

Lecture 7: Logistic Regression Fall 2022

Kai-Wei Chang
CS @ UCLA

kw+cm146@kwchang.net

The instructor gratefully acknowledges Dan Roth, Vivek Srikuar, Sriram Sankararaman, Fei Sha, Ameet Talwalkar, Eric Eaton, and Jessica Wu, Hal Daume whose slides are heavily used, and the many others who made their course material freely available online.

Announcement

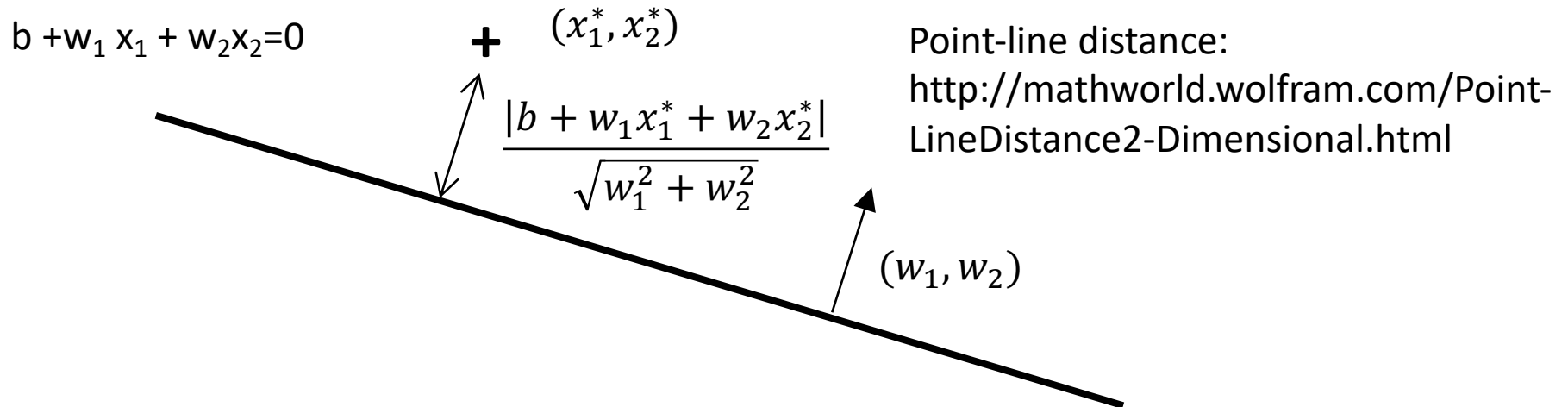
- ❖ Quiz 2 is coming!
- ❖ Suggested readings in BruinLearn
<http://ciml.info/>
- ❖ Please don't put CM146 course materials online

- ❖ Explanation of margin γ

$$y_i(\mathbf{u}^T \mathbf{x}_i + b) \geq \gamma$$

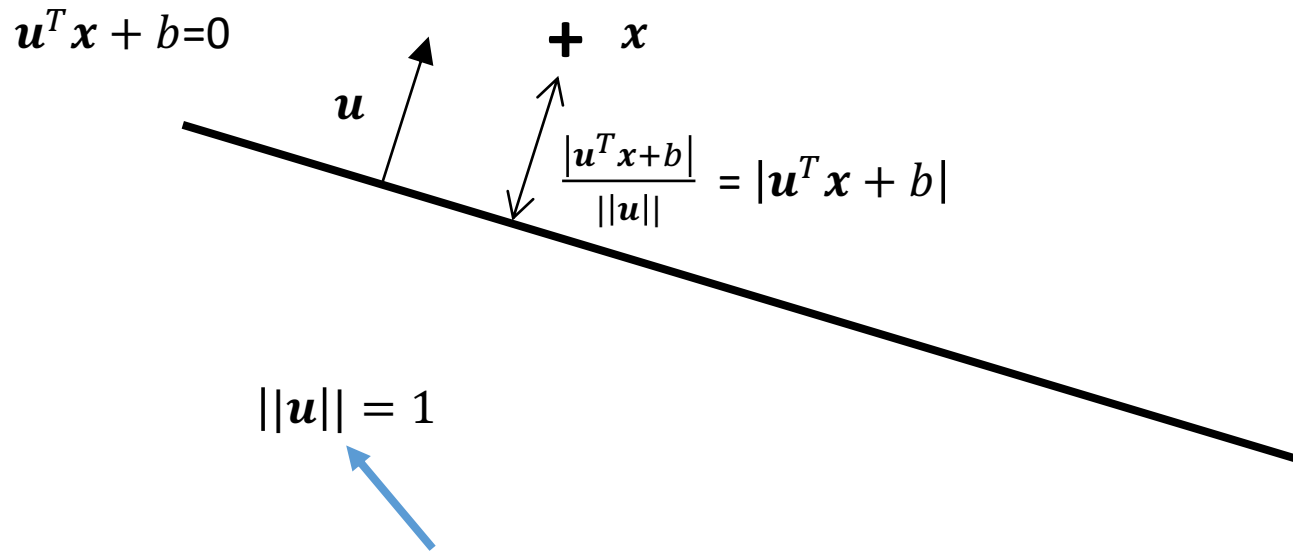
Recall: The geometry of a linear classifier

$$\text{Prediction} = \text{sgn}(b + w_1 x_1 + w_2 x_2)$$



Recall: The geometry of a linear classifier

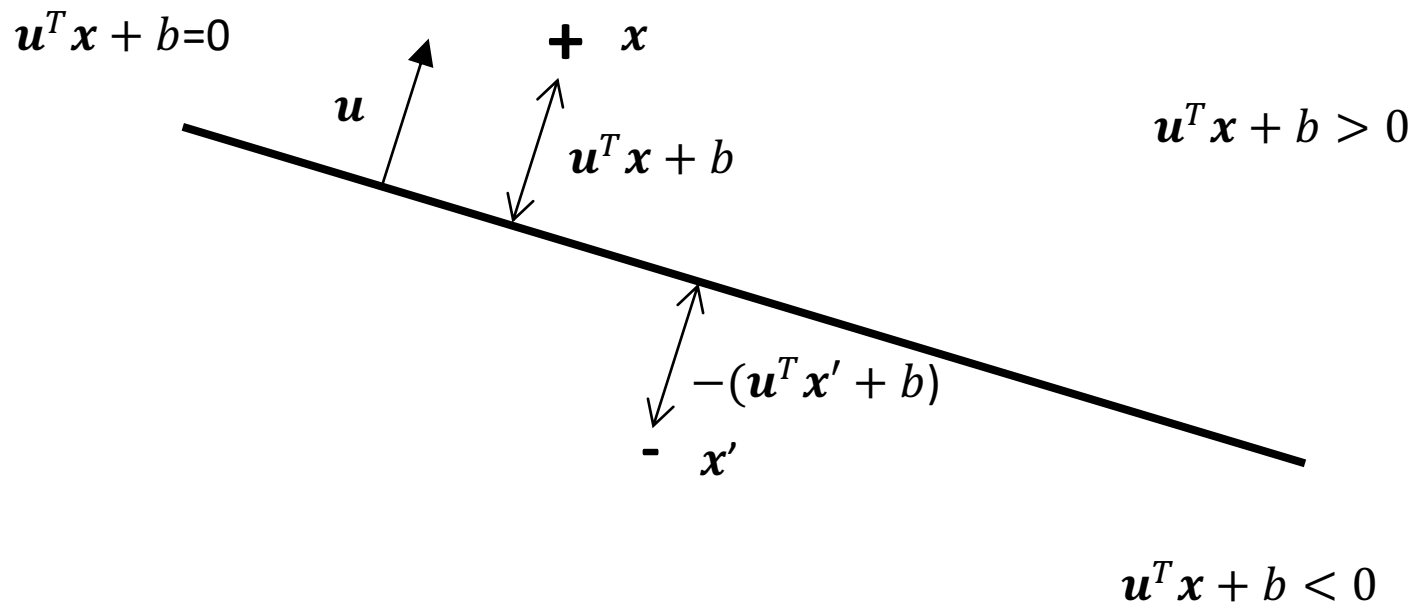
$$\text{Prediction} = \text{sgn}(\mathbf{u}^T \mathbf{x} + b)$$



If we use a **unit** normal vector \mathbf{u} represents the hyperplane, the distance between point \mathbf{x} to plane is $|\mathbf{u}^T \mathbf{x} + b|$ or $y(\mathbf{u}^T \mathbf{x} + b)$

Recall: The geometry of a linear classifier

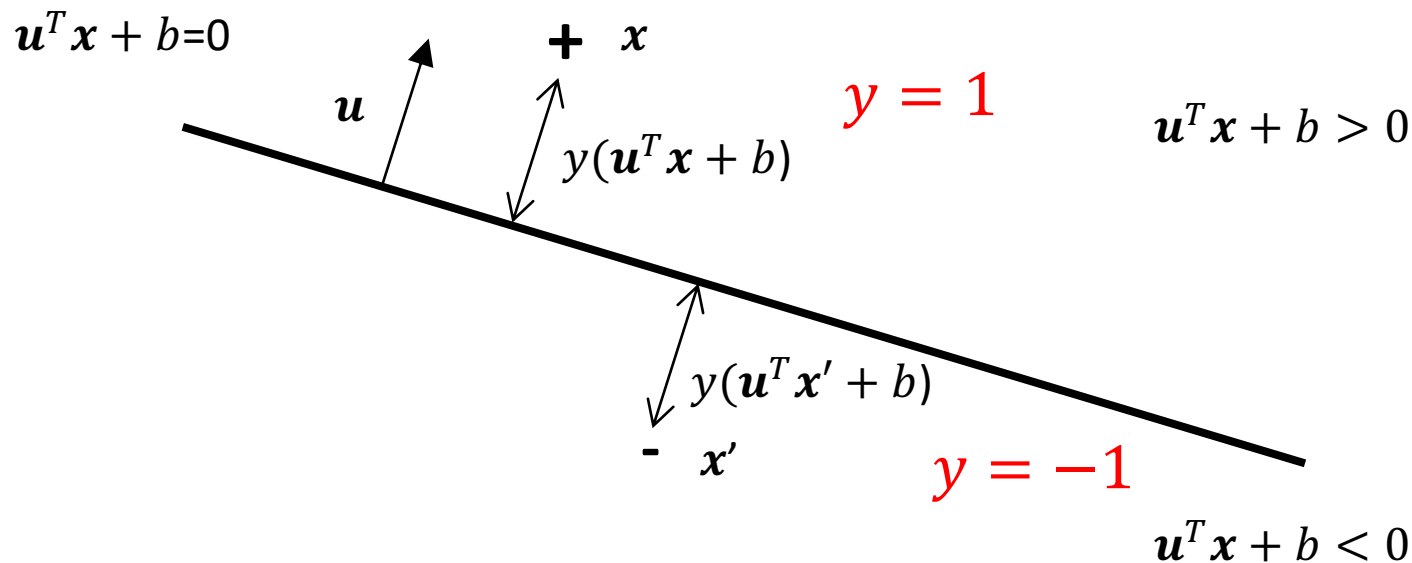
$$\text{Prediction} = \text{sgn}(\mathbf{u}^T \mathbf{x} + b)$$



If we use a **unit** normal vector \mathbf{u} represents the hyperplane, the distance between point \mathbf{x} to plane is $|\mathbf{u}^T \mathbf{x} + b|$ or $y(\mathbf{u}^T \mathbf{x} + b)$

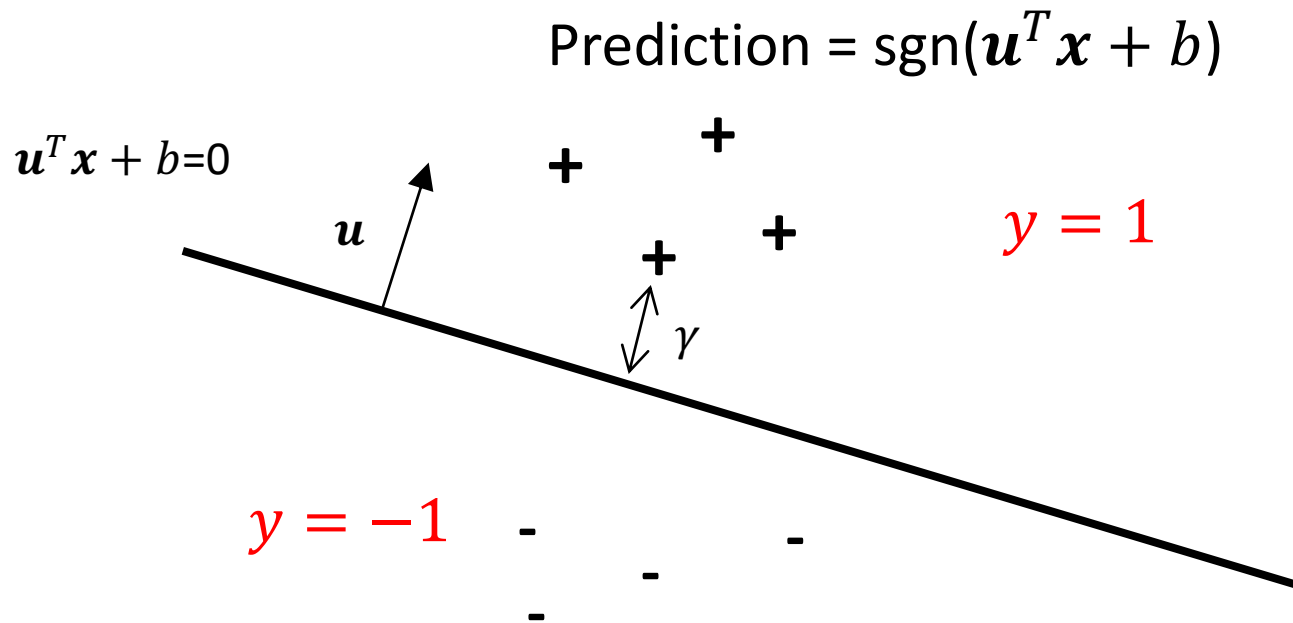
Recall: The geometry of a linear classifier

$$\text{Prediction} = \text{sgn}(\mathbf{u}^T \mathbf{x} + b)$$



If we use a **unit** normal vector \mathbf{u} represents the hyperplane, the distance between point \mathbf{x} to plane is $|\mathbf{u}^T \mathbf{x} + b|$ or $y(\mathbf{u}^T \mathbf{x} + b)$

Recall: The geometry of a linear classifier



If we use a **unit** normal vector \mathbf{u} represents the hyperplane, the distance between point \mathbf{x} to plane is $|\mathbf{u}^T \mathbf{x} + b|$ or $y(\mathbf{u}^T \mathbf{x} + b)$

If the distance between the closest point in dataset D to the plane \mathbf{u} is γ

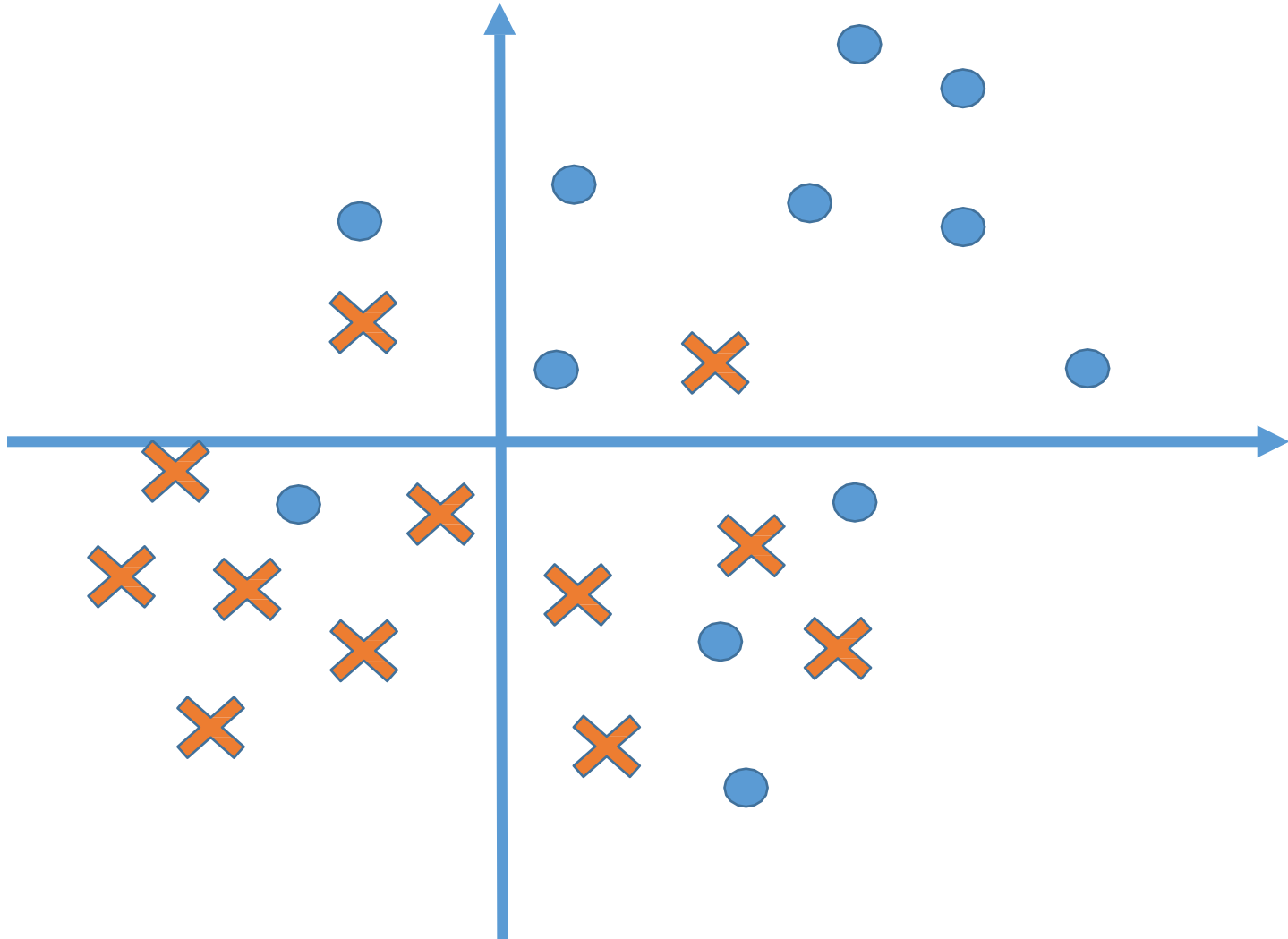
$$y_i(\mathbf{u}^T \mathbf{x}_i + b) \geq \gamma, \forall (x_i, y_i) \in D$$

