

# CS 360: Machine Learning

Prof. Sara Mathieson

Fall 2020



**HVERFORD**  
COLLEGE

# Admin

- Office hours next week:
  - Monday 9:45-11am
  - **Monday** 4:30-6pm
  
- Lab 5 due **October 20**

# Outline for October 16

- Maximum Likelihood Estimation (MLE) in other fields
- Recap handouts and logistic regression so far
- Regularization
- Multi-class logistic regression

# Outline for October 16

- Maximum Likelihood Estimation (MLE) in other fields
- Recap handouts and logistic regression so far
- Regularization
- Multi-class logistic regression

# MLE in other fields

## 1. Chemistry

### **Estimating the number of pure chemical components in a mixture by maximum likelihood**

**E. Levina<sup>1\*</sup>, A.S. Wagaman<sup>1</sup>, A.F. Callender<sup>2</sup>, G.S. Mandair<sup>2</sup> and M.D. Morris<sup>2</sup>**

<sup>1</sup>Department of Statistics, The University of Michigan, Ann Arbor, MI 48109, USA

<sup>2</sup>Department of Chemistry, The University of Michigan, Ann Arbor, MI 48109, USA

Received 31 August 2006; Revised 4 December 2006; Accepted 6 December 2006

## 2. Physics

### **Maximum Likelihood Blood Velocity Estimator Incorporating Properties of Flow Physics**

Malene Schlaikjer and Jørgen Arendt Jensen, *Senior Member, IEEE*

## 3. Biology

### **Maximum-Likelihood Estimation of Admixture Proportions From Genetic Data**

Jinliang Wang

GENETICS *June 1, 2003 vol. 164 no. 2 747-765*

## 4. Economics

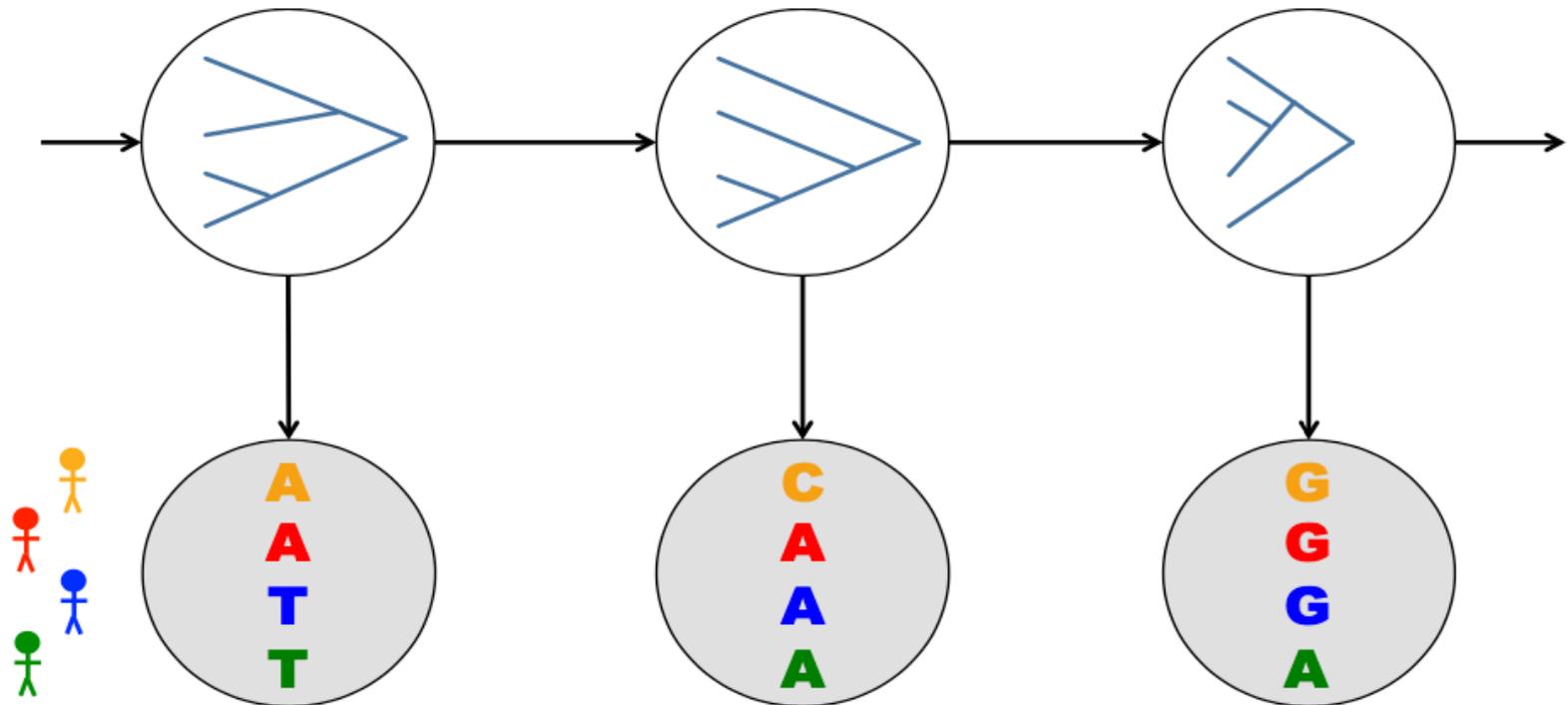
### **Maximum likelihood estimation of mixed C-vines with application to exchange rates**

**Claudia Czado<sup>1</sup>, Ulf Schepsmeier<sup>1</sup> and Aleksey Min<sup>1</sup>**

<sup>1</sup>Zentrum Mathematik, Technische Universität München, Boltzmannstr. 3, 85747 Garching bei München, Germany

# MLE in my own research

- Known: DNA data from multiple individuals
- Unknown: tree of relationships between the individuals



# Outline for October 16

- Maximum Likelihood Estimation (MLE) in other fields
- Recap handouts and logistic regression so far
- Regularization
- Multi-class logistic regression

# Logistic Regression Warmup

1. The output of logistic regression is a model that creates:

- (a) a linear decision boundary
- (b) a logistic decision boundary
- (c) no decision boundary

2. We use logistic regression for:

- (a) classification
- (b) regression
- (c) both

3. Our hypothesis in logistic regression is:

$$h_{\mathbf{w}}(\mathbf{x}) = p(y = 1|\mathbf{x}) = \frac{1}{1 + e^{-\mathbf{w} \cdot \mathbf{x}}}$$

If  $\mathbf{w}$  is the zero vector (as it would be when starting SGD), what is the probability  $y = 1$ ?

