# Empirical Economics - Cheat Sheet

## ① Causal Effects, Experiments & Regression Analysis

Potential-Outcome for binary $D_i \in \{0; 1\}$, Outcome $Y_i$ (potential outcome if treated: $Y_{1i}$, untreated: $Y_{0i}$)

$Y_i = Y_{0i} + (Y_{1i} - Y_{0i}) D_i \Rightarrow \tau_i = Y_{1i} - Y_{0i}$ (causal effect)

Only one known $\Rightarrow$ makes causal inference difficult

ATE: $E[Y_{1i} - Y_{0i}]$    ATT: $E[Y_{1i} - Y_{0i} | D_i = 1]$

Observed difference in means: $E[Y_i | D_i = 1] - E[Y_i | D_i = 0]$
$= ATT + $ Selection bias, where $S_{bias} = E[Y_{0i} | D_i = 1] - E[Y_{0i} | D_i = 0]$

Under random assignment: $E[Y_{0i} | D_i = 1] = E[Y_{0i} | D_i = 0]$
$\Rightarrow$ unbiased estimation. $Y_i = \alpha + \beta D_i + \varepsilon_i$   $E[Y_{1i} - Y_{0i}]$

$\rightarrow$ Control variables can be added to improve the model. Heterogenous TE: $Y_i = \alpha + (\beta + \tau_i) D_i + \varepsilon_i$

## ② OLS Recap: Gauss-Markov Assumptions, given in exam

Effects of Assumptions: 1-6: $\hat{\beta}_j \sim N(\beta_j, \sigma^2_{\hat{\beta}_j})$
1-4: $E[\hat{\beta}_j | X] = \beta_j$   1-5: $Var[\hat{\beta}_j | X] = \frac{\sigma^2}{SST_j (1-R_j^2)} \Rightarrow$ BLUE
Asymptotic properties: 1-4 $\text{plim } \hat{\beta}_j = \beta_j$   1-5: $\hat{\beta}_j \overset{a}{\sim} N(\beta_j, \frac{\sigma^2}{n \hat{\sigma}_j^2})$
Threats to intl. val.: Endogeneity (OVB, ME, simult.), HC (efficiency), multicollinearity

## ③ Heteroskedasticity: $Var(u_i | X) = \sigma_i^2 \Rightarrow$ Wrong standard errors

$\Rightarrow$ No longer BLUE (efficiency) & unbiasedness, consistency remains
Testing for Heteroskedasticity ① Breusch-Pagan-Test ② White Test
i) $\hat{u}^2 = \delta_0 + \hat{\delta}_i x_i + v_i \Rightarrow H_0: \delta_i = 0 \ \forall i \geq 1$   $H_1: \exists i$ such that $\delta_i \neq 0$
Test statistic: $nR_u^2 \sim \chi_k^2$ or $F = \frac{R_u^2/k}{(1-R_u^2)/(n-k-1)}$   $F_{k, n-k-1}$
ii) Use non-linear covariates/regressors: simplified:
$\hat{u}^2 = \delta_0 + \delta_1 \hat{y} + \delta_2 \hat{y}^2 + r_i \Rightarrow H_0: \delta_1 = \delta_2 = 0$   $F_{2, n-k-1}$
Test: $F$-test (see above)

Dealing with Heteroskedasticity: i) WLS/FGLS ii) robust SEs
i) WLS: $Var(u|X) = \sigma^2 h(x_i) \Rightarrow$ divide everything by $\sqrt{h}$ s.t. $h > 0$
ii) FGLS: you don't know the contributing factor $\Rightarrow$ estimate with
$\Rightarrow \ln \hat{u}_i^2 = \delta_0 + \sum \delta_j x_{ij} + r_i \Rightarrow \hat{h}_i = \exp(\hat{\delta}_0 + \sum \hat{\delta}_j x_{ij})$
$\Rightarrow$ follow same procedure with $\sqrt{\hat{h}_i}$ in WLS.
ii) Robust SE: $Var(\hat{\beta}_j) = \frac{\sum r_{ij}^2 \hat{u}_i^2}{SSR_j^2} \Rightarrow$ still inefficient though
HC & LPM: Assume $y \in [0,1]$ and $y_i = \alpha + \beta x + u$
$\Rightarrow Var(y_i | x_i) = p(x_i)[1 - p(x_i)]$  values close to 0.5 large error, close to 0 or 1 small error

