# 10-301/601: Introduction to Machine Learning Lecture 9 – Logistic Regression

Henry Chai & Matt Gormley 9/28/22

#### Q & A:

Where does the  $O(n^3)$  computational complexity for matrix inversion come from?

And that other number you quoted for that matter?

Matrix inversion	One $n  imes n$ matrix	One $n  imes n$ matrix	Gauss-Jordan elimination	$O(n^3)$
			Strassen algorithm	$O(n^{2.807})$
			Coppersmith— Winograd algorithm	$O(n^{2.376})$
			Optimized CW-like algorithms	$O(n^{2.373})$

#### **Front Matter**

- Announcements:
  - HW3 released 9/21, due 9/28 (today!) at 11:59 PM
    - Only two grace days allowed on HW3
    - HW3 exit poll has also been released: you have until one week from the due date to complete it
  - Exam 1 on 10/4 (one week from yesterday) from 6:30 PM 8:30 PM
    - Exam 1 practice problems released on the course website, under <u>Coursework</u>
  - Lecture instead of recitation this Friday (9/30)

### Probabilistic Learning

- Previously:
  - (Unknown) Target function,  $c^*: \mathcal{X} \to \mathcal{Y}$
  - Classifier,  $h: \mathcal{X} \to \mathcal{Y}$
  - Goal: find a classifier, h, that best approximates  $c^*$
- Now:
  - (Unknown) Target distribution,  $y \sim p^*(Y|x)$
  - Distribution, p(Y|x)
  - Goal: find a distribution, p, that best approximates  $p^*$

#### Likelihood

- Given N independent, identically distribution (iid) samples  $\mathcal{D} = \{x^{(1)}, \dots, x^{(N)}\}$  of a random variable X
  - If X is discrete with probability mass function (pmf)  $p(X|\theta)$ , then the *likelihood* of  $\mathcal{D}$  is

$$L(\theta) = \prod_{i=1}^{N} p(x^{(i)}|\theta)$$

• If X is continuous with probability density function (pdf)  $f(X|\theta)$ , then the *likelihood* of  $\mathcal{D}$  is

$$L(\theta) = \prod_{i=1}^{N} f(x^{(i)}|\theta)$$

#### Log-Likelihood

- Given N independent, identically distribution (iid) samples  $\mathcal{D} = \{x^{(1)}, \dots, x^{(N)}\}$  of a random variable X
  - If X is discrete with probability mass function (pmf)  $p(X|\theta)$ , then the log-likelihood of  $\mathcal{D}$  is

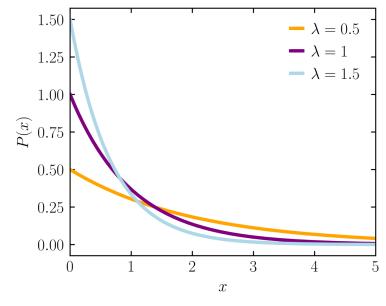
$$\ell(\theta) = \log \prod_{i=1}^{N} p(x^{(i)}|\theta) = \sum_{i=1}^{N} \log p(x^{(i)}|\theta)$$

• If X is continuous with probability density function (pdf)  $f(X|\theta)$ , then the log-likelihood of  $\mathcal{D}$  is

$$\ell(\theta) = \log \prod_{i=1}^{N} f(x^{(i)}|\theta) = \sum_{i=1}^{N} \log f(x^{(i)}|\theta)$$

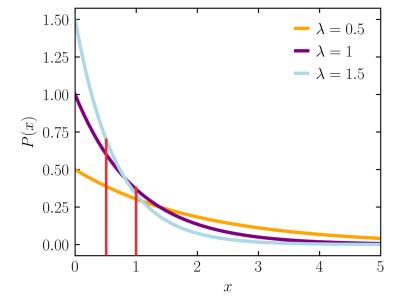
#### Maximum Likelihood Estimation (MLE)

- Insight: every valid probability distribution has a finite amount of probability mass as it must sum/integrate to 1
- Idea: set the parameter(s) so that the likelihood of the samples is maximized
- Intuition: assign as much of the (finite) probability mass to the observed data at the expense of unobserved data
- Example: the exponential distribution



