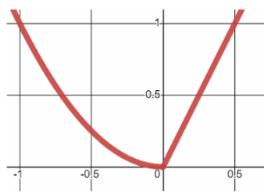1. Let $f(x) = \begin{cases} x^2, & x < 0 \\ 2x, & x \geq 0 \end{cases}$.

   1.1. Is $f$ convex? No need to explain.

   1.2. Propose a sub-derivative function $g$ for $f$. That is, $g \in \partial f$.

   Use the above definition to prove that $g(u) \in \partial f(u), \forall u \in \mathbb{R}$.

   1.3. Set a learning rate of $\eta = 0.25$ and a starting point $x_0 = -1.5$.

   Q

   Running subgradient descent, will the algorithm converge to a minimum?

   Prove your answer by filling the following table like we did in Tutorial 07 using as many rows as needed.

| i | $x_i$ | $f(x_i)$ | $\frac{\partial}{\partial x} f(x_i) = g(x_i)$ |
|---|-------|----------|-----------------------------------------------|
| 0 | −1 | 1 | |
| 1 | | | |
| ⋮ | | | |

   1.4. Repeat 1.3 with $\eta = 1$, $x_0 = -1.5$.

1.1. yes

1.2. $g(u) = \begin{cases} 2u, & u < 0 \\ 2, & u \geq 0 \end{cases}$

$\underline{u<0}$: $f(v) \geq f(u) + 2u(v-u)$
- $\underline{v<0}$: $v^2 \geq u^2 + 2uv - 2u^2$
  
  $v^2 - 2uv + u^2 \geq 0$
  
  $(v-u)^2 \geq 0$ which holds always
- $\underline{v \geq 0}$: $2v \geq u^2 + 2u(v-u)$

  $u^2 + 2v - 2uv \geq u^2 + 2v \geq 0$    which always holds
  
  $\underset{\uparrow}{-uv \geq 0}$    $\underset{\uparrow}{}$ from u,v definition

$\underline{u \geq 0}$: $f(v) \geq f(u) + 2(v-u)$
- $\underline{v<0}$: $v^2 \geq 2u + 2v - 2u$
  
  $v^2 - 2v \geq 0$    which always holds since $v^2 > 0$, $-2v > 0$
- $\underline{v \geq 0}$: $2v \geq 2u + 2v - 2u$
  
  $0 \geq 0$    which is always true

1.3. $f(x) = \begin{cases} x^2, & x < 0 \\ 2x, & x \geq 0 \end{cases}$   $\nabla f(x) = \begin{cases} 2x, & x < 0 \\ 2, & x \geq 0 \end{cases}$

$x_{i+1} = x_i - \eta \cdot f'(x_i)$

$x_0 = -1.5$     ; $f(x_0) = \frac{9}{4}$ ; $\nabla f(x_0) = -3$

$x_1 = -1.5 - \frac{1}{4} \cdot (-3) = -\frac{3}{4}$ ; $f(x_1) = \frac{9}{16}$ ; $\nabla f(x_1) = -\frac{3}{2}$

$x_2 = -\frac{3}{4} - \frac{1}{4} \cdot (-\frac{3}{2}) = -\frac{3}{8}$ ; $f(x_2) = \frac{9}{64}$ ; $\nabla f(x_2) = -\frac{3}{4}$

$x_3 = -\frac{3}{8} - \frac{1}{4}(-\frac{3}{4}) = -\frac{3}{16}$ ; $f(x_3) = \frac{9}{256}$ ; $\nabla f(x_3) = -\frac{3}{8}$

we can see that $x_i \xrightarrow{i \to \infty} 0$ ; $f(x_i) \xrightarrow{i \to \infty} 0$ and $\nabla f(x_i) \xrightarrow{i \to \infty} 0$

1.4   $x_0 = -1.5$     ; $f(x_0) = \frac{9}{4}$ ; $\nabla f(x_0) = -3$

$x_1 = -1.5 - 1 \cdot (-3) = 1.5$ ; $f(x_1) = 3$ ; $\nabla f(x_1) = 2$

$x_2 = 1.5 - 1 \cdot 2 = -\frac{1}{2}$    ; $f(x_2) = \frac{1}{4}$ ; $\nabla f(x_2) = -1$

$x_3 = -\frac{1}{2} - 1 \cdot (-1) = \frac{1}{2}$   ; $f(x_3) = 1$ ; $\nabla f(x_3) = 2$

$x_4 = \frac{1}{2} - 1 \cdot 2 = -1.5$ ; $f(x_4) = \frac{9}{4}$ ; $\nabla f(x_4) = -3$