4. How did we define the cost function for logistic regression? (Bonus: write down the cost function)

- (a) likelihood
- (b) log likelihood
- (c) negative log likelihood

# Logistic Regression Warmup

1. The output of logistic regression is a model that creates:
  - (a) a linear decision boundary
  - (b) a logistic decision boundary
  - (c) no decision boundary

# Logistic Regression Warmup

1. The output of logistic regression is a model that creates:

- (a) a linear decision boundary
- (b) a logistic decision boundary
- (c) no decision boundary

2. We use logistic regression for:

- (a) classification
- (b) regression
- (c) both

# Logistic Regression Warmup

1. The output of logistic regression is a model that creates:

- (a) a linear decision boundary
- (b) a logistic decision boundary
- (c) no decision boundary

2. We use logistic regression for:

- (a) classification
- (b) regression
- (c) both

3. Our hypothesis in logistic regression is:

$$h_{\mathbf{w}}(\mathbf{x}) = p(y = 1|\mathbf{x}) = \frac{1}{1 + e^{-\mathbf{w} \cdot \mathbf{x}}}$$

If  $\mathbf{w}$  is the zero vector (as it would be when starting SGD), what is the probability  $y = 1$ ?

# Logistic Regression Warmup

1. The output of logistic regression is a model that creates:

- (a) a linear decision boundary
- (b) a logistic decision boundary
- (c) no decision boundary

2. We use logistic regression for:

- (a) classification
- (b) regression
- (c) both

3. Our hypothesis in logistic regression is:

$$h_{\mathbf{w}}(\mathbf{x}) = p(y = 1|\mathbf{x}) = \frac{1}{1 + e^{-\mathbf{w} \cdot \mathbf{x}}}$$

If  $\mathbf{w}$  is the zero vector (as it would be when starting SGD), what is the probability  $y = 1$ ?

$\frac{1}{2}$

4. How did we define the cost function for logistic regression? (Bonus: write down the cost function)

- (a) likelihood
- (b) log likelihood
- (c) negative log likelihood

# Logistic Regression Warmup

1. The output of logistic regression is a model that creates:

- (a) a linear decision boundary
- (b) a logistic decision boundary
- (c) no decision boundary

2. We use logistic regression for:

- (a) classification
- (b) regression
- (c) both

3. Our hypothesis in logistic regression is:

$$h_{\mathbf{w}}(\mathbf{x}) = p(y = 1|\mathbf{x}) = \frac{1}{1 + e^{-\mathbf{w} \cdot \mathbf{x}}}$$

If  $\mathbf{w}$  is the zero vector (as it would be when starting SGD), what is the probability  $y = 1$ ?

$\frac{1}{2}$

4. How did we define the cost function for logistic regression? (Bonus: write down the cost function)

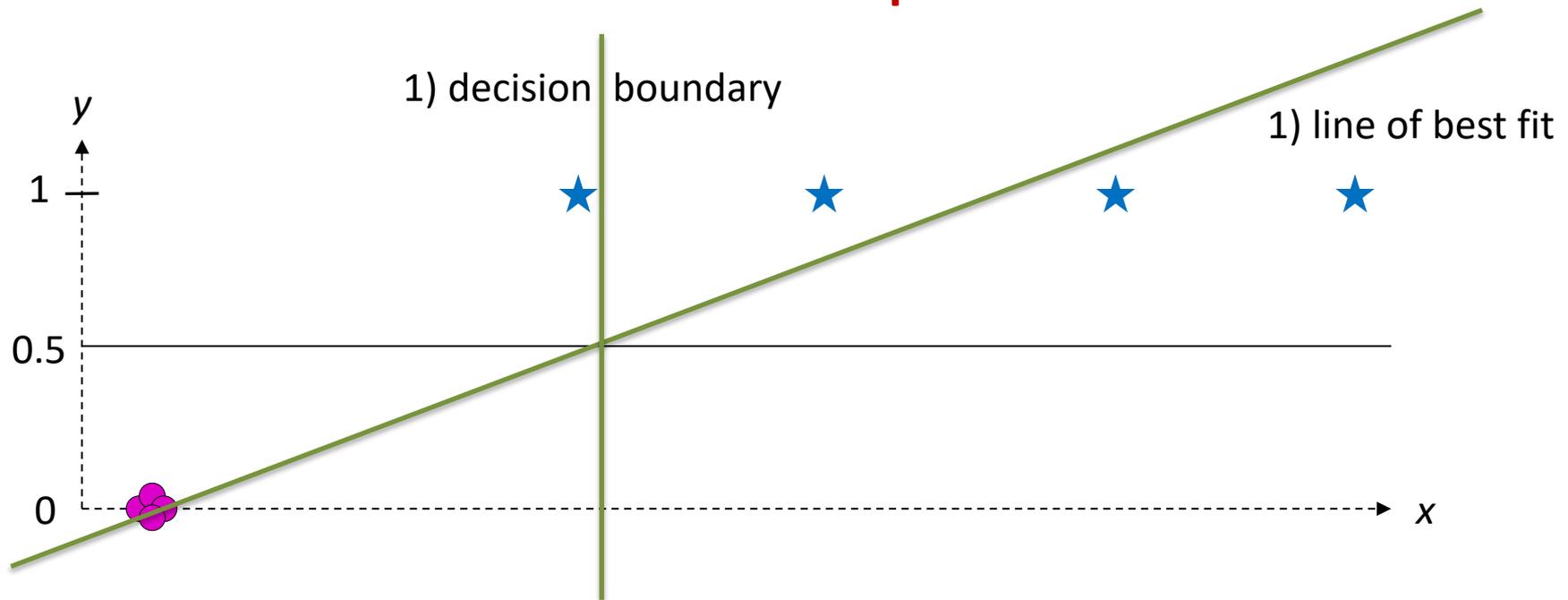
- (a) likelihood
- (b) log likelihood
- (c) negative log likelihood