#### Maximum Likelihood Estimation (MLE)

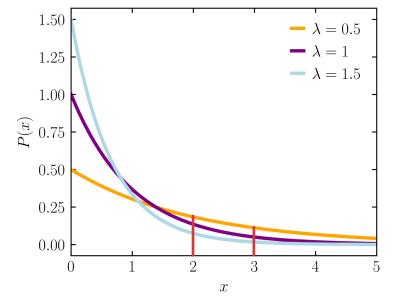
- Insight: every valid probability distribution has a finite amount of probability mass as it must sum/integrate to 1
- Idea: set the parameter(s) so that the likelihood of the samples is maximized
- Intuition: assign as much of the (finite) probability mass to the observed data at the expense of unobserved data
- Example: the exponential distribution



$$\begin{cases} x^{(1)} = 0.5, \\ x^{(2)} = 1 \end{cases}$$

#### Maximum Likelihood Estimation (MLE)

- Insight: every valid probability distribution has a finite amount of probability mass as it must sum/integrate to 1
- Idea: set the parameter(s) so that the likelihood of the samples is maximized
- Intuition: assign as much of the (finite) probability mass to the observed data at the expense of unobserved data
- Example: the exponential distribution



$$\{x^{(1)} = 2 \\ x^{(2)} = 3 \}$$

#### Exponential Distribution MLE

The pdf of the exponential distribution is

$$f(x|\lambda) = \lambda e^{-\lambda x}$$

• Given 
$$N$$
 iid samples  $\{x^{(1)}, \dots, x^{(N)}\}$ , the likelihood is 
$$L(\lambda) = \prod_{i=1}^{N} f(x^{(i)}|\lambda) = \prod_{i=1}^{N} \lambda e^{-\lambda x^{(i)}}$$

### Exponential Distribution MLE

The pdf of the exponential distribution is

$$f(x|\lambda) = \lambda e^{-\lambda x}$$

• Given N iid samples  $\{x^{(1)}, ..., x^{(N)}\}$ , the log-likelihood is

$$\ell(\lambda) = \sum_{i=1}^{N} \log f(x^{(i)}|\lambda) = \sum_{i=1}^{N} \log \lambda e^{-\lambda x^{(i)}}$$

$$= \sum_{i=1}^{N} \log \lambda + \log e^{-\lambda x^{(i)}} = N \log \lambda - \lambda \sum_{i=1}^{N} x^{(i)}$$

#### Okay, but now what?

How do we actually find the best value for the parameter,  $\lambda$ ?

The pdf of the exponential distribution is

$$f(x|\lambda) = \lambda e^{-\lambda x}$$

• Given N iid samples  $\{x^{(1)}, ..., x^{(N)}\}$ , the log-likelihood is

$$\ell(\lambda) = \sum_{i=1}^{N} \log f(x^{(i)}|\lambda) = \sum_{i=1}^{N} \log \lambda e^{-\lambda x^{(i)}}$$

$$= \sum_{i=1}^{N} \log \lambda + \log e^{-\lambda x^{(i)}} = N \log \lambda - \lambda \sum_{i=1}^{N} x^{(i)}$$

### Exponential Distribution MLE

The pdf of the exponential distribution is

$$f(x|\lambda) = \lambda e^{-\lambda x}$$

• Given N iid samples  $\{x^{(1)}, ..., x^{(N)}\}$ , the log-likelihood is

$$\ell(\lambda) = \sum_{i=1}^{N} \log f(x^{(i)}|\lambda) = \sum_{i=1}^{N} \log \lambda e^{-\lambda x^{(i)}}$$

$$= \sum_{i=1}^{N} \log \lambda + \log e^{-\lambda x^{(i)}} = N \log \lambda - \lambda \sum_{i=1}^{N} x^{(i)}$$

Taking the partial derivative and setting it equal to 0 gives

$$\frac{\partial \ell}{\partial \lambda} = \frac{N}{\lambda} - \sum_{i=1}^{N} x^{(i)}$$

### Exponential Distribution MLE

The pdf of the exponential distribution is

$$f(x|\lambda) = \lambda e^{-\lambda x}$$

• Given N iid samples  $\{x^{(1)}, ..., x^{(N)}\}$ , the log-likelihood is

$$\ell(\lambda) = \sum_{i=1}^{N} \log f(x^{(i)}|\lambda) = \sum_{i=1}^{N} \log \lambda e^{-\lambda x^{(i)}}$$

$$= \sum_{i=1}^{N} \log \lambda + \log e^{-\lambda x^{(i)}} = N \log \lambda - \lambda \sum_{i=1}^{N} x^{(i)}$$

Taking the partial derivative and setting it equal to 0 gives

$$\frac{N}{\hat{\lambda}} - \sum_{i=1}^{N} x^{(i)} = 0 \to \frac{N}{\hat{\lambda}} = \sum_{i=1}^{N} x^{(i)} \to \hat{\lambda} = \frac{N}{\sum_{i=1}^{N} x^{(i)}}$$

#### Building a Probabilistic Classifier

- Define a decision rule
  - Given a test data point x', predict its label  $\hat{y}$  using the posterior distribution P(Y = y | x')
  - Common choice:  $\hat{y} = \underset{y}{\operatorname{argmax}} P(Y = y | x')$

