

Machine Learning

Unsupervised Learning

dsai.asia

Asian Data Science and Artificial Intelligence Master's Program



Co-funded by the
Erasmus+ Programme
of the European Union



Readings for these lecture notes:

- Bishop, C. (2006), *Pattern Recognition and Machine Learning*, Springer, Chapter 9.
- Hastie, T., Tibshirani, R., and Friedman, J. (2016), *Elements of Statistical Learning: Data Mining, Inference, and Prediction*, Springer, Chapters 13, 14.
- Ng, A. (2017), *Learning Theory*. Lecture note set 7 for CS229, Stanford University.

These notes contain material © Bishop (2006), Hastie et al. (2016), and Ng (2017).

Outline

- 1 Introduction
- 2 k -means
- 3 EM for Gaussian mixture models
- 4 Other unsupervised learning algorithms

Now we consider the problem of **unsupervised learning** in which, rather than trying to predict some target given x , we want to understand the relationship of x to the examples in our unlabeled training set.

We begin with k -means clustering then move to other models.

Outline

- 1 Introduction
- 2 k -means
- 3 EM for Gaussian mixture models
- 4 Other unsupervised learning algorithms

k -means

The algorithm

In a **clustering** problem, we are given a training set $S = \{x^{(1)}, \dots, x^{(m)}\}$, $x^{(i)} \in \mathbb{R}^n$ and we want to group the $x^{(i)}$ into cohesive clusters.

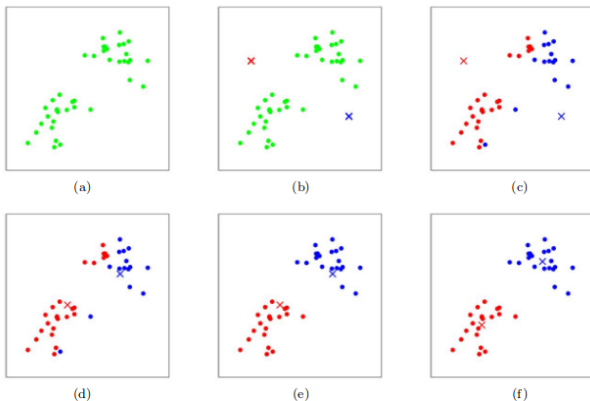
The **k -means** clustering algorithm is as follows:

- 1 Randomly initialize k **cluster centroids** $\mu_1, \dots, \mu_k \in \mathbb{R}^n$.
- 2 Repeat until convergence:
 - 1 For $i \in 1..m$, $c^{(i)} \leftarrow \operatorname{argmin}_j \|x^{(i)} - \mu_j\|^2$.
 - 2 For $j \in 1..k$,

$$\mu_j \leftarrow \frac{\sum_{i=1}^m \delta(c^{(i)} = j) x^{(i)}}{\sum_{i=1}^m \delta(c^{(i)} = j)}.$$

k -means

Example



Ng (2017), CS229 lecture notes

(a) Training data. (b) Initial mean assignment. (c) Iteration 1 cluster assignment. (d) Iteration 1 mean computation. (e) Iteration 2 cluster assignment. (f) Iteration 2 mean computation.

As an exercise, generate data from three Gaussians with different means and covariances and produce an animation of k -means in action.

k-means

Convergence

k-means is equivalent to **coordinate descent** on c and the μ_i 's for cost function or **distortion function**

$$J(c, \theta) = \sum_{i=1}^m \|x^{(i)} - \mu_{c(i)}\|^2$$

with

$$\theta = [\mu_1 \quad \cdots \quad \mu_k] .$$

Prove this by finding the minimum with respect to c holding the μ_i fixed then finding the minimum with respect to the μ_i holding the c fixed.

The method is thus guaranteed to converge in the sense that on every iteration, J will monotonically decrease.

As J is not convex, k-means is not guaranteed to converge to a global minimum. Local minima are possible.

Outline

- 1 Introduction
- 2 k -means
- 3 EM for Gaussian mixture models**
- 4 Other unsupervised learning algorithms

EM for Gaussian mixture models

Setting

Similar to the k -means setting, suppose we are given a training set $S = \{x^{(1)}, \dots, x^{(m)}\}$ without labels.

Rather than simply find cluster centers, however, now our goal is to estimate a joint distribution $p(x^{(i)}, z^{(i)}) = p(x^{(i)} | z^{(i)})p(z^{(i)})$.

The **latent** (hidden) variable $z^{(i)}$ indicates **which of k clusters or components** the example $x^{(i)}$ was drawn from.

