

More Algorithms

Naive Bayes

X, y (supervised) dataset (I.I.D.)

Goal: Approximate $P_{\text{data}}(y | \vec{x}) = \frac{P_{\text{data}}(y, \vec{x})}{P_{\text{data}}(\vec{x})}$

Chain rule for conditional probabilities

$$\begin{aligned}
 P_{\text{data}}(y, \vec{x}) &= p(y, x_1, \dots, x_n) \\
 &= p(x_1 | x_2, \dots, x_n, y) \cdot p(x_2 | x_3, \dots, x_n, y) \\
 &\quad \vdots \\
 &\quad \circ p(x_{n-1} | x_n, y) \\
 &\quad \circ p(x_n | y) \\
 &\quad \circ p(y)
 \end{aligned}$$

Naive assumption: all features are independent

given y .

$$\rightarrow p(y, \vec{x}) = p(y) \prod_{j=1}^n p(x_j | y)$$

$$\therefore p(y|\vec{x}) = \frac{p(y) \prod_{j=1}^n p(x_j|y)}{p(\vec{x})}$$

Naive Bayes Prediction

$$\hat{y} = \operatorname{argmax}_y p(y|\vec{x})$$

$$\hat{y} = \operatorname{argmax}_y p(y) \prod_{j=1}^n p(x_j|y)$$

Maximum
a posteriori
(MAP)
decision rule

Find maximum likelihood estimates
for $p(y)$ and $p(x_j|y)$ $1 \leq j \leq n$

$$p_{\text{model}}(y; \vec{\theta}) \quad p_{\text{model}}(x_j|y; \vec{\theta}^{(j)})$$

Goal: Find $\vec{\theta}_{\text{ML}}$, $\vec{\theta}_{\text{ML}}^{(j)}$

(models depend on the type of data)

E.g.: (Classification)

$$\vec{X} \in \mathbb{R}^{m \times n}$$

X, \vec{y}
continuous discrete

$$y_i \in \{1, \dots, s\} \quad \forall i$$

$$p_{\text{model}}(y=i; \vec{\theta}) = \theta_i$$

$$\sum_{i=1}^{185} \theta_i = 1$$

$$p_{\text{model}}(x_j | y=i; \mu_{i,j}, \sigma_{i,j}^2) = \mathcal{N}(x_j; \mu_{i,j}, \sigma_{i,j}^2)$$

Goal Calculate $\vec{\Theta}_{ML}$ and $(\mu_{i,j}, \sigma_{i,j}^2)_{ML}$

$$\underline{\text{Claim}}: (\vec{\Theta}_{ML})_i = \frac{\#\{k | y_k = i\}}{m}$$

$i = \text{label}$

$j = \text{feature}$

$$(\mu_{i,j})_{ML} = \frac{\sum_{k: y_k = i} X_{k,j}}{\#\{k | y_k = i\}} = \text{Mean}(\{X_{k,j} | y_k = i\})$$

$$(\sigma_{i,j}^2)_{ML} = \text{Var}(\{X_{k,j} | y_k = i\})$$

↑
use $\frac{1}{n-1}$ instead of $\frac{1}{n}$

"Bessel corrected variance"

Then we can predict for a new \vec{x} :

$$\hat{y} = \arg \max_i \Theta_i \prod_{j=1}^n \mathcal{N}(x_j; (\mu_{i,j})_{ML}, (\sigma_{i,j}^2)_{ML})$$

Honorable Mentions

- Support Vector Machines (kernel methods)
- k-nearest neighbors
- Decision Trees / Random Forest

Unsupervised Algorithms (no target \vec{y})

Have dataset X or $\{\vec{x}^{(1)}, \dots, \vec{x}^{(n)}\}$

Some Common Tasks:

- ① Estimate $p_{\text{data}}(\vec{x})$
- ② Draw new samples from $p_{\text{data}}(\vec{x})$
- ★ ③ Denoise X
- ④ Find a manifold which lies near data
- ★ ⑤ Find a "simpler representation"
 - lower dimensional
 - sparse
 - independence of features (remove dependencies)

PCA revisited: $X \in \mathbb{R}^{m \times n}$

WLOG each column has mean 0
 (if not, subtract the mean from each column)
 (so $E[\vec{x}] = 0$)

Consider the unbiased covariance matrix:

$$\text{Var}[\vec{x}] = \frac{1}{m-1} X^T X$$

(i,j) entry measures how much feature i is linearly correlated with feature j. (i,i) entry measures the variance in feature i.

