

Class 8 - Techniques I - Regression

Agenda

- The logic of regression: what, why, when, how (40 minutes)
- Application paper discussion (20 minutes)
- *Break*
- Replication presentation (Group 1-3-5; 40 minutes)
- General discussion (15 minutes)

The logic of regression

Preamble

I will weave in the “core” papers as needed in my remarks later, but we will not discuss them in detail per se

WARNING: I am going deeper into the weeds this week, but we will come back up to a higher level soon

What is a regression?

The multiple linear regression **model** is used to study the relationship between a dependent variable (y) and one or more independent variables (X). The generic form of the linear regression model is (Greene 2012):

$$y = X\beta + \epsilon$$

where β is computed as follows:

$$(\mathbf{X}^T \mathbf{X})^{-1} (\mathbf{X}^T \mathbf{y}) = \beta$$

Let's figure out why...

What is a beta coefficient?

Small aside: Note that in a “simple” (bivariate) regression with just one predictor and an intercept

$$y = \beta_0 + \beta_1 x + \epsilon$$

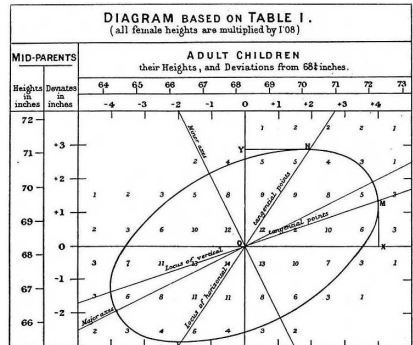
this funky formula is essentially:

$$Cov(x, y)/Var(x) = \beta_1$$

Note that in the multiple regression case, the same logic applies but to covariances and variances that have been “residualized” by accounting for all of the other X variables in the equation (Angrist and Pischke 2008)

What's with the name "regression"?

“Galton observed that extreme characteristics (e.g., height) in parents are not passed on completely to their offspring. Rather, the characteristics in the offspring **regress** toward a mediocre point (a point which has since been identified as the mean)” (Wikipedia)



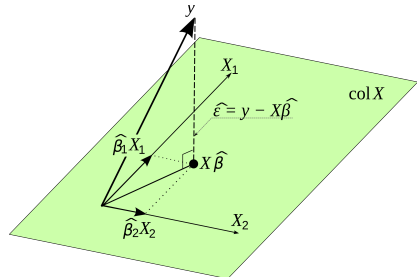
What do we usually mean when we say “regression”?

The term “regression” is often used synonymously with a particular technique: ordinary least squares (OLS).

- Essentially, OLS and its generalizations (e.g., WLS, GLS, discussed later) seek to minimize the distance between observed and predicted values
- It is constrained in this task in that only the data provided by X can be used, and that there are too many observations to simply interpolate

What is OLS doing?

- The data in X forms a (hyper)-plane, and each X variable is **weighted** by β , resulting in a prediction
- The discrepancy between this predicted value $X\beta$ and y is the error term ϵ , which by definition is unrelated to the X variables



A sample data generating process

Let's say we **know** that the relationship between X and Y is as follows:

$$y = (\beta_0 = 0) + (\beta_1 = 2)x + \epsilon$$

and that ϵ is randomly distributed with mean of 0 and a standard deviation of 3.

This is an example **data generating process**, what might such data look like?

A sample data generating process

```
# A tibble: 10 x 2
```

	x	y
	<dbl>	<dbl>
1	1	-0.714
2	2	7.28
3	3	3.00
4	4	6.05
5	5	10.6
6	6	10.2
7	7	13.4
8	1	-2.31
9	2	3.78
10	3	7.14

OLS as a mathematical procedure

The population regression equation applies to each observation (thus the matrix form):

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

and beta is **selected** / **defined** to minimize the sum of the squared errors (i.e., positive and negative errors don't cancel out):

$$\min(\boldsymbol{\epsilon}^T \boldsymbol{\epsilon}) = (\mathbf{y} - \mathbf{X}\boldsymbol{\beta})^T (\mathbf{y} - \mathbf{X}\boldsymbol{\beta})$$

In other words, we ask; what value of beta should we select to minimize the LHS?

