

Topic 2: k -Nearest Neighbours (kNN)

2.1	Dissimilarity/Distance Measure	44
2.1.1	Manhattan Distance	45
2.1.2	Cosine Similarity	45
2.1.3	Other Dissimilarities	46
2.1.4	Application in Recommender System	48
2.2	Bayes Optimal Classifier Theory	49
2.3	Algorithms of kNN Classifier	50
2.4	Algorithms of kNN Regressor	53
2.5	Classifier Boundary	55
2.6	Feature Scaling	57
2.7	Weighted kNN Models	58
2.8	Optimal: Data Partitioning Algorithms	59
2.9	More Performance Evaluation — k-fold cross validation (CV)	60

The k -nearest neighbours (kNN) algorithm is a **non-parametric** method used for classification and regression (https://en.wikipedia.org/wiki/K-nearest_neighbors_algorithm). When $k = 1$, it is called the **nearest neighbour (NN) algorithm**. The strengths of kNN are simple, effective, makes no assumptions about the underlying data distribution and has a fast training phase; The weaknesses of kNN are having slow classification phase and the inability to provide mathematical reasoning on how the features are related to the class [Lantz, 2015, Chapter 3].

The kNN algorithm will not work for data which are uniformly distributed in very high dimensional feature space, i.e. when $p \gg 1$. This is because the points in p dimensional spaces will be distributed “close” to the boundary and approximately equal distance from each other due to the volume $(1 - 2\epsilon)^d \ll 1$ (curse of dimensionality). So kNN only works for data of “intrinsically” low dimension.

The assumption of kNN is “similar inputs have similar outputs”. Based on this assumption, a test input \mathbf{x} should be assigned the most common label amongst its k most similar training inputs.

Application: Use in classical recommender system such as online bookshop, online movie recommendation, etc.

Given a positive integer k (must be smaller than the number of data, n) and an input \mathbf{x} , the kNN algorithm first identifies the k points in the training data (\mathbf{x}_i, y_i) that are “similar” to \mathbf{x} , represented by $N(\mathbf{x})$.

Definition 2.0.1. $N(\mathbf{x}) \subseteq D$ satisfies $|N(\mathbf{x})| = k$ and $\forall(\mathbf{x}', y') \in D \setminus N(\mathbf{x})$,

$$d(\mathbf{x}, \mathbf{x}') \geq \max_{(\mathbf{x}'', y'') \in N(\mathbf{x})} d(\mathbf{x}, \mathbf{x}''),$$

(i.e. every point in D but not in $N(\mathbf{x})$ is at least as far away from \mathbf{x} as the furthest point in $N(\mathbf{x})$).

We can then define

- kNN classifier:

$$\hat{Y} = h(\mathbf{x}) = \text{mode}(\{y'' : (\mathbf{x}'', y'') \in N(\mathbf{x})\}), \quad (2.1)$$

$$\mathbb{P}(Y = j | \mathbf{X} = \mathbf{x}) = \frac{1}{k} \sum_{\mathbf{x}_i \in N(\mathbf{x})} I(y_i = j). \quad (2.2)$$

where $\text{mode}(\cdot)$ means to select the label of the highest occurrence.

- kNN regressor:

$$h(\mathbf{x}) = \frac{1}{k} \sum_{(\mathbf{x}'', y'') \in N(\mathbf{x})} y''. \quad (2.3)$$

In case of a draw, a good solution is to return the result of k -NN with smaller k .

2.1 Dissimilarity/Distance Measure

The kNN algorithm, clustering algorithms, recommender systems fundamentally relies on the **dissimilarity** metric — a function which takes two “vectors” and returns a non-negative number defined below.

Definition 2.1.1. A *dissimilarity function* $d(\mathbf{x}_i, \mathbf{x}_j)$ is a function which satisfies the following conditions: For any $\mathbf{x}_i, \mathbf{x}_j, \mathbf{x}_k$,

- (i) $d(\mathbf{x}_i, \mathbf{x}_j) \geq 0$ and $d(\mathbf{x}_i, \mathbf{x}_j) = 0$ iff $\mathbf{x}_i = \mathbf{x}_j$ (nonnegativity);
- (ii) $d(\mathbf{x}_i, \mathbf{x}_j) = d(\mathbf{x}_j, \mathbf{x}_i)$ (symmetry); and

When the dissimilarity function satisfies

- (iii) $d(\mathbf{x}_i, \mathbf{x}_j) \leq d(\mathbf{x}_i, \mathbf{x}_k) + d(\mathbf{x}_k, \mathbf{x}_j)$ (triangle inequality).

it is called a **distance/metric function**.

For a kNN classifier, the better that metric reflects label similarity, the better the classification will be.

The most common choice is the *Minkowski distance*:

$$d(\mathbf{x}, \mathbf{z}) = \|\mathbf{x} - \mathbf{z}\|_r = \left(\sum_{i=1}^p |x_i - z_i|^r \right)^{\frac{1}{r}}, \quad \mathbf{x}, \mathbf{z} \in \mathbb{R}^p. \quad (2.4)$$

Note that $\|\cdot\|_r$ is called the ℓ^r norm.

When $r = 1$, we have the *Manhattan distance*:

$$\|\mathbf{x} - \mathbf{z}\|_1 = |x_1 - z_1| + |x_2 - z_2| + \cdots + |x_p - z_p|. \quad (2.5)$$

When $r = 2$, we have the *Euclidean distance*:

$$\|\mathbf{x} - \mathbf{z}\|_2 = \sqrt{(x_1 - z_1)^2 + (x_2 - z_2)^2 + \cdots + (x_p - z_p)^2}. \quad (2.6)$$

When $r = \infty$, we have the *Chebyshev distance*:

$$\|\mathbf{x} - \mathbf{z}\|_\infty = \max\{|x_1 - z_1|, |x_2 - z_2|, \dots, |x_p - z_p|\}. \quad (2.7)$$

Note that the Euclidean distance is more sensitive to outliers than the Manhattan distance. When outliers are rare, the Euclidean distance performs very well and is generally preferred. When the outliers are significant, the Manhattan distance is more stable.

The three most common metrics in user-based recommender system are

- Pearson correlation coefficient (PC):

$$PC : \mathbb{R}^p \times \mathbb{R}^p \rightarrow [-1, 1], \quad (\mathbf{x}, \mathbf{z}) \mapsto \frac{\sum_j (x_j - \bar{c}_j)(z_j - \bar{c}_j)}{\sqrt{\sum_j (x_j - \bar{c}_j)^2 \cdot \sum_j (z_j - \bar{c}_j)^2}}.$$

where the vector $\bar{c} = (\bar{c}_1, \dots, \bar{c}_p)$ is a constant vector (usually the mean vector of a set of data).

