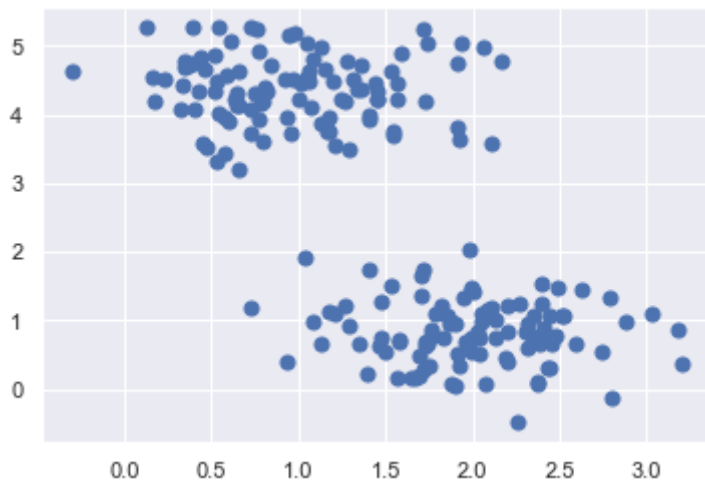


```
In [1]: import numpy as np
import matplotlib.pyplot as plt
import seaborn as sns
sns.set()
```

## Clustering intuition: Putting data into groups *without* labels (i.e., unsupervised learning)

- Abstract idea: "Put like data together in a group" or "like with like"
- Analogy: Consider a small "city" of people.
  - Each point represents a person
  - Friendships are formed entirely based on how close they live to each other
  - Could you put these people into communities?

```
In [2]: from sklearn.datasets import make_blobs
X, y_true = make_blobs(n_samples=200, centers=2,
                        cluster_std=0.50, random_state=0)
plt.scatter(X[:, 0], X[:, 1], s=50);
```

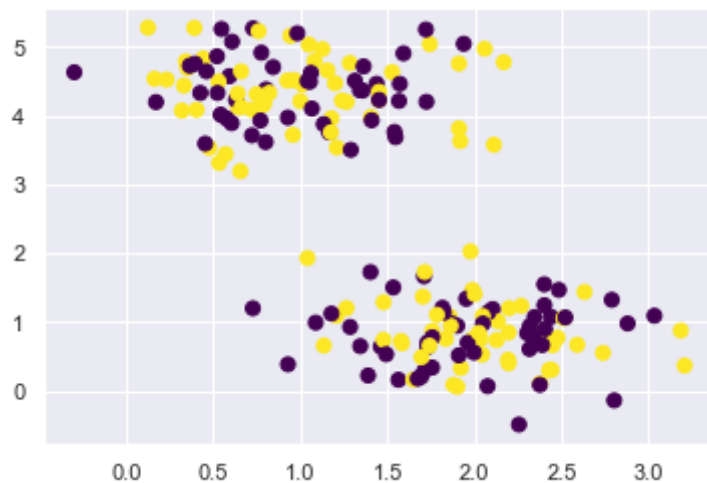


## Math: How do we formalize what we did visually?

- Let's assume for now that we know there are exactly *two* clusters
- How can we assign each point to a cluster?
- Naive idea: Randomly assign points to each cluster

```
In [3]: from sklearn.utils import check_random_state
def get_random_assignment(random_state=None):
    rng = check_random_state(random_state)
    y = rng.randint(2, size=X.shape[0])
    return y
y_rand = get_random_assignment(random_state=0)
plt.scatter(X[:, 0], X[:, 1], c=y_rand, s=50, cmap='viridis')
```

Out[3]: <matplotlib.collections.PathCollection at 0x7fed71ff98e0>



**This clustering "looks" quite bad.**

**How can we formalize whether a particular assignment is good or bad?**

- Intuition from analogy: People in a communities will be as close to each other as possible.
- Take average distance between each point in a cluster to every other point in the **same** cluster.
- Sum over all communities.

