Computer Science

# IB Foundations
# of Data Science

MRes AI4ER

# Data Science part 1

Lecturer

# Dr Damon Wischik

Data Scientist: The Sexies ×

🔒 Secure | https://hbr.org/2012/10/data-scientist-the-sexiest-job-of-the-21st-century

☰ **Harvard Business Review**

Subscribe  Sign

WALTER BU
Academic
formation
Camera
"A W

DON

Business
Bowling
Honor So
Dancing

A, 1, 2,
Future
Record

ARTWORK: TAMAR COHEN, ANDREW J
ON A PAGE FROM A HIGH SCHOOL YEAR

DATA

# Data Scientist: The Sexiest Job of the 21st Century

by Thomas H. Davenport and D.J. Patil

FROM THE OCTOBER 2012 ISSUE

WHAT TO READ NEXT

Still the Sexiest Prof

📄 SUMMARY  ⊞ SAVE  ↗ SHARE  💬 **11** COMMENT  H**H** TEXT SIZE  🖨 PRINT  **$8.95** BUY COPIES

---

The world's most valuab ×

🔒 Secure | https://www.economist.com/news/leaders/21721656-data-economy-demands-new-approac...

38

**Regulating the internet giants**

# The world's most valuable resource is no longer oil, but data

*The data economy demands a new approach to antitrust rules*

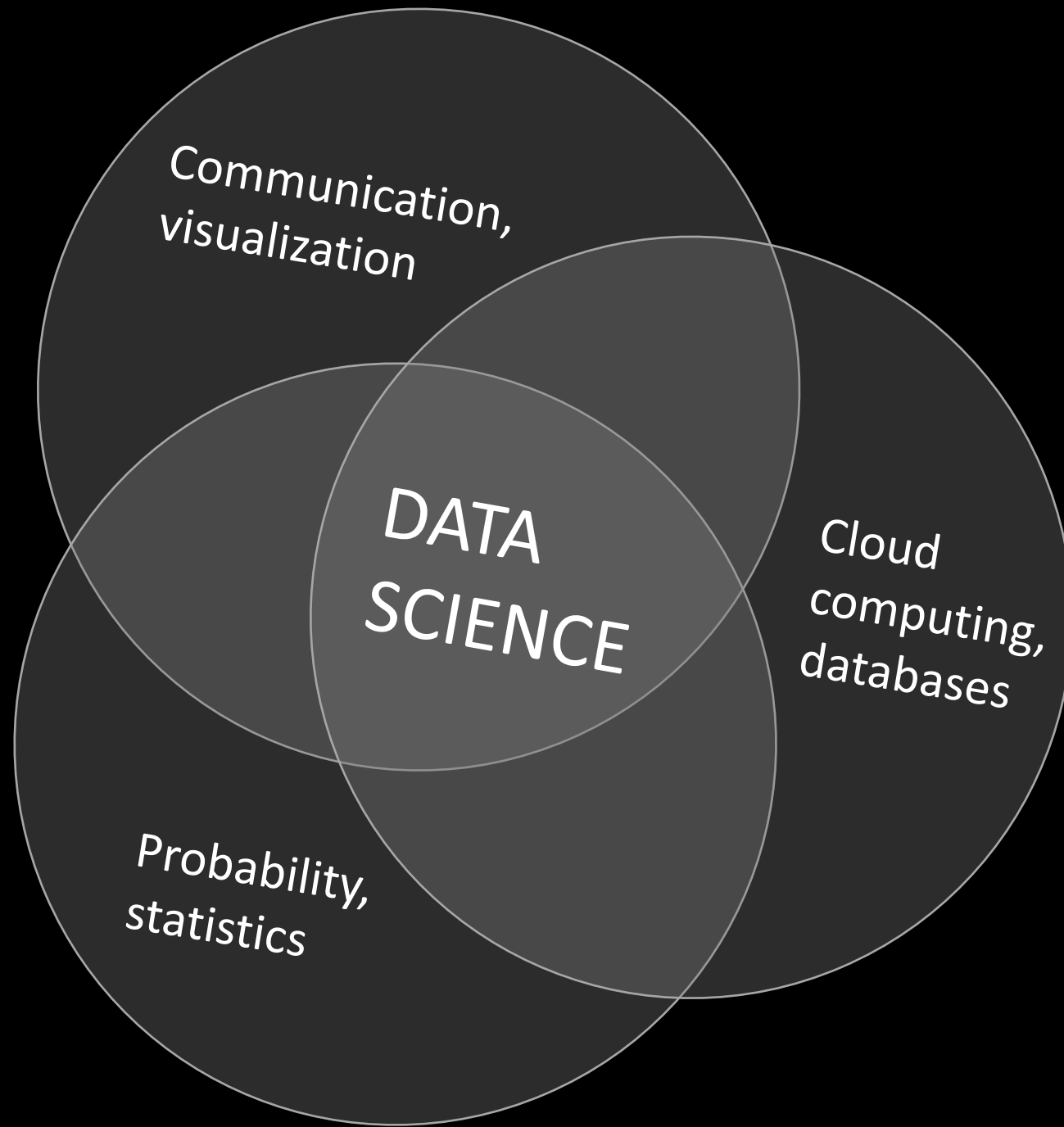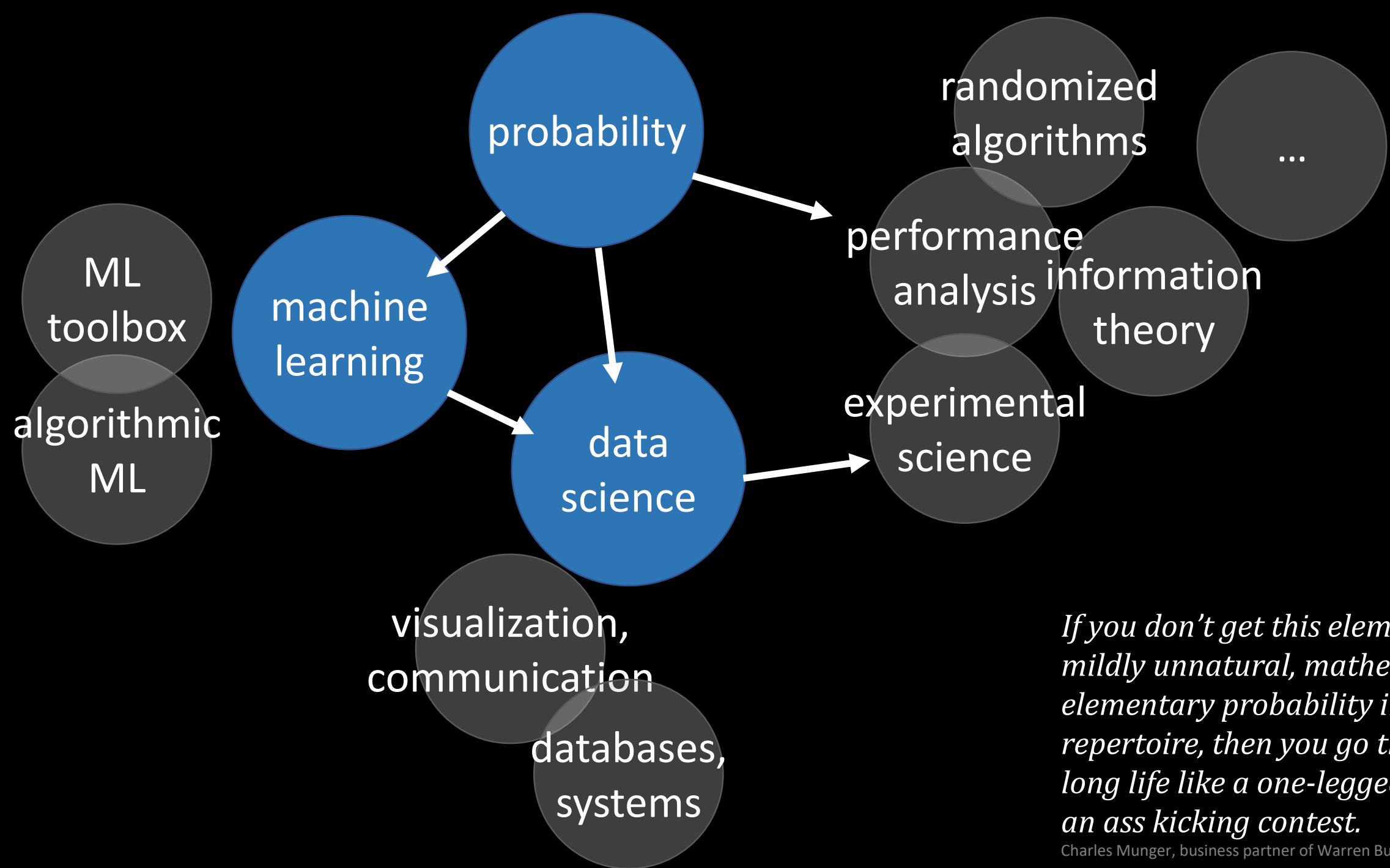amazon  UBER  Microsoft  Google  f  TESLA