Consider the SVD $X = U \Sigma V^T$ and define

$T: \mathbb{R}^n \rightarrow \mathbb{R}^n$ $T(\vec{x}) = V^T \vec{x}$. Consider the new representation of the data:

$$\vec{z} = T(\vec{x})$$

$$\text{Then } \text{Var}[\vec{z}] = \frac{1}{m-1} \begin{bmatrix} -T(\vec{x}_1) - \\ \vdots \\ -T(\vec{x}_m) - \end{bmatrix}^T \begin{bmatrix} -T(\vec{x}_1) - \\ \vdots \\ -T(\vec{x}_m) - \end{bmatrix} = \frac{1}{m-1} \vec{z}^T \vec{z} \quad \text{where}$$

$$= \frac{1}{m-1} \begin{bmatrix} 1 & 1 & \dots & 1 \\ V^T \vec{x}_1 & V^T \vec{x}_2 & \dots & V^T \vec{x}_m \\ \vdots & \vdots & \ddots & \vdots \end{bmatrix} \begin{bmatrix} " \\ " \\ " \\ " \end{bmatrix} \quad \vec{z} = X \vec{V}$$

$$= \frac{1}{m-1} V^T X^T X V \quad \left(V \text{ is an orthogonal eigenvector basis for } X^T X \text{ with eigenvalues } \sigma_1^2, \dots, \sigma_n^2 \right)$$

~~rank 2~~

$$= \frac{1}{m-1} \begin{bmatrix} \sigma_1^2 & 0 & \dots & 0 \\ 0 & \sigma_2^2 & \dots & 0 \\ \vdots & \vdots & \ddots & \vdots \\ 0 & 0 & \dots & \sigma_n^2 \end{bmatrix}$$

\therefore no features of \vec{z} are linearly correlated.

$$\text{Also, } \sigma_1^2 \geq \sigma_2^2 \geq \dots \geq \sigma_n^2 \text{ so}$$

$$\text{Var}(z_1) \geq \text{Var}(z_2) \geq \dots \geq \text{Var}(z_n)$$

so features of \vec{z} are ordered in terms of importance

We also know $\forall k \leq n$

$$T_k: \mathbb{R}^n \rightarrow \mathbb{R}^k, \quad T_k(\vec{x}) = \begin{bmatrix} 1 & 1 & \dots & 1 \\ V_1 & V_2 & \dots & V_k \\ \vdots & \vdots & \ddots & \vdots \\ 1 & 1 & \dots & 1 \end{bmatrix}^T \vec{x} \quad \text{V}_k^T$$

