

Week 7: Model Training and Evaluation

Ekarat Rattagan

July 28, 2025

Which model is the best choice?

$$h_{\theta}(x) = \theta_0 + \theta_1 x_1$$

$$h_{\theta}(x) = \theta_0 + \theta_1 x_1 + \theta_2 x_2^2$$

...

$$h_{\theta}(x) = \theta_0 + \theta_1 x_1 + \theta_2 x_2^2 + \dots + \theta_9 x_9^9$$

$$h_{\theta}(x) = \theta_0 + \theta_1 x_1 + \theta_2 x_2^2 + \dots + \theta_9 x_9^9 + \theta_{10} x_{10}^{10}$$

Two Methods in Training process (1/2)

- **Hold-out Method**

- 80% Train set
- 20% Test set

- Pros:

- Simple
- Fast

- Cons:

- Overfitting

Two Methods in Training process (2/2)

- **Cross validation Method**

- Training set: 60%
- Validation set: 20%
- Test set: 20%

- Pros:

- Variance (Overfitting) is reduced

- Cons:

- Slow

Which one reduces the error caused by overfit or underfit?

- Get more training data
- Try smaller sets of features
- Try getting additional features
- Try adding polynomial features
- Try decreasing λ (Ridge penalty)
- Try increasing λ (Ridge penalty)

Bias vs Variance

- **Underfit:** High Bias (Linear Model: $\theta_0 + \theta_1 x$)
- **Good Fit:** Proper Balance (Quadratic Model: $\theta_0 + \theta_1 x + \theta_2 x^2$)
- **Overfit:** High Variance (Higher-order Polynomial Model:
 $\theta_0 + \theta_1 x + \theta_2 x^2 + \theta_3 x^3 + \theta_4 x^4$)

Diagnosis: Bias vs Variance (Polynomial Degree)

- Plot error (MSE) against polynomial degree
- Low-degree polynomial: Underfitting
- Moderate-degree polynomial: Good Fit
- High-degree polynomial: Overfitting

Diagnosis: Bias vs Variance (Regularization Parameter)

- Regularization function:

$$J(\theta) = \frac{1}{2N} \sum_{i=1}^N (h_{\theta}(x_i) - y_i)^2 + \frac{\lambda}{2} \sum_{j=1}^d \theta_j^2 \quad (1)$$

- Increasing λ reduces variance but increases bias
- Decreasing λ reduces bias but increases variance

Effect of Training Data Size

- Increasing training data reduces variance but not bias
- Underfit: Not enough complexity to capture patterns
- Overfit: Too complex, memorizing data instead of generalizing

- Confusion Matrix
- Precision, Recall, Accuracy
- F1-Score
- ROC, AUC

Confusion Matrix for Classification

	Actual Positive	Actual Negative
Predicted Positive	TP (True Positive)	FP (False Positive)
Predicted Negative	FN (False Negative)	TN (True Negative)

Key Metrics:

- Sensitivity (TPR): $\frac{TP}{TP+FN}$
- Specificity (TNR): $\frac{TN}{TN+FP}$
- False Positive Rate (FPR): $\frac{FP}{FP+TN}$

Accuracy, Recall, and Precision

- Accuracy: $\frac{TP+TN}{TotalPopulation}$
- Recall: $\frac{TP}{TP+FN}$
- Precision: $\frac{TP}{TP+FP}$

Imbalanced Dataset Challenges

- Accuracy alone is misleading for imbalanced datasets
- Need to consider Precision and Recall
- Example of different models and their performance metrics

- Raschka, S. (2018). Model evaluation, model selection, and algorithm selection in machine learning. arXiv preprint arXiv:1811.12808.