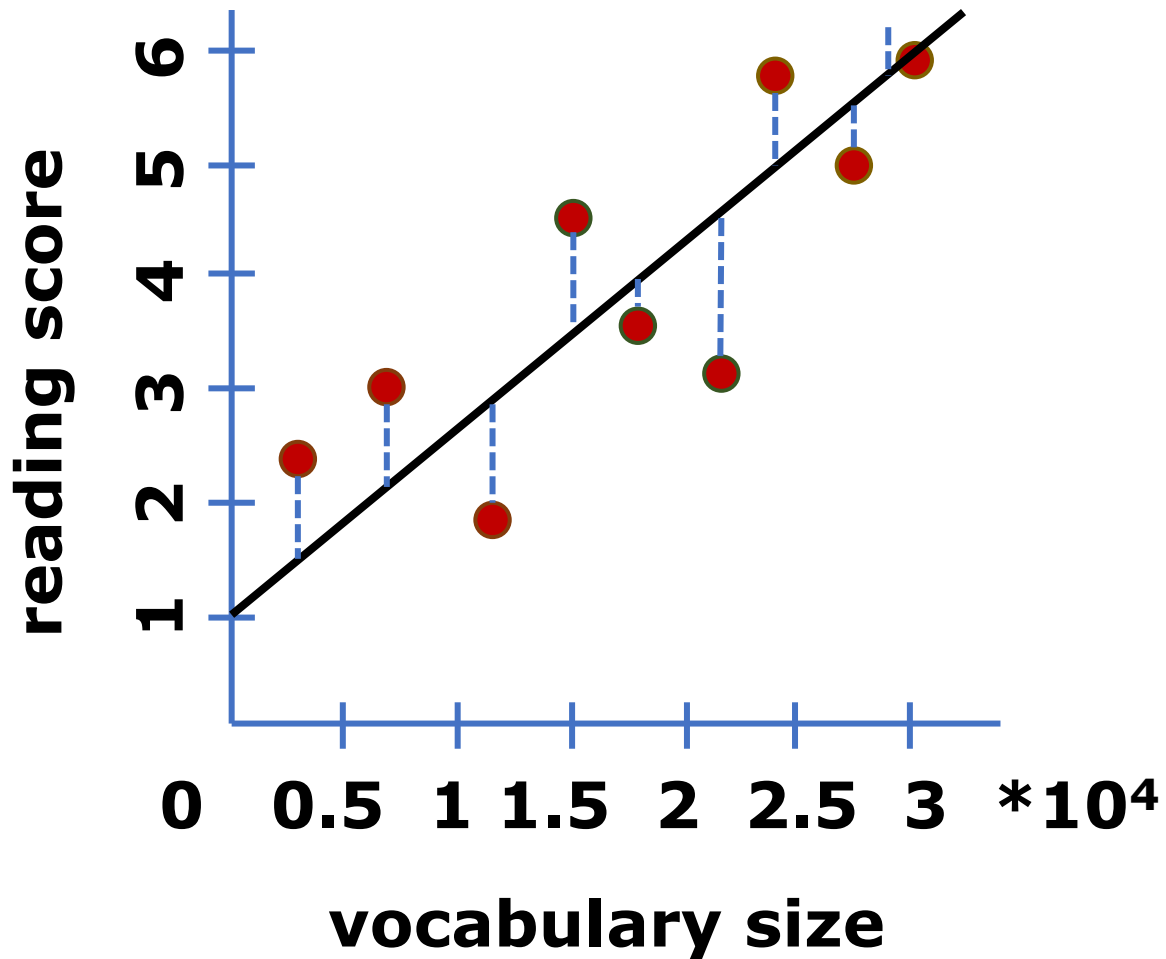# Fundamentals of Statistics for