

# CS228 Homework 5

Instructor: Stefano Ermon – [ermon@stanford.edu](mailto:ermon@stanford.edu)

Available: 03/3/2017; Due: 03/18/2017

---

1. [25 points] (**Bayesian inference**) Let  $X \in \{x^1, \dots, x^K\}$  be a multinomial variable and let  $\theta$  be a parametrization of the distribution of  $X$ , i.e.  $P(X = x^k | \theta) = \theta_k$ . Let  $\mathcal{D} = \{x[1], \dots, x[M]\}$  be a dataset consisting of  $M$  realizations of  $X$ . We would like to infer something interesting about  $\theta$  based on  $\mathcal{D}$ .

Our strategy so far has been to infer a true set of parameters  $\theta^*$  from which the data was generated; we found  $\theta^*$  using the principle of maximum likelihood; this was an example of the so-called *frequentist* approach to statistics. In *Bayesian* statistics, we instead construct a *posterior* distribution  $P(\theta | \mathcal{D})$  that can be used to describe our uncertainty over the parameters given the evidence we observed. Recall that we will construct  $P(\theta | \mathcal{D})$  using Bayes theorem:

$$P(\theta | \mathcal{D}) = \frac{P(\mathcal{D} | \theta) P(\theta)}{P(\mathcal{D})}.$$

This approach has the advantage of better modeling the full uncertainty over the parameters. However, the probabilities no longer correspond to limiting frequencies within a data-generating process described by  $P(\mathcal{D} | \theta)$ . Instead, they can only be interpreted as “beliefs”. Most crucially, these probabilities can be arguably called “subjective” because they depend on a set of arbitrary initial beliefs specified by  $P(\theta)$ . This may raise objections, since we might want our inferences about the world to be independent of any subjective choices by the statistician. It is also not always clear how to specify  $P(\theta)$  and how to translate our prior beliefs into probabilities. Arguments like these are part of the great frequentist vs. Bayesian debate in statistics. Here, we will see a concrete example of how the Bayesian approach can be useful.

- (a) [8 points] Let’s use the posterior to make predictions on new samples. Suppose that the likelihood  $P(\mathcal{D} | \theta) = \prod_{j=1}^M P(x[j] | \theta)$  is a product of categorical distributions (i.e. we assume that the observations are independent of each other given  $\theta$ ) and let’s choose a Dirichlet prior over  $\theta$ , i.e.  $\theta \sim \text{Dirichlet}(\alpha_1, \dots, \alpha_K)$ . Recall that  $\theta \sim \text{Dirichlet}(\alpha_1, \dots, \alpha_K)$  if  $P(\theta) \propto \prod_k \theta_k^{\alpha_k - 1}$ . Show that the Bayesian predictive probability using a Dirichlet prior is

$$P(X[M+1] = x^i | \mathcal{D}) = \frac{M[i] + \alpha_i}{M + \alpha},$$

where  $M[i]$  is the number of times  $x[m] = x^i$  appears in the dataset and  $\alpha = \sum_i \alpha_i$ .  $X[M+1]$  is assumed conditionally independent of  $\mathcal{D}$  given  $\theta$ .

Note that this probability has a very neat interpretation:  $M[i]/M$  by itself is simply the frequency of class  $i$  in our dataset. By adding a Dirichlet prior, we effectively augment our dataset with  $\alpha[i]$  “virtual” data points of class  $i$ : the predictive probability is the same we would’ve had it there were  $\alpha[i]$  extra points of class  $i$  in the actual dataset!

