

# MLDL assignment1

Q1) Explain whether each scenario is a classification or regression problem, and indicate whether we are most interested in inference or prediction. Finally, provide  $n$  and  $p$ .

(a) We collect a set of data on the top 500 firms in the US. For each firm we record profit, number of employees, industry, and the CEO's salary. We are interested in understanding which factors affect CEO salary.

profit + number of employees → CEO salary (Quantitative Response)

Regression

Inference, CEO salary based on features of the firm.

$n = 500$ ,  $p = 3$  (profit, number of employees, industry)

(b) We are considering launching a new product and wish to know whether it will be a success or a failure. We collect data on 20 similar products that were previously launched. For each product, we have recorded whether it was a success or failure, the price charged for the product, the marketing budget, the competition price, and ten other variables.

price charged for the product, the marketing budget, the competition price, and ten other variables → Success or Failure (Qualitative Response)

Classification

Prediction, Success, or Failure of the new product.

$n = 20$ ,  $p = 13$  (price charged for the product, the marketing budget, the competition price, and ten other variables)

(c) We are interested in predicting the % change in the USD/Euro exchange rate in relation to the weekly changes in the world stock markets. Hence we collect weekly data for all of 2012. For each week we record the % change in the USD/Euro, the % change in the US market, the % change in the British market, and the % change in the German market.

USD/Euro exchange rate

Regression

Prediction, % change (Quantitative Response)

$n = 52$  (weekly data for all of 2012),  $p = 3$  (the % change in the US market, the % change in the British market, and the % change in the German market)

Q2) You will now think of some real-life applications for statistical learning

(a) Describe three real-life applications in which classification might be useful. Describe the response, as well as predictors. Is the goal of each application inference or prediction? Explain your answer.

## 1. Email Spam Detection:

- **Response:** The response in this application is binary - classifying emails as either "spam" or "not spam".
- **Predictors:** Predictors are the various features of an email, such as sender, subject, content, and attached links or files.
- **Goal:** The goal in email spam detection is prediction. The model predicts whether an incoming email belongs to the "spam" class or the "not spam" class based on the given features. The primary objective is to correctly classify emails to reduce unwanted spam in users' inboxes.

## 2. Medical Diagnosis:

- **Response:** The response in medical diagnosis can vary depending on the specific case, but it often involves classifying a patient's condition into categories like "healthy," "benign," or "malignant."

- **Predictors:** Predictors include patient medical history, symptoms, test results (e.g., blood tests, imaging), and genetic information.
- **Goal:** The goal in medical diagnosis can be both inference and prediction. Inference is often used to understand the relationship between predictors and outcomes, such as identifying significant risk factors for a disease. Prediction is used to diagnose a patient's current condition based on the available data, assisting healthcare professionals in making informed decisions about treatment and care.

### 3. Sentiment Analysis in Social Media:

- **Response:** The response in sentiment analysis is typically categorical, classifying text data into sentiment categories like "positive," "negative," or "neutral."
- **Predictors:** Predictors are the textual content of social media posts, comments, or reviews, along with additional context like user information, timestamps, or hashtags.
- **Goal:** The primary goal in sentiment analysis is prediction. It aims to predict the sentiment expressed in a given text based on the content and context. For example, businesses use sentiment analysis to gauge public opinion about their products or services. While some inferences about public sentiment trends can be drawn, the primary focus is on predicting the sentiment of individual pieces of content.

(b) Describe three real-life applications in which regression might be useful. Describe the response, as well as the predictors. Is the goal of each application inference or prediction? Explain your answer.

#### 1. Economics - Inflation Prediction:

- **Response Variable:** Inflation rate (e.g., Consumer Price Index - CPI)
- **Predictors:** Various economic indicators like unemployment rate, money supply, GDP growth, and interest rates.
- **Goal:** Prediction
- **Explanation:** In this application, the primary goal is prediction. Economists use regression to model the relationship between inflation (response variable) and various economic indicators (predictors) to forecast future inflation rates. Understanding the relationships between these predictors and inflation is secondary, as the primary aim is to make accurate predictions for economic planning.

#### 2. Medicine - Disease Progression Modeling:

- **Response Variable:** Progression of a disease (e.g., tumor size, blood pressure)
- **Predictors:** Patient demographics (age, gender), genetic factors, lifestyle variables (diet, exercise), and medical treatments.
- **Goal:** Inference and Prediction
- **Explanation:** In this case, regression serves both inference and prediction purposes. Researchers may use regression to understand how different factors (predictors) influence the progression of a disease (inference). Additionally, doctors can use regression models to predict how a disease will progress in an individual patient based on their characteristics and treatment history (prediction).

#### 3. Retail - Sales Forecasting:

- **Response Variable:** Sales revenue (e.g., monthly sales)
- **Predictors:** Historical sales data, marketing expenditures, seasonality, economic conditions, and competitor actions.
- **Goal:** Prediction
- **Explanation:** Sales forecasting is primarily about prediction. Retailers use regression models to analyze historical sales patterns and their relationship with various predictors. By doing so, they can make accurate predictions of future sales, allowing for better inventory management, staffing, and marketing strategy planning. While understanding the impact of marketing spend or other factors on sales is important, the primary objective is to predict future sales volumes.

