## 1. Bias, Variance, Underfitting, and Overfitting

These terms describe the common issues encountered when training machine learning models, particularly concerning the model's ability to learn from the training data and generalize to unseen data.    &#8239;

### i) Bias

**Bias** is the error introduced by approximating a real-world problem, which may be complex, by a much simpler model.

- **Definition**: It measures how far the average prediction of a model is from the correct value. A high-bias model makes strong assumptions about the form of the target function, often leading to **underfitting**.
- **Effect**: High bias means the model is too simple and consistently misses the true relationship between features and the target variable.
- **Example**: Using a **linear regression model** to fit data that has a **non-linear, parabolic relationship**. The straight line will fail to capture the curvature, resulting in a high bias and large error on both training and test data.

### ii) Variance

**Variance** is the model's sensitivity to small fluctuations or changes in the training data.

- **Definition**: It measures how much the predictions for a given data point vary if different training datasets were used. A high-variance model pays too much attention to the training data, including the noise, leading to **overfitting**.
- **Effect**: High variance means the model performs very well on the training data but poorly on unseen test data because it has essentially *memorized* the training set.
- **Example**: Using a **high-degree polynomial regression model** to fit a small dataset. The model may create a complex curve that perfectly passes through every training point, but this complex curve will wildly change shape if a few training points are altered, indicating poor generalization.

### iii) Underfitting and Overfitting

**Underfitting** and **Overfitting** describe how well a model has learned the relationship between features and the target variable. 🔗 🔗

| Feature | Underfitting (High Bias) 🔗 🔗 | Overfitting (High Variance) 🔗 🔗 |
|---|---|---|
| **Model Complexity** | Too Simple 🔗 | Too Complex 🔗 |
| **Training Error** | High 🔗 | Very Low 🔗 |
| **Test/Generalization Error** | High 🔗 | High (significantly higher than training error) 🔗 |
| **Problem** | Model failed to capture the underlying pattern of the data (didn't learn enough). 🔗 🔗 | Model captured noise and random fluctuations in the training data (learned too much). 🔗 🔗 |
| **Analogy** | A student who didn't study the material well enough. 🔗 | A student who memorized only the practice exam answers and can't answer new questions. 🔗 |
| **Solution** | Use a more complex model, add more features, reduce regularization. 🔗 | Use a simpler model, use more training data, add regularization, feature selection. 🔗 |

⊞ Export to Sheets                                                    ⬚

## 2. Differentiate between Lasso Regression and Ridge Regression

Both Lasso (Least Absolute Shrinkage and Selection Operator) and Ridge Regression are types of **regularized linear regression** that address overfitting by adding a penalty term to the standard Ordinary Least Squares (OLS) cost function. 🔗

| Feature | Ridge Regression (L2 Regularization) 🔗 | Lasso Regression (L1 Regularization) 🔗 |
|---|---|---|
| **Penalty Term** | Sum of squares of the coefficients ($\sum_{j=1}^{p} \beta_j^2$) 🔗 | Sum of absolute values of the coefficients ($\sum_{j=1}^{p}$ |
| **Cost Function** | $$J(\beta) = \text{MSE}(\beta) + \lambda \sum_{j=1}^{p} \beta_j^2$$ 🔗 | $J(\beta) = \text{MSE}(\beta) + \lambda \sum_{j=1}^{p}$ |
| **Impact on Coefficients** | Shrinks coefficients towards zero, but never exactly to zero. 🔗 | Shrinks coefficients towards zero and can drive some coefficients *exactly* to zero. 🔗 |
| **Feature Selection** | Does not perform inherent feature selection. 🔗 | Performs automatic **feature selection** by eliminating irrelevant features (setting their |

coefficients to zero).   🔗

| | | |
|---|---|---|
| **Use Case** | Effective when all features are relevant and need to be kept. 🔗 | Effective when there are many features, but only a few are truly important (useful for sparse models). 🔗 |

