



Lecture 18: Clustering

Applied Machine Learning

Volodymyr Kuleshov
Cornell Tech

Part 1: Gaussian Mixture Models

Clustering is a common unsupervised learning problem with numerous applications.

We will start by defining the problem and outlining some models for this problem.

Review: Unsupervised Learning

We will assume that the dataset is sampled from a probability distribution P_{data} , which we will call the *data distribution*. We will denote this as

$$x \sim P_{\text{data}}.$$

The dataset $\mathcal{D} = \{x^{(i)} \mid i = 1, 2, \dots, n\}$ consists of *independent and identically distributed* (IID) samples from P_{data} .

Clustering

Clustering is the problem of identifying distinct components in the data distribution.

- A cluster $C_k \subseteq \mathcal{X}$ is associated with a subset of the x coming from P_{data} .
- Datapoints in a cluster are more similar to each other than to other clusters
- Clusters are usually defined by their centers, and potentially by other shape parameters.

Review: *K*-Means

K-Means is the simplest example of a clustering algorithm.

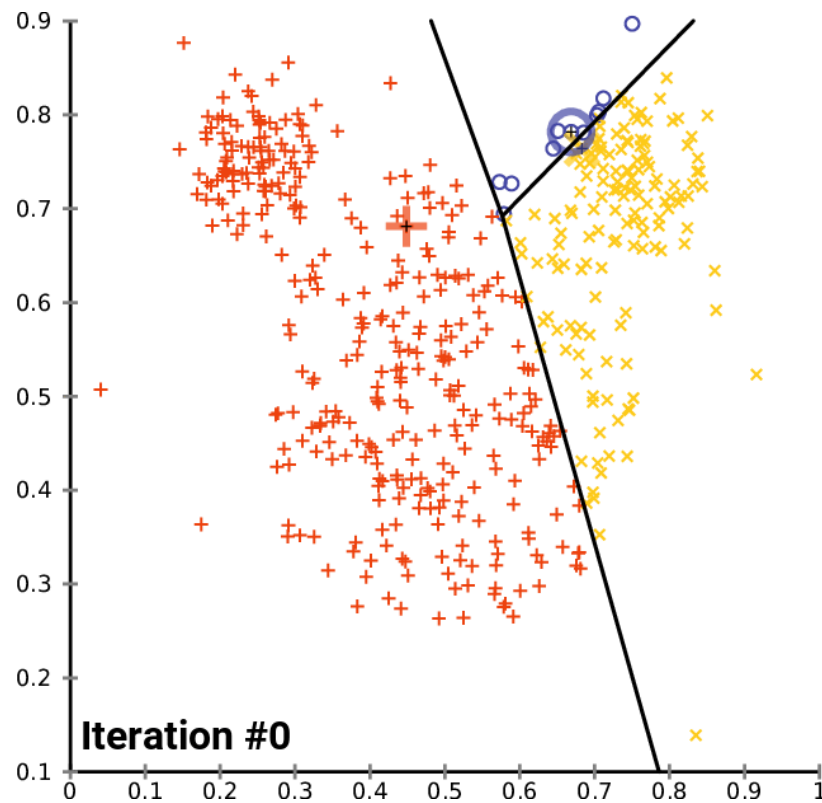
- The algorithm seeks to find *K* hidden clusters in the data.
- Each cluster is characterized by its centroid (its mean).
- The clusters reveal interesting structure in the data.

We seek centroids c_k such that the distance between the points and their closest centroid is minimized:

$$J(\theta) = \sum_{i=1}^n ||x^{(i)} - \text{centroid}(f_{\theta}(x^{(i)}))||,$$

where $\text{centroid}(k) = c_k$ denotes the centroid for cluster k .

