

# **ECE/SIOC 228 Machine Learning for Physical Applications**

## **Lecture 2: Supervised Learning Setup**

**Instructor:** Yuanyuan Shi

Teaching Assistants:

Srinivas Rao Daru, [sdaru@ucsd.edu](mailto:sdaru@ucsd.edu)

Tawaana Homavazir, [thomavaz@ucsd.edu](mailto:thomavaz@ucsd.edu)

Rohin Garg, [rgarg@ucsd.edu](mailto:rgarg@ucsd.edu)

Rishabh Jangir, [rjangir@ucsd.edu](mailto:rjangir@ucsd.edu),

# Logistics

---

Exam:

- Midterm exam will be on May 17th Tuesday (2:00pm - 3:20pm)
- It will be an closed-book exam and you are allowed to bring a A4 paper written note (two-sided)
- Based on what is covered on lectures

Office hour update:

- My office hour: 4-5pm Thursdays
- Jacobs Hall Room 4401

Survey Results

# Logistics

---

## Project

- Many of you have already started to form teams - that is great!!
- For those haven't started - start early :)
- Feel free to discuss your idea at any of the TAs' office hour and my office hour
- Feel free to reach out to your assigned TA for idea discussion if you're not able to attend their office hours; we'll adjust the TA assignment after Week 3 when the teams are formed.
- Please don't share individual docs and writeup with us - the proposal is due on Week 3 and we won't be able to provide feedback on writings before then.

# Today's Lecture

---

- Supervised Learning Setup
- Linear Regression
- Logistic Regression
- Computation Graph

Reference & Acknowledgement: Stanford CS230 Deep Learning Lecture Note, Andrew Ng

# Supervised Learning Setup

---

Suppose we have  $m$  pairs of data points:  $\{(x^{(1)}, y^{(1)}), (x^{(2)}, y^{(2)}), \dots, (x^{(m)}, y^{(m)})\}$ , where for each data  $(x^{(i)}, y^{(i)})$ : Specifically,

- $x^{(i)} \in \mathbb{R}^n$  is the input vector for the  $i^{th}$  sample, and  $\mathbb{R}^n$  is the  $n$ -dimension feature space
- $y^{(i)} \in C$  is the label of the  $i^{th}$  sample and  $C$  is the label space

Example of label spaces:

- $C = \mathbb{R}$  Regression: predicting the electricity consumption of a house given number of rooms, square feet, outside temperature
- $C = \{0, 1\}$  Binary classification: predicting whether this email is a spam (1) or not (0).
- $C = \{1, 2, \dots, K\}$  Multi-class classification: CIFAR-10 (cat, deer, dog, frog, horse, bird, airplane, ...)

# Supervised Learning Setup

Examples of feature vectors:

- Structured data
  - Electricity consumption data.  $x^{(i)} = (x_1^{(i)}, x_2^{(i)}, \dots, x_n^{(i)})$ , where  $x_1^{(i)} = \{0, 1, 2, \dots\}$  may refer to the number of rooms in a house,  $x_2^{(i)}$  could be the square feet, etc.
- Unstructured data **Yawn**



Chinese (Simplified)  English

知识就是力量  
Zhishi jiùshì lìlìng

knowledge is power

Open in Google Translate • Feedback

A screenshot of a translation interface. It shows a Chinese phrase "知识就是力量" (Knowledge is power) on the left and its English translation "knowledge is power" on the right. Both phrases are underlined in red. The interface includes language selection dropdowns, a red "Verified" badge, and standard translation controls like a microphone icon and a "Feedback" link.

# Today's Lecture

---

- Supervised Learning Setup
- Linear Regression
- Logistic Regression ← *binary classification*
- Computation Graph

# Linear Regression

$$\begin{matrix} \underline{x}^{(i)} \rightarrow \underline{y}^{(i)} \in R \\ \in R^n \end{matrix}$$

$$f_w(\underline{x}) = \underline{w}_0 + \underline{w}_1 x_1^{(i)} + \dots + \underline{w}_n x_n^{(i)}$$

parameters (weights)

$$= \underline{w}_0 + \sum_{i=1}^n w_i \underline{x}_i \leftarrow \gamma R^{n+1}$$

$$\underline{x}_0 = 1 \text{ (intercept term)} \quad = \sum_{i=0}^n w_i \underline{x}_i = \underline{w}^T \underline{x}$$

# Linear Regression

$$\{(x^{(1)}, y^{(1)}), (x^{(2)}, y^{(2)}) \dots (x^{(m)}, y^{(m)})\}$$

$$\Rightarrow \hat{y} = \underline{f(x)} \approx y$$

$\Rightarrow$  Cost function:

$$J(\underline{\underline{w}}) = \frac{1}{2} \sum_{i=1}^m (\underline{f_w(x^{(i)})} - y^{(i)})^2$$

ordinary least  
square

# Linear Regression

How to choose  $\underline{w}$  minimize  $J(w)$ ?