David Parkins

📖 Print edition | Leaders ❯

May 6th 2017

🐦 f in ✉ 🖨 💬

A NEW commodity spawns a lucrative, fast-growing industry, prompting antitrust regulators to step in to restrain those who control its flow. A century ago, the resource in question was oil. Now similar concerns are

If you don't get this elementary, but mildly unnatural, mathematics of elementary probability into your repertoire, then you go through a long life like a one-legged man in an ass kicking contest.

Charles Munger, business partner of Warren Buffett.

# NEWS

Find local news 📍

Home    UK    World    Business    Politics    Tech    Science    Health    More ▾
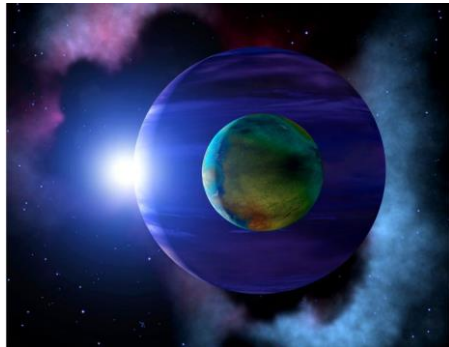
## Science & Environment

# Signal may be from first 'exomoon'

By Paul Rincon
Science editor, BBC News website

🕐 27 July 2017 | Science & Environment

f    🐦    💬    ✉    ⤴ Share



The work by Dr Kipping, his Columbia colleague Alex Teachey and citizen scientist Allan R Schmitt, assigns a confidence level of four sigma to the signal from the distant planetary system. The confidence level describes how unlikely it is that an experimental result is simply down to chance. If you express it in terms of tossing a coin, it's equivalent to tossing 15 heads in row.

But Dr Kipping said this is not the best way to gauge the potential detection.

He told BBC News: "We're excited about it... statistically, formally, it's a very high probability. But do we really trust the statistics? That's something unquantifiable. Until we get the measurements from Hubble, it may as well be 50-50 in my mind."

*If you don't get this elementary, but mildly unnatural, mathematics of elementary probability into your repertoire, then you go through a long life like a one-legged man in an ass kicking contest.*

Charles Munger, business partner of Warren Buffett.

