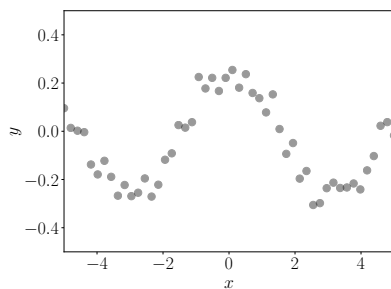


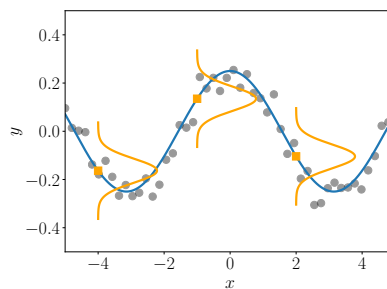
Linear Regression

In the following, we will apply the mathematical concepts from Chapters 2, 5, 6 and 7 to solving linear regression (curve fitting) problems. In *regression*, we want to find a function f that maps inputs $\mathbf{x} \in \mathbb{R}^D$ to corresponding function values $f(\mathbf{x}) \in \mathbb{R}$ given a set of training inputs \mathbf{x}_n and corresponding observations $y_n = f(\mathbf{x}_n) + \epsilon$, where ϵ is a random variable that comprises measurement noise and unmodeled processes. An illustration of such a regression problem is given in Figure 9.1. A typical regression problem is given in Figure 9.1(a): For some input values x we observe (noisy) function values $y = f(x) + \epsilon$. The task is to infer the function f that generated the data. A possible solution is given in Figure 9.1(b), where we also show three distributions centered at the function values $f(x)$ that represent the noise in the data.

Regression is a fundamental problem in machine learning, and regression problems appear in a diverse range of research areas and applications, including time-series analysis (e.g., system identification), control and robotics (e.g., reinforcement learning, forward/inverse model learning), optimization (e.g., line searches, global optimization), and deep-learning applications (e.g., computer games, speech-to-text translation, image recognition, automatic video annotation). Regression is also a key ingredient of classification algorithms.



(a) Regression problem: Observed noisy function values from which we wish to infer the underlying function that generated the data.



(b) Regression solution: Possible function that could have generated the data (blue) with indication of the measurement noise of the function value at the corresponding inputs (orange distributions).

regression

Figure 9.1
(a) Dataset;
(b) Possible solution to the regression problem.

5041 Finding a regression function requires solving a variety of problems,
5042 including

- 5043 • **Choice of the model (type) and the parametrization** of the regres-
5044 sion function. Given a data set, what function classes (e.g., polynomi-
5045 als) are good candidates for modeling the data, and what particular
5046 parametrization (e.g., degree of the polynomial) should we choose?
5047 Model selection, as discussed in Section 8.5, allows us to compare var-
5048 ious models to find the simplest model that explains the training data
5049 reasonably well.
- 5050 • **Finding good parameters.** Having chosen a model of the regression
5051 function, how do we find good model parameters? Here, we will need to
5052 look at different loss/objective functions (they determine what a “good”
5053 fit is) and optimization algorithms that allow us to minimize this loss.
- 5054 • **Overfitting and model selection.** Overfitting is a problem when the
5055 regression function fits the training data “too well” but does not gen-
5056 eralize to unseen test data. Overfitting typically occurs if the underly-
5057 ing model (or its parametrization) is overly flexible and expressive, see
5058 Section 8.5. We will look at the underlying reasons and discuss ways to
5059 mitigate the effect of overfitting in the context of linear regression.
- 5060 • **Relationship between loss functions and parameter priors.** Loss func-
5061 tions (optimization objectives) are often motivated and induced by prob-
5062 abilistic models. We will look at the connection between loss functions
5063 and the underlying prior assumptions that induce these losses.
- 5064 • **Uncertainty modeling.** In any practical setting, we have access to only
5065 a finite, potentially large, amount of (training) data for selecting the
5066 model class and the corresponding parameters. Given that this finite
5067 amount of training data does not cover all possible scenarios, we way
5068 want to describe the remaining parameter uncertainty to obtain a mea-
5069 sure of confidence of the model’s prediction at test time; the smaller the
5070 training set the more important uncertainty modeling. Consistent mod-
5071 eling of uncertainty equips model predictions with confidence bounds.

5072 In the following, we will be using the mathematical tools from Chap-
5073 ters 3, 5, 6 and 7 to solve linear regression problems. We will discuss
5074 maximum likelihood and maximum a posteriori (MAP) estimation to find
5075 optimal model parameters. Using these parameter estimates, we will have
5076 a brief look at generalization errors and overfitting. Toward the end of
5077 this chapter, we will discuss Bayesian linear regression, which allows us to
5078 reason about model parameters at a higher level, thereby removing some
5079 of the problems encountered in maximum likelihood and MAP estimation.

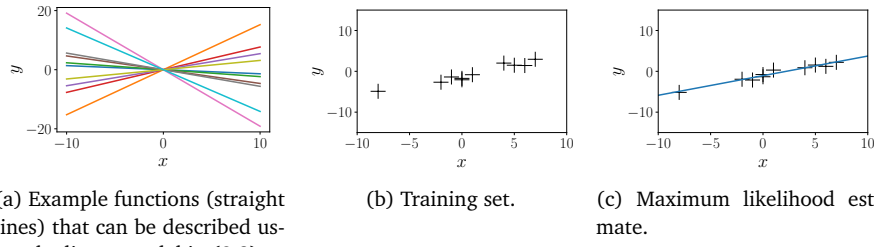


Figure 9.2 Linear regression without features.

(a) Example functions that fall into this category.
(b) Training set.
(c) Maximum likelihood estimate.

9.1 Problem Formulation

We consider the regression problem

$$y = f(\mathbf{x}) + \epsilon, \quad (9.1)$$

where $\mathbf{x} \in \mathbb{R}^D$ are inputs and $y \in \mathbb{R}$ are noisy function values (targets). Furthermore, $\epsilon \sim \mathcal{N}(0, \sigma^2)$ is independent, identically distributed (i.i.d.) measurement noise. In this particular case, ϵ is Gaussian distributed with mean 0 and variance σ^2 . Our objective is to find a function that is close (similar) to the unknown function that generated the data.

In this chapter, we focus on parametric models, i.e., we choose a parametrized function f and find parameters that “work well” for modeling the data. In linear regression, we consider the special case that the parameters appear linearly in our model. An example of linear regression is

$$y = f(\mathbf{x}) + \epsilon = \mathbf{x}^\top \boldsymbol{\theta} + \epsilon, \quad (9.2)$$

where $\boldsymbol{\theta} \in \mathbb{R}^D$ are the parameters we seek, and $\epsilon \sim \mathcal{N}(0, \sigma^2)$ is i.i.d. Gaussian measurement/observation noise. The class of functions described by (9.2) are straight lines that pass through the origin. In (9.2), we chose a parametrization $f(\mathbf{x}) = \mathbf{x}^\top \boldsymbol{\theta}$. For the time being we assume that the noise variance σ^2 is known. The noise model induces the *likelihood*

$$p(y | \mathbf{x}, \boldsymbol{\theta}) = \mathcal{N}(y | \mathbf{x}^\top \boldsymbol{\theta}, \sigma^2), \quad (9.3)$$

which is the probability of observing a target value y given that we know the input location \mathbf{x} and the parameters $\boldsymbol{\theta}$. Note that the only source of uncertainty originates from the observation noise (as \mathbf{x} and $\boldsymbol{\theta}$ are assumed known in (9.3))—without any observation noise, the relationship between \mathbf{x} and y would be deterministic and (9.3) would be a delta distribution.

For $x, \theta \in \mathbb{R}$ the linear regression model in (9.2) describes straight lines (linear functions), and the parameter θ would be the slope of the line. Figure 9.2(a) shows some examples. This model is not only linear in the parameters, but also linear in the inputs x . We will see later that $y = \phi(x)\theta$ for nonlinear transformations ϕ is also a linear regression model because “linear regression” refers to models that are “linear in the parameters”, i.e., models that describe a function by a linear combination of input features.

likelihood

Linear regression refers to models that are linear in the parameters.

In the following, we will discuss in more detail how to find good parameters θ and how to evaluate whether a parameter set “works well”.

9.2 Parameter Estimation

Consider the linear regression setting (9.2) and assume we are given a *training set* \mathcal{D} consisting of N inputs $\mathbf{x}_n \in \mathbb{R}^D$ and corresponding observations/targets $y_n \in \mathbb{R}$, $n = 1, \dots, N$. The corresponding graphical model is given in Figure 9.3. Note that y_i and y_j are conditionally independent given their respective inputs $\mathbf{x}_i, \mathbf{x}_j$, such that the likelihood function factorizes according to

$$p(y_1, \dots, y_N | \mathbf{x}_1, \dots, \mathbf{x}_N) = \prod_{n=1}^N p(y_n | \mathbf{x}_n) = \prod_{n=1}^N \mathcal{N}(y_n | \mathbf{x}_n^\top \theta, \sigma^2). \quad (9.4)$$

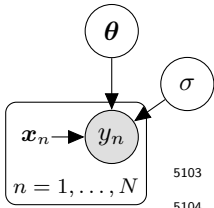
The likelihood and the factors $p(y_n | \mathbf{x}_n)$ are Gaussian due to the noise distribution.

In the following, we are interested in finding optimal parameters $\theta^* \in \mathbb{R}^D$ for the linear regression model (9.2). Once the parameters θ^* are found, we can predict function values by using this parameter estimate in (9.2) so that at an arbitrary test input \mathbf{x}_* we predict the probability for an output y_* as

$$p(y_* | \mathbf{x}_*, \theta^*) = \mathcal{N}(y_* | \mathbf{x}_*^\top \theta^*, \sigma^2). \quad (9.5)$$

In the following, we will have a look at parameter estimation by maximizing the likelihood, a topic that we already covered to some degree in Section 8.2.

Figure 9.3
Probabilistic graphical model for linear regression. Observed random variables are shaded, deterministic/known values are without circles. The parameters θ are treated as unknown/latent quantities.



9.2.1 Maximum Likelihood Estimation

A widely used approach to finding the desired parameters θ_{ML} is *maximum likelihood estimation* where we find parameters θ_{ML} that maximize the likelihood (9.4).

We obtain the maximum likelihood parameters as

$$\theta_{\text{ML}} = \arg \max_{\theta} p(\mathbf{y} | \mathbf{X}, \theta), \quad (9.6)$$

where we define the *design matrix* $\mathbf{X} := [\mathbf{x}_1, \dots, \mathbf{x}_N]^\top \in \mathbb{R}^{N \times D}$ and $\mathbf{y} := [y_1, \dots, y_N]^\top \in \mathbb{R}^N$ as the collections of training inputs and targets, respectively. Note that the n th row in the design matrix \mathbf{X} corresponds to the data point \mathbf{x}_n .

Remark. Note that the likelihood is not a probability distribution in θ : It is simply a function of the parameters θ but does not integrate to 1 (i.e., it is unnormalized), and may not even be integrable with respect to θ . However, the likelihood in (9.6) is a normalized probability distribution in the data \mathbf{y} . \diamond

ce the logarithm
(strictly)
monotonically
creasing function,
optimum of a
ction f is
ntical to the
imum of $\log f$.

To find the desired parameters θ_{ML} that maximize the likelihood, we typically perform gradient ascent (or gradient descent on the negative likelihood). In the case of linear regression we consider here, however, a closed-form solution exists, which makes iterative gradient descent unnecessary. In practice, instead of maximizing the likelihood directly, we apply the log-transformation to the likelihood function and minimize the negative log-likelihood.

