John Smith, Ph.D.

Department of Computer Science
University Name

Email: email@university.edu
Website: www.university.edu

July 19, 2025

# Introduction to Overfitting and Underfitting

## Understanding Model Performance

In machine learning, achieving a well-performing model is crucial, and two common issues that can significantly hinder this performance are **overfitting** and **underfitting**. Understanding these concepts is essential for developing models that generalize well when faced with new, unseen data.

# 1. Overfitting

- **Definition:** Occurs when a model learns the training data too well, capturing noise and random fluctuations instead of the underlying pattern.
- **Characteristics:**
  - High accuracy on training data
  - Poor performance on validation/testing data
- **Example:** A polynomial regression model fitting a dataset with only a few points perfectly by using a high-degree polynomial. While it might show high accuracy on the training set, it performs poorly on test data due to its complexity.

### Illustration

(Imagine a graph where a complex curve perfectly fits all training data points but diverges wildly outside this data range.)

## 2. Underfitting

- **Definition:** Happens when a model is too simple to capture the underlying pattern in the data, leading to poor performance on both training and testing datasets.
- **Characteristics:**
    - Low accuracy on both training and testing data
    - Simplistic model assumptions (e.g., a linear model applied to a non-linear dataset)
- **Example:** Using a simple linear regression model on a dataset that follows a quadratic trend, resulting in large errors.

### Illustration
(Imagine a graph where a straight line tries to fit a parabolic curve but misses the data points entirely.)

# Key Points to Emphasize

- **Trade-Off:** There is a trade-off between overfitting and underfitting known as the **bias-variance trade-off**:
    - **Bias:** Errors due to overly simplistic assumptions in the learning algorithm, leading to underfitting.
    - **Variance:** Errors due to excessive complexity in the model, leading to overfitting.
- **Performance Measurement:** Always evaluate your model's performance using cross-validation techniques to identify overfitting or underfitting.
- **Goal:** The primary goal in machine learning is to find a sweet spot where the model achieves the best generalization on unknown data.

# Conclusion and Next Steps

By mastering the concepts of overfitting and underfitting, you can construct models that fit the training data well while maintaining good predictive accuracy on new datasets, ensuring robust model performance.

**Next Slide Preview:** In the next slide, we will delve deeper into understanding overfitting, its definitions, characteristics, and the scenarios in which it occurs.

# What is Overfitting? - Definition

## Definition

Overfitting occurs in machine learning when a model learns not only the underlying patterns in the training data but also the noise and outliers. Consequently, the model performs exceptionally well on the training dataset but fails to generalize to new, unseen data.

# What is Overfitting? - Characteristics

- **High Training Accuracy**: Very high accuracy on training data.
- **Poor Validation/Test Accuracy**: Underperforms on validation/test datasets.
- **Complex Model Structures**: Often results from overly complex architectures.
- **Sensitivity to Noise**: Adapts too much to every data point, including outliers.

- **Insufficient Training Data**: Limited or non-representative dataset.
- **Excessive Complexity**: Too many parameters relative to data amount.
- **Lack of Regularization**: Absence of techniques to constrain model complexity.

## What is Overfitting? - Key Concepts

- **Bias-Variance Trade-off:** Overfitting increases variance and decreases bias.
- **Validation Importance:** Cross-validation can help identify overfitting.
- **Preventive Techniques:**
    - Pruning (for decision trees)
    - Regularization (penalties for complexity)
    - Early Stopping (halting training at decline in validation accuracy)
    - Data Augmentation (increasing dataset size)

# What is Overfitting? - Example

Consider a polynomial regression model attempting to fit a dataset of points:

- **Underfitting:** Linear model failing to capture trend (high bias).
- **Ideal Fit:** Quadratic model capturing trend (balanced bias and variance).
- **Overfitting:** High-degree polynomial fitting training data perfectly, but performs poorly on new data.

```python
# Example of Overfitting in Python using Scikit-Learn
from sklearn.datasets import make_regression
from sklearn.model_selection import train_test_split
from sklearn.linear_model import LinearRegression, Ridge
import matplotlib.pyplot as plt

# Create dataset
X, y = make_regression(n_samples=100, noise=15)
# Split data
X_train, X_test, y_train, y_test = train_test_split(X, y, test_size=0.2,
    random_state=42)

# Fit a complex model
model = Ridge(alpha=1)  # Using Ridge for simplification
model.fit(X_train, y_train)

# Predict
```

# What is Overfitting? - Conclusion

By understanding overfitting, we can better tune our machine learning models for optimal performance, ensuring they are robust and capable of making accurate predictions on unseen data.

