Text Clustering, K-Means, Gaussian Mixture Models, Expectation-Maximization, Hierarchical Clustering

Sameer Maskey

Week 3, Sept 19, 2012

Topics for Today

- Text Clustering
- Gaussian Mixture Models
- K-Means
- Expectation Maximization
- Hierarchical Clustering

Announcement

Proposal Due tonight (11:59pm) – not graded

Feedback by Friday

- Final Proposal due (11:59pm) next Wednesday
 - □ 5% of the project grade
 - Email me the proposal with the title
 - "Project Proposal : Statistical NLP for the Web"
- Homework 1 is out
 - Due October 4th (11:59pm) Thursday
- Please use courseworks

Course Initial Survey

100.00% 95.45% 95.45% 90.91% 90.48% 90.00% 80.00% 77.27% 77.27% 72.73% 70.00% 63.64% 63.64% Percentage (Yes, No) 59.09% 59.09% 60.00% 54.55% 50.550%00% Yes 50.00% 45.45% ■ No 40.91 40.91% 40.00% 36.36% 36.36% 27.27% 30.00% 22.73 22.73% 20.00% 9.09 9.09% 10.00% 4.55% 4.55% 0.00% NLP SLP ML NLP for ML Adv ML NLP-ML Pace Math Matlab Matlab Excited for Industry Larger Tutorial Project Audience Mentors

Class Survey

Category

Perceptron Algorithm

We are given (x_i, y_i) Initialize wDo until converged if $\operatorname{error}(y_i, sign(w.x_i)) == TRUE$ $w \leftarrow w + y_i x_i$ end if End do

If predicted class is wrong, subtract or add that point to weight vector

Perceptron (cont.)

$$y_j(t) = f[w(t).x_j]$$

 $w_i(t+1) = w_i(t) + \alpha(d_j - y_j(t))x_{i,j}$
Error is either 1, 0 or -1

Input					laitial waighta		Output						Correction	Final weights		
Sen	sor va	alues	Desired output	mua	arve	Per sensor Sum Netwo		Network	LIIO	Conection	r mai weights					
x_0	$ x_1 $	x_2	z	w_0	w_1	w_2	c_0	c_1	c_2	s	n	e	d	w_0	w_1	w_2
				_			$x_0 * w_0$	$x_1 * w_1$	$x_2 * w_2$	$c_0 + c_1 + c_2$	if $s>t$ then 1, else 0	z-n	r * e	$\Delta(x_0 * d$	$\Delta(x_1 * d)$	$\Delta(x_2 * d)$
1	0	0	1	0.4	0	0.1	0.4	0	0	0.4	0	1	+0.1	0.5	0	0.1
1	0	1	1	0.5	0	0.1	0.5	0	0.1	0.6	1	0	0	0.5	0	0.1
1	1	0	1	0.5	0	0.1	0.5	0	0	0.5	0	1	+0.1	0.6	0.1	0.1
1	1	1	0	0.6	0.1	0.1	0.6	0.1	0.1	0.8	1	-1	-0.1	0.5	0	0

Example from Wikipedia

Naïve Bayes Classifier for Text

$$P(Y = y_k | X_1, X_2, ..., X_N) = \frac{P(Y = y_k) P(X_1, X_2, ..., X_N | Y = y_k)}{\sum_j P(Y = y_j) P(X_1, X_2, ..., X_N | Y = y_j)}$$
$$= \frac{P(Y = y_k) \prod_i P(X_i | Y = y_k)}{\sum_j P(Y = y_j) \prod_i P(X_i | Y = y_j)}$$

$$Y \leftarrow argmax_{y_k} P(Y = y_k) \Pi_i P(X_i | Y = y_k)$$

Naïve Bayes Classifier for Text

Given the training data what are the parameters to be estimated?