Remark (Log Transformation). Since the likelihood function is a product of N Gaussian distributions, the log-transformation is useful since a) it does not suffer from numerical underflow, b) the differentiation rules will turn out simpler. Numerical underflow will be a problem when we multiply N probabilities, where N is the number of data points, since we cannot represent very small numbers, such as 10^{-256} . Furthermore, the log-transform will turn the product into a sum of log-probabilities such that the corresponding gradient is a sum of individual gradients, instead of a repeated application of the product rule (5.54) to compute the gradient of a product of N terms. \diamond

To find the optimal parameters θ_{ML} of our linear regression problem, we minimize the negative log-likelihood

$$-\log p(\mathbf{y} | \mathbf{X}, \theta) = -\log \prod_{n=1}^N p(y_n | \mathbf{x}_n, \theta) = -\sum_{n=1}^N \log p(y_n | \mathbf{x}_n, \theta), \quad (9.7)$$

where we exploited that the likelihood (9.4) factorizes over the number of data points due to our independence assumption on the training set.

In the linear regression model (9.2) the likelihood is Gaussian (due to the Gaussian additive noise term), such that we arrive at

$$\log p(y_n | \mathbf{x}_n, \theta) = -\frac{1}{2\sigma^2} (y_n - \mathbf{x}_n^\top \theta)^2 + \text{const} \quad (9.8)$$

where the constant includes all terms independent of θ . Using (9.8) in the negative log-likelihood (9.7) we obtain (ignoring the constant terms)

$$\mathcal{L}(\theta) := -\log p(\mathbf{y} | \mathbf{X}, \theta) = \frac{1}{2\sigma^2} \sum_{n=1}^N (y_n - \mathbf{x}_n^\top \theta)^2 \quad (9.9a)$$

$$= \frac{1}{2\sigma^2} (\mathbf{y} - \mathbf{X}\theta)^\top (\mathbf{y} - \mathbf{X}\theta) = \frac{1}{2\sigma^2} \|\mathbf{y} - \mathbf{X}\theta\|^2, \quad (9.9b)$$

where $\mathbf{X} = [\mathbf{x}_1, \dots, \mathbf{x}_N]^\top \in \mathbb{R}^{N \times D}$.

Remark. There is some notation overloading: We often summarize the set of training inputs in \mathbf{X} , whereas in the design matrix we additionally assume a specific “shape”. \diamond

In (9.9b) we used the fact that the sum of squared errors between the observations y_n and the corresponding model prediction $\mathbf{x}_n^\top \theta$ equals the

The negative log-likelihood function is also called *error function*.

squared distance between \mathbf{y} and $\mathbf{X}\boldsymbol{\theta}$. Remember from Section 3.1 that $\|\mathbf{x}\|^2 = \mathbf{x}^\top \mathbf{x}$ if we choose the dot product as the inner product.

With (9.9b) we have now a concrete form of the negative log-likelihood function we need to optimize. We immediately see that (9.9b) is quadratic in $\boldsymbol{\theta}$. This means that we can find a unique global solution $\boldsymbol{\theta}_{\text{ML}}$ for minimizing the negative log-likelihood \mathcal{L} . We can find the global optimum by computing the gradient of \mathcal{L} , setting it to $\mathbf{0}$ and solving for $\boldsymbol{\theta}$.

Using the results from Chapter 5, we compute the gradient of \mathcal{L} with respect to the parameters as

$$\frac{d\mathcal{L}}{d\boldsymbol{\theta}} = \frac{d}{d\boldsymbol{\theta}} \left(\frac{1}{2\sigma^2} (\mathbf{y} - \mathbf{X}\boldsymbol{\theta})^\top (\mathbf{y} - \mathbf{X}\boldsymbol{\theta}) \right) \quad (9.10a)$$

$$= \frac{1}{2\sigma^2} \frac{d}{d\boldsymbol{\theta}} \left(\mathbf{y}^\top \mathbf{y} - 2\mathbf{y}^\top \mathbf{X}\boldsymbol{\theta} + \boldsymbol{\theta}^\top \mathbf{X}^\top \mathbf{X}\boldsymbol{\theta} \right) \quad (9.10b)$$

$$= \frac{1}{\sigma^2} (-\mathbf{y}^\top \mathbf{X} + \boldsymbol{\theta}^\top \mathbf{X}^\top \mathbf{X}) \in \mathbb{R}^{1 \times D}. \quad (9.10c)$$

As a necessary optimality condition we set this gradient to $\mathbf{0}$ and obtain

$$\frac{d\mathcal{L}}{d\boldsymbol{\theta}} = \mathbf{0} \stackrel{(9.10c)}{\iff} \boldsymbol{\theta}^\top \mathbf{X}^\top \mathbf{X} = \mathbf{y}^\top \mathbf{X} \quad (9.11a)$$

$$\iff \boldsymbol{\theta}^\top = \mathbf{y}^\top \mathbf{X} (\mathbf{X}^\top \mathbf{X})^{-1} \quad (9.11b)$$

$$\iff \boldsymbol{\theta}_{\text{ML}} = (\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{X}^\top \mathbf{y}. \quad (9.11c)$$

We could right-multiply the first equation by $(\mathbf{X}^\top \mathbf{X})^{-1}$ because $\mathbf{X}^\top \mathbf{X}$ is positive definite (if we do not have two identical inputs $\mathbf{x}_i, \mathbf{x}_j$ for $i \neq j$).

Remark. In this case, setting the gradient to $\mathbf{0}$ is a necessary and sufficient condition and we obtain a global minimum since the Hessian $\nabla_{\boldsymbol{\theta}}^2 \mathcal{L}(\boldsymbol{\theta}) = \mathbf{X}^\top \mathbf{X} \in \mathbb{R}^{D \times D}$ is positive definite. \diamond

Example 9.1 (Fitting Lines)

Let us have a look at Figure 9.2, where we aim to fit a straight line $f(x) = \theta x$, where θ is an unknown slope, to a data set using maximum likelihood estimation. Examples of functions in this model class (straight lines) are shown in Figure 9.2(a). For the data set shown in Figure 9.2(b) we find the maximum likelihood estimate of the slope parameter θ using (9.11c) and obtain the maximum likelihood linear function in Figure 9.2(c).

Linear regression⁵¹⁵⁶ refers to “linear-in-the-parameters” regression models, but the inputs can undergo any nonlinear transformation.

Maximum Likelihood Estimation with Features

So far, we considered the linear regression setting described in (9.2), which allowed us to fit straight lines to data using maximum likelihood estimation. However, straight lines are not particularly expressive when it comes to fitting more interesting data. Fortunately, linear regression offers us a way to fit nonlinear functions within the linear regression framework: Since “linear regression” only refers to “linear in the parameters”, we can

perform an arbitrary nonlinear transformation $\phi(\mathbf{x})$ of the inputs \mathbf{x} and then linearly combine the components of the result. The model parameters $\boldsymbol{\theta}$ still appear only linearly. The corresponding linear regression model is

$$y = \boldsymbol{\phi}^\top(\mathbf{x})\boldsymbol{\theta} + \epsilon = \sum_{k=0}^{K-1} \theta_k \phi_k(\mathbf{x}) + \epsilon, \quad (9.12)$$

5157 where $\phi : \mathbb{R}^D \rightarrow \mathbb{R}^K$ is a (nonlinear) transformation of the inputs \mathbf{x} and
 5158 $\phi_k : \mathbb{R}^D \rightarrow \mathbb{R}$ is the k th component of the *feature vector* ϕ .

feature vector

Example 9.2 (Polynomial Regression)

We are concerned with a regression problem $y = \boldsymbol{\phi}^\top(x)\boldsymbol{\theta} + \epsilon$, where $x \in \mathbb{R}$ and $\boldsymbol{\theta} \in \mathbb{R}^K$. A transformation that is often used in this context is

$$\phi(x) = \begin{bmatrix} \phi_0(x) \\ \phi_1(x) \\ \vdots \\ \phi_{K-1}(x) \end{bmatrix} = \begin{bmatrix} 1 \\ x \\ x^2 \\ x^3 \\ \vdots \\ x^{K-1} \end{bmatrix} \in \mathbb{R}^K. \quad (9.13)$$

This means, we “lift” the original one-dimensional input space into a K -dimensional feature space consisting of all monomials x^k for $k = 0, \dots, K-1$. With these features, we can model polynomials of degree $\leq K-1$ within the framework of linear regression: A polynomial of degree $K-1$ is

$$f(x) = \sum_{k=0}^{K-1} \theta_k x^k = \boldsymbol{\phi}^\top(x)\boldsymbol{\theta} \quad (9.14)$$

where ϕ is defined in (9.13) and $\boldsymbol{\theta} = [\theta_0, \dots, \theta_{K-1}]^\top \in \mathbb{R}^K$ contains the (linear) parameters θ_k .

Let us now have a look at maximum likelihood estimation of the parameters $\boldsymbol{\theta}$ in the linear regression model (9.12). We consider training inputs $\mathbf{x}_n \in \mathbb{R}^D$ and targets $y_n \in \mathbb{R}$, $n = 1, \dots, N$, and define the *feature matrix* (design matrix) as

feature matrix
design matrix

$$\Phi := \begin{bmatrix} \boldsymbol{\phi}^\top(\mathbf{x}_1) \\ \vdots \\ \boldsymbol{\phi}^\top(\mathbf{x}_N) \end{bmatrix} = \begin{bmatrix} \phi_0(\mathbf{x}_1) & \cdots & \phi_{K-1}(\mathbf{x}_1) \\ \phi_0(\mathbf{x}_2) & \cdots & \phi_{K-1}(\mathbf{x}_2) \\ \vdots & & \vdots \\ \phi_0(\mathbf{x}_N) & \cdots & \phi_{K-1}(\mathbf{x}_N) \end{bmatrix} \in \mathbb{R}^{N \times K}, \quad (9.15)$$

5159 where $\Phi_{ij} = \phi_j(\mathbf{x}_i)$ and $\phi_j : \mathbb{R}^D \rightarrow \mathbb{R}$.

Example 9.3 (Feature Matrix for Second-order Polynomials)

For a second-order polynomial and N training points $x_n \in \mathbb{R}, n = 1, \dots, N$, the feature matrix is

$$\Phi = \begin{bmatrix} 1 & x_1 & x_1^2 \\ 1 & x_2 & x_2^2 \\ \vdots & \vdots & \vdots \\ 1 & x_N & x_N^2 \end{bmatrix}. \quad (9.16)$$

With the feature matrix Φ defined in (9.15) the negative log-likelihood for the linear regression model (9.12) can be written as

$$-\log p(\mathbf{y} | \mathbf{X}, \boldsymbol{\theta}) = \frac{1}{2\sigma^2} (\mathbf{y} - \Phi\boldsymbol{\theta})^\top (\mathbf{y} - \Phi\boldsymbol{\theta}) + \text{const}. \quad (9.17)$$

Comparing (9.17) with the negative log-likelihood in (9.9b) for the “feature-free” model, we immediately see we just need to replace \mathbf{X} with Φ . Since both \mathbf{X} and Φ are independent of the parameters $\boldsymbol{\theta}$ that we wish to optimize, we arrive immediately at the *maximum likelihood estimate*

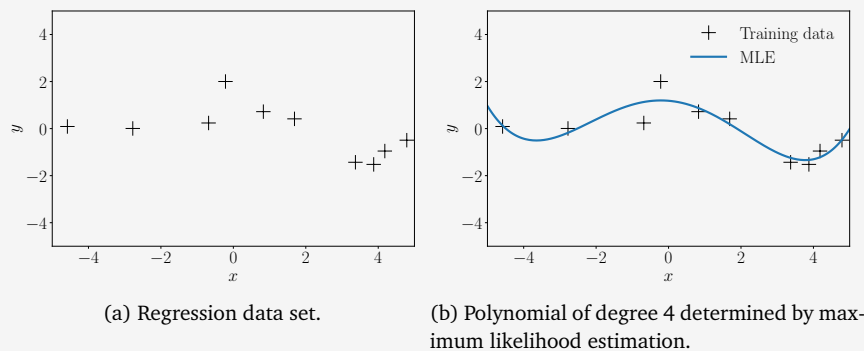
$$\boldsymbol{\theta}_{\text{ML}} = (\Phi^\top \Phi)^{-1} \Phi^\top \mathbf{y} \quad (9.18)$$

for the linear regression problem with nonlinear features defined in (9.12).

