

# cs231n

lec3

# 선형분류(linear classification)

**So far:** Defined a (linear) score function  $f(x, W) = Wx + b$

Example class  
scores for 3  
images for  
some  $W$ :



airplane	-3.45	-0.51	3.42
automobile	-8.87	<b>6.04</b>	4.64
bird	0.09	5.31	2.65
cat	<b>2.9</b>	-4.22	5.1
deer	4.48	-4.19	2.64
dog	8.02	3.58	5.55
frog	3.78	4.49	<b>-4.34</b>
horse	1.06	-4.37	-1.5
ship	-0.36	-2.09	-4.79
truck	-0.72	-2.93	6.14

How can we tell  
whether this  $W$   
is good or bad?

image by Nikita is licensed under CC-BY 2.0  
image is CC0 1.0 public domain  
image is in the public domain

- 이 모델의 분류( $W$  값)이 좋을까?

- 정도를 측정할 도구가 필요  
→ loss function(손실함수)

Given a dataset of examples

$$\{(x_i, y_i)\}_{i=1}^N$$

Where  $x_i$  is image and  
 $y_i$  is (integer) label

Loss over the dataset is a  
sum of loss over examples:

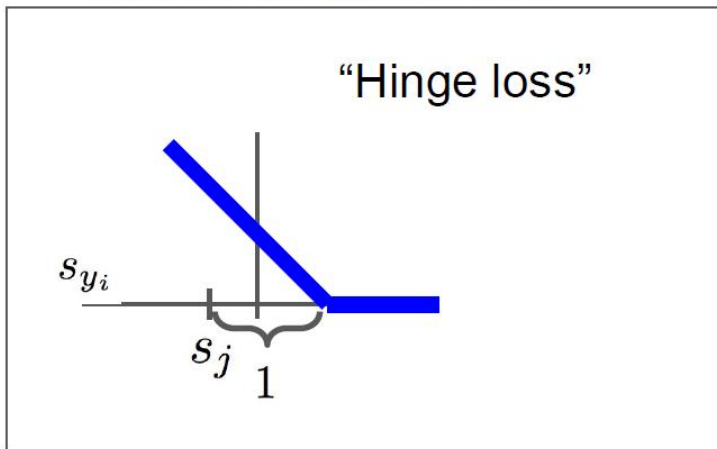
$$L = \frac{1}{N} \sum_i L_i(f(x_i, W), y_i)$$

## 손실 함수 (Loss function)

1. Multiclass SVM loss
2. Softmax Classifier

# Multiclass SVM loss

## Multiclass SVM loss:



$$L_i = \sum_{j \neq y_i} \begin{cases} 0 & \text{if } s_{y_i} \geq s_j + 1 \\ s_j - s_{y_i} + 1 & \text{otherwise} \end{cases}$$

$$= \sum_{j \neq y_i} \max(0, s_j - s_{y_i} + 1)$$

Suppose: 3 training examples, 3 classes.  
With some  $W$  the scores  $f(x, W) = Wx$  are:



cat	<b>3.2</b>	1.3	2.2
car	5.1	<b>4.9</b>	2.5
frog	-1.7	2.0	<b>-3.1</b>

<b>Loss</b>	<b>2.9</b>	<b>0</b>	<b>12.9</b>
$L = \frac{(2.9 + 0 + 12.9)}{3}$ $= 5.7$	$= \max(5.1 - 3.2 + 1, 0)$ $+ \max(-1.7 - 3.2 + 1, 0)$		

# Several question

- Q: What happens to loss if car scores change a bit?
- A: No change
- Q: what is the min/max possible loss?
- A: min = 0 , max = infinite
- Q: At initialization W is small so all  $s \approx 0$ . What is the loss?
- A:  $\frac{1}{N} \sum_i \sum_{j \neq y_i} \max(s_j - s_{y_i} + 1, 0) = \frac{1}{N} \sum_i (N - 1) = \frac{N(N-1)}{N} = N - 1$
- Q: What if the sum was over all classes? (including  $j = y_i$ )
- A: loss value(손실값) + 1
- Q:  $L_i$  값을 구할때 합대신 평균을 사용한다면?
- A: loss value는 작아지지만, 큰의미는 없다
- Q:  $L_i = \sum_{j \neq y_i} \max(s_j - s_{y_i} + 1, 0)^2$  으로 정의한다면
- A: 다른분류가 된다.

# Softmax classifier

scores = unnormalized log probabilities of the classes.