Logistic regression

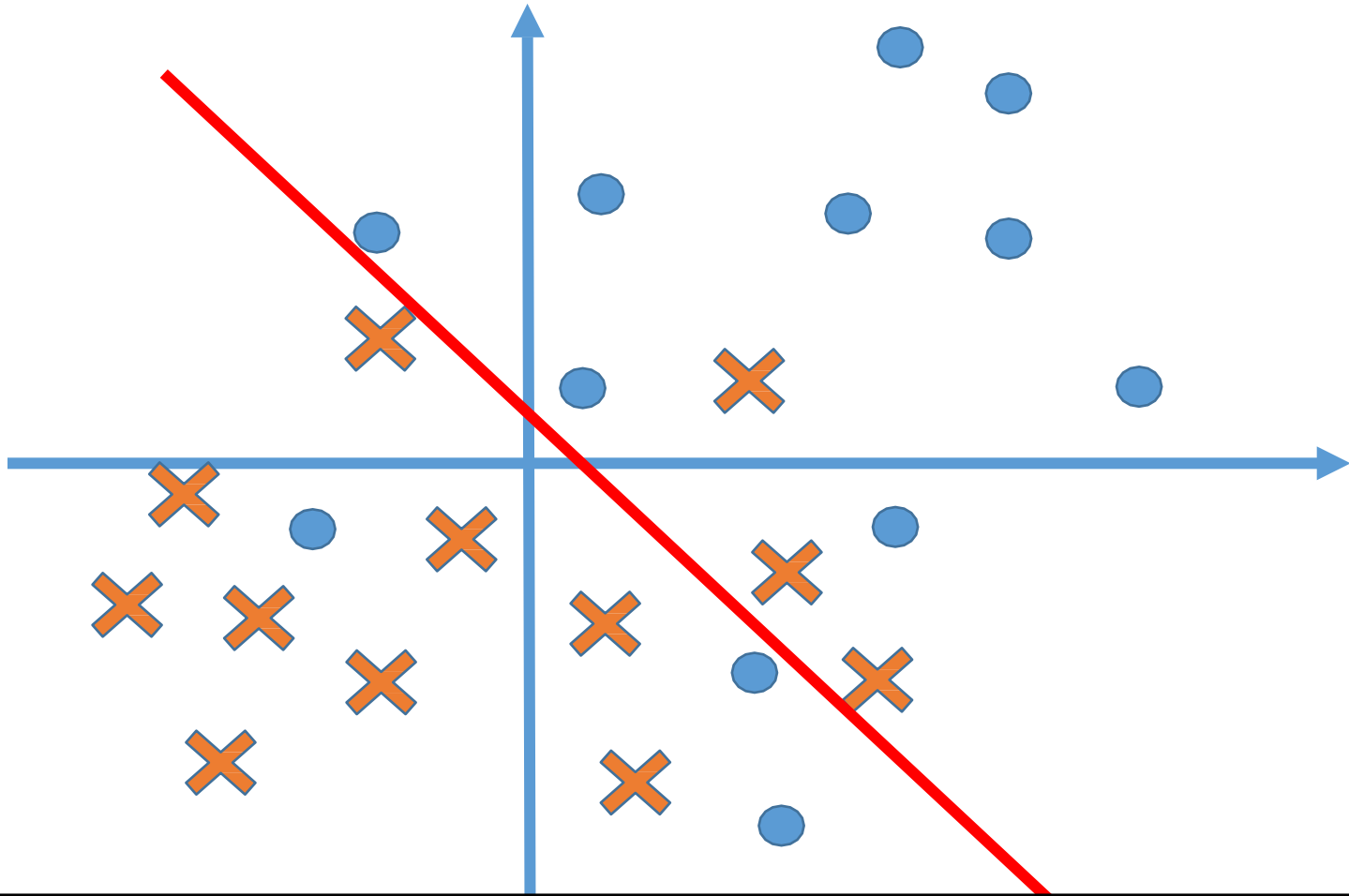
What you will learn today

- ❖ Logistic regression assumption
- ❖ Sigmoid function
- ❖ Maximum likelihood principle
- ❖ Optimization in ML
 - ❖ Stochastic gradient decent

What if data is not linearly separable?

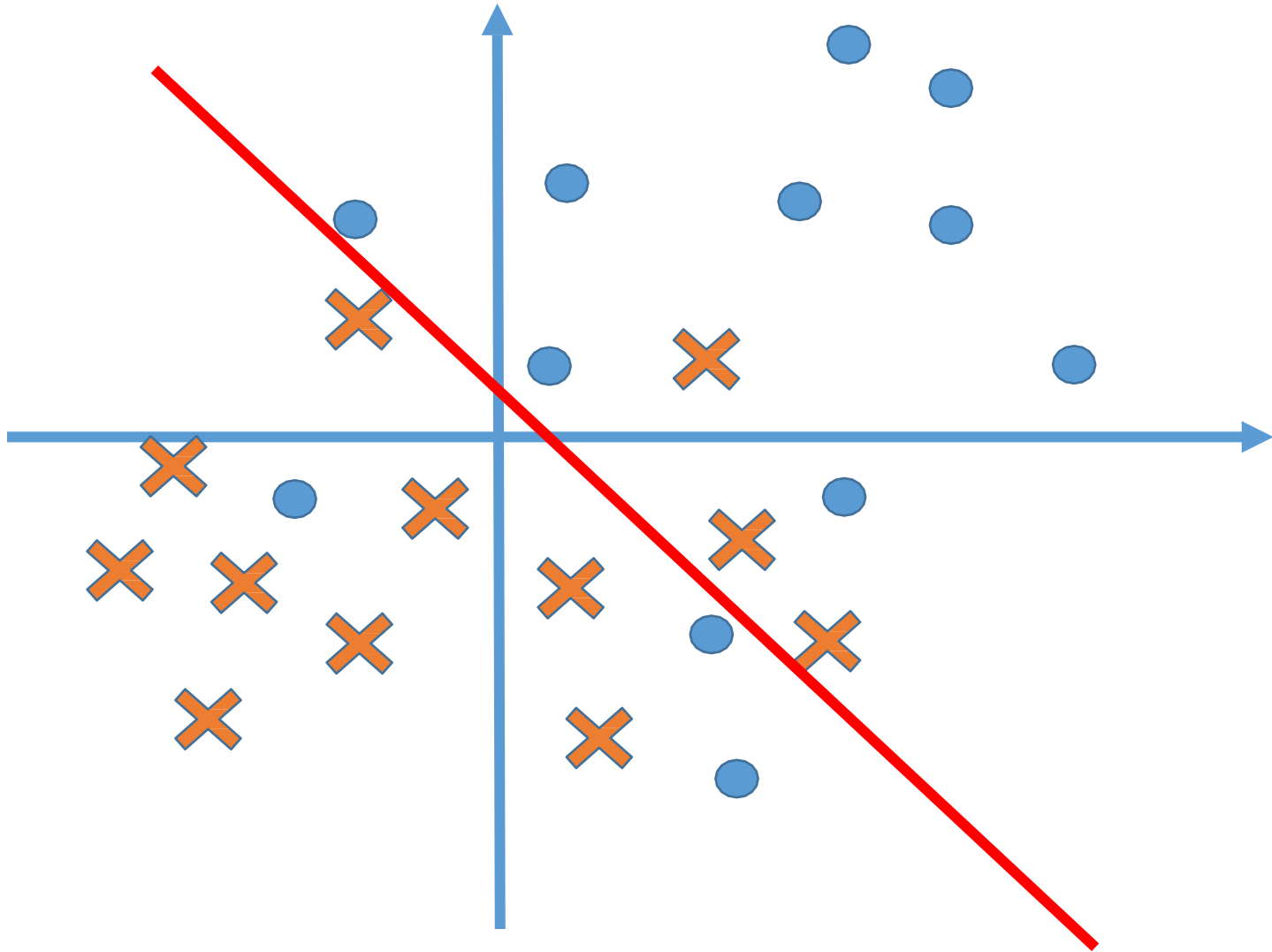


What if data is not linearly separable?



There is no linear model can separate all the data well

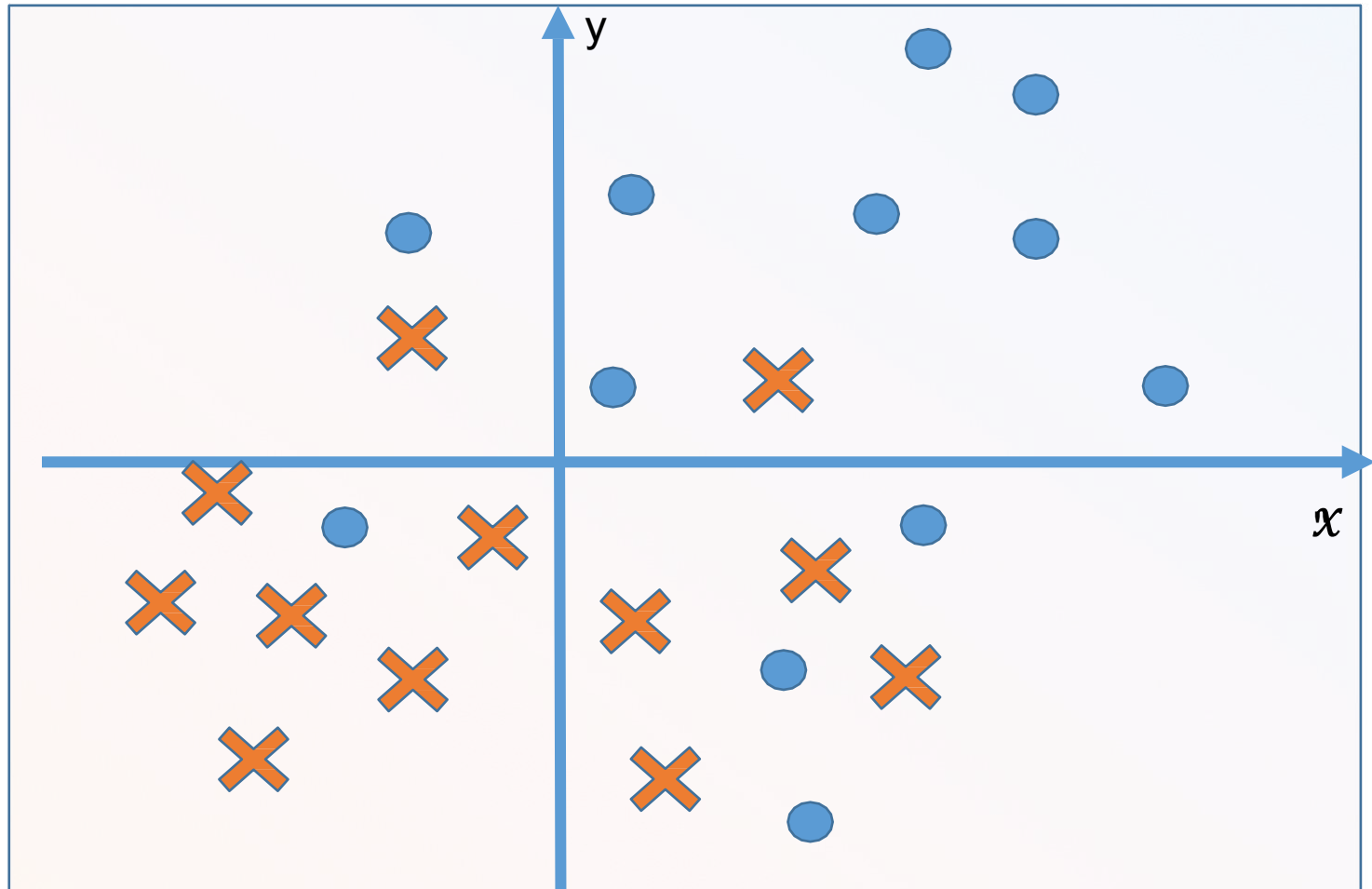
What Making data not linearly separable? [Discussion]



What Making data not linearly separable?

- ❖ Decision boundary is nonlinear
 - ❖ E.g., XOR example
- ❖ Noise in the training data
 - ❖ Outlier due to annotation errors
- ❖ Not enough features
- ❖ The natural of the prediction task
 - ❖ Patients with the same lab test results may not always have the same disease

What if data is not linearly separable?

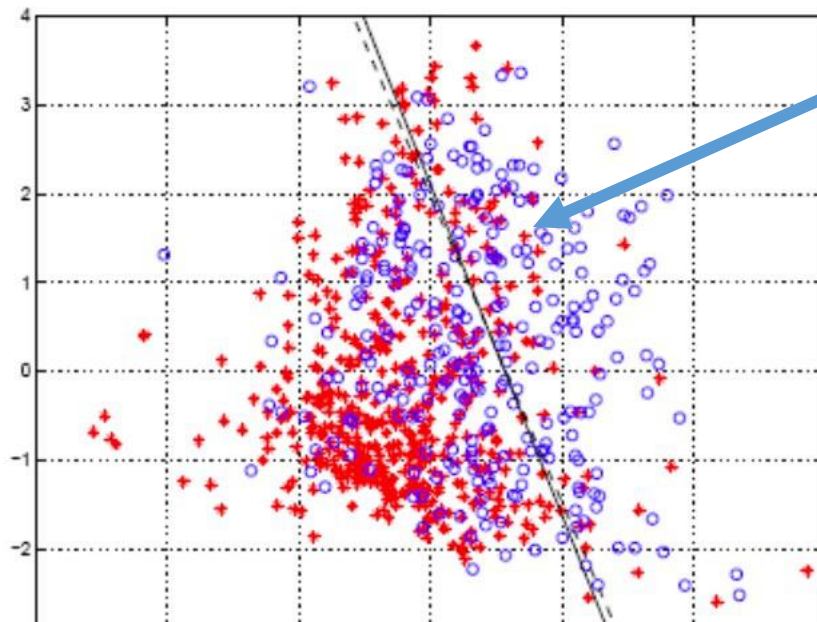


We consider a probabilistic model (today's lecture)

Modeling the Probability

Classification, but...

- ❖ The output y is a discrete value
- ❖ Instead of predicting the output label, let's predict $P(y = 1 \mid \mathbf{x})$

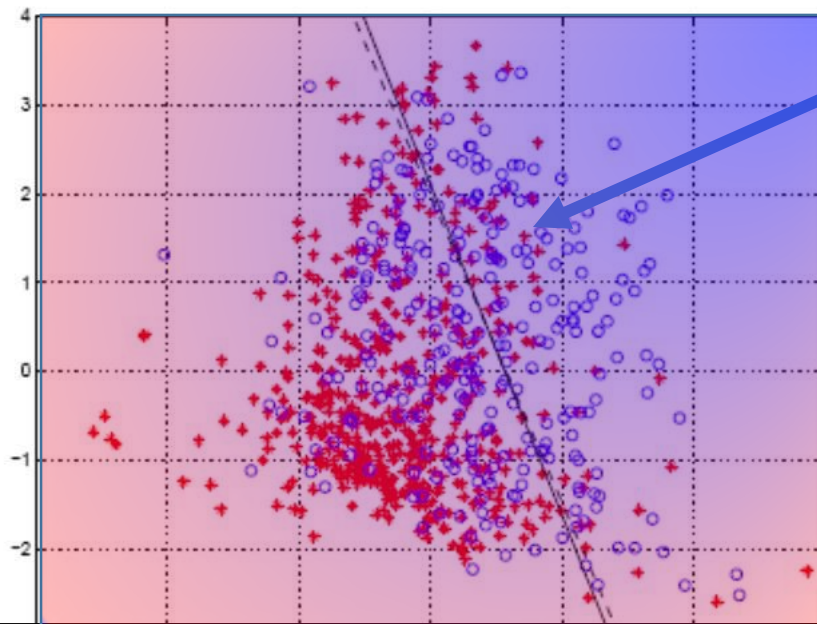


How likely the label $y = 1$ if my feature vector is in this region

Perceptron does not produce probability estimates

Classification, but...

- ❖ The output y is a discrete value
- ❖ Instead of predicting the output label, let us predict $P(y = 1 \mid \mathbf{x})$



How likely the label $y = 1$ if my feature vector is in this region

Perceptron does not produce probability estimates

Predict $P(y = 1 \mid \mathbf{x})$

Input: $x \in \mathbb{R}^d$

Output: $y \in \{1, -1\}$

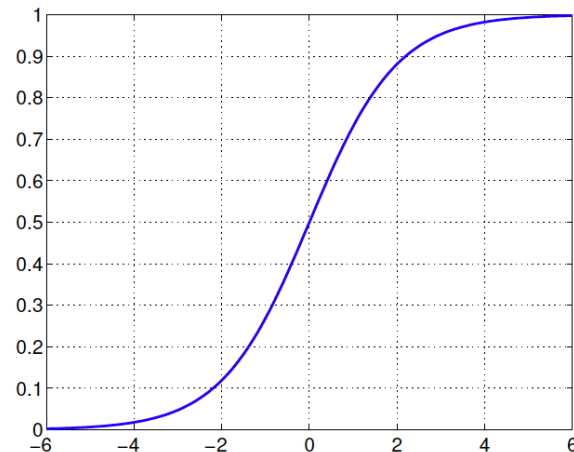
Build a model $h(x)$ such that

$$h(x) = \sigma(w^T x + b) \approx P(y = 1|x)$$

a regression problem

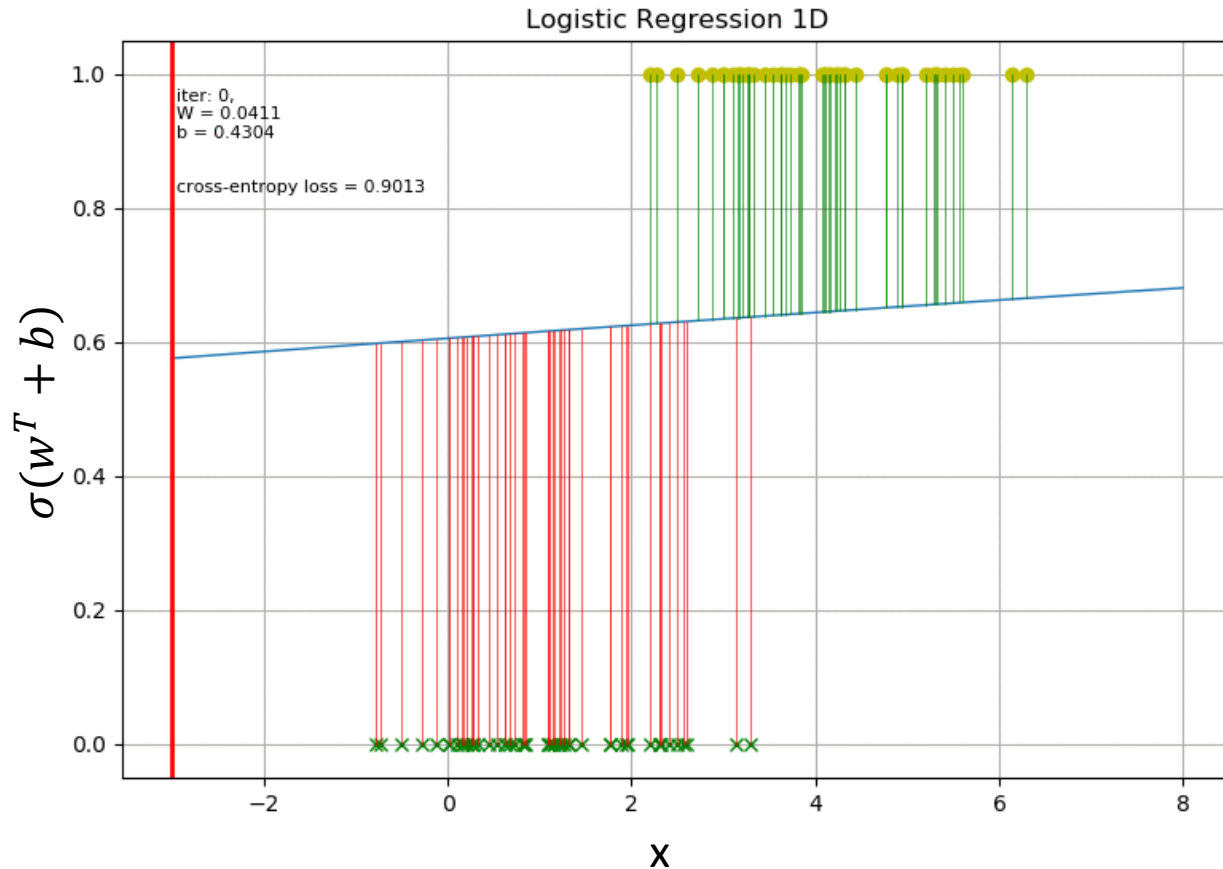
σ is a sigmoid function

$$\begin{aligned}\sigma(z) &= \frac{\exp(z)}{1 + \exp(z)} \\ &= \frac{1}{1 + \exp(-z)}\end{aligned}$$



