

10-605/805 – ML for Large Datasets

Lecture 6: Distributed PCA & Logistic Regression

Henry Chai

9/15/22

Front Matter

- Recitation 3 on 9/16 will go over HW1 solutions
- HW2 released 9/8, due 9/22 at 11:59 PM

Recall: PCA Algorithm

- Input: $\mathcal{D} = \{(\mathbf{x}^{(i)})\}_{i=1}^n, r$
 1. Center the data
 - A. Optionally, normalize the data by features so that all features are of the same scale
 2. Compute the covariance matrix $\mathbf{C}_X = X^T X$
 3. Collect the top r eigenvectors (corresponding to the r largest eigenvalues), $P \in \mathbb{R}^{k \times r}$
 4. Project the data into the space defined by P , $Z = XP$

Recall: Computational Cost of PCA

- Input: $\mathcal{D} = \{(\mathbf{x}^{(i)})\}_{i=1}^n, r$
 1. Center the data
 - A. Optionally, normalize the data by features so that all features are of the same scale
 2. Compute the covariance matrix $\mathbf{C}_X = X^T X$ ($O(nk^2)$)
 3. Collect the top r eigenvectors (corresponding to the r largest eigenvalues), $P \in \mathbb{R}^{k \times r}$ ($O(k^3)$)
 4. Project the data into the space defined by P , $Z = XP$ ($O(nkr)$)

PCA: Large n , Small k

- Assume $O(k^3)$ computation and $O(k^2)$ storage is possible on a single machine
 - ✓ We can store and compute the eigenvalues of $X^T X$
 - We cannot compute $X^T X$
 - We cannot store X
- Approach: basically the same as distributed linear regression
 1. Center the data in a distributed way
 2. Store the rows of X across different machines
 3. Compute $X^T X$ as the sum of outer products

μ

Workers	$\begin{bmatrix} \leftarrow & \mathbf{x}^{(1)T} & \rightarrow \\ \leftarrow & \mathbf{x}^{(4)T} & \rightarrow \\ \vdots & \vdots & \vdots \end{bmatrix}$	$\begin{bmatrix} \leftarrow & \mathbf{x}^{(2)T} & \rightarrow \\ \leftarrow & \mathbf{x}^{(3)T} & \rightarrow \\ \vdots & \vdots & \vdots \end{bmatrix}$	$\begin{bmatrix} \leftarrow & \mathbf{x}^{(5)T} & \rightarrow \\ \leftarrow & \mathbf{x}^{(7)T} & \rightarrow \\ \vdots & \vdots & \vdots \end{bmatrix}$
Reduce	$\mu = \frac{1}{n} \sum_{i=1}^n \mathbf{x}^{(i)}$		
Map	$\tilde{\mathbf{x}}^{(i)} = \mathbf{x}^{(i)} - \mu$	$\tilde{\mathbf{x}}^{(i)} = \mathbf{x}^{(i)} - \mu$	$\tilde{\mathbf{x}}^{(i)} = \mathbf{x}^{(i)} - \mu$

$O(nk)$ distributed storage
(total)

$O(k)$
local work

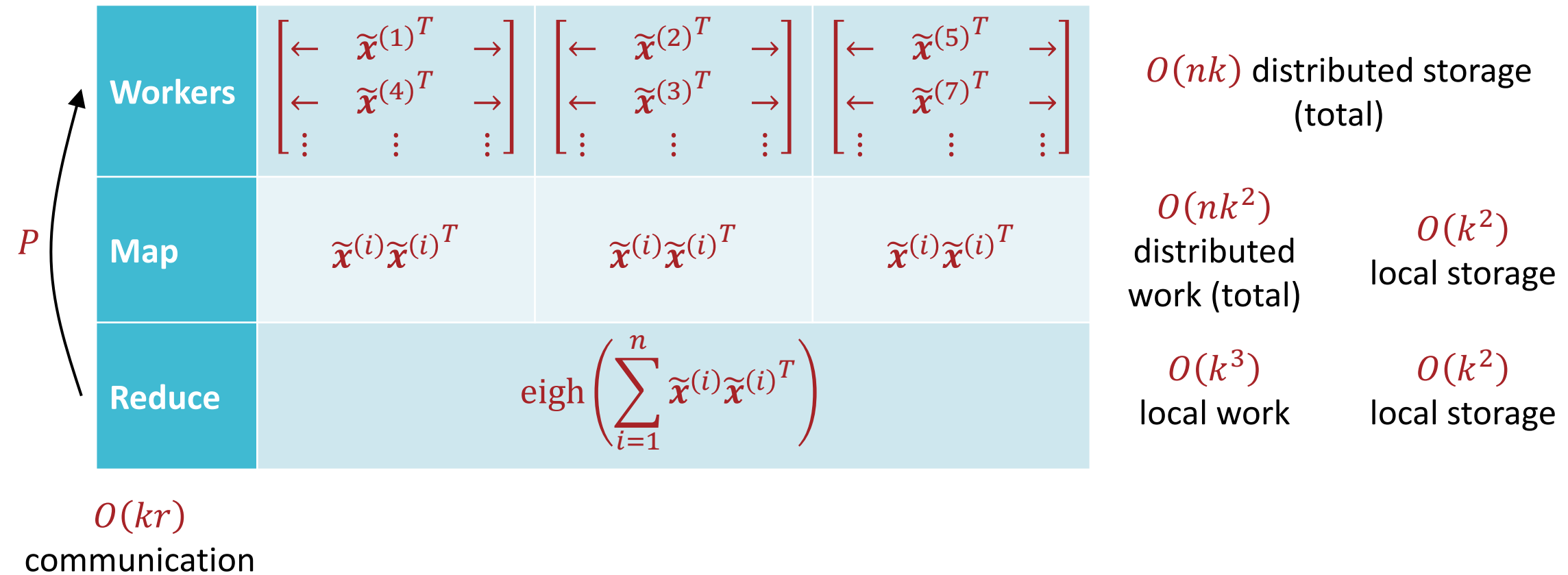
$O(k)$
local storage

$O(nk)$
distributed
work (total)

$O(nk)$
distributed
storage (total)

$O(k)$
communication

Distributed Centering of the Data



Distributed Eigendecomposition of $X^T X$

Workers	$\begin{bmatrix} \leftarrow & \tilde{\mathbf{x}}^{(1)T} & \rightarrow \\ \leftarrow & \tilde{\mathbf{x}}^{(4)T} & \rightarrow \\ \vdots & \vdots & \vdots \end{bmatrix}$	$\begin{bmatrix} \leftarrow & \tilde{\mathbf{x}}^{(2)T} & \rightarrow \\ \leftarrow & \tilde{\mathbf{x}}^{(3)T} & \rightarrow \\ \vdots & \vdots & \vdots \end{bmatrix}$	$\begin{bmatrix} \leftarrow & \tilde{\mathbf{x}}^{(5)T} & \rightarrow \\ \leftarrow & \tilde{\mathbf{x}}^{(7)T} & \rightarrow \\ \vdots & \vdots & \vdots \end{bmatrix}$	$O(nk)$ distributed storage (total)	
Map	$P\tilde{\mathbf{x}}^{(i)}$	$P\tilde{\mathbf{x}}^{(i)}$	$P\tilde{\mathbf{x}}^{(i)}$	$O(nkr)$ distributed work (total)	$O(nr)$ local storage

