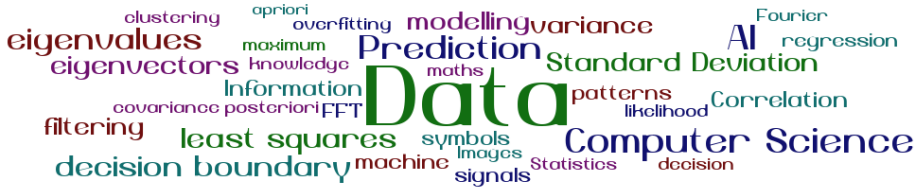


COMS20011 – Data-Driven Computer Science



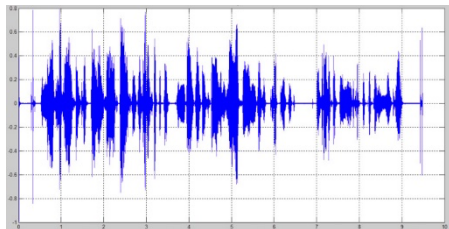
February 2021

Majid Mirmehdi

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

Lecture Video #2

Ex2. Speech Recognition



Data: Analogue speech signals (time series numerical data)

Aim: Convert audio into text (think Echo/Siri...)

1. Pre-processing Digitisation
2. Feature Selection Wave amplitude, frequencies
3. Inference Hidden Markov Models (Viterbi algorithm) [or Deep learning]

Ex3. Spam Filter

Data: Email texts

Aim: Determine whether the email is spam



1. Pre-processing - Normalise words
2. Feature Selection - Presence of words
3. Classification - Naive Bayes classifier

Select subset of words w_i and determine $P(w_i | spam)$ and $P(w_i | \neg spam)$ from frequencies in training data.

For an Email that contains w_1, w_2, \dots, w_n of the subset of words, assume

$$P(email | spam) = P(w_1 | spam)P(w_2 | spam) \dots P(w_n | spam) \quad (1)$$

and

$$P(email | \neg spam) = P(w_1 | \neg spam)P(w_2 | \neg spam) \dots P(w_n | \neg spam) \quad (2)$$

A new Email is spam if

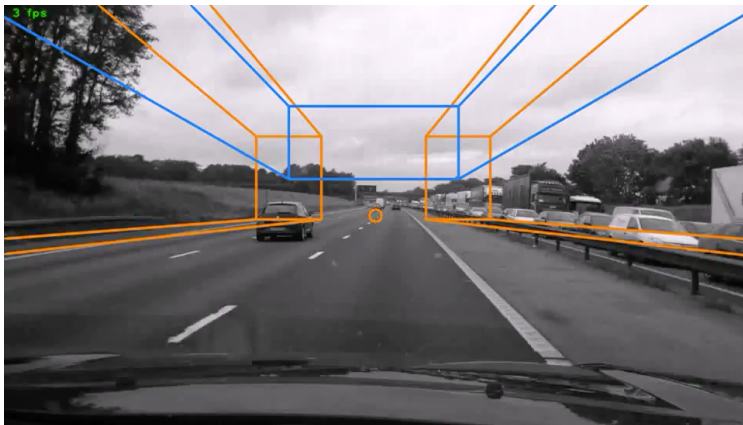
$$P(email | spam) > P(email | \neg spam) \quad (3)$$

Ex4.1 – Towards Autonomous Driving

Data: Video

Aim: Determine knowledge from the road or inside the vehicle

1. Pre-processing (Detect vanishing point)
2. Feature Selection (Use constraints to reduce number and dimensionality)
3. Recognition (Perspective transformations and OCR)



Ex4.2 – Towards Autonomous Driving

1. Pre-processing (Detect vanishing point)
2. Feature Selection (Straight lines)
3. Model Building (Detecting, predicting, decision making)



Ex4.3 – Towards Autonomous Driving

1. Pre-processing (Detect vanishing point)
2. Feature Selection (MSERs, Histogram of Gradients)
3. Classification (Support Vector Machines)



Ex4.4 – Towards Autonomous Driving

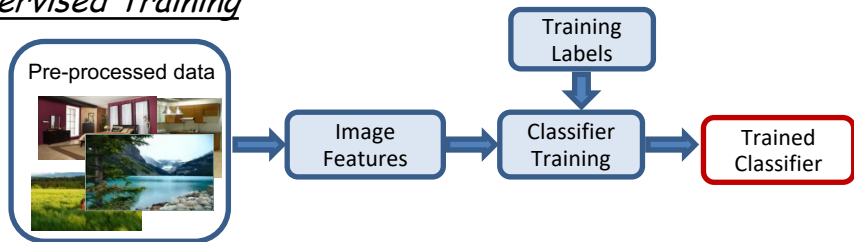
1. Pre-processing (Background subtraction)
2. Feature Selection (hand shapes)
3. Classification (Random Forest classifier)



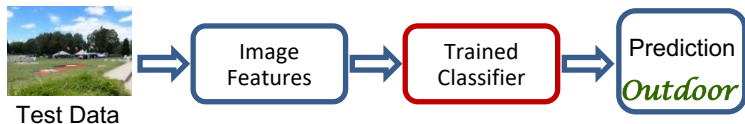
Summary: Typical Data Analysis Problem

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

Supervised Training



Testing



Assessments

- CW – *An Unknown Signal*: report + code (40% of unit) weeks 15-21 [submission in week 21]
- Discuss with others, but submissions are individual
- Assessment for course work is marked in the form of a report – to illustrate what you have understood about the data
- Exam (60% of unit)

Unit pages : https://github.com/LaurenceA/COMS20011_2020

Labs

- Thursdays 15:00 - 16:00 [by timetable]: Group 1
- Thursdays 16:00 - 17:00 [by timetable]: Group 2
- Lab Environment [Jupyter + Python]
- TA support in Teams: **grp-COMS20011_2020**
- Labs are essential for the coursework!

