

Bruce M.
Boghosian

Properties of
Linear Model
estimators

Estimating $\hat{\sigma}^2$

Covariance
and
Correlation

Summary

Regression

Covariance and Correlation

Bruce M. Boghosian



Tufts
UNIVERSITY

School of Arts
and Sciences

Department of Mathematics

Tufts University

Bruce M.
Boghossian

Properties of
Linear Model
estimators

Estimating $\hat{\sigma}^2$

Covariance
and
Correlation

Summary

1 Properties of Linear Model estimators

2 Estimating $\hat{\sigma}^2$

3 Covariance and Correlation

4 Summary

Bruce M.
Boghossian

Properties of
Linear Model
estimators

Estimating $\hat{\sigma}^2$

Covariance
and
Correlation

Summary

- Let Y_1, \dots, Y_n be any set of independent random variables with means μ_1, \dots, μ_n and variances $\sigma_1^2, \dots, \sigma_n^2$, respectively. Let a_1, \dots, a_n be any set of constants. Then $Y = a_1 Y_1 + \dots + a_n Y_n$ is normally distributed with mean $\mu = \sum_i^n a_i \mu_i$ and variance $\sigma^2 = \sum_i^n a_i^2 \sigma_i^2$.
- Proof is straightforward using moment-generating functions.

Bruce M.
Boghossian

Properties of
Linear Model
estimators

Estimating $\hat{\sigma}^2$

Covariance
and
Correlation

Summary

- We are given data $(x_1, Y_1), \dots, (x_n, Y_n)$.
- Linear Model supposes that $f_{Y|x}(y)$ is normal for all x .
- Standard deviation σ assumed independent of x .
- Means are collinear, $E(Y|x) = \beta_0 + \beta_1 x$.
- Results of maximum likelihood estimation

$$\hat{\beta}_0 = \frac{(\sum_{i=1}^n x_i^2) (\sum_{i=1}^n Y_i) - (\sum_{i=1}^n x_i) (\sum_{i=1}^n x_i Y_i)}{n (\sum_{i=1}^n x_i^2) - (\sum_{i=1}^n x_i)^2}$$

$$\hat{\beta}_1 = \frac{n (\sum_{i=1}^n x_i Y_i) - (\sum_{i=1}^n x_i) (\sum_{i=1}^n Y_i)}{n (\sum_{i=1}^n x_i^2) - (\sum_{i=1}^n x_i)^2}$$

- Same form obtained from geometric approach to fitting.

Bruce M.
Boghossian

Properties of
Linear Model
estimators

Estimating $\hat{\sigma}^2$

Covariance
and
Correlation

Summary

- **Thm.:** $\hat{\beta}_0$ and $\hat{\beta}_1$ are both normally distributed.
- **Thm.:** $\hat{\beta}_0$ and $\hat{\beta}_1$ are both unbiased,

$$E(\hat{\beta}_0) = \beta_0$$

$$E(\hat{\beta}_1) = \beta_1$$

- **Thm.:** The variances of the estimators are

$$\text{Var}(\hat{\beta}_0) = \frac{\sigma^2 \sum_i^n x_i^2}{n \sum_i^n (x_i - \bar{x})^2}$$

$$\text{Var}(\hat{\beta}_1) = \frac{\sigma^2}{n \sum_i^n (x_i - \bar{x})^2}$$

Bruce M.
Boghossian

Properties of
Linear Model
estimators

Estimating $\hat{\sigma}^2$

Covariance
and
Correlation

Summary

- **Pf. (for $\hat{\beta}_1$):** Note that

$$\begin{aligned}\hat{\beta}_1 &= \frac{n \sum_i^n x_i Y_i - (\sum_i^n x_i) (\sum_i^n Y_i)}{n (\sum_i^n x_i^2) - (\sum_i^n x_i)^2} = \frac{\sum_i^n x_i Y_i - \left(\frac{1}{n} \sum_i^n x_i\right) (\sum_i^n Y_i)}{\left(\sum_i^n x_i^2\right) - n \left(\frac{1}{n} \sum_i^n x_i\right) \left(\frac{1}{n} \sum_i^n Y_i\right)} \\ &= \frac{\sum_i^n x_i Y_i - \bar{x} (\sum_i^n Y_i)}{(\sum_i^n x_i^2) - n \bar{x}^2} = \frac{\sum_i^n (x_i - \bar{x}) Y_i}{(\sum_i^n x_i^2) - n \bar{x}^2}\end{aligned}$$

- This is a linear combination of normally distributed r.v.s, and thus normally distributed.

