# 26: Logistic Regression

Lisa Yan and Jerry Cain
November 11, 2020

# Quick slide reference

# Background

# 1. Weighted sum

If $\boldsymbol{X} = (X_1, X_2, \ldots, X_m)$:

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

$$= \sum_{j=1}^{m} \theta_j X_j \qquad \text{weighted sum}$$

$$= \theta^T \boldsymbol{X} \qquad \text{dot product}$$

$$\begin{bmatrix} \theta_1 & \theta_2 & & \theta_m \end{bmatrix} \begin{bmatrix} X_1 \\ \\ \\ X_m \end{bmatrix}$$

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

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

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

# 1. Weighted sum

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

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

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

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

$$= \theta^T X \qquad\qquad \text{New } X = (1, X_1, X_2, \ldots, X_m) \quad, \theta = (\theta_0, \theta_1, \ldots, \theta_m)$$
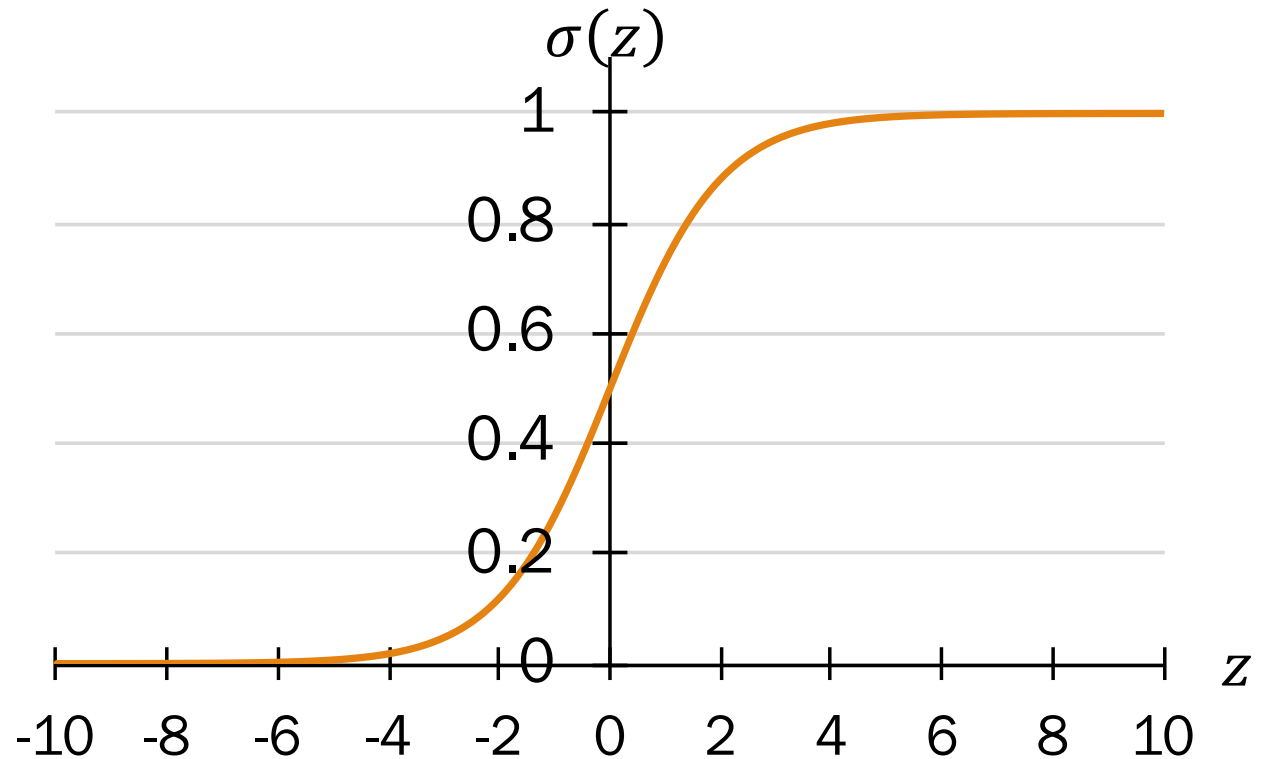
> Prepending $X_0 = 1$ to each feature vector $X$ makes matrix operators more accessible.

