

Linear Models Lecture 2

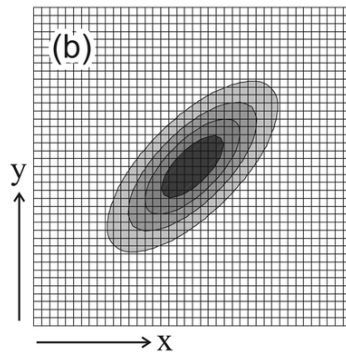
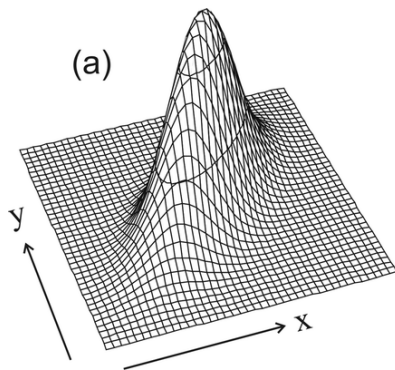
Robert Gulotty

University of Chicago

May 30, 2023

Joint Distributions in Political Science

- We saw the 'roof' distribution, but what joint distributions arise in Political Science?
- If our variables data is discrete and binary, we can have a 2x2 table of probabilities.
- If our data is continuous, we will have a 3d density.



Definition

$\mathbf{X} = (X_1, X_2, \dots, X_N)$ is a vector of random variables. If \mathbf{X} has pdf

$$f(\mathbf{x}) = (2\pi)^{-N/2} \det(\Sigma)^{-1/2} \exp\left(-\frac{1}{2}(\mathbf{x} - \boldsymbol{\mu})' \Sigma^{-1}(\mathbf{x} - \boldsymbol{\mu})\right)$$

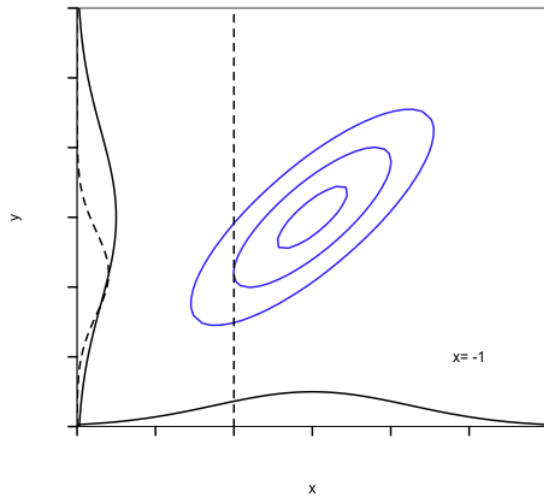
\mathbf{X} is a *Multivariate Normal* Distribution,

$$\mathbf{X} \sim \text{Multivariate Normal}(\boldsymbol{\mu}, \Sigma)$$

$\mathbf{Z} = (X, Y)$ is a vector of jointly distributed random variables.

$$f(\mathbf{z}) = (2\pi)^{-1} \det \begin{bmatrix} \sigma_x^2 & \sigma_{xy} \\ \sigma_{yx} & \sigma_y^2 \end{bmatrix}^{-1/2} \exp \left(-\frac{1}{2} \left(\begin{bmatrix} x \\ y \end{bmatrix} - \begin{bmatrix} \mu_x \\ \mu_y \end{bmatrix} \right)' \begin{bmatrix} \sigma_x^2 & \sigma_{xy} \\ \sigma_{yx} & \sigma_y^2 \end{bmatrix}^{-1} \left(\begin{bmatrix} x \\ y \end{bmatrix} - \begin{bmatrix} \mu_x \\ \mu_y \end{bmatrix} \right) \right)$$

$$f(\mathbf{z}) = \frac{1}{\sqrt{(2\pi)(2\pi)(\sigma_x^2\sigma_y^2 - \sigma_{xy}^2)}} \exp \left(-\frac{1}{2} \begin{bmatrix} x - \mu_x \\ y - \mu_y \end{bmatrix}' \begin{bmatrix} \sigma_x^2 & \sigma_{xy} \\ \sigma_{yx} & \sigma_y^2 \end{bmatrix}^{-1} \begin{bmatrix} x - \mu_x \\ y - \mu_y \end{bmatrix} \right)$$



Prediction

- If we want to predict y with a known pdf. The “best” prediction is μ_Y (minimum mean square error (MSE)).
- Suppose we have a bivariate distribution $f(x, y)$. We will be told X . What should we guess for Y ?
- We can use any function $h(X)$. What is the best, minimizing $(Y - h(X))^2$?
- We will show that only the Conditional Expectation Function $E(Y|X)$ is the best.

Proof that CEF is the best predictor (Goldberger pg 51)

- Goal: Minimize $(Y - h(X))^2$.
- Define $U \equiv Y - h(X)$, $\epsilon \equiv Y - E(Y|X)$, $W \equiv E(Y|X) - h(X)$.

$$U = Y - h(X)$$

$$U = (\epsilon + E(Y|X)) - (E(Y|X) - W)$$

$$U = \epsilon + W$$

- W only depends on X . So for $X = x$, let's call it $W(X = x) = w$.

$$U^2 = \epsilon^2 + 2w\epsilon + w^2$$

$$E(U^2|x) = E(\epsilon^2|x) + 2E(w\epsilon|x) + E(w^2|x)$$

$$= E(\epsilon^2|x) + 2wE(\epsilon|x) + w^2$$

$$= E(\epsilon^2|x) + 2w * 0 + w^2$$

$$E(U^2|x) = E(\epsilon^2|x) + 2w * 0 + w^2$$

from above

$$E(U^2) = E_X[E(U^2|X)]$$

by Law of Iterated Expectations

$$= E_X[E(\epsilon^2|x) + 0 + w^2]$$

Plugging in

$$= E[\sigma_{Y|x}^2] + E[W^2]$$

$E[W^2] \geq 0$, so $E(U^2)$ is minimized if $W = 0$, and recall the definition of W .

$$W = E(Y|X) - h(X)$$

$$0 = E(Y|X) - h(X)$$

$$h(X) = E(Y|X)$$

The Conditional Expectation Function $E(Y|X)$ *is* the function that minimizes $E(U^2)$.

Conditional Expectation as a Prediction

Suppose Y and X are random variables.

$$\mu_{Y|X} \equiv E[Y|X]$$

$\mu_{Y|X}$ is our best guess for Y given X
 ϵ is the amount we are off.

$$\epsilon \equiv Y - \mu_{Y|X}$$

Properties of Estimation error

ϵ is a random variable where:

$$\begin{aligned} E[\epsilon|X] &= E[Y - \mu_{Y|X}|X] \\ &= E[Y|X] - E[\mu_{Y|X}|X] \\ &= E[Y|X] - E[E[Y|X]|X] \\ &= \mu_{Y|X} - \mu_{Y|X} = 0 \end{aligned}$$

By the law of iterated expectations:

$$E[\epsilon] = E[E[\epsilon|X]] = 0$$

The disturbances center on 0 *by construction*.

