

Study Guide

Quiz 2: PCA, Clustering, and Regression

DataSci 347: Machine Learning 1

How to Use This Guide

This study guide covers all concepts needed for Quiz 2. Work through each section carefully, complete the practice problems, and verify your understanding. If you can answer all the checkpoint questions correctly, you should be well-prepared for the quiz.

Contents

1	Principal Component Analysis (PCA)	3
1.1	What is PCA?	3
1.2	Standardization: When and Why	3
1.2.1	Centering	3
1.2.2	Scaling	3
1.3	Interpreting Principal Components	4
1.3.1	Principal Component Loadings	4
1.3.2	Computing PC Scores	4
1.4	Variance and PCA	4
1.5	Proportion of Variance Explained (PVE)	5
1.6	Practice Problems: PCA	5
2	K-Means Clustering	6
2.1	What is K-Means?	6
2.2	How K-Means Works	6
2.3	Python Implementation	6
2.4	Important Notes	6
2.5	Interpreting Cluster Sizes	7
2.6	Practice Problems: Clustering	7
3	Linear Regression	8
3.1	Simple Linear Regression	8
3.2	Interpreting Coefficients	8
3.3	Multiple Linear Regression	9
3.4	Making Predictions	9
3.5	Python Implementation	9
3.6	Model Quality: R^2	10
3.7	F-Test for Overall Model Significance	11

3.8	Practice Problems: Regression	12
4	Model Diagnostics	13
4.1	Residual Plots	13
4.1.1	Residuals vs. Fitted Values Plot	13
4.2	Interpreting Residual Patterns	13
4.3	Visual Examples	14
4.4	Practice Problems: Diagnostics	14
5	Common Pitfalls and Tips	15
5.1	PCA Pitfalls	15
5.2	Clustering Pitfalls	15
5.3	Regression Pitfalls	15
5.4	Exam Strategy	15
6	Quick Reference Formulas	16
6.1	PCA	16
6.2	Regression	16
6.3	Key Python Commands	16

1 Principal Component Analysis (PCA)

1.1 What is PCA?

Principal Component Analysis is a dimensionality reduction technique that transforms correlated variables into a set of uncorrelated variables called principal components (PCs). The key goals are:

- Reduce dimensionality while retaining maximum variance
- Create uncorrelated features from correlated ones
- Identify patterns in high-dimensional data

1.2 Standardization: When and Why

CRITICAL CONCEPT: Before running PCA, you must understand whether to standardize (center and scale) your data.

1.2.1 Centering

Centering means subtracting the mean from each variable:

$$x_{centered} = x - \bar{x} \quad (1)$$

1.2.2 Scaling

Scaling means dividing by the standard deviation after centering:

$$x_{scaled} = \frac{x - \bar{x}}{s_x} \quad (2)$$

When to standardize:

- Variables are measured in different units (e.g., weight in kg, height in cm)
- Variables have vastly different variances
- You want each variable to contribute equally to the analysis

Python Implementation:

```
1 from sklearn.preprocessing import StandardScaler
2 from sklearn.decomposition import PCA
3
4 # Standardize the data
5 scaler = StandardScaler() # This centers AND scales
6 X_scaled = scaler.fit_transform(X)
7
8 # Then run PCA
9 pca = PCA()
10 pca.fit(X_scaled)
```

1.3 Interpreting Principal Components

1.3.1 Principal Component Loadings

The **loadings** (also called rotation matrix or components) tell you how each original variable contributes to each PC.

```

1 # Get the loadings
2 loadings = pca.components_.T # Transpose to get variables x PCs
3 print(loadings)
4 #          PC1      PC2      PC3      PC4
5 # Var1    0.50   -0.50    0.20   -0.70
6 # Var2    0.49    0.51   -0.30    0.10
7 # Var3    0.51    0.48    0.40    0.20
8 # Var4    0.50   -0.49   -0.30    0.67

```

How to read loadings:

- Each column represents one principal component
- Each row shows how much that original variable contributes
- Positive values: variable increases with PC
- Negative values: variable decreases with PC
- Larger absolute values = stronger contribution

1.3.2 Computing PC Scores

EXAM TIP: You must know how to write the formula for PC scores!