Language Sciences LT2206
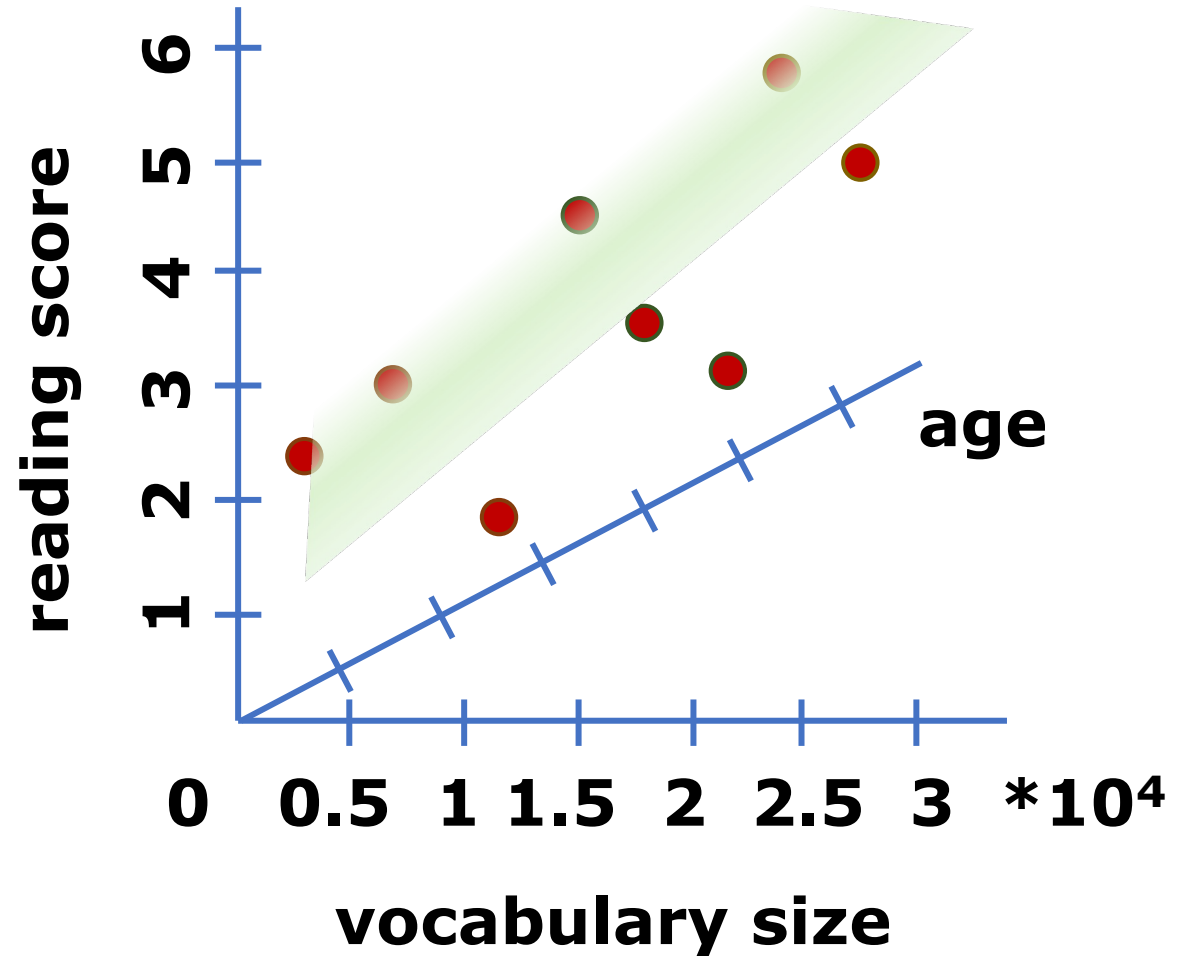
Jixing Li