This is best illustrated visually (from [Wikipedia](https://commons.wikimedia.org/wiki/File:K-means_convergence.gif)
(https://commons.wikimedia.org/wiki/File:K-means_convergence.gif):



K-Means has a number of limitations:

- Clustering can get stuck in local minima
- Measuring clustering quality is hard and relies on heuristics
- Cluster assignment is binary and doesn't estimate confidence

Gaussian Mixture Models

Gaussian mixtures are latent-variable probabilistic models that are useful for clustering. They define a model

$$P_{\theta}(x, z) = P_{\theta}(x|z)P_{\theta}(z)$$

- $z \in \{1, 2, \dots, K\}$ is discrete and follows a categorical distribution
 $P_{\theta}(z = k) = \phi_k$.
- $x \in \mathbb{R}$ is continuous; conditioned on $z = k$, it follows a Normal distribution
 $P_{\theta}(x|z = k) = \mathcal{N}(\mu_k, \Sigma_k)$.

The parameters θ are the μ_k, Σ_k, ϕ_k for all $k = 1, 2, \dots, K$.

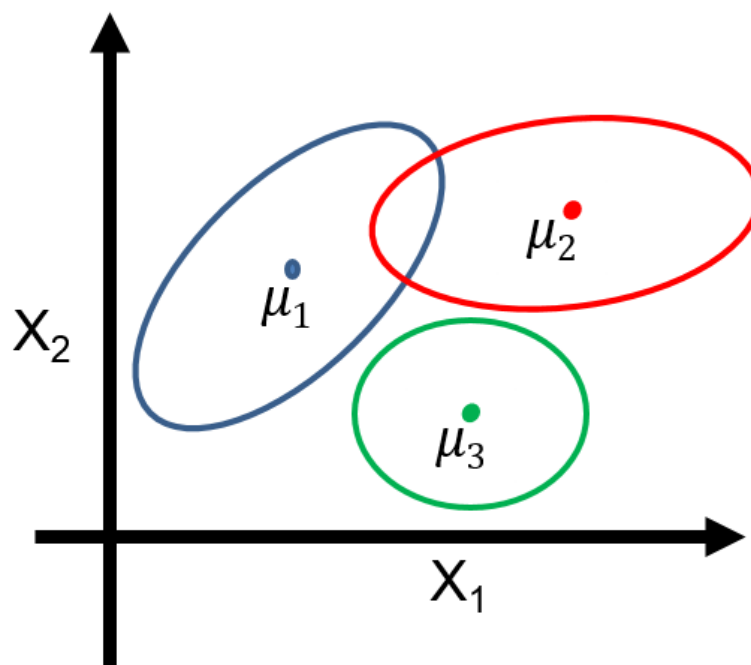
Gaussian Mixtures for Clustering

Gaussian mixtures define a model

$$P_{\theta}(x, z) = P_{\theta}(x|z)P_{\theta}(z)$$

- This model postulates that our observed data is comprised of K clusters with proportions specified by $\phi_1, \phi_2, \dots, \phi_K$
- The points within each cluster follow a Normal distribution
- To generate a new data point, we sample a cluster $z = k$ from $P_{\theta}(z)$ and then x sample from its Gaussian $P_{\theta}(x|z = k)$

This is best understood via a picture.



Mixtures of Gaussians fit more complex distributions than one Gaussian.

Raw data

Single Gaussian

Mixture of Gaussians

Recovering Clusters from GMMs

Given a trained model $P_\theta(x, z) = P_\theta(x|z)P_\theta(z)$, we can look at the *posterior* probability

$$P_\theta(z = k | x) = \frac{P_\theta(z = k, x)}{P_\theta(x)} = \frac{P_\theta(x|z = k)P_\theta(z = k)}{\sum_{l=1}^K P_\theta(x|z = l)P_\theta(z = l)}$$

of a point x belonging to class k .

- The posterior distribution defines a "soft" assignment of x to each class.
- This is in contrast to the hard assignments from K -Means.