# Extra Example



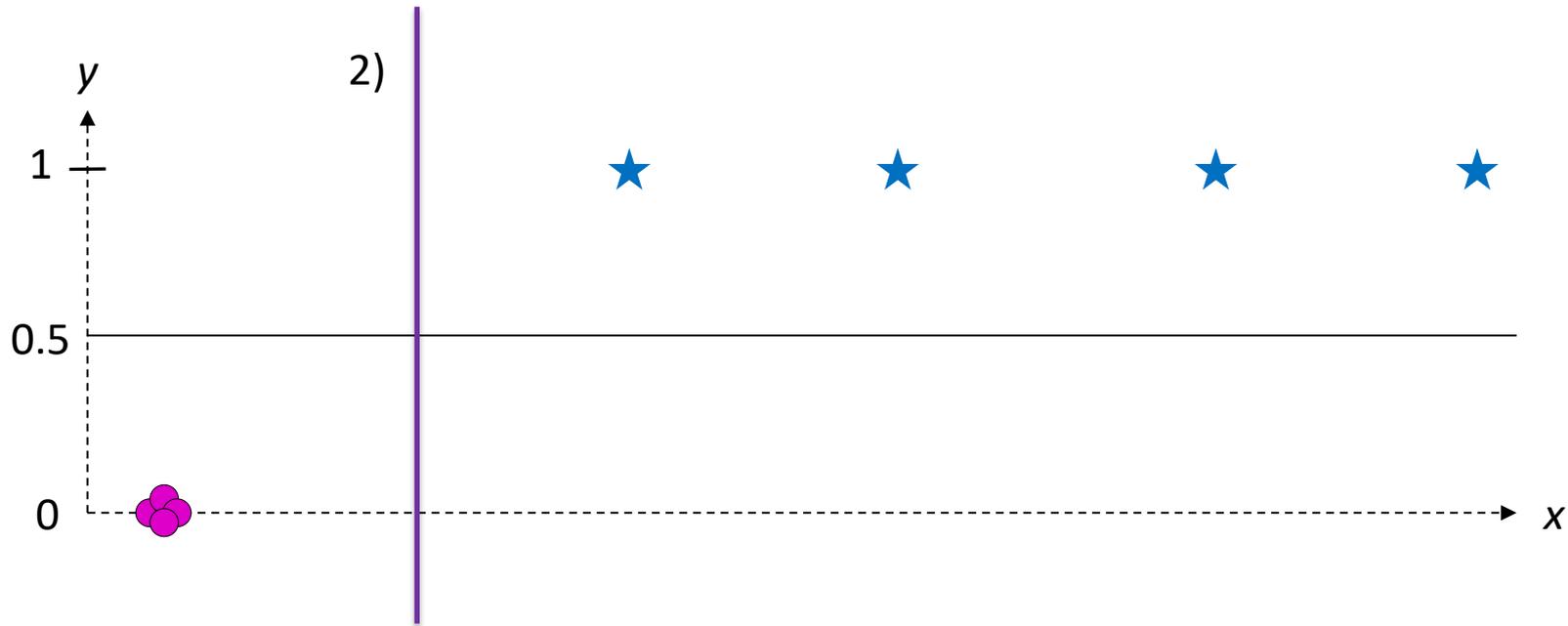
- 1) What line of best fit would be produced by **linear regression**? (roughly)
- 2) What linear decision boundary would be produced by **logistic regression**? (roughly)

# Extra Example



- 1) What line of best fit would be produced by **linear regression**? (roughly)
- 2) What linear decision boundary would be produced by **logistic regression**? (roughly)

# Extra Example



- 1) What line of best fit would be produced by **linear regression**? (roughly)
- 2) What linear decision boundary would be produced by **logistic regression**? (roughly)

# Stochastic Gradient Descent for Logistic Regression (binary classification)

```
set  $w = 0$  vector
```

```
while cost  $J(w)$  still changing:
```

```
    shuffle data points
```

```
    for  $i = 1 \dots n$ :
```

```
         $w \leftarrow w - \alpha(\text{derivative of } J(w) \text{ wrt } x_i)$ 
```

```
    store  $J(w)$ 
```

# 3 important pieces to SGD

- Hypothesis function (prediction)

$$h_w(\mathbf{x}) = p(y = 1|\mathbf{x}) = \frac{1}{1 + e^{-w \cdot \mathbf{x}}}$$

# 3 important pieces to SGD

- Hypothesis function (prediction)

$$h_{\mathbf{w}}(\mathbf{x}) = p(y = 1|\mathbf{x}) = \frac{1}{1 + e^{-\mathbf{w} \cdot \mathbf{x}}}$$

- Cost function (want to minimize)

$$J(\mathbf{w}) = - \sum_{i=1}^n y_i \log h_{\mathbf{w}}(\mathbf{x}_i) + (1 - y_i) \log(1 - h_{\mathbf{w}}(\mathbf{x}_i))$$

# 3 important pieces to SGD

- Hypothesis function (prediction)

$$h_{\mathbf{w}}(\mathbf{x}) = p(y = 1|\mathbf{x}) = \frac{1}{1 + e^{-\mathbf{w} \cdot \mathbf{x}}}$$

- Cost function (want to minimize)

$$J(\mathbf{w}) = - \sum_{i=1}^n y_i \log h_{\mathbf{w}}(\mathbf{x}_i) + (1 - y_i) \log(1 - h_{\mathbf{w}}(\mathbf{x}_i))$$

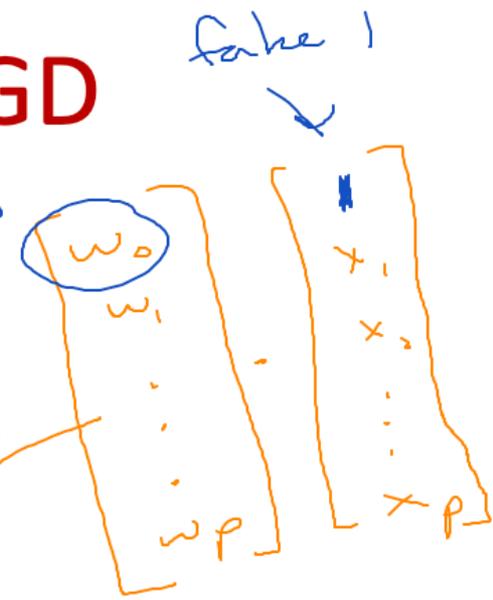
- Gradient of cost wrt single data point  $\mathbf{x}_i$

$$\nabla J_{\mathbf{x}_i}(\mathbf{w}) = (h_{\mathbf{w}}(\mathbf{x}_i) - y_i)\mathbf{x}_i$$

# 3 important pieces to SGD

- Hypothesis function (prediction)

$$h_{\mathbf{w}}(\mathbf{x}) = p(y = 1|\mathbf{x}) = \frac{1}{1 + e^{-\mathbf{w} \cdot \mathbf{x}}}$$

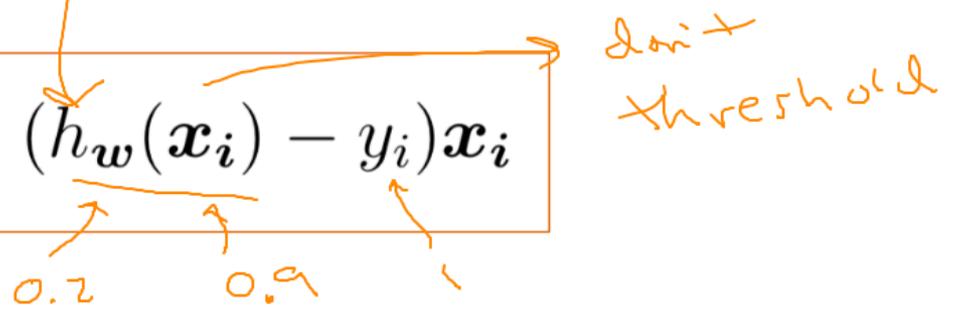


- Cost function (want to minimize)

$$J(\mathbf{w}) = - \sum_{i=1}^n y_i \log h_{\mathbf{w}}(\mathbf{x}_i) + (1 - y_i) \log(1 - h_{\mathbf{w}}(\mathbf{x}_i))$$

- Gradient of cost wrt single data point  $\mathbf{x}_i$

$$\nabla J_{\mathbf{x}_i}(\mathbf{w}) = (h_{\mathbf{w}}(\mathbf{x}_i) - y_i) \mathbf{x}_i$$