# What is Underfitting?

## Definition of Underfitting

Underfitting occurs when a machine learning model is too simple to capture the underlying patterns in the training data. It results in poor performance on both the training data and new, unseen data.

This situation often arises when a model has insufficient complexity to represent the relationships in the dataset.

# Underfitting vs. Overfitting

- **Underfitting:**
    - Model performs poorly on both training and test datasets.
    - Results from overly simplistic models, like linear regression on non-linear data.
    - Indicates a lack of learned patterns.
- **Overfitting:**
    - Model performs well on training data but poorly on unseen data (test dataset).
    - Caused by excessive complexity in the model, such as too many parameters or layers.
    - Indicates that the model has memorized training data rather than learned general patterns.

# Characteristics of Underfit Models

1. **High Bias:** Characterized by high bias, leading to systematic errors due to strong assumptions about the data.
2. **Low Complexity:** Simple models like linear regression on polynomial relationships often suffer from underfitting.
3. **Inadequate Feature Use:** Fails to capture key features that could enhance predictive accuracy.
4. **Poor Performance Metrics:** Evaluation metrics (e.g., Mean Squared Error, accuracy) show disappointing results both on training and validation datasets.
5. **Example:** A linear regression model fitting data points that follow a curved pattern results in a high error rate.

## Key Takeaways

- Aim for a balance between model complexity and data representation. - Address underfitting by increasing model complexity or using more features. - Regularly evaluate model performance using appropriate metrics.

## Understanding Overfitting and Underfitting

- **Overfitting:** Model learns training data too well, capturing noise.
  - Results in poor performance on unseen data.
- **Underfitting:** Model is too simple, failing to capture the underlying trend.
  - Results in poor performance on both training and testing datasets.

## Techniques for Detection

**1** **Performance Metrics**
   - Train/Test Split: Compare performance metrics (e.g., accuracy, loss) on training vs. validation sets.
   - **Example:** Training Accuracy: 95%, Validation Accuracy: 70% (indicative of overfitting).

**2** **Cross-Validation**
   - k-fold method: Train model k times, each time omitting one subset for validation.

**3** **Learning Curves**
   - Plot performance over iterations.
   - Overfitting Indicator: Large gap between training and validation performance.
   - Underfitting Indicator: Low performance on both.

## More Techniques

- **Regularization Techniques**
  - Add penalties for complexity during model training (e.g., L1/L2).
- **Visual Diagnosis**
  - Scatter plots and residual plots assess model fit.
  - Even residuals indicate a well-fitted model; patterns suggest overfitting/underfitting.

## Key Points to Emphasize

- High Training Accuracy and Low Validation Accuracy $\rightarrow$ Overfitting.
- Low Accuracy on both training and validation sets $\rightarrow$ Underfitting.
- Regular use of performance metrics and visualizations is crucial.

## Example Code Snippet

```python
from sklearn.model_selection import train_test_split
from sklearn.metrics import accuracy_score
from sklearn.linear_model import LogisticRegression

# Create data
X_train, X_val, y_train, y_val = train_test_split(X, y, test_size=0.2)
model = LogisticRegression()
model.fit(X_train, y_train)

# Evaluate performance
train_accuracy = accuracy_score(y_train, model.predict(X_train))
val_accuracy = accuracy_score(y_val, model.predict(X_val))

print(f'Train Accuracy: {train_accuracy}, Validation Accuracy: {val_accuracy}')
```

# Introduction to Overfitting

- Overfitting occurs when a model learns the training data too well, capturing noise rather than the underlying patterns.
- Results:
    - Excellent performance on training data
    - Poor performance on unseen data
- Leads to high variance problems.

