# Wrap up Final Lecture: Part I

Nayel Bettache

NI.		D.		***	-1	
INd	/er	D	ະເ	ιa	cr	re

2

<ロト <回ト < 回ト < 回ト -

Linear Model: Assumptions and Interpretations Model:

 $Y = X\beta + \varepsilon$  where  $\varepsilon$  represents random errors.

## Key Assumptions:

- Linearity: The expected value of Y is a linear function of X.
- **Independence:** The errors  $\varepsilon_i$  are independent across observations.
- Homoscedasticity: The errors have constant variance: Var(ε<sub>i</sub>) = σ<sup>2</sup> for all *i*.
- **Normality:** The errors  $\varepsilon_i$  are normally distributed (required mainly for valid *inference*, not for point estimates).
- Low Multicollinearity: Predictors should not be nearly perfectly correlated (to ensure stable estimation of  $\beta$ ).

#### Notes:

- Normality is *not required* for estimating β, but important for hypothesis tests.
- Multicollinearity does not bias estimates, but increases their variance.

Nayel Bettache

# Consequences of Assumption Violations

**Reminder:** Under the Gauss-Markov assumptions, OLS provides the **Best** Linear Unbiased Estimator (BLUE).

If assumptions are violated:

## • Linearity violation:

 $\rightarrow$  Model is misspecified  $\Rightarrow$  Biased estimates, poor predictions.

## Independence violation:

 $\rightarrow$  Standard errors are underestimated  $\Rightarrow$  Inference becomes invalid (e.g., misleading p-values).

## • Homoscedasticity violation:

 $\rightarrow$  Estimates remain unbiased but are inefficient  $\Rightarrow$  Larger standard errors, invalid usual inference (need robust methods).

## Normality violation:

 $\rightarrow$  OLS estimates are still unbiased and consistent, but small-sample inference (e.g., t-tests, F-tests) may be invalid.

## Multicollinearity:

 $\rightarrow$  Estimates remain unbiased, but coefficients are highly sensitive to data perturbations  $\Rightarrow$  Inflated standard errors, unstable predictions.

# Checking Model Assumptions

## Visual Diagnostics:

## • Residual plots:

Check for patterns  $\Rightarrow$  Linearity and homoscedasticity.

## • Q-Q plots (Quantile-Quantile plots):

Check if residuals align with a normal distribution.

## Statistical Tests:

#### • Durbin-Watson test:

Detect autocorrelation in residuals (independence violation).

### • Variance Inflation Factor (VIF):

Quantify multicollinearity among predictors.

**Note:** Always combine visual inspection and statistical tests for robust diagnostics.

く 目 ト く ヨ ト く ヨ ト

# Multiple Linear Regression: Fundamentals Model:

$$Y_i = \beta_0 + \beta_1 X_{i1} + \beta_2 X_{i2} + \dots + \beta_p X_{ip} + \varepsilon_i$$

where  $\varepsilon_i$  are random errors with mean 0 and constant variance  $\sigma^2$ .

Matrix Form:

$$\mathbf{Y} = \mathbf{X} \boldsymbol{eta} + \boldsymbol{arepsilon}$$

#### Least Squares Estimation:

$$\hat{\boldsymbol{\beta}} = (\mathbf{X}^{\mathsf{T}} \mathbf{X})^{-1} \mathbf{X}^{\mathsf{T}} \mathbf{Y}$$

(provided that  $\mathbf{X}^{\mathsf{T}}\mathbf{X}$  is invertible)

Fitted Values:

$$\hat{\mathbf{Y}} = \mathbf{X}\hat{\boldsymbol{\beta}}$$

**Objective:** Minimize the residual sum of squares (RSS):  $\sum_{i=1}^{n} (Y_i - \hat{Y}_i)^2$ .

# Multiple Linear Regression: Interpretation

• **Coefficient**  $\beta_j$ :

Expected change in Y for a one-unit increase in  $X_j$ , holding all other predictors constant.

## • Coefficient of Determination (*R*<sup>2</sup>):

Proportion of variance in Y explained by the model.

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

• Adjusted R<sup>2</sup>:

Adjusts  $R^2$  for the number of predictors. Penalizes overfitting.

## • F-test:

Tests the joint hypothesis  $H_0: \beta_1 = \beta_2 = \cdots = \beta_p = 0.$ 

#### • t-tests:

Test if individual  $\beta_i$  significantly differs from zero.

# Categorical Variables and Interactions

Handling Categorical Variables:

• Dummy variables:

$$X_{ij} = \begin{cases} 1 & \text{if observation } i \text{ belongs to category } j \\ 0 & \text{otherwise} \end{cases}$$

• Reference category:

One category is omitted to avoid the "dummy variable trap" (perfect multicollinearity).

