

Day 22, Oct-24, 2024.

→ Foundations of Deep learning

Understanding Some Probability Related Concepts & terms

# Flipping Coins (log likelihood function)

# Using Input (with  $\text{log}(y|\pi)$ )

# Decision Boundary & Non-linear Decision Boundary

# Adaptability - Learning  $\pi(x)$  or function-

$$\log P(y|\pi) \Rightarrow$$

$$\log \left[ \prod_{i=1}^n \pi^{I(y_i=1)} (1-\pi)^{I(y_i=0)} \right]$$

$$\Rightarrow \sum_{i=1}^n I(y_i=1) \log \pi + I(y_i=0) \log(1-\pi)$$

# log-likelihood function  
for a series of coin  
flips.

- ▶ Imagine  $n$  flips of a coin with probability  $\pi$  of coming up heads
- ▶ Loglikelihood function (independent bernoulli trials)

$$\begin{aligned} \log p(y|\pi) &= \log \left[ \prod_{i=1}^n \pi^{I(y_i=1)} (1-\pi)^{I(y_i=0)} \right] \\ &= \sum_{i=1}^n I(y_i=1) \log \pi + I(y_i=0) \log(1-\pi) \end{aligned}$$

where we use the shorthand  $y = (y_1, \dots, y_n)$ .

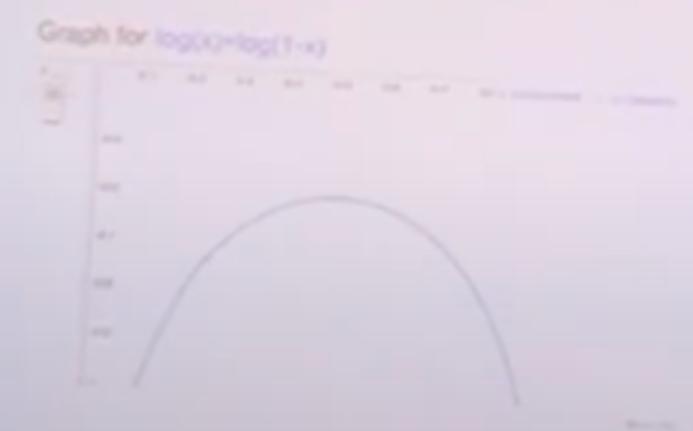


Figure: Logliklihoods over  $\pi$  when we see one head and one tail. Seems likely that  $\pi \approx 0.5$ .

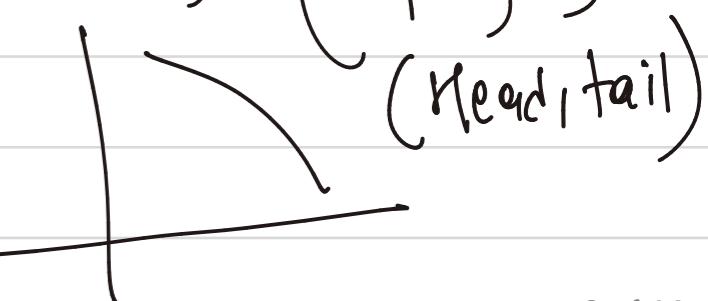
"gf is used to estimate the probability of getting heads ( $\pi$ ) based on the Observed Outcomes ( $y$ ) of the Coin flips."

$$\log p(y|\pi) = \log \left[ \prod_{i=1}^n \pi^{y_i=1} (1-\pi)^{y_i=0} \right]$$