(c) Describe three real-life applications in which cluster analysis might be useful

#### 1. Customer Segmentation in Marketing:

- **Application:** Companies often use cluster analysis to segment their customer base into distinct groups with similar purchasing behavior, demographics, or preferences.
- **Usefulness:** By identifying customer segments, businesses can tailor their marketing strategies, product offerings, and advertising campaigns to better target each group's specific needs and interests. For example, an e-commerce company might discover that it has one segment of budget-conscious shoppers and another segment of luxury buyers. This information allows them to create customized marketing messages and promotions for each group.

## 2. Image and Document Classification:

- **Application:** In computer vision and natural language processing, cluster analysis can be used for image and document classification.
- **Usefulness:** For images, it can group similar images together based on visual features like color, shape, or texture. In document classification, it can group similar documents together based on text content and structure. This is valuable for tasks such as content recommendation, search engines, and organizing large datasets of images or documents. For example, a news website might use cluster analysis to categorize articles into topics like politics, sports, and entertainment.

## 3. Healthcare - Patient Segmentation:

- **Application:** Cluster analysis can be applied to patient data in healthcare to identify groups of patients with similar medical profiles or treatment responses.
- **Usefulness:** Healthcare providers can use patient segmentation for personalized medicine and treatment planning. By clustering patients with similar characteristics, doctors can better understand which treatments are effective for specific groups of patients, predict disease outcomes, and make informed decisions about healthcare resource allocation. For instance, in oncology, clustering patients based on genetic markers can help identify the most appropriate cancer treatments for different groups of patients.

Q3) Describe the difference between a parametric and a non-parametric statistical learning approach.

(a) What are the advantages of a parametric approach to regression or classification (as opposed to a nonparametric approach)? What are its disadvantages

A parametric approach to regression or classification involves making specific assumptions about the functional form of the relationship between the predictors (independent variables) and the response variable (dependent variable). In contrast, a nonparametric approach makes fewer or no assumptions about the functional form of this relationship. Each approach has its own advantages and disadvantages:

### Advantages of a Parametric Approach:

1. **Simplicity and Interpretability:** Parametric models, such as linear regression or logistic regression, have a simple and interpretable form. The coefficients in these models can provide insights into the strength and direction of the relationships between predictors and the response variable.
2. **Efficiency:** Parametric models often require fewer data points to estimate the model parameters accurately compared to nonparametric models. This can be advantageous when dealing with limited data.
3. **Statistical Inference:** Parametric models provide a framework for statistical hypothesis testing and confidence interval estimation. You can test the significance of individual predictors and assess the overall goodness-of-fit of the model.
4. **Reduced Risk of Overfitting:** Because parametric models have a fixed number of parameters, they are less prone to overfitting when the dataset is small or noisy. This makes them useful in situations where overfitting is a concern.

### Disadvantages of a Parametric Approach:

1. **Model Assumptions:** Parametric models rely on assumptions about the functional form of the relationship between predictors and the response variable. If these assumptions are violated, the model's predictions may be inaccurate.
2. **Limited Flexibility:** Parametric models may not capture complex or nonlinear relationships between variables effectively. If the true relationship is highly nonlinear, a parametric model like linear regression may perform poorly.
3. **Bias:** If the chosen parametric form is inappropriate for the data, the model may introduce bias into predictions. This can lead to systematic errors.
4. **Underfitting:** Parametric models may underfit the data when the assumed functional form is too rigid and cannot capture the underlying patterns in the data.

5. **Extrapolation:** Parametric models are generally not well-suited for extrapolation beyond the range of the observed data because they rely on the assumed functional form within that range.

In contrast, nonparametric approaches like decision trees, k-nearest neighbors, or kernel methods make fewer assumptions about the data and can capture complex relationships. However, they often require more data and may be less interpretable. The choice between parametric and nonparametric approaches should depend on the specific characteristics of the data, the research goals, and the trade-offs between simplicity and flexibility. In practice, some machine learning algorithms combine elements of both parametric and nonparametric modeling to strike a balance.

(b) What are the advantages and disadvantages of a very flexible (versus a less flexible) approach for regression or classification? Under what circumstances might a more flexible approach be preferred to a less flexible approach? When might a less flexible approach be preferred?

The flexibility of an approach in regression or classification refers to its ability to capture complex and nonlinear relationships in the data. A very flexible approach can fit the data closely, while a less flexible one may make simpler, more constrained assumptions. Both have their advantages and disadvantages, and the choice depends on the characteristics of the data and the goals of the analysis.

#### **Advantages of a Very Flexible Approach:**

1. **Better Fit:** Very flexible models, such as ensemble methods like Random Forests or deep learning neural networks, can fit the data closely and capture intricate patterns, leading to potentially higher predictive accuracy.
2. **Handling Complex Relationships:** These models can handle highly nonlinear and intricate relationships between predictors and the response variable, making them suitable for data with complex structures.
3. **Feature Importance:** Flexible models can often provide insights into feature importance, helping identify which variables have the most significant impact on the outcome.

#### **Disadvantages of a Very Flexible Approach:**

1. **Overfitting:** Very flexible models are more prone to overfitting, especially when dealing with small datasets or noisy data. They can capture noise in the data and perform poorly on unseen data.
2. **Interpretability:** Highly flexible models are typically less interpretable than simpler models. They may provide excellent predictions but offer little insight into why certain predictions are made.
3. **Computational Complexity:** Training and tuning very flexible models can be computationally intensive and time-consuming, which may not be practical for all applications.

