

Introduction to Machine Learning

Lecture 6 - Regression

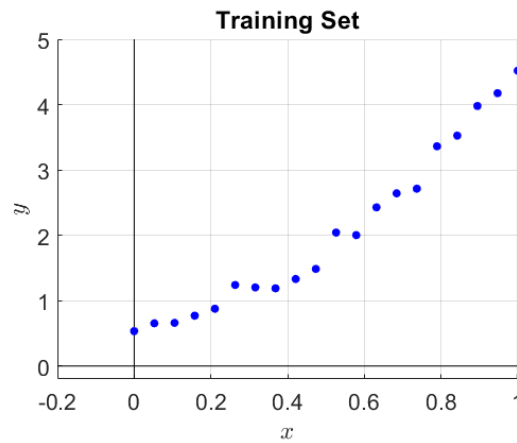
1 Introduction

Consider a training set $\mathcal{D} = \{\mathbf{x}_i, y_i\}_{i=1}^N$ such that:

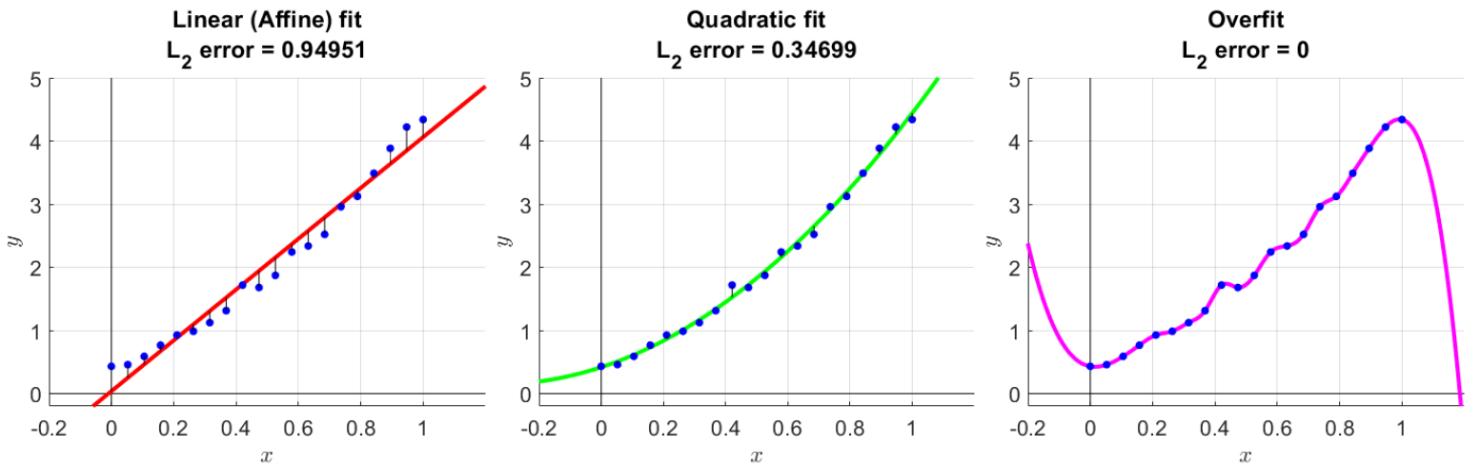
$$y_i = f(\mathbf{x}_i) + \epsilon_i$$

where ϵ_i is assumed to be small and f is an unknown function.

For example:



In the regression problem we search \hat{f} which estimates (fits) the unknown f .
Some possible options to \hat{f} :



The L_2 error is given by:

$$L_2\text{-error} \triangleq \frac{1}{N} \sum_{i=1}^N \left(y_i - \hat{f}(\mathbf{x}_i) \right)^2$$

On the left, we assume a simple model: a linear (affine) function and obtain a relatively high error.

On the right, we assume a complex model which obtains zero error.

However, the quadratic model (middle) which obtains relatively small (but not zero) error seems to be the most reasonable estimation.