Bruce M.
Boghossian

Properties of
Linear Model
estimators

Estimating $\hat{\sigma}^2$

Covariance
and
Correlation

Summary

- **Pf. (for $\hat{\beta}_1$):** Using the same form for $\hat{\beta}_1$ used above,

$$\begin{aligned} E(\hat{\beta}_1) &= E\left(\frac{\sum_i^n (x_i - \bar{x}) Y_i}{(\sum_i^n x_i^2) - n\bar{x}^2}\right) \\ &= \frac{\sum_i^n (x_i - \bar{x}) E(Y_i)}{(\sum_i^n x_i^2) - n\bar{x}^2} \\ &= \frac{\sum_i^n (x_i - \bar{x}) (\beta_0 + \beta_1 x_i)}{(\sum_i^n x_i^2) - n\bar{x}^2} \\ &= \beta_1. \end{aligned}$$

Bruce M.
Boghossian

Properties of
Linear Model
estimators

Estimating σ^2

Covariance
and
Correlation

Summary

- **Pf. (for $\hat{\beta}_1$):** Using the same form for $\hat{\beta}_1$ used above,

$$\begin{aligned}\text{Var}(\hat{\beta}_1) &= \text{Var}\left(\frac{\sum_i^n (x_i - \bar{x}) Y_i}{(\sum_i^n x_i^2) - n\bar{x}^2}\right) \\ &= \sum_i^n \left(\frac{(x_i - \bar{x})}{(\sum_i^n x_i^2) - n\bar{x}^2}\right)^2 \text{Var}(Y_i) \\ &= \frac{\sigma^2}{n \sum_i^n (x_i - \bar{x})^2}\end{aligned}$$

Bruce M.
Boghossian

Properties of
Linear Model
estimators

Estimating $\hat{\sigma}^2$

Covariance
and
Correlation

Summary

- **Thm.:** Let $(x_1, Y_1), \dots, (x_n, Y_n)$ satisfy the assumptions of the Linear Model. Then
 - $\hat{\beta}_1, \bar{Y}$, and $\hat{\sigma}^2$ are mutually independent.
 - $\frac{n}{\sigma^2} \hat{\sigma}^2$ has a χ^2 distribution with $n - 2$ degrees of freedom.
- **Corr.:** Let $\hat{\sigma}^2$ be MLE for σ^2 in the Linear Model. Then
 - $\left(\frac{n}{n-2}\right) \hat{\sigma}^2$ is an unbiased estimator for σ^2 .
 - The random variables \hat{Y}^2 and $\hat{\sigma}^2$ are independent.
- Proofs relegated to appendix of Chapter 11 in Larsen & Marx text.

Bruce M.
Boghossian

Properties of
Linear Model
estimators

Estimating σ^2

Covariance
and
Correlation

Summary

- Calculating $\sum_i^n (y_i - \hat{y}_i)^2 = \sum_i^n (y_i - \hat{\beta}_0 - \hat{\beta}_1 x_i)^2$

$$= \sum_i^n (y_i - \bar{y})^2 - \hat{\beta}_1^2 \sum_i^n (x_i - \bar{x})^2$$

$$= \sum_i^n y_i^2 - \frac{1}{n} \sum_i^n y_i - \frac{[\sum_i^n x_i y_i - \frac{1}{n} (\sum_i^n x_i) (\sum_i^n y_i)]^2}{(\sum_i^n x_i^2) - \frac{1}{n} (\sum_i^n x_i)^2}$$

$$= \sum_i^n y_i^2 - \hat{\beta}_0 \sum_i^n y_i - \hat{\beta}_1 \sum_i^n x_i y_i.$$

- Note that the above are true only for the actual estimators $\hat{\beta}_0$ and $\hat{\beta}_1$, calculated from the data.
- The above are not true if $\hat{\beta}_0$ and $\hat{\beta}_1$ are arbitrary.