$p=1$

# Handout 9

1. Say I train a binary logistic regression model (i.e. outcomes  $\in \{0,1\}$ ) and end up with  $\hat{w} = [\hat{w}_0, \hat{w}_1]^T = [-4, -5]^T$ . What is the decision boundary? Sketch a graph of this logistic model and label the decision boundary. How would you classify a new point  $x_{test} = -2$ ?

$$\vec{w} \cdot \vec{x} > 0 \Rightarrow \text{predict } 1$$

$$-4 - 5x > 0$$

$$-5x > 4$$

$$x < -\frac{4}{5}$$

$\Rightarrow$  predict 1

$$\frac{1}{1 + e^{-\vec{w} \cdot \vec{x}}} > \frac{1}{2}$$

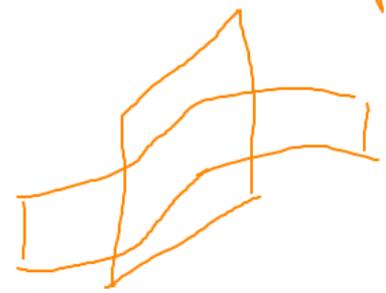
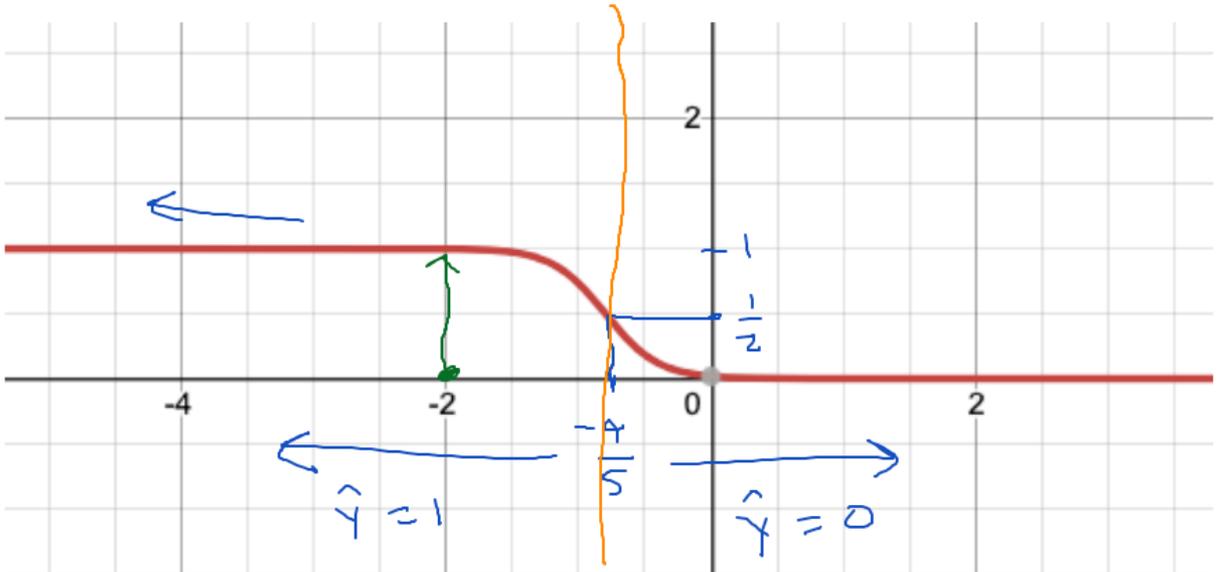
$$\hat{y}_{test} = 1$$

$$\vec{w} \cdot \vec{x} > 0$$

$p=2$

$$w_0 + w_1 x_1 + w_2 x_2 > 0$$

plane

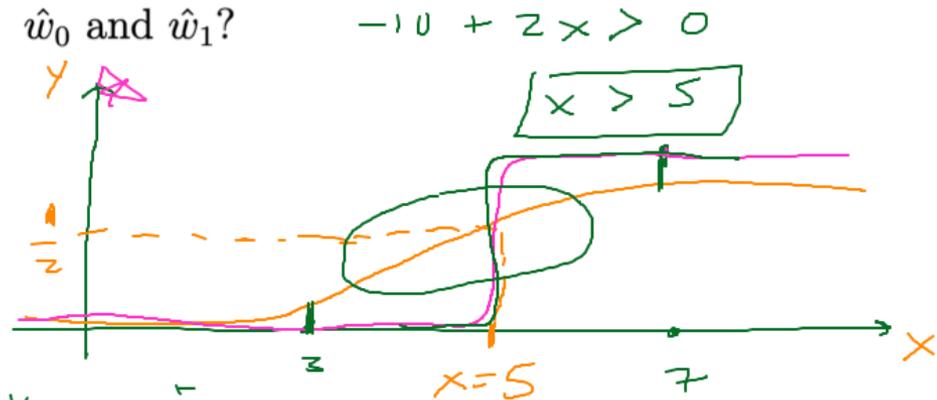


$$\frac{1}{1 + e^{-(-4 - 5x)}}$$

# Lab 5 intro

Say we have  $p = 1$  and two training examples:  $(x_1, y_1) = (3, 0)$  and  $(x_2, y_2) = (7, 1)$ , and we would like to fit a logistic model to this dataset.

1. Draw these two examples on a coordinate system and sketch a logistic function that would fit them (roughly). What is the optimal decision boundary? Does this help us uniquely determine  $\hat{w}_0$  and  $\hat{w}_1$ ?



**NO**

$x > 5$

$Z(-5 + x > 0)$

$w_0 = -5$   
 $w_1 = 1$

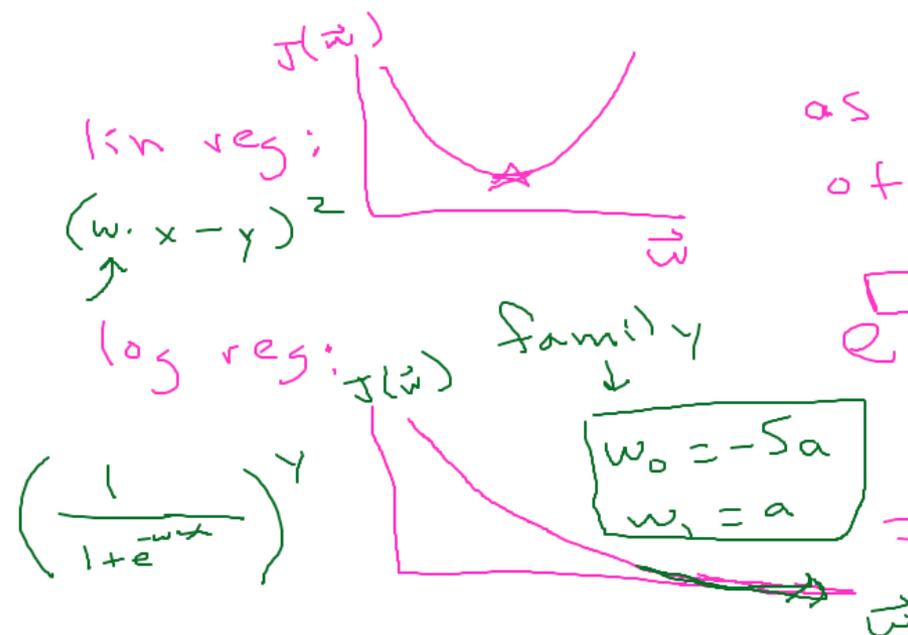
$\Rightarrow h_{\hat{w}}(\vec{x}) = \frac{1}{1 + e^{-(-5+x)}}$

