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

i.e. estimate its parameters using Maximum Likelihood Estimation (mle)

The model typically has unknown parameters.

The <u>likelihood</u> is the probability of seeing the data that we actually saw.

It depends on the parameters.

Let's simply pick the parameters that maximize the likelihood!

#### Exercise 1.3.1 (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

$$X \sim \text{Binom}(n, p)$$

where p is the probability of heads. Estimate p.

#### Likelihood of the observed data:

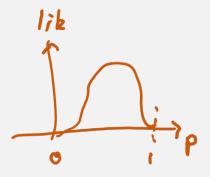
$$lik = P(X=x) \qquad x = 6$$

$$= \binom{h}{x} P^{x} (l-p)^{h-x}$$

Parameter that maximizes it:

$$\frac{d}{dp} lik = {n \choose 2} \left[ \times p^{x-1} \left( r - p \right)^{n-2} - (n-x) p^{x} \left( r - p \right)^{n-x-1} \right]$$

$$\frac{d}{dp} lik = 0 \Rightarrow \hat{r} = \frac{x}{n}$$



## There are standard numerical random variables that you should know:

DISCRETE RANDO	M VARIABLES	
Binomial $X \sim Bin(n, p)$	$\mathbb{P}(X = x) = \binom{n}{x} p^x (1 - p)^{n - x}$ $x \in \{0, 1, \dots, n\}$	For count data, e.g. number of heads in $n$ coin tosses
Poisson $X \sim Pois(\lambda)$	$\mathbb{P}(X = x) = \frac{\lambda^x e^{-\lambda}}{x!}$ $x \in \{0, 1, \dots\}$	For count data, e.g. number of buses passing a spot
Categorical $X \sim Cat([p_1,, p_k])$	$\mathbb{P}(X = x) = p_x$ $x \in \{1, \dots, k\}$	For picking one of a fixed number of choices
CONTINUOUS RA	NDOM VARIABLES	
Uniform $X \sim U[a, b]$	$pdf(x) = \frac{1}{b - a}$ $x \in [a, b]$	A uniformly-distributed floating point value
Normal / Gaussian $X \sim N(\mu, \sigma^2)$	$pdf(x) = \frac{1}{\sqrt{2\pi\sigma^2}} e^{-(x-\mu)^2/2\sigma^2}$ $x \in \mathbb{R}$	For data about magnitudes, e.g. temperature or height
Pareto $X \sim Pareto(\alpha)$	$pdf(x) = \alpha x^{-(\alpha+1)}$ $x \ge 1$	For data about "cascade" magnitudes, e.g. forest fires
Exponential $X \sim \text{Exp}(\lambda)$	$pdf(x) = \lambda e^{-\lambda x}$ $x > 0$	For waiting times, e.g. time until next bus
Beta $X \sim \text{Beta}(a, b)$	$pdf(x) \propto x^{a-1} (1-x)^{b-1}$ $x \in (0,1)$	Arises in Bayesian inference

#### Exercise 1.3.1 (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

 $X \sim \text{Binom}(n, p)$ 

where p is the probability of heads. Estimate p.

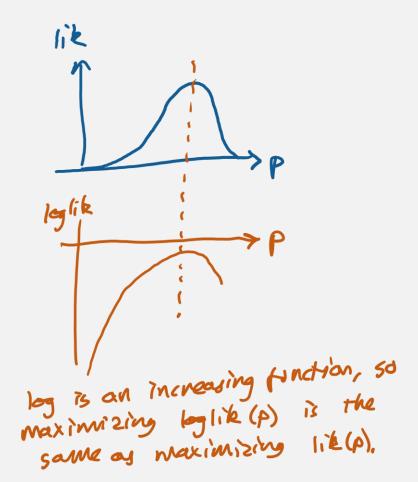
#### Log likelihood of the observed data:

lik = 
$$P(X=x) = {n \choose x} p^x (f-p)^{n-x}$$
  
log lik =  $\log {n \choose x} + x \log p + (n-x) \log f-p$   
corst [it depends on the data x but doesn't depend on p, so for the purposes finding the mle for p it's a constant]