• Idea: model P(Y|x) as some parametric function of x

#### Modelling the Posterior

• Suppose we have binary labels  $y \in \{0,1\}$  and *D*-dimensional inputs  $\mathbf{x} = [1, x_1, ..., x_D]^T \in \mathbb{R}^{D+1}$ • 1 prepended to **x** 

Assume

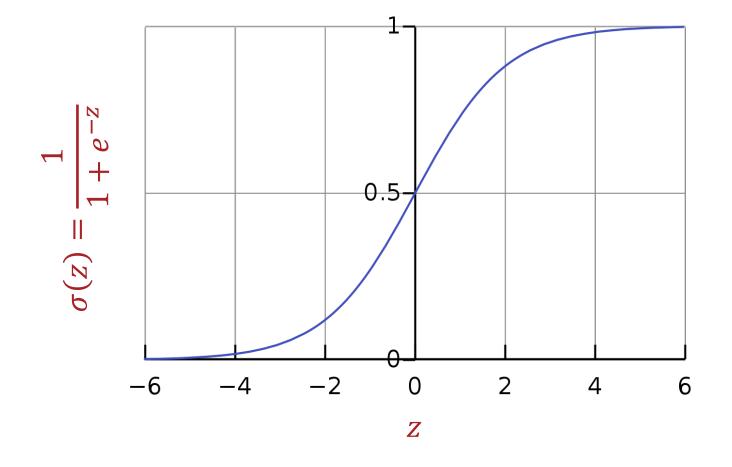
$$P(Y = 1 | \boldsymbol{x}, \boldsymbol{\theta}) = \sigma(\boldsymbol{\theta}^T \boldsymbol{x}) = \frac{1}{1 + \exp(-\boldsymbol{\theta}^T \boldsymbol{x})} = \frac{\exp(\boldsymbol{\theta}^T \boldsymbol{x})}{\exp(\boldsymbol{\theta}^T \boldsymbol{x}) + 1}$$

This implies two useful facts:

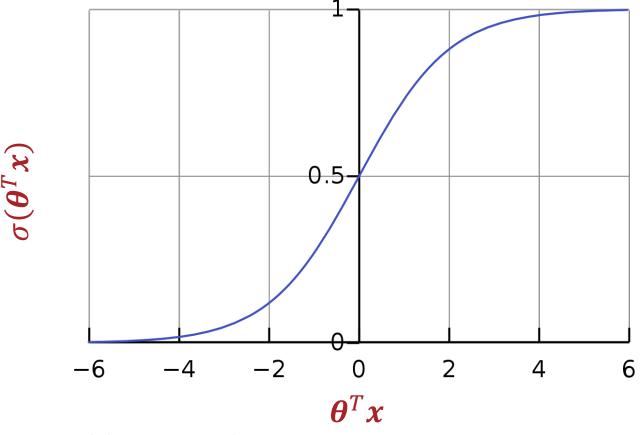
1. 
$$P(Y = 0 | \mathbf{x}, \boldsymbol{\theta}) = 1 - P(Y = 1 | \mathbf{x}, \boldsymbol{\theta}) = \frac{1}{\exp(\boldsymbol{\theta}^T \mathbf{x}) + 1}$$

2. 
$$\frac{P(Y=1|\mathbf{x},\boldsymbol{\theta})}{P(Y=0|\mathbf{x},\boldsymbol{\theta})} = \exp(\boldsymbol{\theta}^T \mathbf{x}) \to \log \frac{P(Y=1|\mathbf{x},\boldsymbol{\theta})}{P(Y=0|\mathbf{x},\boldsymbol{\theta})} = \boldsymbol{\theta}^T \mathbf{x}$$

### Logistic Function



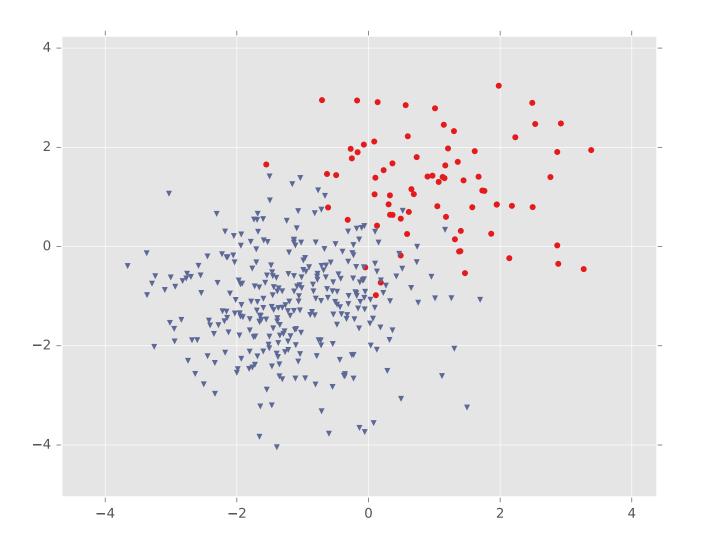
### Why use the Logistic Function?