Hint: Recall from Lecture 15 that the posterior  $P(\theta | x[1], \dots, x[M])$  is given by  $\text{Dirichlet}(\alpha'_1, \dots, \alpha'_K)$ , where

$$\alpha'_k = \alpha_k + \sum_{j=1}^M 1\{x[j] = x^k\}.$$

- (b) **[8 points]** Now we want to compute the Bayesian predictive probability over two samples. Show how to compute

$$P(X[M+1] = x^i, X[M+2] = x^j | \mathcal{D})$$

- (c) **[9 points]** Suppose we decide to use the approximation

$$P(X[M+1] = x^i, X[M+2] = x^j | \mathcal{D}) \approx P(X[M+1] = x^i | \mathcal{D}) \cdot P(X[M+2] = x^j | \mathcal{D})$$

That is, we ignore the dependencies between  $X[M+1]$  and  $X[M+2]$ . Analyze the error in this approximation (the ratio between the approximation and the correct probability). What is the quality of this approximation for small  $M$ ? What is the asymptotic behavior of the approximation when  $M \rightarrow \infty$ .

In general, Bayesian inference may not always be tractable, and often requires approximations such as the one in this question. Finding such approximations is the topic of a large subfield of machine learning which studies the problem of “approximate inference”.

[illegible]

$$(b) \quad P(X_{1M+1} = x^1, X_{1M+2} = x^2 | D) \\ = P(X_{1M+1} = x^1 | D) \cdot P(X_{1M+2} = x^2 | D, X_{1M+1} = x^1) \\ = \frac{x^1 + M E^1}{\alpha + M} \cdot P(X_{1M+2} = x^2 | D, X_{1M+1} = x^1) \\ = \begin{cases} \frac{x^1 + M E^1}{\alpha + M} \cdot \frac{x^2 + M E^2}{\alpha + M + 1} & \text{if } i=j \\ \frac{x^1 + M E^1}{\alpha + M} \cdot \frac{x^2 + M E^2 + 1}{\alpha + M + 1} & \text{if } i \neq j \end{cases}$$

$$\begin{aligned}
 (C). \quad \text{for } i=j \\
 P(X_{1:M} = i, X_{1:M+1:j} | D) &= \frac{\alpha^i + M \cdot i}{\alpha + M} \cdot \frac{\alpha^j + M \cdot j + 1}{\alpha + M + 1} \\
 P(X_{1:M} = i | D) \cdot P(X_{1:M+1:j} | D) &= \frac{\alpha^i + M \cdot i}{\alpha + M} \cdot \frac{\alpha^j + M \cdot j + 1}{\alpha + M} \\
 \therefore P(X_{1:M} = i | D) \cdot P(X_{1:M+1:j} | D) &\leq P(X_{1:M} = i, X_{1:M+1:j} | D) \\
 \frac{P(X_{1:M} = i | D) \cdot P(X_{1:M+1:j} | D)}{P(X_{1:M} = i, X_{1:M+1:j} | D)} &= \frac{\alpha^i + M \cdot i}{\alpha + M} \cdot \frac{\alpha + M + 1}{\alpha^j + M \cdot j + 1} \\
 &= \frac{\alpha^i + pM}{\alpha + M} \cdot \frac{\alpha + M + 1}{\alpha^j + pM + 1} \\
 &= \frac{\alpha^i / M + p}{\alpha / M + 1} \cdot \frac{\alpha / M + 1 + 1/M}{\alpha / M + p + 1/M} \\
 (p \text{ is the true probability true}
 \end{aligned}$$

when  $M$  is small, the ratio depends highly on the choice of  $\alpha$   
if  $\alpha_i = \rho \alpha$ , then the ratio is one

when  $N$  is large,  $\lim_{N \rightarrow \infty} \frac{P(X_1 = i | D) \cdot P(X_2 = j | D)}{P(X_1 = i, X_2 = j | D)} = 1$   
 the more samples we have, the "more independent" samples are given history samples

## 2. [75 points] Programming Assignment <sup>1</sup>

This homework explores parameter learning in latent variable graphical models in the context of a hypothetical problem involving voter registration.