Data without Labels

Data with corresponding Human Scores



Document Clustering

- Previously we classified Documents into Two Classes
 - Diabetes (Class1) and Hepatitis (Class2)
- We had human labeled data
 - Supervised learning
- What if we do not have manually tagged documents
 - Can we still classify documents?
 - Document clustering
 - Unsupervised Learning

Classification vs. Clustering



Supervised Training of Classification Algorithm Unsupervised Training of Clustering Algorithm

Clusters for Classification



Automatically Found Clusters can be used for Classification



Document Clustering

- Cluster the documents in 'N' clusters/categories
- For classification we were able to estimate parameters using labeled data
 - Perceptrons find the parameters that decide the separating hyperplane
 - Naïve Bayes count the number of times word occurs in the given class and normalize
- Not evident on how to find separating hyperplane when no labeled data available
- Not evident how many classes we have for data when we do not have labels

Document Clustering Application

 Even though we do not know human labels automatically induced clusters could be useful
 News Clusters



Document Clustering Application





A Map of Yahoo!, Mappa.Mundi Magazine, February 2000. Map of the Market with Headlines Smartmoney [2]

How to Cluster Documents with No Labeled Data?

- Treat cluster IDs or class labels as hidden variables
- Maximize the likelihood of the unlabeled data
- Cannot simply count for MLE as we do not know which point belongs to which class
 - User Iterative Algorithm such as K-Means, EM

Hidden Variables? What do we mean by this?

Hidden vs. Observed Variables

Assuming our observed data is in R2



How many observed variables?

How many observed variables? How many hidden variables?

lustering

If we have data with labels

data comes from 2 classes

for both classes

Find out μ_i and \sum_i from data for both classes





K-Means Clustering

Let us define Dataset in D dimension $\{x_1, x_2, ..., x_N\}$

We want to cluster the data in K clusters

Let μ_k be D dimension vector representing cluster K

Let us define r_{nk} for each x_n such that $r_{nk} \in \{0, 1\}$ where k = 1, ..., K and $r_{nk} = 1$ if x_n is assigned to cluster k

Distortion Measure

$$J = \sum_{n=1}^{N} \sum_{k=1}^{K} r_{nk} ||x_n - \mu_k||^2$$

Represents sum of squares of distances to mu_k from each data point

We want to minimize J

Estimating Parameters

- We can estimate parameters by doing 2 step iterative process
 - Image: Minimize J with respect to r_{nk}
 - Keep μ_k fixed
 - Minimize J with respect to μ_k
 - Keep r_{nk} fixed



Step 1



- Optimize for each n separately by choosing r_{nk} for k that gives minimum $||x_n - r_{nk}||^2$

$$r_{nk} = 1$$
 if $k = argmin_j ||x_n - \mu_j||^2$
= 0 otherwise

- Assign each data point to the cluster that is the closest
- Hard decision to cluster assignment



- J is quadratic in μ_k . Minimize by setting derivative w.rt. μ_k to zero

$$\mu_k = \frac{\sum_n r_{nk} x_n}{\sum_n r_{nk}}$$

 Take all the points assigned to cluster K and re-estimate the mean for cluster K

Document Clustering with K-means

- Assuming we have data with no labels for Hockey and Baseball data
- We want to be able to categorize a new document into one of the 2 classes (K=2)
- We can extract represent document as feature vectors
 - Features can be word id or other NLP features such as POS tags, word context etc (D=total dimension of Feature vectors)
 - N documents are available
- Randomly initialize 2 class means
- Compute square distance of each point (x_n)(D dimension) to class means (µ_k)
- Assign the point to K for which μ_k is lowest
- Re-compute μ_k and re-iterate

