

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

$$\begin{aligned} A) \left\| \mathbf{x}_i - \sum_{j=1}^k z_{ij} \mathbf{v}_j \right\|^2 &= \left(\mathbf{x}_i - \sum_{j=1}^k z_{ij} \mathbf{v}_j \right)^\top \left(\mathbf{x}_i - \sum_{j=1}^k z_{ij} \mathbf{v}_j \right) \\ &= \mathbf{x}_i^\top \mathbf{x}_i - \mathbf{x}_i^\top \sum_{j=1}^k z_{ij} \mathbf{v}_j - \sum_{j=1}^k z_{ij} \mathbf{v}_j^\top \mathbf{x}_i + \left(\sum_{j=1}^k z_{ij} \mathbf{v}_j \right)^\top \left(\sum_{j=1}^k z_{ij} \mathbf{v}_j \right) \\ &= \mathbf{x}_i^\top \mathbf{x}_i - 2 \sum_{j=1}^k z_{ij} \mathbf{v}_j^\top \mathbf{x}_i + \sum_{j=1}^k z_{ij}^2 \mathbf{v}_j^\top \mathbf{v}_j \quad \text{by bringing in } \mathbf{x}_i^\top \text{ into the summation, } \mathbf{v}_j \text{ becomes } \mathbf{v}_j^\top \\ &= \mathbf{x}_i^\top \mathbf{x}_i - 2 \sum_{j=1}^k \mathbf{v}_j^\top \mathbf{x}_i \mathbf{x}_i^\top \mathbf{v}_j + \sum_{j=1}^k \mathbf{v}_j^\top \mathbf{x}_i \mathbf{x}_i^\top \mathbf{v}_j \quad \text{since } \mathbf{v}_i^\top \mathbf{v}_j = 1 \text{ if } i=j \text{ and 0 otherwise, plus in } z_{ij} = \mathbf{x}_i^\top \mathbf{v}_j \\ &= \mathbf{x}_i^\top \mathbf{x}_i - \sum_{j=1}^k \mathbf{v}_j^\top \mathbf{x}_i \mathbf{x}_i^\top \mathbf{v}_j \quad \text{as we wished} \end{aligned}$$

$$B) J_X = \frac{1}{n} \sum_{i=1}^n (X_i^T X_i - \sum_{j=1}^n v_j^T X_i X_i^T v_j)$$

$$= \frac{1}{n} \sum_{i=1}^n X_i^T X_i - \frac{1}{n} \sum_{i=1}^n v_i^T \left(\sum_{j=1}^n X_i X_i^T v_j \right)$$

$$= \frac{1}{n} \sum_{i=1}^n X_i^T X_i - \sum_{j=1}^n v_j^T \sum_{i=1}^n X_i X_i^T v_j$$

$$= \frac{1}{n} \sum_{i=1}^n X_i^T X_i - \sum_{j=1}^n \lambda_j$$

as we wished

$$\text{as } v_i^T \sum_{j=1}^n v_j = \lambda_i, v_j^T v_j = \lambda_j$$

C) We can utilize the fact that $J_d = 0$ by $\sum_{j=1}^d \lambda_j = \frac{1}{n} \sum_{i=1}^n X_i^T X_i$

Thus,

$$J_X = \frac{1}{n} \sum_{i=1}^n X_i^T X_i - \sum_{j=1}^d \lambda_j + \sum_{j=d+1}^n \lambda_j$$

$$= \sum_{j=d+1}^n \lambda_j$$

as desired.

2 (ℓ_1 -Regularization) Consider the ℓ_1 norm of a vector $\mathbf{x} \in \mathbb{R}^n$:

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

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

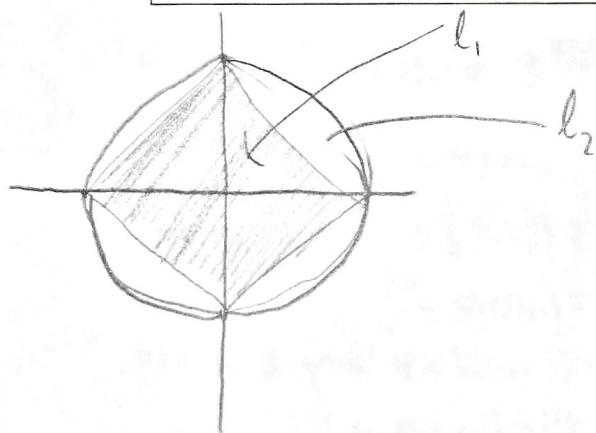
Show that the optimization problem

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

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 ℓ_1 regularization (adding a $\lambda \|\mathbf{x}\|_1$ term to the objective) will give sparser solutions than using ℓ_2 regularization for suitably large λ .