In [ ]:

**Math: Sum of squares per cluster objective function**

- Average sum of squares between each point in the same cluster

$$C_j = \{x \in \mathcal{X} : y = j\}$$

$$\sum_{j=1}^k \frac{1}{2|C_j|} \sum_{x \in C_j, z \in C_j} \|x - z\|_2^2$$

```
In [4]: from sklearn.metrics import pairwise_distances
# Using vectorized and list comprehensions computation
def objective(X, y):
    y_vals = np.unique(y)
    def inner(yv):
        sel = (y==yv) # boolean array
        Xj = X[sel, :]
        n_community = np.sum(sel)
        community_sum = np.sum(pairwise_distances(Xj, Xj)**2)
        return community_sum / (2*n_community)
    return np.sum([inner(yv) for yv in y_vals])

print(objective(X, y_rand))
```

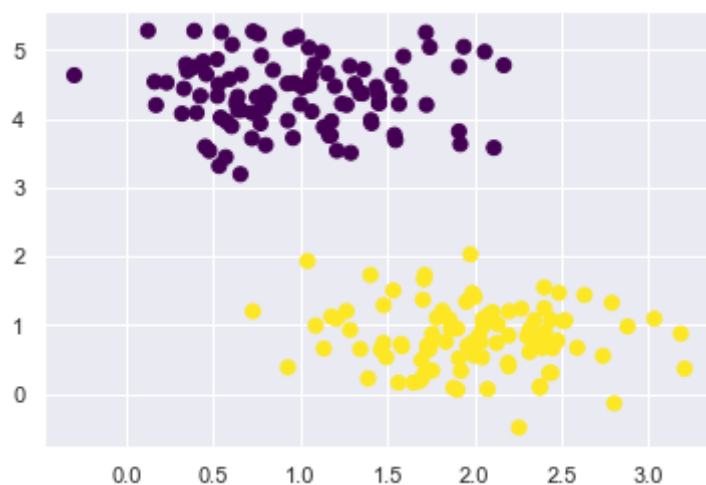
767.2572924351311

## Intuition sanity check, does visual clustering solution have a low value?

```
In [5]: print(objective(X, y_true))
plt.scatter(X[:, 0], X[:, 1], c=y_true, s=50, cmap='viridis')
```

94.67363954089785

Out[5]: <matplotlib.collections.PathCollection at 0x7fed7366e3d0>



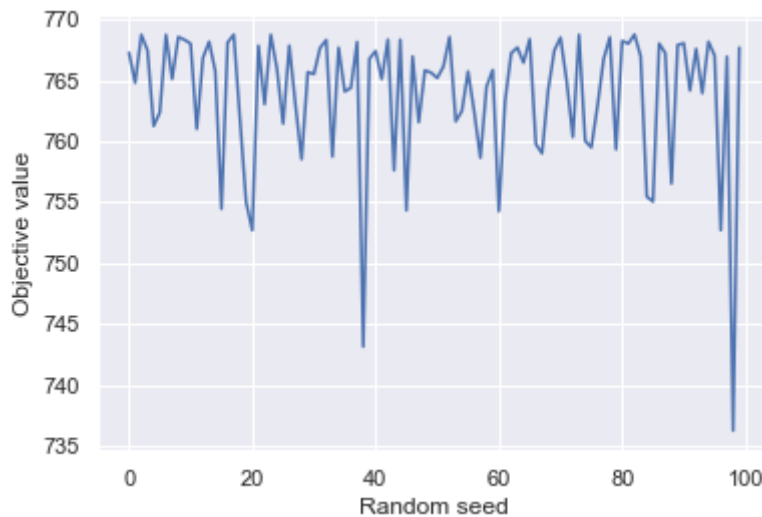
## Optimization: Minimize objective over possible community assignments

$$\arg \min_{C_1, C_2} \sum_{j=1}^k \frac{1}{2|C_j|} \sum_{x \in C_j, z \in C_j} \|x - z\|_2^2$$

- Naively, we could just enumerate all possibilities
- Let's try several random combinations

```
In [6]: rand_obj = np.nan * np.ones(100)
for seed in range(rand_obj.shape[0]):
    y_rand = get_random_assignment(random_state=seed)
    rand_obj[seed] = objective(X, y_rand)
    #print('Seed = %2d, Objective = %g' % (seed, obj))
plt.plot(rand_obj)
plt.xlabel('Random seed')
plt.ylabel('Objective value')
#plt.scatter(X[:, 0], X[:, 1], c=y_rand, s=50, cmap='viridis')
```