"better"

$w_0 = -10$   
 $w_1 = 2$

↓

$w_0 = -100$   
 $w_1 = 20$



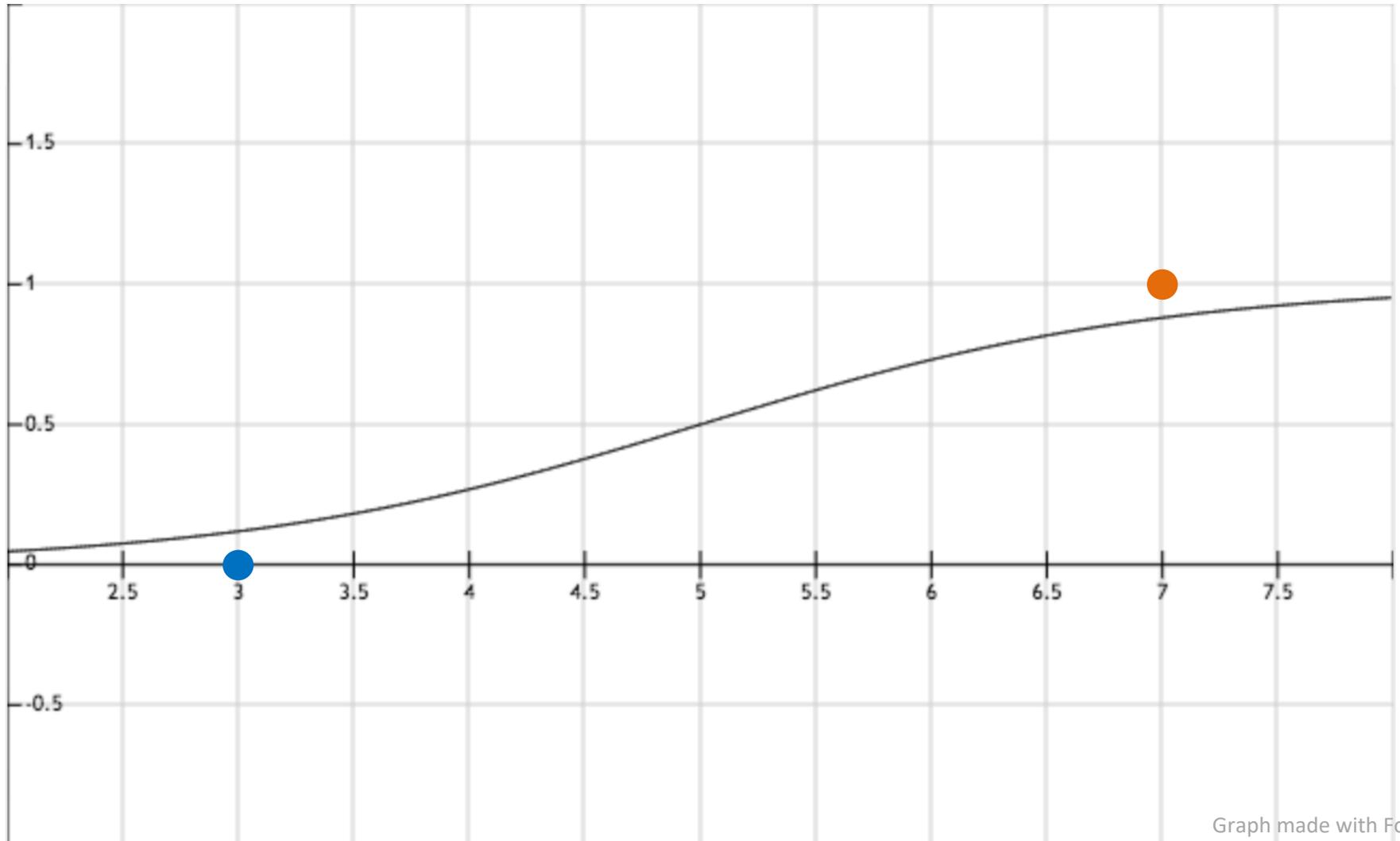
as magnitude of weight ↑

$e \rightarrow 0$   
 $\rightarrow \infty$

$\Rightarrow$  overly confident

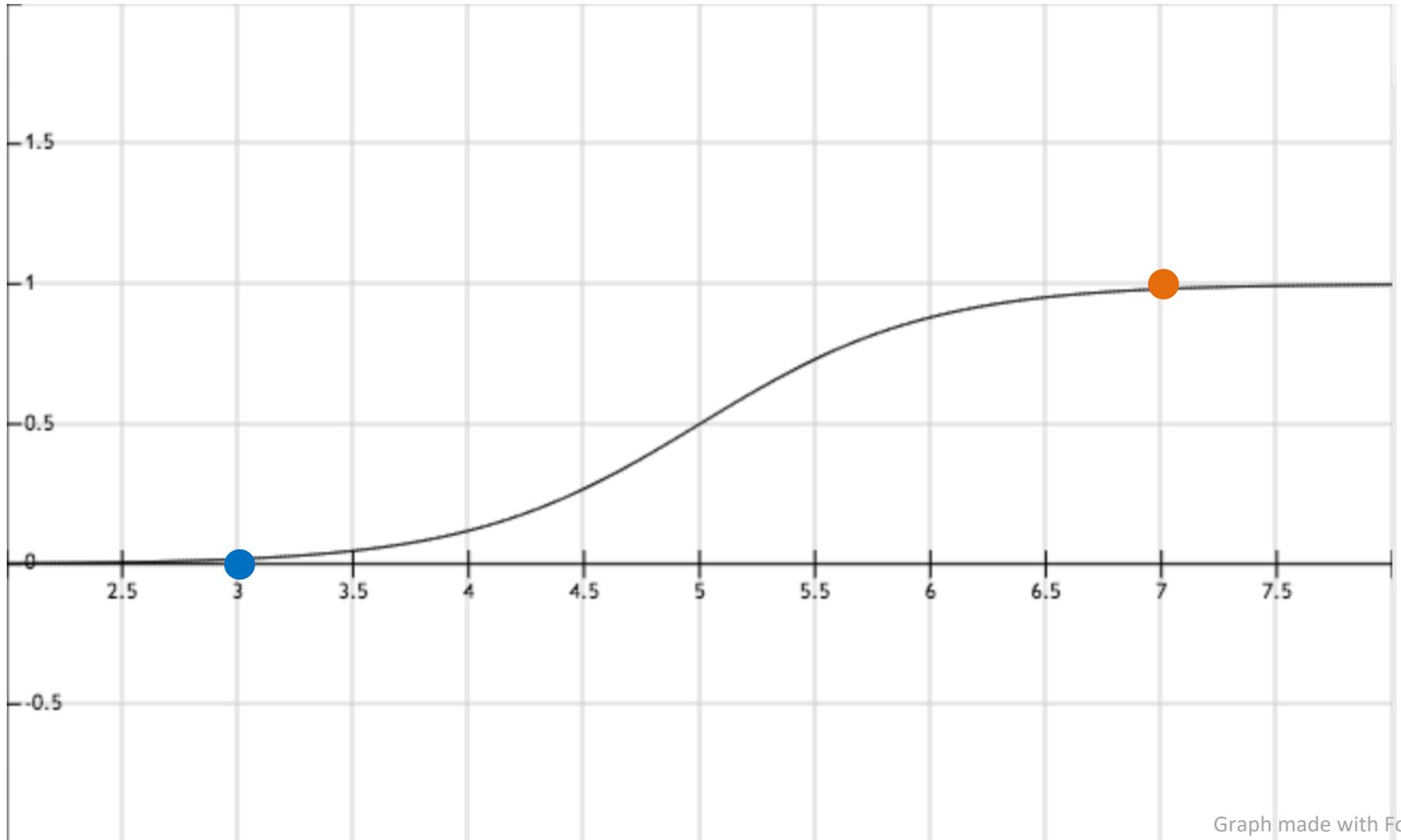
$$w_0 = -5, w_1 = 1$$

$$h_w(x) = 1 / (1 + e^{(5-x)})$$



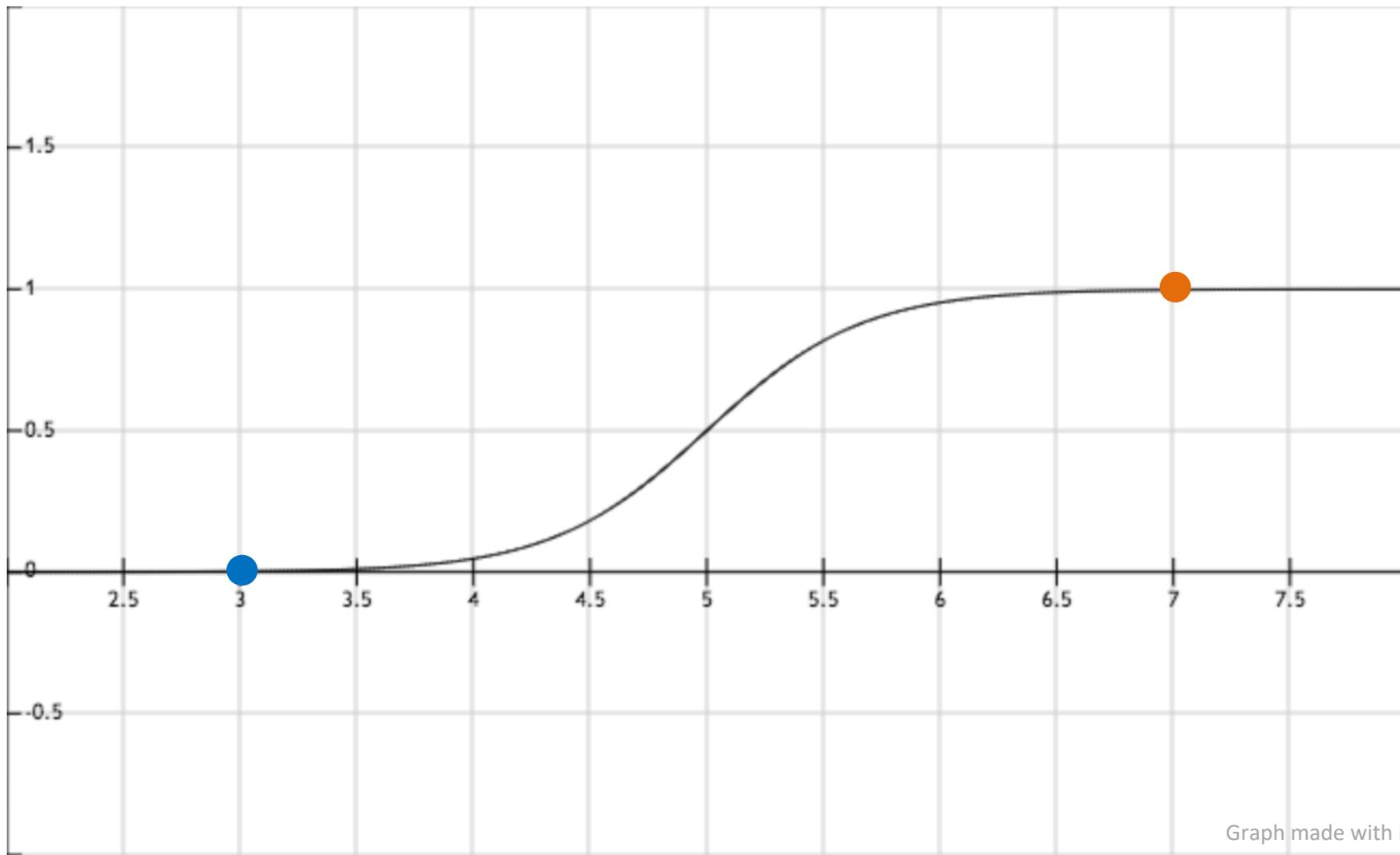
$$w_0 = -10, w_1 = 2$$

$$h_w(x) = 1 / (1 + e^{(10 - 2x)})$$

