

A wide-angle photograph of Mount Bromo volcano at sunset. The volcano's peak is dark and shadowed, while its slopes are illuminated with vibrant orange and yellow hues from the setting sun. In the foreground, there are some green trees and shrubs. The background shows more volcanic peaks under a clear sky.

# Basic Machine Learning: Lecture 1

## Machine Learning Summer School, Indonesia 2020

Daniel Worrall

Gunung Bromo National Park

## Today: Recap of “basics”

You are a diverse class: Some will find this easy, some will find this difficult

Do not worry if you find this hard!

The intention is not for you to understand today but to have exposure.

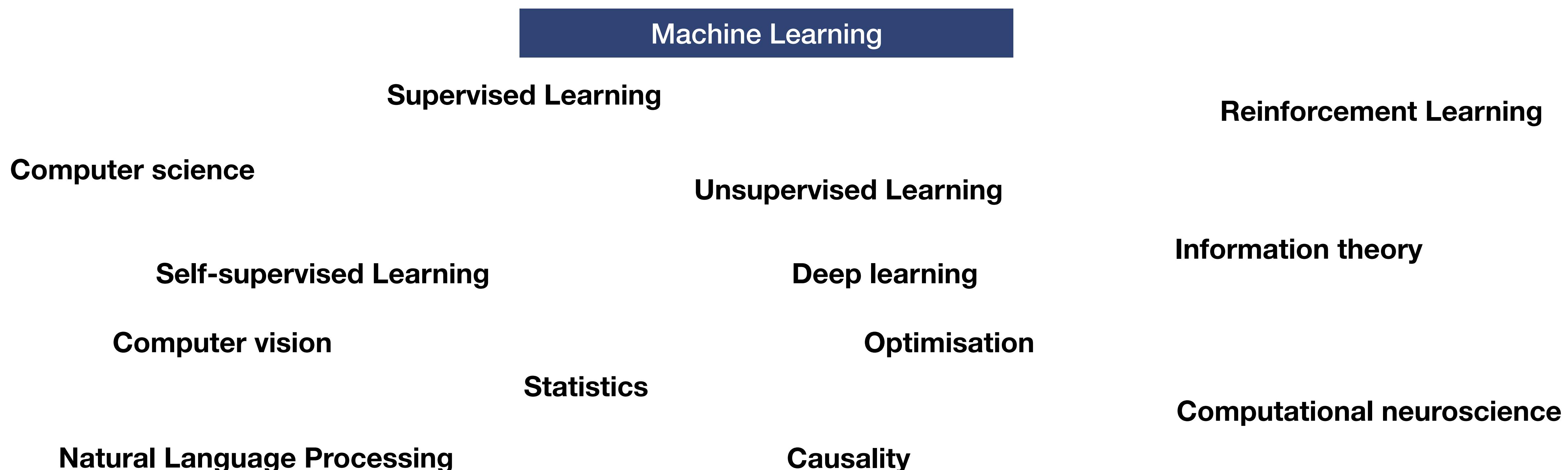
## This lecture: Machine Learning Basics

- What is Machine Learning?
- Probability Theory recap
- Modeling paradigms: Probability theory
- Generative models
- Statistics: Maximum likelihood
- Bayesian Inference
- Prediction
- Conjugacy
- Model Comparison
- Polynomial Regression

# What is Machine Learning?

“Machine learning is the latest in a long line of attempts to distill human knowledge and reasoning into a form that is suitable for constructing machines and engineering automated systems.”

— *Mathematics for machine learning*



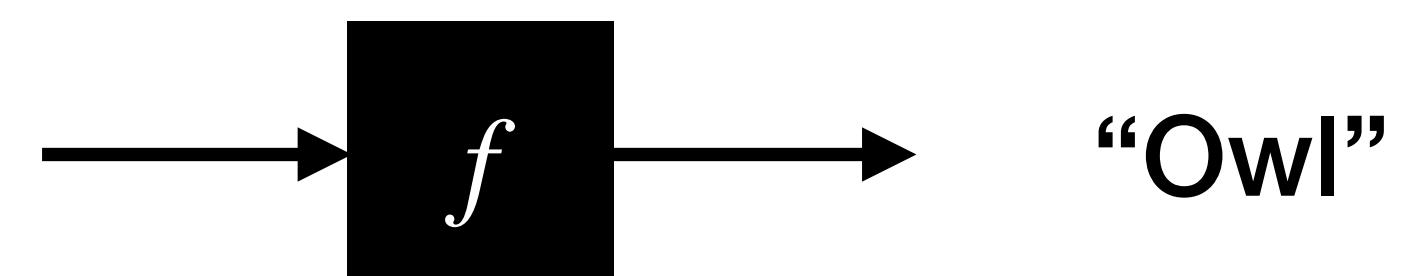
# Supervised learning

We focus on one area of machine learning called *supervised learning*.

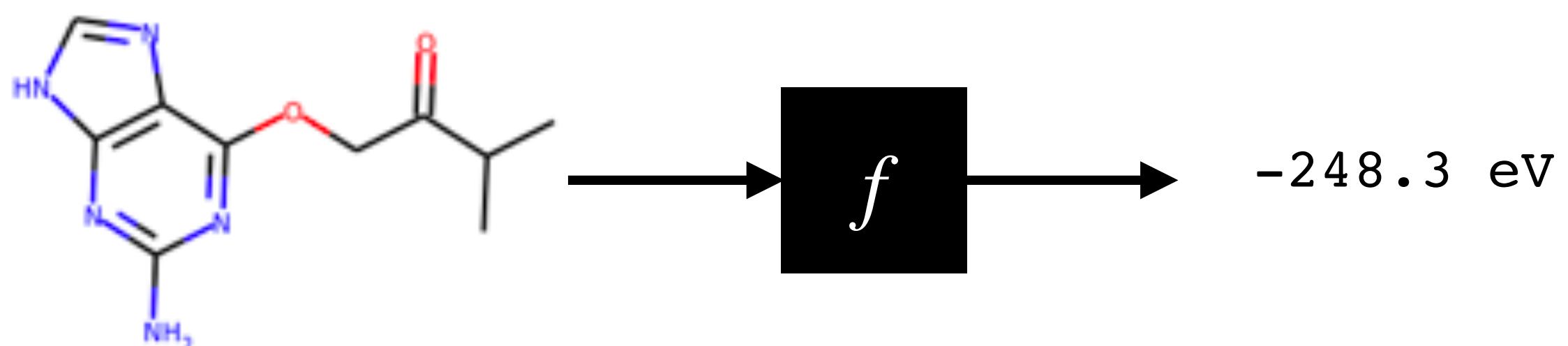
**Supervised learning** problems have 3 parts

- **Data:** inputs  $\mathbf{x}$  and outputs  $\mathbf{y}$
- **Model space:** a collection of *models*  $\mathcal{M}$  which convert inputs into outputs
- **Algorithm:** a method to choose the best model  $m \in \mathcal{M}$  from the model space<sup>1</sup>, which best *fits* the data

**Classification** e.g. image recognition



**Regression** e.g. molecular property prediction

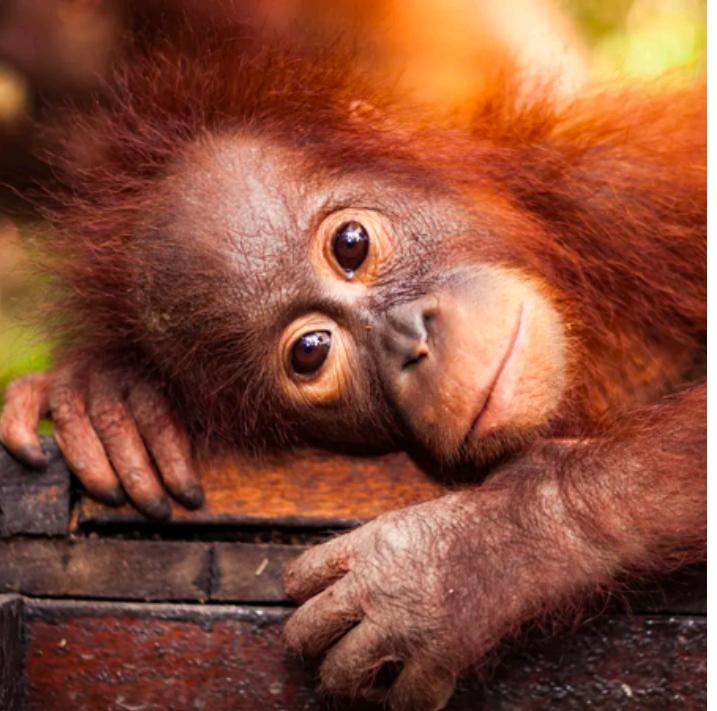
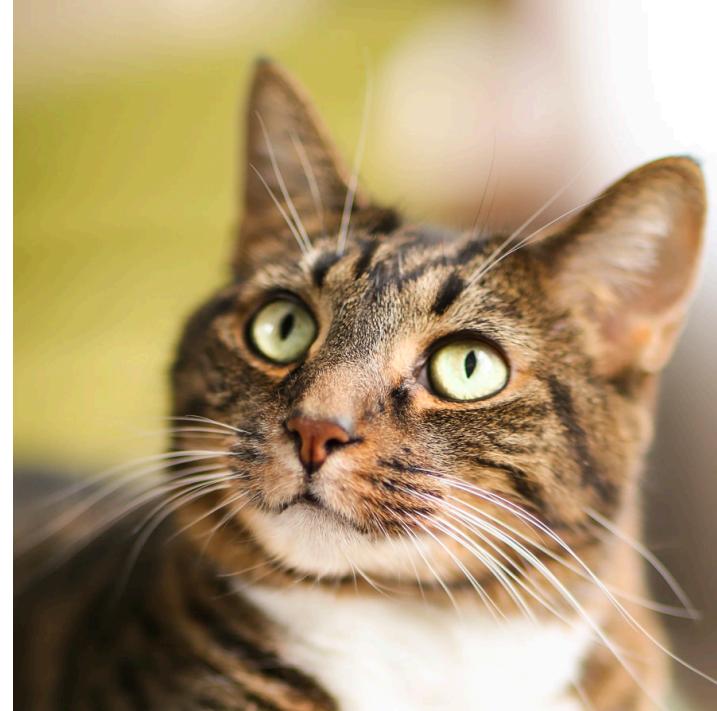
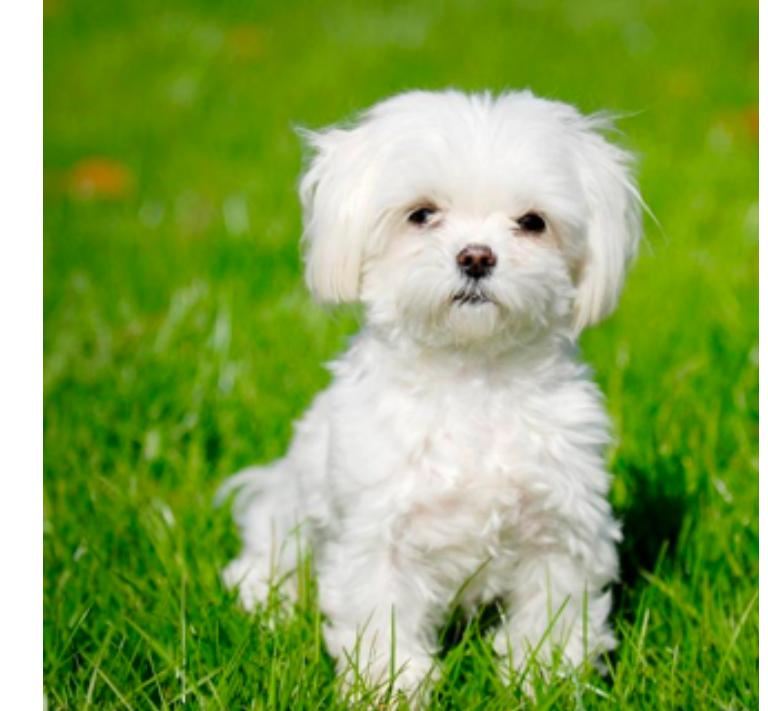


<sup>1</sup> The notation ‘ $\in$ ’ is pronounced *in*, so  $m \in \mathcal{M}$  is spoken ‘m in M’

# Example I - Image Classification

## Image classification

- **Inputs:** 256x256 pixel RGB images
- **Outputs:** labels in {‘dog’, ‘cat’, ‘orangutan’}



- **Model:** ?
  - **Suggestion:** Collect examples of images  $\{x_1, x_2, \dots\}$  with labels  $\{y_1, y_2, \dots\}$  and build a *lookup table*. If new input  $\mathbf{x}_* = \mathbf{x}_j$ , where  $\mathbf{x}_j \in \{x_1, x_2, \dots\}$  then its *predicted label* is  $\mathbf{y}_* = \mathbf{y}_j$ .
  - **Problem 1:** What if we have never seen  $\mathbf{x}_*$  before?
  - **Problem 2:** What if we feed in something which is not a dog/cat/orangutan?
  - **Problem 3:** Are there ways to quantify how good our model is?

# Example II - Machine Translation

'Orang memanggil aku: Minke. Namaku sendiri ... sementara ini tak perlu kusebutkan. Bukan karena gila misteri. Telah aku timbang: belum perlu benar tampilkan diri dihadapan mata orang lain.'

—Pramoedya Ananta Toer (~1980).



'People called me Minke. My own name ... for the time being I need not tell it. Not because I'm crazy for mystery. I've thought about it quite a lot: I don't yet really need to reveal who I am before the eyes of others.'

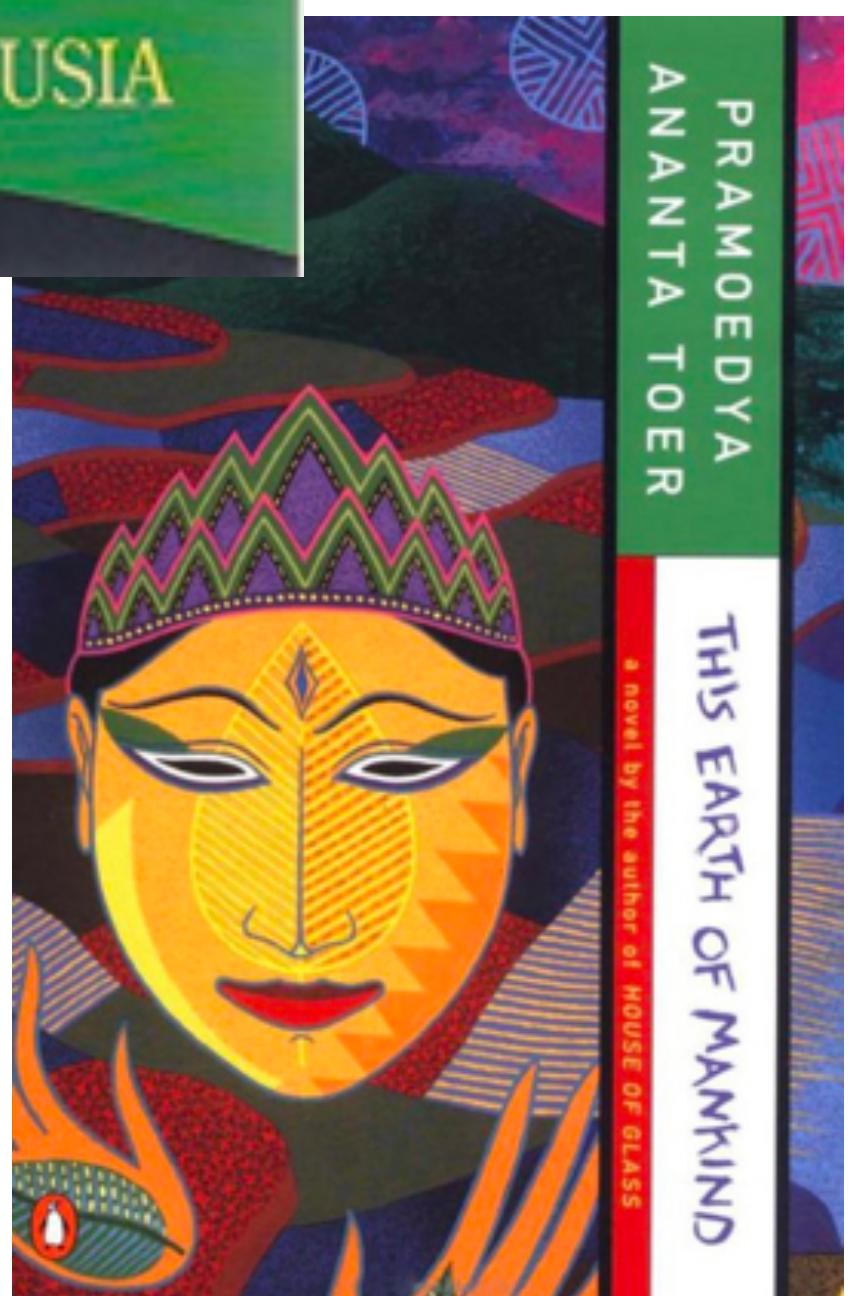
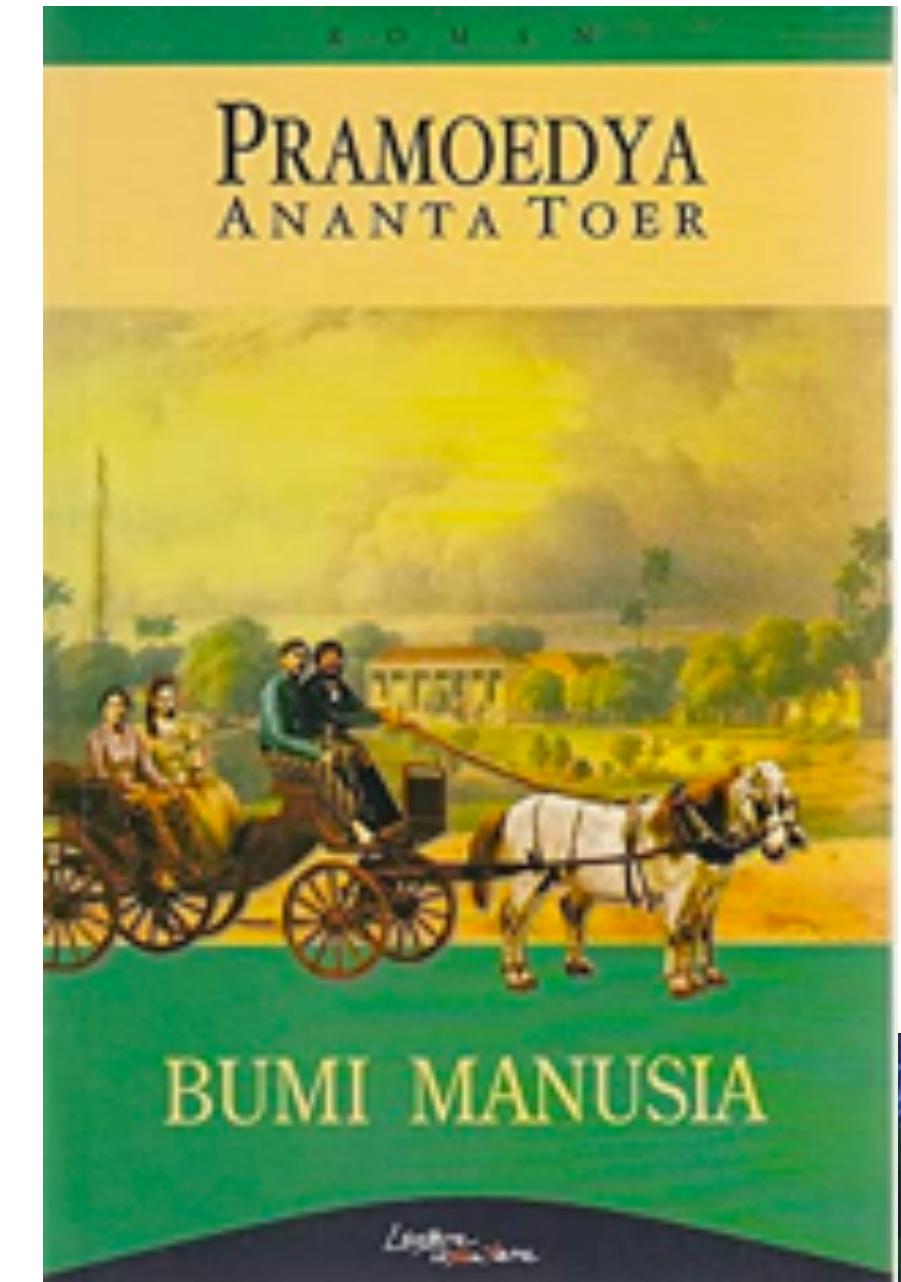
—Max Lane (1981)

'People called me Minke. As for my real name ... for now it doesn't need to be mentioned. Not because I need mystery. I have weighed it up: I needn't yet reveal myself before the eyes of others.'

—Daniel Worrall (2020)

'People call me: Minke. My own name ... meanwhile I don't need to mention it. Not because of a mystery mad. I have weighed: do not need to properly present yourself before the eyes of others.'

—Google Translate (2020)



# Example II - Machine Translation

Machine translation is an input–output task

- **Inputs:** variable length Indonesian strings
- **Outputs:** variable length English strings
- **Problem 1:** Each word has multiple translations
- **Problem 2:** We cannot possibly collect all input–output pairs
- **Problem 3:** What even is a good translation?

Many of the problems we have seen can be addressed (to some extent) by thinking *probabilistically*. To understand what this means though, we need to learn what probability is.

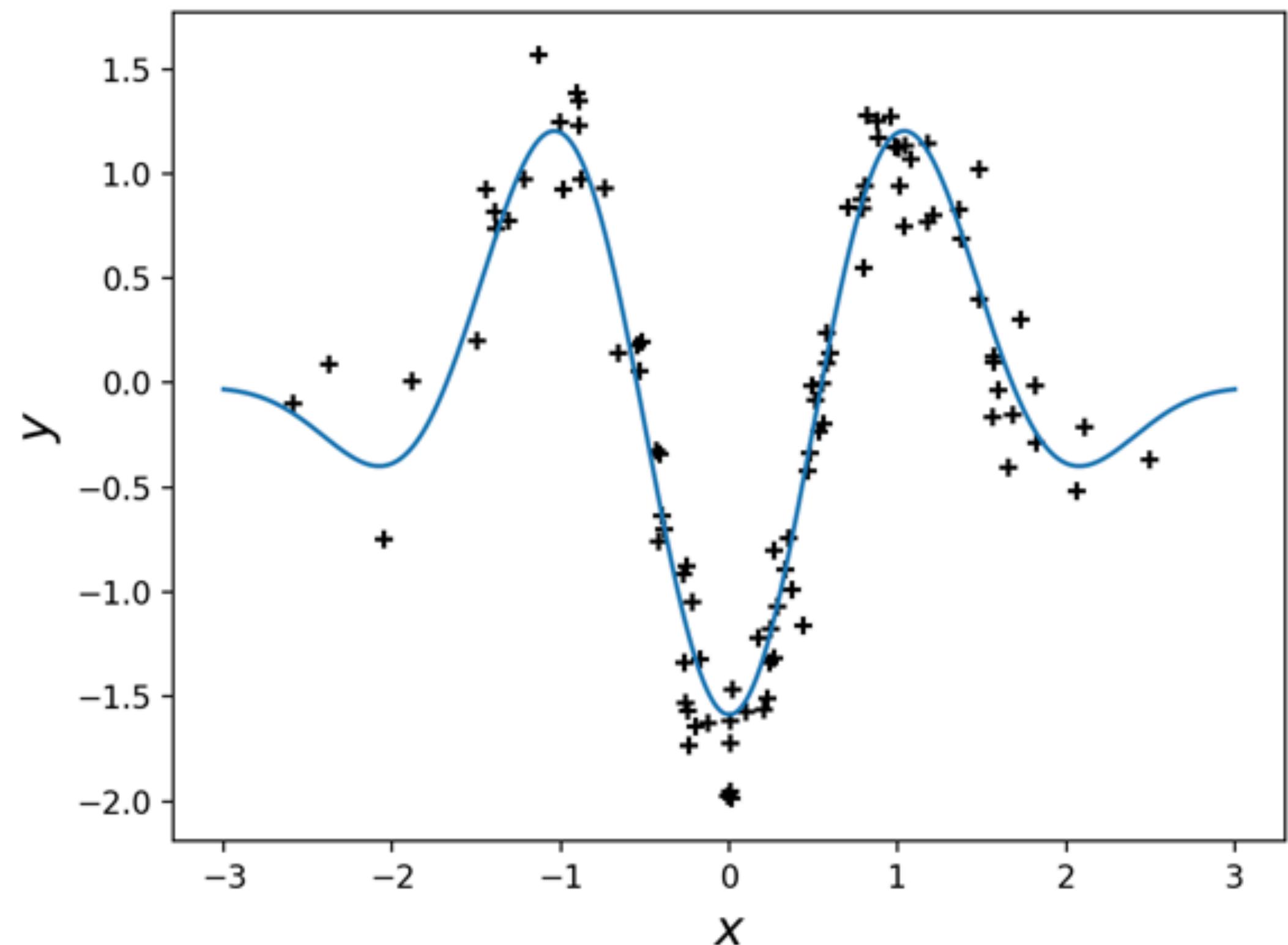
# Example III - Regression

Regression is the canonical input–output task

**Inputs:** real numbers  $x_i \in \mathbb{R}$

**Outputs:** real numbers  $y_i \in \mathbb{R}$

- **Problem 1:** How to handle residual error?
- **Problem 2:** Is a linear model the best we can do?
- **Problem 3:** What about higher dimensions?



The previous two examples are just variants to of regression (in a very liberal sense).

Machine learning is primarily a **conceptual** discipline.

Machine learning is automated machine building\*

Machine learning is...