Lecture 11: Logistic Regression

# Simple linear regression & multiple regression



$$y = b_1 x + b_0$$

$$y = b_1 x_1 + b_2 x_2 + b_0$$

# Model comparison

| Model | F | p | R² | R² adjusted | VIF |
|---|---|---|---|---|---|
| RATL | 67.77 | 3.419e-15 | 0.158 | 0.156 | 1 |
| LPTL | 115 | < 2.2e-16 | 0.242 | 0.24 | 1 |
| LIFG | 38.87 | 1.27e-09 | 0.097 | 0.095 | 1 |
| RATL+LPTL | 95.09 | < 2.2e-16 | 0.346 | 0.343 | 1.03 |
| RATL+LIFG | 34.88 | 1.438e-14 | 0.163 | 0.158 | 1.78 |
| LPTL+LIFG | 81.8 | < 2.2e-16 | 0.313 | 0.309 | 1.01 |
| RATL+LPTL+LIFG | 64.67 | < 2.2e-16 | 0.352 | 0.346 | 1.81,1.03,1.78 |

# Logistic regression

When the dependent variable ($y$) is binary (0 or 1):

e.g., a person's name is male or female?

a movie review if positive or negative?

an email is spam or not?

The task of text classification

- *Input*:
  - a document $x$
  - two classes $C = \{c_1, c_2\}$

- *Output*: a predicted class $\hat{y} \in C$

# Features in logistic regression

**Input vector:** $x = [x_1, x_2, ..., x_n]$

[卓, 琳, Cheuk, Lam, LLA]

Probability of these features in female names:

→ $x = [0.5, 0.7, 0.5, 0.6, 0.8]$

**Weights:** one per feature: $w = [w_1, w_2, ..., w_n]$

→ $w = [0.1, 0.8, -0.1, 0.2, 0.7]$

**Prediction:** $z = w \cdot x + b$

$z = w_1{*}x_1 + w_2{*}x_2 + w_3{*}x_3 + w_4{*}x_4 + w_5{*}x_5 + b$

$\quad = 0.05 + 0.56 + (-0.05) + 0.12 + 0.56 + 0.3$

$\quad = 1.54$

# Transform prediction into probability
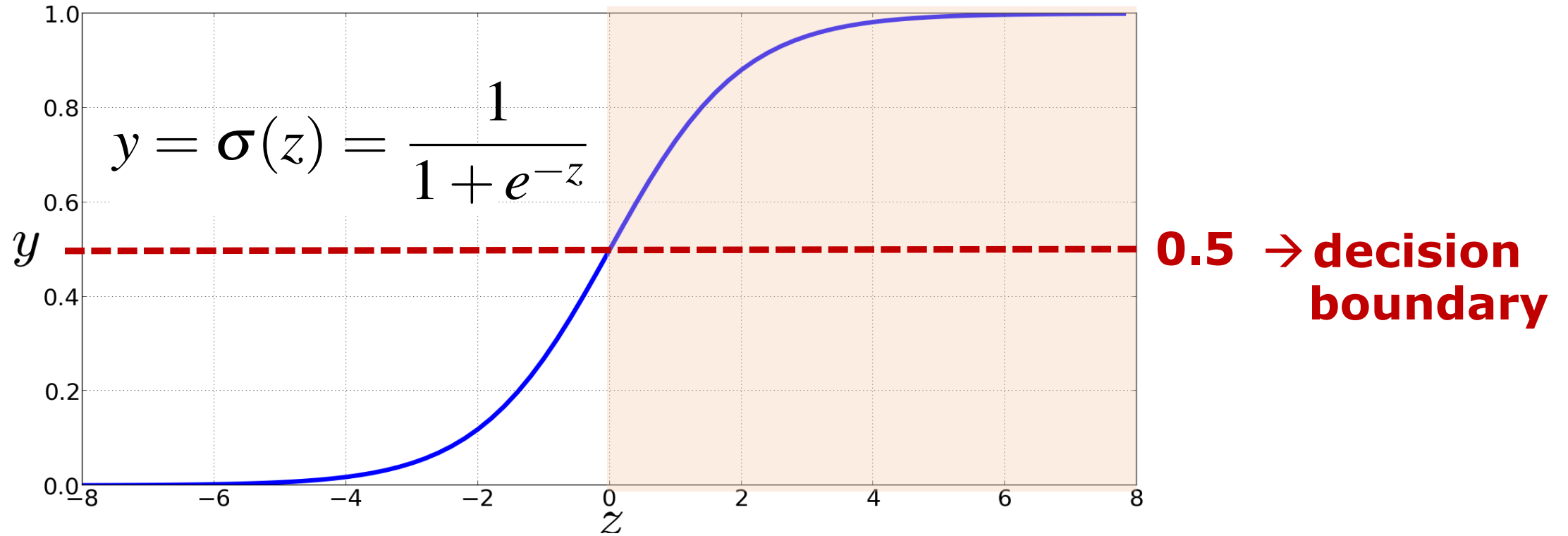
$$z \;=\; w \cdot x + b$$

$z$ is a number, but we We'd like a classifier that gives us a probability

**Solution:** use a function of $z$ that goes from 0 to 1

$$y = \sigma(z) = \frac{1}{1 + e^{-z}} = \frac{1}{1 + \exp(-z)} \qquad \rightarrow \textbf{ the sigmoid function}$$