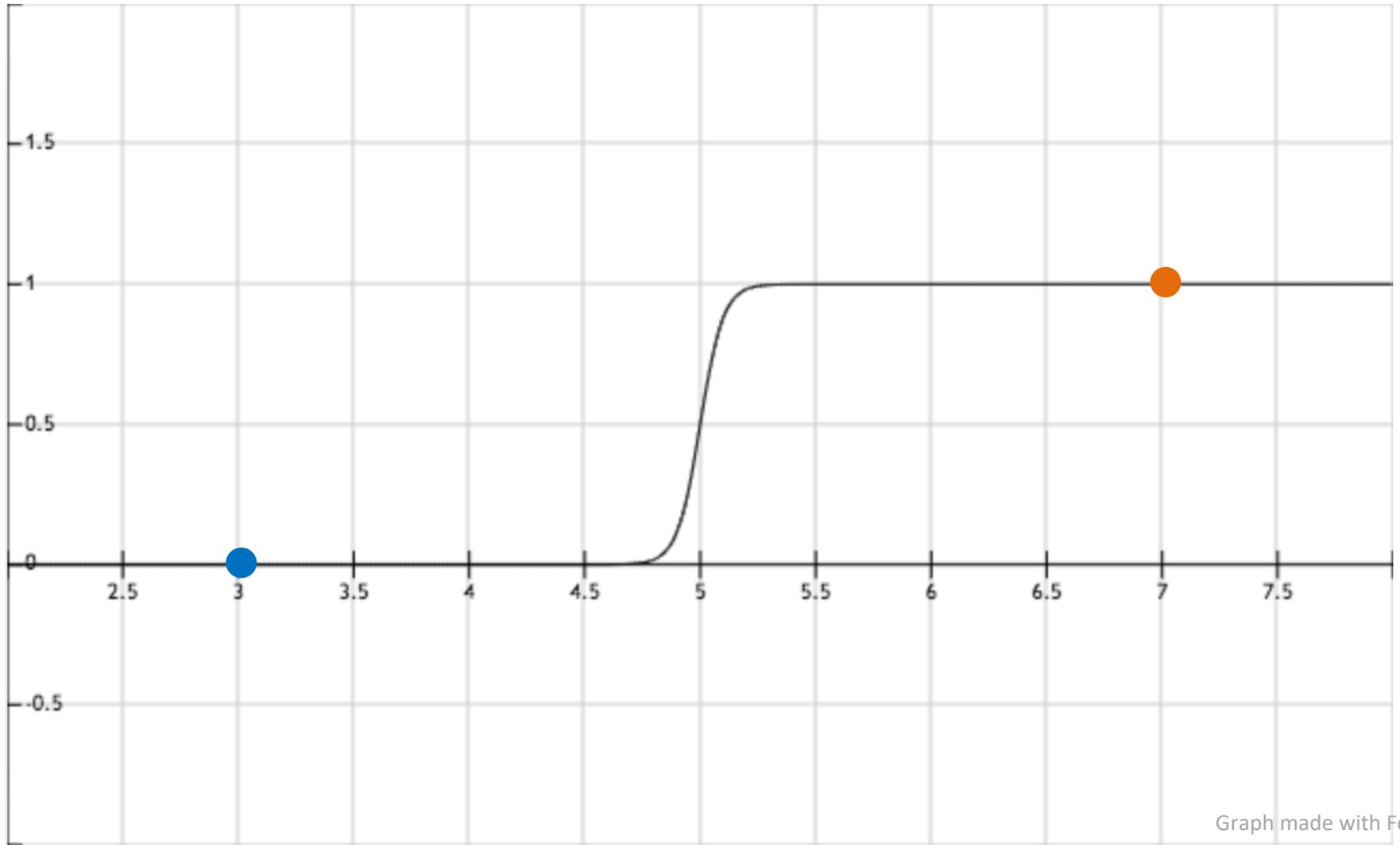


$$w_0 = -15, w_1 = 3$$

$$h_w(x) = 1 / (1 + e^{(15 - 3x)})$$



$$w_0 = -100, w_1 = 20 \quad h_w(x) = 1 / (1 + e^{(100 - 20x)})$$



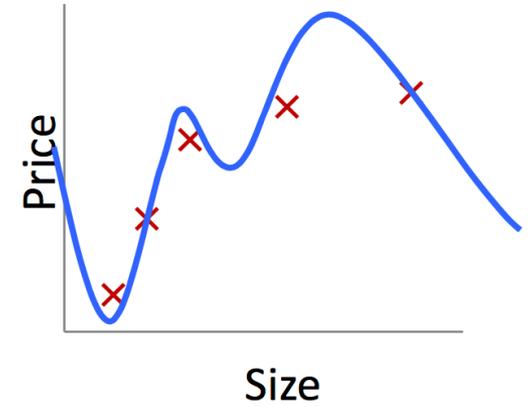
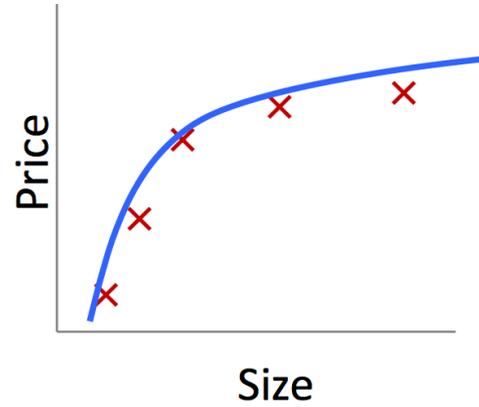
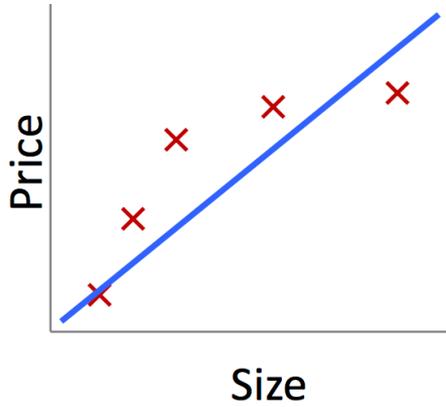


# Outline for October 16

- Maximum Likelihood Estimation (MLE) in other fields
