

Tut 2: kNN

Feb 2026

kNN is discriminative, non-parametric predictive model

- For kNN classifier, the mathematical formulation is

$$\hat{h}(\mathbf{x}) = \underset{j \in \{1, \dots, K\}}{\operatorname{argmax}} \frac{1}{k} \sum_{\mathbf{x}_i \in N(\mathbf{x})} I(y_i = j)$$

- For kNN regressor, the mathematical formulation is

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

One popular choice of **distance** (nonnegative, symmetric and triangle inequality) in kNN is the *Minkowski distance of order/degree r* :

$$d(\mathbf{x}, \mathbf{z}) = \|\mathbf{x} - \mathbf{z}\|_r = \left(\sum_{i=1}^p |x_i - z_i|^r \right)^{\frac{1}{r}}, \quad \mathbf{x}, \mathbf{z} \in \mathbb{R}^p. \quad (2.1)$$

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

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

$$\|\mathbf{x} - \mathbf{z}\|_1 = |x_1 - z_1| + |x_2 - z_2| + \dots + |x_p - z_p|.$$

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

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

There are other distance / dissimilarity functions which are used in specific cases:

- Gower; • Tanimoto; • Jaccard; • Mahalanobis

Distance Measurements

1. (Final Exam Feb 2026 Sem, Q3(a)) State the **three conditions** that define a **distance function** and hence determine whether the Pearson correlation coefficients

$$PC(\mathbf{x}, \mathbf{z}) = \frac{\sum_j (x_j - \bar{c}_j)(z_j - \bar{c}_j)}{\sqrt{\sum_j (x_j - \bar{c}_j)^2 \cdot \sum_j (z_j - \bar{c}_j)^2}}$$

is a distance function or not with proper justification. The vector $(\bar{c}_1, \dots, \bar{c}_p)$ is a constant vector representing the centre of a set of data. (5 marks)

Solution. The three conditions that define a distance function are

- (a) Nonnegativity: $d(\mathbf{x}_i, \mathbf{x}_j) \geq 0$ and $d(\mathbf{x}_i, \mathbf{x}_j) = 0$ iff $\mathbf{x}_i = \mathbf{x}_j$ for any $\mathbf{x}_i, \mathbf{x}_j$; ... [1 mark]
- (b) Symmetry: $d(\mathbf{x}_i, \mathbf{x}_j) = d(\mathbf{x}_j, \mathbf{x}_i)$ for any $\mathbf{x}_i, \mathbf{x}_j$; [1 mark]
- (c) Triangle inequality: $d(\mathbf{x}_i, \mathbf{x}_j) \leq d(\mathbf{x}_i, \mathbf{x}_k) + d(\mathbf{x}_k, \mathbf{x}_j)$ for any $\mathbf{x}_i, \mathbf{x}_j, \mathbf{x}_k$ [1 mark]

The Pearson correlation coefficients is **not a distance function** because it does not satisfy the nonnegativity condition. For example, consider the pair \mathbf{x} and $-\mathbf{x}$ and the centre is a

zero vector $\mathbf{c} = (0, \dots, 0)$, we have $PC(\mathbf{x}, \mathbf{x}) = 1 \neq 0$ and

$$PC(\mathbf{x}, -\mathbf{x}) = \frac{\sum_j (x_j - 0)(-x_j - 0)}{\sqrt{\sum_j (x_j - 0)^2 \cdot \sum_j (-x_j - 0)^2}} = \frac{\sum_j (-x_j^2)}{\sqrt{(\sum_j (x_j)^2)^2}} = \frac{-\sum_j x_j^2}{\sum_j (x_j)^2} = -1 < 0$$

[2 marks]

□

kNN Predictive Models

2. The given table provides a training data set containing six observations, three predictors and one qualitative response variable. Suppose we wish to use this data set to make a prediction for Y when $X_1 = X_2 = X_3 = 0$ using k -nearest neighbours.

Obs.	X_1	X_2	X_3	Y
1	0	3	0	Red
2	2	0	0	Red
3	0	1	3	Red
4	0	1	2	Green
5	-1	0	1	Green
6	1	1	1	Red

- (a) Compute the Euclidean distance between each observation and the test point (TP).

