

Probabilistic machine learning

Karsten W. Jacobsen

Practical stuff

- Notes about probabilistic machine learning available at DTU Learn
- Jupyter/Python notebook available with code necessary to produce all figures - but try to make your own code!
- The lecture slides are also available, and will be updated as we go along
- Switching between lectures and exercises - we have to figure out what works the best
- Lectures:
 - Speak up to ask questions or
 - Use chat to ask questions. Andreas Vishart will keep an eye on the chat and alert me.
- Exercises:
 - Join in groups in “breakout rooms”
 - Use chat to get Andreas Vishart or me to join you and help
 - Use “blue hand” to signal that we can move on

Probabilistic machine learning - overview

- Probability theory:
 - Intro, function fitting, probability theory, sum and product rules (pgs 1-14)
 - Bayes' theorem, models and data, model selection, probability distributions
 - Applications of Bayes' theorem
- Machine learning:
 - Methods of least squares, linear basis function models, regularized least squares, prior and posterior
 - Kernel regression and Gaussian processes
 - GP and hyperparameters
- Small project:
 - Perovskites

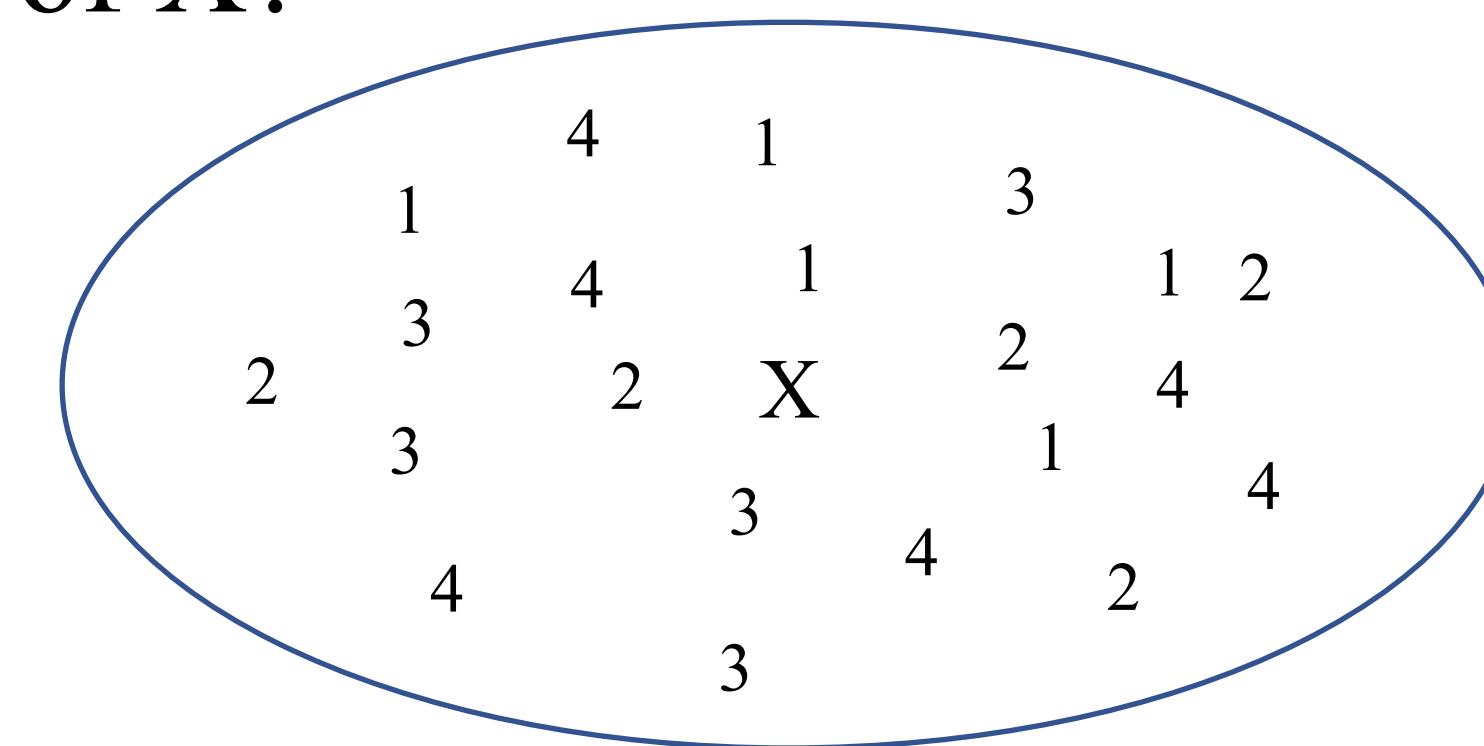
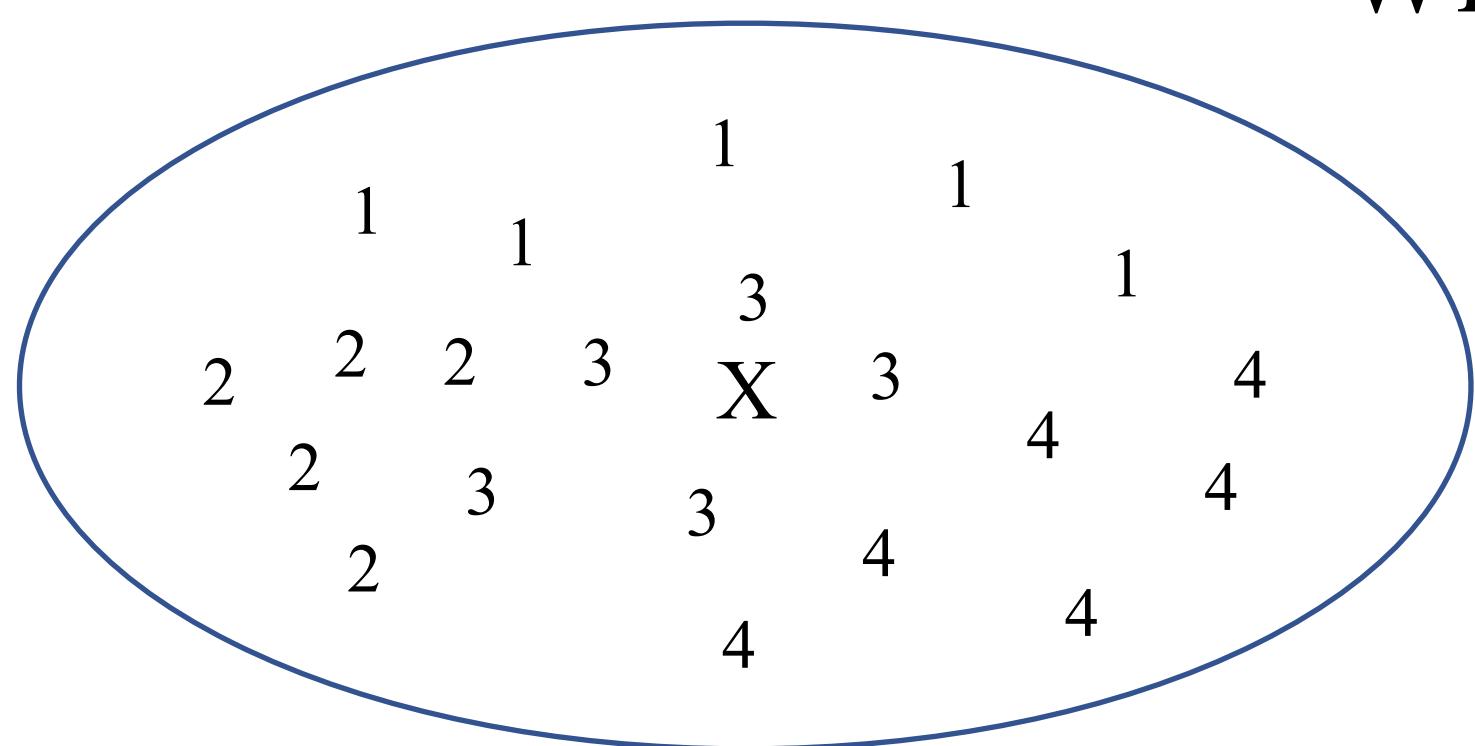
Machine learning

The kind of question we will ask:

I know the light absorption efficiency of 500 materials.

Now I discover a new material. Can I say anything about the light absorption of this material based on my previous knowledge?

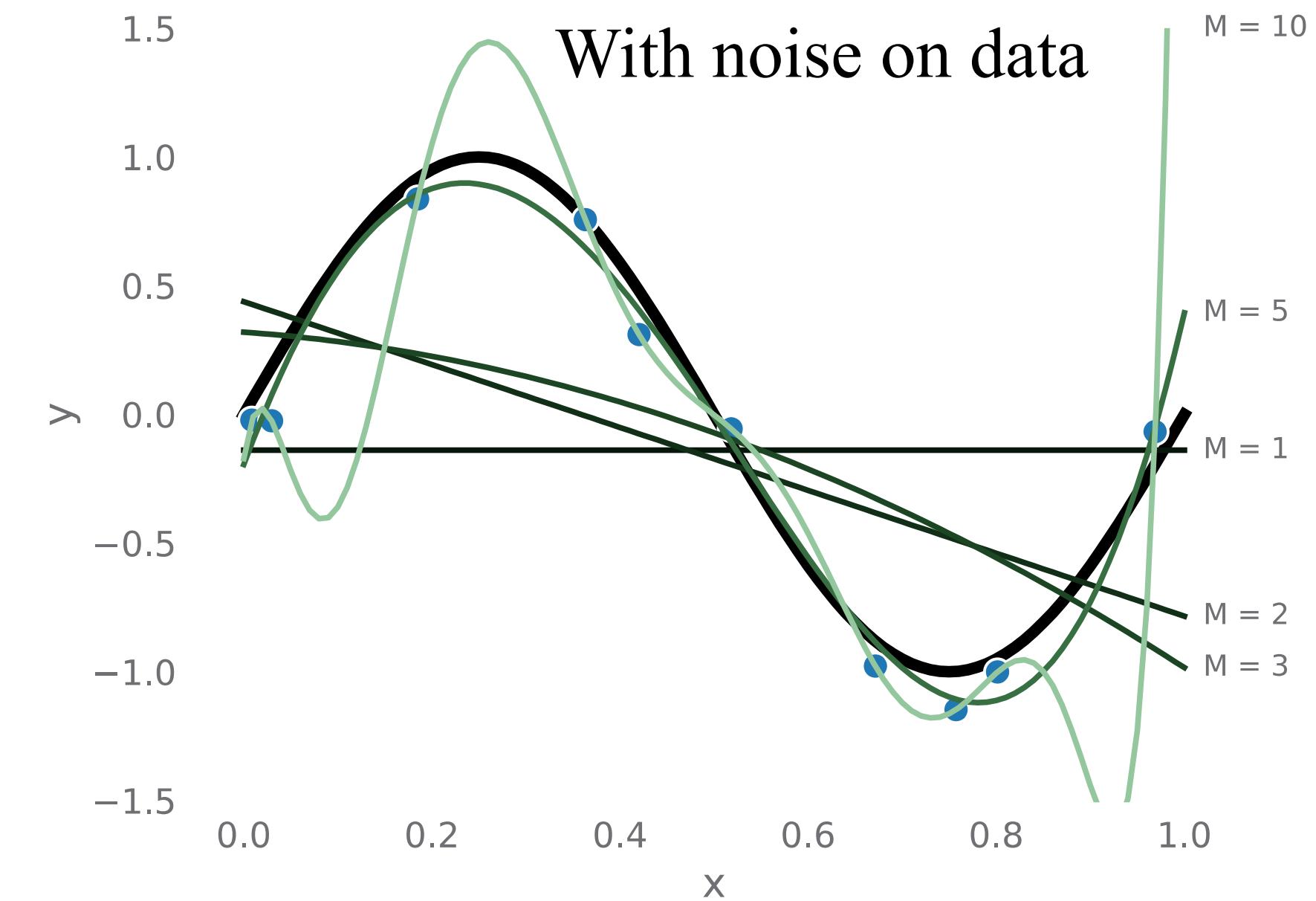
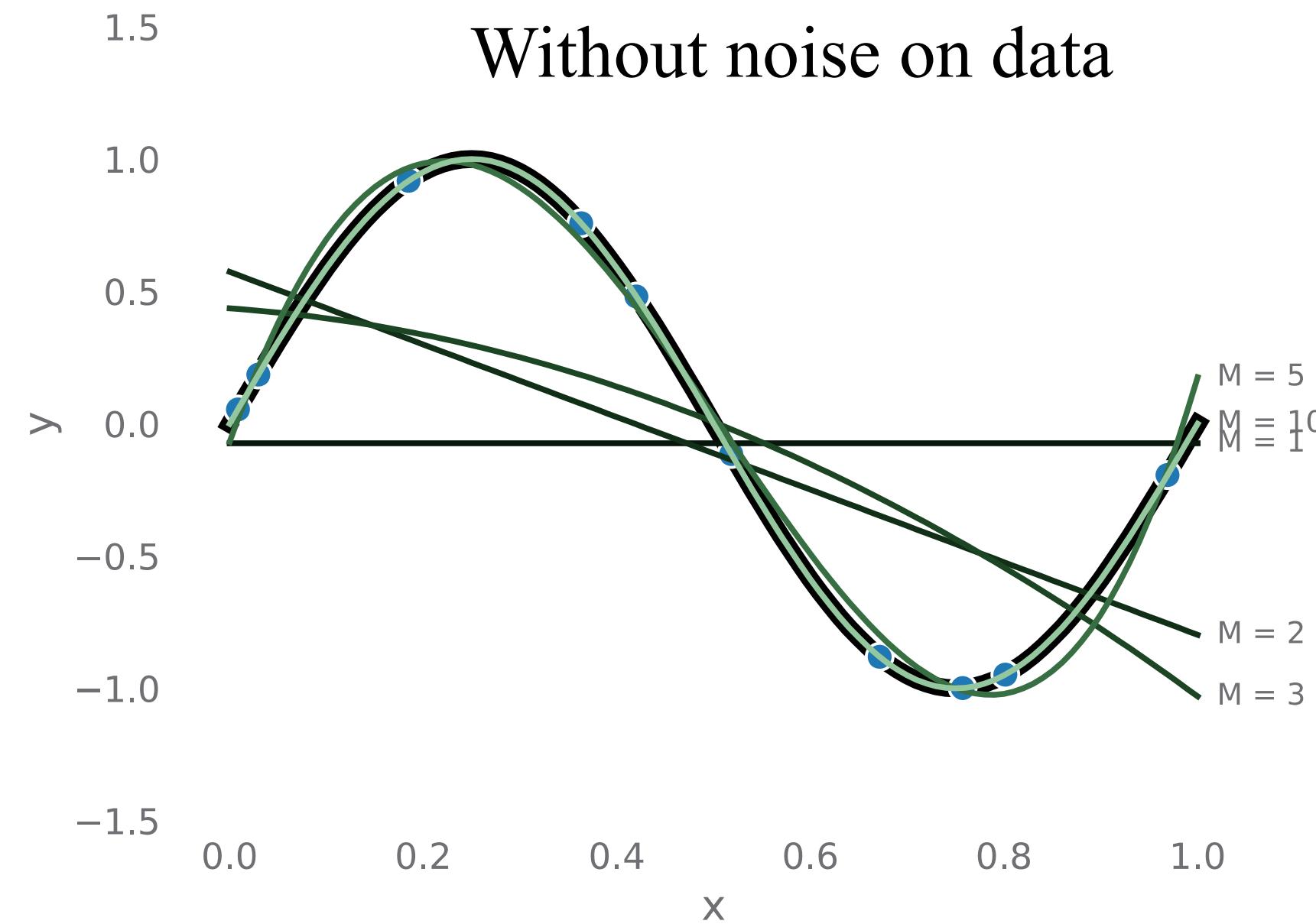
What is the value of X?



A key concept: correlation

The primary example: fitting a function

Polynomial fit: $f_{\text{fit}}(x) = w_0 + w_1x + w_2x^2 + \cdots + w_{M-1}x^{M-1}$

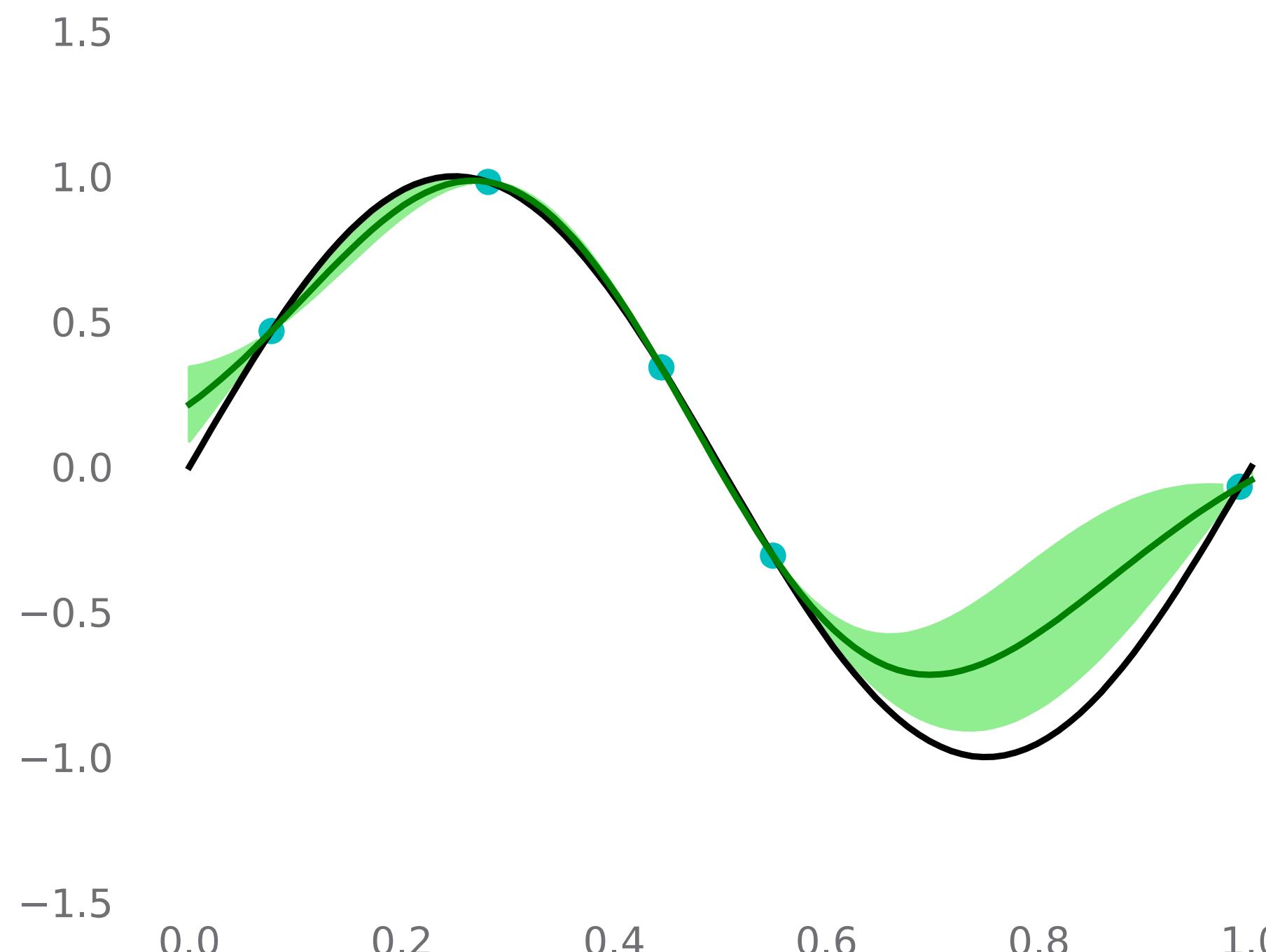


Too simple function (M too small): high *bias*

Too flexible function (M too high): high *variance (overfitting)*

What is the optimal value for M?

The goal: Reliable predictions with uncertainty estimates



Prediction with *Gaussian Process*

Black curve: target function
Green curve: best prediction
Green area: uncertainty

To deal with uncertainties we need *probability theory*. Coming up next ...

Probability theory

We assign probabilities to *propositions*:

Examples:

Y: When you through a die, it will show the number 3.

Y is an *event*, which can be investigated statistically (throw the die many times)

$$\mathcal{P}(Y) = 1/6$$

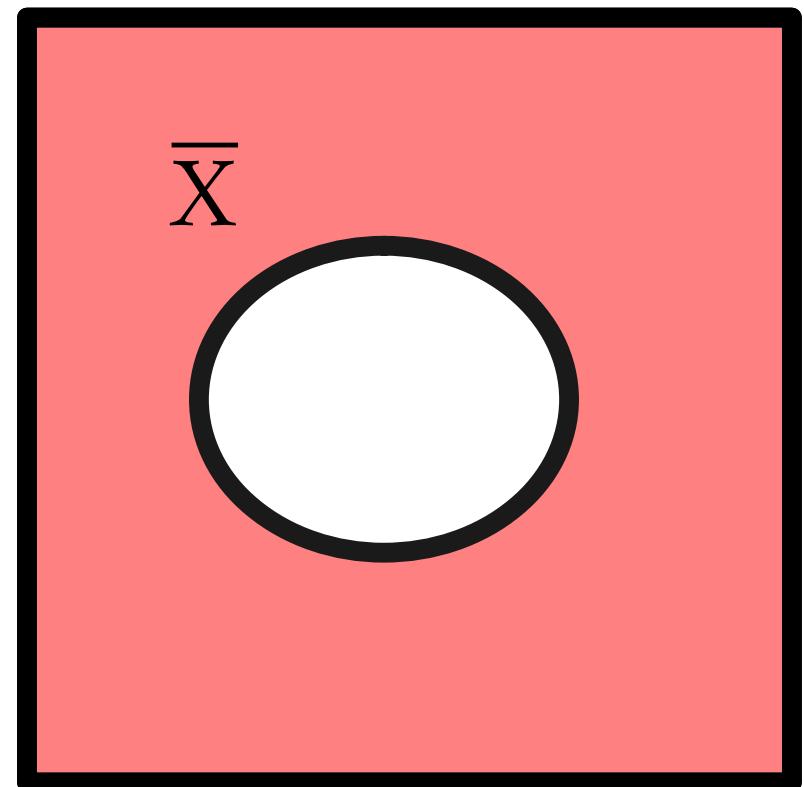
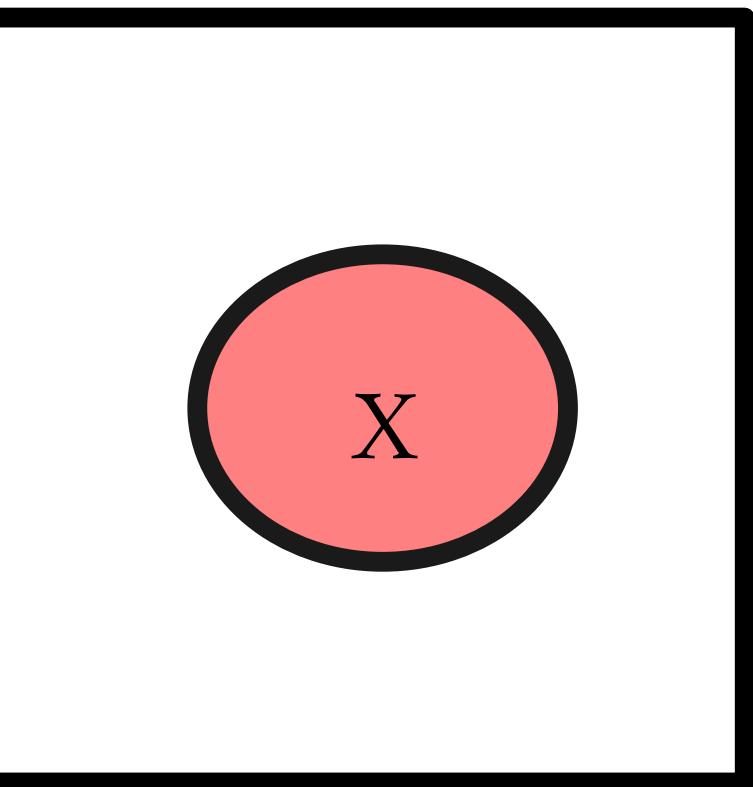
X: The best polynomial fit to the function is of order m=5

Not an event, but still a proposition.

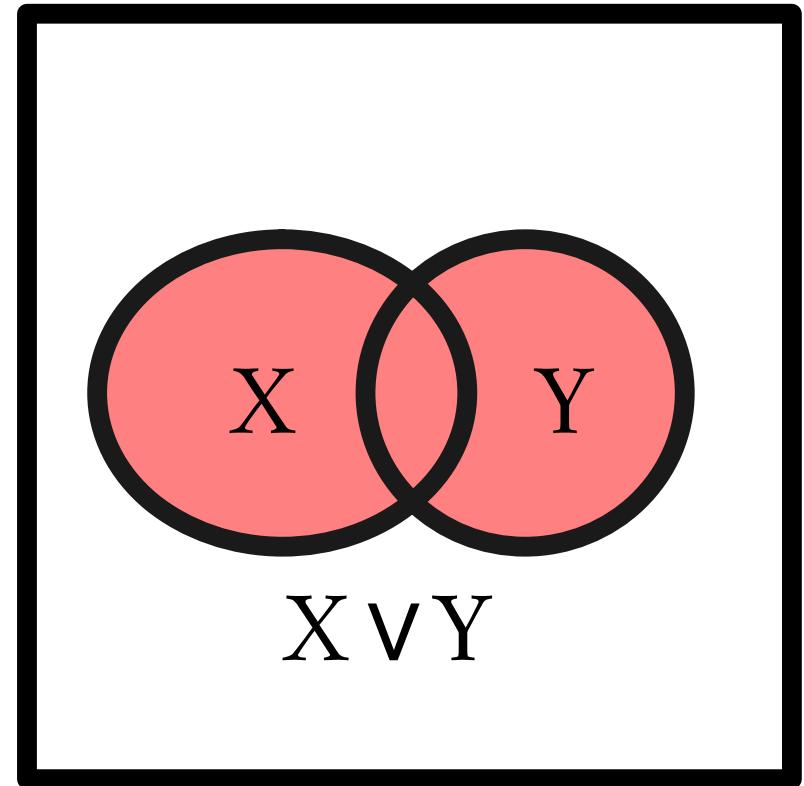
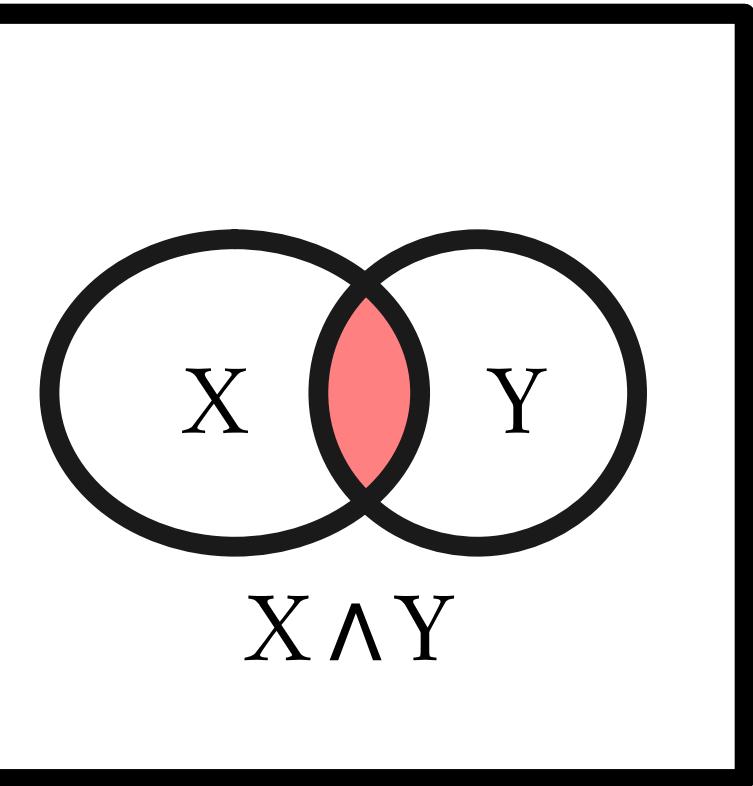
$$\mathcal{P}(X) = ?$$