Remark. When we were working without features, we required $\mathbf{X}^\top \mathbf{X}$ to be invertible, which is the case when the rows of \mathbf{X} are linearly independent. In (9.18), we therefore require $\Phi^\top \Phi$ to be invertible. This is the case if and only if the rows of the feature matrix are linearly independent. Nonlinear feature transformations can make previously linearly dependent inputs \mathbf{X} linearly independent (and vice versa). \diamond

Example 9.4 (Maximum Likelihood Polynomial Fit)

Figure 9.4
Polynomial regression. (a) Data set consisting of (x_n, y_n) pairs, $n = 1, \dots, 10$; (b) Maximum likelihood polynomial of degree 4.



Consider the data set in Figure 9.5(a). The data set consists of $N = 20$

pairs (x_n, y_n) , where $x_n \sim \mathcal{U}[-5, 5]$ and $y_n = -\sin(x_n/5) + \cos(x_n) + \epsilon$, where $\epsilon \sim \mathcal{N}(0, 0.2^2)$.

We fit a polynomial of degree $K = 4$ using maximum likelihood estimation, i.e., parameters θ_{ML} are given in (9.18). The maximum likelihood estimate yields function values $\phi^\top(x_*)\theta_{\text{ML}}$ at any test location x_* . The result is shown in Figure 9.5(b).

Estimating the Noise Variance

Thus far, we assumed that the noise variance σ^2 is known. However, we can also use the principle of maximum likelihood estimation to obtain σ_{ML}^2 for the noise variance. To do this, we follow the standard procedure: we write down the log-likelihood, compute its derivative with respect to $\sigma^2 > 0$, set it to 0 and solve:

$$\log p(\mathbf{y} | \mathbf{X}, \theta, \sigma^2) = \sum_{n=1}^N \log \mathcal{N}(y_n | \theta^\top \phi(\mathbf{x}_n), \sigma^2) \quad (9.19a)$$

$$= \sum_{n=1}^N \left(-\frac{1}{2} \log(2\pi) - \frac{1}{2} \log \sigma^2 - \frac{1}{2\sigma^2} (y_n - \theta^\top \phi(\mathbf{x}_n))^2 \right) \quad (9.19b)$$

$$= -\frac{N}{2} \log \sigma^2 - \frac{1}{2\sigma^2} \underbrace{\sum_{n=1}^N (y_n - \theta^\top \phi(\mathbf{x}_n))^2}_{=:s} + \text{const.} \quad (9.19c)$$

The partial derivative of the log-likelihood with respect to σ^2 is then

$$\frac{\partial \log p(\mathbf{y} | \mathbf{X}, \theta, \sigma^2)}{\partial \sigma^2} = -\frac{N}{2\sigma^2} + \frac{1}{2\sigma^4} s = 0 \quad (9.20a)$$

$$\iff \frac{N}{2\sigma^2} = \frac{s}{2\sigma^4} \quad (9.20b)$$

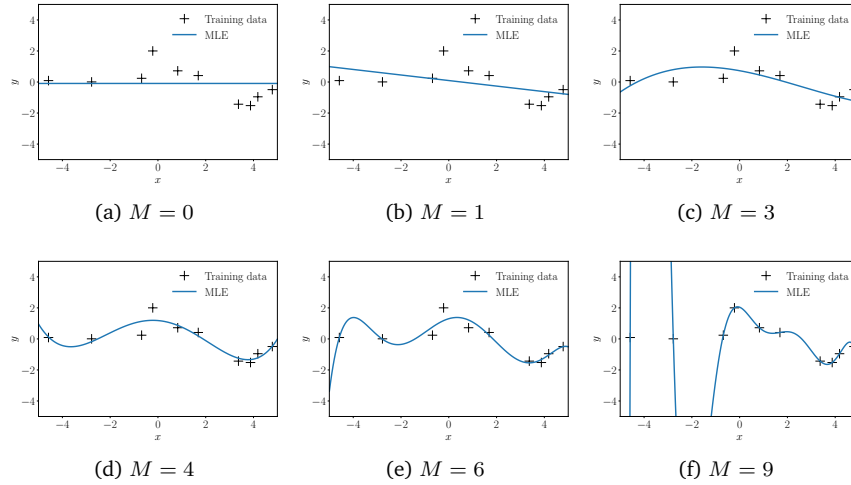
$$\iff \sigma_{\text{ML}}^2 = \frac{s}{N} = \frac{1}{N} \sum_{n=1}^N (y_n - \theta^\top \phi(\mathbf{x}_n))^2. \quad (9.20c)$$

Therefore, the maximum likelihood estimate for the noise variance is the mean squared distance between the noise-free function values $\theta^\top \phi(\mathbf{x}_n)$ and the corresponding noisy observations y_n at \mathbf{x}_n , for $n = 1, \dots, N$.

9.2.2 Overfitting in Linear Regression

We just discussed how to use maximum likelihood estimation to fit linear models (e.g., polynomials) to data. We can evaluate the quality of the model by computing the error/loss incurred. One way of doing this is to compute the negative log-likelihood (9.9b), which we minimized to determine the MLE. Alternatively, given that the noise parameter σ^2 is not

Figure 9.5
Maximum
likelihood fits for
different polynomial
degrees M .



root mean squared
error (RMSE)

a free model parameter, we can ignore the scaling by $1/\sigma^2$, so that we end up with a squared-error-loss function $\|\mathbf{y} - \Phi\boldsymbol{\theta}\|^2$. Instead of using this squared loss, we often use the *root mean squared error (RMSE)*

$$\sqrt{\|\mathbf{y} - \Phi\boldsymbol{\theta}\|^2 / N} = \sqrt{\frac{1}{N} \sum_{n=1}^N (y_n - \phi^\top(x_n)\boldsymbol{\theta})^2}, \quad (9.21)$$

The RMSE is
normalized.

which (a) allows us to compare errors of data sets with different sizes and (b) has the same scale and the same units as the observed function values y_n . For example, assume we fit a model that maps post-codes (\mathbf{x} is given in latitude,longitude) to house prices (y -values are EUR). Then, the RMSE is also measured in EUR, whereas the squared error is given in EUR². If we choose to include the factor σ^2 from the original negative log-likelihood (9.9b) then we end up with a “unit-free” objective.

For model selection (see Section 8.5) we can use the RMSE (or the negative log-likelihood) to determine the best degree of the polynomial by finding the polynomial degree M that minimizes the objective. Given that the polynomial degree is a natural number, we can perform a brute-force search and enumerate all (reasonable) values of M . For a training set of size N it is sufficient to test $0 \leq M \leq N - 1$. For $M \geq N$ we would need to solve an underdetermined system of linear equations so that we would end up with infinitely many solutions.

Figure 9.5 shows a number of polynomial fits determined by maximum likelihood for the dataset from Figure 9.5(a) with $N = 10$ observations. We notice that polynomials of low degree (e.g., constants ($M = 0$) or linear ($M = 1$) fit the data poorly and, hence, are poor representations of the true underlying function. For degrees $M = 3, \dots, 5$ the fits look plausible and smoothly interpolate the data. When we go to higher-degree polynomials, we notice that they fit the data better and better. In the extreme

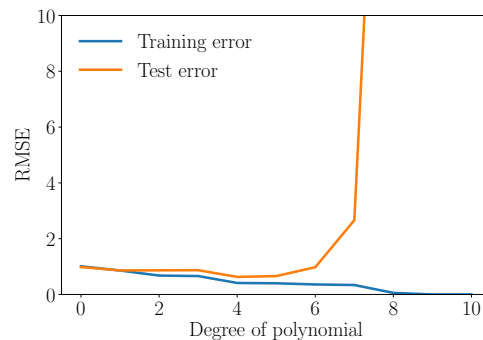


Figure 9.6 Training and test error.

case of $M = N - 1 = 9$, the function will pass through every single data point. However, these high-degree polynomials oscillate wildly and are a poor representation of the underlying function that generated the data, such that we suffer from *overfitting*.

overfitting

Remember that the goal is to achieve good generalization by making accurate predictions for new (unseen) data. We obtain some quantitative insight into the dependence of the generalization performance on the polynomial of degree M by considering a separate test set comprising 200 data points generated using exactly the same procedure used to generate the training set. As test inputs, we chose a linear grid of 200 points in the interval of $[-5, 5]$. For each choice of M , we evaluate the RMSE (9.21) for both the training data and the test data.

Note that the noise variance $\sigma^2 > 0$.

Looking now at the test error, which is a qualitative measure of the generalization properties of the corresponding polynomial, we notice that initially the test error decreases, see Figure 9.6 (orange). For fourth-order polynomials the test error is relatively low and stays relatively constant up to degree 5. However, from degree 6 onward the test error increases significantly, and high-order polynomials have very bad generalization properties. In this particular example, this also is evident from the corresponding maximum likelihood fits in Figure 9.5. Note that the *training error* (blue curve in Figure 9.6) never increases when the degree of the polynomial increases. In our example, the best generalization (the point of the smallest *test error*) is obtained for a polynomial of degree $M = 4$.

training error

test error

9.2.3 Regularization and Maximum A Posteriori Estimation

We just saw that maximum likelihood estimation is prone to overfitting. It often happens that the magnitude of the parameter values becomes relatively big if we run into overfitting (Bishop, 2006). One way to mitigate the effect of overfitting is to penalize big parameter values by a technique called *regularization*. In regularization, we add a term to the log-likelihood that penalizes the magnitude of the parameters θ . A typical example is a

regularization

regularized “loss function” of the form

$$-\log p(\mathbf{y} | \mathbf{X}, \boldsymbol{\theta}) + \lambda \|\boldsymbol{\theta}\|_2^2, \quad (9.22)$$

regularizer

5218 where the second term is the *regularizer*, and $\lambda \geq 0$ controls the “strict-
5219 ness” of the regularization.

LASSO

5220 *Remark.* Instead of the Euclidean norm $\|\cdot\|_2$, we can choose any p -norm
5221 $\|\cdot\|_p$. In practice, smaller values for p lead to sparser solutions. Here,
5222 “sparse” means that many parameter values $\theta_n = 0$, which is also use-
5223 ful for variable selection. For $p = 1$, the regularizer is called *LASSO* (least
5224 absolute shrinkage and selection operator) and was proposed by Tibshi-
5225 rani (1996). \diamond

From a probabilistic perspective, adding a regularizer is identical to placing a prior distribution $p(\boldsymbol{\theta})$ on the parameters and then selecting the parameters that maximize the posterior distribution $p(\boldsymbol{\theta} | \mathbf{X}, \mathbf{y})$, i.e., we choose the parameters $\boldsymbol{\theta}$ that are “most probable” given the training data. The posterior over the parameters $\boldsymbol{\theta}$, given the training data \mathbf{X}, \mathbf{y} , is obtained by applying Bayes’ theorem as

$$p(\boldsymbol{\theta} | \mathbf{X}, \mathbf{y}) = \frac{p(\mathbf{y} | \mathbf{X}, \boldsymbol{\theta})p(\boldsymbol{\theta})}{p(\mathbf{y} | \mathbf{X})}. \quad (9.23)$$

maximum
a-posteriori
MAP

5226 The parameter vector $\boldsymbol{\theta}_{\text{MAP}}$ that maximizes the posterior (9.23) is called
5227 the *maximum a-posteriori (MAP)* estimate.

