

Model Selection & Training

Problem: Multivariate Regression

Preprocessing

Evaluation Model with training Data set

Model

Linear Regression

Decision Tree Regression

RMSE

RMSE

Under fitting

Over fitting

## Linear Regression Goals & Evaluation Method

$$\text{MSE} \quad \frac{1}{n} \sum_{i=1}^n (Y_i - \hat{Y}_i)^2$$

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

$$\text{RMSLE} \quad \log \sqrt{\frac{\sum_{i=1}^n (Y_i - \hat{Y}_i)^2}{n}}$$

$$\text{MAE} \quad \frac{\sum_{i=1}^n |Y_i - \hat{Y}_i|}{n}$$

$$\text{MAMPE} \quad \frac{1}{n} \sum_{i=1}^n \left| \frac{Y_i - \hat{Y}_i}{Y_i} \right|$$

if data set have a lot of noise value

Which method is fit for that data set?

## Linear Regression Goals & Evaluation Method

$$\text{MSE} \quad \frac{1}{n} \sum_{i=1}^n (Y_i - \hat{Y}_i)^2$$

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

$$\text{RMSLE} \quad \log \sqrt{\frac{\sum_{i=1}^n (Y_i - \hat{Y}_i)^2}{n}}$$

$$\text{MAE} \quad \frac{\sum_{i=1}^n |Y_i - \hat{Y}_i|}{n}$$

$$\text{MAMPE} \quad \frac{1}{n} \sum_{i=1}^n \left| \frac{Y_i - \hat{Y}_i}{Y_i} \right|$$

if data set have a lot of noise value

Which method is fit for that data set?

MAE

## Linear Regression Goals & Evaluation Method

$$\text{MSE} \quad \frac{1}{n} \sum_{i=1}^n (Y_i - \hat{Y}_i)^2$$

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

$$\text{RMSLE} \quad \log \sqrt{\frac{\sum_{i=1}^n (Y_i - \hat{Y}_i)^2}{n}}$$

$$\text{MAE} \quad \frac{\sum_{i=1}^n |Y_i - \hat{Y}_i|}{n}$$

$$\text{MAMPE} \quad \frac{1}{n} \sum_{i=1}^n \left| \frac{Y_i - \hat{Y}_i}{Y_i} \right|$$

if data set have a lot of noise value

Which method is fit for that data set?

MAE

WHY?

$$(X + Y)^2 \quad \text{vs} \quad |X + Y|$$

What's difference?

Evaluation Model with training Data set

Model

Linear Regression

Decision Tree Regression

RMSE

RMSE

Under fitting

Over fitting

Cross validation

## Cross validation

Q. What is Goal of machine learning & statistical inference model?



## Cross validation

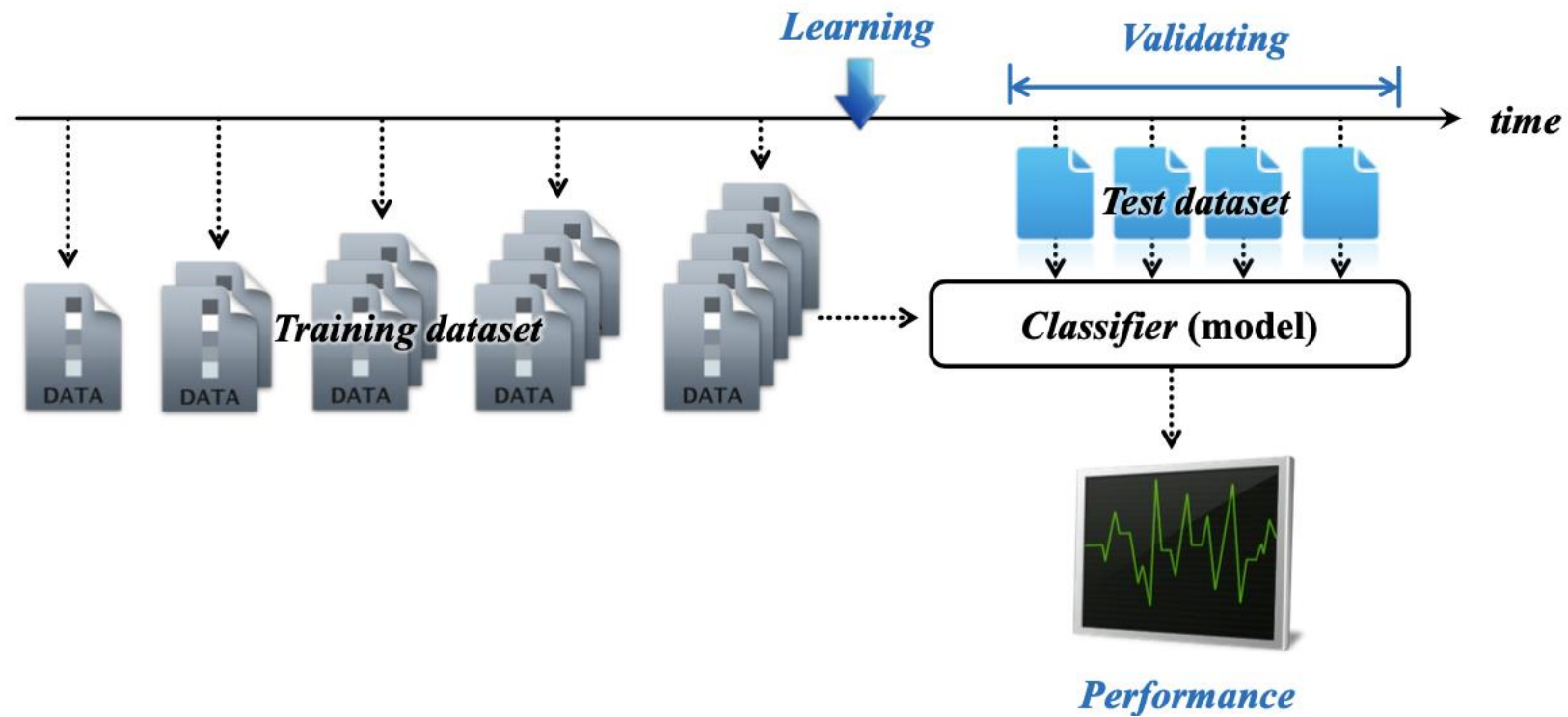
Q. What is Goal of machine learning & statistical inference model?

