

26: Logistic Regression

Lisa Yan and Jerry Cain

November 11, 2020

Quick slide reference

3	Background	26a_background
9	Logistic Regression	26b_logistic_regression
27	Training: The big picture	26c_lr_training
56	Training: The details, Testing	LIVE
59	Philosophy	LIVE
63	Gradient Derivation	26e_derivation

Background

1. Weighted sum

If $\mathbf{X} = (X_1, X_2, \dots, X_m)$:

$$Z = \theta_1 X_1 + \theta_2 X_2 + \dots + \theta_m X_m$$

$$= \sum_{j=1}^m \theta_j X_j$$

weighted sum

$$= \boldsymbol{\theta}^T \mathbf{X}$$

dot product

$$[\theta_1 \quad \theta_2 \quad \dots \quad \theta_m] \begin{bmatrix} X_1 \\ \vdots \\ X_m \end{bmatrix}$$

1. Weighted sum

Dot product/
weighted sum $\theta^T \mathbf{X} = \sum_{j=1}^m \theta_j X_j$

Recall the linear regression model, where $\mathbf{X} = (X_1, X_2, \dots, X_m)$ and $Y \in \mathbb{R}$:

$$\hat{Y} = g(\mathbf{X}) = \theta_0 + \sum_{j=1}^m \theta_j X_j$$

How would you rewrite this expression as a single dot product?



1. Weighted sum

Dot product/
weighted sum $\theta^T \mathbf{X} = \sum_{j=1}^m \theta_j X_j$

Recall the linear regression model, where $\mathbf{X} = (X_1, X_2, \dots, X_m)$ and $Y \in \mathbb{R}$:

$$g(\mathbf{X}) = \theta_0 + \sum_{j=1}^m \theta_j X_j$$

How would you rewrite this expression as a single dot product?

$$g(\mathbf{X}) = \theta_0 X_0 + \theta_1 X_1 + \theta_2 X_2 + \dots + \theta_m X_m \quad \text{Define } X_0 = 1$$

$$= \theta^T \mathbf{X}$$

$$\text{New } \mathbf{X} = (1, X_1, X_2, \dots, X_m), \quad \theta = (\theta_0, \theta_1, \dots, \theta_m)$$

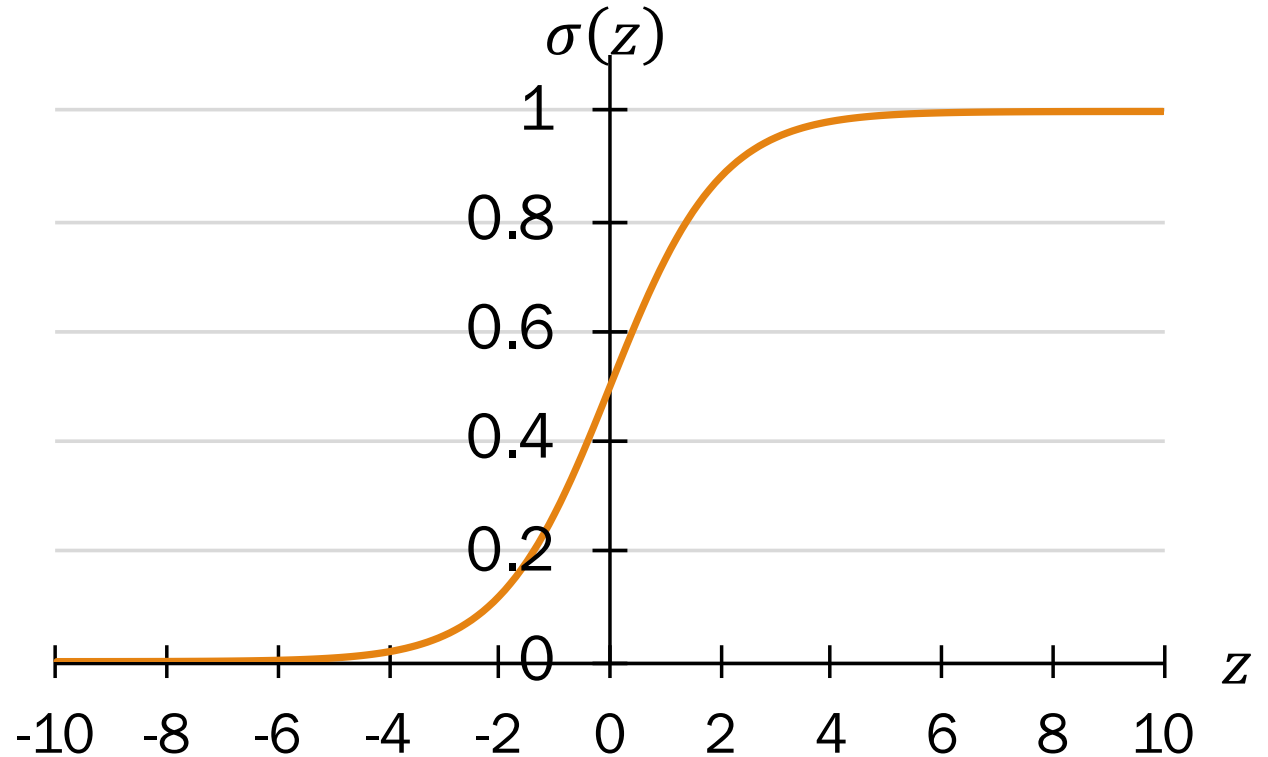
Prepending $X_0 = 1$ to each feature vector \mathbf{X} makes matrix operators more accessible.

2. Sigmoid function $\sigma(z)$

- The sigmoid function:

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

- Sigmoid squashes z to a number between 0 and 1.
- Recall definition of probability:
A number between 0 and 1



$\sigma(z)$ can represent a probability.

3. Conditional likelihood function

Training data (n datapoints):

- $(\mathbf{x}^{(i)}, y^{(i)})$ drawn i.i.d. from a distribution $f(\mathbf{X} = \mathbf{x}^{(i)}, Y = y^{(i)} | \theta) = f(\mathbf{x}^{(i)}, y^{(i)} | \theta)$

$$\theta_{MLE} = \arg \max_{\theta} \prod_{i=1}^n f(y^{(i)} | \mathbf{x}^{(i)}, \theta)$$

conditional likelihood
of training data

$$= \arg \max_{\theta} \sum_{i=1}^n \log f(y^{(i)} | \mathbf{x}^{(i)}, \theta)$$

log conditional likelihood

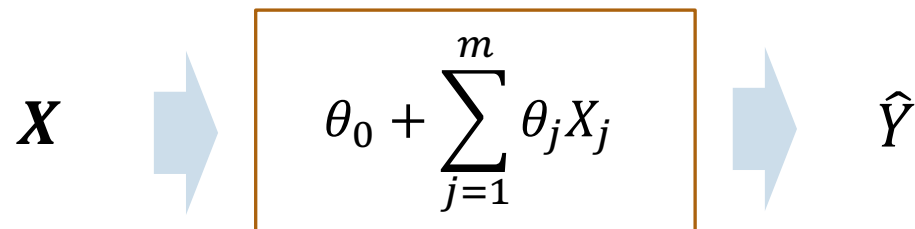
$$= \arg \max_{\theta} LL(\theta)$$

- MLE in this lecture is estimator that maximizes conditional likelihood
- Confusingly, log conditional likelihood is also written as $LL(\theta)$

Logistic Regression

Prediction models so far

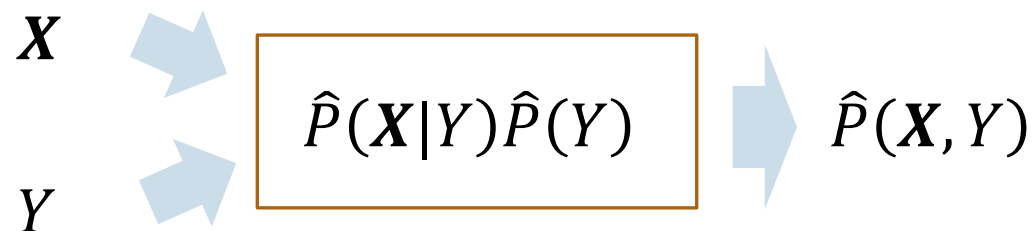
Linear Regression (Regression)



$$\hat{Y} = \theta_0 + \sum_{j=1}^m \theta_j X_j$$

- ✓ \mathbf{X} can be dependent
- 🙋 Regression model ($\hat{Y} \in \mathbb{R}$, not discrete)

Naïve Bayes (Classification)



$$\begin{aligned} \hat{Y} &= \arg \max_{y=\{0,1\}} P(Y | \mathbf{X}) \\ &= \arg \max_{y=\{0,1\}} P(\mathbf{X}|Y)P(Y) \end{aligned}$$

- ✓ Tractable with NB assumption, but...
- ⚠ Realistically, X_j features not necessarily conditionally independent
- 🙋 Actually models $P(\mathbf{X}, Y)$, not $P(Y|\mathbf{X})$?

Introducing Logistic Regression!

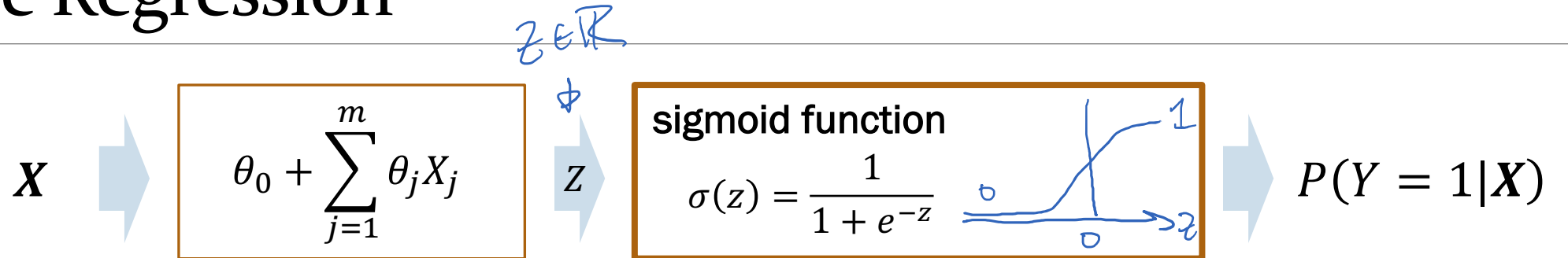


Linear Regression ideas

Classification models

+ *compute power*

Logistic Regression



Logistic Regression
Model:

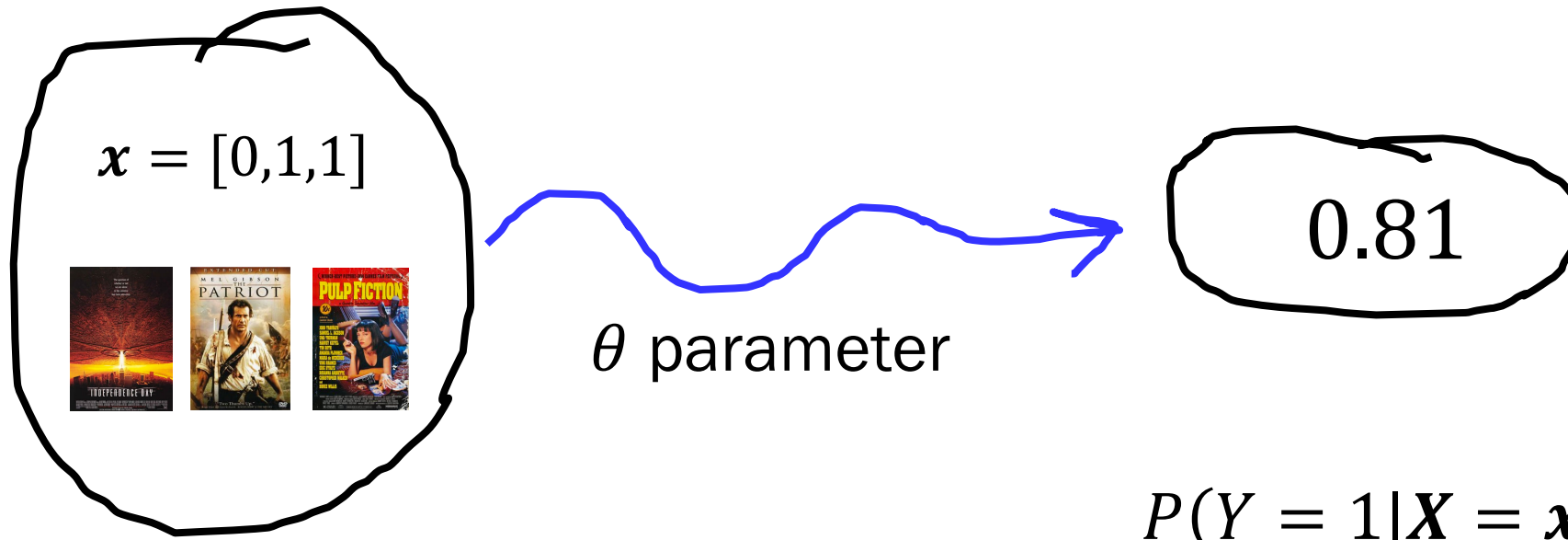
$$P(Y = 1 | \mathbf{X} = \mathbf{x}) = \sigma \left(\theta_0 + \sum_{j=1}^m \theta_j x_j \right)$$

