

CS 360: Machine Learning

Sara Mathieson, Sorelle Friedler

Spring 2024



HVERFORD
COLLEGE

Admin

- **Lab 4** due TONIGHT!
- No office hours today
- **Lab 5** posted tonight
- Review session next week in class (study guide out on Tuesday)

Outline for Feb 22

- Extending logistic regression to multi-class classification (softmax)
- Regularization introduction
- Fairness regularization

Outline for Feb 22

- Extending logistic regression to multi-class classification (softmax)
- Regularization introduction
- Fairness regularization

Sorry, this is s_k !

Multi-class Logistic Regression

$$\vec{y} = \frac{1}{\sum_{k=1}^K e^{\vec{w}^{(k)} \cdot \vec{x}}}$$
$$\left[\begin{array}{c} e^{\vec{w}^{(1)} \cdot \vec{x}} \\ e^{\vec{w}^{(2)} \cdot \vec{x}} \\ \dots \\ e^{\vec{w}^{(K)} \cdot \vec{x}} \end{array} \right]$$

$K = \# \text{ classes}$

$$s_k = \frac{e^{\vec{w}^{(k)} \cdot \vec{x}}}{\sum_{l=1}^K e^{\vec{w}^{(l)} \cdot \vec{x}}} = p(y_i = k | \vec{x}_i)$$

$$\rightarrow [-10, 2, 5] \quad K=3$$

~~$$s_2 = \frac{2}{-10+2+5}$$~~

$$s_2 = \frac{e^2}{e^{-10} + e^2 + e^5}$$

~~$$s_1 = \frac{-10}{-10+2+5}$$~~

Cost function

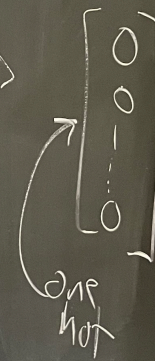
$$J(W) = \sum_{i=1}^n \sum_{k=1}^K y_{ik} \log(p(y_i=k | \vec{x}_i))$$

