

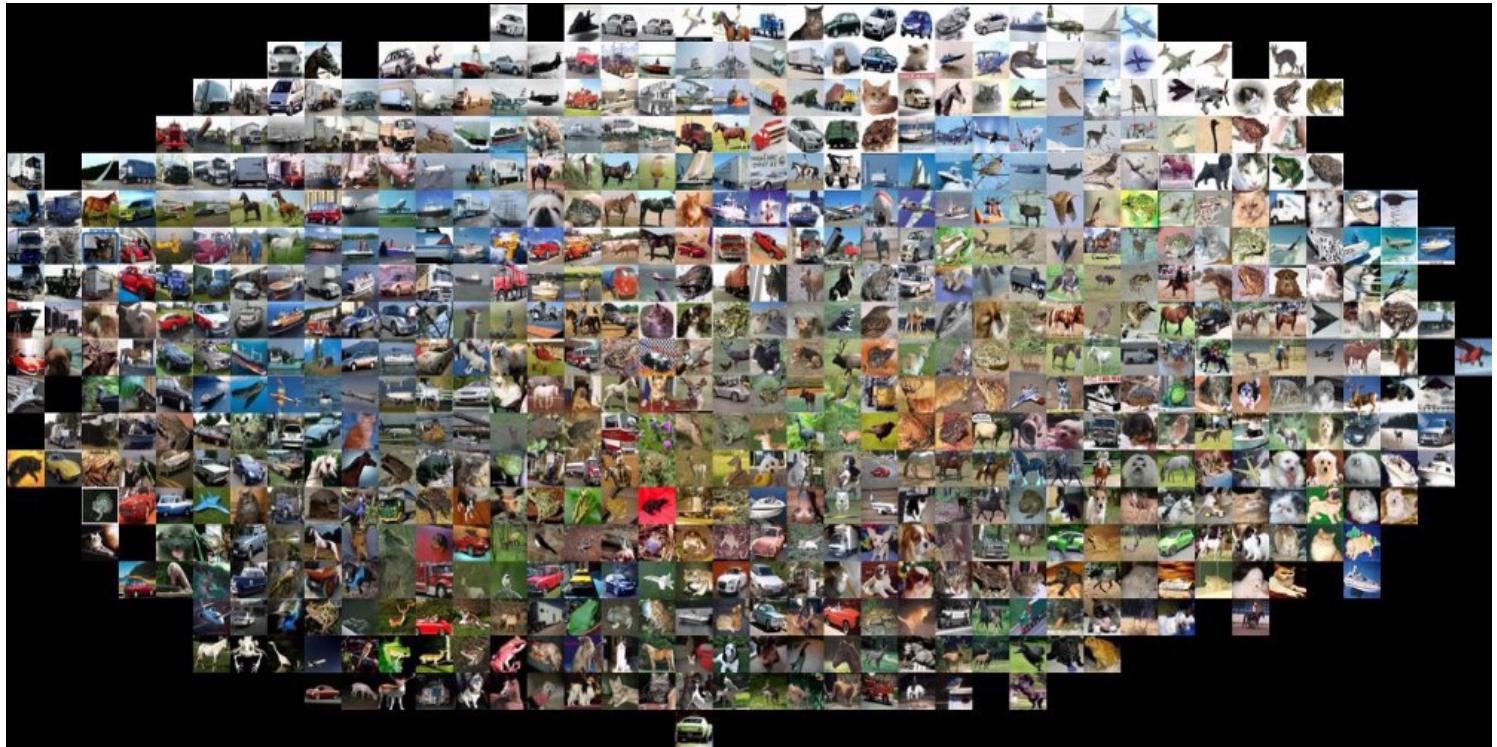


Sergey Korolev

[Follow](#)

Mar 26 · 15 min read

Open Machine Learning Course. Topic 7. Unsupervised Learning: PCA and Clustering



Popular visualization technique t-SNE for CIFAR-10 images using L2 distance. Image source: CS231n course

Welcome to the seventh part of our Open Machine Learning Course!

In this lesson, we will work with unsupervised learning methods such as Principal Component Analysis (PCA) and clustering. You will learn why and how we can reduce the dimensionality of the original data and what the main approaches are for grouping similar data points.

Article outline

1. Introduction

2. Principal Component Analysis

- Intuition, theories, and application issues
- Application examples

3. Cluster analysis

- K-means
- Affinity Propagation
- Spectral clustering
- Agglomerative clustering
- Accuracy metrics

4. Assignment #7

5. Useful links

1. Introduction

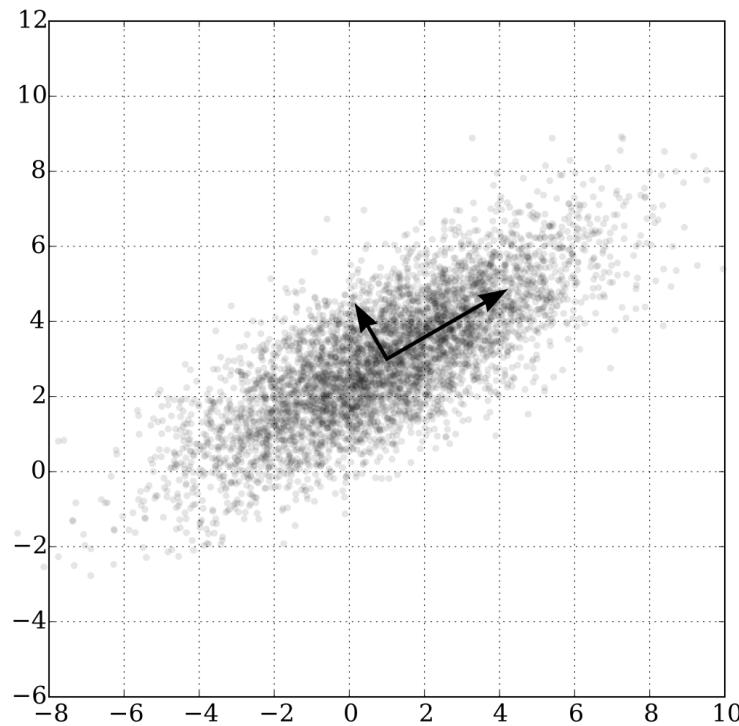
The main feature of unsupervised learning algorithms, when compared to classification and regression methods, is that input data are unlabeled (i.e. no labels or classes given) and that the algorithm learns the structure of the data without any assistance. This creates two main differences. First, it allows us to process large amounts of data because the data does not need to be manually labeled. Second, it is difficult to evaluate the quality of an unsupervised algorithm due to the absence of an explicit goodness metric as used in supervised learning.