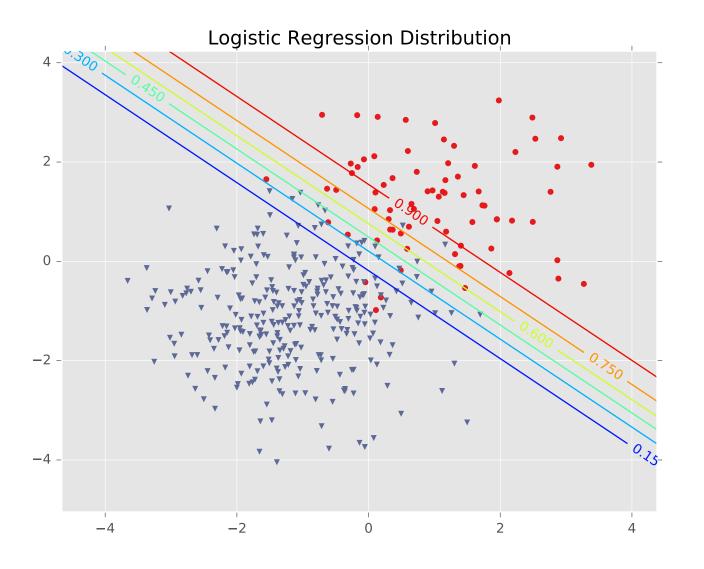
- Differentiable everywhere
- $\sigma$ :  $\mathbb{R} \rightarrow [0, 1]$
- The decision boundary is linear in x!

$$\hat{y} = \begin{cases} 1 \text{ if } P(Y = 1 | \mathbf{x}, \boldsymbol{\theta}) \ge \frac{1}{2} \\ 0 \text{ otherwise.} \end{cases}$$

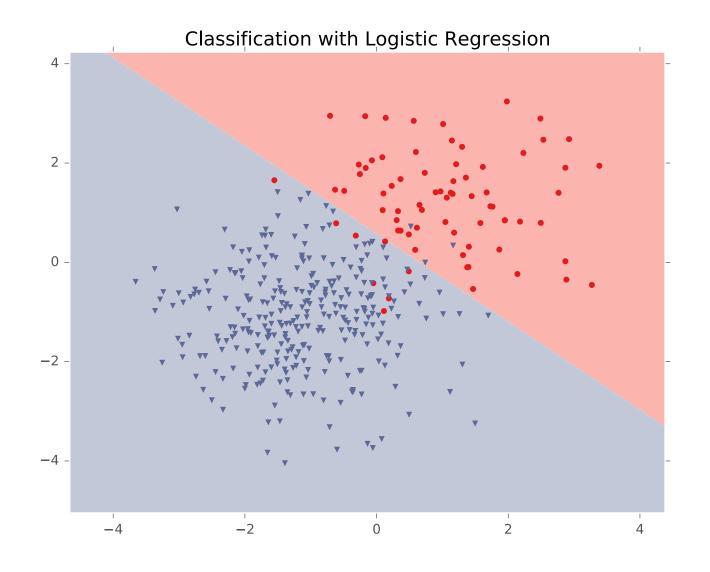
erwise. 
$$P(Y = 1|x) = \sigma(\boldsymbol{\theta}^T x) = \frac{1}{1 + \exp(-\boldsymbol{\theta}^T x)} \ge \frac{1}{2}$$
$$2 \ge 1 + \exp(-\boldsymbol{\theta}^T x)$$
$$1 \ge \exp(-\boldsymbol{\theta}^T x)$$
$$\log(1) \ge -\boldsymbol{\theta}^T x$$
$$0 \le \boldsymbol{\theta}^T x$$



