

25: Logistic Regression

David Varodayan

March 6, 2020

Adapted from slides by Lisa Yan

Announcements

Updated CS109 logistics and policies

We are implementing several changes to CS109 logistics and policies in response to growing concern about Covid-19. These changes affect lectures, section, office hours and the final exam.

Read the new policies in the announcements on the course webpage:

<http://web.stanford.edu/class/cs109/>

We welcome your questions on piazza

Announcements

Problem Set 6

Due: Wednesday 3/11
Covers: Up to Lecture 25
Extra Python Office Hours: Saturday 3/7, 3-5PM

Regrades

Pset 1 to 5 and
Midterm regrades to
close on 3/11 at 1pm

Autograded Coding Problems

Run your code in the command line or
install Pycharm following directions on
Pset 6 webpage

Late Day Reminder

No late days permitted past
last day of the quarter, 3/13

Today's plan

Logistic Regression



- Chapter 0: Background
- Chapter 1: Big Picture
- Chapter 2: Details
- Chapter 3: Philosophy

Background: Weighted sum

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

$$z = \theta^T \mathbf{X} = \sum_{j=1}^m \theta_j X_j$$

Weighted sum
(aka dot product)

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

Weighted sum
with an
intercept term:

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

$$= \theta_0 X_0 + \theta_1 X_1 + \theta_2 X_2 + \dots + \theta_m X_m \quad \text{Define } X_0 = 1$$

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

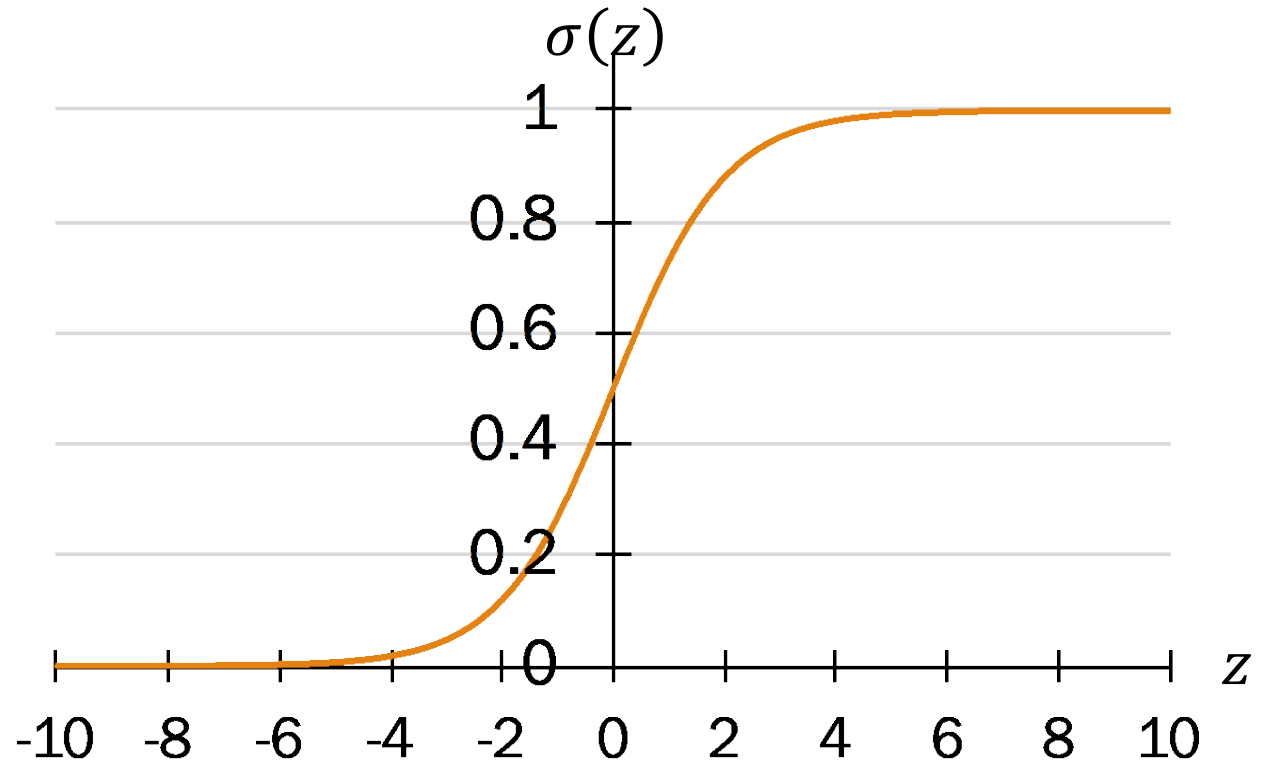
New $\mathbf{X} = (1, X_1, X_2, \dots, X_m)$

Background: 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.

Background: Chain Rule


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

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

Calculus
Chain Rule

Today's plan

Logistic Regression

- Chapter 0: Background
-  • Chapter 1: Big Picture
- Chapter 2: Details
- Chapter 3: Philosophy

From Naïve Bayes to Logistic Regression

Classification goal:

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

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

Predict the Y that is most likely given our observation \mathbf{X}

Naïve Bayes Classifier:

- Estimate $P(\mathbf{X} | Y)$ and $P(Y)$ because $\arg \max_{y=\{0,1\}} P(Y | \mathbf{X}) = \arg \max_{y=\{0,1\}} P(\mathbf{X}|Y)P(Y)$
- Actually modeling $P(\mathbf{X}, Y)$
- Assume $P(\mathbf{X}|Y) = P(X_1, X_2, \dots, X_n|Y) = \prod_{i=1}^m P(X_i|Y)$

Can we model $P(Y | \mathbf{X})$ directly?

- Welcome our friend: Logistic Regression!

Logistic Regression

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

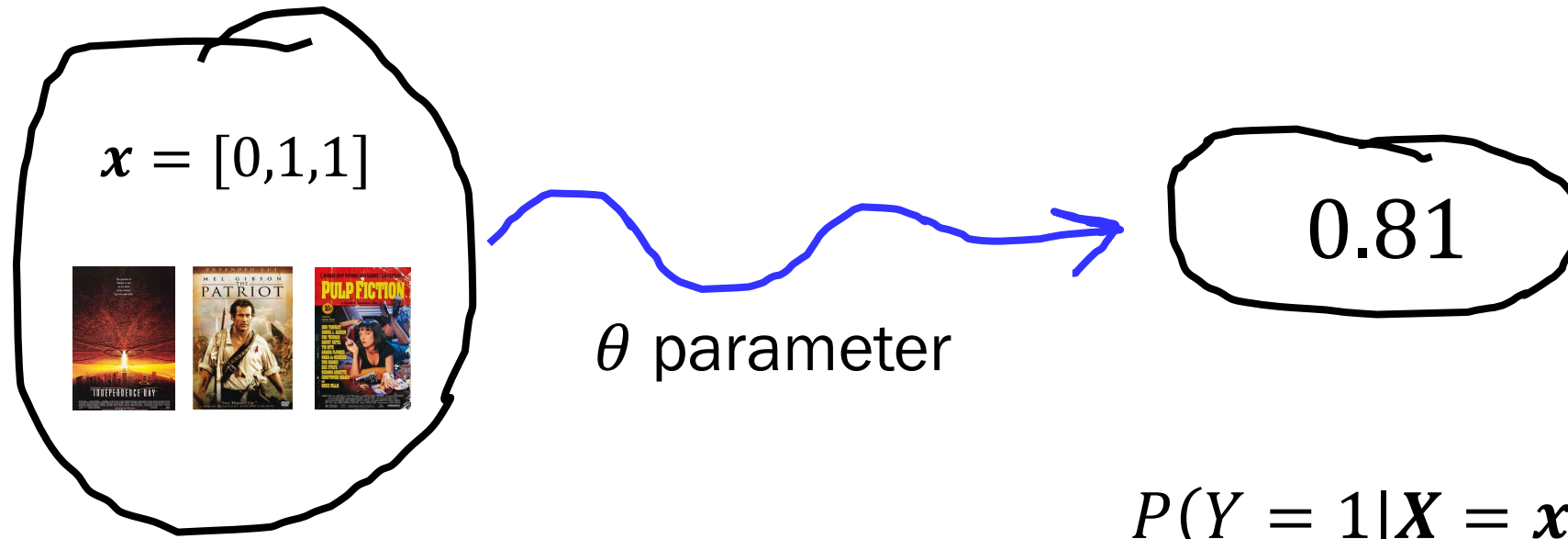
Predict the Y that is most likely given our observation \mathbf{X}