Covariance of Estimator and Disturbance

$$E[\epsilon \mu_{Y|X}] = E[E[\epsilon \mu_{Y|X} | X]] = E[\mu_{Y|X} E[\epsilon | X]] = E[\mu_{Y|X} * 0] = 0$$

because $\mu_{Y|X}$ only depends on X .

$$\text{Cov}(\epsilon, \mu_{Y|X}) = E[\epsilon \mu_{Y|X}] - E[\epsilon]E[\mu_{Y|X}] = 0 - 0 * \mu_{Y|X} = 0$$

The disturbance is uncorrelated with the conditional expectation. This is what allows us to separate the signal from the noise.

Discussion of CEF

- CEF solves the minimum mean squared error (MSE) prediction problem.
- But it depends on knowing the conditional distribution of $Y|x$, because $E[Y|X] = \int y f_{Y|x}(y|x) dy$.
- This is a general problem: MSE minimizing estimators often require knowledge about the population.
- Instead we will restrict attention to linear predictors. We seek to find the "best" linear predictor (BLP).
- Later we will show that OLS regression estimates are BLUE (best linear unbiased estimators) **if** the CEF is linear and residuals are constant.
- However, what if the CEF is not linear?
- Linear regression is the best (minimum MSE) linear approximation to the CEF.

(Best) Linear Predictors

- Suppose we want a predictor $E(Y|X)$ that is **linear**:

$$f(X) = a + bX$$

Our standard for prediction is to minimize the mean-square error (best):

- Whereas the CEF might be infinitely complex, the BLP is characterized just by two numbers, a and b . If the CEF is linear, the BLP is the CEF.
Choose a and b to minimize

$$M = E[(Y - (a + bX))^2]$$

Best Linear Predictor (a)

Take partial derivative with respect to a .

$$\begin{aligned}\frac{\partial E[(Y - (a + bX))^2]}{\partial a} &= E\left[\frac{\partial(Y - (a + bX))^2}{\partial a}\right] \\ &= E[-2(Y - (a + bX))]\end{aligned}$$

setting equal to 0:

$$0 = E[-2(Y - (a + bX))]$$

$$0 = -2E[Y] + 2E(a + bX)$$

$$E(Y) = a + bE(X)$$

$$E(Y) - bE(X) = a$$

Best Linear Predictor (b)

Take partial derivative with respect to b .

$$\begin{aligned}\frac{\partial E[(Y - (a + bX))^2]}{\partial b} &= E[\partial(Y - (a + bX))^2 / \partial b] \\ &= 2E[(Y - (a + bX))(-X)]\end{aligned}$$

setting equal to 0:

$$\begin{aligned}0 &= -2E[(YX) - (a + bX)(X)] \\ 0 &= -2E(YX) + 2E[(a + bX)(X)]\end{aligned}$$

Best Linear Predictor (b continued)

$$0 = -2E[YX] + 2E[(a + bX)(X)]$$

$$E(YX) = aE(X) + bE(X^2)$$

$$E(YX) = [E(Y) - bE(X)]E(X) + bE(X^2)$$

$$E(YX) = E(Y)E(X) - bE(X)E(X) + bE(X^2)$$

$$E(YX) - E(Y)E(X) = b[E(X^2) - E(X)^2]$$

$$b = \frac{E(YX) - E(Y)E(X)}{[E(X^2) - E(X)^2]} = \frac{\text{Cov}(X, Y)}{\text{Var}(X)}$$

Best Linear Predictor

$$E[Y|X] = \beta_0 + \beta_1 X$$
$$Y = \beta_0 + \beta_1 X + \epsilon$$

- Parameters (in Greek)
 - $\epsilon = Y - E[Y|X]$
 - The slope $\beta = \frac{\text{Cov}(X, Y)}{\text{Var}(X)}$
 - The intercept $\alpha = E(Y) - \beta E(X)$
- Y is called the dependent variable, X is called the independent variable.

What is "linear in parameters"?

- The Best Linear Predictor is linear in *parameters* $\theta \in \{\beta_0, \beta_1 \dots\}$.
- Examples of linear in parameters:

$$Y = \beta_0 + \beta_1 X^4 + \epsilon$$

$$Y = \beta_0 + \beta_1 e^X + \epsilon$$

- Examples of nonlinear in parameters:

$$Y = \beta_0 + \frac{1}{\beta_1} X + \epsilon$$

$$Y = \beta_0 + \beta_1^2 X + \epsilon$$

$$Y = \beta_0 + e^{\beta_1 X} + \epsilon$$

Deriving variance of BLP at minimized values

$$\beta = \frac{\text{Cov}(X, Y)}{\text{Var}(X)}$$

$$\beta \text{Var}(X) = \text{Cov}(X, Y)$$

$$\begin{aligned} E[(Y - (\alpha + \beta X))^2] - E[Y - (\alpha + \beta X)]^2 &= \text{Var}(Y - (\alpha + \beta X)) \\ &= \text{Var}(Y - \beta X) \\ &= \text{Var}(Y) + \text{Var}(\beta X) - 2\text{Cov}(\beta X, Y) \\ &= \text{Var}(Y) + \beta^2 \text{Var}(X) - 2\beta \text{Cov}(X, Y) \\ &= \text{Var}(Y) + \beta^2 \text{Var}(X) - 2\beta^2 \text{Var}(X) \\ &= \text{Var}(Y) - \beta^2 \text{Var}(X) \\ \sigma_\epsilon^2 &= \sigma_Y^2 - \beta^2 \sigma_X^2 \end{aligned}$$

