

CS 360: Machine Learning

Prof. Sara Mathieson

Fall 2019



Admin

- No office hours Friday (but feel free to make an appointment)
- **Lab 5 due October 22** (Tuesday after fall break)
 - Code reviews today

Reading Quiz 5

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

Reading Quiz 5

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

Reading Quiz 5

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 $y = 1$?

Reading Quiz 5

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

Reading Quiz 5

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

Outline for October 10

- Recap SGD for logistic regression
- Regularization
- Multi-class logistic regression
- Begin: evaluation metrics

Outline for October 10

- Recap SGD for logistic regression
- Multi-class logistic regression
- Regularization
- Begin: evaluation metrics

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)

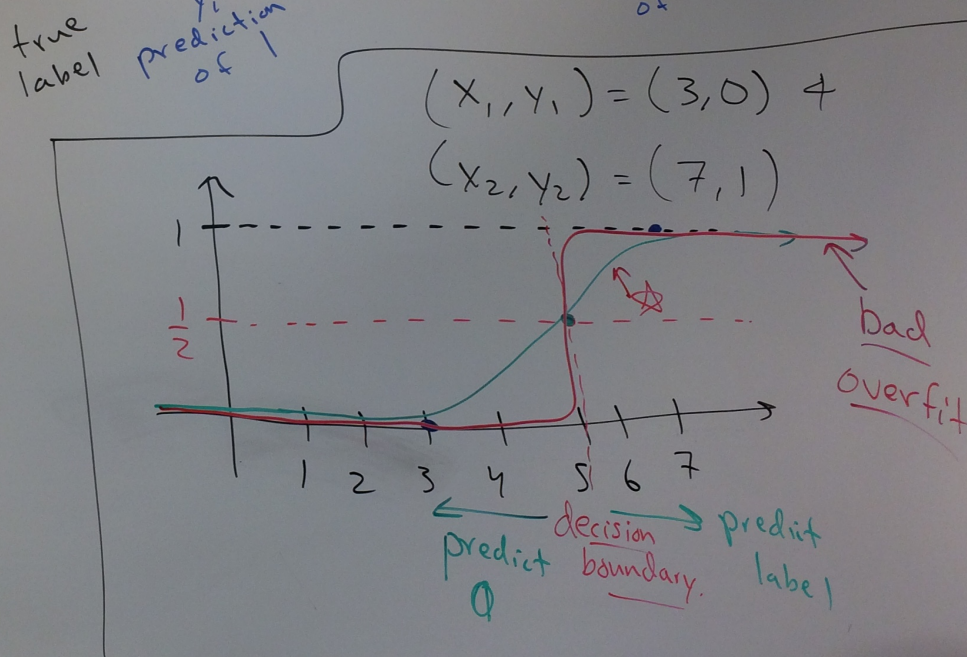
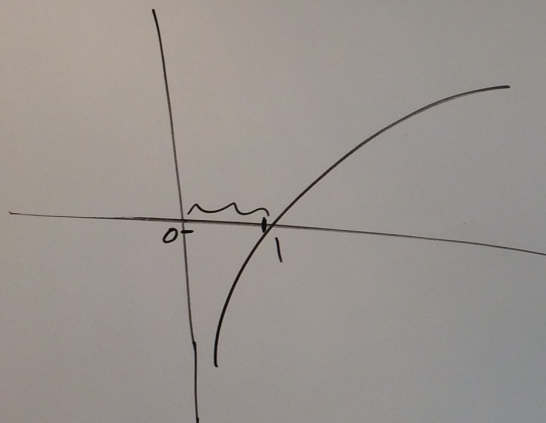
$$J(\mathbf{w}) = - \sum_{i=1}^n y_i \log h_{\mathbf{w}}(\mathbf{x}_i) + (1 - y_i) \log(1 - h_{\mathbf{w}}(\mathbf{x}_i))$$

