

Before & After
The AI Winter

Session #1: Before the AI Winter

Microsoft Student Partners
Evangelist Younjoon Chung

Goals

- Introduction to machine learning
- Basic understanding of machine learning algorithms
 - Linear regression, Logistic regression
 - Gradient Descent Algorithm
 - Sigmoid function
 - Softmax Classifier
 - Overfitting, Regularization

Acknowledgement

- Andrew Ng's ML class
 - <https://class.coursera.org/ml-003/lecture>
 - <http://www.holehouse.org/mlclass> (note)
- 모두의 딥러닝
 - <https://hunkim.github.io/ml>
- The Deep Learning textbook
 - <http://www.deeplearningbook.org>
- 김지환 교수님 System programming
- <http://www.andreykurenkov.com/writing/a-brief-history-of-neural-nets-and-deep-learning>

Machine Learning

- Limitations of explicit programming
 - Spam filter: many rules
 - Automatic driving: too many rules
- Machine learning: "Field of study that gives computers the ability to learn without being explicitly programmed" Arthur Samuel (1959)

What does it mean for machines to learn?

- A computer program is said to learn from experience E with respect to some class of tasks T and performance measure P if its performance at tasks in T , as measured by P , improves with experience E . Tom M. Mitchell (1997)
- E : experience
- T : task
- P : performance measure

The Task, T

- The process of learning itself is not the task.
- Usually described in terms of **how the machine learning system should process a [collection of features]**. ("Deep learning" 99)

What are the tasks?

- **Classification**
 - Learning a function $f: \mathbb{R}^n \rightarrow \{1 \dots k\}$
 - When $y = f(x)$, the model assigns an input described by vector x to a category identified by numeric code y ("Deep learning" 100)
 - Outputs can also be probability distribution over classes.
- Examples
 - Object recognition
 - Spam detector
 - Face recognition -> Viola Jones!
 - MNIST

What are the tasks?

- **Regression**
 - Learning a function $f: \mathbf{R}^n \rightarrow \mathbf{R}$
 - Outputs a numerical value (continuous)
 - Similar to classification; only the format of the output is different
- Examples
 - Weather forecast
 - Predict a test score (if the score is mapped to grades, then it is a classification!)

What are the tasks?

- **Translation**
 - Different from transcription in that the input values are already in a structured format.
 - Machines are asked to convert a sequence of symbols to another language.
- Examples
 - Google Translation

What are the tasks?

- **Anomaly detection**
 - Sifts through a set of events or objects, and flags some of them as being unusual or atypical. ("Deep Learning" 102)
- Examples
 - Credit card fraud detection. By modeling your purchasing habits, a credit card company detect misuse of your cards.

The Performance Measure, P

- The choice of P depends on the nature of the problems.
 - We have to first decide what should be measured.
 - Should we penalize the system more if it frequently makes medium-sized mistakes and rarely makes very large mistakes, and vice versa.
 - how they should be measured
 - 0-1 loss, cross-entropy, MSE, etc.
- Usually we are interested in how well the machine learning algorithm performs on data it has not seen before ("Deep Learning" 104)

The Experience, *E*

- Machine learning algorithms can be broadly categorized as **unsupervised** or **supervised** by what kind of experience they are allowed to have during the learning process. ("Deep Learning" 104)
- Supervised learning algorithms
 - Dataset with labels -> temperature prediction
- Unsupervised learning algorithms
 - Dataset without labels -> hash tag clustering

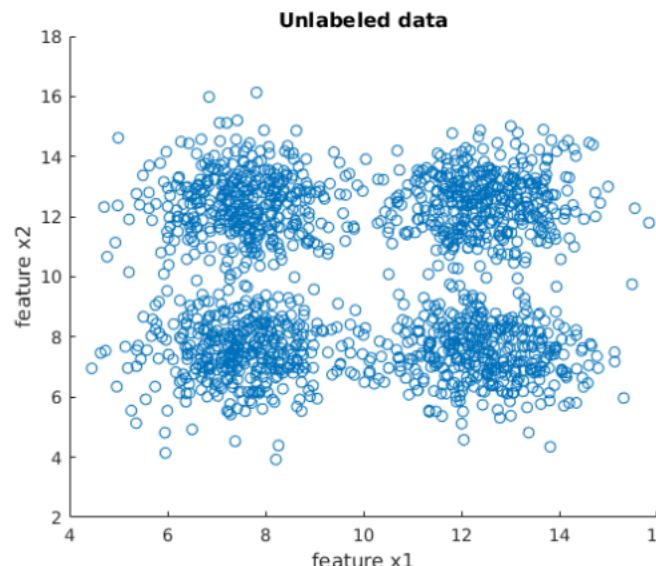
Supervised learning

- Supervised learning
 - learning with labeled examples - training set
- An example training set for four visual categories.

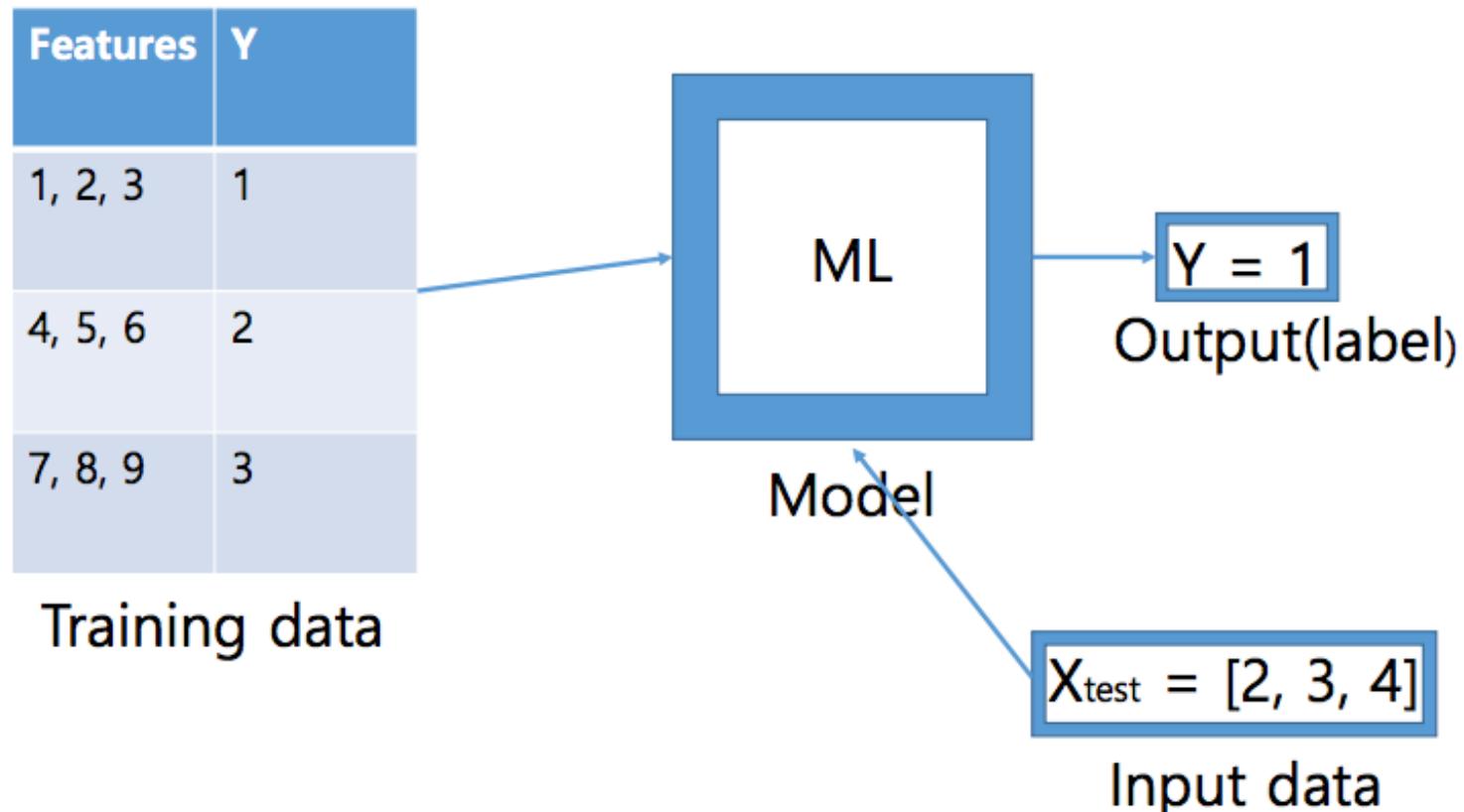


Unsupervised learning

- Unsupervised learning: un-labeled data
 - Google news grouping
 - Word clustering
- An example of clustering



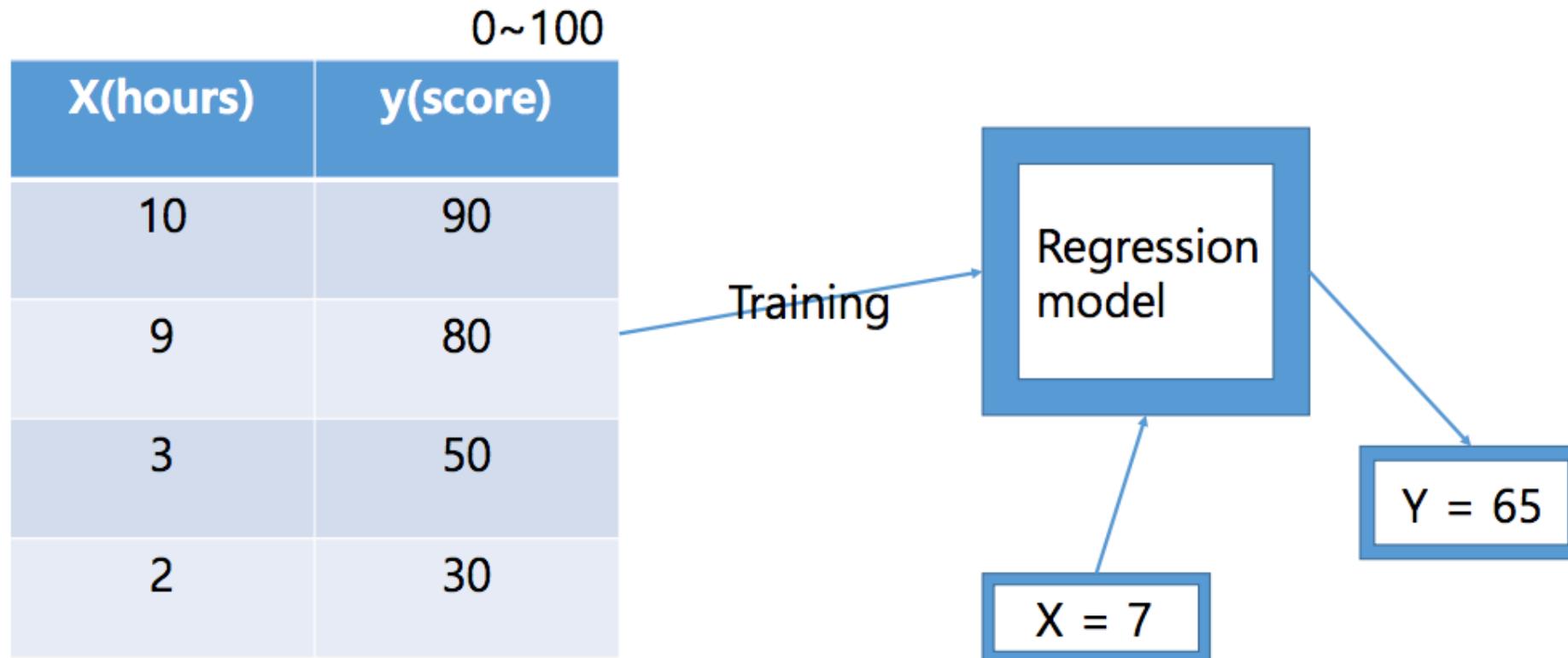
Supervised learning in a nutshell



Types of supervised learning

- Predicting final exam score based on time spent
 - **Regression**
- Pass/non-pass based on time spent
 - **Binary classification**
- Letter grade(A, B, C, D, and F) based on time spent
 - **Multi-label classification**

Predicting exam score: regression

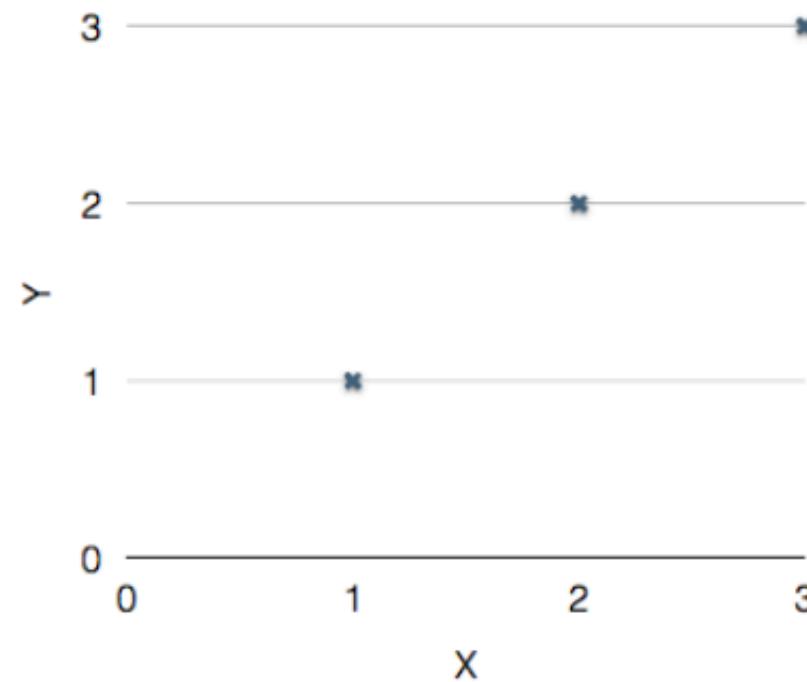


Regression (data)

X	Y
1	1
2	2
3	3

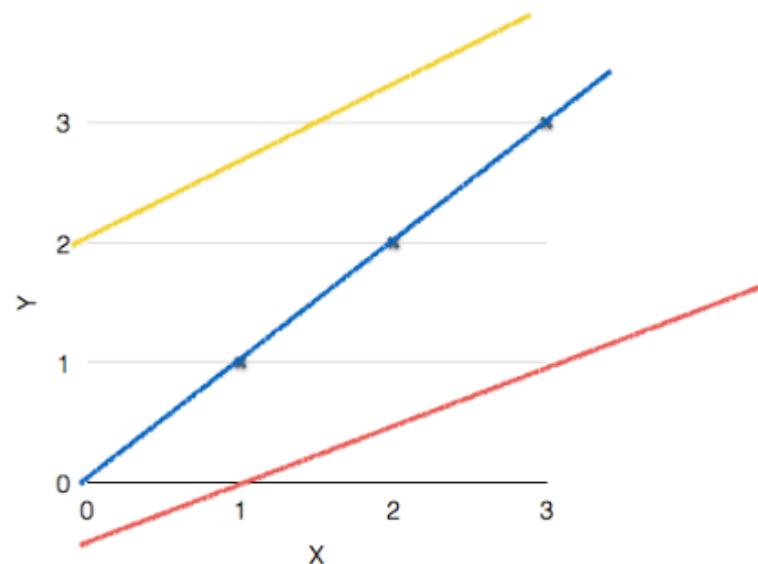
Regression (presentation)

X	Y
1	1
2	2
3	3



(Linear) Hypothesis

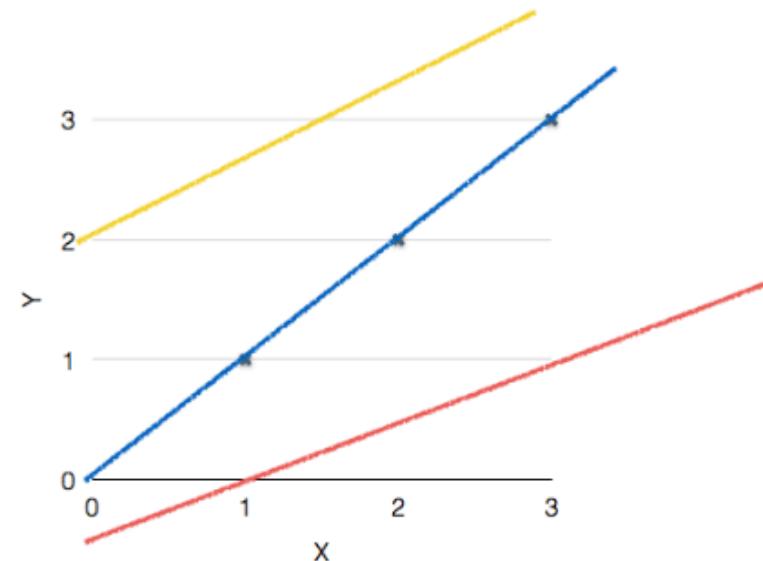
$$H(x) = Wx + b$$



(Linear) Hypothesis

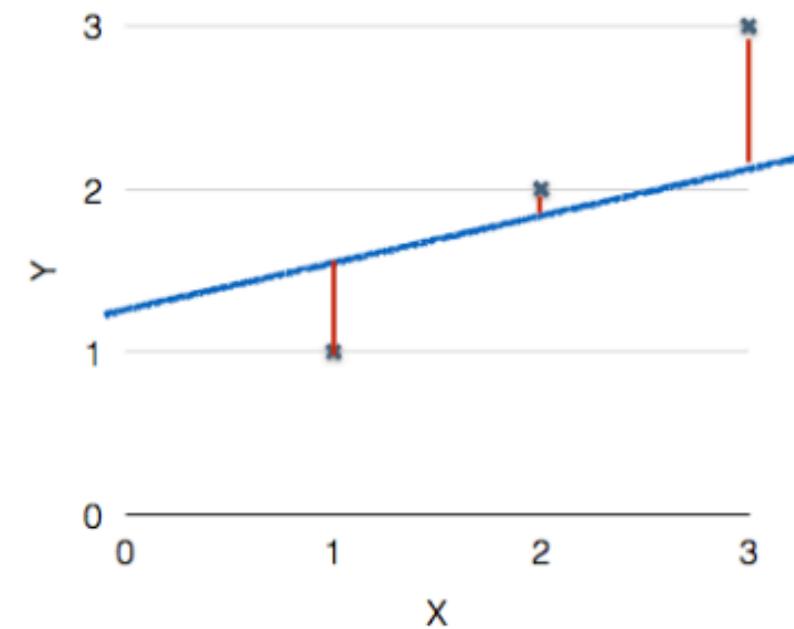
$$H(x) = Wx + b$$

Which hypothesis is better?



Which hypothesis is better?

- We need a measure for deciding which hypothesis is most accurate!
 - > use differences between prediction & real value

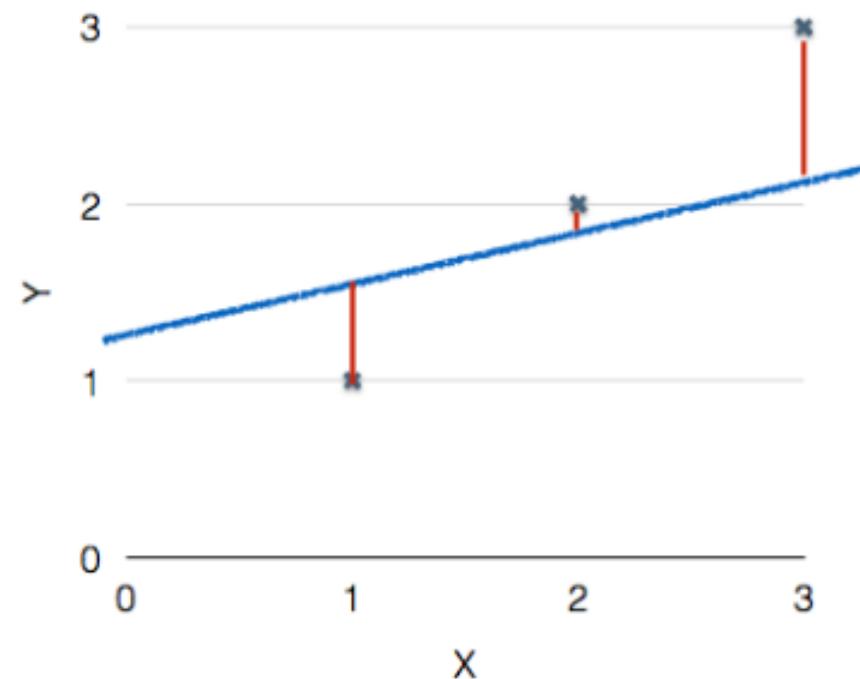


Cost function

- How fit the line to our (training) data

$$\frac{(H(x^{(1)}) - y^{(1)})^2 + (H(x^{(2)}) - y^{(2)})^2 + (H(x^{(3)}) - y^{(3)})^2}{3}$$

$$cost = \frac{1}{m} \sum_{i=1}^m (H(x^{(i)}) - y^{(i)})^2$$



Cost function in a nutshell

- Training Goal:

Cost function: $cost = \frac{1}{m} \sum_{i=1}^m (H(x^{(i)}) - y^{(i)})^2$

$$\underset{W,b}{\text{minimize}} \ cost(W, b)$$

Hypothesis: $H(x) = Wx + b$

Lower the cost,
better the
hypothesis(model)!

$$cost(W, b) = \frac{1}{m} \sum_{i=1}^m (H(x^{(i)}) - y^{(i)})^2$$

Cost function is now a function of
W, b(weight, bias)!

Simplified hypothesis

Eliminated bias for educational purposes

$$H(x) = Wx$$

$$cost(W) = \frac{1}{m} \sum_{i=1}^m (Wx^{(i)} - y^{(i)})^2$$

Now the cost function is a function of only W

What does cost(W) look like?

$$cost(W) = \frac{1}{m} \sum_{i=1}^m (Wx^{(i)} - y^{(i)})^2$$

X	Y
1	1
2	2
3	3

- $W = 1, cost(W) = ?$

What does $\text{cost}(W)$ look like?

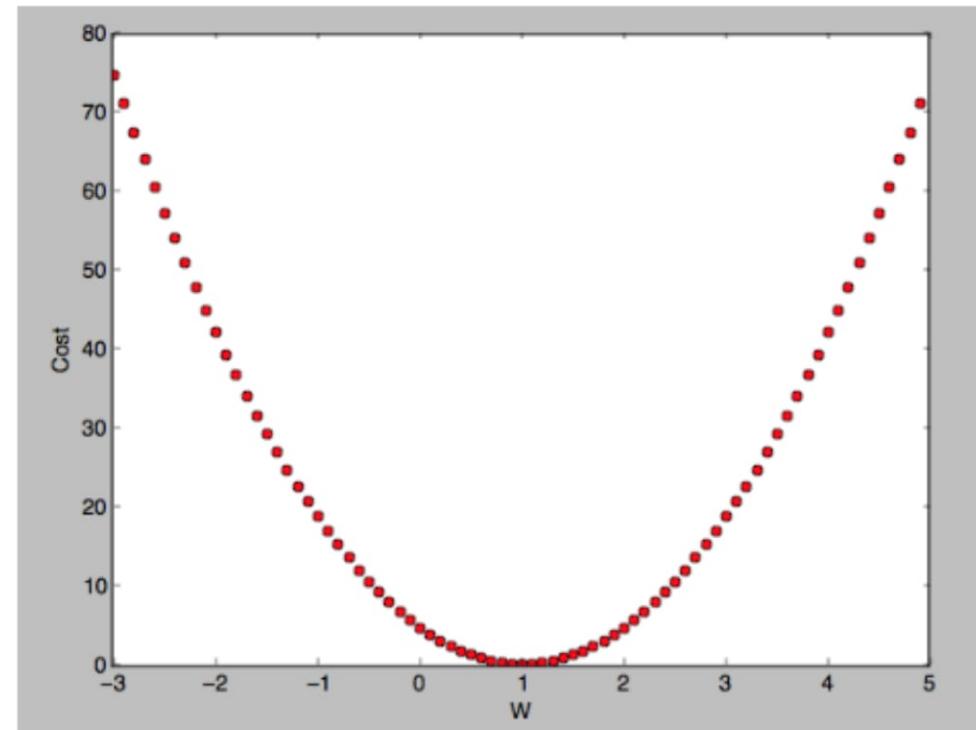
$$\text{cost}(W) = \frac{1}{m} \sum_{i=1}^m (Wx^{(i)} - y^{(i)})^2$$

X	Y
1	1
2	2
3	3

- $W = 1, \text{cost}(W) = 0$
 $\frac{1}{3}((1 * 1 - 1)^2 + (1 * 2 - 2)^2 + (1 * 3 - 3)^2)$
- $W = 0, \text{cost}(W) = 4.67$
- $W = 2, \text{cost}(W) = 4.67$

What does $\text{cost}(W)$ look like?

$$\text{cost}(W) = \frac{1}{m} \sum_{i=1}^m (W x^{(i)} - y^{(i)})^2$$

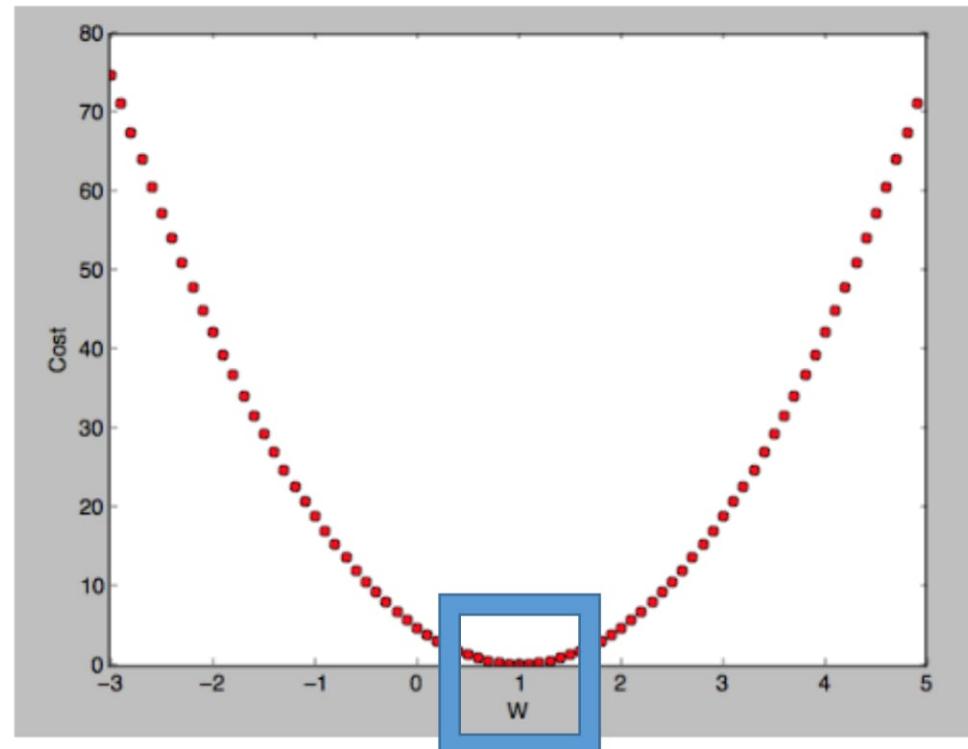


What does $\text{cost}(W)$ look like?

$$\text{cost}(W) = \frac{1}{m} \sum_{i=1}^m (W x^{(i)} - y^{(i)})^2$$

How would you find the lowest point?

Finding root of derivative not computationally feasible with large number of features!



Gradient descent algorithm

- Minimize cost function
- Gradient descent is used many minimization problems
- For a given cost function, $\text{cost}(W, b)$, it will find W, b to minimize cost
- It can be applied to more general function: $\text{cost}(w_1, w_2, \dots)$

How does gradient descent work?

- Start with initial guesses
 - Start at 0,0 (or any other value)
 - Keeping changing W and b a little bit to try and reduce $\text{cost}(W, b)$
- Each time you change the parameters, you select the gradient which reduces $\text{cost}(W, b)$ the most possible
- Repeat
- Do so until you converge to a local minimum
- Has an interesting property
 - Where you start can determine which minimum you end up

Formal definition

$$cost(W) = \frac{1}{m} \sum_{i=1}^m (Wx^{(i)} - y^{(i)})^2$$



$$cost(W) = \frac{1}{2m} \sum_{i=1}^m (Wx^{(i)} - y^{(i)})^2$$

- Updating weights

$$W := W - \alpha \frac{\partial}{\partial W} cost(W)$$

Learning rate,
*determines the
step size*

gradient

Formal definition

$$cost(W) = \frac{1}{2m} \sum_{i=1}^m (Wx^{(i)} - y^{(i)})^2$$

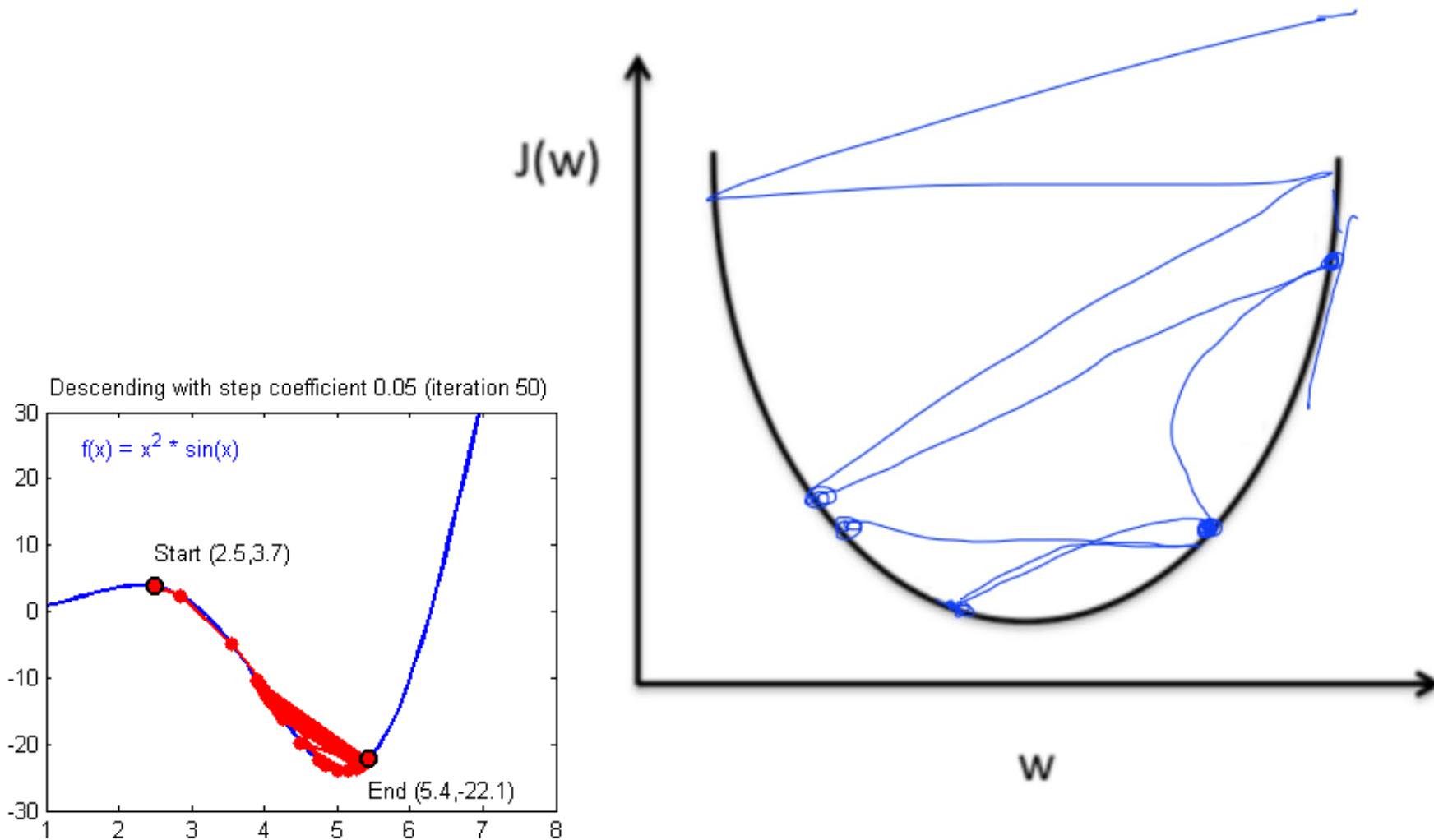
$$W := W - \alpha \frac{\partial}{\partial W} cost(W)$$

$$W := W - \alpha \frac{\partial}{\partial W} \frac{1}{2m} \sum_{i=1}^m (Wx^{(i)} - y^{(i)})^2$$

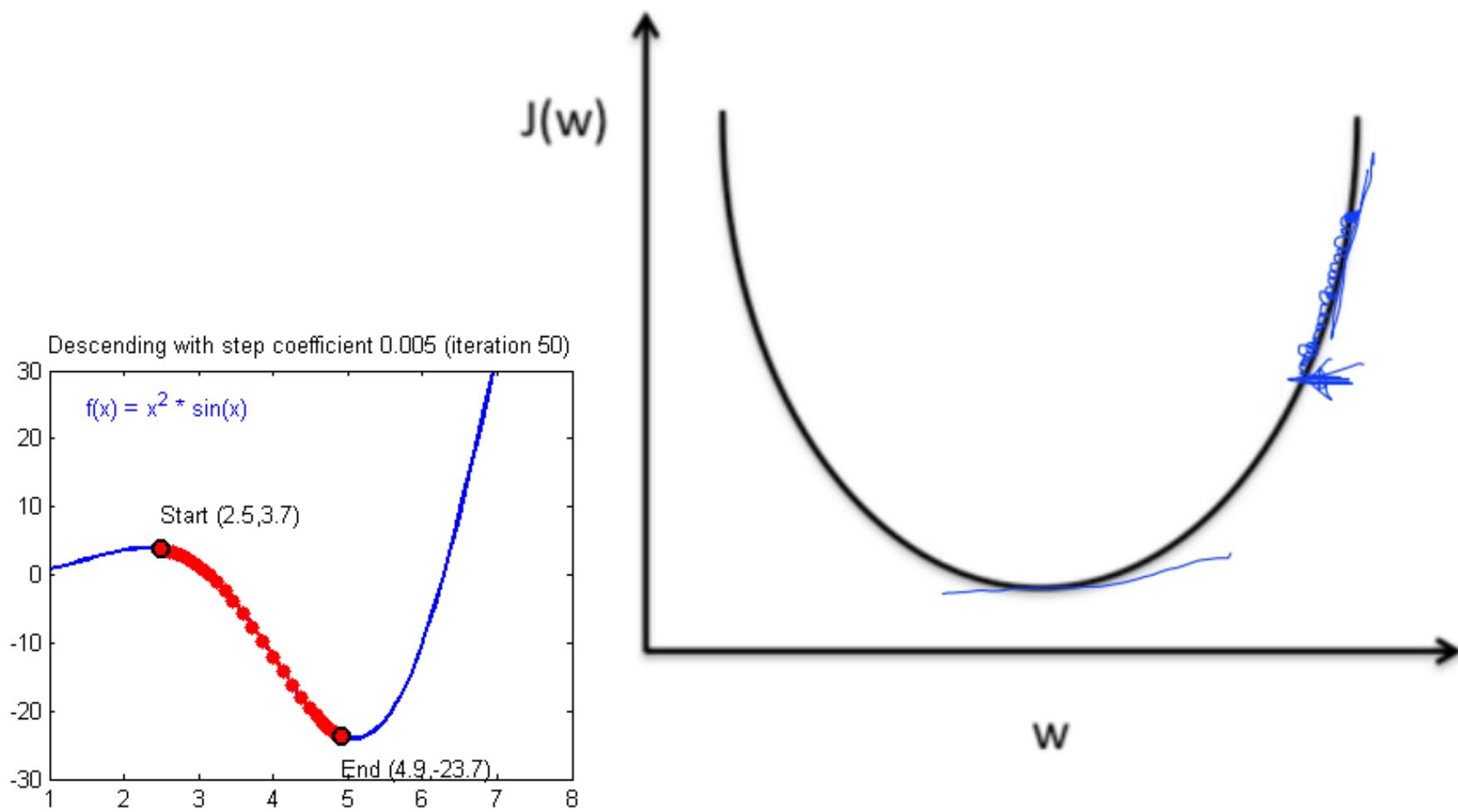
$$W := W - \alpha \frac{1}{2m} \sum_{i=1}^m 2(Wx^{(i)} - y^{(i)})x^{(i)}$$

$$W := W - \alpha \frac{1}{m} \sum_{i=1}^m (Wx^{(i)} - y^{(i)})x^{(i)}$$

Issues: large learning rate - overshooting

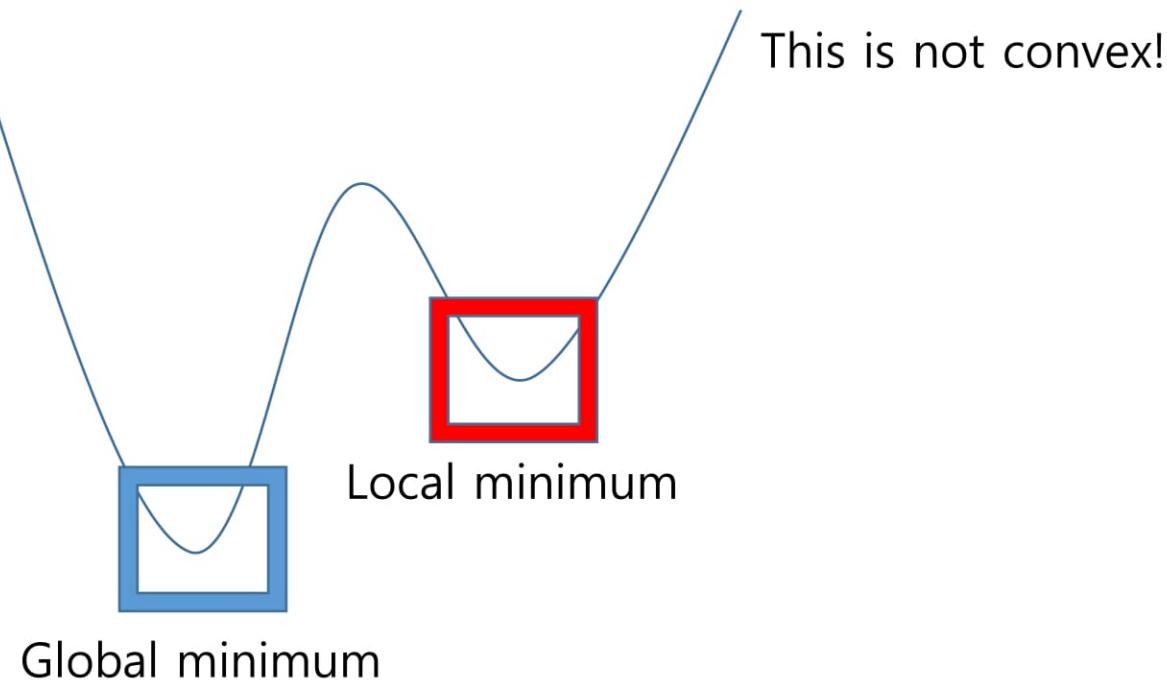


Issues:
small learning rate - takes too long



Issues: starting point is important

- Only works on convex functions!
- Can end up in local minimum depending on the starting point



Issues: only works on differentiable functions

$$W := W - \alpha \frac{\partial}{\partial W} cost(W)$$

We take the derivative of $cost(W)$,
therefore $cost(W)$ must be differentiable

$$cost(W) = \frac{1}{m} \sum_{i=1}^m (H(x^{(i)}) - y^{(i)})^2$$

By chain rule,
the hypothesis must also be differentiable

$$H(x) = Wx$$

Predicting exam score: regression using three inputs (x_1 , x_2 , x_3)

multi-variable/feature

x_1 (quiz 1)	x_2 (quiz 2)	x_3 (midterm 1)	Y (final)
73	80	75	152
93	88	93	185
89	91	90	180
96	98	100	196
73	66	70	142

Test Scores for General Psychology

For multiple features

$$H(x) = Wx + b$$

this is for a single feature

$$cost(W, b) = \frac{1}{m} \sum_{i=1}^m (H(x^{(i)}) - y^{(i)})^2$$

$$H(x_1, x_2, x_3) = w_1 x_1 + w_2 x_2 + w_3 x_3 + b$$

$$cost(W, b) = \frac{1}{m} \sum_{I=1}^m (H(x_1^{(i)}, x_2^{(i)}, x_3^{(i)}) - y^{(i)})^2$$

Recap

- (Linear) Hypothesis

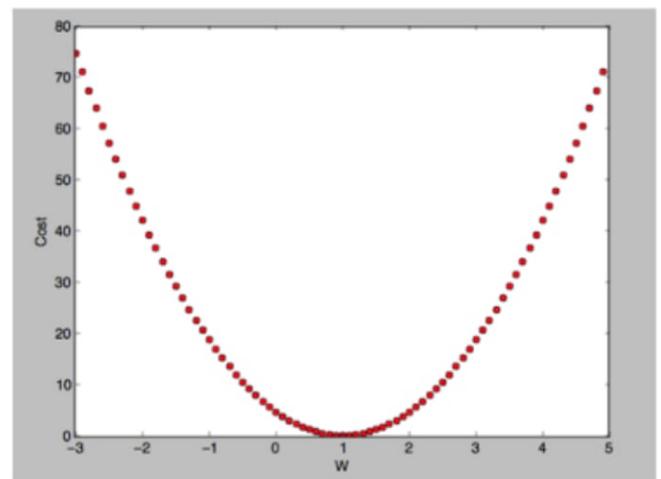
$$H(x_1, x_2, x_3, \dots, x_n) = w_1x_1 + w_2x_2 + w_3x_3 + \dots + w_nx_n + b$$

- Cost function

$$cost(W, b) = \frac{1}{m} \sum_{I=1}^m (H(x_1^{(i)}, x_2^{(i)}, x_3^{(i)}) - y^{(i)})^2$$

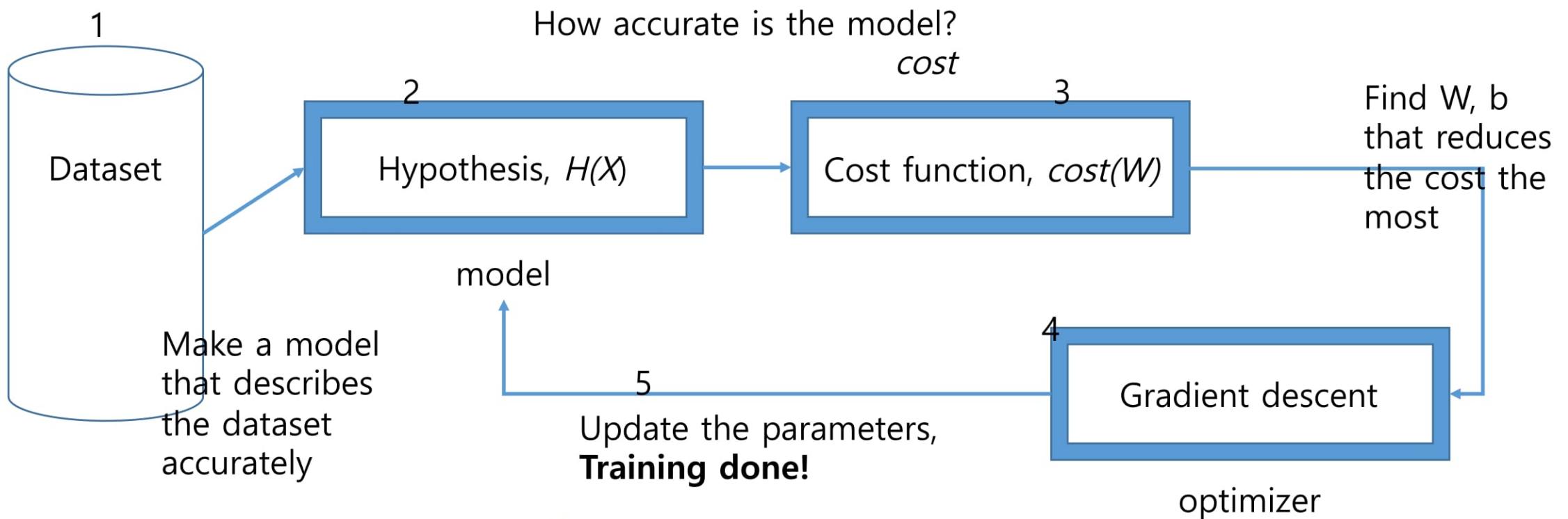
- Gradient descent algorithm

$$W := W - \alpha \frac{\partial}{\partial W} cost(W)$$



Cost function for single variable hypothesis

Training in a nutshell



So module! Many simple!
Wow.



wow

(Review) Types of supervised learning

- Predicting final exam score based on time spent
 - **Regression**
- Pass/non-pass based on time spent
 - **Binary classification**
- Letter grade(A, B, C, D, and F) based on time spent
 - **Multi-label classification**

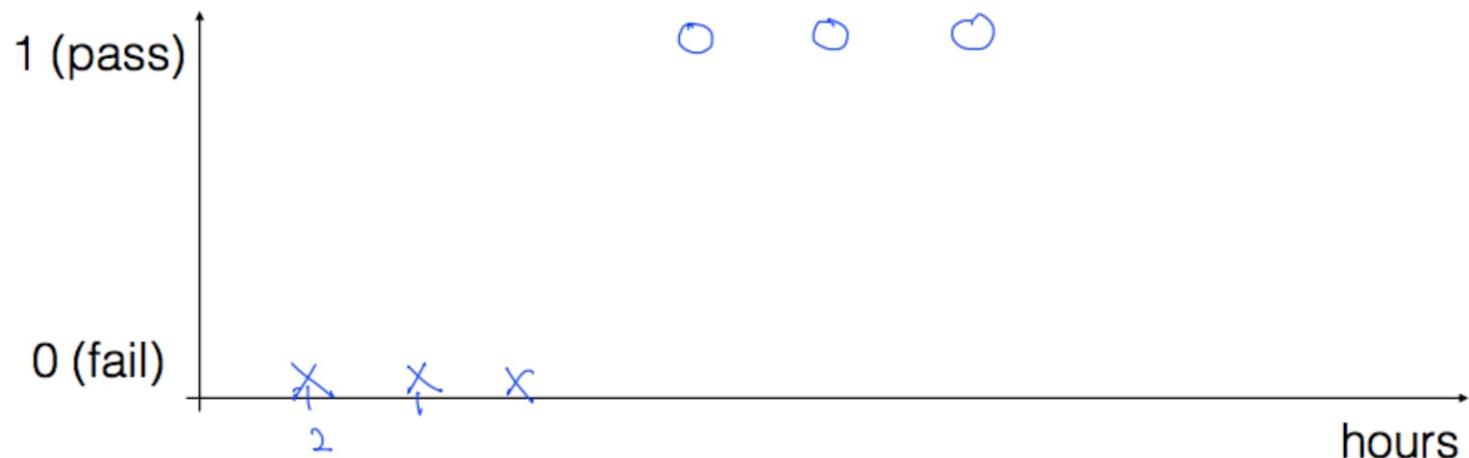
(Binary) Classification

- Spam Detection: Spam or Ham
- Facebook feed: show or hide
- Credit Card Fraudulent Transaction detection: legitimate/fraud

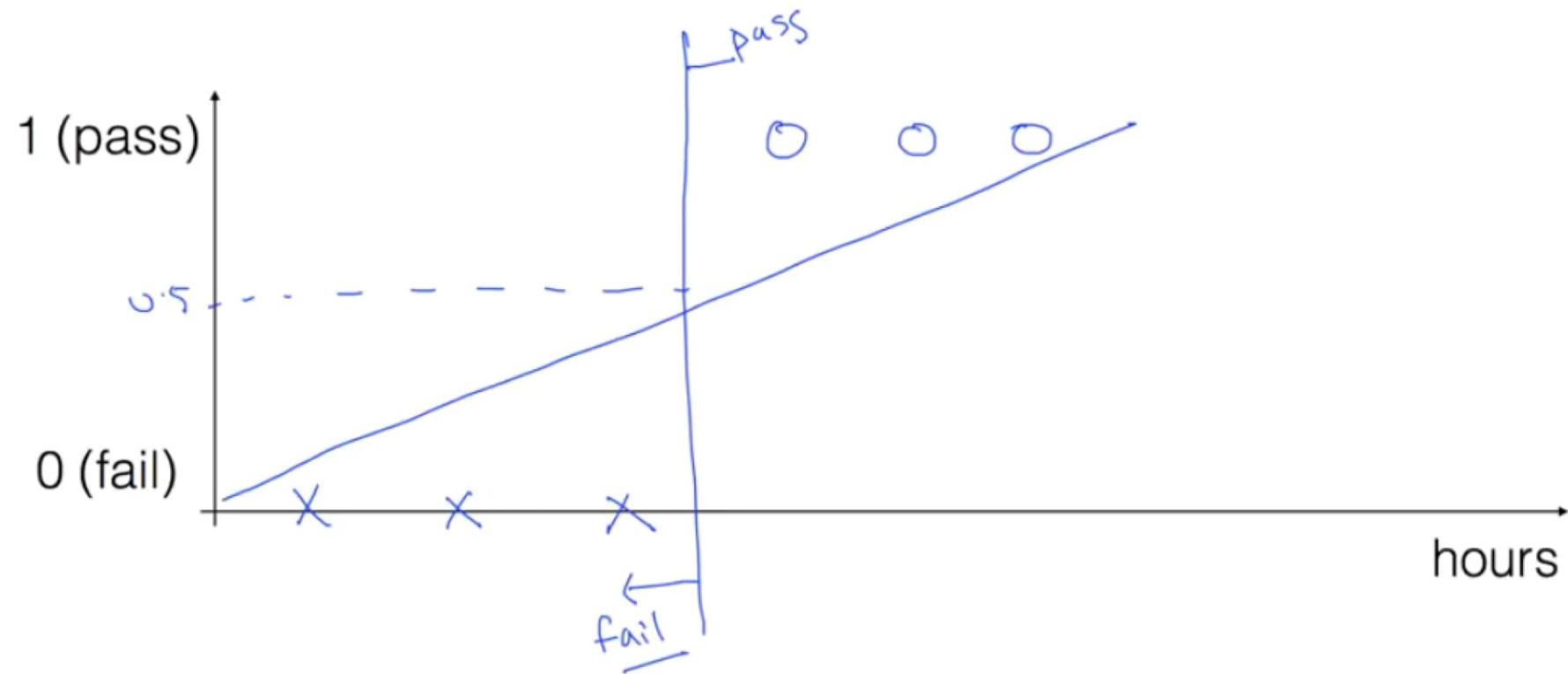
0, 1 encoding

- Spam Detection: Spam (1) or Ham (0)
- Facebook feed: show(1) or hide(0)
- Credit Card Fraudulent Transaction detection: legitimate(0) or fraud(1)

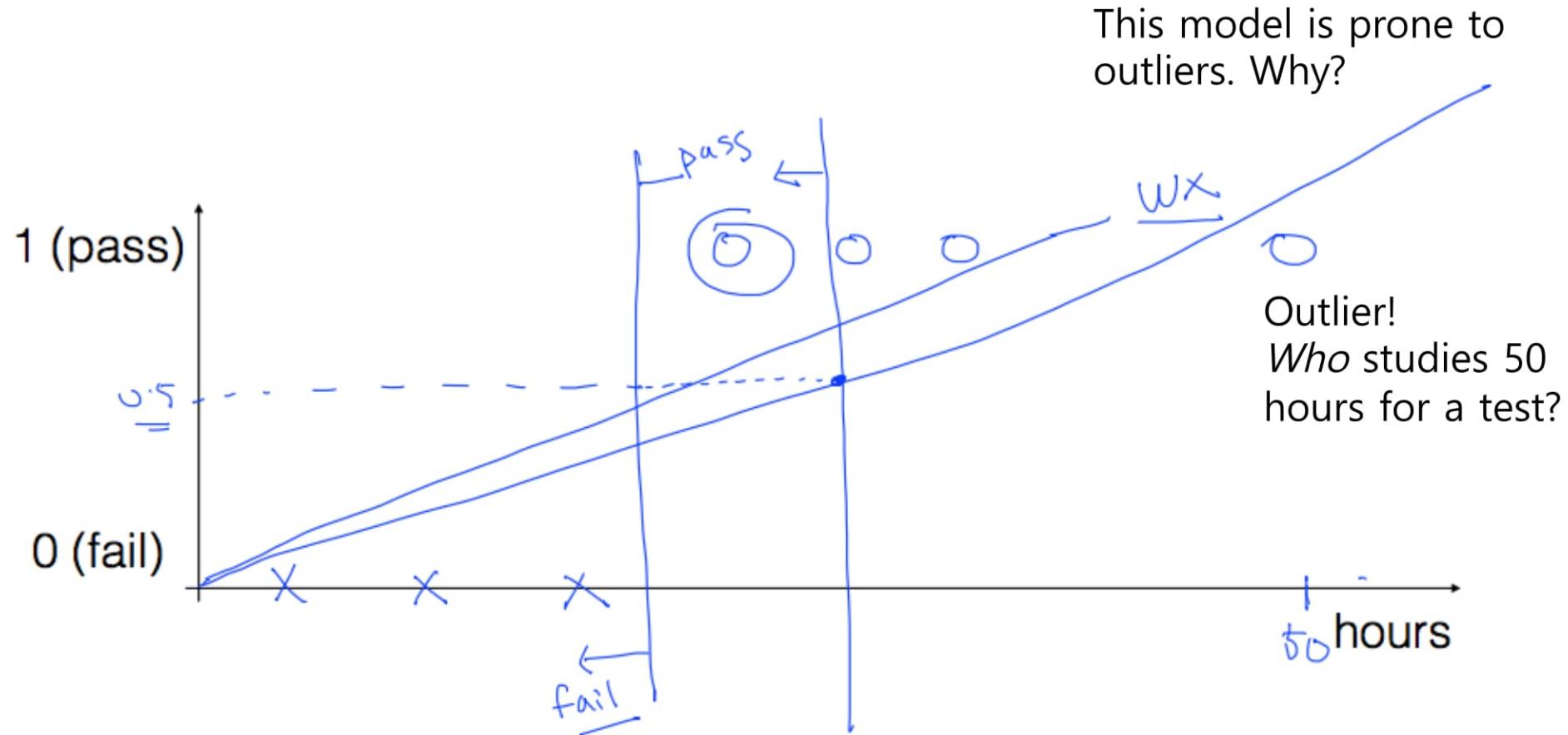
Pass(1)/Fail(0) based on study hours



Linear regression?



Linear regression?



We need conversion!

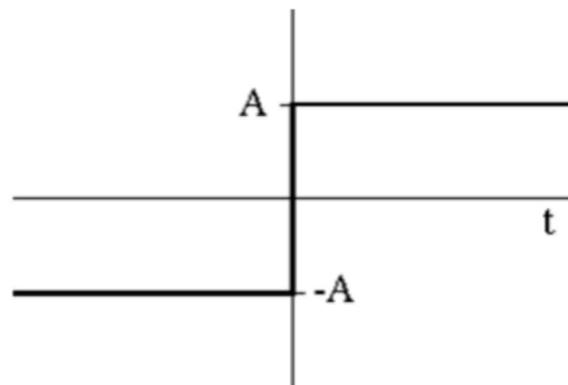
- We know Y is 0 or 1, and $H(X)$ outputs values larger than 1 or less than 0

$$H(x) = Wx + b$$



We want to limit the output range to be within $0 \sim 1$

How about this?



We need conversion!

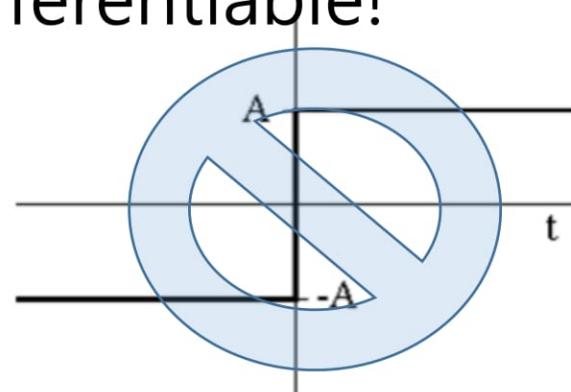
- We know Y is 0 or 1, and $H(X)$ outputs values larger than 1 or less than 0

$$H(x) = Wx + b$$



We want to limit the output range to be within $0 \sim 1$

- It has to be differentiable!

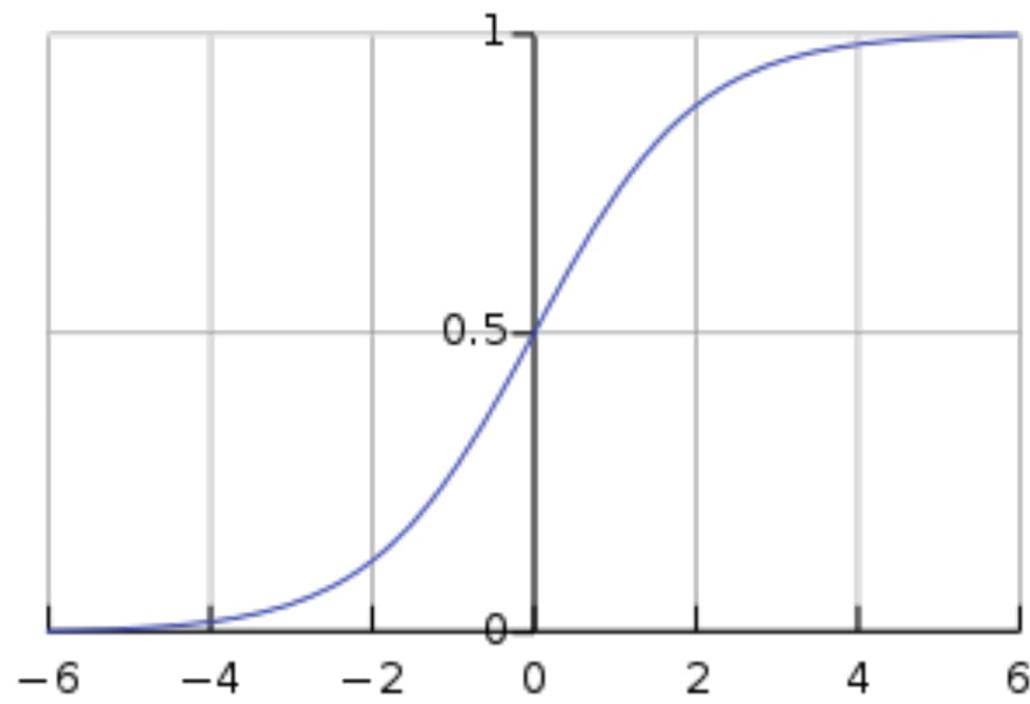
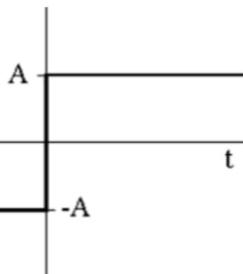


Plus, its derivative is 0 no matter what x. This poses a big problem later on.

Sigmoid

- Curved in two directions, like the letter "S", or the Greek ς (sigma)

$$H(X) = \frac{1}{1 + e^{-W^T X}}$$



We need conversion!

- We know Y is 0 or 1, and $H(X)$ outputs values larger than 1 or less than 0

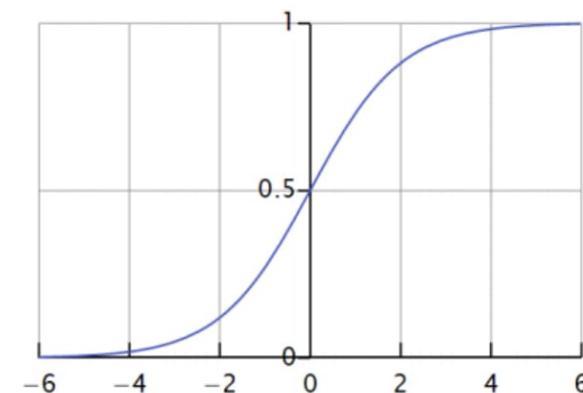
$$H(x) = Wx + b$$



We want to limit the output range to be within $0 \sim 1$

- It has to be differentiable!

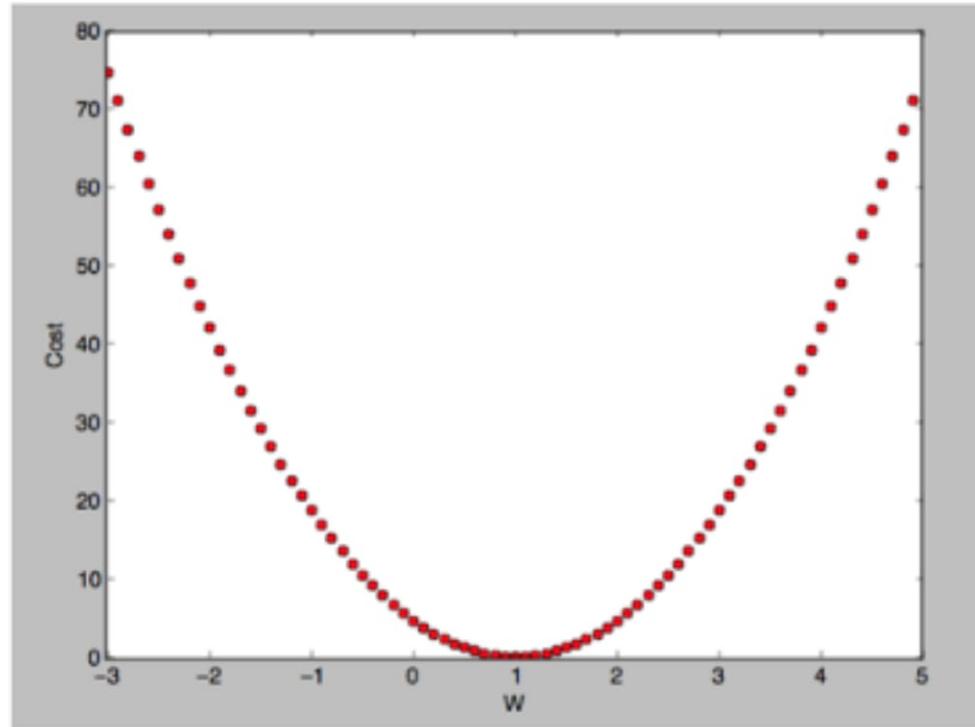
$$H(X) = \frac{1}{1 + e^{-W^T X}}$$



Sigmoid!

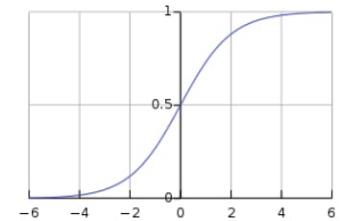
(Review) Cost

$$cost(W) = \frac{1}{m} \sum_{i=1}^m (Wx^{(i)} - y^{(i)})^2 \quad \text{when} \quad H(X) = XW$$



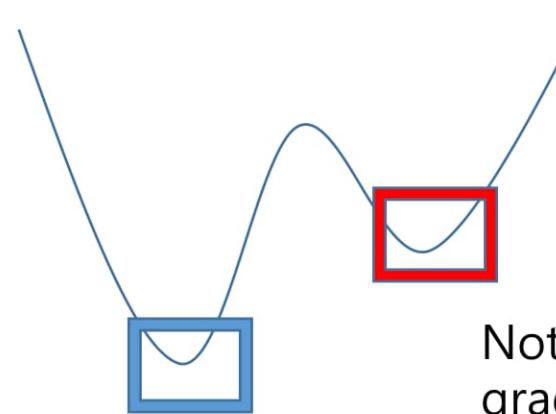
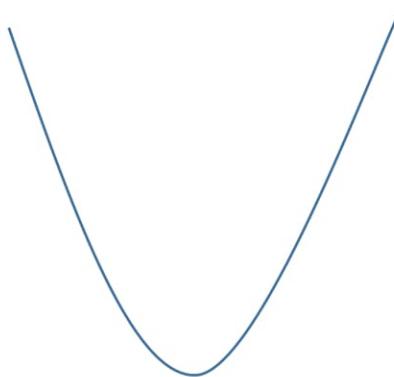
Cost function

$$cost(W, b) = \frac{1}{m} \sum_{i=1}^m (H(x^{(i)}) - y^{(i)})^2$$



$$H(x) = Wx + b$$

$$H(X) = \frac{1}{1 + e^{-W^T X}}$$



Not appropriate for
gradient descent
algorithm!

New cost function for logistic

$$cost(W) = \frac{1}{m} \sum c(H(x), y)$$

$$c(\underline{H(x)}, y) = \begin{cases} -\log(H(x)) & : y = 1 \\ -\log(1 - H(x)) & : y = 0 \end{cases}$$

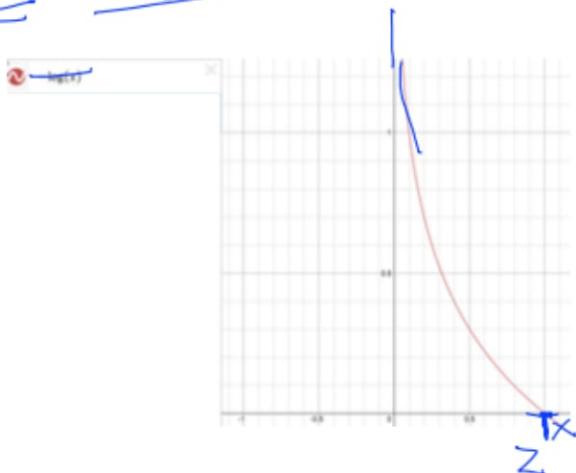
Real value(answer, label)

Understanding cost function

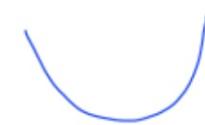
$$C(H(x), y) = \begin{cases} -\log(H(x)) & : y = 1 \\ -\log(1 - H(x)) & : y = 0 \end{cases}$$

Cost $\boxed{y=1}$
 $H(x)=1 \rightarrow \text{cost}(1) = 0,$
 $H(x)=0 \rightarrow \text{cost} = \infty \uparrow$

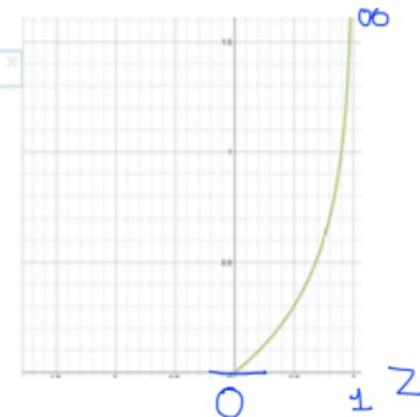
Cost
 $\approx g(z) = -\log(z)$



$y=0$
 $H(x)=0, \text{ cost}=0$
 $H(x)=1, \text{ cost}=\infty \uparrow$



$$-\log(1-z)$$



New cost function for logistic

$$cost(W) = \frac{1}{m} \sum c(H(x), y)$$

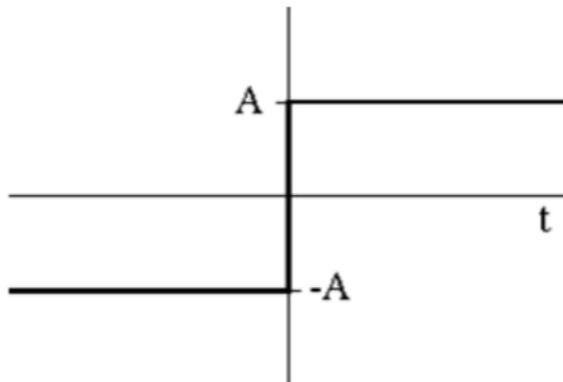
$$c(\underline{H(x)}, y) = \begin{cases} -\log(H(x)) & : y = 1 \\ -\log(1 - H(x)) & : y = 0 \end{cases}$$

$$C(H(x), y) = -y \log(H(x)) - (1 - y) \log(1 - H(x))$$

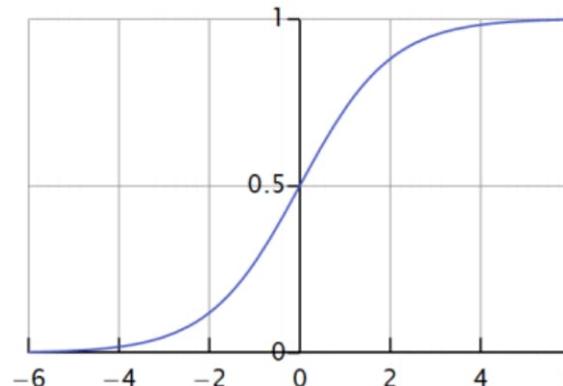
Minimize cost with gradient descent algorithm

$$\underline{cost}(W) = -\frac{1}{m} \sum y \log(H(x)) + (1 - y) \log(1 - H(x))$$

$$W := W - \alpha \frac{\partial}{\partial W} cost(W)$$

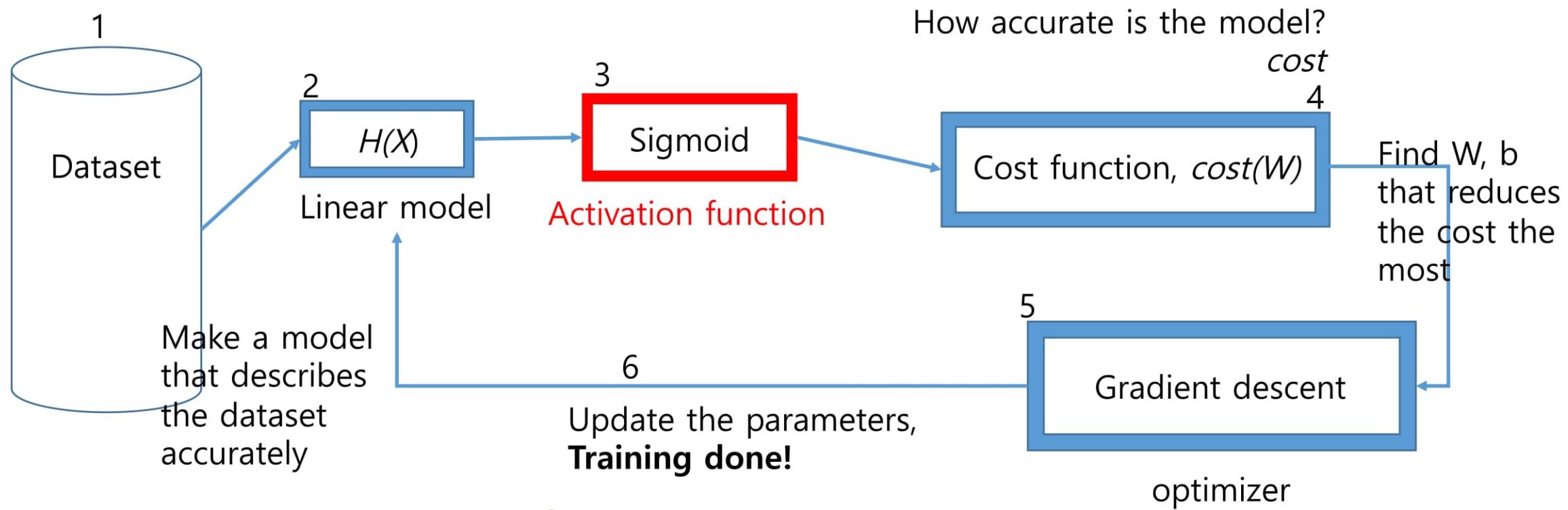
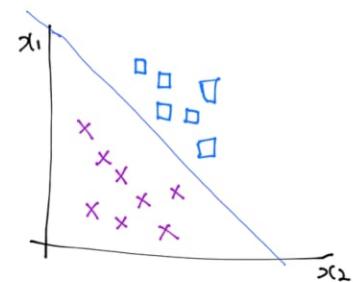


vs



Why sigmoid?

Training in a nutshell



So module! Many simple!
Wow.



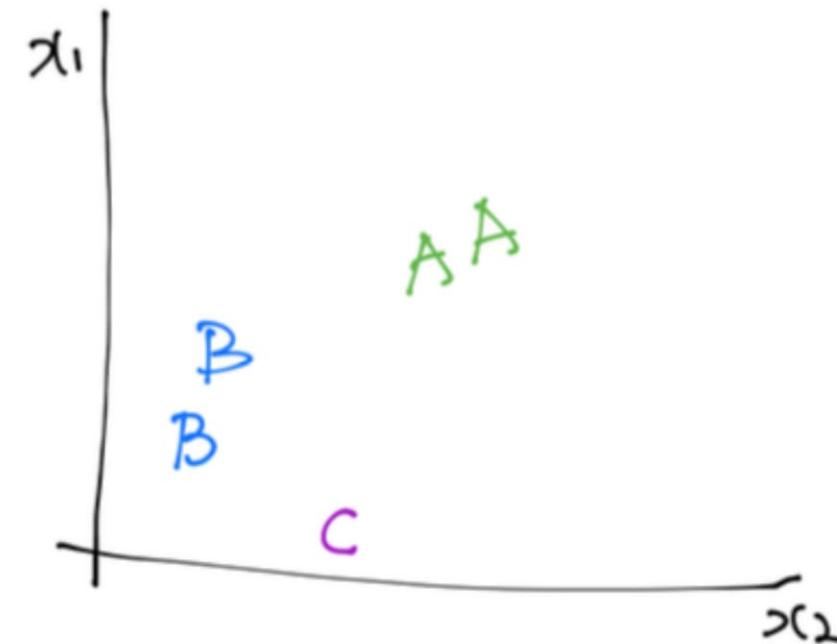
wow

(Review) Types of supervised learning

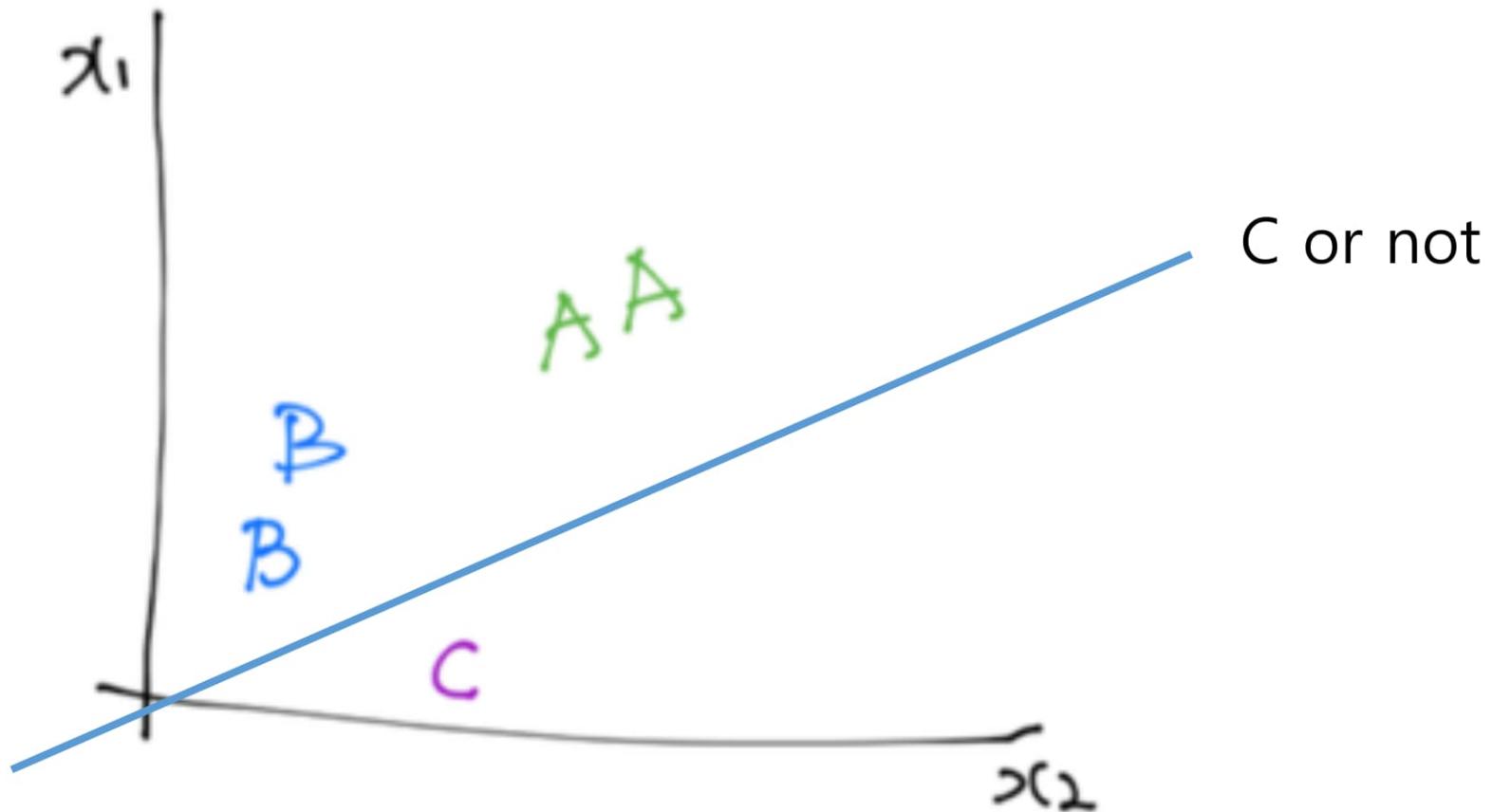
- Predicting final exam score based on time spent
 - **Regression**
- Pass/non-pass based on time spent
 - **Binary classification**
- Letter grade(A, B, C, D, and F) based on time spent
 - **Multi-label classification**

Multinomial classification

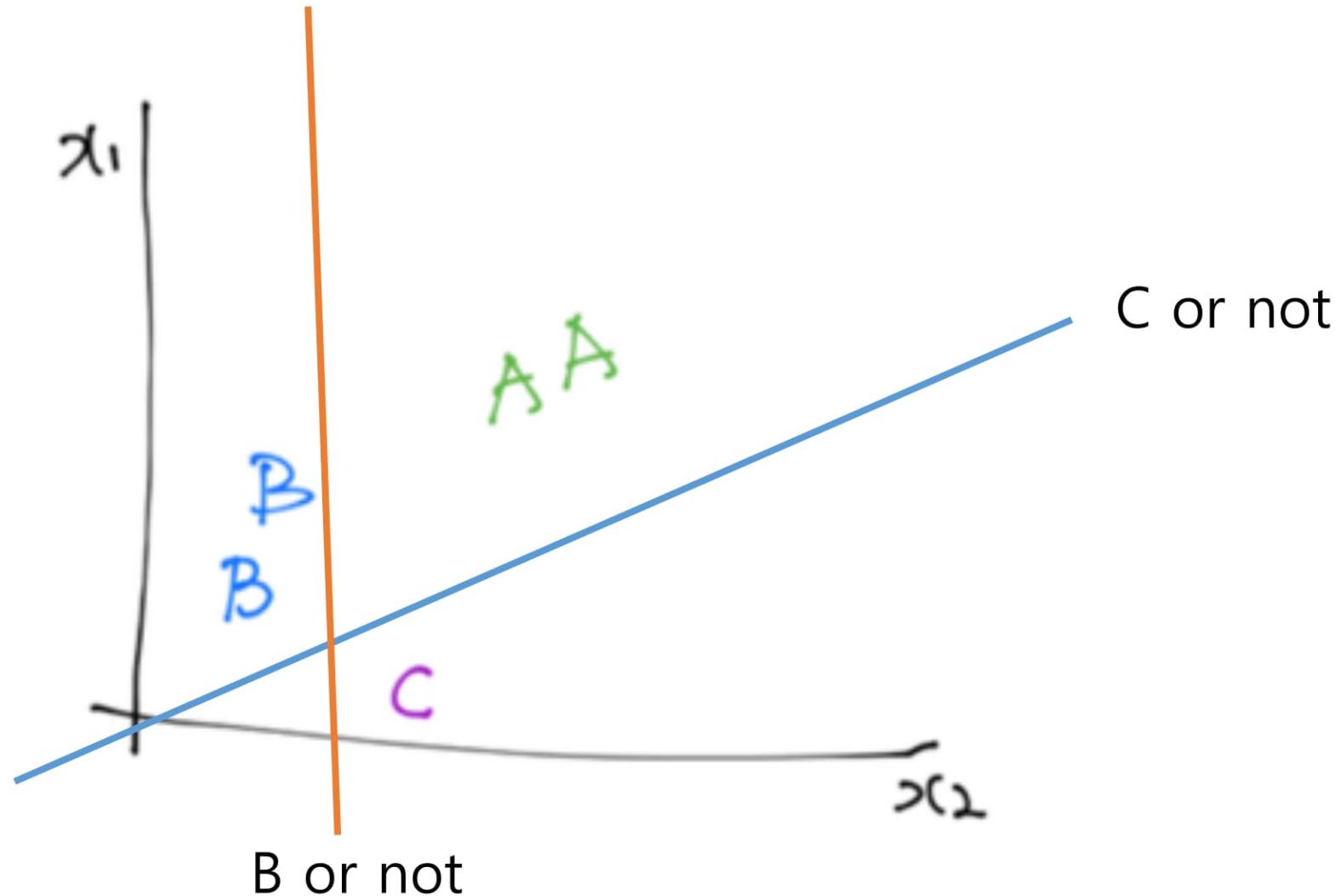
x1 (hours)	x2 (attendance)	y (grade)
10	5	A
9	5	A
3	2	B
2	4	B
11	1	C



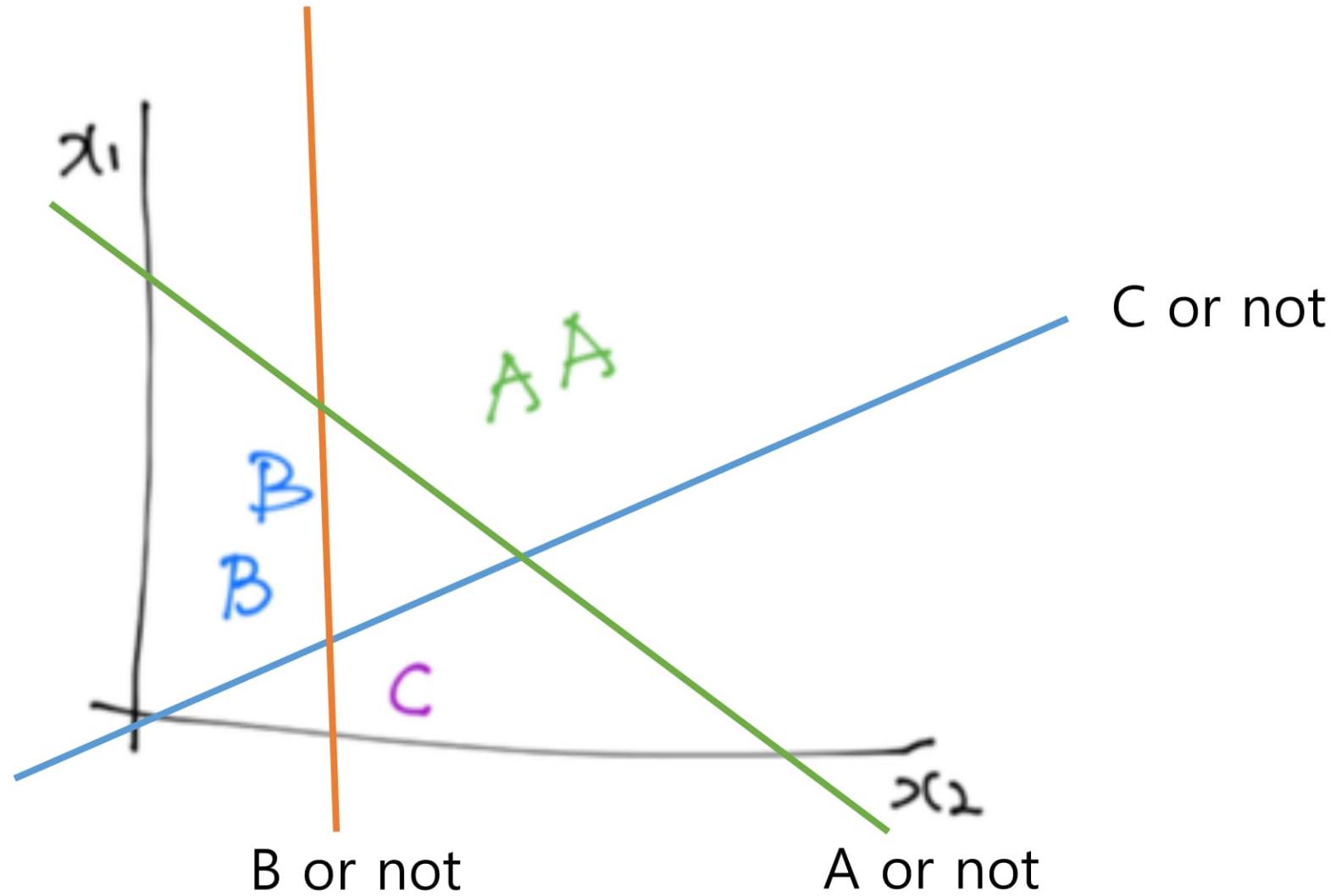
Multinomial classification



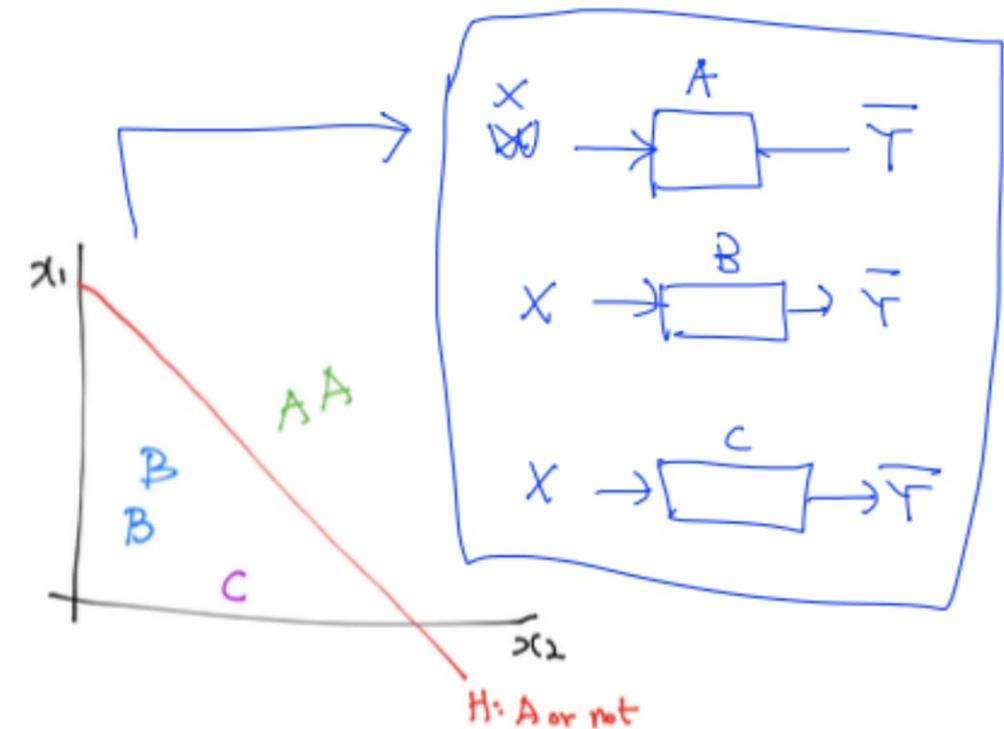
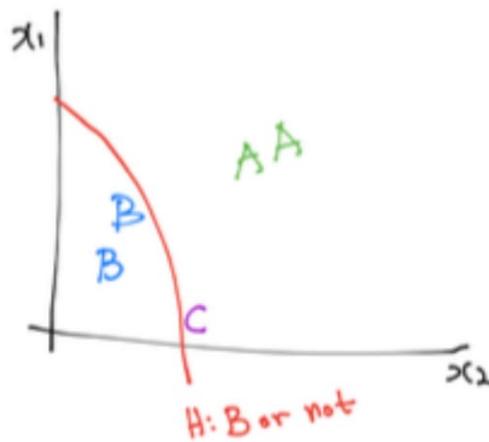
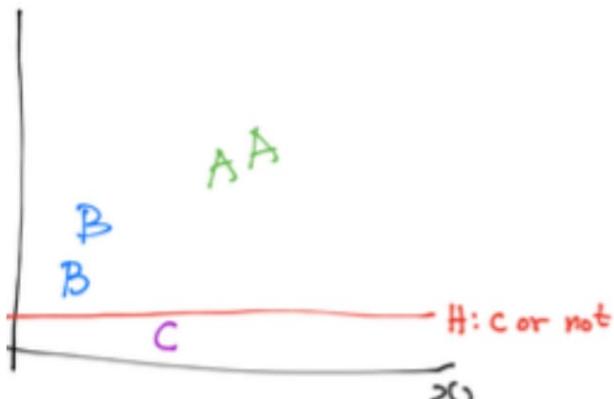
Multinomial classification



Multinomial classification



Multinomial classification

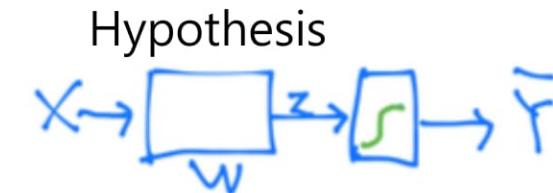


Multinomial classification

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

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

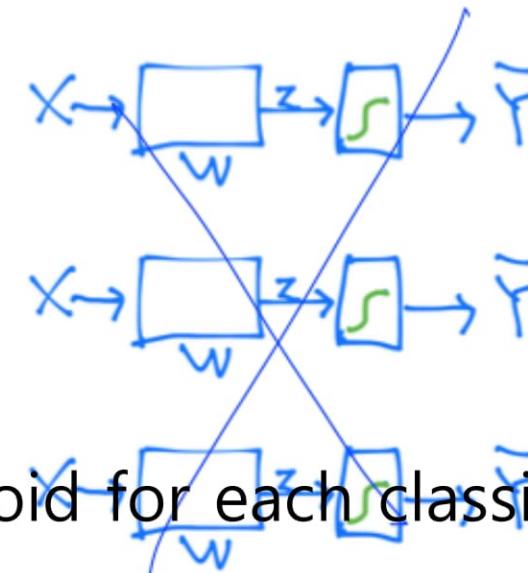
$$[w_1 \ w_2 \ w_3] \begin{bmatrix} x_1 \\ x_2 \\ x_3 \end{bmatrix} = [w_1 x_1 + w_2 x_2 + w_3 x_3]$$



Activation function

Where is sigmoid?

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



1. Sigmoid for each classifier is too much work!

2. Does not add up to 1 (inappropriate for probabilities)

Where is sigmoid?

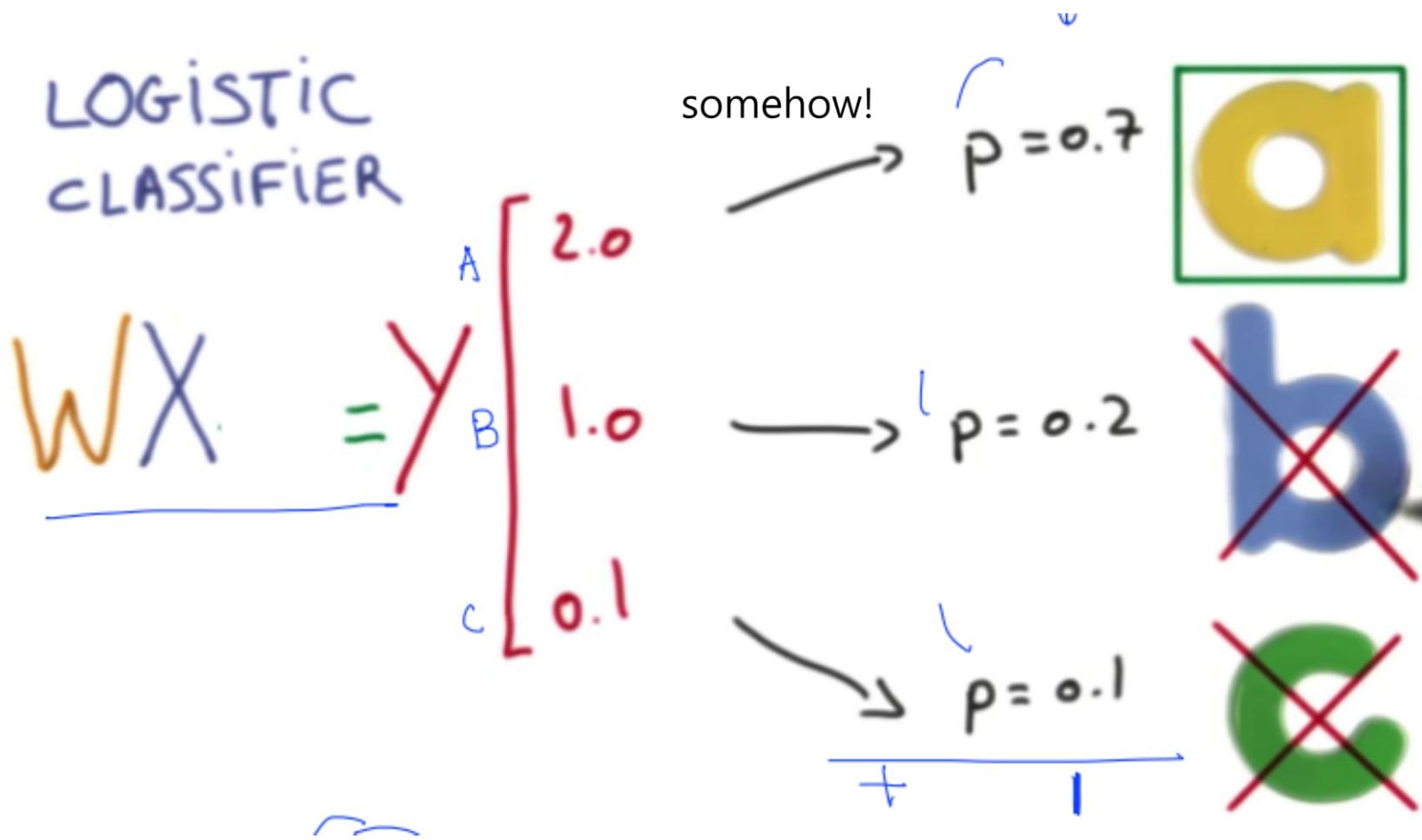
$$\begin{bmatrix} w_{A1} & w_{A2} & w_{A3} \\ w_{B1} & w_{B2} & w_{B3} \\ w_{C1} & w_{C2} & w_{C3} \end{bmatrix} \begin{bmatrix} x_1 \\ x_2 \\ x_3 \end{bmatrix} = \begin{bmatrix} w_{A1}x_1 + w_{A2}x_2 + w_{A3}x_3 \\ w_{B1}x_1 + w_{B2}x_2 + w_{B3}x_3 \\ w_{C1}x_1 + w_{C2}x_2 + w_{C3}x_3 \end{bmatrix}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

0 ~ 1
↓



 We need an activation function that outputs values between $0 \sim 1$, and that add up to 1

No more sigmoid



Softmax regression

SOFTMAX

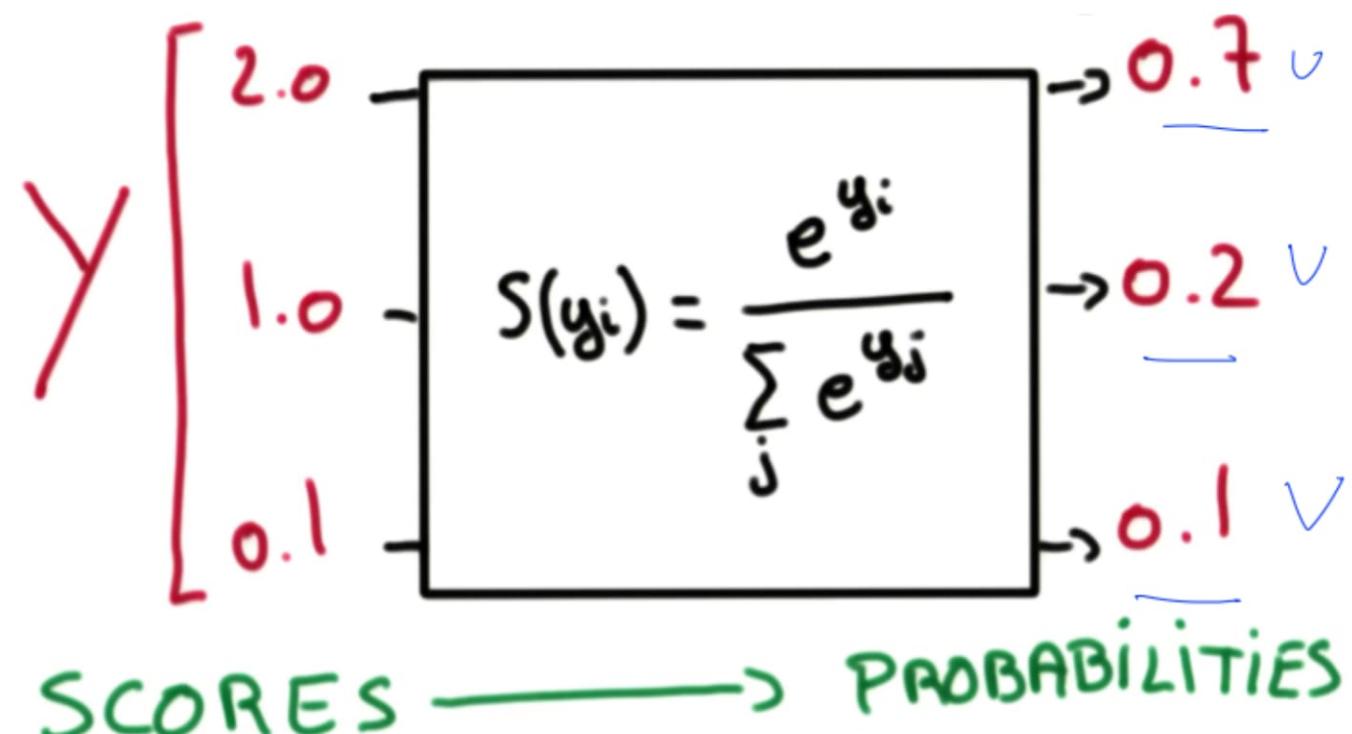
First step:
 $\exp(score)$

Second step:
Normalize

Why "e"?

