# Computational Linguistics LT3233



Jixing Li

Lecture 7: Logistic Regression

Slides adapted from Dan Jurafsky

# Lecture plan

- Naive Bayes and Laplace smoothing review
- Logistic regression
  - feature representation
  - classification function: **sigmoid**
  - loss function: **cross-entropy loss**
  - optimization algorithm: **gradient descent**
- Short break (15 mins)
- Hands-on exercises

# Logistic regression

The task of text classification

- *Input*:
  - a document $x$
  - a fixed set of classes $C = \{c_1, c_2, ..., c_J\}$

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

**Naive Bayes:** Compare P (male|卓琳) P (female|卓琳)
→ **generative classifier**

**Logistic regression:** P (male|卓琳)
→ **discriminative classifier**

# Components of logistic regression

1. **feature representation** of the input. For each input observation $x^{(i)}$, a vector of features $[x_1, x_2, \ldots, x_n]$.

   $x^{(i)} = [卓，琳]$

   $\phantom{x^{(i)}} = [卓，琳, \text{Cheuk, Lam, LLA}]$

2. **classification function** that computes $\hat{y}$, the estimated class: **sigmoid** functions

3. **objective function for learning**: **cross-entropy loss**

4. **algorithm for optimizing** the objective function: **gradient descent**

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

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

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

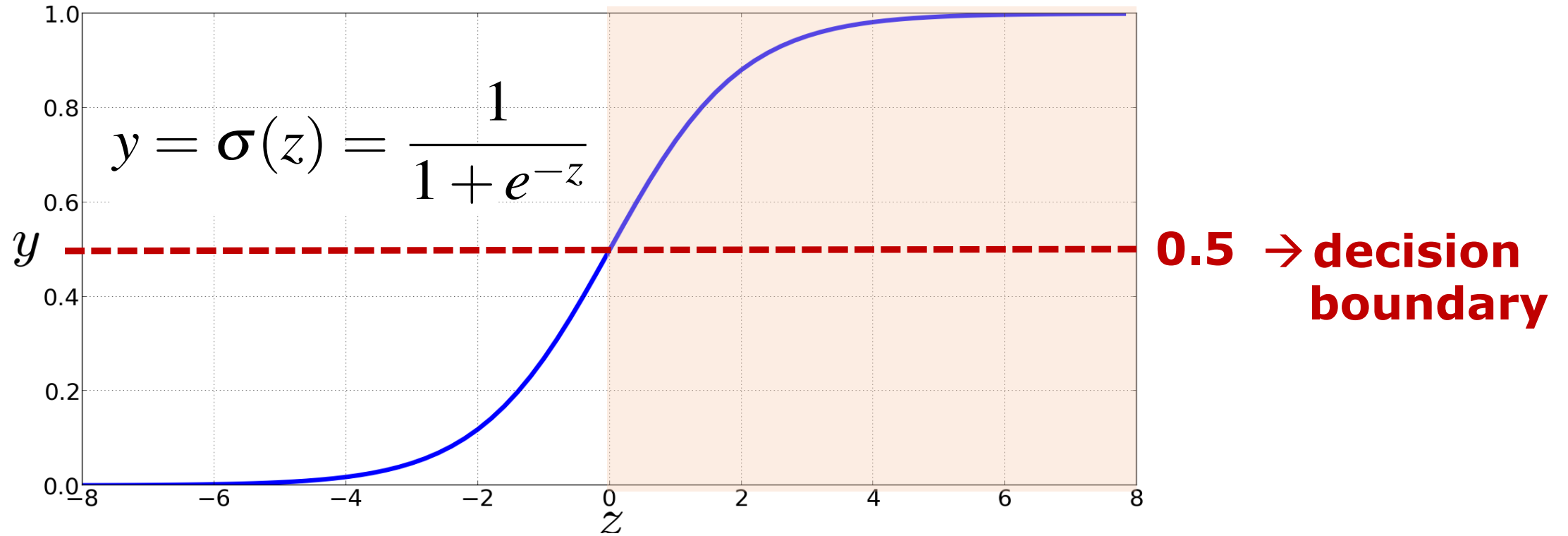$\quad = 1.54$

# 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, just like Naive Bayes did

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

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

© Jixing Li

# The sigmoid function



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

0.5 → decision boundary

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

© Jixing Li

# 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 · x + b

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

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

$= 1.54$

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

© Jixing Li

# 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:

$$y = \text{either 0 or 1}$$

We'll call this difference the loss:

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

# 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:

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

if 卓琳 is male: y = 0:

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

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

# Minimize the loss

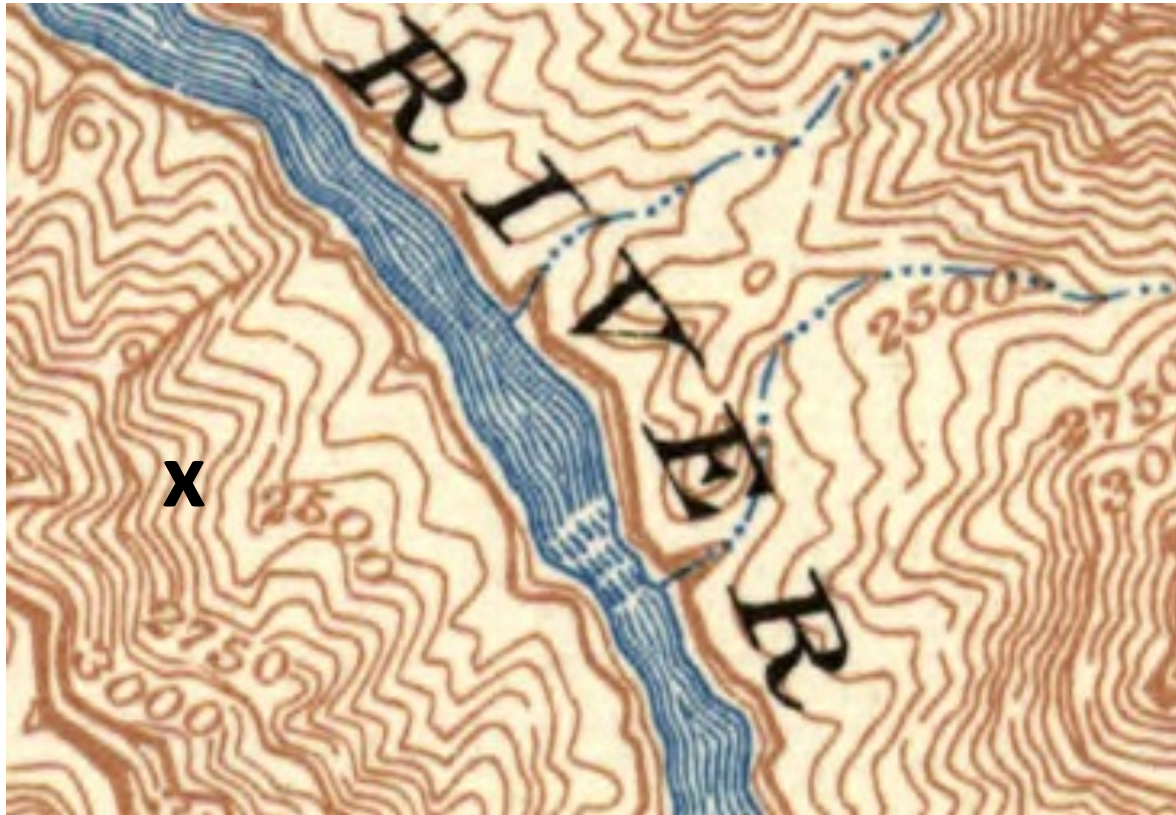Let's make explicit that the loss function is parameterized by weights $\theta=(w,b)$

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

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

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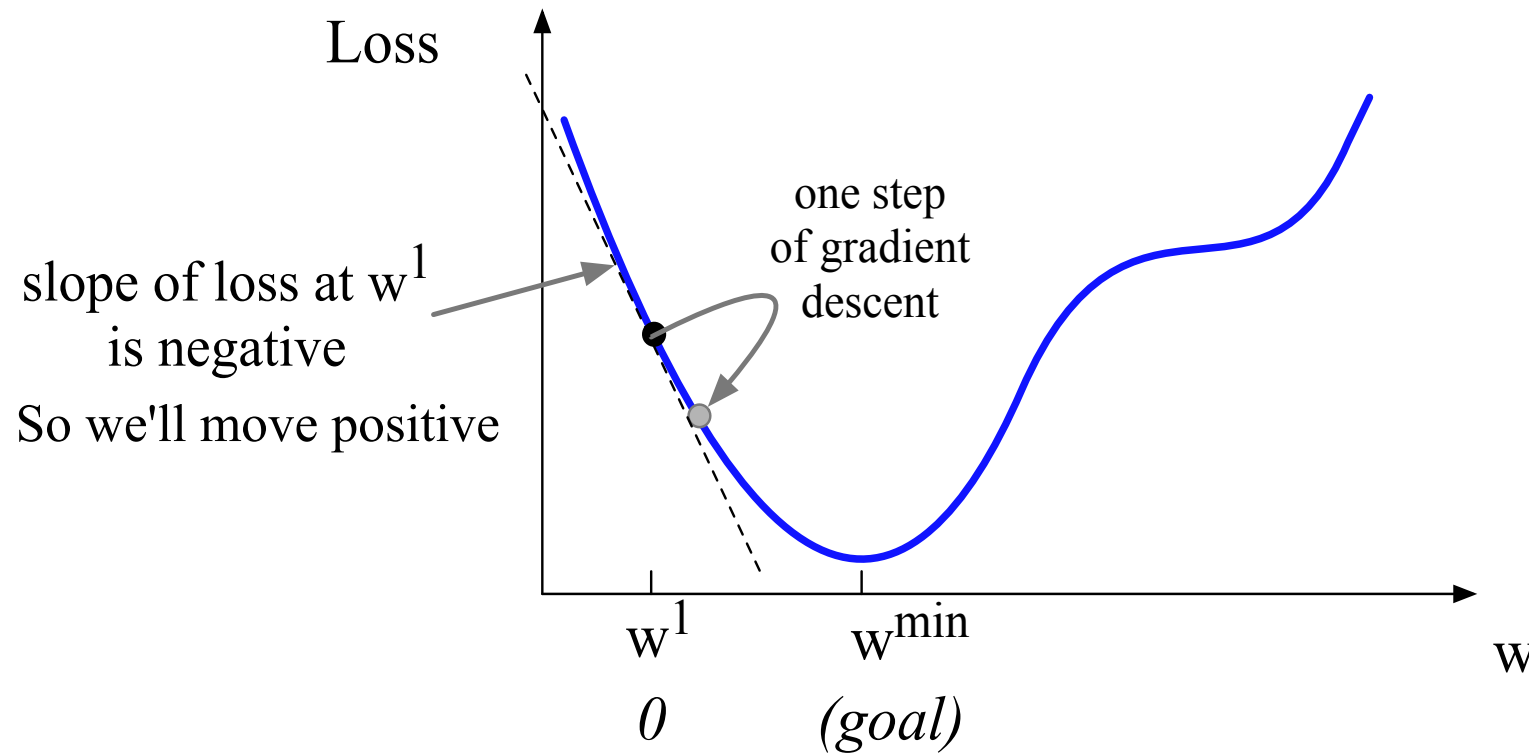
# Gradient descent