One of the most common tasks in unsupervised learning is dimensionality reduction. On one hand, dimensionality reduction may help with data visualization (e.g. t-SNE method) while, on the other hand, it may help deal with the multicollinearity of your data and prepare the data for a supervised learning method (e.g. decision trees).

2. Principal Component Analysis (PCA)

Intuition, theories, and application issues

Principal Component Analysis is one of the easiest, most intuitive, and most frequently used methods for dimensionality reduction, projecting data onto its orthogonal feature subspace.



More generally speaking, all observations can be considered as an ellipsoid in a subspace of an initial feature space, and the new basis set in this subspace is aligned with the ellipsoid axes. This assumption lets us remove highly correlated features since basis set vectors are orthogonal. In the general case, the resulting ellipsoid dimensionality matches the initial space dimensionality, but the assumption that our data lies in a subspace with a smaller dimension allows us to cut off the “excessive” space with the new projection (subspace). We accomplish this in a ‘greedy’ fashion, sequentially selecting each of the ellipsoid axes by identifying where the dispersion is maximal.

“To deal with hyper-planes in a 14 dimensional space, visualize a 3D space and say ‘fourteen’ very loudly. Everyone does it.”—Geoffrey

Hinton

Let's take a look at the mathematical formulation of this process:

In order to decrease the dimensionality of our data from n to k with $k \leq n$, we sort our list of axes in order of decreasing dispersion and take the top- k of them.

We begin by computing the dispersion and the covariance of the initial features. This is usually done with the covariance matrix. According to the covariance definition, the covariance of two features is computed as follows:

$$\text{cov}(X_i, X_j) = E[(X_i - \mu_i)(X_j - \mu_j)] = E[X_i X_j] - \mu_i \mu_j$$

where μ is the expected value of the i -th feature. It is worth noting that the covariance is symmetric, and the covariance of a vector with itself is equal to its dispersion.

Therefore the covariance matrix is symmetric with the dispersion of the corresponding features on the diagonal. Non-diagonal values are the covariances of the corresponding pair of features. In terms of matrices where X is the matrix of observations, the covariance matrix is as follows:

$$\Sigma = E[(\mathbf{X} - E[\mathbf{X}])(\mathbf{X} - E[\mathbf{X}])^T]$$

Quick recap: matrices, as linear operators, have eigenvalues and eigenvectors. They are very convenient because they describe parts of our space that do not rotate and only stretch when we apply linear operators on them; eigenvectors remain in the same direction but are stretched by a corresponding eigenvalue. Formally, a matrix M with eigenvector w and eigenvalue λ satisfy this equation:

$$Mw_i = \lambda_i w_i$$

The covariance matrix for a sample X can be written as a product of a transposed matrix X and X itself. According to the [Rayleigh quotient](#), the maximum variation of our sample lies along the eigenvector of this matrix and is consistent with the maximum eigenvalue. Therefore, the principal components we aim to retain from the data are just the eigenvectors corresponding to the top-k largest eigenvalues of the matrix.

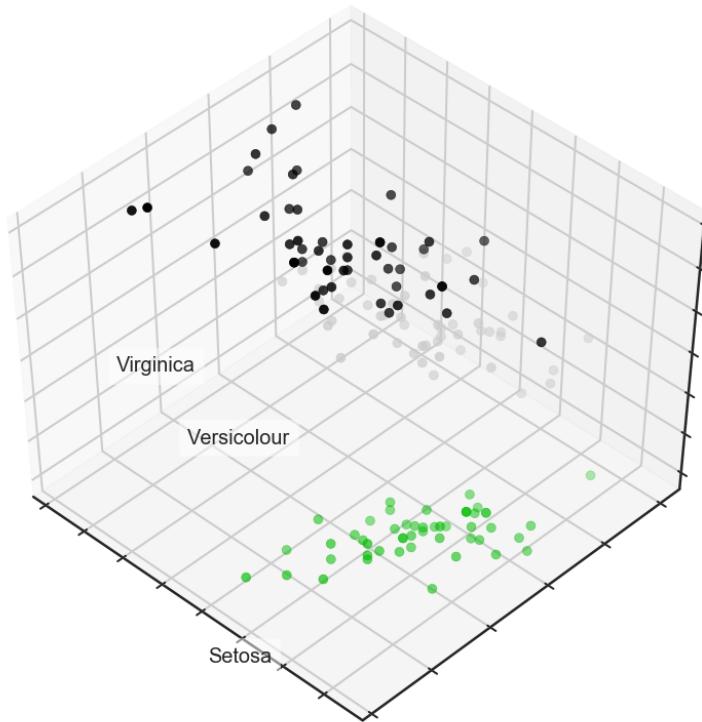
The next steps are easier to digest. We multiply the matrix of our data X by these components to get the projection of our data onto the orthogonal basis of the chosen components. If the number of components was smaller than the initial space dimensionality, remember that we will lose some information upon applying this transformation.

Examples

Fisher's iris dataset

Let's start by uploading all of the essential modules and try out the iris example from the [scikit-learn](#) documentation.

```
1 import numpy as np
2 import matplotlib.pyplot as plt
3 import seaborn as sns; sns.set(style='white')
4 %matplotlib inline
5 %config InlineBackend.figure_format = 'retina'
6 from sklearn import decomposition
7 from sklearn import datasets
8 from mpl_toolkits.mplot3d import Axes3D
9
10 # Loading the dataset
11 iris = datasets.load_iris()
12 X = iris.data
13 y = iris.target
14
15 # Let's create a beautiful 3d-plot
16 fig = plt.figure(1, figsize=(6, 5))
17 plt.clf()
18 ax = Axes3D(fig, rect=[0, 0, .95, 1], elev=48, azim=13
19
20 plt.cla()
21
22 for name, label in [ ('Setosa', 0), ('Versicolour', 1),
```



Now let's see how PCA will improve the results of a simple model that is not able to correctly fit all of the training data:

```
1 from sklearn.tree import DecisionTreeClassifier
2 from sklearn.model_selection import train_test_split
3 from sklearn.metrics import accuracy_score, roc_auc_sc
4
5 # Train, test splits
6 X_train, X_test, y_train, y_test = train_test_split(X,
7 st
8 ra
9
10 # Decision trees with depth = 2
```