Bruce M.
Boghossian

Properties of
Linear Model
estimators

Estimating $\hat{\sigma}^2$

Covariance
and
Correlation

Summary

- **Thm.:** Let $(x_1, Y_1), \dots, (x_n, Y_n)$ satisfy the assumptions of the Linear Model, and let

$$s^2 = \frac{1}{n-2} \sum_i^n (Y_i - \hat{\beta}_0 - \hat{\beta}_1 x_i)^2 = \frac{1}{n-2} \sum_i^n (Y_i - \hat{Y}_i)^2$$

- Then the following has a Student T distribution with $n - 2$ df.

$$T_{n-2} = \frac{\hat{\beta}_1 - \beta_1}{s / \sqrt{\sum_i^n (x_i - \bar{x})^2}}$$

- **Pf.:** Note that the following is normally distributed

$$\hat{\beta}_1 - \beta_1 = \frac{\sum_i^n (x_i - \bar{x})(Y_i - \hat{y}_i)}{\sum_i^n (x_i - \bar{x})^2}$$

- Note that, since $E(Y_i) = \hat{y}_i$, it follows $E(\hat{\beta}_1 - \beta_1) = 0$.
- Then note that

$$\text{Var}(\hat{\beta}_1 - \beta_1) = \frac{\sum_i^n (x_i - \bar{x})^2 \text{Var}(Y_i)}{\left[\sum_i^n (x_i - \bar{x})^2\right]^2} = \frac{\sigma^2}{\sum_i^n (x_i - \bar{x})^2}$$

- Hence the following is distributed as a standard normal

$$Z = \frac{\hat{\beta}_1 - \beta_1}{\frac{\sigma}{\sqrt{\sum_i^n (x_i - \bar{x})^2}}}$$

- **Pf. (continued):** We know that

$$Z = \frac{\hat{\beta}_1 - \beta_1}{\frac{\sigma}{\sqrt{\sum_i^n (x_i - \bar{x})^2}}}$$

is distributed like a standard normal.

- We also know that

$$\frac{n}{\sigma^2} \hat{\sigma}^2 = \frac{n-2}{\sigma^2} S^2$$

is distributed as a χ^2 r.v. with $n-2$ df.

- Hence the following is T distributed with $n-2$ df

$$T_{n-2} = \frac{Z}{\sqrt{\frac{\left(\frac{n-2}{\sigma^2} S^2\right)}{n-2}}} = \frac{\hat{\beta}_1 - \beta_1}{S / \sqrt{\sum_i^n (x_i - \bar{x})^2}} \quad \square$$

- We are given data $(x_1, Y_1), \dots, (x_n, Y_n)$ for Linear Model.
- Calculate the statistic

$$t = \frac{\hat{\beta}_1 - \beta'_1}{s / \sqrt{\sum_i^n (x_i - \bar{x})^2}}$$

where

$$s^2 = \frac{1}{n-2} \sum_i^n (y_i - \hat{\beta}_0 - \hat{\beta}_1 x_i)^2 = \frac{1}{n-2} \sum_i^n (y_i - \hat{y}_i)^2$$

- Hypothesis tests involving $\hat{\beta}_1$:
 - To test $H_0 : \beta_1 = \beta'_1$ versus $H_1 : \beta_1 > \beta'_1$ at α level of significance, reject H_0 if $t \geq +t_{\alpha, n-2}$.
 - To test $H_0 : \beta_1 = \beta'_1$ versus $H_1 : \beta_1 < \beta'_1$ at α level of significance, reject H_0 if $t \leq -t_{\alpha, n-2}$.
 - To test $H_0 : \beta_1 = \beta'_1$ versus $H_1 : \beta_1 \neq \beta'_1$ at α level of significance, reject H_0 if $t \geq +t_{\alpha/2, n-2}$ or $t \leq -t_{\alpha/2, n-2}$.

- We are given data $(x_1, Y_1), \dots, (x_n, Y_n)$ for Linear Model.
- As with hypothesis testing, calculate the statistic

$$t = \frac{\hat{\beta}_1 - \beta_1'}{s / \sqrt{\sum_i^n (x_i - \bar{x})^2}}$$

- Then

$$\left[\hat{\beta}_1 - \frac{t_{\alpha/2, n-2} s}{\sqrt{\sum_i^n (x_i - \bar{x})^2}}, \hat{\beta}_1 + \frac{t_{\alpha/2, n-2} s}{\sqrt{\sum_i^n (x_i - \bar{x})^2}} \right]$$

is a $100(1 - \alpha)\%$ confidence interval for β_1 .

