STA 314: Statistical Methods for Machine Learning I Lecture 7 - Unsupervised Learning, *K*-Means

Chris J. Maddison

University of Toronto

- What can we do without labels?
- How can we even define what learning means without labels?
- Unsupervised learning is the study of learning without labels.
- In some sense, the ML community does not exactly agree on what it means to do unsupervised learning, but intuitively, unsupervised learning is the task of grouping, explaining, and finding structure data.

Motivating Examples

- Some examples of situations where you'd use unupservised learning
 - You want to understand how a scientific field has changed over time. You want to take a large database of papers and model how the distribution of topics changes from year to year. But what are the topics?
 - You're a biologist studying animal behavior, so you want to infer a high-level description of their behavior from video. You don't know the set of behaviors ahead of time.
 - You want to reduce your energy consumption, so you take a time series of your energy consumption over time, and try to break it down into separate components (refrigerator, washing machine, etc.).

• Common themes:

- you have some data, and you want to infer some structure underlying the data.
- you have some data, and you want to be able to make more data that looks similar

- Today we will look at the simplest type of unsupervised learning: clustering.
- Clustering is the task of organizing data into groups or clusters.
- We will study the simplest method for doing this: the *K*-means algorithm.

Clustering

• Sometimes the data form clusters, where samples within a cluster are similar to each other, and samples in different clusters are dissimilar:



- Such a distribution is multimodal, since it has multiple modes, or regions of high probability mass.
- Grouping data points into clusters, with no observed labels, is called clustering. It is an unsupervised learning technique.
- E.g. clustering machine learning papers based on topic (deep learning, Bayesian models, etc.)
 - But topics are never observed (unsupervised).



- Assume the data $\{\mathbf{x}^{(1)}, \dots, \mathbf{x}^{(N)}\}$ lives in a Euclidean space, $\mathbf{x}^{(n)} \in \mathbb{R}^{D}$.
- Assume each data point belongs to one of K clusters
- Assume the data points from same cluster are similar, i.e. close in Euclidean distance.
- How can we identify those clusters (data points that belong to each cluster)? Let's formulate as an optimization problem.

Intro ML (UofT)

K-means Objective



K-means Objective: Find centers $\{\mathbf{m}_k\}_{k=1}^{K}$ and assignments $\{\mathbf{r}^{(n)}\}_{n=1}^{N}$ to minimize the sum of squared distances of data points $\{\mathbf{x}^{(n)}\}$ to their assigned centers.

- Data sample n = 1, ..., N: $\mathbf{x}^{(n)} \in \mathbb{R}^D$ (observed),
- Cluster center k = 1, ..., K: $\mathbf{m}_k \in \mathbb{R}^D$ (not observed),
- Responsibilities: Cluster assignment for sample *n*: $\mathbf{r}^{(n)} \in \mathbb{R}^{K}$ 1-of-K encoding (not observed)

K-means Objective

- K-means Objective: Find cluster centers {m_k}^K_{k=1} and assignments {r⁽ⁿ⁾}^N_{n=1} to minimize the sum of squared distances of data points {x⁽ⁿ⁾} to their assigned centers.
 - ▶ Data sample n = 1, .., N: $\mathbf{x}^{(n)} \in \mathbb{R}^D$ (observed),
 - Cluster center k = 1, ..., K: $\mathbf{m}_k \in \mathbb{R}^D$ (not observed),
 - ▶ Responsibilities: Cluster assignment for sample *n*: $\mathbf{r}^{(n)} \in \mathbb{R}^{K}$ 1-of-K encoding (not observed)
- Mathematically:

$$\min_{\{\mathbf{m}_k\},\{\mathbf{r}^{(n)}\}} \hat{\mathcal{R}}(\{\mathbf{m}_k\},\{\mathbf{r}^{(n)}\}) = \min_{\{\mathbf{m}_k\},\{\mathbf{r}^{(n)}\}} \sum_{n=1}^{N} \sum_{k=1}^{K} r_k^{(n)} ||\mathbf{m}_k - \mathbf{x}^{(n)}||^2$$

where $r_k^{(n)} = \mathbb{I}[\mathbf{x}^{(n)} \text{ is assigned to cluster } k]$, i.e., $\mathbf{r}^{(n)} = [0, .., 1, .., 0]^{\top}$

• Finding an optimal solution is an NP-hard problem!