Predict \hat{Y} as the most likely Y
given our observation $\mathbf{X} = \mathbf{x}$:

$$\hat{Y} = \arg \max_{y \in \{0,1\}} P(Y | \mathbf{X})$$

- Since $Y \in \{0,1\}$, $P(Y = 0 | \mathbf{X} = \mathbf{x}) = 1 - \sigma(\theta_0 + \sum_{j=1}^m \theta_j x_j)$
- Sigmoid function also known as “logit” function

Logistic Regression



X
input features

$P(Y = 1 | \mathbf{X} = \mathbf{x})$
conditional likelihood

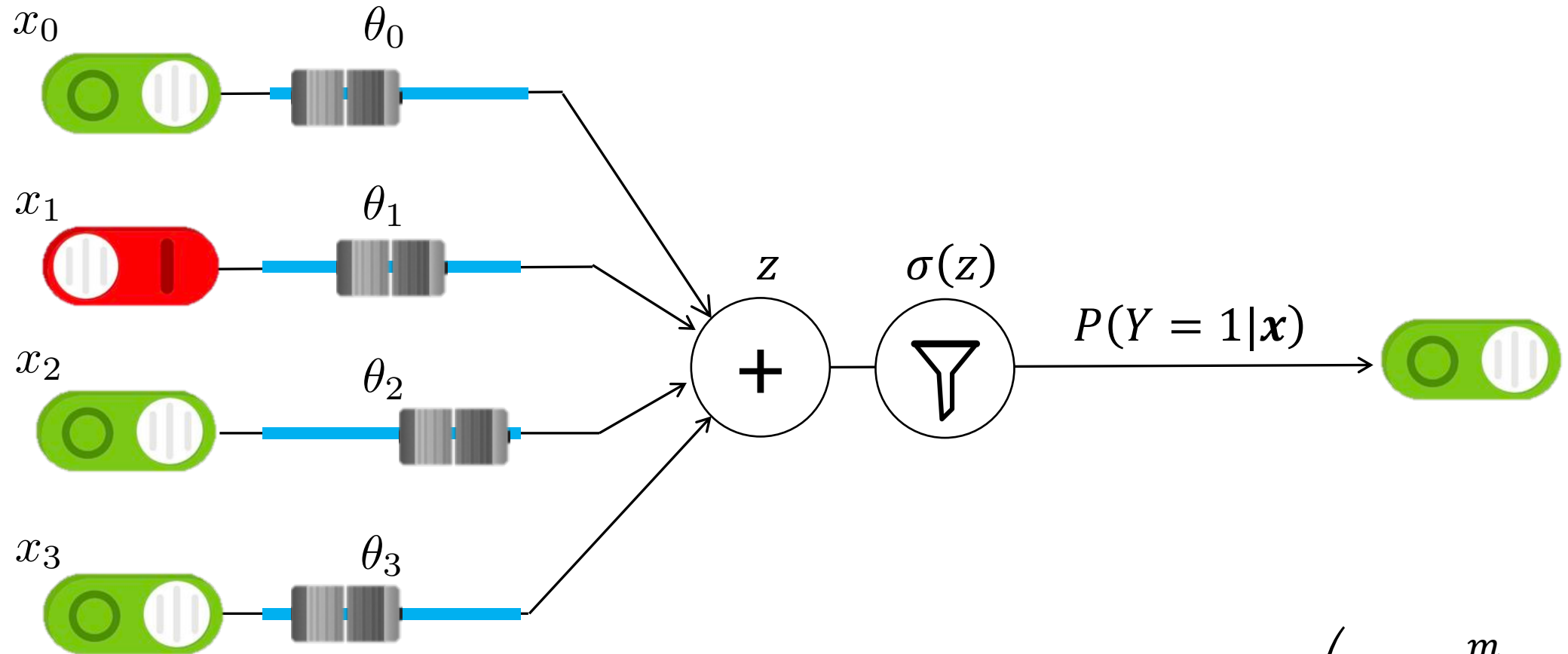
$$P(Y = 1 | \mathbf{X} = \mathbf{x}) = \sigma \left(\theta_0 + \sum_{j=1}^m \theta_j x_j \right)$$

Logistic Regression cartoon



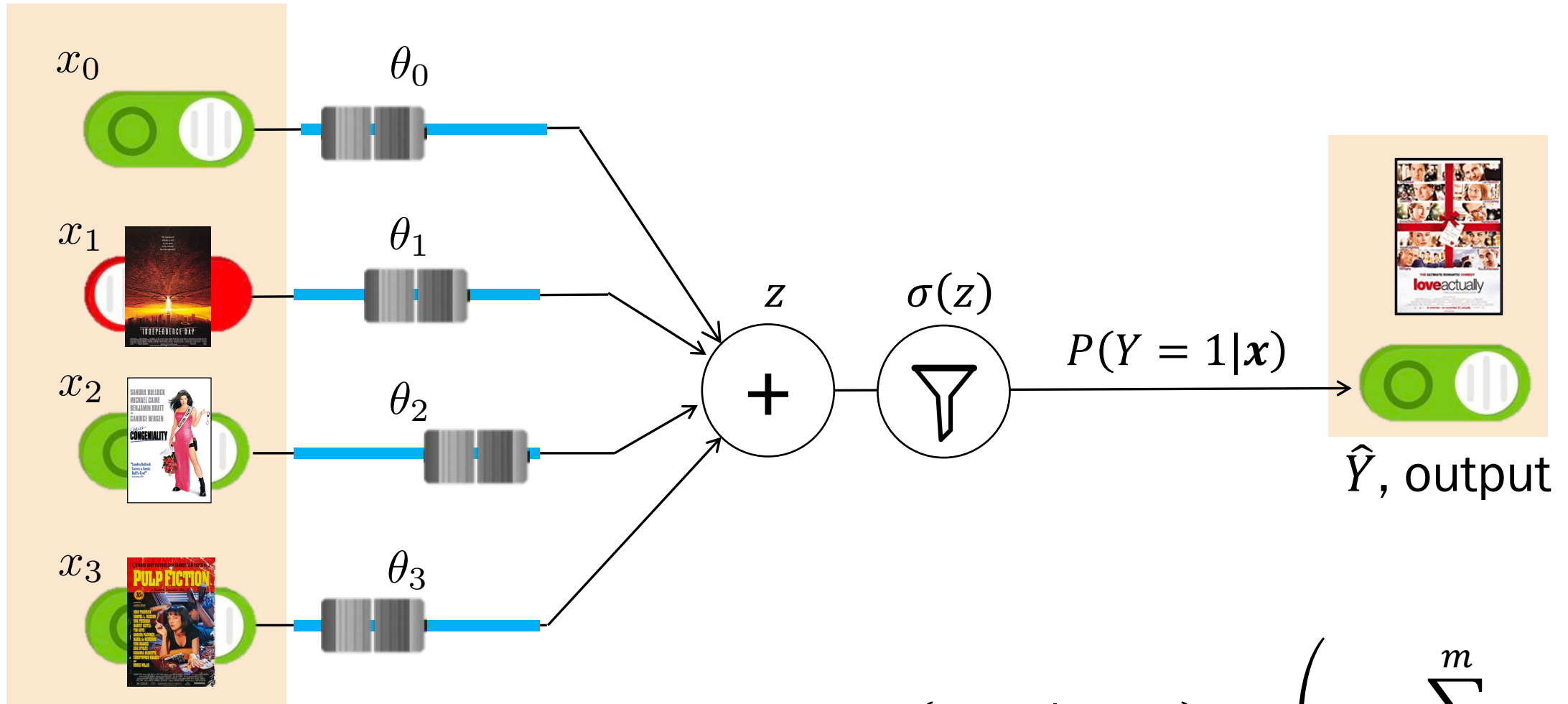
θ parameter

Logistic Regression cartoon



$$P(Y = 1 | \mathbf{X} = \mathbf{x}) = \sigma \left(\theta_0 + \sum_{j=1}^m \theta_j x_j \right)$$

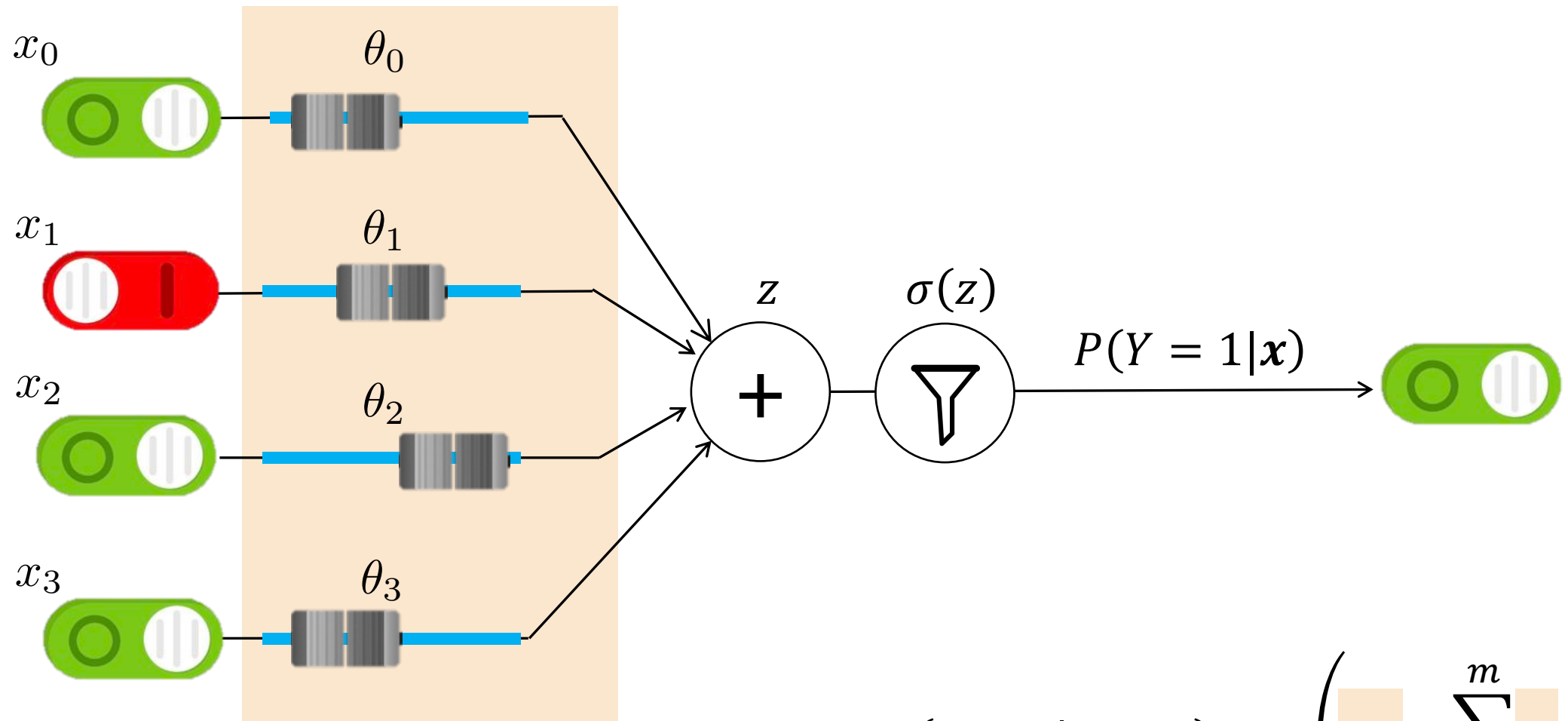
Logistic Regression cartoon



\mathbf{X} , input features
[0,1,1]

$$P(Y = 1 | \mathbf{X} = \mathbf{x}) = \sigma \left(\theta_0 + \sum_{j=1}^m \theta_j x_j \right)$$

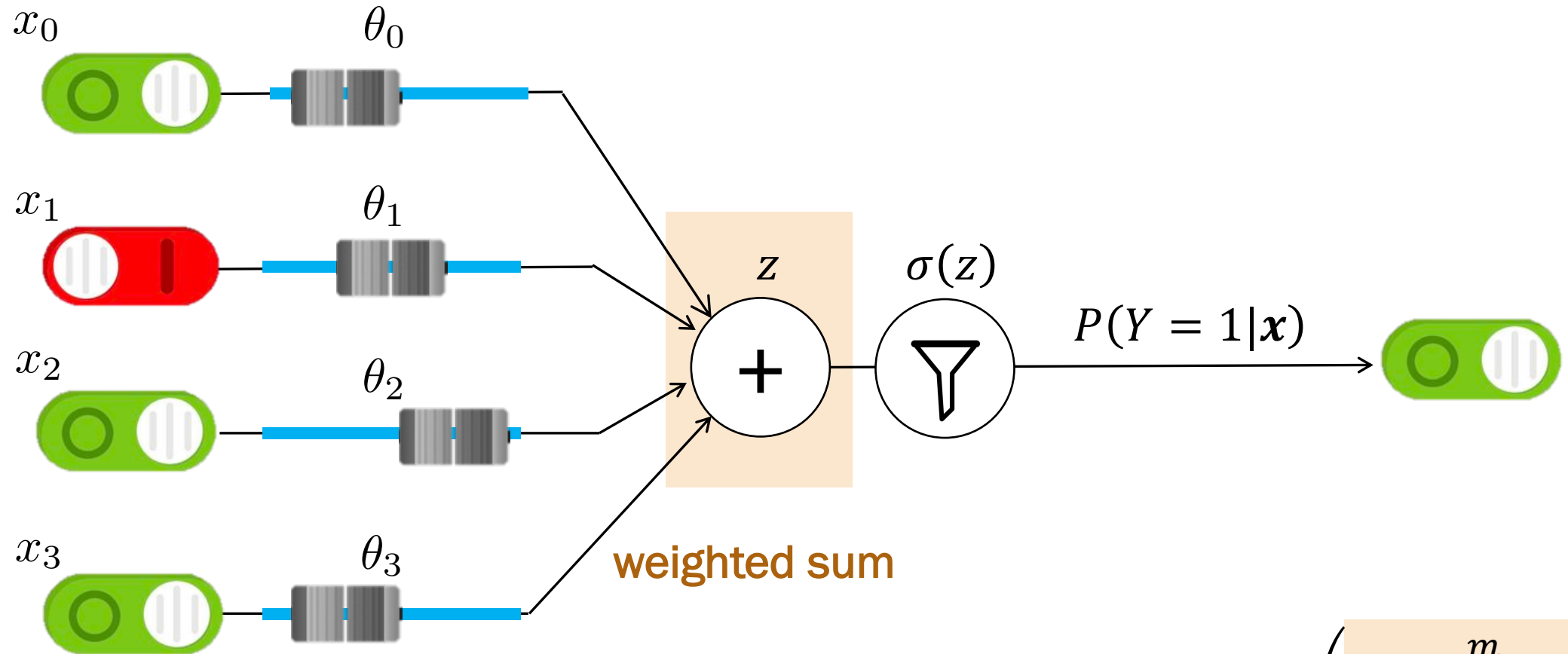
Components of Logistic Regression



θ weights
(aka parameters)

