# MATH 423
# Applied Regression

Optimal prediction for a single r.v. :
Say we have an r.v. $Y$, whose distribution is unknown.

Theorem : The optimal prediction $m^*$ minimizes the expected mean squared error,
$$m^* = \arg\min_m \mathbb{E}((Y-m)^2)$$

The optimal predictor of $Y$ is $\mathbb{E}(Y)$, i.e. $m^* = \mathbb{E}(Y)$.

→ Proof :
$$\frac{\partial \, MSE}{\partial m} = \frac{\partial \, \mathbb{E}((Y-m)^2)}{\partial m} \qquad \text{recall that } Var(X) = \mathbb{E}(X^2) - \mathbb{E}(X)^2$$

$$= \frac{\partial [(\mathbb{E}(Y-m))^2 + Var(Y-m)]}{\partial m}$$

$$= \frac{\partial}{\partial m} (\mathbb{E}(Y) - m)^2 + \frac{\partial}{\partial m} Var(Y)$$

$$= \frac{\partial}{\partial m} (\mathbb{E}(Y)^2 - 2m\,\mathbb{E}(Y) + m^2)$$

$$= -2\mathbb{E}(Y) + 2m$$

$$\frac{\partial \, MSE}{\partial m} := 0 \Rightarrow m^* = \mathbb{E}(Y) \quad ▨$$

In other words, the best 1 number guess we could make for $Y$ is just its expected value.

What if we used something else other than MSE to measure how good the pred. is ?
e.g. mean absolute deviation (MAD)
$$MAD = \mathbb{E}(|Y-m|)$$
$$\Rightarrow m^* = median(Y)$$
The median is more stable than the mean, especially if data has a lot of outliers.
But solving for the $m^*$ is more difficult as this includes taking the derivative of the
absolute value →???

Optimal prediction of an r.v. from other variables
- input variables —model→ output variables
- input vars $(X)$ are vars that are correlated with output, aka features, predictors
- output vars $(Y)$ measure the outcome of interest, aka dep. var.

Note that predictive models do not show causation, only correlation!

Say we have 2 r.v.s, $X, Y$. The joint distribution of $(X,Y)$ is known. We would
like to predict $Y$ using $X$.

Theorem: The optimal prediction function $m^*(\cdot)$ minimizes the expected MSE :
$$m^*(\cdot) = \arg\min_{m(\cdot)} \mathbb{E}_{X,Y}[(Y - m(X))^2]$$

The optimal predictor of $Y$ is $\mathbb{E}(Y|X=x)$ i.e. $m^*(x) = \mathbb{E}(Y|X=x)$.

$\rightarrow$ **Proof**: Denote $\mu(x) = \mathbb{E}_{Y|X}(Y|X=x)$

We want to prove that $m^*(x) = \mu(x)$.

$$\mathbb{E}_{x,y}\left[(Y-m(x))^2\right] = \mathbb{E}_{x,y}\left[(Y-\mu(x)+\mu(x)-m(x))^2\right]$$

$$= \mathbb{E}_{x,y}\left[(Y-\mu(x))^2\right] + 2\,\mathbb{E}_{x,y}\left[(Y-\mu(x))(\mu(x)-m(x))\right]$$

$$+ \mathbb{E}_{x,y}\left[(\mu(x)-m(x))^2\right]$$

$$= \mathbb{E}_{x,y}\left[(Y-\mu(x))^2\right] + \mathbb{E}_{x}\left[(\mu(x)-m(x))^2\right]$$

Notice that $\mathbb{E}_{x,y}\left[(Y-\mu(x))(\mu(x)-m(x))\right] = 0$:

$$\mathbb{E}_{x,y}\left[(Y-\mu(x))(\mu(x)-m(x))\right] = \mathbb{E}_{x}\left[\mathbb{E}_{Y|X}(Y-\mu(x))(\mu(x)-m(x))\,|\,X\right]$$

$$= \mathbb{E}_{x}\left[(\mu(x)-m(x))\,\underbrace{\mathbb{E}_{Y|X}\left[Y-\mu(x)|X\right]}_{=0 \text{ by defn of } \mu(x)}\right]$$

From $\mathbb{E}_{x,y}\left[(Y-m(x))^2\right] = \mathbb{E}_{x,y}\left[(Y-\mu(x))^2\right] + \mathbb{E}_{x}\left[(\mu(x)-m(x))^2\right]$,
it must be that MSE is minimized when $\mathbb{E}_{x}\left[(\mu(x)-m(x))^2\right]$ is minimized.

$$\mathbb{E}\left[(\mu(x)-m(x))^2\right] = 0 \quad \text{when} \quad m(x) = \mu(x):$$

$$m^*(x) = \mu(x)$$

$$= \mathbb{E}_{Y|X}(Y|X=x).$$

Thus the expected MSE-optimal prediction is made by using $\mu(x)$.

We call $\mathbb{E}(Y|X=x)$ the regression function.

# Lecture 3 (Sep 10)    KNN Regression

How might we estimate the regression function?
$$\downarrow$$
$$m^*(x) = \mathbb{E}(Y|X=x)$$

Given $n$ observations: $(x_i, y_i)$, $1 \le i \le n$, $x_1, \ldots, x_n \in \mathbb{R}^p$

$\mathbb{E}(Y|X=x) = ?$

- average? $\hat{m}(x) = $ average $(\{y_i : x_i = x\})$   ← observations

     ie for an $x_i = $ some $x$ value we set, take the mean of the $y_i$ observations that correspond to $x_i = x$.

     → problem with this: what if there's only 1 point at $x_i = x$? Cannot estimate directly
     We can relax our criteria and use data points that are in the vicinity of the $x_i = x$ point.

- average of nearest neighbours:

$$\hat{m}(x) = \text{average } (\{y_i : x_i \text{ equal to or close to } x\})$$
$$= \frac{1}{k} \sum_{x_i \in N_k(x)} y_i$$

   - where $N_k(x)$ is the neighbourhood of $x$ defined by the $k$-closest points $x_i$ in the training sample.   i.e. ↗
   - $k$ is a hyper-parameter that we set    Euclidean distance.

e.g. our dataset $(n=6)$, $X = (x_1, x_2, x_3)$, $y$

| $x_{i1}$ | $x_{i2}$ | $x_{i3}$ | $y_i$ |
|---|---|---|---|
| 1 | 1 | 0 | 3 |
| 2 | 0 | 3 | 2 |
| 0 | 1 | 2 | 0 |
| 0 | 1 | 3 | 1 |
| -1 | 1 | 1 | 4 |
| 1 | 1 | 2 | 3 |

Idea: use KNN regression to make predictions about $y$.

⇒ what is $\hat{m}(x)$?

Let $x_0 = (0, 0, 0)$.   $y_0 = ?$

→ Compute the Euclidean dist. between each obs and $x_0$.
- $\|x_1 - x_0\|_2 = \sqrt{(1-0)^2 + (1-0)^2 + (0-0)^2}$
        $= \sqrt{2} \approx 1.41$
- $\|x_2 - x_0\|_2 = 3.61$
- $\|x_3 - x_0\|_2 = 2.23$
- $\|x_4 - x_0\|_2 = 3.16$
- $\|x_5 - x_0\|_2 = 1.73$
- $\|x_6 - x_0\|_2 = 2.44$

$x_0$ neighbours in order of closeness: $\{x_1, x_5, x_3, x_6, x_4, x_2\}$

our formula is $\hat{m}(x) = \frac{1}{k} \sum_{x_i \in N_k(x)} y_i$

→ KNN prediction when $k=1$?

$N_1(x_0) = \{x_1\}$

$\hat{m}(x_0) = y_1 = \underline{3}$ //

→ KNN prediction when $k=3$?

$N_3(x_0) = \{x_1, x_5, x_3\}$

$\hat{m}(x_0) = \frac{1}{3} \sum_{x_i \in N_3(x_0)}^{3} y_i$

$= \frac{1}{3}(y_1 + y_5 + y_3)$

$= \frac{1}{3}(3 + 0 + 4) = 2.33$

(R) R computation for KNN prediction in handouts

# Lec 4 (Sep 15)    Limitations of KNN

## Theoretical guarantee of KNN:

- $X \in \mathbb{R}^p$, $Y \in \mathbb{R}$
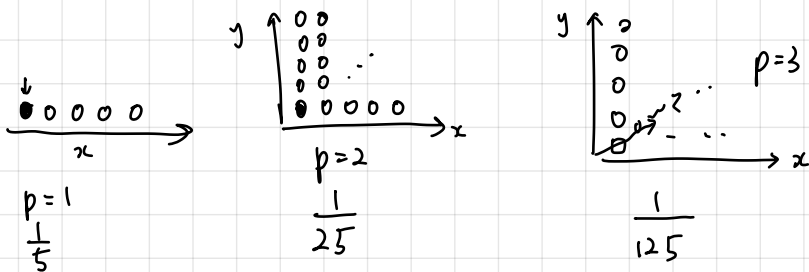  ↳ vector of $p$ variables. our predictors
- KNN can be pretty good for small $p$ (i.e. $p < 4$) and large $N$.

under mild regularity conditions on joint prob. dist. of $\Pr(X, Y)$, one can show that as
$N, k \to \infty$ such that $\frac{k}{N} \to 0$  (i.e. $k$ grows at a slower rate than $N$),
then: $\hat{m}(x) = \frac{1}{k} \sum_{x_i \in N_k(x)} y_i \longrightarrow \mathbb{E}(Y | X = x)$
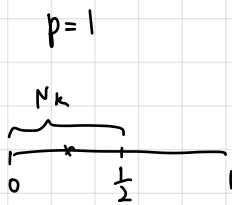
## Curse of dimensionality
- as we increase the dimension of predictors, model performance is affected.
- for small $p$, we can always find a fairly large neighbourhood of observations close to the target $x$ and take their average.

- But for very large $p$, it's more difficult to find a neighbourhood that contains enough data.

$p=1$
$\frac{1}{5}$

$p=2$
$\frac{1}{25}$

$p=3$
$\frac{1}{125}$

The "area" that a neighbourhood takes up in the data space grows smaller and smaller

- Consider the KNN for inputs X uniformly distributed in a $p$-dim unit hypercube
- Suppose we want to predict at a target X, set up a neighborhood around X to capture a fraction $R$ of the observation.

$R = \frac{1}{2}$          $p=1$                    $p=2$



Expected edge length ("radius") of the neighborhood:

$$L_p(R) = R^{\frac{1}{p}}$$

eg  $L_{10}(0.01) = 0.63$ → to capture 1% of data, in a 10-dimensional case, radius is 0.63

$L_{10}(0.1) = 0.8$     * in our eg, it's a unit, so 0.8 radius is
$\uparrow$                almost the whole space. This is not local anymore!
10% of data

To capture 1% or 10% of the data to form a local average, we must cover 63% or 80% respectively of the range of each input variable. Such neighborhoods are not "local" anymore!

# Optimal Linear Prediction

In general, $\hat{m}(x) = \mathbb{E}(Y \mid X = x)$ might be a very complicated form:

| model | $\mathbb{E}(Y \mid X = x)$ |
|---|---|
| KNN | $\frac{1}{k} \sum_{x_i \in N_k(x)} y_i$ |
| linear regression | $x'\beta$ |
| additive model | $f_1(x_1) + \cdots + f_p(x_p)$ ,  $X = (x_1, \ldots, x_p)$ |
| decision tree | $T(x)$ — nonlinear function |
| random forest / gradient boosting | $\sum_{m=1}^{M} \beta_m T_m(x)$ |
| deep learning | $G(\sigma(w_x))$ |
| SVM | $\sum_i \alpha_i K(x, x_i)$ |

To simplify $m(x)$, we restrict $m(x)$:

$$m(x) = \beta_0 + \beta_1 X$$

only 1 predictor $X \in \mathbb{R}$

→ What are the optimal values of $\beta_0$ and $\beta_1$?

Theoretical ans:
- if distribution $(X, Y)$ known ← MSE

$$(\beta_0^*, \beta_1^*) = \underset{(\beta_0, \beta_1)}{\text{argmin}} \; \mathbb{E}_{X,Y}\left[(Y - m(X))^2\right]$$

$$= \underset{(\beta_0, \beta_1)}{\text{argmin}} \; \mathbb{E}_{X,Y}\left[Y - (\beta_0 + \beta_1 X))^2\right]$$

$$\Rightarrow \beta_1^* = \frac{\text{Cov}(X, Y)}{\text{Var}(X)} \quad , \quad \beta_0^* = \mathbb{E}(Y) - \beta_1^* \mathbb{E}(X)$$
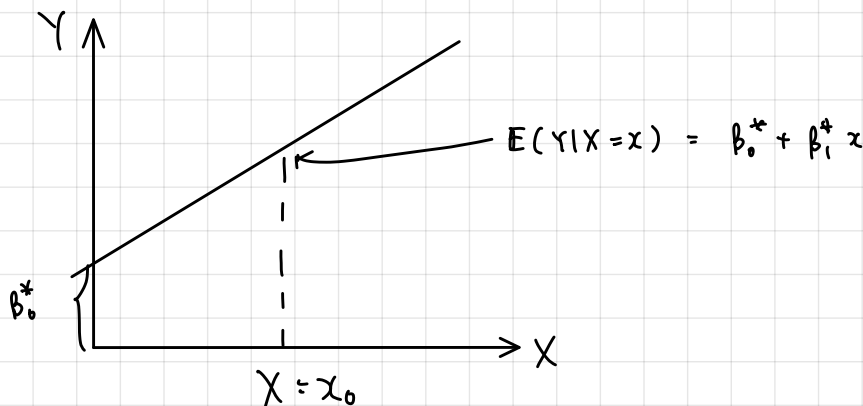
Proof is in handout : 1. decompose MSE
2. set derivs to 0
3. plug and solve

$\Rightarrow$ optimal linear predictor: $m^*(X) = \beta_0 + \beta_1 X$

$$= \mathbb{E}(Y) - \frac{\text{Cov}(X, Y)}{\text{Var}(X)} \mathbb{E}(X) + \frac{\text{Cov}(X, Y)}{\text{Var}(X)} X$$

# Lec 5 (Sep 17)

The line $\beta_0^* + \beta_1^* X$ is called the optimal prediction line
i.e. linear regression function.



$\mathbb{E}(Y|X=x) = \beta_0^* + \beta_1^* x$

$X = X_0$

Some notes about the optimal predictor line:

1. The optimal predictor line passes through $(E(X), E(Y))$

$$m^*(x) = \beta_0^* + \beta_1^* x$$
$$= \underbrace{E(Y) - \beta_1^* E(X)}_{\beta_0^*} + \beta_1^* x$$
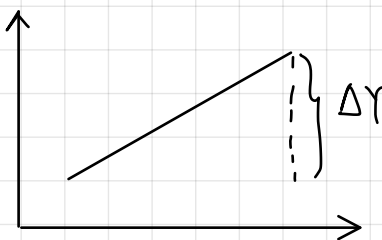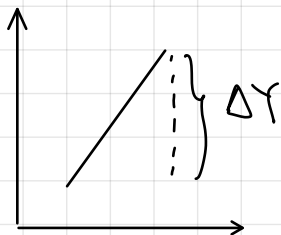
$$\Rightarrow m^*(E(X)) = E(Y)$$

2. If $X$ and $Y$ are "centered", i.e. $E(Y) = E(X) = 0$, the optimal regression line passes through $(0,0)$ since $\beta_0^* = 0$.