Logistic
Regression
Model

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

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

Logistic Regression



X
input features

$P(Y = 1 | X = x)$
conditional likelihood

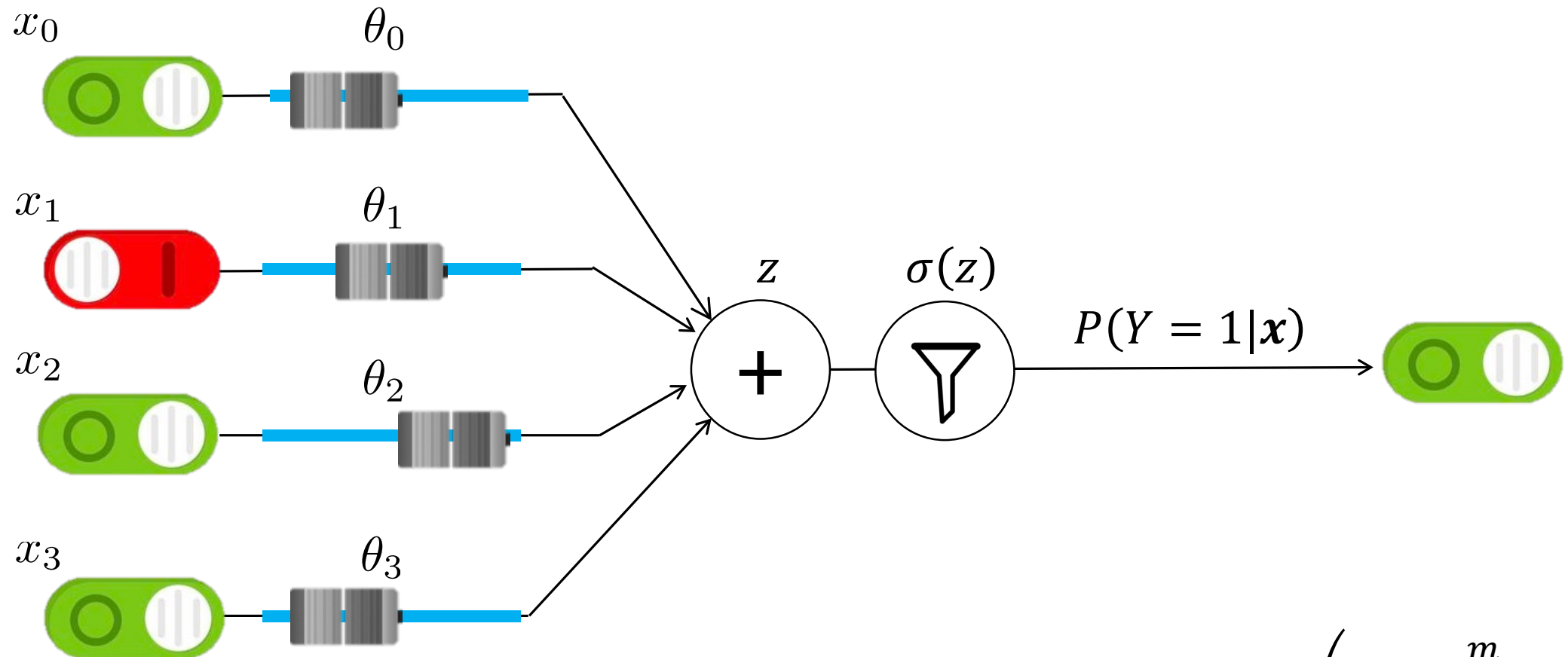
$$P(Y = 1 | X = x) = \sigma \left(\theta_0 + \sum_{j=1}^m \theta_j x_j \right)$$