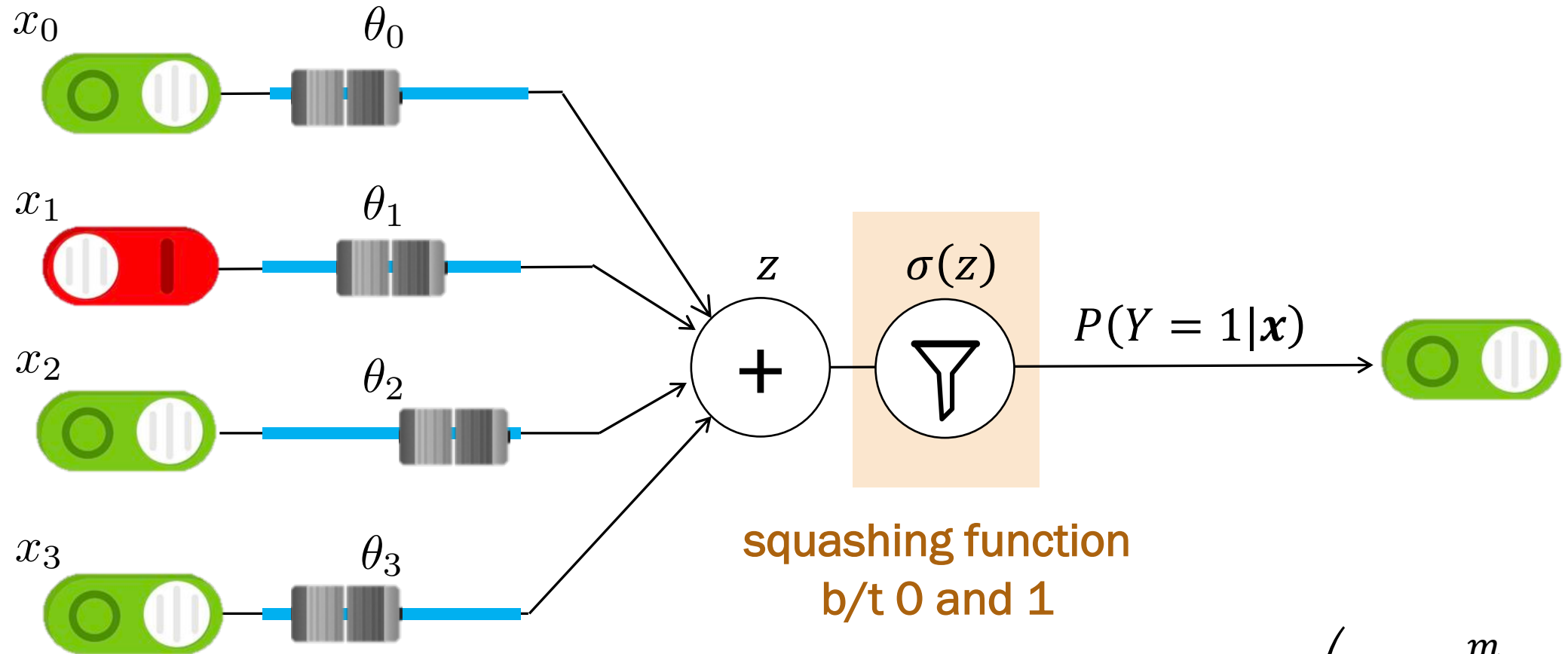
$$P(Y = 1 | \mathbf{X} = \mathbf{x}) = \sigma \left(\theta_0 + \sum_{j=1}^m \theta_j x_j \right)$$

Components of Logistic Regression



$$P(Y = 1 | \mathbf{X} = \mathbf{x}) = \sigma \left(\theta_0 + \sum_{j=1}^m \theta_j x_j \right)$$

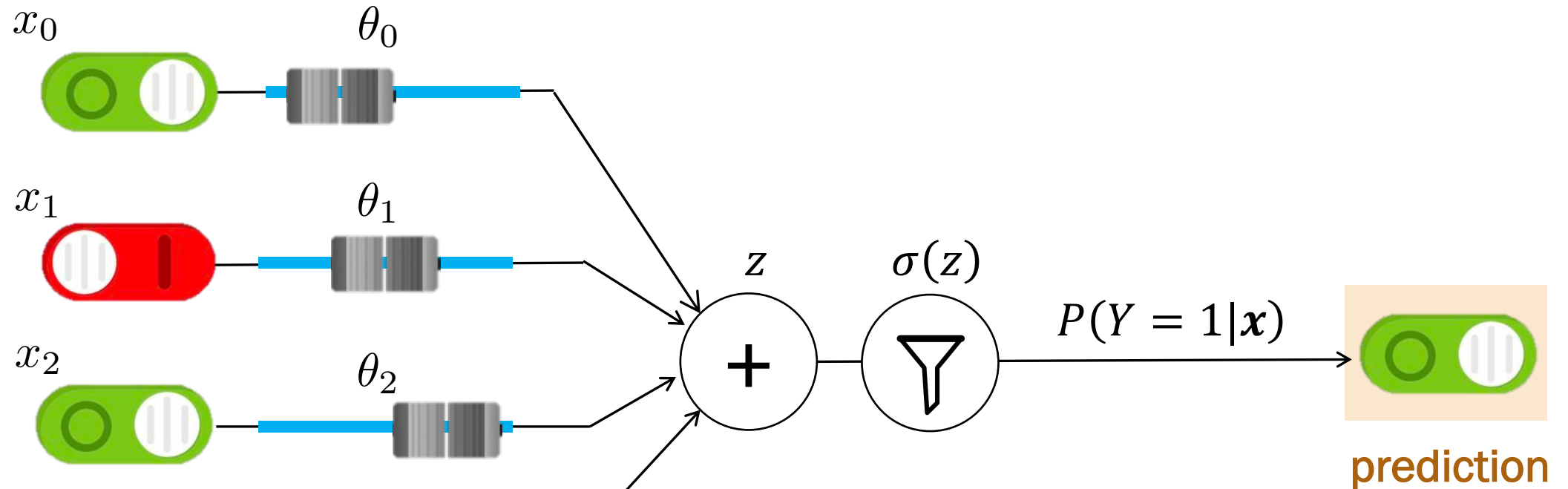
Components of Logistic Regression



squashing function
b/t 0 and 1

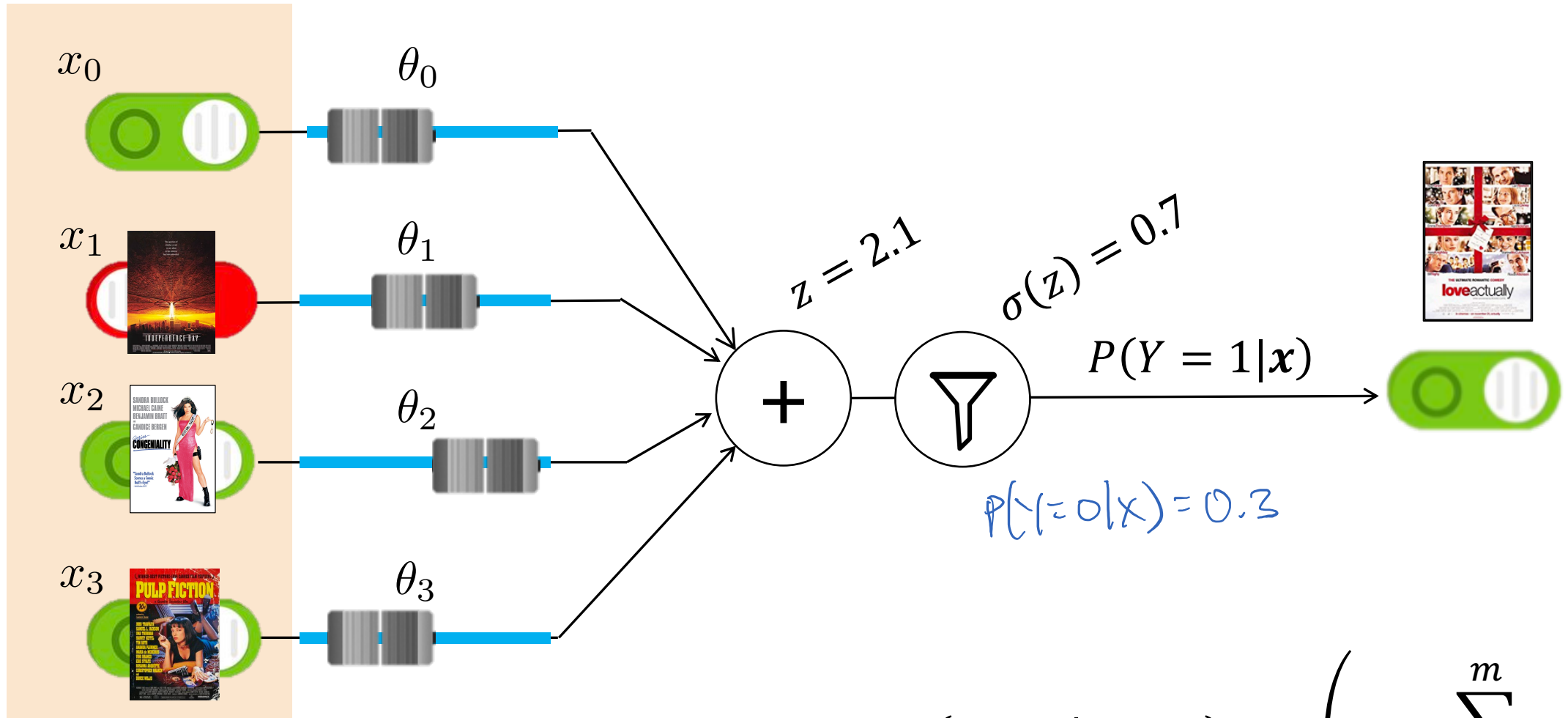
$$P(Y = 1 | \mathbf{X} = \mathbf{x}) = \sigma \left(\theta_0 + \sum_{j=1}^m \theta_j x_j \right)$$

Components of Logistic Regression



$$P(Y = 1 | \mathbf{X} = \mathbf{x}) = \sigma \left(\theta_0 + \sum_{j=1}^m \theta_j x_j \right)$$

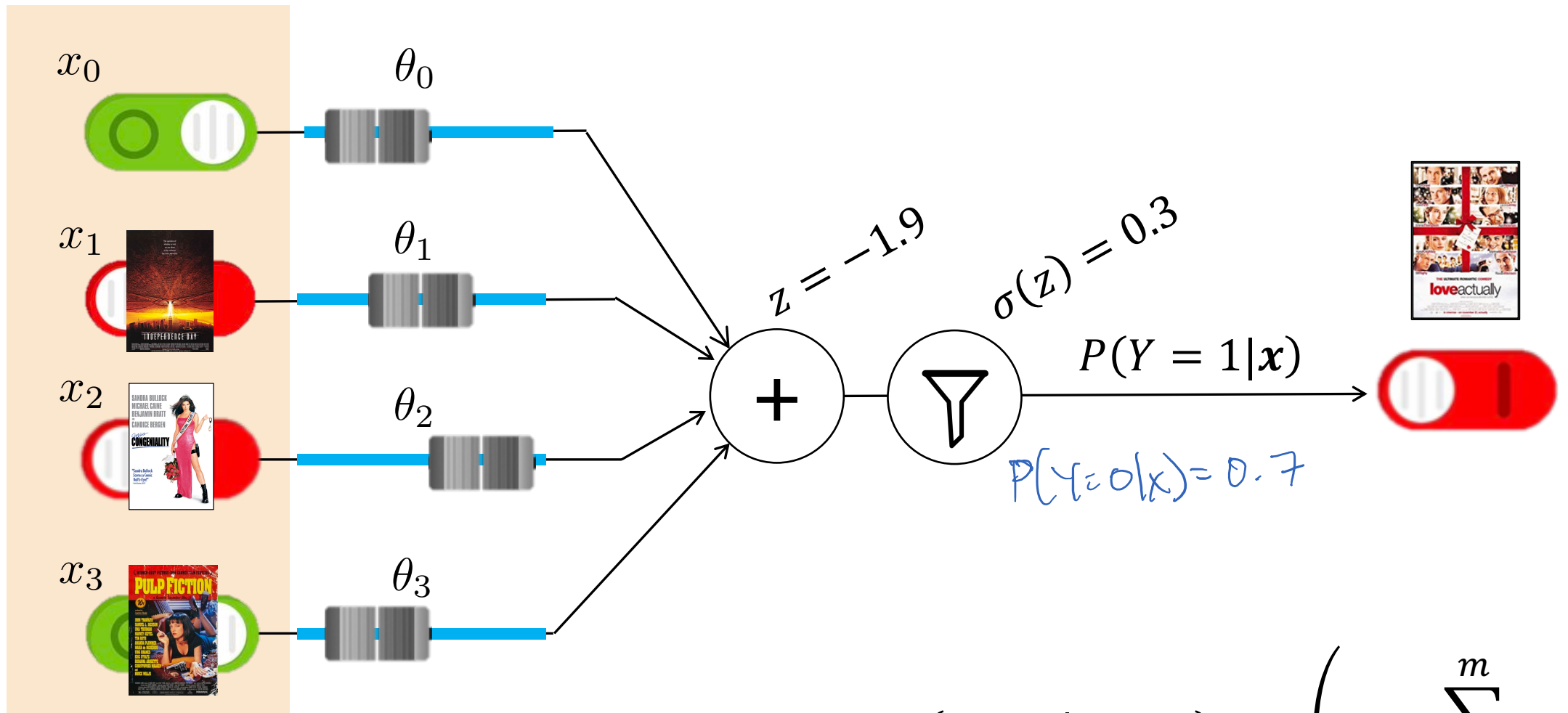
Different predictions for different inputs