Recall that optimal slope $\boxed{\beta_1^* = \dfrac{Cov(X,Y)}{Var(X)}}$ .

3. The optimal slope $\beta_1^*$ increases as $Cov(X,Y)$ increases.

4. The optimal slope $\beta_1^*$ decreases as $Var(X)$ increases.



For the same $\Delta$ in $Y$, larger variance of $X$ leads to flatter slope.

5. The optimal slope $\beta_1^*$ does not change if we use instead $Y-c$ and $X-c$. But the intercept $\beta_0^*$ will change.

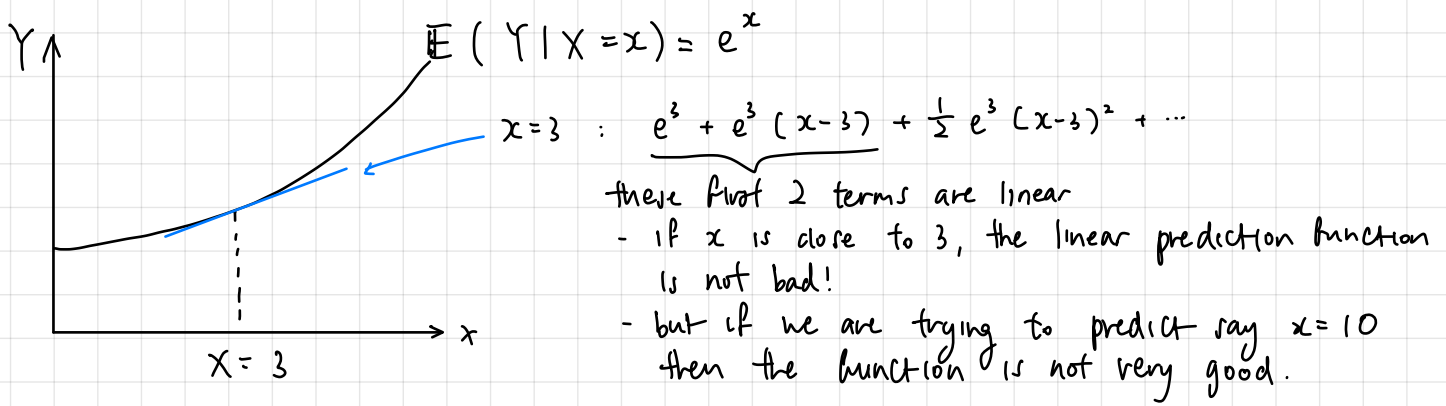6. Non-linear pattern cannot be appropriately modelled.

eg: Imagine true regression function:
$$E(Y|X=x) = e^x \quad \longleftarrow \text{(this is nonlinear)}$$

Taylor's expansion at $X = x_0$
$$e^x = e^{x_0} + \left.\frac{\partial e^x}{\partial x}\right|_{x=x_0} (x - x_0) + \frac{1}{2} \frac{\partial^2 e^x}{\partial x^2} (x-x_0)^2 + \cdots$$

$$= e^{x_0} + e^{x_0}(x - x_0) + \frac{1}{2} e^{x_0}(x-x_0)^2 + \cdots$$

$$E(Y|X=x) = e^x$$

$x=3$ : $\underline{e^3 + e^3(x-3) + \frac{1}{2}e^3(x-3)^2 + \cdots}$

these first 2 terms are linear
- if $x$ is close to 3, the linear prediction function is not bad!
- but if we are trying to predict say $x=10$ then the function is not very good.

If we want to use the linear component as a prediction function, need to make sure that it dominates the quadratic term

i.e. $\frac{1}{2}e^{x_0}|x-x_0|^2 \ll e^{x_0}|x-x_0|$

$$\frac{|x-x_0|^2}{|x-x_0|} < \frac{2e^{x_0}}{e^{x_0}}$$

$$|x-x_0| < 2$$

$\Rightarrow$ In our example, $|x-3| < 2$

$\qquad\qquad\qquad x-3 < 2 \quad$ or $\quad x-3 > -2$

$\qquad\qquad\qquad x < 5 \quad$ or $\quad x > 1$

Linear prediction is good for $1 < x < 5$


## Plug-in estimation (estimating optimal linear prediction using data)

How to estimate optimal linear prediction $m^*(x) = E(Y|X=x) = \beta_0^* + \beta_1^* x$
from $n$ observations of data $(x_1, y_1), \ldots (x_n, y_n)$?

$$\Rightarrow \hat{\beta}_1 = \frac{\widehat{Cov(X,Y)}}{\widehat{Var(X)}}$$

$$= \frac{\sum_{i=1}^{n}(y_i - \bar{y})(x_i - \bar{x})}{\sum_{i=1}^{n}(x_i - \bar{x})^2} , \qquad \bar{y} = \frac{1}{n}\sum_{i=1}^{n}y_i , \quad \bar{x} = \frac{1}{n}\sum_{i=1}^{n}x_i$$

$$\hat{\beta}_0 = \widehat{E(Y)} - \hat{\beta}_1 \widehat{E(X)}$$

$$= \bar{y} - \hat{\beta}_1 \bar{x}$$

The fitted regression line: $\hat{m}(x) = \hat{\beta}_0 + \hat{\beta}_1 x$

- an approximation of the true regression line, calculated with data.

$$\hat{m}(x) = \hat{\beta}_0 + \hat{\beta}_1 x$$

$$= \left[ \bar{y} - \frac{\sum_{i=1}^{n}(y_i - \bar{y})(x_i - \bar{x})}{\sum_{i=1}^{n}(x_i - \bar{x})^2} \right] + \left( \frac{\sum_{i=1}^{n}(y_i - \bar{y})(x_i - \bar{x})}{\sum_{i=1}^{n}(x_i - \bar{x})^2} \right) x$$

see handout for proof.

idea: plug-in sample values, i.e. $E(Y) \approx \bar{y}$ } sample mean
$E(X) \approx \bar{x}$

we could use $\frac{1}{n-1}$ too for sample but it doesn't matter when $n$ is large.

sample variance
$Var(X) = E((X - EX)^2) \approx \frac{1}{n}\sum_{i=1}^{n}(x_i - \bar{x})^2$

$Cov(X,Y) = E((X - EX)(Y - EY))$
$$\approx \frac{1}{n}\sum_{i=1}^{n}(y_i - \bar{y})(x_i - \bar{x})$$
sample covariance

Some notation:

$$\hat{\beta}_1 = \frac{S_{xy}}{S_{xx}},$$ where $S_{xy} = \sum_{i=1}^{n}(y_i - \bar{y})(x_i - \bar{x}) = \sum_{i=1}^{n}(x_i - \bar{x})y_i = \sum_{i=1}^{n}x_i y_i - n\bar{x}\bar{y}$

and $S_{xx} = \sum_{i=1}^{n}(x_i - \bar{x})^2 = \sum_{i=1}^{n}(x_i - \bar{x})x_i = \sum_{i=1}^{n}x_i^2 - n(\bar{x})^2$

$$\hat{\beta}_0 = \bar{y} - \hat{\beta}_1\bar{x} = \bar{y} - \frac{S_{xy}}{S_{xx}}\bar{x}$$

useful result : $\sum_{i=1}^{n}(y_i - \bar{y}) = 0$ , $\sum_{i=1}^{n}(x_i - \bar{x}) = 0$

## Some notes:

1. As $n \to \infty$, $\hat{\beta}_1 \longrightarrow \beta_1^*$ and $\hat{\beta}_0 \to \beta_0^*$

2. The fitted regression line $\hat{m}(x)$ passes through $(\bar{x}, \bar{y})$.

3. If the data is centered, $x_i' = x_i - \bar{x}$ , $y_i' = y_i - \bar{y}$ $\forall i = 1, \dots, n$
$$\bar{x}' = \frac{1}{n}\sum_{i=1}^{n}x_i' = 0 \quad , \quad \bar{y}' = \frac{1}{n}\sum_{i=1}^{n}y_i' = 0,$$
then the fitted line passes through $(0, 0)$.

4. Slope does not change under a shift of the data
$$\text{eg } y_i' = y_i - c \quad , \quad x_i' = x_i - c'$$

# Lec 6 (Sep 22)   <mark>Simple Linear Regression</mark>

"simple": 1 predictor for 1 DV

We are interested in the underlying mechanism / DGP.

Recall 2 goals of regression:

1. prediction: forecast unobserved data   (cause $\longrightarrow$ effect)

2. inference: identify unknown data generating process

   In order to do inference, we have to make <u>assumptions</u> and believe our data is generated from the SLR model.
   $\rightarrow$ we think that the model approximates reality relatively well.

## SLR model specification:

- consider 2 rvs: $Y$ (the response / DV) , $X_1 \swarrow$ only 1 predictor
- our goal here is to predict $Y$ using $X_1$.
- assumptions:

  A1. <mark>arbitrary input</mark>: The distribution of $X_1$ is arbitrary,
      and $X_1$ is even non-random.

  A2. <mark>linear function and additive error</mark>:   $Y = \underbrace{\beta_0 + \beta_1 X_1}_{\text{linear fcn}} + \underbrace{\varepsilon}_{\text{error}}$

      e.g. if $X_1 = x_1$ , then $Y = \beta_0 + \beta_1 x_1 + \varepsilon$ , for unknown
           coefficients $\beta_0$ and $\beta_1$ , and random noise $\varepsilon$.

  A3. <mark>zero mean and constant variance error</mark> (homoskedasticity):
      $\mathbb{E}(\varepsilon) = 0$ ,  $\text{var}(\varepsilon) = \sigma^2 \geqslant 0.$ ,  $\sigma^2$ is unknown

  A4. <mark>independent error</mark>: $\varepsilon$ is independent of $X$.

## some notes about SLR:

- if $E(\varepsilon) = c \neq 0$, we can find an rv $\varepsilon'$ with $E(\varepsilon') = 0$ and $\mathrm{var}(\varepsilon') = \sigma^2$, s.t. $\varepsilon = \varepsilon' + c$.

$$Y = \beta_0 + \beta_1 X_1 + \varepsilon$$
$$= \beta_0 + \beta_1 X_1 + \varepsilon' + c$$
$$= \underbrace{(\beta_0 + c)}_{\beta_0'} + \beta_1 X_1 + \varepsilon'$$

- SLR assumptions actually imply that the true (optimal) regression line is linear and $Y$ has constant variance.

$$E(Y \mid X_1 = x_1) = E(\beta_0 + \beta_1 X_1 + \varepsilon \mid X_1 = x_1)$$
$$= E(\beta_0 + \beta_1 X_1 \mid X_1 = x_1) + \underbrace{E(\varepsilon)}_{=0}$$
$$= \beta_0 + \beta_1 x_1$$

←— this is a realized value

$$\mathrm{Var}(Y \mid X_1 = x_1) = \mathrm{Var}(\beta_0 + \beta_1 X_1 + \varepsilon \mid X_1 = x_1)$$
$$= \underbrace{\mathrm{var}(\beta_0)}_{0} + \underbrace{\mathrm{var}(\beta_1 x_1)}_{0} + \mathrm{var}(\varepsilon)$$
$$= \mathrm{Var}(\varepsilon)$$
$$= \sigma^2$$

$$\Rightarrow Y \mid X_1 = x_1 \overset{d}{\sim} (\beta_0 + \beta_1 x_1, \ \sigma^2)$$

# Lec7 (Sep 24)

The noise variable $\varepsilon$ can represent:
- other factors not considered in the model
- measurement error
- or some combination of both.

We think of the SLR assumptions A1 – A4 as modeling decisions that (we hope) will be useful, rather than as facts about the actual underlying relationship between $Y$ and $X_s$.

## Interpretation of parameters:
- $\beta_0$ = intercept, the expected value of $Y$ when $X_1$ is 0.

$$E(Y \mid X_1 = 0) = \beta_0 + 0 \cdot \beta_1 = \beta_0$$

- $\beta_1$ = slope, diff between exp. value of $Y$ when $x_1$ is shifted by 1.
$$E(Y \mid X_1 = x_1) - E(Y \mid X_1 = x_1 + 1) = \beta_0 + \beta_1 x_1 - \beta_0 - \beta_1 x_1 - \beta_1$$
$$= \beta_1$$

Note that $\beta_1$ does not imply causality! It is only statistical association.

If we select 2 sets of cases from $(X_1, Y)$ distribution where $X_1$ differs by 1, we expect the associated $Y$ to differ by $\beta_1$.

- $\sigma^2$ = error variance, the variance of the noise around the reg. line. It represents a typical distance of a point from the true regression line.

## Model set up for multiple data points:

We assume multiple data points $(X_{11}, Y_1)$, $(X_{21}, Y_2)$, ... , $(X_{n1}, Y_n)$ are generated from the same model.

$$Y_i = \beta_0 + \beta_1 X_{i1} + \varepsilon_i \quad , \quad i = 1, ..., n$$
$$\text{where } E(\varepsilon_i) = 0 , \quad Var(\varepsilon_i) = \sigma^2 \quad \forall i ,$$
$$\text{and} \quad \varepsilon_i \text{ and } \varepsilon_j \text{ are independent for } i \neq j$$

Equivalently, $E(Y_i \mid X_{i1} = x_{i1}) = \beta_0 + \beta_1 x_{i1}$
$$Var(Y_i \mid X_{i1} = x_{i1}) = \sigma^2 \quad , \quad i = 1, ..., n$$

### Question: $Y_i \perp Y_j$ (independent)?

$$Y_i \perp Y_j \mid X_{i1}, X_{i2} \quad \text{(conditionally independent)}?$$

$Y_i$ and $Y_j$ are not unconditionally independent:

$$Cov(Y_i, Y_j) = Cov(\beta_0 + \beta_1 X_i + \varepsilon_i , \beta_0 + \beta_1 X_j + \varepsilon_j)$$
$$= \beta_1^2 Cov(X_i, X_j) + \beta_1 \underbrace{Cov(X_i, \varepsilon_j)}_{=0} + \beta_1 \underbrace{Cov(\varepsilon_i, x_j)}_{=0} + \underbrace{Cov(\varepsilon_i, \varepsilon_j)}_{=0}$$

because $X$ and $\varepsilon$ are indep.

$$= \beta_1^2 Cov(X_i, X_j)$$

If $X_i \perp X_j$, $Cov(X_i, X_j) = 0 \Rightarrow Cov(Y_i, Y_j) = 0$

But if not $X_i \perp X_j$, $Cov(Y_i, Y_j) \neq 0$

$$\text{Cov}(Y_i, Y_j \mid X_i = x_i, X_j = x_j) = \text{Cov}(\varepsilon_i, \varepsilon_j \mid X_i = x_i, X_j = x_j)$$
$$= \text{Cov}(\varepsilon_i, \varepsilon_j)$$
$$= 0$$

↖ as realized values, they become constants and $\text{Cov} = 0$.

