

SDS 387 Linear Models

Fall 2025

Lecture 17 - Tue, Oct 28, 2025

Instructor: Prof. Ale Rinaldo

- Last time: regression modeling. Task of relating a response variable Y to a vector of covariates $X \in \mathbb{R}^d$.

Regression function $x \in \mathbb{R}^d \mapsto \mathbb{E}[Y | X=x]$

↓
↳ projection of Y
on the vector space of
functions of X

↑
target in regression

- We will be focussing on linear regression, i.e. we will approximate the regression function with a linear function

$$x \mapsto x^T \beta \quad \text{some } \beta \in \mathbb{R}^d$$

↓
recall the first coordinate
of the vector x is a 1
to allow for an intercept

- 2 tasks in regression modeling

hypothesis testing
confidence sets

i) statistical inference: carry out inference on
a target parameter. If the model is linear

$$\text{i.e. } \mathbb{E}[Y | X=x] = x^T \beta^* \quad \text{some } \beta^*$$

then our target is β^* .

What if the model is mis-specified (i.e. the regression function is not necessarily linear)

In this case there still exists a natural target parameter - the parameter β^* that gives us the "best" linear approximation to regression function

This is called the projection parameter. It is well defined provided that

$$\mathbb{E}[Y^2] < \infty \quad \text{and} \quad \mathbb{E}[XX^T] \text{ exists (is invertible)}$$

It is equal to

$$\beta^* = \arg\min_{\beta \in \mathbb{R}^d} \mathbb{E}[(\mathbb{E}[Y|X] - x^T \beta)^2]$$

$$= \arg\min_{\beta \in \mathbb{R}^d} \mathbb{E}[(Y - x^T \beta)^2]$$

$$= \underbrace{\left(\mathbb{E}[X X^T] \right)^{-1}}_{d \times d} \underbrace{\mathbb{E}[Y \cdot X]}_{d \times 1}$$

PA/ β^* is the minimizer of

$$\beta \in \mathbb{R}^d \mapsto \mathbb{E}[(X^T \beta)^2] - 2 \mathbb{E}[\mathbb{E}[Y|X] \cdot (X^T \beta)]$$

The gradient at β is

$$\mathbb{E}[2 X X^T \beta] - 2 \mathbb{E}[\mathbb{E}[Y|X] \cdot X]$$

Because this is a strictly convex function it is enough to set gradient = 0 and solve for β . The solution is

$$\begin{aligned} \beta^* &= \left(\mathbb{E}[X X^T] \right)^{-1} \underbrace{\mathbb{E}[\mathbb{E}[Y|X] \cdot X]}_{= \mathbb{E}[Y \cdot X]} \\ &= \mathbb{E}[Y \cdot X] \end{aligned}$$

- Interpretation of β^* : vector of coefficient of the L_2 projection of Y onto the linear span of X (vector space of all linear combinations of X).

↓ it is a parameter expressing linear association btw Y and X

2) prediction: We want to predict a new observation of the response Y_{new} (which we do not observe) based on a new observation of the covariates X_{new} , which we do observe.

Formally, we want to minimize the prediction error:

$$\mathbb{E} \left[(Y_{\text{new}} - X_{\text{new}}^T \beta)^2 \right]$$

over all $\beta \in \mathbb{R}^d$. The minimizer here is of course the projection parameter.

We want to quantify the prediction risk

(let us write Y for Y_{new})
 any $\beta \in \mathbb{R}^d$ X for X_{new}

$$\begin{aligned} \mathbb{E} \left[(Y - X^T \beta)^2 \right] &= \mathbb{E} \left[(Y - X^T \beta^* + X^T \beta^* - X^T \beta)^2 \right] \\ &= \mathbb{E} \left[(Y - X^T \beta^*)^2 \right] + \mathbb{E} \left[(X^T (\beta^* - \beta))^2 \right] \\ &\quad + 2 \mathbb{E} \left[(Y - X^T \beta^*) X^T (\beta^* - \beta) \right] \end{aligned}$$

