

DS 4400

Machine Learning and Data Mining I

Alina Oprea
Associate Professor
Khoury College of Computer Science
Northeastern University

October 8 2020

Announcements

- Released solutions for Homework 1
- Homework 2 is due on Tuesday, Oct. 13 at midnight
 - University holiday on Monday, Oct. 12
- Projects
 - Start thinking about theme, dataset, and topic
 - Look at shared resources and project examples in Piazza
 - Fill in a form with area preferences on Friday
 - Participate in discussion next week

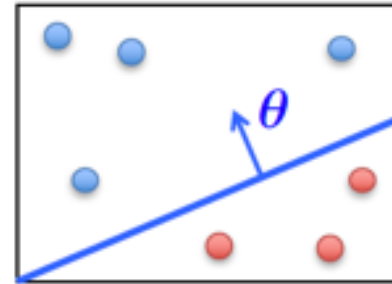
Outline

- Linear classifiers
- Perceptron wrap up
- Logistic regression
 - Classification based on probability
- Maximum Likelihood Estimation
 - Application to logistic regression
 - Cross-entropy objective
- Gradient descent for logistic regression

Linear Classifiers

- **Linear classifiers:** represent decision boundary by hyperplane

$$\theta = \begin{bmatrix} \theta_0 \\ \theta_1 \\ \vdots \\ \theta_d \end{bmatrix} \quad x^\top = \begin{bmatrix} 1 & x_1 & \dots & x_d \end{bmatrix}$$

