

Week_11_Unsupervised_Learning

May 22, 2024

1 Clustering

Intro - Classification vs Clustering

```
[1]: # extra code - this cell generates figure below

import matplotlib.pyplot as plt
from sklearn.datasets import load_iris

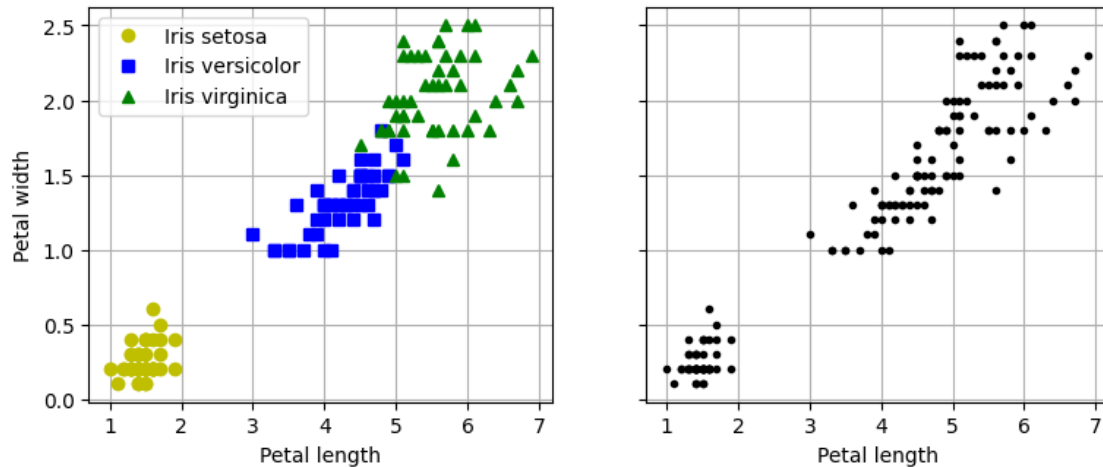
data = load_iris()
X = data.data
y = data.target
data.target_names

plt.figure(figsize=(9, 3.5))

plt.subplot(121)
plt.plot(X[y==0, 2], X[y==0, 3], "yo", label="Iris setosa")
plt.plot(X[y==1, 2], X[y==1, 3], "bs", label="Iris versicolor")
plt.plot(X[y==2, 2], X[y==2, 3], "g^", label="Iris virginica")
plt.xlabel("Petal length")
plt.ylabel("Petal width")
plt.grid()
plt.legend()

plt.subplot(122)
plt.scatter(X[:, 2], X[:, 3], c="k", marker=".")
plt.xlabel("Petal length")
plt.tick_params(labelleft=False)
plt.gca().set_axisbelow(True)
plt.grid()

plt.show()
```



Note: the next cell shows how a Gaussian mixture model (explained later in this topic) can actually separate these clusters pretty well using all 4 features: petal length & width, and sepal length & width. This code maps each cluster to a class. Instead of hard coding the mapping, the code picks the most common class for each cluster using the `scipy.stats.mode()` function:

```
[2]: # extra code

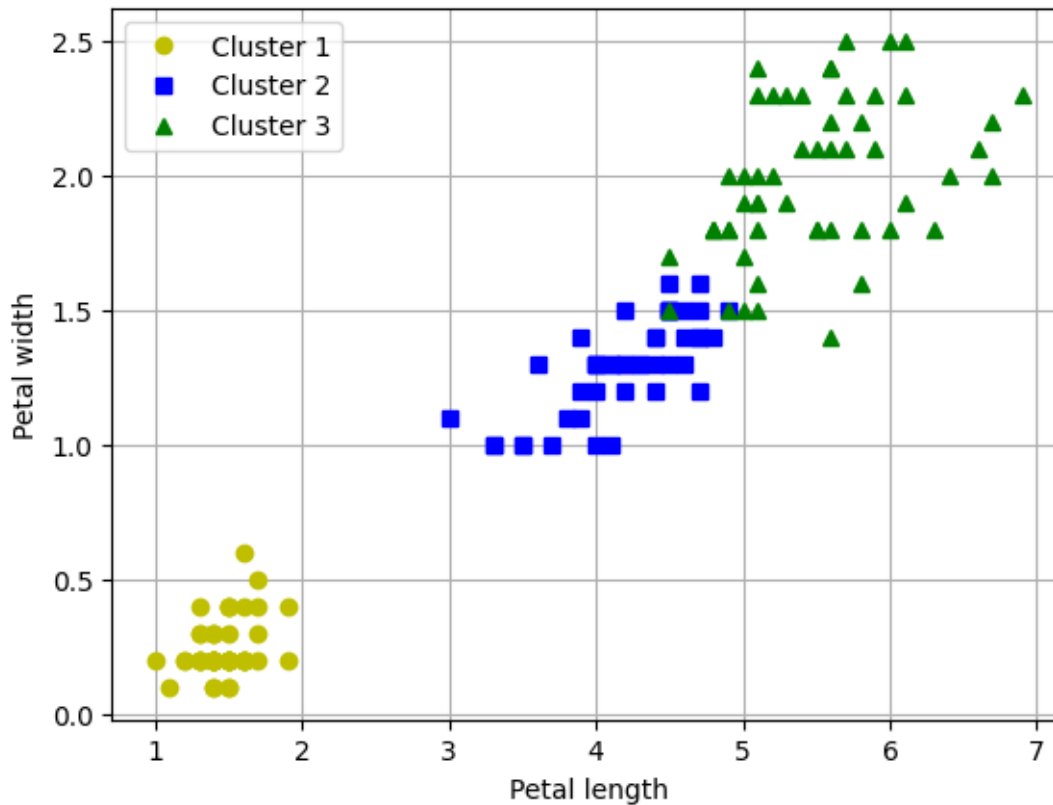
import numpy as np
from scipy import stats
from sklearn.mixture import GaussianMixture

y_pred = GaussianMixture(n_components=3, random_state=42).fit(X).predict(X)

mapping = {}
for class_id in np.unique(y):
    mode, _ = stats.mode(y_pred[y==class_id])
    mapping[mode] = class_id

y_pred = np.array([mapping[cluster_id] for cluster_id in y_pred])

plt.plot(X[y_pred==0, 2], X[y_pred==0, 3], "yo", label="Cluster 1")
plt.plot(X[y_pred==1, 2], X[y_pred==1, 3], "bs", label="Cluster 2")
plt.plot(X[y_pred==2, 2], X[y_pred==2, 3], "g^", label="Cluster 3")
plt.xlabel("Petal length")
plt.ylabel("Petal width")
plt.legend(loc="upper left")
plt.grid()
plt.show()
```



What's the ratio of iris plants we assigned to the right cluster?

```
[3]: (y_pred==y).sum() / len(y_pred)
```

```
[3]: 0.9666666666666667
```

1.1 K-Means

1.1.1 Fit and predict

Let's train a K-Means clusterer on a dataset of blobs. It will try to find each blob's center and assign each instance to the closest blob:

```
[4]: from sklearn.cluster import KMeans
from sklearn.datasets import make_blobs

# extra code - the exact arguments of make_blobs() are not important
blob_centers = np.array([[ 0.2,  2.3], [-1.5,  2.3], [-2.8,  1.8],
                          [-2.8,  2.8], [-2.8,  1.3]])
blob_std = np.array([0.4, 0.3, 0.1, 0.1, 0.1])
X, y = make_blobs(n_samples=2000, centers=blob_centers, cluster_std=blob_std,
                  random_state=7)
```

```
k = 5
kmeans = KMeans(n_clusters=k, n_init=10, random_state=42)
y_pred = kmeans.fit_predict(X)
```

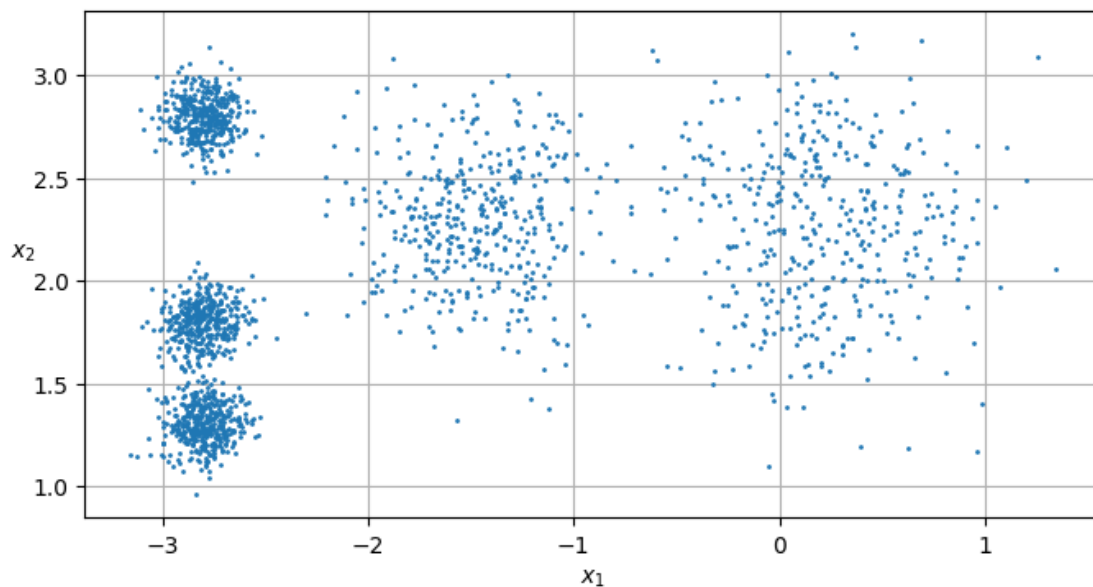
Note: Throughout this notebook, when `n_init` was not set when creating a `KMeans` estimator, I explicitly set it to `n_init=10` to avoid a warning about the fact that the default value for this hyperparameter will change from 10 to “auto” in Scikit-Learn 1.4.

Now let's plot them:

```
[5]: # extra code - this cell generates figure below
```

```
def plot_clusters(X, y=None):
    plt.scatter(X[:, 0], X[:, 1], c=y, s=1)
    plt.xlabel("$x_1$")
    plt.ylabel("$x_2$", rotation=0)

plt.figure(figsize=(8, 4))
plot_clusters(X)
plt.gca().set_axisbelow(True)
plt.grid()
plt.show()
```



Each instance was assigned to one of the 5 clusters:

```
[6]: y_pred
```

```
[6]: array([0, 0, 4, ..., 3, 1, 0])
```

```
[7]: y_pred is kmeans.labels_
```

```
[7]: True
```

And the following 5 centroids (i.e., cluster centers) were estimated:

```
[8]: kmeans.cluster_centers_
```

```
[8]: array([[ -2.80214068,  1.55162671],
          [ 0.08703534,  2.58438091],
          [-1.46869323,  2.28214236],
          [-2.79290307,  2.79641063],
          [ 0.31332823,  1.96822352]])
```

Note that the KMeans instance preserves the labels of the instances it was trained on. Somewhat confusingly, in this context, the label of an instance is the index of the cluster that instance gets assigned to (they are not targets, they are predictions):

```
[9]: kmeans.labels_
```

```
[9]: array([0, 0, 4, ..., 3, 1, 0])
```

Of course, we can predict the labels of new instances:

```
[10]: import numpy as np

X_new = np.array([[0, 2], [3, 2], [-3, 3], [-3, 2.5]])
kmeans.predict(X_new)
```

```
[10]: array([4, 4, 3, 3])
```

