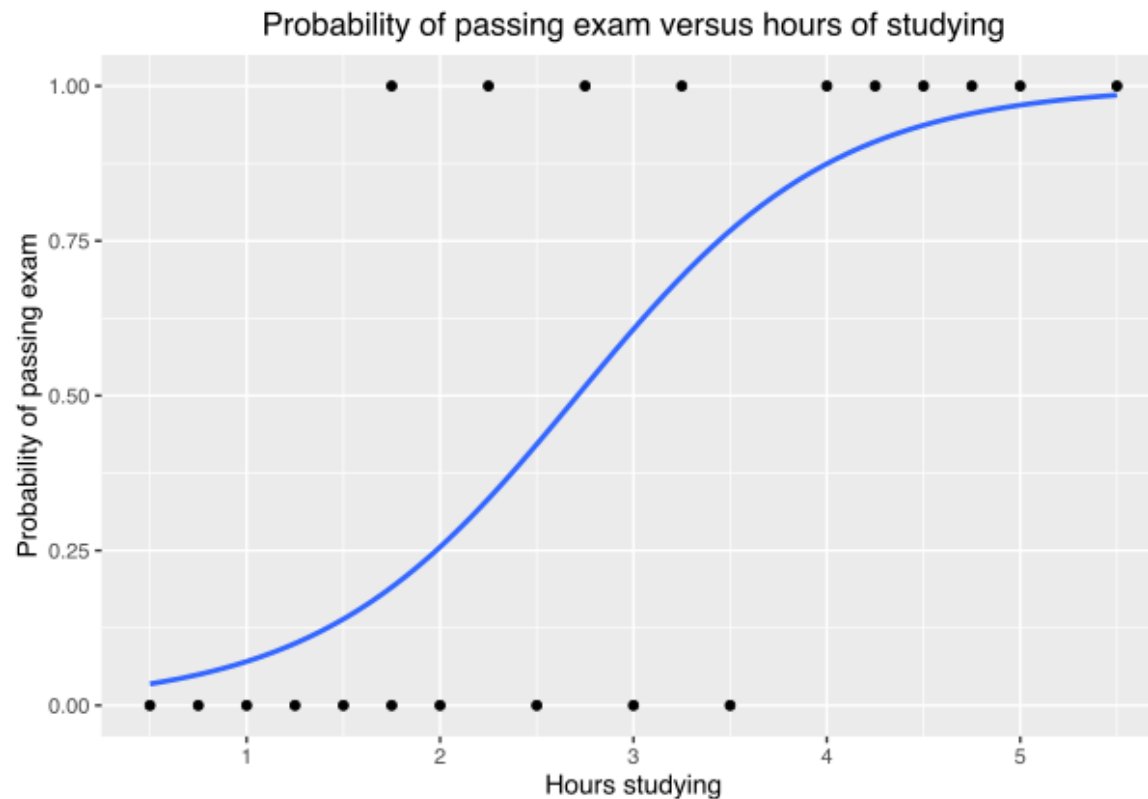


CS440/ECE448 Lecture 8: Logistic Regression

Mark Hasegawa-Johnson, 2/2025

These slides are in the public domain. Re-use, remix, redistribute at will.



CC-SA 4.0, https://commons.wikimedia.org/wiki/File:Exam_pass_logistic_curve.svg

Outline

- Linear regression: Review
- Logistic regression: Output is a probability
- Derivative of the sigmoid
- Derivative of the log sigmoid

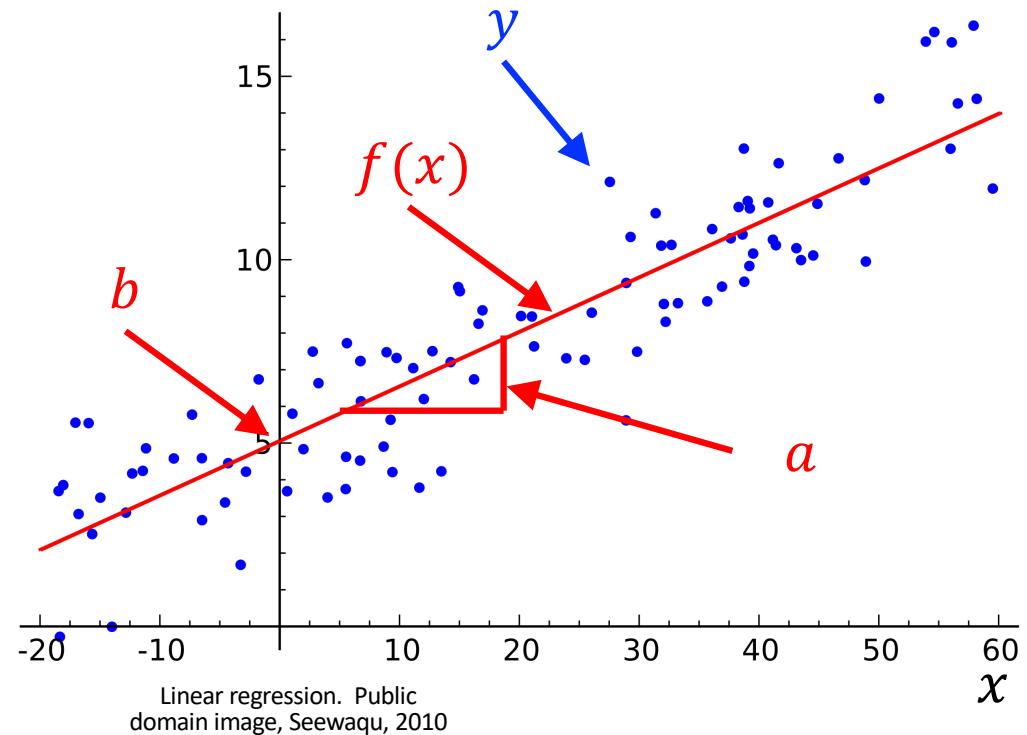
Linear regression

Linear regression is used to estimate a real-valued target variable, y , using a linear function of another variable, x :

$$f(x) = ax + b$$

... so that ...