Logistic Regression in Action

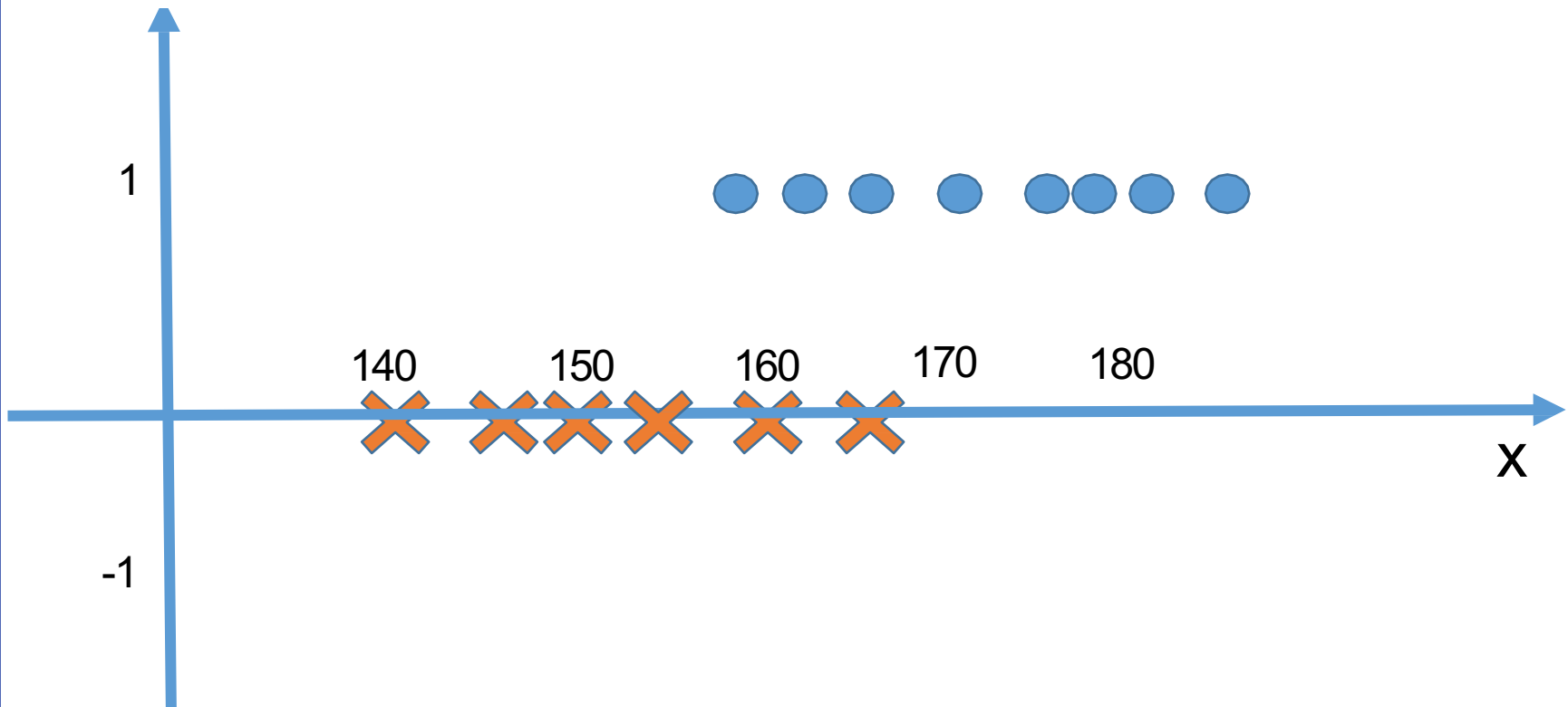
$$\sigma(w^T x + b) \approx P(y = 1|x)$$



<https://medium.com/swlh/from-animation-to-intuition-linear-regression-and-logistic-regression-f641a31e1caf>

Why sigmoid function?

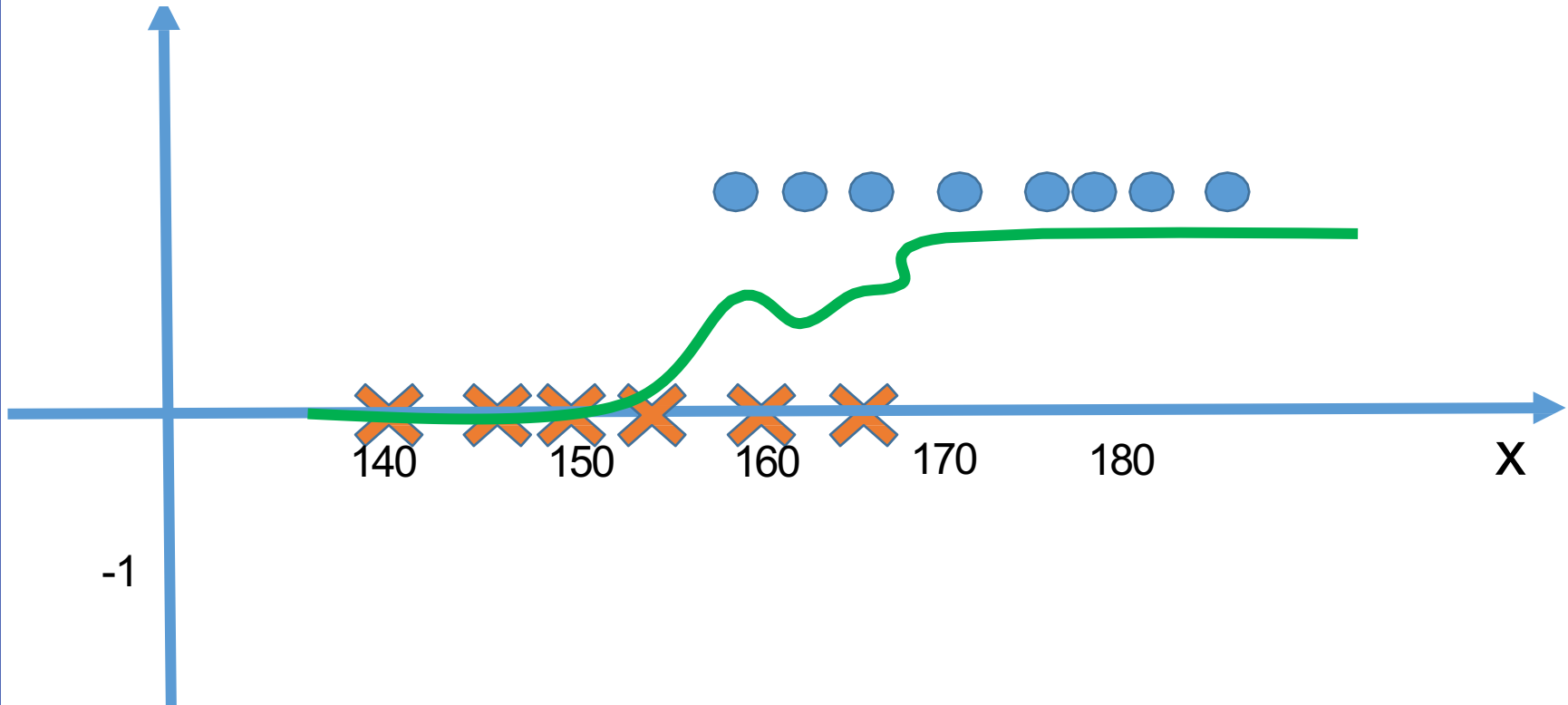
$$P(y = 1|x)$$



What is the σ function

What is the underlying target function?

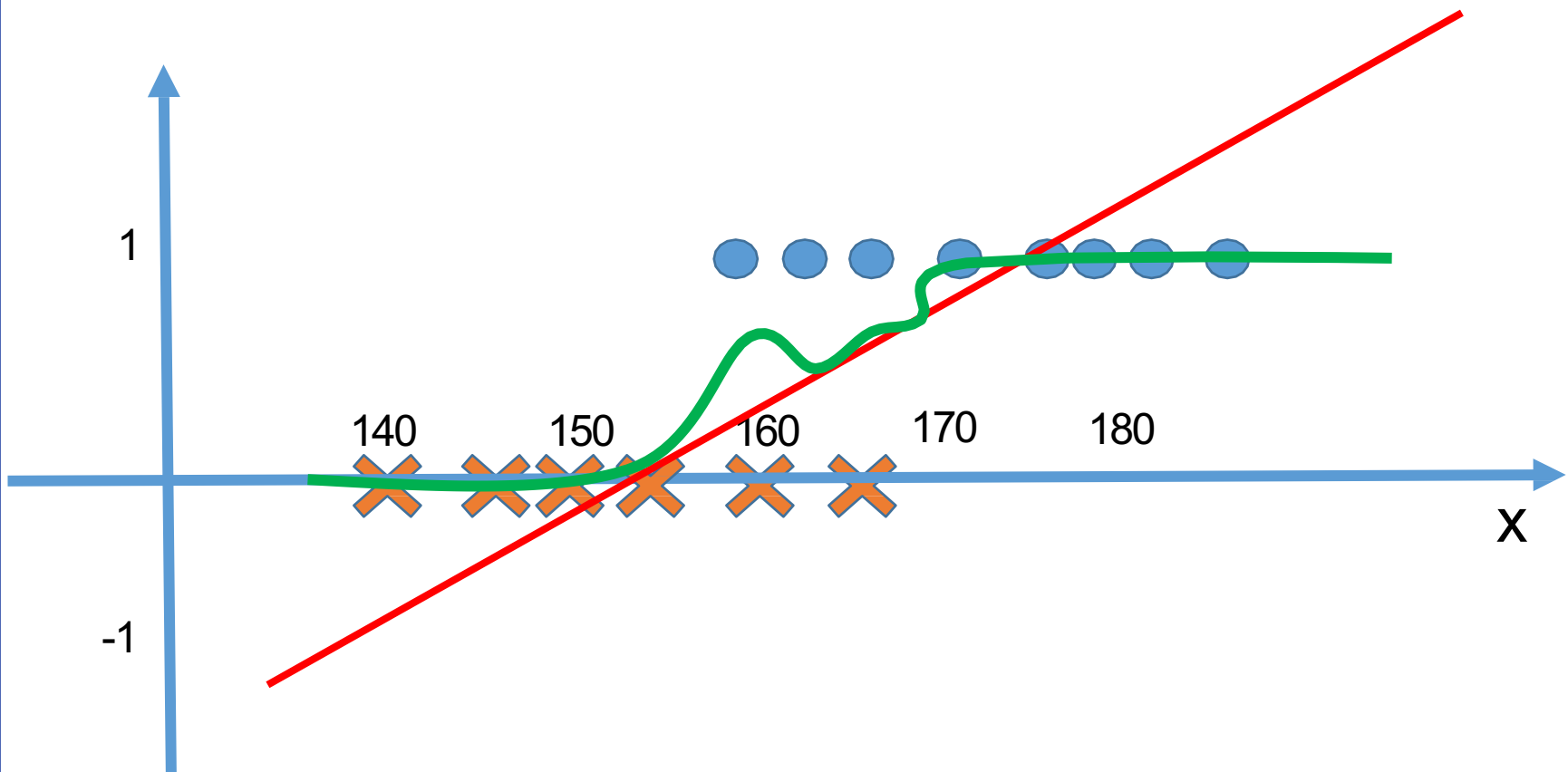
$$P(y = 1|x)$$



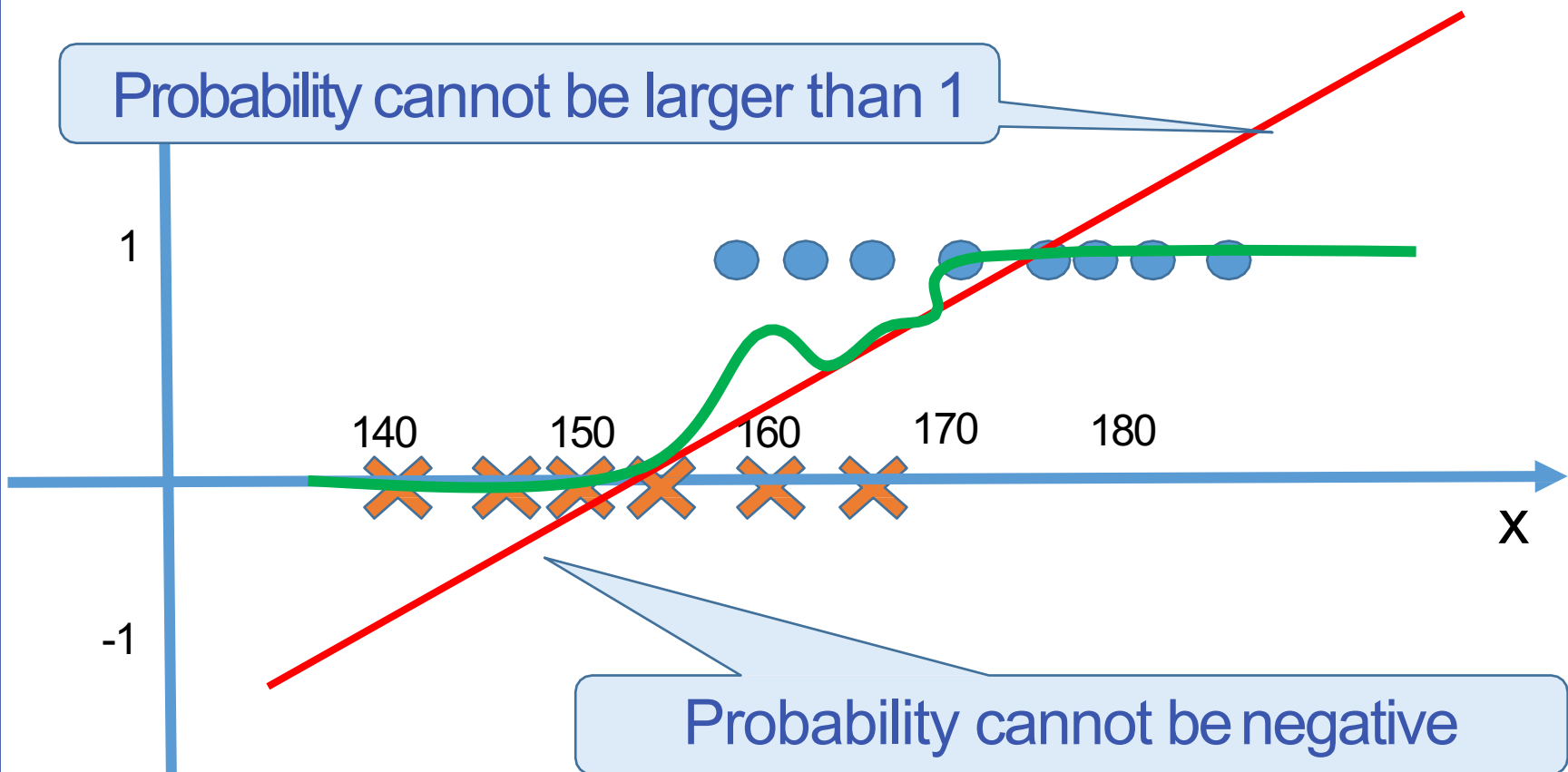
How to fit $P(y = 1|x)$

Can we fit it with a linear function?

$$y = \mathbf{w}^T \mathbf{x} + b$$



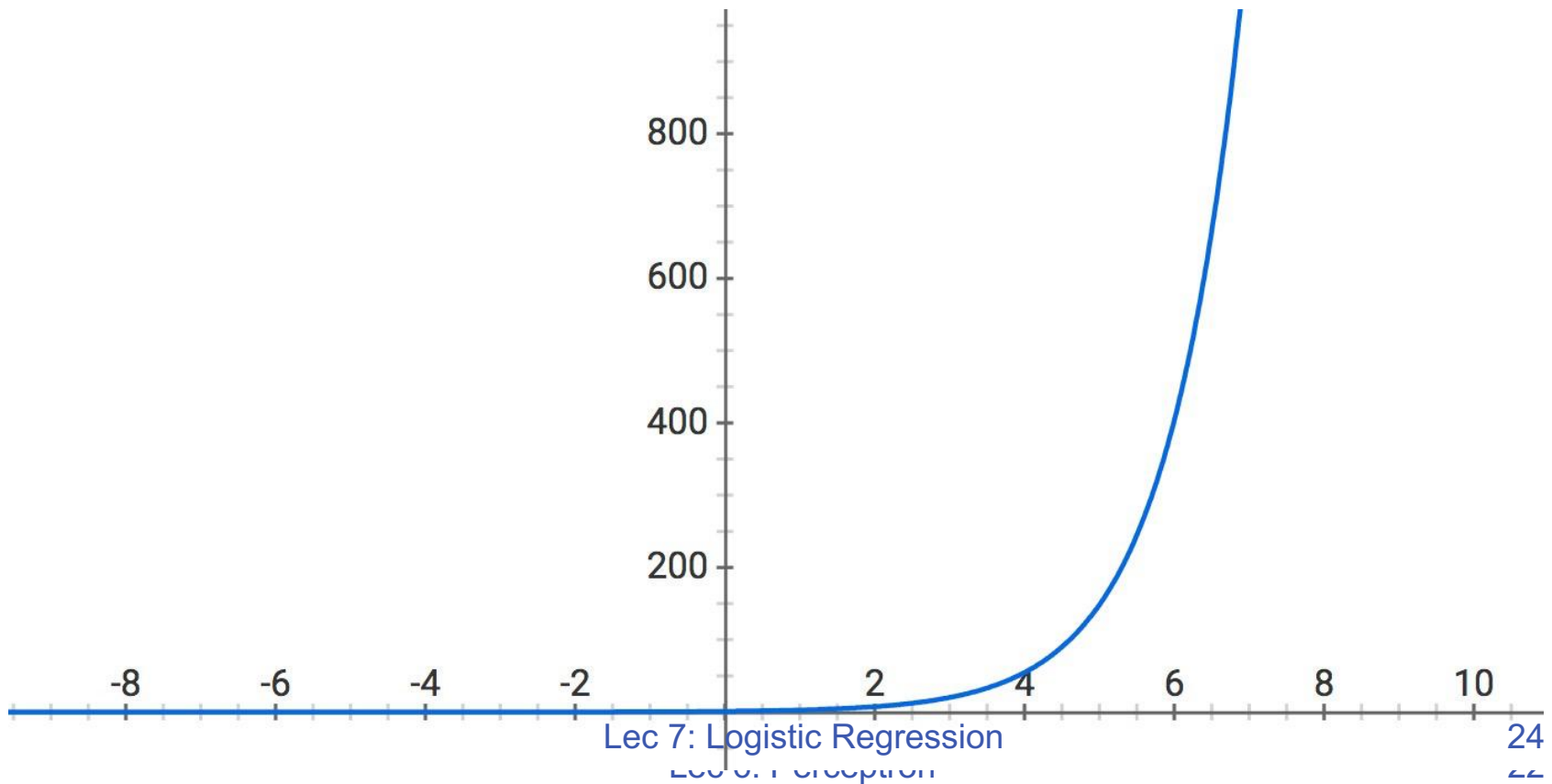
How to fit $P(y = 1|x)$



How can we design such a transformation function?

