

# ECE595 / STAT598: Machine Learning I

## Lecture 14.3: Logistic Regression 1 - Convexity

Spring 2020

Stanley Chan

School of Electrical and Computer Engineering  
Purdue University



# Outline

## Discriminative Approaches

- Lecture 14 Logistic Regression 1
- Lecture 15 Logistic Regression 2

## This lecture: Logistic Regression 1

- From Linear to Logistic
  - Motivation
  - Loss Function
  - Why not L2 Loss?
- Interpreting Logistic
  - Maximum Likelihood
  - Log-odd
- Convexity
  - Is logistic loss convex?
  - Computation

## Convexity of Logistic Training Loss

Recall that

$$J(\theta) = \sum_{n=1}^n - \left\{ y_n \log \left( \frac{h_{\theta}(\mathbf{x}_n)}{1 - h_{\theta}(\mathbf{x}_n)} \right) + \log(1 - h_{\theta}(\mathbf{x}_n)) \right\}$$

- The first term is linear, so it is convex.
- The second term: Gradient:

$$\begin{aligned} \nabla_{\theta}[-\log(1 - h_{\theta}(\mathbf{x}))] &= -\nabla_{\theta} \left[ \log \left( 1 - \frac{1}{1 + e^{-\theta^T \mathbf{x}}} \right) \right] \\ &= -\nabla_{\theta} \left[ \log \frac{e^{-\theta^T \mathbf{x}}}{1 + e^{-\theta^T \mathbf{x}}} \right] = -\nabla_{\theta} \left[ \log e^{-\theta^T \mathbf{x}} - \log(1 + e^{-\theta^T \mathbf{x}}) \right] \\ &= -\nabla_{\theta} \left[ -\theta^T \mathbf{x} - \log(1 + e^{-\theta^T \mathbf{x}}) \right] = \mathbf{x} + \nabla_{\theta} \left[ \log(1 + e^{-\theta^T \mathbf{x}}) \right] \\ &= \mathbf{x} + \left( \frac{-e^{-\theta^T \mathbf{x}}}{1 + e^{-\theta^T \mathbf{x}}} \right) \mathbf{x} = h_{\theta}(\mathbf{x})\mathbf{x}. \end{aligned}$$

## Convexity of Logistic Training Loss

- Gradient of second term is

$$\nabla_{\theta}[-\log(1 - h_{\theta}(\mathbf{x}))] = h_{\theta}(\mathbf{x})\mathbf{x}.$$

- Hessian is:

$$\begin{aligned}\nabla_{\theta}^2[-\log(1 - h_{\theta}(\mathbf{x}))] &= \nabla_{\theta} [h_{\theta}(\mathbf{x})\mathbf{x}] \\ &= \nabla_{\theta} \left[ \left( \frac{1}{1 + e^{-\theta^T \mathbf{x}}} \right) \mathbf{x} \right] \\ &= \left( \frac{1}{(1 + e^{-\theta^T \mathbf{x}})^2} \right) (-e^{-\theta^T \mathbf{x}}) \mathbf{x} \mathbf{x}^T \\ &= \left( \frac{1}{1 + e^{-\theta^T \mathbf{x}}} \right) \left( 1 - \frac{1}{1 + e^{-\theta^T \mathbf{x}}} \right) \mathbf{x} \mathbf{x}^T \\ &= h_{\theta}(\mathbf{x})[1 - h_{\theta}(\mathbf{x})]\mathbf{x} \mathbf{x}^T.\end{aligned}$$

## Convexity of Logistic Training Loss

- For any  $\mathbf{v} \in \mathbb{R}^d$ , we have that

$$\begin{aligned}\mathbf{v}^T \nabla_{\theta}^2 [-\log(1 - h_{\theta}(\mathbf{x}))] \mathbf{v} &= \mathbf{v}^T \left[ h_{\theta}(\mathbf{x}) [1 - h_{\theta}(\mathbf{x})] \mathbf{x} \mathbf{x}^T \right] \mathbf{v} \\ &= (h_{\theta}(\mathbf{x}) [1 - h_{\theta}(\mathbf{x})]) \|\mathbf{v}^T \mathbf{x}\|^2 \geq 0.\end{aligned}$$

