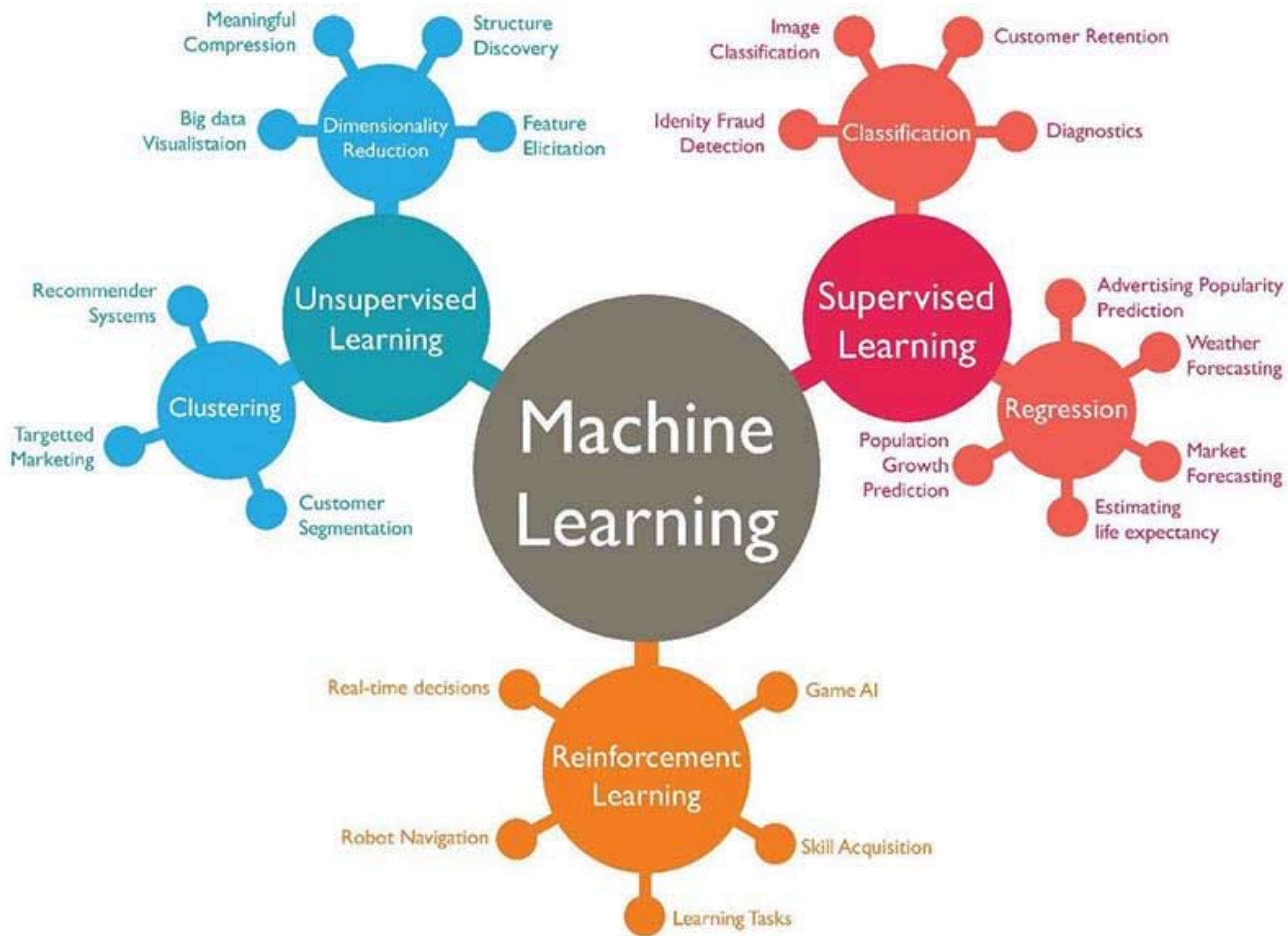


# Machine Learning

# Intro

# Machine Learning?

기계(즉 컴퓨터)를 인간처럼 학습시켜 스스로 규칙을 형성하고자 하는 시도에서 시작되었다



크게 Supervised Learning과 Unsupervised Learning으로 나뉘며,  
이 중 우리는 Supervised Learning, 그 중에서도 Tree 기반의 알고리즘을 집중적으로 다룰 것이다

## Machine Learning Algorithms *(sample)*

	<u>Unsupervised</u>	<u>Supervised</u>
<u>Continuous</u>	<ul style="list-style-type: none"><li>• Clustering &amp; Dimensionality Reduction<ul style="list-style-type: none"><li>○ SVD</li><li>○ PCA</li><li>○ K-means</li></ul></li></ul>	<ul style="list-style-type: none"><li>• Regression<ul style="list-style-type: none"><li>○ Linear</li><li>○ Polynomial</li></ul></li><li>• Decision Trees</li><li>• Random Forests</li></ul>
<u>Categorical</u>	<ul style="list-style-type: none"><li>• Association Analysis<ul style="list-style-type: none"><li>○ Apriori</li><li>○ FP-Growth</li></ul></li><li>• Hidden Markov Model</li></ul>	<ul style="list-style-type: none"><li>• Classification<ul style="list-style-type: none"><li>○ KNN</li><li>○ Trees</li><li>○ Logistic Regression</li><li>○ Naive-Bayes</li><li>○ SVM</li></ul></li></ul>



# Machine Learning의 분류

크게 Supervised Learning과 Unsupervised Learning으로 나뉘며,  
이 중 우리는 Supervised Learning, 그 중에서도 Tree 기반의 알고리즘을 집중적으로 다룰 것이다

## Machine Learning Algorithms (sample)

	<u>Unsupervised</u>	<u>Supervised</u>
<u>Continuous</u>	<ul style="list-style-type: none"><li>• Clustering &amp; Dimensionality Reduction<ul style="list-style-type: none"><li>○ SVD</li><li>○ PCA</li><li>○ K-means</li></ul></li></ul>	<ul style="list-style-type: none"><li>• Regression<ul style="list-style-type: none"><li>○ Linear</li><li>○ Polynomial</li></ul></li><li>• Decision Trees</li><li>• Random Forests</li></ul>
<u>Categorical</u>	<ul style="list-style-type: none"><li>• Association Analysis<ul style="list-style-type: none"><li>○ Apriori</li><li>○ FP-Growth</li></ul></li><li>• Hidden Markov Model</li></ul>	<ul style="list-style-type: none"><li>• Classification<ul style="list-style-type: none"><li>○ KNN</li><li>○ Trees</li><li>○ Logistic Regression</li><li>○ Naive-Bayes</li><li>○ SVM</li></ul></li></ul>

이 부분(+Gradient Boosting Machine)이  
이번 수업의 핵심이 될 예정

# Structured vs Unstructured Dataset

데이터는 구조화된(Strurctured) 데이터와 비구조화된 데이터(Unstructured)로 나뉘어지며, 우리는 이 중 구조화된 데이터에서 가장 좋은 성능이 나오는 알고리즘을 중점적으로 다룰 것이다.

