

LECTURE NOTES

Machine Learning Essentials

TUM, WiSe 2025/26

Prof. Dr. Julija Zavadlav

Last update: October 30, 2025

Disclaimers

The notes are only informally distributed and intended ONLY as a study aid for the TUM students enrolled in the course Machine Learning Essentials WiSe 2025/26. The script may contain errors, so use it with care. Sharing and distribution of any course content, other than between individual students registered in the course, is not permitted without permission.

Supervised learning: $\mathcal{D} = \{(\vec{x}_n, \vec{y}_n)\}_{n=1}^N$

Regression: \vec{y}_n is continuous

Linear Basis Function model: $f(\vec{x}, \vec{w}) = \vec{w} \vec{\phi}(\vec{x})$

Maximum Likelihood: $\vec{\theta}_{ML} = \underset{\vec{\theta}}{\operatorname{arg\,max}} p(D|\vec{\theta})$

LECTURE

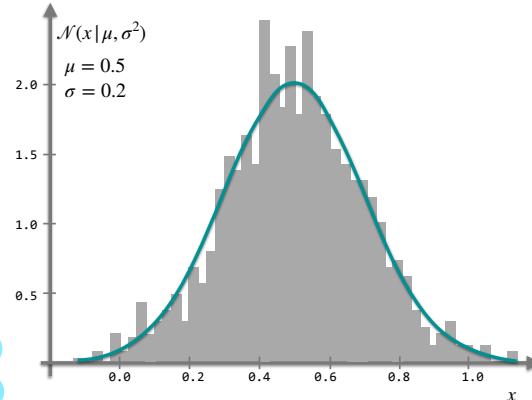
3 Regression Problem

In this lecture, we will look at the regression problem. We will use the three steps we have introduced last time: (1) Specification of the Model Architecture, (2) Parameter Estimation/ Training, and (3) Prediction of Novel Inputs. Our goal is to derive the loss function for a regression problem by using the maximum likelihood "measure of best" (frequentist approach) and assuming Gaussian noise. We will then see how this loss function has a closed-form solution when we further assume a linear basis function model. This solution is, in fact, equal to the Linear Least Squares - a method employed when you run a "fit" command in MATLAB or other software.

3.1 Gaussian Probability Distribution

The Gaussian distribution, also known as the normal distribution, is one of the most widely used distributions in probability theory and appears in numerous contexts. It was invented by Carl Friedrich Gauss, who also demonstrated that least-squares can be derived under the assumption of Gaussian errors (which we will show in this lecture).

Consider, for example, the problem of determining the resistance by measuring the current at different voltages. The measurement device (multimeter) has some intrinsic, random error (inherent to the device itself). Therefore, if we measure the current at the same applied voltage, we will not always obtain the same result. By recording the results in a histogram, we would obtain a distribution that tends toward a Gaussian distribution. Thus, we can often assume that the random noise in the system is Gaussian distributed.



The Gaussian distribution also arises in the central limit theorem (also known as the law of large numbers), which states that the sum of random variables has a distribution that becomes increasingly Gaussian as the number of terms in the sum increases.

Univariate Gaussian

The univariate Gaussian distribution is determined by two parameters μ and σ^2 . The probability density function of the scalar variable x is given by

$$\mathcal{N}(x | \mu, \sigma^2) = \frac{1}{\sqrt{2\pi}\sigma} \exp\left[-\frac{1}{2\sigma^2}(x-\mu)^2\right] \quad (3.1)$$

Univariate Gaussian

We can see that the Gaussian distribution satisfies the properties of a probability distribution function, i.e.,

$$\begin{aligned} \mathcal{N}(x | \mu, \sigma^2) &> 0 \\ \int_{-\infty}^{\infty} \mathcal{N}(x | \mu, \sigma^2) dx &= 1 \quad \leftarrow \text{normalization: factor } \frac{1}{\sqrt{2\pi}\sigma} \end{aligned}$$

The expectation of x under the Gaussian distribution is $\mathbb{E}[f(x)] = \int f(x) p(x) dx$

