

Generalization Error

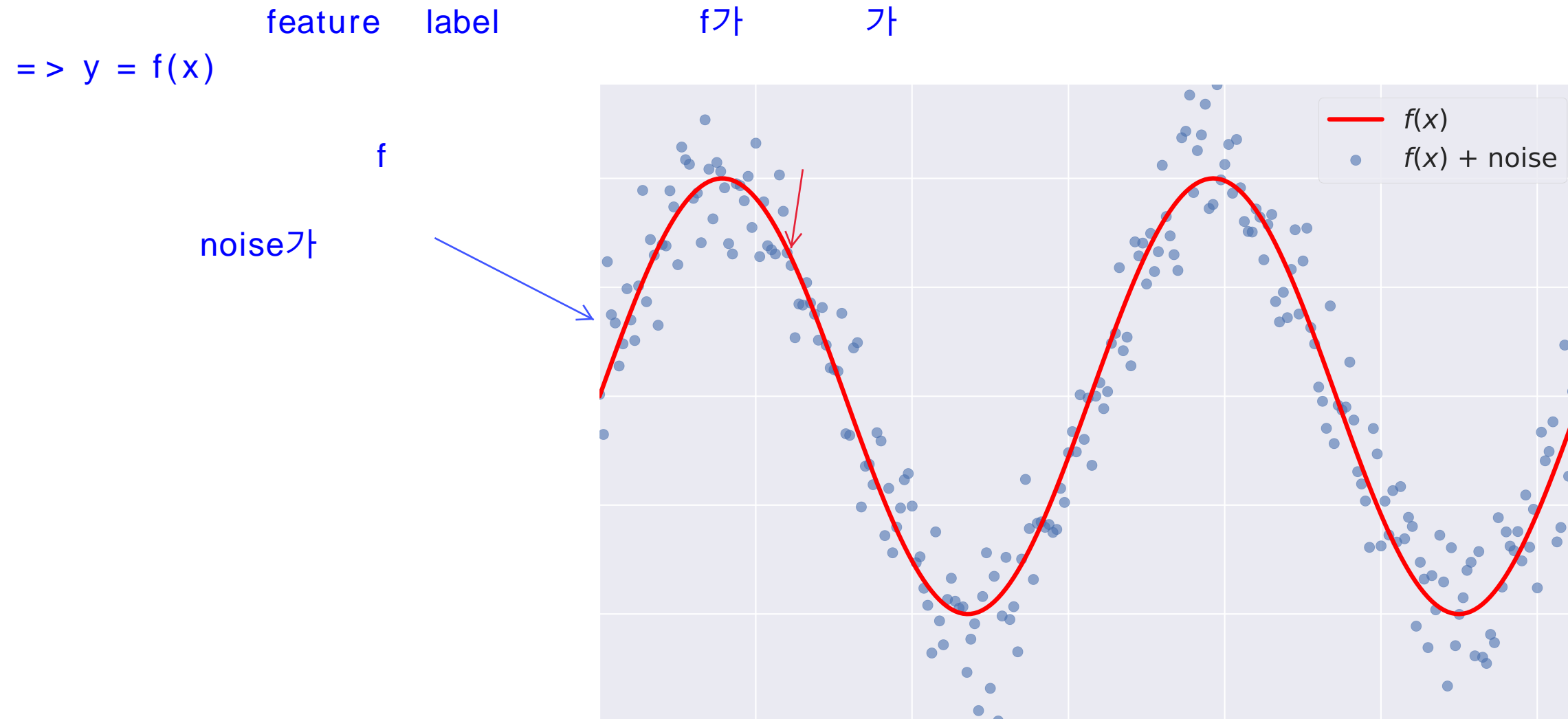
MACHINE LEARNING WITH TREE-BASED MODELS IN PYTHON



Elie Kawerk
Data Scientist

Supervised Learning - Under the Hood

- Supervised Learning: $y = f(x)$, f is unknown.



Goals of Supervised Learning

- Find a model \hat{f} that best approximates f : $\hat{f} \approx f$
- \hat{f} can be Logistic Regression, Decision Tree, Neural Network ...
- Discard noise as much as possible.
- **End goal:** \hat{f} should achieve a low predictive error on unseen datasets.

f 가 가 f_{hat}
 f_{hat} 가 f_{hat} noise가

Difficulties in Approximating f

- **Overfitting:** $\hat{f}(x)$ fits the training set noise.
- **Underfitting:** \hat{f} is not flexible enough to approximate f .

