

CS 274A Final

Yushang Lai

TOTAL POINTS

83 / 90

QUESTION 1

1 Written Pledge 0 / 0

✓ - 0 pts Correct

should be Kd for inverting K diagonal matrices plus NKd for computing the memberships with a diagonal covariance matrix

QUESTION 2

Problem 1 30 pts

2.1 Part 1 10 / 10

✓ - 0 pts Correct

2.2 Part 2 9 / 9

✓ - 0 pts Correct

2.3 Part 3 11 / 11

✓ - 0 pts Correct

4.2 Part 2 10 / 10

✓ - 0 pts Correct

4.3 Part 3 5 / 8

✓ - 3 pts Interesting idea, but weights based on subsets of features could potentially lead to non-convergence of EM. A more straightforward way would be to just compute E and M updates based on a subset B of the data samples at each iteration.

QUESTION 3

Problem 2 30 pts

3.1 Part 1 10 / 10

✓ - 0 pts Correct

3.2 Part 2 10 / 10

✓ - 0 pts Correct

3.3 Part 3 10 / 10

✓ - 0 pts Correct

QUESTION 5

5 Blank Page 0 / 0

✓ - 0 pts Correct

QUESTION 4

Problem 3 30 pts

4.1 Part 1 8 / 12

✓ - 2 pts Part (a): Membership weights w_{ik} are not considered to be parameters. Number of parameters is $2Kd + K-1$

✓ - 2 pts Part (c): incorrect complexity for E-step. It

Take-Home Final Exam, CS 274A, Winter 2020

Time allowed: 10am Wednesday March 18th to noon Thursday March 19th

- Please write your name and sign at the bottom of the page.
- This page must also be accompanied by the following (truthful) pledge written/typed out below and signed by you: “On my honor, I have neither given nor received any unauthorized aid on this examination.”
- Write your exam on these pages, not on other pages. If you run out of page space for a problem, you can add pages to your solution and clearly indicate which problem and part at the top (e.g., “Problem 2.2”). But try if you can to fit your solution on these pages.
- Submit your final solution to Gradescope by noon Thursday March 19th. No need to submit any scratch work.
- You can consult any class-related material: lecture notes, homeworks, materials linked to from the class Web page.
- You are not allowed to use any other resources (e.g., other material on the Web, other textbooks, papers, etc).
- You are not allowed to communicate with anyone (other than the instructor) about the exam until noon on Thursday March 19th.
- If you don’t understand a question please read it carefully and at least twice. If you believe there is a typo or error in the statement of the question please contact the Professor with a **private message** on Piazza. Please do not ask for hints on how to solve the problem.

WRITTEN PLEDGE: *On my honor, I have neither given nor received any unauthorized aid on this examination.*

SIGNATURE:



PRINTED NAME:

YUSHANG LAI

Problem 1: (30 points)

1. Consider a classification problem with d -dimensional input vector \underline{x} and a class variable y taking K possible values. Given only this information define both (i) a lower bound and (ii) an upper bound on the Bayes error rate P_e^* , for this problem. No need to derive your answer, just state the bounds. The bounds should be as tight as possible given the information.

lower bound: 0

upper bound: $(K-1)/K$

2. For the problem in part 1, lets say we added a new variable to the input space (call it x_{d+1}) and we now consider the Bayes error rate of this new problem, relative to the Bayes error rate of the original problem where we only had d variables. Answer TRUE or FALSE to each of the following:

- (a) The Bayes error rate decreases

can

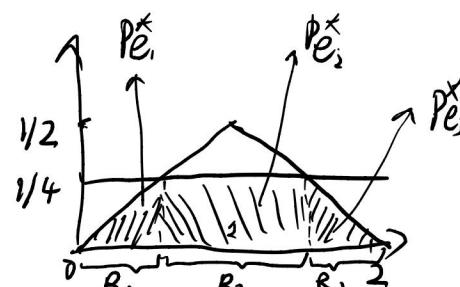
True

- (b) The Bayes error rate remains the same

True

- (c) The Bayes error rate increases

False



3. Consider the following 1-dimensional binary classification problem with $y \in \{1, 2\}$. For class 1, $p(x|y=1)$ is a uniform density $U(0, 2)$. For class 2, $p(x|y=2) = x$ for $x \in [0, 1]$, $p(x|y=2) = 2 - x$ for $x \in [1, 2]$, and $p(x|y=2) = 0$ otherwise. We also know that $p(y=1) = p(y=2)$.

