

Principles of Machine Learning

Lecture 4: Linear Regression

Sharif University of Technology
Dept. of Aerospace Engineering

March 11, 2025



Table of Contents

- 1 Introduction
- 2 Standard Linear Regression
- 3 Maximum Likelihood Estimation
- 4 Least Squares Estimation
- 5 Overfitting and Regularization
- 6 Ridge Regression
- 7 Lasso Regression



Table of Contents

- 1 Introduction
- 2 Standard Linear Regression
- 3 Maximum Likelihood Estimation
- 4 Least Squares Estimation
- 5 Overfitting and Regularization
- 6 Ridge Regression
- 7 Lasso Regression



Linear Regression

- A very widely-used method for predicting a real-valued output (**dependent variable** or **target**) $y \in \mathbb{R}$, given a vector of real-valued input (**independent variables** or **explanatory variables**, or **covariates**) $x \in \mathbb{R}^D$
- Key property: We assume that $\mathbb{E}[y|\mathbf{x}] = \boldsymbol{\theta}^\top \mathbf{x}$



Linear Regression

- **Choice of the model (type) and the parametrization**

Function classes (e.g., polynomials) and particular parametrization (e.g., degree of the polynomial)

- **Finding good parameters**

Loss/objective function and an optimization algorithm to minimize it

- **Overfitting and model selection**

Typically occurs if the underlying model (or its parametrization) is overly flexible and expressive

- **Relationship between loss functions and parameter priors**

Connection between loss functions and the underlying prior assumptions that induce these losses

- **Uncertainty modeling**

Describe the remaining parameter uncertainty to obtain a measure of confidence of the model's prediction at test time



Linear Regression Applications

- 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)
- Deep learning applications (e.g., computer games, speech-to-text translation, image recognition, automatic video annotation)



Introduction

Linear Regression

We aim to find a function f that maps inputs $\mathbf{x} \in \mathbb{R}^D$ to corresponding function values $f(\mathbf{x}) \in \mathbb{R}$.

We are given training inputs \mathbf{x}_n and corresponding noisy observations $y_n = f(\mathbf{x}_n) + \epsilon$, where ϵ is an i.i.d random variable that describes measurement/observation noise.

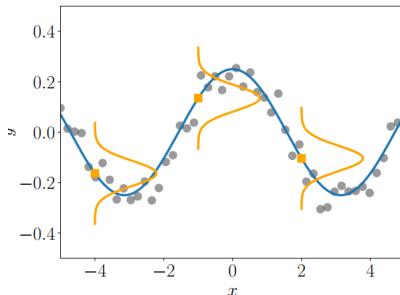


Table of Contents

- 1 Introduction
- 2 Standard Linear Regression**
- 3 Maximum Likelihood Estimation
- 4 Least Squares Estimation
- 5 Overfitting and Regularization
- 6 Ridge Regression
- 7 Lasso Regression



Standard Linear Regression – Problem Formulation

Presence of observation noise \rightarrow Probabilistic approach \rightarrow Explicitly model the noise using a likelihood function:

$$p(y|\mathbf{x}) = \mathcal{N}(y|f(\mathbf{x}), \sigma^2).$$

with the input-output relationship given as:

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

where $\epsilon \sim \mathcal{N}(0, \sigma^2)$ is i.i.d. Gaussian measurement noise with mean 0 and variance σ^2 .

Our objective is to find a function that is similar to the unknown function f that generated the data and that generalizes well.



Standard Linear Regression – Problem Formulation

For now, assume that noise variance σ^2 is known and focus on learning the model parameters θ , where in linear regression, θ appear linearly in our model:

$$p(y|\mathbf{x}, \theta) = \mathcal{N}(y|\mathbf{x}^\top \theta, \sigma^2) \quad (1)$$

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

where $\theta \in \mathbb{R}^D$ are the parameters we seek.

