

Section 5: Mixture models

Mixture models

So far, we've assumed that our data are conditionally exchangeable given their covariates. In other words, for every unique set of covariates there exists a set of parameters, conditioned on which, the data with those covariates are i.i.d. We used various distributions over functions to learn a distribution over these parameters, for all covariate settings.

A common setting was when our data was normally distributed, with mean $\beta^T x_i$ and variance σ^2 . If we did not have the covariate values x_i , our data would no longer be normally distributed.

Exercise 5.1 Download the dataset `restaurants.csv`. This contains profit information for restaurants, based on seating capacity and whether they are open for dinner. Run a Bayesian regression of Profit vs SeatingCapacity and a dummy for DinnerService (you can reuse code from 2.12) (I'd suggest whitening Profit, it will make later prior specification easier). Do the residuals look normal? (e.g. plot histograms, qq plots). Now, let's just look at the raw Profit data: Does it look normal?

Let's assume we're in the situation where we don't know any of these covariate values. For now, let's ignore the continuous-valued covariate (SeatingCapacity), and try to infer the categorical covariate. Let's say we know that half our restaurants are open for dinner. We could assume that each restaurant is associated with a latent indicator variable Z_i , that assigns them to one of two groups, so that

$$Z_i \sim \text{Bernoulli}(\pi)$$

As in the regression setting, conditioned on the latent variable, we will assume that the observed profits are i.i.d. normal. Again, as in the basic regression setting, we will assume the variances of the two normals are the same, but the means are different, i.e.

$$X_i | Z_i = z \sim \text{Normal}(\mu_z, \sigma^2).$$

If we marginalize over these binary indicators, our observations are assumed to be distributed according to a mixture of two Gaussians:

$$X_i \sim 0.5N(\mu_1, \sigma_1^2) + 0.5(\mu_2, \sigma_2^2)$$

We can then look at the posterior distribution over each indicator variable, conditioned on the class probabilities and parameters:

$$\begin{aligned} \mathbf{P}(Z_i = z | X_i, \pi, \mu_1, \sigma^2) &\propto P(Z_i = z | \pi) p(X_i | \mu_z, \sigma^2) \\ \text{so, } \mathbf{P}(Z_i = 1 | X_i, \pi, \mu_1, \sigma^2) &\propto \pi p(X_i | \mu_1, \sigma^2) \\ \mathbf{P}(Z_i = 0 | X_i, \pi, \mu_1, \sigma^2) &\propto (1 - \pi) p(X_i | \mu_2, \sigma^2) \end{aligned}$$

Conditioned on the Z_i , we can update the means of the Gaussians using conjugacy.

Note that we are not guaranteed to find latent clusters that correspond to the covariate we were expecting! If there is a more parsimonious partitioning of the data, then the posterior will tend to favor that partitioning.

Answer: The regression above was completed, with the X and Y data scaled such that the mean of each set was zero. The optimal β value was calculated as 0.4166, and the distributions are included below.

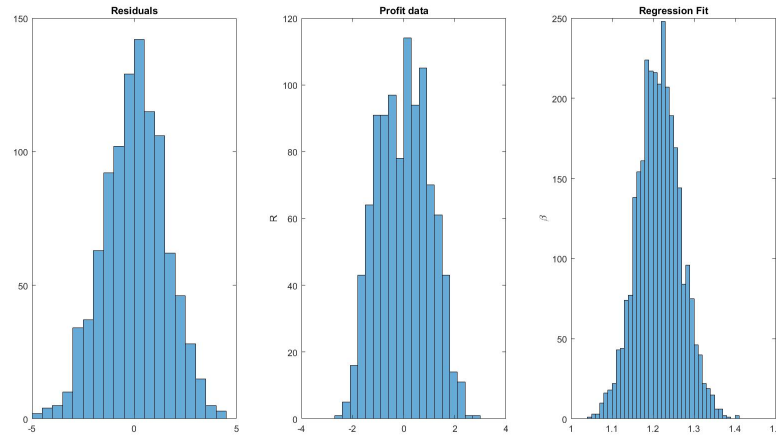


Figure 5.1: Histograms of the residuals, the Y data, and the beta values, with the mean at 0.4166.

