

Reminders

- ↳ Notes (investing time is worth it)
- ↳ Office Hours: M + W 12:30 - 2,
and by appointment
- ↳ Discord
 - DM
 - Pset Questions

Context

Regression : $\underline{x}^{(i)} \in \mathbb{R}^d$ $y^{(i)} \in \mathbb{R}$

Classification: $\underline{x}^{(i)} \in \mathbb{R}^d$ $y^{(i)} \in \{0, \dots, K\}$

↳ spam or not spam (binary)

↳ objects in an image

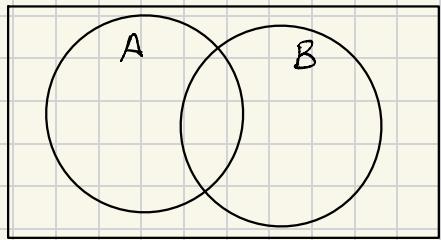
↳ next word

Today: Probability background

Probability

A, B random events

it rains tomorrow
I carry an umbrella



$$\Pr(A)$$

$$\Pr(A \cup B)$$

$$\Pr(A \cap B)$$

$$\Pr(A | B)$$

Properties:

$$0 \leq \Pr(A) \leq 1$$

$$\begin{aligned}\Pr(A \cap B) &= \Pr(A) \Pr(B|A) \\ &= \Pr(B) \Pr(A|B)\end{aligned}$$

$$\Pr(A \cup B) = \Pr(A) + \Pr(B) - \Pr(A \cap B)$$

$$\neg A \leftarrow \text{"complement"} \\ \neg A = \underline{\text{not}} A$$

$$\Pr(\neg A) + \Pr(A) = 1 \iff \Pr(\neg A) = 1 - \Pr(A)$$

$$A, B \text{ indep iff } \Pr(A \cap B) = \Pr(A) \Pr(B)$$

Random Variables:

$X \sim D$ ↗ a distribution
Sampled from

For example, $X \sim$ Dice roll

$$\Pr(X=1) = \Pr(X=2) = \dots = 1/6$$

Bayes Rule:

$$\Pr(A|B) = \frac{\Pr(B|A) \Pr(A)}{\Pr(B)}$$

Proof?

Maximum A Posteriori

Our regression strategy: Choose model that most likely generated data

Classification analogy: Choose label that is most likely given data

Example: $X = \text{email}$, $y = \mathbb{I}[\text{email is spam}]$

Posterior:

$$\Pr(Y=1 | X=x) \quad \text{vs.} \quad \Pr(Y=0 | X=x)$$

Compute using Bayes Rule!

$$\Pr(Y=1 | X=x) = \frac{\Pr(Y=1) \Pr(X=x | Y=1)}{\Pr(X=x)} = \frac{\text{prior}}{\text{evidence}} \frac{\text{likelihood}}{\text{evidence}}$$

Medical Example:

$X \in \{0, 1\}$ outcome of test

$Y \in \{0, 1\}$ disease

- Rare (1% population)
- 5% FPR
- 10% FNR

Name Bayes Classifier (with binary features)

assume independent features

↙ d

$$\text{Naive assumption : } \Pr(X=x | Y=1) = \prod_{i=1}^d \underbrace{\Pr(X_i=x_i | Y=1)}_{p_i^{(1)} \text{ or } (1-p_i^{(1)})}$$

$$\text{For example, } \Pr(X=\begin{bmatrix} 0 \\ 1 \\ 0 \\ 0 \\ 1 \end{bmatrix} | Y=1) = (1-p_1) p_2 (1-p_3) (1-p_4) p_5$$

1. Compute priors $\Pr(Y=1), \Pr(Y=0)$

2. Compute observed probs $p_i^{(1)} = \Pr(X_i=1 | Y=1), p_i^{(0)} = \Pr(X_i=1 | Y=0)$

3. Compute likelihoods

$$\Pr(X=x | Y=1) = \prod_{i=1}^d \Pr(X_i=x_i | Y=1)$$

4. Predict class
with higher posterior

$$\operatorname{argmax}_{i \in \{0, 1\}}$$

$$\frac{\Pr(X=x | Y=i) \Pr(Y=i)}{\Pr(X=x)}$$

9/18/2025

Reminders

↳ Quiz

- going forward, single answer
- learning rate strategies
 - ↳ manual
 - ↳ scheduler
 - ↳ adaptive
- NOT analytic (just problem)
- SGD vs GD
 - ↳ speed (RAM vs memory)

↳ Notes

↳ Student Sessions

- discuss material or psets
- advertise on discord
- incentive: organizer can write 1 quiz question
- selfie = proof

↳ Psets 1 and 4 due Monday

↳ Office hours, let's talk!