f 가

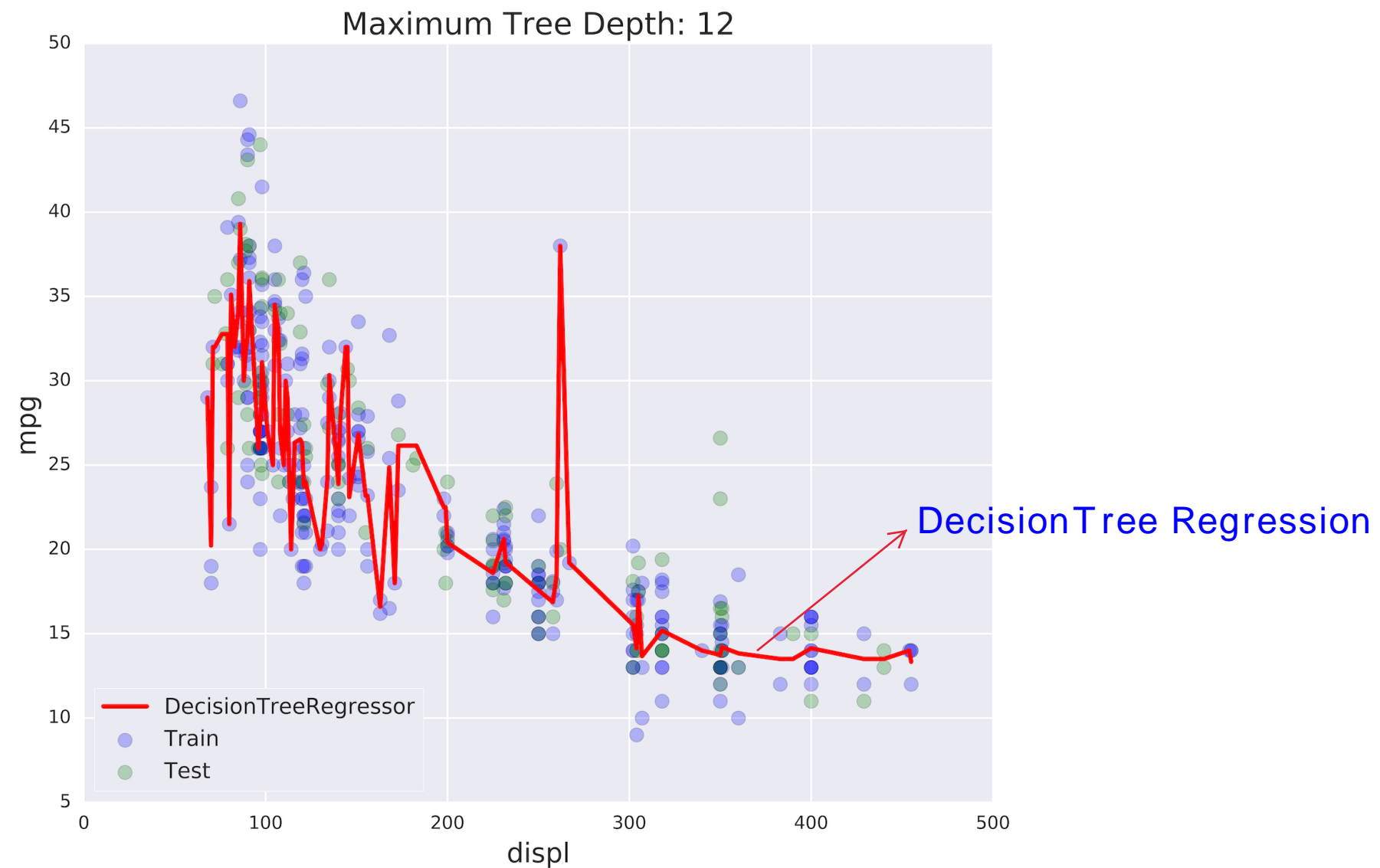
- overfitting : $\hat{f}_{\text{train_set}}$ noise
- underfitting : \hat{f} f