# 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):

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

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

**conditional likelihood**
of training data

$$= \arg\max_{\theta} \sum_{i=1}^{n} \log f\left(y^{(i)} | \boldsymbol{x}^{(i)}, \theta\right)$$
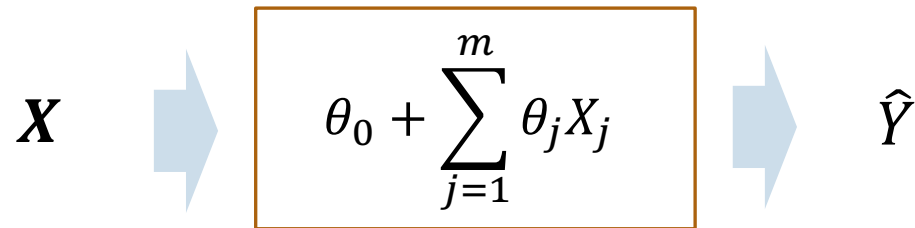
**log conditional likelihood**

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

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

# Logistic Regression
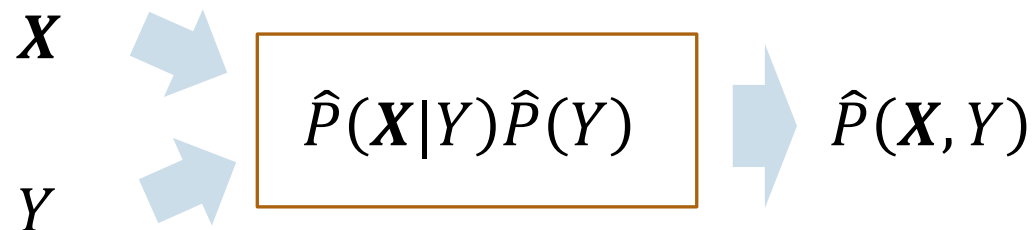
# Prediction models so far

## Linear Regression (Regression)

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

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

✅ $X$ can be dependent

🤷‍♀️ Regression model ($\hat{Y} \in \mathbb{R}$, not discrete)

## Naïve Bayes (Classification)

$$X, Y \Rightarrow \boxed{\hat{P}(X|Y)\hat{P}(Y)} \Rightarrow \hat{P}(X, Y)$$

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

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

✅ Tractable with NB assumption, but...

⚠️ Realistically, $X_j$ features not necessarily conditionally independent

🤷‍♀️ Actually models $P(X, Y)$, not $P(Y|X)$?

# Introducing Logistic Regression!

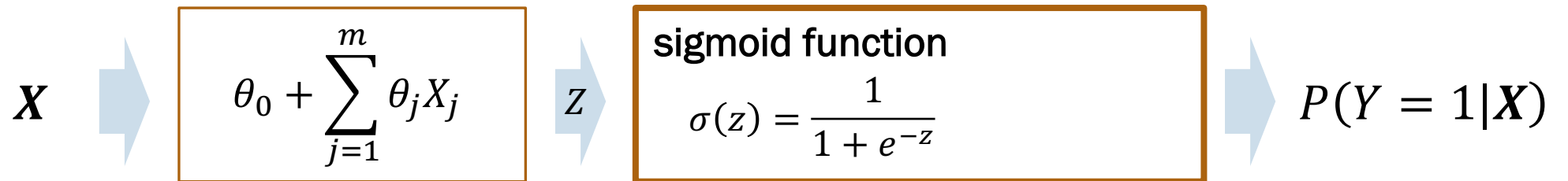Linear Regression ideas          Classification models

*+ compute power*

# Logistic Regression

$$X \implies \boxed{\theta_0 + \sum_{j=1}^{m} \theta_j X_j} \implies z \implies \boxed{\begin{array}{l}\text{sigmoid function} \\ \sigma(z) = \dfrac{1}{1 + e^{-z}}\end{array}} \implies P(Y = 1 | X)$$

