1. The criterion function of a 2-class perceptron algorithm is

$$J(\underline{w}) = - \sum_{n=1}^{N} [\![ \underline{w}^T z_n \underline{x}_n \leq 0 ]\!] \, \underline{w}^T z_n \underline{x}_n$$

This can be written as follows:

$$J(\underline{w}) = \begin{cases} - \sum_{n=1}^{N} \underline{w}^T z_n \underline{x}_n & \text{if } \underline{w}^T z_n \underline{x}_n \leq 0 \\ \\ 0 & \text{if } \underline{w}^T z_n \underline{x}_n > 0 \end{cases}$$

$$= \max \left( - \sum_{n=1}^{N} \underline{w}^T z_n \underline{x}_n , 0 \right)$$

Now, the standard definition of a convex function is as follows: A function $f : \mathbb{R}^n \to \mathbb{R}$ is convex, if the domain of $f$ is convex and $\forall x, y \in$ domain $f$ and $0 \leq \theta \leq 1$ st $f(\theta x + (1-\theta)y) \leq \theta(f(x)) + (1-\theta)f(y)$.

$$J(\underline{w}) = \max \left( - \sum_{n=1}^{N} \underline{w}^T z_n \underline{x}_n , 0 \right) .$$

Let $f(\underline{w_0}) = - \sum_{n=1}^{N} \underline{w}^T z_n \underline{x}_n$ and $g(\underline{w}) = 0 \; \forall \; w_1, w_2 \in \mathbb{R}^n$

and $\theta \in (0,1)$.

Now, $f\left( \theta w_1 + (1-\theta) w_2 \right) \leq \theta f(w_1) + (1-\theta) f(w_2) .$ — ①

$f(w) \leq J(w) = \max \left( f(w), g(w) \right)$

① can be written as:

$$f\left( \theta(w_1) + (1-\theta)w_2 \right) \leq \theta J(\underline{w_1}) + (1-\theta) J(\underline{w_2}) . \quad — ②$$

Similarly, $g\left( \theta w_1 + (1-\theta) w_2 \right) \leq \theta g(w_1) + (1-\theta) g(w_2)$

$$\leq \theta J(w_1) + J(w_2)(1-\theta) . \quad — ③$$

We made an assumption That $f(\underline{w}) = -\sum_{n=1}^{N} \underline{w}^T \underline{z}_n x_n$ and $g(w)$ are

convex.

and so, from ② and ③,

$$J(w_1 + (1-\theta)w_2) \leq \theta J(w_1) + (1-\theta)J(w_2).$$

or, from ② and ③

$$J(\theta w_1 + (1-\theta)w_2) = \max\{f(\theta w_1 + (1-\theta)w_2), g(\theta w_1 + (1-\theta)w_2)\}$$
$$\leq \theta J(w_1) + (1-\theta)J(w_2) \quad\text{//} \qquad \theta \in (0,1).$$

which shows that the max function $J(\underline{w}) = \max(f(w), g(w))$ is convex,