Overfitting

training set overfitting

noise
train_set

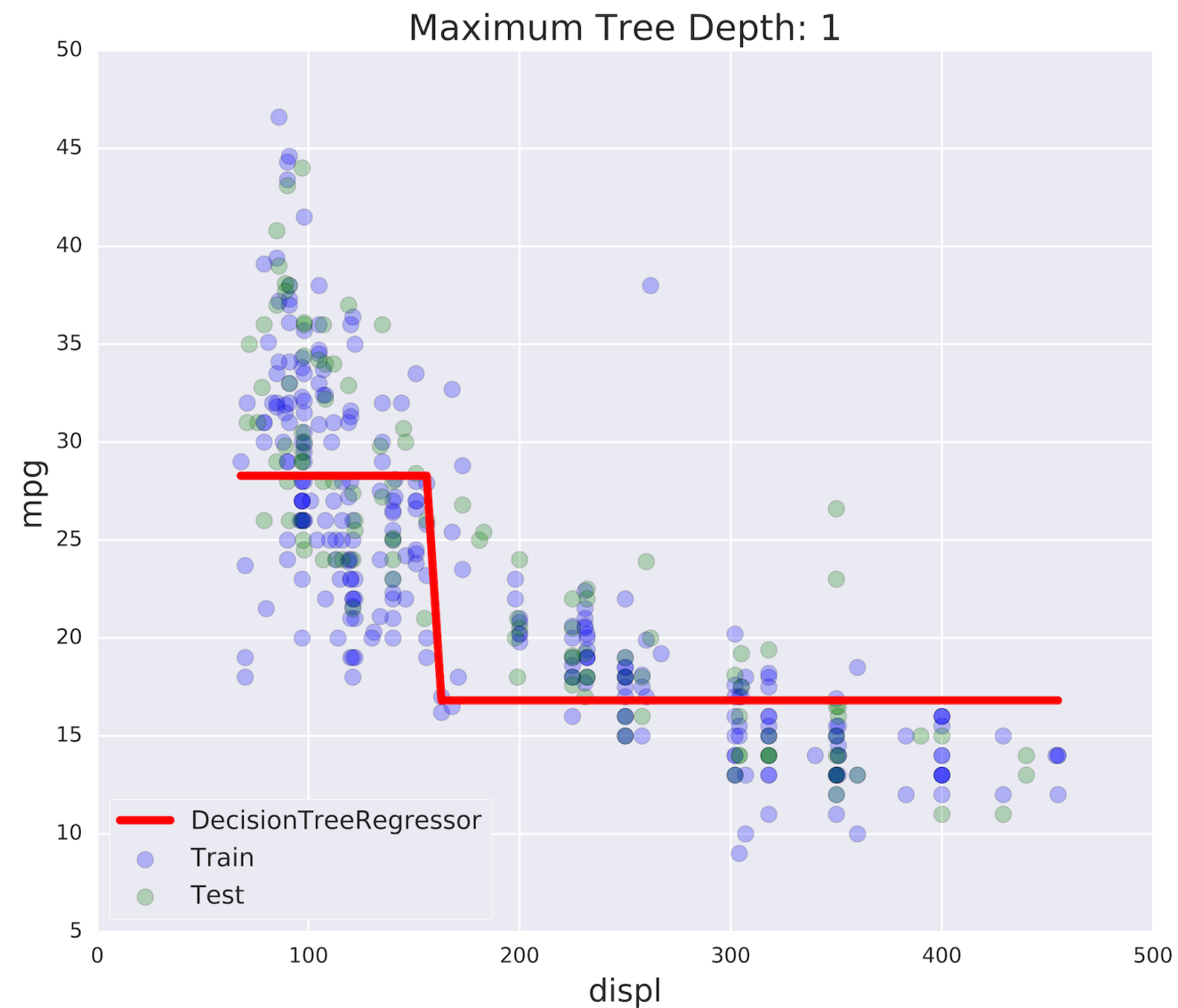
test set



Underfitting

training set test set

train model feature label



Generalization Error

- **Generalization Error of \hat{f} :** Does \hat{f} generalize well on unseen data?
- It can be decomposed as follows: Generalization Error of $\hat{f} = bias^2 + variance + \text{irreducible error}$

model generalization Error model
 bias, variance, irreducible error 가
 (irreducible error() : 가 ())

$$MSE = \text{분산} + \text{편파성}^2 + \text{irreducible error}$$

분산(variance): 모델이 예측한 데이터의 분포정도, 전체 데이터 집합 중 다른 학습데이터를 사용하였을 때 모델f가 변하는 정도

편파성(bias): 모델이 예측한 데이터와 실제 데이터간의 차이의 정도, 학습 알고리즘에서 잘못된 가정을 했을 때 발생하는 오차

irreducible error: 우리가 알수없는 오차(오차항)에 대한 분산

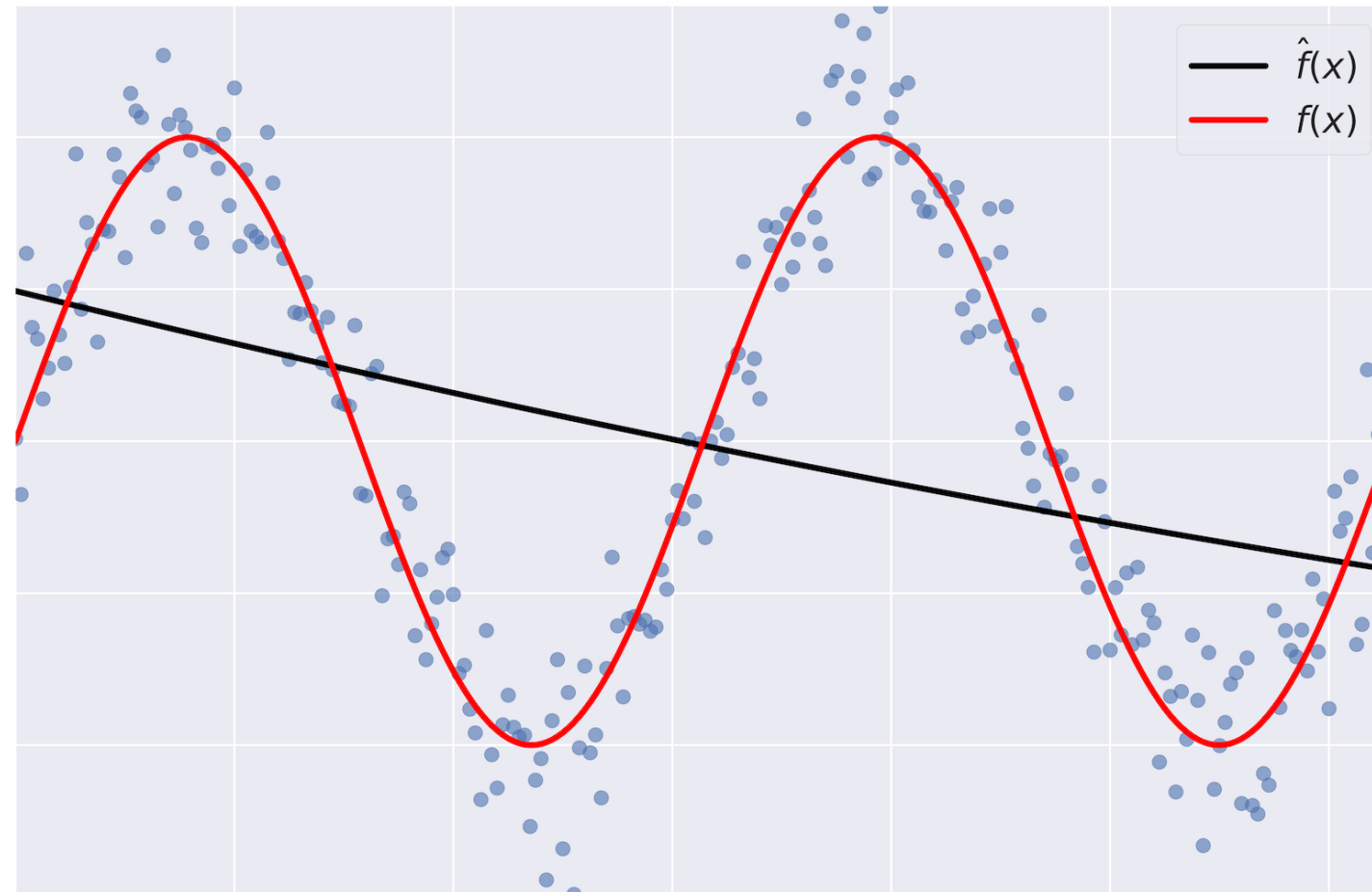
-> 복잡한 모델일수록 분산이 높고, 간단한 모델일수록 편파성이 높다(즉 모델의 복잡도에 따라 분산과 편파성의 정도가 반비례한다.)