K-means Objective

Optimization problem:

$$\min_{\{\mathbf{m}_k\},\{\mathbf{r}^{(n)}\}} \sum_{n=1}^{N} \underbrace{\sum_{k=1}^{K} r_k^{(n)} ||\mathbf{m}_k - \mathbf{x}^{(n)}||^2}_{\text{distance between } x^{(n)}_{\text{end its assigned cluster center}}}$$

Since r_k⁽ⁿ⁾ = I[**x**⁽ⁿ⁾ is assigned to cluster k], i.e., **r**⁽ⁿ⁾ = [0, .., 1, .., 0][⊤]
inner sum is over K terms but only one of them is non-zero.
E.g. say sample **x**⁽ⁿ⁾ is assigned to cluster k = 3, then

$$\mathbf{r}^n = [0, 0, 1, 0, ...]$$

$$\sum_{k=1}^{K} r_k^{(n)} ||\mathbf{m}_k - \mathbf{x}^{(n)}||^2 = ||\mathbf{m}_3 - \mathbf{x}^{(n)}||^2$$

How to optimize?: Alternating Minimization

Optimization problem:

$$\min_{\{\mathbf{m}_k\},\{\mathbf{r}^{(n)}\}} \sum_{n=1}^{N} \sum_{k=1}^{K} r_k^{(n)} ||\mathbf{m}_k - \mathbf{x}^{(n)}||^2$$

- Problem is hard when minimizing jointly over the parameters $\{\mathbf{m}_k\}, \{\mathbf{r}^{(n)}\}$
- But note that if we fix one and minimize over the other, then it becomes easy.
- Doesn't guarantee the same solution!

How to optimize?: Alternating Minimization

Optimization problem:

$$\min_{\{\mathbf{m}_k\},\{\mathbf{r}^{(n)}\}} \sum_{n=1}^{N} \sum_{k=1}^{K} r_k^{(n)} ||\mathbf{m}_k - \mathbf{x}^{(n)}||^2$$

• Fix the centers {**m**_k} then we can easily find the optimal assignments {**r**⁽ⁿ⁾} for each sample *n*

$$\min_{\mathbf{r}^{(n)}} \sum_{k=1}^{K} r_k^{(n)} ||\mathbf{m}_k - \mathbf{x}^{(n)}||^2$$

Assign each point to the cluster with the nearest center

$$r_k^{(n)} = \begin{cases} 1 & \text{if } k = \arg\min_j \|\mathbf{x}^{(n)} - \mathbf{m}_j\|^2\\ 0 & \text{otherwise} \end{cases}$$

Only \hat{k} -th entry is 1

• E.g. if $\mathbf{x}^{(n)}$ is assigned to cluster \hat{k} , $\mathbf{r}^{(n)} = [0, 0, ..., 1, ..., 0]^{\top}$

Intro ML (UofT)

Alternating Minimization

- Likewise, if we fix the assignments $\{\mathbf{r}^{(n)}\}$ then can easily find optimal centers $\{\mathbf{m}_k\}$
 - Set each cluster's center to the average of its assigned data points: For l = 1, 2, ..., K

$$0 = \frac{\partial}{\partial \mathbf{m}_{l}} \sum_{n=1}^{N} \sum_{k=1}^{K} r_{k}^{(n)} ||\mathbf{m}_{k} - \mathbf{x}^{(n)}||^{2}$$
$$= 2 \sum_{n=1}^{N} r_{l}^{(n)}(\mathbf{m}_{l} - \mathbf{x}^{(n)}) \implies \mathbf{m}_{l} = \frac{\sum_{n} r_{l}^{(n)} \mathbf{x}^{(n)}}{\sum_{n} r_{l}^{(n)}}$$

- Let's alternate between minimizing $\hat{\mathcal{R}}(\{\mathbf{m}_k\}, \{\mathbf{r}^{(n)}\})$ with respect to $\{\mathbf{m}_k\}$ and $\{\mathbf{r}^{(n)}\}$
- This is called alternating minimization

High level overview of algorithm:

- Initialization: randomly initialize cluster centers
- The algorithm iteratively alternates between two steps:
 - Assignment step: Assign each data point to the closest cluster
 - Refitting step: Move each cluster center to the mean of the data assigned to it





Figure from Bishop

Simple demo: http://syskall.com/kmeans.js/

Intro ML (UofT)