- The *likelihood* in (1) is the probability density function of y evaluated at $\mathbf{x}^\top \theta$.
- Without the source of uncertainty (noise), (1) would be a Dirac delta.



Standard Linear Regression – Problem Formulation

- The linear regression model (2) is linear in the parameters θ and the inputs \mathbf{x} .
- Note that $y = \phi^\top(\mathbf{x})\theta$ is also a linear regression model for nonlinear transformations ϕ . (Why?)

In the following, we discuss how to find good parameters θ and how to evaluate whether or not they work well.

For now, we assume that σ^2 is known.



Table of Contents

- 1 Introduction
- 2 Standard Linear Regression
- 3 Maximum Likelihood Estimation**
- 4 Least Squares Estimation
- 5 Overfitting and Regularization
- 6 Ridge Regression
- 7 Lasso Regression



Given a *training set* $\mathcal{D} := \{(\mathbf{x}_1, y_1), \dots, (\mathbf{x}_N, y_N)\}$ and noting that y_i and y_j are conditionally independent, the likelihood is:

$$p(\mathcal{Y} | \mathcal{X}, \boldsymbol{\theta}) = p(y_1, \dots, y_N | \mathbf{x}_1, \dots, \mathbf{x}_N, \boldsymbol{\theta}) \quad (3)$$

$$= \prod_{n=1}^N p(y_n | \mathbf{x}_n, \boldsymbol{\theta}) = \prod_{n=1}^N \mathcal{N}(y_n | \mathbf{x}_n^\top \boldsymbol{\theta}, \sigma^2), \quad (4)$$

where $\mathcal{X} := \{\mathbf{x}_1, \dots, \mathbf{x}_N\}$ and $\mathcal{Y} := \{y_1, \dots, y_N\}$.

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



Parameter Estimation

Parameter estimation problem: find optimal parameters $\theta^* \in \mathbb{R}^D$ for the linear regression model (3).

Once the θ^* are found \rightarrow predict function values by using this parameter estimate in (3):

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

where \mathbf{x}_* is an arbitrary test input and y_* is the distribution of the corresponding target at that input.

How to solve this parameter estimation problem / find θ^* ?



Maximum Likelihood Estimation (MLE)

Maximum Likelihood Estimation: finding θ_{ML} that maximize the likelihood (4).

- Maximizing the likelihood \rightarrow maximizing the predictive distribution of the training data given the model params:

$$\theta_{ML} = \arg \max_{\theta} p(\mathcal{Y} | \mathcal{X}, \theta). \quad (5)$$

- The likelihood $p(\mathbf{y} | \mathbf{x}, \theta)$ is not a probability distribution in θ .
- It is simply a function of the parameters θ but does not integrate to 1 and may not even be integrable with respect to θ .
- However, the likelihood in (5) is a normalized probability distribution in \mathbf{y} .



Maximum Likelihood Estimation (MLE)

To find θ_{ML} :

- ① Gradient descent (or GD on the negative likelihood).
- ② Apply log-transformation to the likelihood function & Minimize the negative log-likelihood
 - Logarithm is a (strictly) monotonically increasing function \rightarrow the optimum of f is identical to optimum of $\log f$

MLE in linear regression has closed-form solution \rightarrow Option 2 is preferable.

Log-transformations are useful because:

- does not suffer from numerical underflow
- more simple differentiation rules



Maximum Likelihood Estimation (MLE)

Thus, to find θ_{ML} , we minimize the negative log-likelihood:

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

In the linear regression model (2), the likelihood is Gaussian, such that:

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

where the constant includes all terms independent of θ .