Idea 1: function always output positive value

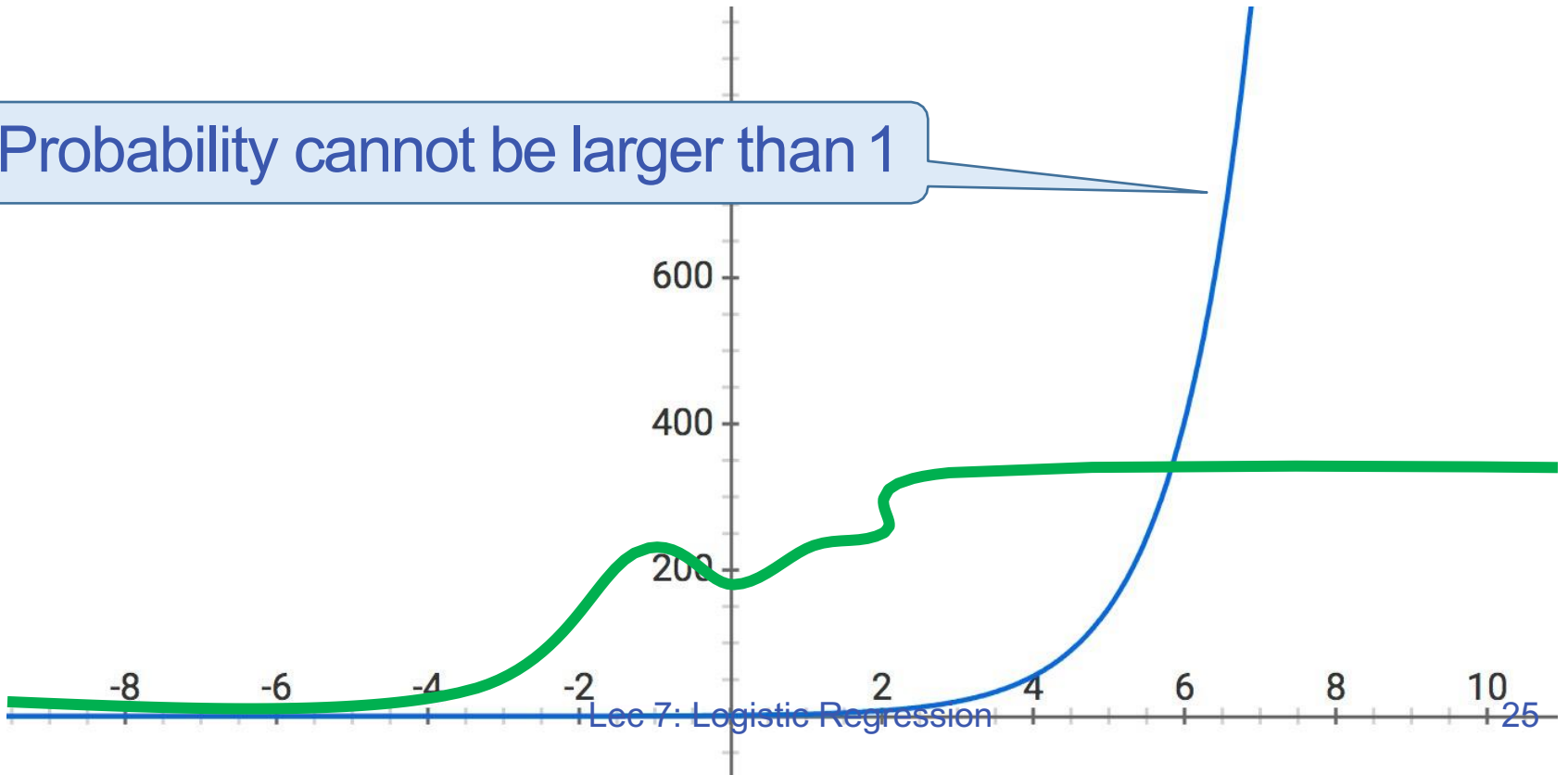
$\exp(\cdot)$ is always positive



How can we design such a transformation function?

- ❖ Idea 1: function always output positive value
 - ❖ $\exp(\cdot)$ is always positive
 - ❖ $\exp(w^T x + b)$ always return positive value

Probability cannot be larger than 1



How can we design such a transformation function?

Idea 2: normalize the value such that it is less than 1

❖ $\exp(w^T x + b) \in (0, \infty)$ grows very fast, so we need to use exp to normalize itself

❖ Let's use

$$\sigma(w^T x + b) = \frac{\exp(w^T x + b)}{1 + \exp(w^T x + b)}$$

❖ When $w^T x + b \rightarrow \infty$,
 $\sigma(w^T x + b) \rightarrow 1$

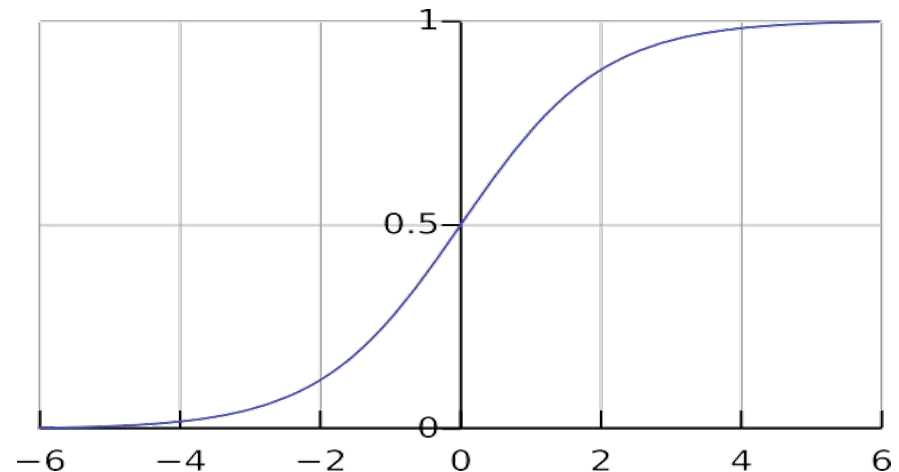
❖ When $w^T x + b \rightarrow -\infty$,
 $\sigma(w^T x + b) \rightarrow 0$



The Sigmoid function

The $\sigma(z)$ function is called sigmoid function (or logistic function)

$$\begin{aligned}\sigma(z) &= \frac{\exp(z)}{1 + \exp(z)} \\ &= \frac{1}{1 + \exp(-z)}\end{aligned}$$



Summary (Modeling)

❖ What is the goal of logistic regression?

❖ Model $P(y = 1|x)$

❖ What is the hypothesis space?

$$H = \{ h \mid h : X \rightarrow P(Y \mid X), h(x) = \sigma(w^T x + b) \}$$
$$\sigma(z) = \frac{1}{1 + \exp(-z)}$$

We want to find $h(x)$ such that

$$h(x) \approx P(y = 1 | x)$$

Decision Boundary of Logistic Regression

Predicting a label

$$P(y = 1|\mathbf{x}; \mathbf{w}) = \sigma(\mathbf{w}^T \mathbf{x}) = \frac{1}{1 + \exp(-\mathbf{w}^T \mathbf{x})}$$

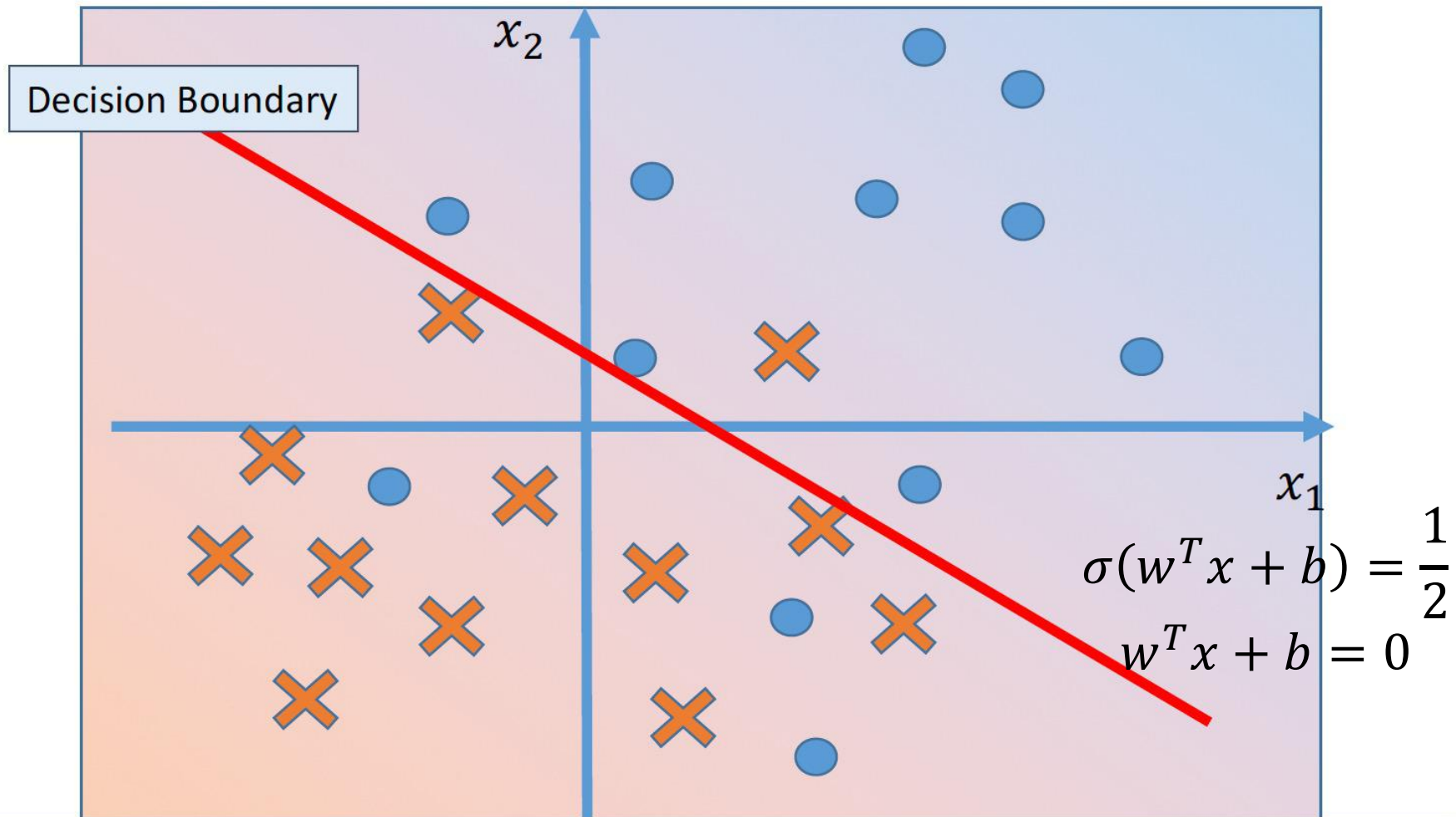
Compute $\sigma(w^T x)$;

If this is greater than half, predict 1
else predict -1

What does this correspond to in terms of $\mathbf{w}^T \mathbf{x}$?

$$w^T x = 0$$

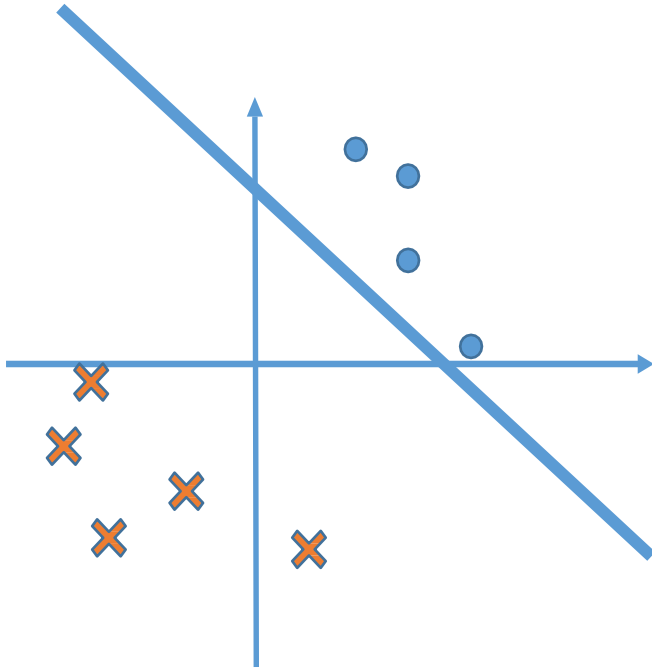
Prediction by logistic regression



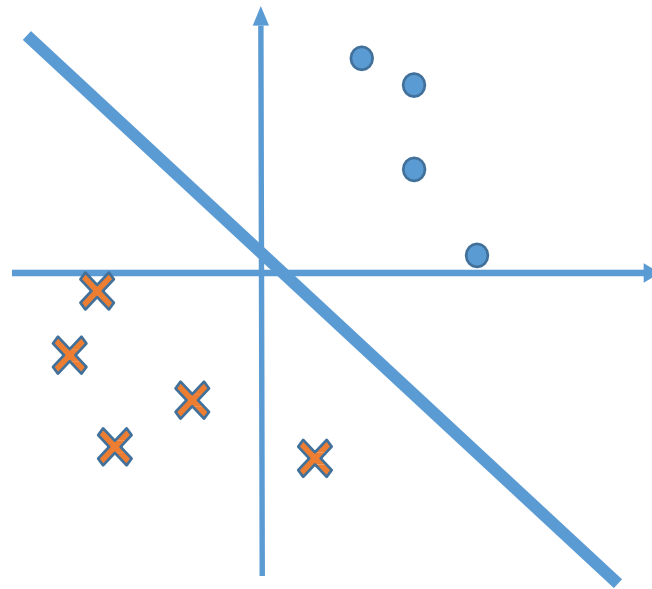
A probabilistic model of modeling $\sigma(w^T x + b) = P(y = 1|x)$

Exercise

- ❖ Which function(s) are likely to be a decision function of logistic regression



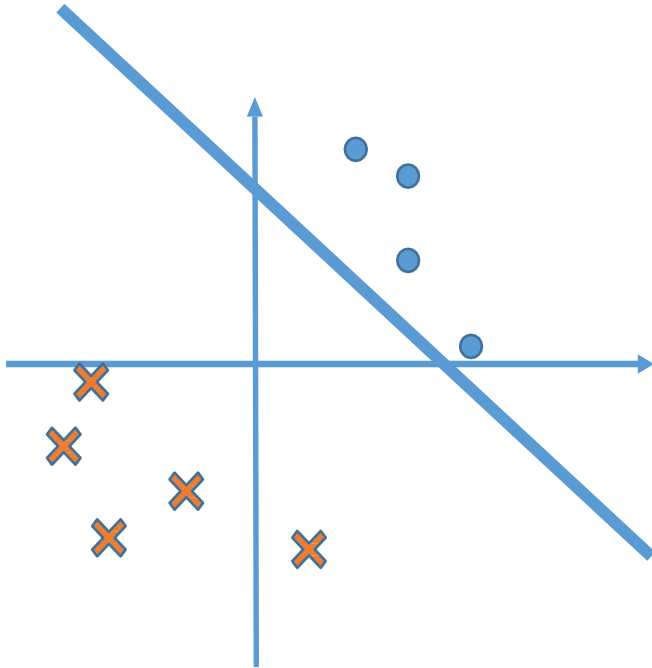
(A)



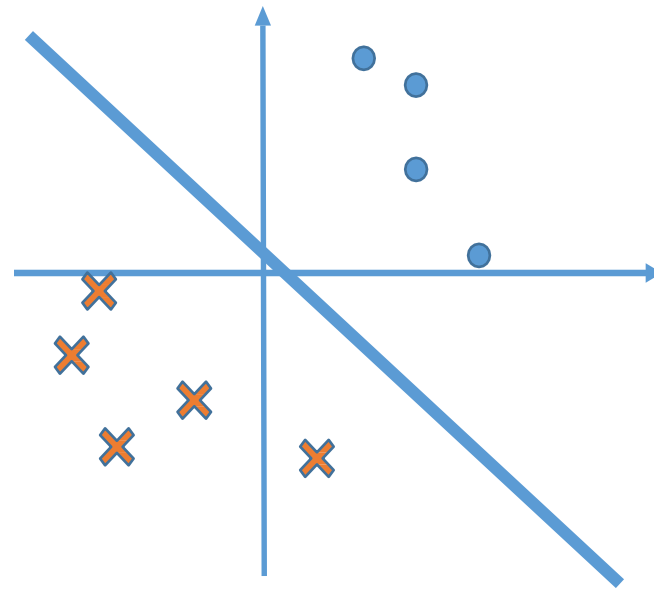
(B)

Exercise

- ❖ Which function(s) are likely to be a decision function of Perceptron



(A)



(B)

How to Train a Logistic Regression Model?

Logistic Regression: Setup

❖ The setting

- ❖ Binary classification

- ❖ Inputs: Feature vectors $x \in R^N$

- ❖ Labels: $y \in \{-1, +1\}$

❖ Training data

- ❖ $S = \{(\mathbf{x}_i, y_i)\}$, m examples

❖ Hypothesis space

$$H = \{ h \mid h : X \rightarrow P(Y \mid X), h(x) = \sigma (w^T x + b) \}$$
$$\sigma(z) = \frac{1}{1 + \exp(-z)}$$

Logistic Regression: Setup

❖ Training data

❖ $S = \{(\mathbf{x}_i, y_i)\}$, m examples

❖ Hypothesis space

$$H = \{ h \mid h : X \rightarrow Y, h(x) = \sigma (w^T x + b) \}$$

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

❖ Learning Goal:

❖ Find an $h \in H$, such that $h(x) \approx P(y = 1 \mid x)$

❖ How to?