$$\Rightarrow \sum_{i=1}^n \pi^{y_i=1} \log \pi + \pi^{y_i=0} \log (1-\pi)$$

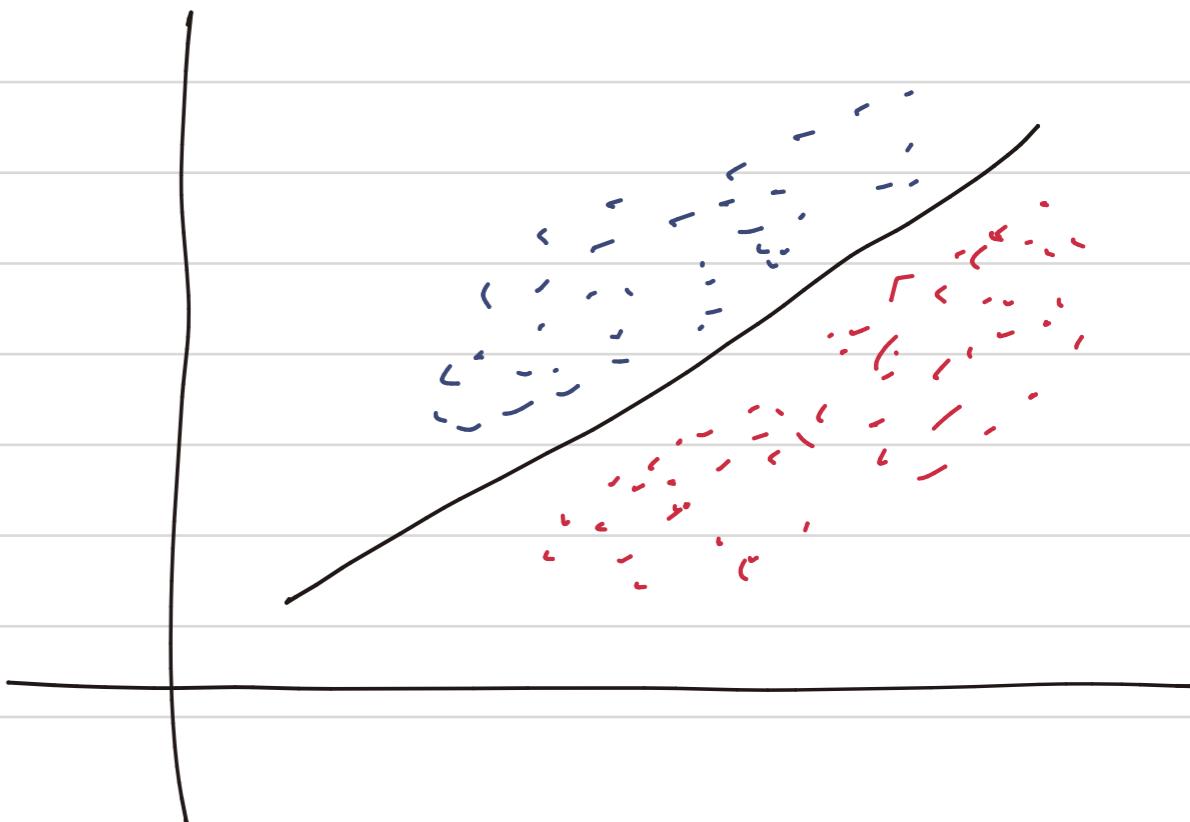
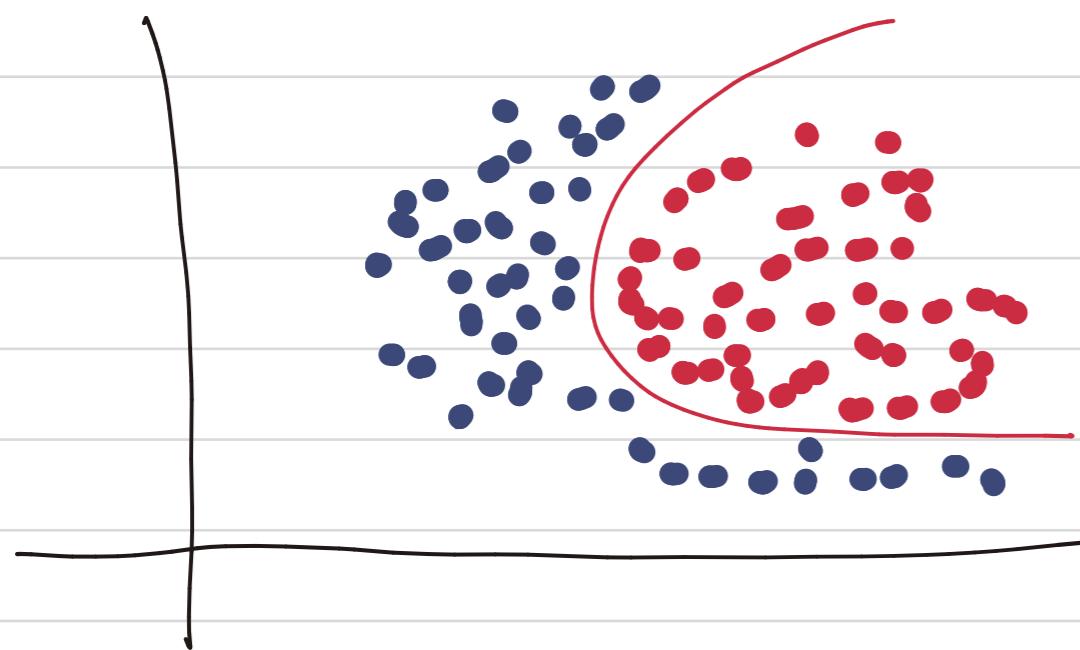
where we use the shorthand  $y = (y_1, y_2, \dots, y_n)$   $(3, 7) \rightarrow$

