

Prediction

with less "error"

notion of distance/norm (e.g. L_1 norm, L_2 norm etc.)

applied/using given data

(true) $Y = \beta_0 + \beta_1 X_i + \epsilon \rightarrow$ random variable
 ↓
 dependent variable of Y , β_0, β_1 and non-stochastic (i.e. non-random)

(prediction) $\hat{Y} = \hat{\beta}_0 + \hat{\beta}_1 X_{\text{new}}$

Estimate $\hat{\beta}_0, \hat{\beta}_1$ (to get $\hat{\beta}_0, \hat{\beta}_1$)

Error: (squared error) or L_2 norm (aka Euclidean norm)

$$Y = (y_1, y_2, \dots, y_n)^T, \text{ then } \|Y\|_2^2 = \sum_i y_i^2 \approx Y^T Y \quad (\text{vector notation})$$

$$\text{suppose, } \tilde{Z} = \begin{pmatrix} z_1 \\ z_2 \\ \vdots \\ z_n \end{pmatrix} \text{ and, } \tilde{Y} = \begin{pmatrix} y_1 \\ y_2 \\ \vdots \\ y_n \end{pmatrix} \text{ then, } \|Z - Y\|_2^2 = \sum_i (y_i - z_i)^2 = (Y - Z)^T (Y - Z)$$

We want to

minimize $S = \sum_{i=1}^n (y_i - (\beta_0 + \beta_1 x_i))^2$ & consequently the optimal values for β_0, β_1 are called the least square estimates ($\hat{\beta}_0, \hat{\beta}_1$)

$$\begin{aligned} \text{Solv: } \frac{\partial S}{\partial \beta_0} &= 0 \Rightarrow 2 \sum_i (y_i - \hat{\beta}_0 - \hat{\beta}_1 x_i) = 0 \Rightarrow \bar{Y} = \hat{\beta}_0 + \hat{\beta}_1 \bar{X} \\ \frac{\partial S}{\partial \beta_1} &= 0 \Rightarrow 2 \sum_i (y_i - \hat{\beta}_0 - \hat{\beta}_1 x_i) x_i = 0 \Rightarrow \sum_i x_i y_i = \hat{\beta}_0 \bar{x} + \hat{\beta}_1 \sum_i x_i^2 \end{aligned}$$

(aka Normal Eqn)

* Linear in parameters (y & variables)

(eq: $Y = \beta_0 + \beta_1 X \rightarrow$ linear model,

$$Y = \beta_0 + \beta_1 X^2 + \beta_2 Xz + \beta_3 z \log x \rightarrow \text{"linear" model}$$

$$\text{or, } Y = (1 \ x^2 \ xz \ z \ log x) \begin{pmatrix} \beta_0 \\ \beta_1 \\ \beta_2 \\ \beta_3 \end{pmatrix}$$

$$\text{or, } Y = X^T \beta \quad (\text{matrix representation})$$

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

But, no reliable measure of error, no testing of hypothesis, no check of model

error, no testing of hypothesis, no check of model
 no comparisons, etc. can be done
 unless we approach the problem
 statistically (via ad-hoc math, as done above)

• We need to develop the foundations rigorously to answer those questions. (which can't be done by merely 2 normal equations!) (In the above ad-hoc set up, ④, ⑤, ⑥ is all that we can do. Nothing more. Point predictions/estimat

* we don't know the distn of $\beta_0, \hat{\beta}_1, e^2$ & hence we can't know any interval estimates (& hence no hyp. testing)

Further, what if we have n -variables (x_1, x_2, \dots, x_n) ; hence a need to impose Linear Algebra structure.

Results from Linear Algebra:

Data: $(x_1, y_1), (x_2, y_2), \dots, (x_n, y_n)$; $y_i = \beta_0 + \beta_1 x_i + \epsilon_i \quad \forall i=1, 2, \dots, n \rightarrow$ Linear s/M

We define $\tilde{Y} = \begin{pmatrix} y_1 \\ y_2 \\ \vdots \\ y_n \end{pmatrix}, \tilde{X} = \begin{pmatrix} 1 & x_1 \\ 1 & x_2 \\ \vdots & \vdots \\ 1 & x_n \end{pmatrix}, \tilde{\beta} = \begin{pmatrix} \beta_0 \\ \beta_1 \end{pmatrix}, \tilde{\epsilon} = \begin{pmatrix} \epsilon_1 \\ \epsilon_2 \\ \vdots \\ \epsilon_n \end{pmatrix}$

* Assumptions:
 (a foresight)

① $\epsilon_i \sim N(0, \sigma^2)$

② $\sigma^2, \beta_0, \beta_1$ are unknown (but indep. of X, Y)

③ X 's are non-stochastic

$$\tilde{Y} = \tilde{X} \tilde{\beta} + \tilde{\epsilon} \quad (\text{Linear s/M in vector-matrix form})$$

where $\tilde{Y} \in \mathbb{R}^n, \tilde{X} \in \mathbb{R}^{n \times 2}, \tilde{\beta} \in \mathbb{R}^{2 \times 1}, \tilde{\epsilon} \in \mathbb{R}^n$
 important note: \tilde{X} is not square. Important to note that \tilde{X} has full rank. (M2.12)

(unconstrained optimization in \mathbb{R}^2 from M1)

Vector space: A vector space V over \mathbb{R} , denoted by $(V, +, \cdot, \mathbb{R})$, has the following properties

$$\forall \alpha, \beta \in \mathbb{R} \text{ and } \underline{x}, \underline{y}, \underline{z} \in V$$

1. $\underline{+}: V \times V \rightarrow V$

b) $(\underline{x} + \underline{y}) + \underline{z} = \underline{x} + (\underline{y} + \underline{z})$

c) $\exists \underline{0} \in V$ such that: $\underline{0} + \underline{x} = \underline{x} + \underline{0} = \underline{x} \quad \forall \underline{x} \in V$

d) $\exists -\underline{x} \quad \forall \underline{x} \in V$ such that: $-\underline{x} + \underline{x} = \underline{0}$

e) $\underline{x} + \underline{y} = \underline{y} + \underline{x} \quad \forall \underline{x}, \underline{y} \in V$

2. $\cdot: \mathbb{R} \times V \rightarrow V$

b) $\alpha \cdot (\beta \cdot \underline{x}) = (\alpha \cdot \beta) \cdot \underline{x}$

c) $1 \cdot \underline{x} = \underline{x}$

d) $(\alpha + \beta) \cdot \underline{x} = \alpha \cdot \underline{x} + \beta \cdot \underline{x}$

e) $\alpha \cdot (\underline{x} + \underline{y}) = \alpha \cdot \underline{x} + \alpha \cdot \underline{y}$

Subspace: Let $(V, \cdot, +, \mathbb{R})$ is a vector space, and $\boxed{S \subseteq V}$ such that $(S, +, \cdot, \mathbb{R})$ is also a vector space

then S is known as subspace of V

Eg: ① \mathbb{R}^2 : (a) origin $(0,0)$

(b) x axis, y axis, all lines passing through origin. } are subspaces of \mathbb{R}^2

② \mathbb{R}^3 : (a) origin $(0,0,0)$

(b) all lines passing thr. origin

(c) all planes passing thr. origin } are subspaces of \mathbb{R}^3

③ IP_n : (a) $S = \text{IP}_r, 1 \leq r \leq n$ are subspaces of IP_n

(i.e. $\text{IP}_0, \text{IP}_1, \text{IP}_2, \dots, \text{IP}_n$)

span: set (defined for a set of vectors $\{v_1, v_2, \dots, v_n\} = A$)

$\text{span}(A) = \text{set of all linear combinations of } v_1, v_2, \dots, v_n = \{x | \exists c_1, c_2, \dots, c_n \in \mathbb{R}, \text{ s.t. } c_1 v_1 + c_2 v_2 + \dots + c_n v_n = x\}$

basis: (defined for a vector space) A basis of a vector space (V) or subspace ($S \subseteq V$) is defined as

a set of linearly independent vectors s.t. they span V or S respectively.

* Need not be unique; but the no. of non-zero vectors in the basis is unique for a given V for a given V

\Rightarrow Dimension of V (i.e. $\dim(V) = \text{cardinality of basis of } V$)

Eg: $V = \{(x, y, 0) | x \in \mathbb{R}, y \in \mathbb{R}\}$

Now, consider $B = \{\underline{x}_1, \underline{x}_2\}$ and clearly \underline{x}_1 & \underline{x}_2 are lin-indep. & span V

$\therefore B$ is a basis for $V \therefore \dim(V) = 2$

(Note: This situation shall repeatedly occur in regression analysis) \leftarrow i.e. $V \subset \mathbb{R}^3, \dim(V) = 2$ (+ "components" in vector of V)

orthogonal vectors: vector $\underline{x}, \underline{y} \in V$ are said to be orthogonal if $\underline{x}^\top \underline{y} = \underline{y}^\top \underline{x} = 0$

* (Notation: $\underline{x} \perp \underline{y}$)

orthogonal basis: (special choice of basis) If the vectors of basis set B are orthogonal to each other then the basis is known as a

* we can obtain an orthogonalized basis from a non-orthogonal one (via, say, Gram-Schmidt orthogonalization)

(In most situations of RTSM, data set given will not be orthogonal, however we can transform them so as to become orthogonal)

Eg:

① \mathbb{R}^n { $(d)m^2$ is n }

② C^n { m^2 is n }

③ IP_n (polynomials of degree $\leq n$)

(but note, the dim is $n+1$)

(\because $n+1$ coefficients)

(example: $\underline{x} = [1, 2, 3]$)

orthogonal complement: If S is a subspace of V , then the orthogonal complement of S is defined as

$$S^\perp = \{ \underline{v} \mid \underline{v} \in V \text{ and } \underline{v}^T \underline{u} = 0 \forall \underline{u} \in S \}$$

* $S^\perp \cap S = \emptyset$ (\neq empty set)

* $\dim(S^\perp) + \dim(S) = \dim(V)$

(Note: we shall use this heavily in RTSM as subspace of error (E) and subspace of data set's prediction will be orthogonal complements & both shall together form entire V ; however, S_2 alone cannot span entire V).

(* These dimensions will give us the "degrees of freedom" in statistical analysis.)

$$(* \text{ corr}(X, Y) = \frac{E(XY) - E(X)E(Y)}{\sqrt{V(X) \cdot V(Y)}}$$

Projection Matrix (because geometrically as we shall see, regression is nothing but a projection)

If S is a subspace of a vector space V , then the "projection matrix" P_S for $S \subseteq V$ satisfies the following: (definition)

① If $\underline{x} \in S$, $P_S \underline{x} = \underline{x}$

② $\forall \underline{v} \in V$, $P_S \underline{v} \in S$

③ If P_S is a projection matrix for $S \subseteq V$, and $(I - P_S)$ is the projection matrix for $S^\perp \subseteq V$, then P_S & $(I - P_S)$ are orthogonal projection matrices of S and S^\perp .

