

Machine Learning for Time Series

(MLTS or MLTS-Deluxe Lectures)

Dr. Dario Zanca

Machine Learning and Data Analytics (MaD) Lab
Friedrich-Alexander-Universität Erlangen-Nürnberg
25.10.2022

Organisational Information

Lectures (online)

A new lecture recording every Tuesday.

Consultation hours from November 15th, h. 8:15 – 9:30

Exercises (online)

Live Zoom Session starting on November 9th

Recordings will be uploaded

Project (online)

Introduction during the first exercise Live Zoom Session (November 9th)

Applications started on Nov 9th

-
- Time series fundamentals and definitions (2 lectures)
 - Bayesian Inference (1 lecture) ←
 - Gaussian processes (2 lectures)
 - State space models (2 lectures)
 - Autoregressive models (1 lecture)
 - Data mining on time series (1 lecture)
 - Deep learning on time series (4 lectures)
 - Domain adaptation (1 lecture)

In this lecture...

1. Bayes Theorem
2. Bayesian Model Selection
3. Prior Distributions
4. Linear Regression (Bayesian treatment)



Bayesian Inference

Bayes' Theorem



The **Bayes' Theorem** was formulated by the English philosopher **Thomas Bayes** (1701 – 1761), whose notes were edited and published posthumously by Richard Price.

Bayes' theorem is stated mathematically as the following equation:

$$P(A|B) = \frac{P(B|A)P(A)}{P(B)}$$

where A and B are events, and $P(B) \neq 0$.

Portrait from: Terence O'Donnell, *History of Life Insurance in Its Formative Years* (Chicago: American Conservation Co., 1936), p. 335





Posterior probability

The probability of event A occurring given that B is true

Likelihood

The probability of event B occurring given that A is true

$$P(A|B) = \frac{P(B|A) P(A)}{P(B)}$$

Marginal probability

The probability of observing B without any given conditions

Prior probability

The probability of observing A without any given conditions

Portrait from: Terence O'Donnell, *History of Life Insurance in Its Formative Years* (Chicago: American Conservation Co., 1936), p. 335

Bayes' Theorem

An example. Iterative application of the Bayes' theorem.

We know that:

- Disease chance: 1%
- Test accuracy: 95%

A: Having the disease

B: Testing positive to the disease

$P(B|A)$: Test accuracy (Likelihood)

$P(B)$: Prob. of a positive test (Marginal)

$P(A)$: Disease chance (Prior)

First positive test

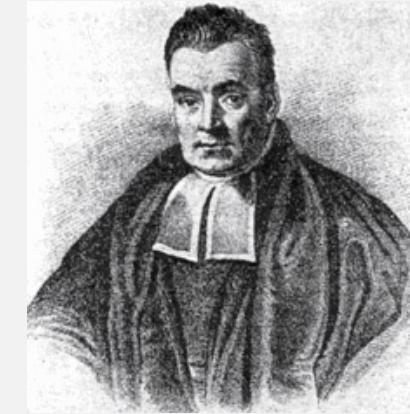
$$P(A|B) = \frac{.95 \times 0.01}{0.95 \cdot 0.01 + (1 - 0.95)(1 - 0.01)} = 0.161$$

Second positive test

$$P(A|B) = \frac{.95 \times 0.161}{0.95 \cdot 0.161 + (1 - 0.95)(1 - 0.161)} = 0.785$$

Third positive test

$$P(A|B) = \frac{.95 \times 0.785}{0.95 \cdot 0.785 + (1 - 0.95)(1 - 0.785)} = 0.986$$



$$P(A|B) = \frac{P(B|A) P(A)}{P(B)}$$

Portrait from: Terence O'Donnell, *History of Life Insurance in Its Formative Years* (Chicago: American Conservation Co., 1936), p. 335

Let \mathcal{D} denote the **observed data**,

$$\mathcal{D} = \{x^{(n)}, y^{(n)}\}$$

with $x^{(n)} \in \mathcal{R}$ represents the input, and $y^{(n)} \in \mathcal{R}$ represents the output (labels).

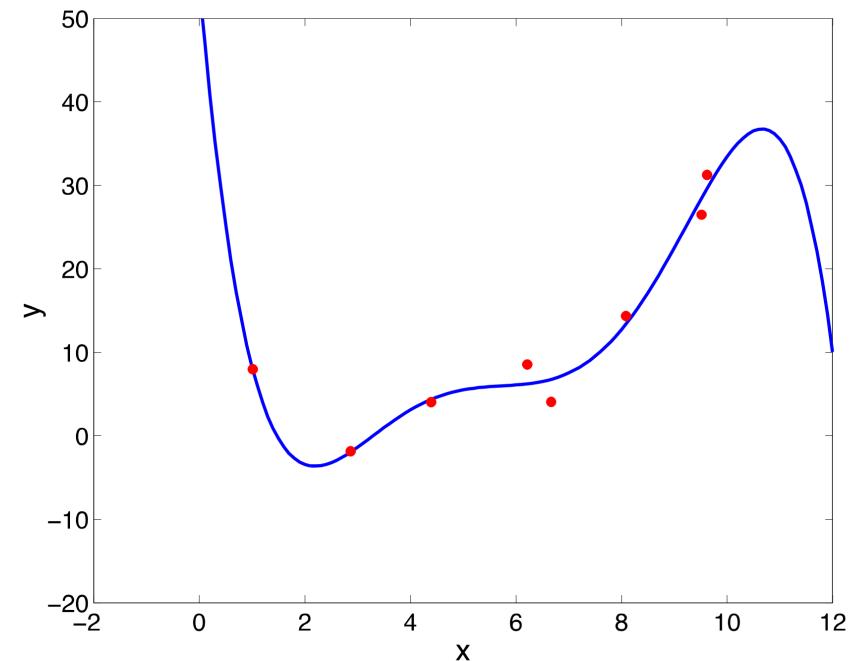
The **model** is defined as

$$y^{(n)} = \omega_0 + \omega_1 x^{(n)} + \omega_2 x^{(n)} \dots + \omega_m x^{(n)} + \epsilon$$

where data noise is Gaussian distributed, i.e., $\epsilon \sim \mathcal{N}(0, \sigma^2)$.