Maximum Likelihood

Which bag of words more likely generate:

aDaaa



Maximum Likelihood

Which bag of words more likely generate:

aDaaa

$$0.7 \times 0.1 \times 0.7 \times 0.7 \times 0.7 \\ = 2.401 \times 10^{-2}$$



$$0.2 \times 0.1 \times 0.2 \times 0.2 \times 0.2 \\ = 1.6 \times 10^{-4}$$



Drawing color cards from the envelope

- ❖ Let say we have several cards in the envelope
- ❖ Assume



$$P(card = yellow) = \theta$$



$$P(card = purple) = 1 - \theta$$

Drawing color cards from the envelope

- ❖ Sample with replacement n times
- ❖ k times we get yellow card, and $n-k$ times, we get purple card
- ❖ The joint probability (likelihood)
$$C_K^N \theta^k (1 - \theta)^{n-k}$$
- ❖ What is the best θ making the joint probability maximal?

Drawing color cards from the envelope

$$\cancel{C_K^N} \theta^k (1 - \theta)^{n-k}$$

❖ Solving $\max_{\theta} \theta^k (1 - \theta)^{n-k}$

❖ Equivalently, we can solve

$$\max_{\theta} \log(\theta^k (1 - \theta)^{n-k})$$

$$\max_{\theta} k \log \theta + (n - k) \log(1 - \theta)$$

❖ At the optimum,

$$\frac{d(k \log \theta + (n - k) \log(1 - \theta))}{d\theta} = 0$$

The usual trick: Convert products to sums by taking log

Recall that this works only because log is an increasing function and the maximizer will not change

Drawing color cards from the envelope

$$\frac{d(k \log \theta + (n-k) \log(1-\theta))}{d\theta} = 0$$

$$\Rightarrow \frac{k}{\theta} - \frac{n-k}{1-\theta} = 0$$

For this simple problem, we have a closed-form solution.

We are not always lucky like this

$$\Rightarrow \theta(n-k) = (1-\theta)k$$

$$\Rightarrow \theta = \frac{k}{n}$$



Maximum Likelihood Estimator (formal definition)

Likelihood function of parameters

Let X_1, \dots, X_N be **IID** (independent and identically distributed) with PDF $p(x|\theta)$ (also written as $p(x; \theta)$). The *likelihood function* is defined by $L(\theta)$,

$$L(\theta) = p(X_1, \dots, X_N; \theta). \qquad = \prod_{i=1}^N p(X_i; \theta).$$

Notes The likelihood function is just the joint density of the data, except that we treat it as a function of the parameter θ .

Maximum Likelihood Estimation

Maximum Likelihood Estimator

Definition: The maximum likelihood estimator (MLE) $\hat{\theta}$, is the value of θ that maximizes $L(\theta)$.

The log-likelihood function is defined by $l(\theta) = \log L(\theta)$. Its maximum occurs at the same place as that of the likelihood function.

Back to Logistic regression

Training data

$S = \{(\mathbf{x}_i, y_i)\}$, m examples

Hypothesis space

$$H = \{ h \mid h : X \rightarrow Y, h(x) = \sigma(w^T x + b) \}$$
$$\sigma(z) = \frac{1}{1 + \exp(-z)}$$

Find an $h \in H$, such that $h(x) \approx P(y = 1 \mid x)$

Back to Logistic regression

Training data

$S = \{(\mathbf{x}_i, y_i)\}$, m examples

Hypothesis space

$$H = \{ h \mid h : X \rightarrow Y, h(x) = \sigma(w^T x + b) \}$$
$$\sigma(z) = \frac{1}{1 + \exp(-z)}$$

Find an $h \in H$, such that $h(x) \approx P(y = 1 \mid x)$

How to? **Maximum Likelihood Estimator**

Likelihood function



$$P(\text{card} = \text{yellow}) = \theta$$



$$P(\text{card} = \text{purple}) = 1 - \theta$$

Find θ by maximizing $L(D; \theta)$

Logistic regression

Find w, b by maximizing $P(S; w, b)$

Likelihood function



$$P(\text{card} = \text{yellow}) = \theta$$



$$P(\text{card} = \text{purple}) = 1 - \theta$$

Find θ by maximizing $L(D; \theta)$

Logistic regression

Find w, b by maximizing $P(S; w, b)$

$$\operatorname{argmax}_{w,b} P(S; w, b) = \operatorname{argmax}_{w,b} \prod_{i=1}^m P(y_i | x_i; w, b)$$

Maximum likelihood estimator for logistic regression

$$\operatorname{argmax}_{w,b} P(S; w, b) = \operatorname{argmax}_{w,b} \prod_{i=1}^m P(y_i | x_i; w, b)$$

Equivalent to solve

$$\operatorname{argmax}_{w,b} \sum_{i=1}^m \log P(y_i | x_i; w, b)$$

Remember our assumption:

$$P(y = 1 | x; w, b) = \sigma(w^T x + b) = \frac{1}{1 + \exp(-(w^T x + b))}$$

$$P(y = -1 | x; w, b) = 1 - \sigma(w^T x + b) = 1 - \frac{1}{1 + \exp(-(w^T x + b))}$$

$$= \frac{\exp(-(w^T x + b))}{1 + \exp(-(w^T x + b))} = \frac{1}{1 + \exp((w^T x + b))} = \sigma(-(w^T x + b))$$

Maximum likelihood estimator for logistic regression

$$\operatorname{argmax}_{w,b} P(S; w, b) = \operatorname{argmax}_{w,b} \prod_{i=1}^m P(y_i | x_i; w, b)$$

Equivalent to solve

$$\operatorname{argmax}_{w,b} \sum_{i=1}^m \log P(y_i | x_i; w, b)$$

Remember our assumption:

$$P(y = 1 | x; w, b) = \sigma(w^T x + b) = \frac{1}{1 + \exp(-(w^T x + b))}$$

$$P(y_i | x_i; w, b) = \begin{cases} \sigma(w^T x_i + b) \\ \sigma(-(w^T x_i + b)) \end{cases} \Rightarrow P(y_i | x_i; w, b) = \sigma(y_i(w^T x_i + b))$$

Maximum likelihood estimator for logistic regression

$$\operatorname{argmax}_{w,b} P(S; w, b) = \operatorname{argmax}_{w,b} \prod_{i=1}^m P(y_i | x_i; w, b)$$

Equivalent to solve

$$\operatorname{argmax}_{w,b} \sum_{i=1}^m \log P(y_i | x_i; w, b)$$

Using $P(y_i | x_i; w, b) = \sigma(y_i(w^T x_i + b))$

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

$$\begin{aligned} & \operatorname{argmax}_{w,b} \sum_{i=1}^m \log \sigma(y_i(w^T x_i + b)) \\ &= - \sum_{i=1}^m \log(1 + \exp(-y_i(w^T x_i + b))) \end{aligned}$$

How to Optimize the Loss?

How to minimizing the loss

- ❖ Optimization methods
 - ❖ Gradient Descent
 - ❖ Stochastic Gradient Descent
 - ❖ Analytic solution
- ❖ ...many other approaches

How to solve it?

$$\operatorname{argmax}_{w,b} - \sum_{i=1}^m \log(1 + \exp(-y_i(w^T x_i + b)))$$

There is no closed-form solution

Max $f(x)$ is equivalent to min $-f(x)$

=> only need to consider minimization problems.

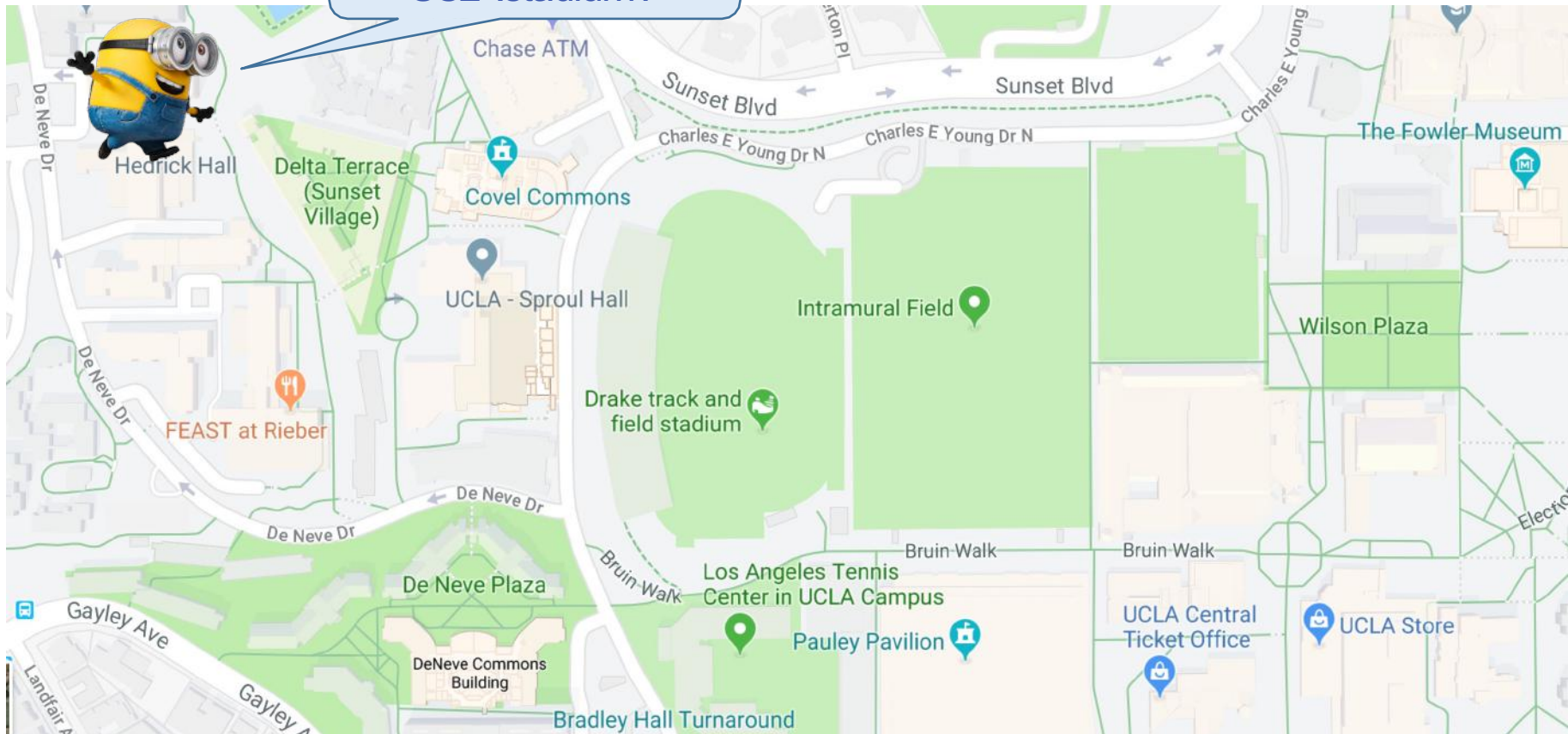
One way to solve it is by **gradient descent**

Gradient Descent

Intuition

Asking direction

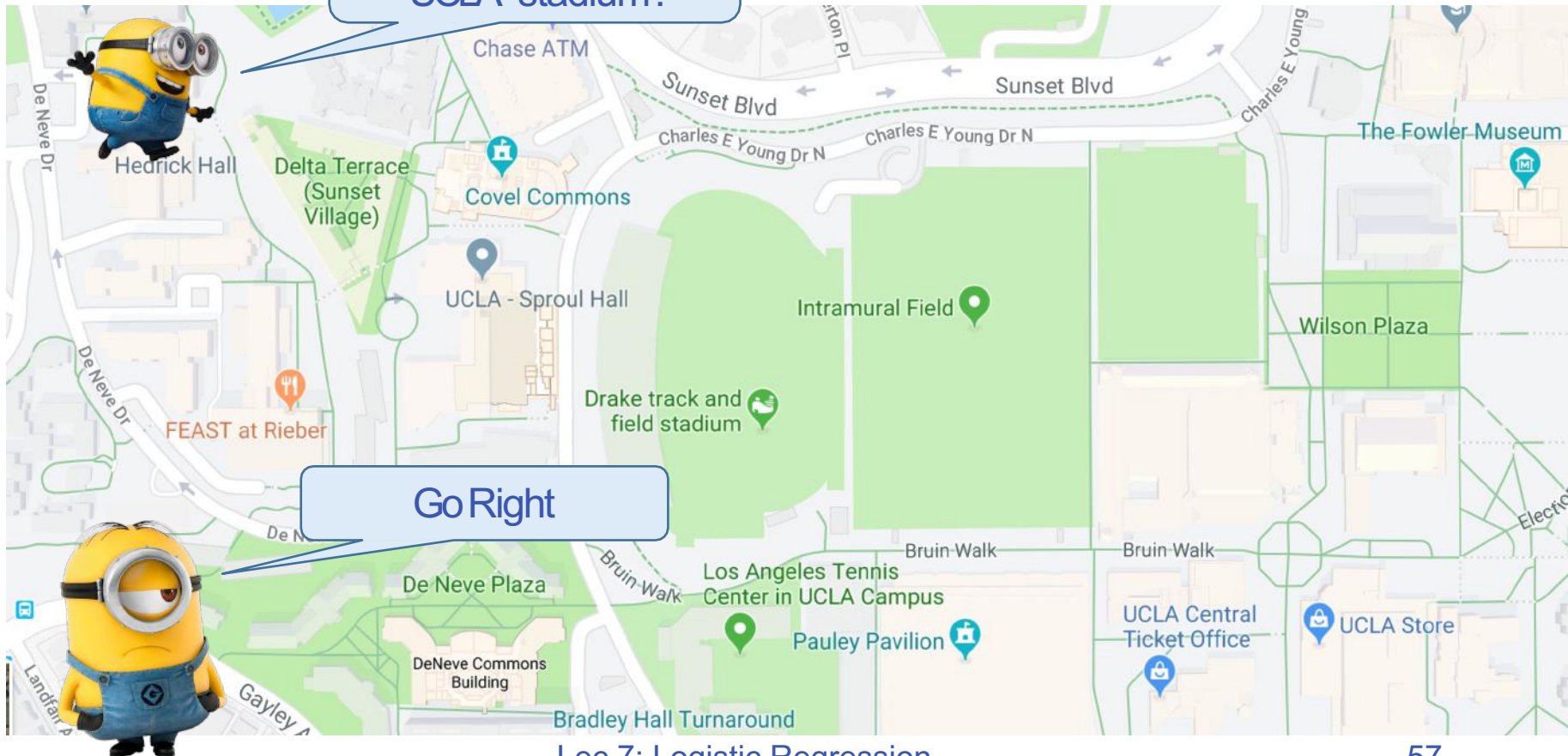
How can I get to
UCLA stadium?



Intuition

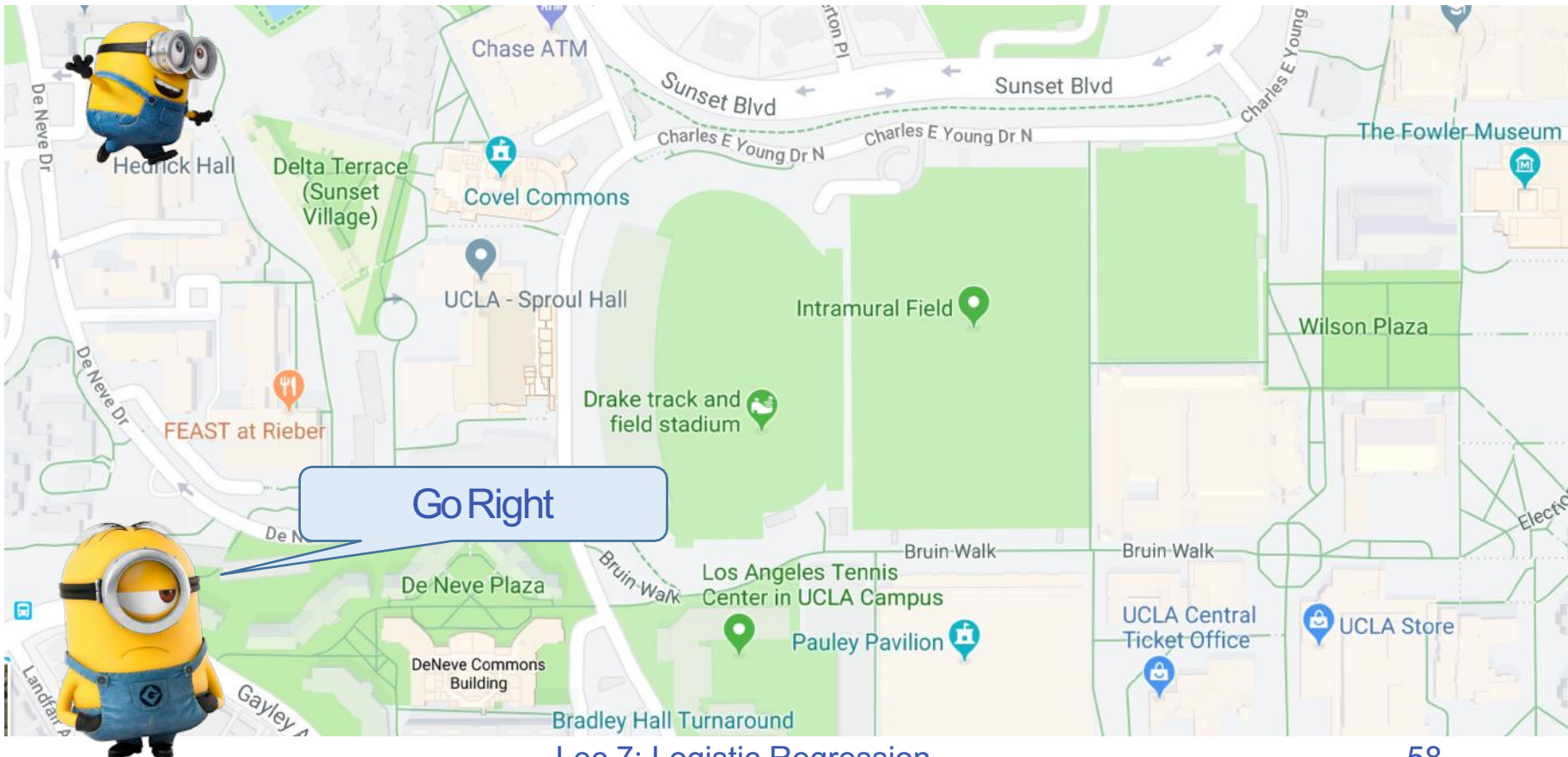
Asking direction

How can I get to
UCLA stadium?



