



CS M146 Discussion: Week 9 Naive Bayes, Clustering (K-Means, Gaussian Mixture Model)

Junheng Hao Friday, 03/05/2021







- Announcement
- Naive Bayes
- K-Means
- Gaussian Mixture Model





- **5:00 pm PST, Mar 5 (Friday):** Weekly Quiz 9 released on Gradescope.
- **11:59 pm PST, Mar 7 (Sunday):** Weekly Quiz 9 closed on Gradescope!
 - Start the quiz before **11:00 pm PST, Mar 7** to have the full 60-minute time
- **Grading update:** Lowest two quiz scores are dropped. The rest 7 quizzes are counted into final grading.
- **Problem set 4** released on CCLE, submission on Gradescope.
 - Please assign pages of your submission with corresponding problem set outline items on GradeScope.
 - You need to submit code and the results required by the problem set
 - Due on **next Friday, 11:59pm PST, Mar 12 (Friday)**

Late Submission of PS will NOT be accepted!



About Quiz 9



- Quiz release date and time: Mar 5, 2021 (Friday) 05:00 PM PST
- Quiz due/close date and time: Mar 7, 2021 (Sunday) 11:59 PM PST
- You will have up to **60 minutes** to take this exam. → Start before **11:00 PM** Sunday
- You can find the exam entry named "Week 9 Quiz" on GradeScope.
- Topics: Naive Bayes, Clustering
- Question Types
 - True/false, multiple choices
 - Some questions may include several subquestions.
- Some light calculations are expected. Some scratch paper and one scientific calculator (physical or online) are recommended for preparation.
- Note: This is the last quiz in this quarter. Highest 7 quiz scores are counted for final grading.

Prof. Sankararaman's post on updated quiz grading: https://campuswire.com/c/GB5E561C3/feed/438





- Open book and open notes, on GradeScope: "quiz"-like exam
- Start attempting the exam from 8:00 am PST on March 15; Submit your exam before
 8:00am PST March 16 (No extensions). → 24h time window
- Exam duration: **3 hours** (time limit after start the exam)
- Type: True/false and multiple choice questions (free text boxes are given for justification)
- The instructors will be available to provide clarifications on CampusWire (visible for everyone) from 8:00am-11:00am on March 15. Later questions on Campuswire may not be answered.
- Some calculations are expected.

MUST READ: Official post about final exam on Campuswire: <u>https://campuswire.com/c/GB5E561C3/feed/437</u>





• Defines a joint distribution

$$P(X = x, Y = c) = P(Y = c) \prod_{d=1}^{D} P(X_d = x_d | Y = c)$$

= $P(Y = c) \prod_k P(k | Y = c)^{z_k} = \pi_c \prod_k \theta_{ck}^{z_k}$

• Learning problem formulation

Training data
$$\mathcal{D} = \{(x_n, y_n)\}_{n=1}^{\mathsf{N}} \to \mathcal{D} = \{(\{z_{nk}\}_{k=1}^{\mathsf{K}}, y_n)\}_{n=1}^{\mathsf{N}}$$
 Constraints $\sum_c \pi_c = 1$ $\sum_k \theta_{ck} = \sum_k P(k|Y=c) = 1$
Objective $\mathcal{L} = \log P(\mathcal{D}) = \log \prod_{n=1}^{\mathsf{N}} \pi_{y_n} P(x_n|y_n) \longrightarrow (\pi_c^*, \theta_{ck}^*) = \arg \max \sum_n \log \pi_{y_n} + \sum_{n,k} z_{nk} \log \theta_{y_nk}$

• Solution

Δ

Engineer Change.

$$\pi_c^* = \frac{\#\text{of data points labeled as c}}{N}$$
$$\theta_{ck}^* = \frac{\#\text{of times word k shows up in data points labeled as c}}{\#\text{data points labeled as c}}$$





A short derivation of the maximum likelihood estimation

To maximize

$$\sum_{n:y_n=c,k} z_{nk} \log \theta_{ck}$$

We use the Lagrangian multiplier

$$\sum_{n:y_n=c,k} z_{nk} \log \theta_{ck} + \lambda \left(\sum_k \theta_{ck} - 1\right)$$

