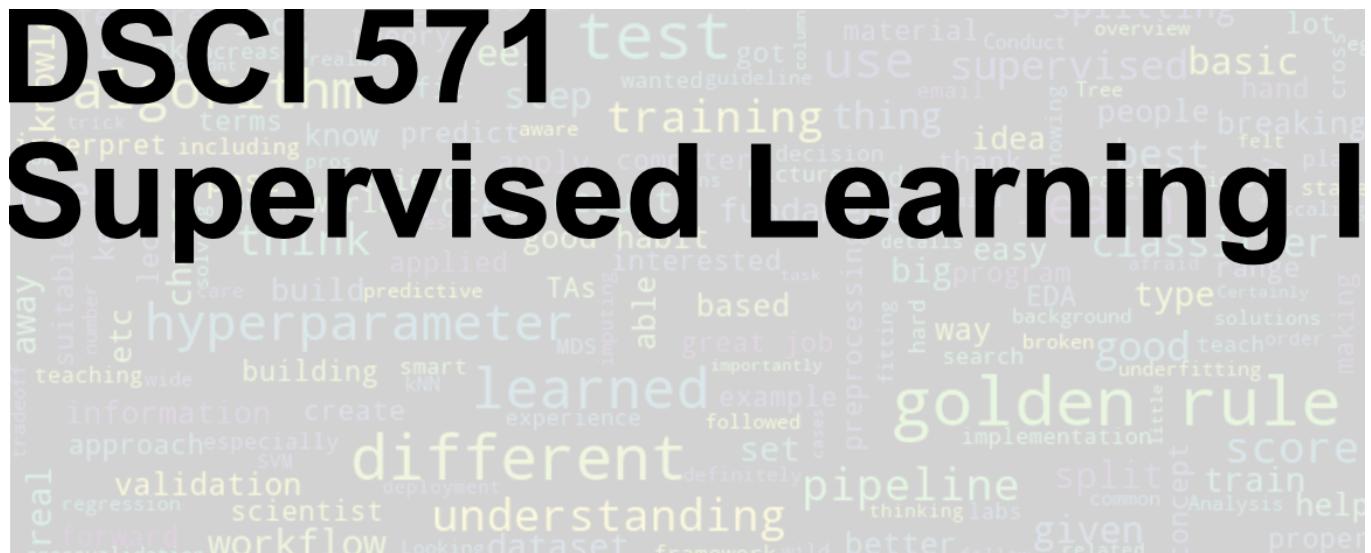


# Lecture 3: k-Nearest Neighbours and SVM RBFs

# Contents

- Imports and LOs
  - Motivation and distances [video]
  - $k$ -Nearest Neighbours ( $k$ -NNs) [video]
  - **?** **?** Questions for you
  - Break (5 min)
  - More on  $k$ -NNs [video]
  - Support Vector Machines (SVMs) with RBF kernel [video]
  - **?** **?** Questions for you
  - (iClicker) Exercise 3.2
  - Summary



UBC Master of Data Science program, 2024-25

If two things are similar, the thought of one will tend to trigger the thought of the other  
– Aristotle

# Imports and LOs

## Imports

```
import sys
import os

import IPython
import matplotlib.pyplot as plt
import numpy as np
import pandas as pd
from IPython.display import HTML

sys.path.append(os.path.join(os.path.abspath("."), "code"))

import ipywidgets as widgets
import mlearn
from IPython.display import display
from ipywidgets import interact, interactive
from plotting_functions import *
from sklearn.dummy import DummyClassifier
from sklearn.model_selection import cross_validate, train_test_split
from utils import *

%matplotlib inline

pd.set_option("display.max_colwidth", 200)
import warnings

warnings.filterwarnings("ignore")
DATA_DIR = DATA_DIR = os.path.join(os.path.abspath("."), "data/")
```

## Learning outcomes

From this lecture, you will be able to

- explain the notion of similarity-based algorithms;
- broadly describe how  $k$ -NNs use distances;
- discuss the effect of using a small/large value of the hyperparameter  $k$  when using the  $k$ -NN algorithm;
- describe the problem of curse of dimensionality;
- explain the general idea of SVMs with RBF kernel;
- broadly describe the relation of `gamma` and `C` hyperparameters of SVMs with the fundamental tradeoff.

## Quick recap

- Why do we split the data?
- What are the 4 types of data splits we discussed in the last lecture?
- What are the benefits of cross-validation?
- What is overfitting?
- What's the fundamental trade-off in supervised machine learning?
- What is the golden rule of machine learning?

### ! Important

If you want to run this notebook you will have to install `ipywidgets`. Follow the installation instructions [here](#).

## Motivation and distances [[video](#)]

## Analogy-based models

- Suppose you are given the following training examples with corresponding labels and are asked to label a given test example.

## Training examples: $X = \text{pictures}$ , $y = \text{names}$



Test example



[source](#)

- An intuitive way to classify the test example is by finding the most "similar" example(s) from the training set and using that label for the test example.

## Analogy-based algorithms in practice

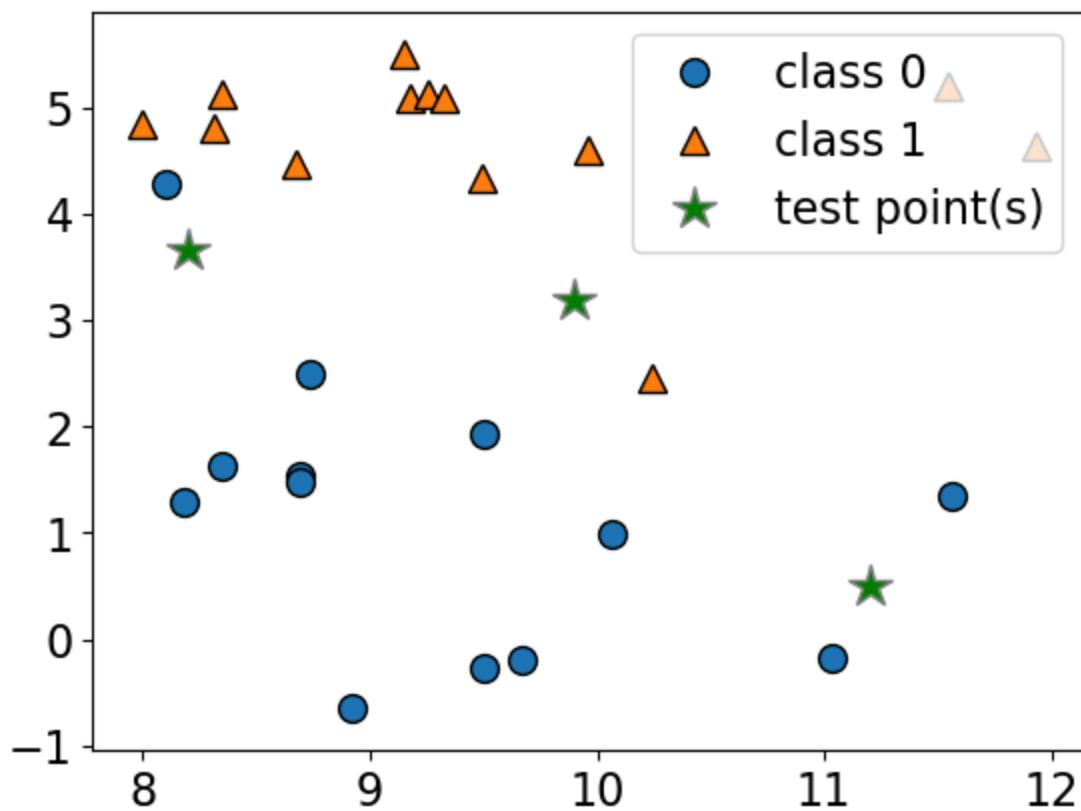
- [Herta's High-tech Facial Recognition](#)
  - Feature vectors for human faces
  - $k$ -NN to identify which face is on their watch list
- Recommendation systems

## General idea of $k$ -nearest neighbours algorithm

- Consider the following toy dataset with two classes.
  - blue circles → class 0
  - red triangles → class 1
  - green stars → test examples

```
X, y = mglearn.datasets.make_forge()
X_test = np.array([[8.2, 3.66214339], [9.9, 3.2], [11.2, 0.5]])
```

```
plot_train_test_points(X, y, X_test)
```



