

# 10-701: Introduction to Machine Learning

## Lecture 7 – Logistic Regression

Henry Chai

2/7/24

# Front Matter

- Announcements:
  - HW2 released 2/7 (today!), due on 2/16 at 11:59 PM
- Recommended Readings:
  - Murphy, [Section 8.1 - 8.3](#)

# Recall: 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  $P^*$

# Recipe for Naïve Bayes

$$\vec{\theta} = [\theta_{1,0}, \theta_{1,1}, \theta_{2,0}, \theta_{2,1}, \dots]$$

- Define a model and model parameters
  - Make the Naïve Bayes assumption
  - Assume independent, identically distributed (iid) data
  - Parameters:  $\pi = P(Y = 1)$ ,  $\theta_{d,y} = P(X_d = 1|Y = y)$
- Write down an objective function
  - Maximize the log-likelihood
- Optimize the objective w.r.t. the model parameters
  - Solve in *closed form*: take partial derivatives, set to 0 and solve

# Bernoulli Naïve Bayes

- Binary label

- $Y \sim \text{Bernoulli}(\pi)$

→  $\hat{\pi} = N_{Y=1} / N$

- $N = \#$  of data points
    - $N_{Y=1} = \#$  of data points with label 1

- Binary features

- $X_d | Y = y \sim \text{Bernoulli}(\theta_{d,y})$

→  $\hat{\theta}_{d,y} = N_{Y=y, X_d=1} / N_{Y=y}$

- $N_{Y=y} = \#$  of data points with label  $y$
    - $N_{Y=y, X_d=1} = \#$  of data points with label  $y$  and feature  $X_d = 1$

What if some  
Word-Label  
pairs don't  
appear in our  
training data?  
Making  
Predictions

- Given a test data point  $x' = [x'_1, \dots, x'_D]^T$

$$P(Y=1|x') \geq P(Y=0|x')$$

$$P(Y=1|x') \propto P(x'|Y=1)P(Y=1)$$

$$= \frac{\hat{\pi}}{\pi} \left( \prod_{d=1}^D \hat{\theta}_{d,1}^{x'_d} (1 - \hat{\theta}_{d,1})^{1-x'_d} \right)$$

$$P(Y=0|x') \propto (1 - \hat{\pi}) \left( \prod_{d=1}^D \hat{\theta}_{d,0}^{x'_d} (1 - \hat{\theta}_{d,0})^{1-x'_d} \right)$$



What if some  
Word-Label  
pair never  
appears in our  
training data?

$x_1$ ("hat")	$x_2$ ("cat")	$x_3$ ("dog")	$x_4$ ("fish")	$x_5$ ("mom")	$x_6$ ("dad")	$y$ (Dr. Seuss)
1	1	0	0	0	0	1
0	0	1	0	0	0	0
0	0	0	1	0	0	1
0	0	0	0	1	0	0

The Cat in the Hat gets a Dog (by ???)

- If some  $\hat{\theta}_{d,y} = 0$  and that word appears in our test data  $\mathbf{x}'$ , then  $P(Y = y|\mathbf{x}') = 0$  even if all the other features in  $\mathbf{x}'$  point to the label being  $y$ !
- The model has been overfit to the training data...
- We can address this with a prior over the parameters!

# Setting the Parameters via MAP

- Binary label
    - $Y \sim \text{Bernoulli}(\pi)$
    - $\hat{\pi} = N_{Y=1} / N$ 
      - $N = \#$  of data points
      - $N_{Y=1} = \#$  of data points with label 1
  - Binary features
    - $X_d | Y = y \sim \text{Bernoulli}(\theta_{d,y})$  and  $\theta_{d,y} \sim \text{Beta}(\alpha, \beta)$
    - $\hat{\theta}_{d,y} = \frac{N_{Y=y, X_d=1} + (\alpha - 1)}{N_{Y=y} + (\alpha - 1) + (\beta - 1)}$ 
      - $N_{Y=y} = \#$  of data points with label  $y$
      - $N_{Y=y, X_d=1} = \#$  of data points with label  $y$  and feature  $X_d = 1$
      - $\alpha$  and  $\beta$  are “pseudocounts” of imagined data points that help avoid zero-probability predictions.
- • Common choice:  $\alpha = \beta = 2$



# Recall: Building a Probabilistic Classifier

- Define a decision rule
  - Given a test data point  $\mathbf{x}'$ , predict its label  $\hat{y}$  using the *posterior distribution*  $P(Y = y|X = \mathbf{x}')$
  - Common choice:  $\hat{y} = \underset{y}{\operatorname{argmax}} P(Y = y|X = \mathbf{x}')$
- Model the posterior distribution
  - Option 1 - Model  $P(Y|X)$  directly as some function of  $X$  (today!)
  - Option 2 - Use Bayes' rule (Monday):

$$\star P(Y|X) = \frac{P(X|Y) P(Y)}{P(X)} \propto P(X|Y) P(Y)$$

# 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**

$$\underline{P(Y = 1|\mathbf{x})} = \text{logit}(\mathbf{w}^T \mathbf{x}) = \frac{1}{1 + \exp(-\mathbf{w}^T \mathbf{x})}$$

$$\boxed{= \frac{\exp(\mathbf{w}^T \mathbf{x})}{\exp(\mathbf{w}^T \mathbf{x}) + 1}}$$

$\frac{\exp(\mathbf{w}^T \mathbf{x})}{\exp(\mathbf{w}^T \mathbf{x}) + 1}$

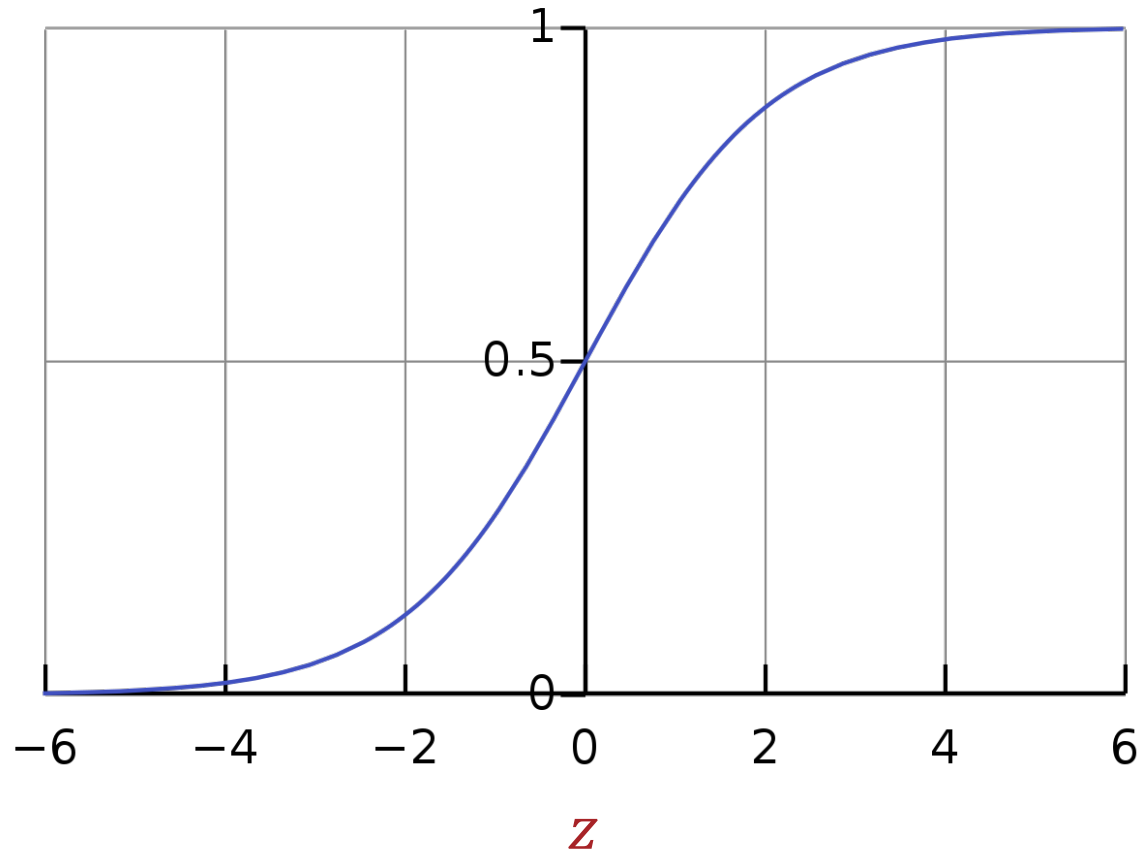
- This implies two useful facts:

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

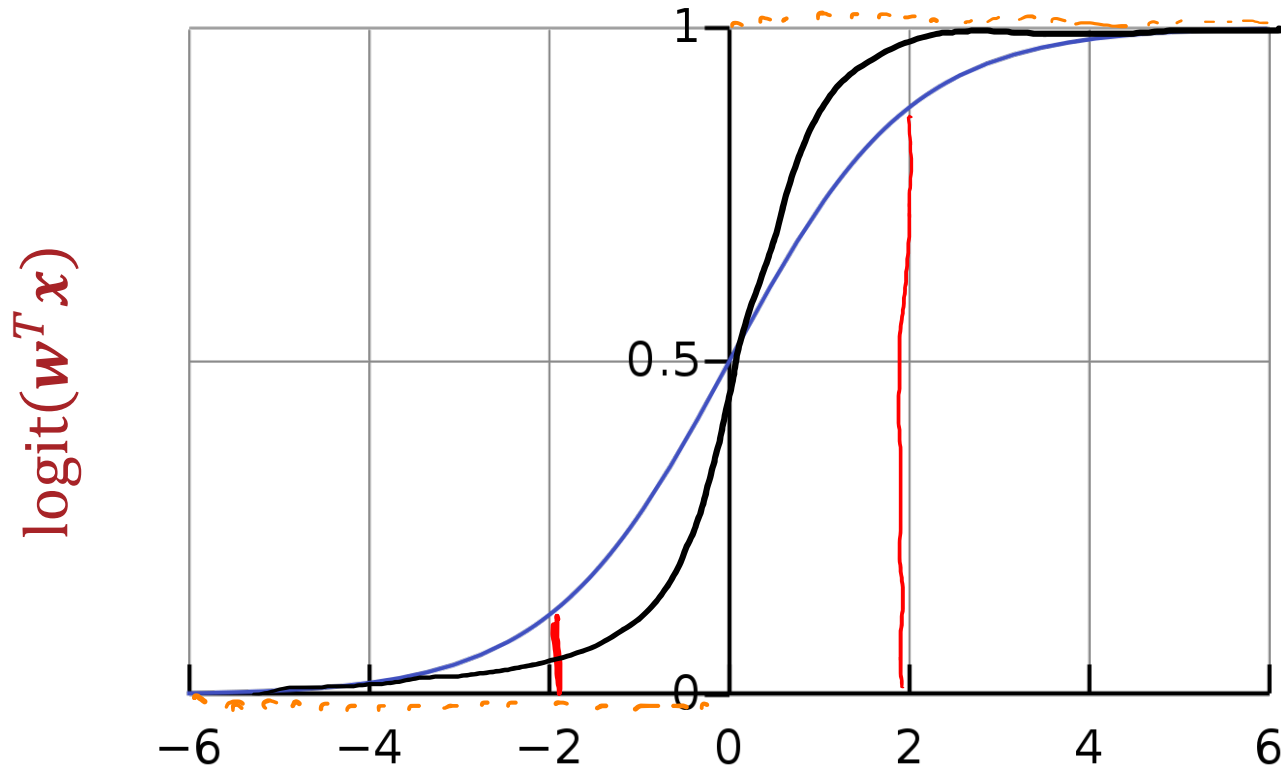
$$2. \frac{P(Y=1|\mathbf{x})}{P(Y=0|\mathbf{x})} = \exp(\mathbf{w}^T \mathbf{x}) \Rightarrow \log \text{ odds} = \mathbf{w}^T \mathbf{x}$$

# Logistic Function

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



# Why use the Logistic Function?



$\text{logit}(w^T x) \in (0, 1) \Rightarrow$  a valid probability  
monotonically increasing  $\Rightarrow$  linear decision boundary  
differentiable everywhere

# Logistic Regression Decision Boundary

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

$$P(Y=1 | \mathbf{x}') = \text{logit}(\hat{\omega}^T \mathbf{x}') = \frac{1}{1 + \exp(-\hat{\omega}^T \mathbf{x}')} \geq \frac{1}{2}$$

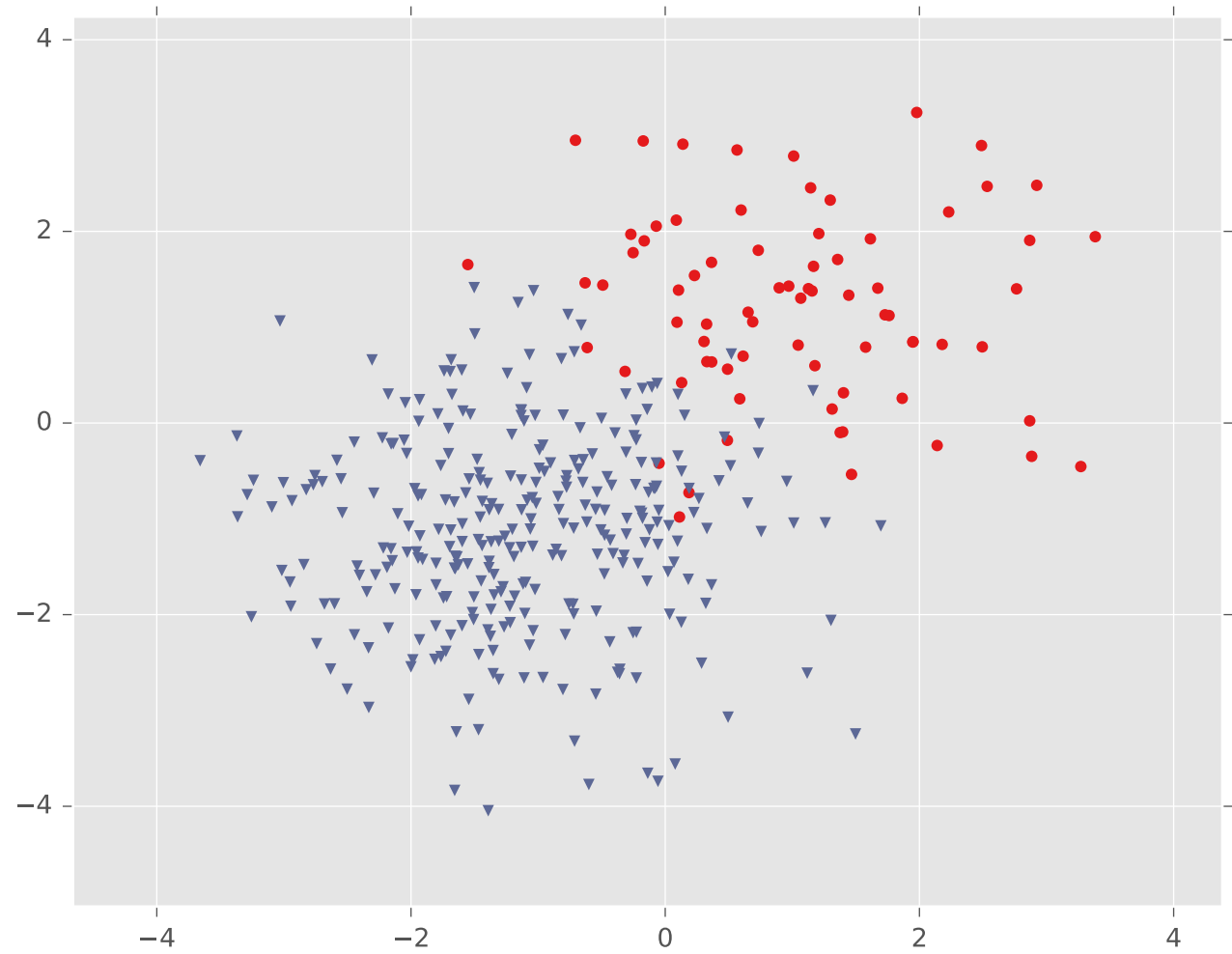
$$\Rightarrow Z \geq 1 + \exp(-\hat{\omega}^T \mathbf{x}')$$

