

1. Assume we have a dataset $\{(x(1), y(1)), (x(2), y(2)), \dots, (x(n), y(n))\}$ of points in the 2D plane and we want to estimate the following **linear regression** model: $\hat{y} = w_1 \cdot x + w_0$.
Compute w_0, w_1 in terms of $x(1), y(1), \dots, x(n), y(n)$. Numerical application $(1, 1), (2, -1), (3, 2)$.

Assume we have a dataset $\{(x^{(1)}, y^{(1)}), (x^{(2)}, y^{(2)}), \dots, (x^{(n)}, y^{(n)})\}$ of points in the 2D plane, and we want to estimate the following linear regression model:

$$\hat{y} = w_1 x + w_0.$$

Compute w_0, w_1 in terms of $x^{(1)}, y^{(1)}, x^{(2)}, y^{(2)}, \dots, x^{(n)}, y^{(n)}$. Numerical application: $(1, 1), (2, -1), (3, 2)$.

Linear Regression

- Assume we have a dataset $\{(x^{(1)}, y^{(1)}), (x^{(2)}, y^{(2)}), \dots, (x^{(m)}, y^{(m)})\}$ of points in the 2D plane and we want to estimate the following linear regression model:

$$\hat{y} = w_1 x + w_0$$

Compute w_0, w_1 in terms of $x^{(1)}, y^{(1)}, x^{(2)}, y^{(2)}, \dots, x^{(m)}, y^{(m)}$. Numerical application: $(1, 1), (2, -1), (3, 2)$.

$$m=3$$

$$x = \begin{bmatrix} 1 & 1 \\ 1 & 2 \\ 1 & 3 \end{bmatrix}, \quad y = \begin{bmatrix} 1 \\ -1 \\ 2 \end{bmatrix}$$

$$w^* = (x^T \cdot x)^{-1} \cdot x^T \cdot y$$

$$x^T \cdot x = \begin{bmatrix} 1 & 1 & 1 \\ 1 & 2 & 3 \end{bmatrix} \cdot \begin{bmatrix} 1 & 1 \\ 1 & 2 \\ 1 & 3 \end{bmatrix} = \begin{bmatrix} 3 & 6 \\ 6 & 14 \end{bmatrix}$$

$$(x^T \cdot x)^{-1}$$

$$A^{-1} = \frac{1}{\det(A)} \cdot A^*$$

$$A^* = \begin{bmatrix} \det(a_{22}) & -\det(a_{12}) \\ -\det(a_{21}) & \det(a_{11}) \end{bmatrix}$$

$$\det(A) = \begin{vmatrix} 3 & 6 \\ 6 & 14 \end{vmatrix} = 42 - 36 = 6$$

$$A^T = \begin{bmatrix} 3 & 6 \\ 6 & 14 \end{bmatrix} \rightarrow A^* = \begin{bmatrix} 14 & -6 \\ -6 & 3 \end{bmatrix}$$

$$A^{-1} = \frac{1}{6} \cdot \begin{bmatrix} 14 & -6 \\ -6 & 3 \end{bmatrix}$$

! la matricele 2×2 A^* - se înverzesc elementele de pe diagonala principală, iar cele de pe diagonala secundară se negăsesc

$$w = \frac{1}{6} \cdot \begin{bmatrix} 14 & -6 \\ -6 & 3 \end{bmatrix} \cdot \begin{bmatrix} 1 & 1 & 1 \\ 1 & 2 & 3 \end{bmatrix} \cdot \begin{bmatrix} 1 \\ -1 \\ 2 \end{bmatrix} = \frac{1}{6} \cdot \begin{bmatrix} -2 \\ 3 \end{bmatrix} = \begin{bmatrix} -\frac{1}{3} \\ \frac{1}{2} \end{bmatrix} \begin{matrix} w_0 \\ w_1 \end{matrix}$$

+ error

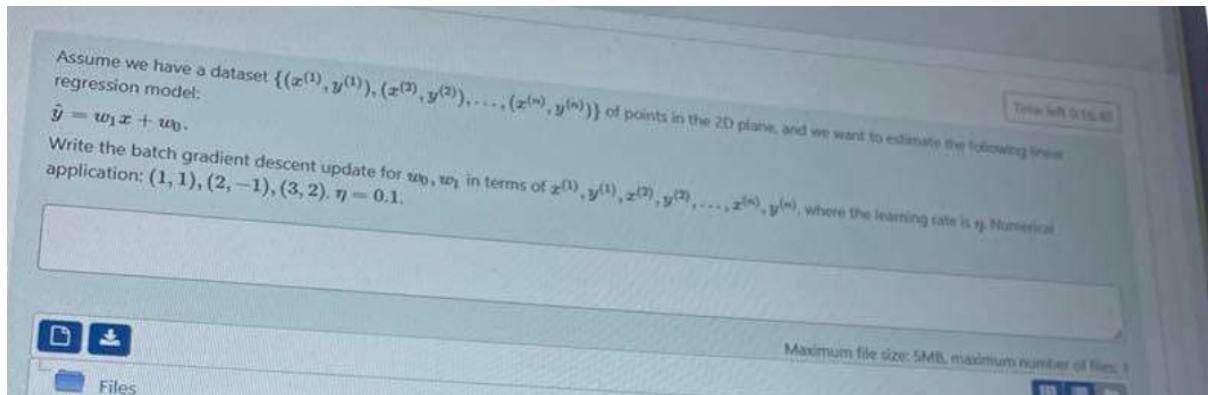
$$\text{err} = \sum (y_i - \hat{y}_i)^2 \quad \leftarrow \text{error}$$

$$\boxed{\hat{y} = x \cdot w} = \begin{pmatrix} 1 & 1 \\ 1 & 2 \\ 1 & 3 \end{pmatrix} \cdot \begin{pmatrix} -\frac{1}{3} \\ \frac{1}{2} \end{pmatrix} = \begin{pmatrix} \frac{1}{6} \\ \frac{2}{3} \\ \frac{7}{6} \end{pmatrix}$$

$$\text{err} = \left(1 - \frac{1}{6}\right)^2 + \left(-1 - \frac{2}{3}\right)^2 + \left(2 - \frac{7}{6}\right)^2 = \frac{25}{6} = 4,166$$

2. Assume we have a dataset $\{(x(1), y(1)), (x(2), y(2)), \dots, (x(n), y(n))\}$ of points in the 2D plane and we want to estimate the following linear regression model: $\hat{y} = w_1 x + w_0$.

Write the batch **gradient descent** update for w_0, w_1 in terms of $x(1), y(1), \dots, x(n), y(n)$ where the following rate is n . Numerical application: $(1, 1), (2, -1), (3, 2)$, $n=0.1$.



Gradient descent

- 2) Assume we have a dataset $(x^{(1)}, y^{(1)}), (x^{(2)}, y^{(2)}), \dots, (x^{(m)}, y^{(m)})$ of points in the 2D plane, and we want to estimate the following linear regression model:

$$\hat{y} = w_1 x + w_0$$

Write the batch gradient descent update for w_0, w_1 in terms of $x^{(1)}, y^{(1)}$, $x^{(2)}, y^{(2)}, \dots, x^{(m)}, y^{(m)}$ where the learning rate is η . Numerical application: $(1, 1), (2, -1), (3, 2)$, $\eta = 0.1$.

$$\frac{1}{m} \cdot x^T \cdot (xw - y) = \nabla_w L(w) \Rightarrow m=3$$

$$x = \begin{bmatrix} 1 & 1 \\ 1 & 2 \\ 1 & 3 \end{bmatrix}, \quad y = \begin{bmatrix} 1 \\ -1 \\ 2 \end{bmatrix}$$

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

$$w \leftarrow w - \eta \nabla_w L(w)$$

$$w \leftarrow \begin{bmatrix} w_0 \\ w_1 \end{bmatrix} - \frac{0.1}{3} \cdot x^T \cdot \begin{pmatrix} w_0 + w_1 \\ w_0 + 2w_1 \\ w_0 + 3w_1 \end{pmatrix} - \begin{pmatrix} 1 \\ -1 \\ 2 \end{pmatrix}$$

$$\begin{bmatrix} w_0 \\ w_1 \end{bmatrix} \leftarrow \begin{bmatrix} w_0 \\ w_1 \end{bmatrix} - \frac{0.1}{3} \cdot x^T \cdot \begin{bmatrix} w_0 + w_1 - 1 \\ w_0 + 2w_1 + 1 \\ w_0 + 3w_1 - 2 \end{bmatrix}$$

$$\begin{bmatrix} w_0 \\ w_1 \end{bmatrix} \leftarrow \begin{bmatrix} w_0 \\ w_1 \end{bmatrix} - \frac{0.1}{3} \begin{bmatrix} 3w_0 + 6w_1 - 2 \\ 6w_0 + 14w_1 - 5 \end{bmatrix}$$

$$\begin{cases} w_0 \leftarrow w_0 - 0.1w_0 - 0.2w_1 + \frac{0.2}{3} \\ w_1 \leftarrow w_1 - 0.2w_0 - \frac{1.4}{3}w_1 + \frac{0.5}{3} \end{cases}$$

3. Assume we have a training set for two-class classification in two dimensions that contains seven sample points: points $(1, 1), (2, 2), (3, 1)$ with label +1 and points $(4, 5), (5, 7), (6, 5), (7, 7)$ with label -1. Which points are the support vectors for a **hard margin SVM**?