Review

$$\underline{x}^{(i)} \in \mathbb{R}^d$$

$y^{(i)} \in \{0, 1, \dots, k-1\}$

$k=2$ for now

↓

classification

Today: Build classification model with parameters

Naive Bayes:

Choose $\underset{y \in \{0, 1\}}{\operatorname{argmax}} \Pr(Y=y | X=x)$

Naive assumption is independence

$$\Pr(X=x | Y=y) = \prod_{i=1}^d \Pr(X_i=x_i | Y=y)$$

↑ supervised (labeled data)

but no parameters

Model & Loss

Idea: Adapt linear regression

Attempt #1:

$$f(x) = \langle w, x \rangle$$

$$\mathcal{L}(w) = |y - \langle w, x \rangle|$$

Attempt #2:

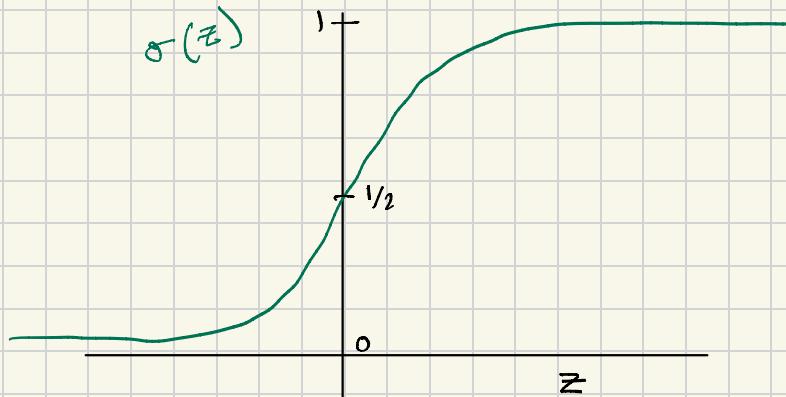
$$f(x) = \langle w, x \rangle$$

$$\mathcal{L}(w) = (y - \langle w, x \rangle)^2$$

Attempt #3:

$$f(x) = \sigma(w, x) \quad \text{where}$$

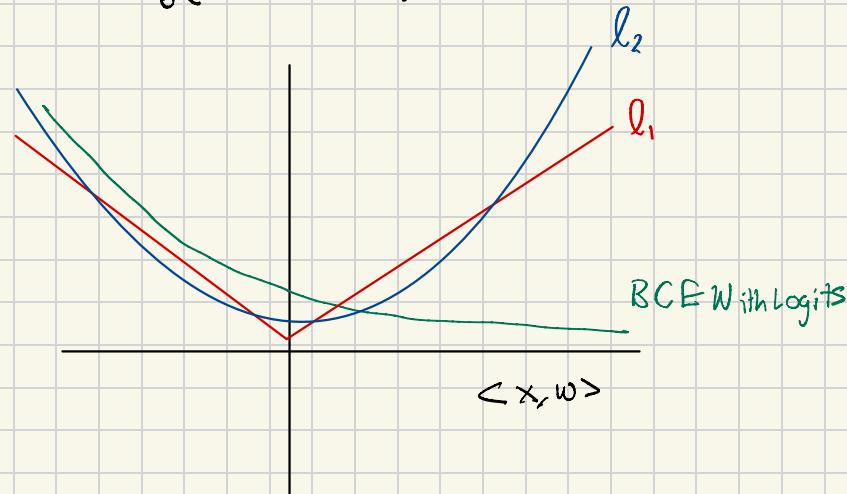
$$\sigma(z) = \frac{1}{1+e^{-z}}$$



Binary Cross Entropy loss

$$\mathcal{L}(w) = - \left[y \log (\sigma(\langle w, x \rangle)) + (1-y) \log (1 - \sigma(\langle w, x \rangle)) \right]$$

$$y=1 \\ = - \log (\sigma(\langle w, x \rangle))$$



Optimization

Goal: Compute $\nabla \mathcal{L}(\omega)$

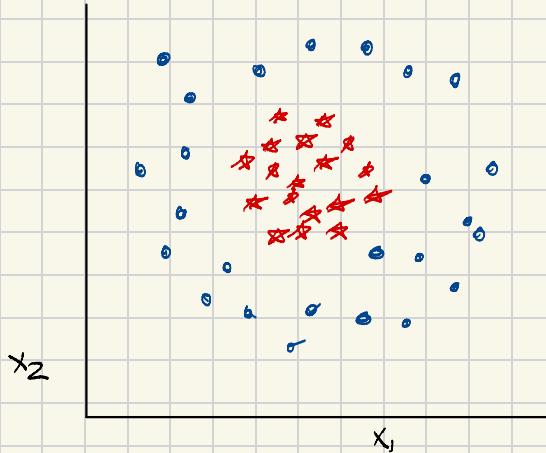
Theorem:

$$\nabla \mathcal{L}(\omega) = X^T (\sigma(X\omega) - y)$$

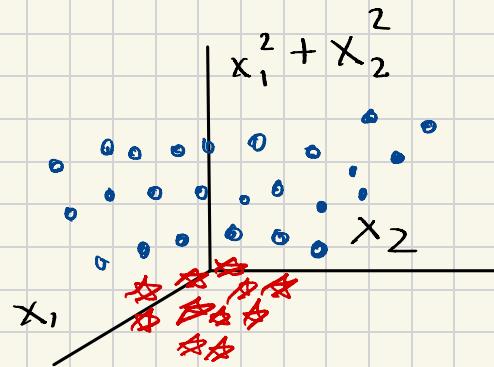
↙ similar to
linear regression?

↙ can we solve for
 ω^* ?

Non-linear Transformation



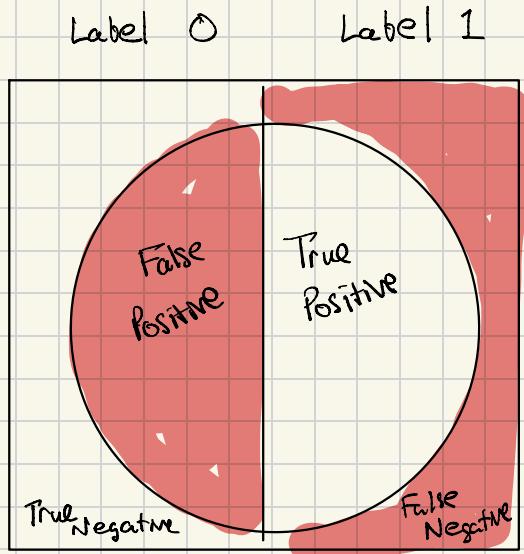
not linearly separable



linearly separable

Measuring Error in Binary Classification

$$\text{error rate} = \frac{\# \text{ incorrect}}{\# \text{ predictions}}$$



$$\text{TPR} = \frac{P}{D}$$

$$\text{FPR} = \frac{Q}{D}$$

$$\text{Precision} = \frac{D}{Q}$$

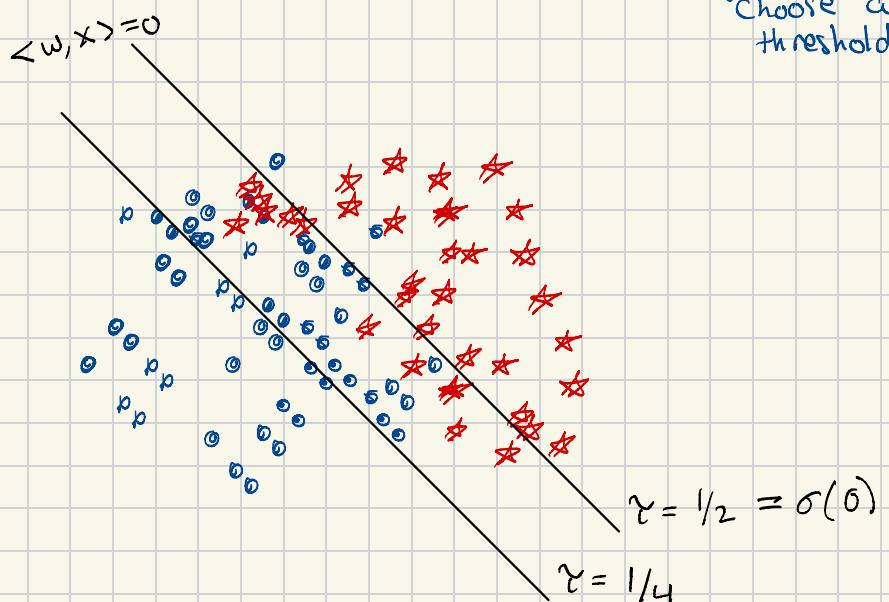
A Hidden Lever...