Suppose you are working as a volunteer on behalf of one of the two major presidential candidates in an evenly divided city in a swing state – let's say, Cleveland, Ohio. Your goal is to register as many voters as possible who are likely to support your candidate. However, your party has a limited number of volunteers, and needs to be wise about how it spends scarce resources in canvassing for voters. Fortunately, a local university recently conducted an extensive survey in which the city was partitioned into fifty precincts, and twenty citizens per precinct were surveyed on their political views. A table has been made publicly available that lists for each respondent:

- The precinct  $i = 1, \dots, N$ . In our city,  $N = 50$ .
- The index  $j = 1, \dots, M$  of the respondent within its precinct; in our case,  $M = 20$ .
- Two summary statistics  $X_{ij} = (X_{ij}^{(1)}, X_{ij}^{(2)})^T$  indicating respectively the overall social conservatism/liberalism and economic conservatism/liberalism of the  $j$ -th respondent in district  $i$ .
- In five of the fifty precincts, we have explicit respondent party preferences  $Z_{ij} \in \{0, 1\}$ .

Your goal will be to use the summary statistics to obtain probabilistic information about party preference for the respondents that have not been explicitly surveyed (missing data).

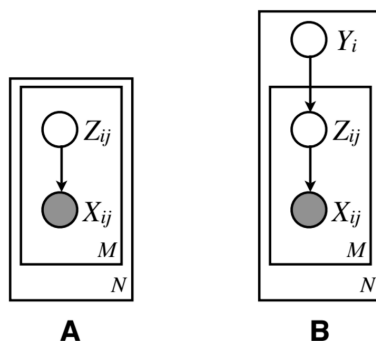


Figure 1: Graphical representation of the models used in this problem.

We will use two models shown in Figure 1 and described below.

### (A) Gaussian mixture model

We model the  $X_{ij}$  as a mixture of two Gaussians  $P(X) = (1 - \pi)\mathcal{N}(X|\mu_0, \Sigma_0) + \pi\mathcal{N}(X|\mu_1, \Sigma_1)$ . This model tries to capture the following generative story: first, each person independently samples a party preference  $Z_{ij}$  from a Bernoulli distribution with parameter  $\pi$ ; then, their sample summary statistics are sampled as  $X_{ij} | z_{ij} \sim \mathcal{N}(\mu_{z_{ij}}, \Sigma_{z_{ij}})$ , where  $\mu_l = (\mu_l^{(1)}, \mu_l^{(2)})^T$  is the class-conditional mean for party  $l$ , and  $\Sigma_l$  is the class-conditional variance/covariance matrix for party  $l$ . Note that precinct membership does not factor into this model, despite its use in indexing the variables.

With regards to the above model, answer the following questions.

- [10 points]** Estimate the parameters  $\pi$ ,  $\mu_0$ ,  $\mu_1$ ,  $\Sigma_0$ , and  $\Sigma_1$  by maximum likelihood using just data from the five precincts in which respondents have reported their party preferences. The data is provided in 'survey-labeled.dat'. Provide an explicit estimator formula for each parameter.
- [15 points]** In this part, we will use the unlabelled dataset 'survey-unlabeled.dat' to learn the model parameters. Implement a Gaussian-Mixtures EM algorithm for model A and use it to estimate  $\theta = \{\pi, \mu_0, \mu_1, \sigma_0, \sigma_1\}$ .

<sup>1</sup>Assignment adapted from Cornell's BTRY 6790, instructed by Adam Siepel

- The  $E$ -step computes  $P(z_{ij}|x_{ij}, \theta)$ , which is the expected “counts” for every respondent given a fixed value of the model parameters.
- The  $M$ -step re-estimates the model parameters,  $\theta$  based on the expected values calculated in the  $E$ -step.

The algorithm should compute and output the log likelihood on every iteration, and should terminate when this quantity increases by less than a value of 0.01 between iterations. Run the algorithm with three different initializations: one equal to your estimates from part 2(A)i and two other (poorer) initializations of your choice. Plot the log likelihood as a function of algorithm iteration for all three cases. Comment on any differences in the local maxima that are found. Report your parameter estimates.