Probability theory: the ingredients

Not $X : \overline{X}$



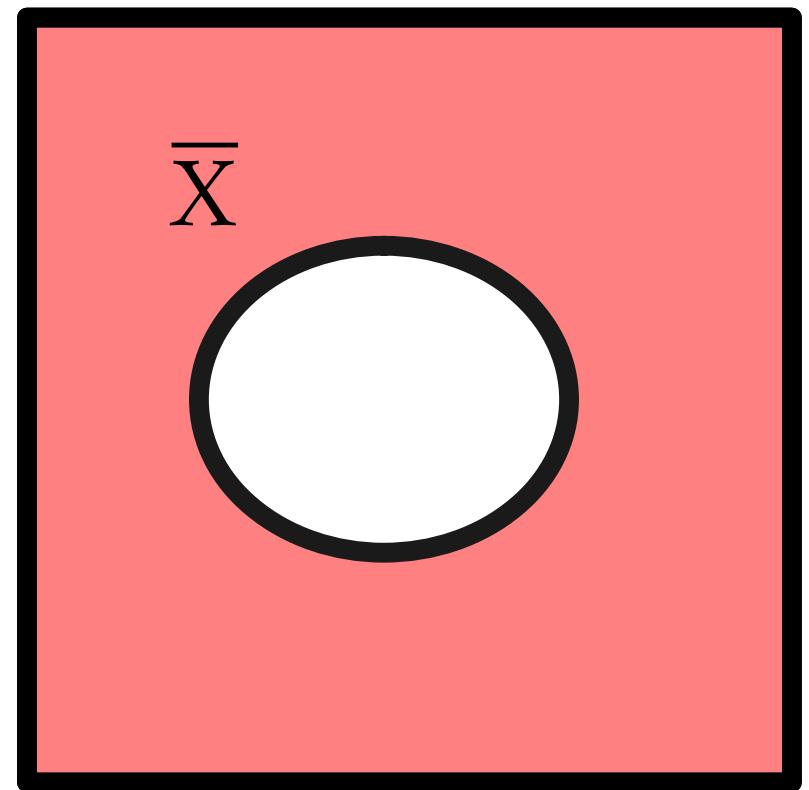
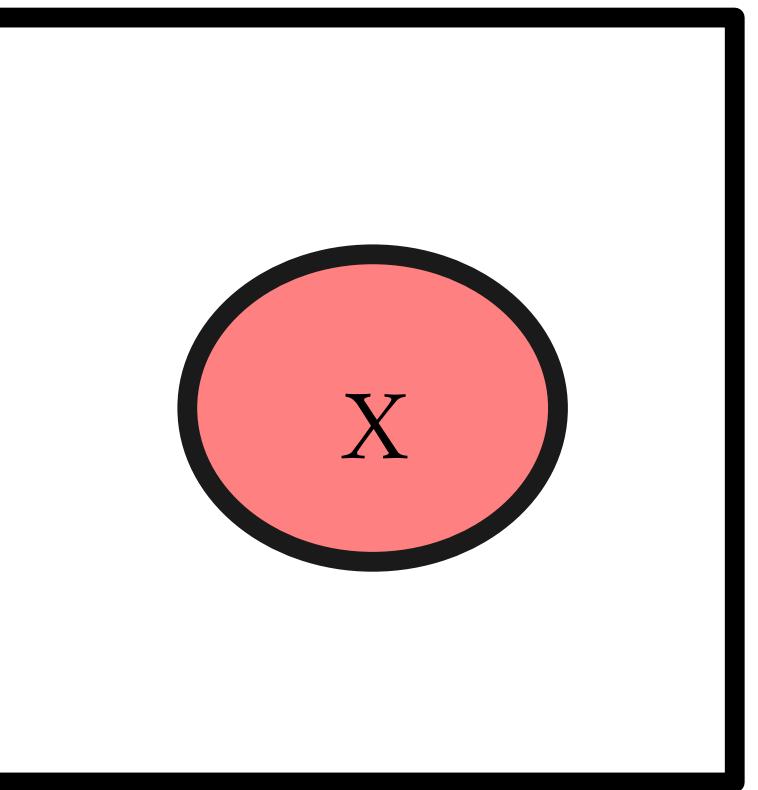
X and $Y : X \wedge Y$ (or X, Y or just XY)



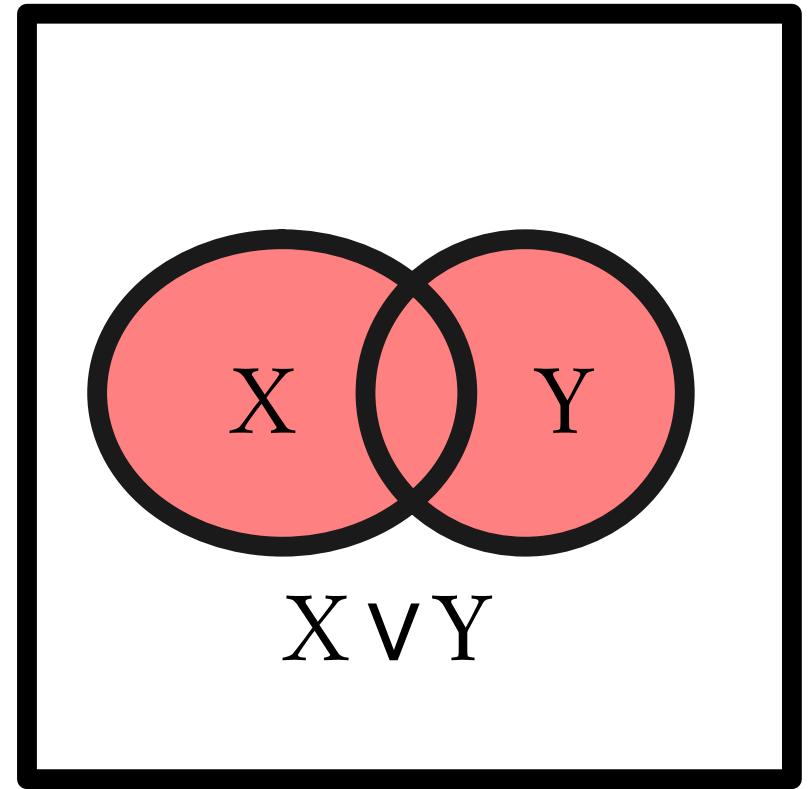
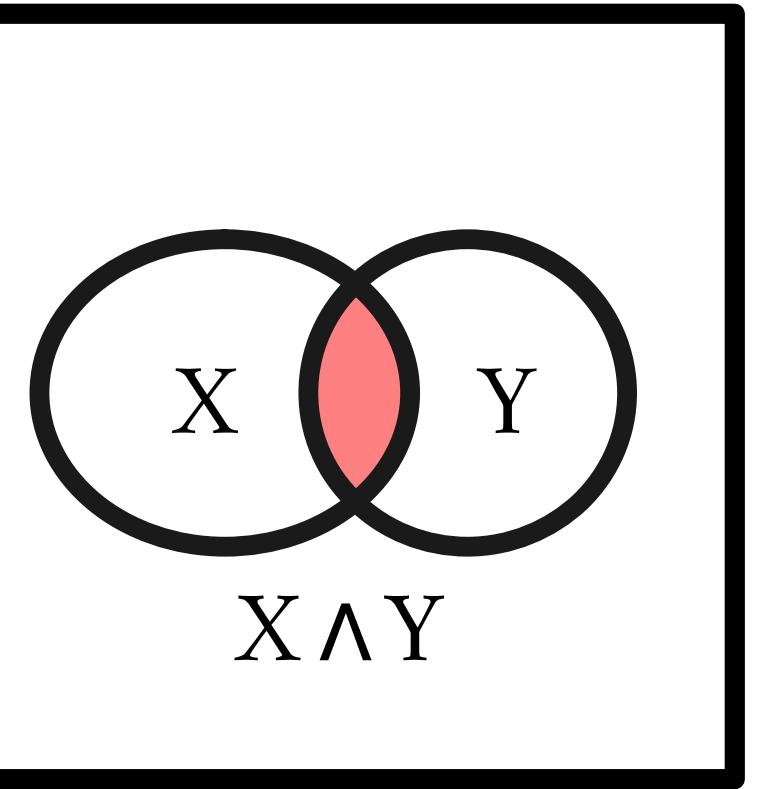
X	Y	\overline{X}	$X \wedge Y$	$X \vee Y$
True	True	False	True	True
True	False	False	False	True
False	True	True	False	True
False	False	True	False	False

Probability theory: the ingredients

Not $X : \overline{X}$



X and $Y : X \wedge Y$ (or X, Y or just XY)



Example: Results of throwing a die

$X =$ the die shows 5 or 6 = $\{5, 6\}$

$Y = \{1, 6\}$

$\overline{X} = \{1, 2, 3, 4\}$

$X \wedge Y = \{6\}$

$X \vee Y = \{1, 5, 6\}$

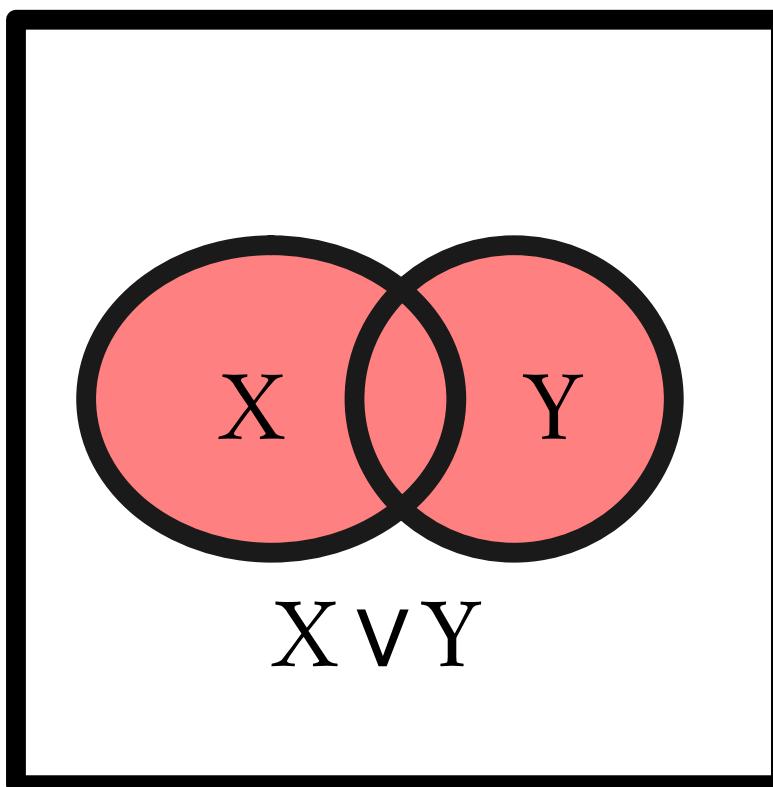
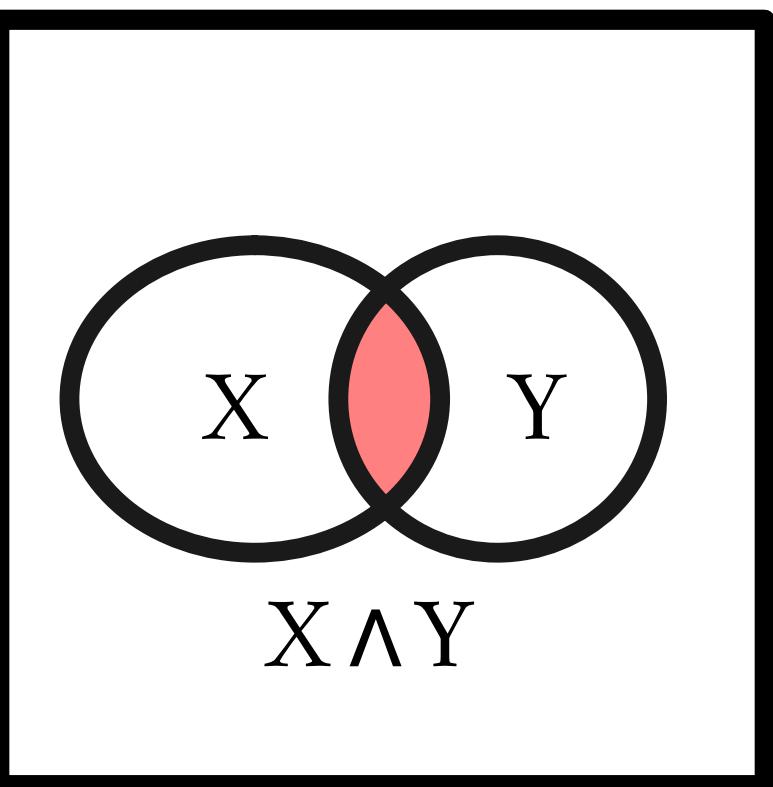
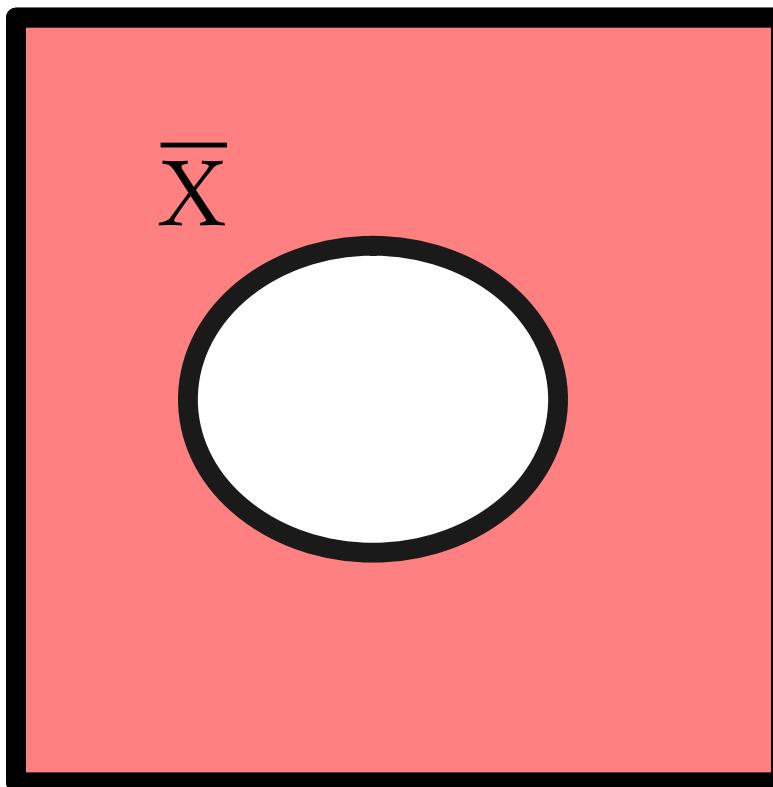
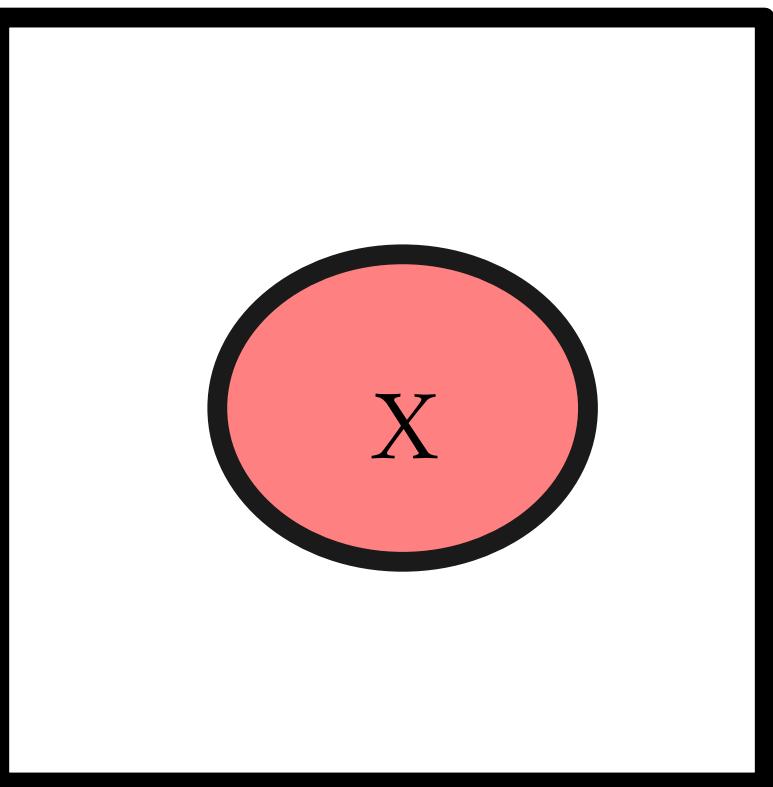
Probability theory: the ingredients

or using a truth table

Show graphically, that $X \vee Y = \overline{\overline{X} \wedge \overline{Y}}$

And that $\overline{X \wedge Y} = \overline{X} \vee \overline{Y}$

And $X \wedge (Y \vee Z) = (X \wedge Y) \vee (X \wedge Z)$



Form groups in the breakout rooms and discuss.
Contact me or Andreas via chat if you have
questions and we shall join the breakout room.

When you are done raise the blue hand in zoom!

Probability theory: the rules of the game

Positivity: $\mathcal{P}(X) \geq 0$

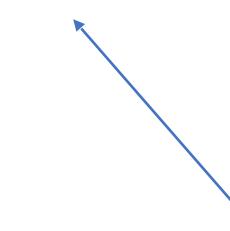
Sum rule: $\mathcal{P}(X) + \mathcal{P}(\overline{X}) = 1$ or $\mathcal{P}(X|I) + \mathcal{P}(\overline{X}|I) = 1$

Product rule: $\mathcal{P}(X, Y) = \mathcal{P}(X|Y)\mathcal{P}(Y)$ or $\mathcal{P}(X, Y|I) = \mathcal{P}(X|Y, I)\mathcal{P}(Y|I)$

Probability of X given I



Probability of X given Y



That is all!

The extended sum rule

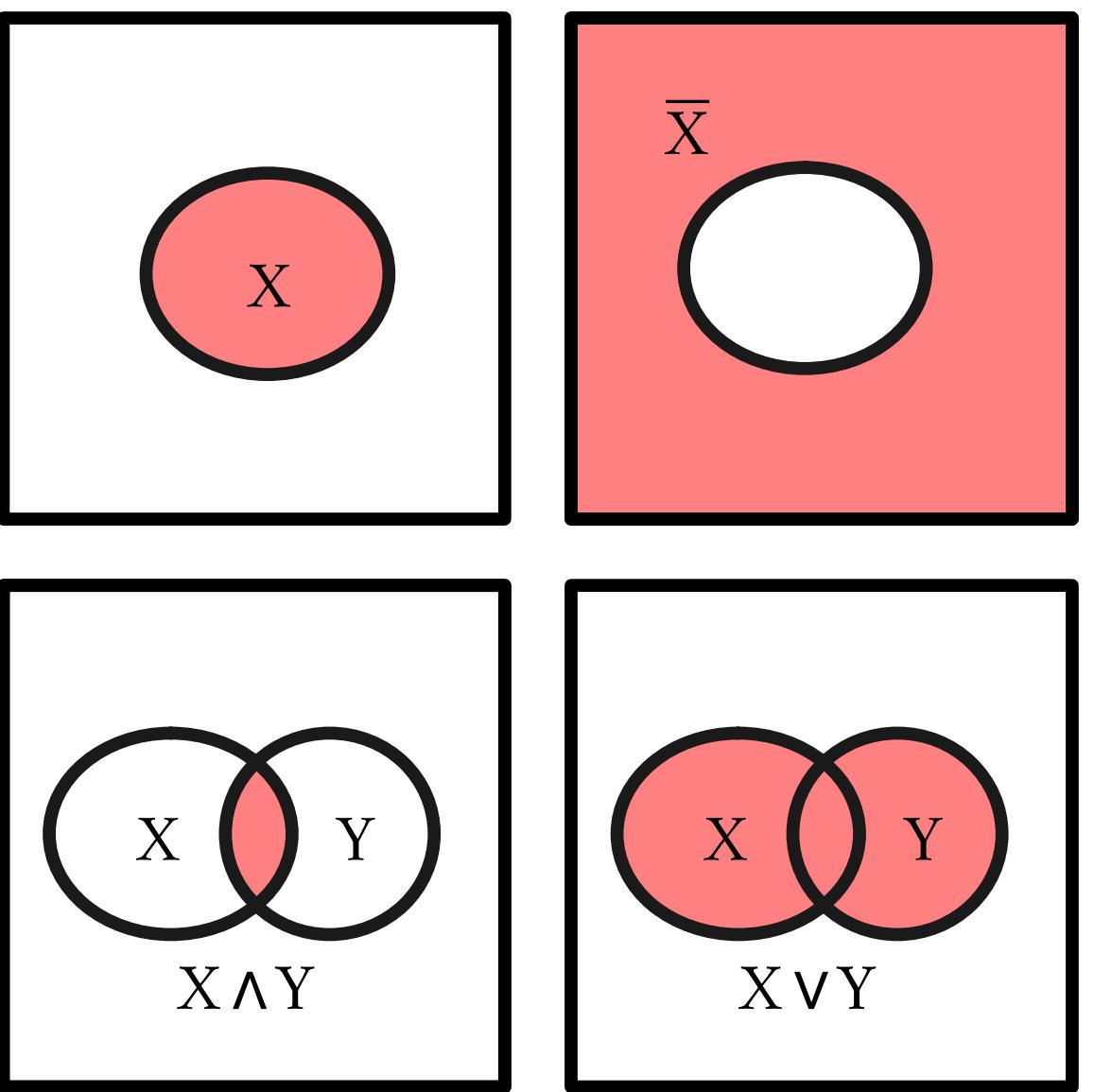
Sum rule: $\mathcal{P}(X) + \mathcal{P}(\bar{X}) = 1$

Product rule: $\mathcal{P}(X, Y) = \mathcal{P}(X|Y)\mathcal{P}(Y)$

$$\mathcal{P}(X, Y) = \mathcal{P}(X) + \mathcal{P}(Y) - \mathcal{P}(X \vee Y)$$

Proof:

$$\begin{aligned}\mathcal{P}(X \vee Y) &= \mathcal{P}(\bar{\bar{X}} \bar{\bar{Y}}) = 1 - \mathcal{P}(\bar{X} \bar{Y}) \\&= 1 - \mathcal{P}(\bar{X}|\bar{Y})\mathcal{P}(\bar{Y}) = 1 - (1 - \mathcal{P}(X|\bar{Y}))\mathcal{P}(\bar{Y}) \\&= 1 - \mathcal{P}(\bar{Y}) + \mathcal{P}(X\bar{Y}) = \mathcal{P}(Y) + \mathcal{P}(\bar{Y}|X)\mathcal{P}(X) \\&= \mathcal{P}(Y) + (1 - \mathcal{P}(Y|X))\mathcal{P}(X) \\&= \mathcal{P}(X) + \mathcal{P}(Y) - \mathcal{P}(XY)\end{aligned}$$



Sum rule: $\mathcal{P}(X) + \mathcal{P}(\overline{X}) = 1$

Marginalization

Let Y_i be a set of mutually exclusive propositions ($\mathcal{P}(Y_i, Y_j) = 0$ for $i \neq j$), which together cover all possible propositions.

Then $\sum_{i=1}^n \mathcal{P}(Y_i) = 1$ or, if Y is continuous $\int \mathcal{P}(Y) dY = 1$

Using the product rule, we get

$$\mathcal{P}(X) = \sum_{i=1}^n \mathcal{P}(Y_i|X)\mathcal{P}(X) = \sum_{i=1}^n \mathcal{P}(Y_i, X) = \sum_{i=1}^n \mathcal{P}(X|Y_i)\mathcal{P}(Y_i) = \left(\int \mathcal{P}(X|Y)\mathcal{P}(Y) dY \right)$$

The process of “integrating out” Y is called *marginalization*.

Independence

A proposition X is said to be *independent* of another proposition Y if

$$\mathcal{P}(X|Y) = \mathcal{P}(X)$$

i.e. the probability of X is unchanged by Y being true or not.

For independent (also called *uncorrelated*) propositions we have

$$\mathcal{P}(X, Y) = \mathcal{P}(X|Y)\mathcal{P}(Y) = \mathcal{P}(X)\mathcal{P}(Y)$$

Exercises

- Use the sum and product rules to solve the exercises in the notes about
 - The product rule
 - The extended sum rule (I just did that)
 - The extended, extended sum rule
 - Independence of propositions
- Come back to main room at 11:45

Bayes' theorem

Reshuffle the equation

$$\mathcal{P}(X, Y) = \mathcal{P}(X|Y)\mathcal{P}(Y) = \mathcal{P}(Y|X)\mathcal{P}(X)$$

and we get *Bayes' theorem*

$$\mathcal{P}(Y|X) = \frac{\mathcal{P}(X|Y)\mathcal{P}(Y)}{\mathcal{P}(X)} \quad \left(\text{or } \mathcal{P}(Y|X, I) = \frac{\mathcal{P}(X|Y, I)\mathcal{P}(Y|I)}{\mathcal{P}(X|I)} \right)$$

This is the most important formula in this course!

Bayes' theorem and rational reasoning

From the book by Jaynes:

Suppose some dark night a policeman walks down a street, apparently deserted. Suddenly he hears a burglar alarm, looks across the street, and sees a jewelry store with a broken window. Then a gentleman wearing a mask comes crawling out through the broken window, carrying a bag which turns out to be full of expensive jewelry. The policeman doesn't hesitate at all in deciding that this gentleman is dishonest. But by what reasoning process does he arrive at this conclusion?

$$\mathcal{P}(\text{thief}|\text{behavior}) = \frac{\mathcal{P}(\text{behavior}|\text{thief})\mathcal{P}(\text{thief})}{\mathcal{P}(\text{behavior})}$$

Example: Alice, Bob, and their two kids

Alice and Bob have two children. What is the probability that they have two girls when you get the following information (four different cases)?

They have at least one girl (“ $\geq 1g$ ”)

The oldest child is a girl (“old”)

They have a girl with blue eyes (“blue”)

They have a girl with the name “Mushroom” (“Mushroom”)

They have at least one girl

$$\mathcal{P}(2g| \geq 1g) = \frac{\mathcal{P}(\geq 1g|2g)\mathcal{P}(2g)}{\mathcal{P}(\geq 1g)} = \frac{\mathcal{P}(2g)}{\mathcal{P}(\geq 1g)} = \frac{1/4}{3/4} = \frac{1}{3}$$

The oldest child is a girl

$$\mathcal{P}(2g|old) = \frac{\mathcal{P}(old|2g)\mathcal{P}(2g)}{\mathcal{P}(old)} = \frac{\mathcal{P}(2g)}{\mathcal{P}(old)} = \frac{1/4}{1/2} = \frac{1}{2}$$

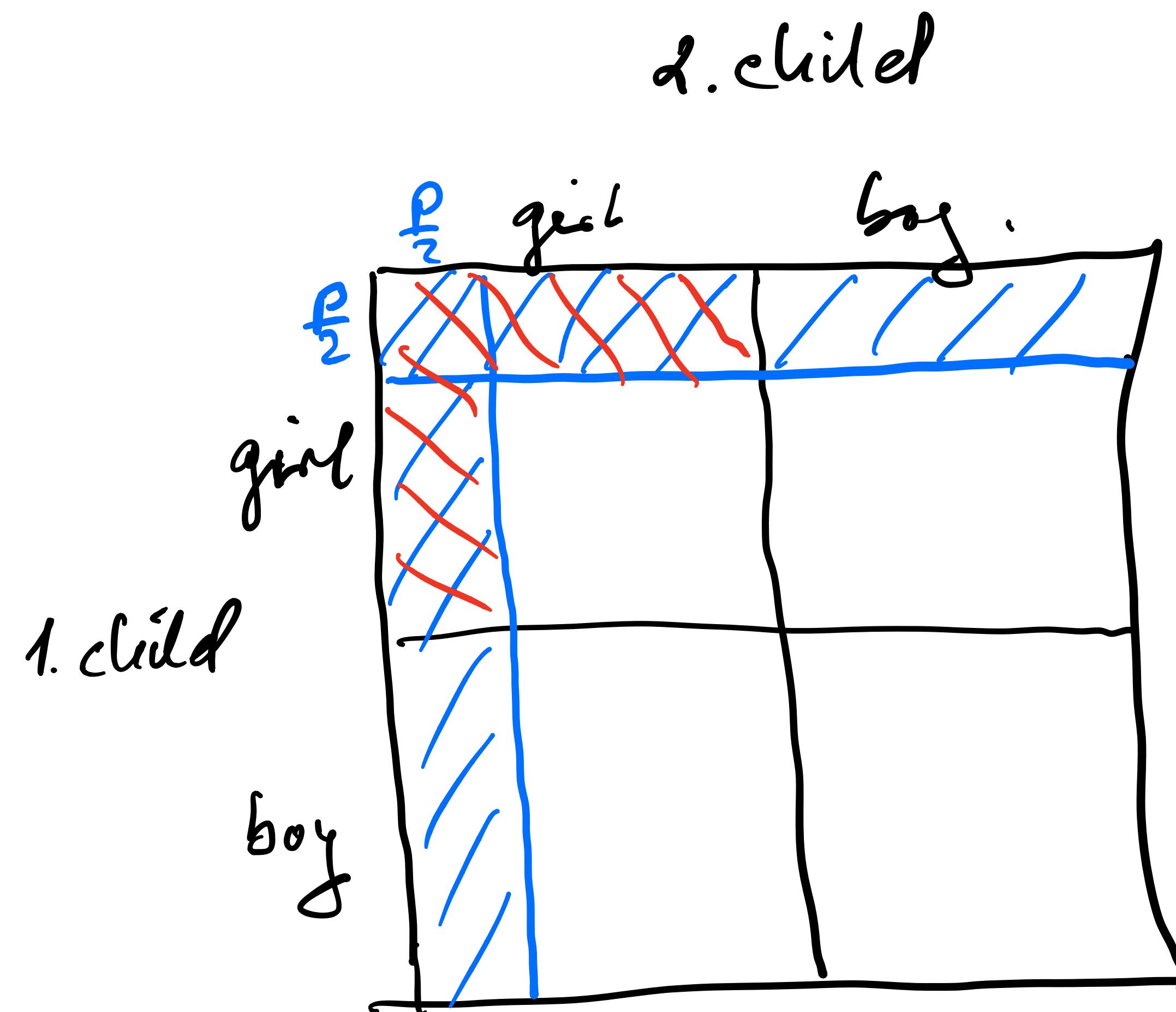
They have a girl with blue eyes
(probability of blue eyes: p)

$$\begin{aligned}\mathcal{P}(2g|blue) &= \frac{\mathcal{P}(blue|2g)\mathcal{P}(2g)}{\mathcal{P}(blue)} = \frac{\mathcal{P}(blue|2g)\mathcal{P}(2g)}{\mathcal{P}(blue|2g)\mathcal{P}(2g) + \mathcal{P}(blue|1g)\mathcal{P}(1g)} \\ &= \frac{(p^2 + 2p(1-p)) \cdot 1/4}{(p^2 + 2p(1-p)) \cdot 1/4 + p \cdot 1/2} = \frac{2-p}{4-p}\end{aligned}$$