-> 즉, 줄일 수 있는 오류를 최대한 줄여야하며, 본산과 편파성의 합이 가장 적은 모델을 설계하도록 해야한다.

Bias

- **Bias:** error term that tells you, on average, how much $\hat{f} \neq f$.

\hat{f} 가 f 가



\hat{f} f

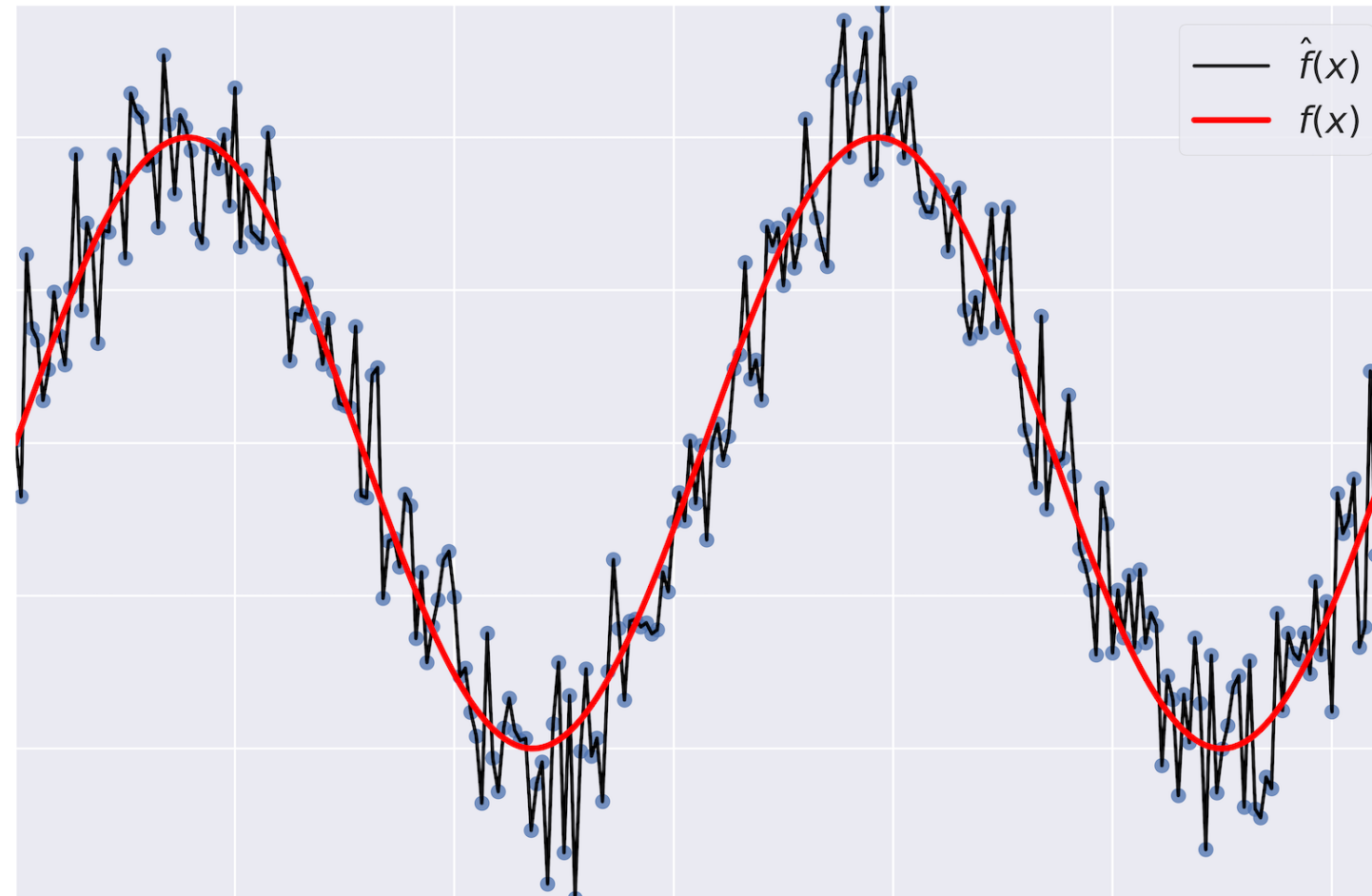
bias model underfitting

Variance

- **Variance:** tells you how much \hat{f} is inconsistent over different training sets.

training set

that



that training data point
가
f

variance() model
underfitting

Model Complexity

- **Model Complexity:** sets the flexibility of \hat{f} . model f
- Example: Maximum tree depth, Minimum samples per leaf, ... ex) 가 decision tree