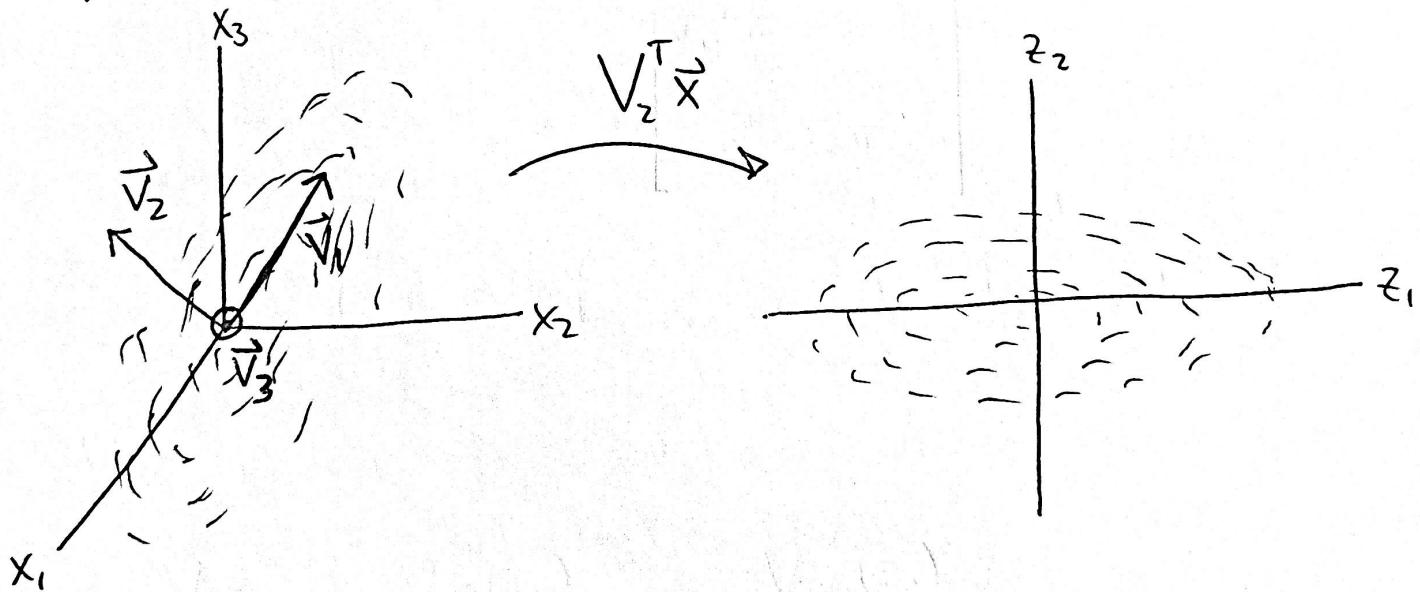
optimally preserves the data (as you proved in HW 1)

\therefore representing the data via

$$Z = X V_k \quad \text{accomplishes:}$$

- ① Maximally preserves data
- ② Uncorrelates features
- ③ Orders features by importance

E.g.



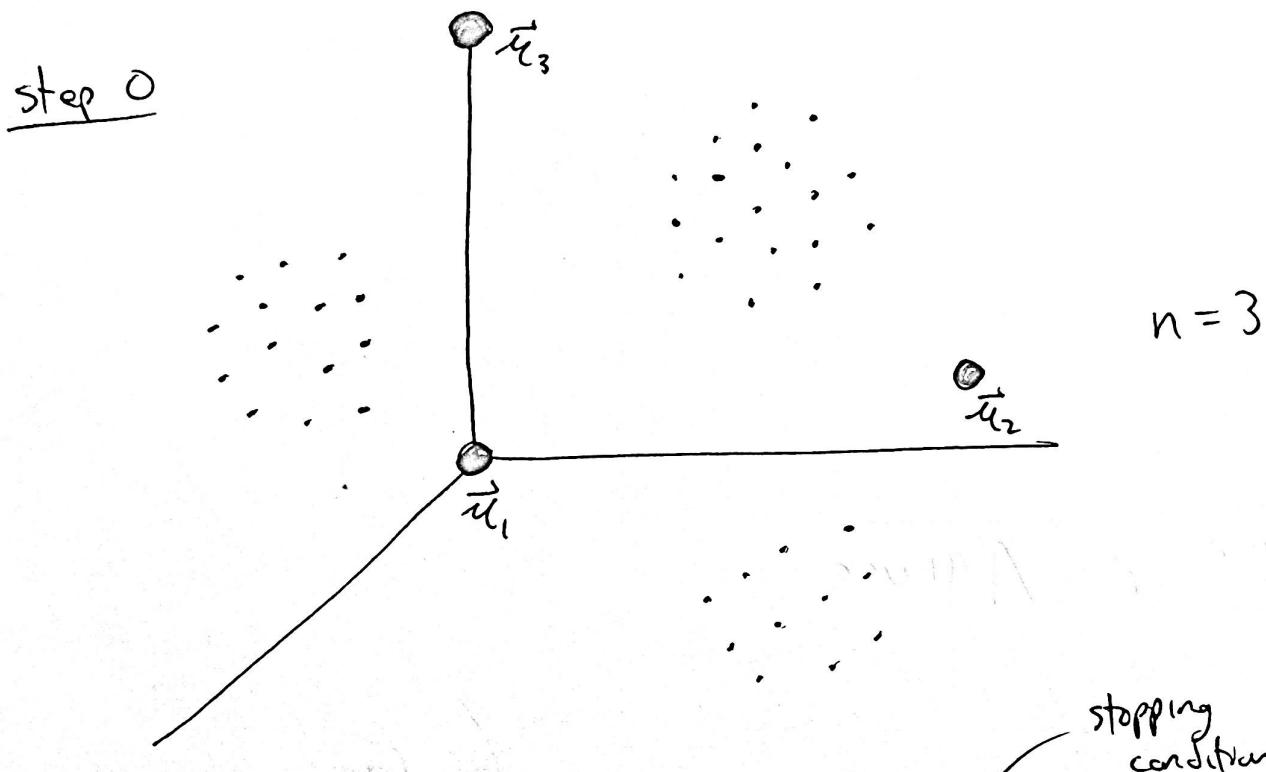
"almost on a plane"