$x_0 = x_4 \Rightarrow$ we'll be in an infinite non-converging loop

Nick   322315318

2. This exercise will investigate the regularization coefficient $\lambda$ as it was presented in the ridge linear regression section of this course. Suppose we are trying to fit a polynomial to the following data:

| X | Y |
|---|---|
| 0 | 0 |
| 1 | 3 |
| 2 | 12 |

Our hypothesis class for this problem will be
$$\mathcal{H} = \{w_0 + w_1 x + w_2 x^2 + w_3 x^3 : (w_0, w_1, w_2, w_3) \in \mathbb{R}^4\}.$$

2.1. Show that we can fit the data with $w_0 = 0, w_1 = 2, w_2 = 0, w_3 = 1$.

2.2. Show that our hypothesis class is too expressive for the problem we're dealing with. In other words, find a simple quadratic polynomial that fits the data perfectly.

**2.1**   $0 + 2x + 0 + x^3 = x^3 + 2x$,

$x = 0 \to y = 0$
$x = 1 \to y = 3$
$x = 2 \to y = 12$

**2.2**   $3x^2$

$x = 0 \to y = 0$
$x = 1 \to y = 3$
$x = 2 \to y = 12$

2.3. Denote the mean squared error (MSE)
$$\mathcal{L}(w) = \frac{1}{m}\|Xw - y\|_2^2,$$
Where $X$ is the appropriate Vandermonde matrix.
Calculate $\mathcal{L}(w)$ for the quadratic model in (2.2) and the cubic model in (2.1).

$$X = \begin{bmatrix} 0^0 & 0^1 & 0^2 & 0^3 \\ 1^0 & 1^1 & 1^2 & 1^3 \\ 2^0 & 2^1 & 2^2 & 2^3 \end{bmatrix} = \begin{bmatrix} 1 & 0 & 0 & 0 \\ 1 & 1 & 1 & 1 \\ 1 & 2 & 4 & 8 \end{bmatrix}$$

$$\mathcal{L}(w_{2.1}) = \left\|\begin{bmatrix} 0 & 3 & 12\end{bmatrix} - \begin{bmatrix} 0 & 3 & 12\end{bmatrix}\right\|_2^2 = 0$$

$$\mathcal{L}(w_{2.2}) = \left\|\begin{bmatrix} 0 & 3 & 12\end{bmatrix} - \begin{bmatrix} 0 & 3 & 12\end{bmatrix}\right\|_2^2 = 0$$

$$W_{2.1} = \begin{bmatrix} 0 \\ 2 \\ 0 \\ 1 \end{bmatrix} \quad W_{2.2} = \begin{bmatrix} 0 \\ 0 \\ 3 \\ 0 \end{bmatrix} \quad y = \begin{bmatrix} 0 & 3 & 12\end{bmatrix}$$

2.4. The best <u>line</u> for fitting the data is $y = 6x - 1$. Calculate $\mathcal{L}(w)$ for this line.

$$W_{2.4} = \begin{bmatrix} -1 \\ 6 \\ 0 \\ 0 \end{bmatrix} \quad \mathcal{L}(w_{2.4}) = \frac{1}{3}\left\|\begin{bmatrix} -1 & 5 & 11\end{bmatrix} - \begin{bmatrix} 0 & 3 & 12\end{bmatrix}\right\|_2^2 = \frac{1}{3}\left\|\begin{bmatrix} -1 & 2 & -1\end{bmatrix}\right\|_2^2 = \frac{1}{3}\cdot 6 = 2$$

2.5. Now denote the MSE with regularization as show in class
$$\mathcal{L}_\lambda(w) = \frac{1}{m}\|Xw - y\|_2^2 + \lambda\|w\|_2^2.$$
Here $\lambda > 0$ is a hyperparameter, which is not given. As we learned in class, the regularization imposes a "cost" on models with large coefficients. Calculate $\mathcal{L}_\lambda(w)$ for each of the three models in (2.1), (2.2) and (2.4).

