



k-Nearest Neighbors (k-NN)

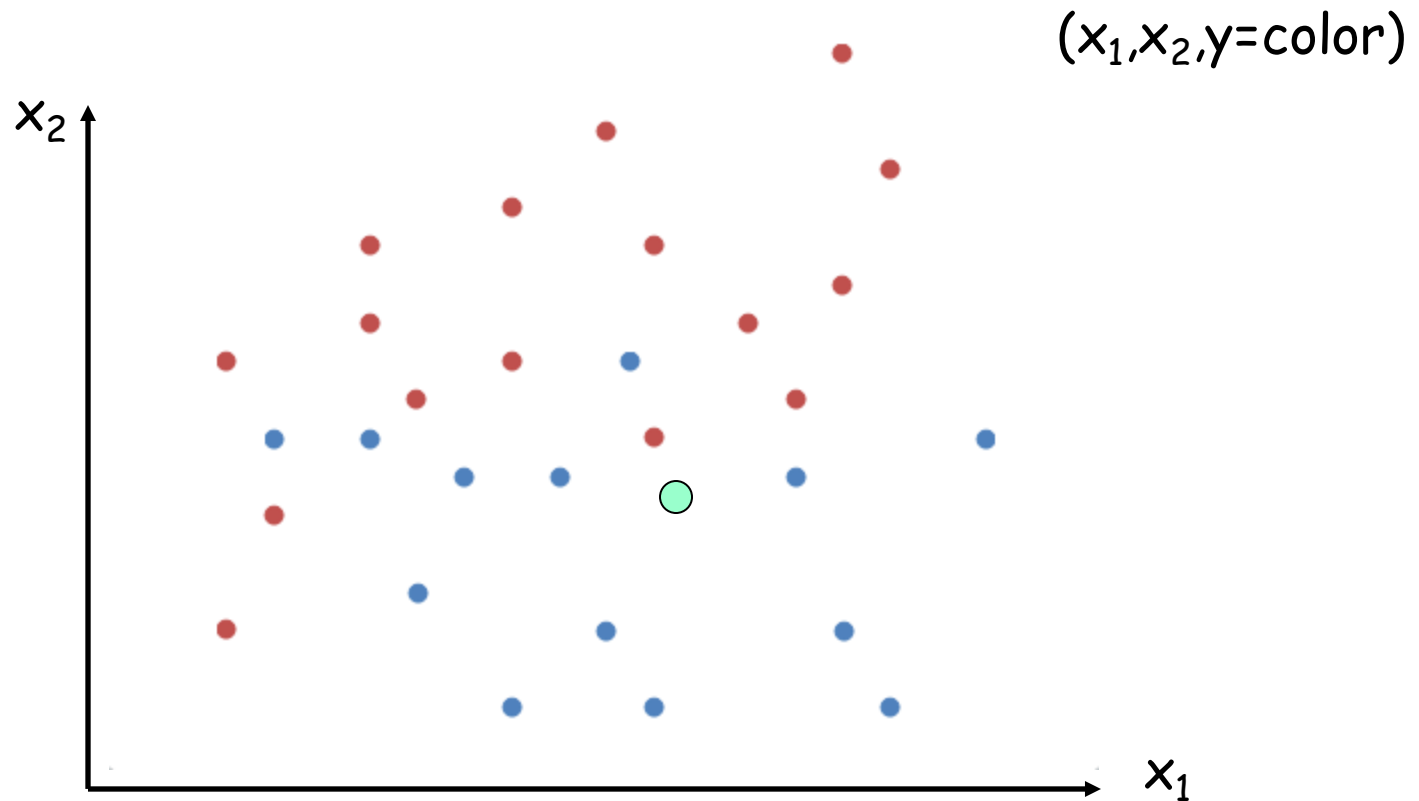
If your friends are RED, you must be RED.

Contents

- **Classification with k-NN**
- **Regression with k-NN**
- **Summary**

Classification

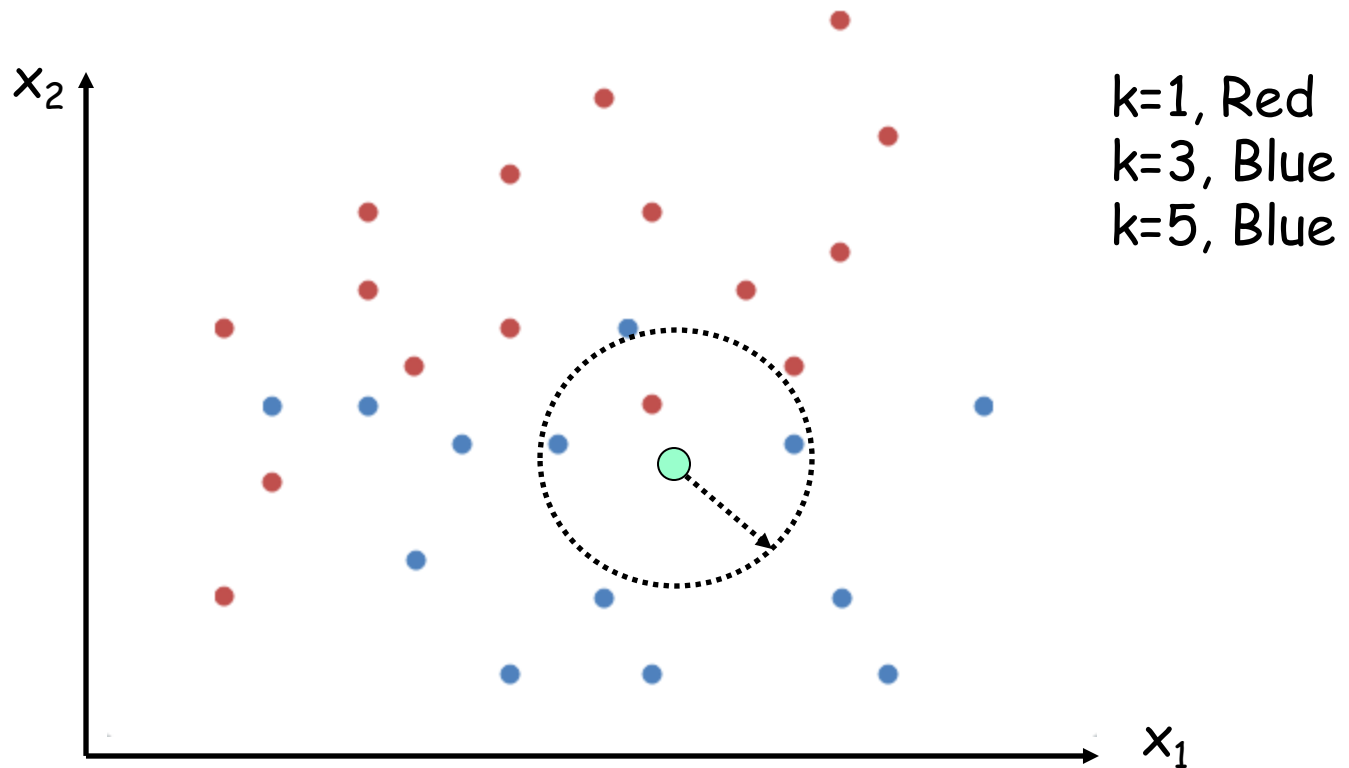
- How to Predict Class of Unknown Data?