- Given a new data point, predict the class of the data point by finding the “closest” data point in the training set, i.e., by finding its “nearest neighbour” or majority vote of nearest neighbours.

```
import matplotlib
import panel as pn
from panel import widgets
from panel.interact import interact

pn.extension()
```

```

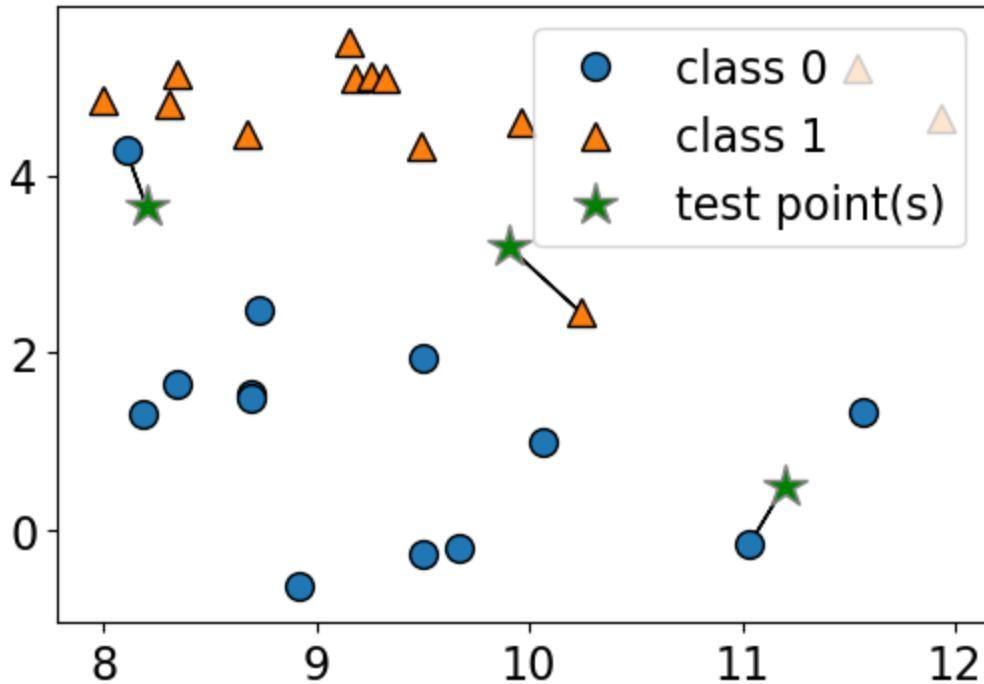
def f(n_neighbors):
    plt.clf()
    fig = plt.figure(figsize=(6, 4))
    plot_knn_clf(X, y, X_test, n_neighbors=n_neighbors)
    plt.close()
    return pn.pane.Matplotlib(fig, tight=True)

n_neighbors_selector = pn.widgets.IntSlider(
    name="n_neighbors", start=1, end=10, value=1
)
# interact(f, n_neighbors=n_neighbors_selector)
interactive_plot = interact(f, n_neighbors=n_neighbors_selector).embed(max_op
interactive_plot

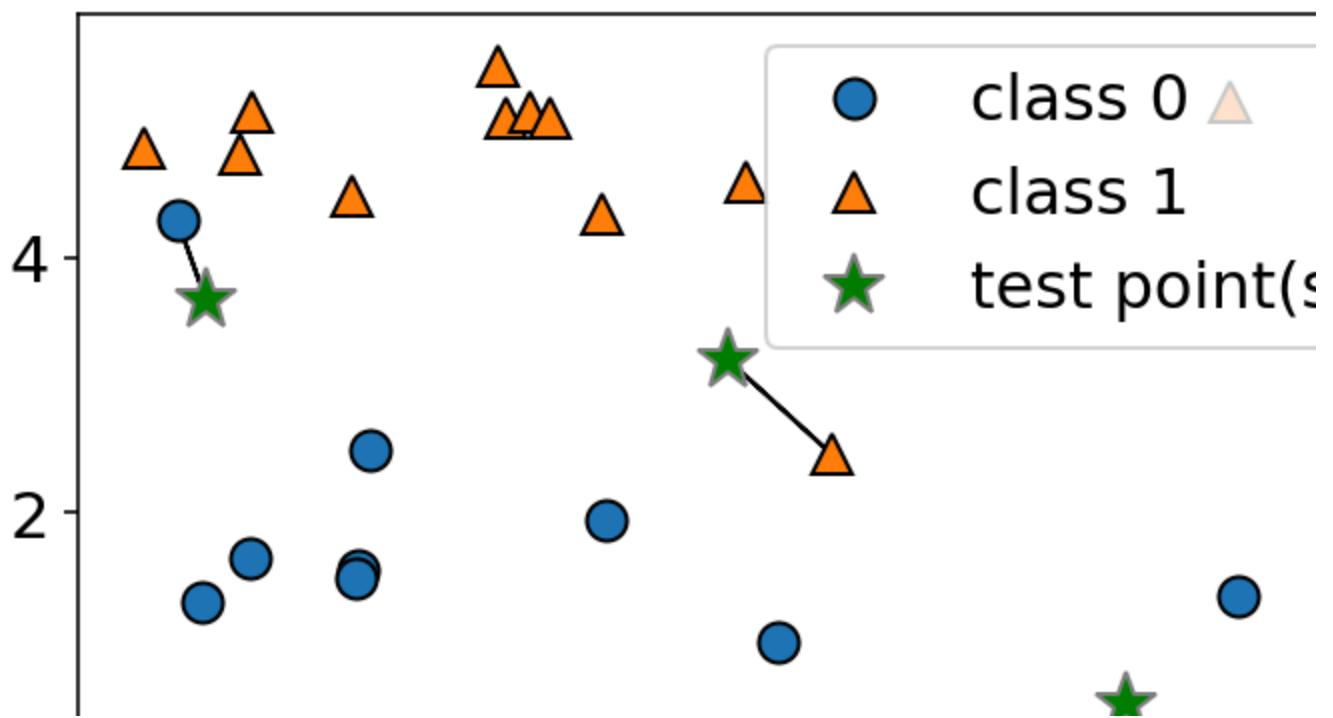
```

n\_neighbors 1

<Figure size 640x480 with 0 Axes>



n\_neighbors: 1



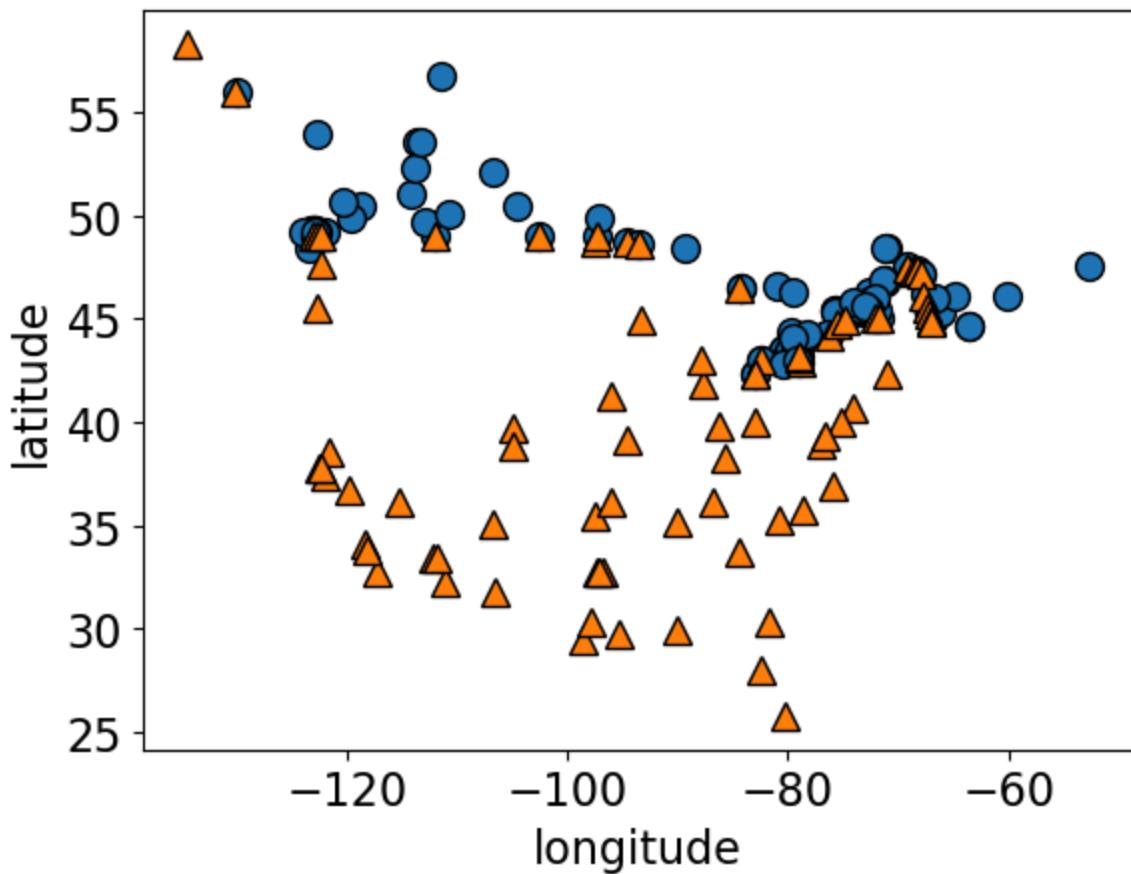
## Geometric view of tabular data and dimensions

- To understand analogy-based algorithms it's useful to think of data as points in a high dimensional space.
- Our  represents the problem in terms of relevant **features** ( $d$ ) with one dimension for each **feature** (column).
- Examples are **points in a  $d$ -dimensional space**.

How many dimensions (features) are there in the cities data?

```
cities_df = pd.read_csv(DATA_DIR + "canada_usa_cities.csv")
X_cities = cities_df[["longitude", "latitude"]]
y_cities = cities_df["country"]
```

```
mglearn.discrete_scatter(X_cities.iloc[:, 0], X_cities.iloc[:, 1], y_cities)
plt.xlabel("longitude")
plt.ylabel("latitude");
```



- Recall the [Spotify Song Attributes](#) dataset from homework 2.
- How many dimensions (features) we used in the homework?

```
spotify_df = pd.read_csv(DATA_DIR + "spotify.csv", index_col=0)
X_spotify = spotify_df.drop(columns=["target", "song_title", "artist"])
print("The number of features in the Spotify dataset: %d" % X_spotify.shape[1])
X_spotify.head()
```

The number of features in the Spotify dataset: 13

	acousticness	danceability	duration_ms	energy	instrumentalness	key	live
0	0.0102	0.833	204600	0.434	0.021900	2	0.
1	0.1990	0.743	326933	0.359	0.006110	1	0.
2	0.0344	0.838	185707	0.412	0.000234	2	0.
3	0.6040	0.494	199413	0.338	0.510000	5	0.
4	0.1800	0.678	392893	0.561	0.512000	5	0.

# Dimensions in ML problems

In ML, usually we deal with high dimensional problems where examples are hard to visualize.

- $d \approx 20$  is considered low dimensional
- $d \approx 1000$  is considered medium dimensional
- $d \approx 100,000$  is considered high dimensional

## Feature vectors

### Feature vector

is composed of feature values associated with an example.

Some example feature vectors are shown below.

```
print(
    "An example feature vector from the cities dataset: %s"
    % (X_cities.iloc[0].to_numpy())
)
print(
    "An example feature vector from the Spotify dataset: \n%s"
    % (X_spotify.iloc[0].to_numpy())
)
```

```
An example feature vector from the cities dataset: [-130.0437  55.9773]
An example feature vector from the Spotify dataset:
[ 1.02000e-02  8.33000e-01  2.04600e+05  4.34000e-01  2.19000e-02
 2.00000e+00  1.65000e-01 -8.79500e+00  1.00000e+00  4.31000e-01
 1.50062e+02  4.00000e+00  2.86000e-01]
```

## Similarity between examples

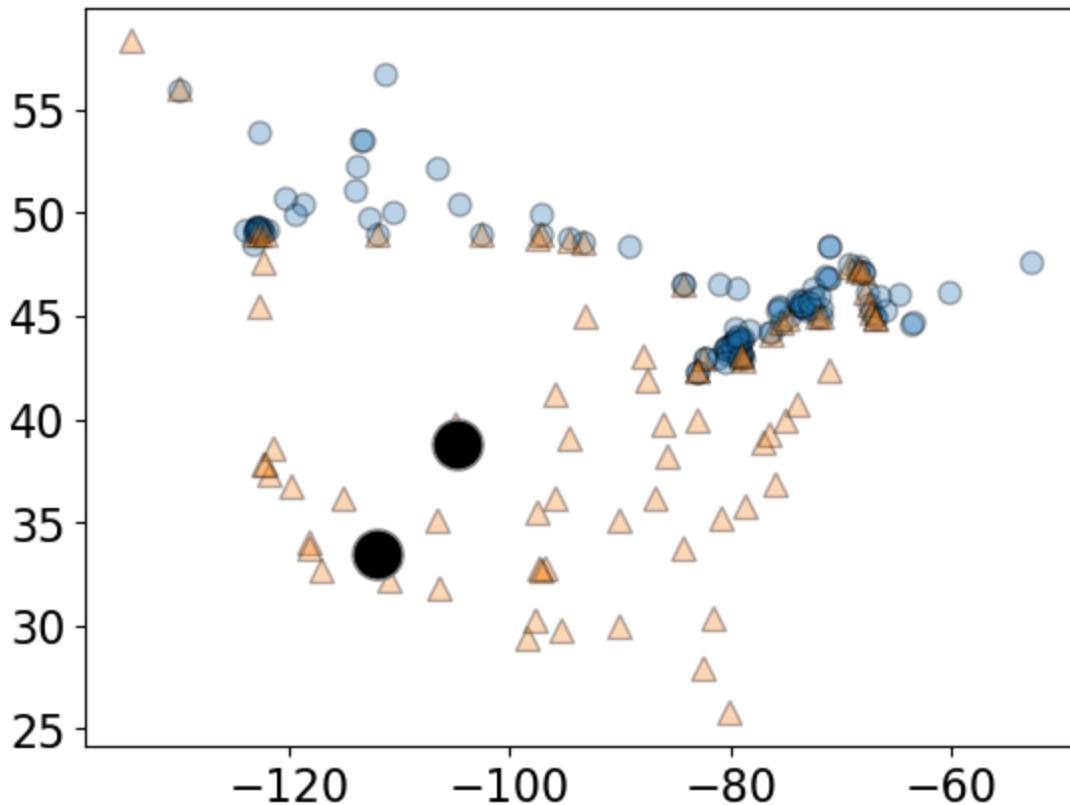
Let's take 2 points (two feature vectors) from the cities dataset.

```
two_cities = X_cities.sample(2, random_state=120)
two_cities
```

	longitude	latitude
69	-104.8253	38.8340
35	-112.0741	33.4484

The two sampled points are shown as big black circles.

```
mglearn.discrete_scatter(
    X_cities.iloc[:, 0], X_cities.iloc[:, 1], y_cities, s=8, alpha=0.3
)
mglearn.discrete_scatter(
    two_cities.iloc[:, 0], two_cities.iloc[:, 1], markers="o", c="k", s=18
);
```