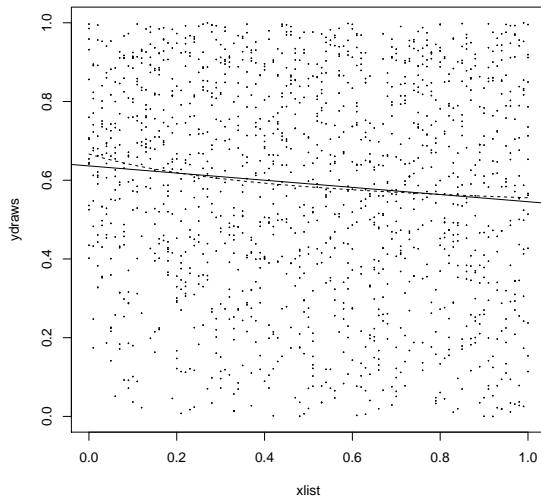
Example Best Linear Predictor, roof

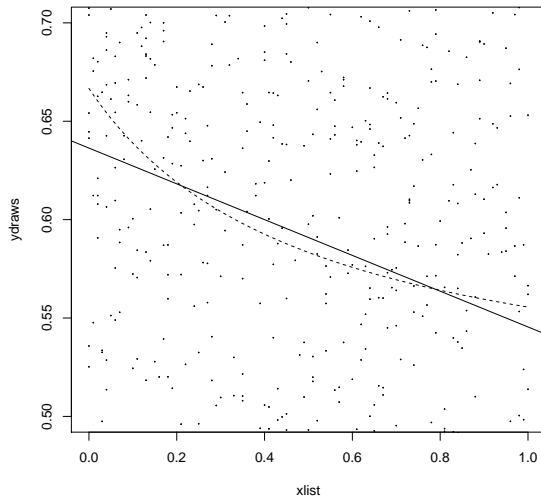
Example: $X + Y$

Suppose X and Y have pdf $x + y$ for $x, y \in [0, 1]$.

$$b = \frac{\text{Cov}(X, Y)}{\text{Var}(X)} = -1/11$$

$$a = E(Y) - bE(X) = 84/132$$





Sample vs Population

- Theoretical "Population" Objects (not observed): Greek Letters: α , β , μ , σ , ϵ , $E()$, $\text{Var}()$, $\text{Cov}()$.
- Empirical Objects (observed): latin letters, a , b , \bar{x} , s^2 , e , $\widehat{E}()$, $\widehat{\text{Cov}}()$.
- Objects of the first group are not exactly equal to their empirical analogues from the second group.
 - $\alpha = \mu_Y - \beta\mu_X$ is correct.
 - $a = \bar{y} - b\bar{x}$ is correct.
 - $E(\bar{x}) = \mu$ is correct.
 - $\bar{x} = \mu$ is not correct.
 - $\text{Cov}(x) = s_x^2$ is not correct.

Sample Linear projection

Population linear projection:

$$E(Y|X) = \alpha + \beta X$$

where

$$\beta = \frac{\sigma_{XY}}{\sigma_X^2}, \quad \alpha = \mu_Y - \beta \mu_X$$

Sample linear projection:

$$\hat{Y} = b_0 + b_1 X$$

where

$$b_1 = \frac{S_{XY}}{S_X^2}, \quad b_0 = \bar{Y} - b_1 \bar{X}$$

Asymptotic Distribution of Sample Slope

Call $X^* = X - \bar{X}$, $\epsilon = Y - (\alpha + \beta X)$, then the Bivariate Delta Method (Goldberg 10.5) tells us

$$b_1 \overset{A}{\sim} N\left(\beta, \frac{E(X^{*2}\epsilon^2)}{(\sigma_X^2)^2}\right)$$

If $E(\epsilon^2|X)$ is a constant σ^2 , then

$$\begin{aligned} E(X^{*2}\epsilon^2) &= E_X[E(X^{*2}\epsilon^2|X)] = E_X[X^{*2}E(\epsilon^2|X)] \\ &= E_X[X^{*2}\sigma^2] \\ &= \sigma^2 E_X[X^{*2}] = \sigma^2 V(X) = \sigma^2 \sigma_X^2 \end{aligned}$$

$$b_1 \overset{A}{\sim} N\left(\beta, \frac{\sigma^2}{\sigma_X^2}\right)$$

Ordinary Least Squares Regression Estimator

- In our effort to approximate the CEF, we simply replaced the expectations, covariances etc. in the BLP with the sample means, covariances etc.
- It turns out that we can algebraically process data with the least squares procedure (OLS) to estimate the BLP parameters.
- That is, we will solve $\min_{b_0, b_1} \sum_i e^2$
- In addition, **if** the CEF function is linear and $E(\epsilon^2|X)$ is a constant σ^2 , then the Gauss Markov Theorem tells us that OLS is a Best Linear Unbiased Estimator (BLUE).

Brief Linear Algebra Interlude

- $\mathbf{x} = \begin{bmatrix} x_1 \\ x_2 \\ \dots \\ x_n \end{bmatrix}$ is $n \times 1$ column vector. $\mathbf{i} = \begin{bmatrix} 1 \\ 1 \\ \dots \\ 1 \end{bmatrix}$.
- $\mathbf{x}' = \mathbf{x}^T = [x_1 \ x_2 \ \dots \ x_n]$ is $1 \times n$ vector.
- If c is a scalar, $c\mathbf{x}' = [cx_1 \ cx_2 \ \dots \ cx_n]$.
- $\mathbf{x}'\mathbf{i} = \mathbf{x} \cdot \mathbf{i} = x_1 * 1 + x_2 * 1 + x_3 * 1 + \dots + x_n * 1 = \sum x_i = n\bar{x}$
- $\mathbf{x}'\mathbf{x} = x_1^2 + x_2^2 + x_3^2 + \dots + x_n^2 = \sum x_i^2 = n * \widehat{Var}(x) + n\bar{x}^2$
- $\mathbf{xx}' = \begin{bmatrix} x_1x_1 & x_1x_2 & \dots \\ x_2x_1 & \dots & \\ \dots & & \end{bmatrix}, \quad \mathbf{I} = \begin{bmatrix} 1 & 0 & \dots \\ 0 & 1 & \dots \\ \dots & & \end{bmatrix}$

Linear Algebra Rules

If \mathbf{x} , \mathbf{y} , and \mathbf{z} are vectors of equal length:

- You can add, subtract and dot product them. No division of vectors.
- Commutative property: If \mathbf{x} and \mathbf{y} are vectors of equal length, $\mathbf{x}'\mathbf{y} = \mathbf{y}'\mathbf{x}$
- Distributive property: $\mathbf{x}'(\mathbf{y} + \mathbf{z}) = \mathbf{x}'\mathbf{y} + \mathbf{x}'\mathbf{z}$

Linear Algebra Geometry (Advanced Topic)

- $\sqrt{\mathbf{x}'\mathbf{x}} = \|\mathbf{x}\|$ is called the Euclidean Norm, from the classic Euclidean distance $\sqrt{a^2 + b^2} = c$ in Descartes' theory of coordinates.
- $\mathbf{u} = \frac{\mathbf{x}}{\sqrt{\mathbf{x}'\mathbf{x}}}$ is called the **unit vector** in the direction of \mathbf{x} .
- If \mathbf{u} is a unit vector in the direction of \mathbf{x} , $\mathbf{y}'\mathbf{u} = \mathbf{u}'\mathbf{y}$ is the *scalar projection* of \mathbf{y} onto \mathbf{x} .
- $(\mathbf{u}'\mathbf{y})\mathbf{u}$ is the *vector projection* of \mathbf{y} onto \mathbf{x} .

$$(\mathbf{u}'\mathbf{y})\mathbf{u} = \left(\frac{\mathbf{x}'\mathbf{y}}{\sqrt{\mathbf{x}'\mathbf{x}}} \right) \frac{\mathbf{x}}{\sqrt{\mathbf{x}'\mathbf{x}}} = \left(\frac{\mathbf{x}'\mathbf{y}}{\mathbf{x}'\mathbf{x}} \right) \mathbf{x}$$

- If $\mathbf{x}'\mathbf{y} = 0$, \mathbf{x} is "orthogonal" to \mathbf{y} .

