

Lecture 5 - *Linear Supervised Learning*

OBJECTIVE: Linear regression is a supervised learning task. It is of great interest because:

- Many real processes can be approximated with linear models.
- Linear regression appears as part of larger problems.
- It can be solved analytically.
- It illustrates many of the approaches to machine learning.

Given the data $\{x_{1:n}, y_{1:n}\}$, with $x_i \in \mathbb{R}^d$ and $y_i \in \mathbb{R}$, we want to fit a hyper-plane that maps x to y .



Mathematically, the linear model is expressed as follows:

$$\hat{y}_i = \theta_0 + \sum_{j=1}^d x_{ij} \theta_j$$

We let $x_{i,0} = 1$ to obtain

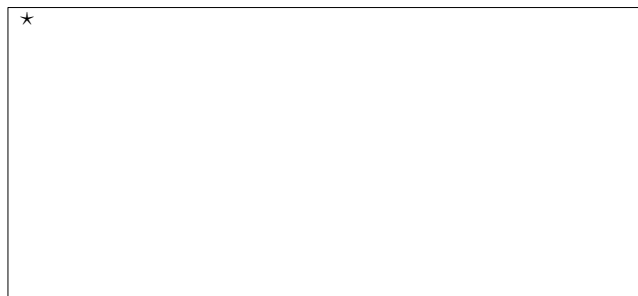
$$\hat{y}_i = \sum_{j=0}^d x_{ij} \theta_j$$

In matrix form, this expression is

$$\hat{Y} = X\theta$$

$$\begin{bmatrix} \hat{y}_1 \\ \vdots \\ \hat{y}_n \end{bmatrix} = \begin{bmatrix} x_{10} & \cdots & x_{1d} \\ \vdots & \vdots & \vdots \\ x_{n0} & \cdots & x_{nd} \end{bmatrix} \begin{bmatrix} \theta_0 \\ \vdots \\ \theta_d \end{bmatrix}$$

If we have several outputs $y_i \in \mathbb{R}^c$, our linear regression expression becomes:



We will present several approaches for computing θ .

◇ OPTIMIZATION APPROACH

Our aim is to minimise the quadratic cost between the output labels and the model predictions

$$C(\theta) = (Y - X\theta)^T(Y - X\theta)$$

★

We will need the following results from matrix differentiation: $\frac{\partial A\theta}{\partial \theta} = A^T$ and $\frac{\partial \theta^T A \theta}{\partial \theta} = 2A^T \theta$

★

$$\frac{\partial C}{\partial \theta} =$$

These are the **normal equations**. The solution (estimate) is:

$$\hat{\theta} =$$

The corresponding predictions are

$$\hat{Y} = HY =$$

where H is the “hat” matrix.

◇ GEOMETRIC APPROACH

★

$$X^T(Y - \hat{Y}) =$$

◇ PROBABILISTIC APPROACH

Univariate Gaussian Distribution

The probability density function of a Gaussian distribution is given by

$$p(x) = \frac{1}{\sqrt{2\pi\sigma^2}}e^{-\frac{1}{2\sigma^2}(x-\mu)^2}.$$

where μ is the mean or center of mass and σ^2 is the variance.

★

Our short notation for Gaussian variables is $X \sim \mathcal{N}(\mu, \sigma^2)$.

Multivariate Gaussian Distribution

Let $x \in \mathbb{R}^n$. The pdf of an n -dimensional Gaussian is given by

$$p(x) = \frac{1}{2\pi^{n/2}|\Sigma|^{1/2}} e^{-\frac{1}{2}(x-\mu)^T \Sigma^{-1}(x-\mu)}$$

where

$$\mu = \begin{pmatrix} \mu_1 \\ \vdots \\ \mu_n \end{pmatrix} = \begin{pmatrix} \mathbb{E}(x_1) \\ \vdots \\ \mathbb{E}(x_n) \end{pmatrix}$$

and

$$\Sigma = \begin{pmatrix} \sigma_{11} & \cdots & \sigma_{1n} \\ \vdots & & \vdots \\ \sigma_{n1} & \cdots & \sigma_{nn} \end{pmatrix} = \mathbb{E}[(X - \mu)(X - \mu)^T]$$

with $\sigma_{ij} = \mathbb{E}[X_i - \mu_i)(X_j - \mu_j)^T]$.