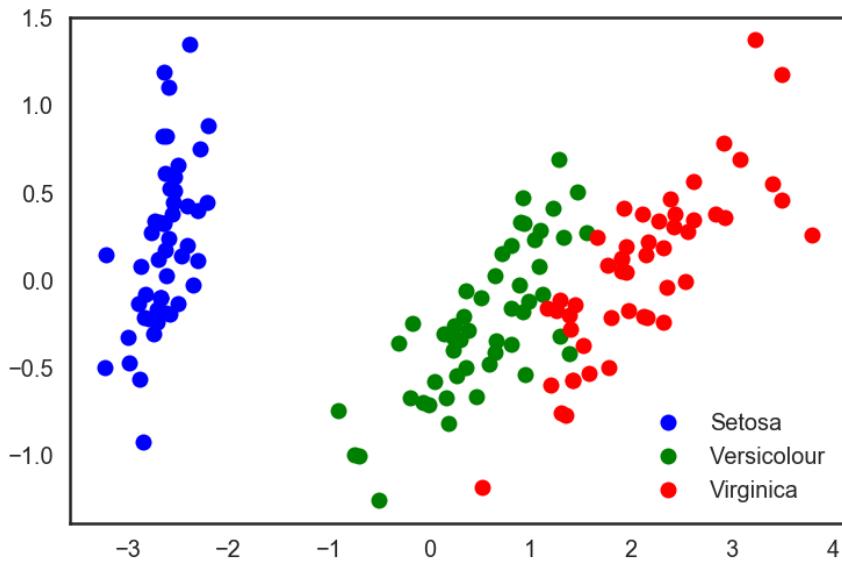
Accuracy: 0.88889

Let's try this again, but, this time, let's reduce the dimensionality to 2 dimensions:

```

1 # Using PCA from sklearn PCA
2 pca = decomposition.PCA(n_components=2)
3 X_centered = X - X.mean(axis=0)
4 pca.fit(X_centered)
5 X_pca = pca.transform(X_centered)
6
7 # Plotting the results of PCA
8 plt.plot(X_pca[y == 0, 0], X_pca[y == 0, 1], 'bo', label='Setosa')
9 plt.plot(X_pca[y == 1, 0], X_pca[y == 1, 1], 'g^', label='Versicolour')
10 plt.plot(X_pca[y == 2, 0], X_pca[y == 2, 1], 'rs', label='Virginica')
11
12 plt.xlabel('PCA Component 1')
13 plt.ylabel('PCA Component 2')
14 plt.legend()
15 plt.show()

```



```

1 # Test-train split and apply PCA
2 X_train, X_test, y_train, y_test = train_test_split(X,
3                                                    test_size=0.2,
4                                                    random_state=42)
5
6 clf = DecisionTreeClassifier(max_depth=2, random_state=42)
7 clf.fit(X_train, y_train)
8
9 accuracy = clf.score(X_test, y_test)
10 print("Accuracy: ", accuracy)

```

Accuracy: 0.91111

The accuracy did not increase significantly in this case, but, with other datasets with a high number of dimensions, PCA can drastically improve the accuracy of decision trees and other ensemble methods.

Now let's check out the percent of variance that can be explained by each of the selected components.

```

1  for i, component in enumerate(pca.components_):
2      print("{} component: {}% of initial variance".format(
3          round(100 * pca.explained_variance_ratio_[i],
4          print(" + ".join("%.3f x %s" % (value, name)
5              for value, name in zip(component,

```

```

1 component: 92.46% of initial variance 0.362 x sepal length
(cm) + -0.082 x sepal width (cm) + 0.857 x petal length (cm)
+ 0.359 x petal width (cm)
2 component: 5.3% of initial variance 0.657 x sepal length
(cm) + 0.730 x sepal width (cm) + -0.176 x petal length (cm)
+ -0.075 x petal width (cm)

```

Handwritten numbers dataset

Let's look at the handwritten numbers dataset that we used before in the [3rd lesson](#).

```

1  digits = datasets.load_digits()
2  X = digits.data
3  y = digits.target

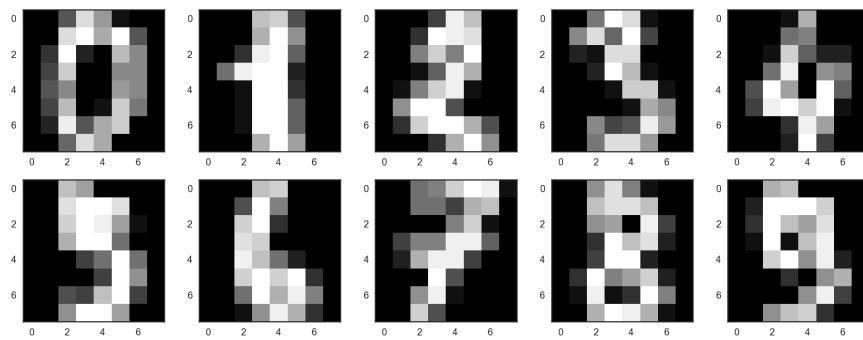
```

Let's start by visualizing our data. Fetch the first 10 numbers. The numbers are represented by 8 x 8 matrixes with the color intensity for each pixel. Every matrix is flattened into a vector of 64 numbers, so we get the feature version of the data.

```

1  # f, axes = plt.subplots(5, 2, sharey=True, figsize=(16, 6))
2  plt.figure(figsize=(16, 6))
3  for i in range(10):
4      plt.subplot(2, 5, i + 1)