- It is also possible to find confidence intervals for β_0 , but usually not as important.

Bruce M.
Boghossian

Properties of
Linear Model
estimators

Estimating $\hat{\sigma}^2$

Covariance
and
Correlation

Summary

- Recall the *covariance* of random variables X and Y

$$\text{Cov}(X, Y) = E[(X - \mu_X)(Y - \mu_Y)]$$

- The covariance depends on the units of the variables.
- Make independent of the units by dividing by σ_X and σ_Y , to obtain the *correlation coefficient*,

$$\rho(X, Y) = \frac{\text{Cov}(X, Y)}{\sigma_X \sigma_Y} = E \left[\left(\frac{X - \mu_X}{\sigma_X} \right) \left(\frac{Y - \mu_Y}{\sigma_Y} \right) \right]$$

- This also has the effect of ensuring $\rho(X, Y) \in [-1, +1]$.

Proof that $\rho(X, Y) \in [-1, +1]$

Bruce M.
Boghossian

Properties of
Linear Model
estimators

Estimating $\hat{\sigma}^2$

Covariance
and
Correlation

Summary

- Define the standardized r.v.s, $X^* = \frac{X - \mu_X}{\sigma_X}$ and $Y^* = \frac{Y - \mu_Y}{\sigma_Y}$
- Hence $E(X^*) = E(Y^*) = 0$ and $\text{Var}(X^*) = \text{Var}(Y^*) = 1$.
- Now consider

$$\begin{aligned} 0 \leq \text{Var}(X^* \pm Y^*) &= E[(X^*)^2] + 2E(X^*Y^*) + E[(Y^*)^2] \\ &= \text{Var}(X^*) \pm 2\text{Cov}(X^*, Y^*) + \text{Var}(Y^*) \\ &= 2 \pm 2\rho(X, Y). \end{aligned}$$

- It follows that $-1 \leq \rho(X, Y) \leq +1$. □

- We proved above earlier using Cauchy-Schwarz inequality.

Bruce M.
Boghossian

Properties of
Linear Model
estimators

Estimating $\hat{\sigma}^2$

Covariance
and
Correlation

Summary

- We have defined the correlation

$$\rho(X, Y) = \frac{E(XY) - E(X)E(Y)}{\sqrt{\text{Var}(X)}\sqrt{\text{Var}(Y)}}$$

- So define the *sample correlation coefficient*,

$$R = \frac{\frac{1}{n} \sum_i^n X_i Y_i - \bar{X} \bar{Y}}{\sqrt{\frac{1}{n} \sum_i^n (X_i - \bar{X})^2} \sqrt{\frac{1}{n} \sum_i^n (Y_i - \bar{Y})^2}}$$

A simple relationship for r^2

Bruce M.
Boghossian

Properties of
Linear Model
estimators

Estimating $\hat{\sigma}^2$

Covariance
and
Correlation

Summary

- Now suppose that we have data $(x_1, y_1), \dots, (x_n, y_n)$.
- Define the *coefficient of determination*

$$r^2 = \frac{\sum_i^n (y_i - \bar{y})^2 - \sum_i^n (y_i - \hat{\beta}_0 - \hat{\beta}_1 x_i)^2}{\sum_i^n (y_i - \bar{y})^2}.$$

- Simple interpretation
 - $\sum_i^n (y_i - \bar{y})^2$ is the *total variability* in y .
 - $\sum_i^n (y_i - \hat{\beta}_0 - \hat{\beta}_1 x_i)^2$ is the variability that can not be explained by linear regression.
 - The numerator of r^2 is the variability that *can* be explained by linear regression.
 - The quantity r^2 is the fraction of the variability that can be explained by regression.

Bruce M.
Boghossian

Properties of
Linear Model
estimators

Estimating $\hat{\sigma}^2$

Covariance
and
Correlation

Summary

- We continued our study of the Linear Model.
- We showed that $\hat{\beta}_1$ is Student T distributed with $n - 2$ degrees of freedom.
 - We used this to do hypothesis testing for $\hat{\beta}_1$.
 - We used this to construct confidence intervals for $\hat{\beta}_1$.
- We defined the *correlation* $\rho(X, Y) \in [-1, +1]$.
- We presented a method of estimating $\rho(X, Y)$ using sample moments.
- We constructed the *Pearson correlation coefficient* R .