Intuition

Asking direction

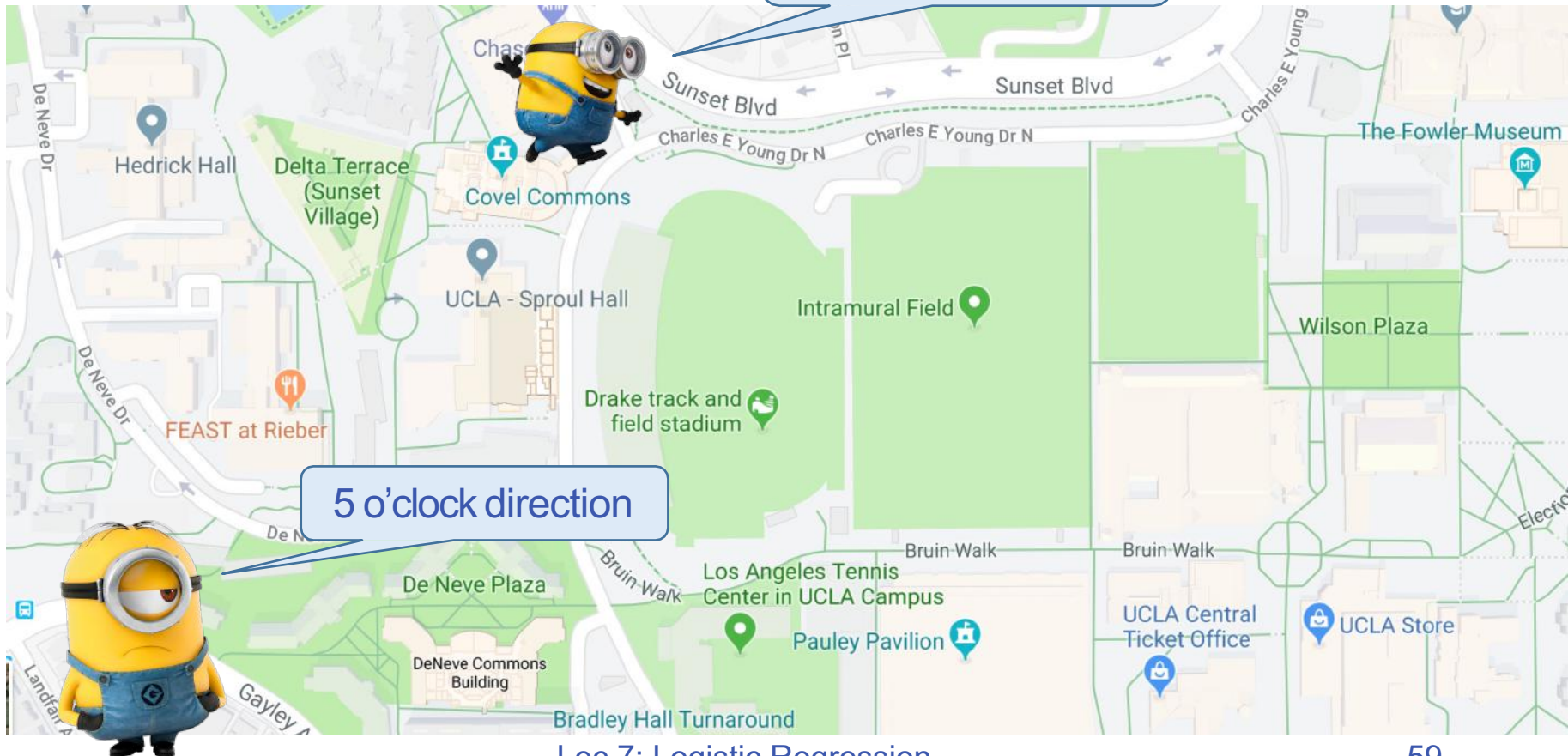


Intuition

Asking direction

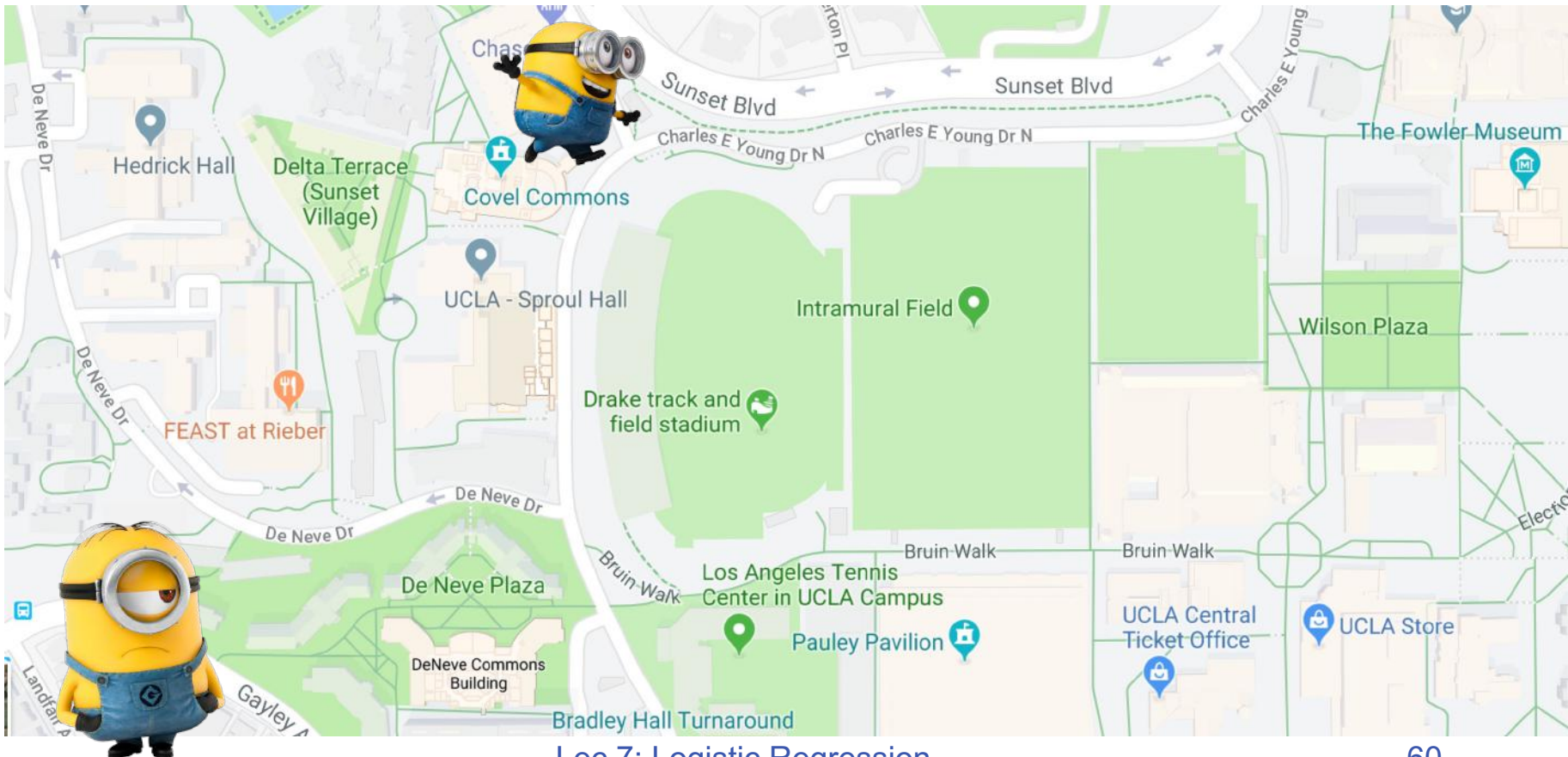
And then?

5 o'clock direction



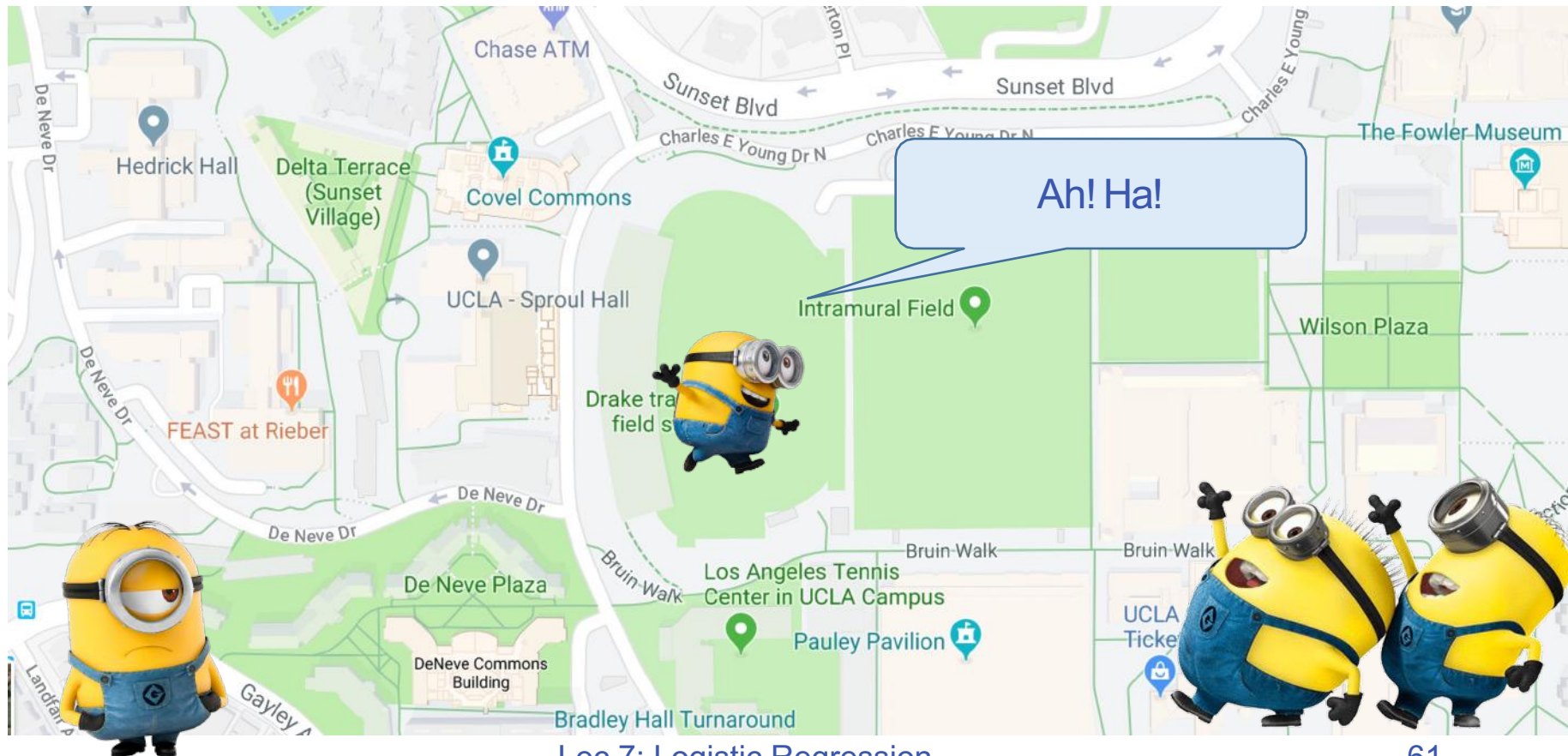
Intuition

Asking direction



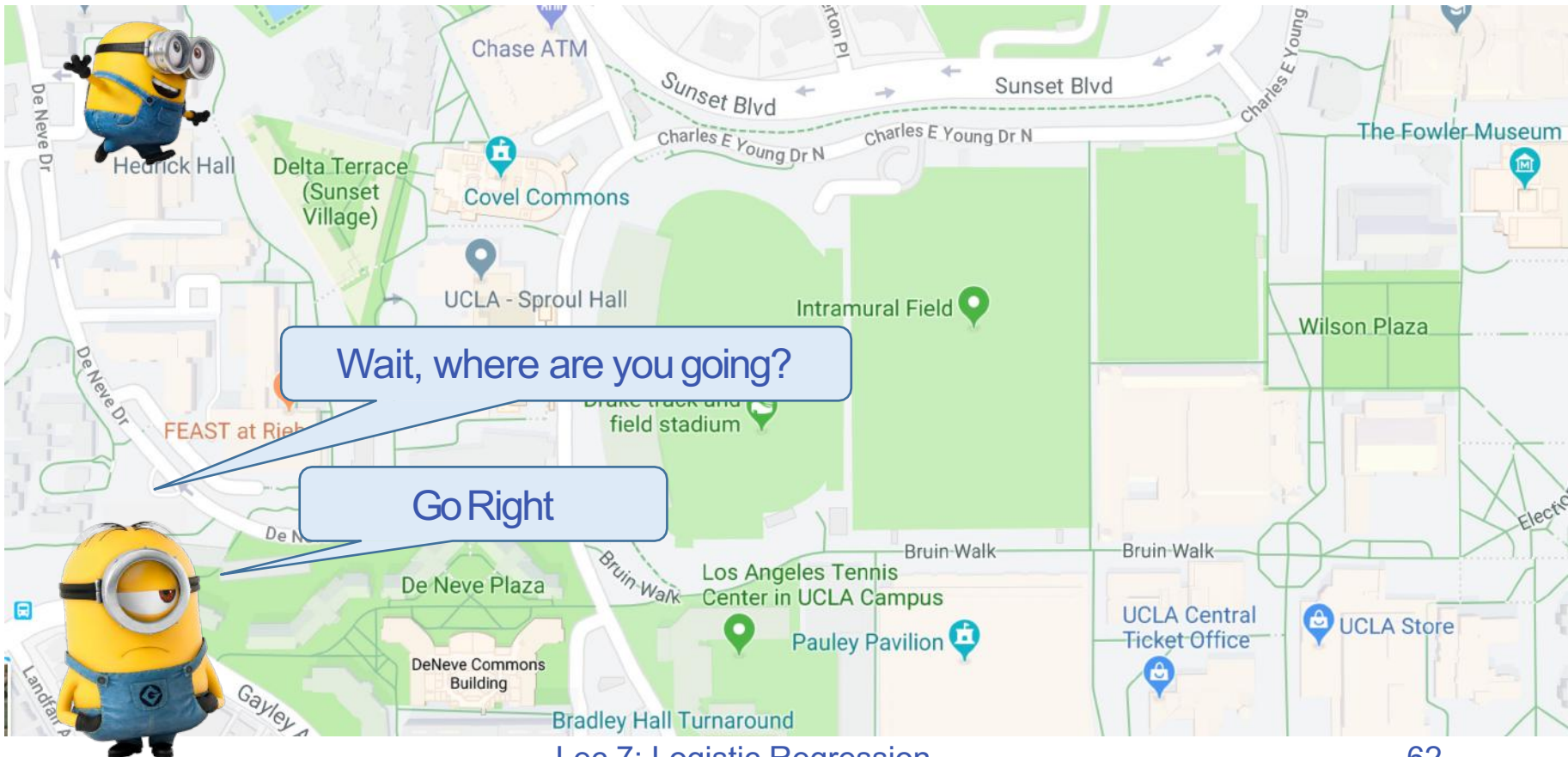
Intuition

Asking direction



Intuition

What may go wrong? Incorrect Step-size

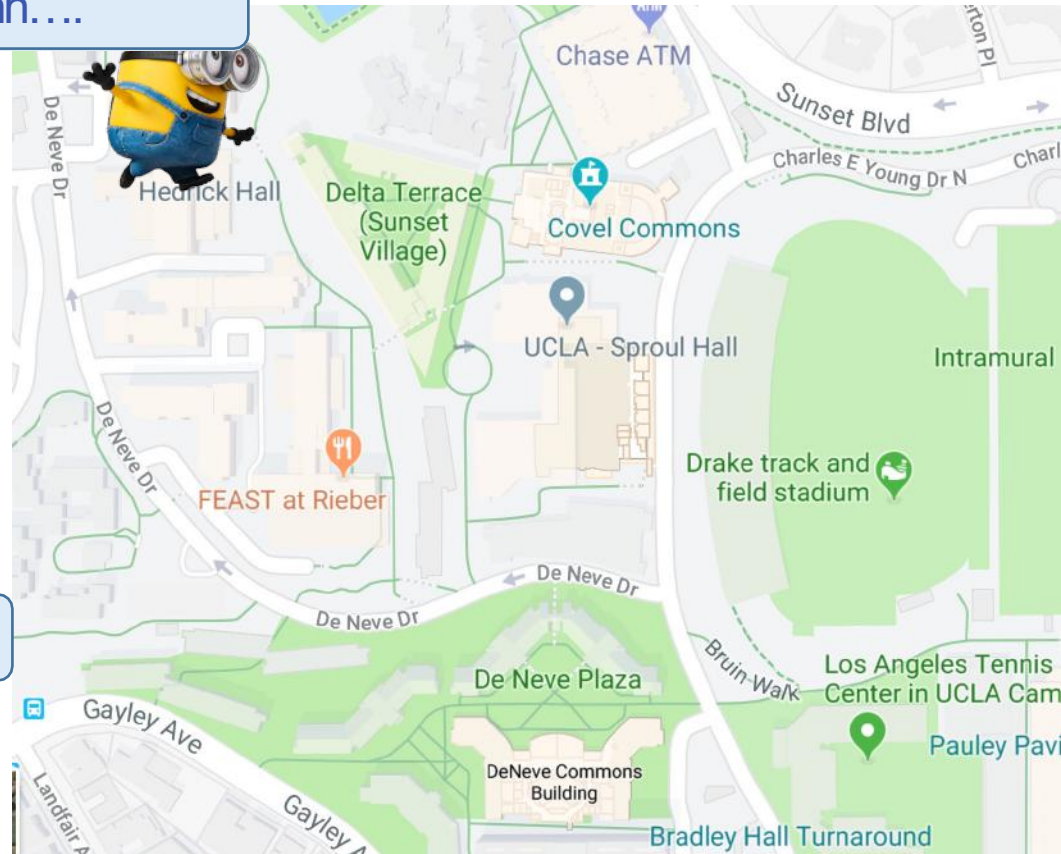


Intuition

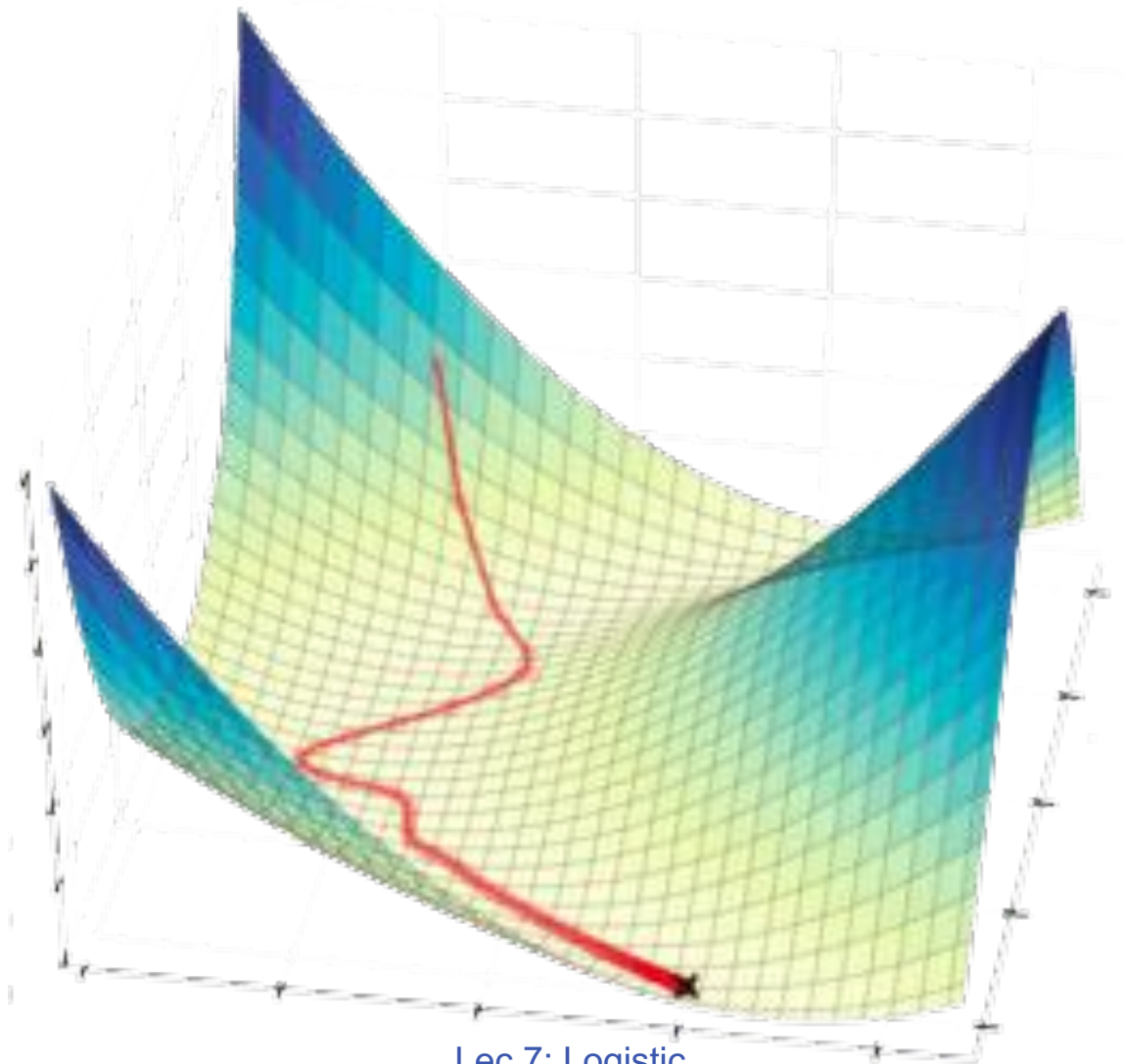
What may go wrong? Incorrect Direction

AHHhhhhhhhh....