```

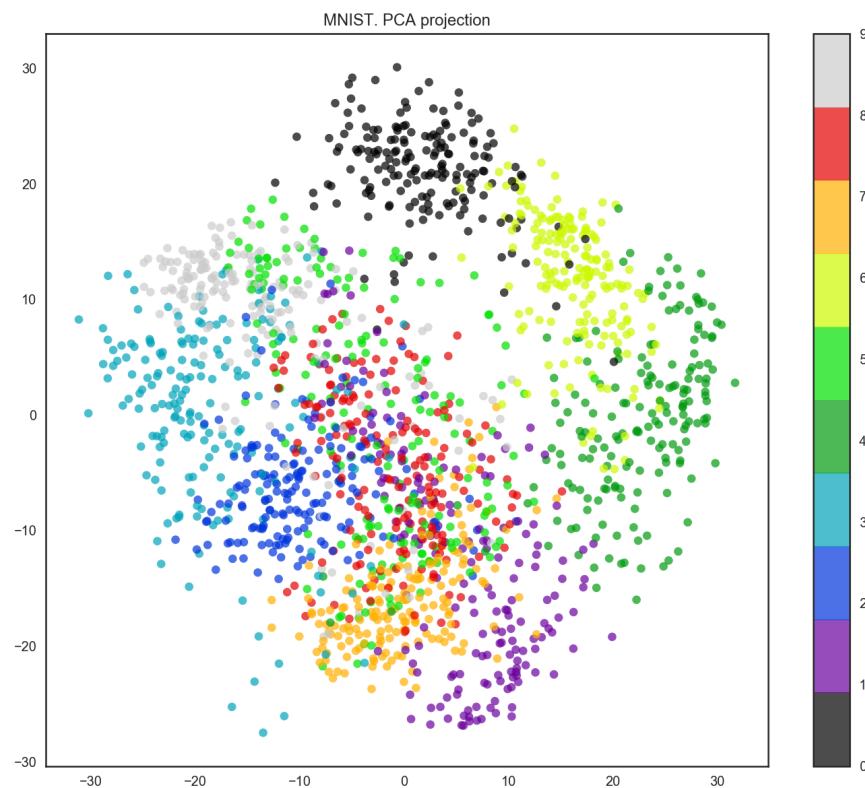


Our data has 64 dimensions, but we are going to reduce it to only 2 and see that, even with just 2 dimensions, we can clearly see that digits separate into clusters.

```

1  pca = decomposition.PCA(n_components=2)
2  X_reduced = pca.fit_transform(X)
3
4  print('Projecting %d-dimensional data to 2D' % X.shape
5
6  plt.figure(figsize=(12,10))
7  plt.scatter(X_reduced[:, 0], X_reduced[:, 1], c=y,
8               edgecolor='none', alpha=0.7, s=40,
```

Projecting 64-dimensional data to 2D

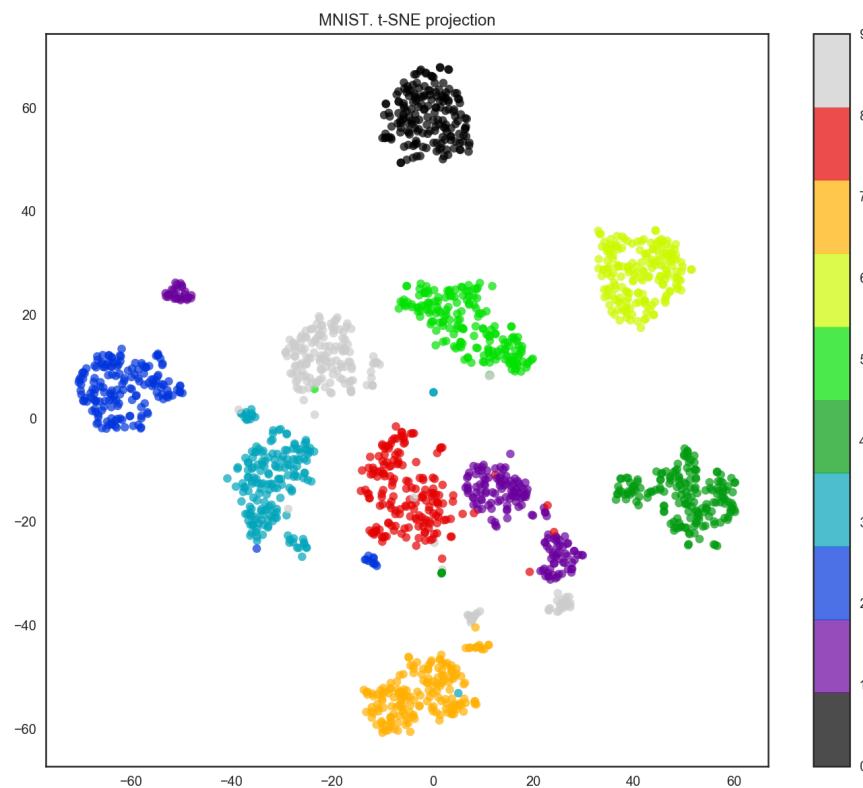


Indeed, with t-SNE, the picture looks better since PCA has a linear constraint while t-SNE does not. However, even with such a small dataset, the t-SNE algorithm takes significantly more time to complete than PCA.

```

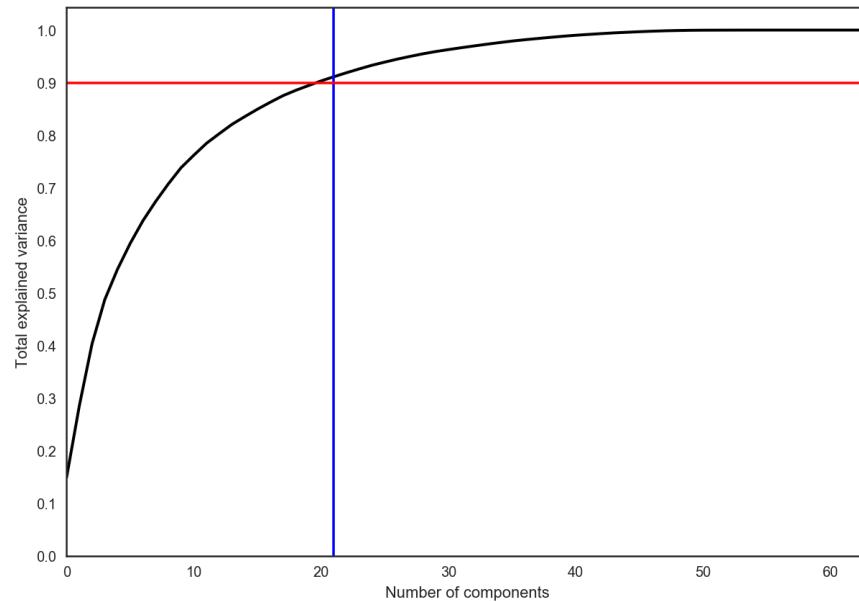
1 %%time
2
3 from sklearn.manifold import TSNE
4 tsne = TSNE(random_state=17)
5
6 X_tsne = tsne.fit_transform(X)
7
8 plt.figure(figsize=(12,10))
9 plt.scatter(X_tsne[:, 0], X_tsne[:, 1], c=y,
```

```
CPU times: user 29.3 s, sys: 2.97 s, total: 32.3 s Wall
time: 32.3 s
```