# Negative Effects of Overfitting

1. **Poor Generalization**:
   - Definition: Ability to perform well on unseen data.
   - Impact: An overfitted model fails to predict accurately on test data.
   - Example: A flower classification model misclassifies new flowers.

2. **Increased Complexity**:
   - Definition: Creating overly complex models with too many parameters.
   - Impact: Greater computational cost and chances of noise misinterpretation.
   - Key Point: Balance model complexity with training data size is crucial.

3. **Unreliable Predictions**:
   - Definition: Overfitted model sensitivity to training data specifics.
   - Impact: Minor input changes cause large output deviations.
   - Example: Misleading medical diagnoses based on spurious patterns.

- Graphical Representation:
  - Training Data Curve: Perfectly fits training data points.
  - Validation/Test Data Curve: Shows poor performance and diverges from training accuracy.

# Preventing Overfitting

- **Strategies**:
  - Cross-Validation: Use k-fold cross-validation for diverse training sets.
  - Regularization: Apply L1 (Lasso) and L2 (Ridge) regularization to prevent complexity.
  - Simplifying the Model: Use feature selection or dimensionality reduction (e.g., PCA).

# Conclusion

- Understanding overfitting is crucial for robust machine learning models.
- Importance of balancing model complexity and training data.
- Always prioritize generalization over training accuracy for reliable predictions.

# Impacts of Underfitting - Understanding Underfitting

## Definition

Underfitting occurs when a machine learning model is too simple to capture the underlying patterns in the data. This results in inadequate learning from the training data, leading to poor predictive performance on both the training and test datasets.

1. **Poor Accuracy:**
   - Low accuracy on both training and test datasets.
   - Example: A linear regression model fitting a quadratic relationship produces suboptimal predictions.

2. **High Bias:**
   - Indicates strong assumptions that do not reflect data complexity.
   - Can lead to systematic errors in predictions.
   - Formula:
   $$\text{Bias}^2 = \text{Expected Prediction} - \text{True Value} \tag{1}$$

3. **Insufficient Complexity:**
   - Occurs when using overly simplistic algorithms.
   - Example: A shallow decision tree failing to capture data variation.

## Example Scenario

A model predicting house prices based solely on the number of bedrooms might ignore important features like location and size, resulting in inaccurate predictions due to underfitting.

## Key Points to Emphasize

- **Feedback Loop:** Poor predictions can lead to incorrect conclusions about data, perpetuating errors.
- **Detection:** Evaluate metrics like mean squared error (MSE) on training and test datasets; both will be high if underfitting is present.

# Techniques to Combat Overfitting - Introduction

## What is Overfitting?

Overfitting occurs when a model learns not only the underlying patterns in the training data but also the noise. This leads to poor performance on unseen data. Combatting overfitting is crucial for developing robust machine learning models.

**1** **Train with More Data**
- **Explanation:** More training examples help the model generalize better.
- **Example:** 10,000 samples vs. 1,000 samples.

**2** **Simpler Models**
- **Explanation:** Choosing less complex models reduces overfitting.
- **Example:** Linear regression vs. 5th-degree polynomial regression.

**3** **Cross-Validation**
- **Explanation:** Ensures model generalization by validating on different subsets.
- **Example:** K-fold cross-validation.

**4** **Early Stopping**
- **Explanation:** Halt training when validation performance starts degrading.
- **Example:** Monitor validation loss during training.

**5** **Regularization Techniques**
- **L1 Regularization (Lasso):**

$$J(\theta) = \text{Loss} + \lambda \sum_{i=1}^{n} |\theta_i| \tag{2}$$

  - **Benefit**: Encourages sparsity in the model.
- **L2 Regularization (Ridge):**

$$J(\theta) = \text{Loss} + \lambda \sum_{i=1}^{n} \theta_i^2 \tag{3}$$

  - **Benefit**: Distributes coefficients and reduces influence of any single feature.

