



Introduction to Data Science

(Lecture 14)

Dr. Mohammad Pourhomayoun
Assistant Professor
Computer Science Department
California State University, Los Angeles





Review:

Evaluating the Accuracy of a Predictive Model Using Randomly Selected Training and Testing Sets

Evaluating The Accuracy Of Our Predictive Model

Here is a simple way to evaluate the accuracy of our predictive model:

- 1- Let's split the dataset **RANDOMLY** into two new datasets: **Training Set** (e.g. 70% of the data samples) and **Testing Set** (30% of the data).
- 2- Let's **pretend** that we do **NOT** know the label of the Testing Set!
- 3- Let's Train the model **ONLY on Training Set**, and then Predict on the Testing Set!
- 4- After prediction, we can compare the **predicted labels** for the Testing Set with the **actual labels** of it to evaluate the accuracy of our prediction!

We will learn more techniques for model evaluation (e.g. **Cross Validation** method) later in this class!



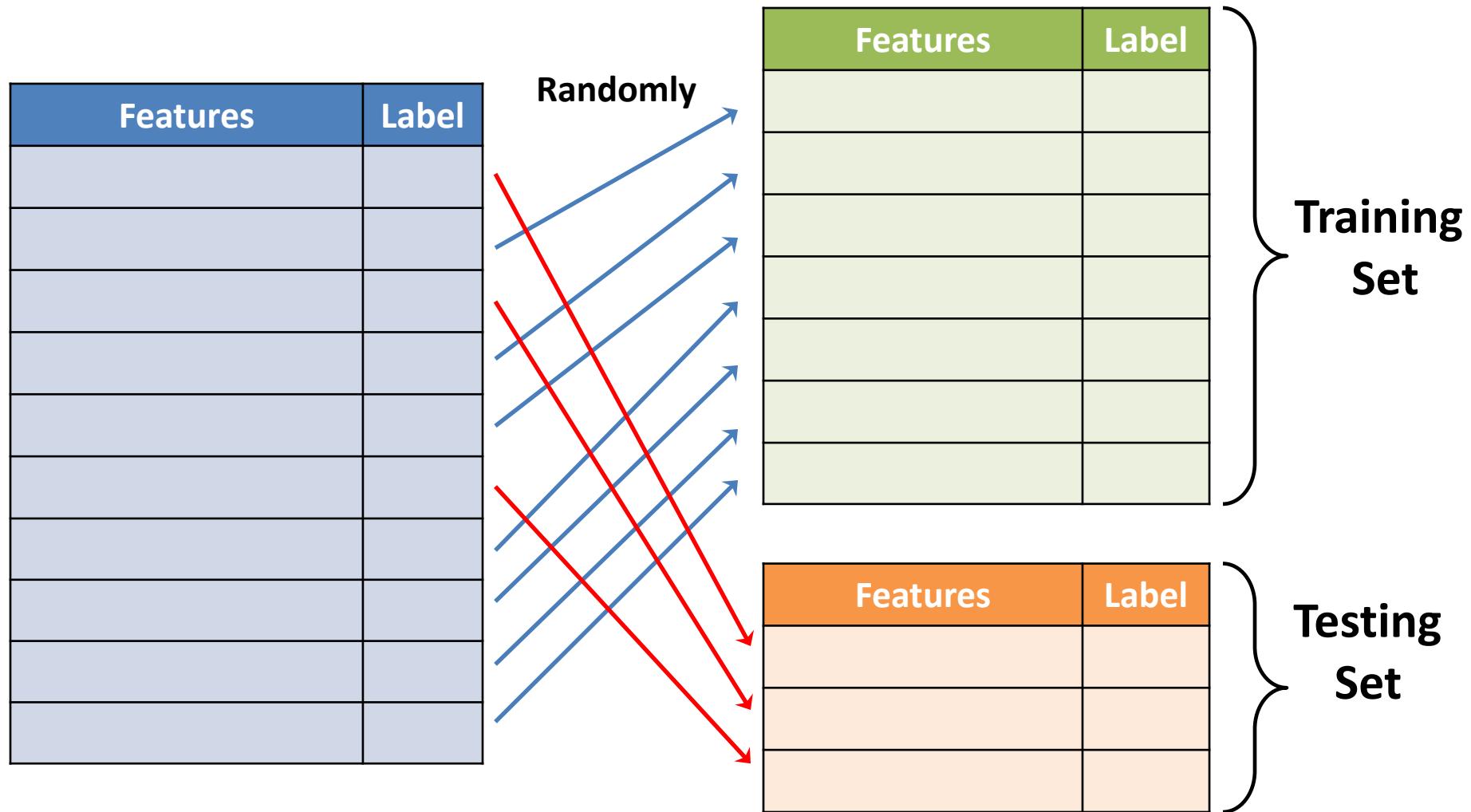
Training and Testing Sets

Features	Label

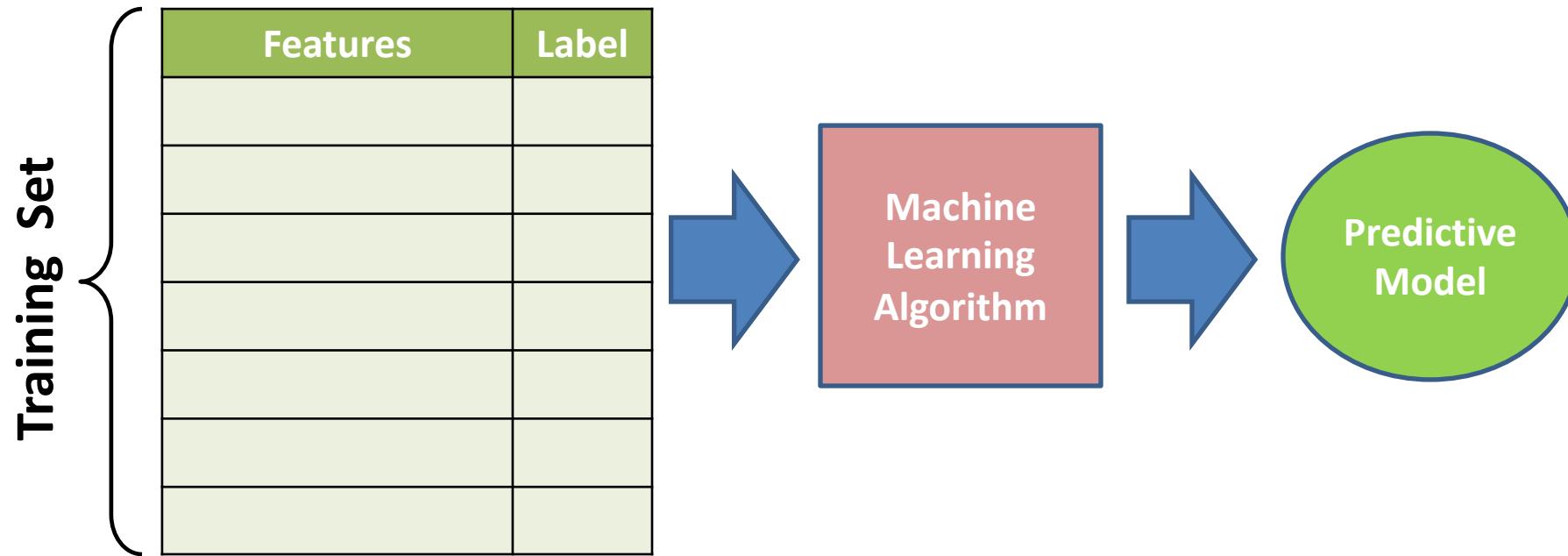
Original Dataset



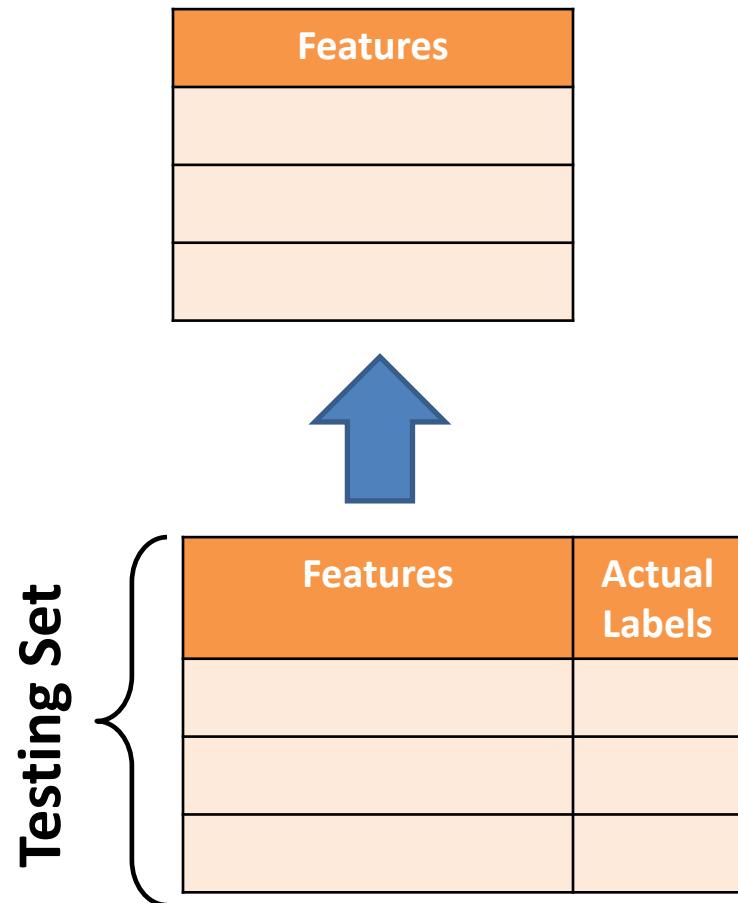
Training and Testing Sets



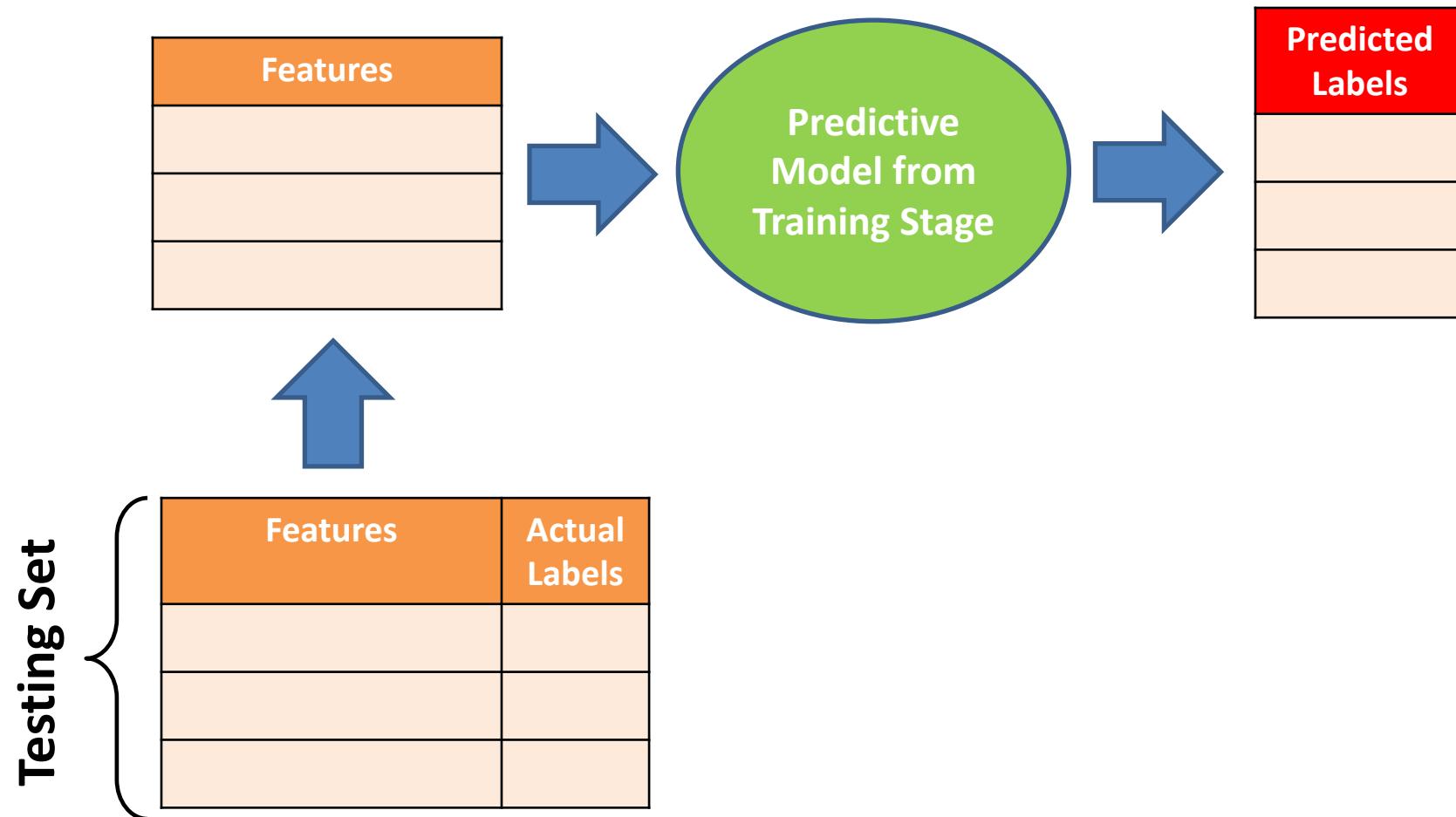
Training Stage



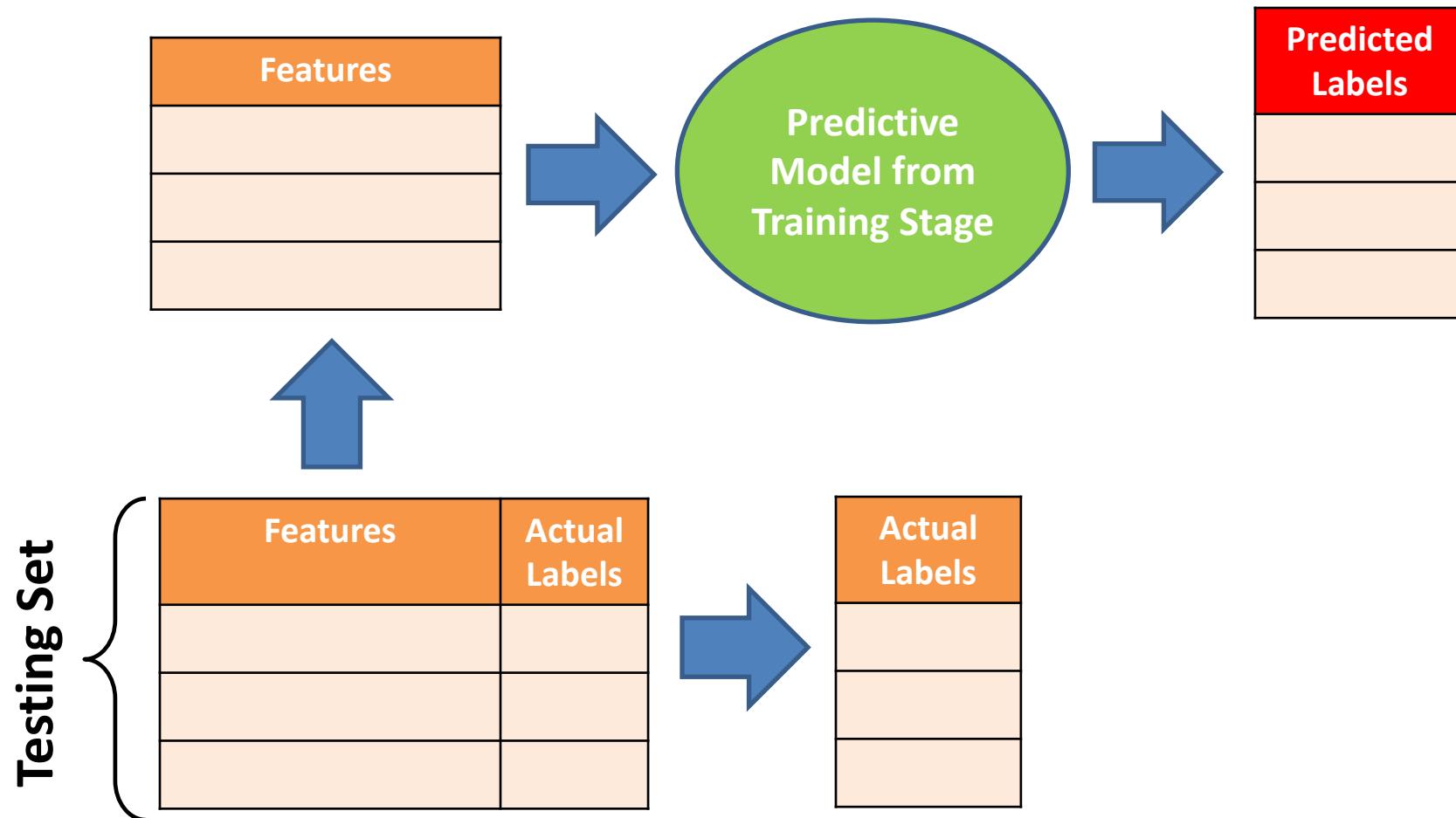
Testing Stage



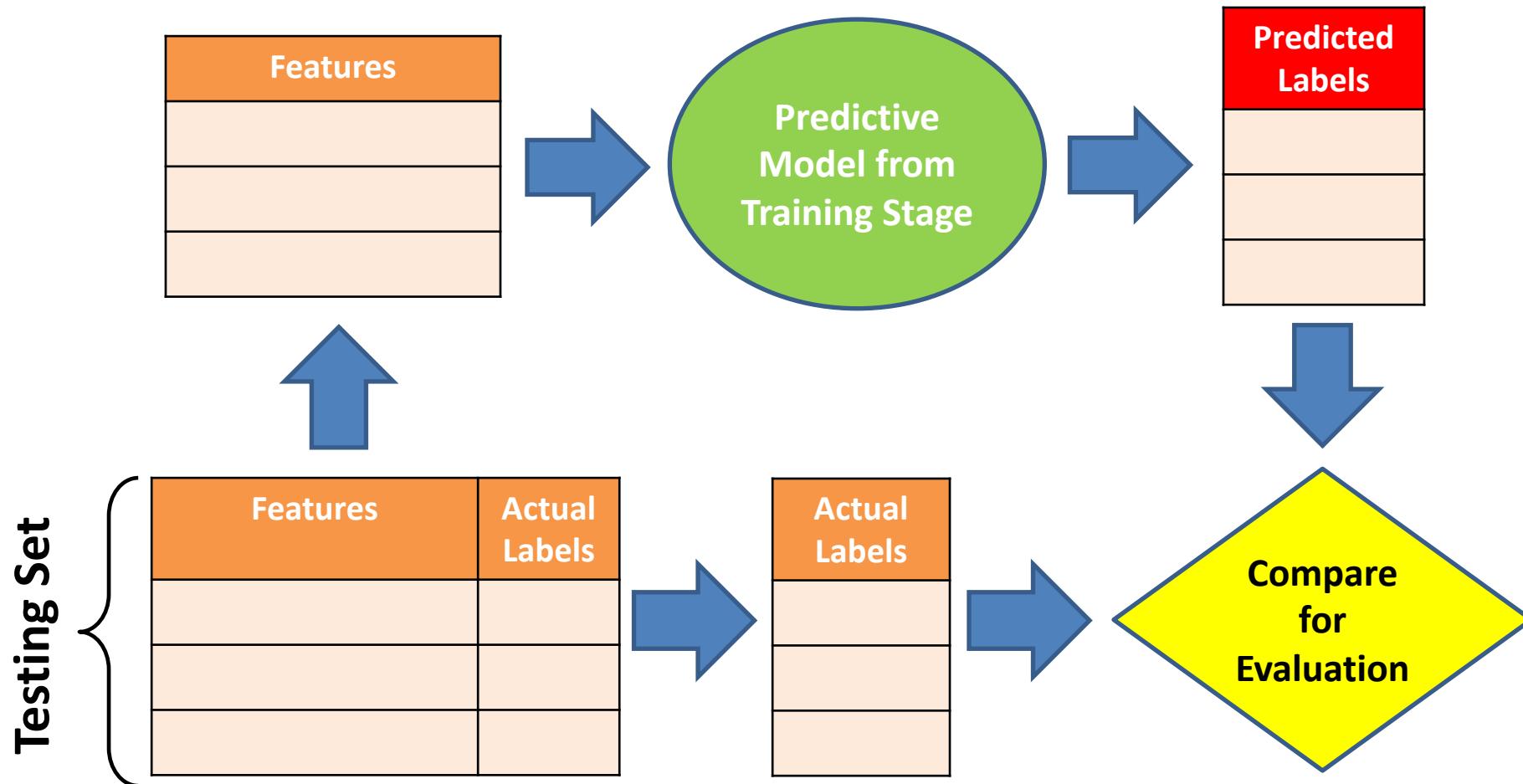
Testing Stage



Testing Stage



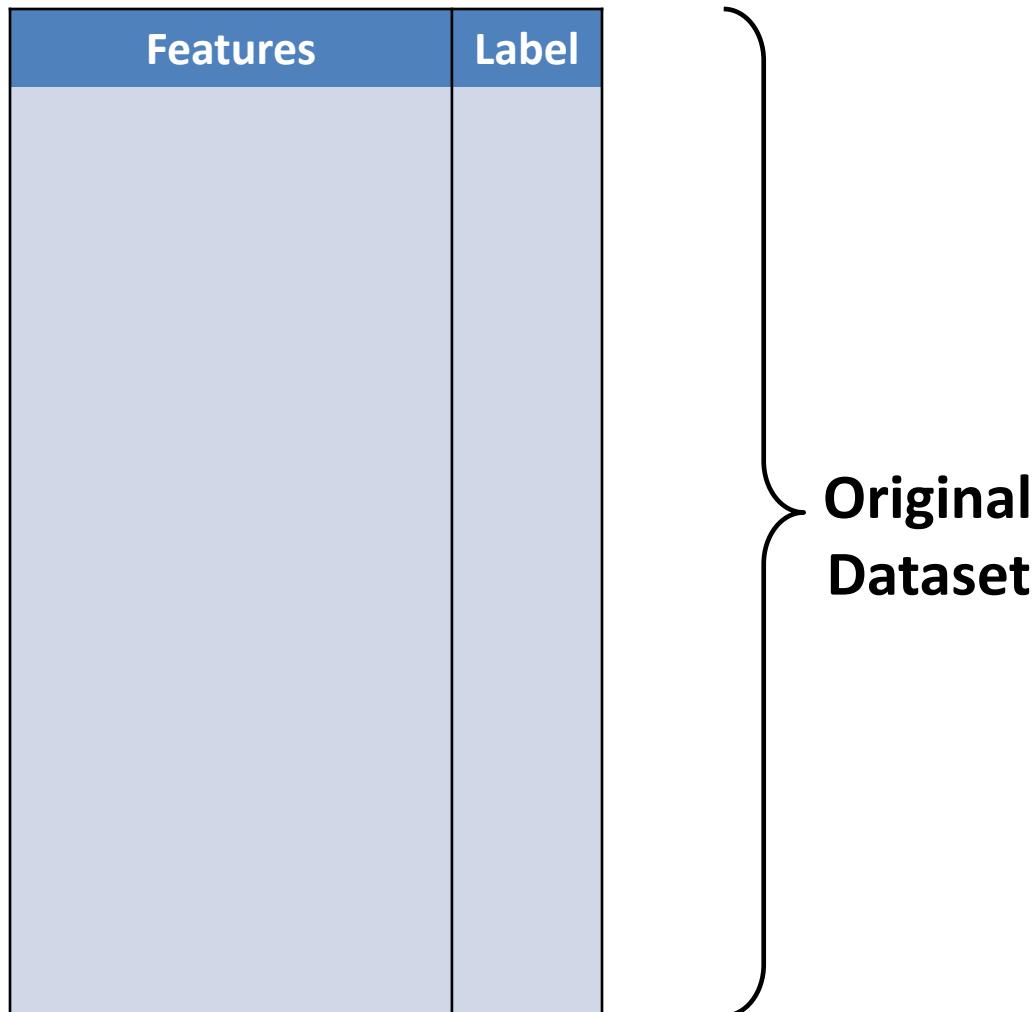
Testing Stage