## **Modeling Interactions:**

- Include products of predictors to allow non-additive effects.
- Interaction model:

$$Y = \beta_0 + \beta_1 X_1 + \beta_2 X_2 + \beta_3 (X_1 \times X_2) + \varepsilon$$

## • Interpretation of $\beta_3$ :

 $\beta_3$  measures how the effect of  $X_1$  on Y changes depending on the level of  $X_2$ .

# Sampling Distributions: Key Concepts

## Definition:

• The **sampling distribution** of a statistic is the probability distribution of that statistic across all possible random samples from the population.

#### **Examples:**

- The sample mean  $\bar{X}$  has a normal distribution (by the Central Limit Theorem) when n is large.
- The *t*-statistic follows a *t*-distribution under the null hypothesis (for small samples).

#### Importance:

- Fundamental for building **confidence intervals** and conducting **hypothesis tests**.
- Explains the variability of estimators from sample to sample.

Key Idea:

• Even though a parameter (like  $\mu$ ) is fixed, the statistic (like  $\bar{X}$ ) is random across samples.

Nayel Bettache

# Central Limit Theorem (CLT)

#### Statement:

- Let  $X_1, X_2, ..., X_n$  be i.i.d. random variables with mean  $\mu$  and variance  $\sigma^2 < \infty$ .
- Then, as  $n \to \infty$ :

$$\sqrt{n}\left(\bar{X}-\mu\right) \stackrel{d}{\longrightarrow} N(0,\sigma^2)$$

• Equivalently:

$$\bar{X} \stackrel{d}{\longrightarrow} N\left(\mu, \frac{\sigma^2}{n}\right)$$

## Implications:

- The sampling distribution of  $\bar{X}$  becomes approximately normal, regardless of the original distribution.
- Justifies normal-based inference (confidence intervals, hypothesis tests) for large *n*.

#### Note:

• Convergence can be slower for heavy-tailed or skewed distributions.

# Important Sampling Distributions

When population variance  $\sigma^2$  is known (or *n* large):

• Sample mean:

$$\bar{X} \sim N\left(\mu, \frac{\sigma^2}{n}\right)$$

• Sample proportion:

$$\hat{p} \sim N\left(p, \frac{p(1-p)}{n}\right)$$

(valid for large *n* by the Central Limit Theorem)

• Difference of two sample means:

$$ar{X}_1 - ar{X}_2 \sim N\left(\mu_1 - \mu_2, rac{\sigma_1^2}{n_1} + rac{\sigma_2^2}{n_2}
ight)$$

• • = • • = •

# Important Sampling Distributions II

## When population variance $\sigma^2$ is unknown:

• Student's *t*-distribution:

If  $\sigma$  is unknown and replaced by the sample standard deviation s, then

$$rac{ar{X}-\mu}{s/\sqrt{n}}\sim t_{n-1}$$

where  $t_{n-1}$  is a *t*-distribution with n-1 degrees of freedom.

#### For variance-related statistics:

• Chi-square distribution:

$$\frac{(n-1)s^2}{\sigma^2} \sim \chi^2_{n-1}$$

(used in confidence intervals and tests for variance)

• F-distribution:

$$rac{\mathsf{Variance}_1/\sigma_1^2}{\mathsf{Variance}_2/\sigma_2^2} \sim \mathcal{F}_{d_1,d_2}$$

(used for comparing two variances)

# Confidence Intervals: Fundamentals General form:

Point estimate  $\pm$  (Critical value  $\times$  Standard error)

#### **Correct Interpretation:**

• If we repeated the sampling procedure infinitely many times, approximately 95% of the constructed confidence intervals would contain the true population parameter.

#### Precision vs Confidence:

- Wider intervals: Higher confidence, but lower precision.
- Narrower intervals: Higher precision, but lower confidence.
- Increasing the sample size *n* leads to narrower intervals without sacrificing confidence.

## Common Misinterpretation (WRONG):

- It is **incorrect** to say that there is a 95% probability that the true parameter lies inside a *realized* interval.
- Once the interval is calculated, the parameter is either inside it or not.