So if $p \sim 1/2$, we get $\mathcal{P}(2g|blue) \sim \frac{3}{7}$

They have a girl with the name “Mushroom”
(Same analysis as for the blue eyes, but now $p \ll 1$)

$$\mathcal{P}(2g|Mushroom) \sim \frac{1}{2}$$



$$P(\text{arg } g \text{ with property } \rho) = \frac{\text{red area}}{\text{blue area}} = \frac{2 \cdot F_2 \cdot \frac{1}{2} - (\frac{F_2}{2})^2}{2 \cdot F_2 - (\frac{F_2}{2})^2} = \frac{2 - \rho}{4 - \rho}$$

Exercise

- Now do the exercise on medical screening
- Let us meet in 15 minutes

Medical screening

$$\begin{aligned}\mathcal{P}(ill|pos) &= \frac{\mathcal{P}(pos|ill)\mathcal{P}(ill)}{\mathcal{P}(pos)} = \frac{\mathcal{P}(pos|ill)\mathcal{P}(ill)}{\mathcal{P}(pos|ill)\mathcal{P}(ill) + \mathcal{P}(pos|h)\mathcal{P}(h)} \\ &= \frac{0.99 \cdot 0.001}{0.99 \cdot 0.001 + 0.01 \cdot 0.999} \sim 0.09\end{aligned}$$

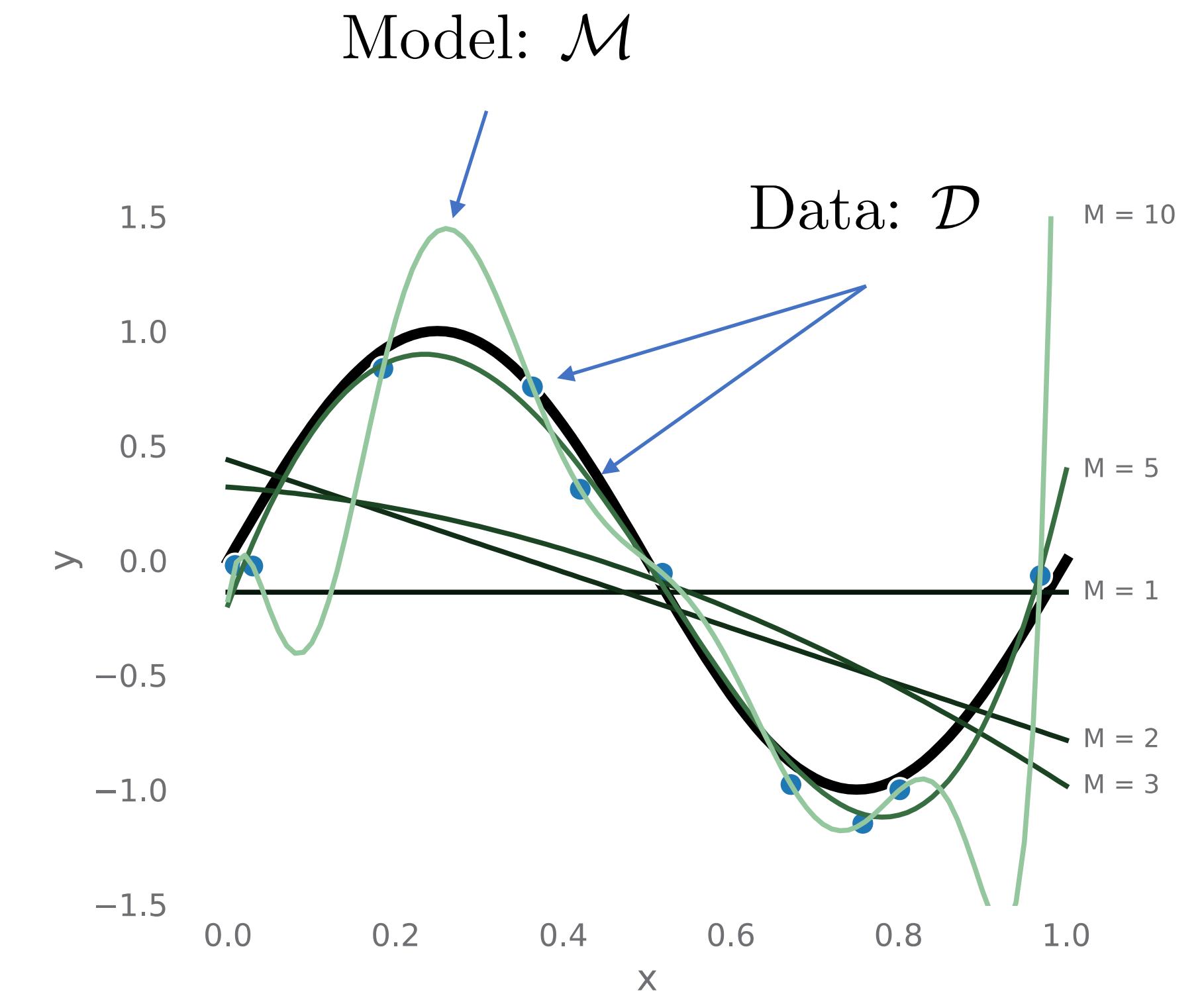
Bayes' theorem: models and data

$$\mathcal{P}(\mathcal{M}|\mathcal{D}) = \frac{1}{\mathcal{P}(\mathcal{D})} \mathcal{P}(\mathcal{D}|\mathcal{M}) \mathcal{P}(\mathcal{M})$$

↑
Posterior probability distribution
Likelihood
↑
Prior probability distribution

Normalization:

$$\mathcal{P}(\mathcal{D}) = \mathcal{P}(\mathcal{D}, \mathcal{M}) + \mathcal{P}(\mathcal{D}, \bar{\mathcal{M}}) = \mathcal{P}(\mathcal{D}|\mathcal{M}) \mathcal{P}(\mathcal{M}) + \mathcal{P}(\mathcal{D}|\bar{\mathcal{M}}) \mathcal{P}(\bar{\mathcal{M}})$$



Why do we believe Newton's second law?

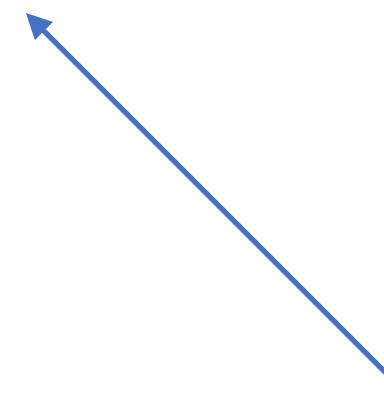
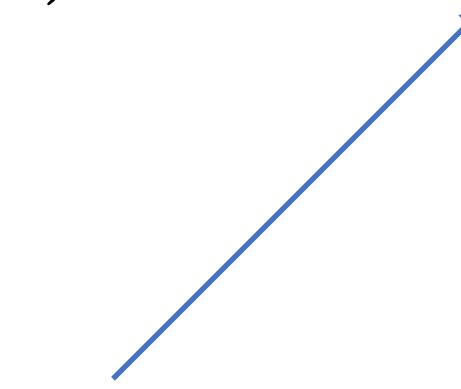
Bayes' theorem - updating with new data

$$\mathcal{P}(\mathcal{M}|\mathcal{D}_{new}, \mathcal{D}_{old}) = \frac{1}{\mathcal{P}(\mathcal{D}_{new}|\mathcal{D}_{old})} \mathcal{P}(\mathcal{D}_{new}|\mathcal{M}, \mathcal{D}_{old}) \mathcal{P}(\mathcal{M}|\mathcal{D}_{old})$$

New posterior distribution

Likelihood of new data

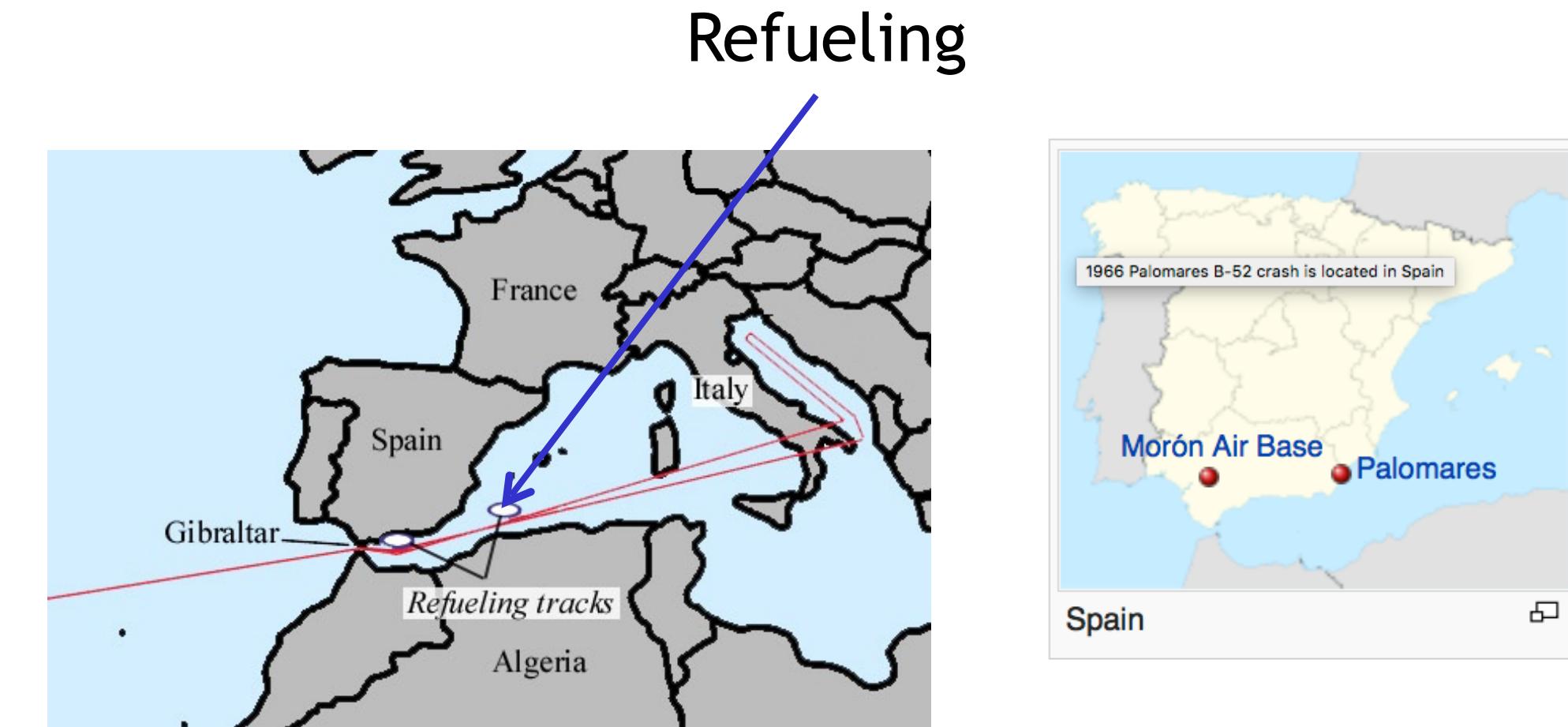
Prior based on old data



Bayesian Search Theory in Practice

The 1966 Palomares B-52 crash

B-52G collided with KC-135 tanker when fueling



4 H-bombs dropped
3 on land
1 in the Mediterranean Sea

Bayesian search theory applied:
Assign probabilities to different areas of the sea
based on available information
(a local fisherman saw the bomb dropping)
Update your probability depending on your search.

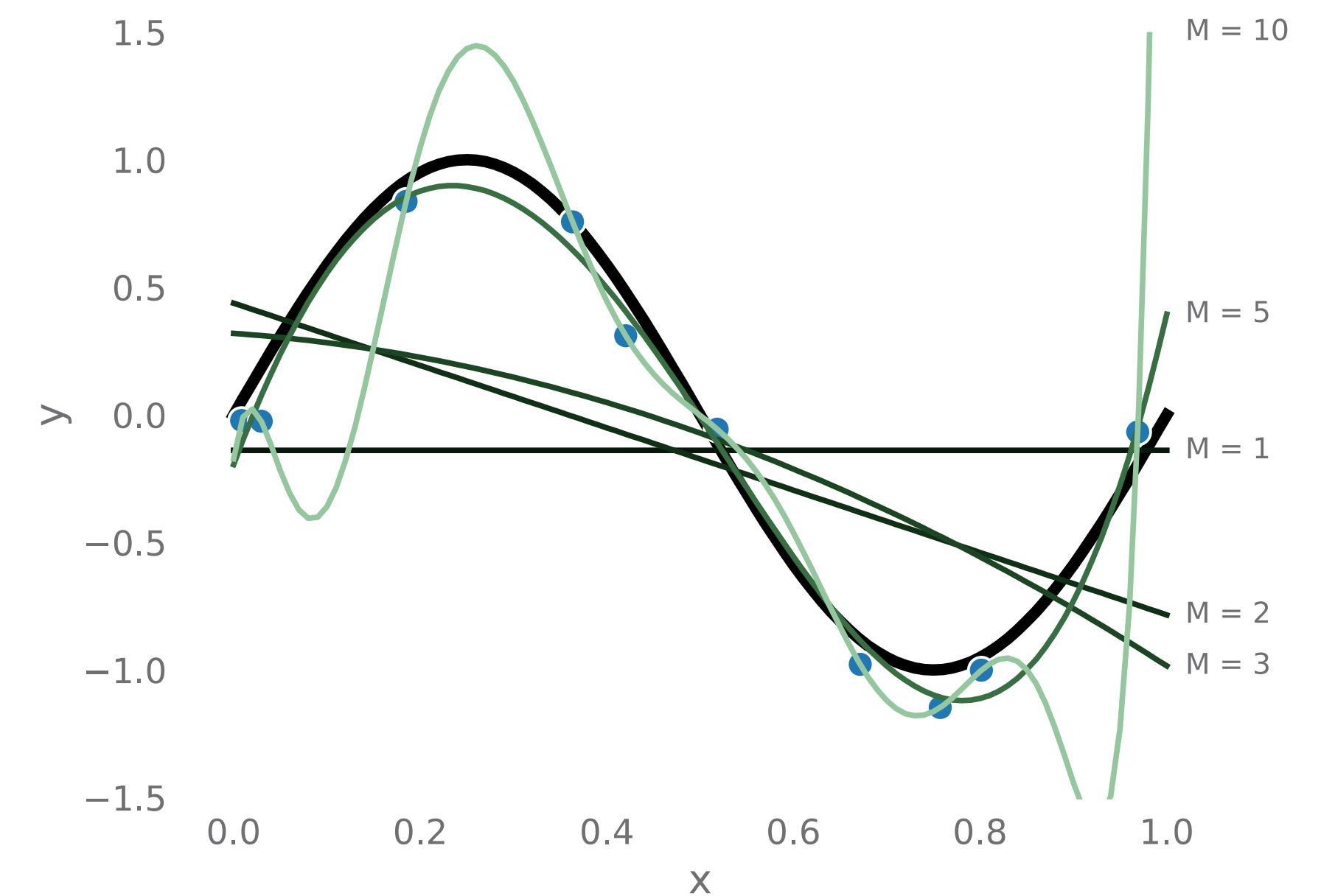


Bomb recovered

Parameter estimation

$$f_{\text{fit}}(x) = w_0 + w_1 x + w_2 x^2 + \cdots + w_{M-1} x^{M-1}$$

$$\begin{aligned}\mathcal{P}(\mathbf{w}|\mathcal{D}, \mathcal{M}) &= \frac{\mathcal{P}(\mathcal{D}|\mathbf{w}, \mathcal{M})\mathcal{P}(\mathbf{w}, \mathcal{M})}{\mathcal{P}(\mathcal{D}, \mathcal{M})} = \frac{\mathcal{P}(\mathcal{D}|\mathbf{w}, \mathcal{M})\mathcal{P}(\mathbf{w}, \mathcal{M})}{\int \mathcal{P}(\mathcal{D}, \mathbf{w}', \mathcal{M})d\mathbf{w}'} \\ &= \frac{\mathcal{P}(\mathcal{D}|\mathbf{w}, \mathcal{M})\mathcal{P}(\mathbf{w}, \mathcal{M})}{\int \mathcal{P}(\mathcal{D}|\mathbf{w}', \mathcal{M})\mathcal{P}(\mathbf{w}', \mathcal{M})d\mathbf{w}'} \\ &= \frac{\mathcal{P}(\mathcal{D}|\mathbf{w}, \mathcal{M})\mathcal{P}(\mathbf{w}|\mathcal{M})\mathcal{P}(\mathcal{M})}{\int \mathcal{P}(\mathcal{D}|\mathbf{w}', \mathcal{M})\mathcal{P}(\mathbf{w}'|\mathcal{M})d\mathbf{w}'\mathcal{P}(\mathcal{M})} \\ &= \frac{\mathcal{P}(\mathcal{D}|\mathbf{w}, \mathcal{M})\mathcal{P}(\mathbf{w}|\mathcal{M})}{\int \mathcal{P}(\mathcal{D}|\mathbf{w}', \mathcal{M})\mathcal{P}(\mathbf{w}'|\mathcal{M})d\mathbf{w}'}\end{aligned}$$



The denominator just ensures $\int \mathcal{P}(\mathbf{w}|\mathcal{D}, \mathcal{M})d\mathbf{w} = 1$

Librarian or farmer? (Tversky and Kahneman)

An individual has been described by a neighbor as follows: “Steve is very shy and withdrawn, invariably helpful but with very little interest in people or in the world of reality. A meek and tidy soul, he has a need for order and structure, and a passion for detail.” Is Steve more likely to be a librarian or a farmer?

$$\frac{\mathcal{P}(\text{farmer}|\text{description})}{\mathcal{P}(\text{librarian}|\text{description})} = \frac{\mathcal{P}(\text{description}|\text{farmer})}{\mathcal{P}(\text{description}|\text{librarian})} \frac{\mathcal{P}(\text{farmer})}{\mathcal{P}(\text{librarian})}$$

In Denmark there are \sim 10000 farms, but only \sim 250 public libraries

Model selection

Comparing two models

$$\begin{aligned}\frac{\mathcal{P}(\mathcal{M}_2|\mathcal{D})}{\mathcal{P}(\mathcal{M}_1|\mathcal{D})} &= \frac{\mathcal{P}(\mathcal{D}|\mathcal{M}_2)}{\mathcal{P}(\mathcal{D}|\mathcal{M}_1)} \cdot \frac{\mathcal{P}(\mathcal{M}_2)}{\mathcal{P}(\mathcal{M}_1)} \\ &= \frac{\int \mathcal{P}(\mathcal{D}|\mathcal{M}_2, \mathbf{w}_2) \mathcal{P}(\mathbf{w}_2|\mathcal{M}_2) d\mathbf{w}_2}{\int \mathcal{P}(\mathcal{D}|\mathcal{M}_1, \mathbf{w}_1) \mathcal{P}(\mathbf{w}_1|\mathcal{M}_1) d\mathbf{w}_1} \cdot \frac{\mathcal{P}(\mathcal{M}_2)}{\mathcal{P}(\mathcal{M}_1)}\end{aligned}$$

Probability distributions - definitions

Probability distribution $\mathcal{P}(x)$

Normalization: $\int \mathcal{P}(x) dx = 1$

Average value of a function:

$$\langle f \rangle = \int f(x) \mathcal{P}(x) dx,$$

Variance:

$$\text{Var}[f] = \langle (f - \langle f \rangle)^2 \rangle = \langle f^2 \rangle - \langle f \rangle^2$$

In particular:

$$\text{Var}[x] = \langle x^2 \rangle - \langle x \rangle^2$$

Several variables: $x = \begin{pmatrix} x_1 \\ x_2 \\ \vdots \\ x_D \end{pmatrix}$

$$xx^T = \begin{pmatrix} x_1 x_1 & x_1 x_2 & \cdots & x_1 x_D \\ x_2 x_1 & x_2 x_2 & \cdots & x_2 x_D \\ \vdots & \vdots & & \vdots \\ x_D x_1 & x_D x_2 & \cdots & x_D x_D \end{pmatrix}$$

Covariance matrix:

$$\text{Cov}[x, x^T] = \langle (x - \langle x \rangle)(x^T - \langle x^T \rangle) \rangle = \langle xx^T \rangle - \langle x \rangle \langle x \rangle^T$$

If x_i and x_j are independent:

$$\text{Cov}[x_i, x_j] = 0$$

The normal (Gaussian) distribution

$$\mathcal{N}(x|\mu, \sigma^2) = \frac{1}{(2\pi\sigma^2)^{(1/2)}} e^{-(x-\mu)^2/2\sigma^2} \quad \begin{aligned} \langle x \rangle &= \mu \\ \text{Var}[x] &= \sigma^2 \end{aligned}$$

Many nice mathematical properties. A product of two Gaussians is for example a new Gaussian:

$$\mathcal{N}(x|\mu, \sigma^2) \propto \mathcal{N}(x|\mu_1, \sigma_1^2) \cdot \mathcal{N}(x|\mu_2, \sigma_2^2) \quad \text{with} \quad \mu = \left(\frac{\mu_1}{\sigma_1^2} + \frac{\mu_2}{\sigma_2^2} \right) \sigma^2 \quad \text{and} \quad \frac{1}{\sigma^2} = \frac{1}{\sigma_1^2} + \frac{1}{\sigma_2^2}$$

The central limit theorem (vaguely formulated): If a stochastic variable y is a sum of many stochastic variables $y = y_1 + y_2 + \dots + y_N$ then as $N \rightarrow \infty$ the distribution $\mathcal{P}(y)$ becomes a normal distribution.

Binomial distribution

Consider a coin with probability p to get heads. The probability of getting n heads, when you throw the coin N times is

$$\mathcal{P}(n|N,p) = \binom{N}{n} p^n (1-p)^{(N-n)}$$
$$\langle n \rangle = pN$$
$$\text{Var}[n] = Np(1-p)$$

For large N we get the normal distribution (to be used in an exercise):

$$\mathcal{P}(n|N,p) \sim \frac{1}{N} \mathcal{N}(x = n/N | p, p(1-p)/N)$$

Exercises about Bayes' theorem

- Now do the exercises on Bayes' theorem:
 - Gaussianly distributed data
 - Tversky and Kahneman
 - Playing with priors - the average height of Norwegians
 - Is the coin fair?
 - (But is the coin really fair?)

Exercise: Gaussianly distributed data

Data x_1, x_2, \dots, x_N which we assume to be Gaussianly distributed with known variance σ^2 . What is the best estimate of the mean of the data?

$$\begin{aligned}\mathcal{P}(\mu|x_1, x_2, \dots, x_N, \sigma^2) &\propto \mathcal{P}(x_1, x_2, \dots, x_N|\mu, \sigma^2)\mathcal{P}(\mu|\sigma^2) \\ &\propto \mathcal{P}(x_1, x_2, \dots, x_N|\mu, \sigma^2) \\ &\propto \prod_{i=1}^N \mathcal{N}(x_i|\mu, \sigma^2) \propto \prod_{i=1}^N e^{-(x_i-\mu)^2/2\sigma^2} \\ &\propto \exp\left(-\sum_{i=1}^N (x_i-\mu)^2/2\sigma^2\right) \\ &\propto \exp\left(-((\mu-\bar{x})^2 + \overline{\Delta x^2})/2\sigma_N^2\right) \\ &\propto \exp\left(-(\mu-\bar{x})^2/2\sigma_N^2\right)\end{aligned}$$

$$\begin{aligned}\bar{x} &= \frac{1}{N} \sum_i x_i \\ \overline{\Delta x^2} &= \frac{1}{N} \sum_i (x_i - \bar{x})^2 \\ \sigma_N &= \sigma/\sqrt{N}\end{aligned}$$

The highest probability is obtained with
 $\mu = \bar{x}$
and the standard deviation (i.e. the square root of the variance) is

$$\sigma/\sqrt{N}$$

Average height of Norwegians

$$\mathcal{P}(h_{\text{Nor}} | h_1, h_2 \dots h_N) \propto \mathcal{P}(h_1, h_2 \dots h_N | h_{\text{Nor}}) \mathcal{P}(h_{\text{Nor}})$$

$$\propto \left(\prod_n \mathcal{N}(h_n | h_{\text{Nor}}, \sigma_{\text{Nor}}^2) \right) \mathcal{N}(h_{\text{Nor}} | h_{\text{DK}}, \sigma_d^2)$$

$$\propto \mathcal{N}(h_{\text{Nor}} | \bar{h}, \sigma_{\text{Nor}}^2/N) \mathcal{N}(h_{\text{Nor}} | h_{\text{DK}}, \sigma_d^2)$$

$$\propto \mathcal{N}(h_{\text{Nor}} | h_0, \sigma_0^2)$$

$$\bar{h} = \frac{1}{N} \sum_{i=1}^N h_i$$

$$\begin{aligned} h_0 &= \frac{\frac{\bar{h}}{\sigma_{\text{Nor}}^2/N} + \frac{h_{\text{DK}}}{\sigma_d^2}}{N/\sigma_{\text{Nor}}^2 + 1/\sigma_d^2} \\ &= \frac{\sum_i h_i + (\frac{\sigma_{\text{Nor}}}{\sigma_d})^2 h_{\text{DK}}}{N + (\frac{\sigma_{\text{Nor}}}{\sigma_d})^2} \end{aligned}$$

$$\frac{\sigma_{\text{Nor}}^2}{\sigma_d^2} \sim \frac{(5 \text{ cm})^2}{(2.5 \text{ cm})^2} \sim 4$$

$$\sigma_0^2 = \frac{\sigma_{\text{Nor}}^2}{N + (\frac{\sigma_{\text{Nor}}}{\sigma_d})^2}$$

“Pseudo-count” of 4 Danes.

Is the coin fair?

You get N_h heads out of N tosses. Is the coin fair?

$$\mathcal{P}(p|N_h, N) \propto \mathcal{P}(N_h|p, N)\mathcal{P}(p) \propto \mathcal{P}(N_h|p, N)$$

$$= \binom{N}{N_h} p^{N_h} (1-p)^{N_t} \propto \mathcal{B}(p|\alpha = N_h + 1, \beta = N_t + 1)$$

The best estimate (the average) of p :

$$\langle p \rangle = \frac{N_h + 1}{(N_h + 1) + (N_t + 1)} = \frac{N_h + 1}{N + 2}$$

Pseudo-count of 1 head and 1 tail

$$\sigma(p) = \sqrt{\text{Var}[p]} = \sqrt{\frac{(N_h + 1)(N_t + 1)}{(N + 2)^2(N + 3)}}$$

Examples:

6 heads out of 10: $\langle p \rangle \pm \sigma(p) = 0.58 \pm 0.14$ Could be fair.

600 heads out of 1000: $\langle p \rangle \pm \sigma(p) = 0.60 \pm 0.02$ Is probably not fair

Probability to get heads: p
Tails: $N_t = N - N_h$

Beta distribution:

$$\mathcal{B}(x|\alpha, \beta) = \frac{1}{B(\alpha, \beta)} x^{\alpha-1} (1-x)^{\beta-1}$$

$$\text{with } B(\alpha, \beta) = \frac{\Gamma(\alpha)\Gamma(\beta)}{\Gamma(\alpha + \beta)}$$

$$\langle x \rangle = \frac{\alpha}{\alpha + \beta}$$

$$\text{Var}[x] = \frac{\alpha\beta}{(\alpha + \beta)^2(\alpha + \beta + 1)}$$

6 heads out of 10

Is the coin fair?

You get N_h heads out of N tosses. Is the coin fair?

