Projects

- Chandrasekar, Arun Kumar, Group 17
- Nearly all group have submitted a proposal
- May 21: Each person gives one slide, 15 min/group.

	First principles	VS	Data driv	en	
Data	Small data		Big data to trai	n	
Domain expertise	High reliance on do expertise	High reliance on domain expertise		Results with little domain knowledge	
Fidelity/ Robustness	Universal link can handle non- linear complex relations		Limited by the range of values spanned by training data		
Adaptability	Complex and time consuming derivation to use new relations		Rapidly adapt to new problems		
Interpretability	Parameters are physical!		Physically agnostic, limited by the rigidity of the functional form		
Perceived Importance.	SIO	SP	Peter	Google	

Machine learning versus knowledge based



Physical/Knowledge Base

Supervised learning



Supervised $\mathbf{y}=\mathbf{w}^{\mathsf{T}}\mathbf{x}$ Training set { $(x^1, y^1), (x^2, y^2), (x^3, y^3)$ } We are given the two classes.



Supervised y=wx Training set $\{(x_1^1, x_2^1), (x_1^2, x_2^2), (x_1^3, x_2^3)\}$

Unsupervised learning

Unsupervised machine learning is inferring a function to describe hidden structure from "unlabeled" data (a classification or categorization is not included in the observations). Since the examples given to the learner are unlabeled, there is no evaluation of the accuracy of the structure that is output by the relevant algorithm—which is one way of distinguishing unsupervised learning from <u>supervised learning</u>.

We are not interested in prediction

Supervised learning: all classification and regression. $Y = W^T X$

Prediction is important.

Unsupervised learning

- Unsupervised learning is more subjective than supervised learning, as there is no simple goal for the analysis, such as prediction of a response.
- But techniques for unsupervised learning are of growing importance in several fields:
 - subgroups of breast cancer patients grouped by their gene expression measurements,
 - groups of shoppers characterized by their browsing and purchase histories,
 - movies grouped by the ratings assigned by movie viewers.
- It is often easier to obtain unlabeled data from a lab instrument or a computer — than labeled data, which can require human intervention.
 - For example it is difficult to automatically assess the overall sentiment of a movie review: is it favorable or not?

Kmeans

- Input: Points $\mathbf{x}_1, \dots, \mathbf{x}_N \in \mathbb{R}^p$; integer K
- **Output**: "Centers", or representatives, $\mu_1, ..., \mu_K \in \mathbb{R}^p$
- Output also $\mathbf{z}_1, \dots, \mathbf{z}_N \in \mathbb{R}^K$

Goal: Minimize average squared distance between points and their nearest representatives:

•
$$cost(\mu_1, ..., \mu_K) = \sum_{n=1}^N \min_j ||x_n - \mu_j||$$



The centers carve \mathbb{R}^p up into k convex regions: μ_j 's region consists of points for which it is the closest center.

K-means

$$J = \sum_{n=1}^{N} \sum_{k=1}^{K} r_{nk} \|\mathbf{x}_n - \boldsymbol{\mu}_k\|^2$$
(9.1)

Solving for r_{nk}

$$r_{nk} = \begin{cases} 1 & \text{if } k = \arg\min_{j} \|\mathbf{x}_{n} - \boldsymbol{\mu}_{j}\|^{2} \\ 0 & \text{otherwise.} \end{cases}$$
(9.2)

Differentiating for μ_k

$$2\sum_{n=1}^{N} r_{nk} (\mathbf{x}_n - \boldsymbol{\mu}_k) = 0$$
(9.3)

.

which we can easily solve for μ_k to give

$$\boldsymbol{\mu}_{k} = \frac{\sum_{n} r_{nk} \mathbf{x}_{n}}{\sum_{n} r_{nk}}.$$
(9.4)



Old Faithful, Kmeans from Murphy



Example



The progress of the K-means algorithm with K=3.