**2.5.1**   $\mathcal{L}_\lambda(w_{2.1}) = 5\lambda$
**2.5.2**   $\mathcal{L}_\lambda(w_{2.2}) = 9\lambda$
**2.5.4**   $\mathcal{L}_\lambda(w_{2.4}) = 2 + 37\lambda$

2.6. As it turns out, $\mathcal{L}_\lambda(w)$ would never prefer the simple quadratic polynomial over the cubic polynomial we found, no matter the value of $\lambda > 0$. Can you explain why?
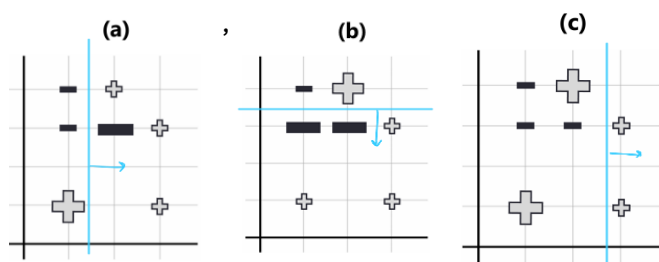
as explained in 2.5, the regularization imposes a "cost" on models with large coefficients therefore, in our example we'd always prefer the cubic polynomial since the $\|\cdot\|_2^2$ of its coefficients is smaller.

2.7. Suggest a way to fix the regularization method to prefer the model we consider to be simpler.

choose $\lambda$ to be a vector that punishes $x^3$ and rewards $x^2$, so we can choose something like: $\lambda = [0,0,0,1]$ and then $\mathcal{L}_\lambda(w) = \frac{1}{m}\|Xw - y\|_2^2 + \sum_{i=0}^{3}\lambda_i w_i^2$

3. Only some of the following figures depict possible distributions that can be obtained after <u>one</u> iteration of AdaBoost. **Which ones?** For each such distribution, propose a weak classifier that can lead to its figure (use a <u>clear</u> drawing or a short description of that classifier).

only figures (a), (b) and (c) can be obtained after one iteration of AdaBoosed

**(a)** , **(b)** **(c)**



(the arrow points towards the side that is supposed to be labeled as positive (+))

4. Show that the error of $h_t$ w.r.t the distribution $D^{t+1}$ is exactly $1/2$. That is, show that $\forall t \in [T]$

$$\sum_{i=1}^{m} D_i^{t+1}\, \mathbb{1}_{[y_i \neq h_t(x_i)]} = 1/2.$$

<u>first way:</u>
prove using induction:

Base: $t=1$: without limitation of generality, let $0 \leq k \leq m$ be the amount of wrongly classified points by the chosen model in $t=0$

$$\varepsilon_1 = \sum_{i=1}^{m} D_{(i)}^1 = \frac{1}{m}\cdot k = \frac{k}{m} \quad ; \quad \alpha_1 = \frac{1}{2}\log\left(\frac{1}{\varepsilon_0}-1\right) = \log\left(\sqrt{\frac{m-k}{k}}\right)$$

if wrongly classified: $D_i^{t+1} \propto \frac{1}{m} e^{\log\left(\sqrt{\frac{m-k}{k}}\right)} = \frac{1}{m}\sqrt{\frac{m-k}{k}} \Longrightarrow \sum_{i=1}^{m} D_i^{t+1}\mathbb{1}_{[y_i \neq t_r(x_i)]} = k\cdot\frac{1}{m}\sqrt{\frac{m-k}{k}} = \frac{\sqrt{k(m-k)}}{m}$

if rightly classified: $D_i^{t+1} \propto \frac{1}{m} e^{\log\left(\sqrt{\frac{k}{m-k}}\right)} = \frac{1}{m}\sqrt{\frac{k}{m-k}} \Longrightarrow \sum_{i=1}^{m} D_i^{t+1}\mathbb{1}_{[y_i=t_+(x_i)]} = (m-k)\cdot\frac{1}{m}\sqrt{\frac{k}{m-k}} = \frac{\sqrt{k(m-k)}}{m}$

as the sum of probabilities sum of both is equal to $1$ therefore each one is $1/2$

Assume correctness for $t'$;

Show correctness for $t+1$: $\sum_{i=1}^{m} D_i^{t+1}\mathbb{1}_{[y_i \neq t_+(x_i)]} = \sum_{i=1}^{m} \frac{D_i^t\, e^{-\frac{1}{2}\log\left(\frac{1}{\sum_i D_i^t \mathbb{1}_{[y_i \neq t_r(x_i)]}}-1\right)}}{\sum_{j=1}^{m} D_j^t\, e^{-\frac{1}{2}\log\left(\frac{1}{\sum_r D_r^t\mathbb{1}_{[y_r \neq t_r(x_r)]}}-1\right)}} \mathbb{1}_{[y_i \neq t_+(x_i)]} = \sum_{i=1}^{m} \frac{D_i^t\, e^0}{\sum_{j=1}^{m} D_j^t e^0}\mathbb{1}_{[y_i \neq t_+(x_i)]} =$

$$= \sum_{i=1}^{m} \underbrace{D_i^t}_{1}\, \mathbb{1}_{[y_i \neq h_+(x_i)]} = \frac{1}{2}$$

