

Chapter 9 – Unsupervised Learning

This notebook contains all the sample code in chapter 9.



Run in Google Colab

Setup

First, let's import a few common modules, ensure Matplotlib plots figures inline and prepare a function to save the figures. We also check that Python 3.5 or later is installed (although Python 2.x may work, it is deprecated so we strongly recommend you use Python 3 instead), as well as Scikit-Learn ≥ 0.20 .

```
In [1]: # Python  $\geq 3.5$  is required
import sys
assert sys.version_info >= (3, 5)

# Scikit-Learn  $\geq 0.20$  is required
import sklearn
assert sklearn.__version__ >= "0.20"

# Common imports
import numpy as np
import os

# to make this notebook's output stable across runs
np.random.seed(42)

# To plot pretty figures
%matplotlib inline
import matplotlib as mpl
import matplotlib.pyplot as plt
mpl.rc('axes', labelsize=14)
mpl.rc('xtick', labelsize=12)
mpl.rc('ytick', labelsize=12)

# Where to save the figures
PROJECT_ROOT_DIR = "."
CHAPTER_ID = "unsupervised_learning"
IMAGES_PATH = os.path.join(PROJECT_ROOT_DIR, "images", CHAPTER_ID)
os.makedirs(IMAGES_PATH, exist_ok=True)

def save_fig(fig_id, tight_layout=True, fig_extension="png", resolution=300):
    path = os.path.join(IMAGES_PATH, fig_id + "." + fig_extension)
    print("Saving figure", fig_id)
    if tight_layout:
        plt.tight_layout()
    plt.savefig(path, format=fig_extension, dpi=resolution)

# Ignore useless warnings (see SciPy issue #5998)
import warnings
warnings.filterwarnings(action="ignore", message="^internal gelsd")
```

Clustering

Introduction – Classification vs Clustering

```
In [2]: from sklearn.datasets import load_iris
```

```
In [3]: data = load_iris()
X = data.data
y = data.target
data.target_names
```

```
Out[3]: array(['setosa', 'versicolor', 'virginica'], dtype='<U10')
```

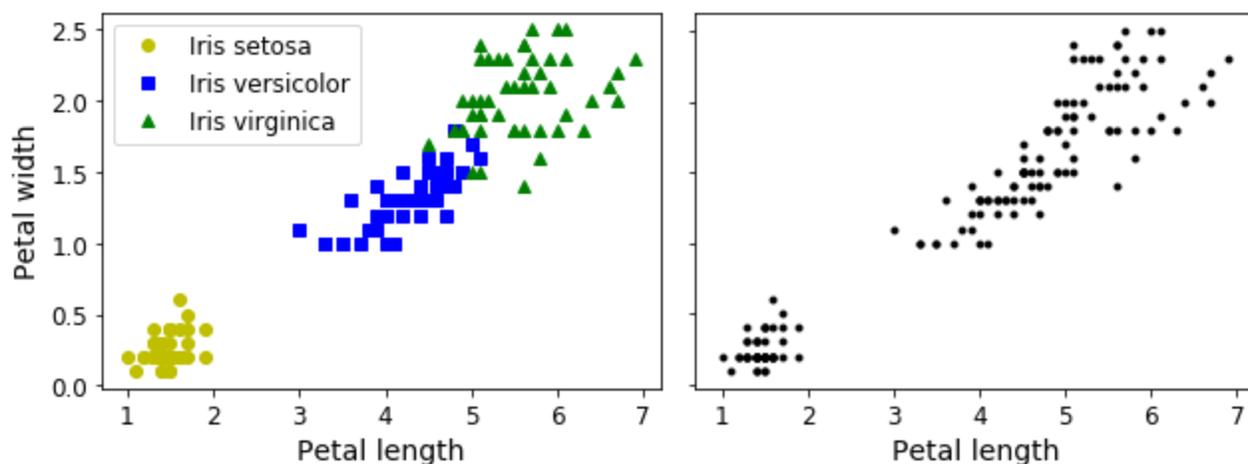
```
In [4]: plt.figure(figsize=(9, 3.5))

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", fontsize=14)
plt.ylabel("Petal width", fontsize=14)
plt.legend(fontsize=12)

plt.subplot(122)
plt.scatter(X[:, 2], X[:, 3], c="k", marker=".")
plt.xlabel("Petal length", fontsize=14)
plt.tick_params(labelleft=False)

save_fig("classification_vs_clustering_plot")
plt.show()
```

Saving figure classification_vs_clustering_plot



K-Means

Let's start by generating some blobs:

```
In [10]: from sklearn.datasets import make_blobs
```

```
In [11]: 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])
```

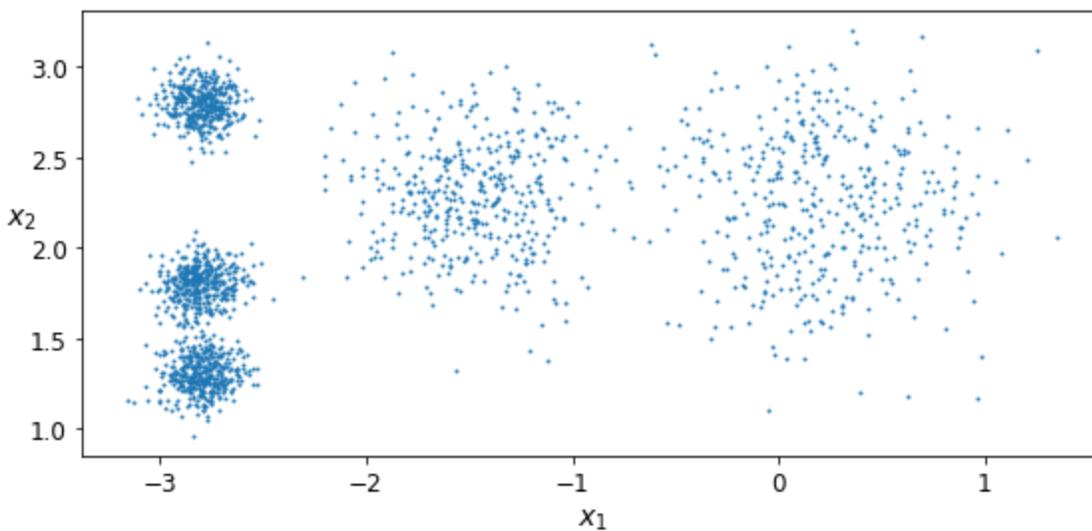
```
In [12]: X, y = make_blobs(n_samples=2000, centers=blob_centers,
                           cluster_std=blob_std, random_state=7)
```

Now let's plot them:

```
In [13]: def plot_clusters(X, y=None):
          plt.scatter(X[:, 0], X[:, 1], c=y, s=1)
          plt.xlabel("$x_1$", fontsize=14)
          plt.ylabel("$x_2$", fontsize=14, rotation=0)
```

```
In [14]: plt.figure(figsize=(8, 4))
          plot_clusters(X)
          save_fig("blobs_plot")
          plt.show()
```

Saving figure blobs_plot



Fit and Predict

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

```
In [15]: from sklearn.cluster import KMeans
```

```
In [16]: k = 5
          kmeans = KMeans(n_clusters=k, random_state=42)
          y_pred = kmeans.fit_predict(X)
```

Each instance was assigned to one of the 5 clusters:

```
In [17]: y_pred
```

```
Out[17]: array([4, 0, 1, ..., 2, 1, 0], dtype=int32)
```

```
In [18]: y_pred is kmeans.labels_
```

```
Out[18]: True
```

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

```
In [19]: kmeans.cluster_centers_
```

```
Out[19]: array([[ -2.80389616,  1.80117999],
                [  0.20876306,  2.25551336],
                [ -2.79290307,  2.79641063],
                [ -1.46679593,  2.28585348],
                [ -2.80037642,  1.30082566]])
```

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:

```
In [20]: kmeans.labels_
```

```
Out[20]: array([4, 0, 1, ..., 2, 1, 0], dtype=int32)
```

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

```
In [21]: X_new = np.array([[0, 2], [3, 2], [-3, 3], [-3, 2.5]])
kmeans.predict(X_new)
```

```
Out[21]: array([1, 1, 2, 2], dtype=int32)
```