#### **When a More Flexible Approach Might Be Preferred:**

1. **Complex Relationships:** When you suspect that the true relationship between predictors and the response variable is highly complex or nonlinear, a more flexible approach can be preferred to capture these nuances.
2. **Large Datasets:** Very flexible models can benefit from larger datasets to help mitigate the risk of overfitting. If you have access to a substantial amount of data, a more flexible approach might be a good choice.
3. **Emphasis on Predictive Accuracy:** In applications where predictive accuracy is the primary concern and model interpretability is less critical, a highly flexible model can be preferred.

#### **When a Less Flexible Approach Might Be Preferred:**

1. **Interpretability:** In cases where model interpretability is crucial for making informed decisions or explaining the results to stakeholders, a less flexible model like linear regression or decision trees might be preferred.
2. **Small Datasets:** When working with small datasets or data with limited features, a simpler model can be less prone to overfitting and may provide more stable results.
3. **Computation Constraints:** If you have limited computational resources or need a model that can be trained quickly, a less flexible approach can be more practical.
4. **Stability:** In situations where you want a model that is less sensitive to small fluctuations in the data or is less likely to capture noise, a less flexible model can provide more stable predictions.

In practice, the choice between a very flexible and a less flexible approach depends on a careful assessment of the specific problem, the available data, computational resources, and the importance of interpretability and generalizability in the context of the application. Often, a balance between flexibility and simplicity may be achieved by using techniques like model ensembles or regularization methods.

Q4) The table below provides a training data set containing six observations, three predictors, and one qualitative response variable. Suppose we wish to use this data set to make a prediction for Y when  $X_1 = X_2 = X_3 = 0$  using K-nearest neighbors

(a) Compute the Euclidean distance between each observation and the test point,  $X_1 = X_2 = X_3 = 0$ .

Obs	X1	X2	X3	Y	Dist
1	0	3	0	R	3
2	2	0	0	R	2
3	0	1	3	R	$\sqrt{10}$ , 3.2
4	0	1	2	G	$\sqrt{5}$ , 2.2
5	-1	0	1	G	$\sqrt{2}$ , 1.4
6	1	1	1	R	$\sqrt{3}$ , 1.7

(b) Green, Because, Obs=5 is 1st nearest point

(c) Red, Because, Obs=2 is 3rd nearest point

(d) If the Bayes decision boundary in this problem is highly nonlinear, then would we expect the best value for K to be large or small?

Small. If K is large, the decision boundary will be linear.

Q5)

(a)

Quantitative: mpg, cylinders, displacement, horsepower, weight, acceleration, year

Qualitative: name, origin

(b)

	mpg	cylinders	displacement	horsepower	weight	acceleration	year
min	11.0	3	68.0	46.0	1649	8.5	70
max	46.6	8	455.0	230.0	4997	24.8	82

(c)

Mean

mpg 23.445918  
 cylinders 5.471939  
 displacement 194.411990  
 horsepower 104.469388  
 weight 2977.584184  
 acceleration 15.541327  
 year 75.979592

Std

mpg 7.805007  
 cylinders 1.705783  
 displacement 104.644004  
 horsepower 38.491160  
 weight 849.402560

```
acceleration  2.758864
year          3.683737
```

(d)

Min, max

```
mpg cylinders displacement horsepower weight acceleration year
min 11.0      3      68.0    46.0 1649      8.5 70
max 46.6      8     455.0   230.0 4997     24.8 82
```

Mean

```
mpg          24.404430
cylinders     5.373418
displacement 187.240506
horsepower   100.721519
weight      2935.971519
acceleration  15.726899
year         77.145570
```

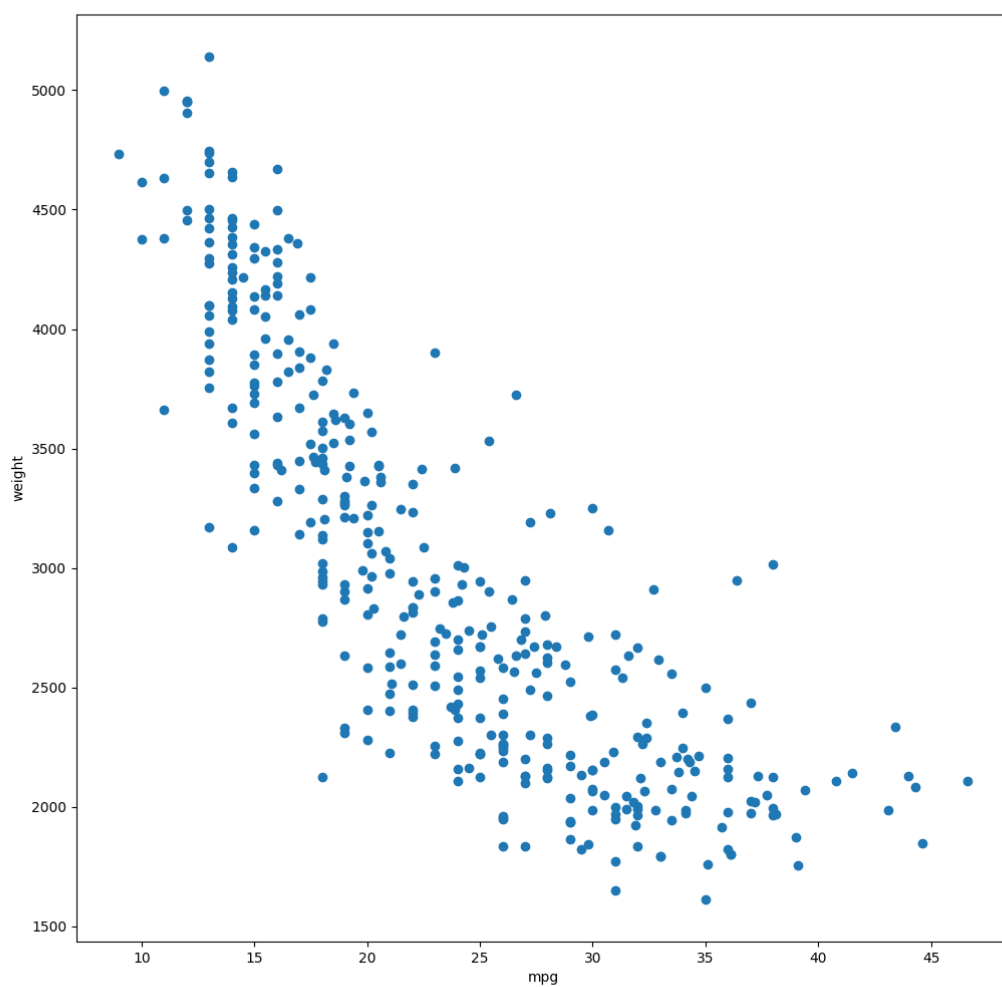
Std

```
mpg          7.867283
cylinders     1.654179
displacement  99.678367
horsepower    35.708853
weight       811.300208
acceleration  2.693721
year         3.106217
```

