# CS 260: Foundations of Data Science

Prof. Sara Mathieson

Fall 2021

HAVERFORD
COLLEGE

# Admin

- I'm giving a talk "at" a conference during our class on Thursday Oct 21
  - There will still be lab though!

- Watch this lecture video instead and do **Handout 14**

- Come prepared to discuss **Handout 14** on Tuesday Oct 26

# Outline for October 21

- Introduction to logistic regression

- Cost function and SGD for logistic regression

- Connection to cross entropy

# Outline for October 21

- Introduction to logistic regression

- Cost function and SGD for logistic regression

- Connection to cross entropy

# Why is linear regression a bad choice for classification?

**Case Study**: you need to identify the medical condition of a patient in the emergency room on the basis of their symptoms.

Possible conditions ($y$) are:
- Stroke
- Drug overdose
- Epileptic seizure

1) If you were forced to use linear regression for this problem, how could you encode $y$ to make it real-valued?

2) What issues arise with making y real-valued?

3) What if you just had two outcomes (i.e. stroke and drug overdose) -- why is linear regression still not a good choice?

# Why is linear regression a bad choice for classification?

**Case Study**: you need to identify the medical condition of a patient in the emergency room on the basis of their symptoms.

Possible conditions ($y$) are:
- Stroke
- Drug overdose
- Epileptic seizure

1) If you were forced to use linear regression for this problem, how could you encode $y$ to make it real-valued?

You could choose stroke=0, drug overdose=1, epileptic seizure=2 (or some permutation)

2) What issues arise with making y real-valued?

3) What if you just had two outcomes (i.e. stroke and drug overdose) -- why is linear regression still not a good choice?

# Why is linear regression a bad choice for classification?

**Case Study**: you need to identify the medical condition of a patient in the emergency room on the basis of their symptoms.

Possible conditions ($y$) are:
- Stroke
- Drug overdose
- Epileptic seizure

1) If you were forced to use linear regression for this problem, how could you encode $y$ to make it real-valued?

You could choose stroke=0, drug overdose=1, epileptic seizure=2 (or some permutation)

2) What issues arise with making y real-valued?

Assumes some *ordering* of the outcomes that is probably not there!

3) What if you just had two outcomes (i.e. stroke and drug overdose) -- why is linear regression still not a good choice?

# Why is linear regression a bad choice for classification?

**Case Study**: you need to identify the medical condition of a patient in the emergency room on the basis of their symptoms.

Possible conditions ($y$) are:
- Stroke
- Drug overdose
- Epileptic seizure

1) If you were forced to use linear regression for this problem, how could you encode $y$ to make it real-valued?

You could choose stroke=0, drug overdose=1, epileptic seizure=2 (or some permutation)

2) What issues arise with making y real-valued?

Assumes some *ordering* of the outcomes that is probably not there!

3) What if you just had two outcomes (i.e. stroke and drug overdose) -- why is linear regression still not a good choice?

The range of a linear function (i.e. $y$ values) is [-∞, ∞], but we want [0, 1]