A. Predict the future by extracting pattern from our data(present)

## Cross validation

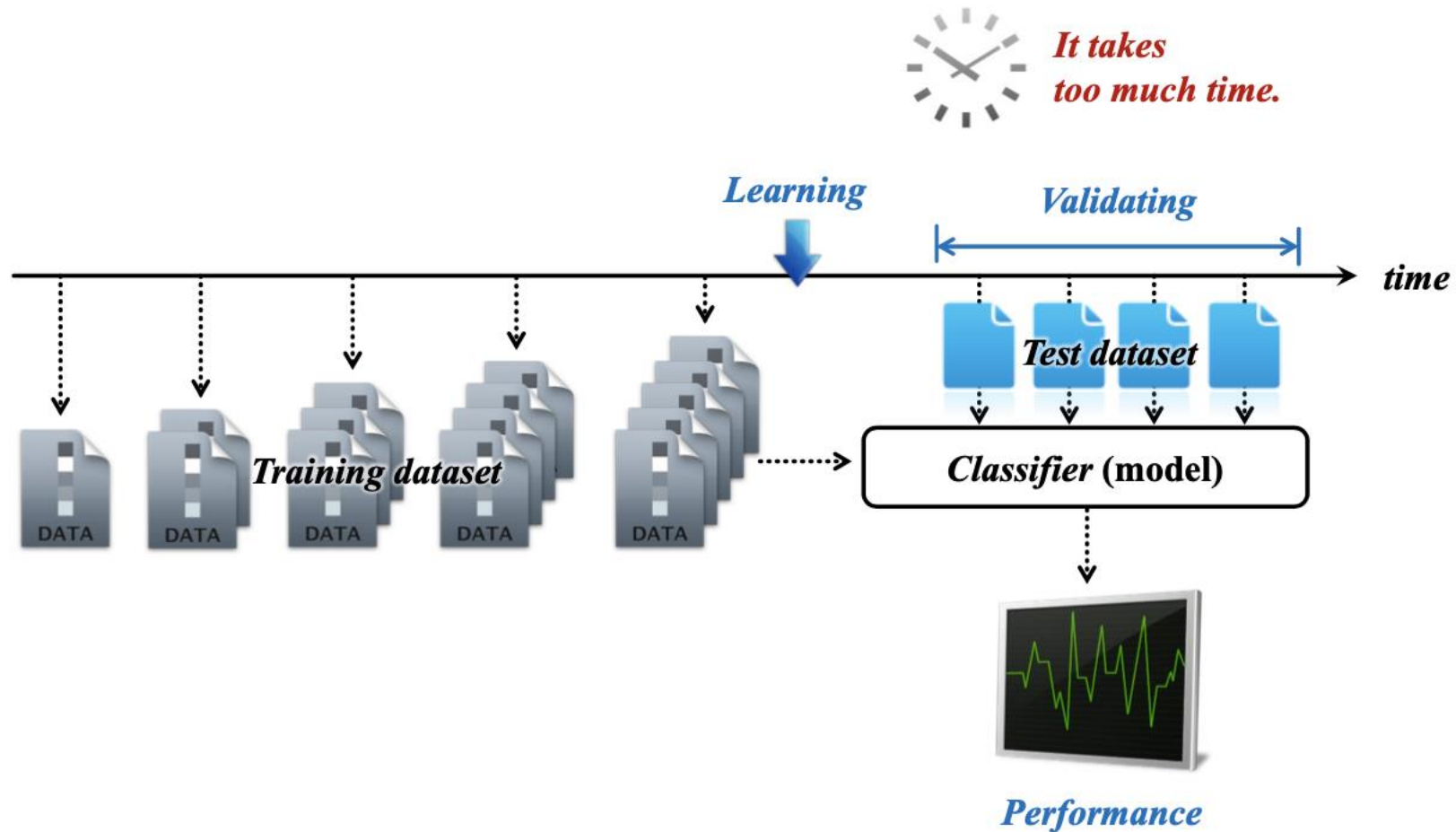
Q. What is Goal of machine learning & statistical inference model?