where

$$f(\underline{w}) = -\sum_{n=1}^{N} w^T z_n x_n \qquad\text{and}\qquad g(\underline{w}) = 0 \quad\text{//}.$$
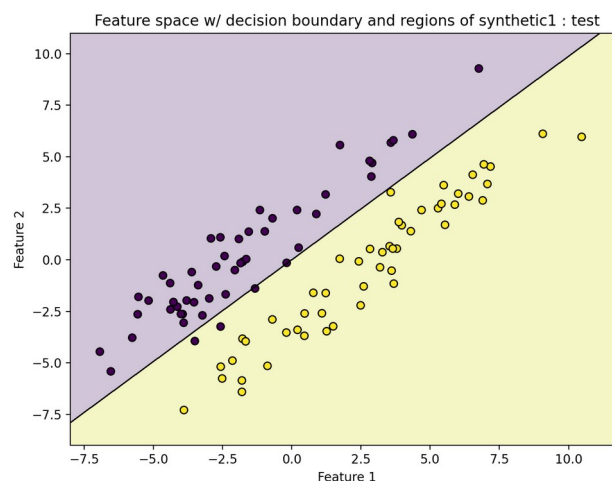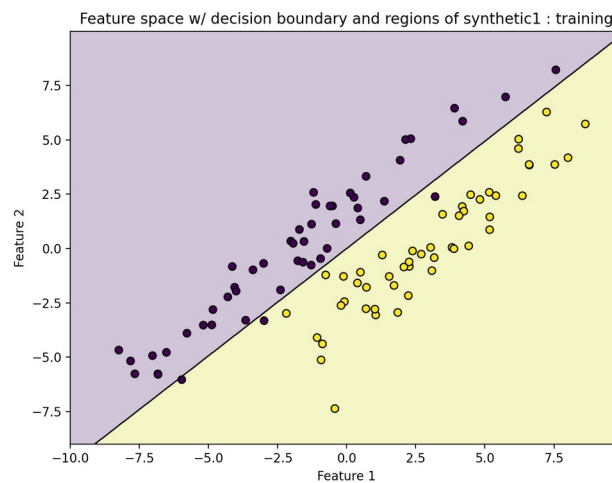
# HOMEWORK 4

2. We train the perceptron classifying algorithm on the training set of each dataset to obtain optimal weights and the final value of the criterion function. Then, the optimal weights are used to make predictions on the test data of the corresponding dataset.

(a) The following are the results on the "synthetic1" dataset:

    (i) The optimal weights $\hat{w}$, i.e the global minima = **[-4.820879, 4.881458]**. The algorithm halted without convergence, with a final **J( $\hat{w}$ ) = 4.19378.**
    (ii) After running the perceptron classifier on the training set and the test set, we obtain classification errors of **0.03 and 0.03**, respectively.
    (iii) The following two figures illustrate the feature space of the training and test data set of synthetic1. Since we're dealing with two features, "feature 1" and"feature 2" are plotted on the x and y axes respectively.  As expected, we can see roughly 3 misclassifications on both the sets.
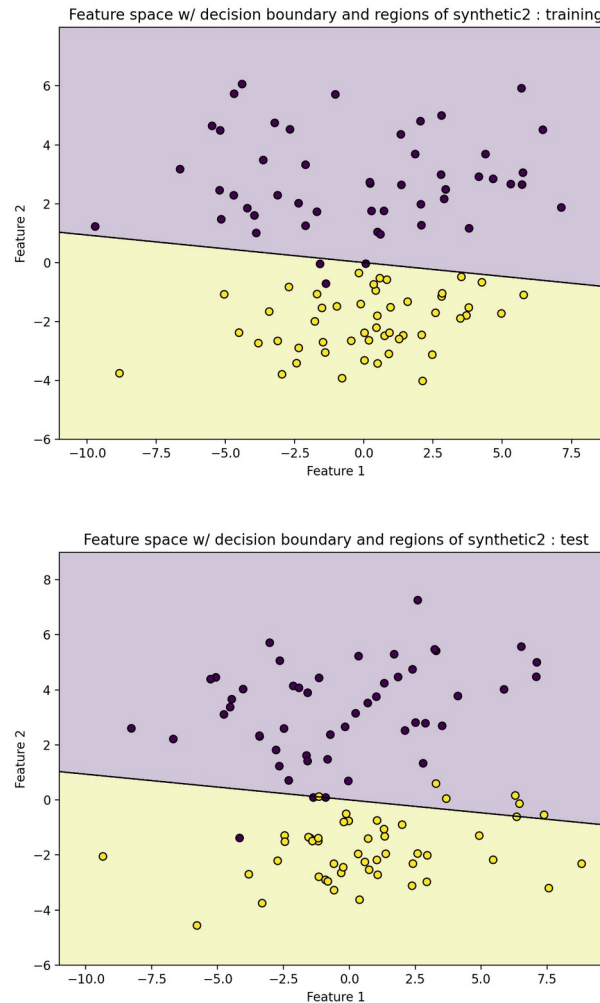


Feature space w/ decision boundary and regions of synthetic1 : training



Feature space w/ decision boundary and regions of synthetic1 : test

(b) The following are the results on the "synthetic2" dataset:

    (i) The optimal weights $\hat{w}$, i.e the global minima = **[0.1579255, 1.688338]**. The algorithm halted without convergence, with a final J( $\hat{w}$ ) = **0.46854**.
    (ii) After running the perceptron classifier on the training set and the test set, we obtain classification errors of **0.03 and 0.08**, respectively.