- Recap handouts and logistic regression so far
- **Regularization**
- Multi-class logistic regression

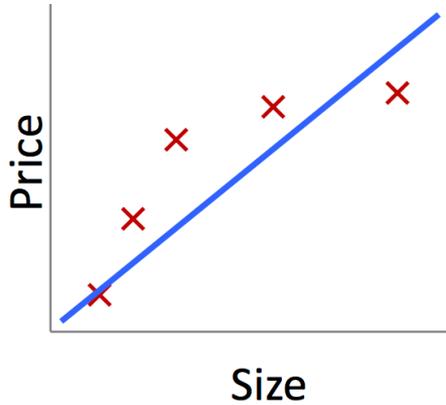
# Generalization error

- Example: price vs. size (i.e. of a house or car)

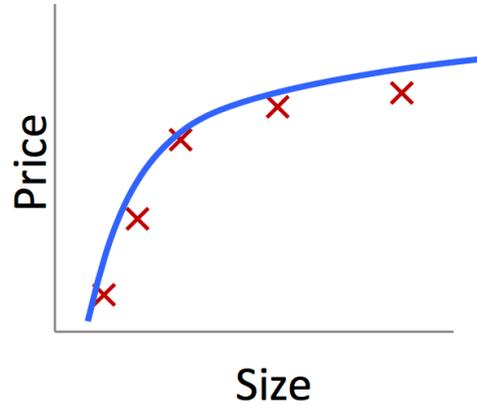


# Generalization error

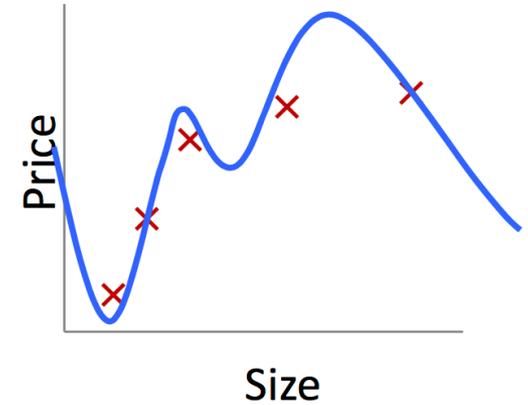
- Example: price vs. size (i.e. of a house or car)



underfitting  
(high bias)