$$f(x) \approx y$$



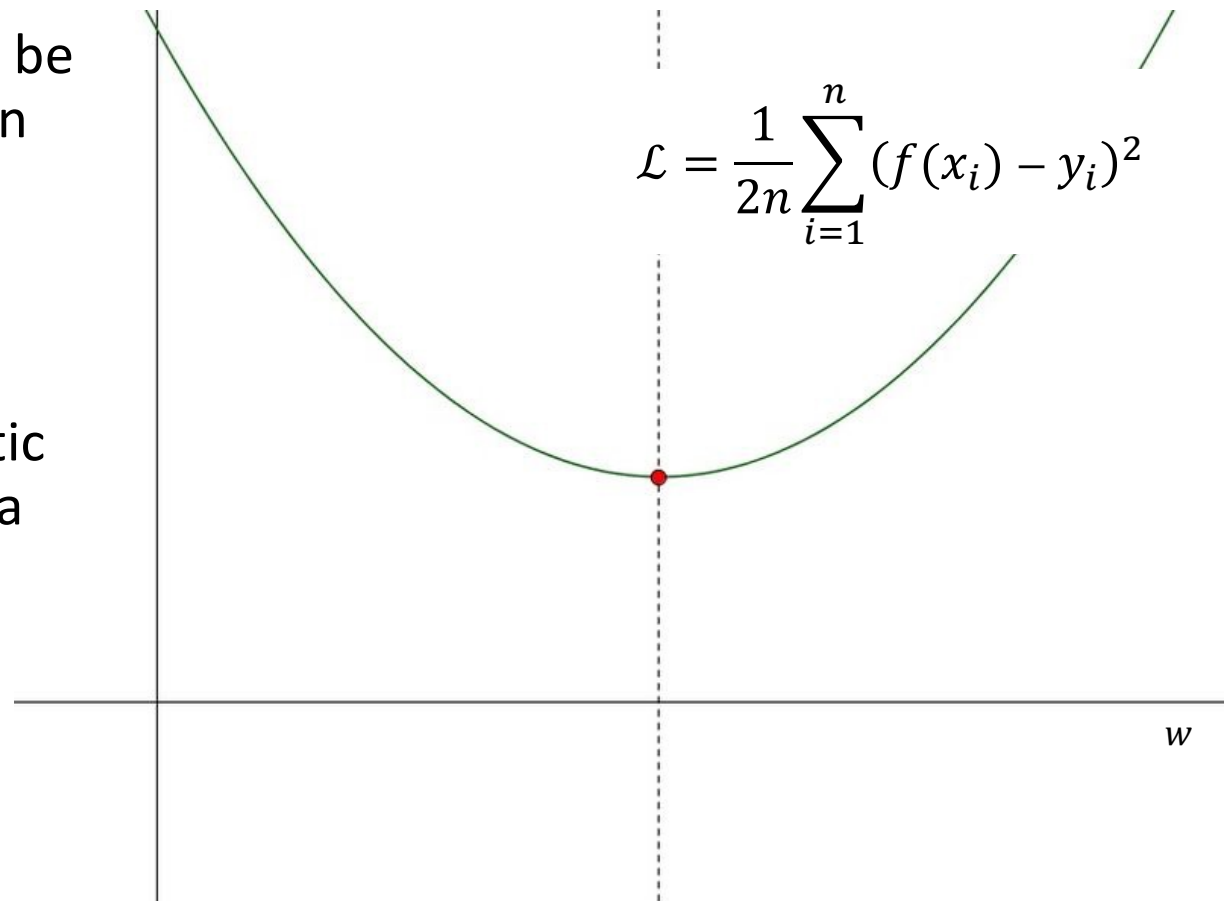
MSE = Parabola

A good closed-form solution can be achieved by minimizing the mean squared error,

$$\mathcal{L} = \frac{1}{2n} \sum_{i=1}^n (f(x_i) - y_i)^2$$

Since it's a non-negative quadratic function of $\mathbf{w} = [a \ b]^T$, it has a unique minimum, given by the closed-form solution

$$\mathbf{w} = (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \mathbf{y}$$



Gradient descent and SGD

Often, closed-form solution is too computationally expensive. In those situations, we choose a random initial guess for the value of \mathbf{w} , and then improve it using either gradient descent:

$$\mathbf{w} \leftarrow \mathbf{w} - \eta \frac{\partial \mathcal{L}}{\partial \mathbf{w}}, \quad \frac{\partial \mathcal{L}}{\partial \mathbf{w}} = \frac{1}{n} \sum_{i=1}^n \epsilon_i \mathbf{x}_i$$

Or stochastic gradient descent:

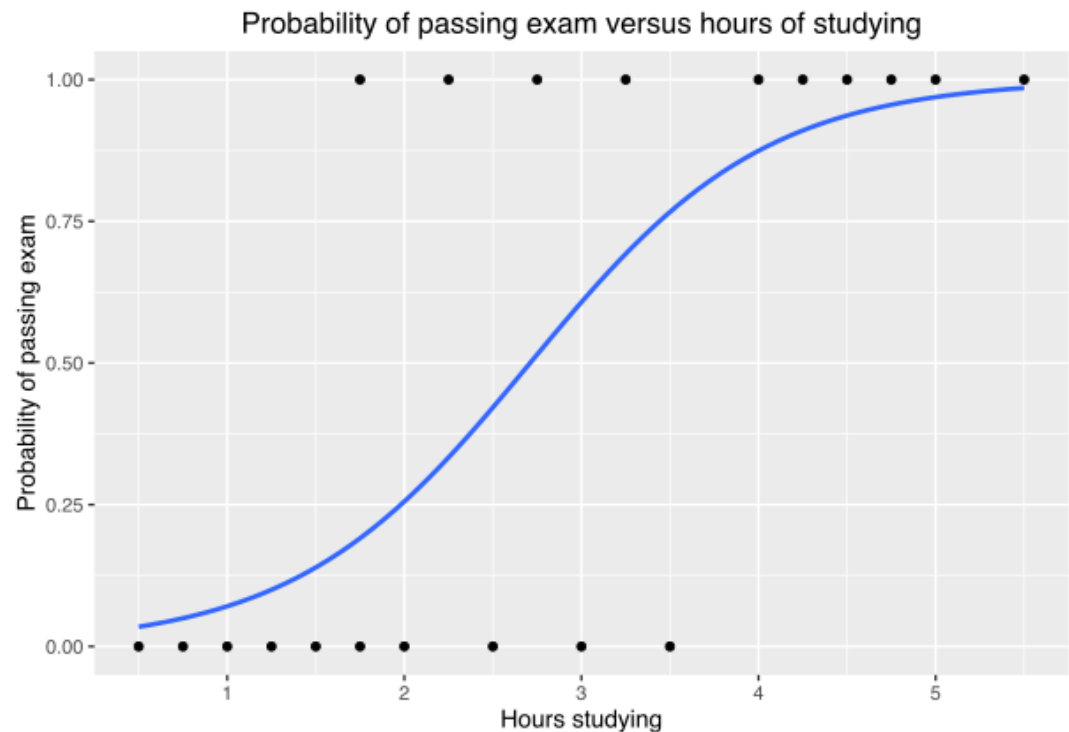
$$\mathbf{w} \leftarrow \mathbf{w} - \eta \frac{\partial \mathcal{L}_i}{\partial \mathbf{w}}, \quad \frac{\partial \mathcal{L}_i}{\partial \mathbf{w}} = \epsilon_i \mathbf{x}_i$$

