- HW3 is out. For the delta method question you need
  to invoke Lemma 2.12 in von der Vaart's book.

  Let $R: \mathbb{R}^d \to \mathbb{R}$ s.t. $R(0) = 0$. Let $\{X_n\} \subseteq \mathbb{R}^d$
  be a sequence of random vectors s.t. $X_n \xrightarrow{P} 0$.
  Then, $\forall p > 0$

$\forall \{h_n\}$  　i) If $R(h_n) = o(\|h_n\|^p)$ then $R(X_n) = o_P(\|X_n\|^p)$

　　　　iv) if $R(h_n) = O(\|h_n\|^p)$ then $R(X_n) = O_P(\|X_n\|^p)$


- Let's finish the proof of $L_2$ projection.

Thm 11.1 (of von der Vaart) $\hat{S}$ is the projection of $T$ onto $S$

　　iff  　i) $\hat{S} \in S$ and  iv) $\mathbb{E}[(T-\hat{S})S] = 0$

　　　　　　　　　　　　　　　　$\langle T - \hat{S}, S \rangle = 0$  　$\forall s \in S$

orthogonality

①

This projection is unique ( in the sense if $\hat{S}$ is another

projection then $P(\hat{S} \neq \hat{S}') = 0$ )

If S contains the constant functions then

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

where T and S are r.v.'s in $L_2$ and S is a

vector space

PF/ Last time we show that if orthogonality holds then

$\hat{S}$ is the unique projection.

Suppose $\hat{S}$ is a projection ( that is

$\hat{S} \in \underset{s \in S}{\text{argmin}} \; \mathbb{E}[(S-T)^2]$ )

Then $\forall \alpha \in \mathbb{R} \quad \forall s \in S$

$$0 \leq \mathbb{E}\left[\underbrace{(T - \hat{S} - \alpha s)^2}_{\in S}\right] - \mathbb{E}\left[(T-\hat{S})^2\right]$$

$$= \alpha^2 \, \mathbb{E}[s^2] - 2\alpha \, \mathbb{E}\left[(T-\hat{S})s\right]$$

This is a parabola in $\alpha$ that has to stay above

the x-axis. The zeros of this parabola are

$$\alpha = 0 \quad \text{and} \quad \alpha = 2 \frac{\mathbb{E}[(T-\hat{S})s]}{\mathbb{E}[s^2]}$$

$\mapsto$ therefore $\mathbb{E}[(T-\hat{S})s] = 0$ and

(2)

since S is generic, the orthogonality condition

is satisfied.

If $S$ contains the constant functions then, by
orthogonality

$$\mathbb{E}\left[(T-\hat{s}) \cdot c\right] = 0 \qquad \forall c \in \mathbb{R}$$

$$\longrightarrow \mathbb{E}[T] = \mathbb{E}[\hat{s}] \qquad \blacksquare$$

Corollary    Pythagora theorem for r.v.'s.

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

Writing    $\|X\|_2^2 = \mathbb{E}\left[X^2\right]$, this gives us

$$\|T\|^2 = \|\hat{s}\|^2 + \|T-\hat{s}\|^2$$

direct sum
decomposition

- The most important type of $L_2$ projection
  is the conditional expectation. Suppose
  $Y$ and $X$ are $L_2$ random variables and
  let
  $$S = \{ f(x), \quad f \text{ is arbitrary} \text{ s.t.}$$
  $$\mathbb{E}\left[f(x)^2\right] < \infty \}$$

  I want to find the function $g$ s.t. $\mathbb{E}\left[g(x)^2\right] < \infty$
  and    $\mathbb{E}\left[(Y-g(x))^2\right] \le \mathbb{E}\left[(Y-f(x))^2\right]$

(3)

for all $f(X) \in S$.

Then $g(X) = \mathbb{E}[Y \mid X]$. This is true because

$X \longmapsto \mathbb{E}[Y \mid X]$ satisfies the orthogonality condition

$$\mathbb{E}\left[\left(Y - \mathbb{E}[Y \mid X]\right) f(X)\right] = \mathbb{E}\left[Y f(X)\right] - \mathbb{E}\left[\mathbb{E}[Y \mid X] f(X)\right]$$

$\downarrow$

any $f(\cdot)$ s.t. $= 0$     by law of iterated

$\mathbb{E}[f^2(X)] < \infty$        expectation

$$\left( \mathbb{E}\left[\left(\mathbb{E}[Y \mid X]\right)^2\right] \leq \mathbb{E}\left[\mathbb{E}[Y^2 \mid X]\right] = \mathbb{E}[Y^2] < \infty \right)$$

by conditional Jensen

- <u>Remark</u> : Conditional expectation is well-defined even
  without a second moment. In this case
  the defining condition is

$$\mathbb{E}\left[Y \, 1_{\{X \in A\}}\right] = \mathbb{E}\left[\mathbb{E}[Y \mid X] \, 1_{\{X \in A\}}\right]$$

all set $A$.