Computational machines need:

Computers

Programming

**Mathematics: Probability, calculus, linear algebra**

**Mathematics is the language of data**



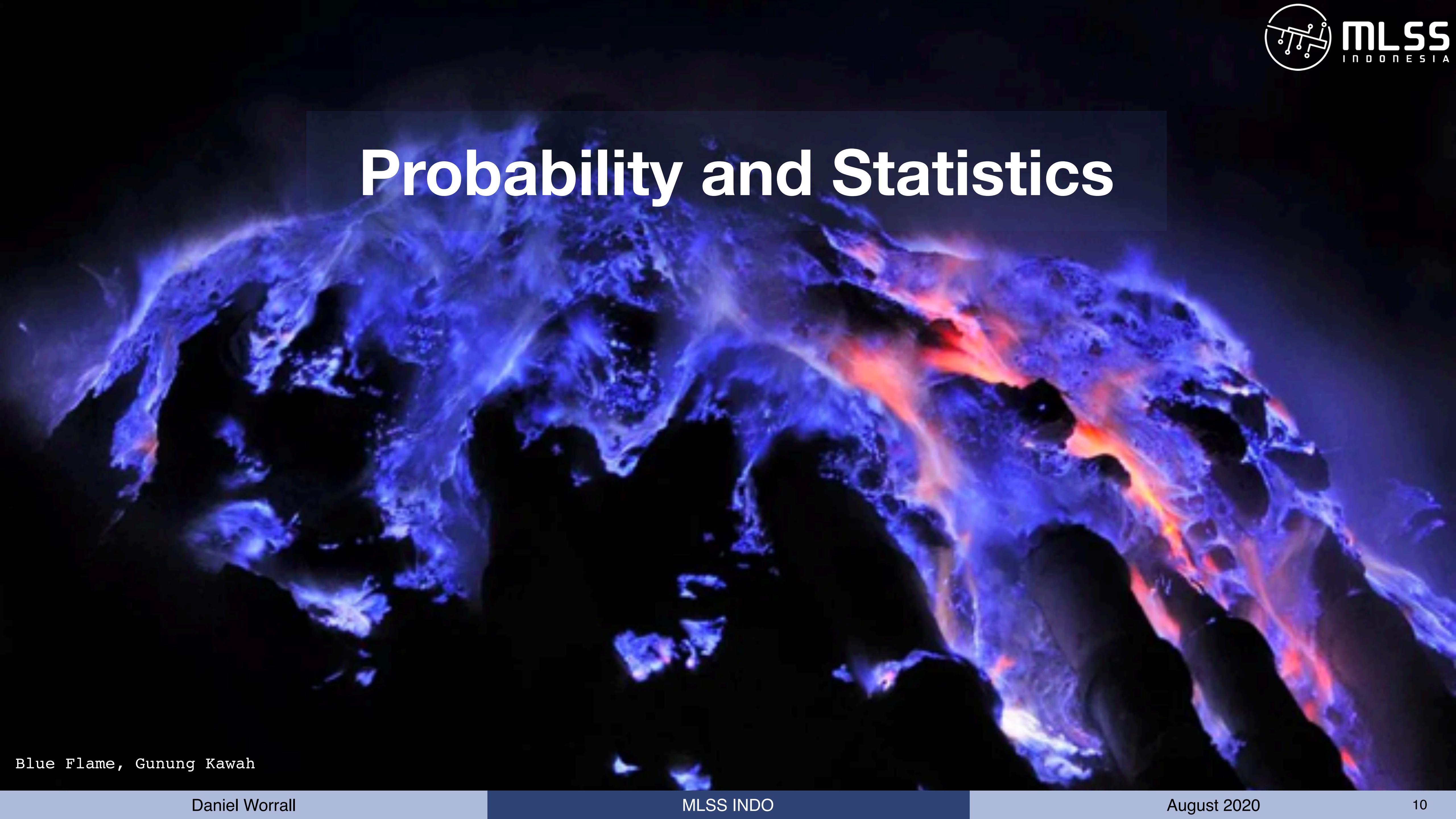
XKCD

Beware: mathematics is just one aspect of ML

Societal awareness: Fairness, Ethics, Decoloniality, ...

\* Some would argue that it is the automation of the scientific method.

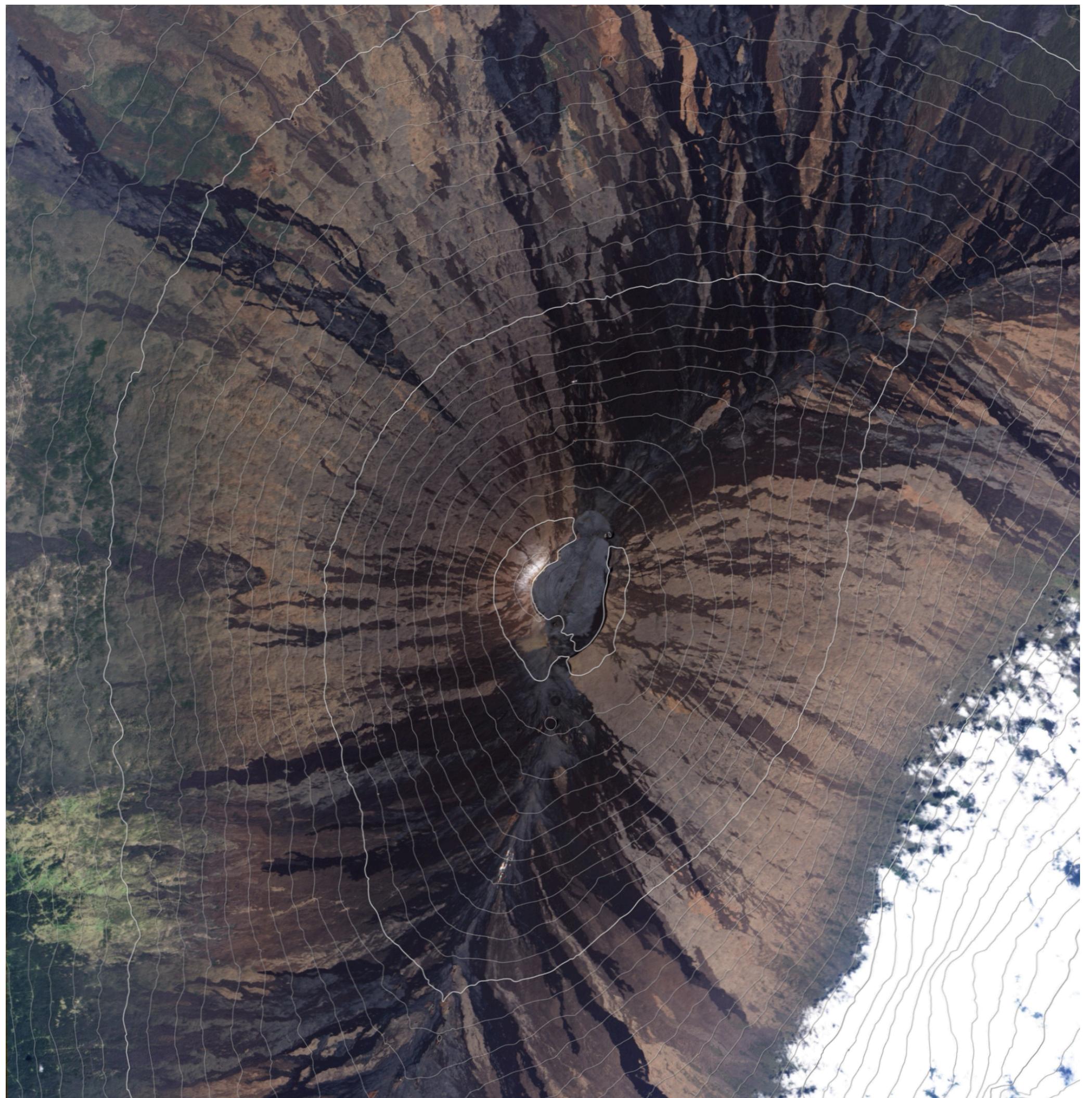
# Probability and Statistics



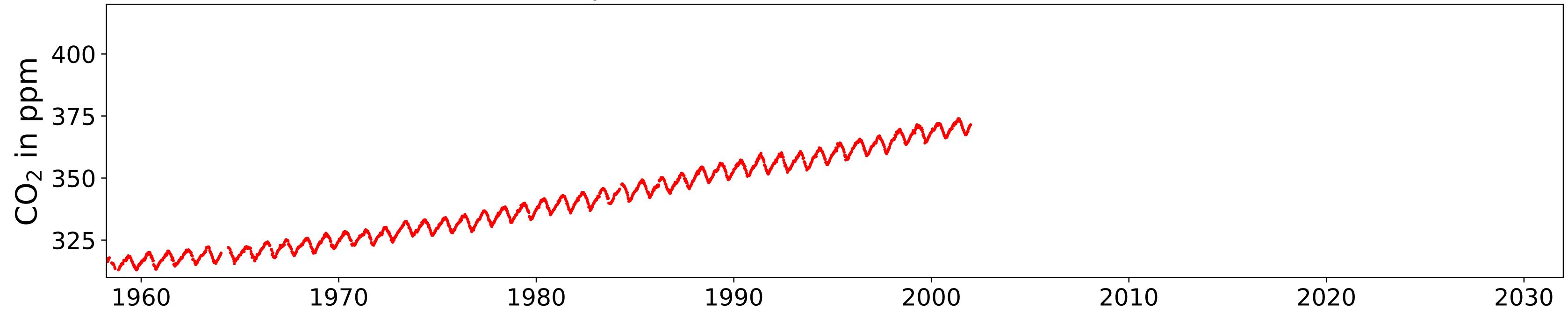
Blue Flame, Gunung Kawah

# Mauna Loa

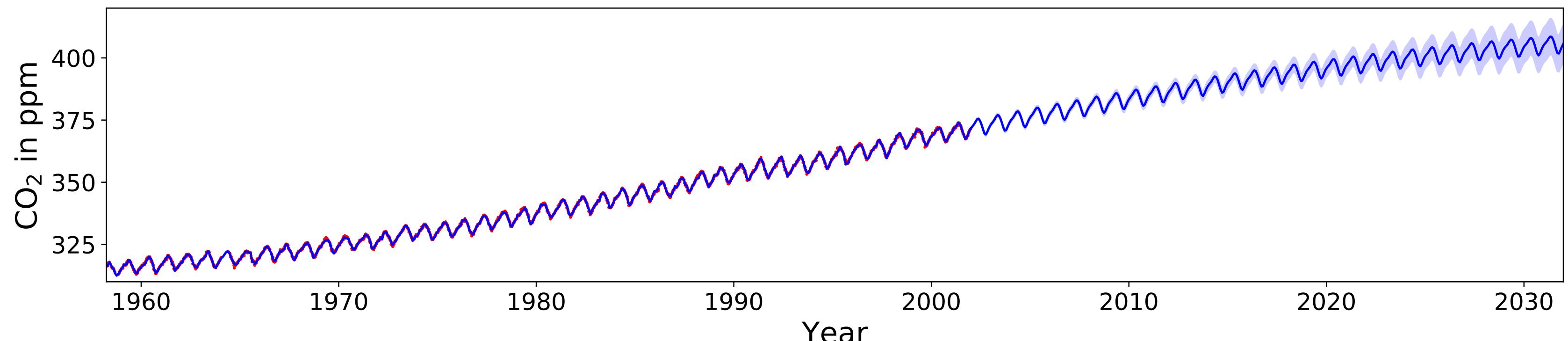
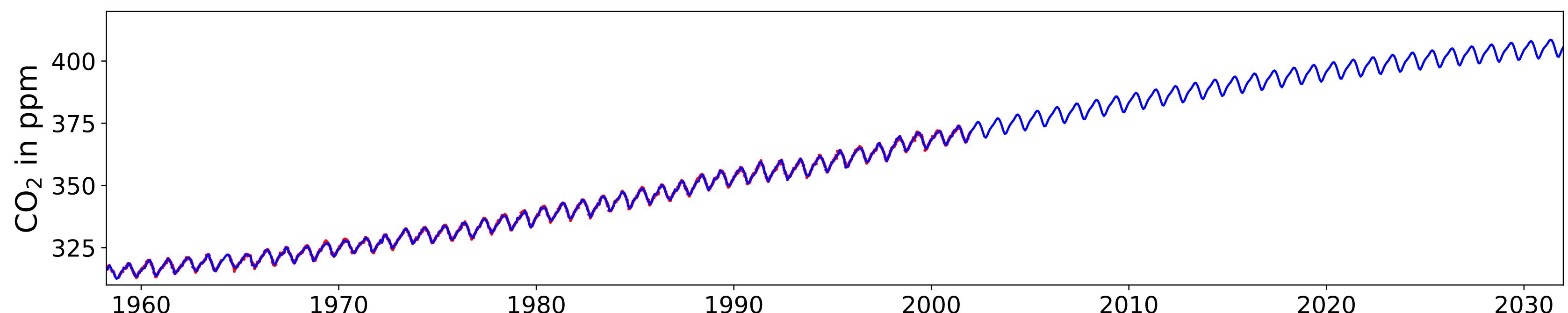
*Mauna Loa* is one of 5 volcanoes forming the Island of Hawaii in the U.S. state of Hawaii in the Pacific Ocean. The largest subaerial volcano in both mass and volume, Mauna Loa has historically been considered the largest volcano on Earth, dwarfed only by Tamu Massif.



Atmospheric CO<sub>2</sub> concentration at Mauna Loa



In machine learning we use past data to make predictions about the future



How would we model the Mauna Loa CO<sub>2</sub> levels as a function of time? Perhaps

$$y = w_1 x + w_2$$

Well this is a straight line, we need an extra periodic component, so how about

$$y = w_1 x + w_2 \cos(2\pi x + w_3) + w_4$$

But this is no good, because  $y$  is negative for certain values of  $x$ ...

## The reality of modelling

“All models are wrong, but some are useful.”

—George Box, 1976

- How do we choose a model?
- Does it matter that we will no doubt make mistakes?
- How do I tell if one model is better than another?

**probability: a mathematical formalism describing uncertain events**

**statistics: the science of collecting and analysing data**

—Carl Rasmussen

*Bayesian statistics* is a branch of statistics loved by machine learners for its computational nature.

- Why is it useful?
- What can we say about uncertain events?
- What be measured?

# Intuitive Probability

Take a coin. Label heads with 1 and tails with 0. Now flip the coin  $N$  times and take the average. Now do this again multiple times.



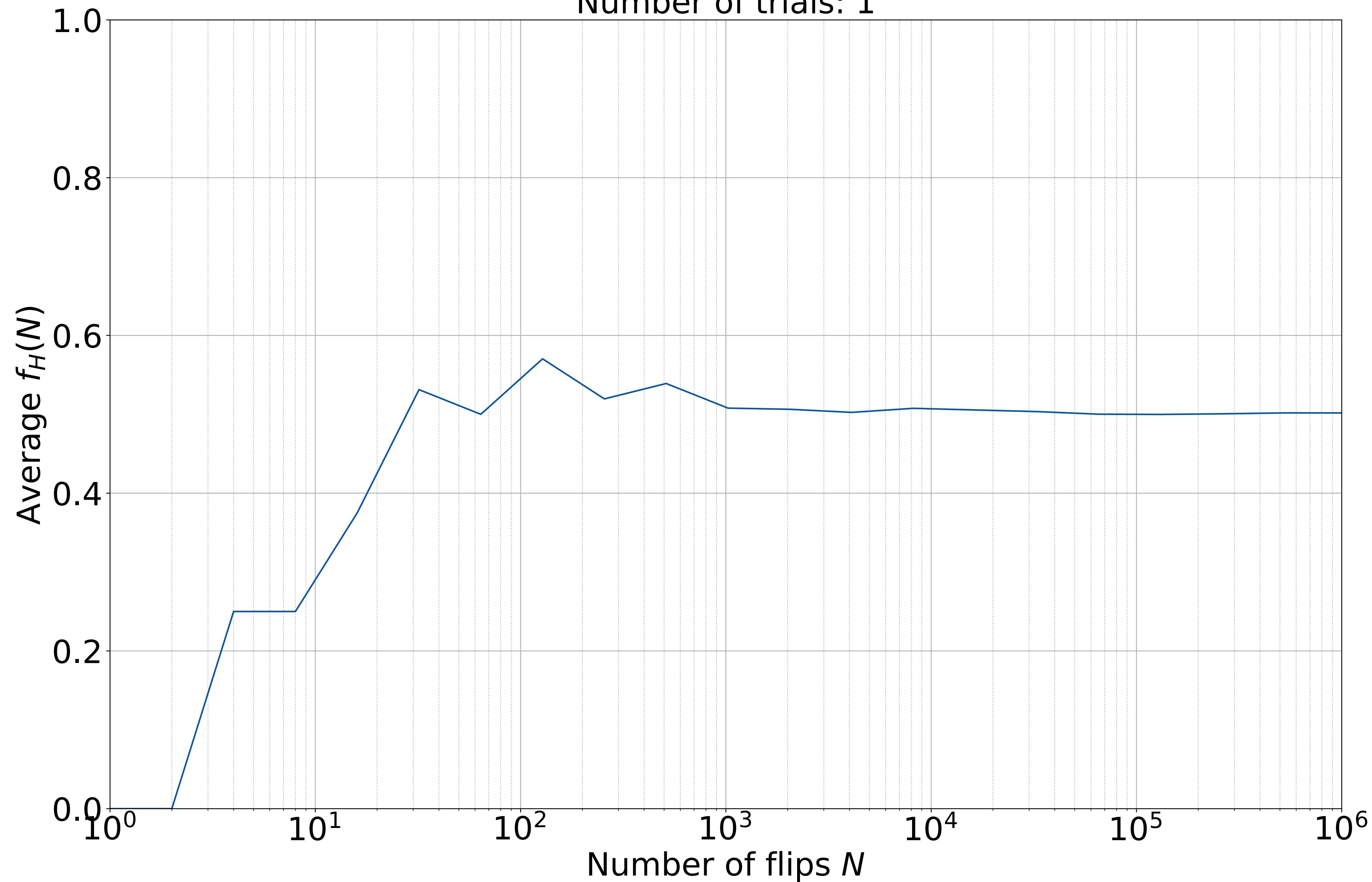
	Trial 1	Trial 2	Trial 3	Trial 4	Trial 5
$N = 10$	0.5000	0.8000	0.6000	0.6000	0.2000
$N = 100$	0.4800	0.4800	0.4800	0.5400	0.5400
$N = 1000$	0.4950	0.5130	0.5080	0.5080	0.4850
$N = 10000$	0.4967	0.5031	0.4980	0.4988	0.4934

Despite the fact that in each trial we get a different result, there is a trend!

As  $N \rightarrow \infty$ , what do you think will happen?

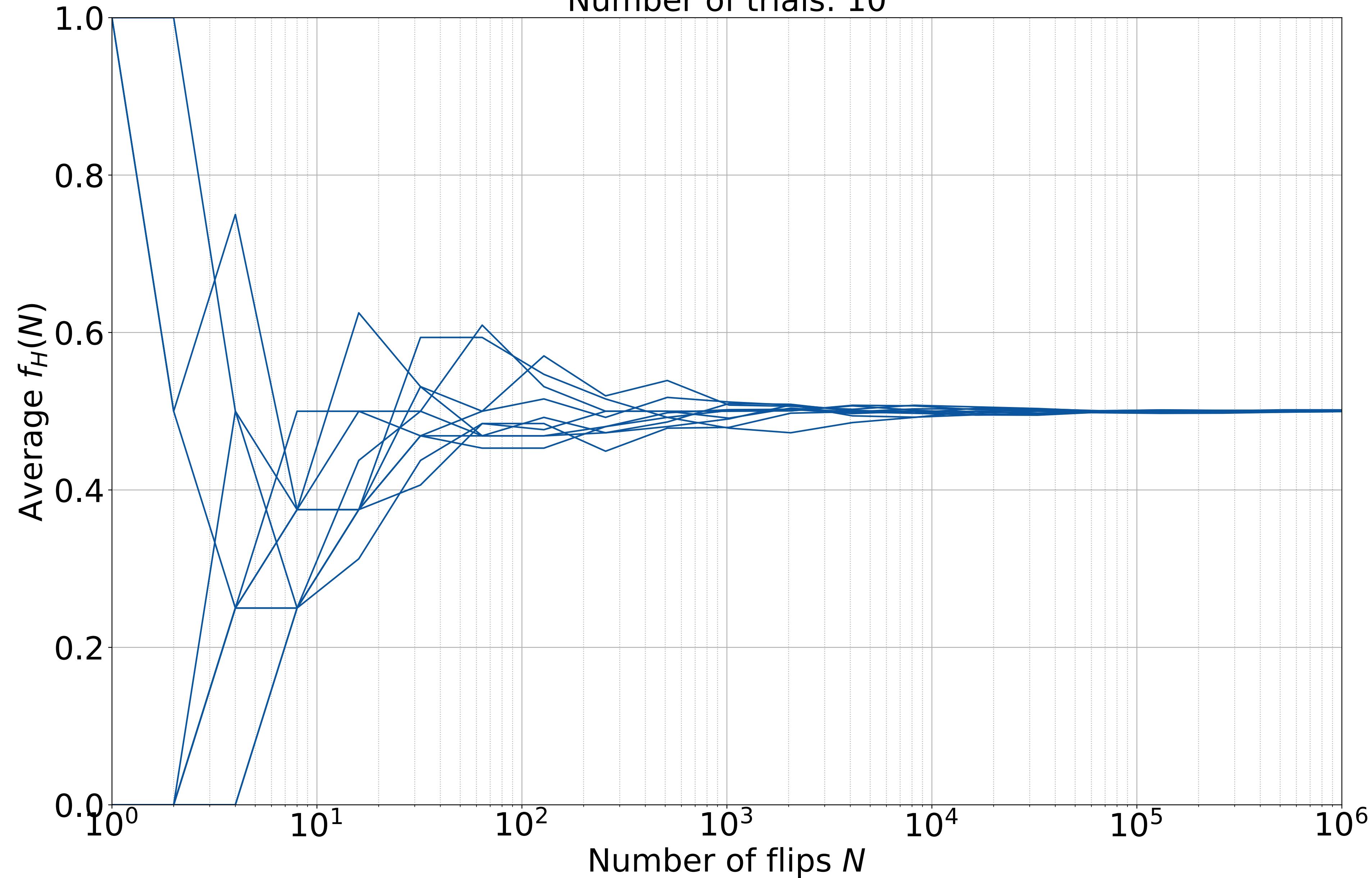
# Frequentist Probability

Number of trials: 1



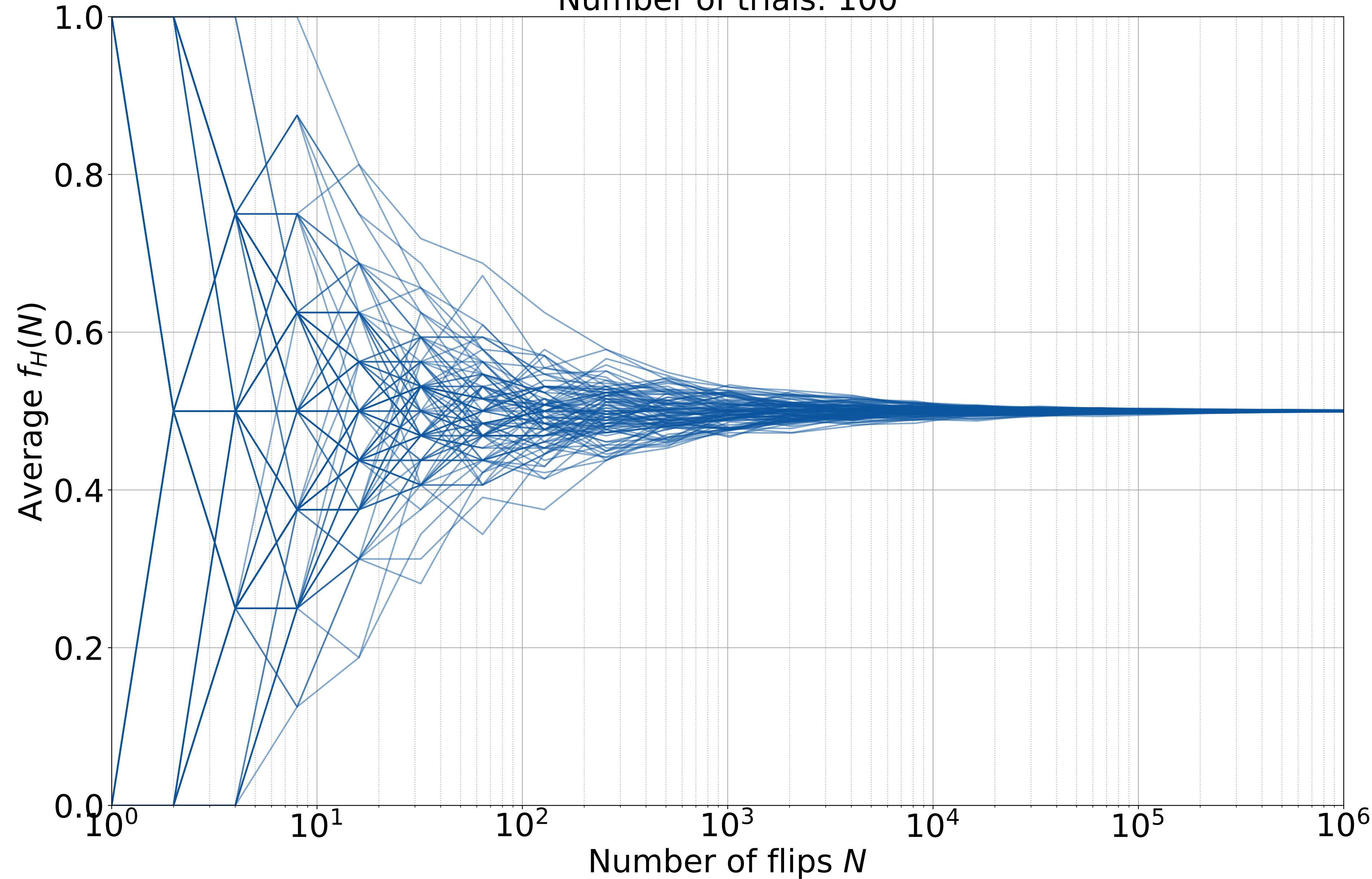
# Frequentist Probability

Number of trials: 10



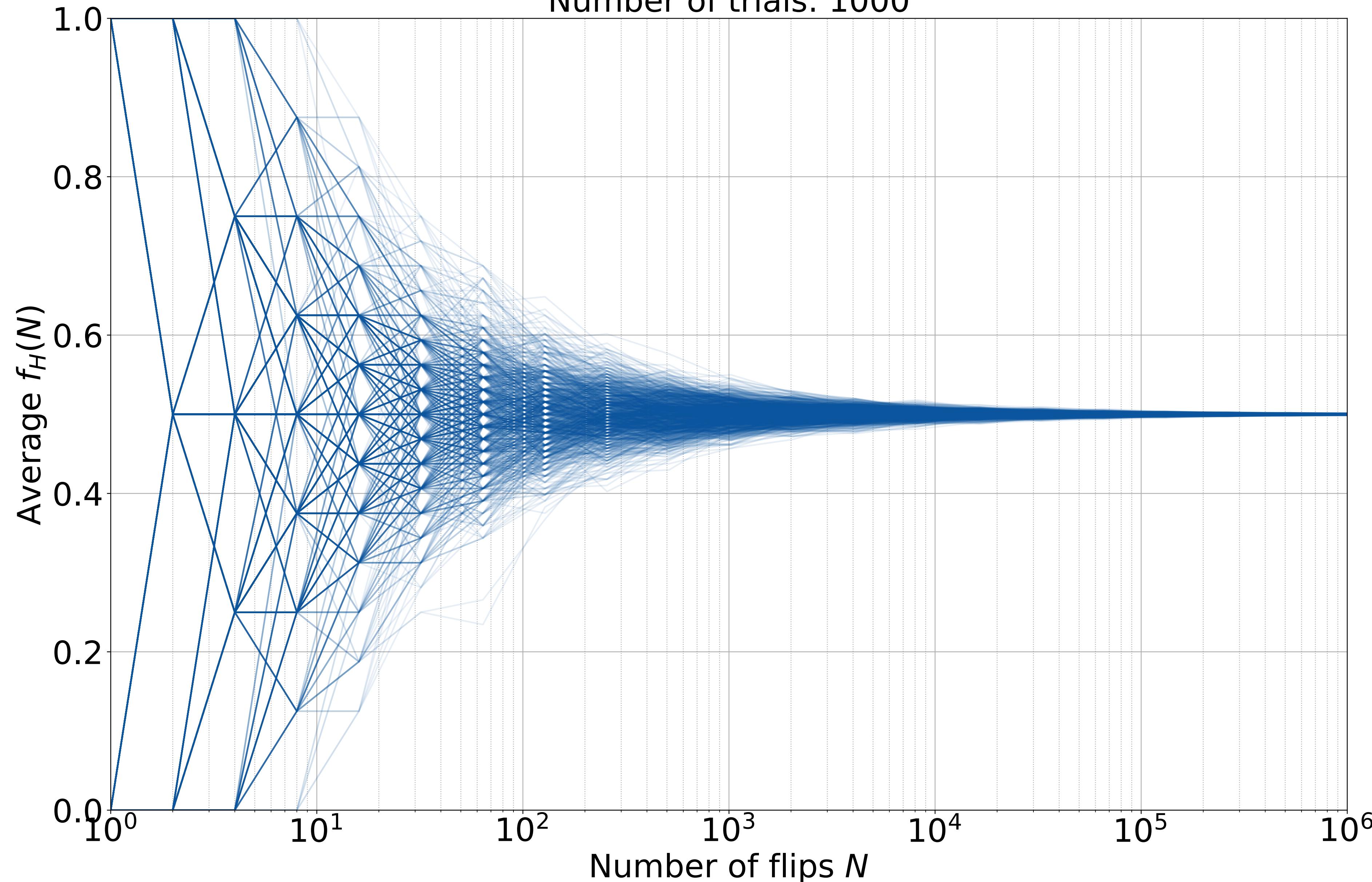
# Frequentist Probability

Number of trials: 100



# Frequentist Probability

Number of trials: 1000



# Frequentist Probability

Probabilities can represent *frequencies*.

e.g. Flip a coin  $N$  times, define  $N(H)$  to be the number of times it lands heads. The *relative frequency*  $f_H(N)$  of landing heads is



$$f_H(N) = \frac{N(H)}{N}$$

The probability of a head, written  $p(H)$  is

$$p(H) := \lim_{N \rightarrow \infty} f_H(N)$$

The symbol  $\lim_{N \rightarrow \infty}$  is called a *limit*. It is the formal way of saying “when  $N$  gets big”.