Ⓡ Generate data from SLR model (simulation)

## Optimal prediction for SLR

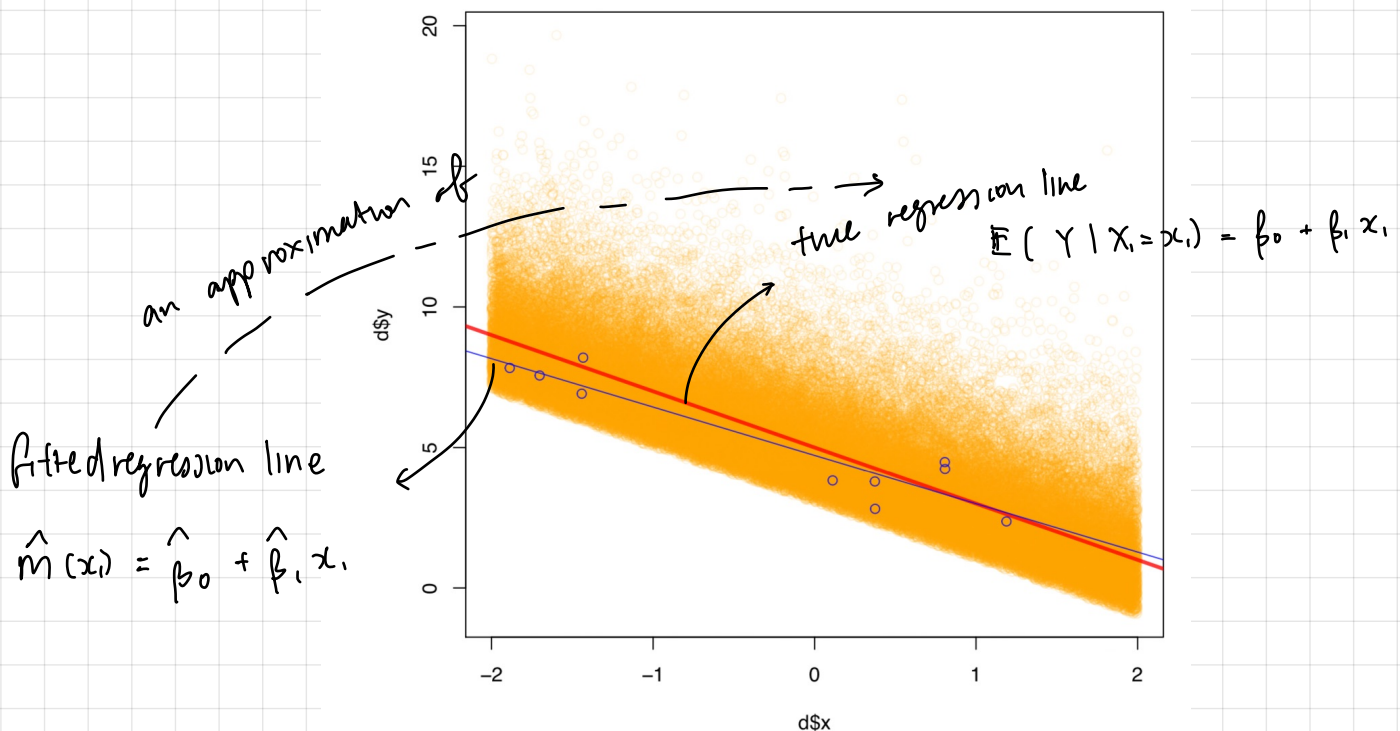Assume $(X_1, Y)$ are generated from the SLR model:
$$Y = \beta_0 + \beta_1 X_1 + \varepsilon, \quad \mathbb{E}(\varepsilon) = 0, \quad \text{Var}(\varepsilon) = \sigma^2, \quad \varepsilon \perp X_1$$

If we want to predict $Y$ using $X_1$, what is the optimal prediction that minimizes the mean squared error?

$$m^*(\cdot) = \underset{m(\cdot)}{\text{argmin}} \; \mathbb{E}_{X,Y}\left[(Y - m(X_1))^2\right]$$

$$\Rightarrow m^*(x_1) = \mathbb{E}(Y \mid X_1 = x_1) \quad \text{is the optimal prediction function}$$
$$= \mathbb{E}(\beta_0 + \beta_1 X_1 + \varepsilon \mid X_1 = x_1)$$
$$= \underline{\beta_0 + \beta_1 x_1}$$
$$\underbrace{\qquad\qquad\qquad}_{\text{true regression line}}$$

# Lec 8 (Sep 29)



an approximation of

true regression line
$$\mathbb{E}(Y \mid X_1 = x_1) = \beta_0 + \beta_1 x_1$$

fitted regression line
$$\hat{m}(x_i) = \hat{\beta}_0 + \hat{\beta}_1 x_i$$

d$y

d$x

# Least Squares estimators

Given $n$ observations $(x_{11}, y_1), (x_{21}, y_2), \ldots, (x_{n1}, y_n)$ to estimate $\beta_0$ and $\beta_1$ in SLR.



$E(Y \mid X_1 = x_1) = \beta_0 + \beta_1 x_1$

Goal: minimize the residual sum of squares

Residual sum of squares: $\frac{1}{n} \sum_{i=1}^{n} e_i^2 = \frac{1}{n} \sum_{i=1}^{n} (y_i - \hat{\beta}_0 - \hat{\beta} x_{i1})^2$

Least squares solve $\hat{\beta}_0, \hat{\beta}_1$ such that $(\hat{\beta}_0, \hat{\beta}_1) = \underset{\beta_0, \beta_1}{\text{argmin}} \; \frac{1}{n} \sum_{i=1}^{n} e_i^2$

$= \underset{\beta_0, \beta_1}{\text{argmin}} \; \frac{1}{n} \sum_{i=1}^{n} (y_i - \hat{\beta}_0 - \hat{\beta}_1 x_{i1})^2$

We can show that:

$$\hat{\beta}_1 = \frac{S_{xy}}{S_{xx}} = \frac{\sum_{i=1}^{n} (y_i - \bar{y})(x_{i1} - \bar{x})}{\sum_{i=1}^{n} (x_{i1} - \bar{x})^2}$$

$$\hat{\beta}_0 = \bar{y} - \hat{\beta}_1 \bar{x}_1$$

$\left. \vphantom{\begin{array}{c} a \\ b \end{array}} \right\}$ the least square estimators are the same as the plug in estimators

Ideally, $(\beta_0^*, \beta_1^*) = \underset{\beta_0, \beta_1}{\text{argmin}} \; \underbrace{E_{X,Y} \left[ (Y - \beta_0 - \beta_1 X_1)^2 \right]}_{\text{"expected risk"}}$

the least square estimators are a sample estimation of the above:

$$(\hat{\beta}_0, \hat{\beta}_1) = \underset{\beta_0, \beta_1}{\text{argmin}} \; \underbrace{\frac{1}{n} \sum_{i=1}^{n} (y_i - \beta_0 - \beta_1 x_{i1})^2}_{\text{"empirical risk"}}$$

$\rightarrow$ $\underline{\text{Proof}}$: $S(\beta_0, \beta_1) = \widehat{MSE}(\beta_0, \beta_1)$

$= \frac{1}{n} \sum_{i=1}^{n} (y_i - \beta_0 - \beta_1 x_{i1})^2$

$\frac{\partial S}{\partial \beta_0} = \frac{1}{n}(-2) \sum_{i=1}^{n} (y_i - \hat{\beta}_0 - \hat{\beta}_1 x_{i1}) := 0$

$\frac{\partial S}{\partial \beta_1} = \frac{1}{n}(-2) \sum_{i=1}^{n} (y_i - \hat{\beta}_0 - \hat{\beta}_1 x_{i1}) x_{i1} := 0$

$\Rightarrow \begin{cases} \bar{y} - \hat{\beta}_0 - \hat{\beta}_1 \bar{x}_1 = 0 \Rightarrow \hat{\beta}_0 = \bar{y} - \hat{\beta}_1 \bar{x}_1 \\ \frac{1}{n} \sum_{i=1}^{n} x_{i1} y_i - \hat{\beta}_0 \bar{x}_1 - \hat{\beta}_1 \frac{1}{n} \sum_{i=1}^{n} x_{i1}^2 = 0 \end{cases}$

Substitute $\hat{\beta}_0 = \bar{y} - \hat{\beta}_1 \bar{x}_1$ to obtain $\hat{\beta}_1 = \frac{S_{xy}}{S_{xx}}$

# SLR in matrix form

$$Y = \begin{bmatrix} Y_1 \\ Y_2 \\ \vdots \\ Y_n \end{bmatrix} \quad , \quad X = \begin{bmatrix} 1 & X_{11} \\ 1 & X_{21} \\ \vdots & \vdots \\ 1 & X_{n1} \end{bmatrix} \quad , \quad \beta = \begin{bmatrix} \beta_0 \\ \beta_1 \end{bmatrix} \quad , \quad \varepsilon = \begin{bmatrix} \varepsilon_1 \\ \varepsilon_2 \\ \vdots \\ \varepsilon_n \end{bmatrix}$$

↑ not observed yet, all are R.V.s

Assumptions (recall SLR A1-A4):

A1. same — arbitrary input

A2. $Y = X\beta + \varepsilon$

A3. $E(\varepsilon) = 0$ , $Var(\varepsilon) = \sigma^2 I_n$

A4. implied by A3 — independence of error ($\varepsilon \perp X$)

$$\begin{bmatrix} \sigma^2 & & 0 \\ & \ddots & \\ 0 & & \sigma^2 \end{bmatrix}$$

variance-covariance matrix:

$$\begin{bmatrix} Cov(\varepsilon_1, \varepsilon_1) & Cov(\varepsilon_1, \varepsilon_2) & \cdots \\ Cov(\varepsilon_2, \varepsilon_1) & Cov(\varepsilon_2, \varepsilon_2) & \cdots \\ \vdots & & \ddots \\ \vdots & & \end{bmatrix}$$

· $Cov(\varepsilon_i, \varepsilon_j) = 0 \quad \forall i \neq j$

  as $\varepsilon_i \perp \varepsilon_j$

· $Cov(\varepsilon_i, \varepsilon_i) = Var(\varepsilon_i) = \sigma^2$

Some notes:

1. $E(Y \mid X) = X\beta + E(\varepsilon)$

   $\qquad = X\beta$

   In regular form, $E(Y_1 \mid X_1 = x_1) = \beta_0 + \beta_1 x_1 \quad$ for $\quad i = 1, \ldots, n$

2. $Y$ has constant variance: $Var(Y \mid X) = \sigma^2 I_n = Var(\varepsilon)$

   $$Var(Y \mid X) = \begin{bmatrix} Var(Y_1 \mid X) & Cov(Y_1, Y_2 \mid X) & \cdots \\ Cov(Y_2, Y_1 \mid X) & Var(Y_2 \mid X) & \\ \vdots & & \ddots \end{bmatrix}$$

   $\Rightarrow Var(Y_i \mid X_{i1} = x_{i1}) = \sigma^2 \qquad$ for $\quad i = 1, \ldots, n$

   $\qquad Cov(Y_j, Y_k \mid X_{j1} = x_{j1}, X_{k1} = x_{k1}) = 0 \qquad$ for $\quad j \neq k$

# Least squares in matrix form

$$\hat{\beta} = \underset{\beta}{argmin} \; \frac{1}{n} \| Y - X\beta \|_2^2 = \frac{1}{n} \sum_{i=1}^{n} (y_i - \beta_0 - \beta_1 x_{i1})^2$$

→ solution is $\hat{\beta} = (X'X)^{-1} X'Y \quad$ (proof in handout)

# Lec 9 (Oct 1)

$$(\beta_0^*, \beta_1^*) = \underset{\beta_0, \beta_1}{\text{argmin}} \; \underbrace{E_{X,Y}\left[(Y - \beta_0 - \beta_1 X_1)^2\right]}_{\text{"expected risk"} = MSE(\beta_0, \beta_1)}$$

$$(\hat{\beta}_0, \hat{\beta}_1) = \underset{\beta_0, \beta_1}{\text{argmin}} \; \underbrace{\frac{1}{n}\sum_{i=1}^{n}(y_i - \beta_0 - \beta_1 x_{i.})^2}_{\text{"empirical risk"} = \widehat{MSE}(\beta_0, \beta_1)}$$

Note: If the data $(x_{1}, y_1)$, $(x_{2}, y_2)$, ... $(x_{n}, y_n)$ are all independent for any fixed, the Law of Large Numbers tells us that:

as $n \to \infty$, $\widehat{MSE}(\beta_0, \beta_1) \to MSE(\beta_0, \beta_1)$

If the SLR assumptions are true, then:

as $n \to \infty$, $(\hat{\beta}_0, \hat{\beta}_1) \to (\beta_0^*, \beta_1^*)$

(i.e. approximation is good when n large)

For SLR, $Y = \beta_0 + \beta_1 X + \varepsilon$, optimal (true) regression function is:

$$E(Y | X = x) = \beta_0 + \beta_1 X$$
$$= \beta_0^* + \beta_1^* X$$

From these results, we see that not only $\beta_0, \beta_1$ are the parameters of the SLR, but also $\beta_0, \beta_1$ can minimize the expected risk (MSE), i.e. $\beta_0^* = \beta_0$, $\beta_1^* = \beta_1$

These symptotic results depend on SLR assumptions being true.

# Statistical properties of least Squares

Estimates vs estimators?
- An <u>estimate</u> is a number (deterministic), the realized value of the estimator.
- An <u>estimator</u> is an r.v. (stochastic)

- For example, we can potentially draw n samples $X_1, X_2, ..., X_n$ from a distribution. The $X_i$'s have not actually been observed yet and can potentially be any value from the distribution.

$$\bar{X} = \frac{X_1 + X_2 + ... + X_n}{n}$$ is also an r.v., whose distribution is a sampling distribution.

↑ this is an estimator of the population mean $\mu$.

If the values are realized, then $\bar{x} = \frac{1}{n}(x_1 + x_2 + ... + x_n)$ is an estimate.