Taking derivatives with respect to θ_{ck} and then find the stationary point

$$\left(\sum_{n:y_n=c} \frac{z_{nk}}{\theta_{ck}}\right) + \lambda = 0 \rightarrow \theta_{ck} = -\frac{1}{\lambda} \sum_{n:y_n=c} z_{nk}$$

Apply constraint $\sum_k \theta_{ck} = 1$, plug in expression above for θ_{ck} , solve for λ , and plug back into expression for θ_{ck} :

$$\theta_{ck} = \frac{\sum_{n:y_n=c} z_{nk}}{\sum_k \sum_{n:y_n=c} z_{nk}}$$





• #docs in Class 1: 25, #docs in Class 2: 75

Index	Word	Count in Class 1	Count in Class 2
1	Ι	1	9
2	like	2	1
3	machine	3	9
4	learning	4	1





• How to obtain the parameters in Naive Bayes classifier? (shown in class)

Index	Word	# in C1	# in C2
1	I	1	9
2	like	2	1
3	machine	3	9
4	learning	4	1





• Predicting new document: {I:3, like:1, machine:5, learning:1} (shown in class)

Index	Word	# in C1	# in C2
1	Ι	1	9
2	like	2	1
3	machine	3	9
4	learning	4	1

One further question: How to predict new document: {I:3, like:1, machine:5, learning:1, love:2} → Label Smoothing





- A linear classifier (same as logistic regression)
- Generative model, modeling joint distribution (probabilities) → What is the model assumption?
- Pros:
 - Fast and simple compared to other complicated algorithms, easy training
 - Works well with high-dimension data such as text classification
- Cons:
 - Strong assumptions (feature independency)
 - Not fit to regression
 - Smoothing is somewhat required for generalization



Generative vs Discriminative Models



- Training classifiers involve estimating $f: X \rightarrow Y$, or P(Y|X)
- Generative classifiers → *"distribution"*
 - Assume some functional form for P(Y), P(X|Y)
 - Estimate parameters of P(X|Y), P(Y) directly from training data
 - \circ Use Bayes rule to calculate P(Y|X)
 - Actually learn the underlying structure of the data
- Discriminative Classifiers → *"boundary"*
 - \circ Assume some functional form for P(Y|X)
 - Estimate parameters of P(Y|X) directly from training data.
 - Learn the mappings directly from the points to the classes

- Generative models
 - Naive Bayes
 - HMM
- Discriminative models
 - Logistic Regression
 - Neural Network / Perceptron
 - SVM

Q: With the aim of classification only, which type of models may less expensive?

Credit:

https://stats.stackexchange.com/questions/ 12421/generative-vs-discriminative





• Compare the learning and prediction procedure on Naive Bayes and Logistic Regression in spam classification example



MLE vs MAP



• From Bayes rule:

Likelihood

 $\mathbf{P}(H \mid e) = \frac{\mathbf{P}(e \mid H) \mathbf{P}(H)}{\mathbf{P}(H)}$

How probable is the evidence given that our hypothesis is true?

Prior

How probable was our hypothesis before observing the evidence?

Posterior

How probable is our hypothesis given the observed evidence? (Not directly computable)

Marginal

How probable is the new evidence under all possible hypotheses? $P(e) = \sum P(e \mid H_i) P(H_i)$

• Comparing MAP and MLE: <u>[Link]</u>

$$egin{aligned} heta_{MLE} &= rg\max_{ heta}\log P(X| heta) \ &= rg\max_{ heta}\log\prod_i P(x_i| heta) \ &= rg\max_{ heta}\sum_i\log P(x_i| heta) \end{aligned}$$

$$egin{aligned} & heta_{MAP} = rg\max_{ heta} P(X| heta) P(heta) \ &= rg\max_{ heta} \log P(X| heta) + \log P(heta) \ &= rg\max_{ heta} \log \prod_i P(x_i| heta) + \log P(heta) \ &= rg\max_{ heta} \sum_i \log P(x_i| heta) + \log P(heta) \end{aligned}$$