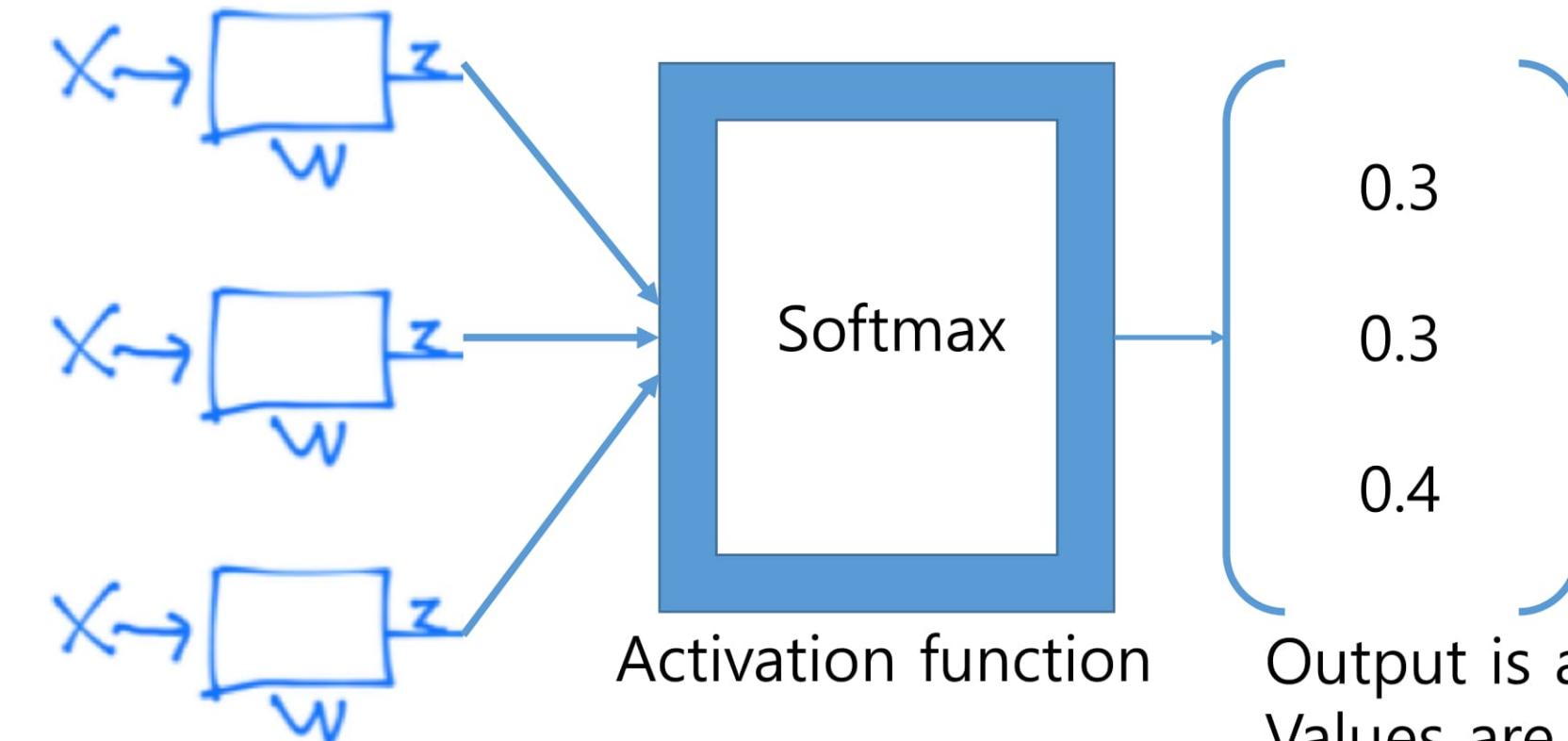
Output values are:

1. between 0 ~ 1
2. normalized



Softmax regression in a nutshell

Dataset

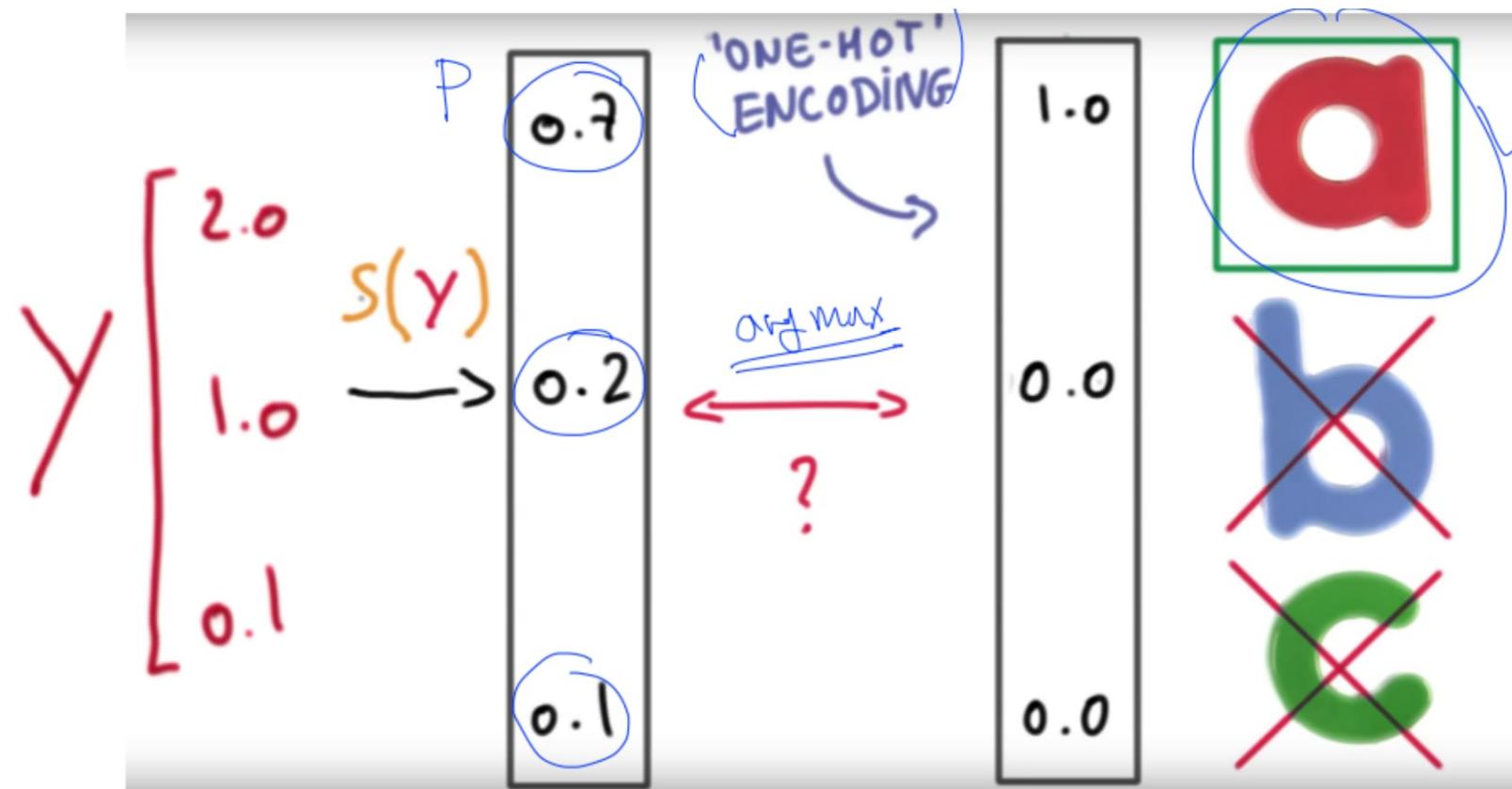


N classifiers
for N
outputs

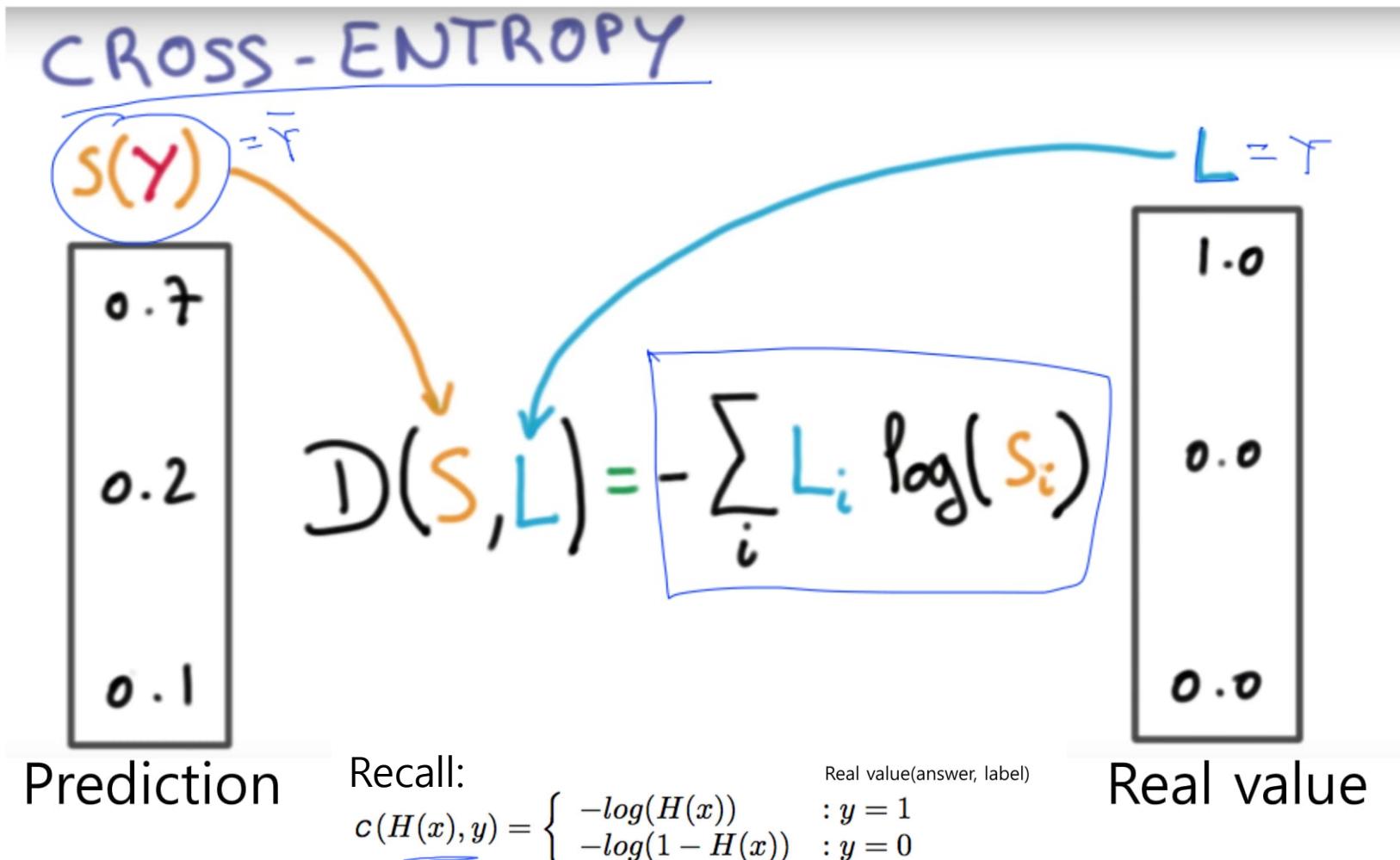
Output is a vector.
Values are

1. between 0 ~ 1
2. normalized

One-hot encoding(optional)

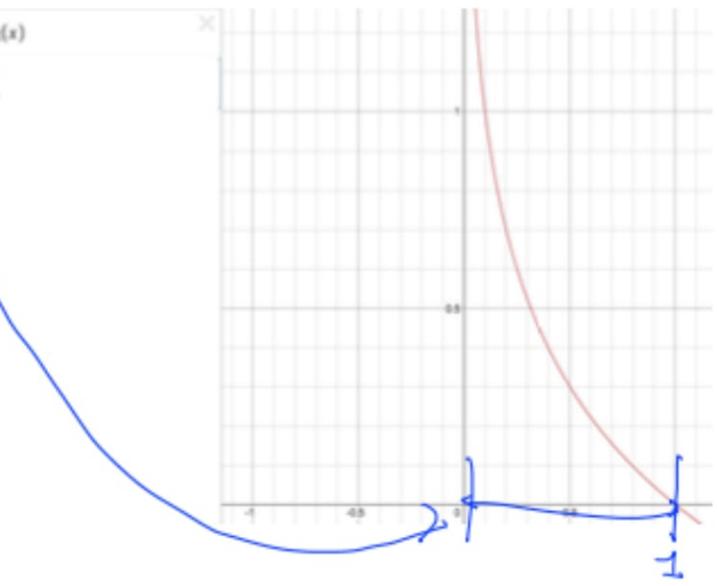


Cost function: cross-entropy



Cross-entropy cost function

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



We are only interested in range (0, 1)
why?

Logistic cost vs. cross entropy

Cost for binary classification

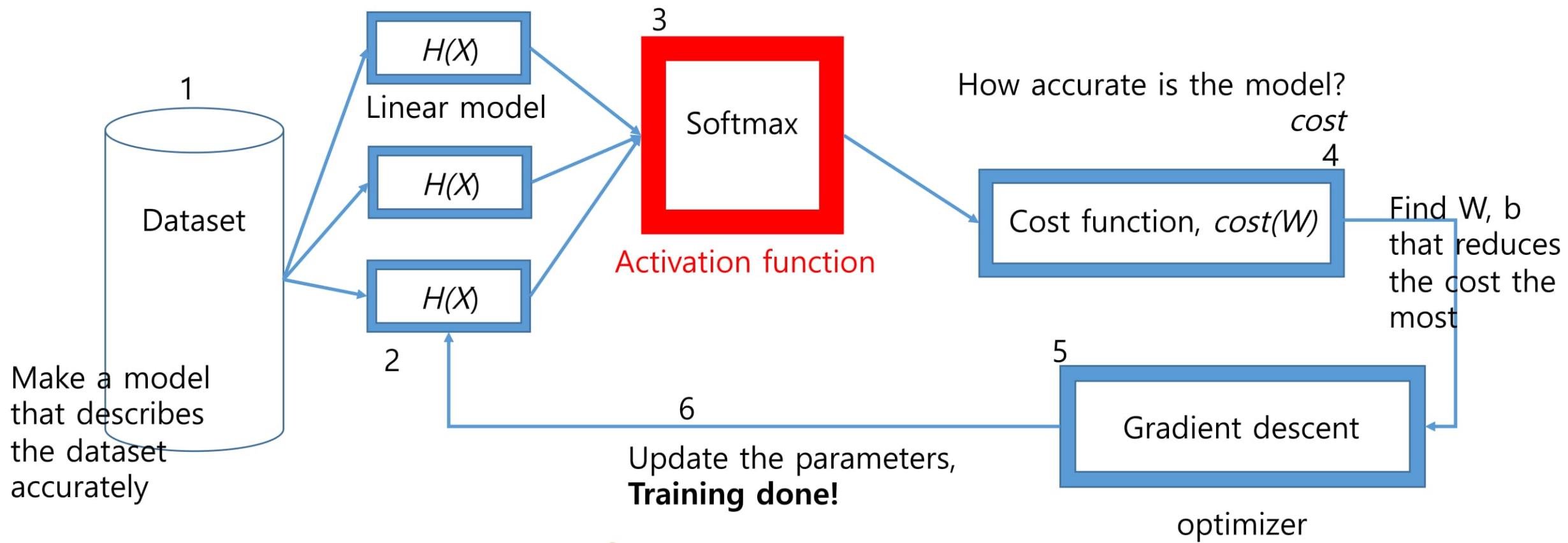
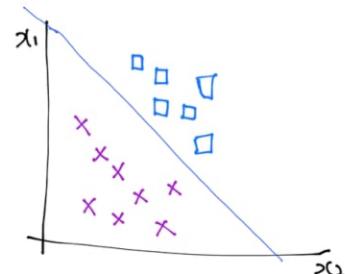
$$C(H(x), y) = y \log(H(x)) - (1 - y) \log(1 - H(x))$$

$$D(S, L) = - \sum_i L_i \log(S_i)$$



Cost for multinomial classification:
generalized logistic cost for n-classification

Training in a nutshell



So module! Many simple!
Wow.



wow