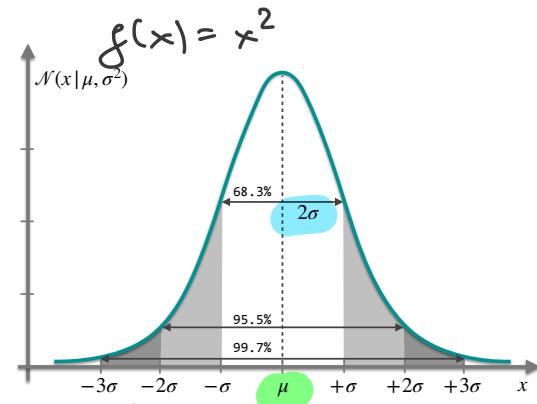
$$\mathbb{E}[x] = \int_{-\infty}^{\infty} \mathcal{N}(x | \mu, \sigma^2) x dx = \mu \quad (3.2)$$

Thus μ is the average value of x and thus referred to as the **mean**. Graphically, the mean corresponds to the position of the peak. If we compute the expectation of x^2 , we obtain

$$\mathbb{E}[x^2] = \int_{-\infty}^{\infty} \mathcal{N}(x | \mu, \sigma^2) x^2 dx = \mu^2 + \sigma^2. \quad (3.3)$$

Thus, the **variance** of x equals

$$\text{var}[x] = \mathbb{E}[x^2] - \mathbb{E}[x]^2 = \sigma^2, \quad (3.4)$$



and the σ^2 is referred to as the **variance** parameter. The square root of variance σ is called the **standard deviation**. Graphically, the standard deviation encodes the width of the curve with full width at half maximum FWHM = 2.355σ .

Multivariate Gaussian

We can generalize the distribution for a D -dimensional vector $\mathbf{x} \in \mathbb{R}^D$ of continuous variables, i.e.,

$$\mathcal{N}(\vec{x} | \vec{\mu}, \Sigma) = \frac{1}{(2\pi)^{D/2}} \frac{1}{\sqrt{|\Sigma|}} \exp\left[-\frac{1}{2} (\vec{x} - \vec{\mu})^\top \Sigma^{-1} (\vec{x} - \vec{\mu})\right] \quad (3.5)$$

Multivariate Gaussian

where the determinant of a matrix Σ is denoted by $|\Sigma|$, and its inverse by Σ^{-1} . $\mu \in \mathbb{R}^D$ is a vector of means and $\Sigma \in \mathbb{R}^{D \times D}$ is the **covariance matrix** since the expected values are now

$$\mathbb{E}[\mathbf{x}] = \mu \quad \text{and} \quad \mathbb{E}[\mathbf{x}\mathbf{x}^\top] = \mu\mu^\top + \Sigma \quad \rightarrow \quad \text{var}[\mathbf{x}] = \Sigma \quad (3.6)$$

The diagonal and off-diagonal components of the covariance matrix are

$$\Sigma_{ii} = \text{var}[x_i] \quad \text{and} \quad \Sigma_{ij} = \text{cov}[x_i, x_j]. \quad (3.7)$$

Geometric Interpretation

The Gaussian distribution is constant on surfaces for which the exponent is constant

$$\Delta^2 = (\mathbf{x} - \boldsymbol{\mu})^T \Sigma^{-1} (\mathbf{x} - \boldsymbol{\mu}) \quad (3.8)$$

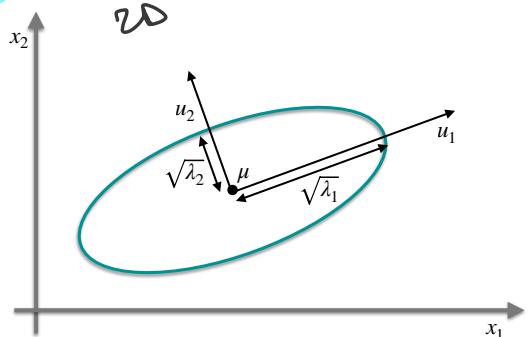
The quantity Δ is called the Mahalanobis distance from \mathbf{x} to $\boldsymbol{\mu}$ and reduces to Euclidean distance when Σ is the identity matrix. To understand how Δ looks in \mathbf{x} space, we can consider the eigenvector equation for the covariance matrix