y_{ik} true label

logistic function

$(P+1) \times K$ cross entropy

$$H(p, q) = - \sum_x p(x) \log q(x)$$

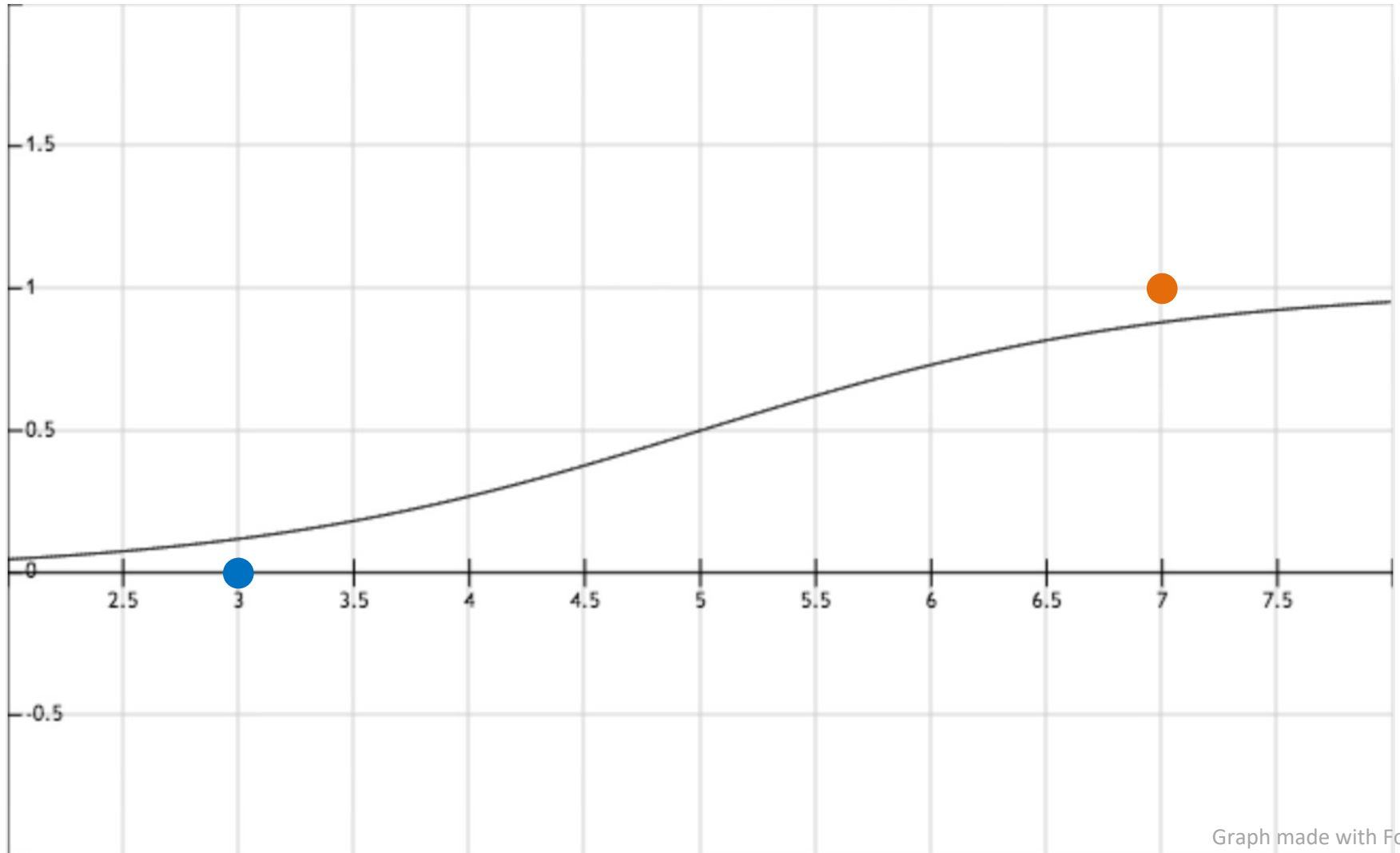


Outline for Feb 22

- Extending logistic regression to multi-class classification (softmax)
- **Regularization introduction**
- Fairness regularization