Maximum Likelihood Estimation (MLE)

The negative log-likelihood function (also called *error function* or *loss function*) is then as follows:

$$\mathcal{L}(\theta) := \frac{1}{2\sigma^2} \sum_{n=1}^N (y_n - \mathbf{x}_n^\top \theta)^2 \times w_n \quad W = \begin{bmatrix} w_1 & & & \\ & \ddots & & 0 \\ & & \ddots & \\ 0 & & & w_N \end{bmatrix} \quad (6)$$

$$= \frac{1}{2\sigma^2} (\mathbf{y} - \mathbf{X}\theta)^\top W (\mathbf{y} - \mathbf{X}\theta) = \frac{1}{2\sigma^2} \|\mathbf{y} - \mathbf{X}\theta\|_W^2, \quad (7)$$

where we define the *design matrix* $\mathbf{X} := [\mathbf{x}_1, \dots, \mathbf{x}_N]^\top \in \mathbb{R}^{N \times D}$ and $\mathbf{y} := [y_1, \dots, y_N]^\top \in \mathbb{R}^N$

- Note that the sum of squared errors between observations y_n and the corresponding model prediction $\mathbf{x}_n^\top \theta$ equals the squared distance between \mathbf{y} and $\mathbf{X}\theta$,



Maximum Likelihood Estimation (MLE)

Using the results from the previous Lecture (Calculus Basics), we compute the gradient of \mathcal{L} w.r.t. the params:

$$\frac{d\mathcal{L}}{d\boldsymbol{\theta}} = \frac{d}{d\boldsymbol{\theta}} \left(\frac{1}{2\sigma^2} (\mathbf{y} - \mathbf{X}\boldsymbol{\theta})^\top (\mathbf{y} - \mathbf{X}\boldsymbol{\theta}) \right) \quad (8)$$

$$= \frac{1}{2\sigma^2} \frac{d}{d\boldsymbol{\theta}} \left(\mathbf{y}^\top \mathbf{y} - 2\mathbf{y}^\top \mathbf{X}\boldsymbol{\theta} + \boldsymbol{\theta}^\top \mathbf{X}^\top \mathbf{X} \boldsymbol{\theta} \right) \quad (9)$$

$$= \frac{1}{\sigma^2} \left(-\mathbf{y}^\top \mathbf{X} + \boldsymbol{\theta}^\top \mathbf{X}^\top \mathbf{X} \right) \in \mathbb{R}^{1 \times D}. \quad (10)$$



Maximum Likelihood Estimation (MLE)

The MLE estimator θ_{ML} solves $\frac{d\mathcal{L}}{d\theta} = \mathbf{0}$:

$$\frac{d\mathcal{L}}{d\theta} = \mathbf{0} \stackrel{(10)}{\iff} \theta_{ML}^T \mathbf{X}^T \mathbf{X} = \mathbf{y}^T \mathbf{X} \quad (11)$$

$$\iff \theta_{ML}^T = \mathbf{y}^T \mathbf{X} (\mathbf{X}^T \mathbf{X})^{-1} \quad (12)$$

$$\iff \boxed{\theta_{ML} = (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \mathbf{y}}. \quad (13)$$

$$\mathbf{y} = \mathbf{X} \theta$$

- Ignoring the possibility of duplicate data $\text{rk}(\mathbf{X}) = D$ if $N \geq D$, i.e., we do not have more parameters than data points.
- We could right-multiply the first equation by $(\mathbf{X}^T \mathbf{X})^{-1}$ because $\mathbf{X}^T \mathbf{X}$ is positive definite if $\text{rk}(\mathbf{X}) = D$.



Maximum Likelihood Estimation (MLE)

$$\frac{d\mathcal{L}}{d\boldsymbol{\theta}} = \mathbf{0} \implies \boxed{\boldsymbol{\theta}_{\text{ML}} = (\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{X}^\top \mathbf{y}}$$

- Setting the gradient to $\mathbf{0}$ is a necessary and sufficient condition to 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.
- The maximum likelihood solution above requires solving a system of linear equations of the form $\mathbf{A}\boldsymbol{\theta} = \mathbf{b}$ with $\mathbf{A} = (\mathbf{X}^\top \mathbf{X})$ and $\mathbf{b} = \mathbf{X}^\top \mathbf{y}$.