\mathbf{X} , input features
[0,1,1]

$$P(Y = 1 | \mathbf{X} = \mathbf{x}) = \sigma \left(\theta_0 + \sum_{j=1}^m \theta_j x_j \right)$$

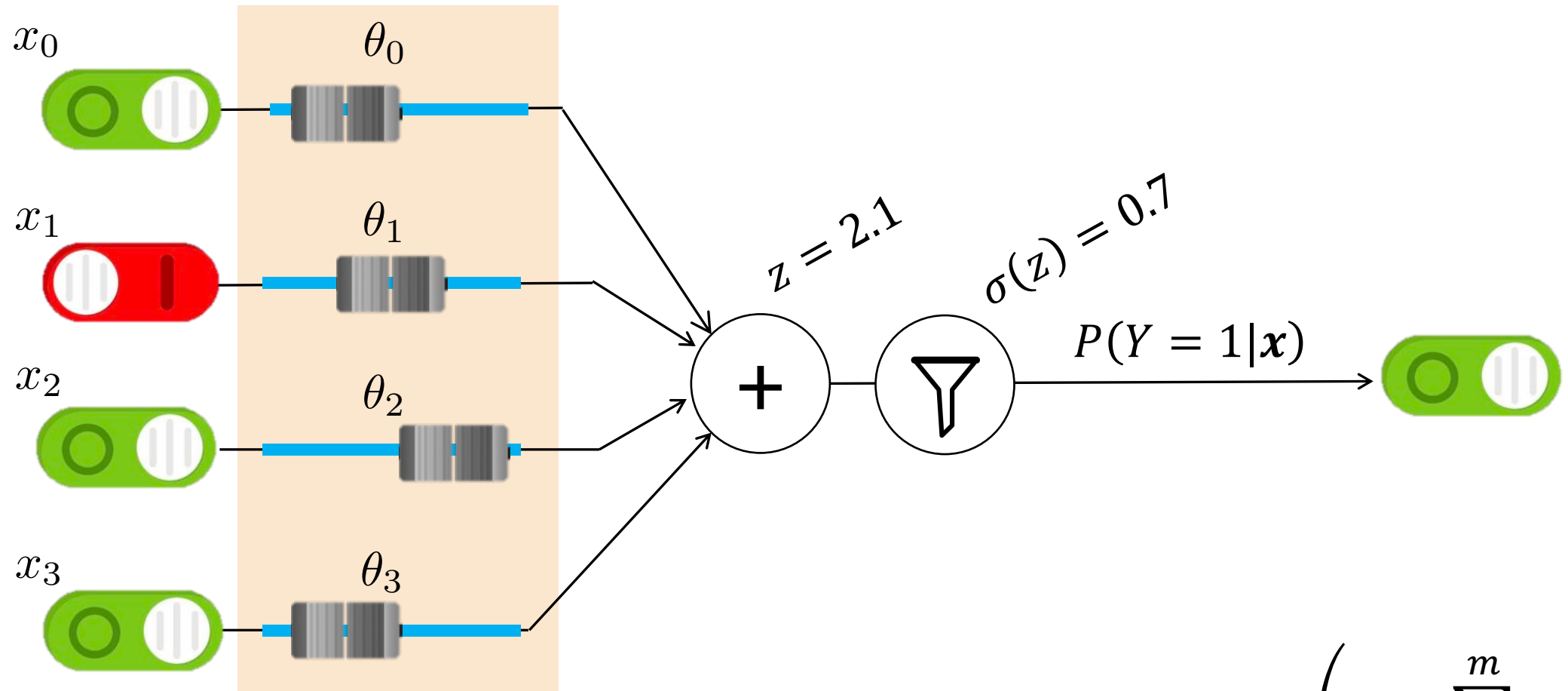
Different predictions for different inputs



\mathbf{X} , input features
[0,0,1]

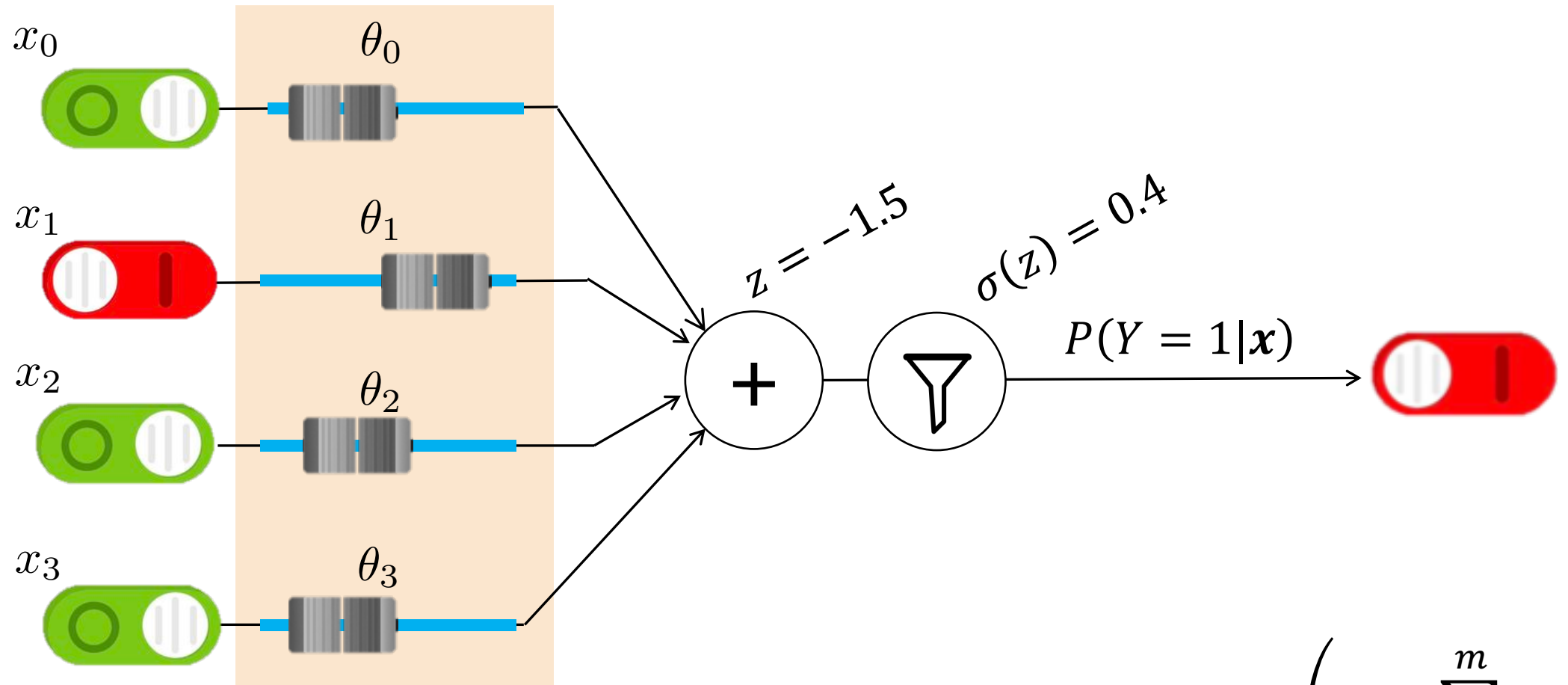
$$P(Y = 1 | \mathbf{X} = \mathbf{x}) = \sigma \left(\theta_0 + \sum_{j=1}^m \theta_j x_j \right)$$

Parameters affect prediction



$$P(Y = 1 | \mathbf{X} = \mathbf{x}) = \sigma \left(\theta_0 + \sum_{j=1}^m \theta_j x_j \right)$$

Parameters affect prediction



$$P(Y = 1 | \mathbf{X} = \mathbf{x}) = \sigma \left(\theta_0 + \sum_{j=1}^m \theta_j x_j \right)$$

For simplicity

$$P(Y = 1 | \mathbf{X} = \mathbf{x}) = \sigma \left(\theta_0 + \sum_{j=1}^m \theta_j x_j \right)$$



$$P(Y = 1 | \mathbf{X} = \mathbf{x}) = \sigma \left(\sum_{j=0}^m \theta_j x_j \right) = \boxed{\sigma(\theta^T \mathbf{x})} \text{ where } x_0 = 1$$

$\theta = (\theta_0, \theta_1, \dots, \theta_m)$
 $\mathbf{x} = (1, x_1, x_2, \dots, x_m)$

Logistic regression classifier

$$\hat{Y} = \arg \max_{y=\{0,1\}} P(Y|\mathbf{X})$$

$$P(Y = 1|\mathbf{X} = \mathbf{x}) = \sigma\left(\sum_{j=0}^m \theta_j x_j\right) = \sigma(\theta^T \mathbf{x})$$

Training

Estimate parameters
from training data

$$\theta = (\theta_0, \theta_1, \theta_2, \dots, \theta_m)$$

Testing

Given an observation $\mathbf{X} = (X_1, X_2, \dots, X_m)$, predict

$$\hat{Y} = \arg \max_{y=\{0,1\}} P(Y|\mathbf{X})$$

Training: The big picture

Logistic regression classifier

$$\hat{Y} = \arg \max_{y=\{0,1\}} P(Y|X)$$

$$P(Y = 1|X = \mathbf{x}) = \sigma\left(\sum_{j=0}^m \theta_j x_j\right) = \sigma(\theta^T \mathbf{x})$$

Training

Estimate parameters
from training data

$$(x^{(i)}, y^{(i)}) \quad i=1, \dots, n$$

$$\theta = (\theta_0, \theta_1, \theta_2, \dots, \theta_m)$$

Choose θ that optimizes some objective:

1. Determine objective function
2. Find gradient with respect to θ
3. Solve analytically by setting to 0, or computationally with gradient ascent

We are modeling $P(Y|X)$ directly, so we maximize the **conditional likelihood** of training data.

Estimating θ

1. Determine objective function

$$\theta_{MLE} = \arg \max_{\theta} \prod_{i=1}^n f(y^{(i)} | \mathbf{x}^{(i)}, \theta)$$

2. Gradient w.r.t. θ_j , for $j = 0, 1, \dots, m$

3. Solve

- No analytical derivation of θ_{MLE} ...
- ...but can still compute θ_{MLE} with gradient ascent!

```
initialize x
repeat many times:
  compute gradient
  x +=  $\eta$  * gradient
```

1. Determine objective function

$$\theta_{MLE} = \arg \max_{\theta} \prod_{i=1}^n f(y^{(i)} | \mathbf{x}^{(i)}, \theta) = \arg \max_{\theta} LL(\theta)$$

$$\begin{aligned} P(Y = 1 | \mathbf{X} = \mathbf{x}) &= \sigma(\sum_{j=0}^m \theta_j x_j) \\ &= \sigma(\theta^T \mathbf{x}) \end{aligned}$$

First: Interpret conditional likelihood with Logistic Regression

Second: Write a differentiable expression for log conditional likelihood

1. Determine objective function (interpret)

$$\theta_{MLE} = \arg \max_{\theta} \prod_{i=1}^n f(y^{(i)} | \mathbf{x}^{(i)}, \theta) = \arg \max_{\theta} LL(\theta)$$

$$P(Y = 1 | \mathbf{X} = \mathbf{x}) = \sigma\left(\sum_{j=0}^m \theta_j x_j\right) = \sigma(\theta^T \mathbf{x})$$

Suppose you have $n = 2$ training datapoints:

$$(\mathbf{x}^{(1)}, \overset{y^{(1)}}{1}), (\mathbf{x}^{(2)}, \overset{y^{(2)}}{0})$$

Consider the following expressions for a given θ :

A. $\sigma(\theta^T \mathbf{x}^{(1)}) \sigma(\theta^T \mathbf{x}^{(2)})$

C. $\sigma(\theta^T \mathbf{x}^{(1)}) (1 - \sigma(\theta^T \mathbf{x}^{(2)}))$

B. $(1 - \sigma(\theta^T \mathbf{x}^{(1)})) \sigma(\theta^T \mathbf{x}^{(2)})$

D. $(1 - \sigma(\theta^T \mathbf{x}^{(1)})) (1 - \sigma(\theta^T \mathbf{x}^{(2)}))$

1. Interpret the above expressions as probabilities.
2. If we let $\theta = \theta_{MLE}$, which probability should be highest?