$$P(X) = \sum P(X|y) = P(X|y=1) + P(X|y=2)$$

(a) Derive from first principles what the decision regions are for this classifier (to minimize 0-1 loss)

(b) Compute the Bayes error rate for this problem (and show clearly how you obtained your answer).

$$(a) 1 = P(y=1) + P(y=2) = 2P(y=1) \Rightarrow P(y=1) = P(y=2) = 1/2$$

$$1 = \int_0^2 P(X|y=1) dx = 2 \cdot P(X|y=1) \Rightarrow P(X|y=1) = 1/2$$

recall R_k = decision region in input x where $g_k(x) > g_j(x) \forall j, j \neq k$ $g_k = P(X|y=k)$

$$g_1(x) = P(y=1|x) P(x) = P(x|y=1) P(y=1) = 1/2 \cdot 1/2 = 1/4$$

$$g_2(x) = P(y=2|x) P(x) = P(x|y=2) P(y=2) = \begin{cases} x/2, & x \in (0, 1) \\ 1-x/2, & x \in (1, 2) \end{cases}$$

on $x \in (0, 1)$, let $g_1(x) = g_2(x) \Rightarrow x = 1/2$, or $x \in (1, 2)$, let $g_1(x) = g_2(x) \Rightarrow x = 3/2$

$$\text{thus } R_1 = (0, 1/2) \cup (3/2, 2) \quad R_2 = (1/2, 3/2)$$

$$(b) P(X) = g_1(x) + g_2(x) = \begin{cases} 1/4 + x/2, & x \in (0, 1) \\ 5/4 - x/2, & x \in (1, 2) \end{cases}$$

$$\begin{aligned} P_{e^*} &= \sum_k \int_{R_k} (1 - P(y=k|x)) P(x) dx = P_{e_1}^* + P_{e_2}^* + P_{e_3}^* \\ &= \int_0^{1/2} (1 - P(y=1|x)) P(x) dx + \int_{1/2}^1 (1 - P(y=2|x)) P(x) dx \\ &\quad + \int_1^{3/2} (1 - P(y=2|x)) P(x) dx + \int_{3/2}^2 (1 - P(y=1|x)) P(x) dx \end{aligned}$$

$$= \int_0^{1/2} (1/4 + x/2 - 1/4) dx + \int_{1/2}^1 (1/4 + x/2 - x/2) dx$$

$$+ \int_1^{3/2} (5/4 - x/2 - 1 + x/2) dx + \int_{3/2}^2 (5/4 - x/2 - 1/4) dx$$

$$= \frac{x^2}{4} \Big|_0^{1/2} + \frac{x}{4} \Big|_{1/2}^1 + \frac{x}{4} \Big|_1^{3/2} + x - \frac{x^2}{4} \Big|_{3/2}^2$$

$$= 1/16 + 1/4 - 1/8 + 3/8 - 1/4 + 2 - 1 - 3/2 + 9/16$$

$$= (1/16 + 1/8) \cdot 2$$

$$= 3/8$$

(Additional space for Problem 1.3, if needed)

Problem 2: (30 points)

1. Say we have a classification problem with K classes and where we have training data $D = \{(\underline{x}_i, y_i)\}, 1 \leq i \leq N$. The $y_i \in \{1, \dots, K\}$ are class labels with unknown class probabilities π_1, \dots, π_K , where $\pi_k = p(y = k)$ and $\sum_{k=1}^K \pi_k = 1$. Each \underline{x}_i represents a document. The data \underline{x}_i for the i th document, $1 \leq i \leq N$, is an M -dimensional vector of counts, $\underline{x}_i = (r_{i1}, \dots, r_{iM})$, where $r_{im} \in \{0, 1, 2, \dots\}$ is the number of times that word m occurs in document i and $m = 1, \dots, M$. The \underline{x}_i are modeled as multinomials, conditioned on the class label. Each multinomial has its own parameters $\theta_k = \{\theta_{k1}, \dots, \theta_{kM}\}$, $1 \leq k \leq K$, and $\sum_{m=1}^M \theta_{km} = 1$. Let $Dir(\alpha_1, \dots, \alpha_M)$ be a Dirichlet prior for each parameter vector θ_k . Documents i are assumed to be conditionally independent of each other and words for each document are assumed to be generated by the corresponding class multinomial.