Learning GMMs

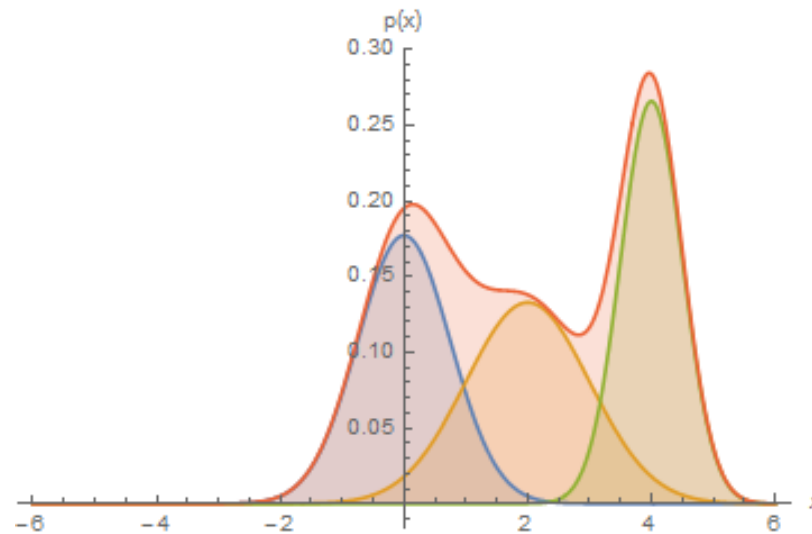
Gaussian mixtures are latent variable models, and we can learn them using maximum marginal log-likelihood:

$$\max_{\theta} \sum_{x \in \mathcal{D}} \log P_{\theta}(x) = \max_{\theta} \sum_{x \in \mathcal{D}} \log \left(\sum_{z \in \mathcal{Z}} P_{\theta}(x, z) \right)$$

- Unlike in GMMs for supervised learning, cluster assignments are latent.
- Hence, there is no closed form solution for θ .
- We will soon see specialized algorithms for this task.

Optimizing the likelihood of latent variable models is hard.

A Gaussian has a single maximum, but a mixture has many and its objective is non-convex (hard to optimize).

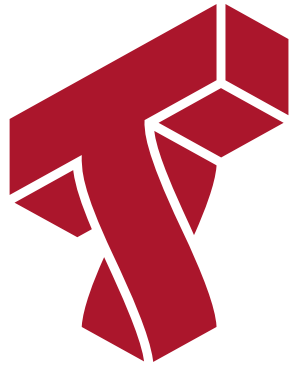


Beyond Gaussian Mixtures

We will focus on Gaussian mixture models in this lecture, but there exist many other kinds of clustering:

- Hierarchical clusters
- Points belonging to multiple clusters (e.g. topics)
- Clusters in graphs

See the `scikit-learn` [guide \(https://scikit-learn.org/stable/modules/clustering.html\)](https://scikit-learn.org/stable/modules/clustering.html) for more!



Part 2: Expectation Maximization

We will now describe expectation maximization (EM), an algorithm that can be used to fit Gaussian mixture models.

Review: Unsupervised Learning

We will assume that the dataset is sampled from a probability distribution P_{data} , which we will call the *data distribution*. We will denote this as

$$x \sim P_{\text{data}}.$$

The dataset $\mathcal{D} = \{x^{(i)} \mid i = 1, 2, \dots, n\}$ consists of *independent and identically distributed* (IID) samples from P_{data} .

Review: Gaussian Mixture Models

Gaussian mixtures are latent-variable probabilistic models that are useful for clustering. They define a model

$$P_{\theta}(x, z) = P_{\theta}(x|z)P_{\theta}(z)$$

- $z \in \{1, 2, \dots, K\}$ is discrete and follows a categorical distribution
 $P_{\theta}(z = k) = \phi_k$.
- $x \in \mathbb{R}$ is continuous; conditioned on $z = k$, it follows a Normal distribution
 $P_{\theta}(x|z = k) = \mathcal{N}(\mu_k, \Sigma_k)$.

The parameters θ are the μ_k, Σ_k, ϕ_k for all $k = 1, 2, \dots, K$.

Review: Learning GMMs

Gaussian mixtures are latent variable models, and we can learn them using maximum marginal log-likelihood:

$$\max_{\theta} \sum_{x \in \mathcal{D}} \log P_{\theta}(x) = \max_{\theta} \sum_{x \in \mathcal{D}} \log \left(\sum_{z \in \mathcal{Z}} P_{\theta}(x, z) \right)$$

- Unlike in GMMs for supervised learning, cluster assignments are latent.
- Hence, there is not a closed form solution for θ .
- We will see specialized algorithm for this task.

