

# 10-301/601: Introduction to Machine Learning

## Lecture 9 – Logistic Regression

Matt Gormley & Henry Chai

9/23/24

# Front Matter

- Announcements:
  - HW3 released 9/16, due 9/23 (today!) at 11:59 PM
    - **Only two grace days allowed on HW3**
  - Exam 1 on 9/30 (next Monday) from 6:30 PM - 8:30 PM
    - If you have a conflict, you must complete the [Exam conflict form](#) by 9/23 (today!) at 1 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

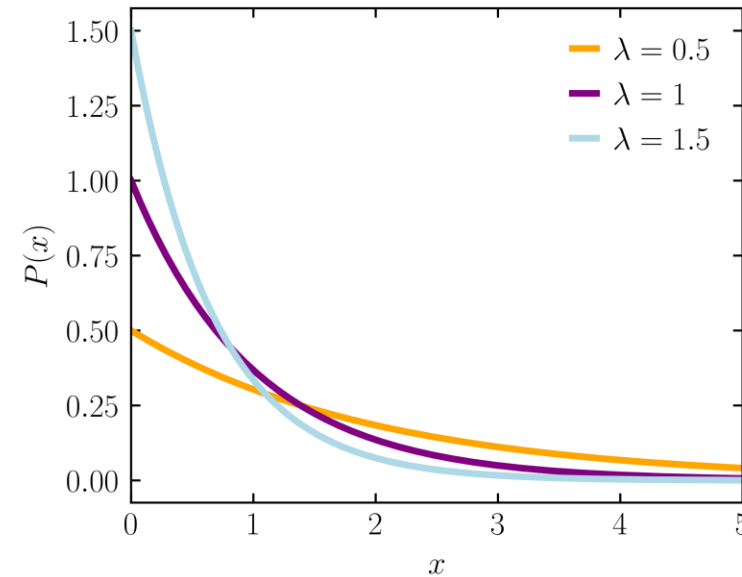
$$\ell(\theta) = \log \prod_{n=1}^N p(x^{(n)}|\theta) = \sum_{n=1}^N \log p(x^{(n)}|\theta)$$

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

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

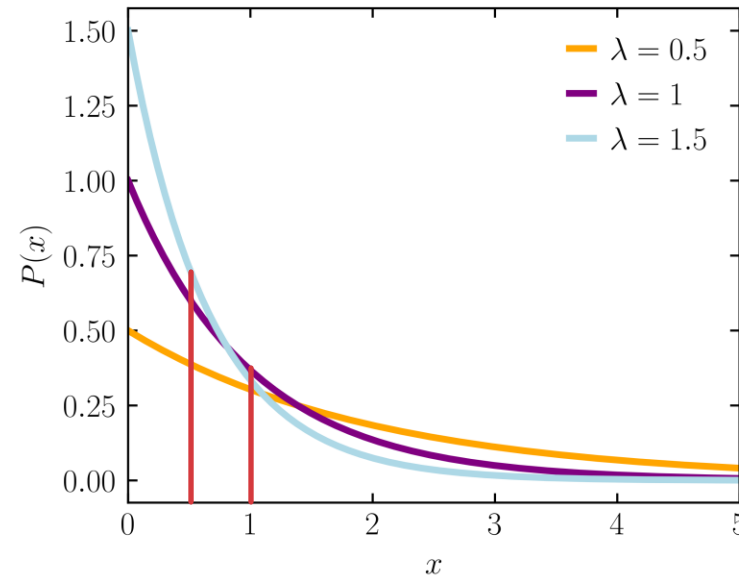
# Maximum Likelihood Estimation (MLE)

- 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)