(e)

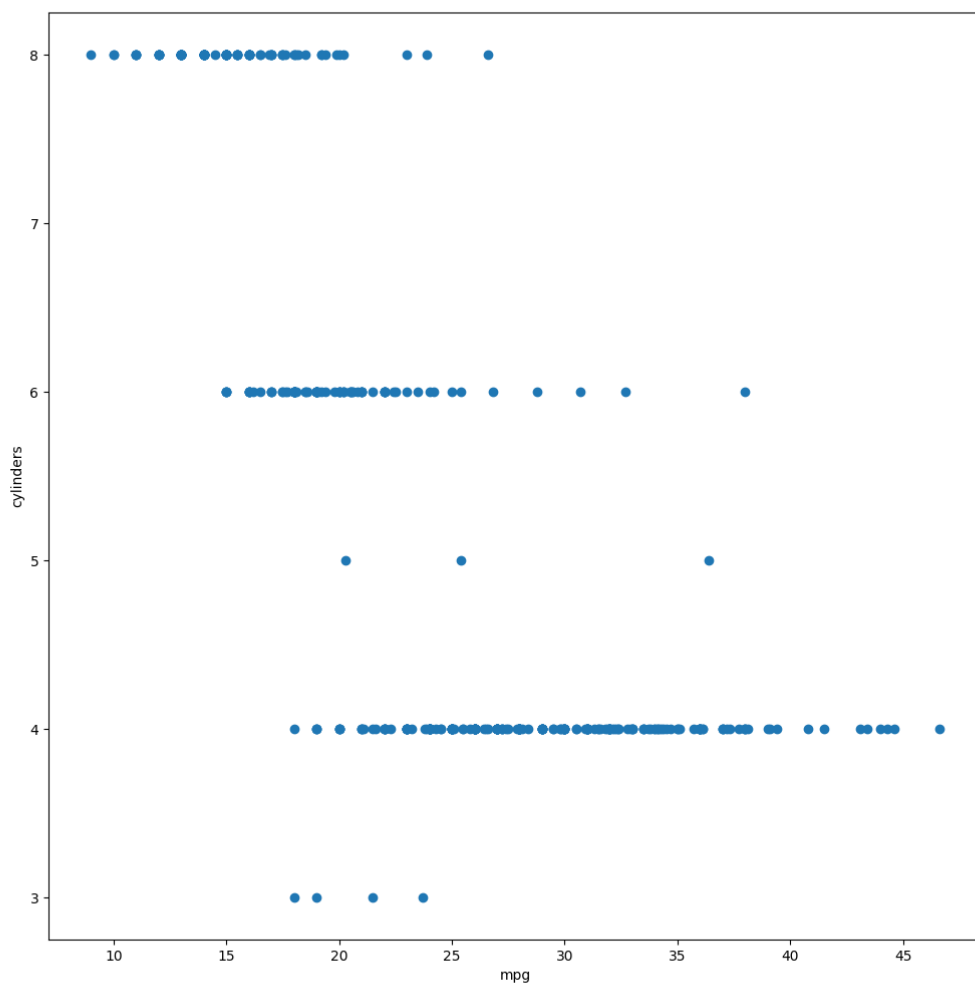
- The heavier the weight, the lower the mpg.

mpg\_vs\_weight



- More cylinders, less mpg.