Maximum Likelihood Estimation (MLE) with Features

We considered the linear regression model as (4) that allowed us to fit straight lines to data using MLE.

However, the straight lines are not sufficiently expressive!

Fitting nonlinear functions within the linear regression framework by performing an arbitrary nonlinear transformation $\phi(\mathbf{x})$ of the inputs \mathbf{x} :

$$p(y|\mathbf{x}, \boldsymbol{\theta}) = \mathcal{N}(y|\phi^\top(\mathbf{x})\boldsymbol{\theta}, \sigma^2) \quad (14)$$

$$\iff y = \phi^\top(\mathbf{x})\boldsymbol{\theta} + \epsilon = \sum_{k=0}^{K-1} \theta_k \phi_k(\mathbf{x}) + \epsilon, \quad (15)$$

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



Maximum Likelihood Estimation (MLE) with Features

Example 1. Polynomial Regression

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

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

Then, a polynomial of degree $K - 1$ is

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

where ϕ is defined as above and $\theta = [\theta_0, \dots, \theta_{K-1}]^\top \in \mathbb{R}^K$ contains the



Maximum Likelihood Estimation (MLE) with Features

Revisiting the design matrix and defining it as the *feature matrix* in this context as:

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

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

Example 2. 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

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



Maximum Likelihood Estimation (MLE) with Features

With the feature matrix Φ defined, the negative log-likelihood for the linear regression model (14,15) can be written as:

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

Since both \mathbf{X} and Φ are independent of the params θ that we wish to optimize, the *maximum likelihood estimate* is

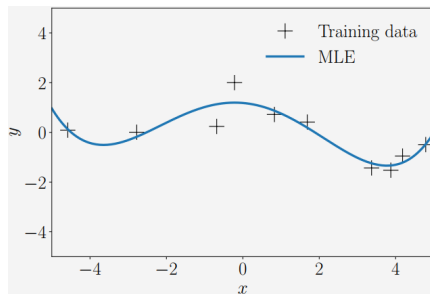
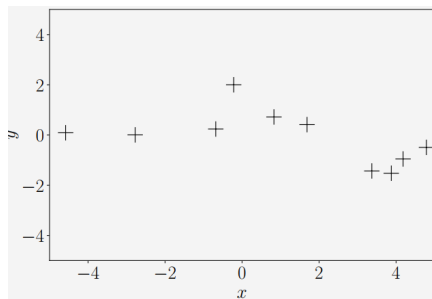
$$\boxed{\theta_{\text{ML}} = (\Phi^\top \Phi)^{-1} \Phi \mathbf{y}} \quad (16)$$

for the linear regression problem with nonlinear features defined in (14,15).



Maximum Likelihood Estimation (MLE) with Features

Example 3. Maximum Likelihood Polynomial Fit



The figure (a) shows the dataset consisting of $N = 10$ 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 4 using MLE, i.e., parameters θ_{ML} are given in (16). The result is shown in figure (b).



Maximum Likelihood Estimation (MLE) with Features

To estimate the noise variance σ^2 :

$$\begin{aligned}\log p(\mathcal{Y} | \mathcal{X}, \boldsymbol{\theta}, \sigma^2) &= \sum_{n=1}^N \log \mathcal{N}(y_n | \boldsymbol{\phi}^\top(\mathbf{x}_n)\boldsymbol{\theta}, \sigma^2) \\ &= \sum_{n=1}^N \left(-\frac{1}{2} \log(2\pi) - \frac{1}{2} \log \sigma^2 - \frac{1}{2\sigma^2} (y_n - \boldsymbol{\phi}^\top(\mathbf{x}_n)\boldsymbol{\theta})^2 \right) \\ &= -\frac{N}{2} \log \sigma^2 - \underbrace{\frac{1}{2\sigma^2} \sum_{n=1}^N (y_n - \boldsymbol{\phi}^\top(\mathbf{x}_n)\boldsymbol{\theta})^2}_{=: s} + \text{const.}\end{aligned}$$