$$X \in \mathbb{R}^{m \times 3}$$

$$V \in \mathbb{R}^{3 \times 3}$$

K-means clustering (Unsupervised)

Dataset $X \in \mathbb{R}^{m \times n}$ datapoint: $X_{i,*} = \vec{x}^{(i)}$



$n = 3$

Pick k ($= \# \text{ clusters}$). Pick $\epsilon > 0$.

Randomly initialize k points $\vec{u}_1, \dots, \vec{u}_k \in \mathbb{R}^n$

$\pi \in \mathbb{R}^m$ $\pi[i] = j$ if $\vec{x}^{(i)}$ is closest to \vec{u}_j

$$J = \frac{1}{m} \sum_{i=1}^m \| \vec{x}^{(i)} - \vec{u}_{\pi[i]} \|_2 ; \text{ costs} = [J] ; i = 0$$

While True:

if $i > 1$ and $| \text{costs}[i] - \text{costs}[i-1] | < \epsilon$:

break

stopping condition

for i in range(k):

$$\vec{\pi}_i = \frac{1}{m_i} \sum_{j: \pi[j]=i} \vec{x}^{(j)}$$

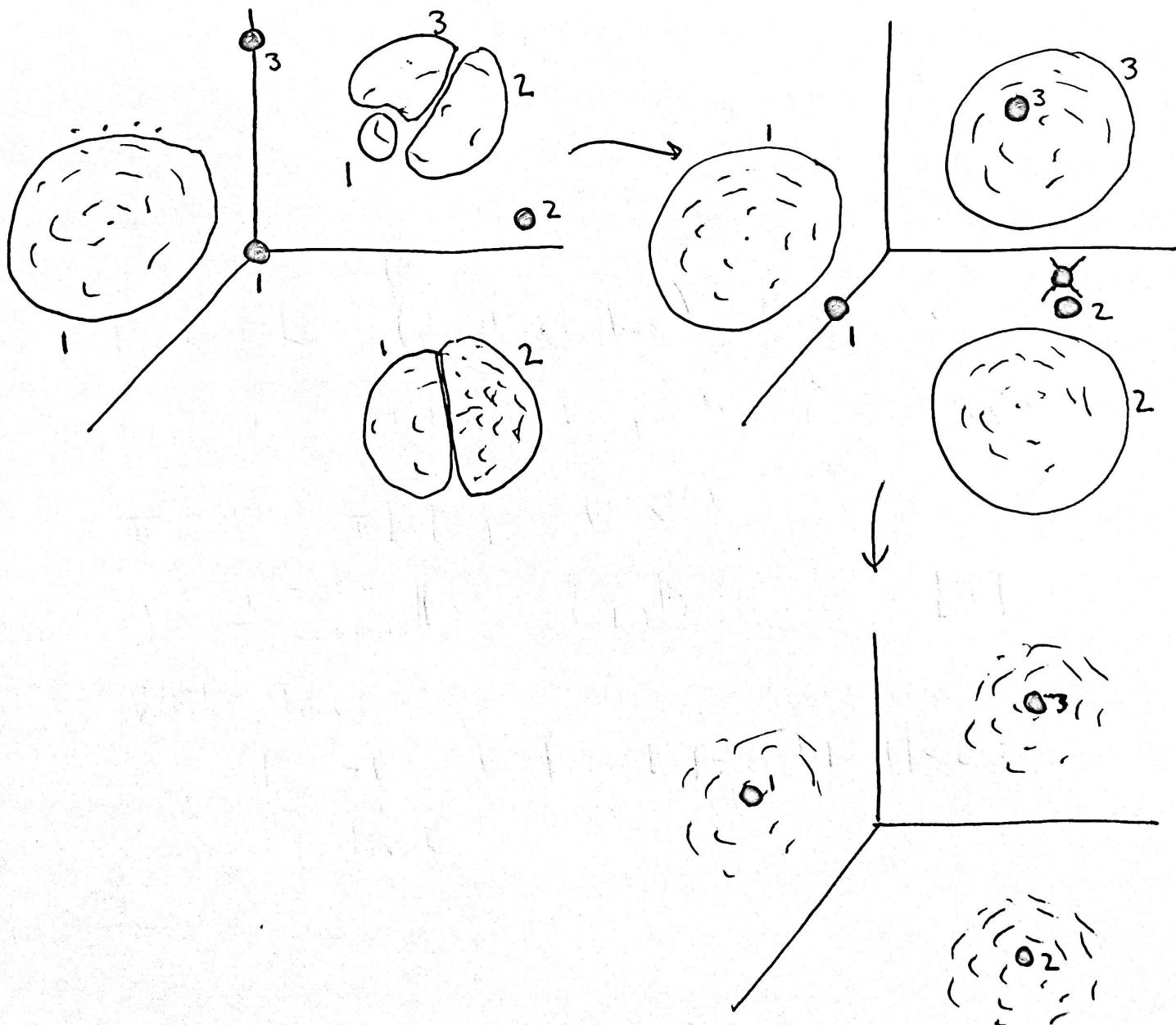
$$m_i = \#\{j \mid \pi[j] = i\}$$

recalculate π, J

i += 1

costs.append(J)

E.g. Continuing from step 0



E.g. Classifying handwritten digits

Given dataset $X \in \mathbb{R}^{m \times 64}$ $X_{i,\ell} =$ straightforward matrix of grayscale values of 8×8 image of digit (with values in)

0-16

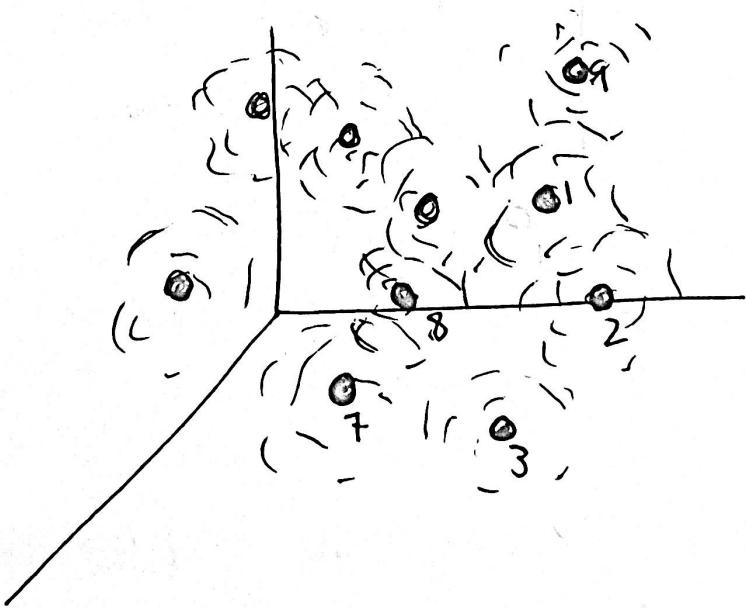
white black

E.g.

