

# Machine Learning

## Summary

---

Carl Henrik Ek - carlhenrik.ek@bristol.ac.uk

November 27, 2017

<http://www.carlhenrik.com>

# Introduction

---



# Phew

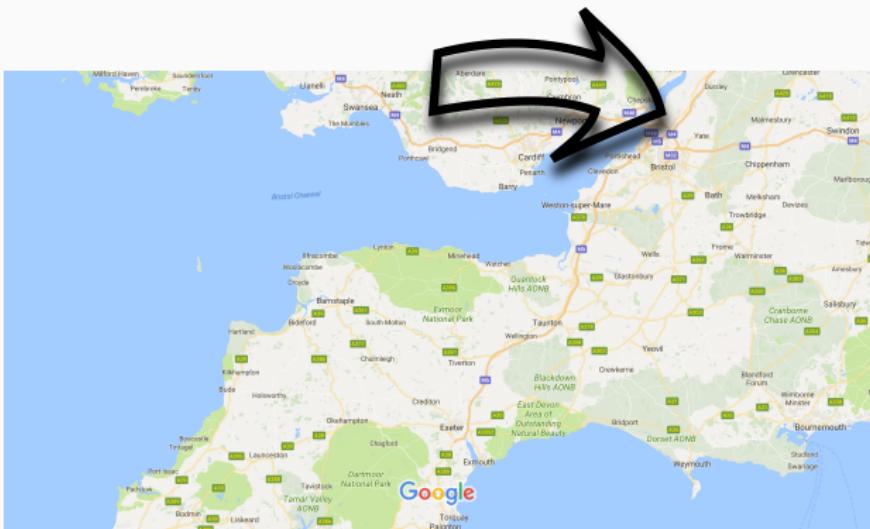
- 17 Lectures
- 17 Labs
- Lots of extra hours
- Just one week more