Distributed Computation of PCA Scores

PCA: Large n , Large k

- Now, $O(k^3)$ computation and $O(k^2)$ storage is *not* possible on a single machine

We cannot store and compute the eigenvalues of $X^T X$

We cannot compute $X^T X$

We cannot store X

- Idea: use a different algorithm!
 - Turn to an iterative method for computing eigenvectors

PCA: Large n , Large k

- Now, $O(k^3)$ computation and $O(k^2)$ storage is *not* possible on a single machine

We cannot store and compute the eigenvalues of $X^T X$

We cannot compute $X^T X$

We cannot store X

- Idea: use a different algorithm!
 - Turn to an iterative method for computing ~~eigenvectors~~ the eigenvector associated with the largest eigenvalue ($r = 1$) \rightarrow power iteration

Power Iteration

- Fact: $A = X^T X$ is “diagonalizable”, i.e., any k -dimensional vector can be written as a linear combination of A ’s eigenvectors:

$$\mathbf{b} = c_1 \mathbf{v}_1 + c_2 \mathbf{v}_2 \dots + c_k \mathbf{v}_k$$

- This follows because $X^T X$ is real and symmetric
- Assume $A = X^T X$ has one eigenvalue that is strictly larger than the others:

$$\lambda_1 > \lambda_2 \geq \dots \geq \lambda_k$$

Power Iteration

- Input: $A = X^T X$
- Initialize $\mathbf{b}^{(0)}$ to all zeros and set $t = 0$
- 1. While TERMINATION CRITERION is not satisfied
 - a. Update the vector \mathbf{b} :
$$\mathbf{b}^{(t+1)} \leftarrow \frac{A\mathbf{b}^{(t)}}{\|A\mathbf{b}^{(t)}\|_2}$$
 - b. Increment t : $t \leftarrow t + 1$
- Output: $\mathbf{b}^{(t)}$, the eigenvector corresponding to the largest eigenvalue of A !

Power Iteration

$$\mathbf{b}^{(0)} = c_1 \mathbf{v}_1 + c_2 \mathbf{v}_2 + \cdots + c_k \mathbf{v}_k$$

$$\begin{aligned} A\mathbf{b}^{(0)} &= c_1 A\mathbf{v}_1 + c_2 A\mathbf{v}_2 + \cdots + c_k A\mathbf{v}_k \\ &= c_1 \lambda_1 \mathbf{v}_1 + c_2 \lambda_2 \mathbf{v}_2 + \cdots + c_k \lambda_k \mathbf{v}_k \end{aligned}$$

$$\begin{aligned} A(A\mathbf{b}^{(0)}) &= (A)^2 \mathbf{b}^{(0)} = c_1 \lambda_1 A\mathbf{v}_1 + c_2 \lambda_2 A\mathbf{v}_2 + \cdots + c_k \lambda_k A\mathbf{v}_k \\ &= c_1 \lambda_1^2 \mathbf{v}_1 + c_2 \lambda_2^2 \mathbf{v}_2 + \cdots + c_k \lambda_k^2 \mathbf{v}_k \end{aligned}$$

$$\begin{aligned} (A)^t \mathbf{b}^{(0)} &= c_1 \lambda_1^t \mathbf{v}_1 + c_2 \lambda_2^t \mathbf{v}_2 + \cdots + c_k \lambda_k^t \mathbf{v}_k \\ &= c_1 \lambda_1^t \left(\mathbf{v}_1 + c_2 \left(\frac{\lambda_2}{\lambda_1} \right)^t \mathbf{v}_2 + \cdots + c_k \left(\frac{\lambda_k}{\lambda_1} \right)^t \mathbf{v}_k \right) \end{aligned}$$

$$\rightarrow c_1 \lambda_1^t (\mathbf{v}_1) \text{ as } t \rightarrow \infty$$

Normalize $(A)^t \mathbf{b}^{(0)}$ to have unit norm: $\mathbf{b}^{(t)} = \frac{(A)^t \mathbf{b}^{(0)}}{\|(A)^t \mathbf{b}^{(0)}\|_2}$

Power Iteration

- Input: $A = X^T X$
- Initialize $\mathbf{b}^{(0)}$ to all zeros and set $t = 0$
- 1. While TERMINATION CRITERION is not satisfied
 - a. Update the vector \mathbf{b} :

$$\mathbf{b}^{(t+1)} \leftarrow \frac{X^T X \mathbf{b}^{(t)}}{\|X^T X \mathbf{b}^{(t)}\|_2}$$

- b. Increment t : $t \leftarrow t + 1$

- Output: $\mathbf{b}^{(t)}$, the eigenvector corresponding to the largest eigenvalue of A !

Power Iteration

- Input: $A = X^T X$
- Initialize $\mathbf{b}^{(0)}$ to all zeros and set $t = 0$
- 1. While TERMINATION CRITERION is not satisfied
 - a. Update the vector \mathbf{b} using some intermediate steps:

$$\mathbf{b}^{(t+1)} \leftarrow X^T X \mathbf{b}^{(t)} = \left(\sum_{i=1}^n \tilde{\mathbf{x}}^{(i)} \tilde{\mathbf{x}}^{(i)T} \right) \mathbf{b}^{(t)}$$

$$\mathbf{b}^{(t+1)} \leftarrow \frac{\mathbf{b}^{(t+1)}}{\|\mathbf{b}^{(t+1)}\|_2}$$

- b. Increment t : $t \leftarrow t + 1$

- Output: $\mathbf{b}^{(t)}$, the eigenvector corresponding to the largest eigenvalue of A !

Power Iteration

- Input: $A = X^T X$
- Initialize $\mathbf{b}^{(0)}$ to all zeros and set $t = 0$
- 1. While TERMINATION CRITERION is not satisfied
 - a. Update the vector \mathbf{b} using some intermediate steps:

$$\mathbf{b}^{(t+1)} \leftarrow X^T X \mathbf{b}^{(t)} = \sum_{i=1}^n \tilde{\mathbf{x}}^{(i)} \left(\tilde{\mathbf{x}}^{(i)T} \mathbf{b}^{(t)} \right)$$

$$\mathbf{b}^{(t+1)} \leftarrow \frac{\mathbf{b}^{(t+1)}}{\|\mathbf{b}^{(t+1)}\|_2}$$

- b. Increment t : $t \leftarrow t + 1$

- Output: $\mathbf{b}^{(t)}$, the eigenvector corresponding to the largest eigenvalue of A !

Power Iteration

- Input: $A = X^T X$
- Initialize $\mathbf{b}^{(0)}$ to all zeros and set $t = 0$
- 1. While TERMINATION CRITERION is not satisfied
 - a. Update the vector \mathbf{b} using some intermediate steps:

$$\mathbf{b}^{(t+1)} \leftarrow X^T X \mathbf{b}^{(t)} = \sum_{i=1}^n \tilde{\mathbf{x}}^{(i)} \left(\beta_i^{(t)} \right)$$

$$\mathbf{b}^{(t+1)} \leftarrow \frac{\mathbf{b}^{(t+1)}}{\|\mathbf{b}^{(t+1)}\|_2}$$

- b. Increment t : $t \leftarrow t + 1$
- Output: $\mathbf{b}^{(t)}$, the eigenvector corresponding to the largest eigenvalue of A !