CONTENT
notes,
example sheets,
practical walkthroughs

EPHEMERA
timetable, slides,
announcements, Q&A

code snippets from lectures

IB Foundations of Data Science

Damon Wischik, Computer Science, Cambridge University

Contents

- Material in the notes is examinable (except sections marked *)

- The notes include a few more examples than we cover in lectures

- I am preparing DRAFT expanded notes; they contain extra material that is not examinable

# Example sheet 0
### Remembering IA Maths for NST
### Foundations of Data Science—DJW—2019/2020

*Foundations of Data Science* builds on the probability theory you learnt in IA *Maths for the Natural Sciences Tripos*. All of the questions below (apart from the last two) are taken from that course. Please look through and make sure you can still answer them! Solutions will be provided.

*For supervisors: it isn't intended that you supervise this example sheet.*

**Question 1.** A card is drawn at random from a pack. Event $A$ is 'the card is an ace', event $B$ is 'the card is a spade', event $C$ is 'the card is either an ace, or a king, or a queen, or a jack, or a 10'. Compute the probability that the card has (i) one of these properties, (ii) all of these properties.

**Question 2.** A biased die has probabilities $p$, $2p$, $3p$, $4p$, $5p$, $6p$ of throwing 1, 2, 3, 4, 5, 6 respectively. Find $p$. What is the probability of throwing an even number?

**Question 3.** Consider drawing 2 balls out of a bag of 5 balls: 1 red, 2 green, 2 blue. What is the probability of the second ball drawn from the bag being blue given that the first ball was blue if (i) the first ball is replaced, (ii) the first ball is not replaced?

**Question 4.** Two cards are drawn from a deck of cards. What is the probability of drawing two queens, given that the first card is not replaced?

**Question 5.** A screening test is 99% effective in detecting a certain disease when a person has the disease. The test yields a 'false positive' for 0.5% of healthy persons tested. Suppose 0.2% of the population has the disease. (i) What is the probability that a person whose test is positive has the disease? (ii) What is the probability that a person whose test is negative actually has the disease after all?

**Question 6.** What is the probability that in a room of $r$ people at least two have the same

- Example sheet 0 review of IA Maths, not for supervision, solutions will be provided

- Example sheets 1, 2, 3 for supervision

- Ideas explained in the example sheets are examinable

SPRINGER TEXTS IN STATISTICS

A Modern
Introduction
to Probability
and Statistics

Understanding Why
and How

F.M. Dekking
C. Kraaikamp
H.P. Lopuhaä
L.E. Meester

Springer

PROBABILITY
MODELS
for Computer Science

DISCRETE CASE

Sheldon M. Ross

Probability and
Computing

Randomization and Probabilistic Techniques
in Algorithms and Data Analysis

Michael Mitzenmacher
and Eli Upfal

SECOND EDITION

reddit

hot | new | rising | controversial | top | gilded | wiki

MY SUBREDDITS ▾   POPULAR · ALL · RANDOM · USERS · ASKREDDIT · PI   MORE »

DATA IS
BEAUTIFUL

hot | new | top | wiki | Want to

MY SUBREDDITS ▾   POPULAR · ALL · RANDOM · USERS · ASKREDDIT · PICS · FU

DATASETS

Search

HOT | NEW | RISING | TOP | GILDED

Probability and Statistics for Programmers

Think
Stats

O'REILLY®

Allen B. Downey

# I. Learning with probability models

Most machine learning and data science tools boil down to

1. write out a probability model

2. learn from data  (fit the model / estimate its parameters)

*Monthly average temperatures [°C] at Cambridge station*



Probability models describe everything that *might have* happened, so you can interpret the significance of what *did* happen.