The RMSE for the regression is 1.5174.

Exercise 5.2 Let's assume (as is the case if our latent variables correspond to the actual *DinnerService* covariate) that the class proportions are roughly equal, and fix $\pi = 0.5$. Using the conditional distributions $P(Z_i|X_i, \pi, \mu_1, \mu_2, \sigma^2)$ and $p(\mu_k|\{X_i : Z_i = k\}, \theta)$, where θ are appropriate (shared) prior parameters for μ_k , implement a Gibbs sampler that samples the means and the latent indicator variables. I'd suggest using the parameters of the initial regression to pick your hyperparameters.

Compare the clustering obtained with the "true" clustering due to the *DinnerService* variable.

Answer: For this exercise, sigma was kept at 0.4, inspired by the solution to the previous exercise. It was found that the results were less accurate when sampling sigma. Prior to solving, the data was sorted and whitened.

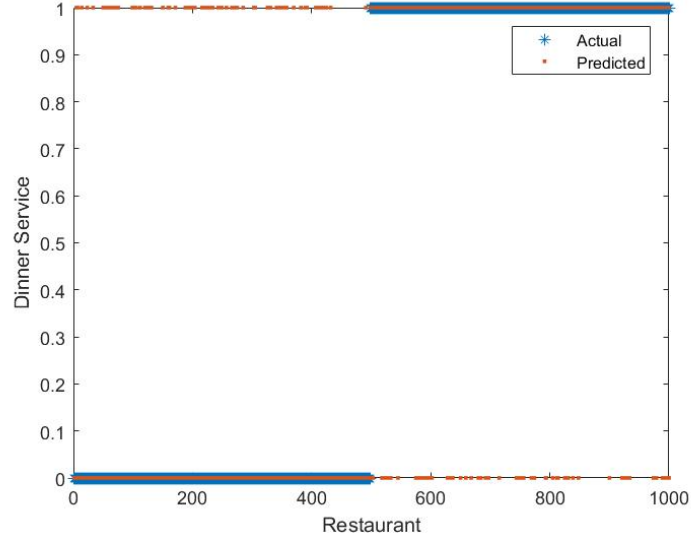


Figure 5.2: Predicted and actual values for the variable Dinner Service.

OK, let's now assume we don't know π , and that the two classes have different values of σ^2 . Let's put a $\text{Beta}(\alpha, \beta)$ prior on π , since it is conjugate to the Bernoulli distribution.

Exercise 5.3 *Let's assume we want to integrate out π . What is the conditional distribution $P(Z_i|Z_{-i}, X_i, \mu_1, \mu_2, \sigma_1, \sigma_2, \alpha, \beta)$, where Z_{-i} means all the values of Z except Z_i ?*

Answer: First, I designate a vector of hyperparameters: $\theta = \{\mu_1, \mu_2, \sigma_1, \sigma_2\}$. The conditional distribution is of the beta-binomial form, which can be derived from the normal likelihood and the posterior.

$$P(Z_i|Z_{-i}, X_i, \theta, \alpha, \beta) = \int P(Z_i|X_i, \pi, \theta) p(\pi|Z_{-i}, \alpha, \beta) d\pi \quad (5.1)$$

$$\propto \int \pi^{Z_i} (1 - \pi)^{1-Z_i} N(X_i|\mu_{Z_i}, \sigma_{Z_i}^2) \text{Beta}(\pi|\alpha + \sum_{j \neq i} Z_j, \beta + n - \sum_{j \neq i} Z_j) d\pi \quad (5.2)$$

$$\propto N(X_i|\mu_{Z_i}, \sigma_{Z_i}^2) \frac{\Gamma(\alpha_n + \beta_n)}{\Gamma(\alpha_n)\Gamma(\beta_n)} \int \pi^{\alpha_n + Z_i} (1 - \pi)^{\beta_n + 1 - Z_i} d\pi \quad (5.3)$$