1. Determine objective function (interpret)

$$\theta_{MLE} = \arg \max_{\theta} \prod_{i=1}^n f(y^{(i)} | \mathbf{x}^{(i)}, \theta) = \arg \max_{\theta} LL(\theta)$$

$$P(Y = 1 | \mathbf{X} = \mathbf{x}) = \sigma(\sum_{j=0}^m \theta_j x_j) = \sigma(\theta^T \mathbf{x})$$

Suppose you have $n = 2$ training datapoints:

$$(\mathbf{x}^{(1)}, 1), (\mathbf{x}^{(2)}, 0)$$

Consider the following expressions for a given θ :

A. $\sigma(\theta^T \mathbf{x}^{(1)}) \sigma(\theta^T \mathbf{x}^{(2)})$
 $P(Y=1 | X=x^{(1)}) P(Y=1 | X=x^{(2)})$

B. $(1 - \sigma(\theta^T \mathbf{x}^{(1)})) \sigma(\theta^T \mathbf{x}^{(2)})$
 $P(Y=0 | X=x^{(1)}) P(Y=1 | X=x^{(2)})$

C. $\sigma(\theta^T \mathbf{x}^{(1)}) (1 - \sigma(\theta^T \mathbf{x}^{(2)}))$
 $P(Y=1 | X=x^{(1)}) P(Y=0 | X=x^{(2)})$

D. $(1 - \sigma(\theta^T \mathbf{x}^{(1)})) (1 - \sigma(\theta^T \mathbf{x}^{(2)}))$
 $P(Y=0 | X=x^{(1)}) P(Y=0 | X=x^{(2)})$

1. Interpret the above expressions as probabilities.
2. If we let $\theta = \theta_{MLE}$, which probability should be highest?

1. Determine objective function (write)

$$\theta_{MLE} = \arg \max_{\theta} \prod_{i=1}^n f(y^{(i)} | \mathbf{x}^{(i)}, \theta) = \arg \max_{\theta} LL(\theta)$$

$$\begin{aligned} P(Y = 1 | \mathbf{X} = \mathbf{x}) &= \sigma(\sum_{j=0}^m \theta_j x_j) \\ &= \sigma(\theta^T \mathbf{x}) \end{aligned}$$

1. What is a differentiable expression for $P(Y = y | \mathbf{X} = \mathbf{x})$?

$$P(Y = y | \mathbf{X} = \mathbf{x}) = \begin{cases} \sigma(\theta^T \mathbf{x}) & \text{if } y = 1 \\ 1 - \sigma(\theta^T \mathbf{x}) & \text{if } y = 0 \end{cases}$$

2. What is a differentiable expression for $LL(\theta)$, log conditional likelihood?

$$LL(\theta) = \log \prod_{i=1}^n f(y^{(i)} | \mathbf{x}^{(i)}, \theta)$$



1. Determine objective function (write)

$$\theta_{MLE} = \arg \max_{\theta} \prod_{i=1}^n f(y^{(i)} | \mathbf{x}^{(i)}, \theta) = \arg \max_{\theta} LL(\theta)$$

$$P(Y = 1 | \mathbf{X} = \mathbf{x}) = \sigma(\sum_{j=0}^m \theta_j x_j) = \sigma(\theta^T \mathbf{x})$$

1. What is a differentiable expression for $P(Y = y | \mathbf{X} = \mathbf{x})$?

$$P(Y = y | \mathbf{X} = \mathbf{x}) = \begin{cases} \sigma(\theta^T \mathbf{x}) & \text{if } y = 1 \\ 1 - \sigma(\theta^T \mathbf{x}) & \text{if } y = 0 \end{cases}$$

Recall $B \sim \text{Ber}(p)$
Bernoulli MLE!
$$P(B=b) = \begin{cases} p & b=1 \\ 1-p & b=0 \end{cases}$$

$$p^b (1-p)^{1-b}$$

$$\sigma(\theta^T \mathbf{x})^y (1 - \sigma(\theta^T \mathbf{x}))^{1-y}$$

2. What is a differentiable expression for $LL(\theta)$, log conditional likelihood?

$$LL(\theta) = \log \prod_{i=1}^n f(y^{(i)} | \mathbf{x}^{(i)}, \theta)$$

$$\sum_{i=1}^n \log \left[\sigma(\theta^T \mathbf{x}^{(i)})^{y^{(i)}} (1 - \sigma(\theta^T \mathbf{x}^{(i)}))^{1-y^{(i)}} \right] = \sum_{i=1}^n \log \sigma(\theta^T \mathbf{x}^{(i)})^{y^{(i)}} + \log (1 - \sigma(\theta^T \mathbf{x}^{(i)}))^{1-y^{(i)}} = \sum_{i=1}^n y^{(i)} \log \sigma(\theta^T \mathbf{x}^{(i)}) + (1-y^{(i)}) \log (1 - \sigma(\theta^T \mathbf{x}^{(i)}))$$

1. Determine objective function (write)

$$\theta_{MLE} = \arg \max_{\theta} \prod_{i=1}^n f(y^{(i)} | \mathbf{x}^{(i)}, \theta) = \arg \max_{\theta} LL(\theta)$$

$$\begin{aligned} P(Y = 1 | \mathbf{X} = \mathbf{x}) &= \sigma(\sum_{j=0}^m \theta_j x_j) \\ &= \sigma(\theta^T \mathbf{x}) \end{aligned}$$

1. What is a differentiable expression for $P(Y = y | \mathbf{X} = \mathbf{x})$?

$$P(Y = y | \mathbf{X} = \mathbf{x}) = (\sigma(\theta^T \mathbf{x}))^y (1 - \sigma(\theta^T \mathbf{x}))^{1-y}$$

2. What is a differentiable expression for $LL(\theta)$, log conditional likelihood?

$$LL(\theta) = \sum_{i=1}^n y^{(i)} \log \sigma(\theta^T \mathbf{x}^{(i)}) + (1 - y^{(i)}) \log (1 - \sigma(\theta^T \mathbf{x}^{(i)}))$$

2. Find gradient with respect to θ

Optimization problem:

$$\theta_{MLE} = \arg \max_{\theta} \prod_{i=1}^n f(y^{(i)} | \mathbf{x}^{(i)}, \theta) = \arg \max_{\theta} LL(\theta)$$

$$LL(\theta) = \sum_{i=1}^n y^{(i)} \log \sigma(\theta^T \mathbf{x}^{(i)}) + (1 - y^{(i)}) \log (1 - \sigma(\theta^T \mathbf{x}^{(i)}))$$

Gradient w.r.t. θ_j , for $j = 0, 1, \dots, m$:

$$\frac{\partial LL(\theta)}{\partial \theta_j} = \sum_{i=1}^n [y^{(i)} - \sigma(\theta^T \mathbf{x}^{(i)})] x_j^{(i)}$$

$\theta = (\theta_0, \theta_1, \theta_2, \dots, \theta_m)$
 $\theta^T \mathbf{x}^{(i)} = \sum_{j=0}^m \theta_j x_j^{(i)}$
(derived later)

How do we interpret the gradient contribution of the i -th training datapoint?



2. Find gradient with respect to θ

Optimization problem:

$$\theta_{MLE} = \arg \max_{\theta} \prod_{i=1}^n f(y^{(i)} | \mathbf{x}^{(i)}, \theta) = \arg \max_{\theta} LL(\theta)$$

$$LL(\theta) = \sum_{i=1}^n y^{(i)} \log \sigma(\theta^T \mathbf{x}^{(i)}) + (1 - y^{(i)}) \log (1 - \sigma(\theta^T \mathbf{x}^{(i)}))$$

Gradient w.r.t. θ_j , for $j = 0, 1, \dots, m$:

$$\frac{\partial LL(\theta)}{\partial \theta_j} = \sum_{i=1}^n [y^{(i)} - \sigma(\theta^T \mathbf{x}^{(i)})] x_j^{(i)} \quad (\text{derived later})$$

↑
scale by j-th feature

$$\theta^T \mathbf{x}^{(i)} = \sum_{j=0}^m \theta_j x_j^{(i)}$$

2. Find gradient with respect to θ

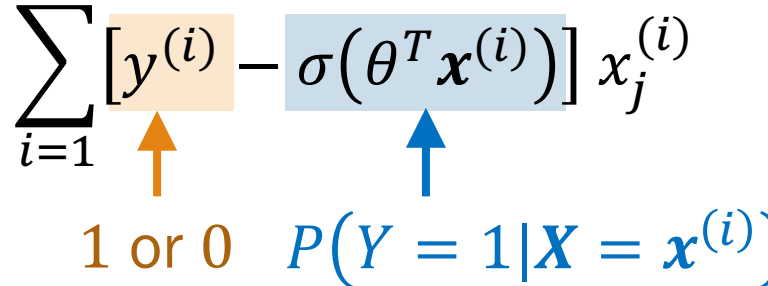
Optimization problem:

$$\theta_{MLE} = \arg \max_{\theta} \prod_{i=1}^n f(y^{(i)} | \mathbf{x}^{(i)}, \theta) = \arg \max_{\theta} LL(\theta)$$

$$LL(\theta) = \sum_{i=1}^n y^{(i)} \log \sigma(\theta^T \mathbf{x}^{(i)}) + (1 - y^{(i)}) \log (1 - \sigma(\theta^T \mathbf{x}^{(i)}))$$

Gradient w.r.t. θ_j , for $j = 0, 1, \dots, m$:

$$\frac{\partial LL(\theta)}{\partial \theta_j} = \sum_{i=1}^n [y^{(i)} - \sigma(\theta^T \mathbf{x}^{(i)})] x_j^{(i)} \quad (\text{derived later})$$



2. Find gradient with respect to θ

Optimization
problem:

$$\theta_{MLE} = \arg \max_{\theta} \prod_{i=1}^n f(y^{(i)} | \mathbf{x}^{(i)}, \theta) = \arg \max_{\theta} LL(\theta)$$

$$LL(\theta) = \sum_{i=1}^n y^{(i)} \log \sigma(\theta^T \mathbf{x}^{(i)}) + (1 - y^{(i)}) \log (1 - \sigma(\theta^T \mathbf{x}^{(i)}))$$

Gradient w.r.t. θ_j , for $j = 0, 1, \dots, m$:

$$\frac{\partial LL(\theta)}{\partial \theta_j} = \sum_{i=1}^n \underbrace{[y^{(i)} - \sigma(\theta^T \mathbf{x}^{(i)})]}_{\text{(derived later)}} x_j^{(i)}$$

Suppose $y^{(i)} = 1$ (the true class label for i -th datapoint):

- If $\sigma(\theta^T \mathbf{x}^{(i)}) \geq 0.5$, correct
- If $\sigma(\theta^T \mathbf{x}^{(i)}) < 0.5$, incorrect → change θ_j more

3. Solve

1. Optimization problem:

$$\theta_{MLE} = \arg \max_{\theta} \prod_{i=1}^n f(y^{(i)} | \mathbf{x}^{(i)}, \theta) = \arg \max_{\theta} LL(\theta)$$

$$LL(\theta) = \sum_{i=1}^n y^{(i)} \log \sigma(\theta^T \mathbf{x}^{(i)}) + (1 - y^{(i)}) \log (1 - \sigma(\theta^T \mathbf{x}^{(i)}))$$

2. Gradient w.r.t. θ_j , for $j = 0, 1, \dots, m$:

$$\frac{\partial LL(\theta)}{\partial \theta_j} = \sum_{i=1}^n [y^{(i)} - \sigma(\theta^T \mathbf{x}^{(i)})] x_j^{(i)}$$

3. Solve

Stay tuned!

