

CS 260: Foundations of Data Science

Prof. Sara Mathieson

Fall 2023



HVERFORD
COLLEGE

Admin

- In lab today:
 - Project meetings with all groups
 - Try to come to the same lab as your partner

Admin

- In lab today:
 - Project meetings with all groups
 - Try to come to the same lab as your partner

- Flexibility for this class:
 - I will drop the lowest lab grade
 - Final project can be simplified as needed
 - Accommodations or flexibility beyond this:
collaboration with class deans

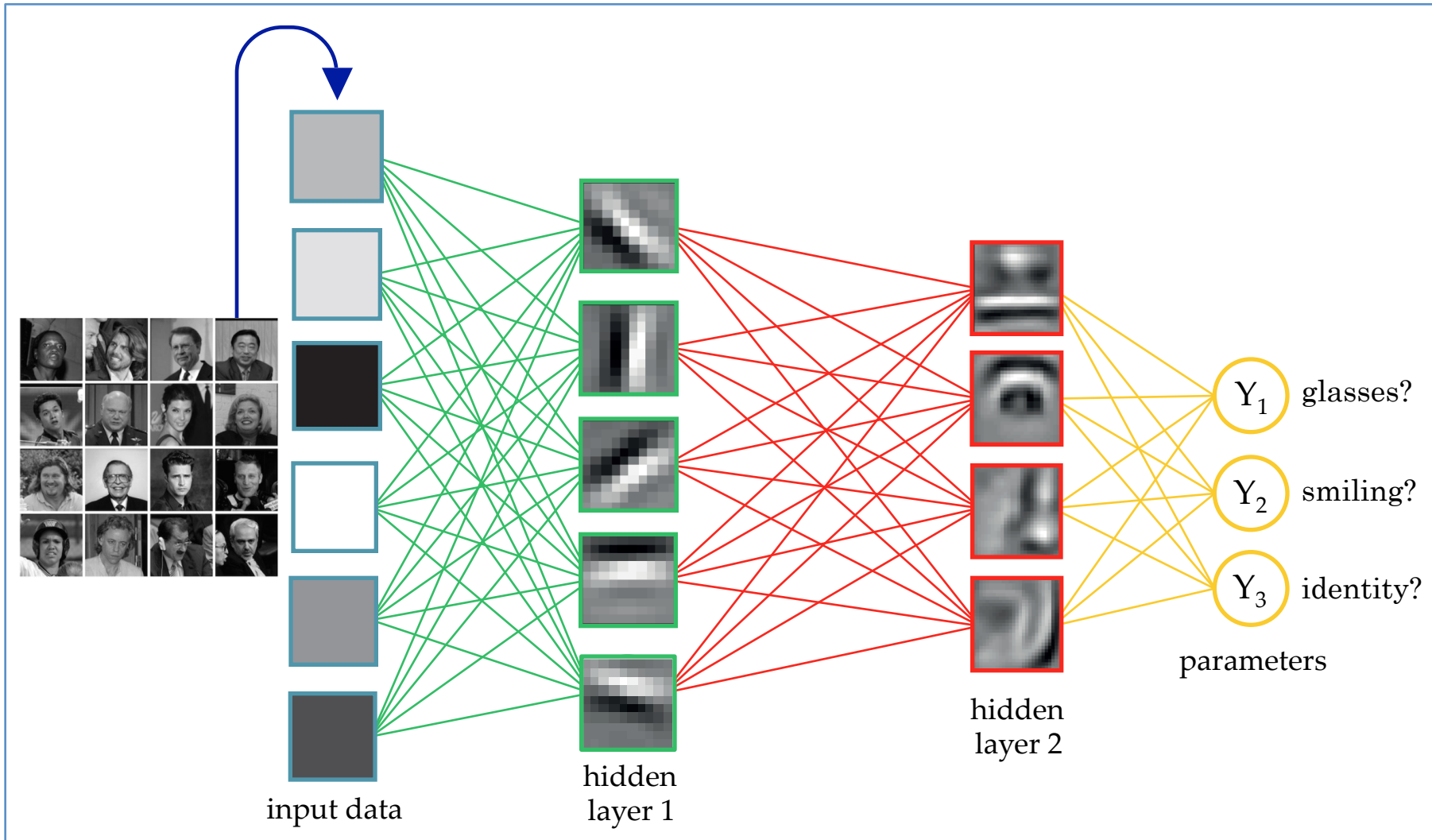
Outline for December 5

- Finish: neural networks
- Advice about git for final project
- Go over Midterm 2