$\mathbf{b}^{(t)}$ $O(k)$ comm.	Workers	$\begin{bmatrix} \leftarrow \tilde{\mathbf{x}}^{(1)T} \rightarrow \\ \leftarrow \tilde{\mathbf{x}}^{(4)T} \rightarrow \\ \vdots \end{bmatrix}$	$\begin{bmatrix} \leftarrow \tilde{\mathbf{x}}^{(2)T} \rightarrow \\ \leftarrow \tilde{\mathbf{x}}^{(3)T} \rightarrow \\ \vdots \end{bmatrix}$	$\begin{bmatrix} \leftarrow \tilde{\mathbf{x}}^{(5)T} \rightarrow \\ \leftarrow \tilde{\mathbf{x}}^{(7)T} \rightarrow \\ \vdots \end{bmatrix}$	$O(nk)$ distributed storage (total)	
	Map	$\beta_i^{(t)} = \tilde{\mathbf{x}}^{(i)T} \mathbf{b}^{(t)}$	$\beta_i^{(t)} = \tilde{\mathbf{x}}^{(i)T} \mathbf{b}^{(t)}$	$\beta_i^{(t)} = \tilde{\mathbf{x}}^{(i)T} \mathbf{b}^{(t)}$	$O(nk)$ distributed work (total)	$O(n)$ distributed storage (total)
	Map	$\tilde{\mathbf{x}}^{(i)} \beta_i^{(t)}$	$\tilde{\mathbf{x}}^{(i)} \beta_i^{(t)}$	$\tilde{\mathbf{x}}^{(i)} \beta_i^{(t)}$	$O(nk)$ distributed work (total)	$O(k)$ local storage
	Reduce	$\mathbf{b}^{(t+1)} = \sum_{i=1}^n \tilde{\mathbf{x}}^{(i)} \beta_i^{(t)} / \left\ \sum_{i=1}^n \tilde{\mathbf{x}}^{(i)} \beta_i^{(t)} \right\ _2$			$O(k)$ local work	$O(k)$ local storage

Distributed Power Iteration

$\mathbf{b}^{(t)}$ $O(k)$ comm.	Workers	$\begin{bmatrix} \leftarrow \tilde{\mathbf{x}}^{(1)T} \rightarrow \\ \leftarrow \tilde{\mathbf{x}}^{(4)T} \rightarrow \\ \vdots \end{bmatrix}$	$\begin{bmatrix} \leftarrow \tilde{\mathbf{x}}^{(2)T} \rightarrow \\ \leftarrow \tilde{\mathbf{x}}^{(3)T} \rightarrow \\ \vdots \end{bmatrix}$	$\begin{bmatrix} \leftarrow \tilde{\mathbf{x}}^{(5)T} \rightarrow \\ \leftarrow \tilde{\mathbf{x}}^{(7)T} \rightarrow \\ \vdots \end{bmatrix}$	$O(nk)$ distributed storage (total)	
	Map	$\beta_i^{(t)} = \tilde{\mathbf{x}}^{(i)T} \mathbf{b}^{(t)}$	$\beta_i^{(t)} = \tilde{\mathbf{x}}^{(i)T} \mathbf{b}^{(t)}$	$\beta_i^{(t)} = \tilde{\mathbf{x}}^{(i)T} \mathbf{b}^{(t)}$	$O(nk)$ distributed work (total)	$O(n)$ distributed storage (total)
	Map	$\tilde{\mathbf{x}}^{(i)} \beta_i^{(t)}$	$\tilde{\mathbf{x}}^{(i)} \beta_i^{(t)}$	$\tilde{\mathbf{x}}^{(i)} \beta_i^{(t)}$	$O(nk)$ distributed work (total)	$O(k)$ local storage
	Reduce	$\mathbf{b}^{(t+1)} = \sum_{i=1}^n \tilde{\mathbf{x}}^{(i)} \beta_i^{(t)} / \left\ \sum_{i=1}^n \tilde{\mathbf{x}}^{(i)} \beta_i^{(t)} \right\ _2$			$O(k)$ local work	$O(k)$ local storage

Disclaimer: to find more than one eigenvector, we need to use similar (but slightly more complicated methods), e.g., PySpark's MLlib uses Krylov subspace methods

Background: Empirical Risk Minimization

- A common framework for supervised learning
- Given:
 - some labelled training dataset $\mathcal{D} = \{(\mathbf{x}^{(i)}, y^{(i)})\}_{i=1}^n$
 - a loss function $\ell: \mathcal{Y} \times \mathcal{Y} \rightarrow \mathbb{R}$
 - a hypothesis class or set of functions \mathcal{F}

the goal is to find

$$\hat{f} = \operatorname{argmin}_{f \in \mathcal{F}} \sum_{i=1}^n \ell(f(\mathbf{x}^{(i)}), y^{(i)})$$

with the hope that

$$\mathbb{E}_{p(\mathbf{x}, y)}[\ell(f(\mathbf{x}), y)] \approx \sum_{i=1}^n \ell(f(\mathbf{x}^{(i)}), y^{(i)})$$

Background: Binary Classification

- A type of supervised learning
- Given:
 - some labelled training dataset $\mathcal{D} = \{(\mathbf{x}^{(i)}, y^{(i)})\}_{i=1}^n$
 - a loss function $\ell: \mathcal{Y} \times \mathcal{Y} \rightarrow \mathbb{R}$ where $\mathcal{Y} = \{0,1\}$
 - a hypothesis class or set of functions \mathcal{F}

the goal is to find

$$\hat{f} = \operatorname{argmin}_{f \in \mathcal{F}} \sum_{i=1}^n \ell(f(\mathbf{x}^{(i)}), y^{(i)})$$

Background: Binary Classification w/ 0-1 Loss

- A type of supervised learning
- Given:
 - some labelled training dataset $\mathcal{D} = \{(\mathbf{x}^{(i)}, y^{(i)})\}_{i=1}^n$
 - $\ell(y, y') = \mathbb{1}(y \neq y')$ for $y, y' \in \{0, 1\}$
 - a hypothesis class or set of functions \mathcal{F}

the goal is to find

$$\hat{f} = \operatorname{argmin}_{f \in \mathcal{F}} \sum_{i=1}^n \ell(f(\mathbf{x}^{(i)}), y^{(i)})$$

- This loss function is difficult to optimize (non-convex)...