Outline

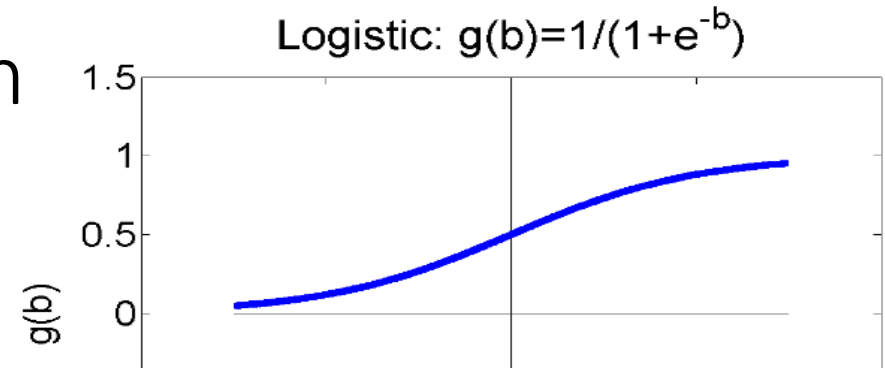
- Linear regression: Review
- Logistic regression: Output is a probability
- Derivative of the sigmoid
- Derivative of the log sigmoid

Logistic regression: Output is a probability

- Logistic regression was invented by psychologists in the early 20th century
- They wanted to model binary outcomes, like “Did student i pass or fail the test?” In other words, every output is either $y_i = 1$ or $y_i = 0$
- Instead of modeling $y_i = \mathbf{x}_i^T \mathbf{w}$ as a real number, it makes more sense to try to model $P(y_i = 1 | \mathbf{x}_i)$ as some kind of function of \mathbf{x}_i .



The logistic sigmoid function



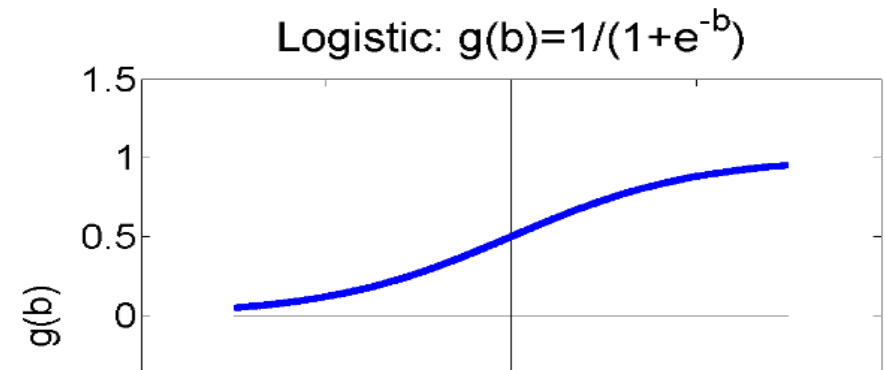
- To model $P(y_i = 1|x_i)$ as some kind of function of x_i , we need some kind of nonlinear function that squashes the linear output $x_i^T \mathbf{w}$ down to the range $0 \leq f(x_i) \leq 1$.
- Psychologists studied many possibilities, but the one most often used today is the logistic sigmoid function:

$$f(x_i) = \sigma(x_i^T \mathbf{w}) = \frac{1}{1 + e^{-x_i^T \mathbf{w}}}$$

This function is called sigmoid because it is S-shaped.

$$\sigma(z) = \begin{cases} 1 & z \rightarrow \infty \\ 0.5 & z = 0 \\ 0 & z \rightarrow -\infty \end{cases}$$

Interpretation as a probability



- Since $0 < f(\mathbf{x}) < 1$, we can interpret $f(\mathbf{x})$ as a probability
- Specifically, we interpret it as $f(\mathbf{x}) = P(Y = 1|X = \mathbf{x})$.

$$f(\mathbf{x}) = \sigma(\mathbf{x}^T \mathbf{w}) = \frac{1}{1 + e^{-\mathbf{x}^T \mathbf{w}}}$$

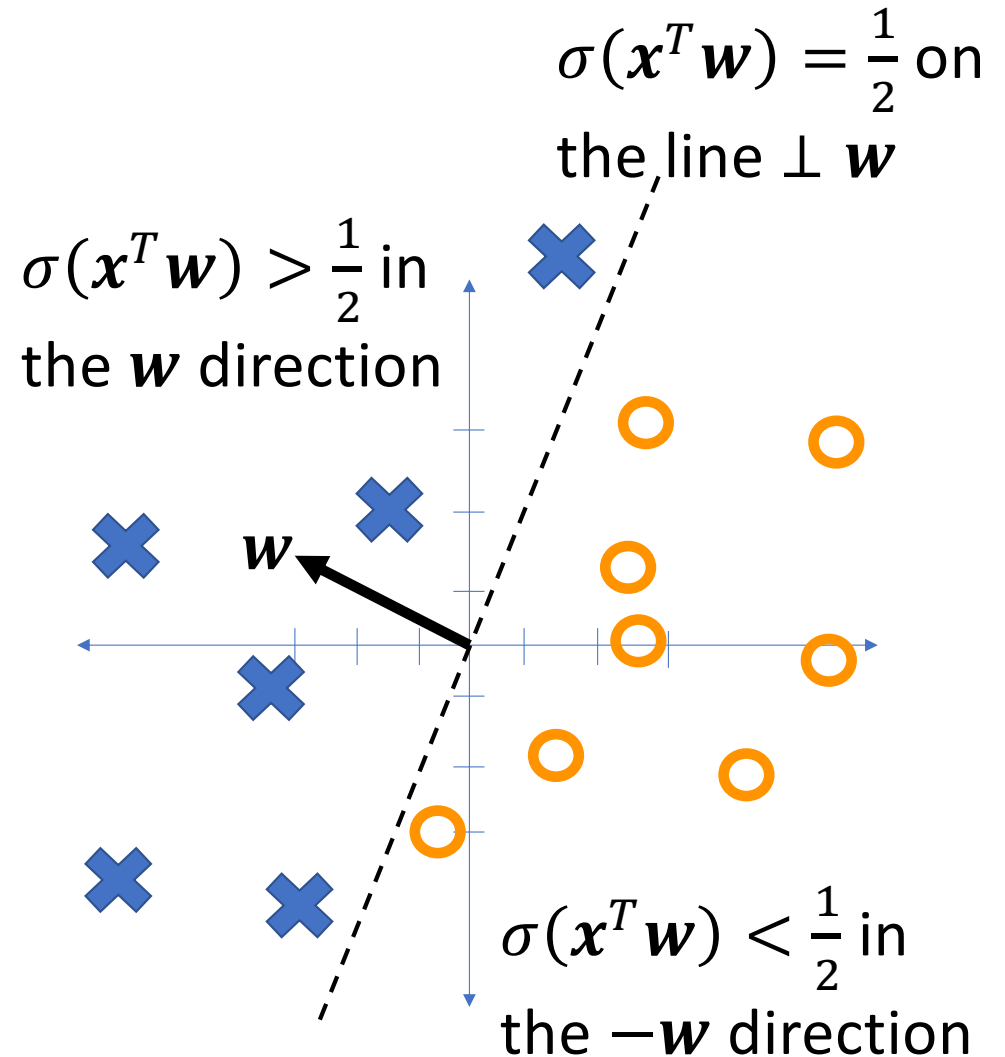
- The argument of the sigmoid, $\mathbf{x}^T \mathbf{w}$, is called the “logit.” Notice that there is a straightforward relationship between the logit and the probability:

$$\sigma(z) = \begin{cases} 1 & z \rightarrow \infty \\ 1/2 & z = 0 \\ 0 & z \rightarrow -\infty \end{cases}$$

Geometric interpretation

Suppose \mathbf{x} is a 2d vector. Then:

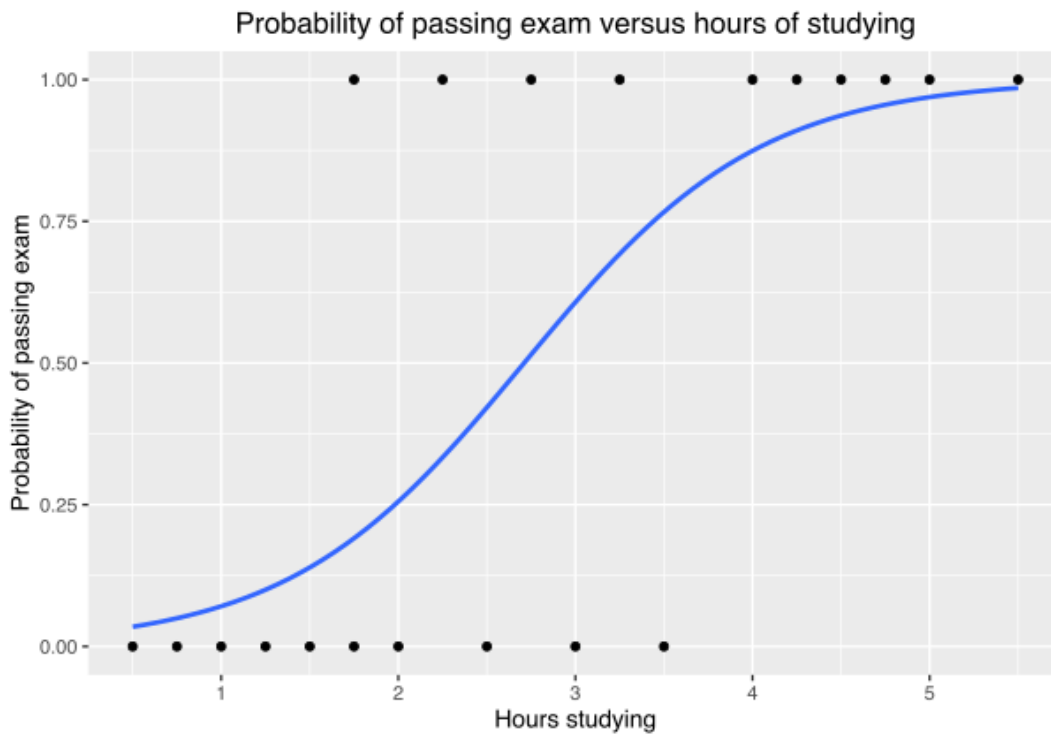
- $\mathbf{x}^T \mathbf{w} = 0$, and $\sigma(\mathbf{x}^T \mathbf{w}) = \frac{1}{2}$, on the line where $\mathbf{x} \perp \mathbf{w}$
- $\mathbf{x}^T \mathbf{w} > 0$, and $\sigma(\mathbf{x}^T \mathbf{w}) > \frac{1}{2}$, in the \mathbf{w} direction
- $\mathbf{x}^T \mathbf{w} < 0$, and $\sigma(\mathbf{x}^T \mathbf{w}) < \frac{1}{2}$, in the $-\mathbf{w}$ direction



Example

For example, the blue line here shows

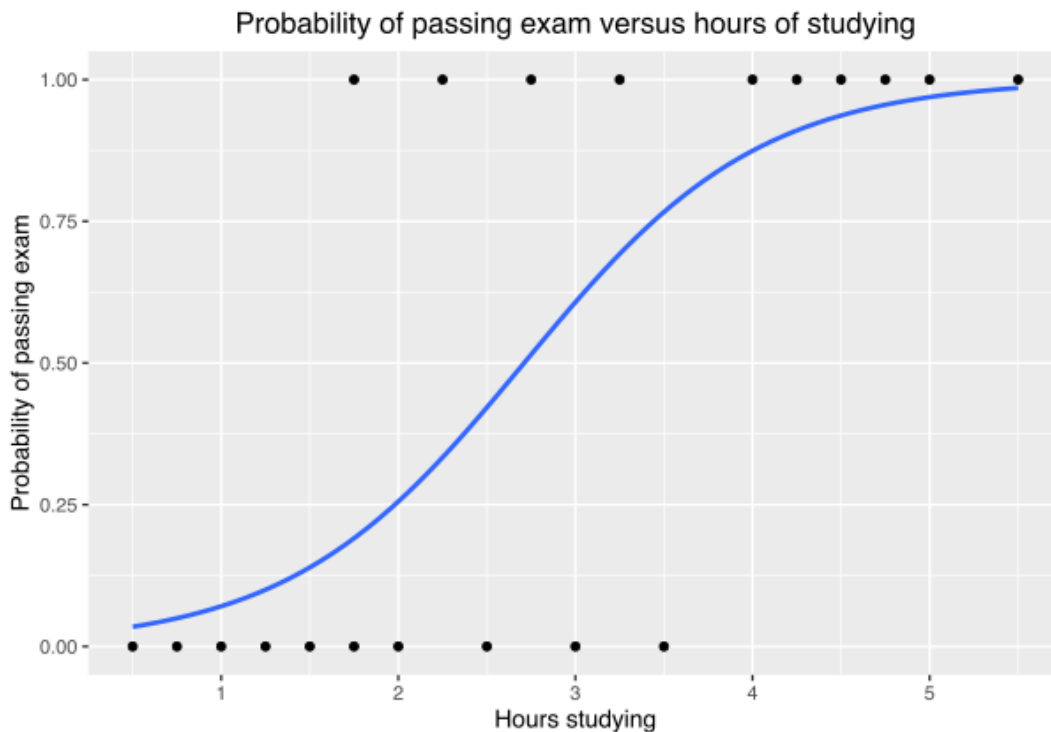
$$P(Y = 1|x) = \sigma(3x - 2.75)$$



Outline

- Linear regression: Review
- Logistic regression: Output is a probability
- Derivative of the sigmoid
- Derivative of the log sigmoid

How is logistic regression trained?