- Clustering: Input / Output / Goal of clustering analysis
 Large amount of unlabeled data in real life
- Supervised learning v.s. unsupervised learning
- Unsupervised learning cases: Clustering and dimension reduction
- Clustering algorithm examples in this course:
 - K-means
 - Gaussian Mixture models
- Dimension reduction algorithm examples in this course:
 - PCA







Classical Machine Learning Task Driven Data Driven [Insupervised Learning Supervised Learning (Pre Categorized Data) ([Inlabelled Data) Predications & Predictive Models Pattern/ Structure Recognition Classification Regression Clustering Association (Divide the (Divide the (Divide by (Identify SOCKS by Color 1 Ties by Length) Similarity | Sequences | Eg. Identity Eg. Market Eg. Targeted Eg. Customer Fraud Detection Forecasting Marketing Recommendation



Applications of ML Categories



Supervised learning

Personalized marketing

Recommendation engines

Insurance / credit underwriting decisions

Fraud detection

Spam filtering

Demand sensing

Predictive maintenance

Sales performance prediction

People analytics

Unsupervised learning

Customer grouping or clustering, e.g. discovering groups of similar visitors to a website or discovering that a group of patients respond to the same treatment

Anomaly detection or finding outliers in the data for better fraud detection or security incident identification

Product affinity/association rule engine, e.g. discovering which two products sell best together

Semi-supervised

Used in applications where labeled data is scarce/ expensive

Speech analytics

Image classification

Web content classification

Medical predictions

Protein sequence classification

Other "learning": Self-supervised learning, reinforcement learning







• Demo 1:

http://stanford.edu/class/ee103/visualizations/kmeans/kmeans.html

• Demo 2:

https://www.naftaliharris.com/blog/visualizing-k-means-clustering/







Until no change

Engineer Change.







• Distortion measure

$$J(\{r_{nk}\}, \{\boldsymbol{\mu}_k\}) = \sum_{n=1}^{N} \sum_{k=1}^{K} r_{nk} \|\boldsymbol{x}_n - \boldsymbol{\mu}_k\|_2^2$$

- Key idea of K-means algorithms:
 - Step 1: Partition into k non-empty subsets (select K points as initial centroids)
 - Step 2: Iteration: Update mean point and assign object to cluster again
 - Step 3: Stop when converge
- Partition-based clustering methods
- Can be considered as a special case of Gaussian Mixture Model (GMM)







- Q1: Will K-means converge?
- Q2: Will different initialization of K-means generate different clustering results?



K-means



- Q1: Will K-means converge?
- A1: Yes. $J = \sum_{j=1}^{k} \sum_{C(i)=j} d(x_i, c_j)^2$
- Q2: Will different initialization of K-means generate different clustering results?
- A2: Yes. Initialization matters!







- Efficiency: **O(tkn)** normally k,t are much smaller than n → efficient
- Can terminate at a local optimum
- Need to specify **k** (or take time to find best **k**)
- Sensitive to noisy data and outliers \rightarrow K-medoids
- Different sizes and variances
- Not suitable to discover clusters with non-convex shapes
- Many variants of K-means:
 - K-means++, Genetics K-means, etc.







*Hierarchical Clustering



- Method
 - Divisive (Top-down)
 - Agglomerative (Bottom-up)

• Distance metrics

- Single linkage
- Complete linkage
- Average linkage
- Centroid
- Medoid





*Hierarchical Clustering



• Single Linkage

• Complete Linkage

• Average Linkage



 $L(r,s) = \min(D(x_{ri}, x_{sj}))$



 $L(r,s) = \max(D(x_{ri}, x_{sj}))$









Probabilistic interpretation of clustering?

We can impose a probabilistic interpretation of our intuition that points

stay close to their cluster centers

How can we model p(x) to reflect this?







Intuition

- We can model each region with a distinct distribution
- Common to use Gaussians, i.e.,
- Gaussian mixture models (GMMs) or mixture of Gaussians (MoGs).
- We don't know cluster assignments (label) or parameters of Gaussians or mixture components







Gaussian mixture models: formal definition

A Gaussian mixture model has the following density function for $oldsymbol{x}$

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

- K: the number of Gaussians they are called (mixture) components
- μ_k and $\mathbf{\Sigma}_k$: mean and covariance matrix of the k-th component
- ω_k: mixture weights they represent how much each component contributes to the final distribution. It satisfies two properties:

$$\forall \; k, \; \omega_k > 0, \quad \text{and} \quad \sum_k \omega_k = 1$$

The properties ensure p(x) is a properly normalized probability density function.





