

# 10-301/601: Introduction to Machine Learning

## Lecture 9 – Logistic Regression

Hoda Heidari, Henry Chai & Matt Gormley

2/14/24

# Front Matter

- Announcements:
  - Exam 1 on 2/19 from 7 PM – 9 PM
  - Exam 1 practice problems released on the course website, under [Coursework](#)

# Probabilistic Learning

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

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

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

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

# 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

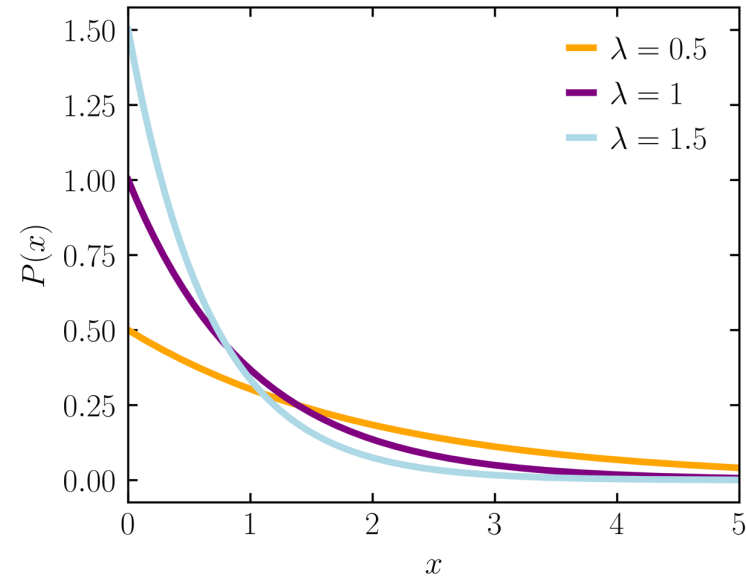
$$\ell(\theta) = \log \prod_{n=1}^N p(x^{(n)}|\theta) = \sum_{n=1}^N \log p(x^{(n)}|\theta)$$

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

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

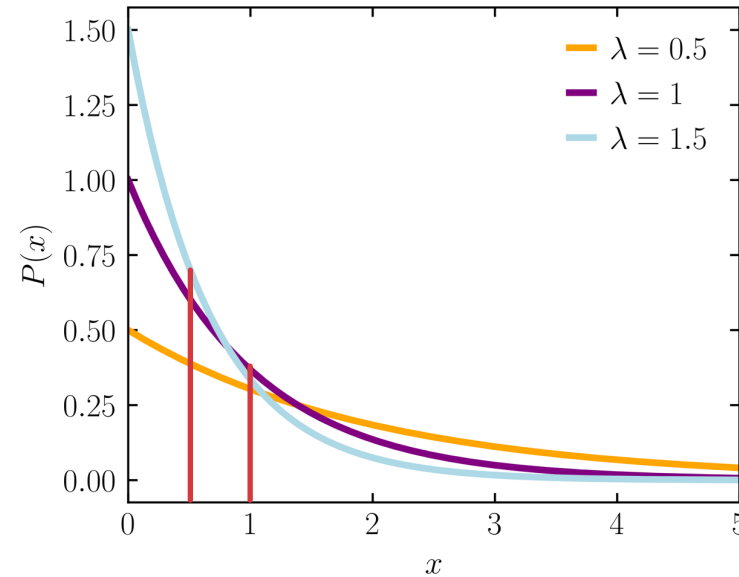
# Maximum Likelihood Estimation (MLE)

- 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)

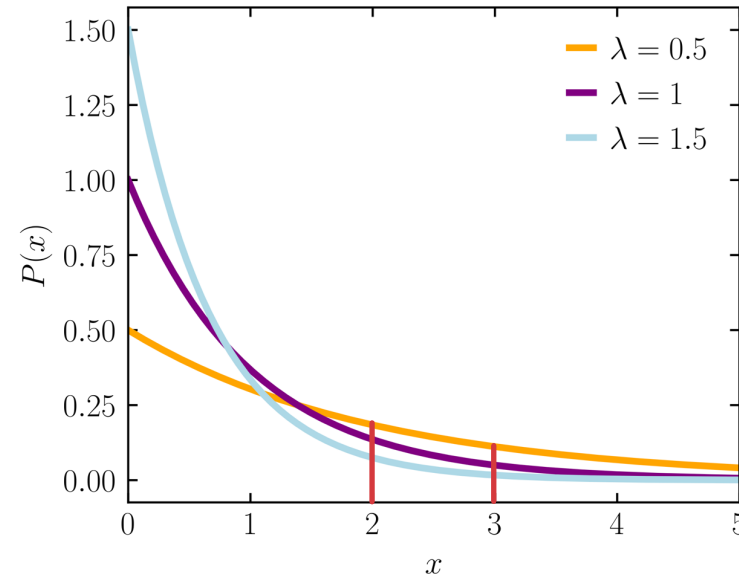
- Insight: every valid probability distribution has a finite amount of probability mass as it must sum/integrate to 1
- Idea: set the parameter(s) so that the likelihood of the samples is maximized
- Intuition: assign as much of the (finite) probability mass to the observed data *at the expense of unobserved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