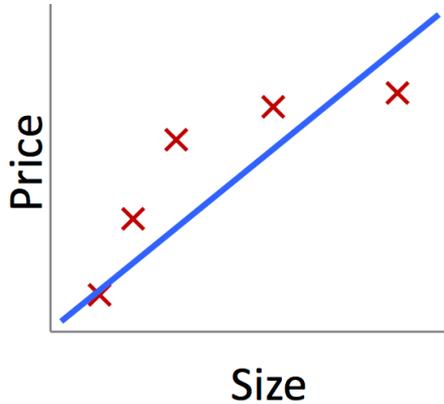
correct fit



overfitting  
(high variance)

# Generalization error

- Example: price vs. size (i.e. of a house or car)

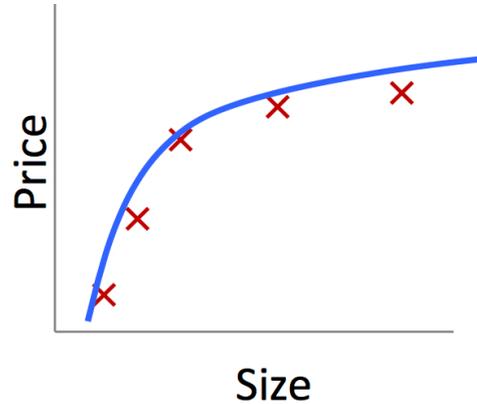


underfitting  
(high bias)

Structural error:

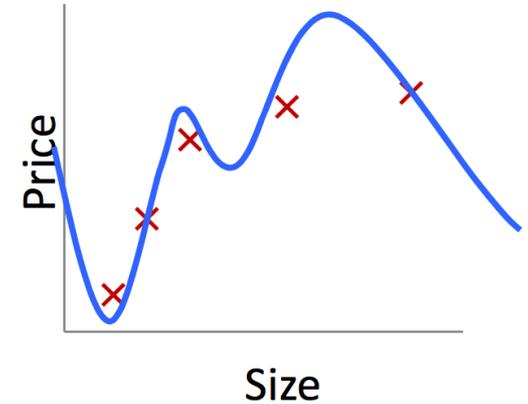
Hypothesis space cannot model true relationship

- More data doesn't help
- Need a more flexible model



correct fit

balance  
↔



overfitting  
(high variance)

Estimation (approximation) error:

Hypothesis space *can* model true relationship, BUT hard to identify correct model due to large hypothesis space, small  $n$ , or noise

- Reduce hypothesis space
- Add more data

# Regularization

What if ...

- we have a limited # of training examples ( $n < p$ ), or
- we want to automatically control the complexity of the learned hypothesis?

# Regularization

What if ...

- we have a limited # of training examples ( $n < p$ ), or
- we want to automatically control the complexity of the learned hypothesis?

Idea: penalize large values of  $w_j$

Why prefer small weights?

- if large weights, small change in feature can result in large change in prediction
- prevent giving too much weight to any one feature
- might prefer zero weight for useless features

# Common Regularizers

$$\|\vec{w}\|_0 = \sum_{j:w_j \neq 0} 1$$

## $L_0$ norm

- Number of non-zero entries
- Minimizing  $L_0$  norm is NP hard

# Common Regularizers

$$\|\vec{w}\|_0 = \sum_{j:w_j \neq 0} 1$$

$L_0$  norm

- Number of non-zero entries
- Minimizing  $L_0$  norm is NP hard

$$\|\vec{w}\|_1 = \sum_{j=1}^p |w_j|$$

$L_1$  norm

- Sum of magnitude of weights
- Not differentiable

# Common Regularizers

$$\|\vec{w}\|_0 = \sum_{j:w_j \neq 0} 1$$

$L_0$  norm

- Number of non-zero entries
- Minimizing  $L_0$  norm is NP hard

$$\|\vec{w}\|_1 = \sum_{j=1}^p |w_j|$$

$L_1$  norm

- Sum of magnitude of weights
- Not differentiable

$$\|\vec{w}\|_2 = \sqrt{\sum_{j=1}^p w_j^2}$$

$L_2$  norm

- Sum of squared weights
- Differentiable

# Adding Regularization to Logistic Regression

$$\min_{\vec{w}} J^R(\vec{w}) = - \left[ \sum_{i=1}^n \underbrace{y_i \log h(\vec{x}_i)}_{y_i=1} + \underbrace{(1-y_i) \log(1-h(\vec{x}_i))}_{y_i=0} \right]$$

$\lambda$  = regularization

abs. value

$\Rightarrow$  small

and

param hyper positive

$[0, 1)$

$$+ \lambda \sum_{j=1}^p w_j^2$$

don't regularize  $(w_0)$   
bias

SGD

$$\vec{w}^{p+1} = \vec{w}^p - \eta \left[ \underbrace{(h(\vec{x}_i) - y_i)}_{\text{error}} \underbrace{\vec{x}_i}_{\text{features}} + \lambda \underbrace{\vec{w}^p}_{\text{weights}} \right]$$



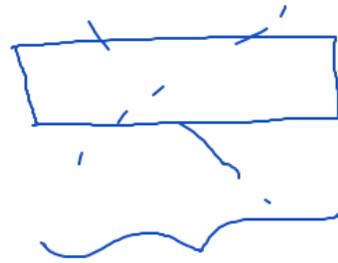
$$\vec{w}^{p+1} = (1 - \eta \lambda) \vec{w}^p - \eta [(h(\vec{x}_i) - y_i) \vec{x}_i]$$

pulls  $\vec{w}$  toward 0

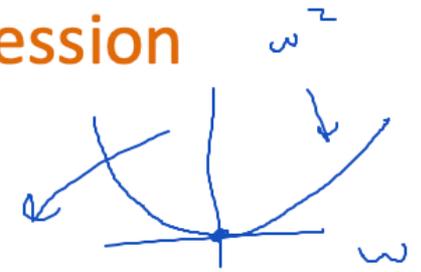
- $\lambda \leftarrow$  regularizer
- $\eta \leftarrow$  learning rate
- $\epsilon \leftarrow$  stopping criteria

# Adding Regularization to Logistic Regression

$$\min J(w)$$



$$+ \frac{\lambda}{2} w^2$$



min with  $w=0$

$$J'(w) = \lambda w = 0$$

$$\Rightarrow \boxed{w=0}$$

# Outline for October 16

- Maximum Likelihood Estimation (MLE) in other fields
- Recap handouts and logistic regression so far
- Regularization
- Multi-class logistic regression

In "Sunday" video!