$$\Sigma \mathbf{u}_i = \lambda_i \mathbf{u}_i \quad \text{or} \quad \Sigma U = U \Lambda, \quad (3.9)$$

where U is a matrix whose columns are eigenvectors \mathbf{u}_i and $\Lambda = \text{diag}(\lambda_1, \dots, \lambda_D)$. Since Σ is symmetric, the eigenvalues are real and eigenvectors can be chosen to form an orthonormal set, i.e., $U^T = U^{-1}$. Right-hand multiplication of the eigenvector equation by U^T gives

$$\Sigma U = U \Lambda / U^T$$

$$\begin{aligned} \cancel{\Sigma} U U^T &= U \Lambda U^T \\ \cancel{\Sigma} &= U \Lambda U^T = \sum_{i=1}^D \lambda_i \mathbf{u}_i \mathbf{u}_i^T \end{aligned}$$



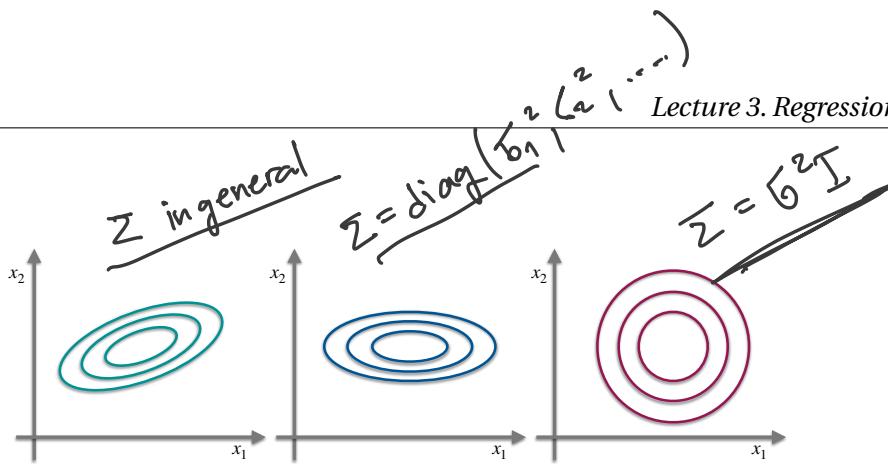
The inverse covariance matrix can be expressed as

$$\Sigma^{-1} = (U \Lambda U^T)^{-1} = (U \Lambda U^{-1})^{-1} = U \Lambda^{-1} U^{-1} = U \Lambda^{-1} U^T = \sum_{i=1}^D \frac{1}{\lambda_i} \mathbf{u}_i \mathbf{u}_i^T \quad (3.10)$$

Substituting this result into the Mahalanobis distance, we obtain

$$\begin{aligned} \Delta^2 &= (\mathbf{x} - \boldsymbol{\mu})^T \Sigma^{-1} (\mathbf{x} - \boldsymbol{\mu}) \\ &= (\mathbf{x} - \boldsymbol{\mu})^T \left(\sum_{i=1}^D \frac{1}{\lambda_i} \mathbf{u}_i \mathbf{u}_i^T \right) (\mathbf{x} - \boldsymbol{\mu}) = \sum_{i=1}^D \frac{(\mathbf{x} - \boldsymbol{\mu})^T \mathbf{u}_i \mathbf{u}_i^T (\mathbf{x} - \boldsymbol{\mu})}{\lambda_i} \\ &= \sum_{i=1}^D \frac{y_i^2}{\lambda_i} \quad \text{where } y_i = \mathbf{u}_i^T (\mathbf{x} - \boldsymbol{\mu}) \quad \text{or } \mathbf{y} = U(\mathbf{x} - \boldsymbol{\mu}), \quad \text{rotation matrix} \end{aligned}$$

