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

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

We have the complete data bog likelihood $l(\mu) = \sum_{i} \sum_{k} r_{ik} \log P(x_i | \theta_k) = \sum_{i} \sum_{k} r_{ik} \sum_{j} \log M_{kj} + (l-x_{ij}) \log (l-M_{kj})$ where i is the datapoint index, k is the compnent, and j is the dimension index of the P dimensional bet vertors. Taking the derivative with respect to M_{kj} , we have $\frac{\partial l}{\partial M_{kj}} = \sum_{i} r_{ik} \left(\frac{x_{ij}}{M_{kj}} - \frac{l-x_{ij}}{l-M_{kj}} \right) = \sum_{i} r_{ik} \left(\frac{x_{ij}-M_{kj}}{M_{kj}(l-M_{kj})} \right) = \frac{\sum_{i} r_{ik}(x_{ij}-M_{kj})}{M_{kj}(l-M_{kj})} = 0.$ This gives the optimality condition $\sum_{i} r_{ik} x_{ij} = M_{kj} \sum_{i} r_{ik}$, which gives the derived verilt.

1

We have the complete data log likelihood plus the log prior (ignoring the π term as no are maximizing without regard to them) $l(\mu): \sum_{i} \sum_{k} r_{ik} \log P(x_{i}|\mu_{k}) + \log P(\mu_{k})$ $= \sum_{i} \sum_{k} r_{ik} (\sum_{i} x_{ij} \log \mu_{kj} + (1-x_{ij}) \log (1-\mu_{kj})) + (a-1) \log \mu_{kj} + (b-1) \log (1-\mu_{kj}).$

Taking derivatives, we have $\frac{\partial L}{\partial L} = \sum_{i} \left(\frac{Y_{i} | X_{ij} + \alpha - I}{M_{ij}} - \frac{Y_{i} | (I - X_{ij}) + b - I}{I - M_{ij}} \right)$ $= \frac{I}{M_{i} | (I - M_{ij})} \sum_{i} Y_{i} | X_{i} | - Y_{i} | M_{i} | + \alpha - I - M_{i} | \alpha + M_{i} | - M_{i} | M_{i} |$ $= \frac{I}{M_{i} | (I - M_{i} | i)} \left[\sum_{i} Y_{i} | | X_{i} | - M_{i} | (\sum_{i} Y_{i} | | | + \alpha + b - 2) + \alpha - I \right] = 0.$

This gives the optimality condition $\sum_{i=1}^{n} x_i L^{n} X_i j + a - 1 = (\sum_{i=1}^{n} x_i L^{n} L^{$

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

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

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).

First, let's wrap the gradient descent step with a thredhold function
$$\begin{cases} \chi_i - \gamma & \chi_i > \gamma \\ 0 & |\chi_i| \leq \gamma \\ \chi_i + \gamma & \chi_i < -\gamma \end{cases}$$

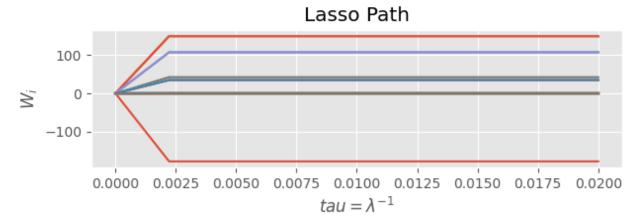
So that each iterate $\chi_{i+1} = pro\chi_{\gamma}(\chi_{i} - \gamma \nabla f(\chi_{i}))$ where γ is the learning rate.

Thus,
$$\frac{\partial ||x||_1}{\partial x_i} = \frac{\partial \sum |x_i|}{\partial x_i} = \text{sign}(x_i)$$
. It follows that $\nabla ||x||_1 = \text{sign}(x_i)$.

We can then see
$$\nabla \|Ax - b\|_{2}^{2} + \lambda \|x\|_{1} = \nabla x^{T} A^{T} A_{x} - 2b^{T} A_{x} + b^{T} b + \lambda \|x\|_{1}$$

= $2A^{T} A_{x} - 2b^{T} A + \lambda sign(x)$.

The plot is attached below. Since we instantiated our weights with the least squares estimate, we can see our losso objective not moving significantly, although if we look at sparsity over time it does in fact increase.



Lasso Objective Convergence: $\lambda = 1e5$

