# Linear Regression Assignment - BoomBikes Analysis
# Subjective Questions and Answers

## Assignment-Based Questions

### Question 1: What can you infer about the effect of categorical variables on bike demand?

Looking at the categorical variables in our bike sharing data, I discovered some really interesting patterns that significantly impact daily rentals:

**1. Seasonal Impact on Demand:**
- **Fall season** is clearly the winner with **5,644 bikes per day** on average.
- **Summer** comes in second place at **4,992 bikes daily**.
- **Winter** actually performs quite well at **4,728 bikes per day**.
- **Spring** is the slowest season with only **2,604 bikes daily**.
- This makes perfect sense because spring weather can be so unpredictable.

**2. Weather Conditions Make a Huge Difference:**
- **Clear, sunny days** drive the highest demand at **4,876 bikes per day**.
- **Misty or cloudy weather** drops demand to **4,035 bikes daily**.
- **Light rain or snow** really hurts business with only **1,803 bikes per day**.
- **Heavy rain and snow** practically kills demand at just **394 rentals daily**.

**3. Strong Business Growth:**
- **2019** showed much higher demand at **5,634 bikes daily**.
- **2018** had **3,842 bikes daily**.
- That's a **47% growth rate**, which is fantastic for the company.

**4. Monthly Patterns Follow the Seasons:**
- **September and October** are the peak months with over **5,700 bikes daily**.
- **January and February** are the slowest at around **2,400 bikes daily**.

**5. Working Days vs. Weekends:**
- Working days show slightly higher demand than weekends.
- The difference isn't huge but it's consistent.

All these patterns were statistically significant when I tested them, so they represent real business effects rather than random chance.

## Question 2: Why is drop_first=True important when creating dummy variables?

Using **drop_first=True** is absolutely essential because it prevents a serious mathematical problem called the **dummy variable trap**:

**1. What happens without drop_first=True:**
- If we have four seasons and create four dummy variables (Spring, Summer, Fall, Winter).
- These four variables will always sum to exactly 1 for every single row.
- If you know three values, you can always perfectly predict the fourth one.
- This creates **perfect multicollinearity** between the variables.

**2. Why this breaks the regression math:**
- Perfect multicollinearity makes the correlation matrix **singular** (impossible to invert).
- Linear regression needs to calculate $(X'X)^{-1}$ for finding coefficients.
- When the matrix can't be inverted, the whole calculation fails.
- You end up with **infinite VIF values** indicating perfect correlation.

**3. How drop_first=True solves this:**
- It removes one category to use as a **reference baseline**.
- For seasons, we keep three dummies (Summer, Fall, Winter) and drop Spring.
- Now **Spring becomes our reference point** for comparison.
- Summer's coefficient tells us how much higher or lower demand is compared to Spring.

**4. Additional benefits:**
- Makes coefficients much easier to interpret.
- Prevents mathematical errors in model calculation.
- Reduces model complexity without losing any information.

## Question 3: Which numerical variable has the highest correlation with bike demand?

**Temperature** has the highest usable correlation with bike demand at **0.627**. While registered users shows a much higher correlation at 0.972, we can't use that for prediction because it's actually part of our target variable (total = casual + registered), which would be data leakage.

# Question 4: How did you validate the linear regression assumptions?

I systematically checked all four key assumptions of linear regression using both visual plots and statistical tests:

**1. Linearity Assumption:**
  • Created a **residuals vs fitted values plot** to check for random scatter around zero.

  • Calculated correlation between residuals and fitted values, which was **0.002** (very close to zero).

  • This confirmed that the relationship between our predictors and target is truly linear.

**2. Independence Assumption:**
  • Used the **Durbin-Watson test** to check if residuals are correlated with each other.

  • Got a test result of **2.034**, which is very close to the ideal value of 2.0.

  • This means our observations are independent of each other.

**3. Homoscedasticity (Constant Variance):**
  • Created a **scale-location plot** showing square root of absolute residuals vs fitted values.

  • Ran the **Breusch-Pagan test** with a p-value greater than 0.05.

  • Confirmed that residual variance stays constant across all prediction levels.

**4. Normality of Residuals:**
  • Used **Q-Q plots** comparing our residuals to what we'd expect from a normal distribution.

  • Applied the **Shapiro-Wilk test** on a sample of residuals.

  • Most points fell reasonably close to the diagonal line, indicating approximate normality.

**5. Additional Multicollinearity Check:**
  • Calculated **VIF scores** for all variables.

  • All values were below 5, indicating no concerning correlation between predictors.

The validation showed that three out of four assumptions were clearly satisfied, with only minor deviations that are totally acceptable for real-world data.


# Question 5: What are the top 3 features explaining bike demand?

Based on my final model using recursive feature elimination, here are the top three features:

**1. Year Variable (coefficient: +1974):**
  • Being in 2019 increases daily bike demand by about **1,974 bikes** compared to 2018.

  • Shows really strong business growth and increasing popularity of bike sharing.

  • Highly statistically significant with p-value less than 0.001.

**2. Temperature (coefficient: +5848):**
  • Each unit increase in normalized temperature increases demand by about **5,848 bikes**.