2 Least Squares (Linear) Regression

2.1 1D

2.1.1 Linear fit

Let $\mathcal{D} = \{x_i, y_i\}_{i=1}^N$ where $x_i, y_i \in \mathbb{R}$ be the training set.

We assume the following model:

$$\hat{f}_{\text{Linear}}(x) = wx, \quad w \in \mathbb{R}$$

The L_2 error (MSE loss) is given by:

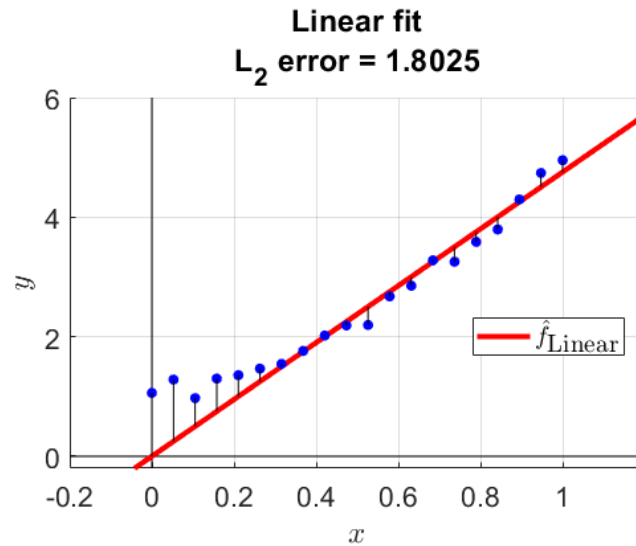
$$L(w) = \frac{1}{N} \sum_{i=1}^N (y_i - \hat{f}(x_i))^2 = \frac{1}{N} \sum_{i=1}^N (y_i - wx_i)^2 = \frac{1}{N} \|\mathbf{y} - w\mathbf{x}\|_2^2$$

where:

$$\mathbf{y} \triangleq \begin{bmatrix} y_1 \\ \vdots \\ y_N \end{bmatrix}, \quad \mathbf{x} \triangleq \begin{bmatrix} x_1 \\ \vdots \\ x_N \end{bmatrix}$$

we can find the optimal w by comparing the derivative to zero:

$$\begin{aligned} \frac{d}{dw} L(w) &= 0 \\ \frac{d}{dw} \left(\frac{1}{N} \|\mathbf{y} - w\mathbf{x}\|_2^2 \right) &= 0 \\ -\frac{2}{N} \mathbf{x}^T (\mathbf{y} - w\mathbf{x}) &= 0 \\ \mathbf{x}^T \mathbf{y} - w \|\mathbf{x}\|_2^2 &= 0 \\ \Rightarrow w &= \frac{\mathbf{x}^T \mathbf{y}}{\|\mathbf{x}\|_2^2} \end{aligned}$$



2.1.2 Affine fit

The linear fit $\hat{f}(x) = wx$ constrains \hat{f} to go through the origin $\hat{f}(0) = 0$. We can use an affine fit to remove this constraint:

$$\hat{f}_{\text{Affine}}(x) = wx + b$$

To obtain the optimal $w \in \mathbb{R}$ and $b \in \mathbb{R}$ (in L_2 error sense) we write:

$$\begin{aligned} \tilde{\mathbf{x}}_i &\triangleq \begin{bmatrix} 1 \\ x_i \end{bmatrix}, & \tilde{\mathbf{w}} &\triangleq \begin{bmatrix} b \\ w \end{bmatrix} \\ \Rightarrow \hat{f}_{\text{Affine}}(x) &= wx + b = \tilde{\mathbf{x}}^T \tilde{\mathbf{w}} \\ \Rightarrow L(w, b) = L(\tilde{\mathbf{w}}) &= \frac{1}{N} \sum_{i=1}^N (y_i - \hat{f}(x_i))^2 = \frac{1}{N} \sum_{i=1}^N (y_i - \tilde{\mathbf{x}}_i^T \tilde{\mathbf{w}})^2 = \frac{1}{N} \|\mathbf{y} - \mathbf{X}\tilde{\mathbf{w}}\|_2^2 \end{aligned}$$

where:

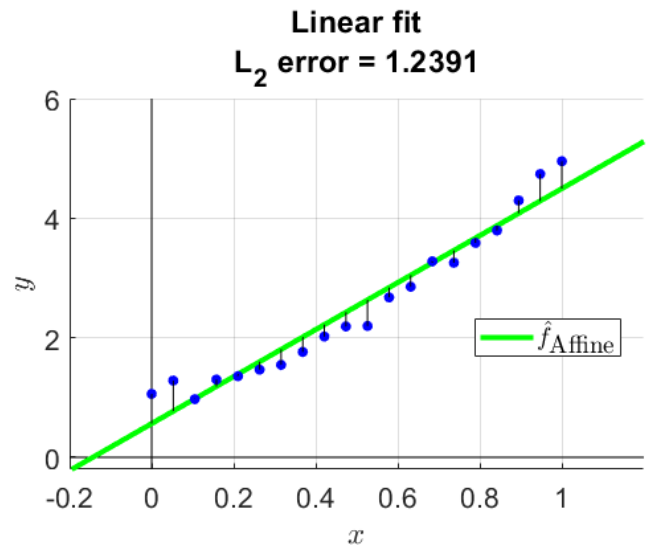
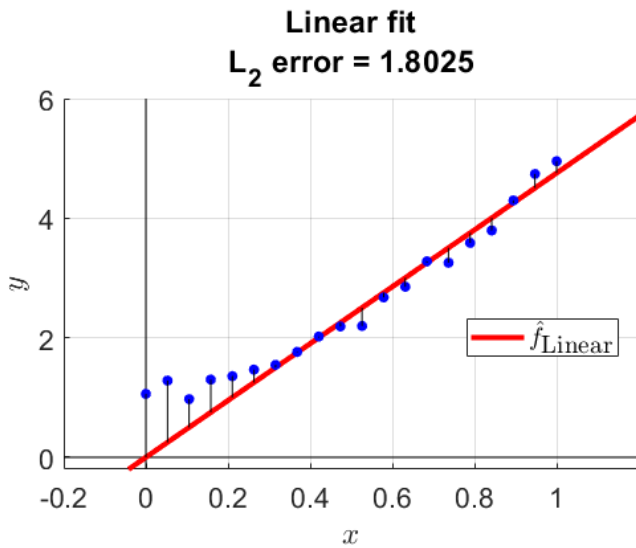
$$\mathbf{X} \triangleq \begin{bmatrix} | & & | \\ \tilde{\mathbf{x}}_1 & \cdots & \tilde{\mathbf{x}}_N \\ | & & | \end{bmatrix}^T = \begin{bmatrix} 1 & x_1 \\ \vdots & \vdots \\ 1 & x_N \end{bmatrix}$$

Comparing the gradient to zero:

$$\begin{aligned} \nabla_{\tilde{\mathbf{w}}} L(\tilde{\mathbf{w}}) &= \mathbf{0} \\ \nabla_{\tilde{\mathbf{w}}} \|\mathbf{y} - \mathbf{X}\tilde{\mathbf{w}}\|_2^2 &= \mathbf{0} \\ \mathbf{X}^T (\mathbf{y} - \mathbf{X}\tilde{\mathbf{w}}) &= \mathbf{0} \\ \mathbf{X}^T \mathbf{X} \tilde{\mathbf{w}} &= \mathbf{X}^T \mathbf{y} \end{aligned}$$

$$\Rightarrow \tilde{\mathbf{w}} = \underbrace{(\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \mathbf{y}}_{\triangleq \mathbf{X}^\dagger} = \mathbf{X}^\dagger \mathbf{y}$$

where \mathbf{X}^\dagger is the Moore–Penrose inverse of \mathbf{X} .