Decision Boundaries

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

```
In [22]: 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=30, linewidths=8,
                color=circle_color, zorder=10, alpha=0.9)
    plt.scatter(centroids[:, 0], centroids[:, 1],
                marker='x', s=50, linewidths=50,
                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("$x_1$", fontsize=14)
else:
    plt.tick_params(labelbottom=False)
if show_ylabels:
    plt.ylabel("$x_2$", fontsize=14, rotation=0)
else:
    plt.tick_params(labelleft=False)

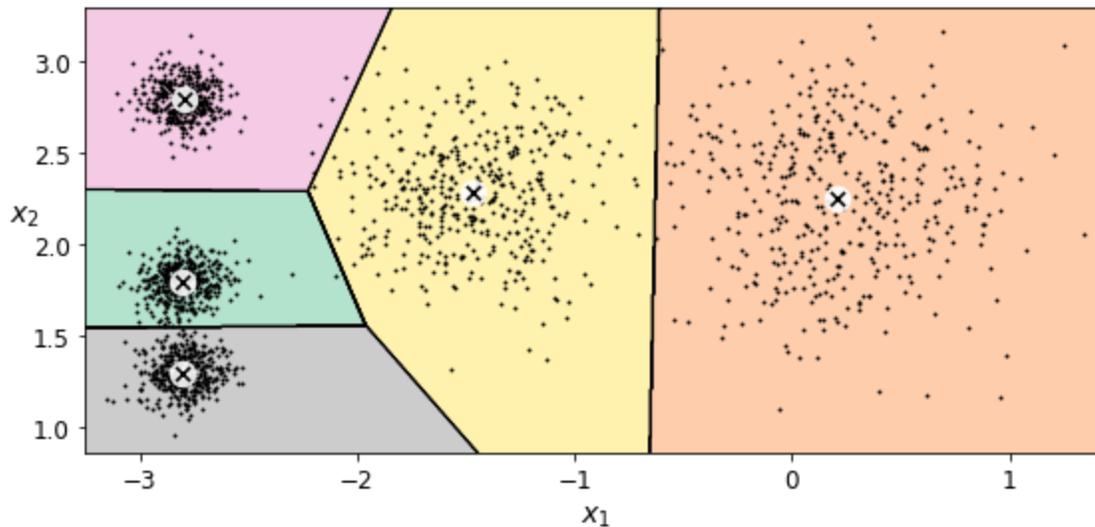
```

```

In [23]: plt.figure(figsize=(8, 4))
plot_decision_boundaries(kmeans, X)
save_fig("voronoi_plot")
plt.show()

```

Saving figure voronoi_plot



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

Hard Clustering vs Soft Clustering

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

```

In [24]: kmeans.transform(X_new)

```

```

Out[24]: array([[2.81093633, 0.32995317, 2.9042344 , 1.49439034, 2.88633901],
 [5.80730058, 2.80290755, 5.84739223, 4.4759332 , 5.84236351],
 [1.21475352, 3.29399768, 0.29040966, 1.69136631, 1.71086031],
 [0.72581411, 3.21806371, 0.36159148, 1.54808703, 1.21567622]])

```

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

```

In [25]: np.linalg.norm(np.tile(X_new, (1, k)).reshape(-1, k, 2) - kmeans.cluster_centers_, axis=

```

```
Out[25]: array([[2.81093633, 0.32995317, 2.9042344 , 1.49439034, 2.88633901],
 [5.80730058, 2.80290755, 5.84739223, 4.4759332 , 5.84236351],
 [1.21475352, 3.29399768, 0.29040966, 1.69136631, 1.71086031],
 [0.72581411, 3.21806371, 0.36159148, 1.54808703, 1.21567622]])
```

K-Means Algorithm

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

- First initialize k centroids randomly: 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 applies an optimized algorithm by default. To get the original K-Means algorithm (for educational purposes only), you must set `init="random"`, `n_init=1` and `algorithm="full"`. These hyperparameters will be explained below.

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

```
In [26]: kmeans_iter1 = KMeans(n_clusters=5, init="random", n_init=1,
                             algorithm="full", max_iter=1, random_state=1)
kmeans_iter2 = KMeans(n_clusters=5, init="random", n_init=1,
                     algorithm="full", max_iter=2, random_state=1)
kmeans_iter3 = KMeans(n_clusters=5, init="random", n_init=1,
                     algorithm="full", max_iter=3, random_state=1)

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

```
Out[26]: KMeans(algorithm='full', copy_x=True, init='random', max_iter=3, n_clusters=5,
               n_init=1, n_jobs=None, precompute_distances='auto', random_state=1,
               tol=0.0001, verbose=0)
```

And let's plot this:

```
In [27]: 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("$x_2$", fontsize=14, rotation=0)
plt.tick_params(labelbottom=False)
plt.title("Update the centroids (initially randomly)", fontsize=14)

plt.subplot(322)
plot_decision_boundaries(kmeans_iter1, X, show_xlabels=False, show_ylabels=False)
plt.title("Label the instances", fontsize=14)

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

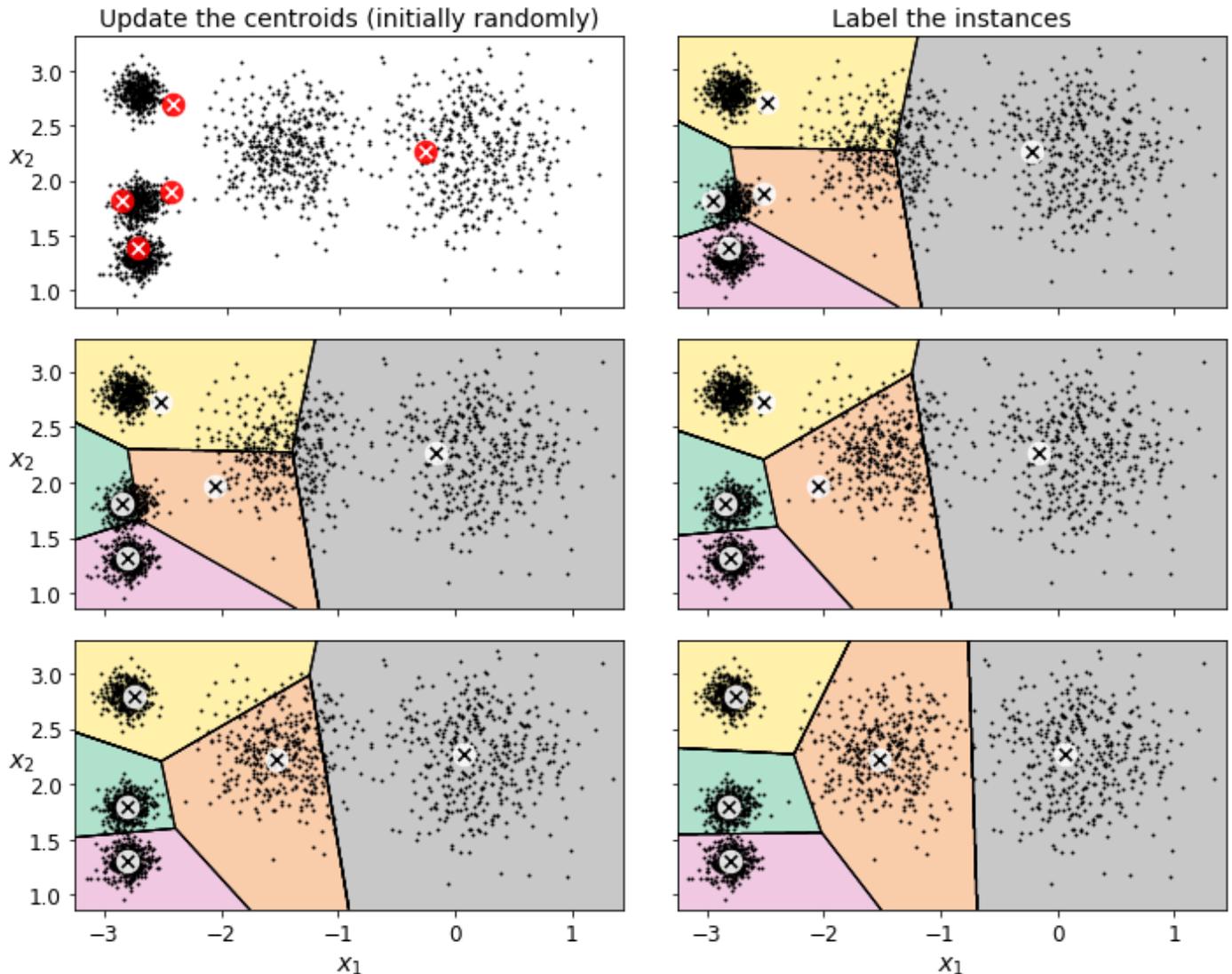
plt.subplot(324)
plot_decision_boundaries(kmeans_iter2, X, show_xlabels=False, show_ylabels=False)
```