- Cosine measure or *cosine similarity*:

$$COS : [0, \infty)^p \times [0, \infty)^p \rightarrow [0, 1], \quad (\mathbf{x}, \mathbf{z}) \mapsto \frac{\sum_j x_j z_j}{\sqrt{\sum_j x_j^2} \cdot \sqrt{\sum_j z_j^2}} =: \frac{\langle \mathbf{x}, \mathbf{z} \rangle}{\|\mathbf{x}\|_2 \cdot \|\mathbf{z}\|_2}.$$

- Mean square difference (similar to MSE except that it compares common rating items only):

$$MSD : \mathbb{R}^p \times \mathbb{R}^p \rightarrow [0, \infty], \quad (\mathbf{x}, \mathbf{z}) \mapsto \frac{\sum_{j \in R_x \cap R_z} (x_j - z_j)^2}{|R_x \cap R_z|}.$$

2.1.1 Manhattan Distance

The Manhattan distance (2.5) is measured along axes at right angles and in R programming, we need to specify its usage by the expression `dist(x, method='manhattan')`. In Python, it is implemented as `cityblock` in `scipy.spatial.distance` and `manhattan_distances` in `sklearn.metrics.pairwise`.

Example 2.1.2 (Final Exam Jan 2019, Q1(a)). State and explain any two types of distance measures. (4 marks)

Solution: Consider two n -dimensional points be $\mathbf{x} = (x_1, \dots, x_n)$, $\mathbf{y} = (y_1, \dots, y_n)$.
 Euclidean distance: $\sqrt{(x_1 - y_1)^2 + \dots + (x_n - y_n)^2}$
 Manhattan distance: $|x_1 - y_1| + |x_2 - y_2| + \dots + |x_n - y_n|$

2.1.2 Cosine Similarity

The cosine similarity is generally used as a metric for measuring distance when the *magnitude of the vectors does not matter*.

The **cosine distance** is a complement of cosine similarity and it is the angular distance between two points, i.e.

$$\text{cos-dist}(\mathbf{x}, \mathbf{y}) = 1 - \text{COS}(\mathbf{x}, \mathbf{y}).$$

It is not implemented in R by default. It is implemented in additional R packages such as `proxy::dist(x, method="cosine")`, `text2vec::sim2(x, method="cosine")`, `lsa::cosine()`, `clv::dot_product()`, `rules::dissimilarity()`. `scipy.spatial.distance.cosine` and `sklearn.metrics.pairwise.cosine_distances` are the implementations in Python.

Cosine distance is a dissimilarity function but **not a distance function**. Consider $A = (1, 0)$, $B = \left(\frac{\sqrt{2}}{2}, \frac{\sqrt{2}}{2}\right)$, $C = (0, 1)$. Note that vectors A and C are orthogonal, so we would get simply 0:

$$\text{COS}(A, C) = \frac{0}{\sqrt{1}\sqrt{1}} = 0.$$

Each pair of vectors A and B as well as B and C would give the same value:

$$\text{COS}(A, B) = \frac{\frac{\sqrt{2}}{2} + 0}{\sqrt{1}\sqrt{1}} = \frac{\sqrt{2}}{2}, \quad \text{COS}(B, C) = \frac{0 + \frac{\sqrt{2}}{2}}{\sqrt{1}\sqrt{1}} = \frac{\sqrt{2}}{2}.$$

Since the cosine distance does not satisfy the triangle inequality and it is therefore not a distance function:

$$\text{cos-dist}(A, C) = 1 - 0 \not\leq \text{cos-dist}(A, B) + \text{cos-dist}(B, C) = \left(1 - \frac{\sqrt{2}}{2}\right) + \left(1 - \frac{\sqrt{2}}{2}\right) = 2 - \sqrt{2}.$$

Despite the cosine distance is not a distance function, it is used in the construction of distance matrix for text data represented by word counts. We could assume that when a word (e.g. science) occurs more frequent in document A than document B, then document A is more related to the topic of science (shorter distance). However, it could also be the case that we are working with documents of uneven lengths. Then, the word ‘science’ probably occurred more in document A just because it was way longer than document B. Cosine similarity corrects for this.

Example 2.1.3. Consider two ‘naive word counts’ for Calculus 1 and Calculus 2 texts with six ‘popular common words’:

	.	=	+	the	f	x
Calculus 1	3626	1446	915	798	552	556
Calculus 2	926	476	317	356	283	146

Calculate the cosine dissimilarity.

Solution:

$$\text{cos-dist}(\text{Calculus 1}, \text{Calculus 2}) = 1 - \frac{4857507}{\sqrt{17326661}\sqrt{1412682}} = 1 - 0.981824 = 0.018176$$

So calculus 1 and calculus 2 text are ‘similar’ based on the 6 words. However, if we include the technical words and remove common words, it may not be the case. Therefore, for text comparison, choosing the right ‘dictionary’ (i.e. the ‘headers’ of the table) is essential.

2.1.3 Other Dissimilarities

Mahalanobis distance (knnGarden::knnMCN):

$$\text{Mahalanobis}(\mathbf{x}, \mathbf{y}) = (\mathbf{x} - \mathbf{y})\Sigma^{-1}(\mathbf{x} - \mathbf{y})^T$$

where \mathbf{x} and \mathbf{y} are row vectors and Σ is the covariance of the data.

The Mahalanobis distance is useful when attributes are correlated, but having different ranges of values and the distribution of the data is approximately Gaussian/normal.

Similarity Measures for Binary Data: Let \mathbf{x} and \mathbf{y} be two objects with binary attributes and p be the number of attributes.

1. Simple Matching Coefficients: $SMC(\mathbf{x}, \mathbf{y}) = \frac{\text{number of matching attribute values}}{p}$

2. Jaccard index / similarity coefficient / Kumar hassebrook distance:

$$J(\mathbf{x}, \mathbf{y}) = \frac{\text{number of matching attribute values}}{\text{number of attribute not involved in 00 matches}}$$

3. Jaccard distance: $d_J(\mathbf{x}, \mathbf{y}) = 1 - J(\mathbf{x}, \mathbf{y})$

Example 2.1.4. Let $\mathbf{x} = (1, 0, 1, 1, 0, 0, 0, 0, 0)$, $\mathbf{y} = (0, 0, 1, 1, 0, 0, 1, 0, 0, 1)$.

$$SMC(\mathbf{x}, \mathbf{y}) = \frac{7 \text{ matches}}{10 \text{ items}} = 0.7$$

$$d_J(\mathbf{x}, \mathbf{y}) = 1 - \frac{2 \text{ '11' pair}}{5 \text{ not '00' pair}} = 1 - 0.4 = 0.6$$

Tanimoto / Soergel score: It is used for data that can take on continuous values. It and Jaccard score are used in Chemoinformatics, plagiarism detection, thesauras extration, market-basket transactional data, anomalies detection in spatio-temporal data:

$$\text{Tanimoto}(\mathbf{x}, \mathbf{y}) = \frac{\sum_i (\max\{x_i, y_i\} - \min\{x_i, y_i\})}{\sum_i \max\{x_i, y_i\}}.$$

For binary data, Tanimoto score is the same as Jaccard distance.