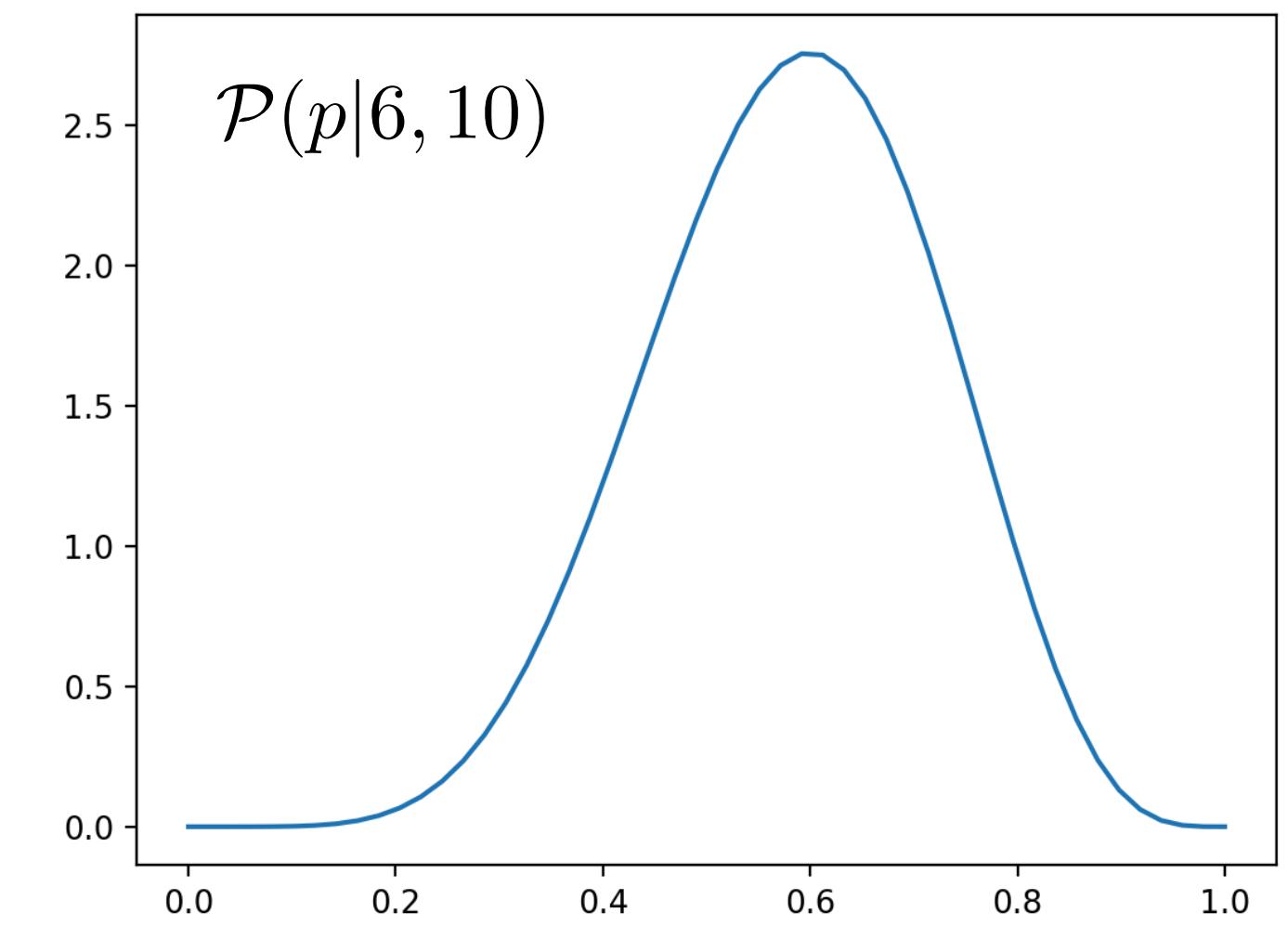
$$\begin{aligned}\mathcal{P}(p|N_h, N) &\propto \mathcal{P}(N_h|p, N)\mathcal{P}(p) \propto \mathcal{P}(N_h|p, N) \\ &= \binom{N}{N_h} p^{N_h} (1-p)^{N_t} \propto \mathcal{B}(p|\alpha = N_h + 1, \beta = N_t + 1)\end{aligned}$$

The best estimate (the average) of p :

$$\langle p \rangle = \frac{N_h + 1}{(N_h + 1) + (N_t + 1)} = \frac{N_h + 1}{N + 2}$$

$$\sigma(p) = \sqrt{\text{Var}[p]} = \sqrt{\frac{(N_h + 1)(N_t + 1)}{(N + 2)^2(N + 3)}}$$

Pseudo-count of 1 head and 1 tail



600 heads out of 1000

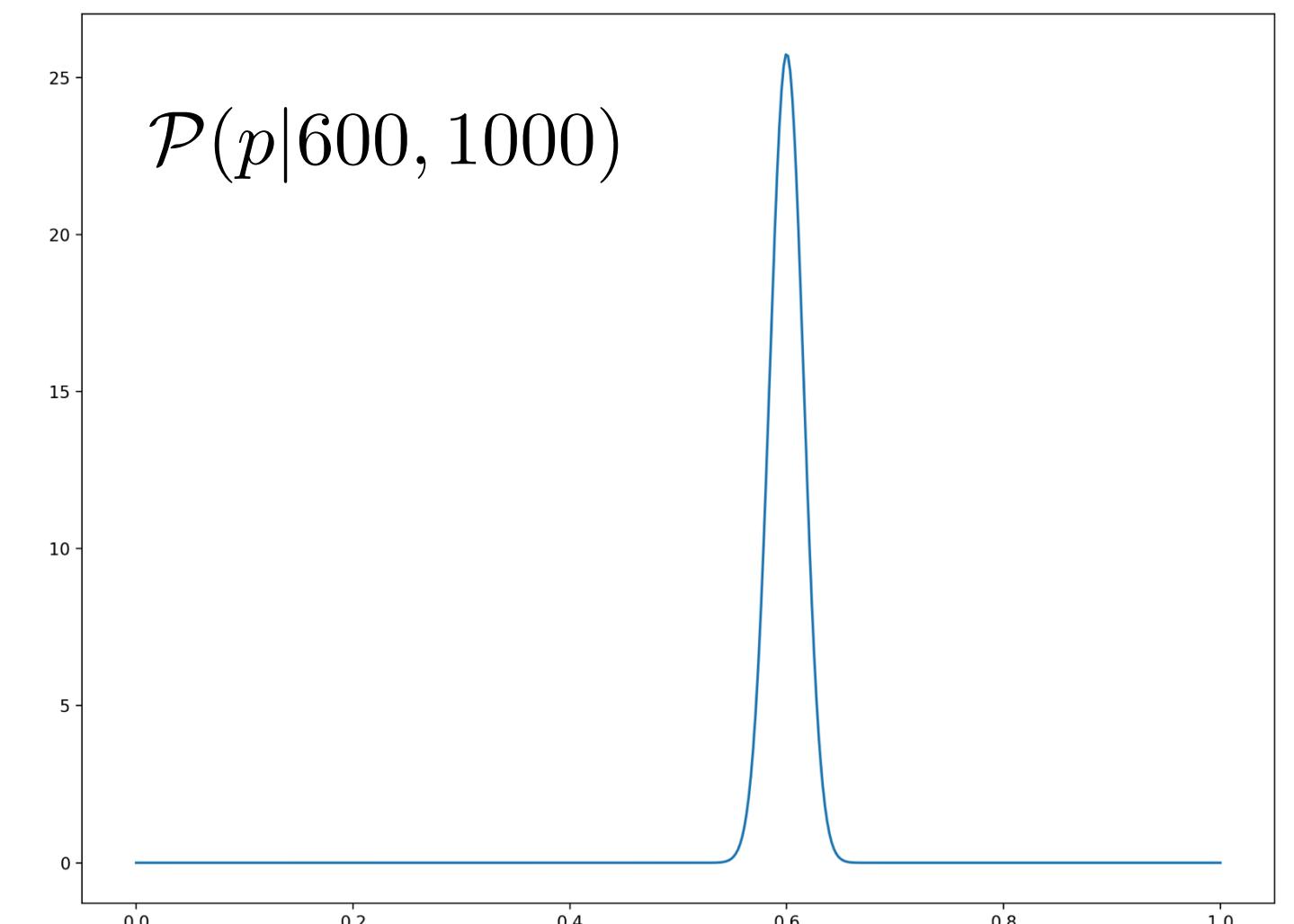
Examples:

6 heads out of 10: $\langle p \rangle \pm \sigma(p) = 0.58 \pm 0.14$

Could be fair.

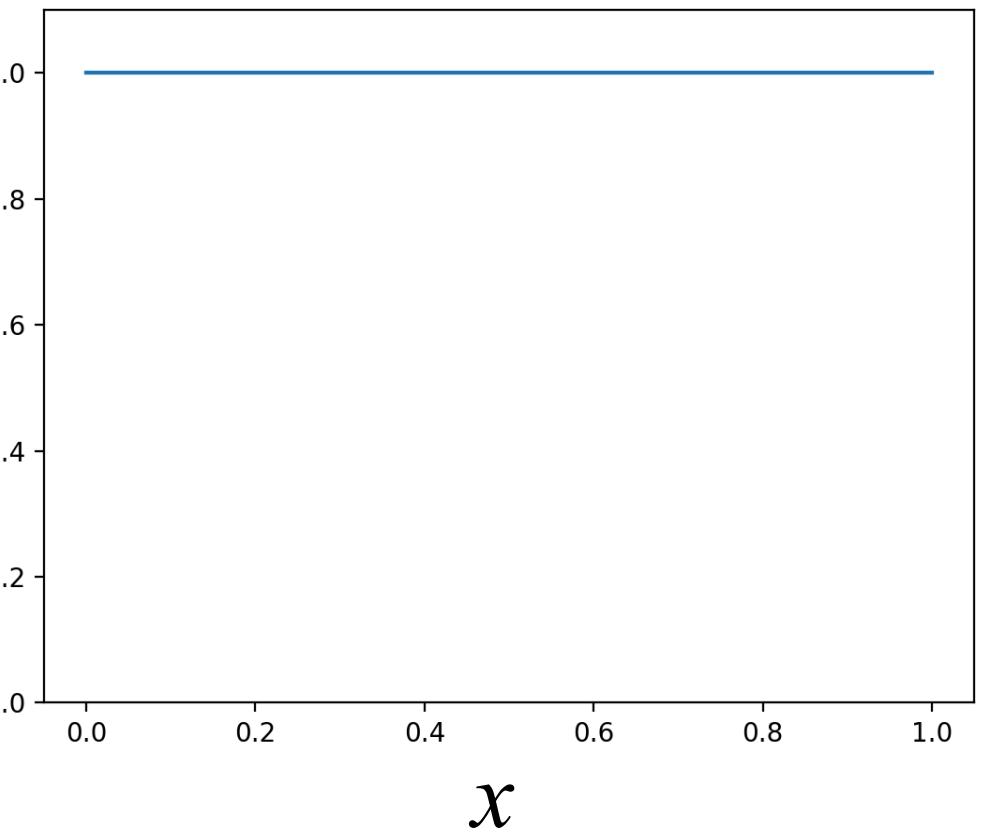
600 heads out of 1000: $\langle p \rangle \pm \sigma(p) = 0.60 \pm 0.02$

Is probably not fair

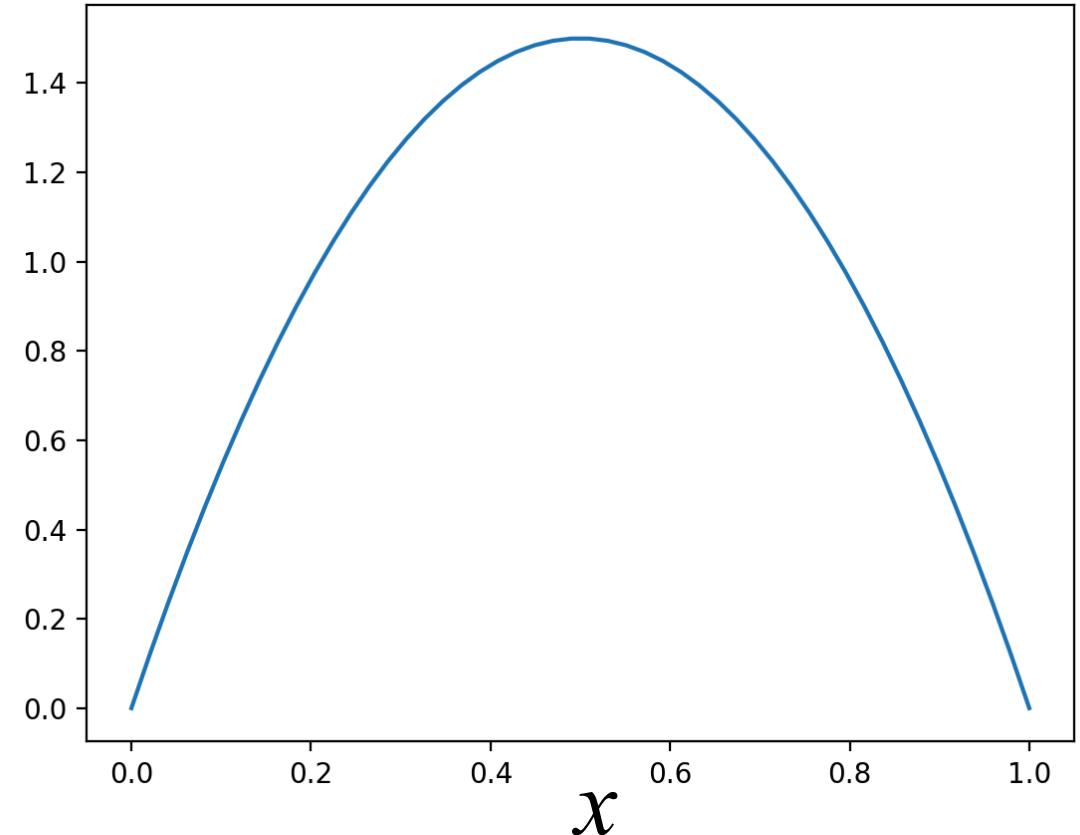


Is the coin fair? - changing the prior

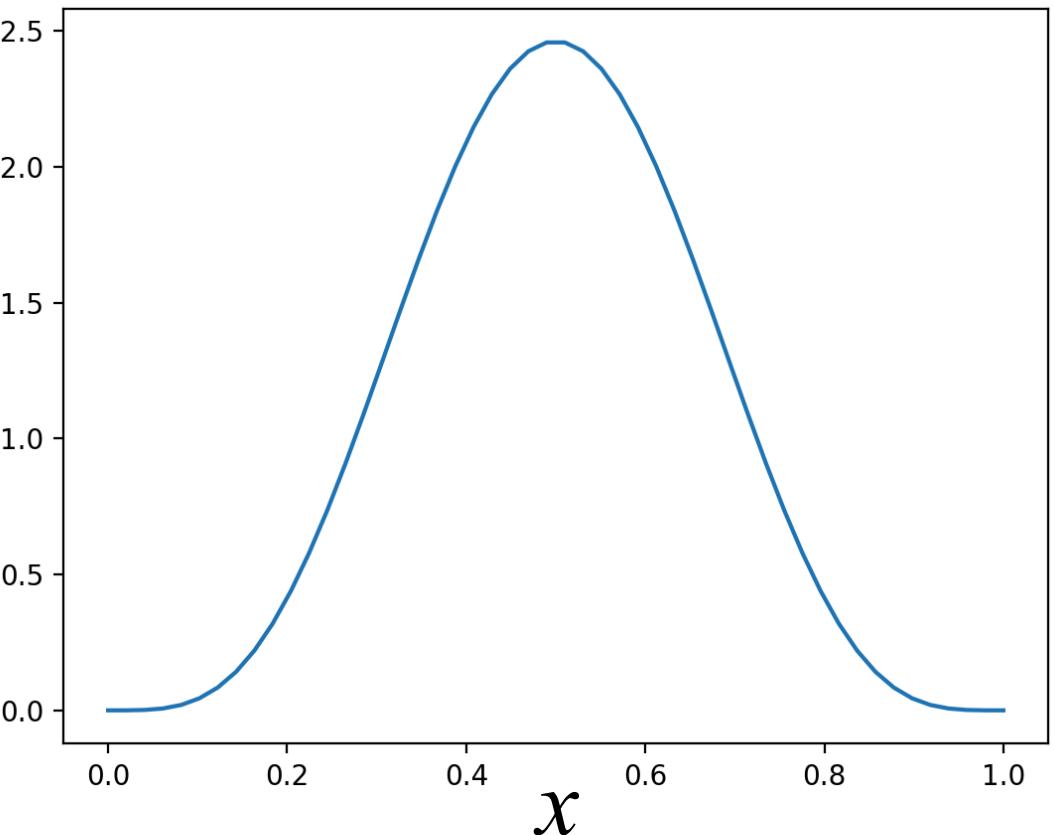
Constant prior: $\alpha = \beta = 1$



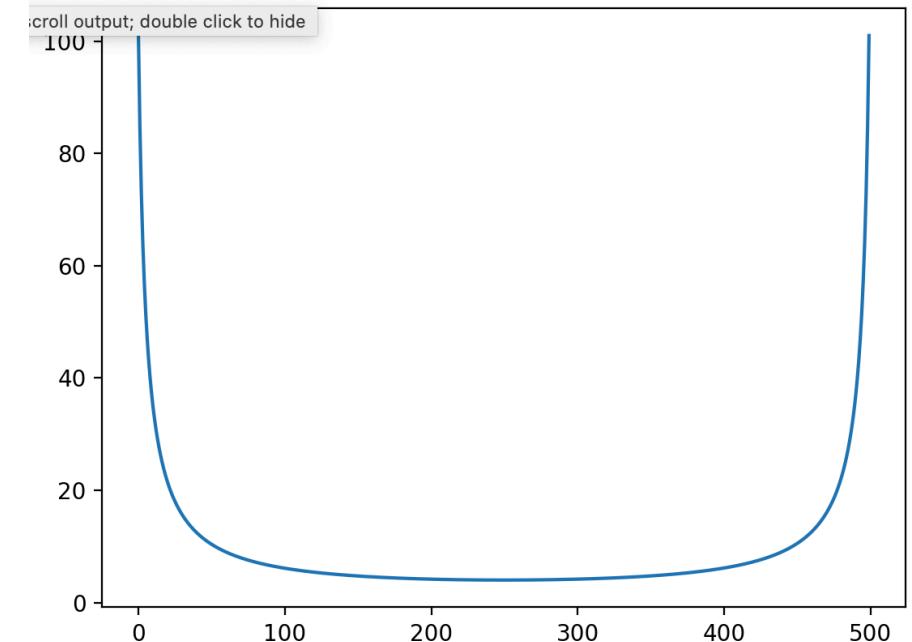
$\alpha = \beta = 2$



$\alpha = \beta = 5$



$\alpha = \beta = 0$



Beta distribution:

$$\mathcal{B}(x|\alpha, \beta) = \frac{1}{B(\alpha, \beta)} x^{\alpha-1} (1-x)^{\beta-1}$$

with
$$B(\alpha, \beta) = \frac{\Gamma(\alpha)\Gamma(\beta)}{\Gamma(\alpha + \beta)}$$

$$\langle x \rangle = \frac{\alpha}{\alpha + \beta}$$

$$\text{Var}[x] = \frac{\alpha\beta}{(\alpha + \beta)^2(\alpha + \beta + 1)}$$

Is the coin fair? - including a beta prior

$$\begin{aligned}\mathcal{P}(p|N_h, N) &\propto \mathcal{P}(N_h|p, N)\mathcal{P}(p) \propto \binom{N}{N_h} p^{N_h} (1-p)^{N_t} \mathcal{P}(p) \\ &\propto \mathcal{B}(p|\alpha = N_h + 1, \beta = N_t + 1)\mathcal{P}(p) \\ &\propto \mathcal{B}(p|\alpha = N_h + 1, \beta = N_t + 1)\mathcal{B}(p|\alpha = a, \beta = a) \\ &\propto \mathcal{B}(p|\alpha = N_h + a, \beta = N_t + a)\end{aligned}$$

This gives a pseudo-count of a .

The best estimate (the average) of p :

$$\langle p \rangle = \frac{N_h + a}{(N_h + a) + (N_t + a)} = \frac{N_h + a}{N + 2a}$$

Pseudo-count of a
heads and a tails

Probability to get heads: p
Tails: $N_t = N - N_h$

Beta distribution:

$$\mathcal{B}(x|\alpha, \beta) = \frac{1}{B(\alpha, \beta)} x^{\alpha-1} (1-x)^{\beta-1}$$

$$\text{with } B(\alpha, \beta) = \frac{\Gamma(\alpha)\Gamma(\beta)}{\Gamma(\alpha + \beta)}$$

$$\langle x \rangle = \frac{\alpha}{\alpha + \beta}$$

$$\text{Var}[x] = \frac{\alpha\beta}{(\alpha + \beta)^2(\alpha + \beta + 1)}$$

Machine learning

Method of least squares

How do we determine the best fit, i.e. the parameters w_i ?

Function values: $f(x_n)$ (not directly known)

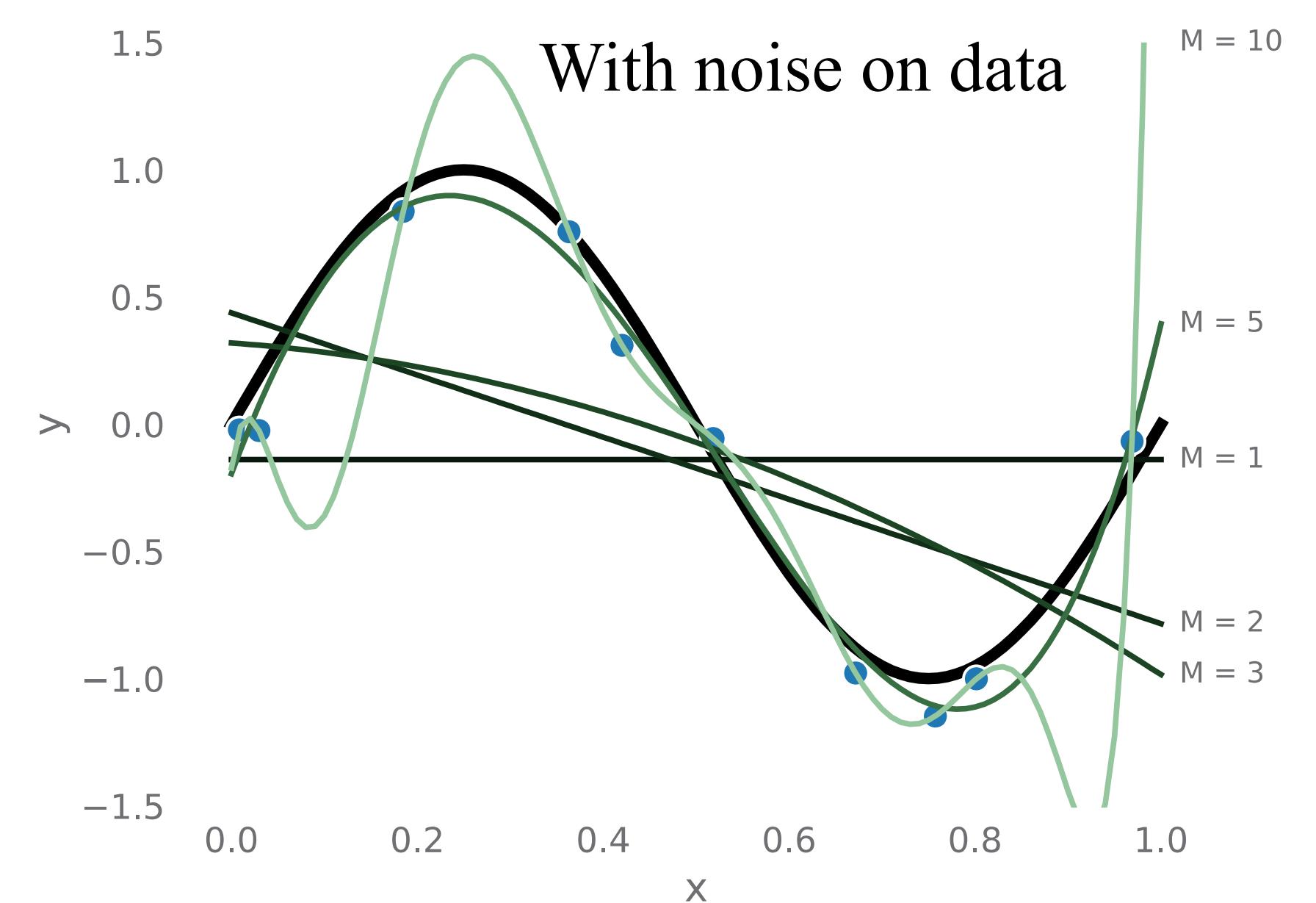
Measured values: $t_n = f(x_n) + \text{noise}$

Method of least squares:

Minimize the *cost function* C :

$$C(\mathbf{w}) = \frac{1}{2} \sum_{n=1}^N (t_n - y(x_n, \mathbf{w}))^2$$

$$y(x, \mathbf{w}) = w_0 + w_1 x + w_2 x^2 + \cdots + w_{M-1} x^{M-1}$$



From Bayes to least squares

$$y(x, \mathbf{w}) = w_0 + w_1 x + w_2 x^2 + \cdots + w_{M-1} x^{M-1}$$

Function values: $f(x_n)$ (not directly known)

Measured values: $t_n = f(x_n) + \text{noise}$

Assume, data is Gaussianly distributed with noise σ :

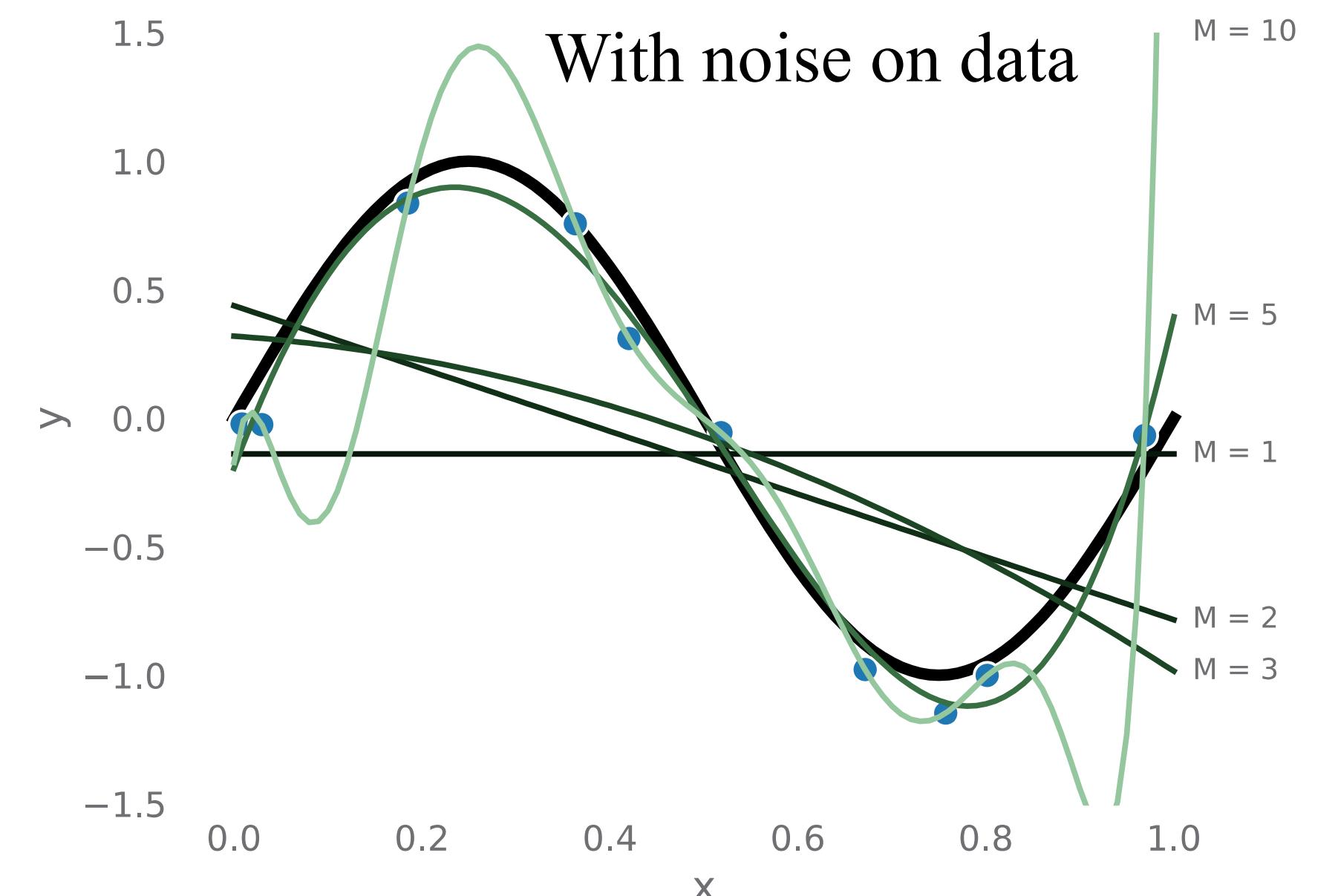
$$\mathcal{N}(t_n | f(x_n), \sigma^2)$$

The probability distribution for the parameters becomes

$$\begin{aligned} \mathcal{P}(\mathbf{w}|\mathbf{t}) \propto \mathcal{P}(\mathbf{t}|\mathbf{w})\mathcal{P}(\mathbf{w}) &= \prod_{n=1}^N \left(\frac{1}{(2\pi\sigma^2)^{(1/2)}} e^{-(t_n - y(x_n, \mathbf{w}))^2 / 2\sigma^2} \right) \mathcal{P}(\mathbf{w}) \\ &= \frac{1}{(2\pi\sigma^2)^{(N/2)}} e^{-\sum_n (t_n - y(x_n, \mathbf{w}))^2 / 2\sigma^2} \mathcal{P}(\mathbf{w}) \end{aligned}$$