Most machine learning and data science tools boil down to

1. write out a probability model

2. learn from data  (fit the model / estimate its parameters)

*fitting ≡ optimization*

**Programming in the 2.0 stack**

Software 1.0 is code we write. Software 2.0 is code written by the *optimization* based on an evaluation criterion (such as "classify this training data correctly"). It is likely that any setting where the program is not obvious but one can repeatedly evaluate the performance of it (e.g. — did you classify some images correctly? do you win games of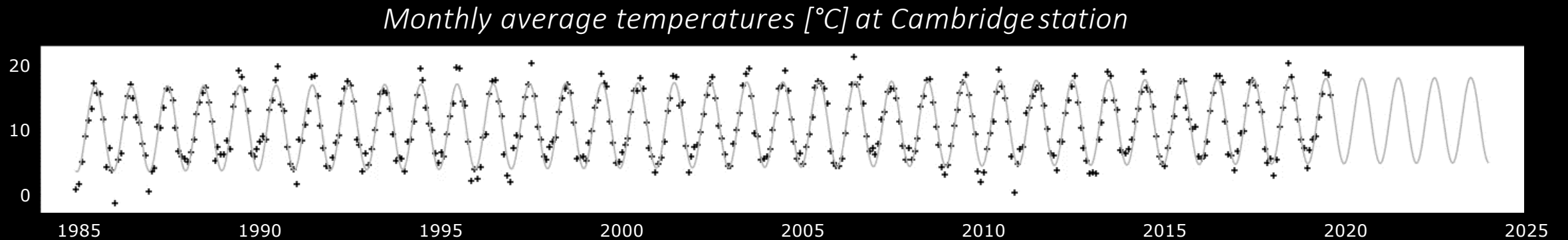 Go?) will be subject to this transition, because the optimization can find much better code than what a human can write.

Andrej Karpathy
@karpathy

Gradient descent can write code better than you. I'm sorry.

3:56 PM - 4 Aug 2017

343 Retweets  1,161 Likes

72     343     1.2K

# 1.1 Maximum likelihood estimation

tl;dr. Assume we have observed data, and we're told the probability model behind the data. Assume also that this probability model has an unknown parameter, which we wish to estimate.
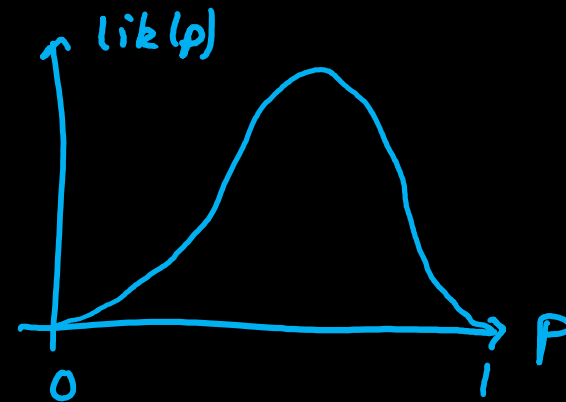
The *likelihood* is the probability of the observed data, viewed as a function of the unknown parameter. The *maximum likelihood estimator* or *mle* is the parameter value that maximizes the likelihood.

Exercise 1.1 (Coin tosses).

Suppose we take a biased coin, and tossed it $n = 10$ times, and observe $x = 6$ heads. Let's use the probability model

$$\mathbb{P}(\text{num.heads} = x) = \binom{n}{x} p^x (1-p)^{n-x}, \qquad x \in \{0,1,\ldots,n\}$$

where $p$ is the probability of heads and $1-p$ is the probability of tails. What is $p$?



$$\text{lik}(p) = \binom{n}{x} p^x (1-p)^{n-x}$$

$$\frac{d}{dp}: \quad \binom{n}{x}\left[ x\, p^{x-1}(1-p)^{n-x} - (n-x)p^x(1-p)^{n-x-1}\right]$$

$$\longrightarrow \quad \hat{p} = \frac{x}{n}$$

Or:

$$\log \text{lik}(p) = k + x \log p + (n-x)\log(1-p), \qquad \text{k doesn't depend on } p$$

$$\frac{d}{dp}: \quad \frac{x}{p} - \frac{n-x}{1-p} = 0 \quad \longrightarrow \quad \hat{p} = \frac{x}{n}$$

often we write $\text{lik}(p\,|\,x)$

Suppose we take a biased coin, and tossed it $n = 10$ times, and observe $x = 6$ heads. Let's use the probability model

$$\mathbb{P}(\text{num.heads} = x) = \binom{n}{x} p^x (1-p)^{n-x}, \qquad x \in \{0,1,\dots,n\}$$

where $p$ is the probability of heads and $1-p$ is the probability of tails.