$$\propto N(X_i|\mu_{Z_i}, \sigma_{Z_i}^2) \frac{\Gamma(\alpha_n + \beta_n)}{\Gamma(\alpha_n)\Gamma(\beta_n)} \frac{\Gamma(\alpha_n + Z_i)\Gamma(\beta_n + 1 - Z_i)}{\Gamma(\alpha_n + \beta_n + 1)} \quad (5.4)$$

$$\propto N(X_i|\mu_{Z_i}, \sigma_{Z_i}^2) \text{Beta} - \text{Binom}(Z_i|1, \alpha_n, \beta_n) \quad (5.5)$$

$$\propto N(X_i|\mu_{Z_i}, \sigma_{Z_i}^2) \alpha_n^{Z_i} \beta_n^{1-Z_i} \quad (5.6)$$

In the above, $\alpha_n = \alpha + \sum_{j \neq i} Z_j$, and $\beta_n = \beta + n - \sum_{j \neq i} Z_j$. The addition of π in the above enables use of the cluster proportion of the sample.

Exercise 5.4 *How about if we want to integrate out all of the continuous variables? What is the conditional distribution $P(Z_i|Z_{-i}, X, \theta)$, where θ is the set of all hyperparameters?*

Answer: This problem is similar to exercise 2.2, and therefore the solution will be used. Updating the model from the previous exercise, where $\lambda = \sigma^2$, we get:

$$y = \pi_1 N(\mu_1, \lambda) + \pi_2 N(\mu_2, \lambda) \quad (5.7)$$

$$\mu_1, \mu_2 \sim N(0, \lambda) \quad (5.8)$$

$$\lambda \sim \text{Gamma}(\alpha, \beta) \quad (5.9)$$

By noting the conjugacy of $\{\lambda_1 \lambda_2, \mu_1, \mu_2\}$, the parameters can be integrated out, by combining this with the solution from exercise 2.2, we get:

$$p(\lambda_1|y) \sim \text{Gamma}\left(\frac{n_1}{2} + \alpha, \beta + \frac{\sum_{i:Z_i=1} y_i^2}{2} - \frac{(\sum_{i:Z_i=1} y_i)^2}{2n_1}\right) \quad (5.10)$$

$$p(Z_i|Z_{-i}, X, \theta) = \frac{\Gamma(\alpha + n_1)\Gamma(\beta + n_2)\Gamma(a_1)\Gamma(a_2)}{\Gamma(\alpha + \beta + n)b_1^{a_1}b_2^{a_2}} \quad (5.11)$$

Where $a_1 = \frac{n_1}{2} + \alpha$, $a_2 = \frac{n_2}{2} + \alpha$, $b_1 = \beta + \frac{\sum_{i:Z_i=1} y_i^2}{2} - \frac{(\sum_{i:Z_i=1} y_i)^2}{2n_1}$, $b_2 = \beta + \frac{\sum_{i:Z_i=0} y_i^2}{2} - \frac{(\sum_{i:Z_i=0} y_i)^2}{2n_1}$.

Exercise 5.5 Implement a Gibbs sampler for this new model where we learn the cluster proportions. You can either implement one of the variants in the previous two exercises, or the fully uncollapsed model where we sample Z , π , μ_1 , μ_2 , σ_1^2 and σ_2^2 .

Answer: To calculate the cluster proportions, the following terms were added onto the likelihoods in problem 5.2:

$$p(z_i = 1|z_{-i}, \alpha) \propto \frac{(N_1 + \alpha)}{N - 1 + \alpha + \beta} N_{pdf}(\mu_1, \sigma^2) \quad (5.12)$$

$$p(z_i = 0|z_{-i}, \alpha) \propto \frac{(N_0 + \alpha)}{N - 1 + \alpha + \beta} N_{pdf}(\mu_0, \sigma^2) \quad (5.13)$$

The error between the actual and predicted counts was: $E_0 = 3.21\%$ and $E_1 = 3.19\%$, the plot shown below tracks the counts for each iteration, showing convergence around iteration 10 for both 1 and 0.