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

# 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

$$L(\lambda) = \prod_{n=1}^N f(x^{(n)}|\lambda) = \prod_{n=1}^N \lambda e^{-\lambda x^{(n)}}$$

# 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 log-likelihood is

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

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

- Taking the partial derivative and setting it equal to 0 gives

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

# 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 log-likelihood is

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

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

- Taking the partial derivative and setting it equal to 0 gives

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

# Building a Probabilistic Classifier

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

# Modelling the Posterior

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

- **Assume**

1 prepended to  $\mathbf{x}$

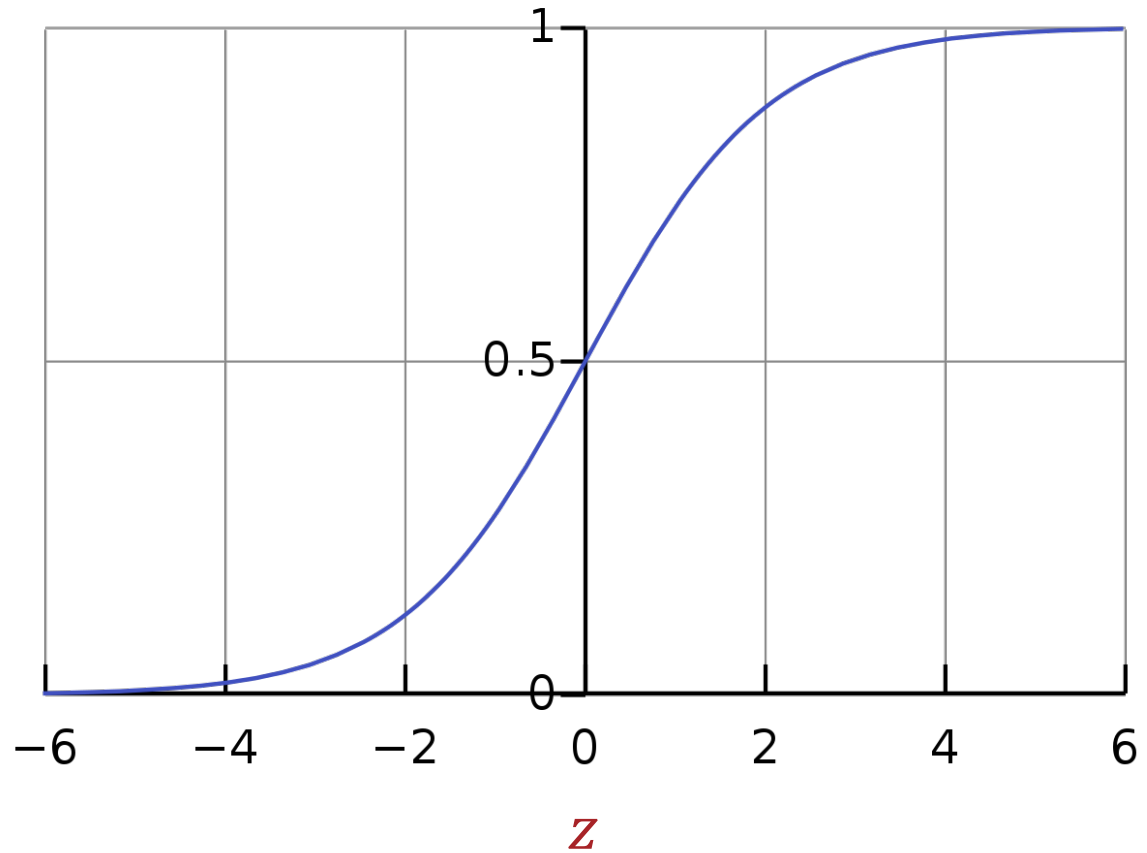
$$P(Y = 1|\mathbf{x}, \boldsymbol{\theta}) = \sigma(\boldsymbol{\theta}^T \mathbf{x}) = \frac{1}{1 + \exp(-\boldsymbol{\theta}^T \mathbf{x})} = \frac{\exp(\boldsymbol{\theta}^T \mathbf{x})}{\exp(\boldsymbol{\theta}^T \mathbf{x}) + 1}$$

- This implies two useful facts:

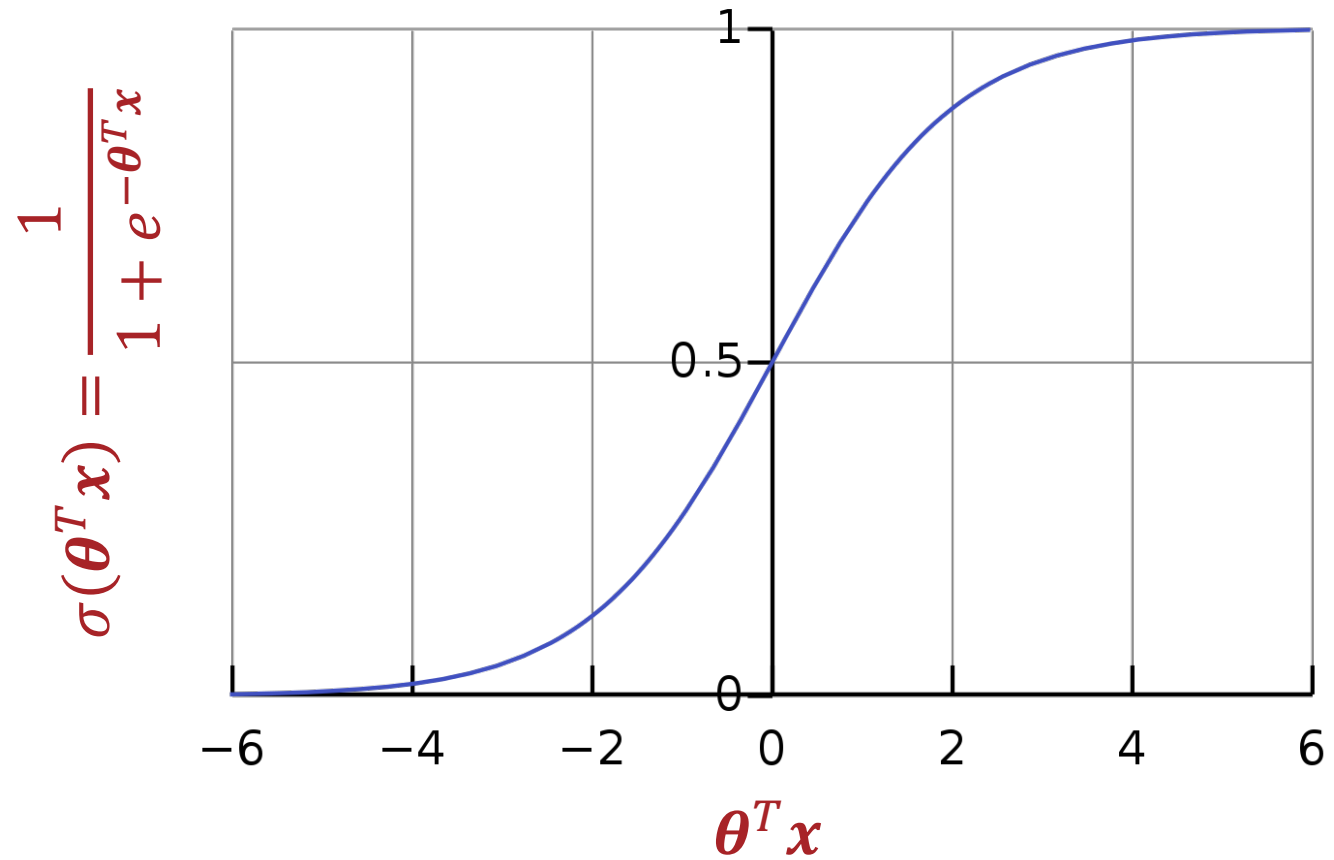
1.  $P(Y = 0|\mathbf{x}, \boldsymbol{\theta}) = 1 - P(Y = 1|\mathbf{x}, \boldsymbol{\theta}) = \frac{1}{\exp(\boldsymbol{\theta}^T \mathbf{x}) + 1}$
2.  $\frac{P(Y = 1|\mathbf{x}, \boldsymbol{\theta})}{P(Y = 0|\mathbf{x}, \boldsymbol{\theta})} = \exp(\boldsymbol{\theta}^T \mathbf{x}) \rightarrow \log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

# Logistic Function

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



# Why use the Logistic Function?



# Logistic Regression Decision Boundary

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

$$P(Y = 1|\mathbf{x}) = \sigma(\boldsymbol{\theta}^T \mathbf{x}) = \frac{1}{1 + \exp(-\boldsymbol{\theta}^T \mathbf{x})} \geq \frac{1}{2}$$