2.1.3 Polyfit

We can assume an M order polynomial model:

$$\hat{f}(x) = w_0 + w_1x + w_2x^2 + \dots + w_Mx^M = \sum_{m=0}^M w_mx^m$$

To obtain the optimal $\{w_m \in \mathbb{R}\}_{m=1}^M$ (in L_2 error sense) we write:

$$\mathbf{w} \triangleq \begin{bmatrix} w_0 \\ w_1 \\ w_2 \\ \vdots \\ w_M \end{bmatrix}, \quad \phi(x) \triangleq \begin{bmatrix} 1 \\ x \\ x^2 \\ \vdots \\ x^M \end{bmatrix}$$

$$\Rightarrow \hat{f}(x) = \phi^T(x) \cdot \mathbf{w}$$

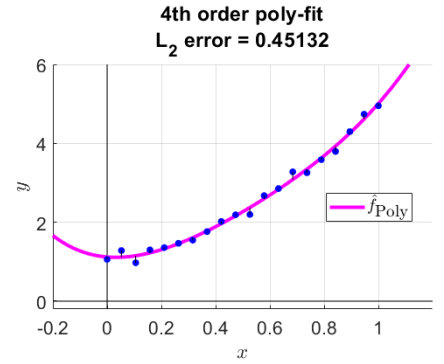
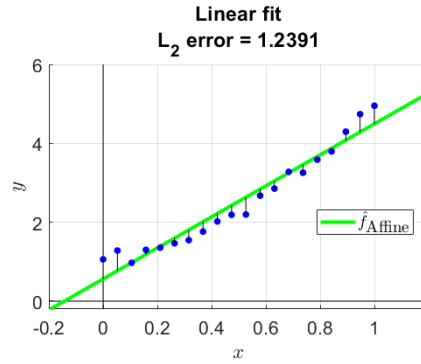
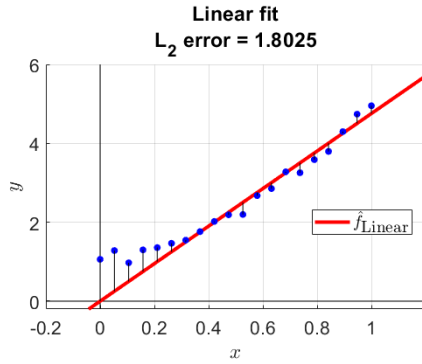
$$\Rightarrow L(\mathbf{w}) = \frac{1}{N} \sum_{i=1}^N (y_i - \hat{f}(x_i))^2 = \frac{1}{N} \sum_{i=1}^N (y_i - \phi^T(x_i) \mathbf{w})^2 = \frac{1}{N} \|\mathbf{y} - \Phi \mathbf{w}\|_2^2$$

where:

$$\Phi \triangleq \begin{bmatrix} \phi(x_1) & \dots & \phi(x_N) \end{bmatrix}^T = \begin{bmatrix} 1 & x_1 & x_1^2 & \dots & x_1^M \\ \vdots & \vdots & \vdots & \dots & \vdots \\ 1 & x_N & x_N^2 & \dots & x_N^M \end{bmatrix}$$

The loss function is the same as before. Thus, the optimal \mathbf{w} is given by:

$$\Rightarrow \mathbf{w} = \left(\Phi^T \Phi \right)^{-1} \Phi^T \mathbf{y} = \Phi^\dagger \mathbf{y}$$



2.1.4 Phase estimation example using feature transform

Consider the following function:

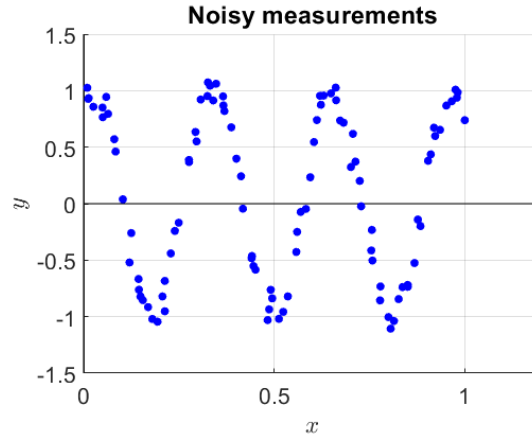
$$f(x) = \sin(\omega_0 x + \theta)$$

where ω_0 is known and θ is the parameter we want to estimate.