If data is **NOT** standardized, PC1 score for observation i is:

$$PC1_i = w_1 \times x_{1i} + w_2 \times x_{2i} + \dots + w_p \times x_{pi} \quad (3)$$

If data **IS** standardized, PC1 score for observation i is:

$$PC1_i = w_1 \times \frac{(x_{1i} - \bar{x}_1)}{s_1} + w_2 \times \frac{(x_{2i} - \bar{x}_2)}{s_2} + \dots + w_p \times \frac{(x_{pi} - \bar{x}_p)}{s_p} \quad (4)$$

where w_j are the loadings for PC1, \bar{x}_j are means, and s_j are standard deviations.

Example: Given loadings [0.5, 0.5, 0.5, 0.5] and standardized data, PC1 equals:

$$PC1 = 0.5 \times \frac{(x_1 - \bar{x}_1)}{s_1} + 0.5 \times \frac{(x_2 - \bar{x}_2)}{s_2} + 0.5 \times \frac{(x_3 - \bar{x}_3)}{s_3} + 0.5 \times \frac{(x_4 - \bar{x}_4)}{s_4} \quad (5)$$

1.4 Variance and PCA

KEY FACT: Principal components are ordered by variance!

- PC1 has the **largest variance** of all possible linear combinations
- PC2 has the second largest variance, uncorrelated with PC1
- PC3 has the third largest variance, uncorrelated with PC1 and PC2

- And so on...

```

1 # Get variance explained by each PC
2 var_explained = pca.explained_variance_
3 print(var_explained)
4 # [3.2, 0.5, 0.2, 0.1] # PC1 has highest variance
5
6 # Total variance
7 total_var = np.sum(var_explained)
8 print(total_var) # Sum of all PC variances = total variance

```

1.5 Proportion of Variance Explained (PVE)

PVE tells you what percentage of total variance each PC captures.

$$PVE_j = \frac{\text{Variance of } PC_j}{\text{Total Variance}} = \frac{\text{Variance of } PC_j}{\sum_{i=1}^p \text{Variance of } PC_i} \quad (6)$$

```

1 # Calculate PVE
2 pve = pca.explained_variance_ratio_
3 print(pve)
4 # [0.80, 0.125, 0.05, 0.025]
5
6 # This means:
7 # PC1 explains 80% of total variance
8 # PC2 explains 12.5% of total variance
9 # PC3 explains 5% of total variance
10 # PC4 explains 2.5% of total variance
11
12 # They sum to 1 (100%)
13 print(np.sum(pve)) # 1.0

```

Interpreting PVE Plots:

- X-axis: Principal component index (1, 2, 3, ...)
- Y-axis: Proportion of variance explained
- Look for "elbow" to determine how many PCs to keep
- First PC always has highest PVE

1.6 Practice Problems: PCA

Problem 1: You run PCA on 5 variables after centering and scaling. The loadings for PC1 are [0.45, 0.45, 0.44, 0.46, 0.45]. Write the formula for computing PC1 scores. Include the means and standard deviations in your formula.

Problem 2: Given PVE values [0.65, 0.20, 0.10, 0.05], what percentage of variance does PC1 capture? What about PC1 and PC2 combined?

Problem 3: True or False: If you run PCA on standardized data and all variables have equal weight (similar loadings) in PC1, then PC1 is approximately an average of the standardized variables.

Problem 4: You have 4 PCs with variances [16, 3, 0.8, 0.2]. Calculate the PVE for each PC. Which PC explains the most variance?

2 K-Means Clustering

2.1 What is K-Means?

K-means is an unsupervised learning algorithm that partitions data into k clusters. Each observation belongs to the cluster with the nearest mean (centroid).

2.2 How K-Means Works

1. Choose number of clusters k
2. Randomly initialize k cluster centroids
3. Assign each point to nearest centroid
4. Recalculate centroids based on assigned points
5. Repeat steps 3-4 until convergence

2.3 Python Implementation

```
1 from sklearn.cluster import KMeans
2
3 # Fit k-means with 2 clusters
4 kmeans = KMeans(n_clusters=2, random_state=42)
5 kmeans.fit(X)
6
7 # Get cluster assignments (0 or 1 for 2 clusters)
8 labels = kmeans.labels_
9 print(labels) # [0, 1, 0, 1, 1, 0, ...]
10
11 # Count observations in each cluster
12 unique, counts = np.unique(labels, return_counts=True)
13 print(counts) # [120, 80] means 120 in cluster 0, 80 in cluster 1
14
15 # Get cluster centers
16 centers = kmeans.cluster_centers_
17 print(centers)
```