Outline for December 5

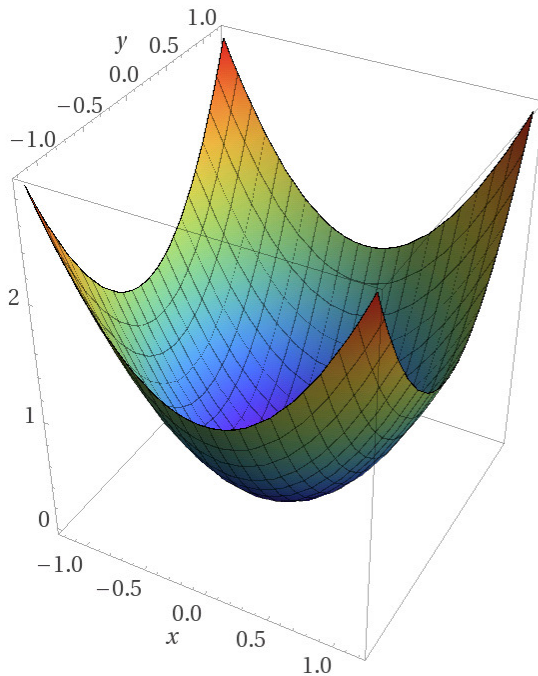
- Finish: neural networks
- Advice about git for final project
- Go over Midterm 2

First fully connected neural networks for images



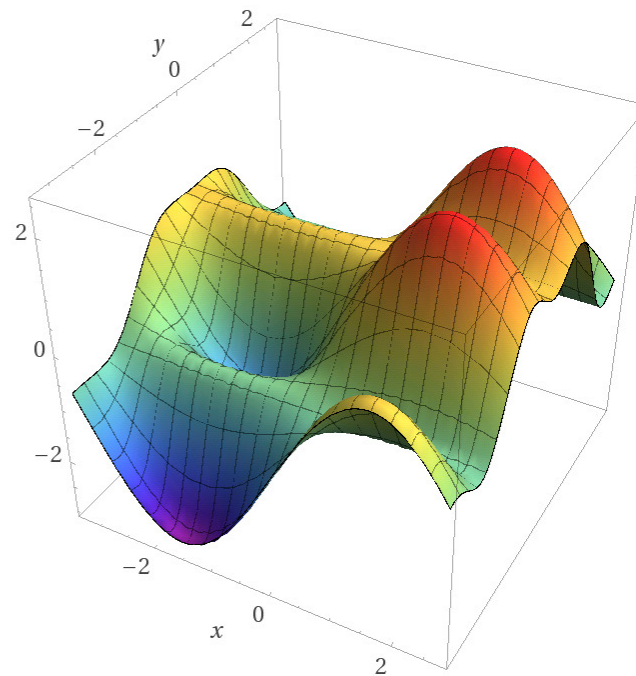
Takeaways from last time

- As the number of parameters grows, a non-convex function often has more and more local minima
- Starting at a “good” point is crucial!



Computed by Wolfram|Alpha

Convex



Computed by Wolfram|Alpha

Non-convex

Takeaways from last time

- Unsupervised pre-training uses latent structure in the data as a starting point for weight initialization
- After this process, the network is “fine-tuned”
- In practice this has been found to increase accuracy on specific tasks (which could be specified after feature learning)

Weight initialization

- We still have to initialize the pre-training
- All 0's initialization is bad! Causes nodes to compute the same outputs, so then the weights go through the same updates during gradient descent
- Need asymmetry! => usually use small random values

Mini-batches

- So far in this class, we have considered *stochastic gradient descent*, where one data point is used to compute the gradient and update the weights
- On the flipside is *batch gradient descent*, where we compute the gradient with respect to all the data, and then update the weights
- A middle ground uses *mini-batches* of examples before updating the weights

Notes about scores and softmax

- The output of the final fully connected layer is a vector of length K (number of classes)