**6** **Dropout (for Neural Networks)**
- **Explanation:** Randomly sets a fraction of input units to zero during training, preventing co-adaptation.
- **Example:** 0.5 dropout rate means half of the neurons are turned off during training.

**7** **Data Augmentation**
- **Explanation:** Enhances training set with modified existing data, increasing diversity.
- **Example:** Applying rotations and shifts to image data.

# Techniques to Combat Overfitting - Conclusion

## Conclusion

Reducing overfitting is key to creating models that perform well on unseen data. Employing a mix of the above techniques can significantly improve model generalization. Understanding and implementing these techniques is a crucial skill in machine learning.

# Overview of Regularization

- Regularization is used in machine learning to prevent overfitting.
- It introduces additional information or constraints into the model, effectively reducing complexity.
- Helps models generalize better to unseen data.

## Key Concepts

- **Overfitting**: Learning noise in training data, resulting in poor performance on new data.
- **Underfitting**: Model is too simple to capture the underlying trends in the data.

# Common Regularization Techniques

1. L1 Regularization (Lasso Regularization)
2. L2 Regularization (Ridge Regularization)

# L1 Regularization (Lasso Regularization)

- **Concept**: Adds a penalty equal to the absolute value of coefficients, shrinking some to zero.
- **Formula**:

$$J(\theta) = \text{Loss} + \lambda \sum_{i=1}^{n} |\theta_i| \tag{4}$$

- **Example**: In a model with multiple features, L1 can identify significant features by reducing less important ones to zero.

# L2 Regularization (Ridge Regularization)

- **Concept**: Adds a penalty equal to the square of the coefficients' magnitudes, reducing their size without driving some to zero.

- **Formula**:

$$J(\theta) = \text{Loss} + \lambda \sum_{i=1}^{n} \theta_i^2 \tag{5}$$

- **Example**: In linear regression, L2 helps smooth the learned function and reduces variance, improving predictive capabilities.

# Comparison of L1 and L2 Regularization

- **Sparsity**:
  - L1 leads to sparse solutions (zero coefficients).
  - L2 generally keeps all features in the model.
- **Performance**:
  - L1 is preferable when feature selection is crucial.
  - L2 is suitable when all features are believed to have influence.

# Key Points to Emphasize

- Choosing the right regularization depends on the problem and the dataset.
- Regularization enhances model generalization and prevents overfitting.
- Tuning the regularization parameter $\lambda$ is essential, often done via cross-validation.

# Code Snippet Example (Python)

```python
from sklearn.linear_model import Lasso, Ridge

# L1 Regularization (Lasso)
lasso_model = Lasso(alpha=0.1)
lasso_model.fit(X_train, y_train)

# L2 Regularization (Ridge)
ridge_model = Ridge(alpha=0.1)
ridge_model.fit(X_train, y_train)
```

# Cross-Validation - Overview

## What is Cross-Validation?

Cross-validation is a statistical technique used to assess how well a model generalizes to an independent dataset. It evaluates the model's ability to predict new data and helps to prevent overfitting.

## Why Use Cross-Validation?

- Prevents overfitting by ensuring the model captures underlying patterns rather than noise.
- Provides a more reliable estimate of model performance than a single train-test split.

## Cross-Validation - Process

### How Does Cross-Validation Work?

1. **Data Splitting**: The dataset is divided into several subsets or "folds".
2. **Training and Validation**: The model is trained on (k-1) folds and validated on the remaining fold.
3. **Repetition**: This process is repeated k times, with each fold used as the validation set once.

### Common Types of Cross-Validation

- **K-Fold Cross-Validation**: Divides the dataset into k subsets. Each fold is used for validation once.
- **Stratified K-Fold**: Ensures each fold has the same proportion of class labels as the original dataset, beneficial for imbalanced datasets.

# Cross-Validation - Example

## Example: K-Fold with k=5

Suppose you have 100 samples: - Split into 5 folds of 20 samples each.
- Train on 4 folds (80 samples), validate on remaining fold (20 samples).

## Cross-Validation Process

1. 1st Iteration: Train on Folds 1 to 4, Validate on Fold 5.
2. 2nd Iteration: Train on Folds 1, 2, 3, 5, Validate on Fold 4.
3. ... Repeat until all folds are validated.

