

Foundations of Data Science & Analytics: Evaluation and Testing

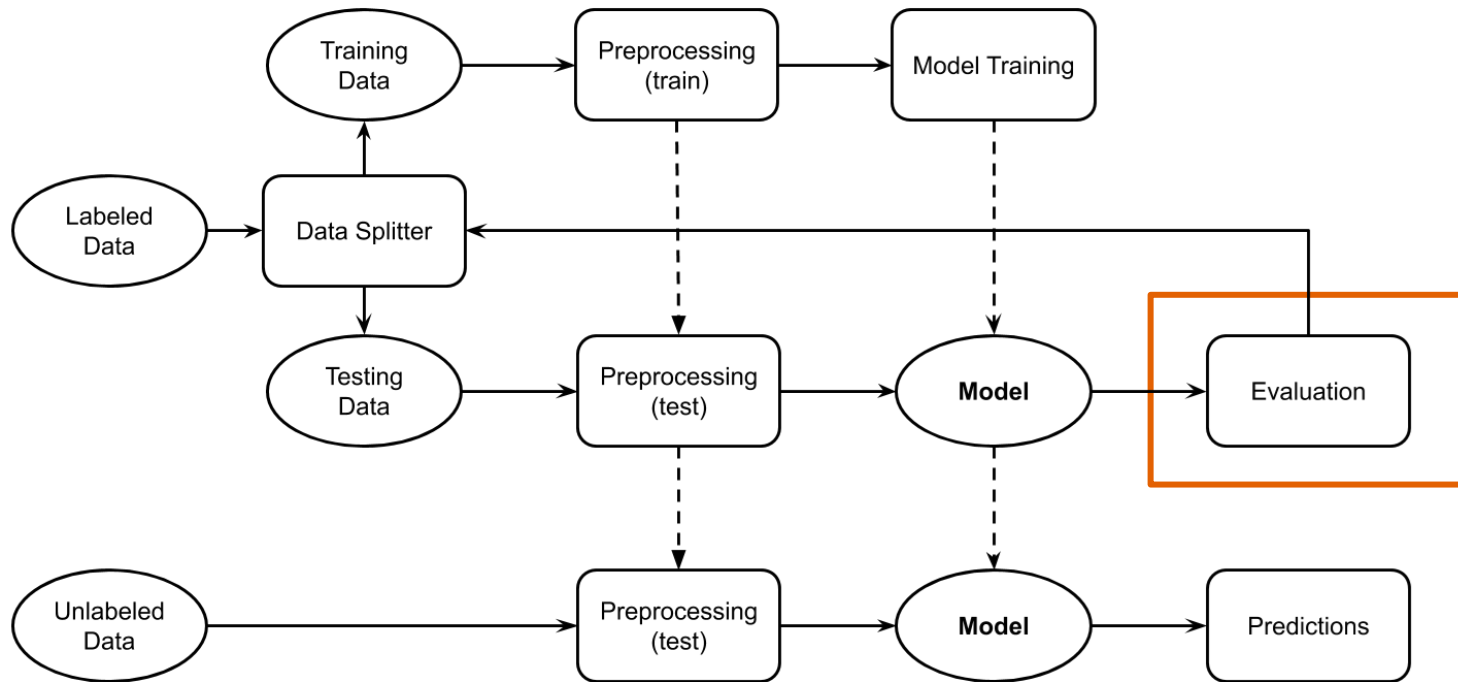
Ezgi Siir Kibris

[Introduction to Data Mining, 2nd Edition](#)

by

Tan, Steinbach, Karpatne, Kumar

Evaluation



Types of Classification Problems

- Binary Classification
- Multi-Class Classification

Binary Classification

Confusion Matrix		Predicted	
		Yes	No
Actual	Yes	TP	FN
	No	FP	TN

Binary Classification

Confusion Matrix		Predicted	
		Yes	No
Actual	Yes	5	2
	No	40	1000

Binary Classification

Confusion Matrix		Predicted	
		Yes	No
Actual	Yes	TP	FN
	No	FP	TN

$$Accuracy = \frac{TP+TN}{TP+TN+FP+FN}$$

$$Error Rate = \frac{FP+FN}{TP+TN+FP+FN}$$
$$= 1 - Accuracy$$

- Pros:
 - Can apply to multi-class
 - A general evaluation
- Cons:
 - Tends to ignore minority classes
 - Need to decide a threshold

