

Naïve Bayes Example

- Reference: <https://web.stanford.edu/~jurafsky/slp3/>, Chapter 3.

Positive tweets
I am happy because I am learning NLP
I am happy, not sad.
Negative tweets
I am sad, I am not learning NLP
I am sad, not happy

word	Pos	Neg
I	3	3
am	3	3
happy	2	1
because	1	0
learning	1	1
NLP	1	1
sad	1	2
not	1	2
N_{class}	13	12

word	Pos	Neg
I	0.24	0.25
am	0.24	0.25
happy	0.15	0.08
because	0.08	0
learning	0.08	0.08
NLP	0.08	0.08
sad	0.08	0.17
not	0.08	0.17

$P(x|pos)$ \rightarrow $P(x|neg)$

Let's classify the following tweets as **positive or negative**:

1. I am not sad.

2. I am learning NLP.

$$P(Pos) = \frac{13}{25}$$

$$P(Neg) = \frac{12}{25}$$

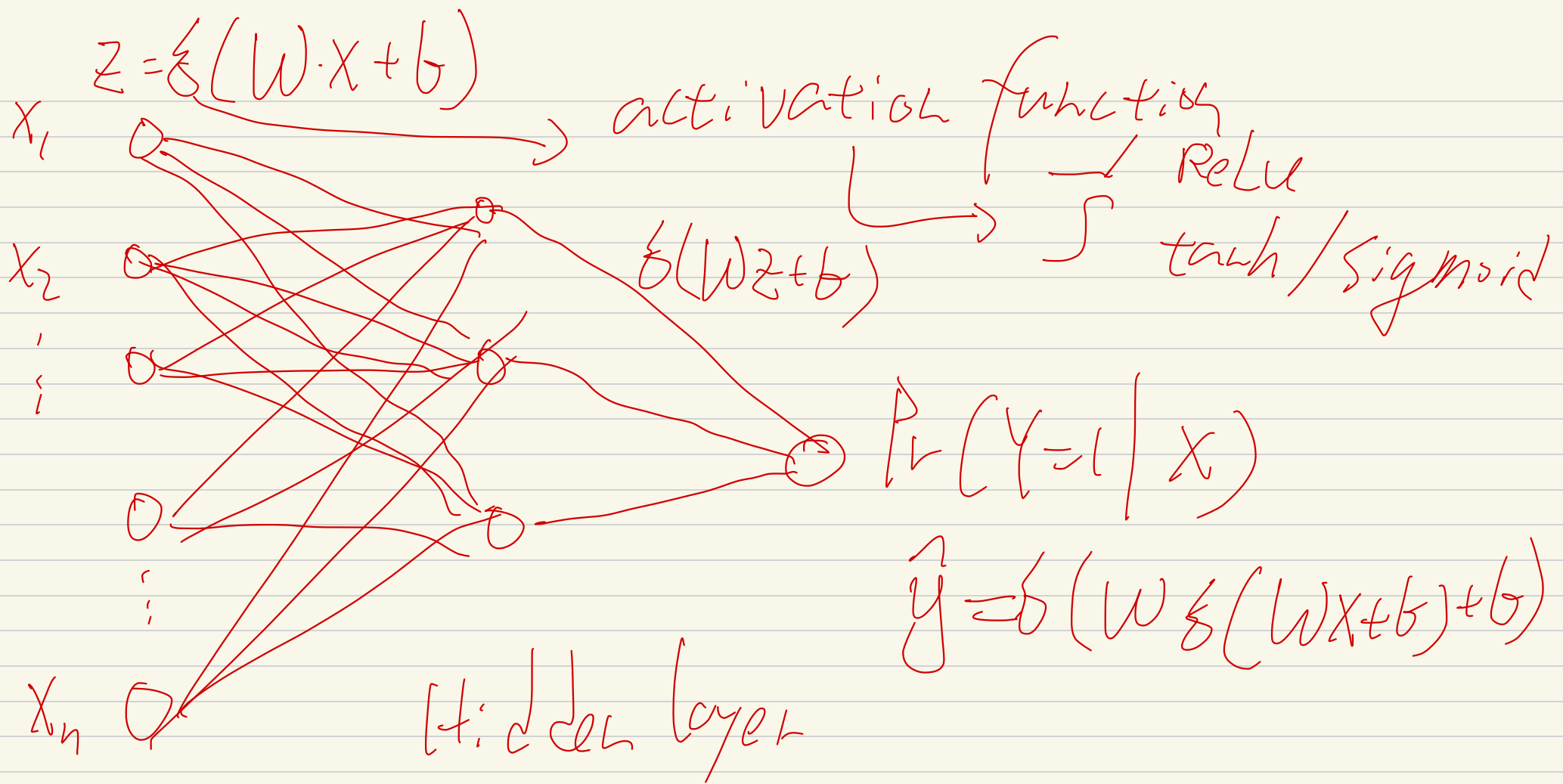
$$P(Pos) \cdot P(I|Pos) \cdot P(am|Pos) \cdot P(sad|Pos)$$

$$P(Neg) \cdot P(I|Neg) \cdot P(am|Neg)$$

$$- P(not|Neg) \cdot P(sad|Neg)$$

$$\frac{13}{25} \cdot \frac{0.15}{0.5} \cdot \frac{0.24}{0.5} = \frac{0.07}{0.17}$$

< 1



① How to obtain $\theta = (W, b)$?

$$\min_{\theta} \frac{1}{n} \sum_{i=1}^n \mathcal{L}(y_i, f(x_i, \theta))$$

↳ DNN.

② What is L ? $\left\{ \begin{array}{l} \text{① regression: } (y - \hat{y})^2 \\ \text{② Classification} \\ \rightarrow y \log \hat{y} - (1 - y) \log (1 - \hat{y}) \end{array} \right.$

③ How to find $\hat{\theta} = \underset{\theta}{\operatorname{argmin}} L(\theta)$

Gradient Descent.

$$\hat{\theta}_{t+1} = \hat{\theta}_t - \alpha \cdot \nabla_{\theta} L(\hat{\theta}_t)$$

④ How to estimate $\nabla_{\theta} \mathcal{L}(\hat{\theta}_t)$

(i) Chain-Rule / Bp

$$y = f(z) \quad z = g(x)$$

$$\frac{dy}{dx} = \frac{df}{dz} \bigg|_{z=g(x)} \cdot \frac{dz}{dx} \bigg|_x$$

(ii) SGD: Shuffle data observations

$$\hat{\theta}_{t+1} = \hat{\theta}_t - \alpha \cdot \nabla_{\theta} \mathcal{L}(Y_i, f(X_i, \hat{\theta}_t))$$

Adams

Default Choice of optimizer

(1) Initialization $(w, b) \sim N(0, \frac{1}{2n_l})$

(2) Normalization: $\frac{x_i - \bar{x}_i}{SD_i}$

(3) Skip Connection

g_1, g_2, \dots, g_M

Vanishing
Gradients

$$f(x) = \underbrace{-x + f(x)} + x$$