Expectation Maximization: Intuition

Expectation maximization (EM) is an algorithm for maximizing marginal log-likelihood

$$\max_{\theta} \sum_{x^{(i)} \in \mathcal{D}} \log \left(\sum_{z \in \mathcal{Z}} P_{\theta}(x^{(i)}, z) \right)$$

that can also be used to learn Gaussian mixtures.

We want to optimize the marginal log-likelihood

$$\max_{\theta} \sum_{x^{(i)} \in \mathcal{D}} \log \left(\sum_{z \in \mathcal{Z}} P_{\theta}(x^{(i)}, z) \right).$$

- If we know the true $z^{(i)}$ for each $x^{(i)}$, we maximize

$$\max_{\theta} \sum_{x^{(i)}, z^{(i)} \in \mathcal{D}} \log \left(P_{\theta}(x^{(i)}, z^{(i)}) \right).$$

and it's easy to find the best θ (use solution for supervised learning).

- If we know θ , we can estimate the cluster assignments $z^{(i)}$ for each i by computing $P_{\theta}(z|x^{(i)})$.

Expectation maximization alternates between these two steps.

1. **(E-Step)** Given an estimate θ_t of the weights, compute $P_{\theta}(z|x^{(i)})$. and use it to “hallucinate” expected cluster assignments $z^{(i)}$.
2. **(M-Step)** Find a new θ_{t+1} that maximizes the marginal log-likelihood by optimizing $P_{\theta}(x^{(i)}, z^{(i)})$ given the $z^{(i)}$ from step 1.

This process increases the marginal likelihood at each step and eventually converges.

Expectation Maximization: Definition

Formally, EM learns the parameters θ of a latent-variable model $P_\theta(x, z)$ over a dataset $\mathcal{D} = \{x^{(i)} \mid i = 1, 2, \dots, n\}$ as follows.

For $t = 0, 1, 2, \dots$, repeat until convergence:

1. **(E-Step)** For each $x^{(i)} \in \mathcal{D}$ compute $P_{\theta_t}(z|x^{(i)})$
2. **(M-Step)** Compute new weights θ_{t+1} as

$$\theta_{t+1} = \arg \max_{\theta} \sum_{i=1}^n \mathbb{E}_{z^{(i)} \sim P_{\theta_t}(z|x^{(i)})} \log P_{\theta}(x^{(i)}, z^{(i)})$$

Since assignments $P_{\theta_t}(z|x^{(i)})$ are "soft", M-step involves an expectation.

Understanding the E-Step

Intuitively, we hallucinate $z^{(i)}$ in the E-Step.

In practice, the $P_{\theta_t}(z|x^{(i)})$ define "soft" assignments, and we compute a vector of class probabilities for each $x^{(i)}$. * The $P_{\theta_t}(z|x^{(i)})$ define "soft" assignments, and we compute a vector of class probabilities for each $x^{(i)}$.

- We compute an expected values over $z^{(i)}$ instead of hallucinating one value. -->

Understanding the M-Step

Since class assignments from E-step are probabilistic, we maximize an expectation:

$$\begin{aligned}\theta_{t+1} &= \arg \max_{\theta} \sum_{i=1}^n \mathbb{E}_{z^{(i)} \sim P_{\theta_t}(z|x^{(i)})} \log P_{\theta}(x^{(i)}, z^{(i)}) \\ &= \arg \max_{\theta} \sum_{i=1}^n \sum_{k=1}^K P_{\theta_t}(z = k|x^{(i)}) \log P_{\theta}(x^{(i)}, z = k)\end{aligned}$$

For many interesting models, this is tractable.

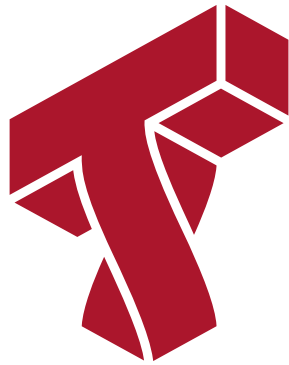
Pros and Cons of EM

EM is a very important optimization algorithm in machine learning.