20



21



#### Setting the **Parameters** via Minimum Negative Conditional (log-)Likelihood **Estimation** (MCLE)

#### Find $\theta$ that minimizes

$$\ell(\boldsymbol{\theta}) = -\log P(y^{(1)}, ..., y^{(N)} | \boldsymbol{x}^{(1)}, ..., \boldsymbol{x}^{(N)}, \boldsymbol{\theta}) = -\log \prod_{i=1}^{N} P(y^{(i)} | \boldsymbol{x}^{(i)}, \boldsymbol{\theta})$$

$$= -\log \prod_{i=1}^{N} P(Y = 1 | \boldsymbol{x}^{(i)}, \boldsymbol{\theta})^{y^{(i)}} \left( P(Y = 0 | \boldsymbol{x}^{(i)}, \boldsymbol{\theta}) \right)^{1-y^{(i)}}$$

$$= -\sum_{i=1}^{N} y^{(i)} \log P(Y = 1 | \boldsymbol{x}^{(i)}, \boldsymbol{\theta}) + (1 - y^{(i)}) \log P(Y = 0 | \boldsymbol{x}^{(i)}, \boldsymbol{\theta})$$

$$= -\sum_{i=1}^{N} y^{(i)} \log \frac{P(Y = 1 | \boldsymbol{x}^{(i)}, \boldsymbol{\theta})}{P(Y = 0 | \boldsymbol{x}^{(i)}, \boldsymbol{\theta})} + \log P(Y = 0 | \boldsymbol{x}^{(i)}, \boldsymbol{\theta})$$

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

$$J(\boldsymbol{\theta}) = \frac{1}{N} \ell(\boldsymbol{\theta}) = -\frac{1}{N} \sum_{i=1}^{N} y^{(i)} \boldsymbol{\theta}^{T} \boldsymbol{x}^{(i)} - \log \left( 1 + \exp(\boldsymbol{\theta}^{T} \boldsymbol{x}^{(i)}) \right)$$

## Minimizing the Negative Conditional (log-)Likelihood

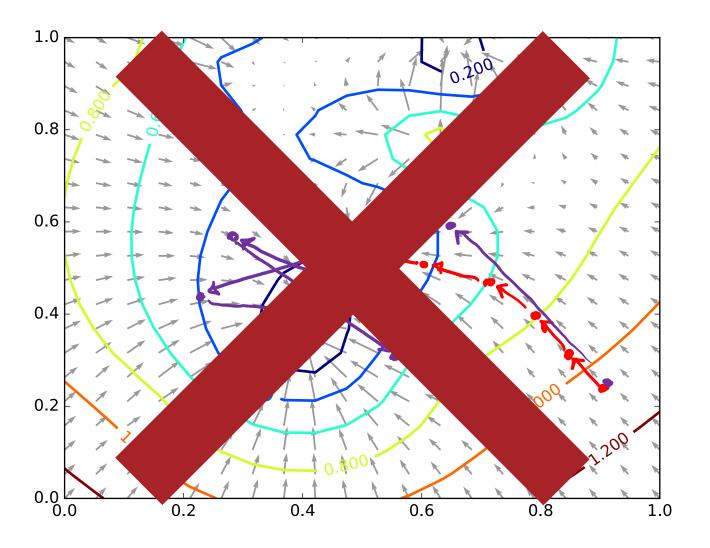
$$J(\boldsymbol{\theta}) = -\frac{1}{N} \sum_{i=1}^{N} y^{(i)} \boldsymbol{\theta}^{T} \boldsymbol{x}^{(i)} - \log \left( 1 + \exp(\boldsymbol{\theta}^{T} \boldsymbol{x}^{(i)}) \right)$$

$$\nabla_{\boldsymbol{\theta}} J(\boldsymbol{\theta}) = -\frac{1}{N} \sum_{i=1}^{N} y^{(i)} \nabla_{\boldsymbol{\theta}} (\boldsymbol{\theta}^{T} \boldsymbol{x}^{(i)}) - \nabla_{\boldsymbol{\theta}} \log \left( 1 + \exp(\boldsymbol{\theta}^{T} \boldsymbol{x}^{(i)}) \right)$$

$$= -\frac{1}{N} \sum_{i=1}^{N} y^{(i)} \boldsymbol{x}^{(i)} - \frac{\exp(\boldsymbol{\theta}^{T} \boldsymbol{x}^{(i)})}{1 + \exp(\boldsymbol{\theta}^{T} \boldsymbol{x}^{(i)})} \boldsymbol{x}^{(i)}$$

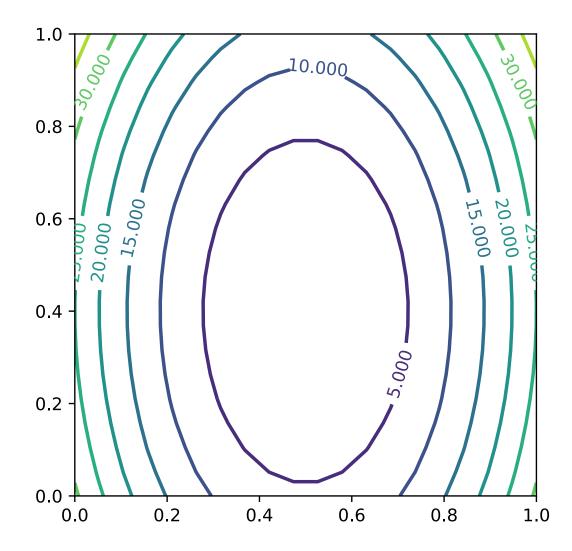
$$= \frac{1}{N} \sum_{i=1}^{N} x^{(i)} (P(Y = 1 | x^{(i)}, \theta) - y^{(i)})$$