where $\mathbf{y} = (y_1, \dots, y_D)^T$. We interpret $\{y_i\}$ as a new coordinate system defined by orthonormal vectors \mathbf{u}_i that are shifted and rotated with respect to the original $\{x_i\}$ coordinates. The Gaussian distribution is constant on ellipsoids (distribution can only be properly normalized with positive eigenvalues) with the center at $\boldsymbol{\mu}$ and axes oriented along \mathbf{u}_i . The scaling factors in the directions of the axes are given by $\sqrt{\lambda_i}$. If the covariance matrix is diagonal, i.e., $\Sigma = \text{diag}(\sigma_1^2, \dots, \sigma_D^2)$, then $U = I$ and constant density contours are axis-aligned ellipsoids. If the covariance matrix is isotropic, $\Sigma = \sigma^2 I$, the contours are spherical.



3.2 Frequentist Approach to Regression

Statistical Model

Suppose we have observed some data \mathcal{D} , which consists of N observations of the form $\{\mathbf{x}_n, y_n\}_{n=1}^N$. For simplicity, we have assumed a vector input \mathbf{x}_n and a scalar output variable y_n , but the problem can be generalized for a vector output. We postulate the following statistical model

$$\underline{\text{assume:}} \quad y = f(\vec{x}, \vec{w}) + \epsilon \quad \text{deterministic function} \quad (3.11)$$

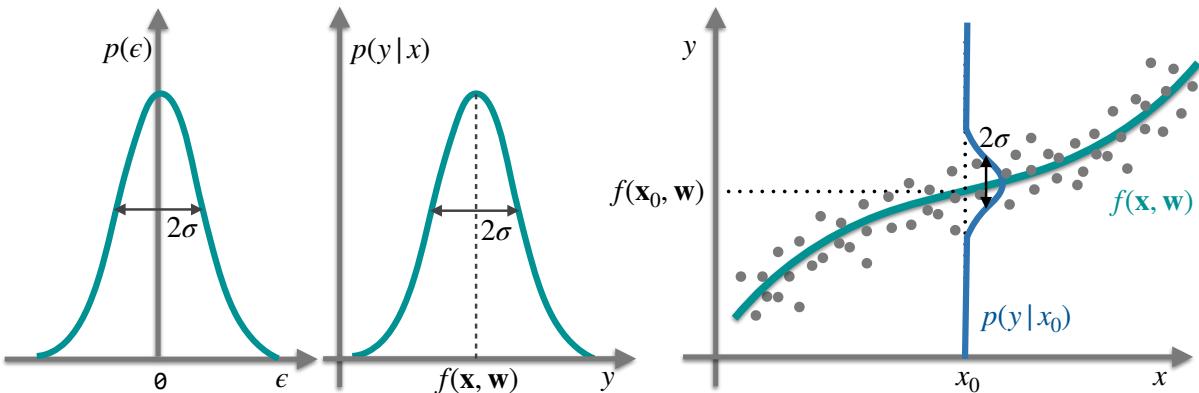
where $f(\mathbf{x}, \mathbf{w})$ is a deterministic function and ϵ is a noise/error term that describes everything that the model cannot capture. We assume that ϵ is a zero-mean Gaussian random variable with variance σ^2 , i.e.,

$$\underline{\text{assume:}} \quad p(\epsilon) = \mathcal{N}(\epsilon | 0, \sigma^2) \quad (3.12)$$

From Eqs 3.11 and 3.12, it follows that the probability distribution $p(y|\mathbf{x})$ is also a Gaussian distribution with a mean equal to $f(\mathbf{x}, \mathbf{w})$ and variance σ^2 , i.e.

$$p(y | \vec{x}) = \mathcal{N}(y | f(\vec{x}, \vec{w}), \sigma^2) \quad (3.13)$$

The parameters of the statistical model are $\vec{\theta} = [\vec{w}, \sigma]$. While a Gaussian random variable is a reasonable choice in many situations, it may not be the best choice in specific settings. Examples where this assumption should be replaced include discrete data, small datasets (Student's t-distribution), and rare events, such as the time until an earthquake occurs (Exponential distribution).