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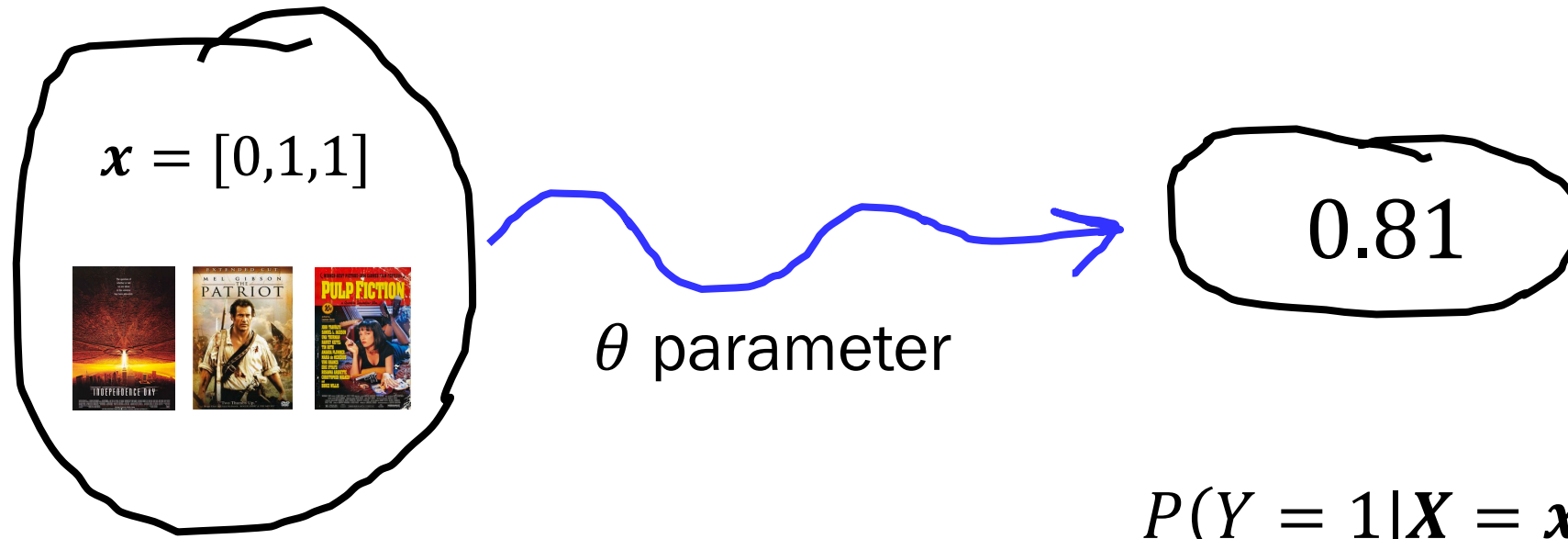
$$P(Y = 1 | X = x) = \sigma\left( \theta_0 + \sum_{j=1}^{m} \theta_j x_j \right)$$

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

$$\hat{Y} = \underset{y=\{0,1\}}{\arg\max} \, P(Y \mid X)$$

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

# Logistic Regression



$$x = [0,1,1]$$

$\theta$ parameter

0.81

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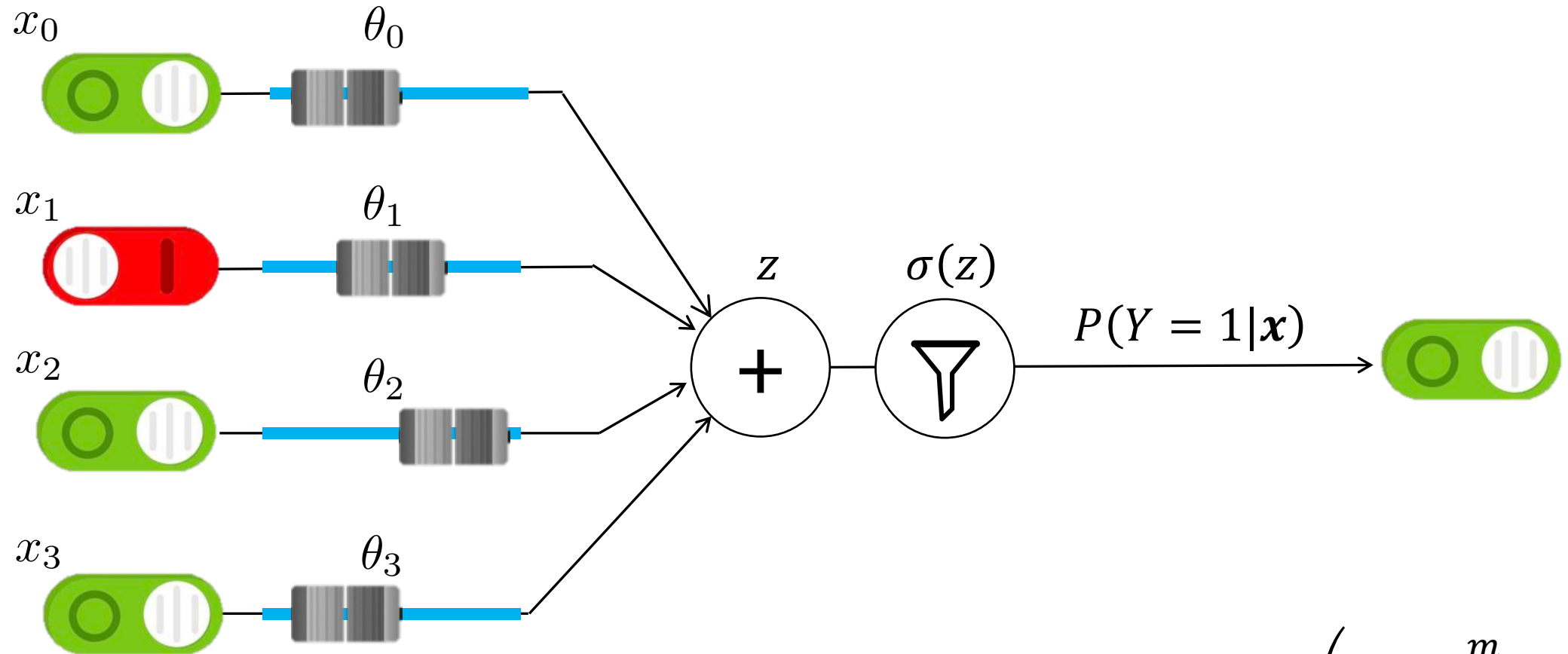
$X$
input features

