

Welcome

Course: High-Dimensional Econometrics

Instructor: [Dr David T. Jacho-Chavez](#)

Start Date: Monday, January 6th, 2025 **End Date:** Friday, January 10th, 2025

Server: <https://server.davidjachochavez.org>

GitHub Repository: <https://github.com/DTJCM/LHD-Metrics>



Date	Schedule
Monday, January 6th	02:00 pm - 06:00 pm
Tuesday, January 7th	08:00 am - 11:45 am 02:15 pm - 06:00 pm
Wednesday, January 8th	08:00 am - 11:45 am 02:15 pm - 06:00 pm
Thursday, January 9th	08:00 am - 11:45 am 02:15 pm - 06:00 pm
Friday, January 10th	08:00 am - 11:45 am 02:15 pm - 06:00 pm

Disclaimer

All the materials in this course are taken from these sources and the original authors hold *all* text, figures, and computer code copyrights:

1. Jeffrey M. Wooldridge's '[Introductory Econometrics: A Modern Approach](#)' textbook.
2. Bruce Hansen's '[Econometrics](#)' online textbook.
3. Sergio Correia's [various](#) Stata packages.
4. Fernando Rios-Avila's [various](#) Stata packages.
5. Asjad Naqvi's [stata-schemepack](#) package.
6. Mingze Gao's [specurve](#) package.
7. Various chapters of '[An Introduction to Statistical Learning: With Applications in R](#)' book by Gareth James, Daniela Witten, Trevor Hastie and Rob Tibshirani.
8. David Drukker and Di Lu's '[An introduction to the lasso in Stata](#)'.
9. StataCorp. 2021. '[Stata Statistical Software: Release 17](#).' College Station, TX: StataCorp LLC.

Exam

Date: Sunday, January 12th, 2025

Time: 09:00 am - 09:00 pm [12 hours]

Format: Electronic - Jupyter Notebook - 100%

- Part 1: Theory - 60% [32% corresponds to questions in [Exercises](#)]
- Part 2: Practice - 40%

Final Score Calculation: $(\text{Exam Score}/100) \times 18 + (\text{Data Camp points})$

Guidelines:

- You can use your own notes, books, the internet, artificial intelligence, etc.
- The top-2 final scores among those who work individually will earn my recommendation letter for graduate school for up to 24 months, i.e., January 2027.
- Submit your `student_code.ipynb` to [Manuel Vasquez](#) [email: manuel.vasquez@bcrp.gob.pe] by 09:00 pm on Sunday, January 12th, 2025.

Emory University

Source: [US News - 2025 Best National Universities](https://www.usnews.com/best-colleges/rankings/national-universities)

The screenshot shows a web browser displaying the US News & World Report website for the 2025 Best National Universities rankings. The page features a header with the US News logo and navigation links. Below the header, there's a section titled "Best National University Rankings" with a brief description of the category. A sidebar on the right lists the top 30 universities. The main content area displays a card for Princeton University, including its rank, name, location, a photo of the campus, and a summary of its history and policies. There are also filters and comparison tools available.

Rank	University
1.	Princeton University
2.	MIT
3.	Harvard University
4.	Stanford University
5.	Yale University
6.	California Institute of Technology
7.	Duke University
8.	Johns Hopkins University
9.	Northwestern University
10.	University of Pennsylvania
11.	Cornell University
12.	University of Chicago
13.	Brown University
14.	Columbia University
15.	Dartmouth College
16.	UCLA
17.	Berkeley
18.	Rice University
19.	University of Notre Dame
20.	Vanderbilt University
21.	Carnegie Mellon University
22.	University of Michigan
23.	Washington University in St. Louis
24.	Emory University
25.	Georgetown University
26.	University of Virginia
27.	University of North Carolina - Chapel Hill
28.	University of Southern California
29.	University of California, San Diego
30.	New York University

1. Princeton University
2. MIT
3. Harvard University
4. Stanford University
5. Yale University
6. California Institute of Technology
7. Duke University
8. Johns Hopkins University
9. Northwestern University
10. University of Pennsylvania
11. Cornell University
12. University of Chicago
13. Brown University
14. Columbia University
15. Dartmouth College
16. UCLA
17. Berkeley
18. Rice University
19. University of Notre Dame
20. Vanderbilt University
21. Carnegie Mellon University
22. University of Michigan
23. Washington University in St. Louis
24. **Emory University**
25. Georgetown University
26. University of Virginia
27. University of North Carolina - Chapel Hill
28. University of Southern California
29. University of California, San Diego
30. New York University

LinkedIn



EMORY

Department
of Economics

Follow us on LinkedIn



emory-university-
department-of-
economics

Regression

1. Find $\mathbb{E}[\mathbb{E}[Y | X_1, X_2, X_3] | X_1, X_2] | X_1]$

2. Suppose that the random variables Y and X only take the values 0 and 1, and have the following joint probability distribution

	$X = 0$	$X = 1$
$Y = 0$.1	.2
$Y = 1$.4	.3

Find $\mathbb{E}[Y | X]$, $\mathbb{E}[Y^2 | X]$ and $\text{var}[Y | X]$ for $X = 0$ and $X = 1$.

3. True or False?

- (a) If $Y = X\beta + e$, $X \in \mathbb{R}$, and $\mathbb{E}[e | X] = 0$, then $\mathbb{E}[X^2e] = 0$.
- (b) If $Y = X\beta + e$, $X \in \mathbb{R}$, and $\mathbb{E}[Xe] = 0$, then $\mathbb{E}[X^2e] = 0$.
- (c) If $Y = X'\beta + e$ and $\mathbb{E}[e | X] = 0$, then e is independent of X .
- (d) If $Y = X'\beta + e$ and $\mathbb{E}[Xe] = 0$, then $\mathbb{E}[e | X] = 0$.
- (e) If $Y = X'\beta + e$, $\mathbb{E}[e | X] = 0$, and $\mathbb{E}[e^2 | X] = \sigma^2$, then e is independent of X .

4. Take the homoskedastic model

$$\begin{aligned} Y &= X'_1\beta_1 + X'_2\beta_2 + e \\ \mathbb{E}[e | X_1, X_2] &= 0 \\ \mathbb{E}[e^2 | X_1, X_2] &= \sigma^2 \\ \mathbb{E}[X_2 | X_1] &= \Gamma X_1 \end{aligned} .$$

Assume $\Gamma \neq 0$. Suppose the parameter β_1 is of interest. We know that the exclusion of X_2 creates omitted variable bias in the projection coefficient on X_2 . It also changes the equation error. Our question is: what is the effect on the homoskedasticity property of the induced equation error? Does the exclusion of X_2 induce heteroskedasticity or not? Be specific.

5. Let $\hat{\beta}_n = (\mathbf{X}'_n \mathbf{X}_n)^{-1} \mathbf{X}'_n \mathbf{Y}_n$ denote the OLS estimate when \mathbf{Y}_n is $n \times 1$ and \mathbf{X}_n is $n \times k$. A new observation (Y_{n+1}, X_{n+1}) becomes available. Prove that the OLS estimate computed using this additional observation is

$$\hat{\beta}_{n+1} = \hat{\beta}_n + \frac{1}{1 + X'_{n+1}(\mathbf{X}'_n \mathbf{X}_n)^{-1} X_{n+1}} (\mathbf{X}'_n \mathbf{X}_n)^{-1} X_{n+1} (Y_{n+1} - X'_{n+1} \hat{\beta}_n).$$

6. Prove that R^2 is the square of the sample correlation between \mathbf{Y} and $\hat{\mathbf{Y}}$.

7. For the intercept-only model $Y_i = \beta + e_i$, show that the leave-one-out prediction error is

$$\tilde{e}_i = \left(\frac{n}{n-1} \right) (Y_i - \bar{Y}).$$

8. The observations are $(Y_i, X_{1i}, X_{2i}), i = 1, \dots, n$. You estimate two least squares regressions.

$$\begin{aligned} Y_i &= X'_{1i} \tilde{\beta}_1 + \tilde{e}_i \\ Y_i &= X'_{1i} \hat{\beta}_1 + X'_{2i} \hat{\beta}_2 + \hat{e}'_i \end{aligned}$$

and calculate the residual variance estimates

$$\begin{aligned}\tilde{\sigma}^2 &= \frac{1}{n} \sum_{i=1}^n \tilde{e}_i^2 \\ \hat{\sigma}^2 &= \frac{1}{n} \sum_{i=1}^n \hat{e}_i^2.\end{aligned}$$

Show that for any $w \in (0,1)$, there is a constant $a \in (0,1)$ such that

$$\frac{1}{n} \sum_{i=1}^n (w\hat{e}_i + (1-w)\tilde{e}_i)^2 = (1-a)\hat{\sigma}^2 + a\tilde{\sigma}^2.$$

Find the constant a .

9. Consider a regression of a $n \times 1$ vector of responses \mathbf{Y} on a $n \times k$ matrix of explanatory variables \mathbf{X} . Let $\mathbf{M} = \mathbf{I}_n - \mathbf{P}$, where \mathbf{P} is the projection matrix $\mathbf{X}(\mathbf{X}'\mathbf{X})^{-1}\mathbf{X}'$, and \mathbf{I}_n is the identity matrix of order n . Suppose one adds a new explanatory variable, \mathbf{Z} , to the regression, so there are now $(k+1)$ regressors. Show that the new residual sum of squares is given by $\mathbf{Y}'\mathbf{M}\mathbf{Y} - b^2\mathbf{Z}'\mathbf{M}\mathbf{Z}$ where $b = \mathbf{Z}'\mathbf{M}\mathbf{Y}/\mathbf{Z}'\mathbf{M}\mathbf{Z}$ is the OLS coefficient on \mathbf{Z} .

Regression (Cont.)

1. For some integer k , set $\mu_k = \mathbb{E}[Y^k]$.

(a) Construct an estimator $\hat{\mu}_k$ for μ_k .

(b) Show that $\hat{\mu}_k$ is unbiased for μ_k .

(c) Calculate the variance of $\hat{\mu}_k$, say $\text{var}[\hat{\mu}_k]$. What assumption is needed for $\text{var}[\hat{\mu}_k]$ to be finite?

(d) Propose an estimator of $\text{var}[\hat{\mu}_k]$.

2. Let $\mu = \mathbb{E}[Y]$, $\sigma^2 = \mathbb{E}[(Y - \mu)^2]$ and $\mu_3 = \mathbb{E}[(Y - \mu)^3]$ and consider the sample mean $\bar{Y} = \frac{1}{n} \sum_{i=1}^n Y_i$. Find $\mathbb{E}[(\bar{Y} - \mu)^3]$ as a function of μ , σ^2 , μ_3 and n .

3. Now assume that

$$\sqrt{n} \begin{pmatrix} \hat{\theta}_1 - \theta_1 \\ \hat{\theta}_2 - \theta_2 \end{pmatrix} \xrightarrow{d} N(0, \Sigma)$$

where $\|\Sigma\| < \infty$, then find the asymptotic distribution of the following statistics

(a) $\hat{\theta}_1 \hat{\theta}_2$

(b) $\exp(\hat{\theta}_1 + \hat{\theta}_2)$

(c) If $\theta_2 \neq 0$, $\hat{\theta}_1 / \hat{\theta}_2^2$

(d) $\hat{\theta}_1^3 + \hat{\theta}_1 \hat{\theta}_2^2$

4. Take a regression model $Y = X\beta + e$ with $\mathbb{E}[e | X] = 0$ and i.i.d. observations (Y_i, X_i) and scalar X . The parameter of interest is $\theta = \beta^2$. Consider the OLS estimators $\hat{\beta}$ and $\hat{\theta} = \hat{\beta}^2$.

(a) Find $\mathbb{E}[\hat{\theta} | X]$ using our knowledge of $\mathbb{E}[\hat{\beta} | X]$ and $V_{\hat{\beta}} = \text{var}[\hat{\beta} | X]$. Is $\hat{\theta}$ biased for θ ?

(b) Suggest an (approximate) biased-corrected estimator $\hat{\theta}^*$ using an estimator $\hat{V}_{\hat{\beta}}$ for $V_{\hat{\beta}}$.

(c) For $\hat{\theta}^*$ to be potentially unbiased, which estimator of $V_{\hat{\beta}}$ is most appropriate?

Under which conditions is $\hat{\theta}^*$ unbiased?

5. Consider an i.i.d. sample $\{Y_i, X_i\}_{i=1, \dots, n}$ where X_i is $k \times 1$. Assume the linear conditional expectation model $Y = X'\beta + e$ with $\mathbb{E}[e | X] = 0$. Assume that $n^{-1}X'X = I_k$ (orthonormal regressors). Consider the OLS estimator $\hat{\beta}$.

(a) Find $V_{\hat{\beta}} = \text{var}[\hat{\beta}]$

(b) In general, are $\hat{\beta}_j$ and $\hat{\beta}_\ell$ for $j \neq \ell$ correlated or uncorrelated?

(c) Find a sufficient condition so that $\hat{\beta}_j$ and $\hat{\beta}_\ell$ for $j \neq \ell$ are uncorrelated.

6. Take the model in vector notation

$$\begin{aligned} Y &= X\beta + e \\ \mathbb{E}[e | X] &= 0 \\ \mathbb{E}[ee' | X] &= \Omega \end{aligned}$$

Assume for simplicity that Ω is known. Consider the OLS and GLS estimators $\hat{\beta} = (X'X)^{-1}(X'Y)$ and $\tilde{\beta} = (X'\Omega^{-1}X)^{-1}(X'\Omega^{-1}Y)$. Compute the (conditional) covariance between $\hat{\beta}$ and $\tilde{\beta}$:

$$\mathbb{E}[(\hat{\beta} - \beta)(\tilde{\beta} - \beta)' | X].$$

Find the (conditional) covariance matrix for $\hat{\beta} - \tilde{\beta}$:

$$\mathbb{E}[(\hat{\beta} - \tilde{\beta})(\hat{\beta} - \tilde{\beta})' | X].$$

7. The model is $Y = X\beta + e$ with $\mathbb{E}[e | X] = 0$ and $X \in \mathbb{R}$. Consider the two estimators

$$\begin{aligned}\hat{\beta} &= \frac{\sum_{i=1}^n X_i Y_i}{\sum_{i=1}^n X_i^2} \\ \tilde{\beta} &= \frac{1}{n} \sum_{i=1}^n \frac{Y_i}{X_i}\end{aligned}$$

(a) Under the stated assumptions are both estimators consistent for β ?

(b) Are there conditions under which either estimator is efficient?

8. Take the linear model $Y = X\beta + e$ with $\mathbb{E}[e | X] = 0$ and $X_i \in \mathbb{R}$. Consider the estimator

$$\hat{\beta} = \frac{\sum_{i=1}^n X_i^3 Y_i}{\sum_{i=1}^n X_i^4}.$$

Find the asymptotic distribution of $\sqrt{n}(\hat{\beta} - \beta)$ as $n \rightarrow \infty$.

9. The model is $y = X'\beta + e$ with $\mathbb{E}[e | X] = 0$. An econometrician is worried about the impact of some unusually large values of the regressors. The model is thus estimated on the subsample for which $|X_i| \leq c$ for some fixed c . Let $\tilde{\beta}$ denote the OLS estimator on this subsample. It equals

$$\tilde{\beta} = \left(\sum_{i=1}^n X_i X_i' \mathbb{1}\{|X_i| \leq c\} \right)^{-1} \left(\sum_{i=1}^n X_i Y_i \mathbb{1}\{|X_i| \leq c\} \right).$$

(a) Show that $\tilde{\beta} \xrightarrow{p} \beta$.

(b) Find the asymptotic distribution of $\sqrt{n}(\tilde{\beta} - \beta)$.

Binary Response

1. For a binary response y , let \bar{y} be the proportion of ones in the sample (which is equal to the sample average of the y_i). Let \hat{q}_0 be the percent correctly predicted for the outcome $y = 0$ and let \hat{q}_1 be the percent correctly predicted for the outcome $y = 1$. If \hat{p} is the overall percent correctly predicted, show that \hat{p} is a weighted average of \hat{q}_0 and \hat{q}_1 :

$$\hat{p} = (1 - \bar{y})\hat{q}_0 + \bar{y}\hat{q}_1.$$

In a sample of 300, suppose that $\bar{y} = .70$, so that there are 210 outcomes with $y_i = 1$ and 90 with $y_i = 0$. Suppose that the percent correctly predicted when $y = 0$ is 80, and the percent correctly predicted when $y = 1$ is 40. Find the overall percent correctly predicted.

2. Emily estimates a probit regression setting her dependent variable to equal $Y = 1$ for a purchase and $Y = 0$ for no purchase. Using the same data and regressors, Jacob estimates a probit regression setting the dependent variable to equal $Y = 1$ if there is no purchase and $Y = 0$ for a purchase. What is the difference in their estimated slope coefficients?
3. Jackson estimates a logit regression where the primary regressor is measured in dollars. Julie estimates a logit regression with the same sample and dependent variable, but measures the primary regressor in thousands of dollars. What is the difference in the estimated slope coefficients?
4. For the logistic distribution $\Lambda(x) = (1 + \exp(-x))^{-1}$ verify that
- (a) $\frac{d}{dx}\Lambda(x) = \Lambda(x)(1 - \Lambda(x))$.
 - (b) $h_{\text{logit}}(x) = \frac{d}{dx}\log\Lambda(x) = 1 - \Lambda(x)$.
 - (c) $H_{\text{logit}}(x) = -\frac{d^2}{dx^2}\log\Lambda(x) = \Lambda(x)(1 - \Lambda(x))$.
 - (d) $|H_{\text{logit}}(x)| \leq 1$.
5. For the normal distribution $\Phi(x)$ verify that
- (a) $h_{\text{probit}}(x) = \frac{d}{dx}\log\Phi(x) = \lambda(x)$ where $\lambda(x) = \phi(x)/\Phi(x)$.
 - (b) $H_{\text{probit}}(x) = -\frac{d^2}{dx^2}\log\Phi(x) = \lambda(x)(x + \lambda(x))$.

Multinomial Response

1.

- a. For estimating the mean of a nonnegative random variable y using a random sample $\{y_i\}_{i=1}^n$, the Poisson quasi-log likelihood for a random draw is $\gamma_i(\mu) = y_i \log(\mu) - \mu, \mu > 0$ (where terms not depending on μ have been dropped). Letting $\mu_0 \equiv \mathbb{E}(y_i)$, we have $\mathbb{E}[\ell_i(\mu)] = \mu_0 \log(\mu) - \mu$. Show that this function is uniquely maximized at $\mu = \mu_0$.
 - b. The gamma (exponential) quasi-log likelihood is $\ell_i(\mu) = -y_i/\mu - \log(\mu), \mu > 0$. Show that $\mathbb{E}[\ell_i(\mu)]$ is uniquely maximized at $\mu = \mu_0$.
2. In a balanced panel data setting, i.e. $\{\{y_{it}, \mathbf{x}_{it}\}_{t=1}^T\}_{i=1}^n$, consider an unobserved effects model for count data with exponential regression function

$$\mathbb{E}(y_{it} | \mathbf{x}_{i1}, \dots, \mathbf{x}_{iT}, c_i) = c_i \exp(\mathbf{x}_{ii} \boldsymbol{\beta}).$$

If $\mathbb{E}(c_i | \mathbf{x}_{i1}, \dots, \mathbf{x}_{iT}) = \exp(\alpha + \bar{\mathbf{x}}_i \boldsymbol{\gamma})$, find $\mathbb{E}(y_{it} | \mathbf{x}_{i1}, \dots, \mathbf{x}_{iT})$ where $\bar{\mathbf{x}}_i = T^{-1}(\mathbf{x}_{i1} + \dots + \mathbf{x}_{iT})$.

3. Let patents be the number of patents applied for by a firm during a given year. Assume that the conditional expectation of patents given sales and RD is

$$\mathbb{E}(\text{patents} | \text{sales}, RD) = \exp[\beta_0 + \beta_1 \log(\text{sales}) + \beta_2 RD + \beta_3 RD^2],$$

where sales is annual firm sales and RD is total spending on research and development over the past 10 years.

- a. How would you estimate the β_j ? Justify your answer by discussing the nature of patents.
- b. How do you interpret β_1 ?
- c. Find the partial effect of RD on $\mathbb{E}(\text{patents} | \text{sales}, RD)$.

High-Dimensional Models

1. Take the model $Y = X'\beta + e$ with $\mathbb{E}[Xe] = 0$. Define the ridge regression estimator

$$\hat{\beta} = \left(\sum_{i=1}^n X_i X_i' + \lambda I_k \right)^{-1} \left(\sum_{i=1}^n X_i Y_i \right),$$

here $\lambda > 0$ is a fixed constant.

- (a) Find $\mathbb{E}[\hat{\beta} | X]$.
 - (b) Is $\hat{\beta}$ biased for β ?
 - (c) Find the probability limit of $\hat{\beta}$ as $n \rightarrow \infty$.
 - (d) Is $\hat{\beta}$ consistent for β ?
 - (e) Derive $\text{var}[\hat{\beta} | X]$.
2. Show that the ridge regression estimator can be computed as least squares applied to an augmented data set. Take the original data (Y, X) . Add p 0's to Y and p rows of $\sqrt{\lambda} I_p$ to X , apply least squares, and show that this equals $\hat{\beta}_{\text{ridge}}$.
3. Which estimator produces a higher regression R^2 , least squares or ridge regression?

Quick Reminders

The **proportionate change** in v moving from v_0 to v_1 is

$$\frac{v_1 - v_0}{v_0} = \Delta v/v_0,$$

assuming of course that $v_0 \neq 0$. The **percentage change** in v moving from v_0 to v_1 is simply 100 times the proportionate change:

$$\% \Delta v = 100(\Delta v/v_0),$$

the notation "% Δv " is read as "the percentage change in v ."

Let v_0 and v_1 be positive values, then, it can be shown that

$$\log(v_1) - \log(v_0) \approx \frac{v_1 - v_0}{v_0} = \Delta v/v_0,$$

for small changes in v . Therefore by writing $\Delta \log(v) = \log(v_1) - \log(v_0)$, then

$$100 \cdot \Delta \log(v) \approx \% \Delta v.$$

Definition (Elasticity)

The **elasticity** of y with respect to x is defined as

$$\frac{\Delta y}{\Delta x} \cdot \frac{x}{y} = \frac{\% \Delta y}{\% \Delta x}.$$

 The elasticity of y with respect to x is the percentage change in y when x increases by 1%.

Consider the relationship

$$\log(y) = \beta_0 + \beta_1 \log(x),$$

Then

$$\beta_1 = \frac{\Delta \log(y)}{\Delta \log(x)} = \frac{100 \times \Delta \log(y)}{100 \times \Delta \log(x)} \approx \frac{\% \Delta y}{\% \Delta x}.$$

 The slope parameter β_1 is the elasticity of y with respect to x [assuming that $x, y > 0$].

Constant Elasticity Demand Function:

If q is quantity demanded and p is price and these variables are related by

$$\log(q) = 4.7 - 1.25 \log(p),$$

then the *price elasticity of demand* is -1.25 . Roughly, a 1% increase in price leads to a 1.25% fall in the quantity demanded.

Suppose that $y > 0$ and

$$\log(y) = \beta_0 + \beta_1 x,$$

Then $\Delta \log(y) = \beta_1 \Delta x$, so

$$100 \times \Delta \log(y) = (100 \times \beta_1) \Delta x$$

$$\% \Delta y \approx (100 \times \beta_1) \Delta x$$

$$\frac{\% \Delta y}{\Delta x} \approx (100 \times \beta_1)$$

 The slope parameter β_1 is the semi-elasticity of y with respect to x [assuming that $y > 0$].

Logarithmic Wage Equation:

Suppose that hourly wage and years of education are related by

$$\log(wage) = 2.78 + .094 educ.$$

Then,

$$\% \Delta wage \approx 100(.094) \Delta educ = 9.4 \Delta educ.$$

It follows that one more year of education increases hourly wage by about 9.4%.

Suppose that $x > 0$ and

$$y = \beta_0 + \beta_1 \log(x),$$

Then $\Delta y = \beta_1 \Delta \log(x)$, so

$$\Delta y = (\beta_1/100) [100 \times \Delta \log(x)]$$

$$\Delta y \approx (\beta_1/100) \% \Delta x$$

Q $(\beta_1/100)$ is the unit change in y when x increases by 1%. [assuming that $x > 0$].

Labor Supply Equation:

Assume that the labor supply of a worker can be described by

$$hours = 33 + 45.1 \log(wage),$$

where $wage$ is hourly wage and $hours$ is hours worked per week. Then,

$$\Delta hours \approx \left(\frac{45.1}{100}\right) (\% \Delta wage) = .451 \% \Delta wage.$$

In other words, a 1 % increase in wage increases the weekly hours worked by about .45, or slightly less than one-half hour.

C If the wage increases by 10%, then $\Delta hours = .451(10) = 4.51$, or about four and one-half hours.

Suppose that $1_A(x)$ represents the indicator function that equals one if $x \in A$ and zero otherwise, and consider the linear relationship

$$y = \beta_0 + \beta_1 1_A(x),$$

Then for pairs $(y_2, x \in A)$ and $(y_1, x \notin A)$, we have $\Delta y = \beta_1$.

Q β_1 is the unit change in y when x belongs to A .

Logarithmic Wage Equation:

Suppose that hourly wage and the indicator for male are related by

$$\log(wage) = 2.78 + .094 male.$$

Then,

$$\% \Delta wage \approx 100(.094) = 9.4.$$

It follows that men earn roughly 9.4% more than women.

Law of Iterated Expectations

Theorem: Simple Law of Iterated Expectations If $\mathbb{E}|Y| < \infty$ then for any random vector X ,

$$\mathbb{E}(\mathbb{E}(Y | X)) = \mathbb{E}(Y)$$

▲ The average of the conditional averages is the unconditional average.

Theorem: Law of Iterated Expectations

If $\mathbb{E}|Y| < \infty$ then for any random vectors X_1 and X_2 ,

$$\mathbb{E}[\mathbb{E}[Y | X_1, X_2] | X_1] = \mathbb{E}[Y | X_1]$$

▲ “The smaller information set wins.”

Theorem: Conditioning Theorem

If $\mathbb{E}|Y| < \infty$ then

$$\mathbb{E}[g(X)Y | X] = g(X)\mathbb{E}[Y | X]$$

If in addition $\mathbb{E}|g(X)| < \infty$ then

$$\mathbb{E}[g(X)Y] = \mathbb{E}[g(X)\mathbb{E}[Y | X]]$$

▲ A property of conditional expectations is that when you condition on a random vector X you can effectively treat it as if it is constant.

Linear Projection Model

Linear Projection Model

$$Y = X'\beta + e$$
$$\mathbb{E}[Xe] = 0$$
$$\beta = (\mathbb{E}[XX'])^{-1}\mathbb{E}[XY]$$

Recall:

Linear CEF Model

$$Y = X'\beta + e$$
$$\mathbb{E}[e | X] = 0$$

Homoskedastic Linear CEF Model

$$Y = X'\beta + e$$
$$\mathbb{E}[e | X] = 0$$
$$\mathbb{E}[e^2 | X] = \sigma^2$$

β is known as the *linear projection coefficient*.

Invertibility & Identification

⚠️ Identification means **uniqueness**!

Invertibility and Identification

The linear projection coefficient $\beta = (\mathbb{E}[XX'])^{-1}\mathbb{E}[XY]$ exists and is unique as long as the $k \times k$ matrix $\mathbf{Q}_{XX} = \mathbb{E}[XX']$ is invertible. The matrix \mathbf{Q}_{XX} is often called the design matrix as in experimental settings the researcher is able to control \mathbf{Q}_{XX} by manipulating the distribution of the regressors X . Observe that for any non-zero $\alpha \in \mathbb{R}^k$

$$\alpha' \mathbf{Q}_{XX} \alpha = \mathbb{E}[\alpha' XX' \alpha] = \mathbb{E}[(\alpha' X)^2] \geq 0$$

so \mathbf{Q}_{XX} by construction is positive semi-definite, conventionally written as $\mathbf{Q}_{XX} \geq 0$. The assumption that it is positive definite means that this is a strict inequality, $\mathbb{E}[(\alpha' X)^2] > 0$. This is conventionally written as $\mathbf{Q}_{XX} > 0$. This condition means that there is no non-zero vector α such that $\alpha' X = 0$ identically. Positive definite matrices are invertible. Thus when $\mathbf{Q}_{XX} > 0$ then $\beta = (\mathbb{E}[XX'])^{-1}\mathbb{E}[XY]$ exists and is uniquely defined. In other words, if we can exclude the possibility that a linear function of X is degenerate, then β is uniquely defined.

We have shown that the linear projection coefficient β is identified (uniquely determined) under Assumption 1. **The key is invertibility of \mathbf{Q}_{XX} .** Otherwise, there is no unique solution to the equation

$$\mathbf{Q}_{XX}\beta = \mathbf{Q}_{XY}$$

when \mathbf{Q}_{XX} is not invertible there are multiple solutions. In this case the coefficient β is not identified as it does not have a unique value.

💡 Although $\mathbf{Q}_{XX}\beta = \mathbf{Q}_{XY}$ has multiple solutions, it turns out that one of these solutions denoted as $\mathbf{Q}_{XX}^\perp \mathbf{Q}_{XY}$ [A^\perp denoting the [Moore-Penrose Inverse](#)] is such that $\|\mathbf{Q}_{XX}^\perp \mathbf{Q}_{XY}\|_2 \leq \|b\|_2$ for all solutions b , i.e., [minimum norm solution to a linear system](#).

Least Squares

The best linear predictor of Y given X for a pair of random variables $(Y, X) \in \mathbb{R} \times \mathbb{R}^k$, is summarized in the linear projection model. We are now interested in *estimating* the parameters of this model, in particular the projection coefficient

$$\beta = (\mathbb{E}(XX'))^{-1}\mathbb{E}(XY).$$

Q We can estimate β from observational data which includes joint measurements on the variables (Y, X) .

Assumption: The variables $\{(Y_1, X_1), \dots, (Y_i, X_i), \dots, (Y_n, X_n)\}$ are identically distributed; they are draws from a common distribution F .

The linear projection model applies to the random observations (Y_i, X_i) . We can write the model as

$$Y_i = X'_i\beta + e_i.$$

where the *linear projection coefficient* β is defined as

$$\beta = \underset{b \in \mathbb{R}^k}{\operatorname{argmin}} S(b),$$

the minimizer of the expected squared error

$$S(\beta) = \mathbb{E}\left((Y_i - X'_i\beta)^2\right),$$

and has the explicit solution

$$\beta = (\mathbb{E}(X_i X'_i))^{-1} \mathbb{E}(X_i Y_i).$$

Moment of functions of random variables: $\mathbb{E}[h(X, Y; \beta)]$.

$$h(X, Y; \beta) = (Y - X'\beta)^2,$$

$$h(X, Y; \beta) = XX',$$

$$h(X, Y; \beta) = XY.$$

Recall that for any random variable, W , its *Cumulative Distribution Function*, $F(t)$, is well-defined as $F(t) = \Pr\{W \leq t\}$. Then, if $\mathbb{E}|h(w)| < \infty$, one has $\mathbb{E}[h(W)] = \int h(W)dF(W)$.

Based on a sample $\{W_i; i = 1, \dots, n\}$, we can estimate this unknown quantity by replacing the unknown CDF with the **empirical distribution function**, $\hat{F}(t) = \frac{1}{n} \sum_{i=1}^n \mathbb{I}(W_i \leq t)$, then a natural 'plug in' estimator of $\mathbb{E}[h(W)]$ is

$$\hat{\mathbb{E}}[h(W)] = \int h(W)d\hat{F}(W).$$

It turns out that through the application of various measurement theory concepts one has

$$\hat{\mathbb{E}}[h(W)] = \frac{1}{n} \sum_{i=1}^n h(W_i).$$



American Economic Review



Econometrica



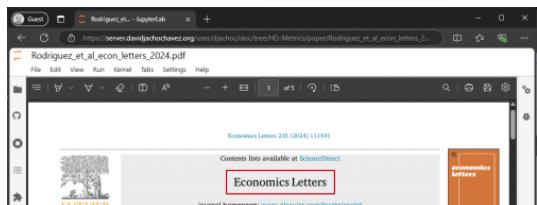
Journal of Political Economy



The Quarterly Journal of Economics



Review of Economic Studies



Economics Letters 235 (2024) 111941
Contents lists available at ScienceDirect
Economics Letters
journal homepage: www.elsevier.com/locate/econle

Abstract readability: Evidence from top-5 economics journals
Belén Rodríguez^{a,*}, Daniel P. Hinsch^b, David T. Jacho-Chávez^c, Leonardo Sánchez-Aragón^c
^aDepartment of Economics, Emory University, 350 Peachtree St., Atlanta, GA 30303-4246, USA
^bBank of Canada, 234 Wellington Ave., Ottawa, ON, K1A 0M2, Canada
^cFacultad de Ciencias Sociales y Humanas, Campus Universitario Avda. 363 3do Piso, 36000 Guadalajara, Jalisco, Mexico

ARTICLE INFO
All classification codes:
CBA
CGE
Economic History
Financial Markets
International Finance
Macroeconomics
Methodology
Microeconomics
Public Economics
Statistics
Latin American Studies

ABSTRACT
Abstract readability is a measure of how easy it is to read a text. Over time, general interest journals have become more technical. This affects how accessible research is to a general audience. Our analysis looks at how readability abstracts are. We study the readability of abstracts of top five economics journals between 2000–2019. We consider two measures of readability: Flesch-Kincaid grade level and log of the number of words per page. We find that higher proportion of women co-authors are more readable. These results are robust to various readability measures and model specification.

1. Introduction
Additionally, we discuss the potential mechanisms of our findings, such as the likelihood for names in macroeconomics and microeconomics.

Rodríguez_et_al_econ_letters_2024.pdf

Data
Rodríguez, Belén; Kim, Eunji; Hinsch, Daniel; Jacho-Chávez, David; Sánchez-Aragón, Leonardo; "Abstract Readability: Evidence from Top-5 Economics Journals"; *Economics Letters*; 111941; Manuscript; 2024-01-18; Materials; GitHub Repository

```
(1): layout panes as pd.set_option('display.max_rows', None)
pd.set_option('display.max_colwidth', 150)
import pyreadstat

# Load the Stats dataset
data, meta = pyreadstat.read_dta("../data/meta.dta")

# Create a DataFrame from the variable names and labels
var_labels = pd.DataFrame({
    "Variable Name": meta.column_names,
    "Variable Label": meta.column_labels
})

# Print the DataFrame
var_labels
```

"The 'JEL' classification system originated with the *Journal of Economic Literature* and is a standard method of classifying scholarly literature in the field of economics. It is used in many of the AEA's published research materials."

Source: [American Economic Association: JEL Guide](#)

```
# Define the list of variables (columns) to include
columns_of_interest = ["log_flesch_kinacid_grade_level", "log_num_authors", "log_num_pages",
                      "both_genders", "prop_women", "journal", "jelcodes", "year", "cluster", "jel_flag"]

# Check if all the columns exist in the dataset
missing_columns = [col for col in columns_of_interest if col not in df.columns]
if missing_columns:
    print(f"The following columns are missing from the DataFrame: {missing_columns}")

# Select the specified columns and sample 5 random rows
random_rows = df[columns_of_interest].sample(n=5, random_state=94)

# Print the result
print(random_rows)

log_flesch_kinacid_grade_level log_num_authors log_num_pages
4707 2.525300 0.000000 3.453987
3695 2.509540 0.000012 3.700390
3697 2.515270 0.000012 3.700390
4551 2.184400 0.450147 2.404987
181 2.059737 0.450147 3.704439

both_genders prop_women journal
4707 0 0.000000 The Society of Economic Theory
3697 0 0.000000 The Quarterly Journal of Economics
4579 1 0.000000 The Review of Economics and Statistics
4551 0 0.000000 The Review of Economic Studies
181 0 0.000000 American Economic Review

cluster year cluster jel_flag
3697 2001 1 1
3697 001.00012644 0123456789 1
4579 2001 1 1
4524 2000 1 1
181 001.00012644 2000 1
```

$$\begin{aligned} X &= [X'_1, X'_2, D', G', d'] \\ X_{4707} &= [0, 3.433987, 0, 0, \text{i}, \text{ECM}, \text{i}, \text{J1}, \text{i}, 2015] \\ G_{4707} &= \text{i}, 1 \\ d_{4707} &= 1 \end{aligned}$$

- i.journal:** corresponds to 4 indicators for ECM, JPE, QJE, RES
- i.jelcodes:** corresponds to 19 indicators for all JEL codes except 'D'
- i.year:** corresponds to 19 indicators for years 2001, 2002,..., 2019
- i.cluster:** corresponds to 214 indicators for clusters 2, 3,..., 215, except cluster '1'
- jel_flag:** indicator of whether the article includes a JEL classification

Y =	log_flesch_kinacid_grade_level
X₁ =	prop_women
X₂ =	[log_num_authors, log_num_page, both_genders, 1]
D =	[i.journal, i.jelcodes, i.year]
G =	i.cluster
d =	jel_flag

Base category: AER, 'D' JEL (Microeconomics), year 2000, for which the JEL was imputed

Least Squares Estimator

For given b , the expected squared error is $\mathbb{E}[(Y_i - X'_i b)^2]$. The *moment estimator* of $S(b)$ is the sample average

$$\hat{S}(b) = \frac{1}{n} \sum_{i=1}^n (Y_i - X'_i b)^2,$$

$$= \frac{1}{n} \text{SSE}(b)$$

Sum-of-Squared-Errors function.

⚠ The estimator is also commonly referred to as the **ordinary least-squares (OLS)** estimator. The following notation is also used:

- $\hat{\beta}_{\text{ols}}$
- $\hat{\beta}_n$

The " \wedge " (caret) symbol over the parameter β indicates that $\hat{\beta}$ is a sample estimate of β . Obviously

- $\hat{\beta}$ is a random variable (since it depends on the sample n).
- β is not random but the finite-dimensional vector of parameter of interest.
- $\hat{\beta} \neq \beta$.

Definition: The **least squares estimator** is

$$\hat{\beta} = \underset{b \in \mathbb{R}^k}{\operatorname{argmin}} \hat{S}(b),$$

where $\hat{S}(b) = \frac{1}{n} \sum_{i=1}^n (Y_i - X'_i b)^2$.

Solving for Least Squares

The sum of squared errors can be written as

$$\text{SSE}(b) = \sum_{i=1}^n Y_i^2 - 2b' \sum_{i=1}^n X_i Y_i + b' \sum_{i=1}^n X_i X'_i b$$

This is a **Quadratic Function in b** .

The FOC are

$$0 = \frac{\partial}{\partial b} \text{SSE}(b) \Big|_{b=\hat{\beta}} = -2 \sum_{i=1}^n X_i Y_i + 2 \sum_{i=1}^n X_i X'_i \hat{\beta}$$

⚠ Notice that $\frac{\partial^2}{\partial b \partial b'} \text{SSE}(b) \Big|_{b=\hat{\beta}} = 2 \sum_{i=1}^n X_i X'_i \geq 0$.

The solution for $\hat{\beta}$ may be found by solving the system of k equations

$$\sum_{i=1}^n X_i X'_i \hat{\beta} = \sum_{i=1}^n X_i Y_i.$$

⚠ If you ever need to *numerically* calculate the solution to a system of equations, you **should** use this equation. In Python you should use the [numpy.linalg.solve\(\)](#).

