Computer Science

IB Data Science

Lecturer

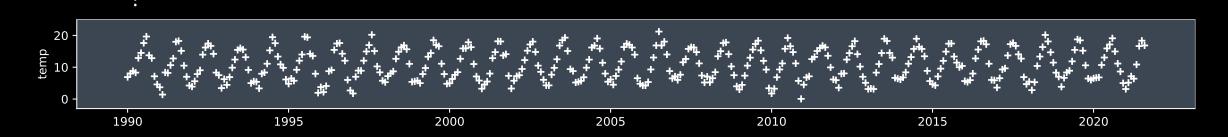
Dr Damon Wischik

Met Office climate dataset

https://www.metoffice.gov.uk/research/climate/maps-and-data/historic-station-data

Monthly readings from 37 weather stations around the country. Let's look at Cambridge, from 1990.

station	уууу	mm	t	af	rain	sun	tmin	tmax	temp
Cambridge	1990	1	1990.00	0	43.8	64.7	4.0	9.8	6.90
Cambridge	1990	2	1990.08	1	71.1	102.0	4.7	11.4	8.05
Cambridge	1990	3	1990.16	3	23.2	153.2	4.7	12.9	8.80

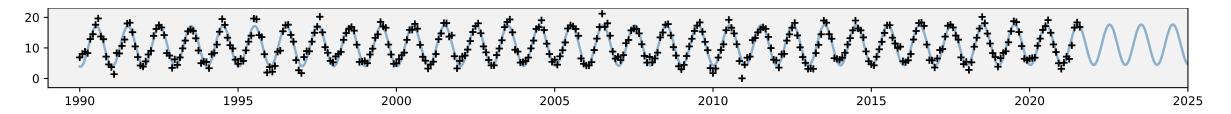


```
What model / formula would you suggest to fit this dataset?

def temp_model(t, ...):
    return ...
```

A SCIENTIST'S DETERMINISTIC MODEL

```
def temp_model(t, \alpha, \phi, c, \gamma):
return c + \alpha * np.sin(2*\pi*(t+\phi)) + \gamma*t
```



why? . to describe the Ada in front of my eys!

to tell my fitting routine how much attention to pay to errors low liers

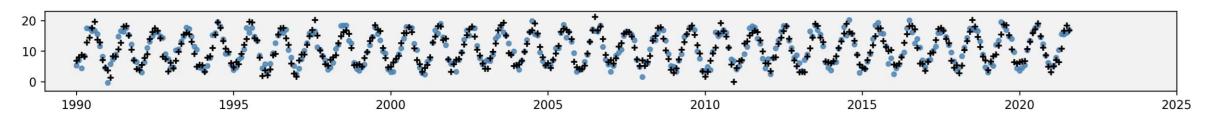
to reason about contridence, e.g. "8=1.4°C (anny ± 0.7°C"

A DATA SCIENTIST'S PROBABILITY MODEL

```
def rtemp(t, \alpha, \phi, c, \gamma, \sigma):

pred = c + \alpha * np.sin(2*\pi*(t+\phi)) + \gamma*t

return np.random.normal(loc=pred, scale=\sigma)
```



All of machine learning is based on a single idea:

- 1. Write out a probability model
- 2. Fit the model from data

This is behind

- A-level statistics formulae
- our climate model
- ChatGPT training

LECTURE NOTES

IB Data Science:

modelling and machine learning with probability

Damon Wischik, Computer Laboratory, Cambridge University









Contents

Introduction											
P	Prerequisites viii										
Ι	Lea	arning with probability models	1								
1	\mathbf{Spe}	cifying and fitting models	3								
	1.1	Specifying a probability model	3								
	1.2	Standard random variables	8								
	1.3	Maximum likelihood estimation	1								
	1.4	Numerical optimization with scipy	8								
	1.5	Likelihood notation	0								
	1.6	Types of model	3								
	1.7	Supervised and unsupervised learning	5								
2	Fea	ture spaces / linear regression 3	1								
	2.1	Fitting a linear model	2								
	2.2	Feature design	5								
		2.2.1 One-hot coding	5								
		2.2.2 Non-linear response	6								
		2.2.3 Comparing groups	6								
		2.2.4 Periodic patterns	7								
		2.2.5 Secular trend	8								
	2.3	Diagnosing a linear model	0								
	2.4	Linear regression and least squares	2								