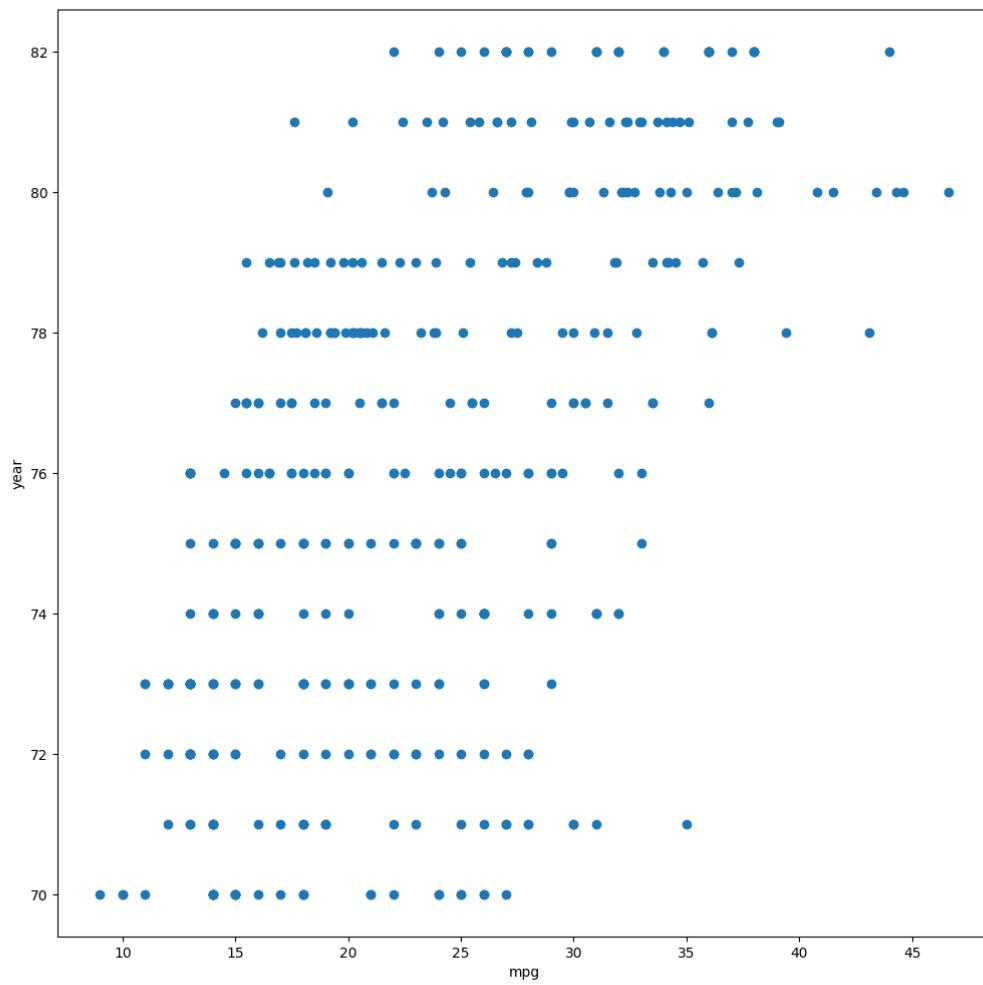
mpg\_vs\_cylinders



- Cars become more efficient over the years.

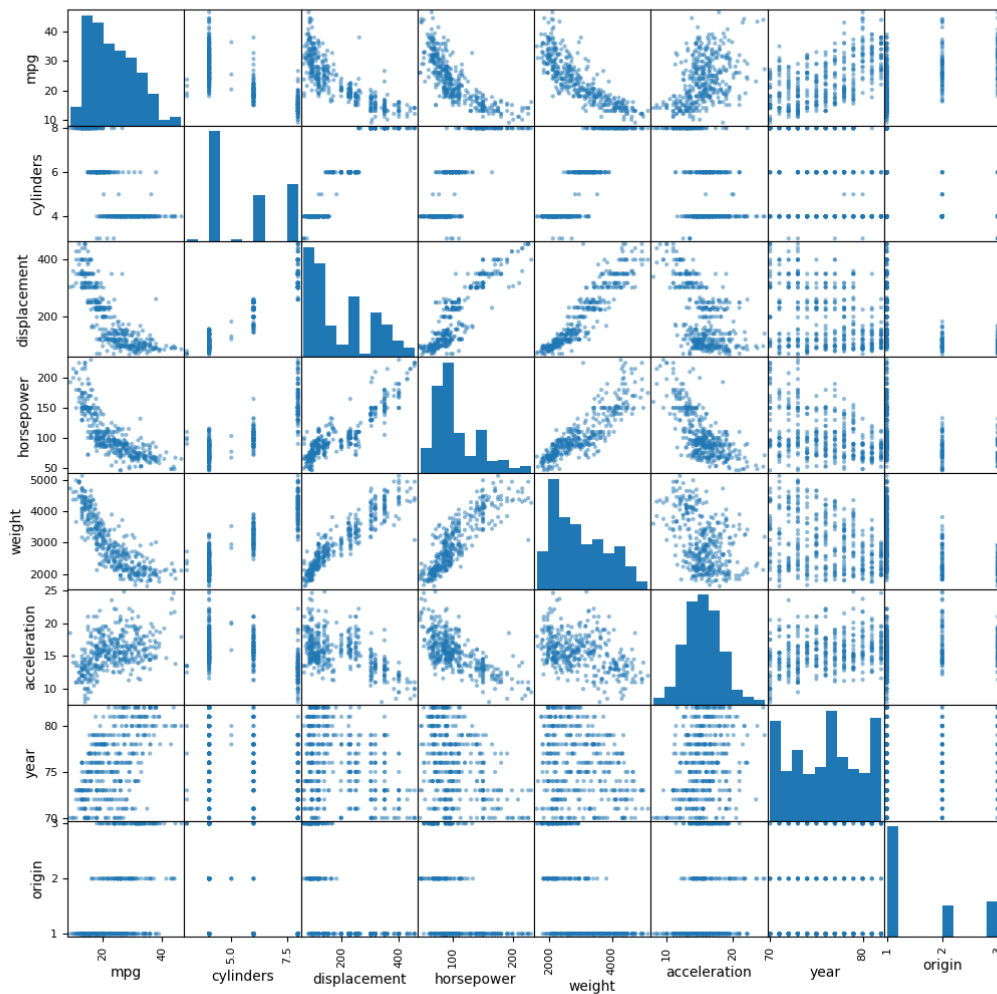
mpg\_vs\_year





(f)

- All predictors correlated with mpg.
- However, for a few observations per name, using 'name' as a predictor is likely to overfit the data and does not generalize well.