A Probabilistic Approach to Binary Classification

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

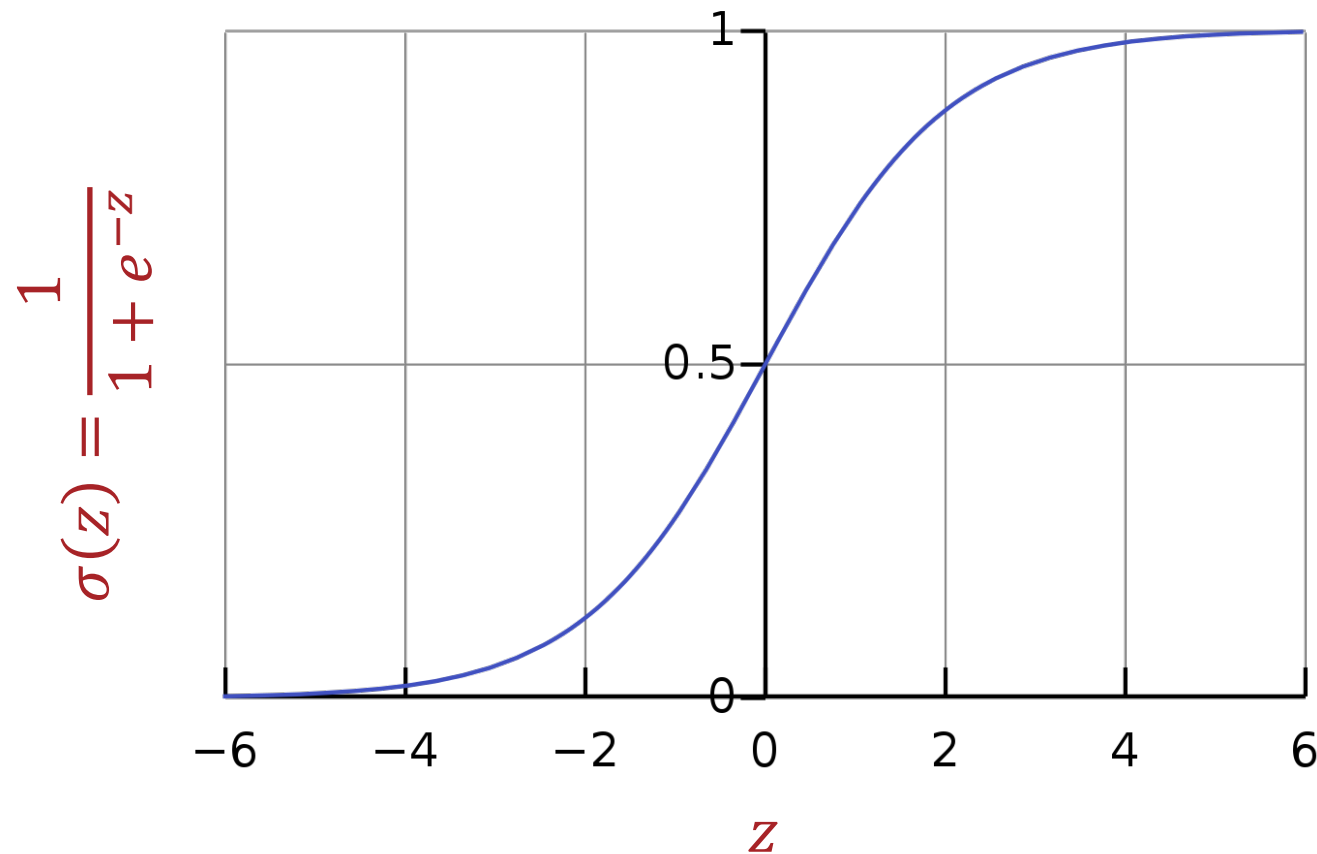
- **Assume**

$$\begin{aligned} P(Y = 1|\mathbf{x}) &= \sigma(\mathbf{w}^T \mathbf{x}) = \frac{1}{1 + \exp(-\mathbf{w}^T \mathbf{x})} \\ &= \frac{\exp(\mathbf{w}^T \mathbf{x})}{\exp(\mathbf{w}^T \mathbf{x}) + 1} \end{aligned}$$

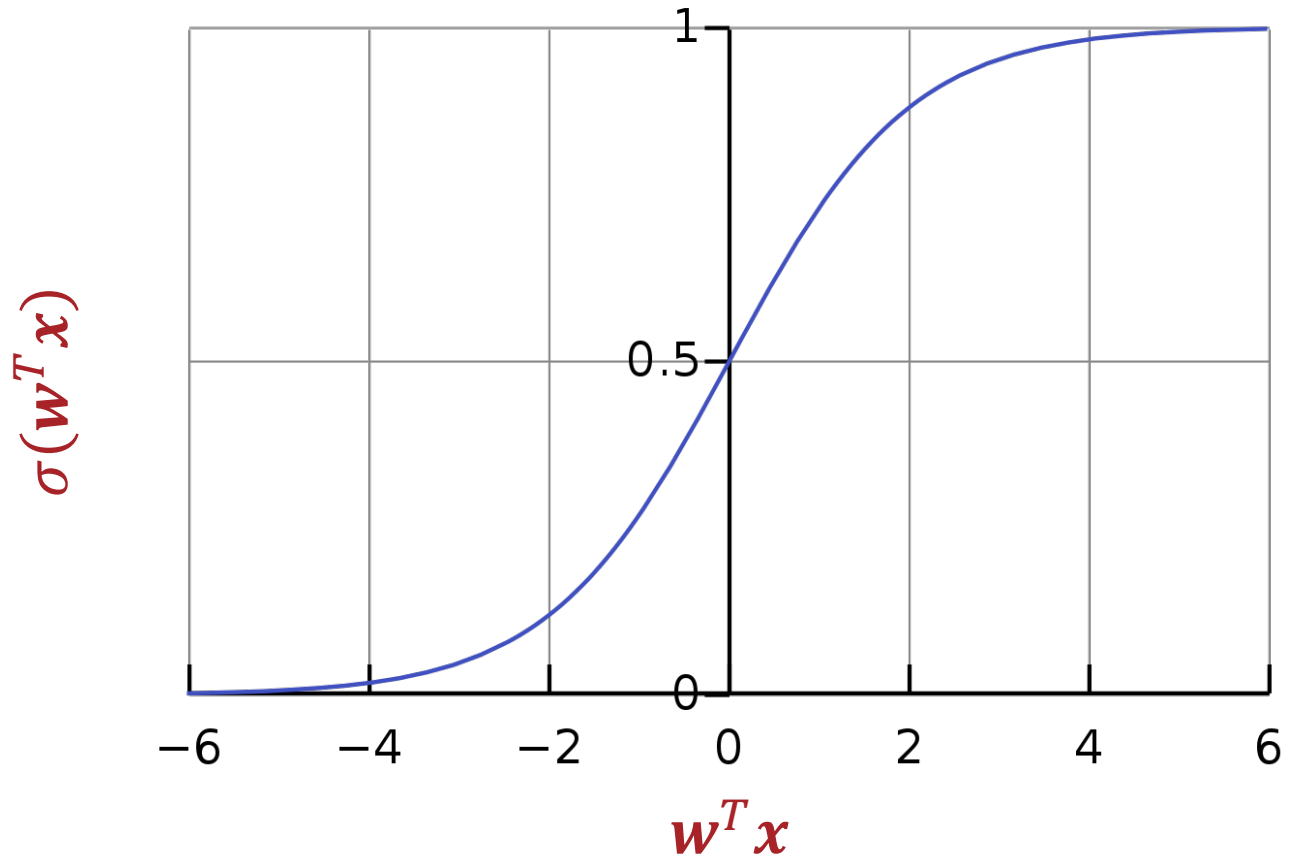
- This implies two useful facts:

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

Logistic Function



Why use the Logistic Function?