Note: we shall split a vector into two orthogonal components or spaces via multiplying it once with P_S & other time via $P(I - P_S)$. $P_S \underline{v} \perp (I - P_S) \underline{v}$

$\left[(P_S \cdot \underline{v}) \cdot (I - P_S) \cdot \underline{v} = 0 \text{ (since } P_S^2 = P_S \text{ idempotent)} \right]$ (hence, $P_S \cdot (I - P_S) = P - P^2 = 0$) $\boxed{P_S^T = P_S}$

* In regression analysis, these 2 components will be prediction and error. We shall / our job will be to simply (!) find the projection matrix.

Thm: Projection matrix is an idempotent matrix. (Proof: $P_S \underline{v} \in S \quad \forall \underline{v} \in V$ (by ②))

Thm: The eigen values of projection matrix are 1 & 0 (Proof: $P_S \underline{v} = \lambda \underline{v}$ or, $P_S(P_S \underline{v}) = \lambda P_S \underline{v}$ or, $P_S^2 \underline{v} = \lambda P_S \underline{v}$ or, $\lambda^2 \underline{v} = \lambda \underline{v} \Rightarrow \lambda^2 = \lambda \Rightarrow \lambda = 0, 1$)

* Thm: If $\{ \underline{v}_1, \underline{v}_2, \dots, \underline{v}_k \}$ is an orthonormal basis of a subspace $S \subseteq V$ then, $P_S = \sum_{i=1}^k \underline{v}_i \underline{v}_i^T$ is an orthogonal projection matrix for $S \subseteq V$

We use this

Theorem to find projection matrix via orthogonal orthonormal basis.

(* To find orthonormal basis, we use Gram-Schmidt orthonormalization)

$$\underline{v} \times \underline{v}^T \underline{v} = (\underline{v} + \underline{y}, \underline{z} + \underline{g}) \underline{v} = (\underline{v}, \underline{v})$$

$$(\underline{v}) \underline{v} = (\underline{v} + \underline{y}) \underline{v}$$

column space: $\mathcal{C}(A) = \text{linear combinations of the columns of a matrix } A = [a_1 \ a_2 \ \dots \ a_n]$

$$\Leftrightarrow \mathcal{C}(A) = \text{span}(a_1, a_2, \dots, a_n)$$

$$= \left\{ \sum_{i=1}^n x_i a_i \mid \forall x_i \in \mathbb{R} \right\}$$

Similarly,

Rowspace

$$= \mathcal{R}(A) = \mathcal{C}(A^T)$$

Properties:

$$\begin{aligned} \textcircled{1} \quad \mathcal{C}(A+B) &= \mathcal{C}(A) + \mathcal{C}(B); \quad \mathcal{C}(AB) \subseteq \mathcal{C}(A) \\ \text{partitioned} &\quad \text{partitioned} \end{aligned}$$

$$\textcircled{2} \quad \dim(\mathcal{C}(A)) = \text{Rank}(A)$$

$$\textcircled{3} \quad \mathcal{C}(AAT) = \mathcal{C}(A) \Rightarrow \text{Rank}(AAT) = \text{Rank}(A)$$

(Proof: i) $\mathcal{C}(AAT) \subseteq \mathcal{C}(A)$ by $\textcircled{1}$

ii) $\mathcal{C}(A) \subseteq \mathcal{C}(AAT)$ → tedious proof.

$$\Leftrightarrow \mathcal{C}(AAT) = \mathcal{C}(A)$$

1) Positive definite (PD) matrix: A is said to be a PD matrix

(analogue of positive nos & negative nos) iff $\mathbf{x}^T A \mathbf{x} > 0 \quad \forall \mathbf{x} \neq 0 \quad (\Leftrightarrow |\mathbf{A}| > 0)$

Positive semi-definite matrix: A is said to be a PSD matrix

2) Generalized inverse: A^{-} has property: $AA^{-}A = A$

(i.e even for non-square matrices)

but NOT unique

3) For any matrix A , the projection matrix of $\mathcal{C}(A)$ is $(AAT)^{-1}$ and the orthogonal projection matrix is $A(A^TA)^{-1}A^T$

→ (but doesn't split into orthogonal components)

Multivariate Analysis // Prob. stats (well mostly for 1 or 2 variables)

$$\mathbf{x} = \begin{pmatrix} x_1 \\ x_2 \\ \vdots \\ x_n \end{pmatrix} = (x_1, x_2, \dots, x_n)^T \quad \mathbf{y} = (y_1, y_2, \dots, y_n)^T$$

$$\textcircled{1} \quad E(\mathbf{x}) = \begin{pmatrix} E(x_1) \\ E(x_2) \\ \vdots \\ E(x_n) \end{pmatrix} = \begin{pmatrix} \mu_1 \\ \mu_2 \\ \vdots \\ \mu_n \end{pmatrix} = \boldsymbol{\mu} \quad (\text{if } E(|x_i|) < \infty \quad \forall i=1, 2, \dots, n) \quad (\text{check of absolute convergence})$$

2) "Dispersion matrix" of \mathbf{x} (vector analogue of variance) denoted by $D(\mathbf{x}) = ((\text{cov}(x_i, x_j)))_{i,j}$

$$\text{variance-covariance matrix} = ((E(x_i x_j) - E(x_i)E(x_j)))_{i,j} = ((E(x_i x_j) - \mu_i \mu_j))_{i,j}$$

i) Note: Diagonal elements are the variances of the corresponding random variables.

$$= E(\mathbf{x} \mathbf{x}^T) - \boldsymbol{\mu} \boldsymbol{\mu}^T$$

ii) Note: Covariance b/w vectors of 2 different dimensions. (p, q)

$$\text{cov}(\mathbf{x}_p, \mathbf{y}_q) = ((\text{cov}(u_i, v_j)))_{i,j} \quad (\text{dimensions: } p \times q) \quad \text{(Generalized)}$$

$$= E(\mathbf{u}_p \mathbf{v}_q^T) - \boldsymbol{\mu}_p \boldsymbol{\mu}_q^T \quad \text{for non-square matrices too.}$$

$$\text{iii) } E(\mathbf{x} + \mathbf{b}) = E(\mathbf{x}) + \mathbf{b}$$

$$\underbrace{pxq}_{pxq}$$

where, \mathbf{b} is a constant/non-random vector.

$$\text{iv) } D(\mathbf{x} + \mathbf{b}) = D(\mathbf{x})$$

$$\text{v) } \text{cov}(\mathbf{x} + \mathbf{b}, \mathbf{y} + \mathbf{c}) = \text{cov}(\mathbf{x}, \mathbf{y})$$

some further results

(vi) $E(\underline{X}^T \underline{X}) = \underline{I}^T \underline{I} \leftarrow$ finally a scalar
(vii) $D(\underline{X}^T \underline{X}) = \underline{I}^T \Sigma \underline{I} \leftarrow$ finally a scalar.

(general result) ~~\underline{X}~~

(viii) $E(A\underline{X}) = A \underline{\mu}$
(ix) $D(A\underline{X}) = A \Sigma A^T$
(x) $\text{cov}(\underline{u}_p, \underline{v}_q) = \gamma$ (gamma symbol), then: $\text{cov}(A\underline{u}, B\underline{v}) = A \Gamma B^T$
 $\Rightarrow \text{cov}(A\underline{X}, B\underline{X}) = A \Sigma B^T$ (special case)

(# Nothing but linearity of expectation)
 $\underline{x} \in E(aX_1 + bX_2) = aE(X_1) + bE(X_2)$
~~which also implies~~ $\text{cov}(a\underline{X}, b\underline{X}) = a \Sigma b^T$ (special case)
 \underline{x} ($\Sigma_{\underline{X}}$)

* Thm: show that $D(\underline{X})$ is a p.s.d (positive definite matrix) $\underline{x} \in \boxed{\begin{matrix} \Sigma \geq 0 \\ \forall \underline{z} \neq 0 \end{matrix}}$

(hint: easily obtained via the fact: $V(X) \geq 0$)
Proof: We need to show that for every $\underline{z} \neq 0$; $\underline{z}^T \Sigma \underline{z} \geq 0$ scalar.
 $\Leftrightarrow \text{var}(\underline{z}^T \underline{X}) \geq 0$ which is true!

* (Very important)

Thm: Let $E(\underline{X}) = \underline{\mu}$, $D(\underline{X}) = \Sigma_{n \times n}$
then, $P((\underline{X} - \underline{\mu}) \in C(\Sigma)) = 1$

(In simple words, the random component of \underline{X} is $(\underline{X} - \underline{\mu})$ and this theorem says that we can always express that random component as a linear combination of column vectors (random vector) (in multivariate) of $\Sigma_{n \times n}$)

3) Theorem: Let \underline{X} be a random vector with n components such that $E(\underline{X}) = \underline{\mu}$ and, $D(\underline{X}) = \Sigma$ and $\text{rank}(\Sigma) = r \leq n$

If we assume that $\Sigma = B B^T$, where B is a $(n \times r)$ matrix and C is a left inverse of B ($B^T C = I_r$) then defining $\underline{Y} = C(\underline{X} - \underline{\mu})$ we get

- (i) $E(\underline{Y}) = \underline{0}$
- (ii) $D(\underline{Y}) = I_r$
- (iii) $\underline{X} = \underline{\mu} + B\underline{Y}$ with probability 1
(for full rank ($r=n$) $B = I_{n \times n}$)

* (Recall SVD: singular value decomp)
For $\Sigma_{n \times n} = \Sigma^T$ (symmetric)

$\Sigma = P D P^T$, where P is an orthogonal matrix

(we shall return here, when we do Principal component regression)

Analogy:
Let $E(Y) = 0$, $V(Y) = I_r$ in \mathbb{R}^r
 $X = \mu + \sigma Y \Rightarrow E(X) = \mu$, $V(X) = \sigma^2 I_n$

$E(\underline{Y}) = \underline{0}$, $D(\underline{Y}) = I_r$ in \mathbb{R}^r
 $X = \underline{\mu} + \Sigma^{1/2} \underline{Y} \Rightarrow E(X) = \underline{\mu}$, $V(X) = \Sigma$

? (what do you mean by " $\Sigma^{1/2}$ " bcoz Σ is a matrix)

$\lambda_1, \lambda_2, \dots, \lambda_n$; $\sum_i \lambda_i = \lambda_i \pi_i$

$\star \Sigma^{1/2} = P D^{1/2} P^T$

$\sqrt{\lambda_1}, \sqrt{\lambda_2}, \dots, \sqrt{\lambda_n}$

4) Multivariate Normal: A random vector \underline{X}_n is said to be following multivariate normal dist if it has a p.d.f:-

$f(\underline{X}) = \frac{1}{(2\pi)^n \sqrt{|D|}} \exp\left\{-\frac{1}{2} (\underline{X} - \underline{\mu})^T \Sigma^{-1} (\underline{X} - \underline{\mu})\right\}$