Assume we have a training set for two-class classification in two dimensions that contains seven sample points: points $(1, 1), (2, 2), (3, 1)$ with label +1 and points $(4, 5), (5, 7), (6, 5), (7, 7)$ with label -1. Which points are the support vectors for a hard-margin SVM?

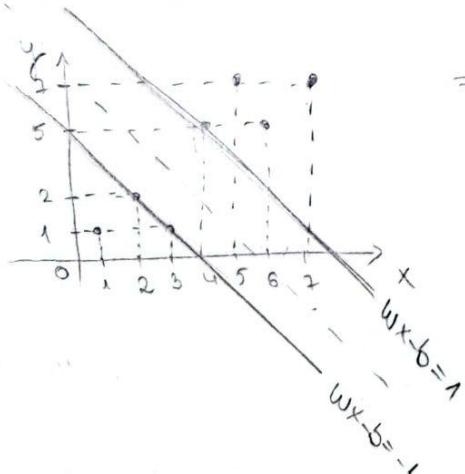
SVM

- 3) Assume we have a training set for two-class classification in two dimensions that contains seven sample points: points $(1, 1), (2, 2), (3, 1)$ with label +1 and points $(4, 5), (5, 7), (6, 5), (7, 7)$ with label -1. Which points are the support vectors for a hard-margin SVM?

- punctele care sunt pe marginie: support vectors (grafic)
 $(3; 3,5) \rightarrow$ mij punctelor $(4,5), (2,2)$

$$\min_{w, w_0} \frac{1}{2} \|w\|_2^2$$

$$y^{(i)} (w^\top x^{(i)} + w_0) \geq 1, \quad 1 \leq i \leq m$$



\Rightarrow punctele care pot fi support vectors
 sunt: $(2, 2), (3, 1)$ și $(4, 5)$

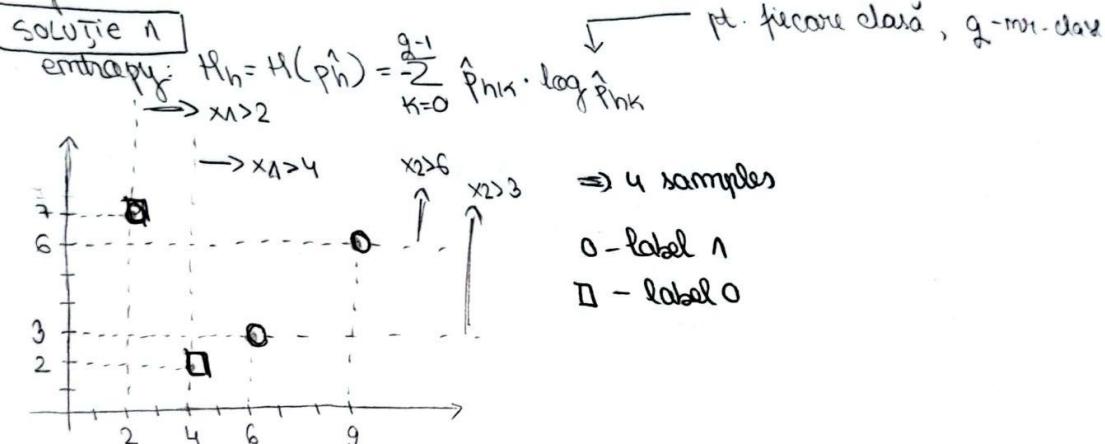
$\Rightarrow (3; 3,5)$

4. Consider a training classification **decision tree** given a design matrix $X = [6 \ 3; 2 \ 7; 9 \ 6; 4 \ 2]$ and labels $y = [1; 0; 1; 0]$. Let x_1 denote feature 1, corresponding to the first column of X and let x_2 denote feature 2, corresponding to the second column of X , Which of the following splits at the root node gives the smallest entropy: $x_1 > 2$, $x_1 > 4$, $x_2 > 3$, $x_2 > 6$? Motivate your answer

Consider training a classification decision tree given a design matrix $X = \begin{bmatrix} 6 & 3 \\ 2 & 7 \\ 9 & 6 \\ 4 & 2 \end{bmatrix}$ and labels $y = \begin{bmatrix} 1 \\ 0 \\ 1 \\ 0 \end{bmatrix}$. Let x_1 denote feature 1, corresponding to the first column of X , and let x_2 denote feature 2, corresponding to the second column of X . Which of the following splits at the root node gives the smallest entropy: $x_1 > 2$, $x_1 > 4$, $x_2 > 3$, $x_2 > 6$? Motivate your answer.

Decision tree

- 4) Consider a binomial classification decision tree given a design matrix $X = \begin{bmatrix} 6 & 3 \\ 2 & 7 \\ 9 & 6 \\ 4 & 2 \end{bmatrix}$ and labels $y = \begin{bmatrix} 1 \\ 0 \\ 1 \\ 0 \end{bmatrix}$. Let x_1 denote feature 1, corresponding to the first column of X and let x_2 denote feature 2, corresponding to the second column of X . Which of the following splits at the root node gives the smallest entropy: $x_1 > 2$, $x_1 > 4$, $x_2 > 3$, $x_2 > 6$? Motivate your answer.



a) $x_1 > 2$

pt. ipămătătoare scrieră: $H_h = -\hat{p}_{11} \cdot \log \hat{p}_{11} - \hat{p}_{00} \cdot \log \hat{p}_{00}$

$$\hat{p}_{hk} = \frac{n}{D_h |I_h|} \quad \sum_{(x,y) \in I_h} 1 \cdot y(i) = K$$

- cete 0 sau 1 sunt în reg. h: 1

\rightarrow - cete 0 sau 1 sunt în reg. h și fac parte din cl. 1: 0

$$\hat{p}_{11} = \frac{0}{4} = 0$$

$$\hat{p}_{00} = \frac{4}{4} = 1$$

$$H_h = -0 \cdot \log 0 - 1 \cdot \log 1 = -1 \cdot 0 = 0$$

pt. jumătatea dreaptă: $H_2 = -\hat{P}_{20} \cdot \log \hat{P}_{20} - \hat{P}_{21} \cdot \log \hat{P}_{21}$

$$\hat{P}_{20} = \frac{1}{3}, \quad \hat{P}_{21} = \frac{2}{3}$$

$$H_2 = -\frac{1}{3} \cdot \underset{\downarrow \text{log}}{\log \frac{1}{3}} - \frac{2}{3} \cdot \underset{\downarrow \text{log}}{\log \frac{2}{3}} = 0,16 + 0,11 = 0,27$$

b) $x_1 > 4$

jum. stânga: $H_3 = -\hat{P}_{30} \cdot \log \hat{P}_{30} - \hat{P}_{31} \cdot \log \hat{P}_{31}$

$$\hat{P}_{30} = \frac{2}{2} = 1$$

$$\hat{P}_{31} = 0 \Rightarrow H_3 = 0$$

jum. dreaptă: $H_4 = -\hat{P}_{40} \cdot \log \hat{P}_{40} - \hat{P}_{41} \cdot \log \hat{P}_{41}$

$$\hat{P}_{40} = 0$$

$$\hat{P}_{41} = 1 \Rightarrow H_4 = 0$$

c) $x_2 > 3$

jum. jos:

$$\hat{P}_{50} = \frac{1}{2}$$

$$\hat{P}_{51} = \frac{1}{2} \Rightarrow H_5 = -\frac{1}{2} \log \frac{1}{2} - \frac{1}{2} \log \frac{1}{2} = 0,3$$

jum. sus:

$$\hat{P}_{60} = \frac{1}{2}$$

$$\hat{P}_{61} = \frac{1}{2} \Rightarrow H_6 = 0,3$$

d) $x_2 > 6$

jum. jos:

$$\hat{P}_{10} = \frac{1}{3}$$

$$\hat{P}_{11} = \frac{2}{3} \Rightarrow H_1 = 0,27$$

jum. sus:

$$\hat{P}_{20} = 1$$

$$\hat{P}_{21} = 0 \Rightarrow H_2 = 0$$

\Rightarrow Concluzie

b) \rightarrow entropia era mai mică

$\Rightarrow x_1 > 4$ cel mai bine reprezentat

Soluție 2

$$x = \begin{pmatrix} 6 & 3 \\ 2 & 7 \\ 9 & 6 \\ 4 & 2 \end{pmatrix} \quad y = \begin{pmatrix} 1 \\ 0 \\ 1 \\ 0 \end{pmatrix}$$

$\uparrow \quad \uparrow$
 $x_1 \quad x_2$ features

Gini index:

$$G_h = \sum_{k=1}^2 p_{hk}^{\hat{1}} (1 - p_{hk}^{\hat{1}}) = \sum_{k=1}^2 p_{hk}^{\hat{1}} - \sum_{k=1}^2 p_{hk}^{\hat{1}} {}^2 = 1 - \sum_{k=1}^2 p_{hk}^{\hat{1}} {}^2,$$

unde $p_{hk}^{\hat{1}}$ = probabilitatea exemplilor din regiunea h care fac parte din clasa k

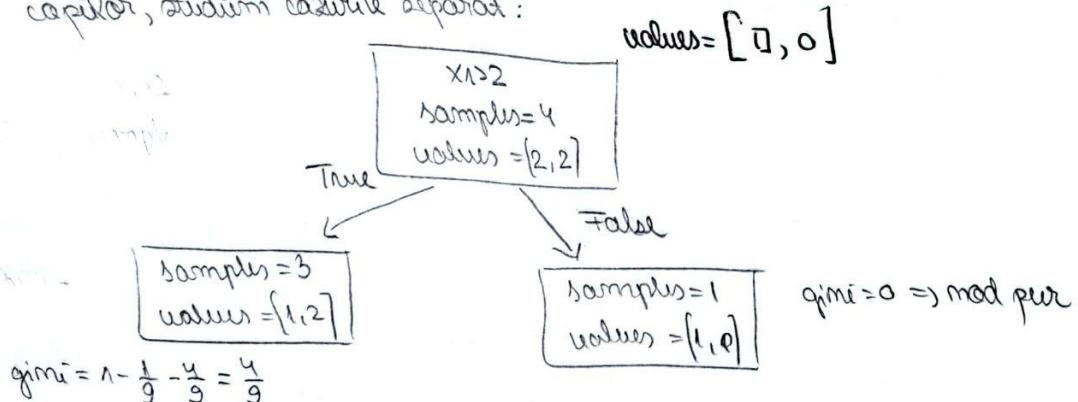
Gini index-ul la rădățea va rămâne aceeași, deoarece toate clasele sunt același.

