

ADVANCED STATISTICS

HW-2

Priya Iddalgi

Rahul Rajput

> Linear Invariance of Regression Coefficients

1) Transformed x_i : $\tilde{x}_i = \frac{x_i - \bar{x}}{s_x}$ (standardized)

Original linear equation: $y_i = \beta_0 + \beta_1 x_i + \epsilon$

Linear eqn on the standardized parameter will be
of the form: $y_i = \alpha_0 + \alpha_1 \tilde{x}_i + \epsilon$

The least sq. estimates can be found as:

$$\alpha_0 = \bar{y}_i - \alpha_1 \bar{x}$$

$$\text{Value of } \alpha_1 = \frac{\sum (y_i - \bar{y}_i)(\tilde{x}_i - \bar{x})}{\sum (\tilde{x}_i - \bar{x})^2}$$

where, \bar{y}_i & \bar{x} are sample means of y_i and \tilde{x}_i

Note that α_0 & α_1 describe changes in y_i with respect to \tilde{x}_i (and not x_i). They can be converted back to their original scale as:

$$\beta_0 = \alpha_0 + \alpha_1 \frac{\bar{x}}{s_x}, \quad \beta_1 = \frac{\alpha_1}{s_x}$$

→ α_0 and α_1 reflect standard deviation of y for a one standard deviation change in x .

2) Relationship b/w \hat{y}_i and \hat{x}_i and α_0, α_1 :

We know that $\hat{y}_i = \alpha_0 + \alpha_1 \hat{x}_i$

$$r = \frac{\sum_{i=1}^{n-1} (\hat{y}_i - \bar{y})(\hat{x}_i - \bar{x})}{\sqrt{\sum_{i=1}^n (\hat{y}_i - \bar{y})^2} \sqrt{\sum_{i=1}^n (\hat{x}_i - \bar{x})^2}} \quad \rightarrow ①$$

$$\text{where } \bar{y} = \frac{\sum_{i=1}^n \hat{y}_i}{n} = \frac{\sum_{i=1}^n (\hat{y}_i - \bar{y}) + \bar{y}}{n} = \bar{y} + \frac{\sum_{i=1}^n (\hat{y}_i - \bar{y})}{n} = \bar{y} + \frac{\sum_{i=1}^n (\hat{y}_i - \bar{y})}{n}$$

$$\text{and } \hat{x}_i = \frac{x_i - \bar{x}}{s_x}$$

Substituting in ①, we have,

$$r = \frac{\sum_{i=1}^{n-1} \left(\frac{(\hat{y}_i - \bar{y})}{s_y} \right) * \left(\frac{(\hat{x}_i - \bar{x})}{s_x} \right)}{\sqrt{\sum_{i=1}^{n-1} \left(\frac{(\hat{y}_i - \bar{y})}{s_y} \right)^2} \sqrt{\sum_{i=1}^{n-1} \left(\frac{(\hat{x}_i - \bar{x})}{s_x} \right)^2}} / n-1$$

$$\text{Also, } \alpha_1 = \frac{\sum_{i=1}^{n-1} (\hat{y}_i - \bar{y})(\hat{x}_i - \bar{x})}{\sum_{i=1}^{n-1} (\hat{x}_i - \bar{x})^2} \quad \rightarrow ②$$

Dividing by $n-1$

$$\text{we have } \boxed{\alpha_1 = r \frac{s_y}{s_x}}$$

3) Sampling variance of α_0 :

$$\text{var}(\alpha_0) = \sigma^2 \left[\frac{1}{n} + \frac{\bar{x}^2}{n-1} \sum_{i=1}^{n-1} (\hat{x}_i)^2 \right]$$

$$\text{var}(\alpha_1) = \frac{\sigma^2}{(n-1)} \sum_{i=1}^{n-1} (\hat{x}_i)^2$$

Sampling covariance b/w slope and intercept:

$$\text{cov}(\alpha_0, \alpha_1) = \frac{\sigma^2 \bar{x}}{(n-1)} \sum_{i=1}^{n-1} (\hat{x}_i)^2$$