Binary Classification

$$Recall = \frac{TP}{TP+FN}$$

$$Precision = \frac{TP}{TP+FP}$$

$$F_1 = 2 \times \frac{precision \times recall}{precision + recall}$$

Confusion Matrix		Predicted	
		Yes	No
Actual	Yes	TP	FN
	No	FP	TN

- Pros:
 - Good indicator for imbalanced classes
- Cons:
 - **Must be used together (precision + recall)**
 - Cannot directly apply to multi-class

Binary Classification

Confusion Matrix		Predicted	
		Yes	No
Actual	Yes	5	2
	No	40	1000

$$Precision = \frac{TP}{TP+FP}$$

$$Recall = \frac{TP}{TP+FN}$$

$$F_1 = 2 \times \frac{precision \times recall}{precision + recall}$$

$$Accuracy = 1005/1047 = 0.960$$

$$Precision = 5/45 = 0.111$$

$$Recall = 5/7 = 0.714$$

$$F1 \text{ score} = 0.192$$

Multi-Class Classification

$$\textit{Accuracy} = \frac{\textit{correct predictions}}{\textit{all data}}$$

$$\begin{aligned}\textit{Error Rate} &= \frac{\textit{incorrect predictions}}{\textit{all data}} \\ &= 1 - \textit{Accuracy}\end{aligned}$$

- Pros:
 - A general evaluation
- Cons:
 - Tend to ignore minority classes
 - Need to decide a threshold

Multi-Class Classification

One vs. Rest		Predicted	
		Yes (Ci)	No (Not Ci)
Actual	Yes (Ci)	TP(Ci)	FN(Ci)
	No (Not Ci)	FP(Ci)	TN(Ci)

Multi-Class Classification

$$\textit{Precision}(C_i) = \frac{TP(C_i)}{TP(C_i) + FP(C_i)}$$

$$\textit{Recall}(C_i) = \frac{TP(C_i)}{TP(C_i) + FN(C_i)}$$

$$F_1(C_i) = 2 \times \frac{\textit{precision}(C_i) \times \textit{recall}(C_i)}{\textit{precision}(C_i) + \textit{recall}(C_i)}$$

Multi-Class Classification

		True/Actual		
		Cat (🐱)	Fish (🐟)	Hen (🐔)
Predicted	Cat (🐱)	4	6	3
	Fish (🐟)	1	2	0
	Hen (🐔)	1	2	6

Multi-Class Classification

Class	Precision	Recall	F1-score
Cat	30.8%	66.7%	42.1%
Fish	66.7%	20.0%	30.8%
Hen	66.7%	66.7%	66.7%

Macro-precision = $(31\% + 67\% + 67\%) / 3 = 54.7\%$

Macro-recall = $(67\% + 20\% + 67\%) / 3 = 51.1\%$

Macro-F1 = $(42.1\% + 30.8\% + 66.7\%) / 3 = 46.5\%$

Minority Classes
are represented
equally in Macro-
metrics

Multi-Class Classification

		True/Actual		
		Cat (🐱)	Fish (🐟)	Hen (🐔)
Predicted	Cat (🐱)	4	6	3
	Fish (🐟)	1	2	0
	Hen (🐔)	1	2	6