(iii) The following two figures illustrate the feature space of the training and test data set of synthetic2. We clearly see that there are 3 misclassifications in the training set and 8 in the test set.



Feature space w/ decision boundary and regions of synthetic2 : training



Feature space w/ decision boundary and regions of synthetic2 : test

(c) The following are the results on the "wine" dataset:

(i) The optimal weights $\hat{w}$, i.e the global minima = **[-3505.025, 2358.549, 592.752, -5012.76, -103.1, -2564.207, -1221.414, 106.042, -737.086, -231.54, -401.67, -142.318, 145.7].** The algorithm halted without convergence, with a final J($\hat{w}$) = **113340.7303.**

(ii) After running the perceptron classifier on the training set and the test set, we obtain classification errors of **0.1538 and 0.2**, respectively.

(iii) **No, the data is not linearly seperable**. Projection of the wine data onto the feature space, would lead to a 13-dimensional space, which is difficult to analyse and to conclude if the data is linearly seperable. And so, this leads us to resorting to dimensionality reduction techniques like Principal Component Analysis (PCA) and t-distributed stochastic neighbour embedding (t-SNE), to model high-dim data.

3. The criterion function $J(\underline{w})$ of a multiclass perceptron is as follows:

$$J(\underline{w}) - \sum_{n=1}^{N} [[ \ell_n \neq k_n ]] \left[ \underline{w}_{k_n}^T \underline{x}_n - \underline{w}_{\ell_n}^T \underline{x}_n \right]$$

where $g_{k_n}(\underline{x}_n) \rightarrow$ discriminant of correct class

$g_{\ell_n}(\underline{x}_n) \rightarrow$ discriminant of predicted class by the ML model

$\ell_n = \underset{m \in \{1,2,\dots c\}}{\arg\max} \{ g_m(\underline{x}_n) \}$.

$c \rightarrow$ total classes

Then, choose the max discriminant function.

(a) The **batch GD** weight update formula is:

$$\underline{w}(i+1) = \underline{w}(i) - \eta \nabla_{\underline{w}} J(\underline{w}).$$

As above,

$\ell_n = \underset{m \in \{1,2,\dots c\}, (m \neq k)}{\arg\max} \{ g_m(\underline{x}_n^{(k)}) \}$

$c \rightarrow$ total classes.

If $k_n \neq \ell_n$, we will get two weight update formula : one for $k_n$ & the other $\ell_n$:

$$\underline{w}_{k_n}^{(i+1)} = \underline{w}_{k_n}^{(i)} - \eta \nabla_{\underline{w}_{k_n}} J(\underline{w}).$$

$$= \underline{w}_{k_n}^{(i)} + \eta \underline{x}_n^{(k)} \; /\!/$$

$$\underline{w}_{\ell_n}^{(i+1)} = \underline{w}_{\ell_n}^{(i)} - \eta \nabla_{\underset{m \in \{1,2,\dots c\} \, m \neq k}{\underline{w} \, \arg\max \, g_m(\underline{x}_n^{(k)})}} J(\underline{w})$$

Find gradients of all $c$ weight vectors and take the maximum value for $\ell_n$ computation.

$$\underline{w}_{\ell_n}^{(i+1)} = \underline{w}_{\ell_n}^{(i)} - \eta \, \underline{x}_n^{(k)} \; /\!/$$