⊞ Export to Sheets       ▢

---

## 3. Explain Gradient Descent Algorithm with Example

The **Gradient Descent** algorithm is an iterative, first-order optimization algorithm used to find the values of a function's parameters (coefficients/weights) that minimize its cost (or error) function. 🔗

### How it Works

1.  **Initialization**: Start with an arbitrary set of parameter values (often initialized randomly). 🔗

2.  **Calculate Gradient**: Compute the gradient (the partial derivatives) of the cost function with respect to each parameter. The gradient indicates the direction of the *steepest ascent* of the cost function. 🔗

3.  **Update Parameters**: Move in the opposite direction of the gradient (the direction of steepest *descent*) by a certain amount, determined by the **learning rate ($\alpha$)**. 🔗

    *   The update rule for a parameter $\theta_j$ is:

    $$\theta_j := \theta_j - \alpha \frac{\partial}{\partial \theta_j} J(\theta)$$

    where $J(\theta)$ is the cost function. 🔗

4.  **Iteration**: Repeat steps 2 and 3 until the cost function converges to a minimum (i.e., the cost is no longer decreasing significantly). 🔗

### Example: Minimizing a Simple Quadratic Function

Consider a simple cost function (parabola) $J(w) = w^2$. We want to find the value of $w$ that minimizes $J(w)$ (which is $w = 0$).

1.  **Gradient**: The derivative is $\frac{d}{dw} J(w) = 2w$.

2.  **Update Rule**: $w := w - \alpha(2w)$.

3.  **Iteration**:

    *   Start with $w = 5$ and $\alpha = 0.1$.

    *   **Iteration 1**:

        *   Gradient is $2(5) = 10$.

- $w = 5 - 0.1(10) = 5 - 1 = 4.$
  - **Iteration 2**:
    - Gradient is $2(4) = 8.$
    - $w = 4 - 0.1(8) = 4 - 0.8 = 3.2.$
  - The value of $w$ continuously decreases, moving closer to the minimum at $w = 0$ with each step. 🔗

---

## 5. Short Note on Training Error and Generalization Error

These two metrics are crucial for evaluating a machine learning model and diagnosing problems like underfitting and overfitting. 🔗

### Training Error

- **Definition**: The error (or loss) of the model when measured on the **training data**—the data the model was actually trained on. 🔗

- **Significance**: A **high training error** often indicates that the model is too simple for the data (**high bias/underfitting**) and has failed to capture the underlying pattern. A very low training error is expected for a well-trained model, but an extremely low error (near zero) might be a sign of overfitting. 🔗

### Generalization Error (Test Error)

- **Definition**: The error of the model when measured on **unseen data** (the test set) that was not used during the training process. It represents how well the model can generalize its learnings to new, real-world data. 🔗

- **Significance**: This is the most important metric for a model's real-world utility. 🔗
  - If **Training Error is low** and **Generalization Error is significantly higher**, the model suffers from **high variance/overfitting**. 🔗
  - If **both errors are high**, the model suffers from **high bias/underfitting**. 🔗

The goal in machine learning is to find a model that has **low training error** *and* a **Generalization Error** that is close to the training error. 🔗

---

## 6. Explain in Brief Importance of Evaluation Metrics

**Evaluation metrics** are quantitative measures used to assess a machine learning model's performance, quality, and effectiveness on a given task. 🔗

### Importance

1. **Model Selection**: Metrics allow for **objective comparison** between different models (e.g., comparing a Linear Regression model's RMSE to a Random Forest model's RMSE) to choose

the best-performing one. 🔗

2. **Performance Assessment**: They provide a single, interpretable score (e.g., accuracy, precision, $R^2$) that summarizes the model's success, making it easy to understand how well the model works. 🔗

3. **Diagnosis of Model Issues**: By analyzing metrics like Training Error vs. Generalization Error (see Q5), one can diagnose issues like **underfitting** or **overfitting** and adjust the model complexity or regularization accordingly. 🔗

4. **Business Alignment**: Choosing the right metric ensures the model is optimized for the **specific business goal**. For example, in fraud detection, **Precision** might be more important than **Accuracy** to minimize false alarms. 🔗