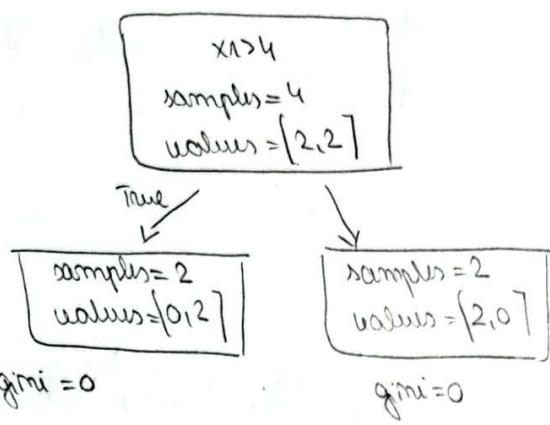
$$\begin{aligned} \text{samples} &= \text{nr. de puncte} \\ &= \text{nr. de linii} \\ &\text{dim } x \end{aligned}$$

$x_i > ?$
 samples = 4
 values = [2, 2]

 $\Rightarrow gini = 1 - \left(\frac{1}{2}\right)^2 - \left(\frac{1}{2}\right)^2 = \frac{1}{2} = 0,5$

Dacă întrebarea se referă la cel mai mic gini dintre cele ale mediilor capilor, studiul cazurilor separați:





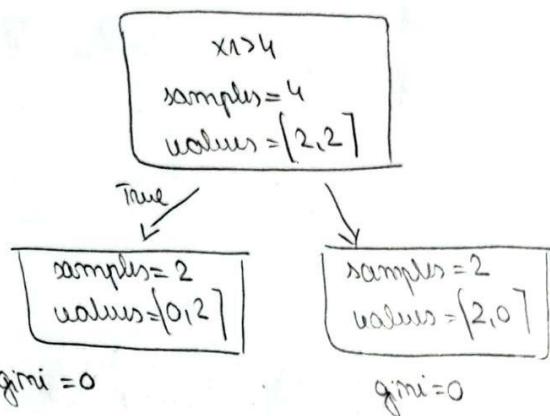
\Rightarrow We are most in control of which variant, because $x_1 > 4$ makes classification perfect, since each branch ends at a single impure class.

5. Assume we have the design matrix $X = [1 \ 0; 0 \ 1; -1 \ -1]$. Use **PCA** to reduce the dimension from 2 to 1 to compute the first principal component.

Assume we have the design matrix

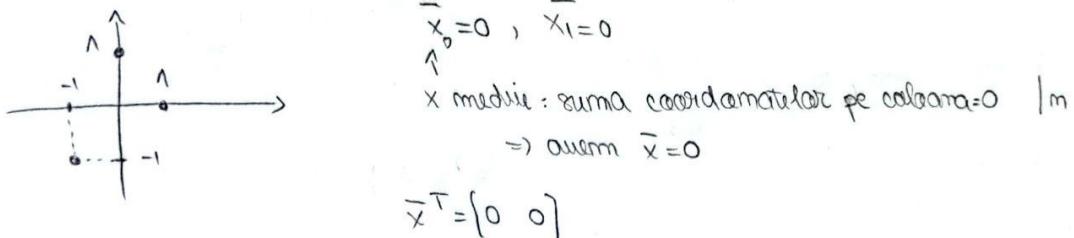
$$X = \begin{bmatrix} 1 & 0 \\ 0 & 1 \\ -1 & -1 \end{bmatrix}.$$

Use PCA to reduce the dimension from 2 to 1, i.e., to compute the first principal component.



⇒ Nu are existat în modul să se întâlnească valori de $x_1 > 4$ și $x_1 \leq 4$. Clasificarea este perfectă, fiindcă nu există exemplare care să fie clasificate corect.

- 5) Assume we have the design matrix $X = \begin{pmatrix} 1 & 0 \\ 0 & 1 \\ -1 & -1 \end{pmatrix}$. Use PCA to reduce the dimension from 2 to 1, i.e., to compute the first principal component.



$$A := \frac{1}{\sqrt{m}} \cdot X$$

$$m=3 \Rightarrow A := \frac{1}{\sqrt{3}} \begin{pmatrix} 1 & 0 \\ 0 & 1 \\ -1 & -1 \end{pmatrix}$$

$$\begin{aligned}
 \Sigma &= \frac{1}{m} \cdot \sum_{i=1}^m (x^{(i)} - \bar{x})(x^{(i)} - \bar{x})^T = \frac{1}{m} \sum_{i=1}^m x^{(i)} \cdot x^{(i)T} = \frac{1}{3} \cdot \sum_{i=1}^3 (x^{(i)} - \bar{x})(x^{(i)} - \bar{x})^T \\
 &= \frac{1}{3} \left[\left(\begin{pmatrix} 1 \\ 0 \\ -1 \end{pmatrix} - \bar{x} \right) \left(\begin{pmatrix} 1 \\ 0 \\ -1 \end{pmatrix} - \bar{x} \right)^T + \left(\begin{pmatrix} 0 \\ 1 \\ -1 \end{pmatrix} - \bar{x} \right) \left(\begin{pmatrix} 0 \\ 1 \\ -1 \end{pmatrix} - \bar{x} \right)^T + \left(\begin{pmatrix} -1 \\ -1 \\ 1 \end{pmatrix} - \bar{x} \right) \left(\begin{pmatrix} -1 \\ -1 \\ 1 \end{pmatrix} - \bar{x} \right)^T \right] \\
 &= \frac{1}{3} \left(\begin{pmatrix} 1 \\ 0 \\ -1 \end{pmatrix} \cdot \begin{pmatrix} 1 & 0 & -1 \end{pmatrix}^T + \begin{pmatrix} 0 \\ 1 \\ -1 \end{pmatrix} \cdot \begin{pmatrix} 0 & 1 & -1 \end{pmatrix}^T + \begin{pmatrix} -1 \\ -1 \\ 1 \end{pmatrix} \cdot \begin{pmatrix} -1 & -1 & 1 \end{pmatrix}^T \right) = \frac{1}{3} \left(\begin{pmatrix} 1 & 0 & -1 \\ 0 & 1 & -1 \\ -1 & -1 & 1 \end{pmatrix} + \begin{pmatrix} 0 & 0 & 0 \\ 0 & 0 & 0 \\ 0 & 0 & 0 \end{pmatrix} + \begin{pmatrix} 1 & 1 & 1 \\ 1 & 1 & 1 \\ 1 & 1 & 1 \end{pmatrix} \right) \\
 &= \frac{1}{3} \begin{pmatrix} 2 & 1 & 1 \\ 1 & 2 & 1 \\ 1 & 1 & 2 \end{pmatrix}
 \end{aligned}$$

rez

$$\text{când } \bar{x}=0 \Rightarrow \Sigma = \frac{1}{n} \cdot \bar{x}^T \cdot x$$

$$A^T \cdot A = \Sigma$$

$$\det(\Sigma - \lambda \cdot J_2) = 0$$

$$\left| \begin{pmatrix} \frac{2}{3} & \frac{1}{3} \\ \frac{1}{3} & \frac{2}{3} \end{pmatrix} - \begin{pmatrix} 2 & 0 \\ 0 & 2 \end{pmatrix} \right| = \left| \begin{pmatrix} \frac{2}{3}-2 & \frac{1}{3} \\ \frac{1}{3} & \frac{2}{3}-2 \end{pmatrix} \right| = \left(\frac{2}{3}-2 \right)^2 - \left(\frac{1}{3} \right)^2 = \left(\frac{2}{3}-2 \right)^2 - \frac{1}{9}$$

$$\Rightarrow \left(\frac{2}{3}-2 \right)^2 - \frac{1}{9} = 0 \Rightarrow \frac{4}{9} - \frac{4}{3}x + x^2 - \frac{1}{9} = 0 \mid :9 \Rightarrow 4 - 12x + 9x^2 - 1 = 0$$

$$\Rightarrow 9x^2 - 12x + 3 = 0 \Rightarrow 3x^2 - 4x + 1 = 0 \Rightarrow (2x-1)(3x-1) = 0$$

$$\Rightarrow \begin{cases} x_1 = 1 \\ x_2 = \frac{1}{3} \end{cases}, \text{ se potrăgăduiesc datorită}$$