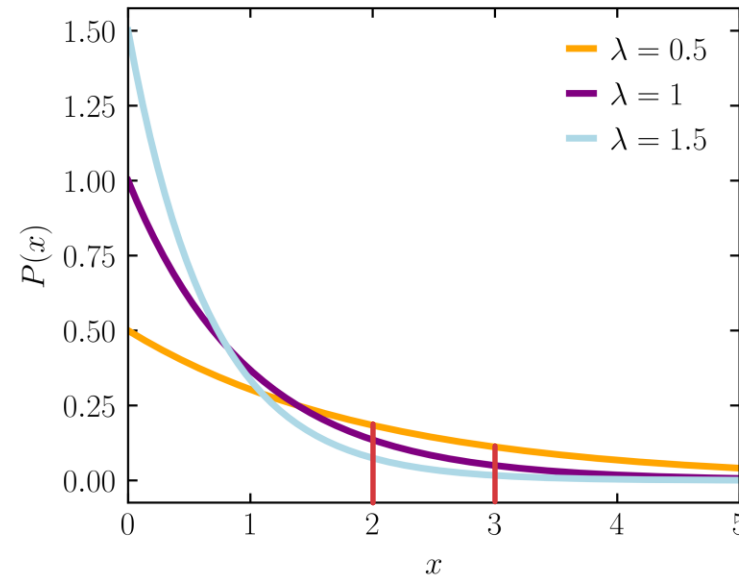
- Insight: every valid probability distribution has a finite amount of probability mass as it must sum/integrate to 1
- Idea: set the parameter(s) so that the likelihood of the samples is maximized
- Intuition: assign as much of the (finite) probability mass to the observed data *at the expense of unobserved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

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

$$\log(ab) = \log a + \log b$$

## 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

$$\begin{aligned} \ell(\lambda) &= \sum_{i=1}^N \log f(x^{(i)}|\lambda) = \sum_{i=1}^N \log \lambda e^{-\lambda x^{(i)}} \\ &= \sum_{i=1}^N (\log \lambda + (-\lambda x^{(i)})) \\ &= N \log \lambda - \lambda \sum_{i=1}^N x^{(i)} \\ \frac{\partial \ell}{\partial \lambda} &= \frac{N}{\lambda} - \sum_{i=1}^N x^{(i)} \Rightarrow \frac{\partial^2 \ell}{\partial \lambda^2} = -\frac{N}{\lambda^2} < 0 \\ \Rightarrow \frac{N}{\lambda} - \sum_{i=1}^N x^{(i)} &= 0 \Rightarrow \hat{\lambda} = \frac{N}{\sum_{i=1}^N x^{(i)}} \end{aligned}$$