26: Logistic Regression (live)

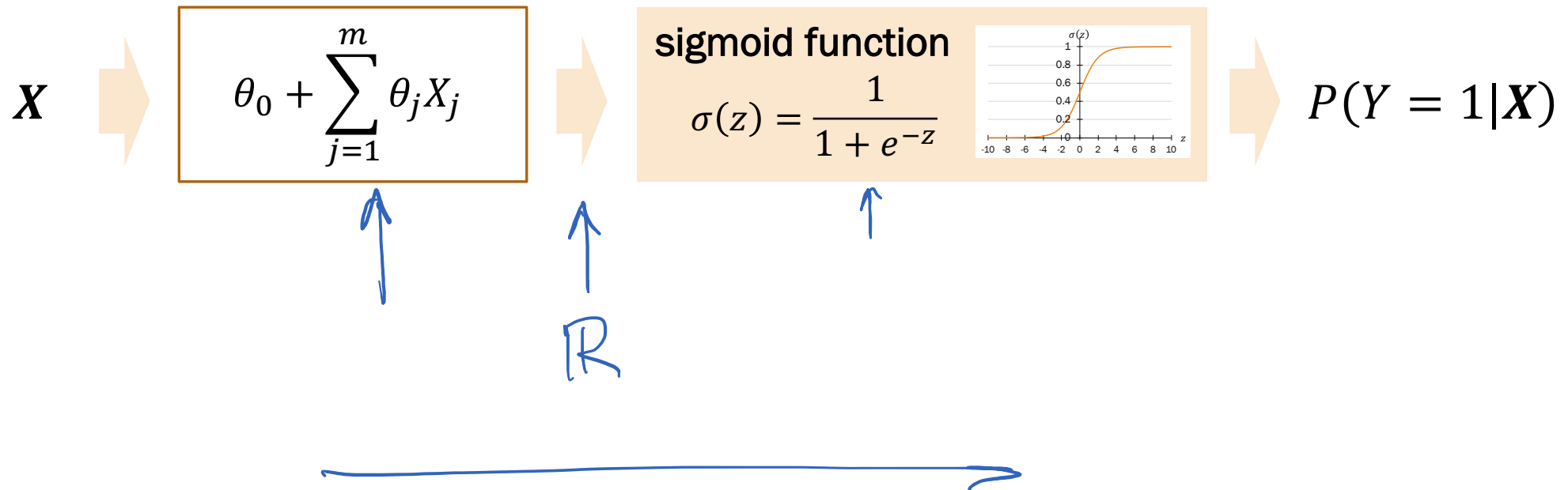
Lisa Yan and Jerry Cain
November 11, 2020

Logistic Regression Model

$$\hat{Y} = \arg \max_{y=\{0,1\}} P(Y|\mathbf{X})$$

$$P(Y = 1|\mathbf{X} = \mathbf{x}) = \sigma\left(\sum_{j=0}^m \theta_j x_j\right) = \sigma(\theta^T \mathbf{x})$$

where $x_0 = 1$



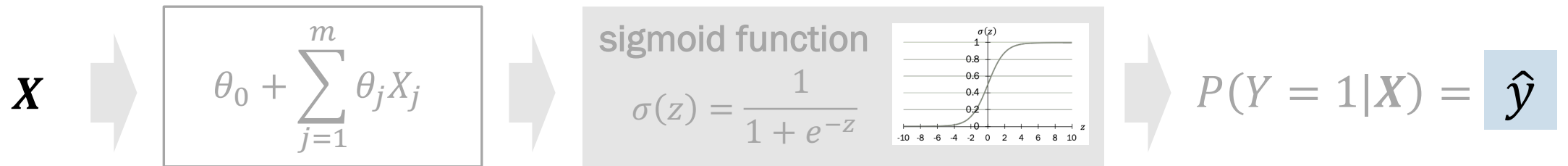
Introducing notation \hat{y}

$$\hat{Y} = \arg \max_{y \in \{0,1\}} P(Y|X)$$

$$P(Y = 1|X = \mathbf{x}) = \sigma\left(\sum_{j=0}^m \theta_j x_j\right) = \sigma(\theta^T \mathbf{x})$$

\hat{Y} is prediction of Y . $\hat{Y} \in \{0,1\}$

where $x_0 = 1$



$$\hat{y} = P(Y = 1|X = \mathbf{x}) = \sigma(\theta^T \mathbf{x})$$

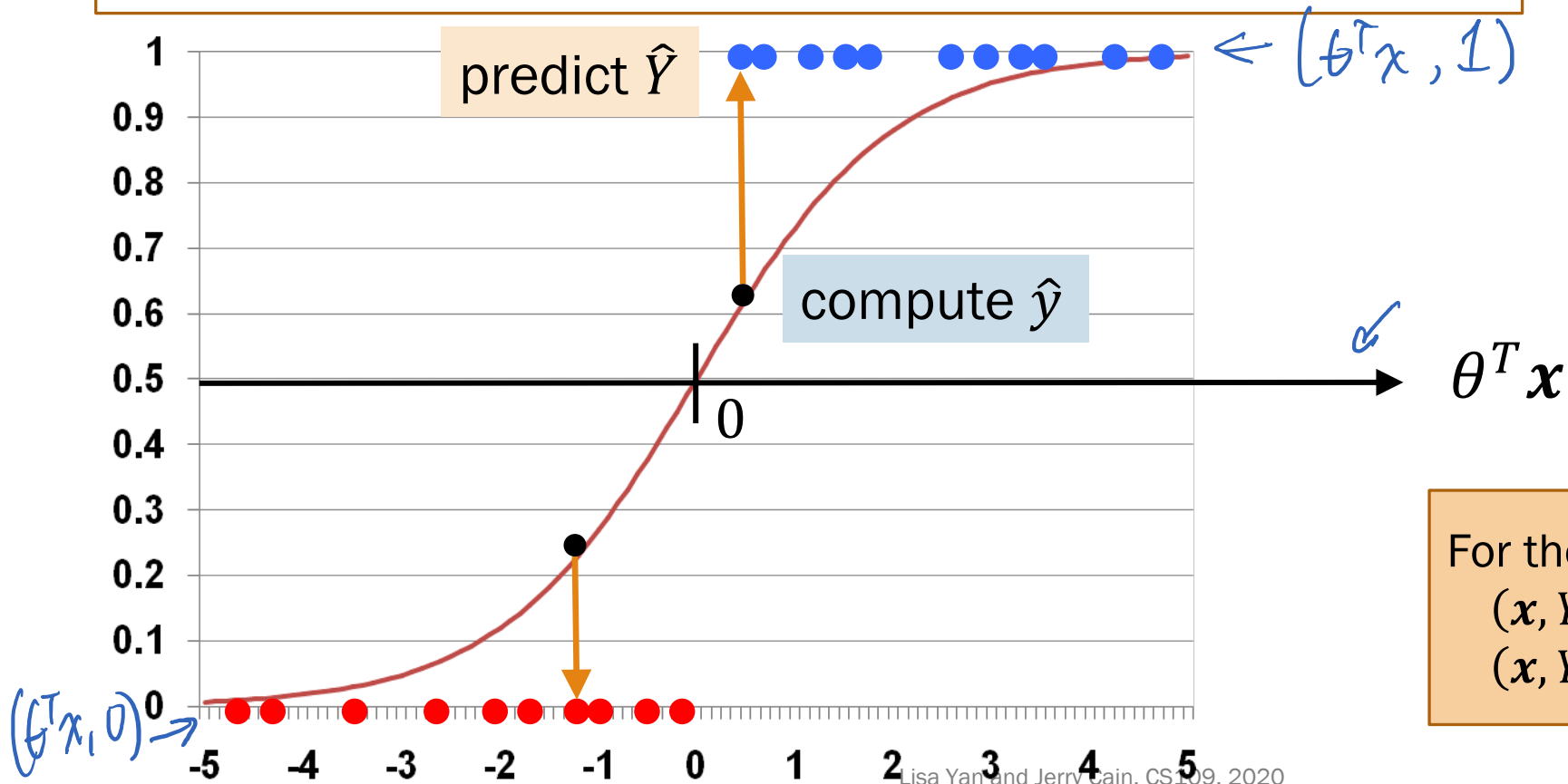
$$P(Y = y|X = \mathbf{x}) = \begin{cases} \hat{y} & \text{if } y = 1 \\ 1 - \hat{y} & \text{if } y = 0 \end{cases}$$

Small \hat{y} is conditional probability of $Y = 1$ given $X = \mathbf{x}$. $\hat{y} \in [0,1]$

Another view of Logistic Regression

$$\hat{Y} = \arg \max_{y=\{0,1\}} P(Y|X)$$

$$\hat{y} = P(Y = 1|X = x) = \sigma\left(\sum_{j=0}^m \theta_j x_j\right) = \sigma(\theta^T x)$$



For the “correct” parameters θ :
 $(x, Y = 1)$ should have $\theta^T x > 0$
 $(x, Y = 0)$ should have $\theta^T x \leq 0$

Today's goals: Logistic Regression

- ✓ At a high level
 - Understand the model
 - Training: Use gradient ascent

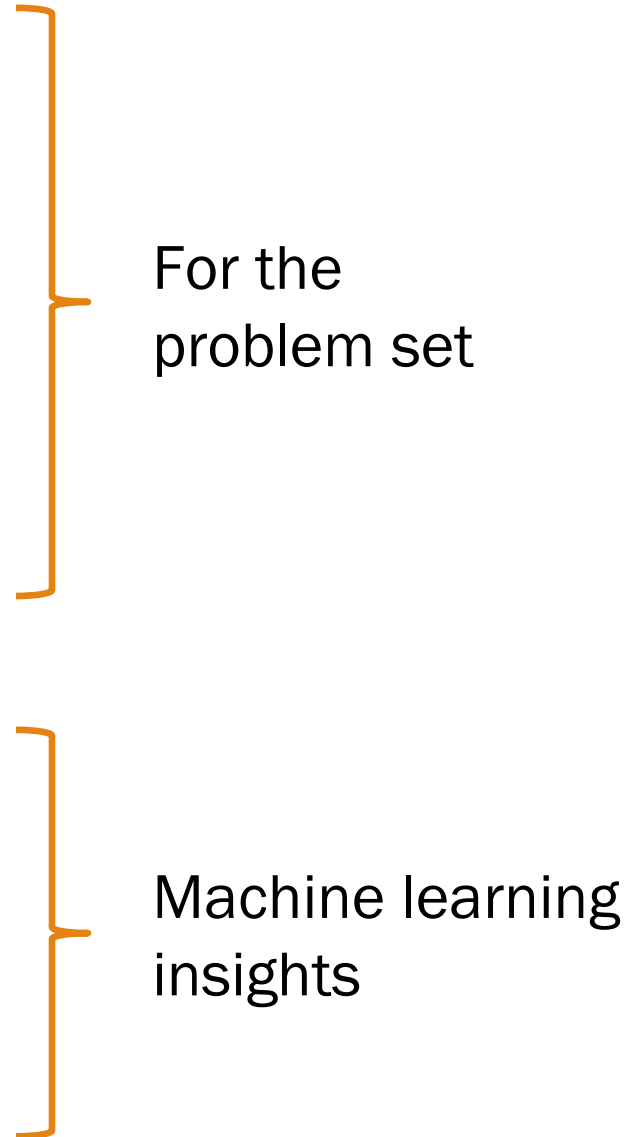
Details

- Gradient ascent pseudocode
- Testing

Philosophy

- Logistic Regression vs Naïve Bayes
- Linearly separable functions

Derivation of gradient (Calculus)



Training: The details

Training: Learning parameters

Review

$$\left(\mathbf{x}^{(i)}, y^{(i)} \right) \\ i=1, \dots, n$$

Training

Learn parameters $\theta = (\theta_0, \theta_1, \dots, \theta_m)$

that maximize log conditional likelihood of training data

Some reminders:

- Log conditional likelihood:

$$LL(\theta) = \sum_{i=1}^n y^{(i)} \log \sigma(\theta^T \mathbf{x}^{(i)}) + (1 - y^{(i)}) \log (1 - \sigma(\theta^T \mathbf{x}^{(i)}))$$

- Gradient with respect to θ :

$$\frac{\partial LL(\theta)}{\partial \theta_j} = \sum_{i=1}^n [y^{(i)} - \sigma(\theta^T \mathbf{x}^{(i)})] x_j^{(i)} \quad \text{for } j = 0, 1, \dots, m \quad \text{(derived at end of lecture)}$$

- No analytical solution; optimize with **gradient ascent**

Training: Gradient ascent step

$$\frac{\partial LL(\theta)}{\partial \theta_j} = \sum_{i=1}^n [y^{(i)} - \sigma(\theta^T \mathbf{x}^{(i)})] x_j^{(i)} \quad \text{for } j = 0, 1, \dots, m$$

repeat many times:

for all thetas: $j = 0, 1, \dots, m$

$$\theta_j^{\text{new}} = \theta_j^{\text{old}} + \eta \cdot \frac{\partial LL(\theta^{\text{old}})}{\partial \theta_j^{\text{old}}}$$

$$= \theta_j^{\text{old}} + \eta \cdot \sum_{i=1}^n [y^{(i)} - \sigma(\theta^{\text{old}T} \mathbf{x}^{(i)})] x_j^{(i)}$$

learning rate (under η)
gradient w.r.t. θ_j (under the sum)

What does this look like in code?

Think

Slide 50 has code to think over by yourself.

Post any clarifications here or in chat!

