

CS5691 - Pattern Recognition and Machine Learning

Jul – Nov, 2023

“Practice Worksheet: SVM, Kernels & Logistic Regression”

November 9, 2023

This worksheet is provided for practice with SVM, Kernels and Logistic Regression only. Students are requested to kindly refer to the Mid Semester and Quiz Worksheets for practice with the other topics included in the End Semester Examination.

1. Consider a SVM Hard Margin problem where the decision boundary is defined by $z(x) = 0$ where $z(x) := w^T x + b$.
 - (i) Derive an expression for the (Euclidean) distance between the margins (margin boundaries).
 - (ii) What is the distance of the origin $(0, 0)$ to the decision boundary $z(x) = 0$?
2. Consider a soft margin SVM where C is the penalty parameter. Explain how the behavior of SVM as a classifier will change as C is increased from a very small value to a very high value.
3. Let $u \in \mathbf{R}^d$ be a point. Let $w \in \mathbf{R}^d, b \in \mathbf{R}$ and the hyperplane given by w, b is $\{x \in \mathbf{R}^d : w^T x + b = 0\}$. Consider the following problem of projection of the point u on to a (hyper)plane given by w, b .

$$\min_{v \in \mathbf{R}^d} \frac{1}{2} \|v - u\|^2$$
$$\text{s.t. } w^T v + b = 0$$

Derive the solution to the above problem via solving the Lagrangian dual (which is an unconstrained quadratic problem, and hence can be easily solved; reviewing the minimax theorem and KKT conditions seen in class may help). Then show that the distance of the point u to the hyperplane given by w, b is $\frac{|w^T u + b|}{\|w\|}$.

4. Let $\{(x_1, y_1), \dots, (x_n, y_n)\}$ be a linearly separable binary classification dataset. Let w^*, b^* be any solution to the problem below:

$$\max_{w \in \mathbf{R}^d, b \in \mathbf{R}} \frac{1}{\|w\|}$$
$$\text{s.t. } y_i(w^T x_i + b) \geq 1$$

Show that $\min_{i \in [n]} y_i(w^T x_i + b) = 1$.

5. Consider a soft margin SVM problem with C set to some constant. Let α^* be the dual solution, and let w^*, b^* be the primal solution. Let the dataset be (x_i, y_i) with i ranging from 1 to n .
 - (i) If $\alpha^* = 0$, what are the possible range of values of $(w^*)^T x_i + b^*$?
 - (ii) If $0 < \alpha^* < C$, what are the possible range of values of $(w^*)^T x_i + b^*$?
 - (iii) If $\alpha^* = C$, what are the possible range of values of $(w^*)^T x_i + b^*$?

(Hint: Use KKT complementary slack conditions and $\beta_i^* = C - \alpha_i^*$)

6. Consider the following 1-dimensional classification dataset with 8 points given by:

$$X^T = [1 \quad 2 \quad 4 \quad 5 \quad 6 \quad 7 \quad 9 \quad 10]$$

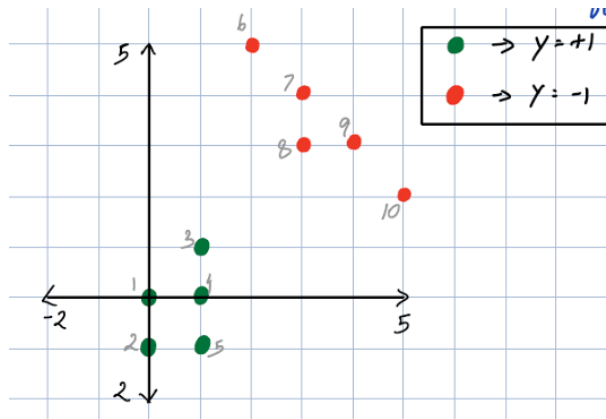
$$y^T = [+1 \quad +1 \quad -1 \quad -1 \quad -1 \quad -1 \quad +1 \quad +1]$$

- (i) Solve Kernel SVM (Hard Margin) (i.e. give the optimal α^*) with $K(u, v) = \exp(-\gamma(u - v)^2)$ where $\gamma = 0.1$.
Hint: $\alpha_2^* = \alpha_7^* = \alpha_3^* = \alpha_6^* > 0$ while all other $\alpha_i^* = 0$.
Find such an α^* and then use KKT conditions to show that it is optimal.
- (ii) Give b^* for all α^* above.
- (iii) Give the decision function $(w^*)^T \phi(x) + b^*$ for $x \in \mathbf{R}$.

7. Consider the following hard margin SVM problem with both w and b . Assume the kernel to be the linear kernel.

- (i) Argue what points are support vectors.
- (ii) Argue what would be the optimal hyperplane and give w^*, b^* .
- (iii) Also argue what α^* should be.

(You can use software to check intuition, and use KKT conditions to verify if a proposed α^* is actually an optimal solution.)



- (iv) Repeat parts (i), (ii), (iii) if point (x_8, y_8) is removed.
- (v) Repeat parts (i), (ii), (iii) with point $(x_8, y_8), (x_9, y_9), (x_7, y_7)$ removed.

(Hint: Optimal α^* need not be unique even if w^* and b^* are.)

8. Consider the following 2-dimensional classification dataset with 5 points given by:

$$X^T = \begin{bmatrix} 1 & 1 & 2 & 4 & 5 \\ 1 & 0 & 5 & 4 & 2 \end{bmatrix}$$

$$y^T = [+1 \quad +1 \quad -1 \quad -1 \quad -1]$$

Consider the hard margin SVM problem with linear kernel $k(u, v) = u^T v$.

- (i) Give the support vectors just by looking at the data. Give reasons.
- (ii) Give the dual solution α^* using the answer to the above part.
- (iii) Check if the entire solution got above is the right answer using KKT conditions. (Thus also checking the first part guessed by “eyeballing”.)
- (iv) Derive the primal solution w^*, b^* from the dual solution α^* and draw a figure illustrating the final solution.

9. Consider the following 2-dimensional binary classification dataset with 10 points given by

$$X^T = \begin{bmatrix} 1 & 1 & 2 & 2 & 4 & 4 & 5 & 5 & 2.9 & 3.1 \\ 0 & 1 & 0 & 1 & 3 & 4 & 3 & 4 & 6 & 6 \end{bmatrix}$$

$$y^T = [-1 \quad -1 \quad -1 \quad -1 \quad +1 \quad +1 \quad +1 \quad +1 \quad -1 \quad +1]$$

Consider the soft-margin linear SVM problem with $C = 0.1, 1, 10, 100$. For each C evaluate the following w, b . By evaluate, we mean you should give the slack variable ξ that make the w, b, ξ feasible, and also give the value of the objective.

- (i) $w = (\frac{1}{2}, 0), b = \frac{-3}{2}$ (ii) $w = (1, 0), b = -3$ (iii) $w = (4, 0), b = -12$ (iv) $w = (16, 0), b = -48$
(v) $w = (64, 0), b = -192$ (vi) $w = (\frac{1}{4}, \frac{1}{4}), b = \frac{-5}{4}$ (vii) $w = (\frac{1}{2}, \frac{1}{2}), b = \frac{-5}{2}$ (viii) $w = (1, 1), b = -5$
(ix) $w = (2, 2), b = -10$ (x) $w = (4, 4), b = -20$

10. (i) Why is logistic regression called regression?

- (ii) Consider the following 2-dimensional classification dataset with 8 points given by:

$$X^T = \begin{bmatrix} -2 & -2 & -1 & -1 & 1 & 1 & 2 & 3 \\ -1 & 2 & 1 & 2 & 1 & 3 & 3 & 2 \end{bmatrix}$$

$$y^T = [+1 \quad +1 \quad +1 \quad -1 \quad -1 \quad +1 \quad -1 \quad -1]$$

Run one iteration of gradient descent with the logistic regression objective by hand. No bias required, only the 2-dimensional weight vector is to be optimised. Choose the step size $\eta = 1$. Initialise at $w = [0, 0]^T$.

11. Recall that the Empirical Logistic Loss Minimisation is given by

$$\hat{R}(w) = \sum_{i=1}^n \log(1 + \exp(-y_i w^T x_i)) = \sum_{i=1}^n \Psi_L(y_i w^T x_i)$$

Prove that $\Psi_L : \mathbf{R} \rightarrow \mathbf{R}$ is convex. Specifically prove that $\Psi_L(u) = \log(1 + \exp(-u))$ is convex in u . Subsequently argue that $\hat{R}(w)$ is convex in w .

12. Let the data instance X be a d -dimensional vector. A function $K : \mathbf{R}^d \times \mathbf{R}^d \rightarrow \mathbf{R}$ is valid kernel function if there exists $\phi : \mathbf{R}^d \rightarrow \mathbf{R}^{d'}$ such that $K(u, v) = \phi(u)^T \phi(v)$. Such a ϕ is called a feature map for the kernel K .

- (i) Let $d = 2, k = 2$. Prove that $K(u, v) = (u^T v)^k$ is a valid kernel. Give the feature map corresponding to this kernel.
(ii) Repeat the above for $d = 3, k = 2$.
(iii) Repeat the above for $d = 2, k = 3$.
(iv) Infer the general form of the feature map ϕ of the kernel $K : (u, v) \rightarrow (u^T v)^k$ for any d, k .
(v) Repeat the four items above for the kernel $K(u, v) = (1 + u^T v)^k$.

13. Let K_1 and K_2 be a valid kernel functions, with feature mapping $\varphi_1 : \mathbb{R}^d \rightarrow \mathbb{R}^{d_1}$ and $\varphi_2 : \mathbb{R}^d \rightarrow \mathbb{R}^{d_2}$.

- (i) Show that $K_3 = K_1 + K_2$ is also a valid kernel. Give the feature mapping φ_3 corresponding to K_3 in terms of φ_1 and φ_2 .
(ii) Show that $K_4 = K_1 \cdot K_2$ is also a valid kernel. Give the feature mapping φ_4 corresponding to K_4 in terms of φ_1 and φ_2 .
(iii) Show that $K_5 = f(u)K_1(u, v)f(v)$ is also a valid kernel for any function $f : \mathbb{R}^d \rightarrow \mathbb{R}$. Give the feature mapping φ_5 corresponding to K_5 in terms of φ_1 and f .
(iv) Show that a Kernel given by $K(u, v) = \exp(2u^T v)$ is a valid kernel. [Hint: Use the results above on a polynomial expansion of $\exp(t)$.]
(v) (Optional) Show that a Kernel given by $K(u, v) = \exp(-\|u - v\|^2)$ is a valid kernel. [Hint: Use the last two parts' results.]