Parameter that maximizes it:

$$\frac{d}{d\rho} \left( \frac{\partial \rho}{\partial \rho} \right) \left( \frac{\partial \rho}{\partial \rho} \right) = 0$$

$$\Rightarrow \frac{\partial \rho}{\partial \rho} \left( \frac{\partial \rho}{\partial \rho} \right) \left( \frac{\partial \rho}{\partial \rho} \right) = 0$$



#### Exercise 1.3.6 (Handling boundaries)

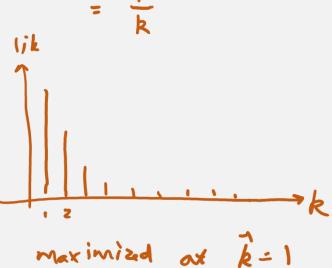
We throw a k-sided dice, and get the answer x=10. Estimate k, using the probability model

$$\mathbb{P}(\text{throw } x) = \frac{1}{k}, \quad x \in \{1, \dots, k\}$$



#### SANITY CHECK

Does our answer depend on the data? In the way we'd expect it to?



But this is daft! How can it be that our estimate of the number of sides (k) doesn't depend on the value we saw (x)? And what even is a 1-sided dice? Dehally, our likelihood function is wrong, because it doesn't capture the constraint in the question, that  $x \in \{1, \dots, k\}$  so  $k > \infty$ .

See betwee notes for a slick way to handle this.

See Ex. 1 93.

For a related mle question.

#### Example sheet 1

Learning with probability models
Data Science—DJW—2025/2026

Some of the questions ask for pseudocode. I suggest you implement your code, and test it using the online tester. There is a notebook with templates for answers and instructions for submission on the course materials webpage.

**Question 1.** Given a dataset  $[x_1, \ldots, x_n]$ , we wish to fit a Poisson distribution. This is a discrete random variable with a single parameter  $\lambda > 0$ , called the rate, and

$$\Pr(x; \lambda) = \frac{\lambda^x e^{-\lambda}}{x!} \quad \text{for } x \in \{0, 1, 2, \dots\}.$$

Show that the maximum likelihood estimator for  $\lambda$  is  $\hat{\lambda} = n^{-1} \sum_{i=1}^{n} x_i$ .

Question 2. Give pseudocode to fit the model of question 1, using scipy.optimize.fmin. [Optional.] If you want to test your code using the online tester, fill in the answer template for PoissonModel.

In practice it answer. But it's problem where yo

Question 3. Gi is unknown. Show

Question 4 (A) B. It has asked y from system A an  $\bar{x} = m^{-1} \sum_{i=1}^{m} x_i$ power of Machine

Suppose the x  $\delta, \sigma^2$ ), and all the estimators for the

#### Hints and comments

**Question 1.** This is a question about fitting parameters from a dataset, as in section 1.3. Formally it's unsupervised learning. There are more examples in section 1.7.

Question 2. See section 1.4. What parameter transform is needed here, to perform a maximization over the restricted domain  $\lambda > 0$ ? How many maxima does the likelihood function have, and how might you choose a sensible starting point for the numerical optimization to make sure it finds a global maximum? Also, if you use numpy, watch out for which variables in your numpy code are vectors and which are scalars.

**Question 3.** This is another question about maximum likelihood estimation on datasets, also known as unsupervised learning, like question 1. You will also need to use the indicator function trick, from section 1.3 exercise 1.3.6

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

i.e. estimate its parameters using Maximum Likelihood Estimation (mle)

The model typically has unknown parameters.

The <u>likelihood</u> is the probability of seeing the data that we actually saw.

It depends on the parameters.

Let's simply pick the parameters that maximize the likelihood!

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

i.e. estimate its parameters using Maximum Likelihood Estimation (mle)

The model typically has unknown parameters.

The <u>likelihood</u> is the probability of seeing the data that we actually saw.