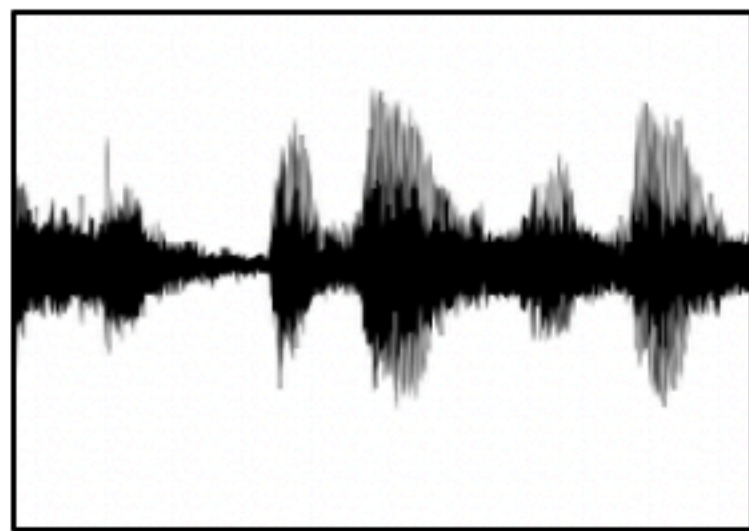
## Supervised Learning

### Structured Data

Size	#bedrooms	...	Price (1000\$s)
2104	3		400
1600	3		330
2400	3		369
⋮	⋮		⋮
3000	4		540

User Age	Ad Id	...	Click
41	93242		1
80	93287		0
18	87312		1
⋮	⋮		⋮
27	71244		1

### Unstructured Data



Audio



Image

Four scores and seven  
years ago...

Text

가장 기본적인 알고리즘인 Decision Tree부터 시작하여,  
가장 최신의 머신러닝 알고리즘과, 이를 실전에 활용하는 방법까지 모두 다루게 될 것.

## Lecture 1: Decision Tree

- Decision Tree의 기본적인 원리
- 이를 파이썬으로 직접 구현하기
- Decision Tree에 있는 Hyperparameter를 배우기

## Lecture 3: Practical usage of GBM

- GBM을 가장 실용적으로 사용할 수 있는 여러 툴을 리뷰하기
- 이 중 LightGBM의 사용 방법을 숙지하기
- LightGBM의 주요 하이퍼패러미터를 살펴보기

## Lecture 2: Ensemble Methods

- Decision Tree의 성능을 올릴 수 있는 앙상블 방식을 배우기
- 배깅(Bagging, Bootstrap Aggregating)를 직접 구현하기
- 그래디언트 부스팅(Gradient Boosting, GBM)을 직접 구현하기

## Lecture 4: Otto Group Product Classification Challenge

- LightGBM을 이용하여 상품을 분류하는 캐글 경진대회에 참여하기
- 순수 Hyperparameter Tuning만으로 상위 10%에 근접하기

# Decision Tree



# Decision Tree

가장 기본적인 머신러닝 알고리즘 중 하나  
사람이 가장 이해하기 쉬운 머신러닝 알고리즘 중 하나이며, 그 성능도 좋은 편에 속한다

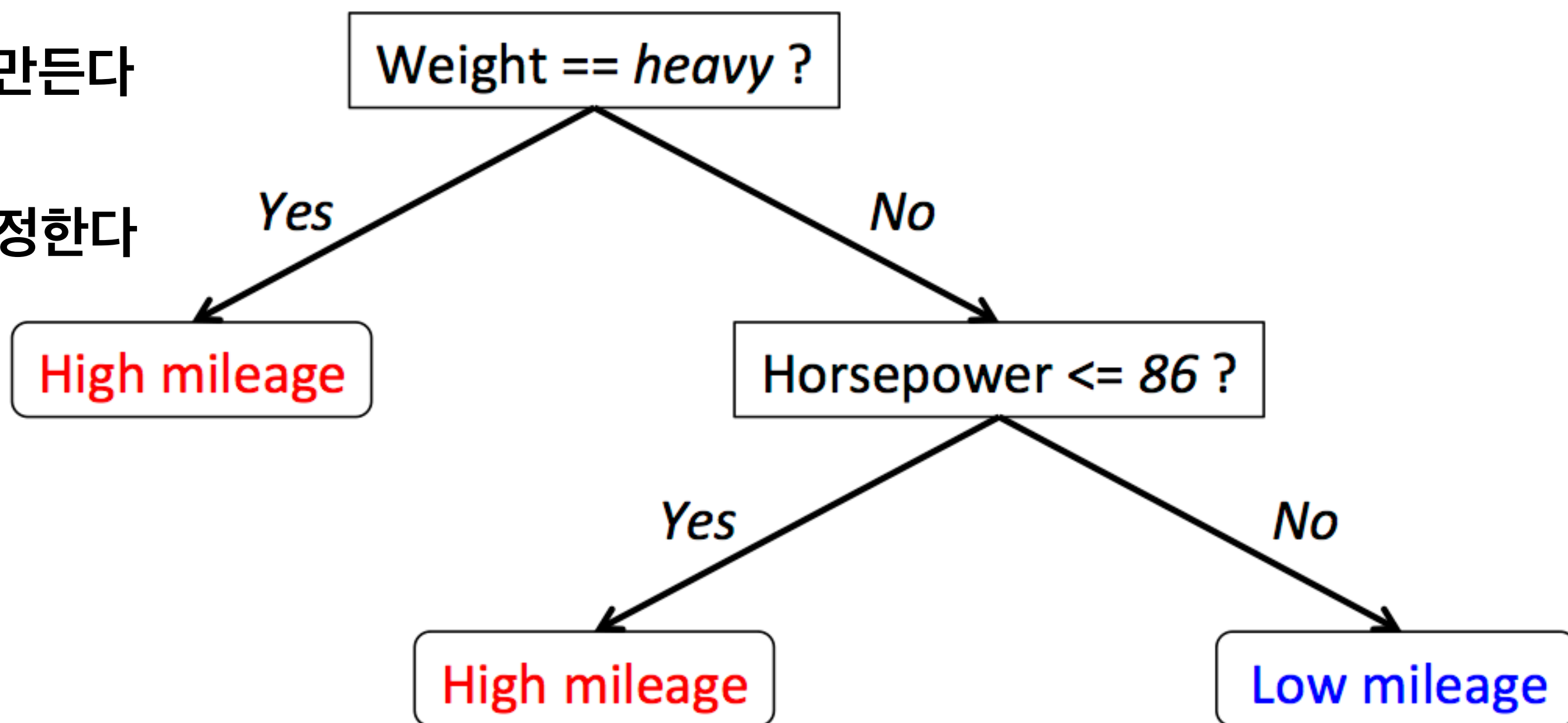
## 동작 원리

1. feature와 label를 받는다
2. feature를 활용해 만들 수 있는 모든 조건(condition)을 만든다
3. 모든 조건마다의 평균 gini impurity를 구한다
4. 평균 gini impurity가 가장 낮은 조건을 좋은 조건이라 가정한다
5. 이 조건을 활용해 나무의 가지를 친다
6. 이를 gini impurity가 더 낮아지지 않을 때 까지 반복한다