# Cross-Validation - Key Benefits

## Key Points to Emphasize

- **Reduces Overfitting**: Validating on different subsets ensures performance isn't due to lucky training on a specific dataset.
- **Better Model Evaluation**: Offers a more reliable estimate of performance compared to a single train-test split.
- **Hyperparameter Tuning**: Useful in tuning hyperparameters by assessing performance across multiple folds.

## Cross-Validation - Code Example

### Code Snippet: K-Fold Cross-Validation in Python

```python
from sklearn.model_selection import KFold
from sklearn.model_selection import cross_val_score
from sklearn.ensemble import RandomForestClassifier
from sklearn.datasets import load_iris

# Load example dataset
data = load_iris()
X, y = data.data, data.target

# Initialize KFold
kf = KFold(n_splits=5)

# Initialize model
model = RandomForestClassifier()
```

# Cross-Validation - Conclusion

## Conclusion

Cross-validation is an essential tool in model evaluation that: - Provides better understanding of model performance.
- Prevents overfitting by testing the model on various unseen samples.
- Guides better model selection and hyperparameter tuning, leading to improved predictive performance.

## Brief Summary

Pruning techniques are employed in decision trees to reduce overfitting by simplifying the model. This involves removing parts of the tree that do not significantly contribute to predictive performance, enhancing generalization on unseen data.

# Understanding Overfitting in Decision Trees

- Overfitting occurs when a model captures noise rather than the underlying data pattern.
- Decision trees are prone to overfitting by becoming too complex and deep.
- Symptoms: High training accuracy but low validation performance.

# Pruning: A Solution to Overfitting

## Definition

Pruning is the process of removing tree sections that add little predictive value, simplifying the decision tree and reducing overfitting.

- Simplifies the model structure
- Enhances performance on unseen data

1. **Pre-Pruning (Early Stopping):**
   - Stop growing the tree based on specific criteria.
   - Example criteria:
     - Maximum depth of the tree
     - Minimum samples required to split a node
     - Minimum impurity decrease
2. **Post-Pruning:**
   - Grow the full tree first and remove less important branches.
   - Utilize validation datasets to evaluate the impact of branch removal.

# Key Points and Summary

- Pruning balances complexity and simplicity, crucial for preventing overfitting.
- An optimal pruned tree usually outperforms an unpruned tree on validation datasets.
- Both pre-pruning and post-pruning can be used based on model needs.

## Conclusion

Implementing pruning enhances decision tree robustness and predictive accuracy on unseen data.

# Feature Selection - Importance

## Importance of Feature Selection

Feature selection is a crucial step in the machine learning pipeline aimed at improving model performance by selecting only the most relevant features from the dataset.

- Helps combat overfitting and underfitting
- Reduces model complexity
- Highlights essential patterns in the data

1. **Overfitting vs. Underfitting**:
   - **Overfitting**: Learning noise instead of trends leads to poor performance on unseen data.
   - **Underfitting**: Too simple of a model fails to capture underlying trends.
2. **Reducing Dimensionality**:
   - Fewer features result in simpler models and reduced likelihood of overfitting.
3. **Noise Reduction**:
   - Irrelevant features introduce noise, making it difficult to uncover relationships.
   - Selecting relevant features minimizes this noise.

# Feature Selection - Techniques

- **Filter Methods**: Score features using statistical techniques.

### Example: Pearson's correlation

Evaluate correlation with the target variable.

```python
import pandas as pd
from sklearn.feature_selection import SelectKBest, f_classif

# Load your dataset
dataframe = pd.read_csv('data.csv')
X = dataframe.iloc[:, :-1] # Features
y = dataframe.iloc[:, -1]   # Target variable

# Select top K features
best_features = SelectKBest(score_func=f_classif, k=10)
fit = best_features.fit(X, y)
```

# Ensemble Methods

## Definition

Ensemble methods combine multiple models to improve prediction accuracy and robustness, effectively addressing overfitting and underfitting.