We denote with θ the **unknown parameters**,

$$\theta = (\omega_0, \dots, \omega_m, \sigma)$$

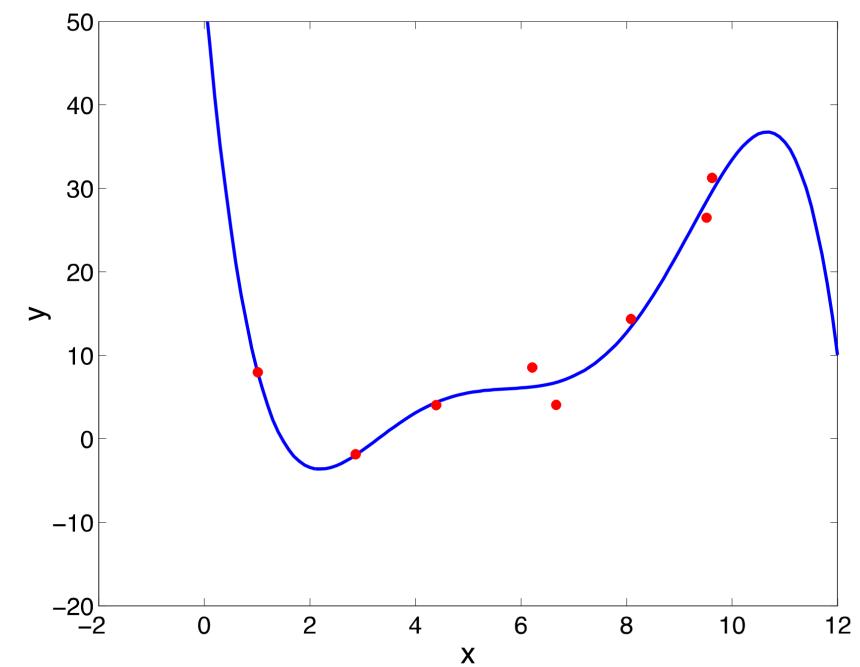


Data: $\mathcal{D} = \{x^{(n)}, y^{(n)}\}$

Model: $y^{(n)} = \omega_0 + \omega_1 x^{(n)} + \omega_2 x^{(n)} \dots + \omega_m x^{(n)} + \epsilon$

Unknown parameters: $\theta = (\omega_0, \dots, \omega_m, \sigma)$

Goal: To infer θ from the data and to predict future outputs $p(y|x, \theta, \mathcal{D})$



$p(\mathcal{D}|\theta)$: likelihood of θ

$p(\theta)$: prior probability of θ

$p(\theta|\mathcal{D})$: posterior of θ given \mathcal{D}

$p(\mathcal{D})$: marginal probability of \mathcal{D}

$p(y|x, \mathcal{D})$: predictive distribution

$p(\mathcal{D}|\theta)$: likelihood of θ

$p(\theta)$: prior probability of θ

$p(\theta|\mathcal{D})$: posterior of θ , given \mathcal{D}

$p(\mathcal{D})$: marginal probability of \mathcal{D}

$p(y|x, \mathcal{D})$: predictive distribution

Bayes' rule:

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

Prediction of a new point:

$$p(y|x, \mathcal{D}) = \int p(y|\theta, x, \mathcal{D})p(\theta|\mathcal{D}) d\theta$$

- In contrast to the maximum likelihood estimation (MLE), in Bayesian learning **we average over possible parameter settings** rather than optimizing over parameter space.
- Bayesian inference gives us a **systematic way to express our uncertainty** about future predictions. Prediction is not just a point estimate (as for MLE) but has a probability form that expresses the uncertainty about the predictions.



Bayesian Inference

Bayesian Model Selection





The **principle of Occam's razor** in its original formulation states that:

"**Entia non sunt multiplicanda praeter necessitatem**"

(In English, "Entities should not be multiplied unnecessarily")

Many scientists have adopted or reformulated the Occam's Razor principle, which is often cited in stronger forms, as in the following statement:

- "If you have two theories that both explain the observed facts, then you should use the simplest until more evidence comes along"
- "One should pick the simplest model that adequately explains the data"



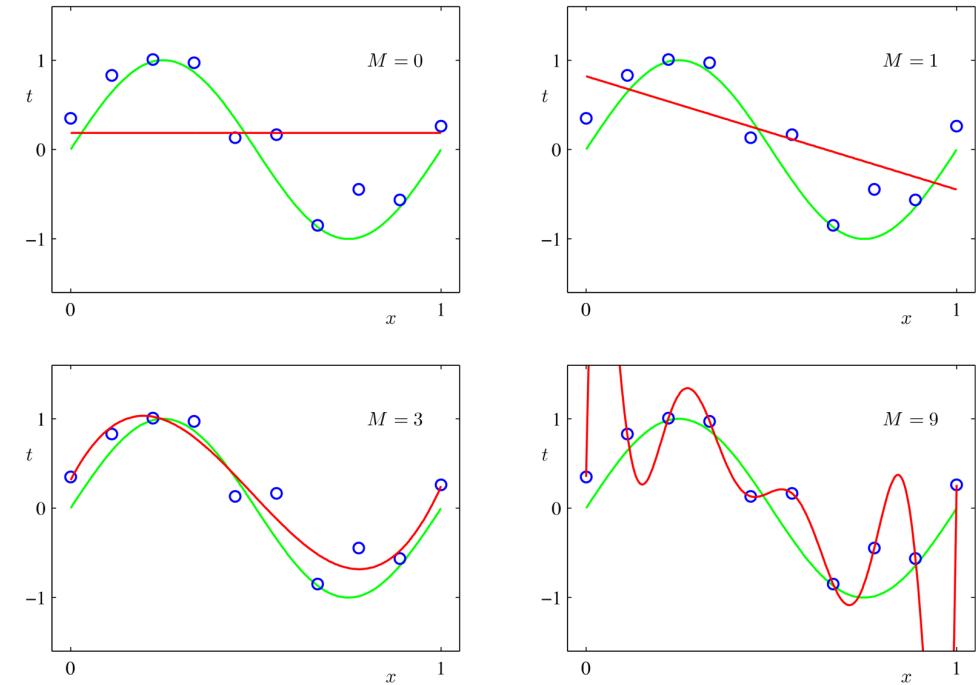
Picture from:
<https://www.britannica.com/topic/Occams-razor>