# The sigmoid function

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

$0.5 \rightarrow$ **decision boundary**

$$\hat{y} = \begin{cases} 1 & \text{if } w{\cdot}x{+}b > 0 \\ 0 & \text{if } w{\cdot}x{+}b \leq 0 \end{cases}$$

# Example

[卓, 琳, Cheuk, Lam, LLA]

$x = [0.5, 0.7, 0.5, 0.6, 0.8]$

$w = [0.1, 0.8, -0.1, 0.2, 0.7]$

$z = w \cdot x + b$

$= w_1*x_1 + w_2*x_2 + w_3*x_3 + w_4*x_4 + w_5*x_5 + b$

$= 0.05 + 0.56 + (-0.05) + 0.12 + 0.56 + 0.3$

$= 1.54$

$\hat{y} = \sigma(z) = \frac{1}{1+e^{-z}} = \frac{1}{1+e^{-1.54}} = 0.82$  $> 0.5$ → female

# How to calculate weights?

**Supervised classification:**

We know the correct label $y$ (either 0 or 1) for each $x$.

But what the system produces is an estimate, $\hat{y}$

We want to know how far is the classifier output:

$\hat{y} = \sigma(w \cdot x + b)$

from the true output:

$y$ = either 0 or 1

We'll call this difference the loss:

$L(\hat{y}, y)$ = how much $\hat{y}$ differs from the true $y$

# Binary cross-entropy loss

**Goal**: **maximize** the probability of the correct label $p(y|x)$

Since there are only 2 outcomes (0 or 1), we can express the probability $p(y|x)$ from our classifier as:

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

if y=1, this simplifies to $\hat{y}$
if y=0, this simplifies to 1- $\hat{y}$

Now take the log of both sides:

$$\log p(y|x) = \log\left[\hat{y}^y (1-\hat{y})^{1-y}\right]$$
$$= y\log\hat{y} + (1-y)\log(1-\hat{y})$$

Now flip sign to turn this into a loss: Something to **minimize**

$$L_{\text{CE}}(\hat{y},y) = -\log p(y|x) = -\left[y\log\hat{y} + (1-y)\log(1-\hat{y})\right]$$

**cross-entropy loss:** negative log likelihood loss

# Example

[卓, 琳, Cheuk, Lam, LLA]

x = [0.5, 0.7, 0.5, 0.6, 0.8]

w = [0.1, 0.8, -0.1, 0.2, 0.7]

b = 0.5

$\hat{y} = \sigma(w \cdot x + b) = 0.82$

if 卓琳 is female: y = 1:

$L_{CE}(\hat{y}, y)$ = -(ylog $\hat{y}$ +(1-y)log(1- $\hat{y}$)) = $-\log(0.82) = 0.2$

if 卓琳 is male: y = 0:

$L_{CE}(\hat{y}, y)$ = -(ylog $\hat{y}$ +(1-y)log(1- $\hat{y}$)) = $-\log(1 - 0.82) = 1.7$