K-Means Example



Clusters

Number of documents clustered together

U.S.	~	Sci/Tech	
Top Stories			Chama to push White House vision for MSA in April
Starred 🙀		15	Reuters - Bernd Debusmann, Jeff Mason - 23 minutes ago
World		A shall	President Barack Obama speaks about heatthcare reform from the East Room of the White
U.S.		diene see	House in Washington March 3, 2010. WASHINGTON (Reuters) President Barack Obama will
Business		BigPond News	President Obama In Florida on April 15, to Elaborate On Vew NASA Initiatives AHN All
Sei/Teeb			Headline News
Sci/Tech			Obama sets conference on future of space program The Associated Press
Entertainment			Baltimore Sur - Wall Street Journal - Sydney Morning Herald - CTV.ca
Sports			
Health		iPad	$ m mathac{10}{10}$ Issues Apple Needs to Address Before Releasing the iPad
Spotlight			eWeek - 2 hours ago
Most Popular			News Analysis: The iPad is now less than a month away from hitting store shelves, but there are
			Jobs: iPad Won't Tether with iPhone. PC Magazine
> All news		New York Times	All about the Apple (Pad (FAQ) CNET
<u>Headlines</u>		(blog)	Wired News - Wall Street Journal - PC World - Computerworld
Images			all 1,49/ news articles » MEmail this story
			☆ Panasonic Announces Lumix DMC-G2 and G10
		an Lange	Slippery Brick - Darrin Olson - 1 hour ago
		4	Panasonic announced on Sunday two new Micro Four-Thirds cameras, the Lumix DMC-G2 and a
		100	less expensive Lumix G10. Both cameras are in Panasonic's line of "smaller" digital cameras in comparison to D SLP's, going without a mirror box or a dedicated
		TrustedReviews	Panasonic's G series gets serious CNET
			Hands On: Panasonic's Micro Four Thirds Touchscreen Camera PC Magazine
			PC World - infoSync World - Digital Photography Review (dpreview.com) - DigitalCameraInfo
			all o'l news articles » Mitmail this story
			☆ Microsoft demos game across PC, mobile, and console platforms
			CNET - Kyle VanHemert - 18 hours ago
		No. 1 Bears	Whoa. During the keynote presentation at TechEd Middle East in Dubai, Microsoft's Eric Rudder
		pitric or	played the same indiana Jones-ish game on a Windows computer, a Windows Phone 7 phone, and an Xhox 360
		SlashGear (blog)	Microsoft Showcases Cross-Platform Gaming for Windows Pho. PC Magazine

Hard Assignment to Clusters

- K-means algorithm assigns each point to the closest cluster
 - Hard decision
 - Each data point affects the mean computation equally
- How does the points almost equidistant from 2 clusters affect the algorithm?
- Soft decision?
 - Fractional counts?

Gaussian Mixture Models (GMMs)





30

Mixtures of 2 Gaussians



Mixture Models



Mixture of Gaussians [1]

- I Gaussian may not fit the data
- 2 Gaussians may fit the data better
- Each Gaussian can be a class category
- When labeled data not available we can treat class category as hidden variable

Mixture Model Classifier

Given a new data point find out posterior probability from each class



Cluster ID/Class Label as Hidden Variables $p(x) = \sum_{z} p(x, z) = \sum_{z} p(z)p(x|z)$

- We can treat class category as hidden variable z
- Z is K-dimensional binary random variable in which $z_k = 1$ and 0 for other elements

z = |00100...|

$$\mathrm{p(z)} = \prod_{k=1}^{K} \pi_k^{z_k}$$

, sum of priors sum to 1 $\sum_{k=1}^{K} \pi_k = 1$

Also, sum of priors sum to 1 $\sum_{k=1}^{7}$

Conditional distribution of x given a particular z can be written as $P(x|z_k=1) = \mathcal{N}(x|\mu_k, \Sigma_k)$ $P(x|\overline{z}) = \prod_{k=1}^{K} \mathcal{N}(x|\mu_k, \sum_k)^{z_k}$