### Recall: Gradient Descent



Good news: the negative conditional log-likelihood is convex! (See HW/recitation)

### Recall: Gradient Descent



Good news: the negative conditional log-likelihood is convex! (See HW/recitation)

#### Gradient Descent

- Input: training dataset  $\mathcal{D} = \left\{ \left( x^{(i)}, y^{(i)} \right) \right\}_{i=1}^N$  and step size  $\gamma$
- 1. Initialize  $\boldsymbol{\theta}^{(0)}$  to all zeros and set t=0
- 2. While TERMINATION CRITERION is not satisfied
  - a. Compute the gradient:

$$\nabla_{\boldsymbol{\theta}} J(\boldsymbol{\theta}^{(t)}) = \frac{1}{N} \sum_{i=1}^{N} \boldsymbol{x}^{(i)} (P(Y = 1 | \boldsymbol{x}^{(i)}, \boldsymbol{\theta}^{(t)}) - \boldsymbol{y}^{(i)})$$

- b. Update  $\boldsymbol{\theta}$ :  $\boldsymbol{\theta}^{(t+1)} \leftarrow \boldsymbol{\theta}^{(t)} \gamma \nabla_{\boldsymbol{\theta}} J(\boldsymbol{\theta}^{(t)})$
- c. Increment  $t: t \leftarrow t + 1$
- Output:  $\boldsymbol{\theta}^{(t)}$

#### Poll Question 1:

What is the computational cost of one iteration of gradient descent for logistic regression?

- A. O(1) (TOXIC) B. O(N) C. O(D) D. O(ND)

- Input: training dataset  $\mathcal{D} = \left\{ \left( x^{(i)}, y^{(i)} \right) \right\}_{i=1}^N$  and step size  $\gamma$
- 1. Initialize  $\theta^{(0)}$  to all zeros and set t=0
- While TERMINATION CRITERION is not satisfied
  - a. Compute the gradient:

$$\nabla_{\boldsymbol{\theta}} J(\boldsymbol{\theta}^{(t)}) = \frac{1}{N} \sum_{i=1}^{N} \boldsymbol{x}^{(i)} (P(Y = 1 | \boldsymbol{x}^{(i)}, \boldsymbol{\theta}^{(t)}) - y^{(i)})$$

- b. Update  $\theta: \theta^{(t+1)} \leftarrow \theta^{(t)} \gamma \nabla_{\theta} I(\theta^{(t)})$
- c. Increment  $t: t \leftarrow t+1$
- Output:  $\boldsymbol{\theta}^{(t)}$

#### Gradient Descent

- Input: training dataset  $\mathcal{D} = \left\{ \left( x^{(i)}, y^{(i)} \right) \right\}_{i=1}^N$  and step size  $\gamma$
- 1. Initialize  $\boldsymbol{\theta}^{(0)}$  to all zeros and set t=0
- 2. While TERMINATION CRITERION is not satisfied
  - a. Compute the gradient:

$$O(ND) \left\{ \nabla_{\boldsymbol{\theta}} J(\boldsymbol{\theta}^{(t)}) = \frac{1}{N} \sum_{i=1}^{N} \boldsymbol{x}^{(i)} (P(Y = 1 | \boldsymbol{x}^{(i)}, \boldsymbol{\theta}^{(t)}) - \boldsymbol{y}^{(i)}) \right\}$$

- b. Update  $\boldsymbol{\theta}$ :  $\boldsymbol{\theta}^{(t+1)} \leftarrow \boldsymbol{\theta}^{(t)} \gamma \nabla_{\boldsymbol{\theta}} J(\boldsymbol{\theta}^{(t)})$
- c. Increment  $t: t \leftarrow t + 1$
- Output:  $\boldsymbol{\theta}^{(t)}$