Q6)

(a) In your own words, describe what bias and variance are. What is the bias-variance tradeoff?

1. **Bias:** Bias refers to the error introduced by approximating a real-world problem (which may be complex) with a simplified model. It represents the model's tendency to systematically underpredict or overpredict the true values. A model with high bias makes strong assumptions about the data and may oversimplify the underlying relationships, leading to poor accuracy and a lack of flexibility to capture complex patterns.
2. **Variance:** Variance measures the model's sensitivity to fluctuations or noise in the training data. It represents the model's ability to fit the training data points closely. A model with high variance is highly flexible and can capture noise in the data, but it may fail to generalize well to new, unseen data because it has essentially memorized the training data.

The **bias-variance tradeoff** is a key concept in machine learning that describes the delicate balance between these two sources of error:

- **High Bias, Low Variance:** Models with high bias and low variance are typically simple and make strong assumptions about the data. They may underfit the training data by oversimplifying the relationships, but they are more stable and tend to generalize better to new data.

- **Low Bias, High Variance:** Models with low bias and high variance are often very flexible and can fit the training data closely. However, they are sensitive to noise and may overfit the training data, performing poorly on unseen data.

The goal of machine learning is to find the right balance between bias and variance to achieve good predictive performance on unseen data. This balance depends on the specific problem and dataset. Methods for achieving this balance include:

1. **Model Selection:** Choosing an appropriate model or algorithm that matches the complexity of the problem at hand. For example, using a linear regression model for a problem with a linear relationship between variables, and using more complex models like random forests or deep neural networks for problems with nonlinear relationships.
2. **Regularization:** Applying techniques like L1 or L2 regularization to penalize overly complex models, helping to reduce variance.
3. **Cross-Validation:** Using techniques like k-fold cross-validation to assess a model's performance on unseen data, which can help identify whether the model is underfitting or overfitting.
4. **Ensemble Methods:** Combining predictions from multiple models (ensemble methods) can often lead to a reduction in variance while maintaining low bias.

In summary, the bias-variance tradeoff is a fundamental concept in machine learning that emphasizes the need to strike a balance between model simplicity and complexity. Finding this balance is essential for building models that generalize well to new data and make accurate predictions.

(b) What happens to bias and variance when the model complexity increases/decreases?

When the model complexity increases.

1. **Bias Decreases:** As you increase the complexity of a model, it becomes more capable of capturing complex patterns and relationships in the data. This means that the model's bias tends to decrease. In other words, it becomes better at fitting the training data, including the intricate details of the data.
2. **Variance Increases:** Conversely, as the model complexity increases, the model becomes more flexible and sensitive to the noise or fluctuations in the training data. This sensitivity results in an increase in variance. In simple terms, a highly complex model is more likely to overfit the training data, capturing not just the underlying patterns but also the random noise present in the data.

The relationship between model complexity, bias, and variance is often visualized as a U-shaped curve known as the bias-variance tradeoff. Here's what the tradeoff looks like:

- **Low Model Complexity (High Bias, Low Variance):** Simple models, like linear regression, have high bias but low variance. They make strong assumptions and may underfit the data, but they are stable and generalize well.
- **Moderate Model Complexity (Balanced Bias and Variance):** Models with moderate complexity strike a balance between bias and variance. They can capture both the underlying patterns and some of the noise, resulting in good generalization to new data.
- **High Model Complexity (Low Bias, High Variance):** Highly complex models, such as deep neural networks, have low bias but high variance. They are capable of fitting the training data closely, but they are prone to overfitting and may not generalize well.

The key challenge in machine learning is to find the right level of model complexity for a given problem and dataset. It often involves experimentation, model selection, and techniques like cross-validation to assess how well a model generalizes to unseen data. Ideally, you want to choose a model complexity that minimizes the total error (sum of bias and variance), resulting in the best predictive performance on new, unseen data. This is why the bias-variance tradeoff is a critical consideration when developing machine learning models.

When the model complexity decreases.

1. **Bias Increases:** As you decrease the model's complexity, it becomes less capable of capturing complex patterns and relationships in the data. This means that the model's bias tends to increase. In other words, it becomes less able to fit the training data, including important underlying patterns.
2. **Variance Decreases:** Conversely, as the model's complexity decreases, it becomes less flexible and less sensitive to the noise or fluctuations in the training data. This reduced sensitivity results in a decrease in variance. In simpler terms, a less

complex model is less likely to overfit the training data because it cannot capture as much of the random noise present in the data.

The relationship between model complexity, bias, and variance is often visualized as a U-shaped curve known as the bias-variance tradeoff. Here's how the tradeoff looks when you decrease model complexity:

- **High Model Complexity (Low Bias, High Variance):** Highly complex models, such as deep neural networks, have low bias but high variance. They are capable of fitting the training data closely but are prone to overfitting.
- **Moderate Model Complexity (Balanced Bias and Variance):** Models with moderate complexity strike a balance between bias and variance. They can capture both the underlying patterns and some of the noise, resulting in good generalization to new data.
- **Low Model Complexity (High Bias, Low Variance):** Simple models, like linear regression or decision trees with limited depth, have high bias but low variance. They make strong assumptions and may underfit the data, but they are stable and generalize reasonably well.