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

# Logistic Regression cartoon



$\theta$ parameter

# Logistic Regression cartoon



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

# Logistic Regression cartoon



$x_0$

$\theta_0$

$x_1$

$\theta_1$

$x_2$

$\theta_2$

$x_3$

$\theta_3$

$z$

$\sigma(z)$

$P(Y = 1 | \boldsymbol{x})$

$\hat{Y}$, output

$\boldsymbol{X}$, input features
[0,1,1]

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

# Components of Logistic Regression



$x_0$

$\theta_0$

$x_1$

$\theta_1$

$z$     $\sigma(z)$

$P(Y = 1 | \boldsymbol{x})$

$x_2$

$\theta_2$

$+$

$x_3$

$\theta_3$

$\theta$ weights
(aka parameters)

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

# Components of Logistic Regression



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

# Components of Logistic Regression



$x_0$

$\theta_0$

$x_1$

$\theta_1$

$x_2$

$\theta_2$

$x_3$

$\theta_3$

$z$

$\sigma(z)$

$P(Y = 1|\boldsymbol{x})$

squashing function
b/t 0 and 1

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

# Components of Logistic Regression



$x_0$

$\theta_0$

$x_1$

$\theta_1$

$x_2$

$\theta_2$

$x_3$

$\theta_3$

$z$

$\sigma(z)$

$P(Y = 1|\boldsymbol{x})$

prediction

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

# Different predictions for different inputs



$x_0$

$\theta_0$

$x_1$

$\theta_1$

$x_2$

$\theta_2$

$x_3$

$\theta_3$

$z = 2.1$

$\sigma(z) = 0.7$

$P(Y = 1 | \boldsymbol{x})$

$\boldsymbol{X}$, input features
$[0,1,1]$

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

# Different predictions for different inputs



$x_0$

$\theta_0$

$x_1$

$\theta_1$

$x_2$

$\theta_2$

$x_3$

$\theta_3$

$z = -1.9$

$\sigma(z) = 0.3$

$P(Y = 1 | \boldsymbol{x})$

$\boldsymbol{X}$, input features $[0,0,1]$

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

# Parameters affect prediction



$x_0$    $\theta_0$

$x_1$    $\theta_1$

$x_2$    $\theta_2$

$x_3$    $\theta_3$

$z = 2.1$

$\sigma(z) = 0.7$

$P(Y = 1|\boldsymbol{x})$

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

# Parameters affect prediction



$x_0$

$\theta_0$

$x_1$

$\theta_1$

$z = -1.5$

$\sigma(z) = 0.4$

$x_2$

$\theta_2$

$P(Y = 1|\boldsymbol{x})$

$x_3$

$\theta_3$

$+$

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

# For simplicity

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

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

# Logistic regression classifier

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

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

**Training**
Estimate parameters from training data

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

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

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

# Training:
# The big picture

# Logistic regression classifier

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

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

**Training**    Estimate parameters from training data

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

Choose $\theta$ that optimizes some objective:

1. Determine objective function
2. Find gradient with respect to $\theta$
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 $\theta$

1. Determine objective function

$$\theta_{MLE} = \arg\max_{\theta} \prod_{i=1}^{n} f\left(y^{(i)} \mid \boldsymbol{x}^{(i)}, \theta\right)$$

2. Gradient w.r.t. $\theta_j$, for $j = 0, 1, \ldots, m$

3. Solve
   - No analytical derivation of $\theta_{MLE}$…
   - …but can still compute $\theta_{MLE}$ with gradient ascent!

```
initialize x
repeat many times:
    compute gradient
    x += η * gradient
```