## Common Confidence Intervals Population Mean (Known Variance $\sigma^2$ ):

$$\bar{x} \pm z_{\alpha/2} \frac{\sigma}{\sqrt{n}}$$

(when  $\sigma$  known,  $z_{\alpha/2}$  from standard normal distribution)

Population Mean (Unknown Variance):

$$ar{x} \pm t_{lpha/2, n-1} rac{s}{\sqrt{n}}$$

(when  $\sigma$  unknown, s sample standard deviation, t-distribution with n-1 degrees of freedom)

**Population Proportion:** 

$$\hat{p} \pm z_{\alpha/2} \sqrt{\frac{\hat{p}(1-\hat{p})}{n}}$$

(for large n, by Central Limit Theorem)

# Confidence Intervals for Regression Coefficients

Confidence interval for each regression coefficient  $\beta_i$ :

$$\hat{\beta}_j \pm t_{\alpha/2, n-p-1} \times SE(\hat{\beta}_j)$$

where:

• 
$$SE(\hat{\beta}_j) = \sqrt{\hat{\sigma}^2 \left[ (\mathbf{X}^T \mathbf{X})^{-1} \right]_{jj}}$$
  
•  $\hat{\sigma}^2 = \frac{1}{n-p-1} \sum_{i=1}^n (Y_i - \hat{Y}_i)^2$  (residual variance estimate)

• Degrees of freedom: n - p - 1 (*n*: observations, *p*: predictors)

#### Notes:

- *t*-distribution used because  $\sigma^2$  is unknown.
- $[(\mathbf{X}^T \mathbf{X})^{-1}]_{jj}$  is the *j*-th diagonal element of the inverse Gram matrix, capturing the variability of  $\hat{\beta}_j$ .

# Hypothesis Testing Framework

Key Elements:

- Null Hypothesis (*H*<sub>0</sub>): Statement of no effect, no difference, or status quo. (Assumed true at the start.)
- Alternative Hypothesis ( $H_1$  or  $H_A$ ): Statement representing an effect, difference, or deviation from  $H_0$ .
- **Test Statistic:** Quantity computed from sample data, with a known distribution under  $H_0$ .
- *p*-value: Probability, assuming  $H_0$  is true, of observing a test statistic as extreme or more extreme than the one observed.
- Significance Level (α): Pre-specified threshold (commonly 0.05) for deciding whether to reject H<sub>0</sub>.

Decision Rule:

- Reject  $H_0$  if *p*-value  $\leq \alpha$ .
- Fail to reject  $H_0$  if *p*-value  $> \alpha$ .

イロト 不得 トイヨト イヨト 二日

# Types of Errors and Power

	H <sub>0</sub> True	$H_0$ False
Reject $H_0$	Type I error ( $lpha$ )	Correct decision (Power)
Fail to reject $H_0$	Correct decision	Type II error ( $\beta$ )

#### Key Concepts:

- Type I error (α): Rejecting H<sub>0</sub> when it is actually true (false positive).
- Type II error (β): Failing to reject H<sub>0</sub> when it is actually false (false negative).
- **Power**  $(1 \beta)$ : Probability of correctly rejecting a false  $H_0$ .

### Factors that Increase Power:

- Larger sample size (n).
- Larger effect size (true difference from  $H_0$ ).
- Higher significance level ( $\alpha$ ).
- Lower variability (smaller  $\sigma^2$ ).

## Common Hypothesis Tests

**One-Sample** *t*-**Test**:

$$t = rac{ar{x} - \mu_0}{s/\sqrt{n}} ~~ \sim ~~ t_{n-1}$$

Testing  $H_0: \mu = \mu_0$ . Requires approximate normality of the population.

**Two-Sample** *t*-**Test** (Independent Samples):

$$t = rac{ar{x}_1 - ar{x}_2}{\sqrt{rac{s_1^2}{n_1} + rac{s_2^2}{n_2}}} ~\sim~ t_{df}$$

Testing  $H_0: \mu_1 = \mu_2$ . Assumes independence between samples. Equal variance assumption may or may not be made (Welch's test if variances unequal).

# Common Hypothesis Tests II

Paired *t*-Test:

$$t = rac{ar{d} - \mu_d}{s_d / \sqrt{n}} ~~ \sim ~~ t_{n-1}$$

Testing  $H_0: \mu_d = 0$  for the mean difference between paired observations. Assumes differences are approximately normally distributed.

F-Test (One-Way ANOVA):

$${\sf F} = rac{MS_{between}}{MS_{within}} ~~ {\sf F}_{k-1,n-k}$$

Testing  $H_0$ : All group means are equal.

- *k* = number of groups
- *MS*<sub>between</sub> = mean square between groups
- *MS<sub>within</sub>* = mean square within groups (residuals)

Hypothesis Tests in Regression: t-test

#### *t*-Test for Individual Coefficients:

$$t = \frac{\hat{\beta}_j - \beta_{j0}}{SE(\hat{\beta}_j)} \quad \sim \quad t_{n-p-1}$$

Testing:

$$H_0: \beta_j = \beta_{j0}$$
 vs.  $H_1: \beta_j \neq \beta_{j0}$ 

(Typically,  $\beta_{j0} = 0$  to test if a predictor is useful.)