# 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  $\mathbf{x}$

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

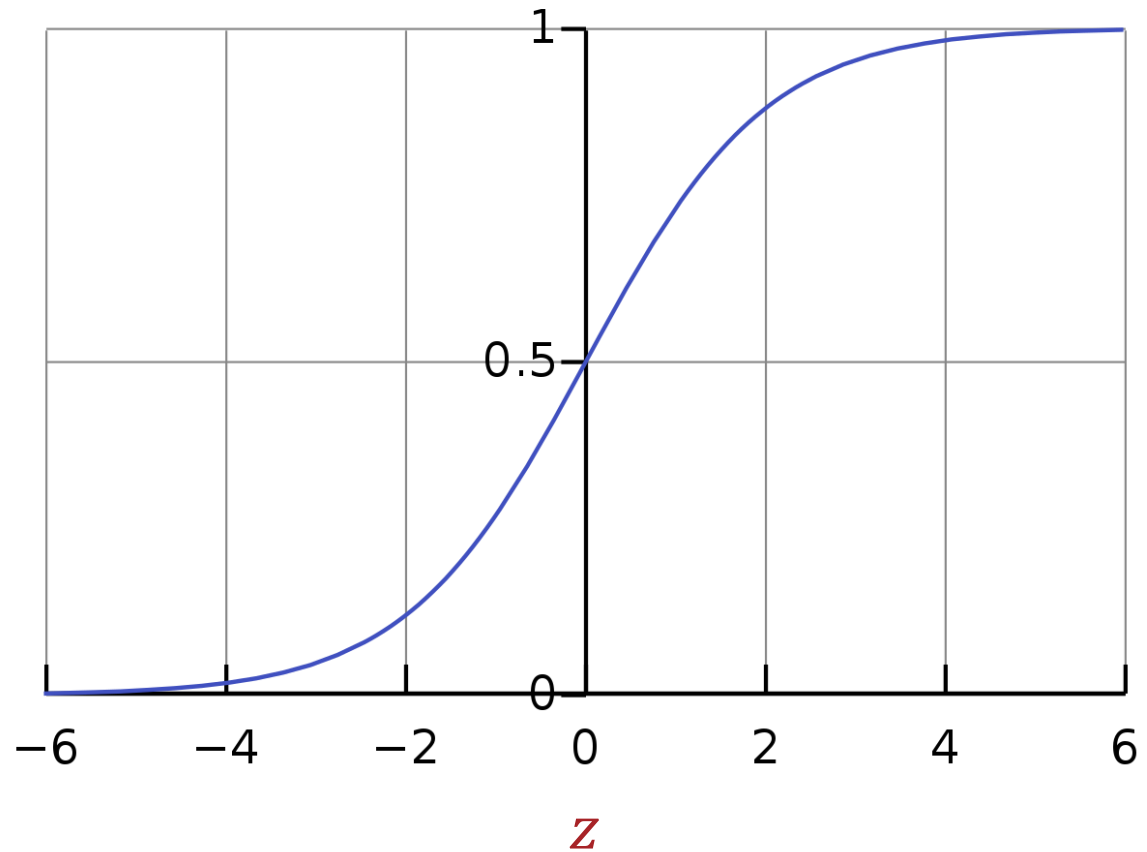
- This implies two useful facts:

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

$$2. \frac{P(Y=1 | \mathbf{x}, \boldsymbol{\theta})}{P(Y=0 | \mathbf{x}, \boldsymbol{\theta})} = \exp(\boldsymbol{\theta}^T \mathbf{x}) \Rightarrow \text{log odds is linear in } \mathbf{x}! \\ \log(\exp(\boldsymbol{\theta}^T \mathbf{x})) = \boldsymbol{\theta}^T \mathbf{x}$$