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

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

cost
= negative
log likelihood

$$J(\vec{w}) = - \sum_{i=1}^n y_i \log \underbrace{h_{\vec{w}}(\vec{x}_i)}_{\substack{\hat{y}_i \\ \text{prediction of 1}}} + (1 - y_i) \log \underbrace{(1 - h_{\vec{w}}(\vec{x}_i))}_{\substack{\text{prediction of 0}}}$$



if $x > 5 \Rightarrow \text{label } \hat{y} = 1$

$x \leq 5 \Rightarrow \text{label } \hat{y} = 0$

positive constant $\rightarrow a[-5 + x > 0]$

$$-5a + ax > 0$$

family
of
solutions!

$$\begin{aligned}\hat{w}_0 &= -5a \\ \hat{w}_1 &= a\end{aligned}$$

$$\lambda = v$$

$$(\leq)$$

$$\vee$$

Outline for October 10

- Recap SGD for logistic regression
- **Regularization**
- Multi-class logistic regression
- Begin: evaluation metrics

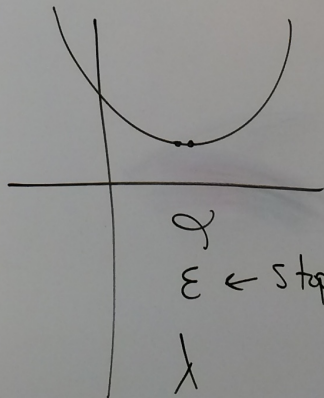
$$\begin{matrix} = 1 \\ y = 0 \end{matrix}$$

Regularization

$$J^R(w) = \left[\sum_{i=1}^n y_i \log h_{\vec{w}}(\vec{x}_i) + (1-y_i) \log (1-h_{\vec{w}}(\vec{x}_i)) \right] + \frac{\lambda}{2} \sum_{j=1}^p w_j^2$$

λ = regularization parameter
(hyper parameter)
(small & positive)

don't
regularize
bias



gradient

$$\nabla_{\vec{x}_i} J^R(\vec{w}) = (h_{\vec{w}}(\vec{x}_i) - y_i) \vec{x}_i + \lambda \vec{w}$$

$$\begin{bmatrix} 0 \\ \vec{w} \end{bmatrix}$$

SGD (regularization)

$$\vec{w} \leftarrow \vec{w} - \eta ((h_{\vec{w}}(\vec{x}_i) - y_i) \vec{x}_i + \lambda \vec{w}^*)$$

$$\leftarrow (1 - \underbrace{\eta \lambda}) \vec{w} - \dots$$

less than
1 (pos)

