy Can Survive Big Data**

New York Times<sup>2</sup>

<sup>1</sup> https://www.ft.com/content/21a6e7d8-b479-11e3-a09a-00144feabdc0#axzz2v0200f0X

<sup>2</sup> https://www.nvtimes.com/2018/03/22/opinion/democracy-survive-data.html



#### Yes and no.

Note that *size* is a flexible term. Here mostly:

Size as in: Number of observations

#### **Big-***n* **setting**

Size as in: Number of variables

#### $\mathbf{Big}$ -p setting

Size as in: Number of observations and variables
Big-n / Big-p setting

Is this all?

#### The Four Vs of Big Data



http://web.archive.org/web/20210506042232/https://www.ibmbigdatahub.com/infographic/four-vs-big-data

Statistics as a science has always been concerned with...

- experimental design or 'how to collect the data'
- modelling of data and underlying assumptions
- ► inference of parameters
- uncertainty quantification in estimated parameters/predictions

Focus is on the last three in this course.

- Increase in sample size often leads to increase in complexity and variety of data (p grows with n)
- More data  $\neq$  less uncertainty
- A lot of classical theory is for fixed p and growing n
- Exploration and visualisation of Big Data can already require statistics
- Probability of extreme values: Unlikely results become much more likely with an increase in n
- Curse of dimensionality: Lot's of space between data points in high-dimensional space

# **Statistical Learning**

#### **Basics about random variables**

- ► We will consider **discrete** and **continuous** random quantities
- ▶ Probability mass function (pmf) p(k) for a discrete variable **Example:** Bernoulli distribution with parameter  $\theta \in (0, 1)$

 $p(0) = \theta, \quad p(1) = 1 - \theta$ 

► Probability density function (pdf)  $p(\mathbf{x})$  for a continuous variables Example: Multivariate normal distribution with mean vector  $\boldsymbol{\mu} \in \mathbb{R}^p$  and covariance matrix  $\boldsymbol{\Sigma} \in \mathbb{R}^{p \times p}$ 

$$p(\mathbf{x}) = |2\pi\Sigma|^{-1/2} \exp\left(-\frac{1}{2}(\mathbf{x}-\boldsymbol{\mu})^{\mathsf{T}}\Sigma^{-1}(\mathbf{x}-\boldsymbol{\mu})\right)$$

#### Two important rules (and a consequence)

#### Marginalisation

For a joint density p(x, y) it holds that

$$p(x) = \sum_{y} p(x, y)$$
 or  $p(x) = \int p(x, y) dy$ 

#### Conditioning

For a joint density p(x, y) it holds that

$$p(x, y) = p(x|y)p(y) = p(y|x)p(x)$$

Both rules together imply Bayes' law

$$p(x|y) = \frac{p(y|x)p(x)}{p(y)}$$

11/29

#### Expectations and variance depend on an underlying pdf/pmf.

#### Notation:

• 
$$\mathbb{E}_{p(x)}[f(x)] = \int f(x)p(x) dx$$
  
•  $\operatorname{Var}_{p(x)}[f(x)] = \mathbb{E}_{p(x)}\left[\left(f(x) - \mathbb{E}_{p(x)}[f(x)]\right)^2\right]$ 

Learn **a model** from **data** by minimizing **expected prediction error** determined by a loss function.

- Model: Find a model that is suitable for the data
- > Data: Data with known outcomes is needed
- Expected prediction error: Focus on quality of prediction (predictive modelling)
- Loss function: Quantifies the discrepancy between observed data and predictions

#### Linear regression - An old friend



14/29

#### **Statistical Learning and Linear Regression**

**Data:** Training data consists of independent pairs

 $(y_i, \mathbf{x}_i), \quad i = 1, \dots, n$ 

Observed response  $y_i \in \mathbb{R}$  for predictors  $\mathbf{x}_i \in \mathbb{R}^p$ 

Model:

$$y_i = \mathbf{x}_i^{\mathsf{T}} \boldsymbol{\beta} + \varepsilon_i$$

where  $\varepsilon_i \sim N(0, \sigma^2)$  independent

Loss function: Squared error loss

 $L(y, \hat{y}) = (y - \hat{y})^2$ 

#### Statistical decision theory for regression (I)

Squared error loss between outcome y and a prediction f(x) dependent on the variable(s) x

$$L(y, f(\mathbf{x})) = (y - f(\mathbf{x}))^2$$

- Assume we want to find the 'best' f that can be learned from training data
- When a new pair of data (y, x) from the same distribution (population) as the training data arrives, expected prediction loss for a given f is

$$J(f) = \mathbb{E}_{p(\mathbf{x}, y)} \left[ L(y, f(\mathbf{x})) \right] = \mathbb{E}_{p(\mathbf{x})} \left[ \mathbb{E}_{p(y|\mathbf{x})} \left[ L(y, f(\mathbf{x})) \right] \right]$$

Define 'best' by:

$$\widehat{f} = \underset{f}{\operatorname{arg\,min}} J(f)$$

#### Statistical decision theory for regression (II)

Can we determine  $\widehat{f}$ ? Focus on inner expectation

$$\begin{split} \mathbb{E}_{p(y|\mathbf{x})}\left[(y-f(\mathbf{x}))^2\right] &= \int (y-\mathbb{E}_{p(y|\mathbf{x})}[y] + \mathbb{E}_{p(y|\mathbf{x})}[y] - f(\mathbf{x}))^2 p(y|\mathbf{x}) \,\mathrm{d}y \\ &= \int (y-\mathbb{E}_{p(y|\mathbf{x})}[y])^2 p(y|\mathbf{x}) \,\mathrm{d}y \\ &+ 2 \int (y-\mathbb{E}_{p(y|\mathbf{x})}[y]) (\mathbb{E}_{p(y|\mathbf{x})}[y] - f(\mathbf{x})) p(y|\mathbf{x}) \,\mathrm{d}y \\ &+ \int (\mathbb{E}_{p(y|\mathbf{x})}[y] - f(\mathbf{x}))^2 p(y|\mathbf{x}) \,\mathrm{d}y \\ &= \operatorname{Var}_{p(y|\mathbf{x})}[y] + (\mathbb{E}_{p(y|\mathbf{x})}[y] - f(\mathbf{x}))^2 \end{split}$$

