Solution Math189R SU17 Homework 5 Monday, June 5, 2017

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.

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

## (a) We know

$$\begin{aligned} \left\| \mathbf{x}_{i} - \sum_{j=1}^{k} z_{ij} \mathbf{v}_{j} \right\|_{2}^{2} &= \left( \mathbf{x}_{i} - \sum_{j=1}^{k} z_{ij} \mathbf{v}_{j} \right)^{\top} \left( \mathbf{x}_{i} - \sum_{j=1}^{k} z_{ij} \mathbf{v}_{j} \right) \end{aligned}$$

$$= \mathbf{x}_{i}^{\top} \mathbf{x}_{i} - \sum_{j=1}^{k} z_{ij} \mathbf{v}_{j}^{\top} \mathbf{x}_{i} - \mathbf{x}_{i}^{\top} \sum_{j=1}^{k} z_{ij} \mathbf{v}_{j} + \left( \sum_{j=1}^{k} z_{ij} \mathbf{v}_{j} \right)^{\top} \left( \sum_{j=1}^{k} z_{ij} \mathbf{v}_{j} \right) \end{aligned}$$

$$= \mathbf{x}_{i}^{\top} \mathbf{x}_{i} - 2 \sum_{j=1}^{k} z_{ij} \mathbf{v}_{j}^{\top} \mathbf{x}_{i} + \left( \sum_{j=1}^{k} z_{ij} \mathbf{v}_{j} \right)^{\top} \left( \sum_{j=1}^{k} z_{ij} \mathbf{v}_{j} \right) \end{aligned}$$

$$= \mathbf{x}_{i}^{\top} \mathbf{x}_{i} - 2 \sum_{j=1}^{k} z_{ij} \mathbf{v}_{j}^{\top} \mathbf{x}_{i} + \sum_{j=1}^{k} \mathbf{v}_{j}^{\top} z_{ij}^{\top} z_{ij} \mathbf{v}_{j} \end{aligned}$$

$$= \mathbf{x}_{i}^{\top} \mathbf{x}_{i} - 2 \sum_{j=1}^{k} z_{ij} \mathbf{v}_{j}^{\top} \mathbf{x}_{i} + \sum_{j=1}^{k} \mathbf{v}_{j}^{\top} \mathbf{x}_{i} \mathbf{x}_{i}^{\top} \mathbf{v}_{j} \qquad \text{(since } \mathbf{v}_{i}^{\top} \mathbf{v}_{j} = 1 \text{ iff } i = j) \end{aligned}$$

$$= \mathbf{x}_{i}^{\top} \mathbf{x}_{i} - 2 \sum_{j=1}^{k} \mathbf{v}_{j}^{\top} \mathbf{x}_{i} \mathbf{x}_{i}^{\top} \mathbf{v}_{j} + \sum_{j=1}^{k} \mathbf{v}_{j}^{\top} \mathbf{x}_{i} \mathbf{x}_{i}^{\top} \mathbf{v}_{j} \qquad \text{(since } z_{ij} \in \mathbb{R}) \end{aligned}$$

$$= \mathbf{x}_{i}^{\top} \mathbf{x}_{i} - 2 \sum_{j=1}^{k} \mathbf{v}_{j}^{\top} \mathbf{x}_{i} \mathbf{x}_{i}^{\top} \mathbf{v}_{j},$$

as desired.

(b) By definition

$$J_{k} = \frac{1}{n} \sum_{i=1}^{n} \left( \mathbf{x}_{i}^{\top} \mathbf{x}_{i} - \sum_{j=1}^{k} \mathbf{v}_{j}^{\top} \mathbf{x}_{i} \mathbf{x}_{i}^{\top} \mathbf{v}_{j} \right)$$
$$= \frac{1}{n} \sum_{i=1}^{n} \mathbf{x}_{i}^{\top} \mathbf{x}_{i} - \sum_{j=1}^{k} \mathbf{v}_{j}^{\top} \frac{1}{n} \left( \sum_{i=1}^{n} \mathbf{x}_{i} \mathbf{x}_{i}^{\top} \right) \mathbf{v}_{j}$$
$$= \frac{1}{n} \sum_{i=1}^{n} \mathbf{x}_{i}^{\top} \mathbf{x}_{i} - \sum_{j=1}^{k} \mathbf{v}_{j}^{\top} \mathbf{\Sigma} \mathbf{v}_{j}$$
$$= \frac{1}{n} \sum_{i=1}^{n} \mathbf{x}_{i}^{\top} \mathbf{x}_{i} - \sum_{j=1}^{k} \lambda_{j}$$

as desired.

(c) Since  $J_d = 0$  we know  $\sum_{j=1}^d \lambda_j = \frac{1}{n} \sum_{i=1}^n \mathbf{x}_i^\top \mathbf{x}_i$ . Then

$$J_k = \frac{1}{n} \sum_{i=1}^n \mathbf{x}_i^\top \mathbf{x}_i - \sum_{j=1}^d \lambda_j + \sum_{j=k+1}^d \lambda_j = \sum_{j=k+1}^d \lambda_j.$$

This is an exciting result. This states that the reconstruction error when using a PCA projection of your data is exactly equal to the sum of the eigenvalues you throw out.