Frequentists: *event probability* = long run frequency in a repeatable experiment.

# Bayesian Probability

Probabilities can represent *beliefs*

e.g. Given the results of a blood test, the probability that Sri Dewi has a nasty disease is  $p\%$ .

e.g. The probability that Gunung Merapi will erupt on 31st December 2020 is  $q\%$ .

Such claims cannot be verified through repeated experimentation. This subjective interpretation or *Bayesian* interpretation expresses *degree of belief*.

## Frequentist and Bayesian probability treated with same theory

Revd. Thomas Bayes



N.B. other interpretations exist: propensity, logical probability, mechanistic, etc.

\*I've heard a rumor that the gentleman pictured above may not actually be the Reverend Thomas Bayes.

# The Probability Axioms

1) Probability of event  $x$  is a non-negative real number

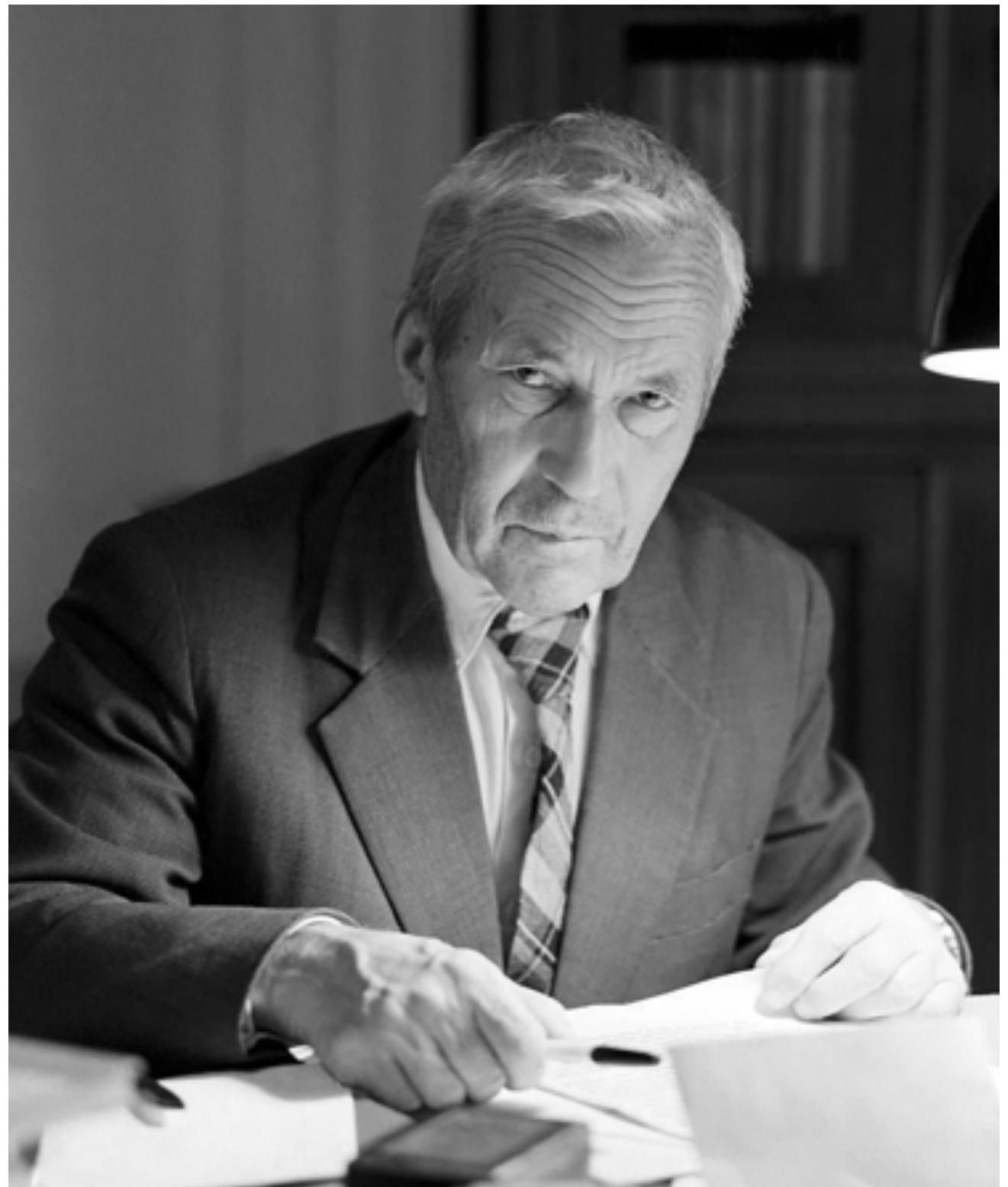
$$p(x) \geq 0 \quad \text{for all } x \subseteq \Omega$$

2) Certain events have unit probability

$$p(\Omega) = 1$$

3) Countable additivity: for disjoint events  $x_1, x_2, \dots, x_N$

$$p(x_1 \cup x_2 \cup \dots \cup x_N) = p(x_1) + p(x_2) + \dots + p(x_N)$$



Other rules:

- **Complement rule:**

$$p(\Omega \setminus x) = 1 - p(x)$$

Andrey Kolmogorov

- **Impossible events:**

$$p(\emptyset) = 0$$

- **Subsets:**

$$x_1 \subseteq x_2 \implies p(x_1) \leq p(x_2)$$

- **Union rule:**

$$p(x_1 \cup x_2) = p(x_1) + p(x_2) - p(x_1 \cap x_2)$$

## Sample space

A *sample space*  $\Omega$  is the *set* of possible outcomes of an experiment.  
 Outcomes are called *samples*.

## Events

An *event*  $E$  is a subset of a sample space  $E \subseteq \Omega$

## Event space

An *event space*<sup>1</sup>  $\Sigma$  is the space of all events  $E \subseteq \Omega$

## Probability Mass Function (PMF)

A *probability mass function*<sup>2</sup>  $p$  assigns a number in  $[0,1]$  to every event in the event space.

- $p(A) = 1$  means that  $A \in \Sigma$  is certain
- $p(A) = 0$  means that  $A \in \Sigma$  will never happen
- If  $p(A) > p(B)$ , then A is more likely than B

<sup>1</sup> In the continuous setting the definition is a bit fiddly

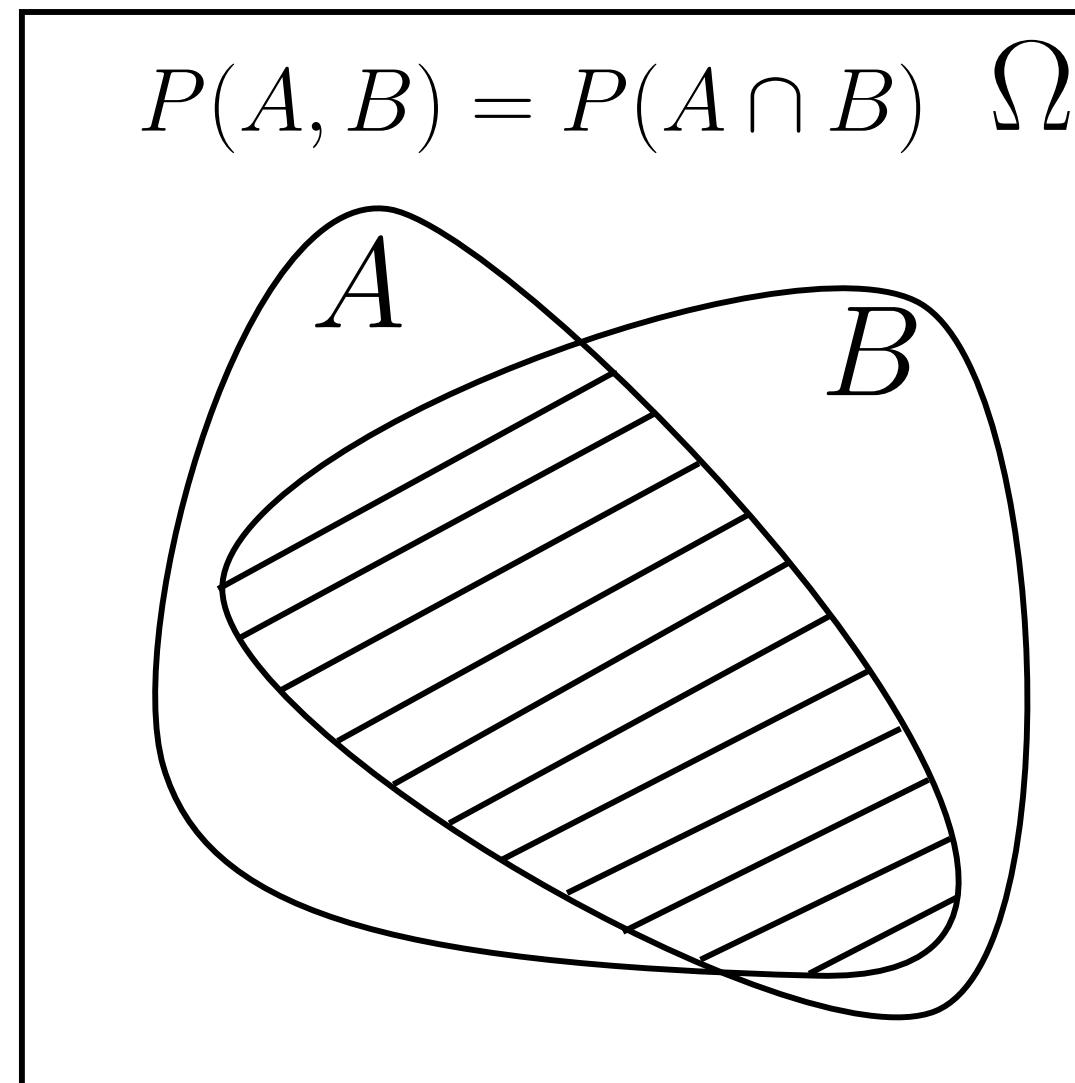
<sup>2</sup> Note in statistics a capital P is used, but in machine learning we are lazy and just use a small p

$i$	$a_i$	$p_i$	
1	a	0.0575	a
2	b	0.0128	b
3	c	0.0263	c
4	d	0.0285	d
5	e	0.0913	e
6	f	0.0173	f
7	g	0.0133	g
8	h	0.0313	h
9	i	0.0599	i
10	j	0.0006	j
11	k	0.0084	k
12	l	0.0335	l
13	m	0.0235	m
14	n	0.0596	n
15	o	0.0689	o
16	p	0.0192	p
17	q	0.0008	q
18	r	0.0508	r
19	s	0.0567	s
20	t	0.0706	t
21	u	0.0334	u
22	v	0.0069	v
23	w	0.0119	w
24	x	0.0073	x
25	y	0.0164	y
26	z	0.0007	z
27	-	0.1928	-

PMF over letters in English

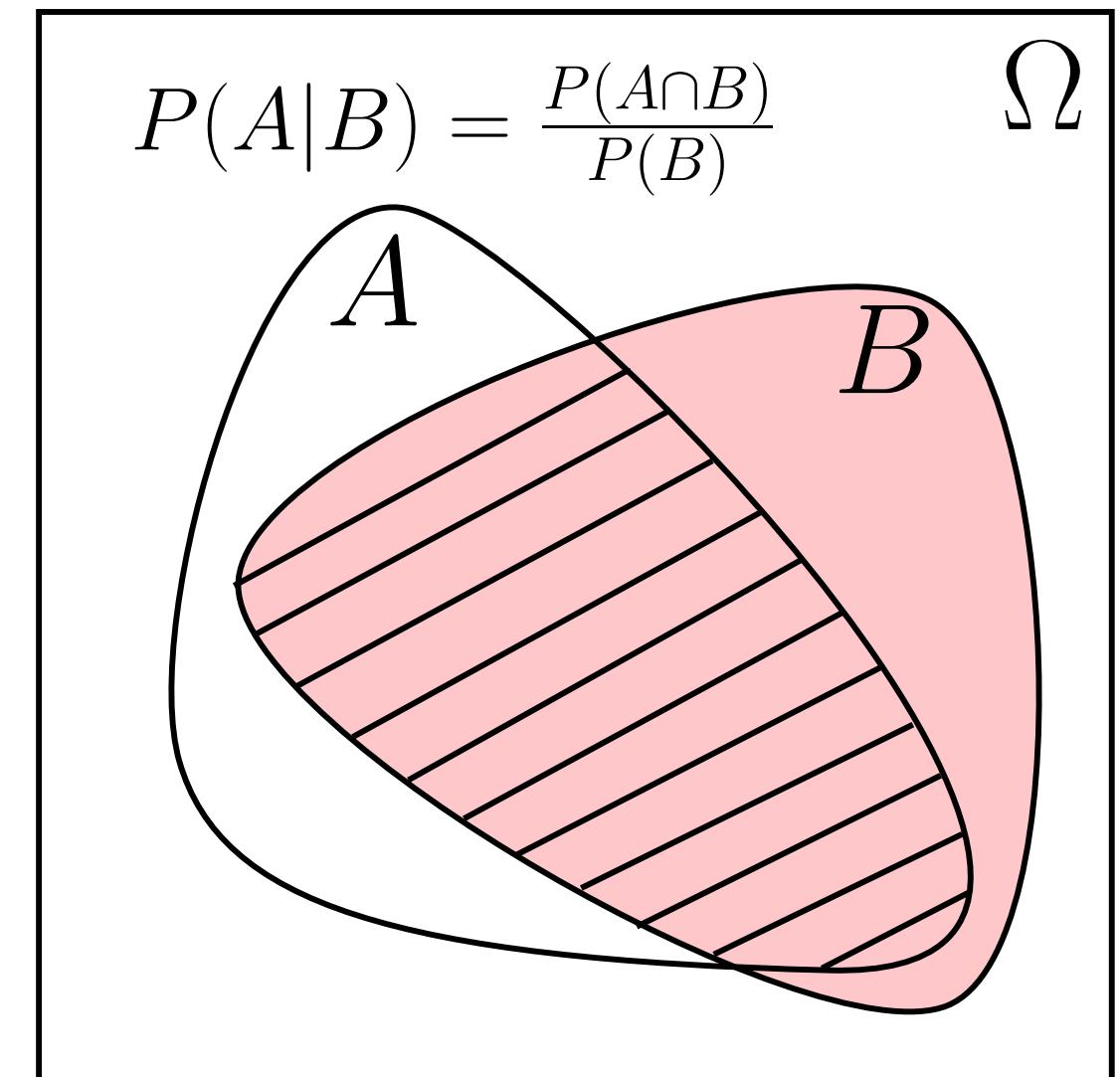
# Conditional and Joint Probability

**Joint probability:**  $A$  and  $B$  co-occur



$$p(A, B) = p(B, A)$$

**Conditional probability:**  $B$  then  $A$



$$p(A|B) \neq p(B|A)$$

**Product rule**  $p(A, B) = p(A|B)p(B)$

**Sum rule**  $p_A(A) = \sum_B p_{A,B}(A, B)$  or  $p_A(A) = \int_B p_{A,B}(A, B) dB$

**Sum rule proof**  $\sum_B p(A, B) = \sum_B p(B|A)p(A) = p(A) \underbrace{\sum_B p(B|A)}_{=1} = p(A)$

Note sometimes we write  $p(x)$  and other times we will write  $p_X(x)$  depending on context

# Some definitions (Google them later if you don't know them)

## Bayes' Rule

$$p(x|y) = \frac{p(y|x)p(x)}{p(y)}$$

Given  $y$  has happened, what is prob. of  $x$ ?

### Proof

$$p(y|x)p(x) = p(x, y) = p(x|y)p(y)$$

## Change of variables formula

$$p_X(x) = p_Z(z) \left| \frac{\partial z}{\partial x} \right| \quad x = f(z)$$

### Proof

$$\begin{aligned} p_X(x) &= \int_Z p(x|z)p_Z(z) dz \\ &\stackrel{\text{e.g.}}{=} \int_Z \delta(x - f(z))p_Z(z) dz \\ &= \int_U \delta(x - u)p_Z(f^{-1}(u)) \left| \frac{\partial z}{\partial u} \right| du \\ &= p_Z(z) \left| \frac{\partial z}{\partial x} \right| \end{aligned}$$

Jacobian tracks volume change

## Expectations

$$\mathbb{E}_{p(\mathbf{x})} [f(\mathbf{x})] = \int f(\mathbf{x})p(\mathbf{x}) d\mathbf{x} \quad \text{or} \quad \sum_{\mathbf{x}} f(\mathbf{x})p(\mathbf{x})$$

### Mean

$$\bar{\mathbf{x}} = \mathbb{E}_{p(\mathbf{x})} [\mathbf{x}]$$

### Covariance

$$\Sigma = \mathbb{E}_{p(\mathbf{x})} [(\mathbf{x} - \bar{\mathbf{x}})(\mathbf{x} - \bar{\mathbf{x}})^{\top}]$$

## Information Theory

### Surprisal

$$I(\mathbf{x}) = -\log p(\mathbf{x})$$

### Entropy = average surprise

$$H(\mathbf{x}) = \mathbb{E}_{p(\mathbf{x})} [-\log p(\mathbf{x})]$$

### Kullback-Leibler divergence

$$D_{KL}(p(\mathbf{x})||q(\mathbf{x})) = \int p(\mathbf{x}) \log \frac{p(\mathbf{x})}{q(\mathbf{x})} d\mathbf{x}$$

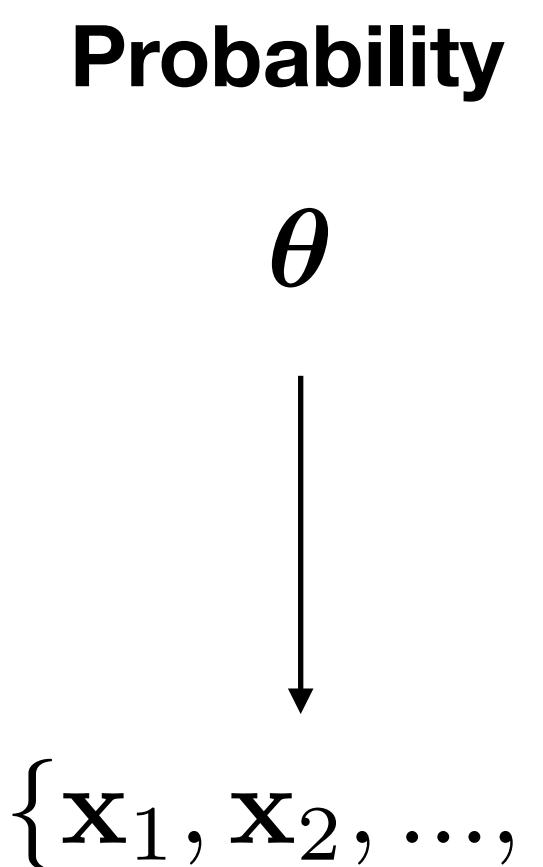
### Mutual information

$$\mathbb{I}[\mathbf{x}; \mathbf{y}] = D_{KL}(p(\mathbf{x}, \mathbf{y})||p(\mathbf{x})p(\mathbf{y}))$$

We are mostly concerned with models which look like

$$p(\mathbf{x} \mid \theta)$$

In many case  $\mathbf{x}$  refers to an *observation* and refers to a set of *parameters*.



# Forward models

You build a model of radioactive decay.

**Model 1:** Particles decay  $x$  cm from the source, following an exponential distribution:

$$p(x|\lambda) = \frac{1}{Z} e^{-\lambda x}$$

Variable      Parameter      Normalization constant

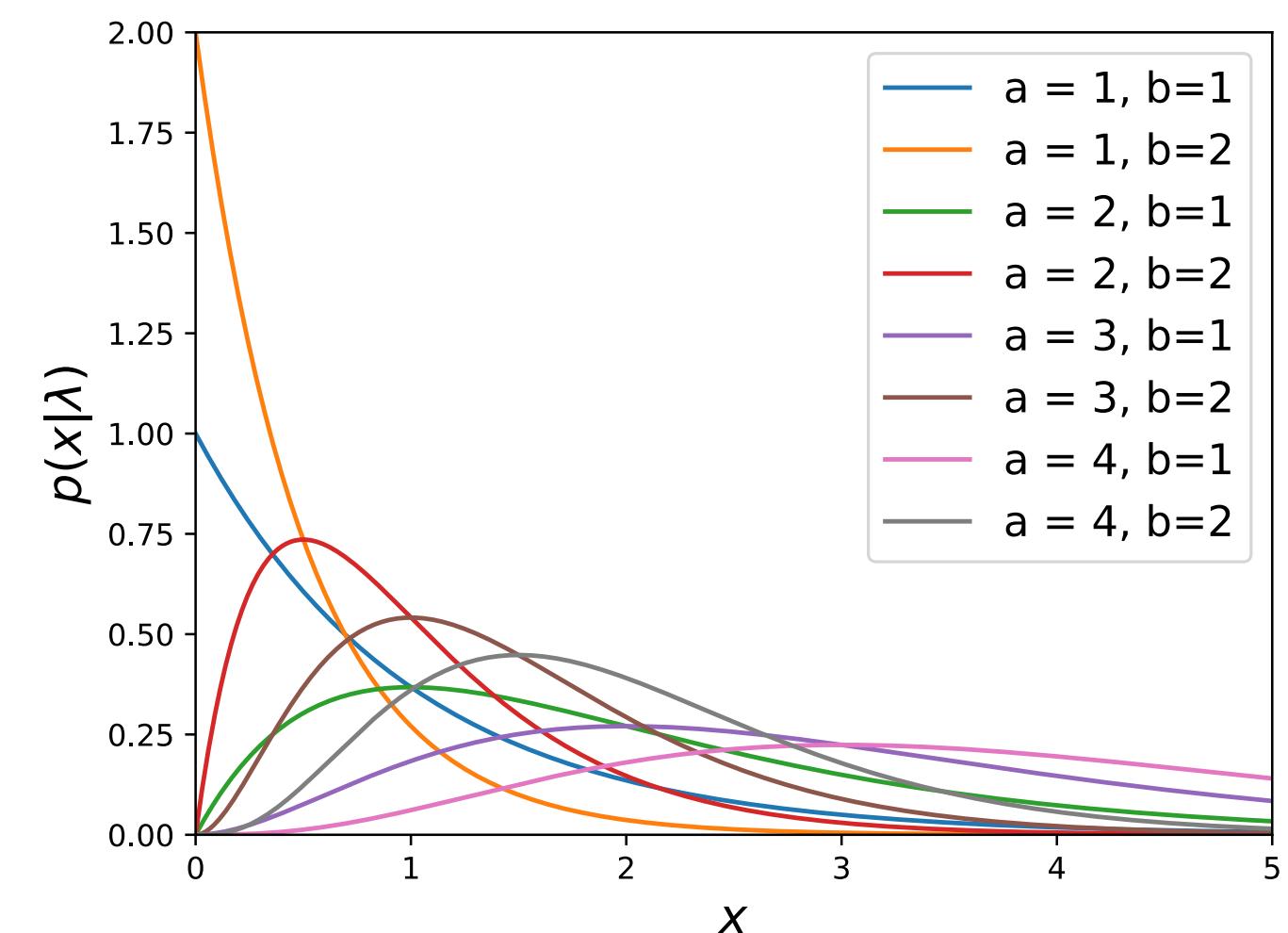
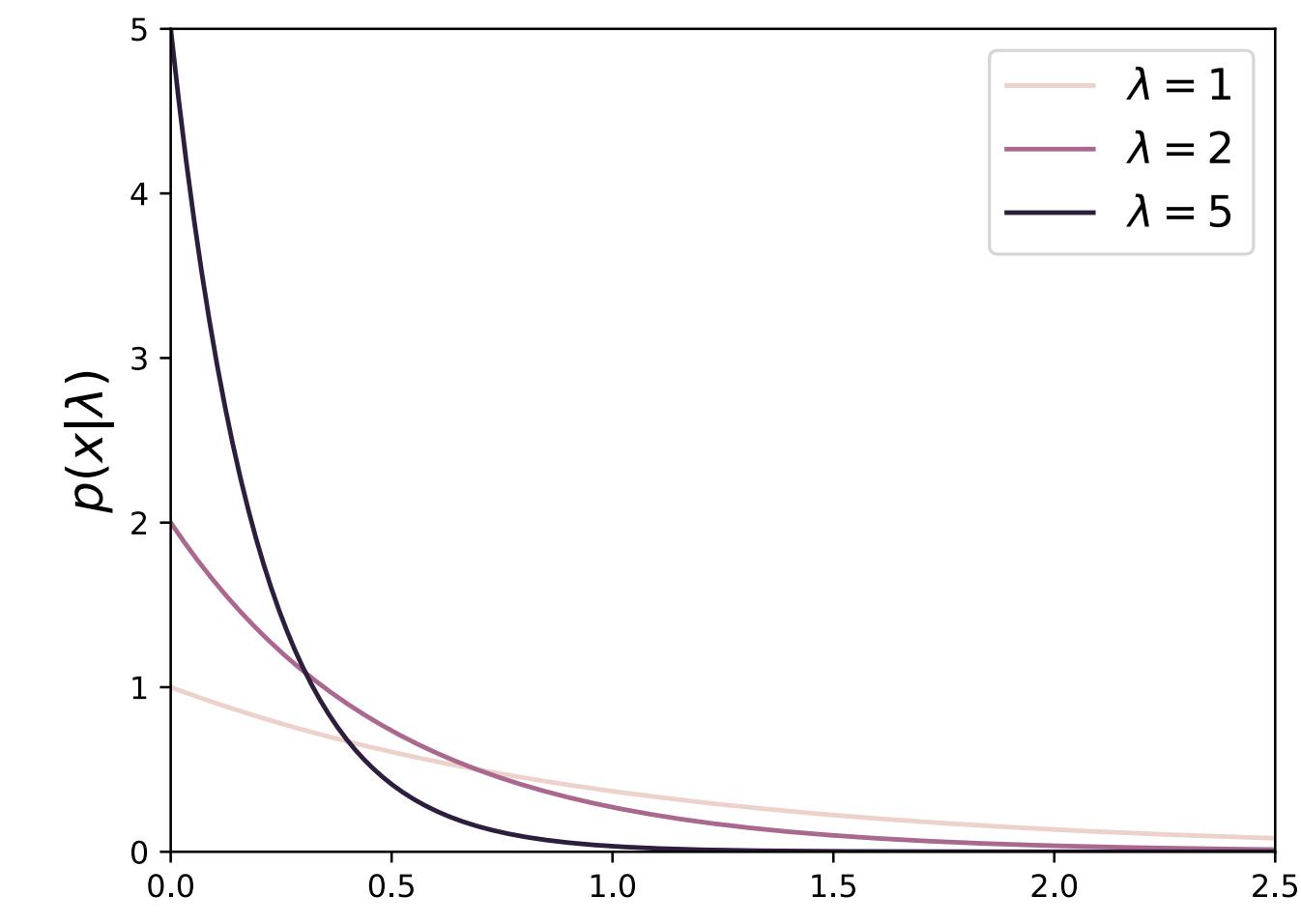
$$\int_0^\infty \frac{1}{Z} e^{-\lambda x} dx = 1 \implies Z = \frac{1}{\lambda}$$

**Model 2:** A gamma distribution

$$p(x|\alpha, \beta) = \frac{1}{Z} x^{\alpha-1} e^{-\beta x}, \quad Z = \frac{\Gamma(\alpha)}{\beta^\alpha}$$

More than 1 parameter

**How to model multiple observations?**



Joint probability of *sequence*  $\{x_1, x_2, \dots, x_N\}$ :

$$p(x_1, x_2, \dots, x_N)$$

## Marginal independence

Probability of each observation *independent* (doesn't depend on other observations) and *identical* (from the same distribution).