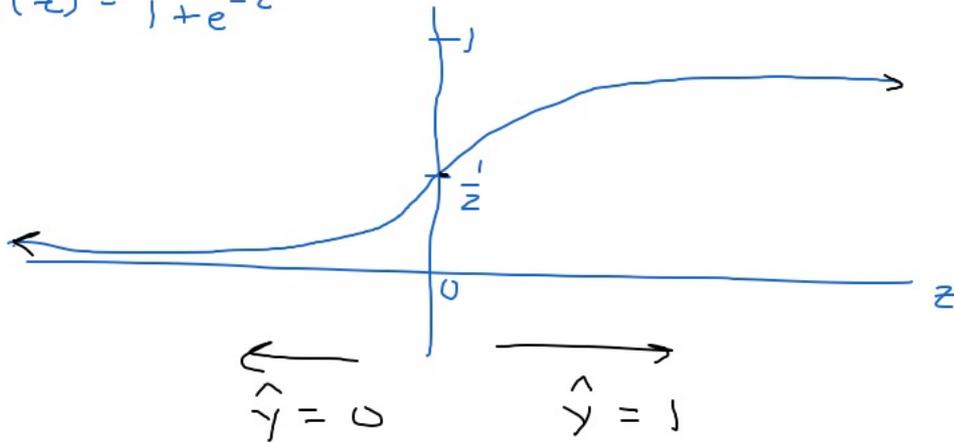# Logistic Regression Intro

$y \in \{0, 1\}$  binary classification

linear regression: weight on each feature ☆

$$X \qquad\qquad Y$$

linear:  $[-\infty, \infty] \longrightarrow [-\infty, \infty]$

logistic:  $[-\infty, \infty] \longrightarrow [0, 1]$  $\}$ can make discrete

$\underbrace{\phantom{[0,1]}}_{\text{prob}}$

$\overbrace{\phantom{p(y=1|\vec{x})}}^{\text{posterior}}$

<u>idea</u>  model:  $h_{\vec{w}}(\vec{x}) = p(y=1 | \vec{x}) = \dfrac{1}{1 + e^{\boxed{\vec{w} \cdot \vec{x}}}}$  linear function

$g(z) = \dfrac{1}{1 + e^{-z}}$



$z \longrightarrow \infty \quad, \quad g(z) \rightarrow 1$

$z \longrightarrow -\infty \quad, \quad g(z) \rightarrow 0$

$z = 0 \quad, \quad g(z) = \dfrac{1}{2}$

# Logistic (sigmoid) function

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

# Logistic Regression Decision Boundaries

if $\boxed{\begin{array}{ll} \vec{w} \cdot \vec{x} > 0 & \Rightarrow \hat{y} = 1 \\ \vec{w} \cdot \vec{x} \leq 0 & \Rightarrow \hat{y} = 0 \end{array}}$

$\vec{w} = \begin{bmatrix} 3 \\ -2 \end{bmatrix}$, $p = 1$

$h(\vec{x}) = \dfrac{1}{1 + e^{-(3 - 2x)}}$ $\Big\}$ model

$= p(y = 1 | \vec{x}) > \dfrac{1}{2}$

Q: what is classified as
$\hat{y} = 1$ ?

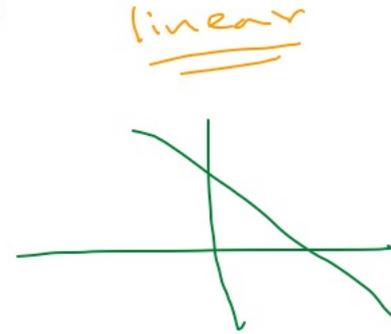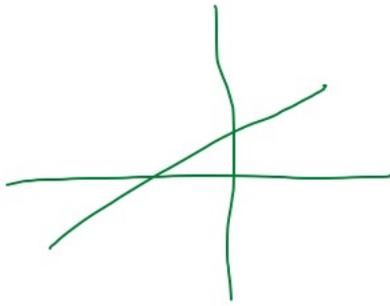$\dfrac{1}{1 + e^{-(3 - 2x)}} > \dfrac{1}{2}$

