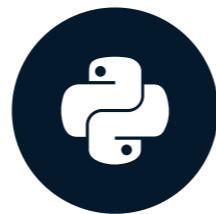


# Decision-Tree for Classification

MACHINE LEARNING WITH TREE-BASED MODELS IN PYTHON



Elie Kawerk

Data Scientist

# Course Overview

- **Chap 1:** Classification And Regression Tree (CART)
- **Chap 2:** The Bias-Variance Tradeoff
- **Chap 3:** Bagging and Random Forests
- **Chap 4:** Boosting
- **Chap 5:** Model Tuning

# Classification-tree

- Sequence of if-else questions about individual features. label
  - **Objective:** infer class labels.
  - Able to capture non-linear relationships between features and labels. linear model
  - Don't require feature scaling (ex: Standardization, ..) ,

가 decision tree  
feature if - else

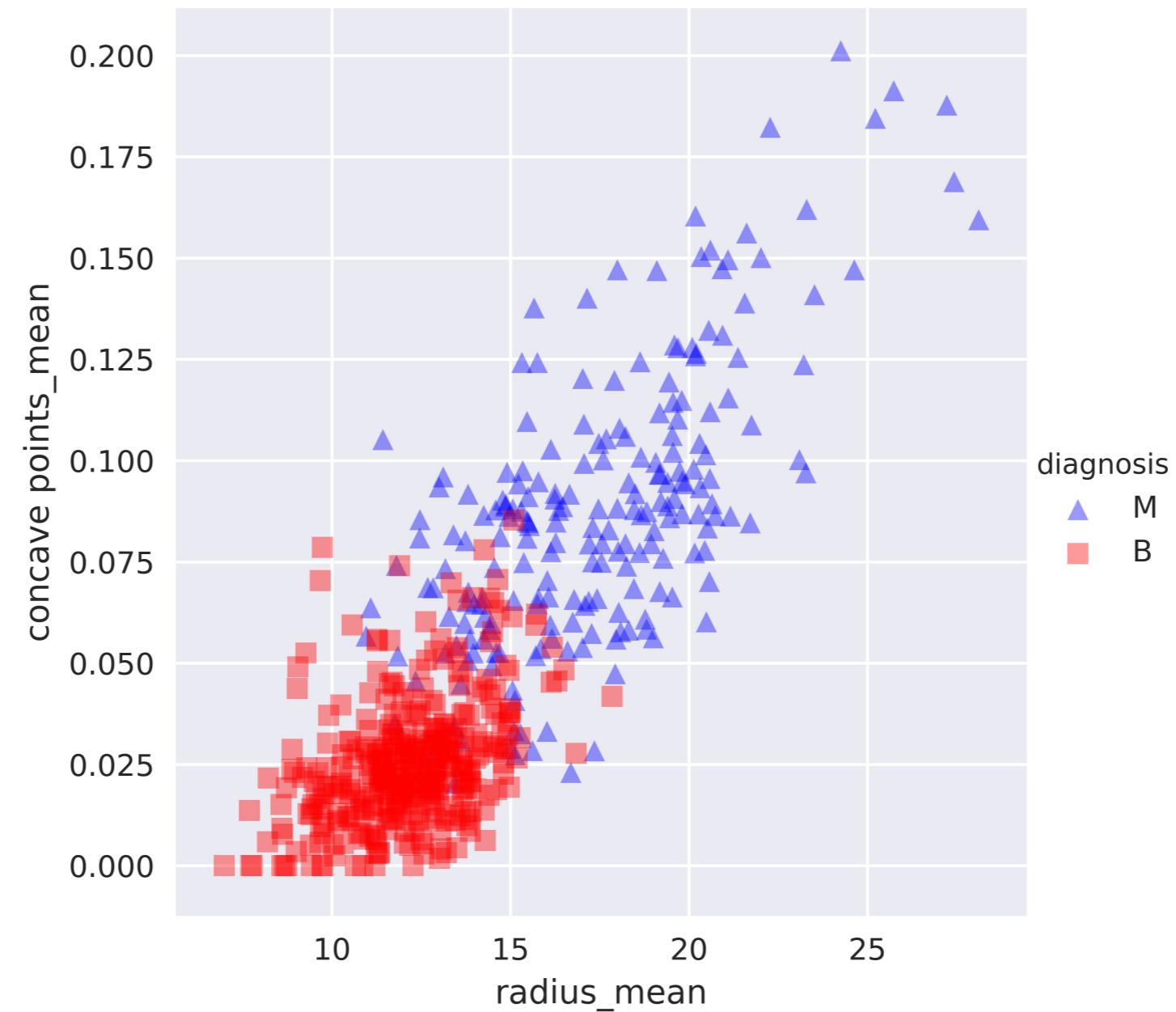
가



# Breast Cancer Dataset in 2D

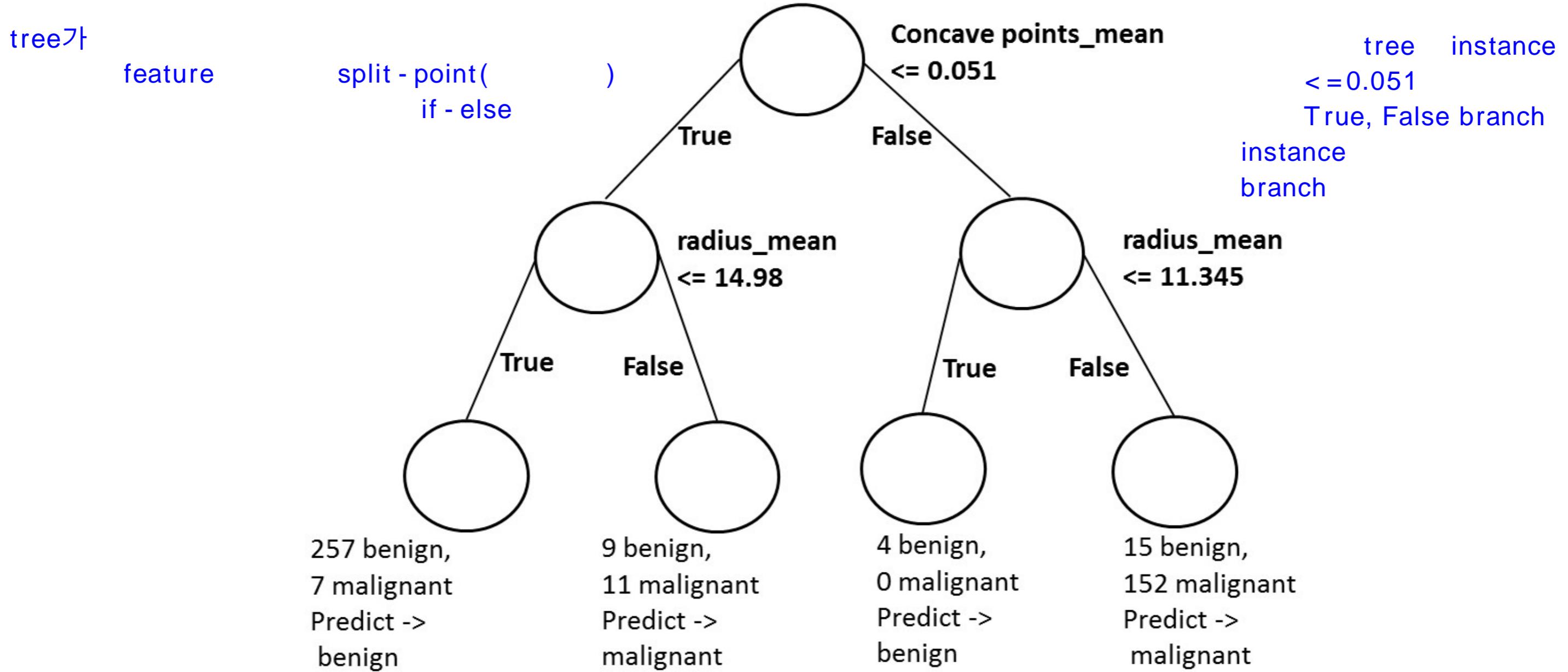
Tree

scatter



blue  
red

# Decision-tree Diagram



# Classification-tree in scikit-learn

```
# Import DecisionTreeClassifier
from sklearn.tree import DecisionTreeClassifier
# Import train_test_split
from sklearn.model_selection import train_test_split