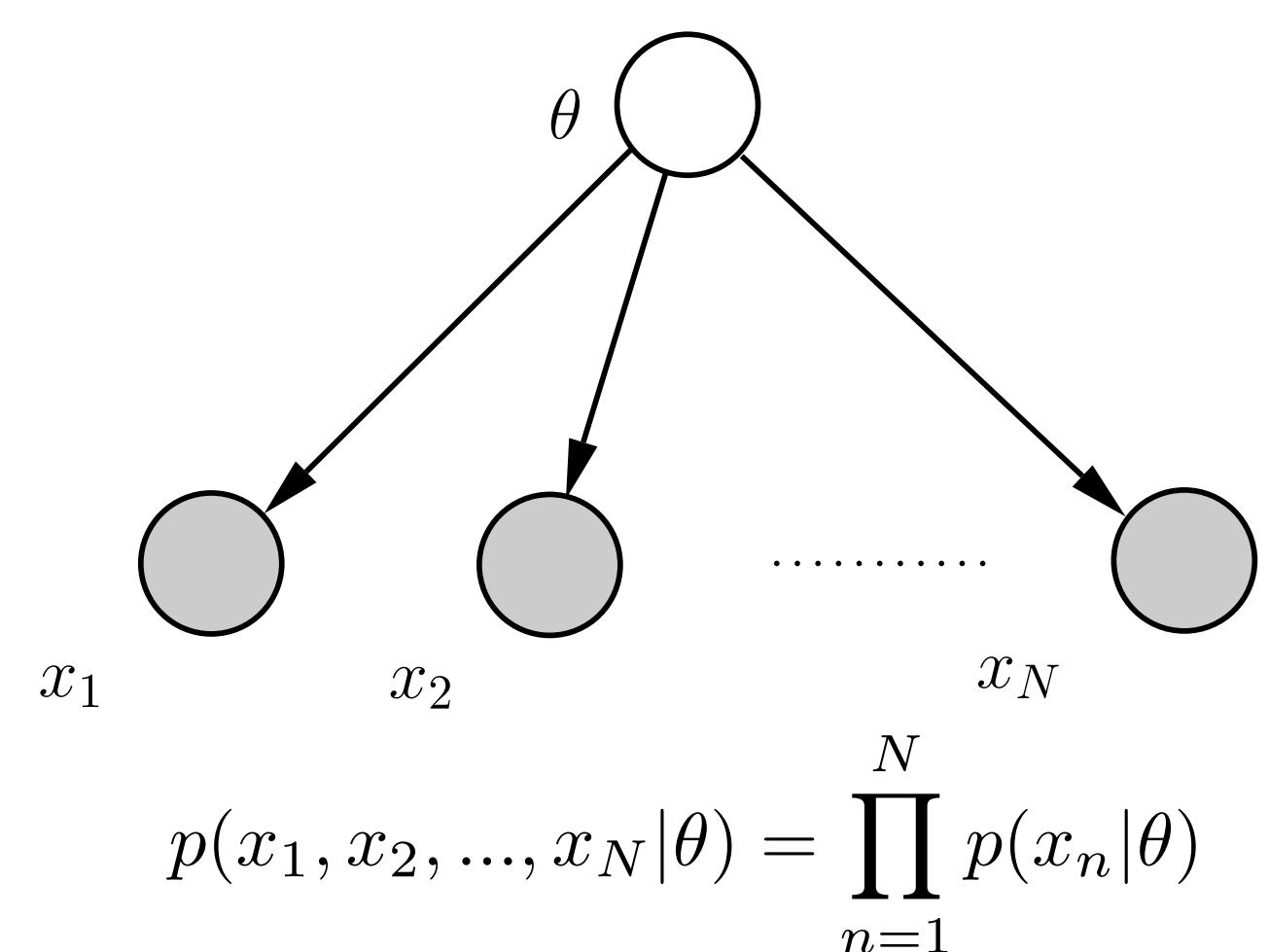
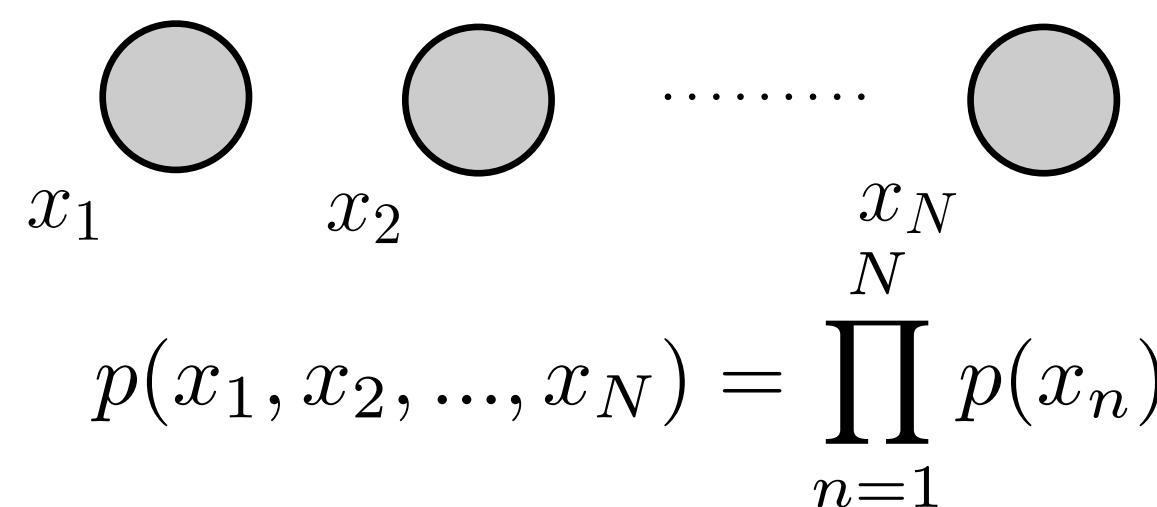
$$p(x_1|x_i) = p(x_1) \quad \text{for all } x_i \neq x_1$$

## Conditional independence

Independence *conditional* on extra information

$$p(x_1, x_2|x_i) = p(x_1|x_i)p(x_2|x_i)$$

for  $x_i \neq x_1, x_i \neq x_2$



# Forward models

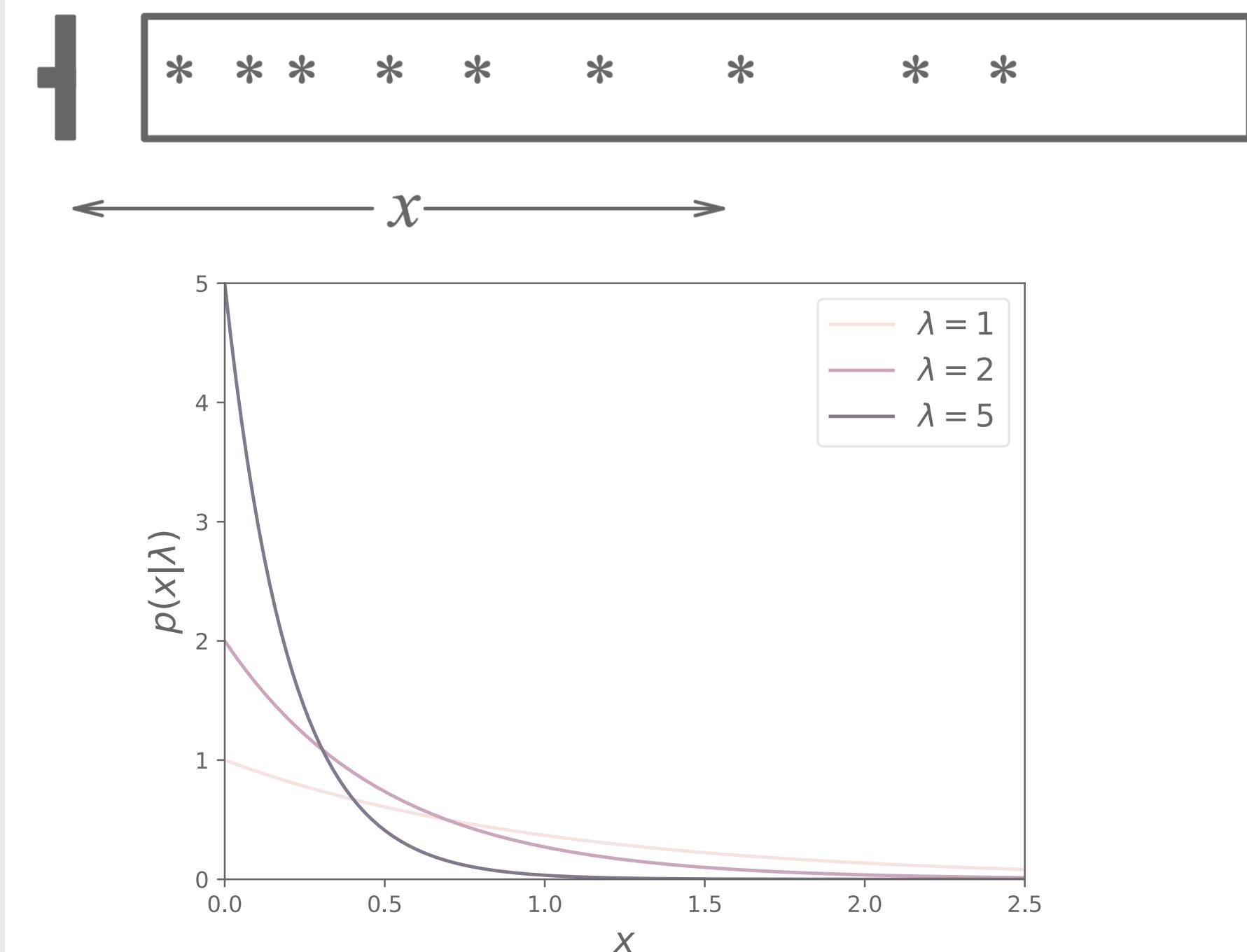
You build a model of radioactive decay.

**Model 1:** Particles decay  $x$  cm from the source, following an exponential distribution:

$$p(x|\lambda) = \frac{1}{Z} e^{-\lambda x}$$

Variable      Parameter      Normalization constant

$$\int_0^\infty \frac{1}{Z} e^{-\lambda x} dx = 1 \implies Z = \frac{1}{\lambda}$$



**How to model multiple observations?**

$$p(x_1, x_2, \dots, x_N | \lambda) = \prod_{n=1}^N p(x_n | \lambda) = \prod_{n=1}^N \frac{1}{Z} e^{-\lambda x_n} = \lambda^N e^{-\lambda \sum_{n=1}^N x_n}$$

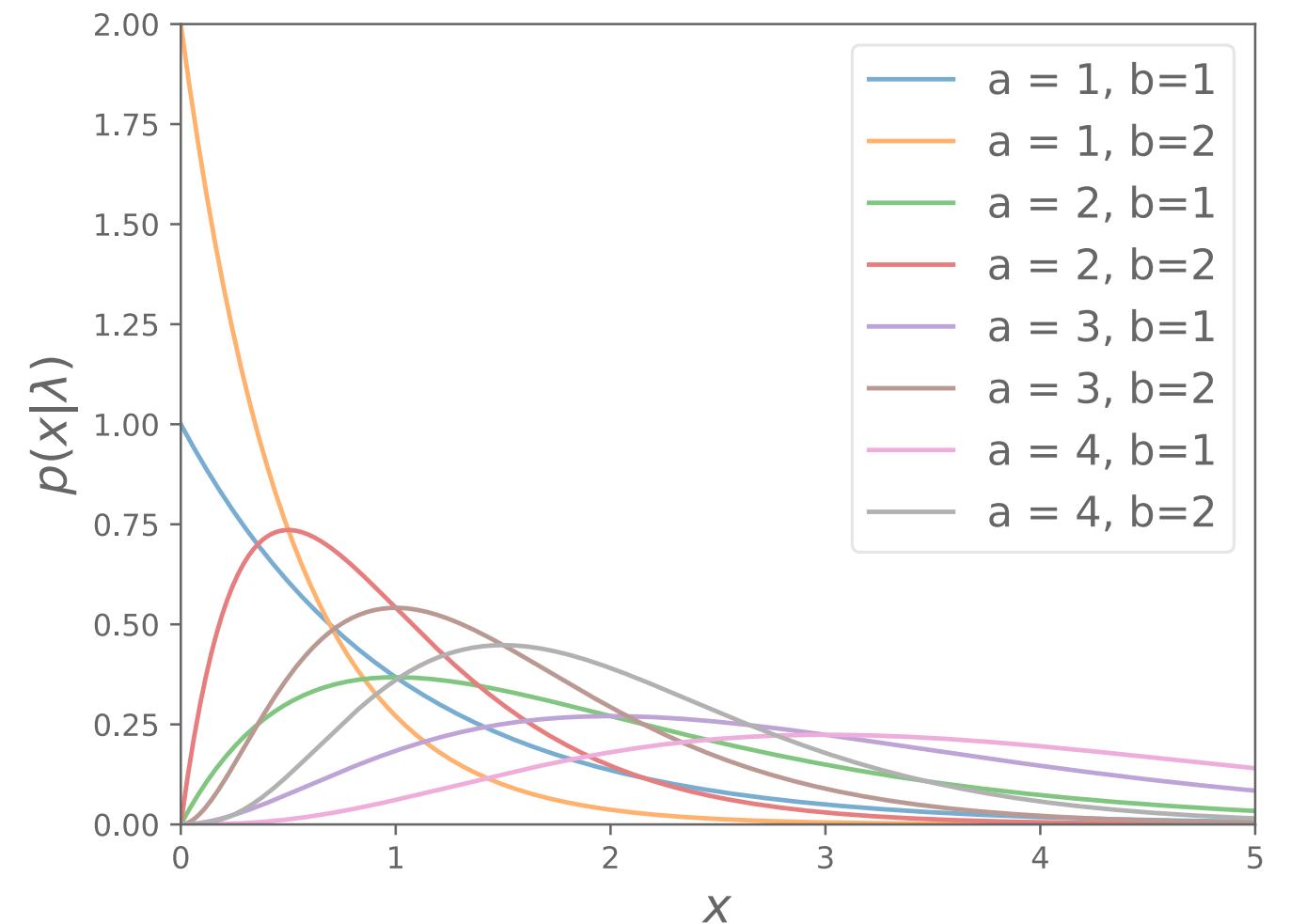
# Forward models

**Model 2:** A gamma distribution

$$p(x|\alpha, \beta) = \frac{1}{Z} x^{\alpha-1} e^{-\beta x}, \quad Z = \frac{\Gamma(\alpha)}{\beta^\alpha}$$

↑  
↑  
More than 1 parameter

**How to model multiple observations?**



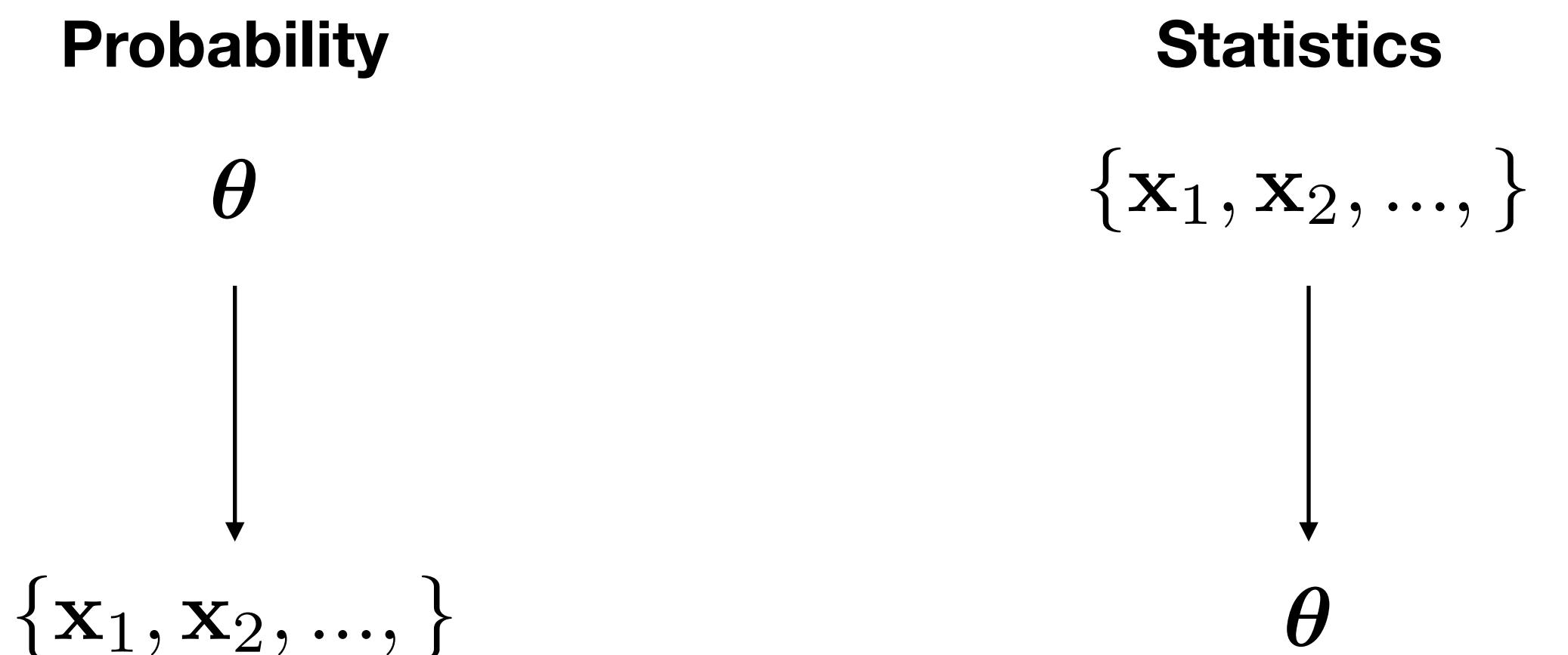
**How to model multiple observations?**

$$p(x_1, x_2, \dots, x_N | \alpha, \beta) = \prod_{n=1}^N p(x_n | \alpha, \beta) = \left( \frac{\beta^\alpha}{\Gamma(\alpha)} \right)^N \prod_{n=1}^N x_n^{\alpha-1} e^{-\beta x_n}$$

We are mostly concerned with models which look like

$$p(\mathbf{x} \mid \theta)$$

In many case  $\mathbf{x}$  refers to an *observation* and refers to a set of *parameters*.



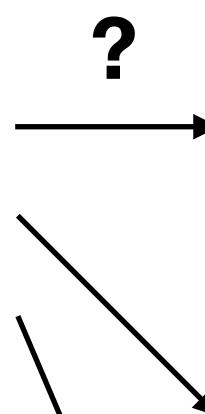
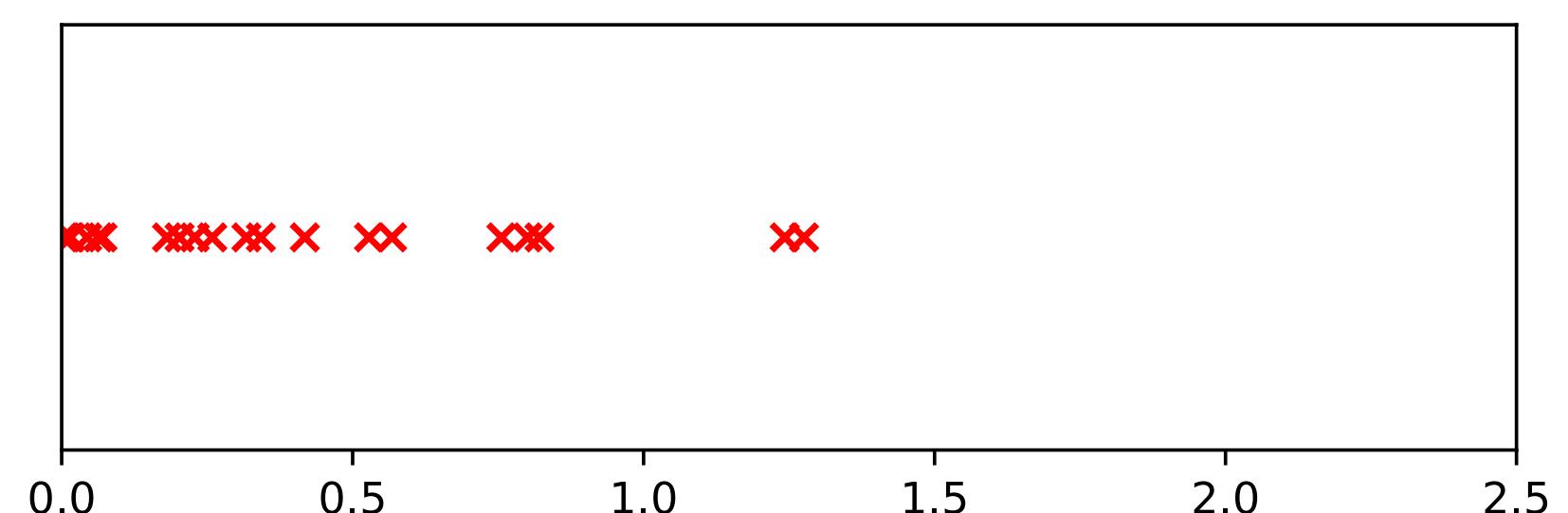
# Maximum likelihood

