

Feel free to work with other students, but make sure you write up the homework and code on your own (no copying homework *or* code; no pair programming). Feel free to ask students or instructors for help debugging code or whatever else, though.

The starter files for problem 2 can be found under the Resource tab on course website. The plot for problem 2 generated by the sample solution has been included in the starter files for reference. Please print out all the graphs generated by your own code and submit them together with the written part, and make sure you upload the code to your Github repository.

1 (Murphy 11.3 - EM for Mixtures of Bernoullis) Show that the M step for ML estimation of a mixture of Bernoullis is given by

$$\mu_{kj} = \frac{\sum_i r_{ik} x_{ij}}{\sum_i r_{ik}}.$$

Show that the M step for MAP estimation of a mixture of Bernoullis with a $\beta(a, b)$ prior is given by

$$\mu_{kj} = \frac{(\sum_i r_{ik} x_{ij}) + a - 1}{(\sum_i r_{ik}) + a + b - 2}.$$

(a) Get complete data log-likelihood:

$$\begin{aligned}\ell(\boldsymbol{\mu}) &= \sum_i \sum_k r_{ik} \log \mathbb{P}(\mathbf{x}_i \mid \boldsymbol{\theta}_k) \\ &= \sum_i \sum_k r_{ik} \sum_j \mathbf{x}_{ij} \log \mu_{kj} + (1 - \mathbf{x}_{ij}) \log (1 - \mu_{kj})\end{aligned}$$

The variables are: i is datapoint index, k is component, j is the dimension index for D dimensional bit vectors.

The derivative with respect to μ_{kj} to get the optimality condition eventually

$$\begin{aligned}
\frac{\partial \ell}{\partial \mu_{kj}} &= \sum_i r_{ik} \left(\frac{\mathbf{x}_{ij}}{\mu_{kj}} - \frac{1 - \mathbf{x}_{ij}}{1 - \mu_{kj}} \right) \\
&= \sum_i r_{ik} \left(\frac{\mathbf{x}_{ij} - \mu_{kj}}{\mu_{kj} (1 - \mu_{kj})} \right) \\
&= \frac{1}{\mu_{kj} (1 - \mu_{kj})} \sum_i r_{ik} (\mathbf{x}_{ij} - \mu_{kj}) = 0.
\end{aligned}$$

$$\sum_i r_{ik} \mathbf{x}_{ij} = \mu_{kj} \sum_i r_{ik}$$

(b) Complete data log likelihood + log prior without π terms

$$\begin{aligned}
\ell(\boldsymbol{\mu}) &= \sum_i \sum_k r_{ik} \log \mathbb{P}(\mathbf{x}_i | \boldsymbol{\mu}_k) + \log \mathbb{P}(\boldsymbol{\mu}_k) \\
&= \sum_i \sum_k r_{ik} \left(\sum_j \mathbf{x}_{ij} \log \mu_{kj} + (1 - \mathbf{x}_{ij}) \log (1 - \mu_{kj}) \right) + \\
&\quad (a - 1) \log \mu_{kj} + (b - 1) \log (1 - \mu_{kj}).
\end{aligned}$$

Taking derivatives to get optimality condition:

$$\begin{aligned}
\frac{\partial \ell}{\partial \mu} &= \sum_i \left(\frac{r_{ik} \mathbf{x}_{ij} + a - 1}{\mu_{kj}} - \frac{r_{ik} (1 - \mathbf{x}_{ij}) + b - 1}{1 - \mu_{kj}} \right) \\
&= \frac{1}{\mu_{kj} (1 - \mu_{kj})} \sum_i r_{ik} \mathbf{x}_{ij} - r_{ik} \mu_{kj} + a - 1 - \mu_{kj} a + \mu_{kj} - \mu_{kj} b + \mu_{kj} \\
&= \frac{1}{\mu_{kj} (1 - \mu_{kj})} \left[\sum_i r_{ik} \mathbf{x}_{ij} - \left(\sum_i r_{ik} + a + b - 2 \right) \mu_{kj} + a - 1 \right] = 0.
\end{aligned}$$

$$\sum_i r_{ik} \mathbf{x}_{ij} + a - 1 = \left(\sum_i r_{ik} + a + b - 2 \right) \mu_{kj}$$

■

2 (Lasso Feature Selection) In this problem, we will use the online news popularity dataset we used in hw2pr3. In the starter code, we have already parsed the data for you. However, you might need internet connection to access the data and therefore successfully run the starter code.

First, ignoring undifferentiability at $x = 0$, take $\frac{\partial |x|}{\partial x} = \text{sign}(x)$. Using this, show that $\nabla \|\mathbf{x}\|_1 = \text{sign}(\mathbf{x})$ where sign is applied elementwise. Derive the gradient of the ℓ_1 regularized linear regression objective

$$\text{minimize: } \|\mathbf{A}\mathbf{x} - \mathbf{b}\|_2^2 + \lambda \|\mathbf{x}\|_1$$

Then, implement a gradient descent based solution of the above optimization problem for this data. Produce the convergence plot (objective vs. iterations) for a non-trivial value of λ . In the same figure (and different axes) produce a 'regularization path' plot. Detailed more in section 13.3.4 of Murphy, a regularization path is a plot of the optimal weight on the y axis at a given regularization strength λ on the x axis. Armed with this plot, provide an ordered list of the top five features in predicting the log-shares of a news article from this dataset (with justification).

$$\text{prox}_\gamma(\mathbf{x})_i = \begin{cases} \mathbf{x}_i - \gamma & \mathbf{x}_i > \gamma \\ 0 & |\mathbf{x}_i| \leq \gamma \\ \mathbf{x}_i + \gamma & \mathbf{x}_i < -\gamma \end{cases}$$

so that each iterate

$$\mathbf{x}_{i+1} = \text{prox}_\gamma(\mathbf{x}_i - \gamma \nabla f(\mathbf{x}_i))$$

where γ is learning rate.

$$\frac{\partial \|\mathbf{x}\|_1}{\partial \mathbf{x}_i} = \frac{\partial \sum |\mathbf{x}_i|}{\partial \mathbf{x}_i} = \text{sign}(\mathbf{x}_i) \text{ similar to } \nabla \|\mathbf{x}\|_1 = \text{sign}(\mathbf{x}_i)$$

$$\begin{aligned} \nabla \|\mathbf{A}\mathbf{x} - \mathbf{b}\|_2^2 + \lambda \|\mathbf{x}\|_1 &= \nabla \mathbf{x}^\top \mathbf{A}^\top \mathbf{A} \mathbf{x} - 2\mathbf{b}^\top \mathbf{A} \mathbf{x} + \mathbf{b}^\top \mathbf{b} + \lambda \|\mathbf{x}\|_1 \\ &= 2\mathbf{A}^\top \mathbf{A} \mathbf{x} - 2\mathbf{b}^\top \mathbf{A} + \lambda \text{sign}(\mathbf{x}). \end{aligned}$$

The plot displayed below shows that since we initialized our weights using the least squares estimate, the lasso objective function value does not change significantly. However, if we examine the sparsity over iterations, we can observe that it actually increases. The most important features are: 'timedelta', 'weekday_is_wednesday', 'weekday_is_thursday', 'weekday_is_friday', and 'weekday_is_saturday'.

■