Out[6]: Text(0, 0.5, 'Objective value')



## How many possible assignments are there?

In terms of the number of samples  $n$  and the number of communities  $k$

```
In [7]: n_samples = X.shape[0]
n_communities = 2
n_assignments = n_communities ** (n_samples-1)
print('For %d samples and %d communities, there are %d possible assignments
      % (n_samples, n_communities, n_assignments))
print('Or in exponential notation: %g possible assignments' % n_assignments)
```

For 200 samples and 2 communities, there are 8034690221294951377709810461  
70581301261101496891396417650688 possible assignments  
Or in exponential notation: 8.03469e+59 possible assignments

**Some perspective: Fastest super computer is 200 petaflops  
=  $2 * 10^{17}$  operations per second**

```
In [8]: ops = 2 * (10 ** 17)
print(ops)
compute_time = n_assignments / ops
compute_time_years = compute_time / 60 / 60 / 24 / 365
print('Years of compute time: %d' % compute_time_years)
```

```
2000000000000000000
Years of compute time: 127389177785625178899305200808361984
```

## Clearly, not a good way to optimize

## Let's consider a *equivalent* optimization

Can you figure out what these two equations mean?

$$\mu_j \equiv \frac{1}{|C_j|} \sum_{x \in C_j} x_i$$

$$\arg \min_{C_1, C_2, \dots, C_k} \sum_{j=1}^k \sum_{x \in C_j} \|x - \mu_j\|_2^2$$

```
In [9]: # Just space holder
```

## Consider an equivalent optimization via community *representatives*

- Analogy
  - Instead of measuring from each person to every other person in the same community, measure between a person and an ideal "representative" of each community, who is at the center of everyone.
  - Representative can move freely.
- If the cluster assignments  $C_j$  are fixed, then the position of the "representative", denoted by  $\mu_j$  is defines as the mean/average point:

$$\mu_j \equiv \frac{1}{|C_j|} \sum_{x \in C_j} x_i$$

- This leads to the following **equivalent** minimization:

$$\arg \min_{C_1, \dots, C_k} \sum_{j=1}^k \frac{1}{2|C_j|} \sum_{x \in C_j, z \in C_j} \|x - z\|_2^2 \equiv \arg \min_{C_1, \dots, C_k} \sum_{j=1}^k \sum_{x \in C_j} \|x - \mu_j\|_2^2$$

$$\equiv \arg \min_{C_1, C_2, \dots, C_k} \sum_{j=1}^k \sum_{x \in C_j} \left\| x - \frac{1}{|C_j|} \sum_{x \in C_j} x_i \right\|_2^2$$

(Derivation of equivalence can be seen at

[https://www.math.ucdavis.edu/~strohmer/courses/180BigData/180lecture\\_kmeans.pdf](https://www.math.ucdavis.edu/~strohmer/courses/180BigData/180lecture_kmeans.pdf)

([https://www.math.ucdavis.edu/~strohmer/courses/180BigData/180lecture\\_kmeans.pdf](https://www.math.ucdavis.edu/~strohmer/courses/180BigData/180lecture_kmeans.pdf).)

## Implement the objective of the equivalent optimization

$$\arg \min_{C_1, C_2, \dots, C_k} \sum_{j=1}^k \sum_{x \in C_j} \|x - \mu_j\|_2^2$$

```
In [10]: def objective2(X, y):
k = len(np.unique(y))
out = 0
for j in range(k):
    sel = (y==j) # boolean array
    Xj = X[sel, :]
    mu_j = np.mean(Xj, axis=0)
    dist_to_mu = np.sqrt(np.sum((Xj - mu_j)**2, axis=0))
    out += np.sum(dist_to_mu**2)
return out

print('Quick sanity check that objective corresponds to visual understanding')
print('Objective random', objective2(X, y_rand))
print('Objective visual', objective2(X, y_true))
```

Quick sanity check that objective corresponds to visual understanding

Objective random 767.679899871254

Objective visual 94.67363954089788

## During optimization, suppose cluster centers can move and are parameters

$$\arg \min_{C_1, \dots, C_k, \mu_1, \dots, \mu_k} \sum_{j=1}^k \sum_{x \in C_j} \|x - \mu_j\|_2^2$$