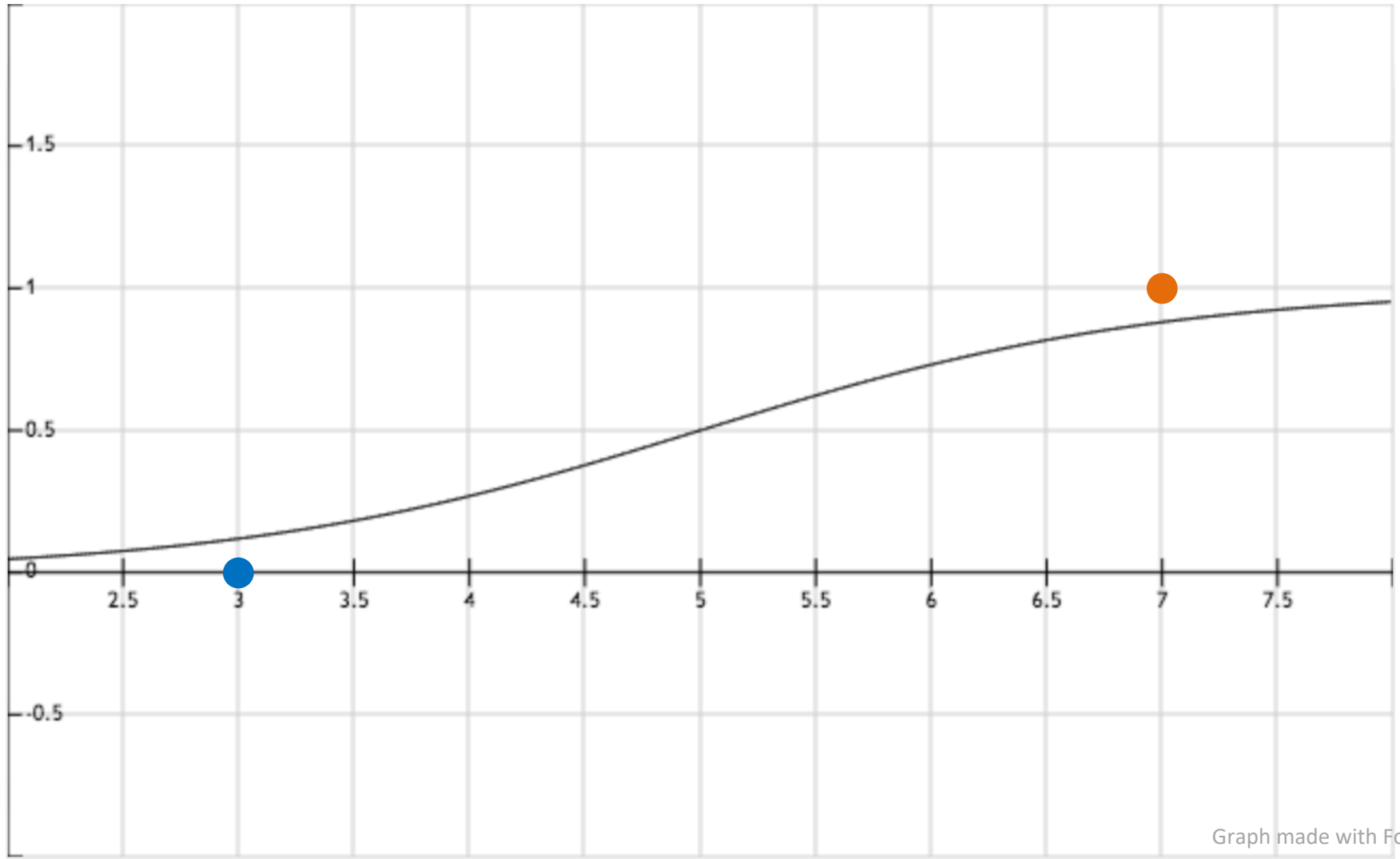
make magnitude
of w 's smaller
all non-bias
terms.

bias

$$w_0 \leftarrow w_0 - \eta \left(\begin{matrix} \text{gradient} \\ \vdots \end{matrix} \right)$$