- ABRIDGED NOTES (contain all examinable material + extras)
- EXTENDED NOTES
 (goes to masters-level, available soon)

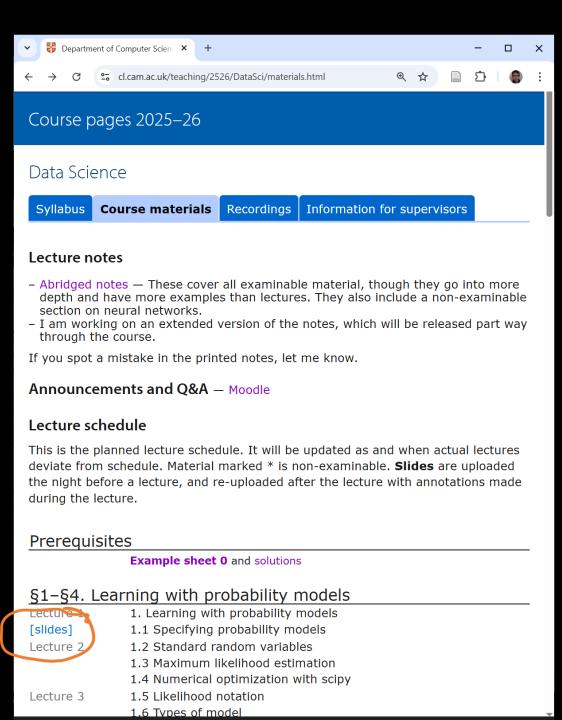
- The handout has more wordy explanations and more examples than lectures
- Use the handout like a textbook and take your own notes during lectures

Do you want printouts?

 $\S x$

Slides for each lecture are on the website
 and most slides say which section they're for

What's examinable?
 Everything in the lecture schedule,
 except for sections marked *



Consent to recordings of live lectures with Panopto

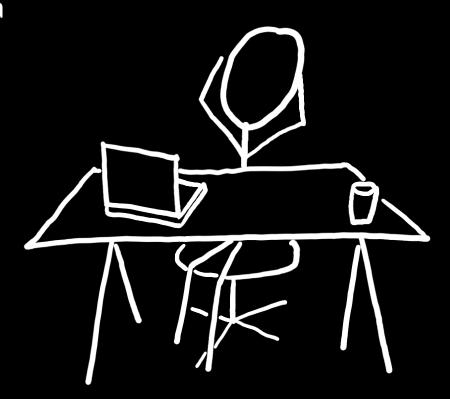
https://www.educationalpolicy.admin.cam.ac.uk/policy-index/recording

For any teaching session where your contribution is mandatory or expected, we must seek your consent to be recorded.

You are not obliged to give this consent, and you have the right to withdraw your consent after it has been given.

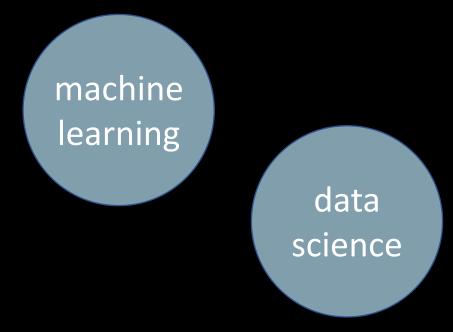
Do you give your consent to recordings?

■ I'll be in the cafe area 1–1.30pm today for office hours.



What is data science? What's the difference between data science and machine learning?

MACHINE LEARNING Building systems that learn from data.



DATA SCIENCE Finding patterns in a dataset.

What is data science? What's the difference between data science and machine learning? course probability algorithmic MACHINE LEARNING Building systems itat learn from data. obabilistic ML systems **Building models** generative Al building and applying data them. building science /isualization, models **communication** fron data statistics

DATA SCIENCE
Finding
patterns in a
dataset
Building
models to
learn about
the dataset.

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

People who try for do ML without probability

Example sheet 0 Prerequisites IB Data Science—DJW—2025/2026

This course assumes that you know how to handle basic probability problems and that you know about random variables, as taught in IA *Introduction to Probability*. It also assumes that you know how to find the maximum or minimum of a function, using calculus, as taught in IA *Maths for NST*. The code snippets in the course are in Python and numpy, and you should be familiar with numpy's way of writing vectorized computations.

This example sheet reviews the material that you need to know. Please look through, and make sure you remember how to answer these questions! Solutions are provided on the course website. For supervisors: this example sheet is not intended for supervision.