GMMs: example



The conditional distribution between $m{x}$ and $m{z}$ (representing color) are

$$p(\boldsymbol{x}|z = red) = N(\boldsymbol{x}|\boldsymbol{\mu}_1, \boldsymbol{\Sigma}_1)$$

 $p(\boldsymbol{x}|z = blue) = N(\boldsymbol{x}|\boldsymbol{\mu}_2, \boldsymbol{\Sigma}_2)$
 $p(\boldsymbol{x}|z = green) = N(\boldsymbol{x}|\boldsymbol{\mu}_3, \boldsymbol{\Sigma}_3)$



$$p(\boldsymbol{x}) = p(red)N(\boldsymbol{x}|\boldsymbol{\mu}_1, \boldsymbol{\Sigma}_1) + p(blue)N(\boldsymbol{x}|\boldsymbol{\mu}_2, \boldsymbol{\Sigma}_2)
onumber \ + p(green)N(\boldsymbol{x}|\boldsymbol{\mu}_3, \boldsymbol{\Sigma}_3)$$





Parameter estimation for GMMs: Incomplete data GMM Parameters

$$oldsymbol{ heta} = \{\omega_k,oldsymbol{\mu}_k,oldsymbol{\Sigma}_k\}_{k=1}^K$$

Incomplete Data

Our data contains observed and unobserved random variables, and hence is incomplete

- Observed: $\mathcal{D} = \{ \boldsymbol{x}_n \}$
- Unobserved (hidden): $\{\boldsymbol{z}_n\}$

Goal Obtain the maximum likelihood estimate of θ :

$$egin{aligned} \widehat{oldsymbol{ heta}} &= rg \max \ell(oldsymbol{ heta}) = rg \max \log p(oldsymbol{x}_n |oldsymbol{ heta}) \ &= rg \max \sum_n \log \sum_{oldsymbol{z}_n} p(oldsymbol{x}_n, oldsymbol{z}_n |oldsymbol{ heta}) \end{aligned}$$

The objective function $\ell(\theta)$ is called the *incomplete* log-likelihood.

Typical EM iterations

- Initialize θ with some values (random or otherwise)
- 2. Repeat
 - **E-Step:** Compute γnk using the current θ
 - b. **M-Step:** Update θ using the γnk we just computed
- 3. Until Convergence

UCLA EM Algorithms: Coin example (only M-step)



a Maximum likelihood



5 sets, 10 tosses per set

	Coin B	Coin A
	5 H, 5 T	
$\hat{\theta}_{A} = \frac{24}{24+6} = 0.80$		9 H, 1 T
ô 9 o 15		8 H, 2 T
$\theta_{B}^{=} = \frac{1}{9 + 11} = 0.45$	4 H, 6 T	
		7 H, 3 T
	9 H, 11 T	24 H, 6 T



EM Algorithms: Coin example (EM)







GMM: E-Step



E-step: Soft cluster assignments

We define γ_{nk} as $p(z_n=k|oldsymbol{x}_n,oldsymbol{ heta})$

- This is the posterior distribution of z_n given $oldsymbol{x}_n$ and $oldsymbol{ heta}$
- Recall that in complete data setting γ_{nk} was binary
- Now it's a "soft" assignment of x_n to k-th component, with x_n assigned to each component with some probability

Given $\boldsymbol{\theta} = \{\omega_k, \boldsymbol{\mu}_k, \boldsymbol{\Sigma}_k\}_{k=1}^K$, we can compute γ_{nk} using Bayes theorem:

$$egin{aligned} &\gamma_{nk} = p(z_n = k | oldsymbol{x}_n) \ &= rac{p(oldsymbol{x}_n | z_n = k) p(z_n = k)}{p(oldsymbol{x}_n)} \ &= rac{p(oldsymbol{x}_n | z_n = k) p(z_n = k)}{\sum_{k'=1}^{K} p(oldsymbol{x}_n | z_n = k') p(z_n = k')} &= rac{\mathcal{N}(oldsymbol{x}_n | oldsymbol{\mu}_k, oldsymbol{\Sigma}_k) \omega_k}{\sum_{k'=1}^{K} \mathcal{N}(oldsymbol{x}_n | oldsymbol{\mu}_{k'}, oldsymbol{\Sigma}_{k'}) \omega_{k'}} \end{aligned}$$