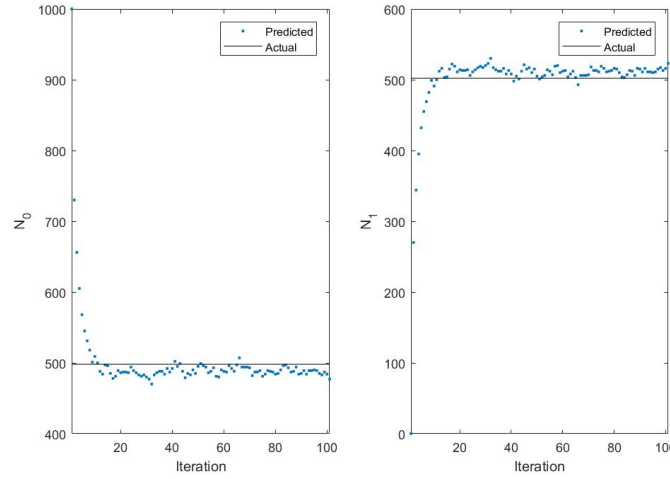


Figure 5.3: Predicted and actual counts for the variable Dinner Service for each iteration.

Let's now consider the case where we have more than two classes. Here, we need to replace our Bernoulli distribution with a multinomial parametrized by some probability vector π , so that:

$$P(Z_i = k) = \pi_k$$

Exercise 5.6 Much as the multinomial is the multivariate generalization of the binomial distribution, the $\text{Dirichlet}(\alpha_1, \dots, \alpha_K)$ distribution, which has pdf

$$\frac{\Gamma(\sum_k \alpha_k)}{\prod_k \Gamma(\alpha_k)} \prod_{k=1}^K \pi_k^{\alpha_k - 1},$$

is the multivariate generalization of the beta distribution. Show that the Dirichlet is conjugate to the multinomial, and derive the posterior predictive distribution

$$P(Z_{n+1} | Z_{1:n}) = \int_{\mathcal{M}} P(Z_{n+1} | \pi) p(\pi) d\pi$$

You may find it helpful to note that, if $\pi \sim \text{Dirichlet}(\alpha_1, \dots, \alpha_K)$, then $E[\pi] = \frac{(\alpha_1, \dots, \alpha_K)}{\sum_k \alpha_k}$.

Exercise 5.7 Modify your previous Gibbs sampler to allow multiple classes, and two-dimensional data. Generate some data according to a Dirichlet mixture of 5 Gaussians in \mathbb{R}^2 , and test your code on it.

Exercise 5.8 OK, let's try a real dataset! We're going to use a set of images from MNIST. Download the dataset `mnist.csv` from the data directory, and transform it to be zero mean, unit variance. Each row contains the vectorized pixel values for an image of a digit. The whole dataset contains 100 copies of each digit, with the first 100 being zeros, the next 100 being ones, etc. You can visualize a data point by reshaping it to be 28×28 :

- R: `image(matrix(X[1,],nrow=28))`
- Python: `import matplotlib.pyplot; plt.imshow(X[0,:].reshape(28,28)); plt.show()`
- Matlab: `imshow(reshape(X(1,:),28,28))`

The data is 784-dimensional; let's reduce this by running PCA and using the first 50 dimensions.

Now, try running your Gibbs sampler with 10 classes, and $\alpha_1 = \alpha_2 = \dots = \alpha_{10} = 1$. This prior corresponds to a uniform distribution on the 9-simplex. It's fine to use a spherical covariance here... in fact it will work fine if you just have a prior on the means, and fix $\sigma^2 = 1$.

Here are some ways you can visualize your output:

- Based on a single sample, plot the recovered clustering vs the ground truth clustering.
- Based on a single sample, visualize the mean image for each cluster, by multiplying the mean embedding with the coefficients obtained using PCA.
- Over multiple samples, create a co-occurrence matrix with entries being the proportion of the times that the two data points are in the same sample.