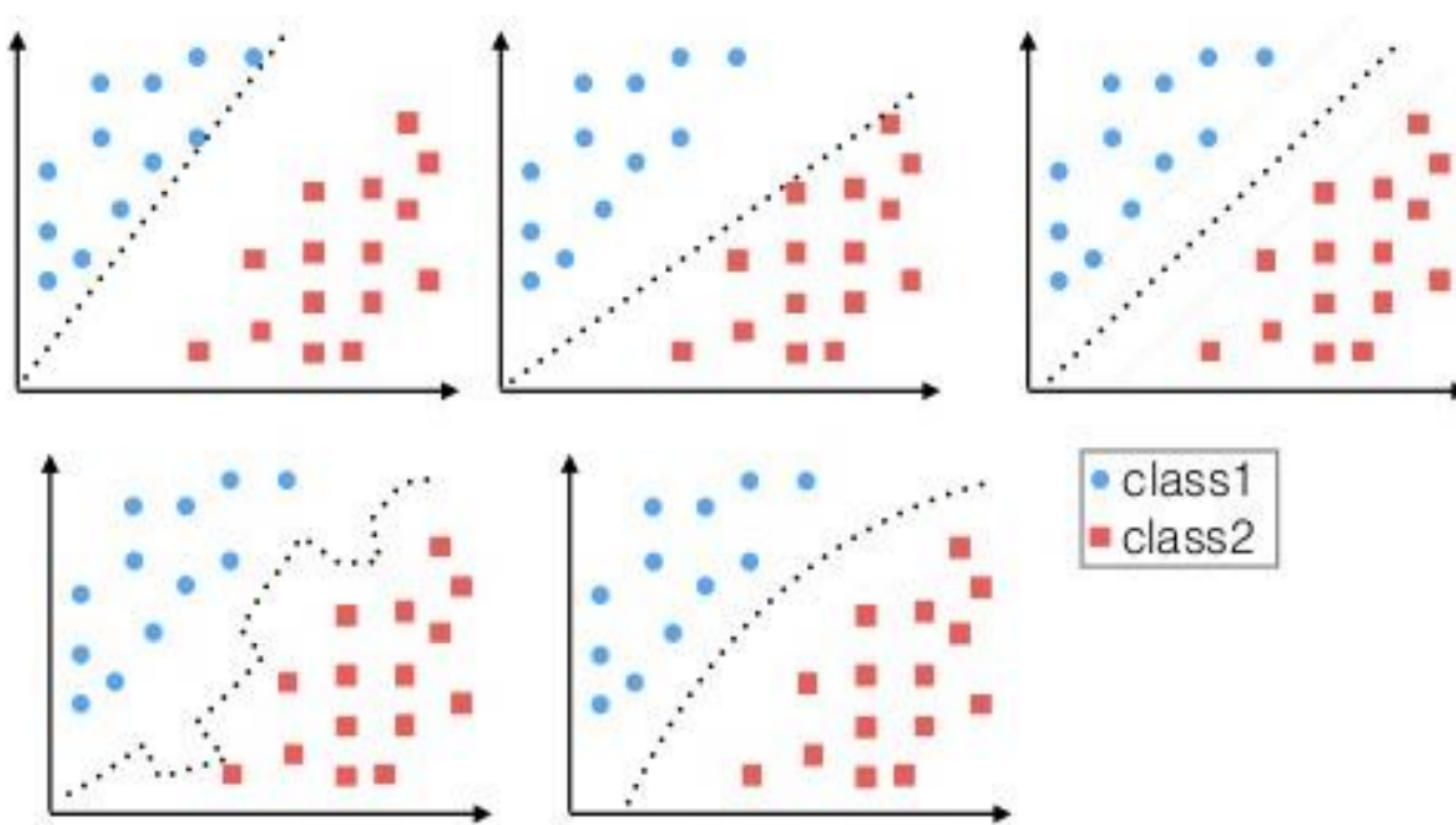


$h_\theta(x) = f(\theta^T x)$ linear function

- If $\theta^T x > 0$ classify “Class A”
- If $\theta^T x < 0$ classify “Class B”

All the points x on the hyperplane satisfy: $\theta^T x = 0$

Linear vs Non-Linear Classifiers



Online Perceptron

Let $\theta \leftarrow [0, 0, \dots, 0]$

Repeat:

Receive training example (x_i, y_i)

If $y_i \theta^T x_i \leq 0$ // prediction is incorrect

$\theta \leftarrow \theta + y_i x_i$

Online learning – the learning mode where the model update is performed each time a single observation is received

Batch learning – the learning mode where the model update is performed after observing the entire training set

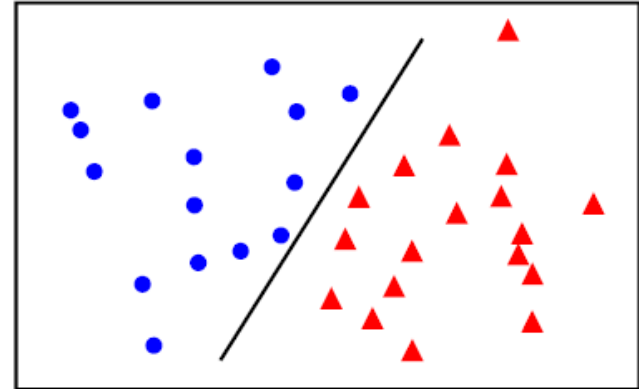
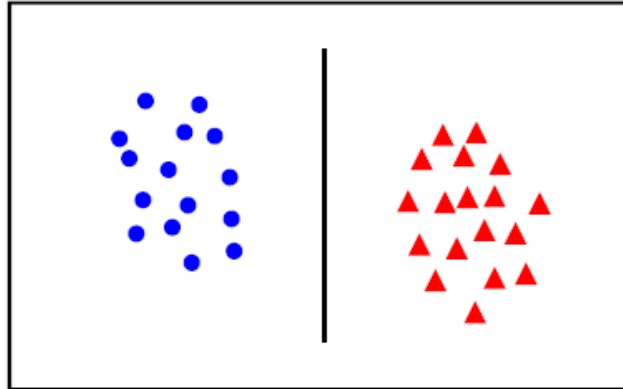
Batch Perceptron

```
Given training data  $\{(x_i, y_i)\}_{i=1}^n$ 
Let  $\theta \leftarrow [0, 0, \dots, 0]$ 
Repeat:
    Let  $\Delta \leftarrow [0, 0, \dots, 0]$ 
    for  $i = 1 \dots n$ , do
        if  $y_i \theta^T x_i \leq 0$  // prediction for  $i^{th}$  instance is incorrect
             $\Delta \leftarrow \Delta + y_i x_i$ 
     $\Delta \leftarrow \Delta / n$  // compute average update
     $\theta \leftarrow \theta + \Delta$ 
Until  $\|\Delta\|_2 < \epsilon$ 
```