## Least Square estimators

Given that the SLR model assumptions A1 - A4 are satisfied, we can show that the least square estimators are unbiased:

$$E(\hat{\beta_1}) = \beta_1, \qquad E(\hat{\beta_0}) = \beta_0$$

and variance of the estimators are:

$$Var(\hat{\beta_1} \mid x_{11}, \ldots, x_{n1}) = \frac{\sigma^2}{S_{xx}}, \qquad var(\hat{\beta_0} \mid x_{11}, \ldots, x_{n1}) = \sigma^2 \left( \frac{1}{n} + \frac{\bar{x}_1^2}{S_{xx}} \right)$$

$\underbrace{\qquad}$ related to var of $\varepsilon$

proof in handout

} proof in handout

Bias and variance in matrix form:

$$\underset{\sim}{Y} = \underset{\sim}{X} \underset{\sim}{\beta} + \underset{\sim}{\varepsilon}$$

$$E(\underset{\sim}{Y} \mid \underset{\sim}{X}) = \underset{\sim}{X}\underset{\sim}{\beta}$$
$$Var(\underset{\sim}{Y} \mid \underset{\sim}{X}) = \sigma^2 \underset{\sim}{I_n}$$

unbiasedness : $E(\hat{\underset{\sim}{\beta}} \mid \underset{\sim}{X}) = \underset{\sim}{\beta}$

variance : $Var(\hat{\underset{\sim}{\beta}} \mid \underset{\sim}{X}) = \sigma^2 (\underset{\sim}{X}'\underset{\sim}{X})^{-1}$

$$\begin{bmatrix} cov(\hat{\beta_0}, \hat{\beta_0}) & cov(\hat{\beta_0}, \hat{\beta_1}) \\ cov(\hat{\beta_1}, \hat{\beta_0}) & cov(\hat{\beta_1}, \hat{\beta_1}) \end{bmatrix}$$

$$= \begin{bmatrix} var(\hat{\beta_0}) & cov(\hat{\beta_0}, \hat{\beta_1}) \\ cov(\hat{\beta_1}, \hat{\beta_0}) & var(\hat{\beta_1}) \end{bmatrix}$$

# Lec 10 ( Oct 6 )

## How to estimate $\sigma^2$ :

Recall $Y = \beta_0 + \beta_1 X_1 + \varepsilon$, $E(\varepsilon) = 0$, $Var(\varepsilon) = \sigma^2$

Although $\sigma^2$ is not used to estimate $\hat{\beta_0}$ and $\hat{\beta_1}$, it is still important as it tells us about :

1. randomness of $Y$

2. $\sigma^2$ is related to $var(\hat{\beta_0} \mid X)$ and $var(\hat{\beta_1} \mid X)$  $\qquad x_{11}, x_{21}, \ldots, x_{n1}$

$$Var(\hat{\beta_0} \mid X) = \sigma^2 \left( \frac{1}{n} + \frac{\bar{x}_1^2}{S_{xx}} \right)$$
$$Var(\hat{\beta_1} \mid X) = \frac{\sigma^2}{S_{xx}}$$

$$\hat{\sigma^2} = \frac{\text{sum of squared residuals}}{n-2} = \text{mean sq. residuals}$$

sum of sq. residuals: $SSR = \sum_{i=1}^{n} \hat{e}_i^2 = \sum_{i=1}^{n} (y_i - \hat{y}_i)^2$

$$= \sum_{i=1}^{n} (y_i - \hat{\beta_0} - \hat{\beta_1} x_{i1})^2$$

mean sq. res: $MSR = \frac{1}{n-2} \sum_{i=1}^{n} (y_i - \hat{\beta_0} - \hat{\beta_1} x_{i1})^2 = \frac{1}{n-2} SSR$

recall that $\widehat{MSE} = \frac{1}{n} SSR$.
When $n$ is large, $MSR \approx \widehat{MSE}$.

Recall that in SLR:
$$\mathbb{E}\left((Y - \beta_0 - \beta_1 X_1)^2\right) = \mathbb{E}(\varepsilon^2)$$
$$= \text{Var}(\varepsilon) + (\mathbb{E}(\varepsilon))^2$$
$$= \sigma^2$$

Plug-in principle: use the observed sample values
- Replace $X_1$ with $(x_{11}, x_{21}, \ldots, x_{n1})$
- Replace $Y$ with $(y_1, y_2, \ldots, y_n)$
- $\mathbb{E}(\cdot)$ replaced by $\frac{1}{n}\sum_{i=1}^{n}$

$$\Rightarrow \mathbb{E}((Y - \beta_0 - \beta_1 X_1)^2) \approx \frac{1}{n}\sum_{i=1}^{n}\left[(y_i - \hat{\beta}_0 - \hat{\beta}_1 x_{1i})^2\right] \quad \longleftarrow \text{which is fully computable}$$

$$\approx \frac{1}{n-2}\sum_{i=1}^{n}\left[(y_i - \hat{\beta}_0 - \hat{\beta}_1 x_{1i})^2\right]$$

$\nwarrow$ this adjustment makes it unbiased.
Does not really matter when $n$ is large,
but makes a diff when $n$ small.

$$\hat{\sigma} = \sqrt{\hat{\sigma}^2} = \text{residual standard error, standard error of regression}$$

$\hat{\sigma}^2$ is unbiased: $\mathbb{E}(\hat{\sigma}^2) = \mathbb{E}\left(\frac{SSR}{n-2}\right) = \sigma^2$    (proof in notes)
$\nwarrow$ an r.v.

If we use $n$ instead of $n-2$, i.e. $\tilde{\sigma}^2 = \frac{1}{n}(SSR)$, $\mathbb{E}(\tilde{\sigma}^2) = \frac{n-2}{n}\sigma^2 \xrightarrow{n\to\infty} \sigma^2$

Alternative formula for SSR: $SSR = \underbrace{\sum_{i=1}^{n} y_i^2 - n\bar{y}^2}_{TSS} - \hat{\beta}_1 S_{xy}$

total sum of squares: $TSS = \sum_{i=1}^{n}(y_i - \bar{y})^2 = \sum_{i=1}^{n} y_i^2 - n\bar{y}^2$

<u>Using $\hat{\sigma}^2$ to estimate $\text{Var}(\hat{\beta}_1 | X)$ and $\text{Var}(\hat{\beta}_0 | X)$</u>

$$\widehat{\text{Var}}(\hat{\beta}_1 | X) = \frac{\hat{\sigma}^2}{S_{xx}} \qquad\qquad \widehat{\text{Var}}(\hat{\beta}_0 | X) = \hat{\sigma}^2\left(\frac{1}{n} + \frac{\bar{x}^2}{S_{xx}}\right)$$

The standard errors (se) are the square roots of the variances.
$$se(\hat{\beta}_1) = \sqrt{\frac{\sigma^2}{S_{xx}}} \qquad\qquad se(\hat{\beta}_0) = \sqrt{\sigma^2\left(\frac{1}{n} + \frac{\bar{x}^2}{S_{xx}}\right)}$$

Replace $\sigma^2$ with $\hat{\sigma}^2$ for the estimates:
$$\widehat{se}(\hat{\beta}_1) = \sqrt{\frac{\hat{\sigma}^2}{S_{xx}}} \qquad\qquad \widehat{se}(\hat{\beta}_0) = \sqrt{\hat{\sigma}^2\left(\frac{1}{n} + \frac{\bar{x}^2}{S_{xx}}\right)}$$

Ⓡ **estimation of $\sigma^2$**
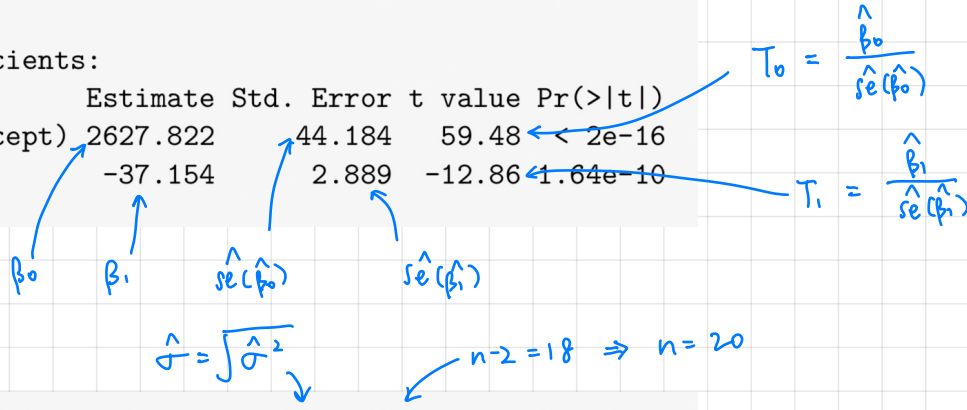
```
Call:
lm(formula = y ~ x)

Residuals:
     Min       1Q  Median      3Q     Max
 -215.98   -50.68   28.74   66.61  106.76

Coefficients:
              Estimate Std. Error t value Pr(>|t|)
(Intercept)  2627.822     44.184   59.48   < 2e-16
x             -37.154      2.889  -12.86  1.64e-10
```

$T_0 = \dfrac{\hat{\beta_0}}{\hat{se}(\hat{\beta_0})}$

$T_1 = \dfrac{\hat{\beta_1}}{\hat{se}(\hat{\beta_1})}$

$\beta_0$   $\beta_1$   $\hat{se}(\hat{\beta_0})$   $\hat{se}(\hat{\beta_1})$

$\hat{\sigma} = \sqrt{\hat{\sigma^2}}$    $n - 2 = 18 \implies n = 20$

```
Residual standard error: 96.11 on 18 degrees of freedom
Multiple R-squared:  0.9018,Adjusted R-squared:  0.8964
F-statistic: 165.4 on 1 and 18 DF,  p-value: 1.643e-10
```

<u>Sampling distributions of $\hat{\beta_0}$, $\hat{\beta_1}$ and $\hat{\sigma^2}$</u>

# Lec 11 (Oct 8)

+ 1 assumption : Gaussian- Noise Simple Linear Regression

We can show that :   $\hat{\beta_1} \sim N(\beta_1, \dfrac{\sigma^2}{S_{xx}})$

$\hat{\beta_0} \sim N(\beta_0, \sigma^2(\dfrac{1}{n} + \dfrac{\bar{x_i}^2}{S_{xx}}))$

$\dfrac{(n-2)\hat{\sigma^2}}{\sigma^2} \sim \chi^2_{n-p}$    $(p = 2)$   — 1 predictor  1 intercept

GN- SLR assumptions:
A1. _____
2. _____   ⎱ same as SLR assumptions.
3.  $\varepsilon \sim N(0, \sigma^2)$   ⎰
4. _____

Note the additional Gaussian assumption on the distribution of $\varepsilon$  (A3).

| | plug in | LS | prediction | unbiased estimators | $\hat{\beta} \sim N$ | t test of F test | CI | $R^2$ |
|---|---|---|---|---|---|---|---|---|
| SLR | ✓ | ✓ | ✓ | ✓ | | | | ✓ |
| GN-SLR | ✓ | ✓ | ✓ | ✓ | ✓ | ✓ | ✓ | ✓ |

The GN-SLR model is strictly stronger than SLR. This means that everything we have done so far directly applies to GN-SLR.

# Sampling distribution of $\frac{\hat{\beta_1} - \beta_1}{se(\hat{\beta_1})}$ and $\frac{\hat{\beta_1} - \beta_1}{\widehat{se}(\hat{\beta_1})}$

---

Results: $\quad \frac{\hat{\beta_1} - \beta_1}{se(\hat{\beta_1})} \sim N(0, 1) \qquad\qquad \frac{\hat{\beta_1} - \beta_1}{\widehat{se}(\hat{\beta_1})} \sim t_{n-p}$

$\nearrow$ # of variables

for $Y = \beta_0 + \beta_1 X_1$, $\quad p = 2$

$\Rightarrow \sim t_{n-2}$

Explanation:

Recall results from prev lecture: $\hat{\beta_1} \sim N\left(\beta_1, \frac{\sigma^2}{S_{xx}}\right)$

$$\hat{\beta_0} \sim N\left(\beta_0, \sigma^2 \left(\frac{1}{n} + \frac{\bar{x}^2}{S_{xx}}\right)\right)$$

For an r.v. $Z$ if $Z \sim N(\mu, \sigma^2)$, then $\frac{Z - \mu}{\sigma} \sim N(0, 1)$. We can similarly standardize $\hat{\beta_1}$:

$$\frac{\hat{\beta_1} - \beta_1}{se(\hat{\beta_1})} \sim N(0, 1) \qquad\qquad \frac{\hat{\beta_0} - \beta_0}{se(\hat{\beta_0})} \sim N(0, 1)$$

$\quad\quad\nwarrow$ this is not an observed value, it's unknown as it is related to the unknown $\sigma^2$

$\quad\quad$ Replace it with estimated se: $\widehat{se}(\hat{\beta_1})$

$$T_1 = \frac{\hat{\beta_1} - \beta_1}{\widehat{se}(\hat{\beta_1})} \sim t_{n-2} \qquad\qquad T_0 = \frac{\hat{\beta_0} - \beta_0}{\widehat{se}(\hat{\beta_0})} \sim t_{n-2}$$

where $\widehat{se}(\hat{\beta_1}) = \sqrt{\frac{\hat{\sigma}^2}{S_{xx}}} \qquad\qquad \widehat{se}(\hat{\beta_0}) = \sqrt{\hat{\sigma}^2 \left(\frac{1}{n} + \frac{\bar{x_1}^2}{S_{xx}}\right)}$

$\qquad\qquad\qquad = \sqrt{\frac{MSR}{S_{xx}}} \qquad\qquad\qquad\quad = \sqrt{MSR\left(\frac{1}{n} + \frac{\bar{x_1}^2}{S_{xx}}\right)}$

## CI for $\beta_1$, $\beta_0$

---

Suppose that $f_T$ is the density function of $T \sim t_{n-2}$.

Let $k > 0$ s.t. $P(-k < T < k) = 1 - \alpha$, $\quad \alpha \in (0, 1)$:



In the textbook,

$k = t_{\frac{\alpha}{2}, \; n-2}$

$\underbrace{\qquad}_{\text{degrees of freedom}}$

Result:

A $100(1 - \alpha)\%$ confidence interval for $\beta_1$ is: $\quad \leftarrow$ not $\hat{\beta_1}$ !!

$$CI(\beta_1) = \left[\hat{\beta_1} - k \cdot \widehat{se}(\hat{\beta_1}) \leq \beta_1 \leq \hat{\beta_1} + k \cdot \widehat{se}(\hat{\beta_1})\right]$$

$\rightarrow$ Proof: $P(\beta_1 \in CI(\beta_1)) = P\left(\hat{\beta_1} - k \cdot \widehat{se}(\hat{\beta_1}) \leq \beta_1 \leq \hat{\beta_1} + k \cdot \widehat{se}(\hat{\beta_1})\right)$

$$= P\left(-k \leq \frac{\hat{\beta_1} - \beta_1}{\widehat{se}(\hat{\beta_1})} \leq k\right)$$

$$= P(-k \leq T_1 \leq k)$$

$$= 1 - \alpha \qquad \text{(by def}^n\text{)}$$

The upper and lower bounds of $CI(\beta_1)$ are random as $\hat{\beta}_1$ can change if a different sample is drawn.
However, $\beta_1$ is fixed / non-random. ← unless you're a Bayesian

Explanation : The random interval $CI(\beta_1)$ traps $\beta_1$ with probability $1-\alpha$.
(i.e. If $\alpha = 0.05$, then 95% of the intervals of $CI(\beta_1)$, obtained by repeatedly sampling the data and finding $\hat{\beta}_1$, will contain $\beta_1$)

How does the width of $CI(\beta_1)$ change?   width $= 2k \cdot \hat{se}(\beta_1)$
___

- $\alpha \downarrow \Rightarrow (1-\alpha) \uparrow \Rightarrow$ width $\uparrow$
  High confidence comes at a price of big margin of error.
- $n \uparrow \Rightarrow$ width $\downarrow$
  Larger sample gives more accurate estimation.
- $\sigma^2 \uparrow \Rightarrow \hat{se}(\hat{\beta}_1) \uparrow \Rightarrow$ width $\uparrow$
  The more noise here is around the true regression line, the less precisely we can measure this line from the data.
- $S_{xx}$ (the variation of $X$) $\uparrow \Rightarrow$ width $\downarrow$

Ⓡ  simulation for constructing CI

# Hypothesis Testing          Lec 12 · Oct 13

$t$ Test (Wald Test) for $\beta_1$ using sampling distribution of $\beta_1$:

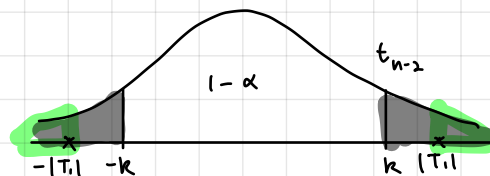Suppose we want to test (2 sided test):
$\quad H_0 : \beta_1 = c$
$\quad H_1 : \beta_1 \neq c$
where $c$ is a specific value e.g $c = 0$

method 1 : Suppose we want to test with a significance level $\alpha$, we compute the
test statistics:   $T_1 = \dfrac{\hat{\beta}_1 - c}{\hat{se}(\hat{\beta}_1)}$  ← all these values are known / can be computed

Reject $H_0$ : $|T_1| \geq k \equiv t_{\frac{\alpha}{2}, n-2}$
$k(n, a)$ is $(1 - \frac{\alpha}{2}) \, 100\%$ percentile of
the distribution



or equivalently, consider the tail probability:
$P(|T| > |T_1|) = \underbrace{P(T < -|T_1| \text{ or } T > |T_1|)}_{p-value} < \alpha$

Why it works:

- We assume GN-SLR assumptions are satisfied.