<i>Solution.</i>	Obs.	X_1	X_2	X_3	Y	Distance	□
	1	0	3	0	Red	3	
	2	2	0	0	Red	2	
	3	0	1	3	Red	$\sqrt{10} \approx 3.1623$	
	4	0	1	2	Green	$\sqrt{5} \approx 2.2361$	
	5	-1	0	1	Green	$\sqrt{2} \approx 1.4142$	
	6	1	1	1	Red	$\sqrt{3} \approx 1.7321$	

- (b) What is our prediction with $k = 1$? Why?

Solution. Green. Observation 5 is the closest neighbour for $k = 1$.

□

- (c) What is our prediction with $k = 3$? Why?

Solution. Red. Observations 2, 5 and 6 are the closest neighbours for $k = 3$, which Y equal to (Red, Green, Red). The probability of Red is two-third, which is larger than

0.5. $\mathbb{P}(Y = \text{Red}) = \frac{2}{3} \geq 0.5$ Hence, the test point will be predicted to be Red.

□

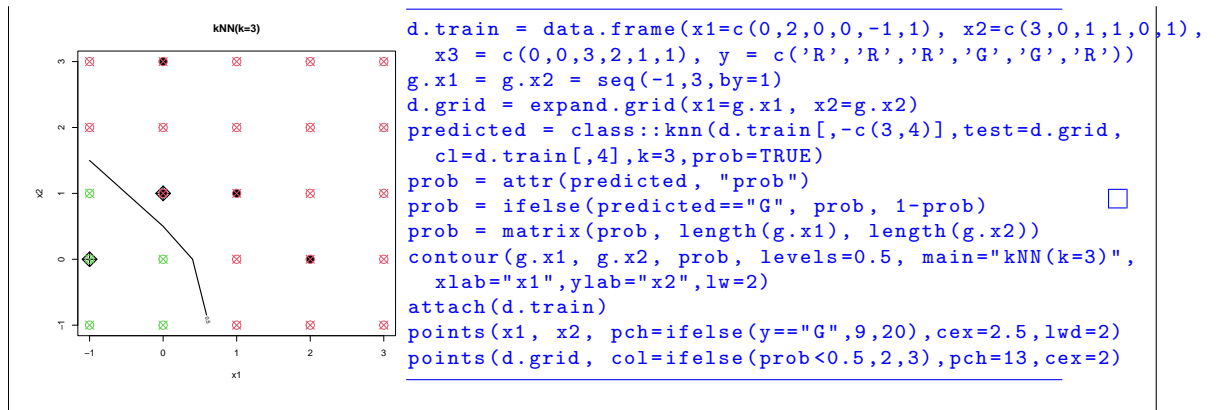
- (d) If the Bayes decision boundary in this problem is highly non-linear, then would we expect the optimum value for k to be large or small? Why?

Solution. Small. A small k would be flexible for a non-linear decision boundary, whereas a large k would try to fit a more linear boundary because it takes more points into consideration.

□

- (e) By considering X_1 and X_2 only, sketch the 3-nearest neighbours decision boundary for range $-1 \leq X_1 \leq 3$ and $-1 \leq X_2 \leq 3$, with the distance measure used in (a). Assume that X_1 and X_2 can only take integer values.

Solution.



3. (Final Exam May 2023 Sem, Q5(a)(i)) Given the training data with features X_1 , X_2 and the label Y in Table 5.1.

Obs.	Petal.Length	Petal.Width	Sepal.Length	Species
1	1.5	0.2	5.0	setosa
2	1.1	0.1	4.3	setosa
3	4.0	1.2	5.8	versicolor
4	3.3	1.0	4.9	versicolor
5	5.4	2.1	6.9	virginica
6	5.1	1.9	5.8	virginica

Table 5.1: Training data with features Petal.Length, Petal.Width, Sepal.Length and the label Species of iris flower.