Assumptions:

- $z^{(i)} \sim \text{Multinomial}(\phi)$ with $\phi_j \geq 0$ and $\sum_{j=1}^k \phi_j = 1$.
- $x^{(i)} | (z^{(i)} = j) \sim \mathcal{N}(\mu_j, \Sigma_j)$.

EM for Gaussian mixture models

Likelihood

As usual, we try to estimate the parameters $\phi, \mu_1, \dots, \Sigma_1, \dots$ via maximum likelihood, first by writing down the (log) likelihood function:

$$\begin{aligned}\ell(\phi, \mu_1, \dots, \Sigma_1, \dots) &= \sum_{i=1}^m \log p(x^{(i)}; \phi, \mu_1, \dots, \Sigma_1, \dots) \\ &= \sum_{i=1}^m \log \sum_{j=1}^k p(x^{(i)} \mid z^{(i)} = j; \mu_j, \Sigma_j) p(z^{(i)} = j; \phi) \\ &= \sum_{i=1}^m \log \sum_{j=1}^k \phi_j \mathcal{N}(x^{(i)}; \mu_j, \Sigma_j)\end{aligned}$$

There is no closed form solution to this maximization problem unfortunately...

EM for Gaussian mixture models

Maximization

But, what if we knew the $z^{(i)}$'s? Could we find the maximum likelihood solution for the μ 's and Σ 's?

The answer is yes...

$$\begin{aligned}\ell(\phi, \mu_1, \dots, \Sigma_1, \dots) &= \sum_{i=1}^m \left[\log p(\mathbf{x}^{(i)} \mid z^{(i)}; \mu_{z^{(i)}}, \Sigma_{z^{(i)}}) + \log p(z^{(i)}; \phi) \right] \\ &= \sum_{i=1}^m \left[\log \phi_{z^{(i)}} + \log \mathcal{N}(\mathbf{x}^{(i)}; \mu_{z^{(i)}}, \Sigma_{z^{(i)}}) \right]\end{aligned}$$

Take the derivatives with respect to each parameter, set them to 0, and solve for the parameters. What do you get?

EM for Gaussian mixture models

Maximization with respect to ϕ_j

If we set up the equation $\frac{\partial \ell}{\partial \phi_j} = 0$, we get close to a solution as $\phi_j \rightarrow \infty$.

This violates the constraint $\sum_{j=1}^k \phi_j = 1$.

Try instead, then, to introduce a Lagrange multiplier for the constraint and find a stationary point of

$$\ell(\phi, \mu_1, \dots, \Sigma_1, \dots) + \lambda \left(\sum_{j=1}^k \phi_j - 1 \right).$$

Maximizing this with the full likelihood expression, you should be able to obtain

$$\sum_{i=1}^m \frac{\mathcal{N}(\mathbf{x}^{(i)}; \mu_{z(i)}, \Sigma_{z(i)})}{\sum_{l=1}^k \phi_l \mathcal{N}(\mathbf{x}^{(i)}; \mu_{z(i)}, \Sigma_{z(i)})} + \lambda = 0.$$

EM for Gaussian mixture models

Maximization with respect to ϕ_j

Here's bit of a magic trick: try multiplying both sides of the equation by ϕ_j then summing over j .

Can you get $\lambda = -m$?

Next, substitute $\lambda = -m$ and use the form of $\ell()$ assuming known $z^{(i)}$ and solve for ϕ_j .

EM for Gaussian mixture models

Maximization with respect to μ_j and Σ_j

Setting up the equation $\frac{\partial \ell}{\partial \mu_j} = 0$ with the likelihood assuming known $z^{(i)}$, and using Equation 81¹ from *The Matrix Cookbook*, we get

$$\sum_{i=1}^m \delta(z^{(i)} = j) \Sigma_j^{-1} (\mathbf{x}^{(i)} - \mu_j) = 0.$$

Doing the same with $\frac{\partial \ell}{\partial \Sigma_j} = 0$ and Equations 72² and 49³ from the *Cookbook*, we get

$$\sum_{i=1}^m \delta(z^{(i)} = j) \Sigma_j^{-1} \left((\mathbf{x}^{(i)} - \mu_j)(\mathbf{x}^{(i)} - \mu_j)^\top - \Sigma_j \right) \Sigma_j^{-1} = 0.$$

¹ $\frac{\partial \mathbf{x}^\top \mathbf{B} \mathbf{x}}{\partial \mathbf{x}} = (\mathbf{B} + \mathbf{B}^\top) \mathbf{x}$