## Key Features

- Leverage strengths of various algorithms.
- Provide generalized predictions beyond individual models.

# Types of Ensemble Methods

1. **Bagging (Bootstrap Aggregating)**
   - Reduces variance by training multiple models on different subsets.
   - Example: Random Forest builds numerous decision trees on bootstrapped datasets.
   - Key Characteristics:
     - Models are trained independently.
     - Predictions are combined by averaging (regression) or voting (classification).

2. **Boosting**
   - Builds a sequential model that corrects errors of previous ones.
   - Example: AdaBoost, XGBoost.
   - Key Characteristics:
     - Models are trained sequentially.
     - Predictions are combined by weighting outputs based on performance.

# Key Points and Formulas

## Key Points

- Reduction of Overfitting: Both methods leverage diversity among models.
- Handling Bias-Variance Trade-off: Balances bias (boosting) and variance (bagging).
- Model Interpretability: Techniques available to comprehend ensemble decisions.

## Formulas

$$\text{Final Prediction} = \text{mode of (predictions } \{y_1, y_2, \ldots, y_n\}) \tag{6}$$

$$F(x) = \sum_{m=1}^{M} \alpha_m h_m(x) \tag{7}$$

## Understanding Underfitting

- **Definition**: Underfitting occurs when a model is too simple to capture underlying patterns in data.
- **Symptoms**: High bias, low variance; low accuracy on training data indicates failure to learn effectively.

**1** **Increase Model Complexity**
- Shift from simpler models (e.g., linear regression) to complex models (e.g., polynomial regression).
- Example: Use a polynomial model $y = ax^2 + bx + c$ instead of $y = mx + c$.
- **Illustration**: A linear fit on a parabolic dataset will underperform compared to a quadratic polynomial.

**2** **Feature Engineering**
- Create new features or modify existing ones to enrich input data.
- Techniques:
  - Interaction Terms: Combine features (e.g., $x_1 \times x_2$).
  - Polynomial Features: Add powers of features up to a certain degree.
- Example: In a housing dataset, create features like "bedrooms $\times$ bathrooms".

# Techniques to Address Underfitting - Part 2

**Increase Training Data**
- More data helps models learn intricate patterns, reducing underfitting.
- Methods:
  - Data Augmentation: For images, apply transformations like rotation, flipping, and zooming.
  - Collecting More Data: Use different sources or surveys for additional data points.
- **Impact**: Increases variability for better generalization.

**Reduce Regularization**
- Adjusting regularization parameters (e.g., in Lasso or Ridge regression) can improve model fitting.
- Effect: Too much regularization can lead to underfitting.
- Example: Lowering alpha in Ridge regression reduces penalty on weights, allowing more complexity.

**Use More Powerful Algorithms**
- Choose algorithms that capture complex patterns.
- Examples:
  - Decision Trees: Model non-linear relationships with deeper trees.
  - SVM with Non-linear Kernels: Find complex decision boundaries.

# Conclusion and Key Takeaways

## Conclusion

Successfully combating underfitting often requires a multi-faceted approach, combining these techniques to find a balance for effective learning and generalization to new data.

- **Key Takeaways:**
  - Identify and understand underfitting: High bias indicates a simplistic model.
  - Iterate on model complexity: Adjust architecture and features to capture data relationships.
  - Regularly evaluate model performance: Use training and validation sets to assess improvements continuously.

# Choosing the Right Model - Overview

## Understanding Overfitting and Underfitting

- **Overfitting:** High accuracy on training data but poor performance on unseen data due to learning noise.
- **Underfitting:** Low accuracy on both training and test datasets due to a model too simple to capture data complexity.

**1 Balance Model Complexity**
- Choose a model complexity that matches data characteristics.
- *Example:* Linear regression for size; polynomial regression for multiple features.

**2 Cross-Validation**
- Use k-fold cross-validation to ensure the model generalizes well.
- A consistently performing model in cross-validation can perform well on unseen data.