* $E(\underline{X}) = \underline{\mu}$
* $D(\underline{X}) = \Sigma_{n \times n}$

Q.) Derive Univariate Normal dist^r ($n=1$) ($X \in \mathbb{R}^1$)
 → Trivial: $f(x) = \frac{1}{\sqrt{2\pi}\sigma} e^{-\frac{1}{2\sigma^2}(x-\mu)^2}$

Q.) Derive Bivariate Normal dist^r ($n=2$) ($X \in \mathbb{R}^2$)

$$\rightarrow f(x_1, x_2) = \frac{\exp\left\{-\frac{1}{2}(x-\mu)^T \Sigma^{-1}(x-\mu)\right\}}{(2\pi)^2 \sqrt{|\Sigma|}}, \quad \Sigma = \begin{bmatrix} \sigma_1^2 & \rho\sigma_1\sigma_2 \\ \rho\sigma_1\sigma_2 & \sigma_2^2 \end{bmatrix}$$

∴ we get $|\Sigma| = \infty$
 (Substitute) and, $\Sigma^{-1} = \infty$

$$= \frac{\exp\left\{-\frac{1}{2(1-\rho^2)}\left[\left(\frac{x_1-\mu_1}{\sigma_1}\right)^2 + \left(\frac{x_2-\mu_2}{\sigma_2}\right)^2 - 2\rho\left(\frac{x_1-\mu_1}{\sigma_1}\right)\left(\frac{x_2-\mu_2}{\sigma_2}\right)\right]\right\}}{(2\pi)^2 \sigma_1 \sigma_2 \sqrt{1-\rho^2}}$$

Q.) Show that

Theorem: (For any distribution) Let \tilde{X} be a random vector in \mathbb{R}^n with $E(\tilde{X}) = \mu$ and
 Then $\rightarrow E(\tilde{X}^T A \tilde{X}) = \text{trace}(A) + \mu^T A \mu$

\downarrow scalar

\downarrow a.k.a "Quadratic Form"

$$(\because \tilde{X}^T A \tilde{X} = \sum_{i=1}^n \sum_{j=1}^n a_{ij} x_i x_j)$$

∴ $\mu, x_1, x_2, \dots, x_n$ iid $N(0, 1)$

Therefore, this is equivalent to

$$\tilde{X} \sim N(\mathbf{0}, I_n)$$

$$\# \text{ then, } E(\Sigma x_i^2) = E(\tilde{X}^T I_n \tilde{X}) = \text{tr}(I_n I_n) + 0 \quad (\because \mu = \mathbf{0})$$

If, $x_i \sim N(\mu_i, 1)$ and indep.

$$\text{then, } E(\Sigma x_i^2) = n + \mu^T \mu$$

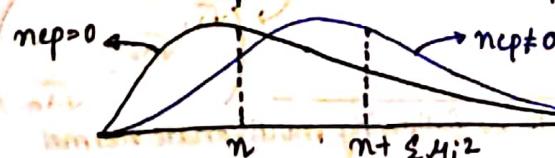
$$\text{or, } E(\Sigma x_i^2) = (n + \Sigma \mu_i^2)$$

If x_i iid $N(0, 1)$ then $\Sigma x_i^2 \sim \chi^2_n$ (already known) (a.k.a "central χ^2 ")

New: x_i iid $N(\mu_i, 1)$, then $\Sigma x_i^2 \sim \chi^2_{df=n}$ (a.k.a "non-central χ^2 ") ($\because \Sigma \mu_i^2 \neq 0$)

(Note: $n \leq n + \mu^T \mu$) (a.k.a "non-central χ^2 ") (more general)

$$\leq E(\chi^2_{df=n, ncp=0}) \leq E(\chi^2_{df=n, ncp=\mu^T \mu})$$



(Remark: whether $\mu_i > 0$ or < 0 , $\Sigma \mu_i^2 > 0$ & hence push the mean rightward & as we shall see later that under null hypo. $ncp=0$, whereas under alternate hyp. $\Sigma \mu_i^2 > 0$ & helps us to reject the null).

Even more general : If $X \sim N(\mu, I_n)$ then

$$\underline{x}^T A \underline{x} \sim \chi^2_{df = \text{rank}(A)}$$

iff A is an idempotent matrix.

ge-jeazzaan

(As we shall see later, A governs strongly the ncp & df, no matter what the M vector's rank & M is are).

Results:

① If $\tilde{\mathbf{x}} \sim N(\boldsymbol{\mu}, \Sigma)$

then $\tilde{X} \sim N(A\tilde{Y}, A\Sigma A^T)$

② If $\tilde{x} \sim N(\mu, \Sigma)$, and if $\exists B$ and its left inverse C , such that $\tilde{y} = C(\tilde{x} - \mu)$

then, $\tilde{x} = \tilde{A} + B\tilde{y}$ (ie \tilde{x} can be represented as a linear transformation of \tilde{y}) follows $N(\mu, \Sigma)$
 (can always be)

③ If A_1, A_2 symmetric and idempotent matrix such that $\alpha = (A_1 - A_2)$ is p.s.d, then $X^T A_1 X$ and $X^T A_2 X$ are independent.

$A_1 = \bar{X} + A_2$ - $181 \geq 0$ - for normally distributed X

④ If $A = AT$ and $CA = 0$ matrix, then $x^T A x$ and Cx are independent

(Special case: $X \sim N(0, I)$) \Rightarrow indep ; then $\frac{X}{\sqrt{Y/n}} \sim t_n$; $\text{mean} \leftarrow \text{definition of } t_n$

→ using ④, we can show that \bar{X} and s^2 are independent. (explanation is given in the next slide.)
 (hence $\frac{\bar{X}-\mu}{s/\sqrt{n}} \sim t_{n-1}$)

(Elaborate: x_i is $N(\mu, \sigma^2)$. Now, $\bar{x} = \frac{1}{n} \sum x_i = (\frac{1}{n} \frac{1}{n} \frac{1}{n} \dots \frac{1}{n}) \begin{pmatrix} x_1 \\ x_2 \\ \vdots \\ x_n \end{pmatrix}$

$$\text{and, } S^2 = \sum (x_i - \bar{x})^2$$

Also, $CA = \tilde{w}^T (I_n - \frac{1}{n} \tilde{w} \tilde{w}^T)$

$$\begin{aligned}
 &= \frac{1}{n} \left(\underbrace{\mathbf{1}^T \mathbf{I}_n}_{\text{all } 1's} - \underbrace{\frac{1}{n} \mathbf{1}^T \mathbf{1}}_{\text{constant}} \mathbf{1} \mathbf{1}^T \right) = \frac{1}{n} \left(\mathbf{1}^T - \frac{1}{n} \mathbf{1}^T \mathbf{1} \mathbf{1}^T \right) \\
 &= \frac{1}{n} \left(\mathbf{1}^T - \frac{1}{n} \mathbf{1}^T \mathbf{n} \right) = \mathbf{0}^T
 \end{aligned}$$

Also, $A = A^T$ (verify yourself) #
 Also, $A = A^2$ (idempotent) #

Note: ④ & ⑤ are enough to show that \bar{x} and s^2 are independent.

Note: $\#$ is needed because we want to show that g_2 will follow g_1 .

Note: (#) is needed because we want $\tilde{A}x$ follows normal distribution.

$$\text{Then } \mathbf{x}^T A^2 \mathbf{x} = \mathbf{x}^T A \mathbf{x}^T A \mathbf{x} = \mathbf{x}^T A \mathbf{x}$$

$$(\text{if } A^T = A) \quad (\text{if } A^2 = A)$$

→ (As $\tilde{X}^T A \tilde{X} = \tilde{\sigma}^T$ and $A^T = A$, then $(\tilde{X}^T \tilde{X} = \tilde{\sigma}^2)$ and, $X^T A X = S^2$) are independently distributed.

$\rightarrow (\bar{x} = \frac{1}{n}^T x \sim N(\frac{1}{n}^T \mu, \sigma^2 \frac{1}{n}^T I n \frac{1}{n})$, where $\mu = \frac{1}{n} \sum_{i=1}^n x_i$

$$\sigma_x \bar{x} \approx N(\mu, \sigma^2/n) - \mathbf{1}^T (\beta x - \mathbf{y})$$

$\rightarrow S^2 = \sum_i (x_i - \bar{x})^2 = x^T A x$; however, this is not χ^2 (because, that required $x_i \sim N(\mu, 1)$)

$$\Gamma X \sim N(\mu, \sigma^2 I_n) \quad \text{or} \quad X_i \sim N(\mu, \sigma^2) \quad \forall i = 1, 2, \dots, n$$

$$Y \sim N\left(\frac{\mu}{\sigma}, 1\right) \Leftarrow \because Y_i \sim N\left(\frac{\mu}{\sigma}, 1\right)$$

because, that required
 $X_i \sim N(\mu_i, \sigma^2)$
 but we have,
 $X_i \stackrel{iid}{\sim} N(\mu, \sigma^2)$
 \nexists (\therefore we need to scale it)

$$\therefore \frac{s^2}{\sigma^2} = \left(\frac{\bar{X}}{\sigma}\right)^T A \left(\frac{\bar{X}}{\sigma}\right) = \bar{Y}^T A \bar{Y}$$

which does follow $\chi^2_{df} = \text{rank}(A)$
 $mcp = (\mu^*)^2 \sum_{i=1}^n A_i^T A_i$

Now, $\text{rank}(A) = n-1$ (trivial)
 $mcp = (\mu^*)^2 \sum_{i=1}^n A_i^T A_i = 0$ (trivial).

