- Reminder: no class on Tue, Oct 22

- HW3, Q2 : it was rewritten and simplified. In the solutions, you will find the following result.

  If $\sqrt{n}\ (\hat{X}_n - \mu) \xrightarrow{d} X$ and $g'(\mu) \neq 0$

  then
  $$n\left[g(\bar{X}_n) - g(\mu)\right] \xrightarrow{d} \frac{g''(\mu)}{2} X^2$$

  Thus, if $\sqrt{n}\ (\hat{X}_n - \mu) \xrightarrow{d} N(0, \sigma^2)$

  $$n\ (\hat{X}_n^2 - \mu^2) \xrightarrow{d} \sigma^2\ \chi_1^2\ (\mu^2)$$

  which is well defined $\forall \mu \in \mathbb{R}$

- $\downarrow$ uses Lemma 2.12 in Van der Vaart

  Let $R: \mathbb{R}^d \to \mathbb{R}$ s.t. $R(0) = 0$. Let $\{X_n\} \subset \mathbb{R}^d$ s.t.

  $X_n \xrightarrow{p} 0$. Then, $\forall p > 0$,

  a) if $R(h) = o(\|h\|^p)$ then $R(X_n) = o_p(\|X_n\|^p)$

$\qquad$ $\text{in)}$ if $R(h) = O\left(\|h\|^{\rho}\right)$ then $R(x_n) = O_p\left(\|X_n\|^{\rho}\right)$

- Last time: linear regression modeling

$\qquad$ $\checkmark$ $\mathbb{E}\left[\,Y\mid X = x\,\right] = x^{\top}\beta$ $\qquad$ some $\beta \in \mathbb{R}^d$

**regression function** *(in blue, left margin)*

$\qquad$ univariate response variable

$\qquad$ $d$-dimensional vector of covariates

Remarks $\quad$ i) $\quad$ linearity here refers to $\beta$. We would call this model:

$$\mathbb{E}[Y\mid X = x] = \phi(x)^{\top}\beta \qquad\qquad \beta \in \mathbb{R}^{d'}$$

also $\quad$ linear, where $\quad \phi : \mathbb{R}^d \to \mathbb{R}^{d'}$ is a _feature vector_ . Example

$$\mathbb{E}[Y\mid X = x] = \alpha_0 + \alpha_1 x + \alpha_2 x^2$$

$\qquad\qquad\qquad\qquad$ is a linear model $\qquad\qquad\qquad\qquad$ (in $(\alpha_0, \alpha_1, \alpha_2)$).

ii) $\quad$ Typically we include an intercept term in the regression function:

$$\mathbb{E}[Y\mid X = x] = \beta_0 + x^{\top}\beta$$

This is important for ANOVA testing and for a correct interpretation of $R^2$ coefficient.

We will always include the intercept, thus we will not write this explicitly. You can think of $X$ or $\phi(X)$ as a vector whose first coordinate

is    non-random and    equal  to    1

• 2   inferential    tasks:
    1)      statistical inference  about  $\beta$

    2)      prediction

## Statistical inference

• If    the model is well-specified  (i.e.  $\mathbb{E}[Y|X=x]=\beta^T x$
  then    $\beta$  is  clearly the  parameter of interest.

• What  if  $\mathbb{E}[Y|X=x]$  is not  linear?  In    this
  mis-specified setting  we  need to first identify the
  target  parameter.    This  can  be  defined  by
  looking  at  the  best  linear  approximation  to
  $\mathbb{E}[Y|X]$:

$$\beta^* = \underset{\beta \in \mathbb{R}^d}{\text{argmin}}  \mathbb{E}\left[\left(\mathbb{E}[Y|X] - X^T\beta\right)^2\right]$$

  This  is  well-defined  and  unique  provided  that
  $Y$ and    $X$  have  2  moments;  in  particular
  $\Sigma_1 = \mathbb{E}[X X^T]$  needs  to  be invertible.
  Also,  $\beta^*$  is also  equal to

$$\underset{\beta \in \mathbb{R}^d}{\text{argmin}}  \mathbb{E}\left[(Y - X^T\beta)^2\right]$$

**Thm** If $\Sigma$ is invertible and $Y$ has 2 moments

$$\beta^* = \Sigma^{-1} \mathbb{E}[Y \cdot X]$$

**Pf** $\beta^*$ is the minimizer of

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

over all $\beta \in \mathbb{R}^d$

Because of moment assumptions we can take the derivative wrt to $\beta$ inside the expectation and obtain the first order optimality condition

$$\mathbb{E}\left[\cancel{2}\, XX^T\beta\right] - \cancel{2}\,\mathbb{E}\left[\mathbb{E}[Y|X] \cdot X\right] = 0$$

Solution is $\mathbb{E}[XX^T]^{-1}\mathbb{E}[Y \cdot X]$

By convexity this is unique! 🔲

**Remark**: What is $\beta^*$? It it is the vector of coefficients of the $L_2$ projection of $Y$ onto the linear span of $X$ ( the vector space of all linear functions of $x$ ) !