If we assume that $H_0: \beta_1 = c$ is true, then we know the sampling dist for
$T_1 = \dfrac{\hat{\beta}_1 - \beta_1}{se(\hat{\beta}_1)}$ will be equal to $T_1 = \dfrac{\hat{\beta}_1 - c}{se(\hat{\beta}_1)}$ and $T_1 \sim t_{n-2}$.

If we observe the unlikely result $|T_1| \geq k$, then there might be 2 explanations:
1. either in fact $\beta_1 = c$ ($H_0$ is true) and therefore you just by chance observe a very rare event that $|T_1| \geq k$ with the probability being only $\alpha$
2. or the assumption is wrong and $\beta_1 \neq c$ ($H_0$ is not true) and $H_0$ should be rejected.

Method 2: We will reject $H_0$ in a hypothesis test with significant level $\alpha$ if the $100(1-\alpha)\%$ confidence interval $CI(\beta_1)$ does not cover $c$, and do not reject $H_0$ if $CI(\beta_1)$ contains $c$.

Test of significance of regression

$$H_0: \beta_1 = 0 = c$$
$$H_1: \beta_1 \neq 0$$

If we fail to reject $H_0$, this implies that there is no linear relationship between X and Y.

Comments on hypothesis testing:
- what $\alpha$ should we use? conventionally $\alpha = 0.05$
  ~ false discovery rate
  If you are very conservative and cant afford false rejection, should use smaller $\alpha$.
- statistical significance:
  If we test $H_0: \beta_1 = 0$ and we reject it, then we say the difference between $\beta_1$ and $0$ is statistically significant with a given significance level $\alpha$.

     — can only test for $c = 0$
       in permutation test!

Permutation test      ($H_0: \beta_1 = 0$ ;  $H_1: \beta_1 \neq 0$)
Steps:
1. compute the observed value of the t statistic:  $T_1 = \dfrac{\hat{\beta}}{se(\hat{\beta}_1)}$
2. randomly permute the data $Y_1, \ldots, Y_n$ while keeping $X_1 \ldots X_n$ unchanged. Recompute the statistic $T_1$ again using permuted data.
3. repeat the previous step B times and let $Z_1 \ldots Z_B$
4. the approximated p-value is p-value $= \dfrac{1}{B} \sum_{j=1}^{n} \mathbb{I}(|Z_j| > |T_1|)$
   reject $H_0: \beta_1 = 0$ if p-value $< \alpha$
                                ↖ indicator function

- small permuted p-value is a strong indication of the rejection of the null hypothesis.

- permuted p-value also indicates false discovery rate

- idea of permutation test is trying to see if there is an association between Y and X. That's why we "scramble" the Ys.

# Analysis of variance · Oct 20

$SS_T$ (total sum of squares) : measures the total variation in Y.
$$SS_T = \sum_{i=1}^{n} (y_i - \bar{y})^2$$

$SS_R$ (regression sum of squares) : measures amount of "systematic variation" in Y due to the $Y \sim X$ linear relationship.
↳ i.e. the variation in Y that can be explained by the regression model.
$$SS_R = \sum_{i=1}^{n} (\hat{y}_i - \bar{y})^2$$
↑ pred. value of y

$SS_{Res}$ (residual sum of squares) : measures amount of "residual variation" in Y. Aggregate measure of misfit of the regression line
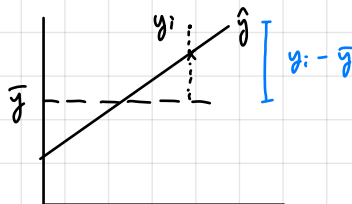$$SS_{Res} = \sum_{i=1}^{n} (y_i - \hat{y}_i) = \sum_{i=1}^{n} e_i^2$$

In general, $SS_T = SS_R + SS_{Res}$
$$\sum_{i=1}^{n} (y_i - \bar{y})^2 = \sum_{i=1}^{n} (\hat{y}_i - \bar{y})^2 + \sum_{i=1}^{n} (y_i - \hat{y}_i)^2$$

→ Proof in handout.

for observation i : $y_i - \bar{y} = \hat{y}_i - y_i + y_i - \bar{y}$

When we have a perfect fit, $\hat{y}_i = y_i$
$$SS_{Res} = 0 \Rightarrow SS_T = SS_R$$

# R² : goodness of fit · Oct 22

If we want to define a global measure of how well $X_1$ predicts Y, we consider the proportion of variation explained by the regression model:
$$R^2 = \frac{SS_R}{SS_T} = \frac{SS_R}{SS_R + SS_{Res}} = 1 - \frac{SS_{Res}}{SS_T}$$

We can show that $SS_R = \hat{\beta}_1 S_{xy}$ (see handout)

$R^2 \in [0, 1]$.

- $R^2 = 0$ when $\hat{\beta}_1 = 0$
  the model cannot explain any variation in $Y$, very bad fit
- $R^2 = 1$ when $SS_{RES} = 0$
- If $R^2$ is near 1 then the predictor $X_1$ can explain a large proportion of the observed variation in $Y$
  → i.e. the predictor $X_1$ included in the model is a "sufficient predictor" of $Y$.

<u>$R^2$ vs $t$-test on $\beta_1$?</u>     ($H_0: \beta_1 = 0$ ; $H_1: \beta_1 \neq 0$)

- $R^2$ tells us how much variation is explained by including $X_1$ (is $X_1$ a sufficient pred. of $Y$?)
- $t$-test tells us about whether $X_1$ is necessary in explaining the variation of $Y$.

We consider 4 scenarios:
1. insignificant p-value and low $R^2$
   - $X_1$ is not useful AND the model does not explain much of the variation (worst case)
   - ⇒ $\beta_1 = 0$

2. insignificant p-value and high $R^2$
   - $X_1$ not useful and the model explains a lot of variation within the data.
     i.e. model without $X_1$ is already sufficient
3. significant p-value and low $R^2$
   - $X_1$ is at least useful but not sufficient, should add more predictors.
4. significant p-value and high $R^2$
   - $X_1$ is useful and sufficient (best case)


<u>Adjusted $R^2$</u>
We have $E_{Y|X}(SS_{RES} | \mathbb{X}) = \sigma^2 (n-p)$        # of variables in model
                                                                         e.g in SLR, $p=2$ ($\beta_0, \beta_1$)
If we increases the model "complexity" to $n$,    i.e. $p \to n$, then
$$SS_{RES} \xrightarrow{p \to n} 0$$
$$\Rightarrow R^2 = 1 - \frac{SS_{RES}}{SS_T} \xrightarrow{p \to n} 1$$

Overfitting!
Alternatively,  $R^2_{adj} = 1 - \dfrac{SS_{RES} / (n-p)}{SS_T / (n-1)}$

· $\uparrow p \Rightarrow \uparrow$ numerator $\Rightarrow \downarrow R^2_{adj}$
$R^2_{adj}$ can account for model complexity.


Relationship between $R^2$ and $R^2_{adj}$ :  $R^2 = 1 - \dfrac{SS_{RES}}{SS_T} = 1 - \dfrac{n-p}{n-1}(1 - R^2_{adj})$
· intercept only model $(p=1)$: $R^2_{adj} = R^2$
·                SLR model $p=2$ :   $R^2_{adj} \neq R^2$
· when we fix $p$, $n \to \infty$,   $R^2_{adj} \xrightarrow{n \to \infty} R^2$
  If $n \gg p$, safe to use $R^2$. There is no overfitting issue when $n$ is large.

## Limitations of $R^2$

As $n \to \infty$, what does $R^2$ converge to?

$$R^2 = \frac{SS_R}{SS_T} = \frac{SS_R}{SS_R + SS_{Res}} = \frac{\hat{\beta}_1^2 \, S_{xx}}{\hat{\beta}_1^2 \, S_{xx} + SS_{Res}}$$

$$= \frac{\hat{\beta}_1^2 \, S_{xx} / n}{\hat{\beta}_1^2 \, S_{xx}/n + SS_{Res}/n}$$

As $n \to \infty$, $\hat{\beta}_1 \to \beta_1$, $\frac{S_{xx}}{n} \to Var(X_1)$, $\frac{SS_{Res}}{n} \to \sigma^2$,

Thus, $R^2 \xrightarrow{n \to \infty} \frac{\beta_1^2 \, var(X_1)}{\beta_1^2 \, var(X_1) + \sigma^2}$   if SLR assumptions hold

- By making $var(X_1)$ small or $\sigma^2$ large, we can drive $R^2$ towards 0 even if the model is correct.
- Conversely, even if the model is wrong, we can make $R^2$ close to 1 by increasing $var(X_1)$ or making $\sigma^2$ small

- $R^2$ can be compared only when different models are fit to the same dataset. We cannot compare $R^2$ across diff datasets!
- $R^2$ is not as useful as a goodness of fit measure on training data, but more useful for testing data.

## F test and ANOVA · Oct 27

The idea is to compare 2 models: (what does data look like?)

$H_0: Y = \beta_0 + \varepsilon$

$H_1: Y = \beta_0 + \beta_1 X_1 + \varepsilon$   (SLR)

equivalently,   $H_0: \beta_1 = 0$
(t test on $\beta_1$)   $H_1: \beta_1 \neq 0$

## Constructing the F test:

idea: If we fit the first model, the OLS estimator is $\hat{\beta}_0 = \bar{y}$.

The idea of the F test is to create a test statistic that measures how much better the 2nd model is compared to the 1st model.

Under GN-SLR assumption, and if $H_0$ is true (i.e. $\beta_1 = 0$), we can show that

$$\frac{SS_T}{\sigma^2} \sim \chi^2_{n-1} \qquad df_T = n-1$$

$$\frac{SS_{Res}}{\sigma^2} \sim \chi^2_{n-p} \qquad \text{In this case, } p=2, \quad \text{so } df_{Res} = n-2$$

↖ we refer to the model in the alt. hyp.
EVEN THOUGH we assume $H_0$ true!

$$\frac{SS_R}{\sigma^2} \sim \chi^2_{p-1}, \qquad df_R = p-1$$

For SLR, $p = 2$. We will generalize this later.

Degrees of freedom are additive:   $df_T = df_R + df_{Res}$

$$n-1 = p-1 + n-p$$
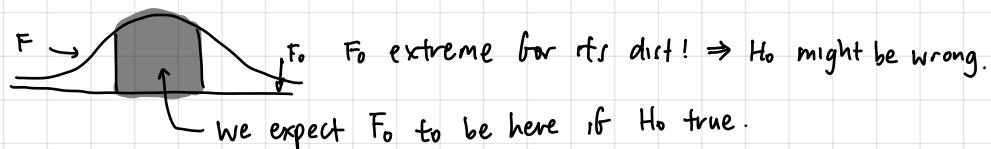
Consider the F-test statistic:

$$F_0 = \frac{SS_R / df_R}{SS_{Res} / df_{Res}} = \frac{SS_R / p-1}{SS_{Res} / n-p} \equiv \frac{MS_R}{MS_{Res}} \quad \overset{\text{mean sq.}}{}$$

Under $H_0 : \beta_1 = 0$, $F_0 \sim F_{(p-1, n-p)} = F_{(1, n-2)}$ for SLR

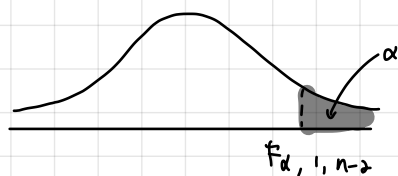Note that $F_0 \geq 0$. So we only need to do a 1-sided test.
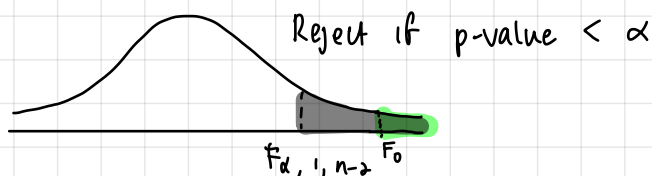We reject $H_0 : \beta_1 = 0$ if $F_0$ is large enough



$F \rightarrow$   $F_0$  $F_0$ extreme for its dist! $\Rightarrow H_0$ might be wrong.

— we expect $F_0$ to be here if $H_0$ true.

— observed statistic

Set $\alpha = 0.05$, $0.01$ etc. Reject $H_0$ if $F_0 > F_{\alpha, p-1, n-p}$.
$F_{\alpha, 1, n-2}$ represents the $1-\alpha$ quantile of the F distribution with df 1 and n-2.

$$\text{i.e.} \quad P[F > F_{\alpha, 1, n-2}] = \alpha$$



$F_{\alpha, 1, n-2}$

Equivalently, we can do F-test using p-value.   p-value $= P(F > F_0)$

Reject if p-value $< \alpha$



$F_{\alpha, 1, n-2}$  $F_0$

We summarize the F-test in an ANOVA table:

|  | SS | df | MS | F |
|---|---|---|---|---|
| Regression | $SS_R$ | $p-1$ | $MS_R = \frac{SS_R}{p-1}$ | $F_0 = \frac{MS_R}{MS_{Res}}$ |
| Residual | $SS_{Res}$ | $n-p$ | $MS_{Res} = \frac{SS_{Res}}{n-p}$ | |
| Total | $SS_T$ | $n-1$ | | |

(R) use anova() function

# Some notes on F-test

1. In deriving the F distribution, it is absolutely vital that all GN-SLR assumptions hold. The test never doubts that the right model is linear.

