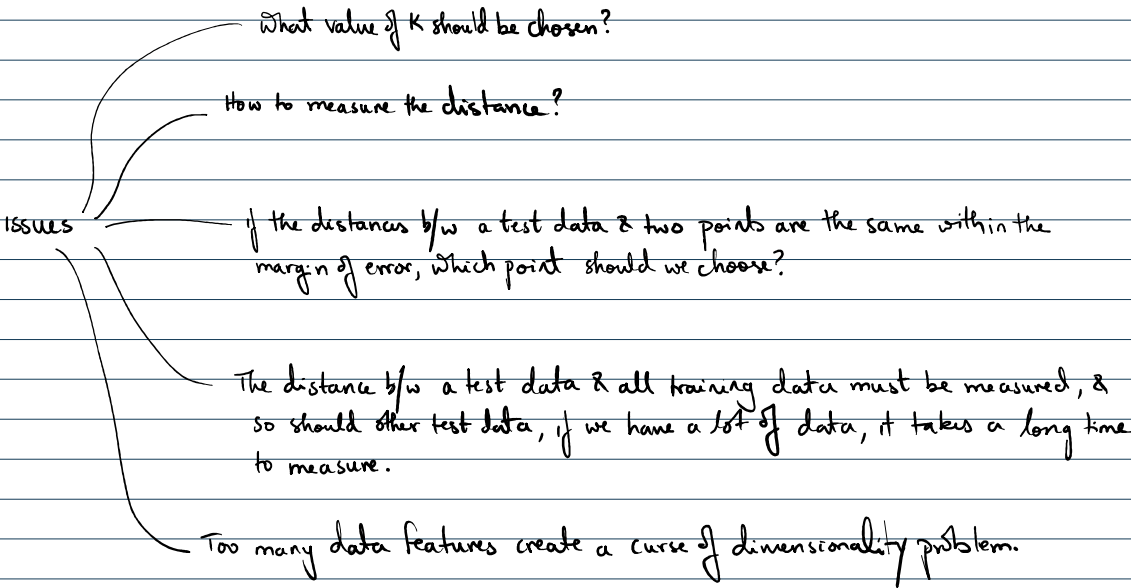


Classification Algorithm

- 1) Measure distances b/w the test data & all training data
- 2) Find K neighborhood closest to the test data
- 3) Find all classes (y) of the neighbors
- 4) Find the most class y & determine this as the class y of the test data.
(Majority Voting)



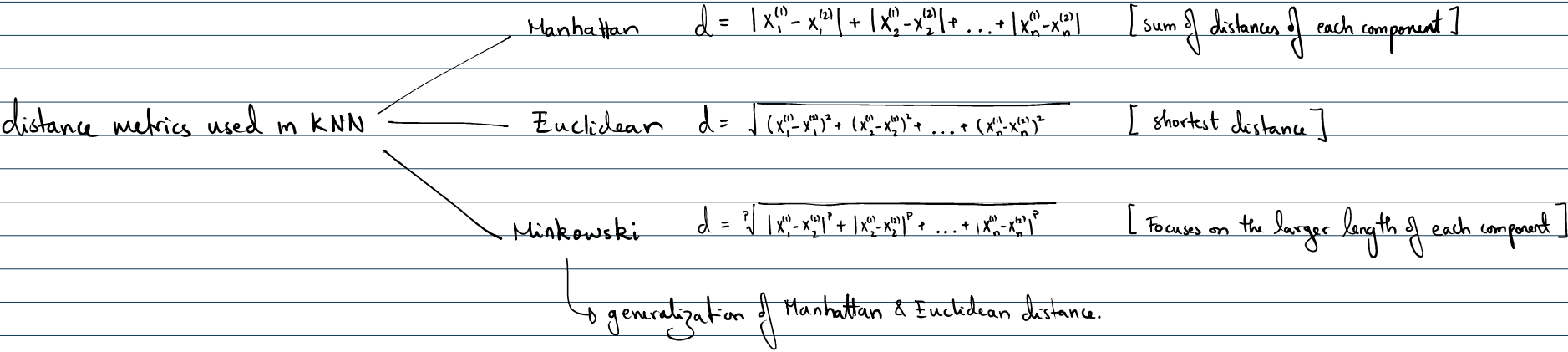
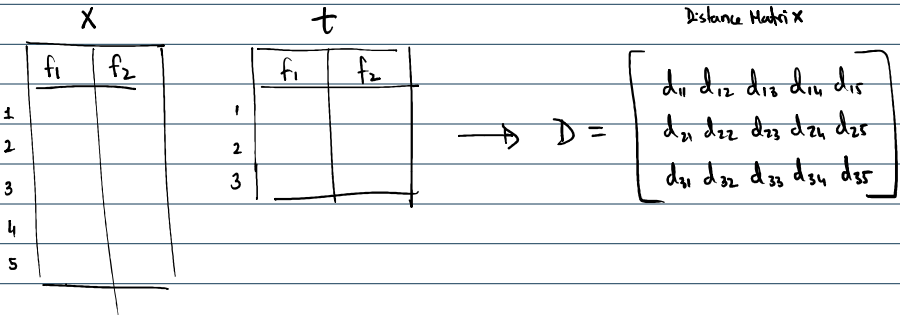
Euclidean distance:

$$\sqrt{(x_1 - a_1)^2 + (x_2 - a_2)^2 + \dots + (x_n - a_n)^2}$$

Matrix Broadcasting

```
x_exp = x[np.newaxis, :, :] # expand D0 axis. (1, 5, 2)
t_exp = t[:, np.newaxis, :] # expand D1 axis. (3, 1, 2)

# Calculating the distance between test and training data.
# Summing data on D2 axis. shape of dist = (3, 5) + distance matrix (D)
dist = np.sqrt(np.sum(np.square(x_exp - t_exp), axis=2))
```