To find the MAP estimate, we follow steps that are similar in flavor to maximum likelihood estimation. We start with the log-transform and compute the log-posterior as

$$\log p(\boldsymbol{\theta} | \mathbf{X}, \mathbf{y}) = \log p(\mathbf{y} | \mathbf{X}, \boldsymbol{\theta}) + \log p(\boldsymbol{\theta}) + \text{const}, \quad (9.24)$$

5228 where the constant comprises the terms that are independent of $\boldsymbol{\theta}$. We see
5229 that the log-posterior in (9.24) is the sum of the log-likelihood $p(\mathbf{y} | \mathbf{X}, \boldsymbol{\theta})$
5230 and the log-prior $\log p(\boldsymbol{\theta})$.

Remark (Relation to Regularization). Choosing a Gaussian parameter prior $p(\boldsymbol{\theta}) = \mathcal{N}(\mathbf{0}, b^2 \mathbf{I})$, $b^2 = \frac{1}{2\lambda}$, the (negative) log-prior term will be

$$-\log p(\boldsymbol{\theta}) = \underbrace{\lambda \boldsymbol{\theta}^\top \boldsymbol{\theta}}_{=\lambda \|\boldsymbol{\theta}\|_2^2} + \text{const}, \quad (9.25)$$

5231 and we recover exactly the regularization term in (9.22). This means that
5232 for a quadratic regularization, the regularization parameter λ in (9.22)
5233 corresponds to twice the precision (inverse variance) of the Gaussian (iso-
5234 tropic) prior $p(\boldsymbol{\theta})$. Therefore, the log-prior in (9.24) reflects the impact
5235 of the regularizer that penalizes implausible values, i.e., values that are
5236 unlikely under the prior. \diamond

To find the MAP estimate $\boldsymbol{\theta}_{\text{MAP}}$, we minimize the negative log-posterior

distribution with respect to θ , i.e., we solve

$$\theta_{\text{MAP}} \in \arg \min_{\theta} \{-\log p(\mathbf{y} | \mathbf{X}, \theta) - \log p(\theta)\}. \quad (9.26)$$

We determine the gradient of the negative log-posterior with respect to θ as

$$-\frac{d \log p(\theta | \mathbf{X}, \mathbf{y})}{d\theta} = -\frac{d \log p(\mathbf{y} | \mathbf{X}, \theta)}{d\theta} - \frac{d \log p(\theta)}{d\theta}, \quad (9.27)$$

where we identify the first term on the right-hand-side as the gradient of the negative log-likelihood given in (9.10c).

More concretely, with a Gaussian prior $p(\theta) = \mathcal{N}(\mathbf{0}, b^2 \mathbf{I})$ on the parameters θ , the negative log-posterior for the linear regression setting (9.12), we obtain the negative log posterior

$$-\log p(\theta | \mathbf{X}, \mathbf{y}) = \frac{1}{2\sigma^2} (\mathbf{y} - \Phi \theta)^\top (\mathbf{y} - \Phi \theta) + \frac{1}{2b^2} \theta^\top \theta + \text{const.} \quad (9.28)$$

Here, the first term corresponds to the contribution from the log-likelihood, and the second term originates from the log-prior. The gradient of the log-posterior with respect to the parameters θ is then

$$-\frac{d \log p(\theta | \mathbf{X}, \mathbf{y})}{d\theta} = \frac{1}{\sigma^2} (\theta^\top \Phi^\top \Phi - \mathbf{y}^\top \Phi) + \frac{1}{b^2} \theta^\top. \quad (9.29)$$

We will find the MAP estimate θ_{MAP} by setting this gradient to $\mathbf{0}$:

$$\frac{1}{\sigma^2} (\theta^\top \Phi^\top \Phi - \mathbf{y}^\top \Phi) + \frac{1}{b^2} \theta^\top = \mathbf{0} \quad (9.30a)$$

$$\iff \theta^\top \left(\frac{1}{\sigma^2} \Phi^\top \Phi + \frac{1}{b^2} \mathbf{I} \right) - \frac{1}{\sigma^2} \mathbf{y}^\top \Phi = \mathbf{0} \quad (9.30b)$$

$$\iff \theta^\top \left(\Phi^\top \Phi + \frac{\sigma^2}{b^2} \mathbf{I} \right) = \mathbf{y}^\top \Phi \quad (9.30c)$$

$$\iff \theta^\top = \mathbf{y}^\top \Phi \left(\Phi^\top \Phi + \frac{\sigma^2}{b^2} \mathbf{I} \right)^{-1} \quad (9.30d)$$

so that we obtain the MAP estimate (by transposing both sides of the last equality)

$$\theta_{\text{MAP}} = \left(\Phi^\top \Phi + \frac{\sigma^2}{b^2} \mathbf{I} \right)^{-1} \Phi^\top \mathbf{y}. \quad (9.31)$$

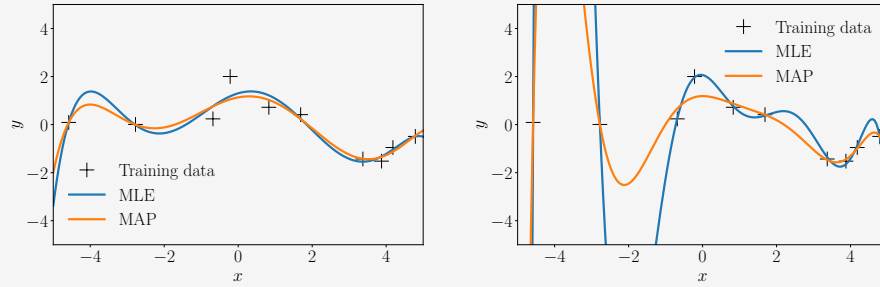
Comparing the MAP estimate in (9.31) with the maximum likelihood estimate in (9.18) we see that the only difference between both solutions is the additional term $\frac{\sigma^2}{b^2} \mathbf{I}$ in the inverse matrix. This term ensures that $\Phi^\top \Phi + \frac{\sigma^2}{b^2} \mathbf{I}$ is symmetric and strictly positive definite (i.e., its inverse exists) and plays the role of the *regularizer*.

$\Phi^\top \Phi$ is symmetric and positive semidefinite and the additional term is strictly positive definite, such that all eigenvalues of the matrix to be inverted are positive.

regularizer

Example 9.5 (MAP Estimation for Polynomial Regression)

Figure 9.7
Polynomial
regression:
Maximum
likelihood and MAP
estimates.



In the polynomial regression example from Section 9.2.1, we place a Gaussian prior $p(\boldsymbol{\theta}) = \mathcal{N}(\mathbf{0}, \mathbf{I})$ on the parameters $\boldsymbol{\theta}$ and determine the MAP estimates according to (9.31). In Figure 9.7, we show both the maximum likelihood and the MAP estimates for polynomials of degree 6 (left) and degree 8 (right). The prior (regularizer) does not play a significant role for the low-degree polynomial, but keeps the function relatively smooth for higher-degree polynomials. However, the MAP estimate can only push the boundaries of overfitting – it is not a general solution to this problem.

5244 In the following, we will discuss Bayesian linear regression where we
5245 average over all plausible sets of parameters instead of focusing on a point
5246 estimate.

9.3 Bayesian Linear Regression

5247

5248 Previously, we looked at linear regression models where we estimated the
5249 model parameters $\boldsymbol{\theta}$, e.g., by means of maximum likelihood or MAP esti-
5250 mation. We discovered that MLE can lead to severe overfitting, in particu-
5251 lar, in the small-data regime. MAP addresses this issue by placing a prior

Bayesian linear
regression

5252 on the parameters that plays the role of a regularizer.
5253 *Bayesian linear regression* pushes the idea of the parameter prior a step
5254 further and does not even attempt to compute a point estimate of the pa-
5255 rameters, but instead the full posterior over the parameters is taken into
5256 account when making predictions. This means we do not fit any param-
5257 eters, but we compute an average over all plausible parameters settings
5258 (according to the posterior).

9.3.1 Model

In Bayesian linear regression, we consider the model

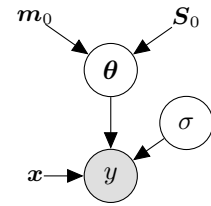
$$\begin{aligned} \text{prior} \quad & p(\boldsymbol{\theta}) = \mathcal{N}(\mathbf{m}_0, \mathbf{S}_0), \\ \text{likelihood} \quad & p(y | \mathbf{x}, \boldsymbol{\theta}) = \mathcal{N}(y | \boldsymbol{\phi}^\top(\mathbf{x})\boldsymbol{\theta}, \sigma^2), \end{aligned} \quad (9.32)$$

where we now explicitly place a Gaussian prior $p(\boldsymbol{\theta}) = \mathcal{N}(\mathbf{m}_0, \mathbf{S}_0)$ on $\boldsymbol{\theta}$, which turns the parameter vector into a latent variable. The full probabilistic model, i.e., the joint distribution of observed and latent variables, y and $\boldsymbol{\theta}$, respectively, is

$$p(y, \boldsymbol{\theta} | \mathbf{x}) = p(y | \mathbf{x}, \boldsymbol{\theta})p(\boldsymbol{\theta}), \quad (9.33)$$

which allows us to write down the corresponding graphical model in Figure 9.8, where we made the parameters of the Gaussian prior on $\boldsymbol{\theta}$ explicit.

Figure 9.8
Graphical model for Bayesian linear regression.



9.3.2 Prior Predictions

In practice, we are usually not so much interested in the parameter values $\boldsymbol{\theta}$. Instead, our focus often lies in the predictions we make with those parameter values. In a Bayesian setting, we take the parameter distribution and average over all plausible parameter settings when we make predictions. More specifically, to make predictions at an input location \mathbf{x}_* , we integrate out $\boldsymbol{\theta}$ and obtain

$$p(y_* | \mathbf{x}_*) = \int p(y_* | \mathbf{x}_*, \boldsymbol{\theta})p(\boldsymbol{\theta})d\boldsymbol{\theta} = \mathbb{E}_{\boldsymbol{\theta}}[p(y_* | \mathbf{x}_*, \boldsymbol{\theta})], \quad (9.34)$$

which we can interpret as the average prediction of $y_* | \mathbf{x}_*, \boldsymbol{\theta}$ for all plausible parameters $\boldsymbol{\theta}$ according to the prior distribution $p(\boldsymbol{\theta})$. Note that predictions using the prior distribution only require to specify the input locations \mathbf{x}_* , but no training data.