# Import accuracy_score
from sklearn.metrics import accuracy_score
# Split the dataset into 80% train, 20% test
X_train, X_test, y_train, y_test= train_test_split(X, y,
                                                    split
                                                    dataset
                                                    test_size=0.2,
                                                    stratify=y,
                                                    random_state=1)
# Instantiate dt
dt = DecisionTreeClassifier(max_depth=2, random_state=1) max_depth      2
                                                    DT
                                                    random_state
```

# Classification-tree in scikit-learn

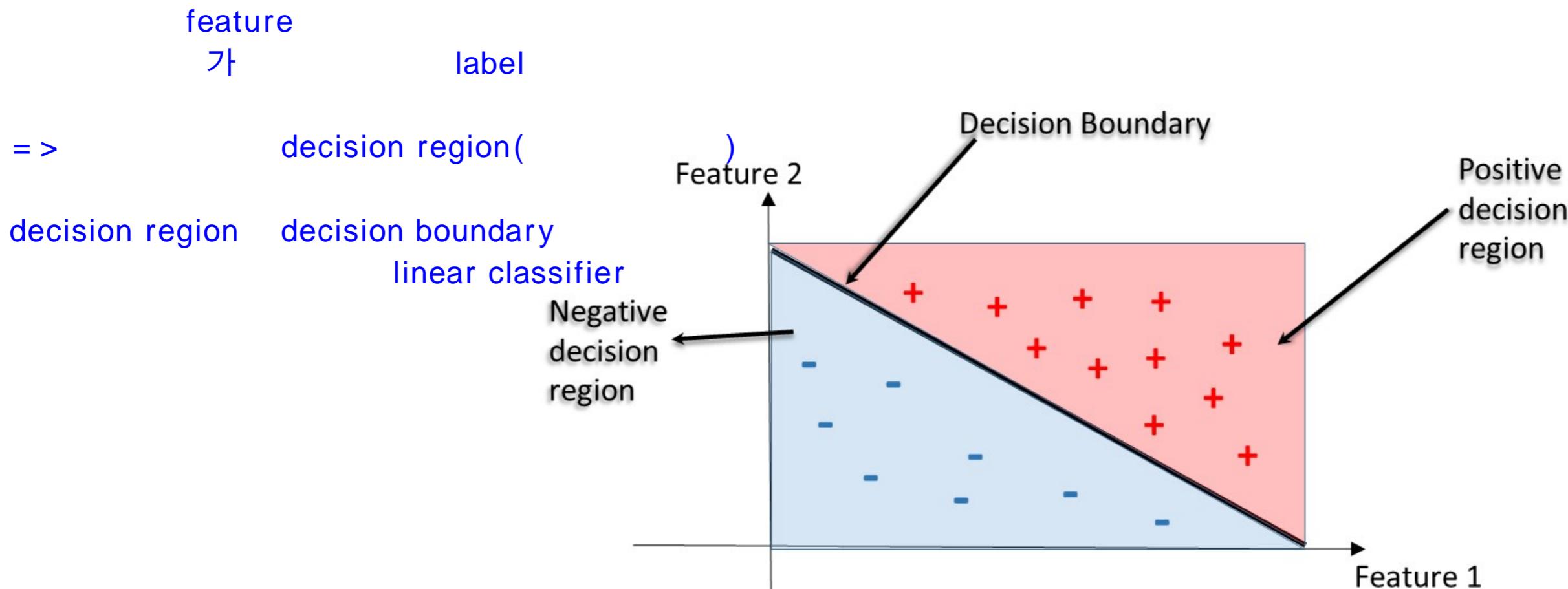
```
# Fit dt to the training set  
dt.fit(X_train,y_train)  
  
# Predict the test set labels  
y_pred = dt.predict(X_test)  
# Evaluate the test-set accuracy  
accuracy_score(y_test, y_pred)
```

```
0.90350877192982459
```