² $\frac{\partial \mathbf{a}^\top \mathbf{x} \mathbf{a}}{\partial \mathbf{a}} = \mathbf{a} \mathbf{a}^\top$

³ $\frac{\partial |\mathbf{X}|}{\partial \mathbf{X}} = |\mathbf{X}| \mathbf{X}^{-\top}$

EM for Gaussian mixture models

Maximization

Finally, we obtain

$$\begin{aligned}\phi_j &= \frac{1}{m} \sum_{i=1}^m \delta(z^{(i)} = j), \\ \mu_j &= \frac{\sum_{i=1}^m \delta(z^{(i)} = j) \mathbf{x}^{(i)}}{\sum_{i=1}^m \delta(z^{(i)} = j)}, \\ \Sigma_j &= \frac{\sum_{i=1}^m \delta(z^{(i)} = j) (\mathbf{x}^{(i)} - \mu_j)(\mathbf{x}^{(i)} - \mu_j)^\top}{\sum_{i=1}^m \delta(z^{(i)} = j)}.\end{aligned}$$

If you look back, you'll see this is very similar to the Gaussian discriminant analysis model...

EM for Gaussian mixture models

EM

OK, that's great, but we don't know the $z^{(i)}$'s!!

To address this issue, we use the EM (Expectation Maximization) algorithm.

EM is a general approach to solving maximum likelihood problems when there are latent (unobserved/hidden) variables.

In the E step, we try to guess values of the hidden variables, then in the M step, we update the parameters of the model based on those guesses.

EM for Gaussian mixture models

EM

Here's the EM algorithm, made specific for our Gaussian mixtures problem:

❶ Repeat until convergence:

❶ (E-step) For each $i \in 1..m, j \in 1..k$,

$$w_j^{(i)} \leftarrow p(z^{(i)} = j \mid \mathbf{x}^{(i)}; \phi, \mu_1, \dots, \Sigma_1, \dots)$$

❷ (M-step) Update the parameters

$$\phi_j \leftarrow \frac{1}{m} \sum_{i=1}^m w_j^{(i)},$$

$$\mu_j \leftarrow \frac{\sum_{i=1}^m w_j^{(i)} \mathbf{x}^{(i)}}{\sum_{i=1}^m w_j^{(i)}},$$

$$\Sigma_j \leftarrow \frac{\sum_{i=1}^m w_j^{(i)} (\mathbf{x}^{(i)} - \mu_j)(\mathbf{x}^{(i)} - \mu_j)^\top}{\sum_{i=1}^m w_j^{(i)}}$$

EM for Gaussian mixture models

More on the E step

The $w_j^{(i)}$'s are our soft guesses as to the membership of each example $x^{(i)}$ in component j of the mixture.

How to calculate them? We have

$$w_j^{(i)} \leftarrow p(z^{(i)} = j \mid x^{(i)}; \phi, \mu_1, \dots, \Sigma_1, \dots).$$

Let's use Bayes rule to turn this into something we know how to calculate:

$$p(z^{(i)} = j \mid x^{(i)}; \phi, \mu_1, \dots, \Sigma_1, \dots) = \frac{p(x^{(i)} \mid z^{(i)} = j; \mu_j, \Sigma_j) p(z^{(i)} = j; \phi)}{\sum_{l=1}^k p(x^{(i)} \mid z^{(i)} = l; \mu_l, \Sigma_l) p(z^{(i)} = l; \phi)}.$$

EM for Gaussian mixture models

More on the E step

By the way, Bishop (2006) calls these $w_j^{(i)}$'s **responsibilities**, denoted γ_{ij} .

γ_{ij} is the degree to which example i is compatible with or could have been generated by Gaussian j .

EM for Gaussian mixture models

Conclusion

You should compare the Gaussian mixture model to k -means.

Just like k -means, EM GMM is susceptible to **local minima**. We should perform several random restarts and retain the best model.

Exercise: implement the Gaussian mixture EM for a small dataset.

As mentioned earlier, this is a specialization of the general EM algorithm for estimating parameters when we have latent variables.

Next we'll look at the general EM algorithm and its guarantees on convergence.

Then we'll look at different forms of unsupervised learning algorithms.

Outline

- 1 Introduction
- 2 k -means
- 3 EM for Gaussian mixture models
- 4 Other unsupervised learning algorithms

Other unsupervised learning algorithms

Types of algorithms

In this brief overview, we've seen k -means and GMMs.