The key challenge in machine learning is to find the right level of model complexity for a given problem and dataset. The goal is to choose a model complexity that minimizes the total error (sum of bias and variance), resulting in the best predictive performance on new, unseen data. This is why the bias-variance tradeoff is a fundamental consideration when developing machine learning models.

(c) True or False: The bias of a model increases as the amount of training data available increases.

False. The statement is false. In general, the bias of a model does not increase as the amount of training data available increases. In fact, having more training data often helps reduce bias, particularly when using more complex models. Here's why:

Bias is related to how well a model fits the training data. A model with high bias makes strong, often simplistic, assumptions about the data, which can lead to underfitting. When you have more training data, it provides the model with a larger and more diverse set of examples to learn from. This additional data can help the model better capture the underlying patterns in the data, reducing bias.

In contrast, having more training data tends to have a regularizing effect on a model, which can help mitigate overfitting (high variance). Overfitting occurs when a model is too complex and starts to capture noise in the data rather than just the true underlying patterns. With more data, the model is less likely to fit the noise and is more likely to generalize well to new, unseen data.

So, in summary, as the amount of training data increases, it is more common for the bias of a model to decrease or remain stable, while the model's ability to generalize to new data (variance) improves.

(d) True or False: The variance of a model decreases as the amount of training data available increases.

True. The statement is generally true: the variance of a model tends to decrease as the amount of training data available increases. Here's why:

1. **Decreased Sensitivity to Noise:** With a larger training dataset, the model has more examples to learn from. This larger dataset provides a more representative sample of the underlying data distribution, making the model less sensitive to random noise in the training data. As a result, the model is less likely to fit the noise and more likely to focus on the true underlying patterns in the data.
2. **Improved Generalization:** When a model has more data to train on, it has a better chance of learning the underlying patterns and relationships that are consistent across a wider range of examples. This improved generalization ability means the model is likely to perform better on new, unseen data, resulting in lower variance.

However, it's important to note that while increasing the amount of training data generally reduces variance, there can be diminishing returns. In some cases, the reduction in variance may not be significant beyond a certain point, especially if the model is already well-regularized and appropriate for the complexity of the task. Additionally, other factors such as the model's architecture, hyperparameter tuning, and data quality also play a role in controlling variance.

(e) True or False: A learning algorithm will always generalize better if we use less features to represent our data

False. It is not always the case that a learning algorithm will generalize better if fewer features are used to represent the data. Whether reducing the number of features improves generalization depends on various factors:

1. **Relevance of Features:** The impact of feature reduction on generalization depends on the relevance of the features being removed. If the features being eliminated contain noise or are irrelevant to the prediction task, then removing them can improve generalization by reducing the risk of overfitting.
2. **Information Loss:** Removing features can result in the loss of valuable information. If important features are removed, the model may lose its ability to capture the true underlying patterns in the data, leading to decreased generalization performance.
3. **Curse of Dimensionality:** In high-dimensional spaces, reducing the number of features can help mitigate the curse of dimensionality, where the data becomes sparse and computational requirements increase. In such cases, feature reduction can improve generalization.
4. **Model Complexity:** The complexity of the learning algorithm plays a role. Some models, such as decision trees or linear models, may benefit from feature reduction, while others, like deep neural networks, may be capable of automatically learning to ignore irrelevant features.
5. **Sample Size:** The amount of available training data also matters. In situations with limited data, reducing the number of features can help avoid overfitting, but in larger datasets, more features may be accommodated without overfitting.
6. **Feature Engineering:** Effective feature engineering, including feature selection and dimensionality reduction techniques, can help identify and retain the most informative features while discarding less useful ones.

In practice, the impact of feature reduction on generalization should be evaluated through cross-validation or other validation techniques. The goal is to strike a balance between retaining enough informative features to capture essential patterns and removing noisy or irrelevant features to prevent overfitting. The choice of which features to keep or remove should be driven by a thorough understanding of the data and the specific machine learning task at hand.

(f) To get better generalization, should we use the train set or the test set to tune our hyperparameters?

To achieve better generalization, you should use the training set to tune your hyperparameters, not the test set. Here's why:

1. **Test Set's Role:** The primary purpose of the test set is to evaluate the final model's performance on data it has never seen before. It serves as an independent measure of how well your model generalizes to unseen examples. If you use the test set for hyperparameter tuning, you risk introducing bias into the evaluation, as the model might have "seen" the test data indirectly through the tuning process.
2. **Overfitting to Test Set:** If you repeatedly evaluate different hyperparameter settings on the test set and choose the one that performs best on the test set, you are effectively tuning your model to perform well on that specific set of data. This can lead to overfitting to the test set, where the model may not generalize well to new, unseen data.
3. **Data Leakage:** Using the test set for hyperparameter tuning can inadvertently leak information about the test set into the model and hyperparameters, compromising the integrity of the test set as an unbiased evaluation metric.

Instead, the typical approach is to split your data into three separate sets: a training set, a validation set, and a test set. Here's how you can use them:

1. **Training Set:** Use the training set to train and fit different models with various hyperparameter settings. This is where the model learns from the data.
2. **Validation Set:** Use the validation set to evaluate the performance of different models with different hyperparameter settings. You can choose the hyperparameters that yield the best performance on the validation set.
3. **Test Set:** After you have chosen the best hyperparameters using the validation set, you should only use the test set once, at the end of the modeling process, to assess how well your final model generalizes to new, unseen data.