Bias-Variance Tradeoff

가

model

diagram model

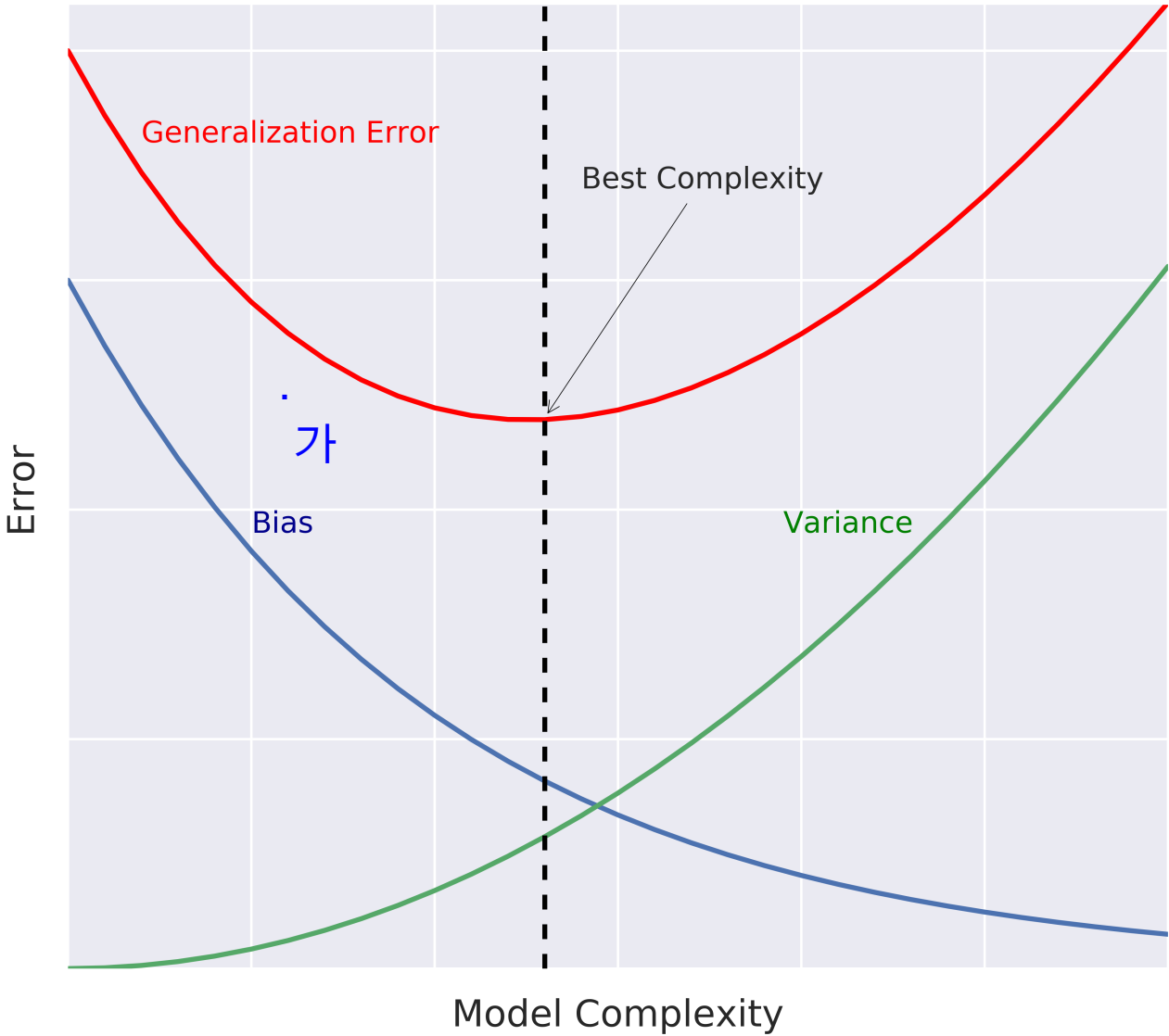
가

가

가

= >

bias - variance tradeoff



Bias-Variance Tradeoff: A Visual Explanation

that f
that
that bias variance
가
that variacne bias ,



This work by Sebastian Raschka is licensed under a Creative Commons Attribution 4.0 International License.

Let's practice!

MACHINE LEARNING WITH TREE-BASED MODELS IN PYTHON

Diagnosing Bias and Variance Problems

MACHINE LEARNING WITH TREE-BASED MODELS IN PYTHON



Elie Kawerk
Data Scientist

Estimating the Generalization Error

- How do we estimate the generalization error of a model? \hat{f} $?$ y \hat{f}
- Cannot be done directly because:
 - f is unknown,
 - usually you only have one dataset,
 - noise is unpredictable.

Estimating the Generalization Error

Solution:

- split the data to training and test sets,
- fit \hat{f} to the training set,
- evaluate the error of \hat{f} on the **unseen** test set.
- generalization error of $\hat{f} \approx$ test set error of \hat{f} . fhat test set fhat

Better Model Evaluation with Cross-Validation

- Test set should not be touched until we are confident about \hat{f} 's performance.
- Evaluating \hat{f} on training set: biased estimate, \hat{f} has already seen all training points.