Rules of probability (IA Probability lecture 1)

Understand what is meant by sample space, written Ω , and know that $\mathbb{P}(\Omega) = 1$. Be able to reason about probabilities of events with Venn diagrams. Know the core definitions and laws ...

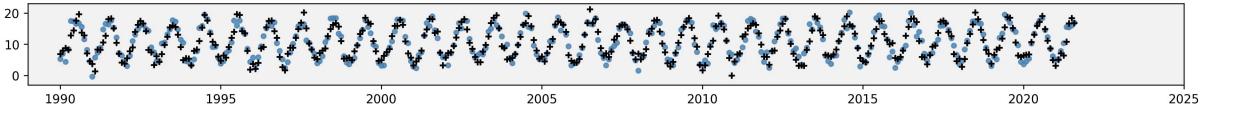
Conditional probability, or equivalently the chain rule:

$$\mathbb{P}(A \mid B) = \frac{\mathbb{P}(A, B)}{\mathbb{P}(B)} \text{ if } \mathbb{P}(B) > 0$$

$$\mathbb{P}(B, A) = \mathbb{P}(B) \mathbb{P}(A \mid B) \quad \text{(chain rule)}$$

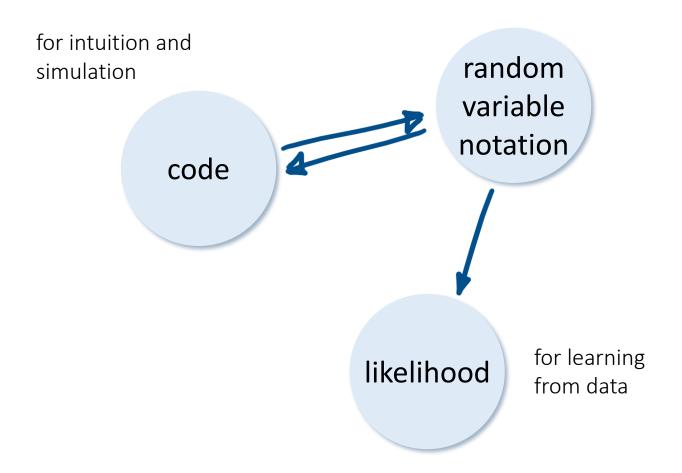
- Example sheet 0 is to remind you about IA Probability, Maths for NST, and Scientific Computing
- It's not for supervision; solutions are provided

How to specify a probability model

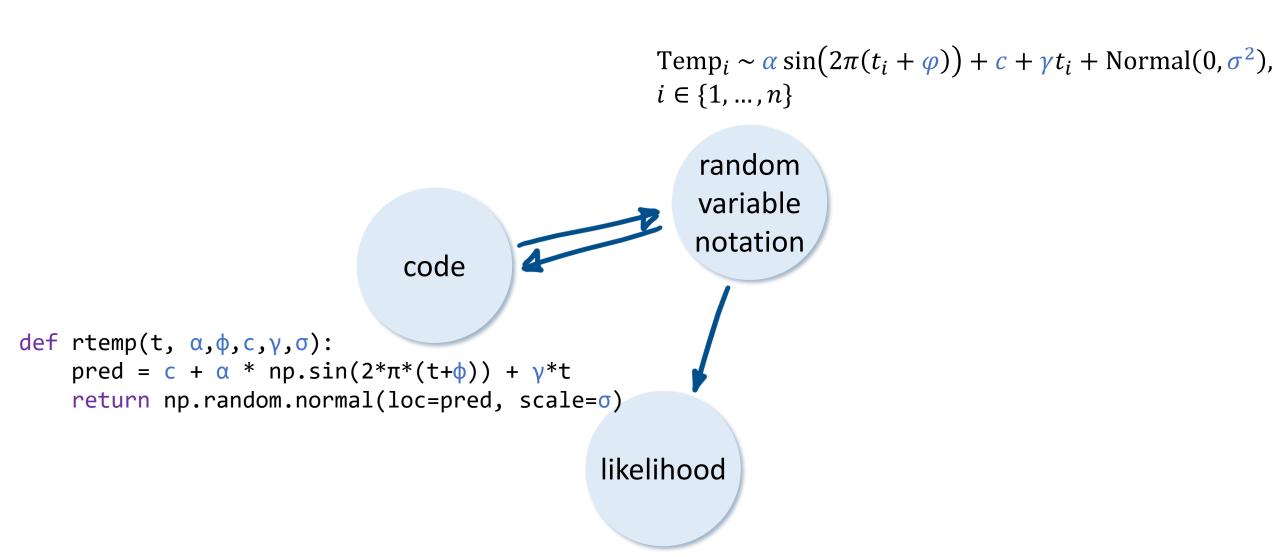


```
def rtemp(t, \alpha=10, \varphi=-0.25, c=11, \gamma=0.035, \sigma=2):
pred = c + \alpha * np.sin(2*\pi*(t+\varphi)) + \gamma*t
return np.random.normal(loc=pred, scale=\sigma)
```