Estimate $h = p/(1-p)$, called the "odds of heads".

HARD WAY. Write the model in terms of $h$, via $h = \frac{p}{1-p} \Leftrightarrow p = \frac{h}{1+h}$

$$\mathbb{P}(\text{numheads} = x) = \binom{n}{x} \left(\frac{h}{1+h}\right)^x \left(1 - \frac{h}{1+h}\right)^{n-x}$$

$$\longrightarrow \quad \hat{h} = \frac{x}{n-x}$$

EASY WAY. We already found $\hat{p} = \frac{x}{n}$.

Plug this into the formula for $h$: $\quad \hat{h} = \frac{\hat{p}}{1-\hat{p}} = \frac{x/n}{1-x/n} = \frac{x}{n-x}$

Suppose we ask $n = 100$ people their views on Brexit, and 37 say Leave, 35 say Remain, and the other 28 don't care. Using the probability model

$$\mathbb{P}(\text{leavers} = x_L, \text{ remainers} = x_R) = \frac{n!}{x_L!\, x_R!\, (n - x_L - x_R)!} p_L^{x_L} p_R^{x_R} (1 - p_L - p_R)^{n - x_L - x_R}$$

find maximum likelihood estimators for $p_L$ and $p_R$.

$$\log \text{lik} (p_L, p_R \mid x_L, x_R) = k + x_L \log p_L + x_R \log p_R + (n - x_L - x_R) \log (1 - p_L - p_R)$$

$$\frac{\partial}{\partial p_L}: \quad \frac{x_L}{p_L} - \frac{n - x_L - x_R}{1 - p_L - p_R} = 0$$

$$\frac{\partial}{\partial p_R}: \quad \frac{x_R}{p_R} - \frac{n - x_L - x_R}{1 - p_L - p_R} = 0$$

$$\hat{p}_L = x_L / n$$

$$\hat{p}_R = x_R / n.$$

Suppose we ask $n = 100$ people their views on Brexit, and 37 say Leave, 35 say Remain, and the other 28 don't care. Using the probability model

$$\mathbb{P}(\text{leavers} = x_L, \text{ remainers} = x_R) = \frac{n!}{x_L! \, x_R! \, (n - x_L - x_R)!} p_L^{x_L} p_R^{x_R} (1 - p_L - p_R)^{n - x_L - x_R}$$

find maximum likelihood estimators for $p_L$ and $p_R$.

⚠️ BAD ANSWER

*An estimator is a function of the observed data. You feed in the data, you get out an estimate.*

The log likelihood is
$$\log \text{lik} = \kappa + x_L \log p_L + x_R \log p_R + (n - x_L - x_R) \log(1 - p_L - p_R)$$

To find the maximum likelihood estimator for $p_L$, set the derivative equal to zero:
$$\frac{d}{d} \log \text{lik} = \frac{x_L}{p_L} - \frac{n - x_L - x_R}{1 - p_L - p_R} = 0$$