- In this "unsettled" state

1. Intuition: People will join the community of their closest *representative*  $\mu_j$ .

Math:

$$y_i = \arg \min_{j \in \{1, 2, \dots, k\}} \|x_i - \mu_j\|_2$$

2. Intuition: The representative will move to the center of it's current community.

Math:

$$\mu_j = \frac{1}{|C_j|} \sum_{x \in C_j} x_i$$

```
In [11]: def objective3(X, y, mu_array):
          k = len(np.unique(y))
          out = 0
          for j in range(k):
              sel = (y==j) # boolean array
              Xj = X[sel, :]
              mu_j = mu_array[j, :]
              dist_to_mu = np.sqrt(np.sum((Xj - mu_j)**2, axis=0))
              out += np.sum(dist_to_mu**2)
          return out
```

## (1) Assign points to closest center / representative

```
In [12]: mu_array = np.array([[0, 1], [1, 0]])
          print(objective3(X, y_rand, mu_array))

          # Assign people
          def best_assignment(X, mu_array):
              y_best = np.argmin(pairwise_distances(X, mu_array), axis=1)
              return y_best

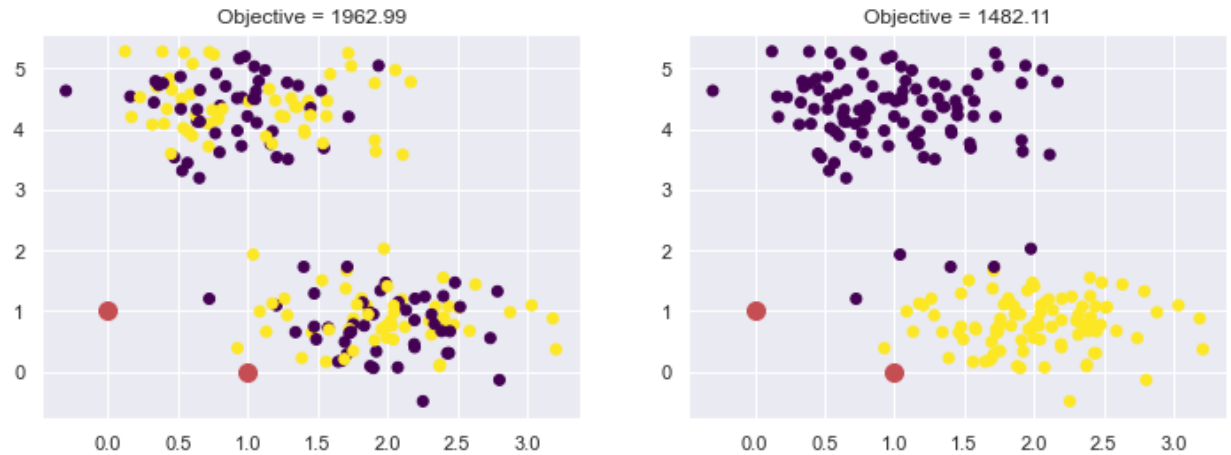
          y_new = best_assignment(X, mu_array)
          print(objective3(X, y_new, mu_array))
```

```
1962.992539917816
1482.1076321431726
```

**Make simple function for plotting (use ax as argument)**

```
In [13]: def plot_clustering(X, y, mu_array, ax=None):
    if ax is None:
        ax = plt.gca()
    ax.plot(mu_array[:, 0], mu_array[:, 1], 'ro', markersize=10)
    ax.scatter(X[:, 0], X[:, 1], c=y, cmap='viridis')
    ax.set_title('Objective = %g' % objective3(X, y, mu_array))

fig, axes = plt.subplots(1, 2, figsize=(12, 4))
for ycur, ax in zip([y_rand, y_new], axes):
    plot_clustering(X, ycur, mu_array, ax=ax)
```