# Today

- Summary
- Exam
- What to do next

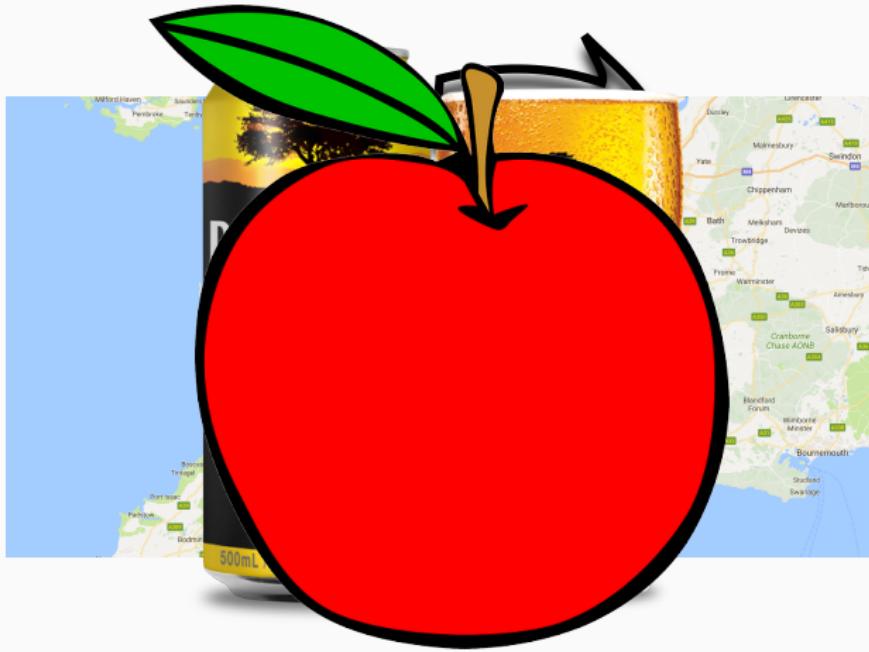
## Summary

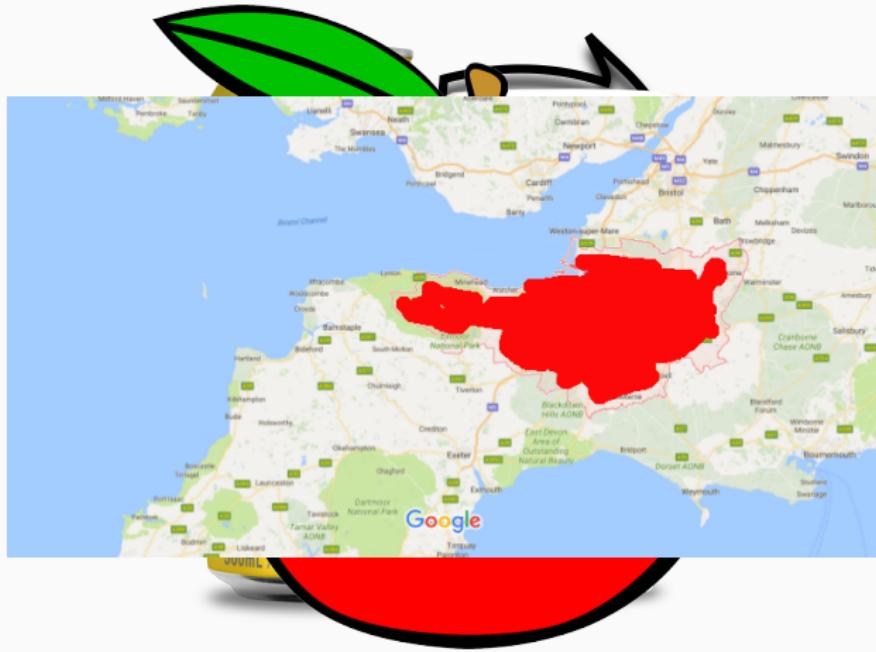
---

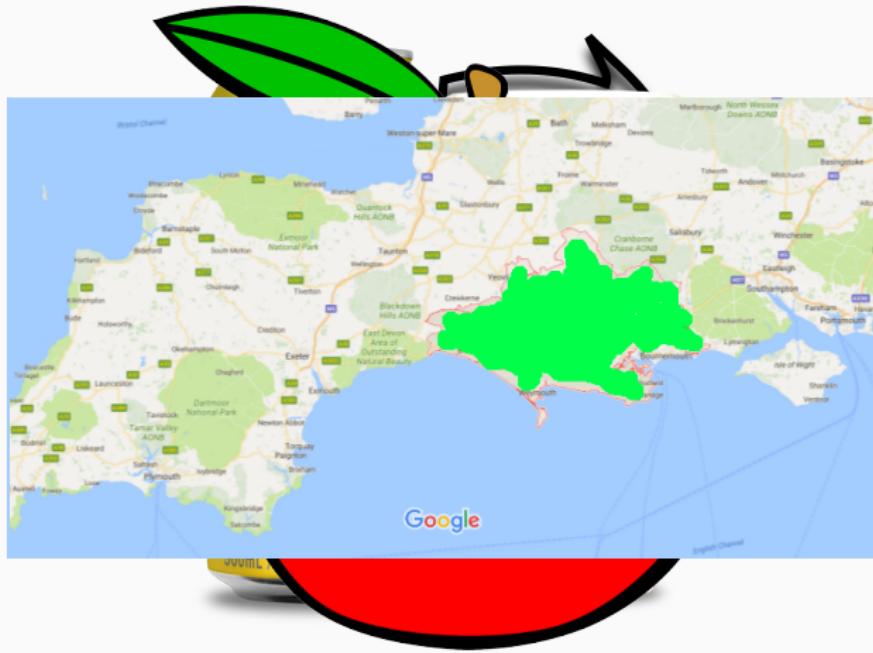


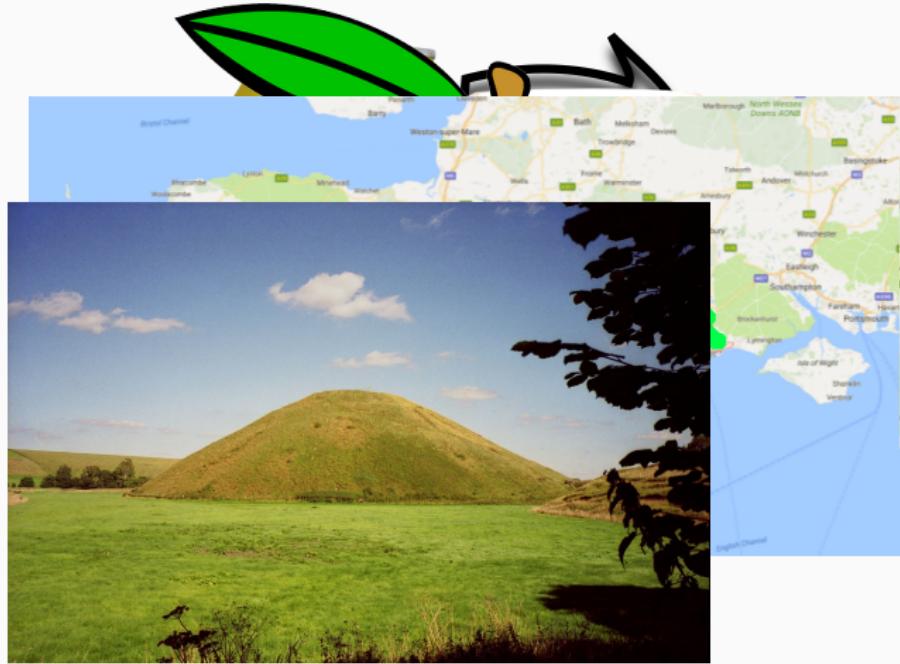




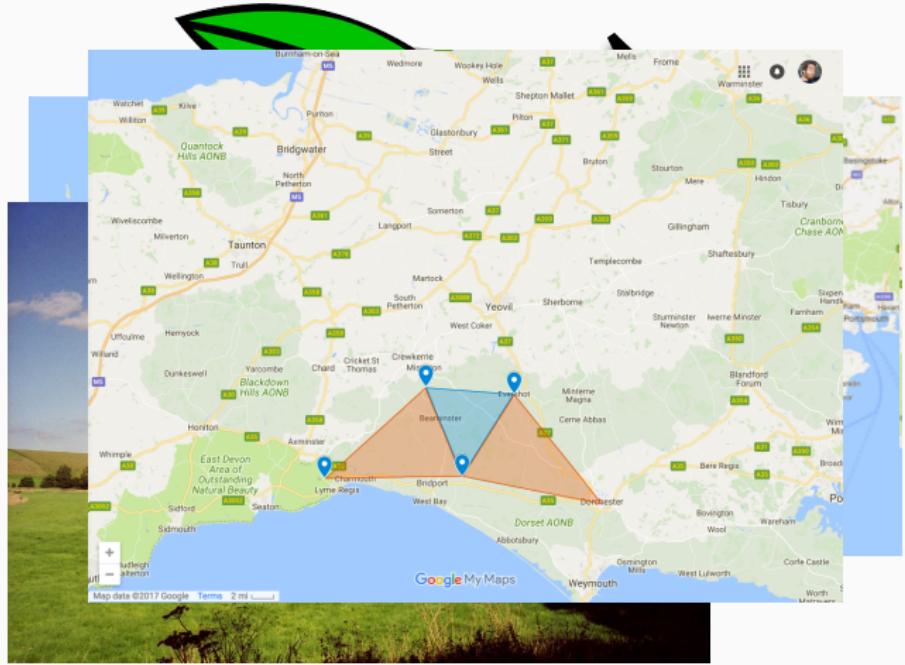