Păstrăm atâtiva 2 căi din dim. cauză menținute $\Rightarrow \boxed{x_2=1}$ și în emună "from 2 to 1"

$$\Rightarrow \text{valo. } \Sigma - x_2 \cdot J = \begin{pmatrix} \frac{2}{3}-1 & \frac{1}{3} \\ \frac{1}{3} & \frac{2}{3}-1 \end{pmatrix} = \begin{pmatrix} -\frac{1}{3} & \frac{1}{3} \\ \frac{1}{3} & -\frac{1}{3} \end{pmatrix}$$

$$\begin{pmatrix} -\frac{1}{3} & \frac{1}{3} \\ \frac{1}{3} & -\frac{1}{3} \end{pmatrix} \begin{pmatrix} w_1 \\ w_2 \end{pmatrix} = \begin{pmatrix} 0 \\ 0 \end{pmatrix} \Rightarrow -\frac{1}{3}w_1 + \frac{1}{3}w_2 = 0 \Rightarrow w_1 = w_2 \Rightarrow \boxed{w_1=w_2} \text{ Ales}$$

$$\tilde{w}_1 = \begin{pmatrix} 1 \\ 1 \end{pmatrix}, w_1 = \begin{pmatrix} w_1 \\ w_2 \end{pmatrix} \quad y = \frac{w_2}{w_1} x$$

$$|\tilde{w}_1| = \sqrt{1^2 + 1^2} = \sqrt{2}$$

\uparrow
 w_1 w_2

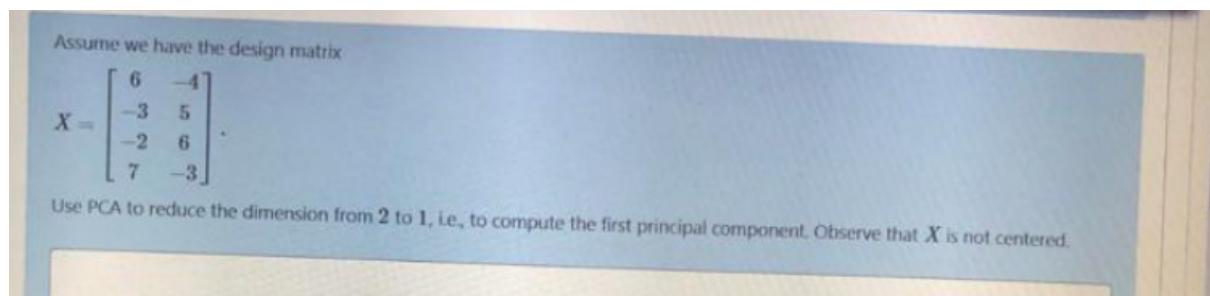
prima comp. principală

$$w_1 = \frac{1}{\sqrt{2}} \begin{pmatrix} 1 \\ 1 \end{pmatrix} = \begin{pmatrix} \frac{1}{\sqrt{2}} \\ \frac{1}{\sqrt{2}} \end{pmatrix}$$

$$\boxed{z = w^T \cdot x} \Rightarrow z^T = w^T \cdot x^T = \begin{pmatrix} \frac{1}{\sqrt{2}} & \frac{1}{\sqrt{2}} \end{pmatrix} \cdot \begin{pmatrix} 1 & 0 & -1 \end{pmatrix} = \begin{pmatrix} \frac{1}{\sqrt{2}} & \frac{1}{\sqrt{2}} & -\frac{1}{\sqrt{2}} \end{pmatrix}$$

$$w = \begin{pmatrix} \frac{1}{\sqrt{2}} \\ \frac{1}{\sqrt{2}} \end{pmatrix} \Rightarrow z = x^T = \begin{pmatrix} \frac{1}{\sqrt{2}} \\ \frac{1}{\sqrt{2}} \\ -\frac{1}{\sqrt{2}} \end{pmatrix}$$

6. Assume we have the design matrix $X = [6 -4; -3 5; -2 6; 7 -3]$. Use **PCA** to reduce the dimension from 2 to 1 to compute the first principal component. Observe that X is not centered.



6) Assume we have the design matrix $X = \begin{bmatrix} 6 & -4 \\ -3 & 5 \\ -2 & 6 \\ 7 & -3 \end{bmatrix}$. Use PCA to reduce the dimension from 2 to 1 to compute the first principal component. Observe that X is not centered.

$$\bar{x} = \begin{bmatrix} 8 & 4 \end{bmatrix} / 4 = \begin{bmatrix} 2 & 1 \end{bmatrix}$$

$$\Sigma = \frac{1}{4} \left[\left(\begin{bmatrix} 6 \\ -4 \end{bmatrix} - \begin{bmatrix} 2 \\ 1 \end{bmatrix} \right) \cdot \left(\begin{bmatrix} 6 \\ -4 \end{bmatrix} - \begin{bmatrix} 2 \\ 1 \end{bmatrix} \right)^T + \left(\begin{bmatrix} -3 \\ 5 \end{bmatrix} - \begin{bmatrix} 2 \\ 1 \end{bmatrix} \right) \cdot \left(\begin{bmatrix} -3 \\ 5 \end{bmatrix} - \begin{bmatrix} 2 \\ 1 \end{bmatrix} \right)^T + \left(\begin{bmatrix} -2 \\ 6 \end{bmatrix} - \begin{bmatrix} 2 \\ 1 \end{bmatrix} \right) \cdot \left(\begin{bmatrix} -2 \\ 6 \end{bmatrix} - \begin{bmatrix} 2 \\ 1 \end{bmatrix} \right)^T + \left(\begin{bmatrix} 7 \\ -3 \end{bmatrix} - \begin{bmatrix} 2 \\ 1 \end{bmatrix} \right) \cdot \left(\begin{bmatrix} 7 \\ -3 \end{bmatrix} - \begin{bmatrix} 2 \\ 1 \end{bmatrix} \right)^T \right] = \frac{1}{4} \left(\begin{bmatrix} 4 & 16 \\ 16 & 64 \end{bmatrix} + \begin{bmatrix} 121 & -11 \\ -11 & 1 \end{bmatrix} + \begin{bmatrix} 100 & -20 \\ -20 & 4 \end{bmatrix} + \begin{bmatrix} 1 & 1 \\ 1 & 9 \end{bmatrix} \right) = \frac{1}{4} \begin{bmatrix} 82 & -20 \\ -20 & 82 \end{bmatrix}$$

pt. write care: $\tilde{x} = x - \bar{x} = \begin{bmatrix} 6 & -4 \\ -3 & 5 \\ -2 & 6 \\ 7 & -3 \end{bmatrix} - \begin{bmatrix} 2 & 1 \\ 2 & 1 \\ 2 & 1 \\ 2 & 1 \end{bmatrix} = \begin{bmatrix} 4 & -5 \\ -5 & 4 \\ -4 & 5 \\ 5 & -4 \end{bmatrix}$

reduce
centered

$$\Sigma = \frac{1}{m} \cdot \tilde{x}^T \cdot \tilde{x}$$

$$\tilde{x} = \begin{bmatrix} 0 & 0 \end{bmatrix}^T$$

$$\det(\mathbb{I} - \lambda \cdot J_2) = 0 \Leftrightarrow \left| \begin{bmatrix} \frac{u_1}{2} & -u_0 \\ -u_0 & \frac{u_1}{2} \end{bmatrix} - \begin{bmatrix} 2 & 0 \\ 0 & 2 \end{bmatrix} \right| = 0 \Leftrightarrow \left| \begin{bmatrix} \frac{u_1}{2} - 2 & -u_0 \\ -u_0 & \frac{u_1}{2} - 2 \end{bmatrix} \right| = 0$$

$$\Leftrightarrow \left(\frac{u_1}{2} - 2 \right)^2 + u_0^2 = 0 \Leftrightarrow \frac{u_1^2}{4} + 2^2 - 4u_1 \cdot 2 + u_0^2 = 0 \mid :4 \Leftrightarrow u_1^2 + 4 \cdot 4 - 4u_1 \cdot 2 + u_0^2 \cdot 4 = 0$$

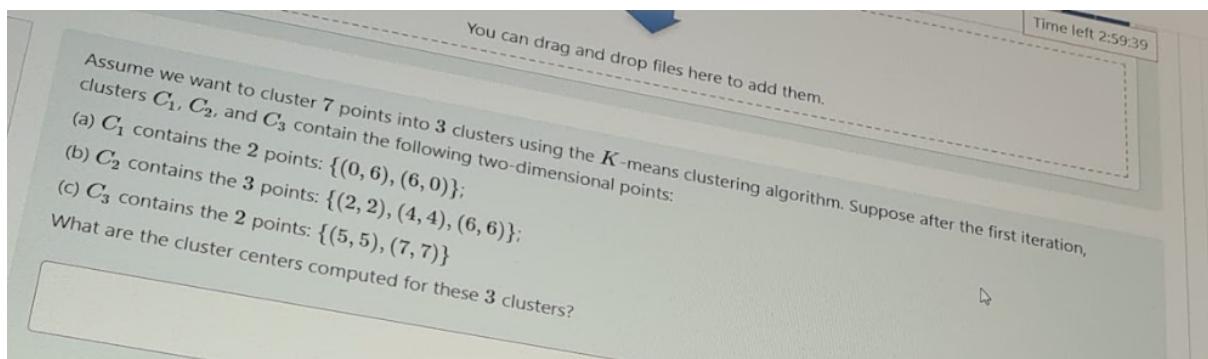
$$u_1^2 - 16u_1 + 4(u_1^2 - u_0^2) = 0 \Leftrightarrow u_1^2 - 16u_1 + 4 \cdot 81 = 0 \mid :4$$

$$\Leftrightarrow 2^2 - 4u_1 + 81 = 0 \Rightarrow \begin{cases} u_1 = \frac{4u_1 + \sqrt{1357}}{2} \\ u_2 = \frac{4u_1 - \sqrt{1357}}{2} \end{cases} \quad (\text{negative discriminant})$$

7. Assume we want to cluster 7 points into 3 clusters using the **K-means clustering algorithm**. Suppose after the first iteration, clusters C₁, C₂ and C₃ contain the following two-dimensional points:

- a. C₁ contains the 2 points: {(0, 6), (6, 0)}
- b. C₂ contains the 3 points: {(2, 2), (4, 4), (6, 6)}
- c. C₃ contains the 2 points: {(5, 5), (7, 7)}

What are the cluster centers computed for these 3 clusters.