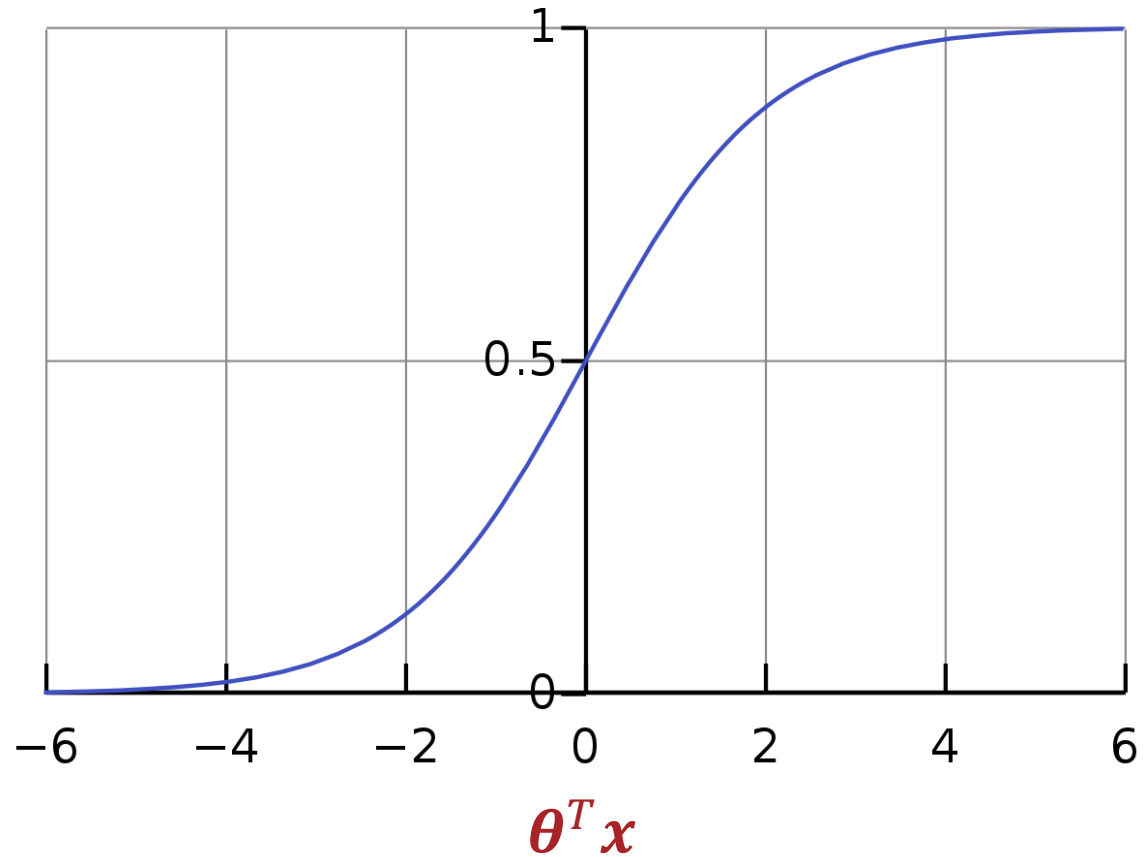
# Logistic Function

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



# Why use the Logistic Function?

$$\sigma(\theta^T x) = \frac{1}{1 + e^{-\theta^T x}}$$



1.  $\sigma: \mathbb{R} \mapsto (0, 1)$
2. differentiable everywhere