Derive the maximum a posteriori estimate for the θ_{km} parameters. Show all steps clearly in your derivation.

$$P(\underline{x}_i | \theta_{ki}) = \prod_{m=1}^M \theta_{ki}^{r_{im}}, \quad P(y_i = k_i | \pi) = \pi_{k_i}, \quad P(\theta_{ki} | \alpha) = \prod_{m=1}^M \theta_{ki}^{d_m - 1}$$

$$\begin{aligned} P(\theta | D, \pi, \alpha) &\propto P(D | \pi, \alpha) P(\theta | \pi, \alpha) \propto \prod_{i=1}^N P(\underline{x}_i | \theta) P(y_i = k_i | \pi) P(\theta_{ki} | \alpha) \\ &= \prod_{i=1}^N \prod_{m=1}^M \theta_{k_i m}^{r_{im} + d_m - 1} \cdot \pi_{k_i} \end{aligned}$$

$$\begin{aligned} \ell(\theta) &= \log(P(\theta | D, \pi, \alpha)) = \sum_{i=1}^N \log \left(\prod_{m=1}^M \theta_{k_i m}^{r_{im} + d_m - 1} \cdot \pi_{k_i} \right) \\ &= \sum_{i=1}^N \sum_{m=1}^M \log \theta_{k_i m}^{r_{im} + d_m - 1} \cdot \pi_{k_i} \\ &= \sum_{i=1}^N \sum_{m=1}^M (r_{im} + d_m - 1) \log \theta_{k_i m} + \log \pi_{k_i} \end{aligned}$$

Note $\sum_{m=1}^M \theta_{km} = 1$

add

$$\ell(\theta) = \left(\sum_{i=1}^N \sum_{m=1}^M (r_{im} + d_m - 1) \log \theta_{k_i m} + \log \pi_{k_i} \right) - \lambda \left(\sum_{m=1}^M \theta_{k_i m} - 1 \right)$$

$$\frac{\partial \ell(\theta)}{\partial \theta_{k_i m}} = \sum_{i \in I_k} (r_{im} + d_m - 1) \frac{1}{\theta_{k_i m}} - \lambda \stackrel{\text{Want}}{=} 0$$

let I_k be the subset of $\{1, \dots, N\}$ st. $k_i = k$ for $\forall i \in I_k$

thus we have $\sum_{i \in I_k} (r_{im} + d_m - 1) \cdot \frac{1}{\theta_{km}} = \lambda$

$$\theta_{km} = \frac{\sum_{i \in I_k} (r_{im} + d_m - 1)}{\lambda}$$

(Additional space for Problem 2.1, if needed)

$$\text{since } \sum_{m=1}^M \theta_{km} = 1 \Rightarrow \sum_{m=1}^M \sum_{i \in I_k} (Y_{im} + \alpha_{m-1}) = \lambda$$

thus MAP $\hat{\theta}_{km} = \frac{\sum_{i \in I_k} (Y_{im} + \alpha_{m-1})}{\sum_{m=1}^M \sum_{i \in I_k} (Y_{im} + \alpha_{m-1})}$

where I_k is a subset
of $\{1, 2, \dots, N\}$
s.t. $k_i = k$ for $\forall i \in I_k$

2. Consider the problem in 2.1, where instead of learning the parameters we assume they are known, i.e., we know θ_k and π_k for each class. Consider a new document \underline{x} with counts (r_1, \dots, r_M) . Define an optimal discriminant function $g_k(\underline{x})$ to classify \underline{x} (assuming 0-1 loss), $k = 1, \dots, K$. Clearly define the discriminant function as a function of the θ_k 's, the π_k 's, and \underline{x} . Can the optimal discriminant function be defined as a linear discriminant? (answer YES or NO). If YES, show how it can be written in this form.

Since θ_k, π_k are known

$$\begin{aligned} \text{define } g_k(\underline{x}) &= \log [P(\underline{x}|y=k) P(y=k)] \\ &= \log \left(\prod_{m=1}^M \theta_{km}^{r_m} \cdot (\pi_k) \right) \\ &= \sum_{m=1}^M r_m \log \theta_{km} + \log \pi_k \\ &= (\log \underline{\theta}_k)^t \underline{x} + \log \pi_k \end{aligned}$$