*) Assume we want to cluster 7 points into 3 clusters using the K-means clustering algorithm. Suppose after the first iteration, clusters C₁, C₂, and C₃ contain the following two-dimensional points:

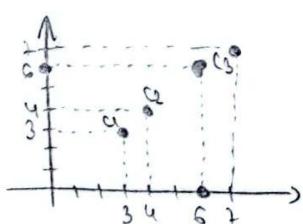
- a) C₁ contains the 2 points: {(0, 6), (6, 0)}
- b) C₂ contains the 3 points: {(2, 2), (4, 4), (6, 6)}
- c) C₃ contains the 2 points: {(5, 5), (7, 7)}

What are the cluster centers computed for these 3 clusters.

$$C_1: \frac{0+6}{2} = 3 \Rightarrow C_1(3, 3)$$

$$C_2: \frac{2+4+6}{3} = \frac{12}{3} = 4 \Rightarrow C_2(4, 4)$$

$$C_3: \frac{5+7}{2} = 6 \Rightarrow C_3(6, 6)$$



xt. inca o ieratice

$$C_3 = (2, 2), (6, 6), (5, 5) \Rightarrow C_3: \frac{2+6+5}{3} = 6 \Rightarrow C_3(6, 6)$$

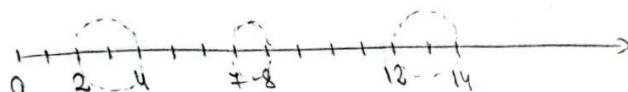
$$C_2 = (4, 4) \Rightarrow C_2(4, 4)$$

$$C_1 = (0, 6), (6, 0), (2, 2) \Rightarrow C_1: \frac{0+6+2}{3} = \frac{8}{3} \Rightarrow C_1(2, 6); \\ (2, 6)$$

8. Consider a dataset containing six one-dimensional points: {2, 4, 7, 8, 12, 14}. After three iterations of hierarchical clustering using **Euclidean distance between points**, we get the 3 clusters: $C_1 = \{2, 4\}$, $C_2 = \{7, 8\}$ and $C_3 = \{12, 14\}$. What is the distance between clusters C_1 and C_2 using single linkage? What is the distance between clusters C_1 and C_2 using complete linkage?

Consider a dataset containing six one-dimensional points: {2, 4, 7, 8, 12, 14}. After three iterations of hierarchical clustering using Euclidean distance between points, we get the 3 clusters: $C_1 = \{2, 4\}$, $C_2 = \{7, 8\}$, and $C_3 = \{12, 14\}$. What is the distance between clusters C_1 and C_2 using single linkage? What is the distance between clusters C_1 and C_2 using complete linkage?

- 8) Consider a dataset containing six one-dimensional points: {2, 4, 7, 8, 12, 14}. After three iterations of hierarchical clustering using Euclidean distance between points, we get the 3 clusters: $C_1 = \{2, 4\}$, $C_2 = \{7, 8\}$ and $C_3 = \{12, 14\}$. What is the distance between clusters C_1 and C_2 using single linkage? What is the distance between clusters C_1 and C_2 using complete linkage?



$C_1 \rightarrow C_2$ single link

$$d(C_1, C_2)_{\text{single}} = 7 - 4 = 3$$

$$d(C_1, C_2)_{\text{complete}} = 8 - 2 = 6$$

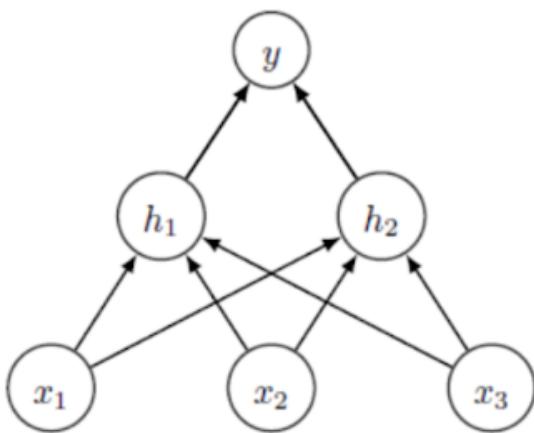
$$d(C_1, C_2)_{\text{average}} = \frac{3+4+5+6}{4} = 4.5$$

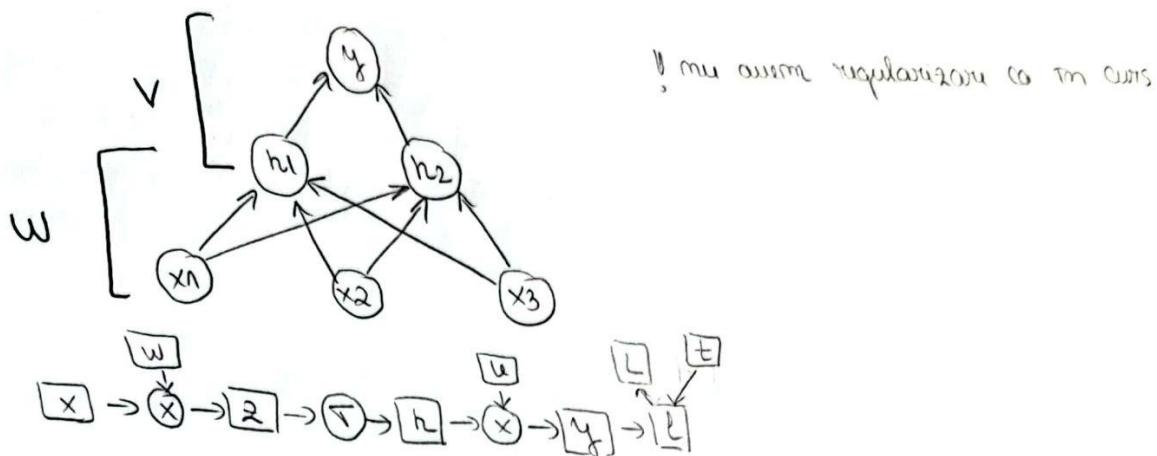
Jmč o iteracije fol. single

$C_1 + C_2 \Rightarrow$ cluster (distancia coa mai mica)

9. The following graph shows the structure of a **simple neural network with a single hidden layer**.

The following graph shows the structure of a simple neural network with a single hidden layer. The input layer consists of three dimensions $x = (x_1, x_2, x_3)$. The hidden layer includes two units $h = (h_1, h_2)$. The output layer includes one unit y . We ignore bias terms for simplicity. We use sigmoid $\sigma(z) = \frac{1}{1+e^{-z}}$ as activation function for the hidden layer and the identity for the output layer. Moreover, denote by $I(y, t) = \frac{1}{2}(y - t)^2$ the loss function. Here, t is the target value for the output unit y . Denote by W and V weight matrices connecting input and hidden layer, and hidden layer and output, respectively. Compute symbolically $\frac{\partial I}{\partial W}$. Numerical application:
 $W = \begin{bmatrix} 1 & 0 & 1 \\ 1 & -1 & 0 \end{bmatrix}, V = \begin{bmatrix} 0 & 1 \end{bmatrix}, x = (1, 2, 1), t = 1.$





$$H = w \star$$

$$y = v \nabla(H) = v \nabla(w \star)$$

$$\frac{\partial L(y, t)}{\partial t} = \frac{1}{2} (y - t)^2 = \frac{1}{2} (v \nabla(w \star) - t)^2 \Rightarrow \frac{\partial^2 L}{\partial y^2} = \frac{1}{2} \cdot 2 \cdot (v \nabla(w \star) - t) \cdot \nabla(w \star) \\ = (v \nabla(w \star) - t) \cdot \nabla(w \star)$$