3. the decision boundary is linear in  $x$ !

# Logistic Regression Decision Boundary

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

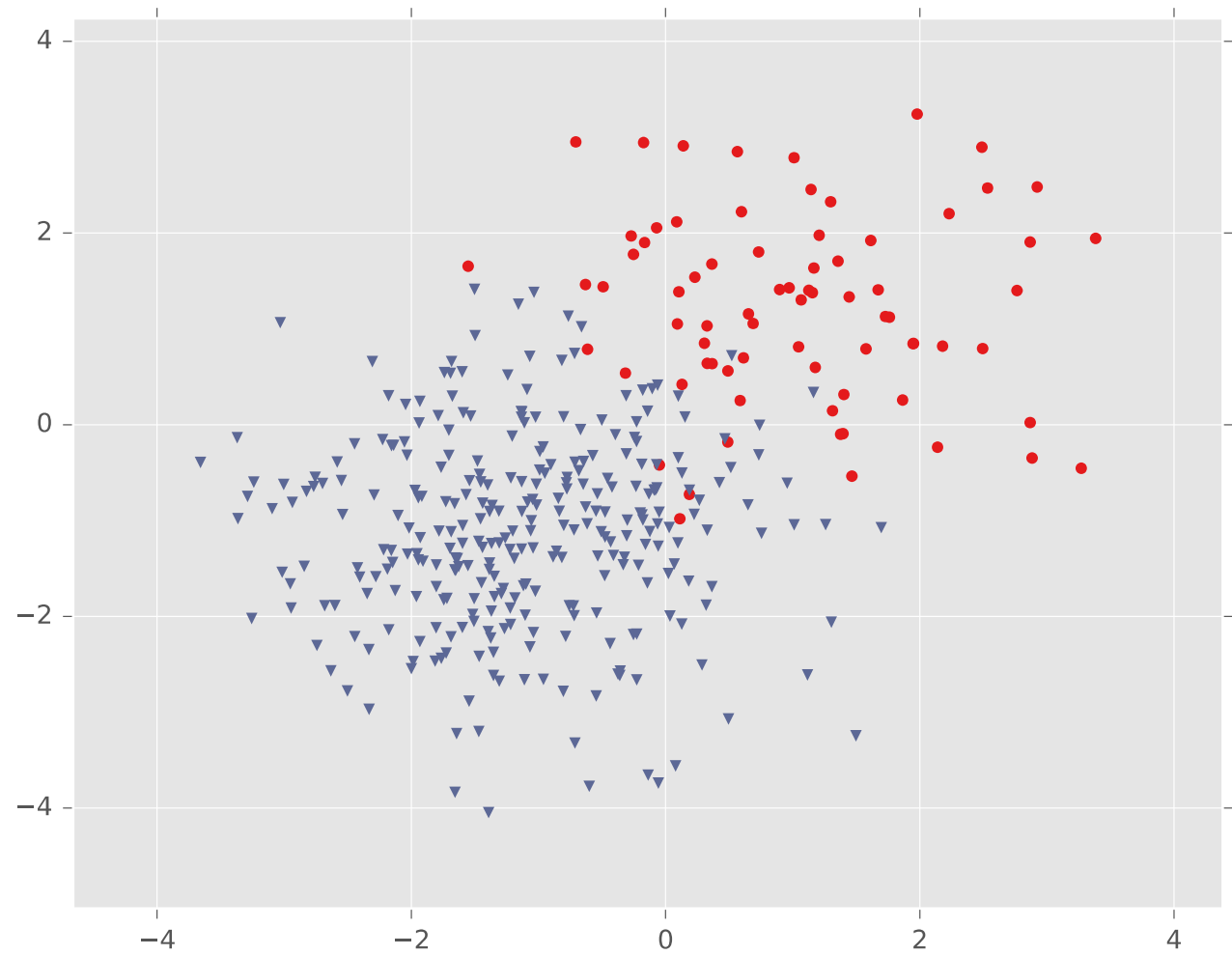
$$P(Y=1 | \mathbf{x}, \boldsymbol{\theta}) = \frac{1}{2} = \frac{1}{1 + \exp(-\boldsymbol{\theta}^T \mathbf{x})}$$