Gradient descent

$\Rightarrow$  initial guess  $\underline{w}$ .

$$\Rightarrow \underline{w} \leftarrow \underline{w} - \beta \nabla_{\underline{w}} J(\underline{w})$$

\* learning rate / step size,

$$w, x \in \mathbb{R}^{n+1} \quad \underline{w}_j \leftarrow \underline{w}_j - \beta \frac{\partial}{\partial w_j} J(\underline{w}), \quad j=0, \dots, n$$

learning

rate

## Linear Regression

Suppose we only have 1 sample

$$\underline{J(w)} = \frac{1}{2} (f(x) - y)^2 = \frac{1}{2} (\underline{w^T x - y})^2$$

$$\frac{\partial}{\partial w_j} J(w) = \frac{\partial J(w)}{\partial z} \cdot \frac{\partial z}{\partial w_j}$$

$\sum_{i=0}^n w_i x_i - y$   
 $w_j x_j$

Chain rule.

$$= (\underbrace{w^T x - y}_{f(x)}) x_j$$

$$w_j \leftarrow w_j - \beta (\underline{f(x) - y}) x_j$$

# Linear Regression

$$\underline{w_j} \leftarrow \underline{w_j} - \beta \sum_{i=1}^m (f(x^{(i)}) - y^{(i)}) x_j^{(i)}$$

\* Why we use the squared loss in  $J(w)$ ?

$$\Rightarrow \underline{\underline{y}^{(i)}} = \underline{\underline{\theta}^T \underline{x}^{(i)}} + (\underline{\varepsilon}^{(i)}) \text{ error term}$$

# Linear Regression

---

$$\underline{\varepsilon^{(i)} \sim N(0, \sigma^2)}, \quad \forall i=1, \dots, m$$

independently, identically distributed.

$$P(\varepsilon^{(i)}) = \frac{1}{\sqrt{2\pi}\sigma} \exp\left(-\frac{(\varepsilon^{(i)})^2}{2\sigma^2}\right)$$

$$\Rightarrow y^{(i)} | \underline{x^{(i)}}; \underline{\theta} \sim N(\theta^T \underline{x^{(i)}}, \sigma^2)$$

# Linear Regression

Given  $\theta$ , and  $m$  data points  $x^{(1)} \dots x^{(m)}$   
what is the joint distribution of  
 $y^{(1)} \dots y^{(m)}$ ?

$$\begin{aligned} L(\theta) &= \prod_{i=1}^m P(y^{(i)} | x^{(i)}, \theta) \\ &= \prod_{i=1}^m \frac{1}{\sqrt{2\pi}\sigma} \exp\left(-\frac{(y^{(i)} - \underline{\theta}^T x^{(i)})^2}{2\sigma^2}\right) \end{aligned}$$

# Linear Regression

---

Find  $\theta$  that maximize the likelihood  
that you see this dataset.

$$\max_{\theta} L(\theta)$$

$$\begin{aligned}\ell(\theta) &= \log L(\theta) \\ &= \sum_{i=1}^m \log \left[ \frac{1}{\sqrt{2\pi}\sigma} \exp \left( -\frac{(y^{(i)} - \theta^T x^{(i)})^2}{2\sigma^2} \right) \right]\end{aligned}$$

# Linear Regression

(continued)

$$= \sum_{i=1}^m \left[ \log \frac{1}{\sqrt{2\lambda}\sigma} - \frac{(y^{(i)} - \theta^T x^{(i)})^2}{2\sigma^2} \right]$$

$$= m \log \frac{1}{\sqrt{2\lambda}\sigma} - \frac{1}{2\sigma^2} \sum_{i=1}^m (y^{(i)} - \theta^T x^{(i)})^2$$

$$\max_{\theta} L(\theta) \iff \min_{\theta} \sum_{i=1}^m (y^{(i)} - \theta^T x^{(i)})^2$$

# Today's Lecture

---

- Supervised Learning Setup
- Linear Regression
- Logistic Regression
- Computation Graph

# Logistic Regression

y take value from 0, 1.  
not  $\frac{1}{1 + e^{-w^T x}}$   $\rightarrow$  spam email

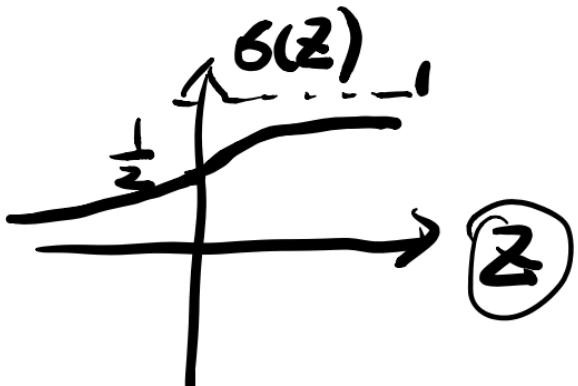
<1> Can we use  $\hat{y} \leftarrow \underline{w^T x}$  • can be larger  
use linear than 1 or  
regression model? smaller than 0