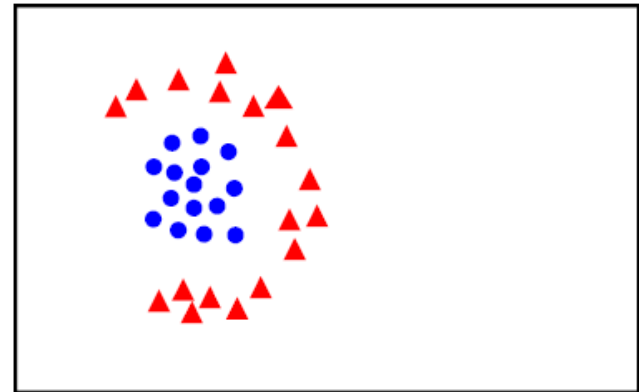
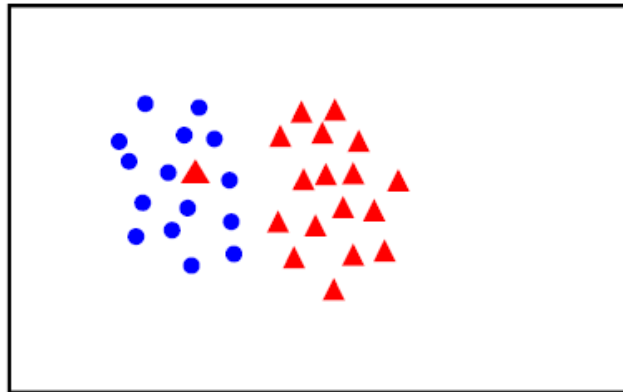
Guaranteed to find separating hyperplane if
data is linearly separable

Linear separability

linearly
separable





not
linearly
separable



- For linearly separable data, can prove bounds on perceptron error (depends on how well separated the data is)

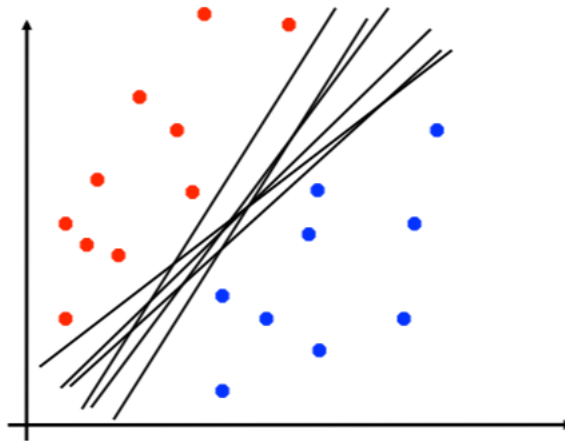
The perceptron

$$h_{\theta}(x) = f(\theta^T x)$$

- Linear classifier
- f is the sign function for the perceptron
- Pros
 - Very compact model (size d) 
- Cons of the perceptron
 - Perceptron depends on the order of training data and it could take many steps for convergence 
 - Only classifies well data that is linearly separable

Perceptron Limitations

- Is dependent on starting point
- It could take many steps for convergence
- Perceptron can overfit
 - Move the decision boundary for every example



Which of this is optimal?