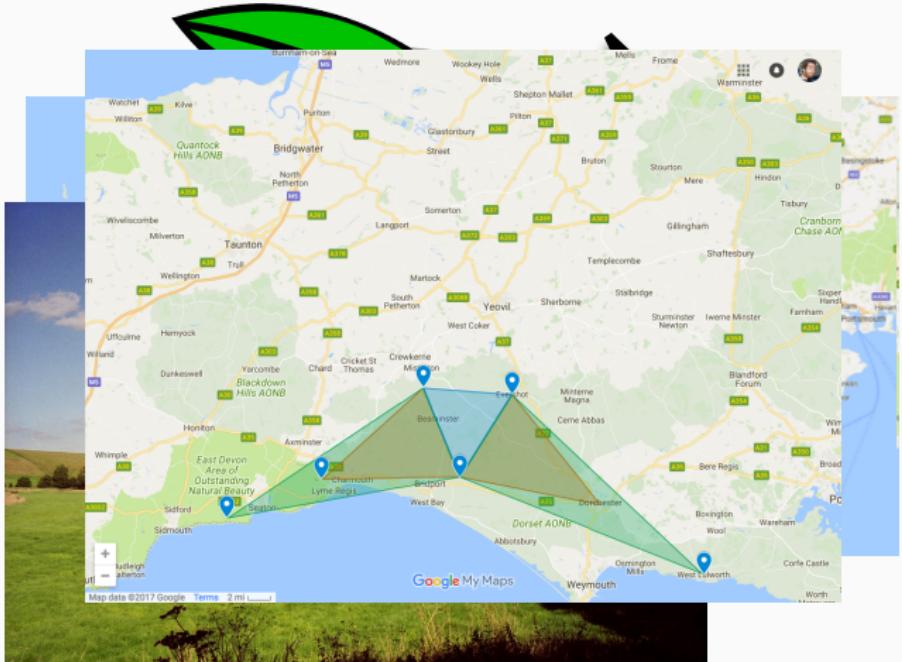




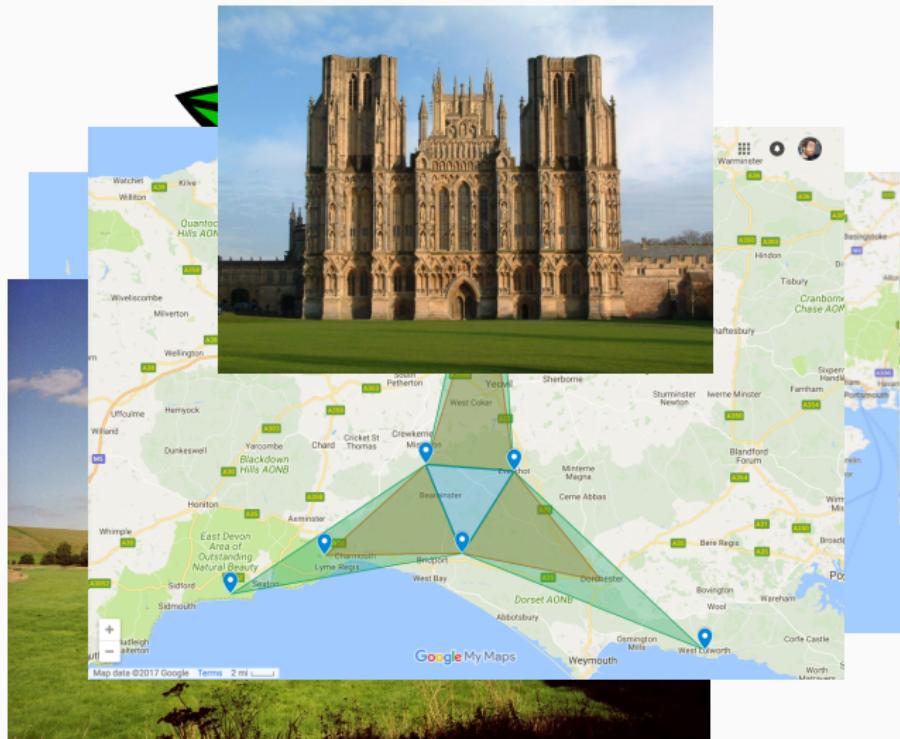


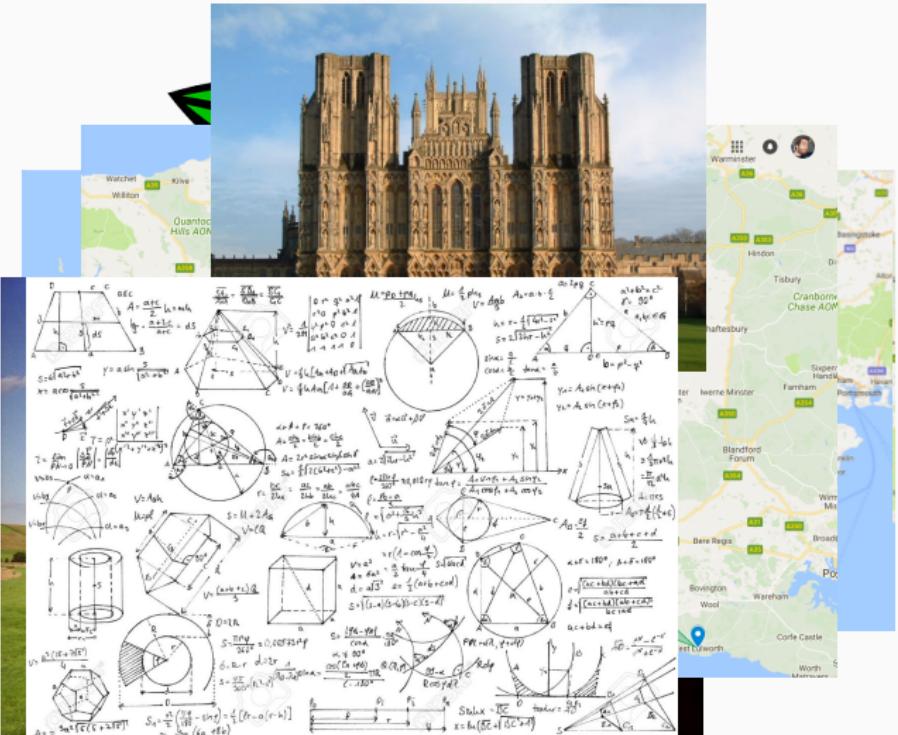


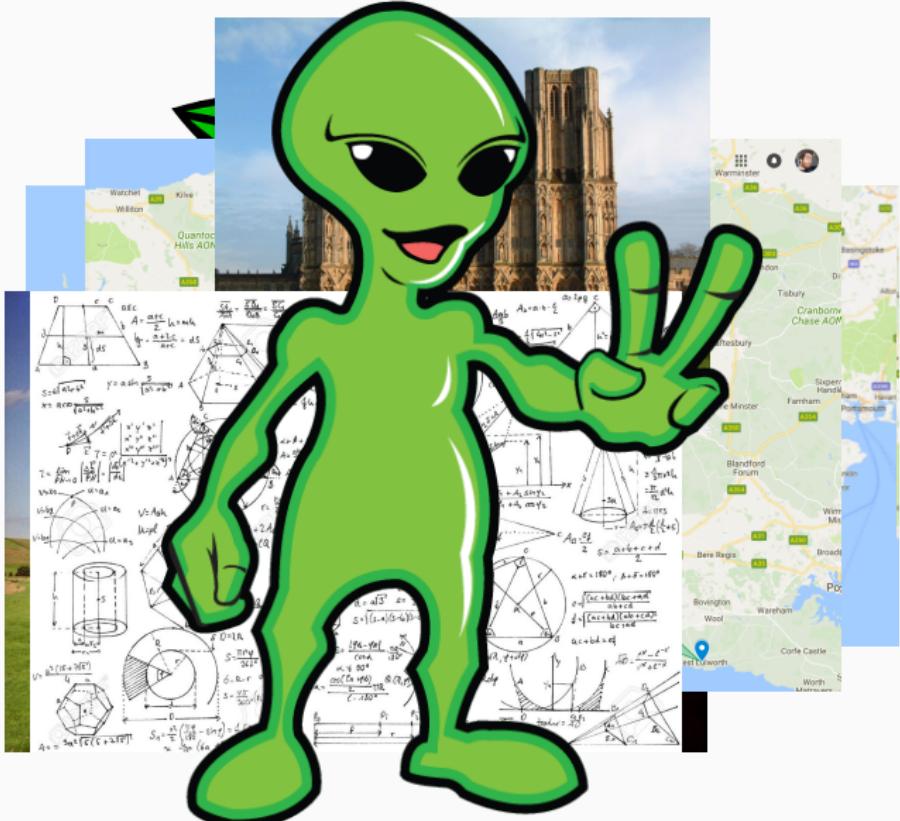






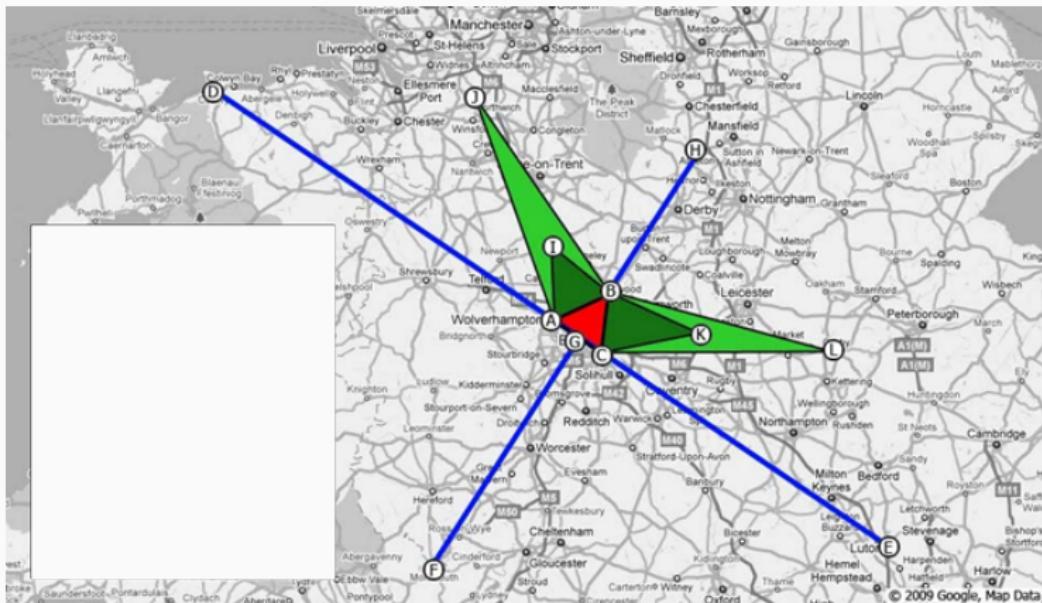






# Tom Brooks







*"We know so little about the ancient Woolworths stores," he explains, "but we do still know their locations. I thought that if we analysed the sites we could learn more about what life was like in 2008 and how these people went about buying cheap kitchen accessories and discount CDs"*— The Guardian<sup>1</sup>