If $k_n = \ell_n$, $\underline{w}_{k_n}(i+1) = \underline{w}_{k_n}(i) \; /\!/$

(b) **Sequential gradient descent :**

For each data point, weight updation is as follows :

$$\underline{w}(i+1) = \underline{w}(i) - \eta \, \nabla_{\underline{w}} J_n(\underline{w}) \qquad \eta(i) \geq 0 .$$

From previous, if $k_n = l_n$ i·e $g_{k_n}(x_n) > g_{j_n}(x_n) \quad \forall j \neq k$,

then is no weight up date

or $\underline{w}_{k_n}(i+1) = \underline{w}_{k_n}(i) \qquad (\because J(\underline{w}) = 0)$ .

If $k_n \neq l_n$ , (in correct classify)
$$w_{k_n}(i+1) = w_{k_n}(i) + \eta \, x_n^{(k)} \qquad \eta(i) \geq 0$$
$$w_{l_n}(i+1) = w_{l_n}(i) - \eta \, x_n^{(k)} \qquad \eta(i) \geq 0 .$$

Hence, the algorithm is as follows :
Randomly shuffle the order of data points .
Initialise $\underline{w}(0)$ i·e $w_1(0) = w_2(0) = \cdots = w_c(0)$ for $C$ classes
Define seq-epoch $m$ as :
  for $n$ in $\{1, 2, \ldots N\}$
    $i = (m-1)N + (n-1)$ ; $k_n \in \{1, 2, \ldots, c\}$ .
    $l_n = \text{argmax} \quad g_m(x_n^{(k)})$ .
      $m \in \{1, 2, \ldots, c\}, m \neq k$ .
    if $l_n \neq k_n$ (incorrect classify) $\qquad k_n \rightarrow$ correct label from
      $w_{k_n}(i+1) = w_{k_n}(i) + \eta \, x_n^{(k)}$ $\qquad$ dataset .
      $w_{l_n}(i+1) = \overline{w_{l_n}(i) + \eta x_n} \quad w_{l_n}(i) - \eta \, x_n^{(k)}$
    else
      $w_{k_n}(i+1) = w_{k_n}(i)$ . $\forall k_n \in \{1, 2, \ldots, c\}$ . //
  until : all points are correctly classified .

(c) This weight-update in (b) is very similar to what was in
lecture class. However, we initialise each class of the weight vector at first,
to compute the gradients (for each class). We make sure $J(\underline{w}) = \sum_{n=1}^{N} J_n(\underline{w})$.

4.

A 2-class perceptron with a margin has the following criterion function :

$$J(\underline{w}) = -\sum_{n=1}^{N} [[\underline{w}^T z_n \underline{x}_n \leq b]] \underline{w}^T z_n \underline{x}_n \; ; \quad \text{where} \quad b \text{ is a marginal value} \in \mathbb{R}.$$

Clearly, there will be a misclassification when $\underline{w}^T z_n \underline{x}_n \leq b$.

We make few assumptions : fixed increment i.e $\eta(i) = \eta > 0$.

$$\nabla_{\underline{w}} J(\underline{w}) = -\sum_{n=1}^{N} [[\underline{w}^T z_n \underline{x}_n \leq b]] z_n \underline{x}_n \quad \text{where} \quad [[\cdot]] \rightarrow \text{indicator function}.$$

At first, we randomly shuffle the data $\rightarrow$ as per the sequential gradient approach.

Initialise $\underline{w}(0) \longrightarrow$ init. weight vector.

$\underline{w}(0) \longrightarrow$ arbitrary.

$$\underline{w}(i+1) = \underline{w}(i) + \eta \sum_{n=1}^{N} [[\underline{w}^T z_n \underline{x}_n \leq b]] z_n \underline{x}_n$$

or $\underline{w}(i+1) = \underline{w}(i) + \eta [[\underline{w}^T z_i \underline{x}_i \leq b]] z_i \underline{x}_i$

$z_i x_i \longrightarrow$ cyclically ordered training data points (over a series of epochs).