An explicit formula for the least-squares estimator is then

$$\hat{\beta} = \left(\sum_{i=1}^n X_i X'_i \right)^{-1} \left(\sum_{i=1}^n X_i Y_i \right).$$

Least Squares Estimator as a *Moment Estimator*

Population		Sample
------------	--	--------

$$\beta = \left(\mathbb{E}(X_i X'_i) \right)^{-1} \mathbb{E}(X_i Y_i)$$

$$\hat{\beta} = \left(\frac{1}{n} \sum_{i=1}^n X_i X'_i \right)^{-1} \left(\frac{1}{n} \sum_{i=1}^n X_i Y_i \right)$$

$$= \left(\sum_{i=1}^n X_i X'_i \right)^{-1} \left(\sum_{i=1}^n X_i Y_i \right)$$

Theorem: If $\sum_{i=1}^n X_i X'_i > 0$, the least squares estimator equals

$$\hat{\beta} = \left(\sum_{i=1}^n X_i X'_i \right)^{-1} \left(\sum_{i=1}^n X_i Y_i \right).$$

The screenshot shows a Jupyter Notebook interface with a Python kernel. The notebook contains the following code:

```
[1]: %%capture
import stata_setup, os
if os.name == 'nt':
    stata_setup.config('C:/Program Files/Stata17/','mp')
else:
    stata_setup.config('/usr/local/stata17','mp')

We load the data, rename the outcome variable, generate the indicator variables for year and cluster and define local Stata variables called journals and jel_imp which collects all relevant indicators.

[2]: %%stata -qui
use "../data/data", clear
rename log_flesch_kincaid_grade_level FKG
quietly tabulate year, generate(y_)
quietly tabulate cluster, generate(c_)

local journals ecm jpe qje res //AER based category

local jel_imp a_imp b_imp c_imp e_imp f_imp g_imp h_imp i_imp j_imp k_imp ///
l_imp m_imp n_imp o_imp p_imp q_imp r_imp y_imp z_imp // D JEL based case
```

Performing the OLS regression of \mathbf{Y} on \mathbf{X} using Stata :

```
[3]: %%stata -qui
#delimit ;
reg FKG log_num_authors log_num_pages both_genders prop_women
`journals' `jel_imp' y_2-y_20 c_2-c_215 jel_flag, vce(cluster cluster);
matrix b_selected = e(b)[1,"log_num_authors"],e(b)[1,"log_num_pages"],
e(b)[1,"both_genders"],e(b)[1,"prop_women"],e(b)[1,"_cons"];
#delimit cr
```

Printing a subset of the OLS estimate $\hat{\beta}$ (originally a 262×1 vector)

```
[4]: %%stata
matrix list b_selected
```

	log_num_authors	log_num_pages	both_genders	prop_women	_cons
y1	-.00397377	.01915903	.00059809	-.01889331	2.7023992

Least Squares Residuals

Fitted Value	+	Residual	=	Outcome
$\hat{Y}_i = X_i' \hat{\beta}$		$\hat{e}_i = Y_i - \hat{Y}_i = Y_i - X_i' \hat{\beta}$		$Y_i = \hat{Y}_i + \hat{e}_i = X_i' \hat{\beta} + \hat{e}_i$

$$\begin{aligned}
\sum_{i=1}^n X_i \hat{e}_i &= \sum_{i=1}^n X_i (Y_i - X_i' \hat{\beta}) \\
&= \sum_{i=1}^n X_i Y_i - \sum_{i=1}^n X_i X_i' \hat{\beta} \\
&= \sum_{i=1}^n X_i Y_i - \sum_{i=1}^n X_i X_i' \left(\sum_{i=1}^n X_i X_i' \right)^{-1} \left(\sum_{i=1}^n X_i Y_i \right) \\
&= \sum_{i=1}^n X_i Y_i - \sum_{i=1}^n X_i Y_i \\
&= \mathbf{0}
\end{aligned}$$

\triangleleft When X contains a constant, then
 $\frac{1}{n} \sum_{i=1}^n \hat{e}_i = 0$.

Model in Matrix Notation

We can stack these n equations together as

$$Y_1 = X'_1 \beta + e_1$$

$$Y_2 = X'_2 \beta + e_2$$

\vdots

$$Y_n = X'_n \beta + e_n$$

Now define

$$\mathbf{Y} = \begin{pmatrix} Y_1 \\ Y_2 \\ \vdots \\ Y_n \end{pmatrix}, \mathbf{X} = \begin{pmatrix} X'_1 \\ X'_2 \\ \vdots \\ X'_n \end{pmatrix}, \mathbf{e} = \begin{pmatrix} e_1 \\ e_2 \\ \vdots \\ e_n \end{pmatrix}$$

$$\mathbf{Y} = \mathbf{X}\beta + \mathbf{e}$$

$$\sum_{i=1}^n X_i X'_i = \mathbf{X}' \mathbf{X}$$

$$\sum_{i=1}^n X_i Y_i = \mathbf{X}' \mathbf{Y}$$

Least Squares Estimator: $\hat{\beta} = (\mathbf{X}' \mathbf{X})^{-1} (\mathbf{X}' \mathbf{Y})$

Theorem: If $\sum_{i=1}^n X_i X'_i > 0$, the least squares estimator equals

$$\hat{\beta} = \left(\sum_{i=1}^n X_i X'_i \right)^{-1} \left(\sum_{i=1}^n X_i Y_i \right)$$

Fitted Values: $\hat{\mathbf{Y}} = \mathbf{X}\hat{\beta}$

Residuals: $\hat{\mathbf{e}} = \mathbf{Y} - \mathbf{X}\hat{\beta}$

Sum-of-Squared Errors: $SSE(b) = (\mathbf{Y} - \mathbf{X}b)'(\mathbf{Y} - \mathbf{X}b)$

Guest 001_ols.ipynb (2) - JupyterLab https://server.davidjachochavez.org/user/djachoc/doc/tree/HD-Metrics/code/001_ols.ipynb

001_ols.ipynb

File Edit View Run Kernel Tabs Settings Help

Predicting the *fitted values* and the *residuals*, then format them to be displayed with up to 4 decimals only along with other variables.

```
[5]: %%stata -qui
predict FKG_hat, xb
predict double e_hat, residuals
format FKG FKG_hat e_hat log_num_authors log_num_pages %5.4f
```

These command randomly sorts the rows of the data set in memory.

```
[6]: %%stata -qui
#delimit ;
set seed 42; tempvar sortorder; gen `sortorder' = runiform(); sort `sortorder';
#delimit cr
```

```
[7]: %stata list FKG FKG_hat e_hat log_num_authors log_num_pages both_genders prop_women in 1/20, table separator(20)
```

	FKG	FKG_hat	e_hat	log_n~rs	log_~ges	both_g~s	prop_w~n
1.	2.8015	2.7041	0.0973	0.6931	3.5835	1	.5
2.	2.7776	2.7294	0.0482	0.6931	3.7136	0	0
3.	2.7829	2.7017	0.0812	1.0986	3.1781	0	0
4.	2.7027	2.7289	-0.0262	0.0000	3.5553	0	0
5.	2.8278	2.7167	0.1111	0.6931	3.3322	0	0
6.	2.3858	2.7210	-0.3352	0.6931	3.2958	0	0
7.	2.1718	2.7226	-0.5509	0.0000	3.8286	0	0
8.	2.5743	2.7851	-0.2108	0.0000	3.4657	0	0
9.	2.8177	2.8233	-0.0056	1.0986	4.1589	0	0
10.	2.7264	2.7048	0.0216	1.3863	3.5264	1	.25
11.	2.7479	2.7620	-0.0141	1.0986	3.6889	0	0
12.	2.7395	2.7227	0.0168	0.6931	3.5553	1	.5
13.	2.6606	2.7556	-0.0950	1.0986	3.1781	0	0
14.	2.8647	2.6974	0.1673	0.0000	3.8286	0	0
15.	3.0438	2.7470	0.2968	0.6931	3.9120	0	0
16.	2.4055	2.7176	-0.3121	0.6931	3.4965	0	0
17.	2.5595	2.7322	-0.1727	0.6931	3.6109	0	0
18.	2.8343	2.6958	0.1386	1.0986	3.9890	0	0
19.	2.8574	2.7618	0.0956	0.0000	3.2189	0	0
20.	2.6740	2.7560	-0.0820	0.6931	3.5264	0	0

+-----+
Y ----- Y ----- e ----- X₁ ----- X₂ -----+

Projection Matrix

Define the matrix

$$\mathbf{P} = \mathbf{X}(\mathbf{X}'\mathbf{X})^{-1}\mathbf{X}'.$$

Observe that

$$\mathbf{P}\mathbf{X} = \mathbf{X}(\mathbf{X}'\mathbf{X})^{-1}\mathbf{X}'\mathbf{X} = \mathbf{X}.$$

This is a property of a **projection matrix**.

Cases:

- The matrix \mathbf{P} creates the fitted values in a least-squares regression:

$$\mathbf{P}\mathbf{Y} = \mathbf{X}(\mathbf{X}'\mathbf{X})^{-1}\mathbf{X}'\mathbf{Y} = \mathbf{X}\hat{\beta} = \hat{\mathbf{Y}}.$$

Because of this, the \mathbf{P} is sometimes called the "hat matrix"

- A special example of a projection matrix occurs when $\mathbf{X} = \mathbf{1}_n$ is a n -vector of ones. Then

$$\mathbf{P} = \mathbf{1}_n(\mathbf{1}'_n\mathbf{1}_n)^{-1}\mathbf{1}'_n = \frac{1}{n}\mathbf{1}_n\mathbf{1}'_n.$$

Note that in this case

$$\begin{aligned}\mathbf{P}\mathbf{Y} &= \mathbf{1}_n(\mathbf{1}'_n\mathbf{1}_n)^{-1}\mathbf{1}'_n\mathbf{Y} \\ &= \mathbf{1}_n\bar{Y}\end{aligned}$$

An n -vector whose elements are the sample mean \bar{Y} of Y_i .

Theorem: The projection matrix $\mathbf{P} = \mathbf{X}(\mathbf{X}'\mathbf{X})^{-1}\mathbf{X}'$ for any $n \times k$ matrix \mathbf{X} with $n \geq k$ has the following algebraic properties

- \mathbf{P} is symmetric ($\mathbf{P}' = \mathbf{P}$).
- \mathbf{P} is idempotent ($\mathbf{P}\mathbf{P} = \mathbf{P}$).
- $\text{tr } \mathbf{P} = k$.
- The eigenvalues of \mathbf{P} are 1 and 0. There are k eigenvalues equaling 1 and $n - k$ equaling 0.
- $\text{rank}(\mathbf{P}) = k$.

$$\begin{aligned}\mathbf{P}' &= (\mathbf{X}(\mathbf{X}'\mathbf{X})^{-1}\mathbf{X}')' \\ &= (\mathbf{X}')'((\mathbf{X}'\mathbf{X})^{-1})'(\mathbf{X}') \\ &= \mathbf{X}((\mathbf{X}'\mathbf{X})')^{-1}\mathbf{X}' \\ &= \mathbf{X}((\mathbf{X}'(\mathbf{X}')')^{-1}\mathbf{X}' \\ &= \mathbf{P}\end{aligned}$$

$$\begin{aligned}\mathbf{P}\mathbf{P} &= \mathbf{P}\mathbf{X}(\mathbf{X}'\mathbf{X})^{-1}\mathbf{X}' \\ &= \mathbf{X}(\mathbf{X}'\mathbf{X})^{-1}\mathbf{X}' \\ &= \mathbf{P}\end{aligned}$$

$$\begin{aligned}\text{tr } \mathbf{P} &= \text{tr}(\mathbf{X}(\mathbf{X}'\mathbf{X})^{-1}\mathbf{X}') \\ &= \text{tr}((\mathbf{X}'\mathbf{X})^{-1}\mathbf{X}'\mathbf{X}) \\ &= \text{tr}(\mathbf{I}_k) \\ &= k\end{aligned}$$

Orthogonal Projection

Define

$$\begin{aligned}\mathbf{M} &= \mathbf{I}_n - \mathbf{P} \\ &= \mathbf{I}_n - \mathbf{X}(\mathbf{X}'\mathbf{X})^{-1}\mathbf{X}'\end{aligned}$$

Note

$$\mathbf{M}\mathbf{X} = (\mathbf{I}_n - \mathbf{P})\mathbf{X} = \mathbf{X} - \mathbf{P}\mathbf{X} = \mathbf{X} - \mathbf{X} = \mathbf{0}. \quad \text{Thus } \mathbf{M} \text{ and } \mathbf{X} \text{ are orthogonal. We call } \mathbf{M} \text{ an orthogonal projection matrix.}$$

Cases:

- \mathbf{M} creates least-squares residuals:

$$\mathbf{M}\mathbf{Y} = \mathbf{Y} - \mathbf{P}\mathbf{Y} = \mathbf{Y} - \mathbf{X}\hat{\beta} = \hat{\mathbf{e}}.$$

- A special example of an orthogonal projection matrix occurs when $\mathbf{X} = \mathbf{1}_n$ is a n -vector of ones. Then

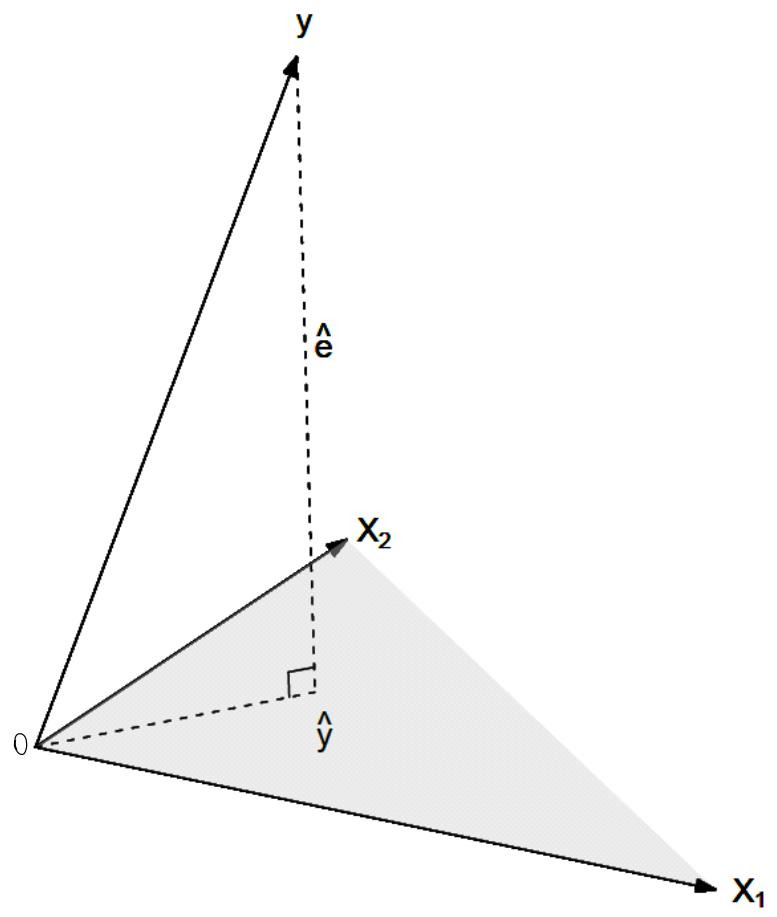
$$\mathbf{M} = \mathbf{I}_n - \mathbf{1}_n(\mathbf{1}'_n \mathbf{1}_n)^{-1} \mathbf{1}'_n.$$

Note that in this case

$$\mathbf{M}\mathbf{Y} = \mathbf{Y} - \mathbf{1}_n\bar{Y}.$$

\mathbf{M} creates demeaned values.

The orthogonal projection matrix \mathbf{M} has similar properties with \mathbf{P} , including that \mathbf{M} is symmetric ($\mathbf{M}' = \mathbf{M}$) and idempotent ($\mathbf{M}\mathbf{M} = \mathbf{M}$). We can calculate $\text{tr } \mathbf{M} = n - k$, and therefore the rank of \mathbf{M} is $n - k$.



Estimation of Error Variance

Substituting $\mathbf{Y} = \mathbf{X}\beta + \mathbf{e}$ into $\hat{\mathbf{e}} = \mathbf{M}\mathbf{Y}$ and using $\mathbf{M}\mathbf{X} = \mathbf{0}$ we find

$$\hat{\mathbf{e}} = \mathbf{M}\mathbf{Y} = \mathbf{M}(\mathbf{X}\beta + \mathbf{e}) = \mathbf{M}\mathbf{e},$$

which is free of dependence on the regression coefficient β .

The error variance $\sigma^2 = \mathbb{E}(e_i^2)$ is a moment, so a natural estimator is a *moment estimator*.

<i>moment estimator</i>	Infeasible		Feasible
Using summations	$\tilde{\sigma}^2 = \frac{1}{n} \sum_{i=1}^n e_i^2$		$\hat{\sigma}^2 = \frac{1}{n} \sum_{i=1}^n \hat{e}_i^2$
Using matrix notation	$\tilde{\sigma}^2 = n^{-1} \mathbf{e}' \mathbf{e}$		$\hat{\sigma}^2 = n^{-1} \hat{\mathbf{e}}' \hat{\mathbf{e}}$

Relationship between Infeasible vs Feasible

$$\begin{aligned}\hat{\sigma}^2 &= n^{-1} \hat{\mathbf{e}}' \hat{\mathbf{e}} \\ &= n^{-1} \mathbf{Y}' \mathbf{M} \mathbf{M} \mathbf{Y} \\ &= n^{-1} \mathbf{Y}' \mathbf{M} \mathbf{Y} \\ &= n^{-1} \mathbf{e}' \mathbf{M} \mathbf{e}\end{aligned}$$

An interesting implication is that

$$\begin{aligned}\tilde{\sigma}^2 - \hat{\sigma}^2 &= n^{-1} \mathbf{e}' \mathbf{e} - n^{-1} \mathbf{e}' \mathbf{M} \mathbf{e} \\ &= n^{-1} \mathbf{e}' \mathbf{P} \mathbf{e} \\ &\geq 0\end{aligned}$$

△ The final inequality holds because \mathbf{P} is positive semi-definite and $\mathbf{e}' \mathbf{P} \mathbf{e}$ is a quadratic form.

Analysis of Variance

$$\mathbf{Y} = \mathbf{PY} + \mathbf{MY} = \hat{\mathbf{Y}} + \hat{\mathbf{e}}$$

This decomposition is **orthogonal**, that is

$$TSS = ESS + RSS$$

$$\hat{\mathbf{Y}}' \hat{\mathbf{e}} = (\mathbf{PY})' (\mathbf{MY}) = \mathbf{Y}' \mathbf{PMY} = 0$$

It follows that

$$\mathbf{Y}' \mathbf{Y} = \hat{\mathbf{Y}}' \hat{\mathbf{Y}} + 2\hat{\mathbf{Y}}' \hat{\mathbf{e}} + \hat{\mathbf{e}}' \hat{\mathbf{e}} = \hat{\mathbf{Y}}' \hat{\mathbf{Y}} + \hat{\mathbf{e}}' \hat{\mathbf{e}}$$

$$\sum_{i=1}^n Y_i^2 = \sum_{i=1}^n \hat{Y}_i^2 + \sum_{i=1}^n \hat{e}_i^2$$

$$\mathbf{Y} - \mathbf{1}_n \bar{Y} = \hat{\mathbf{Y}} - \mathbf{1}_n \bar{Y} + \hat{\mathbf{e}}$$

$$(\hat{\mathbf{Y}} - \mathbf{1}_n \bar{Y})' \hat{\mathbf{e}} = \hat{\mathbf{Y}}' \hat{\mathbf{e}} - \bar{Y} \mathbf{1}_n' \hat{\mathbf{e}} = 0$$

$$(\mathbf{Y} - \mathbf{1}_n \bar{Y})' (\mathbf{Y} - \mathbf{1}_n \bar{Y}) = (\hat{\mathbf{Y}} - \mathbf{1}_n \bar{Y})' (\hat{\mathbf{Y}} - \mathbf{1}_n \bar{Y}) + \hat{\mathbf{e}}' \hat{\mathbf{e}}$$

$$\sum_{i=1}^n (Y_i - \bar{Y})^2 = \sum_{i=1}^n (\hat{Y}_i - \bar{Y})^2 + \sum_{i=1}^n \hat{e}_i^2$$

Analysis-of-variance formula

A commonly reported statistic is the **coefficient of determination or R-squared**

$$R^2 = \frac{\sum_{i=1}^n (\hat{Y}_i - \bar{Y})^2}{\sum_{i=1}^n (Y_i - \bar{Y})^2} = 1 - \frac{\sum_{i=1}^n \hat{e}_i^2}{\sum_{i=1}^n (Y_i - \bar{Y})^2}$$

⌚ It is often described as the fraction of the sample variance of Y_i which is explained by the least-squares fit.

The screenshot shows a Jupyter Notebook interface running on a web browser. The title bar indicates the notebook is titled "001_ols.ipynb - JupyterLab" and is located at "localhost:8888/lab/tree/HD-Metrics/code/001_ols.ipynb". The menu bar includes File, Edit, View, Run, Kernel, Tabs, Settings, and Help. A toolbar below the menu has icons for file operations like New, Open, Save, and Run. The main area displays a code cell [8] containing Stata commands to print TSS and ESS, followed by their respective output values (140.32277 and 127.4019). Below these, another code cell [9] is shown with the command to print RSS. The notebook is running in Python 3 (ipykernel) mode.

```
[8]: %stata display e(mss)+e(rss)
140.32277

Printing the ESS

[9]: %stata display e(mss)

Printing the RSS

[10]: %stata display e(rss)
127.4019

Printing the R2

[11]: %stata display e(mss)/(e(mss)+e(rss))
.09207965
```

Regression Components

Partition

$$\mathbf{X} = [\mathbf{X}_1 \mathbf{X}_2],$$

And

$$\beta = \begin{pmatrix} \beta_1 \\ \beta_2 \end{pmatrix}.$$

Then the regression model can be rewritten as

$$\mathbf{Y} = \mathbf{X}_1\beta_1 + \mathbf{X}_2\beta_2 + \mathbf{e}.$$

The OLS estimator of $\beta = (\beta'_1, \beta'_2)'$ is obtained by regression of \mathbf{Y} on $\mathbf{X} = [\mathbf{X}_1 \mathbf{X}_2]$ and can be written as

$$\mathbf{Y} = \mathbf{X}\hat{\beta} + \hat{\mathbf{e}} = \mathbf{X}_1\hat{\beta}_1 + \mathbf{X}_2\hat{\beta}_2 + \hat{\mathbf{e}}.$$

We are interested in algebraic expressions for $\hat{\beta}_1$ and $\hat{\beta}_2$. The least-squares estimator by definition is found by the joint minimization

$$(\hat{\beta}_1, \hat{\beta}_2) = \underset{b_1, b_2}{\operatorname{argmin}} \text{SSE}(b_1, b_2),$$

where

$$\text{SSE}(b_1, b_2) = (\mathbf{Y} - \mathbf{X}_1 b_1 - \mathbf{X}_2 b_2)' (\mathbf{Y} - \mathbf{X}_1 b_1 - \mathbf{X}_2 b_2).$$

An equivalent expression for $\hat{\beta}_1$ can be obtained by *concentration*.

$$\hat{\beta}_1 = \underset{b_1}{\operatorname{argmin}} \left(\underset{b_2}{\min} \text{SSE}(b_1, b_2) \right)$$

The inner expression
 $\underset{b_2}{\min} \text{SSE}(b_1, b_2)$ minimizes over b_2
while holding b_1 fixed.

This is simply the least squares regression of $\mathbf{Y} - \mathbf{X}_1 b_1$ on \mathbf{X}_2 . This has solution

$$\underset{b_2}{\operatorname{argmin}} \text{SSE}(b_1, b_2) = (\mathbf{X}'_2 \mathbf{X}_2)^{-1} (\mathbf{X}'_2 (\mathbf{Y} - \mathbf{X}_1 b_1))$$

Residuals:

$$\begin{aligned} \mathbf{Y} - \mathbf{X}_1 b_1 - \mathbf{X}_2 (\mathbf{X}'_2 \mathbf{X}_2)^{-1} (\mathbf{X}'_2 (\mathbf{Y} - \mathbf{X}_1 b_1)) &= (\mathbf{M}_2 \mathbf{Y} - \mathbf{M}_2 \mathbf{X}_1 b_1) \\ &= \mathbf{M}_2 (\mathbf{Y} - \mathbf{X}_1 b_1) \end{aligned}$$

The orthogonal projection matrix for \mathbf{X}_2 :
 $\mathbf{M}_2 = \mathbf{I}_n - \mathbf{X}_2 (\mathbf{X}'_2 \mathbf{X}_2)^{-1} \mathbf{X}'_2$ (Idempotent)

$$\begin{aligned} \underset{b_2}{\min} \text{SSE}(b_1, b_2) &= (\mathbf{Y} - \mathbf{X}_1 b_1)' \mathbf{M}_2 \mathbf{M}_2 (\mathbf{Y} - \mathbf{X}_1 b_1) \\ &= (\mathbf{Y} - \mathbf{X}_1 b_1)' \mathbf{M}_2 (\mathbf{Y} - \mathbf{X}_1 b_1) \end{aligned}$$

$$\hat{\beta}_1 = \underset{b_1}{\operatorname{argmin}} (\mathbf{Y} - \mathbf{X}_1 b_1)' \mathbf{M}_2 (\mathbf{Y} - \mathbf{X}_1 b_1)$$

$$= (\mathbf{X}'_1 \mathbf{M}_2 \mathbf{X}_1)^{-1} (\mathbf{X}'_1 \mathbf{M}_2 \mathbf{Y})$$

The orthogonal projection matrix for \mathbf{X}_1 :

$$\mathbf{M}_1 = \mathbf{I}_n - \mathbf{X}_1 (\mathbf{X}'_1 \mathbf{X}_1)^{-1} \mathbf{X}'_1 \text{ (Idempotent)}$$

By similar arguments we can find

$$\hat{\beta}_2 = (\mathbf{X}'_2 \mathbf{M}_1 \mathbf{X}_2)^{-1} (\mathbf{X}'_2 \mathbf{M}_1 \mathbf{Y})$$

Theorem: The least squares estimator $(\hat{\beta}_1, \hat{\beta}_2)$ for $\mathbf{Y} = \mathbf{X}\hat{\beta} + \hat{\mathbf{e}} = \mathbf{X}_1\hat{\beta}_1 + \mathbf{X}_2\hat{\beta}_2 + \hat{\mathbf{e}}$ has the algebraic solution

$$\hat{\beta}_1 = (\mathbf{X}'_1 \mathbf{M}_2 \mathbf{X}_1)^{-1} (\mathbf{X}'_1 \mathbf{M}_2 \mathbf{Y})$$

$$\hat{\beta}_2 = (\mathbf{X}'_2 \mathbf{M}_1 \mathbf{X}_2)^{-1} (\mathbf{X}'_2 \mathbf{M}_1 \mathbf{Y})$$

where \mathbf{M}_1 and \mathbf{M}_2 are defined above.

$$\hat{\beta}_2 = (\mathbf{X}'_2 \mathbf{M}_1 \mathbf{X}_2)^{-1} (\mathbf{X}'_2 \mathbf{M}_1 \mathbf{Y})$$

▲ Thus the coefficient estimate $\hat{\beta}_2$ is algebraically equal to the least-squares regression of $\tilde{\mathbf{e}}_1$ on $\tilde{\mathbf{X}}_2$.

$$= (\mathbf{X}'_2 \mathbf{M}_1 \mathbf{M}_1 \mathbf{X}_2)^{-1} (\mathbf{X}'_2 \mathbf{M}_1 \mathbf{M}_1 \mathbf{Y})$$

$$= (\tilde{\mathbf{X}}'_2 \tilde{\mathbf{X}}_2)^{-1} (\tilde{\mathbf{X}}'_2 \tilde{\mathbf{e}}_1)$$

$$\tilde{\mathbf{X}}_2 = \mathbf{M}_1 \mathbf{X}_2$$

$$\tilde{\mathbf{e}}_1 = \mathbf{M}_1 \mathbf{Y}$$

💡 In panel data, the FWL theorem is used to greatly improve computational time.

Theorem: Frisch-Waugh-Lovell (FWL)

In the model $\mathbf{Y} = \mathbf{X}_1\beta_1 + \mathbf{X}_2\beta_2 + \mathbf{e}$, the OLS estimator of β_2 and the OLS residuals $\hat{\mathbf{e}}$ may be computed by either the OLS regression $\mathbf{Y} = \mathbf{X}\hat{\beta} + \hat{\mathbf{e}} = \mathbf{X}_1\hat{\beta}_1 + \mathbf{X}_2\hat{\beta}_2 + \hat{\mathbf{e}}$ or via the following algorithm:

1. Regress \mathbf{Y} on \mathbf{X}_1 , obtain residuals $\tilde{\mathbf{e}}_1$;
2. Regress \mathbf{X}_2 on \mathbf{X}_1 , obtain residuals $\tilde{\mathbf{X}}_2$;
3. Regress $\tilde{\mathbf{e}}_1$ on $\tilde{\mathbf{X}}_2$, obtain OLS estimates $\hat{\beta}_2$ and residuals $\hat{\mathbf{e}}$.

Leverage Values & LOO Regression

The leverage values for the regressor matrix \mathbf{X} are the diagonal elements of the projection matrix $\mathbf{P} = \mathbf{X}(\mathbf{X}'\mathbf{X})^{-1}\mathbf{X}'$. There are n leverage values, and are typically written as h_{ii} for $i = 1, \dots, n$, i.e.,

$$h_{ii} = \mathbf{x}_i'(\mathbf{X}'\mathbf{X})^{-1}\mathbf{x}_i$$

$$\mathbf{P} = \begin{pmatrix} \mathbf{x}_1' \\ \mathbf{x}_2' \\ \vdots \\ \mathbf{x}_n' \end{pmatrix} (\mathbf{X}'\mathbf{X})^{-1} (\mathbf{x}_1 \mathbf{x}_2 \cdots \mathbf{x}_n)$$

⚠ The leverage value h_{ii} is a normalized length of the observed regressor vector x_i .

Theorem:

1. $0 \leq h_{ii} \leq 1$.
2. $h_{ii} \geq 1/n$ if X includes an intercept.
3. $\sum_{i=1}^n h_{ii} = k$.

```

Leverage Values & LOO Regression
Extracting the leverage values
[1]: Whatare sql
#display ?min
#display ?max
reg fit log_num_authors log_num_pages both_genders prop_comics
#display ?formula
#display ?fit
#display ?min_leverage
#display ?max_leverage
#display ?k
#display ?n

Checking that 0 ≤ h_{ii} ≤ 1
[1]: Whatare
#display ?min_leverage
#display ?max_leverage
#display ?k
#display ?n

Checking that h_{ii} ≥ 1/n by checking that min_{i=1,...,n} h_{ii} ≥ 1/n
[1]: Whatare
local sample_size = r(n)
local reciprocal = 1 / sample_size
local min_leverage = r(min)
display ?min_leverage >> LHS: ("min_leverage" := "reciprocal")

local sample_size = r(n)
local reciprocal = 1 / sample_size
local min_leverage = r(min)
display ?min_leverage >> RHS: ("min_leverage" := "reciprocal")

Checking that sum_{i=1}^n h_{ii} == k
[1]: Whatare
agen total_h = total(h)
local rounded_total_h = call(total_h)
display ?rounded_total_h >> RHS: ("rounded_total_h" := "rounded_total_h")

agen total_h = total(h)
local rounded_total_h = call(total_h)
display ?rounded_total_h >> LHS: ("rounded_total_h" := "rounded_total_h")
#display ?rounded_total_h >> RHS: ("rounded_total_h" := "rounded_total_h")
rounded_total_h == 100
#display ?rounded_total_h

```

Leave-One-Out Regression

Specifically, the leave-one-out least-squares estimator of the regression coefficient β is the least squares estimator constructed using the full sample excluding a single observation i . This can be written

$$\hat{\beta}_{(-i)} = \left(\sum_{j \neq i} X_j X_j' \right)^{-1} \left(\sum_{j \neq i} X_j Y_j \right)$$

$$= (\mathbf{X}'\mathbf{X} - \mathbf{x}_i \mathbf{x}_i')^{-1} (\mathbf{X}'\mathbf{Y} - \mathbf{x}_i \mathbf{Y}_i)$$

$$= (\mathbf{X}'_{(-i)} \mathbf{X}_{(-i)})^{-1} \mathbf{X}'_{(-i)} \mathbf{Y}_{(-i)}$$

- $\mathbf{X}_{(-i)}$ and $\mathbf{Y}_{(-i)}$ are the data matrices omitting the i th row.
- There is a leave-one-out (LOO) estimator for each observation, $i = 1, \dots, n$, so we have n such estimators.

LOO predicted value:	$\tilde{Y}_i = \mathbf{x}_i' \hat{\beta}_{(-i)}$
LOO residual, prediction error, or prediction residual:	$\tilde{e}_i = Y_i - \tilde{Y}_i$

Theorem: The leave-one-out least-squares estimator and prediction error can be calculated as

$$\hat{\beta}_{(-i)} = \hat{\beta} - (\mathbf{X}'\mathbf{X})^{-1} \mathbf{x}_i \tilde{e}_i$$

and

$$\tilde{e}_i = (1 - h_{ii})^{-1} \hat{e}_i$$

where h_{ii} are the leverage values.

The screenshot shows two Jupyter Notebook windows side-by-side. The left window, titled '001_ols.ipynb', contains Stata code. It includes commands to extract the matrix $\hat{\beta}_{-i} - \hat{\beta}$ for the 'prop_women' regressor and the LOO residuals \tilde{e}_i . It also lists observations 1430 and 1438. The right window, also titled '001_ols.ipynb', contains Python code using matplotlib to plot the sets $\{\hat{\beta}_{-i} - \hat{\beta} : i = 1, \dots, n\}$ and $\{\tilde{e}_i : i = 1, \dots, n\}$ against their indexes $i = 1, \dots, n$. The plots show two scatter plots: one for $\hat{\beta}_{-i} - \hat{\beta}$ vs i and another for \tilde{e}_i vs i .

```

[16]: % stata -qui
predict dfbeta, dfbeta(prop_women)
gen double dfbeta_se = dfbeta"se(prop_women)
cvs_regress, generre(e_tilde)

We identify an observation (row 1430) for which there is a perfect fit.

[17]: % stata list FKG FKG_hat h dfbeta_se e_hat e_tilde in 1430, table separator(20)

+-----+
| FKG FKG_hat h dfbeta_se e_hat e_tilde |
+-----+
1430. | 2.6987 2.6987 1 . 0.0000 -.01999181 |
+-----+

[18]: % stata -fouts original,subset
frame put dfbeta_se e_tilde, into(original)
drop in 1430
frame put dfbeta_se e_tilde, into(subset)

. frame put dfbeta_se e_tilde, into(original)
.
. drop in 1430
(1 observation deleted)
. frame put dfbeta_se e_tilde, into(subset)
.

[19]: from sfi import Data
import numpy as np
import pandas as pd
dfbeta_py = subset[:,0]
e_tilde_py = subset[:,1]

Simple 0 1 main Python 3 (ipykernel) | Idle Mode Command Ln 1, Col 1 001_ols.ipynb 0

```

```

[20]: import matplotlib.pyplot as plt
fig, axs = plt.subplots(2)
axs[0].plot(np.linspace(1,e_tilde_py.size,num=e_tilde_py.size).astype(int), dfbeta_py)
axs[0].set_title('$\widehat{\beta}_{-i} - \hat{\beta}$ vs $i=1,\dots,n$')
axs[1].plot(np.linspace(1,e_tilde_py.size,num=e_tilde_py.size).astype(int), e_tilde_py, 'tab:orange')
axs[1].set_title('$\tilde{e}_i$ vs $i=1,\dots,n$')
# Hide x Labels and tick labels for top plots
for ax in axs.flat:
    ax.label_outer()

Simple 0 1 main Python 3 (ipykernel) | Idle Mode Command Ln 1, Col 1 001_ols.ipynb 0

```

Define

$$\begin{aligned} M^* &= (I_n - \text{diag}\{h_{11}, \dots, h_{nn}\})^{-1} \\ &= \text{diag}\{(1 - h_{11})^{-1}, \dots, (1 - h_{nn})^{-1}\} \end{aligned}$$

Then a natural estimator of out-of-sample mean squared error is

$$\tilde{\sigma}^2 = \frac{1}{n} \sum_{i=1}^n \tilde{e}_i^2 = \frac{1}{n} \sum_{i=1}^n (1 - h_{ii})^{-2} \hat{e}_i^2 \quad (\text{sample mean squared prediction error}). \text{ Its root } \tilde{\sigma} = \sqrt{\tilde{\sigma}^2} \text{ is known as the prediction standard error.}$$

OLS: Mean

Assumption 2: linear Regression Model

The observations (Y_i, X_i) satisfy the linear regression equation

$$Y_i = X'_i \beta + e_i, \\ \mathbb{E}(e_i | X_i) = 0.$$

The variables have finite second moments

$$\mathbb{E}(Y_i^2) < \infty, \\ \mathbb{E}||X_i||^2 < \infty,$$

and an invertible design matrix

$$Q_{XX} = \mathbb{E}(X_i X'_i) > 0.$$

Assumption 3: Homoskedastic linear Regression Model

In addition to Assumption 2 ,

$$\mathbb{E}(e_i^2 | X_i) = \sigma^2(X_i) = \sigma^2,$$

is independent of X_i .

We will consider both the general case of **heteroskedastic regression**, where the conditional variance $\mathbb{E}(e_i^2 | X_i) = \sigma^2(X_i) = \sigma_i^2$

Mean of Least-Squares Estimator

$$\mathbb{E}(Y_i | X) = \mathbb{E}(Y_i | X_i) = X'_i \beta$$

The first equality states that the conditional expectation of Y_i given $\{X_1, \dots, X_n\}$ only depends on X_i , since the observations are independent across i . The second equality is the assumption of a linear conditional mean.