<sup>1</sup><https://www.youtube.com/watch?v=XiigTGKZfks>



# Laplace Demon [1]



## Laplace Demon [1]

### Laplace's Demon [1]

*An intelligence which at a given instant knew all the forces acting in nature and the position of every object in the universe - if endowed with a brain sufficiently vast to make all necessary calculations - could describe with a single formula the motions of the largest astronomical bodies and those of the smallest atoms. To such an intelligence, nothing would be uncertain; the future, like the past, would be an open book.*

## Laplace Demon [1]

*All these efforts in the search for truth tend to lead the mind continuously towards the intelligence we have just mentioned, although it will always remain infinitely distant from this intelligence.*

# Zero Sum Games



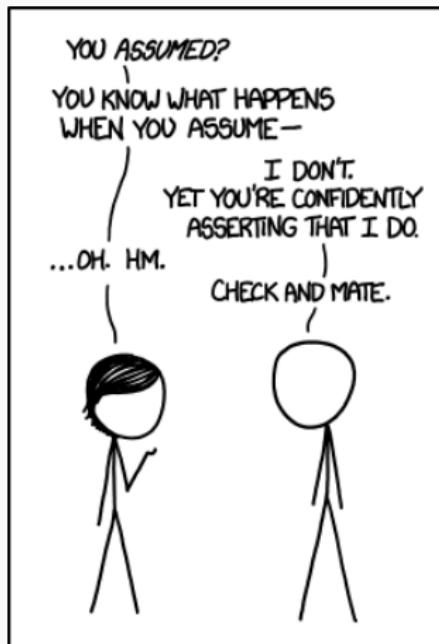
# Assumptions



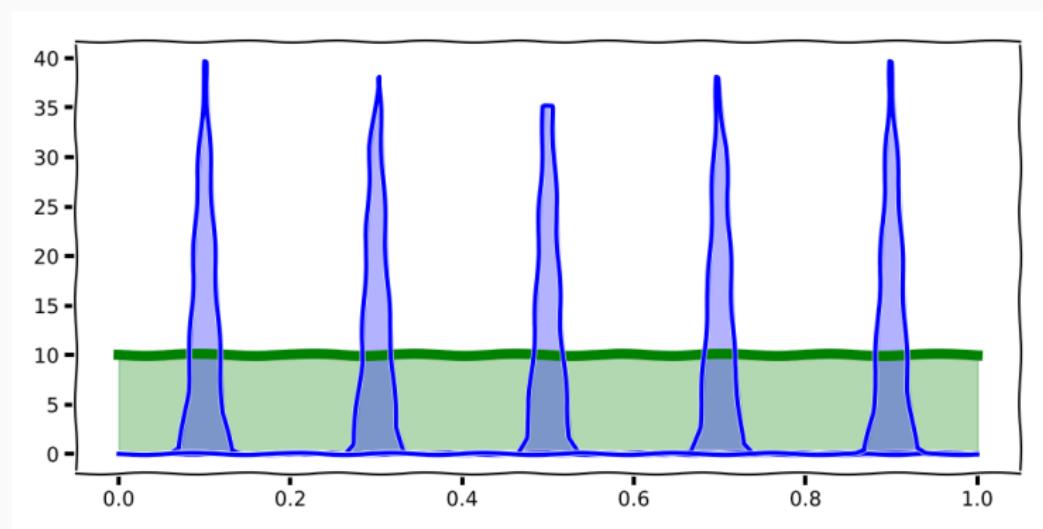
# Assumptions



# Assumptions



# Assumptions



## Ed Jaynes [2]

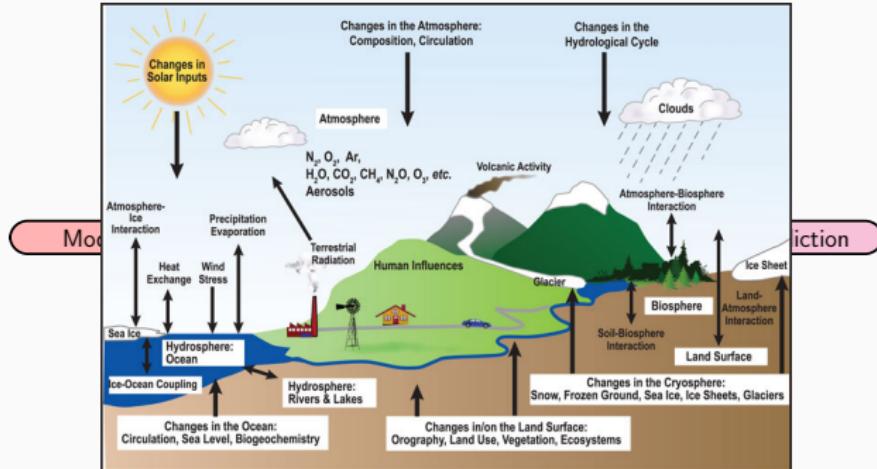


*It was our use of probability theory as logic that has enabled us to do so easily what was impossible for those who thought of probability as a physical phenomenon associated with “randomness”. Quite the opposite; we have thought of probability distributions as carriers of information.*