The blue line is the trained classifier.

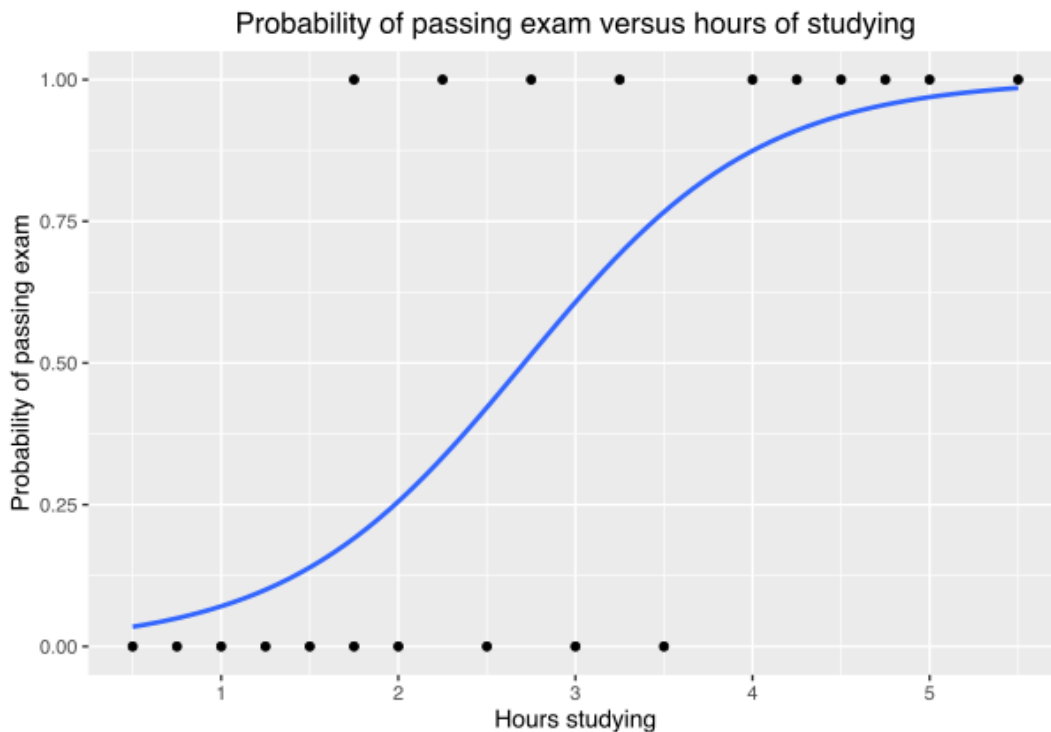
The training data are the black dots:

$$\mathcal{D} = \{(x_1, y_1), \dots, (x_{20}, y_{20})\}$$

$$= \left\{ \begin{array}{l} (0.5,0), (0.75,0), (1.0,0), (1.25,0), \\ (1.5,0), (1.75,0), (2.0,0), (2.5,0), \\ (3.0,0), (3.5,0), (1.75,1), (2.25,1), \\ (2.75,1), (3.25,1), (4.0,1), (4.25,1), \\ (4.5,1), (4.75,1), (5.0,1), (5.5,1) \end{array} \right\}$$

Given that training data, how did we learn the model $P(Y = 1|x) = \sigma(3x - 2.75)$?

First try: SGD



Consider a method based on stochastic gradient descent:

1. Randomly choose a training datum (\mathbf{x}_i, y_i)
2. If $y_i = 1$: adjust \mathbf{w} to increase $\sigma(\mathbf{x}_i^T \mathbf{w})$:

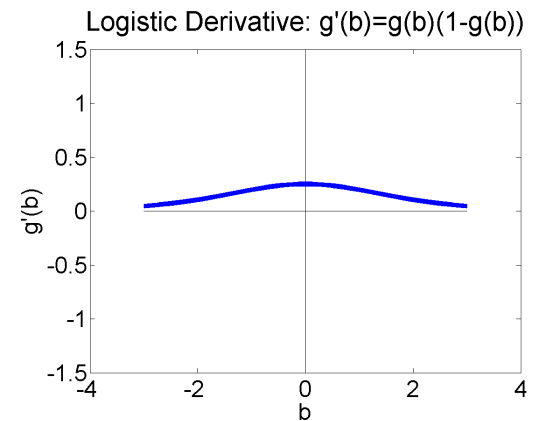
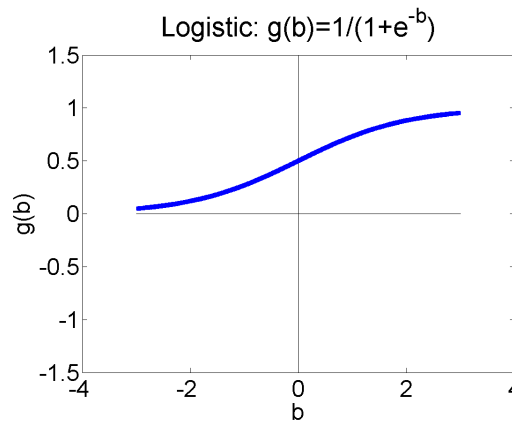
$$\mathbf{w} \leftarrow \mathbf{w} + \eta \frac{\partial \sigma(\mathbf{x}_i^T \mathbf{w})}{\partial \mathbf{w}}$$

3. If $y_i = 0$: adjust \mathbf{w} to decrease $\sigma(\mathbf{x}_i^T \mathbf{w})$:

$$\mathbf{w} \leftarrow \mathbf{w} - \eta \frac{\partial \sigma(\mathbf{x}_i^T \mathbf{w})}{\partial \mathbf{w}}$$

4. Repeat

OK, what's $\frac{\partial \sigma(x_i^T \mathbf{w})}{\partial \mathbf{w}}$?

