

HW 07

down

(a) I work on this homework with Li Yuanming, methods of (b).
We orally discussed the question.

A bit hard for the last question.

(b) I certify that all solution are entirely in my hands and that
I have not looked at another student's solutions. I have not credited any
external sources in this write up. Wang Zijian 8/7

It was a very minimum work at most time left at

HW 07

2. (a) Since $\phi_i > 0$, multiply by ϕ_i just change ξ_i 's weight.

Intuitively, the corresponding unconstrained optimization problem for the SVM with custom margins is:

$$\min_{w,b} \frac{1}{2} \|w\|^2 + C \sum_{i=1}^n \phi_i \max(1 - y_i(w^T x_i - b), 0)$$

$$\text{s.t. } (1 - \xi_i) - y_i(w^T x_i - b) \leq 0 \quad \forall i$$

$$(y_i w)^T X (y_i w)^T X = 0 \quad \text{and} \quad 1 - \xi_i \leq 0^T w \quad (1)$$

(b) Lagrangian for the problem:

$$\mathcal{L}(w, b, \xi, \alpha, \beta) = \frac{1}{2} \|w\|^2 + C \sum_{i=1}^n \phi_i \xi_i + \sum_{i=1}^n \alpha_i ((1 - \xi_i) - y_i(w^T x_i - b)) + \sum_{i=1}^n \beta_i (\xi_i)$$

$$= \frac{1}{2} \|w\|^2 - \sum_{i=1}^n \alpha_i y_i (w^T x_i - b) + \sum_{i=1}^n \alpha_i + \sum_{i=1}^n (C \phi_i - \alpha_i - \beta_i) \xi_i$$

∴ The dual is $\max_{\alpha \geq 0, \beta \geq 0} g(\alpha, \beta)$

$$\text{where } g(\alpha, \beta) = \min_{w, b, \xi} \frac{1}{2} \|w\|^2 - \sum_{i=1}^n \alpha_i y_i (w^T x_i - b) + \sum_{i=1}^n \alpha_i + \sum_{i=1}^n (C \phi_i - \alpha_i - \beta_i) \xi_i$$

Use KKT conditions to find the optimal dual variables. Verify that the primal problem is convex in the primal variables as stated in notes 22. We know that from the stationary conditions, evaluated at the optimal dual values α^* and β^* , and the optimal primal values w^*, b^*, ξ_i^* :

$$\frac{\partial \mathcal{L}}{\partial w_i} = \frac{\partial \mathcal{L}}{\partial b} = \frac{\partial \mathcal{L}}{\partial \xi_i} = 0$$

$$\frac{\partial \mathcal{L}}{\partial w} = w^* - \sum_{i=1}^n \alpha_i^* y_i x_i = 0 \Rightarrow w^* = \sum_{i=1}^n \alpha_i^* y_i x_i$$

$$\frac{\partial \mathcal{L}}{\partial b} = \sum_{i=1}^n \alpha_i^* y_i = 0.$$

$$\frac{\partial \mathcal{L}}{\partial \xi_i} = C \phi_i - \alpha_i^* - \beta_i^* = 0 \Rightarrow 0 \leq \alpha_i^* \leq C \phi_i, \alpha_i^* \text{ are restricted to being less than the hyperparameter } C.$$

Verify that other KKT also hold, indicating strong duality.

$$\Rightarrow \mathcal{L}(w, b, \xi, \alpha^*, \beta^*) = \frac{1}{2} \|w\|^2 - \sum_{i=1}^n \alpha_i^* y_i (w^T x_i) + b \sum_{i=1}^n \alpha_i^* y_i + \sum_{i=1}^n \alpha_i^* + \sum_{i=1}^n (C \phi_i - \alpha_i - \beta_i) \xi_i$$

$$= \frac{1}{2} \|w\|^2 - \sum_{i=1}^n \alpha_i^* y_i (w^T x_i) + \sum_{i=1}^n \alpha_i^*.$$

Since the primal problem is convex, from the KKT conditions we have that the optimal primal variables w^*, b^*, ξ^* minimize $\mathcal{L}(w, b, \xi, \alpha^*, \beta^*)$

$$g(\alpha^*, \beta^*) = \min_{w, b, \xi} \mathcal{L}(w, b, \xi, \alpha^*, \beta^*)$$

$$= \mathcal{L}(w^*, b^*, \xi^*, \alpha^*, \beta^*)$$

$$(\text{same process in notes}) = \frac{1}{2} \left\| \sum_{i=1}^n \alpha_i^* y_i x_i \right\|^2 - \sum_{i=1}^n \alpha_i^* y_i \left(\left(\sum_{j=1}^n \alpha_j^* y_j x_j \right)^T x_i \right) + \sum_{i=1}^n \alpha_i^*$$

$$= \frac{1}{2} \left\| \sum_{i=1}^n \alpha_i^* y_i x_i \right\|^2 - \sum_{i=1}^n \alpha_i^* y_i x_i^T \left(\sum_{j=1}^n \alpha_j^* y_j x_j \right) + \sum_{i=1}^n \alpha_i^*$$

$$= \alpha^{*T} \mathbf{1} - \frac{1}{2} \alpha^{*T} Q \alpha^*, \quad \text{where } Q_{ij} = y_i (x_i^T x_j) y_j$$

$$\text{and } Q = (\text{diag } y) X X^T (\text{diag } y)$$

The dual is only in terms of α and X and y .

To WH

$$\max_{\alpha} \alpha^T 1 - \frac{1}{2} \alpha^T Q \alpha \text{ subject to } 0 \leq \alpha_i \leq C \quad (a)$$

subject to $\sum_i \alpha_i y_i = 0$ (homogeneous constraints) and $\alpha_i \geq 0$ (non-negativity)

$$0 \leq \alpha_i \leq C \quad i=1, \dots, n \quad (\text{margin constraint})$$

$$(y_i(\alpha_i - \alpha_j^T w))_i + \frac{1}{2} \|\alpha\|_2^2 \quad (\text{margin term})$$

$$\text{if } \alpha_i < (d - \alpha_j^T w) \text{ then } -1 \quad (i \neq j)$$

$$(c) \quad \text{We have } \max_{\alpha} \alpha^T 1 - \frac{1}{2} \alpha^T Q \alpha, \text{ where } Q = (\text{diag } y) X X^T (\text{diag } y)$$

Notice that there is our desired term $X X^T$ in Q , so we just kernelize

it by replacing $X X^T$ by gram matrix K ($X = K^T$, $K = X X^T$)

We found that custom margin SVM has only one difference compared to normal (the dual problem of) SVM soft margin sum, which is $0 \leq \alpha_i \leq C \phi_i$.

ϕ_i here limit the value of α_i .

Adjusting ϕ_i could therefore adjust the range of α_i . $\phi_i = (0, b)$ shows

more margin leaving set \mathcal{S} . Smaller b limits α_i to a small set of indices. If b is large, it may leave some margin set empty. So zero or large b is better for leaving set in training.

$$\alpha = \frac{16}{126} = \frac{16}{7+1} = \frac{16}{16} = 1 \quad (\text{if } w \text{ below})$$

$$x_i^T w \frac{1}{\|w\|_2} = w \leftarrow \alpha = x_i^T w \frac{1}{\|w\|_2} - w = 1 \quad w$$

$$\alpha = x_i^T w \frac{1}{\|w\|_2} = \frac{16}{16}$$

$$\text{or position wise, } \alpha_i = x_i^T w - w = \alpha = x_i^T w - w = \frac{16}{16}$$

returning w after α is fixed

minimizes $\|w\|_2^2$ plus also $\frac{1}{2} \alpha^T Q \alpha$ term of loss

$$x_i^T w + \frac{1}{\|w\|_2} + \frac{1}{\|w\|_2} + (x_i^T w)^2 \frac{1}{\|w\|_2} + \|w\|_2^2 = (x_i^T w, \frac{1}{\|w\|_2}, d, w) \perp \rightarrow$$

$$x_i^T w + (x_i^T w)^2 \frac{1}{\|w\|_2} + \|w\|_2^2 =$$

training set and one additional $\frac{1}{\|w\|_2}$ part. Using α in training set

$$(x_i^T w, \frac{1}{\|w\|_2}, d, w) \perp \rightarrow \frac{1}{\|w\|_2} = x_i^T w \perp \rightarrow$$

$$(x_i^T w, \frac{1}{\|w\|_2}, d, w) \perp \rightarrow$$

$$(x_i^T w, \frac{1}{\|w\|_2}, d, w) \perp \rightarrow$$

$$\frac{1}{\|w\|_2} + (x_i^T w)^2 \frac{1}{\|w\|_2} + \|w\|_2^2 = (\text{return in training set})$$

$$\frac{1}{\|w\|_2} + (x_i^T w)^2 \frac{1}{\|w\|_2} + \|w\|_2^2 = \frac{1}{\|w\|_2} + (x_i^T w)^2 \frac{1}{\|w\|_2} + \|w\|_2^2$$

$$\begin{aligned}
 3 \cdot (a) \quad p(y|\theta) &= p(y_1, y_2, y_3, y_4, y_5|\theta) \\
 &= \frac{8!}{y_1! y_2! y_3! y_4! y_5!} \cdot p_1^{y_1} \cdots p_5^{y_5} \\
 \therefore \log p(y|\theta) &= \log \left(\frac{8!}{1 \cdot 1 \cdot 2! \cdot 2! \cdot 3!} \right) + y_1 \log p_1 + y_2 \log(p_2) + y_3 \log(p_3) + y_4 \log(p_4) + y_5 \log(p_5) \\
 &= \log \left(\frac{8!}{2^4} \right) + y_1 \log \left(\frac{1}{2} \right) + y_2 \log \left(\frac{\theta}{4} \right) + y_3 \log \left(\frac{1-\theta}{4} \right) + y_4 \log \left(\frac{\theta}{4} \right) + y_5 \log \left(\frac{1-\theta}{4} \right)
 \end{aligned}$$

where $y_1 + y_2 = 1$

$$p(y|\theta) = \frac{8!}{2^4} \cdot \left(\frac{1}{2}\right)^{y_1} \cdot \left(\frac{\theta}{4}\right)^{y_2} \cdot \left(\frac{1-\theta}{4}\right)^{y_3}$$

(b)

$$\begin{aligned}
 q_{10}^{t+1} &= P_0(Y_1=0, Y_2=1 | y_{\text{sum}}) \\
 &= \frac{P_0(Y_{\text{sum}} | Y_1=0, Y_2=1) P_0(Y_1=0, Y_2=1)}{P_0(Y_{\text{sum}})} \quad (\text{Bayes' Rule}) \\
 &= \frac{1 \cdot P_0(Y_1=0, Y_2=1)}{P_0(Y_1=0, Y_2=1) + P_0(Y_1=1, Y_2=0)} \\
 &= \frac{\frac{8!}{2^4} \cdot \left(\frac{1}{2}\right)^0 \cdot \left(\frac{\theta}{4}\right)^4 \cdot \left(\frac{1-\theta}{4}\right)^4}{\frac{8!}{2^4} \cdot \left(\frac{1}{2}\right)^0 \cdot \left(\frac{\theta}{4}\right)^4 \cdot \left(\frac{1-\theta}{4}\right)^4 + \left(\frac{1}{2}\right)^1 \cdot \left(\frac{\theta}{4}\right)^3 \cdot \left(\frac{1-\theta}{4}\right)^4} \\
 &= \frac{\left(\frac{\theta}{4}\right)^4}{\left(\frac{\theta}{4}\right)^4 + \frac{1}{2} \cdot \left(\frac{\theta}{4}\right)^3}
 \end{aligned}$$

Similarly, $q_{11}^{t+1} = P_0(Y_1=1, Y_2=0 | y_{\text{sum}})$

$$= \frac{\left(\frac{1}{2}\right)^1 \cdot \left(\frac{\theta}{4}\right)^3}{\left(\frac{1}{2}\right)^1 \cdot \left(\frac{\theta}{4}\right)^3 + \left(\frac{1}{2}\right)^0 \cdot \left(\frac{\theta}{4}\right)^4} = \frac{\theta^3}{2 + \theta^4}$$

$$\begin{aligned}
 (c) \quad \theta^{t+1} &= \arg \max_{\theta} E_{\theta} [\log (P_0(Y_1 | y_{\text{sum}}, Y_3, Y_4, Y_5))] \quad , \text{ denote } y_1! y_2! \cdots y_k! \text{ as } k \\
 &= \arg \max_{\theta} [q_{10}^{t+1} \log (P_0(Y_1 | y_{\text{sum}}, Y_3, Y_4, Y_5)) + q_{11}^{t+1} \log (P_0(Y_1 | y_{\text{sum}}, Y_3, Y_4, Y_5))]
 \end{aligned}$$

$$\begin{aligned}
 &= q_{10}^{t+1} \cdot \log \left(\frac{n!}{k} \cdot \left(\frac{1}{2}\right)^0 \cdot \left(\frac{\theta}{4}\right)^{y_1+1} \cdot \left(\frac{1-\theta}{4}\right)^{y_2+y_3+y_4} \right) + q_{11}^{t+1} \log \left(\frac{n!}{k} \left(\frac{1}{2}\right)^1 \cdot \left(\frac{\theta}{4}\right)^{y_1} \cdot \left(\frac{1-\theta}{4}\right)^{y_2+y_3+y_4} \right) \\
 &= q_{10}^{t+1} \left(\log \left(\frac{n!}{k} \right) + (y_5+1) \log \left(\frac{\theta}{4} \right) + (y_3+y_4) \log \left(\frac{1-\theta}{4} \right) \right) + q_{11}^{t+1} \left[\log \left(\frac{n!}{k} \right) + \log \frac{1}{2} + y_5 \log \left(\frac{\theta}{4} \right) + (y_3+y_4) \log \left(\frac{1-\theta}{4} \right) \right]
 \end{aligned}$$

where $k = y_1! y_2! y_3! y_4!$

$$\begin{aligned}
 (f) \theta &= \arg \max_{\theta} P(y_1, \dots, y_5 | \theta, Y_1=0, Y_2=1) \cdot P(y_1, \dots, y_5 | \theta, Y_1=1, Y_2=0) \\
 &= \arg \max_{\theta} \log P(y | \theta, Y_1=0, Y_2=1) + \log P(y | \theta, Y_1=1, Y_2=0) \\
 &= \arg \min_{\theta} -\log \frac{P(y | \theta, Y_1=0, Y_2=1)}{P(y | \theta, Y_1=1, Y_2=0)}
 \end{aligned}$$

where $P(y | \theta, Y_1=0, Y_2=1)$ can be calculated by using the result of part (a).

It's an convex function, so we can find the optimal (minimum) θ .

After obtaining θ , we substitute it back to $P(y | \theta, Y_1=0, Y_2=1)$
and $P(y | \theta, Y_1=1, Y_2=0)$

Pick the one with higher probability and the corresponding Y_1, Y_2 will be selected.

4. (a) Step 1: initialize p_a , p_b with arbitrary value in range (0,1).

Step 2: ^{int stage,} for $\forall i$, compare $P_b^t h_i (p_a^t)^{h_i} (1-p_a^t)^{1-h_i}$ --- $P(1)$
 $(1-h_i) (p_b^t)^{h_i} (1-p_b^t)^{1-h_i}$ --- $P(2)$

Compare these two possibility probability. Assign h_i to class a if $P(1)$ is higher, assign the data to class b if $P(2)$ is higher.
 (Therefore, data are classified into 2 groups)

Step 3: calculate p_a^{t+1} and p_b^{t+1}

$$p_a^{t+1} = \frac{\sum h_i}{\text{number of data in class a}} \quad (\text{average of probability of points in a})$$

$$p_b^{t+1} = \frac{\sum h_i}{\text{number of data in class b}} \quad (h_i \in \text{class b})$$

Repeat step 2,3 until p_a , p_b converge.
 $\therefore \hat{p}_a = p_a^{t+1}$, $\hat{p}_b = p_b^{t+1}$, $\alpha_a = \frac{\text{number of data in class a}}{\text{number of all data}}$

(b)

No.

Try When $n \rightarrow \infty$, the number of heads will be their expected value

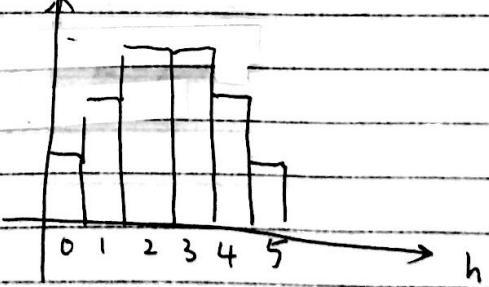
$$\begin{aligned} P(H=0) &= \alpha_a \cdot \binom{5}{0} \cdot p_a^0 \cdot (1-p_a)^5 + \alpha_b \cdot \binom{5}{0} \cdot p_b^0 \cdot (1-p_b)^5 \\ &= 0.5 [0.6^5 + 0.4^5] \\ &= \frac{11}{250} \end{aligned}$$

$$P(H=5) = 0.5 [0.4^5 + 0.6^5] = \frac{11}{250}$$

$$P(H=1) = P(H=4) = 0.5 \cdot [5 \cdot (0.4 \times 0.6^4 + 0.6 \times 0.4^4)] = \frac{21}{125}$$

$$P(H=2) = P(H=3) = 0.5 \cdot [10 \cdot (0.4^2 \times 0.6^3 + 0.6^2 \times 0.4^3)] = \frac{36}{125}$$

- Thus,
 $h=0 \rightarrow \frac{11}{250} n$ data
 $h=1 \rightarrow \frac{21}{125} n$ data
 $h=2 \rightarrow \frac{36}{125} n$ data
 $h=3 \rightarrow \frac{36}{125} n$ data
 $h=4 \rightarrow \frac{21}{125} n$ data
 $h=5 \rightarrow \frac{11}{250} n$ data.



Initial) $P_a = 0.1$, $P_b = 0.9$ from def. of estimation: ϵ gsr (b)

$h=0, 1, 2$ will be assigned to class a (less head)

$h=3, 4, 5$ will be assigned to class b (more head)

$$P_a^1 = \frac{11n}{250} \cdot \frac{1}{5} + \frac{3}{125} \cdot \frac{21n}{5} + \frac{36n}{125} \cdot \frac{2}{5} = 0.2976$$

$$P_b^1 = \frac{11n}{250} \cdot \frac{5}{5} + \frac{21n}{125} \cdot \frac{4}{5} + \frac{36n}{125} \cdot \frac{3}{5} = 0.7024$$

In second iteration, $h=0, 1, 2$ will still be assigned to class a
 $h=3, 4, 5$ will still be assigned to class b

$$\therefore P_a^2 = 0.2976$$

$$P_b^2 = 0.7024$$

$P_a^2 = 0.2976$, $P_b^2 = 0.7024$, they converge but not equal to 0.4 and 0.6.

\therefore "k-means" based estimates will not give the correct estimates.

No matter what initial values we assign to P_a and P_b .

Output histogram right of H is equal to median str., $\epsilon \leftarrow n$ gsr

$$(0.9-1) \cdot 0.9 \cdot \binom{2}{0} \cdot 0.6 + (0.9-1) \cdot 0.9 \cdot \binom{2}{1} \cdot 0.3 = (0 = H) 9$$

$$[0.9 \cdot 0 + 0.3 \cdot 0] \cdot 2.0 =$$

$$\frac{11}{250} =$$

$$\frac{11}{250} = [0.9 \cdot 0 + 0.3 \cdot 0] \cdot 2.0 = (0 = H) 9$$

$$\frac{11}{250} = [(0.9 \cdot 0 + 0.3 \cdot 0) \cdot 2] \cdot 2.0 = (0 = H) 9 = (1 = H) 9$$

$$\frac{11}{250} = [(0.9 \cdot 0 + 0.3 \cdot 0) \cdot 0] \cdot 2.0 = (\epsilon = H) 9 = (\epsilon = H) 9$$

$\bar{x}_0 = 0 = n$ first str. $n \cdot \frac{11}{250}$ str right, $\epsilon = n$

first str. $n \cdot \frac{11}{250} \leftarrow 1 = n$

str. $n \cdot \frac{11}{250} \leftarrow 5 = n$

str. $n \cdot \frac{11}{250} \leftarrow \epsilon = n$

str. $n \cdot \frac{11}{250} \leftarrow \rho = n$

str. $n \cdot \frac{11}{250} \leftarrow 7 = n$

(b) initialize p_a^0, p_b^0, α_a^0

(c) (i) E-step.

$$p^{t+1}(z_i = k | h_i) = \frac{\alpha_k^t p(h_i | z_i = k)}{\alpha_a^t p(h_i | z_i = a) + \alpha_b^t p(h_i | z_i = b)} \quad (z_i \text{ is the label for } h_i, \therefore k = a, b)$$

$$\therefore p^{t+1}(z_i = a | h_i) = \frac{\alpha_a^t p(h_i | z_i = a)}{\alpha_a^t p(h_i | z_i = a) + \alpha_b^t p(h_i | z_i = b)}, \quad \alpha_b^t = 1 - \alpha_a^t$$

$$= \frac{\alpha_a^t (h_i) (p_a^t)^{h_i} (1-p_a^t)^{1-h_i}}{\alpha_a^t (h_i) (p_a^t)^{h_i} (1-p_a^t)^{1-h_i} + \alpha_b^t (h_i) (p_b^t)^{h_i} (1-p_b^t)^{1-h_i}}$$

$$p^{t+1}(z_i = b | h_i) = \frac{\alpha_b^t (h_i) (p_b^t)^{h_i} (1-p_b^t)^{1-h_i}}{\alpha_a^t (h_i) (p_a^t)^{h_i} (1-p_a^t)^{1-h_i} + \alpha_b^t (h_i) (p_b^t)^{h_i} (1-p_b^t)^{1-h_i}}$$

(ii) M step. - denote $p^{t+1}(z_i = a | h_i) = q_{ia}^{t+1}$

$$p^{t+1}(z_i = b | h_i) = q_{ib}^{t+1}$$

The expected complete log likelihood is all data

$$\mathbb{E}_{q^{t+1}} [\ell(\alpha, \theta, z)] \quad \text{let } \Theta \text{ denote } \{\alpha_a, p_a, p_b\},$$

$$\text{Recall the log likelihood } \ell(\theta) = \sum_{i=1}^n \log \left(\sum_{k=a,b} p(h_i | z_i = k) p(z_i = k) \right)$$

The expected complete log likelihood is:

$$\arg \max_{\hat{\alpha}_a, \hat{p}_a, \hat{p}_b} \left(\dots \right)$$

E step: give the probabilities of a coin belongs to class a and class b.

M step: Given observation data, find the optimal $\hat{\alpha}_a, \hat{p}_b, \hat{p}_a$ that maximize the probability of a coin belongs to class a/b.