Given an iris flower with a petal length of 3.9, a petal width of 1.4 and a sepal length of 5.2. Use the Euclidean distance and the supervised learning model kNN ($k=3$) to predict the Species of the iris flower. (7 marks)

Solution. By calculating the Euclidean distance from the point (3.9, 1.4, 5.2) to the points in the training data, we can obtain the following table:

Petal.Length	Petal.Width	Sepal.Length	Species	Distance
1.5	0.2	5.0	setosa	2.6907
1.1	0.1	4.3	setosa	3.2156
4.0	1.2	5.8	versicolor	0.6403
3.3	1.0	4.9	versicolor	0.7810
5.4	2.1	6.9	virginica	2.3728
5.1	1.9	5.8	virginica	1.4318

..... [6 marks]

The 3 nearest neighbours are observations 3, 4 and 6, which correspond to Species versicolor, versicolor and virginica. Therefore, the prediction of the Species of the iris flower is versicolor. [1 mark] ☐

4. (Final Exam May 2024 Sem, Q5(a)) Given the training data with four numeric features “bill length” (unit: mm), “bill depth” (unit: mm), “flipper length” (unit: mm), “body mass” (unit: g) and the label “species” in Table 5.1.

Table 5.1: Training data of the penguin data with three types of penguins — Adelie, Chinstrap and Gentoo.

Obs.	bill length	bill depth	flipper length	body mass	species
A	41.1	19.1	188	4100	Adelie
B	50.6	19.4	193	3800	Chinstrap
C	45.7	17.0	195	3650	Chinstrap
D	43.4	14.4	218	4600	Gentoo
E	44.5	14.7	214	4850	Gentoo
F	35.9	19.2	189	3800	Adelie
G	36.0	17.9	190	3450	Adelie
H	50.0	15.2	218	5700	Gentoo

- (i) Use the supervised learning model kNN (k=3) with the Euclidean distance to predict the species of a penguin with a bill length of 38.9 mm, a bill depth of 17.8 mm, a flipper length of 181 mm and a body mass of 3625 g. You may round the distance to 2 decimal places.

(9 marks)

Solution. By calculating the Euclidean distance from the point (46.5, 14.8, 217, 5200) to the points in the training data, we can obtain the following table:

Obs.	bill length	bill depth	flipper length	body mass	species	Dist
A	41.1	19.1	188	4100	Adelie	475.0584
B	50.6	19.4	193	3800	Chinstrap	175.8080
C	45.7	17.0	195	3650	Chinstrap	29.4598
D	43.4	14.4	218	4600	Gentoo	975.7181
E	44.5	14.7	214	4850	Gentoo	1225.4611
F	35.9	19.2	189	3800	Adelie	175.2140
G	36.0	17.9	190	3450	Adelie	175.2553
H	50.0	15.2	218	5700	Gentoo	2075.3612

..... [8 marks]

The 3 nearest neighbours are observations C, F and G, which correspond to species Chinstrap, Adelie and Adelie. Therefore, the prediction of the species of the penguin is Adelie. [1 mark]

Average: 7.45 / 9 marks in Jan 2024; 17.86% below 4.5 marks. □

- (ii) Based on your calculation in part (i), explain the problem with predictive modelling and what data preparation step is required to resolve the problem. (2 marks)

Solution. Problem: the data are not scaled — “body mass” has a large variation compare to other features [1 mark]

Solution: introduce min-max scaling or standardisation to bring all the features to similar variations. [1 mark]

Average: 0.32 / 2 marks in Jan 2024; 82.14% below 4.5 marks. □

More Performance Evaluation

5. (Final Exam May 2025 Sem, Q1(b)) Given a time-based data (e.g. econometric data), state the data splitting method we use to estimate the generalisation error of a predictive model. Hence determine if k-fold cross validation can be used with justification. (2 marks)

Solution. The time-based data needs to be sorted according to the timestamps and the holdout method with a cut-off time will be used to split the data into training data (historical data) and testing/validation data (future data). [1 mark]

K-fold cross validation is not suitable because “future data” are used in training the predic-

tive models.[1 mark]

Remark: This is related to the **Smarket** data discussed in practical 3. ☐

6. (Jan 2022 Final Q1(d)) Explain the steps in (i) validation set approach and (ii) k -fold cross validation and state each advantages and disadvantages. (5 marks)