- Top left: The observations are shown.
- *Top center:* In Step 1 of the algorithm, each observation is randomly assigned to a cluster.
- Top right: In Step 2(a), the cluster centroids are computed. These are shown as large colored disks. Initially the centroids are almost completely overlapping because the initial cluster assignments were chosen at random.
- *Bottom left:* In Step 2(b), each observation is assigned to the nearest centroid.
- *Bottom center:* Step 2(a) is once again performed, leading to new cluster centroids.
- Bottom right: The results obtained after 10 iterations.

Likely From Hastie book

Different starting values



K-means clustering performed six times on the data from previous figure with K = 3, each time with a different random assignment of the observations in Step 1 of the K-means algorithm.

Above each plot is the value of the objective (4). Three different local optima were obtained, one of which resulted in a smaller value of the objective and provides better separation between the clusters.

Those labeled in red all achieved the same best solution, with an objective value of $235.8\,$

Likely From Hastie book

Vector Quantization VQ

Murphy book Fig 11.12 vqdemo.m

Each pixel \mathbf{x}_i is represented By codebook of K entries μ_k

 $\mathsf{Encode}(\mathbf{x}_{\mathsf{i}}) = \underset{k}{\operatorname{argmin}} \|x_i - \mu_k\|$

Consider N=64k observations, of D=1 (b/w) dimension, C=8 bit

NC=513k

 $Nlog_2$ K+KC bits is needed K=4 gives 128k a factor 4.



Mixtures of Gaussians (1)

Old Faithful geyser:

The time between eruptions has a <u>bimodal distribution</u>, with the mean interval being either 65 or 91 minutes, and is dependent on the length of the prior eruption. Within a margin of error of ±10 minutes, Old Faithful will erupt either 65 minutes after an eruption lasting less than 2 $\frac{1}{2}$ minutes, or 91 minutes after an eruption lasting more than 2 $\frac{1}{2}$ minutes.



Mixtures of Gaussians (2)



Mixtures of Gaussians (3)



• Gaussian mixture - $p(x) = \sum_{k}^{K} \pi_{k} N(x; \mu_{k}, \Sigma_{k})$

- Latent variable:
 - Un-observed
 - Often hidden
- Here $p(z_k) = \pi_k$



$p(\boldsymbol{x}|\boldsymbol{z}_{k} = 1) = N(\boldsymbol{x}; \boldsymbol{\mu}_{k}, \boldsymbol{\Sigma}_{k})$ $p(\boldsymbol{x}|\boldsymbol{z}) =$ $p(\boldsymbol{x}, \boldsymbol{z}) =$ $p(\boldsymbol{x}) =$

Responsibilities

$$\gamma(z_k) = p(z_k = 1 | \boldsymbol{x}) =$$

Mixture of Gaussians

• Mixtures of Gaussians

$$p(\mathbf{x}) = \sum_{k=1}^{K} \pi_k \mathcal{N}(\mathbf{x} | \boldsymbol{\mu}_k, \boldsymbol{\Sigma}_k).$$

Expressed with latent variable z

$$p(\mathbf{x}) = \sum_{\mathbf{z}} p(\mathbf{z}) p(\mathbf{x} | \mathbf{z}) = \sum_{k=1}^{K} \pi_k \mathcal{N}(\mathbf{x} | \boldsymbol{\mu}_k, \boldsymbol{\Sigma}_k)$$

Posterior probability: responsibility



Max Likelihood

- $p(x) = \sum_{k}^{K} \pi_{k} N(x; \mu_{k}, \Sigma_{k})$
- *N* observations *X*
- $\ln[p(\boldsymbol{X}|\boldsymbol{\pi},\boldsymbol{\mu},\boldsymbol{\Sigma})] = \prod_{N} \ln[\sum_{k}^{K} \pi_{k} N(x_{n};\boldsymbol{\mu}_{k},\boldsymbol{\Sigma}_{k})]$





