STA 108 Notes - J. Jiang
Dylan M Ang
February 14, 2022

# Contents

# 1 Simple Linear Regression

$$Y_i = \beta_0 + \beta_1 x_i + \epsilon_i \quad i = 1 \ldots n \quad (1)$$

If assumptions hold true,

- $Y_i$ is normally distributed

$$E(Y_i) = \beta_0 + \beta_1 x_i \quad (2)$$
$$Var(Y_i) = \sigma^2 \quad (3)$$

- Mean: $\beta_0 + \beta_1 x_i$

- Variance: $\sigma^2$

## Assumptions

- $\epsilon_1 \ldots \epsilon_n$

- $E(\epsilon_i) = 0, var(\epsilon_i) = 0$, where $\sigma^2$ is an unknown constant.

- $\epsilon_i$ is normal. (normality assumption)

## 1.1 Least Squares Estimate

$$\Sigma_{i=1}^n (Y_i - \beta_0 - \beta_1 x_i)^2 \quad (4)$$

Take the first order derivative with respect to $\beta_0, \beta_1$ to <mark>minimize</mark> equation (4) to find optimal $\beta_0, \beta_1$.

**LS estimators**

$$\hat{\beta}_1 = \frac{\Sigma_{i=1}^n (x_i - \bar{x})(Y_i - \bar{Y})}{\Sigma_{i=1}^n (x_i - \bar{x})^2} \quad (5)$$

$$\hat{\beta}_0 = \bar{Y} - \hat{\beta}_1 \bar{x} \quad (6)$$

- $\bar{x} = \frac{1}{n}\Sigma_{i=1}^n x_i$ (sample mean of the $x_i$'s)

- $\bar{Y} = \frac{1}{n}\Sigma_{i=1}^n Y_i$ (sample mean of the $Y_i$'s)

- Regression Line: $y = \hat{\beta}_0 + \hat{\beta}_1 x$

**Properties of LS Estimators**

- $E(\hat{\beta}_0) = \beta_0, E(\hat{\beta}_1) = \beta_1$. The average of many sample beta values will approach the true beta values.

<mark>Fitted</mark> (or predicted) values are estimates. The fitted value for $Y_i$ is $\hat{Y}_i = \hat{\beta}_0 + \hat{\beta}_1 x_i$
$\hat{Y}$ is an <mark>unbiased estimator</mark> of $E(Y) = \beta_0 + \beta_1 x$ so $E(\hat{Y}) = E(Y)$

## 1.2 Residuals

<mark>Residuals</mark> : $\hat{\epsilon}_i = Y_i - \hat{Y}_i, i = 1 \ldots n$
Properties of Residuals

- $\Sigma_{i=1}^n \hat{\epsilon}_i = 0$

- The residuals are not independent.

- If one residual is positive, another residual has to compensate.

## 1.3 Variance

Estimation of $\sigma^2$, the variance of the errors (which is the same as the variance of $Y_i$)

$$s^2 = \frac{1}{n-2}\Sigma_{i=1}^n (Y_i - \hat{Y}_i)^2 \quad (7)$$

where $\hat{Y}_i$ is the estimate of $E(Y_i)$.

**Notes**

- $\hat{Y}_i$ is an estimator of $E(Y_i) = \beta_0 + \beta_1 x_i$ in which two parameters are estimated ($\beta_0$ and $\beta_1$) $\implies$ 2 degrees of freedoms are subtracted.

- $E(s^2) = \sigma^2$

When errors are normally distributed, the LS estimators of $\beta_0, \beta_1$ is equal to the MLEs (Maximum Likelihood Estimators) of $\beta_0, \beta_1$, but the MLE of $\sigma^2, \hat{\sigma}^2$, is different from $s^2$

$s^2$ is just (7)

$$\hat{\sigma}^2 = \frac{1}{n}\Sigma_{i=1}^n \hat{\epsilon}_i^2 \quad (8)$$

# 2 Inference in regression and correlation analysis

## 2.1 Inference about $\beta_1$

**For testing $\beta_1$**
$H_0 : \beta_1 = \beta_{10}, \beta_{10}$ is a given value such as 0.
$H_a : \beta_1 \neq \beta_{10}, \beta_1 > \beta_{10},$ or $\beta_1 < \beta_{10}$
<mark>Test statistic</mark> : A statistic whose distribution is known <mark>under the null hypothesis.</mark>

$$t = \frac{\hat{\beta}_1 - \beta_{10}}{s.e.(\hat{\beta}_1)} \quad (9)$$

where $\hat{\beta}_1$ is the LS estimate of $\beta_1$, and

$$s.e.(\hat{\beta}_1) = \sqrt{\frac{MSE}{\Sigma_i (x_i - \bar{x})^2}} \quad (10)$$