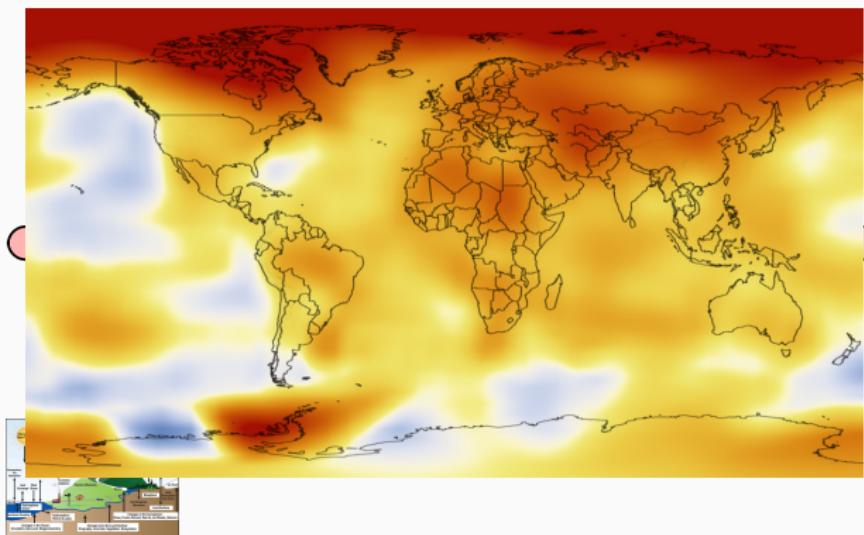
# Machine Learning



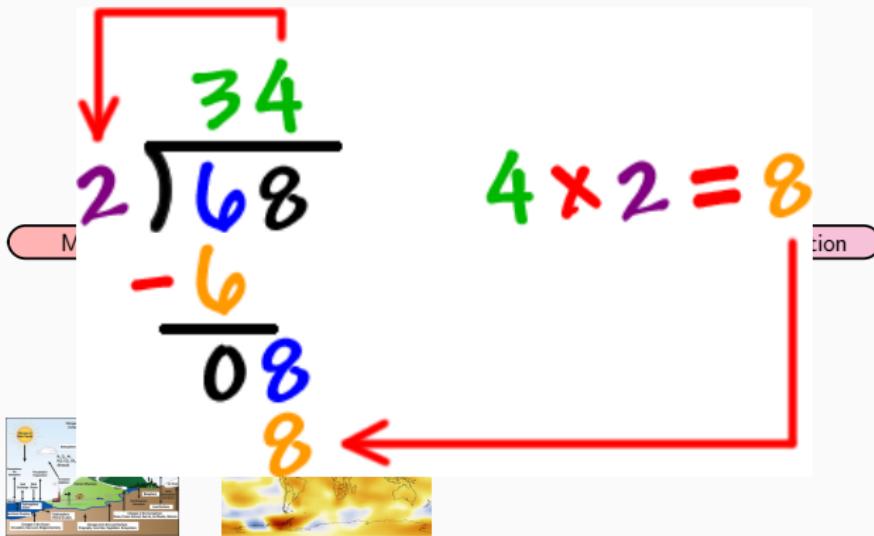
# Machine Learning



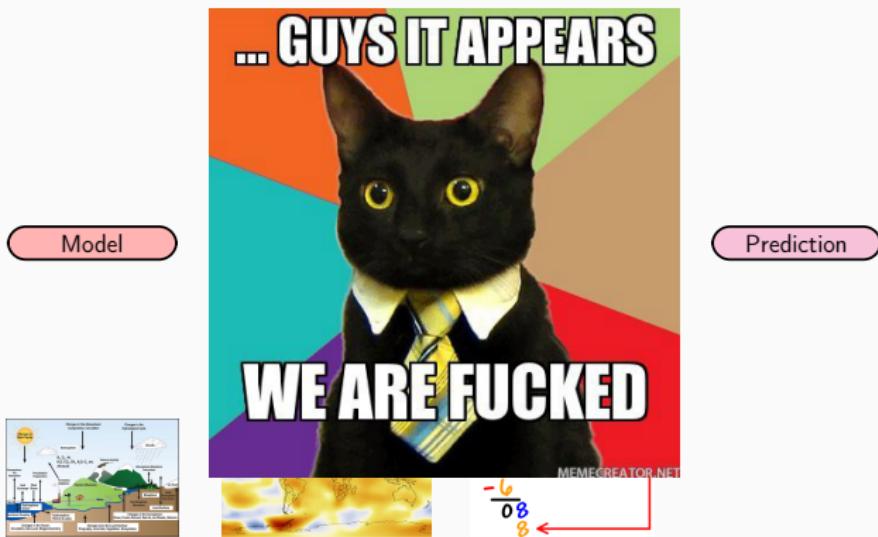
# Machine Learning



# Machine Learning



# Machine Learning



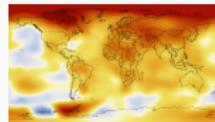
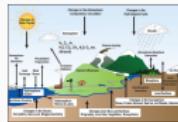
# Machine Learning

Model

Data

Algorithm

Prediction



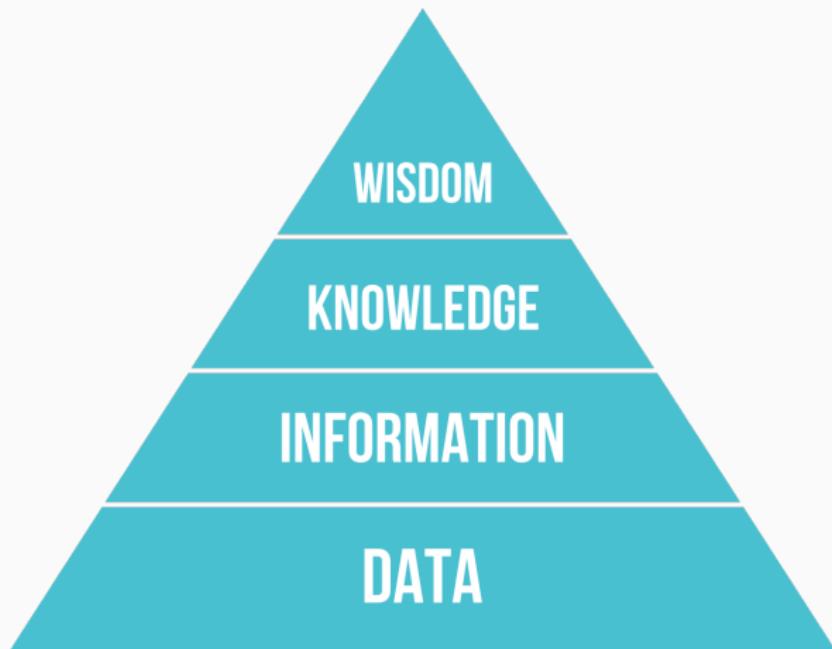
$$\begin{array}{r} 34 \\ \times 2 \\ \hline 68 \end{array}$$

4 x 2 = 8

08

8





# Part I

- Probability theory
- Distributions
- This provides a language

-<2-> *What is the language we have to formulate our assumptions with*

## Part II

- Models
- How can we put together probability distributions to merge several different assumptions
- Examples of different priors, parametric, non-parametric etc.