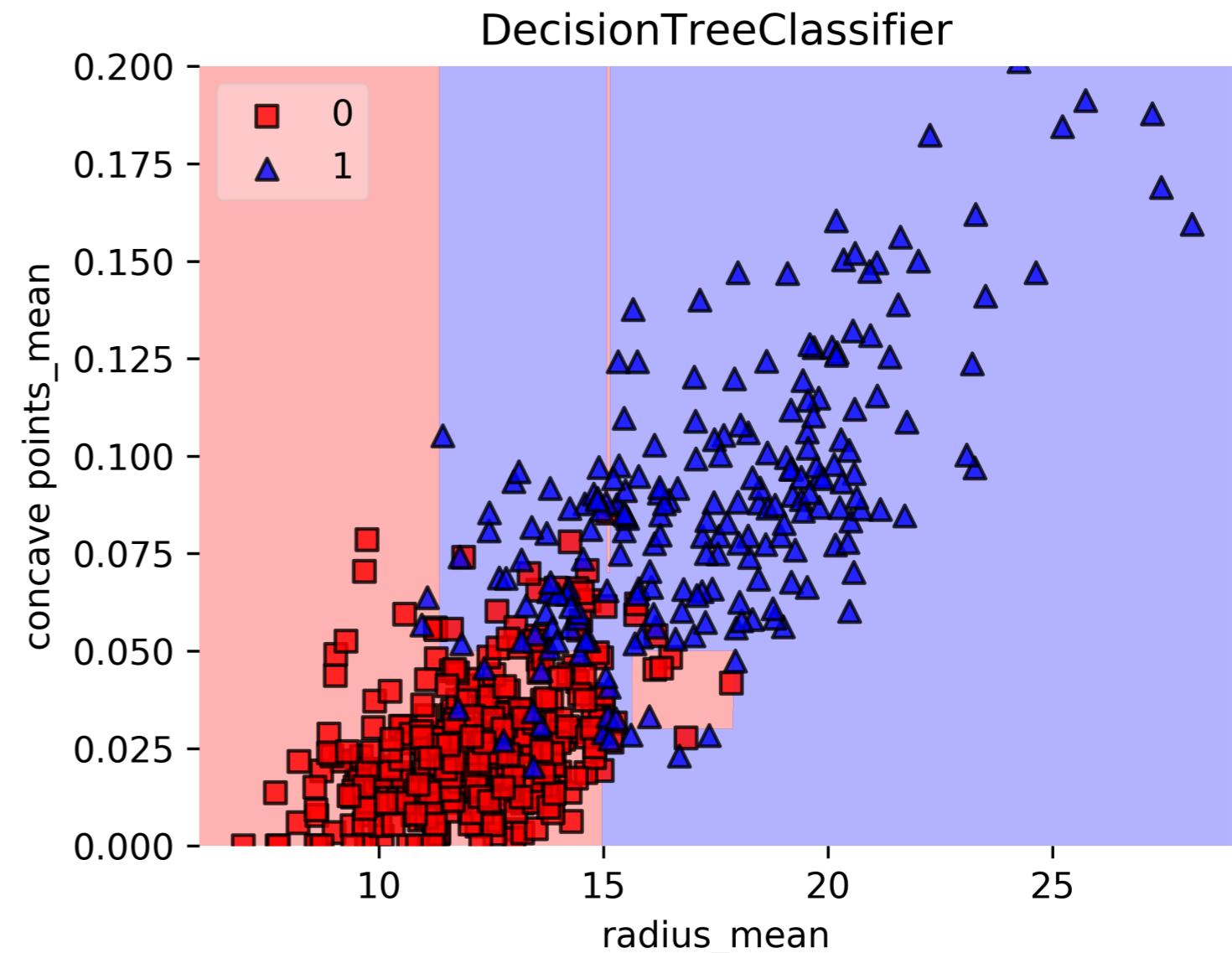
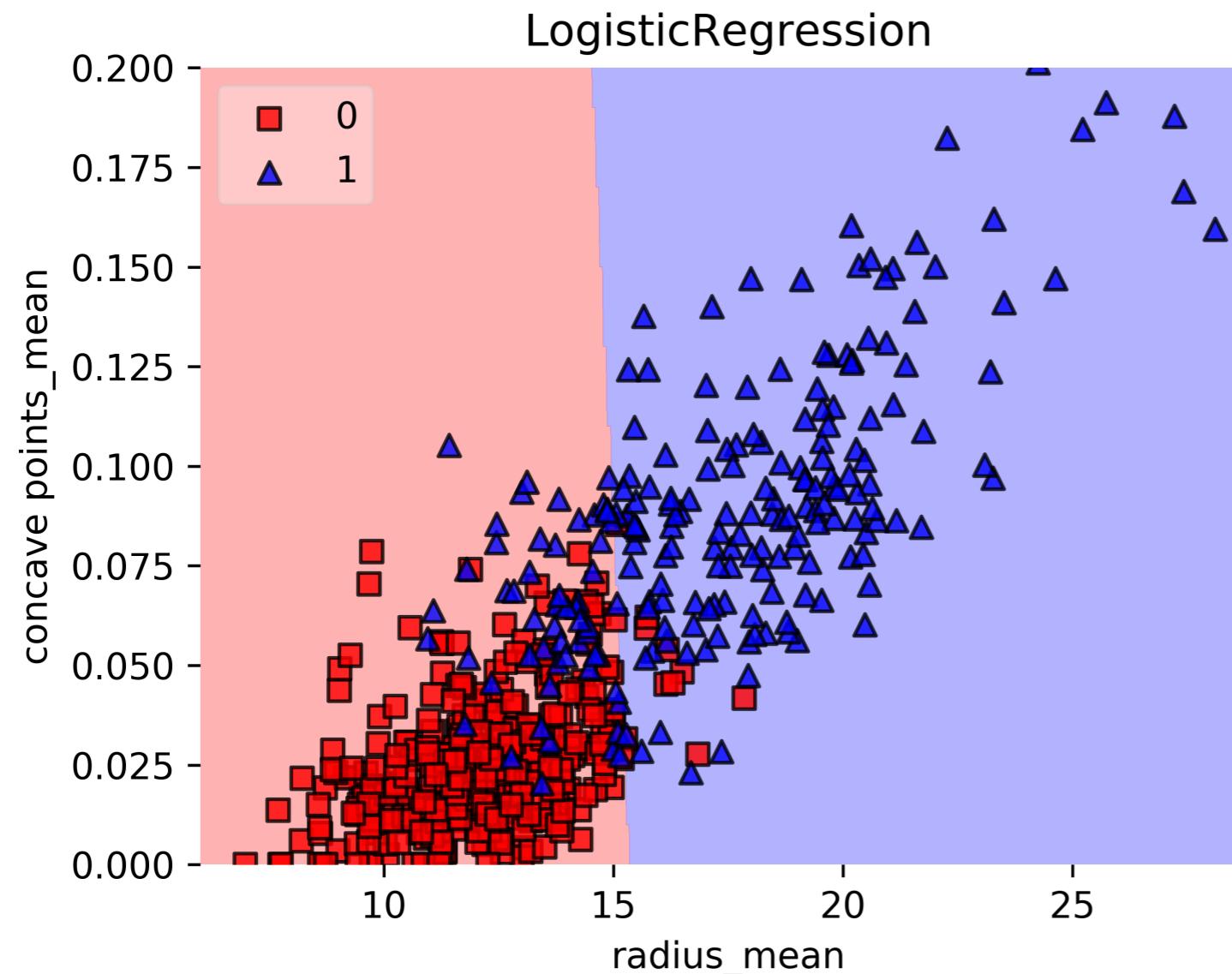
# Decision Regions