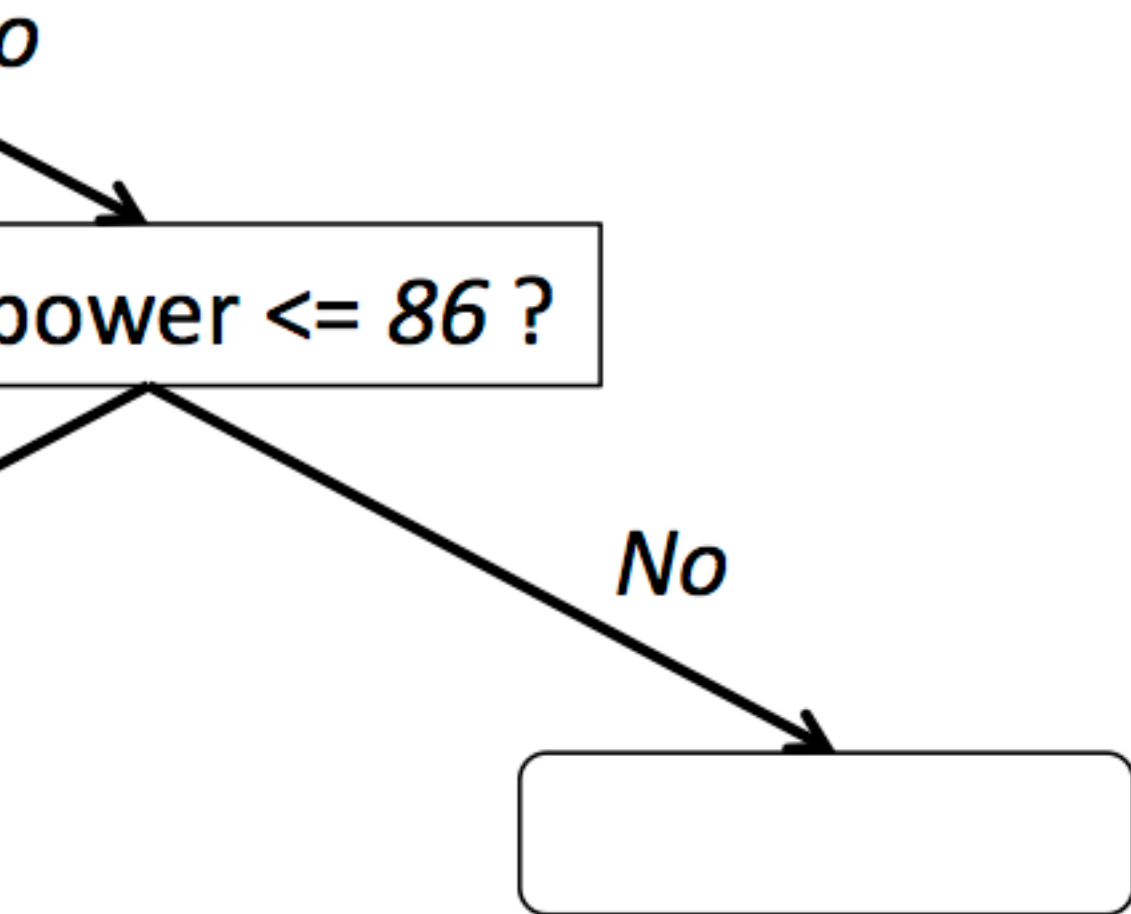
$$I_G(p) = p_t(1 - p_t) + p_f(1 - p_f)$$

gini impurity  
(차후에 자세히 설명)

## Decision Tree Model for Car Mileage Prediction



# Gini Impurity



이 조건에 100개의 데이터가 있을 경우

데이터의 불순도(impurity)를 측정하기 위한 공식이며, 데이터를 가르는 조건이 좋은 조건인지 좋지 않은 조건인지 평가하는 기준이 되는 공식이다

100명이 True, 0명이 False인 경우라면 (굉장히 좋은 경우)

$$I_G(p) = p_t(1 - p_t) + p_f(1 - p_f) = 1.0 * (1 - 1.0) + 0.0 * (1 - 0.0) = 0.0$$

50명이 True, 50명이 False인 경우라면 (굉장히 안 좋은 경우)

$$I_G(p) = p_t(1 - p_t) + p_f(1 - p_f) = 0.5 * (1 - 0.5) + 0.5 * (1 - 0.5) = 0.5$$

**Gini Impurity는 0에 가까울수록 좋은 조건이며, 0.5에 가까울수록 좋지 않은 조건이라 볼 수 있다.**

# Decision Tree의 동작 원리

Feature마다의 Gini Impurity를 계산 후  
이 Gini Impurity가 좋을 수록(=낮을 수록) 더 좋은 Feature(내지는 조건)이라고 가정한다

	짱절미	셀스타그램	우산	like
name				
Kang	True	False	False	True
Kim	False	False	False	False
Choi	False	True	False	True
Park	False	False	False	True
Yoon	False	False	False	False

보유하고 있는 데이터

# Decision Tree의 동작 원리

Feature마다의 Gini Impurity를 계산 후  
이 Gini Impurity가 좋을 수록(=낮을 수록) 더 좋은 Feature(내지는 조건)이라고 가정한다

Features

짱절미	셀스타그램	우산
-----	-------	----

Label

like

name

Kang	True	False	False	True
Kim	False	False	False	False
Choi	False	True	False	True
Park	False	False	False	True
Yoon	False	False	False	False

보유하고 있는 데이터



# Decision Tree의 동작 원리

Feature마다의 Gini Impurity를 계산 후  
이 Gini Impurity가 좋을 수록(=낮을 수록) 더 좋은 Feature(내지는 조건)이라고 가정한다