$$\mathcal{D} = \{x_1, x_2, \dots, x_N\}$$

Forward model: generate **data given parameters**

Inference: find **parameters given data**

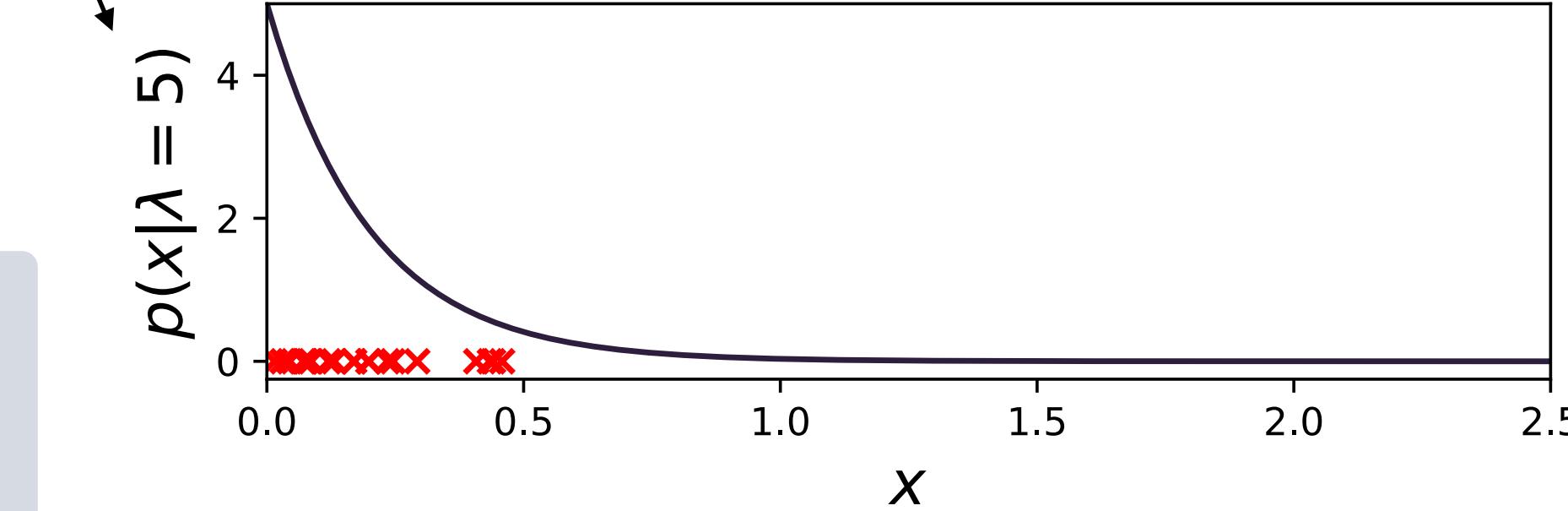
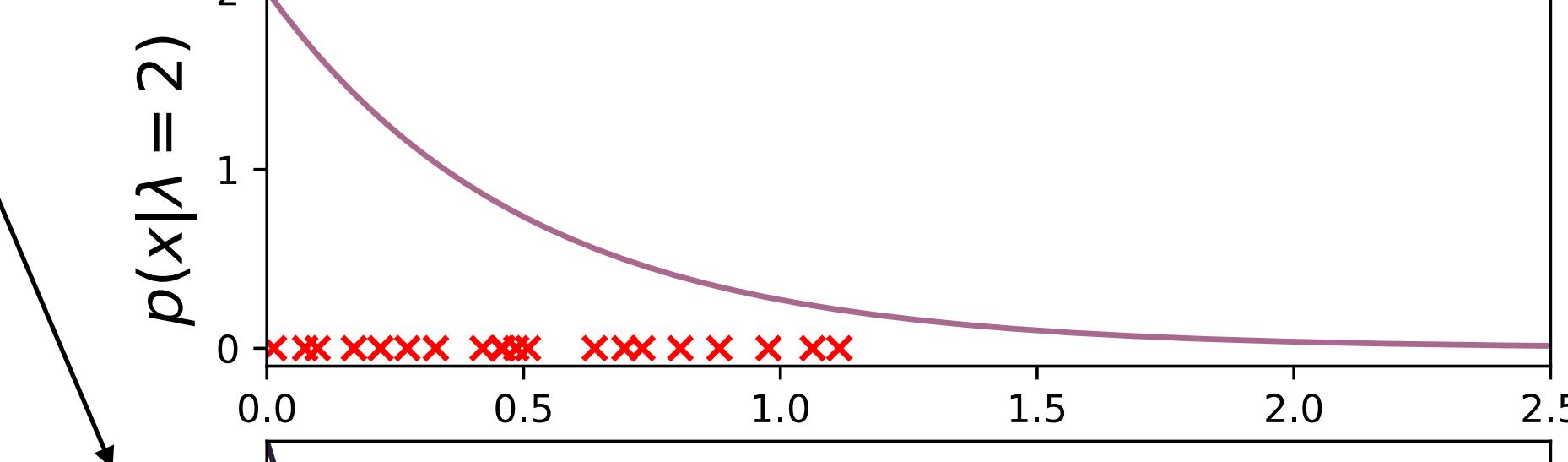
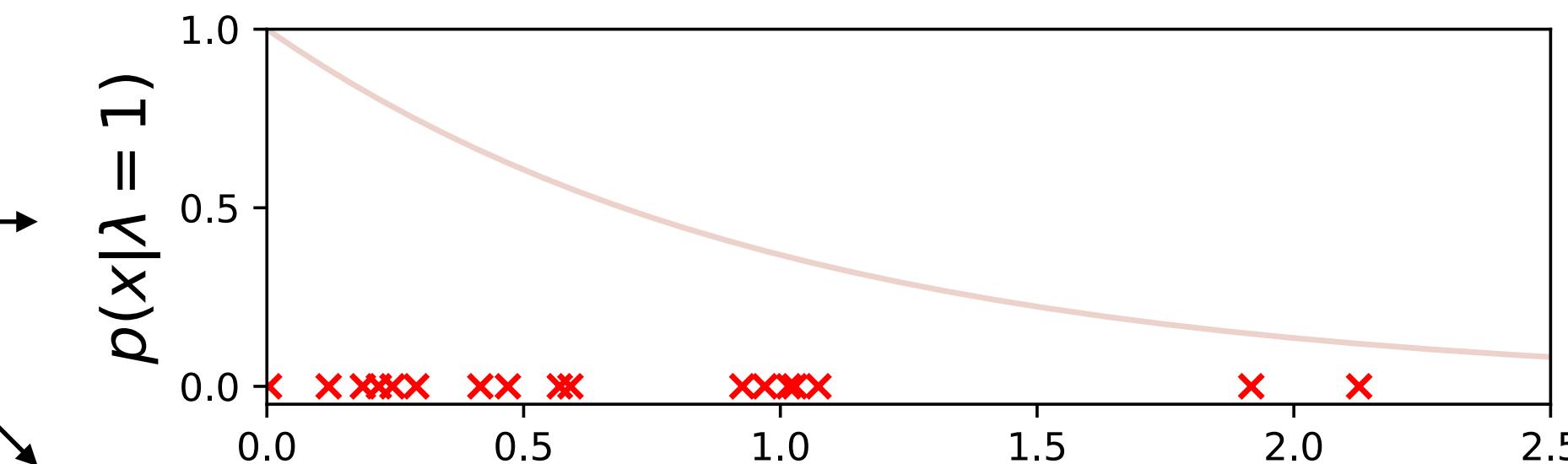
**Idea:** enumerate every possible forward model and see if it matches the data



$$\lambda^* = \arg \max_{\lambda} p(\mathcal{D}|\lambda)$$

**Likelihood:** function of the parameter  $\lambda$

We say the likelihood of  $\lambda$  given  $\mathcal{D}$ , never the likelihood of  $\mathcal{D}$  given  $\lambda$ !



# Maximum log-likelihood

We find the maximum by setting<sup>1</sup> the derivative to 0

$$\lambda^* = \arg \max_{\lambda} p(\mathcal{D}|\lambda) \implies \frac{\partial}{\partial \lambda} p(\mathcal{D}|\lambda^*) = 0$$

We prefer to find the maximum of the log-likelihood

$$\lambda^* = \arg \max_{\lambda} p(\mathcal{D}|\lambda) \iff \lambda^* = \arg \max_{\lambda} \log p(\mathcal{D}|\lambda)$$

$$\begin{aligned} \lambda^* &= \arg \max_{\lambda} \log p(\mathcal{D}|\lambda) \\ &= \arg \max_{\lambda} \log \prod_{n=1}^N p(x_n|\lambda) \\ &= \arg \max_{\lambda} \sum_{n=1}^N \log p(x_n|\lambda) \end{aligned}$$

Why?

- 1) Small likelihoods  $\rightarrow$  numerical underflow
- 2) Derivatives of sums easier than derivatives of products

$$\frac{\partial}{\partial \lambda} [f_1(\lambda)f_2(\lambda)f_3(\lambda)] = \frac{\partial f_1}{\partial \lambda}f_2(\lambda)f_3(\lambda) + f_1(\lambda)\frac{\partial f_2}{\partial \lambda}f_3(\lambda) + f_1(\lambda)f_2(\lambda)\frac{\partial f_3}{\partial \lambda}$$

**Which do you prefer?**

$$\frac{\partial}{\partial \lambda} [\log(f_1(\lambda)f_2(\lambda)f_3(\lambda))] = \frac{\partial \log f_1}{\partial \lambda} + \frac{\partial \log f_2}{\partial \lambda} + \frac{\partial \log f_3}{\partial \lambda}$$

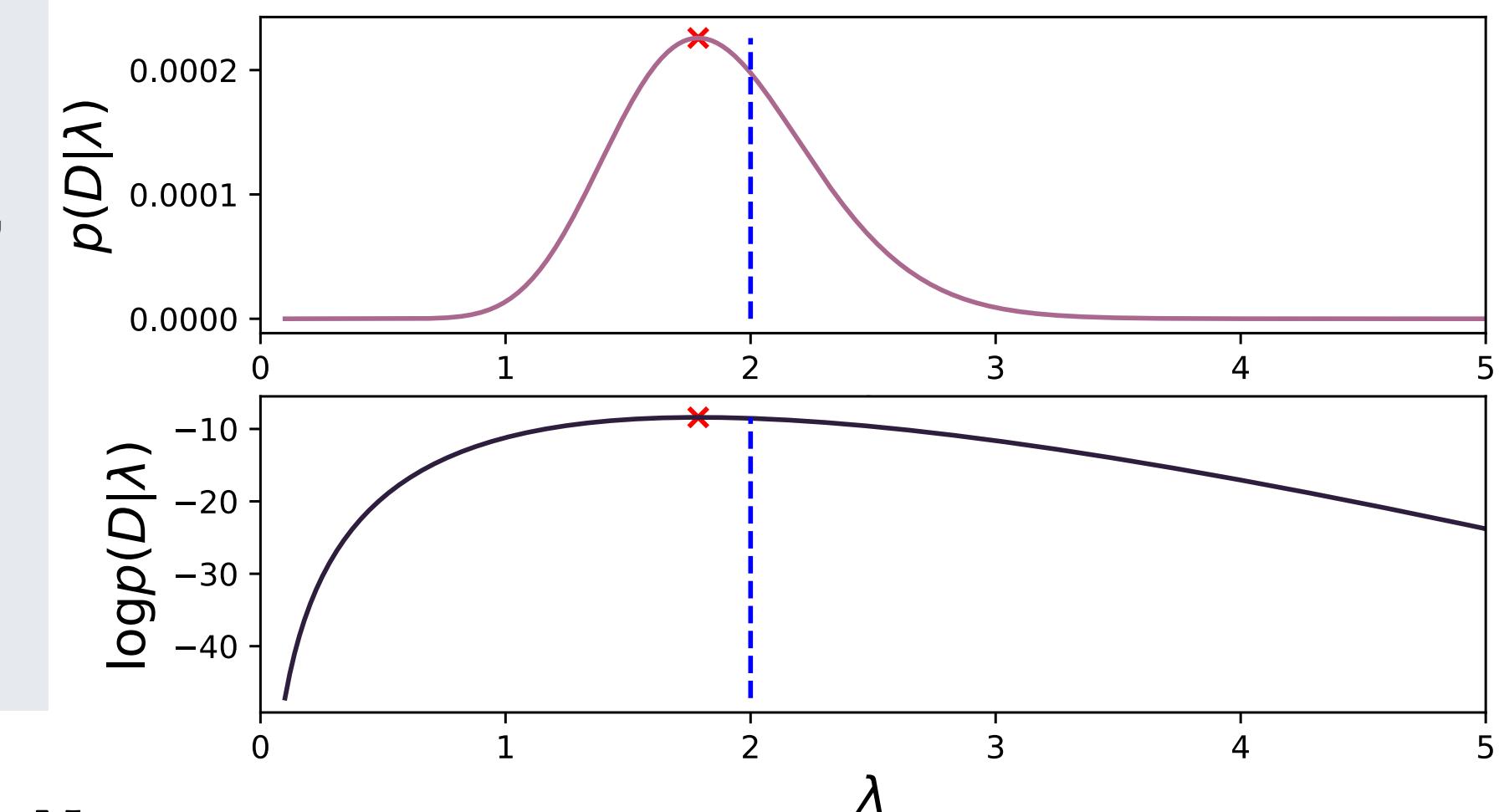
<sup>1</sup> Note we should also check that the Hessian is positive-definite, but for typical distributions this is not necessary, since members of the exponential family are log-concave in the parameters.

# Maximum Likelihood Example

You build a model of radioactive decay.

**Model 1:** Particles decay  $x$  cm from the source, following an exponential distribution:

$$p(x|\lambda) = \frac{1}{Z} e^{-\lambda x}$$



$$\begin{aligned}\lambda^* &= \arg \max_{\lambda} \log p(\mathcal{D}|\lambda) = \arg \max_{\lambda} \sum_{n=1}^N \log p(x_n|\lambda) \\ &= \arg \max_{\lambda} \sum_{n=1}^N \log \lambda e^{-\lambda x_n} = \arg \max_{\lambda} N \log \lambda - \sum_{n=1}^N \lambda x_n\end{aligned}$$

$$\frac{\partial}{\partial \lambda} \left( N \log \lambda - \sum_{n=1}^N \lambda x_n \right) = \frac{N}{\lambda} - \sum_{n=1}^N x_n = 0 \quad \Rightarrow \quad \lambda^* = \frac{1}{\frac{1}{N} \sum_{n=1}^N x_n}$$

# The Exponential Family

Which distributions have analytical maximum likelihood solutions? Most of the distributions we have looked at are fairly similar. They have three main components:

$$p(x|\theta) = \underbrace{\frac{1}{Z(\theta)}}_{\text{normalizer}} \cdot \overbrace{b(x)}^{\text{fnc of } x} \cdot \underbrace{\exp\{\theta^\top \mathbf{t}(x)\}}_{\exp \text{ of linear fnc of } \theta}$$

*Natural parameters*<sup>1</sup>  $\theta$  and *sufficient statistics*  $\mathbf{t}(x)$ .

This may seem like an odd choice, but it has some very handy properties, which allow for **lightning fast** computation. At the ML solution

Model expectation

$$\mathbb{E}_{p(x|\theta)} [\mathbf{t}(x)] = \frac{1}{N} \sum_{n=1}^N \mathbf{t}(x_n)$$

Data expectation

*Method of moments*: Use a mean value mapping to recover the ML parameters, tractably

$$\tau(\theta) = \mathbb{E}_{p(x|\theta)} [\mathbf{t}(x)] \implies \theta^* = \tau^{-1} \left( \frac{1}{N} \sum_{n=1}^N \mathbf{t}(x_n) \right)$$

<sup>1</sup> There are many different names and notations for this, so beware!

# Maximum Likelihood Example

You build a model of radioactive decay.

**Model 1:** Particles decay  $x$  cm from the source, following an exponential distribution:

$$p(x|\lambda) = \frac{1}{Z} e^{-\lambda x}$$

$$p(x|\theta) = \underbrace{\frac{1}{Z(\theta)}}_{\text{normalizer}} \cdot \overbrace{b(x)}^{\text{fnc of } x} \cdot \underbrace{\exp\{\theta^\top t(x)\}}_{\text{exp of linear fnc of } \theta}$$

$$\tau(\theta) = \mathbb{E}_{p(x|\theta)} [t(x)] \implies \theta^* = \tau^{-1} \left( \frac{1}{N} \sum_{n=1}^N t(x_n) \right)$$

Let's apply the method of moments to the exponential distribution example

## 1) Identify parameterisation

$$\theta = -\lambda, \quad t(x) = x, \quad b(x) = 1, \quad Z(\theta) = \frac{1}{\lambda}$$

## 2) Find mean value mapping

$$\tau(\lambda) = \mathbb{E}_{p(x|\lambda)}[x] = \int_0^\infty x \lambda e^{-\lambda x} dx = \frac{1}{\lambda}$$

## 3) Plug into method of moments formula

$$\frac{1}{\lambda^*} = \frac{1}{N} \sum_{n=1}^N x_n$$

We just learnt about maximum likelihood, where we solved

$$\theta_{\text{ML}} = \arg \max_{\theta} p(\mathcal{D}|\theta)$$

- Is this the best we can do?
- **It is almost certainly wrong<sup>1</sup>:**  $p(\theta_{\text{ML}} = \theta_{\text{true}}) = 0$
- Our model is almost certainly wrong as well

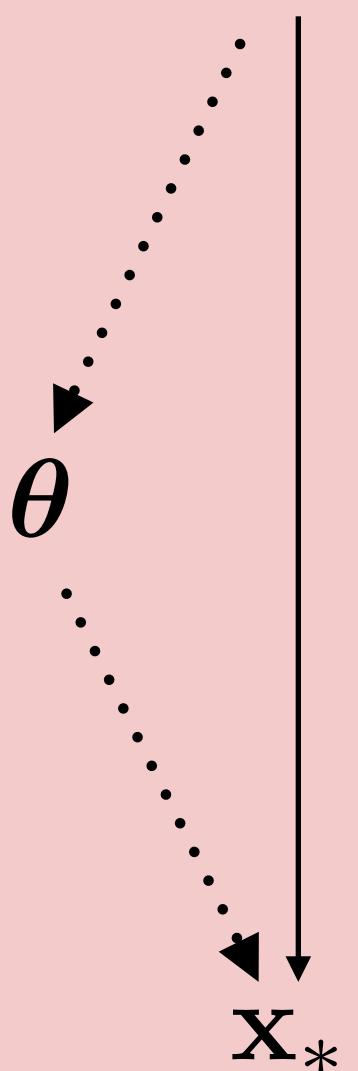
Why do we want  $\theta_{\text{ML}}$  anyway?

## Options

1. We are actually interested in knowing  $\theta$
2. We don't care: want to generate new samples  $x_*$

## Prediction

$$\{\mathbf{x}_1, \mathbf{x}_2, \dots, \}$$



<sup>1</sup> For continuous parameterisations

# Bayesian Inference

Let's think about the distribution  $p(x_*|\mathcal{D})$

We can compute it from *known* quantities:

This approach is called  
*generative modeling*

$$\begin{aligned} p(x_*|\mathcal{D}) &= \int p(x_*, \theta|\mathcal{D}) d\theta \\ &= \int p(x_*|\theta)p(\theta|\mathcal{D}) d\theta \end{aligned}$$

Weighted average

Forward model

Likelihood      Prior

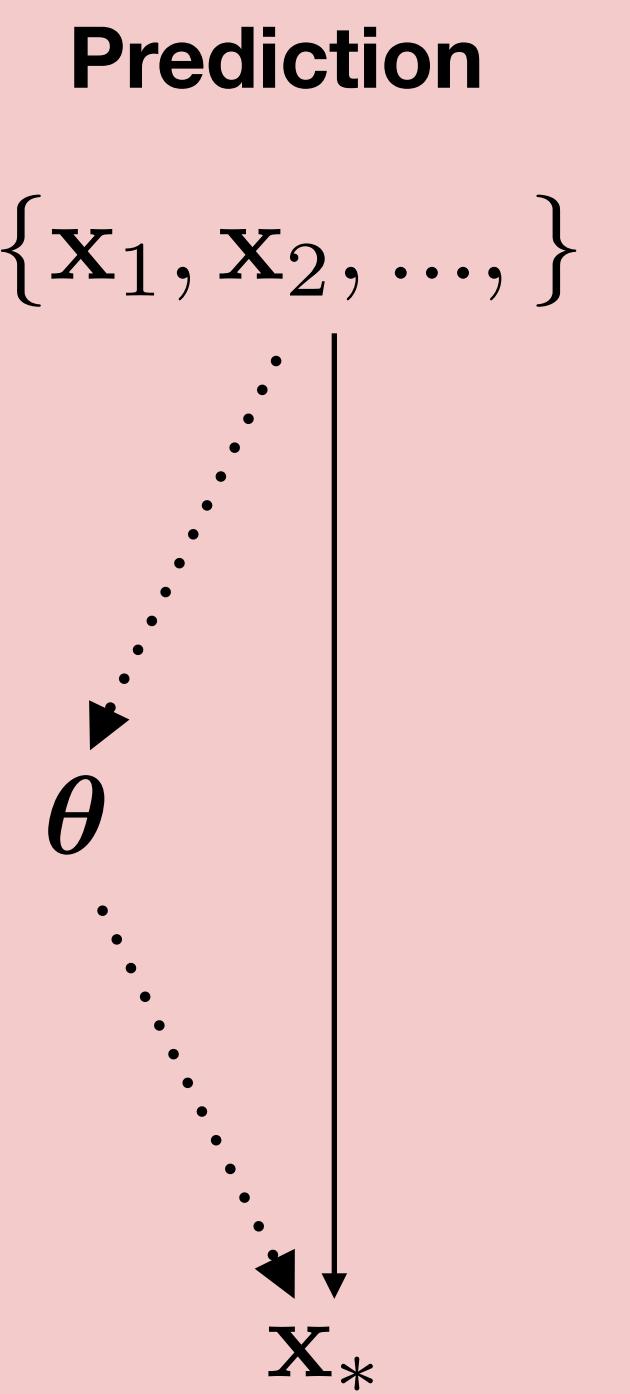
**Posterior distribution**

$$p(\theta|\mathcal{D}) = \frac{p(\mathcal{D}|\theta)p(\theta)}{p(\mathcal{D})} = \frac{p(\mathcal{D}|\theta)p(\theta)}{\int p(\mathcal{D}|\theta)p(\theta) d\theta}$$

Evidence/marginal likelihood

Read, “*the probability of the parameters, given the data*”.

More descriptive than point estimate  $\theta_{ML}$ .

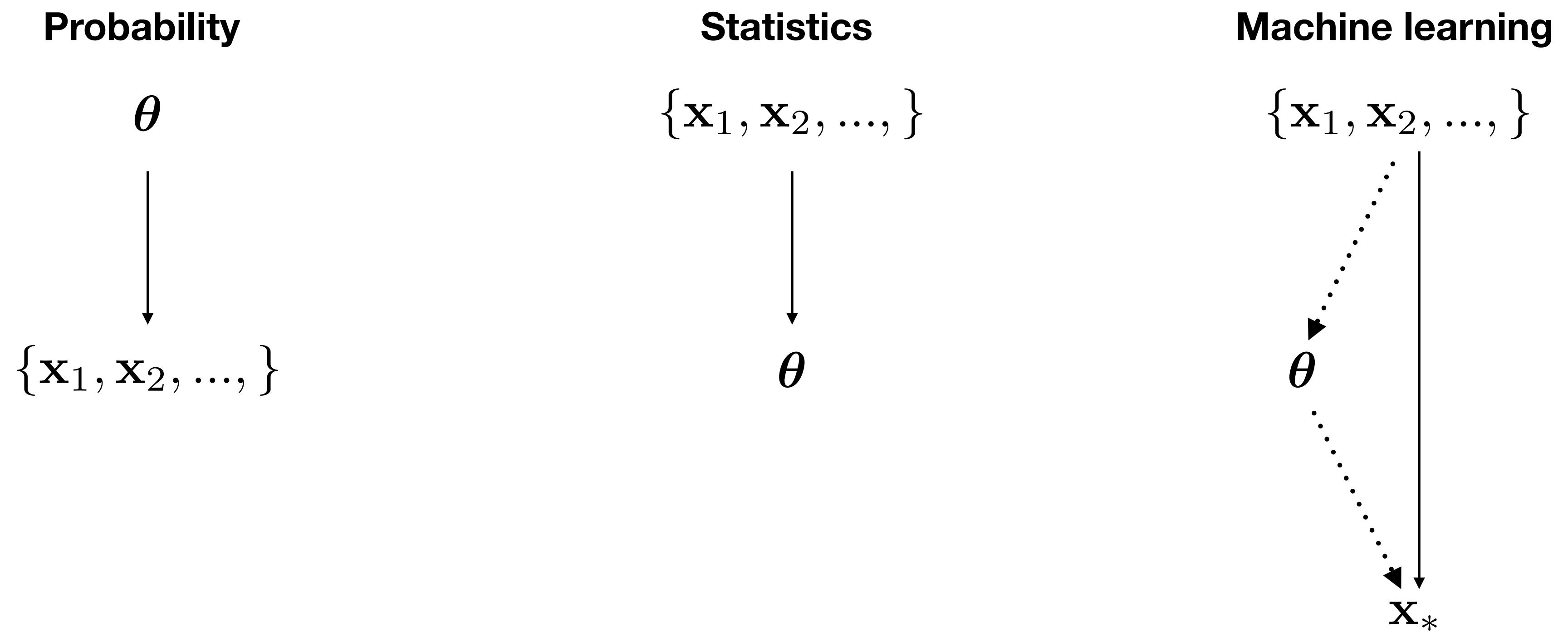


# Probabilistic models

We are mostly concerned with models which look like

$$p(\mathbf{x} \mid \theta)$$

In many case  $\mathbf{x}$  refers to an *observation* and refers to a set of *parameters*.



# Example: The Bent Coin<sup>1</sup>

You are given a bent coin. You flip it  $N$  times. It lands heads  $H$  times.

The probability the coin lands heads is  $\pi$ , what is the posterior  $p(\pi|\mathcal{D})$ ?

$$p(\pi|\mathcal{D}) = \frac{p(\mathcal{D}|\pi)p(\pi)}{p(\mathcal{D})} = \frac{\left[ \prod_{i=1}^N p(x_i|\pi) \right] p(\pi)}{p(\mathcal{D})}$$

We need a prior on  $\pi$ , let's pick a uniform distribution

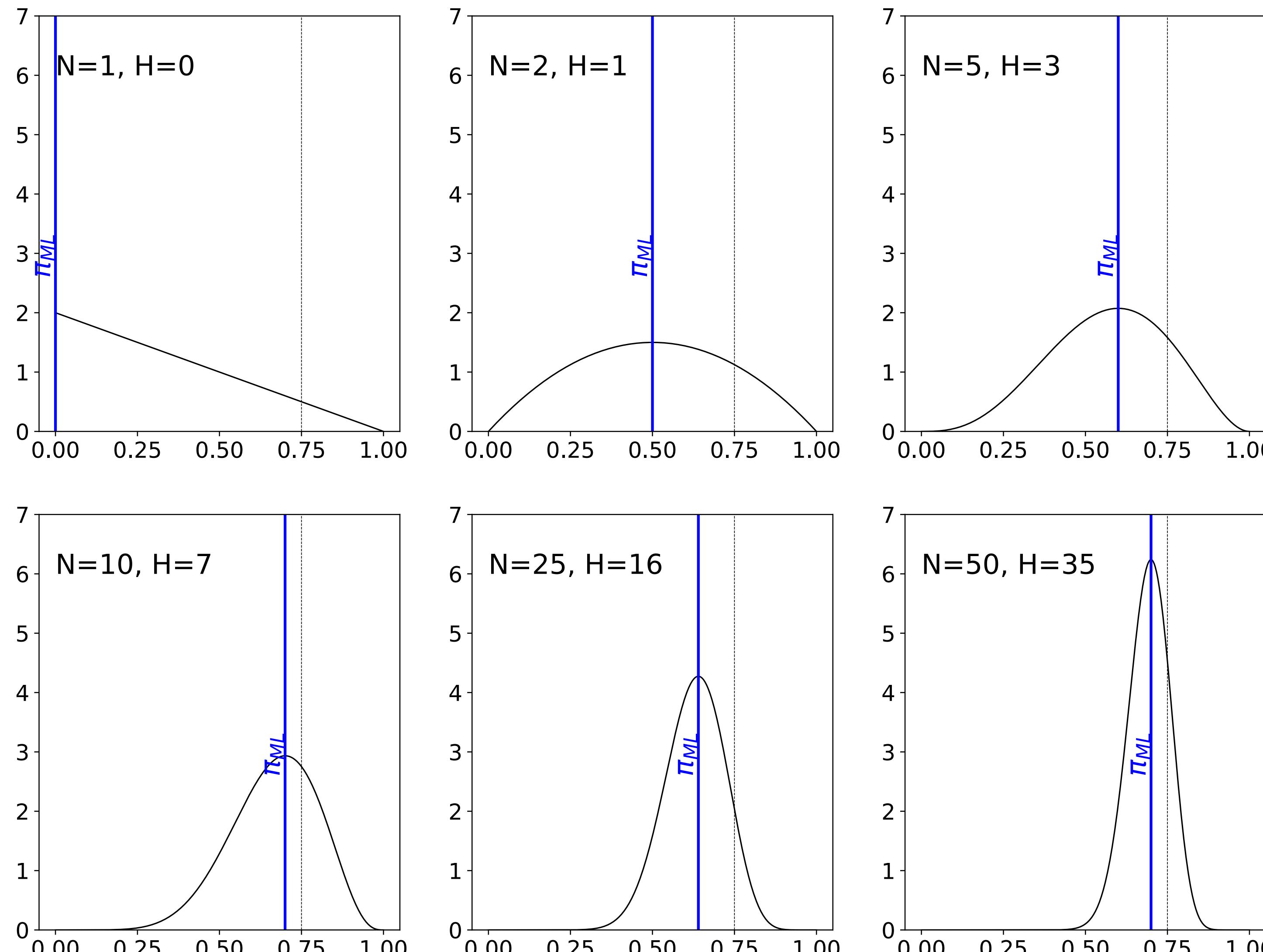
$$\begin{aligned} p(\pi|\mathcal{D}) &= \frac{\left[ \prod_{i=1}^N \pi^{x_i} (1-\pi)^{1-x_i} \right] \mathbb{I}[\pi \in [0, 1]]}{p(\mathcal{D})} \\ &= \frac{\left[ \pi^H (1-\pi)^{N-H} \right] \mathbb{I}[\pi \in [0, 1]]}{p(\mathcal{D})} \\ &= \frac{1}{Z} \pi^H (1-\pi)^{N-H} \end{aligned}$$



<sup>1</sup> This is the original inference problem studied by Thomas Bayes in 1763.