The model Selection problem

We could perform K-fold cross-validation (CV) to estimate the generalization error of all candidates model.

However, it requires fitting each candidate model K times!

→ A more efficient approach is given by Bayesian modelling



Which of the above models represents data the best?

We can compare different models using the marginal likelihood:

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

- Model classes that are **too simple** are unlikely to generate the data set.
- Model classes that are **too complex** can generate many possible data sets, so again, they are unlikely to generate that particular data set \mathcal{D} .

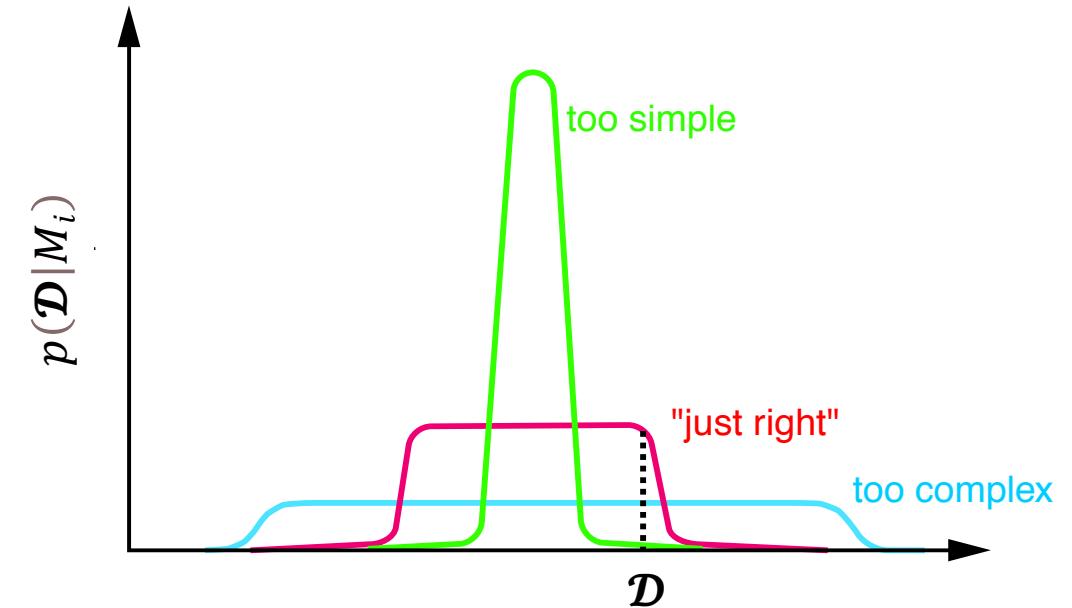


Image from: Rasmussen, C., & Ghahramani, Z. (2000). Occam's razor. Advances in neural information processing systems, 13.

To understand the Bayesian Occam's razor, we notice that:

$$\int_{\mathcal{D}} p(\mathcal{D}|M_i) = 1$$

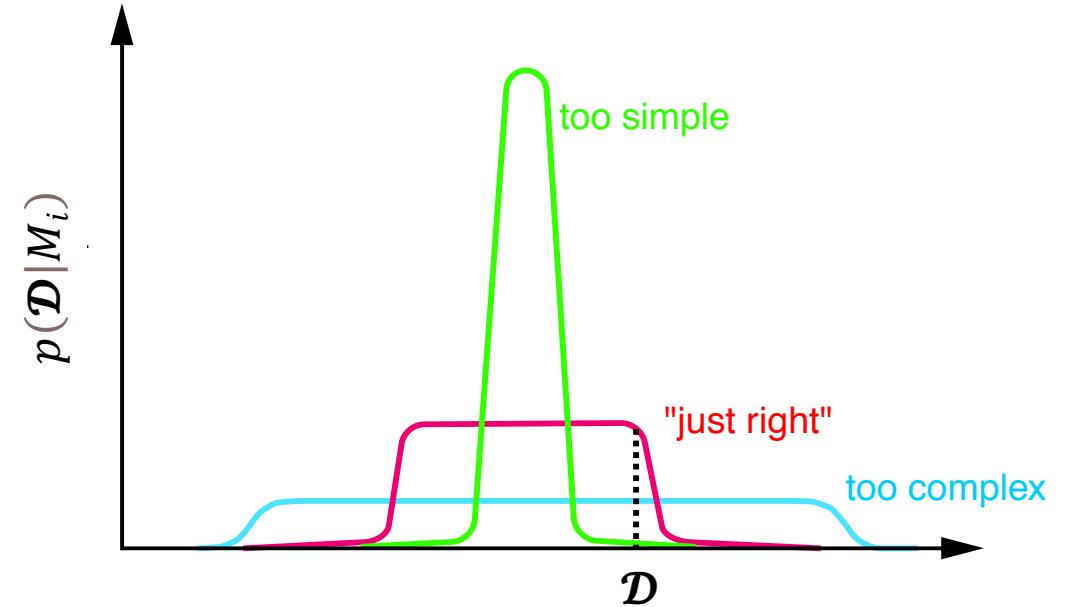


Image from: Rasmussen, C., & Ghahramani, Z. (2000). Occam's razor. Advances in neural information processing systems, 13.