If the data consists of many datapoints  $[x_1, ..., x_n]$  and our model says they're independent, then

$$lik(data) = \prod_{i=1}^{n} lik(x_i)$$

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

i.e. estimate its parameters using Maximum Likelihood Estimation (mle)

The model typically has unknown parameters.

The <u>likelihood</u> is the probability of seeing the data that we actually saw.

 $\mathbb{P}(\text{data})$  if our model is a discrete rand.var. pdf(data) if our model is a continuous rand.var.

If the data consists of many datapoints  $[x_1, ..., x_n]$  and our model says they're independent, then

$$lik(data) = \prod_{i=1}^{n} lik(x_i)$$

#### Exercise 1.3.2 (Exponential sample)

Let the dataset be a list of real numbers,  $x_1, ..., x_n$ , all > 0. Use the probability model that says they're all independent  $\text{Exp}(\lambda)$  random variables, where  $\lambda$  is unknown. Estimate  $\lambda$ .

Log likelihood of the observed data:

$$lik(\frac{dot^{\alpha}}{z_{1}...z_{n}}) = lik(x_{1}) \times ... \times lik(x_{n})$$

$$= (\lambda e^{-\lambda x_{0}}) \times ... \times (\lambda e^{-\lambda x_{n}})$$

$$= \lambda^{n} e^{-\lambda} \xi \times i$$

$$= \lambda^{n} e^{-\lambda} \xi \times i$$

$$log lik = n log \lambda - \lambda \xi \times i$$

WARNING

Watch out copy-pasteitis. Your likelihood must describe the data in the question!

CONTINUOUS RANDOM VARIABLES (real-valued)

Exponential  $X \sim \text{Exp}(\lambda)$ 

$$pdf(x) = \lambda e^{-\lambda x}$$

np.random.exponential(scale= $1/\lambda$ )

Parameter that maximizes it:

$$\frac{d}{d\lambda} \log lik = \frac{n}{\lambda} - \sum_{i} x_{i} = 0 \Rightarrow \hat{\lambda} = \frac{n}{\sum x_{i}}$$

#### Fxercise 1.3.8

Consider a dataset consisting of two collections of real numbers  $x_1, \dots, x_n$ and  $y_1, ..., y_n$  Model the first collection as Normal( $\mu, \sigma^2$ ) and the second as  $Normal(\nu, \sigma^2)$ , where  $\mu, \nu, \sigma$  are all unknown. Estimate  $\sigma$ .

Log likelihood of the observed data:

Log likelihood of the observed data:

$$|g|ik(y) = -\frac{1}{2} |g(2\pi\sigma^2) - \frac{1}{2} \sum_{i=1}^{2} (x_i - y_i)^2 - \frac{1}{2} |g(2\pi\sigma^2) - \frac{1}{2} \sum_{i=1}^{2} (y_i - y_i)^2 - \frac{1}{2} |g(2\pi\sigma^2) - \frac{$$

Parameter that maximizes it:

and all the data?

Tameter that maximizes it:

$$\frac{d}{d\sigma} = 0 \Rightarrow -\frac{2n\sigma}{\sigma^2} + \frac{1}{\sigma^3} (e_1 + e_2) = 0 \Rightarrow \hat{\sigma} = \sqrt{\frac{1}{2n}} (e_1 + e_2)$$

Appendion in the depends on y

Appendion is useless

as an estimator.

SANITY CHECK

Does my answer depend only on the data, or also on unknown parameters?

When there are multiple unknowns, we have to find them all, oven it we've only interested in one:

#### Exercise 1.3.8

Consider a dataset consisting of two collections of real numbers,  $x_1, ..., x_n$  and  $y_1, ..., y_n$ . Model the first collection as  $\operatorname{Normal}(\mu, \sigma^2)$  and the second as  $\operatorname{Normal}(\nu, \sigma^2)$ , where  $\mu, \nu, \sigma$  are all unknown. Estimate  $\mu - \nu$ .