Likelihood function

likelihood

Now, we would like to determine the maximum likelihood estimate $\hat{\theta}_{ML} = \operatorname{argmax}_{\theta} p(\mathcal{D}|\theta)$ of this statistical model. We start by writing the likelihood function

$$\mathcal{D} = \{(\vec{x}_n, y_n)\}_{n=1}^N \quad p(\mathcal{D}|\vec{\theta}) = p(\vec{y} | \vec{X}, \vec{w}, \sigma^2) \quad (3.14)$$

where we have grouped together N observations into a column vector of inputs $\vec{X} = [\vec{x}_1, \dots, \vec{x}_N]^T$ and the corresponding outputs $\vec{y} = [y_1, \dots, y_N]^T$.

We assume that these N observations are independent and identically distributed (i.i.d.). Each observation is distributed according to Eq. (3.13). Since the joint probability of independent events is equal to the product of individual events, we can write the likelihood function as a product of Gaussian functions

$$p(\mathcal{D}|\vec{\theta}) = \prod_{n=1}^N N(y_n | f(\vec{x}_n, \vec{w}), \sigma^2) \quad (3.15)$$

We now switch to the log-likelihood function, as applying the logarithm does not change the location of the maximum (although the value at the maximum will be different). Employing the logarithm product and power rules, we obtain

$$\begin{aligned} \ln p(\mathcal{D}|\vec{\theta}) &= \ln \prod_{n=1}^N N(y_n | f(\vec{x}_n, \vec{w}), \sigma^2) \\ &= \sum_{n=1}^N \ln N(y_n | f(\vec{x}_n, \vec{w}), \sigma^2) \\ &= \sum_{n=1}^N \ln \left\{ \frac{1}{\sqrt{2\pi\sigma^2}} \exp \left[-\frac{1}{2\sigma^2} (y_n - f(\vec{x}_n, \vec{w}))^2 \right] \right\} \\ &= \ln \left(\frac{1}{\sqrt{2\pi\sigma^2}} \right)^N - \frac{1}{2\sigma^2} \sum_{n=1}^N (y_n - f(\vec{x}_n, \vec{w}))^2 \\ &= \frac{N}{2} \ln(2\pi) - N \ln \sigma^2 - \frac{1}{2\sigma^2} \sum_{n=1}^N (y_n - f(\vec{x}_n, \vec{w}))^2 \end{aligned} \quad \left| \begin{array}{l} \text{log-trick!} \\ \ln(a+b) = \ln a + \ln b \\ \ln a^b = b \ln a \end{array} \right. \quad (3.16)$$

To find the $\hat{\vec{w}}_{ML}$, we need to determine the maximum of the log-likelihood function, or equivalently, the minimum of the negative log-likelihood (loss) function.

$$\hat{\vec{w}}_{ML} = \operatorname{argmin}_{\vec{w}} \left[\frac{N}{2} \ln(2\pi) + N \ln \sigma + \underbrace{\frac{1}{2\sigma^2} \sum_{n=1}^N [y_n - f(\vec{x}_n, \vec{w})]^2}_{J(\vec{w})} \right] = \operatorname{argmin}_{\vec{\theta}} L(\vec{w}), \quad (3.17)$$

We see that only the last term depends on \mathbf{w} , thus we can disregard the other two terms. The prefactor $\frac{1}{2\sigma^2}$ will not change the position of the maximum and can be skipped or replaced by 1/2. With these considerations in mind, we end up with the loss function called the sum-of-squares error function

$$\mathcal{L}(\vec{\mathbf{w}}) = \sum_{n=1}^N [y_n - f(\vec{x}_n, \vec{\mathbf{w}})]^2 \quad (3.18)$$