**Decision region:** region in the feature space where all instances are assigned to one class label.

**Decision Boundary:** surface separating different decision regions.



# Decision Regions: CART vs. Linear Model



DT  
feature  
tree

decision - region  
feature

# **Let's practice!**

**MACHINE LEARNING WITH TREE-BASED MODELS IN PYTHON**

# Classification-Tree Learning

MACHINE LEARNING WITH TREE-BASED MODELS IN PYTHON



Elie Kawerk

Data Scientist

# Building Blocks of a Decision-Tree

- **Decision-Tree:** data structure consisting of a hierarchy of nodes.
- **Node:** question or prediction.

DT node

Node:

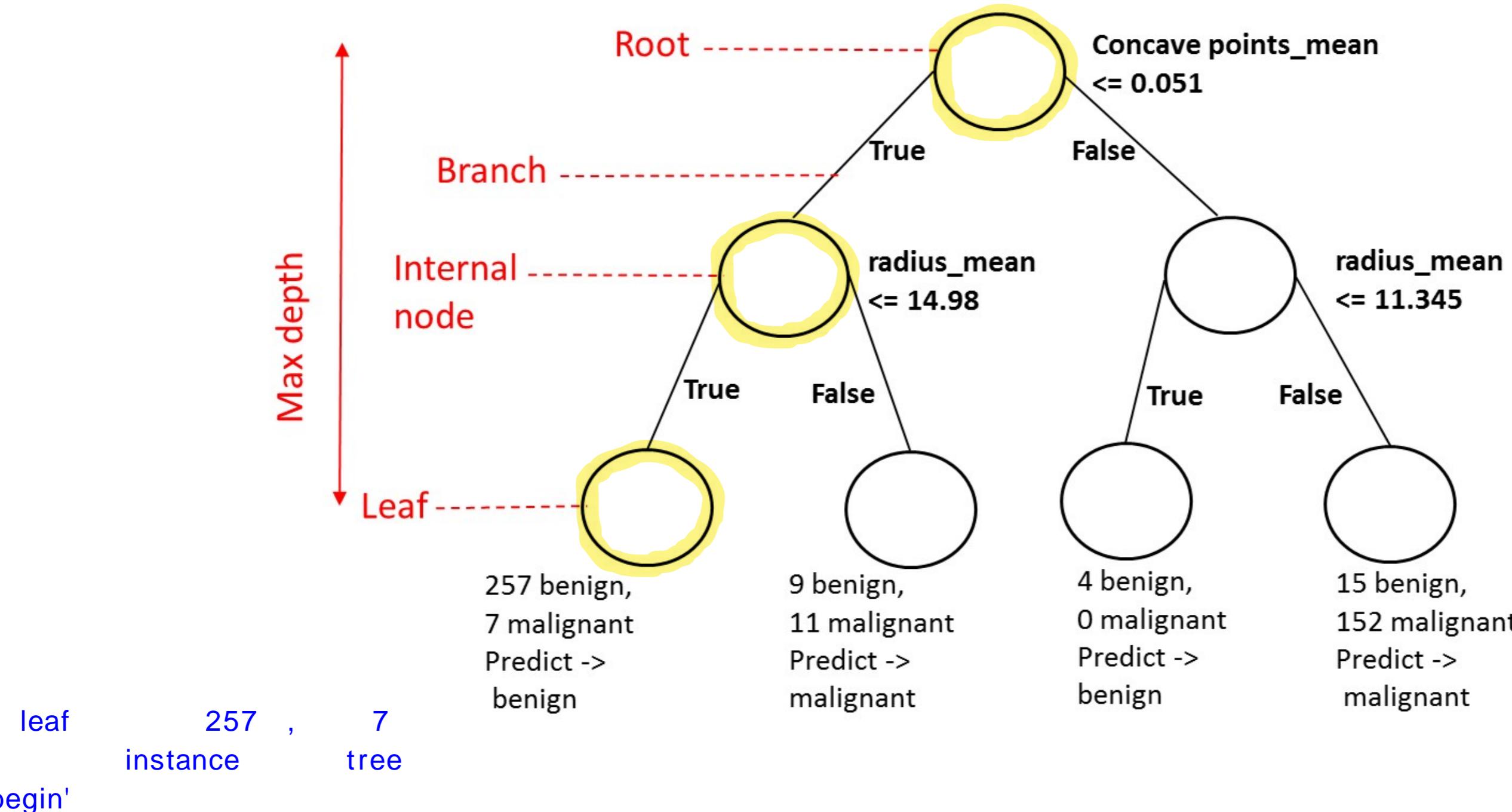