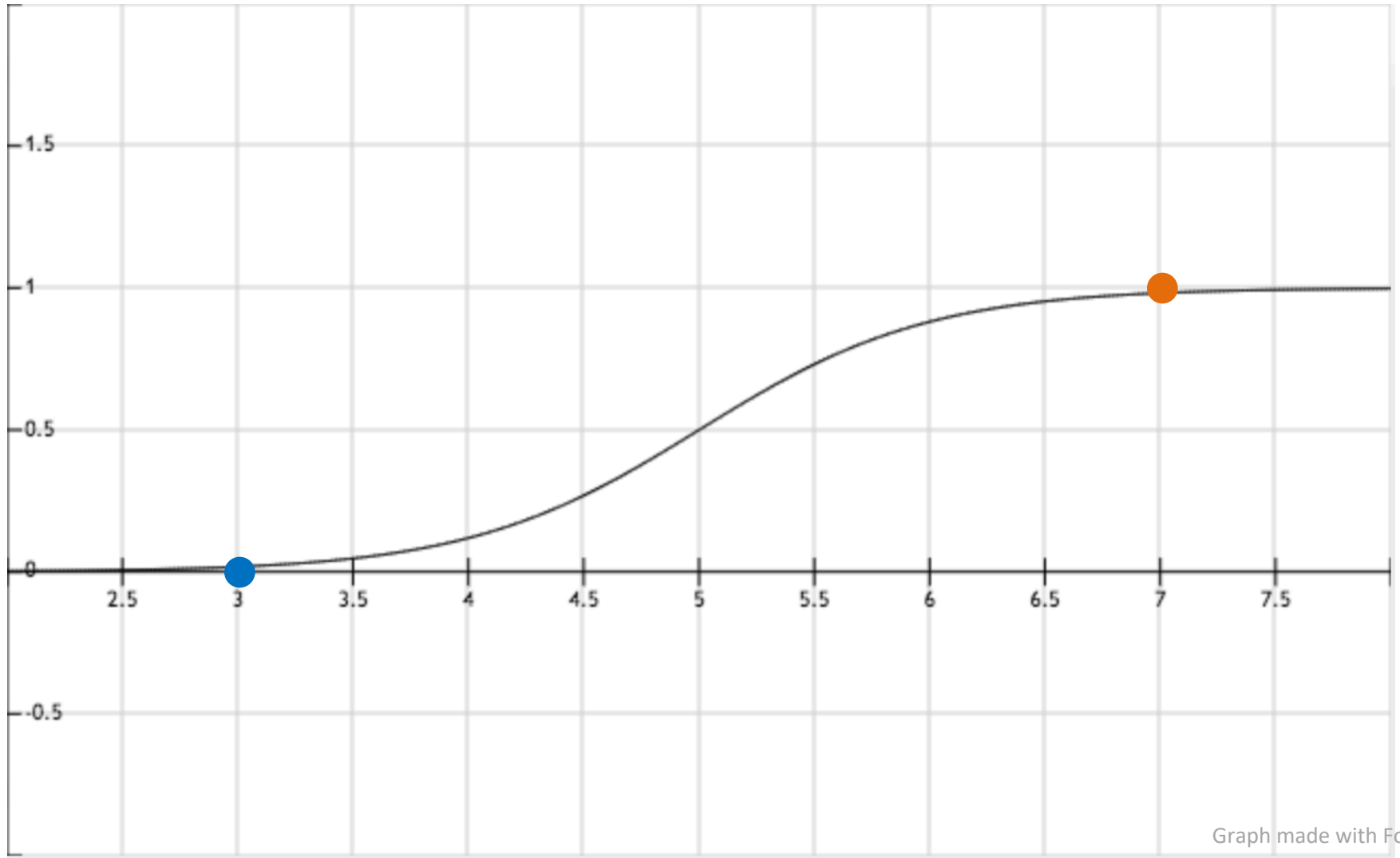
$$w_0 = -5, w_1 = 1$$

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



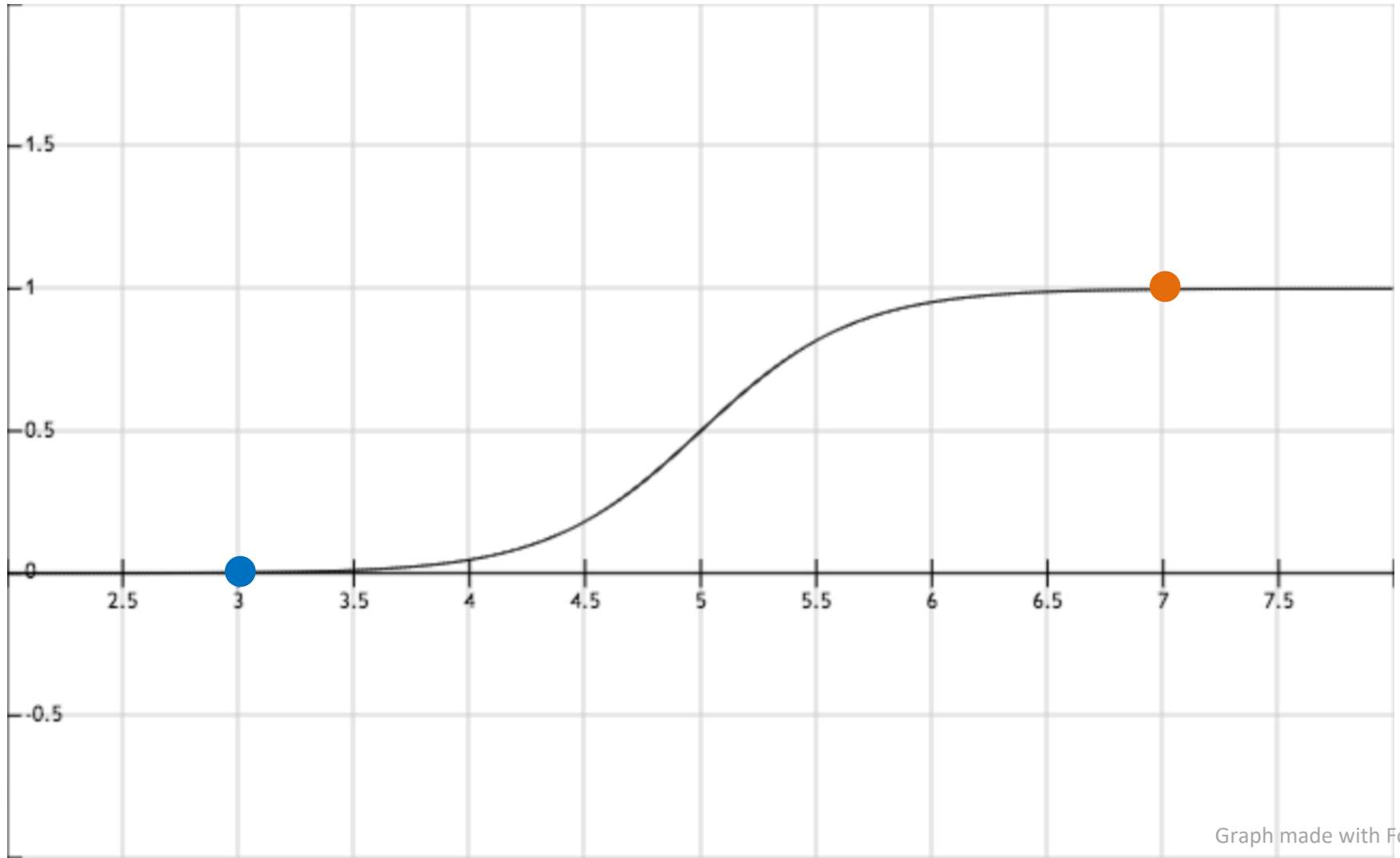
$$w_0 = -10, w_1 = 2$$