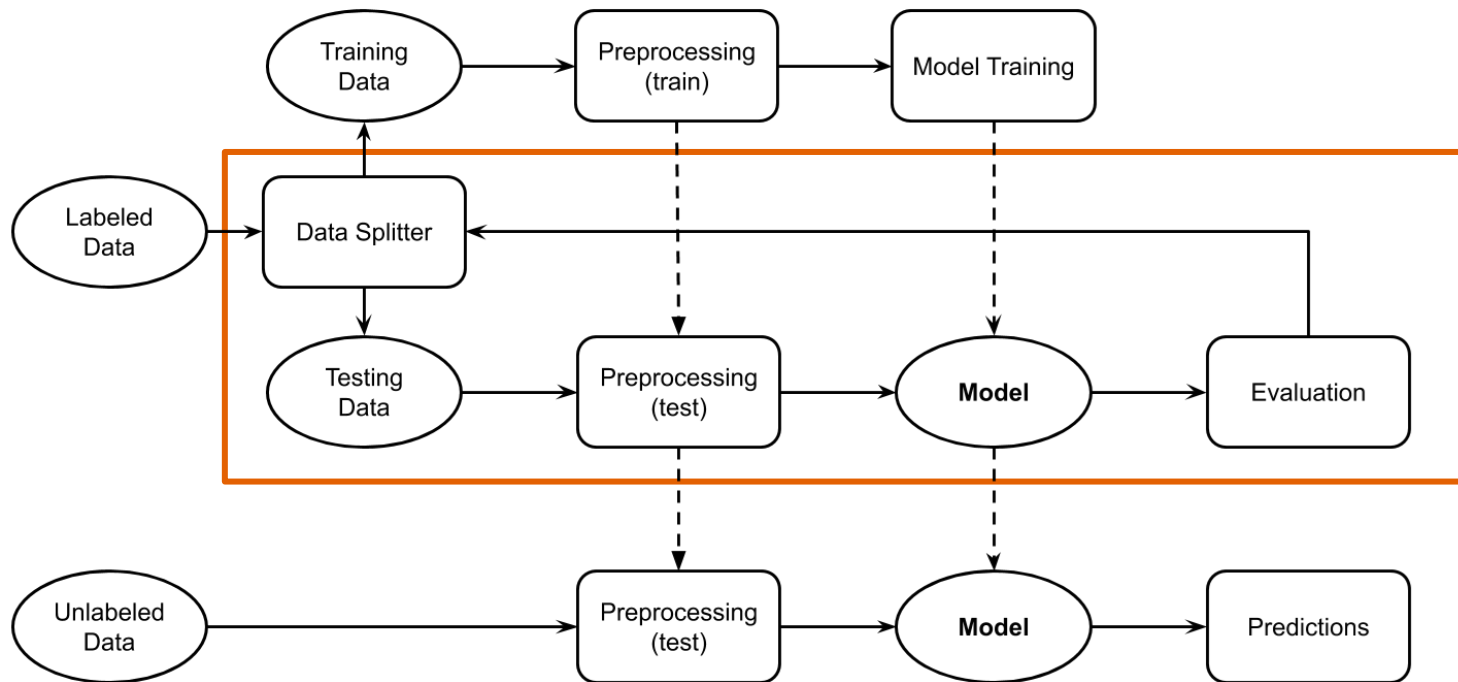
$$\text{Weighted-F1} = (6 \times 42.1\% + 10 \times 30.8\% + 9 \times 66.7\%) / 25 = 46.4\%$$

Multi-Class Classification

		True/Actual		
		Cat (🐱)	Fish (🐟)	Hen (🐔)
Predicted	Cat (🐱)	4	6	3
	Fish (🐟)	1	2	0
	Hen (🐔)	1	2	6

Micro-F1 = Micro-Precision = Micro-Recall = Accuracy = $\frac{\text{Green}}{(\text{Green} + \text{Red})} = \frac{12}{25}$

Testing



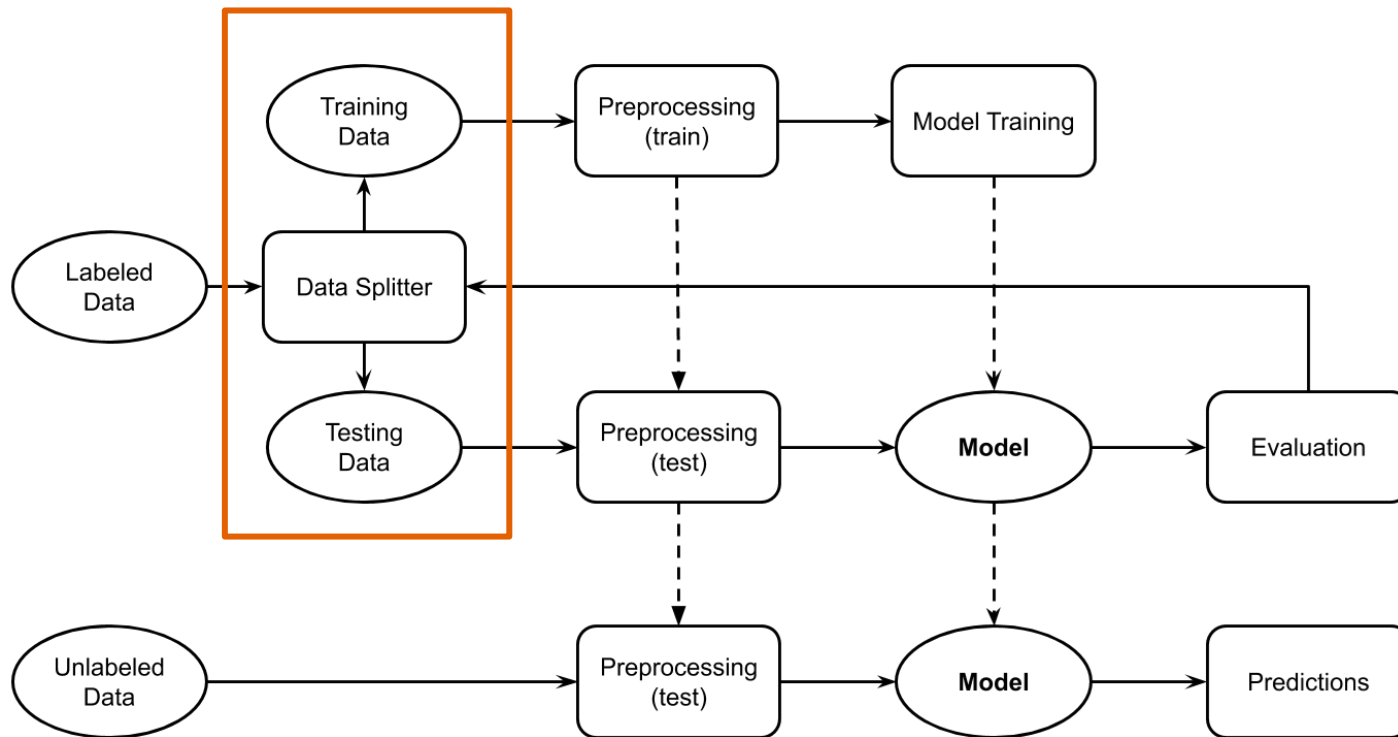
Why is Testing Needed?