#### (B) Geography-aware mixture model

The second model attempts to capture the fact that respondents in the city tend to be geographically separated by party preference, with some precincts showing strong preferences for one party and others showing strong preferences for the other party. In this model, an additional variable  $Y_i \in \{0, 1\}$  is introduced for each precinct  $i$ , representing that precinct’s preferred party. Our new model has the following generative story: first, the  $Y_i$  variables are drawn i.i.d. from a Bernoulli distribution with parameter  $\phi$ . Then, the party preferences  $Z_{ij}$  as sampled according to

$$p(z_{ij}|y_i) = \begin{cases} \lambda & \text{if } z_{ij} = y_i \\ (1 - \lambda) & \text{otherwise} \end{cases}$$

Here,  $\lambda$  is a new parameter that we introduce. Note also that the  $Z_{ij}$  variables are conditionally independent given the  $Y_i$  variables. Finally, the  $X_{ij}$  are sampled as in the previous model.

With regards to the above model, answer the following questions.

- [5 points]** We will first estimate the parameters using data from the five precincts in which respondents have reported their party preferences provided in 'survey-labeled.dat'. Specifically, estimate the parameters  $\phi$  and  $\lambda$  by the approximate method of setting each  $y_i$  to the consensus (majority) of the corresponding  $z_{ij}$  values ( $y_i = I(\sum_{j=1}^M z_{ij} \geq \frac{M}{2})$ ), then acting as as if the  $y_i$ s were observed. Explicitly write out the log likelihood function in terms of  $\phi$ ,  $\lambda$ ,  $\mu_0$ ,  $\mu_1$ ,  $\Sigma_0$ , and  $\Sigma_1$ . Derive maximum likelihood estimators for  $\phi$  and  $\lambda$  in terms of completely observed  $x$ ,  $y$ , and  $z$  variables. Note that the estimates for  $\mu_0$ ,  $\mu_1$ ,  $\Sigma_0$ , and  $\Sigma_1$  will remain unchanged from the ones estimated for part 2(A)i.
- [10 points]** Having estimated your parameters for model B by “supervised” training, use them to analyze the unlabeled data set ('survey-unlabeled.dat') by identifying precincts to be targeted by party 1. Specifically, compute  $p(y_i|x_{i,1:M})$  for each precinct  $i$ , where  $x_{i,1:M} = \{x_{ij} : 1 \leq j \leq M\}$ , and identify those precincts for which this probability exceeds 0.5. Summarize your results by presenting a table with one row per precinct, in ascending order by index, a column with the quantity  $p(y_i|x_{i,1:M})$ , and a mark indicating the precincts that exceed a threshold of 0.5. In addition, compute  $p(z_{ij}|x_{i,1:M})$  for every respondent  $(i, j)$ , and summarize the results by plotting the data points on a two-dimensional plane and coloring them blue if  $p(z_{ij} = 1|x_{i,1:M}) > 0.5$  or red otherwise. Also indicate the positions of the two means.

Hint: From the factorized joint distribution of the precinct,

$$p(y_i = 1|x_{i,1:M}) = \frac{p(y_i = 1) \sum_{z_{i,1:M}} \prod_{j=1}^M p(x_{ij}|z_{ij}) p(z_{ij}|y_i = 1)}{p(x_{i,1:M})}.$$

- [15 points]** Now, we will estimate the parameters for model B using the unlabelled dataset 'survey-unlabeled.dat'. Derive  $E$ - and  $M$ -step updates for model B. Start by writing down the complete log likelihood from part 2(B)i. Let  $\theta = \{\phi, \lambda, \mu_0, \mu_1, \Sigma_0, \Sigma_1\}$  denote the model parameters.

- In the  $E$ -step, derive an expression for  $p(y_i, z_{i,1:M}|x_{i,1:M}, \theta)$ , which is the expected value of the relevant “counts” for a fixed value of the parameters. Can this be represented compactly (in a factored form)?
- The  $M$ -step re-estimate the model parameters,  $\theta$  based on the expected values computed in the  $E$ -step. As a hint, note that the parameter update equations can be easily computed once we have  $p(y_i, z_{ij}|x_{i,1:M}, \theta)$ ,  $p(z_{ij}|x_{i,1:M}, \theta)$ , and  $p(y_{ij}|x_{i,1:M}, \theta)$ .