To understand the Bayesian Occam's razor, we notice that:

$$\int_{\mathcal{D}} p(\mathcal{D}|M_i) = 1$$

Intuitively, complex models which can predict many datasets, must spread their probability mass
→ They don't attribute large probability for any given data set as simpler models.

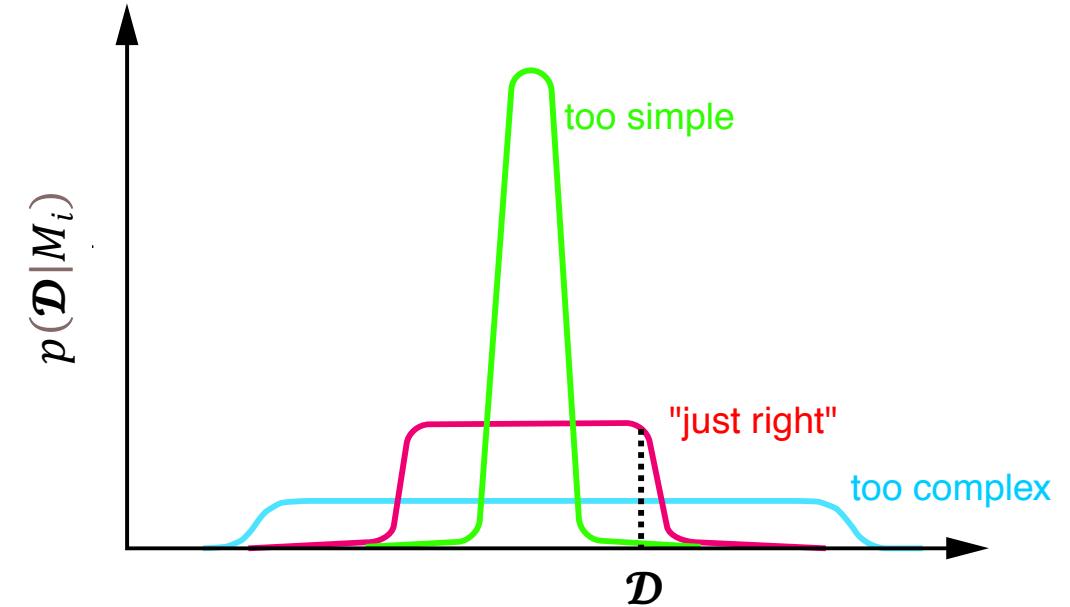


Image from: Rasmussen, C., & Ghahramani, Z. (2000). Occam's razor. Advances in neural information processing systems, 13.