$$\text{minimize } f(\mathbf{x}) \quad \text{subj. to } \|\mathbf{x}\|_p \leq k \quad \text{is the same as}$$

$$\inf_{\mathbf{x}} \sup_{\lambda \geq 0} L(\mathbf{x}, \lambda) = \inf_{\mathbf{x}} \sup_{\lambda \geq 0} (f(\mathbf{x}) + \lambda(\|\mathbf{x}\|_p - k)) \quad \text{flipping inf and sup, we find}$$

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

we know that minimizing $f(\mathbf{x}) + \lambda(\|\mathbf{x}\|_p - k)$ for \mathbf{x} is the same as $\min_{\mathbf{x}} f(\mathbf{x}) + \lambda\|\mathbf{x}\|_p$, as $-\lambda k$ doesn't depend on \mathbf{x} , optimizing for the value of \mathbf{x} is the same as minimizing: $f(\mathbf{x}) + \lambda\|\mathbf{x}\|_p$ for $\lambda \geq 0$.

Considering the plot and our result, the ℓ_1 regularization is essentially projecting the optimal solution onto an ℓ_1 norm-ball. The ℓ_1 norm ball compared to the ℓ_2 norm ball has straight rather than curved edges, thus making the likelihood of landing on an edge and not the face much much larger than the ℓ_2 ball's curved edges. This is because when we rotate the ℓ_2 ball, nothing happens. This isn't the case for the ℓ_1 ball. Hence, the ℓ_1 penalty wants more weights to be 0, rather than the ℓ_2 ball (which helps with simplification of our model, interpretation, and overfitting).

Extra Credit (Lasso) Show that placing an equal zero-mean Laplace prior on each element of the weights θ of a model is equivalent to ℓ_1 regularization in the Maximum-a-Posteriori estimate

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

Note the form of the Laplace distribution is

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

where μ is the location parameter and $b > 0$ controls the variance. Draw (by hand) and compare 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 ℓ_2 regularization).

This problem is equivalent to maximizing $\log \mathbb{P}(\theta|\mathcal{D})$. From the monotonicity of $\log(x)$. Hence, we have

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

$\mathbb{P}(\mathcal{D})$ doesn't depend on θ so we can drop it. Now we have

$$\text{minimize } -\log \mathbb{P}(\mathcal{D}|\theta) - \log \mathbb{P}(\theta). \quad \text{With the prior } \theta_i \sim \text{Lap}(0, b), \text{ we have}$$

$$-\log \mathbb{P}(\theta) = -\log \prod_i e^{-\frac{|\theta_i|}{b}} + z \quad z \text{ is a constant}$$

$$= \frac{1}{b} \sum_i |\theta_i| + z$$

$$= \lambda \|\theta\|_1 + z$$

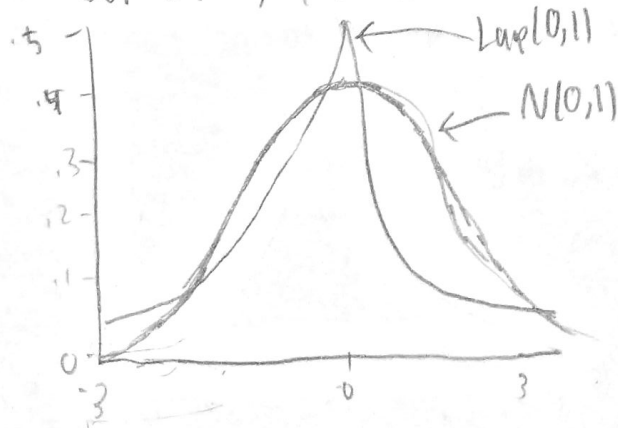
$$\text{where } \lambda = 1/b$$

So, our original problem is the same as

$$\text{minimize: } -\log \mathbb{P}(\mathcal{D}|\theta) + \lambda \|\theta\|_1$$

which is an ℓ_1 regularized maximum likelihood estimate, as we wished.

Our density plots look like



From this plot, we see $\text{Lap}(0,1)$ has much more mass at $x=0$ than the Gaussian. Hence, when we use the Laplace prior instead of a Gaussian prior on the weights, the weights will be pushed to zero more, causing sparser solutions.