#### Stochastic Gradient Descent (SGD)

- Input: training dataset  $\mathcal{D} = \left\{ \left( x^{(i)}, y^{(i)} \right) \right\}_{i=1}^N$  and step size  $\gamma$
- 1. Initialize  $\theta^{(0)}$  to all zeros and set t=0
- While TERMINATION CRITERION is not satisfied
  - a. Randomly sample a data point from  $\mathcal{D}$ ,  $(x^{(i)}, y^{(i)})$
  - b. Compute the pointwise gradient:

$$\nabla_{\boldsymbol{\theta}} J^{(i)}(\boldsymbol{\theta}^{(t)}) = \boldsymbol{x}^{(i)}(P(Y=1|\boldsymbol{x}^{(i)},\boldsymbol{\theta}^{(t)}) - y^{(i)})$$

- c. Update  $\boldsymbol{\theta}$ :  $\boldsymbol{\theta}^{(t+1)} \leftarrow \boldsymbol{\theta}^{(t)} \gamma \nabla_{\boldsymbol{\theta}} J^{(i)} (\boldsymbol{\theta}^{(t)})$
- d. Increment  $t: t \leftarrow t + 1$
- Output:  $\boldsymbol{\theta}^{(t)}$

#### Stochastic Gradient Descent (SGD)

• If the example is sampled uniformly at random, the expected value of the pointwise gradient is the same as the full gradient!

$$E[\nabla_{\boldsymbol{\theta}} J^{(i)}(\boldsymbol{\theta})] = \sum_{i=1}^{N} \left(\text{probability of selecting } \boldsymbol{x}^{(i)}, y^{(i)}\right) \nabla_{\boldsymbol{\theta}} J^{(i)}(\boldsymbol{\theta})$$
$$= \sum_{i=1}^{N} \left(\frac{1}{N}\right) \nabla_{\boldsymbol{\theta}} J^{(i)}(\boldsymbol{\theta}) = \frac{1}{N} \sum_{i=1}^{N} \nabla_{\boldsymbol{\theta}} J^{(i)}(\boldsymbol{\theta}) = \nabla_{\boldsymbol{\theta}} J(\boldsymbol{\theta})$$

 In practice, the data set is randomly shuffled then looped through so that each data point is used equally often

#### Stochastic Gradient Descent (SGD)

- Input: training dataset  $\mathcal{D} = \left\{ \left( x^{(i)}, y^{(i)} \right) \right\}_{i=1}^N$  and step size  $\gamma$
- 1. Initialize  $\boldsymbol{\theta}^{(0)}$  to all zeros and set t=0
- 2. While TERMINATION CRITERION is not satisfied
  - a. For  $i \in \text{shuffle}(\{1, ..., N\})$ 
    - Compute the pointwise gradient:

$$\nabla_{\boldsymbol{\theta}} J^{(i)}(\boldsymbol{\theta}^{(t)}) = \boldsymbol{x}^{(i)}(P(Y=1|\boldsymbol{x}^{(i)},\boldsymbol{\theta}^{(t)}) - y^{(i)})$$

- ii. Update  $\boldsymbol{\theta}$ :  $\boldsymbol{\theta}^{(t+1)} \leftarrow \boldsymbol{\theta}^{(t)} \gamma \nabla_{\boldsymbol{\theta}} J^{(i)} (\boldsymbol{\theta}^{(t)})$
- iii. Increment  $t: t \leftarrow t + 1$
- Output:  $\boldsymbol{\theta}^{(t)}$

### Wait, have we seen something like this before?

- Input: training dataset  $\mathcal{D} = \left\{ \left( x^{(i)}, y^{(i)} \right) \right\}_{i=1}^N$  and step size  $\gamma$
- 1. Initialize  $\boldsymbol{\theta}^{(0)}$  to all zeros and set t=0
- 2. While TERMINATION CRITERION is not satisfied
  - a. For  $i \in \text{shuffle}(\{1, ..., N\})$ 
    - i. Compute the pointwise gradient:

$$\nabla_{\boldsymbol{\theta}} J^{(i)}(\boldsymbol{\theta}^{(t)}) = \boldsymbol{x}^{(i)}(P(Y=1|\boldsymbol{x}^{(i)},\boldsymbol{\theta}^{(t)}) - y^{(i)})$$

- ii. Update  $\boldsymbol{\theta}$ :  $\boldsymbol{\theta}^{(t+1)} \leftarrow \boldsymbol{\theta}^{(t)} \gamma \nabla_{\boldsymbol{\theta}} J^{(i)} (\boldsymbol{\theta}^{(t)})$
- iii. Increment  $t: t \leftarrow t + 1$
- Output:  $\boldsymbol{\theta}^{(t)}$

