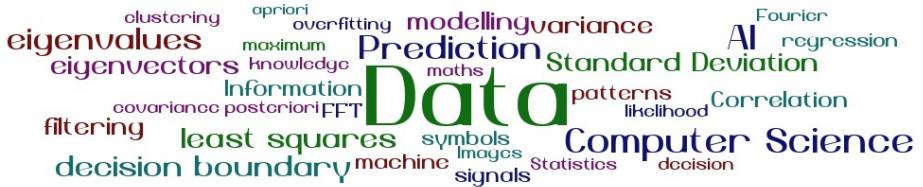# COMS20011 – Data-Driven Computer Science



February 2022
**Majid Mirmehdi**

Some slides in this lecture are adapted from those
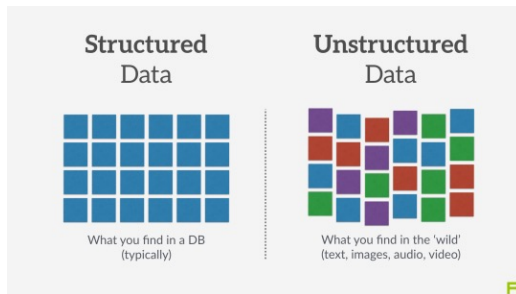authored by **Dima Damen** and **Andrew Calway**

**Lecture Video #1**

# COMS20011 Unit

- ➢ This is a "new" unit that started in the 2020-21 academic year

- ➢ Replaced the 20CP COMS20212 (SPS) unit

- ➢ Exam materials can be used for revision BUT...

- ➢ Use SPS materials with caution...depth, breadth & requirements may differ.

# What is Data?

➤ Data comes in many forms, e.g. symbols, patterns and signals!

➤ Data: *Structured and Unstructured*
  ➤ Numeric (measurements, finance spreadsheets, ...)
  ➤ Textual (emails, social media, web pages, medical records, ...)
  ➤ Visual (images, video, graphics, animations)
  ➤ Auditory (speech, audio)
  ➤ Signals (GPS signals, accelerometer, heart rate, ...)
  ➤ Many others...



Structured Data — What you find in a DB (typically)

Unstructured Data — What you find in the 'wild' (text, images, audio, video)

Image from Garrett Hollander

# This Unit (adapted from COMS20212: Symbols, Patterns and Symbols)

➤ This unit is about doing things with data... *but not*
  - ➤ storing, shuffling, searching (Algorithms I & II)
  - ➤ sending (Computer Systems)
  - ➤ compressing or encrypting (Cryptology)

➤ This unit is about:
  - ➤ extracting knowledge from data
  - ➤ generating data and making predictions
  - ➤ making decisions based on data
  - ➤ Often referred to as:



**DATA SCIENCE**

ANALYSIS STRUCTURE ALGORITHM PROCESS PROGRAMMING SOLVING KNOWLEDGE

Image from https://www.datanami.com

A DAY IN DATA

The exponential growth of data is undisputed, but the numbers behind this explosion – fuelled by internet of things and the use of connected devices – are hard to comprehend, particularly when looked at in the context of one day

Image from Raconteur https://www.weforum.org

5

## Data is the new Oil

**The Largest Companies By Market Cap**

| Year | 1st | 2nd | 3rd | 4th | 5th |
|------|-----|-----|-----|-----|-----|
| 2001 | GE $406B | Microsoft $365B | Exxon $272B | Citi $261B | Walmart $260B |
| 2006 | Exxon $446B | GE $383B | Total $327B | Microsoft $293B | Citi $273B |
| 2011 | Exxon $406B | Apple $376B | PetroChina $277B | Shell $237B | (ICBC) $228B |
| 2016 | Apple $582B | Alphabet $556B | Microsoft $452B | amazon $364B | f $359B |

# Data Science & Analytics



Google Trends: Interest In Data Jobs Over a Decade

Legend: "data scientist jobs" — "data engineering jobs" — "data analytics jobs"

https://onlinedatasciencemasters.virginia.edu/blog/data-science-vs-data-engineering/

# But it's not about the data – it's about the science

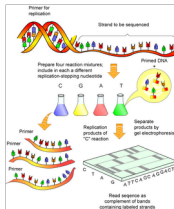Tracking and predicting [disease,mortality,floods,fires, and fun etc.] by Twitter!



## Big Data in Healthcare

CDC-Reported AHD Mortality     Twitter-Predicted AHD Mortality

6/09 – 3/10
826 million tweets
146 million geo-located
Across 1300 counties
Eval for stress & hostility:
     health, attractiveness,
     job, curse words
+CDC, US Census data

Eichstaedt JC et al.,
Psychological Science,
2015;26:159

10  20  30  40  50  60  70  80  90
AHD Mortality (Percentile)

# This Unit

Why is it important for Computer Science?