Handwritten red note: $\frac{1}{\sqrt{2\pi}\sigma} \exp\left[-\frac{(y - \boldsymbol{\phi}^\top\boldsymbol{\theta})^2}{2\sigma^2}\right]$



Maximum Likelihood Estimation (MLE) with Features

$$\frac{\partial \log p(\mathcal{Y} | \mathcal{X}, \boldsymbol{\theta}, \sigma^2)}{\partial \sigma^2} = -\frac{N}{2\sigma^2} + \frac{1}{2\sigma^4} s = 0$$
$$\iff \frac{N}{2\sigma^2} = \frac{s}{2\sigma^4}$$

then

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



Table of Contents

- 1 Introduction
- 2 Standard Linear Regression
- 3 Maximum Likelihood Estimation
- 4 Least Squares Estimation**
- 5 Overfitting and Regularization
- 6 Ridge Regression
- 7 Lasso Regression



LSE – MLE as Orthogonal Projection

Recalling linear regression without features:

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

considering a straight line $f: \mathbb{R} \rightarrow \mathbb{R}$ that go through origin, the parameter θ determines the slope of the line.

Recalling the maximum likelihood estimator:

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

$$\implies \mathbf{X}\theta_{ML} = \mathbf{X} \frac{\mathbf{X}^\top \mathbf{y}}{\mathbf{X}^\top \mathbf{X}} = \frac{\mathbf{X}\mathbf{X}^\top}{\mathbf{X}^\top \mathbf{X}} \mathbf{y}$$

We obtain the approximation with the minimum least-squares error between \mathbf{y} and $\mathbf{X}\theta$.

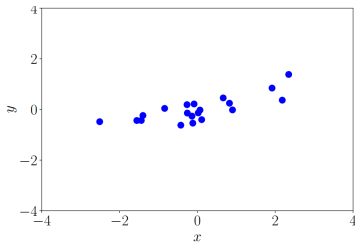


LSE – MLE as Orthogonal Projection

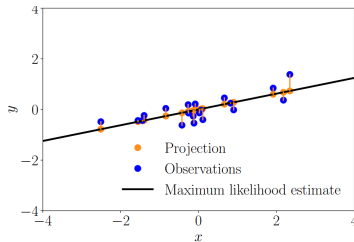
- Linear regression can be thought of as a method for solving systems of linear equations $\mathbf{y} = \mathbf{X}\theta$.
- The ML estimator θ_{ML} effectively does an orthogonal projection of \mathbf{y} onto the 1-D subspace spanned by \mathbf{X} .
- $\frac{\mathbf{X}\mathbf{X}^T}{\mathbf{X}^T\mathbf{X}}$ is the projection matrix.
- θ_{ML} is the coordinates of the projection onto the 1-D 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.



Geometric interpretation of least squares



(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.



LSE – MLE as Orthogonal Projection

In the general linear regression case

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

we again can interpret the maximum likelihood result

$$\mathbf{y} \approx \Phi\boldsymbol{\theta}_{\text{ML}},$$

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

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



LSE – MLE as Orthogonal Projection

If the feature functions ϕ_k that we use to construct the feature matrix Φ are orthonormal, leads to the projection

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

Thus, the maximum likelihood projection is simply the sum of the projections of \mathbf{y} onto the individual basis vectors ϕ_k .

If the basis is not orthogonal \longrightarrow the Gram-Schmidt process



Viewing MLE as an orthogonal projection, the NLL (negative log-likelihood) is equal (up to irrelevant constants) to the **residual sum of squares**:

$$\text{RSS}(\boldsymbol{\theta}) = \frac{1}{2} \sum_{n=1}^N (y_n - \boldsymbol{\theta}^\top \mathbf{x}_n)^2 = \frac{1}{2} \|\mathbf{X}\boldsymbol{\theta} - \mathbf{y}\|^2 = \frac{1}{2} (\mathbf{X}\boldsymbol{\theta} - \mathbf{y})^\top (\mathbf{X}\boldsymbol{\theta} - \mathbf{y}) \quad (17)$$

In the following, methods to optimize the (17) are discussed.