A. Predict the future by extracting pattern from our data(present)



## Cross validation

Q. What is Goal of machine learning & statistical inference model?

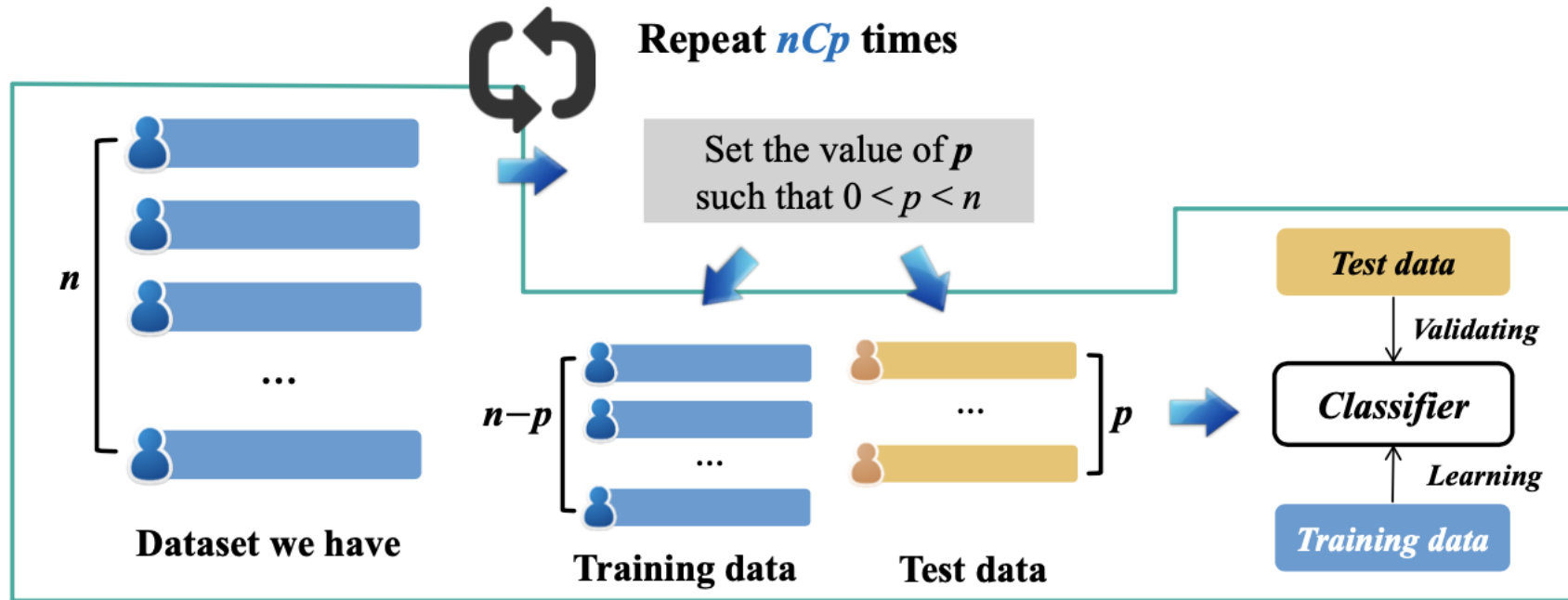


Cross validation

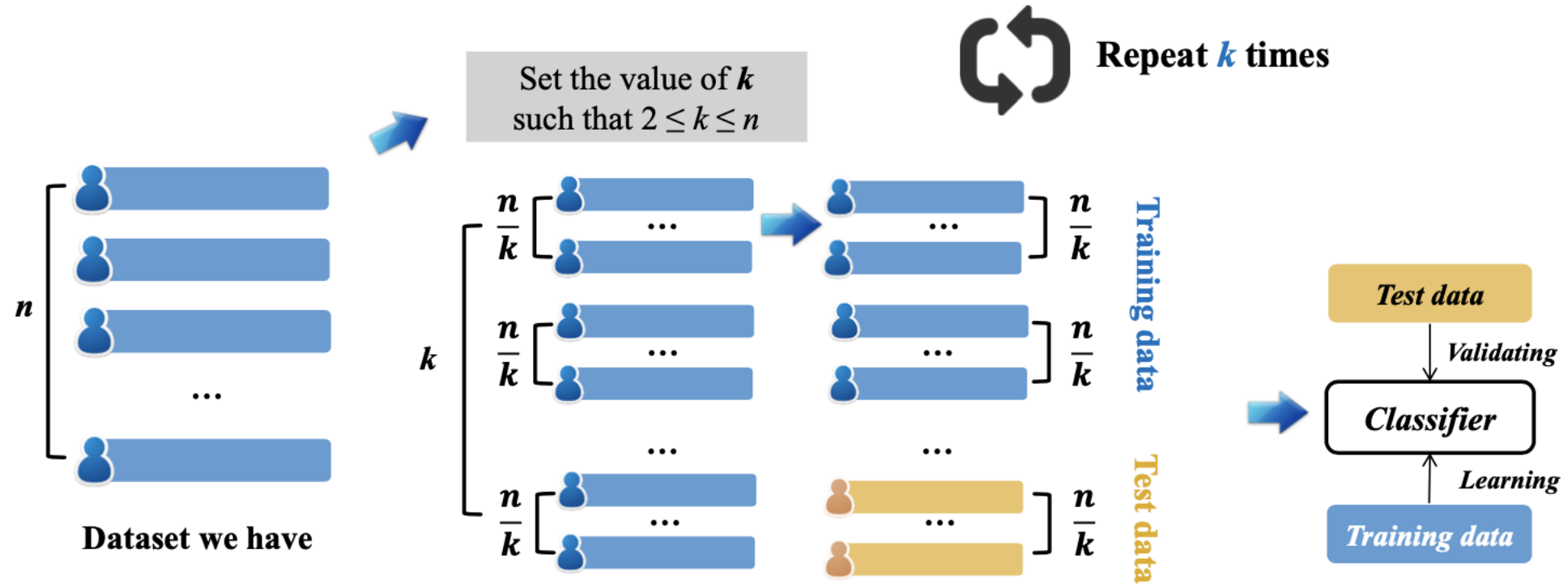
Leave-p-out

K-fold

## Leave – p - out



## K-fold



# Fine Tuning

## Model에게는 Hyper parameters 존재

### Decision Tree hyperparameter

<b>Parameters:</b>	<p><b>criterion : {"gini", "entropy"}, default="gini"</b> The function to measure the quality of a split. Supported criteria are "gini" for the Gini impurity and "entropy" for the information gain.</p> <p><b>splitter : {"best", "random"}, default="best"</b> The strategy used to choose the split at each node. Supported strategies are "best" to choose the best split and "random" to choose the best random split.</p> <p><b>max_depth : int, default=None</b> The maximum depth of the tree. If None, then nodes are expanded until all leaves are pure or until all leaves contain less than min_samples_split samples.</p> <p><b>min_samples_split : int or float, default=2</b> The minimum number of samples required to split an internal node:</p> <ul style="list-style-type: none"><li>• If int, then consider min_samples_split as the minimum number.</li><li>• If float, then min_samples_split is a fraction and <math>\text{ceil}(\text{min\_samples\_split} * \text{n\_samples})</math> are the minimum number of samples for each split.</li></ul> <p><i>Changed in version 0.18: Added float values for fractions.</i></p> <p><b>min_samples_leaf : int or float, default=1</b> The minimum number of samples required to be at a leaf node. A split point at any depth will only be considered if it leaves at least min_samples_leaf training samples in</p>
--------------------	---

1000

1,1,1,1

1,1,1,2

1,1,1,3

Max,max,max,max

2,3,5,1

3,5,6,1

.

.

.