(Conditional) Unbiasedness

$$\begin{aligned} \mathbb{E}(\hat{\beta} | X) &= \mathbb{E}\left(\left(\sum_{i=1}^n X_i X'_i\right)^{-1} \left(\sum_{i=1}^n X_i Y_i\right) | X\right) \\ &= \left(\sum_{i=1}^n X_i X'_i\right)^{-1} \mathbb{E}\left(\sum_{i=1}^n X_i Y_i | X\right) \\ &= \left(\sum_{i=1}^n X_i X'_i\right)^{-1} \sum_{i=1}^n \mathbb{E}(X_i Y_i | X) \\ &= \left(\sum_{i=1}^n X_i X'_i\right)^{-1} \sum_{i=1}^n X_i \mathbb{E}(Y_i | X) \\ &= \left(\sum_{i=1}^n X_i X'_i\right)^{-1} \sum_{i=1}^n X_i X'_i \beta \\ &= \beta \end{aligned}$$

Theorem: Mean of Least-Squares Estimator

In the linear regression model (Assumption 2) and i.i.d. sampling (Assumption 1)

$$\mathbb{E}(\hat{\beta} | X) = \beta.$$

Matrix Notation:

$$\begin{aligned} \mathbb{E}(\hat{\beta} | X) &= \mathbb{E}((X'X)^{-1} X' Y | X) & \mathbb{E}(Y | X) &= \begin{pmatrix} \vdots \\ \mathbb{E}(Y_i | X) \\ \vdots \end{pmatrix} = \begin{pmatrix} \vdots \\ X'_i \beta \\ \vdots \end{pmatrix} = X \beta \\ &= (X'X)^{-1} X' \mathbb{E}(Y | X) \\ &= (X'X)^{-1} X' X \beta \\ &= \beta & \mathbb{E}(e | X) &= \begin{pmatrix} \vdots \\ \mathbb{E}(e_i | X) \\ \vdots \end{pmatrix} = \begin{pmatrix} \vdots \\ \mathbb{E}(e_i | X_i) \\ \vdots \end{pmatrix} = 0 \end{aligned}$$

OLS: Variance

The objective of this section is to derived the **(conditional) covariance matrix** of the regression coefficient estimates, i.e.,

$$V_{\hat{\beta}} \stackrel{\text{def}}{=} \text{var}(\hat{\beta} | X).$$

For any pair (Z, X) the **conditional covariance matrix** is

$$\text{var}(Z | X) = \mathbb{E}((Z - \mathbb{E}(Z | X))(Z - \mathbb{E}(Z | X))' | X)$$

For any $r \times 1$ random vector Z the **covariance matrix** is

$$\begin{aligned} \text{var}(Z) &= \mathbb{E}((Z - \mathbb{E}(Z))(Z - \mathbb{E}(Z))') \\ &= \mathbb{E}(ZZ') - (\mathbb{E}(Z))(\mathbb{E}(Z))' \end{aligned}$$

Firstly, notice that

$$\begin{aligned} \hat{\beta} &= (X'X)^{-1}(X'(\beta + e)) \\ &= (X'X)^{-1}X'\beta + (X'X)^{-1}X'e \\ &= \beta + (X'X)^{-1}X'e \\ \hat{\beta} - \beta &= (X'X)^{-1}X'e \\ \hat{\beta} - \mathbb{E}(\hat{\beta} | X) &= (X'X)^{-1}X'e \end{aligned}$$

$$V_{\hat{\beta}} = \text{var}(\hat{\beta} | X) = \mathbb{E}((\hat{\beta} - \mathbb{E}(\hat{\beta} | X))(\hat{\beta} - \mathbb{E}(\hat{\beta} | X))' | X) = \mathbb{E}((X'X)^{-1}X'ee'X(X'X)^{-1} | X) = (X'X)^{-1}X'\mathbb{E}(ee' | X)X(X'X)^{-1}$$

$$V_{\hat{\beta}} = \text{var}(\hat{\beta} | X) = (X'X)^{-1}X'DX(X'X)^{-1}$$

The i th diagonal element of D is

$$\mathbb{E}(e_i^2 | X) = \mathbb{E}(e_i^2 | X_i) = \sigma_i^2$$

while the ij th off-diagonal element of D is

$$\mathbb{E}(e_i e_j | X) = \mathbb{E}(e_i | X_i) \mathbb{E}(e_j | X_j) = 0$$

Therefore

$$D = \text{diag}(\sigma_1^2, \dots, \sigma_n^2) = \begin{pmatrix} \sigma_1^2 & 0 & \cdots & 0 \\ 0 & \sigma_2^2 & \cdots & 0 \\ \vdots & \vdots & \ddots & \vdots \\ 0 & 0 & \cdots & \sigma_n^2 \end{pmatrix}$$

$$X'DX = \sum_{i=1}^n X_i X'_i \sigma_i^2$$

Theorem: Variance of Least-Squares Estimator In the linear regression model (Assumption 2) and i.i.d. sampling (Assumption 1)

$$\begin{aligned} V_{\hat{\beta}} &= \text{var}(\hat{\beta} | X) \\ &= (X'X)^{-1}(X'DX)(X'X)^{-1} \end{aligned}$$

where D is defined above.

In the homoskedastic linear regression model (Assumption 3) and i.i.d. sampling (Assumption 1)

$$V_{\hat{\beta}} = \sigma^2 (X'X)^{-1} .$$

Residuals

	Residuals		Prediction Errors
Definition:	$\hat{e}_i = Y_i - X'_i \hat{\beta}$		$\tilde{e}_i = Y_i - X'_i \hat{\beta}_{(-i)}$ $\tilde{e}_i = (1 - h_{ii})^{-1} \hat{e}_i$
Matrix Notation:	$\hat{e} = \mathbf{M}\mathbf{e}$ $\mathbf{M} = \mathbf{I}_n - \mathbf{X}(\mathbf{X}'\mathbf{X})^{-1}\mathbf{X}'$		$\tilde{e} = \mathbf{M}^*\hat{e} = \mathbf{M}^*\mathbf{M}\mathbf{e}$ \mathbf{M}^* is a diagonal matrix with i th diagonal element $(1 - h_{ii})^{-1}$
(Conditional) Mean:	$\mathbb{E}(\hat{e} \mathbf{X}) = \mathbb{E}(\mathbf{M}\mathbf{e} \mathbf{X}) = \mathbf{M}\mathbb{E}(\mathbf{e} \mathbf{X}) = \mathbf{0}$		$\mathbb{E}(\tilde{e} \mathbf{X}) = \mathbf{M}^*\mathbf{M}\mathbb{E}(\mathbf{e} \mathbf{X}) = \mathbf{0}$
(Conditional) Variance:	$\text{var}(\hat{e} \mathbf{X}) = \text{var}(\mathbf{M}\mathbf{e} \mathbf{X}) = \mathbf{M}\text{var}(\mathbf{e} \mathbf{X})\mathbf{M} = \mathbf{M}\mathbf{D}\mathbf{M}$		$\text{var}(\tilde{e} \mathbf{X}) = \mathbf{M}^*\mathbf{M}\text{var}(\mathbf{e} \mathbf{X})\mathbf{M}\mathbf{M}^* = \mathbf{M}^*\mathbf{M}\mathbf{D}\mathbf{M}\mathbf{M}^*$

Homoskedasticity ($\mathbf{D} = \mathbf{I}_n \sigma^2$)

	Residuals		Prediction Errors
(Conditional) Variance:	$\text{var}(\hat{e} \mathbf{X}) = \mathbf{M}\sigma^2$		$\text{var}(\tilde{e} \mathbf{X}) = \mathbf{M}^*\mathbf{M}\mathbf{M}^*\sigma^2$
i th Observation:	$\text{var}(\hat{e}_i \mathbf{X}) = \mathbb{E}(\hat{e}_i^2 \mathbf{X}) = (1 - h_{ii})\sigma^2$		$\text{var}(\tilde{e}_i \mathbf{X}) = \mathbb{E}(\tilde{e}_i^2 \mathbf{X}) = (1 - h_{ii})^{-1}(1 - h_{ii})(1 - h_{ii})^{-1}\sigma^2 = (1 - h_{ii})^{-1}\sigma^2$

When the true errors, \mathbf{e} , are *homoskedastic*, then a residual with constant conditional variance can be obtained by rescaling. The **standardized residuals** are

	Standardized Residuals
Definition:	$\bar{e}_i = (1 - h_{ii})^{-1/2} \hat{e}_i$
Matrix Notation:	$\bar{\mathbf{e}} = (\bar{e}_1, \dots, \bar{e}_n)' = \mathbf{M}^{*1/2}\mathbf{M}\mathbf{e}$
(Conditional) Variance:	$\text{var}(\bar{e} \mathbf{X}) = \mathbf{M}^{*1/2}\mathbf{M}\mathbf{M}^{*1/2}\sigma^2$

$$\text{var}(\bar{e}_i | \mathbf{X}) = \mathbb{E}(\bar{e}_i^2 | \mathbf{X}) = \sigma^2$$

Guest 001_ols.ipynb - JupyterLab

localhost:8888/lab/tree/HD-Metrics/code/001_ols.ipynb

File Edit View Run Kernel Tabs Settings Help

001_ols.ipynb +

Notebook Python 3 (ipykernel)

Residuals

```
[21]: %%stata
gen e_bar = e_hat/sqrt(1-h)
format e_hat e_tilde e_bar %5.4f
list e_hat e_tilde e_bar in 1/10, table separator(20)
summarize e_hat e_tilde e_bar

. gen e_bar = e_hat/sqrt(1-h)
(1 missing value generated)

. format e_hat e_tilde e_bar %5.4f

. list e_hat e_tilde e_bar in 1/10, table separator(20)

+-----+
| e_hat   e_tilde   e_bar |
|-----|
1. | 0.0973   0.0981   0.0977 |
2. | 0.0482   0.0485   0.0483 |
3. | 0.0812   0.0821   0.0817 |
4. | -0.0262  -0.0265  -0.0264 |
5. | 0.1111   0.1117   0.1114 |
6. | -0.3352  -0.3381  -0.3366 |
7. | -0.5589  -0.5568  -0.5539 |
8. | -0.2108  -0.2125  -0.2117 |
9. | -0.0056  -0.0066  -0.0061 |
10. | 0.0216   0.0218   0.0217  |
+-----+

. summarize e_hat e_tilde e_bar

Variable |      Obs        Mean     Std. dev.      Min      Max
-----+-----+-----+-----+-----+-----+
    e_hat |    4,988    1.59e-16   .1598337   -.719324   1.287893
    e_tilde |    4,988   -5.42e-06   .1692646   -.78422   1.296832
    e_bar |    4,987   -5.28e-07   .1636692   -.7460606   1.292355
```

Simple

0 \$ 4 Python 3 (ipykernel) | Idle

Mode: Command

Ln 1, Col 33 001_ols.ipynb 1

Covariance Matrix Estimation

1. [Homoskedastic Case](#)
2. [Heteroskedastic Case](#)

Homoskedastic Case

Under homoskedasticity, the covariance matrix takes the relatively simple form

$$\mathbf{V}_{\hat{\beta}}^0 = (\mathbf{X}'\mathbf{X})^{-1}\sigma^2$$

Using the 'plug-in' principle, the classic covariance matrix estimator in the homoskedastic case is

$$\widehat{\mathbf{V}}_{\hat{\beta}}^0 = (\mathbf{X}'\mathbf{X})^{-1}s^2$$

Reminder:

Object Definition:	$\mathbf{V}_{\hat{\beta}}^0 = (\mathbf{X}'\mathbf{X})^{-1}\sigma^2$
Estimator:	$\widehat{\mathbf{V}}_{\hat{\beta}}^0 = (\mathbf{X}'\mathbf{X})^{-1}s^2$
(Conditional) Mean Of the Estimator:	$\begin{aligned} \mathbb{E}(\widehat{\mathbf{V}}_{\hat{\beta}}^0 \mathbf{X}) &= (\mathbf{X}'\mathbf{X})^{-1}\mathbb{E}(s^2 \mathbf{X}) \\ &= (\mathbf{X}'\mathbf{X})^{-1}\sigma^2 \\ &= \mathbf{V}_{\hat{\beta}}^0 \end{aligned}$

✓ $\widehat{\mathbf{V}}_{\hat{\beta}}^0$ is conditionally unbiased for $\mathbf{V}_{\hat{\beta}}^0$ under the assumption of homoskedasticity.

Guest 001_ols.ipynb - JupyterLab

localhost:8888/lab/tree/HD-Metrics/code/001_ols.ipynb

File Edit View Run Kernel Tabs Settings Help

001_ols.ipynb +

Notebook Python 3 (ipykernel)

Covariance Matrix Estimation

Homoskedasticity:

$$\hat{\mathbf{V}}_{\beta}^0 = (\mathbf{X}'\mathbf{X})^{-1}s^2$$

```
[22]: %%stata -qui  
#delimit ;  
quietly reg FKG log_num_authors log_num_pages both_genders prop_women  
`journals' `jel_imp' y_2-y_20 c_2-c_215 jel_flag;  
matrix subV = (e(V)[1,1], e(V)[1,2], e(V)[1,3], e(V)[1,4], e(V)[1,262] \  
e(V)[2,1], e(V)[2,2], e(V)[2,3], e(V)[2,4], e(V)[2,262] \  
e(V)[3,1], e(V)[3,2], e(V)[3,3], e(V)[3,4], e(V)[3,262] \  
e(V)[4,1], e(V)[4,2], e(V)[4,3], e(V)[4,4], e(V)[4,262] \  
e(V)[262,1], e(V)[262,2], e(V)[262,3], e(V)[262,4], e(V)[262,262]);  
#delimit cr
```

```
[23]: %stata matrix list subV
```

	c1	c2	c3	c4	c5
r1	.00004099				
r2	-1.088e-06	.0000722			
r3	-.00002105	1.429e-06	.00006243		
r4	.0000197	-3.541e-06	-.0000547	.00013851	
r5	-7.562e-06	-.000020284	7.309e-06	-.00001451	.00081714

Simple 0 \$ 4 Python 3 (ipykernel) | Idle Mode: Command Ln 3, Col 33 001_ols.ipynb 1

Heteroskedastic Case

Under heteroskedasticity, the general form for the covariance matrix is

$$V_{\beta} = (X'X)^{-1}(X'DX)(X'X)^{-1}$$

Where

$$\begin{aligned} D &= \text{diag}(o_1^2, \dots, o_n^2) \\ &= \mathbb{E}(ee' | X) \\ &= \mathbb{E}(\tilde{D} | X) \quad \tilde{D} = \text{diag}(e_1^2, \dots, e_n^2) \text{ is a conditionally unbiased estimator for } D. \end{aligned}$$

Therefore, if the squared errors e_i^2 were observable, we could construct and unbiased 'ideal' estimator for V_{β} as

Reminder:

Object Definition:	$V_{\beta} = (X'X)^{-1}(X'DX)(X'X)^{-1}$
(ideal) Estimator:	$\hat{V}_{\beta}^{\text{ideal}} = (X'X)^{-1}(X'\tilde{D}X)(X'X)^{-1}$ $= (X'X)^{-1} \left(\sum_{i=1}^n X_i X_i' e_i^2 \right) (X'X)^{-1}$
(Conditional) Mean Of the Estimator:	$\mathbb{E}(\hat{V}_{\beta}^{\text{ideal}} X) = (X'X)^{-1} \left(\sum_{i=1}^n X_i X_i' \mathbb{E}(e_i^2 X) \right) (X'X)^{-1}$ $= (X'X)^{-1} \left(\sum_{i=1}^n X_i X_i' \sigma_i^2 \right) (X'X)^{-1}$ $= (X'X)^{-1}(X'DX)(X'X)^{-1}$ $= V_{\beta}$

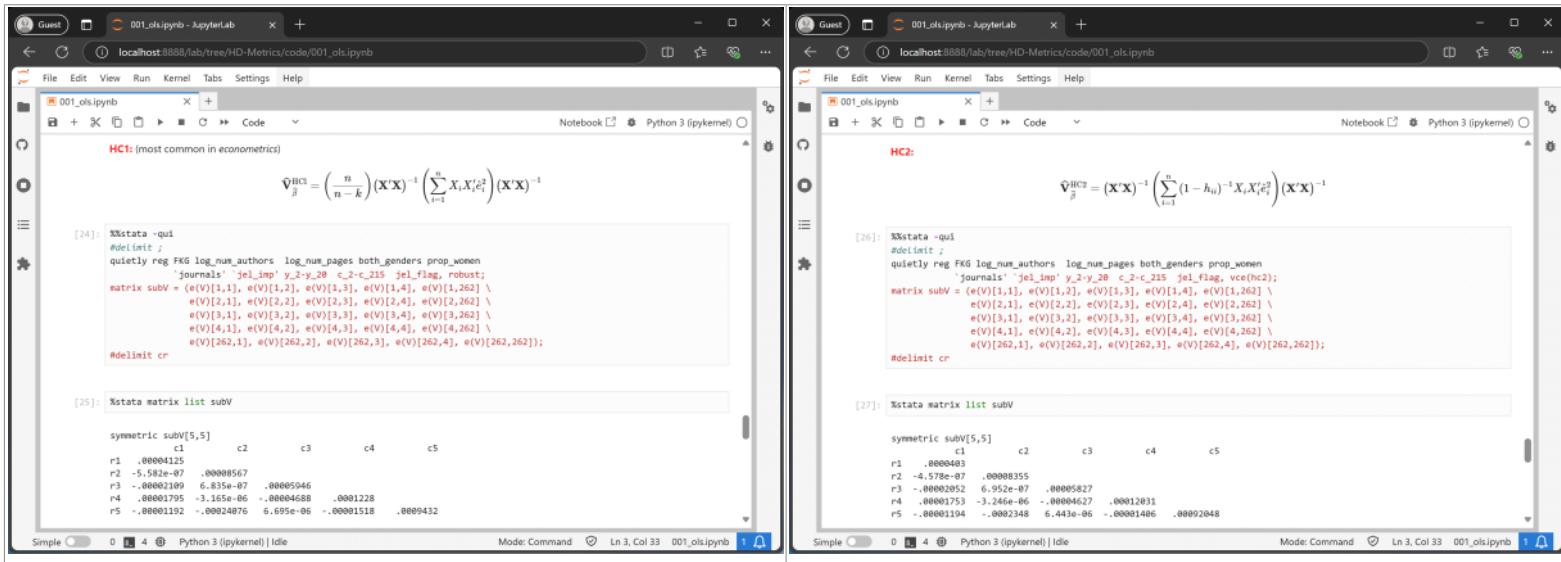
Feasible Heteroskedasticity-Consistent (HC) Estimators

HCO:	$\hat{V}_{\beta}^{\text{HCO}} = (X'X)^{-1} \left(\sum_{i=1}^n X_i X_i' \hat{e}_i^2 \right) (X'X)^{-1}$
HC1:	$\hat{V}_{\beta}^{\text{HC1}} = \left(\frac{n}{n-k} \right) (X'X)^{-1} \left(\sum_{i=1}^n X_i X_i' \hat{e}_i^2 \right) (X'X)^{-1}$
HC2:	$\hat{V}_{\beta}^{\text{HC2}} = (X'X)^{-1} \left(\sum_{i=1}^n X_i X_i' \bar{e}_i^2 \right) (X'X)^{-1}$ $= (X'X)^{-1} \left(\sum_{i=1}^n (1 - h_{ii})^{-1} X_i X_i' \hat{e}_i^2 \right) (X'X)^{-1}$
HC3:	$\hat{V}_{\beta}^{\text{HC3}} = (X'X)^{-1} \left(\sum_{i=1}^n X_i X_i' \hat{e}_i^2 \right) (X'X)^{-1}$ $= (X'X)^{-1} \left(\sum_{i=1}^n (1 - h_{ii})^{-2} X_i X_i' \hat{e}_i^2 \right) (X'X)^{-1}$

⚠️ The default *unchallenged* choice in applied work.

Since $(1 - h_{ii})^{-2} > (1 - h_{ii})^{-1} > 1$
it is straightforward to show that

$$\hat{V}_{\beta}^{\text{HCO}} < \hat{V}_{\beta}^{\text{HC2}} < \hat{V}_{\beta}^{\text{HC3}}$$



```

[24]: % stata -qui
        #delimit ;
quietly reg FKG log_num_authors log_num_pages both_genders prop_women
        `journals' `jel_flag' y_2_y_20 c_2_c_215 `jel_flag', robust;
matrix subV = (e(V)[1,1], e(V)[1,2], e(V)[1,3], e(V)[1,4], e(V)[1,262] \
e(V)[2,1], e(V)[2,2], e(V)[2,3], e(V)[2,4], e(V)[2,262] \
e(V)[3,1], e(V)[3,2], e(V)[3,3], e(V)[3,4], e(V)[3,262] \
e(V)[4,1], e(V)[4,2], e(V)[4,3], e(V)[4,4], e(V)[4,262] \
e(V)[262,1], e(V)[262,2], e(V)[262,3], e(V)[262,4], e(V)[262,262]);
#delimit cr

[25]: % stata matrix list subV
symmetric subV[5,5]
      c1      c2      c3      c4      c5
r1   .00004125
r2   -.5582e-07  .00008567
r3   -.00002109  6.858e-07  .00005946
r4   .00001795  -.3165e-06  -.00004688  .00001228
r5   -.00001192  -.00024076  6.895e-06  -.00001518  .0000432
```



```

[26]: % stata -qui
        #delimit ;
quietly reg FKG log_num_authors log_num_pages both_genders prop_women
        `journals' `jel_flag' y_2_y_20 c_2_c_215 `jel_flag', vce(hc2);
matrix subV = (e(V)[1,1], e(V)[1,2], e(V)[1,3], e(V)[1,4], e(V)[1,262] \
e(V)[2,1], e(V)[2,2], e(V)[2,3], e(V)[2,4], e(V)[2,262] \
e(V)[3,1], e(V)[3,2], e(V)[3,3], e(V)[3,4], e(V)[3,262] \
e(V)[4,1], e(V)[4,2], e(V)[4,3], e(V)[4,4], e(V)[4,262] \
e(V)[262,1], e(V)[262,2], e(V)[262,3], e(V)[262,4], e(V)[262,262]);
#delimit cr

[27]: % stata matrix list subV
symmetric subV[5,5]
      c1      c2      c3      c4      c5
r1   .0000483
r2   -.4578e-07  .00008355
r3   -.00002052  6.952e-07  .00005827
r4   .00001753  -.3146e-06  -.00004627  .00001203
r5   -.00001194  -.0002348  6.443e-06  -.00001406  .00002048
```

The screenshot shows a Jupyter Notebook window with two code cells. Cell [28] contains Stata code for calculating the covariance matrix of the OLS estimator using the HC3 heteroskedasticity-corrected standard errors. Cell [29] contains the resulting symmetric matrix output.

```
[28]: HC3:

$$\hat{V}_{\hat{\beta}}^{HC3} = (\mathbf{X}'\mathbf{X})^{-1} \left( \sum_{i=1}^n (1 - h_{ii})^{-2} X_i X_i' \hat{e}_i^2 \right) (\mathbf{X}'\mathbf{X})^{-1}$$


[29]: K stata -qui
#delimit ;
quietly reg FKG log_num_authors log_num_page both_genders prop_women
`journals' `jel_jmp' y_2-y_20 c_2-c_215 `jel_flag', vce(Hc3);
matrix subv = (e(V)[1,1], e(V)[1,2], e(V)[1,3], e(V)[1,4], e(V)[1,262]) \
e(V)[2,1], e(V)[2,2], e(V)[2,3], e(V)[2,4], e(V)[2,262] \
e(V)[3,1], e(V)[3,2], e(V)[3,3], e(V)[3,4], e(V)[3,262] \
e(V)[4,1], e(V)[4,2], e(V)[4,3], e(V)[4,4], e(V)[4,262] \
e(V)[262,1], e(V)[262,2], e(V)[262,3], e(V)[262,4], e(V)[262,262]);
#delimit cr
```

```
[29]: K stata matrix list subv
symmetric subv[5,5]
      c1        c2        c3        c4        c5
r1  .00004194
r2  -.3135e-07  .00008692
r3  -.00002119  7.647e-07  .00006098
r4  .00001823  -3.6088e-06  .00004889  .00012578
r5  -.00001295  -.00024429  6.464e-06  -.00001345  .00005855
```

Standard Error

Definition: A standard error $s(\hat{\beta})$ for a real-valued estimator $\hat{\beta}$ is an estimator of the standard deviation of the distribution of $\hat{\beta}$.

When β is a vector with estimator $\hat{\beta}$ and covariance matrix estimator $\hat{V}_{\hat{\beta}}$, standard errors for individual elements are the square roots of the diagonal elements of $\hat{V}_{\hat{\beta}}$. That is,

$$s(\hat{\beta}_j) = \sqrt{\hat{V}_{\hat{\beta}_j}} = \sqrt{[\hat{V}_{\hat{\beta}}^2]_{jj}}$$

Where $j \in \{0, \text{HC0}, \text{HC1}, \text{HC2}, \text{HC3}\}$

Measures of Fit

	Definition	Comments
R-Squared (R^2):	$R^2 = 1 - \frac{\sum_{i=1}^n \hat{e}_i^2}{\sum_{i=1}^n (y_i - \bar{y})^2}$	<p>☒ R^2 can be viewed as an estimator of the population parameter</p> $\rho^2 = \frac{\text{var}(X'_i \beta)}{\text{var}(y_i)} = 1 - \frac{\sigma^2}{\sigma_y^2}$ <p>△ It cannot decrease as k increases.</p>
(Adjusted) R-Squared (\bar{R}^2):	$\bar{R}^2 = 1 - \frac{(n-1) \sum_{i=1}^n \hat{e}_i^2}{(n-k) \sum_{i=1}^n (y_i - \bar{y})^2}$	<p>☒ It is called 'adjusted' because it corrects for the terms' degrees-of-freedom.</p> <p>△ \bar{R}^2 is an inappropriate choice for model selection (it tends to select models with too many parameters).</p> <p>However it is commonly reported in applied work.</p>
(Alternative) R-Squared (\tilde{R}^2):	$\tilde{R}^2 = 1 - \frac{\sum_{i=1}^n \tilde{e}_i^2}{\sum_{i=1}^n (y_i - \bar{y})^2}$	<p>✓ Models with high \tilde{R}^2 are better models in terms of expected out of sample squared error, and therefore suitable for model selection.</p>

```

Guest 001_ols.ipynb - JupyterLab
localhost:8888/lab/tree/HD-Metrics/code/001_ols.ipynb
File Edit View Run Kernel Tabs Settings Help
001_ols.ipynb
Notebook Python 3 (ipykernel)
Measures of Fit

R-squared:

$$R^2 = 1 - \frac{\sum_{i=1}^n \hat{e}_i^2}{\sum_{i=1}^n (y_i - \bar{y})^2}$$


[30]: from sfi import Scalar
rsquared=Scalar.getValue('e(r2)')
print(rsquared)
0.09207964664397683

Adjusted R-squared:

$$\bar{R}^2 = 1 - \frac{(n-1) \sum_{i=1}^n \hat{e}_i^2}{(n-k) \sum_{i=1}^n (y_i - \bar{y})^2}$$


[31]: radj=Scalar.getValue('e(r2_a)')
print(radj)
0.041938467586439376

(Alternative) R-squared:

$$\tilde{R}^2 = 1 - \frac{\sum_{i=1}^n \tilde{e}_i^2}{\sum_{i=1}^n (y_i - \bar{y})^2}$$


[32]: from pystata import stata
stata.run(``reg FKG log_num_authors log_num_pages both_genders prop_women ///
`journals' `jel_imp' y_2-y_20 c_2-c_215 jel_flag
``', quietly=True)
mss=Scalar.getValue('e(mss)')
rss=Scalar.getValue('e(rss)')
R2_tilde = 1 - (e_tilde_py**2).sum()/(mss+rss)
print(R2_tilde)

-0.018224014299829117

```

Cluster Sampling

A difficulty with applying the classical regression framework is that articles may share the same co-author across journals and time. Readability of different articles may be affected by writing style of shared coauthors, all of which imply dependence. These concerns, however, do not suggest that readability will be correlated across set of articles that do not share a coauthor, so it seems reasonable to model readability across these sets as mutually independent.

In clustering contexts it is convenient to double index the observations as (Y_{ig}, X_{ig}) where $g = 1, \dots, G$ indexes the cluster and $i = 1, \dots, n_g$ indexes the individual within the g th cluster. The number of observations per cluster n_g may vary across clusters. The number of clusters is G . The total number of observations is $n = \sum_{g=1}^G n_g$. In article readability example, the number of clusters (cluster) in the estimation sample is $G = 215$ and the total number of observations is $n = 4,988$.

While it is typical to write the observations using the double index notation (Y_{ig}, X_{ig}) , it is also useful to use cluster-level notation. Let $\mathbf{Y}_g = (Y_{1g}, \dots, Y_{n_g g})'$ and $\mathbf{X}_g = (X_{1g}, \dots, X_{n_g g})'$ denote the $n_g \times 1$ vector of dependent variables and $n_g \times k$ matrix of regressors for the g th cluster. A linear regression model can be written for the individual observations as

$$Y_{ig} = X'_{ig}\beta + e_{ig}$$

and using cluster notation as

$$\mathbf{Y}_g = \mathbf{X}_g\beta + \mathbf{e}_g$$

where $\mathbf{e}_g = (e_{1g}, \dots, e_{n_g g})'$ is a $n_g \times 1$ error vector. We can also stack the observations into full sample matrices and write the model as

$$\mathbf{Y} = \mathbf{X}\beta + \mathbf{e}$$

```
% stata
sort cluster
list FKG log_num_authors log_num_pages both_genders prop_women cluster ///
in 4304/4322, table separator(20)

sort cluster

list FKG log_num_authors log_num_pages both_genders prop_women cluster ///
> in 4304/4322, table separator(20)

+-----+
| FKG  log_nrs  log_ges  both_g-s  prop_w-n  cluster |
+-----+
4304. | 2.6851  0.6931  2.7726  0   0   1   |
4305. | 2.9652  0.0000  3.4340  0   0   1   |
4306. | 2.7479  1.0986  3.6889  0   0   1   |
4307. | 2.7692  0.6931  3.6189  1   .5   1   |
4308. | 2.9889  0.6931  2.7081  0   0   2   |
4309. | 2.7419  1.0986  2.8904  0   0   2   |
4310. | 2.9387  0.6931  3.6910  0   0   2   |
4311. | 3.0209  0.0000  3.6910  0   0   3   |
4312. | 3.1722  0.6931  2.7726  0   0   3   |
4313. | 2.9696  0.6931  2.3979  0   0   3   |
4314. | 2.8869  0.0000  3.5553  0   0   4   |
4315. | 2.7571  0.0000  2.7726  0   0   4   |
4316. | 3.0897  0.0000  3.6189  0   0   5   |
4317. | 2.6333  1.0986  2.8904  0   0   5   |
4318. | 2.8029  0.6931  3.3322  1   .5   6   |
4319. | 2.7411  0.6931  3.3673  1   .5   6   |
4320. | 2.7576  0.0000  3.2189  0   0   6   |
4321. | 2.7438  0.6931  3.6376  1   .5   6   |
4322. | 2.6561  0.6931  3.8286  1   .5   6   |

```

Least Squares Estimation

$$\begin{aligned} \hat{\beta} &= \left(\sum_{g=1}^G \sum_{i=1}^{n_g} X_{ig} X'_{ig} \right)^{-1} \left(\sum_{g=1}^G \sum_{i=1}^{n_g} X_{ig} Y_{ig} \right) \\ &= \left(\sum_{g=1}^G \mathbf{X}'_g \mathbf{X}_g \right)^{-1} \left(\sum_{g=1}^G \mathbf{X}'_g \mathbf{Y}_g \right) \\ &= (\mathbf{X}' \mathbf{X})^{-1} (\mathbf{X}' \mathbf{Y}) \end{aligned}$$

Individual Level:	$\hat{e}_{ig} = Y_{ig} - X'_{ig}\hat{\beta}$
Group Level:	$\hat{e}_g = \mathbf{Y}_g - \mathbf{X}_g\hat{\beta}$

The model is a linear regression under the assumption

$$\mathbb{E}(\mathbf{e}_g | \mathbf{X}_g) = 0$$

This is the same as assuming that the individual errors are conditionally mean zero

$$\mathbb{E}(e_{ig} | \mathbf{X}_g) = 0$$

Assumption: The clusters (Y_g, X_g) are mutually independent across clusters g .

$$\hat{\beta} - \beta = \left(\sum_{g=1}^G \mathbf{X}'_g \mathbf{X}_g \right)^{-1} \left(\sum_{g=1}^G \mathbf{X}'_g \mathbf{e}_g \right)$$

Theorem: In the clustered linear regression model $\mathbb{E}(\hat{\beta} | \mathbf{X}) = \beta$

(Conditional) Variance

Since the observations are *independent across clusters*

$$\begin{aligned} \text{var} \left(\left(\sum_{g=1}^G \mathbf{X}'_g \mathbf{e}_g \right) | \mathbf{X} \right) &= \sum_{g=1}^G \text{var}(\mathbf{X}'_g \mathbf{e}_g | \mathbf{X}_g) \\ &= \sum_{g=1}^G \mathbf{X}'_g \mathbb{E}(\mathbf{e}_g \mathbf{e}'_g | \mathbf{X}_g) \mathbf{X}_g \\ &= \sum_{g=1}^G \mathbf{X}'_g \Sigma_g \mathbf{X}_g \\ &\stackrel{\text{def}}{=} \Omega_n \end{aligned}$$

$$\mathbf{V}_{\hat{\beta}} = \text{var}(\hat{\beta} | \mathbf{X}) = (\mathbf{X}' \mathbf{X})^{-1} \Omega_n (\mathbf{X}' \mathbf{X})^{-1}$$

Estimator: $\widehat{\mathbf{V}}_{\hat{\beta}} = \text{var}(\hat{\beta} | \mathbf{X}) = (\mathbf{X}' \mathbf{X})^{-1} \Omega_n (\mathbf{X}' \mathbf{X})^{-1}$

$$\begin{aligned}\hat{\Omega}_n &= \sum_{g=1}^G X_g' \hat{e}_g \hat{e}_g' X_g \\ &= \sum_{g=1}^G \sum_{i=1}^{n_g} \sum_{\ell=1}^{n_g} X_{ig} X_{\ell g}' \hat{e}_{ig} \hat{e}_{\ell g} \\ &= \sum_{g=1}^G \left(\sum_{i=1}^{n_g} X_{ig} \hat{e}_{ig} \right) \left(\sum_{\ell=1}^{n_g} X_{\ell g} \hat{e}_{\ell g} \right)'\end{aligned}$$

The factor $G/(G - 1)$ was derived by Chris Hansen (2007) in the context of equal-sized clusters to improve performance when the number of clusters G is small. The factor $(n - 1)/(n - k)$ is an ad hoc generalization which nests the adjustment used in HC1, since when $G = n$ we have the simplification $a_n = n/(n - k)$.

$$a_n = \left(\frac{n-1}{n-k} \right) \left(\frac{G}{G-1} \right)$$

The screenshot shows two side-by-side Jupyter Notebook cells. Both cells contain Stata code for estimating a regression model.

Left Cell (Robust Standard Errors):

```
[34]: % stata -qui  
#delimit ;  
quietly reg FKG log_num_authors log_num_pages both_genders prop_women  
`journals' `jel_im` y_2-y_20 c_2-c_215 jel_flag, robust;  
matrix subV = (e(V)[1,1], e(V)[1,2], e(V)[1,3], e(V)[1,4], e(V)[1,262] \\\  
e(V)[2,1], e(V)[2,2], e(V)[2,3], e(V)[2,4], e(V)[2,262] \\\  
e(V)[3,1], e(V)[3,2], e(V)[3,3], e(V)[3,4], e(V)[3,262] \\\  
e(V)[4,1], e(V)[4,2], e(V)[4,3], e(V)[4,4], e(V)[4,262] \\\  
e(V)[262,1], e(V)[262,2], e(V)[262,3], e(V)[262,4], e(V)[262,262]);  
#delimit cr
```

Right Cell (Clustered Standard Errors):

```
[36]: % stata -qui  
#delimit ;  
quietly reg FKG log_num_authors log_num_pages both_genders prop_women  
`journals' `jel_im` y_2-y_20 c_2-c_215 jel_flag, vce(cluster cluster);  
matrix subV = (e(V)[1,1], e(V)[1,2], e(V)[1,3], e(V)[1,4], e(V)[1,262] \\\  
e(V)[2,1], e(V)[2,2], e(V)[2,3], e(V)[2,4], e(V)[2,262] \\\  
e(V)[3,1], e(V)[3,2], e(V)[3,3], e(V)[3,4], e(V)[3,262] \\\  
e(V)[4,1], e(V)[4,2], e(V)[4,3], e(V)[4,4], e(V)[4,262] \\\  
e(V)[262,1], e(V)[262,2], e(V)[262,3], e(V)[262,4], e(V)[262,262]);  
#delimit cr
```

Output:

Left Cell Output:

```
[35]: % stata matrix list subV  
  
symmetric subV[5,5]  
          c1      c2      c3      c4      c5  
r1  .00004125  
r2  -5.582e-07  .00008567  
r3  -.00002109  6.835e-07  .00005946  
r4  .00001795 -3.165e-06 -.00004688  .0001228  
r5  -.00001192 -.00024076  6.695e-06 -.00001518  .0009432
```

Right Cell Output:

```
[37]: % stata matrix list subV  
  
symmetric subV[5,5]  
          c1      c2      c3      c4      c5  
r1  9.002e-06  
r2  -8.521e-06  .000024048  
r3  2.477e-06  -6.824e-06  .00001387  
r4  .00001121  -.00001846  4.847e-06  .00003053  
r5  7.022e-06  -.00003725  -.00002375  5.047e-06  .00025911
```

A Review of Large Sample Asymptotics

1. [Modes of Convergence](#)
2. [Weak Law of Large Numbers](#)
3. [Central Limit Theorem](#)
4. [Continuous Mapping Theorem](#)
5. [Delta Method](#)
6. [Smooth Function Model](#)
7. [Stochastic Order Symbols](#)

Modes of Convergence

Definition: Convergence in Probability

A random vector $Z_n \in \mathbb{R}^k$ **converges in probability** to Z as $n \rightarrow \infty$, denoted $Z_n \rightarrow_p Z$ or alternatively $\text{plim}_{n \rightarrow \infty} Z_n = Z$, if for all $\delta > 0$,

$$\lim_{n \rightarrow \infty} \mathbb{P}[\|Z_n - Z\| \leq \delta] = 1.$$

We call Z the **probability limit** (or **plim**) of Z_n .

☒ This definition treats random variables and random vectors simultaneously using the vector norm.

☒ It is useful to know that for a random vector, the definition holds if and only if each element in the vector converges in probability to its limit.

Definition: Convergence in Distribution

Let Z_n be a random vector with distribution $F_n(u) = \mathbb{P}[Z_n \leq u]$. We say that Z_n **converges in distribution** to Z as $n \rightarrow \infty$, denoted $Z_n \rightarrow_d Z$ if for all u at which $F(u) = \mathbb{P}[Z \leq u]$ is continuous, $F_n(u) \rightarrow F(u)$ as $n \rightarrow \infty$.

☒ We refer to Z and its distribution $F(u)$ as the asymptotic distribution, large sample distribution, or limit distribution of Z_n .

Weak Law of Large Numbers

Theorem: Weak Law of Large Numbers (WLLN)

If $Y_i \in \mathbb{R}^k$ are i.i.d. and $\mathbb{E} \|Y\| < \infty$, then as $n \rightarrow \infty$

$$\bar{Y} = \frac{1}{n} \sum_{i=1}^n Y_i \xrightarrow{p} \mathbb{E}[Y].$$

↗ The WLLN shows that the sample mean \bar{Y} converges in probability to the true population expectation $\mathbb{E}[Y]$. The result applies to any transformation of a random vector with a finite mean.

Theorem:

If $Y_i \in \mathbb{R}^k$ are i.i.d., $h(y): \mathbb{R}^k \rightarrow \mathbb{R}^q$ and $\mathbb{E} \|h(Y)\| < \infty$, then as $n \rightarrow \infty$

$$\hat{\mu} = \frac{1}{n} \sum_{i=1}^n h(Y_i) \xrightarrow{p} \mu = \mathbb{E}[h(Y)].$$

Definition: An estimator $\hat{\theta}$ of θ is consistent if $\hat{\theta} \xrightarrow{p} \theta$ as $n \rightarrow \infty$.