LSE – Ordinary Least Squares (OLS)

$$\nabla_{\theta} \text{RSS}(\theta) = \mathbf{X}^{\top} \mathbf{X} \theta - \mathbf{X}^{\top} \mathbf{y}$$

$$\xrightarrow{\text{setting the gradient to zero}} \mathbf{X}^{\top} \mathbf{X} \theta = \mathbf{X}^{\top} \mathbf{y}$$

The above equations are known as **normal equations** (why?) and the corresponding solution θ^* is called the **OLS** solution:

$$\theta^* = (\mathbf{X}^{\top} \mathbf{X})^{-1} \mathbf{X}^{\top} \mathbf{y}$$

- Show that the least squares objective has a unique global minimum (show that the Hessian is positive definite).



For the OLS solution as

$$\theta^* = \mathbf{X}^\dagger \mathbf{y} = (\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{X}^\top \mathbf{y},$$

it may not be numerically possible to compute the pseudo-inverse by inverting $\mathbf{X}^\top \mathbf{X}$ since it may be singular.

A better and more general approach \longrightarrow SVD

In `sklearn.linear_model.fit`, it uses `scipy.linalg.lstsq` which calls DGELSD, which is an SVD-based solver.

But ...



If \mathbf{X} is tall and skinny ($N \gg D$) \longrightarrow QR decomposition

Let $\mathbf{X} = \mathbf{QR}$, where $\mathbf{Q}^\top \mathbf{Q} = \mathbf{I}$.

Using QR decomposition, we rewrite solving the linear system $\mathbf{X}\boldsymbol{\theta} = \mathbf{y}$ that minimizes $\|\mathbf{X}\boldsymbol{\theta} - \mathbf{y}\|^2$ (If $N = D$ and \mathbf{X} is full rank, the equations have unique solution and the error will be 0) as

$$\begin{aligned}(\mathbf{QR})\boldsymbol{\theta} &= \mathbf{y} \\ \mathbf{Q}^\top \mathbf{QR}\boldsymbol{\theta} &= \mathbf{Q}^\top \mathbf{y} \\ \boldsymbol{\theta} &= \mathbf{R}^{-1}(\mathbf{Q}^\top \mathbf{y})\end{aligned}$$

Thanks to this approach: backsubstitution instead of matrix inversion



SVD and QR are direct methods based on matrix decomposition.

Alternative approach: iterative solvers such as

- **conjugate gradient** method (assumes \mathbf{X} is symmetric positive definite)
- **GMRES** (generalized minimal residual method) that works for general \mathbf{X} (implemented by `sparse.linalg.gmres` in SciPy)

These are well-suited to problems where \mathbf{X} is sparse or structured.



Weighted LSE

When we want to associate a weight with each example.

An instance: **heteroskedastic regression** (the variance depends on the input):

$$p(y | \mathbf{x}, \boldsymbol{\theta}) = \mathcal{N}(y | \boldsymbol{\theta}^\top \mathbf{X}, \sigma^2(\mathbf{x})) = \frac{1}{\sqrt{2\pi\sigma^2(\mathbf{x})}} \exp\left(-\frac{1}{2\sigma^2(\mathbf{x})}(y - \boldsymbol{\theta}^\top \mathbf{x})^2\right)$$

$$\implies p(\mathbf{y} | \mathbf{x}, \boldsymbol{\theta}) = \mathcal{N}(\mathbf{y} | \mathbf{X}\boldsymbol{\theta}, \boldsymbol{\Lambda}^{-1})$$

where $\boldsymbol{\Lambda} = \text{diag}(1/\sigma^2(\mathbf{x}_n))$.