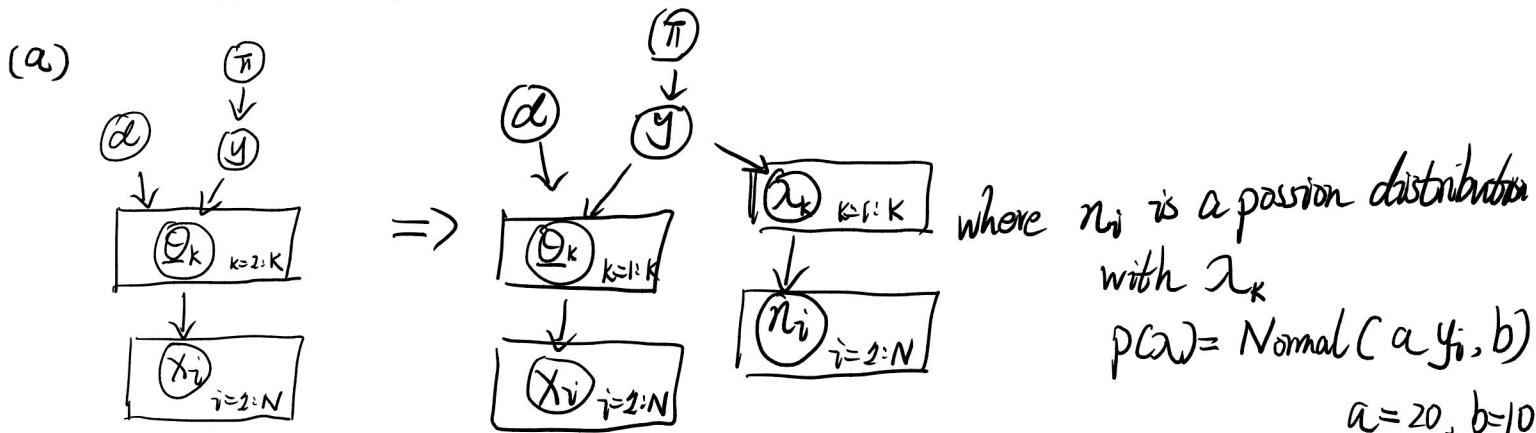
Yes, it can be defined as linear as

$$g_k(\underline{x}) = d_k \underline{x} + \beta_k \quad \text{where } d_k = (\log \underline{\theta}_k)^t$$

$$\beta_k = \log \pi_k$$

3. Say we know that the length of a document, the total number of words, $n_i = \sum_{m=1}^M r_{im}$, might also be related to the class label for a document and we want to extend the multinomial model for documents (from problem 2.1, 2.2) to a model where document length can be included when predicting the class label for a new document x .

- (a) Describe one approach (multiple approaches are possible) that shows how the model from parts 2.1/2.2 could be extended to define a generative model that includes document lengths n_i .
- (b) Show clearly how maximum likelihood could be used to learn any parameters needed for modeling document length in your extended model (and clearly state any assumptions in your derivation). (Ignore the Dirichlet priors in this part of the problem).
- (c) Write the optimal discriminant for your extended model in an additive form.



(b)

$$\begin{aligned}
 p(\theta, \alpha | D) &\propto p(D | \theta, \alpha) p(\theta, \alpha) \\
 &= \prod_{i=1}^N p(x_i | \theta_{k_i}) p(y_{i=k_i} | \pi) p(n_i | \alpha_{k_i}) \cdot p(\theta) p(\alpha)
 \end{aligned}$$

$$= \prod_{i=1}^N \prod_{m=1}^M \theta_{k_i m}^{r_{im}} \cdot \pi_{k_i} \frac{\alpha_{k_i}^{n_i} e^{-\alpha_{k_i}}}{n_i!} \cdot \text{parts not depend on } \theta_{km}, \alpha_k$$

$$\ell(\theta, \alpha) = \log(p(\theta, \alpha | D))$$

$$= \left[\sum_{i=1}^N \sum_{m=1}^M \left(r_{im} \log \theta_{k_i m} + \log \pi_{k_i} + n_i \log \alpha_{k_i} - \alpha_{k_i} - \log(n_i!) \right) \right] - M \left(\sum_{k=1}^M \alpha_k^{-1} \right)$$

$$\frac{\partial \ell(\theta, \alpha)}{\partial \theta_{km}} = \left(\sum_{i \in I_k} r_{im} \frac{1}{\theta_{km}} \right) - \mu \stackrel{\text{want}}{=} 0$$

$$\theta_{km} = \frac{\sum_{i \in I_k} r_{im}}{\mu}$$

$$\text{since } 1 = \sum_{m=1}^M \theta_{km} = \frac{\sum_{m=1}^M \sum_{i \in I_k} r_{im}}{m} \text{ thus } \mu = \sum_{m=1}^M \sum_{i \in I_k} r_{im}$$