$$2 \geq 1 + \exp(-\boldsymbol{\theta}^T \mathbf{x})$$

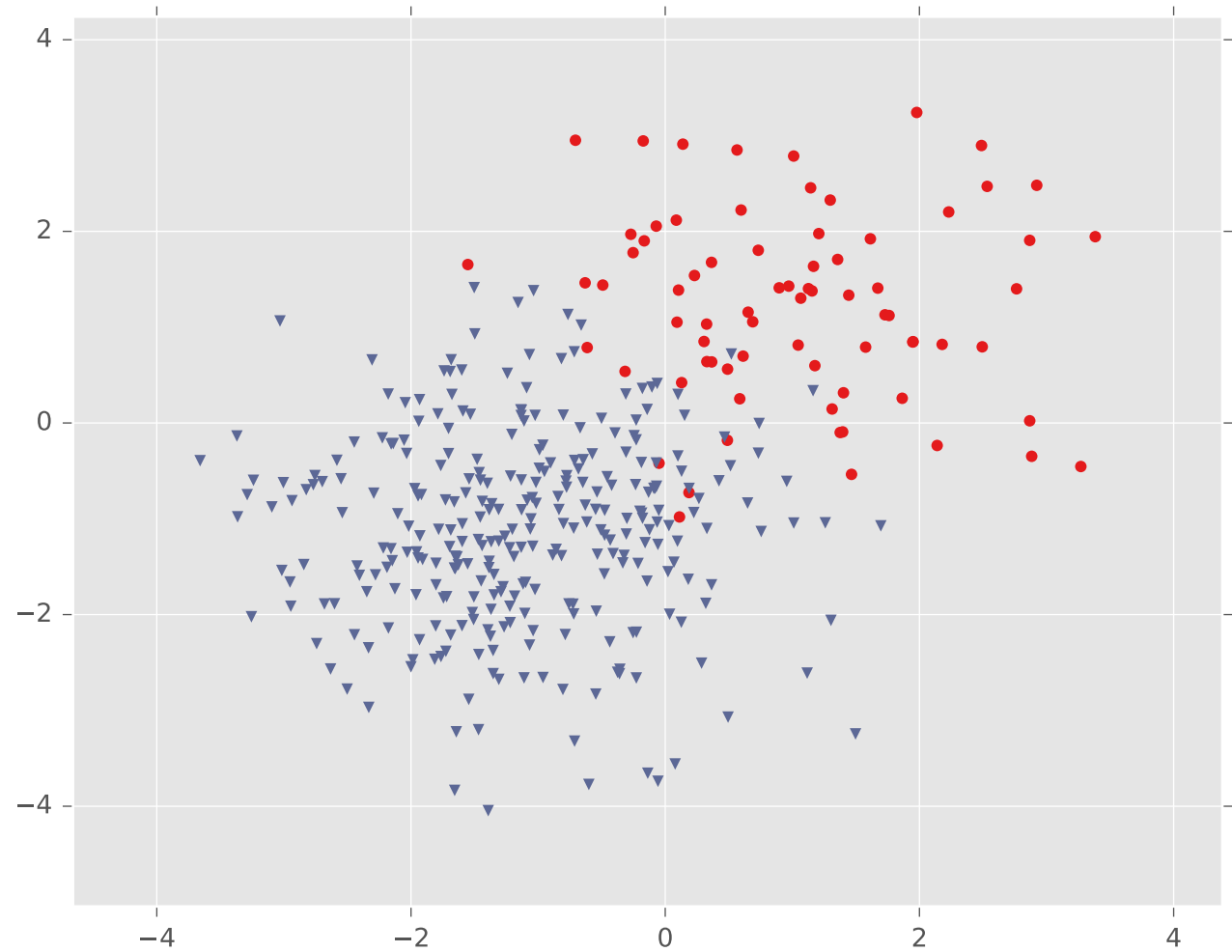
$$1 \geq \exp(-\boldsymbol{\theta}^T \mathbf{x})$$

$$\log(1) \geq -\boldsymbol{\theta}^T \mathbf{x}$$

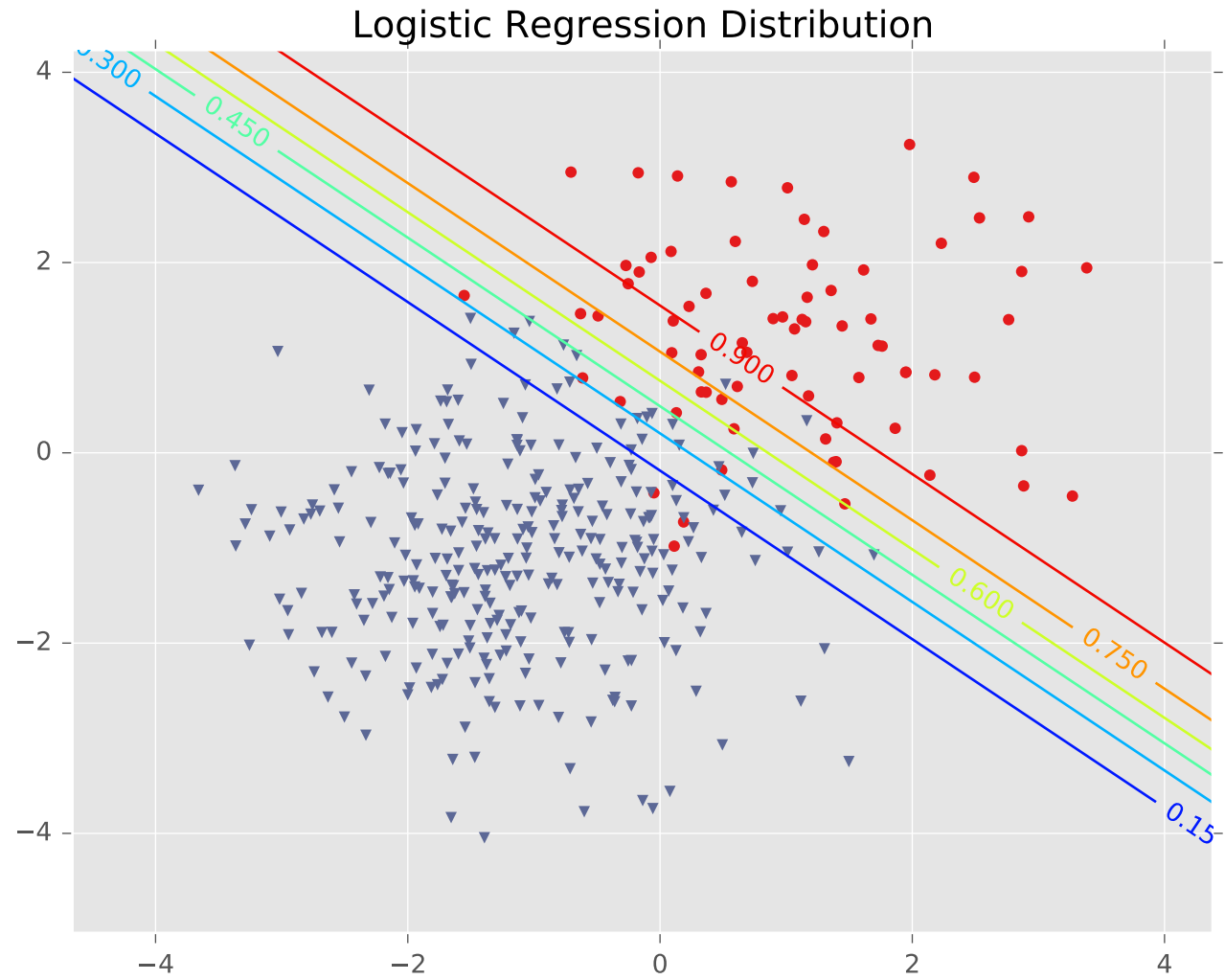
$$0 \leq \boldsymbol{\theta}^T \mathbf{x}$$