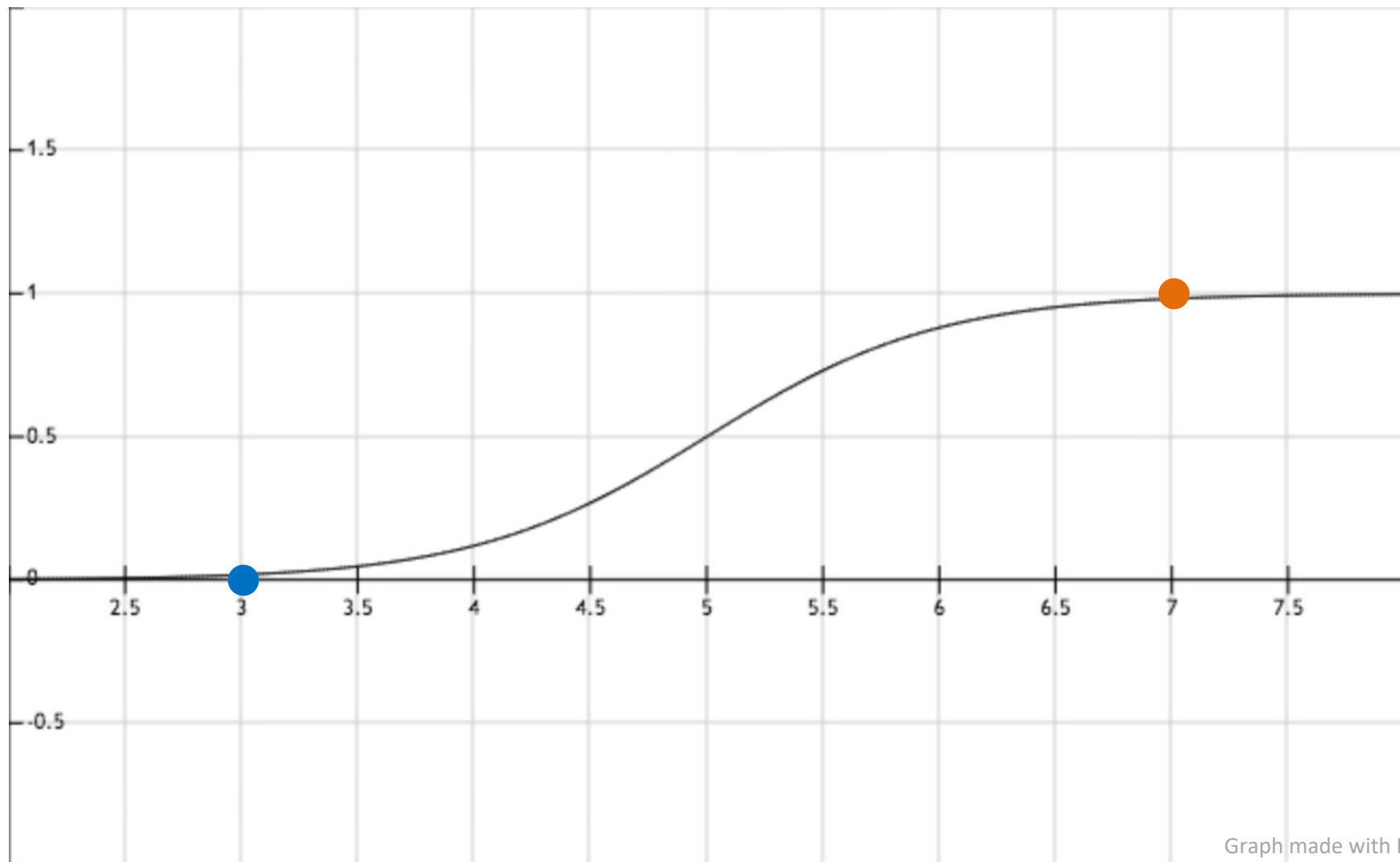
$$w_0 = -5, w_1 = 1$$

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



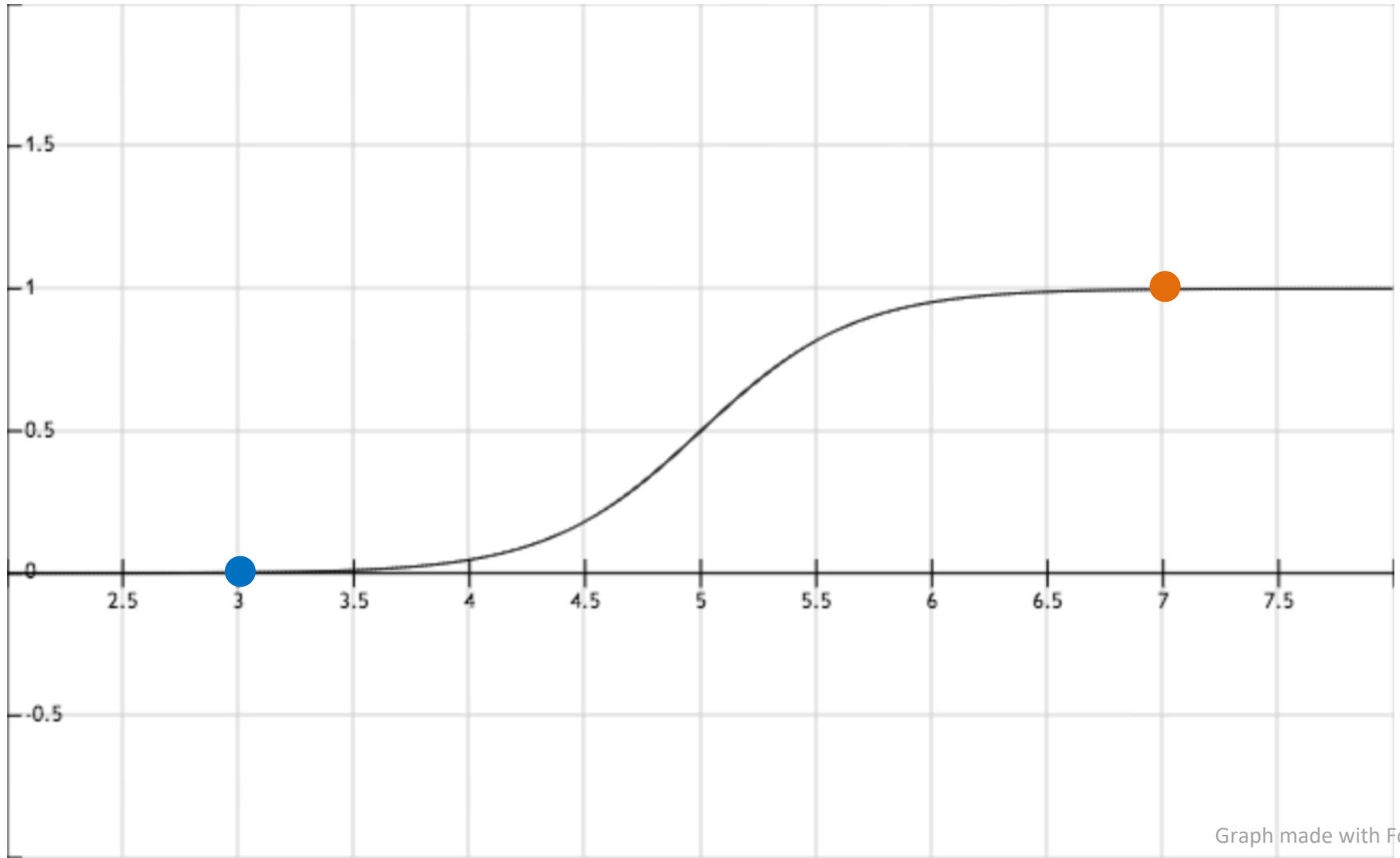
$$w_0 = -10, w_1 = 2$$

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

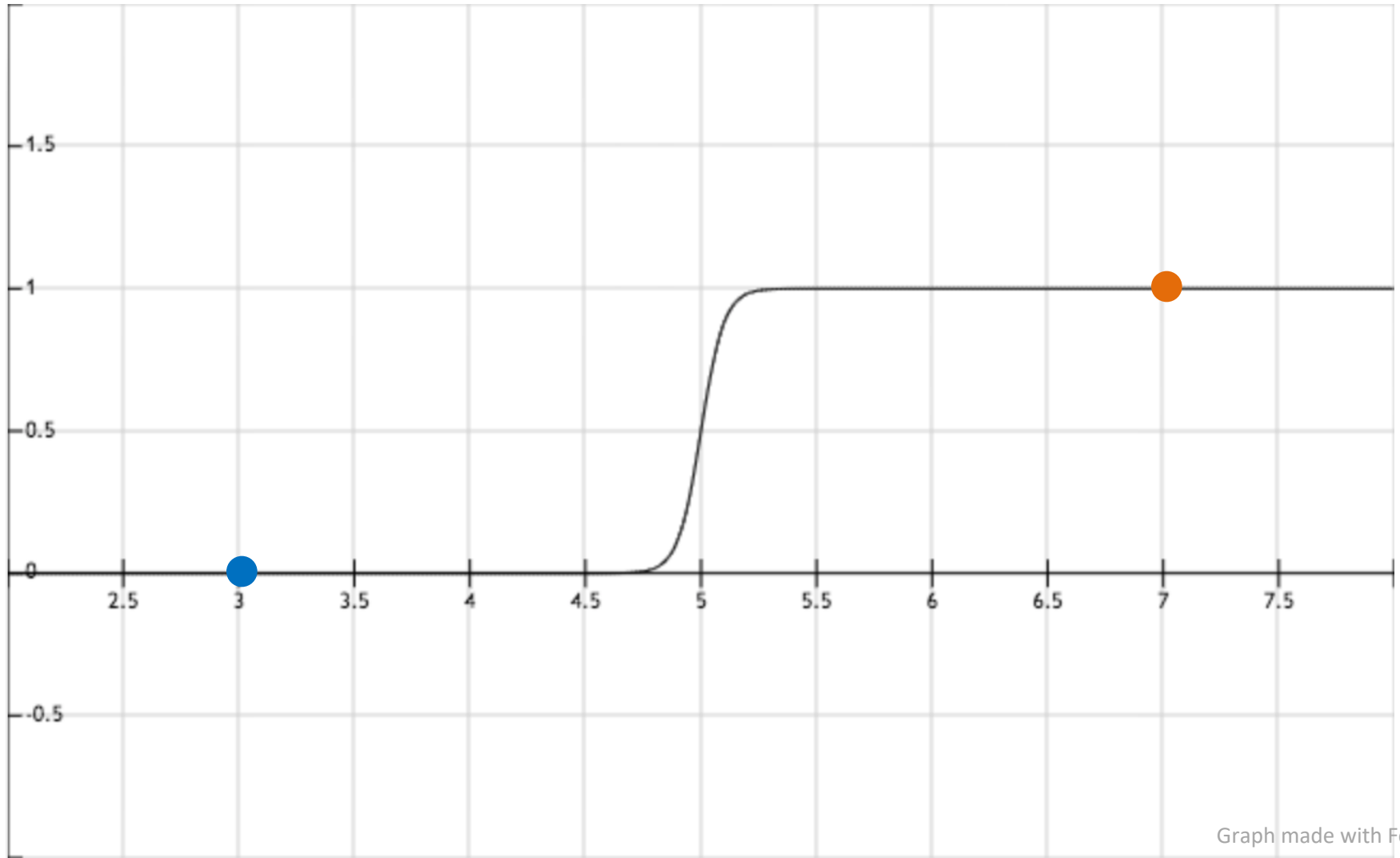


$$w_0 = -15, w_1 = 3$$

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



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



$$\boxed{P=1} \quad h_{\vec{w}}(x) = \frac{1}{1 + e^{-(w_0 + w_1 x)}} \geq 0.5 \Rightarrow \text{predict } \hat{y}=1$$

$$p(\hat{y}=1|x)$$

$$1 \geq 0.5(1 + e^{-(w_0 + w_1 x)})$$

$$\log(e^{-(w_0 + w_1 x)}) \geq \log(1)$$

$$-(w_0 + w_1 x) \geq 0$$

$$\boxed{x \geq -\frac{w_0}{w_1}}$$

Decision Boundary

$$w_0 = -5$$

$$w_1 = 1$$

$$\Rightarrow \boxed{x \geq -\left(\frac{-5}{1}\right) = 5} \Rightarrow \text{predict } 1$$

$$-\left(\frac{-10}{2}\right) = 5$$

Regularization

w 's small

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

regularization parameter

incentivizes small weights

SGD

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

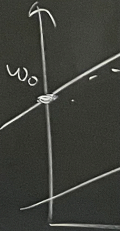
\uparrow \uparrow \uparrow
 $p+1$ $p+1$ $*$