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

Logistic Regression Decision Boundary

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

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

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

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

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

$$0 \leq \mathbf{w}^T \mathbf{x}$$

Why 1/2?

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

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

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

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

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

$$0 \leq \mathbf{w}^T \mathbf{x}$$

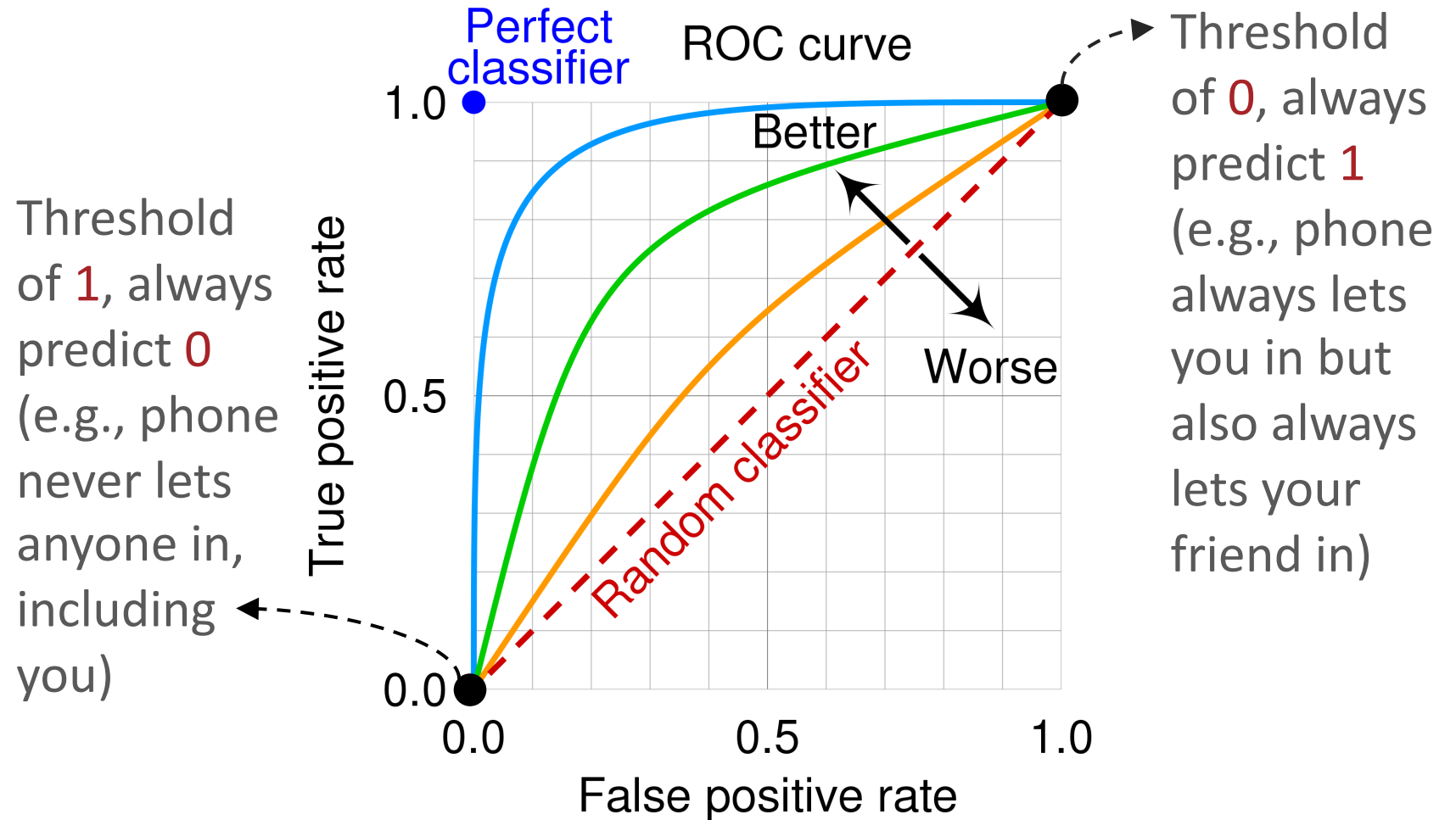
Aside: Different Kinds of Errors

		y	
		0	1
y'	0	True negative (TN)	False negative (FN)
	1	False positive (FP)	True positive (TP)

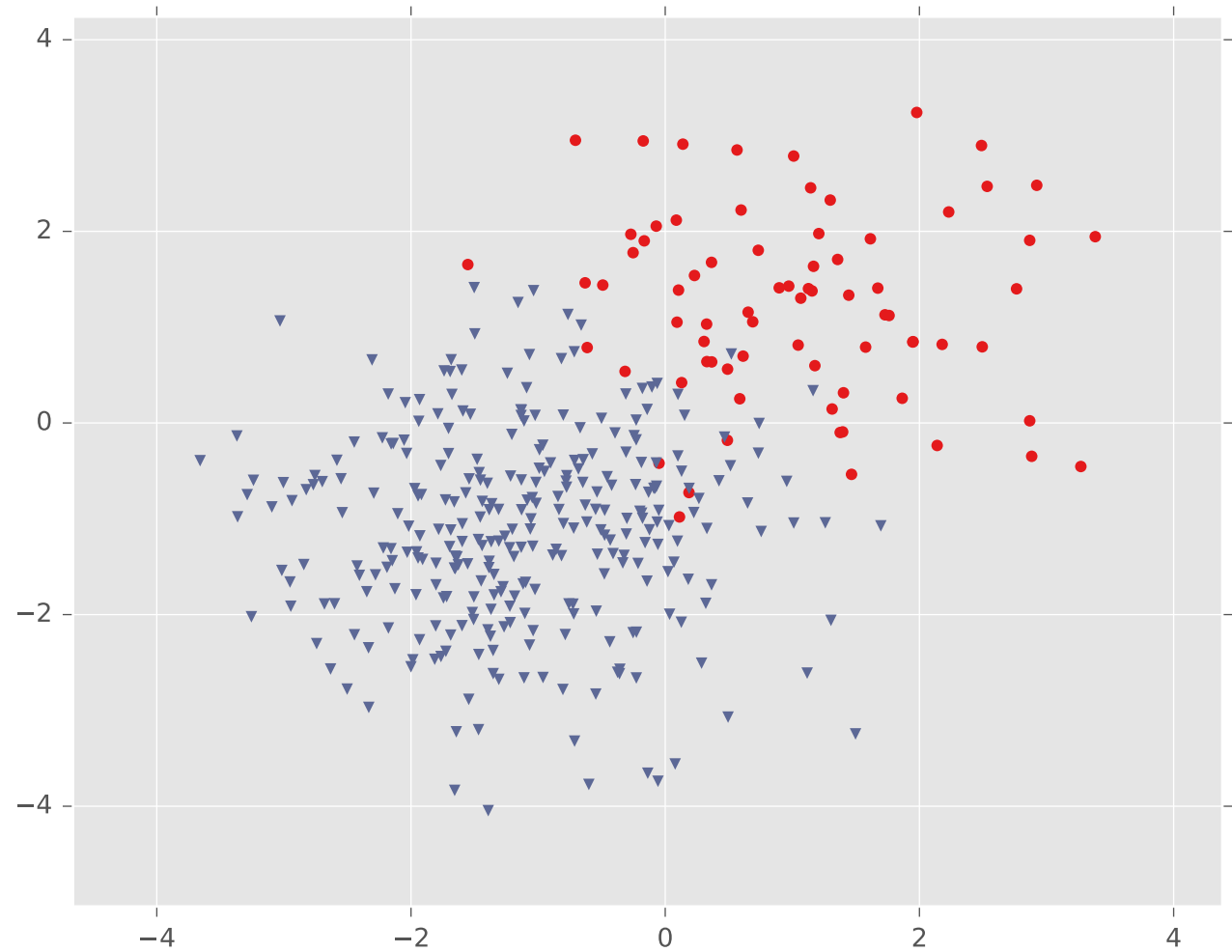
- Depending on the setting, these kinds of errors can have wildly different impacts/costs, e.g., facial recognition for unlocking phones:
 - False negative = minor nuisance, just try again
 - False positive = major security breach
 - Using a threshold $> 1/2$ will tend to decrease the number of false positives