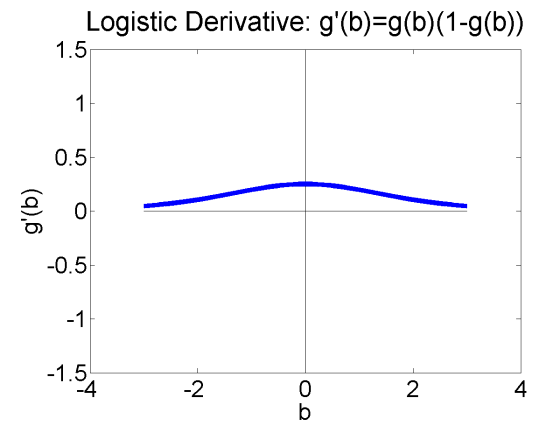
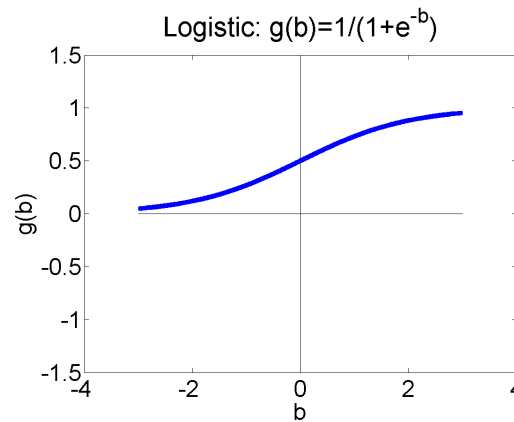


$$\sigma(z) = \frac{1}{1 + e^{-z}}$$

$$\frac{\partial \sigma(z)}{\partial z} = \left(-\frac{1}{(1 + e^{-z})^2} \right) (-e^{-z}) = \left(\frac{1}{1 + e^{-z}} \right) \left(\frac{e^{-z}}{1 + e^{-z}} \right)$$

$$\frac{\partial \sigma(z)}{\partial z} = \sigma(z)(1 - \sigma(z))$$

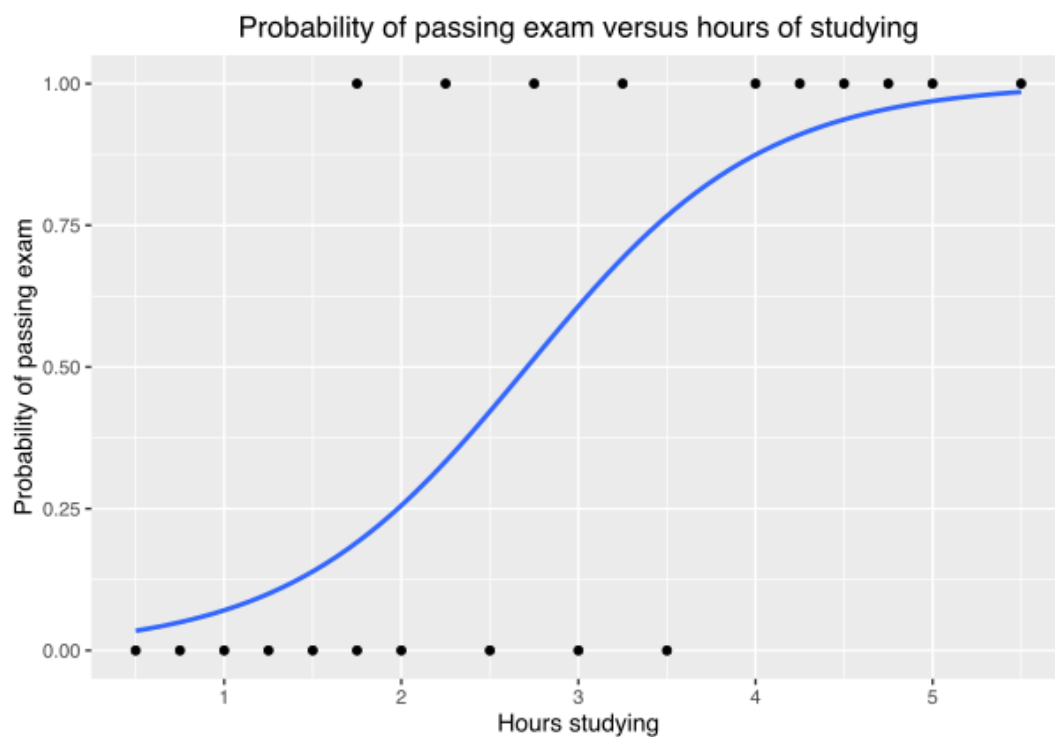
OK, what's $\frac{\partial \sigma(\mathbf{x}_i^T \mathbf{w})}{\partial \mathbf{w}}$?



$$\frac{\partial \sigma(z)}{\partial z} = \sigma(z)(1 - \sigma(z))$$

$$\frac{\partial \sigma(z)}{\partial \mathbf{w}} = \frac{\partial \sigma(z)}{\partial z} \frac{\partial z}{\partial \mathbf{w}} = \sigma(z)(1 - \sigma(z)) \mathbf{x}_i$$

SGD using probabilities



Consider a method based on stochastic gradient descent:

1. Randomly choose a training datum (\mathbf{x}_i, y_i)
2. If $y_i = 1$:
$$\mathbf{w} \leftarrow \mathbf{w} + \eta \sigma(z)(1 - \sigma(z)) \mathbf{x}_i$$
3. If $y_i = 0$:
$$\mathbf{w} \leftarrow \mathbf{w} - \eta \sigma(z)(1 - \sigma(z)) \mathbf{x}_i$$
4. Repeat