Lemma: For $Y \in \mathbb{R}^m$, $\mathbb{E} \|Y\| < \infty$ if and only if $\mathbb{E}|Y_j| < \infty$ for $j = 1, \dots, m$.

⌚ An estimator which converges in probability to the population value is called **consistent**.

Central Limit Theorem

Theorem: Multivariate Lindeberg-Lévy Central limit Theorem ([CLT](#))

If $Y_i \in \mathbb{R}^k$ are i.i.d. and $\mathbb{E} \|Y\|^2 < \infty$, then as $n \rightarrow \infty$

$$\sqrt{n}(\bar{Y} - \mu) \xrightarrow{d} N(0, V),$$

where $\mu = \mathbb{E}[Y]$ and $V = \mathbb{E}[(Y - \mu)(Y - \mu)']$.

The central limit theorem shows that the distribution of the sample mean is approximately normal in large samples. For some applications it may be useful to notice that the Theorem does not impose any restrictions on V other than that the elements are finite.

Continuous Mapping Theorem

- ✓ Continuous functions are limit-preserving.
- ✓ There are two forms of the [continuous mapping theorem](#), for convergence in probability and convergence in distribution.

Theorem: Continuous Mapping Theorem

Let $Z_n \in \mathbb{R}^k$ and $g: \mathbb{R}^k \rightarrow \mathbb{R}^q$. If $Z_n \xrightarrow{p} c$ as $n \rightarrow \infty$ and $g(u)$ is continuous at c then

$$g(Z_n) \xrightarrow{p} g(c) \text{ as } n \rightarrow \infty.$$

Theorem: Continuous Mapping Theorem

If $Z_n \xrightarrow{d} Z$ as $n \rightarrow \infty$ and $g: \mathbb{R}^m \rightarrow \mathbb{R}^k$ has the set of discontinuity points D_g such that $\mathbb{P}[Z \in D_g] = 0$ then

$$g(Z_n) \xrightarrow{d} g(Z) \text{ as } n \rightarrow \infty.$$

Theorem: Slutsky's Theorem

If $Z_n \xrightarrow{d} Z$ and $c_n \xrightarrow{p} c$ as $n \rightarrow \infty$, then

$$\begin{aligned} Z_n + c_n &\xrightarrow{d} Z + c \\ Z_n c_n &\xrightarrow{d} Zc \\ \frac{Z_n}{c_n} &\xrightarrow{d} \frac{Z}{c}, \quad c \neq 0 \end{aligned}$$

Delta Method

Theorem: [Delta Method](#)

Let $\mu \in \mathbb{R}^k$ and $g: \mathbb{R}^k \rightarrow \mathbb{R}^q$. If $\sqrt{n}(\hat{\mu} - \mu) \xrightarrow{d} \xi$, where $g(u)$ is continuously differentiable in a neighborhood of μ , then as $n \rightarrow \infty$

$$\sqrt{n}(g(\hat{\mu}) - g(\mu)) \xrightarrow{d} \mathbf{G}'\xi,$$

where $\mathbf{G}(u) = \frac{\partial}{\partial u} g(u)'$ and $\mathbf{G} = \mathbf{G}(\mu)$.

In particular, if $\xi \sim N(0, \mathbf{V})$ then as $n \rightarrow \infty$

$$\sqrt{n}(g(\hat{\mu}) - g(\mu)) \xrightarrow{d} N(0, \mathbf{G}'\mathbf{V}\mathbf{G}).$$

 Differentiable functions of asymptotically normal random estimators are asymptotically normal.

Example:

Suppose that $\hat{\mu} \equiv [\hat{\mu}_{n,1}, \hat{\mu}_{n,2}]'$ is a random 2×1 vector that satisfies $\sqrt{n}(\hat{\mu}_n - \mu) \xrightarrow{d} N(0, \mathbf{V})$, where $\mu \equiv [1 \ 2]$, and $\mathbf{V} \equiv \begin{bmatrix} 2 & 1 \\ 1 & 1 \end{bmatrix}$ and we want to derive the asymptotic distribution of the sequence $\{\hat{\mu}_{n,1}^2 + \hat{\mu}_{n,2}^3\}$. Obviously, $g(\hat{\mu}) = \hat{\mu}_{n,1}^2 + \hat{\mu}_{n,2}^3$, and $g(\mu) = 1^2 + 2^3 = 9$, and $\mathbf{G}(u) = \frac{\partial}{\partial u} g(u)' = \left[\frac{\partial}{\partial u_1} g(u) \quad \frac{\partial}{\partial u_2} g(u) \right]' = [2u_1 \quad 3u_2^2]',$ so $\mathbf{G} = \mathbf{G}(\mu) = [2 \cdot 1 \quad 3 \cdot 2^2]'$ so $\mathbf{G}'\mathbf{V}\mathbf{G} = [2 \quad 12] \begin{bmatrix} 2 & 1 \\ 1 & 1 \end{bmatrix} \begin{bmatrix} 2 & \\ & 12 \end{bmatrix} = 200$, and we conclude that $\sqrt{n}(\hat{\mu}_{n,1}^2 + \hat{\mu}_{n,2}^3 - 9) \xrightarrow{d} N(0, 200)$.

Note: See other exercises [here](#).

Smooth Function Model

The smooth function model is $\theta = g(\mu)$ where $\mu = \mathbb{E}[h(Y)]$ and $g(\mu)$ is smooth in a suitable sense (see below).

$$\theta = g(\mathbb{E}[h(Y)]).$$

\triangle The parameter of interest θ is a smooth function of a population mean. It is not a population moment so it does not have a direct moment estimator. Instead, we use the 'plug-in' principle and replace μ by its natural moment estimator, i.e.,

$$\begin{aligned}\hat{\mu} &= n^{-1} \sum_{i=1}^n h(Y_i), \\ \hat{\theta} &= g\left(n^{-1} \sum_{i=1}^n h(Y_i)\right).\end{aligned}$$

Theorem: If $Y_i \in \mathbb{R}^m$ are i.i.d., $h(u): \mathbb{R}^m \rightarrow \mathbb{R}^k$, $\mathbb{E} \|h(Y)\| < \infty$ and $g(u): \mathbb{R}^k \rightarrow \mathbb{R}^q$ is continuous at μ , then $\hat{\theta} \xrightarrow{p} \theta$ as $n \rightarrow \infty$.

Theorem: If $Y_i \in \mathbb{R}^m$ are i.i.d., $h(u): \mathbb{R}^m \rightarrow \mathbb{R}^k$, $\mathbb{E} \|h(Y)\|^2 < \infty$, $g(u): \mathbb{R}^k \rightarrow \mathbb{R}^q$, and $G(u) = \frac{\partial}{\partial u} g(u)'$ is continuous in a neighborhood of μ then as $n \rightarrow \infty$

$$\sqrt{n}(\hat{\theta} - \theta) \xrightarrow{d} N(0, V_\theta).$$

where $V_\theta = \mathbb{E}[G(\mu) G(\mu)'] - \mathbb{E}[G(\mu)] \mathbb{E}[G(\mu)']$

Example: Let $\{Y_1, \dots, Y_n\}$ be a random sample taken from the population distribution $F(x) = \Pr\{Y \leq x\}$. Recall that the **Empirical CDF** is defined as $\hat{F}(x) = n^{-1} \sum_{i=1}^n 1\{Y_i \leq x\}$. Using this notation, we simply set

$$g(u) = u$$

$$h(x) = 1\{Y \leq x\}$$

$$\mu = \mathbb{E}[h(x)] = \mathbb{E}[1\{Y \leq x\}] = \int 1\{Y \leq x\} dF(Y) = F(x)$$

$$\hat{\mu} = n^{-1} \sum_{i=1}^n 1\{Y_i \leq x\}.$$

Clearly $\mathbb{E}|1\{Y \leq x\}| \leq 1 < \infty$, so we can conclude by the Theorem that $\hat{\mu} \xrightarrow{p} \mu$ or equivalently that $\hat{F}(x) \xrightarrow{p} F(x)$.

The Empirical CDF is a consistent estimator of the unknown CDF.

Similarly, by the second Theorem we have that

$$\sqrt{n}[\hat{F}(x) - F(x)] \xrightarrow{d} N(0, F(x)[1 - F(x)]) \text{ as } n \rightarrow \infty.$$

\square Consistency requires that $h(Y)$ has a finite expectation; asymptotic normality requires that $h(Y)$ has a finite variance.

\square Consistency requires that $g(u)$ be continuous; asymptotic normality requires that $g(u)$ is continuously differentiable.

Stochastic Order Symbols

Simple symbols for random variables and vectors which converge in probability to zero or are stochastically bounded.

"small oh-P-one"	"big oh-P-one"
<p>Let Z_n and $a_n, n = 1, 2, \dots$ be sequences of random variables and constants respectively. The notation $Z_n = o_p(1)$ ("small oh-P-one") means that $Z_n \xrightarrow{p} 0$ as $n \rightarrow \infty$. We also write $Z_n = o_p(a_n)$ if $a_n^{-1}Z_n = o_p(1)$.</p>	<p>Similarly, the notation $Z_n = O_p(1)$ ("big oh-P-one") means that Z_n is <i>bounded in probability</i>. Precisely, for any $\epsilon > 0$ there is a constant $M_\epsilon < \infty$ such that</p> $\limsup_{n \rightarrow \infty} \mathbb{P}[Z_n > M_\epsilon] \leq \epsilon.$ <p>Furthermore, we write $Z_n = O_p(a_n)$ if $a_n^{-1}Z_n = O_p(1)$.</p>

- ☒ $O_p(1)$ is weaker than $o_p(1)$ in the sense that $Z_n = o_p(1)$ implies $Z_n = O_p(1)$ but not the reverse.
- ☒ However, if $Z_n = O_p(a_n)$ then $Z_n = o_p(b_n)$ for any b_n such that $a_n/b_n \rightarrow 0$.

Algebra of Stochastic Orders

$$o_p(1) + o_p(1) = o_p(1)$$

$$o_p(1) + O_p(1) = O_p(1)$$

$$O_p(1) + o_p(1) = O_p(1)$$

$$o_p(1)o_p(1) = o_p(1)$$

$$o_p(1)O_p(1) = o_p(1)$$

$$O_p(1)o_p(1) = O_p(1)$$

- If $X_n = O_p(f_n)$, and $Y_n = O_p(g_n)$ then

$$X_n Y_n = O_p(f_n g_n),$$

$$X_n + Y_n = O_p(\max(f_n, g_n)).$$

- If $X_n = O_p(f_n)$, and $Y_n = o_p(g_n)$, then $X_n Y_n = o_p(f_n g_n)$.

Example (Sample Mean): Let $\{Z_i\}_{i=1}^n$ be a random sample from a distribution such $\mathbb{E}[Z] = 0$ and $\text{var}[Z] = 1$ and let $X_n = n^{-1} \sum_{i=1}^n Z_i$:

$$X_n = \frac{1}{\sqrt{n}} \left(\frac{1}{\sqrt{n}} \sum_{i=1}^n Z_i \right) = \frac{1}{\sqrt{n}} A_n,$$

$\mathbb{E}[A_n] = 0$, and $\text{var}[A_n] = 1$, it then follows by the [Chebyshev's inequality](#), i.e. $\forall \delta > 0$

$$\begin{aligned} \mathbb{P}\{|A_n| \geq \delta\} &\leq \mathbb{E}(|A_n|^2)/\delta^2 \\ &\leq \delta^{-2}, \\ \mathbb{P}\{|A_n| \geq M_\epsilon\} &\leq \epsilon \end{aligned}$$

where $M_\epsilon \equiv \delta$, and $\epsilon \equiv \delta^{-2}$. Therefore $A_n = O_p(1)$ and

$$X_n = \frac{1}{\sqrt{n}} O_p(1) = O(n^{-1/2}) O_p(1) = O_p(n^{-1/2}).$$

Consistency: OLS Estimator

Theorem: Consistency of the OLS

Under Assumption 1, $\widehat{\mathbf{Q}}_{XX} \xrightarrow{p} \mathbf{Q}_{XX}$, $\widehat{\mathbf{Q}}_{XY} \xrightarrow{p} \mathbf{Q}_{XY}$, $\widehat{\mathbf{Q}}_{XX}^{-1} \xrightarrow{p} \mathbf{Q}_{XX}^{-1}$, $\widehat{\mathbf{Q}}_{Xe} \xrightarrow{p} \mathbf{0}$ and $\widehat{\beta} \xrightarrow{p} \beta$ as $n \rightarrow \infty$.

Proof: Recall that $\widehat{\beta} - \beta = \widehat{\mathbf{Q}}_{XX}^{-1} \widehat{\mathbf{Q}}_{Xe}$ where $\widehat{\mathbf{Q}}_{Xe} = \frac{1}{n} \sum_{i=1}^n X_i e_i$. The WLLN and the fact that β is the linear projection coefficient imply $\widehat{\mathbf{Q}}_{Xe} \xrightarrow{p} \mathbb{E}[Xe] = \mathbf{0}$. Therefore $\widehat{\beta} - \beta = \widehat{\mathbf{Q}}_{XX}^{-1} \widehat{\mathbf{Q}}_{Xe} \xrightarrow{p} \mathbf{Q}_{XX}^{-1} \times \mathbf{0} = \mathbf{0}$.

⚠ The consistency result can be written as $\widehat{\beta} = \beta + o_p(1)$

Assumption 1

1. The variables $(Y_i, X_i), i = 1, \dots, n$, are i.i.d.

2. $\mathbb{E}[Y^2] < \infty$.

These two assumptions make sure that Y and X have finite means, variances, and covariances.

3. $\mathbb{E} \|X\|^2 < \infty$

This implies that the columns of \mathbf{Q}_{XX} are linearly independent and that it is invertible.

4. $\mathbf{Q}_{XX} = \mathbb{E}[XX']$ is positive definite.

Asymptotic Normality

Assumption 2

1. The variables $(Y_i, X_i), i = 1, \dots, n$, are i.i.d.

2. $\mathbb{E}[Y^4] < \infty$.

3. $\mathbb{E} \|X\|^4 < \infty$

4. $\mathbf{Q}_{XX} = \mathbb{E}[XX']$ is positive definite.

Expectation Inequality: For any random vector $Y \in \mathbb{R}^m$ with $\mathbb{E} \|Y\| < \infty$ then
 $\|\mathbb{E}[Y]\| \leq \mathbb{E} \|Y\|$

Matrix norm:
 $\|A\|_F = \|\text{vec}(A)\|$
 $= (\text{tr}(A'A))^{1/2}$
 $= \left(\sum_{i=1}^m \sum_{j=1}^k a_{ij}^2 \right)^{1/2}$
For any $m \times 1$ vector a , one has
 $\|aa'\|_F = \|a\|^2$

Firstly

$$\sqrt{n}(\hat{\beta} - \beta) = \left(\frac{1}{n} \sum_{i=1}^n X_i X_i' \right)^{-1} \left(\frac{1}{\sqrt{n}} \sum_{i=1}^n X_i e_i \right)$$

☞ The product $X_i e_i$ is i.i.d.
(since the observations (Y_i, X_i) are i.i.d.) and mean zero (since $\mathbb{E}[Xe] = 0$).

Theorem:

Under Assumption 2, $\Omega < \infty$, and

$$\frac{1}{\sqrt{n}} \sum_{i=1}^n X_i e_i \xrightarrow{d} N(0, \Omega)$$

as $n \rightarrow \infty$.

In order to show that the elements of Ω are finite we only need to show that $\|\Omega\| < \infty$. By the expectation inequality, the [cauchy-schwarz inequality](#), and Assumption 2

$\|\Omega\| \leq \mathbb{E}\|XX'e^2\| = \mathbb{E}[\|X\|^2 e^2] \leq (\mathbb{E}\|X\|^4)^{1/2} (\mathbb{E}[e^4])^{1/2} < \infty$. Therefore, the finiteness of the covariance matrix means that the CLT applies.

Therefore

$$\begin{aligned} \sqrt{n}(\hat{\beta} - \beta) &\rightarrow \mathbf{Q}_{XX}^{-1} \times N(0, \Omega) \\ &= N(0, \mathbf{Q}_{XX}^{-1} \Omega \mathbf{Q}_{XX}^{-1}) \end{aligned}$$

💡 Linear combinations of normal vectors are normal

Theorem:

Under Assumption 2, as $n \rightarrow \infty$

$$\sqrt{n}(\hat{\beta} - \beta) \xrightarrow{d} N(0, \mathbf{V}_\beta)$$

where $\mathbf{Q}_{XX} = \mathbb{E}[XX']$, $\Omega = \mathbb{E}[XX'e^2]$

and

$$\mathbf{V}_\beta = \mathbf{Q}_{XX}^{-1} \Omega \mathbf{Q}_{XX}^{-1}$$

- $\hat{\beta} = \beta + O_p(n^{-1/2})$

- \mathbf{V}_β is referred as the **asymptotic covariance matrix** of $\hat{\beta}$.

- Recall the *finite-sample* conditional variance in the CEF model

$$\mathbf{V}_{\hat{\beta}} = \text{var}[\hat{\beta} | \mathbf{X}] = (\mathbf{X}'\mathbf{X})^{-1} (\mathbf{X}'\mathbf{D}\mathbf{X}) (\mathbf{X}'\mathbf{X})^{-1}$$

$$n\mathbf{V}_{\hat{\beta}} = \left(\frac{1}{n} \mathbf{X}'\mathbf{X} \right)^{-1} \left(\frac{1}{n} \mathbf{X}'\mathbf{D}\mathbf{X} \right) \left(\frac{1}{n} \mathbf{X}'\mathbf{X} \right)^{-1}$$

As $n \rightarrow \infty$, $n\mathbf{V}_{\hat{\beta}} \xrightarrow{p} \mathbf{V}_\beta$

$$\mathbf{X}'\mathbf{D}\mathbf{X} = \sum_{i=1}^n X_i X_i' \sigma_i^2$$

Heteroskedastic Covariance Matrix Estimation

The moment estimator for Ω is

$$\widehat{\Omega} = \frac{1}{n} \sum_{i=1}^n X_i X_i' \hat{e}_i^2$$

Leading to the plug-in covariance matrix estimator

$$\widehat{\mathbf{V}}_{\beta}^{\text{HCO}} = \widehat{\mathbf{Q}}_{XX}^{-1} \widehat{\Omega} \widehat{\mathbf{Q}}_{XX}^{-1}$$

Remember the HCO covariance matrix estimator:

$$\widehat{\mathbf{V}}_{\widehat{\beta}}^{\text{HCO}} = (\mathbf{X}' \mathbf{X})^{-1} \left(\sum_{i=1}^n X_i X_i' \hat{e}_i^2 \right) (\mathbf{X}' \mathbf{X})^{-1}$$

! Notice that $\widehat{\mathbf{V}}_{\beta}^{\text{HCO}} = n \widehat{\mathbf{V}}_{\widehat{\beta}}^{\text{HCO}}$

Theorem: Under Assumption 2, as $n \rightarrow \infty$, $\widehat{\Omega} \xrightarrow{p} \Omega$ and $\widehat{\mathbf{V}}_{\beta}^{\text{HCO}} \xrightarrow{p} \mathbf{V}_{\beta}$.

Proof: Observe that

$$\begin{aligned} \widehat{\Omega} &= \frac{1}{n} \sum_{i=1}^n X_i X_i' \hat{e}_i^2 \\ &= \frac{1}{n} \sum_{i=1}^n X_i X_i' e_i^2 + \frac{1}{n} \sum_{i=1}^n X_i X_i' (\hat{e}_i^2 - e_i^2) \end{aligned}$$

$$\hat{e}_i^2 = e_i^2 - 2e_i X_i' (\hat{\beta} - \beta) + (\hat{\beta} - \beta)' X_i X_i' (\hat{\beta} - \beta)$$

By the WLLN

$$\frac{1}{n} \sum_{i=1}^n X_i X_i' e_i^2 \xrightarrow{p} \mathbb{E}[XX' e^2] = \Omega$$

Triangle Inequality: $\left| \sum_{j=1}^m x_j \right| \leq \sum_{j=1}^m |x_j|$

Triangle Inequality (matrices): $\| A + B \| \leq \| A \| + \| B \|$

Cauchy-Schwarz Inequality: $|a'b| \leq \| a \| \| b \|$

Hölder Inequality: $\mathbb{E} \| X' Y \| \leq (\mathbb{E} \| X \|^p)^{1/p} (\mathbb{E} \| Y \|)^{1/q}$ for $p > 1$, $q > 1$, and $\frac{1}{p} + \frac{1}{q} = 1$

$$\begin{aligned} \text{Matrix norm:} \\ \| A \|_F &= \|\text{vec}(A) \| \\ &= (\text{tr}(A'A))^{1/2} \\ &= \left(\sum_{i=1}^m \sum_{j=1}^k a_{ij}^2 \right)^{1/2} \end{aligned}$$

For any $m \times 1$ vector a , one has $\| a a' \|_F = \| a \|^2$

We simply need to show that $\frac{1}{n} \sum_{i=1}^n X_i X_i' (\hat{e}_i^2 - e_i^2) \xrightarrow{p} 0$

$$\begin{aligned} \left\| \frac{1}{n} \sum_{i=1}^n X_i X_i' (\hat{e}_i^2 - e_i^2) \right\| &\leq \frac{1}{n} \sum_{i=1}^n \| X_i X_i' (\hat{e}_i^2 - e_i^2) \| \\ &= \frac{1}{n} \sum_{i=1}^n \| X_i \|^2 |\hat{e}_i^2 - e_i^2| \end{aligned}$$

$$\begin{aligned} \left\| \frac{1}{n} \sum_{i=1}^n X_i X_i' (\hat{e}_i^2 - e_i^2) \right\| &\leq 2 \left(\frac{1}{n} \sum_{i=1}^n \| X_i \|^3 |e_i| \right) \| \hat{\beta} - \beta \| + \left(\frac{1}{n} \sum_{i=1}^n \| X_i \|^4 \right) \| \hat{\beta} - \beta \|^2 \\ &= o_p(1) \end{aligned}$$

$$\begin{aligned} |\hat{e}_i^2 - e_i^2| &\leq 2|e_i X_i' (\hat{\beta} - \beta)| + (\hat{\beta} - \beta)' X_i X_i' (\hat{\beta} - \beta) \\ &= 2|e_i| |X_i' (\hat{\beta} - \beta)| + |(\hat{\beta} - \beta)' X_i|^2 \\ &\leq 2|e_i| \| X_i \| \| \hat{\beta} - \beta \| + \| X_i \|^2 \| \hat{\beta} - \beta \|^2 \end{aligned}$$

$$\begin{aligned} \mathbb{E} [\| X \|^3 |e_i|] &\leq \left(\mathbb{E} [(\| X \|^3)^{4/3}] \right)^{3/4} (\mathbb{E}[e^4])^{1/4} \\ &= (\mathbb{E} \| X \|^4)^{3/4} (\mathbb{E}[e^4])^{1/4} < \infty \end{aligned}$$

S.E., t-Statistic, C.I., p -values

Standard Error	t -Statistics	Confidence Intervals
$s(\hat{\beta}) = \sqrt{n^{-1} \hat{V}_\beta^2}$	$T(\beta) = \frac{\hat{\beta} - \beta}{s(\hat{\beta})}$	$\hat{C} = [\hat{\beta} - c \times s(\hat{\beta}), \hat{\beta} + c \times s(\hat{\beta})]$
△ Because this is based on <i>asymptotic</i> theory we call this <u>asymptotic standard errors</u> .	Recall that $\sqrt{n}(\hat{\beta} - \beta) \xrightarrow{d} N(0, V_\beta)$ and $\hat{V}_\beta^2 \xrightarrow{p} V_\beta$, thus $\begin{aligned} T(\beta) &= \frac{\hat{\beta} - \beta}{s(\hat{\beta})} \\ &= \frac{\sqrt{n}(\hat{\beta} - \beta)}{\sqrt{\hat{V}_\beta}} \\ &\xrightarrow{d} \frac{N(0, V_\beta)}{\sqrt{V_\beta}} \\ &= Z \sim N(0, 1) \end{aligned}$	⌚ For $c = 1.965$ we have $\mathbb{P}[\beta \in \hat{C}] \rightarrow 0.95$.
	p -value	
	We are interested in testing $H_0: \beta = 0$ against $H_1: \beta \neq 0$. Then under H_0 one has $T(0) \equiv \hat{\beta}/s(\hat{\beta}) \xrightarrow{d} N(0, 1)$. Then the test taking the form "Reject H_0 if $ T(0) > c$ " has <u>p-value</u> defined as $p = 2(1 - \Phi(T(0)))$. A test with significance level α can be restated as "Reject H_0 if $p < \alpha$ ".	

002_ols_inference.ipynb

```
Performing the OLS regression of Y on X using Stata and saving a sub-vector of the original  $\hat{\beta}$  and its corresponding submatrix  $\hat{V}_\beta$  (Note that row and/or column names are missing in some cases):
```

```
[3]: stata -qui
#delimit ;
reg FGS log_num_authors log_num_pages both_genders prop_women
----- journals` "[jel_impr] y2=y,20 c2_c215 [jel_flag, vce(cluster cluster)];
matrix b = e(b)[1,"log_num_authors"],e(b)[1,"log_num_pages"],
          e(b)[1,"both_genders"],e(b)[1,"prop_women"],e(b)[1,"_cons"];
matrix V = (e(V[1,1]), e(V[1,2]), e(V[1,3]), e(V[1,4]), e(V[1,5]), e(V[1,6]),
           e(V[2,1]), e(V[2,2]), e(V[2,3]), e(V[2,4]), e(V[2,5]), e(V[2,6]),
           e(V[3,1]), e(V[3,2]), e(V[3,3]), e(V[3,4]), e(V[3,5]), e(V[3,6]),
           e(V[4,1]), e(V[4,2]), e(V[4,3]), e(V[4,4]), e(V[4,5]), e(V[4,6]),
           e(V[5,1]), e(V[5,2]), e(V[5,3]), e(V[5,4]), e(V[5,5]), e(V[5,6]));
matrix rnames V = log_num_authors log_num_pages both_genders prop_women _cons;
matrix colnames V = log_num_authors log_num_pages both_genders prop_women _cons;
#delimit cr
```

```
[4]: #state matrix list b
#state matrix list V
```

```
b[1,5]
log_num_authors log_num_pages both_genders prop_women _cons
y1 -.00397377 .01915903 .00059009 -.01889331 2.7023992
```

```
symmetric V[5,5]
log_num_authors log_num_pages both_genders prop_women
log_num_authors 8.062e-06
log_num_pages -.521e-06 .00002404
both_genders 2.477e-06 -.624e-06 .00001387
prop_women .00001121 -.00001846 4.847e-06 .00003053
_cons 7.022e-06 -.00003725 -.00002375 5.047e-06
```

```
_cons
_cons .00025911
```

```
Simple 0 1 @ main Python 3 (ipykernel) | Idle Mode: Command ↵ Ln 1, Col 1 002_ols_inference.ipynb
```

002_ols_inference.ipynb

```
002_ols_inference.ipynb
```

```
File Edit View Run Kernel Git Tabs Settings Help
```

t-Statistics & p-Values

```
Printing the estimation results for these subset of coefficients of interest:
```

```
[5]: stata ereturn post b V
stata ereturn display, 1(90)
```

	Coefficient	Std. err.	z	P> z	[90% conf. interval]
log_num_authors	-.0039738	.0030103	-1.32	0.187	-.0009253 .0009778
log_num_pages	.019159	.0049032	3.91	0.000	.010941 .027224
both_genders	.0005981	.0037246	0.16	0.872	-.0055284 .0067246
prop_women	-.0188933	.0055253	2.14	0.034	-.0279818 .000905
_cons	2.702399	.016097	167.88	0.000	2.673922 2.728876

```
Performing the test of the null hypothesis  $H_0: \beta_{prop\_women} = 0$  against the alternative  $H_1: \beta_{prop\_women} \neq 0$ .
```

```
[6]: stata
capture scalar drop t
scalar T = b[prop_women]/_se[prop_women]
di _n "T(prop_women) = " T
di _n "Prob > |T| = " 2*(1-normal(abs(T)))
```

```
. capture scalar drop t
. scalar T = b[prop_women]/_se[prop_women]
. di _n "T(prop_women) = " T
T(prop_women) = -3.4194328
. di _n "Prob > |T| = " 2*(1-normal(abs(T)))
Prob > |T| = .00062753
```

```
Simple 0 1 @ main Python 3 (ipykernel) | Idle Mode: Command ↵ Ln 1, Col 1 002_ols_inference.ipynb
```

002_ols_inference.ipynb

```
002_ols_inference.ipynb
```

```
File Edit View Run Kernel Git Tabs Settings Help
```

```
[8]: stata ereturn post b V
stata ereturn display, 1(90)
```

	Coefficient	Std. err.	z	P> z	[90% conf. interval]
log_num_authors	-.0039738	.0030103	-1.32	0.187	-.0009253 .0009778
log_num_pages	.019159	.0049032	3.91	0.000	.010941 .027224
both_genders	.0005981	.0037246	0.16	0.872	-.0055284 .0067246
prop_women	-.0188933	.0055253	2.14	0.034	-.0279818 .000905
_cons	2.702399	.016097	167.88	0.000	2.673922 2.728876

```
Manually calculating the 90% = (1 -  $\alpha$ )  $\times 100$  confidence interval for  $\beta_{prop\_women}$  as
 $\hat{C} = [\hat{\beta}_{prop\_women} - c_a \times (\hat{\beta}_{prop\_women}), \hat{\beta}_{prop\_women} + c_a \times (\hat{\beta}_{prop\_women})]$  where  $c_a = F^{-1}(1 - \alpha/2)$  and  $F(\cdot)$  represents the cumulative distribution function of a standard normal distribution function.
```

```
[9]: stata
scalar c_min=_b[prop_women] + invnormal(.05)*_se[prop_women]
scalar c_max=_b[prop_women] + invnormal(.95)*_se[prop_women]
display _n "90% C.I. for b[prop_women]: (" c_min ", " c_max ")"
```

the cumulative distribution function of a standard normal distribution function.

```
[9]: stata  
scalar c_min=b[prop_women] + invnormal(0.05)*_se[prop_women]  
scalar c_max=b[prop_women] + invnormal(0.95)*_se[prop_women]  
display _n "90% C.I. for b[prop_women]: (" c_min ", " c_max ")"  
  
. scalar c_min=b[prop_women] + invnormal(0.05)*_se[prop_women]  
. scalar c_max=b[prop_women] + invnormal(0.95)*_se[prop_women]  
  
. display _n "90% C.I. for b[prop_women]: (" c_min ", " c_max ")"  
90% C.I. for b[prop_women]: (-.0279816, .00000503)
```

Simple 0 1 main Python 3 (ipykernel) | Idle

Mode: Command Ln 1, Col 1 002_ols_inference.ipynb

Motivation

The screenshot shows a web browser displaying a research article. The title is "Abstract readability: Evidence from top-5 economics journals". The authors listed are Belicia Rodriguez, Kim P. Huynh, David T. Jacho-Chávez, and Leonardo Sánchez-Aragón. The journal is Economics Letters, volume 235 (2024) 111541. The Elsevier logo is visible at the top left. The abstract discusses the readability of abstracts from top-5 economics journals between 2000-2019, comparing them across various journals, fields, and years.

(Practical) Implications:

As explained earlier, the information contained in Journals, JEL codes, Years, Cluster, and JEL flag will be included in the model in the form of indicator variables (dummy variables):

Journals: corresponds to 4 indicators for ECM, JPE, QJE, RES

JEL Codes: corresponds to 19 indicators for all JEL codes except 'D'

Years: corresponds to 19 indicators for years 2001, 2002,..., 2019

Cluster: corresponds to 214 indicators for clusters 2, 3,..., 215, except cluster '1'

JEL flag: indicator of whether the article includes a JEL classification

⚠ This implies that we will be estimating 257 parameters multiplying these indicator variables, i.e., 257 fixed effects.

Our empirical strategy is based on the following baseline specification for article a :

$$\log(F - K \text{ grade})_a = \beta_1 \times \log(\text{Number authors})_a + \beta_2 \times \log(\text{Number pages})_a + \beta_3 \times \text{Both genders}_a + \theta \times \text{Proportion Women}_a + \text{cons} + e_a$$

- Some journals (Journals) are more likely to publish more technical articles than others, and more technical articles are on average less readable.

- Research in certain fields (JEL codes) require large groups of authors.

- Editors change over time (Years) and so journal taste for some type of articles over others.

- Articles' readability is affected by co-authors (Cluster).

- The method to impute JEL codes (JEL flag) use characteristics such as co-author information.

⚠ This implies that

$$E[e_a | \log(\text{Number authors})_a, \log(\text{Number pages})_a, \text{Both genders}_a, \text{Proportion Women}_a] \neq 0$$

Because Journals, JEL codes, Years, Cluster, and JEL flag are likely correlated with $\log(\text{Number authors})_a, \log(\text{Number pages})_a, \text{Both genders}_a, \text{Proportion Women}_a$, but they have been omitted from the specification, i.e., omitted variable bias.

(Theoretical) Solution:

$$\begin{aligned} \log(F - K \text{ grade})_a &= \beta_1 \times \log(\text{Number authors})_a + \beta_2 \times \log(\text{Number pages})_a \\ &\quad + \beta_3 \times \text{Both genders}_a + \theta \times \text{Proportion Women}_a \\ &\quad + \text{Journals} + \text{JEL codes} + \text{Cluster} + \text{Years} + \text{JEL flag} + \text{cons} + e_a \end{aligned}$$

where

$$E[e_a | \log(\text{Number authors})_a, \log(\text{Number pages})_a, \text{Both genders}_a, \text{Proportion Women}_a, \text{Journals}, \text{JEL Codes}, \text{Cluster}, \text{Years}, \text{JEL flag}] = 0$$

💡 Notice that all variables we are conditioning on are allowed to be correlated among themselves.

```
[1]: %capture
import stata_setup, os
if os.name == 'nt':
    stata_setup.config('C:/Program Files/Stata17','mp')
else:
    stata_setup.config('/usr/local/stata17','mp')

[2]: from sfi import Data
from pydata import stata
import numpy as np
import pandas as pd

[3]: %stata -qui
use ".../data/data", clear
rename log_fleisch_kincaid_grade_level FKG
quietly tabulate year, generate(y)
quietly tabulate cluster, generate(c_)

local journals ecm jpe qje res //AER based category
local jel_imp a_imp b_imp c_imp e_imp f_imp g_imp h_imp i_imp k_imp l_imp m_imp n_imp o_imp p_imp q_imp r_imp y_imp z_imp // D JEL based case
```

```
[4]: %capture captured_output
state.run('
display "journals"
')

[5]: output = captured_output.stdout
lines = output.splitlines()
journals_output = lines[-4].strip()
journals = stata.read_data_from_file(var=journals_output)
print(journals.sample(n=10, random_state=542))

ecm jpe qje res
4707 0 0 0 1
3807 0 0 1 0
4570 0 0 0 1
4551 0 0 0 1
181 0 0 0 0
4817 0 0 0 1
4932 0 0 0 1
700 0 0 0 0
2711 0 1 0 0
3349 0 0 1 0
```

Four screenshots of JupyterLab notebooks showing data processing steps:

Top Left Notebook:

```
[6]: jel_imp = stata.pandasframe_from_data(var="a_imp b_imp c_imp e_imp f_imp g_imp h_imp i_imp j_imp k_imp l_imp m_imp n_imp")
print(jel_imp.sample(n=10, random_state=542))
```

	a_imp	b_imp	c_imp	e_imp	f_imp	g_imp	h_imp	i_imp	j_imp	k_imp	l_imp	m_imp	n_imp
4787	0	0	0	0	0	0	0	0	1	0	0	0	0
3807	0	0	0	0	0	0	0	0	0	0	0	0	0
4579	0	0	0	0	0	0	0	0	1	0	0	0	0
4551	0	0	1	0	0	0	0	0	0	0	0	0	0
101	0	0	0	0	0	0	0	1	0	0	0	0	0
4817	0	0	1	1	0	0	0	0	0	0	0	0	0
4932	0	0	0	0	0	0	0	0	0	1	0	0	0
700	0	0	0	0	0	0	0	0	0	0	1	0	0
2711	0	0	0	0	0	0	0	0	1	1	0	0	0
3349	0	0	1	0	0	0	0	0	0	0	0	0	0

Top Right Notebook:

```
[7]: %%capture captured_output
stata.run('ds y_2-y_20')

[8]: output = captured_output.stdout
year_output = output.strip()
year = stata.pandasframe_from_data(var=year_output)
print(year.sample(n=10, random_state=542))
```

	y_2	y_4	y_5	y_8	y_10	y_12	y_14	y_16	y_18	y_20	y_3	y_5	y_7
4787	0	0	0	0	0	0	0	0	1	0	0	0	0
3807	0	0	0	0	0	0	0	0	0	0	0	0	0
4579	0	0	0	0	0	0	0	0	0	0	0	0	0
4551	0	0	0	0	0	0	0	0	0	0	0	0	0
101	0	0	0	0	0	0	0	0	0	0	1	0	0
4817	0	0	0	0	0	0	0	0	0	0	1	0	0
4932	0	0	0	0	0	0	0	0	0	0	1	0	0
700	0	0	0	0	0	0	0	0	0	0	0	0	0
2711	0	0	0	1	0	0	0	0	0	1	0	0	0
3349	0	1	0	0	0	0	0	0	0	0	0	0	0

Bottom Left Notebook:

```
[9]: %%capture captured_output
stata.run('ds c_2-c_215')

[10]: output = captured_output.stdout
cluster_output = output.strip()
cluster = stata.pandasframe_from_data(var=cluster_output)
print(cluster.sample(n=10, random_state=542))
```

	c_2	c_22	c_42	c_62	c_82	c_102	c_122	c_142	c_162	c_182	...
4787	1	0	0	0	0	0	0	0	0	0	...
3807	1	0	0	0	0	0	0	0	0	0	...
4579	1	0	0	0	0	0	0	0	0	0	...
4551	0	0	0	0	0	0	0	0	0	0	...
101	1	0	0	0	0	0	0	0	0	0	...
4817	1	0	0	0	0	0	0	0	0	0	...
4932	0	0	0	0	0	0	0	0	0	0	...
700	1	0	0	0	0	0	0	0	0	0	...
2711	1	0	0	0	0	0	0	0	0	0	...
3349	1	0	0	0	0	0	0	0	0	0	...