3×3 case

$$\begin{array}{|c|c|c|}\hline & \text{■} & \text{■} \\ \hline \text{■} & \text{■} & \text{■} \\ \hline \text{■} & \text{■} & \text{■} \\ \hline\end{array} \rightarrow \begin{pmatrix} 0 \\ 0 \\ 0 \\ 15 \\ 14 \\ 7 \\ 0 \\ 0 \\ -1 \end{pmatrix}$$

Run k-means



Label by hand. Make predictions on new digits by checking which is the closest cluster.

New Optimization Algorithm

Stochastic Gradient Descent (SGD)

Usually we have some cost function

$$J(\theta) = \frac{1}{m} \sum_{i=1}^m L(x^{(i)}, y^{(i)}; \theta)$$

per example error
"loss"

$$\text{Then } \nabla_{\theta} J(\theta) = \frac{1}{m} \sum_{i=1}^m \nabla_{\theta} L(x^{(i)}, y^{(i)}; \theta)$$

so $\nabla_{\theta} J(\theta)$ is an average

Idea: Estimate $\nabla_{\theta} J(\theta)$ by picking a subset of $\{1, \dots, m\}$ of size $m' \ll m$ (maybe $m = 1 \text{ billion}$)
 $m' = 100$

This is called a mini-batch:

$$B = \{x^{(i)}, \dots, x^{(im')}\} \subseteq \{x^{(1)}, \dots, x^{(m)}\}$$

Compute $\nabla_{\theta} J(\theta) \approx \underbrace{\frac{1}{m'} \sum_{j=1}^{m'} \nabla_{\theta} L(x^{(ij)}, y^{(ij)}; \theta)}$

Then run similarly to the usual GD

$$\theta = \theta - \alpha g \quad (\alpha = \text{step size})$$

Suggestion Casually read the rest of Chp 5

Now Chp 6 Deep Feedforward Networks

Supplement from
"Deep Learning
Architectures"

or
Feedforward Neural Networks

or
Multilayer Perceptrons (MLPs)

Goal) Approximate a function $y = f^*(\vec{x})$ by finding parameters θ for a model $y = f_{\text{model}}(\vec{x}; \theta)$ by running some optimization algorithm to find suitable θ

Abstract Neurons

Definition An ~~actual~~ abstract neuron (or just neuron) is a quadruple $(\vec{x}, \vec{w}, \varphi, y)$

$$\vec{x} = (x_0, \dots, x_n)$$

"input vector"

$$x_0 = -1$$

$$\vec{w}^T = (w_0, \dots, w_n)$$

"weight vector"