OLS as a mathematical procedure

This minimization problem (with respect to β) leads to the “normal equations” <- click on the link for a full derivation.

$$(\mathbf{X}^T \mathbf{X})^{-1} (\mathbf{X}^T \mathbf{y}) = \beta$$

Comments:

- What would happen the if error/disturbance term was always zero?
- Notice, this equation is agnostic to the form / distribution of the errors

OLS as an “estimator”

Remember, we don't have the whole population, we just have a sample. Thus, the **residuals** we get from an OLS estimation are different from the underlying **errors** in the population due to **sampling error** in our estimate of β

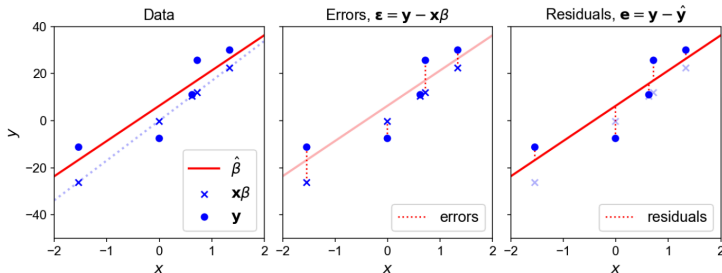


Figure 1. (Left) Ground-truth (un-observed) univariate data $\mathbf{x}\beta$, noisy observations $\mathbf{y} = \mathbf{x}\beta + \epsilon$, and estimated $\hat{\beta}$.

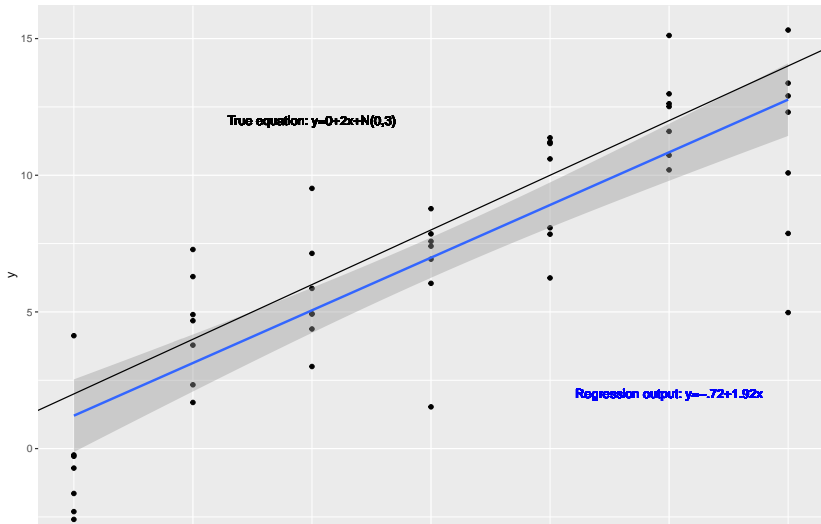
OLS as an “estimator”

Here is another example of this distinction with the data we generated earlier. Recall the “true” parameters:

$$y = (\beta_0 = 0) + (\beta_1 = 2)x + \epsilon$$

	Estimate	Std. Error	t value	Pr(> t)
(Intercept)	-0.7233802	0.8190561	-0.8831876	3.816302e-01
x	1.9272462	0.1831465	10.5229748	6.023326e-14

OLS - Estimation v. reality



Verifying our beta coefficient formula

Let's verify that our simple equation for \hat{b} works out:

$$\hat{Cov}(x, y) / \hat{Var}(x) = \hat{b}_1$$

```
[1] "cov(x,y) 7.8695885474943"
```

```
[1] "var(x) 4.08333333333333"
```

```
[1] "beta-hat 1.92724617489656"
```

What is a regression? - Prediction versus explanation

Notice a few things:

- OLS regression on a sample, by construction, minimizes the sum of the squared residuals
- This means that it seeks to maximize the amount explained by the regression to maximize the chances of correctly predicting the data **in the sample**
- By implication, it is NOT trying to find some true value of β such that we can make good **out of sample** predictions; issues with the sample will contaminate beta
- It also does NOT imply an causal interpretation or explanation

Why linear regression and OLS in particular?