Logistic Regression Cartoon



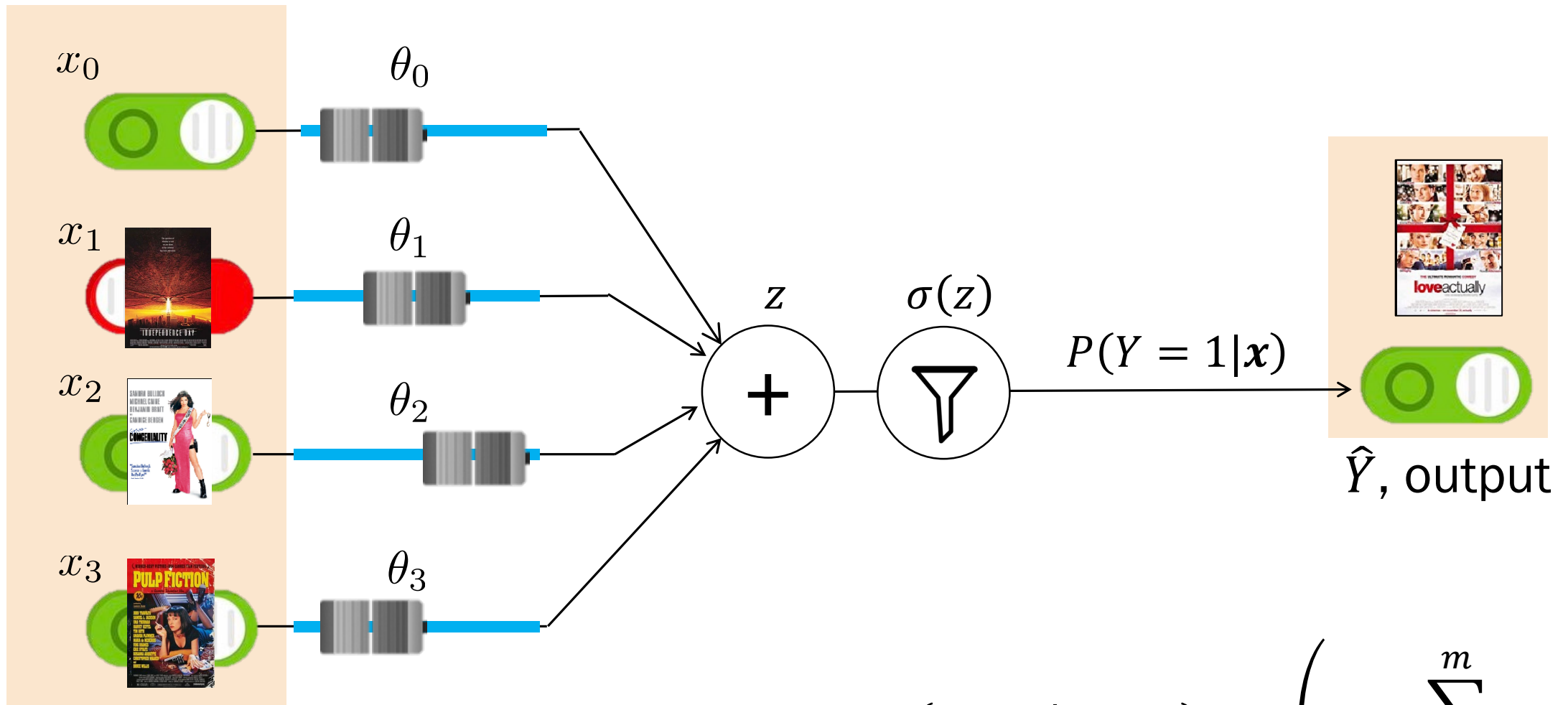
θ parameter

Logistic Regression cartoon



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

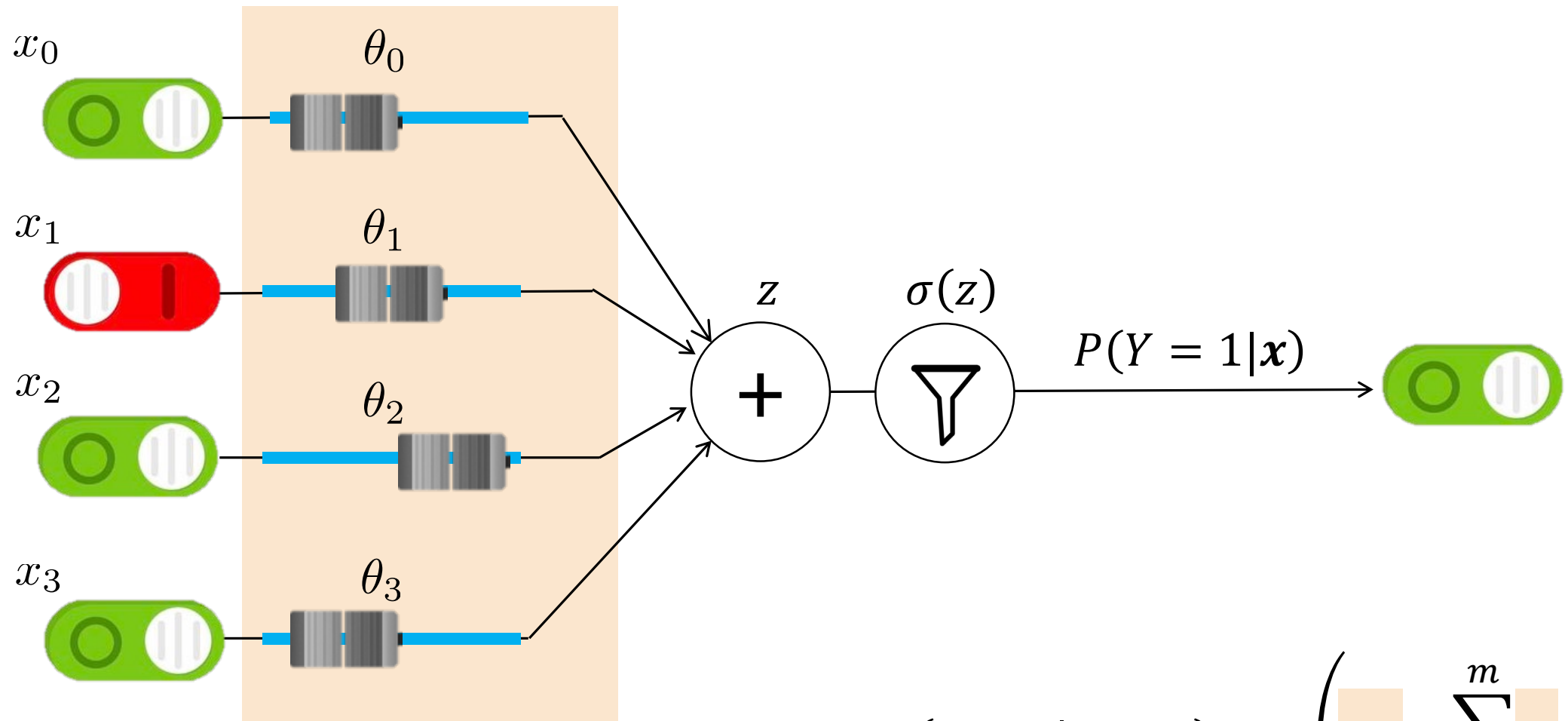
Logistic Regression input/output



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

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

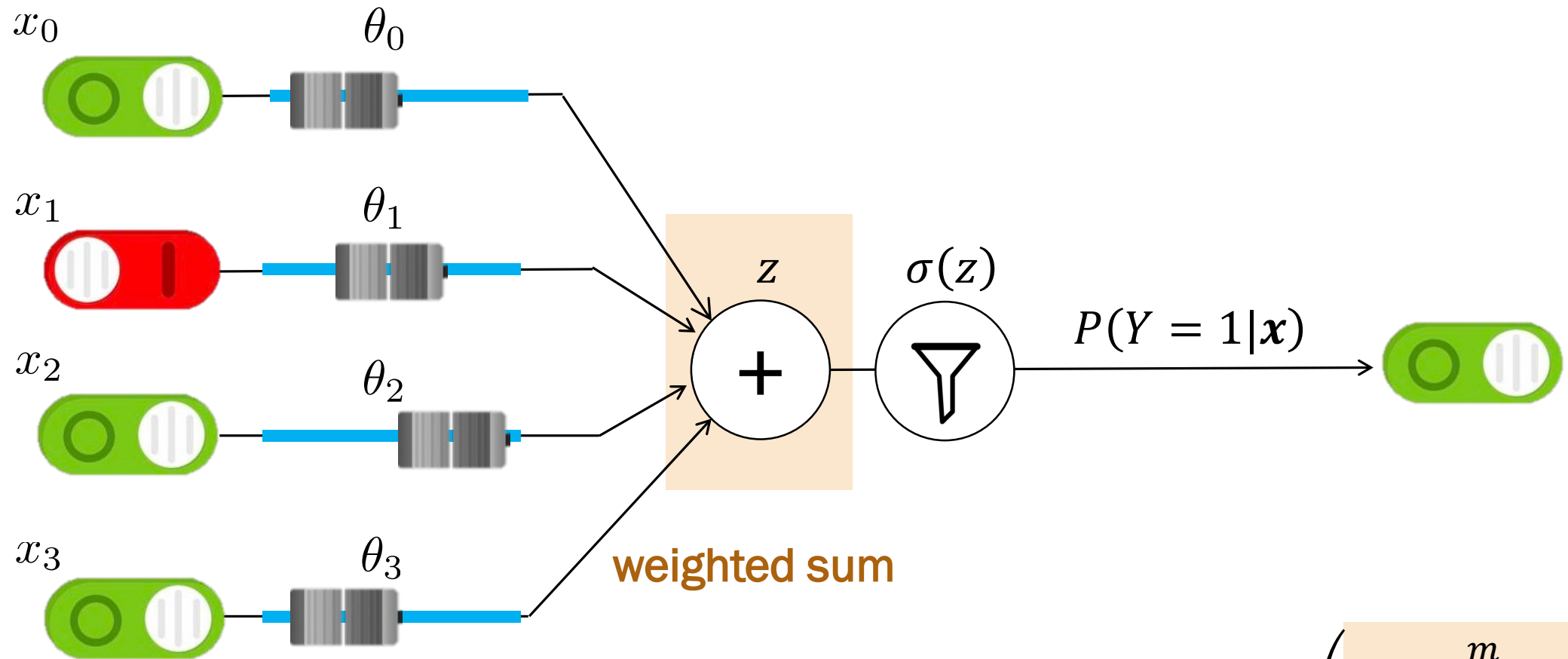
Components of Logistic Regression



θ weights
(aka parameters)

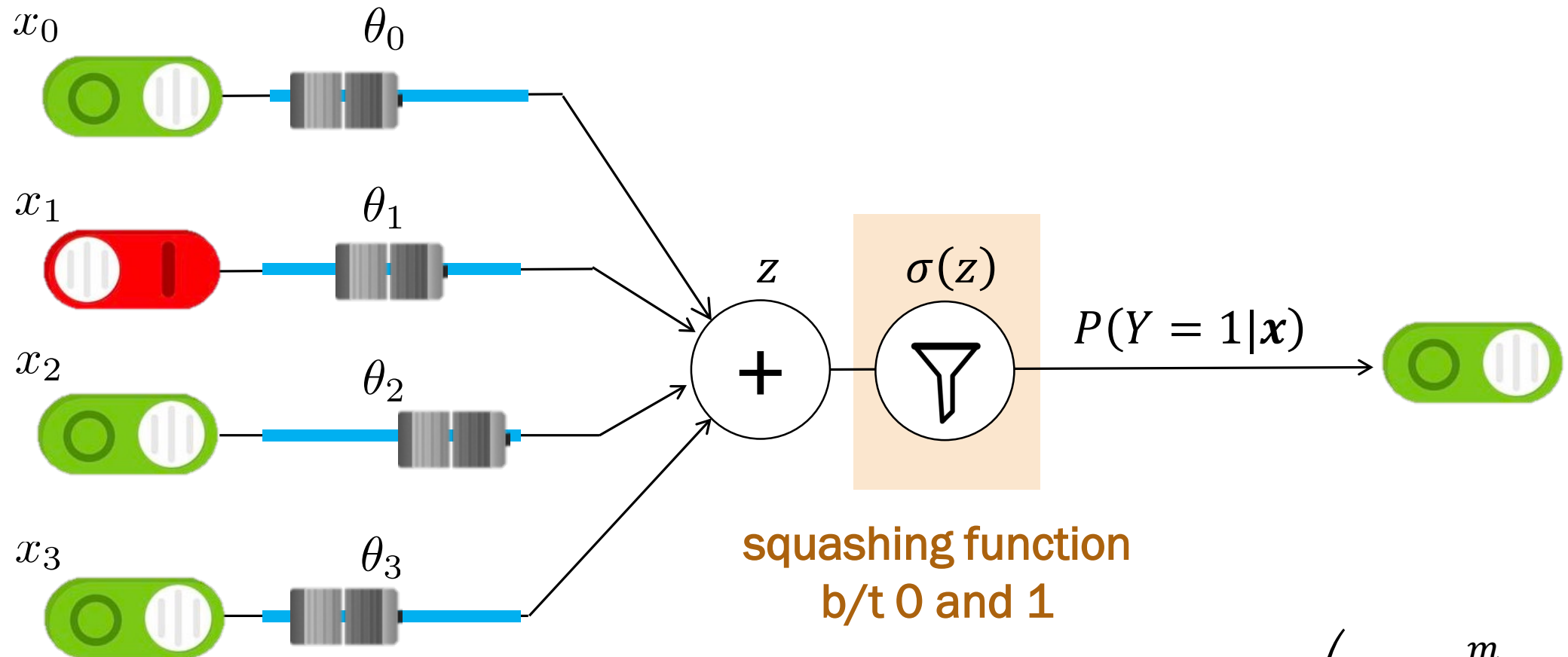
$$P(Y = 1 | \mathbf{X} = \mathbf{x}) = \sigma \left(\theta_0 + \sum_{j=1}^m \theta_j x_j \right)$$