2. If we don't reject $H_0$, we don't find any significant share of variance associated with the regression.

$$F_0 = \frac{SS_R / df_R}{SS_{Res} / df_{Res}} \; \} \text{ larger}$$

The following may be interpreted:
   a) the intercept only model is better
   b) $\beta_1 \neq 0$ but the data doesn't provide enough power to detect departure from the null.
      ↳ this is a rather conservative expl$^n$.
   c) the real relationship between $X_1$ and $Y$ is nonlinear but the best approximation to it has slope zero.
      ↳ but F-test itself does not have power to detect non-linearity

3. If reject $H_0 : \beta_1 = 0$, may interpret the following:
   a) this does not mean that SLR is correct, only that the latter predicts better than the intercept-only model.
   b) SLR might be wrong, with every single one of assumptions violated, and yet better than the intercept-only model.

## F-test ≡ t-test for SLR    (proof in handout)

**Oct 29**

## ANOVA table in R

```
25 Analysis of Variance Table
26
27 Response: y
28           Df  Sum Sq Mean Sq F value    Pr(>F)
29 x          1 1527483 1527483  165.38 1.643e-10 ***
30 Residuals 18  166255    9236
31 ---
32 Signif. codes:  0 *** 0.001 ** 0.01 * 0.05 . 0.1    1
33
```

↗ $SS_R$

↑ same as from lm

$SS_{Res}$

verify $SS_T = SS_R + SS_{Res}$

```
> # To verify that SST = SSR + SSRes
> sum(anova(fit.RP)[,'Sum Sq'])-(n-1)*var(y)
```

$SS_{Res} + SS_R$ should be 0.

$SS_T = n-1 \cdot \frac{\sum_{i=1}^{n}(y_i - \bar{y})^2}{n-1}$

notice that these 2 values are the same! $F \equiv t$ test for SLR

not necessarily true for multiple lin reg.

```
13
14 Coefficients:
15             Estimate Std. Error t value Pr(>|t|)
16 (Intercept) 2627.822    44.184   59.48  < 2e-16 ***
17 x            -37.154     2.889  -12.86 1.64e-10 ***
18 ---
19 Signif. codes:  0 *** 0.001 ** 0.01 * 0.05 . 0.1    1
20
21 Residual standard error: 96.11 on 18 degrees of freedom
22 Multiple R-squared:  0.9018,   Adjusted R-squared:  0.8964
23 F-statistic: 165.4 on 1 and 18 DF,  p-value: 1.643e-10
```

$T_1^2 = F_0$      $p-1$      $n-p$

# Prediction Inference

At an arbitrary value $X_1 = x_{01}$ (not necessarily contained in the training data), we predict that **on average** Y will be:

$$\hat{y} = \hat{m}(x_{01}) = \hat{\beta}_0 + \hat{\beta}_1 x_{01}$$
$$\approx \mathbb{E}(Y | X_1 = x_{01})$$

we don't actually know what y will be! only know what it looks like on average.

The point prediction $\hat{y}$ is called the fitted value of the regression at $X_1 = x_{01}$.
Thus, $\hat{m}(x_{01})$ is an estimate of $\mathbb{E}(Y | X_1 = x_{01})$.
  ↳ random, function of data    ↑ deterministic qty
$\hat{m}(x_{01})$ inherits randomness from the estimators $\hat{\beta}_0$ and $\hat{\beta}_1$, which in turn inherit theirs from Y.

## Confidence interval for conditional mean

A $100(1-\alpha)\%$ CI for the conditional mean at the point $X_1 = x_{01}$ is:

$$CI(m(x_{01})) = \left[ \hat{m}(x_{01}) - k \cdot ese(\hat{m}(x_{01})) \; , \; \hat{m}(x_{01}) + k \cdot ese(\hat{m}(x_{01})) \right]$$

Recall that $k = t_{\frac{\alpha}{2}, n-2}$, it is the $\left[\frac{\alpha}{2} + (1-\alpha)\right]$ quantile of $t$
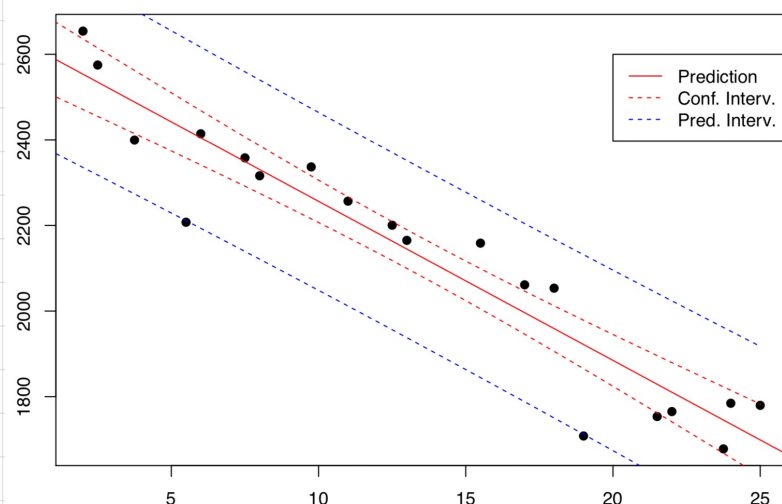→ In R:  $k = qt(1 - \frac{\alpha}{2}, df = n-2)$

How to calculate?
$$\hat{m}(x_{01}) \pm k \cdot ese(\hat{m}(x_{01})) = \hat{m}(x_{01}) \pm k \sqrt{\hat{\sigma}^2 \left(\frac{1}{n} + \frac{(x_{01} - \bar{x})^2}{S_{xx}}\right)}$$

→ proof in handout about CI theory
$\mathbb{E}(\hat{m}(x_{01})) = \beta_0 + \beta_1 x_{01}$   (unbiased)
$Var(\hat{m}(x_{01})) = \sigma^2 \left(\frac{1}{n} + \frac{(x_{01} - \bar{x}_1)^2}{S_{xx}}\right)$

The CI for cond. mean characterizes <u>uncertainty</u> in the prediction.

On the other hand, the prediction interval characterizes <u>uncertainty in the data</u>.

## Prediction interval    · **Nov 3**



CI is narrower than the pred. int.

Using GN-SLR, we know that at $X_1 = x_{01}$ the future observation of $Y_0$ is:

$$Y_0 = m(x_{01}) + \varepsilon = \beta_0 + \beta_1 x_{01} + \varepsilon$$

$x_{11}, x_{12}, \dots x_{1n}$ is the training data.
$x_{01}$ could be any one of $x_{1i}$, but could also be any arbitrary value.

Now we could construct a CI for $Y_0$ instead of $m(x_{01})$, we call it the pred. interval for the future observation of $Y_0$ corresponding to $X_1 = x_{01}$.

$$PI(Y_0) = \left[ \hat{Y}_0 - k \cdot ese(\hat{Y}_0) \ , \ \hat{Y}_0 + k \cdot ese(\hat{Y}_0) \right]$$

$$= \left[ \hat{m}(x_{01}) - k \cdot \sqrt{\hat{\sigma}^2 \left( 1 + \frac{1}{n} + \frac{(x_{01} - \bar{x}_1)^2}{S_{xx}} \right)} \ , \ \hat{m}(x_{01}) + k \cdot \sqrt{\hat{\sigma}^2 \left( 1 + \frac{1}{n} + \frac{(x_{01} - \bar{x}_1)^2}{S_{xx}} \right)} \right]$$

Note:
- A $100(1-\alpha)\%$ CI on conditional mean $m(x_{01})$ is an interval $[a, b]$ where
  $P(a \leq m(x_{01}) \leq b) = 1 - \alpha$
  ↑ true reg function: $\beta_0 + \beta_1 x_{01}$

- A $100(1-\alpha)\%$ PI for a future observation $Y_0 = m(x_{01}) + \varepsilon$ is an interval $[a', b']$:
  $P(a' \leq Y_0 \leq b') = 1 - \alpha$.

<u>Asymptotic behavior when $n \to \infty$</u>

For CI $(m(x_{01}))$:

$$\hat{m}(x_{01}) \pm k \cdot \sqrt{\hat{\sigma}^2 \left( \frac{1}{n} + \frac{(x_{01} - \bar{x}_1)^2}{S_{xx}} \right)} \xrightarrow{n \to \infty} m(x_{01}) = \beta_0 + \beta_1 x_{01}$$

- $\hat{m}(x_{01}) \xrightarrow{n \to \infty} m(x_{01})$ since $\hat{\beta}_0, \hat{\beta}_1 \to \beta_0, \beta_1$ as $n \to \infty$
- $\frac{(x_{01} - \bar{x}_1)^2}{S_{xx}} = \frac{(x_{01} - \bar{x}_1)^2}{\sum_{i=1}^{n}(x_i - \bar{x})^2} \xrightarrow{n \to \infty} 0$

Thus CI$(m(x_{01})) \xrightarrow{n \to \infty} m(x_{01})$, a single point.

For PI$(Y_0)$:

$$\hat{m}(x_{01}) \pm k \cdot \sqrt{\hat{\sigma}^2 \left( 1 + \frac{1}{n} + \frac{(x_{01} - \bar{x}_1)^2}{S_{xx}} \right)} \xrightarrow{n \to \infty} m(x_{01}) \pm k\sigma$$

Thus width of PI$(Y_0) = 2k\sigma$ as $n \to \infty$

$$PI(Y_0) \xrightarrow{n \to \infty} \left[ m(x_{01}) - k\sigma \ , \ m(x_{01}) + k\sigma \right]$$

## Correlation coefficient

$$\rho(X, Y) = \frac{Cov(X, Y)}{\sqrt{Var(X) \, Var(Y)}} = \frac{\sigma_{xy}}{\sigma_x \, \sigma_Y}$$

estimator: $\quad r = \dfrac{\sum_{i=1}^{n} (x_i - \bar{x})(y_i - \bar{y})}{\left[\sum_{i=1}^{n} (x_i - \bar{x})^2 \sum_{i=1}^{n} (y_i - \bar{y})^2\right]^{\frac{1}{2}}} = \dfrac{S_{xy}}{[S_{xx} \, SS_{\tau}]^{\frac{1}{2}}}$

Note that $\hat{\beta}_1 = \dfrac{S_{xy}}{S_{xx}} = \left(\dfrac{SS_{\tau}}{S_{xx}}\right)^{\frac{1}{2}} r$

Furthermore, $\boxed{r^2 = R^2}$ :
$\qquad\qquad\qquad\qquad\qquad \uparrow \sqrt{\dfrac{\text{spread of } Y}{\text{spread of } x}}$

$$\hat{\beta}_1^2 \left(\frac{S_{xx}}{SS_{\tau}}\right) = \frac{\hat{\beta}_1 \, S_{xy}}{SS_{\tau}} = \frac{SS_R}{SS_{\tau}} = R^2$$

## **Multiple Linear Regression** (LR with multiple predictors)

$$\mathbb{E}[Y|X] = \beta_0 + \beta_1 X_1 + \beta_2 X_2 + \cdots + \beta_k X_k = X\beta$$
$$\qquad\qquad\qquad\qquad\qquad\qquad \uparrow$$
$$[1 \quad X_1 \quad X_2 \quad \cdots \quad X_k] \qquad \begin{bmatrix} \beta_0 \\ \vdots \\ \beta_k \end{bmatrix}$$

Polynomial regression: $\mathbb{E}[Y|X_1] = \beta_0 + \beta_1 X_1 + \beta_2 X_1^2 + \cdots + \beta_k X_1^k$
$\qquad\qquad\qquad\qquad$ (single predictor $X_1$, just transformed)
$\to$ shape of the relationship between $X$ and $Y$.

Model with interactions: allow for joint effect of predictors.
$$\mathbb{E}[Y|X_1] = \beta_0 + \beta_1 X_1 + \beta_2 X_2 + \beta_{12} X_1 X_2 = X\beta$$
$$\qquad\qquad\qquad\qquad [1 \; X_1 \; X_2 \; X_1 X_2] \nearrow \quad \nwarrow \begin{bmatrix} \beta_0 \\ \beta_1 \\ \beta_2 \\ \beta_{12} \end{bmatrix}$$

- without interaction: $\mathbb{E}[Y|X] = \beta_0 + \beta_1 X_1 + \beta_2 X_2$
  with $\qquad\qquad$ '' : $\mathbb{E}(Y|X) = \beta_0 + \beta_1 X_1 + (\beta_2 + \beta_{12} X_1) X_2$

Model with transformation: $\mathbb{E}(Y|X) = \beta_0 + \beta_1 \log X_1 + \beta_2 X_2^2 + \beta_3 e^{X_3}$

In all cases we are using $\mathbb{E}(Y|X) = X\beta$
$$\qquad\qquad\qquad\qquad \nwarrow X \in \mathbb{R}^{p} \leftarrow \text{\# of pred.}$$

MLR assumptions:
1. There are $k$ predictors, $X_1, \cdots, X_k$. No assumptions about their distribution (might be random or non-random).
   We denote $X$ (no subscript) as $[X_1 \cdots X_k]$

2. Single response $Y$: $\quad Y = \beta_0 + \beta_1 X_1 + \cdots + \beta_k X_k + \varepsilon = X\beta + \varepsilon$

3. Noise $\varepsilon \perp X$: $\mathbb{E}(\varepsilon) = \mathbb{E}(\varepsilon|X) = 0$, $\quad Var(\varepsilon) = Var(\varepsilon|X) = \sigma^2$
   $\varepsilon_i$'s from different observations are independent $\quad (\varepsilon_i \perp \varepsilon_j)$
   $\Rightarrow Var(\varepsilon | X) = \sigma^2 I_n$

If given $n$ observations, $p = k+1$ parameters,

$$Y = \begin{bmatrix} Y_1 \\ \vdots \\ Y_n \end{bmatrix} \quad \mathbb{X} = \begin{bmatrix} 1 & X_{11} & \cdots & X_{1k} \\ \vdots & \vdots & \cdots & \vdots \\ 1 & X_{n1} & \cdots & X_{nk} \end{bmatrix} \quad \beta = \begin{bmatrix} \beta_0 \\ \vdots \\ \beta_k \end{bmatrix} \quad \varepsilon = \begin{bmatrix} \varepsilon_1 \\ \vdots \\ \varepsilon_n \end{bmatrix}$$

$$(n \times 1) \qquad\qquad (n \times p) \qquad\qquad (p \times 1) \qquad (n \times 1)$$

$Y = \mathbb{X}\beta + \varepsilon$  where:

$\mathbb{E}(\varepsilon \mid \mathbb{X}) = 0_n$

Variance-Cov matrix:

$$\begin{bmatrix} Var(\varepsilon_1) & Cov(\varepsilon_1, \varepsilon_2 \mid \mathbb{X}) & \cdots & Cov(\varepsilon_1, \varepsilon_n \mid \mathbb{X}) \\ Cov(\varepsilon_2, \varepsilon_1 \mid \mathbb{X}) & Var(\varepsilon_2) & & \\ \vdots & & \ddots & \\ Cov(\varepsilon_n, \varepsilon_1 \mid \mathbb{X}) & & & Var(\varepsilon_n) \end{bmatrix} = \begin{bmatrix} \sigma^2 & & 0 \\ & \ddots & \\ 0 & & \sigma^2 \end{bmatrix}$$

$$(n \times n)$$

$$\Rightarrow \mathbb{E}(Y \mid \mathbb{X}) = \mathbb{X}\beta \underset{\nwarrow n \times 1}{}$$

## Parameter Interpretation

$\sigma^2$: variance of noise around the true reg function (hyperplane).
= variance of $Y$ that cannot be explained by $X$.