```
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)

save_fig("kmeans_algorithm_plot")
plt.show()
```

Saving figure kmeans_algorithm_plot



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:

```
In [28]: 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 title:
    plt.title(title1, fontsize=14)

plt.subplot(122)
plot_decision_boundaries(clusterer2, X, show_ylabels=False)
if title2:
    plt.title(title2, fontsize=14)

```

```

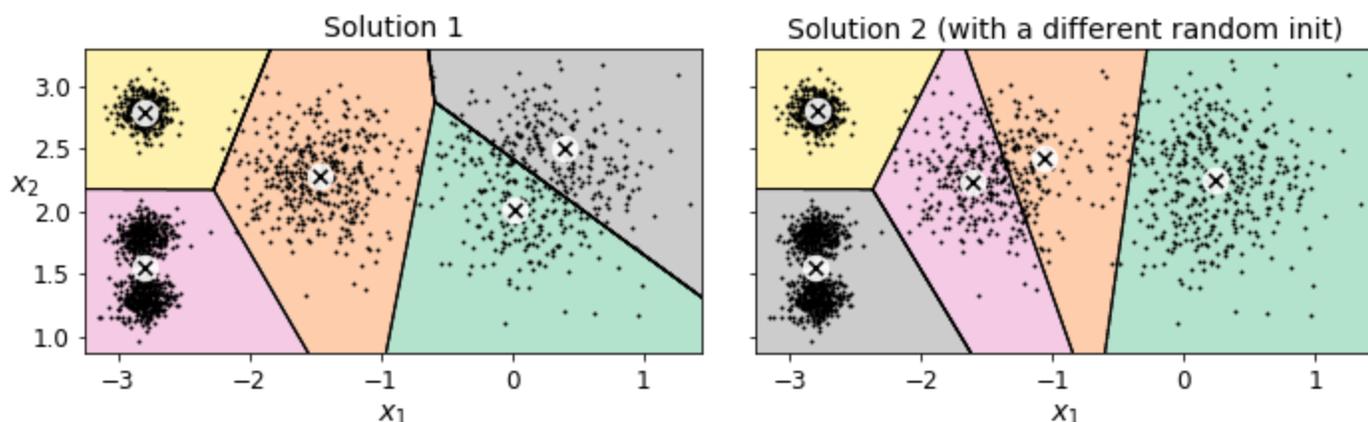
In [29]: kmeans_rnd_init1 = KMeans(n_clusters=5, init="random", n_init=1,
                                algorithm="full", random_state=11)
kmeans_rnd_init2 = KMeans(n_clusters=5, init="random", n_init=1,
                          algorithm="full", random_state=19)

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

save_fig("kmeans_variability_plot")
plt.show()

```

Saving figure kmeans_variability_plot



Inertia

To select the best model, we will need a way to evaluate a K-Means 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:

```
In [30]: kmeans.inertia_
```

```
Out[30]: 211.5985372581684
```

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

```
In [31]: X_dist = kmeans.transform(X)
np.sum(X_dist[np.arange(len(X_dist))], kmeans.labels_**2)
```

```
Out[31]: 211.59853725816856
```

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

```
In [32]: kmeans.score(X)
```

```
Out[32]: -211.59853725816856
```

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. For example, here are the inertias of the two "bad" models shown in the previous figure:

```
In [33]: kmeans_rnd_init1.inertia_
```

```
Out[33]: 223.29108572819035
```

```
In [34]: kmeans_rnd_init2.inertia_
```

```
Out[34]: 237.46249169442845
```

As you can see, they have a higher inertia than the first "good" model we trained, which means they are probably worse.

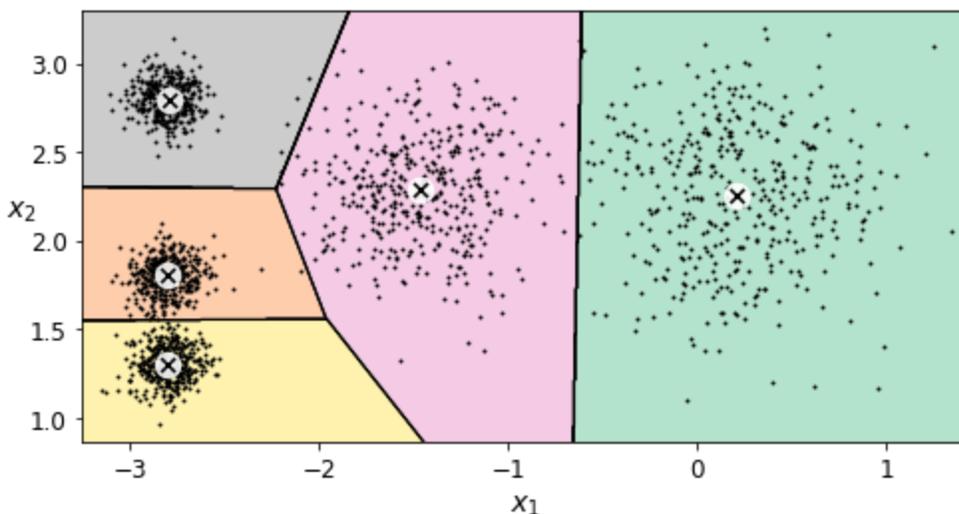
When you set the `n_init` hyperparameter, Scikit-Learn runs the original algorithm `n_init` times, and selects the solution that minimizes the inertia. By default, Scikit-Learn sets `n_init=10`.

```
In [35]: kmeans_rnd_10_inits = KMeans(n_clusters=5, init="random", n_init=10,
                                      algorithm="full", random_state=11)
kmeans_rnd_10_inits.fit(X)
```

```
Out[35]: KMeans(algorithm='full', copy_x=True, init='random', max_iter=300,
                n_clusters=5, n_init=10, n_jobs=None, precompute_distances='auto',
                random_state=11, tol=0.0001, verbose=0)
```

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 $k = 5$).

```
In [36]: plt.figure(figsize=(8, 4))
plot_decision_boundaries(kmeans_rnd_10_inits, X)
plt.show()
```



K-Means++

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 Vassilvskii:

- Take one centroid c_1 , chosen uniformly at random from the dataset.
- Take a new center c_i , choosing an instance \mathbf{x}_i with probability: $D(\mathbf{x}_i)^2 / \sum_{j=1}^m D(\mathbf{x}_j)^2$ where $D(\mathbf{x}_i)$ is the distance between the instance \mathbf{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 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):

```
In [37]: KMeans()
```

```
Out[37]: KMeans(algorithm='auto', copy_x=True, init='k-means++', max_iter=300,
               n_clusters=8, n_init=10, n_jobs=None, precompute_distances='auto',
               random_state=None, tol=0.0001, verbose=0)
```

```
In [38]: 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)
         kmeans.inertia_
```

```
Out[38]: 211.5985372581684
```

Accelerated K-Means

The K-Means algorithm can be significantly 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).

To use Elkan's variant of K-Means, just set `algorithm="elkan"`. Note that it does not support sparse data, so by default, Scikit-Learn uses `"elkan"` for dense data, and `"full"` (the regular K-Means algorithm) for sparse data.