When we fit 
$$p_i, v_i, v_i = v_i = v_i = v_i$$
 when we fit  $p_i, v_i, v_i = v_i = v_i$  when we fit  $p_i, v_i, v_i = v_i = v_i = v_i$  and I want to estimate  $\phi = f(\theta)$ , then the mle for  $\phi$  is  $f(\hat{\theta})$ 

## There are standard numerical random variables that you should know:

DISCRETE RANDO	M VARIABLES	
Binomial $X \sim Bin(n, p)$	$\mathbb{P}(X = x) = \binom{n}{x} p^x (1 - p)^{n - x}$ $x \in \{0, 1, \dots, n\}$	For count data, e.g. number of heads in $n$ coin tosses
Poisson $X \sim Pois(\lambda)$	$\mathbb{P}(X = x) = \frac{\lambda^x e^{-\lambda}}{x!}$ $x \in \{0, 1, \dots\}$	For count data, e.g. number of buses passing a spot
Categorical $X \sim Cat([p_1,, p_k])$	$\mathbb{P}(X = x) = p_x$ $x \in \{1, \dots, k\}$	For picking one of a fixed number of choices
CONTINUOUS RA	NDOM VARIABLES	
Uniform $X \sim U[a, b]$	$pdf(x) = \frac{1}{b - a}$ $x \in [a, b]$	A uniformly-distributed floating point value
Normal / Gaussian $X \sim N(\mu, \sigma^2)$	$pdf(x) = \frac{1}{\sqrt{2\pi\sigma^2}} e^{-(x-\mu)^2/2\sigma^2}$ $x \in \mathbb{R}$	For data about magnitudes, e.g. temperature or height
Pareto $X \sim Pareto(\alpha)$	$pdf(x) = \alpha x^{-(\alpha+1)}$ $x \ge 1$	For data about "cascade" magnitudes, e.g. forest fires
Exponential $X \sim \text{Exp}(\lambda)$	$pdf(x) = \lambda e^{-\lambda x}$ $x > 0$	For waiting times, e.g. time until next bus
Beta $X \sim \text{Beta}(a, b)$	$pdf(x) \propto x^{a-1} (1-x)^{b-1}$ $x \in (0,1)$	Arises in Bayesian inference

## There are standard numerical random variables that you should know:

 $a + b N(0,1) \sim a + N(0, b^2) \sim N(a, b^2)$ for constants a and b

Useful properties of the Normal distribution:

- If we rescale a Normal, we get a Normal
- If we add <u>independent</u> Normals, we get a Normal ~ N(M, T2) + N(7, p2) ~ N(M+0, 52+p2) ((ike the Poisson, binomial)

These two properties are known as "linearity of the Normal distribution".

Normal / Gaussian 
$$pdf(x) = \frac{1}{\sqrt{2\pi\sigma^2}}e^{-(x-\mu)^2/2\sigma^2}$$
 For data about magnitudes, e.g. temperature or height  $x \in \mathbb{R}$ 

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

i.e. estimate its parameters using Maximum Likelihood Estimation (mle)

The model typically has unknown parameters.

The <u>likelihood</u> is the probability of seeing the data that we actually saw.

 $\mathbb{P}(\text{data})$  if our model is a discrete rand.var. pdf(data) if our model is a continuous rand.var.

If the data consists of many datapoints  $[x_1, ..., x_n]$  and our model says they're independent, then

$$lik(data) = \prod_{i=1}^{n} lik(x_i)$$

When there are multiple unknown parameters, we must maximize over all of them simultaneously (even if we're only interested in one).

# §1.4 Numerical optimization

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

i.e. estimate its parameters using Maximum Likelihood Estimation (mle)

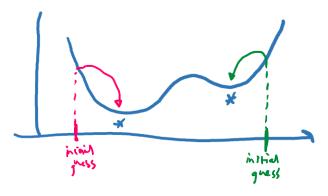
§:

with numerical optimization

(since the likelihood function is usually far too complex for exact optimization)

```
a local minimum
Numerical optimization with Python / scipy
To find the minimum<sup>†</sup> of a smooth function f: \mathbb{R}^K \to \mathbb{R},
      import scipy.optimize
     def f(x):
            return ...
    x_0 = [...] # initial guess

\hat{x} = \text{scipy.optimize.fmin}(f, x_0)
```