We obtain N noisy measurements:

$$y_i = f(x_i) + n_i, \quad i = 1, 2, \dots, N$$

where n_i is some random noise (assume zero mean).



Estimate θ .

Solution:

Note that (trigonometric identity):

$$\sin(\omega_0 x + \theta) = \sin(\omega_0 x) \cos(\theta) + \cos(\omega_0 x) \sin(\theta)$$

and consider the following feature transform:

$$\phi(x) \triangleq \begin{bmatrix} \sin(\omega_0 x) \\ \cos(\omega_0 x) \end{bmatrix}$$

Let us denote:

$$\mathbf{w} \triangleq \begin{bmatrix} \cos(\hat{\theta}) \\ \sin(\hat{\theta}) \end{bmatrix}$$

$$\Rightarrow \hat{f}(x) = \sin(\omega_0 x + \hat{\theta}) = \sin(\omega_0 x) \cos(\hat{\theta}) + \cos(\omega_0 x) \sin(\hat{\theta}) = \phi^T(x) \mathbf{w}$$

As before, the loss function is given by:

$$L(\mathbf{w}) = \frac{1}{N} \sum_{i=1}^N (y_i - \hat{f}(x_i))^2 = \frac{1}{N} \sum_{i=1}^N (y_i - \phi^T(x_i) \mathbf{w})^2 = \frac{1}{N} \|\mathbf{y} - \Phi \mathbf{w}\|_2^2$$

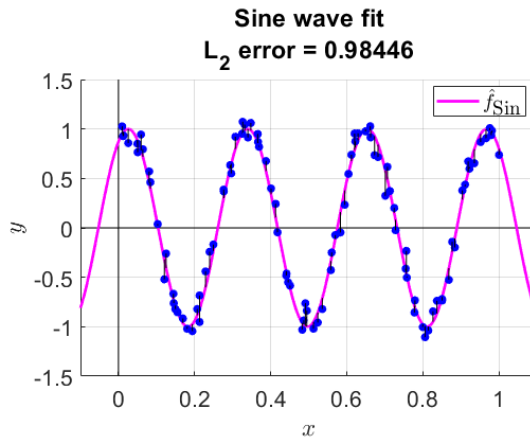
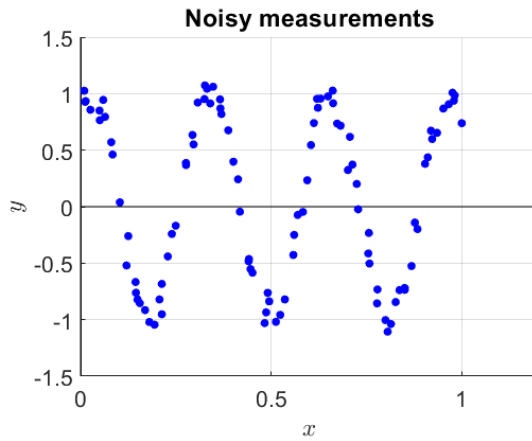
where:

$$\Phi \triangleq \begin{bmatrix} \phi(x_1) & \dots & \phi(x_N) \end{bmatrix}^T = \begin{bmatrix} \sin(\omega_0 x_1) & \cos(\omega_0 x_1) \\ \vdots & \vdots \\ \sin(\omega_0 x_N) & \cos(\omega_0 x_N) \end{bmatrix}$$

Since this is exactly the same loss function, the solution is given by:

$$\begin{bmatrix} \cos(\hat{\theta}) \\ \sin(\hat{\theta}) \end{bmatrix} = \mathbf{w} = \Phi^\dagger \mathbf{y}$$