$\beta_0$: expected value of $Y$ when $X_1 \cdots X_k$ are all $0$
$$\mathbb{E}[Y \mid X_1 = 0, \cdots, X_k = 0] = \beta_0$$

$\beta_j$ for $j = 1, \ldots, k$ :

$$\mathbb{E}(Y \mid X_j = x_j + 1, X_{-j} = x_{-j}) - \mathbb{E}(Y \mid X_j = x_j, X_{-j} = x_{-j})$$
$$= [\beta_0 + \beta_j(x_j + 1) + \cdots] - [\beta_0 + \beta_j x_j + \cdots]$$
$$= \beta_j$$

Interpretation: If we select 2 sets of cases from the distribution of the data where $X_j$ differs by 1, we expect $Y$ to differ by $\beta_j$ on average when all remaining predictors $X_{-j}$ are held constant.

## Least Squares for MLR

Observed data: $(y_1, X_1), (y_2, X_2), \ldots, (y_n, X_n)$

$$X_1 = [1 \ x_{11} \ x_{12} \ \cdots \ x_{1k}]$$
$$\vdots$$
$$X_n = [1 \ x_{n1} \ x_{n2} \ \cdots \ x_{nk}]$$

Goal: estimate $\beta = (\beta_0 \ \beta_1 \ \cdots \ \beta_k)^T$ in $\mathbb{E}(Y \mid \mathbb{X}) = \mathbb{X}\beta$
$\rightarrow$ minimize sums of squared residuals:

$$S(\beta) = \sum_{i=1}^{n}(y_i - X_i\beta)^2 = (Y - \mathbb{X}\beta)^T(Y - \mathbb{X}\beta) = \underbrace{\|Y - \mathbb{X}\beta\|_2^2}_{n \times 1}$$

$$\hat{\beta} = \underset{\beta}{argmin} \|Y - \mathbb{X}\beta\|_2^2, \quad \text{solve by taking derivative wrt } \beta = (\beta_0, \beta_1, \ldots \beta_k)$$

$$\frac{\partial S(\beta)}{\partial \beta_j} = -2\sum_{i=1}^{n} x_{ij}(y_i - X_i\beta) := 0 \qquad \text{for } j = 1, \ldots, k$$

In matrix form: $\nabla_\beta S(\beta) = -2\mathbb{X}^T(Y - \mathbb{X}\beta) = 0_p$

$\Leftrightarrow$ solving the equation $X^T Y - X^T X \beta = 0$

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

but is this invertible?

Some conditions for $(X^T X)$ invertible: $X \in \mathbb{R}^{n \times p}$.
- There must be more data than the number of parameters, $n \geqslant p$.
- In $X$, columns are not linearly dependent

The fitted value of the reg model at $x_i = (1 \ x_{i1} \ x_{i2} \ \dots \ x_{ik})$, $i = 1 \dots n$
$$\hat{y}_i = \beta_0 + \beta_1 x_{i1} + \beta_2 x_{i2} + \dots + \beta_k x_{ik}$$

For $X$: $\hat{Y} = X\hat{\beta} = X (X^T X)^{-1} X^T Y$

Statistical properties of $\hat{\beta}$
- $E(\hat{\beta}) = E(\hat{\beta} | X) = \beta$ unbiased
- $Var(\hat{\beta}) = \sigma^2 (X^T X)^{-1} = \sigma^2 C$
  $(p \times p)$

**NOV 12**

Estimator of $\sigma^2$
$$Var(\varepsilon) = E(\varepsilon^2) - (E(\varepsilon))^2 = E(\varepsilon^2) \approx \frac{1}{n} \sum_{i=1}^{n} \varepsilon_i^2$$
$$= \frac{SS_{Res}}{n-p}$$
$$= \frac{(Y - \hat{Y})^T (Y - \hat{Y})}{n - p}$$
$$= \frac{(Y - X\hat{\beta})^T (Y - X\hat{\beta})}{n - p}$$
$$\therefore \hat{\sigma}^2 = \frac{\| Y - X\hat{\beta} \|_2^2}{n - p} , \quad p = k+1 , \quad k = \# \text{ of predictors.}$$

Note that $n \approx n - p$ when $p$ is fixed and $n$ is large

$E(\hat{\sigma}^2) = \sigma^2$ unbiased

## t-test in MLR
Recall that for SLR, we required Gaussian-Noise assumption to construct t tests.

Therefore we further assume:
4. $\varepsilon \sim N(0_n, \sigma^2 I_n)$    $\varepsilon = \begin{bmatrix} \varepsilon_1 \\ \vdots \\ \varepsilon_n \end{bmatrix}$

5. $\varepsilon \perp X$. It follows that conditional on $X$, $Y$ has a multivariate Gaussian distribution: $Y | X \sim N(X\beta, \sigma^2 I_n)$

We can show that under GN-MLR, $\hat{\beta} \sim N(\beta, \underset{\underset{\text{unknown, use } \hat{\sigma}^2}{\uparrow}}{\sigma^2 (X^T X)^{-1}})$

H follows that $\hat{\beta}_j \sim N(\beta_j, \sigma^2 [(X^T X)^{-1}]_{jj})$.

    $j^{th}$ entry of $\hat{\beta}$ ⨏

$H_0 : \beta_j = 0 \qquad j = 0, 1, \cdots, R$
$H_1 : \beta_j \neq 0$

If we reject $H_0$, it indicates that the predictor $X_j$ is likely to be statistically
detectable / significant in the model.

$$T_j = \frac{\hat{\beta}_j - 0}{ese(\hat{\beta}_j)} \quad , \quad ese(\hat{\beta}_j) = \sqrt{\hat{\sigma}^2 \cdot C_{jj}} = \sqrt{\hat{\sigma}^2 [(X^T X)^{-1}]_{jj}}$$

$$\underset{= \frac{SS_{Res}}{n-p}}{\nwarrow}$$

Under GN-MLR and under $H_0$, $T_j \sim t_{n-p}$.



$P(T < -|T_j| \quad \text{or} \quad T > |T_j|) = \text{p-value}$

Reject $H_0$ if $p < \alpha$

We reject $H_0$ if $|T_j| > k = t_{\frac{\alpha}{2}, n-p} \longleftarrow (1 - \frac{\alpha}{2}) 100\%$ quantile in $t_{n-p}$ dist.

## F-test in MLR

Consider the regression model $Y = X\beta + \varepsilon$, $\beta = (\beta_0, \beta_1, \cdots, \beta_4)$

Test: $H_0 : \beta_2 = \beta_3 = \beta_4 = 0$
      $H_1 :$ at least 1 is not 0.

$\Leftrightarrow H_0 : \quad Y = \beta_0 + \beta_1 X_1 \qquad\qquad\qquad$ (reduced model)

$\qquad H_1 : \quad Y = \beta_0 + \beta_1 X_1 + \beta_2 X_2 + \cdots + \beta_4 X_4 \qquad$ (full model)

Generally, $Y = X\beta + \varepsilon = [X_{(1)} \ X_{(2)}] \begin{bmatrix} \beta_{(1)} \\ \beta_{(2)} \end{bmatrix} + \varepsilon$             **Nov 17**

                              ⤷ partition

$X_{(1)} : (n \times (p-r)) \qquad \beta_{(1)} : ((p-r) \times 1)$
$X_{(2)} : (n \times r) \qquad\quad \beta_{(2)} : (r \times 1)$

$\left. \begin{array}{l} H_0 : \beta_{(2)} = 0_r \\ H_1 : \beta_{(2)} \neq 0_r \end{array} \right\}$ equivalent to $\begin{array}{l} H_0 : Y = X_{(1)} \beta_{(1)} + \varepsilon \\ H_1 : Y = X_{(1)} \beta_{(1)} + X_{(2)} \beta_{(2)} + \varepsilon \end{array}$

The null hypothesis $\beta_{(2)} = 0$ can be tested by F-statistic:

$$F_0 = \underset{MS_{Res}}{\underbrace{\frac{\overline{SS}_R (\beta_{(2)} | \beta_{(1)}) / r}{SS_{Res}(\beta) / (n-p)}}} \quad , \quad \overline{SS}_R (\beta_{(2)} | \beta_{(1)}) = \overline{SS}_R (\beta) - \overline{SS}_R (\beta_{(1)})$$

$$\overline{SS}_R (\beta) = \hat{\beta}^T X^T Y \qquad \overline{SS}_R (\beta_{(1)}) = \hat{\beta}_{(1)}^T X_{(1)}^T Y$$

Recall that $\hat{\beta} = (X^T X)^{-1} X^T Y$, $\hat{\beta}_{(1)} = (X_{(1)}^T X_{(1)})^{-1} X_{(1)}^T Y$

$$MS_{Res} = \frac{SS_{Res}}{n-p} = \frac{Y^T Y - \hat{\beta}^T X^T Y}{n-p}$$

$\overline{SS}_R (\beta_{(2)} | \beta_{(1)})$ has degrees of freedom $p - (p-r) = r$
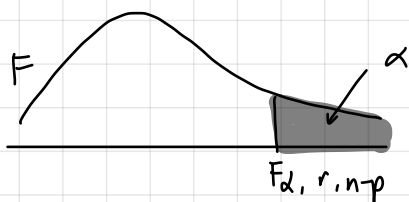It is the extra sums of squares due to $\beta_{(2)}$ when $\beta_{(1)}$ is included in the model.
i.e. extra contribution in $\overline{SS}_R$ due to predictors in $X_{(2)}$.

We can show that under GN-MLR assumptions and under $H_0: \beta_{(2)} = 0$, the random qty $F_0$
follows the F distribution with degrees of freedom $r$, $n-p$.

$$F_0 \sim F_{r, n-p}$$

We reject $H_0$ if: $F_0 > F_{\alpha, r, n-p}$
where $F_{\alpha, r, n-p}$ is $(1-\alpha) 100\%$ quantile of the F distribution with $r$, $n-p$ d.f.



$F$     $\alpha$

$F_{\alpha, r, n-p}$

equivalently,
p-value $= P(F > F_0)$
Reject $H_0$ when p-value $< \alpha$.

Conclusion: If we reject $H_0$, this means that at least 1 of the parameters in $\beta_{(2)}$ is not zero $\Rightarrow$ at least 1 of the predictors in $X_{(2)}$ is not zero.

Ⓑ **Multiple Linear Regression**

Multiple Linear Regression: Find plane of best fit

```
1 > fit.Del12<-lm(y~x1+x2,data=Delivery)
2 > summary(fit.Del12)
3 Coefficients:
4               Estimate Std. Error t value Pr(>|t|)
5 (Intercept)  2.341231   1.096730   2.135  0.044170 *
6 x1           1.615907   0.170735   9.464  3.25e-09 ***
7 x2           0.014385   0.003613   3.981  0.000631 ***
8 ---
```

$ese(\hat{\beta}_i)$

$H_0: \beta_2 = 0$
$H_1: \beta_2 \neq 0$

$\hat{\beta} = [\hat{\beta}_0, \hat{\beta}_1, \hat{\beta}_2]$

```
9 Signif. codes:  0 *** 0.001 ** 0.01 * 0.05 . 0.1   1
10
11 Residual standard error: 3.259 on 22 degrees of freedom
12 Multiple R-squared:  0.9596,   Adjusted R-squared:  0.9559
13 F-statistic: 261.2 on 2 and 22 DF,  p-value: 4.687e-16
```

F test: $H_0: y = \beta_0 + \varepsilon$ (intercept only) $\Big\}$ $H_0: \beta_{(2)} = 0$    $\beta_{(2)} = (\beta_1, \beta_2)$
        $H_1: y = \beta_0 + \beta_1 X_1 + \beta_2 X_2 + \varepsilon$    $H_1: \beta_{(2)} \neq 0$

An equivalent formula for F-statistic

Recall that F statistic:

$$F_0 = \frac{\overline{SS}_R(\beta_{(2)} \mid \beta_{(1)}) / r}{SS_{Res}(\beta) / (n-p)}$$

$$= \frac{(\overline{SS}_R(\beta) - \overline{SS}_R(\beta_{(1)})) / r}{SS_{Res}(\beta) / (n-p)}$$

Equivalently, $F_0 = \dfrac{[SS_{Res}(\beta_{(1)}) - SS_{Res}(\beta)] / r}{SS_{Res}(\beta) / (n-p)}$  $\longleftarrow$ this calculation is used in R

$\rightarrow$ **Proof:** Originally we have $SS_T = SS_{Res} + SS_R$.

$$\sum_{i=1}^{n} (y_i - \bar{y})^2 = \sum_{i=1}^{n} (y_i - \hat{y}_i)^2 + \sum_{i=1}^{n} (\hat{y}_i - \bar{y})^2$$

Now we consider an alternate form:

$$\overline{SS}_T = SS_{Res} + \overline{SS}_R$$

$$\sum_{i=1}^{n} y_i^2 = \sum_{i=1}^{n} (y_i - \hat{y}_i)^2 + \sum_{i=1}^{n} \hat{y}_i^2$$

$$(Y - \hat{Y})^T (X\hat{\beta})$$
$$= \underbrace{(Y - \hat{Y})^T X}_{\text{normal eqn, } = 0.} (X^T X)^{-1} X^T Y$$

Since $\sum_{i=1}^{n} y_i^2 = \sum_{i=1}^{n} (y_i - \hat{y}_i + \hat{y}_i)^2$

$$= \sum_{i=1}^{n} (y_i - \hat{y}_i)^2 + \sum_{i=1}^{n} \hat{y}_i^2 + \underbrace{2 \sum_{i=1}^{n} (y_i - \hat{y}_i)(\hat{y}_i)}_{= 0}$$

$$= \sum_{i=1}^{n} (y_i - \hat{y}_i)^2 + \sum_{i=1}^{n} \hat{y}_i^2$$

Specifically in MLR: $\overline{SS}_R = \sum_{i=1}^{n} \hat{y}_i^2 = (X\hat{\beta})^T (X\hat{\beta})$

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

$$= \hat{\beta}^T X^T Y$$

$$\overline{SS}_T = SS_{Res}(\beta) + \overline{SS}_R(\beta) = SS_{Res}(\beta_{(1)}) + \overline{SS}_R(\beta_{(1)})$$

$$\Rightarrow \overline{SS}_R(\beta) - \overline{SS}_R(\beta_{(1)}) = SS_{Res}(\beta_{(1)}) - SS_{Res}(\beta)$$

## CI in MLR

1. CI for $\beta$
2. CI for conditional mean $E(Y \mid X)$
3. Prediction interval for $Y^0$ (new observation)

see handout

# Model diagnostic

We hope to check the adequacy of our model. The residuals are:
$$e_i = y_i - \hat{y}_i$$
$$= y_i - x_i \hat{\beta} \quad , \quad i = 1, \dots, n$$