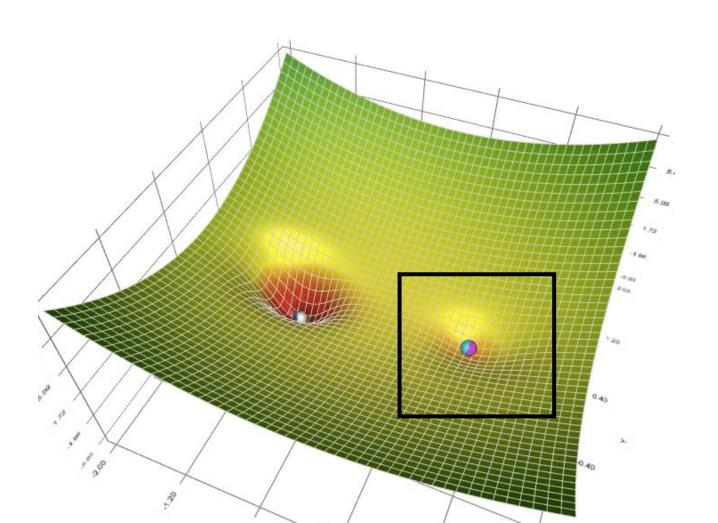


The initial guess will influence which local minimum the fmin ends up finding.

† There is no scipy.optimize.fmax. To maximize f, scipy.optimize.fmin(lambda x: -f(x),  $x_0$ )

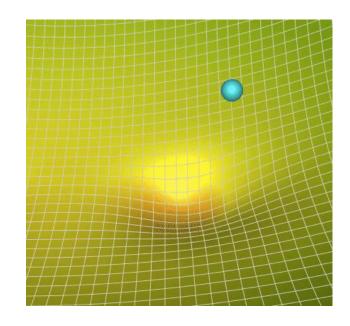
# How does numerical optimization work?

Animations by Lili Jiang, <u>Towards Data Science</u>



#### **GRADIENT DESCENT**

Find the gradient of the function, and take a step in the direction of steepest descent



```
Exercise 1.4.2 (Constraints / softmax transformation) Find the maximum of f(p_1,p_2,p_3) = 0.2\log p_1 + 0.5\log p_2 + 0.3\log p_3 over p_1,p_2,p_3 \in (0,1) such that p_1+p_2+p_3=1.
```

```
fmin is not able to handle constraints.

(only trick: instead of maximizing over Property \in (0,1) st. p_1 + p_2 + p_3 = 1

(only trick: instead of maximizing over (s_1, s_2, s_3) \in \mathbb{R}^3, and set p_i = \frac{e^{s_i}}{e^{s_1} + e^{s_2} + e^{s_3}}.
```