↑ just
substituted
on A
open it up

∴ Finally, $\left(\frac{s^2}{\sigma^2}\right) \sim \chi^2_{n-1} = df$

→ And ultimately, we assimilate: $\bar{X} \sim N(\mu, \frac{\sigma^2}{n}) \Rightarrow \left(\frac{\bar{X}-\mu}{\sigma/\sqrt{n}}\right) \sim N(0,1)$

$\frac{s^2}{\sigma^2} \sim \chi^2_{n-1}$ and $\bar{X}, \frac{s^2}{\sigma^2}$ are independent

$$\therefore t = \frac{\sqrt{n}(\bar{X}-\mu)}{\sqrt{s^2/(n-1)}} = \frac{\bar{X}-\mu}{\sigma/\sqrt{n}} \sim t_{df=n-1}$$

$$\therefore t \sim t_{df=n-1} \quad mcp=0$$

→ (Key result used in ANOVA)

Let $\bar{X} \sim N(\mu, I_n)$ and $\bar{X}^T A \bar{X} = \sum_i \bar{X}^T A_i \bar{X}$

when, $A_i^T = A_i$ and $A_i^2 = A_i$ (can be shown by splitting A)

then A_i 's are independent (shown by splitting A)

and, $\bar{X}^T A_i^T \bar{X} \sim df = \text{rank}(A_i)$

$$= \left(\frac{\bar{X}}{\sigma}\right)^T \frac{1}{\sigma} \sum_i A_i^T A_i \frac{1}{\sigma} \bar{X} = \bar{X}^T \frac{1}{\sigma^2} \sum_i A_i^T A_i \bar{X} = \bar{X}^T \frac{1}{\sigma^2} A \bar{X}$$

$$\left(\underbrace{\frac{1}{\sigma^2} \begin{pmatrix} 1 & 0 & 0 \\ 0 & 1 & 0 \\ 0 & 0 & 1 \end{pmatrix}}_{A} = \begin{pmatrix} 1 & 0 & 0 \\ 0 & 1 & 0 \\ 0 & 0 & 1 \end{pmatrix}\right)$$

$$= A_1 + A_2 + A_3$$

$$\left(\underbrace{\frac{1}{\sigma^2} \begin{pmatrix} 1 & 0 & 0 \\ 0 & 1 & 0 \\ 0 & 0 & 1 \end{pmatrix}}_{A} = \begin{pmatrix} 1 & 0 & 0 \\ 0 & 1 & 0 \\ 0 & 0 & 1 \end{pmatrix}\right)$$

Simple Linear Regression Model $\rightarrow Y_i = \beta_0 + \beta_1 X_i + \epsilon_i$

↓ Added assumption $\epsilon_i \stackrel{iid}{\sim} N(0, \sigma^2)$ (Assumption 1)

Gauss-Markov Model $\therefore \text{Unknowns: } (\beta_0, \beta_1, \sigma^2)$

In matrix notation: $Y = \begin{bmatrix} Y_1 \\ Y_2 \\ \vdots \\ Y_n \end{bmatrix}, \epsilon = \begin{bmatrix} \epsilon_1 \\ \epsilon_2 \\ \vdots \\ \epsilon_n \end{bmatrix}, \beta = \begin{bmatrix} \beta_0 \\ \beta_1 \end{bmatrix}, X = \begin{bmatrix} 1 & X_1 \\ 1 & X_2 \\ \vdots & \vdots \\ 1 & X_n \end{bmatrix}$

$$\therefore Y = X\beta + \epsilon$$

Now, $\epsilon \sim N(0, \sigma^2 I_n)$

$$\Rightarrow Y \sim N(X\beta, \sigma^2 I_n) \quad (\text{from multivariate analysis})$$

Note: ϵ_i are independent & identically distributed

but, Y_i are independent but Not identically distributed (they have different mean $\rightarrow X\beta_i$)

Estimation of $\beta_0, \beta_1, \sigma^2$

1) Least squares method (LS)

2) Maximum Likelihood Method (MLE)

3) L.S. condition is to minimize $S(\beta_0, \beta_1) = \sum_{i=1}^n (Y_i - \beta_0 - \beta_1 X_i)^2$ (a.k.a RSS)

$$= \|Y - X\beta\|_2^2 \quad (\text{square of Euclidean norm})$$

$$= (Y - X\beta)^T (Y - X\beta)$$

$$\left(\begin{array}{c} \hat{\beta}_0 \\ \hat{\beta}_1 \end{array}\right) = \arg \min_{(\beta_0, \beta_1)} S(\beta_0, \beta_1)$$

Now, either we could do individual differentiation $\frac{\partial}{\partial}(\cdot)^2$ or, vector differentiation of last expression

Let's try both!

$$\frac{\partial}{\partial}(\beta_0, \beta_1) \min \rightarrow \frac{\partial}{\partial} S(\beta_0, \beta_1)$$

Method 1: Vector differentiation.

$$\text{Hint: } \frac{\partial \hat{\beta}^T A}{\partial \hat{\beta}} = A^T, \quad \frac{\partial \hat{\beta}^T A \hat{\beta}}{\partial \hat{\beta}} = \hat{\beta}^T (A + A^T)$$

$$\frac{\partial S}{\partial \hat{\beta}} = \frac{\partial (Y^T Y - Y^T X \hat{\beta} - \hat{\beta}^T X^T Y + \hat{\beta}^T X^T X \hat{\beta})}{\partial \hat{\beta}}$$

$$\therefore 0 = 0 - \frac{\partial}{\partial \hat{\beta}} (2 \hat{\beta}^T X^T Y) + \frac{\partial}{\partial \hat{\beta}} (\hat{\beta}^T X^T X \hat{\beta})$$

$$\therefore 0 = -2(X^T Y)^T + \hat{\beta}^T \underbrace{(X^T X + (X^T X)^T)}_{2X^T X}$$

$$\text{or, } (X^T Y)^T = \hat{\beta}^T (X^T X)$$

$$\text{or, } (X^T Y) = \underbrace{I}_{(X^T X)} (X^T X)^T \hat{\beta}$$

($\because X^T X \hat{\beta}$ and $\hat{\beta}^T X^T Y$ are scalars & $(\text{scalar})^T = \text{scalar}$)
 $\therefore (Y^T X \hat{\beta})^T = \hat{\beta}^T (Y^T X)^T = \hat{\beta}^T X^T Y$.

$$\left(\begin{array}{l} \text{For 2 variable case} \\ \text{(i.e. "Simple" Regression)} \end{array} \right) : \quad \begin{aligned} \hat{\beta}_0 &= \bar{y} - \hat{\beta}_1 \bar{x}, \\ \hat{\beta}_1 &= \frac{\sum (x_i - \bar{x})(y_i - \bar{y})}{\sum (x_i - \bar{x})^2} = \left(\frac{S_{xy}}{S_{xx}} \right) \end{aligned} \quad \left(\begin{array}{l} \text{(i.e. symmetric & idemp)} \end{array} \right)$$

Prediction

$$\hat{y}_i = \hat{\beta}_0 + \hat{\beta}_1 x_i \quad \forall i = 1, 2, \dots, n$$

$$\text{or, } \hat{Y} = X \hat{\beta} = X (X^T X)^{-1} (X^T Y).$$

$$\boxed{\text{Observation: We can identify this as the projection of } Y \text{ on the column space of } X.}$$

In fact, this is not just any projection, but an orthogonal projection of Y on the $\mathcal{R}(X)$.

hence, we shall call it as: P_X

$$\therefore \hat{Y} = P_X Y$$

Further, the error vector in the LS prediction is: $e = Y - P_X Y = (I_n - P_X) Y$ (Wow!)

① \hat{e} is orthogonal to \hat{Y}

(Note: orthogonality of e and \hat{Y} implies their uncorrelation. (i.e. an added assumption of normality will imply their independence))

② Although e has covariance matrix $\sigma^2 I_n$, e does not have independent components.

$$\therefore \hat{Y} \sim N(X \hat{\beta}, \sigma^2 I_n)$$

$$\therefore P_X Y \sim N(P_X X \hat{\beta}, \sigma^2 P_X^T P_X) \text{ and, } (I - P_X) Y \sim N(0, \sigma^2 (I_n - P_X))$$

$$\therefore (e, \hat{Y}) \sim N(0, \sigma^2 (I_n - P_X)) \quad \text{(or, } e = Y - \hat{Y})$$

(since, P_X is a projection matrix & hence $P_X^T = P_X = P_X^2$)

(since, projection of X on its own columnspace will be X itself)

Essentially, we just use:-

given: $Y \sim N(X \hat{\beta}, \sigma^2 I_n)$

$\cdot A Y \sim N(A \hat{\beta}, A^T A \sigma^2)$ again & again

keep in mind

$$X = \begin{bmatrix} 1 & x_1 \\ 1 & x_2 \\ \vdots & \vdots \\ 1 & x_n \end{bmatrix}_{n \times 2}$$

$$\text{Also, } \text{cov}(AY, BY) = (A^T A \sigma^2)$$

$$\therefore \hat{\beta} = [(X^T X)^{-1} X^T Y] \sim N\left([(X^T X)^{-1} X^T X \hat{\beta}], \sigma^2 (X^T X)^{-1} X^T I_n X (X^T X)^{-1}\right)$$

Note: $\hat{\beta}$ need not be independent vector

(i.e. $\hat{\beta}_1, \hat{\beta}_2, \hat{\beta}_3, \dots$ may not be indep.)

(i.e. $(X^T X)^{-1}$ need not be diagonal)

* * * (i.e. $\text{cov}(\beta_i, \beta_j)$ need not be zero for all $i \neq j$)

$$\therefore \hat{\beta} \sim N(\hat{\beta}, (X^T X)^{-1} \sigma^2 I_n)$$

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

↳ unbiased ($\therefore \hat{\beta} \sim N(\beta, (X^T X)^{-1} \sigma^2 I_n)$)

$$\therefore E(\hat{\beta}) = \beta$$

Now, squared error/error sum of square/residual sum of square :-

$$RSS \equiv \sum e_i^2 \equiv \mathbf{e}^T \mathbf{e} = [(\mathbf{I} - \mathbf{P}_X) \mathbf{y}]^T [(\mathbf{I} - \mathbf{P}_X) \mathbf{y}] = \mathbf{y}^T (\mathbf{I} - \mathbf{P}_X) \mathbf{y} \quad (\because \mathbf{P}_X^2 = \mathbf{P}_X, \mathbf{P}_X^T = \mathbf{P}_X)$$

$$\therefore RSS \sim \chi^2_{df} \quad (\text{NO})$$

* we need to scale it down

$$\frac{RSS}{\sigma^2} \sim \chi^2_{df} = \text{rank}(\mathbf{I}_n - \mathbf{P}_X)$$

$$nep = \left[\left(\frac{\mathbf{x}_B}{\sigma} \right)^T (\mathbf{I}_n - \mathbf{P}_X) \left(\frac{\mathbf{x}_B}{\sigma} \right) \right]$$

(obviously)

rank = (n-2)

rank = 2

(assuming inverse exists)

Using Result: If A is an idempotent matrix and $\mathbf{z} \sim N(\mathbf{0}, \mathbf{I}_n)$

$$\text{then, } \mathbf{z}^T \mathbf{A} \mathbf{z} \sim \chi^2_{\text{rank}(A)}, \mathbf{z}^T \mathbf{A} \mathbf{z} \downarrow \downarrow$$

Now, $df = n-2$ ✓

df

nep

$$nep = (\mathbf{p}_B^T \mathbf{x}^T (\mathbf{I}_n - \mathbf{P}_X) \mathbf{x} \mathbf{p}_B) / \sigma^2$$

$$= \mathbf{p}_B^T \mathbf{O} \mathbf{p}_B / \sigma^2 = 0 \quad \checkmark$$

For 2 variable case (i.e. $\mathbf{Y} = \beta_0 + \beta_1 \mathbf{x} + \epsilon$)

$$\hat{\beta}_0 \sim N(\beta_0, (\frac{1}{n} + \frac{\mathbf{x}^2}{\mathbf{s}_{xx}}) \sigma^2)$$

$$\hat{\beta}_1 \sim N(\beta_1, \frac{\sigma^2}{\mathbf{s}_{xx}})$$

$$RSS \sim \chi^2_{df} = n-2, nep = 0.$$

$$\hat{\sigma}^2 = \frac{RSS}{n-2} = \frac{\sum (y_i - \hat{y}_i)^2}{n-2}$$

① ϵ and \hat{Y} are orthogonal (trivial) to each other, & ($\because (\mathbf{P}_X \mathbf{y}) \cdot (\mathbf{I}_n - \mathbf{P}_X) \mathbf{y} = \mathbf{y}^T \mathbf{P}_X \mathbf{y}^T (\mathbf{I}_n - \mathbf{P}_X) \mathbf{y} = 0$)