```
In [39]: %timeit -n 50 KMeans(algorithm="elkan").fit(X)
```

```
83.7 ms ± 2.07 ms per loop (mean ± std. dev. of 7 runs, 50 loops each)
```

```
In [40]: %timeit -n 50 KMeans(algorithm="full").fit(X)
```

```
106 ms ± 2.34 ms per loop (mean ± std. dev. of 7 runs, 50 loops each)
```

Mini-Batch K-Means

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

```
In [41]: from sklearn.cluster import MiniBatchKMeans
```

```
In [42]: minibatch_kmeans = MiniBatchKMeans(n_clusters=5, random_state=42)
minibatch_kmeans.fit(X)
```

```
Out[42]: MiniBatchKMeans(batch_size=100, compute_labels=True, init='k-means++',
    init_size=None, max_iter=100, max_no_improvement=10, n_clusters=5,
    n_init=3, random_state=42, reassignment_ratio=0.01, tol=0.0,
    verbose=0)
```

```
In [43]: minibatch_kmeans.inertia_
```

```
Out[43]: 211.93186531476775
```

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 chapter. First let's load MNIST:

```
In [44]: import urllib
from sklearn.datasets import fetch_openml

mnist = fetch_openml('mnist_784', version=1)
mnist.target = mnist.target.astype(np.int64)
```

```
In [45]: from sklearn.model_selection import train_test_split

X_train, X_test, y_train, y_test = train_test_split(
    mnist["data"], mnist["target"], random_state=42)
```

Next, let's write it to a `memmap` :

```
In [46]: filename = "my_mnist.data"
X_mm = np.memmap(filename, dtype='float32', mode='write', shape=X_train.shape)
X_mm[:] = X_train
```

```
In [47]: minibatch_kmeans = MiniBatchKMeans(n_clusters=10, batch_size=10, random_state=42)
minibatch_kmeans.fit(X_mm)
```

```
Out[47]: MiniBatchKMeans(batch_size=10, compute_labels=True, init='k-means++',
    init_size=None, max_iter=100, max_no_improvement=10, n_clusters=10,
    n_init=3, random_state=42, reassignment_ratio=0.01, tol=0.0,
    verbose=0)
```

If your data is so large that you cannot use `memmap`, things get more complicated. Let's start by writing a function to load the next batch (in real life, you would load the data from disk):

```
In [48]: def load_next_batch(batch_size):
    return X[np.random.choice(len(X), batch_size, replace=False)]
```

Now we can train the model by feeding it one batch at a time. We also need to implement multiple initializations and keep the model with the lowest inertia:

```
In [49]: np.random.seed(42)
```

```
In [50]: k = 5
n_init = 10
n_iterations = 100
batch_size = 100
init_size = 500 # more data for K-Means++ initialization
evaluate_on_last_n_iters = 10

best_kmeans = None

for init in range(n_init):
    minibatch_kmeans = MiniBatchKMeans(n_clusters=k, init_size=init_size)
    X_init = load_next_batch(init_size)
    minibatch_kmeans.partial_fit(X_init)

    minibatch_kmeans.sum_inertia_ = 0
    for iteration in range(n_iterations):
        X_batch = load_next_batch(batch_size)
        minibatch_kmeans.partial_fit(X_batch)
        if iteration >= n_iterations - evaluate_on_last_n_iters:
            minibatch_kmeans.sum_inertia_ += minibatch_kmeans.inertia_

    if (best_kmeans is None or
        minibatch_kmeans.sum_inertia_ < best_kmeans.sum_inertia_):
        best_kmeans = minibatch_kmeans
```

```
In [51]: best_kmeans.score(X)
```

```
Out[51]: -211.70999744411483
```

Mini-batch K-Means is much faster than regular K-Means:

```
In [52]: %timeit KMeans(n_clusters=5).fit(X)
```

```
47 ms ± 1.44 ms per loop (mean ± std. dev. of 7 runs, 10 loops each)
```

```
In [53]: %timeit MiniBatchKMeans(n_clusters=5).fit(X)
```

```
26.2 ms ± 2.41 ms per loop (mean ± std. dev. of 7 runs, 10 loops each)
```

That's *much* faster! However, its performance is often lower (higher inertia), and it keeps degrading as k increases. Let's plot the inertia ratio and the training time ratio between Mini-batch K-Means and regular K-Means:

```
In [54]: from timeit import timeit
```

```
In [55]: times = np.empty((100, 2))
inertias = np.empty((100, 2))
for k in range(1, 101):
    kmeans_ = KMeans(n_clusters=k, random_state=42)
    minibatch_kmeans = MiniBatchKMeans(n_clusters=k, random_state=42)
    print("\r{}/{}".format(k, 100), end="")
    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_
```

```
100/100
```

```
In [56]: plt.figure(figsize=(10,4))
```

```

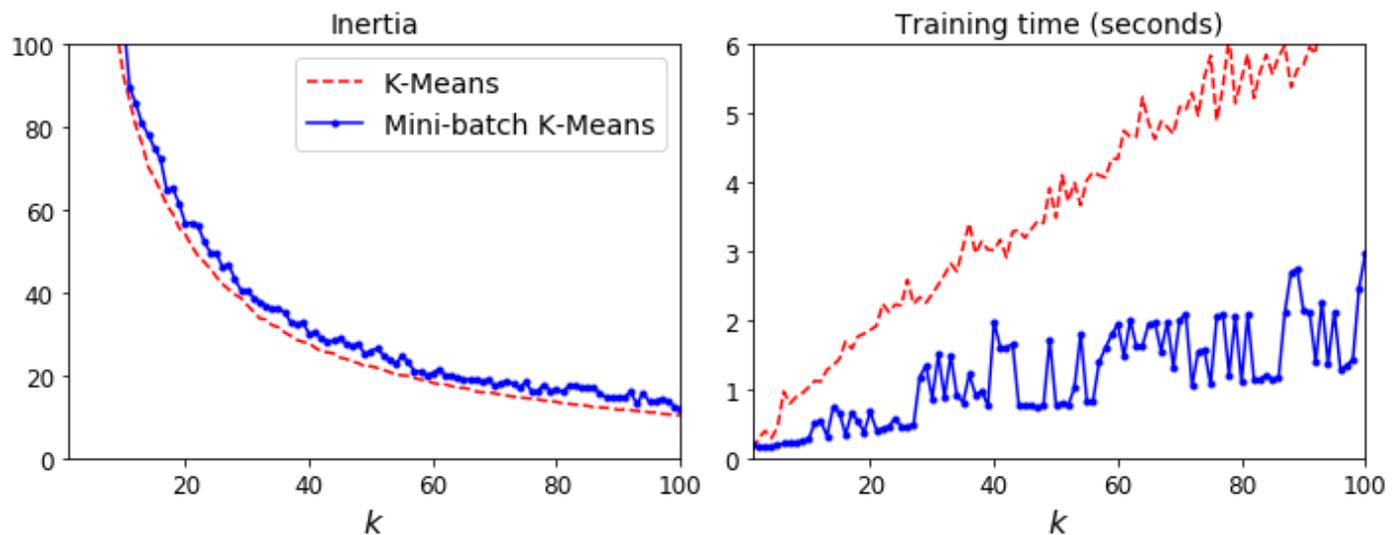
plt.subplot(121)
plt.plot(range(1, 101), inertias[:, 0], "r--", label="K-Means")
plt.plot(range(1, 101), inertias[:, 1], "b.-", label="Mini-batch K-Means")
plt.xlabel("$k$", fontsize=16)
plt.title("Inertia", fontsize=14)
plt.legend(fontsize=14)
plt.axis([1, 100, 0, 100])

plt.subplot(122)
plt.plot(range(1, 101), times[:, 0], "r--", label="K-Means")
plt.plot(range(1, 101), times[:, 1], "b.-", label="Mini-batch K-Means")
plt.xlabel("$k$", fontsize=16)
plt.title("Training time (seconds)", fontsize=14)
plt.axis([1, 100, 0, 6])

save_fig("minibatch_kmeans_vs_kmeans")
plt.show()

```

Saving figure minibatch_kmeans_vs_kmeans



Finding the optimal number of clusters

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

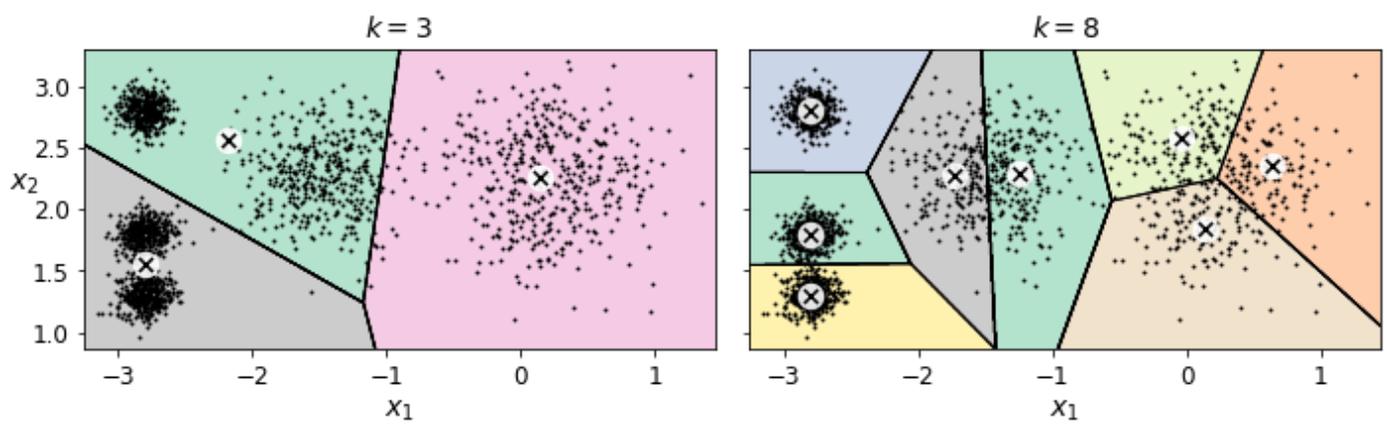
```

