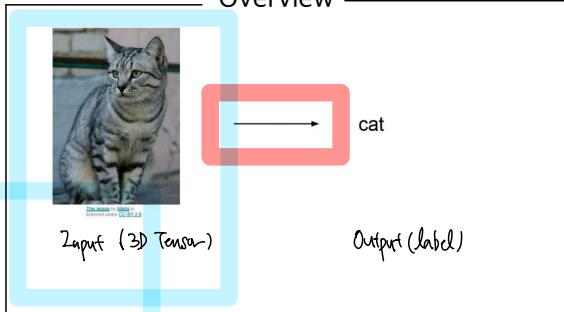




# Example Problem: Image Classification

## Overview



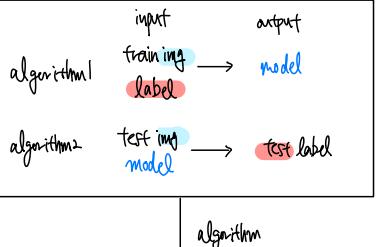
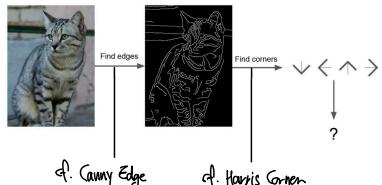
Problem: semantic gap

3D Tensor (오자). 실제적 차원 감소

Challenge

Viewpoint variation	
illumination	
background clutter	
occlusion	
deformation	
intraclass variation	
context	

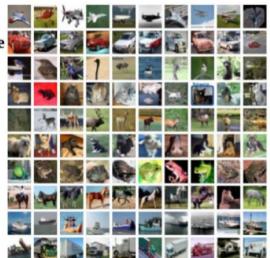
Solutions(Past)  
= hard coding



Solutions(ML)

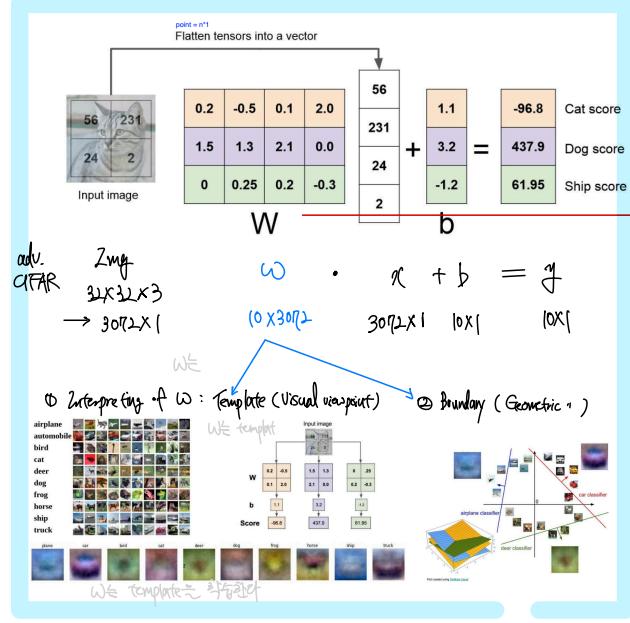
데이터 정리 → train → evaluate  
ML

CIFAR10



# Linear Classifier

Parametric approach:  
Linear Classifier



How about  
non-linear cases?

Class 1: 3 blue squares  
Class 2: 2 red circles  
Class 3: 2 green circles

Loss functions

Def.  $L = \frac{1}{N} \sum_i L_i(\hat{y}_i, y_i)$  s.t.  $f(x_i, y_i) \stackrel{i}{\sim} 1$   
 $= f(x, w)$



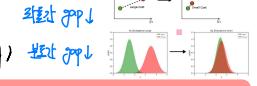
Ans: 3D tensor → 3 (class is integer label)

Do:  $(\hat{y}) \rightarrow \hat{y}_i \neq \hat{y}$   
 model output target

※: wit 액스터이 모델을 보여줄 때 정류하는 틈에 드는 것처럼 예상되는 결과를 보여주는 것

Loss  $\frac{1}{N} \sum_i$  를 만족해라 ① Target = 정답(y-label)

② Target =  $\begin{cases} 1 & (\text{정답인 } f(y)) \\ 0 & (\text{나머지 }) \end{cases}$  Softmax!