$$\Rightarrow 1 \geq \exp(-\hat{\omega}^T \mathbf{x}')$$

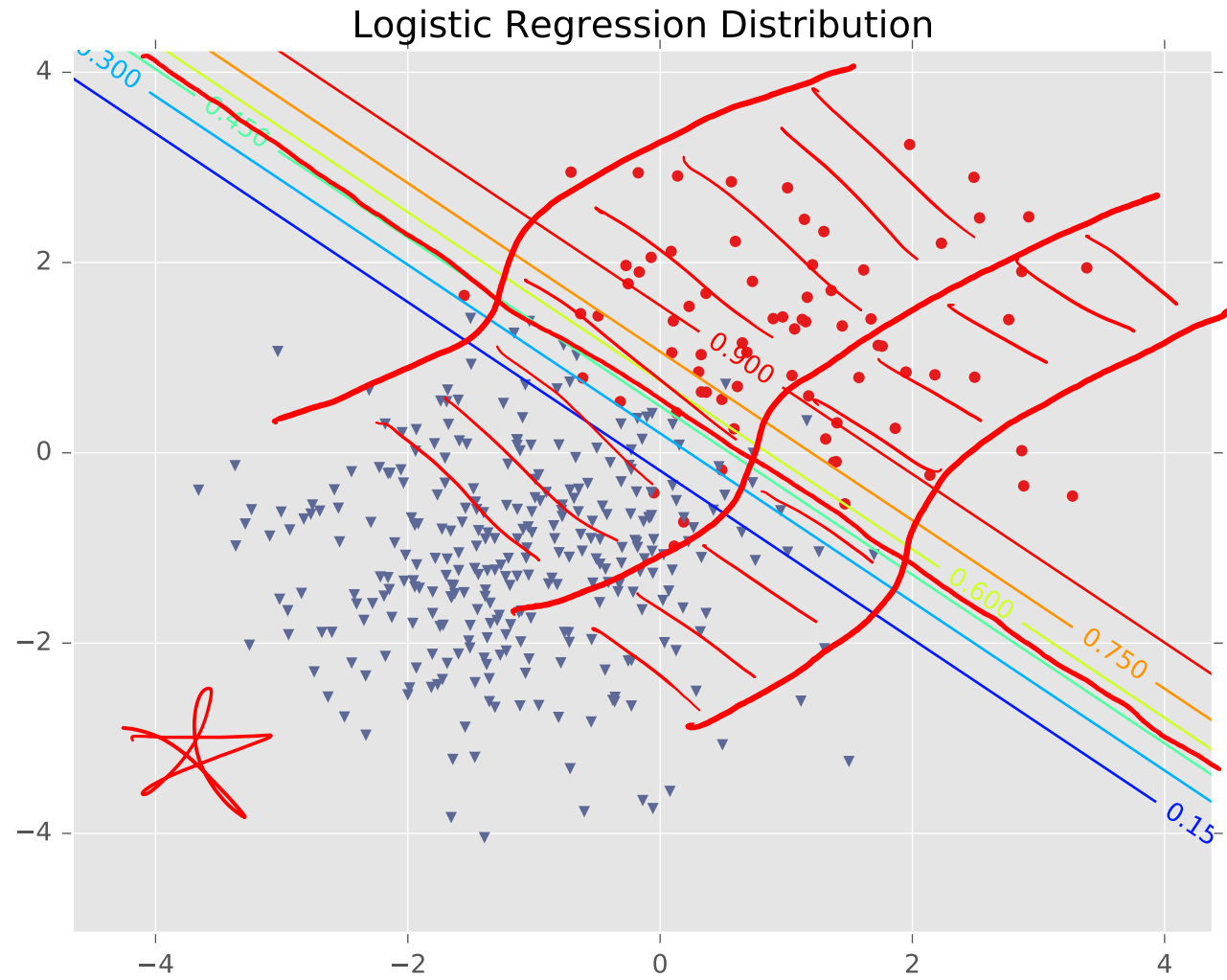
$$\Rightarrow 0 \geq -\hat{\omega}^T \mathbf{x}'$$