Final Exam: CS 274A, Probabilistic Learning: Winter 2020

where I_k be the subset of $\{1, \dots, N\}$ st. $k_i = k$ for $\forall i \in I_k$

(Additional space for Problem 2, 3, if needed)

MAP thus $\hat{\theta}_{km} = \frac{\sum_{i \in I_k} r_{im}}{\sum_{m=1}^M \sum_{i \in I_k} r_{im}}$

$$\frac{\partial \ell(\theta, \lambda)}{\partial \lambda_k} = \sum_{i \in I_k} \sum_{m=1}^M \left(\frac{n_i}{\lambda_k} - 1 \right) \stackrel{\text{want}}{=} 0$$

$$\text{thus } \frac{\sum_{i \in I_k} n_i}{\lambda_k} - |I_k| = 0$$

MAP $\hat{\lambda}_k = \frac{\sum_{i \in I_k} n_i}{|I_k|}$ where I_k is the subset of $\{1, \dots, N\}$
st $k_i = k$ for $\forall i \in I_k$

$|I_k|$ is the number of elms in I_k

$$\begin{aligned} (C) \quad g_k(x, n) &= \log p(x, n | y=k) p(y=k) \\ &= \log p(x | y=k) p(n | y=k) p(y=k) \\ &= (\log \underline{\theta}_k)^t \cdot x + (\log \hat{\lambda}_k)^t \cdot n - \log(n!) + \log \pi_k \end{aligned}$$

$$\underline{\theta}_k = (\hat{\theta}_{k1}, \hat{\theta}_{k2}, \dots, \hat{\theta}_{km})$$

$$\hat{\lambda}_k = \frac{\sum_{i \in I_k} n_i}{|I_k|}$$

where $\underline{\lambda}_k$ is a vector of
 n_i when $k_i = k$

optimal discriminant function

$$\hat{g}_{x,n} = (\log \underline{\theta}_k)^t \cdot x + (\log \hat{\lambda}_k) \cdot n - \log(n!) + \log \pi_k$$

Problem 3: (30 points)

Consider a d -dimensional variable \underline{x} where we have conditionally independent observations $D = \{\underline{x}_1, \dots, \underline{x}_N\}$. We wish to use a mixture model and EM to cluster the data into K Gaussian components, each with mean μ_k and covariance matrix Σ_k , $1 \leq k \leq K$. Let w_{ik} be the membership weights as discussed in lectures, $1 \leq i \leq N, 1 \leq k \leq K$.

1. Say we are told that each Gaussian component should be modeled using a naive Bayes model.

- (a) How many total parameters are in this mixture model?
- (b) Define the M-step for updating the covariance matrix terms for this specific mixture model.
- (c) What is the time complexity of a single EM iteration for this model?

(a) w_{ik} N by K , $\underline{\mu}_k$ d by 1, Σ_k d by 1 since NB thus diagonal d_k k by 1
total parameters, $Nk + kd + kd + k = Nk + 2kd + k$

(b) for $i=1:K$

$$\textcircled{1} \quad N_k = \sum_{i=1}^N w_{ik} \quad \text{effective data pts assigned to compn}$$

$$\textcircled{2} \quad \hat{\alpha}_k = N_k/N$$

$$\textcircled{3} \quad \hat{\mu}_k = \frac{1}{N_k} \sum_{i=1}^{N_k} \underline{x}_{ik}$$

$$\textcircled{4} \quad \Sigma_k = \frac{1}{N_k} \sum_{i=1}^{N_k} \underline{x}_{ik} \begin{pmatrix} (\underline{x}_{i1} - \hat{\mu}_k)^2 & 0 \\ 0 & (\underline{x}_{id} - \hat{\mu}_{kd})^2 \end{pmatrix}$$

$$\text{where } \underline{x}_i = \begin{pmatrix} x_{i1} \\ \vdots \\ x_{id} \end{pmatrix}, \quad \underline{\mu}_k = \begin{pmatrix} \mu_{k1} \\ \vdots \\ \mu_{kd} \end{pmatrix}$$