# 1. Determine objective function

$$\theta_{MLE} = \boxed{\arg\max_\theta \prod_{i=1}^{n} f\left(y^{(i)} \mid \boldsymbol{x}^{(i)}, \theta\right)} = \boxed{\arg\max_\theta \; LL(\theta)}$$

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

First: Interpret conditional likelihood with Logistic Regression

Second: Write a differentiable expression for log conditional likelihood

# 1. Determine objective function (interpret)

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

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

Suppose you have $n = 2$ training datapoints:     $\big(\boldsymbol{x}^{(1)}, 1\big), \big(\boldsymbol{x}^{(2)}, 0\big)$

Consider the following expressions for a given $\theta$:

A.  $\sigma\big(\theta^T \boldsymbol{x}^{(1)}\big)\, \sigma\big(\theta^T \boldsymbol{x}^{(2)}\big)$

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

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

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

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

# 1. Determine objective function (interpret)

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

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

Suppose you have $n = 2$ training datapoints:  $\left(x^{(1)}, 1\right), \left(x^{(2)}, 0\right)$

Consider the following expressions for a given $\theta$:

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

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

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

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

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

# 1. Determine objective function (write)

$$\theta_{MLE} = \arg\max_{\theta} \prod_{i=1}^{n} f\left(y^{(i)} \mid \boldsymbol{x}^{(i)}, \theta\right) = \arg\max_{\theta} LL(\theta)$$

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

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

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

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

$$LL(\theta) = \log \prod_{i=1}^{n} f\left(y^{(i)} \mid \boldsymbol{x}^{(i)}, \theta\right)$$

# 1. Determine objective function (write)

$$\theta_{MLE} = \arg\max_{\theta} \prod_{i=1}^{n} f\left(y^{(i)} \mid \boldsymbol{x}^{(i)}, \theta\right) = \arg\max_{\theta} LL(\theta)$$

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

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

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

Recall
Bernoulli MLE!

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

$$LL(\theta) = \log \prod_{i=1}^{n} f\left(y^{(i)} \mid \boldsymbol{x}^{(i)}, \theta\right)$$

# 1. Determine objective function (write)

$$\theta_{MLE} = \arg \max_{\theta} \prod_{i=1}^{n} f\left(y^{(i)} \mid \boldsymbol{x}^{(i)}, \theta\right) = \arg \max_{\theta} LL(\theta)$$

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

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

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

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

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

# 2. Find gradient with respect to $\theta$

Optimization problem:

$$\theta_{MLE} = \arg\max_\theta \prod_{i=1}^n f\left(y^{(i)} \mid \boldsymbol{x}^{(i)}, \theta\right) = \arg\max_\theta LL(\theta)$$

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

Gradient w.r.t. $\theta_j$, for $j = 0, 1, \ldots, m$:

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

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

# 2. Find gradient with respect to $\theta$

Optimization problem:

$$\theta_{MLE} = \arg \max_{\theta} \prod_{i=1}^{n} f\left(y^{(i)} \mid \boldsymbol{x}^{(i)}, \theta\right) = \arg \max_{\theta} LL(\theta)$$

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

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

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

scale by j-th feature

# 2. Find gradient with respect to $\theta$

Optimization problem:

$$\theta_{MLE} = \arg\max_{\theta} \prod_{i=1}^{n} f\left(y^{(i)} \mid \boldsymbol{x}^{(i)}, \theta\right) = \arg\max_{\theta} LL(\theta)$$

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

Gradient w.r.t. $\theta_j$, for $j = 0, 1, \ldots, m$:

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

1 or 0   $P\left(Y = 1 \mid \boldsymbol{X} = \boldsymbol{x}^{(i)}\right)$

# 2. Find gradient with respect to $\theta$

Optimization problem:

$$\theta_{MLE} = \arg\max_{\theta} \prod_{i=1}^{n} f\left(y^{(i)} \mid \boldsymbol{x}^{(i)}, \theta\right) = \arg\max_{\theta} LL(\theta)$$

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

Gradient w.r.t. $\theta_j$, for $j = 0, 1, \ldots, m$:

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

Suppose $y^{(i)} = 1$ (the true class label for $i$-th datapoint):
- If $\sigma\left(\theta^T \boldsymbol{x}^{(i)}\right) \geq 0.5$, correct
- If $\sigma\left(\theta^T \boldsymbol{x}^{(i)}\right) < 0.5$, incorrect  → change $\theta_j$ more

# 3. Solve

1. Optimization problem:

$$\theta_{MLE} = \arg\max_{\theta} \prod_{i=1}^{n} f\left(y^{(i)} \mid \boldsymbol{x}^{(i)}, \theta\right) = \arg\max_{\theta} LL(\theta)$$

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

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

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

3. Solve

Stay tuned!