$$MSE = s^2 \quad (11)$$

If normal, $T \sim t_{n-1}$

$$T = \frac{\hat{\beta}_1 - \beta_1}{s.e.(\hat{\beta}_1)} \quad (12)$$

Therefore under the $H_0 : \beta_1 = \beta_{10}, t \sim t_{n-2}$
**Decision Rules**

$$H_1 : \beta_1 \neq \beta_{10}, reject \ H_0 \ if \ |t| > t_{n-2,\alpha/2}$$
$$H_1 : \beta_1 > \beta_{10}, reject \ H_0 \ if \ |t| > t_{n-2;\alpha}$$
$$H_1 : \beta_1 < \beta_{10}, reject \ H_0 \ if \ |t| < -t_{n-2;\alpha}$$

Alternatively, Reject $H_0$ if the p-value of t is $\leq \alpha$
**Error**

- <mark>Type I</mark> : Reject $H_0$ when it is true.

- <mark>Type II</mark> : Fail to reject $H_0$ when it is false.

**Level of Significance $\alpha$**
<mark>$\alpha$</mark> is the upper bound for the probability of Type I error.
**P-value**

<mark>p-value</mark> is the observed level of significance: the actual probability that the test statistic is as extreme as observed given $H_0$ is true.
**Power**
<mark>Power</mark> is the probability of rejecting $H_0$ when the alternative holds at a given value.
If $\beta_{10} = 0, \beta_1 = 1, s.d.(\hat{\beta}_1) = 0.5$, we have $\delta = \frac{1}{0.5} = 2$ Let $\alpha = 0.05$. From table B.5 we find the power is 0.48.
**Confidence interval for $\beta_k$**
Assuming normality, a $100(1 - \alpha)\%$ c.i. for $\beta_k$ is

$$\hat{\beta}_k \pm t_{n-2}(1 - \frac{alpha}{2} * s.e.(\hat{\beta})) \quad (13)$$
$$k = 0, 1 \quad (14)$$

where $s.e.(\hat{\beta}_1)$ can be found with eq(10) and $s.e.(\hat{\beta}_0)$ can be found with eq(15).

## 2.2 Inference about $\beta_0$

$$s.e.(\beta_0) = \sqrt{mse(\frac{1}{n} + \frac{\bar{x}^2}{\Sigma_{i=1}^n (x_i - \bar{x})^2})} \quad (15)$$

<mark>Confidence intervals</mark> for $\beta_0$ can be found with (13)

## 2.3 Inference about $\hat{Y}$

<mark>Confidence Interval for $E(Y) = \beta_0 + \beta_1 x$</mark>

$$\hat{Y} \pm t_{n-2}(1 - \frac{alpha}{2}) * s.e.(\hat{Y}) \quad (16)$$
$$s.e.(\hat{Y}) = \sqrt{MSE(\frac{1}{n} + \frac{(x - \bar{x})^2}{\Sigma_{i=1}^n (x_i - \bar{x})^2})} \quad (17)$$

## 2.4 Prediction interval for $\hat{Y}$

A $100(1 - \alpha)\%$ prediction interval for $Y = E(Y) + \epsilon = \beta_0 + \beta_1 x + \epsilon$, where Y is the future observation and $\epsilon$ is the new error:

$$\hat{Y} \pm t_{n-2}(1 - \frac{\alpha}{2}) * p.s.e.(\hat{Y}) \quad (18)$$
$$p.s.e.(\hat{Y}) = \sqrt{MSE(1 + \frac{1}{n} + \frac{(x - \bar{x})^2}{\Sigma_{i=1}^n (x_i - \bar{x})^2})} \quad (19)$$

Where p.s.e. is the percent standard error.
The 1 in the p.s.e is because the variance of $\epsilon = \sigma^2$. If $var(\epsilon) = \frac{\sigma^2}{2}$ change the 1 to $\frac{1}{2}$.