In matrix format: $e = (e_1 \; e_2 \; \cdots \; e_n)^T$

$$e = Y - \hat{Y}$$
$$= Y - X\hat{\beta}$$
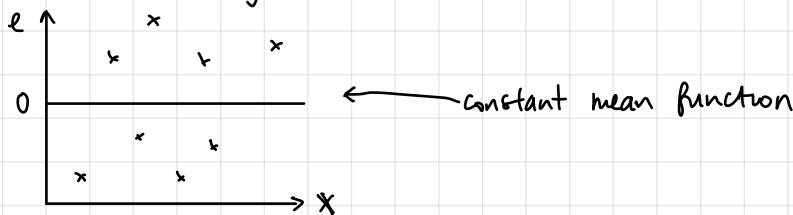$$= Y - \underbrace{X(X^TX)^{-1}X^T}_{H}Y = (I_n - H)Y$$

$H$ is symmetric: $H^T = H$
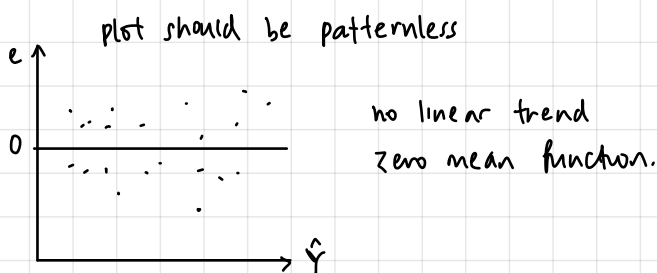idempotent: $HH = H$

**Nov 24**

1. $E(e|X) = O_n$

   since $E(e|X) = E(Y - \hat{Y}|X)$
   $$= E(Y|X) - E(X\hat{\beta}|X)$$
   $$= X\beta - X\beta \qquad \leftarrow \text{unbiased}$$
   $$= O_n$$

   We plot $e$ against any predictor $X_i$ or the linear combination of any predictors:


← constant mean function

2. We already know that $E(e|\hat{Y}) = 0$ 	Since $E(e|\hat{Y}) = E(e|X\hat{\beta})$
   $$= E(e|X) = 0$$

   Also, we have $Cov(e, \hat{Y}) = 0 \in R^{n \times n}$ 	Since $Cov(e, \hat{Y}) = Cov((I_n - H)Y, \; HY)$
   $$= (I_n - H)H \cdot Var(Y)$$
   $$= \sigma^2 (I_n - H)H$$
   $$= \sigma^2 (H - HH)$$
   $$= 0 \qquad \leftarrow H \text{ idempotent}$$

   plot should be patternless


   no linear trend
   zero mean function.

3. Variation of $e$: The variance of residual $e$ 	$Var(e|X) = \sigma^2(I_n - H)$

   since $Var(e|X) = Var((I_n - H)Y | X)$
   $$= (I_n - H)^T (I_n - H) Var(Y|X)$$
   $$= (I_n - H)\sigma^2 \quad \text{which is not a diagonal matrix}$$

   $$[Var(e|X)]_{ii} = (1 - h_{ii})\sigma^2 = Var(e_i|X)$$
   $$\qquad\qquad \uparrow \text{ the } i^{th} \text{ diagonal entry of } H$$

   $$[Var(e|X)]_{ij} = -\sigma^2 h_{ij} \neq 0 \quad \leftarrow$$
   $$= Cov(e_i, e_j | X) \qquad\quad e_i, e_j \text{ are correlated unlike noise terms}$$
   $$\qquad\qquad\qquad\qquad \varepsilon_1 \dots \varepsilon_n \text{ which are independent.}$$

In SLR case, $\text{Var}(e_i | X) = \sigma^2 (1 - h_{ii}) = \sigma^2 (1 - \frac{1}{n} - \frac{(x_i - \bar{x})^2}{S_{xx}}) \approx \underbrace{\sigma^2 (1 - \frac{1}{n})}_{\text{a constant}}$

small! $\uparrow$

$$\frac{(x_i - \bar{x})^2}{S_{xx}} = \frac{(x_i - \bar{x})^2}{(x_1 - \bar{x})^2 + \cdots + (x_n - \bar{x})^2}$$

Therefore the points near the center $\bar{x}$ have larger variance for $e_i$ than $e_i$ at more remote locations.

We should expect to see $\text{Var}(e_i | X) \approx \sigma^2 (1 - \frac{1}{n})$ which is almost constant

The difference in $\frac{(x_i - \bar{x})^2}{S_{xx}}$ gets less pronounced for a larger $n$.

4. For GN model ( SLR and MLR)
   If $Y | X \sim N$, then $e = (e_1, \ldots, e_n)^T$ is also Gaussian.
   $\qquad e_i \sim N(0, \sigma^2)$ for $i = 1, \ldots, n$
   $\qquad\qquad \nwarrow$ approximately true for large $n$.

   Quantile- Quantile (Q-Q) plot : make a histogram of $e$ and compare to Gaussian
   For a continuous dist , CDF: $F(z) = P(Z \leq z)$ has an inverse function:
   $\qquad F^{-1}(p) = z$ s.t. $p = P(Z \leq z)$
   For a Gaussian $Z \sim N(\mu, \sigma^2)$, quantile function: $F^{-1}(p) = \sigma \Phi^{-1}(p) + \mu$
   where $\Phi$ is the quantile function for $N(0,1)$

   If we plot $F^{-1}$ against $\Phi^{-1}$ we should get a straight line.
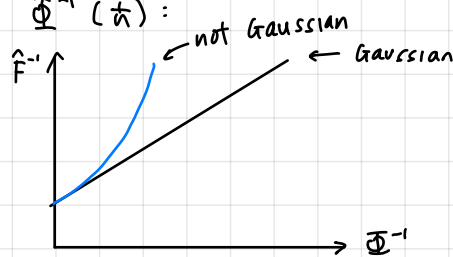   In practice, we replace $F^{-1}$ by the quantile function of $e$, $\hat{F}^{-1}$
   We estimate $\hat{F}^{-1}$ by using residuals $e_1 \ldots e_n$
   Arrange $e_1 \ldots e_n$ in increasing order : $e_{(1)}, e_{(2)}, \ldots, e_{(n)}$
   $\qquad e_{(i)}$ is greater than $\frac{i}{n}$ of the residuals.

   $$\hat{F}^{-1}(\frac{i}{n}) = e_{(i)}$$

   Now we can plot $\hat{F}^{-1}(\frac{i}{n}) = e_{(i)}$ against $\Phi^{-1}(\frac{i}{n})$ :

   $\qquad \hat{F}^{-1}(\frac{1}{n}) \longrightarrow \Phi^{-1}(\frac{1}{n})$
   $\qquad \hat{F}^{-1}(\frac{2}{n}) \longrightarrow \Phi^{-1}(\frac{2}{n})$
   $\qquad\qquad \vdots$
   $\qquad \hat{F}^{-1}(\frac{n}{n}) \longrightarrow \Phi^{-1}(\frac{n}{n})$

Test zero mean, constant variance $\rightarrow$ residual plot

Test Gaussian assumptions of noise $\rightarrow$ plot histogram of residuals, Q-Q plot

Test prediction performance of the model $\rightarrow$ generalization error

## Generalization error

$X = [1 \quad X_1 \quad \cdots \quad X_k]^T$

Response variable $Y$, predictors $X \in \mathbb{R}^P$.
We build a predictive model $\hat{f}(X)$ from the training data $T$.

- In MLR setting:
   $T = (y_1, x_1) \cdots (y_n, x_n) = (Y, X)$
   $\hat{f}(X) = X^T \hat{\beta} = X^T (X^T X)^{-1} X^T Y$

   To evaluate $\hat{f}(X)$, introduce a loss function:
   $$L(Y, \hat{f}(X)) = (Y - \hat{f}(X))^2 \quad \text{or} \quad |Y - \hat{f}(X)|$$
   $\uparrow$ target      $\uparrow$ prediction

   Generalization error:  $Y, X \sim$ test data
   $$\text{Error}_T = \mathbb{E}_{Y,X} (L(Y - \hat{f}(X)) \mid T)$$
   $$\approx \frac{1}{n'} \sum_{i=1}^{n'} L(y_i, \hat{f}(x_i))$$

   $(y_1, x_1) \cdots (y_{n'}, x_{n'})$ are from test data, $n'$ is the size of test data.

Generalization error vs training error:
- training error:
   training data $T = \{(y_i, x_i)\}_{i=1}^n$
   $\text{err} = \mathbb{E}_{Y,X} (L(Y, \hat{f}(X)) \approx \frac{1}{n} \sum_{i=1}^n L(y_i, \hat{f}(x_i))$
        $\uparrow$ from $T$
   In MLR, $\text{err} = \frac{1}{n} \sum_{i=1}^n (y_i - \hat{f}(x_i))^2 = \frac{1}{n} (Y - X\hat{\beta})^T (Y - X\hat{\beta})$

- In general, generalization error > training error.
   Assume data generated from the model $Y = X\beta + \varepsilon$. For simplicity, fix $X$.
   Test: $Y' = X\beta + \varepsilon'$
   Generalization error (test error): $\mathbb{E}_{Y'} [\frac{1}{n} \sum_{i=1}^n (Y_i' - \hat{Y}_i)^2]$
   Training error: $\mathbb{E}_Y [\frac{1}{n} \sum_{i=1}^n (Y_i - \hat{Y}_i)^2]$

   We can show that:
   $$\mathbb{E}_{Y'} [\frac{1}{n} \sum_{i=1}^n (Y_i' - \hat{Y}_i)^2] = \mathbb{E}_Y [\frac{1}{n} \sum_{i=1}^n (Y_i - \hat{Y}_i)^2] + \frac{2p}{n} \sigma^2$$
       (test error)          (training error)

   $p = k + 1$ (total # of variables)
   $\varepsilon, \varepsilon' \sim N(0, \sigma^2)$
   $\uparrow$ # of samples in data

   see handout for proof

Approximation of generalization error:
$$\frac{1}{n} \sum_{i=1}^{n} (Y_i' - \hat{Y}_i)^2 \approx \frac{1}{n} \sum_{i=1}^{n} (Y_i - \hat{Y}_i)^2 + \frac{2p}{n} \sigma^2$$
$$\approx \frac{1}{n} \sum_{i=1}^{n} (Y_i - \hat{Y}_i)^2 + \frac{2p}{n} \hat{\sigma}^2$$

Mallow's $C_p$ statistic:
$$C_p = \frac{1}{n} \sum_{i=1}^{n} (Y_i - \hat{Y}_i)^2 + \frac{2p}{n} \hat{\sigma}^2$$

e.g. 3 pred. models $\hat{f}_1, \hat{f}_2, \hat{f}_3$.
$$E(Y|X) = \hat{\beta}_0 + \hat{\beta}_1 X_1$$
$$E(Y|X) = \hat{\beta}_0$$
$$E(Y|X) = \hat{\beta}_0 + \hat{\beta}_1 X_1 + \hat{\beta}_2 X_2$$

Usually $\hat{\sigma}^2$ is obtained from the largest model i.e. $\hat{f}_3$

Whichever model gives the smallest $C_p$ value is the best in pred. performance

Adjusted $R^2$:
Recall that adjusted $R^2$ is $R_{adj}^2 = 1 - \dfrac{\frac{1}{n} \sum_{i=1}^{n} (Y_i - \hat{Y}_i) \cdot \frac{n}{n-p}}{\frac{1}{n} \sum_{i=1}^{n} (y_i - \bar{y})^2}$   $\overset{err}{\nearrow}$

$$\dfrac{}{\frac{1}{n} S_{yy}}$$

maximization of $R_{adj}^2 \iff$ minimize $err \cdot \frac{n}{n-p}$
$$err \cdot \frac{n}{n-p} = err \cdot \frac{1}{1 - \frac{p}{n}}$$

$$\approx err \left(1 + \frac{p}{n}\right) = err + \frac{p}{n} \cdot err$$
When $n \to \infty$, $err \to \sigma^2$
Thus $err \cdot \frac{n}{n-p} \to err + \frac{p}{n} \sigma^2$ as $n \to \infty$

AIC:
$$\overset{-\frac{err}{2\sigma^2}}{\swarrow}$$
$$AIC = -\frac{2}{n} \cdot \log \text{likelihood} + \frac{2p}{n} \quad (\text{general})$$
in MLR: $AIC = err + \frac{2p}{n} \hat{\sigma}^2 = C_p$

BIC:
$$BIC = err + \frac{\log n}{n} p \hat{\sigma}^2$$

<span style="color:green">**Transformation**</span>                    Dec 1

Transformations to linearize the model: If we detect non-linearity in the scatterplots or residual plots, in some case a nonlinear function can be linearized using a suitable transformation.

e.g.
$$E(Y|X) = \beta_0 + \beta_1 e^{-X}$$
$$E(Y|X) = \beta_0 + \beta_1 \left(\frac{1}{X}\right)$$
$$E(\ln Y|X) = \beta_0 + \beta_1 X_1$$
$$E(Y|X) = \beta_0 + \beta_1 \ln X_1$$
$$E(\ln Y|X) = \beta_0 + \beta_1 \ln X_1$$

e.g. Cobb-Douglas production function: $O_i = e^{\beta_0} L_i^{\beta_1} C_i^{\beta_2} u_i$

Output ↗     labor input ↑     capital input ↑     noise ↖

Take log: $\underbrace{\log O_i}_{Y_i} = \beta_0 + \beta_1 \underbrace{\log L_i}_{X_{i1}} + \beta_2 \underbrace{\log C_i}_{X_{i2}} + \underbrace{\log u_i}_{\varepsilon_i}$

Transformation requires $\varepsilon = \log u \sim N(0, \sigma^2)$ if we hope to assume GN-SLR.

## Polynomial terms

For any non linear function $f(x) = f(x_0) + \underbrace{\frac{f'(x_0)}{1!}(x - x_0)}_{(1^{st} \text{ order})} + \underbrace{\frac{f''(x_0)}{2!}(x - x_0)^2}_{(2^{nd} \text{ order})} + \cdots$

$\underbrace{f(x_0)}_{\text{constant}}$

we hope to approximate the relationship

$E(Y|X) = f(x)$
$\approx \beta_0 + \beta_1 X_1 + \beta_2 X_1^2 + \cdots + \beta_k X_1^k$

The more complicated the model (more terms), the more difficult it is to estimate with the same amount of data.

→ the more data you have, the more you can "afford" a more complicated model without overfitting

## Factor predictors (for categorical)

A factor predictor is a predictor that takes a discrete set of values on a nominal scale (i.e. non-numerical):

e.g. $X_1 = \{ Drug\ A, B, C \}$
$X_1 = \{ male, female \}$

Consider the case where $X_1$ takes $M$ values. We introduce a dummy variable:

$X_1^{(m)} = \begin{cases} 1 & \text{if } X_1 = m \\ 0 & \text{otherwise} \end{cases}$

A categorical variable with $M$ unique categories can be represented by $M-1$ dummies.

For example $X_1$ has $M$ categories. The corresponding model:

$E(Y|X) = \beta_0 + \sum_{m=1}^{M-1} \beta_m X_1^{(m)}$

→ if $(M)$: $X_1^{(1)} \cdots X_1^{(M-1)}$ are all 0.

$E(Y|X) = \beta_0$