EM Gauss Mix

- 1. Initialize the means μ_k , covariances Σ_k and mixing coefficients π_k , and evaluate the initial value of the log likelihood.
- 2. E step. Evaluate the responsibilities using the current parameter values

$$\gamma(z_{nk}) = \frac{\pi_k \mathcal{N}(\mathbf{x}_n | \boldsymbol{\mu}_k, \boldsymbol{\Sigma}_k)}{\sum_{j=1}^K \pi_j \mathcal{N}(\mathbf{x}_n | \boldsymbol{\mu}_j, \boldsymbol{\Sigma}_j)}.$$
(9.23)

3. M step. Re-estimate the parameters using the current responsibilities

$$\boldsymbol{\mu}_{k}^{\text{new}} = \frac{1}{N_{k}} \sum_{n=1}^{N} \gamma(z_{nk}) \mathbf{x}_{n}$$
(9.24)

$$\boldsymbol{\Sigma}_{k}^{\text{new}} = \frac{1}{N_{k}} \sum_{n=1}^{N} \gamma(z_{nk}) \left(\mathbf{x}_{n} - \boldsymbol{\mu}_{k}^{\text{new}} \right) \left(\mathbf{x}_{n} - \boldsymbol{\mu}_{k}^{\text{new}} \right)^{\text{T}} \qquad (9.25)$$

$$\pi_k^{\text{new}} = \frac{N_k}{N} \tag{9.26}$$

where

$$N_k = \sum_{n=1}^N \gamma(z_{nk}). \tag{9.27}$$

4. Evaluate the log likelihood

$$\ln p(\mathbf{X}|\boldsymbol{\mu}, \boldsymbol{\Sigma}, \boldsymbol{\pi}) = \sum_{n=1}^{N} \ln \left\{ \sum_{k=1}^{K} \pi_k \mathcal{N}(\mathbf{x}_n | \boldsymbol{\mu}_k, \boldsymbol{\Sigma}_k) \right\}$$
(9.28)

and check for convergence of either the parameters or the log likelihood. If the convergence criterion is not satisfied return to step 2.

General EM

Given a joint distribution $p(\mathbf{X}, \mathbf{Z}|\boldsymbol{\theta})$ over observed variables \mathbf{X} and latent variables \mathbf{Z} , governed by parameters $\boldsymbol{\theta}$, the goal is to maximize the likelihood function $p(\mathbf{X}|\boldsymbol{\theta})$ with respect to $\boldsymbol{\theta}$.

1. Choose an initial setting for the parameters θ^{old} .

- 2. E step Evaluate $p(\mathbf{Z}|\mathbf{X}, \boldsymbol{\theta}^{\text{old}})$.
- 3. **M step** Evaluate θ^{new} given by

$$\boldsymbol{\theta}^{\text{new}} = \operatorname*{arg\,max}_{\boldsymbol{\theta}} \mathcal{Q}(\boldsymbol{\theta}, \boldsymbol{\theta}^{\text{old}})$$
 (9.32)

where

$$\mathcal{Q}(\boldsymbol{\theta}, \boldsymbol{\theta}^{\text{old}}) = \sum_{\mathbf{Z}} p(\mathbf{Z} | \mathbf{X}, \boldsymbol{\theta}^{\text{old}}) \ln p(\mathbf{X}, \mathbf{Z} | \boldsymbol{\theta}).$$
(9.33)

4. Check for convergence of either the log likelihood or the parameter values. If the convergence criterion is not satisfied, then let

$$\boldsymbol{\theta}^{\mathrm{old}} \leftarrow \boldsymbol{\theta}^{\mathrm{new}}$$
 (9.34)

and return to step 2.