Components of Logistic Regression



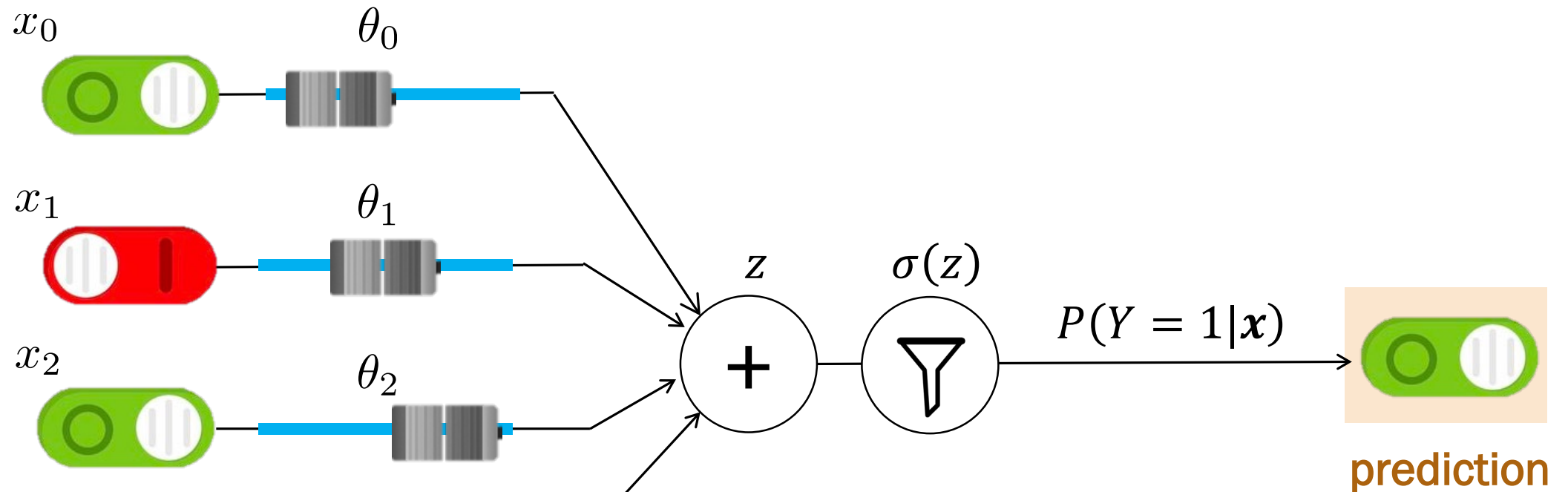
$$P(Y = 1 | \mathbf{X} = \mathbf{x}) = \sigma \left(\theta_0 + \sum_{j=1}^m \theta_j x_j \right)$$

Components of Logistic Regression



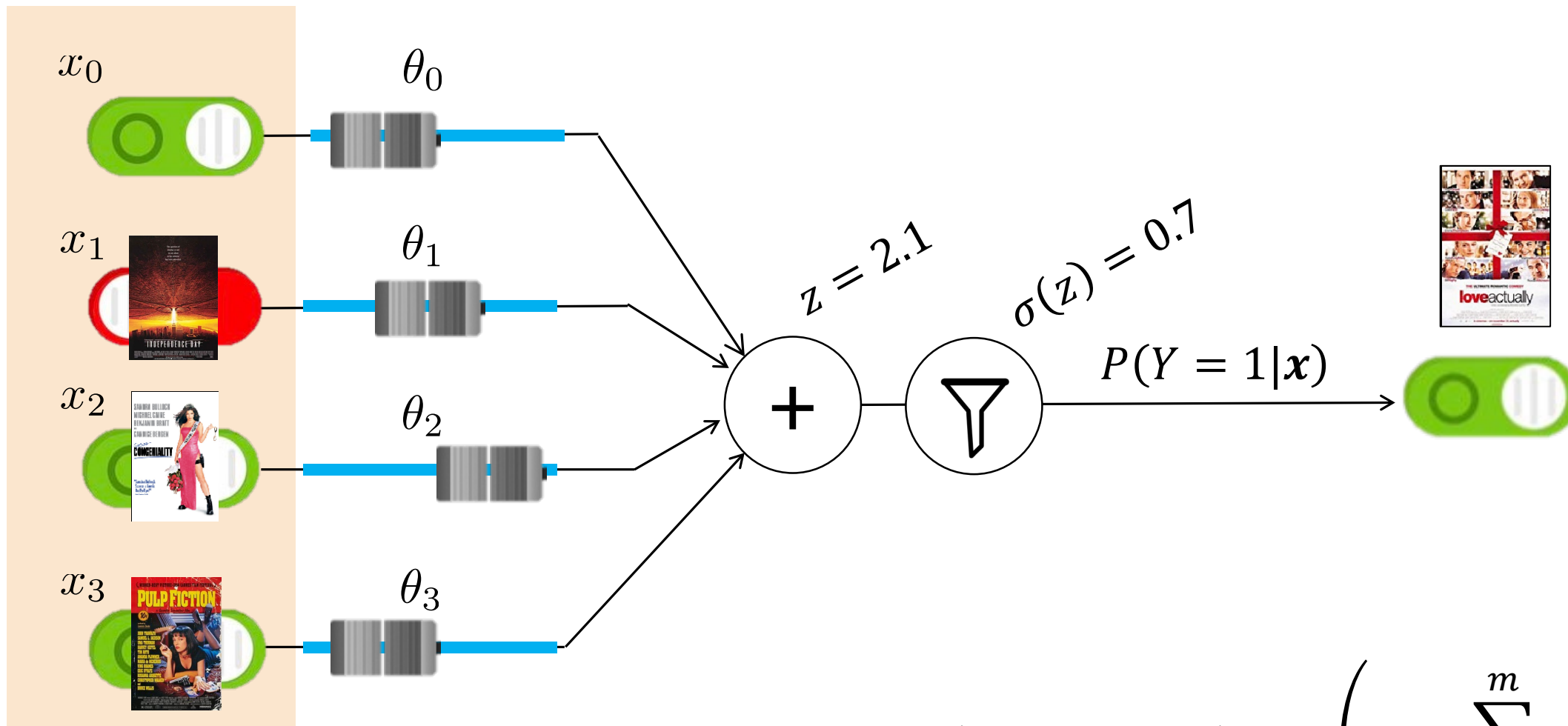
$$P(Y = 1 | \mathbf{X} = \mathbf{x}) = \sigma \left(\theta_0 + \sum_{j=1}^m \theta_j x_j \right)$$