**3 Regularization Techniques**
- Apply L1 (Lasso) or L2 (Ridge) regularization to reduce overfitting.
- *Formula for L2 Regularization:*

$$J(\theta) = \text{Loss}(\theta) + \lambda \sum_{j=1}^{n} \theta_j^2 \tag{8}$$

# Choosing the Right Model - Additional Strategies

**4** **Feature Selection**
- Identify and select relevant features to reduce noise and overfitting.
- *Example:* Backward elimination or forward selection.

**5** **Model Comparison**
- Experiment with different models based on performance metrics.
- Different structures can yield varying results on the same data.

**6** **Learning Curves**
- Plot learning curves to visualize model performance.
- Helps diagnose underfitting or overfitting: both curves converging to high error indicates a need for a more complex model.

**7** **Ensemble Methods**
- Combine predictions from multiple models (e.g., Bagging, Boosting) to enhance accuracy and reduce overfitting variance.

## Conclusion

Choosing the right model is critical to achieving good predictive performance while managing overfitting and underfitting. By understanding your data and applying appropriate strategies, you can enhance the effectiveness and reliability of your models.

# Conclusion - Recap

## Importance of Managing Overfitting and Underfitting

In machine learning, effectively managing both overfitting and underfitting is crucial for developing models that generalize well to unseen data.

# Conclusion - Key Concepts

- **Overfitting**:
  - Definition: A model that is too complex captures noise and performs well on training data but poorly on new data.
  - Example: A decision tree with excessive branches can perfectly classify training data but fails with unseen examples.
- **Underfitting**:
  - Definition: A model that is too simple fails to capture data trends, leading to poor performance on both training and new data.
  - Example: A linear model trying to fit non-linear data will miss essential patterns.

# Conclusion - Managing Techniques and Takeaways

## Techniques for Management

- **Data Handling**: Use cross-validation to validate model performance.
- **Model Complexity**: Select suitable models based on dataset complexity.
- **Regularization Techniques**:
  - **L1 Regularization (Lasso)**:

$$L(\beta) = \frac{1}{N} \sum_{i=1}^{N} (y_i - \hat{y}_i)^2 + \lambda \sum_{j=1}^{p} |\beta_j| \tag{9}$$

  - **L2 Regularization (Ridge)**:

$$L(\beta) = \frac{1}{N} \sum_{i=1}^{N} (y_i - \hat{y}_i)^2 + \lambda \sum_{j=1}^{p} \beta_j^2 \tag{10}$$

# Conclusion - Next Steps

## Q&A Session

Let's discuss any questions or clarifications regarding the concepts of overfitting and underfitting!

# Q&A Session on Overfitting and Underfitting

## Overview of Key Concepts

- **Overfitting**: Model learns the training data too well, capturing noise.
- **Underfitting**: Model is too simple, failing to capture underlying data structure.

# Key Signs of Overfitting and Underfitting

## Overfitting

- High accuracy on training data
- Low accuracy on validation/test data
- Complex models (e.g., deep learning with many parameters)

## Underfitting

- Low accuracy on both training and validation/test data
- Basic model structure (e.g., linear regression on non-linear data)

## Overfitting Example

Consider a polynomial regression model fit to a linear dataset:

$$y = a_0 + a_1 x + a_2 x^2 + \ldots + a_{10} x^{10} \tag{11}$$

A higher degree polynomial will perfectly predict training data but fail on unseen data.

## Underfitting Example

A linear regression on a sinusoidal dataset:

$$y = mx + b \tag{12}$$

Results in a straight line that does not match the underlying curvatures.

# Managing Overfitting and Underfitting

## Techniques to Handle Overfitting

- Regularization (L1 and L2)
- Cross-Validation (k-fold)
- Pruning (for decision trees)

## Strategies for Addressing Underfitting

- Increase Model Complexity
- Feature Engineering
- Adjust Hyperparameters

- Share your experiences with overfitting or underfitting. Can you provide examples?
- What regularization techniques or hyperparameter tuning have you applied? How did they perform?
- Any questions about visualizing model performance and identifying these issues?

# Conclusion

This session aims to clarify and deepen your understanding of overfitting and underfitting, providing you tools to identify, address, and manage these challenges in model training and evaluation. Feel free to ask any questions or seek further clarification!