2.4 Important Notes

EXAM TIP: Pay attention to:

- The `random_state` parameter - affects initialization
- Cluster labels are arbitrary (0, 1, 2, ...) - no inherent ordering
- Cluster sizes are NOT necessarily equal
- Results depend on initialization and can vary

2.5 Interpreting Cluster Sizes

```
1 # If you get output like:
2 unique, counts = np.unique(labels, return_counts=True)
3 print(counts)
4 # [116, 84]
5
6 # This means:
7 # - There are 2 clusters total (length of array)
8 # - Cluster 0 has 116 observations
9 # - Cluster 1 has 84 observations
10 # - Total: 116 + 84 = 200 observations
```

Common Question Format: "There are X subjects in cluster 1 and Y in cluster 2"

- Make sure $X + Y = \text{total number of observations}$
- Check which number corresponds to which cluster
- Labels can be 0-indexed or 1-indexed depending on context

2.6 Practice Problems: Clustering

Problem 1: You run k-means with $k=3$ on 300 observations. The output shows counts [100, 125, 75]. How many observations are in each cluster? Do clusters have equal sizes?

Problem 2: True or False: K-means always produces clusters of equal size.

Problem 3: If you run k-means twice with different random seeds, will you get the same cluster sizes? Why or why not?

3 Linear Regression

3.1 Simple Linear Regression

Simple linear regression models the relationship between one predictor X and response Y :

$$Y = \beta_0 + \beta_1 X + \epsilon \quad (7)$$

where:

- β_0 = intercept
- β_1 = slope (coefficient)
- ϵ = random error

3.2 Interpreting Coefficients

The Slope (β_1):

"On average, for every 1 unit increase in X , Y changes by β_1 units."

CRITICAL DISTINCTION:

- **Average effect:** On average, Y decreases by β_1 when X increases by 1
- **Individual predictions:** For specific observations, actual values can vary around the predicted mean

Example:

```

1 # Output:
2 # Intercept: 35.82
3 # Coefficient: -0.044
4
5 # This means:
6 # Predicted MPG_Hwy = 35.82 - 0.044 * Horsepower

```

Interpretation:

- On average, MPG_Hwy decreases by 0.044 for each 1 unit increase in Horsepower
- A car with Horsepower=200 has predicted MPG = $35.82 - 0.044(200) = 27.02$
- A car with Horsepower=201 has predicted MPG = $35.82 - 0.044(201) = 26.976$

EXAM TIP: Watch out for tricky questions!

TRUE Statement: "On average, MPG decreases by 0.044 when Horsepower increases by 1."

FALSE Statement: "For any two specific cars where one has Horsepower=200 and another has Horsepower=201, the first car is GUARANTEED to have higher MPG."

Why is the second false? Because individual observations have variability around the regression line. The model predicts the **mean**, not individual values exactly.

3.3 Multiple Linear Regression

Multiple regression includes multiple predictors:

$$Y = \beta_0 + \beta_1 X_1 + \beta_2 X_2 + \dots + \beta_p X_p + \epsilon \quad (8)$$

CRITICAL CONCEPT: Interpretation of coefficients in multiple regression

β_1 = "The average change in Y for a 1-unit increase in X_1 , **holding all other variables constant.**"

Example:

```
1 # Output:
2 # Intercept: 46.99
3 # Horsepower: -0.028
4 # Weight: -4.01
5
6 # Model: MPG_Hwy = 46.99 - 0.028*Horsepower - 4.01*Weight
```

How to interpret -0.028 for Horsepower:

- "On average, for each 1 unit increase in Horsepower, MPG_Hwy decreases by 0.028, **when Weight is held constant.**"
- You cannot say "1 unit increase in Horsepower always decreases MPG by 0.028" without the "holding Weight constant" qualifier

3.4 Making Predictions

To predict: Plug values into the equation

```
1 # Given: MPG_Hwy = 46.99 - 0.028*Horsepower - 4.01*Weight
2 # Predict for: Horsepower=240, Weight=3.5
3
4 predicted_MPG = 46.99 - 0.028*240 - 4.01*3.5
5 predicted_MPG = 46.99 - 6.72 - 14.035
6 predicted_MPG = 26.235
```

IMPORTANT: You can only use variables that are in the model!