$2 > 1 + e^{-3 + 2x}$
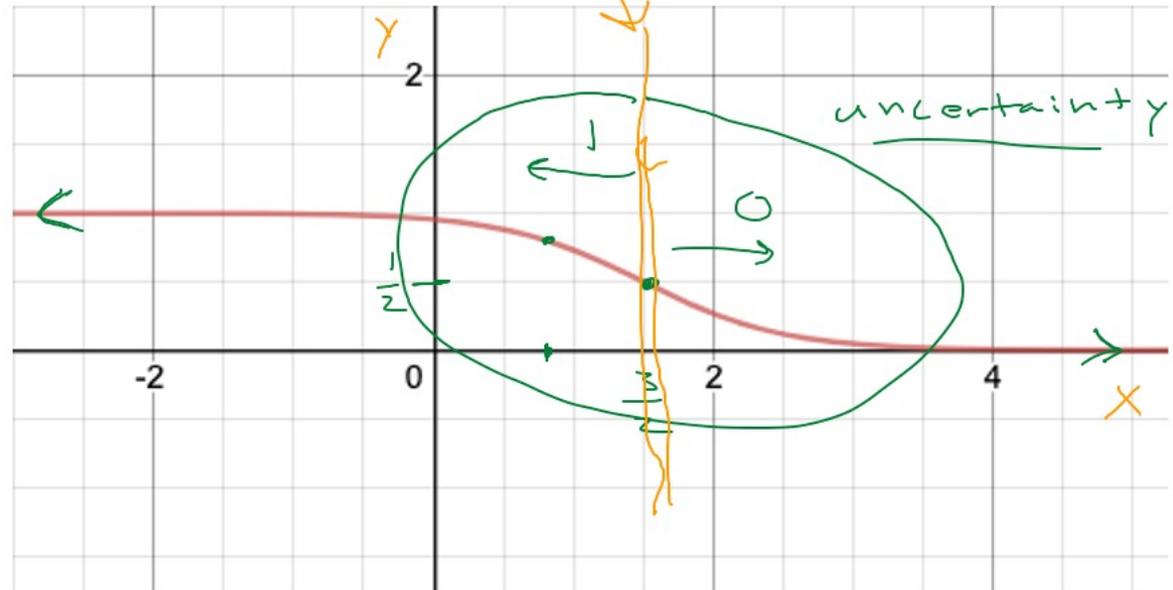
$1 > e^{-3 + 2x}$

$0 > -3 + 2x$

$3 > 2x$

$\boxed{x < \dfrac{3}{2}} \Rightarrow \hat{y} = 1$

# Logistic Regression Decision Boundaries

linear

$x < \dfrac{3}{2}$

$$h_{\vec{w}}(\vec{x}) = \dfrac{1}{\left(1 + e^{-(3-2x)}\right)}$$

y

2

1

0

uncertainty

$\dfrac{1}{2}$

-2

0

$\dfrac{3}{2}$

2

4

x

# Outline for October 21

- Introduction to logistic regression

- Cost function and SGD for logistic regression

- Connection to cross entropy

# Logistic Regression Cost Function

How do we find $\vec{w}$?   need cost function!  $\Rightarrow$ SGD

likelihood: $L(\vec{w}) = \prod\limits_{i=1}^{n} \underbrace{h_{\vec{w}}(\vec{x}_i)}_{\text{prob } y_i=1}^{\overset{\text{if } y_i=1}{\boxed{y_i}}} \underbrace{\left(1 - h_{\vec{w}}(\vec{x}_i)\right)}_{\text{prob } y_i=0}^{\overset{\text{if } y_i=0}{\boxed{1-y_i}}}$
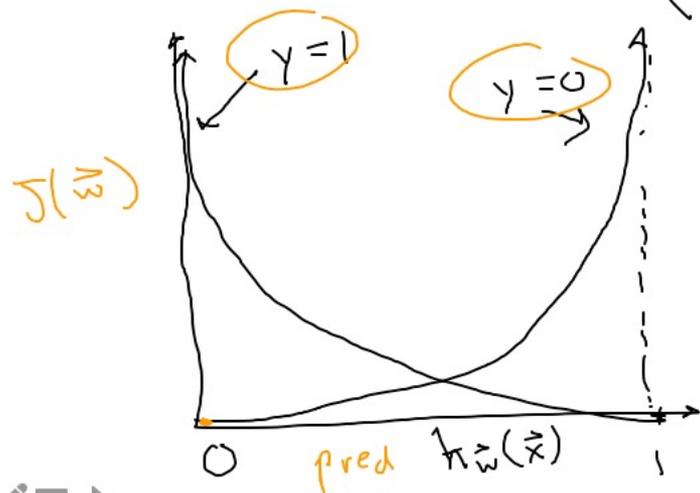
maximize

minimize negative log likelihood = cost

$$J(\vec{w}) = -\log L(\vec{w}) = -\sum\limits_{i=1}^{n} y_i \, h_{\vec{w}}(\vec{x}_i) + (1-y_i) \log\left(1 - h_{\vec{w}}(\vec{x}_i)\right)$$

$\underset{\log}{\wedge}$

single example

$(\vec{x}, y)$

$$J(\vec{w}) = \begin{cases} -\log h_{\vec{w}}(\vec{x}) & \text{if } y = 1 \\ -\log\left(1 - h_{\vec{w}}(\vec{x})\right) & \text{if } y = 0 \end{cases}$$



$y=1$

$y=0$

$J(\vec{w})$

$O$   pred $h_{\vec{w}}(\vec{x})$   $1$

# Stochastic Gradient Descent for Logistic Regression (binary classification)

set w = 0 vector

while cost J(w) still changing:

    shuffle data points

    for i = 1...n:

        w <- w − alpha(derivative of J(w) wrt $x_i$)

    store J(w)

# 3 important pieces to SGD

- Hypothesis function (prediction)

$$h_{\boldsymbol{w}}(\boldsymbol{x}) = p(y = 1 | \boldsymbol{x}) = \frac{1}{1 + e^{-\boldsymbol{w} \cdot \boldsymbol{x}}}$$