## ④ Specification & Data Issues: Misspecification using RESET

Include higher order terms of dependent variable & test:
$y = \beta x + \delta_1 \hat{y}^2 + \delta_2 \hat{y}^3 \Rightarrow F = \frac{R_u^2 - R_R^2}{q} / \frac{1-R_u^2}{n-k-q-1}$   $F_{q, n-k-q-1}$

---

## Nested vs. Non-Nested Models: 1) $\beta x_1$  & $\beta_1 x_1 + \beta_2 x_2$  2) $\beta x_1$ & $\beta \ln x_i$

Testing Non-Nested Models:
CMA: model$_1$ + model$_2$  Test model$_1$: coefficients of model$_2 = 0$ and Use F-test vica versa
Davidson-Mackinnon Test:
$\rightarrow$ T-test: $y = $ model$_1$ + fitted model$_2$  $y = $ model$_2$ + fitted model$_1$
$\rightarrow$ There may be no clear winner (use $R^2$ as a decision)
$\rightarrow$ Rejecting does not validate the other.

Proxy variables (unobservables): $y = \beta_0 + \beta_1 x_1 + \beta_2 x_2 + \beta_3 x_3^* + u_x$
Conditions of a good proxy:  $x_3^* = \delta_0 + \delta_3 x_3 + v_3$
$p(u, x_i) = 0 \ \forall i$   $P(u, x_3^*) = 0$   $P(v_3, x_i) = 0 \ \forall i$

Measurement Errors: CMEA: $E(e_0) = 0$, $P(y, e_0) = 0$
$\rightarrow$ Dependent variable: Remains unbiased & consistent but inefficient
$\rightarrow$ Independent variable: Attenuation bias: $\hat{\beta}_1 \cdot \frac{Var(x_1^*)}{Var(x_1^*) + Var(e_1)}$
$\rightarrow$ biased toward zero

Missing Data & Sample Selection:   $y$ nx omitted
Exogenous: Based on $X \Rightarrow$ OLS remains consistent if equal for all observation
Endogenous: Based on $Y \Rightarrow$ OLS biased $E(u|X) \neq 0$
LAD: $\min \sum |y_i - x_i' \beta|$  $\Rightarrow$ deal with o $y + \Delta i y$, centred around median

## ⑤ Simple Panel Models: Pooled Cross Sections combines multiple

cross-sectional datasets sampled at different points in time
⊕ more observations for precise estimation, allows testing for effects over time
DiD estimation: $y_{it} = \beta_0 + \beta_0 d2_t + \beta_1 D_i + \delta_1 (D_i \cdot d2_t) + u_{it}$
post treatment indicator
Assumptions: Parallel trends   treatment effect   $ATT = (y_{1,post} - y_{1,pre}) - (y_{0,post} - y_{0,pre})$

Fixed effect: $y_{it} = \beta_0 + \beta_1 x_{it} + a_i + u_{it}$   $a_i$: time-invariant ind effect
First differences: $\Delta y_{it} = y_{it} - y_{it-1} = \beta_1 \Delta x_{it} + \Delta u_{it}$  consistent if $p(\Delta x_{it}, \Delta u_{it}) = 0$

Balanced vs. Unbalanced Panel: same number of obs. across time / varying number across individuals
Estimation in First Differences: Conditions for consistency: strict
exogeneity $E[u_{it} | x_{iT}] = 0 \ \forall t$  no residual autocorrelation.
$\rightarrow$ Cannot estimate time-invariant variables, variation in $\Delta x_{it}$
$\rightarrow$ address by interaction term: $\Delta y_{it} = \beta_1 \Delta x_{it} + \gamma_t (Time_t \times a_i)$

Serial Autocorrelation: Regress $\hat{u}_{it}$ on their lags, use FGLS
if significant  Advanced Panel Data Methods: Fixed Effect Estimation
i) within transformation: $y_{it} - \bar{y}_i = \beta_1 (x_{it} - \bar{x}_i) + (u_{it} - \bar{u}_i)$
ii) LSDV: $y_{it} = \beta_0 + \beta x_{it} + \sum \gamma_i D_i + u_{it}$  $\rightarrow$ removes time-invariant heterogeneity
Random Effects Estimator: $y_{it} = \beta_1 x_{it} + a_i + u_{it}$ with $p(x_{it}, a_i) = 0$
FGLS transformation: $y_{it} - \theta \bar{y}_i = \beta_1 (x_{it} - \theta \bar{x}_i) + (1-\theta) \beta_0 + (u_{it} - \theta \bar{u}_i)$
$\theta = 1 - \sqrt{\frac{\sigma_u^2}{\sigma_u^2 + T \sigma_a^2}}$   Model comparison: Hausman-Test for FE & RE