Gauss-Markov Assumptions (classical fixed \mathbf{x})

1 CEF is linear

2 $\mathbf{x} = \begin{bmatrix} x_1 \\ x_2 \\ \dots \\ x_n \end{bmatrix}$ is fixed.

3 $\sum_{i=1}^n (x_i - \bar{x})^2 > 0$ [\mathbf{x} takes on at least two values].

4 $E(\epsilon_i) = 0 \quad \forall i$.

5 $E(\epsilon\epsilon') = \sigma_\epsilon^2 \mathbf{I}$.

- $\text{Var}(\epsilon_i) = \sigma_\epsilon^2 \quad \forall i$ [homoskedasticity, i.i.d.]
- $\text{Cov}(\epsilon_i, \epsilon_j) = 0$

Deriving b_0 , b_1 estimator using OLS

$$\min_{b_0} \min_{b_1} \sum_i e^2 = \min_{b_0} \min_{b_1} \mathbf{e}'\mathbf{e}$$

$$\begin{aligned}\frac{\partial}{\partial b_1} \mathbf{e}'\mathbf{e} &= \frac{\partial}{\partial b_1} (\mathbf{y} - \hat{\mathbf{y}})'(\mathbf{y} - \hat{\mathbf{y}}) \\ &= \frac{\partial}{\partial b_1} (\mathbf{y}'\mathbf{y} - \mathbf{y}'\hat{\mathbf{y}} - \hat{\mathbf{y}}'\mathbf{y} + \hat{\mathbf{y}}'\hat{\mathbf{y}}) \\ &= \frac{\partial}{\partial b_1} (\mathbf{y}'\mathbf{y} - 2\hat{\mathbf{y}}'\mathbf{y} + \hat{\mathbf{y}}'\hat{\mathbf{y}}) \\ &= 0 - 2\frac{\partial}{\partial b_1} \hat{\mathbf{y}}'\mathbf{y} + \frac{\partial}{\partial b_1} \hat{\mathbf{y}}'\hat{\mathbf{y}}\end{aligned}$$

$$\begin{aligned}\frac{\partial}{\partial b_0} \mathbf{e}'\mathbf{e} &= \frac{\partial}{\partial b_0} (\mathbf{y} - \hat{\mathbf{y}})'(\mathbf{y} - \hat{\mathbf{y}}) \\ &= 0 - 2\frac{\partial}{\partial b_0} \hat{\mathbf{y}}'\mathbf{y} + \frac{\partial}{\partial b_0} \hat{\mathbf{y}}'\hat{\mathbf{y}}\end{aligned}$$

Note on calculus with vectors (I)

$$\hat{\mathbf{y}} = \begin{bmatrix} b_0 + b_1 x_1 \\ b_0 + b_1 x_2 \\ b_0 + b_1 x_3 \\ \dots \end{bmatrix} \quad \frac{\partial \hat{\mathbf{y}}}{\partial b_1} = \begin{bmatrix} x_1 \\ x_2 \\ x_3 \\ \dots \end{bmatrix} = \mathbf{x} \quad \frac{\partial \hat{\mathbf{y}}}{\partial b_0} = \begin{bmatrix} 1 \\ 1 \\ 1 \\ \dots \end{bmatrix} = \mathbf{i}$$

$$\hat{\mathbf{y}}' \hat{\mathbf{y}} = \begin{bmatrix} b_0 + b_1 x_1 & b_0 + b_1 x_2 & b_0 + b_1 x_3 & \dots \end{bmatrix} \begin{bmatrix} b_0 + b_1 x_1 \\ b_0 + b_1 x_2 \\ b_0 + b_1 x_3 \\ \dots \end{bmatrix} = \sum_i (b_0 + b_1 x_i)^2$$

Note on calculus with vectors (II)

$$\begin{aligned}\frac{\partial}{\partial b_1} \hat{\mathbf{y}}' \hat{\mathbf{y}} &= \frac{\partial}{\partial b_1} \sum_i (b_0 + b_1 x_i)^2 \\ &= \sum_i 2x_i (b_0 + b_1 x_i) \\ &= 2\mathbf{x}' [b_0 \mathbf{i} + b_1 \mathbf{x}] \\ \frac{\partial}{\partial b_0} \hat{\mathbf{y}}' \hat{\mathbf{y}} &= \sum_i 2(b_0 + b_1 x_i) \\ &= 2\mathbf{i}' [b_0 + b_1 \mathbf{x}]\end{aligned}$$

Deriving b_0 , b_1 estimator using OLS

$$\begin{aligned}\frac{\partial}{\partial b_0} \mathbf{e}'\mathbf{e} &= 0 - 2\frac{\partial}{\partial b_0} \hat{\mathbf{y}}'\mathbf{y} + \frac{\partial}{\partial b_0} \hat{\mathbf{y}}'\hat{\mathbf{y}} \\ &= 0 - 2\mathbf{i}'\mathbf{y} + 2\mathbf{i}'[(b_0 + b_1\mathbf{x})]\end{aligned}$$

$$\mathbf{i}'b_0 = \mathbf{i}'\mathbf{y} - b_1\mathbf{i}'\mathbf{x}$$

$$nb_0 = \mathbf{i}'\mathbf{y} - b_1\mathbf{i}'\mathbf{x}$$

$$nb_0 = n\bar{y} - b_1n\bar{x}$$

$$b_0 = \bar{y} - b_1\bar{x}$$

$$\begin{aligned}\frac{\partial}{\partial b_1} \mathbf{e}'\mathbf{e} &= 0 - 2\frac{\partial}{\partial b_1} \hat{\mathbf{y}}'\mathbf{y} + \frac{\partial}{\partial b_1} \hat{\mathbf{y}}'\hat{\mathbf{y}} \\ &= -2\mathbf{x}'\mathbf{y} + 2\mathbf{x}'[b_0\mathbf{i} + b_1\mathbf{x}] \\ &= -2\mathbf{x}'\mathbf{y} + 2\mathbf{x}'[(\bar{y} - b_1\bar{x})\mathbf{i} + b_1\mathbf{x}] \\ &= -2\mathbf{x}'\mathbf{y} + 2\mathbf{x}'\bar{y}\mathbf{i} + 2\mathbf{x}'b_1(\mathbf{x} - \bar{x}\mathbf{i}) \\ &= -2\mathbf{x}'\mathbf{y} + 2n\bar{x}\bar{y} + 2b_1[\mathbf{x}'\mathbf{x} - n\bar{x}\bar{x}]\end{aligned}$$