## Distance between feature vectors

- For the cities at the two big circles, what is the *distance* between them?
- A common way to calculate the distance between vectors is calculating the **Euclidean distance**.
- The euclidean distance between vectors  $u = \langle u_1, u_2, \dots, u_n \rangle$  and  $v = \langle v_1, v_2, \dots, v_n \rangle$  is defined as:

$$\text{distance}(u, v) = \sqrt{\sum_{i=1}^n (u_i - v_i)^2}$$

## Euclidean distance

`two_cities`

	longitude	latitude
69	-104.8253	38.8340
35	-112.0741	33.4484

- Subtract the two cities
- Square the difference
- Sum them up
- Take the square root

```
# Subtract the two cities
print("Subtract the cities: \n%s\n" % (two_cities.iloc[1] - two_cities.iloc[0]))

# Squared sum of the difference
print(
    "Sum of squares: %0.4f" % (np.sum((two_cities.iloc[1] - two_cities.iloc[0]) ** 2)))

# Take the square root
print(
    "Euclidean distance between cities: %0.4f"
    % (np.sqrt(np.sum((two_cities.iloc[1] - two_cities.iloc[0]) ** 2))))
```

```
Subtract the cities:
longitude   -7.2488
latitude    -5.3856
dtype: float64

Sum of squares: 81.5498
Euclidean distance between cities: 9.0305
```

two\_cities

	longitude	latitude
69	-104.8253	38.8340
35	-112.0741	33.4484

```
# Euclidean distance using sklearn
from sklearn.metrics.pairwise import euclidean_distances

euclidean_distances(two_cities)
```

```
array([[0.          , 9.03049217],
       [9.03049217, 0.        ]])
```

Note: [scikit-learn](#) supports a number of other [distance metrics](#).

## Finding the nearest neighbour

- Let's look at distances from all cities to all other cities

```
dists = euclidean_distances(X_cities)
np.fill_diagonal(dists, np.inf)
dists.shape
```

```
(209, 209)
```

```
pd.DataFrame(dists)
```

	<b>0</b>	<b>1</b>	<b>2</b>	<b>3</b>	<b>4</b>	<b>5</b>	
<b>0</b>	inf	4.955113	9.869531	10.106452	10.449666	19.381676	28.3
<b>1</b>	4.955113	inf	14.677579	14.935802	15.305346	24.308448	33.2
<b>2</b>	9.869531	14.677579	inf	0.334411	0.808958	11.115406	20.5
<b>3</b>	10.106452	14.935802	0.334411	inf	0.474552	10.781004	20.1
<b>4</b>	10.449666	15.305346	0.808958	0.474552	inf	10.306500	19.7
...	...	...	...	...	...	...	...
<b>204</b>	58.289799	63.032656	50.483751	50.150395	49.677629	39.405415	30.0
<b>205</b>	64.183423	68.887343	56.512897	56.179123	55.705696	45.418031	36.0
<b>206</b>	52.426410	57.253724	44.235152	43.904226	43.435186	33.258427	24.0
<b>207</b>	58.033459	62.771969	50.249720	49.916254	49.443317	39.167214	29.7
<b>208</b>	51.498562	56.252160	43.699224	43.365623	42.892477	32.612755	23.2

209 rows × 209 columns

Let's look at the distances between City 0 and some other cities.

```
print("Feature vector for city 0: \n%s\n" % (X_cities.iloc[0]))
print("Distances from city 0 to the first 5 cities: %s" % (dists[0][:5]))
# We can find the closest city with `np.argmin`:
print(
    "The closest city from city 0 is: %d \n\nwith feature vector: \n%s"
    % (np.argmin(dists[0]), X_cities.iloc[np.argmin(dists[0])]))
)
```

```
Feature vector for city 0:
longitude -130.0437
latitude 55.9773
Name: 0, dtype: float64
```

```
Distances from city 0 to the first 5 cities: [ inf  4.95511263  9.869531]
```

The closest city from city 0 is: 81

```
with feature vector:
longitude -129.9912
latitude 55.9383
Name: 81, dtype: float64
```

Ok, so the closest city to City 0 is City 81.

# Question

- Why did we set the diagonal entries to infinity before finding the closest city?

## Finding the distances to a query point

We can also find the distances to a new "test" or "query" city:

```
# Let's find a city that's closest to the a query city
query_point = [[-80, 25]]

dists = euclidean_distances(X_cities, query_point)
dists[0:10]
```

```
array([[58.85545875],
       [63.80062924],
       [49.30530902],
       [49.01473536],
       [48.60495488],
       [39.96834506],
       [32.92852376],
       [29.53520104],
       [29.52881619],
       [27.84679073]])
```

```
# The query point is closest to
print(
    "The query point %s is closest to the city with index %d and the distance
     % (query_point, np.argmin(dists), dists[np.argmin(dists)])
)
```

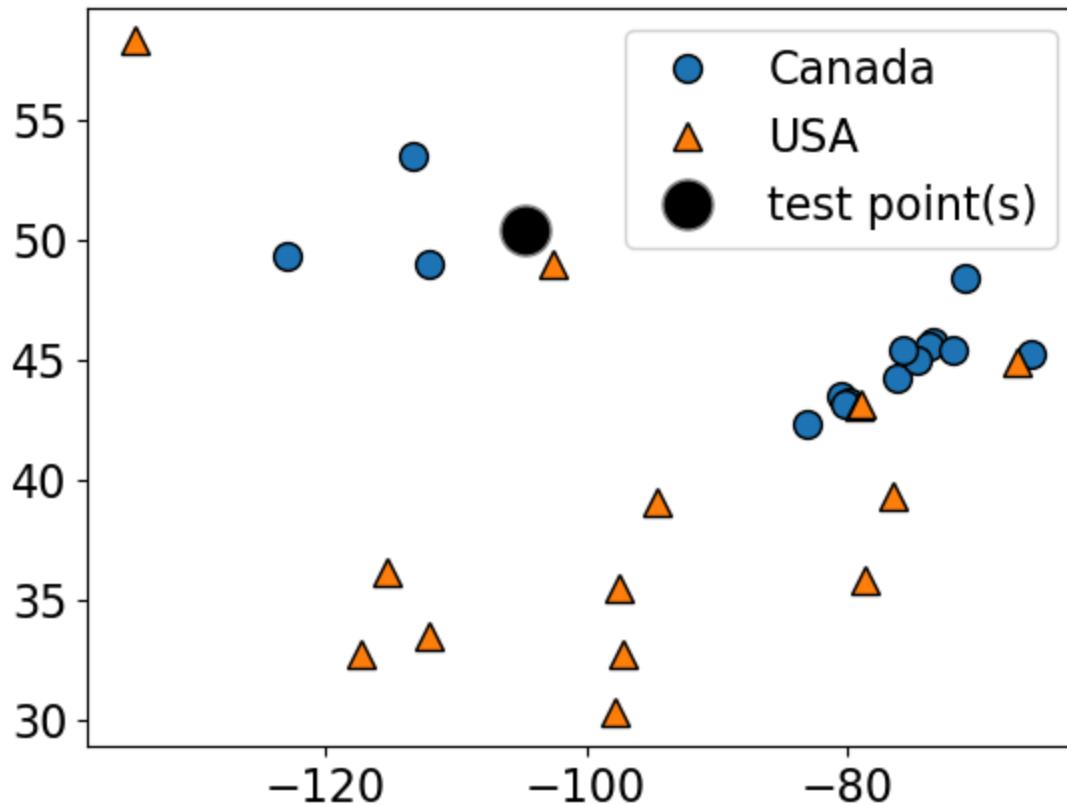
The query point `[-80, 25]` is closest to the city with index 72 and the distance

## $k$ -Nearest Neighbours ( $k$ -NNs) [[video](#)]

```
small_cities = cities_df.sample(30, random_state=90)
one_city = small_cities.sample(1, random_state=44)
small_train_df = pd.concat([small_cities, one_city]).drop_duplicates(keep=False)
```

```
X_small_cities = small_train_df.drop(columns=["country"]).to_numpy()
y_small_cities = small_train_df["country"].to_numpy()
test_point = one_city[["longitude", "latitude"]].to_numpy()
```

```
plot_train_test_points(
    X_small_cities,
    y_small_cities,
    test_point,
    class_names=["Canada", "USA"],
    test_format="circle",
)
```