Features

like

Label

name	짱절미	셀스타그램	우산	like
Kang	True	False	False	True
Kim	False	False	False	False
Choi	False	True	False	True
Park	False	False	False	True
Yoon	False	False	False	False

짱절미 == True

셀스타그램 == True

우산 == True

보유하고 있는 데이터

Feature를 활용해  
가능한 모든 조건(condition)을 만든다

# Decision Tree의 동작 원리

Feature마다의 Gini Impurity를 계산 후  
이 Gini Impurity가 좋을 수록(=낮을 수록) 더 좋은 Feature(내지는 조건)이라고 가정한다

Features

Label

like

name	짱절미	셀스타그램	우산	
Kang	True	False	False	True
Kim	False	False	False	False
Choi	False	True	False	True
Park	False	False	False	True
Yoon	False	False	False	False

보유하고 있는 데이터

짱절미 == True

셀스타그램 == True

우산 == True

Feature를 활용해  
가능한 모든 조건(condition)을 만든다

```
condition = make_condition(binary_condition, "짱절미", True)
gini_impurity = evaluate_average_gini_impurity(data, condition)

("짱절미 == True", f"{gini_impurity:.6f}")

('짱절미 == True', '0.400000')
```

```
condition = make_condition(binary_condition, "셀스타그램", True)
gini_impurity = evaluate_average_gini_impurity(data, condition)

("셀스타그램 == True", f"{gini_impurity:.6f}")

('셀스타그램 == True', '0.400000')
```

```
condition = make_condition(binary_condition, "우산", True)
gini_impurity = evaluate_average_gini_impurity(data, condition)

("우산 == True", f"{gini_impurity:.6f}")

('우산 == True', '0.480000')
```

각 조건마다의 평균 gini impurity를 구한다.  
(해당 조건이 참일때, 거짓일때 gini impurity의 평균)

# Decision Tree의 동작 원리

Feature마다의 Gini Impurity를 계산 후  
이 Gini Impurity가 좋을 수록(=낮을 수록) 더 좋은 Feature(내지는 조건)이라고 가정한다

Features

짱절미 셀스타그램 우산

like

Label

name				
Kang	True	False	False	True
Kim	False	False	False	False
Choi	False	True	False	True
Park	False	False	False	True
Yoon	False	False	False	False

보유하고 있는 데이터

짱절미 == True

셀스타그램 == True

우산 == True

Feature를 활용해  
가능한 모든 조건(condition)을 만든다

```
condition = make_condition(binary_condition, "짱절미", True)
gini_impurity = evaluate_average_gini_impurity(data, condition)

condition = make_condition(binary_condition, "짱절미", True)
gini_impurity = evaluate_average_gini_impurity(data, condition)

condition = make_condition(binary_condition, "셀스타그램", True)
gini_impurity = evaluate_average_gini_impurity(data, condition)

condition = make_condition(binary_condition, "우산", True)
gini_impurity = evaluate_average_gini_impurity(data, condition)
```

각 조건마다의 평균 gini impurity를 구한다.  
(해당 조건이 참일때, 거짓일때 gini impurity의 평균)

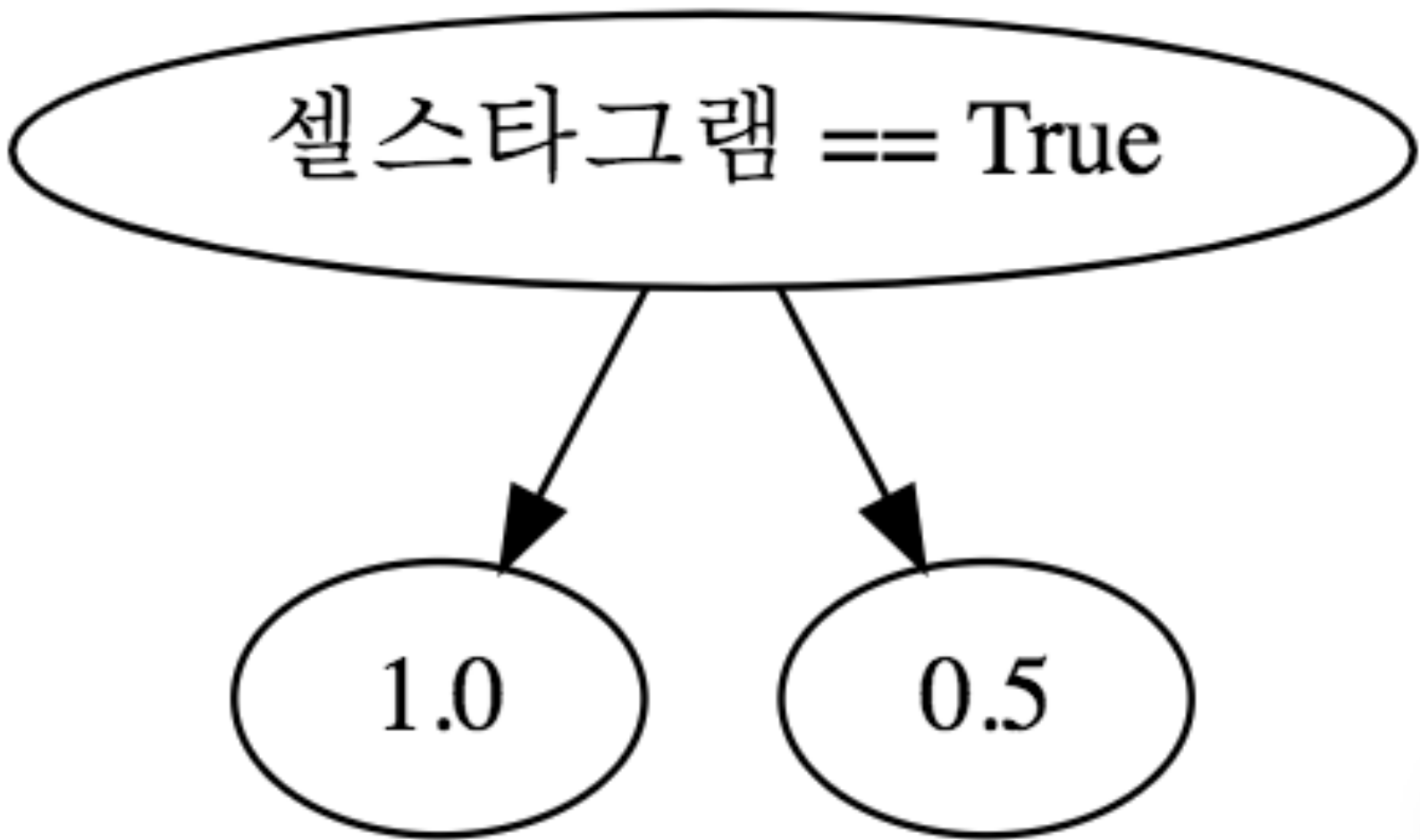
# Decision Tree의 동작 원리

Feature마다의 Gini Impurity를 계산 후  
이 Gini Impurity가 좋을 수록(=낮을 수록) 더 좋은 Feature(내지는 조건)이라고 가정한다