Improving the Perceptron

- The Perceptron produces many θ 's during training
- The standard Perceptron simply uses the final θ at test time
 - This may sometimes not be a good idea!
 - Some other θ may be correct on 1,000 consecutive examples, but one mistake ruins it!
- **Idea:** Use a combination of multiple perceptrons
 - (i.e., neural networks!)
- **Idea:** Use the intermediate θ 's
 - **Voted Perceptron:** vote on predictions of the intermediate θ 's
 - **Averaged Perceptron:** average the intermediate θ 's

Classification Based on Probability

- Instead of just predicting the class, give the *probability of the instance being in that class*
 - Learn $P(Y|X)$
- Consider binary classifier with classes 0 and 1
 - $P(Y = 1|X) + P(Y = 0|X) = 1$
 - Sufficient to learn $P(Y = 1|X)$
- Advantages: interpretability and confidence of output

Logistic Regression

- Setup

- Training data: $\{x_i, y_i\}$, for $i = 1, \dots, N$
- Labels: $y_i \in \{0, 1\}$

- Goals

- Learn $P(Y = 1|X = x)$

- Highlights

- Probabilistic output
- At the basis of more complex models (e.g., neural networks)
- Supports regularization (Ridge, Lasso)
- Can be trained with Gradient Descent

Interpretation of Model Output

$$h_{\theta}(\mathbf{x}) = \text{estimated } P(Y = 1|X; \theta)$$

Example: Cancer diagnosis from tumor size

$$\mathbf{x} = \begin{bmatrix} x_0 \\ x_1 \end{bmatrix} = \begin{bmatrix} 1 \\ \text{tumorSize} \end{bmatrix}$$

$$h_{\theta}(\mathbf{x}) = 0.7$$

→ Tell patient that 70% chance of tumor being malignant

Note that: $P(Y = 0|X; \theta) + P(Y = 1|X; \theta) = 1$

Therefore, $P(Y = 0|X; \theta) = 1 - P(Y = 1|X; \theta)$

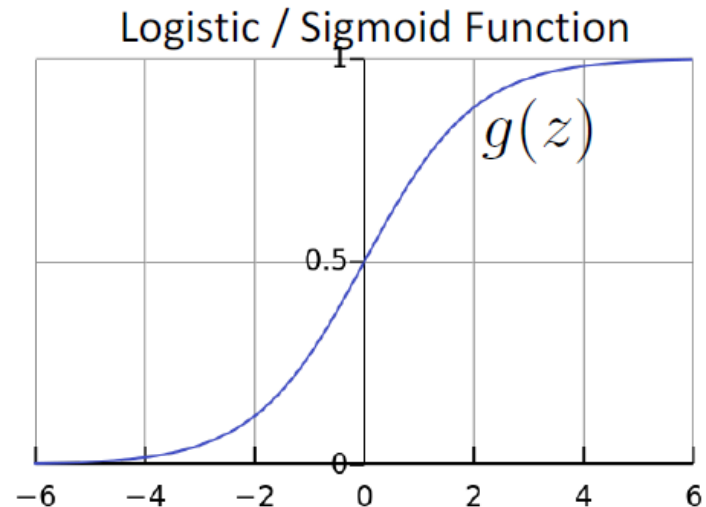
Logistic Regression

- Takes a probabilistic approach to learning discriminative functions (i.e., a classifier)
- $h_{\theta}(x)$ should give $P(Y = 1|X; \theta)$
 - Want $0 \leq h_{\theta}(x) \leq 1$
- Logistic regression model:

$$h_{\theta}(x) = g(\theta^T x)$$

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

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



LR is a Linear Classifier!