By following this approach, you ensure that the test set remains an unbiased and independent measure of your model's performance on unseen data, and you avoid overfitting the model to the test set.

(g) True or False: The training error of a function on the training set provides an overestimate of the true error of that function.

False. The training error of a function on the training set provides an **underestimate** of the true error of that function, not an overestimate. Here's why:

1. **Training Error:** The training error is computed by evaluating how well a model or function fits the data it was trained on. It measures the discrepancy between the model's predictions and the actual target values within the training dataset. In other words, it quantifies how well the model has learned to represent the training data.

2. **True Error:** The true error, often referred to as the generalization error, is an estimate of how well the model will perform on unseen data, which is data that the model has not been trained on. It represents the model's ability to generalize its learned patterns to new, unseen examples.

Since the training error is calculated based on the same data that the model was trained on, it reflects the model's performance on that specific dataset. Because the model has learned to fit the training data closely, the training error is typically lower than the true error. In other words, the training error underestimates the model's performance on unseen data.

The true error accounts for the model's ability to generalize to new, unseen examples and considers how well the model performs beyond the training dataset. Evaluating a model's true error is crucial for assessing its overall performance and its ability to make accurate predictions on data it has not encountered during training.

Q7)

(a) In a K-NN classification problem assume that the distance measure is not explicitly specified to you. Instead, you are given a "black box" where you input a set of instances  $P_1, P_2, \dots, P_n$  and a new example  $Q$ , and the black box outputs the nearest neighbor of  $Q$ , say  $P_i$ , and its corresponding class label  $C_i$ . Is it possible to construct a K-NN classification algorithm (w.r.t. the unknown distance metrics) based on this black box alone? If so, how? If not, why not?

It is not possible to construct a k-nearest neighbors (K-NN) classification algorithm based solely on the black box that provides the nearest neighbor and its corresponding class label without knowledge of the distance metric being used. The reason for this limitation is that the K-NN algorithm relies fundamentally on a distance metric to measure the similarity or distance between data points. Without knowing the distance metric, you cannot effectively determine which data points are the nearest neighbors.

In the standard K-NN algorithm, the steps involved:

1. **Selecting a Distance Metric:** The choice of distance metric, such as Euclidean distance, Manhattan distance, or others, determines how the similarity between data points is calculated.
2. **Calculating Distances:** For a given query point  $Q$ , the algorithm calculates the distance between  $Q$  and all other data points in the dataset.
3. **Sorting by Distance:** The algorithm then sorts the data points based on their distances to  $Q$ , identifying the k-nearest neighbors.
4. **Classifying:** Finally, it assigns a class label to  $Q$  based on the majority class among its k-nearest neighbors.

Without knowing the distance metric used in the black box, you cannot perform the crucial step of calculating distances between the query point and other data points. Therefore, constructing a K-NN classification algorithm based solely on the black box's outputs is not feasible.

The choice of the appropriate distance metric is essential in K-NN, as different metrics are suitable for different types of data and problem domains. The selection of a distance metric depends on the nature of the data, its distribution, and the problem's requirements. Without information about the distance metric, you cannot reliably apply the K-NN algorithm.

(b) If the black box returns the  $j$  nearest neighbors (and their corresponding class labels) instead of the single nearest neighbor (assume  $j \neq k$ ), is it possible to construct a K-NN classification algorithm based on the black box? If so how? If not why not?

Yes. If the black box returns the  $j$  nearest neighbors (and their corresponding class labels), where  $j \neq k$  (where  $k$  is the desired number of nearest neighbors in the K-NN algorithm), you can still construct a K-NN classification algorithm based on the black box, but you need to adjust your approach to handle this difference in the number of neighbors. Here's how you can do it:

1. **Observe the  $j$  Nearest Neighbors:** Use the black box to obtain the  $j$  nearest neighbors ( $P_1, P_2, \dots, P_j$ ) and their corresponding class labels ( $C_1, C_2, \dots, C_j$ ) for a given query point  $Q$ . These neighbors are based on the unknown distance metric used by the black box.
2. **Count Class Frequencies:** Count the frequencies of each class label among the  $j$  nearest neighbors. Keep track of how many neighbors belong to each class.
3. **Majority Voting:** Assign a class label to the query point  $Q$  based on majority voting among the class labels of the  $j$  nearest neighbors. In other words, choose the class label that occurs most frequently among these neighbors.
  - If  $j = k$ , you will be performing traditional K-NN with  $k$  neighbors.
  - If  $j > k$ , you can still use majority voting among the  $j$  neighbors to classify  $Q$ , but you should consider whether this might lead to over-smoothing or excessive influence from a larger number of neighbors.

- If  $j < k$ , you will effectively be using a smaller neighborhood size than  $k$ . Keep in mind that reducing the number of neighbors may lead to increased sensitivity to noise and potentially less robust classification.

It's important to note that this approach adapts the K-NN algorithm to the behavior of the black box that returns  $j$  neighbors. However, the effectiveness of this modified K-NN algorithm will depend on the specific characteristics of the data and the black box's behavior. You may need to experiment with different values of  $j$  and evaluate the algorithm's performance using appropriate validation techniques to determine the optimal number of neighbors for your particular problem.