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

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

Draw the norm-ball  $B_k = \{\mathbf{x} : \|\mathbf{x}\|_1 \le k\}$  for k = 1. On the same graph, draw the Euclidean norm-ball  $A_k = \{\mathbf{x} : \|\mathbf{x}\|_2 \le k\}$  for k = 1 behind the first plot. (Do not need to write any code, draw the graph by hand).

Show that the optimization problem

minimize:  $f(\mathbf{x})$ subj. to:  $\|\mathbf{x}\|_p \le k$ 

is equivalent to

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  $\lambda$ .

We see the norm balls below.



We know the optimization problem

minimize: 
$$f(\mathbf{x})$$
  
subj. to:  $\|\mathbf{x}\|_p \le k$ 

is equivalent to

$$\inf_{\mathbf{x}} \sup_{\lambda \ge 0} \mathcal{L}(\mathbf{x}, \lambda) = \inf_{\mathbf{x}} \sup_{\lambda \ge 0} f(\mathbf{x}) + \lambda(\|\mathbf{x}\|_p - k).$$

In its dual we can flip the inf and sup.

$$\sup_{\lambda \ge 0} \inf_{\mathbf{x}} f(\mathbf{x}) + \lambda(\|\mathbf{x}\|_p - k) = \sup_{\lambda \ge 0} g(\lambda)$$

Since the minimizing value of  $f(\mathbf{x}) + \lambda(||\mathbf{x}||_p - k)$  over  $\mathbf{x}$  is equivalent to the minimizing value of  $f(\mathbf{x}) + \lambda ||\mathbf{x}||_p$  ( $-\lambda k$  doesn't depend on  $\mathbf{x}$ ), we know the the optimizing  $\mathbf{x}$  will solve

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

for some suitable value of  $\lambda \ge 0$ . Looking at the plot and this result, we can consider  $\ell_1$  regularization as projecting the actual optimal solution of your problem onto some suitably sized  $\ell_1$  norm ball. Since the  $\ell_1$  ball has sharper edges, the probability of landing on an edge and not on the face (where both elements of the vector are nonzero) is infinitely larger than the  $\ell_2$  ball. This is due to the rotation invariance of the  $\ell_2$  that certainly doesn't hold for the  $\ell_1$  ball. Generalizing to higher dimensions, we can see that the  $\ell_1$  penalty will encourage more weights to be zero compared to the  $\ell_2$  ball, as desired.

**Extra Credit** (Lasso) Show that placing an equal zero-mean Laplace prior on each element of the weights  $\theta$  of a model is equivelent to  $\ell_1$  regularization in the Maximum-a-Posteriori estimate

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

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

where  $\mu$  is the location parameter and b > 0 controls the variance. Draw (by hand) and compare the density 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).

We know the Maximum-a-Posteriori problem

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

is equivalent to maximizing  $\log \mathbb{P}(\theta | D)$  given the monotonicity of  $\log(x)$ . This gives

maximize: 
$$\log \mathbb{P}(\boldsymbol{\theta}|\mathcal{D}) = \log \mathbb{P}(\mathcal{D}|\boldsymbol{\theta}) + \log \mathbb{P}(\boldsymbol{\theta}) - \log \mathbb{P}(\mathcal{D})$$

Since  $\mathbb{P}(\mathcal{D})$  is a constant not dependent on  $\theta$ , we can drop that term from the problem and flip into a minimization problem, giving

minimize:  $-\log \mathbb{P}(\mathcal{D}|\boldsymbol{\theta}) - \log \mathbb{P}(\boldsymbol{\theta}).$ 

Given a prior  $\boldsymbol{\theta}_i \sim \text{Lap}(0, b)$ ,

$$-\log \mathbb{P}(\boldsymbol{\theta}) = -\log \prod_{i} \exp\left(-\frac{|\boldsymbol{\theta}_{i}|}{b}\right) + Z \qquad \text{(where Z is a constant)}$$
$$= \frac{1}{b} \sum_{i} |\boldsymbol{\theta}_{i}| + Z$$
$$= \lambda \|\boldsymbol{\theta}\|_{1} + Z. \qquad \text{(where } \lambda = 1/b\text{)}$$

It follows that our original problem is equivalent to

minimize:  $-\log \mathbb{P}(\mathcal{D}|\boldsymbol{\theta}) + \lambda \|\boldsymbol{\theta}\|_{1}$ ,

or a  $\ell_1$  regularized maximum likelihood estimate, as desired. Note the plots of the Standard Normal and Laplace Densities.



We can see that Lap(0,1) will place much more mass at x = 0. It follows that when we use a Laplace prior instead of a Gaussian prior on our weights, our weights will be more encouraged to be exactly zero, forcing sparsity.