- Compare the performance of different models
- Check if performance of a certain model is “good enough” on future (unseen) data

How?

- Reserve a subset of the full data set as the test set
- Train model on rest of the data, then evaluate model predictions on the test set.

How?

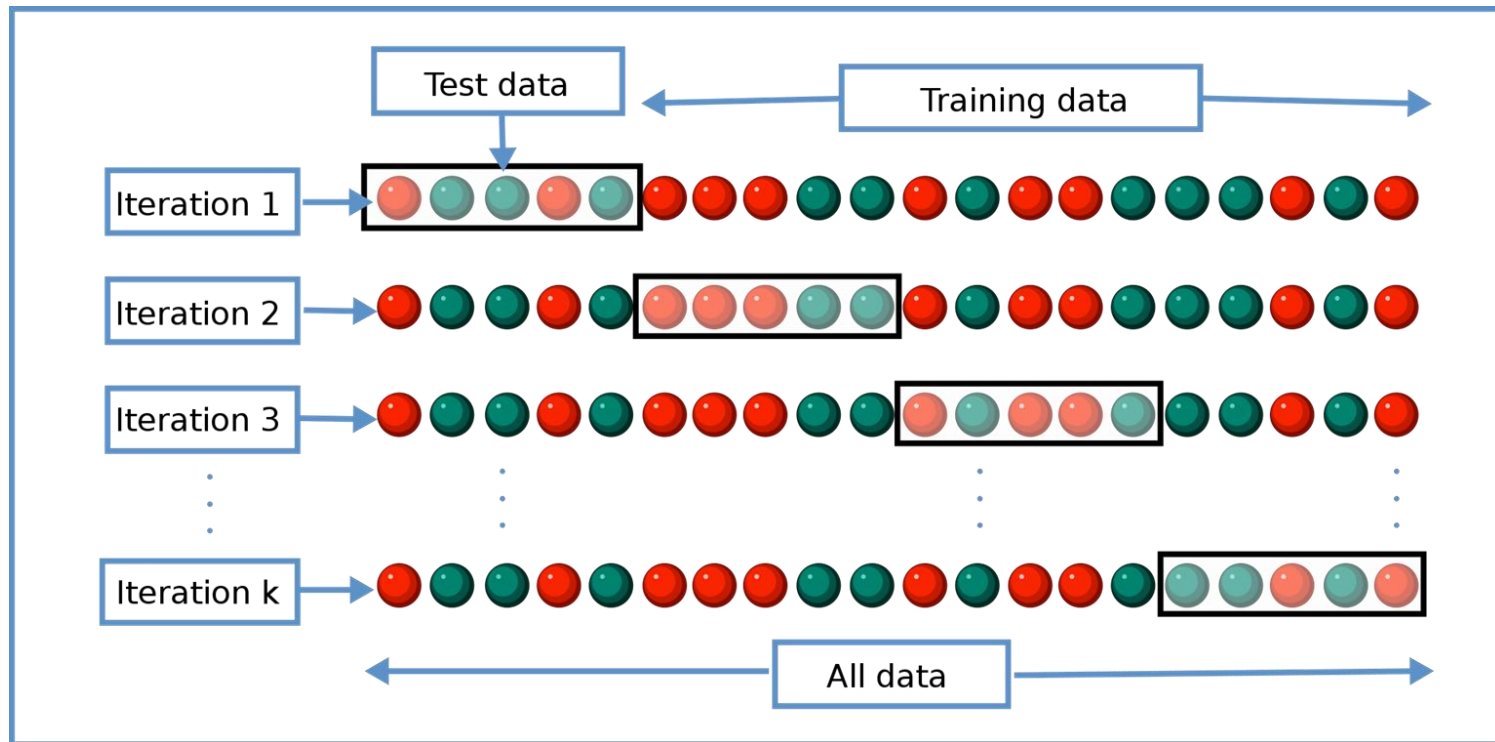


Holdout

- The most basic way to create a training and test set is to use “Holdout”
- Split the data using X% for training and Y% for testing
 - Typical splits tend to be 75/25, 80/20, 90/10
 - Splits should use stratified sampling, based upon the class distribution of the full data set