Solution. (i) Validation set approach shuffles the data and splits it into training set and validation/test set. The training dataset is the sample of data used to fit the model; the validation dataset is the sample of data used to provide an unbiased evaluation of a model fit on the training dataset while tuning model hyperparameters. The evaluation becomes more biased as skew on the validation dataset is incorporated into the model configuration. The validation dataset may also play a role in other forms of model preparation, such as feature selection.[1 mark]

The advantage of this approach is its simplicity in scoring predictive models.[1 mark]

The disadvantage of this approach is its biasness and being too dependent on a particular sampling.[0.5 mark]

(ii) k -fold cross validation shuffles the data and splits it into k groups. Each group will be set as validation/test set and the remainder will be set as training set and be used to score the predictive models.[1 mark]

Advantage: This approach is less prone to biasness problem and less dependent on a particular sampling. It is used to tune model hyperparameters instead of a separate validation dataset.[1 mark]

Disadvantage: When the data size is large and the k is large, the scoring process can be very time consuming.[0.5 mark] ☐

7. What are the advantages of k -fold cross validation relative to

- (a) Validation set approach

Solution. The estimate of the test error rate can be highly variable depending on which observations are included in the training and validation sets.

Secondly, the validation set error rate may tend to overestimate the test error rate for the model fit on the entire data set, which is the overfitting problem. ☐

- (b) Leave-one-out cross validation (LOOCV)

Solution. LOOCV is a special case of k -fold cross-validation with $k = n$. Thus, LOOCV is the most computationally intense method since the model must be fit n times. In addition, LOOCV has **higher variance**, but **lower bias**, than k -fold cross validation. ☐

8. (May 2019 Final Q3)

- (a) Supervised learning includes classification and regression.

- (i) State the difference between classification and regression in term of response variable. (1 mark)

Solution. The response variable for classification is categorical while for regression is numerical. ☐

- (ii) Explain the sampling methods used in splitting data for classification and regression respectively. (4 marks)

Solution. The **classification** uses **stratified sampling** (sklearn's **StratifiedShuffleSplit**). Samples are distributed to different sets according to the proportion of response variable. The **regression** uses **linear sampling** (sklearn's **ShuffleSplit**). Sam-

ples are distributed randomly to different sets. □

- (b) (i) State an issue that comes along with split validation, which can be overcome by using cross validation. (1 mark)

Solution. Overfitting. □

- (ii) Describe the process of a 5-fold cross validation. (4 marks)

Solution. 5-fold cross validation randomly sampled observations into 5 non-overlapping groups with equal size, known as folds. For first iteration, first fold will be treated as validation set, and the remaining four folds act as training set. Five iterations will be run with a different fold is treated as validation set at each iteration, while the other folds served as training set. This will eventually give five estimates of accuracy measures and average will be taken. □

- (c) A sample of 500 males and 800 females had been collected to test on a model of gender prediction. The model resulted that 380 males and 510 females were predicted correctly.

- (i) Assume male as positive class and female as negative class, calculate the count of true positive (TP), true negative (TN), false positive (FP) and false negative (FN) for the model's result. (2 marks)

Solution. Given $TP = 380$; $TN = 510$
 $FN = 500 - TP = 120$; $FP = 800 - TN = 290$. □

- (ii) Construct the confusion matrix for the model. State the classification error, specificity and sensitivity of the model. (4 marks)

Solution. Confusion matrix:

	True +	True -	Precision
Predicted +	380	290	0.5672
Predicted -	120	510	0.8095
Recall	0.7600	0.6375	0.6846

Classification error = $1 - 0.6846 = 0.3154$

Specificity = negative recall = 0.6375

Sensitivity = positive recall = 0.76 □

- (iii) Compare the recall and precision for both male and female. Interpret your results. (4 marks)

Solution. Male (positive class): Recall = 0.7600; Precision = 0.5672.

High recall low precision. The model is wide but generalised in predicting male. This means that the model can capture most males but those predicted males might be incorrect.

Female (negative class): Recall = 0.6375; Precision = 0.8095.

Low recall high precision. The model is small but highly specialised in predicting female. This means that the model can only capture some females but those predicted females are mostly correct. □