3

(日)

## Hypothesis Tests in Regression: F-test *F*-Test for Overall Model Significance:

$$F = \frac{MSR}{MSE} = \frac{\text{Regression Sum of Squares}/p}{\text{Residual Sum of Squares}/(n - p - 1)}$$

Testing:

$$H_0: \beta_1 = \beta_2 = \cdots = \beta_p = 0$$
 vs.  $H_1:$  At least one  $\beta_j \neq 0$ 

- MSR: Mean Square Regression
- MSE: Mean Square Error (Residual)
- p: Number of predictors

#### Partial *F*-Test:

- Compares two nested models (a simpler model inside a more complex one).
- Tests if a subset of predictors improves the model significantly.

# The Multiple Testing Problem

## Key Issue:

• When performing multiple hypothesis tests (*m* tests), the probability of making at least one Type I error increases.

## Definitions:

- Family-Wise Error Rate (FWER): Probability of making at least one Type I error across all tests.
- False Discovery Rate (FDR): Expected proportion of Type I errors among all rejections.

## Without Correction (Independent Tests):

$$FWER = 1 - (1 - \alpha)^m$$
, where:

- $\bullet \ \alpha = {\rm significance}$  level for each individual test
- *m* = number of independent tests

イロト 不得 トイヨト イヨト

# The Multiple Testing Problem II

#### Example:

- m = 20 tests,  $\alpha = 0.05$  per test
- FWER  $\approx 1-(1-0.05)^{20}\approx 0.64$

## Solutions (briefly):

- **FWER control**: Bonferroni correction, Holm's procedure (very conservative).
- **FDR control**: Benjamini-Hochberg procedure (more powerful, used in large-scale testing).

▲圖 医 ▲ 国 医 ▲ 国 医 …

# Multiple Testing Correction Methods

## Bonferroni Correction (Controls FWER):

 $\alpha_{\rm corrected} = \frac{\alpha}{m}$ 

- Test each hypothesis at level  $\alpha/m$ .
- Very simple, but conservative (especially when *m* is large).

#### Holm-Bonferroni Method (Controls FWER, Less Conservative):

- Order *p*-values:  $p_{(1)} \le p_{(2)} \le \cdots \le p_{(m)}$ .
- Solution For each *i*, compare  $p_{(i)}$  to  $\frac{\alpha}{m-i+1}$ :
  - If  $p_{(i)}$  is significant, reject  $H_{0(i)}$  and continue.
  - Stop at the first non-significant *p*-value.
  - Less conservative than Bonferroni while still controlling FWER.

# Multiple Testing Correction Methods II

## Benjamini-Hochberg Procedure (Controls FDR):

- Order *p*-values:  $p_{(1)} \le p_{(2)} \le \cdots \le p_{(m)}$ .
- Find the largest k such that:

$$p_{(k)} \leq \frac{k}{m} \alpha$$

- **3** Reject all  $H_{0(i)}$  for i = 1, 2, ..., k.
  - Controls the expected proportion of false discoveries among all rejections (FDR).
  - Widely used in large-scale testing (e.g., genomics, machine learning).

# Practical Considerations in Multiple Testing

## Key Trade-Off:

- Stricter corrections (e.g., Bonferroni) reduce false positives (Type I errors) but increase false negatives (Type II errors).
- Balancing error types depends on research priorities.

## Independence Assumption:

- Many multiple testing corrections assume independent or weakly dependent tests.
- Violations (strong correlations) may affect FWER/FDR control.

## **Context Matters:**

- **FWER control**: Prefer when false positives are very costly (e.g., clinical trials, regulatory decisions).
- **FDR control**: Prefer when discovery is more important than strict certainty (e.g., genomics, exploratory research).

## **Reporting Standards:**

- Always report:
  - The correction method used (e.g., Bonferroni, BH procedure).
  - The type of error rate controlled (FWER or FDR)

# Understanding Heteroskedasticity

Definitions:

• Homoskedasticity:

$$Var(\varepsilon_i) = \sigma^2$$
 for all  $i$ 

The variance of the errors is constant across observations.

• Heteroskedasticity:

$$Var(\varepsilon_i) = \sigma_i^2$$
 varies with *i*

The variance of the errors differs across observations.