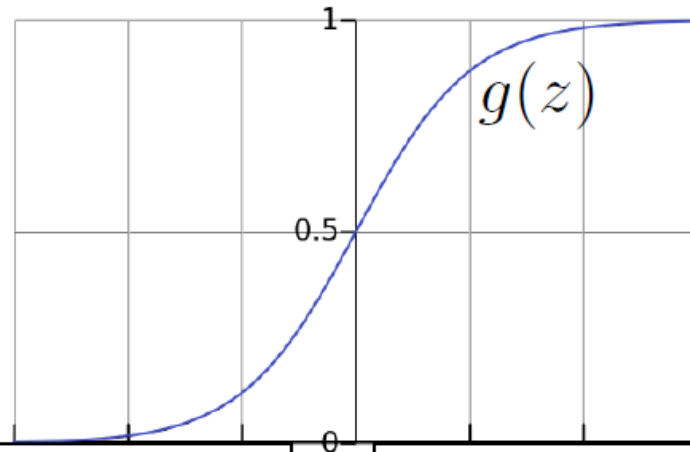
- Predict $Y = 1$ if:
 - $P[Y = 1|X = x; \theta] > P[Y = 0|X = x; \theta]$
 - $P[Y = 1|X = x; \theta] > \frac{1}{2}$
 - $$\frac{1}{1 + e^{-\theta^T x}} > \frac{1}{2}$$
- Equivalent to:
 - $e^{\theta_0 + \sum_{j=1}^d \theta_j x_j} > 1$
 - $\theta_0 + \sum_{j=1}^d \theta_j x_j > 0$

Logistic Regression is a linear classifier!

Logistic Regression

$$h_{\theta}(x) = g(\theta^T x)$$

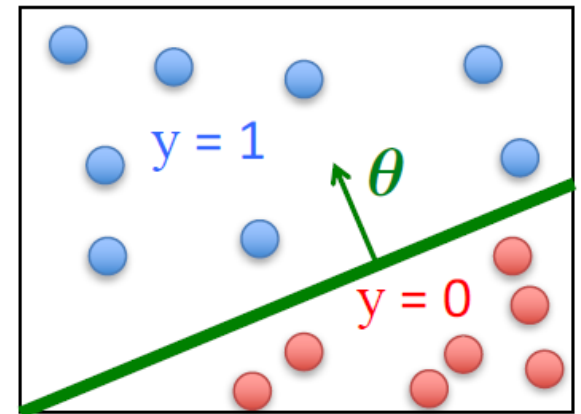
$$g(z) = \frac{1}{1 + e^{-z}}$$



$\theta^T x$ should be large negative values for negative instances

$\theta^T x$ should be large positive values for positive instances

- Assume a threshold and...
 - Predict $Y = 1$ if $h_{\theta}(x) \geq 0.5$
 - Predict $Y = 0$ if $h_{\theta}(x) < 0.5$



Logistic Regression is a linear classifier!

How to Pick Loss Function?

Maximum Likelihood Estimation (MLE)

Given training data $\mathbf{X} = \{x_1, \dots, x_N\}$ with labels $\mathbf{Y} = \{y_1, \dots, y_N\}$

What is the likelihood of training data for parameter θ ?

Define **likelihood function**

$$\text{Max}_{\theta} L(\theta) = P[\mathbf{Y}|\mathbf{X}; \theta]$$

Assumption: training points are independent

$$L(\theta) = \prod_{i=1}^N P[Y = y_i | X = x_i; \theta]$$

General probabilistic method for classifier training

Log Likelihood

- Max likelihood is equivalent to maximizing log of likelihood

$$L(\theta) = \prod_{i=1}^N P[Y = y_i | X = x_i; \theta]$$

$$\log L(\theta) = \sum_{i=1}^N \log P[Y = y_i | X = x_i; \theta]$$

- They both have the same maximum θ_{MLE}

MLE for Logistic Regression

$$P(Y = y_i | X = x_i; \theta) = h_{\theta}(x_i)^{y_i} (1 - h_{\theta}(x_i))^{1-y_i}$$

$$\begin{aligned}\theta_{MLE} &= \operatorname{argmax}_{\theta} \sum_{i=1}^N \log P[Y = y_i | X = x_i; \theta] \\ &= \operatorname{argmax}_{\theta} \sum_{i=1}^N y_i \log h_{\theta}(x_i) + (1 - y_i) \log (1 - h_{\theta}(x_i))\end{aligned}$$