- Therefore the Hessian is positive semi-definite.
- So  $-\log(1 - h_{\theta}(\mathbf{x}))$  is convex in  $\theta$ .
- Conclusion: The training loss function

$$J(\theta) = \sum_{n=1}^n - \left\{ y_n \log \left( \frac{h_{\theta}(\mathbf{x}_n)}{1 - h_{\theta}(\mathbf{x}_n)} \right) + \log(1 - h_{\theta}(\mathbf{x}_n)) \right\}$$

is **convex** in  $\theta$ .

- So we can use convex optimization algorithms to find  $\theta$ .

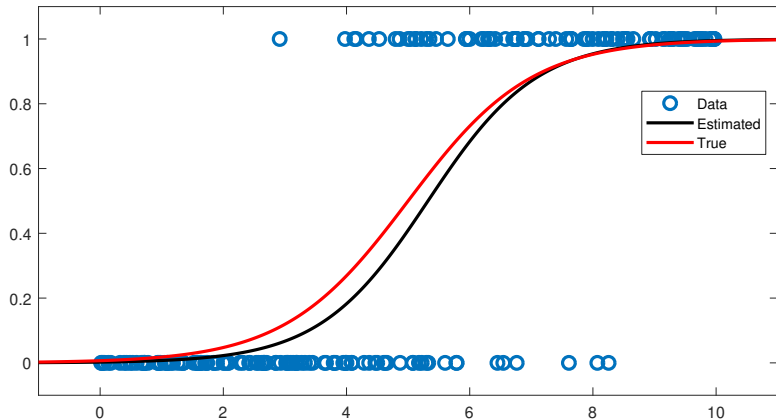
# Convex Optimization for Logistic Regression

- We can use CVX to solve the logistic regression problem
- But it requires some re-organization of the equations

$$\begin{aligned} J(\theta) &= \sum_{n=1}^N -\left\{ y_n \theta^T \mathbf{x}_n + \log(1 - h_{\theta}(\mathbf{x}_n)) \right\} \\ &= \sum_{n=1}^N -\left\{ y_n \theta^T \mathbf{x}_n + \log \left( 1 - \frac{e^{\theta^T \mathbf{x}_n}}{1 + e^{\theta^T \mathbf{x}_n}} \right) \right\} \\ &= \sum_{n=1}^N -\left\{ y_n \theta^T \mathbf{x}_n - \log \left( 1 + e^{\theta^T \mathbf{x}_n} \right) \right\} \\ &= -\left\{ \left( \sum_{n=1}^N y_n \mathbf{x}_n \right)^T \theta - \sum_{n=1}^N \log \left( 1 + e^{\theta^T \mathbf{x}_n} \right) \right\}. \end{aligned}$$

- The last term is a sum of log-sum-exp:  $\log(e^0 + e^{\theta^T \mathbf{x}})$ .

# Convex Optimization for Logistic Regression



# Reading List

## Logistic Regression (Machine Learning Perspective)

- Chris Bishop's *Pattern Recognition*, Chapter 4.3
- Hastie-Tibshirani-Friedman's *Elements of Statistical Learning*, Chapter 4.4
- Stanford CS 229 Discriminant Algorithms  
<http://cs229.stanford.edu/notes/cs229-notes1.pdf>
- CMU Lecture <https://www.stat.cmu.edu/~cshalizi/uADA/12/lectures/ch12.pdf>
- Stanford Language Processing  
<https://web.stanford.edu/~jurafsky/slp3/> (Lecture 5)

## Logistic Regression (Statistics Perspective)

- Duke Lecture <https://www2.stat.duke.edu/courses/Spring13/sta102.001/Lec/Lec20.pdf>
- Princeton Lecture  
<https://data.princeton.edu/wws509/notes/c3.pdf>