

# AdaBoost: From Weak Learners to Strong Classifiers

## Minimizing Exponential Error via Sequential Learning

Nikhil, Jannen

January 9, 2026

# The Power of the Committee

- **The Problem:** Complex data is rarely linearly separable.
- **The Naive Solution:** Build a massive, complex model.
- **The AdaBoost Solution:** Combine simple “Decision boundaries.”

$$H(x) = \text{sign} \left( \sum \alpha_k h_k(x) \right)$$

- **Ensemble Learning**
  - Combines multiple models to improve performance
  - Types: Boosting, Bagging, Stacking
- **Weak Learners**
  - Models that perform slightly better than random guessing
  - Examples: Decision stumps, shallow decision trees

# How AdaBoost Works

- **Algorithm Overview**

- Initialize weights of training instances
- Iteratively train weak learners
- Adjust weights based on misclassifications
- Combine weak learners using weighted majority voting

- **Mathematical Formulation**

- Weight update rule:  $D_{t+1}(i) = D_t(i) \cdot \exp(-\alpha_t \cdot y_i \cdot h_t(x_i)) / Z_t$
- Error calculation:  $\epsilon_t = \sum D_t(i)$  for misclassified instances
- Final model:  $H(x) = \text{sign}(\sum \alpha_t \cdot h_t(x))$

# Advantages of AdaBoost

- **Improved Accuracy**
  - Combines weak learners to form a strong learner
- **Versatility**
  - Compatible with various weak learners
- **Robustness**
  - Resistant to overfitting
- **Ease of Use**
  - Simple implementation and minimal hyperparameter tuning

# Limitations of AdaBoost

- **Sensitivity to Noise**

- Impact of noisy data and outliers

- **Computational Complexity**

- Iterative nature and computational cost

- **Performance on Imbalanced Data**

- Challenges with highly imbalanced datasets

# Example Use Cases

- **Spam Detection**

- Classifying emails as spam or not spam

- **Handwriting Recognition**

- Recognizing handwritten digits and characters

- **Medical Diagnosis**

- Diagnosing diseases based on patient data

- **Fraud Detection**

- Detecting fraudulent transactions in financial data

# Demonstration Project

- **Dataset Selection**

- Titanic Dataset from Kaggle
- Features: Passenger class, age, sex, fare, etc.
- Target Variable: Survival (0 = No, 1 = Yes)

- **Implementation Steps**

- Data preprocessing: Handle missing values, encode categorical variables
- Model training: Use scikit-learn's AdaBoostClassifier
- Evaluation: Compute accuracy, precision, recall, and F1-score
- Visualization: Plot feature importance and confusion matrix

# Conclusion

- **Summary of Key Points**

- AdaBoost combines weak learners to form a strong learner
- It is versatile, robust, and easy to use
- Effective for various classification tasks

- **Future Directions**

- Explore advanced boosting techniques
- Apply AdaBoost to new domains and datasets

- **Q&A**

- Open the floor for questions and discussion

# References

- **Books and Articles**

- Machine Learning by Tom Mitchell
- Pattern Recognition and Machine Learning by Christopher Bishop

- **Online Resources**

- Kaggle: <https://www.kaggle.com>
- scikit-learn documentation: <https://scikit-learn.org>
- LaTeX Beamer documentation: <https://ctan.org/pkg/beamer>