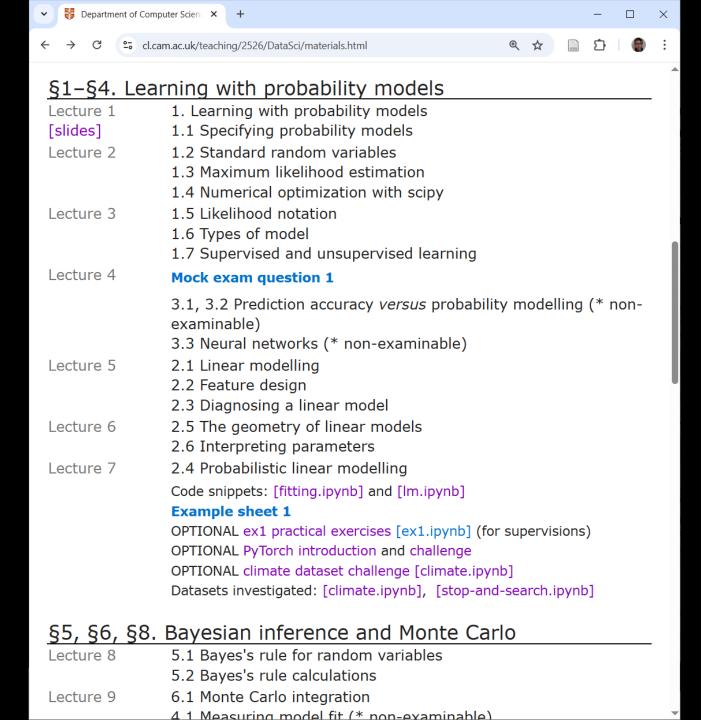
This parameter transformation trick works because (i) every  $(s_1, s_2, s_3)$  yields a valid p, and (ii) every valid  $(p_1, p_2, p_3)$  can be achieved by some s.

```
def f(p):
    p<sub>1</sub>, p<sub>2</sub>, p<sub>3</sub> = p
    return 0.2*np.log(p<sub>1</sub>) + 0.5*np.log(p<sub>2</sub>) + 0.3*np.log(p<sub>3</sub>)

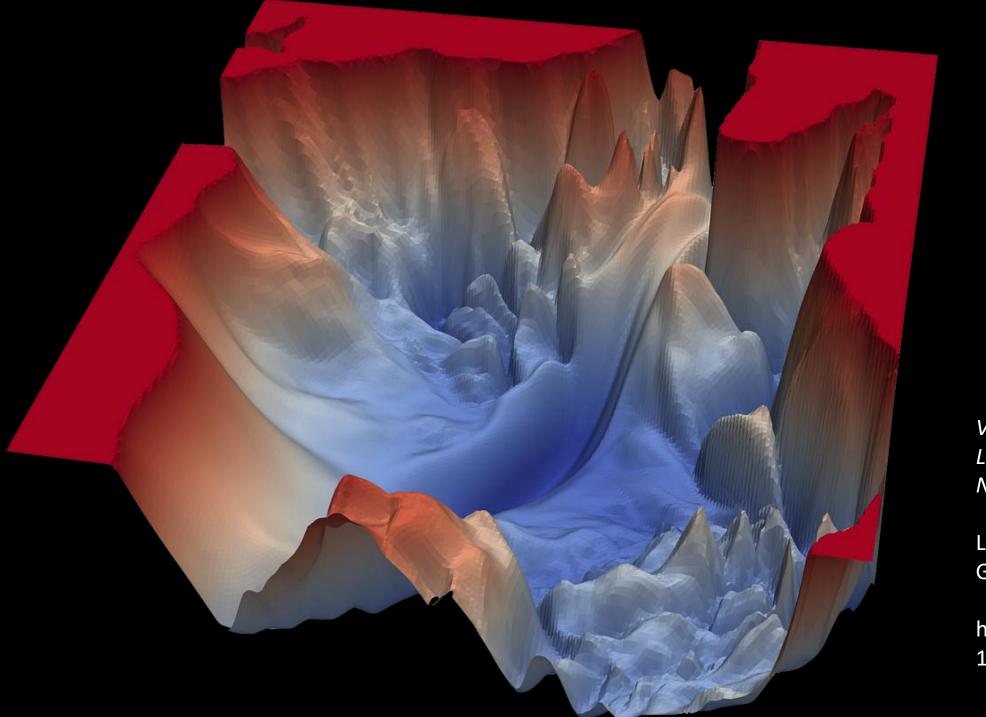
def softmax(s):
    p = np.exp(s)
    return p / np.sum(p)

$\hat{s}$ = scipy.optimize.fmin(lambda s: -f(softmax(s)), [0,0,0])
$\hat{s}$ = softmax(\hat{s})

Optimization terminated successfully. Current function value: 1.02965. Iterations: 63.
Function evaluations: 120
    array([0.19999474, 0.49999912, 0.300000614])
```