-<2-> *How can we factorise a probability distribution such that we get as low entropy model as possible*

## Part III

- Inference
- Learning implies reaching the posterior distribution
- The posterior is often intractable, how can we approximate it

## Part III

- Inference
- Learning implies reaching the posterior distribution
- The posterior is often intractable, how can we approximate it
-

## The other stuff

- Reinforcement Learning and decisions
- Neural networks
- Bayesian Optimisation

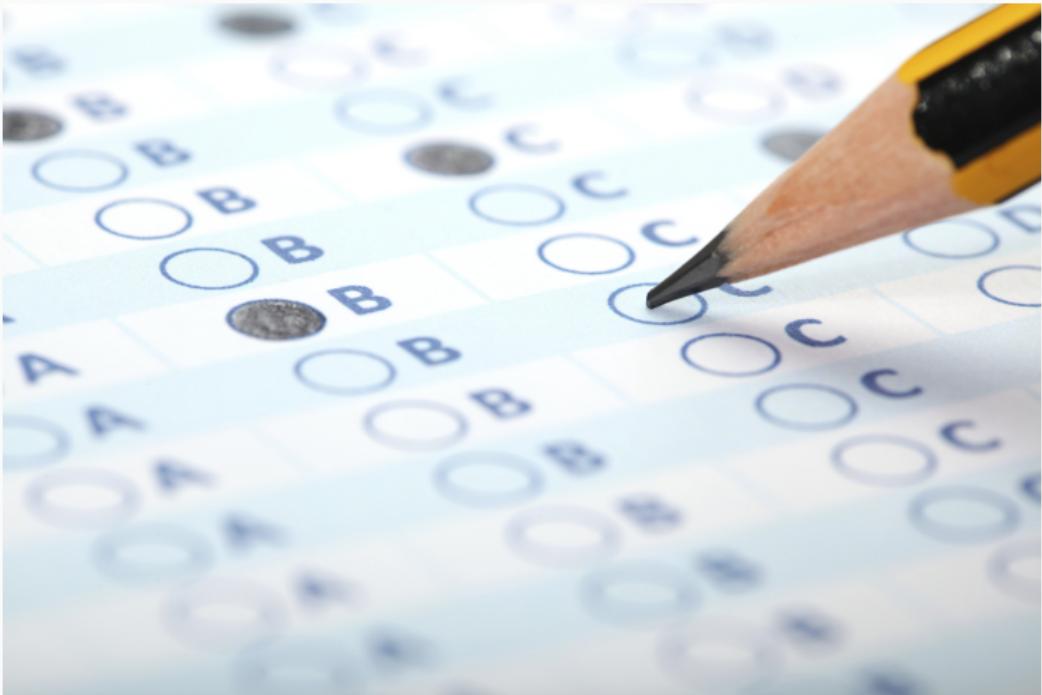
## The other stuff

- Reinforcement Learning and decisions
- Neural networks
- Bayesian Optimisation
- Examples of things that you can do and ways to think

# Exam

---

# Exam



# Exam

- 15 Questions
- Multiple choice
- Conceptual understanding

# Material

- Assignments
  - read through your own reports and the notes, what did you actually do? What were the key conceptual ideas.
- Summary document
  - main thing is the conceptual ideas
- Lecture notes
  - summarise each lecture to yourself, there are a few things in each lecture, understand them but skip the details
  - topic model lecture summarises a lot of part I and II

## Question 1

---

*Can any function be a kernel?*

# Answer

$$||x_i - x_j||_2^2$$

# Answer

$$||x_i - x_j||_2^2 = (x_i - x_j)^T(x_i - x_j)$$

## Answer

$$\|x_i - x_j\|_2^2 = (x_i - x_j)^T (x_i - x_j) = x_i^T x_i + x_j^T x_j - 2x_i^T x_j$$

## Answer

$$\begin{aligned} \|x_i - x_j\|_2^2 &= (x_i - x_j)^T (x_i - x_j) = x_i^T x_i + x_j^T x_j - 2x_i^T x_j \\ &= k(x_i, x_i) + k(x_j, x_j) - 2k(x_i, x_j) \end{aligned}$$

## Answer

$$\begin{aligned}\|x_i - x_j\|_2^2 &= (x_i - x_j)^T(x_i - x_j) = x_i^T x_i + x_j^T x_j - 2x_i^T x_j \\&= k(x_i, x_i) + k(x_j, x_j) - 2k(x_i, x_j) \\||v_i + v_j|| &\leq ||v_i|| + ||v_j||\end{aligned}$$

- The kernel function **needs** to represent an inner-product in a metric space, i.e. the triangle inequality needs to hold

## Question 2

*The conjugate prior to the **scale**  $\theta$  of a Weibull distribution with known shape  $\beta$  is a Inverse-Gamma distribution. Which functional form will the posterior distribution over the scale  $\theta$  take?*

## Answer 2

### Definition (Conjugate Prior)

In Bayesian probability theory, if the posterior distributions  $p(\cdot|x)$  are in the same family as the prior probability distribution  $p(\cdot)$ , the prior and posterior are then called conjugate distributions, and the prior is called a conjugate prior for the likelihood function.