Components of Logistic Regression



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

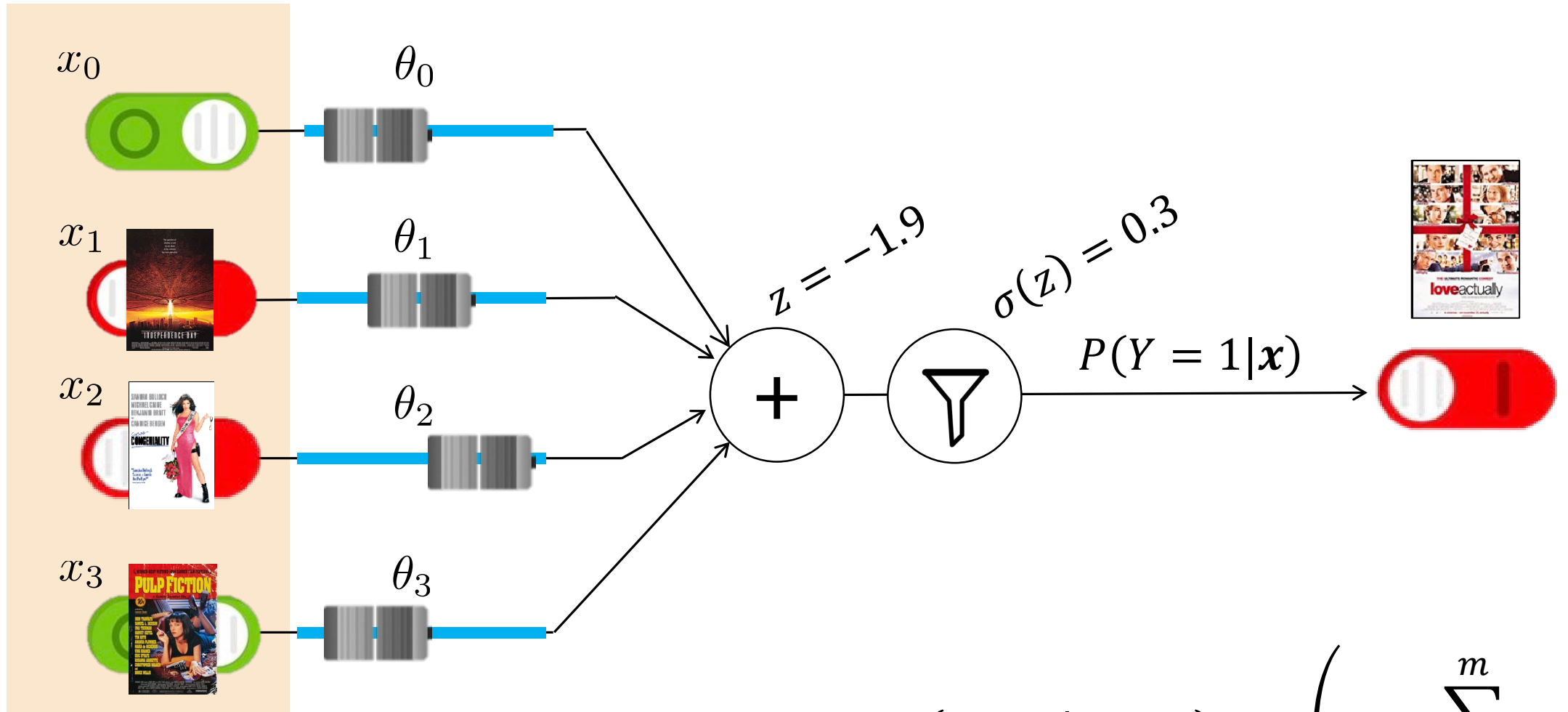
Different predictions for different inputs



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

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

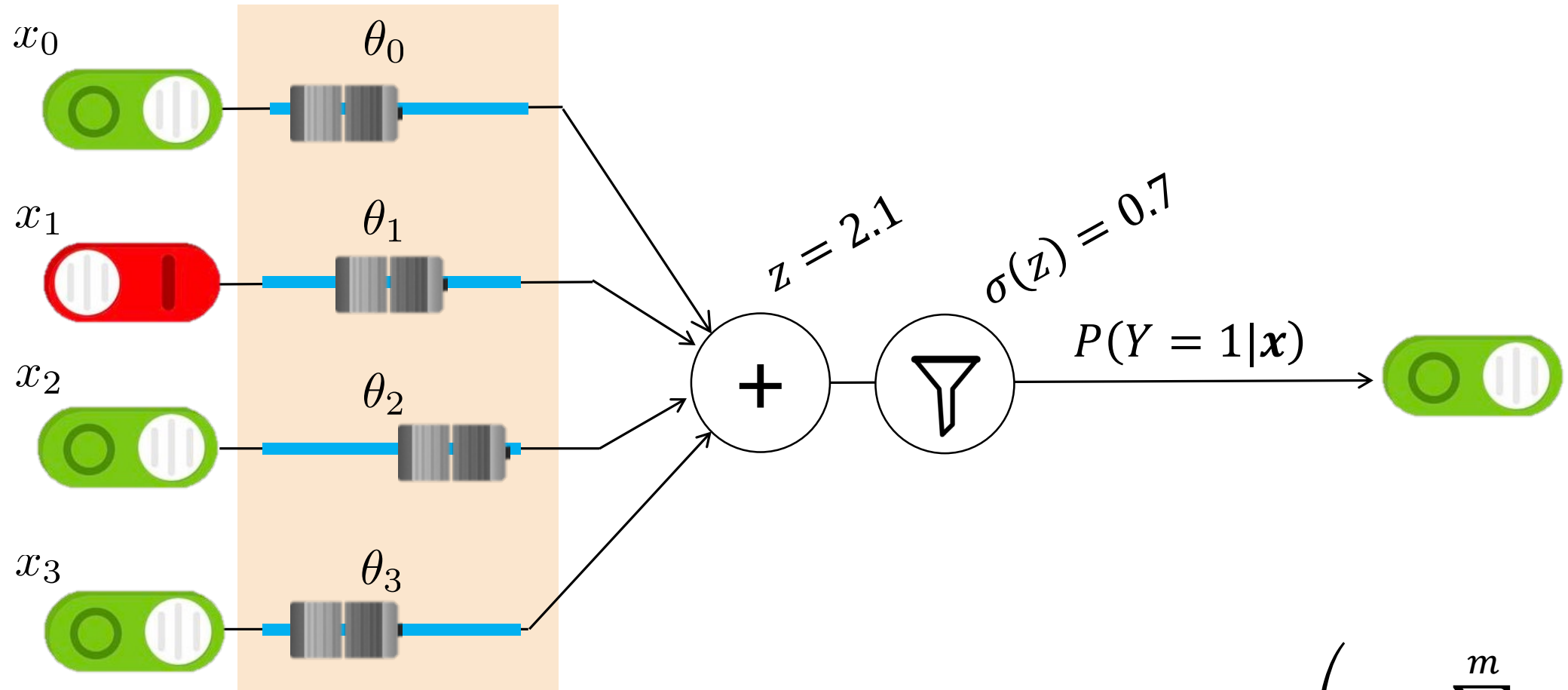
Different predictions for different inputs



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

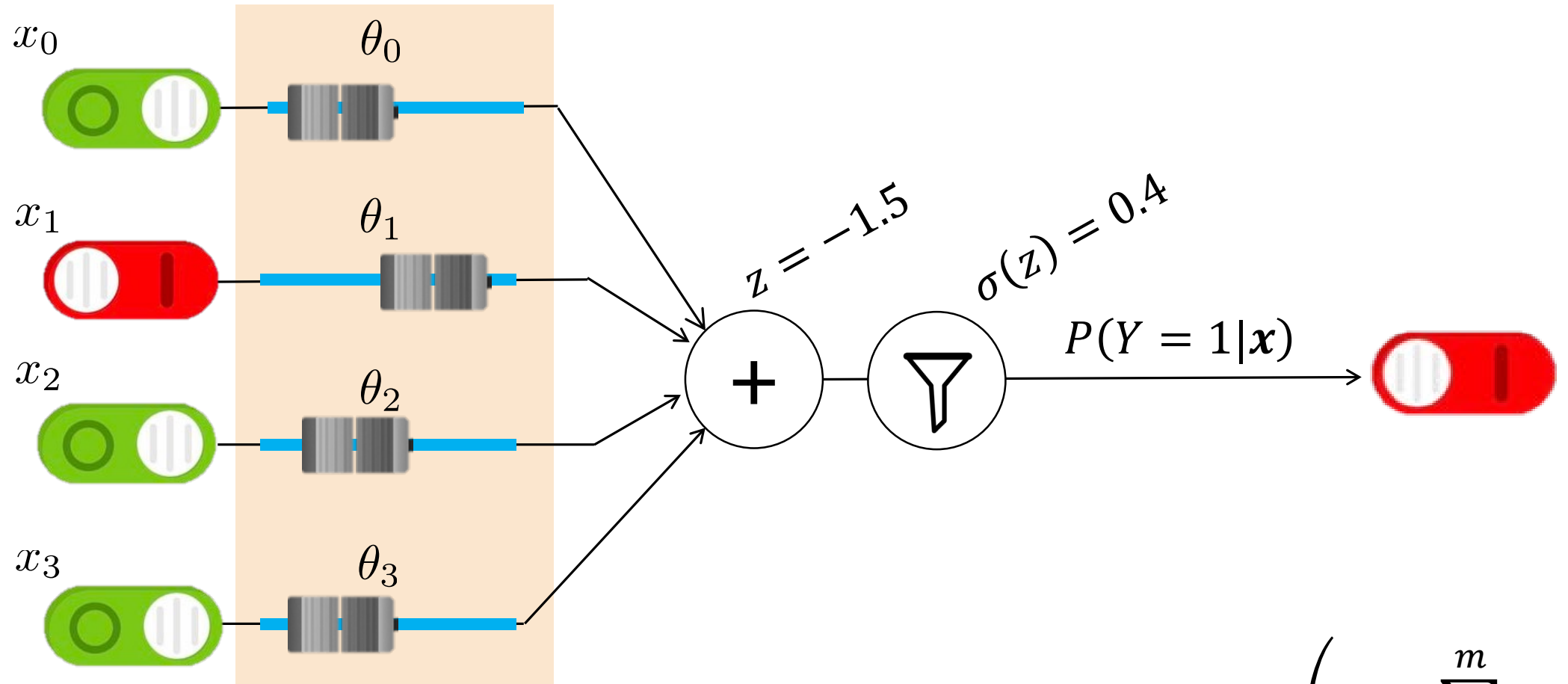
$$P(Y = 1 | \mathbf{X} = \mathbf{x}) = \sigma \left(\theta_0 + \sum_{j=1}^m \theta_j x_j \right)$$

