

Question1:

Q1

a) Pdf of normal distribution  $P(x) = \frac{1}{\sqrt{2\pi}\sigma^2} e^{-\frac{1}{2\sigma^2}(x-\mu)^2}$  . . . . . ①

$$\theta_{MAP} = \underset{\theta \in (-\infty, \infty)}{\operatorname{argmax}} \{P(D|\theta) \cdot P(\theta)\}$$

$$\begin{aligned} \therefore \ln P(D|\theta) &= \ln \prod_{i=1}^n P(x_i|\theta) \quad [\because D \rightarrow \text{iid}] \\ &= n \ln \frac{1}{\sqrt{2\pi}} + n \ln \frac{1}{\sigma_0} - \frac{\sum_{i=1}^n (x_i - \theta)^2}{2\sigma_0^2} \end{aligned}$$

Since  $\theta$  chosen from normal distribution  $N(\mu, \sigma^2)$ .

$$\begin{aligned} \Rightarrow \ln P(\theta) &= \ln L(\theta) \\ &= n \ln \frac{1}{\sqrt{2\pi}} + n \ln \frac{1}{\sigma} - \frac{(\theta - \mu)^2}{2\sigma^2} \ln e \end{aligned}$$

Thus,  $\ln P(\theta|D) \propto \ln P(D|\theta) + \ln P(\theta)$

$$\Rightarrow \ln P(\theta|D) = n \ln \frac{1}{\sqrt{2\pi}} + n \ln \frac{1}{\sigma_0} - \frac{\sum_{i=1}^n (x_i - \theta)^2}{\sigma_0^2} - \frac{\theta - \mu}{\sigma^2} \quad \text{②}$$

Set partial derivative of  $\theta = 0$

$$\frac{\partial}{\partial \theta} \ln P(\theta|D) = 0$$



$$\theta_{MAP} = \frac{\sigma^2 \sum_{i=1}^n x_i + \mu \sigma_0}{\sigma_0 + n \sigma^2}$$

$$\frac{\sum_{i=1}^n (x_i - \theta)}{\sigma_0^2} - \frac{\theta - \mu}{\sigma^2} = 0, \text{ thus } \quad \text{③}$$

Alkemy

Q2

b)

pdf of normal distribution:  $P(x) = \frac{1}{\sqrt{2\pi}\sigma} e^{-\frac{1}{2\sigma^2}(x-\mu)^2}$

$\theta \text{ MAP} = \underset{\theta \in (-\infty, \infty)}{\operatorname{argmax}} \{P(D|\theta) \cdot P(\theta)\}$

$\ln P(D|\theta) = \ln \prod_{i=1}^n P(x_i|\theta) \quad [\because D \rightarrow \text{iid}]$   
 $= n \ln \frac{1}{\sqrt{2\pi}} + n \ln \frac{1}{\sigma} - \frac{\sum_{i=1}^n (x_i - \theta)^2}{2\sigma^2}$

$\theta$  is chosen from  $\mathcal{L}(\mu, b)$ ,  $\mu=0$ , so

$\ln P(\theta) = \ln \left( \frac{1}{\sqrt{2b}} \exp\left(-\frac{|\theta-0|}{b}\right) \right)$   
 $= \ln \frac{1}{\sqrt{2b}} + \frac{|\theta|}{b} \ln e$  (2)

$\ln P(D|\theta) \propto \ln P(D|\theta) + \ln P(\theta) \rightarrow$   
 $\ln P(D|\theta) = n \ln \frac{1}{\sqrt{2\pi}} + n \ln \frac{1}{\sigma} - \frac{\sum_{i=1}^n (x_i - \theta)^2}{2\sigma^2} + \left( \ln \frac{1}{\sqrt{2b}} + \frac{|\theta|}{b} \ln e \right)$

To maximize  $\ln P(D|\theta)$ : derivative.

$\frac{\partial}{\partial \theta} \ln P(D|\theta) = \frac{\sum_{i=1}^n (x_i - \theta)}{\sigma^2} - \frac{1}{2b^2} \frac{\partial}{\partial \theta} |\theta|$  (1)

$\theta$  is a piece wise derivative. It will have a derivative of 1 for  $x > 0$  and a derivative of -1 for  $x < 0$ . Since we don't get a simple value of  $\theta$  from the above expression. so no closed form solution.

Iterative approach typed on PDF.

- b) Iterative approach typed on PDF. By using a gradient descent approach we can solve this problem. We can start with some random value for theta > 1 and take the derivative of theta in Equation 1. After that, we can get the gradient of the MAP function. Our goal is to maximize the function. To achieve that we need to keep decreasing theta by small weight and see if the value of the function in equation 2 increases. Once the gradient for this function approaching 0 and the function reaches it's maxing value. We repeat this process till the time two consecutive theta values become almost equal and the gradient of the function becomes 0. Similarly, we can repeat the same process for a

value of  $\theta < 0$  and repeat the same process until the time the function converges. So, I think there might be 2 values of  $\theta$  for which the function converges.