$$\Rightarrow \hat{\omega}^T \mathbf{x}' \geq 0$$

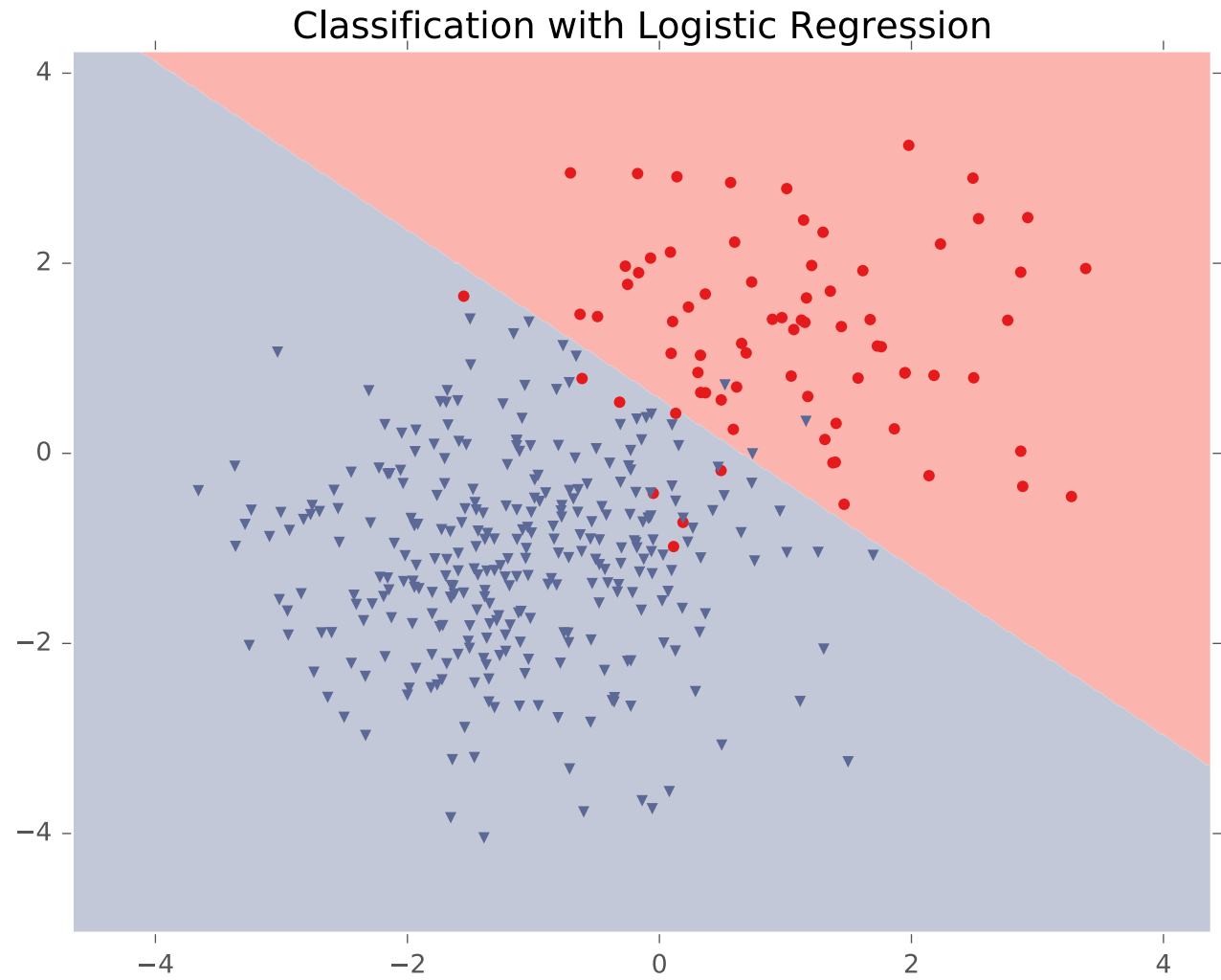
# Logistic Regression Decision Boundary



# Logistic Regression Decision Boundary



# Logistic Regression Decision Boundary





# Recipe for Logistic Regression

- Define a model and model parameters
  - Assume iid data
  - Assume  $P(Y=1|x) = \text{logit}(w^T x)$
  - Parameters  $w = [w_0, w_1, \dots, w_D]^T$
- Write down an objective function
  - Maximum conditional likelihood estimation
  - Minimum negative conditional log-likelihood estimation
- Optimize the objective w.r.t. the model parameters

???

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

$$\log P(D = \{(x^{(1)}, y^{(1)}), (x^{(2)}, y^{(2)}) \dots\})$$

Find  $w$  that minimizes

$$= \log P(\underline{x}^{(1)} \cap \underline{y}^{(1)} \cap \underline{x}^{(2)} \cap \underline{y}^{(2)} \dots)$$

$$\ell_D(w) = -\log P(y^{(1)}, \dots, y^{(N)} | x^{(1)}, \dots, x^{(N)}, w) = -\log \prod_{n=1}^N P(y^{(n)} | x^{(n)}, w)$$

$$\Rightarrow = -\log \prod_{n=1}^N P(Y=1 | x^{(n)}, w)^{y^{(n)}} (P(Y=0 | x^{(n)}, w))^{1-y^{(n)}}$$

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

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

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

$$\rightarrow = -\sum_{n=1}^N y^{(n)} w^T x^{(n)} - \log(1 + \exp(w^T x^{(n)}))$$

# Minimizing the Negative Conditional (log-)Likelihood

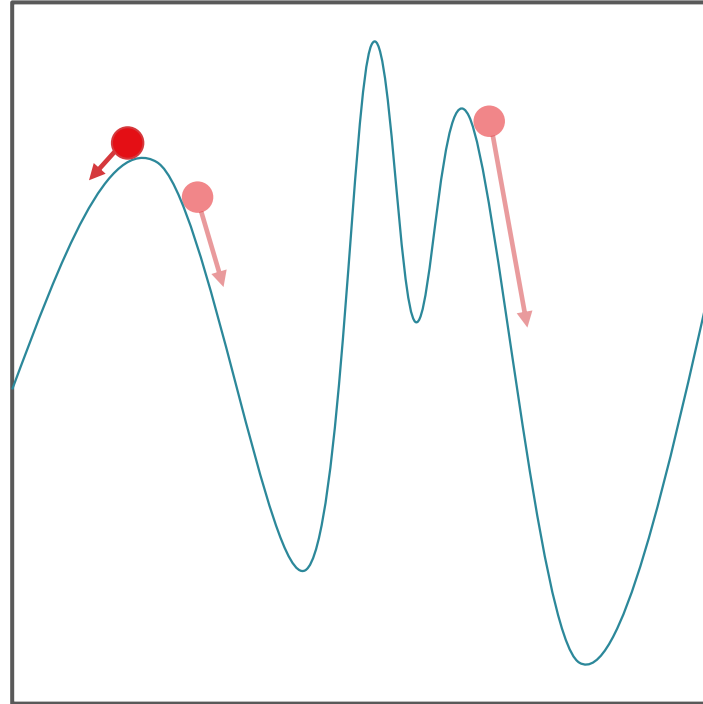
$$w^T x^{(n)} = \sum_{d=1}^D w_d x_d^{(n)}$$

$w_d = ?$

$$\begin{aligned}
 \ell_D(w) &= - \sum_{n=1}^N \left( \frac{y^{(n)} w^T x^{(n)}}{1 + \exp(w^T x^{(n)})} - \log(1 + \exp(w^T x^{(n)})) \right) \\
 \nabla_w \ell_D(w) &= - \sum_{n=1}^N \left( \frac{y^{(n)} x^{(n)}}{1 + \exp(w^T x^{(n)})} - \frac{\exp(w^T x^{(n)})}{1 + \exp(w^T x^{(n)})} x^{(n)} \right) \\
 &= \sum_{n=1}^N x^{(n)} \left( \frac{\exp(w^T x^{(n)})}{1 + \exp(w^T x^{(n)})} - y^{(n)} \right) \\
 &= \sum_{n=1}^N x^{(n)} \left( P(Y=1 | x^{(n)}, w) - y^{(n)} \right)
 \end{aligned}$$

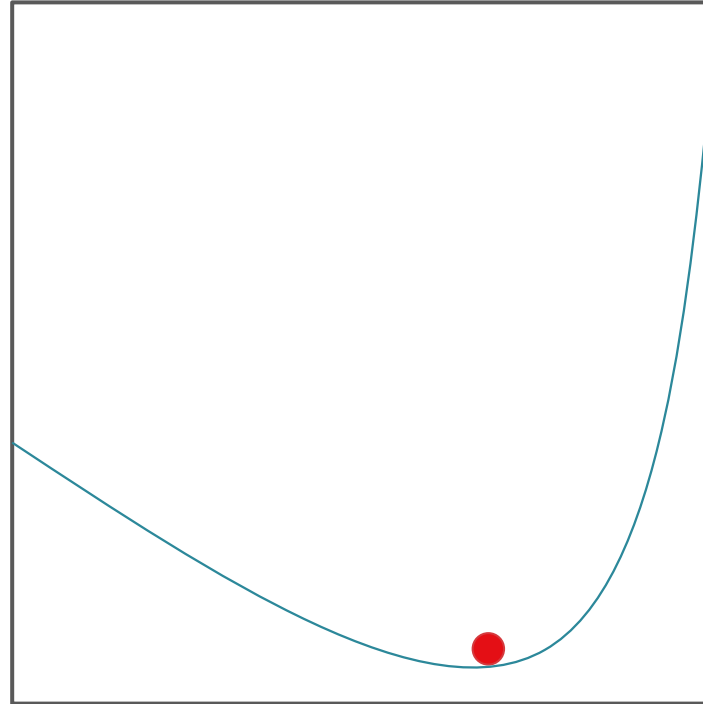
# Recall: Gradient Descent

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



# Recall: Gradient Descent

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



- Good news: the negative conditional log-likelihood, like the squared error, is also convex!

