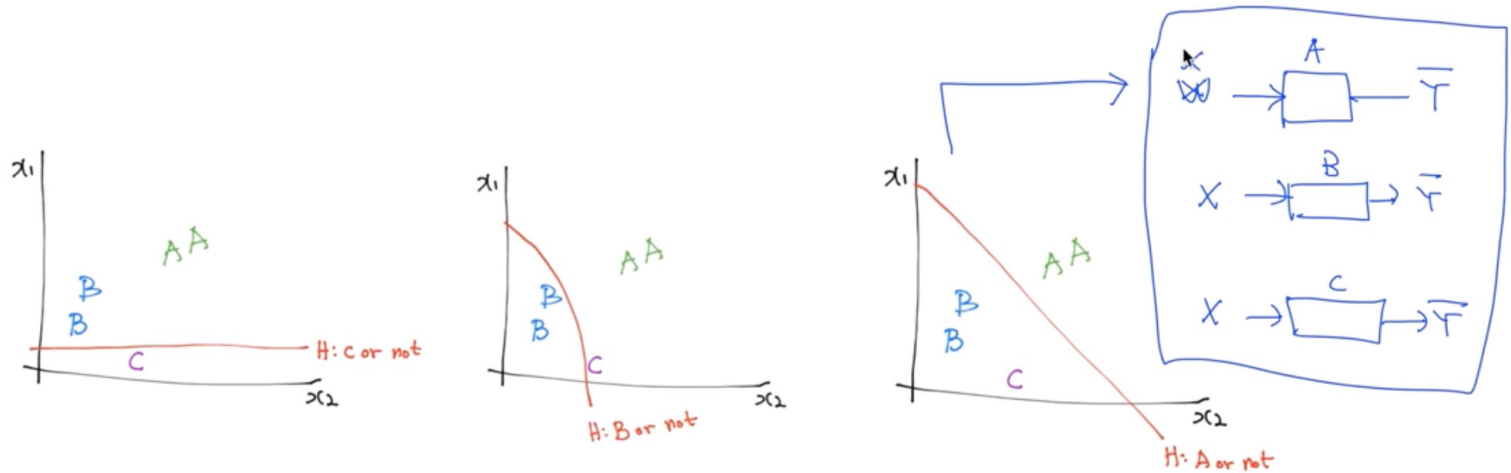# 2주차 머신러닝 발표

김찬원

# Multinomial classification

# Multinomial classification

$$\begin{bmatrix} w_1 & w_2 & w_3 \end{bmatrix} \begin{bmatrix} x_1 \\ x_2 \\ x_3 \end{bmatrix} = \begin{bmatrix} w_1 x_1 + w_2 x_2 + w_3 x_3 \end{bmatrix}$$

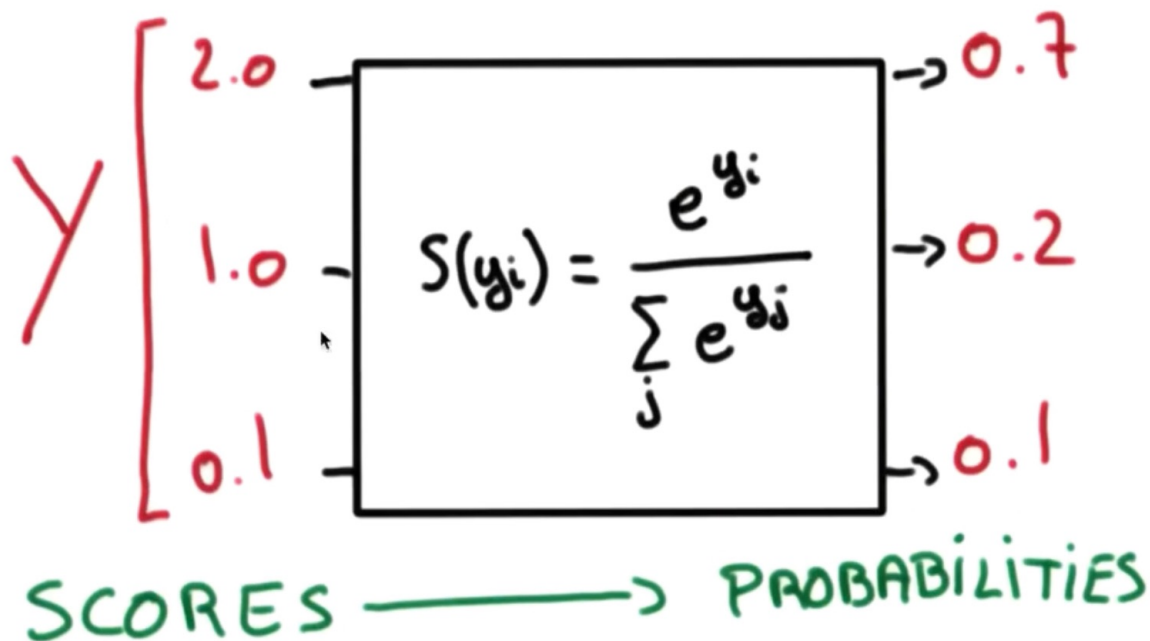$$\begin{bmatrix} w_1 & w_2 & w_3 \end{bmatrix} \begin{bmatrix} x_1 \\ x_2 \\ x_3 \end{bmatrix} = \begin{bmatrix} w_1 x_1 + w_2 x_2 + w_3 x_3 \end{bmatrix}$$

$$\begin{bmatrix} w_1 & w_2 & w_3 \end{bmatrix} \begin{bmatrix} x_1 \\ x_2 \\ x_3 \end{bmatrix} = \begin{bmatrix} w_1 x_1 + w_2 x_2 + w_3 x_3 \end{bmatrix}$$



$$\begin{bmatrix} w_{A1} & w_{A2} & w_{A3} \\ w_{B1} & w_{B2} & w_{B3} \\ w_{C1} & w_{C2} & w_{C3} \end{bmatrix} \begin{bmatrix} x_1 \\ x_2 \\ x_3 \end{bmatrix} = \begin{bmatrix} w_{A1} x_1 + w_{A2} x_2 + w_{A3} x_3 \\ w_{B1} x_1 + w_{B2} x_2 + w_{B3} x_3 \\ w_{C1} x_1 + w_{C2} x_2 + w_{C3} x_3 \end{bmatrix} = \begin{bmatrix} \bar{y}_A \\ \bar{y}_B \\ \bar{y}_C \end{bmatrix} \begin{bmatrix} 2.0 \\ 1.0 \\ 0.1 \end{bmatrix}$$

# SOFTMAX



$Y \begin{bmatrix} 2.0 \\ 1.0 \\ 0.1 \end{bmatrix}$

$$S(y_i) = \frac{e^{y_i}}{\sum_j e^{y_j}}$$

$\rightarrow 0.7$

$\rightarrow 0.2$

$\rightarrow 0.1$

SCORES $\longrightarrow$ PROBABILITIES

# Cross-entropy cost function

$D(S,L)$

$= -\sum_i L_i \log(S_i)$

$-\sum_i L_i \log(\bar{y}_i) = \sum_i L_i * -\log(\bar{y}_i)$


$-\log(x)$

$Y = L = \begin{bmatrix} 0 \\ 1 \end{bmatrix} = \underline{B}$

$-\log$

$Y = \begin{bmatrix} 0 \\ 1 \end{bmatrix} \quad B \, (OK), \quad \begin{bmatrix} 0 \\ 1 \end{bmatrix} \odot -\log \begin{bmatrix} 0 \\ 1 \end{bmatrix} = \begin{bmatrix} 0 \\ 1 \end{bmatrix} \odot \begin{bmatrix} 00 \\ 0 \end{bmatrix} = \begin{bmatrix} 0 \\ 0 \end{bmatrix} \Rightarrow 0,$
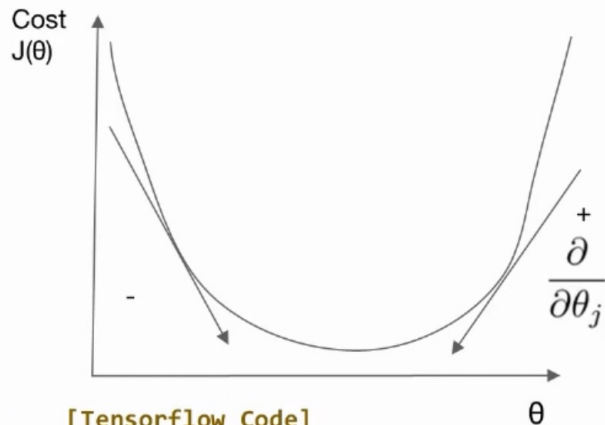
$-\log$

$\bar{Y} = \begin{bmatrix} 1 \\ 0 \end{bmatrix} = A \, (X), \quad \begin{bmatrix} 0 \\ 1 \end{bmatrix} \odot -\log \begin{bmatrix} 1 \\ 0 \end{bmatrix} = \begin{bmatrix} 0 \\ 1 \end{bmatrix} \odot \begin{bmatrix} 0 \\ 00 \end{bmatrix}$

$\frac{1}{8} \text{Loss}: \quad \mathcal{L} = \frac{1}{n} \sum_i D(S(wx_i, b), L_i)$

# Learning rate
## Gradient

Cost
J(θ)

$$\text{Repeat} \left\{ \theta_j := \theta_j - \alpha \frac{\partial}{\partial \theta_j} J(\theta) \right\}$$

+

$$\frac{\partial}{\partial \theta_j}$$

-

Learning rate is a hyper-parameter that controls how much we are adjusting the weights with respect the loss gradient
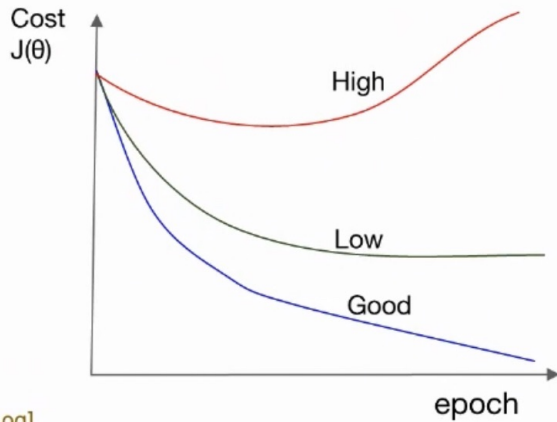
θ

```
[Tensorflow Code]
def grad(hypothesis, labels):
    with tf.GradientTape() as tape:
        loss_value = loss_fn(hypothesis, labels)
    return tape.gradient(loss_value, [W,b])
optimizer = tf.train.GradientDescentOptimizer(learning_rate=0.01)
optimizer.apply_gradients(grads_and_vars=zip(grads,[W,b]))
```
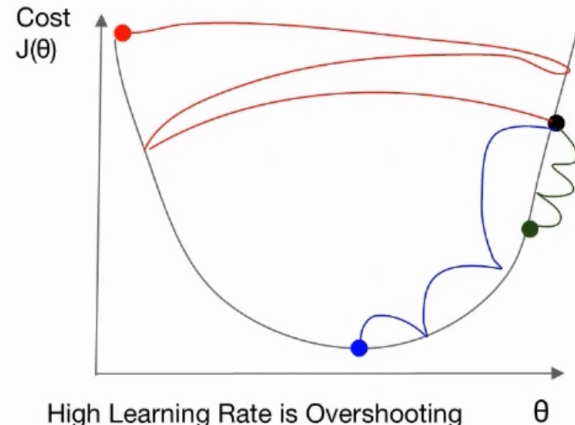
# Learning rate
## Good and Bad

Cost J(θ)

High

Low

Good

epoch

Cost J(θ)

θ

High Learning Rate is Overshooting
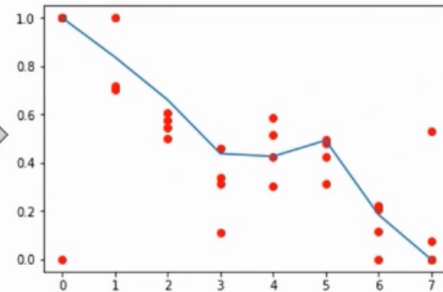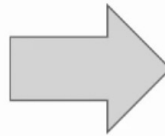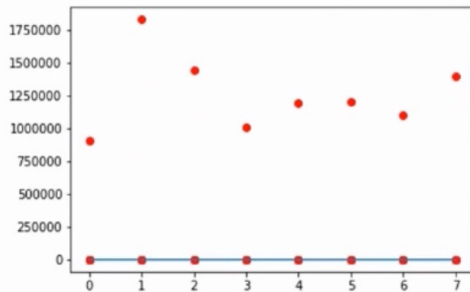Normal Learning Rate is 0.01
3e-4 is the best learning rate for Adam,
hands down (andrej karpathy)

[Tensorflow Code]
```
optimizer = tf.train.GradientDescentOptimizer(learning_rate=0.01)
```

# Data preprocessing
## Feature Scaling



**Standardization**
(Mean Distance)

$$x_{new} = \frac{x - \mu}{\sigma}$$

**Normalization**
(0~1)

$$x_{new} = \frac{x - x_{min}}{x_{max} - x_{min}}$$

```
[Python Code(numpy)]
Standardization = (data - np.mean(data)) / sqrt(np.sum(
(data - np.mean(data))^2 ) / np.count(data))
```

```
Normalization = (data - np.min(data, 0)) / (np.max(data, 0)
- np.min(data, 0))
```
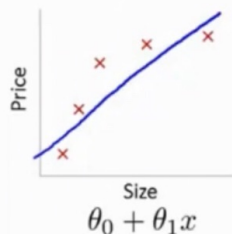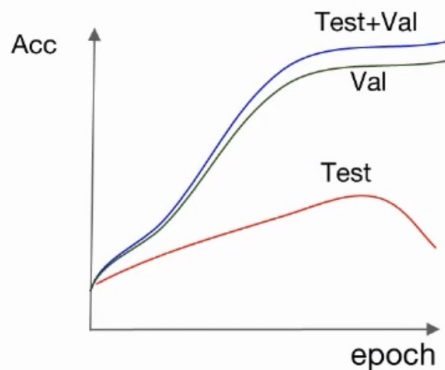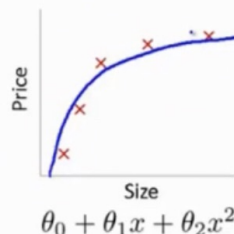
https://github.com/deeplearningzerotoall/TensorFlow/blob/master/lab-07-3-linear_regression_min_max.ipynb
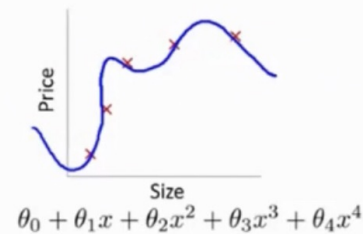
# Overfitting



Acc — Test+Val, Val, Test vs epoch

| | | |
|---|---|---|
| $\theta_0 + \theta_1 x$ | $\theta_0 + \theta_1 x + \theta_2 x^2$ | $\theta_0 + \theta_1 x + \theta_2 x^2 + \theta_3 x^3 + \theta_4 x^4$ |
| High bias (underfit) | "Just right" | High variance (overfit) |

FaceScrub dataset(Aaron Eckhart)

# Overfitting
## Set a features

- Get more training data - more data will actually make a difference, (helps to fix high variance)
- Smaller set of features - dimensionality reduction(PCA) (fixes high variance)
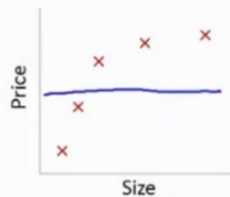- Add additional features - hypothesis is too simple, make hypothesis more specific (fixes high bias)

$$z = U_{reduce}^T x$$

1. $h_\theta(x) = \theta_0 + \theta_1 x$
2. $h_\theta(x) = \theta_0 + \theta_1 x + \theta_2 x^2$
3. $h_\theta(x) = \theta_0 + \theta_1 x + \cdots + \theta_3 x^3$

10. $h_\theta(x) = \theta_0 + \theta_1 x + \cdots + \theta_{10} x^{10}$

$J_{cv}(\theta)$ (cross validation error)

$J_{train}(\theta)$ (training error)
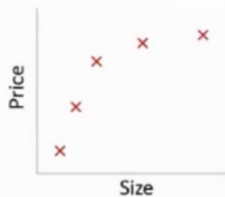
error

degree of polynomial d

[sklearn Code]
```
from sklearn.decomposition import PCA
pca = decomposition.PCA(n_components=3)
pca.fit(X)
X = pca.transform(X)
```
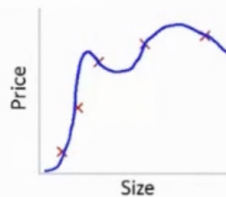
# Overfitting
## Regularization (Add term to loss)



Large $\lambda$
High bias (underfit)
$\lambda = 10000.\ \theta_1 \approx 0, \theta_2 \approx 0, \dots$
$h_\theta(x) \approx \theta_0$

Intermediate $\lambda$
"Just right"

Small $\lambda$
High variance (overfit)
$\lambda \approx 0$

$\lambda$-- : fixes high bias (Under fitting)
$\lambda$ ++ : fixes high variance (overfitting)

**Linear regression with regularization**

Model: $h_\theta(x) = \theta_0 + \theta_1 x + \theta_2 x^2 + \theta_3 x^3 + \theta_4 x^4$

$$J(\theta) = \frac{1}{2m} \sum_{i=1}^{m} (h_\theta(x^{(i)}) - y^{(i)})^2 + \frac{\lambda}{2m} \sum_{j=1}^{m} \theta_j^2$$
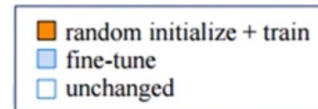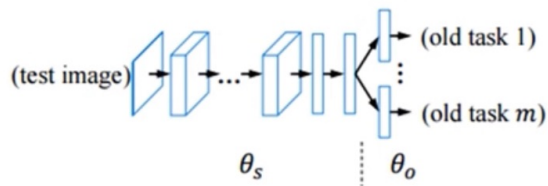
[Tensorflow Code]
```
L2_loss = tf.nn.l2_loss(w) # output = sum(t ** 2) / 2
```
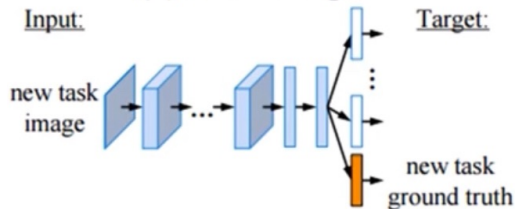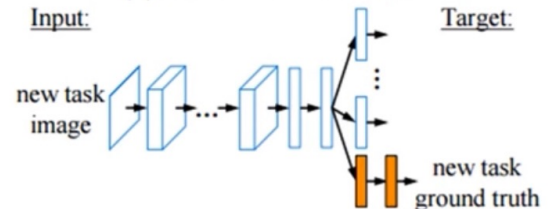
# Learning
## Fine Tuning / Feature Extraction



(a) Original Model

(test image) → ⟶ ⋯ ⟶ → (old task 1) ⋮ (old task $m$)

$\theta_s$   $\theta_o$

■ random initialize + train
□ fine-tune
□ unchanged

(b) Fine-tuning

Input: Target:
new task image → new task ground truth

(c) Feature Extraction

Input: Target:
new task image → new task ground truth

**[Tensorflow Code]**
```
saver = tf.train.import_meta_graph('my-model-1000.meta')
saver.restore(tf.train.latest_checkpoint('./'))
```

Learning without Forgetting : https://arxiv.org/pdf/1606.09282.pdf
Fine tuning : https://goodtogreate.tistory.com/entry/Saving-and-Restoring