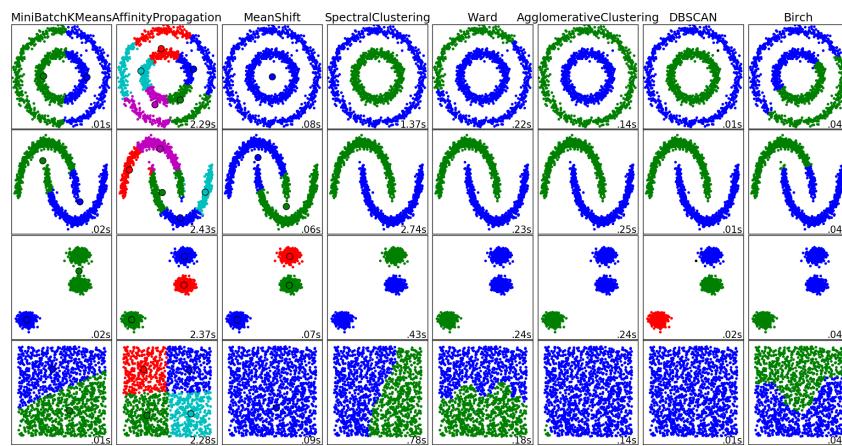
In practice, we would choose the number of principal components such that we can explain 90% of the initial data dispersion (via the `explained_variance_ratio`). Here, that means retaining 21 principal components; therefore, we reduce the dimensionality from 64 features to 21.

```
1 pca = decomposition.PCA().fit(X)
2
3 plt.figure(figsize=(10,7))
4 plt.plot(np.cumsum(pca.explained_variance_ratio_), color='blue')
5 plt.xlabel('Number of components')
6 plt.ylabel('Total explained variance')
7 plt.xlim(0, 63)
8 plt.yticks(np.arange(0, 1.1, 0.1))
```



2. Clustering

The main idea behind clustering is pretty straightforward. Basically, we say to ourselves, “I have these points here, and I can see that they organize into groups. It would be nice to describe these things more concretely, and, when a new point comes in, assign it to the correct group.” This general idea encourages exploration and opens up a variety of algorithms for clustering.



The examples of the outcomes from different algorithms from scikit-learn

The algorithms listed below do not cover all the clustering methods out there, but they are the most commonly used ones.

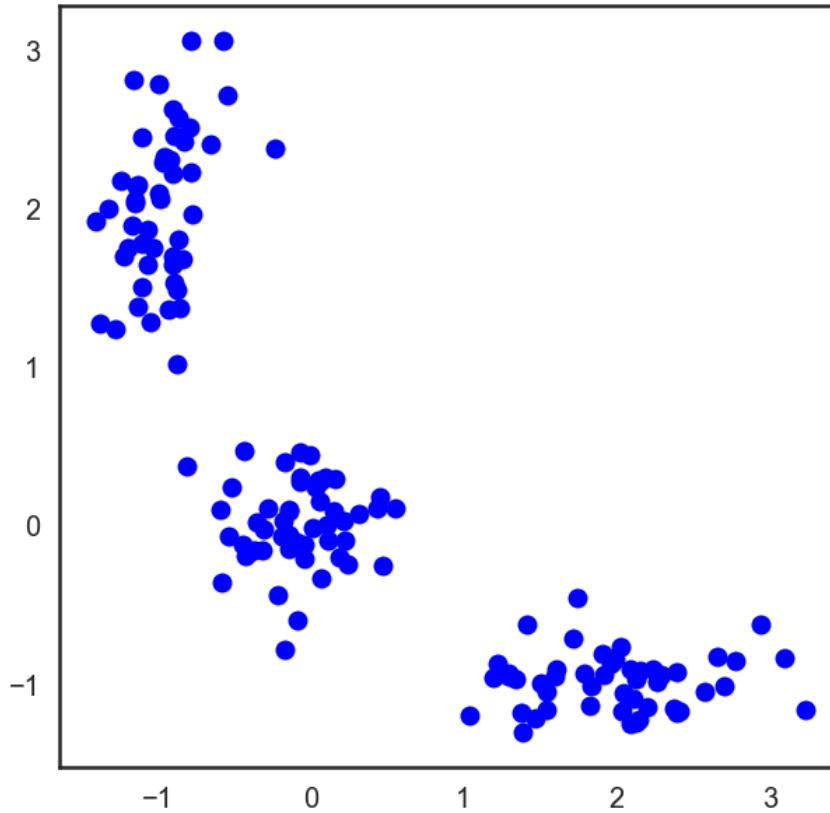
K-means

K-means algorithm is the most popular and yet simplest of all the clustering algorithms. Here is how it works:

1. Select the number of clusters k that you think is the optimal number.
2. Initialize k points as “centroids” randomly within the space of our data.
3. Attribute each observation to its closest centroid.
4. Update the centroids to the center of all the attributed set of observations.
5. Repeat steps 3 and 4 a fixed number of times or until all of the centroids are stable (i.e. no longer change in step 4).

This algorithm is easy to describe and visualize. Let's take a look.

```
1 # Let's begin by allocation 3 cluster's points
2 X = np.zeros((150, 2))
3
4 np.random.seed(seed=42)
5 X[:50, 0] = np.random.normal(loc=0.0, scale=.3, size=50)
6 X[:50, 1] = np.random.normal(loc=0.0, scale=.3, size=50)
7
8 X[50:100, 0] = np.random.normal(loc=2.0, scale=.5, size=50)
9 X[50:100, 1] = np.random.normal(loc=-1.0, scale=.2, size=50)
10
```