Sum-of-Squares Error

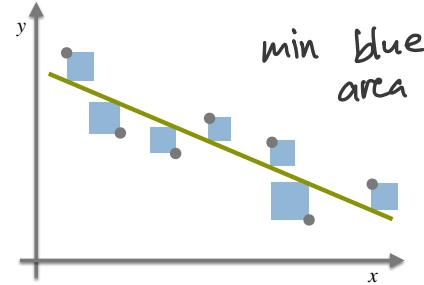
which we need to minimize. We could also normalize the loss by the number of data points, which would again not change the position of the minimum. In this case, we obtain the so-called Mean Squared Error (MSE)

$$\mathcal{L}(\vec{\mathbf{w}}) = \frac{1}{N} \sum_{n=1}^N [y_n - f(\vec{x}_n, \vec{\mathbf{w}})]^2 \quad (3.19)$$

Mean Squared Error (MSE)

Note that we discuss the loss function here because it is used to optimize the parameters. We could also use the same functional form to evaluate the model, and then we would talk about error functions.

The Least Squares method, therefore, is minimizing the squared L2-norm ($\|\mathbf{x}\|_2^2 = \sum_i x_i^2$) of the squared distance between the observed value y_n and the output of the deterministic function $f(\mathbf{x}_n, \mathbf{w})$ (visualized as squares). The optimal parameters are those for which the amount of blue color is minimized.



We can also determine the variance σ^2 parameter of the statistical model. Differentiating Eq. (3.16) with respect to σ and equating to zero we obtain

$$0 = -\frac{N}{\sigma} + \frac{1}{\sigma^3} \sum_{n=1}^N [y_n - f(\mathbf{x}_n, \mathbf{w})]^2$$

$$\hat{\sigma}_{\text{ML}}^2 = \frac{1}{N} \sum_{n=1}^N [y_n - f(\mathbf{x}_n, \hat{\mathbf{w}}_{\text{ML}})]^2.$$

Prediction of Novel Inputs

Once we have determined the parameters of the statistical model $\hat{\theta}_{\text{ML}} = [\hat{\mathbf{w}}_{\text{ML}}, \hat{\sigma}_{\text{ML}}]^T$ we make a prediction for a new input \mathbf{x}_* with

$$y_* = \mathcal{N}(y | f(\mathbf{x}_*, \hat{\mathbf{w}}_{\text{ML}}), \hat{\sigma}_{\text{ML}}^2). \quad (3.20)$$

Note that the prediction is now a probability distribution.

Demo: Least Squares & Outliers

To better understand the Least Squares method, see demo at

<https://mless.pythonanywhere.com/leastsquares> where we investigate the impact of adding an outlier to the training data set. Given a training dataset $\{(x_n, y_n)\}_{n=1}^{N=10}$ we aim to find $f(x, \mathbf{w}) = w_0 + w_1 x + w_2 x^2 + w_3 x^3$ which is an approximation of a true function $g(x) = \sin(2\pi x)$. Data points are drawn from $y = \sin(2\pi x) + \mathcal{N}(0, \sigma^2 = 0.01)$ while the outlier data point is drawn from $y = \sin(2\pi x) + \mathcal{N}(0, \sigma^2 = 0.25)$. Answer the following questions:

- How does the addition of an outlier affect the Least Squares fit? Is (or when is) the Least Squares method sensitive to outliers?

the result is different, critical for small datasets

- How can we modify the "Measure of Best" to obtain a method that is less sensitive to outliers?

$$\|(\hat{\mathbf{w}}) - \sum \|y_n - f(x_n, \mathbf{w})\|_1 \text{-norm}$$

- In what situation could we have an outlier in the dataset?

sim / experiment gone wrong,

- How can we avoid outliers in the dataset?

try to visualise! look at data points with large errors

Summary

- assumptions: gaussian noise, i.i.d data
- we did not say anything about $f(\vec{x}, \vec{w})$ 
- If $f(\vec{x}, \vec{w})$ is LBFM \rightarrow closed-form solution = linear LS method

3.3 Linear Least Squares

If we assume a linear basis function model for the deterministic function $f(\mathbf{x}, \mathbf{w}) = \mathbf{w}^T \phi(\mathbf{x})$, we can find a closed-form solution for the regression problem. For each data sample in the dataset $\mathcal{D} = \{(\mathbf{x}_n, y_n)\}_{n=1}^N$, we have

$$\begin{aligned} y_n &= f(\mathbf{x}_n, \mathbf{w}) + \epsilon_n \quad \forall n = 1, \dots, N \\ &= \mathbf{w}^T \phi(\mathbf{x}_n) + \epsilon_n \\ &= w_0 + \sum_{m=1}^{M-1} w_m \phi_m(\mathbf{x}_n) + \epsilon_n, \end{aligned} \quad \text{N equations}$$

where \mathbf{w} is an M -dimensional vector.