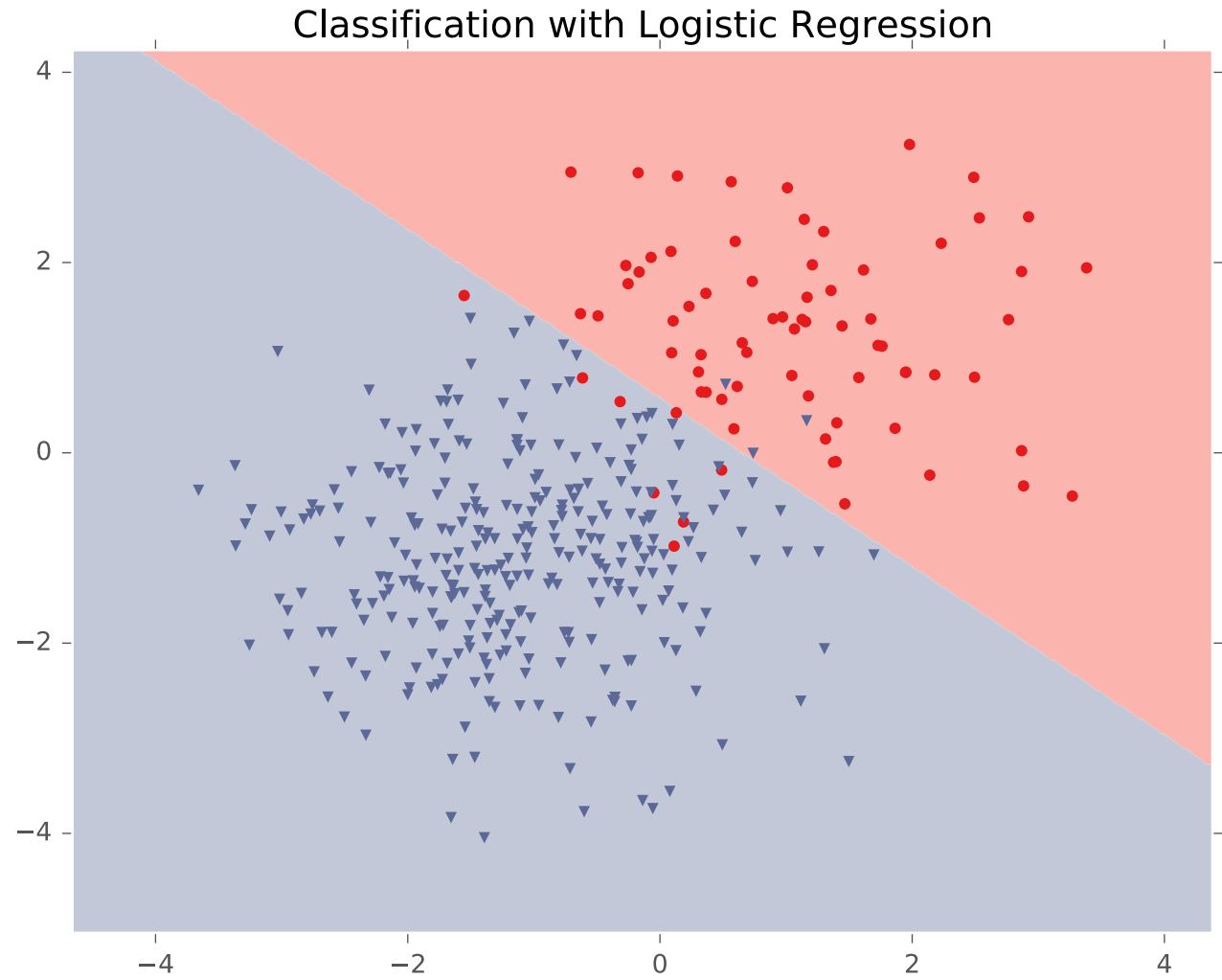
# Logistic Regression Decision Boundary



# Logistic Regression Decision Boundary



# Logistic Regression Decision Boundary



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

- Find  $\boldsymbol{\theta}$  that minimizes

$$\begin{aligned}\ell(\boldsymbol{\theta}) &= -\log P(y^{(1)}, \dots, y^{(N)} | \mathbf{x}^{(1)}, \dots, \mathbf{x}^{(N)}, \boldsymbol{\theta}) = -\log \prod_{n=1}^N P(y^{(n)} | \mathbf{x}^{(n)}, \boldsymbol{\theta}) \\ &= -\log \prod_{n=1}^N P(Y = 1 | \mathbf{x}^{(n)}, \boldsymbol{\theta})^{y^{(n)}} \left( P(Y = 0 | \mathbf{x}^{(n)}, \boldsymbol{\theta}) \right)^{1-y^{(n)}} \\ &= -\sum_{n=1}^N y^{(n)} \log P(Y = 1 | \mathbf{x}^{(n)}, \boldsymbol{\theta}) + (1 - y^{(n)}) \log P(Y = 0 | \mathbf{x}^{(n)}, \boldsymbol{\theta}) \\ &= -\sum_{n=1}^N y^{(n)} \log \frac{P(Y = 1 | \mathbf{x}^{(n)}, \boldsymbol{\theta})}{P(Y = 0 | \mathbf{x}^{(n)}, \boldsymbol{\theta})} + \log P(Y = 0 | \mathbf{x}^{(n)}, \boldsymbol{\theta}) \\ &= -\sum_{n=1}^N y^{(n)} \boldsymbol{\theta}^T \mathbf{x}^{(n)} - \log \left( 1 + \exp(\boldsymbol{\theta}^T \mathbf{x}^{(n)}) \right)\end{aligned}$$

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

## Minimizing the Negative Conditional (log-)Likelihood

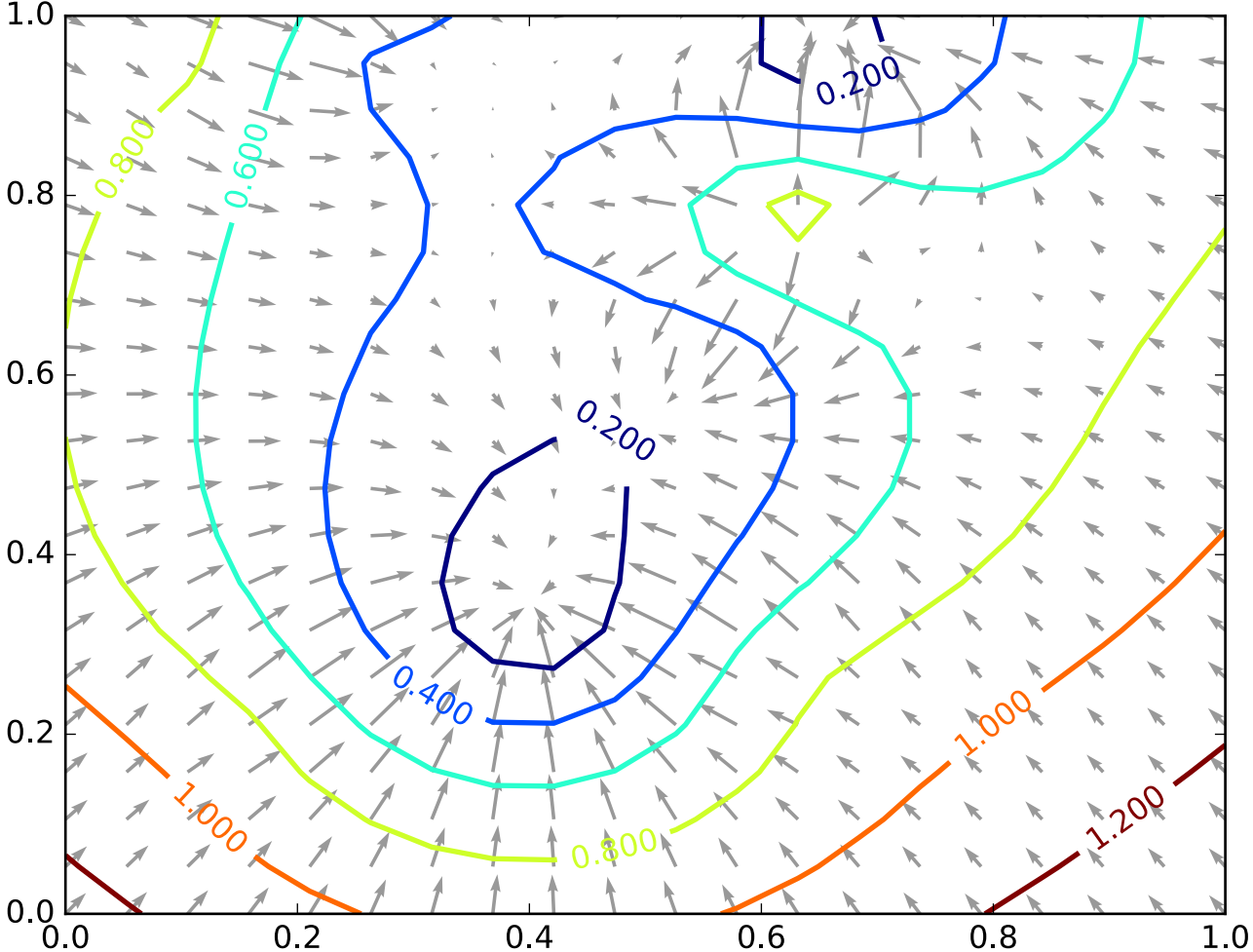
$$J(\boldsymbol{\theta}) = -\frac{1}{N} \sum_{n=1}^N y^{(n)} \boldsymbol{\theta}^T \mathbf{x}^{(n)} - \log \left( 1 + \exp(\boldsymbol{\theta}^T \mathbf{x}^{(n)}) \right)$$