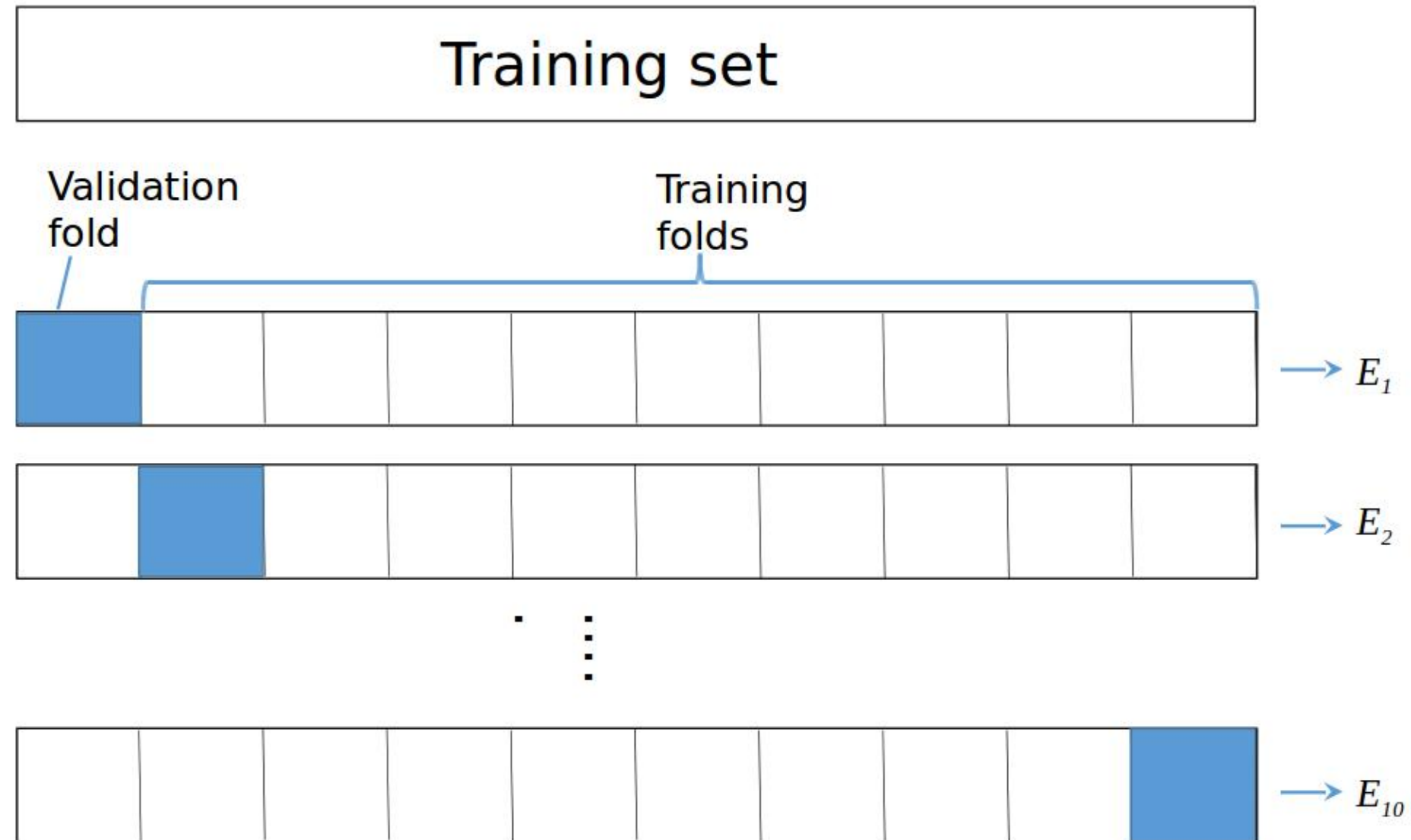
• Solution \rightarrow Cross-Validation (CV):

- K-Fold CV,
- Hold-Out CV.

that
that
가
that
CV
train set
train set
that
CV
K - Fold - CV
hold - out - cv

K-Fold CV

10



K-Fold CV

$$CV \text{ error} = \frac{E_1 + \dots + E_{10}}{10}$$

CV - error 10

Diagnose Variance Problems

- If \hat{f} suffers from **high variance**: CV error of \hat{f} $>$ training set error of \hat{f} .
- \hat{f} is said to overfit the training set. To remedy overfitting:
 - decrease model complexity,
 - for ex: decrease max depth, increase min samples per leaf, ...
 - gather more data, ..

=> \hat{f} variance train_set overfitting

Diagnose Bias Problems

- if \hat{f} suffers from high bias: CV error of $\hat{f} \approx$ ^{training set error} training set error of $\hat{f} \gg$ desired error.
- \hat{f} is said to underfit the training set. To remedy underfitting:
 - increase model complexity () ^{model} ^{feature}
 - for ex: increase max depth, decrease min samples per leaf, ...
 - gather more relevant features

K-Fold CV in sklearn on the Auto Dataset

```
from sklearn.tree import DecisionTreeRegressor
from sklearn.model_selection import train_test_split
from sklearn.metrics import mean_squared_error as MSE
from sklearn.model_selection import cross_val_score
# Set seed for reproducibility
SEED = 123
# Split data into 70% train and 30% test
X_train, X_test, y_train, y_test = train_test_split(X, y,
                                                    test_size=0.3,
                                                    random_state=SEED)
# Instantiate decision tree regressor and assign it to 'dt'
dt = DecisionTreeRegressor(max_depth=4,
                           min_samples_leaf=0.14,
                           random_state=SEED)
```

K-Fold CV in sklearn on the Auto Dataset

```
# Evaluate the list of MSE obtained by 10-fold CV
# Set n_jobs to -1 in order to exploit all CPU cores in computation
MSE_CV = - cross_val_score(dt, X_train, y_train, cv= 10,
                           cross_val_score()
                           scoring='neg_mean_squared_error',
                           scoring
                           n_jobs = -1)
# Fit 'dt' to the training set
dt.fit(X_train, y_train)
# Predict the labels of training set
y_predict_train = dt.predict(X_train)
# Predict the labels of test set
y_predict_test = dt.predict(X_test)
```

10
neg_mean_squared_error
10
CV - MSE
가 CPU

```
# CV MSE
print('CV MSE: {:.2f}'.format(MSE_CV.mean()))
```

```
CV MSE: 20.51
```

```
# Training set MSE  training set error가 CV error          dt가 train set  overfitting          variance
print('Train MSE: {:.2f}'.format(MSE(y_train, y_predict_train)))
```

```
Train MSE: 15.30
```

```
# Test set MSE
print('Test MSE: {:.2f}'.format(MSE(y_test, y_predict_test)))
```

```
Test MSE: 20.92
```

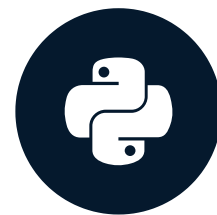
CV test set error가

Let's practice!