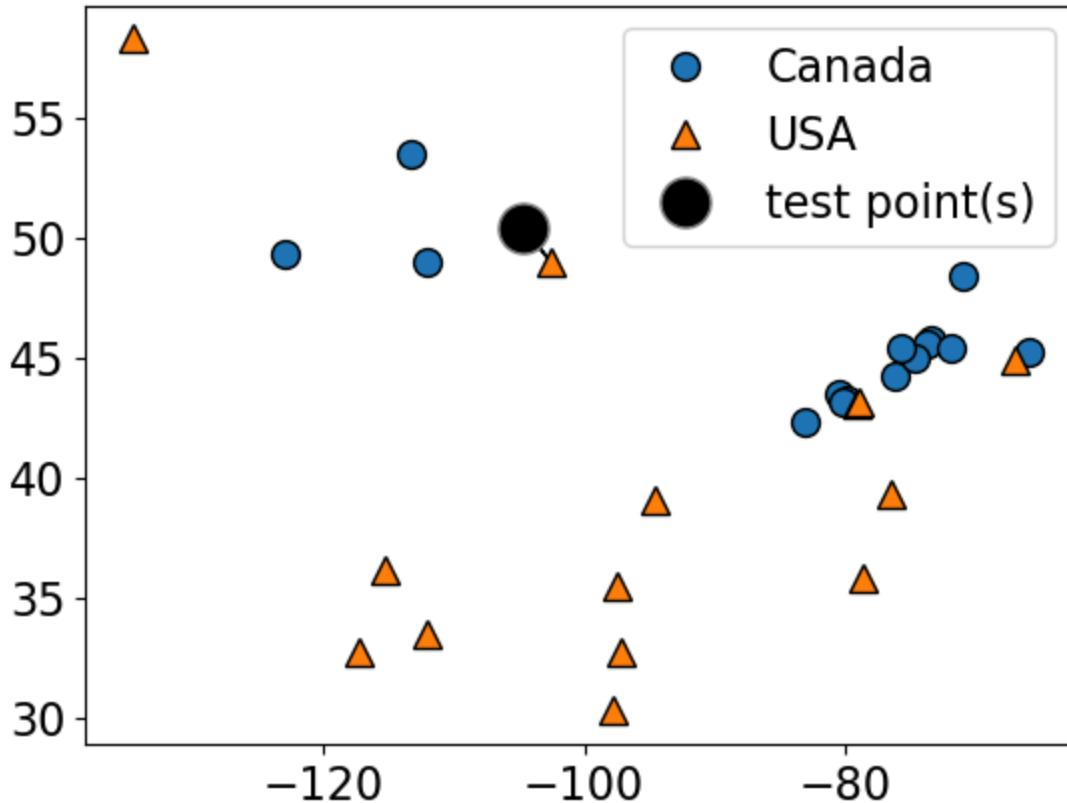
- Given a new data point, predict the class of the data point by finding the “closest” data point in the training set, i.e., by finding its “nearest neighbour” or majority vote of nearest neighbours.

Suppose we want to predict the class of the black point.

- An intuitive way to do this is predict the same label as the “closest” point ( $k = 1$ ) (1-nearest neighbour)
- We would predict a target of **USA** in this case.

```
plot_knn_clf(
    X_small_cities,
    y_small_cities,
    test_point,
    n_neighbors=1,
    class_names=["Canada", "USA"],
    test_format="circle",
)
```

n\_neighbors 1

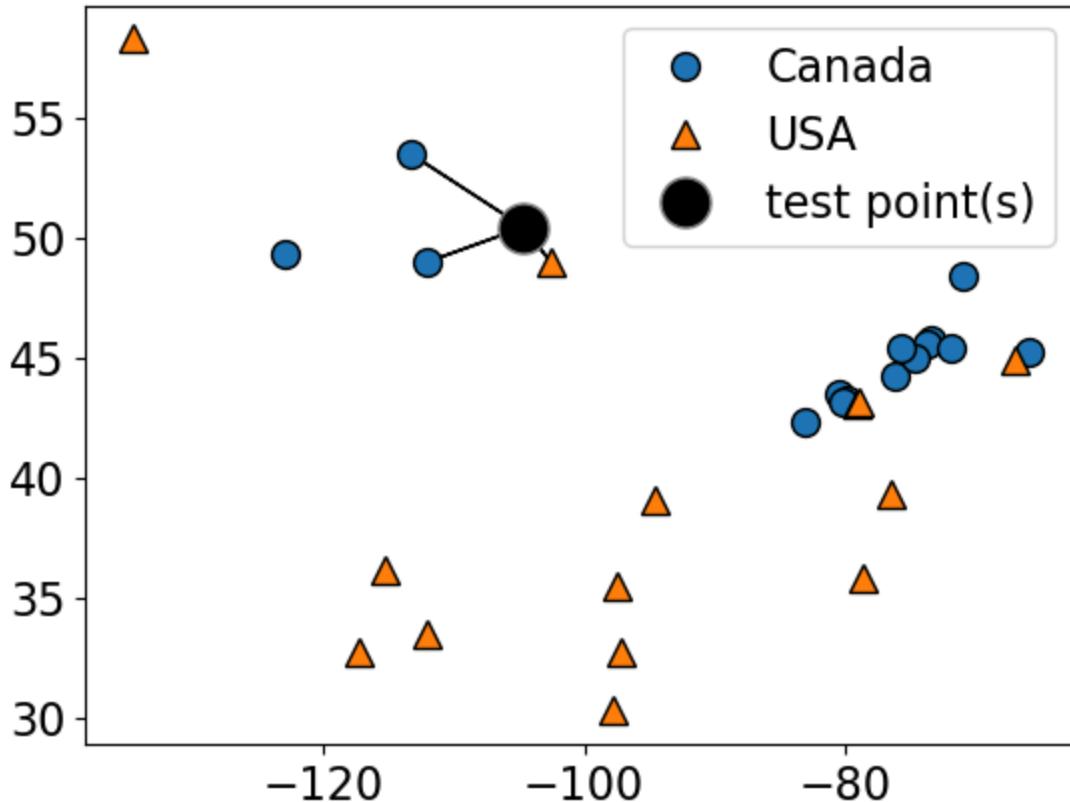


How about using  $k > 1$  to get a more robust estimate?

- For example, we could also use the 3 closest points ( $k = 3$ ) and let them **vote** on the correct class.
- The **Canada** class would win in this case.

```
plot_knn_clf(
    X_small_cities,
    y_small_cities,
    test_point,
    n_neighbors=3,
    class_names=["Canada", "USA"],
    test_format="circle",
)
```

n\_neighbors 3



```
from sklearn.neighbors import KNeighborsClassifier

k_values = [1, 3]

for k in k_values:
    neigh = KNeighborsClassifier(n_neighbors=k)
    neigh.fit(X_small_cities, y_small_cities)
    print(
        "Prediction of the black dot with %d neighbours: %s"
        % (k, neigh.predict(test_point))
    )
```

```
Prediction of the black dot with 1 neighbours: ['USA']
Prediction of the black dot with 3 neighbours: ['Canada']
```

## Choosing `n_neighbors`

- The primary hyperparameter of the model is `n_neighbors` ( $k$ ) which decides how many neighbours should vote during prediction?
- What happens when we play around with `n_neighbors`?
- Are we more likely to overfit with a low `n_neighbors` or a high `n_neighbors`?
- Let's examine the effect of the hyperparameter on our cities data.

```
X = cities_df.drop(columns=["country"])
y = cities_df["country"]

# split into train and test sets
X_train, X_test, y_train, y_test = train_test_split(
    X, y, test_size=0.1, random_state=123
)
```

```
k = 1
knn1 = KNeighborsClassifier(n_neighbors=k)
scores = cross_validate(knn1, X_train, y_train, return_train_score=True)
pd.DataFrame(scores)
```

	<code>fit_time</code>	<code>score_time</code>	<code>test_score</code>	<code>train_score</code>
0	0.002434	0.002767	0.710526	1.0
1	0.000675	0.001394	0.684211	1.0
2	0.000574	0.001393	0.842105	1.0
3	0.000551	0.001191	0.702703	1.0
4	0.000519	0.001220	0.837838	1.0

```
k = 100
knn100 = KNeighborsClassifier(n_neighbors=k)
scores = cross_validate(knn100, X_train, y_train, return_train_score=True)
pd.DataFrame(scores)
```

	<b>fit_time</b>	<b>score_time</b>	<b>test_score</b>	<b>train_score</b>
<b>0</b>	0.000635	0.033191	0.605263	0.600000
<b>1</b>	0.000721	0.001778	0.605263	0.600000
<b>2</b>	0.000579	0.001665	0.605263	0.600000
<b>3</b>	0.000563	0.002264	0.594595	0.602649
<b>4</b>	0.000527	0.001836	0.594595	0.602649

```
plot_knn_decision_boundaries(X_train, y_train, k_values=[1, 11, 100])
```

```
-----
AttributeError                                Traceback (most recent call last)
File ~/miniforge3/envs/jbook/lib/python3.12/site-packages/mglearn/plot_2d_sepa
  85     try:
--> 86         decision_values = classifier.decision_function(X_grid)
  87         levels = [0] if threshold is None else [threshold]
```

`AttributeError: 'KNeighborsClassifier' object has no attribute 'decision_func:`

During handling of the above exception, another exception occurred:

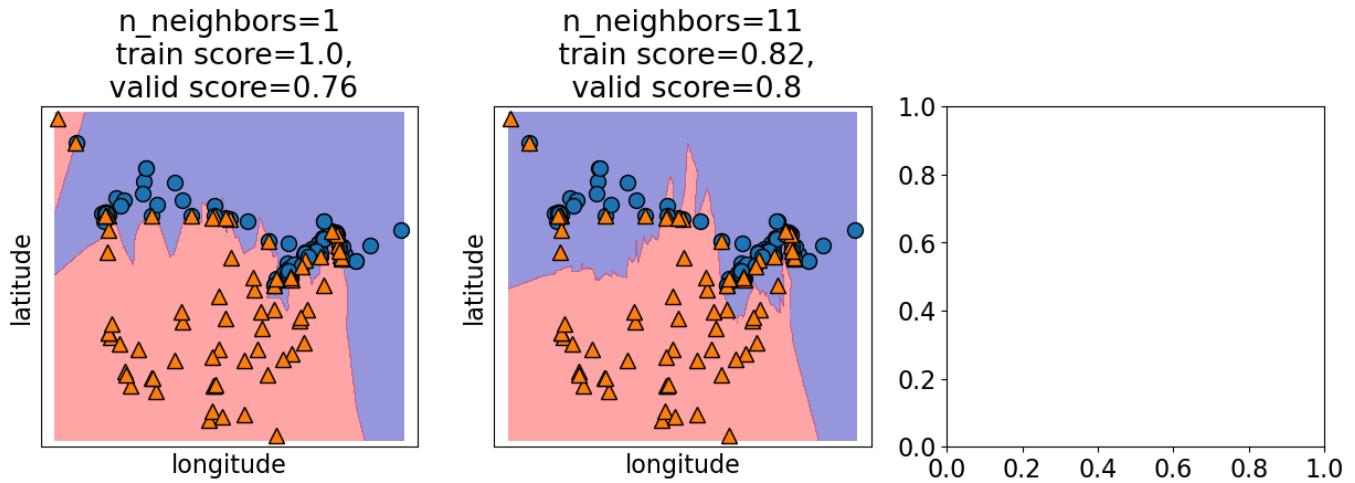
```
KeyboardInterrupt                         Traceback (most recent call last)
Cell In[30], line 1
----> 1 plot_knn_decision_boundaries(X_train, y_train, k_values=[1, 11, 100])

File ~/MDS/2024-25/571/DSCI_571_sup-learn-1_students/lectures/code/plotting_fu
  202 mean_train_score = scores["train_score"].mean()
  203 clf.fit(X_train, y_train)
--> 204 mglearn.plots.plot_2d_separator(
  205     clf, X_train.to_numpy(), fill=True, eps=0.5, ax=ax, alpha=0.4
  206 )
  207 mglearn.discrete_scatter(X_train.iloc[:, 0], X_train.iloc[:, 1], y_train)
  208 title = "n_neighbors={}\n train score={},\n valid score={}".format(
  209     n_neighbors, round(mean_train_score, 2), round(mean_valid_score, 2
  210 )
```

```
File ~/miniforge3/envs/jbook/lib/python3.12/site-packages/mglearn/plot_2d_sepa
  88     fill_levels = [decision_values.min()] + levels + [
  89         decision_values.max()]
  90 except AttributeError:
  91     # no decision_function
--> 92     decision_values = classifier.predict_proba(X_grid)[:, 1]
  93     levels = [.5] if threshold is None else [threshold]
  94     fill_levels = [0] + levels + [1]
```

```
File ~/miniforge3/envs/jbook/lib/python3.12/site-packages/sklearn/neighbors/_c
  395 # a simple ':' index doesn't work right
  396 for i, idx in enumerate(pred_labels.T): # loop is O(n_neighbors)
--> 397     proba_k[all_rows, idx] += weights[:, i]
  399 # normalize 'votes' into real [0,1] probabilities
  400 normalizer = proba_k.sum(axis=1)[:, np.newaxis]
```