Three views of a probability model



Three views of a probability model



```
Generale X from the Uniform dist.
def ry():
                                      X \sim U[0,1]
                                      Y = X^2
    x = random.random()
    y = x ** 2
    return y
                                                       Upper case for random variables
def ri(a,b):
                                      X \sim U[0,1]
                                                       Lower cox for pourameters, constants, destapoints.
                                      I = |aX + b|
    x = random.random()
    i = math.floor(a*x+b)
    return i
```

$$X \sim U[0,1]$$
$$Y = X^2$$

```
X1 and X2 are generated independently.
                                                X_1, X_2 \sim U[0,1]
def rz():
                                               Z = X_1 \log X_2
                                                                        "bearning the value of one rells is nothing about the other".
     x_1 = random.random()
     x_2 = random.random()
      return x_1 * math.log(x_2)
                                                                                                       unless specified
                                                                                                      otherwik.
                                                (Y,Z) \sim Myrandpair
def rmyrandpair():
     x_1 = random.random()
     x_2 = random.random()
     y,z = (x_1+x_2, x_1*x_2)
                                                                    A is lower-cook (not A) so it repers
to a fixed value.
     return (y,z)
                                                When we say "X, and Xz are independent" X_1, X_2 \sim U[0, \lambda] We mean "given all the parameters".
 \lambda = 3
 x_1 = random.uniform(0,\lambda)
 x_2 = random.uniform(0,\lambda)
```

```
x = random.random()
y = 1 - x
```

$$X \sim U[0,1]$$
$$Y = 1 - X$$







```
~ means "has the same distribution
   Y~ U[°1] }

Thek are all true statements.

X~ Y
    X ~ 1-X
       whenever I run the code". Not essignment!
              Y=1-X

X+Y=1 all true statements

X=1-Y
```

In Java, = is for ossignment.

$$X \sim U[0,1]$$
 This is how we inclicate $Y \sim N(X,0.1^2)$ that Y is generated based on X.

It's good to accompany this with a "consal diagram" $X \rightarrow Y$,

(We'll see lots of the conse.)

```
numpy vectorized coole:
def rtemp(t, \alpha=10, \phi=-0.25, c=11, \gamma=0.035, \sigma=2):
     pred = \alpha*np.sin(2*\pi*(t+\phi)) + c + \gamma*t
                                                                       generates 380 independent
     return np.random.normal(loc=pred, scale=o) ≥
                                                                       Normal random raniables
df = pandas.read_csv(...) # data frame, n=380 rows
Temp = rtemp(df.t)
                                                     ach with its own Normal.

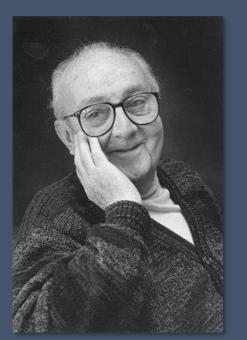
They're all independent, by our notation convention
Temp<sub>i</sub> ~ \alpha \sin(2\pi(t_i + \varphi)) + c + \gamma t_i + \text{Normal}(0, \sigma^2), \quad i \in \{1, ..., n\}
Temp<sub>i</sub> = \alpha \sin(2\pi(t_i + \varphi)) + c + \gamma t_i + E_i, E_i \sim \text{Normal}(0, \sigma^2), i \in \{1, ..., n\}
         What are the r.v.? Tempi, Normal,
         What one the non-1.v.? parameters: \alpha, \phi, c, \delta, \delta
```

- 1. Write out a probability model
- 2. Fit the model from data ____ rest lecture

This is behind

- A-level statistics formulae
- our climate model
- ChatGPT training

A core skill is being able to build useful probability models. During this course you will pick up this skill, through examples.



"All models are wrong but some are useful"

... so don't be shy when building your own. There *is* no "correct model".

§