• $\beta^*$ is vector measuring linear association between $Y$ and $X$

# Prediction

In prediction (the main objective of ML models), we want to predict a new response, say $y^{new}$, using $X^{new}$. Our goal is the to minimize the prediction error, i.e. to solve the problem

$$\min_{\beta \in \mathbb{R}^d} \mathbb{E}\left[ \left(y^{new} - (X^{new})^T \beta\right)^2 \right]$$

↓ prediction MSE

Of course the solution is $\beta^*$. Suppose we instead use a different vector $\beta \in \mathbb{R}^d$. How large is the error that we make by using the wrong $\beta$?

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

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

$$+ 2 \underbrace{\mathbb{E}\left[\left(y - X^T\beta^*\right) X^T(\beta^* - \beta)\right]}_{= 0 \quad \text{by orthogonality of } L_2 \text{ projection}}$$

$$= \underbrace{\mathbb{E}\left[\left(y - X^T\beta^*\right)^2\right]}_{\text{systematic error}} + \underbrace{(\beta^* - \beta)^T \Sigma (\beta^* - \beta)}_{\|\beta^* - \beta\|^2_\Sigma}$$

If $\mathbb{E}[Y|X] = X^T \beta^*$ then the systematic

error is usually written as $\sigma^2$, the variance of $Y$

$$\left( \text{E.g. assuming: } Y = X^T \beta + \varepsilon \quad \text{where} \quad \varepsilon \sim (0, \sigma^2) \quad \varepsilon \perp\!\!\!\perp X \right)$$

and
$$\| \beta^* - \beta \|_\Sigma^2 \quad \text{is a measure of how well we}$$
are estimating the true regression function

# 2  DATA

- Suppose we observe a sample $(Y_1, X_1), ..., (Y_n, X_n)$
  of $n$ points of iid realizations from the
  joint distribution of $Y$ and $X$.

- I will write $Y = \begin{bmatrix} Y_1 \\ \vdots \\ Y_n \end{bmatrix} \in \mathbb{R}^n$ and

  $$\underset{n \times d}{X} = \begin{bmatrix} X_1 & \cdots & X_n \end{bmatrix}^T \quad \text{or} \quad \underset{n \times d}{\Phi} = \begin{bmatrix} \varphi(X_1) & \cdots & \varphi(X_n) \end{bmatrix}^T$$

- To estimate $\beta^*$ we will minimize the empirical
  MSE or predictive risk :

  $$\hat{R}(\beta) = \frac{1}{n} \sum_{i=1}^{n} \left( Y_i - \Phi(X_n)^T \beta \right)^2$$

  $$\overset{\text{expectation}}{\underset{\substack{\text{wrt} \\ \text{empirical} \\ \text{measure}}}{\Longleftarrow}} \quad = \quad \hat{\mathbb{E}}_n \left[ \left( Y - \Phi(X)^T \beta \right)^2 \right]$$

  $$= \quad \frac{1}{n} \| Y - \Phi \beta \|^2$$

- The OLS $\longrightarrow$ ordinary least squares estimator of $\beta^*$ is the minimizer of $\hat{R}(\beta)$ over all $\beta \in \mathbb{R}^d$.

- <u>Thm</u>  Assume that $\Phi$ is of full column rank.
  Then $\left(\text{rank}(\Phi) = d \leq n\right)$

$$\hat{\beta} = \underset{\beta \in \mathbb{R}^d}{\text{argmin}} \; \hat{R}(\beta) = \left(\Phi^T \Phi\right)^{-1} \Phi^T Y$$

$$= \hat{\Sigma}_n^{-1} \frac{\Phi^T Y}{n}$$

where $\hat{\Sigma}_n = \frac{\Phi^T \Phi}{n} = \frac{1}{n} \sum_{i=1}^{n} \Phi_i^T \Phi_i$

$\Phi_i$ $i^{th}$ row of $\Phi$

Notice that $\hat{\beta} = \left(\mathbb{E}_n\left[\Phi(X) \Phi(X)^T\right]\right)^{-1} \mathbb{E}_n\left[Y \, \Phi(X)\right]$

where $\mathbb{E}_n[\cdot]$ expectation wrt empirical measure

- $\hat{\beta}$ is the plug-in estimator for $\beta^*$

Pf/ $\beta \longrightarrow \hat{R}(\beta)$ is strictly convex because $\Phi$ is of full columns rank

[ strict convexity follows because the Hessian of $\hat{R}(\beta)$ is $\frac{\Phi^T \Phi}{n}_{d \times d}$ which is pd by assumption ]

So the minimizer is found by first order optimality condition : $\nabla \hat{R}(\beta) = 0$ solve for $\beta$

This yields:

$$\nabla \hat{R}(\beta) = -\frac{2}{\eta} \Phi^T \left( Y - \Phi\beta \right) = 0$$

↳ Normal equations

Solution is

$$\hat{\beta} = \left( \Phi^T \Phi \right)^{-1} \Phi^T Y$$

invertible by
assumption