OK, then why are regressions ubiquitous?

- 1 Approximates the conditional expectation function (Angrist and Pischke 2008)
- 2 The “best linear unbiased estimator” or BLUE, when the Gauss-Markov assumptions met (Kennedy 2008)
- 3 Equivalent to maximum likelihood estimator (MLE) and maximum a posteriori estimator (MAP) when errors normally distributed (Kennedy (2008), p. 43) and with a uniform prior

Let's consider each in turn.

“Approximates the conditional expectation function”

Recall that much of the information we want to extract from a distribution is summarized in its first two moments: mean and variance

The conditional expectation function [CEF] expresses the expected value (or mean) of a variable (Y) as a function of another variable (X)

The CEF is often not linear, but the OLS estimator is the best linear approximation to the CEF function (Angrist and Pischke 2008, Theorem 3.1.6)

“Approximates the conditional expectation function”

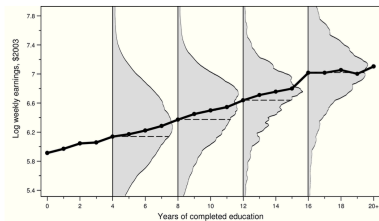
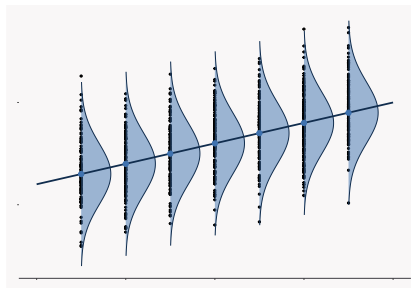


Figure 2: Angrist and Pischke (2008), Ch. 3

Note what this picture is showing you - it highlights the fact that Y has a distribution at each level of X , and what OLS can do is determine how that distribution “shifts” as you change X

“Approximates the conditional expectation function”

The “standard”, unmodified form of OLS assumes that the distribution of Y remains the same across all levels of X (i.e., the distribution is homoskedastic or possesses “spherical errors”), and that the only difference is a shift in where the mean is placed



“Best linear unbiased estimator” (BLUE)

If the Gauss-Markov assumptions are satisfied, then OLS is the:

- **Best** (minimum variance among alternative estimators)
- **Linear**(the estimator is of the form $A(X)y$)
- **Unbiased** (the expected value is equal to the population parameter)
- **Estimator** (employs data from a sample to make inferences about the population)

Recent research seems to suggest that it might actually be the best unbiased estimator if these assumptions hold

“Equivalent to MLE and MAP under certain conditions”

There are other modeling frameworks to help make sense of data.

- Maximum likelihood estimation or MLE: What parameter values make the data most likely, GIVEN an assumed distribution?
- Maximum a posteriori (Bayesian) or MAP estimation: MLE that explicitly incorporates prior beliefs
- Generalized method of moments: A “generalization” of MLE that focuses on moments rather than distributions

Note that MLE and GMM are the primary alternatives in our field to OLS/WLS/GLS: they are used for SEM, HLM, panel data estimators, logit, probit, and other models

“Equivalent to MLE and MAP under certain conditions”

When we have a uniform prior (all parameter values within an interval are equally likely) and when the error term is normally distributed, OLS provides the same estimate as MLE and MAP.