Logistic regression objective

$$\min_{\theta} J(\theta)$$

$$J(\theta) = - \sum_{i=1}^N [y_i \log h_{\theta}(x_i) + (1 - y_i) \log (1 - h_{\theta}(x_i))]$$

Cross-Entropy Objective

$$J(\theta) = - \sum_{i=1}^N [y_i \log h_{\theta}(x_i) + (1 - y_i) \log (1 - h_{\theta}(x_i))]$$

- Cost of a single instance:

$$\text{cost}(h_{\theta}(\mathbf{x}), y) = \begin{cases} -\log(h_{\theta}(\mathbf{x})) & \text{if } y = 1 \\ -\log(1 - h_{\theta}(\mathbf{x})) & \text{if } y = 0 \end{cases}$$

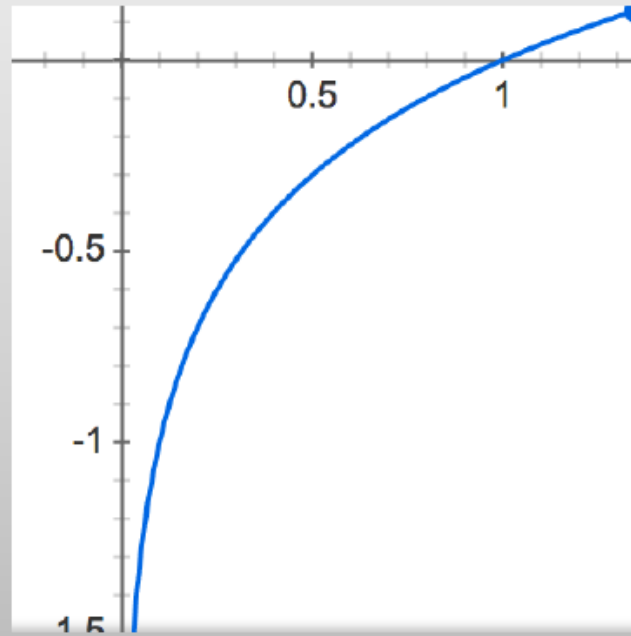
- Can re-write objective function as

$$J(\theta) = \sum_{i=1}^n \underbrace{\text{cost}(h_{\theta}(x_i), y_i)}_{\text{Cross-entropy loss}}$$

Intuition

$$\text{cost}(h_{\theta}(\mathbf{x}), y) = \begin{cases} -\log(h_{\theta}(\mathbf{x})) & \text{if } y = 1 \\ -\log(1 - h_{\theta}(\mathbf{x})) & \text{if } y = 0 \end{cases}$$

Aside: Recall the plot of $\log(z)$

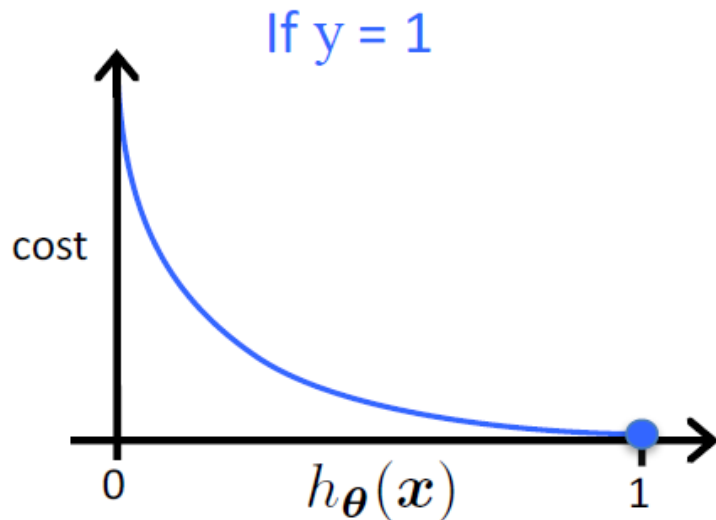


Intuition

$$\text{cost}(h_{\theta}(\mathbf{x}), y) = \begin{cases} -\log(h_{\theta}(\mathbf{x})) & \text{if } y = 1 \\ -\log(1 - h_{\theta}(\mathbf{x})) & \text{if } y = 0 \end{cases}$$