<u>second way:</u> $\sum_{i=1}^{m} D_i^{t+1}\mathbb{1}_{[N]} = \sum_{i=1}^{m} \frac{D_i^t\, e^{-\alpha_t y_i t_r(x_i)}}{Z_+}\mathbb{1}_{[N]} = \sum_{i=1}^{m} \frac{D_i^t\, e^{\alpha_t}}{Z_+}\mathbb{1}_{[N]} = \frac{1}{Z_+}\sum_{i=1}^{m} D_i^t\, e^{\alpha_t}\mathbb{1}_{[N]} = \frac{1}{2\sqrt{\varepsilon_+(1-\varepsilon_+)}}\sqrt{\varepsilon_+(1-\varepsilon_+)} = \frac{1}{2}$ as needed

$Z_+ = \sum_{i=1}^{m} D_i^t\mathbb{1}_{[N]}\, e^{\alpha_t} + \sum_{i=1}^{m} D_i^t\mathbb{1}_{[N]}\, e^{-\alpha_t} = e^{\alpha_t}\sum_{i=1}^{m} D_i^t\mathbb{1}_{[N]} + e^{-\alpha_t}\sum_{i=1}^{m} D_i^t\mathbb{1}_{[N]} = e^{\alpha_t}\varepsilon_+ + e^{-\alpha_t}(1-\varepsilon_+) = e^{\frac{1}{2}\log\left(\frac{1}{\varepsilon_+}-1\right)}\varepsilon_+ + e^{-\frac{1}{2}\log\left(\frac{1}{\varepsilon_+}-1\right)}(1-\varepsilon_+) = \sqrt{\frac{1}{\varepsilon_+}-1}\,\varepsilon_+ + \frac{1}{\sqrt{\frac{1}{\varepsilon_+}-1}}(1-\varepsilon_+) =$

$= \varepsilon_+\sqrt{\frac{1}{\varepsilon_+}-1} + \frac{1-\varepsilon_+}{\sqrt{\frac{1}{\varepsilon_+}-1}} = \sqrt{\varepsilon_+(1-\varepsilon_+)} + \sqrt{\varepsilon_+(1-\varepsilon_+)} = 2\sqrt{\varepsilon_+(1-\varepsilon_+)}$

(***) $\sum_{i=1}^{m} D_i^t\, e^{\alpha_t}\mathbb{1}_{[N]}$      (**) $[N] = [y_i \neq t_+(x_i)]$    $[\bar N] = [y_i = h_+(x_i)]$

5. Prove that $\forall \eta > 0$ the perceptron algorithm will perform the same number of iterations, and will converge to a vector that points to the same direction.

first, we'll show that the perceptron converges to a vector that points in the same direction

let $\eta_1, \eta_2 > 0$. if the data is linearly isepperable, the algorithm won't converge for either of the $\eta$'s

otherwise it'll converge after $n \in \mathbb{N}$ steps; $W_n = W_{n-1} + \eta\, y_{i_n}\cdot x_{i_n} = \ldots = W_0 + \eta\, y_{i_1}\cdot x_{i_1} + \ldots + \eta\, y_{i_n}\cdot x_{i_n} = \eta\sum_{k=1}^{n} y_{i_k} x_{i_k}$

where $i_k$ are the indexes on which the preceptron made a mistake in the labeling

we see that for each $\eta$ we'll get the same answer except the coefficient $\eta$

i.e. $\eta_1\sum_{k=1}^{n} y_{i_k} x_{i_k}$, $\eta_2\sum_{k=1}^{n} y_{i_k} x_{i_k}$ we can see that $\eta$ only scales the sum (its a scalar), therefore it converges

to a vector that points in the same direction

secondly, we'll show that the preceptron algorithm performs the same number of iterations

$\hat{y}_j = \text{Sign}(\langle w, x_j\rangle) \overset{W_{i \leq n}}{=} \text{Sign}\left(\langle \eta\sum y_{i_k} x_{i_k}, x_j\rangle\right) \overset{\eta > 0}{=} \text{Sing}\left(\langle \sum_{k=1}^{n} y_{i_k} x_{i_k}, x_j\rangle\right)$

therefore, we get the same missclassifications, no matter the $\eta$ => we'll have the same number of iterations