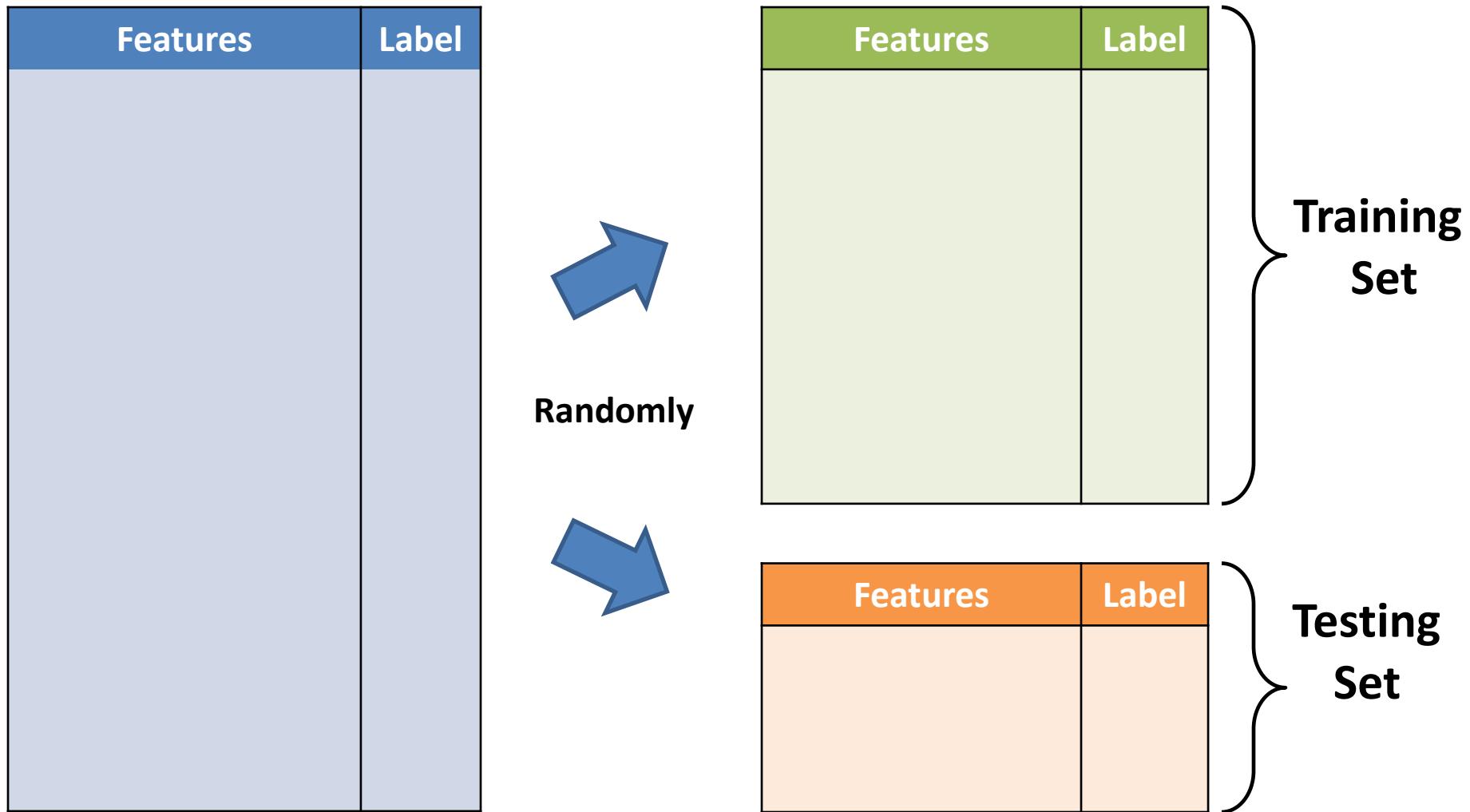


Evaluating the Accuracy of a Predictive Model Using Cross Validation

Training and Testing Sets



Training and Testing Sets



Cross Validation

- We saw how to split the dataset into Training and Testing sets, Fit the model on "training set", and then predict on "testing set" to evaluate the accuracy.
- The problem with this method is that **the results depend on the choice of split**. For example, if you are lucky, some easily predictable samples may happen to be located in the testing set, and then your testing accuracy will be incorrectly high (or vice versa!).
- In order to perform **fair evaluation**, we can repeat the splitting process several times, compute the prediction accuracy for each split, and then average the results.
- **Cross Validation tries to repeat the splitting procedure K times in a smart way such that all data samples will be used in "testing set" one time and in "Training Set" (K-1) times!**