In our model, we chose a conjugate (Gaussian) prior on $\boldsymbol{\theta}$ so that the predictive distribution is Gaussian as well (and can be computed in closed form): With the prior distribution $p(\boldsymbol{\theta}) = \mathcal{N}(\mathbf{m}_0, \mathbf{S}_0)$, we obtain the predictive distribution as

$$p(y_* | \mathbf{x}_*) = \mathcal{N}(\boldsymbol{\phi}^\top(\mathbf{x}_*)\mathbf{m}_0, \boldsymbol{\phi}^\top(\mathbf{x}_*)\mathbf{S}_0\boldsymbol{\phi}(\mathbf{x}_*) + \sigma^2), \quad (9.35)$$

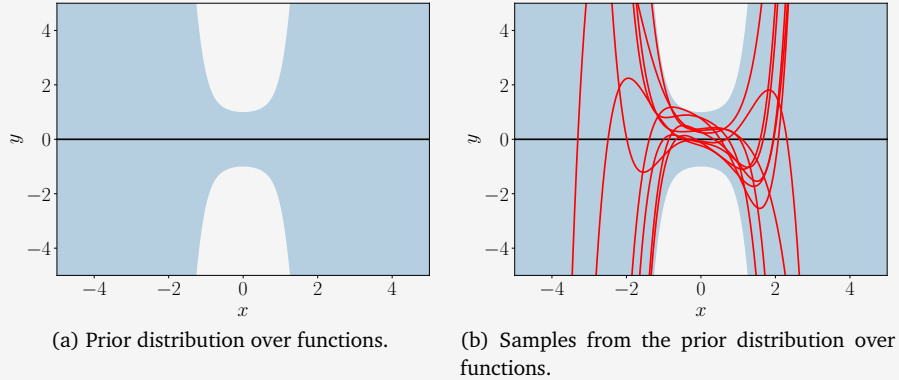
where we used that (i) the prediction is Gaussian due to conjugacy and the marginalization property of Gaussians, (ii), the Gaussian noise is independent so that $\mathbb{V}[y_*] = \mathbb{V}[\boldsymbol{\phi}^\top(\mathbf{x}_*)\boldsymbol{\theta}] + \mathbb{V}[\epsilon]$, (iii) y_* is a linear transformation of $\boldsymbol{\theta}$ so that we can apply the rules for computing the mean and covariance of the prediction analytically by using (6.50) and (6.51), respectively.

In (9.35), the term $\boldsymbol{\phi}^\top(\mathbf{x}_*)\mathbf{S}_0\boldsymbol{\phi}(\mathbf{x}_*)$ in the predictive variance explicitly accounts for the uncertainty associated with the parameters $\boldsymbol{\theta}$, whereas σ^2 is the uncertainty contribution due to the measurement noise.

Example 9.6 (Prior over Functions)

Let us consider a Bayesian linear regression problem with polynomials of degree 5. We choose a parameter prior $p(\boldsymbol{\theta}) = \mathcal{N}(\mathbf{0}, \frac{1}{4}\mathbf{I})$. Figure 9.9 visualizes the distribution over functions induced by this parameter prior, including some function samples from this prior.

Figure 9.9 Prior over functions. (a) Distribution over functions represented by the mean function (black line) and the marginal uncertainties (shaded), representing the 95% confidence bounds; (b) Samples from the prior over functions, which are induced by the samples from the parameter prior.



So far, we looked at computing predictions using the parameter prior $p(\boldsymbol{\theta})$. However, when we have a parameter posterior (given some training data \mathbf{X}, \mathbf{y}), the same principles for prediction and inference hold as in (9.34) – we just need to replace the prior $p(\boldsymbol{\theta})$ with the posterior $p(\boldsymbol{\theta} | \mathbf{X}, \mathbf{y})$. In the following, we will derive the posterior distribution in detail before using it to make predictions.

9.3.3 Posterior Distribution

Given a training set of inputs $\mathbf{x}_n \in \mathbb{R}^D$ and corresponding observations $y_n \in \mathbb{R}$, $n = 1, \dots, N$, we compute the posterior over the parameters using Bayes' theorem as

$$p(\boldsymbol{\theta} | \mathbf{X}, \mathbf{y}) = \frac{p(\mathbf{y} | \mathbf{X}, \boldsymbol{\theta})p(\boldsymbol{\theta})}{p(\mathbf{y} | \mathbf{X})}, \quad (9.36)$$

where \mathbf{X} is the collection of training inputs and \mathbf{y} the collection of training targets. Furthermore, $p(\mathbf{y} | \mathbf{X}, \boldsymbol{\theta})$ is the likelihood, $p(\boldsymbol{\theta})$ the parameter prior and

$$p(\mathbf{y} | \mathbf{X}) = \int p(\mathbf{y} | \mathbf{X}, \boldsymbol{\theta})p(\boldsymbol{\theta})d\boldsymbol{\theta} \quad (9.37)$$

marginal likelihood
evidence

the *marginal likelihood/evidence*, which is independent of the parameters $\boldsymbol{\theta}$ and ensures that the posterior is normalized, i.e., it integrates to 1. We can think of the marginal likelihood as the likelihood averaged over all possible parameter settings (with respect to the prior distribution $p(\boldsymbol{\theta})$).

In our specific model (9.32), the posterior (9.36) can be computed in closed form as

$$p(\boldsymbol{\theta} | \mathbf{X}, \mathbf{y}) = \mathcal{N}(\boldsymbol{\theta} | \mathbf{m}_N, \mathbf{S}_N), \quad (9.38a)$$

$$\mathbf{S}_N = (\mathbf{S}_0^{-1} + \sigma^{-2} \boldsymbol{\Phi}^\top \boldsymbol{\Phi})^{-1}, \quad (9.38b)$$

$$\mathbf{m}_N = \mathbf{S}_N (\mathbf{S}_0^{-1} \mathbf{m}_0 + \sigma^{-2} \boldsymbol{\Phi}^\top \mathbf{y}), \quad (9.38c)$$

where the subscript N indicates the size of the training set. In the following, we will detail how we arrive at this posterior.

Bayes' theorem tells us that the posterior $p(\boldsymbol{\theta} | \mathbf{X}, \mathbf{y})$ is proportional to the product of the likelihood $p(\mathbf{y} | \mathbf{X}, \boldsymbol{\theta})$ and the prior $p(\boldsymbol{\theta})$:

$$\text{posterior} \quad p(\boldsymbol{\theta} | \mathbf{X}, \mathbf{y}) = \frac{p(\mathbf{y} | \mathbf{X}, \boldsymbol{\theta}) p(\boldsymbol{\theta})}{p(\mathbf{y} | \mathbf{X})} \quad (9.39a)$$

$$\text{likelihood} \quad p(\mathbf{y} | \mathbf{X}, \boldsymbol{\theta}) = \mathcal{N}(\mathbf{y} | \boldsymbol{\Phi} \boldsymbol{\theta}, \sigma^2 \mathbf{I}) \quad (9.39b)$$

$$\text{prior} \quad p(\boldsymbol{\theta}) = \mathcal{N}(\boldsymbol{\theta} | \mathbf{m}_0, \mathbf{S}_0) \quad (9.39c)$$

We will discuss two approaches to derive the desired posterior.

Approach 1: Linear Transformation of Gaussian Random Variables

Looking at the numerator of the posterior in (9.39a), we know that the Gaussian prior times the Gaussian likelihood (where the parameters on which we place the Gaussian appears linearly in the mean) is an (un-normalized) Gaussian (see Section 6.6.2). If necessary, we can find the normalizing constant using (6.114). If we want to compute that product by using the results from (6.112)–(6.113) in Section 6.6.2, we need to ensure the product has the “right” form, i.e.,

$$\mathcal{N}(\mathbf{y} | \boldsymbol{\Phi} \boldsymbol{\theta}, \sigma^2 \mathbf{I}) \mathcal{N}(\boldsymbol{\theta} | \mathbf{m}_0, \mathbf{S}_0) = \mathcal{N}(\boldsymbol{\theta} | \boldsymbol{\mu}, \boldsymbol{\Sigma}) \mathcal{N}(\boldsymbol{\theta} | \mathbf{m}_0, \mathbf{S}_0) \quad (9.40)$$

for some $\boldsymbol{\mu}, \boldsymbol{\Sigma}$. With this form we determine the desired product immediately as

$$\mathcal{N}(\boldsymbol{\theta} | \boldsymbol{\mu}, \boldsymbol{\Sigma}) \mathcal{N}(\boldsymbol{\theta} | \mathbf{m}_0, \mathbf{S}_0) \propto \mathcal{N}(\boldsymbol{\theta} | \mathbf{m}_N, \mathbf{S}_N) \quad (9.41a)$$

$$\mathbf{S}_N = (\mathbf{S}_0^{-1} + \boldsymbol{\Sigma}^{-1})^{-1} \quad (9.41b)$$

$$\mathbf{m}_N = \mathbf{S}_N (\mathbf{S}_0^{-1} \mathbf{m}_0 + \boldsymbol{\Sigma}^{-1} \boldsymbol{\mu}). \quad (9.41c)$$

In order to get the “right” form, we need to turn $\mathcal{N}(\mathbf{y} | \boldsymbol{\Phi} \boldsymbol{\theta}, \sigma^2 \mathbf{I})$ into $\mathcal{N}(\boldsymbol{\theta} | \boldsymbol{\mu}, \boldsymbol{\Sigma})$ for appropriate choices of $\boldsymbol{\mu}, \boldsymbol{\Sigma}$. We will do this by using a linear transformation of Gaussian random variables (see Section 6.6), which allows us to exploit the property that linearly transformed Gaussian random variables are Gaussian distributed. More specifically, we will find $\boldsymbol{\mu} = \mathbf{B} \mathbf{y}$ and $\boldsymbol{\Sigma} = \sigma^2 \mathbf{B} \mathbf{B}^\top$ by linearly transforming the relationship $\mathbf{y} = \boldsymbol{\Phi} \boldsymbol{\theta}$ in the likelihood into $\mathbf{B} \mathbf{y} = \boldsymbol{\theta}$ for a suitable \mathbf{B} . We obtain

$$\mathbf{y} = \boldsymbol{\Phi} \boldsymbol{\theta} \xrightarrow{\times \boldsymbol{\Phi}^\top} \boldsymbol{\Phi}^\top \mathbf{y} = \boldsymbol{\Phi}^\top \boldsymbol{\Phi} \boldsymbol{\theta} \xrightarrow{\times (\boldsymbol{\Phi}^\top \boldsymbol{\Phi})^{-1}} \underbrace{(\boldsymbol{\Phi}^\top \boldsymbol{\Phi})^{-1} \boldsymbol{\Phi}^\top}_{=: \mathbf{B}} \mathbf{y} = \boldsymbol{\theta} \quad (9.42)$$

Therefore, we can write $\theta = By$, and by using the rules for linear transformations of the mean and covariance from (6.50)–(6.51) we obtain

$$\mathcal{N}(\theta | By, \sigma^2 BB^\top) = \mathcal{N}(\theta | (\Phi^\top \Phi)^{-1} \Phi^\top y, \sigma^2 (\Phi^\top \Phi)^{-1}) \quad (9.43)$$

after some re-arranging of the terms for the covariance matrix.

If we now look at (9.43) and define its mean as μ and covariance matrix as Σ in (9.41c) and (9.41b), respectively, we obtain the covariance S_N and the mean m_N of the parameter posterior $\mathcal{N}(\theta | m_N, S_N)$ as

$$S_N = (S_0^{-1} + \sigma^{-2} \Phi^\top \Phi)^{-1}, \quad (9.44a)$$

$$m_N = S_N (S_0^{-1} m_0 + \underbrace{\sigma^{-2} (\Phi^\top \Phi)}_{\Sigma^{-1}} \underbrace{(\Phi^\top \Phi)^{-1} \Phi^\top y}_{\mu}) \quad (9.44b)$$

$$= S_N (S_0^{-1} m_0 + \sigma^{-2} \Phi^\top y), \quad (9.44c)$$

The posterior mean equals the MAP estimate.

respectively. Note that the posterior mean m_N equals the MAP estimate θ_{MAP} from (9.31). This also makes sense since the posterior distribution is unimodal (Gaussian) with its maximum at the mean.