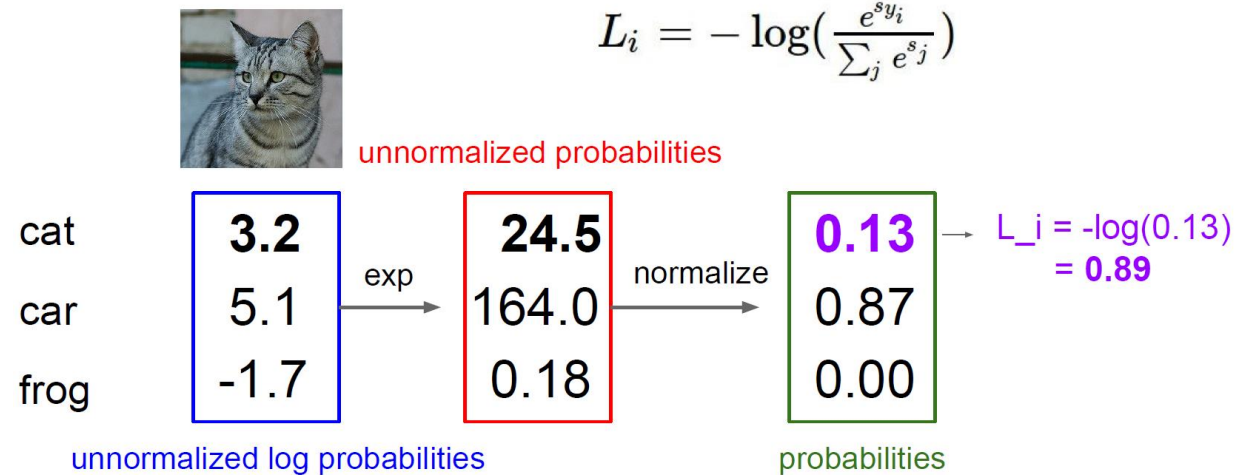
$$P(Y = k|X = x_i) = \frac{e^{s_k}}{\sum_j e^{s_j}} \quad \text{where} \quad s = f(x_i; W)$$

Want to maximize the log likelihood, or (for a loss function) to minimize the negative log likelihood of the correct class:

$$L_i = -\log P(Y = y_i|X = x_i)$$

---

in summary:  $L_i = -\log\left(\frac{e^{s_{y_i}}}{\sum_j e^{s_j}}\right)$



# $w \rightarrow 2w$ , loss is same!

Suppose: 3 training examples, 3 classes.  
With some  $W$  the scores  $f(x, W) = Wx$  are:



cat	<b>3.2</b>	1.3	2.2
car	5.1	<b>4.9</b>	2.5
frog	-1.7	2.0	<b>-3.1</b>
Losses:	2.9	0	

$$L_i = \sum_{j \neq y_i} \max(0, s_j - s_{y_i} + 1)$$

**Before:**

$$\begin{aligned} &= \max(0, 1.3 - 4.9 + 1) \\ &\quad + \max(0, 2.0 - 4.9 + 1) \\ &= \max(0, -2.6) + \max(0, -1.9) \\ &= 0 + 0 \\ &= 0 \end{aligned}$$

**With  $W$  twice as large:**

$$\begin{aligned} &= \max(0, 2.6 - 9.8 + 1) \\ &\quad + \max(0, 4.0 - 9.8 + 1) \\ &= \max(0, -6.2) + \max(0, -4.8) \\ &= 0 + 0 \\ &= 0 \end{aligned}$$

→ 조정이 필요 = regularization

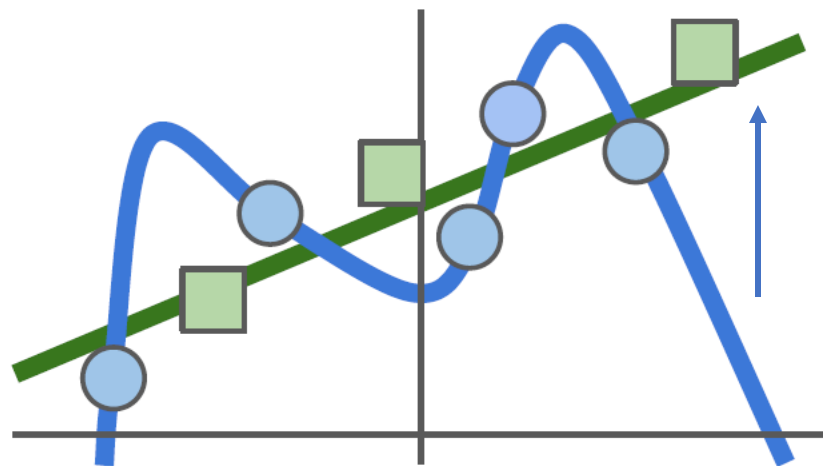
$$L(W) = \underbrace{\frac{1}{N} \sum_{i=1}^N L_i(f(x_i, W), y_i)}_{\text{Data loss: Model predictions should match training data}} + \underbrace{\lambda R(W)}_{\text{Regularization: Model should be "simple", so it works on test data}}$$

