

ECE595 / STAT598: Machine Learning I

Lecture 1.1: Linear Regression

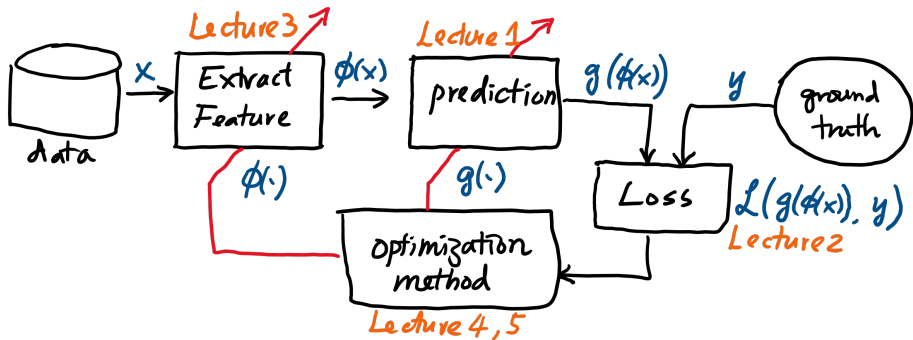
Spring 2020

Stanley Chan

School of Electrical and Computer Engineering
Purdue University



Outline



Outline

Mathematical Background

- Lecture 1: Linear regression: A basic data analytic tool
- Lecture 2: Regularization: Constraining the solution
- Lecture 3: Kernel Method: Enabling nonlinearity

Lecture 1: Linear Regression

- Linear Regression
 - Notation
 - Loss Function
 - Solving the Regression Problem
- Geometry
 - Projection
 - Minimum-Norm Solution
 - Pseudo-Inverse

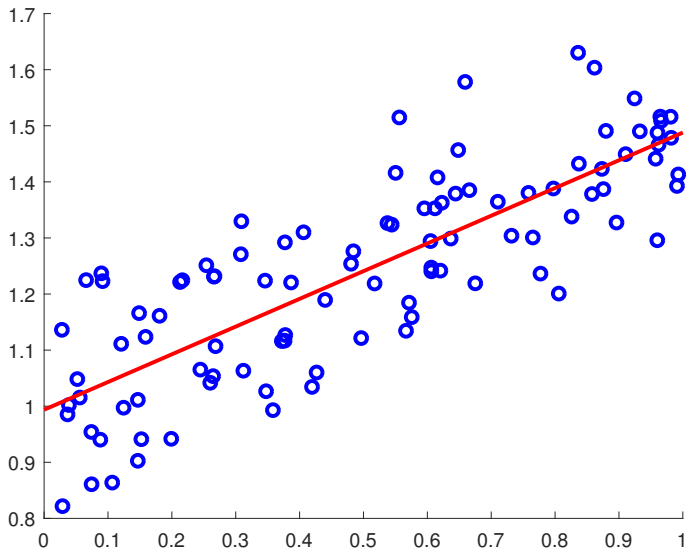
Basic Notation

- Scalar: $a, b, c \in \mathbb{R}$
- Vector: $\mathbf{a}, \mathbf{b}, \mathbf{c} \in \mathbb{R}^d$
- Matrix: $\mathbf{A}, \mathbf{B}, \mathbf{C} \in \mathbb{R}^{N \times d}$; Entries are a_{ij} or $[\mathbf{A}]_{ij}$.
- Rows and Columns

$$\mathbf{A} = \begin{bmatrix} | & | & \dots & | \\ \mathbf{a}_1 & \mathbf{a}_2 & \dots & \mathbf{a}_d \\ | & | & \dots & | \end{bmatrix}, \quad \text{and} \quad \mathbf{A} = \begin{bmatrix} - & (\mathbf{x}^1)^T & - \\ - & (\mathbf{x}^2)^T & - \\ & \vdots & \\ - & (\mathbf{x}^N)^T & - \end{bmatrix}.$$

- $\{\mathbf{a}_j\}$: The j -th feature. $\{\mathbf{x}^n\}$: The n -th sample.
- Identity matrix \mathbf{I}
- All-one vector $\mathbf{1}$ and all-zero vector $\mathbf{0}$
- Standard basis \mathbf{e}_j .

Line Fitting



Linear Regression

The problem of **regression** can be summarized as:

- Given measurements: y^n , (where $n = 1, \dots, N$)
- Given inputs: \mathbf{x}^n
- Given a model: $g_{\theta}(\mathbf{x}^n)$ parameterized by θ
- Determine θ such that $y^n \approx g_{\theta}(\mathbf{x}^n)$.

Linear regression is one type of regression:

- Restrict $g_{\theta}(\cdot)$ to a line:

$$g_{\theta}(\mathbf{x}) = \mathbf{x}^T \theta$$

- The inputs \mathbf{x} and the parameters θ are

$$\mathbf{x} = [x_1, \dots, x_d]^T \quad \text{and} \quad \theta = [\theta_1, \dots, \theta_d]^T$$

- This is equivalent to

$$g_{\theta}(\mathbf{x}) = \mathbf{x}^T \theta = \sum_{j=1}^d x_j \theta_j.$$

Solving the Regression Problem

The (general) regression can be solved via the following logic:

- Define a (Squared-Error) *Loss Function*

$$\begin{aligned} J(\theta) &= \sum_{n=1}^N \mathcal{L}(g_{\theta}(\mathbf{x}^n), y^n) \\ &= \sum_{n=1}^N (g_{\theta}(\mathbf{x}^n) - y^n)^2, \quad \text{e.g., } \mathcal{L}(\clubsuit, \spadesuit) \stackrel{\text{def}}{=} (\clubsuit - \spadesuit)^2 \end{aligned}$$

- Other loss functions can be used, e.g., $\mathcal{L}(\clubsuit, \spadesuit) = |\clubsuit - \spadesuit|$.
- The goal is to solve an optimization

$$\hat{\theta} = \underset{\theta}{\operatorname{argmin}} J(\theta).$$

- The prediction of a new input \mathbf{x}^{new} is $y^{\text{new}} = g_{\hat{\theta}}(\mathbf{x}^{\text{new}}) = \hat{\theta}^T \mathbf{x}^{\text{new}}$.

