

SDS 387 Linear Models

Fall 2024

Lecture 14 - Thu, Oct 15, 2024

Instructor: Prof. Ale Rinaldo

- Announcement: no class on Tue, Oct 22.
- Q2 in HW3: I will clarify the question and update the HW after class

■ Projection of random variables (chapter 16 of Vander Vaart's book on Asymptotic statistics)

Let T and $\{S, S \in \mathcal{S}\}$ be random variables with finite second moments. collection of r.v.'s

We want to project T on \mathcal{S} . A r.v. $\hat{S} \in \mathcal{S}$ is a L_2 -projection of T onto \mathcal{S} when \hat{S}

space of all r.v.'s that have finite variance

minimizes

$$S \in \mathcal{S} \mapsto \mathbb{E}[(T-S)^2]$$

\hat{S} is not necessarily unique!

• Aside : think of L_2 (the space of all r.v.'s with finite 2 moments) as a Hilbert space wrt inner product

$$S_1, S_2 \in L_2 \rightarrow \langle S_1, S_2 \rangle = \mathbb{E}[S_1 S_2]$$

Note: this is not a space of r.v.'s but of equivalence classes of random variables, where 2 r.v.'s are in the same equivalence class when they are identical with prob. 1.

If S is a vector space (closed wrt to addition and scalar multiplication) then \hat{S} is the projection of T onto S iff $T - \hat{S}$ and S are orthogonal
 i.e. $\mathbb{E}[(T - \hat{S})S] = 0 \quad \forall S \in S$:

Thm 11.1 \hat{S} is the projection of T onto S iff
 i) $\hat{S} \in S$ and ii) $\mathbb{E}[(T - \hat{S})S] = 0 \quad \forall S \in S$

The projection is unique (in the sense that if \hat{S}' is another projection then $\mathbb{E}[(\hat{S} - \hat{S}')^2] = 0$).

if S contains the constant functions, then

$$\mathbb{E}[\hat{S}] = \mathbb{E}[T] \quad \text{and} \quad \text{cov}(T - \hat{S}, S) = 0 \quad \forall S \in S.$$

Pf/ The condition $\mathbb{E}[(T - \hat{S})S] = 0$ is called orthogonality

Assume orthogonality. Then, $\forall S \in S$,

$$\mathbb{E}[(T - S)^2] = \mathbb{E}[(T - \hat{S})^2] + 2 \underbrace{\mathbb{E}[(T - \hat{S})(\hat{S} - S)]}_{= 0 \text{ by orthogonality because } \hat{S} - S \in S} + \mathbb{E}[(\hat{S} - S)^2]$$

$$\geq \mathbb{E}[(T - \hat{S})^2]$$

Above, we have an equality iff $\mathbb{E}[(\hat{S} - S)^2] = 0$
 iff $\mathbb{P}(\hat{S} = S) = 1$.

Conversely suppose \hat{S} is a projection. Then $\forall \alpha \in \mathbb{R}$
 $0 \leq \mathbb{E}[(T - \hat{S} - \alpha S)^2] - \mathbb{E}[(T - \hat{S})^2] = \alpha^2 \mathbb{E}[S^2] - 2\alpha \mathbb{E}[(T - \hat{S})S]$

As a function of α , the RHS is a parabola that has to stay above the x -axis. The zeros of this parabola are $\alpha = 0$ and $\alpha = 2 \frac{\mathbb{E}[(T - \hat{S})S]}{\mathbb{E}[S^2]}$

$$\hookrightarrow \mathbb{E}[(T - \hat{S})S] = 0 \quad \forall S \in \mathcal{S}$$

Furthermore, if \mathcal{S} contains the constant r.v.'s then by orthogonality $\mathbb{E}[(T - \hat{S}) \cdot c] = 0 \quad \forall c$

$$\hookrightarrow \mathbb{E}[T] = \mathbb{E}[\hat{S}]$$

Corollary Pythagore theorem for r.v.'s

$$\mathbb{E}[T^2] = \mathbb{E}[\hat{S}^2] + \mathbb{E}[(T - \hat{S})^2]$$

- Arguably the most important type of L_2 -projection is the conditional expectation. Suppose we have 2 r.v.'s, X and Y , with finite 2nd moments. I want to approximate or predict Y using X . Formally, I

want to find the function g s.t. $\mathbb{E}[g^2(X)] < \infty$ and

$$\mathbb{E}[(Y - g(X))^2] \leq \mathbb{E}[(Y - f(X))^2]$$

over all (measurable) functions f s.t. $\mathbb{E}[f^2(X)] < \infty$

In this case $S = \left\{ f(X), f \text{ measurable and } \mathbb{E}[f^2(X)] < \infty \right\}$