At the $i^{th}$ iteration, let the point $z^i x^i$ be misclassified where $z^i x^i$ reflection of the data point $x^i$.     or $\underline{w}^T_{[i]} z^i x^i \leq b$.

and so, the corresponding weight update :

$$w(i+1) = w(i) + \eta \, z^i x^i. \qquad ; \eta > 0 \qquad \text{where} \qquad w(i)^T z^i x^i - b \leq 0 . \forall i$$
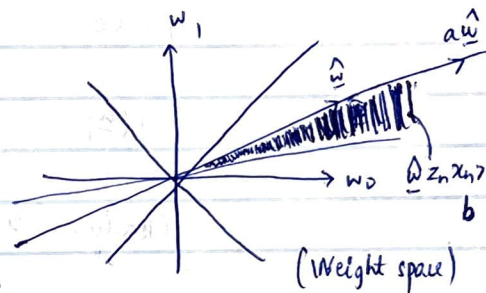
Let $\hat{w}$ be an optimum weight vector solution i.e

$$\hat{w}^T z_n x_n \geq b \quad \forall n$$

then, $a\hat{w}$ be a solution as well:

$$a\hat{w}^T z_n x_n > b \quad \forall n \quad , \quad a > 0$$

Note: The solution region lies within '$b$'.



(Weight space)

Error: $\quad E_w(i) = \| \hat{w}(i) - a\hat{w}(i) \|_2^2$ . where $a$ is adjustable.

Hence, the error must decrease at each iteration and this what we chose to prove.

$\longrightarrow \quad w(i+1) - a\hat{w} = \cancel{w(i)} (w(i) - a\hat{w}) + \eta z^i x^i \quad , a > 0$

Taking squared of $L2$ norm:

or $\| w(i+1) - a\hat{w} \|_2^2 = \| w(i) - a\hat{w} \|_2^2 + \| \eta z^i x^i \|_2^2 + 2 [w(i) - a\hat{w}]^T \eta z^i x^i \quad , a > 0$

Since the increment is fixed, let's drop $\eta$.

~~Now, $w(i) x^i$~~

$w(i) z^i x^i \leq b$.

So, $\| w(i+1) - a\hat{w} \|_2^2 \leq \| w(i) - a\hat{w} \|_2^2 + \| z^i x^i \|_2^2 - 2a\hat{w}^T z^i x^i + 2b$.

but since $\hat{w}$ is optimal weight vector, $\hat{w}^T z^i x^i \geq b$.

Let $k^2 = \max_j \| x_j \|_2^2$ = { Length of longest $x$-vector }$^2$

and $l \triangleq \min_j \{ \hat{w}^T z^i x^i \} > b$.

$\therefore \quad \| w(i+1) - a\hat{w} \|_2^2 \leq \| w(i) - a\hat{w} \|_2^2 + k^2 - 2al + 2b$.

Let $a = \dfrac{k^2}{l}$.

$\therefore \quad \| w(i+1) - a\hat{w} \|_2^2 \leq \| w(i) - a\hat{w} \|_2^2 - k^2 + 2b$.

or $E_w(i+1) - E_w(i) \leq -k^2 + 2b$.

So, the error reduces by $-k^2 + 2b$ when $b$ is the margin.

or $\underline{\varepsilon}(i+1) - \underline{\varepsilon}(i) \leq -k^2 + b$ (approx).

Apply forcing argument :

$$0 \leq \varepsilon_{\underline{w}}(i+1) \leq \varepsilon_{\underline{w}}(i) - k^2 + b \leq 0$$

For some $i_0$, we'd have,

$$0 \leq \underline{\varepsilon}(i_0+1) \leq \underline{\varepsilon}(i_0) - k^2 + b < 0$$

$\longrightarrow$ impossible

Iterations must cease at $i_0-1$ or sooner //

Hence, algorithm converges at a solution, vector at $(i_0-1)^{th}$ iteration

weight

or sooner.

The 2-class perception with margin converges for linearly seperable data (training) //