Gower distance: It can be used to process mixed numeric and categorical data:

- The “sub” Gower distance from a number x to the i th row of a numeric column y_i , $i = 1, \dots, n$ is

$$\frac{|x - y_i|}{\max\{y_i\} - \min\{y_i\}}.$$

- The “sub” Gower distance from a factor x to the i th row of a categorical column y_i , $i = 1, \dots, n$ is 0 if the same, else 1:

$$\text{as.integer}(x \neq y_i).$$

- The Gower distance is the average of all “sub” Gower distance for each column.

Example 2.1.5. Consider the fraud data used in the Practical class.

n	gender	age	status	employment	acclink	supplement	base
1	1	32	2	3	0	1	729.3
2	1	57	1	3	0	0	384.1
3	1	21	3	1	0	0	683.8
4	1	27	1	3	0	0	143.0

Solution: The Gower distance between row 1 to rows 2 to 4 can be calculated in R as follows:

```
1 library(gower)
2 gower_dist(fraud[1,1:7], fraud[2:4,1:7])
3 # => 0.4690316 0.4833087 0.4484127
```

The calculations are shown below:

- Gower distance(row 1, row 2):

$$\frac{((1 \neq 1) + \frac{|32-57|}{57-21} + (2 \neq 1) + (3 \neq 3) + (0 \neq 0) + (1 \neq 0) + \frac{|729.3-384.1|}{729.3-143.0})}{7} = 0.4690316$$

- Gower distance(row 1, row 3):

$$\frac{(1 \neq 1) + \frac{|32-21|}{57-21} + (2 \neq 3) + (3 \neq 1) + (0 \neq 0) + (1 \neq 0) + \frac{|729.3-683.8|}{729.3-143.0}}{7} = 0.4833087$$

- Gower distance(row 1, row 4):

$$\frac{(1 \neq 1) + \frac{|32-27|}{57-21} + (2 \neq 1) + (3 \neq 3) + (0 \neq 0) + (1 \neq 0) + \frac{|729.3-143.0|}{729.3-143.0}}{7} = 0.4484127$$

2.1.4 Application in Recommender System

A simple movie database with just 30 entries, seven genres and their IMDB ratings can be found on the Internet (e.g. <https://github.com/Lozadaa/movie-recommender>). The label column values are all zeroes because this dataset is not used for classification (Section 2.3) or regression (Section 2.4).

There are items (e.g. actors, directors, and themes) among the movies that are not accounted because they are missing from the dataset. So, a simple implementation of the ‘k-nearest neighbour’ (kNN) algorithm will solely be based on the similarity on the movies’ genres and the IMDB ratings.

Consider a **case study**: You are building your own movie recommendation website which uses your kNN (with the Euclidean distance) Recommendation Engine at the back-end. You are going to build this back-end Recommendation Engine. Imagine a user is navigating your recommendation website, and he/she encounters a movie named ‘The Post’. The user is not sure if he/she wants to watch it, but its genres intrigue the user; he/she is curious about other similar movies. The user scrolls down to the “More Like This” section to see what recommendations your recommendation website will make, and the back-end algorithmic gears begin to turn.

Your website sends a request to the back-end for the 5 movies that are most similar to ‘The Post’ with the following information (usually in JSON format these days):

IMDB Rating = 7.2, Biography = Yes, Drama = Yes, Thriller = No, Comedy = No, Crime = No, Mystery = No, History = Yes

The back-end has a recommendation data set exactly like ours. It begins by creating the row representation (better known as a feature vector) for ‘The Post’, then it runs a program similar to the one below to search for the 5 movies that are most similar to ‘The Post’, and finally sends the results back to the user at your website.

An implementation in Python is given at https://github.com/jisilvia/kNN_Recommender_System/blob/main/kNN_Recommender_System.ipynb

Similar implementation in R is listed below.

```
# duplicate entries in movies$Movie.ID, row.names=1 is not allowed.
movies = read.csv('movies_recommendation_data.csv')
# 1 = Movie.ID, 2 = Movie.Name, 11 = Label
X = as.matrix(movies[,-c(1,2,11)])

the_post = data.frame(IMDB.Rating=7.2, Biography=1, Drama=1,
  Thriller=0, Comedy=0, Crime=0, Mystery=0, History=1)

# Find similarities of the movie 'the_post' to 'movies' database
euclid.dist = function(x){sqrt(sum(x**2))}
# You need to practise with Practicals 1 and 2 to understand this:
dst = apply(t(t(X) - unlist(the_post)), 1, euclid.dist)
print(cbind(movies[order(dst), 2:10], dist=dst[order(dst)])[1:5,])
```

After calculating the Euclidean distances to the movie ‘The Post’:

Movie.Name	Rating	Biography	Drama	Thriller	Comedy	Crime	Mystery	History
The Post	7.2	1	1	0	0	0	0	1

we obtain the 5 most similar movies below by ranking:

Movie.Name	Rating	Biography	Drama	Thriller	Comedy	Crime	Mystery	History	dist
12 Years a Slave	8.1	1	1	0	0	0	0	1	0.900000
Hacksaw Ridge	8.2	1	1	0	0	0	0	1	1.000000
Queen of Katwe	7.4	1	1	0	0	0	0	0	1.019804
The Wind Rises	7.8	1	1	0	0	0	0	0	1.166190
A Beautiful Mind	8.2	1	1	0	0	0	0	0	1.414214

2.2 Bayes Optimal Classifier Theory

The best classifier that minimises the *misclassification rate* or the *Bayes error rate* / *Bayes risk* [Devroye et al., 1996, Chapter 3],

$$L(f) = \mathbb{P}(\{\omega : Y \neq f(\mathbf{X})\})$$

is called the *Bayes (optimal) classifier* [Mitchell, 1997, Chapter 6]:

$$h_{\text{opt}}(\mathbf{x}) = \underset{y \in \{1, 2, \dots, K\}}{\text{argmax}} \mathbb{P}(Y = y \mid \mathbf{X} = \mathbf{x}). \quad (2.8)$$

The **error rate** of $h_{\text{opt}}(\mathbf{x})$ is

$$\epsilon_{\text{BayesOpt}} = 1 - \mathbb{P}(y^* \mid \mathbf{x}).$$

It is the **theoretical lower bound of the error rate for any classifier**.

Example 2.2.1. Assume an email \mathbf{x} can either be classified as spam (+1) or ham (-1). For the same email \mathbf{x} the conditional class probabilities are:

$$\mathbb{P}(+1 \mid \mathbf{x}) = 0.8, \quad \mathbb{P}(-1 \mid \mathbf{x}) = 0.2.$$

In this case the Bayes optimal classifier would predict the label $y^* = +1$ as it is most likely, and its error rate would be $\epsilon_{\text{BayesOpt}} = 0.2$.