- This makes temperature the most powerful weather predictor in our model.

- Shows that bike sharing is incredibly weather-dependent.

- Also highly significant with p-value less than 0.001.

**3. Light Rain/Snow Weather (coefficient: -1901):**
- When there's light rain or snow compared to clear weather, demand drops by about **1,901 bikes**.

- Really demonstrates how much adverse weather conditions hurt the business.

- Critical factor for planning weather-responsive operations.

- Highly significant with p-value less than 0.001.

Together, these three features help explain about **84% of the variation** in bike demand, making our model quite accurate for business planning and daily operations.

# General Questions

## Question 1: Explain the linear regression algorithm in detail

Linear regression is one of the most fundamental machine learning algorithms, and it's beautifully simple in concept:

**1. Basic Concept:**
- For one variable, it finds the best straight line $y = mx + b$ that fits through your data points.

- For multiple variables, it extends to $y = b0 + b1*x1 + b2*x2 + ... + bn*xn$.

- The coefficients (b1, b2, etc.) tell you exactly how much each variable affects the outcome.

**2. How the Algorithm Works:**
- It minimizes the **sum of squared errors** between predicted and actual values.

- Takes the difference between what it predicts and what actually happened.

- Squares these differences (to make them all positive) and adds them up.

- Tries to find coefficients that make this total as small as possible.

**3. Solution Methods:**
- **Analytical solution:** Uses matrix math beta = $(X'X)^{-1} * X'Y$ to get exact answer in one step.

- **Gradient descent:** Starts with random coefficients and gradually improves them.

- Analytical works great for smaller datasets, gradient descent for larger ones.

**4. Key Assumptions:**
- The relationship between variables is actually linear.

- Observations are independent of each other.

- Variance of errors remains constant (homoscedasticity).

- Errors are normally distributed.

**5. Main Advantages:**
- Super simple to understand and interpret.

- Very fast to train and make predictions.

- No hyperparameter tuning required.

- Provides statistical significance testing for features.

**6. Main Limitations:**
- Can only capture linear relationships (misses curves and complex patterns).

- Quite sensitive to outliers which can throw off the entire line.

- Requires assumption validation to ensure reliable results.

## Question 2: Explain Anscombes quartet in detail

Anscombes quartet is a brilliant statistical example from 1973 that perfectly demonstrates why you should always visualize your data:

**1. What It Is:**
- Four completely different datasets, each with exactly 11 data points.

- All four have nearly identical statistical properties.

- Same means for x (9.0) and y (7.5), same variances, same correlation coefficient (0.816).

- Even the same linear regression equation (y = 3.00 + 0.500x) and R-squared (0.67).

**2. The Shocking Visual Differences:**
- **Dataset 1:** Nice linear relationship with some scatter - perfect for linear regression.

- **Dataset 2:** Curved, parabolic relationship where linear regression clearly doesnt fit.

- **Dataset 3:** Perfect linear relationship except one extreme outlier ruins everything.

- **Dataset 4:** No relationship at all for most points, but one outlier creates fake correlation.

**3. The Profound Lesson:**
- Identical summary statistics can completely hide different underlying patterns.

- You could look at just the numbers and think all four datasets are identical.

- But they actually require totally different analytical approaches.

- Statistics alone can be incredibly misleading without visualization.

**4. Why This Still Matters Today:**
- Reminds us that exploratory data analysis and visualization are absolutely crucial.

- Shows that evaluation metrics alone dont guarantee your model makes sense.

- Demonstrates how a single outlier can have enormous effects on your results.

- Proves that correlation doesnt necessarily mean theres a meaningful relationship.

This example remains highly relevant in modern data science and serves as a constant reminder that proper analysis requires both statistical rigor and visual intelligence.

## Question 3: What is Pearsons R?

Pearsons R (also called the Pearson correlation coefficient) measures how strongly two variables are linearly related:

**1. Value Range and Interpretation:**
- Always falls between **-1 and +1**.

- **+1** means perfect positive relationship (as one goes up, other goes up proportionally).

- **-1** means perfect negative relationship (as one goes up, other goes down proportionally).

- **0** means no linear relationship whatsoever.

**2. Strength Guidelines:**
   • **Above 0.7:** Strong relationship.

   • **Between 0.3 and 0.7:** Moderate relationship.

   • **Below 0.3:** Weak relationship.

**3. How It Actually Works:**
   • Looks at how much each data point deviates from its variables average.

   • Multiplies these deviations together for each pair of points.

   • Standardizes the result so its independent of units or scale.

   • Gives you a clean measure that works regardless of whether youre measuring temperature or income.

**4. Real Example from Our Project:**
   • Temperature had a Pearsons R of **0.627** with bike demand.

   • This indicates a strong positive relationship, which makes perfect intuitive sense.

   • Warmer weather definitely encourages more people to rent bikes.

**5. Important Limitations to Remember:**
   • Only captures linear relationships (completely misses curves and other patterns).

   • Very sensitive to outliers which can artificially boost or deflate the correlation.

   • **Correlation doesnt imply causation** - just because two things move together doesnt mean one causes the other.