1.1.2 Decision Boundaries

Let's plot the model's decision boundaries. This gives us a Voronoi diagram:

```
[11]: # extra code - this cell generates figure below

def plot_data(X):
    plt.plot(X[:, 0], X[:, 1], 'k.', markersize=2)

def plot_centroids(centroids, weights=None, circle_color='w', cross_color='k'):
    if weights is not None:
        centroids = centroids[weights > weights.max() / 10]
    plt.scatter(centroids[:, 0], centroids[:, 1],
                marker='o', s=35, linewidths=8,
                color=circle_color, zorder=10, alpha=0.9)
    plt.scatter(centroids[:, 0], centroids[:, 1],
                marker='x', s=2, linewidths=12,
```

```

        color=cross_color, zorder=11, alpha=1)

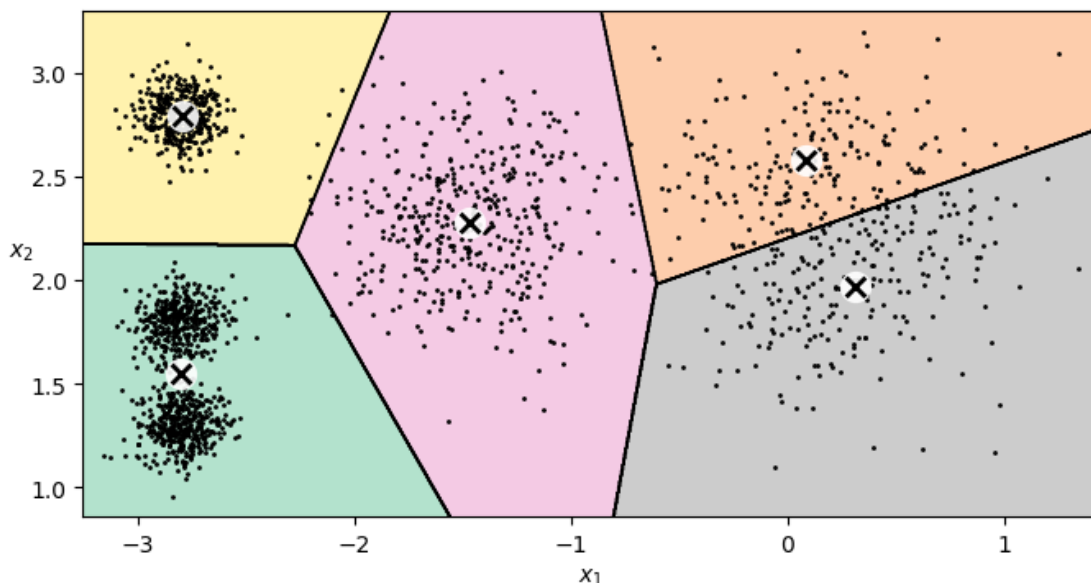
def plot_decision_boundaries(clusterer, X, resolution=1000, show_centroids=True,
                             show_xlabels=True, show_ylabels=True):
    mins = X.min(axis=0) - 0.1
    maxs = X.max(axis=0) + 0.1
    xx, yy = np.meshgrid(np.linspace(mins[0], maxs[0], resolution),
                          np.linspace(mins[1], maxs[1], resolution))
    Z = clusterer.predict(np.c_[xx.ravel(), yy.ravel()])
    Z = Z.reshape(xx.shape)

    plt.contourf(Z, extent=(mins[0], maxs[0], mins[1], maxs[1]),
                 cmap="Pastel2")
    plt.contour(Z, extent=(mins[0], maxs[0], mins[1], maxs[1]),
                linewidths=1, colors='k')
    plot_data(X)
    if show_centroids:
        plot_centroids(clusterer.cluster_centers_)

    if show_xlabels:
        plt.xlabel("$x_1$")
    else:
        plt.tick_params(labelbottom=False)
    if show_ylabels:
        plt.ylabel("$x_2$", rotation=0)
    else:
        plt.tick_params(labelleft=False)

plt.figure(figsize=(8, 4))
plot_decision_boundaries(kmeans, X)
plt.show()

```



Not bad! Some of the instances near the edges were probably assigned to the wrong cluster, but overall it looks pretty good.

1.1.3 Hard Clustering vs Soft Clustering

Rather than arbitrarily choosing the closest cluster for each instance, which is called hard clustering, it might be better to measure the distance of each instance to all 5 centroids. This is what the `transform()` method does:

```
[12]: kmeans.transform(X_new).round(2)
```

```
[12]: array([[2.84, 0.59, 1.5 , 2.9 , 0.31],
            [5.82, 2.97, 4.48, 5.85, 2.69],
            [1.46, 3.11, 1.69, 0.29, 3.47],
            [0.97, 3.09, 1.55, 0.36, 3.36]])
```

You can verify that this is indeed the Euclidian distance between each instance and each centroid:

```
[13]: # extra code
np.linalg.norm(np.tile(X_new, (1, k)).reshape(-1, k, 2)
               - kmeans.cluster_centers_, axis=2).round(2)
```

```
[13]: array([[2.84, 0.59, 1.5 , 2.9 , 0.31],
            [5.82, 2.97, 4.48, 5.85, 2.69],
            [1.46, 3.11, 1.69, 0.29, 3.47],
            [0.97, 3.09, 1.55, 0.36, 3.36]])
```

1.1.4 The K-Means Algorithm

The K-Means algorithm is one of the fastest clustering algorithms, and also one of the simplest:

First initialize k centroids randomly: e.g., k distinct instances are chosen randomly from the dataset and the centroids are placed at their locations. Repeat until convergence (i.e., until the centroids stop moving): Assign each instance to the closest centroid. Update the centroids to be the mean of the instances that are assigned to them.

The `KMeans` class uses an optimized initialization technique by default. To get the original K-Means algorithm (for educational purposes only), you must set `init="random"` and `n_init=1`. More on this later in this topic.

Let's run the K-Means algorithm for 1, 2 and 3 iterations, to see how the centroids move around:

```
[14]: # extra code - this cell generates figure below

kmeans_iter1 = KMeans(n_clusters=5, init="random", n_init=1, max_iter=1,
                      random_state=5)
kmeans_iter2 = KMeans(n_clusters=5, init="random", n_init=1, max_iter=2,
                      random_state=5)
kmeans_iter3 = KMeans(n_clusters=5, init="random", n_init=1, max_iter=3,
                      random_state=5)

kmeans_iter1.fit(X)
kmeans_iter2.fit(X)
kmeans_iter3.fit(X)

plt.figure(figsize=(10, 8))

plt.subplot(321)
plot_data(X)
plot_centroids(kmeans_iter1.cluster_centers_, circle_color='r', cross_color='w')
plt.ylabel("$x_2$", rotation=0)
plt.tick_params(labelbottom=False)
plt.title("Update the centroids (initially randomly)")

plt.subplot(322)
plot_decision_boundaries(kmeans_iter1, X, show_xlabels=False,
                        show_ylabels=False)
plt.title("Label the instances")

plt.subplot(323)
plot_decision_boundaries(kmeans_iter1, X, show_centroids=False,
                        show_xlabels=False)
plot_centroids(kmeans_iter2.cluster_centers_)

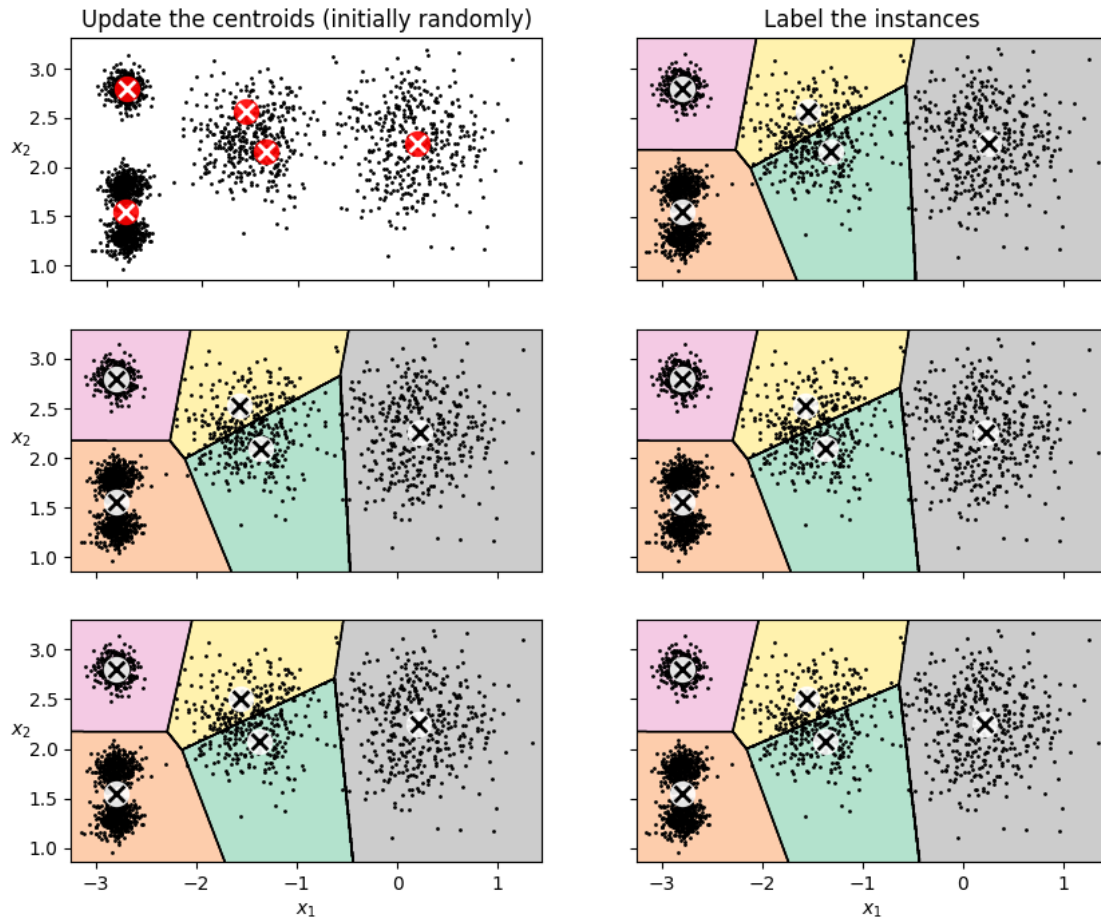
plt.subplot(324)
plot_decision_boundaries(kmeans_iter2, X, show_xlabels=False,
                        show_ylabels=False)
```



```
plt.subplot(325)
plot_decision_boundaries(kmeans_iter2, X, show_centroids=False)
plot_centroids(kmeans_iter3.cluster_centers_)

plt.subplot(326)
plot_decision_boundaries(kmeans_iter3, X, show_ylabels=False)

plt.show()
```



1.1.5 K-Means Variability

In the original K-Means algorithm, the centroids are just initialized randomly, and the algorithm simply runs a single iteration to gradually improve the centroids, as we saw above.

However, one major problem with this approach is that if you run K-Means multiple times (or with different random seeds), it can converge to very different solutions, as you can see below:

```
[15]: # extra code - this cell generates figure below
```

```
def plot_clusterer_comparison(clusterer1, clusterer2, X, title1=None,
                              title2=None):

    clusterer1.fit(X)
    clusterer2.fit(X)

    plt.figure(figsize=(10, 3.2))

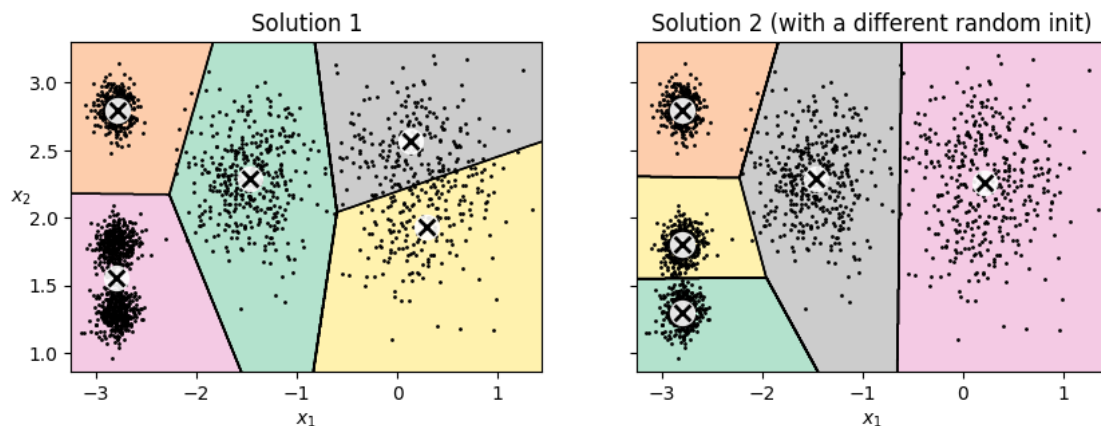
    plt.subplot(121)
    plot_decision_boundaries(clusterer1, X)
    if title1:
        plt.title(title1)

    plt.subplot(122)
    plot_decision_boundaries(clusterer2, X, show_ylabels=False)
    if title2:
        plt.title(title2)

kmeans_rnd_init1 = KMeans(n_clusters=5, init="random", n_init=1, random_state=2)
kmeans_rnd_init2 = KMeans(n_clusters=5, init="random", n_init=1, random_state=9)

plot_clusterer_comparison(kmeans_rnd_init1, kmeans_rnd_init2, X,
                          "Solution 1",
                          "Solution 2 (with a different random init)")

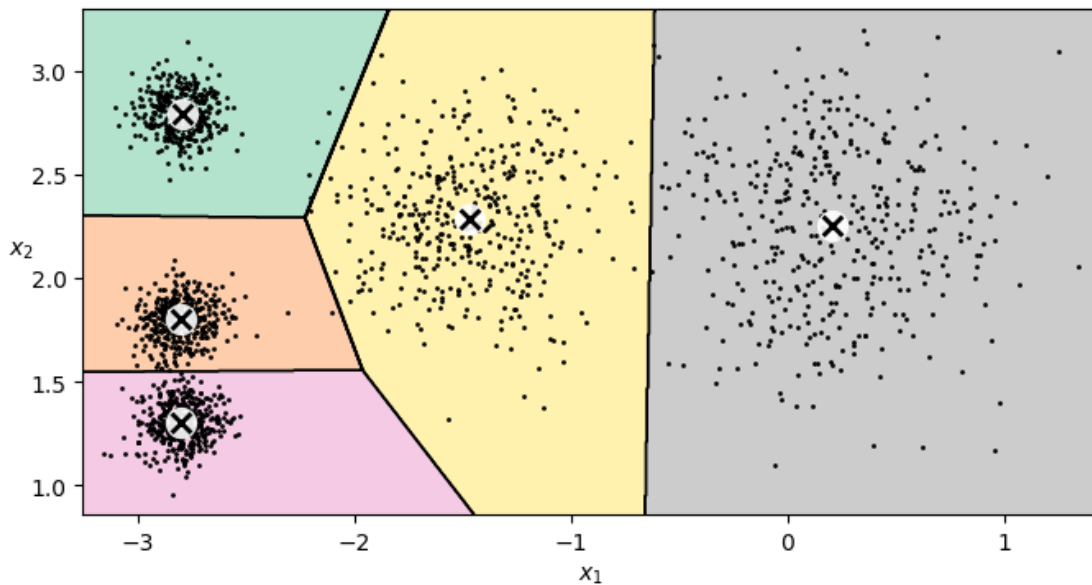
plt.show()
```



```
[16]: good_init = np.array([[ -3,  3], [ -3,  2], [ -3,  1], [ -1,  2], [  0,  2]])
kmeans = KMeans(n_clusters=5, init=good_init, n_init=1, random_state=42)
kmeans.fit(X)
```

```
[16]: KMeans(init=array([[-3,  3],
                        [-3,  2],
                        [-3,  1],
                        [-1,  2],
                        [ 0,  2]]),
        n_clusters=5, n_init=1, random_state=42)
```

```
[17]: # extra code
plt.figure(figsize=(8, 4))
plot_decision_boundaries(kmeans, X)
```