Aside: Receiver Operating Characteristic (ROC) Curves

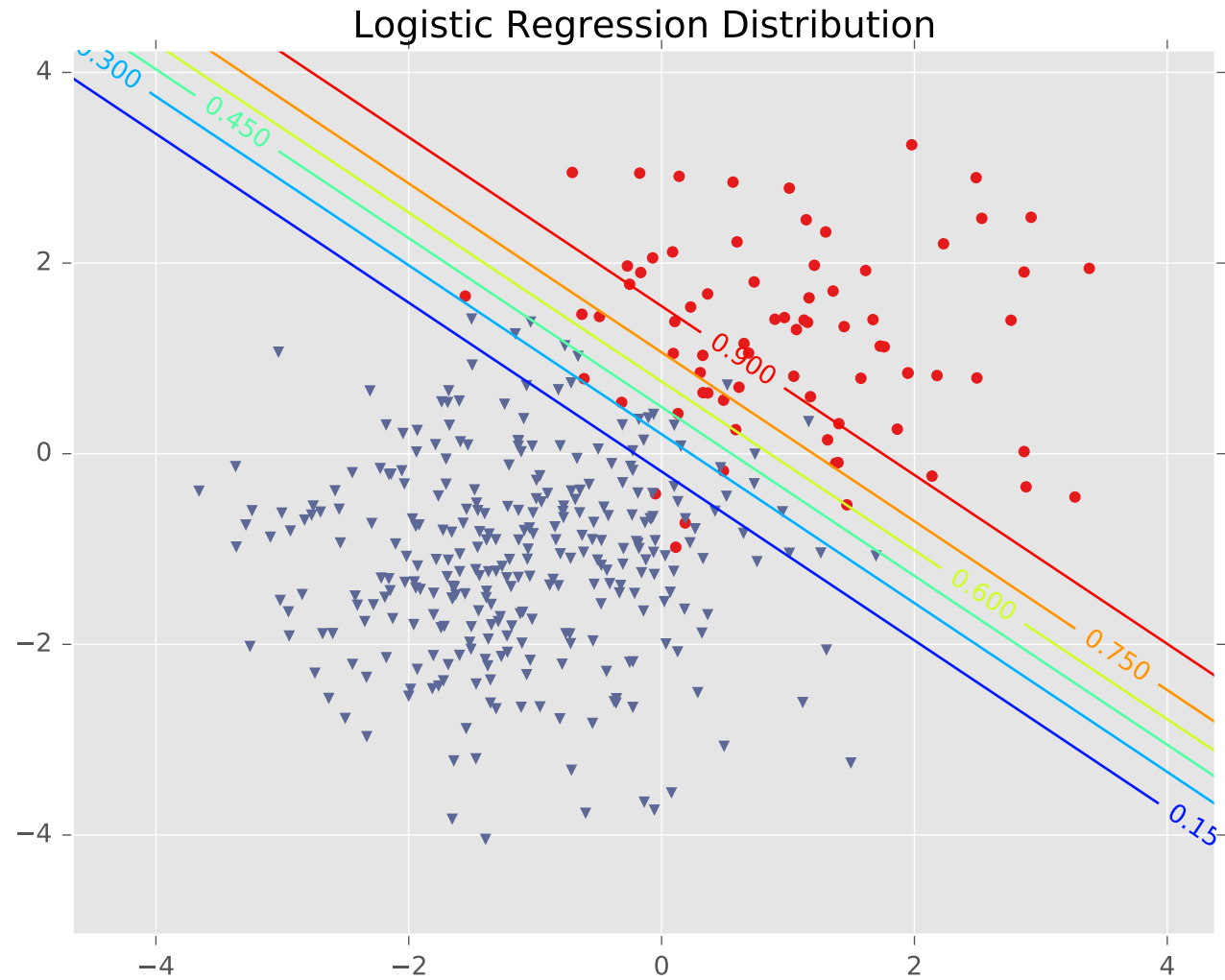
- Plots true positive rate, $TPR = TP / (TP + FN)$, against false positive rate, $FPR = FP / (FP + TN)$