`KeyboardInterrupt:`



## How to choose `n_neighbors`?

- `n_neighbors` is a hyperparameter
- We can use hyperparameter optimization to choose `n_neighbors`.

```
results_dict = {
    "n_neighbors": [],
    "mean_train_score": [],
    "mean_cv_score": [],
    "std_cv_score": [],
    "std_train_score": [],
}
param_grid = {"n_neighbors": np.arange(1, 50, 5)}

for k in param_grid["n_neighbors"]:
    knn = KNeighborsClassifier(n_neighbors=k)
    scores = cross_validate(knn, X_train, y_train, return_train_score=True)
    results_dict["n_neighbors"].append(k)

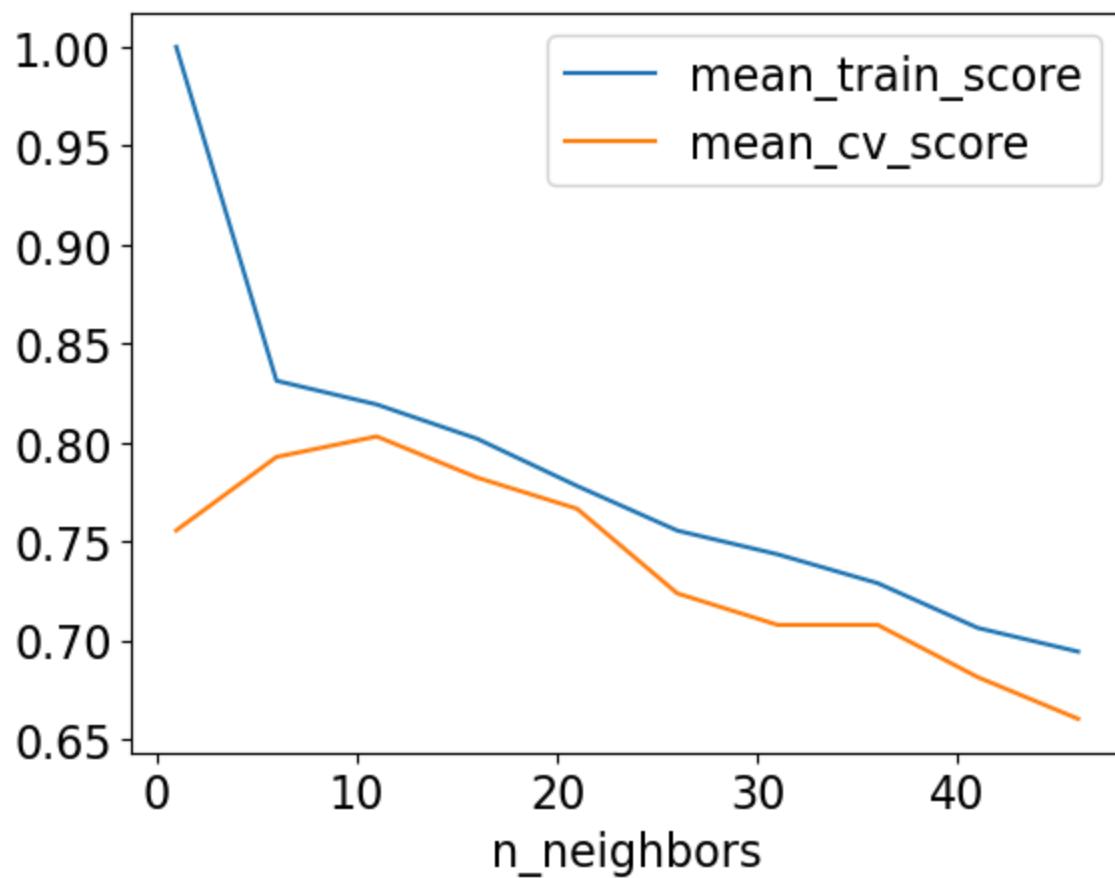
    results_dict["mean_cv_score"].append(np.mean(scores["test_score"]))
    results_dict["mean_train_score"].append(np.mean(scores["train_score"]))
    results_dict["std_cv_score"].append(scores["test_score"].std())
    results_dict["std_train_score"].append(scores["train_score"].std())

results_df = pd.DataFrame(results_dict)
```

```
results_df = results_df.set_index("n_neighbors")
results_df
```

	<b>mean_train_score</b>	<b>mean_cv_score</b>	<b>std_cv_score</b>	<b>std_train_score</b>
<b>n_neighbors</b>				
<b>1</b>	1.000000	0.755477	0.069530	0.000000
<b>6</b>	0.831135	0.792603	0.046020	0.013433
<b>11</b>	0.819152	0.802987	0.041129	0.011336
<b>16</b>	0.801863	0.782219	0.074141	0.008735
<b>21</b>	0.777934	0.766430	0.062792	0.016944
<b>26</b>	0.755364	0.723613	0.061937	0.025910
<b>31</b>	0.743391	0.707681	0.057646	0.030408
<b>36</b>	0.728777	0.707681	0.064452	0.021305
<b>41</b>	0.706128	0.681223	0.061241	0.018310
<b>46</b>	0.694155	0.660171	0.093390	0.018178

```
results_df[["mean_train_score", "mean_cv_score"]].plot();
```



```
best_n_neighbours = results_df.idxmax()["mean_cv_score"]
best_n_neighbours
```

11

Let's try our best model on test data.

```
knn = KNeighborsClassifier(n_neighbors=best_n_neighbours)
knn.fit(X_train, y_train)
print("Test accuracy: %0.3f" % (knn.score(X_test, y_test)))
```

Test accuracy: 0.905

Seems like we got lucky with the test set here.

## ?? Questions for you

### (iClicker) Exercise 3.1

Select all of the following statements which are TRUE.

- (A) Analogy-based models find examples from the test set that are most similar to the query example we are predicting.
- (B) Euclidean distance will always have a non-negative value.
- (C) With  $k$ -NN, setting the hyperparameter  $k$  to larger values typically reduces training error.
- (D) Similar to decision trees,  $k$ -NNs finds a small set of good features.
- (E) In  $k$ -NN, with  $k > 1$ , the classification of the closest neighbour to the test example always contributes the most to the prediction.

💡 V's Solutions!



## Break (5 min)



More on  $k$ -NNs [[video](#)]

Other useful arguments of `KNeighborsClassifier`

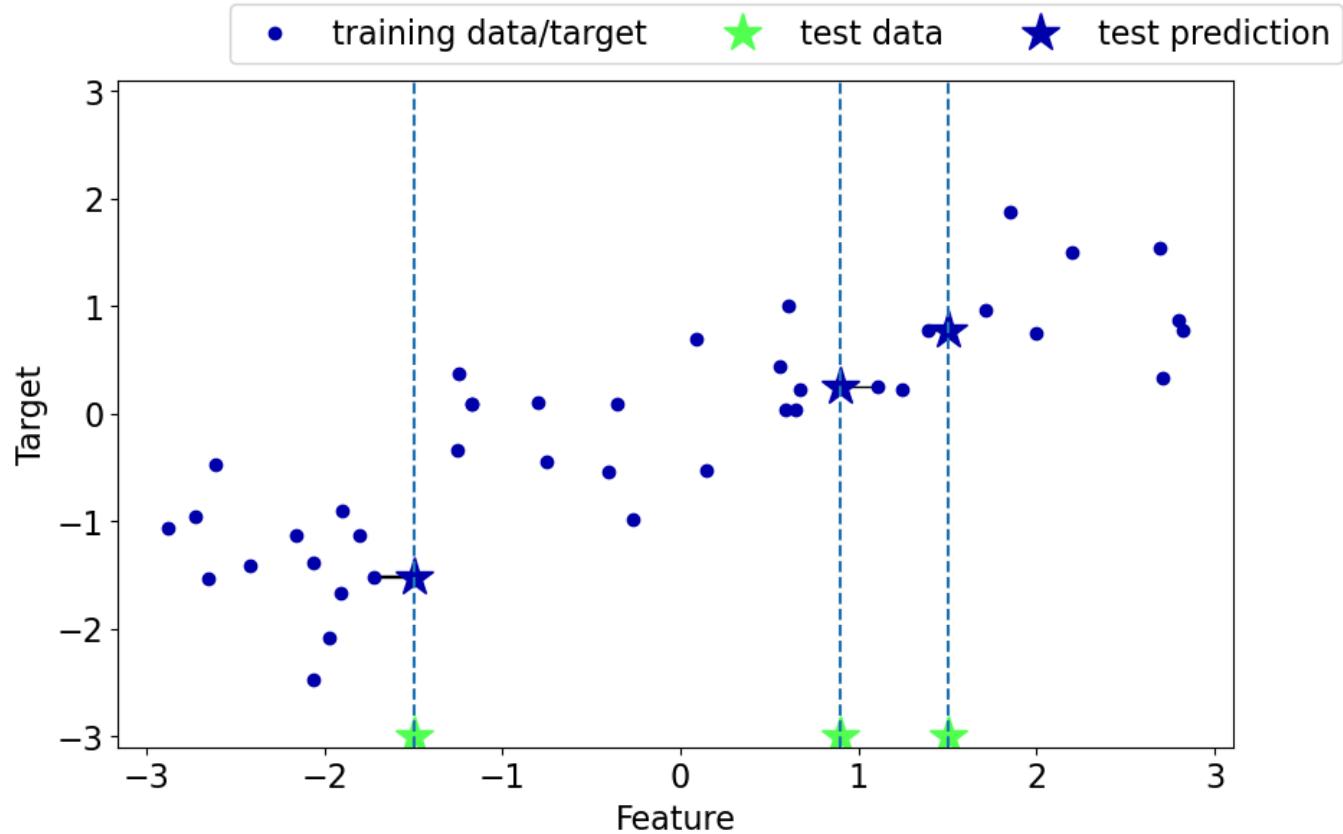
- `weights` → When predicting label, you can assign higher weight to the examples which are closer to the query example.
- Exercise for you: Play around with this argument. Do you get a better validation score?