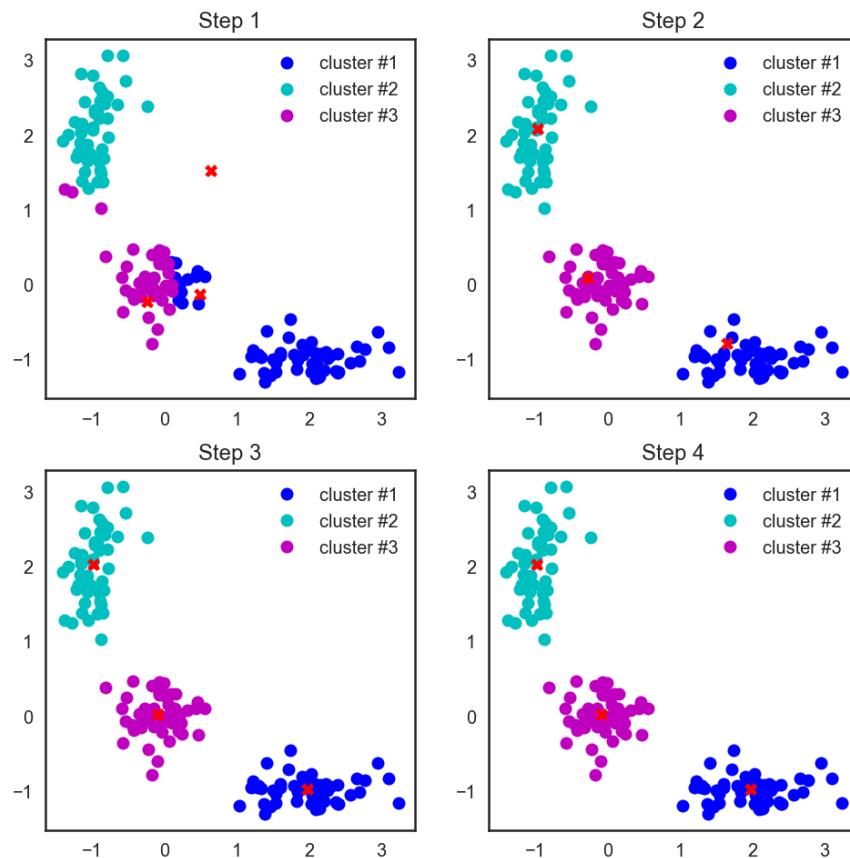


```
1 # Scipy has function that takes 2 tuples and return
2 # calculated distance between them
3 from scipy.spatial.distance import cdist
4
5 # Randomly allocate the 3 centroids
6 np.random.seed(seed=42)
7 centroids = np.random.normal(loc=0.0, scale=1., size=6)
8 centroids = centroids.reshape((3, 2))
9
10 cent_history = []
11 cent_history.append(centroids)
12
13 for i in range(3):
14     # Calculating the distance from a point to a centr
15     distances = cdist(X, centroids)
16     # Checking what's the closest centroid for the poi
```

```

1 # Let's plot K-means
2 plt.figure(figsize=(8, 8))
3 for i in range(4):
4     distances = cdist(X, cent_history[i])
5     labels = distances.argmin(axis=1)
6
7     plt.subplot(2, 2, i + 1)
8     plt.plot(X[labels == 0, 0], X[labels == 0, 1], 'bo')
9     plt.plot(X[labels == 1, 0], X[labels == 1, 1], 'co')

```



Here, we used Euclidean distance, but the algorithm will converge with any other metric. You can not only vary the number of steps or the convergence criteria but also the distance measure between the points and cluster centroids.

Another “feature” of this algorithm is its sensitivity to the initial positions of the cluster centroids. You can run the algorithm several times and then average all the centroid results.

Choosing the number of clusters for K-means

In contrast to the supervised learning tasks such as classification and regression, clustering requires more effort to choose the optimization criterion. Usually, when working with k-means, we optimize the sum of squared distances between the observations and their centroids.

$$J(C) = \sum_{k=1}^K \sum_{i \in C_k} \|x_i - \mu_k\| \rightarrow \min_C$$

where C —is a set of clusters with power K , μ is a centroid of a cluster.

This definition seems reasonable—we want our observations to be as close to their centroids as possible. But, there is a problem—the optimum is reached when the number of centroids is equal to the number of observations, so you would end up with every single observation as its own separate cluster.

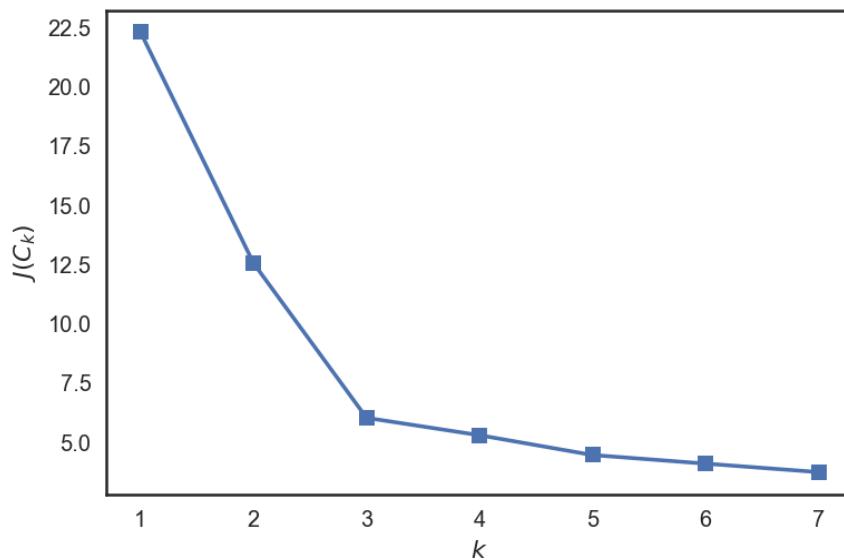
In order to avoid that case, we should choose a number of clusters after which a function $J(C)$ is decreasing less rapidly. More formally,

$$D(k) = \frac{|J(C_k) - J(C_{k+1})|}{|J(C_{k-1}) - J(C_k)|} \rightarrow \min_k$$

Let's look at an example.

```

1  from sklearn.cluster import KMeans
2  inertia = []
3  for k in range(1, 8):
4      kmeans = KMeans(n_clusters=k, random_state=1).fit(x)
5      inertia.append(np.sqrt(kmeans.inertia_))
6  plt.plot(range(1, 8), inertia, marker='s');
```



We see that $J(C)$ decreases significantly until the number of clusters is 3 and then does not change as much anymore. This means that the optimal number of clusters is 3.

Issues

Inherently, K-means is NP-hard. For d dimensions, k clusters, and n observations, we will find a solution in $O(n^{(dk+1)})$ in time. There are some heuristics to deal with this; an example is MiniBatch K-means, which takes portions (batches) of data instead of fitting the whole dataset and then moves centroids by taking the average of the previous steps. Compare the implementation of K-means and MiniBatch K-means in the [scikit-learn documentation](#).