# 3 important pieces to SGD

- Hypothesis function (prediction)

$$h_{\boldsymbol{w}}(\boldsymbol{x}) = p(y = 1|\boldsymbol{x}) = \frac{1}{1 + e^{-\boldsymbol{w} \cdot \boldsymbol{x}}}$$

- Cost function (want to minimize)

$$J(\boldsymbol{w}) = -\sum_{i=1}^{n} y_i \log h_{\boldsymbol{w}}(\boldsymbol{x_i}) + (1 - y_i) \log(1 - h_{\boldsymbol{w}}(\boldsymbol{x_i}))$$

# 3 important pieces to SGD

- Hypothesis function (prediction)

$$h_{\boldsymbol{w}}(\boldsymbol{x}) = p(y = 1|\boldsymbol{x}) = \frac{1}{1 + e^{-\boldsymbol{w} \cdot \boldsymbol{x}}}$$

- Cost function (want to minimize)

$$J(\boldsymbol{w}) = -\sum_{i=1}^{n} y_i \log h_{\boldsymbol{w}}(\boldsymbol{x_i}) + (1 - y_i) \log(1 - h_{\boldsymbol{w}}(\boldsymbol{x_i}))$$

- Gradient of cost wrt single data point $x_i$

$$\nabla J_{\boldsymbol{x_i}}(\boldsymbol{w}) = (h_{\boldsymbol{w}}(\boldsymbol{x_i}) - y_i)\boldsymbol{x_i}$$

# 3 important pieces to SGD

$$\vec{w} \leftarrow \vec{w} - \alpha(h_{\vec{w}}(\vec{x_i}) - y_i)\vec{x_i}$$

pred   truth

- **Hypothesis function (prediction)**

$$h_{\boldsymbol{w}}(\boldsymbol{x}) = p(y = 1|\boldsymbol{x}) = \frac{1}{1 + e^{-\boldsymbol{w}\cdot\boldsymbol{x}}}$$

- **Cost function (want to minimize)**

$$J(\boldsymbol{w}) = -\sum_{i=1}^{n} y_i \log h_{\boldsymbol{w}}(\boldsymbol{x_i}) + (1 - y_i)\log(1 - h_{\boldsymbol{w}}(\boldsymbol{x_i}))$$

- **Gradient of cost wrt single data point x$_i$**

same
form as
linear
reg!

$$\nabla J_{\boldsymbol{x_i}}(\boldsymbol{w}) = (h_{\boldsymbol{w}}(\boldsymbol{x_i}) - y_i)\boldsymbol{x_i}$$

$$\frac{1}{1 + e^{-\vec{w}\cdot\vec{x_i}}}$$

# Outline for October 21

- Introduction to logistic regression

- Cost function and SGD for logistic regression
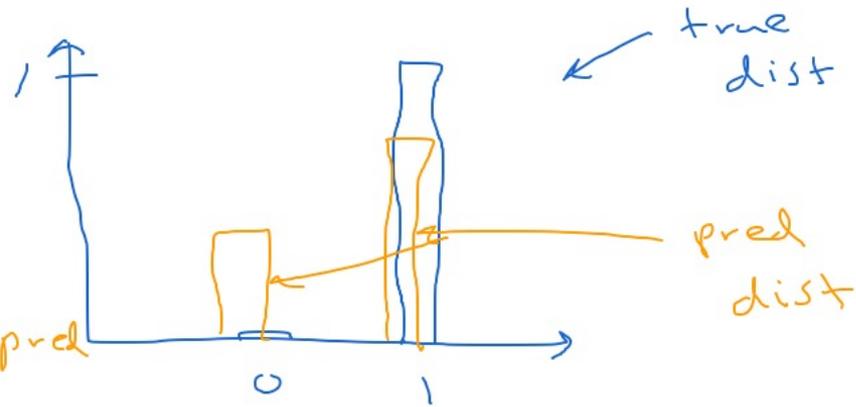
- Connection to cross entropy

# Cross entropy

$$J(\vec{w}) = - y \log h(x) - (1-y) \log(1-h)$$

cost function is

$$H(p,q) = - \sum_x p(x) \log q(x)$$

2 probability distributions

if $y = 1$, $1-y = 0$ } true

$h = 0.6$, $1-h = 0.4$ } pred

high        low

true dist

pred dist

if dist. are exactly the same
$\Rightarrow$ cost / cross-entropy
$= 0$

TODO: Handout 14!