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

$w^1$

$w^{min}$

$w$

*0*

*(goal)*

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.

© Jixing Li

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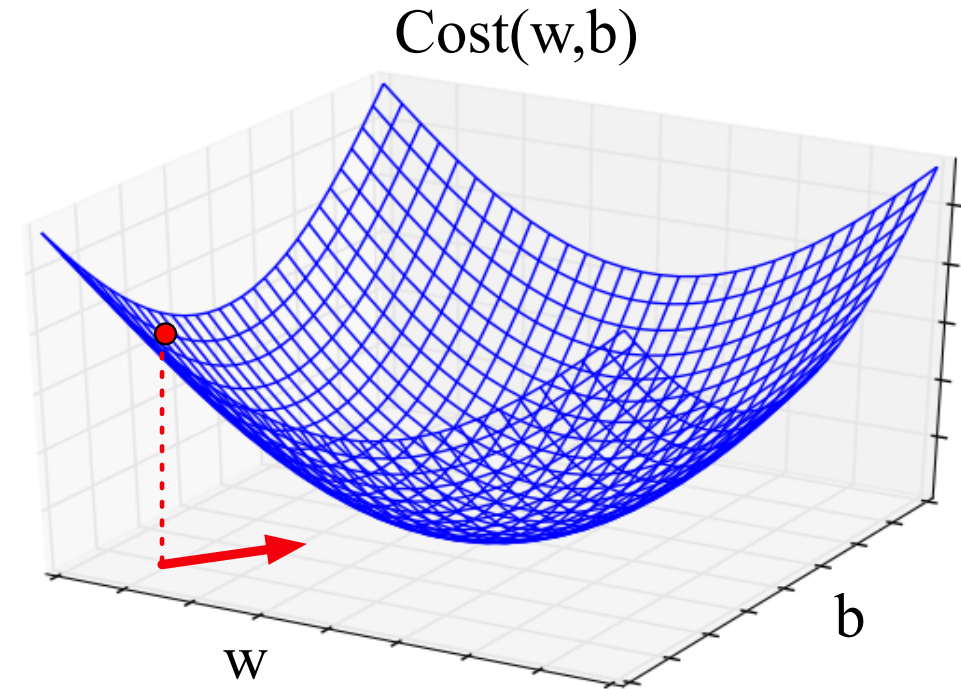
$$w^{t+1} = w^t - \eta \frac{d}{dw} L(f(x; w), y)$$

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

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

is:

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

# 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 $\hat{y}$
$\hat{y} = σ(w \cdot x + b) = 0.5$

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
$w = [0, 0, 0, 0, 0], b = 0, η = 0.1$

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

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 $b$
$w_{t+1} = w_t - η*Gw = [0, 0, 0, 0, 0] - 0.1*[0.025, 0.025, -0.1, 0, -0.2] = [-0.0025, -0.0025, 0.01, 0, 0.02]$, $b_{t+1} = b_t - η*Gb = 0$

# To do

- Optional reading: **SLP** Ch5