# Building Blocks of a Decision-Tree

Three kinds of nodes:

- **Root:** *no* parent node, question giving rise to *two* children nodes.
- **Internal node:** *one* parent node, question giving rise to *two* children nodes.
- **Leaf:** *one* parent node, *no* children nodes --> *prediction*.

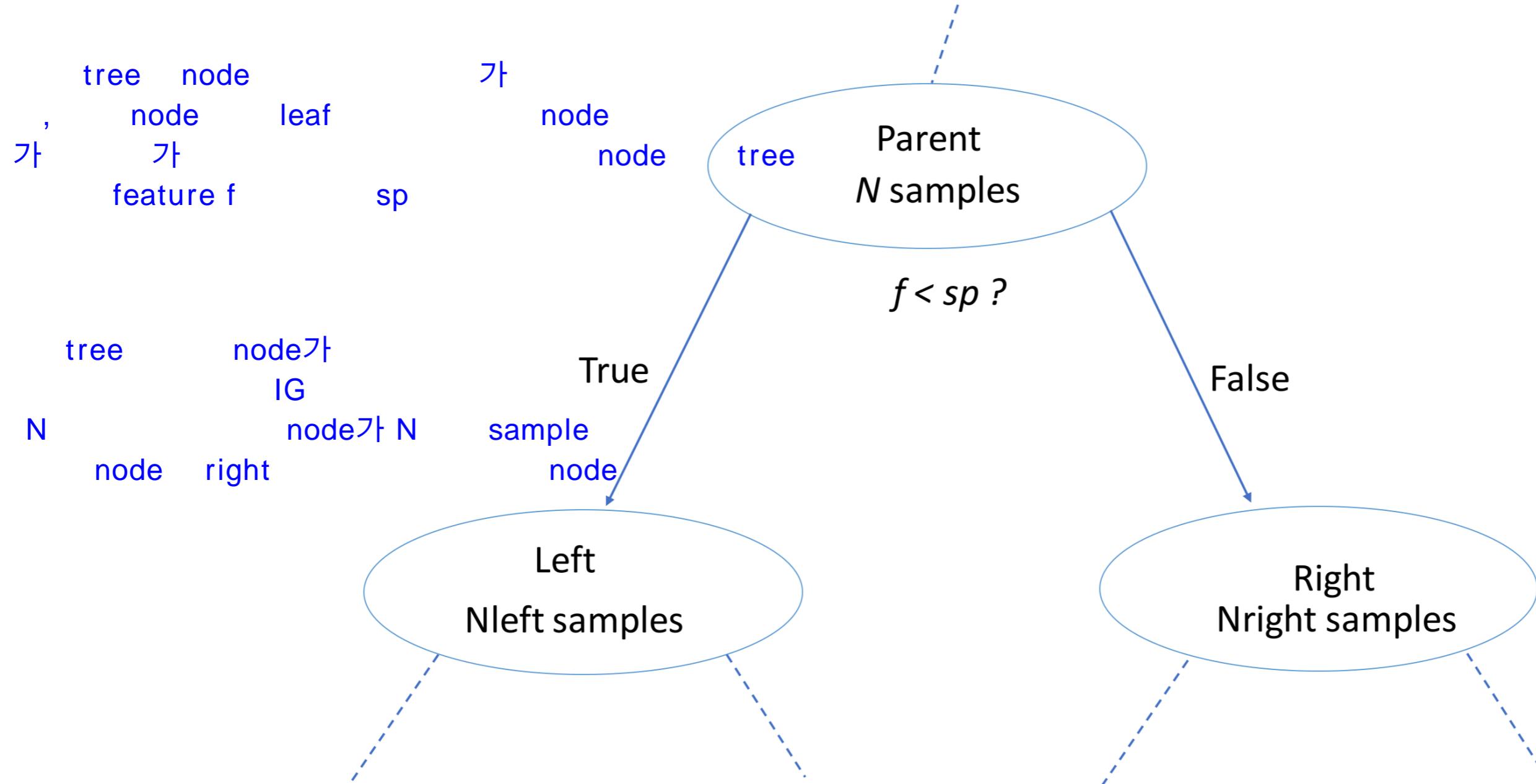


# Prediction



'begin'

# Information Gain (IG)



# Information Gain (IG)

$$IG(\underbrace{f}_{feature}, \underbrace{sp}_{split-point}) = I(parent) - \left( \frac{N_{left}}{N} I(left) + \frac{N_{right}}{N} I(right) \right)$$

Criteria to measure the impurity of a node  $I(node)$ :

- gini index,
- entropy. ...