Cross Validation

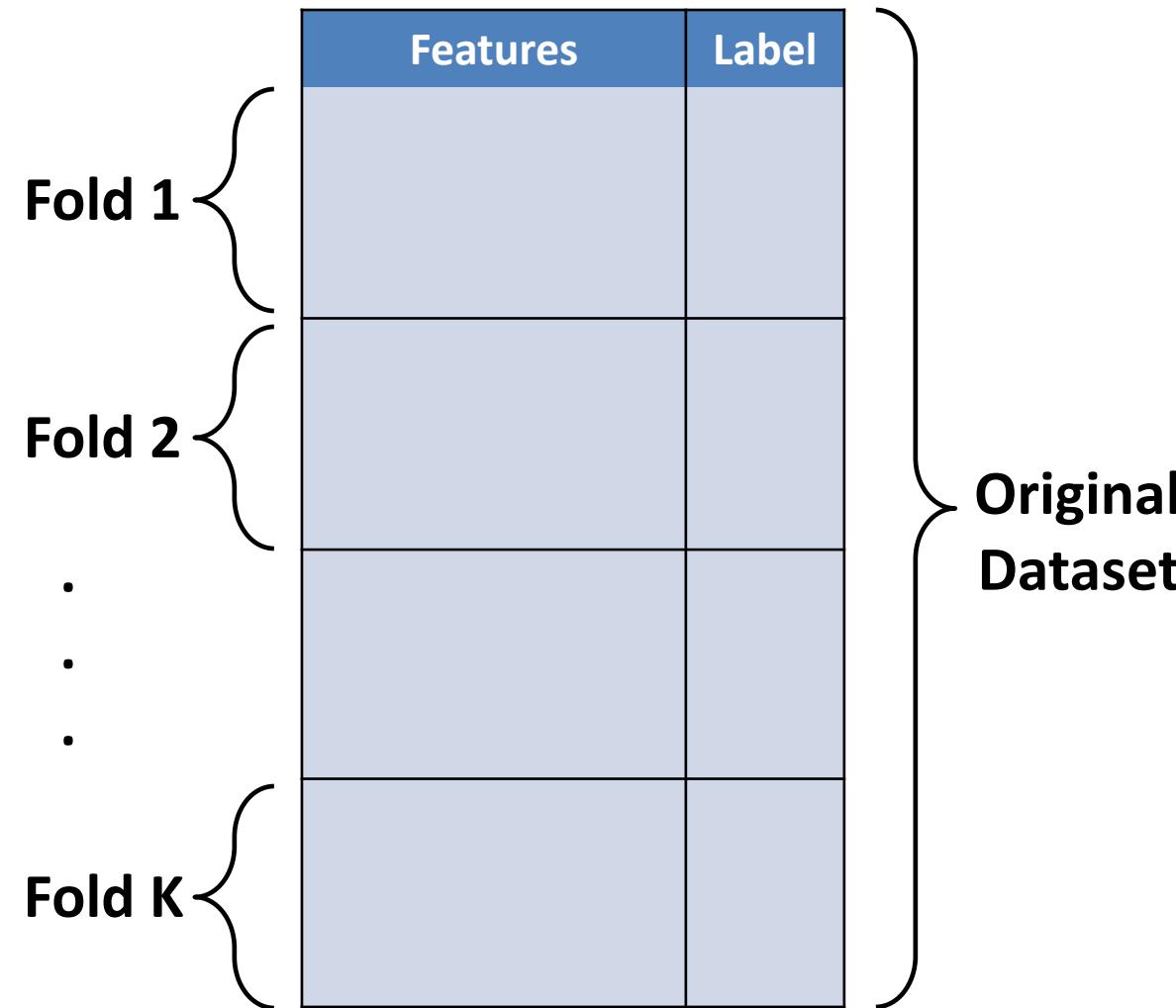
Three main steps for K-fold cross-validation:

1. Partition the dataset Randomly into K equal, non-overlapping sections (called Fold).
2. Use one of the sections as **testing set** at a time and the union of the other (K-1) sections as the **training set**. Perform training stage, testing stage, and compute the accuracy based on the split each time. Repeat this procedure K times, so that each one of the K sections is used as **testing set** one time, and as a part of **training set** (K-1) times.
3. Calculate the average of the accuracies as final result.