if we use  $(0, 1) \rightarrow (\text{Head}, \text{tail})$  then



## 2) Using Input Data

- What if the probability head depend on input data?
- In the ANN, our model predict using input features. So, in the same way some input  $x_i$  we have probability  $\pi(x_i)$  of heads?



→ So the model prediction depend upon the characteristics of data.

→ Using data to influence the probability of an event introduce the idea of conditional probability where the probability of an event is not fixed but depends on additional information (input - data)

# Binomial Distribution:

Where  $P(X=k)$

$$C(n, k) = \frac{n!}{k! \times (n-k)!}$$

$$P(X=k) = C(n, k) \times p^k \times (1-p)^{n-k}$$

Each trial has two possible outcomes

$n \rightarrow$  total number of trials

$k \Rightarrow$  number of successes

$$\log(p(y|\pi)) = \log \left[ \pi \left( \pi^{I(y_i=1)} + (1-\pi)^{I(y_i=0)} \right) \right] \quad - \text{eqn 1}$$

$$= \sum_{i=1}^n (\pi I(y_i=1) + \log \pi + \pi I(y_i=0) + \log(1-\pi))$$

$$\log p(y/x) = \sum_{i=1}^n \pi(y_i=1) \log \pi(x) + \pi(y_i=0) \log(1-\pi(x)) \quad - \text{eqn 14}$$

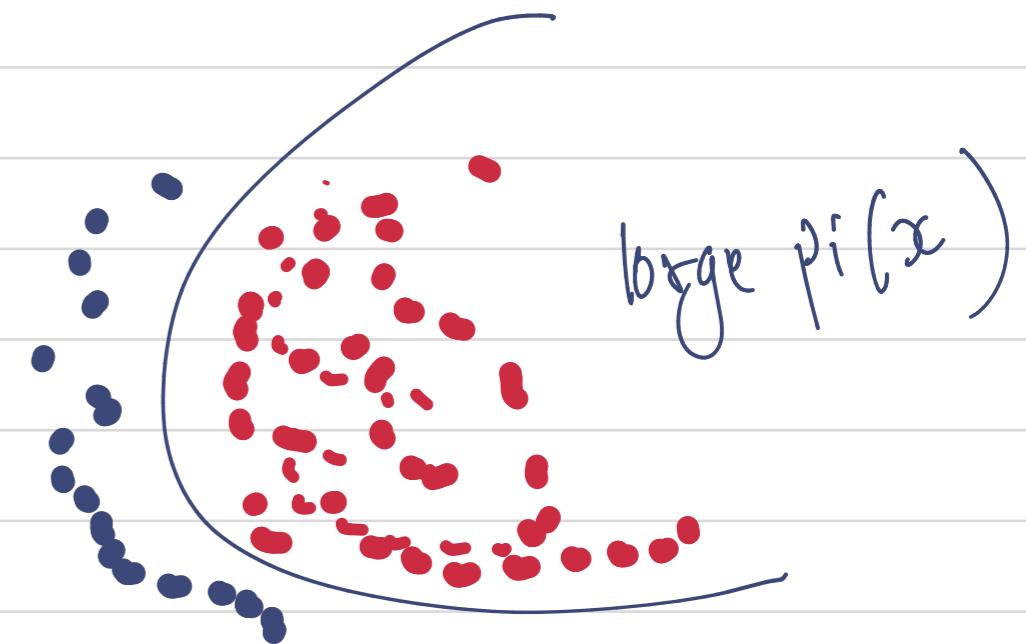
Now Instead of having Same probability for every person

I have probability that depends upon each person.