Hint: Your calculations in the  $M$ -step should roughly mirror the ones in part 2(B)i. In addition, the  $E$ -step will resemble that for model A.

- iv. [10 points] Based on the above updates, implement an EM algorithm for model B. Run the algorithm with three different initializations, plot the log likelihood and report the parameter estimates.
- v. [10 points] Using your best set of parameter estimates (those yielding the highest log likelihood) again compute  $p(y_i = 1|x_{i,1:M})$  for each precinct  $i$ , and identify those precincts for which this probability exceeds 0.5. Report a new version of the table from part 2(B)ii. In addition, again compute  $p(z_{ij}|x_{i,1:M})$  and prepare a new plot with red and blue data points.

### Hints:

- (a) It is important that you understand EM algorithm well in order to do this programming assignment. One good note on EM algorithm is : <https://people.csail.mit.edu/rameshvs/content/gmm-em.pdf>.
- (b) One common source of mistake in part b is how you calculate the log of product-sum. Note that:

$$\log\left(\prod_i \sum_j f(x_{ij})\right) = \sum_i \log\left(\sum_j f(x_{ij})\right)$$

This is not equivalent to:

$$\log\left(\prod_i \sum_j f(x_{ij})\right) = \sum_i \sum_j \log(f(x_{ij}))$$

- (c) Be careful on how you use numpy’s reshape and transpose function when you write code for log-likelihood.



M-step:

$$\max_{\theta} \sum_{ij} \sum_{y \in \mathcal{Y}} \sum_{z \in \mathcal{Z}} Q_{ij}(y \in \mathcal{Y}, z \in \mathcal{Z}) \log [p(y \in \mathcal{Y}) p(z \in \mathcal{Z} | y \in \mathcal{Y}) p(x \in \mathcal{X} | z \in \mathcal{Z})]$$

$$= \max_{\theta} \sum_{ij} \sum_{y \in \mathcal{Y}} \sum_{z \in \mathcal{Z}} Q_{ij}(y \in \mathcal{Y}, z \in \mathcal{Z}) [\log p(y \in \mathcal{Y}) + \log p(z \in \mathcal{Z} | y \in \mathcal{Y}) + \log p(x \in \mathcal{X} | z \in \mathcal{Z}; u, v, \Sigma, \tau)]$$

① max over  $\phi$

$$\max_{\phi} \sum_{ij} \sum_{y \in \mathcal{Y}} \sum_{z \in \mathcal{Z}} Q_{ij}(y \in \mathcal{Y}, z \in \mathcal{Z}) [\log \phi \mathbb{1}\{y \in \mathcal{Y}\} + \log(1-\phi) \mathbb{1}\{y \in \mathcal{Z}\}]$$

$$= \max_{\phi} \log \phi \underbrace{\left[ \sum_{ij} \sum_{y \in \mathcal{Y}} Q_{ij}(1, z \in \mathcal{Z}) \right]}_a + \log(1-\phi) \underbrace{\left[ \sum_{ij} \sum_{z \in \mathcal{Z}} Q_{ij}(0, z \in \mathcal{Z}) \right]}_b$$

$$\nabla_{\phi} = \frac{a}{\phi} - \frac{b}{1-\phi} \stackrel{!}{=} 0$$

$$\phi = \frac{a}{a+b}$$

② max over  $\lambda$

$$\max_{\lambda} \sum_{ij} \sum_{y \in \mathcal{Y}} \sum_{z \in \mathcal{Z}} Q_{ij}(y \in \mathcal{Y}, z \in \mathcal{Z}) [\log \lambda \mathbb{1}\{z \in \mathcal{Z} = y \in \mathcal{Y}\} + \log(1-\lambda) \mathbb{1}\{z \in \mathcal{Z} \neq y \in \mathcal{Y}\}]$$

$$= \max_{\lambda} \log \lambda \underbrace{\left[ \sum_{ij} Q_{ij}(0,0) + Q_{ij}(1,1) \right]}_{\alpha} + \log(1-\lambda) \underbrace{\left[ \sum_{ij} Q_{ij}(0,1) + Q_{ij}(1,0) \right]}_{\beta}$$

$$\lambda = \frac{\alpha}{\alpha + \beta}$$

③ max over  $u$ .

$$\max_u \sum_{ij} \sum_{y \in \mathcal{Y}} \sum_{z \in \mathcal{Z}} Q_{ij}(y \in \mathcal{Y}, z \in \mathcal{Z}) \log p(x \in \mathcal{X} | z \in \mathcal{Z}; u, v, \Sigma, \tau)$$

$$= \max_u \sum_{ij} \sum_{y \in \mathcal{Y}} Q_{ij}(y \in \mathcal{Y}, 0) \cdot \left[ -\frac{1}{2} (x \in \mathcal{X} - u) \Sigma^{-1} (x \in \mathcal{X} - u) \right]$$

$$\nabla_u = \sum_{ij} \sum_{y \in \mathcal{Y}} Q_{ij}(y \in \mathcal{Y}, 0) \Sigma^{-1} (x \in \mathcal{X} - u) \stackrel{!}{=} 0$$

$$u_0 = \frac{\sum_{ij} \sum_{y \in \mathcal{Y}} Q_{ij}(y \in \mathcal{Y}, 0) x \in \mathcal{X}}{\sum_{ij} \sum_{y \in \mathcal{Y}} Q_{ij}(y \in \mathcal{Y}, 0)} = \frac{\sum_{ij} x \in \mathcal{X} \cdot \sum_{y \in \mathcal{Y}} Q_{ij}(y \in \mathcal{Y}, 0)}{\sum_{ij} \sum_{y \in \mathcal{Y}} Q_{ij}(y \in \mathcal{Y}, 0)}$$

$$u_1 = \frac{\sum_{ij} x \in \mathcal{X} \cdot \sum_{y \in \mathcal{Y}} Q_{ij}(y \in \mathcal{Y}, 1)}{\sum_{ij} \sum_{y \in \mathcal{Y}} Q_{ij}(y \in \mathcal{Y}, 1)}$$

④ max over  $\Sigma$ .

$$\max_{\Sigma} \sum_{ij} \sum_{y \in \mathcal{Y}} Q_{ij}(y \in \mathcal{Y}, 0) \cdot \log \frac{1}{(2\pi)^{\frac{1}{2}} |\Sigma|^{1/2}} \exp \left[ -\frac{1}{2} (x \in \mathcal{X} - u) \Sigma^{-1} (x \in \mathcal{X} - u) \right]$$

$$= \max_{\Sigma} \sum_{ij} \sum_{y \in \mathcal{Y}} Q_{ij}(y \in \mathcal{Y}, 0) \left[ -\frac{1}{2} \log \det(\Sigma) - \frac{1}{2} \text{tr}(\Sigma^{-1} (x \in \mathcal{X} - u)(x \in \mathcal{X} - u)^T) \right]$$