- It is easy to implement and is guaranteed to converge.
- It works in a lot of important ML models.

Its limitations include:

- It can get stuck in local optima.
- We may not be able to compute $P_{\theta_t}(z|x^{(i)})$ in every model.



Part 3: Expectation Maximization in Gaussian Mixture Models

Next, let's work through how Expectation Maximization works in Gaussian Mixture Models.

Review: Gaussian Mixture Models

Gaussian mixtures are latent-variable probabilistic models that are useful for clustering. They define a model

$$P_{\theta}(x, z) = P_{\theta}(x|z)P_{\theta}(z)$$

- $z \in \{1, 2, \dots, K\}$ is discrete and follows a categorical distribution
 $P_{\theta}(z = k) = \phi_k$.
- $x \in \mathbb{R}$ is continuous; conditioned on $z = k$, it follows a Normal distribution
 $P_{\theta}(x|z = k) = \mathcal{N}(\mu_k, \Sigma_k)$.

The parameters θ are the μ_k, Σ_k, ϕ_k for all $k = 1, 2, \dots, K$.

Review: Expectation Maximization

Formally, EM learns the parameters θ of a latent-variable model $P_\theta(x, z)$ over a dataset $\mathcal{D} = \{x^{(i)} \mid i = 1, 2, \dots, n\}$ as follows.

For $t = 0, 1, 2, \dots$, repeat until convergence:

1. **(E-Step)** For each $x^{(i)} \in \mathcal{D}$ compute $P_{\theta_t}(z|x^{(i)})$
2. **(M-Step)** Compute new weights θ_{t+1} as

$$\theta_{t+1} = \arg \max_{\theta} \sum_{i=1}^n \mathbb{E}_{z^{(i)} \sim P_{\theta_t}(z|x^{(i)})} \log P_{\theta}(x^{(i)}, z^{(i)})$$

Since assignments $P_{\theta_t}(z|x^{(i)})$ are "soft", M-step involves an expectation.

Deriving the E-Step

In the E-step, we compute the posterior for each data point x as follows

$$P_{\theta}(z = k \mid x) = \frac{P_{\theta}(z = k, x)}{P_{\theta}(x)} = \frac{P_{\theta}(x|z = k)P_{\theta}(z = k)}{\sum_{l=1}^K P_{\theta}(x|z = l)P_{\theta}(z = l)}$$

$P_{\theta}(z \mid x)$ defines a vector of probabilities that x originates from component k given the current set of parameters θ

Deriving the M-Step

At the M-step, we optimize the expected log-likelihood of our model.

$$\begin{aligned} \max_{\theta} \sum_{x \in D} \mathbb{E}_{z \sim P_{\theta_t}(z|x)} \log P_{\theta}(x, z) = \\ \max_{\theta} \left(\sum_{k=1}^K \sum_{x \in D} P_{\theta_t}(z_k|x) \log P_{\theta}(x|z_k) + \sum_{k=1}^K \sum_{x \in D} P_{\theta_t}(z_k|x) \log P_{\theta}(z_k) \right) \end{aligned}$$

As in supervised learning, we can optimize the two terms above separately.

We will start with $P_\theta(x \mid z = k) = \mathcal{N}(x; \mu_k, \Sigma_k)$. We have to find μ_k, Σ_k that optimize

$$\max_{\theta} \sum_{x^{(i)} \in D} P(z = k | x^{(i)}) \log P_\theta(x^{(i)} | z = k)$$

Note that this corresponds to fitting a Gaussian to a dataset whose elements $x^{(i)}$ each have a weight $P(z = k | x^{(i)})$.

Similarly to how we did this the supervised regime, we compute the derivative, set it to zero, and obtain closed form solutions:

$$\begin{aligned}\mu_k &= \frac{\sum_{i=1}^n P(z = k | x^{(i)}) x^{(i)}}{n_k} \\ \Sigma_k &= \frac{\sum_{i=1}^n P(z = k | x^{(i)}) (x^{(i)} - \mu_k)(x^{(i)} - \mu_k)^\top}{n_k} \\ n_k &= \sum_{i=1}^n P(z = k | x^{(i)})\end{aligned}$$

Intuitively, the optimal mean and covariance are the empirical mean and covariance of the dataset \mathcal{D} when each element $x^{(i)}$ has a weight $P(z = k | x^{(i)})$.