<sup>^</sup>  
strictly

# Gradient Descent

• Input:  $\mathcal{D} = \{(\mathbf{x}^{(n)}, y^{(n)})\}_{n=1}^N, \eta^{(0)}$

1. Initialize  $\mathbf{w}^{(0)}$  to all zeros and set  $t = 0$

2. While TERMINATION CRITERION is not satisfied

a. Compute the gradient:

$$O(ND) \left\{ \nabla_{\mathbf{w}} \ell_{\mathcal{D}}(\mathbf{w}^{(t)}) = \sum_{n=1}^N \mathbf{x}^{(n)} (P(Y = 1 | \mathbf{x}^{(n)}, \mathbf{w}^{(t)}) - y^{(n)}) \right.$$

b. Update  $\mathbf{w}$ :  $\mathbf{w}^{(t+1)} \leftarrow \mathbf{w}^{(t)} - \eta^{(0)} \nabla_{\mathbf{w}} \ell_{\mathcal{D}}(\mathbf{w}^{(t)})$

c. Increment  $t$ :  $t \leftarrow t + 1$

• Output:  $\mathbf{w}^{(t)}$

# Stochastic Gradient Descent

- Input:  $\mathcal{D} = \{(\mathbf{x}^{(n)}, y^{(n)})\}_{n=1}^N, \eta_{SGD}^{(0)}$
- 1. Initialize  $\mathbf{w}^{(0)}$  to all zeros and set  $t = 0$
- 2. While TERMINATION CRITERION is not satisfied
  - a. Randomly sample a data point from  $\mathcal{D}$ ,  $(\mathbf{x}^{(n)}, y^{(n)})$
  - b. Compute the pointwise gradient:  
 $\rightarrow \nabla_{\mathbf{w}} \ell^{(n)}(\mathbf{w}^{(t)}) = \mathbf{x}^{(n)} (P(Y = 1 | \mathbf{x}^{(n)}, \mathbf{w}^{(t)}) - y^{(n)})$
  - c. Update  $\mathbf{w}$ :  $\mathbf{w}^{(t+1)} \leftarrow \mathbf{w}^{(t)} - \eta_{SGD}^{(0)} \nabla_{\mathbf{w}} \ell^{(n)}(\mathbf{w}^{(t)})$
  - d. Increment  $t$ :  $t \leftarrow t + 1$
- Output:  $\mathbf{w}^{(t)}$

# Stochastic Gradient Descent

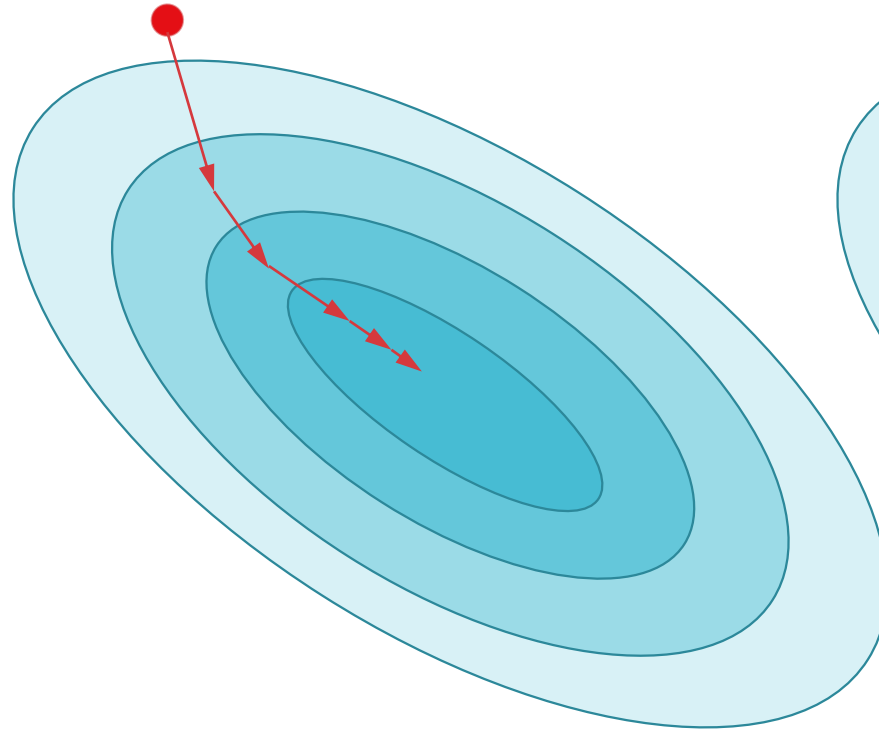
- If the data point is sampled uniformly at random, then the expected value of the pointwise gradient is proportional to the full gradient:

$$\begin{aligned} E \left[ \nabla_{\mathbf{w}} \ell_{\mathbf{x}^{(n)}, \mathbf{y}^{(n)}}(\mathbf{w}^{(t)}) \right] &= \frac{1}{N} \sum_{n=1}^N \nabla_{\mathbf{w}} \ell^{(n)}(\mathbf{w}^{(t)}) \\ &= \frac{1}{N} \sum_{n=1}^N \mathbf{x}^{(n)} (P(Y = 1 | \mathbf{x}^{(n)}, \mathbf{w}^{(t)}) - y^{(n)}) \\ &= \frac{1}{N} \nabla_{\mathbf{w}} \ell_{\mathcal{D}}(\mathbf{w}^{(t)}) \end{aligned}$$

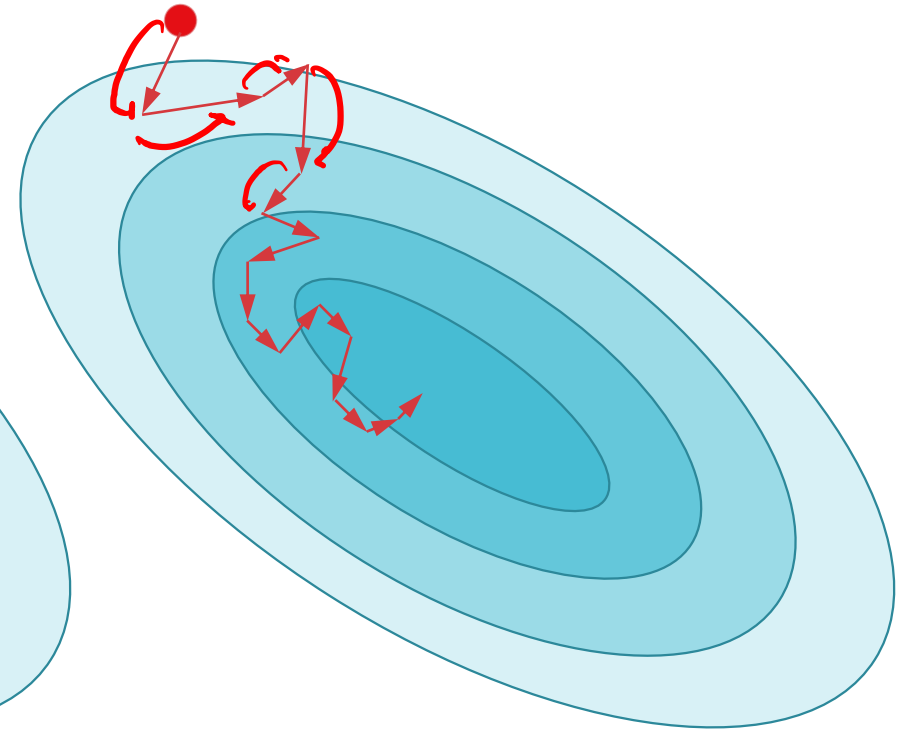
- In practice, the data set is randomly shuffled then looped through so that each data point is used equally often



# Stochastic Gradient Descent vs. Gradient Descent



Gradient Descent



Stochastic Gradient Descent

# Mini-batch Stochastic Gradient Descent

$O(ND)$   $\rightarrow$   $O(BD)$

• Input:  $\mathcal{D} = \{(\mathbf{x}^{(n)}, y^{(n)})\}_{n=1}^N, \eta_{MB}^{(0)}, B$

1. Initialize  $\mathbf{w}^{(0)}$  to all zeros and set  $t = 0$
2. While TERMINATION CRITERION is not satisfied

a. Randomly sample  $B$  data points from  $\mathcal{D}$ :

$$\mathcal{D}_{batch} \{(\mathbf{x}^{(b)}, y^{(b)})\}_{b=1}^B$$

b. Compute the gradient w.r.t. the sampled *batch*:

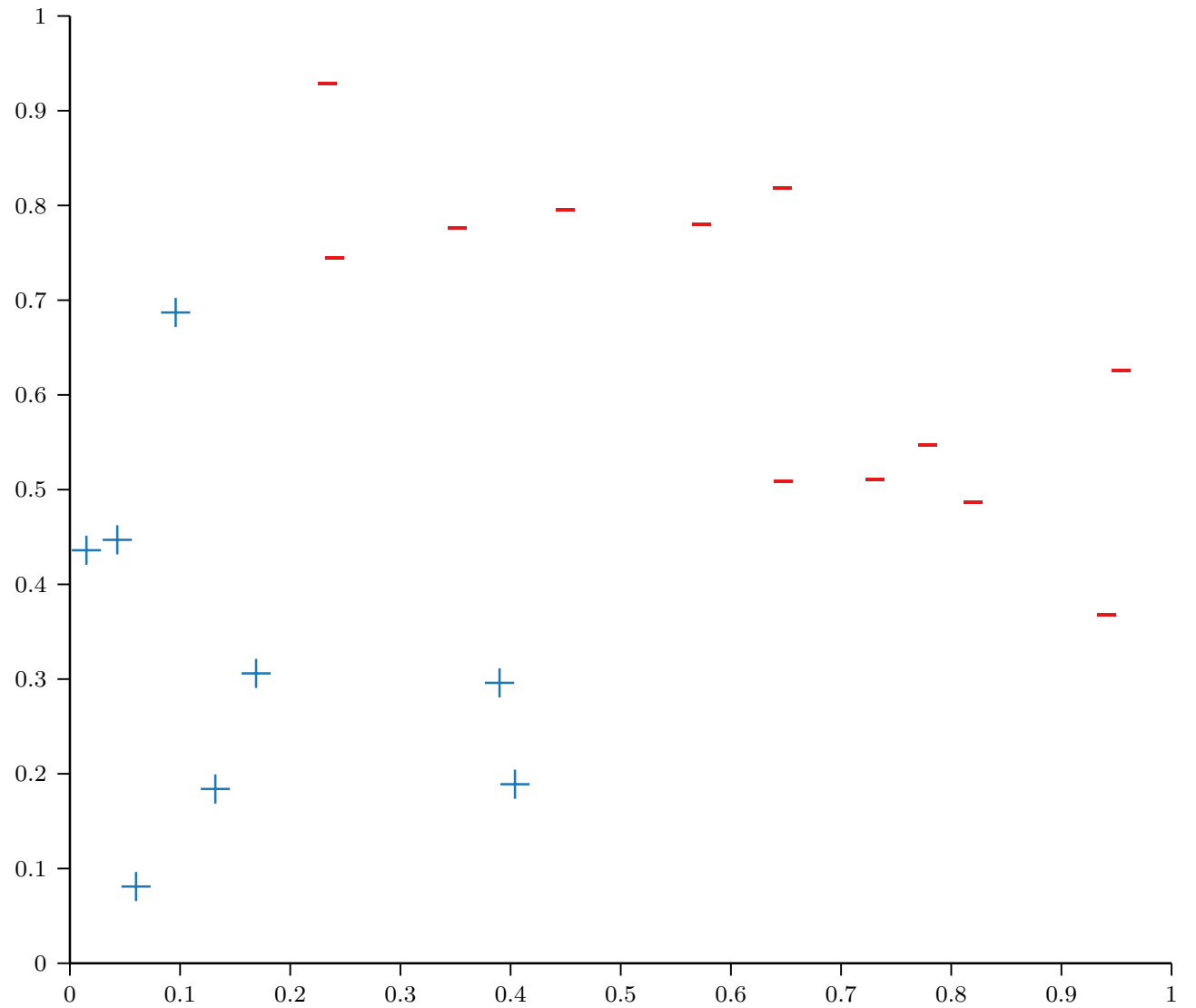
$$\nabla_{\mathbf{w}} \ell_{\mathcal{D}_{batch}}(\mathbf{w}^{(t)}) = \sum_{b=1}^B \mathbf{x}^{(b)} (P(Y = 1 | \mathbf{x}^{(b)}, \mathbf{w}) - y^{(b)})$$

c. Update  $\mathbf{w}$ :  $\mathbf{w}^{(t+1)} \leftarrow \mathbf{w}^{(t)} - \eta_{MB}^{(0)} \nabla_{\mathbf{w}} \ell_{\mathcal{D}_{batch}}(\mathbf{w}^{(t)})$

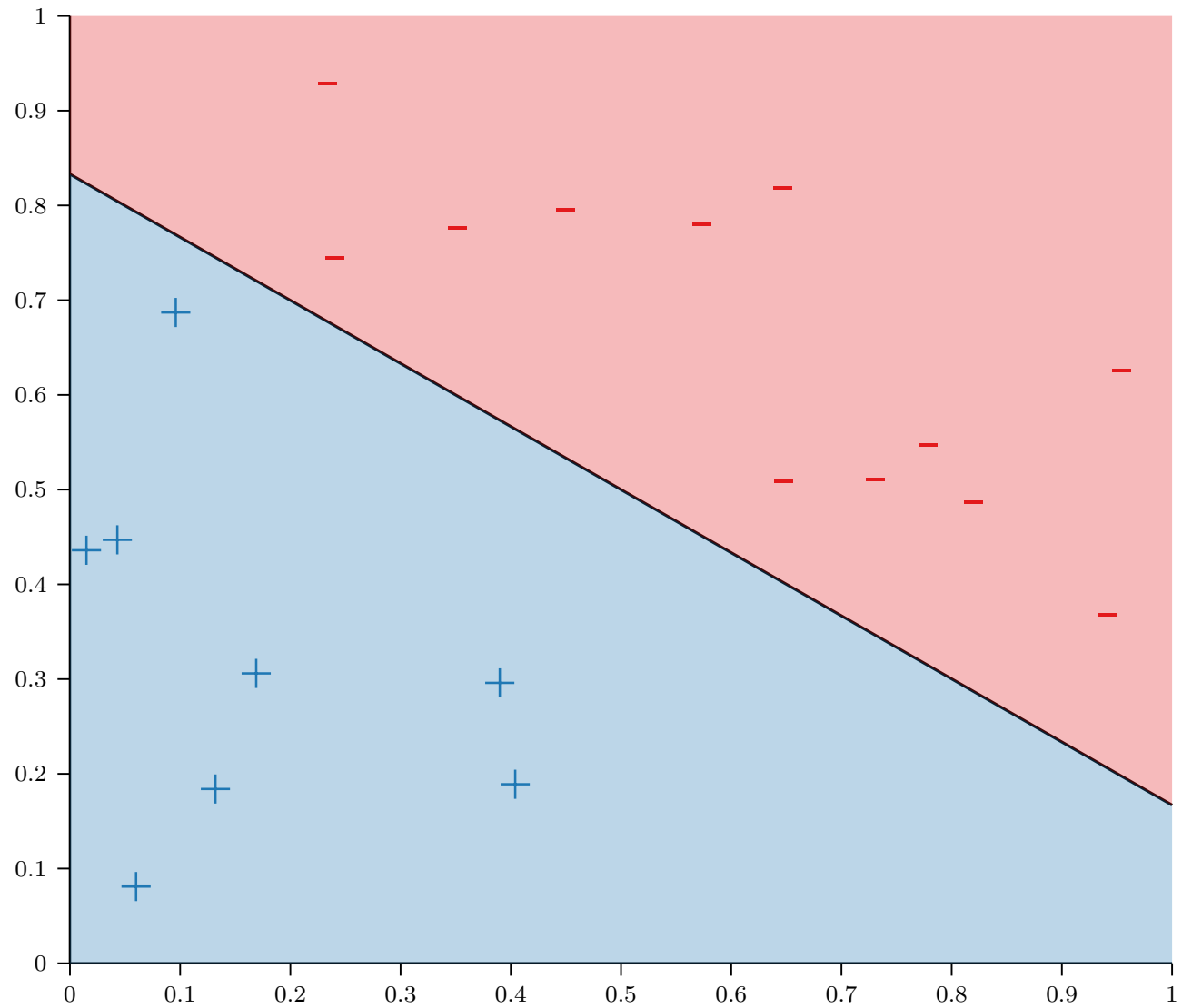
d. Increment  $t$ :  $t \leftarrow t + 1$