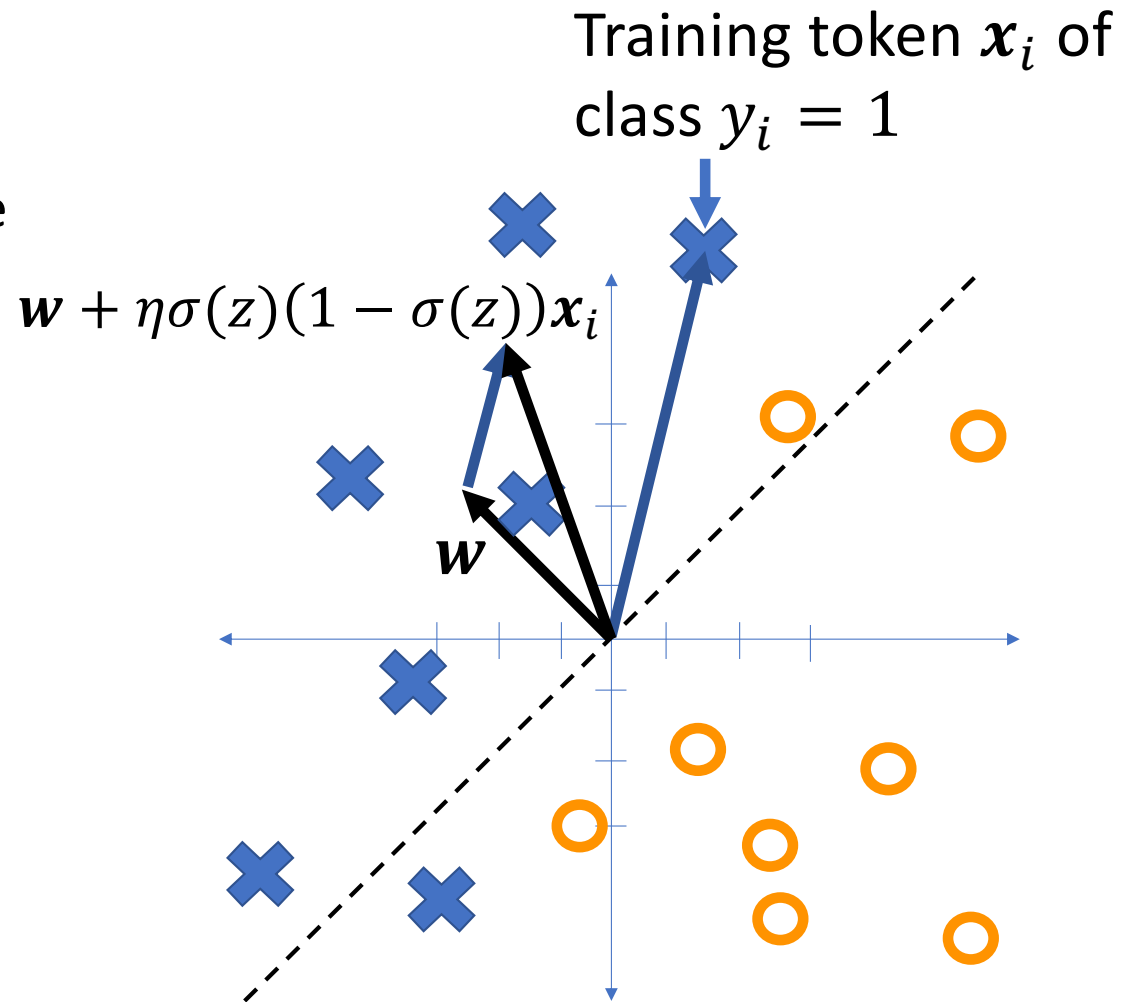
Try the quiz!

- Go to PrairieLearn, try the quiz!

1. Randomly choose a training datum (\mathbf{x}_i, y_i)
2. If $y_i = 1$: adjust \mathbf{w} to increase $\sigma(\mathbf{x}_i^T \mathbf{w})$:

$$\mathbf{w} \leftarrow \mathbf{w} + \eta \frac{\partial \sigma(\mathbf{x}_i^T \mathbf{w})}{\partial \mathbf{w}}$$

$$= \mathbf{w} + \eta \sigma(z)(1 - \sigma(z)) \mathbf{x}_i$$



Outline

- Linear regression: Review
- Logistic regression: Output is a probability
- Derivative of the sigmoid
- Derivative of the log sigmoid

Maximize probability of the training dataset?

- SGD with probabilities is not very systematic
- Here's a systematic approach: Let's try to maximize the probability of the training dataset.
- The training dataset is $\mathcal{D} = \{(\mathbf{x}_1, y_1), \dots, (\mathbf{x}_n, y_n)\}$
- Let's suppose all those examples are chosen independently. Then

$$P(\mathcal{D}) = \prod_{i=1}^n P(Y = y_i | \mathbf{x}_i)$$

- What is $\frac{\partial}{\partial \mathbf{w}} \prod_{i=1}^n P(Y = y_i | \mathbf{x}_i)$? Answer: UGLY!!!!!!

No! Maximize its log probability!

- Notice that whatever value of \mathbf{w} maximizes this:

$$P(\mathcal{D}) = \prod_{i=1}^n P(Y = y_i | \mathbf{x}_i)$$

... will also maximize this:

$$\begin{aligned} \log P(\mathcal{D}) &= \sum_{i=1}^n \log P(Y = y_i | \mathbf{x}_i) \\ &= \sum_{y_i=1} \log \sigma(\mathbf{x}_i^T \mathbf{w}) + \sum_{y_i=0} \log(1 - \sigma(\mathbf{x}_i^T \mathbf{w})) \end{aligned}$$

OK, what's $\frac{\partial \log \sigma(\mathbf{x}_i^T \mathbf{w})}{\partial \mathbf{w}}$?

$$\frac{\partial \log \sigma(\mathbf{x}_i^T \mathbf{w})}{\partial \mathbf{w}} = \frac{1}{\sigma(\mathbf{x}_i^T \mathbf{w})} \frac{\partial \sigma(\mathbf{x}_i^T \mathbf{w})}{\partial \mathbf{w}}$$

$$= \frac{1}{\cancel{\sigma(\mathbf{x}_i^T \mathbf{w})}} \cancel{\sigma(\mathbf{x}_i^T \mathbf{w})} (1 - \sigma(\mathbf{x}_i^T \mathbf{w})) x_i$$

$$= (1 - \sigma(\mathbf{x}_i^T \mathbf{w})) x_i$$

OK, what's $\frac{\partial \log(1 - \sigma(\mathbf{x}_i^T \mathbf{w}))}{\partial \mathbf{w}}$?

$$\frac{\partial \log(1 - \sigma(\mathbf{x}_i^T \mathbf{w}))}{\partial \mathbf{w}} = \frac{1}{(1 - \sigma(\mathbf{x}_i^T \mathbf{w}))} \frac{-\partial \sigma(\mathbf{x}_i^T \mathbf{w})}{\partial \mathbf{w}}$$

$$= \frac{1}{(1 - \cancel{\sigma(\mathbf{x}_i^T \mathbf{w})})} \left(-\cancel{\sigma(\mathbf{x}_i^T \mathbf{w})} (1 - \cancel{\sigma(\mathbf{x}_i^T \mathbf{w})}) \right) \mathbf{x}_i$$

$$= (0 - \sigma(\mathbf{x}_i^T \mathbf{w})) \mathbf{x}_i$$

OK, what's $\frac{\partial \log P(\mathcal{D})}{\partial \mathbf{w}}$?

$$\begin{aligned}\frac{\partial \log P(\mathcal{D})}{\partial \mathbf{w}} &= \sum_{y_i=1} \frac{\partial \log \sigma(\mathbf{x}_i^T \mathbf{w})}{\partial \mathbf{w}} + \sum_{y_i=0} \frac{\partial \log(1 - \sigma(\mathbf{x}_i^T \mathbf{w}))}{\partial \mathbf{w}} \\ &= \sum_{y_i=1} (1 - \sigma(\mathbf{x}_i^T \mathbf{w})) \mathbf{x}_i + \sum_{y_i=0} (0 - \sigma(\mathbf{x}_i^T \mathbf{w})) \mathbf{x}_i\end{aligned}$$

Cross-entropy

This loss function:

$$\mathcal{L}_i = -\ln \Pr(Y = y_i | \mathbf{x}_i)$$

is called cross-entropy. The term comes from physics, where “entropy” is our degree of uncertainty about whether or not something will happen.