Classification

■ K-Nearest Neighbors

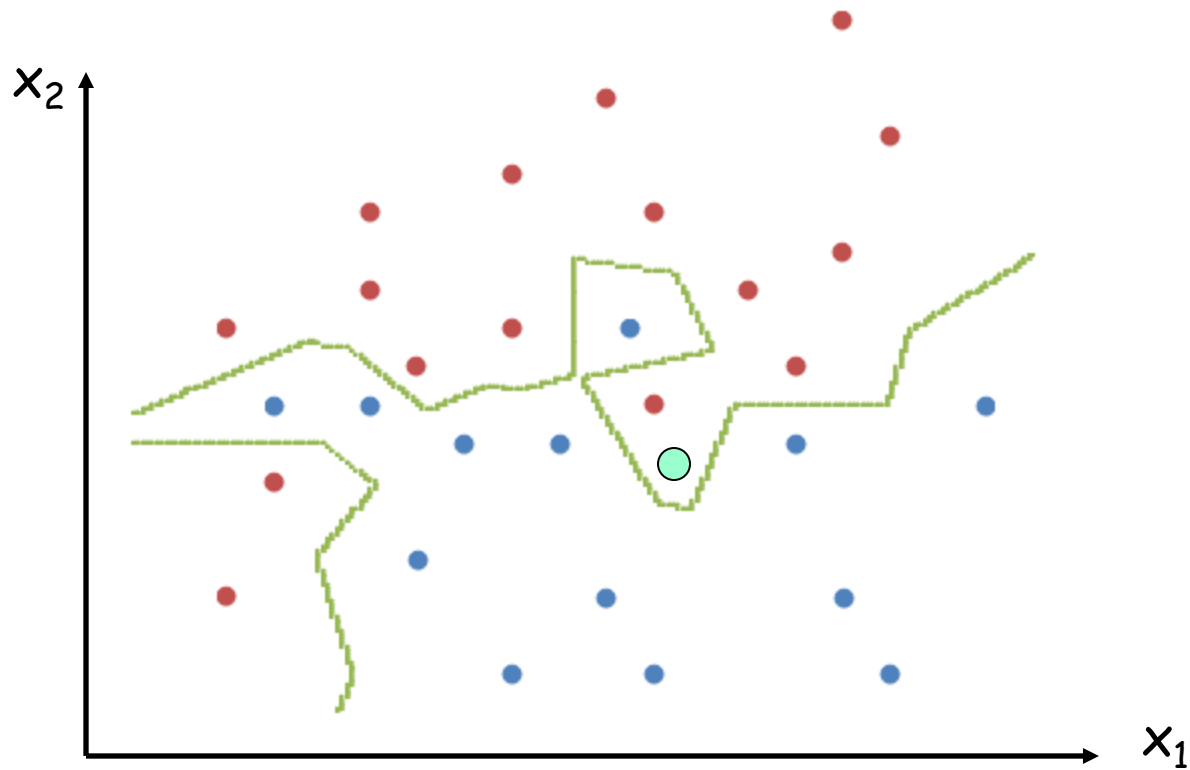
- Choose k nearest neighbors
- Determine the class based on the majority