In [57]: kmeans_k3 = KMeans(n_clusters=3, random_state=42)
kmeans_k8 = KMeans(n_clusters=8, random_state=42)

plot_clusterer_comparison(kmeans_k3, kmeans_k8, X, "$k=3$", "$k=8$")
save_fig("bad_n_clusters_plot")
plt.show()

```

Saving figure bad_n_clusters_plot



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

```
In [58]: kmeans_k3.inertia_
```

```
Out[58]: 653.2167190021553
```

```
In [59]: kmeans_k8.inertia_
```

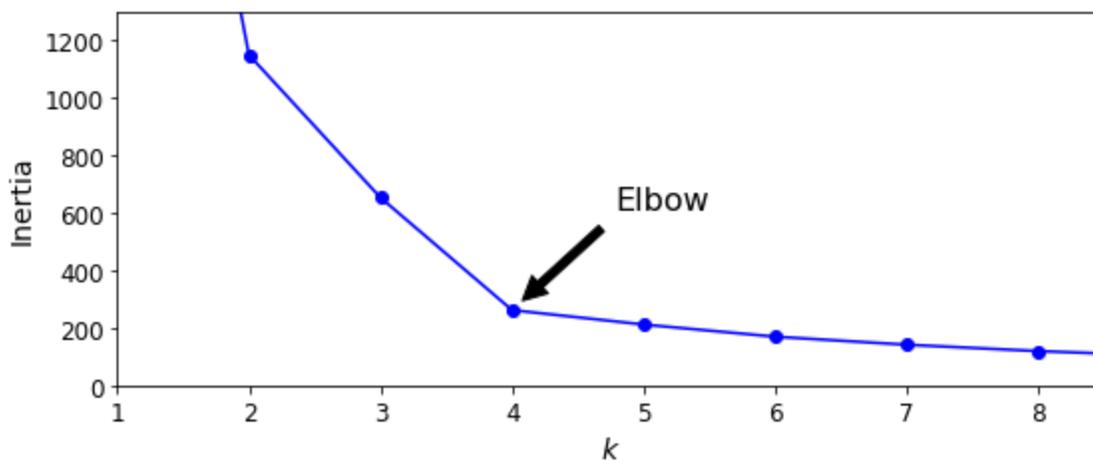
```
Out[59]: 119.11983416102879
```

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:

```
In [60]: kmeans_per_k = [KMeans(n_clusters=k, random_state=42).fit(X)
                        for k in range(1, 10)]
inertias = [model.inertia_ for model in kmeans_per_k]
```

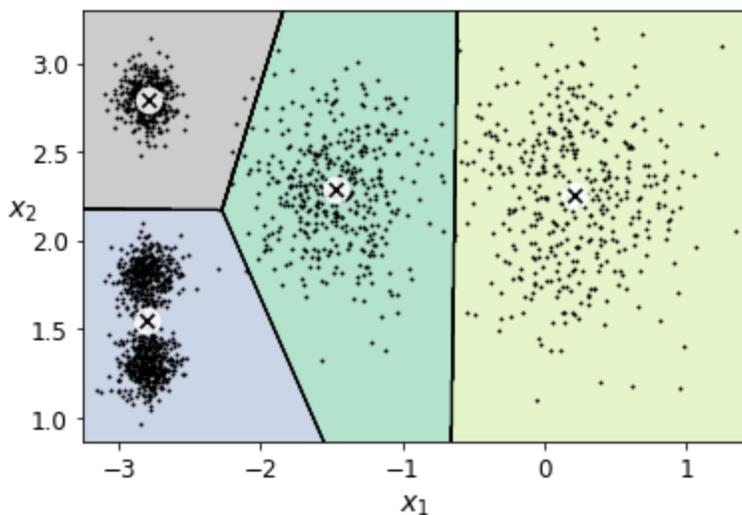
```
In [61]: plt.figure(figsize=(8, 3.5))
plt.plot(range(1, 10), inertias, "bo-")
plt.xlabel("$k$", fontsize=14)
plt.ylabel("Inertia", fontsize=14)
plt.annotate('Elbow',
            xy=(4, inertias[3]),
            xytext=(0.55, 0.55),
            textcoords='figure fraction',
            fontsize=16,
            arrowprops=dict(facecolor='black', shrink=0.1)
            )
plt.axis([1, 8.5, 0, 1300])
save_fig("inertia_vs_k_plot")
plt.show()
```

Saving figure inertia_vs_k_diagram



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.

```
In [62]: 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 :

```
In [63]: from sklearn.metrics import silhouette_score
```

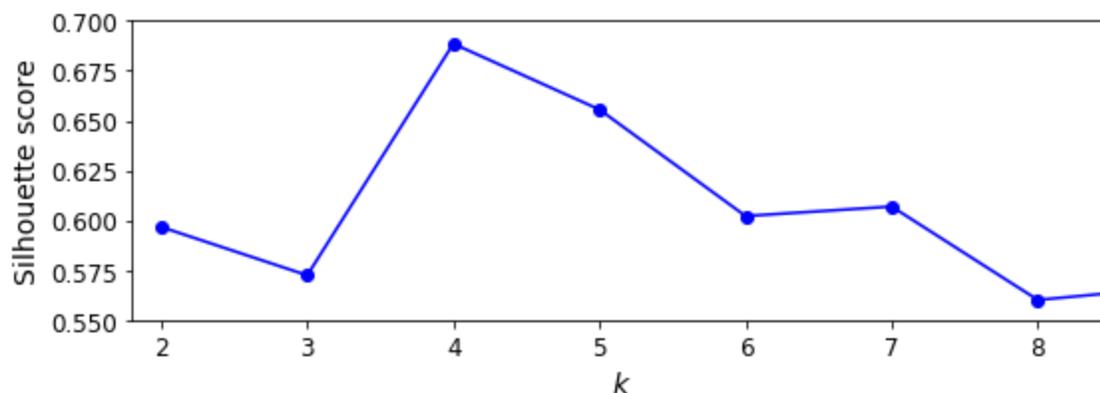
```
In [64]: silhouette_score(X, kmeans.labels_)
```

Out[64]: 0.655517642572828

```
In [65]: silhouette_scores = [silhouette_score(X, model.labels_)  
                             for model in kmeans_per_k[1:]]
```

```
In [66]: plt.figure(figsize=(8, 3))  
plt.plot(range(2, 10), silhouette_scores, "bo-")  
plt.xlabel("$k$", fontsize=14)  
plt.ylabel("Silhouette score", fontsize=14)  
plt.axis([1.8, 8.5, 0.55, 0.7])  
save_fig("silhouette_score_vs_k_plot")  
plt.show()
```