Maximize the likelihood \rightarrow Minimize the cost function

$$C(\mathbf{w}) = \frac{1}{2} \sum_{n=1}^N (t_n - y(x_n, \mathbf{w}))^2$$



Linear basis function models

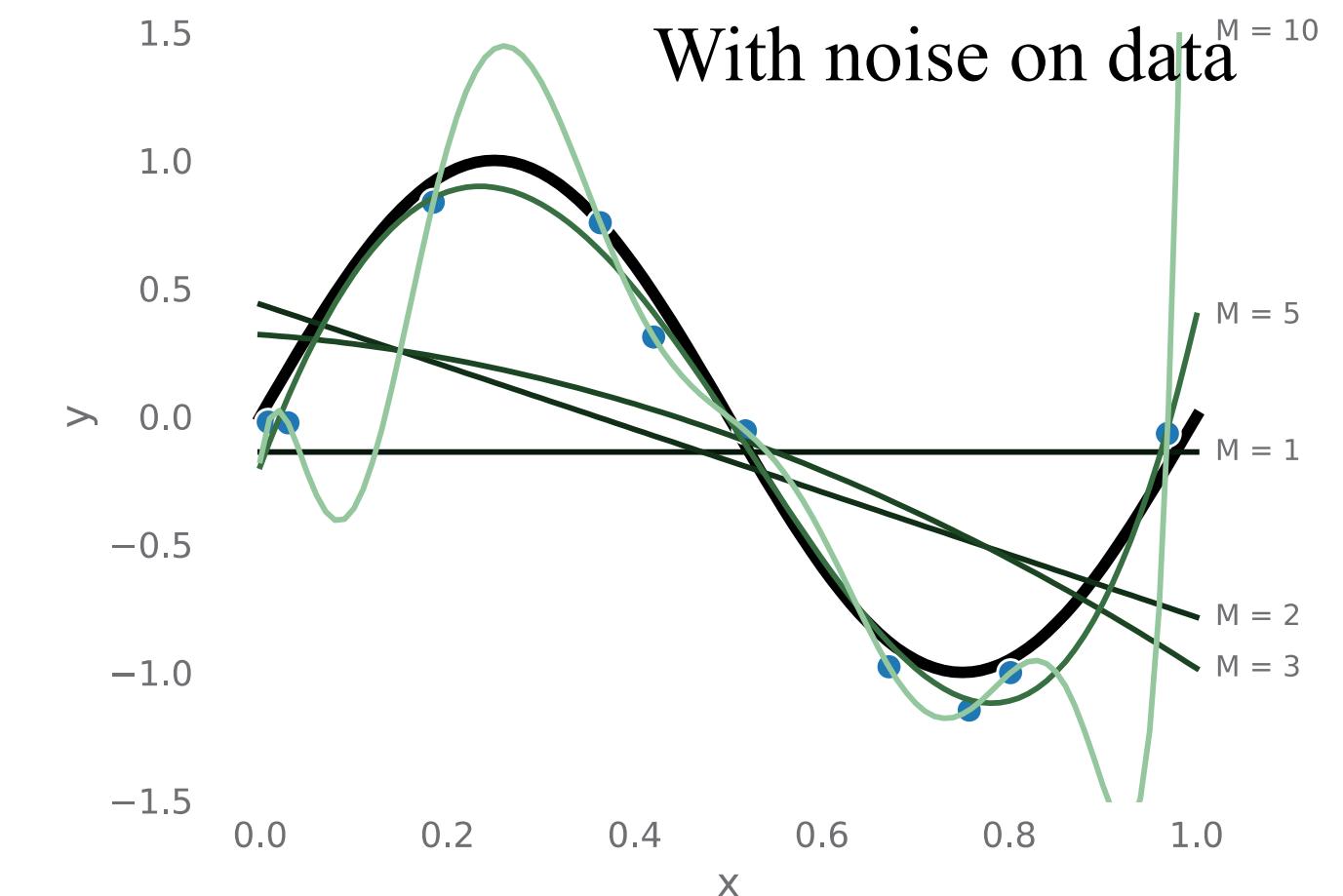
Linear model with basis functions: $y(\mathbf{x}, \mathbf{w}) = \sum_{i=1}^M w_i \phi_i(\mathbf{x}) = \mathbf{w}^T \boldsymbol{\phi}(\mathbf{x})$ $\boldsymbol{\phi} = (\phi_1, \phi_2, \dots, \phi_M)$

Probability distribution of parameters:

$$\begin{aligned}\mathcal{P}(\mathbf{w}|\mathbf{t}) &= \mathcal{P}(\mathbf{t}|\mathbf{w})\mathcal{P}(\mathbf{w}) = \frac{1}{(2\pi\sigma^2)^{(N/2)}} e^{-\sum_n(t_n - y(x_n, \mathbf{w}))^2/2\sigma^2} \mathcal{P}(\mathbf{w}) \\ &= \frac{1}{(2\pi\sigma^2)^{(N/2)}} e^{-C(\mathbf{w})/\sigma^2} \mathcal{P}(\mathbf{w})\end{aligned}$$

with cost function $C(\mathbf{w}) = \frac{1}{2} \sum_n (t_n - \sum_i w_i \phi_i(x_n))^2 = \frac{1}{2} (\mathbf{t} - \Phi \mathbf{w})^2$

$\Phi_{nj} = \phi_j(x_n)$ is the so-called *design matrix*. $N \times M$ matrix



N data points
 M basis functions

First step: Neglect prior, and maximize the likelihood \rightarrow Least squares method

Least squares method - exercise

Linear model with basis functions: $y(\mathbf{x}, \mathbf{w}) = \sum_{i=1}^M w_i \phi_i(\mathbf{x}) = \mathbf{w}^T \boldsymbol{\phi}(\mathbf{x})$ $\boldsymbol{\phi} = (\phi_1, \phi_2, \dots, \phi_M)$

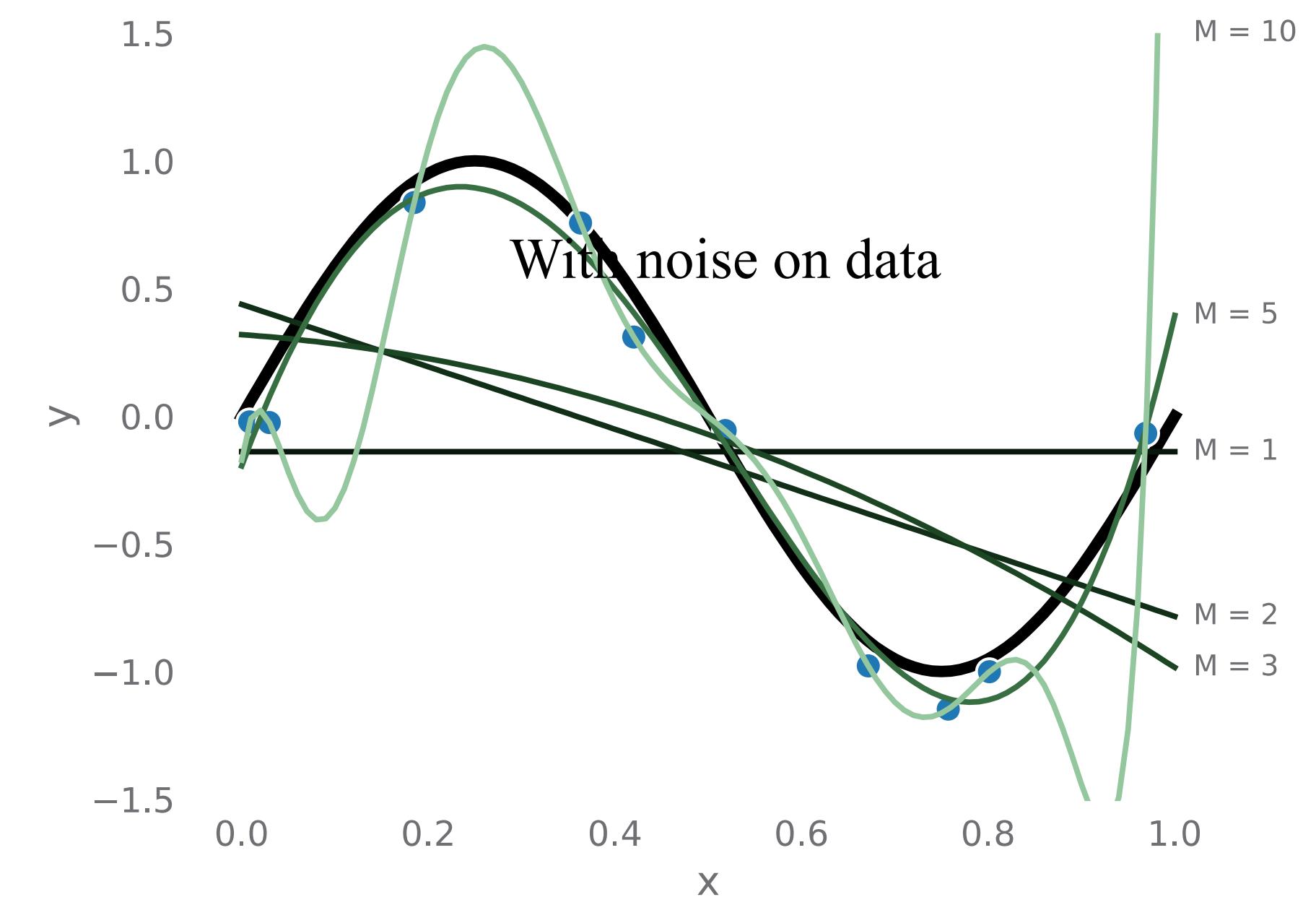
Minimize cost function: $C(\mathbf{w}) = \frac{1}{2} \sum_n (t_n - \sum_i w_i \phi_i(\mathbf{x}_n))^2 = \frac{1}{2} (\mathbf{t} - \mathbf{\Phi} \mathbf{w})^2$ $\Phi_{nj} = \phi_j(x_n)$

$$\frac{\partial C}{\partial \mathbf{w}} = \mathbf{\Phi}^T \mathbf{\Phi} \mathbf{w} - \mathbf{\Phi}^T \mathbf{t} = 0$$

Solution: $\mathbf{w}_0 = (\mathbf{\Phi}^T \mathbf{\Phi})^{-1} \mathbf{\Phi}^T \mathbf{t}$

Moore-Penrose inverse of $\mathbf{\Phi}$.

Do the exercise in Chapter 13 in the notes.



Least squares method - lessons learned

- Very simple to solve by matrix inversion
- Works well for simple models (small M)
- Overfitting for large M and parameters (w) become huge

Regularized least squares - including the prior

We would expect in advance that the coefficients \mathbf{w} should not be huge.
(A Taylor expansion of the sine function will give coefficients of the order one.)

$$\begin{aligned}\mathcal{P}(\mathbf{w}|\mathbf{t}) &\propto \mathcal{P}(\mathbf{t}|\mathbf{w})\mathcal{P}(\mathbf{w}) \propto e^{-\sum_n(t_n-y(x_n,\mathbf{w}))^2/2\sigma^2}\mathcal{P}(\mathbf{w}) \\ &\propto e^{-C(\mathbf{w})/\sigma^2}\mathcal{P}(\mathbf{w}) \propto e^{-C_{\text{reg}}(\mathbf{w})/\sigma^2}\end{aligned}$$

Take

$$\mathcal{P}(\mathbf{w}) \propto e^{-\frac{1}{2}\lambda\mathbf{w}^T\mathbf{w}} = e^{-\frac{1}{2}\lambda\sum_i|w_i|^2} \quad \longrightarrow \quad C_{\text{reg}}(\mathbf{w}) = \frac{1}{2}(\mathbf{t} - \Phi\mathbf{w})^2 + \frac{1}{2}\tilde{\lambda}\mathbf{w}^T\mathbf{w}$$
$$\tilde{\lambda} = \sigma^2\lambda$$

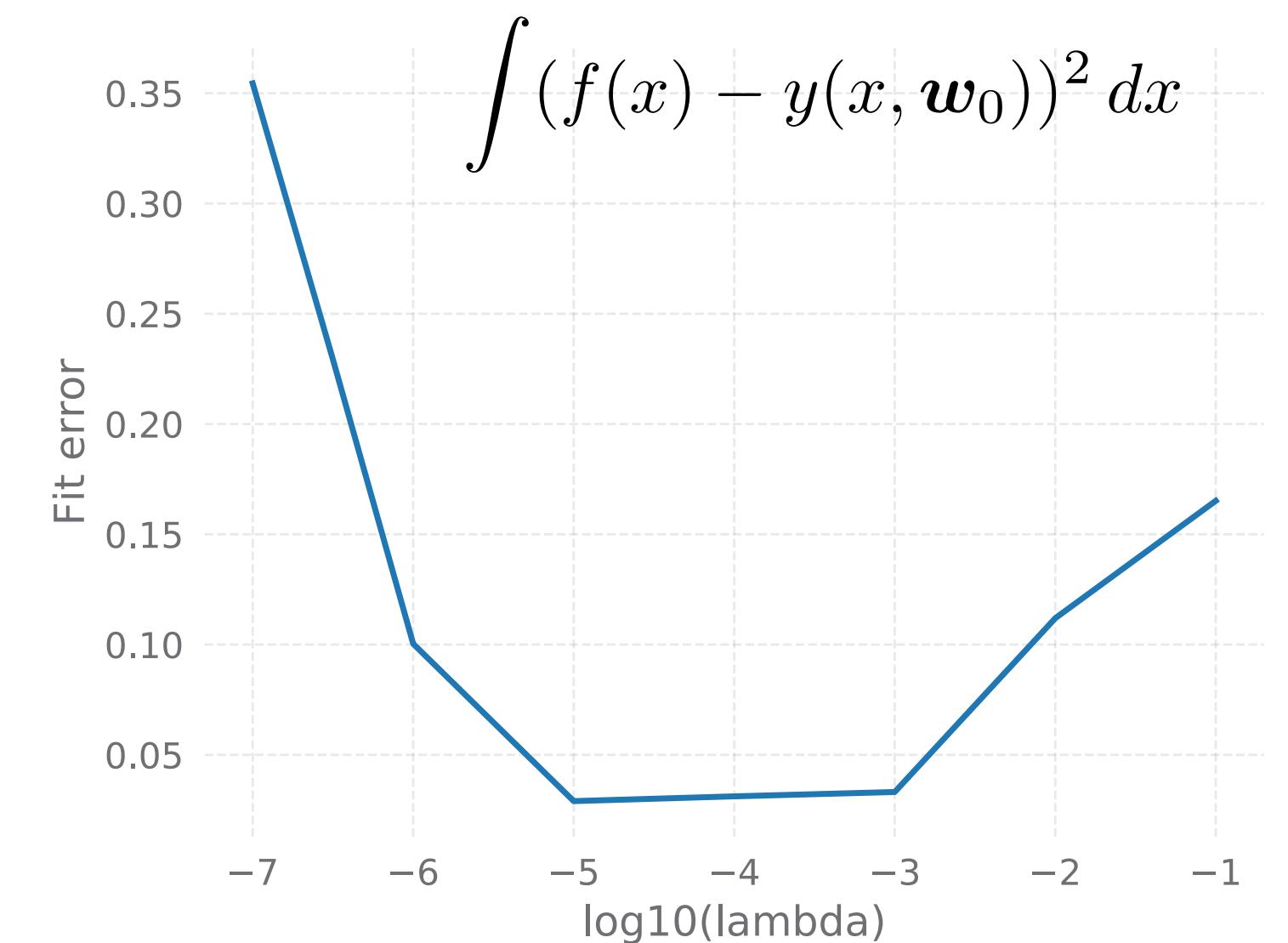
Maximizing the *posterior distribution* corresponds to minimizing the regularized cost function. This gives

$$\mathbf{w}_0 = (\Phi^T\Phi + \tilde{\lambda}\mathbf{I}_M)^{-1}\Phi^T\mathbf{t}$$

Now do the exercise on the regularized least squares method!

Regularized least squares - lessons learned

- The regularization makes it possible to control the size of the parameters, so they do not go crazy.
- This allows for more parameters, because the regularization is controlling the complexity of the model.
- Figure shows calculated error versus $\log_{10}(\tilde{\lambda})$.



Multi-dimensional Gaussian distribution

Consider two variables: y_1, y_2

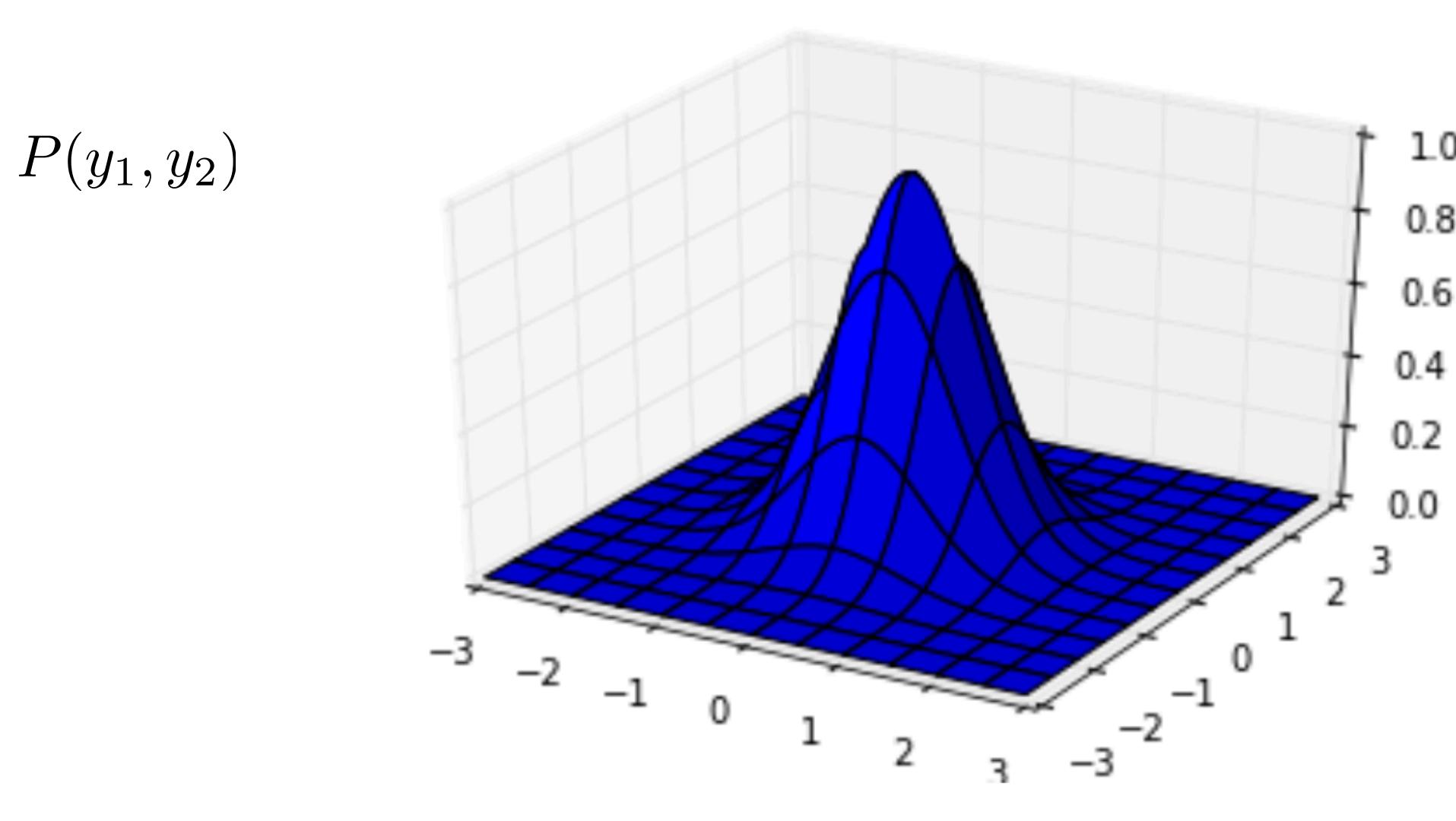
Assume the two variables are Gaussianly distributed with zero mean and some correlation between them:

$$\langle y_1^2 \rangle = \sigma^2, \quad \langle y_2^2 \rangle = \sigma^2, \quad \langle y_1 y_2 \rangle = \tau^2 \quad \text{or} \quad \langle \mathbf{y} \mathbf{y}^T \rangle = \begin{pmatrix} \sigma^2 & \tau^2 \\ \tau^2 & \sigma^2 \end{pmatrix} \equiv \mathbf{K}$$

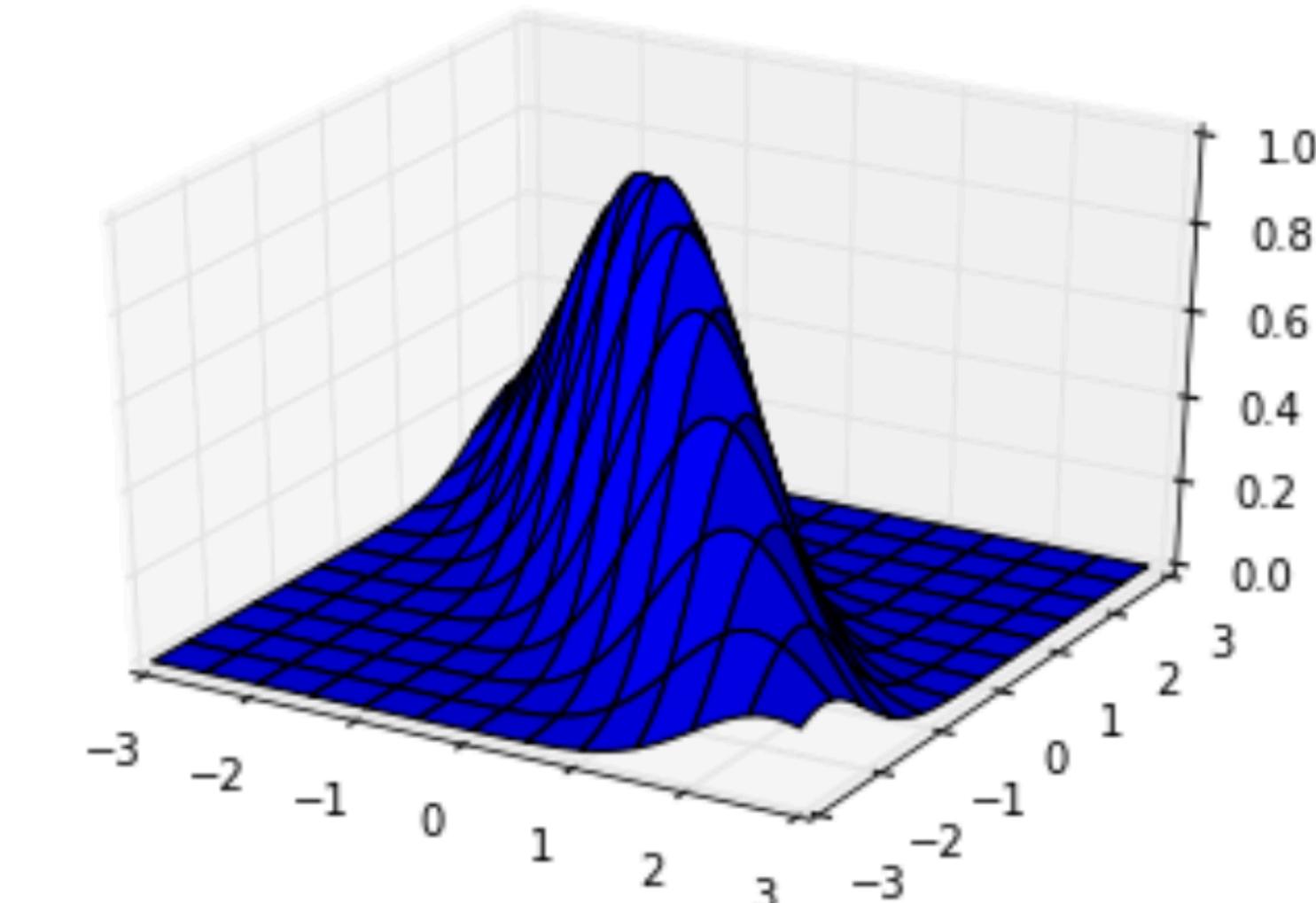
This corresponds to the probability distribution:

$$P_0(y_1, y_2) = \frac{1}{\sqrt{2\pi \det(\mathbf{K})}} \exp\left(-\frac{1}{2} \mathbf{y}^T \mathbf{K}^{-1} \mathbf{y}\right)$$

No correlation $\tau = 0$



Strong correlation $\tau \sim \sigma$



Multi-dimensional Gaussian distribution

Now we get the information that y_1 actually has the value y_1^0 .

What is then the probability distribution for y_2 ?

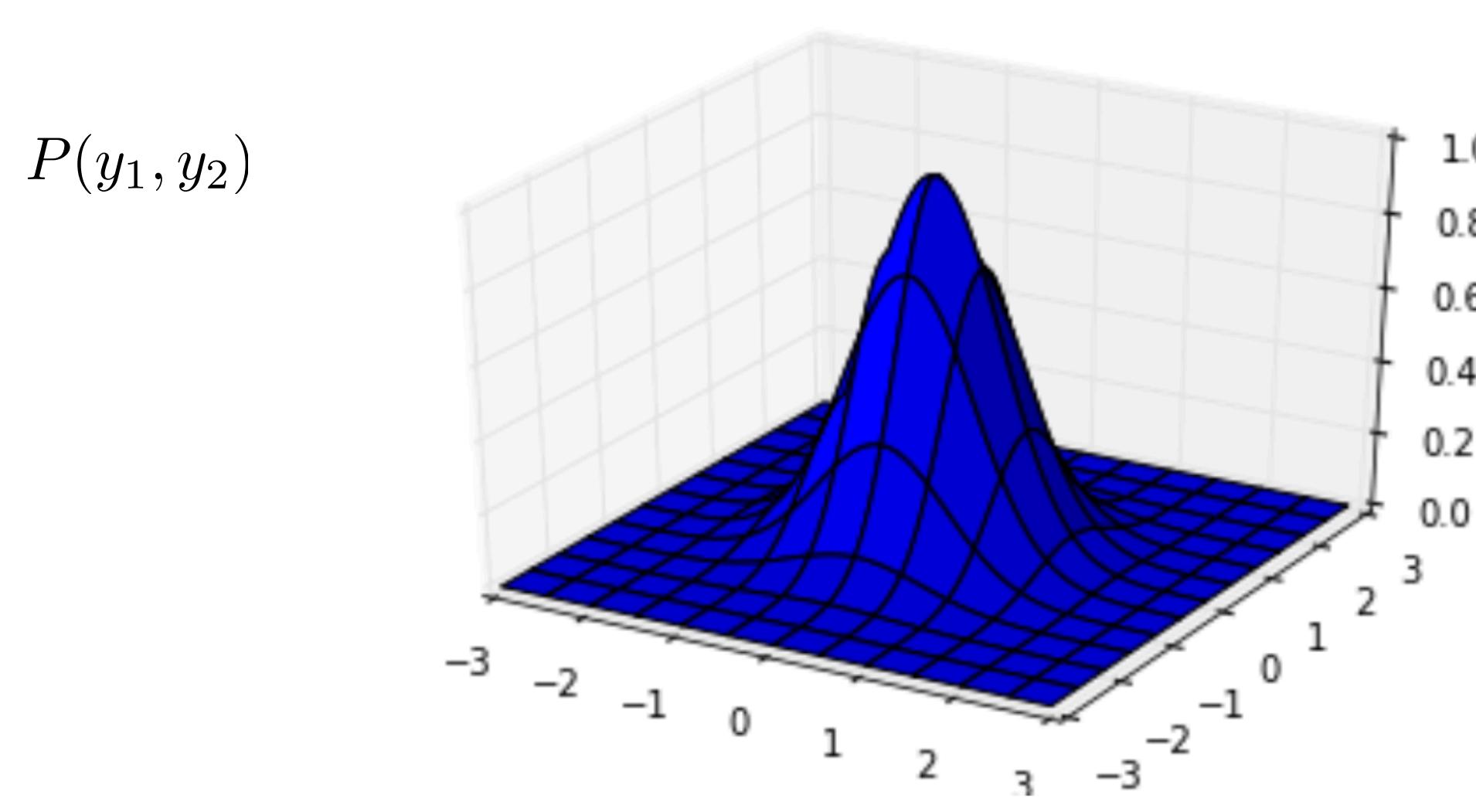
$$P(y_2) \propto P_0(y_1^0, y_2) \propto \exp \left[- \left(y_2 - \left(\frac{\tau}{\sigma} \right)^2 y_1^0 \right)^2 / 2(\sigma^2(1 - (\tau/\sigma)^4)) \right]$$

A new Gaussian!