σ^2 = common variance of errors ϵ_i

4) Standardizing (or scaling) will help predictability

Yes, we can use standardized slope & intercept-estimates ($\hat{\alpha}_0, \hat{\alpha}_1$) for regression of x_i and y_i

$$\text{Original: } y_i = \beta_0 + \beta_1 x_i + \varepsilon \quad \rightarrow \textcircled{1}$$

$$\text{Standardized: } \tilde{y}_i = \frac{x_i - \bar{x}}{s_x} + \hat{\alpha}_0 + \hat{\alpha}_1 \frac{x_i - \bar{x}}{s_x} + \varepsilon \quad \rightarrow \textcircled{2}$$

$$\bar{x} = \bar{x}_0 + \bar{x}_1 \quad \text{now}$$

$$\text{We know that, } \tilde{x}_i = \frac{x_i - \bar{x}}{s_x} \quad | \times \hat{\alpha}_1$$

Multiplying by $\hat{\alpha}_1$,

$$\tilde{x}_i \hat{\alpha}_1 = \hat{\alpha}_1 \left(\frac{x_i - \bar{x}}{s_x} \right) \quad | \text{ (LHS)} = \text{RHS}$$

Substituting in $\textcircled{2}$, $(\bar{x} - \bar{x}_0) = 0$ now

$$y_i = \hat{\alpha}_0 + \hat{\alpha}_1 \left(\frac{x_i - \bar{x}}{s_x} \right) + \varepsilon$$

$$\textcircled{4} \quad y_i = \hat{\alpha}_0 + \underbrace{\hat{\alpha}_1 \tilde{x}_i}_{s_x} - \frac{\hat{\alpha}_1 \bar{x}}{s_x} + \varepsilon = \text{RHS}$$

Thus,
$$\boxed{\beta_1 = \frac{\hat{\alpha}_1}{s_x}}$$
 coefficient (term with x_i)

$$y_i = \hat{\alpha}_0 - \frac{\hat{\alpha}_1}{s_x} \bar{x} + \beta_1 x_i + \varepsilon \quad | \text{ (LHS) now}$$

then
$$\boxed{\beta_0 = \hat{\alpha}_0 s_x - \hat{\alpha}_1 \bar{x}}$$
 $| \text{ (LHS) now}$

thus, we can use $\hat{\alpha}_0$ and $\hat{\alpha}_1$ by changing back their scale to predict y_i from x_i :

5. If each variable gets multiplied by 100, there should be no change to α_0 and α_1 as they are "standardized estimates".
Consider,

$$\alpha_0 = \bar{y} - \beta_1 \bar{x} \quad \rightarrow ①$$

$$\beta_1 = \frac{\sum (x_i - \bar{x})(y_i - \bar{y})}{\sum (x_i - \bar{x})^2} \quad \rightarrow ②$$

$$\text{Now } x \rightarrow x'_i = 100x_i \quad \& \quad \bar{x}' = 100\bar{x}$$

Substituting,

$$\alpha'_0 = \bar{y} - \beta'_1 \bar{x}' / s_x$$

$$\beta'_1 = \bar{y} - \beta'_1 (100\bar{x}') / 100 s_x$$

$$\beta'_1 = \frac{\sum (100x_i - 100\bar{x})(y_i - \bar{y}) + 100 s_x}{\sum (100x_i - 100\bar{x}_i)}$$

$$\therefore \alpha'_0 = \bar{y} - \beta'_1 \bar{x}' / s_x$$

$$\beta'_1 = \frac{\sum (x_i - \bar{x})(y_i - \bar{y}) s_x}{\sum (x_i - \bar{x})^2}$$

$$\therefore \underline{\alpha'_0 = \alpha_0} \quad \& \quad \underline{\alpha'_1 = \alpha_1}$$

Thus, scaling will not change α_0 or α_1 .

Practical Implications:

i) Improved Interpretability: Standardizing all predictors will make their scale same & effect of each predictor can be compared.

ii) Robustness to Outliers & Improved Stability:
Scaling can make regression estimates

sensitive to outliers. Standardizing can help remove this sensitivity. Standardizing also makes the predictors more stable.

$$\text{OLS} \rightarrow \hat{\beta}_0 = \bar{y} - \hat{\beta}_1 \bar{x} = 0.9$$

$$\text{OLS} \rightarrow \frac{(\hat{\beta}_0 - \bar{y})(\bar{x} - \bar{w})}{\sqrt{(\hat{\beta}_0 - \bar{y})^2}} = 1.9$$