The **upper bound on the error** is given by the constant classifier.

In practise, the distribution of the random variables \mathbf{X} is unknown and so the classifier has to be determined from the observations $\{(\mathbf{x}_i, y_i)\}$ by minimising the expected 0-1 loss:

$$\widehat{L}(f) = \mathbb{E} \left[\frac{1}{n} \sum_{i=1}^n \mathbf{1}(\hat{y}_i \neq y_i) \right]. \quad (2.9)$$

The kNN Classifier's 'ability to predict' is characterised by the following theory.

Theorem 2.2.2. *As $n \rightarrow \infty$, the 1-NN error is no more than twice the error of the Bayes Optimal classifier:*

$$\epsilon_{\text{BayesOpt}} \leq \epsilon_{\text{NN}} \leq 2\epsilon_{\text{BayesOpt}}$$

Similar guarantees hold for $k > 1$.

According to Thomas Cover and Peter Hart, "Nearest neighbor pattern classification", IEEE Transactions on Information Theory, 1967, 13(1):21-27, a **Proof** goes like:

Let \mathbf{x}_{NN} be the nearest neighbour of our test point \mathbf{x}_t . As $n \rightarrow \infty$, $d(\mathbf{x}_{\text{NN}}, \mathbf{x}_t) \rightarrow 0$, i.e. $\mathbf{x}_{\text{NN}} \rightarrow \mathbf{x}_t$. This means the nearest neighbour is identical to \mathbf{x}_t . You return the label of \mathbf{x}_{NN} . What is the probability that this is not the label of \mathbf{x}_t ? It is the probability of drawing two different label of \mathbf{x} :

$$\begin{aligned} \epsilon_{\text{NN}} &= \mathbb{P}(y^* \mid \mathbf{x}_t)(1 - \mathbb{P}(y^* \mid \mathbf{x}_{\text{NN}})) + \mathbb{P}(y^* \mid \mathbf{x}_{\text{NN}})(1 - \mathbb{P}(y^* \mid \mathbf{x}_t)) \\ &\leq (1 - \mathbb{P}(y^* \mid \mathbf{x}_{\text{NN}})) + (1 - \mathbb{P}(y^* \mid \mathbf{x}_t)) = 2(1 - \mathbb{P}(y^* \mid \mathbf{x}_t)) = 2\epsilon_{\text{BayesOpt}}, \end{aligned}$$

where the inequality follows from $\mathbb{P}(y^* \mid \mathbf{x}_t) \leq 1$ and $\mathbb{P}(y^* \mid \mathbf{x}_{\text{NN}}) \leq 1$. We also used that $\mathbb{P}(y^* \mid \mathbf{x}_t) = \mathbb{P}(y^* \mid \mathbf{x}_{\text{NN}})$.

2.3 Algorithms of kNN Classifier

A pure Python implementation of kNN is given in Grus [2015, Chapter 12]. If numpy library is used, a simple implementation (involving brute force search) can be stated below.

```

1 # https://github.com/joelgrus/data-science-from-scratch
2 from scipy.spatial.distance import euclidean
3 import logging
4 logging.basicConfig(level=logging.DEBUG)
5
6 def knn_classify(k: int, X, y, x_new, distance=euclidean) -> str:
7     import numpy as np, collections
8     dist_list = [distance(x, x_new) for x in X]
9     logging.debug('distances={}'.format([round(d,4) for d in dist_list]))
10    order = np.argsort(dist_list)
11    k_nearest_labels = y[order[:k]]
12    logging.debug('knn={}'.format(k_nearest_labels))
13    vote_counts = collections.Counter(k_nearest_labels)
14    winner, winner_count = vote_counts.most_common(1)[0]
15    return winner, winner_count/k
16
17 if __name__=='__main__':
18     from sklearn import datasets
19     iris_df = datasets.load_iris()
20     X = iris_df['data']; y = iris_df['target']
21     knn_classify(1,X,y,[6,3,5,2]); knn_classify(1,X,y,[5,3,1,0])

```

Python's `sklearn.neighbors.KNeighborsClassifier` implements a more sophisticated algorithm involving KD tree and ball tree.

In R, the `knn` classifier is available in the `class` library which only supports Euclidean distance and has the following form.

```
knn(train, test, cl, k = 1, l = 0, prob = FALSE, use.all = TRUE)
```

Here,

- `train`: matrix or data frame of training set cases.
- `test`: matrix or data frame of test set cases.
- `cl`: factor of true classifications of training set
- `k`: number of neighbours considered.

Let us practise the algorithm with an example. More examples are found in the tutorial.

Example 2.3.1. A sport school would like to group their new enrolled students into 2 groups, as according to the existing students' weight and height. The weight and height of 7 existing students with group are shown in the table below.

Student	Weight (kg)	Height (cm)	Group
A	29	118	A
B	53	137	B
C	38	127	B
D	49	135	B
E	28	111	A
F	24	111	A
G	30	121	A

- (a) Perform and use $k = 3$ -NN method (with Euclidean distance) to predict which group the following students will be grouped into, based on a **cut-off of 0.7** and a default cut-off on group A.

Student	Weight (kg)	Height (cm)
H	35	120
I	47	131
J	22	115
K	38	119
L	31	136

Solution: First, we construct a distance table from the testing data to the training data and take note that $k = 3$:

		H	I	J	K	L
Student	Group	Distance	Distance	Distance	Distance	Distance
A	A	6.3246	22.2036	7.6158	9.0554	18.1108
B	B	24.7588	8.4853	38.0132	23.4307	22.0227
C	B	7.6158	9.8489	20.0000	8.0000	11.4018
D	B	20.5183	4.4721	33.6006	19.4165	18.0278
E	A	11.4018	27.5862	7.2111	12.8062	25.1794
F	A	14.2127	30.4795	4.4721	16.1245	25.9615
G	A	5.0990	19.7231	10.0000	8.2462	15.0333
$P(Y = A)$		0.6667	0.0000	1.0000	0.6667	0.3333
Cut-off = 0.5	\hat{y}	A	B	A	A	B
Cut-off = 0.7	\hat{y}	B	B	A	B	B

How do we get $d(A, H) = 6.3246$? The calculation is as follows:

$$d(A, H) = \sqrt{(29 - (35))^2 + (118 - (120))^2} = \sqrt{40} = 0.5673$$

The rest is similar.

- (b) The actual groups of the students are {A, B, A, B, B} for students {H, I, J, K, L} respectively. Construct a confusion matrix and calculate the accuracy measurements for a cut-off of 0.5 and a cut-off of 0.7.
 (Answer: For a cut-off 0.5: ACC=0.8, TPR=1, TNR=0.6767, PPV=0.6667, NPV=1; For a cut-off 0.7: ACC=0.8, TPR=0.5, TNR=1, PPV=1, NPV=0.75)

- (c) Write a Python script to produce the above calculations using the simple implementation by lecturer in Section 2.3 and also use sklearn's implementation.