The K-means Algorithm

- Initialization: Set K cluster means $\mathbf{m}_1, \ldots, \mathbf{m}_K$ to random values
- Repeat until convergence (until assignments do not change):
 - Assignment: Optimize w.r.t. {r}: Each data point x⁽ⁿ⁾ assigned to nearest center

$$\hat{k}^{(n)} = \arg\min_{k} ||\mathbf{m}_k - \mathbf{x}^{(n)}||^2$$

and Responsibilities (1-hot or 1-of-K encoding)

$$r_k^{(n)} = \mathbb{I}[\hat{k}^{(n)} = k]$$
 for $k = 1,..,K$

▶ Refitting: Optimize $\hat{\mathcal{R}}$ w.r.t. {m}: Each center is set to mean of data assigned to it

$$\mathbf{m}_k = \frac{\sum_n r_k^{(n)} \mathbf{x}^{(n)}}{\sum_n r_k^{(n)}}.$$

K-means for Vector Quantization



Figure from Bishop

- Given image, construct "dataset" of pixels represented by their RGB pixel intensities
- Run k-means, replace each pixel by its cluster center

Intro ML (UofT)

K-means for Image Segmentation



- Given image, construct "dataset" of pixels, represented by their RGB pixel intensities and grid locations
- Run k-means (with some modifications) to get superpixels

Intro ML (UofT)

- Why does update set **m**_k to mean of assigned points?
- What if we used a different distance measure?
- How can we choose the best distance?
- How to choose *K*?
- Will it converge?

Hard cases - unequal spreads, non-circular spreads, in-between points

Why K-means Converges

- K-means algorithm reduces the cost at each iteration.
 - Whenever an assignment is changed, the sum squared distances R̂ of data points from their assigned cluster centers is reduced.
 - Whenever a cluster center is moved, $\hat{\mathcal{R}}$ is reduced.
- Test for convergence: If the assignments do not change in the assignment step, we have converged (to at least a local minimum).
- This will always happen after a finite number of iterations, since the number of possible cluster assignments is finite



• K-means cost function after each assignment step (blue) and refitting step (red). The algorithm has converged after the third refitting step.

- The objective *Â* is non-convex (so coordinate descent on *Â* is not guaranteed to converge to the global minimum)
- There is nothing to prevent k-means getting stuck at local minima.
- We could try many random starting points

A bad local optimum



- Instead of making hard assignments of data points to clusters, we can make soft assignments. One cluster may have a responsibility of .7 for a datapoint and another may have a responsibility of .3.
 - Allows a cluster to use more information about the data in the refitting step.
 - How do we decide on the soft assignments?
 - We already saw this in multi-class classification:
 - ▶ 1-of-K encoding vs softmax assignments

Soft K-means Algorithm

- Initialization: Set K means $\{\mathbf{m}_k\}$ to random values
- Repeat until convergence (measured by how much $\hat{\mathcal{R}}$ changes):
 - Assignment: Each data point n given soft "degree of assignment" to each cluster mean k, based on responsibilities

$$r_k^{(n)} = \frac{\exp[-\beta \|\mathbf{m}_k - \mathbf{x}^{(n)}\|^2]}{\sum_j \exp[-\beta \|\mathbf{m}_j - \mathbf{x}^{(n)}\|^2]}$$
$$\implies \mathbf{r}^{(n)} = \operatorname{softmax}(-\beta \{\|\mathbf{m}_k - \mathbf{x}^{(n)}\|^2\}_{k=1}^K)$$

Refitting: Model parameters, means, are adjusted to match sample means of datapoints they are responsible for:

$$\mathbf{m}_k = \frac{\sum_n r_k^{(n)} \mathbf{x}^{(n)}}{\sum_n r_k^{(n)}}$$

Some remaining issues

- How to set β ?
- Clusters with unequal weight and width?

These aren't straightforward to address with K-means. Instead, in the sequel, we'll reformulate clustering using a generative model.

As $\beta \to \infty$, soft k-Means becomes k-Means! (Exercise)

- We now turn to the second unsupervised learning algorithm for this course: principal component analysis (PCA)
- Dimensionality reduction: map data to a lower dimensional space
 - Save computation/memory
 - Reduce overfitting, achieve better generalization
 - Visualize in 2 dimensions
- PCA is a linear model. It's useful for understanding lots of other algorithms.
 - Autoencoders
 - Matrix factorizations (next week)
- PCA is finding linear structure in data.

Gaussian distribution

• First, we will review the Gaussian, or normal, distribution:

$$\mathcal{N}(x;\mu,\sigma^2) = rac{1}{\sqrt{2\pi}\sigma} \exp\left(-rac{(x-\mu)^2}{2\sigma^2}
ight)$$

- The Central Limit Theorem says that sums of lots of independent random variables are approximately Gaussian.
- In machine learning, we use Gaussians a lot because they make the calculations easy.