34

Mixture of Gaussians with Hidden Variables



$$p(x) = \sum_{k=1}^{K} \pi_k \frac{1}{(2\pi)^{D/2} \sqrt{(|\sum_k|}} e^{xp(-\frac{1}{2}(x-\mu_k)^T \sum_{k=1}^{-1} (x-\mu_k))} e^{xp(-\frac{1}{2}(x-\mu_k)^T \sum_{k=1}^{-1} (x-\mu_k))}$$

- Mixture models can be linear combinations of other distributions as well
- Mixture of binomial distribution for example

Conditional Probability of Label Given Data

- Mixture model with parameters mu, sigma and prior can represent the parameter
- We can maximize the data given the model parameters to find the best parameters
- If we know the best parameters we can estimate

$${}^{\circ}(z_k) \,\, ' \,\, \mathbf{p}(\mathbf{z}_k = 1 | x) = \frac{p(z_k = 1)p(x | z_k = 1)}{\sum_{j=1}^{K} p(z_j = 1)p(x | z_j = 1)}$$
$$= \frac{\pi_k \mathcal{N}(x | \mu_k, \sum_k)}{\sum_{j=1}^{K} \pi_j \mathcal{N}(x | \mu_j, \sum_j)}$$

This essentially gives us probability of class given the data i.e label for the given data point

Maximizing Likelihood

 If we had labeled data we could maximize likelihood simply by counting and normalizing to get mean and variance of Gaussians for the given classes

(55, 2)

$$l = \sum_{n=1}^{N} \log p(x_n, y_n | \pi, \mu, \Sigma)$$

$$l = \sum_{n=1}^{N} \log \pi_{y_n} \mathcal{N}(x_n | \mu_{y_n}, \sum_{y_n}) |_{(30, 1)}$$

If we have two classes C1 and C2

Let's say we have a feature x
x = number of words 'field'
And class label (y)
y = 1 hockey or 2 baseball documents $N(\mu_1, \sum_1)$ Find out μ_i and \sum_i from data $N(\mu_2, \sum_2)$

Maximizing Likelihood for Mixture Model with Hidden Variables

 For a mixture model with a hidden variable representing 2 classes, log likelihood is

$$l = \sum_{n=1}^{N} logp(x_n | \pi, \mu, \Sigma)$$
$$l = \sum_{n=1}^{N} log \sum_{y=0}^{1} \mathcal{N}(x_n, y | \pi, \mu, \Sigma)$$

$$= \sum_{n=1}^{N} \log \left(\pi_0 \mathcal{N}(x_n | \mu_0, \sum_0) + \pi_1 \mathcal{N}(x_n | \mu_1, \sum_1) \right)$$

Log-likelihood for Mixture of Gaussians

$\log p(X|\pi, \mu, \Sigma) = \sum_{n=1}^{N} \log \left(\sum_{k=1}^{k} \pi_k \mathcal{N}(x|\mu_k, \Sigma_k) \right)$

- We want to find maximum likelihood of the above loglikelihood function to find the best parameters that maximize the data given the model
- We can again do iterative process for estimating the loglikelihood of the above function
 - □ This 2-step iterative process is called Expectation-Maximization

Explaining Expectation Maximization



Expectation Maximization

An expectation-maximization (EM) algorithm is used in statistics for finding maximum likelihood estimates of parameters in probabilistic models, where the model depends on unobserved hidden variables.

EM alternates between performing an <u>expectation (E) step, which</u> <u>computes an expectation of the likelihood by including the latent</u> <u>variables as if they were observed</u>, and a <u>maximization (M) step,</u> <u>which computes the maximum likelihood estimates of the</u> <u>parameters by maximizing the expected likelihood found on the E</u> <u>step.</u> The parameters found on the M step are then used to begin another E step, and the process is repeated.

The EM algorithm was explained and given its name in a classic 1977 paper by A. Dempster and D. Rubin in the Journal of the Royal Statistical Society.