1.1.6 Inertia

To select the best model, we will need a way to evaluate a K-Mean model's performance. Unfortunately, clustering is an unsupervised task, so we do not have the targets. But at least we can measure the distance between each instance and its centroid. This is the idea behind the inertia metric:

```
[18]: kmeans.inertia_
```

```
[18]: 211.59853725816836
```

```
[19]: kmeans_rnd_init1.inertia_ # extra code
```

```
[19]: 219.58201503602294
```

```
[20]: kmeans_rnd_init2.inertia_ # extra code
```

```
[20]: 211.59853725816836
```

As you can easily verify, inertia is the sum of the squared distances between each training instance and its closest centroid:

```
[21]: # extra code
X_dist = kmeans.transform(X)
(X_dist[np.arange(len(X_dist)), kmeans.labels_] ** 2).sum()
```

```
[21]: 211.5985372581688
```

The `score()` method returns the negative inertia. Why negative? Well, it is because a predictor's `score()` method must always respect the “greater is better” rule.

```
[22]: kmeans.score(X)
```

```
[22]: -211.59853725816836
```

1.1.7 Multiple Initializations

So one approach to solve the variability issue is to simply run the K-Means algorithm multiple times with different random initializations, and select the solution that minimizes the inertia.

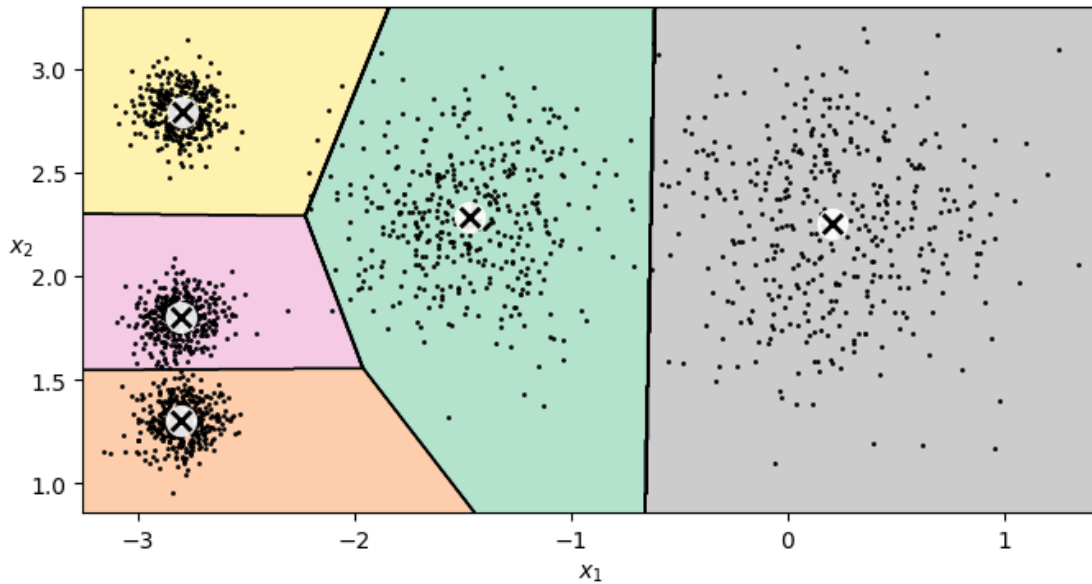
When you set the `n_init` hyperparameter, Scikit-Learn runs the original algorithm `n_init` times, and selects the solution that minimizes the inertia. By default, Scikit-Learn sets `n_init=10`.

```
[23]: # extra code
kmeans_rnd_10_inits = KMeans(n_clusters=5, init="random", n_init=10,
                             random_state=2)
kmeans_rnd_10_inits.fit(X)
```

```
[23]: KMeans(init='random', n_clusters=5, n_init=10, random_state=2)
```

As you can see, we end up with the initial model, which is certainly the optimal K-Means solution (at least in terms of inertia, and assuming `k=5`).

```
[24]: # extra code
plt.figure(figsize=(8, 4))
plot_decision_boundaries(kmeans_rnd_10_inits, X)
plt.show()
```



```
[25]: kmeans_rnd_10_inits.inertia_
```

```
[25]: 211.59853725816836
```

1.1.8 Centroid Initialization methods

Instead of initializing the centroids entirely randomly, it is preferable to initialize them using the following algorithm, proposed in a 2006 paper by David Arthur and Sergei Vassilvitskii:

Take one centroid c_1 , chosen uniformly at random from the dataset. Take a new center c_i , choosing an instance x_i with probability: $D(x_i)^2 / \sum_{j=1}^m D(x_j)^2$ where $D(x_i)$ is the distance between the instance x_i and the closest centroid that was already chosen. This probability distribution ensures that instances that are further away from already chosen centroids are much more likely to be selected as centroids. Repeat the previous step until all k centroids have been chosen. The rest of the K-Means++ algorithm is just regular K-Means. With this initialization, the K-Means algorithm is much less likely to converge to a suboptimal solution, so it is possible to reduce `n_init` considerably. Most of the time, this largely compensates for the additional complexity of the initialization process.

To set the initialization to K-Means++, simply set `init="k-means++"` (this is actually the default):

1.1.9 Accelerated K-Means

The K-Means algorithm can sometimes be accelerated by avoiding many unnecessary distance calculations: this is achieved by exploiting the triangle inequality (given three points A , B and C , the distance AC is always such that $AC \leq AB + BC$) and by keeping track of lower and upper bounds for distances between instances and centroids (see this 2003 paper by Charles Elkan for more details).

For Elkan's variant of K-Means, use `algorithm="elkan"`. For regular KMeans, use `algorithm="full"`.

The default is “auto”, which uses the full algorithm since Scikit-Learn 1.1 (it used Elkan’s algorithm before that).

1.1.10 Mini-Batch K-Means

Scikit-Learn also implements a variant of the K-Means algorithm that supports mini-batches (see this paper):

```
[26]: from sklearn.cluster import MiniBatchKMeans

minibatch_kmeans = MiniBatchKMeans(n_clusters=5, n_init=3, random_state=42)
minibatch_kmeans.fit(X)
```

```
[26]: MiniBatchKMeans(n_clusters=5, n_init=3, random_state=42)
```

Note: Throughout this notebook, when `n_init` was not set when creating a `MiniBatchKMeans` estimator, I explicitly set it to `n_init=3` to avoid a warning about the fact that the default value for this hyperparameter will change from 3 to “auto” in Scikit-Learn 1.4.

```
[27]: minibatch_kmeans.inertia_
```

```
[27]: 211.6589937457431
```

Using `MiniBatchKMeans` along with `memmap`

If the dataset does not fit in memory, the simplest option is to use the `memmap` class, just like we did for incremental PCA in the previous topic. First let’s load MNIST:

```
[28]: from sklearn.datasets import fetch_openml

mnist = fetch_openml('mnist_784', as_frame=False, parser="auto")
```

Let’s split the dataset:

```
[29]: X_train, y_train = mnist.data[:60000], mnist.target[:60000]
X_test, y_test = mnist.data[60000:], mnist.target[60000:]
```

Next, let’s write the training set to a `memmap`:

```
[30]: filename = "my_mnist.mmap"
X_memmap = np.memmap(filename, dtype='float32', mode='write',
                     shape=X_train.shape)
X_memmap[:] = X_train
X_memmap.flush()
```

```
[31]: from sklearn.cluster import MiniBatchKMeans

minibatch_kmeans = MiniBatchKMeans(n_clusters=10, batch_size=10,
                                   n_init=3, random_state=42)
minibatch_kmeans.fit(X_memmap)
```

```
[31]: MiniBatchKMeans(batch_size=10, n_clusters=10, n_init=3, random_state=42)
```

Let's plot the inertia ratio and the training time ratio between Mini-batch K-Means and regular K-Means:

```
[32]: # extra code - this cell generates figure below

from timeit import timeit

max_k = 100
times = np.empty((max_k, 2))
inertias = np.empty((max_k, 2))
for k in range(1, max_k + 1):
    kmeans_ = KMeans(n_clusters=k, algorithm="lloyd", n_init=10,
↳ random_state=42)
    minibatch_kmeans = MiniBatchKMeans(n_clusters=k, n_init=10, random_state=42)
    print(f"\r{k}/{max_k}", end="") # \r returns to the start of line
    times[k - 1, 0] = timeit("kmeans_.fit(X)", number=10, globals=globals())
    times[k - 1, 1] = timeit("minibatch_kmeans.fit(X)", number=10,
                             globals=globals())
    inertias[k - 1, 0] = kmeans_.inertia_
    inertias[k - 1, 1] = minibatch_kmeans.inertia_

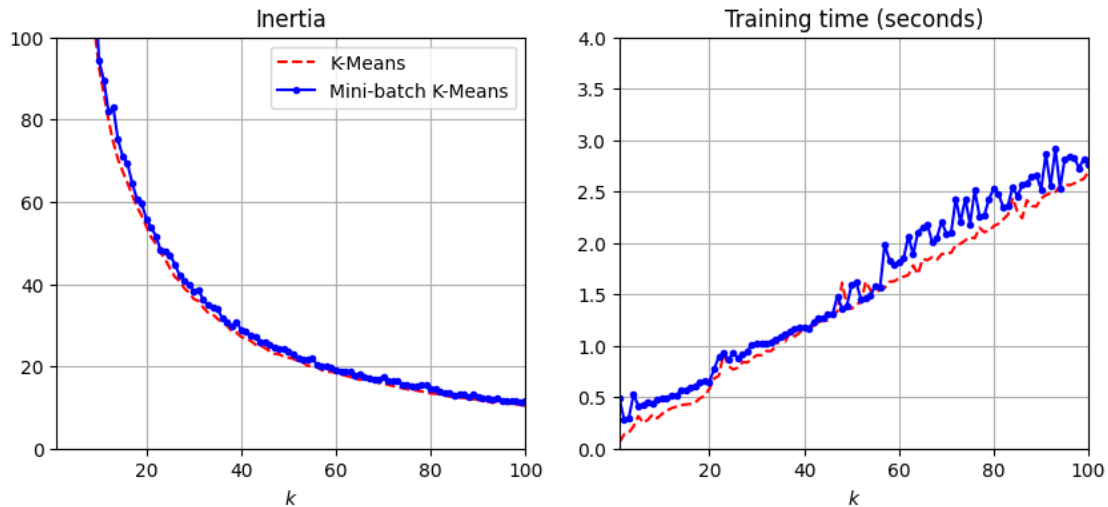
plt.figure(figsize=(10, 4))

plt.subplot(121)
plt.plot(range(1, max_k + 1), inertias[:, 0], "r--", label="K-Means")
plt.plot(range(1, max_k + 1), inertias[:, 1], "b.-", label="Mini-batch K-Means")
plt.xlabel("$k$")
plt.title("Inertia")
plt.legend()
plt.axis([1, max_k, 0, 100])
plt.grid()

plt.subplot(122)
plt.plot(range(1, max_k + 1), times[:, 0], "r--", label="K-Means")
plt.plot(range(1, max_k + 1), times[:, 1], "b.-", label="Mini-batch K-Means")
plt.xlabel("$k$")
plt.title("Training time (seconds)")
plt.axis([1, max_k, 0, 4])
plt.grid()

plt.show()
```