Remark. The posterior precision (inverse covariance)

$$S_N^{-1} = S_0^{-1} + \frac{1}{\sigma^2} \Phi^\top \Phi \quad (9.45)$$

$\Phi^\top \Phi$ accumulates contributions from the data.

of the parameters θ (see (9.44a)) contains two terms: S_0^{-1} is the prior precision and $\frac{1}{\sigma^2} \Phi^\top \Phi$ is a data-dependent (precision) term. Both terms (matrices) are symmetric and positive definite. The data-dependent term $\frac{1}{\sigma^2} \Phi^\top \Phi$ grows as more data is taken into account. This means (at least) two things:

- The posterior precision grows as more and more data is taken into account; therefore, the covariance, and with it the uncertainty about the parameters, shrinks.
- The relative influence of the parameter prior vanishes for large N .

Therefore, for $N \rightarrow \infty$ the prior plays no role, and the parameter posterior tends to a point estimate, the MAP estimate. \diamond

Approach 2: Completing the Squares

Instead of looking at the product of the prior and the likelihood, we can transform the problem into log-space and solve for the mean and covariance of the posterior by completing the squares.

The sum of the log-prior and the log-likelihood is

$$\log \mathcal{N}(y | \Phi \theta, \sigma^2 I) + \log \mathcal{N}(\theta | m_0, S_0) \quad (9.46a)$$

$$= -\frac{1}{2} (\sigma^{-2} (y - \Phi \theta)^\top (y - \Phi \theta) + (\theta - m_0)^\top S_0^{-1} (\theta - m_0)) + \text{const} \quad (9.46b)$$

where the constant contains terms independent of θ . We will ignore the constant in the following. We now factorize (9.46b), which yields

$$-\frac{1}{2}(\sigma^{-2}\mathbf{y}^\top\mathbf{y} - 2\sigma^{-2}\mathbf{y}^\top\Phi\theta + \theta^\top\sigma^{-2}\Phi^\top\Phi\theta + \theta^\top\mathbf{S}_0^{-1}\theta - 2\mathbf{m}_0^\top\mathbf{S}_0^{-1}\theta + \mathbf{m}_0^\top\mathbf{S}_0^{-1}\mathbf{m}_0) \quad (9.47a)$$

$$= -\frac{1}{2}(\theta^\top(\sigma^{-2}\Phi^\top\Phi + \mathbf{S}_0^{-1})\theta - 2(\sigma^{-2}\Phi^\top\mathbf{y} + \mathbf{S}_0^{-1}\mathbf{m}_0)^\top\theta) + \text{const}, \quad (9.47b)$$

where the constant contains the black terms in (9.47a), which are independent of θ . The orange terms are terms that are linear in θ , and the blue terms are the ones that are quadratic in θ . By inspecting (9.47b), we find that this equation is quadratic in θ . The fact that the unnormalized log-posterior distribution is a (negative) quadratic form implies that the posterior is Gaussian, i.e.,

$$p(\theta | \mathbf{X}, \mathbf{y}) = \exp(\log p(\theta | \mathbf{X}, \mathbf{y})) \propto \exp(\log p(\mathbf{y} | \mathbf{X}, \theta) + \log p(\theta)) \quad (9.48a)$$

$$\propto \exp\left(-\frac{1}{2}(\theta^\top(\sigma^{-2}\Phi^\top\Phi + \mathbf{S}_0^{-1})\theta - 2(\sigma^{-2}\Phi^\top\mathbf{y} + \mathbf{S}_0^{-1}\mathbf{m}_0)^\top\theta)\right), \quad (9.48b)$$

5309 where we used (9.47b) in the last expression.

The remaining task is it to bring this (unnormalized) Gaussian into the form that is proportional to $\mathcal{N}(\theta | \mathbf{m}_N, \mathbf{S}_N)$, i.e., we need to identify the mean \mathbf{m}_N and the covariance matrix \mathbf{S}_N . To do this, we use the concept of *completing the squares*. The desired log-posterior is

completing the squares

$$\log \mathcal{N}(\theta | \mathbf{m}_N, \mathbf{S}_N) = -\frac{1}{2}((\theta - \mathbf{m}_N)^\top \mathbf{S}_N^{-1}(\theta - \mathbf{m}_N)) + \text{const} \quad (9.49a)$$

$$= -\frac{1}{2}(\theta^\top \mathbf{S}_N^{-1}\theta - 2\mathbf{m}_N^\top \mathbf{S}_N^{-1}\theta + \mathbf{m}_N^\top \mathbf{S}_N^{-1}\mathbf{m}_N). \quad (9.49b)$$

Here, we factorized the quadratic form $(\theta - \mathbf{m}_N)^\top \mathbf{S}_N^{-1}(\theta - \mathbf{m}_N)$ into a term that is quadratic in θ alone (blue), a term that is linear in θ (orange), and a constant term (black). This allows us now to find \mathbf{S}_N and \mathbf{m}_N by matching the colored expressions in (9.47b) and (9.49b), which yields

$$\mathbf{S}_N^{-1} = \Phi^\top \sigma^{-2} \mathbf{I} \Phi + \mathbf{S}_0^{-1} \iff \mathbf{S}_N = (\sigma^{-2} \Phi^\top \Phi + \mathbf{S}_0^{-1})^{-1}, \quad (9.50)$$

$$\mathbf{m}_N^\top \mathbf{S}_N^{-1} = (\sigma^{-2} \Phi^\top \mathbf{y} + \mathbf{S}_0^{-1} \mathbf{m}_0)^\top \iff \mathbf{m}_N = \mathbf{S}_N (\sigma^{-2} \Phi^\top \mathbf{y} + \mathbf{S}_0^{-1} \mathbf{m}_0). \quad (9.51)$$

5310 This is identical to the solution in (9.44a)–(9.44c), which we obtained by
5311 linear transformations of Gaussian random variables.

Remark (Completing the Squares—General Approach). If we are given an equation

$$\mathbf{x}^\top \mathbf{A} \mathbf{x} - 2\mathbf{a}^\top \mathbf{x} + \text{const}_1, \quad (9.52)$$

where \mathbf{A} is symmetric and positive definite, which we wish to bring into the form

$$(\mathbf{x} - \boldsymbol{\mu})^\top \boldsymbol{\Sigma} (\mathbf{x} - \boldsymbol{\mu}) + \text{const}_2, \quad (9.53)$$

we can do this by setting

$$\boldsymbol{\Sigma} := \mathbf{A}, \quad (9.54)$$

$$\boldsymbol{\mu} := \boldsymbol{\Sigma}^{-1} \mathbf{a} \quad (9.55)$$

5312 and $\text{const}_2 = \text{const}_1 - \boldsymbol{\mu}^\top \boldsymbol{\Sigma} \boldsymbol{\mu}$. \diamond

We can see that the terms inside the exponential in (9.48b) are of the form (9.52) with

$$\mathbf{A} := \sigma^{-2} \boldsymbol{\Phi}^\top \boldsymbol{\Phi} + \mathbf{S}_0^{-1}, \quad (9.56)$$

$$\mathbf{a} := \sigma^{-2} \boldsymbol{\Phi}^\top \mathbf{y} + \mathbf{S}_0^{-1} \mathbf{m}_0. \quad (9.57)$$

5313 Since \mathbf{A}, \mathbf{a} can be difficult to identify in equations like (9.47a), it is of-
5314 ten helpful to bring these equations into the form (9.52) that decouples
5315 quadratic term, linear terms and constants, which simplifies finding the
5316 desired solution.

9.3.4 Posterior Predictions

5317

In (9.34), we computed the predictive distribution of y_* at a test input \mathbf{x}_* using the parameter prior $p(\boldsymbol{\theta})$. In principle, predicting with the parameter posterior $p(\boldsymbol{\theta} | \mathbf{X}, \mathbf{y})$ is not fundamentally different given that in our conjugate model the prior and posterior are both Gaussian (with different parameters). Therefore, by following the same reasoning as in Section 9.3.2 we obtain the (posterior) predictive distribution

$$p(y_* | \mathbf{X}, \mathbf{y}, \mathbf{x}_*) = \int p(y_* | \mathbf{x}_*, \boldsymbol{\theta}) p(\boldsymbol{\theta} | \mathbf{X}, \mathbf{y}) d\boldsymbol{\theta} \quad (9.58a)$$

$$= \int \mathcal{N}(y_* | \boldsymbol{\phi}^\top(\mathbf{x}_*) \boldsymbol{\theta}, \sigma^2) \mathcal{N}(\boldsymbol{\theta} | \mathbf{m}_N, \mathbf{S}_N) d\boldsymbol{\theta} \quad (9.58b)$$

$$= \mathcal{N}(y_* | \boldsymbol{\phi}^\top(\mathbf{x}_*) \mathbf{m}_N, \boldsymbol{\phi}^\top(\mathbf{x}_*) \mathbf{S}_N \boldsymbol{\phi}(\mathbf{x}_*) + \sigma^2) \quad (9.58c)$$

5318 The term $\boldsymbol{\phi}^\top(\mathbf{x}_*) \mathbf{S}_N \boldsymbol{\phi}(\mathbf{x}_*)$ reflects the posterior uncertainty associated
5319 with the parameters $\boldsymbol{\theta}$. Note that \mathbf{S}_N depends on the training inputs \mathbf{X} ,
5320 see (9.44a). The predictive mean coincides with the MAP estimate.

Remark (Mean and Variance of Noise-Free Function Values). In many cases, we are not interested in the predictive distribution $p(y_* | \mathbf{X}, \mathbf{y}, \mathbf{x}_*)$ of a (noisy) observation. Instead, we would like to obtain the distribution

of the (noise-free) latent function values $f(\mathbf{x}_*) = \phi^\top(\mathbf{x}_*)\boldsymbol{\theta}$. We determine the corresponding moments by exploiting the properties of means and variances, which yields

$$\begin{aligned}\mathbb{E}[f(\mathbf{x}_*) | \mathbf{X}, \mathbf{y}] &= \mathbb{E}_\theta[\phi^\top(\mathbf{x}_*)\boldsymbol{\theta} | \mathbf{X}, \mathbf{y}] = \phi^\top(\mathbf{x}_*)\mathbb{E}_\theta[\boldsymbol{\theta} | \mathbf{X}, \mathbf{y}] \\ &= \phi^\top(\mathbf{x}_*)\mathbf{m}_N = \mathbf{m}_N^\top\phi(\mathbf{x}_*),\end{aligned}\quad (9.59)$$

$$\begin{aligned}\mathbb{V}_\theta[f(\mathbf{x}_*) | \mathbf{X}, \mathbf{y}] &= \mathbb{V}_\theta[\phi^\top(\mathbf{x}_*)\boldsymbol{\theta} | \mathbf{X}, \mathbf{y}] \\ &= \phi^\top(\mathbf{x}_*)\mathbb{V}_\theta[\boldsymbol{\theta} | \mathbf{X}, \mathbf{y}]\phi(\mathbf{x}_*) \\ &= \phi^\top(\mathbf{x}_*)\mathbf{S}_N\phi(\mathbf{x}_*)\end{aligned}\quad (9.60)$$

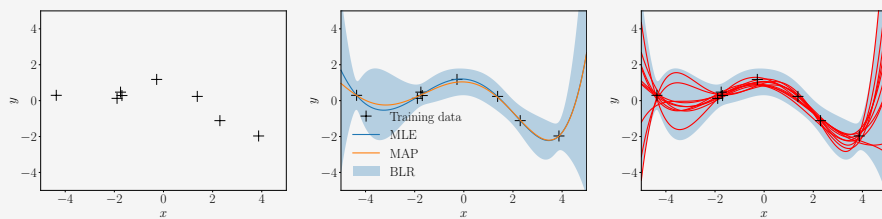
We see that the predictive mean is the same as the predictive mean for noisy observations as the noise has mean 0, and the predictive variance only differs by σ^2 , which is the variance of the measurement noise: When we predict noisy function values, we need to include σ^2 as a source of uncertainty, but this term is not needed for noise-free predictions. Here, the only remaining uncertainty stems from the parameter posterior. \diamond

Remark (Distribution over Functions). The fact that we integrate out the parameters $\boldsymbol{\theta}$ induces a distribution over functions: If we sample $\boldsymbol{\theta}_i \sim p(\boldsymbol{\theta} | \mathbf{X}, \mathbf{y})$ from the parameter posterior, we obtain a single function realization $\boldsymbol{\theta}_i^\top\phi(\cdot)$. The *mean function*, i.e., the set of all expected function values $\mathbb{E}_\theta[f(\cdot) | \boldsymbol{\theta}, \mathbf{X}, \mathbf{y}]$, of this distribution over functions is $\mathbf{m}_N^\top\phi(\cdot)$. The (marginal) variance, i.e., the variance of the function $f(\cdot)$, are given by $\phi^\top(\cdot)\mathbf{S}_N\phi(\cdot)$. \diamond

Integrating out parameters induces a distribution over functions.

mean function

Example 9.7 (Posterior over Functions)



(a) Training data.