If incorrect predictions, change threshold

Let $p^{(i)} = f(x^{(i)}) = \sigma(\langle w, x^{(i)} \rangle)$

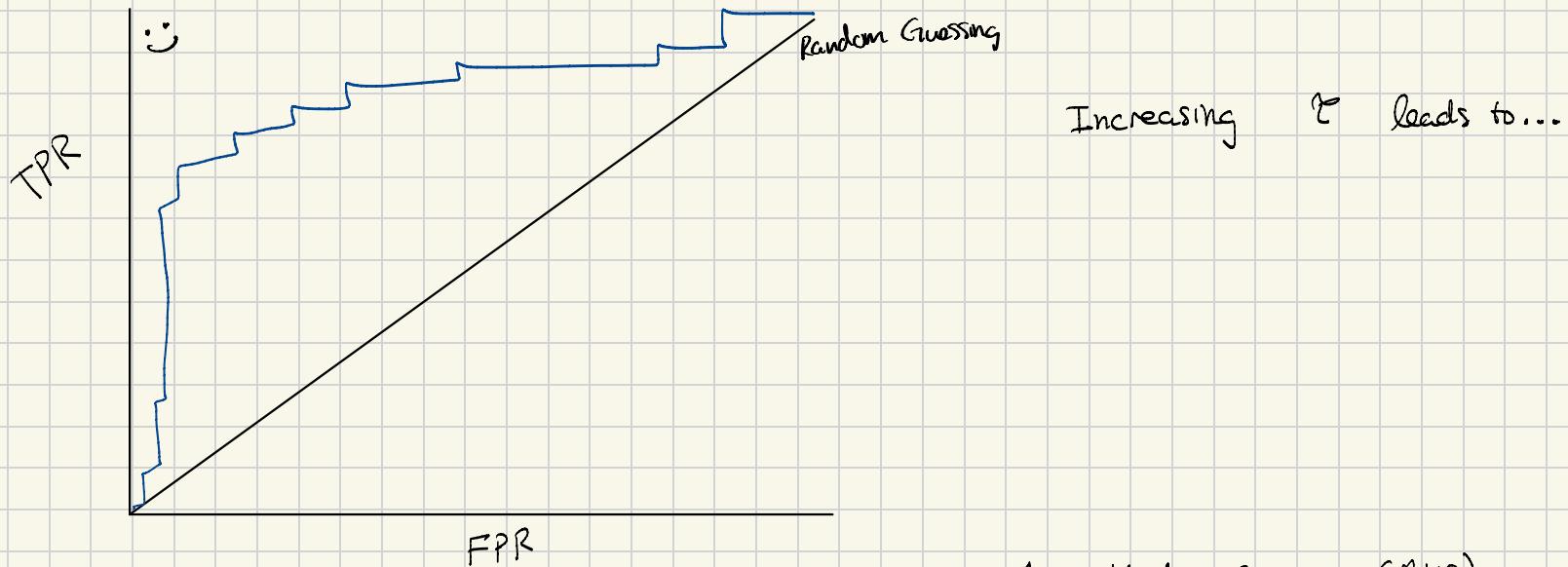
Generally, $\hat{y}^{(i)} = p^{(i)} > 1/2$

↑
choose any
threshold!



Each τ gives a different
FPR and TPR!

Receiver Operating Characteristic Curve



Area Under Curve (AUC)
gives single number
of performance

Multiple Classes

$$\underline{x}^{(i)} \in \mathbb{R}^d$$

$$y^{(i)} \in \{0, 1, \dots, K-1\}$$

$$f: \mathbb{R}^d \rightarrow [0, 1]^K$$

$$\begin{bmatrix} \langle w_1, x \rangle \\ \langle w_2, x \rangle \\ \vdots \\ \langle w_K, x \rangle \end{bmatrix}$$

make probabilities

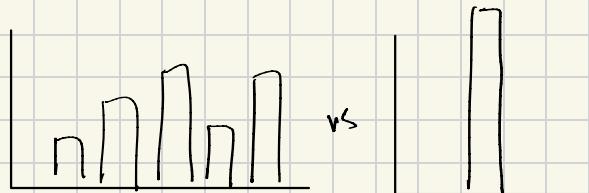
Attempt #1: Sigmoid

Attempt #2: Softmax

$$\begin{bmatrix} z_1 \\ z_2 \\ \vdots \\ z_K \end{bmatrix} \rightarrow \begin{bmatrix} e^{z_1} / \sum_{k=1}^K e^{z_k} \\ e^{z_2} / \sum_{k=1}^K e^{z_k} \\ \vdots \\ e^{z_K} / \sum_{k=1}^K e^{z_k} \end{bmatrix}$$

Cross Entropy Loss

$$\mathcal{L}(w) = - \sum_{l=1}^K \mathbb{I}[y=l] \log [f_j(x)]$$



Computing the Gradient ... :-)