100/100



1.1.11 Finding the optimal number of clusters

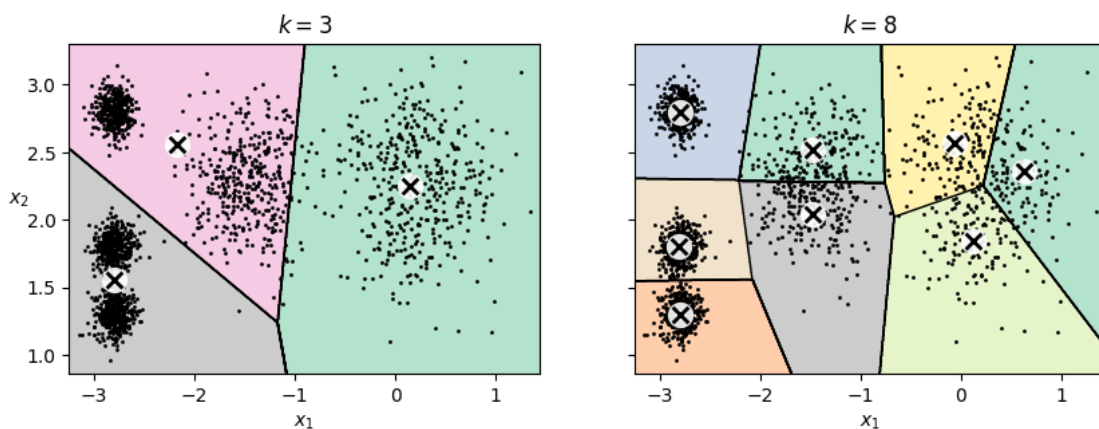
What if the number of clusters was set to a lower or greater value than 5?

```
[33]: # extra code - this cell generates figure below

kmeans_k3 = KMeans(n_clusters=3, n_init=10, random_state=42)
kmeans_k8 = KMeans(n_clusters=8, n_init=10, random_state=42)

plot_clusterer_comparison(kmeans_k3, kmeans_k8, X, "$k=3$", "$k=8$")

plt.show()
```



Ouch, these two models don't look great. What about their inertias?


```
[34]: kmeans_k3.inertia_
```

```
[34]: 653.2167190021553
```

```
[35]: kmeans_k8.inertia_
```

```
[35]: 119.22484592677125
```

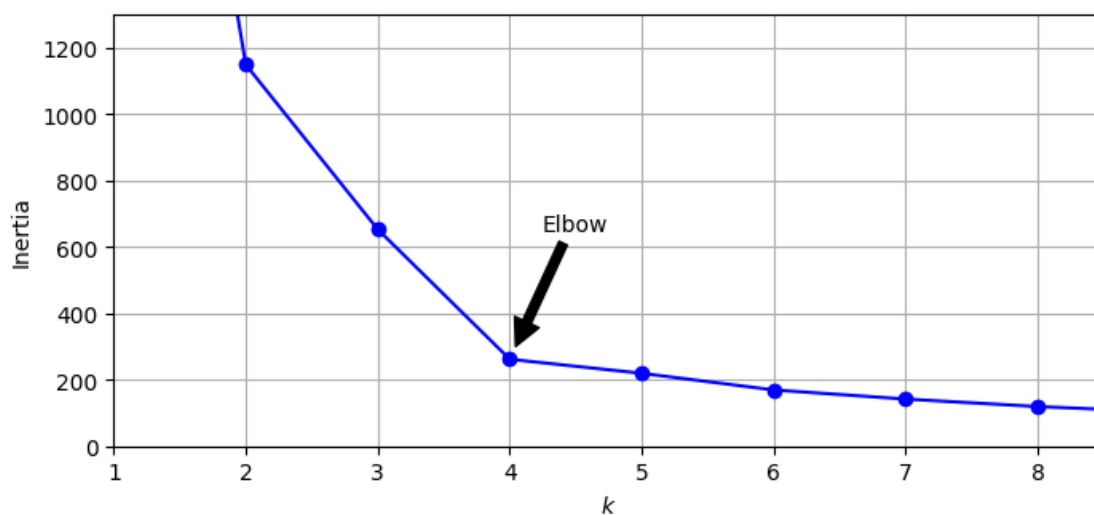
No, we cannot simply take the value of k that minimizes the inertia, since it keeps getting lower as we increase k . Indeed, the more clusters there are, the closer each instance will be to its closest centroid, and therefore the lower the inertia will be. However, we can plot the inertia as a function of k and analyze the resulting curve:

```
[36]: # extra code - this cell generates figure below
```

```
kmeans_per_k = [KMeans(n_clusters=k, n_init=10, random_state=42).fit(X)
                  for k in range(1, 10)]
inertias = [model.inertia_ for model in kmeans_per_k]

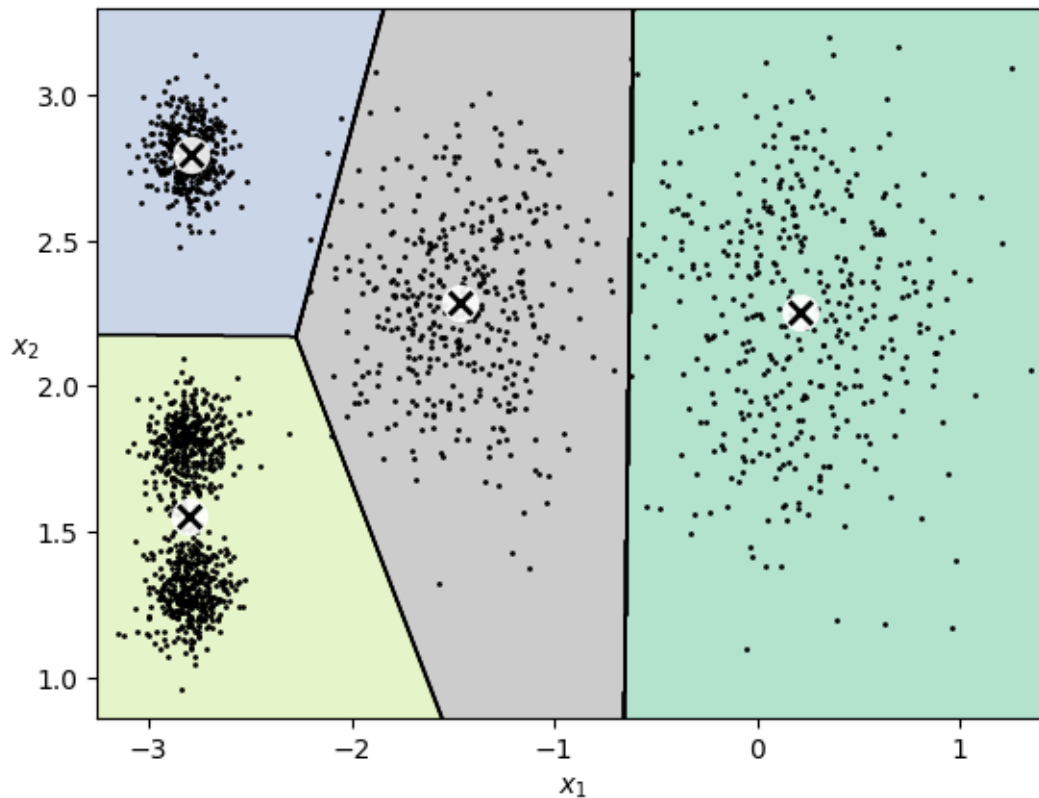
plt.figure(figsize=(8, 3.5))
plt.plot(range(1, 10), inertias, "bo-")
plt.xlabel("$k$")
plt.ylabel("Inertia")
plt.annotate("", xy=(4, inertias[3]), xytext=(4.45, 650),
              arrowprops=dict(facecolor='black', shrink=0.1))
plt.text(4.5, 650, "Elbow", horizontalalignment="center")
plt.axis([1, 8.5, 0, 1300])
plt.grid()

plt.show()
```



As you can see, there is an elbow at $k=4$, which means that less clusters than that would be bad, and more clusters would not help much and might cut clusters in half. So $k=4$ is a pretty good choice. Of course in this example it is not perfect since it means that the two blobs in the lower left will be considered as just a single cluster, but it's a pretty good clustering nonetheless.

```
[37]: # extra code
      plot_decision_boundaries(kmeans_per_k[4 - 1], X)
      plt.show()
```



Another approach is to look at the silhouette score, which is the mean silhouette coefficient over all the instances. An instance's silhouette coefficient is equal to $(b - a) / \max(a, b)$ where a is the mean distance to the other instances in the same cluster (it is the mean intra-cluster distance), and b is the mean nearest-cluster distance, that is the mean distance to the instances of the next closest cluster (defined as the one that minimizes b , excluding the instance's own cluster). The silhouette coefficient can vary between -1 and +1: a coefficient close to +1 means that the instance is well inside its own cluster and far from other clusters, while a coefficient close to 0 means that it is close to a cluster boundary, and finally a coefficient close to -1 means that the instance may have been assigned to the wrong cluster.

Let's plot the silhouette score as a function of k :

```
[38]: from sklearn.metrics import silhouette_score
```

```
[39]: silhouette_score(X, kmeans.labels_)
```

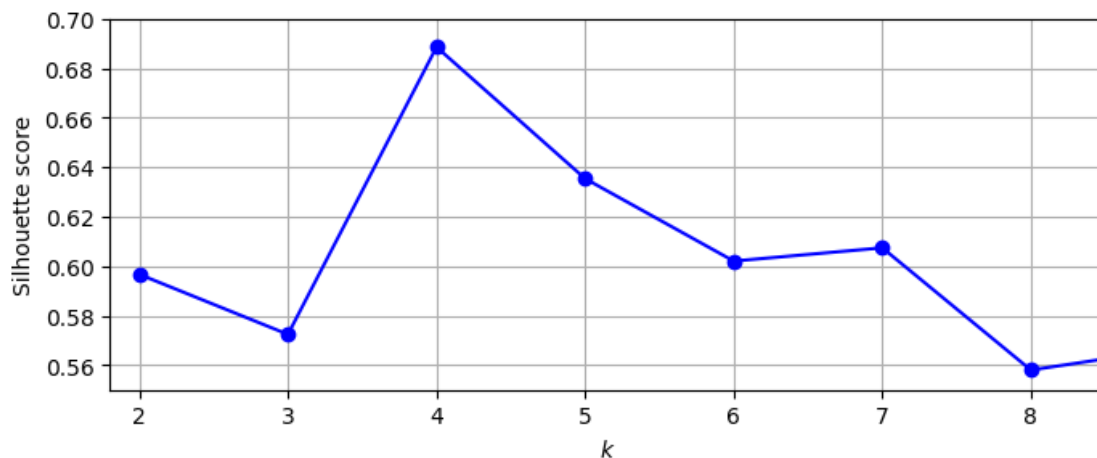
```
[39]: 0.655517642572828
```

```
[40]: # extra code - this cell generates figure below

silhouette_scores = [silhouette_score(X, model.labels_)
                      for model in kmeans_per_k[1:]]

plt.figure(figsize=(8, 3))
plt.plot(range(2, 10), silhouette_scores, "bo-")
plt.xlabel("$k$")
plt.ylabel("Silhouette score")
plt.axis([1.8, 8.5, 0.55, 0.7])
plt.grid()

plt.show()
```



As you can see, this visualization is much richer than the previous one: in particular, although it confirms that $k=4$ is a very good choice, but it also underlines the fact that $k=5$ is quite good as well.

An even more informative visualization is given when you plot every instance's silhouette coefficient, sorted by the cluster they are assigned to and by the value of the coefficient. This is called a silhouette diagram:

```
[41]: # extra code - this cell generates figure below

from sklearn.metrics import silhouette_samples
from matplotlib.ticker import FixedLocator, FixedFormatter
```

```

plt.figure(figsize=(11, 9))

for k in (3, 4, 5, 6):
    plt.subplot(2, 2, k - 2)

    y_pred = kmeans_per_k[k - 1].labels_
    silhouette_coefficients = silhouette_samples(X, y_pred)

    padding = len(X) // 30
    pos = padding
    ticks = []
    for i in range(k):
        coeffs = silhouette_coefficients[y_pred == i]
        coeffs.sort()

        color = plt.cm.Spectral(i / k)
        plt.fill_between(np.arange(pos, pos + len(coeffs)), 0, coeffs,
                        facecolor=color, edgecolor=color, alpha=0.7)
        ticks.append(pos + len(coeffs) // 2)
        pos += len(coeffs) + padding

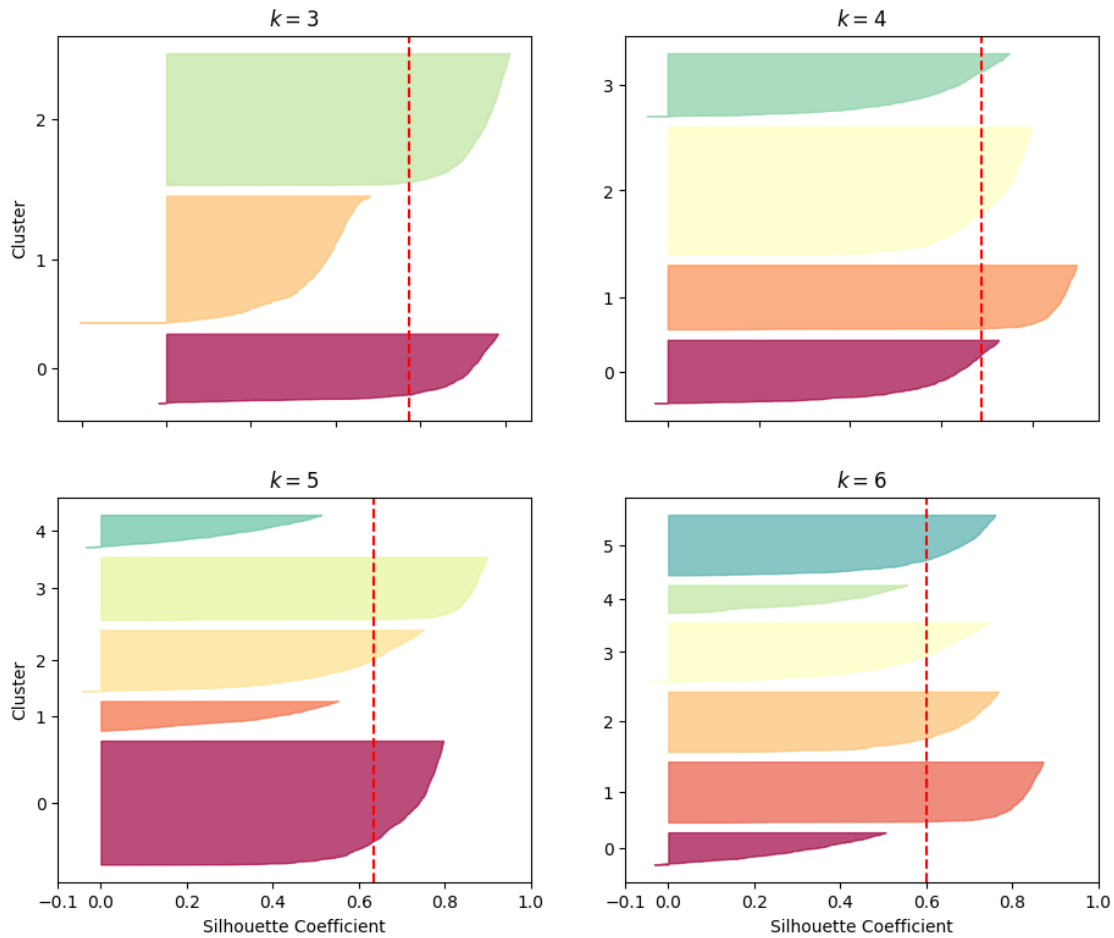
    plt.gca().yaxis.set_major_locator(FixedLocator(ticks))
    plt.gca().yaxis.set_major_formatter(FixedFormatter(range(k)))
    if k in (3, 5):
        plt.ylabel("Cluster")

    if k in (5, 6):
        plt.gca().set_xticks([-0.1, 0, 0.2, 0.4, 0.6, 0.8, 1])
        plt.xlabel("Silhouette Coefficient")
    else:
        plt.tick_params(labelbottom=False)

    plt.axvline(x=silhouette_scores[k - 2], color="red", linestyle="--")
    plt.title(f"$k={k}$")

plt.show()