(b) Posterior over functions represented by the marginal uncertainties (shaded) showing the 95% predictive confidence bounds, the maximum likelihood estimate (MLE) and the MAP estimate (MAP), which is identical to the posterior mean function.

(c) Samples from the posterior over functions, which are induced by the samples from the parameter posterior.

Figure 9.10 Bayesian linear regression and posterior over functions. (a) Training data; (b) posterior distribution over functions; (c) Samples from the posterior over functions.

Let us revisit the Bayesian linear regression problem with polynomials of degree 5. We choose a parameter prior $p(\boldsymbol{\theta}) = \mathcal{N}(\mathbf{0}, \frac{1}{4}\mathbf{I})$. Figure 9.9

visualizes the prior over functions induced by the parameter prior and sample functions from this prior.

Figure 9.10 shows the posterior over functions that we obtain via Bayesian linear regression. The training dataset is shown in Figure 9.11(a); Figure 9.11(b) shows the posterior distribution over functions, including the functions we would obtain via maximum likelihood and MAP estimation. The function we obtain using the MAP estimate also corresponds to the posterior mean function in the Bayesian linear regression setting. Figure 9.11(c) shows some plausible realizations (samples) of functions under that posterior over functions.

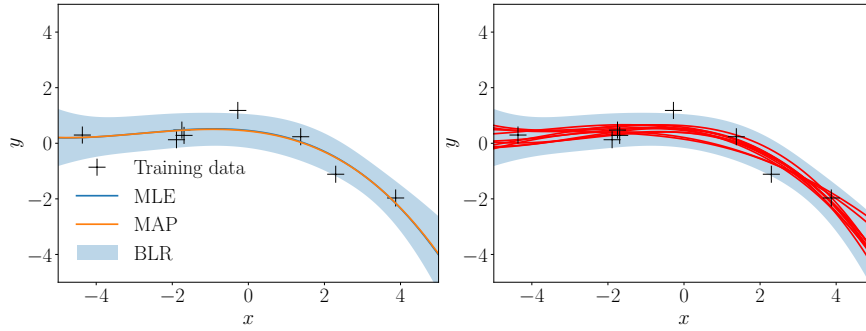
Figure 9.11 shows some examples of the posterior distribution over functions induced by the parameter posterior. For different polynomial degrees M the left panels show the maximum likelihood estimate, the MAP estimate (which is identical to the posterior mean function) and the 95% predictive confidence bounds, represented by the shaded area. The right panels show samples from the posterior over functions: Here, we sampled parameters θ_i from the parameter posterior and computed the function $\phi^\top(\mathbf{x}_*)\theta_i$, which is a single realization of a function under the posterior distribution over functions. For low-order polynomials, the parameter posterior does not allow the parameters to vary much: The sampled functions are nearly identical. When we make the model more flexible by adding more parameters (i.e., we end up with a higher-order polynomial), these parameters are not sufficiently constrained by the posterior, and the sampled functions can be easily visually separated. We also see in the corresponding panels on the left how the uncertainty increases, especially at the boundaries. Although for a 7th-order polynomial the MAP estimate yields a reasonable fit, the Bayesian linear regression model additionally tells us that the posterior uncertainty is huge. This information can be critical when we use these predictions in a decision-making system, where bad decisions can have significant consequences (e.g., in reinforcement learning or robotics).

9.3.5 Computing the Marginal Likelihood

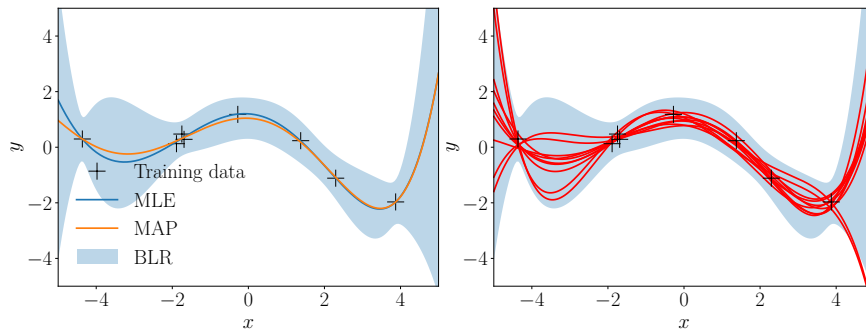
In Section 8.5.2, we highlighted the importance of the marginal likelihood for Bayesian model selection. In the following, we compute the marginal likelihood for Bayesian linear regression with a conjugate Gaussian prior on the parameters, i.e., exactly the setting we have been discussing in this chapter. Just to re-cap, we consider the following generative process:

$$\boldsymbol{\theta} \sim \mathcal{N}(\mathbf{m}_0, \mathbf{S}_0) \quad (9.61a)$$

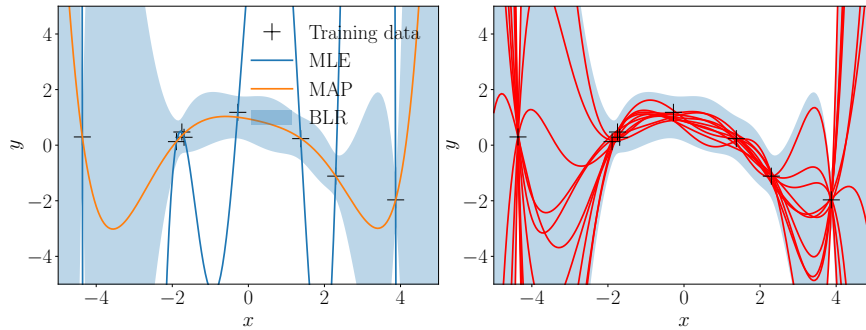
$$y_n | \mathbf{x}_n, \boldsymbol{\theta} \sim \mathcal{N}(\mathbf{x}_n^\top \boldsymbol{\theta}, \sigma^2), \quad (9.61b)$$



(a) Posterior distribution for polynomials of degree $M = 3$ (left) and samples from the posterior over functions (right).



(b) Posterior distribution for polynomials of degree $M = 5$ (left) and samples from the posterior over functions (right).



(c) Posterior distribution for polynomials of degree $M = 7$ (left) and samples from the posterior over functions (right).

Figure 9.11 Bayesian linear regression. Left panels: Shaded areas indicate the 95% predictive confidence bounds. The mean of the Bayesian linear regression model coincides with the MAP estimate. The predictive uncertainty is the sum of the noise term and the posterior parameter uncertainty, which depends on the location of the test input. Right panels: Sampled functions from the posterior distribution.

$n = 1, \dots, N$. The marginal likelihood is given by

$$p(\mathbf{y} | \mathbf{X}) = \int p(\mathbf{y} | \mathbf{X}, \boldsymbol{\theta}) p(\boldsymbol{\theta}) d\boldsymbol{\theta} \quad (9.62a)$$

$$= \int \mathcal{N}(\mathbf{y} | \mathbf{X}\boldsymbol{\theta}, \sigma^2 \mathbf{I}) \mathcal{N}(\boldsymbol{\theta} | \mathbf{m}_0, \mathbf{S}_0) d\boldsymbol{\theta}, \quad (9.62b)$$

The marginal likelihood can be interpreted as the expected likelihood under the prior, i.e., $\mathbb{E}_{\boldsymbol{\theta}}[p(\mathbf{y} | \mathbf{X}, \boldsymbol{\theta})]$.

where we integrate out the model parameters $\boldsymbol{\theta}$. We compute the marginal likelihood in two steps: First, we show that the marginal likelihood is

5358 Gaussian (as a distribution in \mathbf{y}); Second, we compute the mean and co-
5359 variance of this Gaussian.

- 5360 1. The marginal likelihood is Gaussian: From Section 6.6.2 we know that
5361 (i) the product of two Gaussian random variables is an (unnormal-
5362 ized) Gaussian distribution, (ii) a linear transformation of a Gaussian
5363 random variable is Gaussian distributed. In (9.62b), we require a linear
5364 transformation to bring $\mathcal{N}(\mathbf{y} | \mathbf{X}\boldsymbol{\theta}, \sigma^2 \mathbf{I})$ into the form $\mathcal{N}(\boldsymbol{\theta} | \boldsymbol{\mu}, \boldsymbol{\Sigma})$ for
5365 some $\boldsymbol{\mu}, \boldsymbol{\Sigma}$. Once this is done, the integral can be solved in closed form.
5366 The result is the normalizing constant of the product of the two Gaus-
5367 sians. The normalizing constant itself has Gaussian shape, see (6.114).
2. Mean and covariance. We compute the mean and covariance matrix of the marginal likelihood by exploiting the standard results for means and covariances of affine transformations of random variables, see Section 6.4.4. The mean of the marginal likelihood is computed as

$$\mathbb{E}_{\boldsymbol{\theta}}[\mathbf{y} | \mathbf{X}] = \mathbb{E}_{\boldsymbol{\theta}}[\mathbf{X}\boldsymbol{\theta} + \boldsymbol{\epsilon}] = \mathbf{X}\mathbb{E}_{\boldsymbol{\theta}}[\boldsymbol{\theta}] = \mathbf{X}\mathbf{m}_0. \quad (9.63)$$

Note that $\boldsymbol{\epsilon} \sim \mathcal{N}(\mathbf{0}, \sigma^2 \mathbf{I})$ is a vector of i.i.d. random variables. The covariance matrix is given as

$$\text{Cov}_{\boldsymbol{\theta}}[\mathbf{y}] = \text{Cov}[\mathbf{X}\boldsymbol{\theta}] + \sigma^2 \mathbf{I} = \mathbf{X} \text{Cov}_{\boldsymbol{\theta}}[\boldsymbol{\theta}] \mathbf{X}^{\top} + \sigma^2 \mathbf{I} \quad (9.64a)$$

$$= \mathbf{X} \mathbf{S}_0 \mathbf{X}^{\top} + \sigma^2 \mathbf{I} \quad (9.64b)$$

Hence, the marginal likelihood is

$$\begin{aligned} p(\mathbf{y} | \mathbf{X}) &= (2\pi)^{-\frac{N}{2}} \det(\mathbf{X} \mathbf{S}_0 \mathbf{X}^{\top} + \sigma^2 \mathbf{I})^{-\frac{1}{2}} \\ &\quad \times \exp\left(-\frac{1}{2}(\mathbf{y} - \mathbf{X}\mathbf{m}_0)^{\top}(\mathbf{X} \mathbf{S}_0 \mathbf{X}^{\top} + \sigma^2 \mathbf{I})^{-1}(\mathbf{y} - \mathbf{X}\mathbf{m}_0)\right). \end{aligned} \quad (9.65)$$