$$b = w_0 = \text{bias}$$

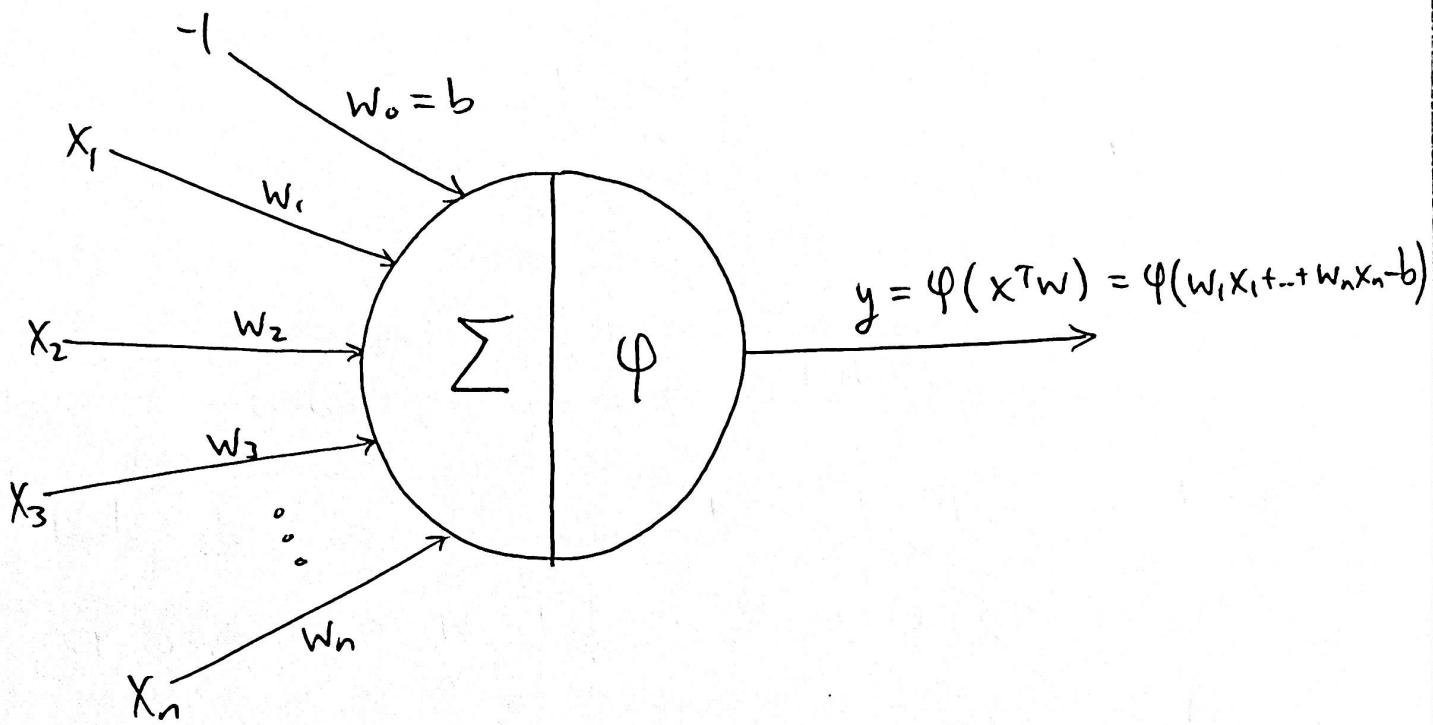
$$\varphi : \mathbb{R} \rightarrow \mathbb{R}$$

"activation function"

$$y = \varphi(x^T w)$$

"output"

Idea: A neuron tries to "learn" a target z by tuning \vec{w} so that $\varphi(x^T w) \approx z$



Def: A perceptron is a neuron with inputs $\vec{x} \in \{0,1\}^n$ and with activation function

$$\varphi(x) = \begin{cases} 0 & x < 0 \\ 1 & x \geq 0 \end{cases} \quad (\text{Heaviside function})$$

E.g. $f^*: \{0,1\} \times \{0,1\} \rightarrow \{0,1\}$

$0 = \text{False}$

$1 = \text{True}$

$$f^*(x,y) = x \text{ and } y$$

Use a perceptron with $w = \begin{pmatrix} w_0 \\ w_1 \\ w_2 \end{pmatrix}$

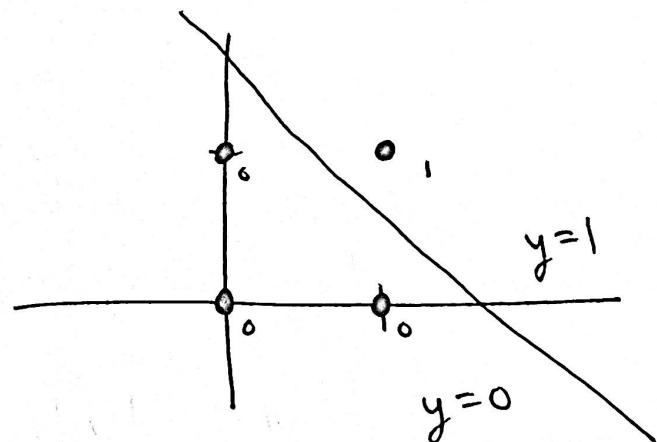
We want $\varphi(w^T x) = f^+(x)$ $\forall x \in \{0,1\}^2$