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

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

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

# Recall: Gradient Descent



# Gradient Descent

- Input: training dataset  $\mathcal{D} = \{(\mathbf{x}^{(i)}, y^{(i)})\}_{i=1}^N$  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:
$$\nabla_{\boldsymbol{\theta}} J(\boldsymbol{\theta}^{(t)}) = \frac{1}{N} \sum_{i=1}^N \mathbf{x}^{(i)} (P(Y = 1 | \mathbf{x}^{(i)}, \boldsymbol{\theta}^{(t)}) - y^{(i)})$$
    - 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  $\mathcal{D} = \{(\mathbf{x}^{(i)}, y^{(i)})\}_{i=1}^N$  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:

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

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:  $\boldsymbol{\theta}^{(t)}$



# Gradient Descent

- Input: training dataset  $\mathcal{D} = \{(\mathbf{x}^{(i)}, y^{(i)})\}_{i=1}^N$  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:

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

- 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  $\mathcal{D} = \{(\mathbf{x}^{(i)}, y^{(i)})\}_{i=1}^N$  and step size  $\gamma$ 
  1. Initialize  $\boldsymbol{\theta}^{(0)}$  to all zeros and set  $t = 0$
  2. While TERMINATION CRITERION is not satisfied
    - a. Randomly sample a data point from  $\mathcal{D}$ ,  $(\mathbf{x}^{(i)}, y^{(i)})$
    - b. Compute the pointwise gradient:
$$\nabla_{\boldsymbol{\theta}} J^{(i)}(\boldsymbol{\theta}^{(t)}) = \mathbf{x}^{(i)} (P(Y = 1 | \mathbf{x}^{(i)}, \boldsymbol{\theta}^{(t)}) - y^{(i)})$$
    - 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!

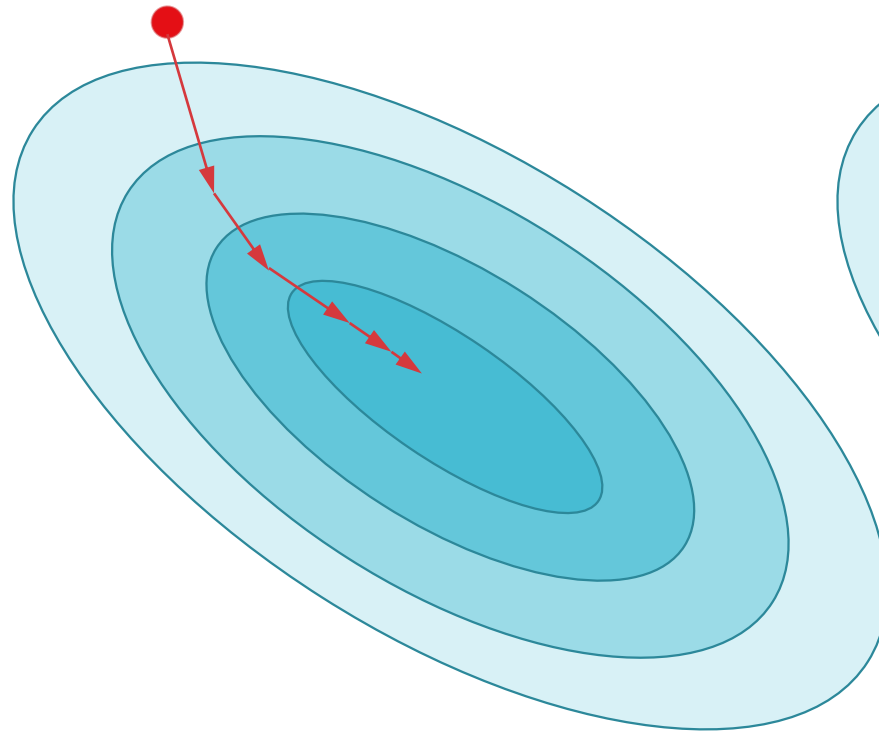
$$\begin{aligned} E[\nabla_{\theta} J^{(i)}(\theta)] &= \sum_{i=1}^N (\text{probability of selecting } \mathbf{x}^{(i)}, y^{(i)}) \nabla_{\theta} J^{(i)}(\theta) \\ &= \sum_{i=1}^N \left(\frac{1}{N}\right) \nabla_{\theta} J^{(i)}(\theta) = \frac{1}{N} \sum_{i=1}^N \nabla_{\theta} J^{(i)}(\theta) = \nabla_{\theta}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