Solution:

```

1 import pandas as pd
2 d_train = pd.DataFrame({
3     'weight': [29,53,38,49,28,24,30],
4     'height': [118,137,127,135,111,111,121],
5     'group' : ['A','B','B','B','A','A','A']
6 }, index = list('ABCDEFG'))
7 d_test = pd.DataFrame({
8     'weight': [35, 47, 22, 38, 31],
9     'height': [120, 131, 115, 119, 136]
10 }, index = list('HIJKL'))
11
12 from knn_naive import *
13 X = d_train[['weight','height']].values
14 y = d_train['group']
15 X_test = d_test[['weight','height']].values
16 for new_x in X_test:
17     print(knn_classify(3,X,y,new_x))
18
19 from sklearn.neighbors import KNeighborsClassifier
20 model = KNeighborsClassifier(n_neighbors=3, p=2)
21 model.fit(X,y)
22 print(model.predict(X_test)) # cut-off = 0.5
23 print(model.predict_proba(X_test))

```

- (d) Write an R script to produce the above calculations.

Solution:

```

1 X_y.train = data.frame(row.names = LETTERS[1:7],
2     weight = c(29,53,38,49,28,24,30),
3     height = c(118,137,127,135,111,111,121),
4     group = c('A','B','B','B','A','A','A'))
5 X.test = data.frame(row.names = LETTERS[8:12],
6     weight = c(35, 47, 22, 38, 31),
7     height = c(120, 131, 115, 119, 136))
8
9 dist.table = data.frame(row.names = LETTERS[1:7])
10 norm2 = function(x){return(sqrt(sum(x^2)))}
11 X = X_y.train[,1:2]
12 print(X)

```

```

13 for(i in 1:nrow(X.test)) {
14     dist.table[,i] = apply(t(X)-t(X.test)[,i],2,norm2)
15 }
16 names(dist.table) = row.names(X.test)
17 print(round(dist.table,4))
18
19 library(class) # for knn
20 M = ncol(X_y.train)
21 X.test$group = knn(train=X_y.train[,-M], test=X.test,
22                   cl =X_y.train[, M], k=3, prob=TRUE)
23 X.test$prob = attr(X.test$group,"prob")
24 X.test$prob = ifelse(X.test$group=="A", X.test$prob, 1-X.test$prob)
25 cut.off = 0.7
26 X.test$group = ifelse(X.test$prob>cut.off, "A", "B")
27 print(cbind(round(X.test[,1:2],4),prob.A=round(X.test[,4],4),
28            predicted=X.test$group, actual=c('A','B','A','B','B'))))

```

2.4 Algorithms of kNN Regressor

The algorithms for the *kNN regressor* are similar to the kNN classifier except that the “mean” is used to estimate the output as in (2.3) [Altman, 1992].

In Python, `sklearn.neighbors.KNeighborsRegressor` implements the kNN regressor with Minkowski distance; In R, FNN’s `knn.reg` provides an implementation with the Euclidean distance.

Example 2.4.1 (Exam SRM Study Manual, p225, Q15.10). A continuous variable Y is modelled as a function of X using kNN with $k = 3$. With the following data:

X	5	8	15	22	30
Y	4	1	10	16	30

Calculate the fitted value of Y at $X = 12$. Try write a script using R (or Python) to perform the calculation for you.

(Answer: $\hat{Y} = 5$)

Example 2.4.2 (Final Exam May 2024 Sem, Q5(a)). Given the training data with two numeric features “age”, “experience” and the response “income” in Table 5.1.

Table 5.1: Training data of the income data.

Obs.	age	experience	income
A	44	9	44190
B	43	10	47830
C	25	1	30450
D	30	3	35670
E	51	7	41630
F	28	5	41340
G	37	10	48700
H	54	5	36720

- (i) Use the supervised learning model kNN ($k=3$) with the Euclidean distance to predict the income for a 47 years old person with 2 years experience and a 41 years old person with 6 years experience. (10 marks)

(Answer: $Y_{kNN}(P1)=40846.67$, $Y_{kNN}(P2)=46906.67$)

- (ii) Use the cosine distance formula $\text{cos-dist}(\mathbf{x}, \mathbf{z}) = 1 - \frac{\sum_j x_j z_j}{\sqrt{\sum_j x_j^2} \sqrt{\sum_j z_j^2}}$ to calculate $\text{cos-dist}(\text{age}, \text{experience})$ for the data in Table 5.1. (2 marks)

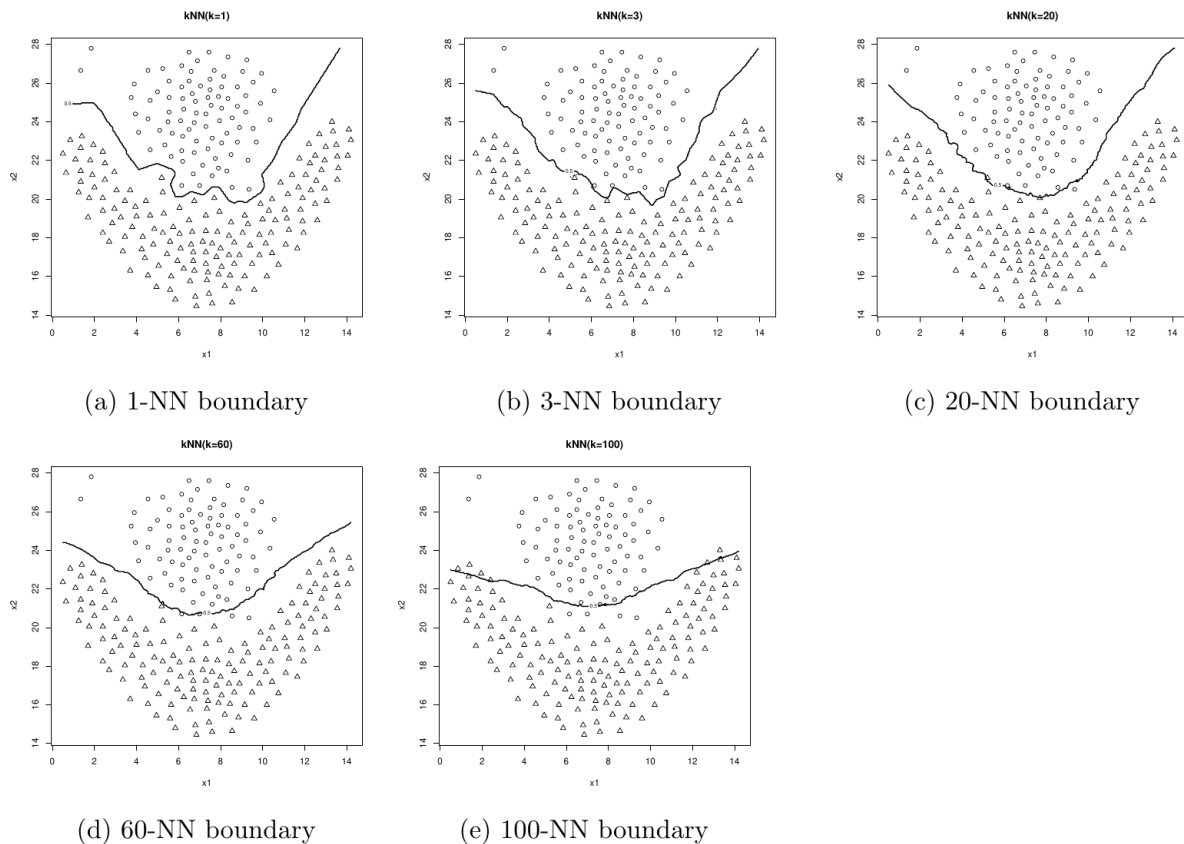
(Answer: 0.076417)