Exercise 5.9 OK, let's try a different likelihood. Let's consider modeling documents. A common modeling assumption is to treat a document as a “bag-of-words” – assuming that all the information is in the words, and none of it is in the ordering. Under this assumption, an appropriate distribution is a multinomial distribution over words, with a Dirichlet prior. Concretely, let:

$$\begin{aligned}\pi &\sim \text{Dirichlet}_K(\alpha) \\ \eta_k &\sim \text{Dirichlet}_V(\beta), \quad k = 1, \dots, K \\ z_i &\sim \text{Discrete}(\pi), \quad i = 1, \dots, N \\ \mathbf{w}_i &\sim \text{Multinomial}(\eta_{z_i})\end{aligned}$$

where N is the number of documents, V is the number of words in the dictionary, K is the number of clusters, and \mathbf{w}_i is a V -dimensional count vector representing the i th document.

Write out the conditional distributions for a collapsed (i.e. integrating out π and the η_k) Gibbs sampler for this model.

Exercise 5.10 Implement the code. Generate a test set by generating data from a mixture of two multinomials, one with probabilities $(1, 1, 1, 1, 9, 9, 9, 9)/40$ and the other with probabilities $(9, 9, 9, 9, 1, 1, 1, 1)/40$. Test your code on this dataset, and compare a single sample's clustering pattern with the ground truth values.

Once you've got it to work on the toy data, try it on some real data! The file `cora.csv` on Github contains a bag-of-words representation of a collection of 2410 scientific documents from the Cora search engine (taken from the R package `lda`). Each row corresponds to a document, each column to a word, each element is the number of times that word appears in that document. The list of words is at `cora_vocab.csv`. Try clustering them into say 10 clusters. The NIPS dataset on Github contains the text of NIPS papers. Try clustering them into say 10 clusters. Based on a single sample for each cluster, report the 10 most frequently occurring words.

5.0.1 Admixture models

A mixture model for text isn't massively realistic. Consider the NIPS papers: is it really reasonable to separate multiple documents into distinct clusters? It is more likely that two papers share some aspects in common, but differ on others.

We can use a hierarchical Bayesian formulation to model each document using a mixture model, with a shared prior on the mixing components. Concretely, let

$$\begin{aligned}\theta_i &\sim \text{Dirichlet}_K(\alpha), & i = 1, \dots, N \\ \eta_k &\sim \text{Dirichlet}_V(\beta), & k = 1, \dots, K \\ z_{i,j} &\sim \text{Discrete}(\theta_i), & j = 1, \dots, M_i \\ w_{i,j} &\sim \text{Discrete}(\eta_{z_{i,j}}),\end{aligned}$$

where M_i is the number of words in the j th document. This model is commonly known as Latent Dirichlet Allocation ?; it is an example of an *admixture* model.

This means that each document is associated with a distribution θ_i over clusters, and each word is associated with a single cluster.

Exercise 5.11 We can construct a collapsed Gibbs sampler for this model by integrating out the θ_i and the η_k . Derive the predictive distributions $p(z_{i,j}|\{z_{\neg i,j}\}, \alpha)$ and $p(w_{i,j}|z_{i,j}, z_{\neg i,j}, w_{\neg i,j}, \beta)$, and hence the conditional distribution $p(z_{i,j}|\text{rest})$

Exercise 5.12 I'm not going to make you implement this one (although if you want to, feel free!). Instead, let's use the R package `lda` (sorry Python/R folk! it should be fairly easy to use). The documentation is here: <https://cran.r-project.org/web/packages/lda/lda.pdf>. Run the Gibbs sampler on the built-in document dataset `cora`, and report the 5 words with highest probability for each cluster (hint: look at the example under `top.topic.words` – note that you might need more iterations than is given in the example, R has a rule that examples have to run quickly, hence the low number in the example). Why is this sort of model commonly called a topic model?

5.1 Bayesian nonparametric models

When we were modeling the MNIST dataset, we used 10 clusters. This seems reasonable, right – there are 10 digits! However, if you look at the data, there is a lot of variation within each digit. Maybe we'd be better off using more clusters... but how many?

One answer to this question is to allow *infinitely* many clusters *a priori*. Each data point can only belong to a single cluster, so there will only be at most N occupied clusters. By allowing infinitely many clusters, we can allow N data points to occupy a random number of clusters. Further, if we see more data, we are not restricted to the previously occupied clusters.