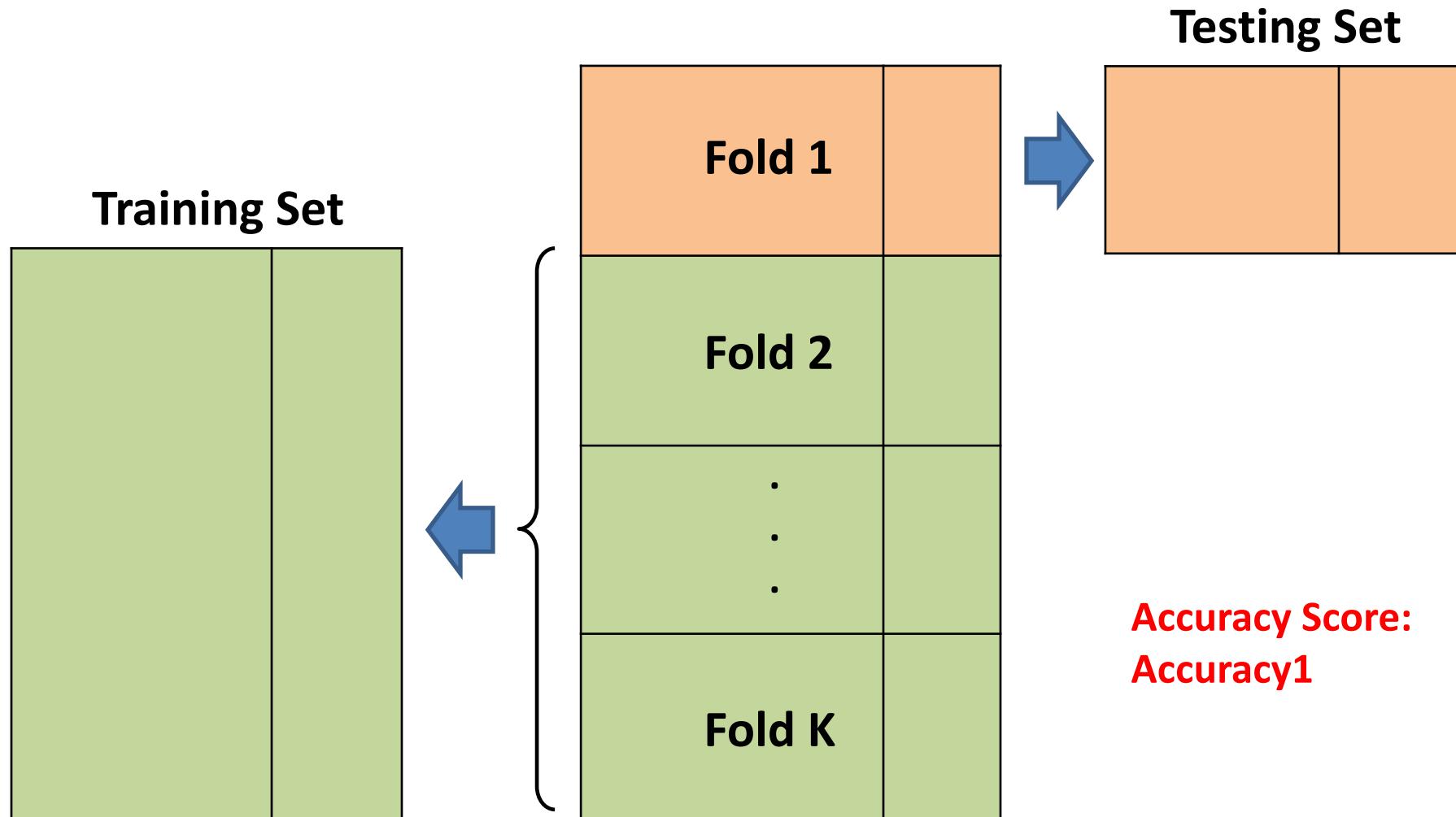
Note: K is arbitrary, but Using K=10 (10-fold cross-validation) is very common and recommended in machine learning.