# (live)
# 26: Logistic Regression

Lisa Yan and Jerry Cain

November 11, 2020

# Logistic Regression Model

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

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

$$\boldsymbol{X} \Rightarrow \boxed{\theta_0 + \sum_{j=1}^{m} \theta_j X_j} \Rightarrow$$

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



$$\Rightarrow P(Y = 1|\boldsymbol{X})$$

# Introducing notation $\hat{y}$

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

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

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

where $x_0 = 1$

$\boldsymbol{X}$ ⟹ $\theta_0 + \sum_{j=1}^{m} \theta_j X_j$ ⟹ sigmoid function $\sigma(z) = \dfrac{1}{1 + e^{-z}}$ ⟹ $P(Y = 1|\boldsymbol{X}) = \hat{y}$

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

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

# Another view of Logistic Regression

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

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

predict $\hat{Y}$

compute $\hat{y}$

$\theta^T \boldsymbol{x}$

For the "correct" parameters $\theta$:
  $(\boldsymbol{x}, Y = 1)$ should have $\theta^T \boldsymbol{x} > 0$
  $(\boldsymbol{x}, Y = 0)$ should have $\theta^T \boldsymbol{x} \leq 0$

Stanford University

# Today's goals: Logistic Regresison

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

Details
- Gradient ascent pseudocode
- Testing

For the problem set

Philosophy
- Logistic Regression vs Naïve Bayes
- Linearly separable functions

Derivation of gradient (Calculus)

Machine learning insights

# Training: The details

# Training: Learning parameters

**Training**

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

that maximize log conditional likelihood of training data

Some reminders:

- Log conditional likelihood:

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

- Gradient with respect to $\theta$:

$$\frac{\partial LL(\theta)}{\partial \theta_j} = \sum_{i=1}^{n} \left[y^{(i)} - \sigma(\theta^T \boldsymbol{x}^{(i)})\right] x_j^{(i)} \qquad \text{for } j = 0, 1, \ldots, m$$

(derived at end of lecture)

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

# Training: Gradient ascent step

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

```
repeat many times:
    for all thetas:
```

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

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

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

Gradient Ascent Step
for $j = 0, 1, \ldots, m$:
$$\theta_j^{\text{new}} = \theta_j^{\text{old}} + \eta \cdot \sum_{i=1}^{n} \left[ y^{(i)} - \sigma\left( \theta^{\text{old}^T} \boldsymbol{x}^{(i)} \right) \right] x_j^{(i)}$$

```
initialize θ_j = 0 for 0 ≤ j ≤ m
repeat many times:
    gradient[j] = 0 for 0 ≤ j ≤ m
    // TODO: your code here

    // compute all gradient[j]'s
    // based on n training examples

    θ_j += η  * gradient[j] for all 0 ≤ j ≤ m
```

(by yourself)

# Training: Gradient Ascent
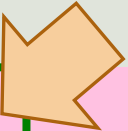
inner loop → for $j = 0, 1, \ldots, m$:

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

compute

outer loop

```
initialize θⱼ = 0 for 0 ≤ j ≤ m
repeat many times:

    gradient[j] = 0 for 0 ≤ j ≤ m

    for each training example (𝒙,𝑦):

        for each 0 ≤ j ≤ m:

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

    θⱼ += η  * gradient[j] for all 0 ≤ j ≤ m
```

# Training: Gradient Ascent

inner loop ➡

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

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

compute

outer loop

```
initialize θⱼ = 0 for 0 ≤ j ≤ m
repeat many times:

    gradient[j] = 0 for 0 ≤ j ≤ m

    for each training example (𝒙,𝑦):

        for each 0 ≤ j ≤ m:
```

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

```
    θⱼ += η * gradient[j] for all 0 ≤ j ≤ m
```

Some important details...