Exercise 5.13 To get a feel for this, we can “approximate” a model with infinitely many clusters with a model with a large number of clusters. Let's start with a Dirichlet prior on cluster membership, with 100 clusters.

Sample $\pi \sim \text{Dirichlet}_{100}(10, 10, \dots, 10)$, and then sample 10 cluster indicators $z_i \sim \pi$. Record the list of cluster indicators, e.g. $\{1, 10, 11, 11, \dots\}$. Do this 5 times, with a different π each time.

Repeat this with $\alpha = (1, 1, \dots, 1)$, $\alpha = (0.1, 0.1, \dots, 0.1)$ and $\alpha = (0.01, 0.01, \dots, 0.01)$.

Comment on how the value of α affects your clustering behavior.

OK, now let's explore some further properties of the Dirichlet distribution. First, we note an important relationship between the Dirichlet distribution and the gamma distribution: If

$$\gamma_i \stackrel{\text{iid}}{\sim} \text{Gamma}(\alpha_i, \beta)$$

then

$$Z = \sum_{i=1}^K \gamma_i \sim \text{Gamma}\left(\sum_{i=1}^K \alpha_i, \beta\right)$$

and

$$\pi = \left(\frac{\gamma_1}{Z}, \dots, \frac{\gamma_K}{Z}\right) \sim \text{Dirichlet}(\alpha_1, \alpha_K)$$

Exercise 5.14 Using the change-of-variable technique with the transform $(\gamma_1, \dots, \gamma_K) \rightarrow (\pi_1, \dots, \pi_{K-1}, Z)$, prove the above result.

You will probably find this relationship helpful in proving the following

Exercise 5.15 (Agglomeration property) Show that, if $(\pi_1, \dots, \pi_K) \sim \text{Dirichlet}(\alpha_1, \dots, \alpha_K)$, then $(\pi_1 + \pi_2, \dots, \pi_K) \sim \text{Dirichlet}(\alpha_1 + \alpha_2, \alpha_3, \dots, \alpha_K)$.

Exercise 5.16 Let $\pi \sim \text{Dirichlet}_K\left(\frac{\alpha}{K}, \dots, \frac{\alpha}{K}\right)$, and assign weight π_k to the interval $\left[\frac{k-1}{K}, \frac{k}{K}\right)$. Show that, for any partition with breaks at multiples of $\frac{1}{k}$, the distribution over the weights associated with the blocks in the partition will be Dirichlet distributed.

The Dirichlet process extends this idea to arbitrary partitions. Concretely, the Dirichlet process is a distribution over measures¹ on some space \mathcal{X} , parametrized by some probability distribution H on \mathcal{X} and some positive scalar α such that for any partition A_1, \dots, A_K of \mathcal{X} , the masses assigned to A_1, \dots, A_K are distributed according to a Dirichlet $(\alpha H(A_1), \dots, \alpha H(A_K))$ distribution. The resulting probability distribution D will have its probability concentrated on infinitely many singletons $D = \sum_{i=1}^{\infty} \pi_i \delta_{\theta_i}$ —what is known as an atomic probability distribution.

We can construct a finite dimensional approximation to the Dirichlet process by sampling $\pi \sim \text{Dirichlet}_K\left(\frac{\alpha}{K}, \dots, \frac{\alpha}{K}\right)$ for some large α , and associating each probability π_k with a location $\theta_k \sim H$. This distribution will converge weakly to the Dirichlet process as $K \rightarrow \infty$.

Exercise 5.17 Return to the MNIST mixture model, and replace your 10-dimensional Dirichlet distribution with a 100-dimensional Dirichlet with parameters $\alpha/100$ for, say, $\alpha = 1$. How many clusters does it use (look at a distribution over multiple samples)? Based on a single sample, what do those clusters look like?

¹If you're not familiar with measure theory, a measure on some space is just a function that assigns a positive number to every subset of that space. So, a probability is a measure. Area is a measure.