Before we did maximum likelihood for  $\pi$  now we need to learn some function  $\pi_i(x)$  that explains the data well.

The only difference between them  
eqn(i) and eqn(ii) is  $\pi_i(x)$

$\pi$  of  $x$ .



Learning  $\pi(x)$ : figuring out the best  
way to calculate the probability of Heads based on the Input data.

- ① Log-likelihood of flipping a coin is calculated, estimating  $\pi$   
(Probability of heads)
- ② Probability of Heads depends on input data ( $x_i$ ), shown as Classification.
- ③ Log-likelihood is adapted to include input data, learning  $\pi(x)$ .

## Flipping Coins

- ▶ Imagine  $n$  flips of a coin with probability  $\pi$  of coming up heads
- ▶ Loglikelihood function (independent bernoulli trials)

$$\begin{aligned}\log p(y|\pi) &= \log \left[ \prod_{i=1}^n \pi^{\mathbb{I}(y_i=1)} (1-\pi)^{\mathbb{I}(y_i=0)} \right] \\ &= \sum_{i=1}^n \mathbb{I}(y_i=1) \log \pi + \mathbb{I}(y_i=0) \log (1-\pi)\end{aligned}$$

where we use the shorthand  $y = (y_1, \dots, y_n)$ .

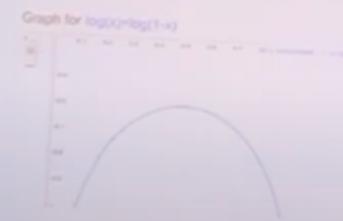


Figure: Loglikelihods over  $\pi$  when we see one head and one tail. Seems likely that  $\pi \approx 0.5$ .

## Using Input Data

- ▶ What if the probability depended on input data?
- ▶ Depending on some input  $x_i$ , we have probability  $\pi(x_i)$  of heads



Figure: Example of a classification problem when  $x_i$  are two dimensional, heads are green, and tails are red.

## Using Input Data

### Using Input Data

- ▶ What if the probability depended on input data?
- ▶ Depending on some input  $x_i$ , we have probability  $\pi(x_i)$  of heads



Figure: Example with a nonlinear boundary

- ▶ Can adapt previous loglikelihood to this setting

$$\log p(y|x) = \sum_{i=1}^n \mathbb{I}(y_i=1) \log \pi(x) + \mathbb{I}(y_i=0) \log (1 - \pi(x))$$

- ▶ Need to learn  $\pi(x)$ , just like we learned  $\pi$  before

# # Gradient Descent of Logistic Regression

(1) We need

derivative of

Sigmoid Function

(2) Partial Derivatives

of weights  $w_j$ .

w.r.t. Cost  $J(\vec{w}, b)$

(3) Gradient with respect of  $w_j$  & Gradient w.r.t.  $b$

## Gradient descent

cost

$$J(\vec{w}, b) = -\frac{1}{m} \sum_{i=1}^m [y^{(i)} \log(f_{\vec{w}, b}(\vec{x}^{(i)})) + (1 - y^{(i)}) \log(1 - f_{\vec{w}, b}(\vec{x}^{(i)}))]$$

repeat {

$$\underbrace{w_j = w_j - \alpha \frac{\partial}{\partial w_j} J(\vec{w}, b)}$$

$$b = b - \alpha \underbrace{\frac{\partial}{\partial b} J(\vec{w}, b)}$$

}

$$\frac{\partial}{\partial w_j} J(\vec{w}, b) = \frac{1}{m} \sum_{i=1}^m (f_{\vec{w}, b}(\vec{x}^{(i)}) - y^{(i)}) \underbrace{x_j^{(i)}}$$

$$\frac{\partial}{\partial b} J(\vec{w}, b) = \frac{1}{m} \sum_{i=1}^m (f_{\vec{w}, b}(\vec{x}^{(i)}) - y^{(i)})$$

$$\textcircled{1} \quad f_{\vec{w}, b}(\vec{x}) = \frac{1}{1 + e^{-(\vec{w} \cdot \vec{x} + b)}}$$

$$\textcircled{2} \quad J(\vec{w}, b) = -\frac{1}{m} \sum_{i=1}^m \left[ y^{(i)} \cdot \log(f_{\vec{w}, b}(\vec{x}^{(i)})) + (1-y^{(i)}) \cdot \log(1-f_{\vec{w}, b}(\vec{x}^{(i)})) \right]$$