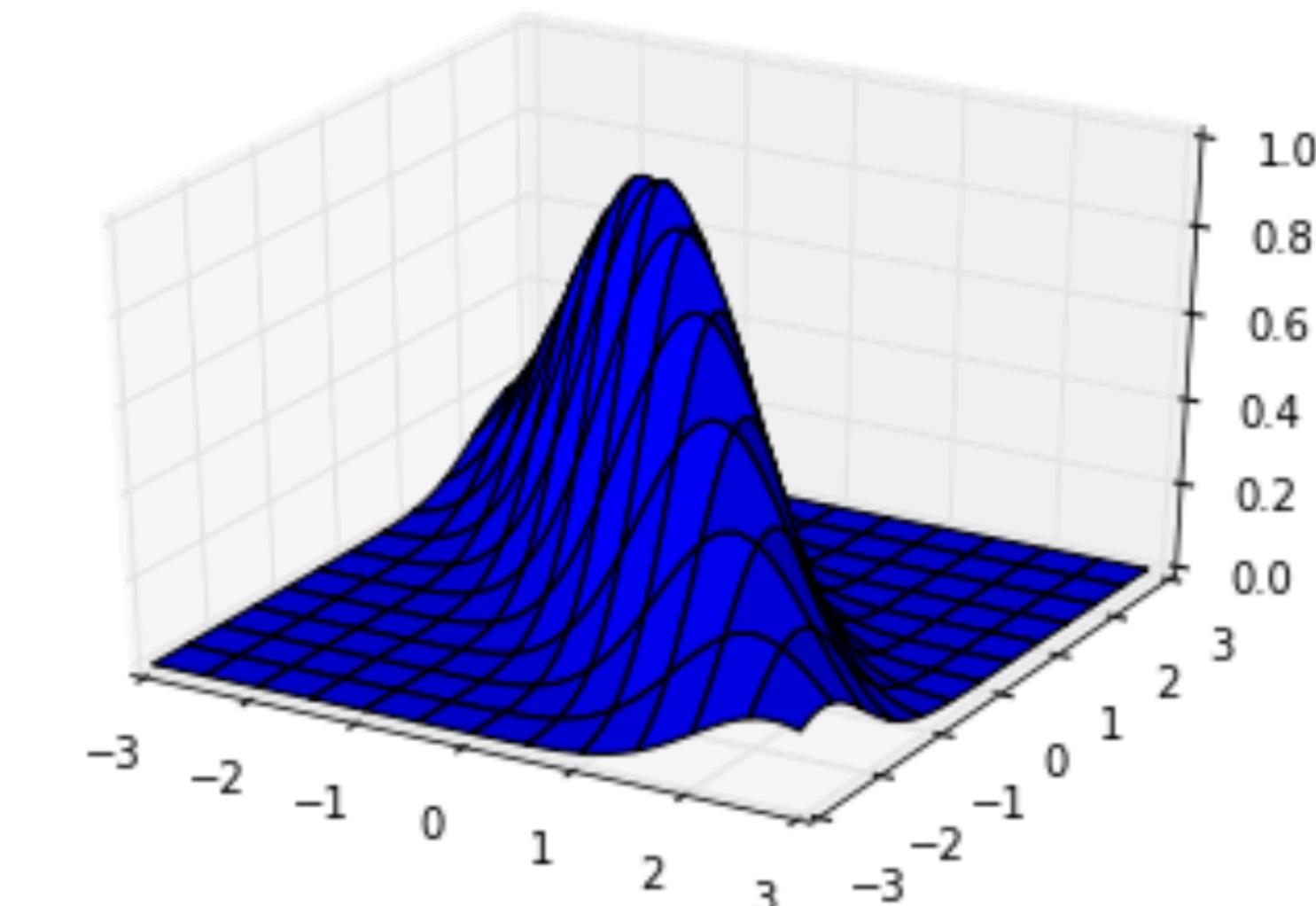
Mean: $\langle y_2 \rangle = \left(\frac{\tau}{\sigma} \right)^2 y_1^0$

Variance: $\langle (y_2 - \langle y_2 \rangle)^2 \rangle = \sigma^2(1 - (\tau/\sigma)^4)$

No correlation $\tau = 0$



Strong correlation $\tau \sim \sigma$



See appendix in notes

Multi-dimensional Gaussian distribution

$$\mathcal{N}(\mathbf{x}|\boldsymbol{\mu}, \boldsymbol{\Sigma}) = \frac{1}{(2\pi)^{D/2}} \frac{1}{\det(\boldsymbol{\Sigma})^{1/2}} \exp\left(-\frac{1}{2}(\mathbf{x} - \boldsymbol{\mu})^T \boldsymbol{\Sigma}^{-1} (\mathbf{x} - \boldsymbol{\mu})\right) \quad \mathbf{x} = \begin{pmatrix} x_1 \\ x_2 \\ \vdots \\ x_D \end{pmatrix}$$

$$\langle \mathbf{x} \rangle = \boldsymbol{\mu} \quad \text{and} \quad \text{Cov}[\mathbf{x}, \mathbf{x}^T] = \boldsymbol{\Sigma}$$

Convolution:

$$\mathcal{P}(\mathbf{t}) = \int \mathcal{N}(\mathbf{t}|\mathbf{y}, \boldsymbol{\Sigma}) \mathcal{N}(\mathbf{y}|\mathbf{y}_p, \mathbf{K}) d\mathbf{y} = \mathcal{N}(\mathbf{t}|\mathbf{y}_p, \mathbf{C}) \quad \text{where} \quad \mathbf{C} = \mathbf{K} + \boldsymbol{\Sigma}$$

Product:

$$\mathcal{N}(\mathbf{y}|\tilde{\mathbf{y}}, \mathbf{K}) \propto \mathcal{N}(\mathbf{y}|\mathbf{y}_1, \mathbf{K}_1) \mathcal{N}(\mathbf{y}|\mathbf{y}_2, \mathbf{K}_2) \quad \text{where} \quad \mathbf{K}^{-1} = \mathbf{K}_1^{-1} + \mathbf{K}_2^{-1} \quad \text{and} \quad \tilde{\mathbf{y}} = \mathbf{K} (\mathbf{K}_1^{-1} \mathbf{y}_1 + \mathbf{K}_2^{-1} \mathbf{y}_2)$$

Fully Bayesian: the posterior distribution

$$\mathcal{P}(\mathbf{w}|\mathbf{t}) = \mathcal{P}(\mathbf{t}|\mathbf{w})\mathcal{P}(\mathbf{w}) \propto e^{-\sum_n(t_n - y(x_n, \mathbf{w}))^2/2\sigma^2} e^{-\frac{1}{2}\lambda \mathbf{w}^T \mathbf{w}} = e^{-C_{reg}(\mathbf{w})/\sigma^2}$$

$$C_{\text{reg}}(\mathbf{w}) = \frac{1}{2}(\mathbf{t} - \Phi \mathbf{w})^2 + \frac{1}{2}\tilde{\lambda} \mathbf{w}^T \mathbf{w}$$

$\tilde{\lambda} = \sigma^2 \lambda$

Maximize the posterior: $\mathbf{w}_0 = (\Phi^T \Phi + \tilde{\lambda} \mathbf{I}_M)^{-1} \Phi^T \mathbf{t}$

Rewrite:

$$\begin{aligned} C_{\text{reg}}(\mathbf{w}) &= \frac{1}{2}(\mathbf{t} - \Phi \mathbf{w})^2 + \frac{1}{2}\tilde{\lambda} \mathbf{w}^T \mathbf{w} = \frac{1}{2}(\mathbf{t}^T - \mathbf{w}^T \Phi^T)(\mathbf{t} - \Phi \mathbf{w}) + \frac{1}{2}\tilde{\lambda} \mathbf{w}^T \mathbf{w} \\ &= \frac{1}{2}\mathbf{w}^T (\Phi^T \Phi + \tilde{\lambda}) \mathbf{w} + \text{ linear in } \mathbf{w} + \text{constant} \\ &= \frac{1}{2}(\mathbf{w} - \mathbf{w}_0)^T \Omega^{-1}(\mathbf{w} - \mathbf{w}_0) + \text{constant} \quad \text{with} \quad \Omega = \sigma^2(\Phi^T \Phi + \tilde{\lambda} \mathbf{I}_M)^{-1} \end{aligned}$$

So, the posterior can be written

$$\mathcal{P}(\mathbf{w}|\mathbf{t}) \propto \exp\left(-\frac{1}{2}(\mathbf{w} - \mathbf{w}_0)^T \Omega^{-1}(\mathbf{w} - \mathbf{w}_0)\right) = \mathcal{N}(\mathbf{w}|\mathbf{w}_0, \Omega)$$

Fully Bayesian: the posterior probability distribution

Posterior:

$$\text{with } \Omega = \sigma^2(\Phi^T \Phi + \tilde{\lambda} I_M)^{-1}$$

$$\mathcal{P}(\mathbf{w}|\mathbf{t}) \propto \exp\left(-\frac{1}{2}(\mathbf{w} - \mathbf{w}_0)^T \Omega^{-1} (\mathbf{w} - \mathbf{w}_0)\right) = \mathcal{N}(\mathbf{w}|\mathbf{w}_0, \Omega)$$

$$\mathbf{w}_0 = (\Phi^T \Phi + \tilde{\lambda} I_M)^{-1} \Phi^T \mathbf{t} = \frac{1}{\sigma^2} \Omega \Phi^T \mathbf{t}$$

So, we get for the average and the covariance: $\langle \mathbf{w} \rangle = \mathbf{w}_0$ and $\text{Cov}[\mathbf{w}^T, \mathbf{w}] = \Omega$

The predicted functions are given by: $y(x, \mathbf{w}) = \mathbf{w}^T \phi(x)$

So, the best fit function is: $y_0(x, \mathbf{w}) = \mathbf{w}_0^T \phi(x) = \sum_{i=1}^M w_{0,i} \phi_i(x)$

and the variance (i.e. the (square of) uncertainty of the prediction becomes

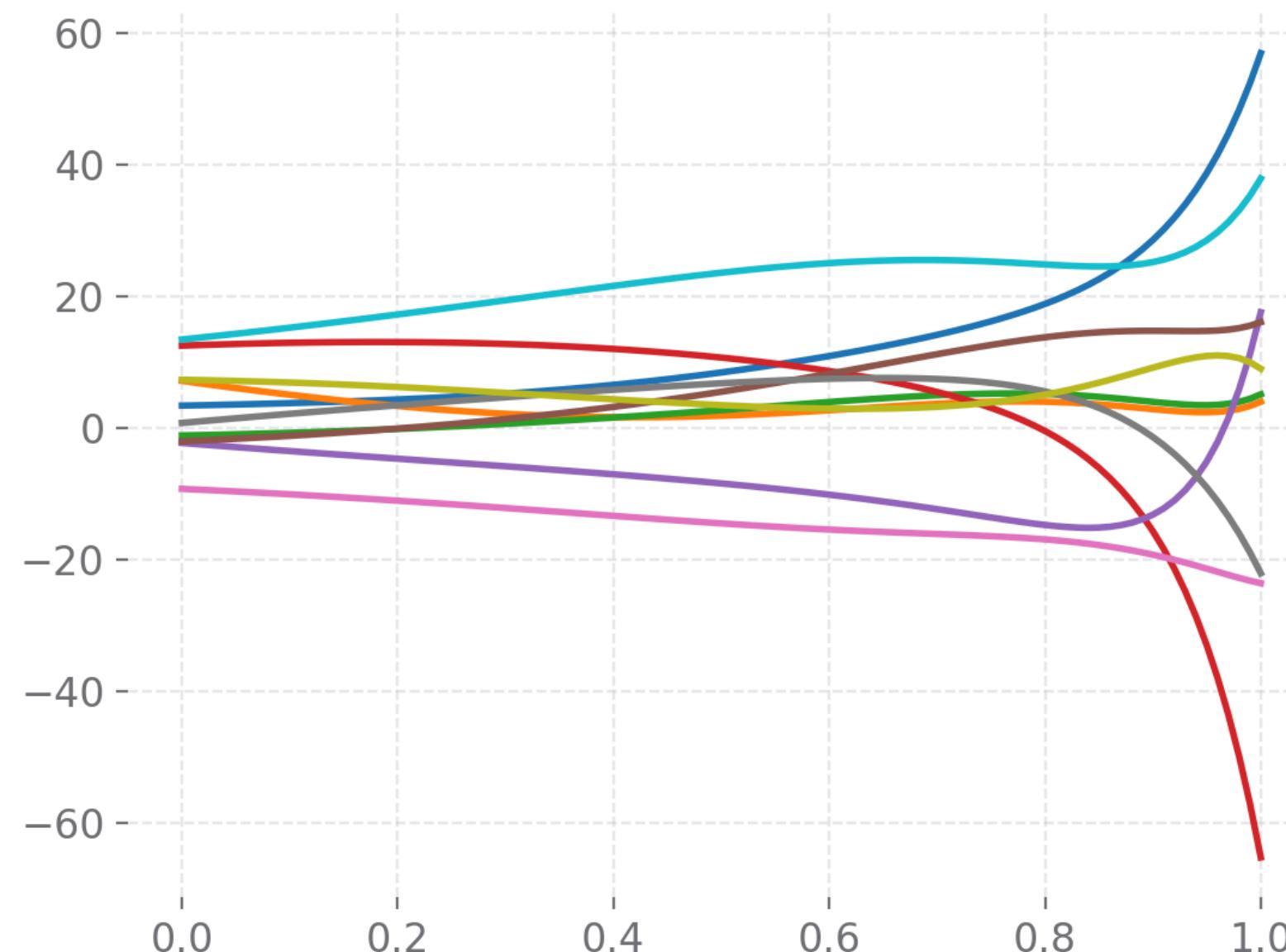
$$\text{Var}[y(x, \mathbf{w})] = \text{Cov}[y(x, \mathbf{w}), y(x, \mathbf{w})] = \text{Cov}[\phi^T(x) \mathbf{w}, \mathbf{w}^T \phi(x)] = \phi^T(x) \text{Cov}[\mathbf{w}, \mathbf{w}^T] \phi(x) = \phi^T(x) \Omega \phi(x)$$

Now do the exercise about Bayesian analysis of a linear basis function model!

Prior distribution

$$\mathcal{P}(\mathbf{w}) \propto e^{-\frac{1}{2}\lambda \mathbf{w}^T \mathbf{w}} = e^{-\frac{1}{2}\lambda \sum_i |w_i|^2}$$
$$\propto \mathcal{N}(\mathbf{w}|0, 1/\lambda * \mathbf{I}_M)$$

$$y(\mathbf{x}, \mathbf{w}) = \sum_{i=1}^M w_i \phi_i(\mathbf{x}) = \mathbf{w}^T \phi(\mathbf{x})$$



```
ngrid = 101
xg = np.linspace(0,1,num=ngrid)

M = 20
lamb = 10**(-2) # Do not use lambda here! That is a reserved word in Python.
Nens = 10

def basisf(j,x):
    return x**j

def fitf(x):
    return sum([w[j]*basisf(j,x) for j in range(M)])

for i in range(Nens):
    w=np.random.multivariate_normal(np.zeros(M),np.identity(M)/lamb)
    plt.plot(xg, fitf(xg))
```

Posterior distribution

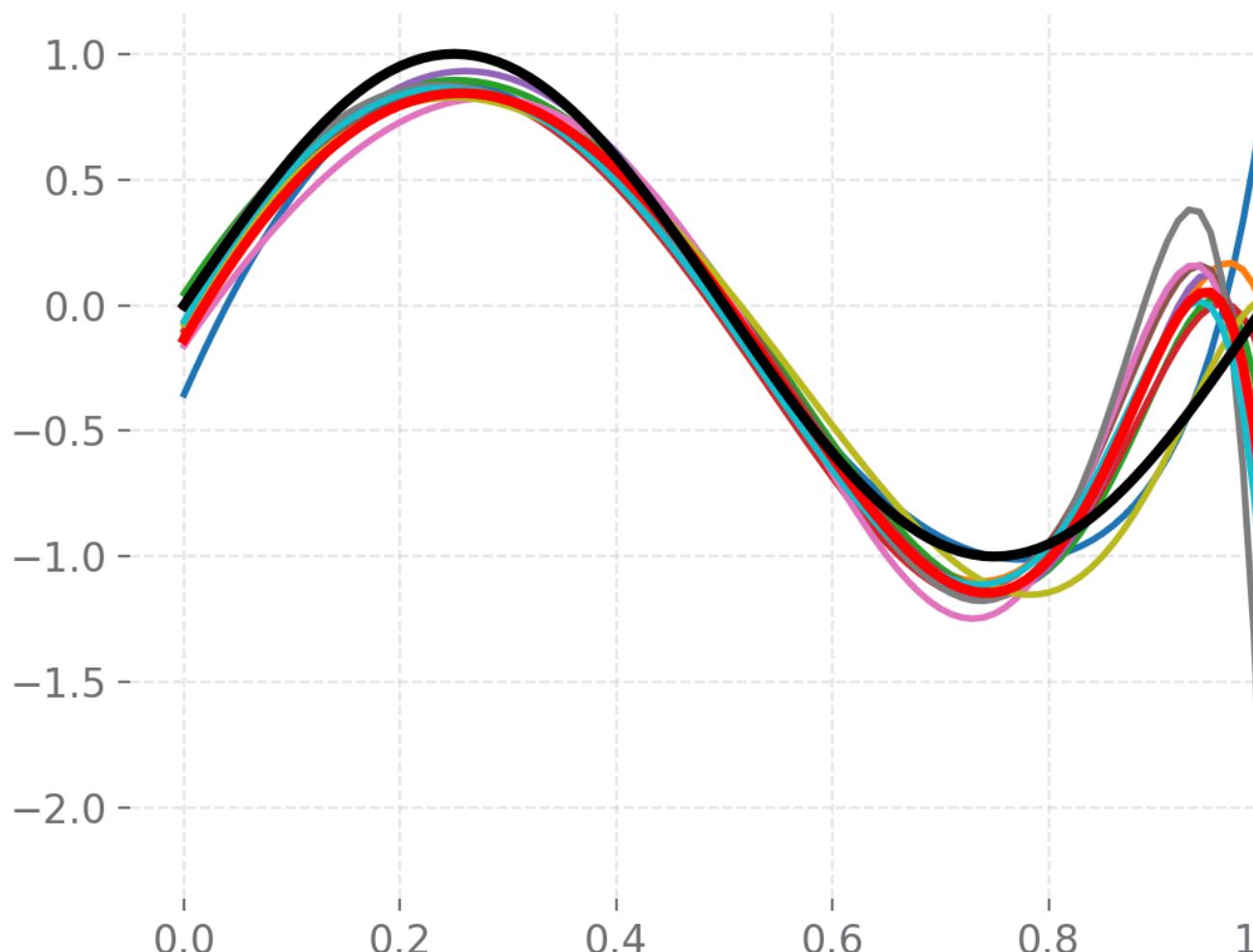
$$\begin{aligned}\mathcal{P}(\mathbf{w}|\mathbf{t}) &\propto \exp\left(-\frac{1}{2}(\mathbf{w} - \mathbf{w}_0)^T \boldsymbol{\Omega}^{-1} (\mathbf{w} - \mathbf{w}_0)\right) \\ &= \mathcal{N}(\mathbf{w}|\mathbf{w}_0, \boldsymbol{\Omega})\end{aligned}$$

$$\Phi_{nj} = \phi_j(x_n)$$

$$\mathbf{w}_0 = \frac{1}{\sigma^2} \boldsymbol{\Omega} \Phi^T \mathbf{t}$$

$$\boldsymbol{\Omega} = \sigma^2 (\Phi^T \Phi + \tilde{\lambda} \mathbf{I}_M)^{-1}$$

$$y(\mathbf{x}, \mathbf{w}) = \sum_{i=1}^M w_i \phi_i(\mathbf{x}) = \mathbf{w}^T \boldsymbol{\phi}(\mathbf{x})$$



```
initwithnoise()
M = 2*N
lambtildelde = 10**(-4)
Nens = 10
sigma = 0.1
ngrid = 101
xg = np.linspace(0,1,num=ngrid)

def basisf(j,x):
    return x**j

def fitf(x):
    return sum([w[j]*basisf(j,x) for j in range(M)])

phi = np.zeros((N,M))
for n, x in enumerate(xp):
    for j in range(M):
        phi[n,j] = basisf(j,x)
phit = phi.transpose()

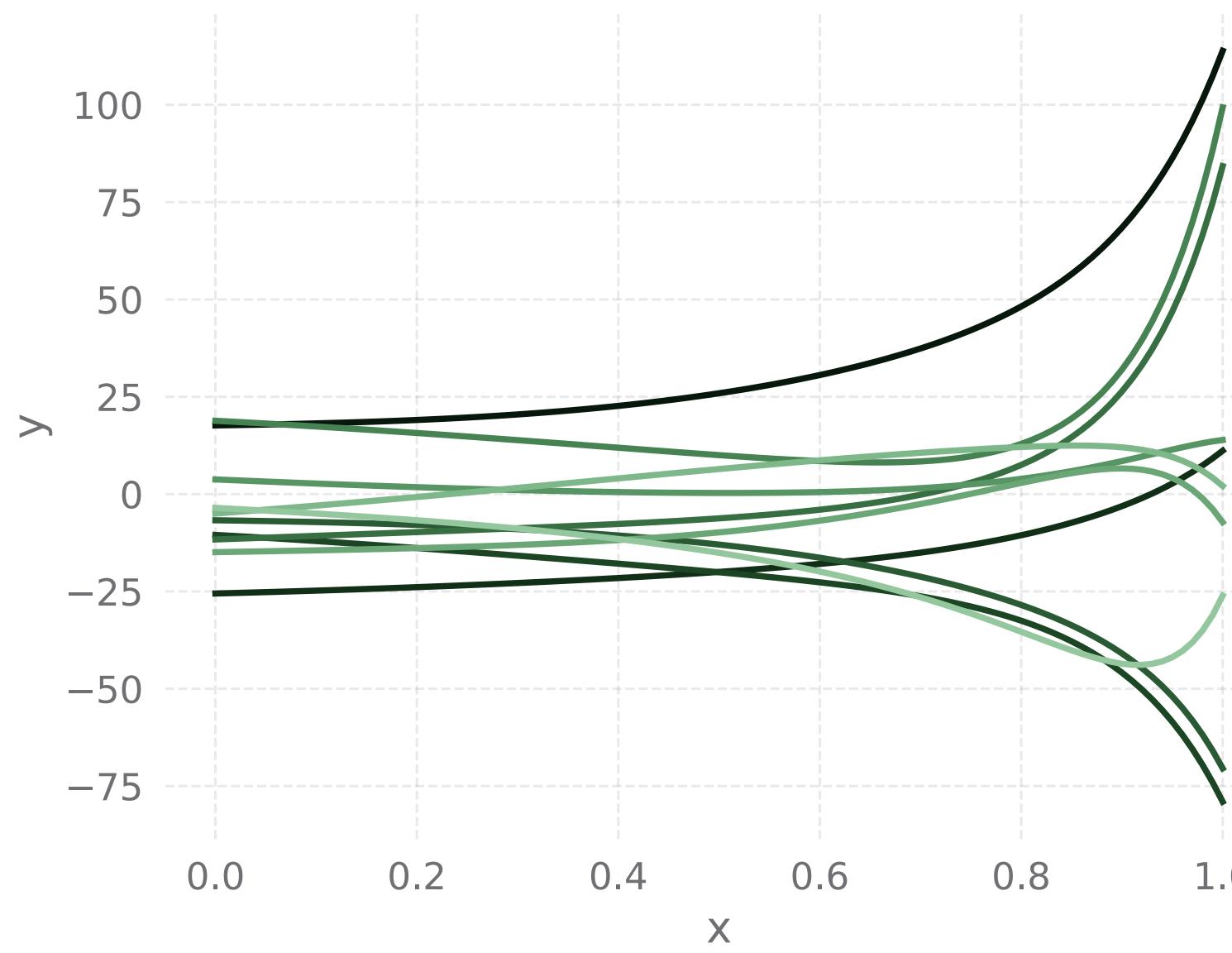
omega = sigma**2*np.linalg.inv(phit @ phi + lambtildelde*np.identity(M))

w0 = np.dot(omega/sigma**2, np.dot(phi.transpose(),tp))

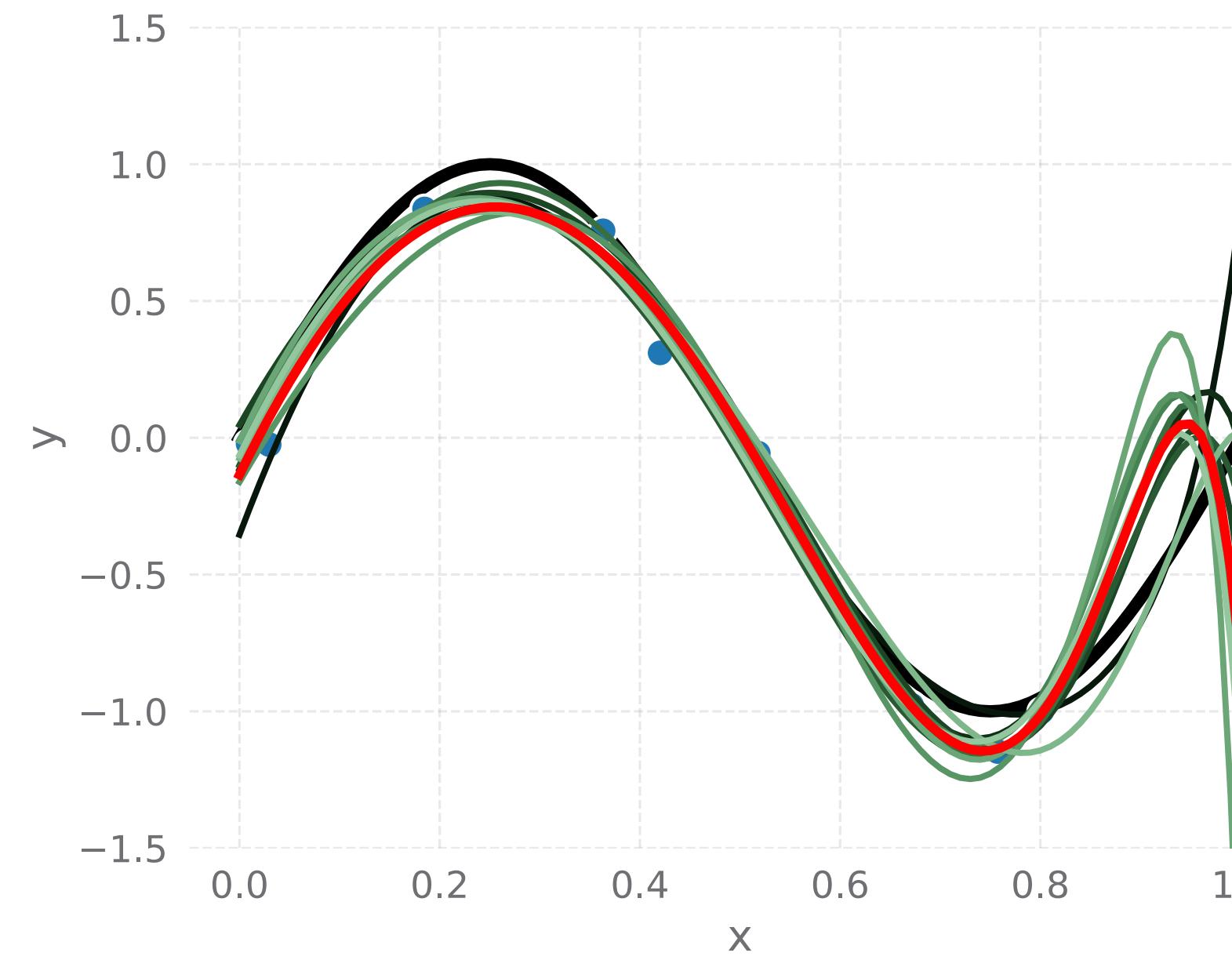
for i in range(Nens):
    w=np.random.multivariate_normal(w0,omega)
    plt.plot(xg, fitf(xg))
w=w0
plt.plot(xg, fitf(xg), color='r', lw=3)
plt.plot(xg,np.sin(2*np.pi*xg), color='k', lw=3)
```

Exercise about Bayesian analysis of a linear basis function model - lessons learned

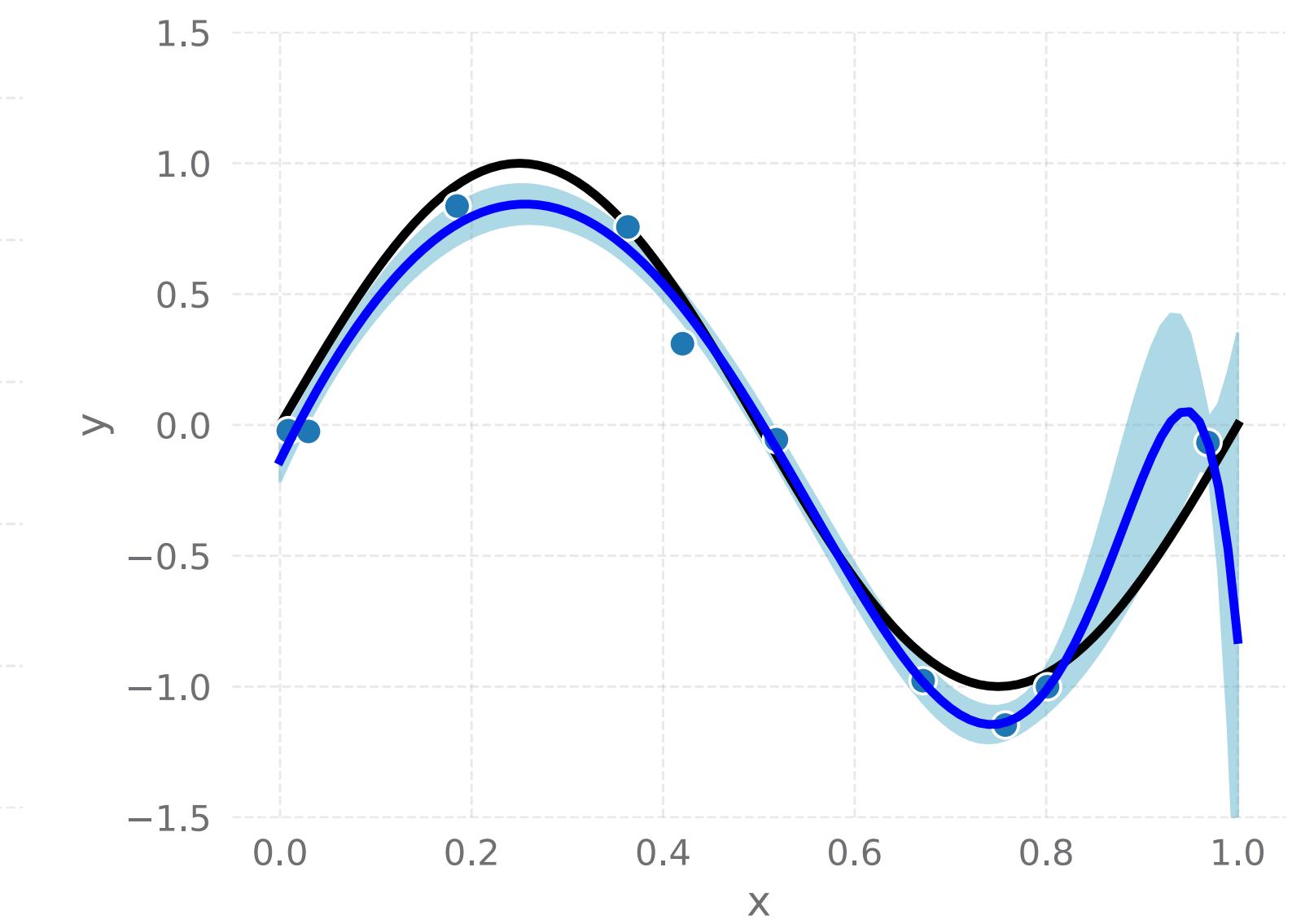
Ensemble of functions drawn from the prior distribution



Ensemble of functions drawn from the posterior distribution



Best fit function with error bar (one standard deviation)



And now for something completely different: Kernel regression

We now choose the basis functions as a Gaussian function *at each data point*.