Now we want to make it robust

$$\text{robust} \rightarrow \hat{\beta}_0 = \bar{y} = 1.9$$

$$\text{robust}(\hat{\beta}_0 = 1.9) = \frac{(\hat{\beta}_0 - \bar{y})(\bar{x} - \bar{w})}{\sqrt{(\hat{\beta}_0 - \bar{y})^2}} = 1.9$$

$$\text{robust}(\hat{\beta}_0 = 1.9) = \hat{\beta}_1 = 1.9$$

$$\text{robust}(\hat{\beta}_0 = 1.9)(\hat{\beta}_1 = 1.9) = \frac{(\hat{\beta}_0 - \bar{y})(\bar{x} - \bar{w})}{\sqrt{(\hat{\beta}_0 - \bar{y})^2}} = 1.9$$

so this is robust

and it's robust for flat outliers, that's unhelpful though

using the residuals to distinguish between flat and non-flat outliers is a good idea

but it's better to use a different method

6.)

Examples where x_i & ϵ_i may not be independent:

i) Autocorrelation / Time series data:

In many data collected over time, the errors ϵ_i are correlated with previous errors. Sequences are often correlated to themselves, for e.g., advertising affecting sales for each month.

ii) Measurement Errors:

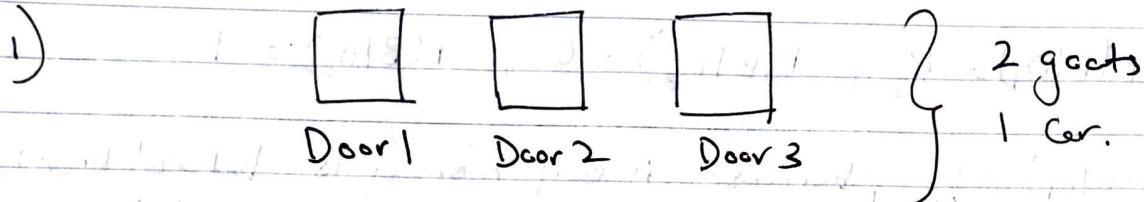
If independent variables (x_i) are measured with error, these errors can be correlated with other errors, violating independence assumption.

iii) Spatial / Cluster Data:

If data is collected from different locations, ϵ_i maybe correlated with x_i . When data are collected from groups of individuals (clusters), the ϵ_i within the same group might be correlated with x_i .

In all these cases, we may get biased and inefficient estimates.

Probability Review



→ At the start, we assume that there is equal probability of the car being behind each door.

$$\Rightarrow P(D_1) = P(D_2) = P(D_3) = \frac{1}{3} \quad [\text{Prob of car being behind a given door}]$$

→ This is Unconditional probability.

$$\rightarrow \text{We choose Door 1, } \Rightarrow P(D_1) = \frac{1}{3}$$

$$\text{but, } P(\sim D_1) = \frac{2}{3}$$

→ Host, has prior knowledge, opens door 3, which has a goat.

Event 2 = 'B', Host opening door 3 to reveal a goat.

→ We want to calculate probability of car being behind Door 1, given the Host opened Door 3 to reveal a goat, i.e., Event 'B'.

$\Rightarrow P(D_1|B) = ?$ Using Bayes Theorem:

$$P(D_1|B) = \frac{P(B|D_1)P(D_1)}{P(B|D_1)P(D_1) + P(B|\sim D_1)P(\sim D_1)} =$$

$$\frac{P(B|D_1)P(D_1)}{P(B|D_1)P(D_1) + P(B|D_2)P(D_2) + P(B|D_3)P(D_3)}$$

$$P(D_1) = P(D_2) = P(D_3) = \frac{1}{3}$$

$$P(B|D_1) = \frac{1}{2}, \quad P(B|D_3) = 0, \quad P(B|D_2) = 1$$

$\rightarrow P(B|D_1) = \frac{1}{2}$, because if money is behind Door 1, then the host could have only opened Door 2 or Door 3. Since Door 3 does not have the car, it is either behind Door 1 or Door 2.

$\rightarrow P(B|D_3) = 0$, because we now know it has a goat.

$\rightarrow P(B|D_2) = 1$, because Prob of opening Door 3 given car is behind Door 2. Host has only 2 options Door 2 & Door 3, since ~~we have chosen Door 1~~ & host cannot open that. Door 3 is opened to ~~releas~~ reveal goat, which only leaves Door 2 as the definite option.

$$\Rightarrow P(D_1|B) = \frac{\frac{1}{2} \times \frac{1}{3}}{\frac{1}{2} \times \frac{1}{3} + 1 \times \frac{1}{3} + 0 \times \frac{1}{3}} = \frac{\frac{1}{6}}{\frac{3}{6}} = \frac{1}{3}$$

\Rightarrow Probability of winning if we stay = $\frac{1}{3}$ OR,

Probability of winning if we switch = $\frac{2}{3}$

MLE Estimator of Poisson Distribution.

$x_1, x_2, x_3, \dots, x_n \sim \text{Poisson}(\lambda)$, where

$$P(x|\lambda) = \frac{e^{-\lambda} \lambda^x}{x!}, \quad x \in \{0, 1, 2, \dots\}$$

1) Mean, Expected Value = $E(x) = \sum_{i=1}^n x \cdot \frac{e^{-\lambda} \lambda^x}{x!}$

$$= e^{-\lambda} \sum_{i=1}^n x \cdot \lambda^x$$

$$= \lambda e^{-\lambda} \sum_{i=1}^n \frac{\lambda^{x-1}}{(x-1)!}$$

$$\text{let } z = x-1$$

$$= \lambda e^{-\lambda} \sum_{i=1}^n \frac{\lambda^z}{z!} \quad [\text{using Taylor Series Expansion}]$$

$$= \lambda e^{-\lambda} \times e^\lambda$$

$$= \underline{\lambda}$$

$$\text{Var}(x) = E(x^2) - (E(x))^2, \quad (E(x))^2 = \lambda^2$$

$$E(x^2) = \sum_{i=1}^n x^2 \frac{e^{-\lambda} \lambda^x}{x!}$$

$$\begin{aligned}
 &= \lambda e^{-\lambda} \sum_{c=1}^{\hat{x}} x \cdot \frac{\lambda^{x-1}}{(x-1)!} \\
 &= \lambda e^{-\lambda} \left[\sum_{c=1}^{\hat{x}} \frac{(x-1) \cdot \lambda^{x-1}}{(x-1)!} + \sum_{c=1}^{\hat{x}} \frac{1}{(x-1)!} \lambda^{x-1} \right] \\
 &= \lambda e^{-\lambda} \left[\lambda \sum_{c=1}^{\hat{x}} \frac{\lambda^{x-2}}{(x-2)!} + \sum_{c=1}^{\hat{x}} \frac{\lambda^{x-1}}{(x-1)!} \right]
 \end{aligned}$$

Substituting $(x-2)$ & $(x-1)$ with j & k & applying Taylor Series expansion,

$$\begin{aligned}
 &= \lambda e^{-\lambda} (\lambda e^\lambda + e^\lambda) \\
 &= \lambda^2 + \lambda \\
 \Rightarrow \text{Var}(x) &= (\lambda^2 + \lambda) - \lambda^2 = \underline{\underline{\lambda}}
 \end{aligned}$$

2) MLE Estimator.

$$P(x|\lambda) = e^{-\lambda} \frac{\lambda^x}{x!}, \quad x \in \{0, 1, 2, 3, \dots\}$$

→ MLE is ~~is~~ a way to assess the distribution of the ~~original~~ data from which the random sample is generated. We try to estimate the value of the parameter of the original data distribution, such that the probability of getting the same values as the observed values is the highest.

Likelihood function

$$f(x_1, x_2, \dots, x_n | \lambda) = \prod_{i=1}^n \frac{e^{-\lambda}}{x_i!} \cdot x_i! \left(p(x_i | \lambda) \cdot x_i \right)$$

$$\Rightarrow \ln f(x_1, x_2, \dots, x_n | \lambda) = \ln \left[\frac{e^{-n\lambda} \times \lambda^{\sum x_i}}{(x_1! x_2! \dots x_n!)} \right]$$

$$\Rightarrow \ln f(x_1, x_2, \dots, x_n | \lambda) = \ln(e^{-n\lambda}) + \left[\ln\left(\lambda^{\sum x_i}\right) - \ln(x_1! x_2! \dots x_n!) \right]$$

$$\Rightarrow \ln f(x_1, x_2, \dots, x_n | \lambda) = -n\lambda + \sum_{i=1}^n x_i \cdot \ln \lambda - \ln(x_1! x_2! \dots x_n!)$$

Now, to find maxima, $\frac{d \ln f(x_1, x_2, \dots, x_n | \lambda)}{d\lambda} = 0$.

$$\Rightarrow -n + \frac{1}{\lambda} \sum_{i=1}^n x_i = 0$$

$$\Rightarrow \lambda = \frac{1}{n} \sum_{i=1}^n x_i$$

$$\Rightarrow \boxed{\lambda = \bar{x}}$$

\therefore MLE for $\lambda = \bar{x}$ or $\hat{\lambda} = \bar{x}$

3) Now, $E(\hat{\lambda}) = E(\bar{x}) = E[\bar{y}] = \lambda$,

Since $E(\hat{\lambda}) = \lambda$, unbiased.

MLE and Inadmissible Estimators.

$X \rightarrow$ binary random variable, 0 or 1.
 prob $\rightarrow p$ for 0, $p-1$ for 1.

$$1) f(x_i | p) = p^{x_i} \cdot (1-p)^{1-x_i}$$

$$\text{Likelihood} = L(p) = \prod_{i=1}^n f(x_i | p) = p^{x_1} (1-p)^{1-x_1} \cdots p^{x_n} (1-p)^{1-x_n}$$

$$\text{or, } L(p) = p^{\sum_{i=1}^n x_i} (1-p)^{n - \sum_{i=1}^n x_i}$$

Taking log,

$$\ln L(p) = \left(\sum_{i=1}^n x_i \right) \ln p + \left(n - \sum_{i=1}^n x_i \right) \ln (1-p)$$

Differentiating & setting derivative to 0 to locate maxima.

$$\frac{d}{dp} \ln L(p) = \frac{\sum_{i=1}^n x_i}{p} - \frac{\left(n - \sum_{i=1}^n x_i \right)}{1-p} = 0$$

$$\Rightarrow \sum_{i=1}^n x_i - p \sum_{i=1}^n x_i - np + p \sum_{i=1}^n x_i = 0$$

$$\Rightarrow \sum_{i=1}^n x_i - np = 0$$

$$\Rightarrow p = \frac{\sum_{i=1}^n x_i}{n}$$

$$\text{MLE of } p \Rightarrow \hat{p} = \frac{\sum_{i=1}^n x_i}{n} = \bar{x}_i$$

$$\begin{aligned}
 2) E[\hat{p}] &= E[\bar{x}] = E\left(\frac{x_1 + x_2 + \dots + x_n}{n}\right) \\
 &= E\left(\frac{\underbrace{E(x_1)}_{\text{Bern}} + E(x_2) + \dots + E(x_n)}{n}\right) \\
 &= \left(\frac{p + p + \dots + p}{n}\right)_n = \frac{np}{n} = p
 \end{aligned}$$

Since $E(\hat{p}) = p$, the MLE estimator, \hat{p} , is unbiased.

3) Expected Square Error.

$$\begin{aligned}
 \text{MSE} &= E(\hat{p} - p)^2 = E(\hat{p}^2) + E(p^2) - 2pE(\hat{p}) = \text{Var}(\hat{p}) + \\
 &= E\left[\left(\frac{x}{n} - p\right)^2\right] = E\left[\frac{x^2}{n^2} - 2\frac{px}{n} + p^2\right]
 \end{aligned}$$

$$= \frac{1}{n^2} E[x^2] - \frac{2p}{n} E[x] + p^2$$

$$\begin{aligned}
 \text{for Bernoulli Distr, } E(x) &= np \quad \& \quad E[x^2] = \text{Var}(x) + E[x]^2 \\
 &= np(1-p) + np^2
 \end{aligned}$$

$$= \frac{1}{n^2} (np(1-p) + np^2) - \frac{2p}{n} \cdot np + p^2$$

$$= \boxed{\frac{p(1-p)}{n}}$$

4) from the prev part, we get the expected MSE (for a Bernoulli distribution as,

$$\frac{p(1-p)}{n} \rightarrow \text{***}$$

\Rightarrow if $p = 0$ or 1 , the MSE = 0.

i.e., in a distribution where there can only be success or failure, the MSE = 0.

Now if $p \in \left[\frac{1}{4}, \frac{3}{4}\right]$, $n = 3$

$$\Rightarrow \text{if } p = 1/4, \text{ MSE} = \frac{1/4 \times 3/4}{3} = 1/16$$

$$\text{if } p = 3/4, \text{ MSE} = \frac{(1-3/4) \times 1/4}{3} = \frac{1}{16}$$

for $p \in [1/4, 3/4]$, since $R(\hat{p}, p)$ is the same for

both values of p & $E(\hat{p}) = p$, ~~the MSE~~ there is

~~no~~ minima for MSE among all values of \hat{p} .

\therefore we can claim that the MLE estimator is inadmissible.