giving

$$\hat{p}_L = (1 - p_R) \frac{x_L}{n - x_R}$$

Similarly,

$$\hat{p}_R = (1 - p_L) \frac{x_R}{n - x_L}$$

*This is not a valid estimator because it involves an unknown parameter $p_R$.*

# 1.2 Numerical optimization

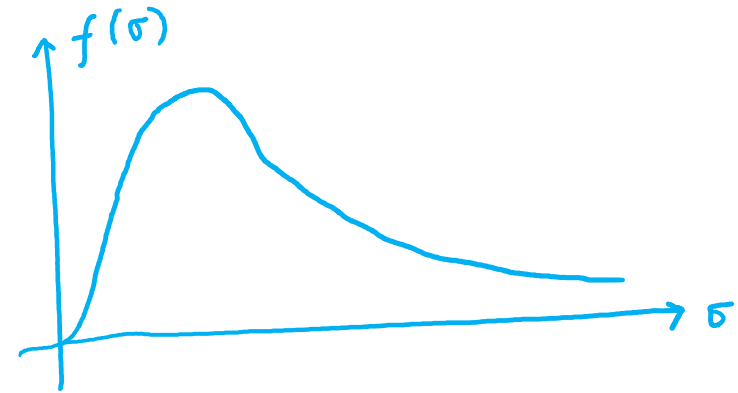tl;dr. To find the minimum of a function $f: \mathbb{R}^K \to \mathbb{R}$,

```
1    import scipy.optimize
2
3    def f(x):
4        return …
5
6    x₀ = […]  # initial guess
7    x̂ = scipy.optimize.fmin(f, x₀)
```

- Choose $x_0$ wisely
- This function finds a local minimum, perhaps not a global minimum
- See the documentation to control number of iterations, …

**Exercise 1.4 (Constrained optimization).**
Find the maximum over $\sigma > 0$ of
$$f(\sigma) = \frac{1}{\sqrt{2\pi\sigma^2}} e^{-3/2\sigma^2}$$

*Hint. Instead of maximizing over $\sigma > 0$, maximize over $\tau \in \mathbb{R}$ using the transform $\sigma = e^\tau$.*



```
1    import scipy.optimize
2    import numpy
3
4    π = numpy.pi
5
6    def f(σ):
7        return numpy.exp(-3/2/numpy.power(σ,2))  \
8                / numpy.sqrt(2*π*numpy.power(σ,2))
9
10   (τ̂,) = scipy.optimize.fmin(
11               lambda τ: -f(numpy.exp(τ)),
12               numpy.log(σ0))
13   σ̂ = numpy.exp(τ̂)

Optimization terminated successfully.
    Current function value: -0.139702
    Iterations: 13
    Function evaluations: 26

1.7320498691939412
```

**Exercise 1.5 (Softmax transformation).**

Find the maximum of

$$f(x_1, x_2, x_3) = 0.2 \log x_1 + 0.5 \log x_2 + 0.3 \log x_3$$

over $x_1, x_2, x_3 \in [0,1]$ such that $x_1 + x_2 + x_3 = 1$.

This is called the softmax transformation, and it's widespread in machine learning models.

Instead of optimizing over $(x_1, x_2, x_3)$, we'll optimize over $(\xi_1, \xi_2) \in \mathbb{R}^2$ with

$$x_1 = \frac{e^{\xi_1}}{e^{\xi_1} + e^{\xi_2} + 1} \qquad x_2 = \frac{e^{\xi_2}}{e^{\xi_1} + e^{\xi_2} + 1} \qquad x_3 = \frac{1}{e^{\xi_1} + e^{\xi_2} + 1}$$

The exponentiation ensures we get positive values, even for negative $\xi$.

The normalization ensures $x_1 + x_2 + x_3 = 1$.

```
1    def f(ξ):
2        ξ1,ξ2 = ξ
3        x = numpy.exp([ξ1,ξ2,0])
4        x1,x2,x3 = x / sum(x)
5        return 0.2*numpy.log(x1) + 0.5*numpy.log(x2) + 0.3*numpy.log(x3)
6
7    ξ1,ξ2 = scipy.optimize.fmin(lambda ξ: -f(ξ), [0,0])
8    x = numpy.exp([ξ1,ξ2,0])
9    x = x / numpy.sum(x)
```

*Optimization terminated successfully. Current function value: 1.02965. Iterations: 63. Function evaluations: 120*
*array([0.19999474, 0.49999912, 0.30000614])*