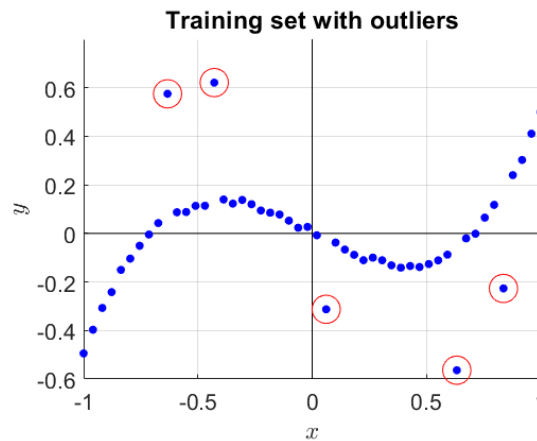
$$\Rightarrow \hat{\theta} = \begin{cases} \arctan\left(\frac{\sin(\hat{\theta})}{\cos(\hat{\theta})}\right) & \cos(\hat{\theta}) \geq 0 \\ \arctan\left(\frac{\sin(\hat{\theta})}{\cos(\hat{\theta})}\right) + \pi & \cos(\hat{\theta}) < 0 \end{cases}$$



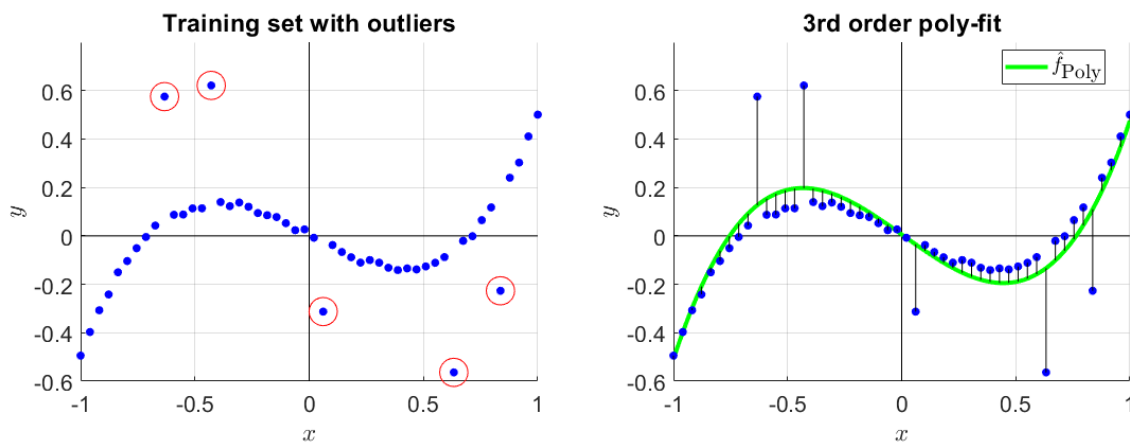
2.2 Outliers and regularization

2.2.1 Outliers

Consider the following training set with few outliers:



If we ignore the outliers the data seems to originate by a 3rd order polynomial. Thus, we can try to apply poly-fit (of order 3):



The fitted curve is minimizing the squared error:

$$L(\mathbf{w}) = \sum_{i=1}^N \left(y_i - \hat{f}(x_i) \right)^2 = \|\mathbf{y} - \Phi\mathbf{w}\|_2^2$$

Thus, the squared errors of the outliers causing the fitted curve to be too twisted (i.e. the values of \mathbf{w} are too big)

2.2.2 Tichonov regularization (Ridge regression)

Regularization can reduce the influence of the outliers.