$$\Rightarrow \left. \begin{array}{l} \varphi(w^T \begin{pmatrix} -1 \\ 0 \\ 0 \end{pmatrix}) = 0 \\ \varphi(w^T \begin{pmatrix} -1 \\ 1 \\ 0 \end{pmatrix}) = 0 \\ \varphi(w^T \begin{pmatrix} -1 \\ 0 \\ 1 \end{pmatrix}) = 0 \\ \varphi(w^T \begin{pmatrix} -1 \\ 1 \\ 1 \end{pmatrix}) = 1 \end{array} \right\} \Leftrightarrow \begin{array}{l} w^T \begin{pmatrix} -1 \\ 0 \\ 0 \end{pmatrix} < 0 \\ w^T \begin{pmatrix} -1 \\ 1 \\ 0 \end{pmatrix} < 0 \\ w^T \begin{pmatrix} -1 \\ 0 \\ 1 \end{pmatrix} < 0 \\ w^T \begin{pmatrix} -1 \\ 1 \\ 1 \end{pmatrix} \geq 0 \end{array}$$

A solution: $w = \begin{pmatrix} 1.5 \\ 1 \\ 1 \end{pmatrix}$

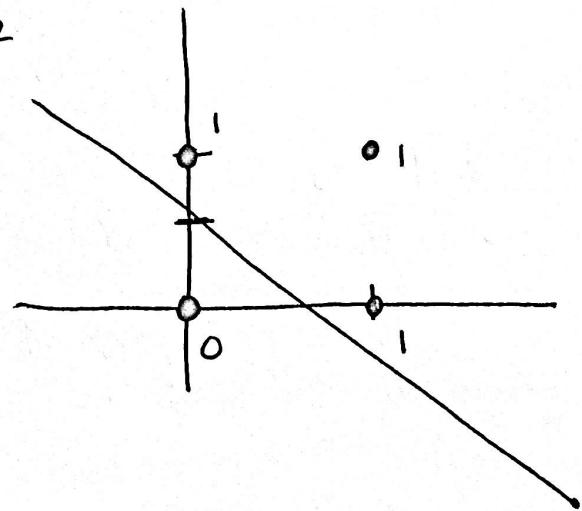
$$w^T x = -1.5 + x_1 + x_2 = 0$$

$$x_2 = 1.5 - x_1$$

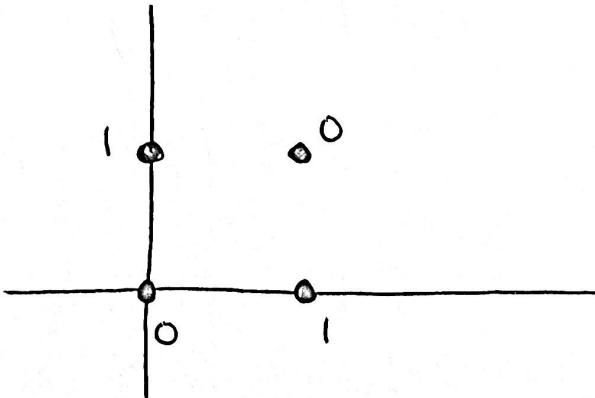


E.g. $f^+(x_1, x_2) = x_1 \text{ or } x_2$

$$w = \begin{pmatrix} 0.5 \\ 1 \\ 1 \end{pmatrix}$$



E.g. $f^*(x_1, x_2) = x_1 \text{ XOR } x_2$ (one or the other
but not both)



Perceptron cannot fit this function

(no line in \mathbb{R}^2 separates)
1's from 0's