Numeric: $H = w \star = \begin{pmatrix} 1 & 0 & 1 \\ 1 & -1 & 0 \end{pmatrix} \cdot \begin{pmatrix} 1 \\ 2 \\ 1 \end{pmatrix} = \begin{pmatrix} 2 \\ -1 \end{pmatrix}$

$$\nabla(H) = \begin{pmatrix} \frac{1}{1+e^{-2}} \\ \frac{1}{1+e} \end{pmatrix}$$

$$\Rightarrow y = (0 \ 1) \cdot \begin{pmatrix} \frac{1}{1+e^{-2}} \\ \frac{1}{1+e} \end{pmatrix} = \frac{1}{1+e} = [0, 269]$$

$$v^\top = \begin{bmatrix} 0 \\ 1 \end{bmatrix}$$

$$\frac{\partial L}{\partial w} = (v^\top \cdot ([0, 269] \wedge)) \odot \nabla'(2) \cdot x^\top$$

$$= ([0] \cdot [-0,731] \odot \nabla'(2)) \cdot x^\top$$

$$\nabla'(2) = \frac{e^{-2}}{(1+e^{-2})^2} = \begin{pmatrix} 0,105 \\ 0,192 \end{pmatrix}$$

$$\frac{\partial L}{\partial w} = \left(\begin{bmatrix} 0 \\ -0,731 \end{bmatrix} \odot \begin{pmatrix} 0,105 \\ 0,192 \end{pmatrix} \right) \cdot \begin{bmatrix} 1 \\ 2 \\ 1 \end{bmatrix} = \begin{bmatrix} 0 \\ -0,144 \end{bmatrix} \cdot \begin{bmatrix} 1 \\ 2 \\ 1 \end{bmatrix} = \begin{pmatrix} 0 \\ -0,144 \\ -0,288 \\ -0,144 \end{pmatrix}$$

10. Assume we have the following **CNN classifier**:

INPUT, CONV-9-32, POOL-2, CONV-5-64, POOL-2, FC-3

For each layer, calculate the number of weights, number of biases and the size of the associated feature maps. The size of INPUT is $128 \times 128 \times 3$. The notation follows the convention:

- CONV-K-N denotes a convolutional layer with N filters, each of them of size KxK, padding and stride parameters are always 0 and 1, respectively.
- POOL-K indicates a KxK pooling layer with stride K and padding 0.
- FC-N stands for a fully-connected layer with N neurons.

Assume we have the following CNN classifier:

INPUT, CONV-9-32, POOL-2, CONV-5-64, POOL-2, CONV-5-64, POOL-2, FC-3.

For each layer, calculate the number of weights, number of biases, and the size of the associated feature maps. The size of INPUT is $128 \times 128 \times 3$. The notation follows the convention:

- CONV-K-N denotes a convolutional layer with N filters, each of them of size KxK, padding and stride parameters are always 0 and 1, respectively.
- POOL-K indicates a KxK pooling layer with stride K and padding 0.
- FC-N stands for a fully-connected layer with N neurons.

CNN classifier

10) Assume we have the following CNN classifier:

INPUT, CONV-9-32, POOL-2, CONV-5-64, POOL-2, FC-3.

For each layer, calculate the number of weights, number of bias and the size of the associated feature maps. The size of INPUT is $128 \times 128 \times 3$. The notation follows the convention:

- CONV-K-N denotes a convolutional layer with N filters, each of kernel of size $K \times K$, padding and stride parameters are always 0 and 1
- POOL-K indicates a $K \times K$ pooling layer with stride K and padding 0
- FC-N stands for a fully-connected layer with N neurons.

W: ~~128x128x3~~

b: 3

$$(24-2+0+2)/2 = 12$$

$$(28-5+0+1)/1 = 11$$

$$11 \cdot 102400 = 1126400$$

b: 64 - output

$$(56-2+0+2)/2 = 28$$

$$(80-5+0+1)/1 = 64$$

weights: ~~64-32-5-1~~

$$\text{bias: } 64 = 5200$$

$$(120-2+0+2)/2 = 60$$

weights: ~~60-32-2-2~~

$$\text{bias: } 60 = 7680$$

input size

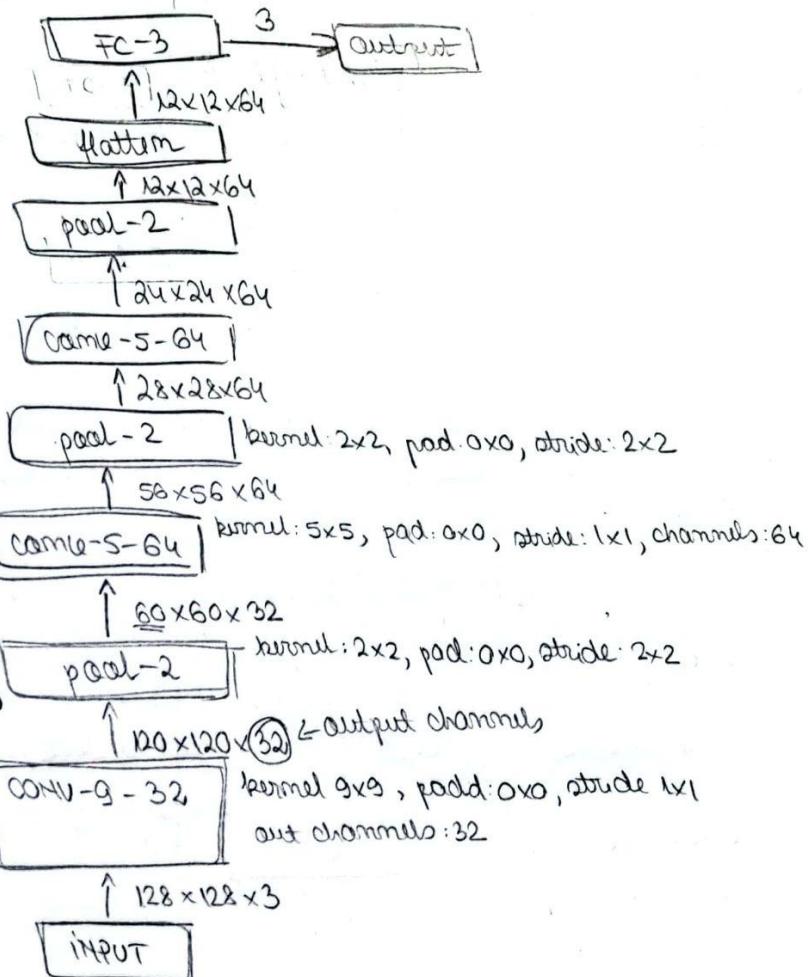
$$(32 \times 3) \times 9 \times 9 = 7776$$

kernel kernel

$$\uparrow \quad \uparrow$$

weights: ~~223~~ 7776

$$\text{bias: } 32$$

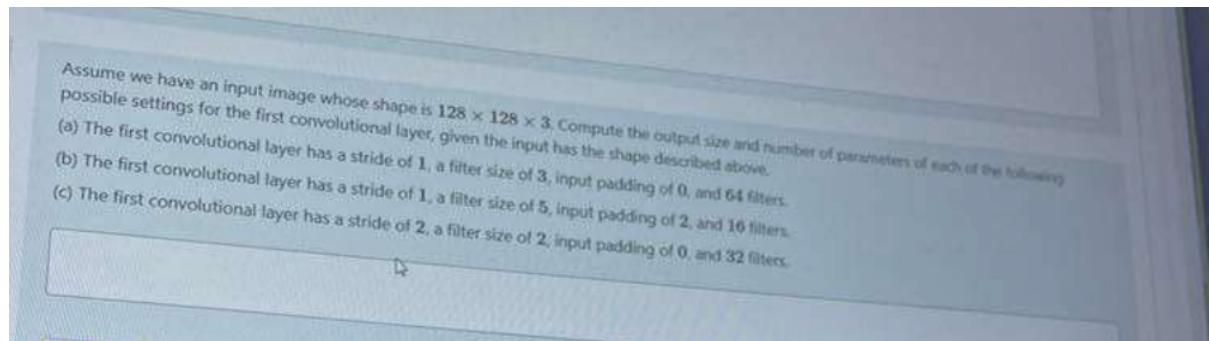


! Pt. pad m. limii output latime output

$$\left\{ \frac{(m_h - kh + ph + sh)}{sh} \right\} \times \left\{ \frac{(m_w - kw + pw + sw)}{sw} \right\}$$

$$128 - 9 + 0 + 1 / 1 = 120 \times (128 - 9 + 0 + 1 / 1 = 120)$$

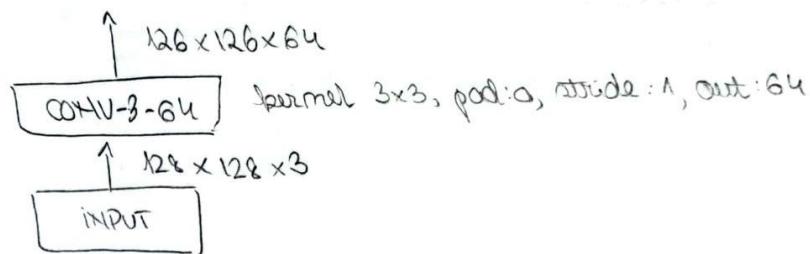
11. Assume we have an input **image** whose shape is $128 \times 128 \times 3$. Compute the output size and number of parameters of each of the following possible settings for the first convolutional layer, given the input has the shape described above:
- The first convolutional layer has a stride of 1, a filter size of 3, input padding 0 and 64 filters.
 - The first convolutional layer has a stride of 1, a filter size of 5, input padding 2 and 16 filters.
 - The first convolutional layer has a stride of 2, a filter size of 2, input padding 0 and 32 filters.



- M) Assume we have an input image whose shape is $128 \times 128 \times 3$. Compute the output size and number of parameters of each of the following possible settings for the first convolutional layer, given the input has the shape described above.
- The first convolutional layer has a stride of 1, a filter size of 3×3 , input padding of 0 and 64 filters.
 - The first convolutional layer has a stride of 1, a filter size of 5×5 , input padding of 2 and 16 filters.
 - The first convolutional layer has a stride of 2, a filter size of 2×2 , input padding of 0 and 32 filters.