# ☰ LINEAR REGRESSION

- <u>General regression settings</u>:

response or dependent variable

$Y$     univariate r.v. of interest

- covariates
- independent variable
- features
- explanatory variables

$X$     random vector in $\mathbb{R}^d$

- Our goal is to "model" or "learn" $Y$ using $X$

- Assuming $Y$ and $X$ have finite second moments

$(\ \mathbb{E}[Y^2] < \infty \quad$ and

$$\underline{\underline{\Sigma}} = \mathbb{E}\left[\left(X - \mathbb{E}[X]\right)\left(X - \mathbb{E}[X]\right)^T\right]$$

$$= \mathbb{E}\left[XX^T\right] - \mathbb{E}[X]\left(\mathbb{E}[X]\right)^T$$

exist $)$     this can be cast as the

problem of finding or modeling a function

$f: \mathbb{R}^d \to \mathbb{R} \quad$ s.t.

$$\mathbb{E}\left[\left(Y - f(X)\right)^2\right] \quad \text{is "small".}$$

⑤

- Hierarchy of modeling assumption

Agnostic or model free approach:

$$Y = \mathbb{E}[Y|X] + \underbrace{(Y - \mathbb{E}[Y|X])}_{\varepsilon \quad \text{error}}$$

$$\downarrow$$

best approximation
of $Y$ using $X$

The function $x \in \mathbb{R}^d \longmapsto \mathbb{E}[Y|X=x]$ is the regression function. So letting $f^{opt}(x) = \mathbb{E}[Y|X]$

$$Y = f^{opt}(X) + \varepsilon$$

signal     error

Properties: i) $\mathbb{E}[\varepsilon] = \mathbb{E}[(Y - \mathbb{E}[Y|X])] = 0$

$$\downarrow$$

natural desirable
property of $\varepsilon$

ii) $\mathbb{E}[f^{opt}(X)\,\varepsilon] = \text{Cov}[f^{opt}(X)\,\varepsilon] = 0$

by orthogonality

$\hookrightarrow$ **Remark**   This does not mean

$$\varepsilon \perp\!\!\!\perp f^{opt}(X)$$

(6)

- <u>Signal + independent noise assumption</u>     non-parametric regression

$$Y = f(X) + \varepsilon \qquad \text{where} \qquad \mathbb{E}[\varepsilon] = 0$$
$$\varepsilon \perp\!\!\!\perp X$$

where $f$ belongs a $\left(\text{large}\right)$ class of functions

satisfying certain properties. Typically these

functions are assumed to be "smooth".



- <u>Parametric regression</u>

$$Y = f(x) + \varepsilon \qquad \mathbb{E}[\varepsilon] = 0$$
$$X \perp\!\!\!\perp \varepsilon = 0$$

where $f$ is a function belonging to a parametric

class, i.e.     each $f$ in this class is

parametrized by a parameter $\theta \in \Theta \subseteq \mathbb{R}^m$

            so    regression function if $f_{\theta^a}$ some $\theta^* \in \Theta$



- <u>Linear regression</u>:

$$Y = f_{\theta^*}(X) + \varepsilon \qquad \mathbb{E}[\varepsilon] = 0 \qquad X \perp\!\!\!\perp \varepsilon$$

⑦

where $\qquad f_{\theta^*}(X) = X^T \theta^*$, with

$\theta^*$ unknown but in a set $\Theta \subseteq \mathbb{R}^d$

Remark : to allow for an intercept term, we assume
that the first coordinate of $X$ is a
constant, say $1$. In this case

$$f_{\theta^*}(X) = \theta_1^* + \sum_{j=2}^{d} \theta_j^* X_j$$

- We usually need to have an intercept in your
  model, to make sense of ANOVA decomposition
  and $R^2$ statistic.

- <u>Mis-specified modeling</u> : You are in an agnostic setting
  Fit a parametric model and carry out inference
  on some well-defined "projection parameters"

- The fixed-X setting: the covariates are deterministic

- 2 tasks :       1) Prediction : we want to predict a
                     new instance of response variables,
                     say $Y_{new}$, using a new instance

of the covariates $X_{new}$.

We observe $X_{new}$ and would like to predict $Y_{new}$

2) statistical inference on $\theta^*$ the true parameter indexing the regression function, or the projection parameter if the model is mis-specified