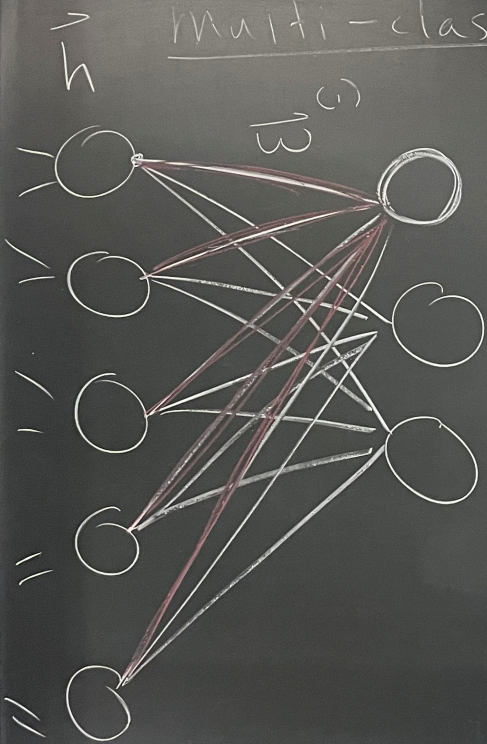
Notes about scores and softmax

- The output of the final fully connected layer is a vector of length K (number of classes)
- The raw scores are transformed into probabilities using the *softmax function*: (let s_k be the score for class k)

$$\hat{y}_k = \frac{e^{s_k}}{\sum_{j=1}^K e^{s_j}}$$

- Then we apply *cross-entropy loss* to these probabilities

Multi-class classification



$y=1$ (dog)

$y=2$ (cat)

$y=3$ (bird)

(score)

$$s_1 = \vec{w}^{(1)} \cdot \vec{h}$$

$$s_2 = \vec{w}^{(2)} \cdot \vec{h}$$

$$s_3 = \vec{w}^{(3)} \cdot \vec{h}$$

$y = \begin{bmatrix} 0 \\ 1 \\ 0 \end{bmatrix}$

linear functions

$$\begin{bmatrix} s_1 \\ s_2 \\ s_3 \end{bmatrix} = \begin{bmatrix} -7.2 \\ 10.9 \\ 0.71 \end{bmatrix}$$

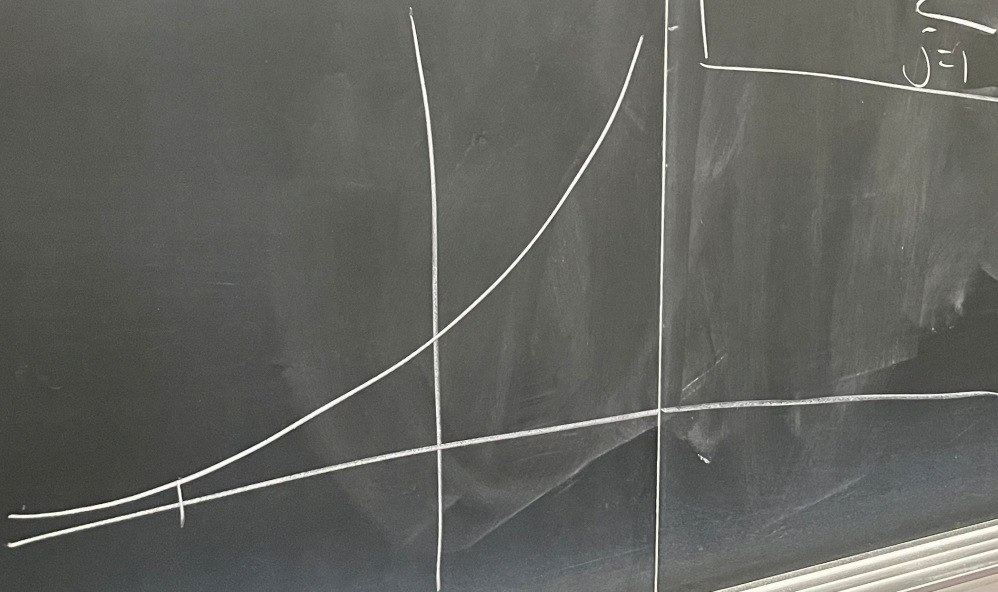
argmax

$$\hat{y} = 2 \text{ (cat)}$$

Softmax

prob?