	짱절미	셀스타그램	우산	like	Label
name					
Kang	True	False	False	True	
Kim	False	False	False	False	
Choi	False	True	False	True	
Park	False	False	False	True	
Yoon	False	False	False	False	

보유하고 있는 데이터

Gini Impurity가 가장 낮은(==0.4)  
“셀스타그램 == True” 조건으로 나눈다



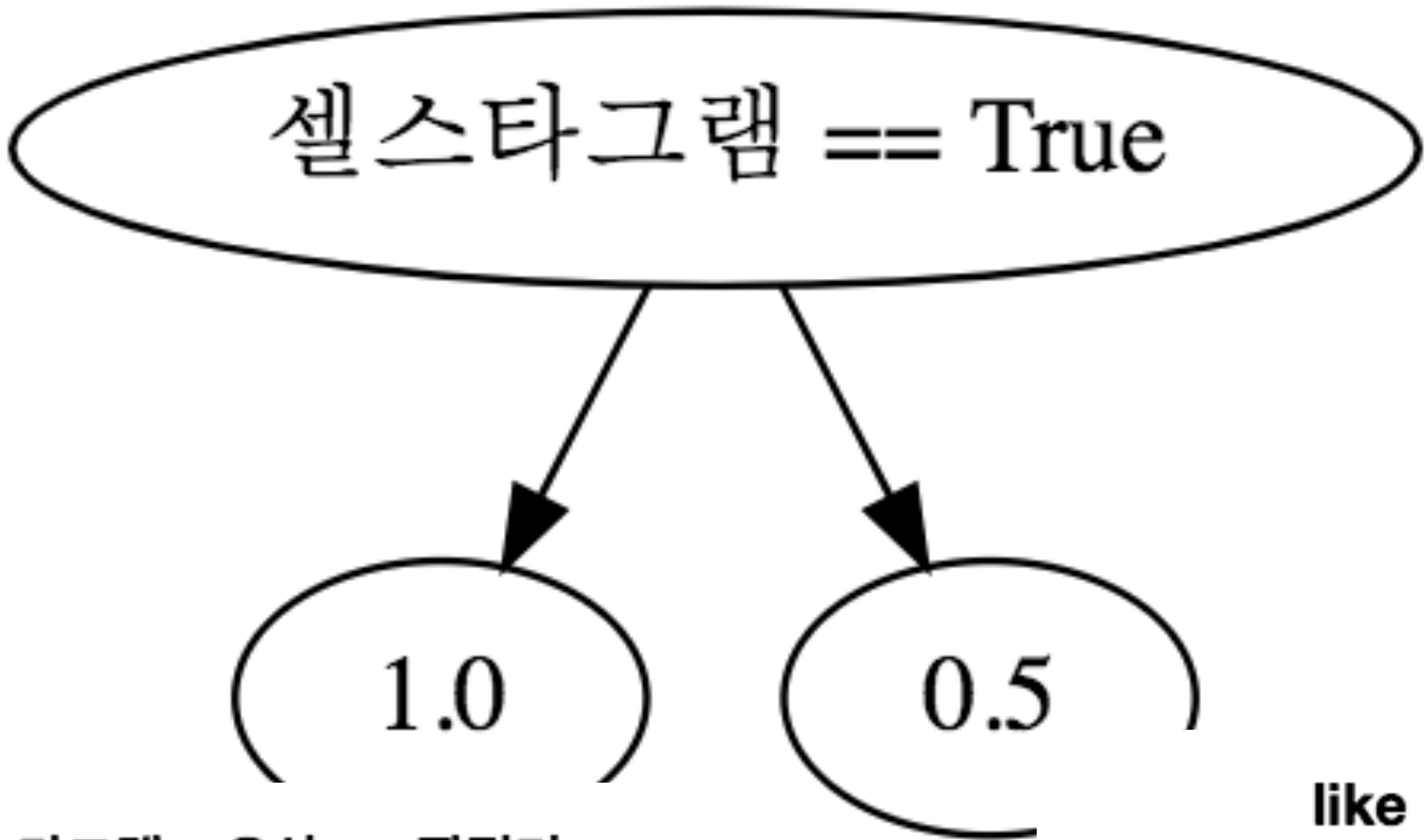


# Decision Tree의 동작 원리

Feature마다의 Gini Impurity를 계산 후  
이 Gini Impurity가 좋을 수록(=낮을 수록) 더 좋은 Feature(내지는 조건)이라고 가정한다

Gini Impurity가 가장 낮은(==0.4)  
“셀스타그램 == True” 조건으로 나눈다

	짱절미	셀스타그램	우산	like
Label				
name				
Kang	True	False	False	True
Kim	False	False	False	False
Choi	False	True	False	True
Park	False	False	False	True
Yoon	False	False	False	False



보유하고 있는 데이터

name	like	셀스타그램	우산	짱절미
Choi	True	True	False	False

name	like	셀스타그램	우산	짱절미
Kang	True	False	False	True
Kim	False	False	False	False
Park	True	False	False	False
Yoon	False	False	False	False