# Training: Gradient Ascent

```
initialize θ_j = 0 for 0 ≤ j ≤ m
repeat many times:
    gradient[j] = 0 for 0 ≤ j ≤ m
    for each training example (x,y):
        for each 0 ≤ j ≤ m:
            gradient[j] +=
```

$$\left[ y - \frac{1}{1 + e^{-\theta^T \boldsymbol{x}}} \right] x_j$$

```
    θ_j += η  * gradient[j] for all 0 ≤ j ≤ m
```

- Finish computing gradient with $\theta^{\mathrm{old}}$ prior to any $\theta$ update

# Training: Gradient Ascent

```
initialize θⱼ = 0 for 0 ≤ j ≤ m
repeat many times:
    gradient[j] = 0 for 0 ≤ j ≤ m
    for each training example (𝒙,y):
        for each 0 ≤ j ≤ m:
```

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

```
    θⱼ += η  * gradient[j] for all 0 ≤ j ≤ m
```

- Finish computing gradient with $\theta^{\text{old}}$ prior to any $\theta$ update
- Learning rate $\eta$ is a constant you set before training

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

```
initialize θ_j = 0 for 0 ≤ j ≤ m
repeat many times:
    gradient[j] = 0 for 0 ≤ j ≤ m
    for each training example (𝒙,𝑦):
        for each 0 ≤ j ≤ m:
```

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

```
    θ_j += η  * gradient[j] for all 0 ≤ j ≤ m
```

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

# Training: Gradient Ascent

```
initialize θ_j = 0 for 0 ≤ j ≤ m
repeat many times:
    gradient[j] = 0 for 0 ≤ j ≤ m
    for each training example (x,y):
        for each 0 ≤ j ≤ m:
```

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

```
    θ_j += η  * gradient[j] for all 0 ≤ j ≤ m
```

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

# Training: Gradient Ascent

```
initialize θⱼ = 0 for 0 ≤ j ≤ m
repeat many times:
    gradient[j] = 0 for 0 ≤ j ≤ m
    for each training example (𝒙,y):
        for each 0 ≤ j ≤ m:
```

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

```
    θⱼ += η * gradient[j] for all 0 ≤ j ≤ m
```

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

# Testing.

# Testing: Classification with Logistic Regression

Training

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

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

Testing

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

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

⚠️ Parameters $\theta_j$ are **not** updated during testing phase

# 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

# Today's goals: Logistic Regression

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

Details
✅
- Gradient ascent pseudocode
- Testing

For the problem set

Philosophy
- Logistic Regression vs Naïve Bayes
- Linearly separable functions

Derivation of gradient (Calculus)

Machine learning insights

# 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

(by yourself) 🤔

# Naïve Bayes        vs        Logistic Regression

$$\theta$$

$X$

$$\hat{P}(\boldsymbol{X}|Y)\hat{P}(Y)$$

$$\hat{P}(\boldsymbol{X}, Y)$$

$Y$

$$\hat{Y} = \underset{y=\{0,1\}}{\arg\max} P(Y \mid \boldsymbol{X}) = \underset{y=\{0,1\}}{\arg\max} P(\boldsymbol{X}|Y)P(Y)$$

$X$

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



$$P(Y = 1|\boldsymbol{X})$$

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

Compare/contrast:

1. What **distributions** are we modeling?

2. After learning our parameters, could we randomly **generate** a new datapoint $(\boldsymbol{x}, y)$?

3. Could we model a **continuous** $X_j$ feature (e.g., $X_j \sim$ Normal, or $X_j \sim$ Unknown)?

4. Could we model a non-binary **discrete** $X_j$ (e.g., $X_j \in \{1, 2, \dots, 6\}$)?

(by yourself)

# Tradeoffs:  Naïve Bayes  Logistic Regression

| | Naïve Bayes | Logistic Regression |
|---|---|---|
| **1.** Modeling goal | $P(\boldsymbol{X}, Y)$ | $P(Y \| \boldsymbol{X})$ |
| **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 $P(\boldsymbol{X}, Y)$) |
| **3.** Continuous input features | ⚠️ Needs parametric form (e.g., Gaussian) or discretized buckets (for multinomial features) | ✅ Yes, easily |
| **4.** 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 $\{1, 2, 3\}$ |

# Linearly separable data

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



$$\theta^T \boldsymbol{x} = 0$$

- 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).

$$\hat{P}(\boldsymbol{X}|Y) = \prod_{j=1}^{m} \hat{P}(X_j|Y)$$

# 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

$$f(x) = f\big(z(x)\big)$$

aka decomposition
of composed functions

# Are you ready?



Quora    Home    Answer    Spaces    Notifications    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

# Our goal

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

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

log conditional likelihood

Two "pre-processing" steps to prepare for chain rule
1. Rewrite $LL(\theta)$ with $\hat{y}$
2. Compute gradient of $\hat{y}$

# 1. Rewriting $LL(\theta)$ with $\hat{y}$

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

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

log conditional likelihood

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

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

# 2. Compute gradient of $\hat{y} = \sigma(\theta^T \boldsymbol{x})$

Aside: Sigmoid has a
beautiful derivative!

