

Goodness of Fit: $R^2 = \text{SSE}/\text{SST}$. For SLR, $R^2 = (r)^2$ | **DF for sample variance and MLR (equations):** We lose 1DF for every x incl. β_0 . For sample variance of simple regression, square root of $\text{SSR}/(n-2)$. For sample variance in MLR, variance of error/ $\text{SST}(1-R^2)$, $\text{SST} = \sum(x_j - \text{avg}(x))^2$

Figure 1- Transformations (Level-log etc.)

Interaction term: Using interaction terms, we can quantify how multiple features work together to influence the price. This helps us identify which combinations of features drive the most value – valuable information for crafting pricing strategies and promotions. $Y = \beta_0 + \beta_1 x_1 + \beta_2 x_2 + \beta_3(x_1 * x_2) + \epsilon$, If $\beta_3 > 0$, it indicates a **synergistic effect**: x_1 and x_2 reinforce each other's impact on Y. If $\beta_3 < 0$, it indicates a **dampening effect**: x_1 and x_2 weaken each other's impact on Y. If $\beta_3 = 0$, it means there is no interaction between x_1 and x_2 , and their effects are additive.

Model	Dependent Variable	Independent Variable	Interpretation of β_1
Level-level	y	x	$\Delta y = \beta_1 \Delta x$
Level-log	y	$\log(x)$	$\Delta y = (\beta_1/100)\% \Delta x$
Log-level	$\log(y)$	x	$\% \Delta y = (100\beta_1) \Delta x$
Log-log	$\log(y)$	$\log(x)$	$\% \Delta y = \beta_1 \% \Delta x$

Changing units: Changing the units of measurement of one or more variables does not change the regression itself. It changes only the magnitude of the coefficients. If model has log in it, no changes as it is %. SE changes tho!

Calculations: Predictions: Plug in b's . Quadratic terms (finding turning point): $Y^* = \beta_0 + \beta_1 X + \beta_2 X^2$, turning point: $X^* = -(\beta_1/2\beta_2)$. **Changing units in log-level model:** $\ln(Y) = \beta_0 + \beta_1 X + \epsilon$, a 1-unit change in X leads to a $\beta_1 \times 100\%$ change in Y. **Residual calculation:** $e_i = Y_i - Y^*$

[**Bias:** MLR1, MLR2, MLR3, MLR4. | **Inference:** MLR5, MLR6 (in addition to MLR1–MLR4)]