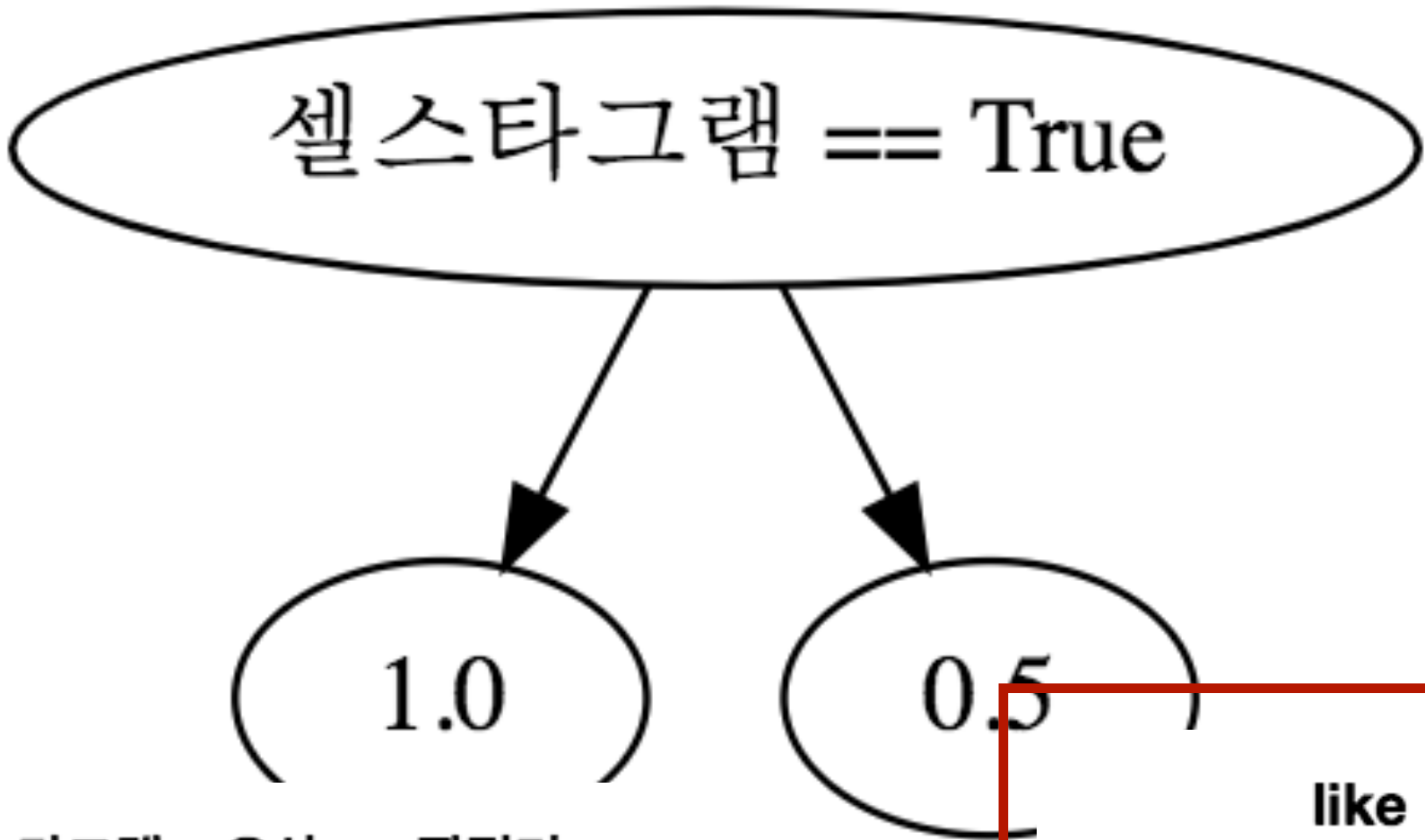
# Decision Tree의 동작 원리

Feature마다의 Gini Impurity를 계산 후  
이 Gini Impurity가 좋을 수록(=낮을 수록) 더 좋은 Feature(내지는 조건)이라고 가정한다

	짱절미	셀스타그램	우산	like
name				
Kang	True	False	False	True
Kim	False	False	False	False
Choi	False	True	False	True
Park	False	False	False	True
Yoon	False	False	False	False

보유하고 있는 데이터

name	like	셀스타그램	우산	짱절미
Choi	True	True	False	False

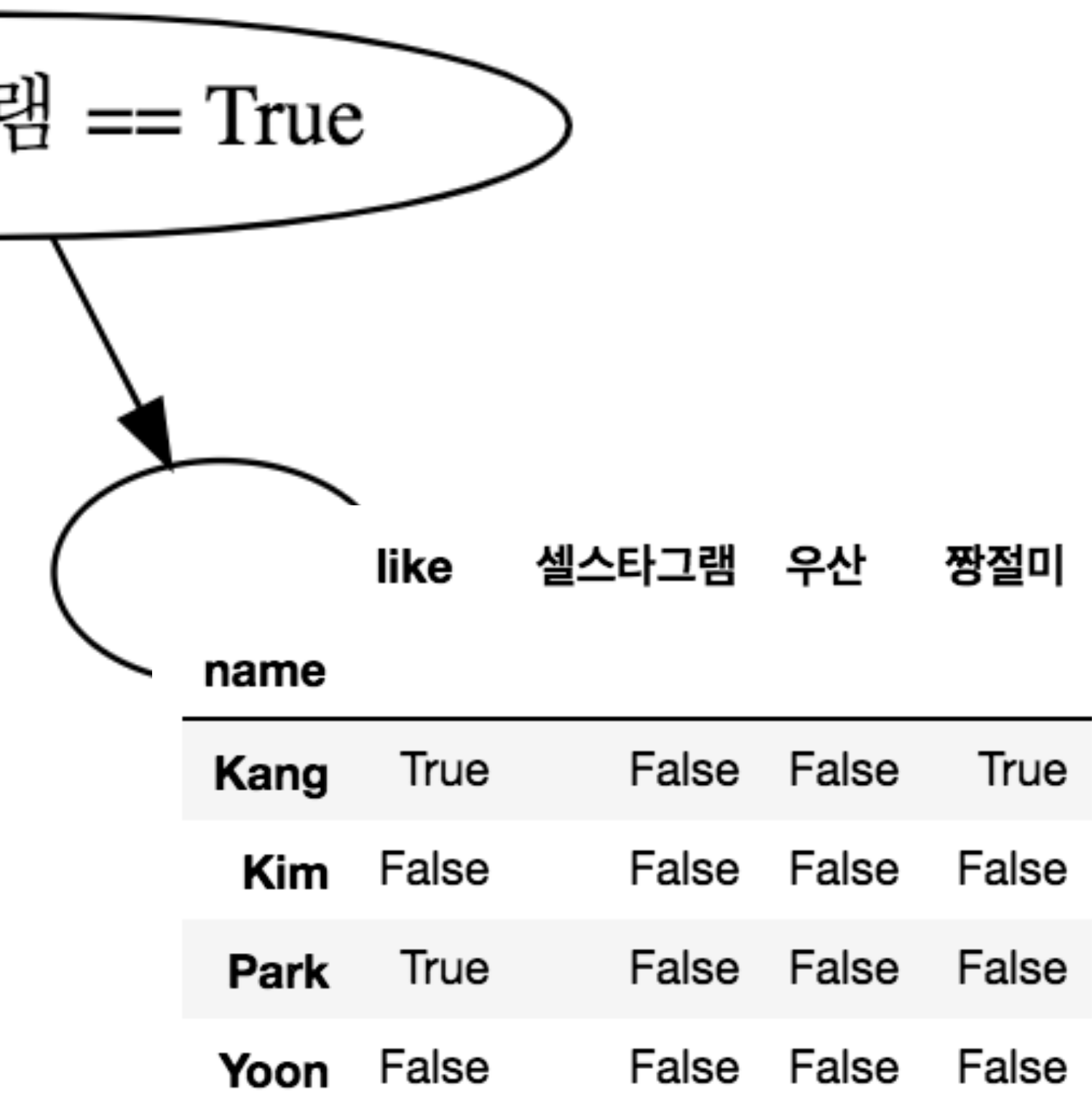


오른쪽 데이터는 한 번 더 나눌 수 있다

name	like	셀스타그램	우산	짱절미
Kang	True	False	False	True
Kim	False	False	False	False
Park	True	False	False	False
Yoon	False	False	False	False