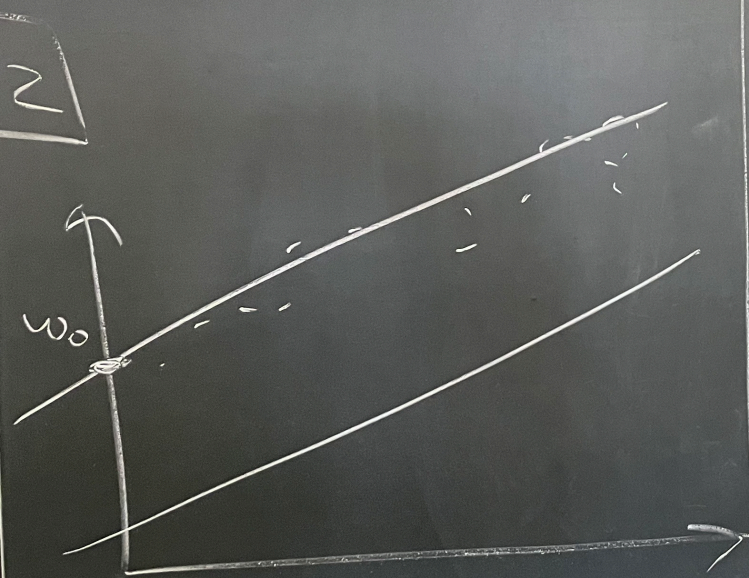
$0 < \lambda < 1$

ex: $\lambda = 0.02$

$\begin{bmatrix} 0 \\ 3 \\ 3 \\ \dots \\ 3 \end{bmatrix}$

$\vec{w} = [w_0, w_1, w_2, \dots, w_p]$
 \uparrow
 not regularizing

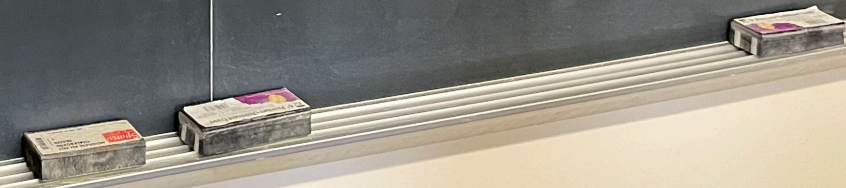




$$\vec{w} \leftarrow \vec{w} - \alpha \left[(h_{\vec{w}}(\vec{x}_i) - y_i) \vec{x}_i + \lambda \vec{w}^* \right]$$

$$\leftarrow (1 - \underbrace{\alpha \lambda}_{\substack{\text{small} \\ \& \text{positive}}}) \vec{w} + \dots$$

Handout 10
 Q1



Handout 10, Q1

$$g'(z) = +1(1 + e^{-z})^{-2} e^{-z}$$

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

\neq

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

adding 0

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

Outline for Feb 22

- Extending logistic regression to multi-class classification (softmax)
- Regularization introduction
- **Fairness regularization**

Fairness Regularization setup

Fairness Regularization

- features $\Rightarrow X$
- labels $\Rightarrow Y \in \{0, 1\}$
not hired \uparrow hired
admitted
- sensitive attributes
 $A \in \{0, 1\}$ race, sex
unprotected \uparrow protected
- prediction $\hat{Y} \in \{0, 1\}$
 \rightarrow accurate wrt Y
 \rightarrow fair wrt A

demographic parity

$$P(\hat{Y}=1 | A=1)$$

$$P(\hat{Y}=1 | A=0)$$

≈ 1
goal
(fairness)

equalized odds

$$P(\hat{Y}=1 | A=1, Y=y)$$

$$P(\hat{Y}=1 | A=0, Y=y)$$

≈ 1
goal

for
 $y \in \{0, 1\}$

$Y=0 \Rightarrow$ FPR per demographic

$Y=1 \Rightarrow$ TPR " "

regularization term (demographic parity)

$$R(h_w, D) = 1 - P(\hat{Y}=1 | A=1)$$

log reg data (train)

want small

$$J^R(\vec{w}) = \underbrace{\text{cost}}_{\text{accuracy}} + \underbrace{R(h_w, D)}_{\text{fairness}}$$

$$1 - \frac{1}{|D_0|} \sum_{\vec{x} \in D_0} p(\hat{y}=1 | \vec{x})$$

given $A=1$ $p(\hat{y}=1)$

$D_0 =$ all examples
in the protected group
($A=1$)

$$1 - \frac{1}{|D_0|} \sum_{\vec{x} \in D_0} h_{\vec{w}}(\vec{x})$$

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

$$\vec{w} \leftarrow \vec{w} - \eta \left[(h - y) \vec{x} \cdot \frac{1}{|D_0|} h(1-h) \vec{x} \right]$$

if $A=1$
for \vec{x}_i

gradient

$$\nabla h_{\vec{w}}(\vec{x}) = h_{\vec{w}}(\vec{x}) (1 - h_{\vec{w}}(\vec{x})) \vec{x}$$

SGD add this on! to the gradient

Handout 10
Q 243

Handout 10, Q2 & Q3

$A=1$ non-men

$A=0$ men

$Y=0$ not hired

$Y=1$ hired

$$P(\hat{Y}=1 | A=1) = \frac{170 + 56}{542 + 170 + 23 + 56}$$
$$= 0.29$$

$$P(\hat{Y}=1 | A=0) = \frac{430 + 190}{\# \text{men}} \approx 1$$
$$= 0.24$$