# Example: The Bent Coin<sup>1</sup>

$D = [0 \ 1 \ 0 \ 1 \ 1 \ 1 \ 0 \ 1 \ 1 \ 1 \ 1 \ 0 \ 0 \ 1 \ 1 \ 0 \ 0 \ 1 \ 1 \ 1 \ 1 \ 0 \ 0 \ 1 \ 1 \ 1 \ 1 \ 1 \ 1 \ 1 \ 1 \ 0 \ 1 \ 1 \ 1 \ 1 \ 1 \ 1 \ 0 \ 0 \ 0 \ 1 \ 0 \ 0 \ 1 \ 1]$



# Example: The Bent Coin<sup>1</sup>

Notice how the posterior is ‘less temperamental’ than the likelihood function.

Next we need to figure out the *marginal likelihood*

$$Z = p(\mathcal{D}) = \int p(\mathcal{D}, \pi) d\pi = \int \underbrace{p(\mathcal{D}|\pi)}_{\text{likelihood}} \underbrace{p(\pi)}_{\text{prior}} d\pi.$$

$$\begin{aligned} p(\pi|\mathcal{D}) &= \frac{\left[ \prod_{i=1}^N \pi^{x_i} (1-\pi)^{1-x_i} \right] \mathbb{I}[\pi \in [0, 1]]}{p(\mathcal{D})} \\ &= \frac{\left[ \pi^H (1-\pi)^{N-H} \right] \mathbb{I}[\pi \in [0, 1]]}{p(\mathcal{D})} \\ &= \frac{1}{Z} \pi^H (1-\pi)^{N-H} \end{aligned}$$

The marginal likelihood is an instance of the famous Beta integral<sup>1</sup>

$$p(\mathcal{D}) = \int_0^1 \pi^H (1-\pi)^{N-H} d\pi = B(H+1, N-H+1) = \frac{H!(N-H)!}{(N+1)!}$$

$$p(\pi|\mathcal{D}) = \frac{(N+1)!}{H!(N-H)!} \pi^H (1-\pi)^{N-H}$$

Don’t worry if this integral scares you. It frightens me too! Resources such as Wolfram Alpha, Wikipedia, the Bishop book, and the MacKay book are handy.

<sup>1</sup> $B(x, y) = \int_0^1 t^{x-1} (1-t)^{y-1} dt$

# Example: The Bent Coin<sup>1</sup>

The posterior has the form of a *Beta distribution*

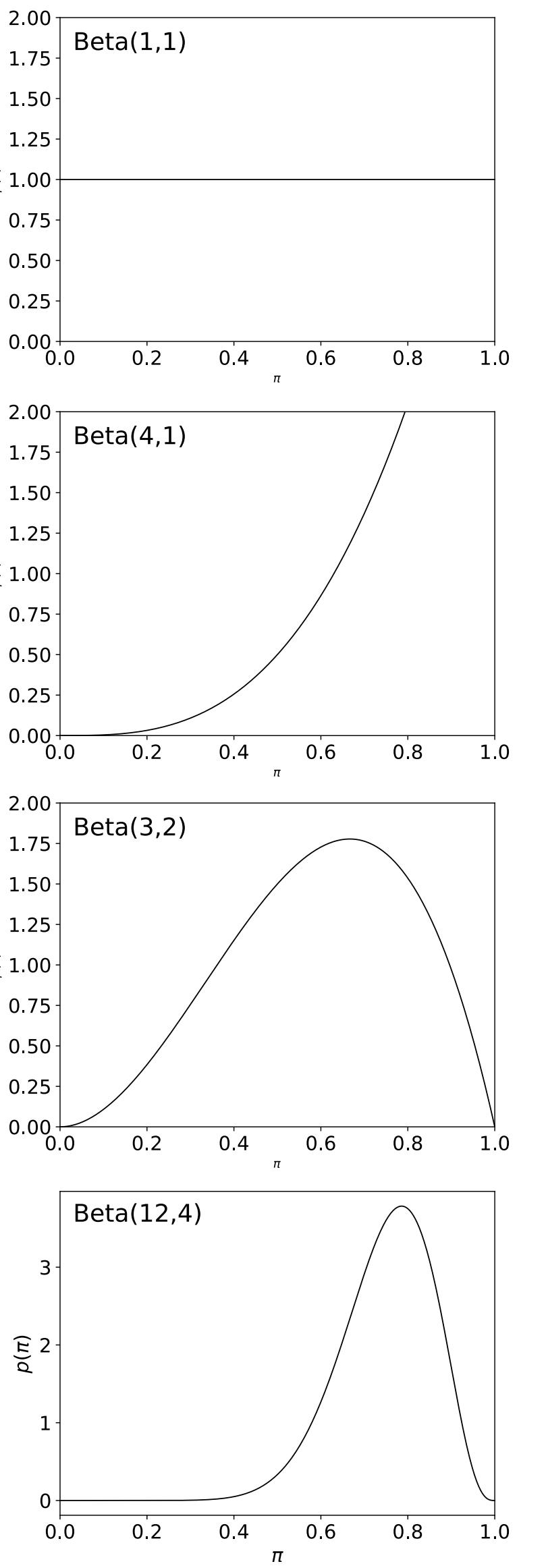
$$\text{Beta}(\pi|\alpha, \beta) = \frac{1}{Z(\alpha, \beta)} \pi^{\alpha-1} (1-\pi)^{\beta-1}, \quad Z(\alpha, \beta) = \frac{\Gamma(\alpha)\Gamma(\beta)}{\Gamma(\alpha+\beta)}$$

- The Beta distribution is a probability distribution over probabilities.
- The two parameters  $\alpha$  and  $\beta$  control the shape of the distribution.
- The *Gamma* function<sup>1</sup> satisfies  $\Gamma(\alpha) = (\alpha - 1)!$  and  $\Gamma(\alpha + 1) = \alpha\Gamma(\alpha)$

**Posterior predictive distribution:** *Laplace's rule of succession*

$$\begin{aligned} p(x_* = \text{head} | \mathcal{D}) &= \int \underbrace{p(x_* = \text{head} | \pi)}_{\text{forward likelihood}} \underbrace{p(\pi | \mathcal{D})}_{\text{posterior}} d\pi \\ &= \int_0^1 \pi \cdot \frac{\pi^H (1-\pi)^{N-H}}{p(\mathcal{D})} d\pi \\ &= \frac{H+1}{N+2} \end{aligned}$$

<sup>1</sup> $\Gamma(\alpha) := \int_0^\infty x^{\alpha-1} e^{-x} dx$



Sometimes, instead of the ML estimate people take the *maximum a posteriori*

$$\begin{aligned}\theta_{\text{MAP}} &= \arg \max_{\theta} p(\theta | \mathcal{D}) = \arg \max_{\theta} \log p(\theta | \mathcal{D}) \\ &= \arg \max_{\theta} \underbrace{\log p(\mathcal{D} | \theta)}_{\text{log-likelihood}} + \underbrace{\log p(\theta)}_{\text{log prior}} - \log p(\mathcal{D})\end{aligned}$$

As data goes to infinity, MAP  $\rightarrow$  ML

$$\theta_{\text{MAP}} = \arg \max_{\theta} \sum_{n=1}^N \log p(x_n | \theta) + \log p(\theta)$$

Infinite data limit: ML = MAP = Bayes'

$$p(x_* | \mathcal{D}) = \int p(x_* | \theta) p(\theta | \mathcal{D}) d\theta \stackrel{N \rightarrow \infty}{=} \int p(x_* | \theta) \delta(\theta - \theta_{\text{ML}}) d\theta = \int p(x_* | \theta) \delta(\theta - \theta_{\text{MAP}}) d\theta$$

<sup>1</sup> The word “conjugate” comes from conjugal, meaning the relationship of a married couple

When the posterior and prior have the same form, we call it *conjugacy*<sup>1</sup>. The *exponential family* admits conjugate pairs:

## Likelihood

$$p(\mathcal{D}|\boldsymbol{\theta}) = \frac{1}{Z(\boldsymbol{\theta})^N} \cdot \exp \left\{ \boldsymbol{\theta}^\top \sum_{n=1}^N \mathbf{t}(x_n) \right\} \cdot \prod_{n=1}^N b(x_n)$$

## Prior

$$p(\boldsymbol{\theta}) = \frac{1}{Z(\boldsymbol{\tau}, \nu)} \cdot \frac{1}{Z(\boldsymbol{\theta})^\nu} \cdot \exp \left\{ \boldsymbol{\theta}^\top \boldsymbol{\tau} \right\}$$

$$p(\boldsymbol{\theta}|\mathcal{D}) \propto p(\mathcal{D}|\boldsymbol{\theta})p(\boldsymbol{\theta})$$

$$= \underbrace{\frac{1}{Z(\boldsymbol{\theta})^N} \exp \left\{ \boldsymbol{\theta}^\top \sum_{n=1}^N \mathbf{t}(x_n) \right\}}_{\text{likelihood}} \left[ \prod_{n=1}^N b(x_n) \right] \cdot \underbrace{\frac{1}{Z(\boldsymbol{\tau}, \nu)} \frac{1}{Z(\boldsymbol{\theta})^\nu} \exp \left\{ \boldsymbol{\theta}^\top \boldsymbol{\tau} \right\}}_{\text{prior}}$$

$$\propto \frac{1}{Z(\boldsymbol{\theta})^{N+\nu}} \exp \left\{ \boldsymbol{\theta}^\top \left( \boldsymbol{\tau} + \sum_{n=1}^N \mathbf{t}(x_n) \right) \right\}$$

Drop terms not containing  $\boldsymbol{\theta}$

Normalisation is easy: just compare with prior

$$\nu \rightarrow N + \nu \quad \boldsymbol{\tau} \rightarrow \boldsymbol{\tau} + \sum_{n=1}^N \mathbf{t}(x_n) \quad Z(\boldsymbol{\tau}, \nu) \rightarrow Z \left( \boldsymbol{\tau} + \sum_{n=1}^N \mathbf{t}(x_n), N + \nu \right)$$

Lightning fast computation:  $O(N)$

<sup>1</sup> The word “conjugate” comes from conjugal, meaning the relationship of a married couple



# Model Comparison

Kabupaten Tana Toraja,

Say I have some data, how do I pick a likelihood and a prior, aka models? Pick a few different *models*, and then find the posterior distribution over the models given the data.

$$p(\mathcal{M}_i | \mathcal{D}) \propto p(\mathcal{D} | \mathcal{M}_i) p(\mathcal{M}_i)$$

Typically, we just want **one** model → MAP inference

Furthermore, the *model prior* is usually flat → MAP = ML

$$\arg \max_{\mathcal{M}_i} p(\mathcal{D} | \mathcal{M}_i) = \arg \max_{\mathcal{M}_i} \int p(\mathcal{D} | \boldsymbol{\theta}) p(\boldsymbol{\theta} | \mathcal{M}_i) d\boldsymbol{\theta}$$

But hang on, this is just the marginal likelihood/evidence!

$$p(\boldsymbol{\theta} | \mathcal{D}, \mathcal{M}_i) = \frac{p(\mathcal{D} | \boldsymbol{\theta}, \mathcal{M}_i) p(\boldsymbol{\theta} | \mathcal{M}_i)}{p(\mathcal{D} | \mathcal{M}_i)}$$

Best model has highest evidence!

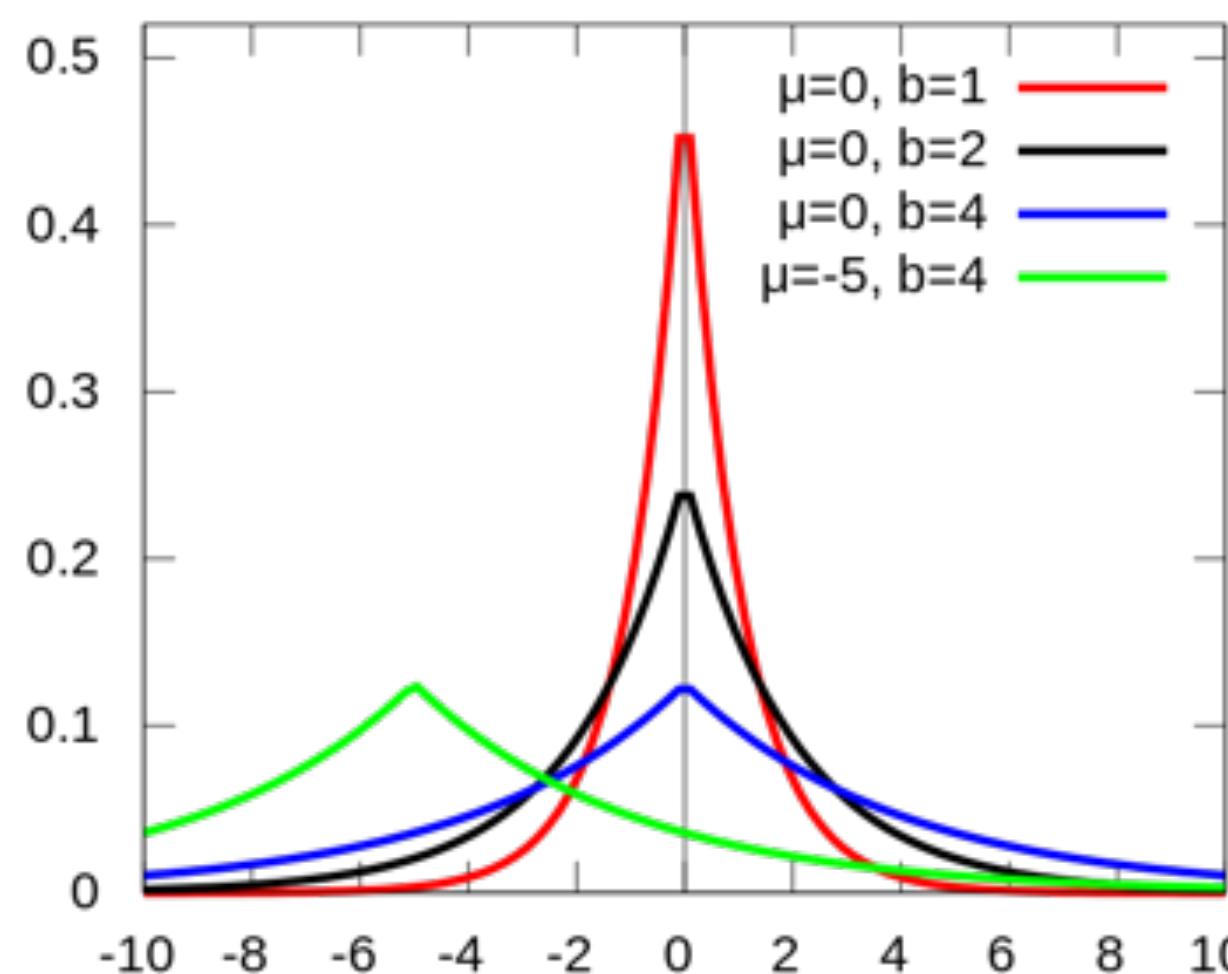
# Example: Laplace versus Gauss

We have zero mean and unit variance data  $\{x_1, \dots, x_N\}$ .

Laplacian

$$p(x|\mathcal{M}_1) = \frac{1}{\sqrt{2}} \exp \left\{ -\sqrt{2}|x| \right\}$$

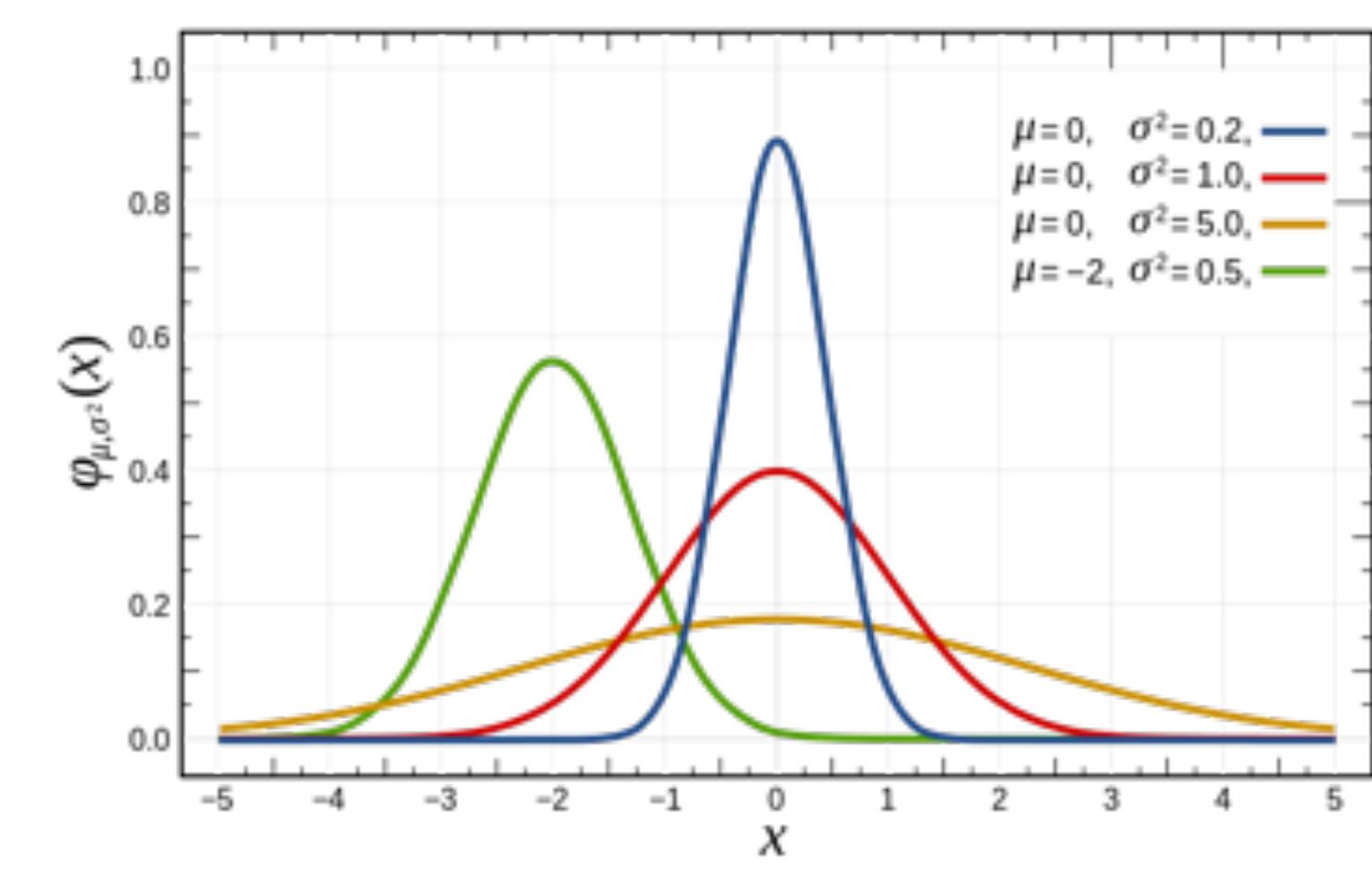
$$\begin{aligned} \log p(\mathcal{D}|\mathcal{M}_1) &= \sum_{n=1}^N \log \left( \frac{1}{\sqrt{2}} \exp \left\{ -\sqrt{2}|x_n| \right\} \right) \\ &= N \log \frac{1}{\sqrt{2}} - \sqrt{2} \sum_{i=1}^N |x_i| \end{aligned}$$



Gaussian

$$p(x|\mathcal{M}_2) = \frac{1}{\sqrt{2\pi}} \exp \left\{ -\frac{x^2}{2} \right\}$$

$$\begin{aligned} \log p(\mathcal{D}|\mathcal{M}_2) &= \sum_{i=1}^N \log \left( \frac{1}{\sqrt{2\pi}} \exp \left\{ -\frac{x_i^2}{2} \right\} \right) \\ &= N \log \frac{1}{\sqrt{2\pi}} - \frac{1}{2} \sum_{i=1}^N x_i^2 \end{aligned}$$