Bottom Right Notebook:

```
[11]: jel_flag = stata.pandasframe_from_data(var="jel_flag")
print(jel_flag.sample(n=10, random_state=542))
```

	jel_flag
4787	1
3807	1
4579	1
4551	0
101	1
4817	1
4932	0
700	1
2711	0
3349	0

High Dimensional Fixed Effects

We want to compute the least squares estimates $\hat{\beta}$ of

$$\mathbf{Y} = \mathbf{X}\boldsymbol{\beta} + \mathbf{D}\boldsymbol{\alpha} + \mathbf{e}.$$

- $\mathbf{D} = [\mathbf{D}_1 \ \mathbf{D}_2 \ \dots \ \mathbf{D}_F]$ consists of F indicator matrices
- If $F = 1$, this collapses to a standard fixed effect regression (`xtreg`, `areg`)
- Can't use dummies because $[\mathbf{D}_2 \dots \mathbf{D}_F]$ is too large

Solution Strategy:

1. Compute the residuals of \mathbf{y} and \mathbf{X} against \mathbf{D} :

$$\begin{aligned}\tilde{\mathbf{Y}} &= \mathbf{M}_{\mathbf{D}}\mathbf{Y} \\ \tilde{\mathbf{X}} &= \mathbf{M}_{\mathbf{D}}\mathbf{X}\end{aligned}$$

2. Apply the Frisch-Waugh-Lovell Theorem:

$$\hat{\boldsymbol{\beta}} = (\tilde{\mathbf{X}}'\tilde{\mathbf{X}})^{-1}\tilde{\mathbf{X}}'\tilde{\mathbf{Y}}.$$

Thus, we can just focus on one variable at a time: $\hat{\mathbf{Y}}$

To obtain $\hat{\mathbf{Y}} = \mathbf{M}_{\mathbf{D}}\mathbf{Y}$, find an $\hat{\boldsymbol{\alpha}}$ that satisfies the normal equations

$$\mathbf{D}'\mathbf{e} = 0, \mathbf{e} \stackrel{\text{def}}{=} \mathbf{Y} - \mathbf{D}\hat{\boldsymbol{\alpha}}.$$

In plain English:

For every level g of every fixed effect f the mean of the residuals must be zero:

$$\bar{e}_i = 0, i \in \mathcal{J}(f, g).$$

Note: We don't care if $\hat{\boldsymbol{\alpha}}$ is unique.

The screenshot shows a Jupyter Notebook interface with a Stata code cell and its corresponding output. The code performs a two-step estimation. First, it uses `reghdfe` to absorb journal fixed effects and then performs a second stage regression. The output provides summary statistics for the regression, including the number of observations, F-statistic, R-squared, and Root MSE. It also displays a coefficient table for the FKG model, showing coefficients, standard errors, t-values, p-values, and 95% confidence intervals for various variables like log_num_authors, log_num_pages, both_genders, prop_women, and cons.

```
%>%%stata
. egen journal1 = group(journal)
. #delimit ;
. reghdfe FKG log_num_authors log_num_pages both_genders prop_women,
>    absorb(journal1 a_imp b_imp c_imp e_imp f_imp g_imp h_imp
>    i_imp j_imp k_imp l_imp m_imp n_imp o_imp p_imp q_imp r_imp y_
>    year ib0.cluster jel_flag) vce(cluster cluster);
. #delimit cr

. egen journal1 = group(journal)

. #delimit ;
delimited now ;
. reghdfe FKG log_num_authors log_num_pages both_genders prop_women,
>    absorb(journal1 a_imp b_imp c_imp e_imp f_imp g_imp h_imp
>    i_imp j_imp k_imp l_imp m_imp n_imp o_imp p_imp q_imp r_imp y_
>    year ib0.cluster jel_flag) vce(cluster cluster);
(dropped 1 singleton observations)
(MWFE estimator converged in 9 iterations)

HDFE Linear regression
Number of obs      =     4,987
Absorbing 23 HDFE groups
F(  4,    214) =      5.90
Statistics robust to heteroskedasticity
Prob > F        =     0.0002
R-squared        =     0.0921
Adj R-squared   =     0.0421
Within R-sq.    =     0.0019
Number of clusters (cluster) =      215
Root MSE         =     0.1642

(Std. err. adjusted for 215 clusters in cluster)

-----| Robust
      | Coefficient  std. err.      t    P>|t|    [95% conf. interval]
-----+-----log_num_aus | -.0039738  .0029441   -1.35   0.179   -.009777  .0018294
log_num_pa-s |  .019159  .0047953   4.00  0.000   .009707  .0286111
both_genders |  .0005981  .0036427   0.16  0.870   -.0065821  .0077782
prop_women |  -.0188933  .0054037   -3.50  0.001   -.0295447  -.0082419
_cons |  2.666831  .0147896  180.32  0.000   2.637679  2.695982
```


Specification Curve Analysis

Recall that we want to compute the least squares estimates $\hat{\beta}$ of

$$\mathbf{Y} = \mathbf{X}\boldsymbol{\beta} + \mathbf{D}\boldsymbol{\alpha} + e,$$

where $\mathbf{D} = [\mathbf{D}_1 \quad \mathbf{D}_2 \quad \cdots \quad \mathbf{D}_F]$ consists of F indicator matrices.

```

1 # Comment lines starting with a hashbang will be ignored
2 # Each choice is independent from others
3 # Choices
4 # - Independent Variable: # reserved keyword
5 #   - Ln(DATE_CODE): log_date_code
6 #   - Ln(K_GRADE): log_flaesch_kincaid_grade_level
7 #   - Ln(COLEMAN_LIN INDEX): log_coleman_linc_index
8 #   - Ln(Automated readability index): log_automated_readability_index
9
10 # Focal Variables: # reserved keyword
11 #   - Only female: only_female
12 #   - Proportion of women: prop_women
13 #   - Only male: only_male
14 #   - Only gender: both_genders
15 #   - Gender: # reserved keyword
16 #     - Baseline: both_genders log_num_authors log_num_pages jel_flag
17 #     - And Union: age grade collgrad whs_rtt_exp tenure hours whs_work_union
18 # Fixed Effects: # reserved keyword AER , D, 2000 THE base cases
19 # Journals: journal a_imp b_imp c_imp d_imp f_imp g_imp h_imp i_imp j_imp k_imp l_imp m_imp n_imp o_imp p_imp q_imp r_imp y_imp z_imp
20 # Year: year
21 # Cluster: lit_cluster
22 # Journals, Year: journal year
23 # Journals, Cluster: journal lit_cluster
24 # Year, Cluster: year lit_cluster
25 # Year, Cluster: year lit_cluster a_imp b_imp c_imp d_imp f_imp g_imp h_imp i_imp j_imp k_imp l_imp m_imp n_imp o_imp p_imp q_imp r_imp y_imp z_imp
26 # Journals, Year, Cluster: journal year lit_cluster a_imp b_imp c_imp d_imp f_imp g_imp h_imp i_imp j_imp k_imp l_imp m_imp n_imp o_imp p_imp q_imp r_imp y_imp z_imp
27 # Year, Cluster: year lit0_cluster
28 # Journals, lit0_code, Year: journal year a_imp b_imp c_imp d_imp f_imp g_imp h_imp i_imp j_imp k_imp l_imp m_imp n_imp o_imp p_imp q_imp r_imp y_imp z_imp
29 # Journals, lit0_code, Cluster: journal lit0_cluster a_imp b_imp c_imp d_imp f_imp g_imp h_imp i_imp j_imp k_imp l_imp m_imp n_imp o_imp p_imp q_imp r_imp y_imp z_imp
30 # Journals, lit0_code, Year, Cluster: journal year lit0_cluster a_imp b_imp c_imp d_imp f_imp g_imp h_imp i_imp j_imp k_imp l_imp m_imp n_imp o_imp p_imp q_imp r_imp y_imp z_imp
31 # Journals, lit0_code, Year, Cluster: journal year lit0_cluster a_imp b_imp c_imp d_imp f_imp g_imp h_imp i_imp j_imp k_imp l_imp m_imp n_imp o_imp p_imp q_imp r_imp y_imp z_imp
32 # Journals, lit0_code, Years, Cluster: journal year lit0_cluster a_imp b_imp c_imp d_imp f_imp g_imp h_imp i_imp j_imp k_imp l_imp m_imp n_imp o_imp p_imp q_imp r_imp y_imp z_imp
33 Standard Error Clustering: # reserved keyword
34 Conditionality: # reserved keyword cluster
35 Subsample: # reserved keyword
36 All sample: num_words_90_flag==1 | num_words_90_flag==2
37 # Abstracts with > 90 words: num_words_90_flag==2
38

```

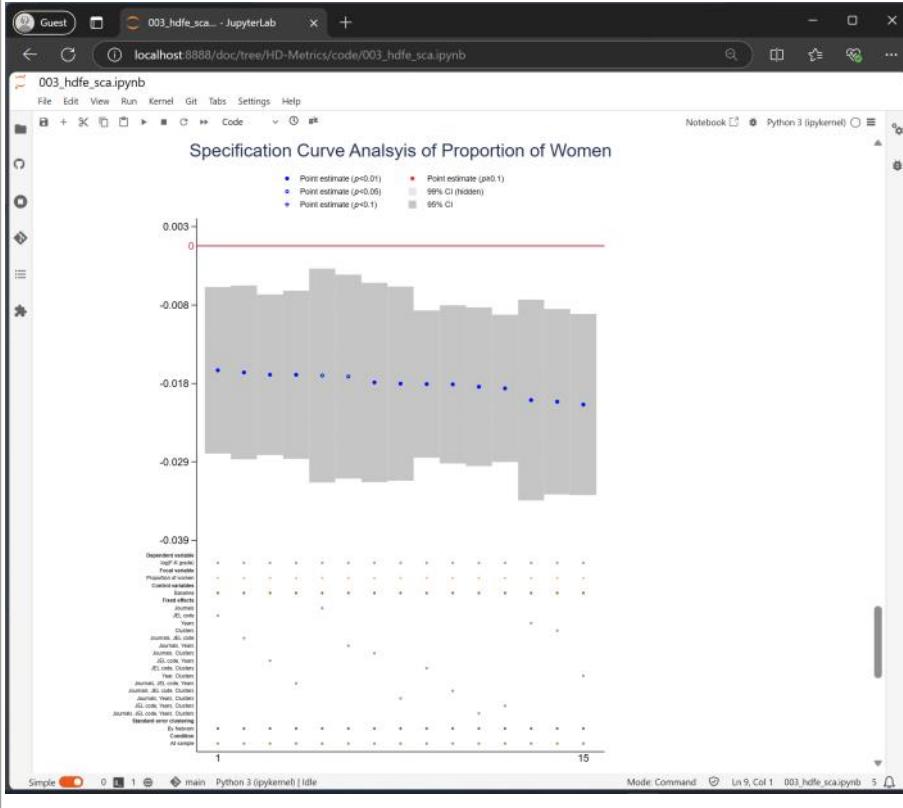
Let's rewrite this model as $\mathbf{Y} = X_1\beta_1 + X_2\beta_2 + \mathbf{D}_1\boldsymbol{\alpha}_1 + \cdots + \mathbf{D}_F\boldsymbol{\alpha}_F + e$, and we are only interested in the effect of the *focal* variable, X_1 , and how robust this effect is to the inclusion or not of $\mathbf{D}_1, \mathbf{D}_2, \dots, \mathbf{D}_F$ individually or grouped.

- X_1 is called the focal variable and kept in the model always, e.g., `prop_women`.
- X_2 are called control variables and kept in the model always, e.g., `log_num_authors`, `log_num_pages`, `both_genders`, `jel_flag`, and the `_const`.

Specification Curve Analysis entails setting:

- $\boldsymbol{\alpha}_l = \mathbf{0}$ for $l = 1, \dots, F$ one at the time.
- $\boldsymbol{\alpha}_l = \mathbf{0}, \boldsymbol{\alpha}_m = \mathbf{0}$ for $l, m = 1, \dots, F; l \neq m$.
- $\boldsymbol{\alpha}_l = \mathbf{0}, \boldsymbol{\alpha}_m = \mathbf{0}, \boldsymbol{\alpha}_n = \mathbf{0}$ for $l, m, n = 1, \dots, F; l \neq m \neq n$.
- ...

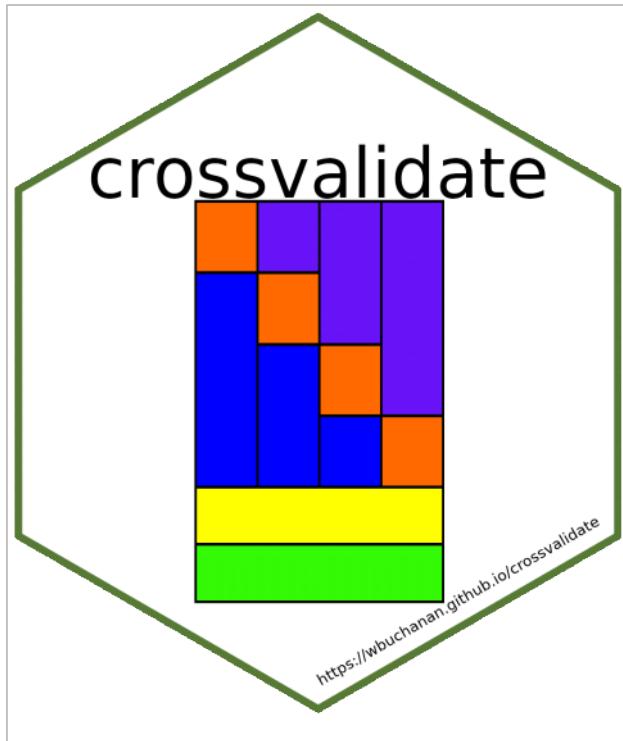
and see how the magnitude and statistical significance of $\hat{\beta}_1$ is affected by these different specifications.



Resampling Methods

There are two most commonly used *resampling methods*:

- ✓ *Cross-validation*,
- ✓ *Bootstrap*.



Source: <https://github.com/wbuchanan/crossvalidate>
Capabilities: <https://411steven.github.io/stataConference2024/>

☞ Package is still under development.

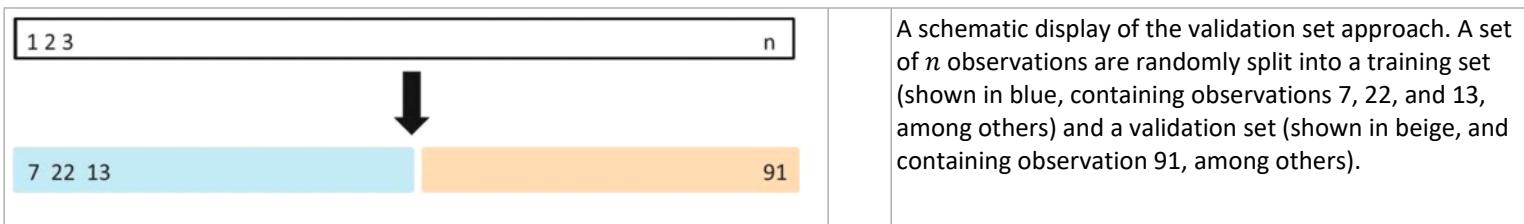
Cross-Validation

These methods are used to do two things

1. *Model Assessment* - the process of evaluating a model's performance.
2. *Model Selection* - the process of selecting the proper level of flexibility for a model.

Since models are 'trained' using training data sets, they are by construction suitable to fit data in these training data sets only (they were fit by minimizing in-sample fitted errors). Since the validation set was not used to fit the model, these set of observations can be used to assess the performance of the model and therefore will allow us to do model selection.

Validation Set Approach



☞ Usually $m/n \approx 0.8$ or 80% and $(n - m)/n \approx 0.2$ or 20%.

Imagine our sample $\{Y_1, Y_2, Y_3, \dots, Y_m, Y_{m+1}, \dots, Y_n\} = \{7, 22, 13, \dots, 91\}$ is such that we also observe a sample of $k \times 1$ vector of realized covariates $\{X_1, X_2, X_3, \dots, X_m, X_{m+1}, \dots, X_n\}$ so that

Training Set:	$\{(Y_i, X_i)\}_{i=1}^m$	Using these m observations, calculate the OLS estimator, $\hat{\beta}$, using this sample.
Validation Set:	$\{(Y_j, X_j)\}_{j=m+1}^n$	Using these $n - m$ observations, calculate predictions $\{X'_{m+1}\hat{\beta}, \dots, X'_n\hat{\beta}\} = \{\hat{Y}_{m+1}, \dots, \hat{Y}_n\}$

Two common out-of-sample goodness-of-fit measures:

Mean Squared Errors:	$MSE = \frac{1}{n-m} \sum_{j=m+1}^n (Y_j - \hat{Y}_j)^2$
R-squared:	$R^2 = 1 - \sum_{j=m+1}^n (Y_j - \hat{Y}_j)^2 / \sum_{j=m+1}^n (Y_j - \bar{Y})^2$

☞ The baseline is the outcome mean of your training data set, i.e., $\bar{Y} = \frac{1}{m} \sum_{i=1}^m Y_i$

Guest 006_resampl... - JupyterLab x +

localhost:8888/doc/tree/HD-Metrics/code/006_resampling.ipynb

006_resampling.ipynb

File Edit View Run Kernel Git Tabs Settings Help

Code git Notebook Python 3 (ipykernel)

Validation Set Approach

```
[3]: %%stata -qui  
splitsample , generate(sample) split(.80 .20) rseed(42)  
label define slabel 1 "Training" 2 "Validation"  
label values sample slabel
```

```
[4]: %stata tabulate sample
```

sample	Freq.	Percent	Cum.
Training	3,990	79.99	79.99
Validation	998	20.01	100.00
Total	4,988	100.00	

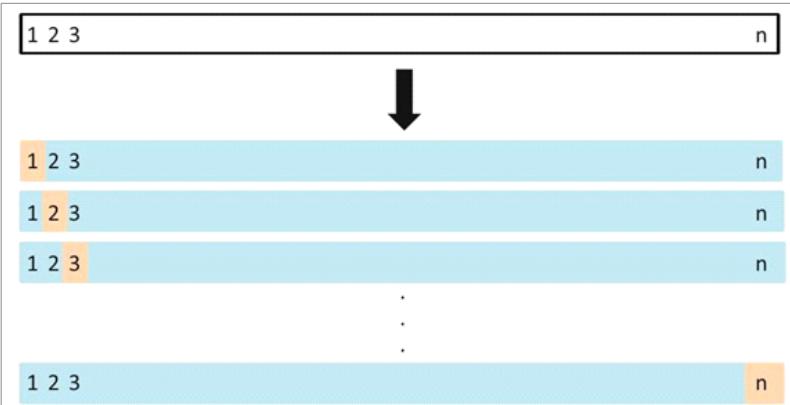
```
[5]: %%stata -qui  
#delimit ;  
qui reg FK6 log_num_authors log_num_pages both_genders prop_women  
`journals' `jel_imp' y_2-y_20 c_2-c_215 jel_flag  
if sample==1;  
#delimit cr  
estimates store ols
```

```
[6]: %stata lassogof ols, over(sample)
```

Penalized coefficients				
Name	sample	MSE	R-squared	Obs
ols				
	Training	.0261068	0.0990	3,990
	Validation	.0254966	0.0304	998

Simple 0 5 main Python 3 (ipykernel) | Idle Mode: Command Ln 1, Col 1 006_resampling.ipynb 0

Leave-One-Out Cross-Validation



- Since (X_1, Y_1) was not used in the fitting process, $MSE_1 = (Y_1 - \hat{Y}_1)^2$ provides an approximately unbiased estimate for the test error.
- We can repeat the procedure by selecting (X_2, Y_2) for the validation data, training the statistical learning procedure on the $n - 1$ observations $\{(X_1, Y_1), (X_3, Y_3), \dots, (X_n, Y_n)\}$, and computing $MSE_2 = (Y_2 - \hat{Y}_2)^2$.

Leave-one-out cross-validation (LOOCV) is closely related to the validation set approach

Like the validation set approach, LOOCV involves splitting the set of observations into two parts.

However, instead of creating two subsets of comparable size, a single observation (X_1, Y_1) is used for the validation set, and the remaining observations $\{(X_2, Y_2), \dots, (X_n, Y_n)\}$ make up the training set.

Repeating this approach n times produces n squared errors, MSE_1, \dots, MSE_n . The LOOCV estimate for the test MSE is the average of these n test error estimates:

$$CV_{(n)} = \frac{1}{n} \sum_{i=1}^n MSE_i.$$

```
%stata -qui
#delimit ;
qui reg FKG log_num_authors log_num_pages both_genders prop_women
`journals' `jel_imp' y_2-y_20 c_2-c_215 jel_flag;
#delimit cr

%stata cv_regress

-----  

Method | Value  

-----  

Root Mean Squared Errors | 0.1692  

Log Mean Squared Errors | -3.5528  

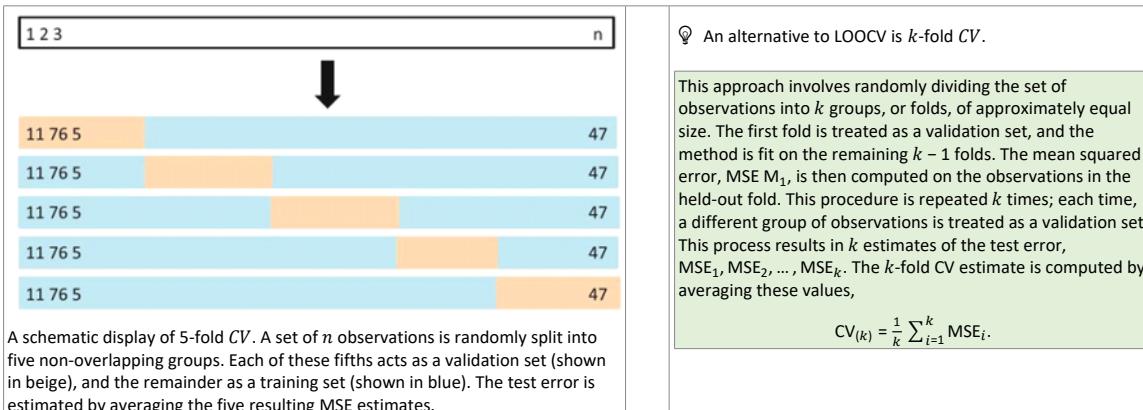
Mean Absolute Errors | 0.1305  

Pseudo-R2 | 0.82028
```

Given the original sample $\{Y_1, \dots, Y_n\}$ and the loocv predictions $\{\hat{Y}_1, \dots, \hat{Y}_n\}$, then

Root Mean Squared Errors =	$\sqrt{n^{-1} \sum_{i=1}^n (Y_i - \hat{Y}_i)^2}$
Mean Absolute Errors =	$n^{-1} \sum_{i=1}^n Y_i - \hat{Y}_i $
Pseudo-R2 =	$\text{corr}(Y_i, \hat{Y}_i)^2$

k-Fold Cross-Validation



```
[9]: %stata
#delimit ;
set seed 42;
crossfold reg FKG log_num_authors log_num_pages both_genders prop_women
    journals `jel_flag' y_2-y_20 c_2-c_215 jel_flag;
k(5) stub(fold);

#delimit cr
```
set seed 42;
crossfold reg FKG log_num_authors log_num_pages both_genders prop_women
 journals `jel_flag' y_2-y_20 c_2-c_215 jel_flag;
k(5) stub(fold) r2;
```
#delimit cr
```
df_r2 = pd.DataFrame(sum(Matrix.get('r(fold)'),[]))

Export to result with datframe format
result = pd.concat([df_rmse, df_mae, df_r2],axis=1)
result.columns = ['RMSE', 'MAE', 'pseudo R2'];
result.index = rows
```
print(result)

      RMSE     MAE  pseudo R2
fold1  0.172111  0.132485  0.898257
fold2  0.169455  0.138320  0.892947
Fold3  0.175192  0.138475  0.814250
fold4  0.172340  0.128885  0.832889
fold5  0.159737  0.124862  0.822878
```
In this case $\sqrt{CV^2}$ equals

[12]: import math as math
import statistics as st
print(math.sqrt(st.mean(result['RMSE']**2)))
0.16985879515275103
```

# $p$ -norms

We will be making extensive use of the 1-norm and 2-norm, so it is useful to review the definition of the general  $p$ -norm.

For a vector  $a = (a_1, \dots, a_k)'$  the  $p$ -norm ( $p \geq 1$ ) is

$$\|a\|_p = \left( \sum_{j=1}^k |a_j|^p \right)^{1/p}.$$

Important special cases include:

| 1-norm                           | 2-norm                                                | sup-norm                                        |
|----------------------------------|-------------------------------------------------------|-------------------------------------------------|
| $\ a\ _1 = \sum_{j=1}^k  a_j $ . | $\ a\ _2 = \left( \sum_{j=1}^k a_j^2 \right)^{1/2}$ . | $\ a\ _\infty = \max_{1 \leq j \leq k}  a_j $ . |

We also define the "**0-norm**"

$$\|a\|_0 = \sum_{j=1}^k \mathbb{1}\{a_j \neq 0\}.$$

↳ The number of non-zero elements.

⚠ This is only heuristically labeled as a "norm."

## Useful results

- ✓ The  $p$ -norm satisfies the following additivity property. If  $a = (a_0, a_1)$  then

$$\|a\|_p^p = \|a_0\|_p^p + \|a_1\|_p^p.$$

- ✓ The [Hölder inequality](#) for  $1/p + 1/q = 1$  is

$$|a'b| \leq \|a\|_p \|b\|_q.$$

- ✓ The case  $p = 1$  and  $q = \infty$  is

$$|a'b| \leq \|a\|_1 \|b\|_\infty.$$

- ✓ The [Minkowski inequality](#) for  $p \geq 1$  is

$$\|a+b\|_p \leq \|a\|_p + \|b\|_p.$$

- ✓ The  $p$ -norms for  $p \geq 1$  satisfy norm monotonicity, i.e.,

$$\|a\|_1 \geq \|a\|_2 \geq \|a\|_\infty.$$

- ✓ Applying Hölder's, we also have the inequality

$$\|a\|_1 = \sum_{j=1}^k |a_j| \mathbb{1}\{a_j \neq 0\} \leq \|a\|_2 \|a\|_0^{1/2}.$$

# Ridge Regression

[Hoerl and Kennard \(1970\)](#) proposed the [Ridge Regression](#) estimator

$$\hat{\beta} = (\mathbf{X}'\mathbf{X} + \lambda \mathbf{I}_p)^{-1} \mathbf{X}'\mathbf{Y}.$$

where  $\lambda > 0$  is called the ridge parameter (tuning parameter).

☞ This estimator has the property that it is well-defined and does not suffer from multicollinearity.

☞ This even holds if  $p > n$  ! the ridge regression estimator is well-defined even when the number of regressors exceeds the sample size.

## Why it works?

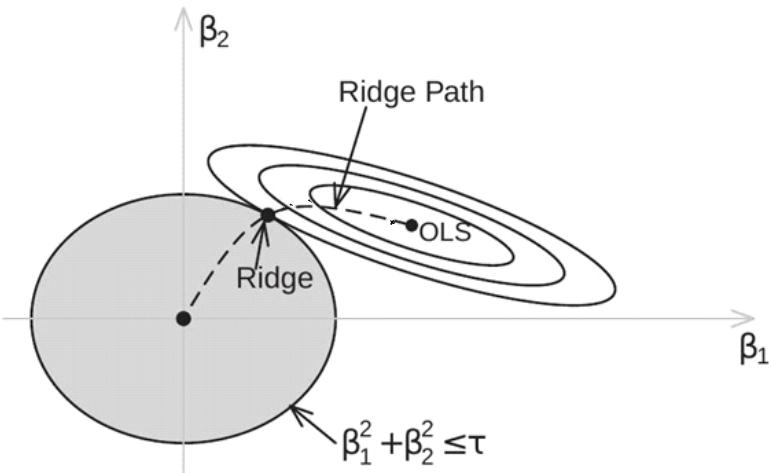
use the spectral decomposition to write  $\mathbf{X}'\mathbf{X} = \mathbf{H}'\mathbf{D}\mathbf{H}$  where  $\mathbf{H}$  is orthonormal and  $\mathbf{D} = \text{diag}\{r_1, \dots, r_p\}$  is a diagonal matrix with the eigenvalues  $r_j$  of  $\mathbf{X}'\mathbf{X}$  on the diagonal. Set  $\Lambda = \lambda \mathbf{I}_p$ . We can write

$$\mathbf{X}'\mathbf{X} + \lambda \mathbf{I}_p = \mathbf{H}'\mathbf{D}\mathbf{H} + \lambda \mathbf{H}'\mathbf{H} = \mathbf{H}'(\mathbf{D} + \Lambda)\mathbf{H},$$

which has strictly positive eigenvalues  $r_j + \lambda > 0$ . Thus all eigenvalues are bounded away from zero so  $\mathbf{X}'\mathbf{X} + \lambda \mathbf{I}_p$  is full rank and well-conditioned.

# Definition

| Definition: <u>Penalized Least Squares</u>                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                     | Definition: <u>Constrained Minimization</u>                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                           |
|--------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------|-------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------|
| <p>Consider the sum of squared errors penalized by the squared 2-norm of the coefficient vector</p> $\text{SSE}_2(b, \lambda) = (\mathbf{Y} - \mathbf{X}b)'(\mathbf{Y} - \mathbf{X}b) + \lambda b'b = \ \mathbf{Y} - \mathbf{X}b\ _2^2 + \lambda \ b\ _2^2.$ <p>The minimizer of <math>\text{SSE}_2(b, \lambda)</math> is a regularized least squares estimator.</p> <p>The first order condition for minimization of <math>\text{SSE}_2(b, \lambda)</math> over <math>b</math> is</p> $-2\mathbf{X}'(\mathbf{Y} - \mathbf{X}b) + 2\lambda b = 0.$ <p>The solution is <math>\hat{\beta}_{\text{ridge}}</math>.</p> <p>You specify the ridge parameter <math>\lambda</math></p> | <p>This problem is</p> $\min_{b'b \leq \tau} (\mathbf{Y} - \mathbf{X}b)'(\mathbf{Y} - \mathbf{X}b).$ <p>for some <math>\tau &gt; 0</math>. To see the connection, the Lagrangian for the constrained problem is</p> $\min_b (\mathbf{Y} - \mathbf{X}b)'(\mathbf{Y} - \mathbf{X}b) + \lambda(b'b - \tau),$ <p>where <math>\lambda</math> is a Lagrange multiplier. The first order condition over <math>b</math> is</p> $-2\mathbf{X}'(\mathbf{Y} - \mathbf{X}b) + 2\lambda b = 0.$ <p>The solution is <math>\hat{\beta}_{\text{ridge}}</math>.</p> <p>You specify the parameter <math>\tau</math></p> |



They are connected because the values of  $\lambda$  and  $\tau$  satisfy the relationship  $\hat{\beta}_{\text{ridge}}'\hat{\beta}_{\text{ridge}} - \tau = 0$

$$\mathbf{Y}'\mathbf{X}(\mathbf{X}'\mathbf{X} + \lambda I_p)^{-1}(\mathbf{X}'\mathbf{X} + \lambda I_p)^{-1}\mathbf{X}'\mathbf{Y} = \tau.$$

To find  $\lambda$  given  $\tau$  it is sufficient to (numerically) solve this equation.

007\_ridge.ipynb · JupyterLab

## 007\_ridge.ipynb

File Edit View Run Kernel Git Tabs Settings Help

Code git Python 3 (ipykernel) ⚙️

### Ridge Regression

The ridge regression estimator

$$\hat{\beta}_{\text{ridge}} = (\mathbf{X}'\mathbf{X} + \lambda\mathbf{I}_p)^{-1}\mathbf{X}'\mathbf{Y}$$

where  $\lambda > 0$  is called the *ridge* parameter.

```
[3]: %%stata -qui -eret steret
#delimit ;
quietly elasticnet linear FKG log_num_authors log_num_pages both_genders prop_women
`journals' `jel_imp' y_2-y_20 c_2-c_215 jel_flag, alpha(0) lambda(5.34) nolog;
#delimit cr
```

```
[4]: import numpy as np
np.set_printoptions(formatter={'float': '{: 0.3f}'.format})
steret['e(b)']
```

```
[4]: array([[-0.001, 0.002, -0.001, -0.003, 0.005, -0.003, -0.001, 0.002,
-0.005, -0.003, 0.004, -0.001, -0.001, -0.001, -0.000, 0.000,
-0.002, -0.003, -0.004, -0.000, 0.000, 0.005, 0.003,
-0.001, 0.004, -0.004, 0.003, 0.000, -0.001, -0.002, -0.002,
-0.001, -0.002, -0.001, -0.001, -0.003, 0.001, 0.003, -0.000,
0.002, -0.001, 0.002, 0.003, -0.001, -0.002, 0.026, 0.052,
0.015, 0.015, 0.002, 0.007, -0.001, 0.000, 0.002, 0.039,
0.050, -0.003, 0.002, 0.034, -0.014, -0.031, -0.025, -0.012,
0.010, -0.014, -0.016, -0.001, 0.024, -0.007, 0.009, 0.002,
0.026, -0.015, -0.005, -0.026, 0.027, -0.022, -0.056, -0.002,
-0.021, 0.020, 0.007, 0.005, 0.037, 0.006, -0.004, -0.015,
0.009, -0.008, -0.010, 0.009, -0.056, -0.004, -0.026, -0.013,
-0.020, -0.001, -0.031, -0.025, -0.003, -0.004, 0.025, -0.006,
0.000, 0.030, 0.024, 0.012, 0.043, 0.005, -0.000, 0.007,
-0.017, -0.008, -0.018, 0.050, -0.024, 0.017, -0.003, 0.016,
-0.021, 0.013, 0.018, -0.013, 0.004, 0.005, -0.011, 0.016,
-0.047, 0.027, 0.021, -0.013, -0.002, 0.009, 0.012, -0.003,
-0.023, 0.016, 0.007, -0.000, 0.039, 0.001, -0.009, 0.021,
-0.022, 0.016, 0.004, 0.004, -0.002, -0.036, -0.009, -0.001,
-0.007, -0.004, -0.011, -0.014, -0.016, 0.030, -0.018, -0.002,
-0.008, 0.024, -0.009, -0.005, -0.014, -0.029, -0.028, -0.014,
-0.004, -0.015, -0.010, -0.013, -0.015, 0.007, -0.020, 0.012,
0.008, -0.040, 0.005, 0.004, -0.011, -0.000, 0.000, 0.035,
0.017, -0.028, 0.026, -0.004, 0.013, -0.040, 0.008, 0.014,
0.000, 0.009, 0.007, -0.018, -0.015, 0.012, 0.010, 0.018,
0.013, -0.010, 0.001, 0.053, 0.036, 0.003, -0.000, -0.013,
-0.010, 0.025, 0.002, -0.005, -0.003, -0.005, -0.018, 0.007,
0.000, 0.028, -0.047, -0.007, -0.021, -0.049, -0.023, 0.046,
0.029, 0.015, 0.019, -0.032, -0.016, -0.030, 0.036, 0.007,
-0.010, 0.014, 0.029, 0.013, -0.020, 0.005, -0.013, 0.026,
0.012, 0.015, -0.003, -0.016, 0.033, -0.011, -0.030, 0.014,
-0.031, 0.032, -0.002, 2.724]])
```

Simple 0 0 3 main Python 3 (ipykernel) | Busy Mode: Command ⚙️ Ln 5, Col 15 007\_ridge.ipynb 0 0

# Cross Validation

The most popular method to select the ridge parameter  $\lambda$  is cross validation. The leave-one-out ridge regression estimator, prediction errors, and CV criterion are

$$\begin{aligned}\hat{\beta}_{-i}(\lambda) &= \left( \sum_{j \neq i} X_j X_j' + \Lambda \right)^{-1} \left( \sum_{j \neq i} X_j Y_j \right) \\ \tilde{e}_i(\lambda) &= Y_i - X_i' \hat{\beta}_{-i}(\lambda) \\ \text{CV}(\lambda) &= \sum_{i=1}^n \tilde{e}_i(\lambda)^2\end{aligned}$$

ⓘ The CV-selected ridge parameter  $\lambda_{cv}$  minimizes  $\text{CV}(\lambda)$ .

ⓘ The cross-validation ridge estimator is calculated using  $\lambda_{cv}$ .

The screenshot shows a Jupyter Notebook interface with a Python 3 (ipykernel) kernel. The code cell contains a Stata command for elastic net linear model selection:

```
[5]: %%stata
#delimit ;
elasticnet linear FKG log_num_authors log_num_pages both_genders prop_women
`journals' `jel_flag' y_2-y_20 c_2-c_215 jel_flag,
alpha(0) lambda(1.3(.1).15) folds(4988) nolog;
#delimit cr
```

The output displays the results of the elastic net linear model selection:

```
Elastic net linear model
No. of obs = 4,988
No. of covariates = 260
Selection: Cross-validation
No. of CV folds = 4,988
```

| alpha | ID  | Description     | lambda | No. of nonzero coef. | Out-of-sample R-squared | CV mean error |
|-------|-----|-----------------|--------|----------------------|-------------------------|---------------|
| 0.000 | 1   | first lambda    | 1.5    | 260                  | 0.0221                  | .0275093      |
|       | * 2 | selected lambda | 1.4    | 260                  | 0.0222                  | .0275089      |
|       | 3   | last lambda     | 1.3    | 260                  | 0.0221                  | .0275106      |

\* alpha and lambda selected by cross-validation.

```
. #delimit cr
delimiter now cr
.
```

At the bottom right, there are three checkmarks with corresponding text:

- ✓ # of folds = n
- ✓ implements the 'leave-one-out' strategy
- ✓ Minimization is always done over a grid of values for  $\lambda$ .

# Statistical Properties

Take the linear regression model

$$\begin{aligned} Y &= X'\beta + e \\ \mathbb{E}[e | X] &= 0 \end{aligned}$$

The bias of the ridge estimator with fixed  $\lambda$  is

$$\text{bias} [\hat{\beta}_{\text{ridge}} | X] = -\lambda(X'X + \lambda I_p)^{-1}\beta.$$

Under random sampling its covariance matrix is

$$\text{var} [\hat{\beta}_{\text{ridge}} | X] = (X'X + \lambda I_p)^{-1}(X'DX)(X'X + \lambda I_p)^{-1}.$$

where  $D = \text{diag}\{\sigma^2(X_1), \dots, \sigma^2(X_n)\}$  and  $\sigma^2(X) = \mathbb{E}[e^2 | X]$ .

We can measure estimation efficiency by the mean squared error (MSE) matrix

$$\text{mse}[\hat{\beta}_{\text{ridge}} | X] = \mathbb{E}[(\hat{\beta} - \beta)(\hat{\beta} - \beta)'] | X.$$

Define  $\underline{\sigma}^2 = \min_{x \in \mathcal{X}} \sigma^2(x)$  where  $\mathcal{X}$  is the support of  $X$ .

## Theorem

In the linear regression model, if  $0 < \lambda < 2\underline{\sigma}^2/\beta'\beta$ ,

$$\text{mse}[\hat{\beta}_{\text{ridge}} | X] < \text{mse}[\hat{\beta} | X].$$

- Theorem shows that the ridge estimator dominates the least squares estimator, if  $\lambda$  satisfies a specific range of values. This holds regardless of the dimension of  $\beta$ .
- Since the upper bound  $2\underline{\sigma}^2/\beta'\beta$  is unknown, however, it is unclear if feasible ridge regression dominates least squares.
- The upper bound does not give practical guidance for selection of  $\lambda$ .

## Variance Estimation

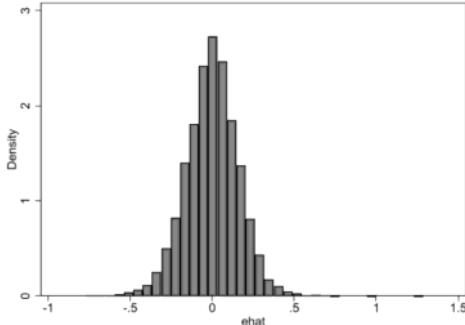
It is straightforward to construct estimators of  $V_{\hat{\beta}} = \text{var}[\hat{\beta} | X]$ . I suggest the [HC3](#)

$$\tilde{V}_{\hat{\beta}} = (X'X + \lambda I_p)^{-1} \left( \sum_{i=1}^n X_i X_i' \hat{e}_i(\lambda)^2 \right) (X'X + \lambda I_p)^{-1}$$

where  $\hat{e}_i(\lambda)$  are the leave-one-out ridge regression prediction errors

$$\hat{e}_i(\lambda) = (1 - X_i'(X'X + \lambda)^{-1}X_i)^{-1} \hat{e}_i(\lambda).$$

where  $\hat{e}_i(\lambda) = Y_i - X_i'\hat{\beta}_{\text{ridge}}(\lambda)$  are the ridge regression residuals.



$\triangle$  The estimators and confidence intervals are valid for the pseudo-true projections, e.g.  $\beta^* = (X'X + \lambda I_p)^{-1}X'\beta$ , not the coefficients  $\beta$  themselves.

