

Lab Sheet

Nearest Neighbour Classifiers

0. Work through the “02 k-NN” notebook.
1. Three examples are shown below from the “penguins” dataset. Each example is represented by a vector of 4 numeric features. Example x_1 has been manually labelled as belonging to “Class

Example: x_1		Example: x_2		Query: q	
<i>Bill length</i>	4.4	<i>Bill length</i>	5.6	<i>Bill length</i>	6.1
<i>Bill depth</i>	2.9	<i>Bill depth</i>	3.0	<i>Bill depth</i>	3.0
<i>Flipper length</i>	1.4	<i>Flipper length</i>	4.5	<i>Flipper length</i>	4.6
<i>Body Mass</i>	0.2	<i>Body Mass</i>	1.5	<i>Body Mass</i>	1.4
<i>Class</i>	A	<i>Class</i>	B	<i>Class</i>	???

A”, while Example x_2 has been labelled as belonging to “Class B”.

- a) What type of distance function might be appropriate for comparing the examples above?
- b) Use this distance function to calculate the distances between the query example q and the two labelled examples. Which class label would a 1-NN classifier assign to the query based on the distances?

2. The table below shows three examples from a system for predicting whether a person is over or under the drink driving limit. The 5 input features for this system are:

- Gender: categorical feature {male, female}
- Weight: numeric, with range [50,150]
- Amount of alcohol in units: numeric, with range [1,16]
- Meal type: ordinal feature {None, Snack, Lunch, Full}
- Duration of drinking session: numeric, range [20,230]

Example: x_1		Example: x_2		Query: q	
Gender	female	Gender	male	Gender	male
Weight	60	Weight	75	Weight	70
Amount	4	Amount	2	Amount	1
Meal	full	Meal	full	Meal	snack
Duration	90	Duration	60	Duration	30
Class	over	Class	under	Class	???

- Normalise all numeric features to the range [0,1]
- Propose an appropriate global distance function for comparing examples such as the above.
-
- Use your proposed distance function to calculate the distances between the query example q and the two labelled examples. Which class label would a 1-NN classifier assign to the query based on the distances?

3. The table below reports the pairwise distances between a set of 9 labelled training examples and a new query example q , for the system described in Question 2.

Example	Class	Distance to q
$x1$	over	1.5
$x2$	under	2.8
$x3$	over	1.8
$x4$	under	2.9
$x5$	under	2.2
$x6$	under	3.0
$x7$	under	2.4
$x8$	over	3.2
$x9$	over	3.6

- a) What class label would a 3-NN classifier assign to q ?
- b) What class label would a 4-NN classifier assign to q ?
- c) What class label would a weighted 4-NN classifier assign to q ?

4. Case-based Reasoning (CBR) is a reasoning approach that uses k-NN to retrieve the most similar examples to query cases and uses these to make decisions about the query case.

(For information on CBR see Aamodt and Plaza 2001 seminal paper available in Brightspace under the Reading Unit).

Two different examples from a CBR system for estimating the price of second-hand cars are shown in the tables below. Each example is described by 6 features.

Example: <i>x1</i>	
<i>Manufacturer</i>	Ford
<i>Model</i>	Fiesta
<i>Engine Size</i>	1,100
<i>Fuel</i>	Petrol
<i>Mileage</i>	65,000
<i>Condition</i>	Excellent
<i>Price</i>	€3,100

Example: <i>x2</i>	
<i>Manufacturer</i>	Citroen
<i>Model</i>	BX
<i>Engine Size</i>	1,800
<i>Fuel</i>	Diesel
<i>Mileage</i>	37,000
<i>Condition</i>	Fair
<i>Price</i>	€4,500

- a) Normalise all numeric features to the range [0,1]. Assume that the feature ranges are:
- Engine Size 1,000 to 3,000
 - Mileage 1,000 to 100,000
- b) Propose a suitable global distance function that might be used in a *k*-Nearest Neighbour case retrieval system for this data. Assume that "Condition" is an ordinal feature that has the possible values {Poor, Fair, Good, Excellent},
- c) Use the proposed global distance function to calculate the distance between the examples *x1* and *x2* above.

5. The data below shows households classified by how budget is allocated ('Household.csv').

The notebook '02 kNN Lab Sheet' contains code to load in this dataset. Add in code to classify the query example using 1-NN and Euclidean distance. In this example households are classified based on how budget is allocated, correlation would be a better measure of similarity. Modify this code so that correlation is used rather than Euclidean distance.

Household	Groceries	Education	Travel	Category
H1	2000	4000	500	C1
H2	3000	6000	1000	C1
H3	2000	2000	2000	C2
H4	3000	3000	3000	C2
query	2500	3500	2000	?

6. In the Data Normalisation example in the "02-kNN" Notebook replace the $N(0,1)$ scaler with a min-max scaler. Are there any differences?
7. Download the zip file '02-BYO kNN-Python focus'. It contains a notebook that takes you through building your own kNN classifier in Python with a significant focus on writing good Python code. Work through this notebook if you are interested.