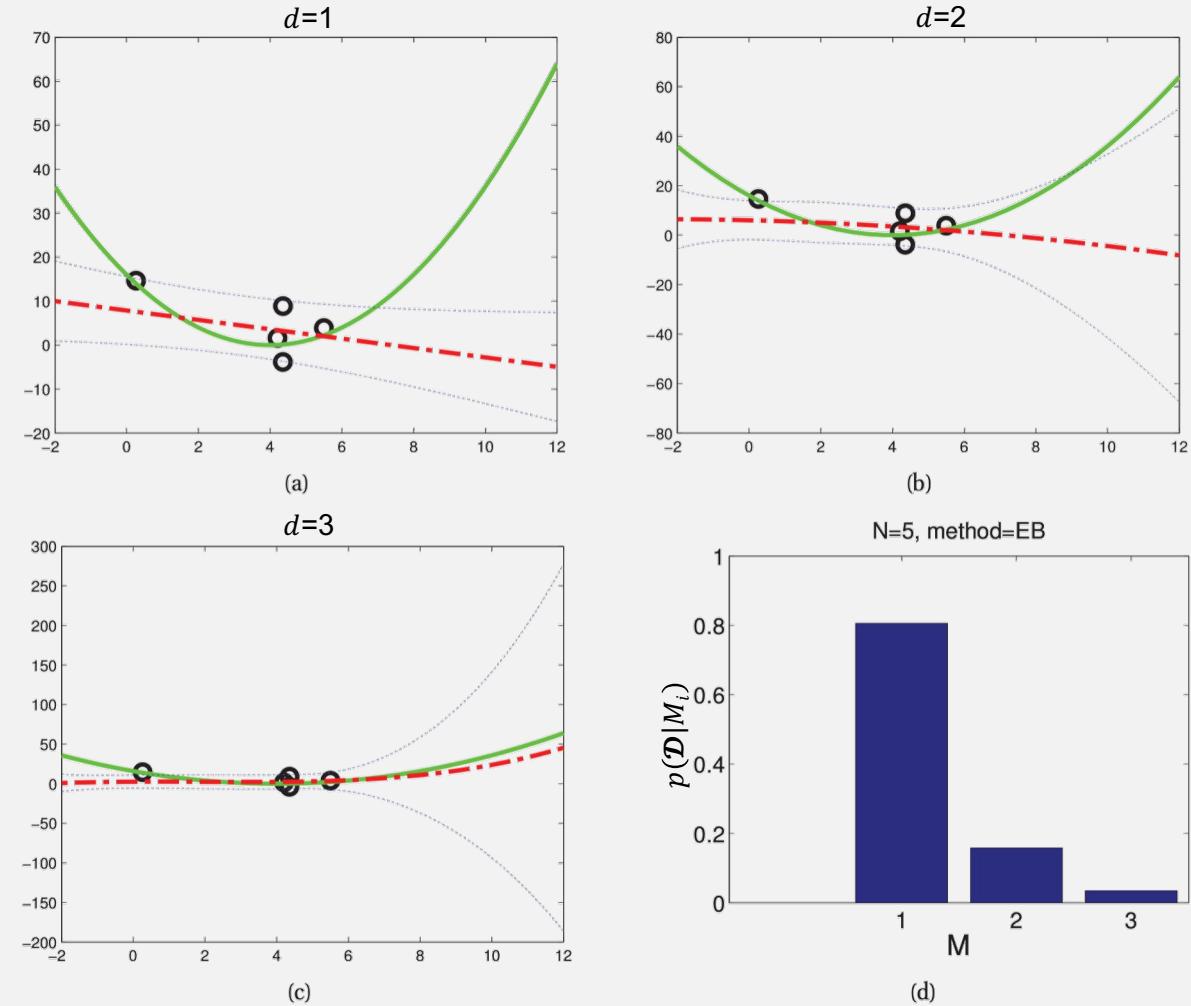
Bayesian Occam's razor

A concrete example

We plot polynomials of degrees 1, 2 and 3 fit to N=5 data points using (empirical) Bayes.

- True function
- - - Prediction
- $\pm\sigma$ around the mean

There is not enough data to justify a complex model, so the best model is $d = 1$.



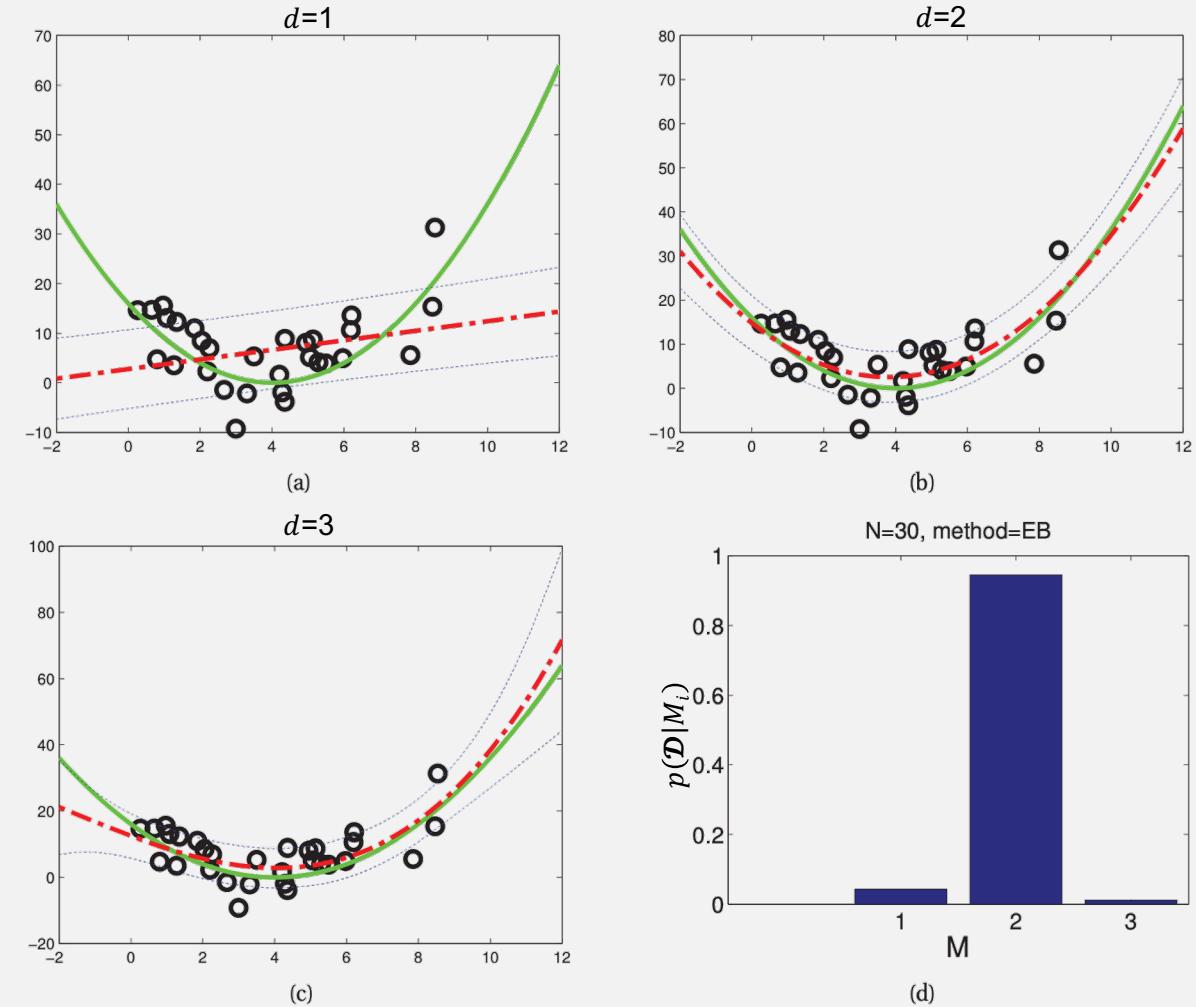
Bayesian Occam's razor

A concrete example

We plot polynomials of degrees 1, 2 and 3 fit to N=30 data points using (empirical) Bayes.

- True function
- - - Prediction
- $\pm\sigma$ around the mean

When more data is available, $d = 2$ is the right model





Bayesian Inference

Prior Distributions







p: Prob. purring

1-p: Prob. grumpy

What is the best guess
for the probability p?



How can I update my
belief on p ?