Matrix Formulation

We can rewrite the above N equations jointly using a matrix formulation

$$\mathbf{y} = \Phi \mathbf{w} + \boldsymbol{\epsilon} \quad (3.21)$$

$$y_1 = \underbrace{\phi_0(x_1) w_0}_1 + \underbrace{\phi_1(x_1) w_1}_2 + \dots + \underbrace{\phi_{M-1}(x_1) w_{M-1}}_7 + \epsilon_1$$

or

$$\underbrace{\mathbf{y} \in \mathbb{R}^N}_{\text{each row is 1 data example}} = \underbrace{\Phi \in \mathbb{R}^{N \times M}}_{\substack{\text{columns basis functions} \\ \text{basis functions}}} \underbrace{\mathbf{w} \in \mathbb{R}^M}_{\text{fiting basis function}} + \underbrace{\boldsymbol{\epsilon} \in \mathbb{R}^N}_{\epsilon_1 \epsilon_2 \dots \epsilon_N}, \quad (3.22)$$

where Φ is called a design matrix. Here, each row corresponds to one data sample.

Example: Matrix Formulation for Quadratic Polynomial

Let us consider a polynomial basis set, i.e., $\phi_m(x) = x^m$ and $M = 3$. Furthermore, we consider the following dataset:

	x_n	y_n
n=1	1	2
n=2	2	9

$$f(x, \vec{w}) = \overline{w_0 + w_1 x + w_2 x^2} \quad \begin{cases} \phi_0(x) = 1 \\ \phi_1(x) = x \\ \phi_2(x) = x^2 \end{cases} \quad M=3$$

$$\begin{bmatrix} 2 \\ 9 \end{bmatrix} = \begin{bmatrix} 1 & 1 & 1 \\ 1 & 2 & 4 \end{bmatrix} \begin{bmatrix} w_0 \\ w_1 \\ w_2 \end{bmatrix} + \begin{bmatrix} \epsilon_1 \\ \epsilon_2 \end{bmatrix}$$

Closed-form Solution

We want to find the Maximum Likelihood Estimate for this problem. We saw that assuming Gaussian noise and independent and identically distributed (i.i.d.) observations boils down to minimizing the

sum-of-squares error function, i.e.

$$\begin{aligned}
 E(\mathbf{w}) &= \sum_{n=1}^N [y_n - f(\mathbf{x}_n, \mathbf{w})]^2 \\
 &= \sum_{n=1}^N [y_n - \mathbf{w}^T \Phi(\mathbf{x}_n)]^2 \\
 &= (\mathbf{y} - \Phi \mathbf{w})^T (\mathbf{y} - \Phi \mathbf{w}) \\
 &= (\mathbf{y}^T - \mathbf{w}^T \Phi^T)(\mathbf{y} - \Phi \mathbf{w}) \\
 &= [\mathbf{y}^T \mathbf{y} - \underbrace{\mathbf{y}^T \Phi \mathbf{w}}_{\mathbf{w}^T \Phi^T \mathbf{y}} - \underbrace{\mathbf{w}^T \Phi^T \mathbf{y}}_{\mathbf{w}^T \Phi^T \Phi \mathbf{w}} + \mathbf{w}^T \Phi^T \Phi \mathbf{w}] \\
 &= [\mathbf{y}^T \mathbf{y} - 2\mathbf{w}^T \Phi^T \mathbf{y} + \mathbf{w}^T \Phi^T \Phi \mathbf{w}].
 \end{aligned}
 \tag{3.23}$$