Saving figure silhouette_score_vs_k_plot



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*:

```
In [67]: 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 = mpl.cm.Spectral(i / k)  
        plt.fill_betweenx(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")

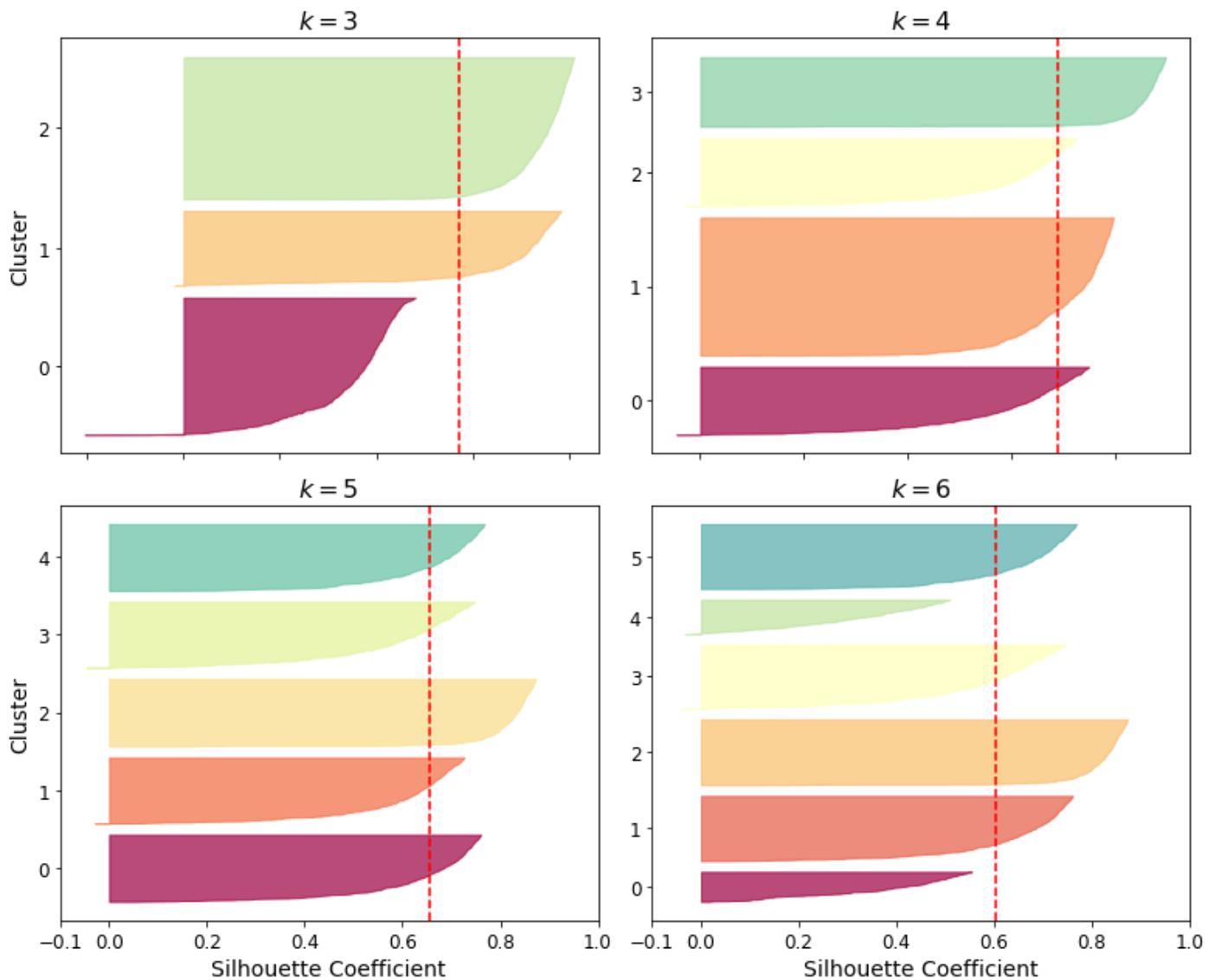
if k in (5, 6):
    plt.gca().set_xticks([-0.1, 0, 0.2, 0.4, 0.6, 0.8, 1])
    plt.xlabel("Silhouette Coefficient")
else:
    plt.tick_params(labelbottom=False)

plt.axvline(x=silhouette_scores[k - 2], color="red", linestyle="--")
plt.title("$k={}$".format(k), fontsize=16)

save_fig("silhouette_analysis_plot")
plt.show()

```

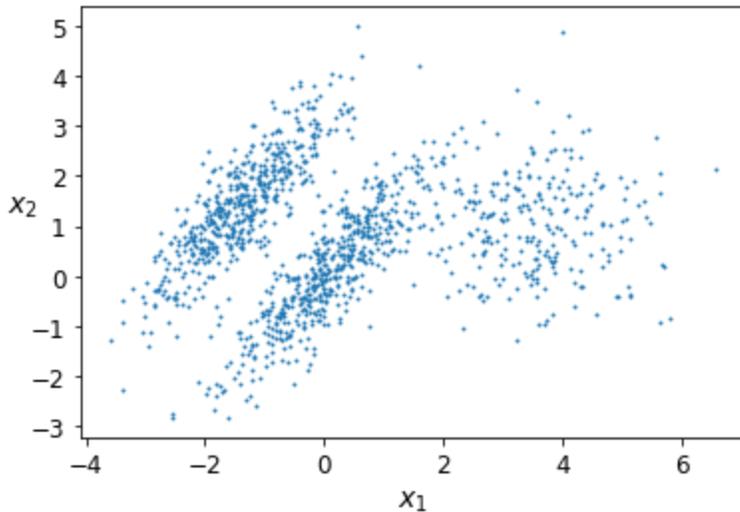
Saving figure silhouette_analysis_diagram



Limits of K-Means

```
In [68]: 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]
```

```
In [69]: plot_clusters(X)
```



```
In [70]: kmeans_good = KMeans(n_clusters=3, init=np.array([[ -1.5, 2.5], [ 0.5, 0], [ 4, 0]]), n_init=10, random_state=42)
kmeans_bad = KMeans(n_clusters=3, random_state=42)
kmeans_good.fit(X)
kmeans_bad.fit(X)
```

```
Out[70]: KMeans(algorithm='auto', copy_x=True, init='k-means++', max_iter=300, n_clusters=3, n_init=10, n_jobs=None, precompute_distances='auto', random_state=42, tol=0.0001, verbose=0)
```

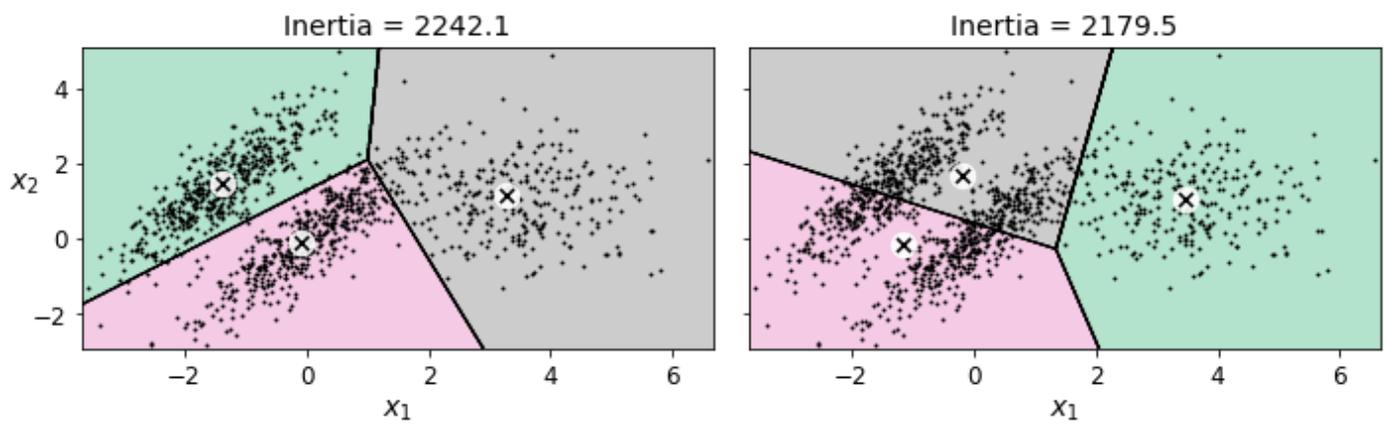
```
In [71]: plt.figure(figsize=(10, 3.2))

plt.subplot(121)
plot_decision_boundaries(kmeans_good, X)
plt.title("Inertia = {:.1f}".format(kmeans_good.inertia_), fontsize=14)

plt.subplot(122)
plot_decision_boundaries(kmeans_bad, X, show_ylabels=False)
plt.title("Inertia = {:.1f}".format(kmeans_bad.inertia_), fontsize=14)

save_fig("bad_kmeans_plot")
plt.show()
```

Saving figure bad_kmeans_diagram



Using clustering for image segmentation

```
In [72]: # Download the ladybug image
images_path = os.path.join(PROJECT_ROOT_DIR, "images", "unsupervised_learning")
os.makedirs(images_path, exist_ok=True)
DOWNLOAD_ROOT = "https://raw.githubusercontent.com/ageron/handson-ml2/master/"
filename = "ladybug.png"
print("Downloading", filename)
url = DOWNLOAD_ROOT + "images/unsupervised_learning/" + filename
urllib.request.urlretrieve(url, os.path.join(images_path, filename))
```

```
In [73]: from matplotlib.image import imread
image = imread(os.path.join(images_path, filename))
image.shape
```