## Recall: Perceptron Learning Algorithm

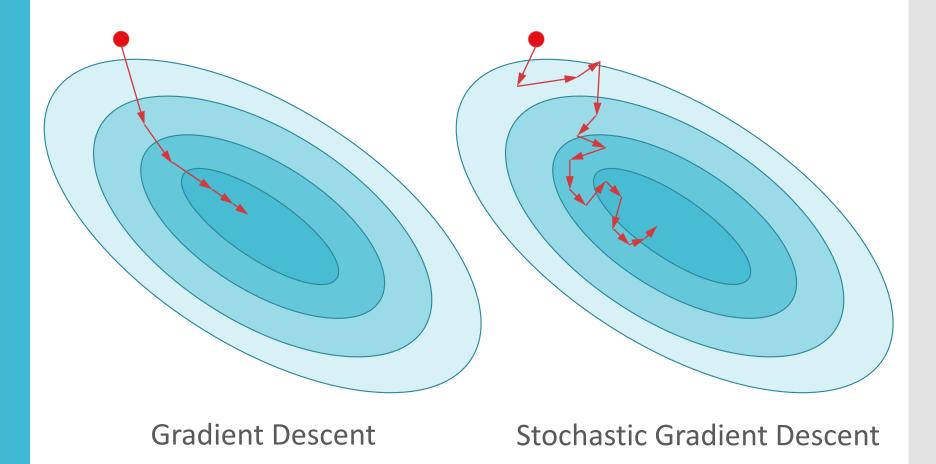
- Input: training dataset  $\mathcal{D} = \{(x^{(i)}, y^{(i)})\}_{i=1}^{N}$
- 1. Initialize  $\boldsymbol{\theta}^{(0)}$  to all zeros and set t=0
- While NOT CONVERGED
  - a. For  $t \in \text{shuffle}(\{1, ..., N\})$ 
    - i. Predict the label of  $x^{(t)}$ ,  $\hat{y} = \text{sign}(\theta^T x^{(t)})$
    - ii. Observe its true label,  $y^{(t)}$
    - iii. If we misclassified  $x^{(t)}$  ( $y^{(t)} \neq \hat{y}$ ):

$$\boldsymbol{\theta} \leftarrow \boldsymbol{\theta} + y^{(t)} \boldsymbol{x}^{(t)}$$

This is the negative gradient of the *hinge loss* 

$$J^{(i)}(\boldsymbol{\theta}) = \max(0, 1 - y^{(i)}\boldsymbol{\theta}^T \boldsymbol{x}^{(i)})$$

Stochastic
Gradient
Descent vs.
Gradient
Descent



## Stochastic Gradient Descent vs. Gradient Descent

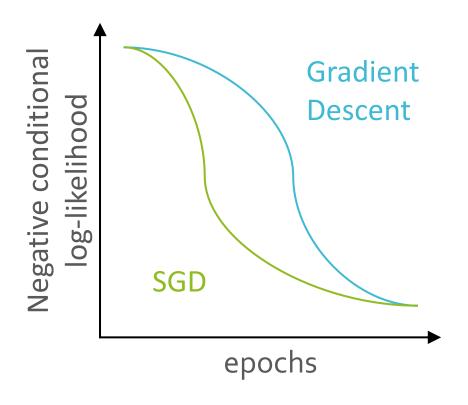
- An *epoch* is a single pass through the entire training dataset
  - Gradient descent updates the parameters once per epoch
  - SGD updates the parameters N times per epoch
- Theoretical comparison:
  - Define convergence to be when  $J(\boldsymbol{\theta^{(t)}}) J(\boldsymbol{\theta^*}) < \epsilon$

Method	Steps to Convergence	Computation per Step		
Gradient descent	$O(\log 1/\epsilon)$	O(ND)		
SGD	$O(1/\epsilon)$	O(D)		

(with high probability under certain assumptions)

## Stochastic Gradient Descent vs. Gradient Descent

- An *epoch* is a single pass through the entire training dataset
  - Gradient descent updates the parameters once per epoch
  - SGD updates the parameters N times per epoch



Empirically, SGD reduces the negative conditional log-likelihood much faster than gradient descent

9/28/22 **37** 

### Optimization for ML Learning Objectives

You should be able to...

- Apply gradient descent to optimize a function
- Apply stochastic gradient descent (SGD) to optimize a function
- Apply knowledge of zero derivatives to identify a closed-form solution (if one exists) to an optimization problem
- Distinguish between convex, concave, and nonconvex functions
- Obtain the gradient (and Hessian) of a (twice) differentiable function

## Logistic Regression Learning Objectives

You should be able to...

- Apply the principle of maximum likelihood estimation (MLE) to learn the parameters of a probabilistic model
- Given a discriminative probabilistic model, derive the conditional log-likelihood, its gradient, and the corresponding Bayes Classifier
- Explain the practical reasons why we work with the log of the likelihood
- Implement logistic regression for binary (and multiclass) classification
- Prove that the decision boundary of binary logistic regression is linear