Similarly, we can show that the class priors are

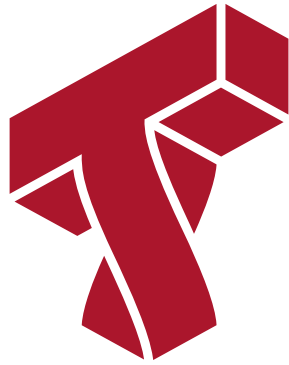
$$\phi_k = \frac{n_k}{n}$$
$$n_k = \sum_{i=1}^n P(z = k | x^{(i)})$$

EM in Gaussian Mixture Models

EM learns the parameters θ of a Gaussian mixture model $P_\theta(x, z)$ over a dataset $\mathcal{D} = \{x^{(i)} \mid i = 1, 2, \dots, n\}$ as follows.

For $t = 0, 1, 2, \dots$, repeat until convergence:

1. **(E-Step)** For each $x^{(i)} \in \mathcal{D}$ compute $P_{\theta_t}(z|x^{(i)})$
2. **(M-Step)** Compute parameters μ_k, Σ_k, ϕ_k using the above formulas



Part 4: Generalization in Probabilistic Models

Let's now revisit the concepts of overfitting and underfitting in GMMs.

Review: Data Distribution

We will assume that the dataset is sampled from a probability distribution \mathbb{P} , which we will call the *data distribution*. We will denote this as

$$x \sim \mathbb{P}.$$

The dataset $\mathcal{D} = \{x^{(i)} \mid i = 1, 2, \dots, n\}$ consists of *independent and identically distributed* (IID) samples from \mathbb{P} .

Review: Gaussian Mixture Models

Gaussian mixtures are latent-variable probabilistic models that are useful for clustering. They define a model

$$P_{\theta}(x, z) = P_{\theta}(x|z)P_{\theta}(z)$$

- $z \in \{1, 2, \dots, K\}$ is discrete and follows a categorical distribution
 $P_{\theta}(z = k) = \phi_k$.
- $x \in \mathbb{R}$ is continuous; conditioned on $z = k$, it follows a Normal distribution
 $P_{\theta}(x|z = k) = \mathcal{N}(\mu_k, \Sigma_k)$.

The parameters θ are the μ_k, Σ_k, ϕ_k for all $k = 1, 2, \dots, K$.

Review: Generalization

In machine learning, **generalization** is the property of predictive models to achieve good performance on new, heldout data that is distinct from the training set.

How does generalization apply to probabilistic unsupervised models like GMMs?

An Unsupervised Learning Dataset

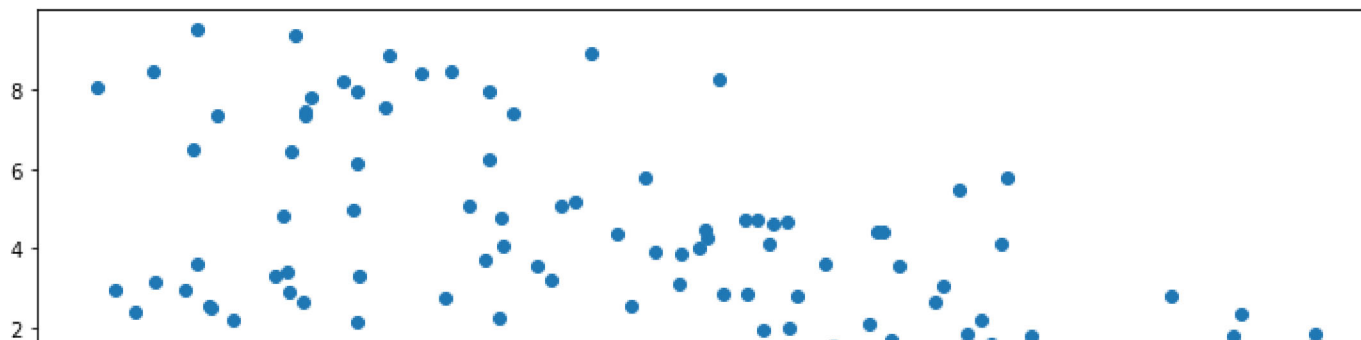
Consider the following dataset, consisting of a mixture of Gaussians.

```
In [49]: import numpy as np
from sklearn import datasets
from matplotlib import pyplot as plt
plt.rcParams['figure.figsize'] = [12, 4]

# generate 150 random points
np.random.seed(0)
X_all, y_all = datasets.make_blobs(150, centers=4)

# use the first 100 points as the main dataset
X, y = X_all[:100], y_all[:100]
plt.scatter(X[:,0], X[:,1])
```