---

## 7. Find Regression Line and Estimate Cost

Given the data for Memory Capacity $X$ (in GB) and Cost $Y$ (in ) :| $X$ (GB) | 2 | 4 | 8 || :--- | :--- | :--- | :--- | :--- || $Y$ ($) | 12 | 16 | 28 |

The linear regression line is $Y = aX + b$. 🔗

**i) Find Regression line $Y = aX + b$ using the least square method.**

The formulas for the coefficients $a$ and $b$ using the least square method are:

$$a = \frac{n \sum(XY) - \sum X \sum Y}{n \sum(X^2) - (\sum X)^2} \quad \text{and} \quad b = \bar{Y} - a\bar{X}$$

where $n$ is the number of data points. Here $n = 3$. 🔗

**Calculation Table**

| $X$ | $Y$ | $XY$ | $X^2$ |
|---|---|---|---|
| 2 | 12 | 24 | 4 |
| 4 | 16 | 64 | 16 |
| 8 | 28 | 224 | 64 |
| $\sum$ | 14 | 56 | 312 |

⊞ Export to Sheets 🗔

$\sum X = 14, \sum Y = 56, \sum(XY) = 312, \sum(X^2) = 84.$

**Calculate $a$ (Slope)**

$$a = \frac{3(312) - (14)(56)}{3(84) - (14)^2} = \frac{936 - 784}{252 - 196} = \frac{152}{56} \approx 2.714$$

**Calculate $b$ (Y-intercept)**

First, calculate the means: $\bar{X} = \frac{14}{3} \approx 4.667$, $\bar{Y} = \frac{56}{3} \approx 18.667$.

$$b = \bar{Y} - a\bar{X} = 18.667 - 2.714(4.667) = 18.667 - 12.664 \approx 6.003$$

Rounding to three decimal places, $a \approx 2.714$ and $b \approx 6.003$.

The Regression Line is:

$$Y = 2.714X + 6.003$$

🔗

**ii) Estimate the cost of 32 GB RAM using the line equation.**

Substitute $X = 32$ into the regression equation:

$$Y = 2.714(32) + 6.003$$

$$Y = 86.848 + 6.003$$

$$Y = 92.851$$

🔗

The estimated cost of 32 GB RAM is **$92.85**. 🔗

---

## 8. Describe the Bias-Variance Trade-Off and its relationship to Underfitting and Overfitting

**Bias-Variance Trade-Off**

The **Bias-Variance Trade-Off** is a central concept in machine learning that states that the effort to decrease one (Bias or Variance) will generally lead to an increase in the other. The total expected error of a model can be decomposed into three parts: 🔗

$$\text{Error} = \text{Bias}^2 + \text{Variance} + \text{Irreducible Error}$$

The **Irreducible Error** is the noise in the data itself and cannot be reduced by any model. The goal is to find a model complexity (or sweet spot) that minimizes the sum of **Bias** and **Variance**. 🔗

**Relationship to Underfitting and Overfitting**

The trade-off directly relates model complexity to the errors:

1. **Low Complexity Models (High Bias):**
   - **Result: Underfitting.** 🔗
   - **State**: The model is too simple. It has high bias because it cannot represent the underlying pattern and low variance because it's consistent across different datasets. 🔗
   - **Error**: High error on both training and test data. 🔗

2. **High Complexity Models (High Variance):**
   - **Result: Overfitting.** 🔗
   - **State**: The model is too complex. It has low bias because it fits the training data almost perfectly, but high variance because its predictions change drastically with a small change in the training data. 🔗
   - **Error**: Very low training error but high test/generalization error. 🔗

3. **Optimal Model (Balanced):**
   - **Result**: Optimal generalization. 🔗
   - **State**: A model of intermediate complexity that achieves the best balance between bias and variance, minimizing the total error on unseen data. 🔗