Out[73]: (533, 800, 3)

```
In [74]: X = image.reshape(-1, 3)
kmeans = KMeans(n_clusters=8, random_state=42).fit(X)
segmented_img = kmeans.cluster_centers_[kmeans.labels_]
segmented_img = segmented_img.reshape(image.shape)
```

```
In [75]: segmented_imgs = []
n_colors = (10, 8, 6, 4, 2)
for n_clusters in n_colors:
    kmeans = KMeans(n_clusters=n_clusters, random_state=42).fit(X)
    segmented_img = kmeans.cluster_centers_[kmeans.labels_]
    segmented_imgs.append(segmented_img.reshape(image.shape))
```

```
In [76]: plt.figure(figsize=(10,5))
plt.subplots_adjust(wspace=0.05, hspace=0.1)

plt.subplot(231)
plt.imshow(image)
plt.title("Original image")
plt.axis('off')

for idx, n_clusters in enumerate(n_colors):
    plt.subplot(232 + idx)
    plt.imshow(segmented_imgs[idx])
    plt.title("{} colors".format(n_clusters))
    plt.axis('off')

save_fig('image_segmentation_diagram', tight_layout=False)
plt.show()
```



Using Clustering for Preprocessing

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.

```
In [77]: from sklearn.datasets import load_digits
```

```
In [78]: X_digits, y_digits = load_digits(return_X_y=True)
```

Let's split it into a training set and a test set:

```
In [79]: from sklearn.model_selection import train_test_split
```

```
In [80]: X_train, X_test, y_train, y_test = train_test_split(X_digits, y_digits, random_state=42)
```

Now let's fit a Logistic Regression model and evaluate it on the test set:

```
In [81]: from sklearn.linear_model import LogisticRegression
```

```
In [82]: log_reg = LogisticRegression(multi_class="ovr", solver="lbfgs", max_iter=5000, random_st
log_reg.fit(X_train, y_train)
```

```
Out[82]: LogisticRegression(C=1.0, class_weight=None, dual=False, fit_intercept=True,
intercept_scaling=1, max_iter=5000, multi_class='ovr',
n_jobs=None, penalty='l2', random_state=42, solver='lbfgs',
tol=0.0001, verbose=0, warm_start=False)
```

```
In [83]: log_reg.score(X_test, y_test)
```

```
Out[83]: 0.9688888888888889
```

Okay, that's our baseline: 96.89% accuracy. Let's see if we can do better by using K-Means as a preprocessing step. We will create a pipeline that will first cluster the training set into 50 clusters and replace the images with their distances to the 50 clusters, then apply a logistic regression model:

```
In [84]: from sklearn.pipeline import Pipeline
```

```
In [85]: pipeline = Pipeline([
    ("kmeans", KMeans(n_clusters=50, random_state=42)),
    ("log_reg", LogisticRegression(multi_class="ovr", solver="lbfgs", max_iter=5000, ran
])
pipeline.fit(X_train, y_train)
```

```
Out[85]: Pipeline(memory=None,
    steps=[('kmeans', KMeans(algorithm='auto', copy_x=True, init='k-means++', max_iter
=300,
    n_clusters=50, n_init=10, n_jobs=None, precompute_distances='auto',
    random_state=42, tol=0.0001, verbose=0)), ('log_reg', LogisticRegression(C=1.0, cla
ss_weight=None, dual=False, fit_intercept=True,
    intercept_scaling=1, max_iter=5000, multi_class='ovr',
    n_jobs=None, penalty='l2', random_state=42, solver='lbfgs',
    tol=0.0001, verbose=0, warm_start=False))])
```

```
In [86]: pipeline.score(X_test, y_test)
```

```
Out[86]: 0.9777777777777777
```

```
In [87]: 1 - (1 - 0.977777) / (1 - 0.968888)
```

```
Out[87]: 0.28570969400874346
```

How about that? We reduced the error rate by over 28%! But we chose the number of clusters k completely arbitrarily, we can surely do better. Since K-Means is just a preprocessing step in a classification pipeline, finding a good value for k is much simpler than earlier: there's no need to perform silhouette analysis or minimize the inertia, the best value of k is simply the one that results in the best classification performance.

```
In [88]: from sklearn.model_selection import GridSearchCV
```

```
In [89]: param_grid = dict(kmeans__n_clusters=range(2, 100))
grid_clf = GridSearchCV(pipeline, param_grid, cv=3, verbose=2)
grid_clf.fit(X_train, y_train)
```

Fitting 3 folds for each of 98 candidates, totalling 294 fits

```
[CV] kmeans__n_clusters=2 .....
[CV] ..... kmeans__n_clusters=2, total= 0.1s
[CV] kmeans__n_clusters=2 .....
```

```
[Parallel(n_jobs=1)]: Using backend SequentialBackend with 1 concurrent workers.
[Parallel(n_jobs=1)]: Done 1 out of 1 | elapsed: 0.1s remaining: 0.0s
```



```
[CV] kmeans__n_clusters=99 .....  
[CV] ..... kmeans__n_clusters=99, total= 3.5s  
[Parallel(n_jobs=1)]: Done 294 out of 294 | elapsed: 13.9min finished
```

```
Out[89]: GridSearchCV(cv=3, error_score='raise-deprecating',  
    estimator=Pipeline(memory=None,  
    steps=[('kmeans', KMeans(algorithm='auto', copy_x=True, init='k-means++', max_iter  
=300,  
    n_clusters=50, n_init=10, n_jobs=None, precompute_distances='auto',  
    random_state=42, tol=0.0001, verbose=0)), ('log_reg', LogisticRegression(C=1.0, cla  
ss_weight=None, ...penalty='l2', random_state=42, solver='lbfgs',  
    tol=0.0001, verbose=0, warm_start=False))]),  
    fit_params=None, iid='warn', n_jobs=None,  
    param_grid={'kmeans__n_clusters': range(2, 100)},  
    pre_dispatch='2*n_jobs', refit=True, return_train_score='warn',  
    scoring=None, verbose=2)
```

```
In [90]: grid_clf.best_params_
```

```
Out[90]: {'kmeans__n_clusters': 99}
```

```
In [91]: grid_clf.score(X_test, y_test)
```

```
Out[91]: 0.9822222222222222
```

The performance improved most with $k = 99$, so 99 it is.

Clustering for Semi-supervised Learning

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

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

```
In [92]: n_labeled = 50
```

```
In [93]: log_reg = LogisticRegression(multi_class="ovr", solver="lbfgs", random_state=42)  
log_reg.fit(X_train[:n_labeled], y_train[:n_labeled])  
log_reg.score(X_test, y_test)
```

```
Out[93]: 0.8333333333333334
```

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:

```
In [94]: k = 50
```

```
In [95]: kmeans = KMeans(n_clusters=k, random_state=42)  
X_digits_dist = kmeans.fit_transform(X_train)  
representative_digit_idx = np.argmin(X_digits_dist, axis=0)  
X_representative_digits = X_train[representative_digit_idx]
```

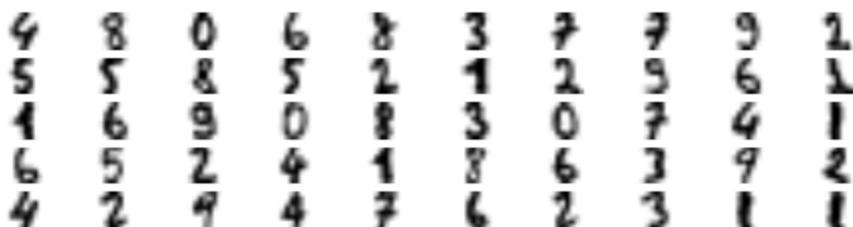
Now let's plot these representative images and label them manually:

```
In [96]: 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')
```

```
save_fig("representative_images_diagram", tight_layout=False)
plt.show()
```

Saving figure representative_images_diagram



```
In [97]: y_representative_digits = np.array([
    4, 8, 0, 6, 8, 3, 7, 7, 9, 2,
    5, 5, 8, 5, 2, 1, 2, 9, 6, 1,
    1, 6, 9, 0, 8, 3, 0, 7, 4, 1,
    6, 5, 2, 4, 1, 8, 6, 3, 9, 2,
    4, 2, 9, 4, 7, 6, 2, 3, 1, 1])
```

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:

```
In [98]: log_reg = LogisticRegression(multi_class="ovr", solver="lbfgs", max_iter=5000, random_state=42)
log_reg.fit(X_representative_digits, y_representative_digits)
log_reg.score(X_test, y_test)
```

```
Out[98]: 0.9222222222222223
```