Softmax



"Soft" (probability)

$$P_i = \frac{-7.2}{-7.2 + 10.9 + 0.71}$$

Negative prob!

$$\hat{y}_k = \frac{e^{s_k}}{\sum_{j=1}^k e^{s_j}}$$

$$\frac{e^{-7.2}}{e^{-7.2} + e^{10.9} + e^{0.71}}$$

$$\frac{1}{1 + e^{-x}}$$

Motivation for moving away from FC architectures

- For a $32 \times 32 \times 3$ image (very small!) we have $p=3072$ features in the input layer
- For a $200 \times 200 \times 3$ image, we would have $p=120,000$! *doesn't scale*

Motivation for moving away from FC architectures

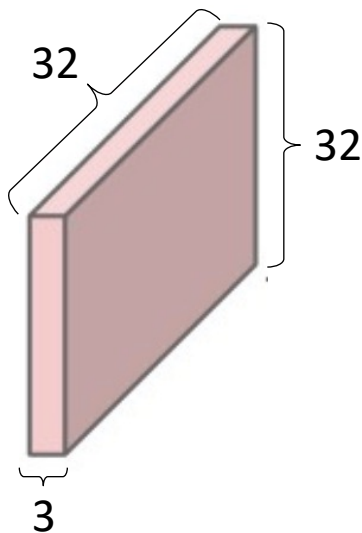
- For a 32x32x3 image (very small!) we have $p=3072$ features in the input layer
- For a 200x200x3 image, we would have $p=120,000!$ *doesn't scale*
- FC networks do not explicitly account for the structure of an image and the correlations/relationships between nearby pixels

Idea: 3D volumes of neurons

- Do not “flatten” image, keep it as a volume with *width*, *height*, and *depth*

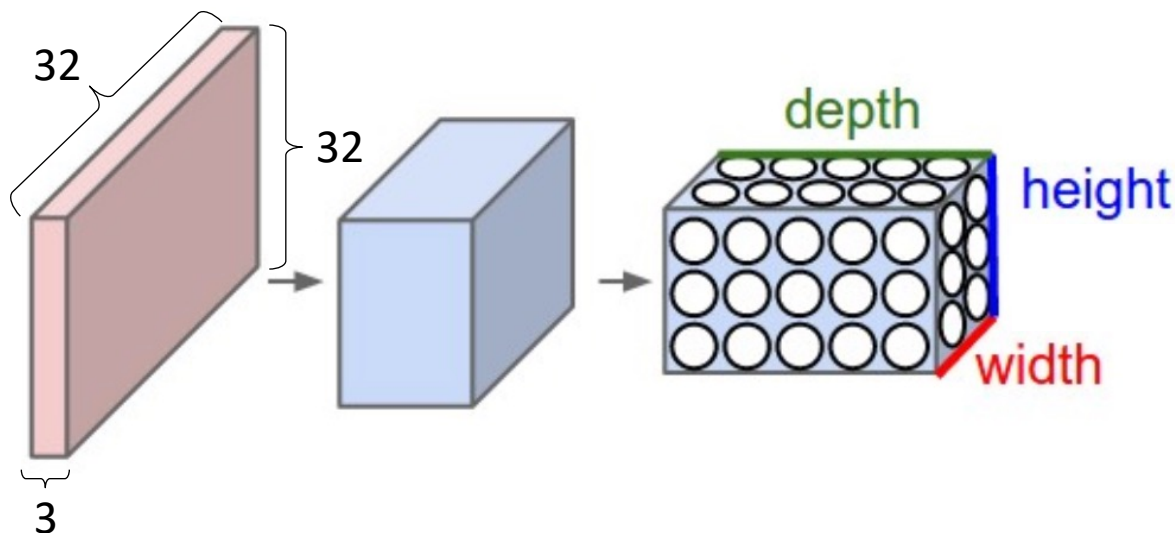
Idea: 3D volumes of neurons

- Do not “flatten” image, keep it as a volume with *width*, *height*, and *depth*
- For **CIFAR-10**, we would have:
 - Width=32, Height=32, Depth=3



