

Instructions

Submission: Assignment submission will be via courses.usciden.net. By the submission date, there will be a folder named 'Theory Assignment 3' set up in which you can submit your files. Please be sure to follow all directions outlined here.

You can submit multiple times, but only *the last submission* counts. That means if you finish some problems and want to submit something first and update later when you finish, that's fine. In fact you are encouraged to do this: that way, if you forget to finish the homework on time or something happens (remember Murphy's Law), you still get credit for whatever you have turned in.

Problem sets must be typewritten or neatly handwritten when submitted. In both cases, your submission must be a single PDF. It is strongly recommended that you typeset with \LaTeX . There are many free integrated \LaTeX editors that are convenient to use (e.g. [Overleaf](#), [ShareLaTeX](#)). Choose the one(s) you like the most. This tutorial [Getting to Grips with LaTeX](#) is a good start if you do not know how to use \LaTeX yet.

Please also follow the rules below:

- The file should be named as `firstname_lastname_USCID.pdf` e.g., `Don_Quijote_de_la_Mancha_8675309045.pdf`.
- Do not have any spaces in your file name when uploading it.
- Please include your name and USCID in the header of your report as well.

Collaboration: You may discuss with your classmates. However, you need to write your own solutions and submit separately. Also in your report, you need to list with whom you have discussed for each problem. Please consult the syllabus for what is and is not acceptable collaboration. Review the rules on academic conduct in the syllabus: a single instance of plagiarism can adversely affect you significantly more than you could stand to gain.

Notes on notation:

- Unless stated otherwise, scalars are denoted by small letter in normal font, vectors are denoted by small letters in bold font and matrices are denoted by capital letters in bold font.
- $\|\cdot\|$ means L2-norm unless specified otherwise i.e. $\|\cdot\| = \|\cdot\|_2$

Problem 1 Principle Component Analysis

(25 points)

In this problem, we use proof by induction to show that the M -th principle component corresponds to the M -th eigenvector of $X^T X$ sorted by the eigenvalue from largest to smallest. Here X is the centered data matrix and we denote the sorted eigenvalues as $\lambda_1 \geq \lambda_2 \geq \dots \geq \lambda_d$. In the lecture, the results was proven for $M = 1$. Now suppose the result holds for a value M , and you are going to show that it holds for $M + 1$. Note that the $M + 1$ principle component corresponds to the solution of the following optimization problem:

$$\max_v v^T X^T X v \quad (1)$$

$$\text{s.t. } \|v\|_2 = 1 \quad (2)$$

$$v^T v_i = 0, i = 1, \dots, M \quad (3)$$

where v_i is the i -th principle component. Write down the Lagrangian of the optimization problem above, and show that the solution v_{M+1} is an eigenvector of $X^T X$. Then show that the quantity in (1) is maximized when the v_{M+1} is the eigenvector with eigenvalue λ_{M+1} .

Ans: The Lagrangian of 1 is:

$$\mathcal{L}(v; \eta_1, \dots, \eta_M, \eta) = v^T X^T X v + \eta(1 - v^T v) + \sum_{i=1}^M \eta_i v^T v_i \quad (4 \text{ points})$$

Taking derivative w.r.t v and set it to 0, we get:

$$2X^T X v - 2\eta v + \sum_{i=1}^M \eta_i v_i = 0 \quad (4 \text{ points})$$

Left Multiply by v_j on both sides, we get:

$$2v_j^T X^T X v - 2\eta v_j^T v + \sum_{i=1}^M \eta_i v_j^T v_i = 0 \quad (3 \text{ points})$$

$$2v_j^T X^T X v + \eta_j = 0 \quad (\text{by } v_1, \dots, v_M, v \text{ are orthonormal}) \quad (3 \text{ points})$$

$$2\lambda_j v_j^T v + \eta_j = 0 \quad (v_j \text{ is eigenvector of } X^T X) \quad (3 \text{ points})$$

$$\eta_j = 0 \quad (1 \text{ points})$$

Therefore, $\eta_j = 0, j = 1, \dots, M$, and we have: $X^T X v = \eta v$ (2 points)

Thus v is an eigenvector of $X^T X$, and the value of equation 1 is $v^T X^T X v = \lambda v^T v = \lambda$, where λ is the eigenvalue of v . (2 points)