## Regression with $k$ -nearest neighbours ( $k$ -NNs)

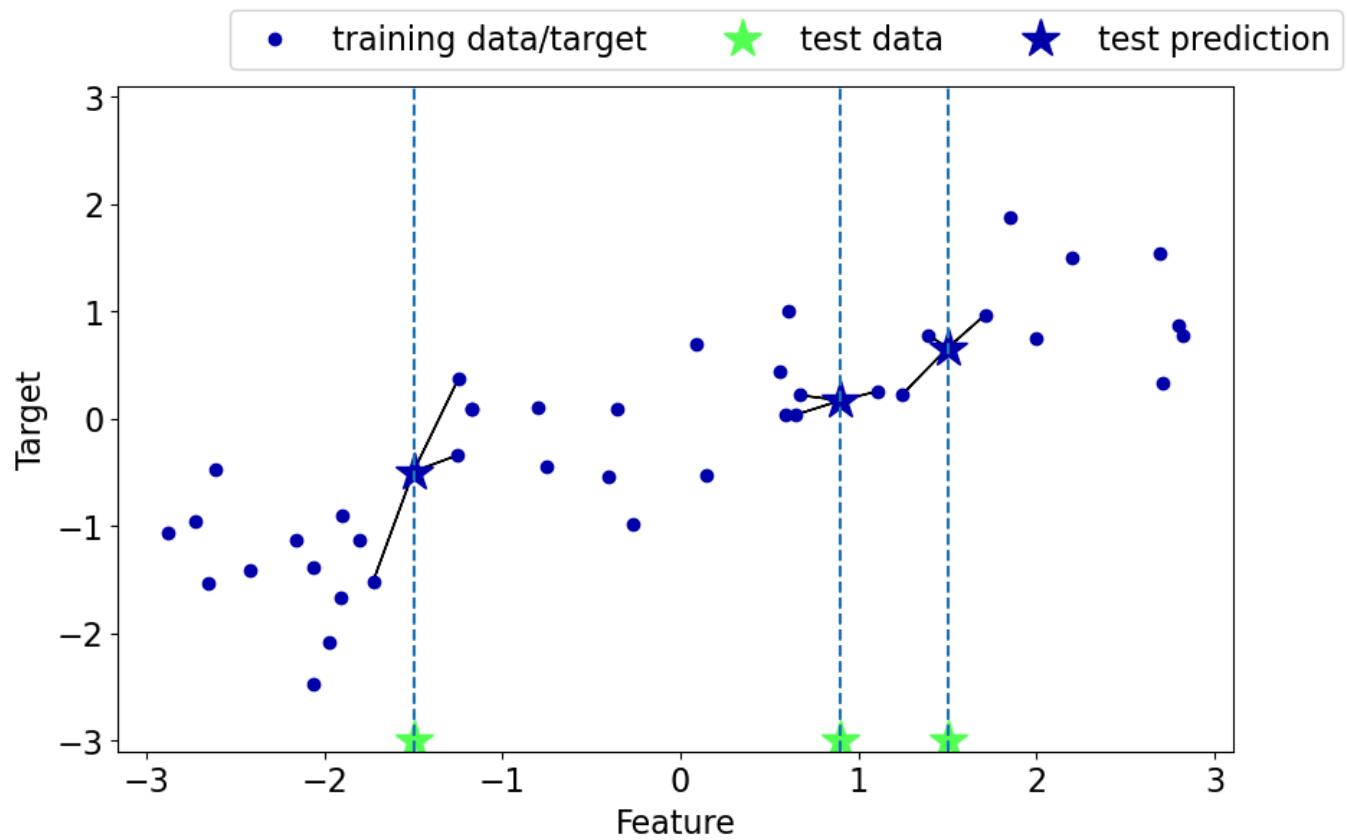
- Can we solve regression problems with  $k$ -nearest neighbours algorithm?
- In  $k$ -NN regression we take the average of the  $k$ -nearest neighbours.
- We can also have weighted regression.

See an example of regression in the lecture notes.

```
mglearn.plots.plot_knn_regression(n_neighbors=1)
```



```
mglearn.plots.plot_knn_regression(n_neighbors=3)
```



## Pros of $k$ -NNs for supervised learning

- Easy to understand, interpret.
- Simple hyperparameter  $k$  (`n_neighbors`) controlling the fundamental tradeoff.
- Can learn very complex functions given enough data.
- Lazy learning: Takes no time to `fit`

## Cons of $k$ -NNs for supervised learning

- Can be potentially be VERY slow during prediction time, especially when the training set is very large.
- Often not that great test accuracy compared to the modern approaches.
- It does not work well on datasets with many features or where most feature values are 0 most of the time (sparse datasets).

### ⚠️ Attention

For regular  $k$ -NN for supervised learning (not with sparse matrices), you should scale your features. We'll be looking into it soon.

## (Optional) Parametric vs non parametric

- You might see a lot of definitions of these terms.
- A simple way to think about this is:
  - do you need to store at least  $O(n)$  worth of stuff to make predictions? If so, it's non-parametric.
- Non-parametric example:  $k$ -NN is a classic example of non-parametric models.
- Parametric example: decision stump
- If you want to know more about this, find some reading material [here](#), [here](#), and [here](#).
- By the way, the terms "parametric" and "non-parametric" are often used differently by statisticians, see [here](#) for more...

### ℹ️ Note

$\mathcal{O}(n)$  is referred to as big  $\mathcal{O}$  notation. It tells you how fast an algorithm is or how much storage space it requires. For example, in simple terms, if you have  $n$  examples and you need to store them all you can say that the algorithm requires  $\mathcal{O}(n)$  worth of stuff.

## Curse of dimensionality

- Affects all learners but especially bad for nearest-neighbour.
- $k$ -NN usually works well when the number of dimensions  $d$  is small but things fall apart quickly as  $d$  goes up.
- If there are many irrelevant attributes,  $k$ -NN is hopelessly confused because all of them contribute to finding similarity between examples.
- With enough irrelevant attributes the accidental similarity swamps out meaningful similarity and  $k$ -NN is no better than random guessing.

```
from sklearn.datasets import make_classification

nfeats_accuracy = {"nfeats": [], "dummy_valid_accuracy": [], "KNN_valid_accur
for n_feats in range(4, 2000, 100):
    X, y = make_classification(n_samples=2000, n_features=n_feats, n_classes=2
    X_train, X_test, y_train, y_test = train_test_split(
        X, y, test_size=0.2, random_state=123
    )
    dummy = DummyClassifier(strategy="most_frequent")
    dummy_scores = cross_validate(dummy, X_train, y_train, return_train_score=
        True)
    knn = KNeighborsClassifier()
    scores = cross_validate(knn, X_train, y_train, return_train_score=True)
    nfeats_accuracy["nfeats"].append(n_feats)
    nfeats_accuracy["KNN_valid_accuracy"].append(np.mean(scores["test_score"]))
    nfeats_accuracy["dummy_valid_accuracy"].append(np.mean(dummy_scores["test_
```

```
pd.DataFrame(nfeats_accuracy)
```

	<b>nfeats</b>	<b>dummy_valid_accuracy</b>	<b>KNN_valid_accuracy</b>
<b>0</b>	4	0.503750	0.867500
<b>1</b>	104	0.501875	0.755625
<b>2</b>	204	0.506250	0.695625
<b>3</b>	304	0.500000	0.614375
<b>4</b>	404	0.505625	0.671250
<b>5</b>	504	0.510625	0.623750
<b>6</b>	604	0.505625	0.648125
<b>7</b>	704	0.511250	0.615000
<b>8</b>	804	0.511250	0.604375
<b>9</b>	904	0.506875	0.553125
<b>10</b>	1004	0.500625	0.609375
<b>11</b>	1104	0.502500	0.564375
<b>12</b>	1204	0.504375	0.638750
<b>13</b>	1304	0.505000	0.620625
<b>14</b>	1404	0.501250	0.615000
<b>15</b>	1504	0.505625	0.575625
<b>16</b>	1604	0.502500	0.587500
<b>17</b>	1704	0.501250	0.582500
<b>18</b>	1804	0.499375	0.560000
<b>19</b>	1904	0.513125	0.591875

## Support Vector Machines (SVMs) with RBF kernel [[video](#)]

- Very high-level overview
- Our goals here are

- Use `scikit-learn`'s SVM model.
- Broadly explain the notion of support vectors.
- Broadly explain the similarities and differences between  $k$ -NNs and SVM RBFs.
- Explain how `C` and `gamma` hyperparameters control the fundamental tradeoff.

(Optional) RBF stands for radial basis functions. We won't go into what it means in this video. Refer to [this video](#) if you want to know more.

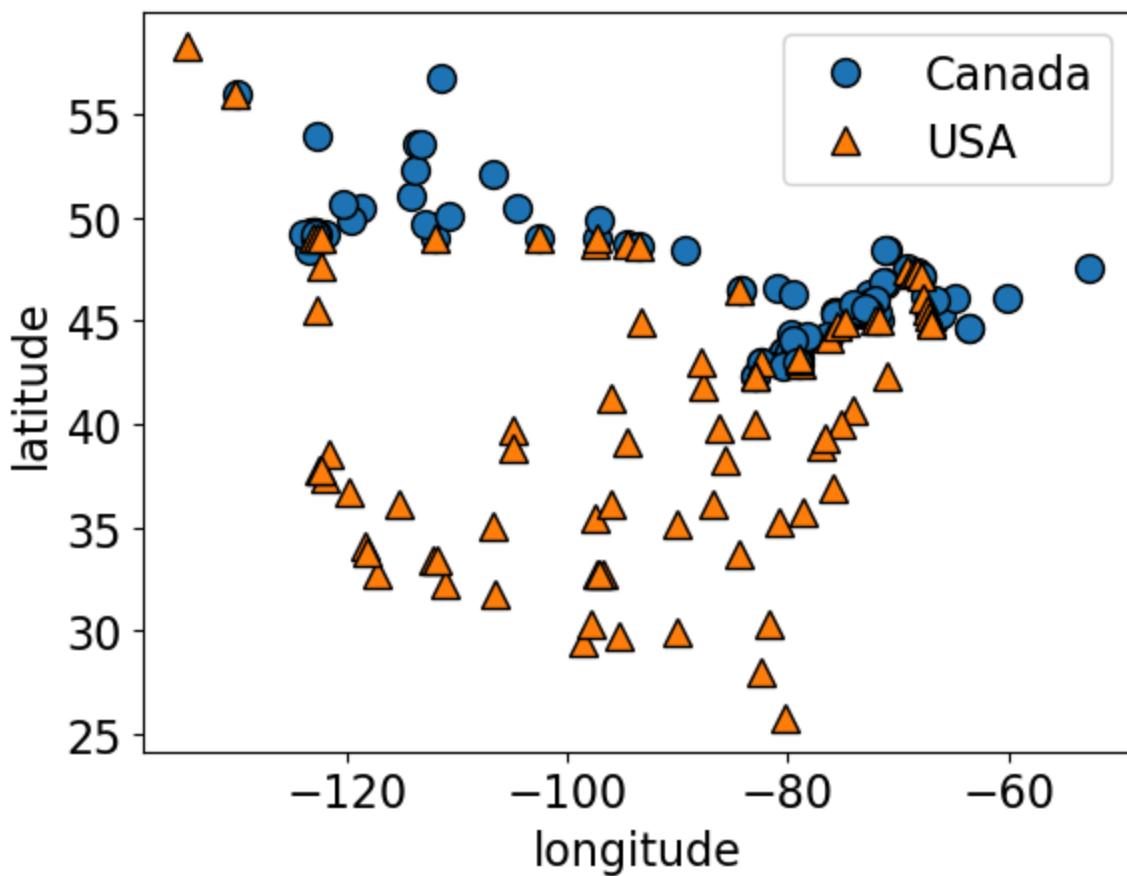
## Overview

- Another popular similarity-based algorithm is Support Vector Machines with RBF Kernel (SVM RBFs)
- Superficially, SVM RBFs are more like weighted  $k$ -NNs.
  - The decision boundary is defined by **a set of positive and negative examples and their weights** together with **their similarity measure**.
  - A test example is labeled positive if on average it looks more like positive examples than the negative examples.
- The primary difference between  $k$ -NNs and SVM RBFs is that
  - Unlike  $k$ -NNs, SVM RBFs only remember the key examples (support vectors).
  - SVMs use a different similarity metric which is called a "kernel". A popular kernel is Radial Basis Functions (RBFs)
  - They usually perform better than  $k$ -NNs!