The [implementation](#) of the algorithm using `scikit-learn` has its benefits such as the possibility to state the number of initializations with the `n_init` function parameter, which enables us to identify more robust centroids. Moreover, these runs can be done in parallel to decrease the computation time.

Affinity Propagation

Affinity propagation is another example of a clustering algorithm. As opposed to K-means, this approach does not require us to set the number of clusters beforehand. The main idea here is that we would like to cluster our data based on the similarity of the observations (or how they “correspond” to each other).

Let's define a similarity metric such that $s(x, y) > s(x, z)$ if an observation x is more similar to observation y and less similar to observation z . A simple example of such a similarity metric is a negative square of distance $s(x, y) = -|x-y|^2$.

Now, let's describe "correspondence" by making two zero matrices. One of them, r , determines how well the k -th observation is as a "role model" for the i -th observation with respect to all other possible "role models". Another matrix, a , determines how appropriate it would be for i -th observation to take the k -th observation as a "role model". This may sound confusing, but it becomes more understandable with some hands-on practice.

The matrices are updated sequentially with the following rules:

$$\begin{aligned} r_{i,k} &\leftarrow s(x_i, x_k) - \max_{k' \neq k} \{a_{i,k'} + s(x_i, x'_{k'})\} \\ a_{i,k} &\leftarrow \min \left(0, r_{k,k} + \sum_{i' \notin \{i,k\}} \max(0, r_{i',k}) \right), \quad i \neq k \\ a_{k,k} &\leftarrow \sum_{i' \neq k} \max(0, r_{i',k}) \end{aligned}$$

Spectral clustering

Spectral clustering combines some of the approaches described above to create a stronger clustering method.

First of all, this algorithm requires us to define the similarity matrix for observations called the adjacency matrix. This can be done in a similar fashion as in the Affinity Propagation algorithm, so that matrix A hosts negative square of the distances between the corresponding points. This matrix describes a full graph with the observations as vertices and the estimated similarity value between a pair of observations as edge weights for that pair of vertices. For the metric defined above and two-dimensional observations, this is pretty intuitive—two observations are similar if the edge between them is shorter. We'd like to split up the graph into two subgraphs in such a way that each observation in each subgraph would be similar to another observation in that subgraph.

Formally, this is a Normalized cuts problem; for more details, we recommend reading [this paper](#).

Agglomerative clustering

The following algorithm is the simplest and easiest to understand among all the the clustering algorithms without a fixed number of clusters.

The algorithm is fairly simple:

1. We start by assigning each observation to its own cluster
2. Then sort the pairwise distances between the centers of clusters in descending order
3. Take the nearest two neighbor clusters and merge them together, and recompute the centers
4. Repeat steps 2 and 3 until all the data is merged into one cluster

The process of searching for the nearest cluster can be conducted with different methods of bounding the observations:

1. Single linkage—distance between centroids is the minimum distance between all the pairs of points where one is from the first cluster, and the second is from the other.
2. Complete linkage—distance between centroids is the maximum distance between all the pairs of points where one is from the first cluster, and the second is from the other.
3. Average linkage—distance between centroids is the adjusted sum of all of the distances between all the pairs of points where one is from the first cluster, and the second is from the other.
4. Centroid linkage—distance between centroids is the distance between their geometric centers of clusters.

The 3rd one is the most effective in computation time since it does not require recomputing the distances every time the clusters are merged.

The results can be visualized as a beautiful cluster tree (dendrogram) to help recognize the moment the algorithm should be stopped to get

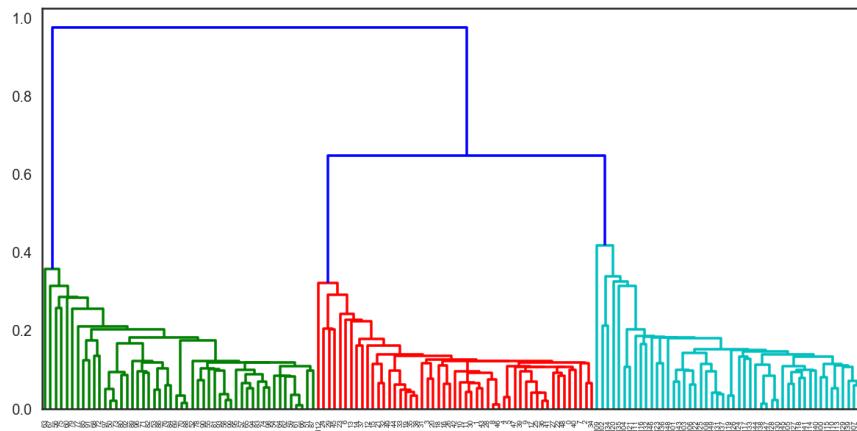
optimal results. There are plenty of Python tools to build these dendograms for agglomerative clustering.

Let's consider an example with the clusters we got from K-means:

```

1  from scipy.cluster import hierarchy
2  from scipy.spatial.distance import pdist
3
4  X = np.zeros((150, 2))
5
6  np.random.seed(seed=42)
7  X[:50, 0] = np.random.normal(loc=0.0, scale=.3, size=50)
8  X[:50, 1] = np.random.normal(loc=0.0, scale=.3, size=50)
9
10 X[50:100, 0] = np.random.normal(loc=2.0, scale=.5, size=50)
11 X[50:100, 1] = np.random.normal(loc=-1.0, scale=.2, size=50)
12
13 X[100:150, 0] = np.random.normal(loc=-1.0, scale=.2, size=50)
14 X[100:150, 1] = np.random.normal(loc=0.5, scale=.3, size=50)

```



Accuracy metrics

As opposed to classification, it is difficult to assess the quality of results from clustering. Here, a metric cannot depend on the labels but only on the goodness of split. Secondly, we do not usually have true labels of the observations when we use clustering.

There are *internal* and *external* goodness metrics. External metrics use the information about the known true split while internal metrics do not use any external information and assess the goodness of clusters based only on the initial data. The optimal number of clusters is usually defined with respect to some internal metrics.

All the metrics described below are implemented in `sklearn.metrics`.

Adjusted Rand Index (ARI)

Here, we assume that the true labels of objects are known. This metric does not depend on the labels' values but on the data cluster split. Let N be the number of observations in a sample. Let a to be the number of observation pairs with the same labels and located in the same cluster, and let b to be the number of observations with different labels and located in different clusters. The Rand Index can be calculated using the following formula: $RI = 2(a + b)/n(n - 1)$. In other words, it evaluates a share of observations for which these splits (initial and clustering result) are consistent. The Rand Index (RI) evaluates the similarity of the two splits of the same sample. In order for this index to be close to zero for any clustering outcomes with any n and number of clusters, it is essential to scale it, hence the Adjusted Rand Index: $ARI = (RI - E(RI)) / (\max(RI) - E(RI))$.