Sigmoid function:

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

Derivative:

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

# Think

Slide 72 has code to think over by yourself.

Post any in chat!

Think by yourself: 2 min

(by yourself) 🤔

# 2. Compute gradient of $\hat{y} = \sigma(\theta^T x)$

Sigmoid function:

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

Derivative:

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

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

A. $\sigma(x_j)[1 - \sigma(x_j)]x_j$

B. $\sigma(\theta^T x)[1 - \sigma(\theta^T x)]x$

C. $\sigma(\theta^T x)[1 - \sigma(\theta^T x)]x_j$

D. $\sigma(\theta^T x)x_j[1 - \sigma(\theta^T x)x_j]$

(by yourself)

E. None/other

# 2. Compute gradient of $\hat{y} = \sigma(\theta^T x)$

Sigmoid function:

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

Derivative:

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

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

A. $\sigma(x_j)[1 - \sigma(x_j)]x_j$

B. $\sigma(\theta^T x)[1 - \sigma(\theta^T x)]x$

C. $\sigma(\theta^T x)[1 - \sigma(\theta^T x)]x_j$

D. $\sigma(\theta^T x)x_j[1 - \sigma(\theta^T x)x_j]$

E. None/other

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

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

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

# Compute gradient of log conditional likelihood

$$\frac{\partial LL(\theta)}{\partial \theta_j} = \sum_{i=1}^{n} \frac{\partial}{\partial \theta_j} \left[ y^{(i)} \log(\hat{y}^{(i)}) + (1 - y^{(i)}) \log(1 - \hat{y}^{(i)}) \right] \qquad \text{Let } \hat{y}^{(i)} = \sigma(\theta^T \boldsymbol{x}^{(i)})$$

$$= \sum_{i=1}^{n} \frac{\partial}{\partial \hat{y}^{(i)}} \left[ y^{(i)} \log(\hat{y}^{(i)}) + (1 - y^{(i)}) \log(1 - \hat{y}^{(i)}) \right] \cdot \frac{\partial \hat{y}^{(i)}}{\partial \theta_j} \qquad \text{(Chain Rule)}$$

$$= \sum_{i=1}^{n} \left[ y^{(i)} \frac{1}{\hat{y}^{(i)}} - (1 - y^{(i)}) \frac{1}{1 - \hat{y}^{(i)}} \right] \cdot \hat{y}^{(i)} (1 - \hat{y}^{(i)}) x_j^{(i)} \qquad \text{(calculus)}$$

$$= \sum_{i=1}^{n} \left[ y^{(i)} - \hat{y}^{(i)} \right] x_j^{(i)} \qquad = \sum_{i=1}^{n} \left[ y^{(i)} - \sigma(\theta^T \boldsymbol{x}^{(i)}) \right] x_j^{(i)} \qquad \text{(simplify)}$$

# Compute gradient of log conditional likelihood

$$\frac{\partial LL(\theta)}{\partial \theta_j} = \sum_{i=1}^{n} \frac{\partial}{\partial \theta_j} \left[ y^{(i)} \log(\hat{y}^{(i)}) + (1 - y^{(i)}) \log(1 - \hat{y}^{(i)}) \right] \qquad \text{Let } \hat{y}^{(i)} = \sigma(\theta^T \boldsymbol{x}^{(i)})$$

$$= \sum_{i=1}^{n} \frac{\partial}{\partial \hat{y}^{(i)}} \left[ y^{(i)} \log(\hat{y}^{(i)}) + (1 - y^{(i)}) \log(1 - \hat{y}^{(i)}) \right] \cdot \frac{\partial \hat{y}^{(i)}}{\partial \theta_j} \qquad \text{(Chain Rule)}$$

$$= \sum_{i=1}^{n} \left[ y^{(i)} \frac{1}{\hat{y}^{(i)}} - (1 - y^{(i)}) \frac{1}{1 - \hat{y}^{(i)}} \right] \cdot \hat{y}^{(i)} (1 - \hat{y}^{(i)}) x_j^{(i)} \qquad \text{(calculus)}$$

$$= \sum_{i=1}^{n} \left[ y^{(i)} - \hat{y}^{(i)} \right] x_j^{(i)} \qquad\qquad = \sum_{i=1}^{n} \left[ y^{(i)} - \sigma(\theta^T \boldsymbol{x}^{(i)}) \right] x_j^{(i)} \quad 🎉 \qquad \text{(simplify)}$$

# Interlude for jokes

# Probability as college students



(A useful construct that connects discrete PMF to continuous PDF)