EM in general

$$p(\mathbf{X}|\boldsymbol{\theta}) = \sum_{\mathbf{Z}} p(\mathbf{X}, \mathbf{Z}|\boldsymbol{\theta}).$$
(9.69)

$$\ln p(\mathbf{X}|\boldsymbol{\theta}) = \mathcal{L}(q, \boldsymbol{\theta}) + \mathrm{KL}(q||p)$$
(9.70)

where we have defined

$$\mathcal{L}(q, \theta) = \sum_{\mathbf{Z}} q(\mathbf{Z}) \ln \left\{ \frac{p(\mathbf{X}, \mathbf{Z} | \theta)}{q(\mathbf{Z})} \right\}$$
(9.71)

$$\operatorname{KL}(q||p) = -\sum_{\mathbf{Z}} q(\mathbf{Z}) \ln \left\{ \frac{p(\mathbf{Z}|\mathbf{X}, \boldsymbol{\theta})}{q(\mathbf{Z})} \right\}.$$
(9.72)

$$\ln p(\mathbf{X}, \mathbf{Z}|\boldsymbol{\theta}) = \ln p(\mathbf{Z}|\mathbf{X}, \boldsymbol{\theta}) + \ln p(\mathbf{X}|\boldsymbol{\theta})$$
(9.73)

$$\mathcal{L}(q, \boldsymbol{\theta}) = \sum_{\mathbf{Z}} p(\mathbf{Z} | \mathbf{X}, \boldsymbol{\theta}^{\text{old}}) \ln p(\mathbf{X}, \mathbf{Z} | \boldsymbol{\theta}) - \sum_{\mathbf{Z}} p(\mathbf{Z} | \mathbf{X}, \boldsymbol{\theta}^{\text{old}}) \ln p(\mathbf{Z} | \mathbf{X}, \boldsymbol{\theta}^{\text{old}})$$

$$= \mathcal{Q}(\boldsymbol{\theta}, \boldsymbol{\theta}^{\text{old}}) + \text{const}$$
(9.74)

Gaussian Mixtures









Hierarchical Clustering

- K-means clustering requires us to pre-specify the number of clusters K. This can be a disadvantage (later we discuss strategies for choosing K)
- *Hierarchical clustering* is an alternative approach which does not require that we commit to a particular choice of *K*.
- In this section, we describe *bottom-up* or *agglomerative* clustering. This is the most common type of hierarchical clustering, and refers to the fact that a dendrogram is built starting from the leaves and combining clusters up to the trunk.

Hierarchical Clustering Algorithm

The approach in words:

- Start with each point in its own cluster.
- Identify the closest two clusters and merge them.
- Repeat.
- Ends when all points are in a single cluster.



An Example



45 observations generated in 2-dimensional space. In reality there are three distinct classes, shown in separate colors. However, we will treat these class labels as unknown and will seek to cluster the observations in order to discover the classes from the data.

Example



45 observations generated in 2-dimensional space. In reality there are three distinct classes, shown in separate colors. However, we will treat these class labels as unknown and will seek to cluster the observations in order to discover the classes from the data.



- *Left:* Dendrogram obtained from hierarchically clustering the data from previous slide, with complete linkage and Euclidean distance.
- *Center:* The dendrogram from the left-hand panel, cut at a height of 9 (indicated by the dashed line). This cut results in two distinct clusters, shown in different colors.
- *Right:* The dendrogram from the left-hand panel, now cut at a height of 5. This cut results in three distinct clusters, shown in different colors. Note that the colors were not used in clustering, but are simply used for display purposes in this figure

•NOT USED

K-means clustering



NOT

INTERESTING

A simulated data set with 150 observations in 2-dimensional space. Panels show the results of applying K-means clustering with different values of K, the number of clusters. The color of each observation indicates the cluster to which it was assigned using the K-means clustering algorithm. Note that there is no ordering of the clusters, so the cluster coloring is arbitrary. These cluster labels were not used in clustering; instead, they are the outputs of the clustering procedure.