• Output:  $\mathbf{w}^{(t)}$

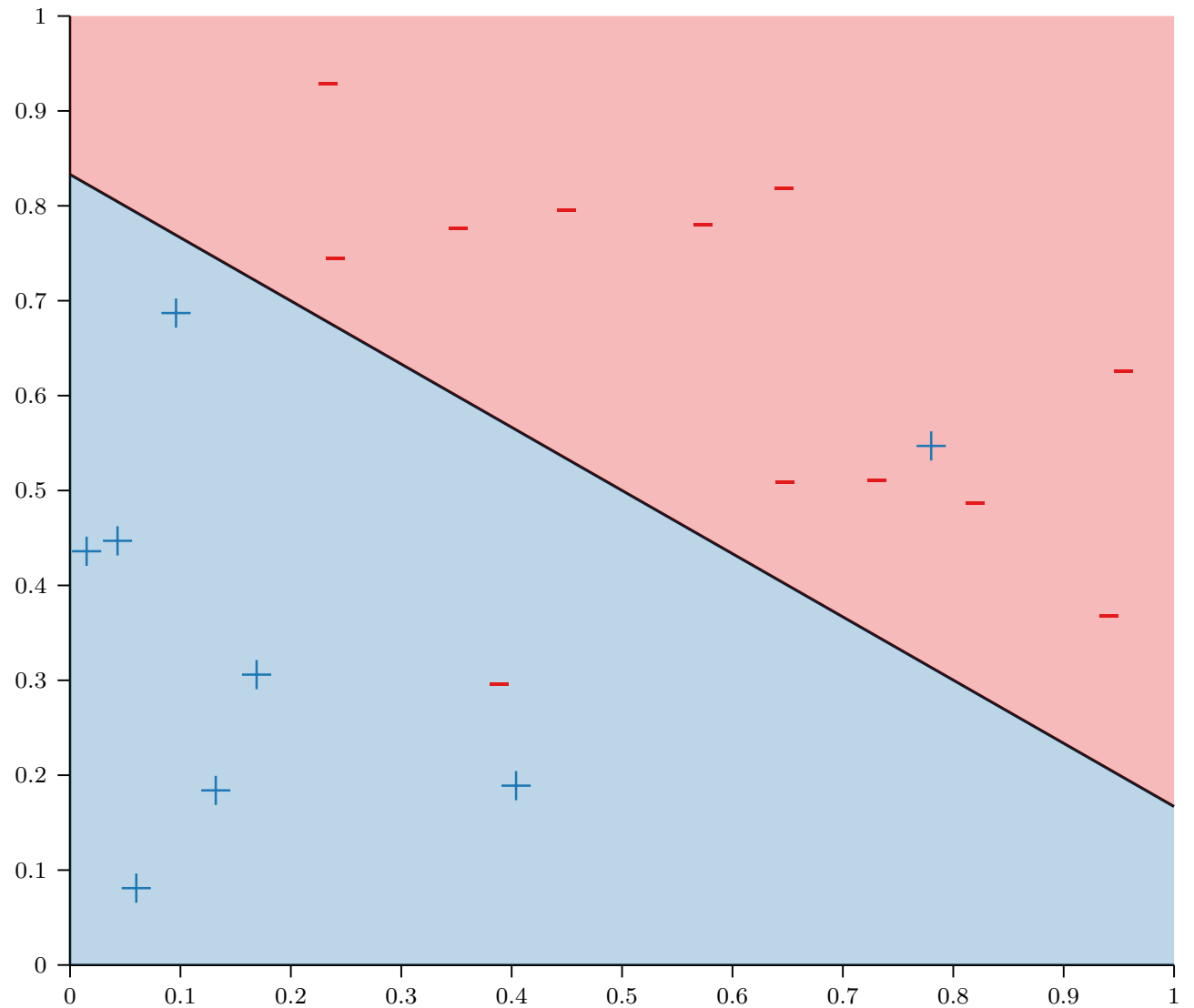
# Linear Models



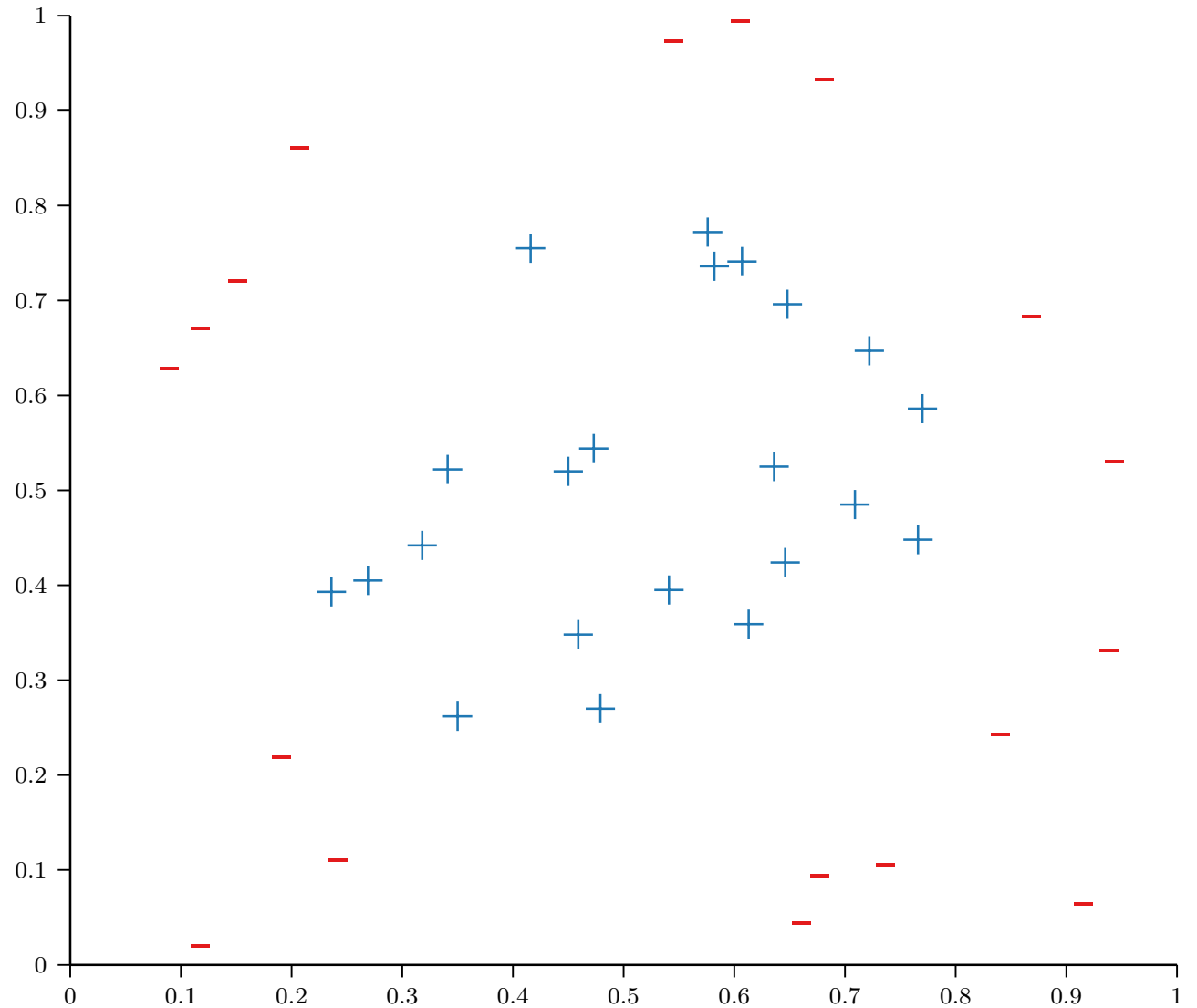
# Linear Models



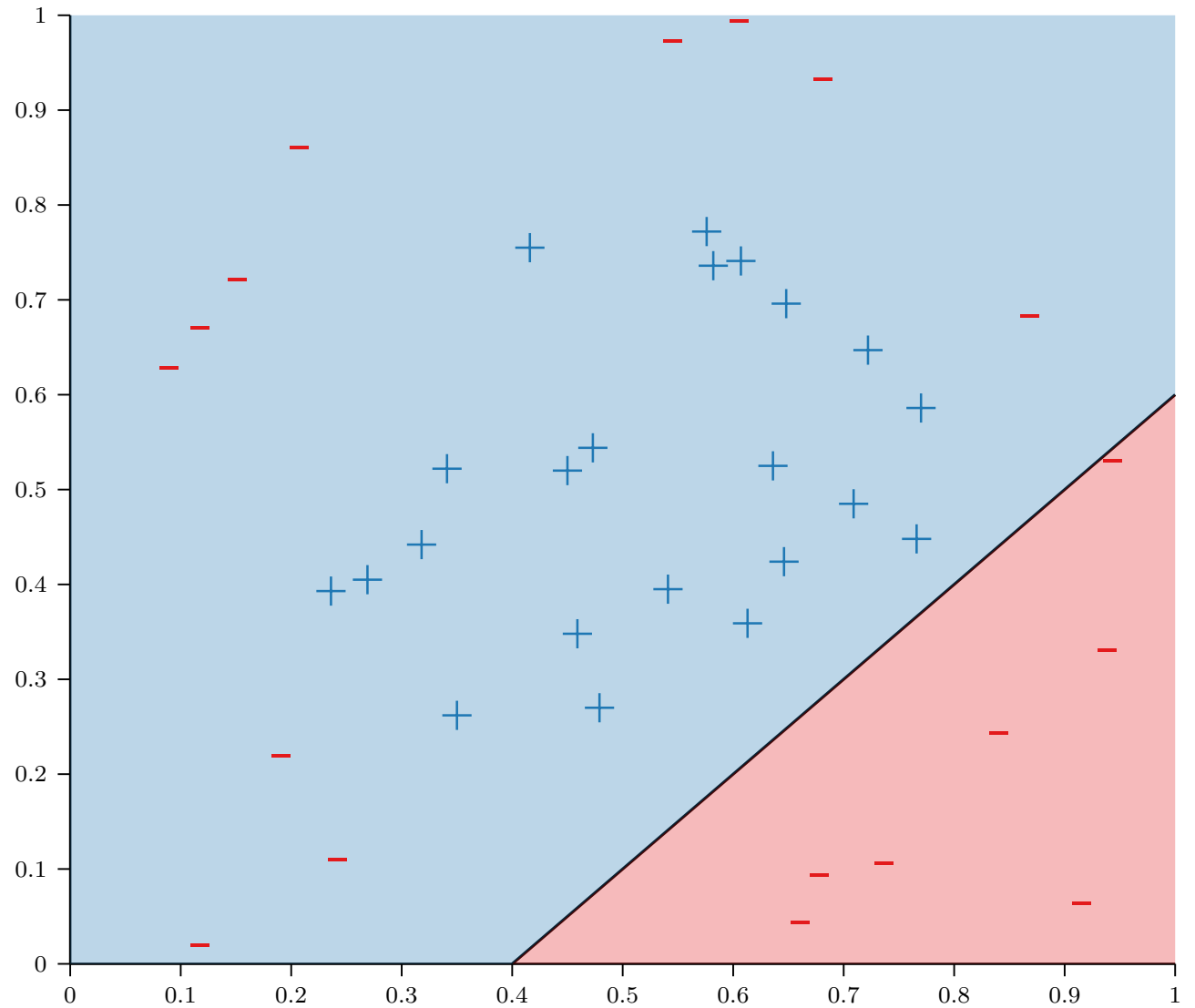