```



As you can see, $k=5$ looks like the best option here, as all clusters are roughly the same size, and they all cross the dashed line, which represents the mean silhouette score.

[42]: *### Limits of K-Means*

Let's generate a more difficult dataset, with elongated blobs and varying densities, and show that K-Means struggles to cluster it correctly:

[43]: *# extra code - this cell generates figure below*

```
X1, y1 = make_blobs(n_samples=1000, centers=((4, -4), (0, 0)), random_state=42)
X1 = X1.dot(np.array([[0.374, 0.95], [0.732, 0.598]]))
X2, y2 = make_blobs(n_samples=250, centers=1, random_state=42)
X2 = X2 + [6, -8]
X = np.r_[X1, X2]
y = np.r_[y1, y2]

kmeans_good = KMeans(n_clusters=3,
```

```

        init=np.array([[ -1.5, 2.5], [0.5, 0], [4, 0]]),
        n_init=1, random_state=42)
kmeans_bad = KMeans(n_clusters=3, n_init=10, random_state=42)
kmeans_good.fit(X)
kmeans_bad.fit(X)

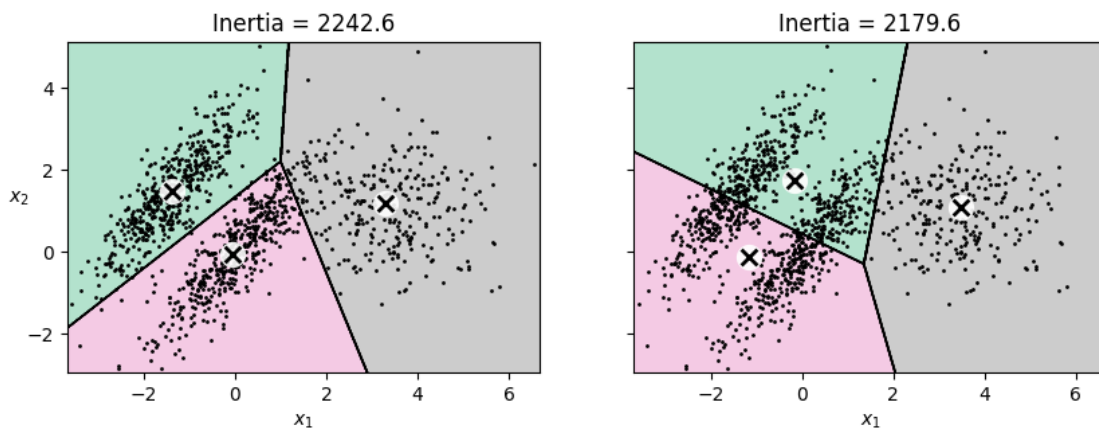
plt.figure(figsize=(10, 3.2))

plt.subplot(121)
plot_decision_boundaries(kmeans_good, X)
plt.title(f"Inertia = {kmeans_good.inertia_:.1f}")

plt.subplot(122)
plot_decision_boundaries(kmeans_bad, X, show_ylabels=False)
plt.title(f"Inertia = {kmeans_bad.inertia_:.1f}")

plt.show()

```



1.2 Using Clustering for Semi-Supervised Learning

Another use case for clustering is semi-supervised learning, when we have plenty of unlabeled instances and very few labeled instances.

Let's tackle the digits dataset which is a simple MNIST-like dataset containing 1,797 grayscale 8×8 images representing digits 0 to 9.

```

[44]: from sklearn.datasets import load_digits

X_digits, y_digits = load_digits(return_X_y=True)
X_train, y_train = X_digits[:1400], y_digits[:1400]
X_test, y_test = X_digits[1400:], y_digits[1400:]

```

Let's look at the performance of a logistic regression model when we only have 50 labeled instances:

```
[45]: from sklearn.linear_model import LogisticRegression
```

```
n_labeled = 50
log_reg = LogisticRegression(max_iter=10_000)
log_reg.fit(X_train[:n_labeled], y_train[:n_labeled])
```

```
[45]: LogisticRegression(max_iter=10000)
```

```
[46]: log_reg.score(X_test, y_test)
```

```
[46]: 0.7581863979848866
```

```
[47]: # extra code - measure the accuracy when we use the whole training set
log_reg_full = LogisticRegression(max_iter=10_000)
log_reg_full.fit(X_train, y_train)
log_reg_full.score(X_test, y_test)
```

```
[47]: 0.9093198992443325
```

It's much less than earlier of course. Let's see how we can do better. First, let's cluster the training set into 50 clusters, then for each cluster let's find the image closest to the centroid. We will call these images the representative images:

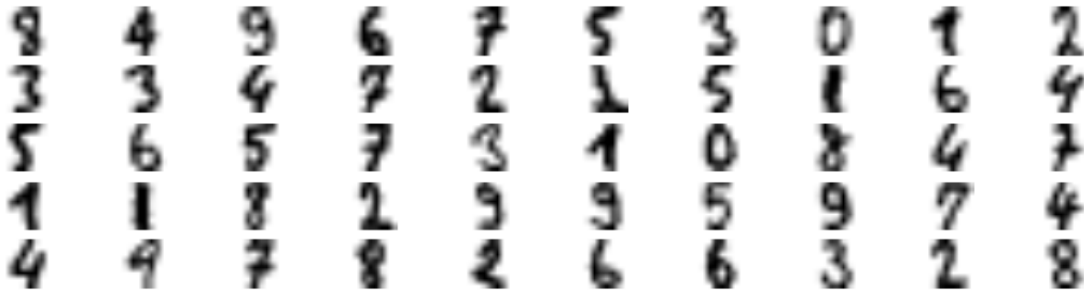
```
[48]: k = 50
kmeans = KMeans(n_clusters=k, n_init=10, random_state=42)
X_digits_dist = kmeans.fit_transform(X_train)
representative_digit_idx = X_digits_dist.argmax(axis=0)
X_representative_digits = X_train[representative_digit_idx]
```

Now let's plot these representative images and label them manually:

```
[49]: # extra code - this cell generates figure below

plt.figure(figsize=(8, 2))
for index, X_representative_digit in enumerate(X_representative_digits):
    plt.subplot(k // 10, 10, index + 1)
    plt.imshow(X_representative_digit.reshape(8, 8), cmap="binary",
               interpolation="bilinear")
    plt.axis('off')

plt.show()
```



```
[50]: y_representative_digits = np.array([
    1, 3, 6, 0, 7, 9, 2, 4, 8, 9,
    5, 4, 7, 1, 2, 6, 1, 2, 5, 1,
    4, 1, 3, 3, 8, 8, 2, 5, 6, 9,
    1, 4, 0, 6, 8, 3, 4, 6, 7, 2,
    4, 1, 0, 7, 5, 1, 9, 9, 3, 7
])
```

Now we have a dataset with just 50 labeled instances, but instead of being completely random instances, each of them is a representative image of its cluster. Let's see if the performance is any better:

```
[51]: log_reg = LogisticRegression(max_iter=10_000)
log_reg.fit(X_representative_digits, y_representative_digits)
log_reg.score(X_test, y_test)
```

```
[51]: 0.06801007556675064
```

Wow! We jumped from 74.8% accuracy to 84.9%, although we are still only training the model on 50 instances. Since it's often costly and painful to label instances, especially when it has to be done manually by experts, it's a good idea to make them label representative instances rather than just random instances.

But perhaps we can go one step further: what if we propagated the labels to all the other instances in the same cluster?

```
[52]: y_train_propagated = np.empty(len(X_train), dtype=np.int64)
for i in range(k):
    y_train_propagated[kmeans.labels_ == i] = y_representative_digits[i]
```

```
[53]: log_reg = LogisticRegression(max_iter=10_000)
log_reg.fit(X_train, y_train_propagated)
```

```
[53]: LogisticRegression(max_iter=10000)
```

```
[54]: log_reg.score(X_test, y_test)
```



```
[54]: 0.07808564231738035
```

We got another significant accuracy boost! Let's see if we can do even better by ignoring the 1% instances that are farthest from their cluster center: this should eliminate some outliers:

```
[55]: percentile_closest = 99

X_cluster_dist = X_digits_dist[np.arange(len(X_train)), kmeans.labels_]
for i in range(k):
    in_cluster = (kmeans.labels_ == i)
    cluster_dist = X_cluster_dist[in_cluster]
    cutoff_distance = np.percentile(cluster_dist, percentile_closest)
    above_cutoff = (X_cluster_dist > cutoff_distance)
    X_cluster_dist[in_cluster & above_cutoff] = -1

partially_propagated = (X_cluster_dist != -1)
X_train_partially_propagated = X_train[partially_propagated]
y_train_partially_propagated = y_train[partially_propagated]
```

```
[56]: log_reg = LogisticRegression(max_iter=10_000)
log_reg.fit(X_train_partially_propagated, y_train_partially_propagated)
log_reg.score(X_test, y_test)
```

```
[56]: 0.07808564231738035
```

Wow, another accuracy boost! We have even slightly surpassed the performance we got by training on the fully labeled training set!

Our propagated labels are actually pretty good: their accuracy is about 97.6%:

```
[57]: (y_train_partially_propagated == y_train[partially_propagated]).mean()
```

```
[57]: 0.09481481481481481
```

You could now do a few iterations of active learning:

Manually label the instances that the classifier is least sure about, if possible by picking them in distinct clusters. Train a new model with these additional labels.