=> gini      entropy  
                        : "node"      가?"

# Classification-Tree Learning

- Nodes are grown recursively.
- At each node, split the data based on:
  - feature  $f$  and split-point  $sp$  to maximize  $IG(\text{node})$ .
- If  $IG(\text{node}) = 0$ , declare the node a leaf. ...

tree가

node가

가

non - leaf node  
sp

node  
- >

IG가 null  
node  
tree

ex)  
null

2  
2  
node  
node가 leaf  
IG가

```
# Import DecisionTreeClassifier
from sklearn.tree import DecisionTreeClassifier
# Import train_test_split
from sklearn.model_selection import train_test_split
# Import accuracy_score
from sklearn.metrics import accuracy_score
# Split dataset into 80% train, 20% test
X_train, X_test, y_train, y_test= train_test_split(X, y,
                                                    test_size=0.2,
                                                    stratify=y,
                                                    random_state=1)
# Instantiate dt, set 'criterion' to 'gini'
dt = DecisionTreeClassifier(criterion='gini', random_state=1)
```

# Information Criterion in scikit-learn

```
# Fit dt to the training set  
dt.fit(X_train,y_train)  
  
# Predict test-set labels  
y_pred= dt.predict(X_test)  
  
# Evaluate test-set accuracy  
accuracy_score(y_test, y_pred)
```