## Let's explore SVM RBFs

Let's try SVMs on the cities dataset.

```
mglearn.discrete_scatter(X_cities.iloc[:, 0], X_cities.iloc[:, 1], y_cities)
plt.xlabel("longitude")
plt.ylabel("latitude")
plt.legend(loc=1);
```



```
X_train, X_test, y_train, y_test = train_test_split(
    X_cities, y_cities, test_size=0.2, random_state=123
)
```

```
knn = KNeighborsClassifier(n_neighbors=best_n_neighbours)
scores = cross_validate(knn, X_train, y_train, return_train_score=True)
print("Mean validation score %0.3f" % (np.mean(scores["test_score"])))
pd.DataFrame(scores)
```

Mean validation score 0.803

	fit_time	score_time	test_score	train_score
0	0.000830	0.002402	0.794118	0.819549
1	0.000653	0.001354	0.764706	0.819549
2	0.000587	0.001280	0.727273	0.850746
3	0.000568	0.001274	0.787879	0.828358
4	0.000557	0.001255	0.939394	0.783582

```
from sklearn.svm import SVC

svm = SVC(gamma=0.01) # Ignore gamma for now
scores = cross_validate(svm, X_train, y_train, return_train_score=True)
print("Mean validation score %0.3f" % (np.mean(scores["test_score"])))
pd.DataFrame(scores)
```

Mean validation score 0.820

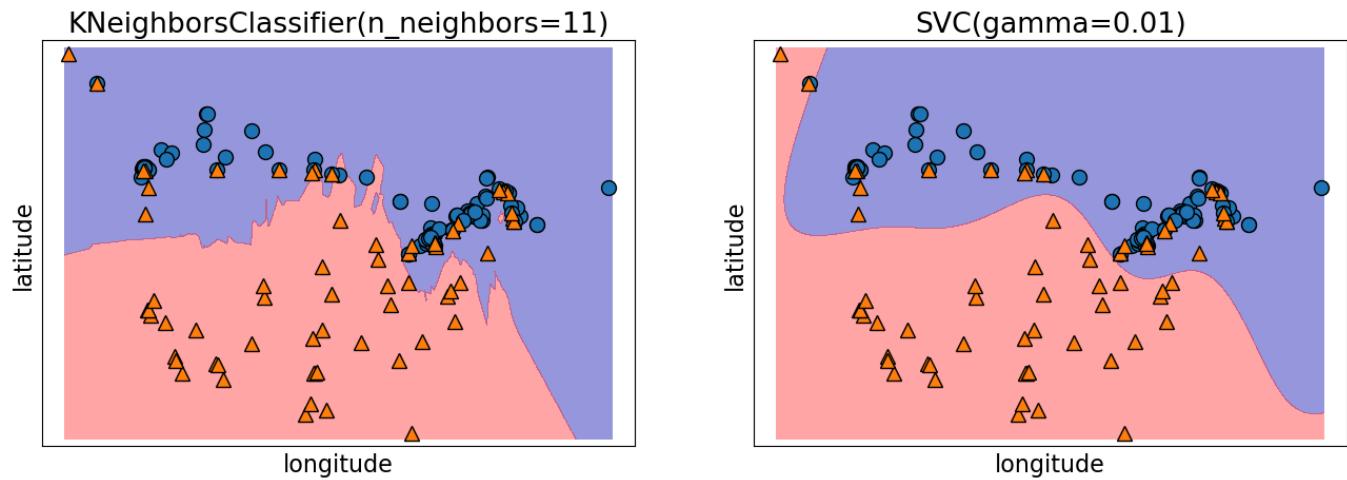
	fit_time	score_time	test_score	train_score
0	0.001008	0.000590	0.823529	0.842105
1	0.000803	0.000535	0.823529	0.842105
2	0.000761	0.000515	0.727273	0.858209
3	0.000848	0.000537	0.787879	0.843284
4	0.000749	0.000881	0.939394	0.805970

## Decision boundary of SVMs

- We can think of SVM with RBF kernel as “smooth KNN”.

```
fig, axes = plt.subplots(1, 2, figsize=(16, 5))

for clf, ax in zip([knn, svm], axes):
    clf.fit(X_train, y_train)
    mglearn.plots.plot_2d_separator(
        clf, X_train.to_numpy(), fill=True, eps=0.5, ax=ax, alpha=0.4
    )
    mglearn.discrete_scatter(X_train.iloc[:, 0], X_train.iloc[:, 1], y_train,
    ax.set_title(clf)
    ax.set_xlabel("longitude")
    ax.set_ylabel("latitude")
```



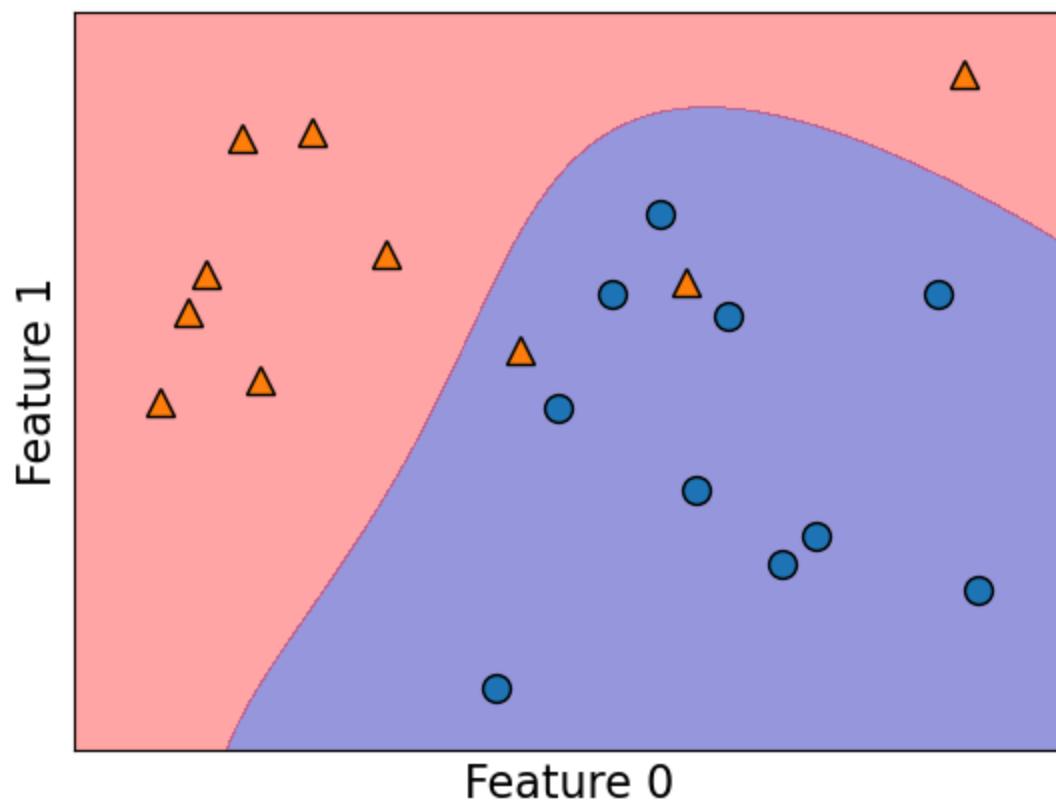
## Support vectors

- Each training example either is or isn't a "support vector".
  - This gets decided during `fit`.
- **Main insight: the decision boundary only depends on the support vectors.**
- Let's look at the support vectors.

```
from sklearn.datasets import make_blobs

n = 20
n_classes = 2
X_toy, y_toy = make_blobs(
    n_samples=n, centers=n_classes, random_state=300
) # Let's generate some fake data
```

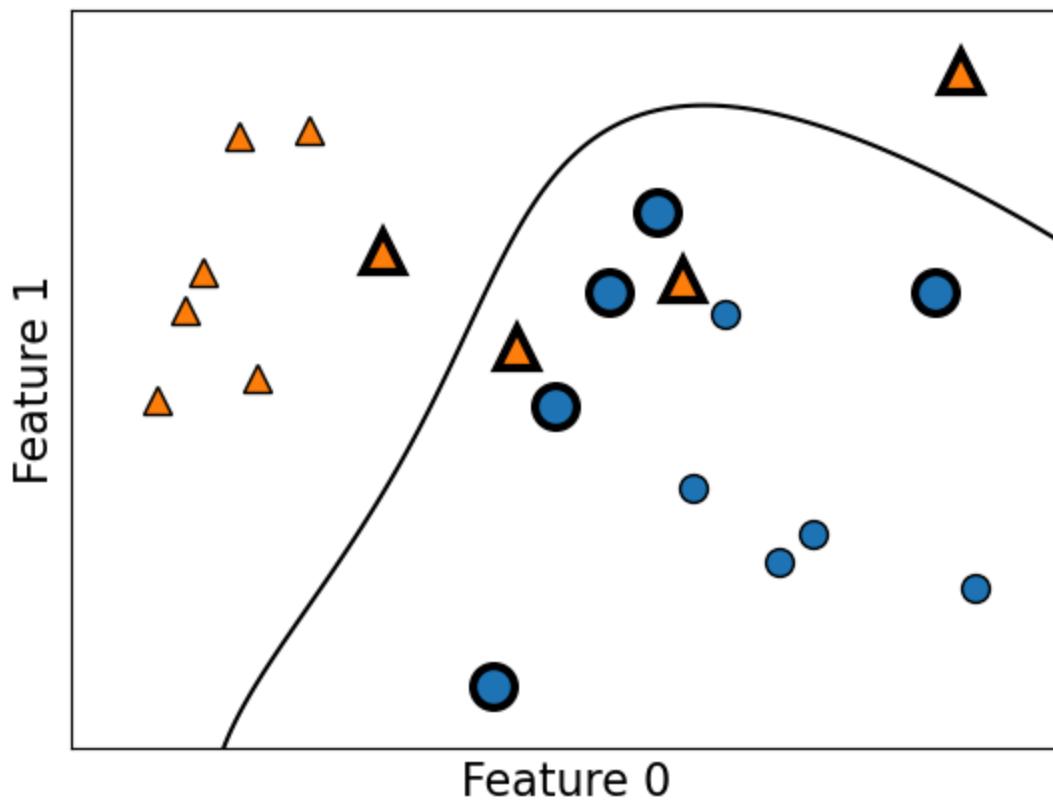
```
mglearn.discrete_scatter(X_toy[:, 0], X_toy[:, 1], y_toy)
plt.xlabel("Feature 0")
plt.ylabel("Feature 1")
svm = SVC(kernel="rbf", C=10, gamma=0.1).fit(X_toy, y_toy)
mglearn.plots.plot_2d_separator(svm, X_toy, fill=True, eps=0.5, alpha=0.4)
```



```
svm.support_
```

```
array([ 3,  8,  9, 14, 19,  1,  4,  6, 17], dtype=int32)
```

```
plot_support_vectors(svm, X_toy, y_toy)
```



The support vectors are the bigger points in the plot above.