$$\mathbf{x}'\mathbf{y} - n\bar{x}\bar{y} = b_1[\mathbf{x}'\mathbf{x} - n\bar{x}\bar{x}]$$

$$\widehat{Cov}(\mathbf{x}, \mathbf{y}) = b_1 \widehat{Var}(\mathbf{x})$$

$$\frac{\widehat{Cov}(\mathbf{x}, \mathbf{y})}{\widehat{Var}(\mathbf{x})} = b_1 \quad \text{By assumption 3}$$

```
> x <- 1:50
> y <- 8 + .5 * x + rnorm(50)

> coef(lm(y~x))
(Intercept)      x 
8.1679311    0.4982869

> mean(y) - cov(x,y) / var(x) * mean(x)
[1] 8.167931
> cov(x,y)/var(x)
[1] 0.4982869
```

Interpretation

- b_1 is the slope: one unit change in x is associated with a b_1 unit change in Y .

$$\hat{\mathbf{y}} = b_0 + \frac{\widehat{\text{Cov}}(\mathbf{x}, \mathbf{y})}{\widehat{\text{Var}}(\mathbf{x})} \mathbf{x}$$
$$\mathbf{y} = \hat{\mathbf{y}} + \mathbf{e}$$

- b_0 is the intercept: What is the predicted value of y when $x = 0$.
- e_i is the residual: observed y minus predicted y .

Decomposition of Variance

$$\sigma_{\epsilon}^2 = \sigma_Y^2 - \beta^2 \sigma_X^2$$

$$\widehat{Var}(\mathbf{e}) = \widehat{Var}(\mathbf{y}) - b_1^2 \widehat{Var}(\mathbf{x})$$

$$b_1^2 \widehat{Var}(\mathbf{x}) + \widehat{Var}(\mathbf{e}) = \widehat{Var}(\mathbf{y})$$

$$\underbrace{b_1^2 \sum (x_i - \bar{x})^2}_{\text{Explained Sum of Squares}} + \underbrace{\sum e_i^2}_{\text{Residual Sum of Squares}} = \underbrace{\sum (y_i - \bar{y})^2}_{\text{Total Sum of Squares}}$$

$$ESS + RSS = TSS$$

$$\frac{ESS}{TSS} + \frac{RSS}{TSS} = 1$$

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

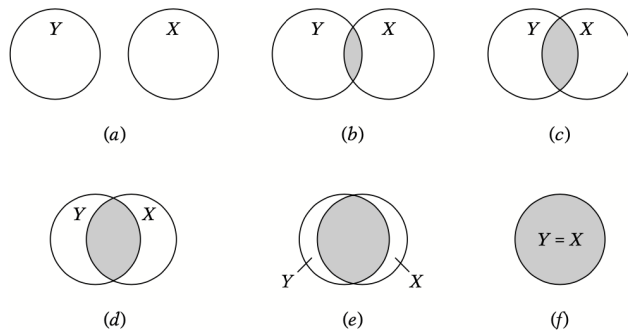


FIGURE 3.9 The Ballentine view of r^2 : (a) $r^2 = 0$; (f) $r^2 = 1$.

Sample Correlation Coefficient

$$r_{xy} \equiv \frac{\widehat{Cov}(\mathbf{x}, \mathbf{y})}{\sqrt{\widehat{Var}(\mathbf{x})}\sqrt{\widehat{Var}(\mathbf{y})}} = \frac{S_{XY}}{S_X S_Y} = \frac{\frac{1}{n-1} \sum (x_i - \bar{x})(y_i - \bar{y})}{\sqrt{\frac{1}{n-1} \sum (x_i - \bar{x})^2} \sqrt{\frac{1}{n-1} \sum (y_i - \bar{y})^2}}$$

- The correlation coefficient (r_{xy}) is the normalized linear relationship.
- The OLS slope coefficient (b_1) is the unnormalized linear relationship.

OLS coefficient vs Correlation

$$b_1 = \frac{\widehat{\text{Cov}}(\mathbf{x}, \mathbf{y})}{\sqrt{\widehat{\text{Var}}(\mathbf{x})}} \frac{1}{\sqrt{\widehat{\text{Var}}(\mathbf{x})}} \quad (\text{Definition of } b_1)$$

$$\begin{aligned} \frac{b_1}{\sqrt{\widehat{\text{Var}}(\mathbf{y})}} &= \frac{\widehat{\text{Cov}}(\mathbf{x}, \mathbf{y})}{\sqrt{\widehat{\text{Var}}(\mathbf{x})} \sqrt{\widehat{\text{Var}}(\mathbf{y})}} \frac{1}{\sqrt{\widehat{\text{Var}}(\mathbf{x})}} \\ &= r_{xy} \frac{1}{\sqrt{\widehat{\text{Var}}(\mathbf{x})}} \\ b_1 &= r_{xy} \frac{\sqrt{\widehat{\text{Var}}(\mathbf{y})}}{\sqrt{\widehat{\text{Var}}(\mathbf{x})}} = r_{xy} \frac{S_Y}{S_X} \end{aligned}$$

```
> x <- 1:50
> y <- 8 + .5 * x + rnorm(50)

> coef(lm(y~x))
(Intercept)          x
8.1679311      0.4982869

> cor(x,y)
[1] 0.9897854

> cor(x,y)*sd(y)/sd(x)
[1] 0.4982869
```

Obtaining the Sample Correlation Coefficient from a Regression

$$\hat{\mathbf{y}} = b_0 + b_1 \mathbf{x}$$

$$\hat{\mathbf{y}} = \bar{y} - b_1 \bar{x} + b_1 \mathbf{x}$$

$$\hat{\mathbf{y}} - \bar{y} = b_1 (\mathbf{x} - \bar{x})$$

$$\hat{\mathbf{y}} - \bar{y} = r_{xy} \frac{S_Y}{S_X} (\mathbf{x} - \bar{x})$$

$$\frac{\hat{\mathbf{y}} - \bar{y}}{S_Y} = r_{xy} \frac{(\mathbf{x} - \bar{x})}{S_X}$$

$$y^* = \rho_0 + \rho_1 x^*$$

```
> x <- 1:50
> y <- 8 + .5 * x + rnorm(50, sd=3)
>
> cor(x, y)
[1] 0.9139259
>
> yhatnorm <- (predict(lm(y~x)) - mean(y)) / sd(y)
> xhatnorm <- (x - mean(x)) / sd(x)
>
> coef(lm(yhatnorm ~ xhatnorm - 1 ))
      xhatnorm
0.9139259
```