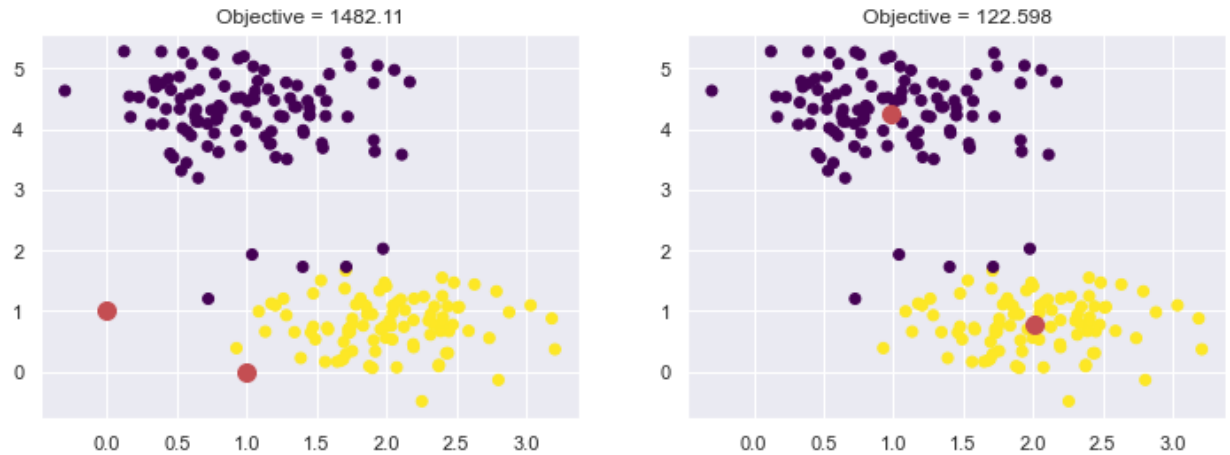
## (2) Update center based on new cluster assignments

(Analogy: Move representative to the center of it's cluster.)



```
In [14]: def recenter(X, y):
    return np.array([
        np.mean(X[y==yv, :], axis=0)
        for yv in np.unique(y)
    ])
mu_array_new = recenter(X, y_new)

fig, axes = plt.subplots(1, 2, figsize=(12, 4))
for m, ax in zip([mu_array, mu_array_new], axes):
    plot_clustering(X, y_new, m, ax=ax)
```



## What do you think you should do next?

```
In [15]: # Program kmeans
def kmeans_alg(X, maxiter=100, random_state=None):
    rng = check_random_state(random_state)

    # Initialize with random points in X
    rand_idx = rng.permutation(X.shape[0])
    mu_array = X[rand_idx[:2], :]
    y = get_random_assignment(random_state=rng)

    for i in range(maxiter):
        # Get new best assignment
        y_old = y # Save old assignment matrix
        y = best_assignment(X, mu_array)

        # Recenter / compute cluster mean
        mu_array = recenter(X, y)

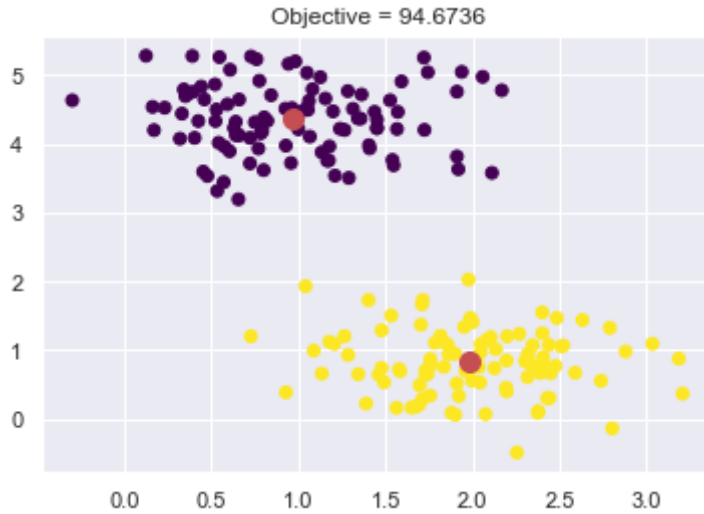
        # Check convergence
        if y_old is not None and np.all(y == y_old):
            print('Converged after %d iteration' % i)
            break
    return y, mu_array
```

```
In [16]: y_kmeans, mu_kmeans = kmeans_alg(X, maxiter=100, random_state=0)

plt.plot(mu_kmeans[:, 0], mu_kmeans[:, 1], 'ro', markersize=10)
plt.scatter(X[:, 0], X[:, 1], c=y_kmeans, cmap='viridis')
plt.title('Objective = %g' % objective3(X, y_kmeans, mu_kmeans))
```