Classification

- **K-Nearest Neighbors**

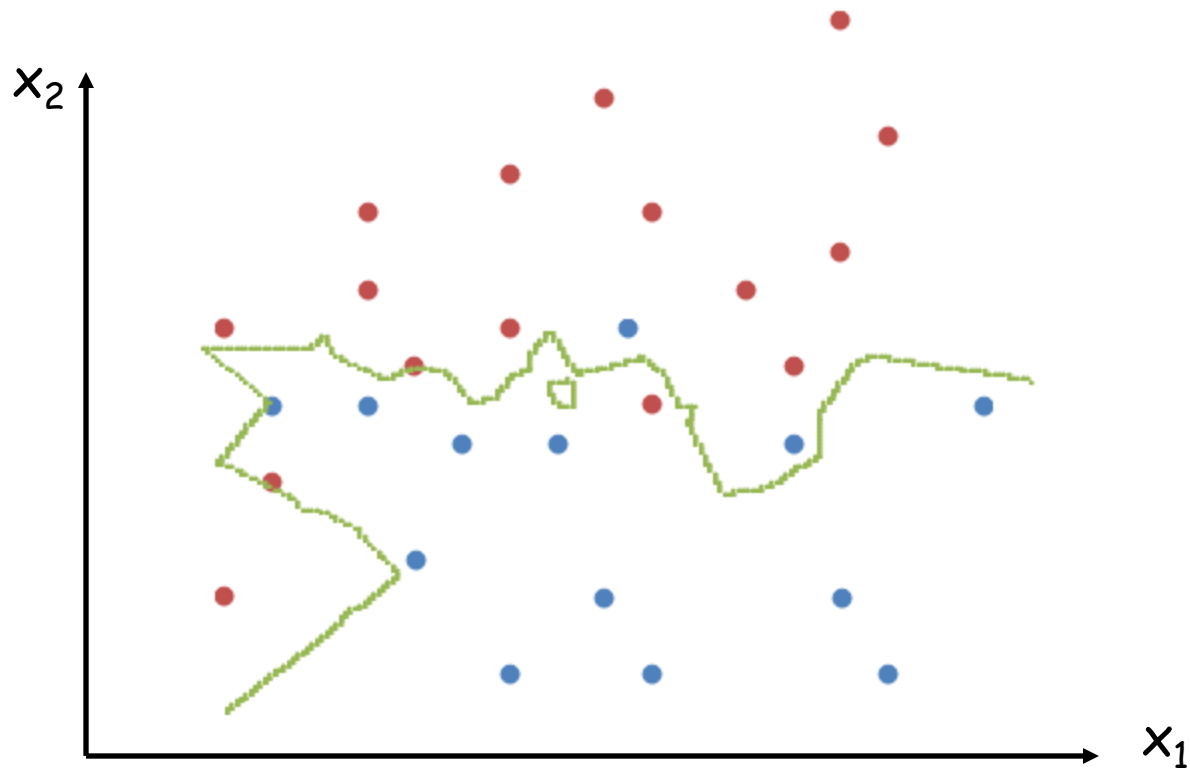
- $K = 1$



Classification

- **K-Nearest Neighbors**

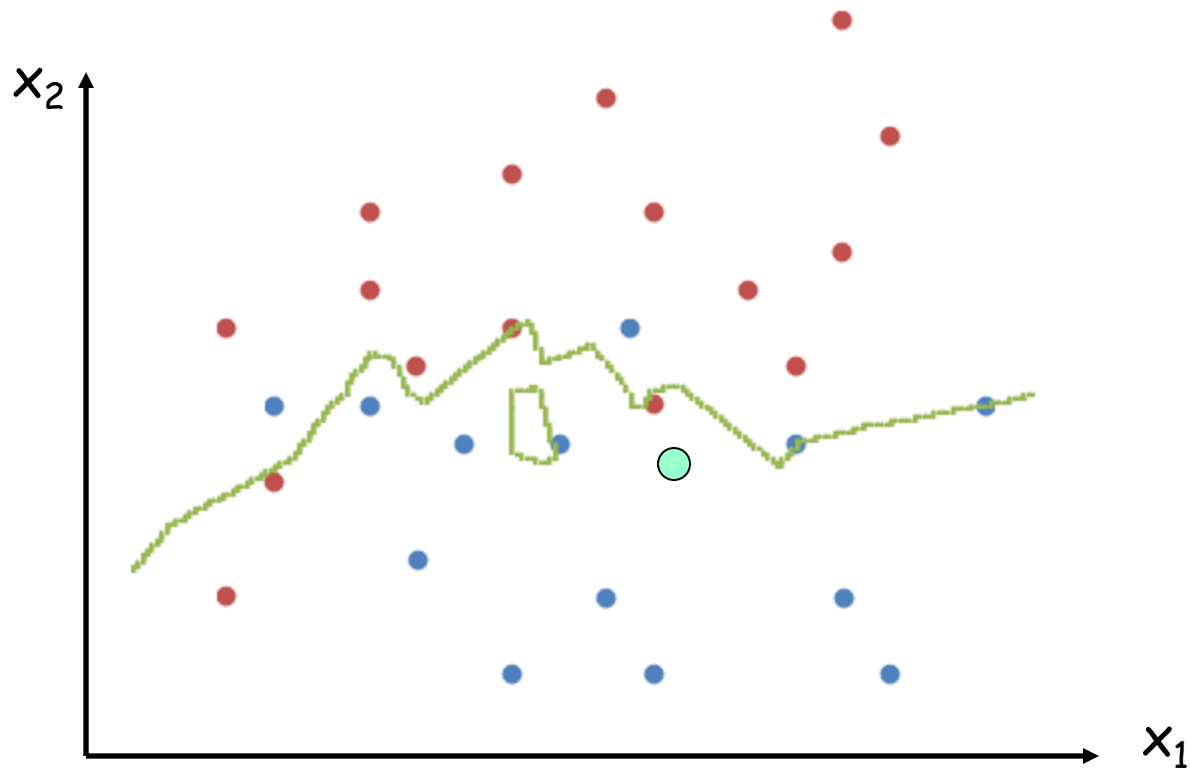
- $K = 3$



Classification

- **K-Nearest Neighbors**

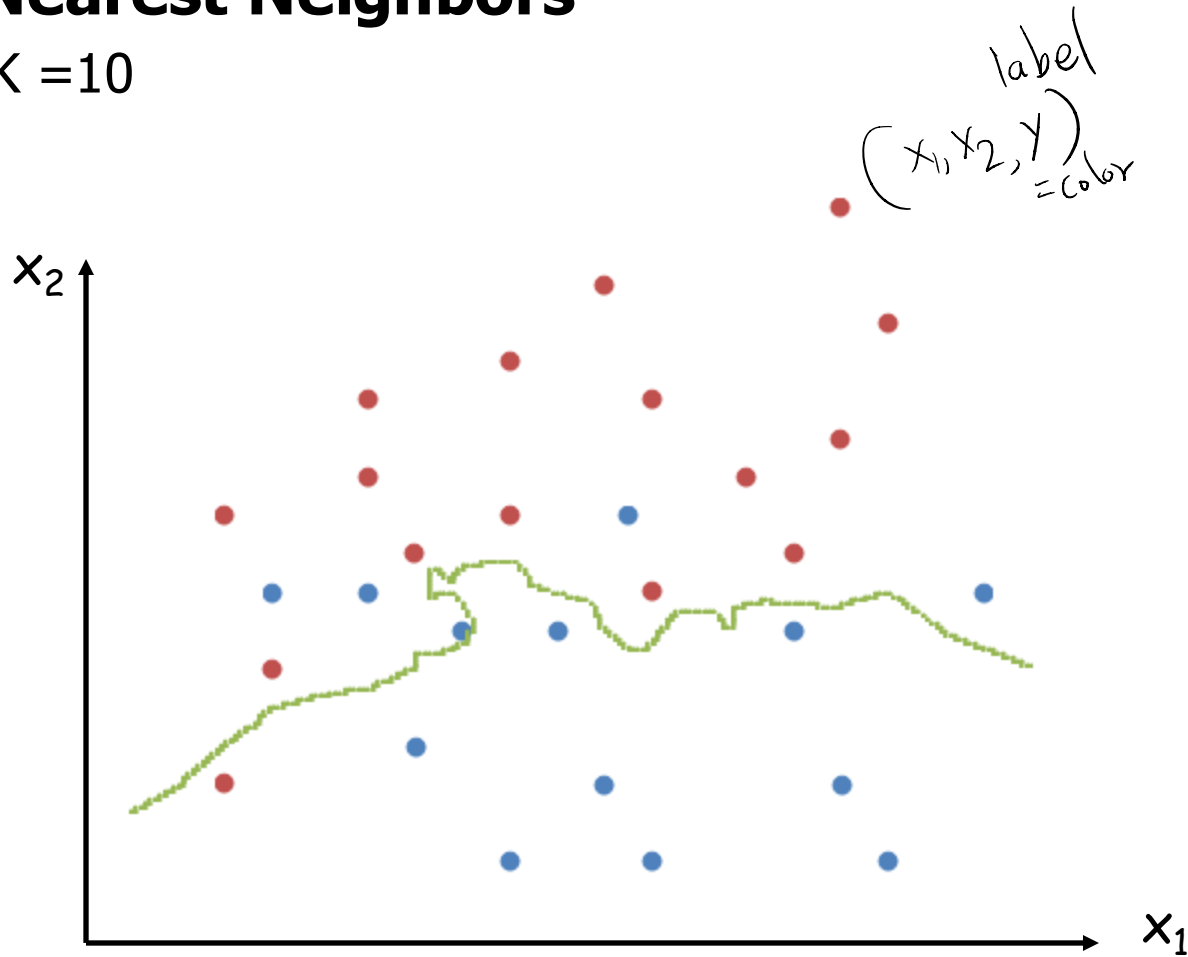
- $K = 5$



Classification

- **K-Nearest Neighbors**

- $K = 10$



Classification

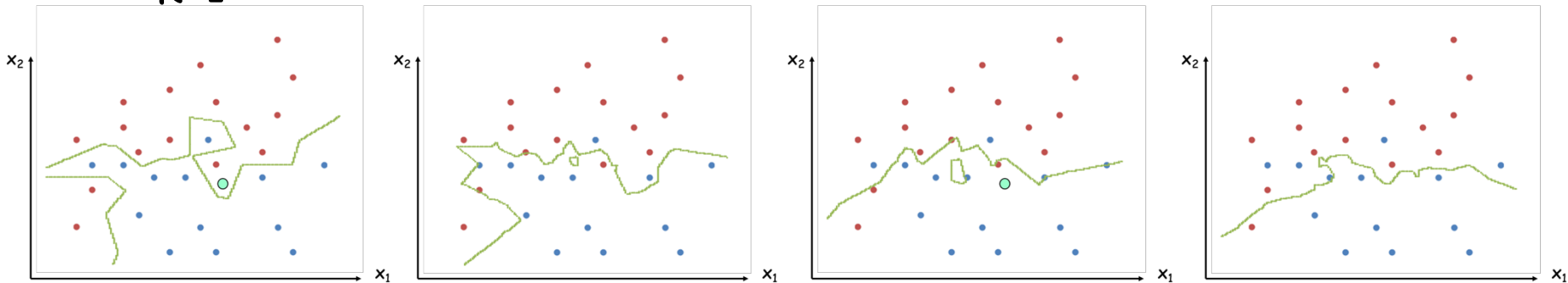