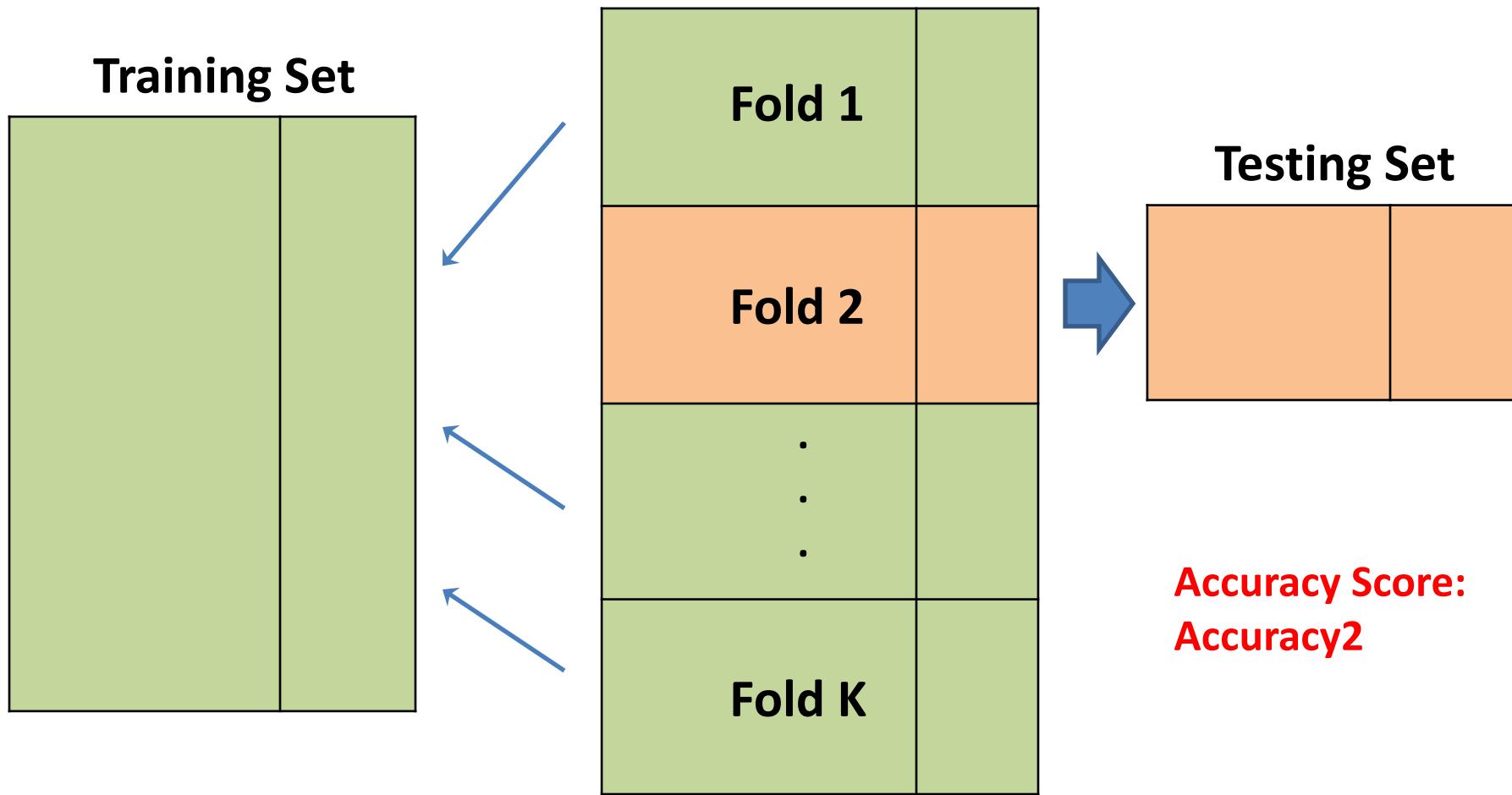
Cross Validation



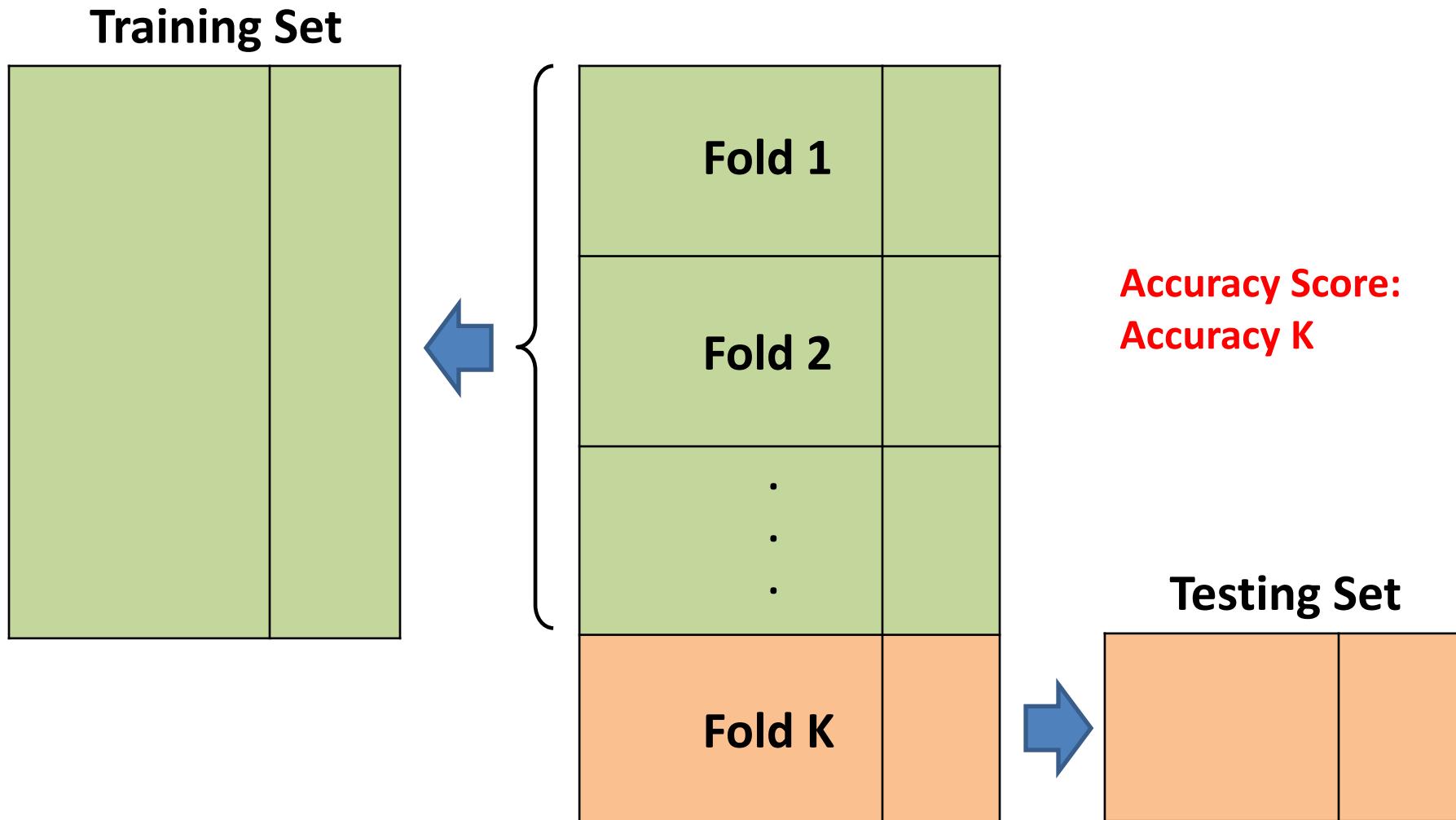
Cross Validation – Round 1



Cross Validation – Round 2



Cross Validation – Round K



Cross Validation

- **Accuracy_Score_Total =**
(Accuracy 1 + Accuracy 2 + Accuracy 3 + ... + Accuracy K) / K



Evaluation for Regression

- So far we talked about accuracy for classifiers: For classifiers we compare "the predicted labels" against "the actual labels" and calculate the accuracy as the percentage of correctly classified samples.
- In regression, the target is continuous valued. So, we need to find the error as the average difference between the "predicted target value" and the "actual value".
- The most popular metric to quantify this difference (error) is **Root Mean Square Error** or **RMSE**:

$$\text{RMSE} = \sqrt{\frac{1}{n} \sum_{i=1}^n (y_i - \hat{y}_i)^2}$$

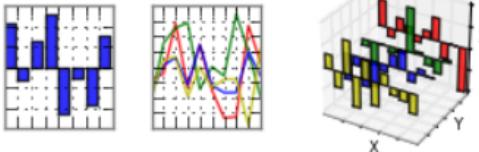
- " y_i " is the actual value. " \hat{y}_i " is the predicted value.
- Notice that **RMSE** is error. So, unlike the accuracy for classifiers, **the lower RMSE, the better!**



Data Science with Python

IP[y]: IPython
Interactive Computing

pandas
 $y_{it} = \beta' x_{it} + \mu_i + \epsilon_{it}$



scikits
learn
machine learning in Python

NumPy

SciPy.org Sponsored By ENTHOUGHT

matplotlib



Data Science in Practice

- Let's open file ***CS4661-PythonDataScience-Lab4.ipynb*** in Jupyter notebook to start the tutorial Lab4.
- In this Tutorial we will cover Logistic Regression, Linear Regression, Model Evaluation, and Cross Validation