We can interpret each component of x , for example, as a feature of an image such as colour or texture. The term $\frac{1}{2}(x-\mu)^T \Sigma^{-1}(x-\mu)$ is called the **Mahalanobis distance**. Conceptually, it measures the distance between x and μ .

Maximum Likelihood

If our errors are Gaussian distributed, we can use the model

$$Y = X\theta + \mathcal{N}(0, \sigma^2 I)$$

Note that the mean of Y is $X\theta$ and that its variance is $\sigma^2 I$. So we can equivalently write this expression using the probability density of Y **given** X , θ and σ :

$$p(Y|X, \theta, \sigma) = (2\pi\sigma^2)^{-n/2} e^{-\frac{1}{2\sigma^2}(Y-X\theta)^T(Y-X\theta)}$$

The maximum likelihood (ML) estimate of θ is obtained by taking the derivative of the log-likelihood, $\log p(Y|X, \theta, \sigma)$. The idea of maximum likelihood learning is to maximise the likelihood of seeing some data Y by modifying the parameters (θ, σ) .

The ML estimate of θ is:

★

Proceeding in the same way, the ML estimate of σ is:

★

Lecture 6 - *Ridge Regression*

OBJECTIVE: Here we learn a cost function for linear supervised learning that is more stable than the one in the previous lecture. We also introduce the very important notion of **regularization**.

All the answers so far are of the form

$$\hat{\theta} = (X^T X)^{-1} X^T Y$$

They require the inversion of XX^T . This can lead to problems if the system of equations is poorly conditioned. A solution is to add a small element to the diagonal:

$$\hat{\theta} = (X^T X + \delta^2 I_d)^{-1} X^T Y$$

This is the ridge regression estimate. It is the solution to the following **regularised quadratic cost function**

$$C(\theta) = (Y - X\theta)^T (Y - X\theta) + \delta^2 \theta^T \theta$$

★ Proof:

It is useful to visualise the quadratic optimisation function and the constraint region.

★

That is, we are solving the following **constrained optimisation** problem:

$$\min_{\theta: \theta^T \theta \leq t} \{(Y - X\theta)^T(Y - X\theta)\}$$

Large values of θ are penalised. We are **shrinking** θ towards zero. This can be used to carry out **feature weighting**. **An input $x_{i,d}$ weighted by a small θ_d will have less influence on the output y_i .**

Spectral View of LS and Ridge Regression

Again, let $X \in \mathbb{R}^{n \times d}$ be factored as

$$X = U\Sigma V^T = \sum_{i=1}^d u_i \sigma_i v_i^T,$$

where we have assumed that the rank of X is d .

★ The least squares prediction is:

$$\hat{Y}_{LS} = \sum_{i=1}^d u_i u_i^T Y$$

★ Likewise, for ridge regression we have:

$$\hat{Y}_{ridge} = \sum_{i=1}^d \frac{\sigma_i^2}{\sigma_i^2 + \delta^2} u_i u_i^T Y$$

The filter factor

$$f_i = \frac{\sigma_i^2}{\sigma_i^2 + \delta^2}$$

penalises small values of σ^2 (they go to zero at a faster rate).

★

Also, by increasing δ^2 we are penalising the weights:

★

Small eigenvectors tend to be wobbly. The Ridge filter factor f_i gets rid of the wobbly eigenvectors. Therefore, the predictions tend to be more stable (smooth, regularised).

The smoothness parameter δ^2 is often estimated by cross-validation or Bayesian hierarchical methods.

Minimax and cross-validation

Cross-validation is a widely used technique for choosing δ . Here's an example:

