

# Linear Regression

## Gradient Descent

- **Model**  $y = w \cdot x + b$ 
  - Model  $\rightarrow$  **training data**  $\rightarrow$  **testing data**  $\rightarrow$  **Overfitting**
- **Loss**  $L(w, b) = \sum_{i=1}^n (y_i - (w \cdot x_i + b))^2$
- **Minimization**
  - $w_1 = w_0 - \alpha \frac{\partial L}{\partial w}|_{w=w_0}$
  - $b_1 = b_0 - \alpha \frac{\partial L}{\partial b}|_{b=b_0}$
  - $\alpha$  **Learning rate**

## Regularization

- **Loss function**
  - $L(w, b) = \sum_{i=1}^n (y_i - (w \cdot x_i + b))^2 + \lambda \sum (w_i)^2$
  - $\lambda$  **Loss function**
    - $\lambda$   $\rightarrow$  0  $\rightarrow$  **Underfitting**
    - $\lambda$   $\rightarrow$  **Overfitting**

# Classification

- **1. 機械学習**
    - **条件付き確率**
    - $P(y=k|x)$   $\rightarrow$  **条件付き確率**
    - **全確率の定理**
- $$P(y=k|x) = \frac{P(x | y = k) P(y = k)}{P(x)}$$
- ここで
- $P(y=k|x) \propto P(x | y = k) \propto P(y = k)$
  - $P(x | y=k) \propto P(y=k | x) \propto P(x)$
  - $P(y=k) \propto 1$

- 
- **2. ロジスティック回帰**
    - $P(y=k|x)$   $\rightarrow$   $P(x|y=k)$   $\rightarrow$   $P(y=k)$   $\rightarrow$  **条件付き確率**
    - **条件付き確率**
    - **条件付き確率**

- 
- **3. ロジスティック**
    1.  $P(y=k)$
    2.  $P(x|y=k)$ 
      - $x \cdot P(y=k)$

3. \*\* $P(y=k|x)$ \*\*

- **MLE**  $\mu, \Sigma$   $\hat{\mu}, \hat{\Sigma}$ 
    - $L(\mu, \Sigma) = f_{\mu, \Sigma}(x^1)f_{\mu, \Sigma}(x^2)f_{\mu, \Sigma}(x^3)f_{\mu, \Sigma}(x^4)\dots f_{\mu, \Sigma}(x^{79})$
    - $L(\mu, \Sigma) \propto \mu_1, \mu_2, \Sigma$
    - $\mu_1, \mu_2, \Sigma$   $\hat{\mu}, \hat{\Sigma}$
  - $P(y=1|x) \propto P(y=2|x)$

**5.**

- $$\frac{P(C_1|x)}{1+P(C_2|x)} = \frac{P(x|C_1)P(C_1)}{P(x|C_1)P(C_1) + P(x|C_2)P(C_2)} = \frac{1}{1 + e^{-z}}$$
    - $\sigma(z)$  sigmoid function
  - $z = wx + b$ 
    - $w^T = (\mu^1 - \mu^2)\Sigma^{-1}$
    - $b = -\frac{1}{2}(\mu^1)^T\Sigma^{-1}\mu^1 + \frac{1}{2}(\mu^2)^T\Sigma^{-1}\mu^2 + \ln\frac{N_1}{N_2}$ 
      - $\mu^1 = \Sigma^{-1}\mu^1$
      - $\Sigma^{-1} = \Sigma^{-1}$

# Logistic Regression

## Loss function

$\text{f}_{\{w,b\}}(X) = P_{\{w,b\}}(C=1|x)$

- 

□□□□	$x^1$	$x^2$	$x^3$	$\dots$
□□	$C_1$	$C_1$	$C_2$	$\dots$
□□	$\hat{y} = 1$	$\hat{y} = 1$	$\hat{y} = 0$	$\dots$

- $\text{f}(w, b) = w_1 x_1 + w_2 x_2 + \dots + w_N x_N + b$
  - $f(w, b)(x) = P(w, b)(C_1 | x)$

██████████████████ \$C\_1\$ □ \$C\_2\$ ██████████ \$C\_2\$ ██████████

$$P_{\{w,b\}}(C=2|x) = 1 - P_{\{w,b\}}(C=1|x) = 1 - f_{\{w,b\}}(x)$$

•  $\text{argmax}^{\star}(w, b) \equiv \text{argmax}(w, b) \equiv \text{arg}$

- $\$-\ln L(w, b) = \lnf_{\{w,b\}}(x^1) + \lnf_{\{w,b\}}(x^2) + \ln(1 - f_{\{w,b\}}(x^3)) = \sum_n -[\hat{y}^n \lnf_{\{w,b\}}(x^n) + (1 - \hat{y}^n) \ln(1 - f_{\{w,b\}}(x^n))]$

- $C(f(x^n), \hat{y}^n) = \sum_n -[\hat{y}^n \ln f_{w,b}(x^n) + (1 - \hat{y}^n) \ln(1 - f_{w,b}(x^n))]$

## Gradient Descent

**Cross entropy**  $C(f(x^n), \hat{y}^n) = \sum_n -[\hat{y}^n \ln f_{w,b}(x^n) + (1 - \hat{y}^n) \ln(1 - f_{w,b}(x^n))]$

- $\frac{\partial}{\partial w_i} \frac{\partial \ln(1 - \sigma(z))}{\partial w_i} = \frac{\partial \ln(1 - \sigma(z))}{\partial z} \frac{\partial z}{\partial w_i} = -\sigma(z)x_i$ 
  - $\frac{\partial}{\partial z} \frac{\partial \ln(1 - \sigma(z))}{\partial z} = -\frac{1}{1 - \sigma(z)}\sigma'(z)$
  - $\frac{\partial}{\partial z} \frac{\partial z}{\partial w_i} = x_i$
- $\frac{\partial}{\partial b} \frac{\partial \ln(1 - \sigma(z))}{\partial b} = \frac{\partial \ln(1 - \sigma(z))}{\partial z} \frac{\partial z}{\partial b} = -(1 - \sigma(z))x_i$
- $\frac{\partial}{\partial w_i} \frac{\partial \ln(-\ln(w, b))}{\partial w_i} = \sum_n -(\hat{y}^n - f_{w,b}(x^n))x_i^n$ 
  - $\frac{\partial}{\partial w_i} \frac{\partial \ln(-\ln(w, b))}{\partial w_i} = \sum_n -(\hat{y}^n - f_{w,b}(x^n))x_i^n$

## Square Error

- **Loss function**  $L(f) = \frac{1}{2} \sum_n (f_{w,b}(x^n) - \hat{y}^n)^2$
- $\frac{\partial}{\partial w_i} \frac{\partial L(f)}{\partial w_i} = \text{normalize } 2(f_{w,b}(x) - \hat{y})f_{w,b}(x)(1 - f_{w,b}(x))x_i$ 
  - $\hat{y} = 0 \Rightarrow f_{w,b}(x) = 1 \Rightarrow 2(f_{w,b}(x) - \hat{y})f_{w,b}(x)(1 - f_{w,b}(x))x_i = 0$
  - $\hat{y} = 1 \Rightarrow f_{w,b}(x) = 0 \Rightarrow 2(f_{w,b}(x) - \hat{y})f_{w,b}(x)(1 - f_{w,b}(x))x_i = 0$

## Multi-class Classification

### weight · bias

- $C_1: w^1, b_1 \quad z_1 = w^1 \cdot x + b_1$
  - $C_2: w^2, b_2 \quad z_2 = w^2 \cdot x + b_2$
  - $C_3: w^3, b_3 \quad z_3 = w^3 \cdot x + b_3$
1. softmax ဆိတ်ချက်
    - $f(z) = \frac{e^z}{\sum_{i=1}^n e^{z_i}}$
    - $f(z) \in [0, 1]^n$
  2. Cross entropy ဆိတ်ချက်
    - $-\sum_{i=1}^n \hat{y}_i \ln y_i$

## Limitation of Logistic Regression

- အမြန် ရှင်းခွင့်
- $\sigma(z) = 0.5 \Rightarrow w \cdot x + b = 0$

### Feature transformation

## Discriminative and Generative

### Discriminative

• ပုံစံမျက်နှာတွင်  $x$  အတွက်

- $P(y|x)$

- ပုံစံမျက်နှာတွင်  $y$  အတွက်

## Generative ပုံစံမျက်နှာ

• ပုံစံမျက်နှာတွင်  $P(x|y)$

- ပုံစံမျက်နှာတွင်  $x$  အတွက်  $w \cdot b$

- ပုံစံမျက်နှာတွင်  $y$  အတွက်

## Deep Learning

### Three steps for Deep Learning

#### 1. Define a set of function

- ပုံစံမျက်နှာတွင်

- ပုံစံ

- ပုံစံမျက်နှာတွင်  $\sigma$  အတွက်

- ပုံစံမျက်နှာတွင် ReLU・Sigmoid・Tanh အဲ

- ပုံစံမျက်နှာတွင်  $\text{Adam}$  အတွက်

- ပုံစံမျက်နှာတွင်  $\text{SGD}$  အတွက်

ပုံစံ

- ပုံစံမျက်နှာတွင် Sigmoid・ReLU・ReLU・ReLU ပုံစံမျက်နှာ

#### Sigmoid ပုံစံ

- $\sigma(x) = (0, 1)$

- ပုံစံမျက်နှာတွင် 0 အတွက်

- ပုံစံမျက်နှာတွင်

#### ReLU ပုံစံ

- $\text{ReLU}(x) = [0, \infty)$

- ပုံစံမျက်နှာတွင် sigmoid ပုံစံမျက်နှာ

- ပုံစံမျက်နှာတွင်

#### □ ReLU ပုံစံ

- ReLU ပုံစံမျက်နှာတွင် 0 အတွက်

#### 2. Goodness of function

- ပုံစံမျက်နှာတွင်

- ပုံစံမျက်နှာတွင်  $\sigma$  အတွက်

- ပုံစံမျက်နှာတွင်  $k$  အတွက်

- ပုံစံမျက်နှာတွင်  $\text{ReLU}$  အတွက်

- ပုံစံမျက်နှာတွင်  $\text{SGD}$  အတွက်

### 3. Pick the best function

- $\text{MSE} = \frac{1}{n} \sum_{i=1}^n (y_i - \hat{y}_i)^2$
- $\text{MSE}$  ကို  $k$  ပုံမှန်လေ့လာတဲ့ အကြောင်း
- $\text{MSE}$  ကို  $k$  ပုံမှန်လေ့လာတဲ့ အကြောင်း
- $\text{MSE}$  ကို  $k$  ပုံမှန်လေ့လာတဲ့ အကြောင်း

## Gradient Descent

- Learning Rate $\alpha$ 
  - $\alpha$  မြင်
  - $\alpha$  မြတ်

ဗိုလ်တော် Taylor Series

## Adagrad

ဗိုလ်တော်

## Stochastic Gradient Descent

$L = \sum_n (\hat{y}^n - (b + \sum w_i x_i^n))^2$

- Gradient Descent $\theta^i = \theta^{i-1} - \eta \nabla L(\theta^{i-1})$
- Stochastic Gradient Descent $\theta^i = \theta^{i-1} - \eta \nabla L(\theta^{i-1})$  example
  - ဗိုလ်တော်

## Feature Scaling

ဗိုလ်တော်

- $x_i \leftarrow \frac{x_i - m_i}{\sigma_i}$
  - $x_i \leftarrow \frac{x_i - m_i}{\sigma_i}$
- ဗိုလ်
- $$x_i \leftarrow \frac{x_i - m_i}{\sigma_i}$$
- $m_i$ \$ မြတ်  
◦  $\sigma_i$ \$ မြတ်
  - 0 မြတ် 1

## Backpropagation

- Forward pass Backword pass
  - ဗိုလ်တော်
  - ဗိုလ်တော်
  - ဗိုလ်တော်



## Forward pass

## Backword pass

## BERT

BERT

## Unsupervised Learning

### Self-supervised Learning

#### Word Embedding

Word Embedding

Word

##### 1. One-hot encoding

- 一个词的独热编码\*\**(one-hot encoding)*\*\*是将这个词表示为一个向量(维度等于词汇量)
- 向量维度通常在几百到几千之间(约50-300)

##### 2. Word2Vec

- 通过上下文信息学习词向量
- 通过语义相似度学习词向量

### Prediction-based

#### Prediction-based

预测

预测

##### 1. 通过独热编码进行one-hot encoding

- "狗" "猫" "鸟" "狗" "鸟" "狗" "鸟" "狗" "鸟"

##### 2. 通过矩阵乘法进行

- " $\begin{pmatrix} \text{狗} \\ \text{猫} \end{pmatrix} \times \begin{pmatrix} \text{狗} & \text{猫} & \text{鸟} \end{pmatrix} \rightarrow \begin{pmatrix} \text{狗} \\ \text{猫} \end{pmatrix}$ "

##### • " $\begin{pmatrix} \text{狗} \\ \text{猫} \end{pmatrix} \times \begin{pmatrix} \text{狗} & \text{猫} & \text{鸟} \end{pmatrix} \rightarrow \begin{pmatrix} \text{狗} \\ \text{猫} \end{pmatrix}$ "

- 狗的向量与狗的向量相似

- 狗的向量与鸟的向量不相似

- 猫的向量与狗的向量不相似

##### 4. 通过神经网络进行

- 狗的向量与狗的向量相似

- ပုဂ္ဂန်များ

## 5. အောက်တို့တော်းခြင်း

- ပုဂ္ဂန်များ
- ပုဂ္ဂန်များ

အောက်တို့တော်းခြင်း

### 1. အောက်တို့တော်းခြင်း

- ပုဂ္ဂန်များ
- ပုဂ္ဂန်များ

### 2. အောက်တို့တော်းခြင်း

- ပုဂ္ဂန်များ
- ပုဂ္ဂန်များ

## Seq2Seq

**Sequence to Sequence Model**အောက်တို့တော်းခြင်း

အောက်

Seq2Seqအောက်တို့တော်းခြင်း

1. အောက်Encoderအောက်တို့တော်းခြင်း
2. အောက်Decoderအောက်တို့တော်းခြင်း

**AT(Auto-regressive) VS NAT(Non-Autoregressive Transformer)**အောက်

အောက်**AR/AT**

- အောက်Encoderအောက်တို့တော်းခြင်း
- အောက်Decoderအောက်တို့တော်းခြင်း
- အောက်GPTအောက်Transformerအောက်

အောက်**NAT**

- အောက်Encoderအောက်တို့တော်းခြင်း
- အောက်Decoderအောက်တို့တော်းခြင်း
- အောက်Transformerအောက်
- အောက်GPTအောက်