```
% stata -qui
#delimit ;
elasticnet linear FKG log_num_authors log_num_pages both_genders prop_women
`journals' `jel_imps' y_2-y_20 c_2-c_215 jel_flag,
alpha(0) lambda(1.4) nolog;
#delimit cr

[7]: x stata matrix list r(table)

r(table)[9,247]
 log_num_authors log_num_pages both_genders prop_women ecm
 b -.0038518 .01776194 .00007089 -.01735091 .02353282
 se .00638722 .00848627 .00789672 .01175006 .00891076
 t -.6030483 2.0930212 .00897685 -.14766652 2.6409455
 pvalue .54650541 .03640028 .99283798 .13983181 .00829462
 ll -.01637372 .00112491 -.01541034 -.0403865 .0060636
 ul .00867011 .03439898 .01555212 .00568467 .041090204
 df 4741 4741 4741 4741 4741
 crit 1.9604645 1.9604645 1.9604645 1.9604645 1.9604645
 eform 0 0 0 0 0

 jpe qje res a_imp b_imp
 b -.01346722 -.00686827 .01397137 -.02455158 -.00913661
 se .00863138 .00885875 .00768578 .03430158 .04662379
 t -1.5602629 -.77530933 1.817821 -.71575652 -.1959645
 pvalue .11876451 .43819556 .06915459 .47417696 .84464635
 ll -.03038872 -.02423553 -.00109633 -.0917986 -.10054088
 ul .00345429 .01049899 .02903906 .04269545 .08226767
 df 4741 4741 4741 4741 4741
 crit 1.9604645 1.9604645 1.9604645 1.9604645 1.9604645
 eform 0 0 0 0 0

 c_imp e_imp f_imp g_imp h_imp
 b .01212098 -.00436706 -.00090114 -.00635589 .00245118
 se .00679658 .00649453 .00833898 .00715904 .00738302
 t 1.7833934 -.67242013 -.10806375 -.88781378 .33200172
 pvalue .07458624 .50134904 .9139497 .37468597 .73990264
 ll -.00120348 -.01709936 -.01724942 -.02039094 -.01202297
 ul .02545454 .00836525 .01544714 .00767915 .01692532
 df 4741 4741 4741 4741 4741
 crit 1.9604645 1.9604645 1.9604645 1.9604645 1.9604645
 eform 0 0 0 0 0

 i_imp j_imp k_imp l_imp m_imp
 b .00836472 -.01248318 -.00977392 -.01175778 .01430392
 se .00832813 .00657588 .0121837 .00702818 .01095725
 t 1.0043936 -1.898329 -.80221277 -.6729465 -.1.30543
 pvalue .3152402 .05771339 .4224701 .09440381 .19189982
 ll -.00796228 -.02537495 -.03635962 -.02553626 -.03578523
 ul .02469172 .00048859 .01411179 .00202073 .00717738
 df 4741 4741 4741 4741 4741
 crit 1.9604645 1.9604645 1.9604645 1.9604645 1.9604645
 eform 0 0 0 0 0
```

# Lasso

[Tibshirani \(1996\)](#) proposed the [Lasso](#) estimator as the minimizer

$$\hat{\beta} = \underset{b}{\operatorname{argmin}} \text{SSE}_1(b, \lambda).$$

where the least squares criterion with a 1-norm penalty is

$$\text{SSE}_1(b, \lambda) = (\mathbf{Y} - \mathbf{X}b)'(\mathbf{Y} - \mathbf{X}b) + \lambda \sum_{j=1}^p |b_j| = \| \mathbf{Y} - \mathbf{X}b \|_2^2 + \lambda \| b \|_1.$$

☞ An important property is that when  $\lambda > 0$  the Lasso estimator is well defined even if  $p > n$ .

# Definition

Definition: Constrained Minimization

The minimization problem

$$\hat{\beta} = \underset{\|\beta\|_1 \leq \tau}{\operatorname{argmin}} \text{SSE}_1(\beta).$$

To see that the two problems are the same observe that the constrained minimization problem has the Lagrangian

$$\min_b (\mathbf{Y} - \mathbf{X}b)'(\mathbf{Y} - \mathbf{X}b) + \lambda \left( \sum_{j=1}^p |b_j| - \tau \right),$$

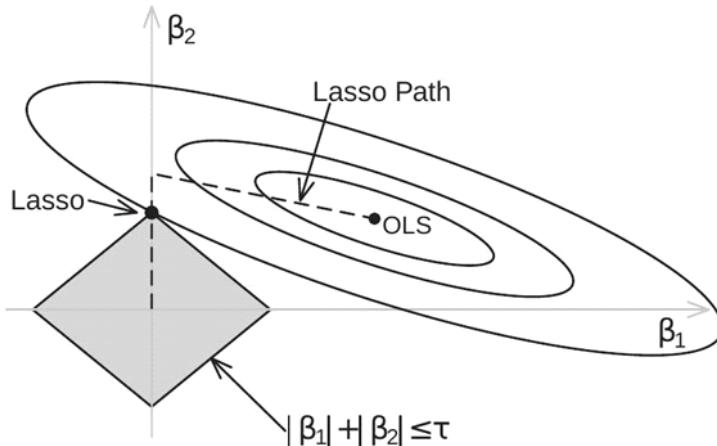
which has 'first order' conditions

$$-2\mathbf{X}'(\mathbf{Y} - \mathbf{X}b) + \lambda \operatorname{sgn}(b_j) = 0.$$

This is the same as those for minimization of the penalized criterion. Thus the solutions are identical.

The sign function of a real number  $x$  is a piecewise function which is defined as follows:

$$\operatorname{sgn} x := \begin{cases} -1 & \text{if } x < 0 \\ 0 & \text{if } x = 0 \\ 1 & \text{if } x > 0 \end{cases}$$



The Lasso path is drawn with the dashed line. This is the sequence of solutions obtained as the constraint set is varied. The solution path has the property that it is a straight line from the least squares estimator to the  $y$ -axis (in this example), at which point  $\beta_2$  is set to zero, and then the solution path follows the  $y$ -axis to the origin. In general, the solution path is linear on segments until a coefficient hits zero, at which point that coefficient is eliminated. In this particular example the solution path shows  $\beta_2$  increasing while  $\beta_1$  decreases. Thus while Lasso is a shrinkage estimator it does not shrink individual coefficients monotonically.

The Lasso is not invariant to the scaling of the regressors. If you rescale a regressor then the penalty has a different meaning. Consequently, it is important to scale the regressors appropriately before applying Lasso. It is conventional to scale all the variables to have mean zero and unit variance.

Lasso is also not invariant to rotations of the regressors. For example, Lasso on  $(\mathbf{X}_1, \mathbf{X}_2)$  is not the same as Lasso on  $(\mathbf{X}_1 - \mathbf{X}_2, \mathbf{X}_2)$  despite having identical least squares solutions. This is troubling as typically there is no default specification.

008\_lasso.ipynb - JupyterLab

## 008\_lasso.ipynb

File Edit View Run Kernel Git Tabs Settings Help

Code git Python 3 (ipykernel)

### Lasso Regression

```
[3]: %%stata -qui -eret steret
#delimit ;
lasso linear FKG log_num_authors log_num_pages both_genders prop_women
`journals' `jel_imp' y_2-y_20 c_2-c_215 jel_flag,
lambda(0.004) nolog;
estimates store lasso;
#delimit cr
```

```
[4]: %stata ereturn display
```

|            | FKG   Coefficient |
|------------|-------------------|
| .0023597s  |                   |
| prop_women | -.0004395         |
| ecm        | .0214475          |
| res        | .0067892          |
| c_imp      | .0100399          |
| j_imp      | -.0002742         |
| l_imp      | -.0040454         |
| m_imp      | -.0058155         |
| p_imp      | .0102725          |
| z_imp      | .0082719          |
| y_2        | -.0020283         |
| y_19       | .0044503          |
| c_2        | -.003098          |
| c_3        | .0054662          |
| c_4        | .1700485          |
| c_12       | .0498898          |
| c_13       | .1115755          |
| c_16       | .0594836          |
| c_17       | -.0158877         |
| c_19       | -.0137163         |
| c_30       | .0670438          |
| c_36       | -.1482595         |
| c_38       | -.0000305         |
| c_42       | .0425364          |
| c_45       | -.0157608         |
| c_50       | -.1565083         |
| c_52       | -.0043919         |
| c_54       | -.0229245         |
| c_60       | .0412154          |

"When we are doing lasso for prediction, we are not supposed to care about the values of the coefficients or look at them"

[Lasso - Stata Manual](#) (page 17)

Simple 0 1 main Python 3 (ipykernel) | Idle Mode: Command Ln 1, Col 29 008\_lasso.ipynb 2

# Group Lasso

Take the model

$$Y = X_1' \beta_1 + \dots + X_G' \beta_G + e,$$

and minimize the SSE  $\equiv \| Y - Xb \|_2^2$  subject to the penalty

$$\lambda_1 \sqrt{b_1' b_1} + \dots + \lambda_G \sqrt{b_G' b_G}.$$

The penalized criterion function is

$$\text{SSE}_3(b, \lambda_1, \dots, \lambda_G) = (Y - Xb)'(Y - Xb) + \sum_{g=1}^G \lambda_g \sqrt{b_g' b_g} = \| Y - Xb \|_2^2 + \sum_{g=1}^G \lambda_g \| b_g \|_2.$$

- One important special case is  $\lambda_1 = 0$ , thus one group of coefficients are not penalized.
- With  $G = 2$  this partitions the coefficients into two groups: penalized and non-penalized.

## Relation with the Lasso

Set  $\lambda_1 = \dots = \lambda_G \equiv \lambda$ , then

$$\text{SSE}_3(b, \lambda) = \| Y - Xb \|_2^2 + \lambda \sum_{g=1}^G \| b_g \|_2.$$

Q When  $G = p$ , one has that the Group Lasso becomes Lasso.

Imagine we 'required' that  $X_1$  must be in the model, but wish to perform model selection among the remaining regressors  $X_2, \dots, X_p$ , then we would minimize

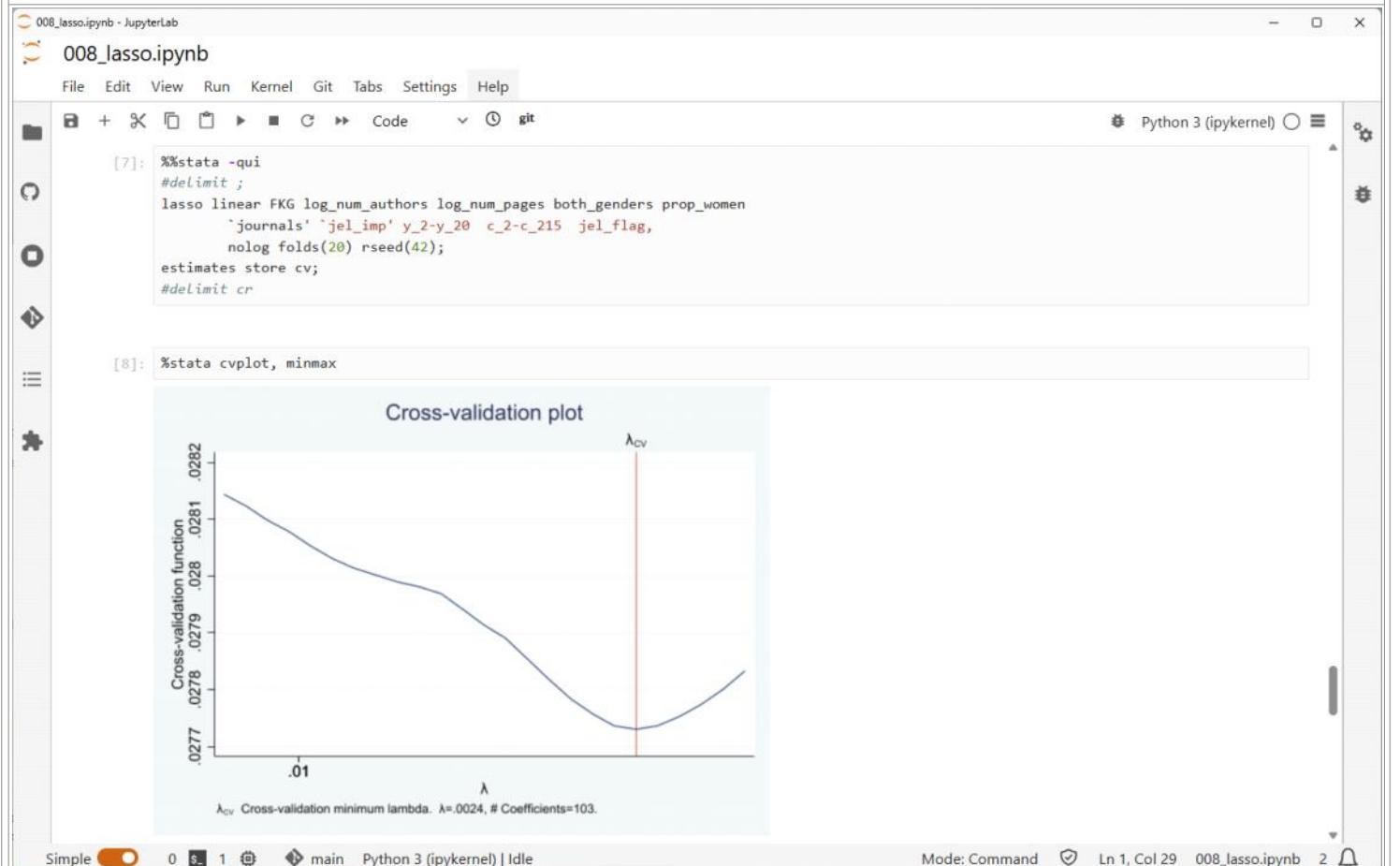
$$(Y - Xb)'(Y - Xb) + \lambda \sum_{j=2}^p |b_j| = \| Y - Xb \|_2^2 + \lambda \sum_{j=2}^p \| b_j \|_1.$$

# Cross Validation

- ⌚ Critically important for Lasso estimation is the penalty  $\lambda$ .
- For  $\lambda$  close to zero the estimates are close to least squares.
  - As  $\lambda$  increases the number of selected variables falls.
  - 🧐 Picking  $\lambda$  induces a trade-off between complexity and parsimony.

Asymptotic consistency of CV selection for Lasso estimation has been demonstrated by [Chetverikov, Liao, and Chernozhukov \(2021\)](#).

## K-fold cross validation



## "1se" Rule

Another popular choice is called the "1se" rule, which is the  $\lambda$  which yields the most parsimonious model for  $\lambda$  values within one standard error of the minimum. The idea is to select a model similar but more parsimonious than the CV-minimizing choice.

# Large Sample Asymptotics

The model is the high-dimensional projection framework:

$$\begin{aligned} Y &= X'\beta + e \\ \mathbb{E}[Xe] &= 0 \end{aligned},$$

where  $X$  is  $p \times 1$  with  $p \gg n$ .

## Sparsity:

The true coefficient vector  $\beta$  is assumed to be sparse in the sense that only a subset of the elements of  $\beta$  are non-zero.

For some  $\lambda$  let  $\hat{\beta}$  be the Lasso estimator which minimizes  $\text{SSE}_1(\beta, \lambda)$ . Define the scaled design matrix  $\mathbf{Q}_n = n^{-1}\mathbf{X}'\mathbf{X}$  and the regression fit



$$(\hat{\beta}_{\text{Lasso}} - \beta)' \mathbf{Q}_n (\hat{\beta}_{\text{Lasso}} - \beta) = \frac{1}{n} \sum_{i=1}^n (X_i' (\hat{\beta}_{\text{Lasso}} - \beta))^2.$$

- ☞ When  $p > n$  the matrix  $\mathbf{Q}_n$  is singular.
- ☞ The theory, however, requires that it not be "too singular".
  - What is required is non-singularity of all sub-matrices of  $\mathbf{Q}_n$  corresponding to the non-zero coefficients.
  - Not "too many" of the zero coefficients.

If  $A$  is a  $m \times n$  matrix, and  $B$  is a  $n \times k$  matrix, then  $\text{rank}(AB) \leq \min(\text{rank}(A), \text{rank}(B))$

Partition  $\beta = (\beta_0, \beta_1)$  where the elements of  $\beta_0$  are all 0 and the elements of  $\beta_1$  are non-zero. (This partition is a theoretical device and unknown to the econometrician.)

Let  $b = (b_0, b_1) \in \mathbb{R}^p$  be partitioned conformably. Define the cone  $B = \{b \in \mathbb{R}^p : \|b_0\|_1 \leq 3\|b_1\|_1\}$ . This is the set of vectors  $b$  such that the sub-vector  $b_0$  is not "too large" relative to the sub-vector  $b_1$ .

## Restricted Eigenvalue Condition (REC)

With probability approaching 1 as  $n \rightarrow \infty$

$$\min_{b \in B} \frac{b' \mathbf{Q}_n b}{b'b} \geq c^2 > 0.$$

☞ Notice that if the minimum is taken without restriction, i.e., over  $\mathbb{R}^p$ , it equals the smallest eigenvalue of  $\mathbf{Q}_n$ . Thus when  $p < n$  a sufficient condition for the REC is  $\lambda_{\min}(\mathbf{Q}_n) \geq c^2 > 0$ . Instead, the minimum is calculated only over  $B$ . In this sense this calculation is similar to a "restricted eigenvalue" which is the source of its name.

**Theorem:** Suppose the high-dimensional projection framework holds with  $p > 1$  and the REC Assumption holds. Assume that each regressor has been standardized so that  $n^{-1}\mathbf{X}'\mathbf{X}_j = 1$  before applying the Lasso. Suppose  $e | X \sim N(0, \sigma^2(X))$  where  $\sigma^2(x) \leq \bar{\sigma}^2 < \infty$ . For some  $C$  sufficiently large set

$$\lambda = C\sqrt{n \log p}.$$

Then there is  $D < \infty$  such that with probability arbitrarily close to 1,

- These rates depend on the number of non-zero coefficients  $\|\beta\|_0$ , the number of variables  $p$ , and the sample size  $n$ .
- Suppose that  $\|\beta\|_0$  is fixed. Then the bounds are  $o(1)$  if  $\log p = o(n)$ . This shows that Lasso estimation is consistent even for an exponentially large number of variables.
- The rates, however, allow the number of non-zero coefficients  $\|\beta\|_0$  to increase with  $n$  at the cost of slowing the allowable rate of increase of  $p$ .

$$\begin{aligned} (\widehat{\beta} - \beta) Q_n (\widehat{\beta} - \beta) &\leq D \|\beta\|_0 \sqrt{\frac{\log p}{n}}, \\ \|\widehat{\beta} - \beta\|_1 &\leq D \|\beta\|_0 \sqrt{\frac{\log p}{n}}, \end{aligned}$$

and

$$\|\widehat{\beta} - \beta\|_2 \leq D \|\beta\|_0^{1/2} \sqrt{\frac{\log p}{n}}.$$

# Postlasso

The Lasso estimator  $\hat{\beta}$  simultaneously selects variables and shrinks coefficients.  $\rightarrow$  Shrinkage introduces bias into estimation.  $\rightarrow$  This bias can be reduced by applying least squares after Lasso selection.  $\leftarrow$  This is known as the Post-Lasso estimator.

The procedure takes two steps:

- (1) Estimate the model  $Y = X'\beta + e$  by Lasso. Let  $X_S$  denote the variables in  $X$  which have non-zero coefficients in  $\hat{\beta}$ . Let  $\beta_S$  denote the corresponding coefficients in  $\hat{\beta}$ .
- (2) the coefficient  $\beta_S$  is estimated by least squares, thus  $\hat{\beta}_S = (X'_S X_S)^{-1} (X'_S Y)$ . This is the Post-Lasso least squares estimator.

The screenshot shows a Jupyter Notebook interface with the title '008\_lasso.ipynb'. The code cell [9] contains R code for estimating a linear model using the lasso package. The output shows the Lasso linear model results, including the number of observations (4,988), covariates (258), and CV folds (100). It also displays the selected lambda values and their corresponding coefficients. The code then moves to Step 1, where it uses the steret package to extract selected variables. Step 2 involves using the b\_postselection function from the steret package to estimate the Post-Lasso least squares coefficients. The bottom status bar indicates the mode is Command, and the file is 008\_lasso.ipynb.

```
[9]: % stata
#delimit ;
lasso linear FKG log_num_authors log_num_pages both_genders prop_women
`journals' `jel_imp' y_2-y_20 c_2-c_215 jel_flag,
nolog sel(cv,serule) folds(100) rseed(42);
#delimit cr

. #delimit ;
delimiter now ;
. lasso linear FKG log_num_authors log_num_pages both_genders prop_women
> `journals' `jel_imp' y_2-y_20 c_2-c_215 jel_flag,
> nolog sel(cv,serule) folds(100) rseed(42);

Lasso linear model
No. of obs = 4,988
No. of covariates = 258
Selection: Cross-validation one s.e. rule No. of CV folds = 100

| No. of Out-of- CV mean
| nonzero sample prediction
ID | Description lambda coef. R-squared error

* 1 | selected lambda .0137899 0 -0.0004 .0281436
2 | lambda after .0125649 1 0.0005 .0281188
25 | last lambda .0014786 158 0.0134 .0277561

* lambda selected by cross-validation one s.e. rule.

. #delimit cr
delimiter now cr
.
Step 1

[10]: steret['e(allvars_sel)']

[10]: 'log_num_pages prop_women ecm res c_imp j_imp l_imp m_imp p_imp z_imp y_2 y_19 c_2 c_3 c_4 c_12 c_13 c_16 c_17 c_19 c_30 c_36 c_38 c_42 c_45 c_50 c_52 c_54 c_60 c_66 c_73 c_75 c_86 c_94 c_98 c_102 c_108 c_139 c_145 c_151 c_165 c_166 c_180 c_183 c_186 c_187 c_190 c_193 c_194 c_198 c_209 c_215'
Step 2

[11]: steret['e(b_postselection)']

[11]: array([[0.01619588, -0.01271079, 0.02942181, 0.01834324, 0.01403871,
-0.00966217, -0.00981967, -0.01423334, 0.03180837, 0.03040684,
-0.01660127, 0.01998842, -0.0132574 , 0.1700719 , 0.33975728,
0.26341346, 0.30595317, 0.22756923, -0.08819735, -0.15223669,
0.13557891, -0.35171181, -0.14866494, 0.24885568, -0.10186321,
-0.35858156, -0.16314212, -0.13988524, 0.15083912, 0.28664969,
0.30666366, 0.09434489, -0.32964873, -0.15622047, 0.23506107,
0.11780536, -0.22804017, -0.26717355, 0.21081936, -0.27031132,
0.30174348, 0.21232284, -0.29287214, -0.29956206, 0.29179373,
0.17902457, -0.2273103 , -0.19565426, 0.1993974 , 0.18060527,
0.18037651, 0.19166047, 2.67286643]])
```

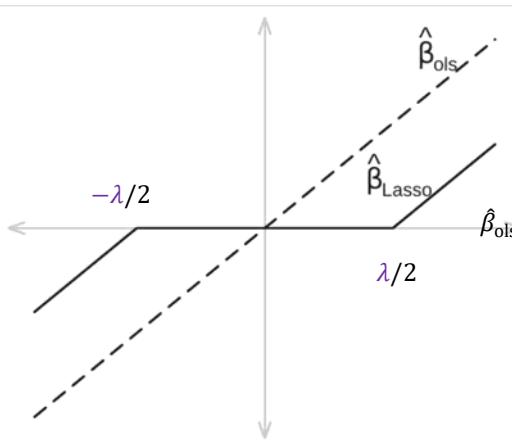
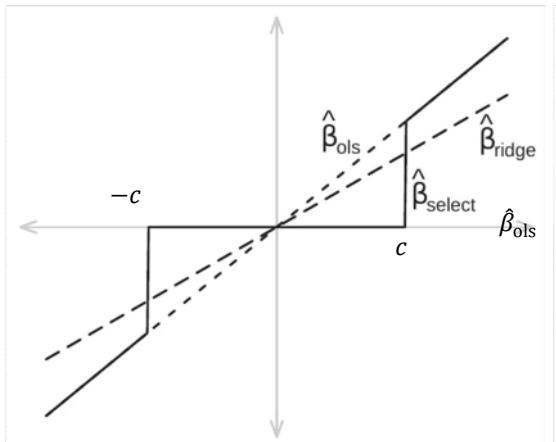
# Lasso | Ridge | Selection

One case where we can explicitly calculate the Lasso estimates is when the regressors are orthogonal, e.g.,  $\mathbf{X}'\mathbf{X} = \mathbf{I}_p$ . Then the first order condition for minimization simplifies to

$$-2(\hat{\beta}_{ols,j} - \hat{\beta}_{Lasso,j}) + \lambda \times \text{sgn}(\hat{\beta}_{Lasso,j}) = 0,$$

which has the explicit solution

| Estimator                                                                                                                                                                                                                               | Remarks                                                                                                                                                                                                                                                                                                                                                                 |
|-----------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------|-------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------|
| $\hat{\beta}_{Lasso,j} = \begin{cases} \hat{\beta}_{ols,j} - \lambda/2 & \hat{\beta}_{ols,j} > \lambda/2 \\ 0 &  \hat{\beta}_{ols,j}  \leq \lambda/2 \\ \hat{\beta}_{ols,j} + \lambda/2 & \hat{\beta}_{ols,j} < -\lambda/2 \end{cases}$ | <ul style="list-style-type: none"> <li>This shows that the Lasso estimate is a continuous transformation of the least squares estimate.</li> <li>For small values of the least squares estimate the Lasso estimate is set to zero.</li> <li>For all other values the Lasso estimate moves the least squares estimate towards zero by <math>\lambda/2</math>.</li> </ul> |
| $\hat{\beta}_{ridge} = (1 + \lambda)^{-1} \hat{\beta}_{ols}$                                                                                                                                                                            | <ul style="list-style-type: none"> <li>It shrinks the OLS coefficient towards zero by a common factor.</li> </ul>                                                                                                                                                                                                                                                       |
| $\hat{\beta}_{select} = \mathbb{1}\{ \hat{\beta}_{ols,j}  > c\} \hat{\beta}_{ols,j}$                                                                                                                                                    |  A selection estimator <ul style="list-style-type: none"> <li>For simplicity consider selection based on a homoskedastic t-test with <math>\sigma^2 = 1</math> and critical value <math>c</math>.</li> </ul>                                                                           |



- The Lasso and ridge estimators are continuous functions while the selection estimator is a discontinuous function.
- The Lasso ("soft") and selection ("hard") estimators are thresholding functions, meaning that the function equals zero for a region about the origin.
- Thresholding estimators are selection estimators because they equal zero when the least squares estimator is sufficiently small.

Hard thresholding rules tend to have high variance due to the discontinuous transformation. Consequently, we expect the Lasso to have reduced variance relative to selection estimators.

# Elastic Net

The difference between Lasso and ridge regression is that the Lasso uses the 1-norm penalty while ridge uses the 2-norm penalty.

💡 Since both procedures have advantages it seems reasonable that further improvements may be obtained by a compromise.

[Zou and Hastie \(2005\)](#) showed that taking a weighted average of the penalties we obtain the [Elastic Net](#) criterion

$$\text{SSE}(b, \lambda, \alpha) = (\mathbf{Y} - \mathbf{X}b)'(\mathbf{Y} - \mathbf{X}b) + \lambda((1 - \alpha) \| b \|_2^2 + \alpha \| b \|_1),$$

with weight  $0 \leq \alpha \leq 1$ .

- ✓ This includes Lasso ( $\alpha = 1$ ) and
- ✓ ridge regression ( $\alpha = 0$ ) as special cases.

Typically the parameters ( $\alpha, \lambda$ ) are selected by joint minimization of the  $k$ -fold cross-validation criterion. Since the elastic net penalty is linear-quadratic the solution is computationally similar to Lasso.

```
% stata
#delimit;
elasticnet linear FKG log_num_authors log_num_pages both_genders prop_women
`journals' `jel_imp' y_2-y_20 c_2-c_215 jel_flag,
alpha(.0001 (.0001) .0005) nolog folds(20) rseed(42);
#delimit cr

. #delimit;
delimiter now ;
. elasticnet linear FKG log_num_authors log_num_pages both_genders prop_women
> `journals' `jel_imp' y_2-y_20 c_2-c_215 jel_flag,
> alpha(.0001 (.0001) .0005) nolog folds(20) rseed(42);

Elastic net linear model
No. of obs = 4,988
No. of covariates = 252
Selection: Cross-validation
No. of CV folds = 20

| No. of Out-of- CV mean
| nonzero sample prediction
alpha ID Description lambda coef. R-squared error
-----+-----
0.001 |
1 | first lambda 137.8993 0 -0.0005 .0281457
61 | last lambda .6674678 235 0.0187 .027607
-----+-----
0.000 |
62 | first lambda 137.8993 0 -0.0005 .0281457
122 | last lambda .6674678 237 0.0184 .0276154
-----+-----
0.000 |
123 | first lambda 137.8993 0 -0.0005 .0281457
182 | last lambda .7325451 238 0.0190 .0275986
-----+-----
0.000 |
183 | first lambda 137.8993 0 -0.0005 .0281457
242 | last lambda .7325451 243 0.0187 .0276068
-----+-----
0.000 |
243 | first lambda 137.8993 0 -0.0005 .0281457
294 | lambda before 1.541937 241 0.0217 .0275225
* 295 | selected lambda 1.404956 242 0.0217 .0275207
296 | lambda after 1.280143 243 0.0217 .0275227
301 | last lambda .8039673 249 0.0192 .0275911
-----+-----
* alpha and lambda selected by cross-validation.
-----+-----
```

Simple 0 3 main Python 3 (ipykernel) | Idle Mode: Command Ln 5, Col 6 009\_elasticnet.ipynb 0

# Prediction

The image shows two side-by-side JupyterLab notebooks, both titled "010\_prediction.ipynb".

**Left Notebook (Stata):**

- Cell 1:** Stata code to generate sample data and split it into training and validation sets.
- Cell 2:** Stata code for OLS regression, including specification of dependent and independent variables, and estimation of coefficients.
- Cell 3:** Stata code for Ridge regression, including specification of dependent and independent variables, and estimation of coefficients.
- Cell 4:** Stata code for Lasso regression, including specification of dependent and independent variables, and estimation of coefficients.
- Cell 5:** Stata code for Elastic Net regression, including specification of dependent and independent variables, and estimation of coefficients.

**Right Notebook (Python):**

- Cell 1:** Python code to run OLS, Ridge, Lasso, and elasticnet regressions using the statsmodels library.
- Cell 2:** Output table showing Penalized coefficients for each model across Training and Validation sets. The table includes columns for Name, sample, MSE, R-squared, and Obs.
- Text:** A note about Postselection coefficients, stating they should not be used with elasticnet and, in particular, with ridge regression. Ridge works by shrinking the coefficient estimates, and these are the estimates that should be used for prediction. Because postselection coefficients are OLS regression coefficients for the selected coefficients and because ridge always selects all variables, postselection coefficients after ridge are OLS regression coefficients for all potential variables, which clearly we do not want to use for prediction.

# Inference

Consider the linear model

$$Y = D\theta + X'\beta + e$$
$$\mathbb{E}[e | D, X] = 0$$

The goal is inference on  $\theta$ .

where  $Y$  and  $D$  are scalar and  $X$  is  $p \times 1$ .

- ☞ The variable  $D$  is the main focus of the regression;
- ☞ the variable  $X$  are controls.

Suppose you estimate model by group post-Lasso, only penalizing  $\beta$ . This performs selection on the variables  $X$ , resulting in a least squares regression of  $Y$  on  $D$  and the selected variables in  $X$ .

- ⌚ The coverage probabilities for  $\theta$  are downward biased, and the distortions are serious.
- ⌚ The distortions are primarily affected by (and increasing in) the correlation between  $D$  and  $X$ .

[Belloni, Chernozhukov, and Hansen \(2014b\)](#) deduce that improved coverage accuracy can be achieved if the variable  $X$  is included in the regression whenever  $X$  and  $D$  are correlated.

# Double Selection Lasso

[Belloni, Chernozhukov, and Hansen \(2014b\)](#) proposed to perform what they call **double-selection**. We start by specifying an auxiliary equation for  $D$  :

$$\frac{D}{\mathbb{E}[V | X]} = X' \textcolor{blue}{y} + \textcolor{blue}{V}$$

$$= 0.$$

Substituting this into our structural model, we obtain a reduced form for  $Y$  :

$$\frac{Y}{\mathbb{E}[U | X]} = X' \textcolor{blue}{\eta} + \textcolor{blue}{U}$$

$$= 0.$$

The proposed double-selection algorithm applies model selection (e.g., Lasso selection) separately to each equations, takes the union of the selected regressors, and then estimates structural equation by least squares using the selected regressors. This method ensures that a variable  $X$  is included if it is relevant for the structural regression or if it is correlated with  $D$ .

The **double-selection** estimator as recommended by [Belloni, Chernozhukov, and Hansen \(2014b\)](#) is:

- 1 Estimate  $\textcolor{blue}{y}$  by Lasso. Let  $X_1$  be the selected variables from  $X$ .
- 2 Estimate  $\textcolor{blue}{\eta}$  by Lasso. Let  $X_2$  be the selected variables from  $X$ .
- 3 Let  $\tilde{X} = X_1 \cup X_2$  be the union of the variables in  $X_1$  and  $X_2$ .
- 4 Regress  $Y$  on  $(D, \tilde{X})$  by OLS to obtain the double-selection coefficient estimate  $\hat{\theta}_{DS}$ .
- 5 Calculate a conventional (heteroskedastic) standard error for  $\hat{\theta}_{DS}$ .

$$\begin{aligned}\textcolor{blue}{\eta} &= \beta + \gamma \theta, \\ \textcolor{brown}{U} &= e + \textcolor{blue}{V} \theta\end{aligned}$$

Provided the models in steps 1 and 2 satisfy an approximate sparsity structure (so that the regressions are well approximated by a finite set of regressors) then the double-selection estimator  $\hat{\theta}_{DS}$  and its  $t$ -ratio are asymptotically normal so conventional inference methods are valid.

The essential idea is that because  $\tilde{X}$  includes the variables in  $X_2$ , the estimator  $\hat{\theta}_{DS}$  is asymptotically equivalent to the regression where  $D$  is replaced with the error  $V$ . Since  $V$  is uncorrelated with the regressors  $X$  the estimator and  $t$ -ratio satisfy the conventional non-selection asymptotic distribution.

011\_dsregress... - JupyterLab  
011\_dsregress.ipynb

File Edit View Run Kernel Git Tabs Settings Help

Code Notebook Python 3 (ipykernel)

```
[1]: %%capture
import stata_setup, os
if os.name == 'nt':
 stata_setup.config('C:/Program Files/Stata17/','mp')
else:
 stata_setup.config('/usr/local/stata17','mp')
```

## Preparing the data

```
[2]: %%stata -qui

use "../data/data", clear
rename log_flesch_kincaid_grade_level FKG
quietly tabulate year, generate(y_)
quietly tabulate cluster, generate(c_)

local journals ecm jpe qje res //AER based category

local jel_imp a_imp b_imp c_imp e_imp f_imp g_imp h_imp i_imp j_imp k_imp ///
l_imp m_imp n_imp o_imp p_imp q_imp r_imp y_imp z_imp // D JEL based case
```

## Double Selection Lasso

The Double Selection Lasso is implemented by the `dsregress` Stata command. It'll utilize a 'plugin' value for the  $\lambda$ s parameters by default, but they can also use standard cross-validation instead.

```
*[3]: %%stata -qui
#delimit ;
dsregress FKG log_num_authors log_num_pages both_genders prop_women,
 controls('journals' `jel_imp' y_2-y_20 c_2-c_215 jel_flag)
 vce(cluster cluster) rseed(42);
estimates store ds_plugin;
dsregress FKG log_num_authors log_num_pages both_genders prop_women,
 controls('journals' `jel_imp' y_2-y_20 c_2-c_215 jel_flag)
 vce(cluster cluster) selection(cv) rseed(42);
estimates store ds_cv;
#delimit cr
```

Simple 0 \$ 8 main Python 3 (ipykernel) | Idle Mode: Command Ln 4, Col 62 011\_dsregress.ipynb 0

011\_dsregress... - JupyterLab  
011\_dsregress.ipynb

File Edit View Run Kernel Git Tabs Settings Help

Code Notebook Python 3 (ipykernel)

```
[4]: %%stata
#delimit ;
lassocoef (ds_plugin, for(FKG)) (ds_cv, for(FKG))
(ds_plugin, for(prop_women)) (ds_cv, for(prop_women));
#delimit cr

. #delimit ;
delimiter now ;
. lassocoef (ds_plugin, for(FKG)) (ds_cv, for(FKG))
> (ds_plugin, for(prop_women)) (ds_cv, for(prop_women));

-----+-----+-----+-----+
| ds_plugin ds_cv ds_plugin ds_cv
| FKG FKG prop_women prop_women
-----+-----+-----+-----+
c_4 | x
c_36 | x
c_50 | x
c_165 | x
c_183 | x
c_51 | x
c_100 | x
c_108 | x
c_174 | x
c_196 | x
c_imp | x
i_imp | x
k_imp | x
l_imp | x
o_imp | x
z_imp | x
y_10 | x
_cons | x x x x
```

Legend:

- b - base level
- e - empty cell
- o - omitted
- x - estimated

```
. #delimit cr
delimiter now cr
```

The first two columns of x's show which controls were selected from the lassos for the dependent variable, `log_flesch_kincaid_grade_level`—the first column for the plugin method and the second for cross-validation. The third and fourth columns of x's show which controls were selected by the lassos for the covariate of interest, `prop_women`.

011\_dsregress... - JupyterLab  
011\_dsregress.ipynb

File Edit View Run Kernel Git Tabs Settings Help

Notebook Python 3 (ipykernel)

```
[5]: %%stata
estimates restore ds_cv
dsregress

. estimates restore ds_cv
(results ds_cv are active now)

. dsregress

Double-selection linear model Number of obs = 4,988
 Number of controls = 257
 Number of selected controls = 37
 Wald chi2(4) = 13.61
 Prob > chi2 = 0.0087

(Std. err. adjusted for 215 clusters in cluster)
-----| Robust
FKG | Coefficient std. err. z P>|z| [95% conf. interval]
-----+-----log_num_au-s | -.0053689 .0042331 -1.27 0.205 -.0136657 .0029278
log_num_pa-s | .0161836 .0052931 3.86 0.002 .0058093 .026558
both_genders | .0085261 .0042799 0.12 0.902 -.0078624 .0089146
prop_women | -.0172381 .0058155 -2.96 0.003 -.0286363 -.0058398

Note: Chi-squared test is Wald test of the coefficients of the variables
of interest jointly equal to zero. Lassos select controls for model
estimation. Type lassoinfo to see number of selected variables in each
lasso.