---

Choosing between FE & FD: FE preferred if $u_{it}$ serially auto-correlated.
FD preferred if $u_{it}$ follows a random walk

## ⑦ Instrumental variable estimation: Model $y = \beta x + u$ exhibits

$p(x, u) \neq 0 \Rightarrow$ Instrumental variable: $cov(z, x) \neq 0$ (relevance),
$cov(z, u) = 0$ (Exogeneity) $\Rightarrow \hat{\beta}_1^{IV} = \frac{cov(z, y)}{cov(z, x)}$  Wald Estimator
Two Stage Least Squares
Multiple Instruments: First stage: $\hat{\beta}_1^{IV} = \frac{\bar{y}_{z=1} - \bar{y}_{z=0}}{\bar{x}_{z=1} - \bar{x}_{z=0}}$
$x = \pi_0 + \pi_1 z_1 + \ldots + \pi_m z_m + v$  all exogenous & instrumental
Second stage: $y = \beta_0 + \beta_1 \hat{x} + $ other exog. vars $+ e$  variables
Variance and Inference:
$Var(\hat{\beta}_1^{IV}) = \frac{\sigma^2}{n Var(x) \rho[z, x]^2}$   $\sigma^2 = Var(u)$

Weak Instruments: if $\rho(z, x)$ small $\Rightarrow$ large SE
$R^2$ is different as OLS and can be negative.

## ⑧ Regression Discontinuity Design:

Sharp RD: deterministic  $D_i = \{1$ if $x_i > x_0$ cutoff, 0 otherwise
Model: $Y_i = \alpha + \beta x_i + \rho D_i + \eta_i$, where $\rho = Y_{1i} - Y_{0i}$
Identification: $\rho = \lim_{x \to x_0} E[Y_i | x_i = x] - \lim_{x \to x_0} E[Y_i | x_i = x]$
Fuzzy RD: $D_i = \gamma_0 + \gamma_1 I(x_i \geq x_0) + g(x_i) + r_i$ (probabilistic)
Wald Estimator: $\rho = \frac{\text{Reduced Form Discontinuity}}{\text{First Stage Discontinuity}} = \frac{E[Y | x_0^+] - E[Y | x_0^-]}{E[D | x_0^+] - E[D | x_0^-]}$
2SLS Specification: $\alpha + \beta x_i + \rho D_i + \eta_i$ (with $D_i$ instrumented by $I(x_i \geq x_0)$)
Estimation Approaches: Parametric: polynomial approximation
$Y_i = \alpha + Z^j \beta_k (x_i - x_0)^k + \rho D_i + \varepsilon_i$  Non-parametric: local linear reg:
$\hat{\rho} = \arg\min \sum K(\frac{x_i - x_0}{\Delta})(Y_i - \alpha - \beta(x_i - x_0) - \rho D_i)^2$   kernel function
Validity checks: Continuity of covariates check: Pre-treatment
covariates should not jump at $x_0$. Placebo Test: no discontinuity in outcome at false cutoffs

## ⑨ Limited Dependent Variable Models & Sample Selection Correction

LPM: $P(y=1|x) = \beta_0 + x\beta$  Logit Model: $P(y=1|x) = \Lambda(x\beta) = \frac{\exp(x\beta)}{1 + \exp(x\beta)}$
Probit Model: $P(y=1|x) = \Phi(x\beta) = \int \phi(t) dt$
Latent Variable Framework: $y^* = x\beta + \varepsilon$   $y = I(y^* > 0)$
where $\varepsilon \sim$ Logistic (logit) or $\varepsilon \sim N(0,1)$ (probit)
Partial Effects: $\frac{\partial P(y=1|x)}{\partial x_j} = g(x\beta)\beta_j$, where $g(\cdot)$ PDF
(logistic/normal). For discrete $x_j$, compute: $\Delta P = G(x\beta + \beta_j) - G(x\beta)$  MLE: $\ell(\beta) = \sum [y_i \log G(x_i\beta) + (1-y_i) \log(1 - G(x_i\beta))]$  Other Models: Poisson: $\lambda_i = \exp(x_i\beta) \Rightarrow P(y_i = k) = \frac{e^{-\lambda_i} \lambda_i^k}{k!}$  Assumption & Tests: Logit/Probit $\Rightarrow$ correct distribution