### Consequences of Heteroskedasticity:

- OLS estimators remain unbiased but are no longer efficient.
- Estimated standard errors are biased  $\rightarrow$  Hypothesis tests and confidence intervals are invalid.
- Predictions have suboptimal (non-minimal) variance.
- Violates a Gauss-Markov assumption: OLS is no longer the **Best** Linear Unbiased Estimator (BLUE).

Nayel Bettache

# Detecting Heteroskedasticity

Visual Methods:

- **Residual vs. Fitted Values Plot:** Look for patterns a funnel shape suggests heteroskedasticity.
- **Residual vs. Predictor Plots:** Examine key predictors individually for changing spread.
- Scale-Location Plot: Plot  $\sqrt{|e_i|}$  versus  $\hat{y}_i$  to better detect non-constant variance.

## Statistical Tests:

## Breusch-Pagan Test:

- Regress squared residuals on predictors.
- *H*<sub>0</sub>: Homoskedasticity (constant variance).
- White Test:
  - General test allowing for nonlinear forms of heteroskedasticity.
  - H<sub>0</sub>: Homoskedasticity.

## Goldfeld-Quandt Test:

- Test for monotonic changes in variance across ordered data.
- Compare variance across different subsamples.

# Addressing Heteroskedasticity

## Transformation Approaches:

- Apply a transformation to the response variable to stabilize variance (e.g., log, square root).
- Box-Cox Transformation:

$$y^{(\lambda)} = \begin{cases} \frac{y^{\lambda} - 1}{\lambda} & \text{if } \lambda \neq 0\\ \log(y) & \text{if } \lambda = 0 \end{cases}$$

where  $\lambda$  is estimated from the data.

## Weighted Least Squares (WLS):

• Minimize a weighted sum of squared residuals:

$$\sum_{i=1}^n w_i (y_i - x_i^T \beta)^2$$

• Use weights  $w_i = 1/\hat{\sigma}_i^2$ , giving more importance to observations with smaller variance.

# Addressing Heteroskedasticity II

#### **Robust Standard Errors:**

- Use heteroskedasticity-consistent (HC) standard errors (e.g., White's estimator).
- Coefficient estimates stay the same; only the standard errors are adjusted for valid inference.

#### Summary:

- If you care about improving model fit  $\rightarrow$  use transformations or WLS.
- If you mainly care about valid inference ightarrow use robust standard errors.

# Bootstrapping: Fundamentals

## Basic Principle:

- Approximate the sampling distribution of a statistic by resampling from the observed data.
- Does not require assumptions about the underlying population distribution.

#### **Bootstrap Process:**

- Draw a random sample **with replacement** of size *n* from the observed dataset.
- Compute the statistic of interest (e.g., mean, median, regression coefficient) from the resample.
- Repeat the process a large number of times (e.g., 1000 or more).
- Use the distribution of bootstrap replicates to estimate:
  - Standard error
  - Confidence intervals
  - Bias (if needed)

イロト 不得 トイヨト イヨト 二日

# Bootstrapping: Advantages

#### Advantages of Bootstrapping:

- **Distribution-free**: No assumptions about the form of the population.
- Handles complex statistics for which theoretical formulas are difficult or unknown.
- Easy to implement with modern computing.

# Types of Bootstrap Methods

#### Nonparametric Bootstrap:

- Resample directly with replacement from the observed data.
- No assumptions about the underlying population.

#### **Block Bootstrap:**

- Used when data are dependent (e.g., time series).
- Resample blocks of consecutive observations to preserve local dependence structure.

## Residual Bootstrap (for Regression Models):

- Assumes the model form is correct but allows for random errors.
- Process:
  - If the regression model and calculate residuals.
  - 2 Resample residuals and generate new response values:  $\hat{y}_i + e_i^*$ .
  - In the model using these new responses.

<ロト <部ト <注入 < 注入 = 二 =

# Bootstrap Confidence Intervals

Idea:

- Use the empirical distribution of bootstrap estimates to construct confidence intervals.
- No assumptions about normality or standard errors needed.

### Percentile Method:

- $\bullet\,$  Take the  $\alpha/2$  and  $1-\alpha/2$  quantiles from the bootstrap distribution.
- Confidence interval:

$$\left[\hat{ heta}^*_{(lpha/2)},\hat{ heta}^*_{(1-lpha/2)}
ight]$$

## **Basic Bootstrap Method:**

- Symmetric adjustment around the original estimate  $\hat{\theta}$ .
- Confidence interval:

$$\left[2\hat{\theta}-\hat{\theta}^*_{(1-lpha/2)},\,2\hat{\theta}-\hat{\theta}^*_{(lpha/2)}
ight]$$

・ 同 ト ・ ヨ ト ・ ヨ ト …