- If Seats and Length are not in the model, you don't need their values to predict
- Having extra information doesn't hurt - just ignore variables not in the model
- You cannot make predictions if required variables are missing

3.5 Python Implementation

```
1 from sklearn.linear_model import LinearRegression
2 import numpy as np
3
4 # Simple regression: Y ~ X1
5 X = car_data[['Horsepower']] # Need double brackets for DataFrame
6 y = car_data['MPG_Hwy']
7
8 model = LinearRegression()
```

```

9  model.fit(X, y)
10
11  print(f"Intercept: {model.intercept}")
12  print(f"Coefficient: {model.coef_[0]}")
13
14  # Multiple regression: Y ~ X1 + X2
15  X_multi = car_data[['Horsepower', 'Weight']]
16  y = car_data['MPG_Hwy']
17
18  model2 = LinearRegression()
19  model2.fit(X_multi, y)
20
21  print(f"Intercept: {model2.intercept}")
22  print(f"Coefficients: {model2.coef}")
23  # model2.coef_[0] is coefficient for Horsepower
24  # model2.coef_[1] is coefficient for Weight
25
26  # Make predictions
27  new_data = np.array([[240, 3.5]]) # Horsepower=240, Weight=3.5
28  prediction = model2.predict(new_data)
29  print(f"Predicted MPG: {prediction[0]}")

```

3.6 Model Quality: R^2

R^2 (**R-squared**) measures the proportion of variance in Y explained by the model.

$$R^2 = 1 - \frac{SS_{residual}}{SS_{total}} = 1 - \frac{\sum (y_i - \hat{y}_i)^2}{\sum (y_i - \bar{y})^2} \quad (9)$$

Interpretation:

- $R^2 = 0.465$ means "The model explains 46.5% of the variance in Y "
- Range: 0 to 1 (0% to 100%)
- Higher R^2 = better fit
- R^2 always increases when adding more variables

Adjusted R^2 : Penalizes adding unhelpful variables

$$R^2_{adj} = 1 - \frac{(1 - R^2)(n - 1)}{n - p - 1} \quad (10)$$

where n = number of observations, p = number of predictors

```

1  from sklearn.metrics import r2_score
2
3  # Calculate R-squared
4  y_pred = model.predict(X)
5  r2 = r2_score(y, y_pred)
6  print(f"R-squared: {r2}")
7
8  # Or use model's score method
9  r2 = model.score(X, y)
10 print(f"R-squared: {r2}")

```

3.7 F-Test for Overall Model Significance

The F-test tests whether **at least one** predictor is useful:

$$H_0 : \beta_1 = \beta_2 = \dots = \beta_p = 0 \text{ (no predictors are useful)} \quad (11)$$

$$H_A : \text{At least one } \beta_j \neq 0 \text{ (at least one predictor is useful)} \quad (12)$$

F-statistic formula:

$$F = \frac{(R^2/p)}{((1 - R^2)/(n - p - 1))} \quad (13)$$

where p = number of predictors, n = sample size

EXAM TIP: Relationship between F-test and R^2

- Larger R^2 leads to larger F-statistic
- But you need the F-test to determine if R^2 is "large enough" to be statistically significant
- Large R^2 alone doesn't prove significance - you need the p-value from the F-test

How to reject H_0 :

- If F-statistic p-value $\leq \alpha$ (significance level), reject H_0
- Example: p-value = 0.0001, $\alpha = 0.001$, so $0.0001 \leq 0.001 \rightarrow$ reject H_0
- "p-value $< 2e-16$ " means p-value is extremely small (essentially 0)

```

1 import scipy.stats as stats
2
3 # Calculate F-statistic
4 n = len(y)
5 p = X_multi.shape[1] # number of predictors
6
7 # Get R-squared
8 y_pred = model2.predict(X_multi)
9 r2 = r2_score(y, y_pred)
10
11 # Calculate F-statistic
12 f_stat = (r2 / p) / ((1 - r2) / (n - p - 1))
13
14 # Get p-value
15 f_pvalue = 1 - stats.f.cdf(f_stat, p, n - p - 1)
16
17 print(f"F-statistic: {f_stat}")
18 print(f"p-value: {f_pvalue}")
19
20 # Interpretation:
21 if f_pvalue < 0.001:
22     print("Reject H0 at alpha=0.001")
23     print("At least one predictor is significant")

```

3.8 Practice Problems: Regression

Problem 1: Given the model: $\hat{Y} = 50 - 0.05X$

- What is the predicted Y when $X=100$?
- What is the predicted Y when $X=101$?
- On average, what happens to Y when X increases by 1?
- If person A has $X=100$ and person B has $X=101$, is Y guaranteed to be higher for person A?