$$\text{let } \Sigma_0^{-1} = J.$$

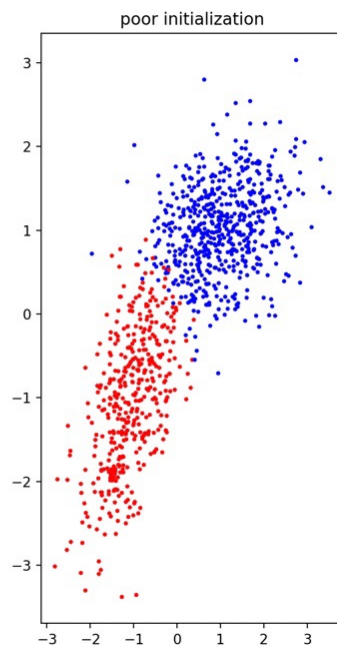
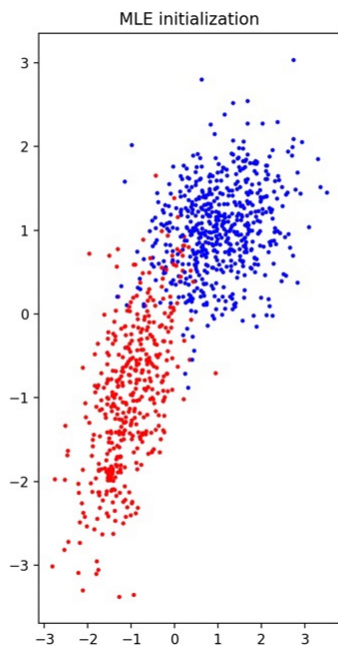
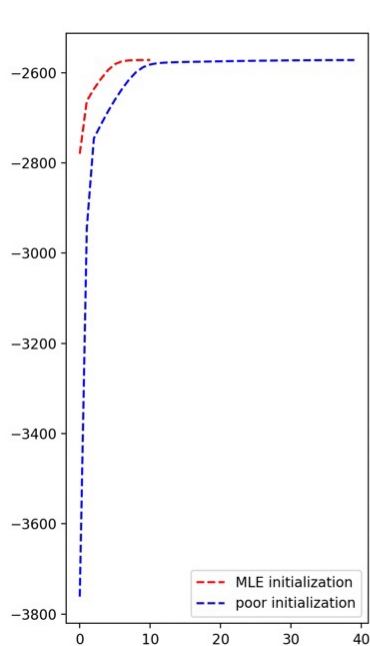
$$= \max_J \sum_{ij} \sum_{y \in \mathcal{Y}} Q_{ij}(y \in \mathcal{Y}, 0) \cdot \log \det(J_0) - \sum_{ij} \sum_{y \in \mathcal{Y}} Q_{ij}(y \in \mathcal{Y}, 0) \text{tr}(J_0 (x \in \mathcal{X} - u)(x \in \mathcal{X} - u)^T)$$

$$= \max_J \log \det(J_0) \sum_{ij} \sum_{y \in \mathcal{Y}} Q_{ij}(y \in \mathcal{Y}, 0) - \text{tr} \left( J_0, \sum_{ij} \sum_{y \in \mathcal{Y}} Q_{ij}(y \in \mathcal{Y}, 0) (x \in \mathcal{X} - u)(x \in \mathcal{X} - u)^T \right)$$

$$\nabla_{J_0} = J_0^{-1} \sum_{ij} \sum_{y \in \mathcal{Y}} Q_{ij}(y \in \mathcal{Y}, 0) - \sum_{ij} (x \in \mathcal{X} - u)(x \in \mathcal{X} - u)^T \cdot \sum_{y \in \mathcal{Y}} Q_{ij}(y \in \mathcal{Y}, 0) \stackrel{!}{=} 0$$

$$\Sigma_0 = J_0^{-1} = \frac{\sum_{ij} (x \in \mathcal{X} - u)(x \in \mathcal{X} - u)^T \cdot \sum_{y \in \mathcal{Y}} Q_{ij}(y \in \mathcal{Y}, 0)}{\sum_{ij} \sum_{y \in \mathcal{Y}} Q_{ij}(y \in \mathcal{Y}, 0)}$$

$$\Sigma_1 = \frac{\sum_{ij} (x \in \mathcal{X} - u_1)(x \in \mathcal{X} - u_1)^T \cdot \sum_{y \in \mathcal{Y}} Q_{ij}(y \in \mathcal{Y}, 1)}{\sum_{ij} \sum_{y \in \mathcal{Y}} Q_{ij}(y \in \mathcal{Y}, 1)}$$



EM depends highly on the initialization