This should be somewhat comforting, since it implies:

- The OLS estimate b is the most likely value for β , given the data and assuming the distribution is normal
- The OLS estimate is what we should update our beliefs to be if we did not have an informative prior

Aside: An example of how likelihood functions work

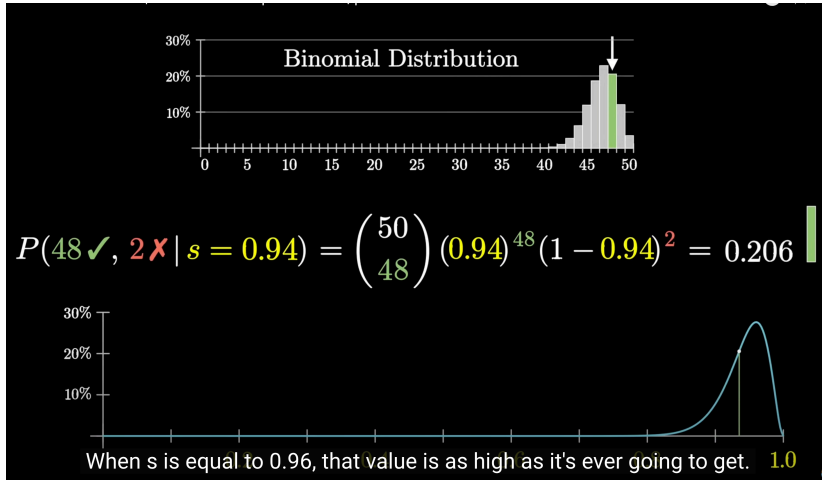


Figure 3: An example of a likelihood function

When is linear regression appropriate and / or optimal?

When the assumptions are satisfied!

And what are those assumptions? There are two “flavors”:

- Fixed design
- Random design

I've tried to align the numbering so that they correspond, where possible.

When is linear regression appropriate and / or optimal?

Table 3.1 The assumptions of the CLR model.

Assumption	Mathematical expression		Violations	Chapter in which discussed
	Bivariate	Multivariate		
1. Dependent variable a linear function of a specific set of independent variables, plus a disturbance	$y_t = \beta_0 + \beta_1 x_t + \varepsilon_t$ $t = 1, \dots, N$	$Y = X\beta + \varepsilon$	Wrong regressors Nonlinearity Changing parameters	6
2. Expected value of disturbance term is zero	$E\varepsilon_t = 0$, for all t	$E\varepsilon = 0$	Biased intercept	7
3. Disturbances have uniform variance and are uncorrelated	$E\varepsilon_t \varepsilon_r = 0$, $t \neq r$ $= \sigma^2$, $t = r$	$E\varepsilon \varepsilon' = \sigma^2 I$	Heteroskedasticity Autocorrelated errors	8
4. Observations on independent variables can be considered fixed in repeated samples	x_t fixed in repeated samples	X fixed in repeated samples	Errors in variables Autoregression Simultaneous equations	10 11
5. No exact linear relationships between independent variables and more observations than independent variables	$\sum_{t=1}^N (x_t - \bar{x})^2 \neq 0$	Rank of $X = K \leq N$	Perfect multicollinearity	12

The mathematical terminology is explained in the technical notes to this section. The notation is as follows: Y is a vector of observations on the dependent variable; X is a matrix of observations on the independent variables; ε is a vector of disturbances;

When is linear regression appropriate and / or optimal?

“Random design”

- 1 Linearity (model is correctly specified)
- 2
- 3 Spherical errors $V[\epsilon|\mathbf{X}] = \sigma^2 \mathbf{I}$
 - implies: Homoskedasticity $V[\epsilon|\mathbf{X}] = \sigma^2$
 - implies: No serial correlation $V[\epsilon_i \epsilon_j | \mathbf{X}] = 0; i \neq j$
- 4 Strict exogeneity $E[\epsilon|\mathbf{X}] = 0$
- 5 No multicollinearity (\mathbf{X} is invertible)
- 6 Normality (optional) - allows us to make inferences about the sampling distribution of b

When is linear regression appropriate and / or optimal?