2.5 Classifier Boundary

Classifier boundary = Points \mathbf{x} in the domain which do not belong to any output class because $P(Y|X = \mathbf{x})$ is not clearly defined.

Example 2.5.1. Analyse the ‘flame’ dataset from <https://cs.joensuu.fi/sipu/datasets/flame.txt> using the kNN classifier with $k = 1$, $k = 3$, $k = 20$, $k = 60$, $k = 100$.

Solution: The boundaries of 1-NN (Figure 2.1a) and 3-NN (Figure 2.1b) are rather “nonlinear” and “rough”. As k is increased to 100, the boundary becomes “smoother” and usually “straighter”: when $k = 20$, the boundary seems to be “smoother” in Figure 2.1c while the boundary becomes “straighter” changing when k is increased to $k = 60$ (Figure 2.1d) to $k = 100$ (Figure 2.1e).

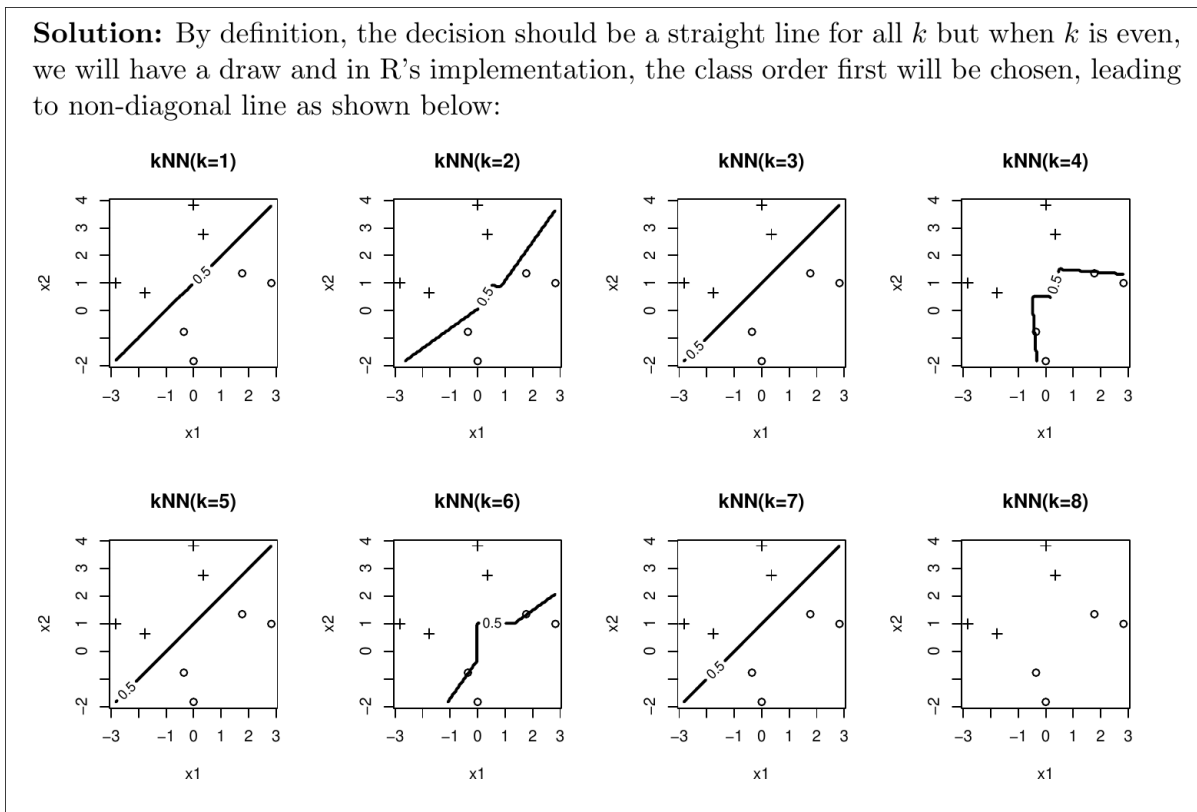


Example 2.5.2 (Symmetric Data Can Lead to Bias due to Class Order). Consider the data

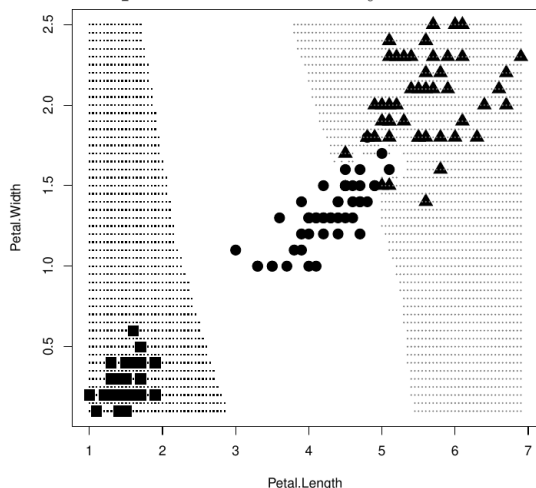
x_1	x_2	y
0.0000	3.8284	+
0.3536	2.7678	+
-2.8284	1.0000	+
-1.7678	0.6464	+
2.8284	1.0000	o
1.7678	1.3536	o
0.0000	-1.8284	o
-0.3536	-0.7678	o

Analyse the decision boundary for $k = 1$ to $k = 8$.

Solution: By definition, the decision should be a straight line for all k but when k is even, we will have a draw and in R's implementation, the class order first will be chosen, leading to non-diagonal line as shown below:



Example 2.5.3 (Three-class Output). The iris data has 4 variables, but if we pick just two variables, we can “visualise” the decision boundary using ‘point labels’ or ‘colours’ rather than the contour plot as in the binary classifications in the earlier examples.