# Logistic Regression

$$(2) \hat{y}^{(i)} = \underline{\sigma}(w^T x^{(i)})$$

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

↙ sigmoid function  
↙ logistic function.



$$z \rightarrow +\infty \quad \sigma(z) \rightarrow 1$$

$$z \rightarrow -\infty \quad \sigma(z) \rightarrow 0$$

# Logistic Regression

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

$$h(x) = \frac{f(x)}{g(x)}$$
$$h'(x) = \frac{g(x)f'(x) - f(x)g'(x)}{g(x)^2}$$

$$\begin{aligned}\sigma'(z) &= \frac{0 - (-e^{-z})}{(1+e^{-z})^2} \\ &= \frac{e^{-z}}{(1+e^{-z})^2} = \frac{1}{1+e^{-z}} \cdot \frac{e^{-z}}{1+e^{-z}} \\ &\quad \text{circled terms: } \frac{1}{1+e^{-z}}, \frac{e^{-z}}{1+e^{-z}}\end{aligned}$$

# Logistic Regression

---

$$\hat{y} = \sigma(w^T \underline{x}) = f(x) \in [0, 1]$$

$$\begin{cases} P(y=1 | x) = \sigma(w^T x) = f(x) \\ P(y=0 | x) = 1 - \sigma(w^T x) = 1 - f(x) \end{cases}$$

$$\underline{P(y|x)} = \underline{[f(x)]^y} \underline{[1-f(x)]^{1-y}}$$

# Logistic Regression

---

$$\begin{aligned} L(w) &= \prod_{i=1}^m \underline{P(y^{(i)} | x^{(i)})} \\ &= \underline{\frac{\prod_{i=1}^m f(x^{(i)})^{y^{(i)}} [1 - f(x)]^{1-y^{(i)}}}{\prod_{i=1}^m}} \end{aligned}$$

$$\underline{l(w)} = \underline{\log L(w)} = \sum_{i=1}^m \left( \log \frac{\underline{f(x^{(i)})^{y^{(i)}}}}{\underline{[1-f(x^{(i)})]^{1-y^{(i)}}}} \right)$$

# Logistic Regression

$$= \sum_{i=1}^m [y^{(i)} \log f(x^{(i)}) + (1-y^{(i)}) \log (1-f(x^{(i)}))]$$

cross-entropy loss

$$\mathbf{w} \leftarrow \mathbf{w} + \frac{\beta}{\text{learning}} \nabla_{\mathbf{w}} L(\mathbf{w})$$

gradient

# Logistic Regression

---

$$w_j \leftarrow w_j + \beta \cdot \frac{\partial L(w)}{\partial w_j}$$

Chain rule

$$\frac{\partial L(w)}{\partial w_j} = \frac{\partial L(w)}{\partial f} \cdot \frac{\partial f}{\partial z} \cdot \frac{\partial z}{\partial w_j}$$

## Logistic Regression

(continued)  $= \sum_{i=1}^m \left( \underbrace{y^{(i)} \frac{1}{f(x^{(i)})}}_{\delta(z) \cdot (1 - \delta(z))} - (1 - y^{(i)}) \frac{1}{1 - f(x^{(i)})} \right)$

$$\underline{\delta(z) \cdot (1 - \delta(z))} x_j^{(i)}$$

$$= \sum_{i=1}^m y^{(i)} [1 - f(x^{(i)})] \\ - (1 - y^{(i)}) f(x^{(i)}) x_j^{(i)}$$

# Logistic Regression

$$\underset{\text{+}}{\cancel{-}} \sum_{i=1}^m \left[ y^{(i)} - f(x^{(i)}) \right] x_j^{(i)}$$

linear regression,

$$f(x^{(i)}) = w^\top x^{(i)}$$

Logistic regression

$$f(x^{(i)}) = \sigma(w^\top x^{(i)})$$

# Today's Lecture

---

- Supervised Learning Setup
- Linear Regression
- Logistic Regression
- Computation Graph

# Computation Graph

$$J(a,b,c) = 2(a + bc)$$

$$= 2 \times (1 + 6) = 14$$

$$\begin{array}{l} a=1 \\ b=2 \\ c=3 \end{array}$$

$$\begin{aligned} u &= bc - \\ \boxed{v = a + bc} &= a + u \\ J &= 2v \end{aligned}$$

$$\begin{array}{c} a=1 \\ b=2 \\ c=3 \end{array} \quad \boxed{u = bc} \quad \boxed{\frac{dv}{dV} = 2} \quad \boxed{V = a + u} \quad \boxed{J = 2V}$$

•  $\frac{du}{du} = 1$

$\frac{du}{du} = \frac{\partial I}{\partial u} \cdot \frac{\partial u}{\partial u} = 1 \cdot 1 = 1$

•  $\frac{db}{du} = \frac{\partial I}{\partial b} = \frac{\partial I}{\partial u} \cdot \frac{\partial u}{\partial b} = 2 \cdot c = 2 \cdot 3 = 6$