Estimating Parameters

$$o(z_{nk}) = E(z_{nk}|x_n) = p(z_k = 1|x_n)$$

$$^{\mathbf{o}}(z_{nk}) = \frac{\pi_k \mathcal{N}(x_n | \mu_k, \sum_k)}{\sum_{j=1}^K \pi_j \mathcal{N}(x_n | \mu_j, \sum_j)}$$

Estimating Parameters

M-step

$$\mu'_{k} = \frac{1}{N_{k}} \sum_{n=1}^{N} {}^{\circ}(z_{nk}) x_{n}$$

$$\sum_{k}' = \frac{1}{N_{k}} \sum_{n=1}^{N} {}^{\circ}(z_{nk}) (x_{n} - \mu'_{k}) (x_{n} - \mu'_{k})^{T}$$

$$\pi'_{k} = \frac{N_{k}}{N}$$

where $N_{k} = \sum_{n=1}^{N} {}^{\circ}(z_{nk})$
terate until convergence of log likelihood

 $\log p(X|\pi, \mu, \Sigma) = \sum_{n=1}^{N} \log \left(\sum_{k=1}^{k} \mathcal{N}(x|\mu_k, \Sigma_k) \right)$

EM Iterations



EM iterations [1]

Clustering Documents with EM

- Clustering documents requires representation of documents in a set of features
 - Set of features can be bag of words model
 - Features such as POS, word similarity, number of sentences, etc
- Can we use mixture of Gaussians for any kind of features?
- How about mixture of multinomial for document clustering?
- How do we get EM algorithm for mixture of multinomial?

Clustering Algorithms

We just described two kinds of clustering algorithms

- K-means
- Expectation Maximization
- Expectation-Maximization is a general way to maximize log likelihood for distributions with hidden variables

□ For example, EM for HMM, state sequences were hidden

 For document clustering other kinds of clustering algorithm exists

Hierarchical Clustering

- Build a binary tree that groups similar data in iterative manner
- K-means
 - distance of data point to center of the gaussian
- EM
 - Posterior of data point w.r.t to the gaussian
- Hierarchical
 - Similarity : ?
 - Similarity across groups of data

Types of Hierarchical Clustering

Agglomerative (bottom-up):

- Assign each data point as one cluster
- Iteratively combine sub-clusters
- Eventually, all data points is a part of 1 cluster

Divisive (top-down):

- Assign all data points to the same cluster.
- Eventually each data point forms its own cluster

One advantage : Do not need to define K, number of clusters before we begin clustering

Hierarchical Clustering Algorithm

- Step 1
 - Assign each data point to its own cluster

Step 2
Compute similarity between clusters
Step 3

Merge two most similar cluster to form one less cluster

Hierarchical Clustering Demo



Animation source [4]

Similar Clusters?

- How do we compute similar clusters?
 - Distance between 2 points in the clusters?
 - Distance from means of two clusters?
 - Distance between two closest points in the clusters?
- Different similarity metric could produce different types of cluster
- Common similarity metric used
 - Single Linkage
 - Complete Linkage
 - Average Group Linkage





Complete Linkage



Average Group Linkage



Hierarchical Cluster for Documents



Figure : [Ho, Qirong, et. al]

Hierarchical Document Clusters

- Highlevel multi view of the corpus
- Taxonomy useful for various purposes
 - Q&A related to a subtopic
 - Finding broadly important topics
 - Recursive drill down on topics
 - Filter irrelevant topics

Summary

- Unsupervised clustering algorithms
 - K-means
 - Expectation Maximization
 - Hierarchical clustering
- EM is a general algorithm that can be used to estimate maximum likelihood of functions with hidden variables
- Similarity Metric is important when clustering segments of text

References

- [1] Christopher Bishop, "Pattern Recognition and Machine Learning," 2006
- [2] <u>http://www.smartmoney.com/map-of-the-market/</u>
- [3] Ho, Qirong, et. al, Document Hierarchies from Text and Links, 2012