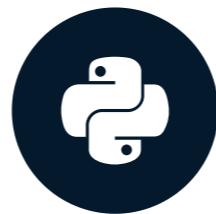
```
0.92105263157894735
```

# **Let's practice!**

**MACHINE LEARNING WITH TREE-BASED MODELS IN PYTHON**

# Decision-Tree for Regression

MACHINE LEARNING WITH TREE-BASED MODELS IN PYTHON



Elie Kawerk

Data Scientist

, model output

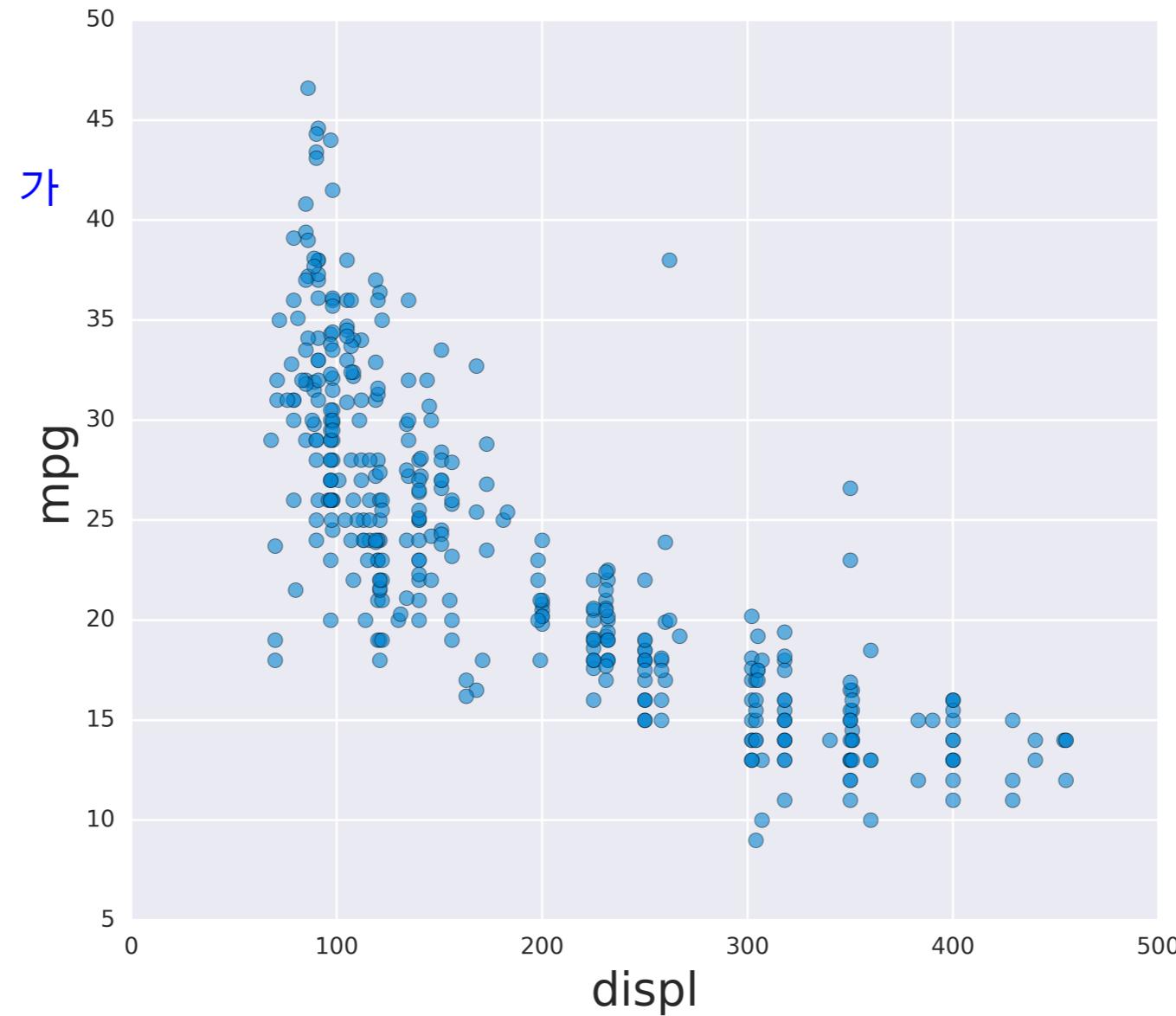
# Auto-mpg Dataset

dataset  
mpg  
=> : 6가 feature MPG

	mpg	displ	hp	weight	accel	origin	size
0	18.0	250.0	88	3139	14.5	US	15.0
1	9.0	304.0	193	4732	18.5	US	20.0
2	36.1	91.0	60	1800	16.4	Asia	10.0
3	18.5	250.0	98	3525	19.0	US	15.0
4	34.3	97.0	78	2188	15.8	Europe	10.0
5	32.9	119.0	100	2615	14.8	Asia	10.0

# Auto-mpg with one feature

mpg    displ    2D scatter    mpg



# Regression-Tree in scikit-learn

```
# Import DecisionTreeRegressor
from sklearn.tree import DecisionTreeRegressor
# Import train_test_split
from sklearn.model_selection import train_test_split
# Import mean_squared_error as MSE
from sklearn.metrics import mean_squared_error as MSE
# Split data into 80% train and 20% test
X_train, X_test, y_train, y_test= train_test_split(X, y,
                                                    test_size=0.2,
                                                    random_state=3)
# Instantiate a DecisionTreeRegressor 'dt'
dt = DecisionTreeRegressor(max_depth=4,
                           min_samples_leaf=0.1,
                           random_state=3)
```

파라미터 명	설명
min_samples_split	- 노드를 분할하기 위한 최소한의 샘플 데이터수 → 과적합을 제어하는데 사용 - Default = 2 → 작게 설정할 수록 분할 노드가 많아져 과적합 가능성 증가
min_samples_leaf	- 리프노드가 되기 위해 필요한 최소한의 샘플 데이터수 - min_samples_split과 함께 과적합 제어 용도 - 불균형 데이터의 경우 특정 클래스의 데이터가 극도로 작을 수 있으므로 작게 설정 필요
max_features	- 최적의 분할을 위해 고려할 최대 feature 개수 - Default = None → 데이터 세트의 모든 피처를 사용 - int형으로 지정 → 피처 갯수 / float형으로 지정 → 비중 - sqrt 또는 auto : 전체 피처 중 $\sqrt{피처개수}$ 만큼 선정 - log : 전체 피처 중 $\log_2(\text{전체 피처 개수})$ 만큼 선정
max_depth	- 트리의 최대 깊이 - default = None → 완벽하게 클래스 값이 결정될 때 까지 분할 또는 데이터 개수가 min_samples_split보다 작아질 때 까지 분할 - 깊이가 깊어지면 과적합될 수 있으므로 적절히 제어 필요
max_leaf_nodes	리프노드의 최대 개수

leaf가  
10%

# Regression-Tree in scikit-learn

```
# Fit 'dt' to the training-set  
dt.fit(X_train, y_train)  
# Predict test-set labels  
y_pred = dt.predict(X_test)      test_set  
# Compute test-set MSE  
mse_dt = MSE(y_test, y_pred)      가  
# Compute test-set RMSE  
rmse_dt = mse_dt**(1/2)          1/2  
# Print rmse_dt  
print(rmse_dt)                  dt  test set rmse      5.1
```

5.1023068889

# Information Criterion for Regression-Tree

$$I(\text{node}) = \underbrace{\text{MSE}(\text{node})}_{\text{mean-squared-error}} = \frac{1}{N_{\text{node}}} \sum_{i \in \text{node}} (y^{(i)} - \hat{y}_{\text{node}})^2$$