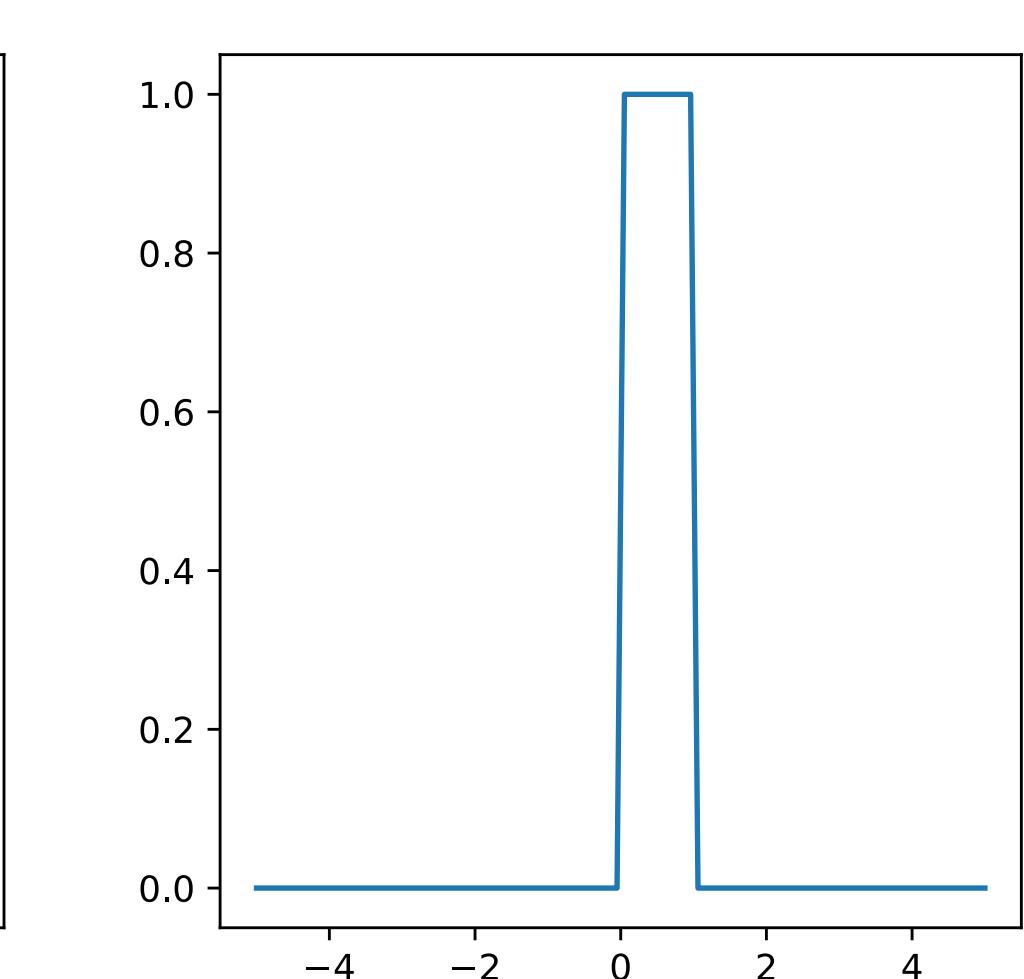
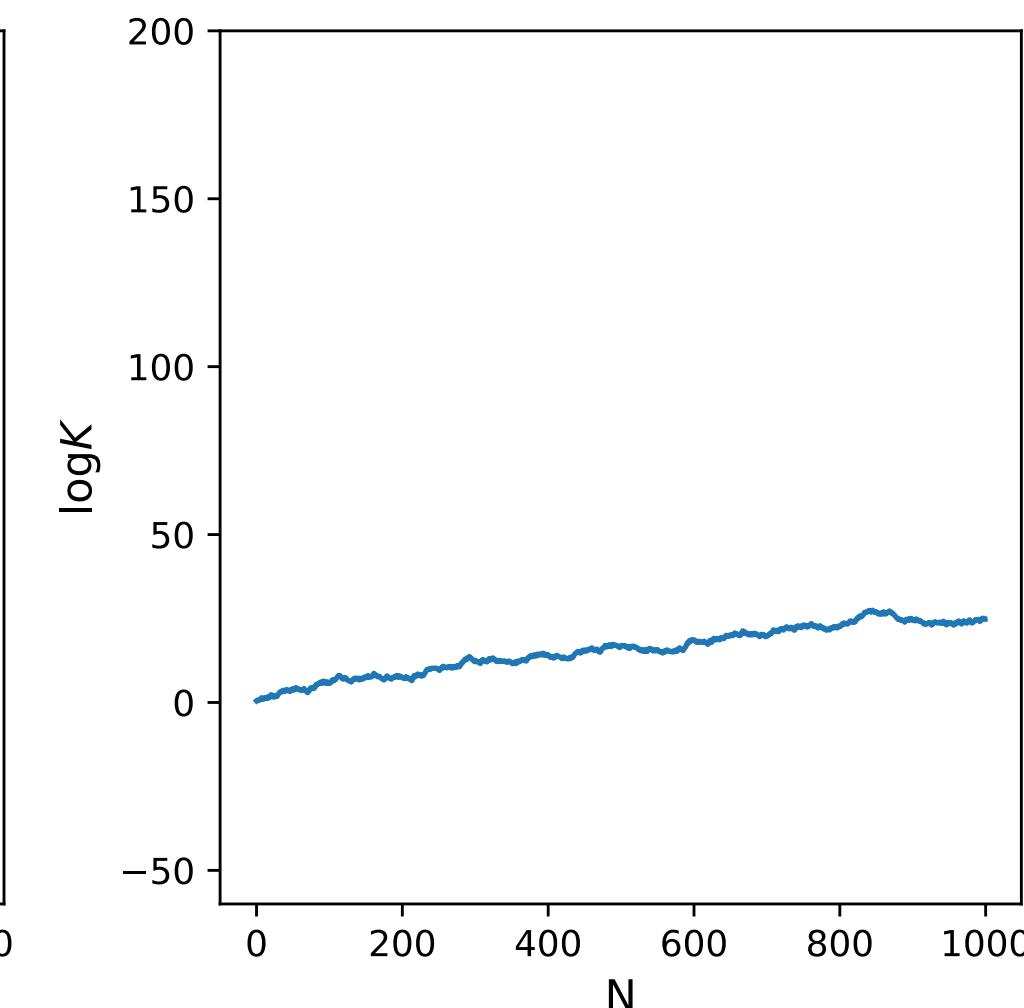
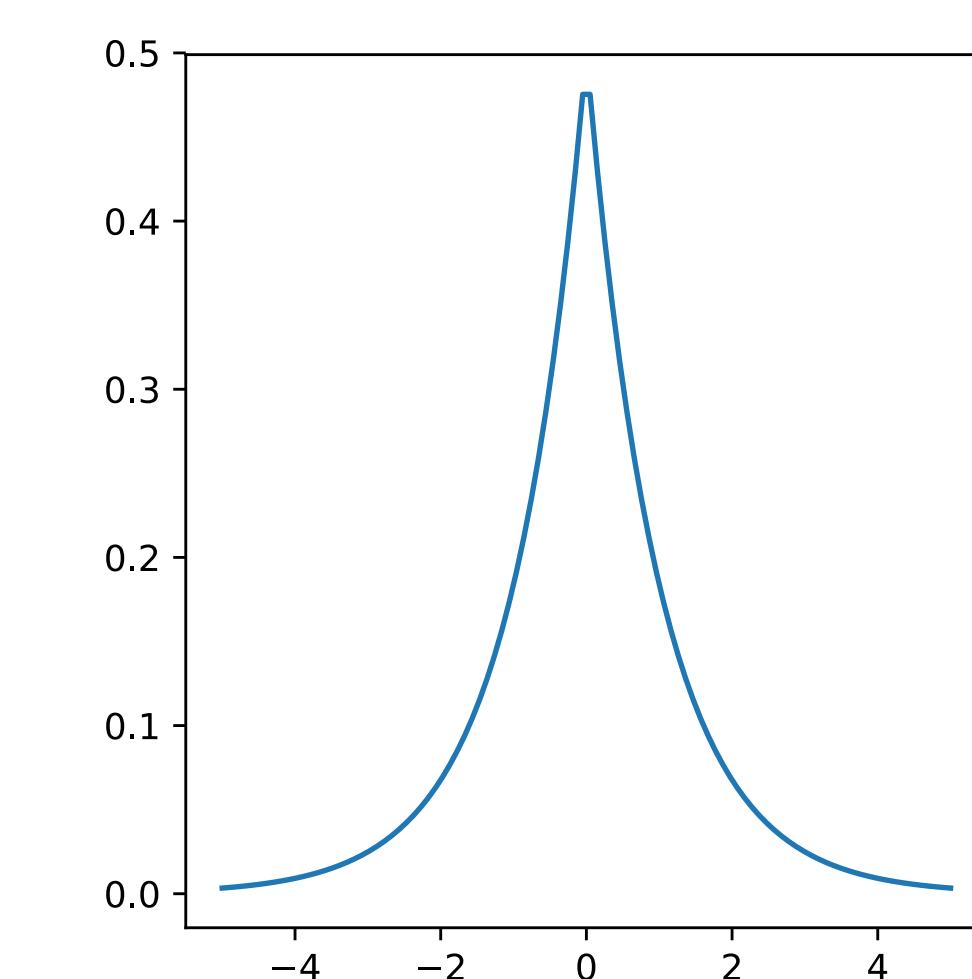
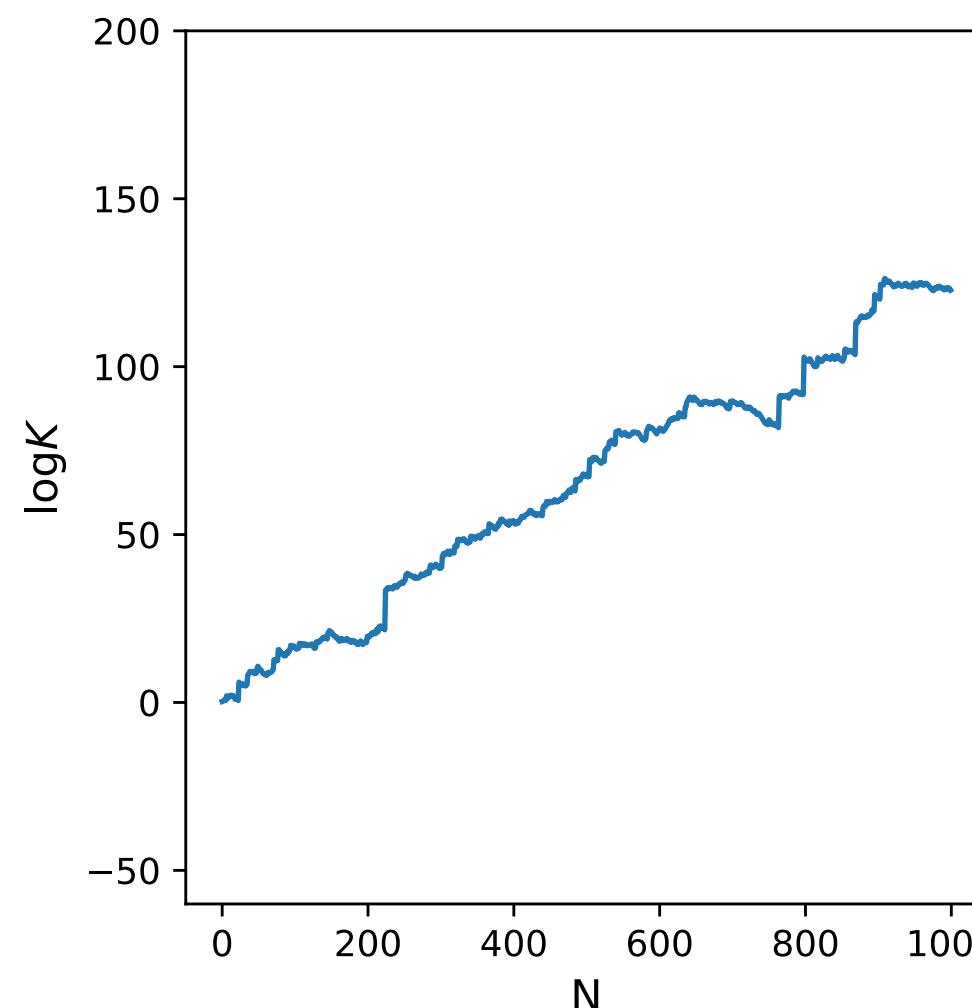
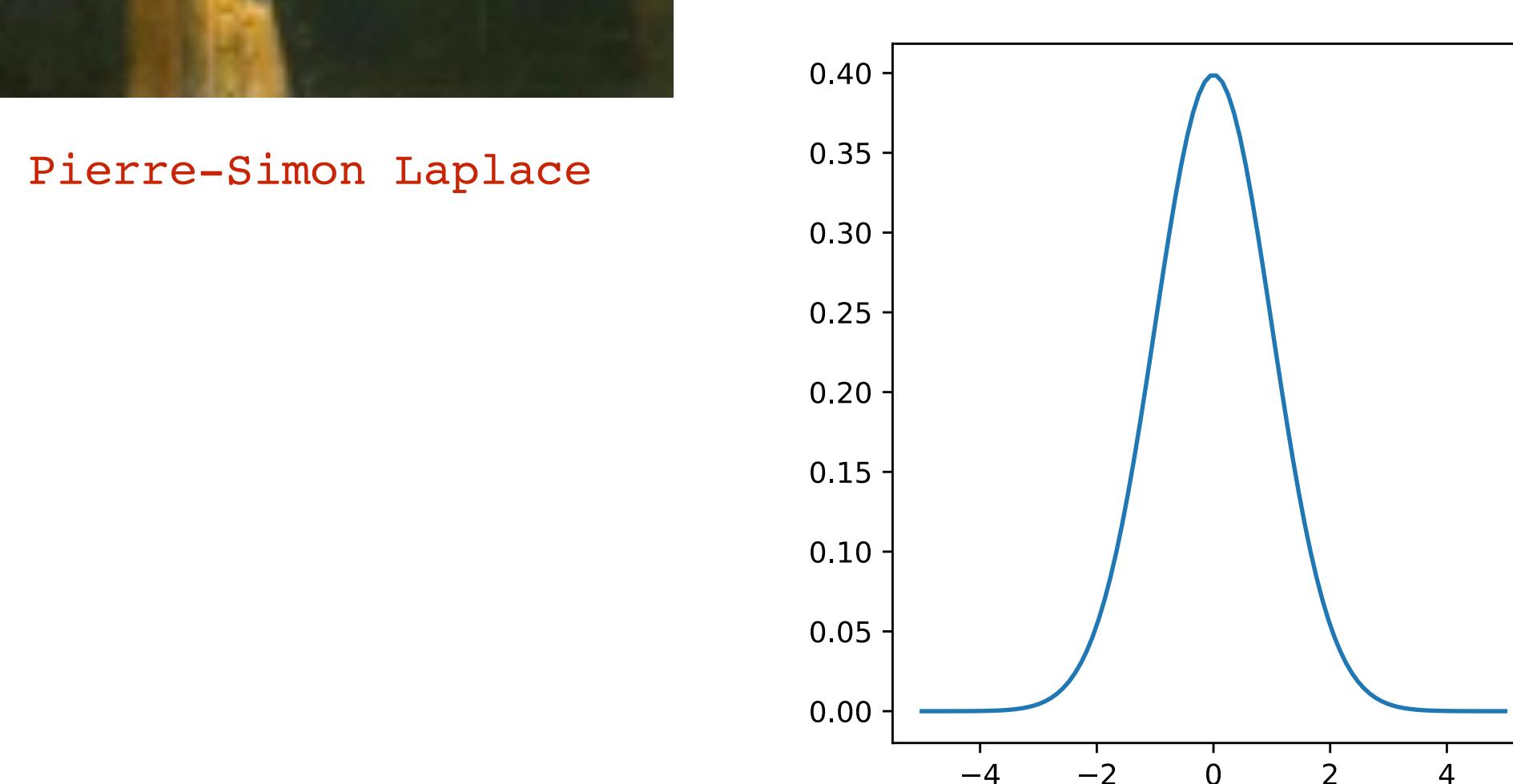
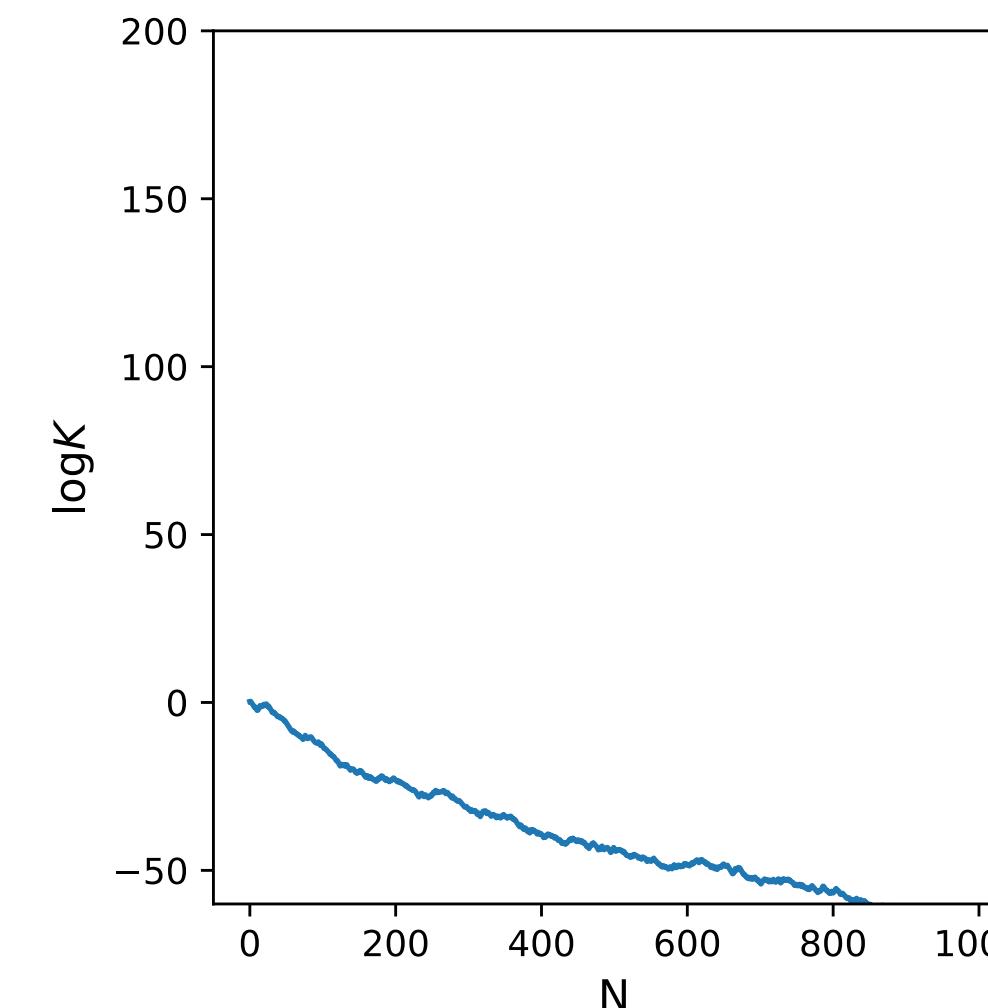
# Example: Laplace versus Gauss

Log Bayes' factor

$$\log K = \log \frac{p(\mathcal{D}|\mathcal{M}_1)}{p(\mathcal{D}|\mathcal{M}_2)}$$



Pierre-Simon Laplace



Carl Friedrich Gauss

# Type-II Maximum likelihood

Can optimise model *hyperparameters* too

$$p(\boldsymbol{\theta}|\mathcal{D}, \boldsymbol{\tau}) = \frac{p(\mathcal{D}|\boldsymbol{\theta}, \boldsymbol{\tau})p(\boldsymbol{\theta}|\boldsymbol{\tau})}{p(\mathcal{D}|\boldsymbol{\tau})}.$$

$$\boldsymbol{\tau}^* = \arg \max_{\boldsymbol{\tau}} p(\boldsymbol{\theta}|\mathcal{D}, \boldsymbol{\tau})$$

This goes by the name of *Type-II maximum likelihood*, *empirical Bayes*, or the *evidence approximation*

It can be difficult to compute  $\boldsymbol{\tau}^*$  in closed-form and the optimization landscape is typically highly multimodal. We will see an example of this in linear regression

# Regression

Candi Borobudur

# Linear regression

In high school you learnt the equation of a straight line

$$f(\mathbf{x}) = w\mathbf{x} + b = [w \quad b] \begin{bmatrix} \mathbf{x} \\ 1 \end{bmatrix} = \mathbf{w}^\top \mathbf{x}$$

$w$  is the gradient and  $b$  is the  $y$ -intercept.

In general, a function  $f$  is *linear* if

$$f(ax_1 + x_2) = af(x_1) + f(x_2)$$

Our model should account for *noise*

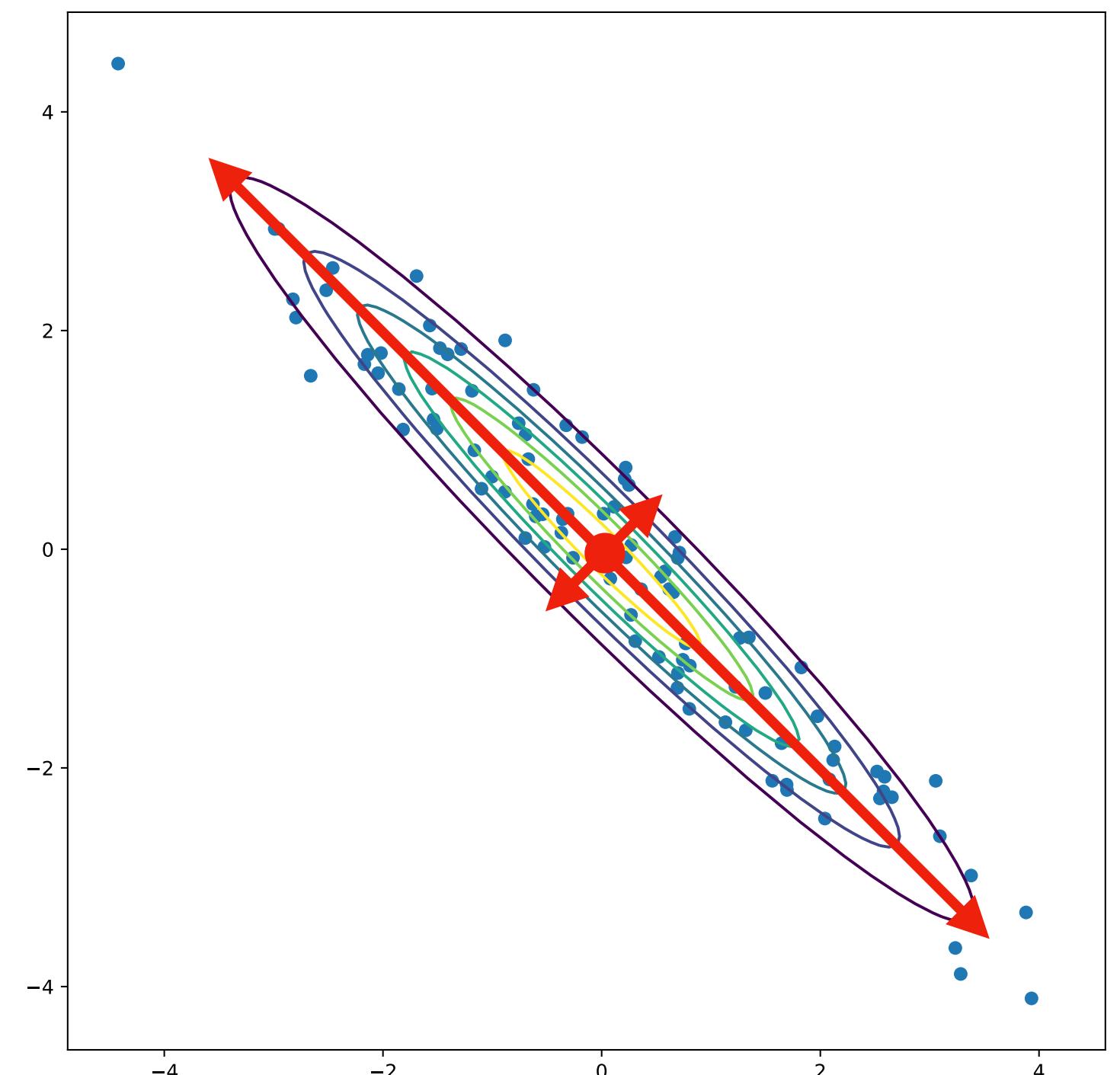
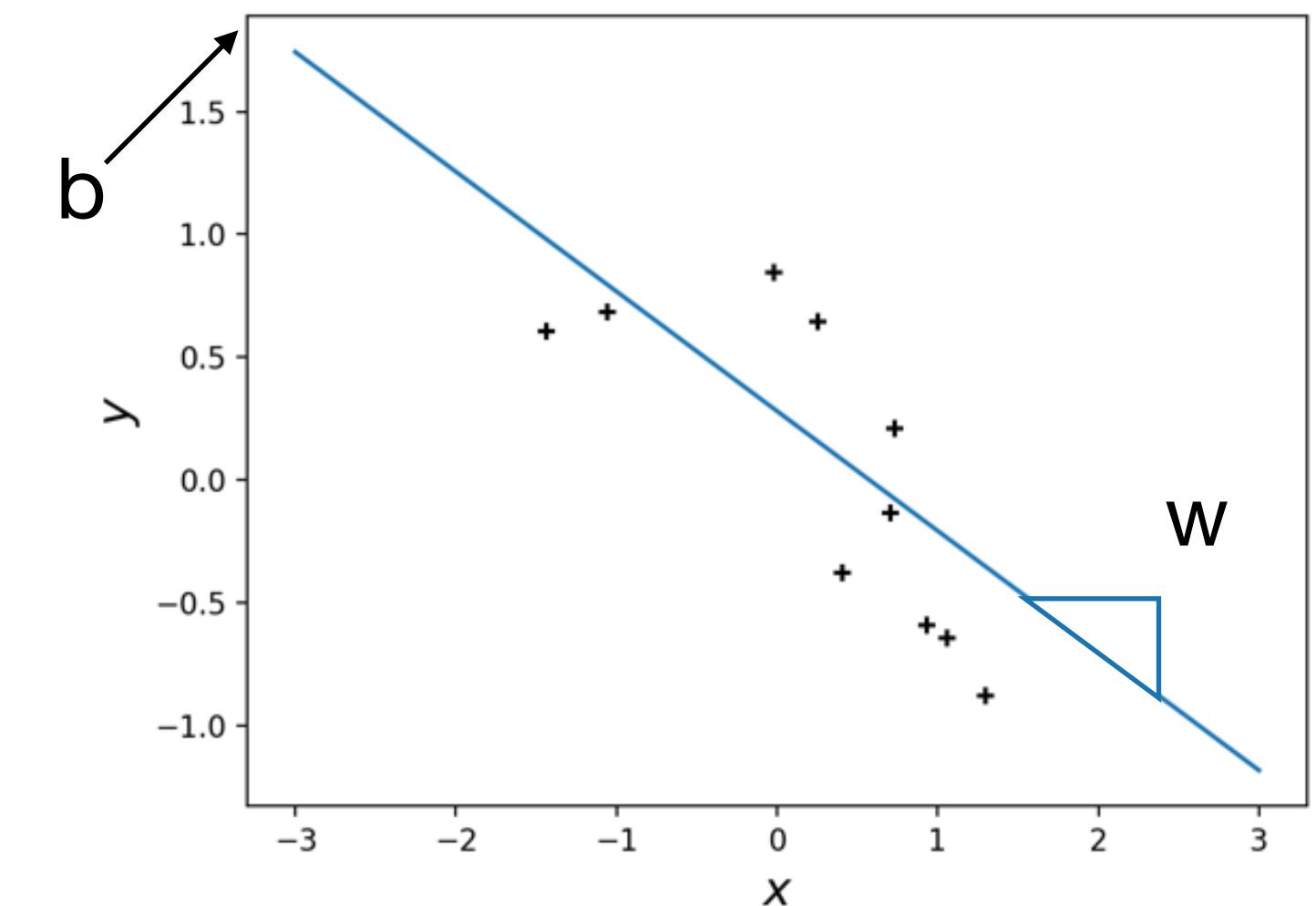
$$y = f(\mathbf{x}) + \epsilon, \quad \epsilon \sim \mathcal{N}(\epsilon \mid 0, \sigma^2)$$

Mean

Covariance

The *Normal/Gaussian distribution*

$$\mathcal{N}(\mathbf{x} \mid \boldsymbol{\mu}, \boldsymbol{\Sigma}) = |2\pi\boldsymbol{\Sigma}|^{-\frac{1}{2}} \exp \left\{ -\frac{1}{2}(\mathbf{x} - \boldsymbol{\mu})^\top \boldsymbol{\Sigma}^{-1} (\mathbf{x} - \boldsymbol{\mu}) \right\}$$



# Linear regression

Dataset

$$\mathbf{y} = [y_1, y_2, \dots, y_N]^\top \quad \mathbf{X} = [\mathbf{x}_1, \mathbf{x}_2, \dots, \mathbf{x}_N]$$

Likelihood

$$p(\mathbf{y} | \mathbf{X}, \mathbf{w}, \boldsymbol{\tau}) = \prod_{n=1}^N \mathcal{N}(y_n | x_n, \mathbf{w}, \boldsymbol{\tau}) = \prod_{n=1}^N \frac{1}{\sqrt{2\pi\sigma^2}} \exp \left\{ -\frac{(y_n - \mathbf{x}_n^\top \mathbf{w})^2}{2\sigma^2} \right\}$$

$$= \frac{1}{(2\pi\sigma^2)^{N/2}} \exp \left\{ -\frac{1}{2\sigma^2} \|\mathbf{y} - \mathbf{X}^\top \mathbf{w}\|^2 \right\} = \mathcal{N}(\mathbf{y} | \mathbf{X}^\top \mathbf{w}, \sigma^2 \mathbf{I})$$

Hyperparameters

Posterior

Likelihood

Prior

Conjugacy

$$p(\mathbf{w} | \mathbf{y}, \mathbf{X}, \boldsymbol{\tau}) = \frac{\mathcal{N}(\mathbf{y} | \mathbf{w}^\top \mathbf{X}, \sigma^2 \mathbf{I}) \mathcal{N}(\mathbf{w} | \mathbf{0}, \boldsymbol{\Sigma}_0)}{\int \mathcal{N}(\mathbf{y} | \mathbf{w}^\top \mathbf{X}, \sigma^2 \mathbf{I}) \mathcal{N}(\mathbf{w} | \mathbf{0}, \boldsymbol{\Sigma}_0) d\mathbf{w}} = \mathcal{N}(\mathbf{w} | \boldsymbol{\mu}_N, \boldsymbol{\Sigma}_N)$$

$$\boldsymbol{\mu}_N = \frac{1}{\sigma^2} \boldsymbol{\Sigma}_N \mathbf{X} \mathbf{y}, \quad \boldsymbol{\Sigma}_N = (\boldsymbol{\Sigma}_0^{-1} + \sigma^{-2} \mathbf{X} \mathbf{X}^\top)^{-1}$$

Covariance weighted mean

Data-weighted mean of precisions

Precision = Covariance<sup>-1</sup>

# The PPD and Type-II ML

Let's find the posterior predictive distribution: it just happens to be Gaussian

$$p(y_* | \mathbf{x}_*, \mathcal{D}, \tau) = \int \mathcal{N}(y_* | \mathbf{x}_*^\top \mathbf{w}, \sigma^2) \mathcal{N}(\mathbf{w} | \boldsymbol{\mu}_N, \boldsymbol{\Sigma}_N) d\mathbf{w} = \mathcal{N}(y_* | \mathbf{x}_*^\top \boldsymbol{\mu}_N, \sigma^2 + \mathbf{x}_*^\top \boldsymbol{\Sigma}_N \mathbf{x}_*)$$

The PPD uses the mean of the posterior distribution. The variance is a sum of  $\sigma^2$  and  $\mathbf{x}_*^\top \boldsymbol{\Sigma}_N \mathbf{x}_*$ . The variance gets larger the larger the data term is.