One can show that the MLE is given by

$$\boldsymbol{\theta}^* = (\mathbf{X}^\top \boldsymbol{\Lambda} \mathbf{X})^{-1} \mathbf{X}^\top \boldsymbol{\Lambda} \mathbf{y}$$

which is known as the **weighted least squares** estimate.



Table of Contents

- 1 Introduction
- 2 Standard Linear Regression
- 3 Maximum Likelihood Estimation
- 4 Least Squares Estimation
- 5 Overfitting and Regularization**
- 6 Ridge Regression
- 7 Lasso Regression



Goodness of Fit is about how well a regression model fits the data.

- **Residual plots**
- **Prediction accuracy and R^2**
 - RSS
 - RMSE
 - R^2



Measuring Goodness of Fit – Residual Plots

For 1d inputs,

Check the *reasonableness* of the model by plotting the *residuals*,

$$r_n = y_n - y_n^*$$

vs the input x_n .

The model assumes that

- The residuals have a $\mathcal{N}(0, \sigma^2)$ distribution

Thus, the residual plots are of points more or less equally above and below the horizontal line at 0 with no obvious trends.



Measuring Goodness of Fit – Prediction accuracy and R^2

Rather than plotting against x_n , measure the quality of a regression's model predictions by plotting y_n^* vs y_n .

To assess the fit quantitatively using the **residual sum of squares**:

$$\text{RSS}(\theta) = \sum_{n=1}^N (y_n - \theta^\top \mathbf{x}_n)^2$$

Another measure, **root mean squared error**:

$$\text{RMSE}(\theta) \triangleq \sqrt{\frac{1}{N} \text{RSS}(\theta)}$$

A more interpretable measure, **coefficient of determination**:

$$R^2 \triangleq 1 - \frac{\sum_{n=1}^N (y_n^* - y_n)^2}{\sum_{n=1}^N (\bar{y} - y_n)^2} = 1 - \frac{\text{RSS}}{\text{TSS}}$$



Coefficient of Determination R^2

$$R^2 \triangleq 1 - \frac{\sum_{n=1}^N (y_n^* - y_n)^2}{\sum_{n=1}^N (\bar{y} - y_n)^2} = 1 - \frac{\text{RSS}}{\text{TSS}}$$

where

- $\bar{y} = \frac{1}{N} \sum_{n=1}^N y_n$ is the empirical mean of the response
- $\text{RSS} = \frac{1}{N} \sum_{n=1}^N (y_n - y_n^*)^2$ is the residual sum of squares
- $\text{TSS} = \frac{1}{N} \sum_{n=1}^N (y_n - \bar{y})^2$ is the total sum of squares

Note that

- R^2 measures *the variance in the prediction* relative to a simple constant prediction of $y_n^* = \bar{y}$
- $0 \leq R^2 \leq 1$ (why?)
 - larger values \longrightarrow a greater reduction in variance \longrightarrow better fit



Overfitting

For **model selection**,

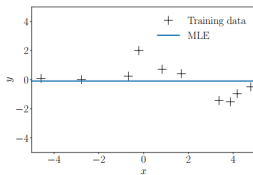
- Use RMSE (or NLL) \rightarrow find the best degree of the polynomial M that minimizes the objective
- Brute-force search \rightarrow enumerate all (reasonable) values of M
- For a training set of size N it is sufficient to test $0 \leq M \leq N - 1$
 - For $M < N$, the maximum likelihood estimator is unique.
 - For $M \geq N$, infinitely many possible maximum likelihood estimators.
 - The case of $M = N - 1$ is extreme.

Note that

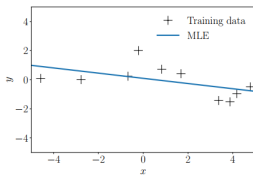
- $M = N - 1$ is extreme in the sense that otherwise the null space of the corresponding linear system would be non-trivial \rightarrow infinitely many optimal solutions to the linear regression problem.
- For $M \geq N$,
more params than data points \rightarrow undetermined system \rightarrow the $\Phi^T \Phi$ in slide 34 is no longer invertible.