1.3 DBSCAN

```
[58]: from sklearn.cluster import DBSCAN
from sklearn.datasets import make_moons

X, y = make_moons(n_samples=1000, noise=0.05, random_state=42)
dbscan = DBSCAN(eps=0.05, min_samples=5)
dbscan.fit(X)
```

```
[58]: DBSCAN(eps=0.05)
```

```
[59]: dbscan.labels_[:10]
```

```
[59]: array([ 0,  2, -1, -1,  1,  0,  0,  0,  2,  5], dtype=int64)
```

```
[60]: dbscan.core_sample_indices_[:10]
```

```
[60]: array([ 0,  4,  5,  6,  7,  8, 10, 11, 12, 13], dtype=int64)
```

```
[61]: dbscan.components_
```

```
[61]: array([[ -0.02137124,  0.40618608],
          [ -0.84192557,  0.53058695],
          [ 0.58930337, -0.32137599],
          ...,
          [ 1.66258462, -0.3079193 ],
          [-0.94355873,  0.3278936 ],
          [ 0.79419406,  0.60777171]])
```

```
[62]: # extra code - this cell generates figure below
```

```
def plot_dbscan(dbscan, X, size, show_xlabels=True, show_ylabels=True):
    core_mask = np.zeros_like(dbscan.labels_, dtype=bool)
    core_mask[dbscan.core_sample_indices_] = True
    anomalies_mask = dbscan.labels_ == -1
    non_core_mask = ~(core_mask | anomalies_mask)

    cores = dbscan.components_
    anomalies = X[anomalies_mask]
    non_cores = X[non_core_mask]

    plt.scatter(cores[:, 0], cores[:, 1],
                c=dbscan.labels_[core_mask], marker='o', s=size, cmap="Paired")
    plt.scatter(cores[:, 0], cores[:, 1], marker='*', s=20,
                c=dbscan.labels_[core_mask])
    plt.scatter(anomalies[:, 0], anomalies[:, 1],
                c="r", marker="x", s=100)
    plt.scatter(non_cores[:, 0], non_cores[:, 1],
                c=dbscan.labels_[non_core_mask], marker=".")
    if show_xlabels:
        plt.xlabel("$x_1$")
    else:
        plt.tick_params(labelbottom=False)
    if show_ylabels:
        plt.ylabel("$x_2$", rotation=0)
    else:
        plt.tick_params(labelleft=False)
    plt.title(f"eps={dbscan.eps:.2f}, min_samples={dbscan.min_samples}")
    plt.grid()
```

```
plt.gca().set_axisbelow(True)

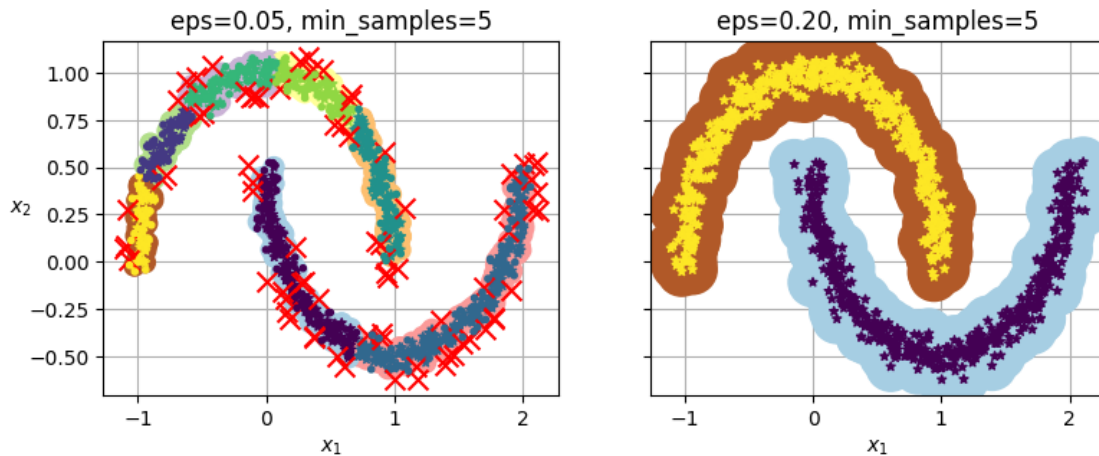
dbscan2 = DBSCAN(eps=0.2)
dbscan2.fit(X)

plt.figure(figsize=(9, 3.2))

plt.subplot(121)
plot_dbscan(dbscan, X, size=100)

plt.subplot(122)
plot_dbscan(dbscan2, X, size=600, show_ylabels=False)

plt.show()
```



```
[63]: dbscan = dbscan2  # extra code - the text says we now use eps=0.2
```

```
[64]: from sklearn.neighbors import KNeighborsClassifier

knn = KNeighborsClassifier(n_neighbors=50)
knn.fit(dbscan.components_, dbscan.labels_[dbscan.core_sample_indices_])
```

```
[64]: KNeighborsClassifier(n_neighbors=50)
```

```
[65]: X_new = np.array([[ -0.5,  0], [ 0,  0.5], [ 1, -0.1], [ 2,  1]])
knn.predict(X_new)
```

```
[65]: array([1, 0, 1, 0], dtype=int64)
```

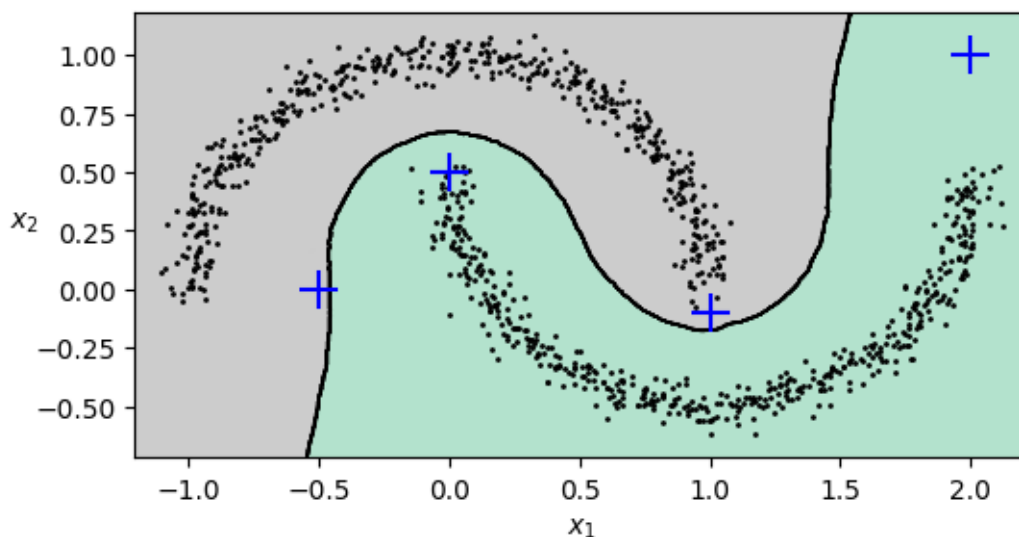
```
[66]: knn.predict_proba(X_new)
```

```
[66]: array([[0.18, 0.82],
            [1.  , 0.  ],
            [0.12, 0.88],
            [1.  , 0.  ]])
```

```
[67]: # extra code - this cell generates figure below
```

```
plt.figure(figsize=(6, 3))
plot_decision_boundaries(knn, X, show_centroids=False)
plt.scatter(X_new[:, 0], X_new[:, 1], c="b", marker="+", s=200, zorder=10)

plt.show()
```



```
[68]: y_dist, y_pred_idx = knn.kneighbors(X_new, n_neighbors=1)
y_pred = dbscan.labels_[dbscan.core_sample_indices_][y_pred_idx]
y_pred[y_dist > 0.2] = -1
y_pred.ravel()
```

```
[68]: array([-1,  0,  1, -1], dtype=int64)
```

1.4 Other Clustering Algorithms

1.4.1 Spectral CLustering

```
[69]: from sklearn.cluster import SpectralClustering
```

```
[70]: sc1 = SpectralClustering(n_clusters=2, gamma=100, random_state=42)
sc1.fit(X)
```

```
[70]: SpectralClustering(gamma=100, n_clusters=2, random_state=42)
```

```
[71]: sc1.affinity_matrix_.round(2)
```

```
[71]: array([[1. , 0. , 0. , ..., 0. , 0. , 0. ],
          [0. , 1. , 0.3, ..., 0. , 0. , 0. ],
          [0. , 0.3, 1. , ..., 0. , 0. , 0. ],
          ...,
          [0. , 0. , 0. , ..., 1. , 0. , 0. ],
          [0. , 0. , 0. , ..., 0. , 1. , 0. ],
          [0. , 0. , 0. , ..., 0. , 0. , 1. ]])
```

```
[72]: sc2 = SpectralClustering(n_clusters=2, gamma=1, random_state=42)
      sc2.fit(X)
```

```
[72]: SpectralClustering(gamma=1, n_clusters=2, random_state=42)
```

```
[73]: def plot_spectral_clustering(sc, X, size, alpha, show_xlabels=True,
                                show_ylabels=True):
    plt.scatter(X[:, 0], X[:, 1], marker='o', s=size, c='gray', alpha=alpha)
    plt.scatter(X[:, 0], X[:, 1], marker='o', s=30, c='w')
    plt.scatter(X[:, 0], X[:, 1], marker='.', s=10, c=sc.labels_, cmap="Paired")

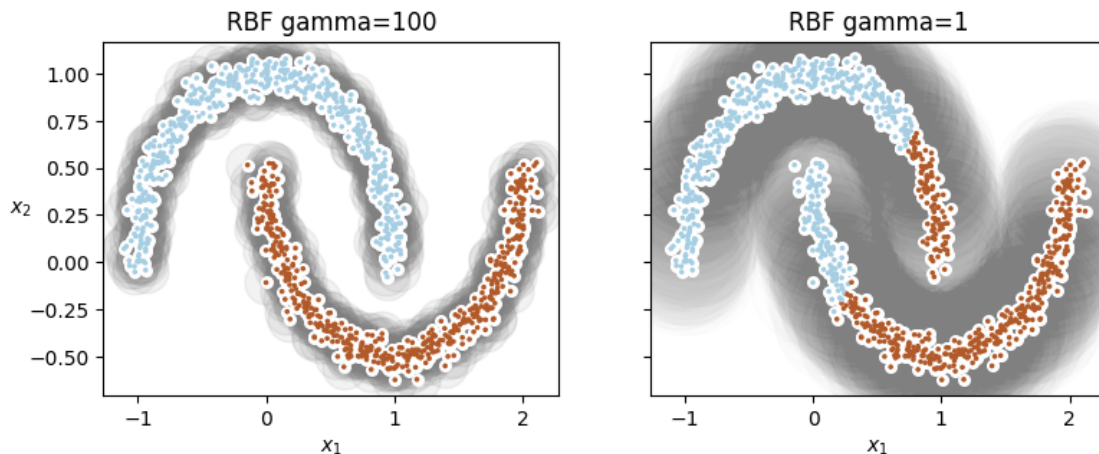
    if show_xlabels:
        plt.xlabel("$x_1$")
    else:
        plt.tick_params(labelbottom=False)
    if show_ylabels:
        plt.ylabel("$x_2$", rotation=0)
    else:
        plt.tick_params(labelleft=False)
    plt.title(f"RBF gamma={sc.gamma}")
```

```
[74]: plt.figure(figsize=(9, 3.2))

      plt.subplot(121)
      plot_spectral_clustering(sc1, X, size=500, alpha=0.1)

      plt.subplot(122)
      plot_spectral_clustering(sc2, X, size=4000, alpha=0.01, show_ylabels=False)

      plt.show()
```



```
[75]: ### Agglomerative Clustering
```

```
[76]: from sklearn.cluster import AgglomerativeClustering
```

```
[77]: X = np.array([0, 2, 5, 8.5]).reshape(-1, 1)
agg = AgglomerativeClustering(linkage="complete").fit(X)
```

```
[78]: def learned_parameters(estimator):
    return [attrib for attrib in dir(estimator)
            if attrib.endswith("_") and not attrib.startswith("_")]
```

```
[79]: learned_parameters(agg)
```

```
[79]: ['children_',
      'labels_',
      'n_clusters_',
      'n_connected_components_',
      'n_features_in_',
      'n_leaves_']
```

```
[80]: agg.children_
```

```
[80]: array([[0, 1],
      [2, 3],
      [4, 5]])
```

2 Gaussian Mixtures

Let's generate the same dataset as earlier with three ellipsoids (the one K-Means had trouble with):

```
[81]: X1, y1 = make_blobs(n_samples=1000, centers=((4, -4), (0, 0)), random_state=42)
      X1 = X1.dot(np.array([[0.374, 0.95], [0.732, 0.598]]))
      X2, y2 = make_blobs(n_samples=250, centers=1, random_state=42)
      X2 = X2 + [6, -8]
      X = np.r_[X1, X2]
      y = np.r_[y1, y2]
```

Let's train a Gaussian mixture model on the previous dataset:

```
[82]: from sklearn.mixture import GaussianMixture
```

```
[83]: gm = GaussianMixture(n_components=3, n_init=10, random_state=42)
      gm.fit(X)
```

```
[83]: GaussianMixture(n_components=3, n_init=10, random_state=42)
```

Let's look at the parameters that the EM algorithm estimated:

```
[84]: gm.weights_
```

```
[84]: array([0.40005972, 0.20961444, 0.39032584])
```

```
[85]: gm.means_
```

```
[85]: array([[ -1.40764129,  1.42712848],
          [ 3.39947665,  1.05931088],
          [ 0.05145113,  0.07534576]])
```

```
[86]: gm.covariances_
```

```
[86]: array([[[ 0.63478217,  0.72970097],
            [ 0.72970097,  1.16094925]],

          [[ 1.14740131, -0.03271106],
            [-0.03271106,  0.95498333]],

          [[ 0.68825143,  0.79617956],
            [ 0.79617956,  1.21242183]]])
```

Did the algorithm actually converge?

```
[87]: gm.converged_
```

```
[87]: True
```

Yes, good. How many iterations did it take?

```
[88]: gm.n_iter_
```

```
[88]: 4
```

You can now use the model to predict which cluster each instance belongs to (hard clustering) or the probabilities that it came from each cluster. For this, just use `predict()` method or the `predict_proba()` method:

```
[89]: gm.predict(X)
```

```
[89]: array([2, 2, 0, ..., 1, 1, 1], dtype=int64)
```

```
[90]: gm.predict_proba(X).round(3)
```

```
[90]: array([[0.    , 0.023, 0.977],
          [0.001, 0.016, 0.983],
          [1.    , 0.    , 0.    ],
          ...,
          [0.    , 1.    , 0.    ],
          [0.    , 1.    , 0.    ],
          [0.    , 1.    , 0.    ]])
```

This is a generative model, so you can sample new instances from it (and get their labels):

```
[91]: X_new, y_new = gm.sample(6)
      X_new
```

```
[91]: array([[ -2.32491052,  1.04752548],
          [-1.16654983,  1.62795173],
          [ 1.84860618,  2.07374016],
          [ 3.98304484,  1.49869936],
          [ 3.8163406 ,  0.53038367],
          [ 0.38079484, -0.56239369]])
```

```
[92]: y_new
```

```
[92]: array([0, 0, 1, 1, 1, 2])
```

Notice that they are sampled sequentially from each cluster.