**Data loss:** Model predictions should match training data

**Regularization:** Model should be “simple”, so it works on test data

● training data  
■ real data

Overfitting





# Regularization

$\lambda$  = regularization strength  
(hyperparameter)

$$L = \frac{1}{N} \sum_{i=1}^N \sum_{j \neq y_i} \max(0, f(x_i; W)_j - f(x_i; W)_{y_i} + 1) + \lambda R(W)$$

In common use:

**L2 regularization**

$$R(W) = \sum_k \sum_l W_{k,l}^2$$

L1 regularization

$$R(W) = \sum_k \sum_l |W_{k,l}|$$

Elastic net (L1 + L2)  $R(W) = \sum_k \sum_l \beta W_{k,l}^2 + |W_{k,l}|$

Max norm regularization (might see later)

Dropout (will see later)

Fancier: Batch normalization, stochastic depth

## L2 Regularization (Weight Decay)

$$x = [1, 1, 1, 1]$$

$$w_1 = [1, 0, 0, 0]$$

$$w_2 = [0.25, 0.25, 0.25, 0.25]$$

$$R(W) = \sum_k \sum_l W_{k,l}^2$$

(If you are a Bayesian: L2 regularization also corresponds MAP inference using a Gaussian prior on  $W$ )

$$w_1^T x = w_2^T x = 1$$

$$R(W) = \sum_k \sum_l |W_{k,l}|$$

L1 Regularization 은  $w_1$  선호

L2 Regularization 은  $w_2$  선호

# Recap

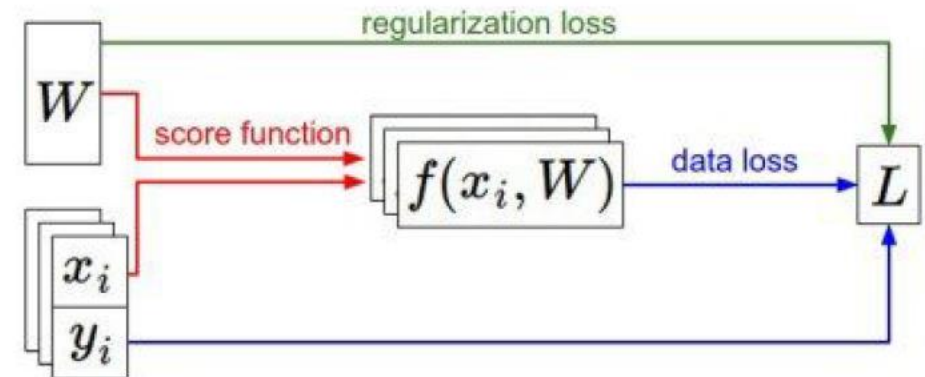
## How do we find the best $W$ ?

- We have some dataset of  $(x, y)$
- We have a **score function**:  $s = f(x; W) \stackrel{\text{e.g.}}{=} Wx$
- We have a **loss function**:

$$L_i = -\log\left(\frac{e^{s_{y_i}}}{\sum_j e^{s_j}}\right) \quad \text{Softmax}$$

$$L_i = \sum_{j \neq y_i} \max(0, s_j - s_{y_i} + 1) \quad \text{SVM}$$

$$L = \frac{1}{N} \sum_{i=1}^N L_i + R(W) \quad \text{Full loss}$$



# Optimization (적절한 $w$ 찾기)

- Random search
- → 시간낭비, 정답에대한 보장이 없음.
- Following Slope 기울기 따라가기

**current W:**

[0.34,  
-1.11,  
0.78,  
0.12,  
0.55,  
2.81,  
-3.1,  
-1.5,  
0.33,...]