The Gauss-Markov theorem proves that, given these assumptions, the OLS estimator is the BLUE

- It becomes the best unbiased estimator if normality is assumed (it reaches the Cramer-Rao lower bound)
- The assumption of a spherical errors can be relaxed via Aitken's theorem, known as weighted (WLS) or generalized (GLS) least squares
- GLS often performs a transformation of the raw data to take into consideration the variance structure

Aside: alternative estimators when OLS/GLS not applicable

- Maximum likelihood estimation (MLE)
- Bayesian / maximum a posteriori (MAP) techniques
- Generalized method of moments (GMM) estimators
- Quantile regression
- Machine learning (e.g., random forest models)
- LASSO (least absolute shrinkage and selection operation)
- Markov Chain Monte Carlo simulation and bootstrapping
- K-L divergence based metrics (e.g., Expectation-Maximization algorithms)

Running a regression

How to properly perform a regression analysis?

- 1 We determine the appropriateness of regression for our model and data from two angles:
 - Column (variable) perspective
 - Row (observation) perspective
- 2 Once we are satisfied that there are no issues (or that they have been addressed), run a “final, publication ready” regression.
- 3 Interpret the results using the tools of statistical inference (discussed later today).

The column (variable) perspective

Conditional mean structure

- 1 Do we have the right variables? (DAGs, theory)
- 2 Is the model in the correct functional form specified? (RESET, Chow tests)
- 3 Are the variables exogenous? (Sargen-Hansen, Durbin-Wu-Hausman)
- 4 Are your models well-conditioned? (VIF)

Variance structure

- 5 Is the mean structure correct? (see above)
- 6 Tests of heteroskedasticity and auto-correlation (e.g., Durbin-Watson, White tests)

The row (observation) perspective

Conditional mean structure

- 1 Are there outliers present? (Cook's distance)
- 2 Are there influential observations? (hat matrix)
- 3 Are certain observations missing or partially missing?
(Censoring and truncation)
- 4 Is the model dynamic in nature? (Arellano-Bond estimators)

Variance structure

- 5 Is there a nested structure to the data? (Level of analysis)

Deep dive: Do we have the right variables?

Note that we don't really have a “test” to apply here, but the assumption of proper specification is critical (including the right variables)

Why? Because if it is violated we don't only lose the “best” part of BLUE, we also lose the “unbiased” part, and then we are just “LE”

And furthermore, the estimator isn't even “consistent” -> i.e., gets closer to the true value as we increase the size of the sample, which is a “fall-back” position we can take if we don't have an unbiased estimator for smaller samples

The issue: Omitted variable bias

Omitted variable bias causes **ALL** of the beta estimates in a model to be biased in proportion the correlation between the regressors and the omitted variable left in the error term

This bias can be expressed in terms of “short” and “long” regression estimates, where short means a relevant variable b_2 has been omitted and the long is where it has been included (and it is assumed that the error term is truly uncorrelated in the case of the long equation):

In the simple case, the upshot is the short regression coefficient b_1 will be biased in the following manner:

$$b_{1short} = b_{1long} + b_{2omitted}(cov(b_1, b_2)/var(b_1))$$

The issue: Omitted variable bias

The problem: remember that OLS **assumes by construction** that the error term and the X variables are uncorrelated, so we can't rely on a quick test of their correlation to draw conclusions. You need to **know from theory or experience** that you are missing the variable and then can run the test of including it to see if it changes the results (assuming it is something you can measure!)

What can you do if this variable is unobservable? [See our endogeneity day!]

The flip side: The “illusion of statistical control”

By mapping practices to these purposes, we demonstrate why current CV practice struggles to accomplish any of them effectively while potentially reducing the interpretability of results. Empirically, we examine correlations between CVs and independent variables (IVs)—relationships that receive little consideration in standard CV practice—demonstrating these relationships can and do influence the magnitude and sign of regression coefficients, sometimes dramatically—just not very often. In fact, the CVs in the studies we reviewed most often have little if any impact on research findings or interpretations, creating the illusion of statistical control when little control actually occurs. (Carlson and Wu 2012, 415)

What if we fail to meet one or more of the assumptions?