Wow! We jumped from 83.3% accuracy to 92.2%, 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?

```
In [99]: y_train_propagated = np.empty(len(X_train), dtype=np.int32)
for i in range(k):
    y_train_propagated[kmeans.labels_==i] = y_representative_digits[i]
```

```
In [100]: log_reg = LogisticRegression(multi_class="ovr", solver="lbfgs", max_iter=5000, random_state=42)
log_reg.fit(X_train, y_train_propagated)
```

```
Out[100]: LogisticRegression(C=1.0, class_weight=None, dual=False, fit_intercept=True,
    intercept_scaling=1, max_iter=5000, multi_class='ovr',
    n_jobs=None, penalty='l2', random_state=42, solver='lbfgs',
    tol=0.0001, verbose=0, warm_start=False)
```

```
In [101]: log_reg.score(X_test, y_test)
```

```
Out[101]: 0.9333333333333333
```

We got a tiny little accuracy boost. Better than nothing, but we should probably have propagated the labels only to the instances closest to the centroid, because by propagating to the full cluster, we have certainly included some outliers. Let's only propagate the labels to the 20th percentile closest to the centroid:

```
In [102... percentile_closest = 20

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
```

```
In [103... partially_propagated = (X_cluster_dist != -1)
X_train_partially_propagated = X_train[partially_propagated]
y_train_partially_propagated = y_train_propagated[partially_propagated]
```

```
In [104... log_reg = LogisticRegression(multi_class="ovr", solver="lbfgs", max_iter=5000, random_st
log_reg.fit(X_train_partially_propagated, y_train_partially_propagated)
```

```
Out[104... LogisticRegression(C=1.0, class_weight=None, dual=False, fit_intercept=True,
intercept_scaling=1, max_iter=5000, multi_class='ovr',
n_jobs=None, penalty='l2', random_state=42, solver='lbfgs',
tol=0.0001, verbose=0, warm_start=False)
```

```
In [105... log_reg.score(X_test, y_test)
```

```
Out[105... 0.94
```

Nice! With just 50 labeled instances (just 5 examples per class on average!), we got 94% performance, which is pretty close to the performance of logistic regression on the fully labeled *digits* dataset (which was 96.9%).

This is because the propagated labels are actually pretty good: their accuracy is very close to 99%:

```
In [106... np.mean(y_train_partially_propagated == y_train[partially_propagated])
```

```
Out[106... 0.9896907216494846
```

You could now do a few iterations of *active learning*:

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

DBSCAN

```
In [107... from sklearn.datasets import make_moons
```

```
In [108... X, y = make_moons(n_samples=1000, noise=0.05, random_state=42)
```

```
In [109... from sklearn.cluster import DBSCAN
```

```
In [110... dbscan = DBSCAN(eps=0.05, min_samples=5)
dbscan.fit(X)
```

```
Out[110... DBSCAN(algorithm='auto', eps=0.05, leaf_size=30, metric='euclidean',
metric_params=None, min_samples=5, n_jobs=None, p=None)
```

```
In [111... dbscan.labels_[:10]
```

```
Out[111... array([ 0,  2, -1, -1,  1,  0,  0,  0,  2,  5])
```

```
In [112... len(dbscan.core_sample_indices_)
```

```
Out[112... 808
```

```
In [113... dbscan.core_sample_indices_[:10]
```

```
Out[113... array([ 0,  4,  5,  6,  7,  8, 10, 11, 12, 13])
```

```
In [114... dbscan.components_[:3]
```

```
Out[114... array([[ -0.02137124,  0.40618608],
          [-0.84192557,  0.53058695],
          [ 0.58930337, -0.32137599]])
```

```
In [115... np.unique(dbscan.labels_)
```

```
Out[115... array([-1,  0,  1,  2,  3,  4,  5,  6])
```

```
In [116... dbscan2 = DBSCAN(eps=0.2)
dbscan2.fit(X)
```

```
Out[116... DBSCAN(algorithm='auto', eps=0.2, leaf_size=30, metric='euclidean',
metric_params=None, min_samples=5, n_jobs=None, p=None)
```

```
In [117... 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], marke
    if show_xlabels:
        plt.xlabel("$x_1$", fontsize=14)
    else:
        plt.tick_params(labelbottom=False)
    if show_ylabels:
        plt.ylabel("$x_2$", fontsize=14, rotation=0)
    else:
        plt.tick_params(labelleft=False)
    plt.title("eps={:.2f}, min_samples={}".format(dbscan.eps, dbscan.min_samples), fonts
```

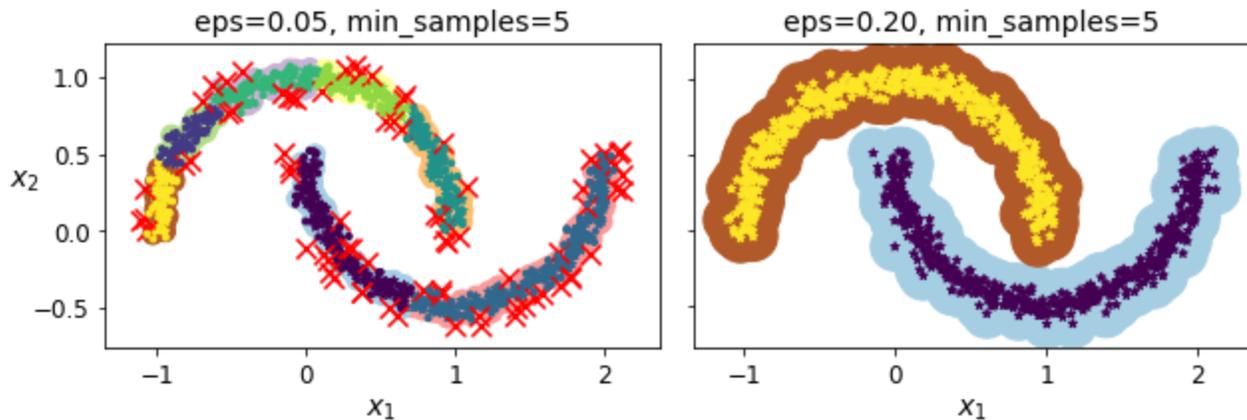
```
In [118... plt.figure(figsize=(9, 3.2))

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

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

save_fig("dbscan_plot")
plt.show()
```

Saving figure dbscan_diagram



```
In [119... dbscan = dbscan2
```

```
In [120... from sklearn.neighbors import KNeighborsClassifier
```

```
In [121... knn = KNeighborsClassifier(n_neighbors=50)
knn.fit(dbscan.components_, dbscan.labels_[dbscan.core_sample_indices_])
```

```
Out[121... KNeighborsClassifier(algorithm='auto', leaf_size=30, metric='minkowski',
metric_params=None, n_jobs=None, n_neighbors=50, p=2,
weights='uniform')
```

```
In [122... X_new = np.array([[-0.5, 0], [0, 0.5], [1, -0.1], [2, 1]])
knn.predict(X_new)
```

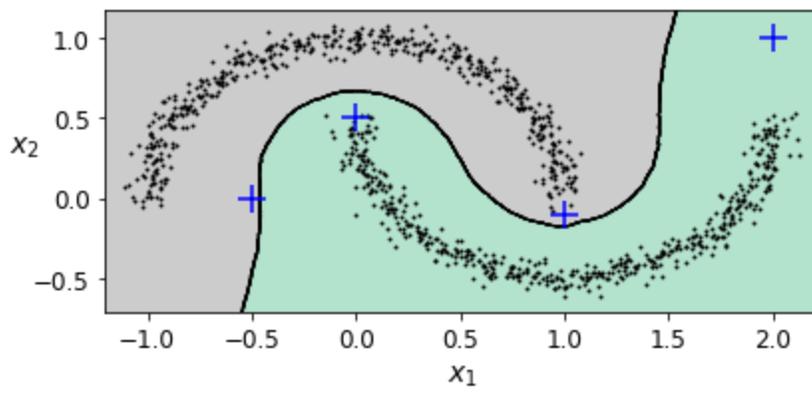
```
Out[122... array([1, 0, 1, 0])
```

```
In [123... knn.predict_proba(X_new)
```

```
Out[123... array([[0.18, 0.82],
[1. , 0. ],
[0.12, 0.88],
[1. , 0. ]])
```

```
In [124... 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)
save_fig("cluster_classification_plot")
plt.show()
```

Saving figure cluster_classification_diagram



```
In [125... 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()
```

```
Out[125... array([-1,  0,  1, -1])
```