This metric is symmetric and does not depend in the label permutation. Therefore, this index is a measure of distances between different sample splits. ARI takes on values in the $[-1, 1]$ range. Negative values indicate the independence of splits, and positive values indicate that these splits are consistent (they match $ARI = 1$).

Adjusted Mutual Information (AMI)

This metric is similar to ARI. It is also symmetric and does not depend on the labels' values and permutation. It is defined by the [entropy] ([https://en.wikipedia.org/wiki/Entropy_\(information_theory\)](https://en.wikipedia.org/wiki/Entropy_(information_theory))) function and interprets a sample split as a discrete distribution (likelihood of assigning to a cluster is equal to the percent of objects in it). The MI index is defined as the mutual information for two distributions, corresponding to the sample split into clusters. Intuitively, the mutual information measures the share of information

common for both clustering splits i.e. how information about one of them decreases the uncertainty of the other one.

Similarly to the ARI, the AMI is defined. This allows us to get rid of the MI index's increase with the number of clusters. The AMI lies in the [0, 1] range. Values close to zero mean the splits are independent, and those close to 1 mean they are similar (with complete match at AMI = 1).

Homogeneity, completeness, V-measure

Formally, these metrics are also defined based on the entropy function and the conditional entropy function, interpreting the sample splits as discrete distributions:

$$h = 1 - \frac{H(C | K)}{H(C)}, c = 1 - \frac{H(K | C)}{H(K)},$$

where K is a clustering result and C is the initial split. Therefore, h evaluates whether each cluster is composed of same class objects, and c measures how well the same class objects fit the clusters. These metrics are not symmetric. Both lie in the [0, 1] range, and values closer to 1 indicate more accurate clustering results. These metrics' values are not scaled as the ARI or AMI metrics are and thus depend on the number of clusters. A random clustering result will not have metrics' values closer to zero when the number of clusters is big enough and the number of objects is small. In such a case, it would be more reasonable to use ARI. However, with a large number of observations (more than 100) and the number of clusters less than 10, this issue is less critical and can be ignored.

V-measure is a combination of h, and c and is their harmonic mean: $v = (2hc)/(h + c)$. It is symmetric and measures how consistent two clustering results are.

Silhouette

In contrast to the metrics described above, this coefficient does not imply the knowledge about the true labels of the objects. It lets us estimate the quality of the clustering using only the initial, unlabeled sample and the clustering result. To start with, for each observation,

the silhouette coefficient is computed. Let a be the mean of the distance between an object and other objects within one cluster and b be the mean distance from an object to an object from the nearest cluster (different from the one the object belongs to). Then the silhouette measure for this object is $s = (b - a) / \max(a, b)$.

The silhouette of a sample is a mean value of silhouette values from this sample. Therefore, the silhouette distance shows to which extent the distance between the objects of the same class differ from the mean distance between the objects from different clusters. This coefficient takes values in the $[-1, 1]$ range. Values close to -1 correspond to bad clustering results while values closer to 1 correspond to dense, well-defined clusters. Therefore, the higher the silhouette value is, the better the results from clustering.

With the help of silhouette, we can identify the optimal number of clusters k (if we don't know it already from the data) by taking the number of clusters that maximizes the silhouette coefficient.

To conclude, let's take a look at how these metrics perform with the MNIST handwritten numbers dataset:

```

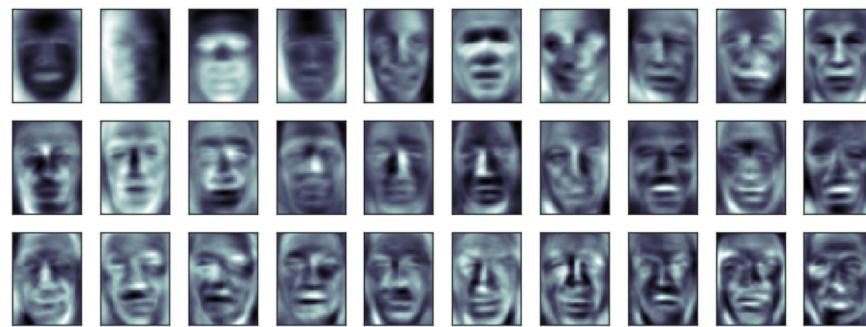
1  from sklearn import metrics
2  from sklearn import datasets
3  import pandas as pd
4  from sklearn.cluster import KMeans, AgglomerativeClust
5
6
7  data = datasets.load_digits()
8  X, y = data.data, data.target
9
10 algorithms = []
11 algorithms.append(KMeans(n_clusters=10, random_state=1))
12 algorithms.append(AffinityPropagation())
13 algorithms.append(SpectralClustering(n_clusters=10, ra
14                                     affinity='nearest')
15 algorithms.append(AgglomerativeClustering(n_clusters=1)
16
17 data = []
18 for algo in algorithms:
19     algo.fit(X)
20     data.append({
21         'ARI': metrics.adjusted_rand_score(y, algo.labels_),
22         'AMI': metrics.adjusted_mutual_info_score(y, a

```

	ARI	AMI	Homogeneity	Completeness	V-measure	Silhouette
K-means	0.662295	0.732799	0.735448	0.742972	0.739191	0.182097
Affinity	0.175174	0.451249	0.958907	0.486901	0.645857	0.115197
Spectral	0.752639	0.827818	0.829544	0.876367	0.852313	0.182195
Agglomerative	0.794003	0.856085	0.857513	0.879096	0.868170	0.178497

4. Assignment #7

In this assignment, you'll apply `scikit-learn` dimensionality reduction and clustering methods to the "faces" dataset to get such crazy pictures like this:



Instructions are given in [this](#) Jupyter notebook.

Deadline: April 4, 23.59 UTC+3.

5. Useful links

- Overview of clustering methods in the [scikit-learn doc](#).
 - [Q&A](#) for PCA with examples
 - [Notebook](#) on k-means and the EM-algorithm

• • •

Author: Sergey Korolev, Software Engineer at Snap Inc. Translated and edited by Anna Golovchenko, Sergey Korolev, Yury Kashnitsky, and Yuanyuan Pao.