Prior distributions

The importance of priors in Bayesian Inference

$p(\mathcal{D}|\theta)$: likelihood of θ

$p(\theta)$: prior probability of θ

$p(\theta|\mathcal{D})$: posterior of θ , given \mathcal{D}

$p(\mathcal{D})$: marginal probability of \mathcal{D}

$p(y|x, \mathcal{D})$: predictive distribution

Bayes' rule:

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

$p(\mathcal{D}|\theta)$: likelihood of θ

$p(\theta)$: prior probability of θ

$p(\theta|\mathcal{D})$: posterior of θ , given \mathcal{D}

$p(\mathcal{D})$: marginal probability of \mathcal{D}

$p(y|x, \mathcal{D})$: predictive distribution

Bayes' rule:

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

A **prior probability distribution** of an uncertain quantity is the probability distribution that would express one's belief, before some evidence is taken into account.

→ For example, a prior could represent the distribution of votes coming from an opinion poll, prior to the election.

A **subjective prior** expresses the modeler's subjective belief.

- We formulate our (subjective) assumptions about modeling the data in terms of priors
- We have to work hard to understand the system under study in order to formulate our assumptions

An **objective prior** constrain prior beliefs to be “uninformative” about the parameters.

- The objective Bayes view is that formulating our assumptions is too difficult, especially in complex models

If we don't have strong beliefs about what θ should be, it is common to use an “uninformative” priors → “**Let the data speak for itself!**”

An **informative prior** expresses a specific information about a variable.

- For example, a reasonable informative prior about the temperature at noon tomorrow could be given by a normal distribution with expected value equal to today's noon temperature and variance equal to the daily variance of the temperature.

An **uninformative prior** is designed to express vague or general information about a variable.

- For example, when tossing a coin, we assign the probability of 0.5 to both heads and tails.

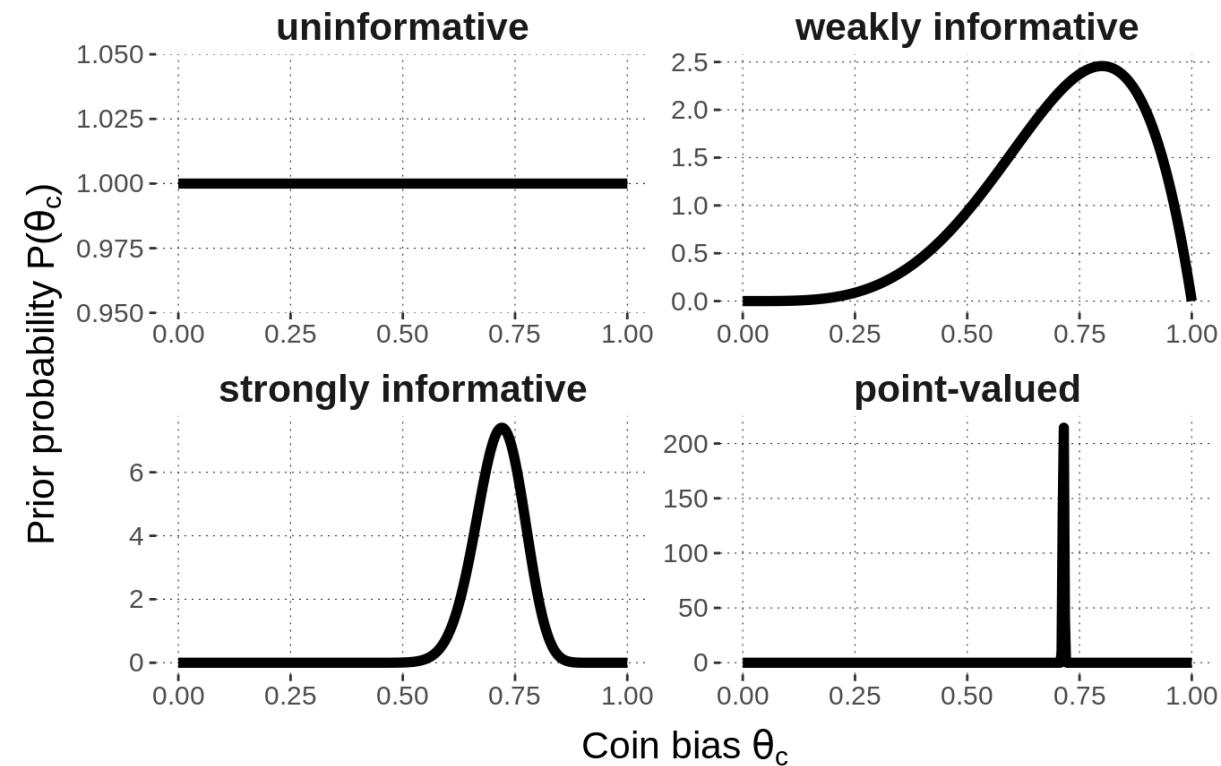


Image from: <https://michael-franke.github.io/intro-data-analysis/Chap-03-03-models-parameters-priors.html>

A key requirement for computing the posterior $p(\theta|\mathcal{D})$ is the specification of a prior $p(\theta|\alpha)$, where alpha denotes a set of hyperparameters.

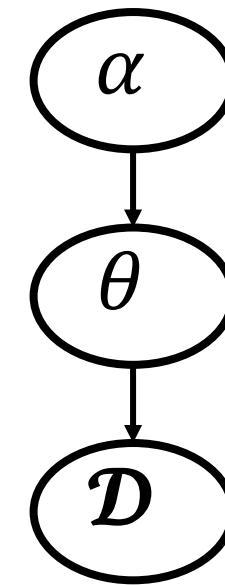
We have multiple levels of priors:

$$\alpha \rightarrow \theta \rightarrow \mathcal{D}$$

$$p(\theta) = \int p(\theta|\alpha) p(\alpha) d\alpha$$

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

The multi-level model



Consider the problem of predicting the cancer mortality rates in various cities. We measure the number of people N_i in various cities, as well as the number of people who died of cancer x_i in those cities. We assume that the mortality follows:

$$x_i \sim Bin(N_i, \theta_i).$$

A reasonable approach to estimate θ_i is that of assuming that they are drawn from the same distribution $\theta_i \sim Beta(a, b)$, where $\alpha = (a, b)$ are hyper-parameters in our model.