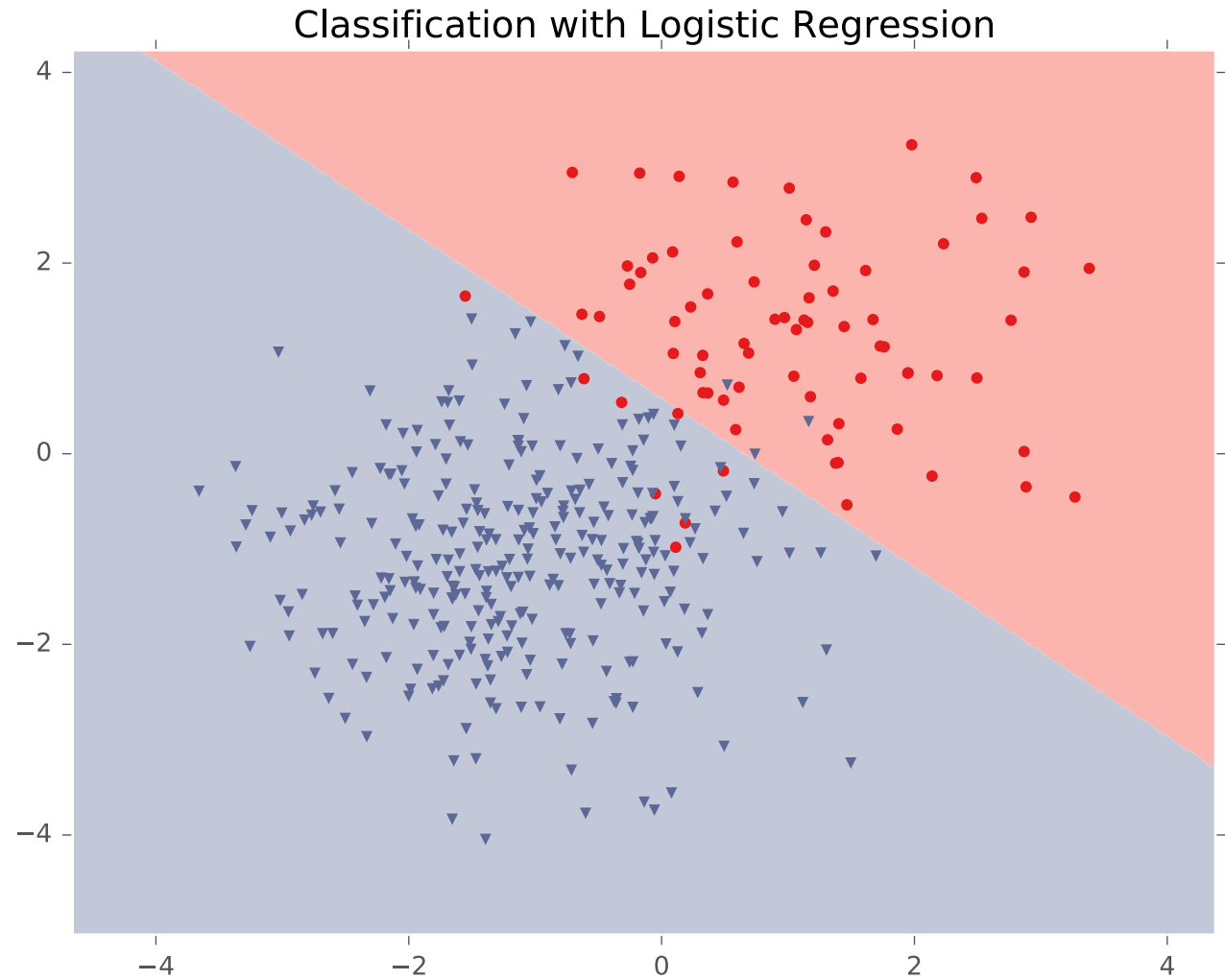
Logistic Regression: Decision Boundary



Logistic Regression: Decision Boundary



Logistic Regression: Decision Boundary



Logistic Regression: Setting the Parameters

- Goal: find the \mathbf{w} that *maximizes* the (conditional) probability of the training dataset:

$$\prod_{i=1}^n P(y^{(i)} | \mathbf{x}^{(i)}, \mathbf{w})$$

- This is equivalent to finding the \mathbf{w} that *minimizes* the negative log of this probability:

$$\begin{aligned} L_{\mathcal{D}}(\mathbf{w}) &= -\log \prod_{i=1}^n P(y^{(i)} | \mathbf{x}^{(i)}, \mathbf{w}) \\ &= -\log \prod_{i=1}^n P(Y = 1 | \mathbf{x}^{(i)}, \mathbf{w})^{y^{(i)}} \left(P(Y = 0 | \mathbf{x}^{(i)}, \mathbf{w}) \right)^{1-y^{(i)}} \\ &= - \underbrace{\sum_{i=1}^n y^{(i)} \log P(Y = 1 | \mathbf{x}^{(i)}, \mathbf{w}) + (1 - y^{(i)}) \log P(Y = 0 | \mathbf{x}^{(i)}, \mathbf{w})}_{\text{log loss or cross-entropy loss}} \end{aligned}$$

Logistic Regression as Empirical Risk Minimization

- A type of supervised learning
- Given:
 - some labelled training dataset $\mathcal{D} = \{(\mathbf{x}^{(i)}, y^{(i)})\}_{i=1}^n$
 - the log loss
 - \mathcal{F} = the set of all linear decision boundaries

the goal is to find

$$\begin{aligned} & \operatorname{argmin}_{\mathbf{w}} - \sum_{i=1}^n (y^{(i)} \log P(Y = 1 | \mathbf{x}^{(i)}, \mathbf{w})) + (1 - y^{(i)}) \log P(Y = 0 | \mathbf{x}^{(i)}, \mathbf{w}) \\ &= \operatorname{argmin}_{\mathbf{w}} - \sum_{i=1}^n y^{(i)} \log \frac{P(Y = 1 | \mathbf{x}^{(i)}, \mathbf{w})}{P(Y = 0 | \mathbf{x}^{(i)}, \mathbf{w})} + \log P(Y = 0 | \mathbf{x}^{(i)}, \mathbf{w}) \\ &= \operatorname{argmin}_{\mathbf{w}} - \sum_{i=1}^n y^{(i)} \mathbf{w}^T \mathbf{x}^{(i)} - \log (1 + \exp(\mathbf{w}^T \mathbf{x}^{(i)})) \end{aligned}$$

Logistic Regression as Empirical Risk Minimization

$$L_{\mathcal{D}}(\mathbf{w}) = - \sum_{i=1}^n y^{(i)} \mathbf{w}^T \mathbf{x}^{(i)} - \log \left(1 + \exp(\mathbf{w}^T \mathbf{x}^{(i)}) \right)$$

$$\nabla_{\mathbf{w}} L_{\mathcal{D}}(\mathbf{w}) = - \sum_{i=1}^n y^{(i)} \nabla_{\mathbf{w}} \mathbf{w}^T \mathbf{x}^{(i)} - \nabla_{\mathbf{w}} \log \left(1 + \exp(\mathbf{w}^T \mathbf{x}^{(i)}) \right)$$

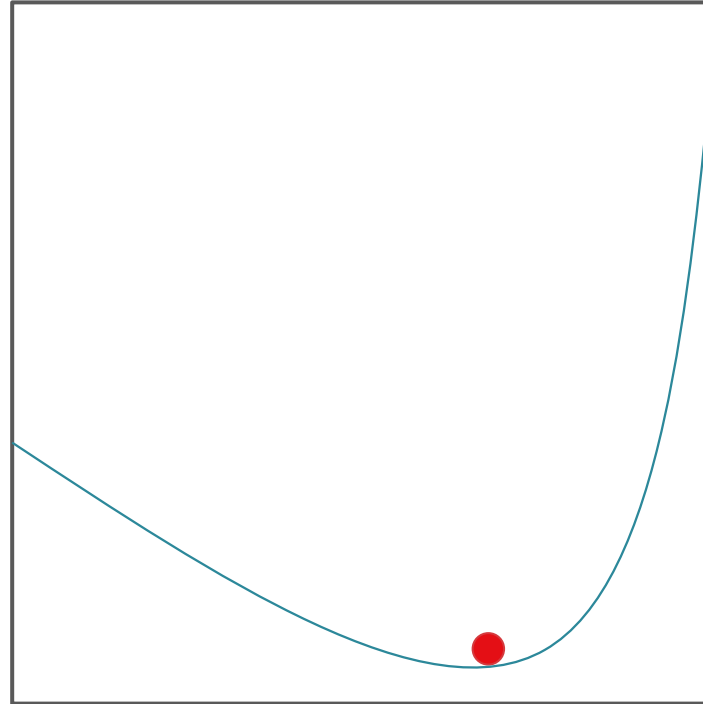
$$= - \sum_{i=1}^n y^{(i)} \mathbf{x}^{(i)} - \frac{\exp(\mathbf{w}^T \mathbf{x}^{(i)})}{1 + \exp(\mathbf{w}^T \mathbf{x}^{(i)})} \mathbf{x}^{(i)}$$

$$= \sum_{i=1}^n \mathbf{x}^{(i)} \frac{1}{\exp(-\mathbf{w}^T \mathbf{x}^{(i)}) + 1} - y^{(i)} \mathbf{x}^{(i)}$$

$$= \sum_{i=1}^n \mathbf{x}^{(i)} (\sigma(\mathbf{w}^T \mathbf{x}^{(i)}) - y^{(i)})$$

Recall: Gradient Descent

- An iterative method for minimizing functions
- Requires the gradient to exist everywhere



- Good news: the log loss for logistic regression is convex!