Converged after 3 iteration

```
Out[16]: Text(0.5, 1.0, 'Objective = 94.6736')
```



**Let's inspect the underlying operation by splitting the iteration**

```

In [17]: # Program kmeans
def kmeans_alg(X, maxiter=100, random_state=None):
    rng = check_random_state(random_state)

    # Initialize with random points in X
    rand_idx = rng.permutation(X.shape[0])
    mu_array = X[rand_idx[:2], :]
    y = get_random_assignment(random_state=rng)

    for i in range(int(2*maxiter)): #CHANGED
        if i % 2 == 0: #CHANGED
            # Get new best assignment
            y_old = y # Save old assignment matrix
            y = best_assignment(X, mu_array)
        else: #CHANGED
            # Recenter / compute cluster mean
            mu_array = recenter(X, y)

        # Check convergence
        if y_old is not None and np.all(y == y_old):
            print('Converged after %d iteration' % (i/2)) #CHANGED
            break
    return y, mu_array

```

```

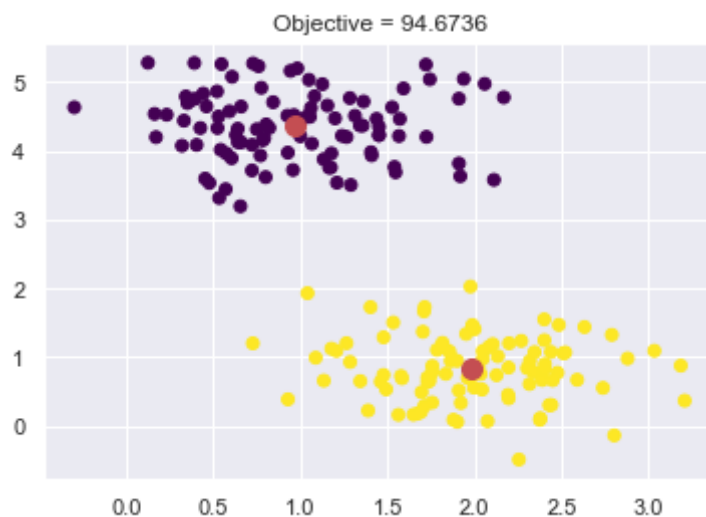
In [18]: y_kmeans, mu_kmeans = kmeans_alg(X, maxiter=4, random_state=0)

plt.plot(mu_kmeans[:, 0], mu_kmeans[:, 1], 'ro', markersize=10)
plt.scatter(X[:, 0], X[:, 1], c=y_kmeans, cmap='viridis')
plt.title('Objective = %g' % objective3(X, y_kmeans, mu_kmeans))

```

Converged after 3 iteration

Out[18]: Text(0.5, 1.0, 'Objective = 94.6736')



**Introducing scikit-learn's  
sklearn.cluster.KMeans**

- Documentation: <https://scikit-learn.org/stable/modules/generated/sklearn.cluster.KMeans.html> (<https://scikit-learn.org/stable/modules/generated/sklearn.cluster.KMeans.html>) (some nice examples at the bottom of the documentation)
- See Python handbook for nice examples of kmeans <https://jakevdp.github.io/PythonDataScienceHandbook/05.11-k-means.html> (<https://jakevdp.github.io/PythonDataScienceHandbook/05.11-k-means.html>)

```
In [19]: from sklearn.datasets import make_blobs
X, y_true = make_blobs(n_samples=200, centers=2,
                       cluster_std=0.50, random_state=0)

from sklearn.cluster import KMeans
kmeans = KMeans(n_clusters=2, random_state=0) # 0 and 2 give opposite clust

kmeans.fit(X)
y_kmeans = kmeans.labels_
mu_array = kmeans.cluster_centers_
plot_clustering(X, y_kmeans, mu_array)
```



**This looks great! But isn't this an NP-Hard problem?**

**First caveat: Does not always converge to the optimal/best solution.**

```
In [20]: # Example from Python handbook
from sklearn.datasets import make_blobs
X2, y_true2 = make_blobs(n_samples=300, centers=4,
                        cluster_std=0.60, random_state=0)

kmeans = KMeans(n_clusters=4, init='random', n_init=1, random_state=104) #
kmeans.fit(X2)
plot_clustering(X2, kmeans.labels_, kmeans.cluster_centers_)
```