- ➤ Fundamental to many application areas:
    - ➤ Artificial Intelligence, Machine Learning, Deep Learning
    - ➤ Image Processing and Pattern Recognition
    - ➤ Graphics, Animation and Virtual Reality
    - ➤ Computer Vision and Robotics
    - ➤ Speech and Audio Processing.
    - ➤ With growing applications in: neuroscience, literature, agriculture, etc.
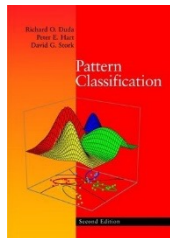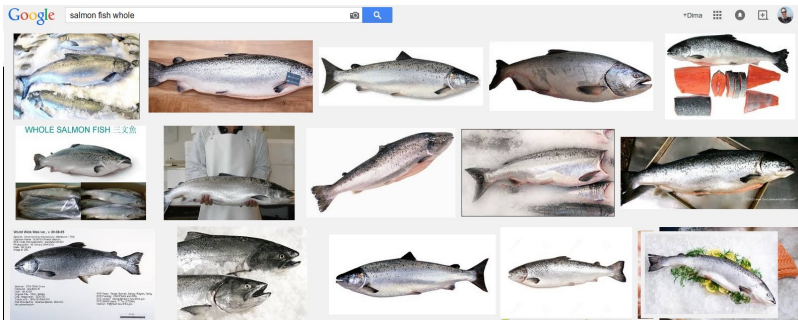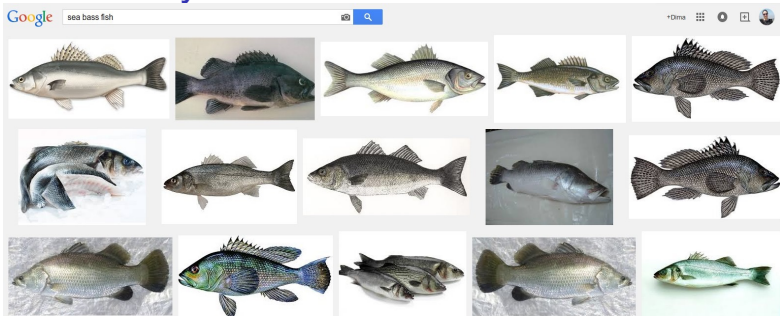- ➤ Hence, preparation for application units in years 3 and 4.

# Ex1. A Fishy Problem



**Data:** images of fish

**Aim:** distinguish between sea bass and salmon

From: Pattern Classification by *Duda, Hart and Stork*,
2nd Edition, Wiley Interscience

# Ex1. A Fishy Problem

# Fishing for a Solution

Steps:

1. Pre-processing
2. Feature Selection
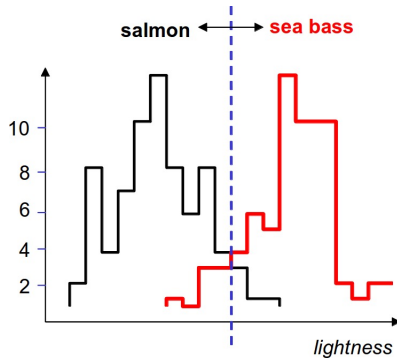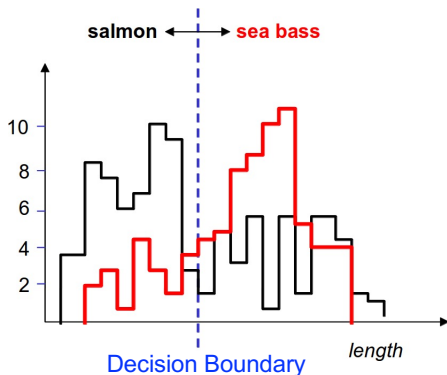3. Classification

# Fishing for a Solution

Steps:

1. Pre-processing e.g. Rotate and align, Segment fish from background
2. Feature Selection e.g. Measure length
3. Classification e.g. Find a threshold
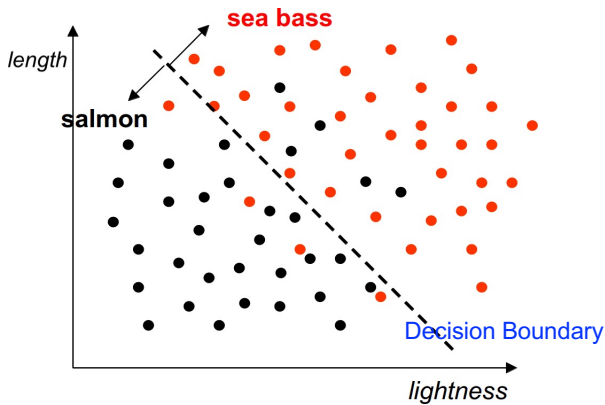
# Fishing for a Solution

Steps:
1. Pre-processing e.g. Rotate and align, Segment fish from background
2. Feature Selection e.g. Measure length or lightness
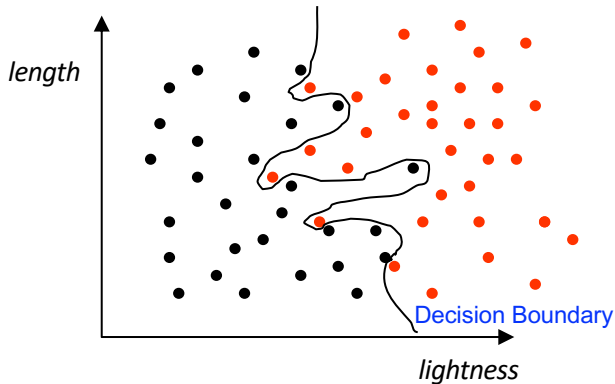3. Classification e.g. Find a threshold

# Fishing for a Solution

Multiple features could be selected, resulting in a multi-dimensional feature vector.

# Fishing for a Solution

Complex decision model



*length*

Decision Boundary

*lightness*

# Typical Data Analysis Problem

Steps:
1. Pre-processing  [Unit - Part 1] → Majid Mirmehdi (~10%)
2. Feature Selection  [Unit - Part 3] → Majid Mirmehdi (~40%)
3. Modelling & Classification  [Unit - Part 2] → Laurence Aitchison **[UD]** (~50%)

# Next Video

More example applications...