\textcircled{3} Gradient with Respect to  $w_j$ :

$$\frac{\partial}{\partial w_j} J(\vec{w}, b) = \frac{1}{m} \sum_{i=1}^m (f_{\vec{w}, b}(\vec{x}^{(i)}) - y^{(i)}) x_j^{(i)}$$

where  $f_{\vec{w}, b}(\vec{x}^{(i)} - y^{(i)}) \rightarrow$  difference between (Actual - Prediction)

Q

Gradient with respect to  $b$ .

$$\frac{\partial}{\partial b} J(\vec{w}, b) = \frac{1}{m} \sum_{i=1}^m (f_{\vec{w}, b}(\vec{x}^{(i)}) - y^{(i)})$$

Q

So, update rule -

$$w_j := w_j - \alpha \frac{\partial}{\partial w_j} J(\vec{w}, b)$$

$$w_j := w_j - \alpha \frac{1}{m} \sum_{i=1}^m (f_{\vec{w}, b}(\vec{x}^{(i)}) - y^{(i)}) x_j^{(i)}$$

$$b := b - \alpha \frac{\partial}{\partial b} J(\vec{w}, b)$$
$$b := b - \alpha \frac{1}{m} \sum_{i=1}^m (f_{\vec{w}, b}(\vec{x}^{(i)}) - y^{(i)})$$

# Derivative of Sigmoid function

①

$$F_{\vec{w}, b}(\vec{x}^{(i)}) \Rightarrow \frac{1}{1 + e^{-(\vec{w} \cdot \vec{x}^{(i)} + b)}}$$

②

$$f(z) \Rightarrow \frac{1}{1 + e^{-z}}$$

Where,  $z = \vec{w} \cdot \vec{x}^{(i)} + b$

So,  $\frac{d}{dz} f(z) \Rightarrow f(z) \left( \frac{1}{1 - f(z)} \right)$  eqn ①

③

Derivative of  $z$  w.r.t.  $w_j^o \rightarrow z = \vec{w} \cdot \vec{x}^{(i)} + b$  w.r.t  $w_j^o$

$$z = w_1 x_1^{(i)} + w_2 x_2^{(i)} + \dots + w_j^o x_j^{(i)} + \dots + b$$

Derivative of  $z$  wrt.  $w_j^o$   $\left[ \frac{\partial z}{\partial w_j^o} \Rightarrow x_j^{(i)} \right]$

Simplly where  $z$  is linear combination of the weights and features. eqn ②



$$\frac{\partial}{\partial w_j} f_{\vec{w}, b}(\vec{x}^{(i)}) \Rightarrow \frac{d}{dz} f(z) \cdot \frac{\partial z}{\partial w_j} \longrightarrow \text{eqn III}$$

Substituting eqn ⑨ and eqn ⑩ — in eqn ⑪

$$\frac{\partial}{\partial w_j} f_{\vec{w}, b}(\vec{x}^{(i)}) \Rightarrow f(z)(1 - f(z)) \cdot x_j^{(i)}$$

But  $f(z)(1 - f(z))$  — eqn ④ has  $f(z)$

and  $f(z)$  can be written as

$$f(z) \Rightarrow f_{\vec{w}, b}(\vec{x}^{(i)})$$

So, the eqn (iii) becomes,

$$\frac{\partial}{\partial w_j} f_{\vec{w}, b}(\vec{x}^{(i)}) \Rightarrow f_{\vec{w}, b}(\vec{x}^{(i)}) (1 - f_{\vec{w}, b}(\vec{x}^{(i)})) \cdot x_j^{(i)}$$

which becomes eqn (iv)

eqn (iv) helps update the weight  $w_j$  during each step of gradient descent.

Since, Gradient descent is optimization algorithm that aim to find the minima and minimize the cost function and optimize the model parameters to reduce the prediction error.