- This process should be repeated multiple times to mitigate the effects of “luck”

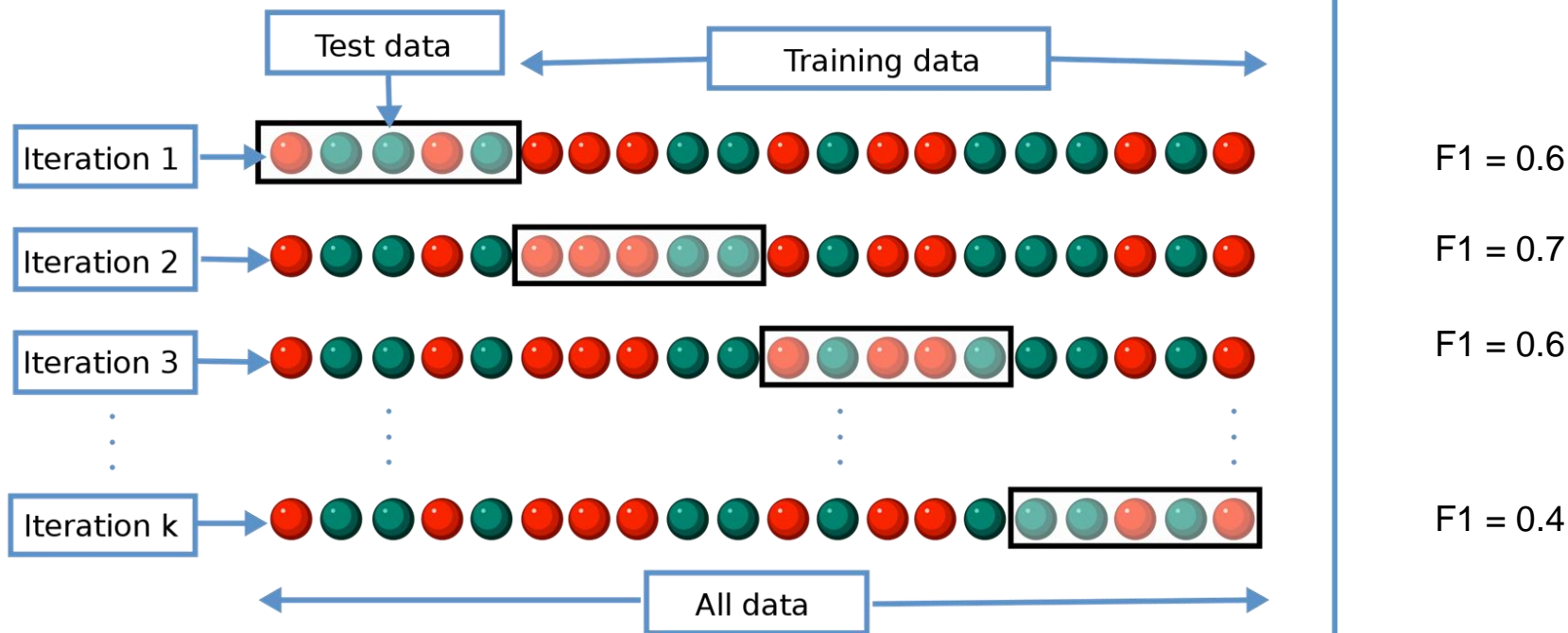
Cross-Validation

Use all data for testing



Cross-Validation

- **K-folds:**
 - split data into K subsets
 - each subset will be used as validation set once
 - performances are averaged across all K subsets



$$F1_{\text{cross}} = (0.6 + 0.7 + 0.6 + 0.4) / 4 = 0.575$$

Cross-Validation

- **L × K-fold:**
 - Shuffle data **L** times randomly
 - Each time perform a K-fold cross-validation
 - Collect performance
 - End up with **L** performance metrics
- **Compare the L metrics of each treatment to decide which is the best**

Avoid
lucky wins