$$\mathbf{g}(x)^T = (g(x, x_1), g(x, x_2), \dots, g(x, x_N))$$

$$g(x, x') = \exp(-(x - x')^2 / 2\ell^2)$$

$$y(x, \mathbf{a}) = \sum_{n=1}^N g(x, x_n) a_n = \mathbf{g}(x)^T \mathbf{a}$$

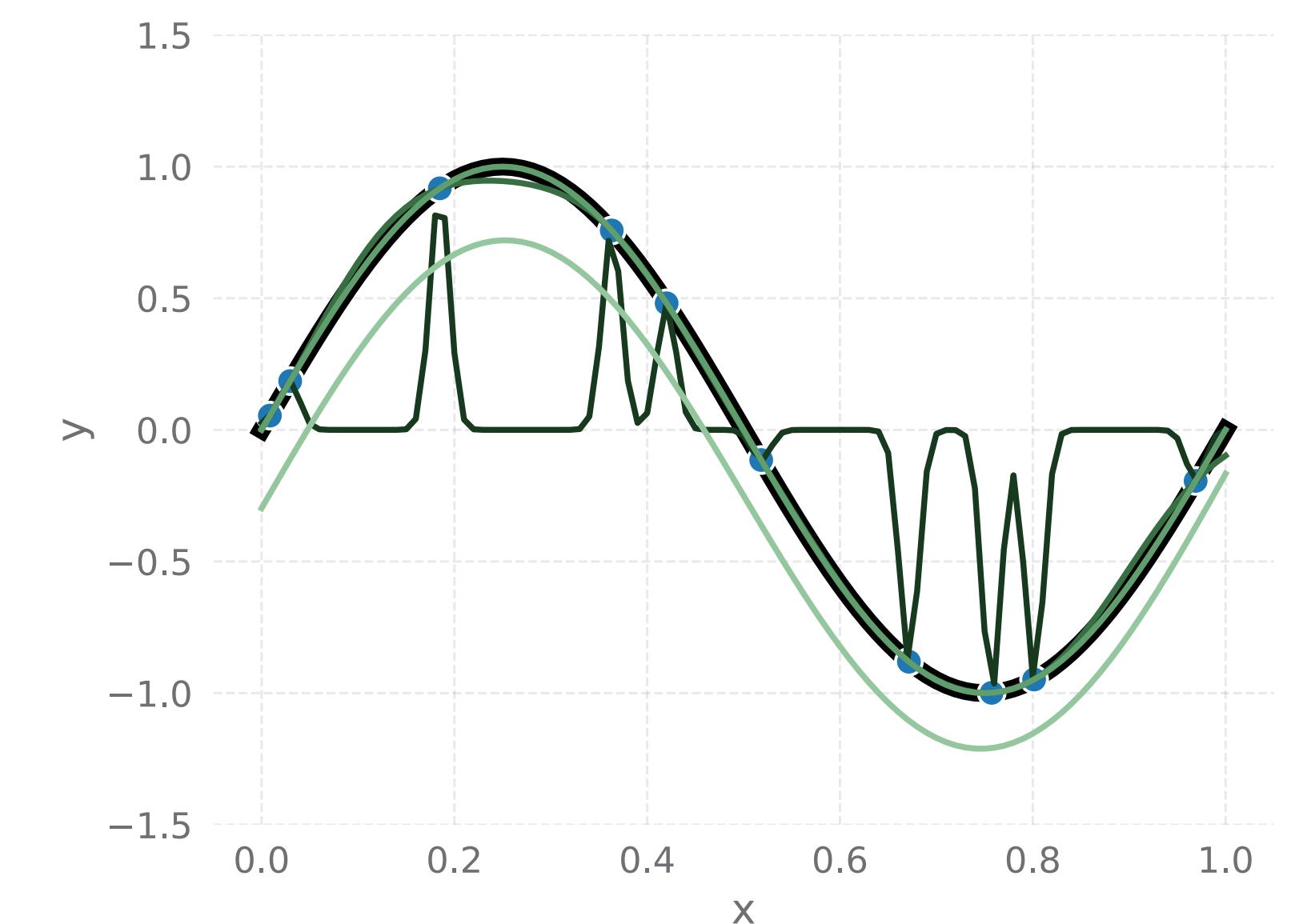
Same number of data points and parameters!

$$G_{nm} = g(x_n, x_m)$$

$$\mathbf{y} = G\mathbf{a}_0 \implies \mathbf{a}_0 = G^{-1}\mathbf{y}$$

Prediction:

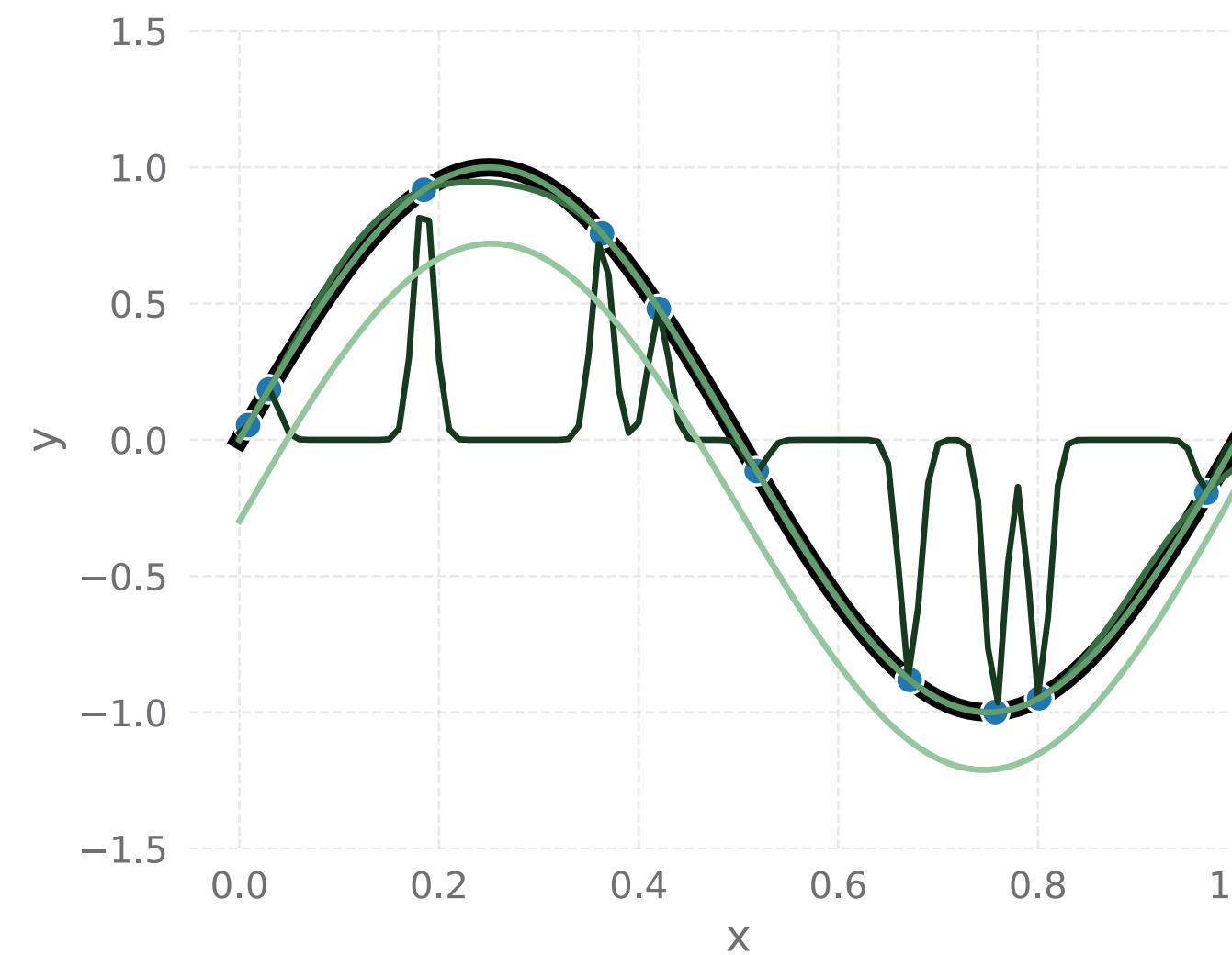
$$y(x, \mathbf{a}_0) = \mathbf{g}(x)^T \mathbf{a}_0 = \mathbf{g}(x)^T G^{-1} \mathbf{y}$$



Exercise on kernel regression - lessons learned

$$g(x, x') = \exp(-(x - x')^2/2\ell^2)$$

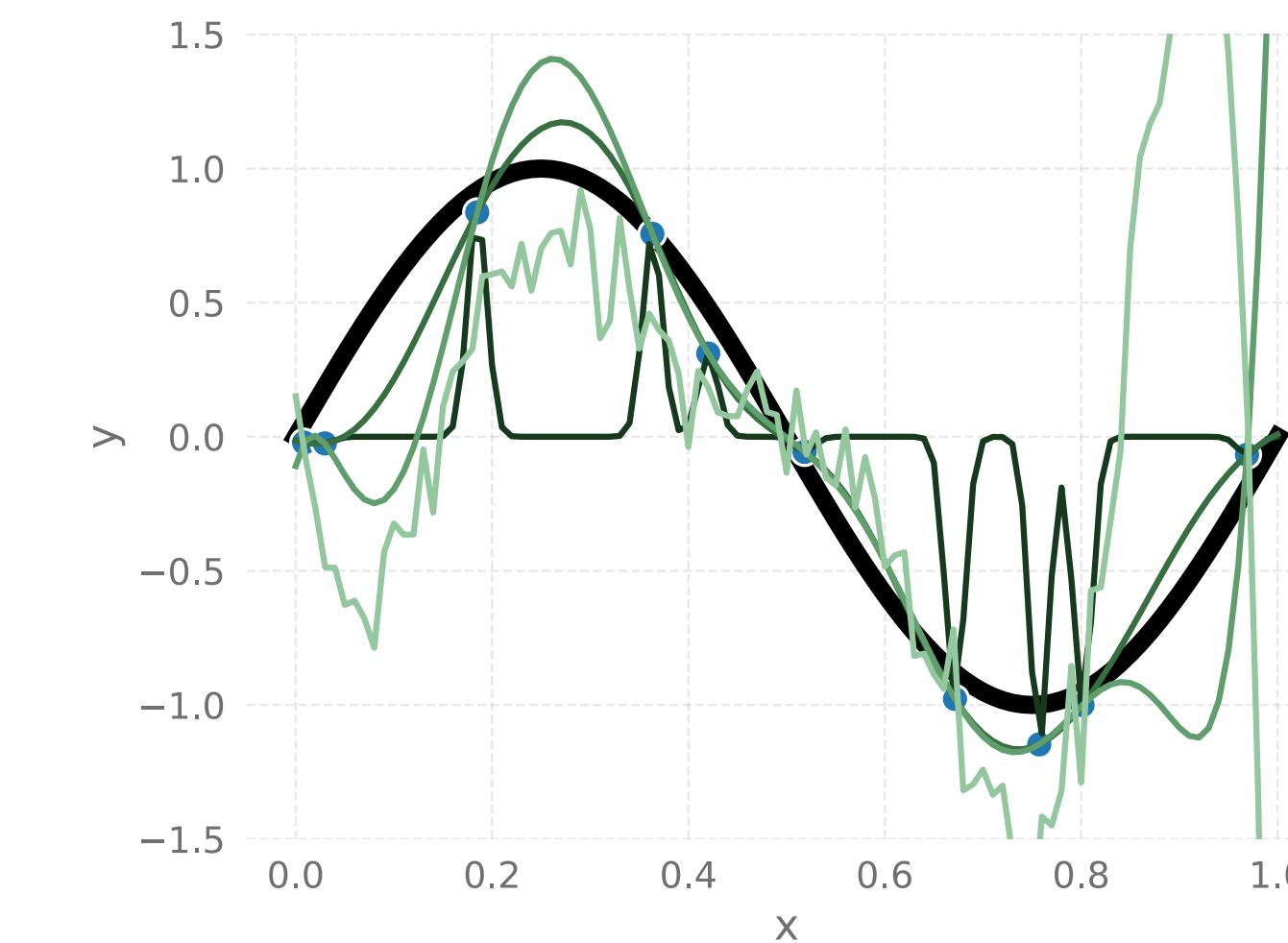
Without noise



$$\ell = 0.01, 0.1, 0.5 \text{ and } 1.0$$

Length scale very important
for quality of fit!

With noise



All noise is fitted. Severe problems
for both short and long scales

Kernel ridge regression

What if we have noise on the data $\mathbf{t} = \mathbf{y} + \text{noise}$?

The solution from before reproduces all noise!

$$y(x, \mathbf{a}_0) = \mathbf{g}(x)^T \mathbf{a}_0 = \mathbf{g}(x)^T \mathbf{G}^{-1} \mathbf{t}$$

Invent a regularized cost function:

$$C_{\text{reg}} = (\mathbf{t} - \mathbf{G}\mathbf{a})^2 + \lambda' \mathbf{a}^T \mathbf{G} \mathbf{a}$$

New solution:

$$\mathbf{a}_0 = (\mathbf{G} + \lambda' \mathbf{I}_M)^{-1} \mathbf{t}$$

New prediction:

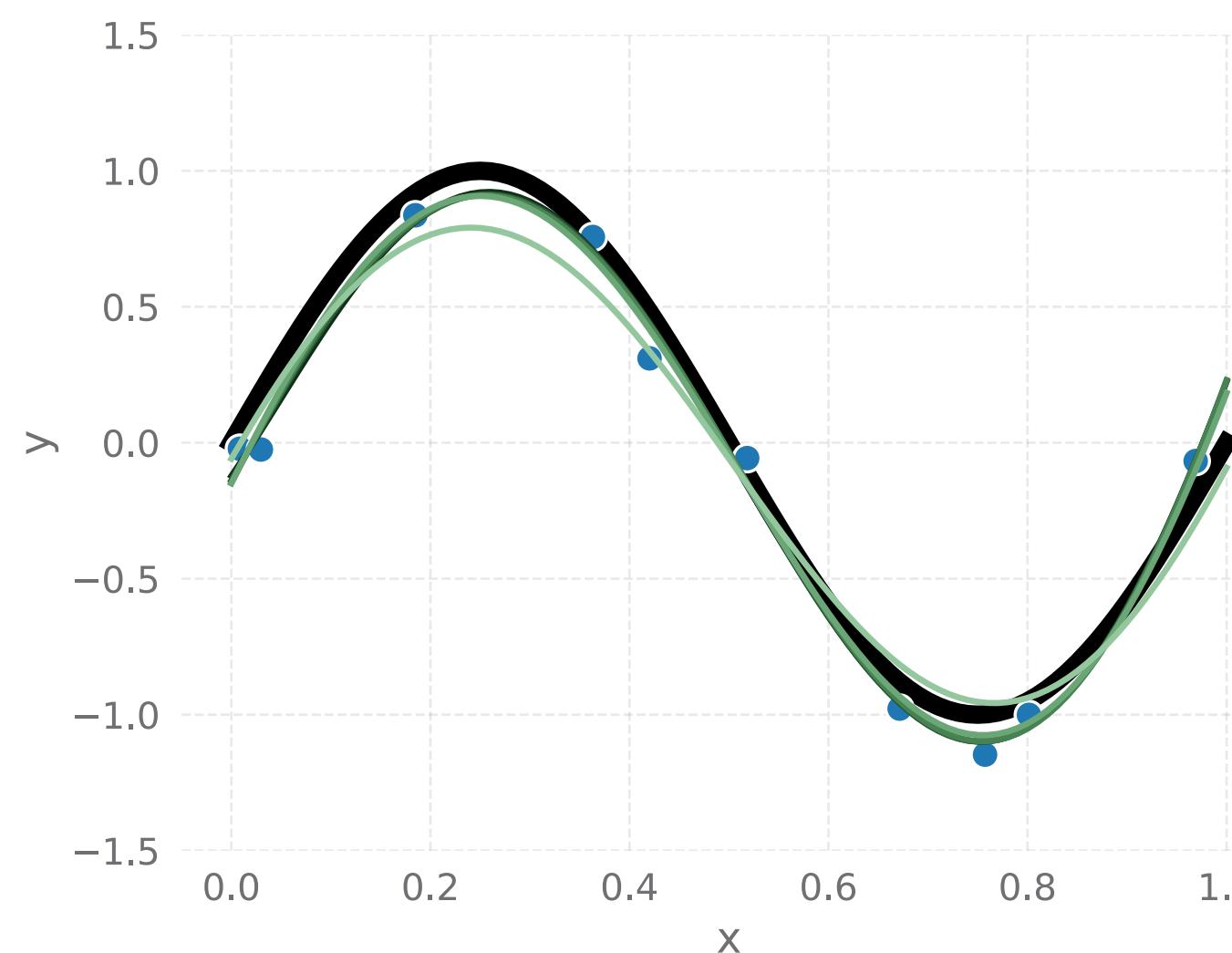
$$y(x, \mathbf{a}_0) = \mathbf{g}(x)^T \mathbf{a}_0 = \mathbf{g}(x)^T (\mathbf{G} + \lambda' \mathbf{I}_N)^{-1} \mathbf{t}$$

Now do the exercises on kernel regression and kernel ridge regression!

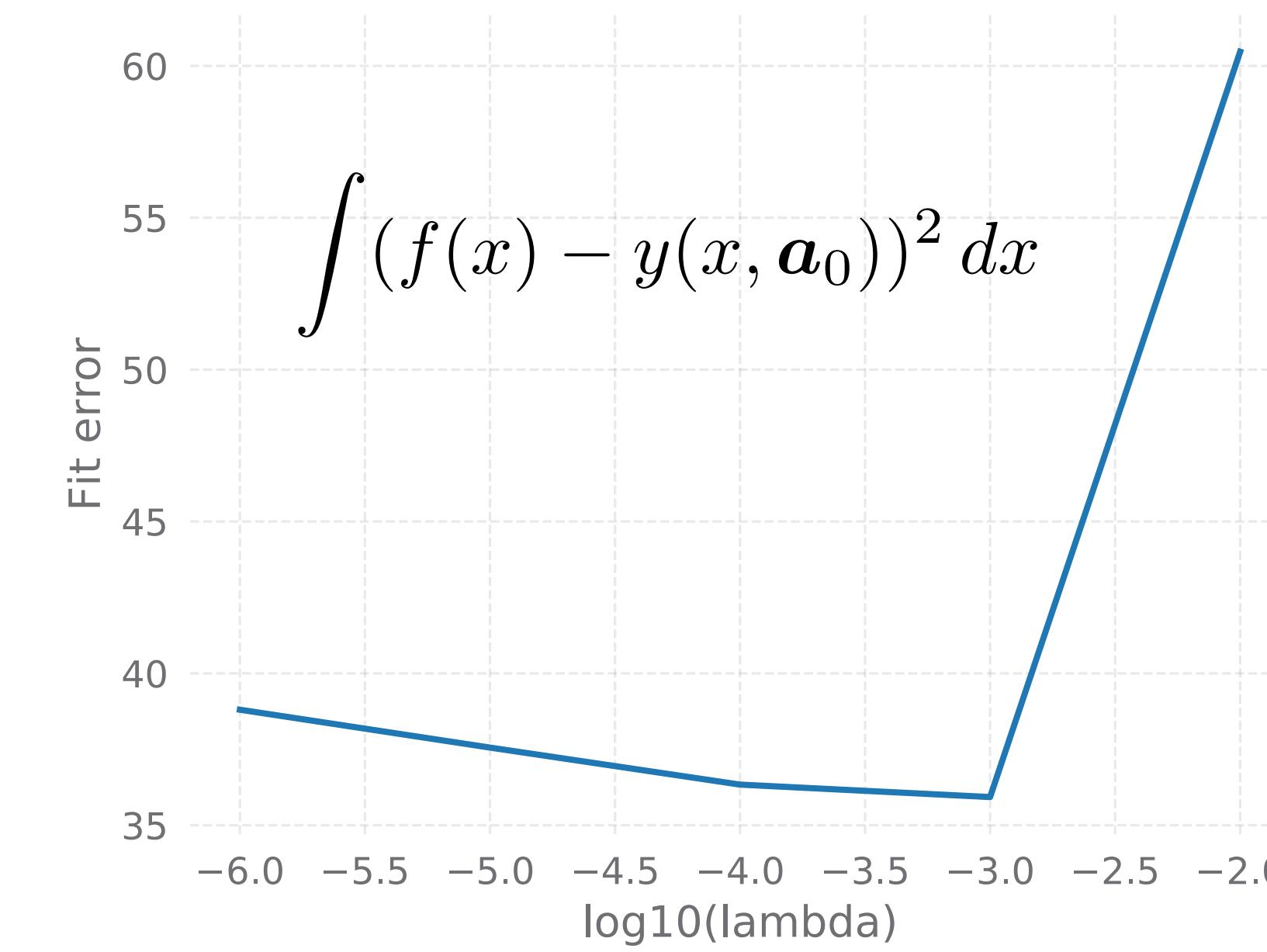
Exercise on kernel ridge regression - lessons learned

$$\ell = 0.5$$

$$\lambda' = 10^{-6} - 10^{-2}$$



Error as a function of regularization



Really nice and stable fits!

Today

- 13:00 Gaussian processes - theory
- 13:30 Exercises on GP
- 15:00 Hyperparameters - theory
- 15:30 Exercises on hyperparameters
-

Gaussian processes - the real thing!

- Includes the linear basis functions models, but are more general
- Can reproduce the results of kernel ridge regression
- Includes full probability distributions, so can predict uncertainties

Gaussian processes - derivation from linear basis function models

Work with parameters \mathbf{w}

Basis functions: $\phi_i(\mathbf{x})$

$$y(\mathbf{x}, \mathbf{w}) = \phi(\mathbf{x})^T \mathbf{w}$$

Prior distribution: $\mathcal{P}(\mathbf{w}) \propto \exp(-\lambda \mathbf{w}^T \mathbf{w}/2)$

Average: $\langle \mathbf{w} \rangle = 0$

Covariance: $\text{Cov}[\mathbf{w}, \mathbf{w}^T] = \frac{1}{\lambda} \mathbf{I}_M$

Prior distribution:

$$\mathcal{P}(\mathbf{w}) = \mathcal{N}(\mathbf{w} | 0, \frac{1}{\lambda} \mathbf{I}_M) \propto \exp(-\lambda \mathbf{w}^T \mathbf{w}/2)$$

Work directly with function $y(\mathbf{x}, \mathbf{w}) = \phi(\mathbf{x})^T \mathbf{w}$

Kernel: $k(\mathbf{x}, \mathbf{x}') = \frac{1}{\lambda} \sum_{i=1}^M \phi_i(\mathbf{x}) \phi_i(\mathbf{x}')$

Prior distribution: $\mathcal{P}([y(\mathbf{x})])?$

Average: $\langle y(\mathbf{x}) \rangle = \phi(\mathbf{x})^T \langle \mathbf{w} \rangle = 0$

Covariance:

$$\begin{aligned} \text{Cov}[y(\mathbf{x}), y(\mathbf{x}')] &= \langle y(\mathbf{x}) y(\mathbf{x}') \rangle = \phi(\mathbf{x})^T \langle \mathbf{w} \mathbf{w}^T \rangle \phi(\mathbf{x}') \\ &= \frac{1}{\lambda} \phi(\mathbf{x})^T \mathbf{I}_M \phi(\mathbf{x}') = \frac{1}{\lambda} \sum_{i=1}^M \phi_i(\mathbf{x}) \phi_i(\mathbf{x}') = k(\mathbf{x}, \mathbf{x}') \end{aligned}$$

Prior distribution:

$$\mathcal{P}([y(\mathbf{x})]) = \mathcal{N}([y(\mathbf{x})] | 0, k(\mathbf{x}, \mathbf{x}'))$$

Gaussian process - the prior distribution

Prior distribution:

Kernel:

$$k(\mathbf{x}, \mathbf{x}') = \frac{1}{\lambda} \sum_{i=1}^M \phi_i(\mathbf{x}) \phi_i(\mathbf{x}')$$

but now we will forget about the basis functions!

$$\mathcal{P}([y(\mathbf{x})]) = \mathcal{N}([y(\mathbf{x})] | 0, k(\mathbf{x}, \mathbf{x}'))$$

$$\propto \exp \left(-\frac{1}{2} \int y(\mathbf{x}) k^{-1}(\mathbf{x}, \mathbf{x}') y(\mathbf{x}') d\mathbf{x} d\mathbf{x}' \right)$$

“Inverse” function defined like matrix inverse

$$\int k^{-1}(\mathbf{x}, \mathbf{x}'') k(\mathbf{x}'', \mathbf{x}') d\mathbf{x}'' = \delta(\mathbf{x} - \mathbf{x}')$$

We shall now look at the Gaussian (“squared-exponential”) kernel:

$$k(\mathbf{x}, \mathbf{x}') = \exp(-(x - x')^2 / 2\ell^2)$$

Interpretation of kernel as covariance matrix!

$y(\mathbf{x})$ and $y(\mathbf{x}')$ are *correlated* if \mathbf{x} and \mathbf{x}' are close and *uncorrelated* or *independent* if they are far away from each other.

$$\text{Cov}[y(\mathbf{x}), y(\mathbf{x}')] = \langle y(\mathbf{x}) y(\mathbf{x}') \rangle = k(\mathbf{x}, \mathbf{x}')$$

(Do the exercise on the prior distribution with Gaussian kernel function.)

Prior distribution:

$$\mathcal{P}([y(\mathbf{x})]) = \mathcal{N}([y(\mathbf{x})] | 0, k(\mathbf{x}, \mathbf{x}'))$$

$$\propto \exp\left(-\frac{1}{2} \int y(\mathbf{x}) k^{-1}(\mathbf{x}, \mathbf{x}') y(\mathbf{x}') d\mathbf{x} d\mathbf{x}'\right)$$

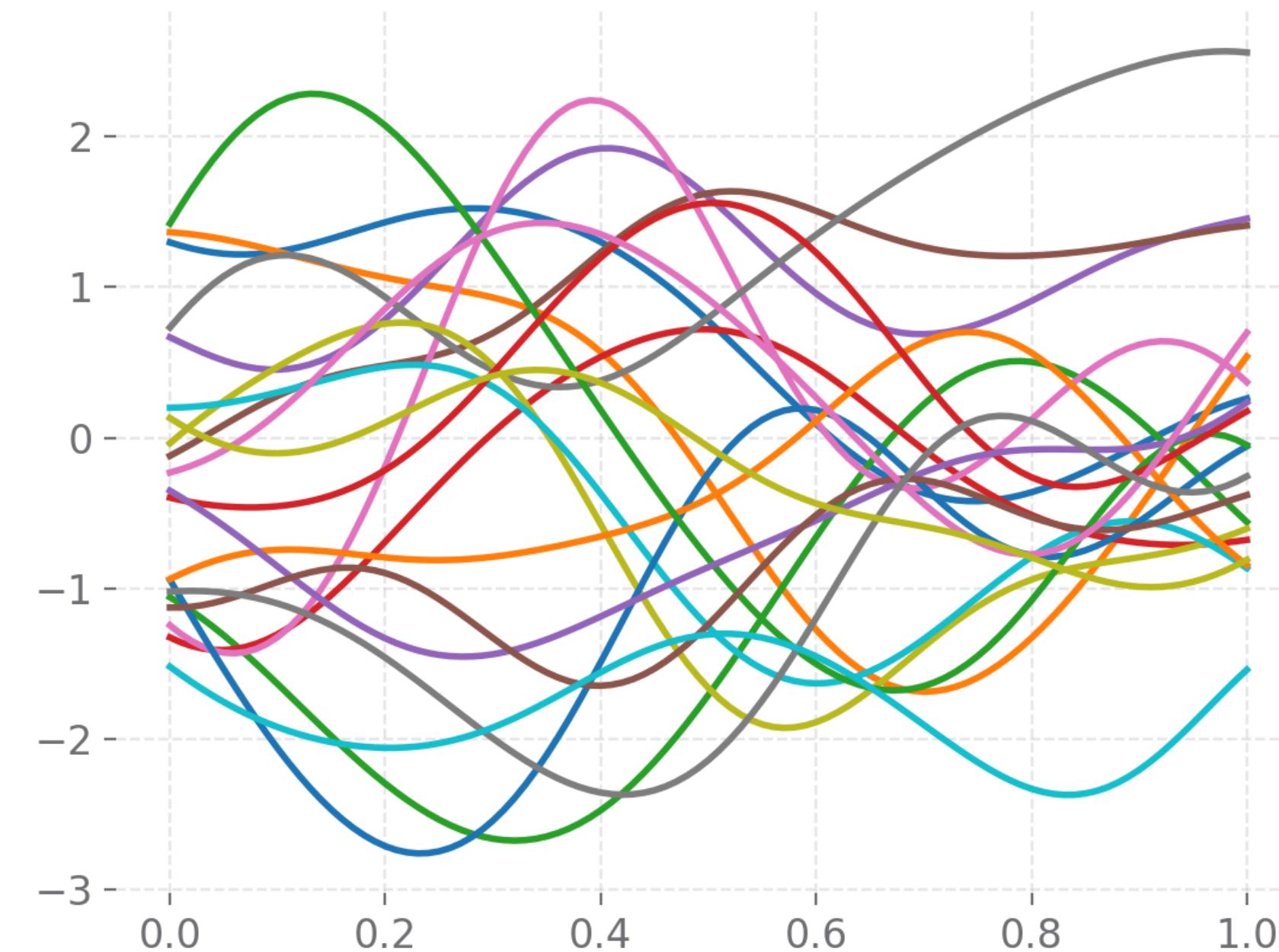
Kernel: $k(\mathbf{x}, \mathbf{x}') = \exp(-(x - x')^2 / 2\ell^2)$

$$\text{Cov}[y(\mathbf{x}), y(\mathbf{x}')] = \langle y(\mathbf{x})y(\mathbf{x}') \rangle = k(\mathbf{x}, \mathbf{x}')$$

```
xg = np.linspace(0,1,101)
ll = 0.2
def kernel(x,xx):
    return np.exp(-(x-xx)**2/(2*ll**2))

Nens = 20
k = np.array([[kernel(x,xx) for x in xg] for xx in xg])

for i in range(Nens):
    y = np.random.multivariate_normal(np.zeros(ngrid),k)
    plt.plot(xg,y)
```

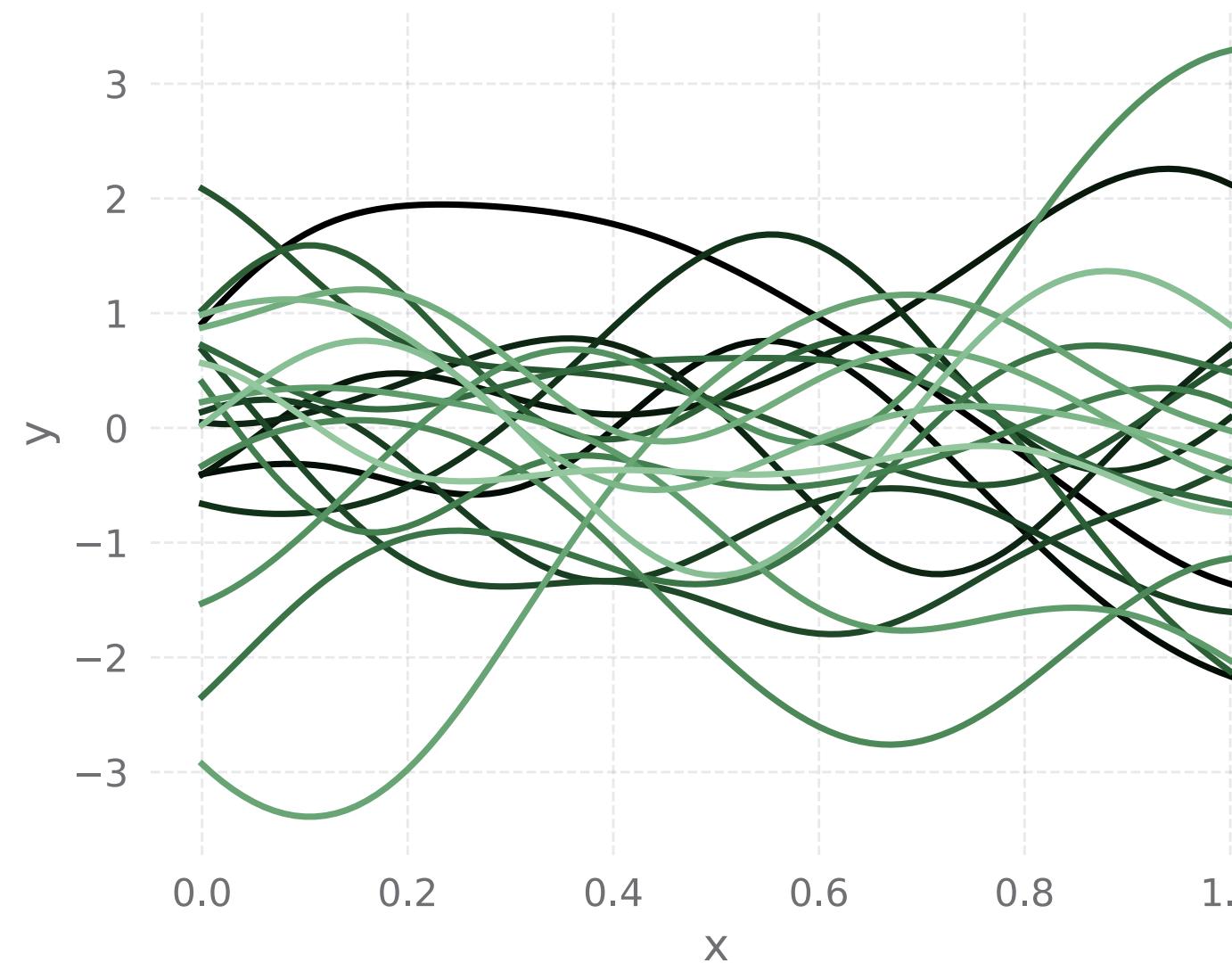


Do the exercise on the prior distribution with Gaussian kernel function.