MACHINE LEARNING WITH TREE-BASED MODELS IN PYTHON

Ensemble Learning

MACHINE LEARNING WITH TREE-BASED MODELS IN PYTHON



Elie Kawerk
Data Scientist

Advantages of CARTs

- Simple to understand.
- Simple to interpret.
- Easy to use.
- Flexibility: ability to describe non-linear dependencies.
- Preprocessing: no need to standardize or normalize features, ...

Limitations of CARTs

- Classification: can only produce orthogonal decision boundaries.
- Sensitive to small variations in the training set. train set point가 CART
가
- High variance: unconstrained CARTs may overfit the training set. train set
- Solution: ensemble learning.

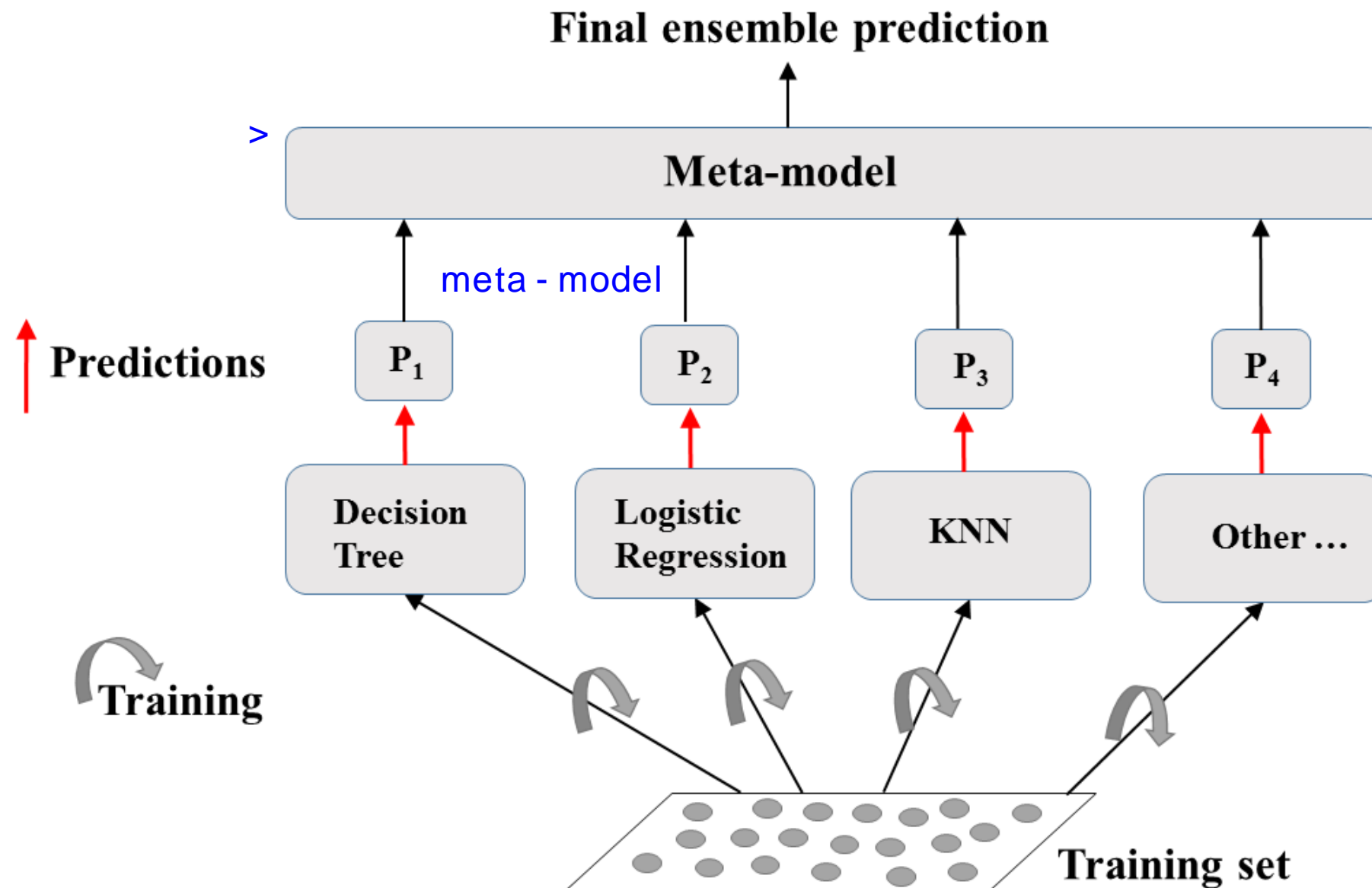
Ensemble Learning

- Train different models on the same dataset.
- Let each model make its predictions.
- Meta-model: aggregates predictions of individual models. meta - model model
- Final prediction: more robust and less prone to errors. 가
- Best results: models are skillful in different ways.

(meta - model)

Ensemble Learning: A Visual Explanation

- < 1. train set
- 2.
- 3.