We are fortunate that the three classes of iris flowers are ‘distinguishable’. Otherwise, the boundary will be mixed and very difficult to identify.

Conclusion: the choice of k depends upon the data. In general, larger values of k reduces effect of the noise on the classification, but make boundaries between classes less distinct. A good k can be selected by various heuristic techniques (e.g. based of prediction accuracy). When $k = 1$, the decision boundary is over “flexible”. This corresponds to a classifier that has low bias but very high variance. As k grows, the classifier becomes less flexible and produces a decision boundary that is closed to linear. This corresponds to a low-variance but high bias classifier.

2.6 Feature Scaling

- Features with large variations may shadow features with small variations in Euclidean distance measures.
- Scaling such as min-max scaling and standardisation (using R's `scale`) can help putting all features to the same ground.

Example 2.6.1 (Standardisation). For the training data and testing data from Example 2.3.1, apply the standardisation and then perform $k = 3$ -NN (with Euclidean distance) to predict which group the following students will be grouped into, based on a **cut-off of 0.5** and a **cut-off of 0.7** on group A.

Solution:

Student	Weight (kg)	Height	Std_Wgt	Std_Hgt	Group
A	29	118	-0.6113	-0.4587	A
B	53	137	1.5284	1.3355	B
C	38	127	0.1910	0.3912	B
D	49	135	1.1717	1.1467	B
E	28	111	-0.7005	-1.1197	A
F	24	111	-1.0571	-1.1197	A
G	30	121	-0.5222	-0.1754	A
Mean	35.8571	122.8571			
Std. dev	11.2165	10.5898			

Important Note: The training data **MUST** used the standardisation used in the training data because the kNN classifier is **TRAINED** with the training data.

$$I: \frac{47 - 35.8571}{11.2165} \approx 0.9934, \quad \frac{131 - 122.8571}{10.5898} \approx 0.7689$$

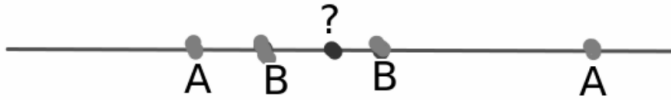
Student	Weight (kg)	Height (cm)	Std_Wgt	Std_Hgt
H	35	120	-0.0764	-0.2698
I	47	131	0.9934	0.7689
J	22	115	-1.2354	-0.7420
K	38	119	0.1910	-0.3642
L	31	136	-0.4330	1.2411

By performing Euclidean distance calculation, we obtain $d(A, H) = \sqrt{(-0.0764 - (-0.6113))^2 + (-0.2698 - (-0.4587))^2} = \sqrt{0.3218} = 0.5673$, etc. leading to the table below.

		H	I	J	K	L
	Group	Distance	Distance	Distance	Distance	Distance
A	A	0.5673	2.0205	0.6854	0.8079	1.7091
B	B	2.2699	0.7792	3.4575	2.1628	1.9637
C	B	0.7131	0.8869	1.8218	0.7554	1.0544
D	B	1.8879	0.4177	3.0596	1.8013	1.6076
E	A	1.0544	2.5370	0.6548	1.1686	2.3759
F	A	1.2977	2.7878	0.4177	1.4590	2.4419
G	A	0.4557	1.7857	0.9109	0.7378	1.4193
$P(Y = A)$		0.6667	0.0000	1.0000	0.6667	0.3333
Cut-off = 0.5	\hat{y}	A	B	A	A	B
Cut-off = 0.7	\hat{y}	B	B	A	B	B

2.7 Weighted kNN Models

Some academicians propose an extension to kNN by considering the following scenario:



Here, a new point marked with ‘?’ has 4 similar/closest points. Suppose the kNN with $k = 4$ is used. The conditional probability

$$P(Y = A|X = ?) = P(Y = B|X = ?) = 0.5$$

is obtained. If we use dictionary ordering, A is before B , so $kNN(k = 4)$ will give $Y = A$. This may not be reasonable since B is closer to ‘?’.

Weighted kNN (wkNN, available in R in `kknn` library) introduces weights (called kernel) to solve this problem [Hechenbichler and Schliep, 2004].

1. Let $N(\mathbf{x}, k + 1)$ be the $k + 1$ nearest neighbours to \mathbf{x} according to a distance function $d(\mathbf{x}, \mathbf{x}_i)$.
2. The $(k + 1)$ th neighbour is used for “standardisation” of the k smallest distances via

$$D_i = \frac{d(\mathbf{x}, \mathbf{x}_i)}{d(\mathbf{x}, \mathbf{x}_{k+1})}.$$

3. A weighted majority of the k nearest neighbour

$$\hat{y} = \max_j \left\{ \sum_{i=1}^k K(D_i) I(y_i = j) \right\}.$$

In order for the ‘weight’ to make sense, we need it to satisfy some requirements. The ‘weight’ K is called a *kernel (function)* if it satisfies

- (1) $K(x) \geq 0$ for all $x \in \mathbb{R}$;
- (2) $K(x)$ is maximum when $x = 0$;
- (3) $K(x)$ descends monotonously when $x \rightarrow \pm\infty$.

Available kernels in R’s `kknn` or Python’s `sklearn`:

- rectangular/uniform: $K(x) = \frac{1}{2} I(|x| \leq 1)$;
- inv/distance: $K(x) = 1/|x|$
- triangular: $K(x) = (1 - |x|) \cdot I(|x| \leq 1)$.
- optimal: The number of neighbours used for this kernel should be $\left(\frac{2(d+4)}{d+2}\right)^{d/(d+4)} k \in [1.2, 2]$, where k is used
- Others: `cos`, `gaussian`, `rank`, `epanechnikov` (or `beta(2,2)`), `biweight` (or `beta(3,3)`), `triweight` (or `beta(4,4)`).

The wkNN is implemented in Python’s `KNeighborsClassifier` and `KNeighborsRegressor` and R’s `kknn`’s `kknn`. By default, Python’s kNN algorithms sets the “weights” parameter to “uniform”, which is the same as the usual kNN.

However, R’s `kknn` sets the “kernel” to “optimal” by default, which leads to a different calculation from the usual kNN. To set it to be the usual kNN, the “kernel” needs to be set to “rectangular”. For each row of the test set, the k nearest training set vectors (according to Minkowski distance) are found, and the classification is done via the maximum of summed kernel densities. In addition, even ordinal and continuous variables can be predicted.

2.8 Optimal: Data Partitioning Algorithms

To **make kNN faster during prediction**, Python sklearn's kNNs implement more sophisticated methods to “find” the nearest neighbours.

brute-force method. It finds all distances and sort in order.

kd_tree method. It partitions the feature space by (a) dividing the data into two halves, e.g. left and right, along one feature. (b) For each training input, remember the half it lies in.

How can this partitioning speed up prediction? Let's think about it for the one neighbour case.