5368 The marginal likelihood can now be used for Bayesian model selection as
5369 discussed in Section 8.5.2.

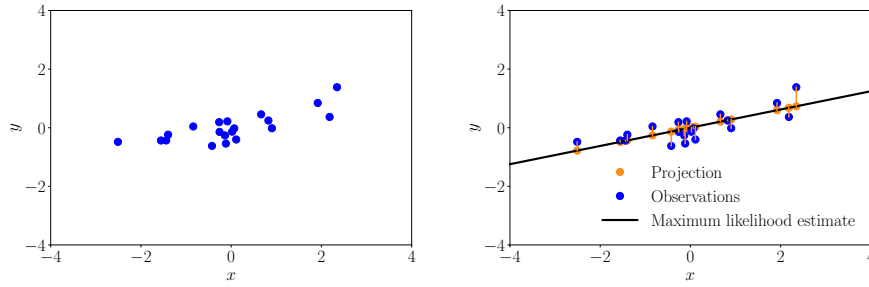
5370 9.4 Maximum Likelihood as Orthogonal Projection

Having crunched through much algebra to derive maximum likelihood and MAP estimates, we will now provide a geometric interpretation of maximum likelihood estimation. Let us consider a simple linear regression setting

$$y = x\theta + \epsilon, \quad \epsilon \sim \mathcal{N}(0, \sigma^2), \quad (9.66)$$

5371 in which we consider linear functions $f : \mathbb{R} \rightarrow \mathbb{R}$ that go through the
5372 origin (we omit features here for clarity). The parameter θ determines the
5373 slope of the line. Figure 9.12(a) shows a one-dimensional dataset.

With a training data set $\mathbf{X} = [x_1, \dots, x_N]^{\top} \in \mathbb{R}^N$, $\mathbf{y} = [y_1, \dots, y_N]^{\top} \in \mathbb{R}^N$



(a) Regression dataset consisting of noisy observations y_n (blue) of function values $f(x_n)$ at input locations x_n .

(b) The orange dots are the projections of the noisy observations (blue dots) onto the line $\theta_{\text{ML}}x$. The maximum likelihood solution to a linear regression problem finds a subspace (line) onto which the overall projection error (orange lines) of the observations is minimized.

Figure 9.12
Geometric interpretation of least squares. (a) Dataset; (b) Maximum likelihood solution interpreted as a projection.

\mathbb{R}^N , we recall the results from Section 9.2.1 and obtain the maximum likelihood estimator for the slope parameter as

$$\theta_{\text{ML}} = (\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{X}^\top \mathbf{y} = \frac{\mathbf{X}^\top \mathbf{y}}{\mathbf{X}^\top \mathbf{X}} \in \mathbb{R}. \quad (9.67)$$

This means for the training inputs \mathbf{X} we obtain the optimal (maximum likelihood) reconstruction of the training data, i.e., the approximation with the minimum least-squares error

$$\mathbf{X}\theta_{\text{ML}} = \mathbf{X} \frac{\mathbf{X}^\top \mathbf{y}}{\mathbf{X}^\top \mathbf{X}} = \frac{\mathbf{X}\mathbf{X}^\top}{\mathbf{X}^\top \mathbf{X}} \mathbf{y}. \quad (9.68)$$

As we are basically looking for a solution of $\mathbf{y} = \mathbf{X}\theta$, we can think of linear regression as a problem for solving systems of linear equations. Therefore, we can relate to concepts from linear algebra and analytic geometry that we discussed in Chapters 2 and 3. In particular, looking carefully at (9.68) we see that the maximum likelihood estimator θ_{ML} in our example from (9.66) effectively does an orthogonal projection of \mathbf{y} onto the one-dimensional subspace spanned by \mathbf{X} . Recalling the results on orthogonal projections from Section 3.7, we identify $\frac{\mathbf{X}\mathbf{X}^\top}{\mathbf{X}^\top \mathbf{X}}$ as the projection matrix, θ_{ML} as the coordinates of the projection onto the one-dimensional subspace of \mathbb{R}^N spanned by \mathbf{X} and $\mathbf{X}\theta_{\text{ML}}$ as the orthogonal projection of \mathbf{y} onto this subspace.

Therefore, the maximum likelihood solution provides also a geometrically optimal solution by finding the vectors in the subspace spanned by \mathbf{X} that are “closest” to the corresponding observations \mathbf{y} , where “closest” means the smallest (squared) distance of the function values y_n to $x_n\theta$. This is achieved by orthogonal projections. Figure 9.12(b) shows the orthogonal projection of the noisy observations onto the subspace that

Linear regression can be thought of as a method for solving systems of linear equations. Maximum likelihood linear regression performs an orthogonal projection.

5391 minimizes the squared distance between the original dataset and its pro-
 5392 jection, which corresponds to the maximum likelihood solution.

In the general linear regression case where

$$y = \phi^\top(x)\theta + \epsilon, \quad \epsilon \sim \mathcal{N}(0, \sigma^2) \quad (9.69)$$

with vector-valued features $\phi(x) \in \mathbb{R}^K$, we again can interpret the maximum likelihood result

$$\mathbf{y} \approx \Phi \theta_{\text{ML}}, \quad (9.70)$$

$$\theta_{\text{ML}} = \Phi(\Phi^\top \Phi)^{-1} \Phi^\top \mathbf{y} \quad (9.71)$$

5393 as a projection onto a K -dimensional subspace of \mathbb{R}^N , which is spanned
 5394 by the columns of the feature matrix Φ , see Section 3.7.2.

If the feature functions ϕ_k that we use to construct the feature matrix Φ are orthonormal (see Section 3.6), we obtain a special case where the columns of Φ form an orthonormal basis (see Section 3.5), such that $\Phi^\top \Phi = \mathbf{I}$. This will then lead to the projection

$$\Phi(\Phi^\top \Phi)^{-1} \Phi \mathbf{y} = \Phi \Phi^\top \mathbf{y} = \left(\sum_{k=1}^K \phi_k \phi_k^\top \right) \mathbf{y} \quad (9.72)$$

5395 so that the coupling between different features has disappeared and the
 5396 maximum likelihood projection is simply the sum of projections of \mathbf{y} onto
 5397 the individual basis vectors ϕ_k , i.e., the columns of Φ . Many popular basis
 5398 functions in signal processing, such as wavelets and Fourier bases, are
 5399 orthogonal basis functions. When the basis is not orthogonal, one can
 5400 convert a set of linearly independent basis functions to an orthogonal basis
 5401 by using the Gram-Schmidt process (Strang, 2003).

9.5 Further Reading

In this chapter, we discussed linear regression for Gaussian likelihoods and conjugate Gaussian priors on the parameters of the model. This allowed for closed-form Bayesian inference. However, in some applications we may want to choose a different likelihood function. For example, in a binary *classification* setting, we observe only two possible (categorical) outcomes, and a Gaussian likelihood is inappropriate in this setting. Instead, we can choose a Bernoulli likelihood that will return a probability of the predicted label to be 1 (or 0). We refer to the books by Bishop (2006); Murphy (2012); Barber (2012) for an in-depth introduction to classification problems. A different example where non-Gaussian likelihoods are important is count data. Counts are non-negative integers, and in this case a Binomial or Poisson likelihood would be a better choice than a Gaussian. All these examples fall into the category of *generalized linear models*, a flexible generalization of linear regression that allows for response variables that have error distribution models other than a Gaussian

distribution. The GLM generalizes linear regression by allowing the linear model to be related to the observed values via a smooth and invertible function $\sigma(\cdot)$ that may be nonlinear so that $y = \sigma(f)$, where $f = \theta^\top \phi(x)$ is the linear regression model from (9.12). We can therefore think of a generalized linear model in terms of function composition $y = \sigma \circ f$ where f is a linear regression model and σ the activation function. Note, that although we are talking about “generalized linear models” the outputs y are no longer linear in the parameters θ . In *logistic regression*, we choose the *logistic sigmoid* $\sigma(f) = \frac{1}{1+\exp(-f)} \in [0, 1]$, which can be interpreted as the probability of observing a binary output $y = 1$ of a Bernoulli random variable. The function $\sigma(\cdot)$ is called *transfer function* or *activation function*, its inverse is called the *canonical link function*. From this perspective, it is also clear that generalized linear models are the building blocks of (deep) feedforward neural networks: If we consider a generalized linear model $y = \sigma(Ax + b)$, where A is a weight matrix and b a bias vector, we identify this generalized linear model as a single-layer neural network with activation function $\sigma(\cdot)$. We can now recursively compose these functions via

$$\begin{aligned} x_{k+1} &= f_k(x_k) \\ f_k(x_k) &= \sigma_k(A_k x_k + b_k) \end{aligned} \quad (9.73)$$

for $k = 0, \dots, K - 1$ where x_0 are the input features and $x_K = y$ are the observed outputs, such that $f_{K-1} \circ \dots \circ f_0$ is a K -layer deep neural network. Therefore, the building blocks of this deep neural network are the generalized linear models defined in (9.73). A great post on the relation between GLMs and deep networks is available at <https://tinyurl.com/glm-dnn>. Neural networks (Bishop, 1995; Goodfellow et al., 2016) are significantly more expressive and flexible than linear regression models. However, maximum likelihood parameter estimation is a non-convex optimization problem, and marginalization of the parameters in a fully Bayesian setting is analytically intractable.

We briefly hinted at the fact that a distribution over parameters induces a distribution over regression functions. *Gaussian processes* (Rasmussen and Williams, 2006) are regression models where the concept of a distribution over function is central. Instead of placing a distribution over parameters a Gaussian process places a distribution directly on the space of functions without the “detour” via the parameters. To do so, the Gaussian process exploits the *kernel trick* (Schölkopf and Smola, 2002), which allows us to compute inner products between two function values $f(x_i), f(x_j)$ only by looking at the corresponding input x_i, x_j . A Gaussian process is closely related to both Bayesian linear regression and support vector regression but can also be interpreted as a Bayesian neural network with a single hidden layer where the number of units tends to infinity (Neal, 1996; Williams, 1997). An excellent introduction to Gaus-

logistic regression

logistic sigmoid

transfer function

activation function

canonical link

function

For ordinary linear

regression the

activation function

would simply be the

identity.

Generalized linear

models are the

building blocks of

deep neural

networks.

Gaussian processes

kernel trick

sian processes can be found in (MacKay, 1998; Rasmussen and Williams, 2006).

We focused on Gaussian parameter priors in the discussions in this chapters because they allow for closed-form inference in linear regression models. However, even in a regression setting with Gaussian likelihoods we may choose a non-Gaussian prior. Consider a setting where the inputs are $\mathbf{x} \in \mathbb{R}^D$ and our training set is small and of size $N \ll D$. This means that the regression problem is under-determined. In this case, we can choose a parameter prior that enforces sparsity, i.e., a prior that tries to set as many parameters to 0 as possible (*variable selection*). This prior provides a stronger regularizer than the Gaussian prior, which often leads to an increased prediction accuracy and interpretability of the model. The Laplace prior is one example that is frequently used for this purpose. A linear regression model with the Laplace prior on the parameters is equivalent to linear regression with L1 regularization (*LASSO*) (Tibshirani, 1996). The Laplace distribution is sharply peaked at zero (its first derivative is discontinuous) and it concentrates its probability mass closer to zero than the Gaussian distribution, which encourages parameters to be 0. Therefore, the non-zero parameters are relevant for the regression problem, which is the reason why we also speak of “variable selection”.

variable selection

LASSO