What else is out there? Most unsupervised learning algorithms have one or more of these aims:

- Clustering, for purposes of coding, discretizing, or detecting anomalies.
- Probability density estimation, for purposes of positive detection, anomaly detection, or interpolation/synthesis of new data distributed similarly to the training set.
- Latent space discovery, for purposes of dimensionality reduction, positive detection, or interpolation/synthesis of new data distributed similarly to the training set.

Other unsupervised learning algorithms

Clustering

We already saw how k -means and GMM can perform data clustering.

This style of clustering is sometimes called **vector quantization**. The goal is coding of inputs as members of a discrete set according to spatial locality.

A second major group of clustering algorithms are **pairwise clustering methods** that use pairwise similarity or distance (affinity) to group inputs.

Two types:

- **Hierarchical clustering** is top-down.
- **Agglomerative clustering** is bottom-up.

Other unsupervised learning algorithms

Density estimation

We saw the GMM as an example of **probability density estimation**.

Density estimators can be used for **detection** (any input whose probability density is larger than some value is classified as a positive).

Density estimators can be used for **anomaly detection** (any input whose probability density is below some value is classified as anomalous).

The GMM is a **parametric density estimator** and is a **universal approximator**.

Non-parametric density estimators like kernel density estimation do not try to model the exact form of the probability density but instead estimate density directly based on sparsity/density of training data near the target input.

Other unsupervised learning algorithms

Latent space discovery

Latent space methods map the input space to a lower-dimensional representation, called a **latent** or **semantic** representation.

Principle components analysis (PCA) models the latent space as Gaussian distribution in the k -dimensional linear subspace of the original space in which the input data have the most variance.

Locally-linear embedding (LLE) and other nonlinear dimensionality reduction methods model the data as generated from a non-linear submanifold of the input space. LLE uses a piecewise linear model.

Latent Dirichlet allocation (LDA) is a **topic model** in NLP that imposes a prior on the data distribution in the latent space.

Other unsupervised learning algorithms

GANs

Generative adversarial networks (GANs) represent the mapping from the latent space to the data space by a neural network, usually a deconvolutional neural network. This generative model is trained alongside a discriminative adversary.

Seminal paper: Goodfellow, I., Pouget-Abadie, J., Mirza, M., Xu, B., Warde-Farley, D., Ozair, S., Courville, A., and Bengio, Y. (2014), Generative adversarial nets, *Advances in Neural Information Processing Systems (NeurIPS 2014)*.

GANs are **generative** models, not discriminative. Recall the differences.

For unsupervised learning tasks, we want to model $P(x)$ given a training set $x^{(1)}, \dots, x^{(m)}$.

Other unsupervised learning algorithms

GANs

Historically, generative models were hard to estimate due to latent variables (like the $z^{(i)}$'s in the GMM) making maximum likelihood estimation difficult.

GANs take a completely new approach involving playing a two-player minimax game with value function

$$\min_G \max_D V(D, G) = \mathbb{E}_{x \sim p_{\text{data}}(x)} [\log D(x)] + \mathbb{E}_{z \sim p_z(z)} [\log(1 - D(G(z)))]$$

G and D are neural networks.

In practice, it works slightly better to have G maximize $\log D(G(z))$ rather than minimizing $\log(1 - D(G(z)))$ and it accomplishes the same thing but provides stronger gradients early in learning.

Other unsupervised learning algorithms

The GAN training algorithm

for number of training iterations do

for k steps do

Sample minibatch of m noise samples $\{z^{(1)}, \dots, z^{(m)}\}$ from noise prior $p_g(z)$.

Sample minibatch of m examples $\{x^{(1)}, \dots, x^{(m)}\}$ from the data generating distribution $p_{\text{data}}(x)$.

Update the discriminator using stochastic gradient ascent:

$$\nabla_{\theta_d} \frac{1}{m} \sum_{i=1}^m \left[\log D(x^{(i)}) + \log (1 - D(G(z^{(i)}))) \right]$$

Sample minibatch of m noise samples $\{z^{(1)}, \dots, z^{(m)}\}$ from noise prior $p_g(z)$.

Update the generator using stochastic gradient descent:

$$\nabla_{\theta_g} \frac{1}{m} \sum_{i=1}^m \log (1 - D(G(z^{(i)})))$$

Other unsupervised learning algorithms

GAN exercise

In-class exercise: Using PyTorch, build and train a GAN to model the synthetic 2D mixture of Gaussians from Lab 12.

Don't use any of PyTorch's specialized GAN classes. Just use ordinary fully connected multilayer perceptrons for G and D.

Other unsupervised learning algorithms

Summary

Unsupervised learning is an extremely rich area, deserving an entire course on its own!