Then, the full joint distribution is written as

$$p(\mathcal{D}, \theta, \alpha | N) = p(\alpha) \prod_{i=1}^N Bin(x_i | N_i, \theta_i) Beta(\theta_i | \alpha).$$

Note: It is crucial to infer α from the data itself.

Example taken from: Machine Learning: A Probabilistic Perspective, Ch. 5.51

In hierarchical models, we need to **compute the posterior on multiple levels of variables**. For example,

$$p(\alpha, \theta | \mathcal{D}) \propto p(\mathcal{D} | \theta) p(\theta | \alpha) p(\alpha).$$

In some case, we can simplify the problems by **marginalizing over θ** . Then, we just need to compute:

$$p(\alpha | D).$$

As a computational “shortcut”, we can approximate the posterior on the hyper-parameters α with a point estimate. Since α is generally of a much smaller dimensionality than θ , it is less prone to overfitting and we safely assume a uniform prior on α .

$$\hat{\alpha} = \operatorname{argmax}_{\alpha} p(\mathcal{D}, \alpha) = \operatorname{argmax}_{\alpha} \int p(\mathcal{D}, \theta) p(\theta | \alpha) d\theta$$

A prior $p(\theta)$ is a **conjugate prior** for a particular likelihood $p(y|\theta)$ if the resulting posterior $p(\theta|y)$ has the same algebraic form.

Conjugate priors are widely used because they provide advantages:

- they usually allow us to derive a closed-form expression for the posterior distribution;
- they are easy to interpret,

Note: Conjugate priors simplify the computation, but are often not flexible enough to encode our prior knowledge → We can also use mixture of conjugate priors.

Likelihood (Binomial): $p(\mathcal{D}|\theta) = Bin(\theta) = \binom{n}{x} \theta^x (1-\theta)^{n-x}$

Prior (Beta): $p(\theta) = Beta(\alpha, \beta) = \frac{1}{B(\alpha, \beta)} \theta^{\alpha-1} (1-\theta)^{\beta-1}$

We plug them into the Bayes' formula to derive the posterior distribution:

$$\begin{aligned} p(\theta|\mathcal{D}) &= \frac{p(\mathcal{D}|\theta) p(\theta)}{\int p(\mathcal{D}|\theta) p(\theta) d\theta} \\ &= \frac{\binom{n}{x} \theta^x (1-\theta)^{n-x} \frac{1}{B(\alpha, \beta)} \theta^{\alpha-1} (1-\theta)^{\beta-1}}{\int \binom{n}{x} \theta^x (1-\theta)^{n-x} \frac{1}{B(\alpha, \beta)} \theta^{\alpha-1} (1-\theta)^{\beta-1} d\theta} \\ &= \frac{\frac{n! C_x}{B(\alpha, \beta)} \theta^{x+\alpha-1} (1-\theta)^{n-x+\beta-1}}{\frac{n! C_x}{B(\alpha, \beta)} \int \theta^{x+\alpha-1} (1-\theta)^{n-x+\beta-1} d\theta} = Beta(x + \alpha, n - x + \beta) \end{aligned}$$

Prior: Beta(2, 2)



Prior: Beta(2, 2)

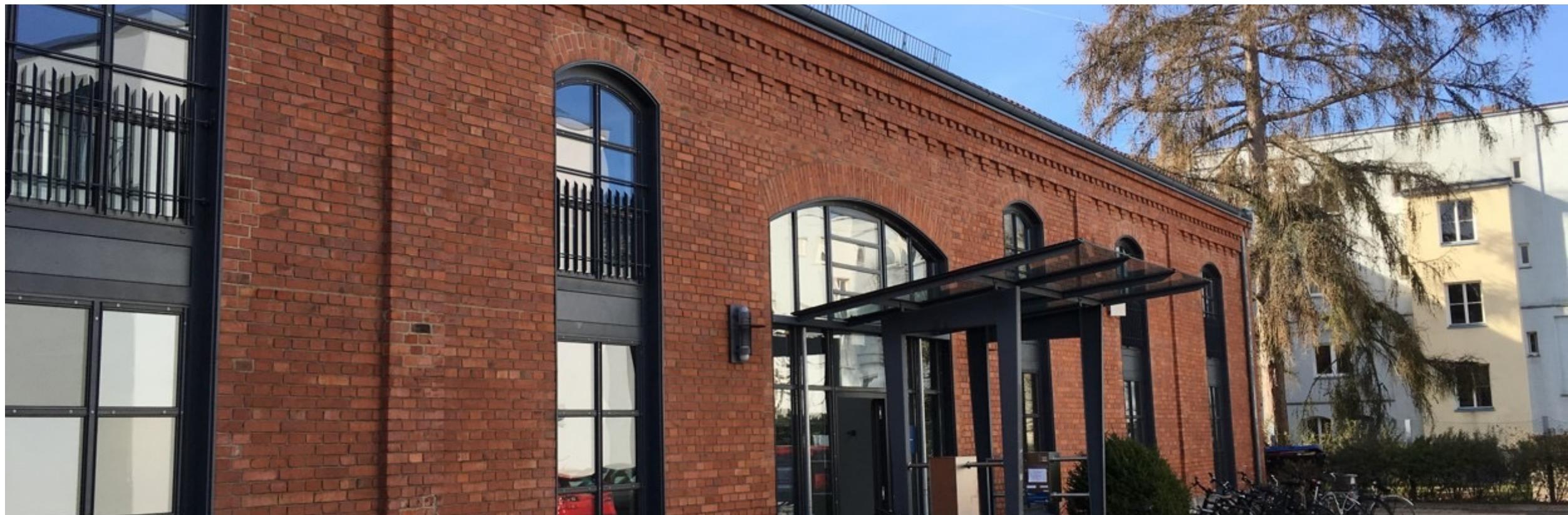


Posterior: Beta($2+2, 4+2$) = Beta(4, 6)