- Multiple measurements (sensors)
- D inputs/features/attributes
- N instances/observations/examples

$$\mathbf{X} = \begin{bmatrix} [\mathbf{x}^{(1)}]^{\top} \\ [\mathbf{x}^{(2)}]^{\top} \\ \vdots \\ [\mathbf{x}^{(N)}]^{\top} \end{bmatrix} = \begin{bmatrix} x_1^{(1)} & x_2^{(1)} & \cdots & x_D^{(1)} \\ x_1^{(2)} & x_2^{(2)} & \cdots & x_D^{(2)} \\ \vdots & \vdots & \ddots & \vdots \\ x_1^{(N)} & x_2^{(N)} & \cdots & x_D^{(N)} \end{bmatrix}$$

Multivariate Mean and Covariance

 $oldsymbol{\mu} = \mathbb{E}[\mathbf{x}] = egin{pmatrix} \mu_1 \ dots \ \mu_d \end{pmatrix}$

Covariance

Mean

$$\boldsymbol{\Sigma} = \operatorname{Cov}(\mathbf{x}) = \mathbb{E}[(\mathbf{x} - \boldsymbol{\mu})(\mathbf{x} - \boldsymbol{\mu})^{\top}] = \begin{pmatrix} \sigma_1^2 & \sigma_{12} & \cdots & \sigma_{1D} \\ \sigma_{12} & \sigma_2^2 & \cdots & \sigma_{2D} \\ \vdots & \vdots & \ddots & \vdots \\ \sigma_{D1} & \sigma_{D2} & \cdots & \sigma_D^2 \end{pmatrix}$$

- The statistics (μ and Σ) uniquely define a multivariate Gaussian (or multivariate Normal) distribution, denoted N(μ, Σ) or N(x; μ, Σ)
 - This is not true for distributions in general!

Intro ML (UofT)

Multivariate Gaussian Distribution

• Normally distributed variable $\mathbf{x} \sim \mathcal{N}(\boldsymbol{\mu}, \boldsymbol{\Sigma})$ has distribution:

$$p(\mathbf{x}) = rac{1}{(2\pi)^{d/2} |\mathbf{\Sigma}|^{1/2}} \exp\left[-rac{1}{2} (\mathbf{x} - \boldsymbol{\mu})^T \mathbf{\Sigma}^{-1} (\mathbf{x} - \boldsymbol{\mu})
ight]$$





Gaussian Intuition: (Multivariate) Shift + Scale

- Recall that in the univariate case, all normal distributions are shaped like the standard normal distribution
- The densities are related to the standard normal by a shift (μ), a scale (or stretch, or dilation) σ , and a normalization factor



Gaussian Intuition: (Multivariate) Shift + Scale

- The same intuition applies in the multivariate case.
- We can think of the multivariate Gaussian as a shifted and "scaled" version of the standard multivariate normal distribution.
 - ullet The standard multivariate normal has $m{\mu}=m{0}$ and $m{\Sigma}=m{I}$
- ullet Multivariate analog of the shift is simple: it's a vector $m \mu$
- But what about the scale?
 - ▶ In the univariate case, the scale factor was the square root of the variance: $\sigma = \sqrt{\sigma^2}$
 - But in the multivariate case, the covariance Σ is a matrix! Does Σ^{1/2} exist, and can we scale by it?

Multivariate Scaling (Intuitive)

We call a matrix "positive definite" if it is symmetric and scales the space in **orthogonal** directions. The univariate analog is positive scalar $\alpha > 0$. Consider, e.g., how these two matrices transform the orthogonal vectors:

Consider
matrix:
$$\begin{pmatrix} 2 & 0 \\ 0 & 0.5 \end{pmatrix}$$
 $\begin{pmatrix} 1 & 0.5 \\ 0.5 & 1 \end{pmatrix}$ Consider
action on: $\begin{pmatrix} 1 \\ 0 \end{pmatrix} \perp \begin{pmatrix} 0 \\ 1 \end{pmatrix}$ $\begin{pmatrix} 1 \\ 1 \end{pmatrix} \perp \begin{pmatrix} 1 \\ -1 \end{pmatrix}$

Draw action on slide:

Notice: both matrices are symmetric!

Multivariate Scaling (Formal) (details optional)

We summarize some definitions/results from linear algebra (without proof). Knowing them is optional, but they may help with intuition (esp. for PCA).