You can also estimate the log of the probability density function (PDF) at any location using the `score_samples()` method:

```
[93]: gm.score_samples(X).round(2)
```

```
[93]: array([-2.61, -3.57, -3.33, ..., -3.51, -4.4 , -3.81])
```

Let's check that the PDF integrates to 1 over the whole space. We just take a large square around the clusters, and chop it into a grid of tiny squares, then we compute the approximate probability that the instances will be generated in each tiny square (by multiplying the PDF at one corner of the tiny square by the area of the square), and finally summing all these probabilities). The result is very close to 1:


```
[94]: # extra code - bonus material

resolution = 100
grid = np.arange(-10, 10, 1 / resolution)
xx, yy = np.meshgrid(grid, grid)
X_full = np.vstack([xx.ravel(), yy.ravel()]).T

pdf = np.exp(gm.score_samples(X_full))
pdf_probas = pdf * (1 / resolution) ** 2
pdf_probas.sum()
```

```
[94]: 0.9999999999225089
```

Now let's plot the resulting decision boundaries (dashed lines) and density contours:

```
[95]: # extra code - this cells generates figure below

from matplotlib.colors import LogNorm

def plot_gaussian_mixture(clusterer, X, resolution=1000, show_ylabels=True):
    mins = X.min(axis=0) - 0.1
    maxs = X.max(axis=0) + 0.1
    xx, yy = np.meshgrid(np.linspace(mins[0], maxs[0], resolution),
                          np.linspace(mins[1], maxs[1], resolution))
    Z = -clusterer.score_samples(np.c_[xx.ravel(), yy.ravel()])
    Z = Z.reshape(xx.shape)

    plt.contourf(xx, yy, Z,
                 norm=LogNorm(vmin=1.0, vmax=30.0),
                 levels=np.logspace(0, 2, 12))
    plt.contour(xx, yy, Z,
               norm=LogNorm(vmin=1.0, vmax=30.0),
               levels=np.logspace(0, 2, 12),
               linewidths=1, colors='k')

    Z = clusterer.predict(np.c_[xx.ravel(), yy.ravel()])
    Z = Z.reshape(xx.shape)
    plt.contour(xx, yy, Z,
               linewidths=2, colors='r', linestyle='dashed')

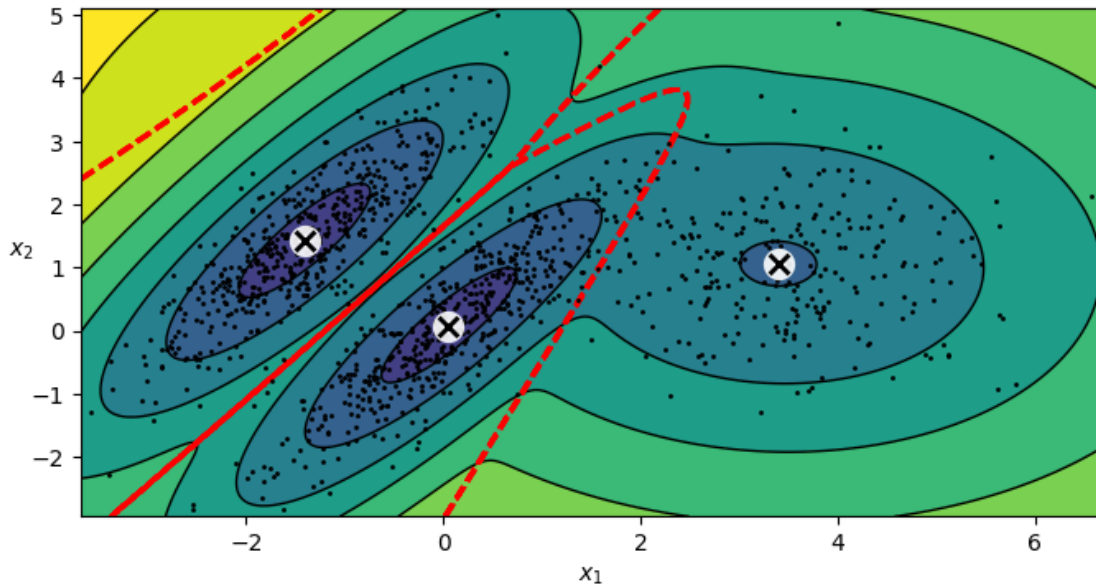
    plt.plot(X[:, 0], X[:, 1], 'k.', markersize=2)
    plot_centroids(clusterer.means_, clusterer.weights_)

    plt.xlabel("$x_1$")
    if show_ylabels:
        plt.ylabel("$x_2$", rotation=0)
    else:
        plt.tick_params(labelleft=False)
```

```
plt.figure(figsize=(8, 4))

plot_gaussian_mixture(gm, X)

plt.show()
```



You can impose constraints on the covariance matrices that the algorithm looks for by setting the `covariance_type` hyperparameter:

“spherical”: all clusters must be spherical, but they can have different diameters (i.e., different variances). “diag”: clusters can take on any ellipsoidal shape of any size, but the ellipsoid’s axes must be parallel to the axes (i.e., the covariance matrices must be diagonal). “tied”: all clusters must have the same shape, which can be any ellipsoid (i.e., they all share the same covariance matrix). “full” (default): no constraint, all clusters can take on any ellipsoidal shape of any size.

```
[96]: # extra code - this cell generates figure below

gm_full = GaussianMixture(n_components=3, n_init=10,
                           covariance_type="full", random_state=42)
gm_tied = GaussianMixture(n_components=3, n_init=10,
                           covariance_type="tied", random_state=42)
gm_spherical = GaussianMixture(n_components=3, n_init=10,
                                covariance_type="spherical", random_state=42)
gm_diag = GaussianMixture(n_components=3, n_init=10,
                            covariance_type="diag", random_state=42)

gm_full.fit(X)
gm_tied.fit(X)
```

```

gm_spherical.fit(X)
gm_diag.fit(X)

def compare_gaussian_mixtures(gm1, gm2, X):
    plt.figure(figsize=(9, 4))

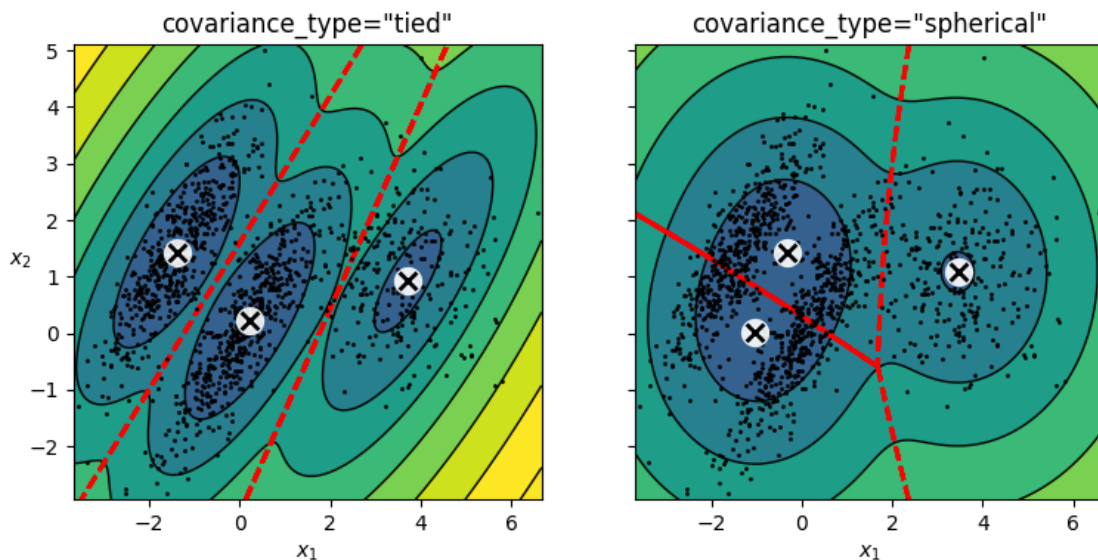
    plt.subplot(121)
    plot_gaussian_mixture(gm1, X)
    plt.title(f'covariance_type="{gm1.covariance_type}"')

    plt.subplot(122)
    plot_gaussian_mixture(gm2, X, show_ylabels=False)
    plt.title(f'covariance_type="{gm2.covariance_type}"')

compare_gaussian_mixtures(gm_tied, gm_spherical, X)

plt.show()

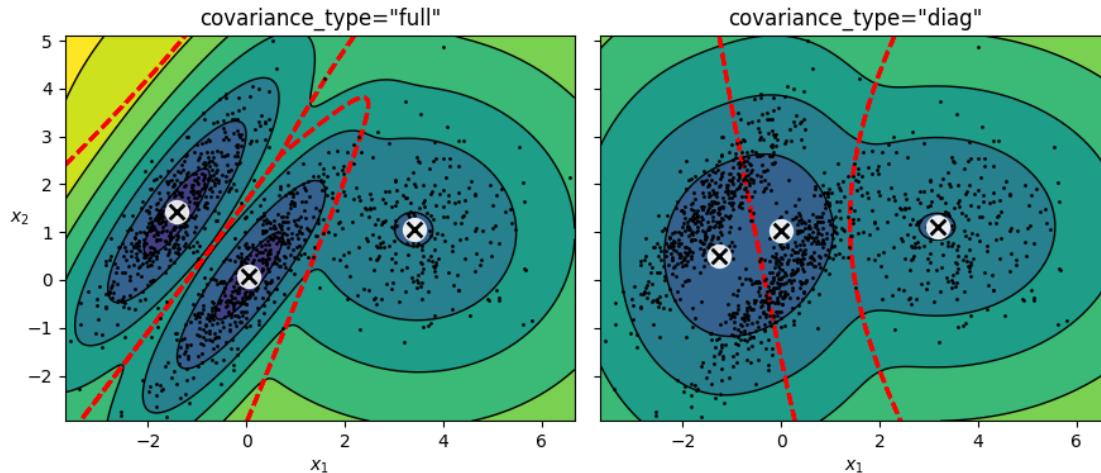
```



```

[97]: # extra code - comparing covariance_type="full" and covariance_type="diag"
compare_gaussian_mixtures(gm_full, gm_diag, X)
plt.tight_layout()
plt.show()

```



2.0.1 Anomaly Detection Using Gaussian Mixtures

Gaussian Mixtures can be used for anomaly detection: instances located in low-density regions can be considered anomalies. You must define what density threshold you want to use. For example, in a manufacturing company that tries to detect defective products, the ratio of defective products is usually well-known. Say it is equal to 2%, then you can set the density threshold to be the value that results in having 2% of the instances located in areas below that threshold density:

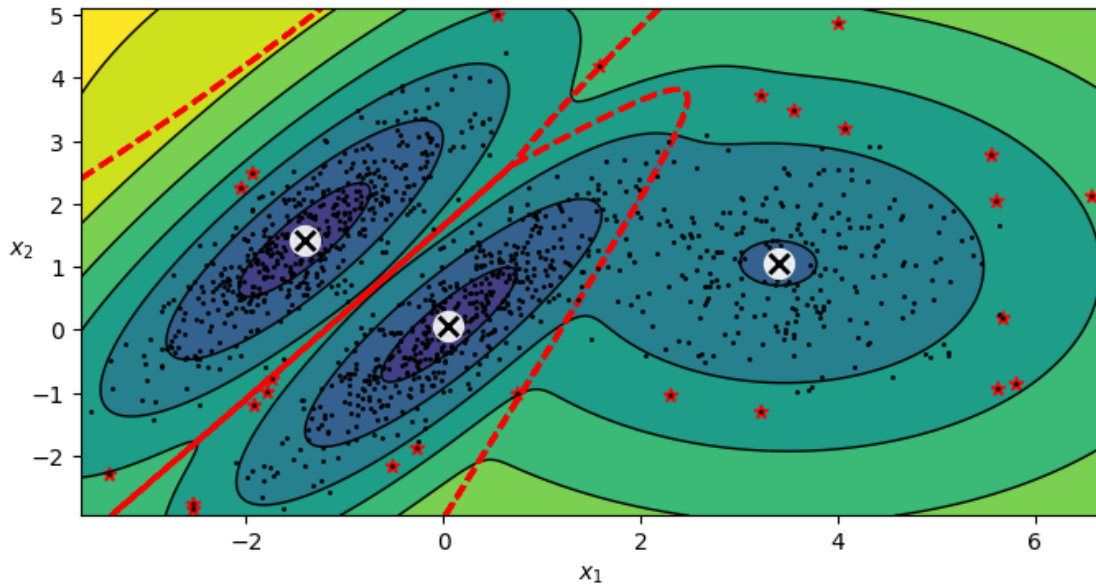
```
[98]: densities = gm.score_samples(X)
      density_threshold = np.percentile(densities, 2)
      anomalies = X[densities < density_threshold]
```

```
[99]: # extra code - this cell generates figure below

plt.figure(figsize=(8, 4))

plot_gaussian_mixture(gm, X)
plt.scatter(anomalies[:, 0], anomalies[:, 1], color='r', marker='*')
plt.ylim(top=5.1)

plt.show()
```