Go left, hehehe...

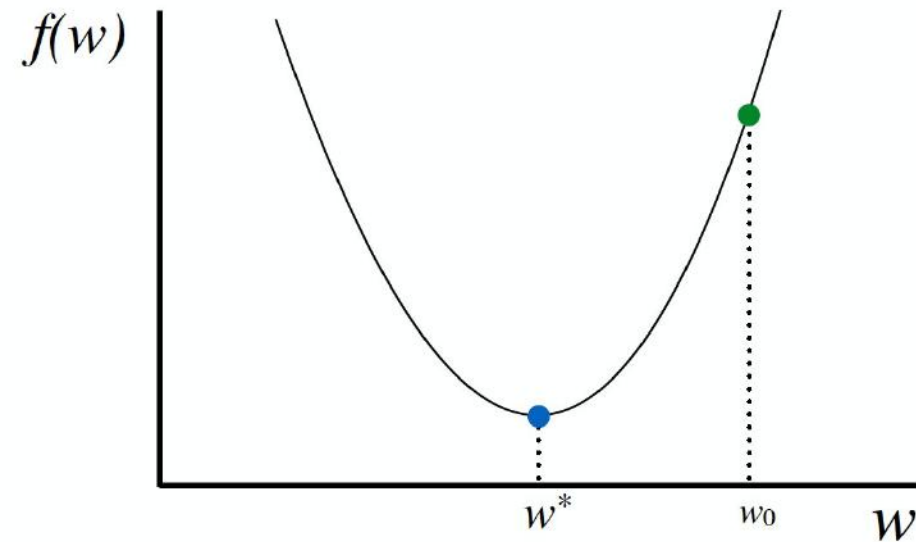


Gradient Descent



Gradient descent

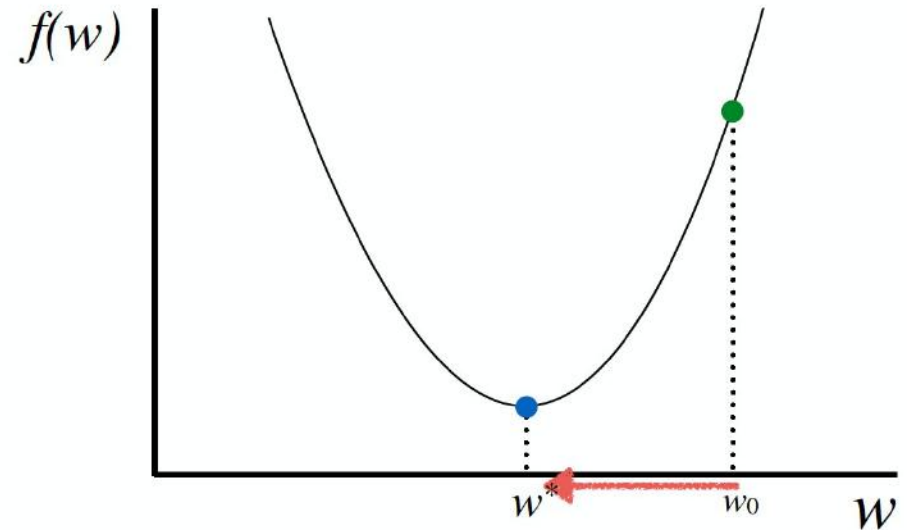
Start at a random point



Gradient descent

Start at a random point

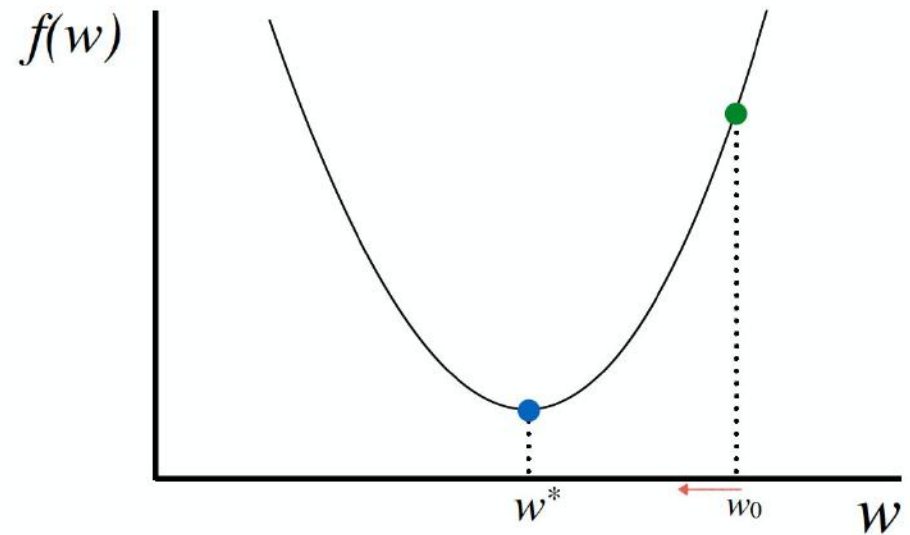
Determine a descent direction



Gradient descent

Start at a random point

Determine a descent direction
Choose a step size



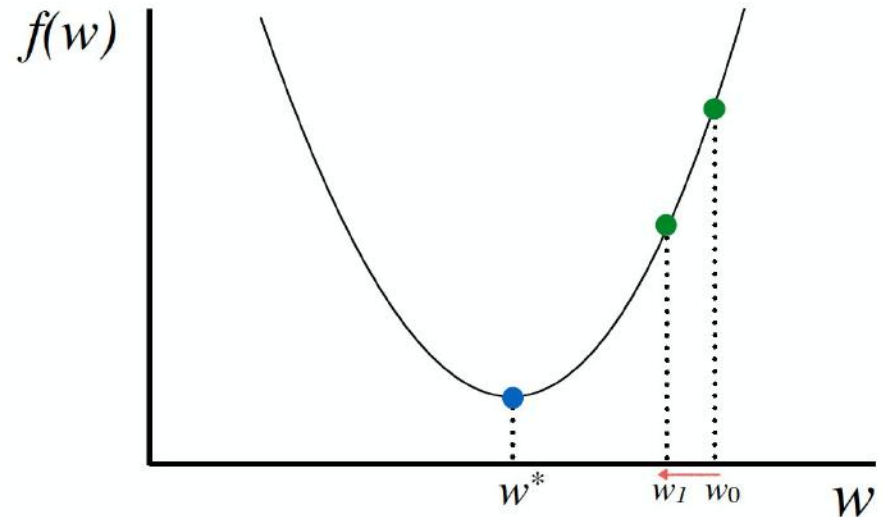
Gradient descent

Start at a random point

Determine a descent direction

Choose a step size

Update



Gradient descent

Start at a random point

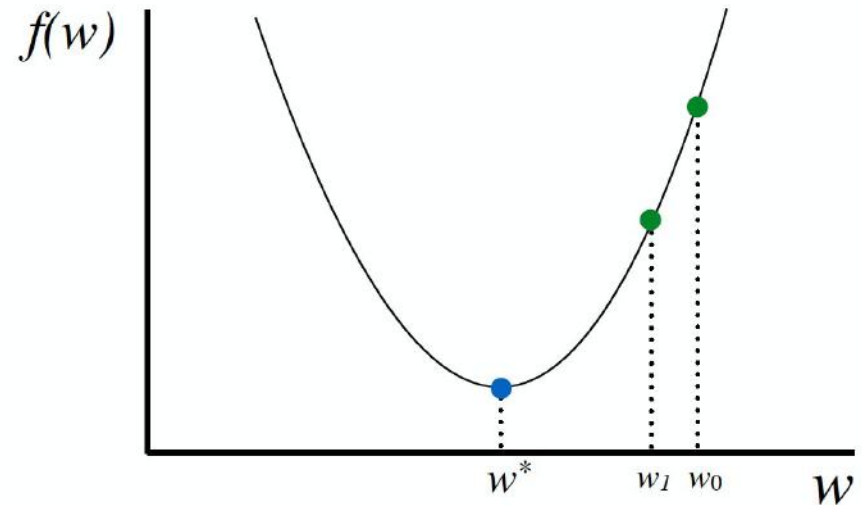
Repeat

Determine a descent direction

Choose a step size

Update

Until stopping criterion is satisfied



Gradient descent

Start at a random point

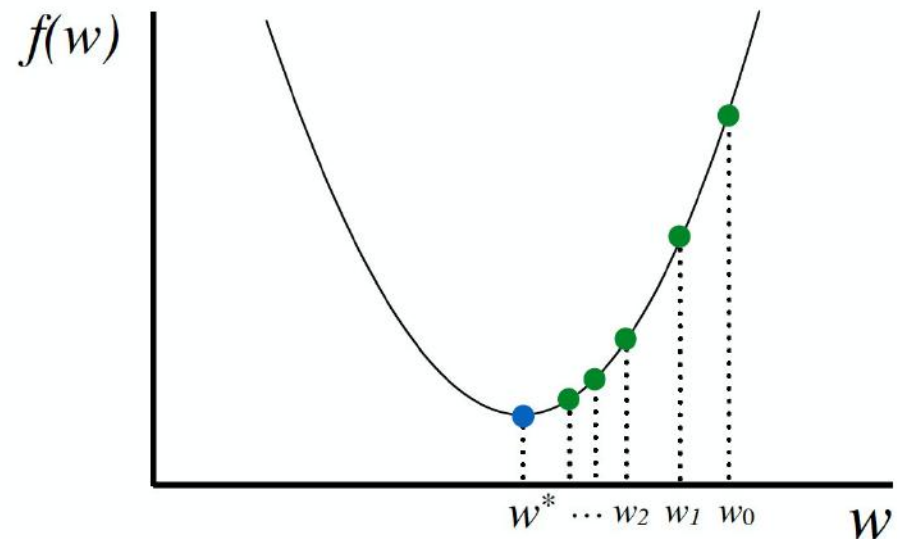
Repeat

Determine a descent direction

Choose a step size

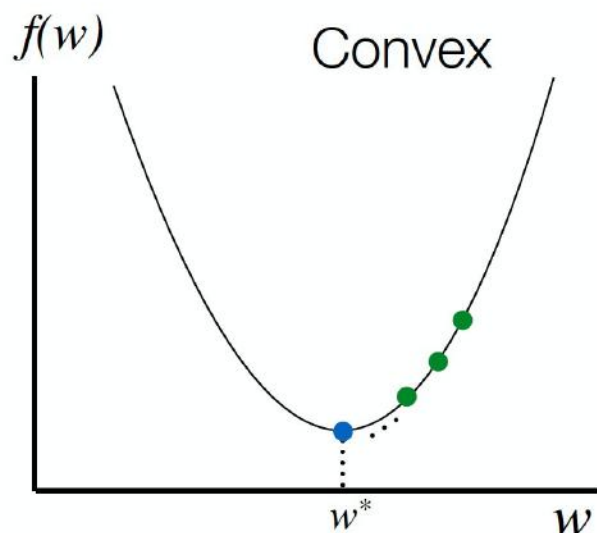
Update

Until stopping criterion is satisfied

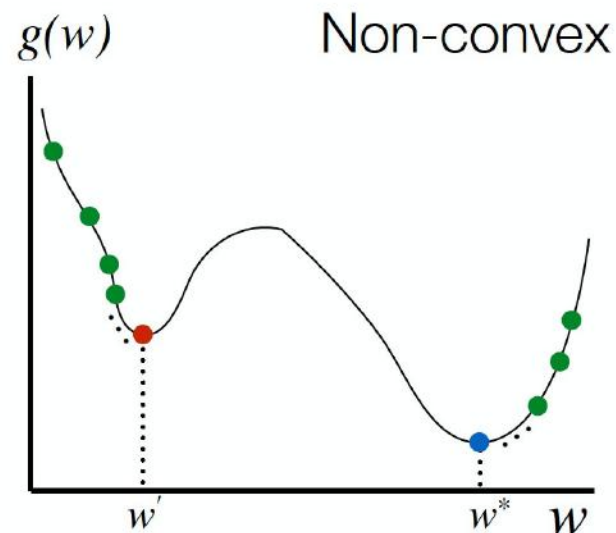


Where will we converge?

If function is convex, it converges to the global optimum (need proper choice of step-size)



Any local minimum is a global minimum



Multiple local minima may exist

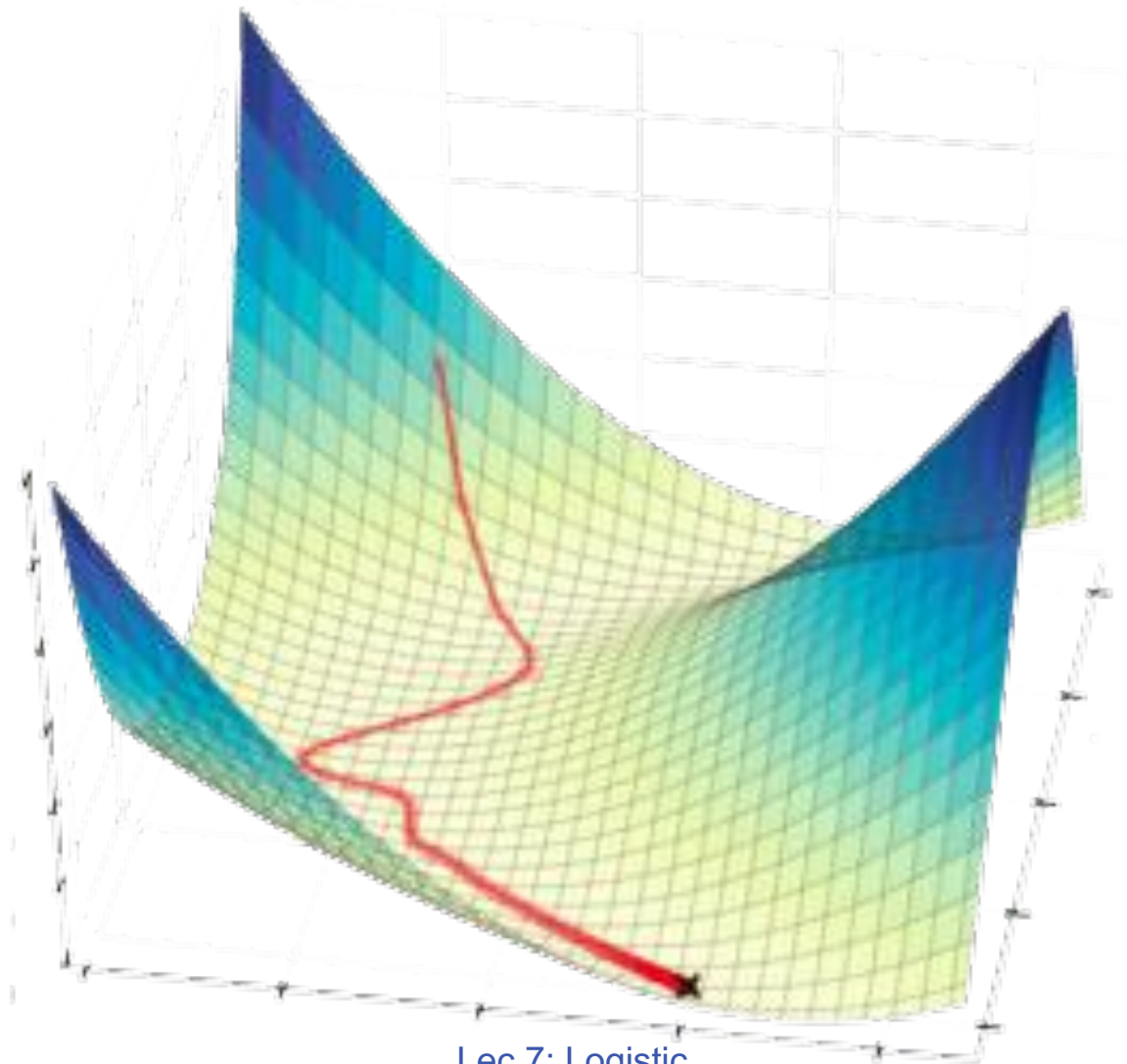
Gradient Descent

Algorithm 1 Gradient Descent (J)

```
1:  $t \leftarrow 0$ 
2: Initialize  $\theta^{(0)}$ 
3: repeat
4:    $\theta^{(t+1)} \leftarrow \theta^{(t)} - \eta \nabla J(\theta^{(t)})$ 
5:    $t \leftarrow t + 1$ 
6: until convergence
7: Return final value of  $\theta$ 
```

Need to compute the gradient for the negative log likelihood

Gradient Descent



Example

$$\min f(\boldsymbol{\theta}) = 0.5(\theta_1^2 - \theta_2)^2 + 0.5(\theta_1 - 1)^2$$