(c) E step

$$\text{for } i=1:N; \text{ for } k=1:K \quad w_{ik} = E[Z_{ik} | \underline{x}_i, \theta] = P(Z_{ik}=1 | \underline{x}_i, \theta)$$

$$= \frac{P(\underline{x}_i | Z_{ik}=1, \theta_k) P(Z_{ik} | \theta)}{\sum_{k=1}^K (\text{above part})} \quad \textcircled{1}$$

$$\frac{P(\underline{x}_i | \theta)}{\sum_{k=1}^K P(\underline{x}_i | Z_{ik}=1, \theta_k)} \quad \textcircled{2}$$

$$\mathbb{E} \textcircled{1} \quad N \times K \times d \times d \quad \textcircled{2} \quad K \times d \quad \text{inverse matrix}$$

$$M \textcircled{1} \quad KN \quad \textcircled{2} \quad K \quad \textcircled{3} \quad K \quad \textcircled{4} \quad KNd$$

total for one iteration $O(NKd^2)$

In parts 3.2 and 3.3 of this problem below assume that we are now using unconstrained covariance matrices Σ_k (i.e., no naive Bayes assumption).

2. Say that m of the N datapoints (where $m < N$) are labeled, i.e., for these m datapoints (only) the data consists of pairs (\underline{x}_i, y_i) , with $y_i \in \{1, \dots, K\}$ indicating the true label. Describe below how you would do maximum likelihood estimation of the μ_k 's and Σ_k 's for such a model. A high-level sketch of the key idea for your proposed algorithm is fine.

I will use EM Algorithm ^{Gaussian Mixture} to do maximum likelihood estimation of μ_k 's and Σ_k 's.

let $\Theta_k = \{\mu_k, \Sigma_k\}$ $\Theta = \{\Theta_1, \dots, \Theta_K\}$

let $D_x = \{\underline{x}_i\}$ $D_z = \{\underline{z}_i\}$ where $\underline{z}_i = (z_{i1}, \dots, z_{ik})$ $z_{ik}=1$ if x_i true
 $z_{ik}=0$ otherwise

since we have m data with labels, we fix those lines W_{ik} in there true columns and not update its membership weights in each iteration's E step.

then based on Jensen's Inequality

we maximize the likelihood of μ_k 's and Σ_k 's by ascent of the lower bound of $\ell(\Theta) = \sum_{i=1}^n \log P(x_i; \Theta)$

Since time complexity linear depend
on N , hard to speed up by variable N .
¹²

3. Consider now a standard Gaussian mixture problem with N unlabeled data points. In class we discussed stochastic gradient for speeding up gradient-based training. Clearly describe some method (e.g., step-by-step, pseudocode, etc) that could use an approach similar to stochastic gradient to speed up the EM algorithm for Gaussian mixtures. Your method need not have a guarantee that the likelihood is non-decreasing at each iteration. If the mini-batch size is B , how much speed-up per iteration would your proposed method gain?

One iteration

E step

random order feature d , split into subset of size B , with the last one w.l.o.g. say G subsets, $\{N_1, \dots, N_G\}$
may less than B

$$w_{ik} = 1 \text{ for } \forall i \in N, k \in K$$

for $i=1:N$.

for $g=1:G$

for $k=1, \dots, K$

X_{ik}^g mean that
we only use the
features in subset
 $N_g \subset \{1, \dots, d\}$

$$w_{ik}^* = \frac{P(X_{ik}^g | z_{ik}=1, \theta_k) P(z_{ik}=1 | \theta_k)}{\sum_k \text{terms on top}}$$

the M step

for $g=1:G$

for $i=1:k$

$$N_{gk} = \sum_{i=1}^{N_g} w_{ik}$$

$$\hat{\alpha}_k = N_{gk}/N_g$$

$$\hat{\mu}_k = \frac{1}{N_{gk}} \sum_{i=1}^{N_g} X_{ik}$$

$$\Sigma_k = \frac{1}{N_{gk}} \sum_{i=1}^{N_g} w_{ik} (X_{ik} - \hat{\mu}_k) (X_{ik} - \hat{\mu}_k)^t$$

speed up per iteration: previous E step $O(Kd^3 + Nd^2)$ M step $O(Nkd^2)$

$$\begin{aligned} \text{new E step} &= O(KB^3 \cdot \frac{d}{B} + NKB^2 \cdot \frac{d}{B}) \\ &= O(KB^2 d + NKBD) \\ &= O(NKBd) \end{aligned}$$

based on E step performance may speed up $(\frac{d}{B})$ times