Data normalization

- When the data is as follows, x_1 and x_2 have different scales. The scale of x_2 is much larger than that of x_1 .
- In this case, even if x_2 with a large scale changes slightly, the distance changes significantly. To prevent this phenomenon, the scale of features must be matched. This is called data standardization (or data normalization).
- The most commonly used methods are the Min-Max and Z-score methods.
- The Min-max method converts data into values between 0 and 1, and the Z-score method converts the distribution of data to mean = 0 and standard deviation = 1.
- When normalizing validation or test data, use the values calculated from the training data (min, max, mean, std).

i	feature		class
	x_1	x_2	y
0	0.16	10.01	0
1	-0.09	-20.03	0
2	0.06	-10.08	0
3	0.58	-13.1	2
4	0.46	-11.07	2
5	0.23	12.54	1
6	0.55	15.07	2
7	0.1	21.42	1
8	0.18	17.63	1
9	-0.03	-19.09	0
10	0.59	-23.14	2
11	0.32	29.06	2
12	-0.08	-18.03	0
13	0.27	16.52	1

Min-Max normalization

$$x'_k(i) = \frac{x_k^{(i)} - \min(x_k)}{\max(x_k) - \min(x_k)}$$

Z-score normalization

$$x'_k(i) = \frac{x_k^{(i)} - \text{mean}(x_k)}{\text{std}(x_k)}$$

Decision boundary depends on the value of K. $K \downarrow$ overfitting $K \uparrow$ underfitting

↳ Hyperparameter

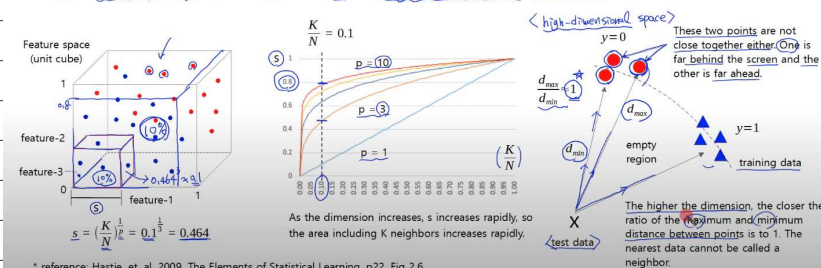
$$K = \sqrt{N}$$

→ Since KNN is an algorithm that uses distances, performance can degrade when features increase. This is called COD. KNN is particularly sensitive to this problem as it uses the concept of neighbors.

→ In higher dimensions distances are not reliable/misleading.

Curse of Dimensionality

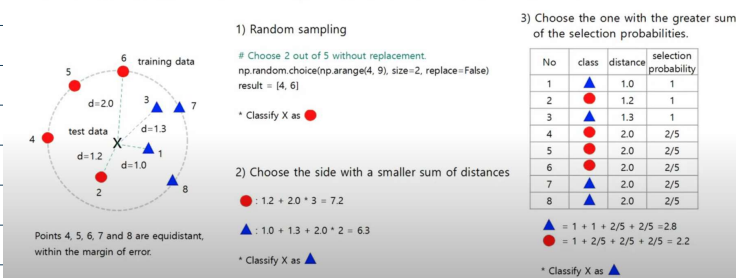
- Since KNN is an algorithm that uses distance, performance can be degraded when features increase. This phenomenon is called the **curse of dimensionality**. Most machine learning algorithms also have this problem, but KNN, which uses the concept of neighbors, is particularly sensitive to this problem.
- If the number of features in the dataset is 3 and the feature values range from 0 to 1, a cube with side length 0.464 is required to keep the number of training data around the test data (10% of the total data, $p=3$, $K/N=0.1$) the length of a side $s = 0.1^{1/3} = 0.464$ (N : the number of training data, K : K value in KNN)
- If there are 10 features ($p=10$), it increases to $s = 0.1^{1/10} = 0.8$, so the area including K neighbors increases rapidly.



KNN → lazy learner → learns every time it estimates without creating a learning model in advance. Takes a long time to estimate. [No learning process]

Dealing with same distance

- If $K=5$, which one should I choose between ● and ▲ for the point X in the following situation?
- Points 1, 2, and 3 are selected unconditionally because they are less than $K=5$, and points 4, 5, 6, 7, and 8 have similar distances, so two of them must be selected. There are three ways to determine the class of x in this situation as follows:



Weighted KNN (WKNN)

- If $K=5$, which one should I choose between ● and ▲ for the point X in the following situation?
- In vanilla KNN, X is classified as ● because ● has three majorities. But X is closer to ▲. Giving higher weight to samples that are closer is called weighted KNN.
- The closer the distance, the greater the weight, and vice versa. (inverse weighting).



- Categorical data has no concept of distance like Manhattan, Euclidean, or Minkowski. Distance (or similarity) for categorical features must be redefined in a different way.