If $y = 1$

- Cost = 0 if prediction is correct
- As $h_{\theta}(\mathbf{x}) \rightarrow 0$, cost $\rightarrow \infty$
- Captures intuition that larger mistakes should get larger penalties
 - e.g., predict $h_{\theta}(\mathbf{x}) = 0$, but $y = 1$

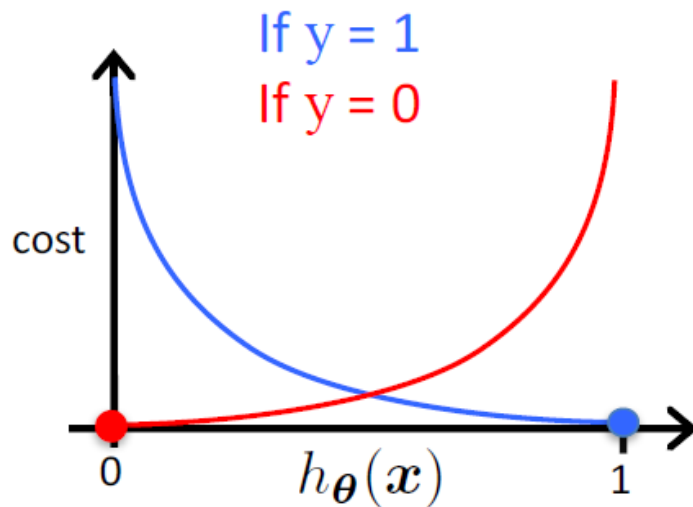


Intuition

$$\text{cost}(h_{\theta}(\mathbf{x}), y) = \begin{cases} -\log(h_{\theta}(\mathbf{x})) & \text{if } y = 1 \\ -\log(1 - h_{\theta}(\mathbf{x})) & \text{if } y = 0 \end{cases}$$

If $y = 0$

- Cost = 0 if prediction is correct
- As $(1 - h_{\theta}(\mathbf{x})) \rightarrow 0$, $\text{cost} \rightarrow \infty$
- Captures intuition that larger mistakes should get larger penalties



Cross-Entropy Objective

$$P(Y = y_i | X = x_i; \theta) = h_{\theta}(x_i)^{y_i} (1 - h_{\theta}(x_i))^{1-y_i}$$

$$\begin{aligned}\theta_{MLE} &= \operatorname{argmax}_{\theta} \sum_{i=1}^N \log P[Y = y_i | X = x_i; \theta] \\ &= \operatorname{argmax}_{\theta} \sum_{i=1}^N y_i \log h_{\theta}(x_i) + (1 - y_i) \log (1 - h_{\theta}(x_i))\end{aligned}$$

Logistic regression objective

$$\min_{\theta} J(\theta)$$

$$J(\theta) = - \sum_{i=1}^N [y_i \log h_{\theta}(x_i) + (1 - y_i) \log (1 - h_{\theta}(x_i))]$$

Logistic Regression

Lab Example

Review

- Perceptron is the first example of linear classifier
 - Online and batch learning
 - Has several limitations
- Logistic regression is a linear classifier that predicts class probability
- Maximum Likelihood Estimation is a method to estimate model parameters
 - Derive cross-entropy loss function
- Logistic regression can be trained with GD

Acknowledgements

- Slides made using resources from:
 - Andrew Ng
 - Eric Eaton
 - David Sontag
- Thanks!