② they are distributionally independent (using, $\text{cov}(\mathbf{A}\mathbf{Y}, \mathbf{B}\mathbf{Y}) = (\mathbf{A}\mathbf{I}\mathbf{B}^T)\sigma^2$)

(Proof: here, $\text{cov}(\hat{Y}, \epsilon) =$ ↘

$$= \text{cov}(\mathbf{P}_X \mathbf{y}, (\mathbf{I}_n - \mathbf{P}_X) \mathbf{y}) = [\mathbf{P}_X \mathbf{I} (\mathbf{P}_X^T - \mathbf{I})] \sigma^2$$

$$= (\mathbf{P}_X \mathbf{P}_X^T - \mathbf{P}_X) \sigma^2 \quad (\because \mathbf{P}_X^T = \mathbf{P}_X = \mathbf{P}_X^2) \\ = (\mathbf{P}_X^2 - \mathbf{P}_X) \sigma^2 = 0$$

Prediction for new x_0 (or, x_0 for 2 var. case)

Recall, we have data $\{(x_i, y_i) | i=1, 2, \dots, n\}$ (2 variable case)

From there we get $\hat{\beta}_0, \hat{\beta}_1$ as LS estimate

so the prediction line or regression line is $\hat{y} = \hat{\beta}_0 + \hat{\beta}_1 x$

Now,

For a new value x_0 under the same model, we can predict the value of y_0 as follows:-

(that is, the same model)

from which the data (x_i, y_i) was obtained

$$\hat{y}_0 = \hat{\beta}_0 + \hat{\beta}_1 x_0 \quad (\text{Not same as } y_0 = \beta_0 + \beta_1 x)$$

* Apart from the $\text{var}(\epsilon) = \sigma^2$, we also have the new sources of uncertainty which impact the distribution of \hat{y}_0 .

(Although, $y_0 \sim N(\beta_0 + \beta_1 x_0, \sigma^2)$);

\hat{y}_0 will have different variance)

So, what is the distribution of \hat{y}_0 ?? (larger)

Ans: $\hat{y}_0 = \hat{\beta}_0 + \hat{\beta}_1 x_0$ ↘ constant

(Normal(Normal) linear comb.)

∴ \hat{y}_0 will follow Normal dist. So, just find $E(\hat{y}_0)$ and $\text{var}(\hat{y}_0)$ & you are done.

- $E(\hat{Y}_0) = E(\hat{\beta}_0 + \hat{\beta}_1 x_0) = (\beta_0 + \beta_1 x_0)$
- $V(\hat{Y}_0) = V(\hat{\beta}_0) + V(\hat{\beta}_1) \cdot x_0^2 + 2\text{cov}(\hat{\beta}_0, \hat{\beta}_1) \cdot x_0$
 $= \sigma^2 \left[\frac{1}{n} + \frac{\bar{x}^2}{S_{xx}} + \frac{x_0^2}{S_{xx}} - 2 \frac{\bar{x}x_0}{S_{xx}} \right]$
 $= \sigma^2 \left[\frac{1}{n} + \frac{(\bar{x} - x_0)^2}{S_{xx}} \right]$

$\therefore \hat{Y}_0 \sim N(\beta_0 + \beta_1 x_0, \sigma^2 \left(\frac{1}{n} + \frac{(\bar{x} - x_0)^2}{S_{xx}} \right))$

As promised,
we have
different variance
(in fact larger).

\rightarrow Not same as σ^2
($\because V(Y_0)$)

only self terms will have non-zero covariance ($\because \text{cov}(Y_i, Y_j) = 0 \ (i \neq j)$)
 $\therefore \text{cov}(\hat{\beta}_0, \hat{\beta}_1) = \sum k_i g_i \sigma^2$
 $= \sigma^2 \left(\sum \frac{k_i}{n} - \bar{x} \sum k_i \bar{x}^2 \right) = \frac{\sigma^2 \bar{x}}{S_{xx}}$

$\hat{\beta}_1 = \sum \left(\frac{x_i - \bar{x}}{S_{xx}} \right) y_i = (k_1 y_1 + k_2 y_2 + \dots)$
 $\hat{\beta}_0 = \sum \left(\frac{1}{n} - \frac{\bar{x}(x_i - \bar{x})}{S_{xx}} \right) y_i = (g_1 y_1 + g_2 y_2 + \dots)$

Testing of hypothesis.

$H_0 \rightarrow \beta_0 = b_0$
 $H_1 \rightarrow \beta_0 \neq b_0$

Now, $\hat{\beta}_0 \sim N(\beta_0, \sigma^2 \left(\frac{1}{n} + \frac{\bar{x}^2}{S_{xx}} \right))$

$\rightarrow \frac{\hat{\beta}_0 - \beta_0}{\sqrt{\sigma^2 \left(\frac{1}{n} + \frac{\bar{x}^2}{S_{xx}} \right)}} \sim N(0, 1)$

\rightarrow Under H_0 : $\frac{\hat{\beta}_0 - \beta_0}{\sqrt{\sigma^2 \left(\frac{1}{n} + \frac{\bar{x}^2}{S_{xx}} \right)}} \sim N(0, 1)$

* Problem: We don't know the true value of σ^2
But, we just know an estimate $\hat{\sigma}^2$

$\hat{\sigma}^2 = \frac{1}{n-2} \sum (y_i - \bar{y})^2$

$\rightarrow \frac{\hat{\beta}_0 - \beta_0}{\sqrt{\hat{\sigma}^2 \left(\frac{1}{n} + \frac{\bar{x}^2}{S_{xx}} \right)}} \sim t_{n-2}$

\rightarrow Under H_0 : $\frac{\hat{\beta}_0 - \beta_0}{\sqrt{\hat{\sigma}^2 \left(\frac{1}{n} + \frac{\bar{x}^2}{S_{xx}} \right)}} \sim t_{n-2}$

($= S_{yy} - \frac{S_{xy}^2}{S_{xx}}$ (final result, computationally shortwrt))

\therefore Now,

if $|t_{\text{computed}}| > t_{\frac{\alpha}{2}, n-2}$,
 $\quad \quad \quad$ (or t_{stat})

we reject [null hypothesis (H_0)] at level of significance α .

(In favour of H_1)

consequently, we can quickly obtain the $(1-\alpha)$ confidence interval for β_0 .

Similarly,

$H_0 \rightarrow \beta_1 = b_1$
 $H_1 \rightarrow \beta_1 \neq b_1$

Now, $\hat{\beta}_1 \sim N(\beta_1, \sigma^2 \left(\frac{1}{S_{xx}} \right))$

but σ^2 is unknown,
so, we rather use: $\frac{\hat{\beta}_1 - \beta_1}{\sqrt{\hat{\sigma}^2 \left(\frac{1}{S_{xx}} \right)}} \sim t_{n-2}$

\rightarrow Under H_0 : $\frac{\hat{\beta}_1 - \beta_1}{\sqrt{\hat{\sigma}^2 \left(\frac{1}{S_{xx}} \right)}} \sim t_{n-2}$

\therefore Now,

if $|t_{\text{stat}}| > t_{\frac{\alpha}{2}, n-2}$

we reject [null hypothesis (H_0)] at level of significance α .

Once again, we can obtain the $(1-\alpha)$ confidence interval of β_1 as follows:-

$$\frac{\hat{\beta}_1 - \beta_1}{\sqrt{\hat{\sigma}^2 / S_{xx}}} \sim t_{n-2} \Rightarrow P(-t_{\frac{\alpha}{2}, n-2} < T < t_{\frac{\alpha}{2}, n-2}) = 1-\alpha \Rightarrow P(\hat{\beta}_1 - t_{\frac{\alpha}{2}, n-2} \sqrt{\frac{\hat{\sigma}^2}{S_{xx}}} < \beta_1 < \hat{\beta}_1 + t_{\frac{\alpha}{2}, n-2} \sqrt{\frac{\hat{\sigma}^2}{S_{xx}}}) = 1-\alpha$$

* (Interval is a random variable)

$$\boxed{P(\hat{\beta}_1 - t_{\frac{\alpha}{2}, n-2} \sqrt{\hat{\sigma}^2 / s_{xx}} < \beta_1 < \hat{\beta}_1 + t_{\frac{\alpha}{2}, n-2} \sqrt{\hat{\sigma}^2 / s_{xx}}) = 1-\alpha}$$

Note: For a given sample, we have a given realization of this interval & in that context, the probability that β_1 is contained in it or not is either 1 or 0; but for the interval estimator the probability is $1-\alpha$ (ie if you keep taking out random samples & keep realizing 1 or 0, the frequency of 1 on avg./long run will be $1-\alpha$)

Informally,

* similarly,

$$\boxed{P(\hat{\beta}_0 - t_{\frac{\alpha}{2}, n-2} \sqrt{\hat{\sigma}^2 (\frac{1}{n} + \frac{\bar{x}^2}{s_{xx}})} < \beta_0 < \hat{\beta}_0 + t_{\frac{\alpha}{2}, n-2} \sqrt{\hat{\sigma}^2 (\frac{1}{n} + \frac{\bar{x}^2}{s_{xx}})}) = 1-\alpha}$$

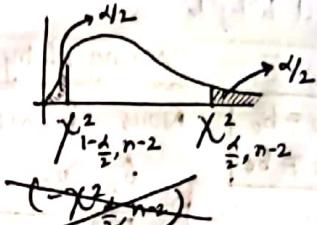
* What about σ^2 ? (Is $(1-\alpha)$ CI for σ^2 ?)