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

$$\Rightarrow -\boldsymbol{\theta}^T \mathbf{x} = \log(2 - 1)$$

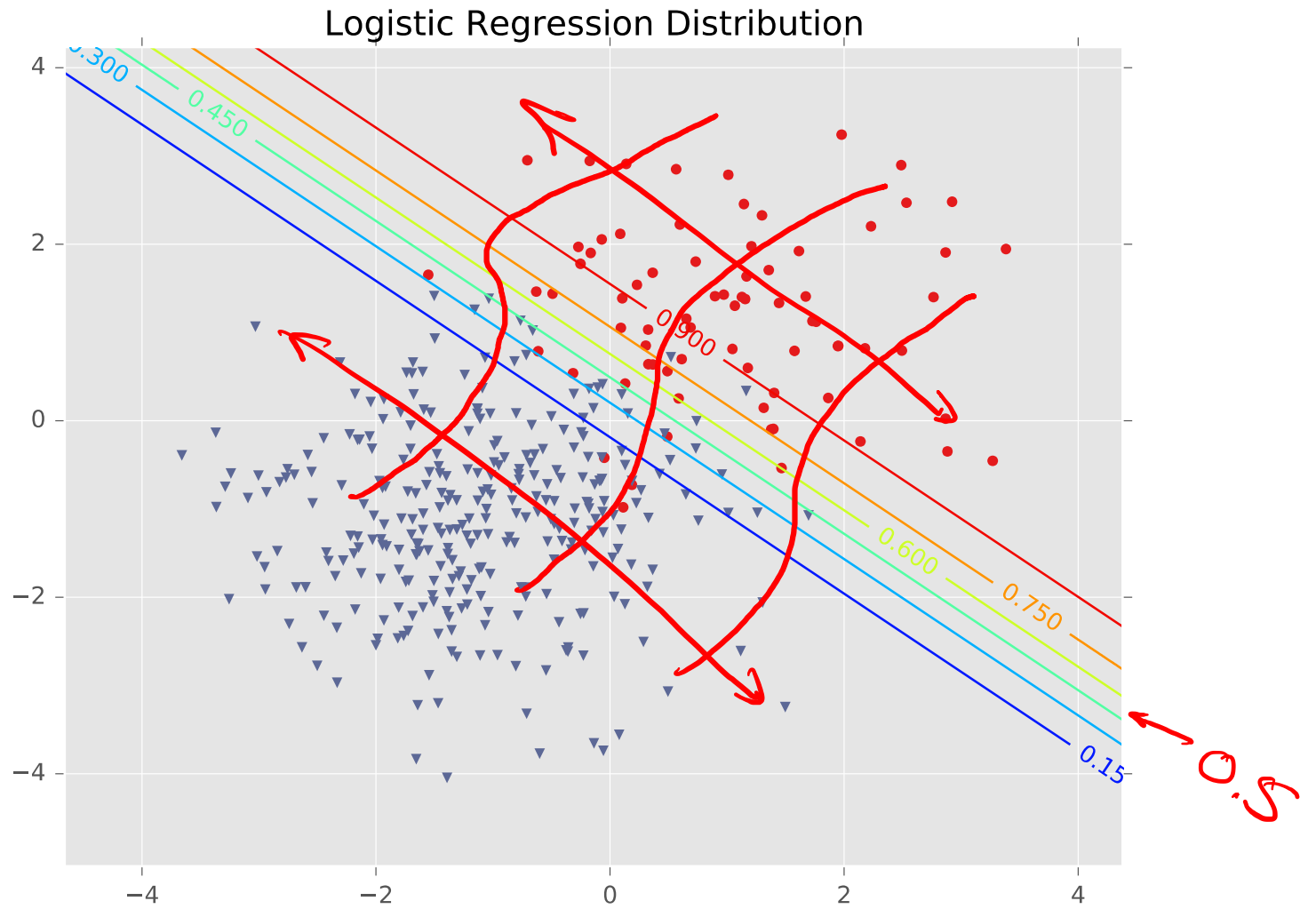
$$\Rightarrow \boldsymbol{\theta}^T \mathbf{x} = 0$$

# Logistic Regression Decision Boundary

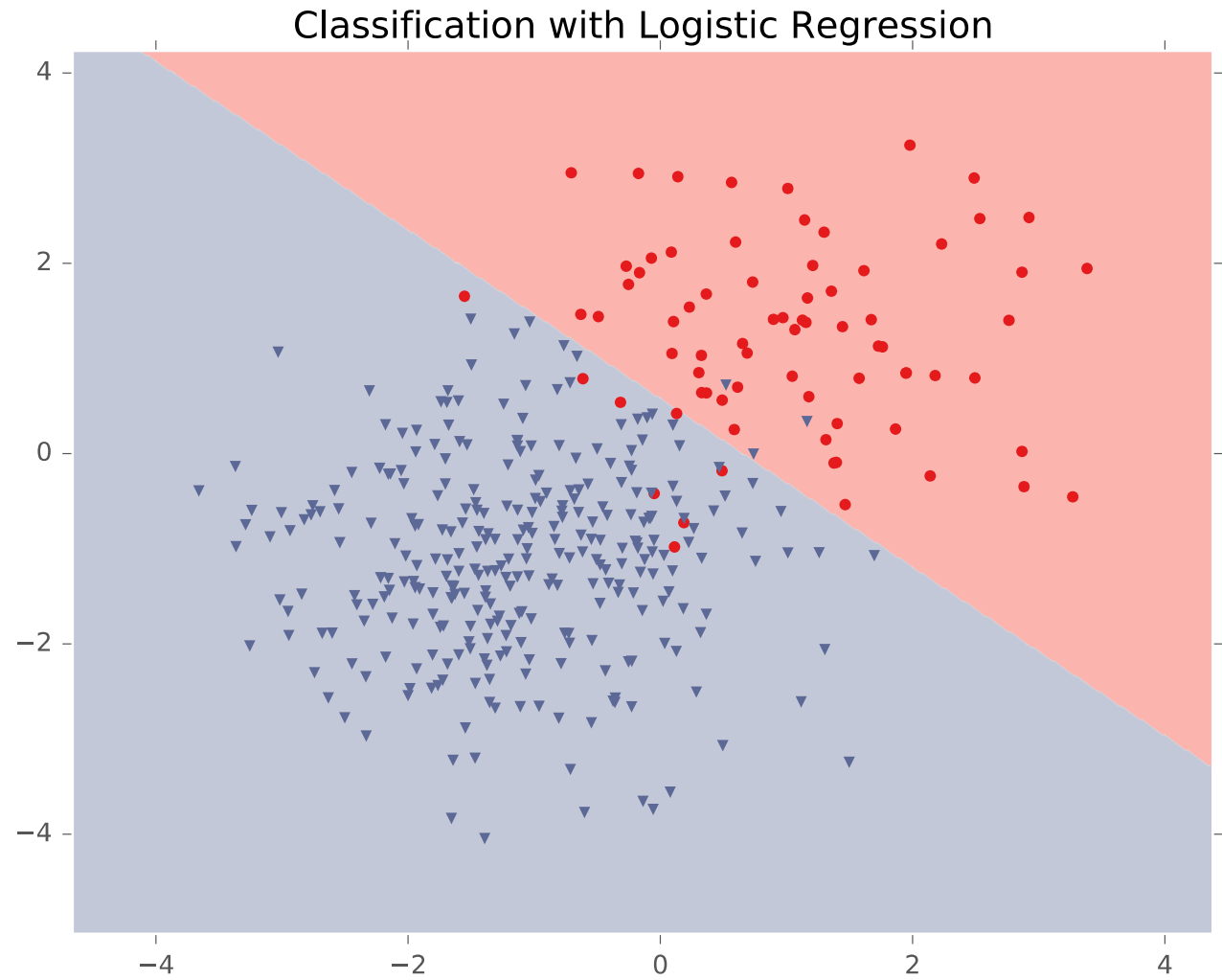




# Logistic Regression Decision Boundary



# Logistic Regression Decision Boundary



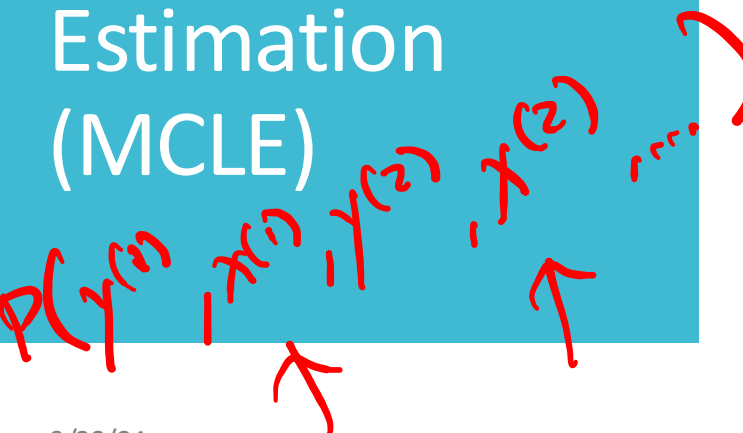
$$\log(a^b c^d)$$

$$= b \log a + d \log c$$

Setting the  
Parameters  
via Minimum  
Negative

**Conditional**

(log-)Likelihood  
Estimation  
(MCLE)



$$\ell(\theta) = -\log \mathbb{P}(y^{(1)}, y^{(2)}, \dots, y^{(N)} | x^{(1)}, x^{(2)}, \dots, x^{(N)}, \theta)$$

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

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

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

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

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

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

Key Takeaway:  
 This objective function is convex but we cannot minimize it in closed form