---

## 9. Explain the advantages of RMSE over MSE as an evaluation metric

Both Root Mean Square Error (**RMSE**) and Mean Square Error (**MSE**) are measures of the average magnitude of the errors (the difference between the actual and predicted values). 🔗

| Metric | Formula |
| --- | --- |
| **MSE** | $$\text{MSE} = \frac{1}{n}\sum_{i=1}^{n}(Y_i - \hat{Y}_i)^2$$ 🔗 |
| **RMSE** | $$\text{RMSE} = \sqrt{\frac{1}{n}\sum_{i=1}^{n}(Y_i - \hat{Y}_i)^2}$$ 🔗 |

▦ Export to Sheets      🗔

RMSE is generally preferred over **MSE** because of the following key advantage:

**Advantage: Interpretability (Same Units)**

- **RMSE is easier to interpret** because it is in the **same units as the target variable ($Y$).**
  Since the errors are squared in MSE, the resulting unit is the *square* of the target variable's
  unit, which is less intuitive.  🔗

  - **Example:** If the target variable $Y$ is cost in dollars ($\$$), MSE will be in dollars-squared ($\$^2$
    ), while RMSE will be in dollars ($\$$). This means the RMSE value can be directly compared
    to the average magnitude of the predictions or the target values.  🔗

**Shared Feature (Sensitivity to Outliers)**

Both RMSE and MSE involve squaring the errors, which gives disproportionately **more weight to
large errors** (outliers). This makes both metrics good for scenarios where large errors are
undesirable.  🔗

---

## 10. What do you mean by least square method? Explain least square method in the context of linear regression.

**Least Square Method**

The **Least Square Method (LSM)** is a standard statistical technique used to find the best-fitting
line (or curve) for a set of data points. The term "least squares" comes from its objective: **to
minimize the sum of the squares of the differences** between the observed data points and
the values predicted by the model.  🔗

**LSM in Context of Linear Regression**

In **Simple Linear Regression**, we model the relationship between one independent variable ($X$)
and one dependent variable ($Y$) using a straight line: $\hat{Y} = aX + b$.  🔗

- **Error (Residual):** For any data point ($X_i, Y_i$), the difference between the actual value $Y_i$
  and the predicted value $\hat{Y}_i$ is the residual, $e_i = Y_i - \hat{Y}_i$.  🔗

- **Objective Function:** The LSM defines the best-fit line as the one that minimizes the **Sum of
  Squared Residuals (SSR)** or **Cost Function** $J(a, b)$:

$$J(a, b) = \sum_{i=1}^{n} e_i^2 = \sum_{i=1}^{n}(Y_i - \hat{Y}_i)^2 = \sum_{i=1}^{n}(Y_i - (aX_i + b))^2$$

🔗

- **Finding the Coefficients:** To find the optimal values for the slope ($a$) and the y-intercept ($b$
  ), we use calculus. We take the partial derivative of the cost function $J(a, b)$ with respect to
  $a$ and $b$, set them both to zero, and solve the resulting system of normal equations. This
  yields the standard formulas for $a$ and $b$:

$$a = \frac{n\sum(XY) - \sum X \sum Y}{n\sum(X^2) - (\sum X)^2} \quad \text{and} \quad b = \bar{Y} - a\bar{X}$$

The line derived from the least square method is called the **line of best fit** because it results in the smallest possible sum of squared vertical distances from the data points to the line. 🔗

---

## 11. Define different regression models.

**Regression models** are a set of supervised learning techniques used to predict a **continuous** target variable based on one or more independent predictor variables. They model the relationship between the features and the target. 🔗

Here are some different types of regression models:

1. **Linear Regression:**

   - **Description**: The simplest form, which assumes a **linear relationship** between the input variable(s) ($X$) and the output variable ($Y$). The equation is $Y = a_1 X_1 + a_2 X_2 + \cdots + b$. 🔗

   - **Types**: Simple (one $X$) and Multiple (multiple $X$'s). 🔗

2. **Polynomial Regression:**

   - **Description**: Used when the relationship between $X$ and $Y$ is **non-linear**. It models the relationship using an $n^{th}$-degree polynomial. The equation is $Y = b_0 + b_1 X + b_2 X^2 + \cdots + b_n X^n$. 🔗