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



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

Logistic Regression Model

$$\hat{Y} = \arg \max_{y=\{0,1\}} P(Y | \mathbf{X}) \quad \text{where} \quad P(Y = 1 | \mathbf{X} = \mathbf{x}) = \sigma \left(\theta_0 + \sum_{j=1}^m \theta_j x_j \right)$$

Predict the Y that is most likely given our observation \mathbf{X}

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

- $\sigma(z) = \frac{1}{1+e^{-z}}$, the sigmoid function

- For simplicity, define $x_0 = 1$: $P(Y = 1 | \mathbf{X} = \mathbf{x}) = \sigma(\theta^T \mathbf{x})$

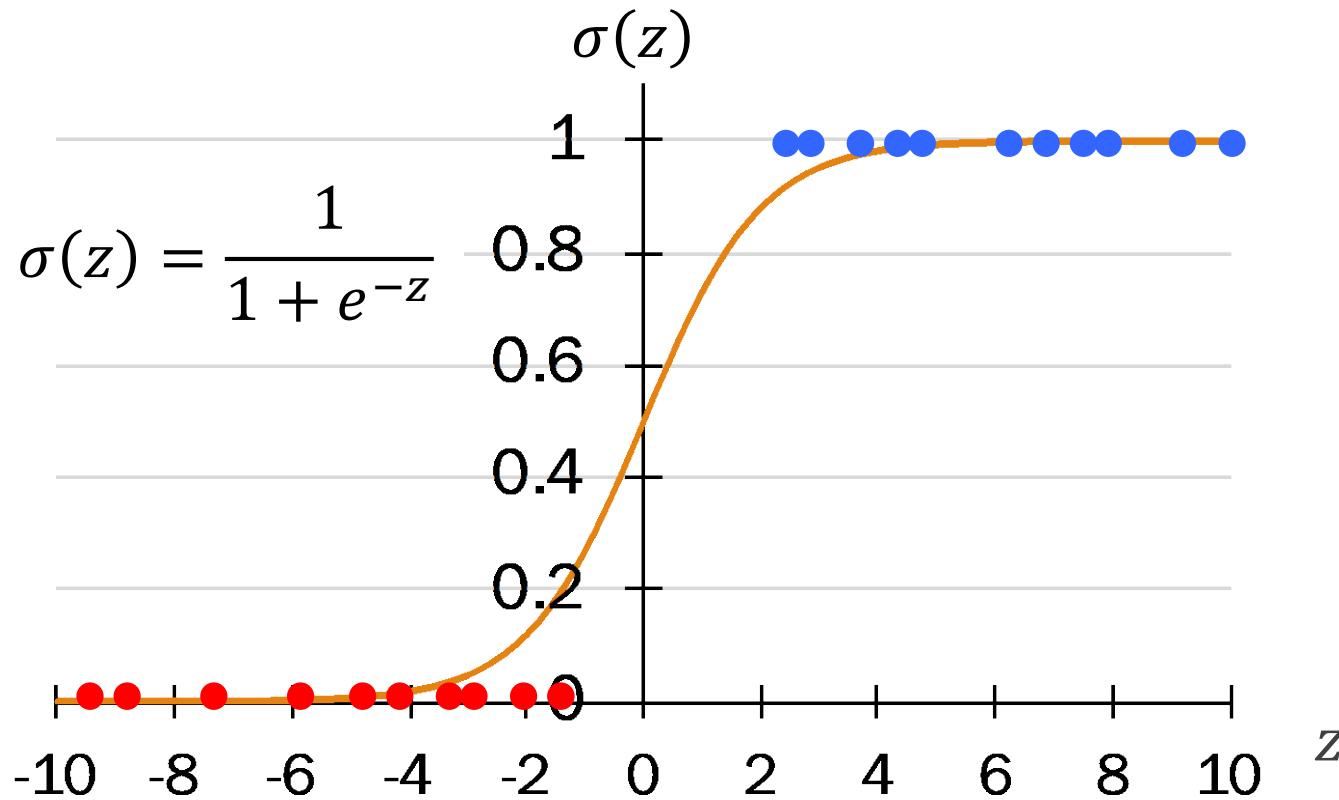
- Since $P(Y = 1 | \mathbf{X} = \mathbf{x}) + P(Y = 0 | \mathbf{X} = \mathbf{x}) = 1$:

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

Classifying using the sigmoid function

Logistic
Regression
Model

$$\hat{Y} = \arg \max_{y=\{0,1\}} P(Y | \mathbf{X}) \quad \text{where} \quad P(Y = 1 | \mathbf{X} = \mathbf{x}) = \sigma \left(\theta_0 + \sum_{j=1}^m \theta_j x_j \right)$$