Idea: 3D volumes of neurons

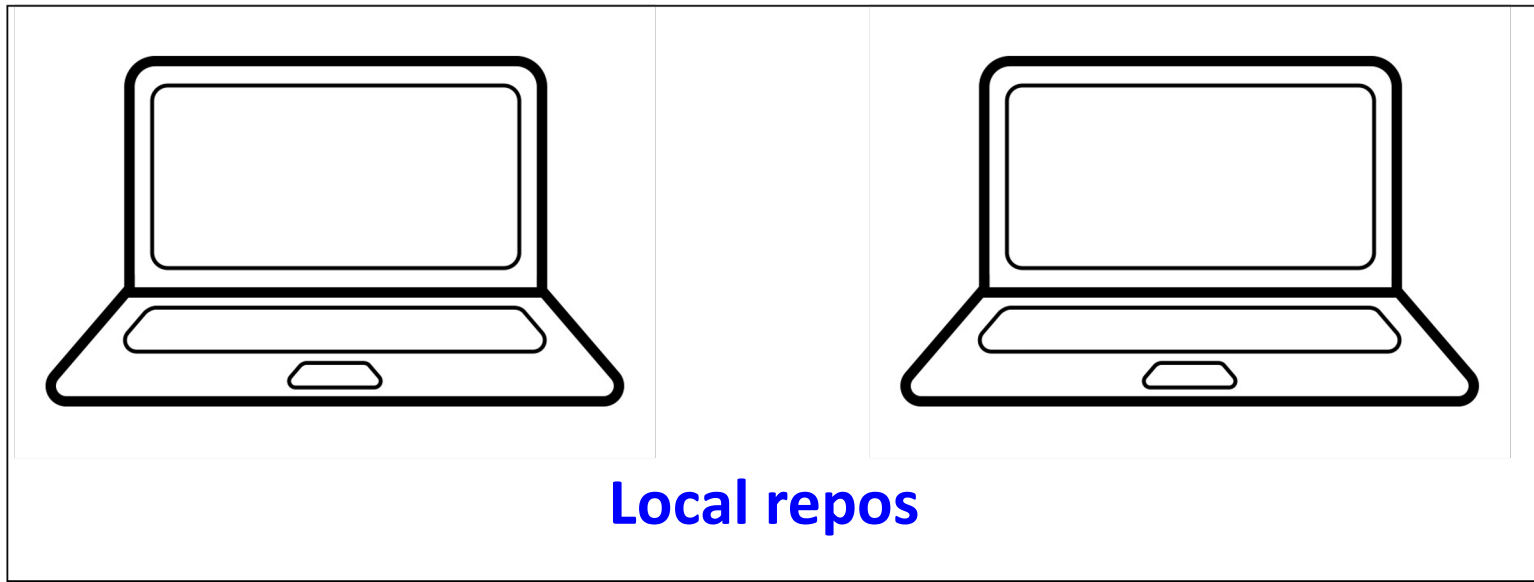
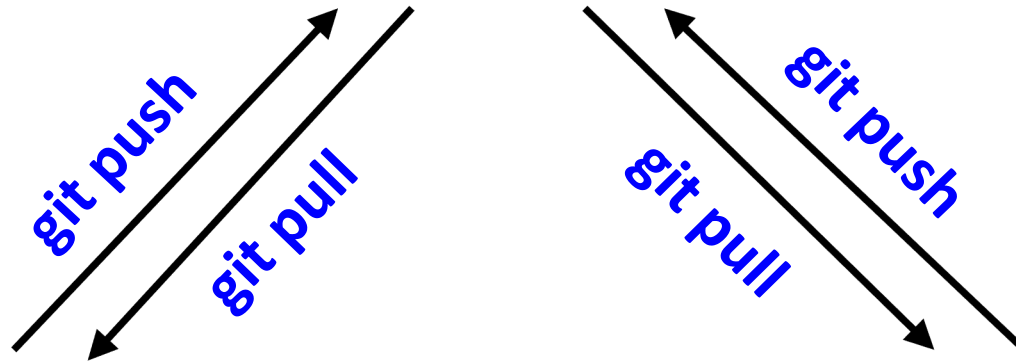
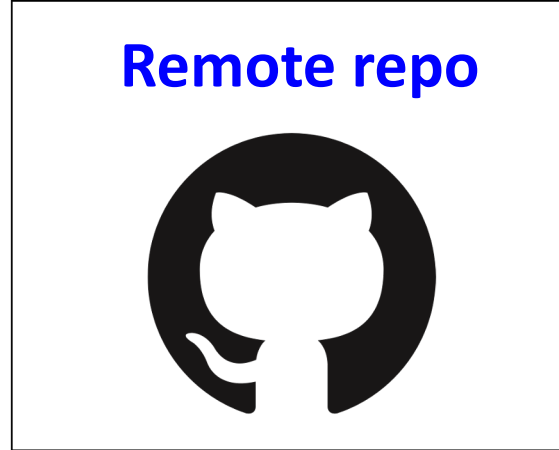
- Do not “flatten” image, keep it as a volume with *width*, *height*, and *depth*
- For **CIFAR-10**, we would have:
 - Width=32, Height=32, Depth=3
- Each layer is also a 3 dimensional volume