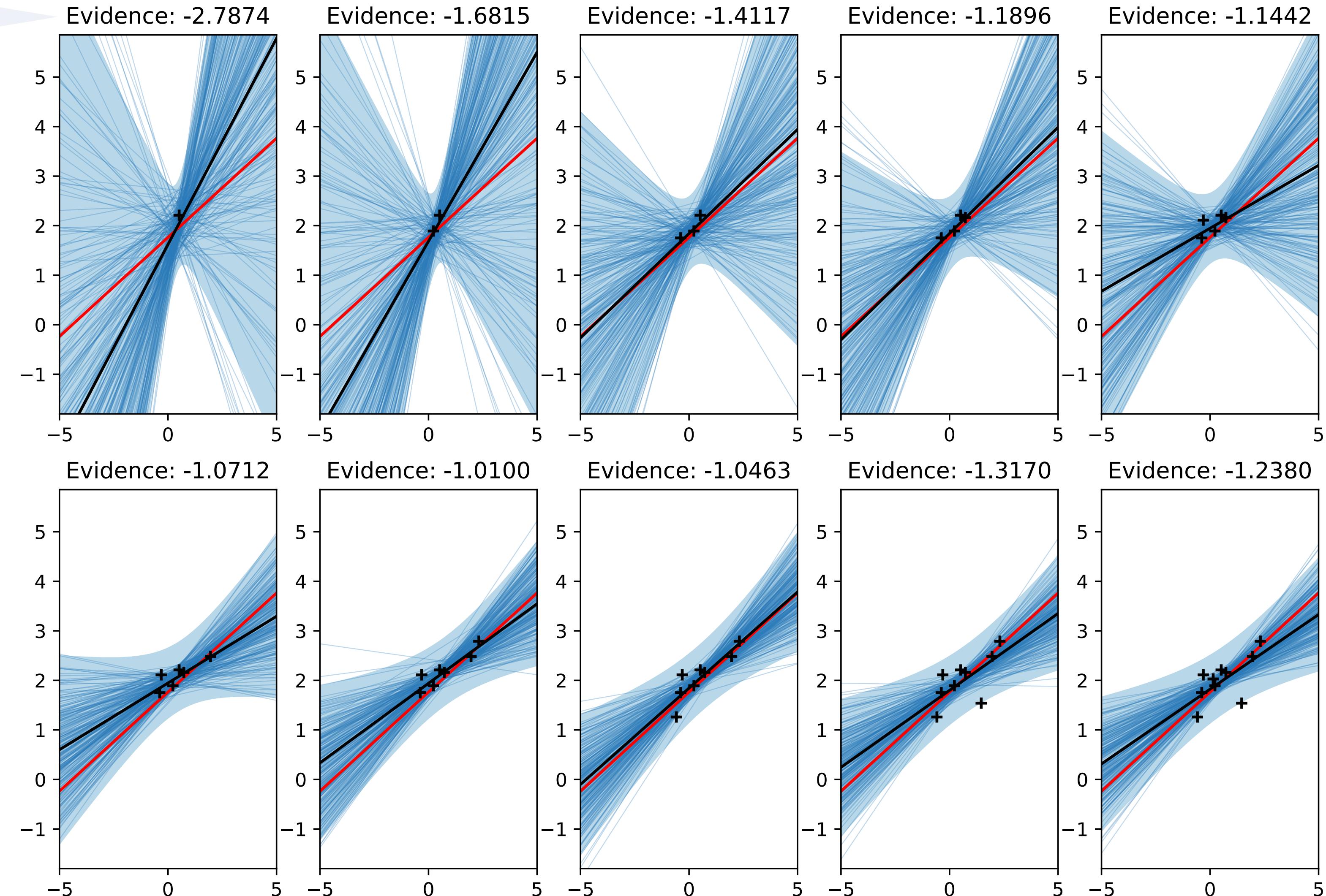
## Type-II Maximum Likelihood

Notice how the PPD is a function of  $\tau$ . We can optimise  $\tau$ : this is called *Type-II Maximum likelihood, the evidence approximation, or empirical Bayes*

$$\begin{aligned} \log p(\mathcal{D} | \sigma^2, \boldsymbol{\Sigma}_0) &= \log p(\mathbf{y} | \mathbf{X}, \sigma^2, \boldsymbol{\Sigma}_0) = \log \int \mathcal{N}(\mathbf{y} | \mathbf{X}^\top \mathbf{w}, \sigma^2 \mathbf{I}) \mathcal{N}(\mathbf{w} | \mathbf{0}, \boldsymbol{\Sigma}_0) d\mathbf{w} \\ &= \log \mathcal{N}(\mathbf{y} | \mathbf{0}, \underbrace{\sigma^2 \mathbf{I} + \mathbf{X}^\top \boldsymbol{\Sigma}_0 \mathbf{X}}_{\mathbf{K}}) \\ &= -\frac{1}{2} \mathbf{y}^\top \mathbf{K}^{-1} \mathbf{y} - \frac{1}{2} \log |\mathbf{K}| - \frac{N}{2} \log 2\pi \end{aligned}$$

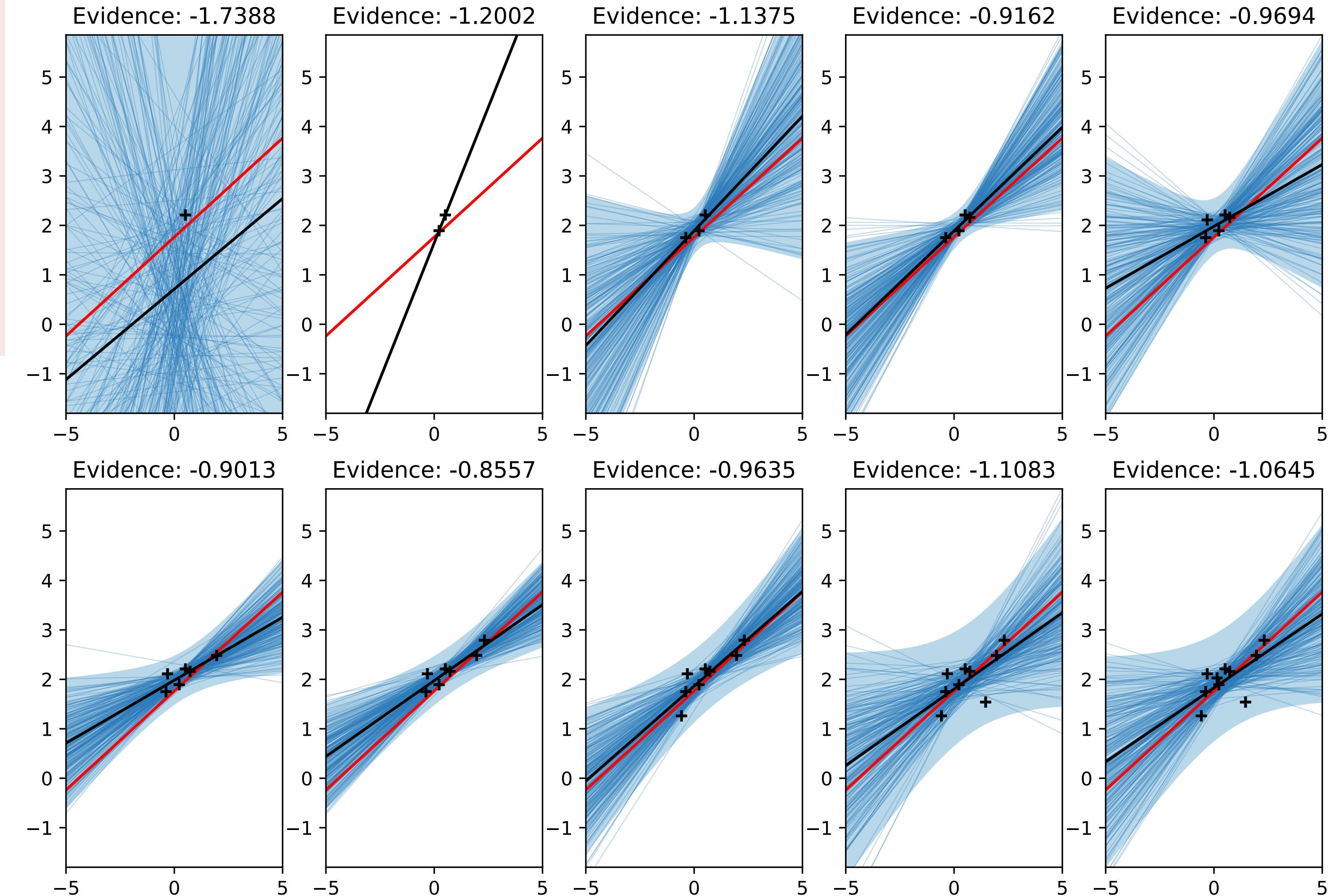
# Example: $\sigma^2 = 0.1$ , $\Sigma_0 = \mathbf{I}$

Sample normalized  
evidence



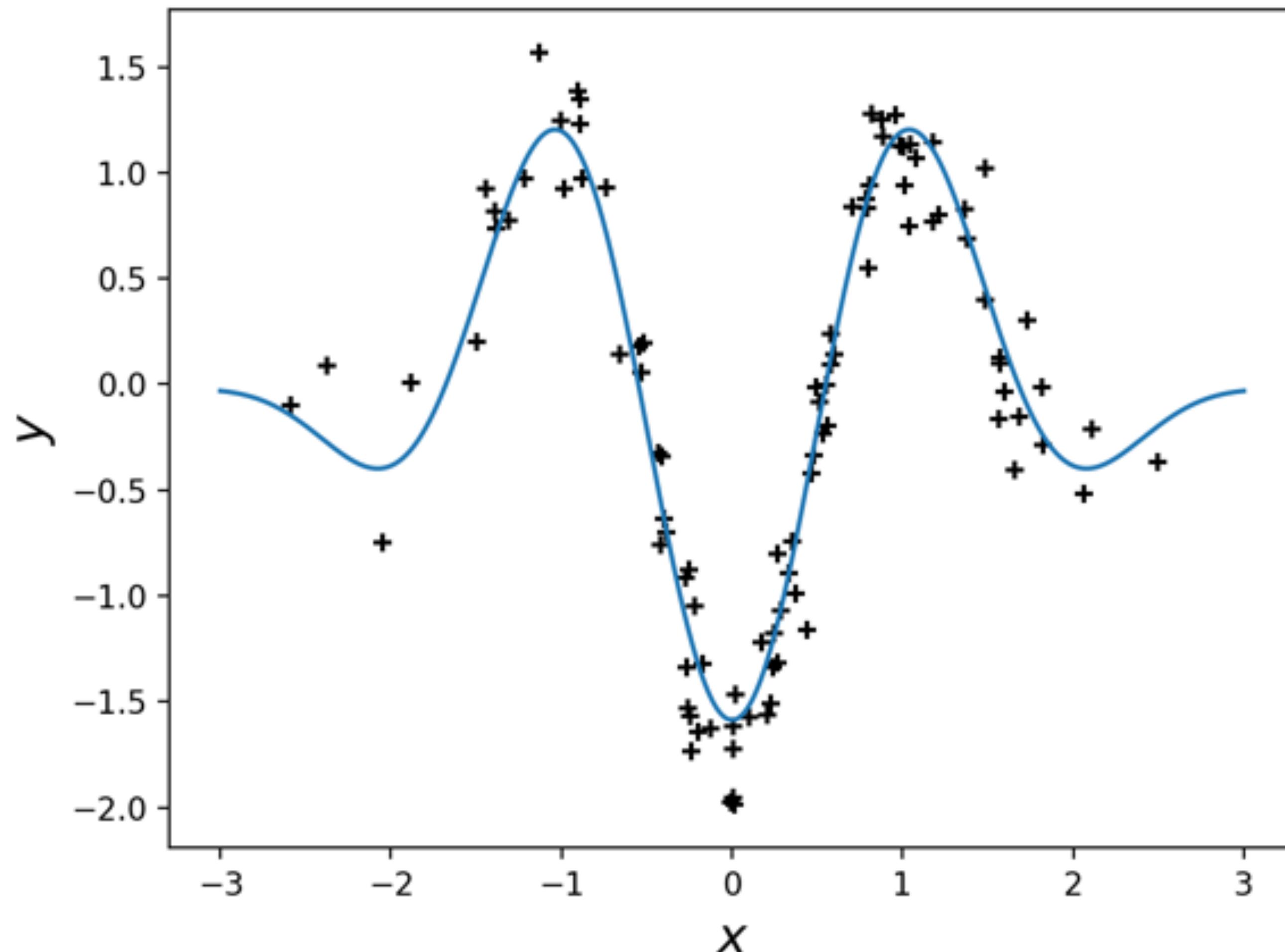
# Example: evidence approximation

The evidence approximation is maximum likelihood at the end of the day, so it can be unstable for small datasets



# Featurisation

Could we describe non-straight curves using linear maths?



## Featurisation

$$\begin{aligned}f(x) &= w_0 + w_1x + w_2x^2 + \dots + w_Mx^M \\&= \sum_{m=0}^M w_m \phi_m(x) \\&= \boldsymbol{\phi}(\mathbf{x})^\top \mathbf{w}\end{aligned}$$

Nonlinear in  $x$ , linear in  $\mathbf{w}$

# Featurisation

Featurisation makes everything easy, because our method remains mostly unchanged. All we have to do is substitute  $\mathbf{X} \mapsto \Phi$

## Likelihood

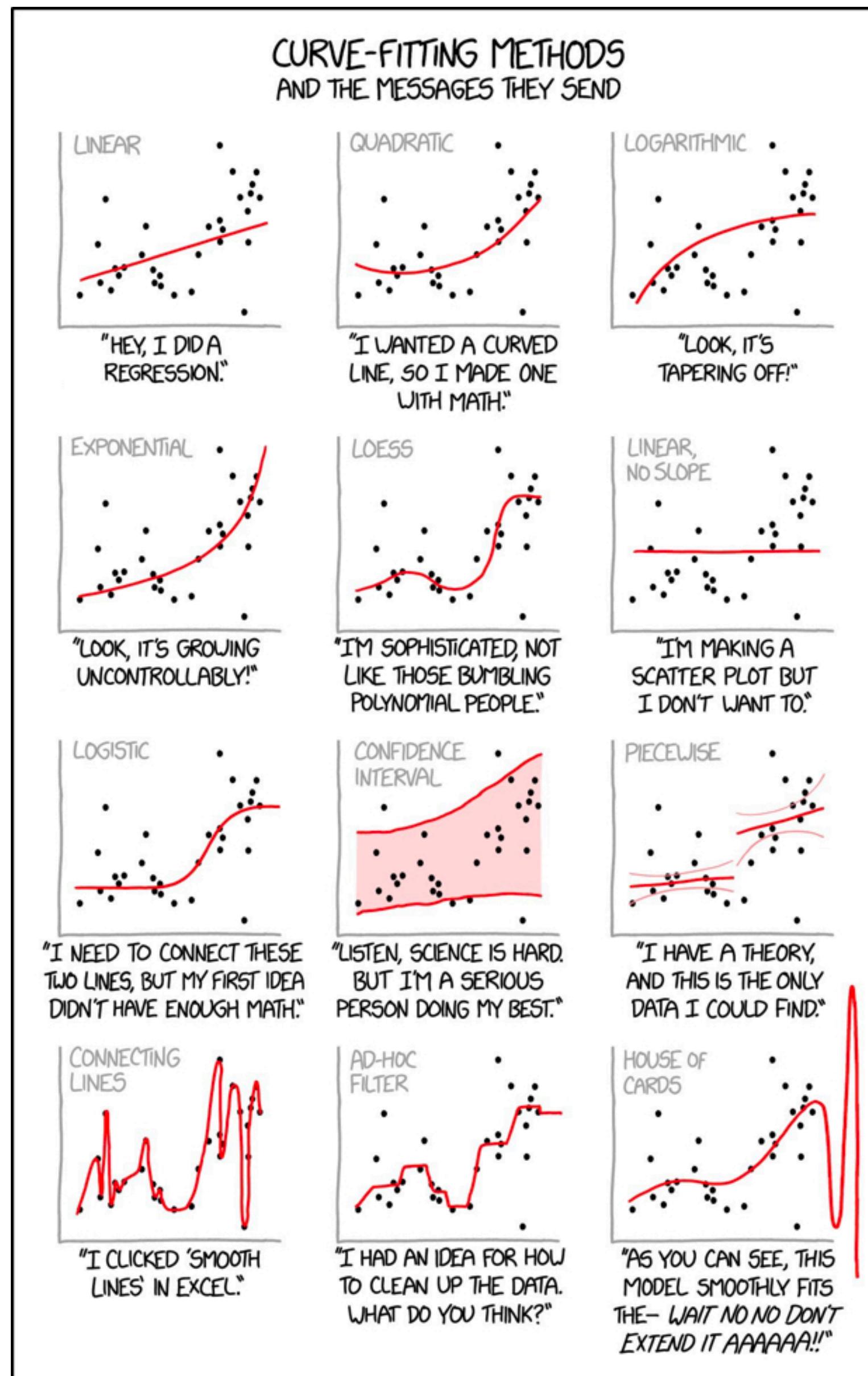
$$p(\mathbf{y} | \mathbf{X}, \mathbf{w}, \tau) = \prod_{n=1}^N \frac{1}{\sqrt{2\pi\sigma^2}} \exp \left\{ -\frac{(y_n - \mathbf{x}_n^\top \mathbf{w})^2}{2\sigma^2} \right\} = \mathcal{N}(\mathbf{y} | \mathbf{X}^\top \mathbf{w}, \sigma^2 \mathbf{I})$$

$$p(\mathbf{y} | \mathbf{X}, \mathbf{w}, \tau) = \prod_{n=1}^N \frac{1}{\sqrt{2\pi\sigma^2}} \exp \left\{ -\frac{(y_n - \phi(x_n)^\top \mathbf{w})^2}{2\sigma^2} \right\} = \mathcal{N}(\mathbf{y} | \Phi^\top \mathbf{w}, \sigma^2 \mathbf{I})$$

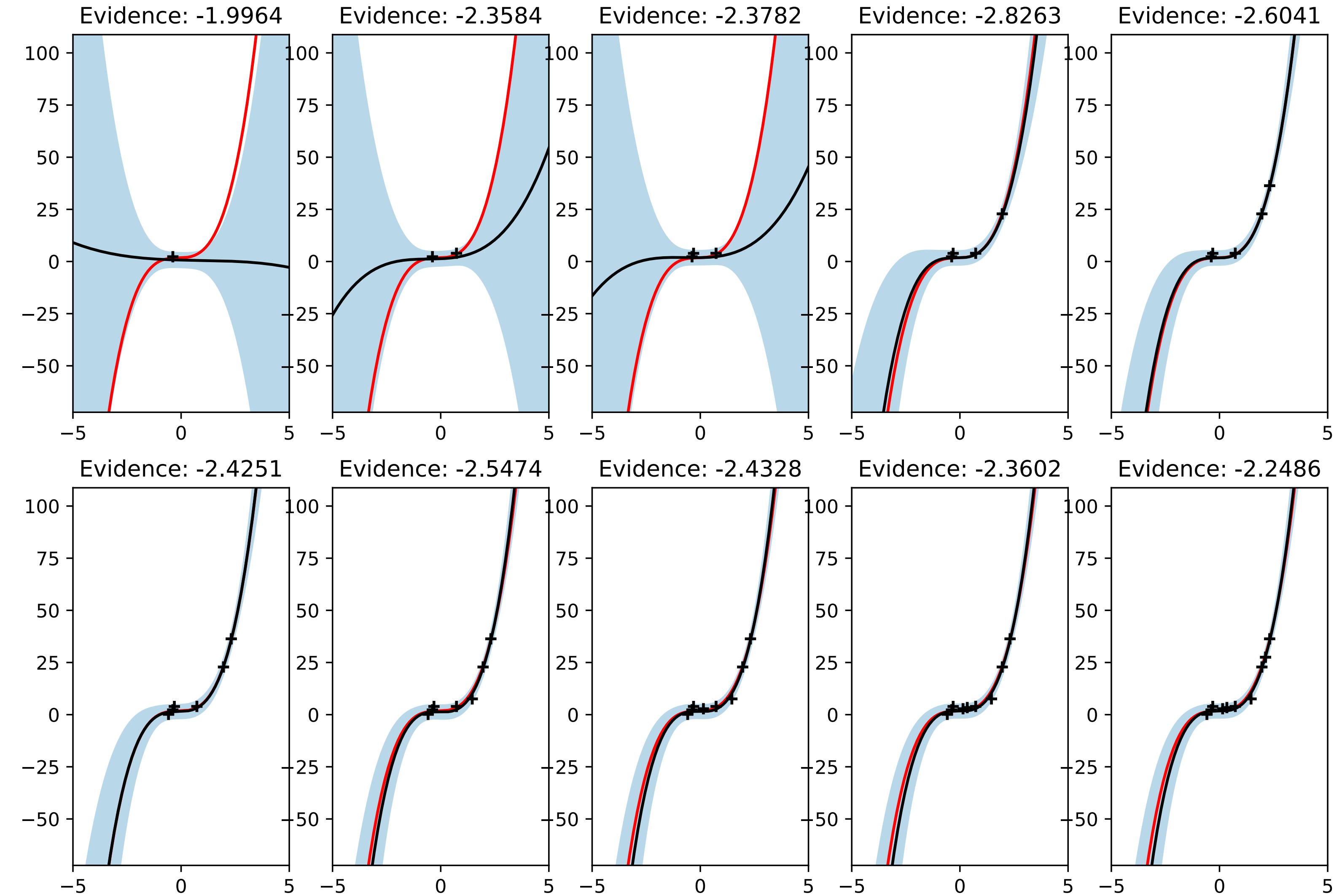
Designing features was a huge research area a decade ago. These are called **handcrafted features**. In Lecture 2, we will show how to learn these..

$$\phi(x) = [1, \cos(x), \sin(x), \dots, \cos(Mx), \sin(Mx)]$$

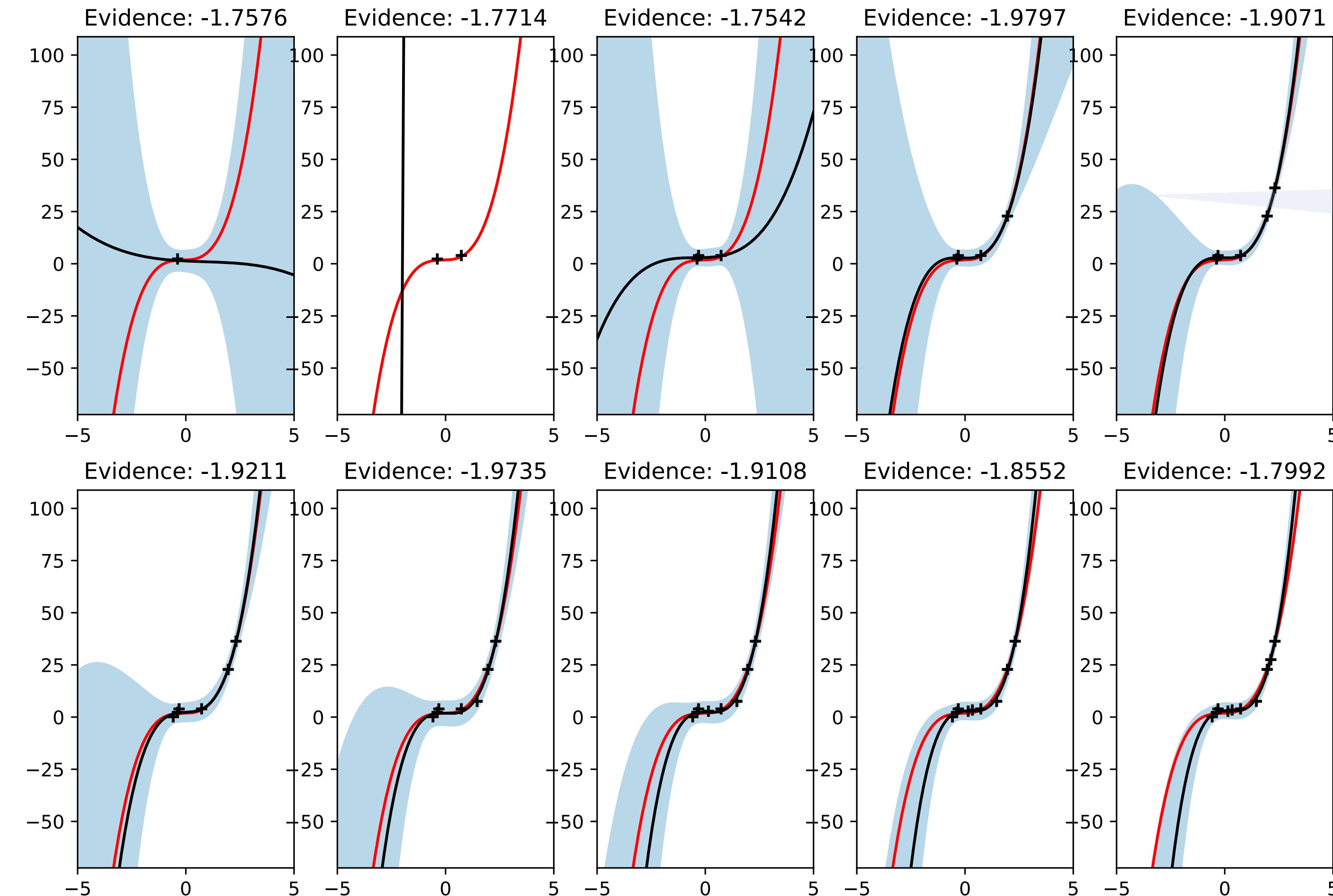
$$\phi(x) = [\exp(-\gamma|x - x_0|^2), \exp(-\gamma|x - x_1|^2), \dots, \exp(-\gamma|x - x_M|^2)]$$



# Example: Cubic Bayes



# Example: Cubic evidence approx.



Notice increased uncertainty far from the data

# Summary

## This lecture: Machine Learning Basics

- What is Machine Learning?
- Probability Theory recap
- Modeling paradigms: Probability theory
- Generative models
- Statistics: Maximum likelihood
- Bayesian Inference
- Prediction
- Conjugacy
- Model Comparison
- Polynomial Regression

## Next lecture: Deep Learning Basics

- Learning Theory
  - Overfitting, Generalization, Bias—variance decomposition
- Classification
  - Logistic regression
  - Gradient-based learning
    - Steepest Descent, Newton's Method
- Deep Learning
  - Backpropagation, Stochastic Gradient Descent
  - The Exploding—Vanishing Gradients problem
  - Initialization, Activation Normalization
- Cross-validation
- Neural Architectures
  - CNN, LSTMs, Graph NNs