An econometrics textbook can be characterized as a catalog of which estimators are most desirable in what estimating situations. Thus, a researcher facing a particular estimating problem simply turns to the catalog to determine which estimator is most appropriate for him or her to employ in that situation. The purpose of this chapter is to explain how this catalog is structured. (Kennedy 2008, 40)

What if we fail to meet one or more of the assumptions?

But note that:

If more than one of the CLR model assumptions is violated at the same time, econometricians often find themselves in trouble because their catalogs usually tell them what to do if only one of the CLR model assumptions is violated. Much recent econometric research examines situations in which two assumptions of the CLR model are violated simultaneously. These situations will be discussed when appropriate. (Kennedy 2008, 44)

Statistical inference

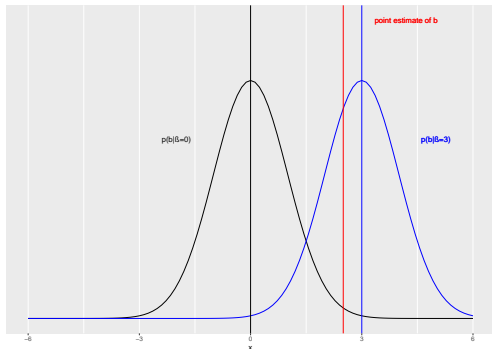
Are our results “statistically significant”?

Once we have an estimate for β (known as \hat{b}), here are two approaches to statistical inference (Kennedy 2008):

- Explicitly assume that the error term is distributed normally (t-tests and F-tests are appropriate straight away, even for small samples)
- Rely on the central limit theorem and large-N asymptotics (beta is an expected value, so its distribution will convergence to a normal distribution - see this video to understand why)

The distribution of \hat{b}

When making an inference, we ask where our point estimate (what we got from the regression) sits within the distribution of the estimator conditional on some assumed true value of β



Determining appropriate standard errors

Questions about the variance structure also need to be settled:

- Are errors spherical?
- If not, do you want to incorporate the error structure information into the estimates (using GLS) or use OLS and “robust” statistics”?
- Two schools of thought: Feasible GLS has no guarantees of improvement, but “Simply computing a robust covariance matrix for an otherwise inconsistent estimator does not give it redemption” (Greene 2012, 692)

Statistical inferences you can draw

The simplest test:

- Is $\beta \neq 0$?: Use t-tests

Intuition: We assume a null distribution that $\beta = 0$, and then see where our estimated value falls within that distribution. The more extreme it is, the more likely the assumption that $\beta = 0$ is false

Statistical inferences you can draw

A more powerful set of tests: linear hypotheses

- Is $\beta_1 = \beta_2$?
- Is $\beta_1 = -\beta_2 + 2\beta_3$?
- Is $\beta_1 = \beta_2 = \beta_3 = 0$?

A restriction matrix is constructed to run an F-test (Wald test), Likelihood Ratio (LR) or LaGrange Multiplier (LM) test (Kennedy 2008, ch. 4)

Intuition: Unconstrained optimization (remember we are minimizing least squares) is easier than constrained. We ask whether there is a significant difference in model fit after imposing the restriction: If yes, then the restriction is probably false.

Examples:

“Linear hypothesis tests fail to reject the assertion that growth and reductions in either team size ($\chi^2(1) = 0.39$, $p > 0.10$) or within-industry experiences ($\chi^2(1) = 0.17$, $p > 0.10$) have equal and opposite effects” (Fox, Simsek, and Heavey 2023, 17)

“A Wald test indicates that the hypothesis that both variables are simultaneously zero is rejected ($\chi^2(2) = 7.50$, $p = .02$)” (Fox, Simsek, and Heavey 2022, 18, FN)

Aside: Statistical inference and p-hacking

- “Searching for asterisks” is all about demonstrating that a beta coefficient is “statistically significant” and differs from 0.
- It does not consider economic / practical effect size
- There are many tools in the toolbox to “manage” statistical significance
- More useful things to focus on are effect size, confidence intervals, and robustness to different specifications of mean and variance structures

Statistical inferences and p-hacking

On the flip side, being careful when performing your analyses can stop you from making faulty inferences - even if it is painful in the short run.