## 2.5 ANOVA and F-test

$$SSTO = SSR + SSE \quad (20)$$
$$= \Sigma_{i=1}^n (Y_i - \bar{Y})^2 \quad (21)$$
$$SSR = \Sigma_{i=1}^n (\hat{Y}_i - \bar{Y})^2 \quad (22)$$
$$SSE = \Sigma_{i=1}^n (Y_i - \hat{Y})^2 \quad (23)$$
$$(24)$$

Sum of Squares of Regression (SSR) explains the variability in $Y$ due to the regression model compared to the baseline model. Sum of Squares of Errors (SSE) is the remaining unexplained variability of Y found from SSTO - SSR.
**Degrees of Freedom**

$$SSRdf = 1$$
$$SSEdf = n - 2$$
$$SSTOdf = n - 2 + 1 = n - 1$$

**Mean Squares**
Mean squares is SS divided by its degrees of freedom.

$$MSR = \frac{SSR}{1} \quad (25)$$
$$MSE = \frac{SSE}{n - 2} \quad (26)$$
$$(27)$$

**F-Statistic**

$$F = \frac{MSR}{MSE} = \frac{SSR*(n-2)}{SSE} \qquad (28)$$

ANOVA table : Analysis of variance.
The distribution of F under the null hypothesis $H_0 : \beta_1 = 0$ is $F_{1,n-2}$.

| Source | SS | df | MS | F |
|---|---|---|---|---|
| Regression | SSR | 1 | MSR | F |
| Error | SSE | n-2 | MSE | |
| Total | SSTO | n-1 | | |

## 2.6 Inference about $\rho$

$R^2$ : a measure of goodness of fit, which is the proportion of variation in $Y$ explained by the regression (i.e. by x).

$$R^2 = \frac{SSR}{SSTO} = 1 - \frac{SSE}{SSTO} \qquad (29)$$

Coefficient of correlation:

$$r = \pm\sqrt{R^2} = \begin{cases} +\sqrt{R^2} & if \; \hat{\beta}_1 > 0 \\ -\sqrt{R^2} & if \; \hat{\beta}_1 < 0 \end{cases} \qquad (30)$$

$$r = \frac{\Sigma_i(Y_i - \bar{Y})(x_i - \bar{x})}{\sqrt{\Sigma_i(Y_i - \bar{Y})^2 \Sigma_i(x_i - \bar{x})^2}} \qquad (31)$$

**Properties of $R^2$ and $r$**

- $0 \le R^2 \le 1 \quad -1 \le r \le 1$

- $R^2 \approx 1$ or $r \approx \pm 1$, if there is a strong linear association between $x$ and $Y$.

- $R^2 \approx 0$, or $r \approx 0$, if there is a weak or no linear association between $x$ and $Y$.

- Both $R^2$ and $r$ are measures of linear association only.

**Covariance and correlation between two random variables**

$$cov(X,Y) = E\{(X - \mu_X)(Y - \mu_Y)\} \qquad (32)$$
$$= E(XY) - E(X)E(Y) \qquad (33)$$
$$cor(X,Y) = \frac{cov(X,Y)}{sd(X)sd(Y)} \qquad (34)$$

where $\mu_X = E(X), sd(X) = \sqrt{var(X)}$,etc.
Special case: $(X,Y)$ has a bivariate normal distribution.
**Testing for $\rho$**
Assume that the bivariate normal distribution holds for $(X,Y)$.
$H_0 : \rho = 0$
$H_a : \rho \ne 0 (\; or \; \rho > 0 \; or \; \rho < 0)$

$$t^* = \frac{r\sqrt{n-2}}{\sqrt{1-r^2}} \sim t_{n-2} \text{ under } H_0 \qquad (35)$$

# 3 Diagnostics

The goal of diagnostics is to examine the departures from the simple linear regression model with normal errors. Typical departures and corresponding diagnostic plots/tests are:

- The regression is not linear - residual plots(residual against the predictor variable, or against the fitted values), lack of fit test.

- The error terms are not normally distributed - histogram, boxplot/dot plot of residuals, normal probability plot (aka QQ plot), Shapiro-Wilk's test, correlation test for normality.

- The error terms do not have constant variance - residual plots, Brown - Forsythe (BF) test.

- The error terms are not independent - residual against time.

- The model fits all but one or a few outlier observations - (semistudentized) residual plots, box plots, dot plots, stem and leaf plots.

- Some important predictors are missing - residual plots (residual against other possibly important predictors).

## 3.1 Residual Plots

Residuals can be used to check whether

- The regression function is not linear.

- The variance of the errors is not constant.

- The errors are not independent.

- Outliers

- The errors are not normal.

- Some important predictors are missing.

**Scatter Plot**

- Check linearity - residuals normally disributed.

- Check constant variance - residuals are random and dont follow a cone pattern.

**Box Plot and Dot Plot**

- Normality - residuals should be centered and symmetric about 0.

**Normality probability plot - QQ Plot**

- QQ plot is linear $\implies$ normal residuals.

- QQ plot is nonlinear $\implies$ non normal residuals.

## 3.2 Diagnostic Tests

**Shapiro Wilk's test**
$H_0 \; data \sim N()$
$H_a$ : data not normal.
$p - val \le \alpha$ reject normality assumption.
**Correlation test for normality**
Step 1. Compute the coefficient of correlation between the ordered residuals and their expected values. The latter are given by

$$\sqrt{MSE}z(\frac{k-0.375}{n+0.25}), \quad k = 1, \ldots, n \qquad (36)$$

where $z(p)$ is the pth quantile of the standard normal distribution, that is, $P[Z \le z(p)] = p$, where $Z$ has the standard normal distribution.
Step 2. Compare the coefficient of correlation on I with the critical value from Table B.6, if the coefficient of correlation exceeds the critical value, accept the normality assumption.
**BF test for constant variance**
1. Divide the residuals into two parts according to residual pattern (or no pattern)
Let $\hat{\epsilon}_{i1} = 1, \ldots, n_1$ be the residuals for the first part, and $\hat{\epsilon}_{i2}, i = 1, \ldots, n_2$ be the residuals for the second part, where $n_1 + n_2 = n$.

Compute $m(\hat{\epsilon}_1) = \; median \; of \; \hat{\epsilon}_{i1}, i = 1, \ldots, n_1$ and $m(\hat{\epsilon}_2)$.
2. Compute $d_{i1} = |\hat{\epsilon}_{i1} - m(\hat{\epsilon}_1)|, i = 1 \ldots n_1$ and $d_{i2} = |\hat{\epsilon}_{i2} - m(\hat{\epsilon}_2)|, i = 1 \ldots n_2$
3. Compute t score.

$$t_{BF} = \frac{\bar{d}_1 - \bar{d}_2}{s\sqrt{n_1^{-1} + n_2^{-1}}} \qquad (37)$$

$$s^2 = \frac{\Sigma_{i=1}^{n_1}(d_{i1} - \bar{d}_1)^2 + \Sigma_{i=1}^{n_2}(d_{i2} - \bar{d}_2)^2}{n-2} \qquad (38)$$

4. Test $H_0 : \sigma_1^2 = \sigma_2^2$ vs $H_a : \sigma_1^2 \ne \sigma_2^2$
$t_{BF} \sim t_{n-2}$ under $H_0$. Given $\alpha$, use the critical value (or p-value) to test $H_0$.
**F-test for lack of fit**
Regression model: $Y_{ij} = \beta_0 + \beta_1 x_j + \epsilon_{ij}, j = 1 \ldots c, i = 1 \ldots n_j$ where $x_j$ is the $j$th value of $x$, $c$ is the number of different $x$ values, and $Y_{ij}, i = 1 \ldots n_j$ are the Y values corresponding to the same $x_j$.
Full model: $Y_{ij} = \mu_j + \epsilon_{ij}, j = 1 \ldots c, i = 1 \ldots n_j$
F-statistic:

$$F = \frac{SSE(R) - SSE(F)}{df_R - df_F}\{\frac{SSE(F)}{df_F}\}^{-1} \qquad (39)$$

where

$$SSE(R) = \Sigma_j \Sigma_i(Y_{ij} - \hat{Y}_{ij})^2 \qquad (40)$$
$$SSE(F) = \Sigma_j \Sigma_i(Y_{ij} - \hat{\mu}_j)^2 \qquad (41)$$

with $\hat{Y}_{ij} = \hat{\beta}_0 + \hat{\beta}_1 x_j$ and $\hat{\mu}_j = \bar{Y}_j - n_j^{-1}\Sigma_{i=1}^{n_j} Y_{ij}, df_R = n-2$
with $n = \Sigma_{j=1}^c n_j$ and $df_F = n - c$.
Under $H_0$ : The assumed model is correct, $F \sim F_{c-2,n-c}$.

## 3.3 Remedial Measures

**Transformation of x**: for nonlinear association.
**Transformation of Y**: for nonnormality/unequal variance.
**Box Cox transformation**
This is a collection of transformations depending on a "tuning parameter", $\lambda$.

$$Y_i' = \begin{cases} K_1(Y_i^\lambda - 1), & \lambda \ne 0 \\ K_2 log(Y_i), & \lambda = 0 \end{cases} \qquad (42)$$

where $K_1, K_2$ are two numbers computed from the data.

$$K_2 = (Y_1 Y_2 \ldots Y_n)^{\frac{1}{n}} = e^{\overline{\log Y}} \qquad (43)$$

$$K_1 = \frac{1}{\lambda K_2^{\lambda-1}} \qquad (44)$$