## Hyperparameters of SVM

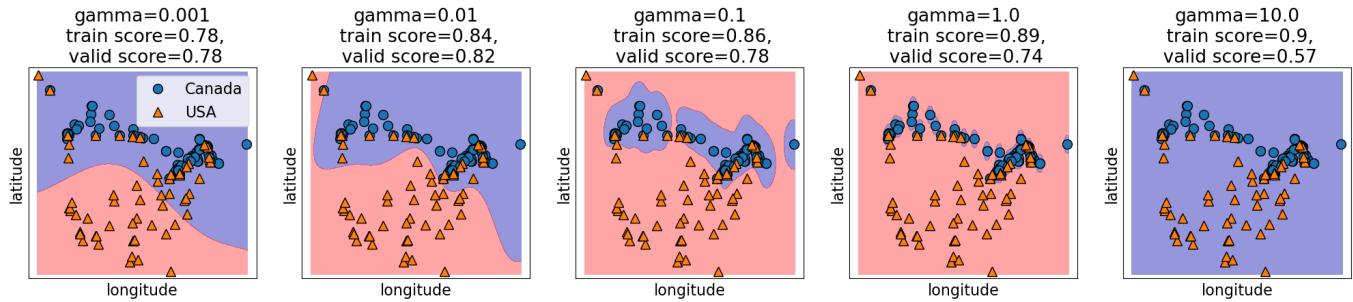
- Key hyperparameters of `rbf` SVM are
  - `gamma`
  - `C`
- We are not equipped to understand the meaning of these parameters at this point but you are expected to describe their relation to the fundamental tradeoff.

See [scikit-learn](#)'s explanation of RBF SVM parameters.

## Relation of `gamma` and the fundamental trade-off

- `gamma` controls the complexity (fundamental trade-off), just like other hyperparameters we've seen.
  - larger `gamma` → more complex
  - smaller `gamma` → less complex

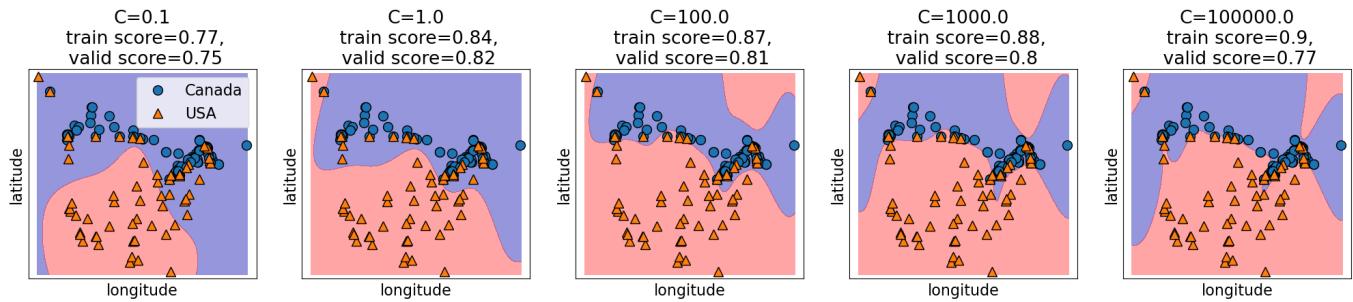
```
gamma = [0.001, 0.01, 0.1, 1.0, 10.0]
plot_svc_gamma(
    gamma,
    X_train.to_numpy(),
    y_train.to_numpy(),
    x_label="longitude",
    y_label="latitude",
)
```



## Relation of $C$ and the fundamental trade-off

- $C$  also affects the fundamental tradeoff
  - larger  $C$  → more complex
  - smaller  $C$  → less complex

```
C = [0.1, 1.0, 100.0, 1000.0, 100000.0]
plot_svc_C(
    C, X_train.to_numpy(), y_train.to_numpy(), x_label="longitude", y_label="l
)
```



## Search over multiple hyperparameters

- So far you have seen how to carry out search over a hyperparameter

- In the above case the best training error is achieved by the most complex model (large `gamma`, large `C`).
- Best validation error requires a hyperparameter search to balance the fundamental tradeoff.
  - In general we can't search them one at a time.
  - More on this next week. But if you cannot wait till then, you may look up the following:
    - [`sklearn.model\_selection.GridSearchCV`](#)
    - [`sklearn.model\_selection.RandomizedSearchCV`](#)

## SVM Regressor

- Similar to KNNs, you can use SVMs for regression problems as well.
- See [`sklearn.svm.SVR`](#) for more details.

## ?? Questions for you

### (iClicker) Exercise 3.2

Select all of the following statements which are TRUE.

- (A)  $k$ -NN may perform poorly in high-dimensional space (say,  $d > 1000$ ).
- (B) In sklearn's SVC classifier, large values of gamma tend to result in higher training score but probably lower validation score.
- (C) If we increase both gamma and C, we can't be certain if the model becomes more complex or less complex.



V's Solutions!



# Playground

In this interactive playground, you will investigate how various algorithms create decision boundaries to distinguish between Iris flower species using their sepal length and width as features. By adjusting the parameters, you can observe how the decision boundaries change, which can result in either overfitting (where the model fits the training data too closely) or underfitting (where the model is too simplistic).

- With **k-Nearest Neighbours (k-NN)**, you'll determine how many neighboring flowers to consult. Should we rely on a single nearest neighbor? Or should we consider a wider group?
- With **Support Vector Machine (SVM)** using the RBF kernel, you'll tweak the hyperparameters `C` and `gamma` to explore the tightrope walk between overly complex boundaries (that might overfit) and overly broad ones (that might underfit).
- With **Decision trees**, you'll observe the effect of `max_depth` on the decision boundary.

Observe the process of crafting and refining decision boundaries, one parameter at a time! Be sure to take breaks to reflect on the results you are observing.

```

from matplotlib.figure import Figure

from sklearn.model_selection import train_test_split
from sklearn.datasets import load_iris
from ipywidgets import interact, FloatLogSlider, IntSlider
import mglearn

# Load dataset and preprocessing
iris = load_iris(as_frame=True)
iris_df = iris.data
iris_df['species'] = iris.target
iris_df = iris_df[iris_df['species'] > 0]
X, y = iris_df[['sepal length (cm)', 'sepal width (cm)']], iris_df['species']
X_train, X_test, y_train, y_test = train_test_split(X, y, test_size=0.4, rando

# Common plot settings
def plot_results(model, X_train, y_train, title, ax):
    mglearn.plots.plot_2d_separator(model, X_train.values, fill=True, alpha=0.
    mglearn.discrete_scatter(
        X_train["sepal length (cm)", X_train["sepal width (cm)", y_train, s=
    );
    ax.set_xlabel("sepal length (cm)", fontsize=12);
    ax.set_ylabel("sepal width (cm)", fontsize=12);
    train_score = np.round(model.score(X_train.values, y_train), 2)
    test_score = np.round(model.score(X_test.values, y_test), 2)
    ax.set_title(
        f"{title}\n train score = {train_score}\ntest score = {test_score}", f
    );
    pass

# Widgets for SVM, k-NN, and Decision Tree
c_widget = pn.widgets.FloatSlider(
    value=1.0, start=1, end=5, step=0.1, name="C (log scale)"
)
gamma_widget = pn.widgets.FloatSlider(
    value=1.0, start=-3, end=5, step=0.1, name="Gamma (log scale)"
)
n_neighbors_widget = pn.widgets.IntSlider(
    start=1, end=40, step=1, value=5, name="n_neighbors"
)
max_depth_widget = pn.widgets.IntSlider(
    start=1, end=20, step=1, value=3, name="max_depth"
)

# The update function to create the plots
def update_plots(c, gamma=1.0, n_neighbors=5, max_depth=3):
    c_log = round(10**c, 2) # Transform C to logarithmic scale
    gamma_log = round(10**gamma, 2) # Transform Gamma to logarithmic scale

    fig = Figure(figsize=(8, 2))
    axes = fig.subplots(1, 3)

```

```

models = [
    SVC(C=c_log, gamma=gamma_log, random_state=42),
    KNeighborsClassifier(n_neighbors=n_neighbors),
    DecisionTreeClassifier(max_depth=max_depth, random_state=42),
]
titles = [
    f"SVM (C={c_log}, gamma={gamma_log})",
    f"k-NN (n_neighbors={n_neighbors})",
    f"Decision Tree (max_depth={max_depth})",
]
for model, title, ax in zip(models, titles, axes):
    model.fit(X_train.values, y_train)
    plot_results(model, X_train, y_train, title, ax);
# print(c, gamma, n_neighbors, max_depth)
return pn.pane.Matplotlib(fig, tight=True);

# Bind the function to the panel widgets
interactive_plot = pn.bind(
    update_plots,
    c=c_widget.param.value_throttled,
    gamma=gamma_widget.param.value_throttled,
    n_neighbors=n_neighbors_widget.param.value_throttled,
    max_depth=max_depth_widget.param.value_throttled,
)

# Layout the widgets and the plot
dashboard = pn.Column(
    pn.Row(c_widget, n_neighbors_widget),
    pn.Row(gamma_widget, max_depth_widget),
    interactive_plot,
)

# Display the interactive dashboard
dashboard

```

## Summary

- We have KNNs and SVMs as new supervised learning techniques in our toolbox.
- These are analogy-based learners and the idea is to assign nearby points the same label.
- Unlike decision trees, all features are equally important.

- Both can be used for classification or regression (much like the other methods we've seen).

## Coming up:

Lingering questions:

- Are we ready to do machine learning on real-world datasets?
- What would happen if we use  $k$ -NNs or SVM RBFs on the spotify dataset from hw2?
- What happens if we have missing values in our data?
- What do we do if we have features with categories or string values?