Gradient Descent for Logistic Regression

- Input: $\mathcal{D} = \{(\mathbf{x}^{(i)}, y^{(i)})\}_{i=1}^n, \alpha$
- 1. Initialize $\mathbf{w}^{(0)}$ to all zeros and set $t = 0$
- 2. While TERMINATION CRITERION is not satisfied
 - a. Compute the gradient:

$$\nabla_{\mathbf{w}} L_{\mathcal{D}}(\mathbf{w}^{(t)}) = \sum_{i=1}^n \mathbf{x}^{(i)} \left(\text{logit}(\mathbf{w}^{(t)T} \mathbf{x}^{(i)}) - y^{(i)} \right)$$

- b. Update \mathbf{w} : $\mathbf{w}^{(t+1)} \leftarrow \mathbf{w}^{(t)} - \alpha \nabla_{\mathbf{w}} L_{\mathcal{D}}(\mathbf{w}^{(t)})$
 - c. Increment t : $t \leftarrow t + 1$
- Output: $\mathbf{w}^{(t)}$

Idea: distribute $\mathbf{x}^{(i)}$ and compute summands in parallel

- Input: $\mathcal{D} = \{(\mathbf{x}^{(i)}, y^{(i)})\}_{i=1}^n, \alpha$
- 1. Initialize $\mathbf{w}^{(0)}$ to all zeros and set $t = 0$
- 2. While TERMINATION CRITERION is not satisfied
 - a. Compute the gradient:

$$\nabla_{\mathbf{w}} L_{\mathcal{D}}(\mathbf{w}^{(t)}) = \sum_{i=1}^n \mathbf{x}^{(i)} \left(\text{logit}(\mathbf{w}^{(t)T} \mathbf{x}^{(i)}) - y^{(i)} \right)$$

- b. Update \mathbf{w} : $\mathbf{w}^{(t+1)} \leftarrow \mathbf{w}^{(t)} - \alpha \nabla_{\mathbf{w}} L_{\mathcal{D}}(\mathbf{w}^{(t)})$
 - c. Increment t : $t \leftarrow t + 1$
- Output: $\mathbf{w}^{(t)}$

$\mathbf{w}^{(t)}$	Workers	$\begin{bmatrix} \leftarrow \mathbf{x}^{(1)T} & \rightarrow \\ \leftarrow \mathbf{x}^{(4)T} & \rightarrow \\ \vdots & \vdots & \vdots \end{bmatrix}$	$\begin{bmatrix} \leftarrow \mathbf{x}^{(2)T} & \rightarrow \\ \leftarrow \mathbf{x}^{(3)T} & \rightarrow \\ \vdots & \vdots & \vdots \end{bmatrix}$	$\begin{bmatrix} \leftarrow \mathbf{x}^{(5)T} & \rightarrow \\ \leftarrow \mathbf{x}^{(7)T} & \rightarrow \\ \vdots & \vdots & \vdots \end{bmatrix}$	$O(nk)$ distributed storage (total)	
	Map	$\mathbf{x}^{(i)} \left(\sigma \left(\mathbf{w}^{(t)T} \mathbf{x}^{(i)} \right) - y^{(i)} \right)$	$\mathbf{x}^{(i)} \left(\sigma \left(\mathbf{w}^{(t)T} \mathbf{x}^{(i)} \right) - y^{(i)} \right)$	$\mathbf{x}^{(i)} \left(\sigma \left(\mathbf{w}^{(t)T} \mathbf{x}^{(i)} \right) - y^{(i)} \right)$	$O(nk)$ distributed work (total)	$O(k)$ local storage
	Reduce	$\mathbf{w}^{(t+1)} \leftarrow \mathbf{w}^{(t)} - \frac{\alpha}{n\sqrt{t}} \sum_{i=1}^n \mathbf{x}^{(i)} \left(\sigma \left(\mathbf{w}^{(t)T} \mathbf{x}^{(i)} \right) - y^{(i)} \right)$			$O(k)$ local work	$O(k)$ local storage

$O(k)$ communication

Distributed Gradient Descent for Logistic Regression

What do linear regression, PCA and logistic regression all have in common?

- They're all linear methods...
- ... that can be kernelized to achieve nonlinear models!
 - What is a kernel and how can we apply it at scale?
Come back on Tuesday and find out!

Key Takeaways

- Distributed PCA
 - Very similar to distributed linear regression:
 - If k is small, simply distribute the storage and computation of $X^T X$
 - If k is large, use iterative methods (e.g., power iteration) to compute eigenvector-eigenvalue pairs
- Distributed logistic regression
 - Distribute the summands of gradient descent and compute in parallel

Positive Semidefinite Matrices

- For a symmetric matrix $A \in \mathbb{R}^{n \times n}$, the following statements are equivalent:
 1. A is positive semidefinite
 2. $\forall x \in \mathbb{R}^n, x^T A x \geq 0$
 3. All eigenvalues of A are greater than or equal to 0
 4. $\exists U$ s.t. $A = U^T U$

Block Matrices

- Sometimes, it may be convenient to express a matrix A as comprised of several smaller matrices:

$$A = \begin{bmatrix} A_{11} & A_{12} \\ A_{21} & A_{22} \end{bmatrix}$$

- The Schur complement of A_{22} is defined as

$$C = A_{11} - A_{12}A_{22}^{-1}A_{21}$$

- Using the Schur complement, we can succinctly express quantities like

$$\det(A) = \det(A_{22}) * \det(C)$$

$$A^{-1} = \begin{bmatrix} I & 0 \\ -A_{22}^{-1}A_{21} & I \end{bmatrix} \begin{bmatrix} C & 0 \\ 0 & A_{22}^{-1} \end{bmatrix} \begin{bmatrix} I & -A_{12}A_{22}^{-1} \\ 0 & I \end{bmatrix}$$

Orthogonal Matrices

- An orthogonal matrix is a real symmetric matrix whose columns and rows are orthonormal vectors.
- If A is an orthogonal matrix, then this implies that

$$A^T A = A A^T = I$$

or equivalently

$$A^T = A^{-1}$$

Singular Value Decomposition

- Any matrix real $A \in \mathbb{R}^{m \times n}$ can be expressed as the product of three matrices:

$$A = U\Sigma V^T$$

where

- $U \in \mathbb{R}^{m \times m}$ is an orthogonal matrix whose columns are the eigenvectors of AA^T
- $V \in \mathbb{R}^{n \times n}$ is also an orthogonal matrix whose columns are the eigenvectors of $A^T A$
- $\Sigma \in \mathbb{R}^{m \times n}$ is a diagonal rectangular matrix whose non-zero elements are the square roots of the eigenvalues of AA^T (equivalently, $A^T A$)