3. **Ridge Regression (L2 Regularization):**

   - **Description**: A form of linear regression that adds an **L2 penalty term** to the cost function to prevent overfitting by shrinking the coefficient values towards zero. 🔗

4. **Lasso Regression (L1 Regularization):**

   - **Description**: Similar to Ridge Regression, but it adds an **L1 penalty term**. This term has the additional ability to force the coefficients of irrelevant features to become *exactly* zero, thus performing **automatic feature selection**. 🔗

5. **Elastic Net Regression:**

   - **Description**: A hybrid model that combines both L1 (Lasso) and L2 (Ridge) penalty terms in its cost function. It is useful for getting the regularization power of Ridge and the feature selection capability of Lasso. 🔗

---

## 12. Predict Expenditure for 6th Month

Given the data for month $X$ and expenditure $Y$ ($n = 5$):

| $X$ (month) | $Y$ (expenditure) |
| --- | --- |
| 1 | 12 |
| 2 | 19 |

| 3 | 29 |
|---|----|
| 4 | 37 |
| 5 | 45 |

⊞ Export to Sheets   ▢

We will use **Simple Linear Regression** ($Y = aX + b$) to predict the expenditure for the 6th month ($X = 6$). 🔗

**Calculation Table**

| $X$ | $Y$ | $XY$ | $X^2$ |
|-----|-----|------|-------|
| 1 | 12 | 12 | 1 |
| 2 | 19 | 38 | 4 |
| 3 | 29 | 87 | 9 |
| 4 | 37 | 148 | 16 |
| 5 | 45 | 225 | 25 |
| $\Sigma$ | 15 | 142 | 510 |

⊞ Export to Sheets   ▢

$$\sum X = 15, \sum Y = 142, \sum(XY) = 510, \sum(X^2) = 55.$$

**Calculate $a$ (Slope)**

$$a = \frac{n \sum(XY) - \sum X \sum Y}{n \sum(X^2) - (\sum X)^2}$$

$$a = \frac{5(510) - (15)(142)}{5(55) - (15)^2} = \frac{2550 - 2130}{275 - 225} = \frac{420}{50} = 8.4$$

**Calculate $b$ (Y-intercept)**

Means: $\bar{X} = \frac{15}{5} = 3, \bar{Y} = \frac{142}{5} = 28.4$.

$$b = \bar{Y} - a\bar{X} = 28.4 - 8.4(3) = 28.4 - 25.2 = 3.2$$

The Regression Line is:

$$Y = 8.4X + 3.2$$

Prediction for 6th Month

Prediction for 6th Month

Substitute $X = 6$ into the equation: &#x1f517;

$$Y_6 = 8.4(6) + 3.2$$

$$Y_6 = 50.4 + 3.2$$

$$Y_6 = 53.6$$

The predicted expenditure for the 6th month is **53.6**. &#x1f517;

---

## 13. Write a Short Note on Stochastic Gradient Descent Algorithms

**Stochastic Gradient Descent (SGD)**

**Stochastic Gradient Descent (SGD)** is an efficient and often preferred variation of the standard Gradient Descent (Batch Gradient Descent) algorithm. It is particularly useful when dealing with **large datasets**. &#x1f517;

**Key Differences from Standard Gradient Descent**

- **Batch Size**: Standard Gradient Descent calculates the gradient using **all** $n$ training examples (the entire dataset or "batch") in each iteration. SGD, in contrast, updates the model parameters by calculating the gradient using **only a single randomly selected training example** in each iteration. &#x1f517;

- **Update Frequency**: SGD performs **more frequent updates** than standard Gradient Descent (one update per data point vs. one update per full dataset). &#x1f517;

- **Convergence**:

  - The path to the minimum in SGD is **noisy/stochastic** because it's based on only one sample's error. It does not precisely converge to the global minimum but rather **oscillates** around it. &#x1f517;

  - Despite the oscillation, this process is much **faster** for large datasets, as it avoids processing the entire dataset for every parameter update. &#x1f517;

**Advantages**

- **Computational Efficiency**: Significantly reduces computation time for very large datasets. &#x1f517;