Out[49]: <matplotlib.collections.PathCollection at 0x12b583780>



We know the true labels of these clusters, and we can visualize them.

We will also keep 50 points as a holdout set.

Underfitting in Unsupervised Learning

Underfitting happens when we are not able to fully learn the signal hidden in the data.

In the context of GMMs, this means not capturing all the clusters in the data.

Let's fit a GMM on our toy dataset.

```
In [52]: # fit a GMM  
from sklearn import mixture  
model = mixture.GaussianMixture(n_components=2)  
model.fit(X)
```

```
Out[52]: GaussianMixture(n_components=2)
```

The model finds two distinct components in the data, but they fail to capture the true structure.

We can also measure the value of our objective (the log-likelihood) on the training and holdout sets.

Consider now what happens if we further increase the number of clusters.

Overfitting in Unsupervised Learning

Overfitting happens when we fit the noise, but not the signal.

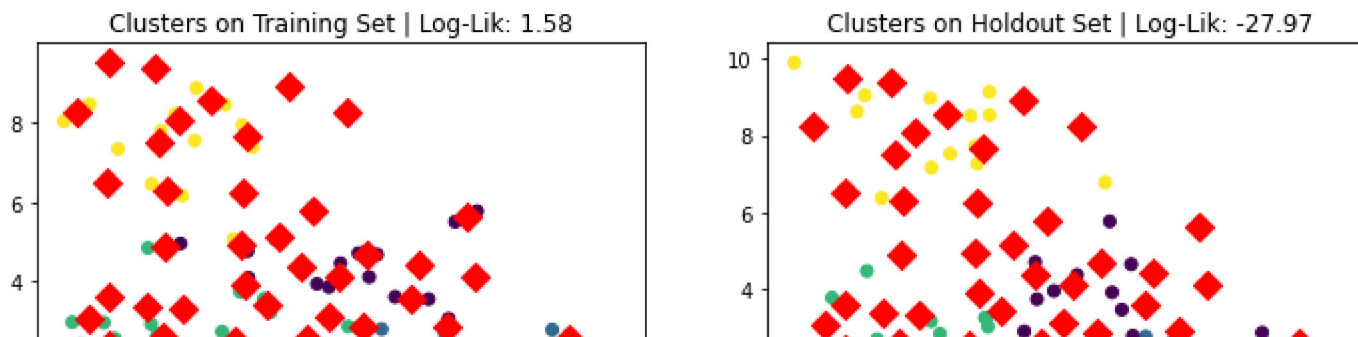
In our example, this means fitting small, local noise clusters rather than the true global clusters.

```
In [55]: model = mixture.GaussianMixture(n_components=50)
model.fit(X)

plt.subplot(121)
plt.title('Clusters on Training Set | Log-Lik: %.2f' % model.score(X))
plt.scatter(X[:,0], X[:,1], c=y)
plt.scatter(model.means_[0], model.means_[1], marker='D', c='r', s=100)

plt.subplot(122)
plt.title('Clusters on Holdout Set | Log-Lik: %.2f' % model.score(X_holdout))
plt.scatter(X_holdout[:,0], X_holdout[:,1], c=y_holdout)
plt.scatter(model.means_[0], model.means_[1], marker='D', c='r', s=100)
```

Out[55]: <matplotlib.collections.PathCollection at 0x12aeb95c0>



Measuring Generalization Using Log-Likelihood

Probabilistic unsupervised models optimize an objective that can be used to detect overfitting and underfitting by comparing performance between training and holdout sets.