CC-SA 4.0,
https://en.wikipedia.org/wiki/File:Ultra_slow-motion_video_of_glass_tea_cup_smashed_on_concrete_floor.webm

Gradient descent

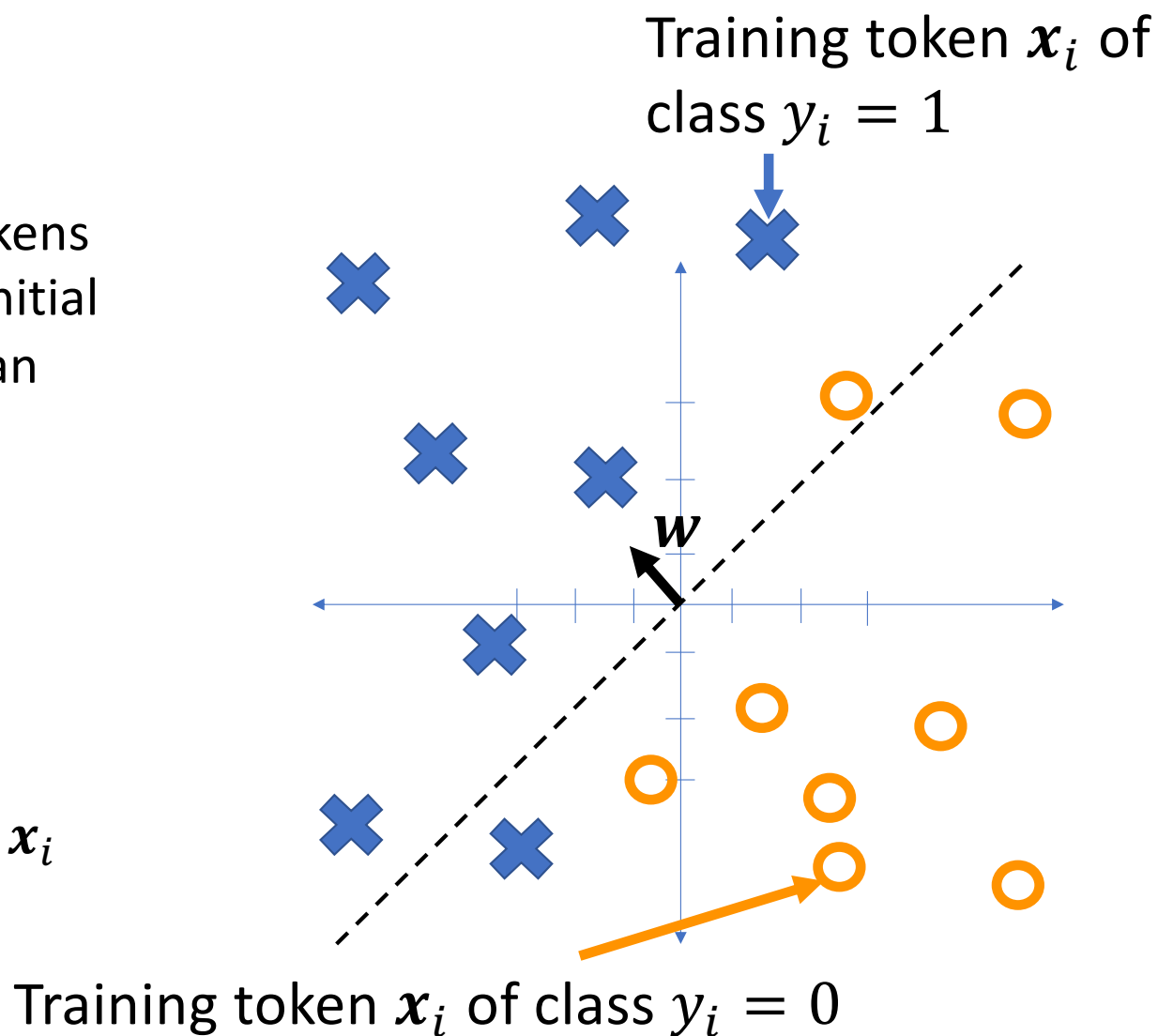
Suppose we have training tokens (\mathbf{x}_i, y_i) , and we have some initial weight vector \mathbf{w} . Then we can update it as

$$\mathbf{w} \leftarrow \mathbf{w} - \eta \frac{\partial \mathcal{L}}{\partial \mathbf{w}}$$

...where...

$$\mathcal{L} = -\log P(\mathcal{D})$$

$$\frac{\partial \mathcal{L}}{\partial \mathbf{w}} = \sum_{i=1}^n (\sigma(\mathbf{x}_i^T \mathbf{w}) - y_i) \mathbf{x}_i$$



Conclusions

- Linear regression

$$\mathbf{w} \leftarrow \mathbf{w} - \eta \frac{\partial \mathcal{L}}{\partial \mathbf{w}}, \quad \frac{\partial \mathcal{L}}{\partial \mathbf{w}} = \frac{1}{n} \sum_{i=1}^n \epsilon_i \mathbf{x}_i$$

- Logistic regression

$$f(\mathbf{x}_i) = \sigma(\mathbf{x}_i^T \mathbf{w}) = \frac{1}{1 + e^{-\mathbf{x}_i^T \mathbf{w}}}$$

- Derivative of the sigmoid

$$\frac{\partial \sigma(z)}{\partial z} = \sigma(z)(1 - \sigma(z))$$

- Derivative of the negative log sigmoid

$$\mathcal{L} = - \sum_{y_i=1} \log \sigma(\mathbf{x}_i^T \mathbf{w}) - \sum_{y_i=0} \log(1 - \sigma(\mathbf{x}_i^T \mathbf{w}))$$