- **Handles Redundancy**: If data points are redundant, SGD still makes progress by using only one. &#x1f517;

- **Potential for Global Minimum**: The oscillation can sometimes help escape shallow local minima in complex, non-convex cost functions. &#x1f517;

**Mini-Batch Gradient Descent** is a compromise between Batch GD and SGD, where the gradient

is calculated using a small subset (or mini-batch) of $k$ training examples $(1 < k < n)$. 🔗

---

## 14. What are different techniques to reduce underfitting?

**Underfitting** occurs when a model is too simple to capture the underlying patterns in the training data (high bias). Techniques to reduce underfitting generally involve increasing the model's complexity or its ability to learn from the data. 🔗

1. **Increase Model Complexity**:

   - **Use a More Complex Model**: Switch from a simpler model (e.g., Linear Regression) to a more powerful model (e.g., Polynomial Regression, Support Vector Machines, Neural Networks with more layers/neurons). 🔗

   - **Increase Degrees of Freedom**: For polynomial regression, increase the degree of the polynomial to better fit non-linear data. 🔗

2. **Add More Features (Feature Engineering)**:

   - The model may be underfitting because the current set of features is insufficient to explain the relationship with the target variable. 🔗

   - **Generate new features** (e.g., combining existing features, taking logarithms, or generating interaction terms like $X_1 \cdot X_2$). 🔗

3. **Reduce Regularization**:

   - **Decrease the regularization parameter ($\lambda$)** in techniques like Ridge or Lasso Regression. High regularization forces coefficients towards zero, simplifying the model; reducing it allows the model to become more complex and better fit the training data. 🔗

4. **Increase Training Time (for iterative models)**:

   - For models trained iteratively (like Neural Networks), the model may simply need to be trained for **more epochs** (more passes over the data) to allow the optimization algorithm to converge to a better solution. 🔗

5. **Remove Noise from the Data**:

   - While not always a common solution, ensuring the data is clean and does not contain gross errors or outliers that mislead the simple model can help it learn the core relationship more effectively. 🔗

---

## 15. Short Note on MAE, RMSE, and $R^2$

These are common **evaluation metrics** used for **regression models**. 🔗

### i) MAE (Mean Absolute Error)

- **Definition**: The average of the absolute differences between the predicted values ($\hat{Y}$) and the actual values ($Y$). 🔗

$$MAE = \frac{1}{n} \sum_{i=1}^{n} |Y_i - \hat{Y}_i|$$

- **Advantages**: Highly **interpretable** as it is in the same units as the target variable. It is also less sensitive to outliers compared to MSE/RMSE because it does not square the errors. 🔗

## ii) RMSE (Root Mean Square Error)

- **Definition**: The square root of the average of the squared differences between predicted and actual values. 🔗

$$RMSE = \sqrt{\frac{1}{n} \sum_{i=1}^{n} (Y_i - \hat{Y}_i)^2}$$

- **Advantages**: Also in the same units as the target variable, making it interpretable. Because it squares the errors, it **penalizes large errors (outliers) heavily**, making it a good metric when avoiding large prediction errors is crucial. 🔗

## iii) $R^2$ (Coefficient of Determination)

- **Definition**: Measures the **proportion of the variance** in the dependent variable ($Y$) that is predictable from the independent variable(s) ($X$). It's a measure of the **goodness of fit**. 🔗

$$R^2 = 1 - \frac{\text{Sum of Squared Residuals (SSR)}}{\text{Total Sum of Squares (SST)}}$$

- **Interpretation**: $R^2$ values range from 0 to 1. An $R^2$ of 0.80 means that 80% of the variance in $Y$ is explained by the model. A higher $R^2$ indicates a better fit. 🔗

- **Limitation**: $R^2$ always increases when new features are added, even if those features are irrelevant, leading to the use of **Adjusted $R^2$** for multiple regression. 🔗

---

Would you like a deeper explanation of any specific model or concept, such as **Stochastic Gradient Descent** or the **Bias-Variance Trade-Off**?