Proof that OLS unbiased $\text{Bias}(\hat{\theta}) \equiv E(\hat{\theta}) - \theta$

- Assume specification : $\mathbf{y} = [\beta_0 \mathbf{i} + \beta_1 \mathbf{x} + \epsilon]$, $\mathbf{i}' \mathbf{y} = [\mathbf{i}' \beta_0 \mathbf{i} + \mathbf{i}' \beta_1 \mathbf{x} + \mathbf{i}' \epsilon]$,
 $\bar{y} = \beta_0 + \beta_1 \bar{x} + \bar{\epsilon}$

$$\begin{aligned}
 b_1 &= \frac{\mathbf{x}' \mathbf{y} - n \bar{x} \bar{y}}{\widehat{\text{Var}}(x)} \\
 b_1 &= \frac{\mathbf{x}' [\beta_0 \mathbf{i} + \beta_1 \mathbf{x} + \epsilon] - n \bar{x} [\beta_0 + \beta_1 \bar{x} + \bar{\epsilon}]}{\widehat{\text{Var}}(x)} \\
 &= \frac{n \bar{x} \beta_0 - n \bar{x} \beta_0 + \mathbf{x}' \beta_1 \mathbf{x} - n \bar{x} \beta_1 + \mathbf{x}' \epsilon - n \bar{x} \bar{\epsilon}}{\widehat{\text{Var}}(x)} \\
 &= \frac{\beta_1 [\mathbf{x}' \mathbf{x} - n \bar{x}] + \mathbf{x}' \epsilon - n \bar{x} \bar{\epsilon}}{\widehat{\text{Var}}(x)} \\
 &= \beta_1 + \frac{\mathbf{x}' \epsilon - n \bar{x} \bar{\epsilon}}{\widehat{\text{Var}}(x)}
 \end{aligned}$$

Proof that OLS unbiased

$$\blacksquare \text{Bias}(\hat{\theta}) \equiv E(\hat{\theta}) - \theta$$

$$b_1 = \beta_1 + \frac{\mathbf{x}'\epsilon - n\bar{x}\bar{\epsilon}}{\widehat{\text{Var}}(x)}$$

$$E[b_1] = E[\beta_1] + E\left[\frac{\mathbf{x}'\epsilon - n\bar{x}\bar{\epsilon}}{\widehat{\text{Var}}(x)}\right]$$

$$E[b_1] = E[\beta_1] + \frac{\mathbf{x}'E[\epsilon] - n\bar{x}E[\bar{\epsilon}]}{\widehat{\text{Var}}(x)} \quad \text{by assumption (2)}$$

$$E[b_1] = E[\beta_1] + \frac{\mathbf{x}'0 - n\bar{x}0}{\widehat{\text{Var}}(x)}$$

Change of Notation

- The notation simplifies if we center our variables $\tilde{\mathbf{x}} = \mathbf{x} - \bar{x}$, $\tilde{\mathbf{y}} = \mathbf{y} - \bar{y}$,
- $b_1 = \frac{\tilde{\mathbf{x}}' \tilde{\mathbf{y}}}{\tilde{\mathbf{x}}' \tilde{\mathbf{x}}}$
- Note: $b_1 \tilde{\mathbf{x}} = \frac{\tilde{\mathbf{x}}' \tilde{\mathbf{y}}}{\tilde{\mathbf{x}}' \tilde{\mathbf{x}}} \tilde{\mathbf{x}}$ is the vector projection of \mathbf{y} in $\tilde{\mathbf{x}}$.
- Define weights $\mathbf{w} \equiv \frac{\tilde{\mathbf{x}}}{\tilde{\mathbf{x}}' \tilde{\mathbf{x}}}$
- Note $\mathbf{w}' \mathbf{w} = \frac{\tilde{\mathbf{x}}' \tilde{\mathbf{x}}}{(\tilde{\mathbf{x}}' \tilde{\mathbf{x}})' \tilde{\mathbf{x}}' \tilde{\mathbf{x}}} = \frac{1}{(\tilde{\mathbf{x}}' \tilde{\mathbf{x}})'}$
- $b_1 = \mathbf{w}' \tilde{\mathbf{y}}$
- Recall under the assumptions above, $b_1 = \beta_1 + \mathbf{w}' \epsilon$. $b_1 - \beta_1 = \mathbf{w}' \epsilon$

Variance of the slope estimator

$$\begin{aligned} \text{Var}(b_1) &= E[(b_1 - E(b_1))^2] \\ &= E[(b_1 - \beta_1)^2] \\ &= E[(\mathbf{w}'\epsilon)^2] \\ &= E[(\mathbf{w}'\epsilon)'(\mathbf{w}'\epsilon)] \\ &= E[(\mathbf{w}'\epsilon\epsilon'\mathbf{w})] \\ &= (\mathbf{w}'E[\epsilon\epsilon']\mathbf{w}) \rightarrow \frac{\sigma^2 I}{\tilde{\mathbf{x}}'\tilde{\mathbf{x}}} \quad \text{by assumption (4).} \end{aligned}$$

Gauss Markov BLUE Step 1

- The OLS weights \mathbf{w} give us: $b_1 = \mathbf{w}'\tilde{\mathbf{y}} = \mathbf{w}'\mathbf{y}$
- Suppose we had some other unbiased linear estimator with weights $\mathbf{c} = \mathbf{w} + \mathbf{d}$

$$b_* = \mathbf{c}'\tilde{\mathbf{y}} = \mathbf{c}'\mathbf{y}$$

$$b_* = \mathbf{c}'[\beta_0\mathbf{i} + \beta_1\tilde{\mathbf{x}} + \boldsymbol{\epsilon}]$$

$$b_* = \beta_0\mathbf{c}'\mathbf{i} + \beta_1\mathbf{c}'\tilde{\mathbf{x}} + \mathbf{c}'\boldsymbol{\epsilon}$$

$$E[b_*] = \beta_0(\mathbf{w} + \mathbf{d})'\mathbf{i} + \beta_1(\mathbf{w} + \mathbf{d})'\tilde{\mathbf{x}} + (\mathbf{w} + \mathbf{d})'E[\boldsymbol{\epsilon}]$$

$$= \beta_0(\mathbf{w} + \mathbf{d})'\mathbf{i} + \beta_1\frac{\tilde{\mathbf{x}}'}{\tilde{\mathbf{x}}'\tilde{\mathbf{x}}}\tilde{\mathbf{x}} + \beta_1\mathbf{d}'\tilde{\mathbf{x}} + (\mathbf{w} + \mathbf{d})'E[\boldsymbol{\epsilon}] = \beta_1$$

only if

$$(\mathbf{w} + \mathbf{d})'\mathbf{i} = 0 \text{ and } \mathbf{d}'\tilde{\mathbf{x}} = 0$$