■ Then, which “k”?

K=1

K=3

K=5

K=10



Complex Boundary Simple

높은 (정확도)
Fidelity to data

High

Low

민감도
Sensitivity to noisy data

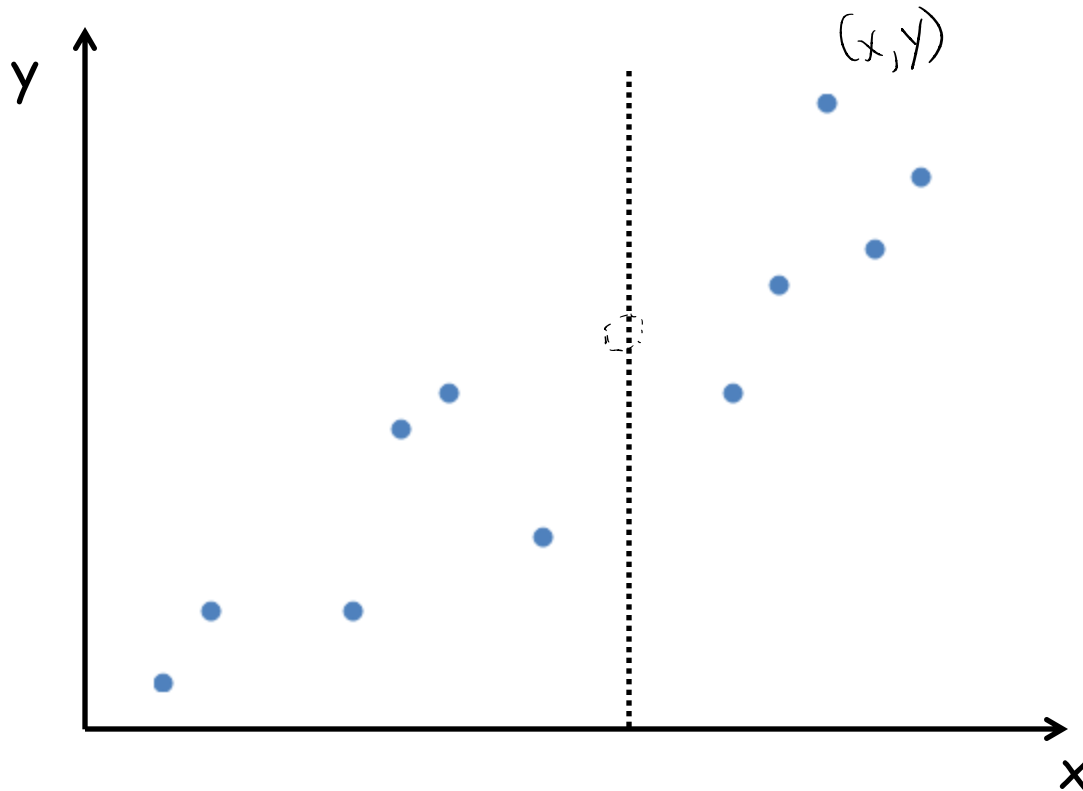
High

Low

Stable

Regression

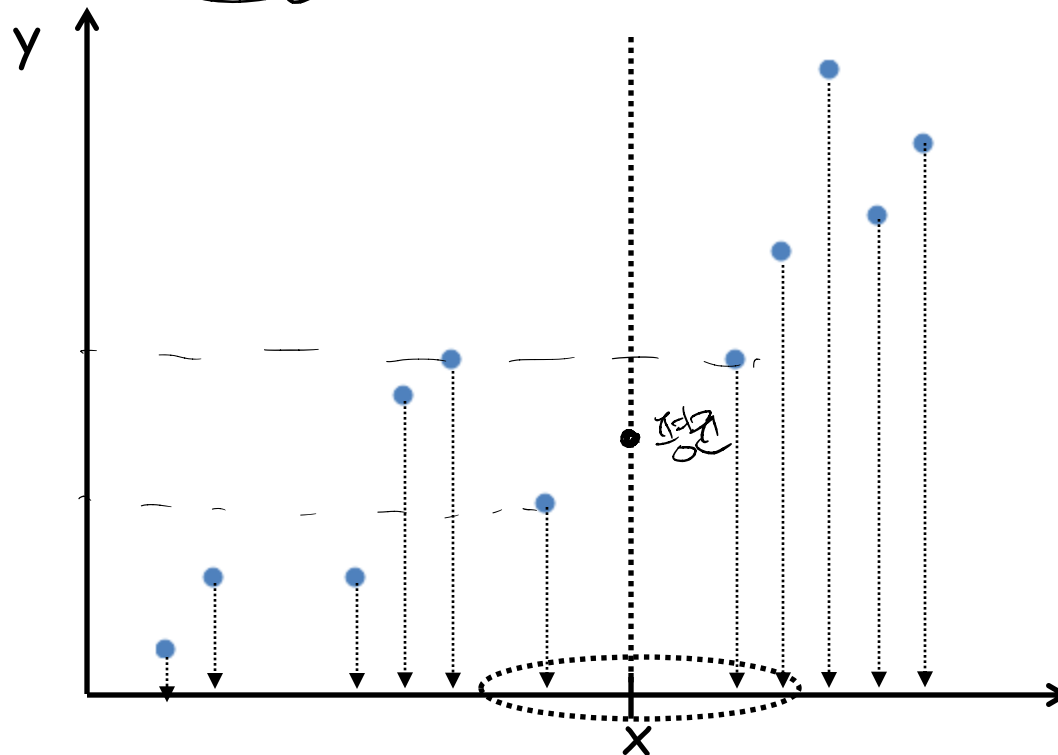
- How to predict “ y ” value of unknown data



Regression

■ K-Nearest Neighbors

- Choose k nearest neighbors in X axis
- Predict the average of their "y"