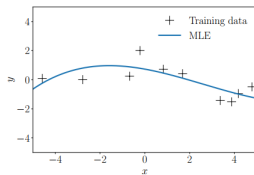
Overfitting



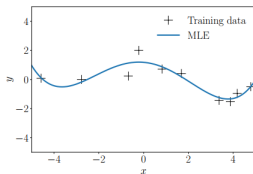
(a) $M = 0$



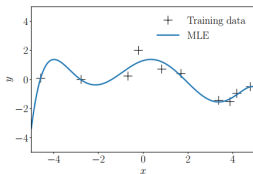
(b) $M = 1$



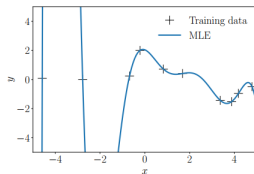
(c) $M = 3$



(d) $M = 4$



(e) $M = 6$



(f) $M = 9$

Figure: Maximum likelihood fits for different polynomial degrees M with $N = 10$ observations



Overfitting

In the figure in previous slide,

- The polynomials of low degree (e.g., constants ($M = 0$) or linear ($M = 1$)) fit the data poorly \rightarrow poor representations of the true underlying function
- For degrees $M = 3, \dots, 6$, the fit looks plausible and smoothly interpolate the data
- For higher degrees, they fit the data better and better
- For the extreme case $M = N - 1 = 9$, the function will pass through every single data point

BUT, the high-degree polynomials (a) oscillate wildly and (b) are a poor representation of the underlying function that generated the data.

That's *overfitting*.



Let's test!

To gain quantitative insight into the dependence of the generalization performance on the polynomials of degree M

- A separate test set of 200 data points generated using exactly the same procedure used to generate the training set
- a linear grid of 200 points in the interval of $[-5, 5]$ as test inputs
- Evaluate the RMSE for both training and test data



Overfitting

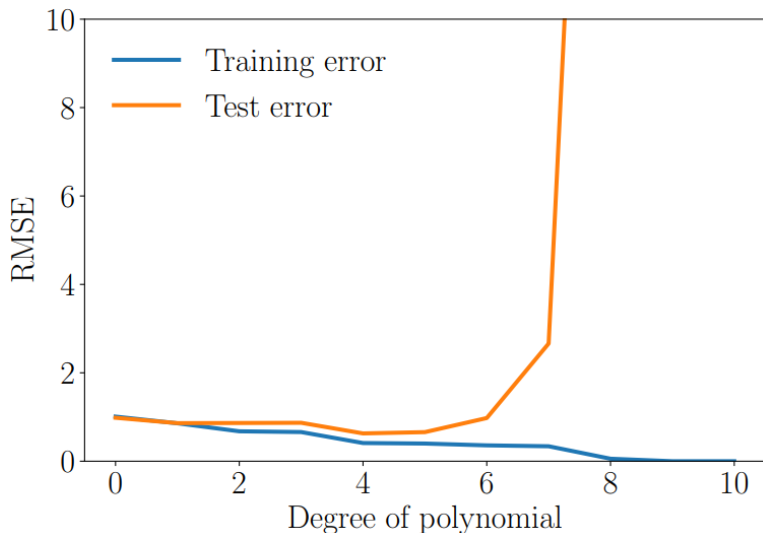


Figure: Training and test errors



Investigating the “training and test errors” figure,

- Test error
 - is a qualitative measure of the generalization properties of the corresponding polynomial
 - initially decreases
 - is relatively low for fourth-order polynomials and stays relatively constant up to degree 5
 - increases significantly from degree 6 onward → high-order polynomials have very bad generalization
- Training error never increases with the increase in the degree of the polynomial
- The best generalization → the point of the smallest test error → $M = 4$