Outline for December 5

- Finish: neural networks
- **Advice about git for final project**
- Go over Midterm 2

Github workflow with your partner



Outline for December 5

- Finish: neural networks
- Advice about git for final project
- **Go over Midterm 2**

1e) $A = \begin{bmatrix} \boxed{1} \\ \vdots \\ \vdots \\ \boxed{1} \end{bmatrix}_{p \times p}$

$$\Rightarrow \sum_{i=1}^n (f_{i1} - \bar{f}_1)(f_{i2} - \bar{f}_2)$$

$\Rightarrow \boxed{O(np^2)}$

3b) cumulative prob $\boxed{0.625}$

$$0. \boxed{1} \cdot Z^{-1} + \boxed{0} \cdot Z^{-2} + \boxed{1} \cdot Z^{-3}$$

$$\begin{array}{r} 0.625 \\ - 0.5 \\ \hline 0.125 \end{array} \quad \begin{array}{r} \rightarrow \cancel{0.125} \\ - \cancel{0.25} \\ \hline \end{array} \quad \begin{array}{r} \rightarrow 0.125 \\ - 0.125 \\ \hline 0 \end{array}$$

$\Rightarrow \boxed{0.101000\dots}$

$$\textcircled{3e} \quad 0.375 \cdot 2 + (0.25 \cdot 2) \cdot 2 + 0.125 \cdot 3 = \boxed{2.125}$$

$\textcircled{4}$ P = pos H = healthy
 N = neg D = disease

$$x = P(P|D) = \frac{P(P)P(D|P)}{P(D)} = \frac{\textcircled{P(P)} \cdot 0.8}{\frac{1}{500}}$$

$$\frac{1}{500} x = \left(\frac{499}{500} (1-x) + \frac{1}{500} x \right) \frac{4}{5}$$

$$\begin{aligned}
 P(P) &= P(P, H) + P(P, D) \\
 &= P(H)P(P|H) + P(D)P(P|D) \\
 &= \frac{499}{500} (1-x) + \frac{1}{500} x
 \end{aligned}$$

$$\begin{aligned}
 &P(P|H) + \overset{P(N|H)}{P(P|D)} = 1 \\
 &\boxed{P(N|D) + P(P|D) = 1}
 \end{aligned}$$

$$x = 0.9995 \dots$$

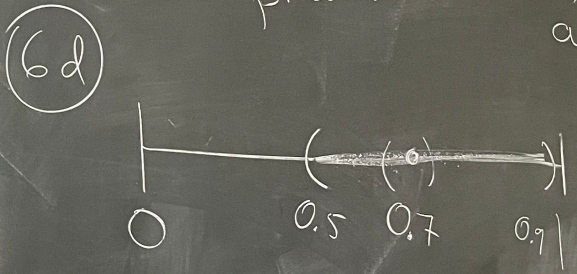
Part 5 f

$$\begin{aligned}
 P(y=k | \vec{x}) &\propto p(y=k) p(x_1, x_2 | y=k) \\
 &= p(x_1 | y=k) p(x_2 | x_1, y=k)
 \end{aligned}$$

$$\begin{aligned}
 &\rightarrow \frac{N_{x_1=v_1, x_2=v_2, y=k+1}}{N_{x_1=v_1, y=k} + |f_2|}
 \end{aligned}$$

$$\begin{aligned}
 &\rightarrow \frac{N_{x_1=v_1, x_2=v_2, y=k+1}}{N_{y=k+1} |f_1| |f_2|}
 \end{aligned}$$

(6c) classifier that always predicts female \Rightarrow acc=75%



(6e-9)

$$\bar{z} = \frac{1}{n} \sum_{i=1}^n x_i - y_i$$

(9) $\vec{x} = [15, 7, 10, 18, 11]$
 $\vec{y} = [12, 4, 7, 15, 8]$

[12	7	10	18
[15	4	7	15
}	}	}	}

8
[]
=
3

$\rightarrow \approx 0$