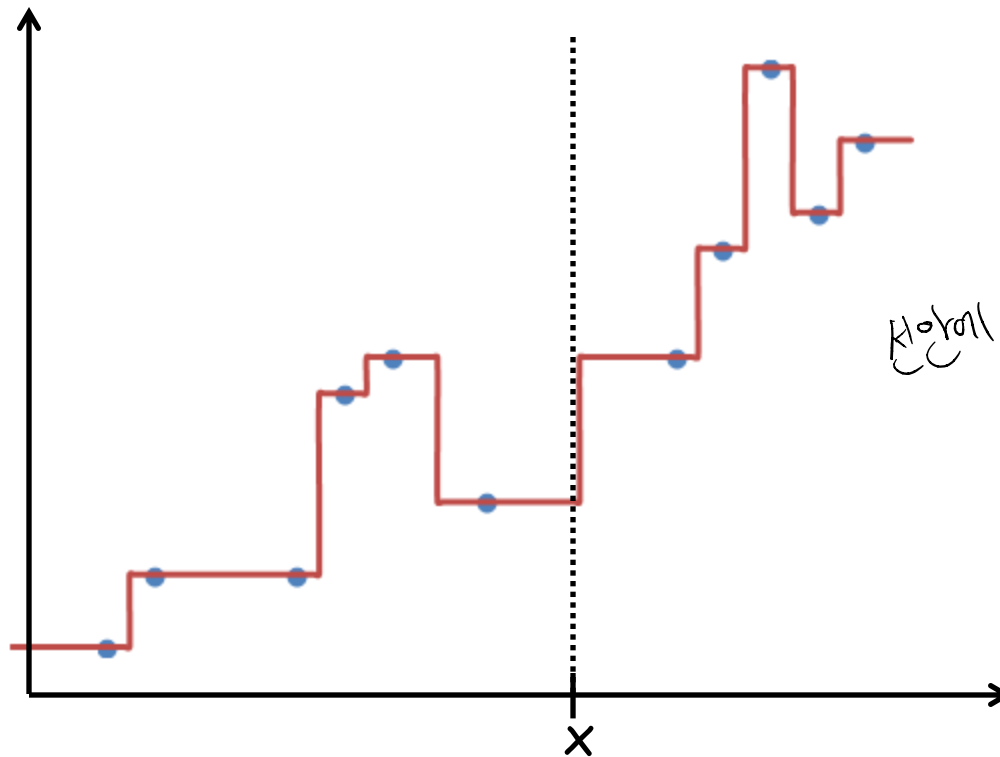


(x_1, y_1)
 (x_2, y_2)
 (x_3, y_3)
 \vdots
 (x_n, y_n)