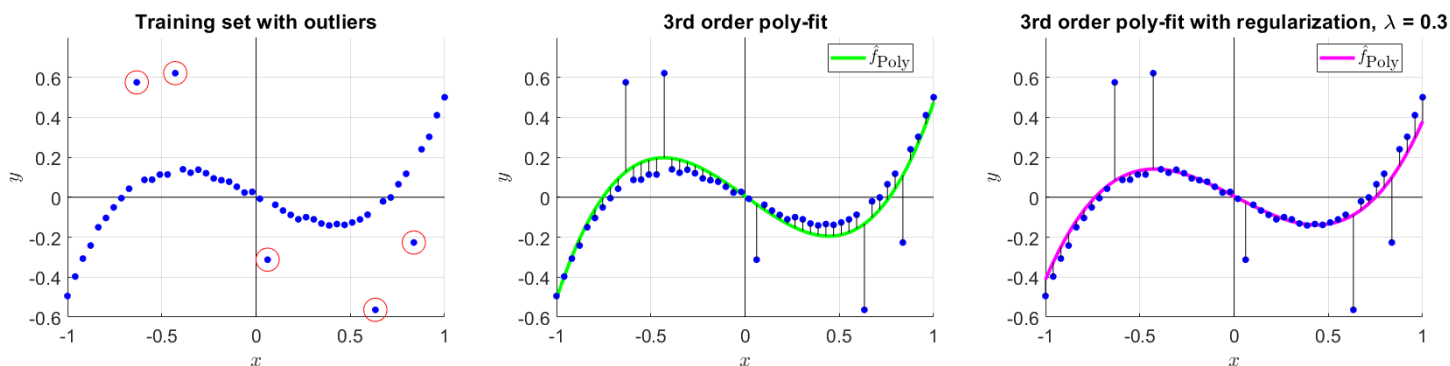
The loss function with a squared norm regularization (Tichonov regularization) is given by:

$$\mathcal{L}(\mathbf{w}) \triangleq L(\mathbf{w}) + \lambda \sum_{m=1}^M w_m^2 = \|\mathbf{y} - \Phi\mathbf{w}\|_2^2 + \lambda \|\mathbf{w}\|_2^2$$

where the parameter λ controls the ratio between the original loss (fidelity) and the regularization term.

To obtain the optimal \mathbf{w} , we compare the gradient to zero:

$$\begin{aligned} \nabla_{\mathbf{w}} \mathcal{L}(\mathbf{w}) &= \mathbf{0} \\ -2\Phi^T(\mathbf{y} - \Phi\mathbf{w}) + 2\lambda\mathbf{w} &= \mathbf{0} \\ \Phi^T\Phi\mathbf{w} + \lambda\mathbf{w} &= \Phi^T\mathbf{y} \\ (\Phi^T\Phi + \lambda\mathbf{I})\mathbf{w} &= \Phi^T\mathbf{y} \\ \Rightarrow \mathbf{w} &= (\Phi^T\Phi + \lambda\mathbf{I})^{-1} \Phi^T\mathbf{y} \end{aligned}$$



2.2.3 General quadratic regularization

We can put different weights for each parameter of the model \mathbf{w} :

$$\mathcal{L}(\mathbf{w}) = L(\mathbf{w}) + \sum_{m=1}^M \lambda_m w_m^2 = \|\mathbf{y} - \Phi \mathbf{w}\|_2^2 + \|\Lambda \mathbf{w}\|_2^2$$

where:

$$\Lambda \triangleq \begin{bmatrix} \lambda_1^{1/2} & 0 & 0 \\ 0 & \ddots & 0 \\ 0 & 0 & \lambda_M^{1/2} \end{bmatrix}$$

In this case the optimal parameters are given by:

$$\Rightarrow \mathbf{w} = \left(\Phi^T \Phi + \Lambda^T \Lambda \right)^{-1} \Phi^T \mathbf{y}$$

Note: Λ does not have to be a diagonal matrix.

2.2.4 ℓ_1 - regularization (LASSO)

Another common choice is using the ℓ_1 regularization:

$$\mathcal{L}(\mathbf{w}) = \|\mathbf{y} - \Phi \mathbf{w}\|_2^2 + \|\mathbf{w}\|_1$$

In this case there is no closed form solution but numerical iterative algorithm can provide the optimal value of \mathbf{w} . This problem is known as LASSO - Least Absolute Shrinkage and Selection Operator.