→ **The loss is greater when the prediction is wrong**

# Minimize the loss

Let's make explicit that the loss function is parameterized by weights *θ=(w,b)*

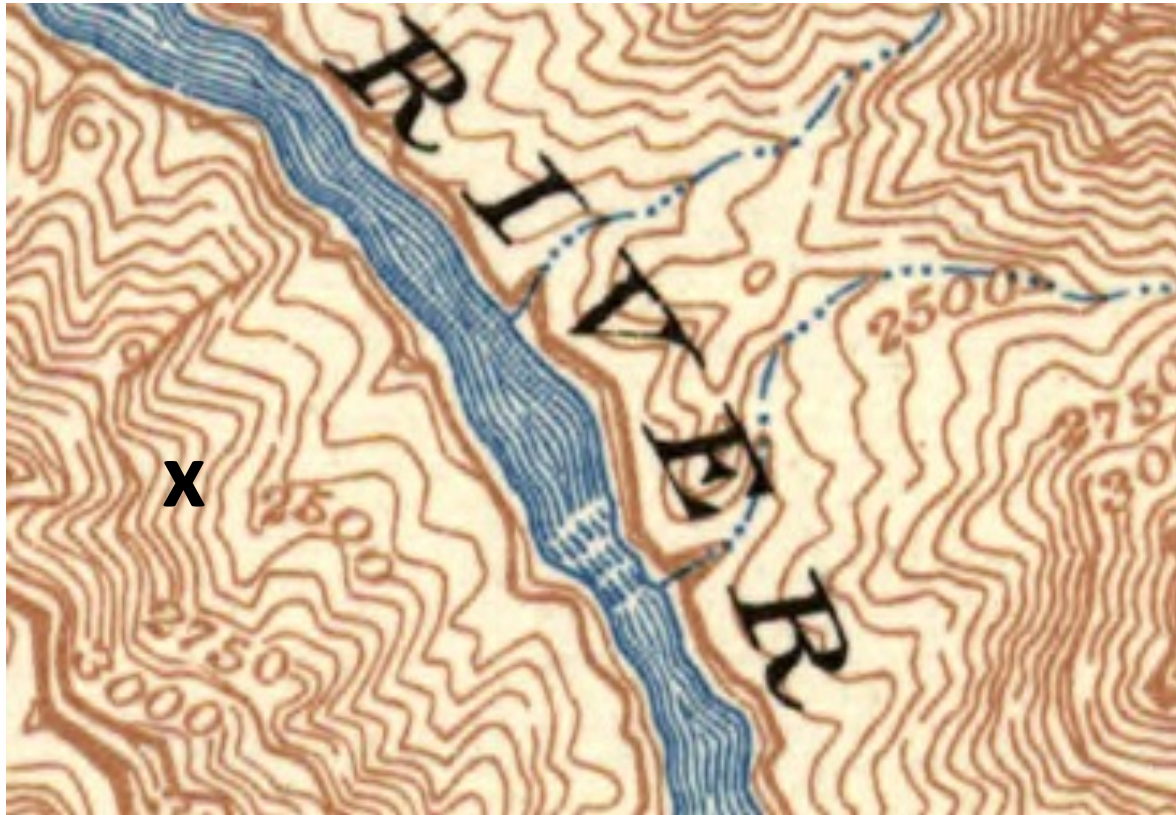And we'll represent $\hat{y}$ as *f(x;θ)* to make the dependence on θ more obvious

We want the weights that minimize the loss, averaged over all examples:

$$\hat{\theta} \;=\; \operatorname*{argmin}_{\theta} \frac{1}{m} \sum_{i=1}^{m} L_{\mathrm{CE}}(f(x^{(i)};\boldsymbol{\theta}),y^{(i)})$$
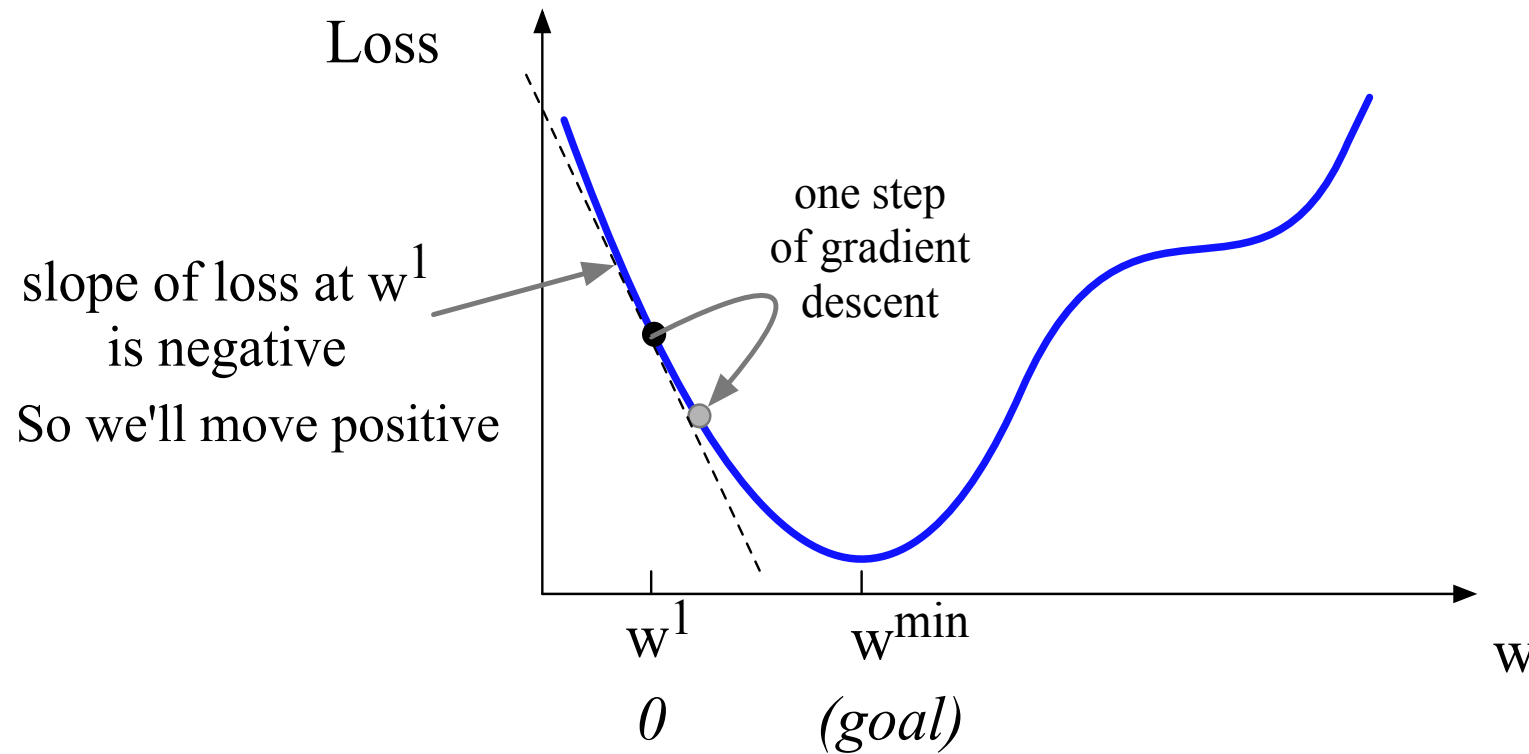
# Gradient descent

How do I get to the bottom of this river canyon?



Look around me 360°

Find the direction of steepest slope down

Go that way

# Gradient descent for a single scaler

**Minimize loss:** Given the current $w$, Move $w$ in the reverse direction from the slope of the function



Loss

slope of loss at $w^1$
is negative

So we'll move positive

one step
of gradient
descent

$w^1$

$w^{min}$

$0$

*(goal)*

$w$

The **gradient** of a function of many variables is a vector pointing in the direction of the greatest increase in a function.

**Gradient descent:** Find the gradient of the loss function at the current point and move in the opposite direction.

# Gradient descent

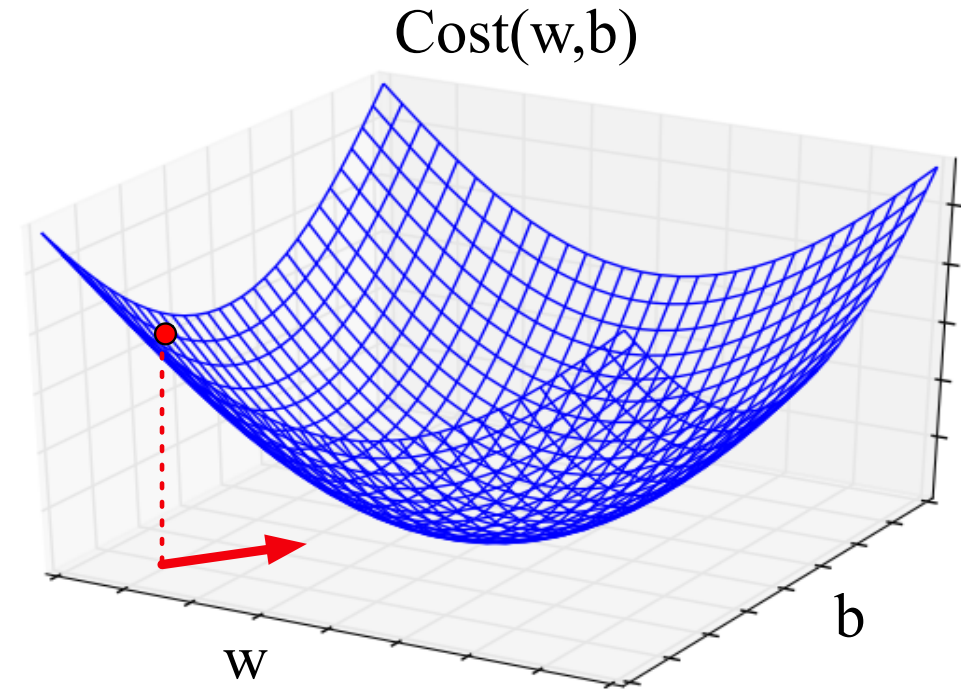The new weight $w^{t+1}$ is the old weight $w^t$ minus the value of the gradient weighted by a learning rate η