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

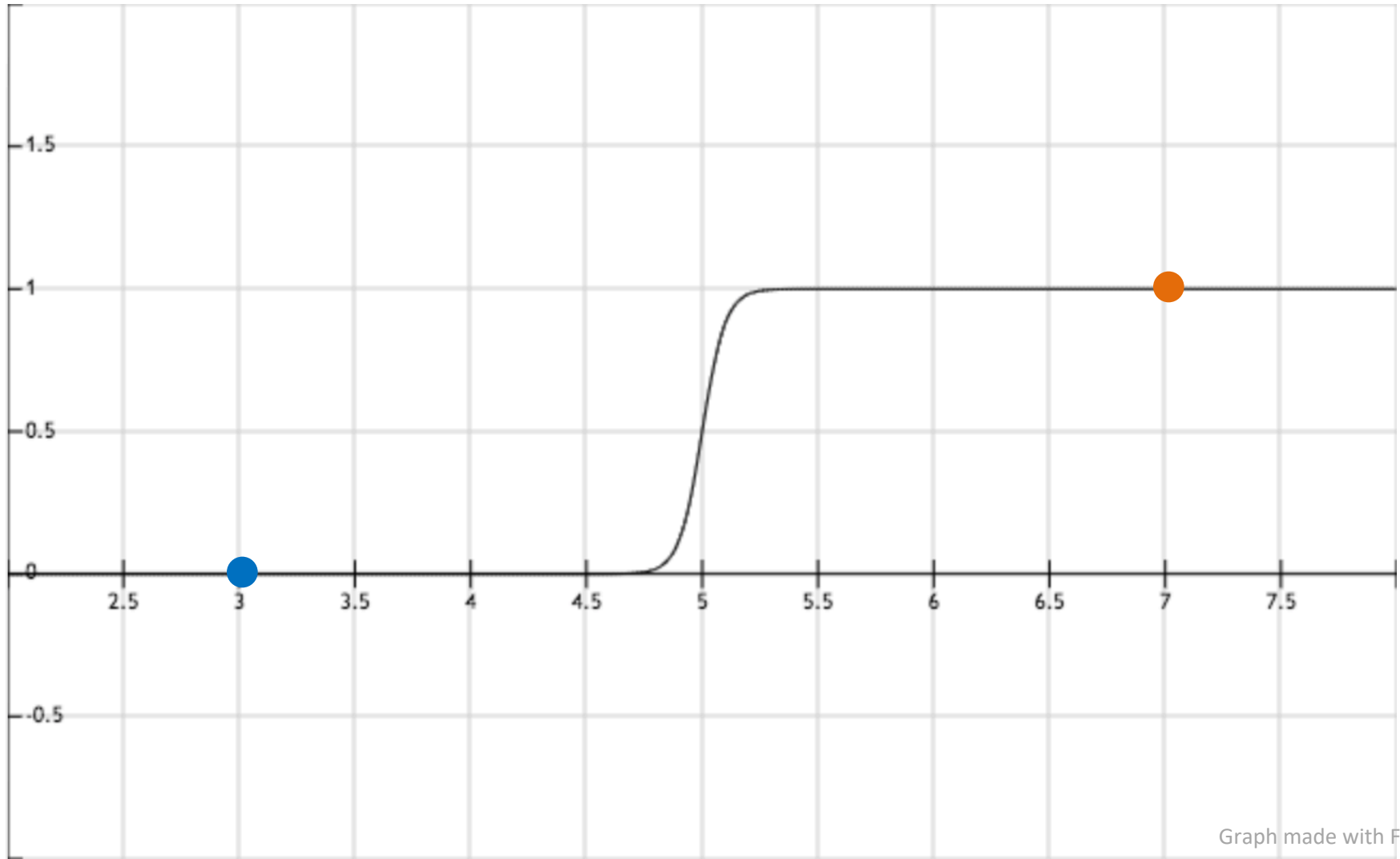


$$w_0 = -15, w_1 = 3$$

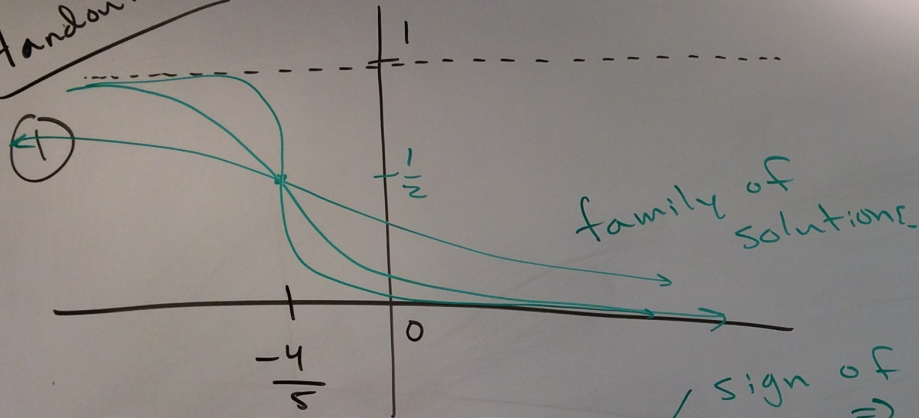
$$h_w(x) = 1/(1+e^{(15-3x)})$$



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



Handout 11



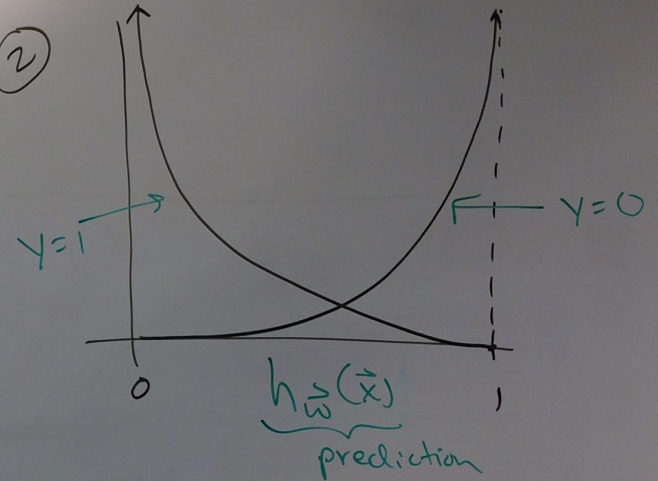
$$\begin{aligned}\hat{w}_0 &= -4 \cdot a \\ \hat{w}_1 &= -5 \cdot a\end{aligned}$$