Prior: Beta(2, 2)



Posterior: Beta($2+2, 4+2$) = Beta(4, 6)



Bayesian Inference

Linear Regression (Bayesian treatment)



Given the observed data $\mathcal{D} = \{x^{(n)}, y^{(n)}\}$, we assume to know the noise variance σ^2 .

We would like to compute the posterior over the parameters, i.e,

$$p(w|\mathcal{D}, \sigma^2).$$

(We assume throughout a Gaussian likelihood model).

In linear regression **the likelihood is given by:**

$$\begin{aligned} p(y|X, w, \mu, \sigma^2) &= \mathcal{N}(y|\mu + Xw, \sigma^2 I_N) \\ &\propto \exp\left(-\frac{1}{2\sigma^2} (y - \mu - Xw)^T (y - \mu - Xw)\right) \end{aligned}$$

where μ is an offset term.

The conjugate prior of a Gaussian likelihood is also Gaussian*, which we will denote by

$$p(w) = \mathcal{N}(w|w_0, V_0).$$

Using the Bayes rule for Gaussian*, the posterior is given by

$$p(w|X, y, \sigma^2) \propto \mathcal{N}(w|w_0, V_0) \mathcal{N}(y|Xw, \sigma^2 I_N) = \mathcal{N}(w|w_N, V_N)$$

where

$$w_N = V_N V_0^{-1} w_0 + \frac{1}{\sigma^2} V_N X^T y$$

$$V_N = \sigma^2 (\sigma^2 V_0^{-1} + X^T X)^{-1}$$

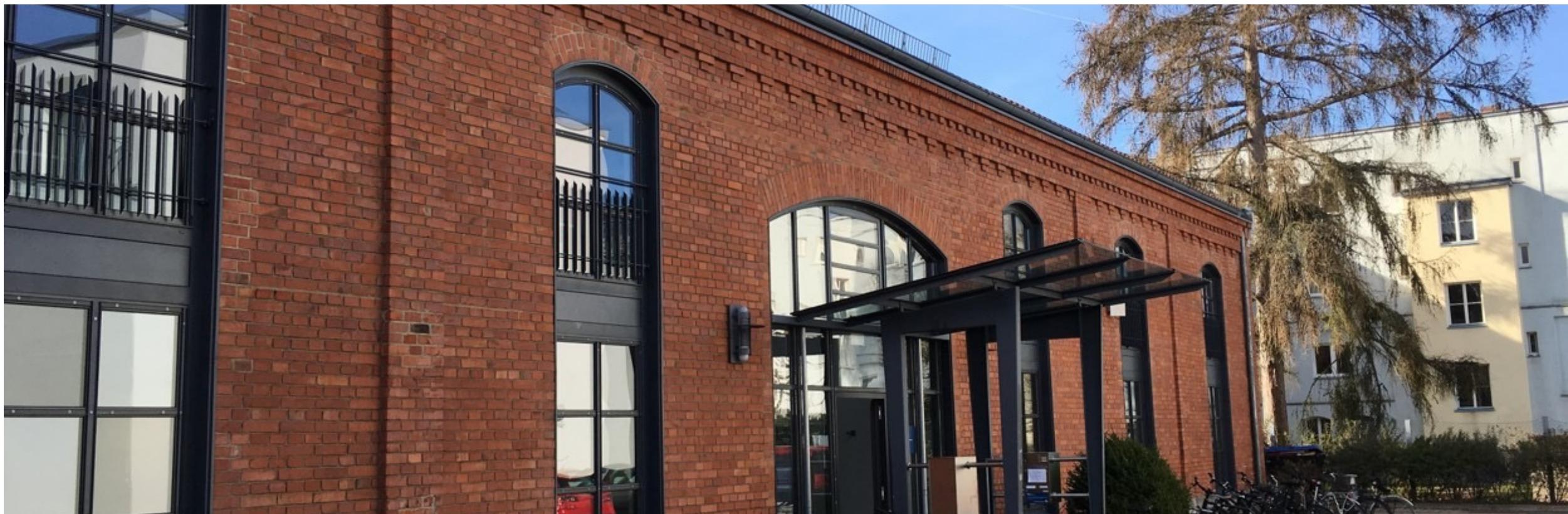
* See: Murphy K., „Machine Learning: A Probabilistic Perspective“ (2012)

The posterior predictive distribution at a test point x is given by *

$$\begin{aligned} p(y|x, \mathcal{D}, \sigma^2) &= \int \mathcal{N}(y|x^T w, \sigma^2) \mathcal{N}(w|w_N, V_N) dw \\ &= \mathcal{N}(y|w_N^T x, \sigma_N^2(x)) \end{aligned}$$

where $\sigma_N^2(x) = \sigma^2 + x^T V_N x$.

The variance in this prediction depends on the variance of the observation noise, σ^2 , and the variance in the parameters, V_N .



Bayesian Inference

Recap



- Bayesian modelling
 - Prior
 - Posterior
 - Likelihood
 - Priors
 - Informative vs Uninformative
 - Conjugate priors
 - Linear regression with Bayesian treatment
- Bayesian modelling requires integration over parameters
 - For complex models it could be not tractable! (We cannot compute the integral analytically)