= 0 by orthogonality of L_2 projection of Y onto linear span of X

$$= \mathbb{E}[(Y - X^T \beta^*)^2] + \mathbb{E}[(\beta^* - \beta)^T X X^T (\beta^* - \beta)]$$

$$= \mathbb{E}[(Y - X^T \beta^*)^2] + \|\beta^* - \beta\|_{\Sigma}^2$$

where $\Sigma = \mathbb{E}[X X^T]$
 component of the
 prediction risk that
 depends on β

Next

$$\mathbb{E}[(Y - X^T \beta^*)^2] = \mathbb{E}[(Y - \mathbb{E}[Y|X] + \mathbb{E}[Y|X] - X^T \beta^*)^2]$$

by orthogonality

$$= \mathbb{E}[(Y - \mathbb{E}[Y|X])^2] + \mathbb{E}[(\mathbb{E}[Y|X] - X^T \beta^*)^2]$$

σ^2
 unavoidable
 variance

non-linearity
 term η^2

So the prediction risk is

$$\mathbb{E}[(Y_{\text{new}} - X_{\text{new}}^T \beta)^2] = \underbrace{\|\beta - \beta^*\|_{\Sigma}^2}_{\text{estimation error}} + \underbrace{\sigma^2}_{\text{irreducible variance}} + \underbrace{\eta^2}_{\text{non-linearity}}$$

↓

by minimizing estimation error we minimize the
 prediction error.

- Suppose we observe a sample of n iid pairs

$$(y_1, x_1), \dots, (y_n, x_n) \in \mathbb{R} \times \mathbb{R}^d$$

from some unknown
joint distribution of (Y, X) .

We will concatenate these observations:

$$Y = \begin{bmatrix} y_1 \\ \vdots \\ y_n \end{bmatrix} \in \mathbb{R}^n$$

$$X = \begin{bmatrix} x_1^T \\ \vdots \\ x_n^T \end{bmatrix}$$

$n \times d$

or following Bach's book

$$\Phi = \begin{bmatrix} \phi(x_1)^T \\ \vdots \\ \phi(x_n)^T \end{bmatrix} = \begin{bmatrix} \Phi_1^T \\ \vdots \\ \Phi_n^T \end{bmatrix}$$

$n \times d$

↓
feature matrix

$\phi(\cdot)$ a feature function
↳ each $\Phi_i \in \mathbb{R}^d$

- We need to estimate β^* (either projection parameter or the actual parameter in a linear model)

We do this by minimizing the empirical MSE

$$\beta \in \mathbb{R}^d \mapsto \hat{R}(\beta) = \frac{1}{n} \sum_{i=1}^n (y_i - \Phi_i^T \beta)^2$$

$$\stackrel{\substack{\text{expected} \\ \text{value wrt to} \\ \text{empirical measure}}}{=} \hat{\mathbb{E}}_n \left[(Y_i - \Phi_i^T \beta)^2 \right]$$

probability measure that puts $\frac{1}{n}$ -mass on each observation

$$= \frac{\|Y - \Phi\beta\|^2}{n}$$

(6)

- The minimizer, $\hat{\beta}$, of $\hat{R}(\beta)$ over all $\beta \in \mathbb{R}^d$ is called the Ordinary Least Squares estimator
OLS

Assume that $\Phi_{n \times d}$ is of full column rank

$$\text{rank}(\Phi) = d \leq n.$$

Then $\hat{\beta} = \arg \min \hat{R}(\beta) = (\underbrace{\Phi^T \Phi}_{\text{well-defined by rank assumption}})^{-1} \Phi^T Y$

$$= (\hat{\mathbb{E}}_n[\Phi \Phi^T])^{-1} \hat{\mathbb{E}}_n[Y \Phi]$$

← same expression for the projection parameter, except we are now using the empirical measure

↳ Plug-in estimator of β^*