Assumption	Cause	Diagnose	Handle	R Code
MLR1: Linearity. The model is linear in parameters. (<i>Linearity = residuals have mean zero for every level of the fitted values and predictors</i>)	(1) Incorrect model specification, (2) missing interaction terms or (3) non-linear relationships	(1) Plot residuals vs. fitted values- should show no pattern. (2) Use scatterplots of independent vs. dependent variables. (3) RESET test (F-test compares long model to smaller models & adds powers of the fitted values (e.g., squared, cubed) to check for functional form misspecification. A significant result indicates nonlinearity. "Ho: The model is well-specified."	(1) Check other assumptions (2) Return to data processing (3) Log transformations / quadratic terms (4) Interaction terms (5) Alternative models	(1) plot(model_name, which=1) (2) resettest(model_name, power = 2:3, type = "fitted")
MLR2: Random Sample. The data is randomly sampled from the population, ensuring that the sample represents the population.	(1) Non-representative or biased sample (2) Sampling errors or convenience sampling, (3) Missing values, outliers, non-random samples Outliers: Types: (1) Extreme values (far outside typical range, often in Y variable), means high residuals, (2) Leverage points (unusual X values). (3) Influential points (observations significantly affecting model estimates (predicted values & parameter estimates) Extreme + Leverage	(1) Inspect data collection methods, (2) Check descriptive statistics for sample characteristics. (3) N/A's at random (no bias but, lose data, ↑ uncertainty, means 1 SE), (4) N/A's NOT at random means bias (endogenous). <i>"In small samples, outliers can have a large effect on model estimates, distorting predictions and coefficient values. Impact of outliers diminishes as N increases."</i> (5) Studentized residuals (standardizing residuals) to identify outliers. Or external res stud (w/o high leverage observation	(1) Over/under representation: Weighting or post stratification (2) N/A's: (a) Listwise delete (stats power down), (b) mean/median imputation (↓ variance->bias), (c) regression imputation (when NAs not random, overstates confidence in predictions), (d) flagging NA's (combines imputation + binary variable when NA's carry value, may introduce bias), (e) Missing category (for categorical variables -> unknown/undisclosed). (3) Outliers: (a) cap extreme values (winsorizing), (3) add, drop, transform variables, (4) fix data entry issues, (5) remove observations outside intended population.	(1) dataset %>% mutate(res_stud = rstudent(m1), res_stud_large = ifelse(between(res_stud,-3, 3), "Normal", "Extreme")) (2) Leverage: Cooks distance (residuals vs leverage) : plot(m1, 5)
MLR3: No perfect collinearity. The independent variables are not perfectly correlated, so the model can estimate coefficients uniquely.	(1) If independent variables are perfectly correlated, (2) if N is too small in relation to parameters, we need at least k + 1 observations to estimate k + 1, (3) dummy variable trap (e.g., including all categories of a categorical variable). <i>"High multicollinearity reduces the reliability of individual predictors, making it hard to assess their true relationship with the outcome variable."</i>	(1) Calculate Variance Inflation Factor (VIF). High VIF (e.g., >10) signals potential multicollinearity. Measures how much variance of a regression coefficient is inflated due to multicollinearity with other predictors. (2) Look at correlation matrices.	(1) Combine highly correlated variables / drop one, (2) increase sample size to help disentangle variable effects, (3) A high degree of linear relationship between x1 and x2 can lead to large variances for the OLS estimators. <i>"What ultimately matters is how big "β_j is in relation to its standard error."</i>	(1) numeric_vars <- data %>% select_if(is.numeric) correlation_matrix <- round(cor(numeric_vars, use = "complete.obs"), 2) print(correlation_matrix) (2) vif(model_name)
MLR4: Zero conditional mean. The error term has an expected value of zero conditional on the independent variables. This is critical for unbiasedness, as it implies no omitted variable bias or endogeneity. <i>Biggest threat to causal modelling: But time periods do not have to be adjacent in panel data</i>	(1) Omitted variable bias, (2) Measurement errors (a) in Y, if nonrandom, different slope, (b) in X, worse(different slope either way!), (3) Reverse causality (simultaneity) -> if x and y are simultaneously correlated are correlated with error term. <i>"Correlation between 1 'x' and 'e' results in ALL estimators being biased. Overspecifying: not unbiased, but variance of estimators increases."</i>	(1) Perform a residual analysis, (2) use theoretical reasoning to identify potential omitted variables, (3) CEV assumption: e is unrelated to true x, $\text{Cov}(x^*, ME) = 0$. CEV deals with random ME in X.	(1) Only think and explain if. + or - bias, (2) determine if CEV assumption holds, then attenuation bias: β_{xj} is pushed towards 0, we underestimate magnitude. $\text{Cov}(x, u - \beta_1 e) = -\beta_1 \text{Cov}(x, e) = -\beta_1 \sigma_e^2$: Covariance between x and composite error.	(1) plot(model, which = 1) (2) numeric_vars <- data_set[sapply(data_set, is.numeric)] # Select only numeric columns cor_matrix <- cor(numeric_vars, use = "complete.obs") # Compute correlation matrix print(cor_matrix)