T-stat looks too good. Use standard errors- significance gone. -Keisuke Hirano (Angrist and Pischke 2008, 6)

The flip side of significance: Statistical power

Don't forget to give yourself a fighting chance to find the effect you are looking for!

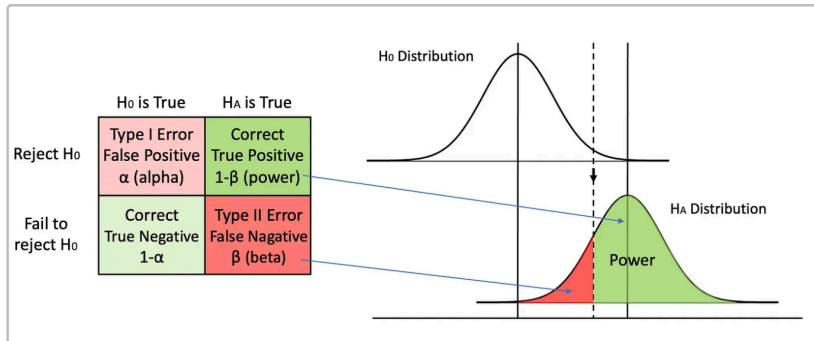


Figure 5: Source: Yao 2021

Applications

Application readings

Let's level-set people's familiarity with these pieces.

- Katila, R., & Ahuja, G. 2002. Something Old, Something New: A Longitudinal Study of Search Behavior and New Product Introduction. *Academy of Management Journal*, 45(6), 1183-1194.
- Simsek, Z., Fox, B., & Heavey, C. 2021. Systematicity in Organizational Research Literature Reviews: A Framework and Assessment. *Organizational Research Methods*, 109442812110086.

Break



COFFEE BREAK

Replication Presentation

- Replication: Simsek, Z., Fox, B., & Heavey, C. 2021. Systematicity in Organizational Research Literature Reviews: A Framework and Assessment. Organizational Research Methods, 109442812110086.

General Discussion

Preparation for next class

Next class

Concept check 2 will be available this evening - please complete it before next class.

Next class

Techniques II: Moderation

- 1 Dawson, J. F. 2014. Moderation in Management Research: What, Why, When, and How. Journal of Business and Psychology, 29(1), 1-19.
- 2 Hitt, M. A., Beamish, P. W., Jackson, S. E., & Mathieu, J. E. 2007. Building Theoretical and Empirical Bridges Across Levels: Multilevel Research in Management. Academy of Management Journal, 50(6), 1385-1399.

Next class

Techniques II: Moderation

Applications:

- 3 Replication: Heavey, C., Simsek, Z., & Fox, B. C. 2015. Managerial Social Networks and Ambidexterity of SMEs: The Moderating Role of a Proactive Commitment to Innovation. *Human Resource Management*, 54(S1).
- 4 Wolfson, M. A., & Mathieu, J. E. 2018. Sprinting to the finish: Toward a theory of Human Capital Resource Complementarity. *J Appl Psychol*, 103(11), 1165-1180.

References

- Angrist, J. D., and J. S Pischke. 2008. *Mostly Harmless Econometrics: An Empiricist's Companion*. Princeton, NJ: Princeton University Press.
- Carlson, Kevin D., and Jinpei Wu. 2012. "The Illusion of Statistical Control." *Organizational Research Methods* 15 (3): 413–35.
- Fox, Brian C., Zeki Simsek, and Ciaran Heavey. 2022. "Top Management Team Experiential Variety, Competitive Repertoires, and Firm Performance: Examining the Law of Requisite Variety in the 3D Printing Industry (1986–2017)." *Academy of Management Journal* 65 (2): 545–76.
- . 2023. "Venture Team Membership Dynamics and New Venture Innovation." *Strategic Entrepreneurship Journal*.
- Greene, William H. 2012. "Econometric Analysis."
- Kennedy, Peter. 2008. *A Guide to Econometrics*. Malden, MA: