

CS228 Homework 5

Instructor: Stefano Ermon – ermon@stanford.edu

Available: 02/27/2018; Due: 03/13/2018

please note: points may be deducted from answers that are needlessly complicated or contain extraneous information

-
1. [25 points] (**Bayesian inference**) Let $X \in \{x^1, \dots, x^K\}$ be a multinomial variable and let θ 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 θ based on \mathcal{D} .

Our strategy so far has been to infer a true set of parameters θ^* from which the data was generated; we found θ^* 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 θ) and let's choose a Dirichlet prior over θ , 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 θ .

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 i.e., the predictive probability is the same we would've had if 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

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

Additionally, you may use the formula for the mean, variance, or mode of any standard probability distribution such as the Dirichlet distribution.

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

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

Notice that the parameters of the model will be updated as we observe new samples.

- (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”.

2. [75 points] Programming Assignment ¹

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).

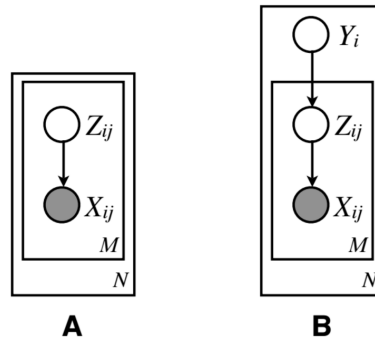


Figure 1: Graphical representation of the models used in this problem.

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 π ; 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 Σ_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 π , μ_0 , μ_1 , Σ_0 , and Σ_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.
- [20 points]** In this part, we will use the unlabelled dataset `survey-unlabeled.dat` to learn the model parameters. Implement a Gaussian-Mixture EM algorithm for model A and use it to estimate $\theta = \{\pi, \mu_0, \mu_1, \Sigma_0, \Sigma_1\}$.

¹Assignment adapted from Cornell's BTRY 6790, instructed by Adam Siepel

- The E -step computes $P(z_{ij} \mid x_{ij}, \theta)$, which is the expected “counts” for every respondent given a fixed value of the model parameters.
- The M -step re-estimates the model parameters, θ 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 ϕ . Then, the party preferences Z_{ij} are sampled according to

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

Here, λ 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`. You can set 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 if the y_i ’s were observed (so all the variables x_{ij}, z_{ij}, y_i are fully observed). Explicitly write out the log-likelihood function in terms of $\phi, \lambda, \mu_0, \mu_1, \Sigma_0$, and Σ_1 . Find maximum likelihood estimators for ϕ and λ in terms of completely observed x, y , and z variables. Write the estimators in analytic form and also give their numerical values calculated from the provided data. Note that the estimates for μ_0, μ_1, Σ_0 , and Σ_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 $p(y_i \mid x_{i,1:M})$ 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 $p(y_i \mid x_{i,1:M})$, and a mark indicating the precincts that exceed a threshold of 0.5. In addition, write down the explicit expression for computing $p(z_{ij} \mid x_{i,1:M})$ for each respondent (i, j) , and summarize the results by plotting the data points on a two-dimensional plane and coloring them blue if $p(z_{ij} = 1 \mid x_{i,1:M}) > 0.5$ or red otherwise. Also indicate the positions of the two means.

Hint: You can start your derivation from the factorized joint distribution of the precinct,

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

Note: for numerical stability you should use log probability instead of the probability and use the log-sum-exp trick, where for any real numbers $A, B > 0$ and $M = \max(A, B)$,

$$\log(e^{\log A} + e^{\log B}) = M + \log(e^{\log A - M} + e^{\log B - M})$$

iii. [25 points] Now, we will estimate the parameters for model B using the unlabelled dataset `survey-unlabeled.dat`. Find E - and M -step updates for model B. 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)?
- Given your results in the E -step, write out an expression for computing $\log p(X)$; your expression should be computable in time $O(NM)$.
- The M -step updates the model parameters θ based on the expected values computed in the E -step by maximizing the likelihood of the observed variables. As a hint, the log-likelihood and the parameter updates can be easily computed once we have $p(y_i, z_{ij} | x_{i,1:M}, \theta)$, $p(z_{ij} | x_{i,1:M}, \theta)$, $p(y_i | x_{i,1:M}, \theta)$, and $p(z_{ij} | x_{i,1:M}, y_i, \theta)$. Note that all four can be computed if we first compute $p(y_i, z_{ij} | x_{i,1:M}, \theta)$.

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.

iv. [5 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:

- 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>.
- One common source of mistake in part b is how you calculate the log of product-sum. Note that:

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

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