Regression

- K-Nearest Neighbors

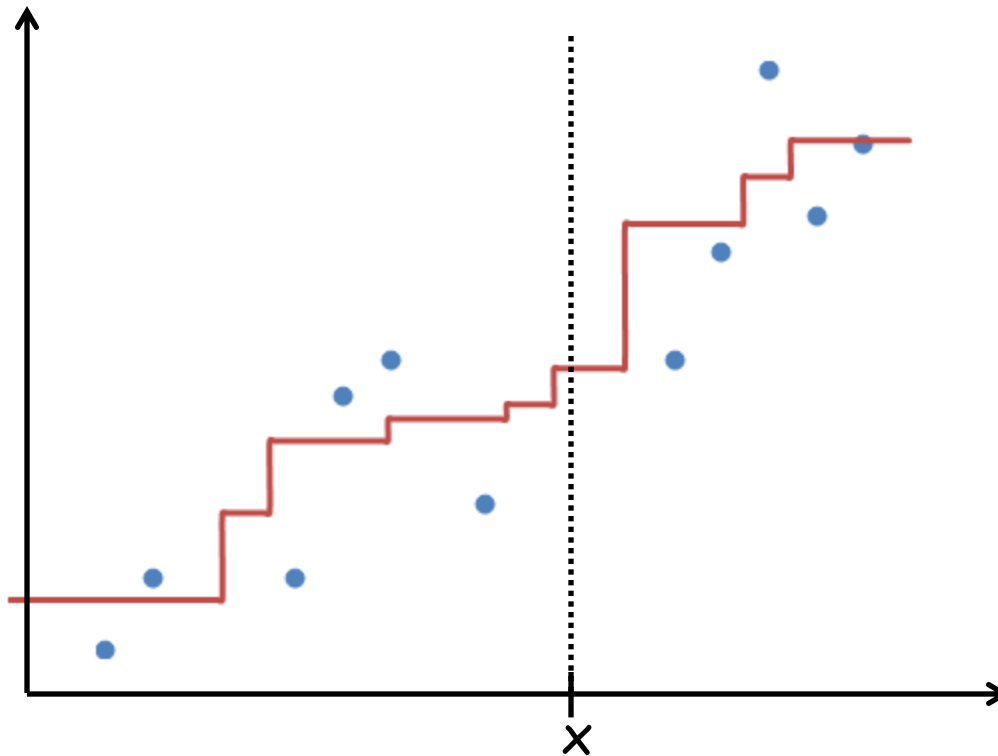
- K=1



선상에 점이 존재할 것

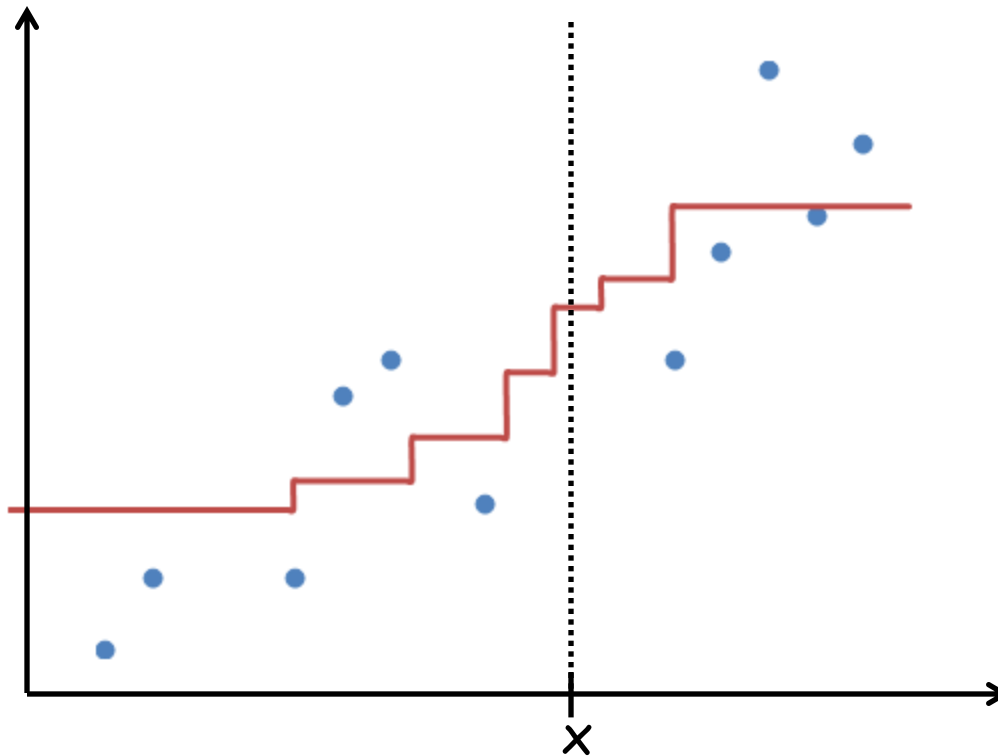
Regression

- **K-Nearest Neighbors**
 - $K=3$



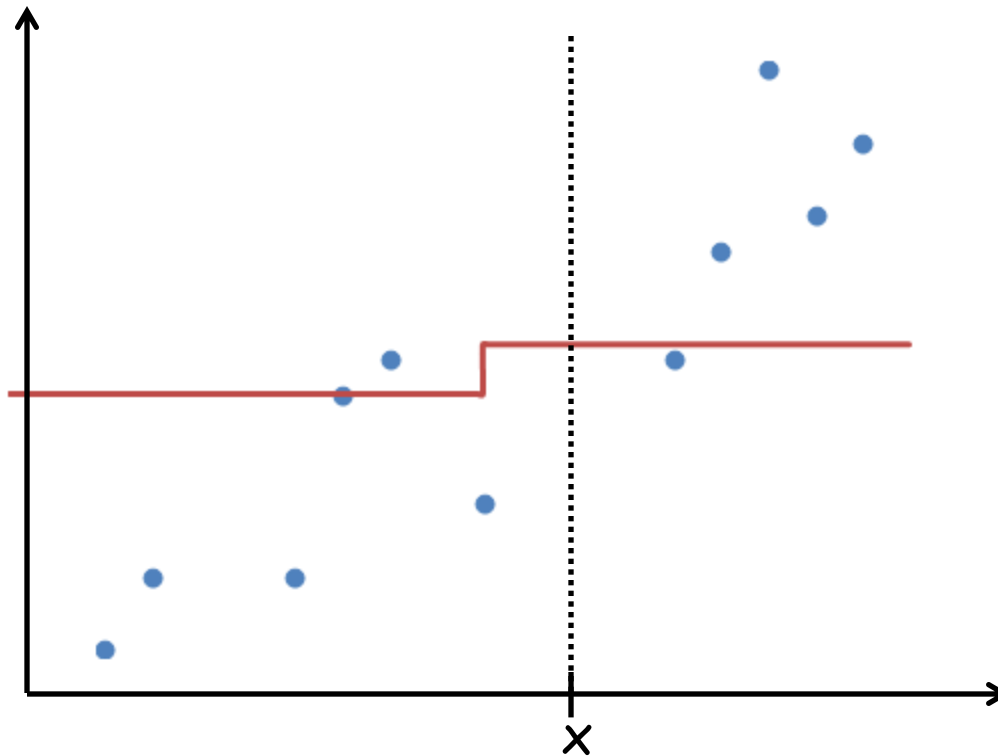
Regression

- **K-Nearest Neighbors**
 - $K=5$



Regression

- **K-Nearest Neighbors**
 - $K=10$

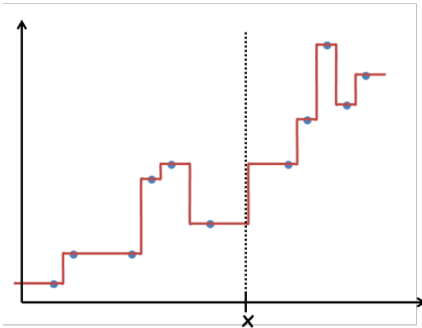


Regression ~~Classification~~

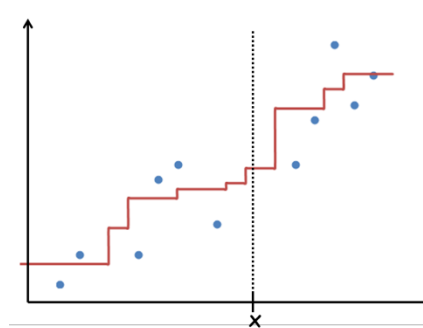
■ Then, which "k"?

값이 높을수록 성능이 좋아!

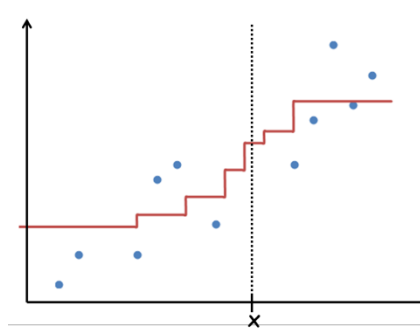
K=1



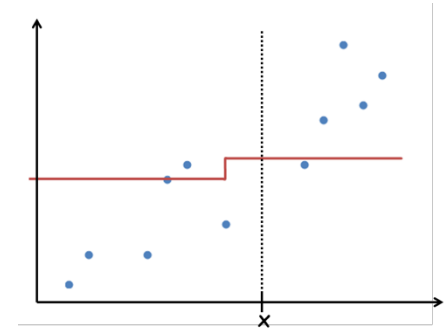
K=3



K=5



K=10



Model output

Complex

Simple

Fidelity to data

High

Low

Sensitivity to noisy data

High

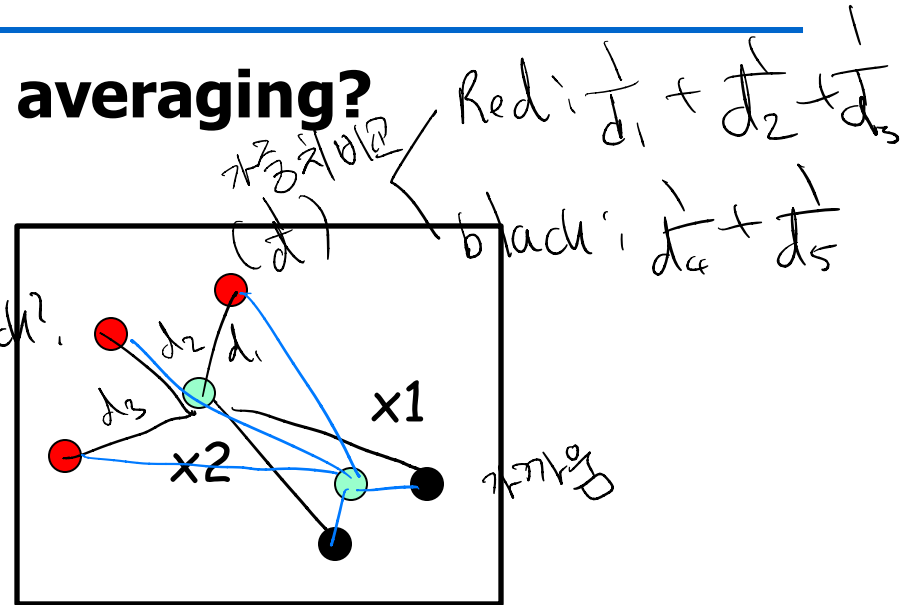
Low

Variation

Why "just" counting or averaging?