**loss 1.25347**

**W + h (second dim):**

[0.34,  
-1.11 + **0.0001**,  
0.78,  
0.12,  
0.55,  
2.81,  
-3.1,  
-1.5,  
0.33,...]

**loss 1.25353**

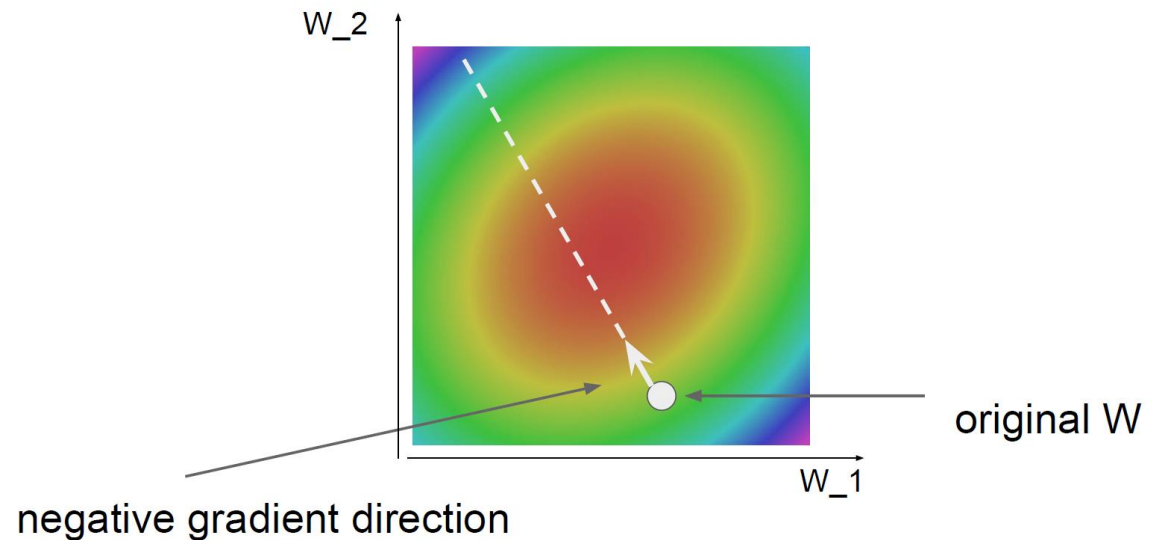
**gradient dW:**

[-2.5,  
?,  
?,  
?,  
?,  
?,  
?,  
?,  
?,...]

$$\frac{df(x)}{dx} = \lim_{h \rightarrow 0} \frac{f(x+h) - f(x)}{h}$$

# Using calculus

- Faster than numerical way
- But some time numerical gradient useful to debug
- Go plat area(red)

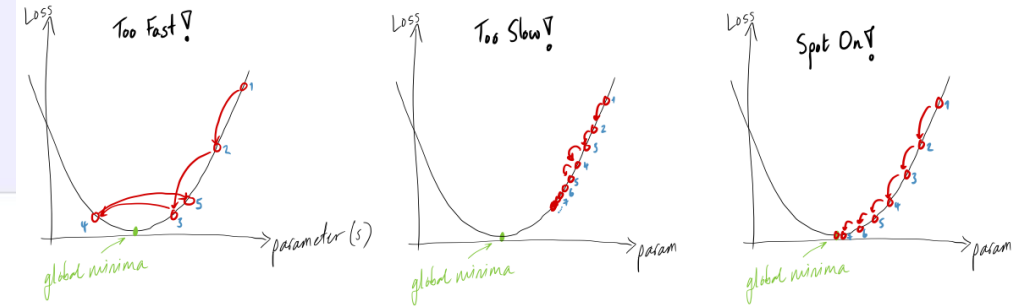


# Gradient Descent

Step\_size = learning rate.  
→ Large size or small size?

```
# Vanilla Gradient Descent

while True:
    weights_grad = evaluate_gradient(loss_fun, data, weights)
    weights += - step_size * weights_grad # perform parameter update
```



## Stochastic Gradient Descent (SGD)

$$L(W) = \frac{1}{N} \sum_{i=1}^N L_i(x_i, y_i, W) + \lambda R(W)$$

$$\nabla_W L(W) = \frac{1}{N} \sum_{i=1}^N \nabla_W L_i(x_i, y_i, W) + \lambda \nabla_W R(W)$$

Full sum expensive  
when N is large!

Approximate sum  
using a **minibatch** of  
examples  
32 / 64 / 128 common

Thank you