Below, we visualize the performance (measured via negative log-likelihood) on training and holdout sets as K increases.

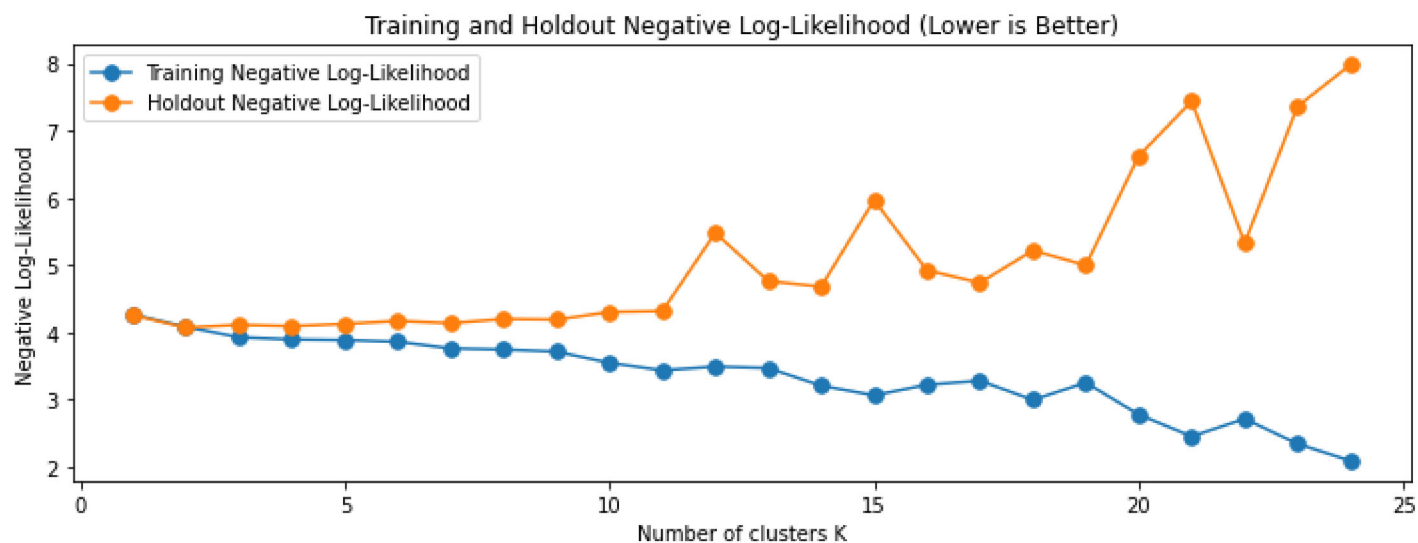
```

In [56]: Ks, training_objs, holdout_objs = range(1,25), [], []
for k in Ks:
    model = mixture.GaussianMixture(n_components=k)
    model.fit(X)
    training_objs.append(-model.score(X))
    holdout_objs.append(-model.score(X_holdout))

plt.plot(Ks, training_objs, '.-', markersize=15)
plt.plot(Ks, holdout_objs, '.-', markersize=15)
plt.xlabel("Number of clusters K")
plt.ylabel("Negative Log-Likelihood")
plt.title("Training and Holdout Negative Log-Likelihood (Lower is Better)")
plt.legend(['Training Negative Log-Likelihood', 'Holdout Negative Log-Likelihood'
])

```

Out[56]: <matplotlib.legend.Legend at 0x12c463320>



Warning: This process doesn't work as well as in supervised learning

For example, detecting overfitting with larger datasets will be paradoxically harder (try it!)

Summary

- Generalization is important for supervised and unsupervised learning.
- A probabilistic model can detect overfitting by comparing the likelihood of training data vs. that of holdout data.
- We can reduce overfitting by making the model less expressive.