Classification (k=5)

- X1: Red or Black? Red? Black?
- X2: Red or Black? Red

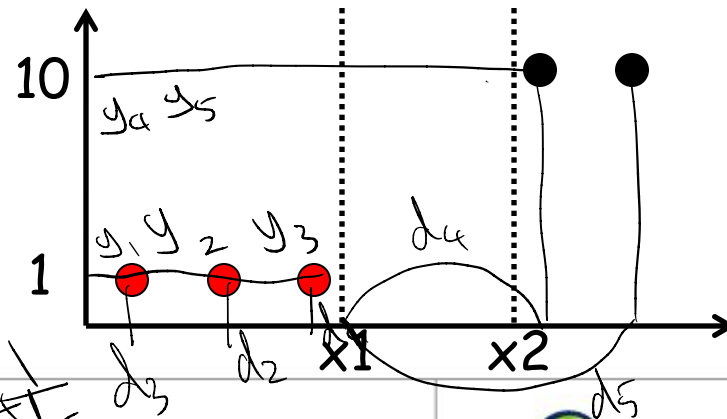


Regression (k=5)

- X1: close to 10 or 1
- X2: close to 10 or 1

$$x_1 = x_2 = \frac{y_1 + y_2 + y_3 + y_4 + y_5}{5}$$

$$= \frac{\frac{1}{d_1} + \frac{1}{d_2} + \dots + \frac{1}{d_5}}{\frac{1}{d_1} + \frac{1}{d_2} + \dots + \frac{1}{d_5}}$$



Variation

- **More weight to closer one**

- Different weight depending on the distance from \mathbf{x}'

- **Classification: Not just counting**

$$S(\mathbf{x}', R) = \sum_{\mathbf{x} \in N(\mathbf{x}', R)} w(\mathbf{x})$$

$$S(\mathbf{x}', B) = \sum_{\mathbf{x} \in N(\mathbf{x}', B)} w(\mathbf{x})$$

if $S(\mathbf{x}', R) > S(\mathbf{x}', B)$ then \mathbf{x}' is R

else \mathbf{x}' is B

$N(\mathbf{x}', R)$: the set of *Red* data among the nearest neighbors of \mathbf{x}'

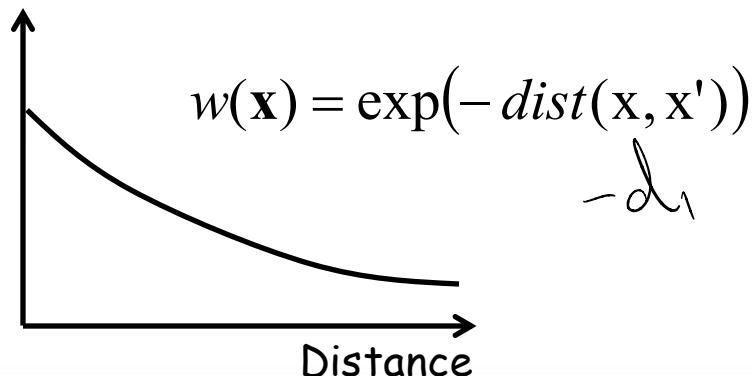
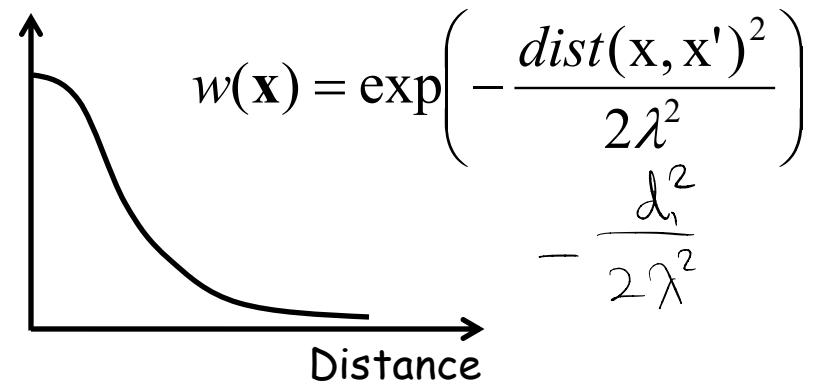
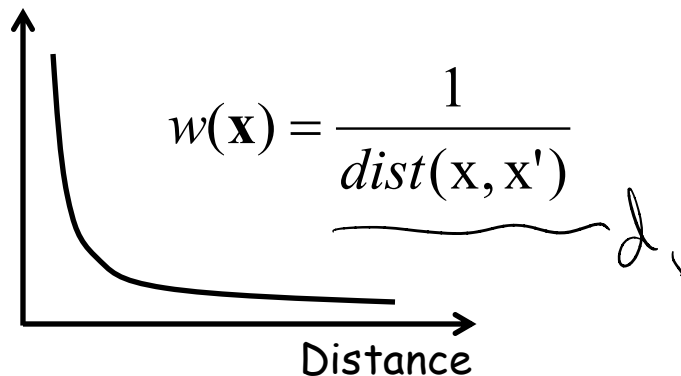
- **Regression: Not just averaging**

$$f(\mathbf{x}') = \frac{\sum_{\mathbf{x} \in N(\mathbf{x}')} w(\mathbf{x}) \cdot f(\mathbf{x})}{\sum_{\mathbf{x} \in N(\mathbf{x}')} w(\mathbf{x})}$$

$N(\mathbf{x}')$: the nearest neighbors of \mathbf{x}'