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

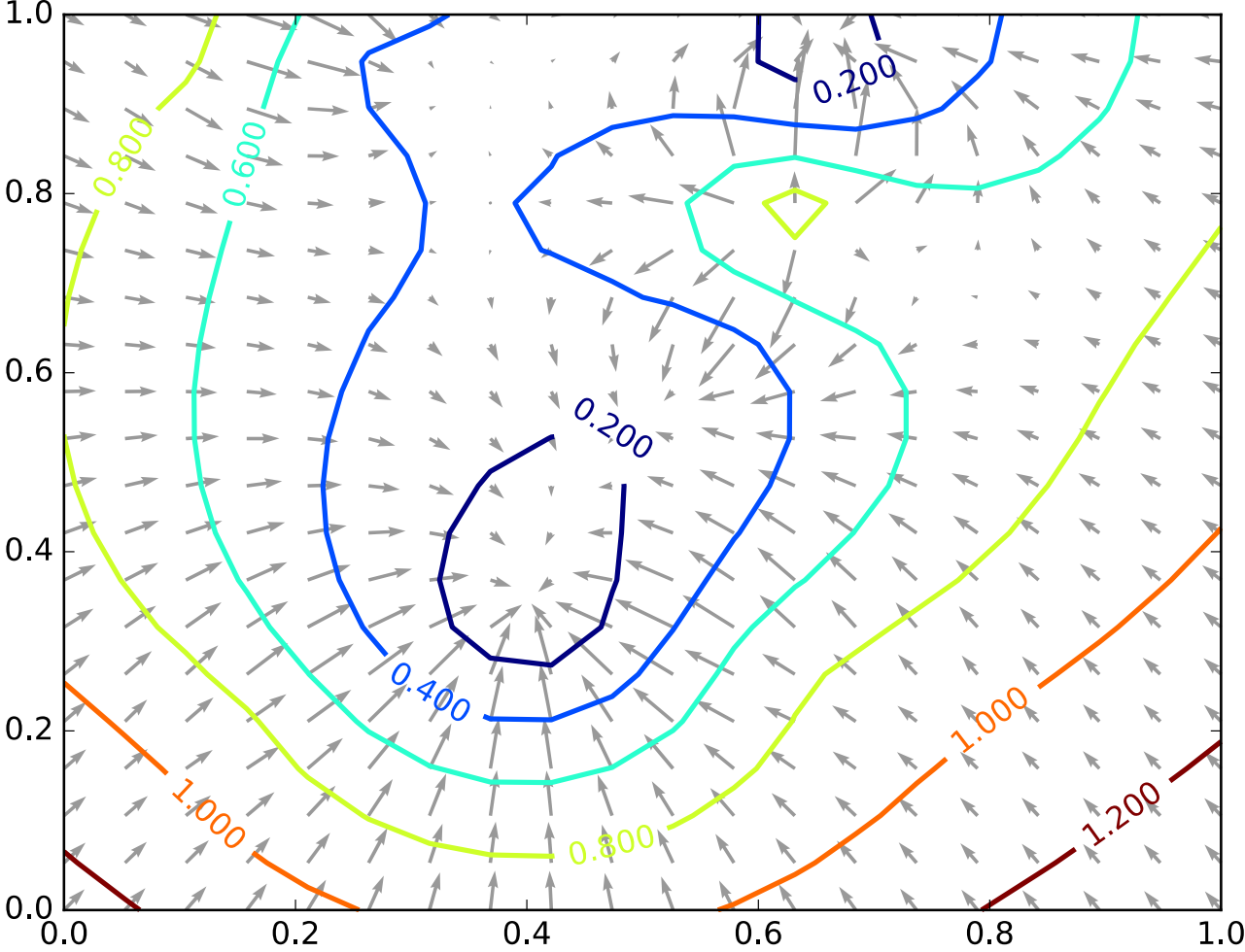
$$\nabla_{\theta} J(\theta) = -\frac{1}{N} \sum_{n=1}^N y^{(n)} x^{(n)} - \frac{\exp(\theta^T x^{(n)})}{1 + \exp(\theta^T x^{(n)})} x^{(n)}$$

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

$$P(Y=1 | x^{(n)}, \theta)$$

$$\nabla_{\theta} J(\hat{\theta}) = 0$$

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

# 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