1. Identify which side the test point lies in, e.g. the right side.
2. Find the nearest neighbour x_{NN}^R of x_t in the same side. The R denotes that our nearest neighbour is also on the right side.
3. Compute the distance between x_y and the dividing "wall". Denote this as d_w . If $d_w > d(x_t, x_{\text{NN}}^R)$ you are done, and we get a 2x speedup.

In other words: if the distance to the partition is larger than the distance to our closest neighbour, we know that none of the data points inside that partition can be closer. We can avoid computing the distance to any of the points in that entire partition. We can prove this formally with the triangular inequality.

Let $d(x_t, x)$ denote the distance between our test point x_t and a candidate x . We know that x lies on the other side of the wall, so this distance is dissected into two parts $d(x_t, x) = d_1 + d_2$, where d_1 is the part of the distance on x_t 's side of the wall and d_2 is the part of the distance on x 's side of the wall. Also let d_w denote the shortest distance from x_t to the wall. We know that $d_1 > d_w$ and therefore it follows that

$$d(x_t, x) = d_1 + d_2 \geq d_w + d_2 \geq d_w.$$

This implies that if d_w is already larger than the distance to the current best candidate point for the nearest neighbor, we can safely discard x as a candidate.

The tree construction involves two steps: (a) Split data recursively in half on exactly one feature. (b) Rotate through features by picking the feature with maximum variance.

The pros of the tree construction is its “exactness” and its easiness to build. The cons are (1) The curse of dimensionality makes KD-Trees ineffective for higher number of dimensions; (2) All splits are axis aligned.

ball_tree method. It uses hyper-spheres (balls) instead of the boxes in KD-tree. As in `kd_tree` method, we can dissect the distance and use the triangular inequality

$$d(x_t, x) = d_1 + d_2 \geq d_b + d_2 \geq d_b \tag{2.10}$$

If the distance to the ball, d_b , is larger than distance to the currently closest neighbor, we can safely ignore the ball and all points within. The ball structure allows us to partition the data along an underlying manifold that our points are on, instead of repeatedly dissecting the entire feature space (as in KD-Trees).

Ball-trees are slower than KD-Trees in low dimensions ($d \leq 3$) but a lot faster in high dimensions. Both are affected by the curse of dimensionality, but Ball-trees tend to still work if data exhibits local structure (e.g. lies on a low-dimensional manifold).

The kNN is slow during prediction because it does a lot of unnecessary work. KD-trees partition the feature space so we can rule out whole partitions that are further away than our closest k neighbours. However, the splits are axis aligned which does not extend well to higher dimensions. Ball-trees partition the manifold the points are on, as opposed to the whole space. This allows it to perform much better in higher dimensions.

The default algorithm in Python is `auto`, which attempts to decide the most appropriate algorithm (out of `brute`, `ball_tree`, `kd_tree`) based on the values passed to `fit` method.

2.9 More Performance Evaluation — *k*-fold cross validation (CV)

The development of CRISP-DM until now:

- Business understanding
- **Data understanding**: structured and unstructured data, EDA (Section 1.5)
- **Data preparation / preprocessing**: feature scaling
- **Modelling**: kNN, wkNN
- **Evaluation**: MSE, R^2 for regression problems; confusion matrix / contingency table, ROC for (binary) classification problems; holdout/validation set/train-test-split method
- Deployment

A few issues regarding the holdout/validation set method from Topic 1:

- **The model is not trained with all available data** since the available data is split into training set and testing set.
- **The result depends on the sampling of the data.** Note that in Practical 3, there is a `set.seed(123)` which fixes a particular **sampling**. If we change the number to 124 or 125, we will get different results for the performance due to **different sampling**. If we apply stratified sampling with odd index, we get an accuracy of 0.8011988 and with even index, we get an accuracy of 0.7654691

Since holdout method is “sampling” dependent, to try to learn the ‘real’ performance of the (kNN) model, we can cover all possible sampling. However,

for a data with 1000 rows, if we choose 700 rows for training and 300 rows for testing using linear sampling, there are

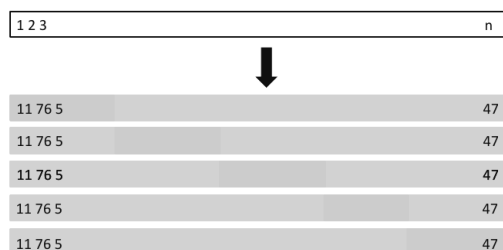
$$\frac{1000!}{700!300!} \approx 5.428250 \times 10^{263}$$

possible combinations! The number is just **too large**.

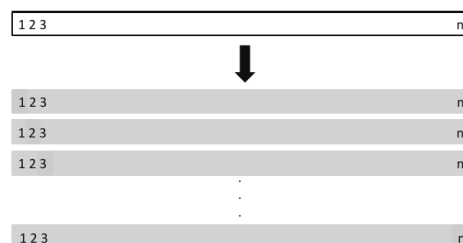
To have a better confidence on the performance than the holdout method we can use the **cross-validation methods** which are found to be acceptable by statisticians.

- Randomly divides the set of observations into k “equal” “folds”;
- First fold = validation set & remaining $k - 1$ folds = training set.
- Second fold = validation set & remaining $k - 1$ folds = training set.
- etc.

Example 2.9.1 (5-fold).



Example 2.9.2 (LOOCV).



leave-one-out cross validation (LOOCV) is a special case k -fold CV with $k = n$.

Software support:

R: `caret::createDataPartition`, `createResamples` (for bootstrapping), `createFolds` (for *k*-fold cross validation), `createMultiFolds` (for repeated cross-validation), `createTimeSlices` (for dealing with time-series).

Python: `train_test_split` (simple random split), `KFold` or `StratifiedKFold` (*k*-fold CV), `LeaveOneOut` (LOOCV) from `sklearn.cross_validation`.

If we don't have the right to install additional package like the `caret` library, we can still use R's basic commands to perform *k*-fold cross validation. The *k* = 10 example is shown below.

<https://stats.stackexchange.com/questions/61090/how-to-split-a-data-set-to-do-10-fold-cross-validation>

```
#Randomly shuffle the data & create 10 fold
d.f = d.f[sample(nrow(d.f)),]
folds = cut(seq(1,nrow(d.f)),breaks=10,labels=FALSE)
#Perform 10 fold cross validation
for(i in 1:10){
  #Segment your data by fold using which()
  testIndexes = which(folds==i,arr.ind=TRUE)
  testData = d.f[testIndexes, ]
  trainData = d.f[-testIndexes, ]
  ...
}
```

Using `caret` library (it is a complex R packages with many dependencies):

```
set.seed(123)
train.control <- trainControl(method="cv", number=5)
model <- train(y ~., data=train.data, method="lm",
              trControl=train.control)
print(model)
```

Note that *k*-fold CV is still sampling dependent but the averaging has reduced the bias problem compare to the holdout method.