Gauss Markov BLUE Step 2

- This gives us that our alternative estimator is:

$$b_* = \beta_0 \mathbf{c}' \mathbf{i} + \beta_1 \mathbf{c}' \tilde{\mathbf{x}} + \mathbf{c}' \boldsymbol{\epsilon} = \beta_1 + \mathbf{c}' \boldsymbol{\epsilon}$$

$$b_* - \beta_1 = \mathbf{c}' \boldsymbol{\epsilon}$$

$$\begin{aligned} \text{Var}(b_*) &= E((b_* - \beta_1)^2) \\ &= E(\mathbf{c}' \boldsymbol{\epsilon} \boldsymbol{\epsilon}' \mathbf{c}) \\ &= \sigma_{\epsilon}^2 (\mathbf{c}' \mathbf{c}) \\ &= \sigma_{\epsilon}^2 (\mathbf{w} + \mathbf{d})' (\mathbf{w} + \mathbf{d}) \\ &= \sigma_{\epsilon}^2 (\mathbf{w}' \mathbf{w} + \mathbf{d}' \mathbf{w} + \mathbf{w}' \mathbf{d} + \mathbf{d}' \mathbf{d}) \\ &= \sigma_{\epsilon}^2 \mathbf{w}' \mathbf{w} + \sigma_{\epsilon}^2 \mathbf{d}' \mathbf{w} + \sigma_{\epsilon}^2 \mathbf{w}' \mathbf{d} + \sigma_{\epsilon}^2 \mathbf{d}' \mathbf{d} \\ &= \sigma_{\epsilon}^2 \mathbf{w}' \mathbf{w} + \sigma_{\epsilon}^2 \mathbf{d}' \mathbf{d} > \text{Var}(b_1) \end{aligned}$$

Discussion of Gauss Markov

- Gauss Markov shows that b_1 is the best (lowest variance) among linear unbiased estimators of β_1 .
- We used all of the assumptions to get this result.
 - We can easily dispense with the assumption that \mathbf{x} is fixed.
 - The assumption that $\widehat{Var}(x_i)$ is not zero is 1) not problematic and 2) testable.
 - If $E(\epsilon_j) = 0$ doesn't hold, we have bias,
 - If $Var(\epsilon_j)$ is not a constant, we have heteroskedasticity,
 - If $Cov(\epsilon_i, \epsilon_j)$ is not zero, we have serial dependence.

Expectation of Sum of Squared Residuals

$$\begin{aligned}E[\sum e_i^2] &= E[\mathbf{e}'\mathbf{e}] \\&= E[TSS - ESS] \\&= E[\tilde{\mathbf{y}}'\tilde{\mathbf{y}}] - E[b_1^2\tilde{\mathbf{x}}'\tilde{\mathbf{x}}] \\&= E[(\beta\tilde{\mathbf{x}} + \tilde{\epsilon})'(\beta\tilde{\mathbf{x}} + \tilde{\epsilon})] - E[b_1^2\tilde{\mathbf{x}}'\tilde{\mathbf{x}}] \\&= [\beta^2\tilde{\mathbf{x}}'\tilde{\mathbf{x}} + 2\beta E[\tilde{\mathbf{x}}'\tilde{\epsilon}] + E[\tilde{\epsilon}'\tilde{\epsilon}]] - [Var(b_1\tilde{\mathbf{x}}) + E(b_1\tilde{\mathbf{x}})^2] \\&= [\beta^2\tilde{\mathbf{x}}'\tilde{\mathbf{x}} + (n-1)\sigma^2] - [Var(b_1) + E(b_1)^2]\tilde{\mathbf{x}}'\tilde{\mathbf{x}} \\&= [\beta^2\tilde{\mathbf{x}}'\tilde{\mathbf{x}} + (n-1)\sigma^2] - [\frac{\sigma^2}{\tilde{\mathbf{x}}'\tilde{\mathbf{x}}} + \beta_1^2]\tilde{\mathbf{x}}'\tilde{\mathbf{x}} \\&= (n-2)\sigma^2\end{aligned}$$

Matrix Formulation

- We can condense notation even further $\boldsymbol{\beta} = \begin{bmatrix} \beta_0 \\ \beta_1 \end{bmatrix}$, $\mathbf{X} = [\mathbf{i} \quad \mathbf{x}] = \begin{bmatrix} 1 & x_1 \\ 1 & x_2 \\ \dots & \dots \end{bmatrix}$,

$$\mathbf{y} = \beta_0 \mathbf{i} + \beta_1 \mathbf{x} + \boldsymbol{\epsilon}$$

$$\mathbf{y} = \mathbf{X}\boldsymbol{\beta} + \boldsymbol{\epsilon}$$

- $\mathbf{X}'\mathbf{X} = \begin{bmatrix} \mathbf{i}' \\ \mathbf{x}' \end{bmatrix} [\mathbf{i} \quad \mathbf{x}] = \begin{bmatrix} \mathbf{i}'\mathbf{i} & \mathbf{i}'\mathbf{x} \\ \mathbf{x}'\mathbf{i} & \mathbf{x}'\mathbf{x} \end{bmatrix}$.

- $\mathbf{X}'\mathbf{y} = \begin{bmatrix} \mathbf{i}' \\ \mathbf{x}' \end{bmatrix} [\mathbf{y}] = \begin{bmatrix} \mathbf{i}'\mathbf{y} \\ \mathbf{x}'\mathbf{y} \end{bmatrix}$.

Matrix Terms

Suppose $\mathbf{x} = (x_1, x_2, \dots, x_N)$

$$E[\mathbf{x}] = \boldsymbol{\mu}$$

$$\text{Cov}(\mathbf{x}) = \boldsymbol{\Sigma}$$

So that,

$$\boldsymbol{\Sigma} = \begin{pmatrix} \text{Var}(x_1) & \text{Cov}(x_1, x_2) & \dots & \text{Cov}(x_1, x_N) \\ \text{Cov}(x_2, x_1) & \text{Var}(x_2) & \dots & \text{Cov}(x_2, x_N) \\ \vdots & \vdots & \ddots & \vdots \\ \text{Cov}(x_N, x_1) & \text{Cov}(x_N, x_2) & \dots & \text{Var}(x_N) \end{pmatrix}$$