Note: Lassos are performed accounting for clusters in cluster.
```

Simple 0 8 main Python 3 (ipykernel) | Idle Mode: Command Ln 4, Col 38 011\_dsregress.ipynb 0

💡 We interpret the coefficient estimates just as we would for a standard linear regression.

⚠️ What we lose with the inferential lasso estimators is the ability to interpret any other coefficients.

💡 Our point estimate for the effect of the proportion of women co-authors on abstract readability is -0.017, meaning that going from no women co-authors (prop\_women=0) to all women co-authors (prop\_women=1) will decrease the number of schooling needed to read the abstract by 1.724%. This value is statistically different from 0 well beyond the 5% level, in fact, beyond the 0.1% level. Our 95% confidence interval is -0.0286 to -0.0059.

# Post-Regularization Lasso

## ☒ [Partialing-Out] Lasso

[Chernozhukov, Hansen, and Spindler \(2015\)](#) proposed to transform the structural equation to eliminate the *high-dimensional* component. Take the expected value of  $Y = D\theta + e$  conditionally on  $X$ , and subtract from each side. This leads to the equation

$$Y - \mathbb{E}[Y | X] = (D - \mathbb{E}[D | X])\theta + e.$$

💡 Notice that this eliminates the regressor  $X$  and the high-dimensional coefficient  $\beta$ !

The authors specify  $\mathbb{E}[Y | X]$  and  $\mathbb{E}[D | X]$  as linear functions of  $X$ . Substituting these expressions we obtain

$$Y - X'\eta = (D - X'y)\theta + e.$$

☒ If  $\eta$  and  $y$  were known the coefficient  $\theta$  could be estimated by least squares. However, as  $\eta$  and  $y$  are unknown they need to be estimated. [Chernozhukov, Hansen, and Spindler \(2015\)](#) recommend estimation by Lasso or post-Lasso, separately for  $Y$  and  $D$  as follows:

1. Estimate a regression of  $Y$  on  $X$  by Lasso or post-Lasso (Stata) with Lasso parameter  $\lambda_1$ . Let  $\hat{\eta}$  be the coefficient estimator and  $\hat{U}_i = Y_i - X'_i\hat{\eta}$  the residual.
2. Estimate a regression of  $D$  on  $X$  by Lasso or post-Lasso (Stata) with Lasso parameter  $\lambda_2$ . Let  $\hat{y}$  be the coefficient estimator and  $\hat{V}_i = D_i - X'_i\hat{y}$  the residual.
3. Let  $\hat{\theta}_{PR}$  be the OLS coefficient from the regression of  $\hat{U}$  on  $\hat{V}$ .
4. Calculate a conventional (heteroskedastic) standard error for  $\hat{\theta}_{PR}$ .

### Theorem

Suppose

$$\begin{aligned} D &= X'y + V && \text{Define the scaled design matrix} \\ \mathbb{E}[V | X] &= 0. && Q_n = n^{-1}X'X, \text{ with probability} \\ Y &= X'\eta + U && \text{approaching 1 as } n \rightarrow \infty \\ \mathbb{E}[U | X] &= 0. && \min_{b \in B} \frac{b'Q_nb}{b'b} \geq c^2 > 0. \end{aligned}$$

Assume that each regressor has been standardized so that  $n^{-1}\mathbf{X}'_j\mathbf{X}_j = 1$ . Suppose  $e | X \sim N(0, \sigma_e^2(X))$  and  $V | X \sim N(0, \sigma_V^2(X))$  where  $\sigma_e^2(x) \leq \bar{\sigma}_e^2 < \infty$  and  $\sigma_V^2(x) \leq \bar{\sigma}_V^2 < \infty$ . For some  $C_1$  and  $C_2$  sufficiently large the Lasso parameters satisfy  $\lambda_1 = C_1\sqrt{n \log p}$  and  $\lambda_2 = C_2\sqrt{n \log p}$ . Assume  $p \rightarrow \infty$  and

$$(\|\beta\|_0 + \|\gamma\|_0) \frac{\log p}{\sqrt{n}} = o(1).$$

Then

$$\sqrt{n}(\hat{\theta}_{PR} - \theta) \xrightarrow{d} N\left(0, \frac{\mathbb{E}[V^2e^2]}{(\mathbb{E}[V^2])^2}\right).$$

- The original result does not need the assumption of normal errors.
- Theorem shows that the post-regularization (partialing-out) Lasso estimator has a conventional asymptotic distribution, allowing conventional inference for the coefficient  $\theta$ , i.e., the standard variance estimator for  $\hat{\theta}_{PR}$  is consistent for the asymptotic variance.

☒ The advantage of the post-regularization estimator  $\hat{\theta}_{PR}$  over the double-selection estimator  $\hat{\theta}_{DS}$  is efficiency.

☒ The post-regularization estimator uses only the relevant components of  $X$  to separately demean  $Y$  and  $D$ , leading to greater parsimony. Different components of  $X$  may be relevant to  $D$  and  $Y$ .

☒ The double-selection estimator uses the union of the two regressor sets for estimation of  $\theta$ , leading to a less parsimonious specification.



As a consequence, an advantage of the double-selection estimator is reduced bias and robustness.

012\_poregres... - JupyterLab

File Edit View Run Kernel Git Tabs Settings Help

012\_poregres.ipynb X +

Notebook Python 3 (ipykernel)

```
[2]: %%stata -qui
use "../data/data", clear
gen prop_women_100 = prop_women*100
rename log_flesch_kincaid_grade_level FKG
quietly tabulate year, generate(y_)
quietly tabulate cluster, generate(c_)

local journals ecm jpe qje res //AER based category

local jel_imp a_imp b_imp c_imp e_imp f_imp g_imp h_imp i_imp j_imp k_imp ///
l_imp m_imp n_imp o_imp p_imp q_imp r_imp y_imp z_imp // D JEL based case
```

We multiply the variable of interest `prop_women` by 100 so a 'unit' change is a 1% change.

```
[3]: %%stata
#delimit ;
poregress FKG log_num_authors log_num_pages both_genders prop_women_100,
 controls(`journals' `jel_imp' y_2-y_20 c_2-c_215 jel_flag)
 vce(cluster cluster) rseed(42);
estimates store po_plugin;
poregress FKG log_num_authors log_num_pages both_genders prop_women_100,
 controls(`journals' `jel_imp' y_2-y_20 c_2-c_215 jel_flag)
 vce(cluster cluster) selection(cv) rseed(42);
estimates store po_cv;
#delimit cr
```

```
[4]: %%stata
#delimit ;
lassocoef (po_plugin, for(FKG)) (po_cv, for(FKG))
 (po_plugin, for(prop_women_100)) (po_cv, for(prop_women_100));
#delimit cr
```

012\_poregres... - JupyterLab

File Edit View Run Kernel Git Tabs Settings Help

012\_poregres.ipynb X +

Notebook Python 3 (ipykernel)

```
. #delimit ;
delimiter now ;
. lassocoef (po_plugin, for(FKG)) (po_cv, for(FKG))
> (po_plugin, for(prop_women_100)) (po_cv, for(prop_women_100));

-----+----| po_plugin po_cv po_plugin po_cv
 | FKG FKG prop_women_100 prop_women_100
-----+----| x x x x
x_4 | x
c_36 | x
c_50 | x
c_165 | x
c_183 | x
c_56 | x
c_91 | x
c_100 | x
c_174 | x
c_196 | x
c_imp | x
i_imp | x
k_imp | x
l_imp | x
o_imp | x
z_imp | x
y_10 | x
_cons | x x x x
```

Legend:

- b - base level
- e - empty cell
- o - omitted
- x - estimated

```
. #delimit cr
delimiter now cr
.
```

```

%%stata
estimates restore po_cv
poregress

. estimates restore po_cv
(results po_cv are active now)

. poregress

Partialing-out linear model Number of obs = 4,988
 Number of controls = 257
 Number of selected controls = 37
 Wald chi2(4) = 11.16
 Prob > chi2 = 0.0248

(Std. err. adjusted for 215 clusters in cluster)
-----| Robust
 | Coefficient std. err. z P>|z| [95% conf. interval]
-----| log_num_aus-s | -.0029956 .0051732 -0.58 0.563 -.0131349 .0071438
log_num_pa-s | .0161292 .0055794 2.89 0.004 .0051939 .0270646
both_genders | -.0021614 .0040552 -0.53 0.594 -.0101093 .0057866
prop_wom-100 | -.0001582 .0000725 -2.18 0.029 -.0003004 -.000016

Note: Chi-squared test is a Wald test of the coefficients of the variables
of interest jointly equal to zero. Lassos select controls for model
estimation. Type lassoinfo to see number of selected variables in each
lasso.

Note: Lassos are performed accounting for clusters in cluster.

```

✓ These use post-lasso estimation to construct the residuals  $\hat{U}$  and  $\hat{V}$ .

💡 Since the dependent variable, `log_flesch_kincaid_grade_level` is measured in logarithm, the point estimate `-0.0001582` is the elasticity of the Flesch-kincaid Grade index with respect to the proportion of women co-authors. This is statistically significant at 10% & 5%.

# Double/Debiased Machine Learning

[Chernozhukov, Chetverikov, Demirer, Duflo, Hansen, Newey, and Robins \(2018\)](#) proposed the Double/Debiased machine learning (DML) estimator.

As presented previously, the post-regularization estimator first estimates the coefficients  $\gamma$  and  $\eta$  in the models

$$\begin{aligned} \frac{D}{\mathbb{E}[V | X]} &= X' \gamma + V \\ &= 0. \\ \frac{Y}{\mathbb{E}[U | X]} &= X' \eta + U \\ &= 0. \end{aligned}$$

and then estimates the coefficient  $\theta$ .

The DML estimator performs these estimation steps using separate samples based on K-fold partitioning. The estimation algorithm is as follows:

- 1 Randomly partition the sample into  $K$  independent folds  $A_k, k = 1, \dots, K$ , of roughly equal size  $n/K$ .
- 2 Write the data matrices for each fold as  $(Y_k, D_k, X_k)$ .
- 3 For  $k = 1, \dots, K$ 
  - (a) Use all observations except for fold  $k$  to estimate the coefficients  $\gamma$  and  $\eta$  by Lasso or post-Lasso (Stata). Write these leave-fold-out estimators as  $\hat{\gamma}_{-k}$  and  $\hat{\eta}_{-k}$ .
  - (b) Set  $\hat{V}_k = D_k - X_k \hat{\gamma}_{-k}$  and  $\hat{U}_k = Y_k - X_k \hat{\eta}_{-k}$ . These are the estimated values of  $V$  and  $U$  for observations in the  $k^{\text{th}}$  fold using the leave-fold-out estimators.
- 4 Set  $\hat{\theta}_{\text{DML}} = (\sum_{k=1}^K \hat{V}_k \hat{V}_k')^{-1} (\sum_{k=1}^K \hat{V}_k \hat{U}_k)$ . Equivalently, stack  $\hat{V}_k$  and  $\hat{U}_k$  into  $n \times 1$  vectors  $\hat{V}$  and  $\hat{U}$  and set  $\hat{\theta}_{\text{DML}} = (\hat{V}' \hat{V})^{-1} (\hat{V}' \hat{U})$ .
- 5 Construct a conventional (heteroskedastic) standard error for  $\hat{\theta}_{\text{DML}}$ .

☞ The authors label this estimator the "DML2" estimator.

☞ An alternative they label "DML1" is  $\hat{\theta}_{\text{DML1}} = K^{-1} \sum_{k=1}^K (\hat{V}_k \hat{V}_k')^{-1} (\hat{V}_k \hat{U}_k)$ .

💡 Although  $\hat{\theta}_{\text{DML1}}$  is asymptotically equivalent to  $\hat{\theta}_{\text{DML}}$ , the latter is preferred.

## Advantage:

- The advantage of the DML estimator over the post-regularization estimator is that the sample splitting eliminates the dependence between the two estimation steps, thereby reducing post-model-selection bias.

## Theorem

Under the assumptions of the previous Theorem,

$$\sqrt{n}(\hat{\theta}_{\text{DML}} - \theta) \xrightarrow{d} N\left(0, \frac{\mathbb{E}[V^2 e^2]}{(\mathbb{E}[V^2])^2}\right).$$

Furthermore, the standard variance estimator for  $\hat{\theta}_{\text{DML}}$  is consistent for the asymptotic variance.

## Disadvantages:

- First, the estimator is random due to the sample splitting. Two researchers with the same data set but making different random splits will obtain two distinct estimators. This randomness can be reduced by using a larger value of  $K$ , but this increases computation cost.
- Another disadvantage of sample-splitting is that estimation of  $\gamma$  and  $\eta$  is performed using smaller samples which reduces estimation efficiency, though this effect is minor if  $K \geq 10$ . Regardless, these considerations suggest that DML may be most appropriate for settings with large  $n$  and  $K \geq 10$ .

**Preparing the data**

```
[2]: %%stata -qui
use "../data/data", clear
rename log_flesch_kincaid_grade_level FKG
quietly tabulate year, generate(y_)
quietly tabulate cluster, generate(c_)

local journals ecm jpe qje res //AER based category
local jel_imp a_imp b_imp c_imp e_imp f_imp g_imp h_imp i_imp j_imp k_imp ///
l_imp m_imp n_imp o_imp p_imp q_imp r_imp y_imp z_imp // D JEL based case
```

**Post-Regularization Lasso**

The Double/Debiased Machine Learning estimator is implemented by the `xporegress` Stata command. It'll utilize a 'plugin' value for the `A`s parameters by default, but they can also use standard cross-validation instead.

```
[3]: %%stata -qui
#delimit ;
xporegress FKG log_num_authors log_num_pages both_genders prop_women,
 controls('journals' `jel_imp' y_2-y_20 c_2-c_215 jel_flag)
 vce(cluster cluster) rseed(42);
estimates store xpo_plugin;
xporegress FKG log_num_authors log_num_pages both_genders prop_women,
 controls('journals' `jel_imp' y_2-y_20 c_2-c_215 jel_flag)
 vce(cluster cluster) selection(cv) rseed(42);
estimates store xpo_cv;
#delimit cr
```

☞ `xporegress` fits 10 lassos for the dependent variable and 10 more lassos for each covariate of interest! That is the default; you can request more. Or you can request fewer, but that is not recommended.

☞ So, `xporegress` is orders of magnitude slower than `poregress` and `dsregress`. And it has orders of magnitude more lassos to explore.

```
[4]: %%stata
estimates restore xpo_cv
xporegress

. estimates restore xpo_cv
(results xpo_cv are active now)

. xporegress

Cross-fit partialing-out Number of obs = 4,988
linear model Number of controls = 257
 Number of selected controls = 41
 Number of folds in cross-fit = 10
 Number of resamples = 1
 Wald chi2(4) = 7.80
 Prob > chi2 = 0.0993

 (Std. err. adjusted for 215 clusters in cluster)

-----| Robust
FKG | Coefficient std. err. z P>|z| [95% conf. interval]
-----+-----log_num_aus | -.0026534 .0068324 -0.39 0.698 -.0160447 .010738
log_num_pas | .0165392 .0100478 1.65 0.100 -.0031541 .0362326
both_genders | -.001836 .0062027 -0.38 0.767 -.013993 .0103211
prop_women | -.016443 .0084284 -1.95 0.051 -.0329623 .0000763

Note: Chi-squared test is a Wald test of the coefficients of the variables
of interest jointly equal to zero. Lassos select controls for model
estimation. Type lassoinfo to see number of selected variables in each
lasso.
Note: Lassos are performed accounting for clusters in cluster.
```

☞ Our point estimate for the effect of the proportion of women co-authors on abstract readability is  $-0.0164$ , meaning that going from no women co-authors (`prop_women=0`) to all women co-authors (`prop_women=1`) will decrease the number of schooling needed to read the abstract by  $1.644\%$ .

# What to Use?

⌚ It is recommended you should start with the Double/Debiased Machine Learning estimator (xporegress) first.

## Why?

⌚ It is safer if you think that the process that generated your data has lots of covariates relative to your sample size.

⌚ Similarly, it is also safer if you want to explore lots of potential controls.

- The number of potential controls is not as problematic as the number of true covariates because it is the natural log of the potential control that counts. For example, needing 10 additional true covariates is the same as requesting just over 22,000 potential controls! The jargon term for this is *sparsity*.
- xporegress has a weaker sparsity requirement than do Post-Regularization Lasso (poregress) and Double Selection Lasso (dsregress).

⌚ However ...

If your model is weakly identified by the data (dropping one observation or regressor would change estimates dramatically, et c.), then the Double Selection Lasso (dsregress) can be more stable than the others because it uses a union of all the selected controls from all the lassos for all of its computations after selection. ⚡ ... more stable doesn't mean better !

# Lasso | Logit & Poisson

Recall that the [Lasso](#) estimator for the *linear* model is defined as the minimizer

$$\hat{\beta} = \underset{b}{\operatorname{argmin}} \text{SSE}_1(b, \lambda).$$

where the least squares criterion with a 1-norm penalty is

$$\text{SSE}_1(b, \lambda) = (\mathbf{Y} - \mathbf{X}b)'(\mathbf{Y} - \mathbf{X}b) + \lambda \sum_{j=1}^p |b_j| = \|\mathbf{Y} - \mathbf{X}b\|_2^2 + \lambda \|b\|_1.$$

Equivalently, we could write

$$\hat{\beta} = \underset{b}{\operatorname{argmin}} Q_n(b) + \lambda \|b\|_1.$$

where  $Q_n(b) \equiv \sum_{i=1}^n (Y_i - X_i'b)^2$  when the model is *linear*.

| Model   | $Q_n(b)$ .                                                                         |
|---------|------------------------------------------------------------------------------------|
| Linear  | $\sum_{i=1}^n (Y_i - X_i'b)^2$                                                     |
| Logit   | $\sum_{i=1}^n -Y_i(X_i'b) + \log\{1 + \exp(X_i'b)\}$                               |
| Probit  | $\sum_{i=1}^n -Y_i \times \log\{\Phi(X_i'b)\} - (1 - Y_i) \log\{1 - \Phi(X_i'b)\}$ |
| Poisson | $\sum_{i=1}^n -Y_i(X_i'b) + \exp\{X_i'b\}$                                         |

014\_lasso\_co... - JupyterLab

File Edit View Run Kernel Tabs Settings Help

014\_lasso\_cont.ipynb +

Code git Notebook Python 3 (ipykernel)

## Lasso - Logistic

```
[3]: %stata -qui -eret steret
#delimit ;
lasso logit FKG_01 prop_women $baseline $mfe, lambda(0.018) nolog;
#delimit cr
```

```
[4]: %stata ereturn display
```

|                     | FKG_01   Coefficient |
|---------------------|----------------------|
| - .0062653m         |                      |
| 0.c_4   -1.185463   |                      |
| 0.c_13   -.0127349  |                      |
| 0.c_42   -.0158715  |                      |
| 0.c_98   -.1500967  |                      |
| 0.c_165   -.0096062 |                      |
| 0.c_186   -.0158715 |                      |
| _cons   -.803255    |                      |

## Lasso - Poisson

```
[5]: %stata -qui -eret steret
#delimit ;
lasso poisson excess_sentences prop_women $baseline $mfe, lambda(.07) nolog;
#delimit cr
```

```
[6]: %stata ereturn display
```

|                      | excess_sen~s   Coefficient |
|----------------------|----------------------------|
| - log_num_pa~s       | .3637272                   |
| 0.ecm   -.2927387    |                            |
| 0.qje   -.174802     |                            |
| 0.res   -.3358273    |                            |
| 0.c_imp   -.0362706  |                            |
| 0.c_86   -.235316    |                            |
| 0.c_114   -.4203244  |                            |
| 0.c_136   -.1.154179 |                            |
| 0.c_142   -.100394   |                            |
| 0.c_174   -.1172538  |                            |

# Elastic Net | Logit & Poisson

Recall that the [Elastic Net](#) estimator for the *linear* model is defined as the minimizer of

$$\text{SSE}(b, \lambda, \alpha) = (\mathbf{Y} - \mathbf{X}b)'(\mathbf{Y} - \mathbf{X}b) + \lambda((1 - \alpha) \| b \|_2^2 + \alpha \| b \|_1),$$

with weight  $0 \leq \alpha \leq 1$ .

- ✓ This includes Lasso ( $\alpha = 1$ ) and
- ✓ ridge regression ( $\alpha = 0$ ) as special cases.

Equivalently, we could write

$$\hat{\beta} = \underset{b}{\operatorname{argmin}} Q_n(b) + \lambda((1 - \alpha) \| b \|_2^2 + \alpha \| b \|_1).$$

where  $Q_n(b) \equiv \sum_{i=1}^n (Y_i - X_i'b)^2$  when the model is *linear*.

| Model   | $Q_n(b)$ .                                                                         |
|---------|------------------------------------------------------------------------------------|
| Linear  | $\sum_{i=1}^n (Y_i - X_i'b)^2$                                                     |
| Logit   | $\sum_{i=1}^n -Y_i(X_i'b) + \log\{1 + \exp(X_i'b)\}$                               |
| Probit  | $\sum_{i=1}^n -Y_i \times \log\{\Phi(X_i'b)\} - (1 - Y_i) \log\{1 - \Phi(X_i'b)\}$ |
| Poisson | $\sum_{i=1}^n -Y_i(X_i'b) + \exp\{X_i'b\}$                                         |

015\_elasticnet... - JupyterLab

File Edit View Run Kernel Git Tabs Settings Help

015\_elasticnet\_cont.ipynb X +

Notebook Python 3 (ipykernel)

## Elastic Net - Logistic

```
[3]: %stata -qui -eret steret
#delimit ;
elasticnet logit FKG_01 prop_women $baseline $mfe, alpha(0.99) lambda(0.02) nolog;
#delimit cr
```

```
[4]: %stata ereturn display
```

|       | FKG_01   Coefficient |
|-------|----------------------|
| c_4   | -.3811678            |
| 1     | .3843454             |
| _cons | -1.816232            |

## Elastic Net - Poisson

```
[5]: %stata -qui -eret steret
#delimit ;
elasticnet poisson excess_sentences prop_women $baseline $mfe, alpha(0.7) lambda(.1) nolog;
#delimit cr
```

```
[6]: %stata ereturn display
```

|              | excess_sentences   Coefficient |
|--------------|--------------------------------|
| log_num_pas~ | .3588161                       |
| ecm          |                                |
| 0            | -.1451147                      |
| 1            | .1451147                       |
| qje          |                                |
| 0            | -.0868885                      |
| 1            | .0868885                       |
| res          |                                |
| 0            | -.1660911                      |
| 1            | .1660911                       |

# Binary Response

When  $Y$  is **binary**, meaning that it takes two values. Without loss of generality these are taken as zero and one, thus  $Y$  has support  $\{0,1\}$ . In econometrics we typically call this class of models **binary choice**.

## Examples:

- ✓ Purchase of a single item
- ✓ Market entry
- ✓ Approval of an application
- ✓ Participation in a program

☞ The dependent variable may be recorded as Yes/No, True/False, or 1/-1, but can always be written as 1/0.

## Binary Choice Models

Let  $(Y, X')$  be random with  $Y \in \{0,1\}$  and  $X \in \mathbb{R}^k$ . The response probability of  $Y$  with respect to  $X$  is

$$P(x) = \Pr[Y = 1 | X = x] = \mathbb{E}[Y | X = x].$$

The response probability completely describes the conditional distribution of  $Y|X$ . The marginal effect ( $X$  continuous) is

$$\frac{\partial}{\partial x} P(x) = \frac{\partial}{\partial x} \Pr[Y = 1 | X = x] = \frac{\partial}{\partial x} \mathbb{E}[Y | X = x].$$

This is the regression derivative.

## ⚠ Linear Regression ⚠

The variables satisfy the regression framework

$$\begin{aligned} Y &= P(X) + e \\ \mathbb{E}[e | X] &= 0 \end{aligned}$$

The error  $e$  is not "classical." It has the two-point conditional distribution

$$e = \begin{cases} 1 - P(X), & \text{with probability } P(X) \\ -P(X), & \text{with probability } 1 - P(X). \end{cases}$$

It is also highly heteroskedastic with conditional variance

$$\text{var}[e | X] = P(X)(1 - P(X)).$$

004\_logit.ipynb - JupyterLab

004\_logit.ipynb

File Edit View Run Kernel Git Tabs Settings Help

Code git

Notebook Python 3 (ipykernel)

[2]: %%stata -qui  
use "../data/data", clear  
rename log\_flesch\_kincaid\_grade\_level FKG  
quietly tabulate year, generate(y\_)  
quietly tabulate cluster, generate(c\_)  
local journals ecm jpe qje res //AER based category  
local jel\_imp a\_imp b\_imp c\_imp e\_imp f\_imp g\_imp h\_imp i\_imp j\_imp k\_imp ///  
l\_imp m\_imp n\_imp o\_imp p\_imp q\_imp r\_imp y\_imp z\_imp // D JEL based case

## Logistic Regression

Let `FKG_01` be a binary variable indicating that an article's abstract readability is among the 10% highest in the Flesch Kincaid grade level scale.

[3]: %%stata -qui  
summarize flesch\_kincaid\_grade\_level, detail  
local cutoff = r(p90)  
gen int FKG\_01 = (flesch\_kincaid\_grade\_level >= `cutoff')

[4]: %stata list FKG\_01 log\_num\_authors log\_num\_pages both\_genders prop\_women jel\_flag in 1/10, clean

|     | FKG_01 | log_n~rs | log_~ges | both_g~s | prop_w~n | jel_flag |
|-----|--------|----------|----------|----------|----------|----------|
| 1.  | 0      | .6931472 | 2.639057 | 0        | 0        | 1        |
| 2.  | 1      | .6931472 | 2.70805  | 0        | 0        | 1        |
| 3.  | 0      | 1.098612 | 3.258096 | 0        | 0        | 1        |
| 4.  | 1      | 0        | 3.091043 | 0        | 0        | 1        |
| 5.  | 1      | 0        | 3.496508 | 0        | 0        | 1        |
| 6.  | 0      | 0        | 2.772589 | 0        | 0        | 1        |
| 7.  | 0      | 1.098612 | 2.890372 | 0        | 0        | 1        |
| 8.  | 1      | .6931472 | 3.295837 | 0        | 0        | 1        |
| 9.  | 0      | .6931472 | 2.833213 | 0        | 0        | 1        |
| 10. | 0      | .6931472 | 3.091043 | 0        | 0        | 1        |

Simple 0 0 1 main Python 3 (ipykernel) | Idle Mode: Command Ln 1, Col 1 004\_logit.ipynb 0

# Models

**Linear Probability Model:**  $P(x) = x'\beta$  where  $\beta \in \mathbb{R}^k$  is a coefficient vector.

Advantages:

- The coefficients  $\beta$  equal the marginal effects (when  $X$  does not include nonlinear transformations)
- Linearity means that estimation is simple as least squares can be used to estimate the coefficients.

Disadvantage:

- It does not respect the [0,1] boundary. Fitted and predicted values from estimated linear probability models frequently violate these boundaries producing nonsense results.

**Index Model:**  $P(x) = G(x'\beta)$  where  $G(u)$  is a *link function* and  $\beta$  is a coefficient vector.

☞ This framework is also called a **single index** model where  $x'\beta$  is a **linear index** function.

☞ In binary choice models,  $G(u)$  is a **distribution function** which respects the probability bounds  $0 \leq G(u) \leq 1$ .

☞  $G(\cdot)$  is *symmetric*, i.e.,  $G(-u) = 1 - G(u)$ .

☞  $g(u) = \frac{\partial}{\partial u} G(u)$  denote the density function of  $G(u)$ .

**Probit Model:**  $P(x) = \Phi(x'\beta)$  where  $\Phi(u)$  is the standard normal distribution function.

**Logit Model:**

$P(x) = \Lambda(x'\beta)$  where  $\Lambda(u) = (1 + \exp(-u))^{-1}$  is the logistic distribution function.

# Latent Variable Interpretation

An index model can be interpreted as a latent variable model. Consider

$$Y^* = X'\beta + e$$

$$e \sim G(e)$$

$$Y = \mathbb{1}\{Y^* > 0\} = \begin{cases} 1 & \text{if } Y^* > 0 \\ 0 & \text{otherwise} \end{cases}$$

The event  $Y = 1$  is the same as  $Y^* > 0$ , which is the same as

$$X'\beta + e > 0.$$

This means that the response probability is

$$P(X) = \Pr[e > -X'\beta | X] = 1 - G(-X'\beta) = G(X'\beta).$$

You may have noticed that we have discussed cases where the error  $e$  is either *standard normal* or *standard logistic*, that is, their scale is fixed.



This is because the scale of the error distribution is not identified.



To see this, suppose that  $e = \sigma \varepsilon$  where  $\varepsilon$  has a distribution  $G^*(\varepsilon)$  with unit variance. Then the response probability is

$$\Pr[Y = 1 | X = x] = \Pr[\sigma \varepsilon > x'\beta | X = x] = G^*\left(\frac{x'\beta}{\sigma}\right) = G^*(x'\beta^*),$$

where  $\beta^* = \beta/\sigma$ . This is an index model with coefficient  $\beta^*$ .



This means that  $\beta$  and  $\sigma$  are not separately identified.

While the coefficient  $\beta$  is not identified the following parameters are identified:

1. Scaled coefficients:  $\beta^* = \beta/\sigma$ .
  2. Ratios of coefficients:  $\beta_1/\beta_2 = \beta_1^*/\beta_2^*$ .
  3. Marginal effects:  $\frac{\partial}{\partial x} P(x) = \frac{\beta}{\sigma} g\left(\frac{x'\beta}{\sigma}\right) = \beta^* g(x'\beta^*)$ .
- These only depend on  $\beta^*$  so are identified.

```

004_logit.ipynb - JupyterLab
004_logit.ipynb
File Edit View Run Kernel Git Tabs Settings Help
Notebook Python 3 (ipykernel)
Variable Preparation
We now utilize Stata capabilities to automatically identify continuous as well as dummy variables in the data set.

[5]: % stata
#delimit ;
vi set log_num_authors log_num_pages both_genders prop_women
"journals" "jel_imp" y_2-y_20 c_2-c_215 jel_flag
, dummy clear nonotes;
delimit cr

. #delimit ;
delimiter now ;
. vi set log_num_authors log_num_pages both_genders prop_women
> "journals" "jel_imp" y_2-y_20 c_2-c_215 jel_flag
> , dummy clear nonotes;

Macro contents
Macro | # Vars Description
System
$vldummy | 258 0/1 variables
$vlcategorical | 0 categorical variables
$vlcontinuous | 3 continuous variables
$vluncertain | 0 perhaps continuous, perhaps categorical variables
$vlother | 0 all missing or constant variables

. delimit cr

```

```

004_logit.ipynb - JupyterLab
004_logit.ipynb
File Edit View Run Kernel Git Tabs Settings Help
Notebook Python 3 (ipykernel)
[6]: %%stata
#delimit ;
vi create fe = vldummy ~ (both_genders jel_flag);
vi substitute mfe = i.fe;
vi create controls = vlcategorical ~ (prop_women);
vi create controls_dummy = (both_genders jel_flag);
vi substitute baseline = i.controls_dummy controls;
vi rebuild;
#delimit cr

. display "$baseline"
i.both_genders i.jel_flag log_num_pages log_num_authors

. display "$fe"
i.eom i.joe i.eje i.res i.a.imp i.b.imp i.c.imp i.f.imp i.g.imp i.h.imp
> i.i.imp i.j.imp i.k.imp i.l.imp i.n.imp i.o.imp i.p.imp i.q.imp i.r.
> i.y_11 i.y_12 i.y_13 i.y_14 i.y_15 i.y_16 i.y_17 i.y_18 i.y_19 i.y_20 i.c_2 i.c
> _3 i.c_4 i.c_5 i.c_6 i.c_7 i.c_8 i.c_9 i.c_10 i.c_11 i.c_12 i.c_13 i.c_14 i.c
> _15 i.c_16 i.c_17 i.c_18 i.c_19 i.c_20 i.c_21 i.c_22 i.c_23 i.c_24 i.c_25 i.c
> _26 i.c_27 i.c_28 i.c_29 i.c_30 i.c_31 i.c_32 i.c_33 i.c_34 i.c_35 i.c_36 i.c
> _37 i.c_38 i.c_39 i.c_40 i.c_41 i.c_42 i.c_43 i.c_44 i.c_45 i.c_46 i.c_47 i.c
> _48 i.c_49 i.c_50 i.c_51 i.c_52 i.c_53 i.c_54 i.c_55 i.c_56 i.c_57 i.c_58 i.c_59 i.c
> _60 i.c_61 i.c_62 i.c_63 i.c_64 i.c_65 i.c_66 i.c_67 i.c_68 i.c_69 i.c_70 i.c_71 i.c_72 i.c_73 i.c_74 i.c_75 i.c_76 i.c_77 i.c_78 i.c_79 i.c_80 i.c
> _81 i.c_82 i.c_83 i.c_84 i.c_85 i.c_86 i.c_87 i.c_88 i.c_89 i.c_90 i.c_91 i.c
> _92 i.c_93 i.c_94 i.c_95 i.c_96 i.c_97 i.c_98 i.c_99 i.c_100 i.c_101 i.c_102
> i.c_103 i.c_104 i.c_105 i.c_106 i.c_107 i.c_108 i.c_109 i.c_110 i.c_111 i.c_1
> _12 i.c_113 i.c_114 i.c_115 i.c_116 i.c_117 i.c_118 i.c_119 i.c_120 i.c_121 i.
> _c_122 i.c_123 i.c_124 i.c_125 i.c_126 i.c_127 i.c_128 i.c_129 i.c_130 i.c_131
> i.c_132 i.c_133 i.c_134 i.c_135 i.c_136 i.c_137 i.c_138 i.c_139 i.c_140 i.c_
> _141 i.c_142 i.c_143 i.c_144 i.c_145 i.c_146 i.c_147 i.c_148 i.c_149 i.c_150 i
> _c_151 i.c_152 i.c_153 i.c_154 i.c_155 i.c_156 i.c_157 i.c_158 i.c_159 i.c_160
> i.c_161 i.c_162 i.c_163 i.c_164 i.c_165 i.c_166 i.c_167 i.c_168 i.c_169 i.c_170 i.c
> _171 i.c_172 i.c_173 i.c_174 i.c_175 i.c_176 i.c_177 i.c_178 i.c_179
> i.c_180 i.c_181 i.c_182 i.c_183 i.c_184 i.c_185 i.c_186 i.c_187 i.c_188 i.c_1
> _189 i.c_190 i.c_191 i.c_192 i.c_193 i.c_194 i.c_195 i.c_196 i.c_197 i.c_198 i.c
> _199 i.c_200 i.c_201 i.c_202 i.c_203 i.c_204 i.c_205 i.c_206 i.c_207 i.c_208
```

# Estimation - Likelihood

Recall that if  $Y$  is Bernoulli, such that  $\Pr[Y = 1] = p$  and  $\Pr[Y = 0] = 1 - p$ , then  $Y$  has the probability mass function

$$\pi(y) = p^y(1-p)^{1-y}, \quad y \in \{0,1\}.$$

In the index model  $\Pr[Y = 1 | \mathbf{X}] = G(\mathbf{X}'\beta)$ ,  $Y$  is **conditionally** Bernoulli, so its **conditional** probability mass function is

$$\pi(Y | \mathbf{X}) = G(\mathbf{X}'\beta)^Y (1 - G(\mathbf{X}'\beta))^{1-Y} = G(\mathbf{X}'\beta)^Y G(-\mathbf{X}'\beta)^{1-Y} = G(Z'\beta),$$

where

$$Z = \begin{cases} \mathbf{X} & ; Y = 1 \\ -\mathbf{X} & ; Y = 0 \end{cases}$$

Taking natural logarithm and summing across observations we obtain:

|                          |                                                                                                                                                                              | Probit                                                                                                                                                                                           | Logit                                                                        |
|--------------------------|------------------------------------------------------------------------------------------------------------------------------------------------------------------------------|--------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------|------------------------------------------------------------------------------|
| Population:              | $\ell(\beta) = \mathbb{E}[\log G(Z'\beta)]$                                                                                                                                  | $\ell^{\text{probit}}(\beta) = \mathbb{E}[\log \Phi(Z'\beta)]$                                                                                                                                   | $\ell^{\text{logit}}(\beta) = \mathbb{E}[\log \Lambda(Z'\beta)]$             |
| log-likelihood function: | $\ell_n(\beta) = n^{-1} \sum_{i=1}^n \log G(Z_i'\beta)$                                                                                                                      | $\ell_n^{\text{probit}}(\beta) = n^{-1} \sum_{i=1}^n \log \Phi(Z_i'\beta)$                                                                                                                       | $\ell_n^{\text{logit}}(\beta) = n^{-1} \sum_{i=1}^n \log \Lambda(Z_i'\beta)$ |
| First derivative:        | $h(x) = \frac{d}{dx} \log G(x)$ $S_n(\beta) = \frac{\partial}{\partial \beta} \ell_n(\beta)$ $= n^{-1} \sum_{i=1}^n Z_i h(Z_i'\beta)$                                        | $h_{\text{probit}}(x) = \frac{\phi(x)}{\Phi(x)} \stackrel{\text{def}}{=} \lambda(x)$ <br>Inverse Mills Ratio | $h_{\text{logit}}(x) = 1 - \Lambda(x)$                                       |
| Second derivative:       | $H(x) = -\frac{d^2}{dx^2} \log G(x)$ $\mathcal{H}_n(\beta) = -\frac{\partial^2}{\partial \beta \partial \beta'} \ell_n(\beta)$ $= n^{-1} \sum_{i=1}^n X_i X_i' H(Z_i'\beta)$ | $H_{\text{probit}}(x) = \lambda(x)(x + \lambda(x))$                                                                                                                                              | $H_{\text{logit}}(x) = \Lambda(x)(1 - \Lambda(x))$                           |

The MLE is the value which maximizes  $\ell_n(\beta)$ . We write this as

$$\begin{aligned} \hat{\beta}^{\text{probit}} &= \underset{\beta}{\operatorname{argmax}} \ell_n^{\text{probit}}(\beta) \\ \hat{\beta}^{\text{logit}} &= \underset{\beta}{\operatorname{argmax}} \ell_n^{\text{logit}}(\beta) \end{aligned}$$

The *pseudo* true values are defined as

$$\begin{aligned} \beta^{\text{probit}} &= \underset{\beta}{\operatorname{argmax}} \ell^{\text{probit}}(\beta) \\ \beta^{\text{logit}} &= \underset{\beta}{\operatorname{argmax}} \ell^{\text{logit}}(\beta) \end{aligned}$$

# Estimation - Asymptotic Distribution

Let  $B$  be the parameter space for  $\beta$

| Probit                                                                                                                                                                                                                                                                                                                                                                                                         | Logit                                                                                                                                                                                                                                                                                                                                                                                                   |
|----------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------|---------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------|
| <b>Theorem (Consistency):</b><br>If (1) $(Y_i, X_i)$ are i.i.d.; (2) $\mathbb{E} \ X\ ^2 < \infty$ ; (3) $Q_{\text{probit}} > 0$ ; and (4) $B$ is <u>compact</u> ; then $\hat{\beta}_{\text{probit}} \xrightarrow[p]{} \beta_{\text{probit}}$ as $n \rightarrow \infty$ .                                                                                                                                      | <b>Theorem (Consistency):</b><br>If (1) $(Y_i, X_i)$ are i.i.d.; (2) $\mathbb{E} \ X\  < \infty$ ; (3) $Q_{\text{logit}} > 0$ ; and (4) $B$ is <u>compact</u> ; then $\hat{\beta}_{\text{logit}} \xrightarrow[p]{} \beta_{\text{logit}}$ as $n \rightarrow \infty$ .                                                                                                                                    |
| $Q_{\text{probit}} \stackrel{\text{def}}{=} \mathbb{E} [XX' H_{\text{probit}}(X'\beta_{\text{probit}})] > 0,$<br>$\Omega_{\text{probit}} = \mathbb{E} [XX' \Lambda(X'\beta_{\text{probit}})^2],$                                                                                                                                                                                                               | $Q_{\text{logit}} \stackrel{\text{def}}{=} \mathbb{E} [XX' \Lambda(X'\beta_{\text{logit}})(1 - \Lambda(X'\beta_{\text{logit}}))] > 0,$<br>$\Omega_{\text{logit}} = \mathbb{E} [XX' (Y - \Lambda(X'\beta_{\text{logit}}))^2].$                                                                                                                                                                           |
| <b>Theorem (Asy. Normality):</b><br>If the conditions of consistency hold plus $\mathbb{E} \ X\ ^4 < \infty$ and $\beta_{\text{probit}}$ is in the <u>interior</u> of $B$ ; then as $n \rightarrow \infty$<br>$\sqrt{n}(\hat{\beta}_{\text{probit}} - \beta_{\text{probit}}) \xrightarrow[d]{} N(0, V),$<br>where $V_{\text{probit}} = Q_{\text{probit}}^{-1} \Omega_{\text{probit}} Q_{\text{probit}}^{-1}$ . | <b>Theorem (Asy. Normality):</b><br>If the conditions of consistency hold plus $\mathbb{E} \ X\ ^4 < \infty$ and $\beta_{\text{logit}}$ is in the <u>interior</u> of $B$ ; then as $n \rightarrow \infty$<br>$\sqrt{n}(\hat{\beta}_{\text{logit}} - \beta_{\text{logit}}) \xrightarrow[d]{} N(0, V),$<br>where $V_{\text{logit}} = Q_{\text{logit}}^{-1} \Omega_{\text{logit}} Q_{\text{logit}}^{-1}$ . |

## Standard Errors

The asymptotic standard errors of  $\hat{\beta}_{\text{probit}}$  and  $\hat{\beta}_{\text{logit}}$  are defined as the squared-root of the main diagonal element of  $n^{-1}\hat{V}_{\text{probit}}$  and  $n^{-1}\hat{V}_{\text{logit}}$ , where

|                                                                                                                                                                                                                     |                                                                                                                                                                           |
|---------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------|---------------------------------------------------------------------------------------------------------------------------------------------------------------------------|
| $\hat{V}_{\text{probit}} = \hat{Q}_{\text{probit}}^{-1} \hat{\Omega}_{\text{probit}} \hat{Q}_{\text{probit}}^{-1}$                                                                                                  | $\hat{V}_{\text{logit}} = \hat{Q}_{\text{logit}}^{-1} \hat{\Omega}_{\text{logit}} \hat{Q}_{\text{logit}}^{-1}$                                                            |
| $\hat{Q}_{\text{probit}} = \frac{1}{n} \sum_{i=1}^n X_i X_i' \hat{\lambda}_i (\hat{\mu}_i + \hat{\lambda}_i)$ , where $\hat{\mu}_i = X_i' \hat{\beta}_{\text{probit}}$ , $\hat{\lambda}_i = \lambda(\hat{\mu}_i)$ , | $\hat{Q}_{\text{logit}} = \frac{1}{n} \sum_{i=1}^n X_i X_i' \hat{\Lambda}_i (1 - \hat{\Lambda}_i)$ , where $\hat{\Lambda}_i = \Lambda(X_i' \hat{\beta}_{\text{logit}})$ , |
| $\hat{\Omega}_{\text{probit}} = \frac{1}{n} \sum_{i=1}^n X_i X_i' \hat{\lambda}_i^2$ .                                                                                                                              | $\hat{\Omega}_{\text{logit}} = \frac{1}{n} \sum_{i=1}^n X_i X_i' (Y_i - \hat{\Lambda}_i)^2$ .                                                                             |

004\_logit.ipynb - JupyterLab

004\_logit.ipynb

File Edit View Run Kernel Git Tabs Settings Help

Code git Notebook Python 3 (ipykernel)