- Definition. Symmetric matrix A is positive semidefinite if x[⊤]Ax ≥ 0 for all non-zero x. It is positive definite if x[⊤]Ax > 0 for all non-zero x.
 - Any positive definite matrix is positive semidefinite.
 - Positive definite matrices have positive eigenvalues, and positive semidefinite matrices have non-negative eigenvalues.
 - ▶ For any matrix **X**, **X**^T**X** and **XX**^T are positive semidefinite.
- Theorem (Unique Positive Square Root). Let A be a positive semidefinite real matrix. Then there is a unique positive semidefinite matrix B such that A = B^TB = BB. We call A^{1/2} ≜ B the positive square root of A.
- **Theorem** (*Spectral Theorem*). If $\mathbf{A} \in \mathbb{R}^{d \times d}$ is symmetric, then
 - 1. \mathbb{R}^{D} has an orthonormal basis consisting of the eigenvectors of **A**.
 - 2. There exists orthonormal matrix ${f Q}$ and diagonal matrix ${f \Lambda}$ such that
 - $\mathbf{A} = \mathbf{Q} \mathbf{\Lambda} \mathbf{Q}^{\mathsf{T}}$. This is called the spectral decomposition of \mathbf{A} .
 - ► The columns of **Q** are (unit) eigenvectors of **A**.

Key properties of Σ :

- 1. $\boldsymbol{\Sigma}$ is positive semidefinite (and therefore symmetric).
- 2. For a distribution with density, $\pmb{\Sigma}$ is positive definite.

Other properties (optional / for reference):

3.
$$\mathbf{\Sigma} = \mathbb{E}[\mathbf{x}\mathbf{x}^{ op}] - \boldsymbol{\mu}\boldsymbol{\mu}^{ op}$$
 (generalizes $\operatorname{Var}(x) = \mathbb{E}[x^2] - \mu^2$))

4.
$$\mathsf{Cov}(\mathbf{A}\mathbf{x} + \mathbf{b}) = \mathbf{A} \mathbf{\Sigma} \mathbf{A}^{ op}$$
 (generalizes $\mathsf{Var}(ax + b) = a^2 \mathsf{Var}(x)$)

So here is the "scale" intuition:

- For positive definite Σ , consider its unique positive square root $\Sigma^{\frac{1}{2}}$.
- $\Sigma^{\frac{1}{2}}$ is also positive definite, so by the Real Spectral Theorem, it "scales" the space in orthogonal directions (its eigenvectors) by its eigenvalues.

Intro ML (UofT)





Intro ML (UofT)



Probability density function



Intro ML (UofT)

$$\begin{split} \mathbf{\Sigma} &= \begin{pmatrix} 1 & 0 \\ 0 & 1 \end{pmatrix} \qquad \mathbf{\Sigma} = \begin{pmatrix} 1 & 0.5 \\ 0.5 & 1 \end{pmatrix} \qquad \mathbf{\Sigma} = \begin{pmatrix} 1 & 0.8 \\ 0.8 & 1 \end{pmatrix} \\ &= \mathbf{Q}_1 \begin{pmatrix} 1.5 & 0. \\ 0. & 0.5 \end{pmatrix} \mathbf{Q}_1^\top \qquad = \mathbf{Q}_2 \begin{pmatrix} 1.8 & 0. \\ 0. & 0.2 \end{pmatrix} \mathbf{Q}_2^\top \end{split}$$



$$\boldsymbol{\Sigma} = \begin{pmatrix} 1 & 0 \\ 0 & 1 \end{pmatrix} \qquad \boldsymbol{\Sigma} = \begin{pmatrix} 1 & 0.5 \\ 0.5 & 1 \end{pmatrix} \qquad \boldsymbol{\Sigma} = \begin{pmatrix} 1 & -0.5 \\ -0.5 & 1 \end{pmatrix}$$
$$= \boldsymbol{Q}_1 \begin{pmatrix} 1.5 & 0. \\ 0. & 0.5 \end{pmatrix} \boldsymbol{Q}_1^\top \qquad = \boldsymbol{Q}_2 \begin{pmatrix} \boldsymbol{\lambda}_1 & 0. \\ 0. & \boldsymbol{\lambda}_2 \end{pmatrix} \boldsymbol{Q}_2^\top$$



37 / 40



 $Cov(x_1, x_2)=0, Var(x_1)=Var(x_2)$

 $Cov(x_1, x_2)=0, Var(x_1)>Var(x_2)$







Intro ML (UofT)

- Next weel: PCA.
- Dimensionality reduction: map data to a lower dimensional space
- PCA is finding linear structure in data.