The pdf of multivariate gaussian variables:

$$P(x, \theta, \Sigma) = \frac{1}{\sqrt{2\pi}|\Sigma|} e^{-\frac{1}{2}(x-\theta)^T \Sigma^{-1}(x-\theta)}$$

Now we need the MAP estimate of  $\theta$  parameter, which can be estimated as

$$\theta_{MAP} = \underset{\theta \in (-\infty, \infty)}{\operatorname{argmax}} \{P(D|\theta) \cdot P(\theta)\}$$

$$\begin{aligned} \ln P(D|\theta) &= \ln \left[ \prod_{i=1}^n P(x_i|\theta) \right] \quad [D \rightarrow \text{ind}] \\ &= \ln \left( \frac{1}{\sqrt{2\pi}|\Sigma|} e^{-\frac{1}{2} \sum_{i=1}^n (x_i - \theta)^T \Sigma^{-1} (x_i - \theta)} \right) \\ &= n \ln \left( \frac{1}{\sqrt{2\pi}} \right) + n \ln \left( \frac{1}{\sqrt{|\Sigma|}} \right) - \frac{1}{2} \sum_{i=1}^n (x_i - \theta)^T \Sigma^{-1} (x_i - \theta) \end{aligned}$$

Gaussian  $N(\mu=0, \Sigma=\sigma^2 I)$ :

$$\begin{aligned} \ln P(\theta) &= \ln \left( \frac{1}{\sqrt{2\pi}|\Sigma|} e^{-\frac{1}{2} \theta^T \Sigma^{-1} \theta} \right) \\ &= \ln \frac{1}{\sqrt{2\pi}} + \ln \frac{1}{|\Sigma|} - \frac{1}{2} \theta^T \Sigma^{-1} \theta \end{aligned}$$

$$\begin{aligned} \ln P(\theta|D) &\propto \ln P(D|\theta) + \ln P(\theta) \\ \Rightarrow \ln P(\theta|D) &= n \ln \left( \frac{1}{\sqrt{2\pi}} \right) + n \ln \left( \frac{1}{\sqrt{|\Sigma|}} \right) - \frac{1}{2} \sum_{i=1}^n (x_i - \theta)^T \Sigma^{-1} (x_i - \theta) \\ &\quad + \ln \frac{1}{\sqrt{2\pi}} + \ln \frac{1}{|\Sigma|} - \frac{1}{2} \theta^T \Sigma^{-1} \theta \end{aligned}$$

To maximize  $\ln P(\theta|D)$  we take the derivative.

$$\begin{aligned} \frac{\partial}{\partial \theta} \ln P(\theta|D) &= \frac{1}{2} \sum_{i=1}^n \frac{\partial}{\partial \theta} (x_i - \theta)^T \Sigma^{-1} (x_i - \theta) - \frac{1}{2} \frac{\partial}{\partial \theta} \theta^T \Sigma^{-1} \theta \\ &= \frac{1}{2} \sum_{i=1}^n 2(x_i - \theta) - \frac{1}{2} 2\theta^T \Sigma^{-1} = 0 \quad \text{and } \Sigma = \sigma^2 I \end{aligned}$$

Now we set derivative obtained to 0 and we get.

$$\begin{aligned} \frac{\partial}{\partial \theta} \ln P(\theta|D) &= 0 \\ \therefore \sum_{i=1}^n (x_i - \theta) - \theta &= 0 \\ \Rightarrow \sum_{i=1}^n x_i - n\theta - \theta &= 0 \\ \theta_{MAP} &= \frac{\sum_{i=1}^n x_i}{n+1} \end{aligned}$$

Question2:

- a) By increasing the number of selected features to 385, we raised an error of "numpy.linalg.LinAlgError: Singular matrix". The error might be due to the issue that we used the formula  $w = (X^T X)^{-1} X^T y$  in our implementation. However, some matrices from the dataset might have a determinant of zero because of identicalness, similarity,

or they are linearly dependent features. This implies that it is a singular matrix; thus, the inverse of the matrix does not exist.

To remedy this problem, we can use singular value decomposition, or just adding the regularization weight decay.

- c) The model trained by Ridge Regression have better result when comparing to the FSLinearRegression from (a). We had average error for FSLR(385features) around 11.53 +-0.37 and average error for RLR(lambda = 0.01) around 8.51 +- 0.19. The reason for RLR error is lower than FSLR error might be the idea of Ridge Regression regularized the model, so the model generally performs better with the test data (help deal with overfitting). Another assumption might be the default lambda of FSLR given code= 0.5(to avoid singular matrix error) is just way too big.
- f) Stochastic gradient descent (SGD) typically reaches convergence much faster than batch gradient descent since it updates weight more frequently. Unlike the batch gradient descent, which computes the gradient using the whole dataset, because the SGD, also known as incremental gradient descent, tries to find minimums or maximums by iteration from a single randomly picked training example, the error is typically noisier than in gradient descent.

In terms of the number of times, the entire training set is processed. SGD processed the entire training set for epochs times and epochs=1000 for my implementation, according to Algorithm 3 from Note section 6.3. BGD processed the entire training set for epochs \* maxIteration times and epochs=1000, maxIteration = 10e5 for my implementation, according to Algorithm 2 and Algorithm 1 from Note section 6.3.

```
Results:
Best parameters for StochasticGradientDescent: {'regwgt': 0.01}
Average error for StochasticGradientDescent: 9.269436681646129 +- 0.0
Standard error for StochasticGradientDescent: 0.0
Average runtime for StochasticGradientDescent: 68.80631349999749

Best parameters for BatchGradientDescent: {'regwgt': 0.01}
Average error for BatchGradientDescent: 8.926717568337368 +- 0.0
Standard error for BatchGradientDescent: 0.0
Average runtime for BatchGradientDescent: 384.89999340000213

zeeman@Zeeman-PC: /mnt/c/Users/akazf/Desktop/a2barebones$
```

Above is my runtimes result for the SGD and BGD. Therefore error versus epochs: SGD(9.27/1000) and BGD (8.93/1000). For error versus runtime: SGD(9.27/68.81) and BGD (8.93/384.90).