Formula for Matrix Inverse

$$\begin{aligned}
 (\mathbf{X}'\mathbf{X})^{-1} &= \begin{bmatrix} \mathbf{i}'\mathbf{i} & \mathbf{i}'\mathbf{x} \\ \mathbf{x}'\mathbf{i} & \mathbf{x}'\mathbf{x} \end{bmatrix}^{-1} \\
 &= \frac{1}{\mathbf{x}'\mathbf{x}\mathbf{i}'\mathbf{i} - \mathbf{i}'\mathbf{x}\mathbf{x}'\mathbf{i}} \begin{bmatrix} \mathbf{x}'\mathbf{x} & -\mathbf{i}'\mathbf{x} \\ -\mathbf{x}'\mathbf{i} & \mathbf{i}'\mathbf{i} \end{bmatrix} \\
 &= \frac{1}{N(\mathbf{x}'\mathbf{x} - N\bar{x}^2)} \begin{bmatrix} \mathbf{x}'\mathbf{x} & -N\bar{x} \\ -N\bar{x} & N \end{bmatrix} \\
 &= \begin{bmatrix} \frac{\frac{1}{N}\mathbf{x}'\mathbf{x}}{(\mathbf{x}'\mathbf{x} - N\bar{x}^2)} & \frac{-\bar{x}}{(\mathbf{x}'\mathbf{x} - N\bar{x}^2)} \\ \frac{-\bar{x}}{(\mathbf{x}'\mathbf{x} - N\bar{x}^2)} & \frac{1}{(\mathbf{x}'\mathbf{x} - N\bar{x}^2)} \end{bmatrix}
 \end{aligned}$$

$$\text{If } \bar{x} = 0, \quad (\mathbf{X}'\mathbf{X})^{-1} = \begin{bmatrix} \frac{1}{N} & 0 \\ 0 & \frac{1}{(\mathbf{x}'\mathbf{x})} \end{bmatrix}$$

$$\mathbf{y} = \beta_0 \mathbf{i} + \beta_1 \mathbf{x} + \boldsymbol{\epsilon}$$

$$\mathbf{y} = \mathbf{X}\boldsymbol{\beta} + \boldsymbol{\epsilon}$$

$$\mathbf{b} = \begin{bmatrix} b_0 \\ b_1 \end{bmatrix} = \begin{bmatrix} \frac{1}{\mathbf{i}'\mathbf{i}} [\mathbf{i}'\mathbf{y} - \frac{\tilde{\mathbf{x}}'\mathbf{y}}{\tilde{\mathbf{x}}'\tilde{\mathbf{x}}} \mathbf{i}'\mathbf{x}] \\ \frac{\tilde{\mathbf{x}}'\mathbf{y}}{\tilde{\mathbf{x}}'\tilde{\mathbf{x}}} \end{bmatrix} \quad \text{if } \bar{x} = 0, = \begin{bmatrix} \frac{1}{N} [\mathbf{i}'\mathbf{y}] \\ \frac{\mathbf{x}'\mathbf{y}}{\mathbf{x}'\mathbf{x}} \end{bmatrix}$$

$$\mathbf{b} = (\mathbf{X}'\mathbf{X})^{-1} \mathbf{X}'\mathbf{y}$$

$$\mathbf{y} = \mathbf{X}(\mathbf{X}'\mathbf{X})^{-1} \mathbf{X}'\mathbf{y} + \mathbf{e}$$

$$\mathbf{y} - \mathbf{X}(\mathbf{X}'\mathbf{X})^{-1} \mathbf{X}'\mathbf{y} = \mathbf{e}$$

$$(\mathbf{I} - \mathbf{X}(\mathbf{X}'\mathbf{X})^{-1} \mathbf{X}')\mathbf{y} = \mathbf{e}$$

$$(\mathbf{I} - \mathbf{P})\mathbf{y} = \mathbf{e}$$

$$\mathbf{M}\mathbf{y} = \mathbf{e}$$

Define the projection matrix

Define the annihilator matrix

P , M Rules

- $P = X(X'X)^{-1}X'$ is called the projection matrix or hat matrix.
- $I - P = M$ is called the annihilator matrix or the residual maker.
- $MM' = MM = M$, that is, M is symmetric and idempotent.
- $PP' = PP = P$, that is, P is symmetric and idempotent.
- $PX = X$, that is, X is invariant under P
- $My = e$
- $E(e) = ME(y)$
- $Var(e) = MVar(y)M' = \sigma^2 M$

Trace Tricks

- The trace is the sum of the diagonals of a matrix.
- $tr(c) = c$, which means given a vector \mathbf{y} , $\mathbf{y}'\mathbf{y} = tr(\mathbf{y}'\mathbf{y})$.
- The trace of $P = \mathbf{X}(\mathbf{X}'\mathbf{X})^{-1}\mathbf{X}' = rank(\mathbf{X})$.
- Given a vector of random variables \mathbf{x} ,

$$E[\mathbf{x}'\mathbf{x}] = tr[V(\mathbf{x})] + E(\mathbf{x})'E(\mathbf{x})$$

Expectation of Sum of Squares

$$\begin{aligned}E[\mathbf{e}'\mathbf{e}] &= E[\text{tr}(\mathbf{e}\mathbf{e}')] \\&= \text{tr}[E(\mathbf{e}\mathbf{e}')] \\&= \text{tr}[\text{Var}(\mathbf{e}) + E(\mathbf{e})'E(\mathbf{e})] \\&= \text{tr}[\text{Var}(\mathbf{e})] + E(\mathbf{e})'E(\mathbf{e}) \\&= \text{tr}[\text{Var}(\mathbf{e})] + 0 \quad E(\mathbf{e}) = 0 \\&= \text{tr}[\sigma^2 \mathbf{M}] \\&= \sigma^2 \text{tr}[(\mathbf{I} - \mathbf{P})] \\&= \sigma^2 [N - \text{rank}(\mathbf{P})] \\&= \sigma^2 (N - 2)\end{aligned}$$

trace is a linear operator

Estimating variance of errors

$$E[\sum e_i^2] = E[\mathbf{e}'\mathbf{e}] = \sigma^2(N-2)$$

$$s_e^2 = \frac{\sum_{i=1}^N e_i^2}{(N-2)}$$

This is an unbiased estimator:

$$E[s_e^2] = E\left[\frac{\sum e_i^2}{(N-2)}\right] = \frac{E[\mathbf{e}'\mathbf{e}]}{(N-2)} = \frac{\sigma^2(N-2)}{(N-2)}$$

```
> x <- 1:50
> y <- 8 + .5 * x + rnorm(50, sd=3)

> vcov(lm(y~x))
              (Intercept)                x
(Intercept)  0.69874415  -0.0207547766
x            -0.02075478   0.0008139128

> xtilde <- x - mean(x)

> sum(resid(lm(y~x))^2)/48 / xtilde%*%xtilde
      [,1]
[1,] 0.0008139128
```

Next time: Statistics