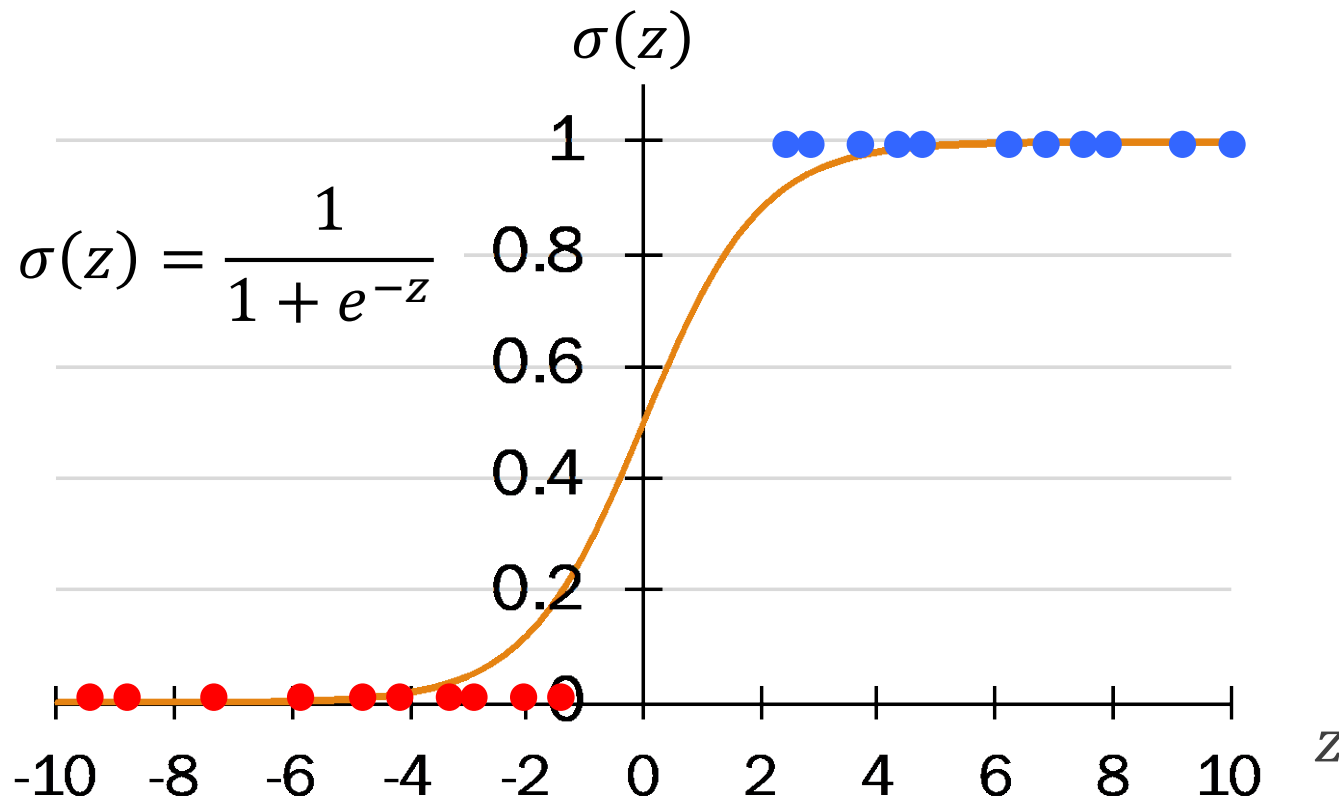


Logistic Regression uses the sigmoid function to try and distinguish $y = 1$ (blue) points from $y = 0$ (red) points.

Classifying using the sigmoid function

Logistic
Regression
Model

$$\hat{Y} = \arg \max_{y=\{0,1\}} P(Y | \mathbf{X}) \quad \text{where} \quad P(Y = 1 | \mathbf{X} = \mathbf{x}) = \sigma \left(\theta_0 + \sum_{j=1}^m \theta_j x_j \right)$$



When do we predict $\hat{Y} = 1$?

- A. If $\sigma(\theta^T \mathbf{x}) > 1 - \sigma(\theta^T \mathbf{x})$
- B. If $\sigma(\theta^T \mathbf{x}) > 0.5$
- C. If $\theta^T \mathbf{x} > 0$
- D. All are valid, but C is easiest
- E. None/Other

Naming algorithms

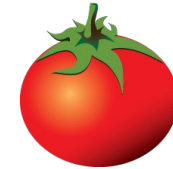
Regression Algorithms

Linear Regression

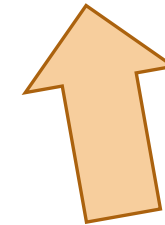


Classification Algorithms

Naïve Bayes



Logistic Regression



Awesome classifier,
terrible name

Training: Learning the parameters

Logistic regression gets its **intelligence** from its parameters $\theta = (\theta_0, \theta_1, \dots, \theta_m)$.

- Logistic Regression Model:

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

- Want to predict training data as correctly as possible:

$$\arg \max_{y \in \{0,1\}} P(Y = y | \mathbf{X} = \mathbf{x}^{(i)}) = y^{(i)} \quad \text{as often as possible}$$

- Therefore, choose θ that maximizes the **conditional likelihood** of observing i.i.d. training data:

$$L(\theta) = \prod_{i=1}^n P(Y = y^{(i)} | \mathbf{X} = \mathbf{x}^{(i)}, \theta)$$

During training, find the θ that maximizes log-conditional likelihood of the training data. Use MLE!

Training: Learning the parameters via MLE

0. Add $x_0^{(i)} = 1$ to each $\mathbf{x}^{(i)}$

1. Logistic Regression model:

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

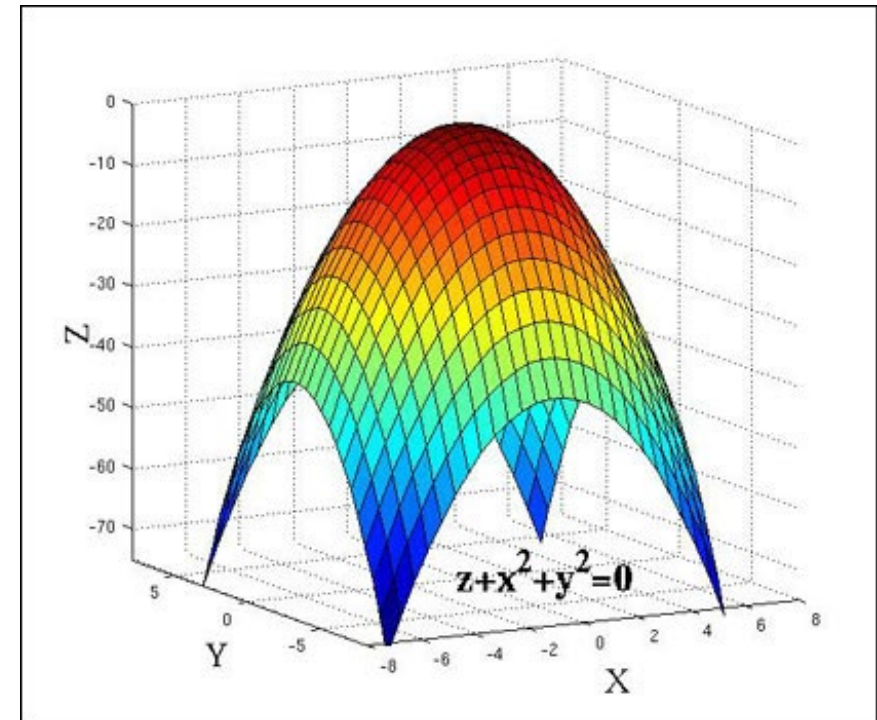
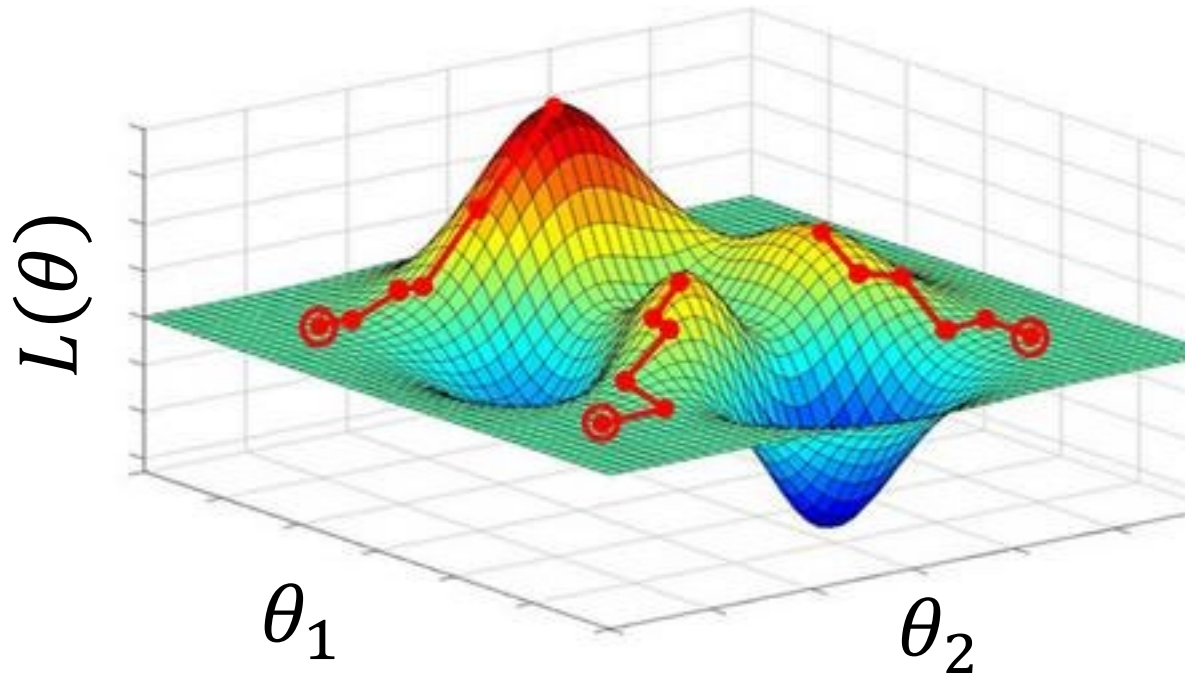
2. Compute log-likelihood of training data:

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

3. Compute derivative of log-likelihood with respect to each $\theta_j, j = 0, 1, \dots, m$:

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

Walk uphill and you will find a local maxima
(if your step is small enough).