GMM: M-Step



M-step: Maximimize complete likelihood

Recall definition of complete likelihood from earlier:

$$\sum_{n} \log p(\boldsymbol{x}_{n}, z_{n}) = \sum_{k} \sum_{n} \gamma_{nk} \log \omega_{k} + \sum_{k} \left\{ \sum_{n} \gamma_{nk} \log \mathcal{N}(\boldsymbol{x}_{n} | \boldsymbol{\mu}_{k}, \boldsymbol{\Sigma}_{k})
ight\}$$

Previously γ_{nk} was binary, but now we define $\gamma_{nk} = p(z_n = k | \boldsymbol{x}_n)$ (E-step)

We get the same simple expression for the MLE as before!

$$egin{aligned} & \omega_k = rac{\sum_n \gamma_{nk}}{\sum_k \sum_n \gamma_{nk}}, \quad oldsymbol{\mu}_k = rac{1}{\sum_n \gamma_{nk}} \sum_n \gamma_{nk} oldsymbol{x}_n \ & oldsymbol{\Sigma}_k = rac{1}{\sum_n \gamma_{nk}} \sum_n \gamma_{nk} (oldsymbol{x}_n - oldsymbol{\mu}_k) (oldsymbol{x}_n - oldsymbol{\mu}_k)^{ ext{T}} \end{aligned}$$

Intuition: Each point now contributes some fractional component to each of the parameters, with weights determined by γ_{nk}





Consider clustering ID data with a mixture of **2** gaussian. You're given the 1-D data points **x** = **[1 2 20 40]**. Suppose the E step is the following matrix :

[0.5	0.5	
	0.2	0.8	
	0	1	
	1	0	1

→ What's the mixing weights after M-step?

 \rightarrow What's the new values of means after M-step?



GMM: Cheatsheet



Expectation (E) Step: Calculate $\forall i, k$ $\hat{\gamma}_{ik} = \frac{\hat{\phi}_k \mathcal{N}(x_i \mid \hat{\mu}_k, \hat{\sigma}_k)}{\sum_{j=1}^K \hat{\phi}_j \mathcal{N}(x_i \mid \hat{\mu}_j, \hat{\sigma}_j)},$ where $\hat{\gamma}_{ik}$ is the probability that x_i is generated by component C_k . Thus, $\hat{\gamma}_{ik} = p(C_k | x_i, \hat{\phi}, \hat{\mu}, \hat{\sigma}).$ Maximization (M) Step:

Using the $\hat{\gamma}_{ik}$ calculated in the expectation step, calculate the following in that order orall k :

•
$$\hat{\phi}_k = \sum_{i=1}^{N} \frac{\hat{\gamma}_{ik}}{N}$$

• $\hat{\mu}_k = \frac{\sum_{i=1}^{N} \hat{\gamma}_{ik} x_i}{\sum_{i=1}^{N} \hat{\gamma}_{ik}}$
• $\hat{\sigma}_k^2 = \frac{\sum_{i=1}^{N} \hat{\gamma}_{ik} (x_i - \hat{\mu}_k)^2}{\sum_{i=1}^{N} \hat{\gamma}_{ik}}$





- How does GMM relate to K-means? What are the similarities and differences?
- Will the GMM optimization process converge? (connected to K-means)





- The EM algorithm is used to find **(local) maximum likelihood parameters** of a statistical model in cases where the equations cannot be solved directly.
- Typically these models involve latent variables in addition to unknown parameters and known data observations.
- Example applied cases: K-Means, GMM
- Reading: <u>http://ai.stanford.edu/~chuongdo/papers/em_tutorial.pdf</u>





Thank you!





Reminder: You have until **Saturday, March 13 8:00 AM PST** to complete confidential evaluations for CSM146 and Dis 1C (Junheng).

Evaluation of Instruction