- CODE SNIPPETS from lectures
- CODING EXERCISES on example sheets



Visualizing the Loss Landscape of Neural Nets

Li, Xu, Taylor, Studer, Goldstein (2018)

https://arxiv.org/abs/ 1712.09913

### Getting started with PyTorch

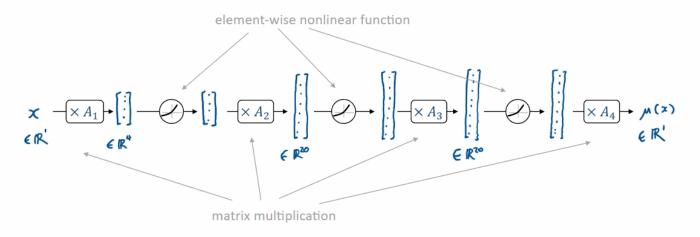
This is a brief introduction to PyTorch, designed to complement the IB Data Science course. It assumes you're familiar with the idea of maximum likelihood estimation. We'll use PyTorch to represent a probability model for regression, and fit it.

There are many other tutorials on PyTorch, including <u>the tutorial in the official documentation</u>. They typically present PyTorch as a software library, and go into much more depth on tensors and GPUs and so on. But they often don't give much guidance on how to use PyTorch for data science.

In this tutorial we'll work with the data behind the <u>xkcd 2048 comic on curve-fitting</u>. First, load the dataset:

#### STEP 7. Using a neural network

A neural network is just another function! We can swap out the  $\mu(x) = a + bx + cx^2$  function and replace it by a neural network, i.e. a sequence of linear maps and nonlinear element-wise operations.



# Exercise The observed data is $[\text{temp}_1,...,\text{temp}_n]$ . Estimate $\gamma$ using the model $\text{Temp}_i \sim c + \alpha \sin(2\pi(t_i + \varphi)) + \gamma t_i + \text{Normal}(0,\sigma^2), \quad i \in \{1,...,n\}$

Temp: ~ N (predi, 
$$\sigma^2$$
) where predi =  $c + \alpha \sin(2\pi(t_i, \sigma^2)) + \delta t_i$   
by linearity of the Normal distribution: pred + N(0, $\sigma^2$ ) ~ N(pred, $\sigma^2$ )

log sik = 
$$-\frac{9}{2}\log(2\pi\sigma^2)$$
  $-\frac{1}{2\sigma^2}$   $\sum_{i} (temp_i - pred_i)^2$ .

by looking up the pdf, and avoiding copy-paste-itis

```
Tempi ~ N (predi, \sigma^2) where predi = c + \alpha \sin(2\pi (t_i, \sigma^2)) + \delta t_i

In the = -\frac{2}{2} \log(2\pi\sigma^2) - \frac{1}{2\sigma^2} \sum (tempi - predi)^2
```

```
url = 'https://www.cl.cam.ac.uk/teaching/current/DataSci/data/climate 202510.csv'
climate = pandas.read csv(url)
climate['t'] = (climate.yyyy + (climate.mm-1)/12)
climate['temp'] = (climate.tmin + climate.tmax) / 2
df = climate.loc[(climate.station=='Cambridge') & (climate.yyyy>=1985)]
# Transform parameters: Let \sigma = e^{\tau} for \tau \in \mathbb{R}
def loglik(\theta,t,temp):
    (c,\alpha,\phi,\gamma,\tau) = \theta
    \sigma 2 = np.exp(\tau) ** 2
    pred = c + \alpha * np.sin(2*\pi*(t+\phi)) + \gamma*t
    n = len(t)
    return - n/2 * np.log(2*\pi*\sigma2) - 1/(2*\sigma2) * np.sum((temp - pred)**2)
init_guess = [12, 7, -0.2, 0, np.log(1.5)]
\thetahat = scipy.optimize.fmin(lambda p: -loglik(df.t,df.temp,\theta), init guess, maxiter=5000)
# View the fitted parameters, compared to initial guess
pandas.DataFrame(\{'par': ['c', '\alpha', '\phi', '\gamma', '\tau', '\sigma'], 
                     'init guess': list(init_guess) + [np.exp(init_guess[-1])],
                     'fitted': list(phat) + [np.exp(phat[-1])]})
```

The fitted value for c comes out to be -67.5, which is surprising compared to our guess of I2. But when we plot the fitted predictions, it looks right. Can you figure out why we got c = -67.5?