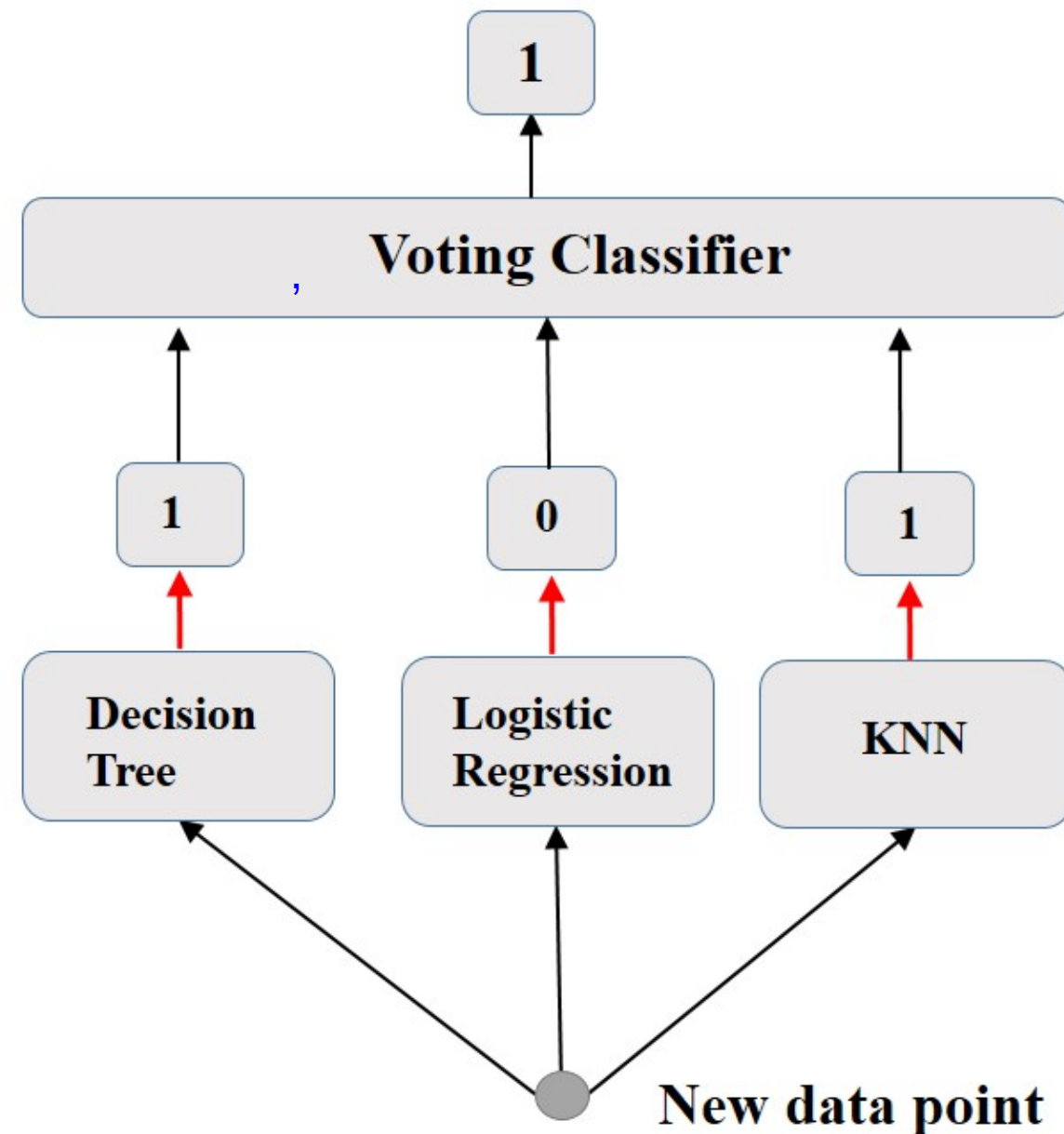
Ensemble Learning in Practice: Voting Classifier

- Binary classification task.
- N classifiers make predictions: P_1, P_2, \dots, P_N with $P_i = 0$ or 1 .
- Meta-model prediction: hard voting. meta - model hard voting

Hard Voting

diagram 3 train
1 3 data point 1 label
2 label 0
=> 1 2 , 2 1
=>

1
Predictions



Voting Classifier in sklearn (Breast-Cancer dataset)

```
# Import functions to compute accuracy and split data
from sklearn.metrics import accuracy_score
from sklearn.model_selection import train_test_split

# Import models, including VotingClassifier meta-model
from sklearn.linear_model import LogisticRegression
from sklearn.tree import import DecisionTreeClassifier
from sklearn.neighbors import import KNeighborsClassifier as KNN
from sklearn.ensemble import import VotingClassifier

# Set seed for reproducibility
SEED = 1
```

Voting Classifier in sklearn (Breast-Cancer dataset)

```
# Split data into 70% train and 30% test
X_train, X_test, y_train, y_test = train_test_split(X, y,
                                                    test_size= 0.3,
                                                    random_state= SEED)

# Instantiate individual classifiers
lr = LogisticRegression(random_state=SEED)
knn = KNN()
dt = DecisionTreeClassifier(random_state=SEED)
# Define a list called classifier that contains the tuples (classifier_name, classifier)
classifiers = [('Logistic Regression', lr),
               ('K Nearest Neighbours', knn),
               ('Classification Tree', dt)]
```

```
# Iterate over the defined list of tuples containing the classifiers
for clf_name, clf in classifiers:
    #fit clf to the training set
    clf.fit(X_train, y_train)

    # Predict the labels of the test set
    y_pred = clf.predict(X_test)

    # Evaluate the accuracy of clf on the test set
    print('{:s} : {:.3f}'.format(clf_name, accuracy_score(y_test, y_pred)))
```

```
Logistic Regression: 0.947
K Nearest Neighbours: 0.930
Classification Tree: 0.930
```

Voting Classifier in sklearn (Breast-Cancer dataset)

```
# Instantiate a VotingClassifier 'vc'
vc = VotingClassifier(estimators=classifiers)

# Fit 'vc' to the training set and predict test set labels
vc.fit(X_train, y_train)
y_pred = vc.predict(X_test)

# Evaluate the test-set accuracy of 'vc'
print('Voting Classifier: {:.3f}'.format(accuracy_score(y_test, y_pred)))
```

```
Voting Classifier: 0.953
```

=> 가

Let's practice!

MACHINE LEARNING WITH TREE-BASED MODELS IN PYTHON