Problem 2: Given: $\hat{Y} = 30 + 2X_1 - 5X_2$

- Interpret the coefficient of X_1
- Interpret the coefficient of X_2
- Predict Y when $X_1 = 10, X_2 = 3$
- Can you predict Y if you're given $X_1 = 10, X_2 = 3, X_3 = 5$?

Problem 3: A model has $R^2 = 0.65$ and $p=2$ predictors, $n=200$ observations. Calculate the F-statistic.

Problem 4: True or False: "If $R^2 = 0.8$, we can automatically reject $H_0 : \beta_1 = \beta_2 = 0$ at $\alpha = 0.001$." Explain.

4 Model Diagnostics

4.1 Residual Plots

Residuals are the differences between observed and predicted values:

$$e_i = y_i - \hat{y}_i \quad (14)$$

4.1.1 Residuals vs. Fitted Values Plot

This plot helps check the **linearity assumption**.

```

1 # Create residual plot
2 y_pred = model.predict(X)
3 residuals = y - y_pred
4
5 plt.scatter(y_pred, residuals)
6 plt.axhline(y=0, color='r', linestyle='--')
7 plt.xlabel('Fitted values')
8 plt.ylabel('Residuals')
9 plt.title('Residuals vs Fitted')
10 plt.show()
```

What to look for:

- **Good:** Random scatter around horizontal line at 0
- **Bad:** Clear patterns, curves, or systematic deviations

4.2 Interpreting Residual Patterns

Pattern 1: Underestimation

- If residuals are mostly **positive** (above 0) for certain fitted values
- This means $y_i - \hat{y}_i > 0$, so $y_i > \hat{y}_i$
- The actual values are **higher** than predictions
- We are **underestimating** (predicting too low)

Pattern 2: Overestimation

- If residuals are mostly **negative** (below 0) for certain fitted values
- This means $y_i - \hat{y}_i < 0$, so $y_i < \hat{y}_i$
- The actual values are **lower** than predictions
- We are **overestimating** (predicting too high)

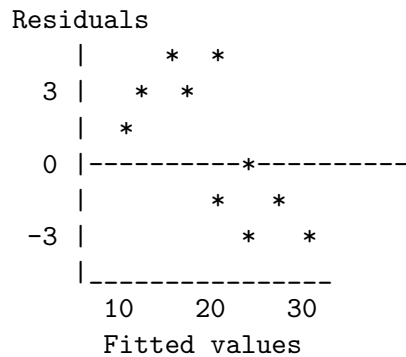
EXAM TIP: Common question pattern

"For cars with smaller MPG.Hwy (left side of plot), residuals are positive. This suggests:"

- Positive residuals = actual > predicted
- We are predicting too low
- We are **underestimating**
- Linearity might be a problem

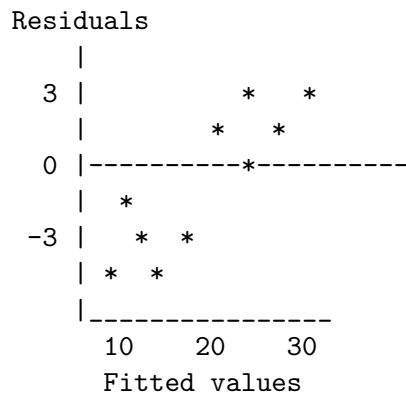
4.3 Visual Examples

Scenario 1: Left side of plot shows mostly positive residuals



Interpretation: Small fitted values (left) have positive residuals → underestimating

Scenario 2: Left side of plot shows mostly negative residuals



Interpretation: Small fitted values (left) have negative residuals → overestimating

4.4 Practice Problems: Diagnostics

Problem 1: You fit a model predicting house prices. The residual plot shows that for expensive houses (right side, high fitted values), most residuals are negative. Are you overestimating or underestimating expensive houses?

Problem 2: For a model predicting test scores, you observe positive residuals for students with low predicted scores. What does this suggest about the model's predictions for low-performing students?

Problem 3: True or False: "If all residuals on the left side of the plot are above 0, the model is overestimating for small values."