In the last line we have used the fact that $\mathbf{y}^T \Phi \mathbf{w}$ is a $(1 \times N)(N \times M)(M \times 1)$ so (1×1) matrix which is always symmetric. To obtain a minimum, we compute the derivative with respect to \mathbf{w}

$$\begin{aligned}
 \overset{\text{A}^T \rightarrow A}{\frac{\partial}{\partial \mathbf{w}}} [\mathbf{y}^T \mathbf{y} - 2\mathbf{w}^T \Phi^T \mathbf{y} + \mathbf{w}^T \Phi^T \Phi \mathbf{w}] &= 0 \\
 -2\Phi^T \mathbf{y} + [\Phi^T \Phi + (\Phi^T \Phi)^T] \mathbf{w} &= 0 \\
 -\Phi^T \mathbf{y} + \Phi^T \Phi \mathbf{w} &= 0,
 \end{aligned}$$

where we have used the matrix differentiation rules $\frac{\partial \mathbf{x}^T \mathbf{a}}{\partial \mathbf{x}} = \mathbf{a}$ and $\frac{\partial \mathbf{x}^T A \mathbf{x}}{\partial \mathbf{x}} = (A + A^T) \mathbf{x}$. We obtain the so-called normal equations

$$\Phi^T \Phi \mathbf{w} = \Phi^T \mathbf{y}. \tag{3.24}$$

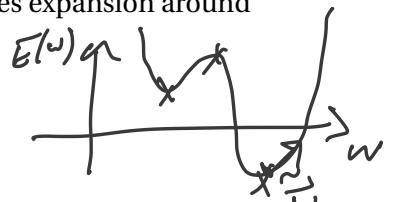
If the basis functions ϕ are linearly independent $\Phi^T \Phi$ is a symmetric, positive definite matrix¹. Thus, it can be inverted, and we can write the least squares solution $\hat{\mathbf{w}}$ as

$$\boxed{\hat{\mathbf{w}}_{ML} = (\Phi^T \Phi)^{-1} \Phi^T \mathbf{y}}
 \tag{3.25}$$

Least Squares Solution

Is the obtained solution truly a minimum? Mathematically, we check if an extremum is a minimum by checking if the second derivative is positive. To see this, we write the Taylor series expansion around $\hat{\mathbf{w}}$:

$$\begin{aligned}
 E(\hat{\mathbf{w}} + \tilde{\mathbf{w}}) &= E(\hat{\mathbf{w}}) + \tilde{\mathbf{w}}^T \cancel{\frac{\partial E(\hat{\mathbf{w}})}{\partial \mathbf{w}}} + \frac{1}{2} \tilde{\mathbf{w}}^T \cancel{\frac{\partial^2 E(\hat{\mathbf{w}})}{\partial^2 \mathbf{w}}} \tilde{\mathbf{w}} \\
 &= E(\hat{\mathbf{w}}) + \frac{1}{2} \tilde{\mathbf{w}}^T \Phi^T \Phi \tilde{\mathbf{w}} > E(\hat{\mathbf{w}}).
 \end{aligned}$$



We used the fact that $\frac{\partial E(\hat{\mathbf{w}})}{\partial \mathbf{w}} = 0$ and $\frac{\partial^2 E(\hat{\mathbf{w}})}{\partial^2 \mathbf{w}} = \Phi^T \Phi$. We conclude that the least square solution $\hat{\mathbf{w}}$ is a minimum since adding a small vector $\tilde{\mathbf{w}}$ will increase E .

¹Symmetry is straight forward $(\Phi^T \Phi)^T = \Phi^T (\Phi^T)^T = \Phi^T \Phi$. Recall that a matrix B is positive definite if for any nonzero vector \mathbf{y} : $\mathbf{y}^T B \mathbf{y} > 0$. To see that this is true for matrix $\Phi^T \Phi$ we write: $\mathbf{y}^T \Phi^T \Phi \mathbf{y} = (\Phi \mathbf{y})^T (\Phi \mathbf{y}) = \|\Phi \mathbf{y}\|_2^2 > 0$.