It turns out that

$$g(X) = \mathbb{E}[Y|X]$$

you can see this by verifying orthogonality:

$$\mathbb{E}[(Y - \mathbb{E}[Y|X]) f(X)] = \mathbb{E}[Y \cdot f(X)] - \mathbb{E}[\mathbb{E}[Y|X] f(X)]$$

$\forall f$

$= 0$

by law of iterated expectation

aka tower property of conditional expectation

• Aside: a more direct way to see this, without orthogonality

is to use the fact that

$$\arg \min_{c \in \mathbb{R}} \mathbb{E}[(X - c)^2] = \mathbb{E}[X].$$

So,

$$\mathbb{E}[(Y - f(X))^2] \geq \mathbb{E}[(Y - \mathbb{E}[Y|X])^2].$$

$\forall f$

exercise!

• Remark: the conditional expectation is well-defined even without a second moment!

(4)

This is the more general measure-theoretic definition of conditional expectation

LINEAR MODELS

(I will follow Bach's book draft for the next several lectures)

- General regression setting: let Y is a univariate random variable called the response variable. Let $X \in \mathbb{R}^d$ a random vector of covariates or features or explanatory variables

Our goal is to "learn" about Y using X .

- In its most general form "learn" refers to learning the regression function: \rightarrow not random

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

- Assuming that Y and X have finite second moments this can be cast as the problem of minimizing

$$\mathbb{E}[(Y - f(X))^2] \quad \text{over all } f \text{ (with 2nd moment)}$$

\downarrow
MSE (mean squared error)

\downarrow
prediction task!

- There are specific instances of this problem, introduced below in increasing order of generality.

- linear regression models

$$\mathbb{E}[Y | X=x] = \alpha + \beta^T x \quad \text{some } \alpha \in \mathbb{R}, \beta \in \mathbb{R}^d \quad (6)$$

Note that

$$\mathbb{E}[Y|X=x] = \alpha + \beta^T \varphi(x) \text{ also linear}$$

where $\varphi: \mathbb{R}^d \rightarrow \mathbb{R}^d$

is a feature mapping

For example, polynomial regression:

$$\mathbb{E}[Y|X=x] = \alpha + \beta_1 x + \beta_2 x^2 + \beta_3 x^3$$

• non-parametric regression: \rightarrow error or noise

$$Y = f(X) + \varepsilon \quad \text{where} \quad \mathbb{E}[\varepsilon|X] = 0$$

\leftarrow
some unknown f

(in fact, after
times $\varepsilon \perp\!\!\!\perp X$)

$$\mathbb{E}[\varepsilon] = 0$$

without any assumptions on f this
is a hopeless task. Typically f

is assumed to belong to a class of well-behaved
(i.e. smooth) functions.

• most general form of regression (assumption free setting)

$$Y = \mathbb{E}[Y|X] + \underbrace{Y - \mathbb{E}[Y|X]}_{\varepsilon}$$

Note: $\mathbb{E}[\varepsilon|X] = 0$

but, of course
 $\varepsilon \not\perp\!\!\!\perp X$

While you may not be able to estimate

$\mathbb{E}[Y|X]$ you may be interested in approximating
it using a simple, e.g. linear, model!

\downarrow
mis-specified
model (7)

To summarize, the model can be

linear $E[Y | X=x] = \alpha + \beta x$

non-linear $E[Y | X=x] = f(x)$

mis-specified approximate $E[Y | X]$ with a simple model

The covariate X can be treated as

- deterministic: X is not random
- random: X is random

The error term $\varepsilon = Y - E[Y | X]$ can be

- parametric $\varepsilon \sim N_d(0, \sigma^2 I_d) \rightarrow$ homoskedastic
 $N_d(0, \Sigma)$
 \hookrightarrow heteroskedastic

• $\varepsilon \perp\!\!\!\perp X$

• $\varepsilon \not\perp\!\!\!\perp X$