Properties of the Algorithm

• This algorithm is guaranteed to decrease the value of the objective (4) at each step. *Why?* Note that

$$\frac{1}{|C_k|} \sum_{i,i' \in C_k} \sum_{j=1}^p (x_{ij} - x_{i'j})^2 = 2 \sum_{i \in C_k} \sum_{j=1}^p (x_{ij} - \bar{x}_{kj})^2,$$

where $\bar{x}_{kj} = \frac{1}{|C_k|} \sum_{i \in C_k} x_{ij}$ is the mean for feature j in cluster C_k .

• however it is not guaranteed to give the global minimum. Why not?

K-Means Clustering Algorithm

- 1. Randomly assign a number, from 1 to K, to each of the observations. These serve as initial cluster assignments for the observations.
- 2. Iterate until the cluster assignments stop changing:
 - 2.1 For each of the K clusters, compute the cluster *centroid*. The kth cluster centroid is the vector of the p feature means for the observations in the kth cluster.
 - 2.2 Assign each observation to the cluster whose centroid is closest (where *closest* is defined using Euclidean distance).

Clustering

- *Clustering* refers to a very broad set of techniques for finding *subgroups*, or *clusters*, in a data set.
- We seek a partition of the data into distinct groups so that the observations within each group are quite similar to each other,
- It make this concrete, we must define what it means for two or more observations to be *similar* or *different*.
- Indeed, this is often a domain-specific consideration that must be made based on knowledge of the data being studied.

Mixture of Experts



Figure 11.6 (a) Some data fit with three separate regression lines. (b) Gating functions for three different "experts". (c) The conditionally weighted average of the three expert predictions. Figure generated by mixexpDemo.

 The key idea is that each expert focus on predicting the right answer for cases where they are already doing better than other experts.



{X,Z}:Complete, {X}: Incomplete, responsibilities







Hierarchical Clustering

- K-means clustering requires us to pre-specify the number of clusters K. This can be a disadvantage (later we discuss strategies for choosing K)
- *Hierarchical clustering* is an alternative approach which does not require that we commit to a particular choice of *K*.
- In this section, we describe *bottom-up* or *agglomerative* clustering. This is the most common type of hierarchical clustering, and refers to the fact that a dendrogram is built starting from the leaves and combining clusters up to the trunk.

Hierarchical Clustering Algorithm

The approach in words:

- Start with each point in its own cluster.
- Identify the closest two clusters and merge them.
- Repeat.
- Ends when all points are in a single cluster.



An Example



45 observations generated in 2-dimensional space. In reality there are three distinct classes, shown in separate colors. However, we will treat these class labels as unknown and will seek to cluster the observations in order to discover the classes from the data.

Example



45 observations generated in 2-dimensional space. In reality there are three distinct classes, shown in separate colors. However, we will treat these class labels as unknown and will seek to cluster the observations in order to discover the classes from the data.



- *Left:* Dendrogram obtained from hierarchically clustering the data from previous slide, with complete linkage and Euclidean distance.
- *Center:* The dendrogram from the left-hand panel, cut at a height of 9 (indicated by the dashed line). This cut results in two distinct clusters, shown in different colors.
- *Right:* The dendrogram from the left-hand panel, now cut at a height of 5. This cut results in three distinct clusters, shown in different colors. Note that the colors were not used in clustering, but are simply used for display purposes in this figure



Mixture of Experts



Figure 11.8 (a) Some data from a simple forwards model. (b) Some data from the inverse model, fit with a mixture of 3 linear regressions. Training points are color coded by their responsibilities. (c) The predictive mean (red cross) and mode (black square). Based on Figures 5.20 and 5.21 of (Bishop 2006b). Figure generated by mixexpDemoOneToMany.

Two clustering methods

- In *K*-means clustering, we seek to partition the observations into a pre-specified number of clusters.
- In *hierarchical clustering*, we do not know in advance how many clusters we want; in fact, we end up with a tree-like visual representation of the observations, called a *dendrogram*, that allows us to view at once the clusterings obtained for each possible number of clusters, from 1 to n.