Logistic regression $LL(\theta)$
is convex

Training: Gradient ascent step

4. Optimize.

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

Repeat many times:

For all thetas:

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

What does this look like in code?

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  $0 \leq j \leq m$ 
```

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

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

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

Training: Gradient Ascent

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

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

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

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 $0 \leq j \leq m$

What are important implementation details?

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 $0 \leq j \leq m$

- x_j is j -th feature of input var $x = (x_1, \dots, x_m)$

Training: Gradient Ascent

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

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 $0 \leq j \leq m$

- x_j is j -th feature of input var $x = (x_1, \dots, x_m)$
- Insert $x_0 = 1$ before training

Training: Gradient Ascent

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

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 $0 \leq j \leq m$

- x_j is j -th feature of input var $x = (x_1, \dots, x_m)$
- Insert $x_0 = 1$ before training
- Finish computing gradient before updating any part of θ

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 $0 \leq j \leq m$:

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

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

- x_j is j -th feature of input var $x = (x_1, \dots, x_m)$
- Insert $x_0 = 1$ before training
- Finish computing gradient before updating any part of θ
- Learning rate η is a constant you set before training

Training: Gradient Ascent

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

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 $0 \leq j \leq m$

- x_j is j -th feature of input var $x = (x_1, \dots, x_m)$
- Insert $x_0 = 1$ before training
- Finish computing gradient before updating any part of θ
- Learning rate η is a constant you set before training

Testing: Classification with Logistic Regression

Training

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

via gradient
ascent:

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

Testing

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

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

Parameters θ_j are not updated during testing phase

Today's plan

Logistic Regression

- Chapter 0: Background
- Chapter 1: Big Picture
- Chapter 2: Details
- Chapter 3: Philosophy



Introducing notation \hat{y}

Logistic
Regression
model:

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



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

Prediction:

$$\hat{Y} = \arg \max_{y=\{0,1\}} P(Y | \mathbf{X} = \mathbf{x}) = \begin{cases} 1 & \text{if } \hat{y} > 0.5 \\ 0 & \text{otherwise} \end{cases}$$

Training: Learning the parameters via MLE

0. Add $x_0^{(i)} = 1$ to each $\mathbf{x}^{(i)}$

1. Logistic Regression model:

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

2. Compute log-likelihood of training data:

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

3. Compute derivative of **How did we get this log-likelihood function?**

log-likelihood with respect to each $\theta_j, j = 0, 1, \dots, m$:

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

Log-likelihood of data

Logistic
Regression
model:

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

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

(see Bernoulli
MLE PMF)

Likelihood
of training data:

$$L(\theta) = \prod_{i=1}^n P(Y = y^{(i)} | \mathbf{X} = \mathbf{x}^{(i)}, \theta)$$

Notes:

- Actually **conditional likelihood**
- Still correctly gets correct θ_{MLE} since \mathbf{X}, θ independent
- See lecture notes

Log-likelihood of data

Logistic Regression model:

$$P(Y = y | \mathbf{X} = \mathbf{x}) = \begin{cases} \hat{y} & \text{if } y = 1 \\ 1 - \hat{y} & \text{if } y = 0 \end{cases} \quad \text{where } \hat{y} = \sigma(\theta^T \mathbf{x})$$
$$= (\hat{y})^y (1 - \hat{y})^{1-y} \quad \text{(see Bernoulli MLE PMF)}$$

Likelihood of training data:

$$L(\theta) = \prod_{i=1}^n P(Y = y^{(i)} | \mathbf{X} = \mathbf{x}^{(i)}, \theta) = \prod_{i=1}^n (\hat{y}^{(i)})^{y^{(i)}} (1 - \hat{y}^{(i)})^{1-y^{(i)}}$$

Log-likelihood:

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

Training: Learning the parameters via MLE

0. Add $x_0^{(i)} = 1$ to each $\mathbf{x}^{(i)}$

1. Logistic Regression model:

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

2. Compute log-likelihood of training data:

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

3. Compute derivative of log-likelihood with respect to each $\theta_j, j = 0, 1, \dots, m$:

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

How did we get this gradient?

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

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

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

Compute gradient of log-conditional likelihood

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

where

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

Compute gradient of log-likelihood

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

Today's plan

Logistic Regression

- Chapter 0: Background
- Chapter 1: Big Picture
- Chapter 2: Details
- Chapter 3: Philosophy



Intuition about Logistic Regression

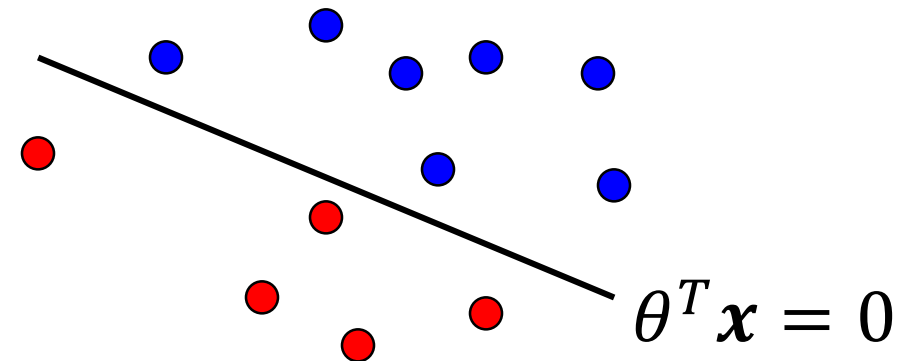
Logistic
Regression
Model

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

where

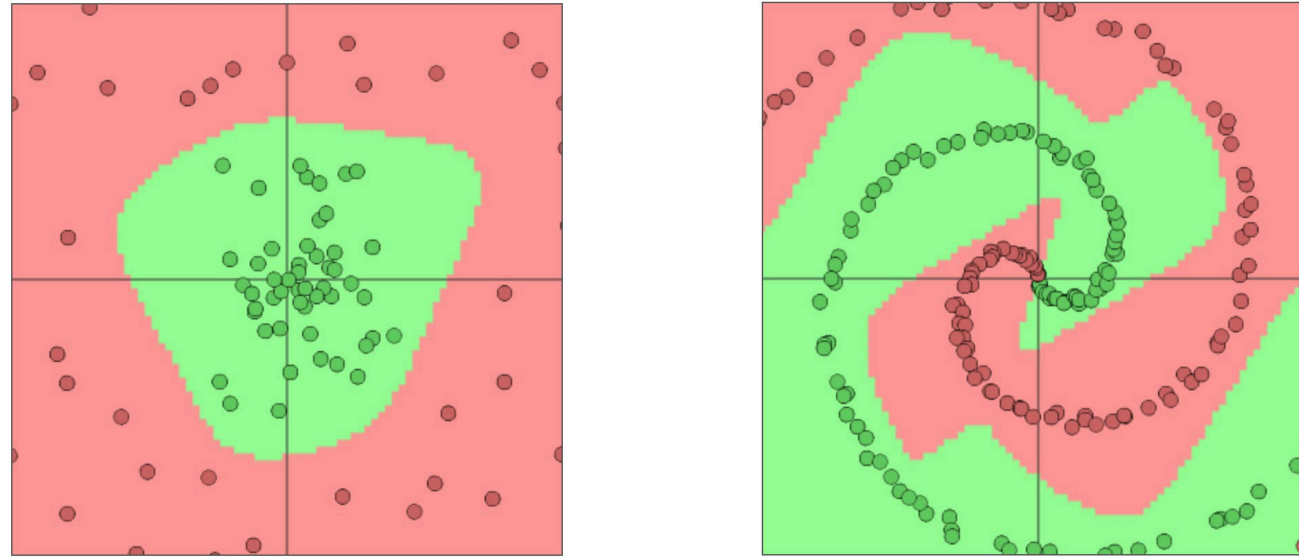
$$\theta^T \mathbf{x} = \sum_{j=0}^m \theta_j x_j$$

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 no interaction between different 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

Many tradeoffs in choosing an algorithm

	Naïve Bayes	Logistic Regression
Modeling goal	$P(\mathbf{X}, Y)$	$P(Y \mathbf{X})$
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 $P(\mathbf{X}, Y)$)
Continuous input features?	Needs parametric form (e.g., Gaussian) or discretized buckets (for multinomial features)	Yes, easily
Discrete input features?	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 $\{1, 2, 3\}$)