Variation

- How to determine weight considering distance

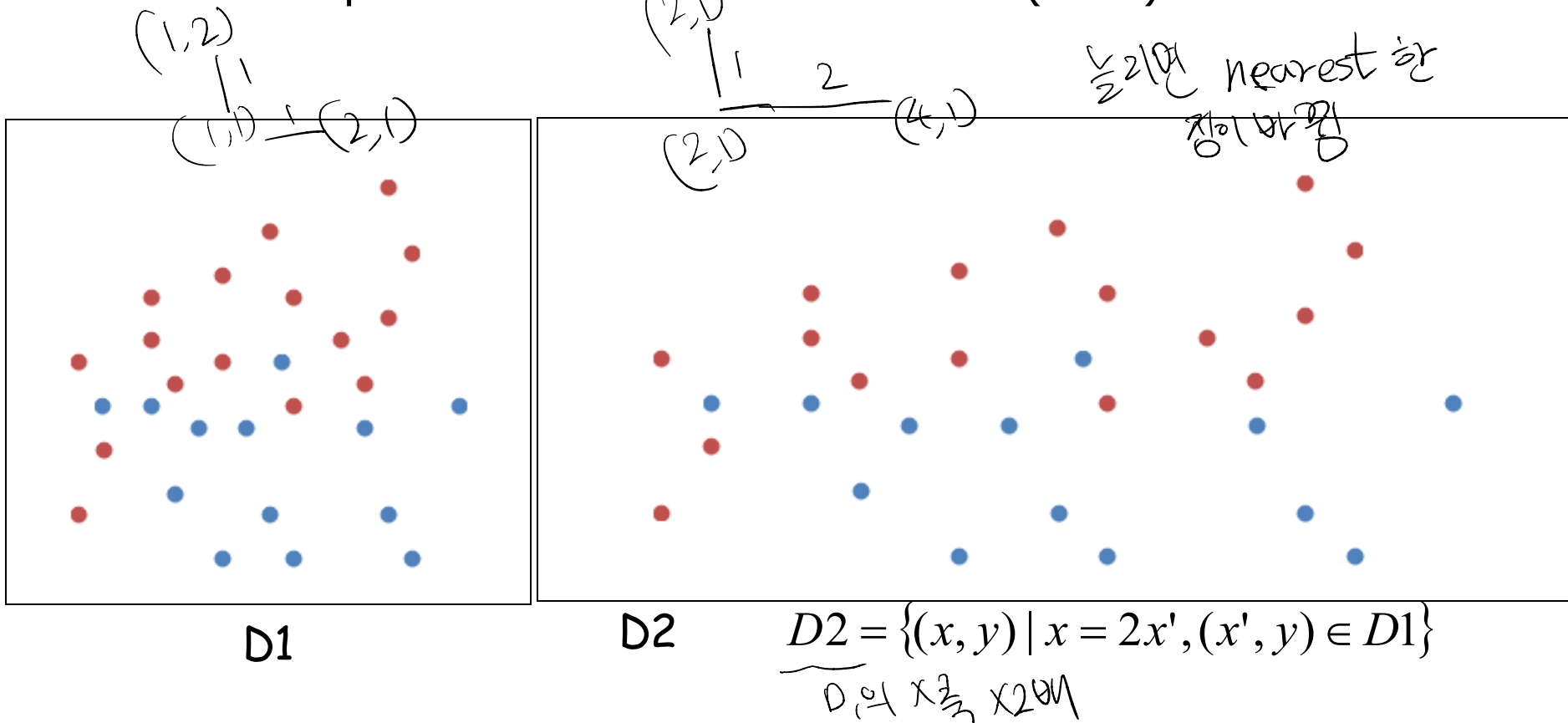


Do we need to choose k NNs?
-Yes, if you want
-Not necessarily

Distance Measure

Two data sets

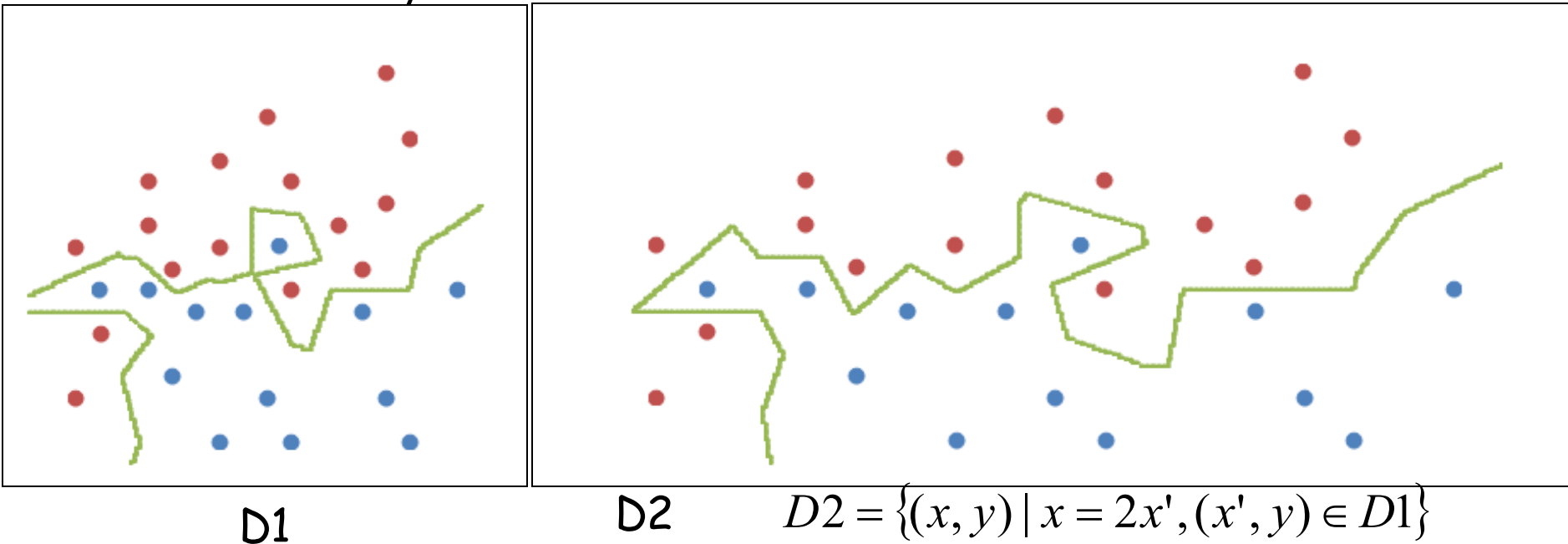
- Compare the boundaries in D1 and D2 (1-NN)



Distance Measure

■ Two data sets

- Even though the data are linearly scaled, the boundary changes!! *변화한다*
- You may need to choose other distance measures



Summary

- **Which k is better?**

- Small k : higher variance (less stable) -> possibly overfitted
- Large k : higher bias (less precise) -> possibly underfitted

- **Proper choice of k**

- Depending on the data
- Use Cross-validation : *의 표준편차를 선택가능*

Summary

- **Advantage** of k -NN
 - No training (Only inference step)
 - Complexity of target functions do not matter
 - No loss of information
 - **Disadvantage**
 - Have to keep all data -> Memory space
 - Sensitive to noise
 - If training data is imbalanced, major class may dominate
 - Need to calculate the distance from all training data -> Time
 - Especially in high dimensional space, expensive
- a lot of Time, space

Summary

■ Reducing Computational Cost

- Finding k nearest neighbors is expensive: $O(nd)$
- Space partitioning
 - quad-tree, locality sensitive hashing, etc.
- Preprocessing
 - Reduce dimensions: Remove less important features, Vector quantization
 - Reduce size of data: Sampling, Clustering

K-NN in Scikit-Learn

```
import numpy as np
from sklearn.datasets import make_blobs
from sklearn.neighbors import KNeighborsClassifier
import matplotlib.pyplot as plt

x,y = make_blobs(50, n_features=2, centers=2, random_state=0)
knn = KNeighborsClassifier(n_neighbors = 3, \
                           weights = 'uniform', \
                           metric = 'euclidean')

knn.fit(x, y)

test_data = np.array([[1.4, 0.2], [1.4,0.5],[0.9, 4.0], \
                       [-0.1, 3.0], [2.5, 0.1]])

prediction = knn.predict(test_data)
print(prediction)
```