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

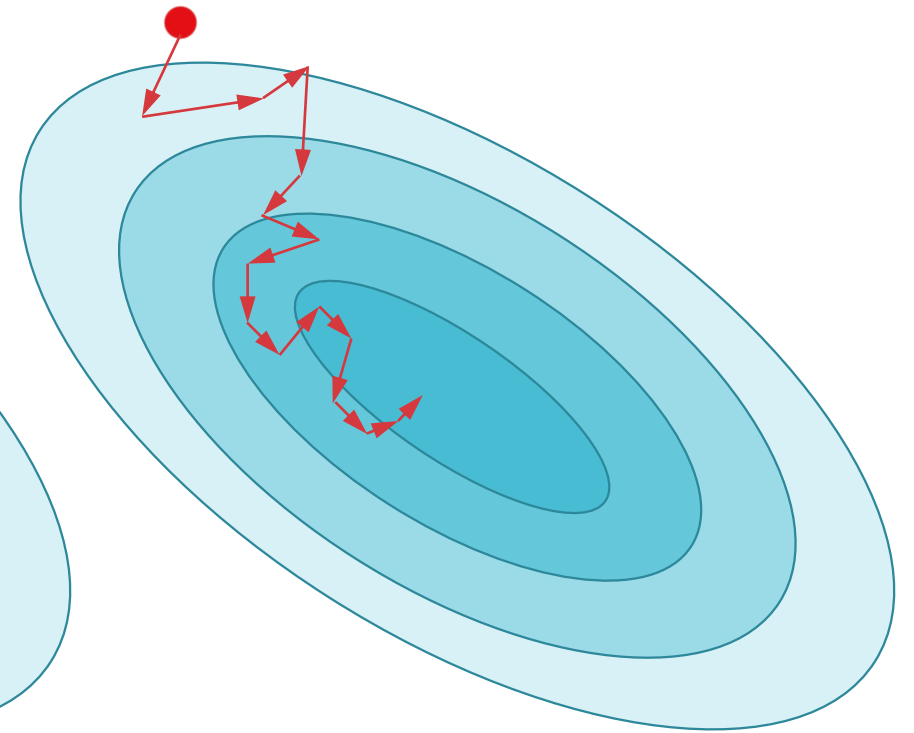
# Stochastic Gradient Descent (SGD)

- Input: training dataset  $\mathcal{D} = \{(\mathbf{x}^{(i)}, y^{(i)})\}_{i=1}^N$  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  $i \in \text{shuffle}(\{1, \dots, N\})$ 
      - i. Compute the pointwise gradient:
$$\nabla_{\boldsymbol{\theta}} J^{(i)}(\boldsymbol{\theta}^{(t)}) = \mathbf{x}^{(i)} (P(Y = 1 | \mathbf{x}^{(i)}, \boldsymbol{\theta}^{(t)}) - y^{(i)})$$
      - 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:  $\boldsymbol{\theta}^{(t)}$

# Stochastic Gradient Descent vs. Gradient Descent



Gradient Descent



Stochastic Gradient Descent

# Stochastic Gradient Descent vs. Gradient Descent

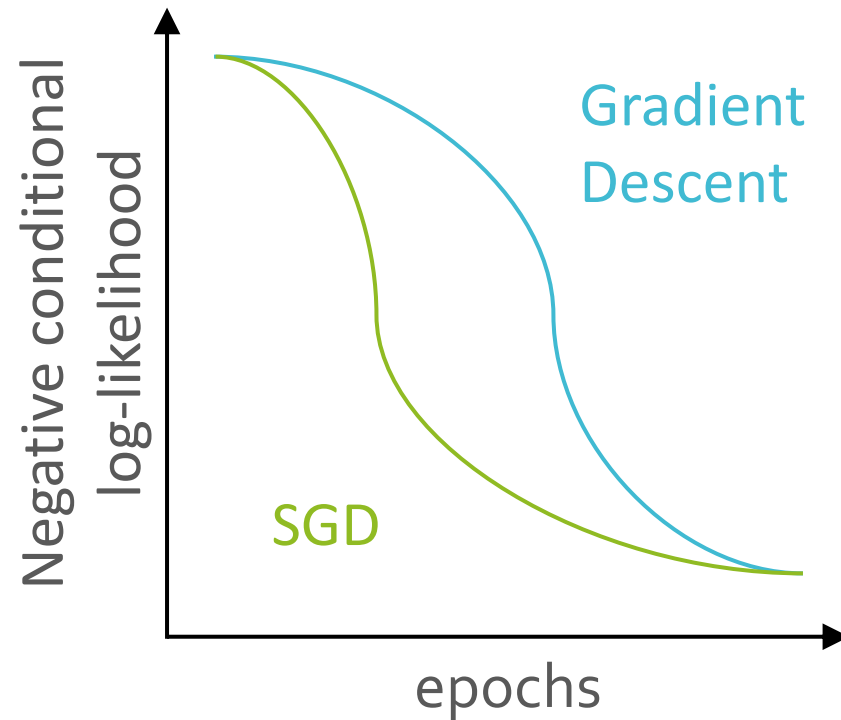
- An *epoch* is a single pass through the entire training dataset
  - Gradient descent updates the parameters once per epoch
  - SGD updates the parameters  $N$  times per epoch
- Theoretical comparison:
  - Define convergence to be when  $J(\boldsymbol{\theta}^{(t)}) - J(\boldsymbol{\theta}^*) < \epsilon$

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

# 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