# Linear Models



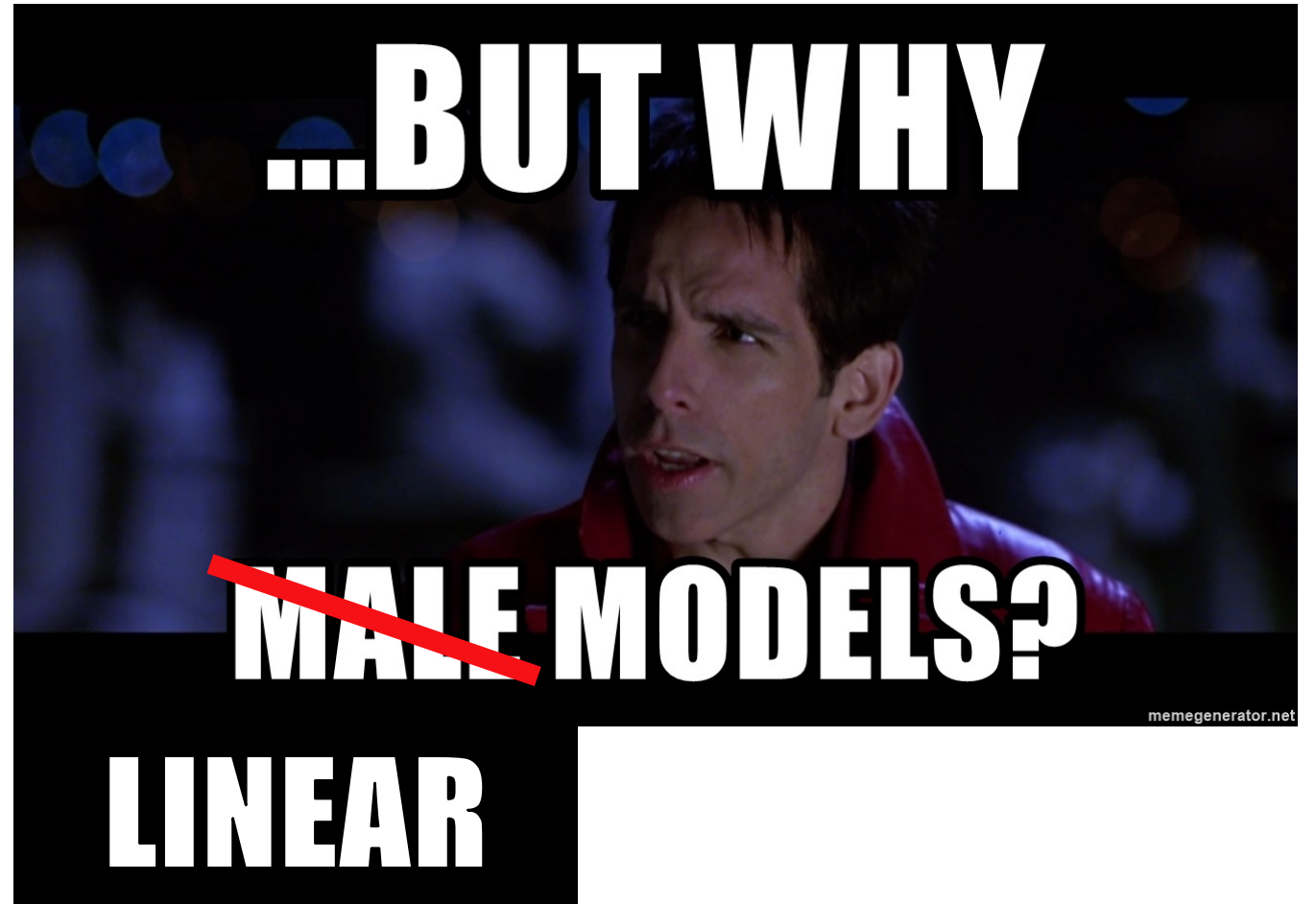
# Linear Models?



# Linear Models?



# Linear Models?





# Nonlinear Models

