

Homework 5 *

Ishan Upadhyaya

February 22, 2016

1 Bias Variance Tradeoff

1.1 Part A

The closed form solution for the given optimization problem can be written as

$$\hat{\beta}_\lambda = (X^T X + n\lambda I_p)^{-1} X^T \mathbf{y} ; \text{ Where } \mathbf{y} \text{ is the output vector}$$

Using the given linear model, $y_i = x_i^T \beta^* + \epsilon_i$, we can write the above equation as

$$\hat{\beta}_\lambda = \underbrace{(X^T X + n\lambda I_p)^{-1} X^T}_{A_\lambda} (X \beta^* + \mathcal{E}) = A_\lambda X \beta^* + A_\lambda \mathcal{E}$$

Here, \mathcal{E} is a vector given by $\mathcal{E} = [\epsilon_1 \quad \epsilon_2 \quad \cdots \quad \epsilon_n]^T$ and ϵ_i is the noise for i^{th} data point

Using properties of affine transformation of Gaussian random vectors, we can say that $\hat{\beta}_\lambda$ is a Gaussian random vector with mean and variance given below

$$\mu_{\hat{\beta}_\lambda} = A_\lambda X \beta^* \text{ and } C_{\hat{\beta}_\lambda} = A_\lambda^T C_{\mathcal{E}} A_\lambda$$

Here $C_{\mathcal{E}}$ is the covariance matrix of \mathcal{E}

1.2 Part B

Using linearity property of Expectation, we get that $\mathbb{E}[x^T \hat{\beta}_\lambda] = x^T \mathbb{E}[\hat{\beta}_\lambda]$. Using the above expression for given problem we get

$$\begin{aligned} \mathbb{E}[x^T \hat{\beta}_\lambda] - x^T \beta^* &= x^T \{\mathbb{E}[\hat{\beta}_\lambda] - \beta^*\} = x^T \{A_\lambda X \beta^* - \beta^*\} \\ &= x^T \underbrace{\{A_\lambda X - I_p\}}_Q \beta^* = b_\lambda \end{aligned}$$

$$\begin{aligned} Q &= A_\lambda X - I_p = (X^T X + n\lambda I_p)^{-1} X^T X - I_p \\ &= (I_p - n\lambda (X^T X)^{-1})^{-1} - I_p = I_p - n\lambda (X^T X + n\lambda I_p)^{-1} - I_p \\ &= -n\lambda (X^T X + n\lambda I_p)^{-1} \end{aligned}$$

*COMP SCI 260 - Machine Learning Algorithms ; Fall 2015

1.3 Part C

The expression $\mathbb{E}[\{x^T \hat{\beta}_\lambda - \mathbb{E}[x^T \hat{\beta}_\lambda]\}^2]$ can be written as $\underbrace{\mathbb{E}[\{x^T \{\hat{\beta}_\lambda - \mathbb{E}[\hat{\beta}_\lambda]\}\}^2]}_P$. Simplyfying P further we get

$$\begin{aligned} P &= \mathbb{E}[\{x^T (A_\lambda X \beta^* + A_\lambda \mathcal{E} - A_\lambda X \beta^*)\}^2] \\ &= \mathbb{E}[\{x^T A_\lambda \mathcal{E}\}^2] = \mathbb{E}[x^T A_\lambda \mathcal{E} \mathcal{E}^T A_\lambda^T x] = x^T A_\lambda \mathbb{E}[\mathcal{E} \mathcal{E}^T] A_\lambda^T x = x^T A_\lambda C_\mathcal{E} A_\lambda^T x \end{aligned}$$

1.4 Part D

Consider behavior of P and A_λ as $\lambda \rightarrow 0$ and $\lambda \rightarrow \infty$.

$$\begin{aligned} \lim_{\lambda \rightarrow 0} P &= 0 \text{ and } \lim_{\lambda \rightarrow \infty} P = -n I_p \\ \lim_{\lambda \rightarrow 0} A_\lambda &= X^{-1} \text{ and } \lim_{\lambda \rightarrow \infty} A_\lambda = \frac{1}{n\lambda} X^T = \mathbf{0}_{p \times n} \end{aligned}$$

Therefore we get

$$\begin{aligned} \lim_{\lambda \rightarrow 0} b_\lambda &= 0 \text{ and } \lim_{\lambda \rightarrow \infty} b_\lambda = -x^T n I_p \beta^* = -n x^T \beta^* \\ \lim_{\lambda \rightarrow 0} P &= x^T X^{-1} C_\mathcal{E} (X^{-1})^T x \text{ and } \lim_{\lambda \rightarrow \infty} P = 0 \end{aligned}$$

In words, the bias becomes zeros as $\lambda = 0$ while the variance term goes to zero as $\lambda \rightarrow \infty$

2 Kernilized Perceptron

2.1 Part A

We show that \mathbf{w} is a linear weighted sum of feature vectors using an argument similar to that of Induction. We know that \mathbf{w} is intialized as 0, thus the first time that \mathbf{w} is updated, we have $\mathbf{w} = y_k \varphi(x_k)$, where 'k' is the first sample for which \mathbf{w} misclassifies. Consequently, we can write this \mathbf{w} as $\sum_{i=1}^n \alpha_i \varphi(x_i)$, where $\alpha_i = 0 \forall i \neq k$ and $\alpha_k = y_k$. Now, suppose the immediately subsequent update to \mathbf{w} occurs for j^{th} sample point, then we write $\mathbf{w} = \sum_{i=1}^n \alpha_i \varphi(x_i) + y_j \varphi(x_j) = \sum_{i=1}^n \hat{\alpha}_i \varphi(x_i)$ such that $\hat{\alpha}_i = \alpha_i \forall i \neq j$ and $\hat{\alpha}_i = \alpha_i + y_j$ for $i = k$. Thus we have shown that \mathbf{w} is a linear sum of weighted feature vectors after two updates. Thus we claim, that it is a weighted sum after 'm' updates. Now consider the $(m+1)^{th}$ update at say p^{th} data point. Using, the same logic as above, we can write $\mathbf{w} = \sum_{i=1}^n \alpha_i \varphi(x_i) + y_p \varphi(x_p) = \sum_{i=1}^n \hat{\alpha}_i \varphi(x_i)$ such that $\hat{\alpha}_i = \alpha_i \forall i \neq p$ and $\hat{\alpha}_i = \alpha_i + y_p$ for $i = p$. Hence, we have proved that if the weight vector is initialized as a zero vector, the we can write it as a linear weighted sum of feature vectors, i.e. $\mathbf{w} = \sum_{i=1}^n \alpha_i \varphi(x_i)$.

2.2 Part B

The prediction of perceptron is given by $\hat{y}_k = \text{sign}(\mathbf{w}^T \varphi(x_k))$. Expanding this equation further we get,

$$\begin{aligned}\hat{y}_k &= \text{sign}(\mathbf{w}^T \varphi(x_k)) \\ &= \text{sign}\left\{\left(\sum_{i=1}^n \alpha_i \varphi(x_i)\right)^T \varphi(x_k)\right\} = \text{sign}\left\{\sum_{i=1}^n \alpha_i \underbrace{\varphi(x_i)^T \varphi(x_k)}_{k(x_i, x_k)}\right\} \\ &= \text{sign} \sum_{i=1}^n \alpha_i k(x_i, x_k)\end{aligned}$$

Hence, the prediction rule for perceptron can be written in terms of a kernel function.

2.3 Part C

From part A we know that $\mathbf{w} = \sum_{i=1}^n \alpha_i \varphi(x_i)$. During the training phase, for any update we have $\mathbf{w} \leftarrow \mathbf{w} + y_k \varphi(x_k) = \sum_{i=1}^n \alpha_i \varphi(x_i) + y_k \varphi(x_k)$. This equation can be modified as $\alpha_k \leftarrow \alpha_k + y_k$ to bring about the same change in the weight vector. Also notice that the prediction of a new data point $\varphi(x_k)$ is given by $\text{sign} \sum_{i=1}^n \alpha_i k(x_i, x_k)$, which is independent of the explicit form of weight vector \mathbf{w} . Thus, during the training time, we only need access to α_i 's and x_i 's to update the weight vector.

During testing time, we are just predicting the label of a data point x_k using the $\text{sign} \sum_{i=1}^n \alpha_i k(x_i, x_k)$, which just requires access to x_i 's and α_i 's to make prediction. Thus proving that we don't need access to explicit form of \mathbf{w} during testing time as well.

The algorithm for kernelized matrix is given as follows.

1. Initialize $\alpha_i = 0 \forall i$.
2. Set counter $k=0$
3. For current data point 'k', predict label $\hat{y}_k = \text{sign} \sum_{i=1}^n \alpha_i k(x_i, x_k)$
4. If $\hat{y}_k \neq y_k$ update α_k as $\alpha_k \leftarrow \alpha_k + y_k$
5. $k \leftarrow k+1$
6. Iterate till convergence of \mathbf{w}

3 Kernels

Let $\mathbf{v} \in \mathbb{R}^n$

3.1 Part A

For matrix $K_3 = a_1 K_1 + a_2 K_2$, the term $\mathbf{v}^T K_3 \mathbf{v}$ is given by

$$\begin{aligned}\mathbf{v}^T K_3 \mathbf{v} &= a_1 \mathbf{v}^T K_1 \mathbf{v} + a_2 \mathbf{v}^T K_2 \mathbf{v} \geq 0 \\ \because a_1, a_2 &\geq 0 \text{ and } \mathbf{v}^T K_1 \mathbf{v}, \mathbf{v}^T K_2 \mathbf{v} \geq 0 \text{ As } k_1 \text{ and } k_2 \text{ are positive semi definite}\end{aligned}$$

$\therefore K_3$ is positive semidefinite

3.2 Part B

The matrix generated by the function $k_4(x_1, x_2)$ is shown below

$$K_4 = \begin{bmatrix} f(x_1)f(x_1) & f(x_1)f(x_2) & \cdots & f(x_1)f(x_n) \\ f(x_2)f(x_1) & f(x_2)f(x_2) & \cdots & f(x_2)f(x_n) \\ \vdots & \vdots & \ddots & \vdots \\ f(x_n)f(x_1) & f(x_n)f(x_2) & \cdots & f(x_n)f(x_n) \end{bmatrix} = \begin{bmatrix} f(x_1) & f(x_2) & \cdots & f(x_n) \end{bmatrix} \begin{bmatrix} f(x_1) \\ f(x_2) \\ \vdots \\ f(x_n) \end{bmatrix} = \mathbf{f}\mathbf{f}^T$$

Therefore, $\mathbf{v}^T K_4 \mathbf{v}$ is given by

$$\mathbf{v}^T K_4 \mathbf{v} = \mathbf{v}^T \mathbf{f} \mathbf{f}^T \mathbf{v} = \{\mathbf{f}^T \mathbf{v}\}^2 \geq 0$$

Therefore, K_4 is positive semi definite.

3.3 Part C

Given that $k_5 = k_1(x_1, x_2)k_2(x_1, x_2)$, where k_1 and k_2 are kernel functions, we can alternatively express k_1 and k_2 as

$$k_1(x_1, x_2) = \varphi_1(x_1)^T \varphi_1(x_2) = \sum_{i=1}^n \varphi_1(x_1)_i * \varphi_1(x_2)_i$$

$$k_2(x_1, x_2) = \varphi_2(x_1)^T \varphi_2(x_2) = \sum_{i=1}^m \varphi_2(x_1)_i * \varphi_2(x_2)_i$$

Using the above equations we can write k_5 as

$$\begin{aligned} k_5(x_1, x_2) &= \left(\sum_{i=1}^n \varphi_1(x_1)_i * \varphi_1(x_2)_i \right) \left(\sum_{j=1}^m \varphi_2(x_1)_j * \varphi_2(x_2)_j \right) \\ &= \sum_{i=1}^n \sum_{j=1}^m (\varphi_1(x_1)_i \varphi_2(x_1)_j) * (\varphi_1(x_2)_i \varphi_2(x_2)_j) \\ &= \sum_{i=1}^n \sum_{j=1}^m \varphi_3(x_1)_{i,j} * \varphi_3(x_2)_{i,j} \\ &= \varphi_3(x_1)^T \varphi_3(x_2) ; \text{ here } \varphi_3(\cdot) \text{ is a } (nm) \times 1 \text{ dimensional vector} \end{aligned}$$

Thus we can write $k_5(\cdot, \cdot)$ as an innerproduct of two vectors i.e. $k_5(x_1, x_2) = \varphi_3(x_1)^T \varphi_3(x_2)$. Now consider $\mathbf{v}^T K_5 \mathbf{v}$

$$\begin{aligned} \mathbf{v}^T K_5 \mathbf{v} &= \sum_{i,j} \mathbf{v}_i < \varphi_3(x_i), \varphi_3(x_j) > \mathbf{v}_j \\ &= < \sum_i \mathbf{v}_i \varphi_3(x_i), \sum_j \varphi_3(x_j) \mathbf{v}_j > = \left\| \sum_j \varphi_3(x_j) \mathbf{v}_j \right\|_2^2 \\ &>> 0 \end{aligned}$$

Hence, K_5 is positive semidefinite.

4 Soft Margin Hyperplanes

4.1 Part A

Lagrangian of the modified optimization problem is given by

$$\begin{aligned}\mathcal{L}_p &= C \sum_n \xi_n + \frac{1}{2} \mathbf{w}^T \mathbf{w} - \sum_n \lambda_n \xi_n + \sum_n \alpha_n [1 - y_n (\mathbf{w}^T \phi(x_n) + b) - \xi_n] \\ \frac{\partial \mathcal{L}_p}{\partial \xi_n} &= p C \xi_n^{p-1} - \lambda_n - \alpha_n = 0 \\ \frac{\partial \mathcal{L}_p}{\partial \mathbf{w}} &= \mathbf{w} - \sum_n y_n \alpha_n \phi(x_n) = 0 \\ \frac{\partial \mathcal{L}_p}{\partial b} &= \sum_n \alpha_n y_n = 0\end{aligned}$$

Substituting the conditions back into lagrangian we get

$$\begin{aligned}\mathcal{G}_p &= \sum_n \{C \xi_n^{p-1} - \lambda_n - \alpha_n\} \xi_n + \frac{1}{2} \left\| \sum_n y_n \alpha_n \phi(x_n) \right\|_2^2 + \sum_n \alpha_n \\ &\quad - \sum_n y_n \alpha_n b - \sum_n y_n \alpha_n \left(\sum_m y_m \alpha_m \phi(x_m) \right)^T \phi(x_n) \\ \mathcal{G}_p &= \sum_n (1-p) C \xi_n^p + \sum_n \alpha_n - \frac{1}{2} \sum_{n,m} y_m y_n \alpha_n \alpha_m \phi(x_n)^T \phi(x_m) \\ &= \underbrace{\frac{1-p}{C^{\frac{1}{p-1}}} \sum_n \left(\frac{\lambda_n + \alpha_n}{p} \right)^{\frac{p}{p-1}}}_{S_p} + \sum_n \alpha_n - \frac{1}{2} \sum_{n,m} y_m y_n \alpha_n \alpha_m \phi(x_n)^T \phi(x_m)\end{aligned}$$

Hence, the dual formation of the problem is given as

$$\begin{aligned}\boldsymbol{\alpha} &= \max_{\boldsymbol{\alpha}} \mathcal{G}_p \\ \alpha_n, \lambda_n &\geq 0 \\ \sum_n y_n \alpha_n &= 0\end{aligned}$$

4.2 Part B

The general formulation is more complex than the one discussed in class because of the additional S_p term. Also, the number of variables involved is double than that of the simpler formulation, as we have the additional λ_n terms.

5 Programming

5.1 Data Preprocessing

Mean and standard deviation before pre-processing

1. 3rd feature: Mean = 2.5260, Standard Deviation = 1.0856

2. 10th feature: Mean = 2.5526, Standard Deviation = 1.1222

Mean and standard deviation after pre-processing

1. 3rd feature: Mean = -0.013, Standard Deviation = 1.0099
2. 10th feature: Mean = 0.0228, Standard Deviation = 0.9979

5.2 Cross validation for Linear SVM

5.2.1 A

Variation of Cross validation accuracy and training time are shown below. After a certain value of C, accuracy flattens out as can be observed in the graph, while the training time generally increases with value of C.

5.2.2 B

As it can be observed from the graphs, maximum cross-validation accuracy is obtained at $C = 4^{-3}$ with an accuracy of 80.9%

5.2.3 C

Test accuracy if the model is trained using $C = 4^{-3}$ is 84.64%.

5.3 Linear SVM using LibSVM

Plots for variation of cross-validation accuracy and training time with C are shown below

As can be observed, the accuracy plot is similar to the one we obtained using quadprog, while the time plot follows a similar trend of increasing training time with value of C. In general, training time of LibSVM is slightly more than the one observed from our implementation.

The maximum cross-validation accuracy achieved was 80.4% with value of C as 4^{-5} . This value is slightly different from the one obtained using our implementation. This might be due to differing optimization algorithm used by quadprog and LibSVM or because of different forms (primal and dual) being optimized by them or a combination of the above two reasons.

5.4 SVM using Kernels

Cross validation and training times for Polynomial kernel function are shown below for different degrees. Please note that the X-axis is in log scale.

Cross validation and training times for RBF kernel function are shown below for different values of gamma. Please note that the X-axis is in log scale.

As can be observed, the maximum cross-validation accuracy is obtained with RBF kernel, with $C = 4$ and $\gamma = 4^{-3}$. The corresponding accuracy = 88%.

The test accuracy for these parameters is 90.4368%.