MLR5: Homoskedasticity. The variance of the error term is constant across all levels of the independent variables. This ensures valid standard errors for hypothesis testing and confidence intervals. If it breaks: <i>OLS ≠ BLUE(best linear unbiased estimator)</i> , large N does not help.	(1) Presence of heteroskedasticity (non-constant error variance) -does not affect the unbiasness or consistency of OLS estimators, nor the interpretation of R ² . Hypothesis testing does not hold. Variance of estimators becomes biased. SE and CI are not reliable, invalid t and f stats, (2) Outliers or data skewness.	(1) Breusch-Pagan (BP) or White test, (2) plot residuals vs. fitted values (look for patterns) and (3) scale-location plot (horizontal red line indicates homoscedasticity) <i>"Fan or funnel shape suggests violation of MLR5, curved patterns or distant points indicate issues with model fit or variance consistency."</i>	(1) Robust standard errors, (2) transformations to reduce variance (3) remove Outlier / cap a predictor	(1) BP test or white test (for linear and nonlinear): bptest(model, ~fitted(model) + I(fitted(model)^2)) (2) Robust SEs: cov1 <- vcovHC(model, type = "HC3") robust_se <- sqrt(diag(cov1))
MLR6: Error term is normally distributed. The error term follows a normal distribution. Important for small samples to ensure that hypothesis testing, and confidence intervals are valid (by the Central Limit Theorem, normality is less critical for large samples).	(1) Small sample size. (2) Extreme outliers. (3) Skewed or non-normal error distribution.	(1) Plot a Q-Q plot of residuals: compares distribution of the residuals from regression model to normal distribution. Points are roughly along 45-degree line. Curves suggest skewness. S shaped mean heavier or lighter tails than Dnorm, (3) Perform a Shapiro-Wilk or Kolmogorov-Smirnov test. (if residuals deviate from Dnorm	(1) transformation, (2) larger sample	shapiro.test(rstandard(model)) H ₀ = Residuals follow a DNorm. *test is appropriate for n <= 5000*

Pooled Cross Section two or more random samples from a population of interest at different points in time. Because samples are collected independently of each other and randomly they need not be of equal size, and they will typically contain different statistical units. **Pooled Time Series** multiple statistical unit of analysis N= {1,...,n} observed in several successive periods of time T= {1,...,t}(e.g. the price of several stocks over time). This differs from panel data in that there are more time periods than statistical units (t>n).

Aspect	Pooled Time Series (Panel)	Pooled Cross Section
Data Structure	Same entities tracked over time	Different entities sampled at each time point.
Time Dependency	Time dependency within entities	No dependency between time points
Example Analysis	Panel data models, fixed/random effects	Difference in means tests, trends in averages
Use Case	Studying individual level trends	Studying population level changes

DIFF in DIFF: A treatment group affected by an intervention, a control group unaffected by the intervention, pre-and post-intervention data for both groups.

Strengths: Provides **intuitive interpretation** of causal effects. Allows use of **observational data** if assumptions (e.g., parallel trends) hold. Works with both **individual- and group-level data**. Focuses on **changes over time**, not absolute levels, accounting for pre-existing differences. Adjusts for **time-varying factors** unrelated to the intervention. **Limitations:** Requires **baseline data** and a valid **control group**. Treatment allocation depends on the **baseline outcome**. Groups exhibit **different outcome trends** over time. Group **composition changes** pre/post intervention. **Note:** Units do not self-select into treated/not treated groups. Parallel trends assumptions: In the absence of treatment, the difference in outcomes between the treatment and control groups would remain constant over time. **Unbalanced data** is no problem, if it is not systematic (MLR2). **Idiosyncratic errors** (v_{it}) are individual-specific errors in a regression model that vary randomly across observations and over time, capturing factors not explained by the model. **Heterogeneity bias** -> correlation between unobserved effects and explanatory variables.

FD: One period of data drops automatically. For T=2, FD reduces to simple cross section. When panel is balanced and T=2, FD = FE. For T>2, FE for long term effects, FD for short term impact on changes. $y_{it} = \beta_0 + \beta_1 x_{it} + a_i + v_{it}$ -> **unobserved effects (a_i) model**. FD: $(y_{i(t-1)} - y_{it}) = \Delta y_{it} = \beta_1 \Delta x_{it} + \Delta v_{it}$, and we assume $E[\Delta v_{it} | \Delta x_{it}] = 0$. Thus, degrees of freedom are reduced by 1 and equal to N (T - 1) - k.

FE: Approaches: (1) Demean the data (subtract mean of each variable for each city from its observations), (2) Include dummy variable for each city (city specific effects constant over time). $\Delta y_{it} = \beta_1 \Delta x_{it} + \Delta v_{it}$ (there is a ~ on top of each to indicate demean). FE estimators require a correction to the calculation of SEs. Because we demean, degrees of freedom are N(T-1) - k.