Minimal for  $f(\mathbf{x}) = \mathbb{E}_{p(y|\mathbf{x})}[y]$ 

#### Statistical decision theory for regression (III)

We just derived that

$$\widehat{f}(\mathbf{x}) = \mathbb{E}_{p(y|\mathbf{x})}[y]$$

the expectation of y given that  $\mathbf{x}$  is fixed (conditional mean)

- Regression methods approximate the conditional mean
- For many observations y with identical  $\mathbf{x}$  we could use

$$\mathbb{E}_{p(y|\mathbf{x})}[y] \approx \frac{1}{|\{y_i : \mathbf{x}_i = \mathbf{x}\}|} \sum_{\mathbf{x}_i = \mathbf{x}} y_i$$

▶ Probably more realistic to look for the *k* closest neighbours of **x** in the training data  $N_k(\mathbf{x}) = {\mathbf{x}_{i_1}, ..., \mathbf{x}_{i_k}}$ . Then

$$\mathbb{E}_{p(y|\mathbf{x})}[y] \approx \frac{1}{k} \sum_{\mathbf{x}_{i_l} \in N_k(\mathbf{x})} y_{i_l}$$

#### Average of *k* neighbours



k — 2 — 5

Linear regression is a **model-based approach** and assumes that the dependence of y on x can be written as a weighted sum, i.e.

 $y = \mathbf{x}^{\mathsf{T}} \boldsymbol{\beta} + \boldsymbol{\varepsilon}$ 

where  $\varepsilon \sim N(0, \sigma^2)$ . The mean of y given **x** is therefore

 $\mathbb{E}_{p(y|x)}[y] = \mathbf{x}^{\mathsf{T}} \boldsymbol{\beta}.$ 

Note that in practice this equality will only hold approximately.

## Classification

#### A simple example of classification



How do we classify a pair of new coordinates  $\mathbf{x} = (x_1, x_2)$ ?

▶ Find the *k* predictors

$$N_k(\mathbf{x}) = \{\mathbf{x}_{i_1}, \dots, \mathbf{x}_{i_k}\}$$

in the training sample, that are closest to  ${\bf x}$  in the Euclidean norm.

Majority vote: Assign x to the class that most predictors in N<sub>k</sub>(x) belong to (highest frequency)

#### kNN and its decision boundaries



#### Classification

Learn a rule  $c(\mathbf{x})$  from data which maps observed features  $\mathbf{x}$  to classes  $\{1, \dots, K\}$ .

#### **Remember:**

#### **Statistical Learning**

Learn a model from data by minimizing expected prediction error determined by a loss function.

Here: rule  $\simeq$  model, and observed classes give us the required outcomes for learning. What is a suitable loss?

#### Statistical decision theory for classification

0-1 misclassification loss: Let i be the actual class of an object and c(x) is a rule that returns the class for the variable(s) x, then

$$L(i, c(\mathbf{x})) = \begin{cases} 0 & i = c(\mathbf{x}), \\ 1 & i \neq c(\mathbf{x}) \end{cases} = \mathbb{1}(i \neq c(\mathbf{x}))$$

Expected prediction error

$$J(c) = \mathbb{E}_{p(\mathbf{x})} \left[ \mathbb{E}_{p(i|\mathbf{x})} [\mathbb{1}(i \neq c(\mathbf{x}))] \right]$$

Minimizing expected prediction error leads to the rule

$$\hat{c}(\mathbf{x}) = \underset{1 \le i \le K}{\arg \max} p(i|\mathbf{x})$$

This is called **Bayes' rule**.

Again, focus on inner expectation

$$\mathbb{E}_{p(i|\mathbf{x})}[\mathbbm{1}(i \neq c(\mathbf{x}))] = \sum_{i=1}^{K} \mathbbm{1}(i \neq c(\mathbf{x}))p(i|\mathbf{x})$$
$$= \sum_{i \neq c(\mathbf{x})} p(i|\mathbf{x})$$
$$= 1 - p(c(\mathbf{x})|\mathbf{x})$$

Minimal for  $\hat{c}(\mathbf{x}) = \arg \max_{1 \le i \le K} p(i|\mathbf{x})$ 

- kNN solves the classification problem by approximating p(i|x) with the frequency of class i among the k closest neighbours of x.
- Given data  $(i_l, \mathbf{x}_l)$  for l = 1, ..., n it holds that

$$\hat{c}(\mathbf{x}) = \operatorname*{arg\,max}_{1 \le i \le K} \frac{1}{k} \sum_{\mathbf{x}_l \in N_k(\mathbf{x})} \mathbb{1}(i_l = i)$$

There are two choices to make when implementing a kNN method

- 1. The metric to determine a neighbourhood
  - ▶ e.g. Euclidean/ $\ell_2$  norm, Manhattan/ $\ell_1$  norm, max norm, ...
- 2. The number of neighbours, i.e. k

The choice of metric changes the underlying local model of the method while k determines the size of this local model.

- Big Data is complex and is multi-faceted
- Regression and classification can be formulated in the framework of Statistical Learning
- In both cases, focus is on prediction