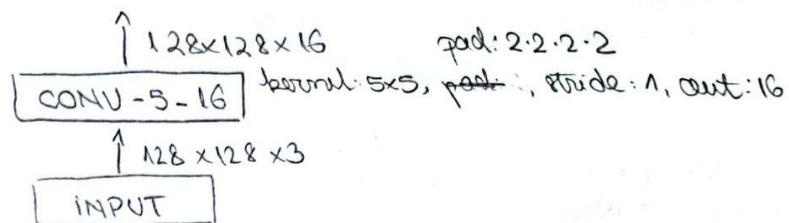
a)

$$\begin{aligned} & (128 - 3 + 0 + 1) / 1 \\ & = 126 \\ W &= (3 \times 3) \times 64 \times 3 \\ b &= 64 \end{aligned}$$



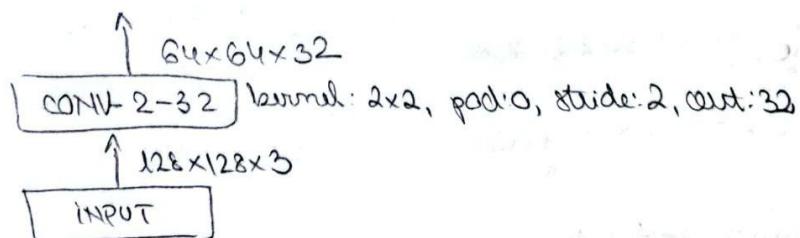
b)

$$\begin{aligned} & W = (5 \times 5) \times 16 \times 3 \\ & = 1200 \\ b &= 16 \end{aligned}$$



c)

$$\begin{aligned} & (128 - 2 + 0 + 2) / 2 \\ & = 64 \\ W &= 2 \times 2 \times 32 \times 3 = 384 \\ b &= 32 \end{aligned}$$



12. Compute symbolically the total number of parameters of a **GRU layer**, knowing that the number of examples is n , the number of input is d and the number of hidden units is h . Numerical application $n=2, d=3, h=4$.

Compute symbolically the total number of parameters of a GRU layer, knowing that the number of examples is n , the number of inputs is d , and the number of hidden units is h . Numerical application: $n = 2, d = 3, h = 4$.

12) Compute symbolically the total number of parameters of a GRU layer, knowing that the number of examples is m , the number of input is d and the number of hidden units is h .

Numerical application: $m=2, d=3, h=4$

$$z_t = \sigma(x_t w_{xz} + h_{t-1} w_{zh} + b_z)$$

$$r_t = \sigma(x_t w_{xr} + h_{t-1} w_{hr} + b_r)$$

$$w_{xz}, w_{xh} \in \mathbb{R}^{d \times h} \Rightarrow 3 \cdot 4 = 12 \\ + 12 = 24 \quad \left. \begin{array}{l} \\ \end{array} \right\} \Rightarrow 24 + 32 + 8 = 64$$

$$w_{hr}, w_{zh} \in \mathbb{R}^{h \times h} \Rightarrow 4 \cdot 4 = 16 \\ + 16 = 32$$

$$b_z, b_r \in \mathbb{R}^{1 \times h} \Rightarrow 4 + 4 = 8$$

$$\tilde{h}_t = \tanh(x_t w_{xh} + (R_t \odot h_{t-1}) w_{hh} + b_h)$$

$$w_{xh} \in \mathbb{R}^{d \times h} = 12 \quad \left. \begin{array}{l} \\ 16+12+4=32 \end{array} \right\}$$

$$w_{hh} \in \mathbb{R}^{h \times h} = 16$$

$$b_h \in \mathbb{R}^{1 \times h} = 4$$

$$64 + 32 = 96$$

$$2 \cdot dh + 2 \cdot hh + 2 \cdot h + dh + hh + h = 3dh + 3hh + 3h = 3(dh + hh + h) = \\ = 3h(d+h+1) = 3 \cdot 4(3+4+1) = \underline{\underline{96}}$$

$$\text{state-vec initial } h_0: 4(h) \Rightarrow 96 + 4 = 100 \quad \left. \begin{array}{l} \\ \rightarrow 115 \end{array} \right\}$$

$$\text{output layer: } hd + d = 4 \cdot 3 + 3 = 15$$

$$o_t = H_t + W_{hq} + b_q$$

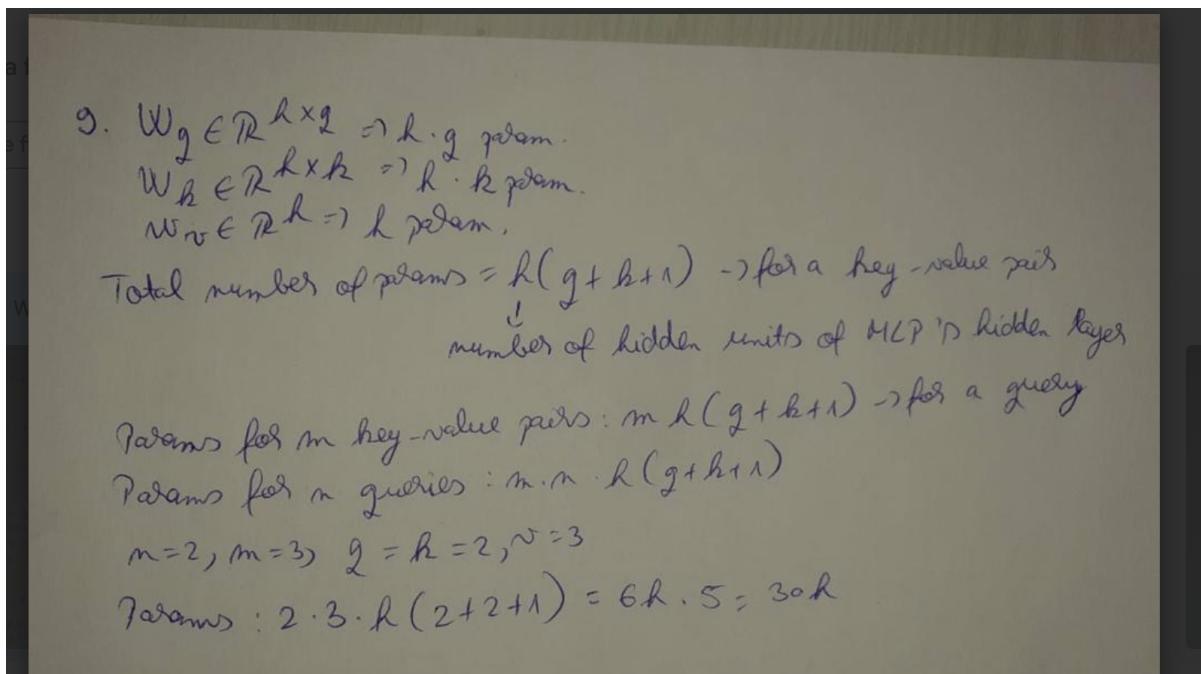
$$W_{hq} \in \mathbb{R}^{h \times q}, q = d$$

$$b_q \in \mathbb{R}^{1 \times q}$$

13. Compute the total number of parameters for **additive attention** for n queries and m key-value pairs, where queries are of length q , keys are of length k and values are of length v . Give the answer symbolically. Numerical application: $n = 2$, $m = 3$, $q = k = 2$, $v = 3$.

Compute the total number of parameters for additive attention for n queries and m key-value pairs, where queries are of length q , keys are of length k , and values are of length v . Give the answer symbolically. Numerical application:

$$n = 2, m = 3, q = k = 2, v = 3.$$



Common Andoni - Alexandru

9.

u queries \rightarrow length q
 m key-value pairs
 \downarrow
 $length$

v \downarrow
 $length$

$M = 2$

$m = 3$
 $q = k = 2$
 $v = 3$

$$a(q, k) = w_v^\top \tanh(w_q q + w_k k)$$

$$J = J(q, (k_1, v_1), \dots, (k_m, v_m)) = \sum_{i=1}^m \alpha(q, k_i) \cdot v_i, v_i \in \mathbb{R}^v$$

$$\alpha(q, k_i) = \frac{\exp(a(q, k_i))}{\sum_{j=1}^m \exp(a(q, k_j))}$$

$$\begin{aligned} w_q \in \mathbb{R}^{h \times q} &\rightarrow h \times q \text{ params} \\ w_k \in \mathbb{R}^{h \times k} &\rightarrow h \times k \text{ params} \\ w_v \in \mathbb{R}^h &\rightarrow h \text{ params} \end{aligned} \quad \left. \begin{array}{l} \Rightarrow \text{function } a = h_q + h_k + h \\ = h(q+k+1) \text{ params} \end{array} \right\} a$$

$$\text{Function } \alpha = \frac{\exp(a)}{\sum_{i=1}^m \exp(a_i)} \Rightarrow a + m \cdot a = (m+1)(q+k+1) \cdot h$$

$$J = m \cdot h \cdot (m+1)(q+k+1) \leftarrow \text{nr of parameters total per query}$$

$$\text{Cu valori: } \Rightarrow 2 \cdot 3 \cdot 2 \cdot 2 \cdot 2 \cdot 2 \Rightarrow 2 \cdot J$$

$$2 \cdot 3 \cdot h \cdot (3+1) \cdot (2+2+1)$$

$$6 \cdot h \cdot 4 \cdot 5 = 120h \text{ - nr of parameters}$$

h - nr. hidden layers

14. Assume a scalar multiplication and a scalar addition both account for one **FLOP** and exponentiation accounts for 10 FLOPS. All other scalar operations account for 0 FLOPS. How many FLOPS does the scaled dot-product attention require for n queries and m key-value pairs, where queries and keys are of length d and values are of length v . Give the answer symbolically. Numerical application: $n = 2, m = 3, d = 2, v = 3$.