Linear Regression Solution

The linear regression problem is a special case which we can solve analytically.

- Restrict $g_{\theta}(\cdot)$ to a line:

$$g_{\theta}(\mathbf{x}^n) = \boldsymbol{\theta}^T \mathbf{x}^n$$

- Then the loss function becomes

$$J(\boldsymbol{\theta}) = \sum_{n=1}^N (\boldsymbol{\theta}^T \mathbf{x}^n - y^n)^2 = \|\mathbf{A}\boldsymbol{\theta} - \mathbf{y}\|^2.$$

- The matrix and vectors are defined as

$$\mathbf{A} = \begin{bmatrix} \text{---} & (\mathbf{x}^1)^T & \text{---} \\ \text{---} & (\mathbf{x}^2)^T & \text{---} \\ & \vdots & \\ \text{---} & (\mathbf{x}^N)^T & \text{---} \end{bmatrix}, \quad \boldsymbol{\theta} = \begin{bmatrix} \theta_1 \\ \vdots \\ \theta_d \end{bmatrix}, \quad \text{and} \quad \mathbf{y} = \begin{bmatrix} y_1 \\ y_2 \\ \vdots \\ y_N \end{bmatrix}$$

- $\|\cdot\|^2$ stands for the ℓ_2 -norm square. See Tutorial on Linear Algebra.

Linear Regression Solution

Theorem

For a linear regression problem

$$\hat{\boldsymbol{\theta}} = \underset{\boldsymbol{\theta}}{\operatorname{argmin}} J(\boldsymbol{\theta}) \stackrel{\text{def}}{=} \|\mathbf{A}\boldsymbol{\theta} - \mathbf{y}\|^2,$$

the minimizer is

$$\hat{\boldsymbol{\theta}} = (\mathbf{A}^T \mathbf{A})^{-1} \mathbf{A}^T \mathbf{y}.$$

- Take derivative and setting to zero: (See Tutorial on “Linear Algebra”.)

$$\begin{aligned}\nabla_{\boldsymbol{\theta}} J(\boldsymbol{\theta}) &= \nabla_{\boldsymbol{\theta}} \{\|\mathbf{A}\boldsymbol{\theta} - \mathbf{y}\|^2\} \\ &= 2\mathbf{A}^T (\mathbf{A}\boldsymbol{\theta} - \mathbf{y}) = \mathbf{0}.\end{aligned}$$

- So solution is $\hat{\boldsymbol{\theta}} = (\mathbf{A}^T \mathbf{A})^{-1} \mathbf{A}^T \mathbf{y}$, assuming $\mathbf{A}^T \mathbf{A}$ is invertible.

Examples

Example 1: Second-order polynomial fitting

$$y_n = ax_n^2 + bx_n + c$$

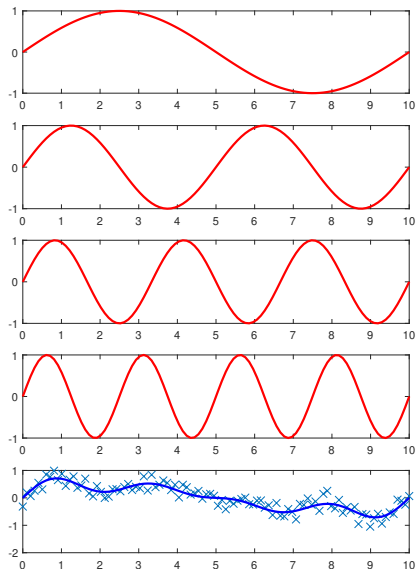
$$\mathbf{A} = \begin{bmatrix} x_1^2 & x_1 & 1 \\ \vdots & \vdots & \vdots \\ x_N^2 & x_N & 1 \end{bmatrix}, \quad \mathbf{y} = \begin{bmatrix} y_1 \\ \vdots \\ y_N \end{bmatrix}, \quad \boldsymbol{\theta} = \begin{bmatrix} a \\ b \\ c \end{bmatrix}$$

Example 2: Auto-regression

$$y_n = ay_{n-1} + by_{n-2}$$

$$\mathbf{A} = \begin{bmatrix} y_2 & y_1 \\ y_3 & y_2 \\ \vdots & \vdots \\ y_{N-1} & y_{N-2} \end{bmatrix}, \quad \mathbf{y} = \begin{bmatrix} y_3 \\ y_4 \\ \vdots \\ y_N \end{bmatrix}, \quad \boldsymbol{\theta} = \begin{bmatrix} a \\ b \end{bmatrix}$$

Generalized Linear Regression



- Eg 1: Fourier series

$$\mathbf{x}^n = \begin{bmatrix} x_1^n \\ x_2^n \\ \vdots \\ x_d^n \end{bmatrix} = \begin{bmatrix} \sin(\omega_0 t_n) \\ \sin(2\omega_0 t_n) \\ \vdots \\ \sin(K\omega_0 t_n) \end{bmatrix}$$

$$y^n = \boldsymbol{\theta}^T \mathbf{x}^n = \sum_{k=1}^d \theta_k \sin(k\omega_0 t_n)$$

- θ_k : k -th Fourier coefficient
- $\sin(k\omega_0 t_n)$: k -th Fourier basis at time t_n