## Question 3

*In Bayesian optimisation we use a surrogate model to optimise an intractable function. If we design an acquisition function that does not explore sufficiently, is our optimisation likely to become more or less dependent on initialisation? If our function over explores what is likely to be the effect?*

## Answer 3

1. *More dependent on initialisation.*
2. *The optimisation will require more iterations.*

## Question 4

*Marginalisation is the process of integrating out our belief in a variable. This is useful because the uncertainty in the variable disappears.*

## Answer 4

*Wrong, very wrong, we take the expectation of the variable so we include all our knowledge in the variable*

## Question 5

*When using a mean-field approximation to a variational distribution, what is our main assumption?*

## Answer 5

$$q(\mathbf{x}) = \prod_i^D q_i(x_i) \approx p(\mathbf{x}|\mathbf{y})$$

*That we model the marginals of the posterior distribution rather than the join.*

## Summary

- Focus on understanding of the conceptual material

## Summary

- Focus on understanding of the conceptual material
- Go through the assignment, what where the actual things that you needed to think of, the choices that enabled you to solve the tasks, the assumptions that you did

## Summary

- Focus on understanding of the conceptual material
- Go through the assignment, what where the actual things that you needed to think of, the choices that enabled you to solve the tasks, the assumptions that you did
  - are you clear with that linear regression and PCA are **exactly** the same thing?

# Summary

- Focus on understanding of the conceptual material
- Go through the assignment, what where the actual things that you needed to think of, the choices that enabled you to solve the tasks, the assumptions that you did
  - are you clear with that linear regression and PCA are **exactly** the same thing?
- Go through the lecture notes/summary document, what is the key conceptual messages for each lecture, understand these

# Summary

- Focus on understanding of the conceptual material
- Go through the assignment, what where the actual things that you needed to think of, the choices that enabled you to solve the tasks, the assumptions that you did
  - are you clear with that linear regression and PCA are **exactly** the same thing?
- Go through the lecture notes/summary document, what is the key conceptual messages for each lecture, understand these
- 2h exam, 15 questions, i.e. 4 minutes per questions, if it is complicated, you are on the wrong track

## Back yourself

$$p(\text{ML}|\text{COMS30007}) = \frac{p(\text{COMS30007}|\text{ML})p(\text{ML})}{p(\text{COMS30007})}$$

## What to do next?

---

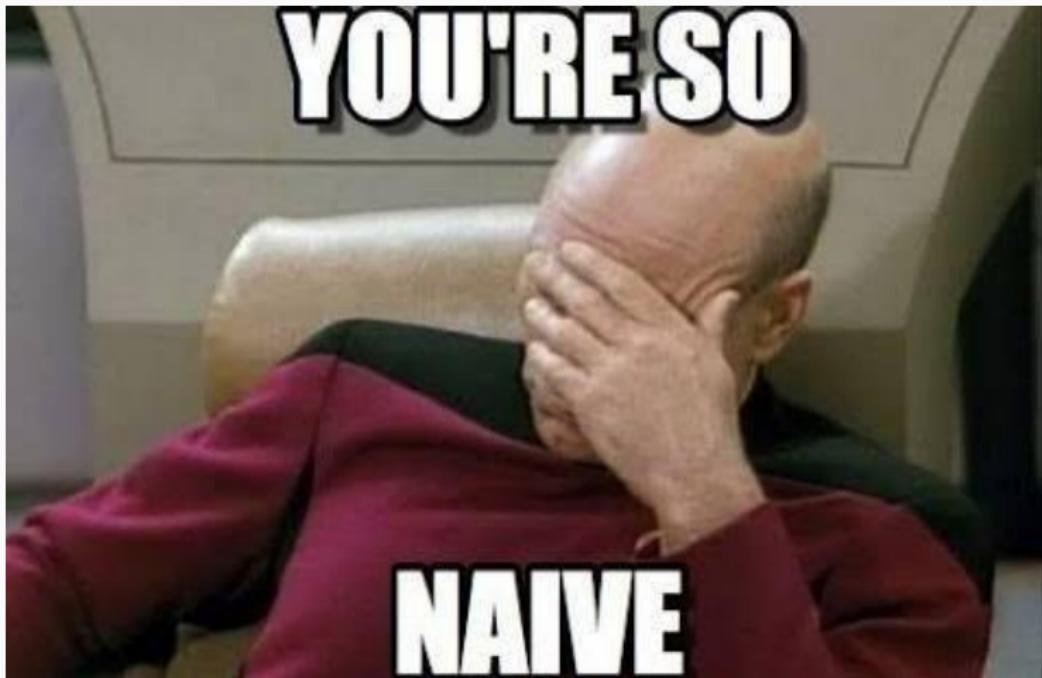
# Data



# Methods



- always remember and think about the trade-off between principle and what it allows you to do



## Right of Explanation<sup>2</sup>



### Article 22. Automated individual decision making, including profiling

*The data subject shall have the right not to be subject to a decision based solely on automated processing, including profiling, which produces legal effects concerning him or her or similarly significantly affects him or her.*

<sup>2</sup>Parliament and Council of the European Union (2016). General Data Protection Regulation.

# Communication



- A machine learner does not know anything but can do everything if we can communicate

# Machine Learning





CONSUMER  
CREATOR

# Ellen Key



*"Knowledge is what is left when you have forgotten  
everything that you have learned"*

*– Ellen Key*

## Summary

---

## Summary

- Machine Learning is the science of formulating assumptions and updating them with data
- Probability theory provides us with a useful language to make assumptions explicit
- You are the best source of knowledge, try to make how you function/reason into explicit formulations

eof

## References

---

 Pierre Simon Laplace.

A philosophical essay on probabilities, 1814.

 E T Jaynes.

*Probability theory: The logic of science.*

Cambridge university press, June 2003.