Sol: Now, $\frac{RSS}{\sigma^2} \sim \chi^2_{n-2}$

(caution: χ^2 is Not a symmetric distriⁿ (unlike t distriⁿ))

$$\therefore P(\chi^2_{1-\frac{\alpha}{2}, n-2} < \frac{RSS}{\sigma^2} < \chi^2_{\frac{\alpha}{2}, n-2}) = 1-\alpha$$

$$\frac{(n-2)\hat{\sigma}^2}{\sigma^2} \text{ or, } P\left(\frac{1}{\chi^2_{1-\frac{\alpha}{2}, n-2}} > \frac{\sigma^2}{RSS} > \frac{1}{\chi^2_{\frac{\alpha}{2}, n-2}}\right) = 1-\alpha$$



$$\therefore P\left(\frac{RSS}{\chi^2_{\frac{\alpha}{2}, n-2}} < \sigma^2 < \frac{RSS}{\chi^2_{1-\frac{\alpha}{2}, n-2}}\right) = 1-\alpha$$

(Note: $RSS = \sum (y_i - \hat{y}_i)^2$)

$$\left(\equiv P\left(\frac{(n-2)\hat{\sigma}^2}{\chi^2_{\frac{\alpha}{2}, n-2}} < \sigma^2 < \frac{(n-2)\hat{\sigma}^2}{\chi^2_{1-\frac{\alpha}{2}, n-2}}\right) = 1-\alpha \right)$$

$\hat{\sigma}^2 = RSS/(n-2)$

* Finally, PREDICTION INTERVAL

(we don't call it by the generic term "confidence interval" to distinguish that we are no longer constructing an interval for a parameter (ie \hat{y}_0 is not a parameter but a random variable itself))

$$\hat{y}_0 \sim N(\hat{y}_0, \sigma^2 (1 + \frac{1}{n} + \frac{(x_0 - \bar{x})^2}{s_{xx}}))$$

$$\text{or, } \hat{y}_0 - y_0 \sim N(0, \sigma^2 \hat{y}_0)$$

$$\text{or, } \frac{\hat{y}_0 - y_0}{\sqrt{\sigma^2 \hat{y}_0}} \sim N(0, 1) \text{ but } \sigma^2 \text{ is unknown (hence } \sigma^2 \hat{y}_0 \text{ is unknown)}$$

$$\therefore \text{we rather use: } \frac{\hat{y}_0 - y_0}{\sqrt{\hat{\sigma}^2 \hat{y}_0}} \sim t_{n-2} \quad (\hat{\sigma}^2 \hat{y}_0 = \hat{\sigma}^2 (1 + \frac{1}{n} + \frac{(x_0 - \bar{x})^2}{s_{xx}}))$$

$$\Rightarrow P(-t_{\frac{\alpha}{2}, n-2} < \frac{\hat{y}_0 - y_0}{\sqrt{\hat{\sigma}^2 \hat{y}_0}} < t_{\frac{\alpha}{2}, n-2}) = 1-\alpha$$

$$\Rightarrow P(\hat{y}_0 - t_{\frac{\alpha}{2}, n-2} \sqrt{\hat{\sigma}^2 \hat{y}_0} < y_0 < \hat{y}_0 + t_{\frac{\alpha}{2}, n-2} \sqrt{\hat{\sigma}^2 \hat{y}_0}) = 1-\alpha$$

carefully note:
 $\sigma^2 \hat{y}_0$ (v/s $\sigma^2 y_0$)
&
 $\hat{\sigma}^2 \hat{y}_0$ (v/s $\hat{\sigma}^2 y_0$)

Scanned by CamScanner

Maximum Likelihood Estimation of regression parameters.

Here, 'even for point estimation' we need the distributional assumption;

$$Y_i = \beta_0 + \beta_1 x_i + \epsilon_i, \quad \epsilon_i \sim \text{iid } N(0, \sigma^2)$$

$\Rightarrow Y_i \sim N(\beta_0 + \beta_1 x_i, \sigma^2)$ and independent. (already known)

Now,

$$\begin{aligned} (\hat{\beta}_0, \hat{\beta}_1, \hat{\sigma}^2)_{\text{MLE}} &= \underset{\beta_0, \beta_1, \sigma^2}{\arg \max.} \prod_{i=1}^n \frac{e^{-\frac{1}{2} \frac{(Y_i - \beta_0 - \beta_1 x_i)^2}{\sigma^2}}}{\sqrt{2\pi} \sigma} \\ &= \underset{\beta_0, \beta_1, \sigma^2}{\arg \max.} \frac{e^{-\frac{1}{2\sigma^2} \sum_i (y_i - \beta_0 - \beta_1 x_i)^2}}{(\sqrt{2\pi} \sigma)^n} \end{aligned}$$

observe: (if we take logarithm to simplify)

maximizing log likelihood = minimizing $\sum_i (y_i - \beta_0 - \beta_1 x_i)^2 / \text{RSS}$ = OLS estimate.

(In above case)

$\hat{\beta}_{LS} = \hat{\beta}_{MLE}$, $\hat{\beta}_{LS}^2 = \hat{\beta}_{MLE}^2$; but $\hat{\sigma}_{LS}^2 = \text{RSS}/n-2$ (whereas $\hat{\sigma}_{MLE}^2 = \text{RSS}/n$).
 ↳ unbiased. ↳ ~~biased~~ biased

When X too, is a random variable.

Suppose, (X, Y) is a paired random variable with Bivariate Normal dist^r $N(\mu_X, \mu_Y, \sigma_X^2, \sigma_Y^2, \rho_{XY})$. Given, data set: $D = \{(x_i, y_i), i=1, 2, \dots, n\}$

1) The regression model $\rightarrow E(Y|X=x)$

Now,

$$f(x, y) = e^{-\frac{1}{2(1-\rho^2)} \left\{ \left(\frac{x-\mu_X}{\sigma_X} \right)^2 + \left(\frac{y-\mu_Y}{\sigma_Y} \right)^2 - 2\rho \left(\frac{x-\mu_X}{\sigma_X} \right) \left(\frac{y-\mu_Y}{\sigma_Y} \right) \right\}} / 2\pi \sigma_X \sigma_Y \sqrt{1-\rho^2}$$

2) Similarly, we could've chosen to study the other potential regression model $\rightarrow E(X|Y=y)$

// Anyways, the results will be symmetric.

we want $E(Y|X=x)$: [So, take out $e^{-\frac{1}{2(1-\rho^2)} \left\{ \left(\frac{x-\mu_X}{\sigma_X} \right)^2 \right\}}$ and complete remaining square" in terms of $e^{-\frac{1}{2(1-\rho^2)} \left\{ y - \text{something} \right\}^2}$ & obtain $E(Y|x=x)$]

$$\text{or, } f(x, y) = \left[e^{-\frac{1}{2(1-\rho^2)} \left[\left(\frac{y-\mu_Y}{\sigma_Y} \right) - \rho \left(\frac{x-\mu_X}{\sigma_X} \right) \right]^2} \right] \left[e^{-\frac{1}{2} \left(\frac{x-\mu_X}{\sigma_X} \right)^2} \right] / \sqrt{2\pi} (\sqrt{1-\rho^2} \sigma_Y)$$

Now, we need conditional dist^r of $(Y|X=x)$

$$\begin{aligned} f_{Y|X=x}(y|x) &= \frac{\text{joint q } X \& Y}{\text{marginal of } X} = \frac{e^{-\frac{1}{2(1-\rho^2)} \left[\left(\frac{y-\mu_Y}{\sigma_Y} \right) - \rho \left(\frac{x-\mu_X}{\sigma_X} \right) \right]^2}}{\sqrt{2\pi} (\sigma_Y \sqrt{1-\rho^2})} \\ &= e^{-\frac{1}{2(1-\rho^2)\sigma_Y^2} \left[y - (\mu_Y + \rho \frac{\sigma_Y}{\sigma_X} (x-\mu_X)) \right]^2} \end{aligned}$$

$$\therefore (Y|X=x) \sim N\left(\mu_Y + \rho \frac{\sigma_Y}{\sigma_X} (x-\mu_X), \sigma_Y^2 (1-\rho^2)\right)$$

$$\hookrightarrow E(Y|X=x) = \mu_Y + \rho \frac{\sigma_Y}{\sigma_X} (x-\mu_X) \quad \text{and, } V(Y|X=x) = \sigma_Y^2 (1-\rho^2).$$

$$\text{or, } E(Y|X=x) = \underbrace{\left(\mu_Y - \rho \frac{\sigma_Y}{\sigma_X} \mu_X \right)}_{\beta_0} + \underbrace{\left(\rho \frac{\sigma_Y}{\sigma_X} \right) x}_{\beta_1}$$

(conditional variance)
 ✪ (\leq uncond^b variance)

Now, testing

whether x is significant or not

= whether $\beta_1 = 0$ or not

= whether $\rho = 0$ or not

$$\begin{array}{l} H_0: \rho = 0 \\ H_1: \rho \neq 0 \end{array}$$

Now, $\hat{\rho} \leq r = \frac{\sum (x_i - \bar{x})(y_i - \bar{y})}{\sqrt{\sum (x_i - \bar{x})^2} \sqrt{\sum (y_i - \bar{y})^2}}$

and, under H_0 :

$$T = \frac{r \sqrt{n-2}}{\sqrt{1-r^2}} \sim t_{n-2, n \rho=0}$$

(Generalized result:

$$\begin{array}{l} H_0: \rho = \rho_0 \\ H_1: \rho \neq \rho_0 \end{array}$$

Now, define $Z = \frac{1}{2} \log \left(\frac{1+r}{1-r} \right) (= \tanh^{-1}(r))$

$$\mu_Z = \frac{1}{2} \log \left(\frac{1+\rho_0}{1-\rho_0} \right) (= \tanh^{-1}(\rho_0))$$

$$\sigma_Z^2 = (n-3)^{-1}$$

and, $\left(\frac{Z - \mu_Z}{\sigma_Z} \right) \sim N(0, 1)$ when $n \rightarrow \infty$

[based on the famous variance stabilization formula]

(& its special cases

Box-Cox method)

$$\text{then, } \sqrt{n} (T_n - \mu) \xrightarrow{d} N(0, \sigma^2)$$

$$\text{then, } \sqrt{n} (g(T_n) - g(\mu)) \xrightarrow{d} N(0, \sigma^2 (g'(\mu))^2)$$

"We reject H_0 in favour of H_1 " vs ("we reject H_0 ") → meaningless

There is no such thing as "Rejecting H_0 "

It's always "in relative/ reference/ favour of H_1 "

MULTIPLE LINEAR REGRESSION

$y = \begin{pmatrix} y_1 \\ y_2 \\ y_3 \\ \vdots \\ y_n \end{pmatrix}$ vs. & intercept

$$x = \begin{pmatrix} 1 & x_{11} & x_{12} & x_{13} & \dots & x_{1k} \\ 1 & x_{21} & x_{22} & x_{23} & \dots & x_{2k} \\ 1 & x_{31} & x_{32} & x_{33} & \dots & x_{3k} \\ \vdots & \vdots & \vdots & \vdots & \ddots & \vdots \\ 1 & x_{n1} & x_{n2} & x_{n3} & \dots & x_{nk} \end{pmatrix}$$

$$\beta = \begin{pmatrix} \beta_0 \\ \beta_1 \\ \vdots \\ \beta_k \end{pmatrix}, e = \begin{pmatrix} e_1 \\ e_2 \\ \vdots \\ e_n \end{pmatrix}$$

$$y = x\beta + e, e \sim N(0, \sigma^2 I_n)$$

LS condition: $S^2 = (y - x\beta)^T (y - x\beta)$
 $* \hat{\beta} = \arg \min_{\beta \in \mathbb{R}^{k+1}} S^2(\beta)$

Now, $H_0: \beta_1 = \beta_2 = \dots = \beta_k = 0$ vs. $H_1: \text{at least one is non-zero}$
 $\hat{\beta}_{LS} = (x^T x)^{-1} x^T y$ we get $\hat{\beta}_{LS} \neq (x^T x)^{-1} x^T y$ (assume $x^T x \neq 0$)
 \leftrightarrow whether we should build a model or not??