We compute the gradients

$$\frac{\partial f}{\partial \theta_1} = 2(\theta_1^2 - \theta_2)\theta_1 + \theta_1 - 1$$

$$\frac{\partial f}{\partial \theta_2} = -(\theta_1^2 - \theta_2)$$

Example $\min f(\boldsymbol{\theta}) = 0.5(\theta_1^2 - \theta_2)^2 + 0.5(\theta_1 - 1)^2$

❖ Use the following iterative procedure for gradient descent

$$\nabla f(\boldsymbol{\theta}) = \begin{bmatrix} 2(\theta_1^2 - \theta_2)\theta_1 + \theta_1 - 1 \\ -(\theta_1^2 - \theta_2) \end{bmatrix}$$

1 Initialize $\theta_1^{(0)}$ and $\theta_2^{(0)}$, and $t = 0$

2 do

Type equation here.

$$\theta_1^{(t+1)} \leftarrow \theta_1^{(t)} - \eta \left[2(\theta_1^{(t)^2} - \theta_2^{(t)})\theta_1^{(t)} + \theta_1^{(t)} - 1 \right]$$

$$\theta_2^{(t+1)} \leftarrow \theta_2^{(t)} - \eta \left[-(\theta_1^{(t)^2} - \theta_2^{(t)}) \right]$$

$$t \leftarrow t + 1$$

3 until $f(\boldsymbol{\theta}^{(t)})$ *does not change much*

Remarks

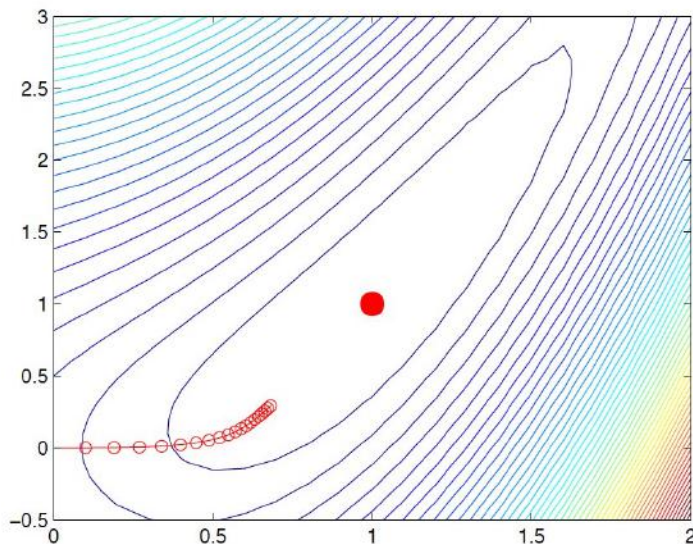
- ❖ η is often called step size or learning rate -- how far our update will go along the the direction of the negative gradient
- ❖ With a **suitable** choice of η , the iterative procedure converges to a stationary point where

$$\frac{\partial f}{\partial \theta} = 0$$

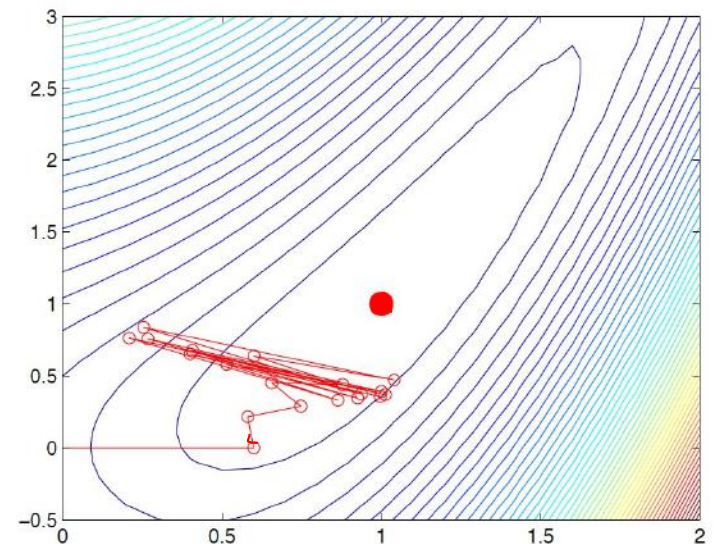
- ❖ A stationary point is only necessary for being the minimum

Choosing the right η is important

small η is too slow?

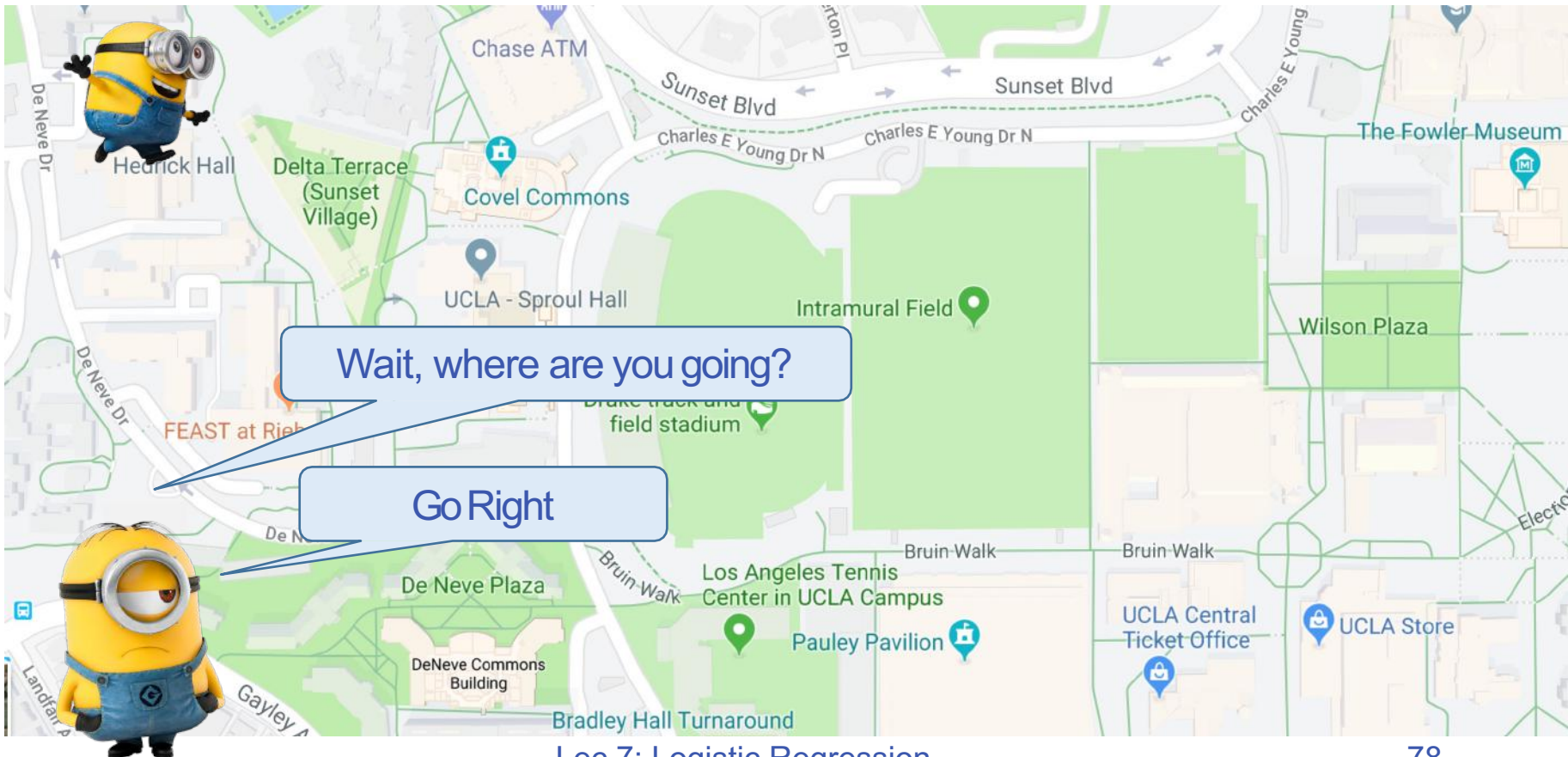


large η is too unstable?



Intuition

What may go wrong? Incorrect Step-size



Iterative optimization

Algorithm 2 Gradient Descent (J)

- 1: $t \leftarrow 0$
 - 2: Initialize $\theta^{(0)}$
 - 3: **repeat**
 - 4: $\theta^{(t+1)} \leftarrow \theta^{(t)} - \eta \nabla J(\theta^{(t)})$
 - 5: $t \leftarrow t + 1$
 - 6: **until** convergence
 - 7: Return final value of θ
-

How to compute the gradient?

Need to compute the gradient for the linear regression cost function (residual sum of squares RSS)

Stochastic Gradient Descent

Incremental/Stochastic gradient descent

Repeat for each example (\mathbf{x}_i, y_i)

Use this example to calculate the
gradient and update the model

Contrast with *batch gradient descent* which
makes one update to the weight vector for
every pass over the data

Stochastic gradient descent

$$\text{If } f(w) = \frac{1}{|D|} \sum_i^{|D|} f_i(w)$$

$$\nabla f(w) = \frac{1}{|D|} \sum_i \nabla f_i(w) = E_{i \sim D} \nabla f_i(w)$$

- ❖ Approximate the true gradient by a gradient at a single example at a time

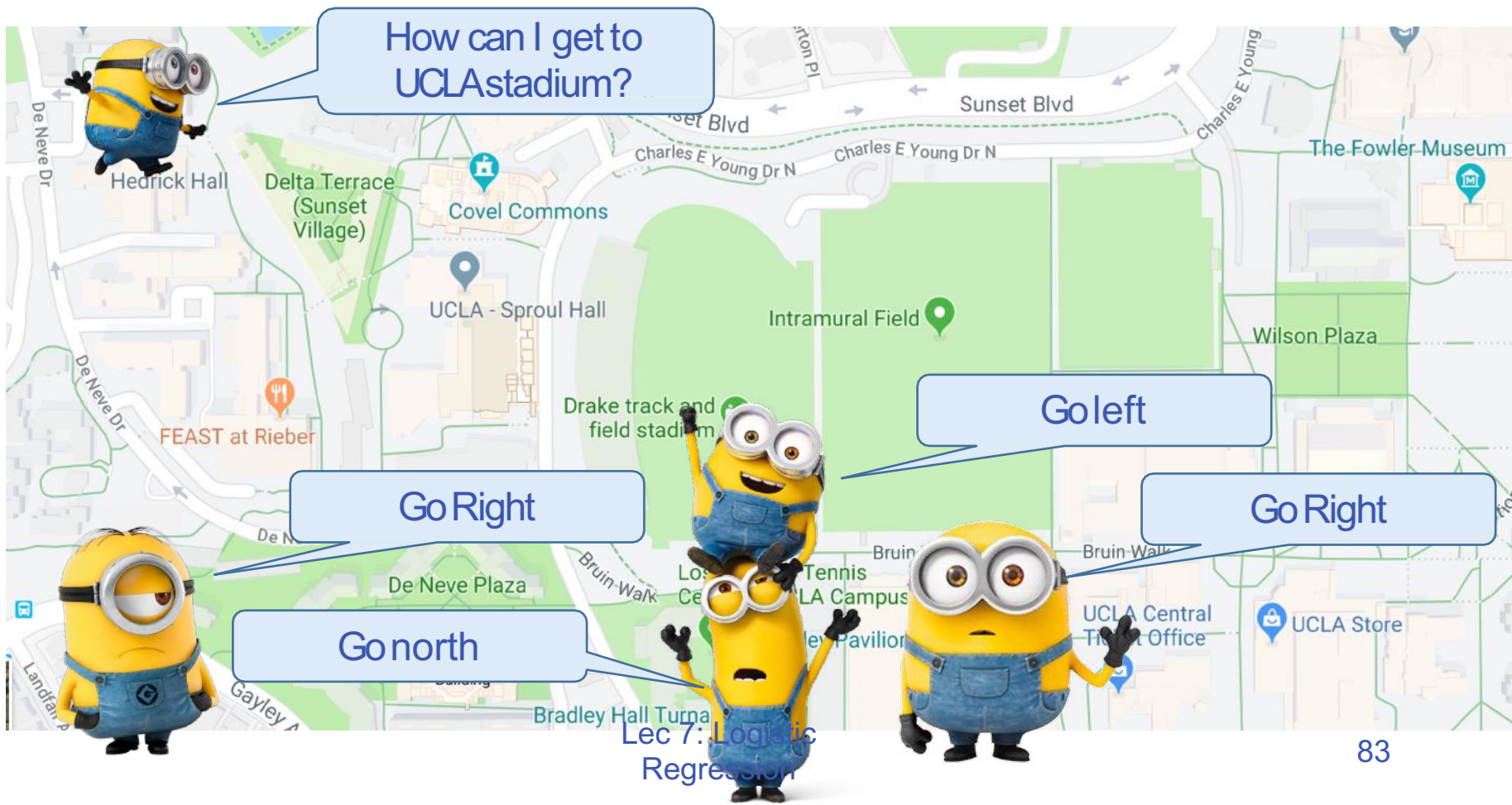
Repeat until converge:

Randomly pick one sample (x_i, y_i)

Update $w \leftarrow w - \eta \nabla f_i(w)$

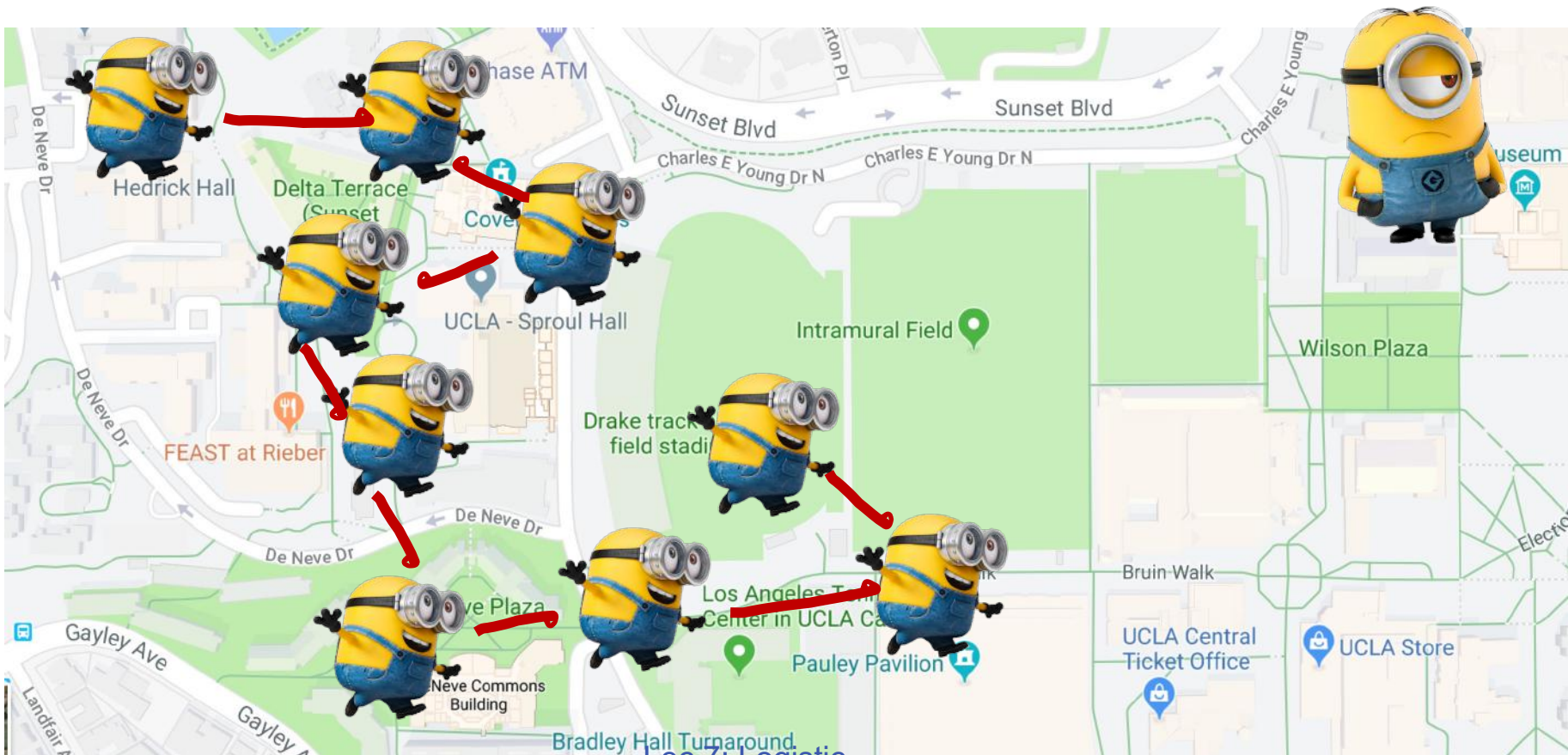
Intuition

Asking direction. Gradient descent:
compute gradient of all instances.



Intuition

Asking direction. Stochastic Gradient descent:
compute approximate gradient by one instance



Stochastic gradient Descent

Given a training set $\mathcal{D} = \{(\mathbf{x}, y)\}$

1. Initialize $\mathbf{w} \leftarrow \mathbf{0} \in \mathbb{R}^n$
2. For epoch $1 \dots T$:
3. For (\mathbf{x}, y) in \mathcal{D} :
4. Update $\mathbf{w} \leftarrow \mathbf{w} - \eta \nabla f(\mathbf{w})$
5. Return \mathbf{w}

The Perceptron Algorithm [Rosenblatt 1958]

Given a training set $\mathcal{D} = \{(\mathbf{x}, y)\}$

1. Initialize $\mathbf{w} \leftarrow \mathbf{0} \in \mathbb{R}^n$
2. For epoch $1 \dots T$:
3. For (\mathbf{x}, y) in \mathcal{D} :
4. if $y(\mathbf{w}^\top \mathbf{x}) < 0$
5. $\mathbf{w} \leftarrow \mathbf{w} + \eta y \mathbf{x}$
6. Return \mathbf{w}

Prediction: $y^{\text{test}} \leftarrow \text{sgn}(\mathbf{w}^\top \mathbf{x}^{\text{test}})$

The Perceptron Algorithm [Rosenblatt 1958]

Given a training set $\mathcal{D} = \{(\mathbf{x}, y)\}$

1. Initialize $\mathbf{w} \leftarrow \mathbf{0} \in \mathbb{R}^n$
2. For epoch $1 \dots T$:
3. For (\mathbf{x}, y) in \mathcal{D} :
4. if $y(\mathbf{w}^\top \mathbf{x}) < 0$
5. $\mathbf{w} \leftarrow \mathbf{w} + \eta y \mathbf{x}$
6. Return \mathbf{w}

Prediction: y^{test}

Perceptron effectively minimizing:

$$\sum_i \max(0, 1 - y_i(\mathbf{w}^\top \mathbf{x}_i))$$

Summary

- ❖ Maximum Likelihood Estimation
- ❖ Gradient Descent
- ❖ Stochastic Gradient Descent

- ❖ We will see more examples in later lectures