5 Common Pitfalls and Tips

5.1 PCA Pitfalls

1. **Forgetting to mention standardization** when writing PC score formulas
2. **Confusing loadings with PC scores** - loadings are weights, scores are the transformed data
3. **Misreading PVE plots** - make sure you understand what the y-axis represents
4. **Thinking PCs are in any order other than by variance** - PC1 ALWAYS has highest variance

5.2 Clustering Pitfalls

1. **Assuming equal cluster sizes** - K-means does NOT force equal sizes
2. **Forgetting clusters depend on initialization** - random seed matters
3. **Confusing cluster labels** - the numbers (0, 1, 2) are arbitrary

5.3 Regression Pitfalls

1. **Confusing "on average" with "always"** - regression predicts means, not guarantees
2. **Forgetting "holding other variables constant"** in multiple regression interpretation
3. **Thinking high R^2 alone proves significance** - need F-test p-value
4. **Mixing up over- and underestimation** - positive residuals = underestimating
5. **Trying to use variables not in the model** for prediction

5.4 Exam Strategy

1. **Read carefully:** Look for words like "always," "guaranteed," "on average"
2. **Check your assumptions:** Is data standardized? What variables are in the model?
3. **Show your work:** For calculations, write out the formula first
4. **Think about interpretation:** Don't just memorize - understand what statistics mean
5. **Watch for trick answers:** Some will be technically correct but misleading

6 Quick Reference Formulas

6.1 PCA

$$\text{Standardized variable: } \frac{x - \bar{x}}{s} \quad (15)$$

$$\text{PC score (unstandardized): } PC = w_1x_1 + w_2x_2 + \dots + w_px_p \quad (16)$$

$$\text{PC score (standardized): } PC = w_1 \frac{(x_1 - \bar{x}_1)}{s_1} + \dots + w_p \frac{(x_p - \bar{x}_p)}{s_p} \quad (17)$$

$$\text{PVE: } \frac{\text{Var}(PC_j)}{\sum_{i=1}^p \text{Var}(PC_i)} \quad (18)$$

6.2 Regression

$$\text{Simple: } Y = \beta_0 + \beta_1 X + \epsilon \quad (19)$$

$$\text{Multiple: } Y = \beta_0 + \beta_1 X_1 + \dots + \beta_p X_p + \epsilon \quad (20)$$

$$\text{Residual: } e_i = y_i - \hat{y}_i \quad (21)$$

$$R^2 : 1 - \frac{\sum (y_i - \hat{y}_i)^2}{\sum (y_i - \bar{y})^2} \quad (22)$$

$$\text{F-statistic: } \frac{(R^2/p)}{((1 - R^2)/(n - p - 1))} \quad (23)$$

6.3 Key Python Commands

```

1  # PCA
2  from sklearn.decomposition import PCA
3  from sklearn.preprocessing import StandardScaler
4  scaler = StandardScaler()
5  X_scaled = scaler.fit_transform(X)
6  pca = PCA()
7  pca.fit(X_scaled)
8  loadings = pca.components_.T
9  pve = pca.explained_variance_ratio_
10
11 # K-Means
12 from sklearn.cluster import KMeans
13 kmeans = KMeans(n_clusters=k, random_state=seed)
14 kmeans.fit(X)
15 labels = kmeans.labels_
16 counts = np.unique(labels, return_counts=True)[1]
17
18 # Regression
19 from sklearn.linear_model import LinearRegression
20 model = LinearRegression()
21 model.fit(X, y)
22 intercept = model.intercept_
23 coefficients = model.coef_
24 predictions = model.predict(X_new)
25 r2 = model.score(X, y)

```


Final Checklist

Before the quiz, make sure you can:

- ☐ Write the formula for PC scores with and without standardization
- ☐ Explain what PVE means and how to interpret a PVE plot
- ☐ Identify which PC has the largest variance
- ☐ Interpret cluster sizes from k-means output
- ☐ Distinguish between "on average" and "always" in regression
- ☐ Interpret regression coefficients in multiple regression
- ☐ Make predictions using regression equations
- ☐ Calculate R^2 and understand what it means
- ☐ Explain when to use the F-test and how to interpret p-values
- ☐ Identify over- vs. underestimation from residual plots
- ☐ Know when holding other variables constant matters

Study smart, think carefully, and you've got this!