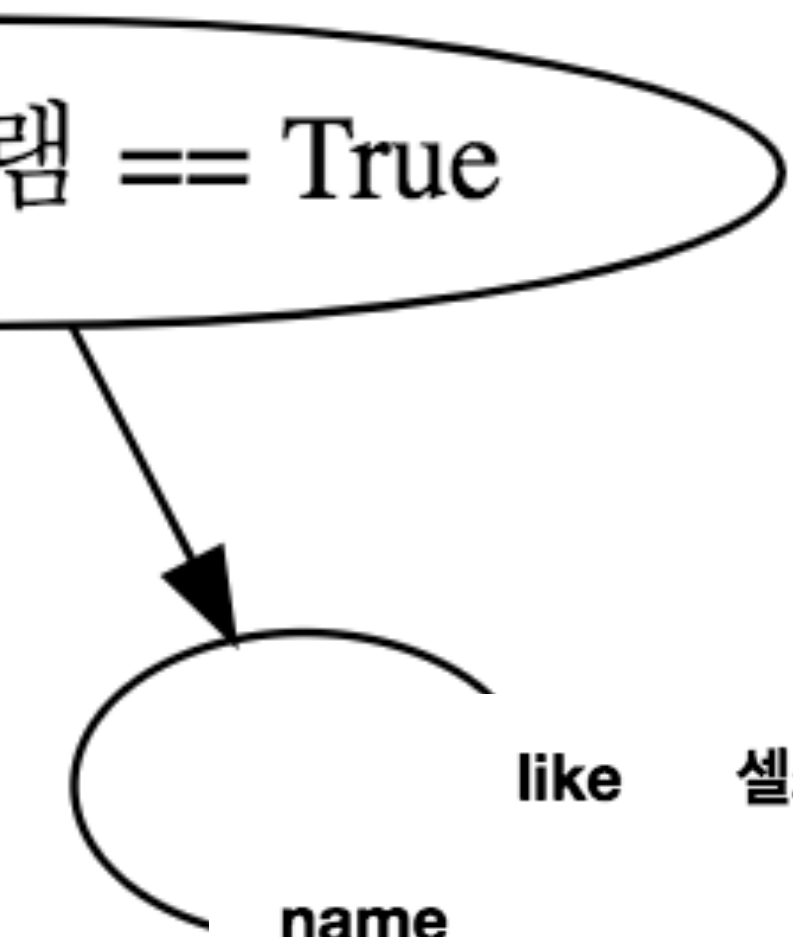
# Decision Tree의 동작 원리

Feature마다의 Gini Impurity를 계산 후  
이 Gini Impurity가 좋을 수록(=낮을 수록) 더 좋은 Feature(내지는 조건)이라고 가정한다



트리의 우측 데이터

# Decision Tree의 동작 원리



name	like	셀스타그램	우산	짱절미
Kang	True	False	False	True
Kim	False	False	False	False
Park	True	False	False	False
Yoon	False	False	False	False

트리의 우측 데이터

Feature마다의 Gini Impurity를 계산 후 이 Gini Impurity가 좋을 수록(=낮을 수록) 더 좋은 Feature(내지는 조건)이라고 가정한다

```
condition = make_condition(binary_condition, "짱절미",  
gini_impurity = evaluate_average_gini_impurity(right_  
"짱절미 == True", f"{gini_impurity:.6f}")  
( '짱절미 == True', '0.333333' )
```

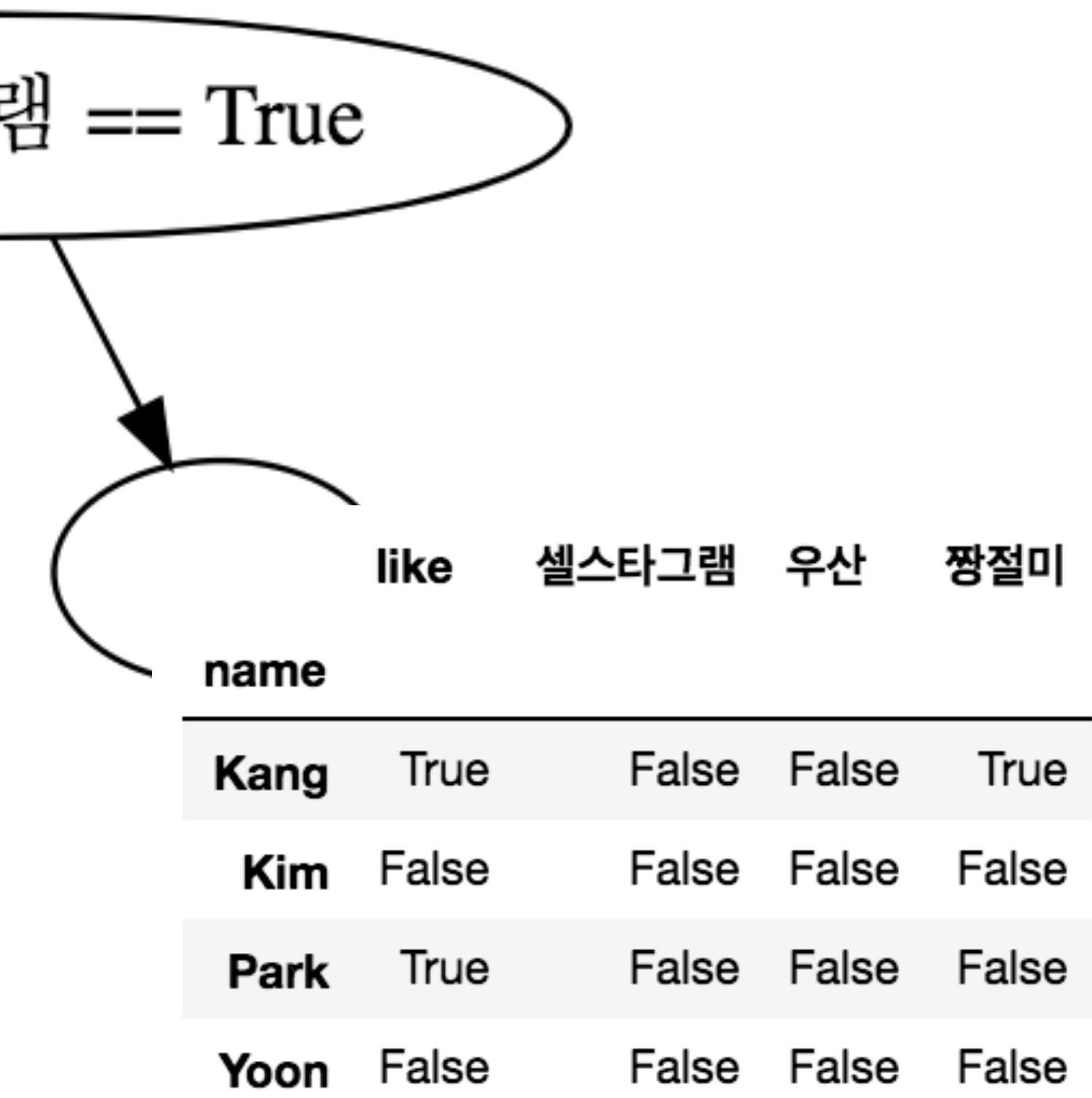
```
condition = make_condition(binary_condition, "우산", T  
gini_impurity = evaluate_average_gini_impurity(right_  
"우산 == True", f"{gini_impurity:.6f}")  
( '우산 == True', '0.500000' )
```