Therefore, to optimize the value, we should let v to be the eigenvector with eigenvalue λ_{M+1} , since v is perpendicular to the first M principle components, which are eigenvectors with eigenvalues $\lambda_1, \dots, \lambda_M$. (3 points)

Problem 2 Support Vector Regression

(30 points)

In this problem, we derive an extension of support vector machine to regression problem, called Support Vector Regression (SVR). Define the regressor $f(x) = \mathbf{w}^T \boldsymbol{\phi}(x) + b$, and given a dataset $\{(x_n, y_n)\}_{n=1}^N, y_n \in \mathbb{R}$. Intuitively, we want to find a regressor that has small weight \mathbf{w} and also ensure small approximation error to $\{(x_n, y_n)\}_{n=1}^N$. The intuition can be formulated as the following optimization problem:

$$\begin{aligned} \min_{\mathbf{w}, b} \quad & \frac{1}{2} \|\mathbf{w}\|_2^2 \\ \text{s.t.} \quad & |\mathbf{w}^T \boldsymbol{\phi}(x_n) + b - y_n| \leq \epsilon \end{aligned}$$

For an arbitrary dataset, the ϵ -close constraint may not be feasible, Therefore, we optimize the “soft” version of the loss above:

$$\min_{\mathbf{w}, b} \quad \frac{1}{2} \|\mathbf{w}\|_2^2 + C \sum_{n=1}^N E_\epsilon(y_n - f(x_n)) \quad (4)$$

E_ϵ is the ϵ -insensitive error function which gives zero error if the difference between prediction and ground truth is smaller than ϵ and incurs linear penalty otherwise. It is defined as follow:

$$E_\epsilon(x) = \begin{cases} 0 & |x| \leq \epsilon \\ |x| - \epsilon & |x| > \epsilon \end{cases}$$

Question 1 Reformulate the unconstrained optimization problem in equation 4 as a constraint optimization problem by introducing slack variables for each data points. Hint: For each data point, introduce slack variables $\xi_n \geq 0, \xi'_n \geq 0$ such that $-\epsilon - \xi'_n \leq y_n - f(x_n) \leq \epsilon + \xi_n$. Then replace E_ϵ with ξ_n, ξ'_n . (12 points)

Ans: Minimizing the loss function in equation 4 is equivalent to the following constraint optimization problem:

$$\begin{aligned} \min_{\mathbf{w}, b} \quad & \frac{1}{2} \|\mathbf{w}\|_2^2 + C \sum_{n=1}^N (\xi_n + \xi'_n) & (6 \text{ points}) \\ \text{s.t.} \quad & \epsilon + f(x_n) - y_n + \xi_n \geq 0 & (2 \text{ points}) \\ & \epsilon + y_n - f(x_n) + \xi'_n \geq 0 & (2 \text{ points}) \\ & \xi_n \geq 0, \xi'_n \geq 0 & (2 \text{ points}) \end{aligned}$$

Question 2 Write down the Lagrangian of the constrained optimization derived in Question 1, then minimize the Lagrangian by taking derivative w.r.t $\mathbf{w}, b, \xi_n, \xi'_n$ and set the gradient to 0, and simplify expressions. Hint: there are no b, ξ_n, ξ'_n in the final expressions. (18 points)

Ans: Introduce $\xi = [\xi_1, \dots, \xi_N]^T, \xi' = [\xi'_1, \dots, \xi'_N]^T$. The Lagrangian is as follows:

$$\mathcal{L}(\mathbf{w}, b, \xi, \xi') = \frac{1}{2} \|\mathbf{w}\|_2^2 + C \sum_{n=1}^N (\xi_n + \xi'_n) - \sum_{n=1}^N (b_n \xi_n + b'_n \xi'_n) - \sum_{n=1}^N a_n (\epsilon + \xi_n + f(x_n) - y_n) - \sum_{n=1}^N a'_n (\epsilon + \xi'_n - f(x_n) + y_n) \quad (6 \text{ points})$$

Taking derivative w.r.t $\mathbf{w}, b, \xi_n, \xi'_n$ then gives:

$$\frac{\partial \mathcal{L}}{\partial \mathbf{w}} = 0 \quad \implies \quad \mathbf{w} = \sum_{n=1}^N (a_n - a'_n) \boldsymbol{\phi}(x_n) \quad (3 \text{ points})$$

$$\frac{\partial \mathcal{L}}{\partial b} = 0 \quad \implies \quad \sum_{n=1}^N (a_n - a'_n) = 0 \quad (3 \text{ points})$$

$$\frac{\partial \mathcal{L}}{\partial \xi_n} = 0 \quad \implies \quad a_n + b_n = C \quad (3 \text{ points})$$

$$\frac{\partial \mathcal{L}}{\partial \xi'_n} = 0 \quad \implies \quad a'_n + b'_n = C \quad (3 \text{ points})$$

Problem 3 Support Vector Machine

(25 points)

Consider the dataset consisting of points (x, y) , where x is a real value, and $y \in \{-1, 1\}$ is the class label. There are only three points $(x_1, y_1) = (0, 1)$, $(x_2, y_2) = (\frac{\pi}{2}, -1)$, $(x_3, y_3) = (\pi, 1)$. Let the feature mapping $\phi(x) = [\cos x, \sin x]^T$, corresponding to the kernel function $k(x, y) = \cos(x - y)$.

Question 1 Write down the primal and dual formulations of SVM for this dataset in the transformed two-dimensional feature space based on $\phi(\cdot)$. Note that we assume the data points are separable and set the hyperparameter C to be $+\infty$, which forces all slack variables (ξ) in the primal formulation to be 0 (and thus can be removed from the optimization). (12 points)

Ans: General primal formulation of SVM for separable data is:

$$\begin{aligned} \min_{\mathbf{w}, b} \quad & \frac{1}{2} \|\mathbf{w}\|_2^2 \\ \text{s.t.} \quad & y_n [\mathbf{w}^T \phi(x_n) + b] \geq 1, \forall n \end{aligned}$$

Plugging in the specific dataset gives:

$$\begin{aligned} \min_{w_1, w_2, b} \quad & \frac{1}{2} (w_1^2 + w_2^2) & (3 \text{ points}) \\ \text{s.t.} \quad & w_1 + b \geq 1 & (1 \text{ points}) \\ & -w_2 - b \geq 1 & (1 \text{ points}) \\ & -w_1 + b \geq 1 & (1 \text{ points}) \end{aligned}$$

(Remark: Give 2 points for partial credit if the points obtained in plugging in specific dataset is less than 2, but the general primal formulation is present.)

General dual formulation of SVM is:

$$\begin{aligned} \max_{\alpha} \quad & \sum_n \alpha_n - \frac{1}{2} \sum_{m, n} y_m y_n \alpha_m \alpha_n k(x_m, x_n) \\ \text{s.t.} \quad & \alpha_n \geq 0, \forall n \\ & \sum_n \alpha_n y_n = 0 \end{aligned}$$

Plugging in the specific dataset gives:

$$\begin{aligned} \max_{\alpha_1, \alpha_2, \alpha_3 \geq 0} \quad & \alpha_1 + \alpha_2 + \alpha_3 - \frac{1}{2} \alpha_1^2 - \frac{1}{2} \alpha_2^2 - \frac{1}{2} \alpha_3^2 + \alpha_1 \alpha_3 & (2 \text{ points}) \\ \text{s.t.} \quad & \alpha_2 = \alpha_1 + \alpha_3 & (2 \text{ points}) \\ & \alpha_i \geq 0 & (2 \text{ points}) \end{aligned}$$

(Remark: Give 2 points for partial credit if the points obtained in plugging in specific dataset is less than 2, but the general dual formulation is present.)

Question 2 Next, solve the dual formulation. Based on that, derive the primal solution. (13 points)

Ans: Eliminating the dependence on α_2 using the constraint $\alpha_1 + \alpha_3 = \alpha_2$, we arrive at the objective

$$\max_{\alpha_1, \alpha_2 \geq 0} \quad 2\alpha_1 - \alpha_1^2 + 2\alpha_2 - \alpha_2^2. \quad (3 \text{ points})$$

Clearly, we can maximize over α_1 and α_2 separately, which gives $\alpha_1^* = \alpha_3^* = 1$ (2 points) and thus $\alpha_2^* = 2$ (2 points).

The primal solution can be found by

$$(w_1^*, w_2^*)^T = \sum_{n=1}^3 y_n \alpha_n^* \phi(x_n) = [0, -2]^T \quad (3 \text{ points})$$

$$b^* = y_1 - \mathbf{w}^{*T} \phi(x_1) = 1 \quad (\text{using any example works in this case}) \quad (3 \text{ points})$$

Problem 4 Boosting

(20 points)

Recall the procedure of AdaBoost algorithm described in class:

Algorithm 1: Adaboost

1 **Given:** A training set $\{(x_n, y_n \in \{+1, -1\})\}_{n=1}^N$, and a set of classifier \mathcal{H} , where each $h \in \mathcal{H}$ takes a feature vector as input and outputs $+1$ or -1 .

2 **Goal:** Learn $H(x) = \text{sign}\left(\sum_{t=1}^T \beta_t h_t(x)\right)$

3 **Initialization:** $D_1(n) = \frac{1}{N}, \forall n \in [N]$.

4 **for** $t = 1, 2, \dots, T$ **do**

5 Find $h_t = \arg \min_{h \in \mathcal{H}} \sum_{n: y_n \neq h(x_n)} D_t(n)$.

6 Compute

$$\epsilon_t = \sum_{n: y_n \neq h_t(x_n)} D_t(n) \quad \text{and} \quad \beta_t = \frac{1}{2} \log \frac{1 - \epsilon_t}{\epsilon_t}.$$

7 Compute

$$D_{t+1}(n) = \frac{D_t(n) e^{-\beta_t y_n h_t(x_n)}}{\sum_{n'=1}^N D_t(n') e^{-\beta_t y_{n'} h_t(x_{n'})}}$$

for each $n \in [N]$

Question 1 We discussed in class that AdaBoost minimizes the exponential loss greedily. In particular, Adaboost seeks the optimal β_t that minimizes

$$\epsilon_t (e^{\beta_t} - e^{-\beta_t}) + e^{-\beta_t}$$

where ϵ_t is the weighted classification error of h_t and is fixed. Show that $\beta^* = \frac{1}{2} \ln \left(\frac{1 - \epsilon_t}{\epsilon_t} \right)$ is the minimizer.

(8 points)

Ans: Set the derivative to 0

(1 points):

$$\epsilon_t (e^{\beta_t} + e^{-\beta_t}) - e^{-\beta_t} = 0$$

(3 points)

Multiplying both sides by e^{β_t} and rearranging gives

$$e^{2\beta_t} = \frac{1}{\epsilon_t} - 1$$

(2 points)

Solving for β_t finishes the proof.

(2 points)

(**Remark:** correctly calculate the gradient (3 points). Setting gradient to 0 (1 point) and get the final form of β_t (4 points).)

Question 2 Recall that at round t of AdaBoost, a classifier h_t is obtained and the weighting over the training set is updated from D_t to D_{t+1} . Prove that h_t is only as good as random guessing in terms of classification error weighted by D_{t+1} . That is **(12 points)**

$$\sum_{n:h_t(\mathbf{x}_n) \neq y_n} D_{t+1}(n) = \frac{1}{2}.$$

Hint: you can somehow ignore the denominator of $D_{t+1}(n)$ to simplify calculation.

Ans: By the update algorithm, we have

$$\sum_{n:h_t(\mathbf{x}_n) \neq y_n} D_{t+1}(n) \propto \sum_{n:h_t(\mathbf{x}_n) \neq y_n} D_t(n) e^{\beta_t} = \epsilon_t e^{\beta_t} = \sqrt{\epsilon_t(1 - \epsilon_t)}$$

(3 points)

and similarly

$$\sum_{n:h_t(\mathbf{x}_n) = y_n} D_{t+1}(n) \propto \sum_{n:h_t(\mathbf{x}_n) = y_n} D_t(n) e^{-\beta_t} = (1 - \epsilon_t) e^{-\beta_t} = \sqrt{(1 - \epsilon_t)\epsilon_t}$$

(3 points)

Note that $\sum_{n:h_t(\mathbf{x}_n) \neq y_n} D_{t+1}(n) + \sum_{n:h_t(\mathbf{x}_n) = y_n} D_{t+1}(n) = 1$.

(3 points)

Therefore,

$$\sum_{n:h_t(\mathbf{x}_n) \neq y_n} D_{t+1}(n) = \sum_{n:h_t(\mathbf{x}_n) = y_n} D_{t+1}(n) = \frac{1}{2}$$

(3 points)