Assume a scalar multiplication and a scalar addition both account for one FLOP and exponentiation accounts for 10 FLOPS. All other scalar operations account for 0 FLOPS. How many FLOPS does the scaled dot-product attention require for n queries and m key-value pairs, where queries and keys are of length d and values are of length v . Give the answer symbolically. Numerical application: $n = 2, m = 3, d = 2, v = 3$.

⑨ Scalar multiplication $\Rightarrow 1 \text{ FLOP}$
 Scalar addition $\Rightarrow 1 \text{ FLOP}$
 Exponentiation $\Rightarrow 10 \text{ FLOPs}$.
 All other scalar operations $\Rightarrow 0 \text{ FLOPs}$.

How many FLOPs for the scaled dot-product attention require for n queries and m key-value pairs?

{ queries, keys $\Rightarrow d$ Numerical: $n=2$ $v=3$
 values $\Rightarrow v$ $m=3$
 $d=2$

Solution

The scaled dot-product attention scoring function:

$$a(q, k) = q^T k / \sqrt{d}$$

↓ query ↓ key

$\Rightarrow Q \in \mathbb{R}^{n \times d}$ queries, $K \in \mathbb{R}^{m \times d}$ keys,
 $V \in \mathbb{R}^{m \times v}$ values

\Rightarrow We have to compute FLOPS for:

$$\text{softmax} \left(\frac{Q^T K}{\sqrt{d}} \right) V \in \mathbb{R}^{m \times v} (\in \mathbb{R}^{m \times v})$$

$$Q \in \mathbb{R}^{n \times d} \Rightarrow Q^T \in \mathbb{R}^{d \times n}$$

$$K \in \mathbb{R}^{m \times d}$$

Studiem pt un query "i" un key:

$$\alpha(Q, k) = Q^T k / \sqrt{d}$$

$$Q \in \mathbb{R}^{1 \times d}, k \in \mathbb{R}^{1 \times d}$$

$$\Rightarrow Q^T k \in \mathbb{R}^{d \times d} \Rightarrow Q^T k \text{ va genera } d^2 \text{ FLOPS.}$$

(fiecare element Q se va înmulți cu fiecare el k)

\Rightarrow Se vor genera înălătura d^2 FLOPS prin
multiplicarea fiecărui el. cu $\frac{1}{\sqrt{d}}$.

$$\Rightarrow 2d^2 \text{ generat pt. } \alpha(Q, k)$$

\Rightarrow Se va genera $m \cdot n \cdot 2d^2$ Flop's prin
aplicarea aceliei funcții.

$\Rightarrow m \cdot n \cdot 2d^2 (d-1)$ fiind cont de
înmulțirea matricelor

15. Assume we have a training set for **two-class classification in one dimension** that contains three sample points: point $x(1) = 3$ with label $y(1) = 1$, point $x(2) = 1$ with label $y(2) = 1$ and point $x(3) = -1$ with label $y(3) = -1$. What are the values of w and b given by a hard-margin SVM?

Assume we have a training set for two-class classification in one dimension that contains three sample points: point $x(1) = 3$ with label $y(1) = 1$, point $x(2) = 1$ with label $y(2) = 1$, and point $x(3) = -1$ with label $y(3) = -1$. What are the values of w and b given by a hard-margin SVM?

Maximum file size: 5MB, maximum number of attachments: 5

DRINCIANU ALEXANDRU-MIHAI
CTI-RO, AN 3

2. $x^{(1)} = 3, y^{(1)} = 1$
 $x^{(2)} = 1, y^{(2)} = 1$
 $x^{(3)} = -1, y^{(3)} = -1$

$\min_{w, w_0} \frac{1}{2} \|w\|_2^2, y^{(i)}(w^T x^{(i)} + w_0) \geq 1, i \in \{1, 2, 3\}$

$w^* = (x^T \cdot x)^{-1} \cdot x^T \cdot y$

$x = \begin{bmatrix} 1 & x^{(1)} \\ 1 & x^{(2)} \\ 1 & x^{(3)} \end{bmatrix} \Rightarrow w^* = \left[\begin{pmatrix} 1 & x^{(1)} \\ 1 & x^{(2)} \\ 1 & x^{(3)} \end{pmatrix}^{-1} \begin{pmatrix} 1 & 1 \\ 1 & 1 \\ 1 & 1 \end{pmatrix} \right] = \begin{pmatrix} y_1 & y_2 & y_3 \\ x_1^{(1)} & x_2^{(1)} & x_3^{(1)} \\ x_1^{(2)} & x_2^{(2)} & x_3^{(2)} \\ x_1^{(3)} & x_2^{(3)} & x_3^{(3)} \end{pmatrix}^{-1} \begin{pmatrix} y_1 & y_2 & y_3 \\ 1 & 1 & 1 \\ 1 & 1 & 1 \end{pmatrix}$

$w^* = \begin{pmatrix} 3 & x^{(1)}+x^{(2)}+x^{(3)} \\ x^{(1)}+x^{(2)}+x^{(3)} & x_1^{(1)} \cdot x_1^{(2)} \cdot x_1^{(3)} \end{pmatrix}^{-1} \cdot \begin{pmatrix} y_1 & y_2 & y_3 \\ x_1^{(1)} \cdot g_1 & x_2^{(1)} \cdot g_2 & x_3^{(1)} \cdot g_3 \end{pmatrix} =$

$w^* = \begin{pmatrix} 3 & 3 \\ 3 & 11 \end{pmatrix}^{-1} \cdot \cancel{\begin{pmatrix} 1 \\ 5 \end{pmatrix}} \cdot \begin{pmatrix} 1 \\ 5 \end{pmatrix} =$

$= \begin{pmatrix} \frac{11}{2} & -\frac{3}{2} \\ -\frac{3}{2} & \frac{3}{2} \end{pmatrix} \cdot \begin{pmatrix} 1 \\ 5 \end{pmatrix} = \begin{pmatrix} \frac{21}{2} - \frac{15}{2} \\ -\frac{15}{2} + \frac{15}{2} \end{pmatrix} = \begin{pmatrix} -2 \\ 0 \end{pmatrix} \Rightarrow w = -2, w_0 = 0$

$y^{(1)}(w^T \cdot x^{(1)} + w_0) \geq 1$
 $y^{(2)}(w^T \cdot x^{(2)} + w_0) \geq 1$
 $y^{(3)}(w^T \cdot x^{(3)} + w_0) \geq 1$

16. Compute symbolically the total number of parameters (weights and biases) of an **LSTM** layer, knowing that the number of examples is n , the number of inputs is d and the number of hidden units is h . Numerical application: $n = 2, d = 3, h = 4$.

You can drag and drop files here to add them.

Time left 2:59:29

Compute symbolically the total number of parameters (weights and biases) of an LSTM layer, knowing that the number of examples is n , the number of inputs is d , and the number of hidden units is h . Numerical application: $n = 2, d = 3, h = 4$.

(8)

Total number of parameters of an LSTM layer.

$\begin{cases} n - \text{examples} \\ d - \text{inputs} \\ h - \text{hidden units} \end{cases}$

Numerical: $\begin{cases} n=2 \\ d=3 \\ h=4 \end{cases}$

Solution

LSTM layer

$$I_t = \sigma(x_t W_{xi} + H_{t-1} W_{hi} + b_i)$$

$$F_t = \sigma(x_t W_{xf} + H_{t-1} W_{hf} + b_f)$$

$$O_t = \sigma(x_t W_{xo} + H_{t-1} W_{ho} + b_o)$$

$$\tilde{C}_t = \tanh(x_t W_{xc} + H_{t-1} W_{hc} + b_c)$$

$W_{xi}, W_{xf}, W_{xo}, W_{xc} \in \mathbb{R}^{d \times h} \Rightarrow d \times h \text{ parametri'}$

$W_{hi}, W_{hf}, W_{ho}, W_{hc} \in \mathbb{R}^{h \times h} \Rightarrow h \times h \text{ parametri'}$

$b_i, b_f, b_o, b_c \in \mathbb{R}^{1 \times h} \Rightarrow 1 \times h \text{ parametri'}$

$$\Rightarrow \text{Numărul total de parametri: } 4(dh + h^2 + h)$$
$$\Rightarrow \underline{4h(d+2+1)}$$

$$\text{Numeric: } 4 \cdot 4(3+4+1) = 16 \cdot 8 = 128 \text{ parametri}$$

$$b \in \mathbb{R}^{12}$$
$$\Rightarrow 3 \cdot (3 \cdot 4) + 3 \cdot (4 \cdot 4) + 3 \cdot 4 + 4 \cdot 4 = 112 \text{ weight parameters}$$
$$3 \cdot (1 \cdot 4) + 1 \cdot 4 = 16 \text{ bias parameters}$$