   • Works best when both variables are roughly normally distributed.


## Question 4: What is scaling and why is it important?

Scaling transforms your numerical features so theyre all on similar ranges, preventing any single variable from dominating just because of its scale:

**1. The Problem Without Scaling:**
   • Variables with bigger numbers unfairly bias the algorithm.

   • Classic example: Age (ranges 20-80) vs Income (ranges 20,000-200,000).

   • Income would completely dominate the analysis simply because its numbers are much larger.

   • Even though both variables might be equally important for your prediction.

**2. Key Benefits of Scaling:**
   • All features get to contribute equally to the model.

   • Algorithms converge much faster (especially gradient-based ones).

   • Better numerical stability in calculations.

   • Fair comparison of feature importance scores.

**3. Two Main Scaling Types:**
  • **Normalization (Min-Max Scaling):** Formula: (x - minimum) / (maximum - minimum).

  • Transforms everything to fall between 0 and 1.

  • **Standardization (Z-Score):** Formula: (x - mean) / standard deviation.

  • Results in mean=0 and standard deviation=1.

**4. When Scaling is Especially Critical:**
  • Algorithms using distance calculations (k-nearest neighbors, SVM, neural networks).

  • Gradient descent optimization (helps it converge faster).

  • Principal component analysis and clustering.

  • Any algorithm that compares feature magnitudes.

**5. Critical Best Practice:**
  • Always fit the scaler on training data only.

  • Apply the same transformation to test data.

  • This prevents data leakage and ensures your model works on new, unseen data.

# Question 5: Why might VIF be infinite?

VIF (Variance Inflation Factor) becomes infinite when you have perfect multicollinearity in your regression variables, which creates serious mathematical problems:

**1. Understanding VIF:**
- VIF measures how much multicollinearity exists between your variables.
- Formula: VIF = 1 / (1 - R-squared), where R-squared comes from regressing one variable against all others.
- When R-squared = 1 (perfect correlation), you get VIF = 1/0 = infinite.

**2. Most Common Cause - Dummy Variable Trap:**
- Happens when you create dummy variables for all categories of a categorical variable.
- Example: Four season dummies (Spring, Summer, Fall, Winter) always sum to exactly 1.
- If you know three values, you can perfectly predict the fourth one.
- Creates perfect linear dependence between variables.

**3. Other Common Causes:**
- **Mathematical dependencies:** Including variables that are exact combinations of others.
- **Duplicate variables:** Accidentally including the same variable twice.
- **Constant multiples:** Having temperature in Celsius and the same temperature times 2.

**4. Why Infinite VIF Breaks Everything:**
- Makes the correlation matrix singular (mathematically non-invertible).
- Linear regression needs to calculate $(X'X)^{-1}$ for finding coefficients.
- When matrix inversion fails, the entire regression calculation collapses.
- Results in computational errors and completely unreliable results.

**5. Simple Solutions:**
- **For dummy variables:** Use drop_first=True to remove one category as reference.
- **For mathematical dependencies:** Remove the redundant variables.
- **For duplicates:** Keep only one copy of each variable.
- **Prevention:** Check correlation matrix and VIF scores before modeling.


# Question 6: What is a Q-Q plot and its importance in linear regression?

A Q-Q (quantile-quantile) plot is a powerful graphical tool that helps you check whether your data follows a particular distribution, usually the normal distribution:

**1. How Q-Q Plots Work:**

• Compares the quantiles of your actual data against what youd expect from a perfect normal distribution.

• If your data is normally distributed, the points will fall approximately on a straight diagonal line.

• Any systematic deviations from this line indicate departures from normality.

**2. Reading Different Patterns:**

• **Points on straight line:** Your data is normally distributed (excellent!).

• **Curve upward at both ends:** Distribution has heavier tails (more extreme values than normal).

• **Curve downward at both ends:** Distribution has lighter tails (fewer extreme values).

• **Consistent upward or downward curve:** Data is skewed in one direction.

**3. Critical Importance in Linear Regression:**

• Tests the crucial assumption that residuals should be normally distributed.

• This assumption directly affects the validity of confidence intervals and hypothesis tests.

• If residuals arent normal, your statistical inference might be completely wrong.

**4. Why This Assumption Matters So Much:**

• Confidence intervals for coefficients rely on the normal distribution assumption.

• Hypothesis testing (t-tests, F-tests) assumes residuals are normally distributed.

• Prediction intervals accuracy depends heavily on residual normality.

**5. How We Used Q-Q Plots in Our Analysis:**

• Created Q-Q plots of our model residuals to verify approximate normality.

• Points fell reasonably close to the diagonal line, which was reassuring.

• This gave us confidence that our statistical tests and confidence intervals were reliable.

**6. What to Do When Problems Are Found:**

• Consider data transformations (log, square root, Box-Cox).

• Investigate and potentially treat or remove outliers.

• Might need alternative modeling approaches that dont assume normality.

• Sometimes the problem is minor enough to proceed with caution.

Q-Q plots provide an intuitive visual check that perfectly complements statistical tests, helping ensure your model assumptions are reasonable and your results are trustworthy.