Serial Correlation occurs when e_{it} and $e_{i(t-1)}$ are correlated. Positive serial correlation is a problem, may underestimate SE estimates (negative correlation less problematic). In R, we use pgbttest(pmodel, h0="fd/fe", type=HC3) and pfdtest(pmodel). If p-value <0.05, serial correlation is present H₀=no serial correlation. If idiosyncratic errors in levels are uncorrelated (FE hypothesis), first-differenced (FD) errors show serial correlation of 0.5. *If errors in levels follow a random walk (FD hypothesis), FD errors have no serial correlation.* Both FE and FD are consistent, but FE is more efficient under "fe" and FD under "fd". **Summary:** FD & FE & dummy inclusion control for time-constant unobserved factors, however no control over omitted variables varying over entities and time. Cannot include time-invariant factors (e.g. age), differencing will result in 0, but interaction possible. Missing data: FE handles seamlessly, FD loses two observations for each missing period, lowers sample size). Attrition occurs when units leave panel (they die, companies liquidate), if related to variables in model, leads to sample selection bias. FE can handle attrition correlated with a_i (time-constant error) reducing impact.

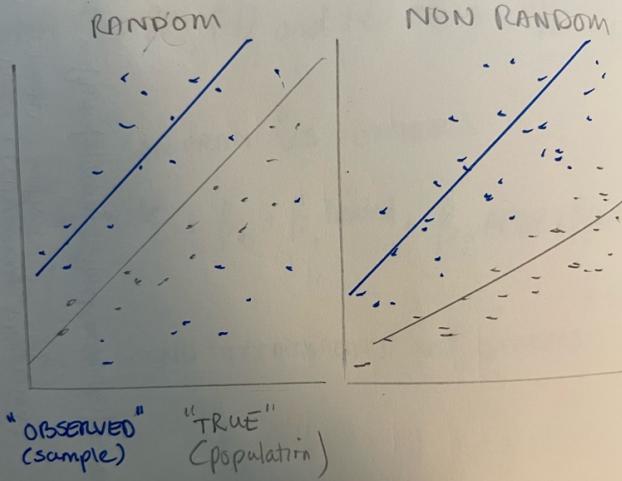
Linear Probability Model (LPM) is so named because probability is linear in the parameters. $P(\text{approve}=1|x) = \beta_0 + \beta_1 \text{white}$. Where β_1 is the average, all else constant increase/decrease probability of success added to β_0 . **Advantage:** coefficients represent change in probability for unit change in predictor, easy computation thanks to OLS, comparable to logistic regression when predictors are categorical. **Disadvantage:** predicted probabilities can fall outside of [0,1] range (nonsensical). Error variance depends on predicted probability, violating MLR5. Errors are not normally distributed (MLR6 violated). MLR5 -> Robust errors. MLR6 -> In LPMs, residuals are inherently not normal because y is bound to 0 and 1. Not critical if sample size is large, but for small sample inference it may be problematic. Robust SEs recommended. R² assumes continuous outcomes, so misleads interpretation. Instead include percent correctly predicted or naive benchmark (predicting majority class). Interpretation: 4% to 6% is a 2% percentage point change but a 50% percent change.

Non-linear (logistic & standard normal cumulative distribution): $P(y=1|x) = G(\beta_0 + \beta_1 x_1 + \dots + \beta_k x_k) = G(\beta_0 + x\beta)$, where G is a cumulative distribution function. The difference between probit and logit models is small in terms of estimates and structure. Preferred over LPM because (1) Constrain predicted probabilities to lie between 0 and 1, (2) Ensure marginal effects are sensible across the range of explanatory variables. LPM provides simple and reasonable estimates and effective for statistical significance, however unsuitable for predicting individual probabilities.

CHEAT SHEET (MLR4)

2.12.2024
[FSAR]
check

Measurement Error in Y



FORMULA

$$y = \beta_0 + \beta_1 x_1 + \dots + \beta_n x_n + u + e_0$$

measurement error (me)

(true y)
 $y^* = y - e_0$

Measurement Error in X