```
[8]: %%stata -qui
#delimit ;
logit FKG_01 prop_women $baseline $mfe, vce(cluster cluster) nolog;
#delimit cr
```

```
[9]: %%stata
estimates table, keep(log_num_authors log_num_pages 1.both_genders prop_women) b(%5.4f) star varwidth(25)
```

| Variable        | Active     |
|-----------------|------------|
| log_num_authors | -0.2194*** |
| log_num_pages   | 0.4182***  |
| 1.both_genders  | 0.2597*    |
| prop_women      | -0.5166**  |

Legend: \* p<0.05; \*\* p<0.01; \*\*\* p<0.001

```
[10]: %stata estat class
```

Logistic model for FKG\_01

| True       |     |      | Total |
|------------|-----|------|-------|
| Classified | D   | ~D   |       |
| +          | 25  | 16   | 41    |
| -          | 463 | 4045 | 4508  |
| Total      | 488 | 4061 | 4549  |

Classified + if predicted Pr(D) >= .5  
 True D defined as FKG\_01 != 0

|                               | Pr( +  D)  | 5.12%  |
|-------------------------------|------------|--------|
| Sensitivity                   | Pr( -  ~D) | 99.61% |
| Specificity                   | Pr( D  +)  | 60.98% |
| Positive predictive value     | Pr( ~D  -) | 89.73% |
| Negative predictive value     | Pr( +  ~D) | 0.39%  |
| False + rate for true ~D      | Pr( -  D)  | 94.88% |
| False - rate for true D       | Pr( ~D  +) | 39.02% |
| False + rate for classified + | Pr( D  -)  | 10.27% |
| Correctly classified          |            | 89.47% |

Simple 0 s 1 main Python 3 (ipykernel) | Idle

Mode: Command Ln 1, Col 1 004\_logit.ipynb 0

# Marginal Effects

Recall that  $\Pr[Y = 1 | X = x] = G(x'\beta) = G(\beta_1x_1 + \dots + \beta_jx_j + \dots + \beta_{k-1}x_{k-1} + \beta_k)$

|                          | $x_j$ is continuous                                                     | $x_j$ is binary                                                                                                                                                                 |
|--------------------------|-------------------------------------------------------------------------|---------------------------------------------------------------------------------------------------------------------------------------------------------------------------------|
| Marginal Effects:        | $\delta_j(x) = \beta_j g(x'\beta)$                                      | $\delta_j(x) = G(\beta_1x_1 + \dots + \beta_j \cdot 1 + \dots + \beta_{k-1}x_{k-1} + \beta_k) - G(\beta_1x_1 + \dots + \beta_j \cdot 0 + \dots + \beta_{k-1}x_{k-1} + \beta_k)$ |
| Average Marginal Effects | $\text{AME} = \mathbb{E}[\delta_j(X)] = \beta_j \mathbb{E}[g(X'\beta)]$ | $\text{AME} = \mathbb{E}[\delta_j(X)]$                                                                                                                                          |

Estimator:  $\widehat{\text{AME}} = \frac{1}{n} \sum_{i=1}^n \delta_j(X_i)$ , where  $\delta_j(X_i) = \hat{\beta}_j g(X_i'\hat{\beta})$  or  $\delta_j(X_i) = G(\hat{\beta}_1 X_{1,i} + \dots + \hat{\beta}_j \cdot 1 + \dots + \hat{\beta}_{k-1} X_{k-1,i} + \hat{\beta}_k) - G(\hat{\beta}_1 X_{1,i} + \dots + \hat{\beta}_j \cdot 0 + \dots + \hat{\beta}_{k-1} X_{k-1,i} + \hat{\beta}_k)$ .

```

004_logit.ipynb - JupyterLab
004_logit.ipynb
File Edit View Run Kernel Git Tabs Settings Help
Marginal Effects

Recall that the base category is an article published in the AER, in the field of 'Microeconomics', in 2000, that belongs to the first cluster. For this type of articles, we are further interested in calculating the marginal effect of the prop_women among 4 male (both_genders=0) co-authors (log_num_authors=0.60205999132) at 0, 0.25, 0.5, 0.75, and 1 across the whole spectrum of articles' length (number of pages).

[12]: % stata -qui
 adopath init;
 cd "../data";
 margins, dydx(prop_women) at(prop_women = (0@.25)1) log_num_pages=(1(0.5)4.5)
 log_num_authors=0.60205999132 (base) _factor
 saving(predictions, replace);
 #delimit cr

[13]: % stata -qui
 use predictions, clear
 rename _at1 prop_women
 rename _at4 log_num_pages
 rename _margin delta_hat
 save predictions, replace

```

Simple 0 1 main Python 3 (ipykernel) | idle Mode: Command ⚙ Ln 29, Col 36 004\_logit.ipynb 0

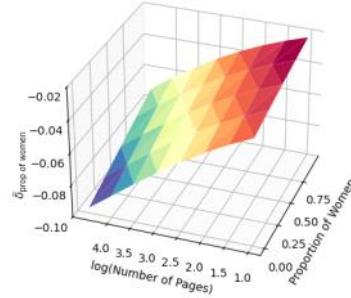
⚠ Please note that  $\hat{\delta}_{\text{prop of women}}(\text{Proportion Women}, \log(\text{Number pages}))$  is an estimate of  $\delta_{\text{prop of women}}(\text{Proportion Women}, \log(\text{Number pages}))$  and therefore one also needs to report s.e. ( $\delta_{\text{prop of women}}(\text{Proportion Women}, \log(\text{Number pages}))$ ) which are calculated via the [Delta Method](#).

$\Pr[\text{FKG\_01}_a | \log(\text{Number authors})_a, \log(\text{Number pages})_a, \text{Both genders } a, \text{Proportion Women } a, \text{Journals, JEL Codes, Cluster, Years, JEL flag}] = G(\beta_1 \times \log(\text{Number authors})_a + \beta_2 \times \log(\text{Number pages})_a + \beta_3 \times \text{Both genders } a + \theta \times \text{Proportion Women } a + \text{Journals} + \text{JEL codes} + \text{Cluster} + \text{Years} + \text{JEL flag} + \text{cons})$

We are interested in the marginal effect of education on the probability of participating in the labor force.

After a Logit regression, one obtains  $[\hat{\beta}_1, \hat{\beta}_2, \hat{\beta}_3, \hat{\theta}, \dots, \hat{\beta}_{262}]$ , we proceed to set  $\log(\text{Number authors}) = 0.60205999132$  (Number authors = 4), Both genders = 0, and all the sets of dummies in `Journals, JEL Codes, Cluster, Years, and JEL flag` equal to 0, i.e., an article published in 2000, in the field of 'Microeconomics,' without author-provided JEL codes, published in the AER. Then the *marginal effect of prop\_women* is

$$\hat{\delta}_{\text{prop of women}}(\text{Proportion Women}, \log(\text{Number pages})) = \hat{\theta}g(\hat{\beta}_1 \times 0.60205999132 + \hat{\beta}_2 \log(\text{Number pages}) + \hat{\theta} \times \text{Proportion Women} + \hat{\beta}_{262}) + \hat{\theta}$$



```

004_logit.ipynb - JupyterLab
004_logit.ipynb
File Edit View Run Kernel Git Tabs Settings Help
Average Marginal Effects

[11]: % stata
 margins, dydx(log_num_authors log_num_pages both_genders prop_women)

 Average marginal effects
 Number of obs = 4,549
 Model VCE: Robust
 Expression: Pr(FKG_01), predict()
 dy/dx wrt: prop_women 1.both_genders log_num_pages log_num_authors

 Delta-method
 | dy/dx std. err. z P>|z| [95% conf. interval]

 -.04657823 | .01120000 0.000 -.08212000 -.01210000
 1.log_num_authors | .00400268 -.01100000 2.27 0.024 -.00390000 .0065437
 log_num_pages | .03787021 .00783800 4.83 0.000 .025862 .0533242
 log_num_authors | -.0198674 .0045374 -4.38 0.000 -.0287685 -.0189742

 Note: dy/dx for factor levels is the discrete change from the base level.


```

Simple 0 1 main Python 3 (ipykernel) | Busy Mode: Command ⚙ Ln 29, Col 36 004\_logit.ipynb 0

⌚ On average if the proportion of women authors increases by 1, the probability of it to be very difficult to read decreases by -.047.

⌚ On average the probability for an article to be very difficult to read when it is written by mixed-gender co-authors is 0.025 higher than single-gender authors holding everything else constant.

# Multinomial Response

A **multinomial** random variable  $Y$  takes values in a finite set, typically written as  $Y \in \{0,1,2,\dots,J\}$ . We typically describe the pair  $(Y, X')$  as **multinomial response** when  $Y$  is multinomial and  $X \in \mathbb{R}^k$  are regressors. The conditional distribution of  $Y$  given  $X$  is summarized by the response probability

$$P_j(x) = \Pr[Y = j \mid X = x].$$

# Poisson Regression

**Count data** refers to situations where the dependent variable is the number of "events" recorded as positive integers  $Y \in \{0,1,2, \dots\}$

|                                                                                                                                                                                    |                                                                                                                                                                                                                                                     |
|------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------|-----------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------|
| <u>Examples:</u>                                                                                                                                                                   |  A count data model specifies the response probabilities $P_j(X) = \Pr[Y = j   X]$ for $j = 0, 1, 2, \dots$ , with the property $\sum_{j=0}^{\infty} P_j(X) = 1$ . |
| <ul style="list-style-type: none"> <li>✓ Number of doctor visits</li> <li>✓ Number of accidents</li> <li>✓ Number of patent registrations</li> <li>✓ Number of absences</li> </ul> |                                                                                                                                                                                                                                                     |

005\_poison.ipynb - JupyterLab

Poisson Regression

We are interested in modelling the number of sentences in an abstract in excess of 3 (5% of articles have at most 3 sentences). i.e. our model is

$$E[Excess\ Sentences|\log(\text{Number authors}), \log(\text{Number pages}), \text{Both genders}, \text{Proportion women}] = \exp(\beta_1 \log(\text{Number authors}) + \beta_2 \log(\text{Number pages}) + \beta_3 \text{Both genders} + \beta_4 \text{Proportion women} + \beta_5 \text{Journals} + \beta_6 \text{JEL codes} + \beta_7 \text{Cluster} + \beta_8 \text{Years} + \beta_9 \text{JEL flag} + \beta_{10})$$

```
[3]: %% stata -qui
summarize num_sentences, detail
local cutoff_sentences = r(p5)
gen int excess_sentences = num_sentences - `cutoff_sentences'
replace excess_sentences = 0 if excess_sentences < 0
```

```
[4]: %% stata
set scheme s1mono
histogram excess_sentences, discrete freq width(1)
```

```
. set scheme s1mono
. histogram excess_sentences, discrete freq width(1)
(start=0, width=1)
```

Frequency

excess\_sentences

| excess_sentences | Frequency (approx.) |
|------------------|---------------------|
| 0                | 500                 |
| 1                | 1200                |
| 2                | 1000                |
| 3                | 800                 |
| 4                | 600                 |
| 5                | 400                 |
| 6                | 200                 |
| 7                | 100                 |
| 8                | 50                  |
| 9                | 20                  |
| 10               | 10                  |
| 11               | 5                   |
| 12               | 2                   |
| 13               | 1                   |
| 14               | 1                   |
| 15               | 1                   |
| 16               | 1                   |
| 17               | 1                   |
| 18               | 1                   |
| 19               | 1                   |
| 20               | 1                   |

## Poisson Regression:

$$P_j(X) = \frac{\exp(-\lambda(X))\lambda(X)^j}{j!}$$

$$\lambda(X) = \exp(X'\beta)$$

Recall that the [Poisson distribution](#) has the property that its mean and variance equal the Poisson parameter  $\lambda$ . Thus (conditionally on  $X$ )

$$\text{var}[Y | X] = \exp(X'\beta)$$

|                          |                                                                                                                           |
|--------------------------|---------------------------------------------------------------------------------------------------------------------------|
| Estimator:               | The MLE $\hat{\beta}$ is the value $\beta$ which maximizes $\ell_n(\beta)$ .                                              |
| Log-likelihood Function: | $\ell_n(\beta) = \sum_{i=1}^n \log P_{Y_i}(X_i   \beta) = \sum_{i=1}^n (-\exp(X'_i \beta) + Y_i X'_i \beta - \log(Y_i!))$ |

1st & 2nd Derivatives:

$$\frac{\partial}{\partial \beta} \ell_n(\beta) = \sum_{i=1}^n X_i (Y_i - \exp(X_i' \beta))$$

$$\frac{\partial^2}{\partial \beta \partial \beta'} \ell_n(\beta) = - \sum_{i=1}^n X_i X_i' \exp(X_i' \beta)$$

### Coefficients Interpretation (Continuous Controls)

Let us rewrite

$$\mathbb{E}[Y | X] = \exp(X' \beta) = \exp(\beta_1 X_1 + \dots + \beta_{k-1} X_{k-1} + \beta_k),$$

then by taking logs one has

$$\log\{\mathbb{E}[Y | X_1, \dots, X_{k-1}]\} = \beta_1 X_1 + \dots + \beta_{k-1} X_{k-1} + \beta_k.$$

Let  $X_j \in \mathbb{R}$ , then

$$\log\{\mathbb{E}[Y | X_1, \dots, X_j + \Delta X_j, \dots, X_{k-1}]\} - \log\{\mathbb{E}[Y | X_1, \dots, X_j, \dots, X_{k-1}]\} = \beta_j \Delta X_j,$$

$$\% \Delta \mathbb{E}[Y | X_1, \dots, X_{k-1}] \approx (100 \beta_j) \Delta X_j.$$

$\beta_j$  is roughly the percentage change in  $\mathbb{E}[Y | X]$ , given a one-unit increase in  $X_j$ .

If  $X_j = \log(Z_j)$ , for some variable  $Z_j > 0$ , then  $\Delta X_j = \log(Z_{j;1}) - \log(Z_{j;0}) = \Delta \log(Z_j)$ , so  
 $\% \Delta \mathbb{E}[Y | X_1, \dots, X_{k-1}] \approx (100 \beta_j) \Delta X_j = \beta_j (100 \cdot \Delta \log(Z_j)) \approx \beta_j \% \Delta Z_j$  and therefore

$$\beta_j \approx \% \Delta \mathbb{E}[Y | X_1, \dots, \log(Z_j), \dots, X_{k-1}] / \% \Delta Z_j.$$

$\beta_j$  is roughly the elasticity of the expected value of  $Y$  with respect to  $Z_j$ .

The proportionate change in  $v$  moving from  $v_0$  to  $v_1$  is  

$$\frac{v_1 - v_0}{v_0} = \Delta v/v_0,$$

assuming of course that  $v_0 \neq 0$ . The percentage change in  $v$  moving from  $v_0$  to  $v_1$  is simply 100 times the proportionate change:

$$\% \Delta v = 100(\Delta v/v_0),$$

the notation "% $\Delta v$ " is read as "the percentage change in  $v$ ". Let  $v_0$  and  $v_1$  be positive values, then, [it can be shown](#) that

$$\log(v_1) - \log(v_0) \approx \frac{v_1 - v_0}{v_0} = \Delta v/v_0,$$

for small changes in  $v$ . Therefore by writing  $\Delta \log(v) = \log(v_1) - \log(v_0)$ , then

$$100 \cdot \Delta \log(v) \approx \% \Delta v.$$

### Coefficient Interpretation (Discrete Controls)

Let  $X_j \in \{0,1\}$ , then the proportionate change in the expected value of  $Y$  from  $X_j = 0$  to  $X_j = 1$  is

$$\frac{\mathbb{E}[Y | X_1, \dots, X_j = 1, \dots, X_{k-1}]}{\mathbb{E}[Y | X_1, \dots, X_j = 0, \dots, X_{k-1}]} - 1 = \frac{\exp(\beta_1 X_1 + \dots + \beta_j \cdot 1 + \dots + \beta_{k-1} X_{k-1} + \beta_k)}{\exp(\beta_1 X_1 + \dots + \beta_j \cdot 0 + \dots + \beta_{k-1} X_{k-1} + \beta_k)} - 1 = \exp(\beta_j) - 1$$

$100 \cdot [\exp(\beta_j) - 1]$  is the percentage change in  $\mathbb{E}[Y | X]$  from  $X_j = 0$  to  $X_j = 1$ .

005\_poisson... - JupyterLab

## 005\_poisson.ipynb

File Edit View Run Kernel Git Tabs Settings Help

Markdown git excess\_sentences

Python 3 (ipykernel)

```
[5]: % stata -qui
#delimit;
poisson excess_sentences log_num_authors log_num_pages both_genders prop_women
`journals'`jel_imp' y_2-y_20 c_2-c_215 jel_flag, vce(cluster cluster);
#delimit cr
```

```
[6]: % stata
estimates table, keep(log_num_authors log_num_pages both_genders prop_women jel_flag) b(%5.4f) star varwidth(50) varlabel
```

| Variable                                           | Active    |
|----------------------------------------------------|-----------|
| Logarithm of the number of authors                 | -0.0118   |
| Logarithm of the number of pages                   | 0.4093*** |
| Indicator if author team includes both genders     | 0.0004    |
| Proportion of women among the authors              | 0.0416*   |
| Indicator of whether the article includes a JEL cl | 0.0349*   |

Legend: \* p<0.05; \*\* p<0.01; \*\*\* p<0.001

```
[7]: % stata
nlcom (_b[log_num_pages]) (100*_b[prop_women]) (100*(exp(_b[jel_flag])-1))
```

| excess_sen_s | Coefficient | Std. err. | z     | P> z  | [95% conf. interval] |
|--------------|-------------|-----------|-------|-------|----------------------|
| _n1_1        | .4093298    | .0162724  | 25.15 | 0.000 | .3774364 .4412231    |
| _n1_2        | 4.157834    | 2.077484  | 2.00  | 0.045 | .0860391 8.229628    |
| _n1_3        | 3.54721     | 1.444164  | 2.46  | 0.014 | .7166994 6.37772     |

Other factors equal:

- The elasticity of the expected number of sentences in excess of 3 for an article's abstract with respect to the article's number of pages is 0.41.
- If an article's authors are all women (`prop_women` goes from 0 to 1), the abstract's expected number of sentences in excess of 3 increase by 4.16%.
- The expected number of sentences in excess of 3 for an article's abstract is roughly 3.55% higher when the JEL is imputed *ceteris paribus*.

Simple 0 5 main Python 3 (ipykernel) | Idle Mode: Command Ln 5, Col 106 005\_poisson.ipynb

# HDFE & SCA

We are interested in estimating  $\beta \equiv [\theta, \beta_2']'$  in the model

$$\mathbb{E}[Y|X_1, X_2, D_1, D_2, \dots, D_F] = \exp(\theta X_1 + X_2' \beta_2 + D_1' \alpha_1 + D_2' \alpha_2 + \dots + D_F' \alpha_F),$$

where  $[D'_1, D'_2, \dots, D'_F]$  represents a row of  $\mathbf{D} = [\mathbf{D}_1 \quad \mathbf{D}_2 \quad \dots \quad \mathbf{D}_F]$  which consists of  $F$  indicator matrices.

Manuscript: <https://arxiv.org/abs/1903.01690>  
Source Code: <https://github.com/sergiocorreia/ppmlhdf>

PPMLHDFE: Fast Poisson Estimation with  
High-Dimensional Fixed Effects

Sergio Correia<sup>1</sup>, Paulo Guimarães<sup>2</sup>, and Tom Zylkin<sup>3</sup>  
<sup>1</sup>Federal Reserve Board, sergio.a.correia@frb.gov  
<sup>2</sup>Banco de Portugal, pguimaraes@bportugal.pt  
<sup>3</sup>University of Richmond, tzylin@richmond.edu

August 5, 2019

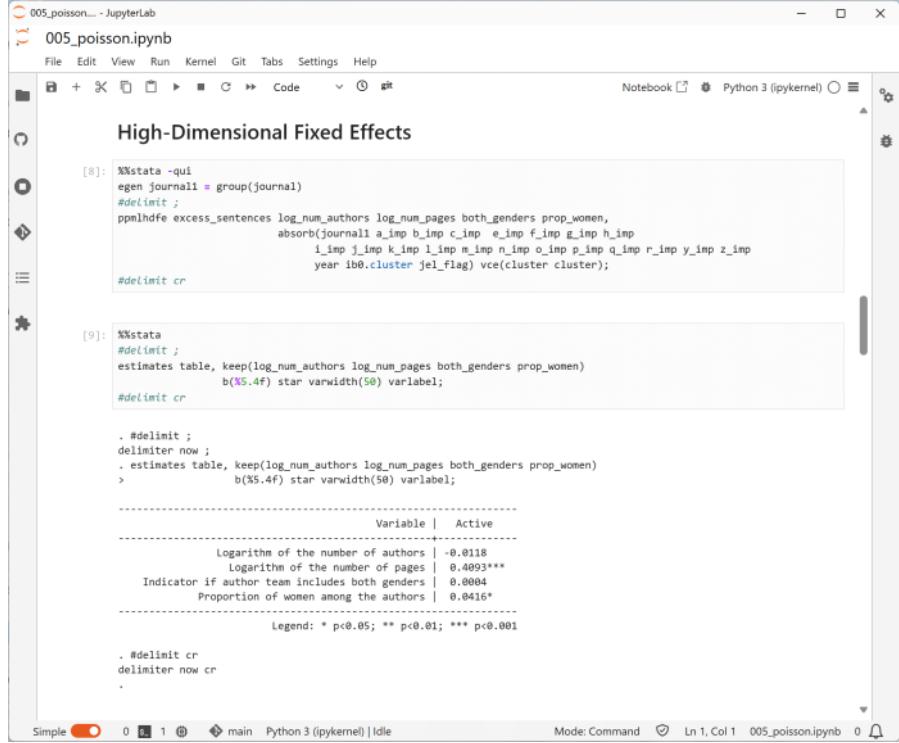
## Abstract

In this paper we present `ppmlhdf`, a new Stata command for estimation of (pseudo) Poisson regression models with multiple high-dimensional fixed effects (HDFE). Estimation is implemented using a modified version of the iteratively reweighted least squares (IRLS) algorithm, which is well known to converge rapidly in the presence of HDFE. Because the code is built around the `regdfe` package, it has similar syntax, supports many of the same functionalities, and benefits from `regdfe`'s fast convergence properties for estimating high-dimensional least squares (LSS) Poisson regression. Compared to `regdfe`, `ppmlhdf` also offers a fast alternative for accelerating HDFE-IRLS estimation specifically. `ppmlhdf` also implements a novel and more robust approach to check for the existence of (pseudo) maximum likelihood estimates.

Keywords: `ppmlhdf`, `regdfe`, Poisson regression, high-dimensional fixed-effects

## 1 Introduction

Poisson regression is now well established as the standard approach to model count data. However, it is also gaining popularity as a viable alternative for estimation of multiplicative models where the dependent variable is nonnegative. Specifically, these models are estimated by regressing the log of the dependent variable. Due to the fact that, as with ordinary least squares (OLS), the only assumption required for consistency of the Poisson regression estimator is the correct specification of the conditional mean of the dependent variable (Gourieroux et al., 1984). In this setting, Poisson regression becomes Poisson pseudo maximum likelihood (PPML) regression. Gourieroux et al.'s results greatly extend the realm of application of Poisson regression because there is no



```
005_poisson.... - JupyterLab
005_poisson.ipynb
File Edit View Run Kernel Git Tabs Settings Help
Notebook Python 3 (ipykernel)
High-Dimensional Fixed Effects

[8]: %%stata -qui
egen journali = group(journal)
#delimit ;
ppmlhdf excess_sentences log_num_authors log_num_pages both_genders prop_women,
absorb(journali a_imp b_imp c_imp e_imp f_imp g_imp h_imp
i_imp j_imp k_imp l_imp m_imp n_imp o_imp p_imp q_imp r_imp y_imp z_imp
year ib0.cluster jel_flag) vce(cluster cluster);
#delimit cr

[9]: %%stata
#delimit ;
estimates table, keep(log_num_authors log_num_pages both_genders prop_women)
b(%5.4f) star varwidth(50) varlabel;
#delimit cr

. #delimit ;
delimiter now ;
estimates table, keep(log_num_authors log_num_pages both_genders prop_women)
> b(%5.4f) star varwidth(50) varlabel;

----- Variable | Active -----
Logarithm of the number of authors | -0.0118
Logarithm of the number of pages | 0.4093***
Indicator if author team includes both genders | 0.0004
Proportion of women among the authors | 0.0416*

Legend: * p<0.05; ** p<0.01; *** p<0.001

. #delimit cr
delimiter now cr
.
```

005\_poisson.... - JupyterLab

005\_poisson.ipynb

File Edit View Run Kernel Git Tabs Settings Help

Notebook Python 3 (ipykernel)

## Specification Curve Analysis

Clearing the data set in memory, reloading, and creating the `journal1` variable

```
[10]: %%stata -qui
clear all
use "../data/data", clear
egen journal1 = group(journal)

summarize num_sentences, detail
local cutoff_sentences = r(p5)
gen int excess_sentences = num_sentences - `cutoff_sentences'
replace excess_sentences = 0 if excess_sentences < 0
```

Commenting out the `- log(F-K grade): log_flesch_kincaid_grade_level` line in the `.yml` file.

```
[11]: import yaml

File name
file_path = '../code/readability_graph_speccurve.yml'

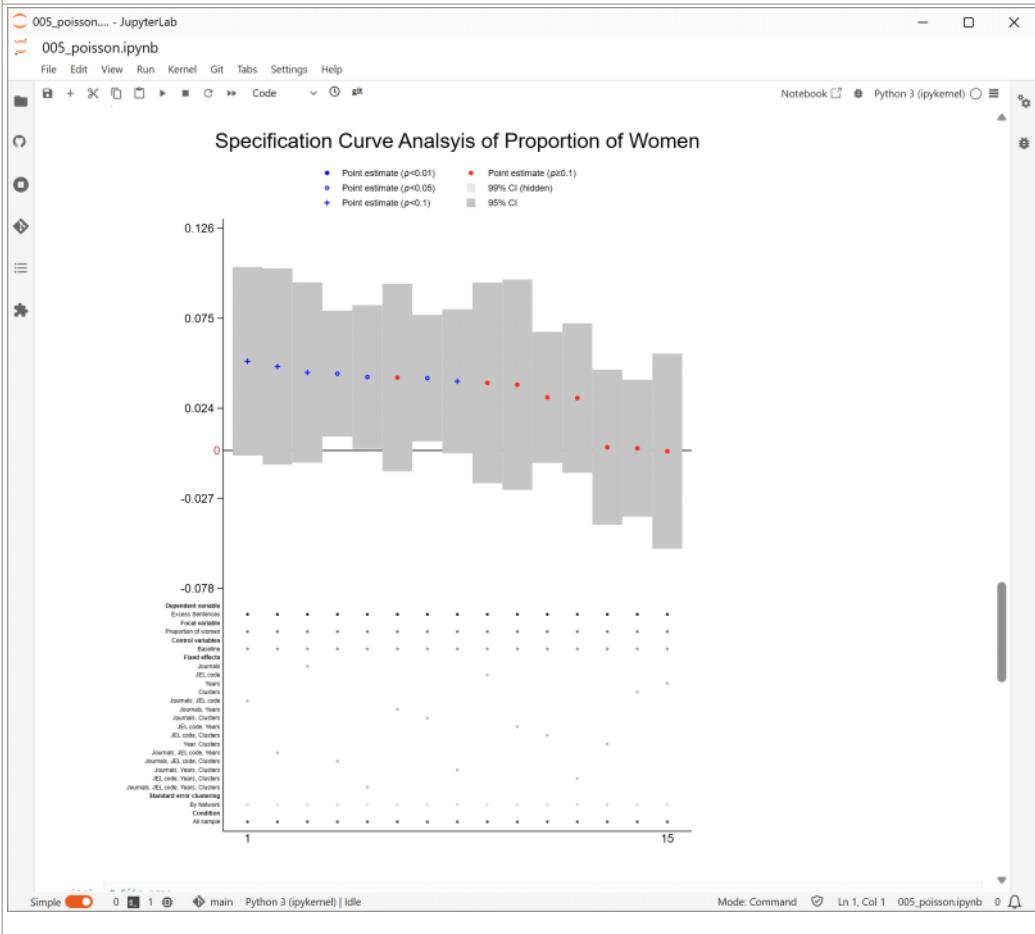
Read, modify, and save the file
with open(file_path, 'r') as file:
 lines = file.readlines()

Modify the target line
modified_lines = []
for line in lines:
 if line.strip() == "- log(F-K grade): log_flesch_kincaid_grade_level":
 modified_lines.append(f"# {line}") # Comment out the line
 else:
 modified_lines.append(line)

Write the modified content back to the file
with open(file_path, 'w') as file:
 file.writelines(modified_lines)

print("Line commented successfully!")
```

Simple 0 1 main Python 3 (ipykernel) | Idle Mode: Command Ln 1, Col 1 005\_poisson.ipynb 0

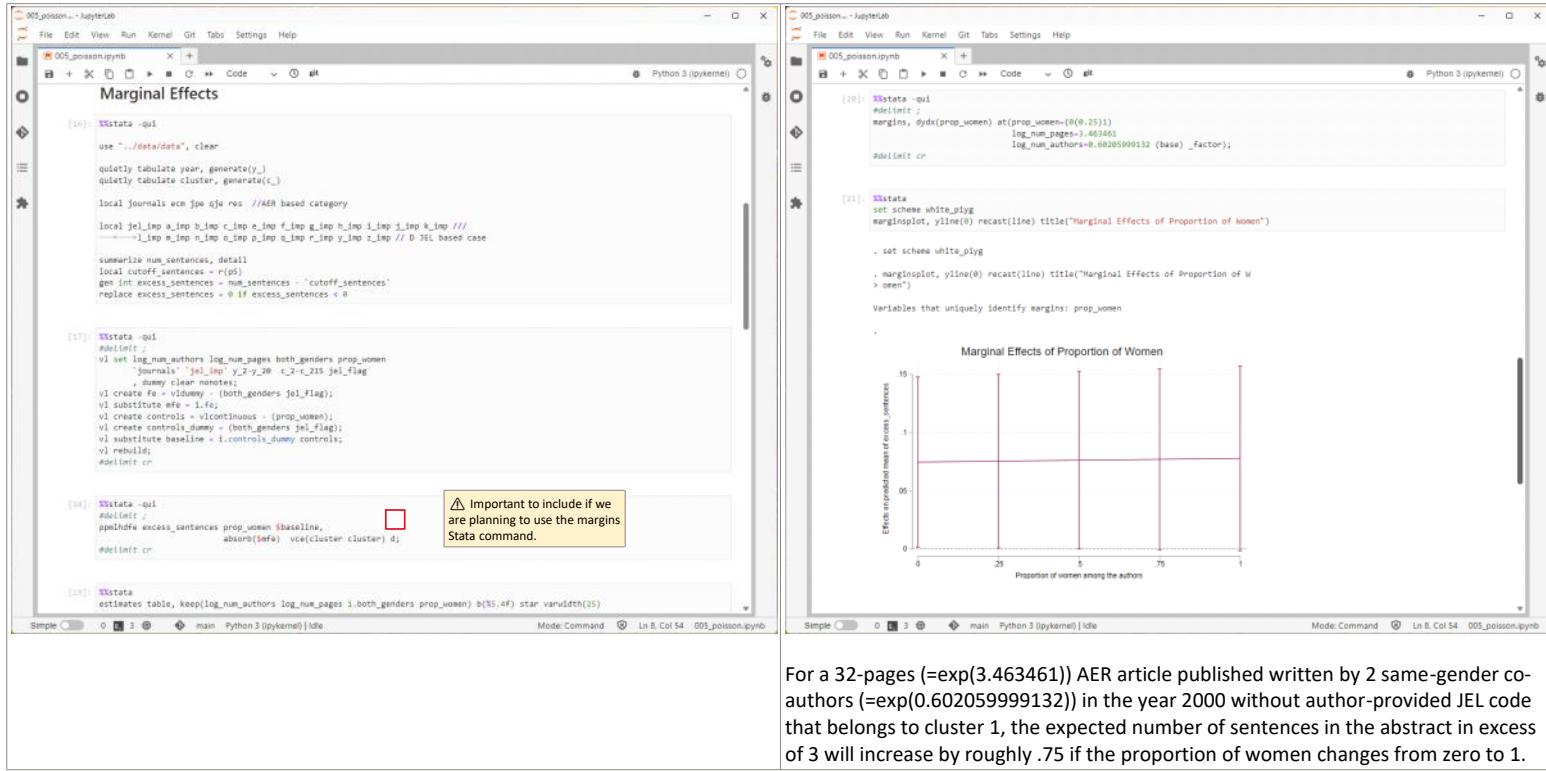


# Marginal Effects

Recall that  $\mathbb{E}[Y | X = x] = \exp(x'\beta) = \exp(\beta_1 x_1 + \dots + \beta_j x_j + \dots + \beta_{k-1} x_{k-1} + \beta_k)$

|                          | $x_j$ is continuous                                                 | $x_j$ is binary                                                                                                                                                                           |
|--------------------------|---------------------------------------------------------------------|-------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------|
| Marginal Effects:        | $\delta_j(x) = \beta_j \exp(x'\beta)$                               | $\delta_j(x) = \exp(\beta_1 x_1 + \dots + \beta_j \cdot 1 + \dots + \beta_{k-1} x_{k-1} + \beta_k) - \exp(\beta_1 x_1 + \dots + \beta_j \cdot 0 + \dots + \beta_{k-1} x_{k-1} + \beta_k)$ |
| Average Marginal Effects | $AME = \mathbb{E}[\delta_j(X)] = \beta_j \mathbb{E}[\exp(X'\beta)]$ | $AME = \mathbb{E}[\delta_j(X)]$                                                                                                                                                           |

Estimator:  $\widehat{AME} = \frac{1}{n} \sum_{i=1}^n \hat{\delta}_j(X_i)$ , where  $\hat{\delta}_j(X_i) = \hat{\beta}_j \exp(X_i'\hat{\beta})$  or  
 $\hat{\delta}_j(X_i) = \exp(\hat{\beta}_1 X_{1,i} + \dots + \hat{\beta}_j \cdot 1 + \dots + \hat{\beta}_{k-1} X_{k-1,i} + \hat{\beta}_k) - G(\hat{\beta}_1 X_{1,i} + \dots + \hat{\beta}_j \cdot 0 + \dots + \hat{\beta}_{k-1} X_{k-1,i} + \hat{\beta}_k)$ .



For a 32-pages ( $=\exp(3.463461)$ ) AER article published written by 2 same-gender co-authors ( $=\exp(0.602059999132)$ ) in the year 2000 without author-provided JEL code that belongs to cluster 1, the expected number of sentences in the abstract in excess of 3 will increase by roughly .75 if the proportion of women changes from zero to 1.