Under H_0 : $y_i \sim N(\beta_0, \sigma^2)$

$$(NOT \quad \beta_0 + \beta_1 x_1 + \beta_2 x_2 + \dots + \beta_k x_k)$$

$$((x-x)^T)^{-1} x^T y = (x^T x)^{-1} x^T y$$

$$((x-x)^T)^{-1} x^T y = ((x^T x)^{-1} x^T x)^T y = I y$$

* The famous ANOVA identity.

$$TSS = ESS + RSS$$

$$\sum (y_i - \bar{y})^2 = \sum (\hat{y}_i - \bar{y})^2 + \sum (y_i - \hat{y}_i)^2$$

$$\sum (y_i - \bar{y})^2 = \sum (y_i - \hat{y}_i)^2 + \sum (\hat{y}_i - \bar{y})^2$$

(Proof) \Rightarrow

In matrix notation:

$$TSS = y^T (I_n - \frac{1}{n} \mathbf{1} \mathbf{1}^T) y$$

$$ESS = y^T (P_x - \frac{1}{n} \mathbf{1} \mathbf{1}^T) y$$

Approach (to test this):

ANOVA

(does the model explain any variability in response variable(y))

$$\sum (y_i - \bar{y})^2$$

Total sum of squares

$$\sum (\hat{y}_i - \bar{y})^2$$

Explained sum of squares

$$\text{and, } RSS = y^T (I_n - P_x) y$$

$$(\text{show that } TSS = ESS + RSS) : ESS + RSS = \tilde{Y}^T (P_X - \frac{1}{n} \mathbf{1} \mathbf{1}^T) \tilde{Y} + \tilde{Y}^T (I_n - P_X) \tilde{Y}$$

$$= \tilde{Y}^T (P_X - \frac{1}{n} \mathbf{1} \mathbf{1}^T) \tilde{Y} + I_n - P_X \tilde{Y}$$

NOW, we are ready to apply Cochran's thm $= \tilde{Y}^T (I_n - \frac{1}{n} \mathbf{1} \mathbf{1}^T) \tilde{Y} = TSS$ (hence proved)
 * (Recall we had said, "Key result in ANOVA").

- $\left(\frac{RSS}{\sigma^2} \right) \sim \chi^2_{n-(k+1)}, ncp=0$

(Note: $A = A_1 + A_2 + \dots + A_k$,

then,

$$\text{rank}(A) = \sum \text{rank}(A_i)$$

only when A_i are

symmetric
idempotent

- $\left(\frac{ESS}{\sigma^2} \right) \sim \chi^2_k, ncp=1$

- $\left(\frac{TSS}{\sigma^2} \right) \sim \chi^2_{n-1}, ncp=1$

what is λ ?

$$ncp = \lambda = \frac{1}{\sigma^2} (\tilde{X} \beta_2)^T (P_X - \frac{1}{n} \mathbf{1} \mathbf{1}^T) (\tilde{X} \beta_2)$$

$$= \frac{1}{\sigma^2} \beta_2^T (\tilde{X} P_X \tilde{X} - \frac{1}{n} \mathbf{1} \mathbf{1}^T \mathbf{1} \mathbf{1}^T) \beta_2 = \frac{1}{\sigma^2} \beta_2^T (\tilde{X} \tilde{X} - \frac{1}{n} \mathbf{1} \mathbf{1}^T \mathbf{1} \mathbf{1}^T) \beta_2 \quad (\because P_X \tilde{X} = \tilde{X})$$

To conveniently solve this, we partition \tilde{X} as $\tilde{X} = (\mathbf{1} : X_R)$ ($\mathbf{1} \beta_2 = \begin{pmatrix} \beta_0 \\ \beta_R \end{pmatrix}$)

$$= \frac{1}{\sigma^2} \beta_2^T \left(\begin{bmatrix} \mathbf{1}^T \mathbf{1} & \mathbf{1}^T X_R \\ X_R^T \mathbf{1} & X_R^T X_R \end{bmatrix} - \frac{1}{n} \begin{bmatrix} \mathbf{1}^T \\ X_R^T \end{bmatrix} (\mathbf{1} \mathbf{1}^T) (\mathbf{1} X_R) \right) \beta_2$$

$$= \frac{1}{\sigma^2} \beta_2^T \left(\begin{bmatrix} \mathbf{1}^T \mathbf{1} & \mathbf{1}^T X_R \\ X_R^T \mathbf{1} & X_R^T X_R \end{bmatrix} - \frac{1}{n} \begin{bmatrix} \mathbf{1}^T \\ X_R^T \end{bmatrix} \mathbf{1}^T \right) \beta_2$$

$$= \frac{1}{\sigma^2} \beta_2^T \left(\begin{bmatrix} \mathbf{1}^T \mathbf{1} & \mathbf{1}^T X_R \\ X_R^T \mathbf{1} & X_R^T X_R \end{bmatrix} - \begin{bmatrix} \mathbf{1}^T \mathbf{1} & \mathbf{1}^T X_R \\ X_R^T \mathbf{1} & \frac{1}{n} X_R^T \mathbf{1} \mathbf{1}^T X_R \end{bmatrix} \right) \beta_2$$

$$= \frac{1}{\sigma^2} \beta_2^T \left(\begin{bmatrix} 0 & \mathbf{0}^T \\ \mathbf{0}^T & X_R^T X_R - \frac{1}{n} X_R^T \mathbf{1} \mathbf{1}^T X_R \end{bmatrix} \right) \beta_2$$

$$= \frac{1}{\sigma^2} \beta_2^T (X_R^T X_R - \frac{1}{n} X_R^T \mathbf{1} \mathbf{1}^T X_R) \beta_2 \quad (\text{i.e. independent of } \beta_0)$$

$$= \frac{1}{\sigma^2} \beta_2^T X_R^T \beta_2 (I_n - \frac{1}{n} \mathbf{1} \mathbf{1}^T) X_R \beta_2$$

Now, the $ncp = 0$ iff all $\beta_1, \beta_2, \dots, \beta_k = 0$ (i.e. $\beta_R = \mathbf{0}$)

(A) (Not β_0)

which is nothing but the null hypothesis of F test ($H_0: \beta_1 = \beta_2 = \dots = \beta_k = 0$)

Moving forward,

To test: $H_0: \beta_1 = \beta_2 = \dots = \beta_k = 0$
 $H_1: \text{at least one is non-zero}$

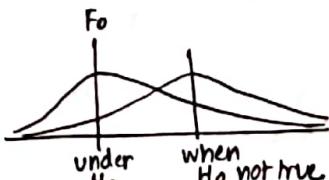
Now, $\frac{ESS}{\sigma^2} \sim \chi^2_k, ncp=\lambda$

but, we don't know σ^2

so, instead we have: $\hat{\sigma}^2 = RSS/m-(k+1)$

Now, $\frac{ESS/k}{RSS/m-(k+1)} \sim F_{k, n-k-1}, ncp=\lambda$.

And, under H_0 : $F_0 = \left[\frac{ESS/k}{RSS/m-(k+1)} \right] \sim F_{k, n-k-1}$
 $(\lambda=0 \Leftrightarrow \text{central F dist} \Leftrightarrow \text{noncentral F})$



(i.e. central F distribution)
 (v/s. noncentral F).

∴ We reject null hypothesis in favour of H_1 iff $F_0 > F_{\alpha, k, n-k-1}$ at α level of significance

(Hence, we are studying 'mean' through variance, hence the name is "Analysis of variance for mean")

↳ heavily employed in design & Analysis of Experiments

Now, (testing individual variable's significance) $\rightarrow C$

$$\begin{array}{l} H_0: \beta_j = b_j \\ \text{v/s} \\ H_1: \beta_j \neq b_j \end{array}$$

where, $\hat{\beta}_n \sim N(\beta_n, (X^T X)^{-1} \sigma^2)$

we need to extract out the distribution of $\hat{\beta}_j$

$$\Rightarrow \hat{\beta}_j \sim N(\beta_j, \sigma^2 c_{j+1,j+1})$$

$$X = \begin{bmatrix} 1 & 2 & 3 & \dots & k \end{bmatrix}$$

$$\text{and, under } H_0: T = \frac{\hat{\beta}_j - b_j}{\hat{\sigma} \sqrt{c_{j+1,j+1}}} \sim t_{n-k-1}, \text{ where } \hat{\sigma}^2 = \frac{RSS}{n-k-1}$$

\therefore we reject H_0 in favour of H_1 , if $|T| > t_{\frac{\alpha}{2}, n-k-1}$

Let's test more general hypotheses!

$$\begin{array}{l} H_0: \beta_j - \theta \beta_i = b \\ \text{v/s} \\ H_1: \beta_j - \theta \beta_i \neq b \end{array}$$

$$\begin{array}{l} H_0: \beta_j - \theta \beta_i - k \beta_k = b \\ \text{v/s} \\ H_1: \beta_j - \theta \beta_i - k \beta_k \neq b \end{array}$$

or, in general
any linear combination of β_i 's

$$\begin{array}{l} H_0: \lambda^T \beta_n = b \\ \text{v/s} \\ H_1: \lambda^T \beta_n \neq b \end{array}$$

(In fact ①, ⑩, ⑪ are all special cases), where λ is chosen accordingly;

$$\rightarrow \text{①: } \lambda = \begin{pmatrix} 0 \\ \vdots \\ 1 \\ \vdots \\ 0 \end{pmatrix} \xleftarrow{\text{j-th element}}$$

$$\rightarrow \text{⑩: } \lambda = \begin{pmatrix} 0 \\ \vdots \\ 1 \\ \vdots \\ 0 \end{pmatrix} \xleftarrow{\text{j-th element}} \quad \text{so on.}$$

Now, $\hat{\beta}_n \sim N(\beta_n, \sigma^2 (X^T X)^{-1})$

$$\therefore \lambda^T \hat{\beta}_n \sim N(\lambda^T \beta_n, \sigma^2 \lambda^T (X^T X)^{-1} \lambda)$$

#(Using: $\mathbf{Y} \sim N(\mu, \Sigma)$)

The result in multivariate distribution analysis $\therefore \lambda^T Y \sim N(\lambda^T \mu, \lambda^T \Sigma \lambda)$

Something crucial:- generalized inverse

$$\hat{\beta}_n \sim N(\beta_n, (X^T X)^{-1} \sigma^2)$$

If $|X^T X| \neq 0$, then any linear parametric function of β in the form $\beta^T \beta$ is ESTIMABLE.