$$w^{t+1} = w^t - \eta \frac{d}{dw} L(f(x;w), y)$$

learning rate: Higher learning rate means move w faster → a **hyperparameter** not learned by algorithm from supervision, but are chosen by algorithm designer.

gradient (a vector of the derivatives with respect to the weight w)

# Gradient in N-dimensional space

The gradient expresses the directional components of the sharpest slope along each of the N dimensions. For each dimension $w_i$, we express the slope as a partial derivative $\partial$ of the loss $\partial w_i$



Cost(w,b)

b

w

The derivative of

$$L_{\mathrm{CE}}(\hat{y}, y) = -\left[ y \log \sigma(w \cdot x + b) + (1-y) \log(1 - \sigma(w \cdot x + b)) \right]$$

is:

$$\frac{\partial L_{\mathrm{CE}}(\hat{y}, y)}{\partial w_j} = \left[ \sigma(w \cdot x + b) - y \right] x_j$$

# Example

[卓, 琳, Cheuk, Lam, LLA]
x = [0.5, 0.7, 0.5, 0.6, 0.8]

1. initialize w and b, set η
w = [0, 0, 0, 0, 0], b = 0, η = 0.1

2. compute $\hat{y}$
$\hat{y}$ = σ(w · x + b) = 0.5

3. compute the gradients for w and b
Gw = (0.5− y)x = -0.5x = [-0.25, -0.35, -0.25, -0.3, -0.4]
Gb = 0.5− y = -0.5

4. update w and b
$w_{t+1}$ = $w_t$ − η*Gw = [0, 0, 0, 0, 0] − 0.1*[-0.25, -0.35, -0.25, -0.3, -0.4] = [0.025, 0.035, 0.205, 0.03, 0.04]
$b_{t+1}$ = $b_t$ − η*Gb = 0- 0.1*(-0.5) = 0.05

# Calculate gradient descent over all examples

[卓, 琳, Cheuk, Lam, LLA] $x_1 = [0.5, 0.7, 0.5, 0.6, 0.8]$
[承, 璋, Shing Cheung, LLA] $x_2 = [-0.6, -0.8, -0.1, -0.6, 0.8]$

1. initialize w and b, set η
$w = [0, 0, 0, 0, 0], b = 0, \eta = 0.1$

2. compute $\hat{y}$
$\hat{y}_1 = \sigma(w \cdot x + b) = 0.5, \hat{y}_2 = \sigma(w \cdot x + b) = 0.5$

3. compute the gradients for w and b
$Gw = \frac{1}{2}((0.5 - y)x_1 + (0.5 - y)x_2) = \frac{1}{2}(-0.5x_1 - 0.5x_2) = [0.025, 0.025, -0.1, 0, -0.2]$
$Gb = \frac{1}{2}((0.5 - y_1) + (0.5 - y_2)) = 0$

4. update w and b
$w_{t+1} = w_t - \eta*Gw = [0, 0, 0, 0, 0] - 0.1*[0.025, 0.025, -0.1, 0, -0.2] = [-0.0025, -0.0025, 0.01, 0, 0.02]$, $b_{t+1} = b_t - \eta*Gb = 0$