<https://us.edstem.org/courses/2678/discussion/171556>

Think by yourself: 2 min




(by yourself)

Training: Gradient Ascent

for $j = 0, 1, \dots, m$:

Gradient Ascent Step $\theta_j^{\text{new}} = \theta_j^{\text{old}} + \eta \cdot \sum_{i=1}^n [y^{(i)} - \sigma(\theta^{\text{old}T} \mathbf{x}^{(i)})] x_j^{(i)}$



initialize $\theta_j = 0$ for $0 \leq j \leq m$
repeat many times:

```
gradient[j] = 0 for  $0 \leq j \leq m$ 
```

```
// TODO: your code here
```

```
// compute all gradient[j]'s
```

```
// based on n training examples
```

```
 $\theta_j$  +=  $\eta$  * gradient[j] for all  $0 \leq j \leq m$ 
```



Training: Gradient Ascent

inner loop

for $j = 0, 1, \dots, m$:

Gradient
Ascent Step

$$\theta_j^{\text{new}} = \theta_j^{\text{old}} + \eta \cdot \sum_{i=1}^n [y^{(i)} - \sigma(\theta^{\text{old}T} \mathbf{x}^{(i)})] x_j^{(i)}$$

outer loop

compute

```
initialize  $\theta_j = 0$  for  $0 \leq j \leq m$   
repeat many times:
```

```
  gradient[j] = 0 for  $0 \leq j \leq m$ 
```

```
  for each training example  $(x, y)$ :
```

```
    for each  $0 \leq j \leq m$ :
```

```
      // update gradient[j] for  
      // current  $(x, y)$  example
```

```
   $\theta_j += \eta * \text{gradient}[j]$  for all  $0 \leq j \leq m$ 
```

Training: Gradient Ascent

inner loop

for $j = 0, 1, \dots, m$:

Gradient
Ascent Step

$$\theta_j^{\text{new}} = \theta_j^{\text{old}} + \eta \cdot \sum_{i=1}^n \left[y^{(i)} - \sigma(\theta^{\text{old}T} \mathbf{x}^{(i)}) \right] x_j^{(i)}$$

outer loop

compute

initialize $\theta_j = 0$ for $0 \leq j \leq m$
repeat many times:

gradient[j] = 0 for $0 \leq j \leq m$

for each training example (x, y) :

for each $0 \leq j \leq m$:

$$\text{gradient}[j] += \left[y - \frac{1}{1 + e^{-\theta^T x}} \right] x_j$$

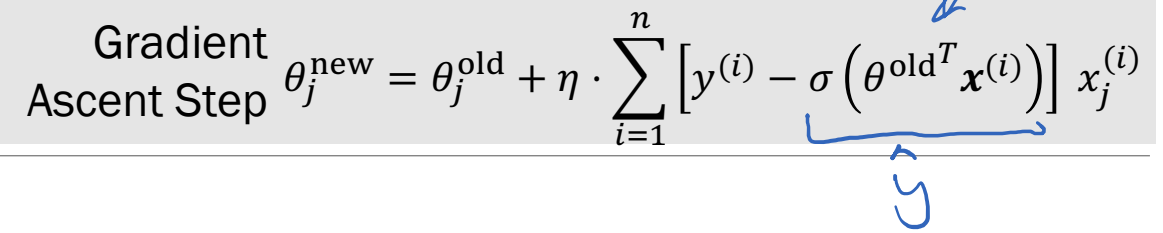
$\theta_j += \eta * \text{gradient}[j]$ for all $0 \leq j \leq m$

$\hat{y} = \sigma(\theta^T x)$

Some important
details...

Training: Gradient Ascent

Gradient Ascent Step $\theta_j^{\text{new}} = \theta_j^{\text{old}} + \eta \cdot \sum_{i=1}^n [y^{(i)} - \sigma(\theta^{\text{old}T} \mathbf{x}^{(i)})] x_j^{(i)}$



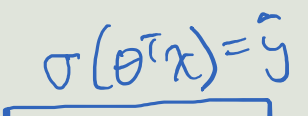
initialize $\theta_j = 0$ for $0 \leq j \leq m$
repeat many times:

gradient[j] = 0 for $0 \leq j \leq m$

for each training example (\mathbf{x}, y) :

for each $0 \leq j \leq m$:

gradient[j] += $\left[y - \frac{1}{1 + e^{-\theta^T \mathbf{x}}} \right] x_j$



$\theta_j += \eta * \text{gradient}[j]$ for all $0 \leq j \leq m$



- Finish computing gradient with θ^{old} prior to any θ update

Training: Gradient Ascent

$$\text{Gradient Ascent Step } \theta_j^{\text{new}} = \theta_j^{\text{old}} + \eta \cdot \sum_{i=1}^n [y^{(i)} - \sigma(\theta^{\text{old}T} \mathbf{x}^{(i)})] x_j^{(i)}$$

initialize $\theta_j = 0$ for $0 \leq j \leq m$
repeat many times:

gradient[j] = 0 for $0 \leq j \leq m$

for each training example (x, y) :

for each $0 \leq j \leq m$:

$$\text{gradient}[j] += \left[y - \frac{1}{1 + e^{-\theta^T x}} \right] x_j$$

$\theta_j += \eta * \text{gradient}[j]$ for all $0 \leq j \leq m$

\uparrow
 $\propto 1 \times 10^{-6}$

- Finish computing gradient with θ^{old} prior to any θ update
- Learning rate η is a constant you set before training

Training: Gradient Ascent

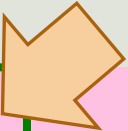
$$\text{Gradient Ascent Step } \theta_j^{\text{new}} = \theta_j^{\text{old}} + \eta \cdot \sum_{i=1}^n [y^{(i)} - \sigma(\theta^{\text{old}T} \mathbf{x}^{(i)})] x_j^{(i)}$$

initialize $\theta_j = 0$ for $0 \leq j \leq m$
repeat many times:

gradient[j] = 0 for $0 \leq j \leq m$

for each training example (\mathbf{x}, y) :

for each $0 \leq j \leq m$:

$$\text{gradient}[j] += \left[y - \frac{1}{1 + e^{-\theta^T \mathbf{x}}} \right] x_j$$


$\theta_j += \eta * \text{gradient}[j]$ for all $0 \leq j \leq m$

- Finish computing gradient with θ^{old} prior to any θ update
- Learning rate η is a constant you set before training
- x_j is j -th feature of input $\mathbf{x} = (x_1, \dots, x_m)$

Training: Gradient Ascent

$$\text{Gradient Ascent Step } \theta_j^{\text{new}} = \theta_j^{\text{old}} + \eta \cdot \sum_{i=1}^n [y^{(i)} - \sigma(\theta^{\text{old}T} \mathbf{x}^{(i)})] x_j^{(i)}$$

insert $x_0 = 1$ into all training data x $(x_1, x_2, \dots, x_m) \rightarrow (1, x_1, x_2, \dots, x_m)$ $\theta^T x = \theta_0 + \sum_{j=1}^m \theta_j x_j$

initialize $\theta_j = 0$ for $0 \leq j \leq m$
repeat many times:

gradient[j] = 0 for $0 \leq j \leq m$

for each training example (x, y) :

for each $0 \leq j \leq m$:

$$\text{gradient}[j] += \left[y - \frac{1}{1 + e^{-\theta^T x}} \right] x_j$$

$\theta_j += \eta * \text{gradient}[j]$ for all $0 \leq j \leq m$

- Finish computing gradient with θ^{old} prior to any θ update
- Learning rate η is a constant you set before training
- x_j is j -th feature of input $\mathbf{x} = (x_1, \dots, x_m)$
- Insert $x_0 = 1$ before training

Training: Gradient Ascent

$$\text{Gradient Ascent Step } \theta_j^{\text{new}} = \theta_j^{\text{old}} + \eta \cdot \sum_{i=1}^n [y^{(i)} - \sigma(\theta^{\text{old}T} \mathbf{x}^{(i)})] x_j^{(i)}$$

$$(x^{(1)}, y^{(1)}), (x^{(2)}, y^{(2)}), \dots, (x^{(n)}, y^{(n)})$$

initialize $\theta_j = 0$ for $0 \leq j \leq m$
repeat many times:

gradient[j] = 0 for $0 \leq j \leq m$

for each training example (x, y) :

for each $0 \leq j \leq m$:

$$\text{gradient}[j] += \left[y - \frac{1}{1 + e^{-\theta^T x}} \right] x_j$$

$\theta_j += \eta * \text{gradient}[j]$ for all $0 \leq j \leq m$

- Finish computing gradient with θ^{old} prior to any θ update
- Learning rate η is a constant you set before training
- x_j is j -th feature of input $\mathbf{x} = (x_1, \dots, x_m)$
- Insert $x_0 = 1$ before training



Testing

Testing: Classification with Logistic Regression

Training

Learn parameters $\theta = (\theta_0, \theta_1, \dots, \theta_m)$

via gradient ascent:

$$\theta_j^{\text{new}} = \theta_j^{\text{old}} + \eta \cdot \sum_{i=1}^n [y^{(i)} - \sigma(\theta^{\text{old}T} \mathbf{x}^{(i)})] x_j^{(i)}$$

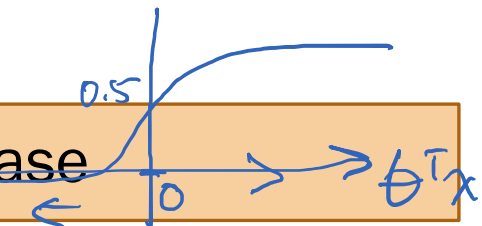
Testing

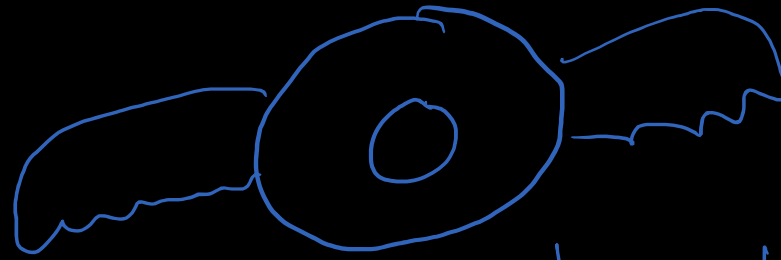
- Compute $\hat{y} = P(Y = 1 | X = \mathbf{x}) = \sigma(\theta^T \mathbf{x}) = \frac{1}{1 + e^{-\theta^T \mathbf{x}}}$
- Classify instance as: 1 if $P(Y=1|X=\mathbf{x}) > P(Y=0|X=\mathbf{x})$

$$\hat{y} = \begin{cases} 1 & \hat{y} > 0.5, \text{ equivalently } \theta^T \mathbf{x} > 0 \\ 0 & \text{otherwise} \end{cases}$$



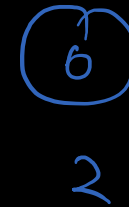
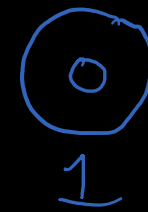
Parameters θ_j are not updated during testing phase





plane bagel

schmar
schmar
campuzn



Interlude for jokes/announcements

<https://www.bagelbakerygainesville.com/top-8-bagel-jokes-of-all-time/>

Announcements

Quiz #3

Time frame: Wednesday 11/18 2:00pm – Friday 11/20 12:59pm PT
Covers: Up to and including logistic regression
Info and practice: [Quizzes page](#)

Next week: Last section

Review session for Quiz #3

Probability Reference ([Overleaf](#))

Updated to include all of Quiz 3-relevant material (sampling defs, MLE/MAP, classifiers)

Interesting probability news

The Time Everyone “Corrected” the World’s Smartest Woman



<https://priceconomics.com/the-time-everyone-corrected-the-worlds-smartest/>

Today's goals: Logistic Regression

- ✓ At a high level
 - Understand the model
 - Training: Use gradient ascent

Details

- ✓
 - Gradient ascent pseudocode
 - Testing

Philosophy

- Logistic Regression vs Naïve Bayes
- Linearly separable functions

Derivation of gradient (Calculus)



Philosophy

Think

Slide 64 asks you to think over by yourself.

Post any clarifications here or in chat!

<https://us.edstem.org/courses/2678/discussion/171556>

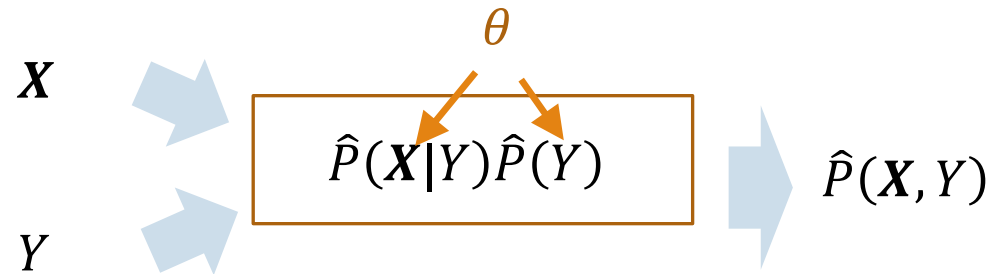
Think by yourself: 2 min



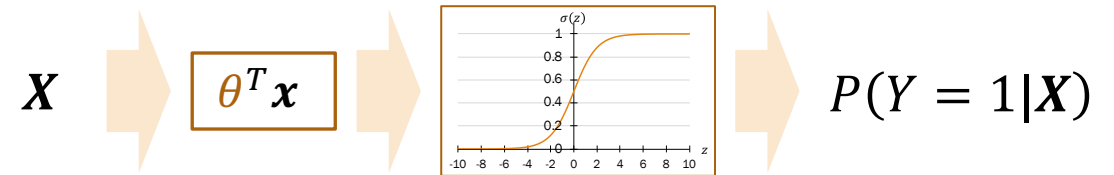
Naïve Bayes

vs

Logistic Regression



$$\hat{Y} = \arg \max_{y=\{0,1\}} P(Y | \mathbf{X}) = \arg \max_{y=\{0,1\}} P(\mathbf{X}|Y)P(Y)$$



$$\hat{Y} = \arg \max_{y=\{0,1\}} P(Y|\mathbf{X})$$

Compare/contrast:

1. What **distributions** are we modeling?
2. After learning our parameters, could we randomly **generate** a new datapoint (x, y) ?
3. Could we model a **continuous** X_j feature (e.g., $X_j \sim \text{Normal}$, or $X_j \sim \text{Unknown}$)?
4. Could we model a non-binary **discrete** X_j (e.g., $X_j \in \{1, 2, \dots, 6\}$)?



Tradeoffs:

Naïve Bayes

Logistic Regression

1. Modeling goal

$$P(\mathbf{X}, Y)$$

$$P(Y|\mathbf{X})$$

2. Generative or discriminative?

Generative: could use joint distribution to generate new points (⚠️ but you might not need this extra effort)

Discriminative: just tries to discriminate $y = 0$ vs $y = 1$ (❌ cannot generate new points b/c no $P(\mathbf{X}, Y)$)

3. Continuous input features

$$X \sim N(\mu, \sigma^2)$$

⚠️ Needs parametric form (e.g., Gaussian) or discretized buckets (for multinomial features)

✅ Yes, easily

$$\underline{\underline{\theta^T x}}$$

4. Discrete input features

0 - 100

{apple, banana, orange}

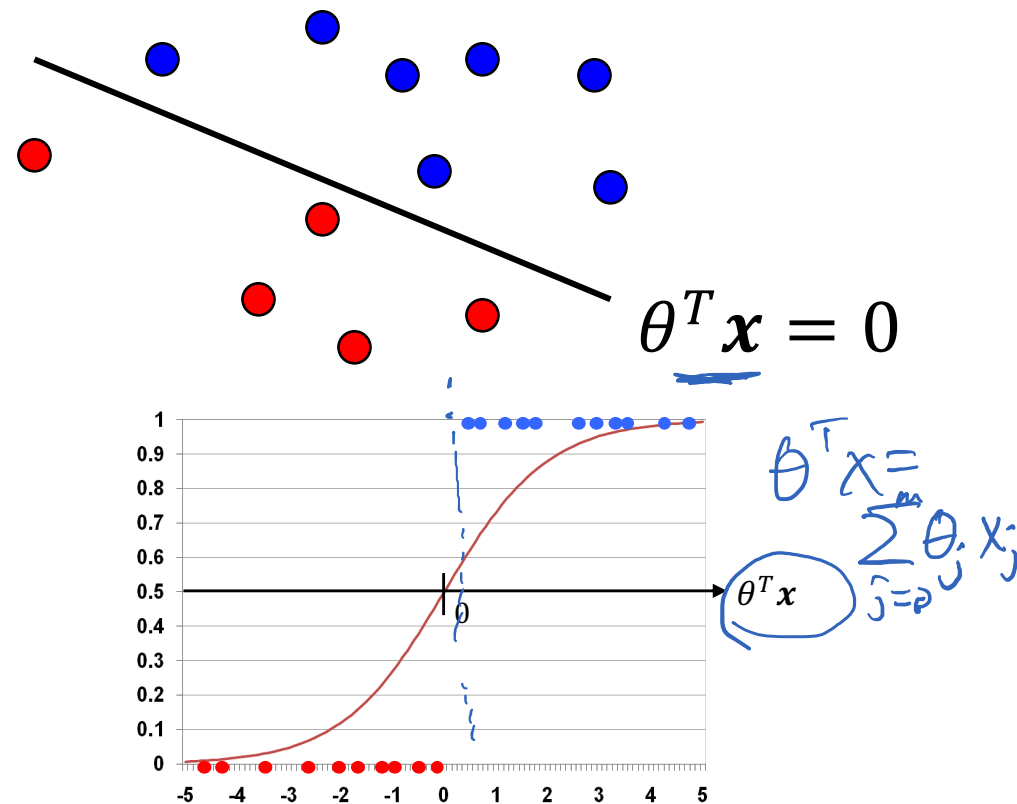
✅ Yes, multi-value discrete data = multinomial $P(X_i|Y)$

⚠️ Multi-valued discrete data hard (e.g., if $X_i \in \{A, B, C\}$, not necessarily good to encode as $\{1, 2, 3\}$)

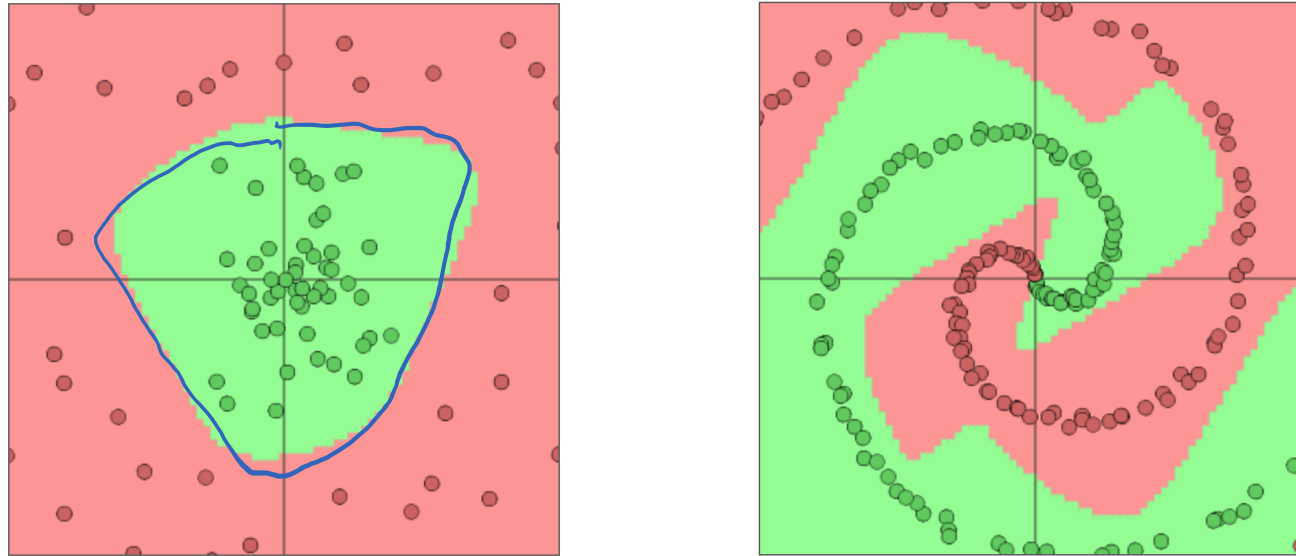
Linearly separable data

Logistic Regression is trying to fit a line that separates data instances where $y = 1$ from those where $y = 0$:

- We call such data (or functions generating the data) linearly separable.
- Naïve Bayes is linear too, because there is one parameter for each feature (and no parameters that involve multiple features).



Data is often not linearly separable



- Not possible to draw a line that successfully separates all the $y = 1$ points (green) from the $y = 0$ points (red)
- Despite this fact, Logistic Regression and Naive Bayes still often work well in practice

Gradient Derivation

Background: Calculus

Calculus refresher #1:

Derivative(sum) =
sum(derivative)

$$\frac{\partial}{\partial x} \sum_{i=1}^n f_i(x) = \sum_{i=1}^n \frac{\partial f_i(x)}{\partial x}$$

Calculus refresher #2:

Chain rule 🌟🌟🌟

$$\frac{\partial f(x)}{\partial x} = \frac{\partial f(z)}{\partial z} \frac{\partial z}{\partial x}$$

Calculus Chain Rule

$$f(x) = f(z(x))$$

aka decomposition
of composed functions

Are you ready?

The screenshot shows the Quora website interface. At the top, the Quora logo is on the left, followed by navigation links: Home, Answer, Spaces, and Notifications (with a red badge showing '1'). A search bar is on the right. Below the navigation is a horizontal menu with categories: Moments, Personal Experiences, Important Life Lessons, and a '+5' link with a pencil icon. The main question is "What is your best 'I've never been more ready in my life' moment?". Below the question are interaction options: Answer, Follow (with a '- 2' count), and Request. There are also icons for comments, downvotes, Facebook, Twitter, and a share icon. Below the question, it says "1 Answer". At the bottom of the question area, there are options for Upvote (with a '- 1' count) and Share, along with downvote, share, and more options icons.

Quora Home Answer Spaces Notifications 1 Search

Moments Personal Experiences Important Life Lessons +5

What is your best "I've never been more ready in my life" moment?

Answer Follow - 2 Request

1 Answer

Right now!!!

12 views · View Upvoters

Upvote · 1 Share

Downvote Share More

Our goal

Find: $\frac{\partial LL(\theta)}{\partial \theta_j}$ where

$$LL(\theta) = \sum_{i=1}^n y^{(i)} \log \sigma(\theta^T \mathbf{x}^{(i)}) + (1 - y^{(i)}) \log (1 - \sigma(\theta^T \mathbf{x}^{(i)}))$$

log conditional likelihood

Think

Slide 72 has code to think over by yourself.

Post any clarifications here or in chat!

<https://us.edstem.org/courses/2678/discussion/171556>

Think by yourself: 2 min



(by yourself)

Aside: Sigmoid has a beautiful derivative

Sigmoid function:

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

Derivative:

$$\frac{d}{dz} \sigma(z) = \sigma(z)[1 - \sigma(z)]$$

What is $\frac{\partial}{\partial \theta_j} \sigma(\theta^T \mathbf{x})$?

- A. $\sigma(x_j)[1 - \sigma(x_j)]x_j$
- B. $\sigma(\theta^T \mathbf{x})[1 - \sigma(\theta^T \mathbf{x})]\mathbf{x}$
- C. $\sigma(\theta^T \mathbf{x})[1 - \sigma(\theta^T \mathbf{x})]x_j$
- D. $\sigma(\theta^T \mathbf{x})x_j[1 - \sigma(\theta^T \mathbf{x})x_j]$
- E. None/other



Aside: Sigmoid has a beautiful derivative

Sigmoid function:

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

Derivative:

$$\frac{d}{dz} \sigma(z) = \sigma(z)[1 - \sigma(z)]$$

What is $\frac{\partial}{\partial \theta_j} \sigma(\theta^T \mathbf{x})$?

$$\text{Let } z = \theta^T \mathbf{x} = \sum_{k=0}^m \theta_k x_k.$$

- A. $\sigma(x_j)[1 - \sigma(x_j)]x_j$
- B. $\sigma(\theta^T \mathbf{x})[1 - \sigma(\theta^T \mathbf{x})]x$
- C.** $\sigma(\theta^T \mathbf{x})[1 - \sigma(\theta^T \mathbf{x})]x_j$
- D. $\sigma(\theta^T \mathbf{x})x_j[1 - \sigma(\theta^T \mathbf{x})x_j]$
- E. None/other

$$\frac{\partial}{\partial \theta_j} \sigma(\theta^T \mathbf{x}) = \frac{\partial}{\partial z} \sigma(z) \cdot \frac{\partial z}{\partial \theta_j} \quad (\text{Chain Rule})$$

$$= \sigma(\theta^T \mathbf{x})[1 - \sigma(\theta^T \mathbf{x})]x_j$$

Our goal: Re-introducing notation \hat{y}

Find: $\frac{\partial LL(\theta)}{\partial \theta_j}$ where

$$LL(\theta) = \sum_{i=1}^n y^{(i)} \log \sigma(\theta^T \mathbf{x}^{(i)}) + (1 - y^{(i)}) \log(1 - \sigma(\theta^T \mathbf{x}^{(i)}))$$

log conditional likelihood



$$LL(\theta) = \sum_{i=1}^n y^{(i)} \log \hat{y}^{(i)} + (1 - y^{(i)}) \log(1 - \hat{y}^{(i)})$$

$$\text{Let } \hat{y}^{(i)} = \sigma(\theta^T \mathbf{x}^{(i)})$$

Compute gradient of log conditional likelihood

$$\begin{aligned}\frac{\partial LL(\theta)}{\partial \theta_j} &= \sum_{i=1}^n \frac{\partial}{\partial \theta_j} [y^{(i)} \log(\hat{y}^{(i)}) + (1 - y^{(i)}) \log(1 - \hat{y}^{(i)})] && \text{Let } \hat{y}^{(i)} = \sigma(\theta^T \mathbf{x}^{(i)}) \\ &= \sum_{i=1}^n \frac{\partial}{\partial \hat{y}^{(i)}} [y^{(i)} \log(\hat{y}^{(i)}) + (1 - y^{(i)}) \log(1 - \hat{y}^{(i)})] \cdot \frac{\partial \hat{y}^{(i)}}{\partial \theta_j} && \text{(Chain Rule)} \\ &= \sum_{i=1}^n \left[y^{(i)} \frac{1}{\hat{y}^{(i)}} - (1 - y^{(i)}) \frac{1}{1 - \hat{y}^{(i)}} \right] \cdot \hat{y}^{(i)} (1 - \hat{y}^{(i)}) x_j^{(i)} && \text{(calculus)} \\ &= \sum_{i=1}^n [y^{(i)} - \hat{y}^{(i)}] x_j^{(i)} && = \sum_{i=1}^n [y^{(i)} - \sigma(\theta^T \mathbf{x}^{(i)})] x_j^{(i)} && \text{(simplify)}\end{aligned}$$

Compute gradient of log conditional likelihood

$$\begin{aligned}\frac{\partial LL(\theta)}{\partial \theta_j} &= \sum_{i=1}^n \frac{\partial}{\partial \theta_j} [y^{(i)} \log(\hat{y}^{(i)}) + (1 - y^{(i)}) \log(1 - \hat{y}^{(i)})] && \text{Let } \hat{y}^{(i)} = \sigma(\theta^T \mathbf{x}^{(i)}) \\ &= \sum_{i=1}^n \frac{\partial}{\partial \hat{y}^{(i)}} [y^{(i)} \log(\hat{y}^{(i)}) + (1 - y^{(i)}) \log(1 - \hat{y}^{(i)})] \cdot \frac{\partial \hat{y}^{(i)}}{\partial \theta_j} && \text{(Chain Rule)} \\ &= \sum_{i=1}^n \left[y^{(i)} \frac{1}{\hat{y}^{(i)}} - (1 - y^{(i)}) \frac{1}{1 - \hat{y}^{(i)}} \right] \cdot \hat{y}^{(i)} (1 - \hat{y}^{(i)}) x_j^{(i)} && \text{(calculus)} \\ &= \sum_{i=1}^n [y^{(i)} - \hat{y}^{(i)}] x_j^{(i)} && = \sum_{i=1}^n [y^{(i)} - \sigma(\theta^T \mathbf{x}^{(i)})] x_j^{(i)} \quad \text{(simplify)}\end{aligned}$$