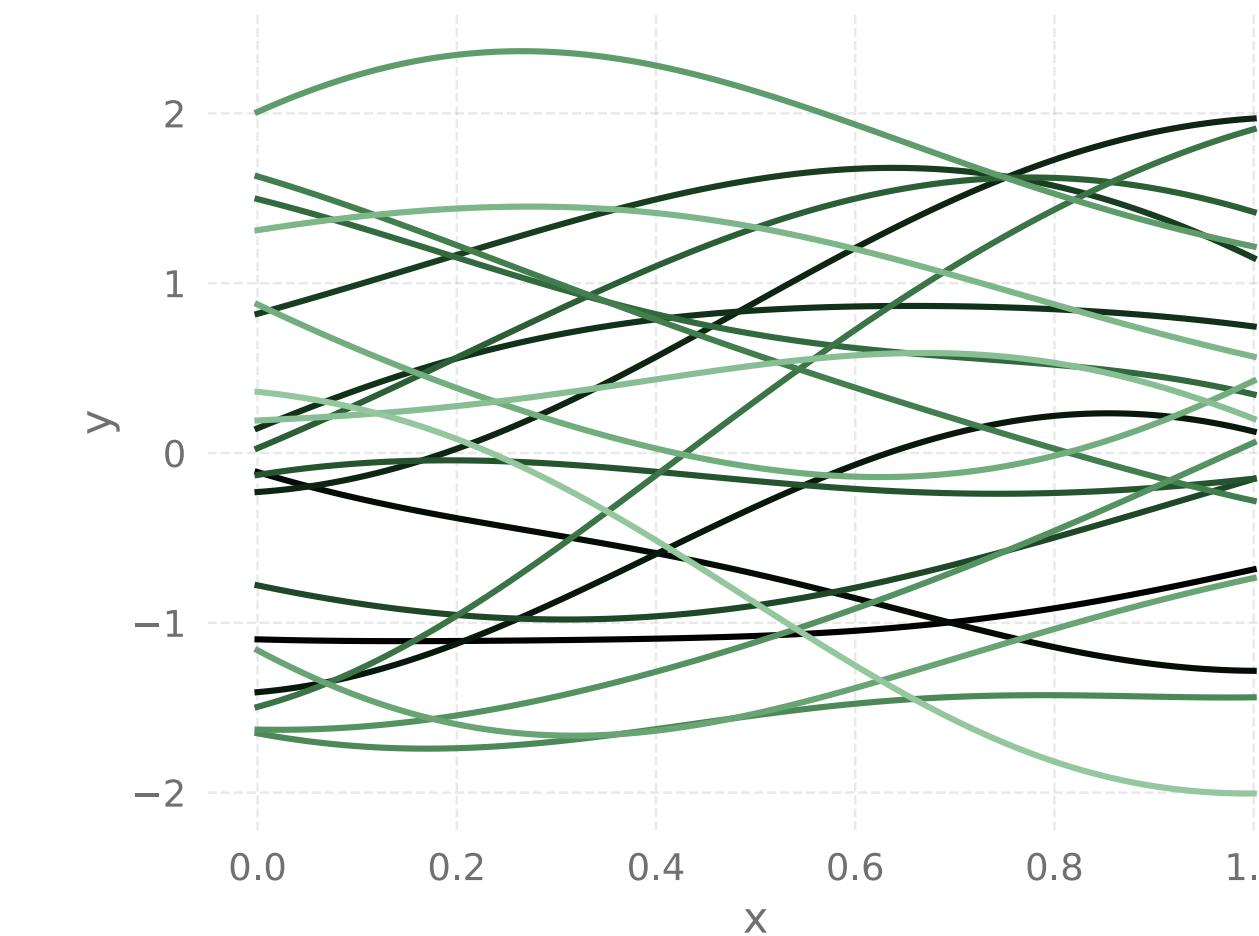
Gaussian process - the prior distribution - lessons learned

$$k(x, x') = \exp(-(x - x')^2 / 2\ell^2)$$

Ensemble of functions with $\ell = 0.2$



Ensemble of functions with $\ell = 0.6$



The length scale determines our expectations about the functions.
Much nicer functions than the ones we obtained with the polynomial basis.

Gaussian process - prediction

From the linear basis function models, we have the prediction

$$y_0(\mathbf{x}) := \langle y(\mathbf{x}) \rangle = \boldsymbol{\phi}(\mathbf{x})^T \langle \mathbf{w} \rangle = \boldsymbol{\phi}(\mathbf{x})^T \mathbf{w}_0 = \boldsymbol{\phi}(\mathbf{x})^T (\boldsymbol{\Phi}^T \boldsymbol{\Phi} + \tilde{\lambda} \mathbf{I}_M)^{-1} \boldsymbol{\Phi}^T \mathbf{t}$$

↑
basis functions ↑
design matrix ↑
data

This can be rewritten (see the notes)

$$y_0(\mathbf{x}) = \mathbf{k}(\mathbf{x})^T \mathbf{C}^{-1} \mathbf{t}$$

$K_{nm} = k(\mathbf{x}_n, \mathbf{x}_m)$

↑
 $k_n(\mathbf{x}) = k(\mathbf{x}, \mathbf{x}_n)$ $\mathbf{C} = \mathbf{K} + \sigma^2 \mathbf{I}_N$

We do not need the basis functions. The prediction is expressed using the kernel function. Prediction looks exactly like kernel ridge regression, if the function g is identified with the kernel.

Gaussian process - variance

From the linear basis function models, we have

$$\sigma^2(\mathbf{x}) := \langle (y(\mathbf{x}, \mathbf{w}) - y_0(\mathbf{x}))^2 \rangle = \boldsymbol{\phi}^T(\mathbf{x}) \sigma^2 \left(\boldsymbol{\Phi}^T \boldsymbol{\Phi} + \tilde{\lambda} \mathbf{I}_M \right)^{-1} \boldsymbol{\phi}(\mathbf{x})$$

This can be rewritten (see the notes)

$$\sigma^2(\mathbf{x}) = \langle (y(\mathbf{x}, \mathbf{w}) - y_0(\mathbf{x}))^2 \rangle = k(\mathbf{x}, \mathbf{x}) - \mathbf{k}(\mathbf{x})^T \mathbf{C}^{-1} \mathbf{k}(\mathbf{x})$$

$$k_n(\mathbf{x}) = k(\mathbf{x}, \mathbf{x}_n)$$

$$K_{nm} = k(\mathbf{x}_n, \mathbf{x}_m)$$

$$\mathbf{C} = \mathbf{K} + \sigma^2 \mathbf{I}_N$$

Also the variance is expressed exclusively using the kernel!

Gaussian process

Gaussian process. Given the kernel function $k(\mathbf{x}, \mathbf{x}')$ we can define the vector function $\mathbf{k}(\mathbf{x})$ by $k_n(\mathbf{x}) = k(\mathbf{x}, \mathbf{x}_n)$ and the kernel matrix \mathbf{K} by $K_{nm} = k(\mathbf{x}_n, \mathbf{x}_n)$. We furthermore define the matrix $\mathbf{C} = \mathbf{K} + \sigma^2 \mathbf{I}_N$. The prior average function is denoted $y_p(\mathbf{x})$ with the corresponding vector $(\mathbf{y}_p)_n = y_p(\mathbf{x}_n)$. With these definitions we get the prediction for the average function $y_0(\mathbf{x})$ and the variance $\sigma^2(\mathbf{x})$ as

$$y_0(\mathbf{x}) = y_p(\mathbf{x}) + \mathbf{k}(\mathbf{x})^T \mathbf{C}^{-1} (\mathbf{t} - \mathbf{y}_p) \quad (17.16)$$

$$\sigma^2(\mathbf{x}) = k(\mathbf{x}, \mathbf{x}) - \mathbf{k}(\mathbf{x})^T \mathbf{C}^{-1} \mathbf{k}(\mathbf{x}) \quad (17.17)$$

What happens, if we try to make a prediction far away from the data points?

Gaussian process

Kernel: $k(x, x') = k_0 \exp(-(x - x')^2 / 2\ell^2)$

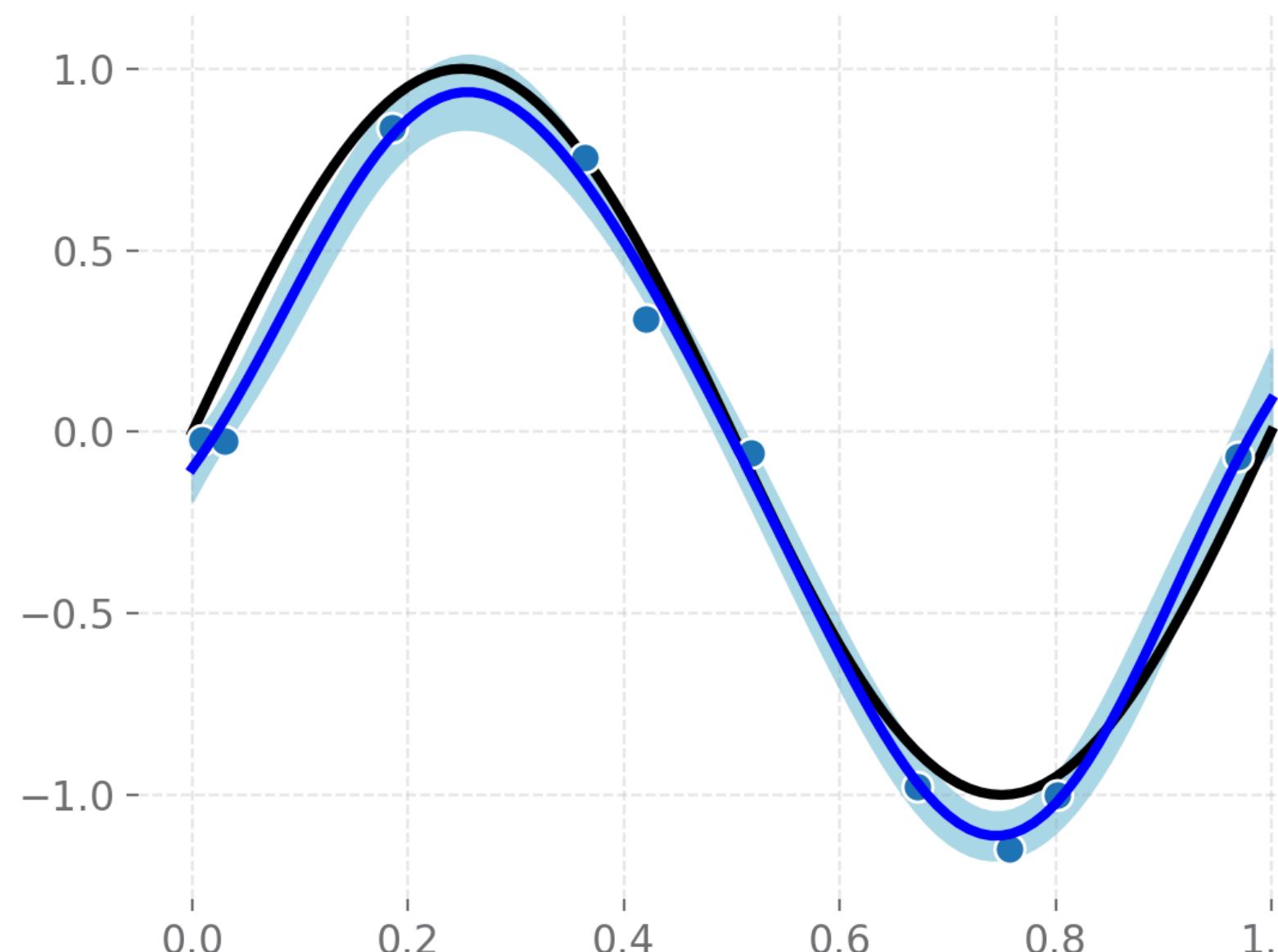
$$k_n(\mathbf{x}) = k(\mathbf{x}, \mathbf{x}_n)$$

$$K_{nm} = k(\mathbf{x}_n, \mathbf{x}_m)$$

$$\mathbf{C} = \mathbf{K} + \sigma^2 \mathbf{I}_N$$

Prediction: $y_0(\mathbf{x}) = \mathbf{k}(\mathbf{x})^T \mathbf{C}^{-1} \mathbf{t}$

Variance: $\sigma^2(\mathbf{x}) = \langle (y(\mathbf{x}, \mathbf{w}) - y_0(\mathbf{x}))^2 \rangle$
 $= k(\mathbf{x}, \mathbf{x}) - \mathbf{k}(\mathbf{x})^T \mathbf{C}^{-1} \mathbf{k}(\mathbf{x})$



```

initwithnoise()
# The x-values of the data are in the array "xp"
# and the t-values are in the array tp

ngrid = 101
xg = np.linspace(0,1,ngrid)
sigma = 0.1

ll = 0.2
k0 = 1

def kernel(x,xx):
    return k0*np.exp(-(x-xx)**2/(2*ll**2))

def kvec(x):
    return np.array([kernel(x,xx) for xx in xp])

K = np.array([[kernel(x,xx) for x in xp] for xx in xp])
C = K + sigma**2*np.identity(N)
Cinv = np.linalg.inv(C)
Cinvt = np.dot(Cinv,tp)

def fitf(x):
    return np.dot(kvec(x),Cinvt)

def varx(x):
    return kernel(x,x)- np.dot(kvec(x),np.dot(Cinv,kvec(x)))

plt.plot(xg,np.sin(2*np.pi*xg), color='k', lw=3) # Plot sine function
plt.plot(xp,tp,'o') # Plot data points

fitg = np.array([fitf(x) for x in xg]) # Fit function on the grid
width = np.array([np.sqrt(varx(x)) for x in xg]) # Width on the grid
plt.plot(xg, fitg, color='b', lw=3)
plt.fill_between(xg, fitg-width, fitg+width, color='lightblue')

```

Now do the exercise on the Gaussian process for the sine function.

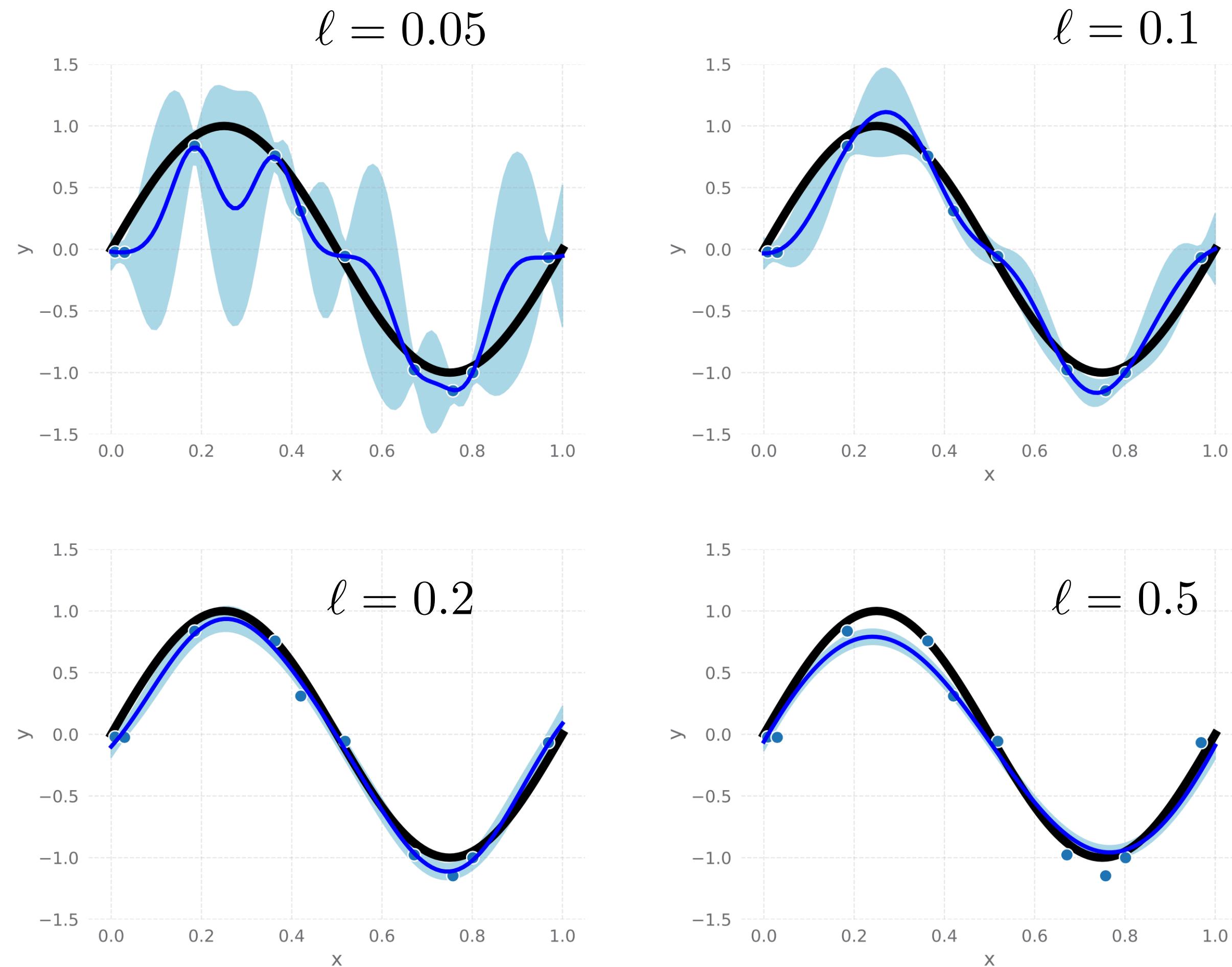
Gaussian process for the sine function - lessons learned

Noise $\sigma = 0.1$

Kernel $k(x, x') = k_0 \exp(-(x - x')^2 / 2\ell^2)$

$$k_0 = 1 = \langle y(x)^2 \rangle_{\text{prior}}$$

Too small ℓ : poor interpolation between points, large variance
Too large ℓ : model too “stiff”, underestimation of errors



Hyperparameters

$$k(x, x') = k_0 \exp(-(x - x')^2 / 2\ell^2)$$

Examples: noise σ , length scale ℓ , prefactor k_0

$$\mathbf{C} = \mathbf{K} + \sigma^2 \mathbf{I}_N$$

Different strategies to determine the hyperparameters

$$k_0 = 1 = \langle y(x)^2 \rangle_{\text{prior}}$$

1. Prior information

- Noise might be known
- The scale ℓ of the problem might be known
- Prefactor may come from expectations for the variance $k_0 = \langle y(x)^2 \rangle_{\text{prior}}$

2. Cross validation

- Only use part of the data for “training” and validate on the rest
- 10-fold cross validation: divide data in 10 sets, train on 9 and test on the last one
- Pick the values of the hyperparameters that do the best in the cross validation

3. Bayesian analysis (yes!)

Bayesian analysis of hyperparameters for a Gaussian process

Hyperparameters θ

Bayes' theorem:

$$\begin{aligned} \mathcal{P}(\theta|t) &\propto \mathcal{P}(t|\theta)\mathcal{P}(\theta) = \int \mathcal{P}(t, y|\theta) dy \mathcal{P}(\theta) \\ &= \int \mathcal{P}(t|y)\mathcal{P}(y|\theta) dy \mathcal{P}(\theta) \\ &= \int \mathcal{N}(t|y, \sigma^2 I_N) \mathcal{N}(y|y_p, K) dy \mathcal{P}(\theta) \\ &= \int \mathcal{N}(t|y_p, C) dy \mathcal{P}(\theta) \end{aligned}$$

$C = K + \sigma^2 I_N$

marginalization over y

Too heavy, to work with the full posterior distribution, so we maximize the log-likelihood:

$$\begin{aligned} \log \mathcal{P}(t|\theta) &= \log \mathcal{N}(t|y_p, C) \\ &= -\frac{1}{2} \log (\det(C)) - \frac{1}{2} (t - y_p)^T C^{-1} (t - y_p) - \frac{N}{2} \log(2\pi) \end{aligned}$$

$$\mathcal{N}(t|y_p, C) = \frac{1}{(2\pi)^{N/2}} \frac{1}{\det(C)^{1/2}} \exp \left(-\frac{1}{2} (t - y_p)^T C^{-1} (t - y_p) \right)$$

Bayesian analysis of hyperparameters for a Gaussian process

$$\mathcal{N}(\mathbf{t}|\mathbf{y}_p, \mathbf{C}) = \frac{1}{(2\pi)^{N/2}} \frac{1}{\det(\mathbf{C})^{1/2}} \exp\left(-\frac{1}{2}(\mathbf{t} - \mathbf{y}_p)^T \mathbf{C}^{-1} (\mathbf{t} - \mathbf{y}_p)\right)$$

$$\begin{aligned}\log \mathcal{P}(\mathbf{t}|\boldsymbol{\theta}) &= \log \mathcal{N}(\mathbf{t}|\mathbf{y}_p, \mathbf{C}) \\ &= -\frac{1}{2} \log (\det(\mathbf{C})) - \frac{1}{2} (\mathbf{t} - \mathbf{y}_p)^T \mathbf{C}^{-1} (\mathbf{t} - \mathbf{y}_p) - \frac{N}{2} \log(2\pi)\end{aligned}$$

Maximize by setting derivative to zero:

$$\frac{\partial}{\partial \theta_i} \log \mathcal{P}(\mathbf{t}|\boldsymbol{\theta}) = -\frac{1}{2} \text{Tr}\left(\mathbf{C}^{-1} \frac{\partial \mathbf{C}}{\partial \theta_i}\right) + \frac{1}{2} (\mathbf{t} - \mathbf{y}_p)^T \mathbf{C}^{-1} \frac{\partial \mathbf{C}}{\partial \theta_i} \mathbf{C}^{-1} (\mathbf{t} - \mathbf{y}_p) = 0 \quad \text{consider fixed}$$

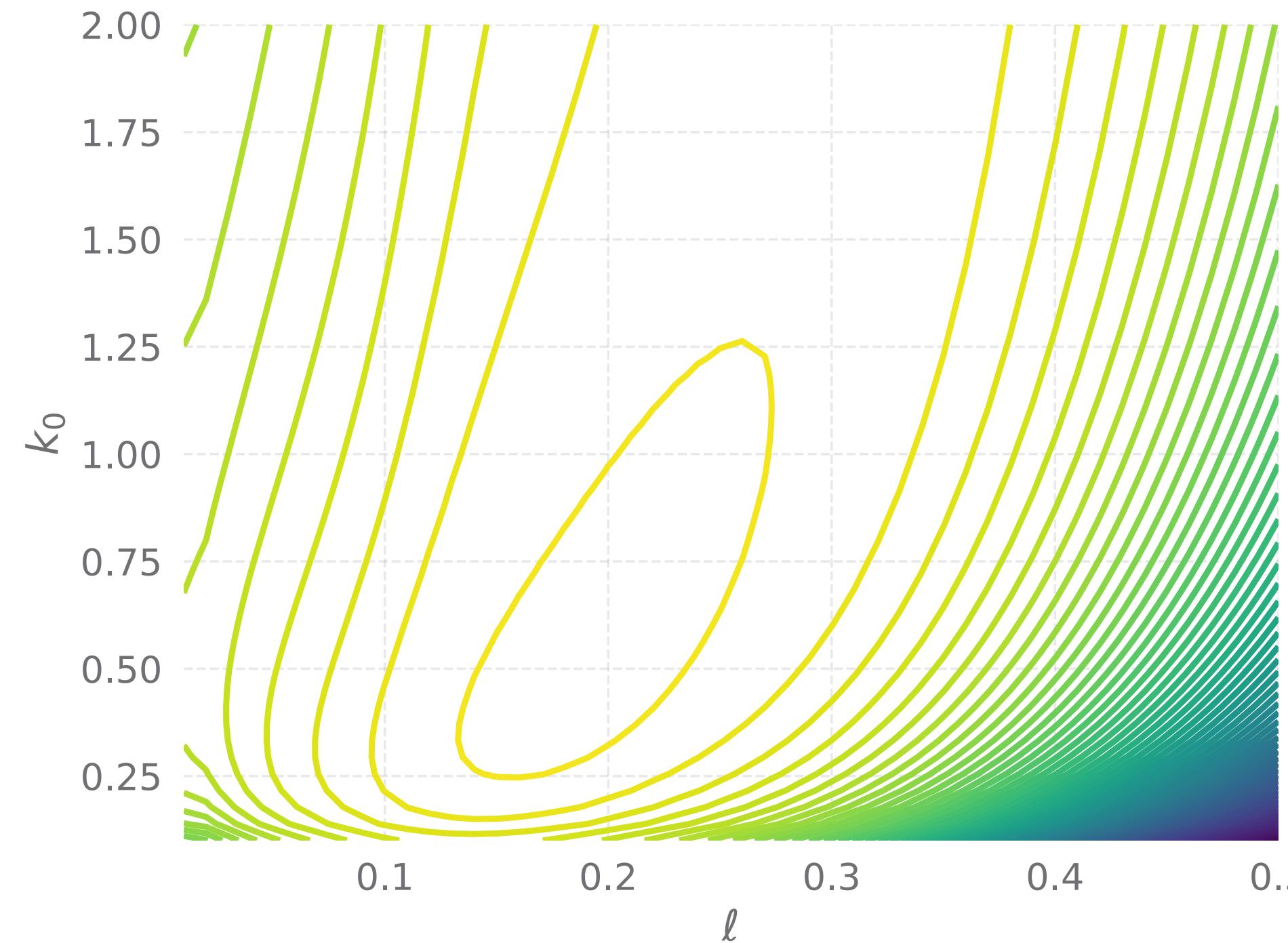
Example: kernel prefactor

$$\mathbf{C} = \mathbf{K} + \sigma^2 \mathbf{I}_N = k_0 \mathbf{K}_0 + \sigma^2 \mathbf{I}_N = k_0 \left(\mathbf{K}_0 + \frac{\sigma^2}{k_0} \mathbf{I}_N \right) = k_0 \mathbf{C}_0$$

Maximizing the likelihood gives: $k_0 = \frac{1}{N} (\mathbf{t} - \mathbf{y}_p)^T \mathbf{C}_0^{-1} (\mathbf{t} - \mathbf{y}_p)$

Do the exercise on the Bayesian determination of hyperparameters.

Optimizing hyperparameters for sine problem



Too large ℓ requires larger k_0 .