우측 데이터의 평균 Gini Impurity를 구한다



# Decision Tree의 동작 원리

Feature마다의 Gini Impurity를 계산 후 이 Gini Impurity가 좋을 수록(=낮을 수록) 더 좋은 Feature(내지는 조건)이라고 가정한다



트리의 우측 데이터

```
condition = make_condition(binary_condition, "짱절미",  
gini_impurity = evaluate_짱절미 == True 조건이 Gini Impurity가 높다  
"짱절미 == True", f"{gini_impurity:.6f}")  
( '짱절미 == True', '0.333333' )
```

```
condition = make_condition(binary_condition, "우산", T  
gini_impurity = evaluate_average_gini_impurity(right_  
"우산 == True", f"{gini_impurity:.6f}")  
( '우산 == True', '0.500000' )
```

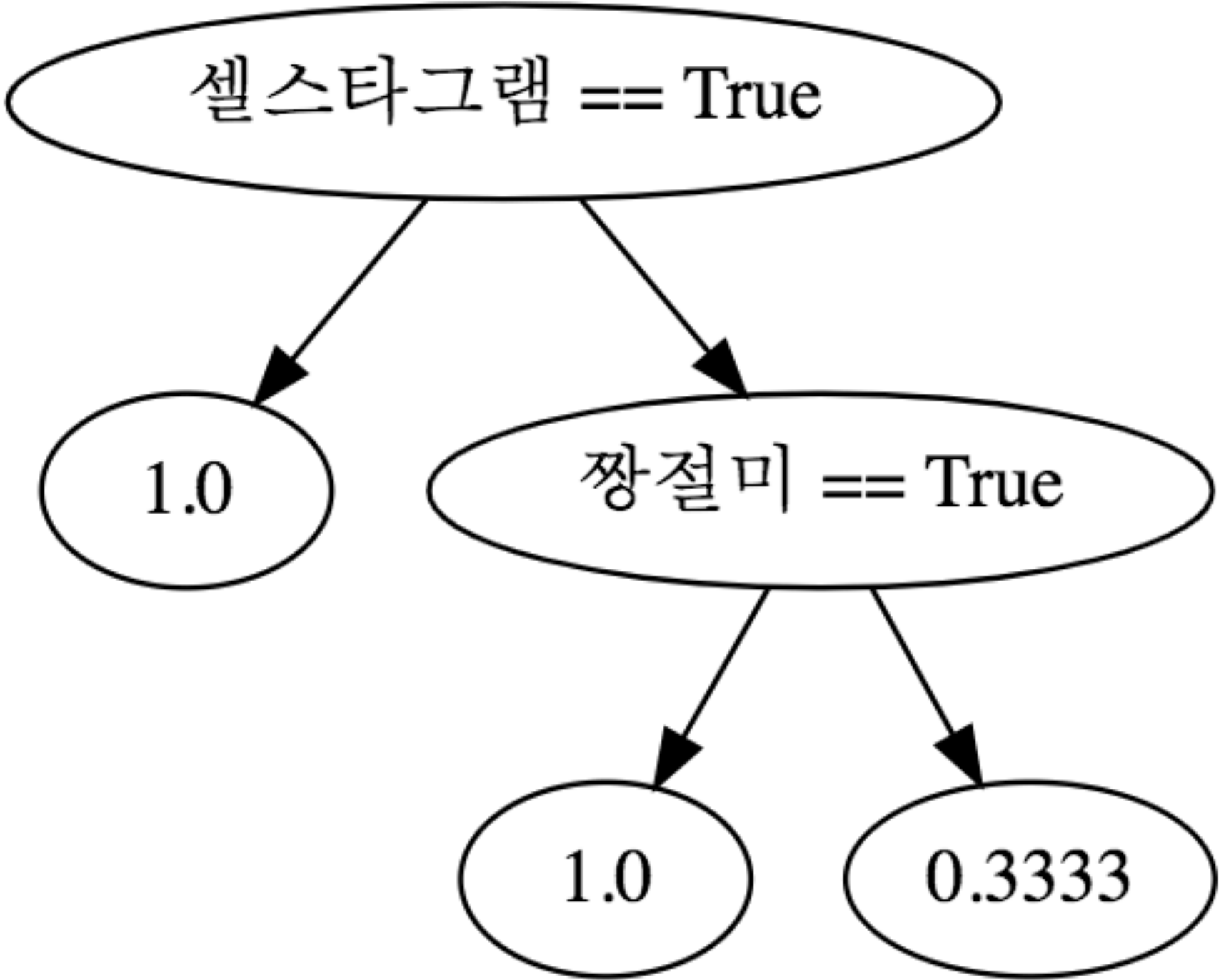
우측 데이터의 평균 Gini Impurity를 구한다

# Decision Tree의 동작 원리

Feature마다의 Gini Impurity를 계산 후 이 Gini Impurity가 좋을 수록(=낮을 수록) 더 좋은 Feature(내지는 조건)이라고 가정한다

```
condition = make_condition(binary_condition, "짱절미",  
gini_impurity = evaluate  
"짱절미 == True", f"{gini_impurity:.6f}"  
( '짱절미 == True', '0.333333' )  
  
condition = make_condition(binary_condition, "우산", T  
gini_impurity = evaluate_average_gini_impurity(right_  
"우산 == True", f"{gini_impurity:.6f}"  
( '우산 == True', '0.500000' )
```

우측 데이터의 평균 Gini Impurity를 구한다



조건을 추가하여 Tree를 구축한다

**Q & A**