

There are 6 problems in this set. You need to do 3 problems (due in class on Monday) every week for 2 weeks. Note that this means you must eventually complete all problems. 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. When implementing algorithms you may not use any library (such as sklearn) that already implements the algorithms but you may use any other library for data cleaning and numeric purposes (numpy or pandas). Use common sense. Problems are in no specific order.

**1. (Conditioning a Gaussian)** Note that from Murphy page 113. “Equation 4.69 is of such importance in this book that we have put a box around it, so you can easily find it.” That equation is important. Read through the proof of the result. Suppose we have a distribution over random variables  $\mathbf{x} = (\mathbf{x}_1, \mathbf{x}_2)$  that is jointly Gaussian with parameters

$$\boldsymbol{\mu} = \begin{bmatrix} \mu_1 \\ \mu_2 \end{bmatrix} \quad \boldsymbol{\Sigma} = \begin{bmatrix} \boldsymbol{\Sigma}_{11} & \boldsymbol{\Sigma}_{12} \\ \boldsymbol{\Sigma}_{21} & \boldsymbol{\Sigma}_{22} \end{bmatrix},$$

where

$$\mu_1 = \begin{bmatrix} 0 \\ 0 \end{bmatrix}, \quad \mu_2 = 5, \quad \boldsymbol{\Sigma}_{11} = \begin{bmatrix} 6 & 8 \\ 8 & 13 \end{bmatrix}, \quad \boldsymbol{\Sigma}_{21}^\top = \boldsymbol{\Sigma}_{12} = \begin{bmatrix} 5 \\ 11 \end{bmatrix}, \quad \boldsymbol{\Sigma}_{22} = [14].$$

Compute

- (a) The marginal distribution  $p(\mathbf{x}_1)$ . Plot the density in  $\mathbb{R}^2$ .
- (b) The marginal distribution  $p(\mathbf{x}_2)$ . Plot the density in  $\mathbb{R}^1$ .
- (c) The conditional distribution  $p(\mathbf{x}_1|\mathbf{x}_2)$
- (d) The conditional distribution  $p(\mathbf{x}_2|\mathbf{x}_1)$

**2. ( $\ell_1$ -Regularization)** Consider the  $\ell_1$  norm of a vector  $\mathbf{x} \in \mathbb{R}^n$ :

$$\|\mathbf{x}\|_1 = \sum_i |\mathbf{x}_i|.$$

Plot the norm-ball  $B_k = \{\mathbf{x} : \|\mathbf{x}\|_1 \leq k\}$  for  $k = 1$ . On the same plot, plot the Euclidean norm-ball  $A_k = \{\mathbf{x} : \|\mathbf{x}\|_2 \leq k\}$  for  $k = 1$  behind the first plot. Show that the optimization problem

$$\begin{array}{ll} \text{minimize:} & f(\mathbf{x}) \\ \text{subj. to:} & \|\mathbf{x}\|_p \leq k \end{array}$$

is equivalent to

$$\text{minimize: } f(\mathbf{x}) + \lambda \|\mathbf{x}\|_p$$

(hint: create the Lagrangian). With this knowledge, and the plots given above, argue why using  $\ell_1$  regularization (adding a  $\lambda \|\mathbf{x}\|_1$  term to the objective) will give sparser solutions than using  $\ell_2$  regularization for suitably large  $\lambda$ .

**3. (Lasso)** Show that placing an equal zero-mean Laplace prior on each element of the weights  $\boldsymbol{\theta}$  of a model is equivalent to  $\ell_1$  regularization in the Maximum-a-Posteriori estimate

$$\text{maximize: } \mathbb{P}(\boldsymbol{\theta}|\mathcal{D}) = \frac{\mathbb{P}(\mathcal{D}|\boldsymbol{\theta})\mathbb{P}(\boldsymbol{\theta})}{\mathbb{P}(\mathcal{D})}.$$

Note the form of the Laplace distribution is

$$\text{Lap}(x|\mu, b) = \frac{1}{2b} \exp\left(-\frac{|x - \mu|}{b}\right)$$

where  $\mu$  is the location parameter and  $b > 0$  controls the variance. Plot the density  $\text{Lap}(x|0, 1)$  and the standard normal  $\mathcal{N}(x|0, 1)$  and suggest why this would lead to sparser solutions than a Gaussian prior on each elements of the weights (which correspond to  $\ell_2$  regularization).

**4. (Lasso Feature Selection)** 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  $\ell_1$  regularized linear regression objective

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

Now consider the shares dataset we used in problem 1 of homework 1 ([https://math189r.github.io/hw/data/online\\_news\\_popularity/online\\_news\\_popularity.txt](https://math189r.github.io/hw/data/online_news_popularity/online_news_popularity.txt)). 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  $\lambda$ . 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  $\lambda$  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). We can see a more detailed analysis of this at [https://en.wikipedia.org/wiki/Proximal\\_gradient\\_methods\\_for\\_learning](https://en.wikipedia.org/wiki/Proximal_gradient_methods_for_learning) and [https://web.stanford.edu/~boyd/papers/pdf/prox\\_algs.pdf](https://web.stanford.edu/~boyd/papers/pdf/prox_algs.pdf) but you will have to wrap the gradient descent step with a threshold function

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

so that each iterate

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

where  $\gamma$  is your learning rate. Tip: you can reuse most of your code from the first homework.

**5. (SVD Image Compression)** Load the image of a scary clown at <http://i.imgur.com/X017qGH.jpg> into a matrix/array. Plot the progression of the 100 largest singular values for the original image and a randomly shuffled version of the same image (all on the same plot). In a single figure plot a grid of four images: the original image, and a rank  $k$  truncated SVD approximation of the original image for  $k \in \{2, 10, 20\}$ .

**6. (Murphy 12.5 - Deriving the Residual Error for PCA)** It may be helpful to reference section 12.2.2 of Murphy.

(a) Prove that

$$\left\| \mathbf{x}_i - \sum_{j=1}^k z_{ij} \mathbf{v}_j \right\|^2 = \mathbf{x}_i^{\top} \mathbf{x}_i - \sum_{j=1}^k \mathbf{v}_j^{\top} \mathbf{x}_i \mathbf{x}_i^{\top} \mathbf{v}_j.$$

Hint: first consider the case when  $k = 2$ . Use the fact that  $\mathbf{v}_i^{\top} \mathbf{v}_j$  is 1 if  $i = j$  and 0 otherwise. Recall that  $z_{ij} = \mathbf{x}_i^{\top} \mathbf{v}_j$ .

(b) Now show that

$$J_k = \frac{1}{n} \sum_{i=1}^n \left( \mathbf{x}_i^{\top} \mathbf{x}_i - \sum_{j=1}^k \mathbf{v}_j^{\top} \mathbf{x}_i \mathbf{x}_i^{\top} \mathbf{v}_j \right) = \frac{1}{n} \sum_{i=1}^n \mathbf{x}_i^{\top} \mathbf{x}_i - \sum_{j=1}^k \lambda_j.$$

Hint: recall that  $\mathbf{v}_j^{\top} \Sigma \mathbf{v}_j = \lambda_j \mathbf{v}_j^{\top} \mathbf{v}_j = \lambda_j$ .

(c) If  $k = d$  there is no truncation, so  $J_d = 0$ . Use this to show that the error from only using  $k < d$  terms is given by

$$J_k = \sum_{j=k+1}^d \lambda_j.$$

Hint: partition the sum  $\sum_{j=1}^d \lambda_j$  into  $\sum_{j=1}^k \lambda_j$  and  $\sum_{j=k+1}^d \lambda_j$ .