**Second caveat: Choosing the number of clusters is not obvious**

```
In [21]: # Example from Python handbook
from sklearn.datasets import make_blobs
X2, y_true2 = make_blobs(n_samples=300, centers=4,
                          cluster_std=0.60, random_state=0)

kmeans = KMeans(n_clusters=9, init='random', n_init=1, random_state=0)
kmeans.fit(X2)
plot_clustering(X2, kmeans.labels_, kmeans.cluster_centers_)
```



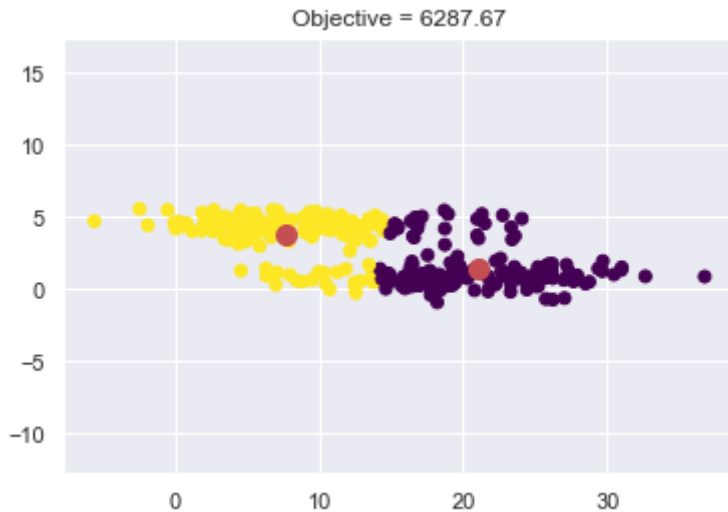
**Third caveat: Scaling of variables and clusters matters**

```
In [22]: from sklearn.datasets import make_moons
X3, y_true = make_blobs(n_samples=300, centers=2,
                        cluster_std=0.60, random_state=0)
X3[:, 0] = X3[:, 0]*10

kmeans = KMeans(n_clusters=2, random_state=0).fit(X3)

plot_clustering(X3, kmeans.labels_, kmeans.cluster_centers_)
plt.axis('equal')
```

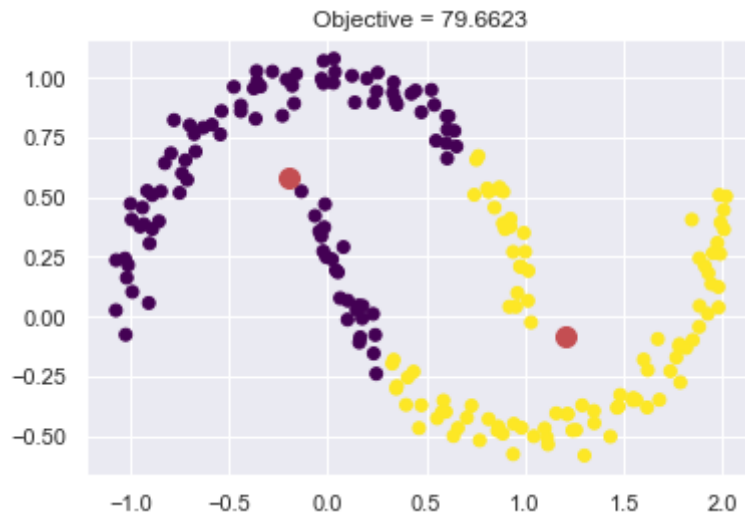
```
Out[22]: (-7.6695009915020105,
          38.84428241167275,
          -1.2536474338023191,
          5.866108305124282)
```



**Fourth caveat: Only linear boundaries between clusters**

```
In [23]: from sklearn.datasets import make_moons
X4, y_true4 = make_moons(200, noise=.05, random_state=0)

kmeans = KMeans(n_clusters=2, random_state=0).fit(X4)
plot_clustering(X4, kmeans.labels_, kmeans.cluster_centers_)
```



In [ ]: