

Question 4 - AdaBoost

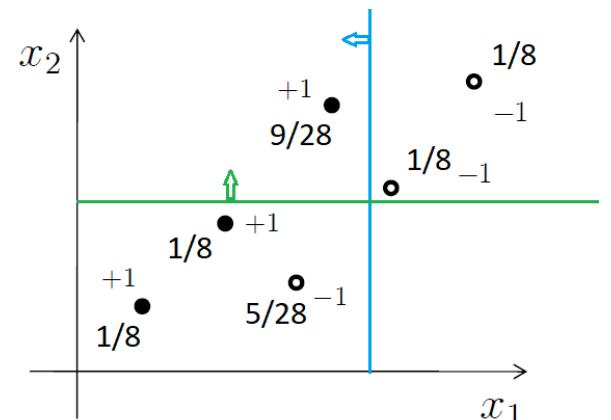
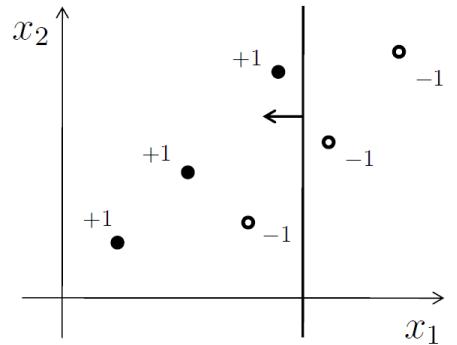
AdaBoost (Adaptive Boosting) is another approach to the ensemble method field.

It always uses the entire data and features (unlike before) and aims to create T weighted classifiers (unlike before, where each classifier had same influence). The new classification will be decided by linear combination of all the classifiers, by:

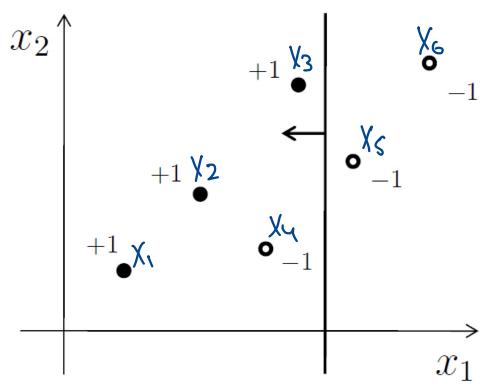
$$g(x) = \text{sign} \left(\sum_{t=1}^T \alpha_t f_t(x) \right), \alpha_t \geq 0$$

Consider the following dataset in \mathbb{R}^2 :

- 1) The first decision stump is already drawn, the arrow points in the positive direction. Calculate the classifier error (ε_1) and weight (α_1).
- 2) Calculate the new weights of the samples (and normalize them to get valid distribution).
- 3) Draw the second decision stump. Reminder: the decision stump (our classifiers) are parallel to x/y axis.
- 4) Without calculations, which classifier's weight is larger, α_1 or α_2 ? Explain why.
- 5) In the right image, there is the dataset and the weights for each point, after finding the third decision stump and calculating the new weights. Which of the following (green or blue) is the correct third decision stump?
- 6) Given $\alpha_2 = 1.1$, $\alpha_3 = 0.62$, draw the full classifier, like in slide 13.
What is the train accuracy?



- 1) The first decision stump is already drawn, the arrow points in the positive direction. Calculate the classifier error (ε_1) and weight (α_1).



השאלה מבקשת חישוב שגיאה של החלטה סימפלית. נסמן x_i כהיפך של נקודה (x_{i1}, x_{i2}) , ו- y_i כערך החלטה של נקודה. השגיאה היא $\varepsilon_i = |y_i - \text{ההחלטה}|$. נסמן $s = \frac{1}{6} \sum \varepsilon_i$.

$$\varepsilon_i = \frac{1}{6}$$

$$\alpha_1 = \frac{1}{2} \ln \left(\frac{\varepsilon_i - 1}{\varepsilon_i} \right) = \frac{1}{2} \ln \left(\frac{\frac{1}{6} - 1}{\frac{1}{6}} \right) = \frac{1}{2} \ln (5) \approx 0.8$$

- 2) Calculate the new weights of the samples (and normalize them to get valid distribution).

השאלה מבקשת חישוב משקלים חדשים $w_i^{(2)}$ ש�� יתנו שגיאה מינימלית:

$$w_i^{(2)} = w_i^{(1)} \cdot e^{(-\alpha_1 y_i f_1(x_i))}$$

המשמעות של פונקציית $f_1(x_i)$ היא שמיון נזק ב- x_i יתבטא ב- y_i ביחס ל- $f_1(x_i)$.

$$w_i^{(2)} = w_i^{(1)} e^{(-\alpha_1 \cdot 1 \cdot 1)} = w_i^{(1)} e^{-0.8}$$

אנו מודדים שכך:

$$\Rightarrow d_1(x_1) = \frac{1}{6} \cdot e^{-0.8} \approx 0.074$$

$$w_i^{(2)} = w_i^{(1)} \cdot e^{(-\alpha_1 \cdot 1 \cdot (-1))} = w_i^{(1)} e^{0.8}$$

אנו מודדים שכך:

$$\Rightarrow d_1(x_2) = \frac{1}{6} \cdot e^{0.8} \approx 0.37$$

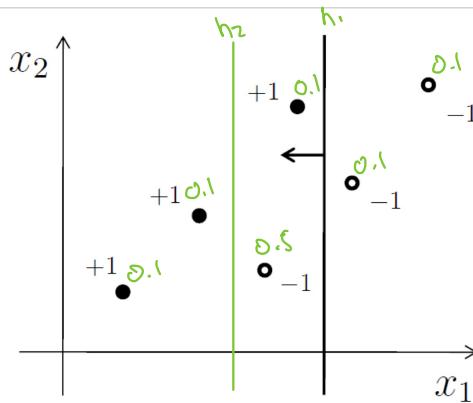
$$z = 0.37 \cdot 1 + 0.074 \cdot 5 \approx 0.74$$

מגניטו גלאי

$$\|\alpha_1(x_C)\| = \frac{0.074}{0.74} = \underline{0.1}$$

$$\|\alpha_1(x_F)\| = \frac{0.37}{0.74} = \underline{0.5}$$

- 3) Draw the second decision stump. Reminder: the decision stump (our classifiers) are parallel to x/y axis.



- 4) Without calculations, which classifier's weight is larger, α_1 or α_2 ? Explain why.

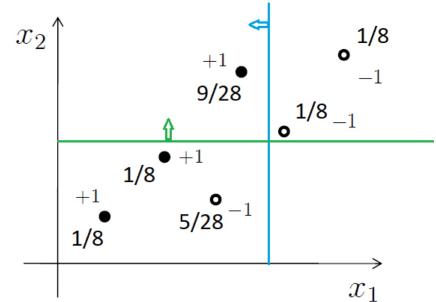
השאלה מבקשת בזבזון נזק. נזכיר כי החלטת ה- α_1 מוגדרת כ $x_1 > 0$.

$\alpha_2 > \alpha_1$ - כי $\epsilon_2 = 0.166$ ו- $\epsilon_1 = 0.1$.

- 5) In the right image, there is the dataset and the weights for each point, after finding the third decision stump and calculating the new weights. Which of the following (green or blue) is the correct third decision stump?

$$E_{\text{blue}} = \frac{5}{28}$$

$$E_{green} = 0.125 \cdot 4 = 0.5$$



ולא מילא ה-Decision stump ב-CART. זו או זו Decision tree ה-Decision stump ב-CART.
בCART מילא ה-Decision stump ב-CART בDecision tree בCART.
הDecision stump ב-CART מילא ה-Decision tree ב-CART.

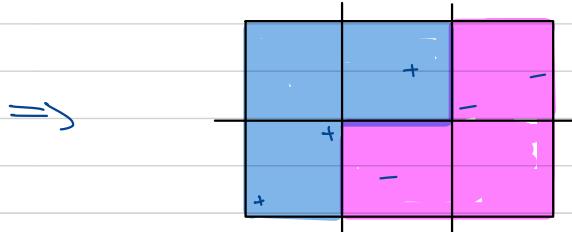
- 6) Given $\alpha_2 = 1.1$, $\alpha_3 = 0.62$, draw the full classifier, like in slide 13.

What is the train accuracy?

$\therefore g(x) \leq N \leq 100$, $\alpha_3 = 0.62$, $\alpha_2 = 1.1$: 178

$$g(x) = \text{sign} \left(\sum_{t=1}^3 \alpha_t f_t(x) \right)$$

$$\Rightarrow \text{Sign}(0.8f_1(x) + 1.1f_2(x) + 0.6f_3(x))$$



Question 5 – Kernel PCA (Bonus 10 points)

PCA is a useful tool to lower the dimensions of the data and get much less feature that yet represent the data. Even though, sometimes the original data is not good enough, so we want to map the data into higher (or same) dimension!

- Recall that in PCA, we require the samples to be mean-centered, $\frac{1}{n} \sum_{i=1}^n x_i = 0$.

We now want to do the same thing, but after x was mapped into $\varphi(x)$. Denote the following, mean-centered version of $\varphi(x)$:

$$v_i = \varphi(x_i) - \frac{1}{n} \sum_{t=1}^n \varphi(x_t)$$

However, we won't always have access to φ / costs a lot of computation (and won't be able to compute v_1, \dots, v_n).

Therefore, we will use the kernel trick. Denote K' as the kernel matrix for v_1, \dots, v_n . Show how K' can be calculated using only the original kernel matrix K on x_1, \dots, x_n .

Reminder: $K'_{i,j} = \langle v_i, v_j \rangle$ (since the new kernel is linear in terms of $\varphi(x)$).

- Now, for the rest of the question we assume that $\frac{1}{n} \sum_{i=1}^n \varphi(x_i) = 0$.

We would like to apply PCA to the vectors $\varphi(x_i)$. Denote by $u_1, \dots, u_k \in R^{d'}$ the first k principal components, the eigenvectors of the scatter matrix S .

Show that u_j (for $j=1, \dots, k$) is linear combination of $\varphi(x_1), \dots, \varphi(x_n)$.

Moreover, show who are the $\alpha_{j,1}, \dots, \alpha_{j,k}$ such that $u_j = \sum_i \alpha_i \varphi(x_i)$.

- Use the expression for $\alpha_{j,1}, \dots, \alpha_{j,n}$ and find an algorithm to obtain the entire $\alpha_j = [\alpha_{j,1}, \dots, \alpha_{j,n}]$ using K .

Hint: substitute u_j you found in the expression for $\alpha_{i,j}$ you found.

- Since we don't really know φ , we can't use u_j . But we also don't need them!

Their only purpose was to project $\varphi(x)$ onto the new dimension, using the dot product, $z_j = \langle u_j, \varphi(x) \rangle$ (and $z^i = [z_1^i \dots z_k^i]$).

Show how to calculate z_j without using φ .

5. Now, use [sklearn](#) and apply Kernel PCA with 'cosine' kernel (doesn't have hyperparameters) and classify with KNN using 125 neighbors, on the same data from question 3.

Did it do better than the regular PCA (on the test set)?

MAKE SURE TO SAVE IT TO COLAB

1. Recall that in PCA, we require the samples to be mean-centered, $\frac{1}{n} \sum_{i=1}^n x_i = 0$.

We now want to do the same thing, but after x was mapped into $\varphi(x)$. Denote the following, mean-centered version of $\varphi(x)$:

$$v_i = \varphi(x_i) - \frac{1}{n} \sum_{t=1}^n \varphi(x_t)$$

However, we won't always have access to φ / costs a lot of computation (and won't be able to compute v_1, \dots, v_n).

Therefore, we will use the kernel trick. Denote K' as the kernel matrix for v_1, \dots, v_n . Show how K' can be calculated using only the original kernel matrix K on x_1, \dots, x_n .

Reminder: $K'_{i,j} = \langle v_i, v_j \rangle$ (since the new kernel is linear in terms of $\varphi(x)$).

We will show how the kernel matrix K' can be calculated using only the original kernel matrix K .

firstly, we will compute K using $K_{i,j} = \langle \varphi(x_i), \varphi(x_j) \rangle$.

after that, we compute the mean of $\varphi(x)$: $\mu = \frac{1}{n} \sum_{t=1}^n \varphi(x_t)$

then, to find the mean-centered data v_i , we express it as $\varphi(x_i) - \mu$

from the Reminder : $K'_{i,j} = \langle v_i, v_j \rangle$

then we expand the inner product, we get that:

$$K'_{i,j} = \langle \varphi(x_i), \varphi(x_j) \rangle - \langle \varphi(x_i), \mu \rangle - \langle \mu, \varphi(x_j) \rangle + \langle \mu, \mu \rangle$$

Now, will replace the values from K :

*: this is the original kernel matrix element, $K_{i,j}$

**: $\langle \varphi(x_i), \frac{1}{n} \sum_{t=1}^n \varphi(x_t) \rangle = \frac{1}{n} \sum_{t=1}^n \langle \varphi(x_i), \varphi(x_t) \rangle = \frac{1}{n} \sum_{t=1}^n K_{i,t}$

***: $\frac{1}{n} \sum_{t=1}^n K_{t,j}$

****: this will be the mean with itself, $\frac{1}{n^2} \sum_{t=1}^n \sum_{s=1}^n K_{s,t}$

So, we get,

$$K'_{i,j} = K_{i,j} - \frac{1}{n} \sum_{t=1}^n K_{i,t} - \frac{1}{n} \sum_{t=1}^n K_{t,j} + \frac{1}{n^2} \sum_{t=1}^n \sum_{s=1}^n K_{s,t}$$

Finally, we will express this in matrix form :

$$k' = k - \frac{1}{n} 1 k - \frac{1}{n} k 1 + \frac{1}{n^2} 1 k 1$$

Thus, we derived the mean-centered kernel matrix k' using only the original matrix k .

2. Now, for the rest of the question we assume that $\frac{1}{n} \sum_{i=1}^n \varphi(x_i) = 0$.

We would like to apply PCA to the vectors $\varphi(x_i)$. Denote by $u_1, \dots, u_k \in R^{d'}$ the first k principal components, the eigenvectors of the scatter matrix S .

Show that u_j (for $j=1, \dots, k$) is linear combination of $\varphi(x_1), \dots, \varphi(x_n)$.

Moreover, show who are the $\alpha_{j,1}, \dots, \alpha_{j,k}$ such that $u_j = \sum_i \alpha_{j,i} \varphi(x_i)$.

here we assume that $\varphi(x_i)$ is already mean-centered.

we will show that each u_j is a linear combination of $\varphi(x_1), \dots, \varphi(x_n)$, and to identify the coefficients $\alpha_{j,1}, \dots, \alpha_{j,k}$ such that $u_j = \sum_{i=1}^n \alpha_{j,i} \varphi(x_i)$.

firstly, we represent the scatter matrix S in the feature space as: $\sum_{i=1}^n \varphi(x_i) \varphi(x_i)^T$.

then we consider the eigenvalue equation for the scatter matrix: $S u_j = \lambda_j u_j$

we will replace this expression into the eigenvalue equation, we get:

$$S \sum_{i=1}^n \alpha_{j,i} \varphi(x_i) = \lambda_j \sum_{i=1}^n \alpha_{j,i} \varphi(x_i)$$

next, we expand the scatter matrix S , we get:

$$\sum_{i=1}^n \sum_{l=1}^n \alpha_{j,l} \varphi(x_l) \underbrace{\langle \varphi(x_l), \varphi(x_i) \rangle}_{k_{il}} = \lambda_j \sum_{i=1}^n \alpha_{j,i} \varphi(x_i)$$

So, we get: $\sum_{i=1}^n k_{il} \alpha_{j,i} = \lambda_j \alpha_{j,l}$, and at the matrix form: $k \alpha_j = \lambda_j \alpha_j$

Thus, we show that each principal component u_j is a linear combination of $\varphi(x_1), \dots, \varphi(x_n)$ with coefficients α_j being the eigenvectors of k .

we can express $u_j = \sum_{i=1}^n \alpha_{ji} e(x_i)$, where the α_{ji} are the eigenvectors corresponding to the eigenvalue λ_j .

3. Use the expression for $\alpha_{j,1}, \dots, \alpha_{j,n}$ and find an algorithm to obtain the entire $\alpha_j = [\alpha_{j,1}, \dots, \alpha_{j,n}]$ using K.

Hint: substitute u_j you found in the expression for $\alpha_{i,j}$ you found.

To find an algorithm for obtaining the entire $\alpha_j = [\alpha_{j,1}, \dots, \alpha_{j,n}]$ using K, we will use the previous relationship.

from part 2, each u_j can be expressed as: $\sum_{i=1}^n \alpha_{ji} e(x_i)$ *

then as part 2, we can get that $K\alpha_j = \lambda_j \alpha_j$

To find α_j we should solve the eigenvalue problem for K.

Algorithm:

1 - calculate $k_{i,j} = \langle e(x_i), e(x_j) \rangle$ for all data points.

2 - Decompose the kernel matrix to obtain its eigenvalues λ_i and corresponding eigenvectors α_i .

3 - α_j that we get represent the coefficients $\alpha_{j,1}, \dots, \alpha_{j,n}$

4 - then, use the coefficients to express each principal component u_j as *

4. Since we don't really know φ , we can't use u_j . But we also don't need them!

Their only purpose was to project $\varphi(x)$ onto the new dimension, using the dot product, $z_j = \langle u_j, \varphi(x) \rangle$ (and $z^i = [z_1^i \dots z_k^i]$).

Show how to calculate z_j without using φ .

Our goal is to compute $z_j = \langle u_j, \varphi(x) \rangle$ without using $\varphi(x)$.
we will use the kernel trick and α_j

firstly, we know that: $u_j = \sum_{i=1}^n \alpha_{j,i} \varphi(x_i)$

next, we need to calculate the inner product $z_j = \langle u_j, \varphi(x) \rangle$, after replacing u_j we get:

$$z_j = \left\langle \sum_{i=1}^n \alpha_{j,i} \varphi(x_i), \varphi(x) \right\rangle = \sum_{i=1}^n \alpha_{j,i} \langle \varphi(x_i), \varphi(x) \rangle$$

↑
linearity of inner product

we know that: $\langle \varphi(x_i), \varphi(x) \rangle$ is equivalent to $k(x_i, x)$, we get:

$$z_j = \sum_{i=1}^n \alpha_{j,i} k(x_i, x) = \alpha_j^\top k(x)$$

↑
in matrix notation

Thus, we have shown how to calculate z_j without using φ .

5. Now, use sklearn and apply Kernel PCA with 'cosine' kernel (doesn't have hyperparameters) and classify with KNN using 125 neighbors, on the same data from question 3.

Did it do better than the regular PCA (on the test set)?

MAKE SURE TO SAVE IT TO COLAB

After applying kernel PCA with 'cosine' kernel and classifying with KNN on the dataset, it performed better than regular PCA on the test set. This indicates that kernel PCA with the 'cosine' kernel captured the data's structure more effectively, improving classification accuracy.