(i.e. only then, there exists a linear combination of Y (e.g. $\lambda^T Y$) such that $E(\lambda^T Y) = \beta^T \beta$)

// Motivation: For testing, we first need to do estimation, (i.e. develop estimator(s)) & estimators exist only if the parametric function is ESTIMABLE) unbiased

Best Linear Unbiased Estimator (BLUE)

An unbiased estimator of $\beta^T \beta$ is said to be the BLUE iff it has minimum variance among all linear unbiased estimators of $\beta^T \beta$.

$$U = \{ J^T Y \mid E(J^T Y) = \beta^T \beta \quad \forall \beta \in \mathbb{R}^{k+1} \}$$

$$V(J^T Y) \leq V(J^T Z) \quad \forall J, J \in U$$

Theorem: A linear function of Y is BLUE of its expectation iff it is uncorrelated with all linear zero function (LZF).

$$E(J^T Y) = \beta^T \beta$$

iff $\text{cov}(J^T Y, J^T Z) = 0$ where, $J^T Y$ is a LZF

* Very important result: If $J^T Y$ is a LUE of $\beta^T \beta$ then the corresponding BLUE is $J^T P_X Y$ (gives us a "recipe" to go from LUE \rightarrow BLUE) (Proof: skipped!)

* Every estimable linear parametric function has a unique BLUE

(Proof: assume $\exists J_1^T Y$ and $J_2^T Y$ which are BLUE of $\beta^T \beta$; & then prove by contradiction)

Polynomial Regression

- ① Eq: $y_i = \beta_0 + \beta_1 x_i + \beta_2 x_i^2 + \beta_3 x_i^3 + \epsilon_i$, $\epsilon_i \sim N(0, \sigma^2)$ in variables, but NOT parameters. (*i.e.* $y = \beta_0 + \beta_1 x + \beta_2 x^2 + \epsilon$)
- ② Eq: $y_i = \beta_0 + \beta_1 x_{1i} + \beta_2 x_{2i} + \beta_3 x_{1i} x_{2i} + \beta_4 x_{1i}^2 + \beta_5 x_{2i}^2 + \epsilon_i$, $\epsilon_i \sim N(0, \sigma^2)$

Note, however that both are "linear" in parameters.

& hence can be analyzed in the same framework as Linear Regression.

just their X matrix shall be different. Eq ①: $X = \begin{bmatrix} 1 & x_1 & x_1^2 & x_1^3 \\ 1 & x_2 & x_2^2 & x_2^3 \\ \vdots & \vdots & \vdots & \vdots \\ 1 & x_n & x_n^2 & x_n^3 \end{bmatrix}$

Note: 1) If all x_{ij} s are equal then $X^T X$ shall not be invertible.

$$\text{Eq. ②: } X = \begin{bmatrix} 1 & x_{11} & x_{21} & x_{11}x_{21} & x_{11}^2 & x_{21}^2 \\ 1 & x_{12} & x_{22} & x_{12}x_{22} & x_{12}^2 & x_{22}^2 \\ \vdots & \vdots & \vdots & \vdots & \vdots & \vdots \\ 1 & x_{1n} & x_{2n} & x_{1n}x_{2n} & x_{1n}^2 & x_{2n}^2 \end{bmatrix}$$

- 2) If $\sum |x_{ij}| < \epsilon$ small no. then with higher degrees of x is the columns will become numerically zero!
- 3) At most $(n-1)$ degree polynomial, will allow estimation of $\beta_0, \beta_1, \dots, \beta_{n-1}$. At most $(n-2)$ degree polynomial, will allow estimation of $\underbrace{\dots}_{k+2 \text{ parameters}} \& \sigma^2$.

• $k+2$ parameters
(slopes, intercept & σ^2)

$$\therefore \text{At most: } k+2 \leq n \\ k \leq n-2$$

cofficients (coefficients of the terms)
every other term is zero
the of terms in sum will incorporate with the sum
• Non-zero terms to the final

Model adequacy checking

$E(Y) = X\beta$, β is model parameters

$\tilde{Y} \sim N(X\beta, I\sigma^2)$, σ^2 unknown.

errors are uncorrelated,

errors are normally distributed & errors are independent

coefficient of determination (R^2) = $\frac{\text{Explained variation by model (ESS)}}{\text{Total variation in response (TSS)}} \quad (= \frac{\text{ESS}}{\text{TSS}})$

$$= \frac{\sum (\hat{y}_i - \bar{y})^2}{\sum (y_i - \bar{y})^2} = \frac{\sum (\hat{y}_i - \bar{y})^2}{\sum (y_i - \bar{y})^2} \quad (\because \bar{y} = \bar{y}, \text{ Under OLS})$$

- Now, because

$$\text{TSS} = \text{ESS} + \text{RSS}$$

$$\text{or, } 1 = \frac{\text{ESS}}{\text{TSS}} + \frac{\text{RSS}}{\text{TSS}}$$

$$\text{or, } 1 = R^2 + \frac{\text{RSS}}{\text{TSS}}$$

$$\text{b/w } 0 \text{ to } 1 \quad \therefore R^2 \in [0, 1]$$

$$(R^2 = 1 - \frac{\text{RSS}}{\text{TSS}})$$

- If we increase the number of regressors, (such that $X^T X$ remains invertible)

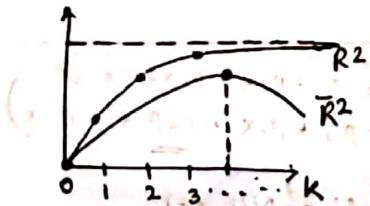
" R^2 cannot decrease" (ie either constant or increase)

[! the very objective of OLS estimation is $\text{Min. RSS} = \frac{\text{RSS}}{\text{TSS}}$ (! TSS constant for given sample; irrespective of prediction or parameter estimates)]

becomes 1 when $k=n-1$

\therefore we cannot use R^2 as a means to decide the adequate no. of regressor variables (k).

Adjusted R^2 (\bar{R}^2) = $1 - \frac{\text{RSS}/df(\text{RSS})}{\text{TSS}/df(\text{TSS})} = 1 - \frac{\text{RSS}/(n-k-1)}{\text{TSS}/(n-1)} \leq R^2 \quad (\because k \geq 0)$



Even though this (so to say) gives us an 'optimal' number of regressors, but how to "pick" which shall be those regressors from a set of several potential explanatory vars? (Umm.. Try out all $P_C k$ possibilities (say, p no. of potential regressors) & compare their \bar{R}^2 ?)

Error Analysis / Residual Analysis

We try to get the best idea about the 'error' from our known 'residual':

* Assumptions made about errors

can't be observed

uncorrelated

\therefore Let's use 'residuals' as their proxy & verify these assumptions on the sample:

$$\hat{e} = \tilde{Y} - \hat{Y} = (I - P_X)\tilde{Y} \sim N(0, \sigma^2(I - P_X))$$

$$\frac{\hat{e}^T \hat{e}}{\sigma^2} \sim \chi_{n-k-1}^2$$

$$e_i = (y_i - \hat{y}_i) \sim N(0, \sigma^2(1-h_{ii})) ; \text{ assuming } H = P_X$$

$$\text{corr}(e_i, e_j) = \text{cov}(y_i - \hat{y}_i, y_j - \hat{y}_j)$$

$$= \sigma^2 \text{corr}(1-h_{ii}, 1-h_{jj}) = \sigma^2 (I-H)_{ij}$$

\therefore estimated errors are correlated in general.

$\Rightarrow e_i$ will follow $N(0, 1)$ when $n \rightarrow \infty$ because $\hat{\sigma}^2 \rightarrow \sigma^2$ with prob. 1

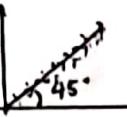
If we plot the histogram then we are expected to get pdf of standard Normal.

Define: $\tau_i = \frac{e_i}{\hat{\sigma}}$
(standardized) $\sqrt{\hat{\sigma}^2(1-h_{ii})}$

we also can do "q-q plot" for r_i 's (ie quantile-quantile plot)

If the errors are normally distributed, then only q-q plot will give a straight line passing through origin with slope 1.

$$q_i = x \text{ such that } P(Z \leq x) = p \\ r_i \text{ such that } \frac{\#\{r_i < x\}}{n} = p$$



or, we could do a "p-p" plot

(probability - probability plot)

• If the plot is different at all, then this indicates Non-normality (ie violation of normality assumption)

Now, ~~error approximation~~

suppose we define $e_{(i)} = y_i - \hat{y}_i$, prediction error of y_i but the same (y_i, x_i) has been removed to get \hat{y}_i .

* If we remove (y_i, x_i) from the dataset, run the OLS, get $\hat{\beta}_2$ and hence $\hat{y}_{(i)}$

$\hat{y}_{(i)}$ = predicted value of y_i based on $(n-1)$ observations.

$\sum_{i=1}^n e_{(i)}^2$ = Predicted Residual error sum of squares

* $\sum_i (y_i - \hat{y}_{(i)})^2 = \text{PRESS}$. (a.k.a "Leave one out" or "Jackknife")

We can show that, $e_{(i)} = \frac{e_i}{1-h_{ii}} \sim N(0, \frac{\sigma^2}{1-h_{ii}})$

standardize $e_{(i)}$

$$\frac{e_{(i)}}{\sqrt{\sigma^2/1-h_{ii}}} = \frac{e_i/1-h_{ii}}{\sqrt{\sigma^2/(1-h_{ii})}} = \frac{e_i}{\sqrt{\sigma^2(1-h_{ii})}}$$

$$\frac{e_i}{\sqrt{\sigma^2(1-h_{ii})}}$$

but not really

(\because here, we have to use ~~to use~~ to get $\hat{\sigma}^2$ based on $(n-1)$ obsⁿ v/s (n) obsⁿ in here)

(previous page's)

$$r_i = \frac{e_i}{\sqrt{\hat{\sigma}^2(1-h_{ii})}}$$

* Q: Do we really have to estimate the model n times

to obtain PRESS?

Anc: Obviously NOT!

The results beside come to our aid

(of course, you need to prove the statement "we can show that, $e_{(i)} = \frac{e_i}{1-h_{ii}} \sim N(0, \frac{\sigma^2}{1-h_{ii}})$ ")

Long matrix based proof!!! (of course, I've skipped it)

and, similarly we can show that:

$$\hat{\sigma}^2 = \frac{(n-k-1)\text{MSRes} - \frac{e_i^2}{1-h_{ii}}}{n-k-2}$$

where MSRes is obviously $\text{SSRes} \left(\frac{n-k-1}{n-k-1} \right)$

The original SSRes using n obsⁿ

(Q) How would you test if i th observation is an outlier?

$$\rightarrow H_0: \mu = 0$$

$$H_1: \mu \neq 0$$

$$y_i - \hat{y}_i = e_{(i)} \sim N(0, \frac{\sigma^2}{1-h_{ii}})$$

$$e_{(i)}$$

$$\frac{e_{(i)}}{\sqrt{\sigma^2/1-h_{ii}}} \sim N(0,1) \text{ under } H_0.$$

$$\Rightarrow \frac{e_{(i)}}{\sqrt{s_{(i)}^2/1-h_{ii}}} \sim t_{n-k-2}$$