tree가 dataset  
node

node

$\hat{y}_{\text{node}}$

$$= \frac{1}{N_{\text{node}}} \sum_{i \in \text{node}} y^{(i)}$$

leaf

가 가

mean-target-value

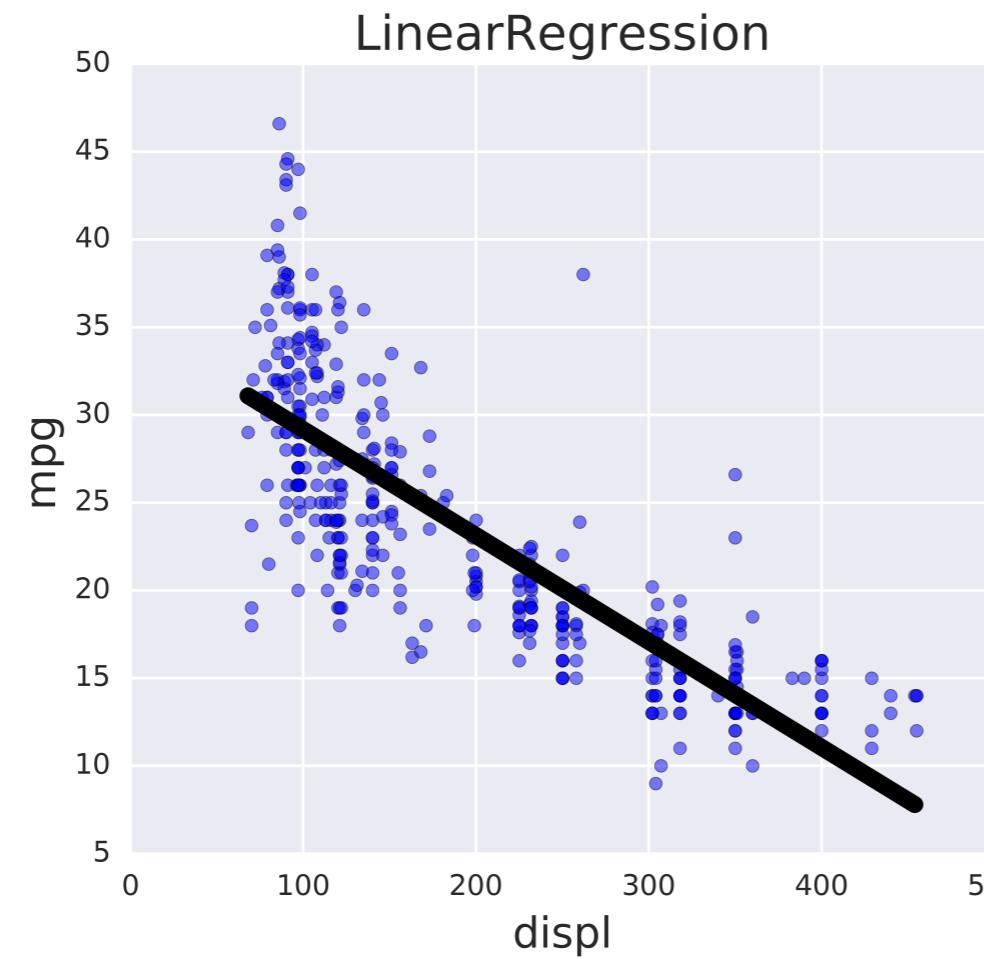
, 가 가 가

# Prediction

$$\hat{y}_{pred}(\text{leaf}) = \frac{1}{N_{\text{leaf}}} \sum_{i \in \text{leaf}} y^{(i)}$$

instance 가 tree 가 leaf 'y'  
leaf

# Linear Regression vs. Regression-Tree

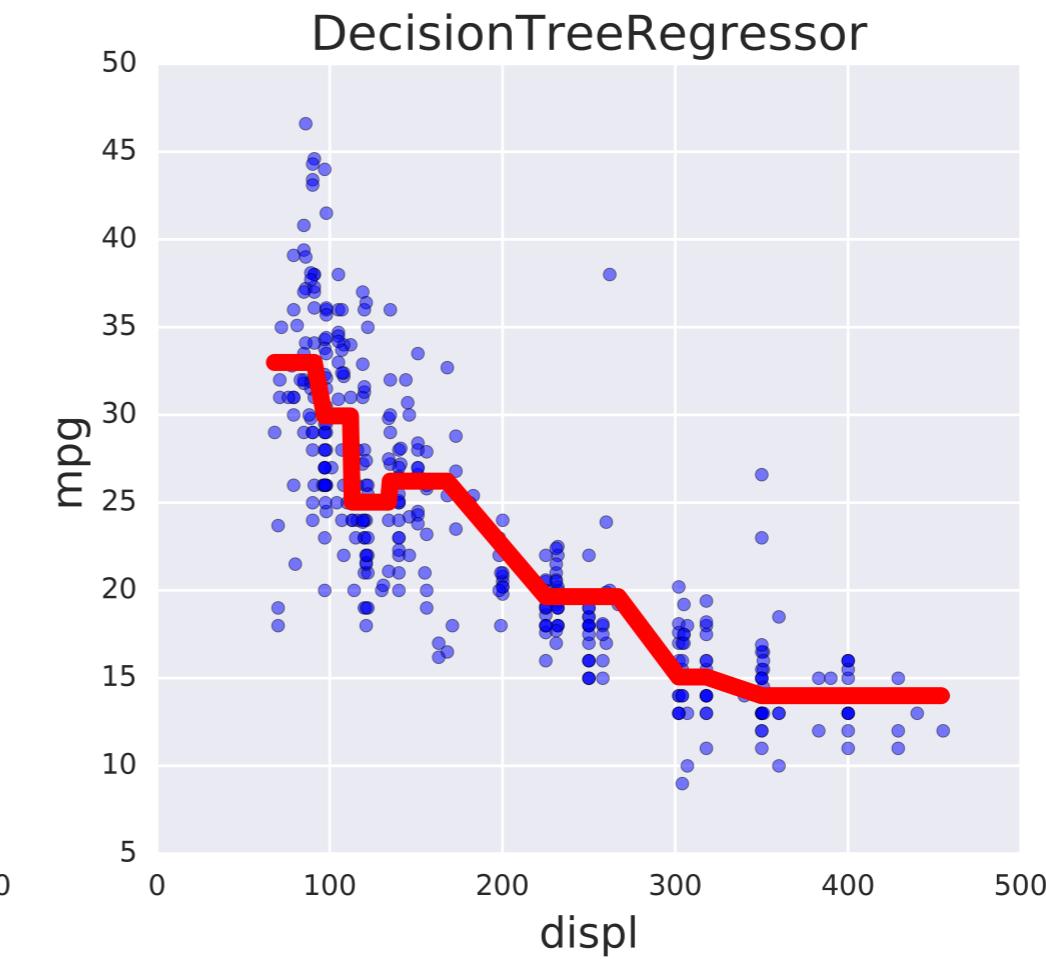


tree

data scatter

(

regression tree



)

scatter

# **Let's practice!**

**MACHINE LEARNING WITH TREE-BASED MODELS IN PYTHON**