2.0.2 Selecting the Number of Clusters

We cannot use the inertia or the silhouette score because they both assume that the clusters are spherical. Instead, we can try to find the model that minimizes a theoretical information criterion such as the Bayesian Information Criterion (BIC) or the Akaike Information Criterion (AIC):

$$\text{BIC} = \log(m)p - 2\log(\hat{L})$$

$$\text{AIC} = 2p - 2\log(\hat{L})$$

m is the number of instances. p is the number of parameters learned by the model. \hat{L} is the maximized value of the likelihood function of the model. This is the conditional probability of the observed data X , given the model and its optimized parameters. Both BIC and AIC penalize models that have more parameters to learn (e.g., more clusters), and reward models that fit the data well (i.e., models that give a high likelihood to the observed data).

```
[100]: # extra code - this cell generates figure below

from scipy.stats import norm

x_val = 2.5
std_val = 1.3
x_range = [-6, 4]
x_proba_range = [-2, 2]
stds_range = [1, 2]

xs = np.linspace(x_range[0], x_range[1], 501)
stds = np.linspace(stds_range[0], stds_range[1], 501)
Xs, Stds = np.meshgrid(xs, stds)
```

```

Z = 2 * norm.pdf(Xs - 1.0, 0, Stds) + norm.pdf(Xs + 4.0, 0, Stds)
Z = Z / Z.sum(axis=1)[: , np.newaxis] / (xs[1] - xs[0])

x_example_idx = (xs >= x_val).argmax() # index of the first value >= x_val
max_idx = Z[:, x_example_idx].argmax()
max_val = Z[:, x_example_idx].max()
s_example_idx = (stds >= std_val).argmax()
x_range_min_idx = (xs >= x_proba_range[0]).argmax()
x_range_max_idx = (xs >= x_proba_range[1]).argmax()
log_max_idx = np.log(Z[:, x_example_idx]).argmax()
log_max_val = np.log(Z[:, x_example_idx]).max()

plt.figure(figsize=(8, 4.5))

plt.subplot(2, 2, 1)
plt.contourf(Xs, Stds, Z, cmap="GnBu")
plt.plot([-6, 4], [std_val, std_val], "k-", linewidth=2)
plt.plot([x_val, x_val], [1, 2], "b-", linewidth=2)
plt.ylabel(r"$\theta$", rotation=0, labelpad=10)
plt.title(r"Model $f(x; \theta)$")

plt.subplot(2, 2, 2)
plt.plot(stds, Z[:, x_example_idx], "b-")
plt.plot(stds[max_idx], max_val, "r.")
plt.plot([stds[max_idx], stds[max_idx]], [0, max_val], "r:")
plt.plot([0, stds[max_idx]], [max_val, max_val], "r:")
plt.text(stds[max_idx] + 0.01, 0.081, r"$\hat{\theta}$")
plt.text(stds[max_idx] + 0.01, max_val - 0.006, r"$Max$")
plt.text(1.01, max_val - 0.008, r"$\hat{\mathcal{L}}$")
plt.ylabel(r"$\mathcal{L}$", rotation=0, labelpad=10)
plt.title(fr"$\mathcal{L}(\theta|x=\{x\_val\}) = f(x=\{x\_val\}; \theta)$")
plt.grid()
plt.axis([1, 2, 0.08, 0.12])

plt.subplot(2, 2, 3)
plt.plot(xs, Z[s_example_idx], "k-")
plt.fill_between(xs[x_range_min_idx:x_range_max_idx+1],
                 Z[s_example_idx, x_range_min_idx:x_range_max_idx+1], alpha=0.2)
plt.xlabel(r"$x$")
plt.ylabel("PDF")
plt.title(fr"PDF $f(x; \theta=\{std\_val\})$")
plt.grid()
plt.axis([-6, 4, 0, 0.25])

plt.subplot(2, 2, 4)
plt.plot(stds, np.log(Z[:, x_example_idx]), "b-")
plt.plot(stds[log_max_idx], log_max_val, "r.")

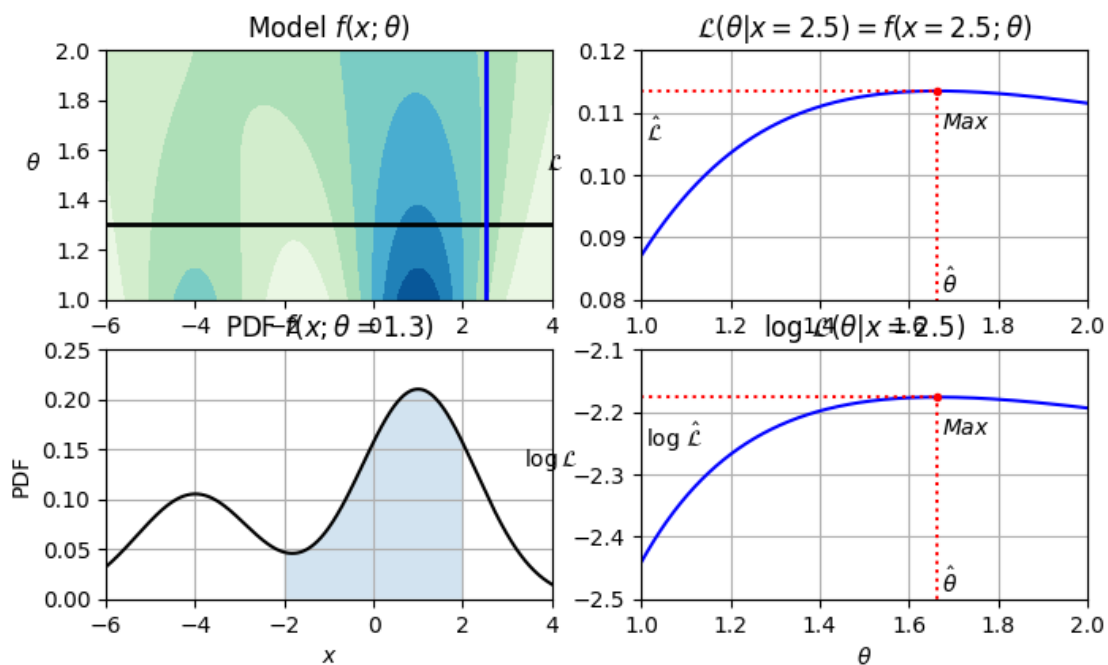
```

```

plt.plot([stds[log_max_idx], stds[log_max_idx]], [-5, log_max_val], "r:")
plt.plot([0, stds[log_max_idx]], [log_max_val, log_max_val], "r:")
plt.text(stds[log_max_idx]+ 0.01, log_max_val - 0.06, r"$Max$")
plt.text(stds[log_max_idx]+ 0.01, -2.49, r"$\hat{\theta}$")
plt.text(1.01, log_max_val - 0.08, r"$\log \, , \, \hat{\mathcal{L}}$")
plt.xlabel(r"$\theta$")
plt.ylabel(r"$\log \mathcal{L}$", rotation=0, labelpad=10)
plt.title(fr"$\log \, , \, \hat{\mathcal{L}}(\theta|x=\{x\_val\})$")
plt.grid()
plt.axis([1, 2, -2.5, -2.1])

plt.show()

```



```
[101]: gm.bic(X)
```

```
[101]: 8189.733705221636
```

```
[102]: gm.aic(X)
```

```
[102]: 8102.508425106598
```

We could compute the BIC manually like this:

```

[103]: # extra code - bonus material
n_clusters = 3

```

```

n_dims = 2
n_params_for_weights = n_clusters - 1
n_params_for_means = n_clusters * n_dims
n_params_for_covariance = n_clusters * n_dims * (n_dims + 1) // 2
n_params = n_params_for_weights + n_params_for_means + n_params_for_covariance
max_log_likelihood = gm.score(X) * len(X) # log(L~)
bic = np.log(len(X)) * n_params - 2 * max_log_likelihood
aic = 2 * n_params - 2 * max_log_likelihood
print(f"bic = {bic}")
print(f"aic = {aic}")
print(f"n_params = {n_params}")

```

```

bic = 8189.733705221636
aic = 8102.508425106598
n_params = 17

```

There's one weight per cluster, but the sum must be equal to 1, so we have one degree of freedom less, hence the -1. Similarly, the degrees of freedom for an $n \times n$ covariance matrix is not n^2 , but $1 + 2 + \dots + n = n(n+1)/2$.

Let's train Gaussian Mixture models with various values of k and measure their BIC:

```

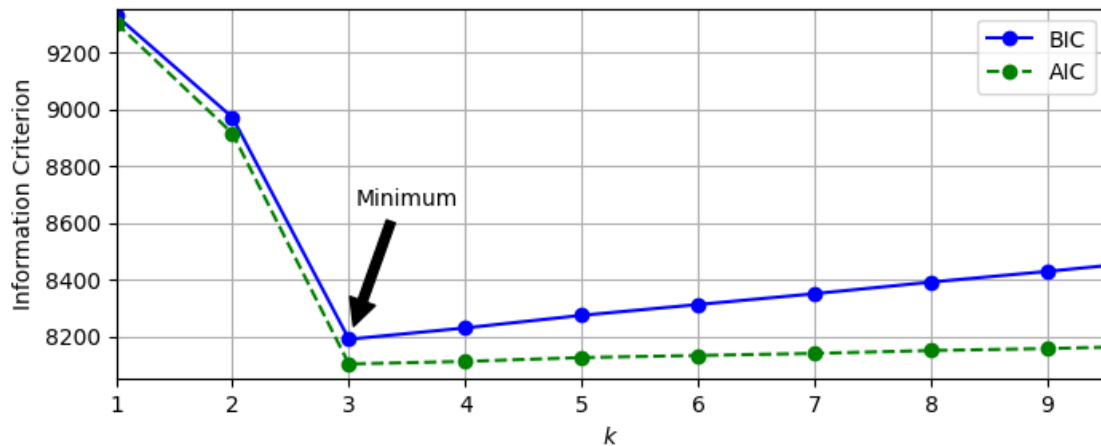
[104]: # extra code - this cell generates figure below

gms_per_k = [GaussianMixture(n_components=k, n_init=10, random_state=42).fit(X)
              for k in range(1, 11)]
bics = [model.bic(X) for model in gms_per_k]
aics = [model.aic(X) for model in gms_per_k]

plt.figure(figsize=(8, 3))
plt.plot(range(1, 11), bics, "bo-", label="BIC")
plt.plot(range(1, 11), aics, "go--", label="AIC")
plt.xlabel("$k$")
plt.ylabel("Information Criterion")
plt.axis([1, 9.5, min(aics) - 50, max(aics) + 50])
plt.annotate("", xy=(3, bics[2]), xytext=(3.4, 8650),
             arrowprops=dict(facecolor='black', shrink=0.1))
plt.text(3.5, 8660, "Minimum", horizontalalignment="center")
plt.legend()
plt.grid()

plt.show()

```

2.0.3 Bayesian Gaussian Mixture Models

Rather than manually searching for the optimal number of clusters, it is possible to use instead the `BayesianGaussianMixture` class which is capable of giving weights equal (or close) to zero to unnecessary clusters. Just set the number of components to a value that you believe is greater than the optimal number of clusters, and the algorithm will eliminate the unnecessary clusters automatically.

```
[105]: from sklearn.mixture import BayesianGaussianMixture

bgm = BayesianGaussianMixture(n_components=10, n_init=10, random_state=42)
bgm.fit(X)
bgm.weights_.round(2)
```

```
[105]: array([0.4 , 0.21, 0. , 0.4 , 0. , 0. , 0. , 0. , 0. , 0. ])
```

```
[106]: The algorithm automatically detected that only 3 components are needed!
```

```
Cell In[106], line 1
```

```
The algorithm automatically detected that only 3 components are needed!
```

```
SyntaxError: invalid syntax
```

```
[ ]: # extra code
plt.figure(figsize=(8, 5))
plot_gaussian_mixture(bgm, X)
plt.show()
```

```
[ ]: # extra code - this cell generates figure below

X_moons, y_moons = make_moons(n_samples=1000, noise=0.05, random_state=42)

bgm = BayesianGaussianMixture(n_components=10, n_init=10, random_state=42)
bgm.fit(X_moons)

plt.figure(figsize=(9, 3.2))

plt.subplot(121)
plot_data(X_moons)
plt.xlabel("$x_1$")
plt.ylabel("$x_2$", rotation=0)
plt.grid()

plt.subplot(122)
plot_gaussian_mixture(bgm, X_moons, show_ylabels=False)

plt.show()
```

Oops, not great... instead of detecting 2 moon-shaped clusters, the algorithm detected 8 ellipsoidal clusters. However, the density plot does not look too bad, so it might be usable for anomaly detection.