$$\boxed{-4 - 5x > 0} \quad \left. \begin{array}{l} \text{sign of } w_1 \\ \Rightarrow \text{slope} \end{array} \right\} \text{pred 1}$$

$$-5x > 4$$

$$\star \boxed{x < -\frac{4}{5}} \Rightarrow \text{pred 1}$$

②



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

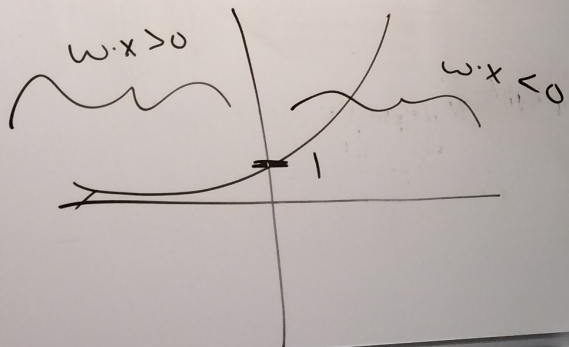
Scalar

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

$$\log(1) > \log(e^{-\vec{w} \cdot \vec{x}})$$

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

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



#3

$$g'(z) = -1 (\text{inside})^{-2} \text{ (derivative of inside)}$$

$$= \frac{e^{-z}}{(1 + e^{-z})^2}$$

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

$$\textcircled{4} \begin{bmatrix} 0.05 \\ 0.35 \end{bmatrix}$$

$$\textcircled{c} \begin{bmatrix} -0.025 \\ 0.125 \end{bmatrix}$$

$$= g(z) (1 - g(z))$$

Outline for October 10

- Recap SGD for logistic regression
- Regularization
- **Multi-class logistic regression**
- Begin: evaluation metrics

Multiclass Logistic Regression

K classes (political party, blood groups)

2 classes: $h_{\vec{w}}(\vec{x}) = \frac{1}{1 + e^{-\vec{w} \cdot \vec{x}}} = \frac{e^{\vec{w} \cdot \vec{x}}}{e^{\vec{w} \cdot \vec{x}} + 1}$

weight on class (under $e^{\vec{w} \cdot \vec{x}}$)
weight on class (under 1)

K classes

$$\hat{y} = h_w(\vec{x}) = \begin{bmatrix} p(y=1|\vec{x}) \\ p(y=2|\vec{x}) \\ \vdots \\ p(y=K|\vec{x}) \end{bmatrix} = \frac{1}{\sum_{k=1}^K e^{\vec{w}^{(k)} \cdot \vec{x}}} \begin{bmatrix} e^{\vec{w}^{(1)} \cdot \vec{x}} \\ e^{\vec{w}^{(2)} \cdot \vec{x}} \\ \vdots \\ e^{\vec{w}^{(K)} \cdot \vec{x}} \end{bmatrix}$$

sum to 1 (bracketed around the denominator)

Weights are a matrix

$$W = \begin{bmatrix} | & | & \dots & | \\ \vec{w}^{(1)} & \vec{w}^{(2)} & \dots & \vec{w}^{(K)} \\ | & | & \dots & | \end{bmatrix}$$

$(p+1) \times K$

cost

$$J(W) = - \underbrace{\sum_{i=1}^n \sum_{k=1}^K y_{ik} \log p(y_i = k | \vec{x}_i)}_{\hat{y}_{ik}}$$

cross entropy

$\vec{y}_i = \begin{bmatrix} 0 \\ 0 \\ \vdots \\ 1 \\ \vdots \\ 0 \end{bmatrix}$ ← index of true label k

$y_i = 4$

$K=5 \rightarrow \begin{bmatrix} 0 \\ 0 \\ 0 \\ 0 \\ 1 \\ 0 \end{bmatrix}$

Red arrows point from y_{i4} to the 4th row and from y_{i1} to the 1st row.