## Multiclass SVM Loss

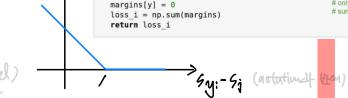
$y_{ij} \geq y_{ij} + 1$

$$L_i = \sum_{j \neq y_i} \max(0, y_{ij} - y_{iy} + 1)$$

$\therefore y_{ij} \geq y_{iy} \Rightarrow$   $y_{ij} \geq y_{iy}$   $\forall j \neq y_i$   $\forall i$

 $\therefore L_i = 0$

$y_{ij} \geq y_{ij} + 1$



exercise



cat	3.2	1.3	2.2
car	5.1	4.9	2.5
frog	-1.7	2.0	-3.1

cat:  $L_i = \max(0, 5 - 3.2 + 1) + \max(0, 7 - 3.2 + 1)$   
 $= \max(0, 2.9) + \max(0, -3.9)$   
 $= 2.9$

car:  $L_i = 0 \because (\text{모든 } j \neq 1 \text{일 때 } y_{ij} \leq 1)$

frog:  $L_i = (5 - 3 + 1) + (7 - 6 + 1) = 2.9$

$\therefore L = \frac{1}{3} (L_1 + L_2 + L_3) = \frac{1}{3} (5.9) \approx 1.97$

trivial question

Q1. car loss of -0.5?  $\rightarrow$  Losses must be positive

Q2. SVM Loss  $L_i = \min / \max$

Q3.  $\approx 0$  in all Loss

Q4.  $y_{ij} = y_i$ : case when all margins are positive

Q5.  $\approx$  card mean  $\approx$  margin size

Q6. What if we used

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

A1.  $\min X$   
 $\max(0, \cdot) = 0 \therefore (\text{모든 } j \neq 1 \text{일 때 } y_{ij} \leq 1)$

A2.  $\min = 0, \max = \infty$

A3.  $C - 1$   
 $\therefore C - 1$

A4. Just adding 1  $\therefore \sum \max(0, y_j - y_i + 1) = \max(0, 1) = 1$

A5. Just Rescaling

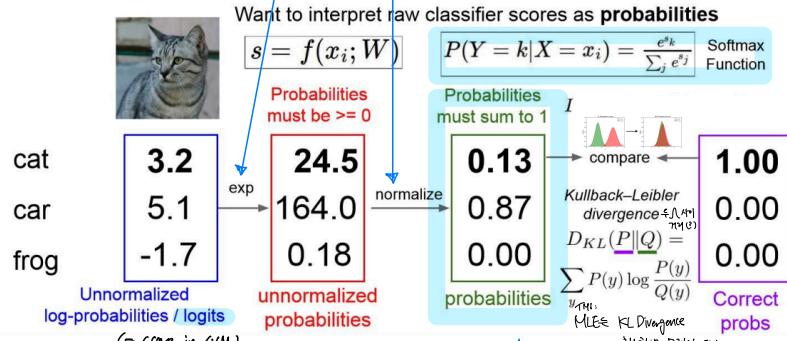
A6.  $\approx$  card loss + Ekt (squared hinge loss)

## 2. Softmax classifier a.k.a. multinomial logistic regression

### Motivation

장수는 확률로 표기

In statistics, by def.  $\sum_i s_i = 1$ ,  $s_i \in [0, 1]$



Putting it all together:

$$L_i = -\log \left( \frac{e^{s_i}}{\sum_j e^{s_j}} \right)$$

\*!  $-\log$ 은 확률이 아님. ① 정상 확률 확정  $\rightarrow$  loss 최소화

$-\log p$ 은 확률 gradient (loss 정의에 있음)  
 $\neq$  sensitive gradient, 확률  $\rightarrow$  확률

Def. Putting it all together:

$$L_i = -\log \left( \frac{e^{s_i}}{\sum_j e^{s_j}} \right)$$

trivial questions.

Q1. softmax loss  $L_i = -\ln / \max$

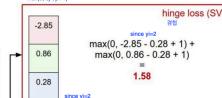
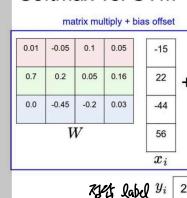
Q2. if  $s_i$  is equal, C classes of  $L_i$  are same

A1. 0/∞

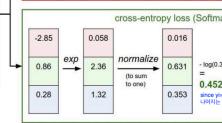
A2.  $-\log \left( \frac{1}{C} \right) = \log C$

SVM vs softmax.

Softmax vs. SVM



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



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

questions

assume scores:

[10, -2, 3]

[10, 9, 9]

[10, -100, -100]

and  $y_i = 0$

Q1. SVM loss vs. Softmax loss  $\frac{1}{N} \sum L_i$   $\approx$  1.14.

$$\text{SVM loss } L = \frac{1}{N} \sum L_i = \frac{1}{3} (0+0+0) = 0$$

$$\text{Softmax } L = \frac{1}{N} \sum L_i = \frac{1}{3} \left( -\log \left( \frac{e^{10}}{e^{10} + e^{-2} + e^3} \right), -\log \left( \frac{e^{-2}}{e^{10} + e^{-2} + e^3} \right), -\log \left( \frac{e^3}{e^{10} + e^{-2} + e^3} \right) \right) \approx 0$$

$$\begin{aligned} \text{assume scores:} \\ [10, -2, 3] &\rightarrow \frac{e^{10} e^3 e^{-2}}{e^{10} + e^{-2} + e^3} \\ [10, 9, 9] &\rightarrow \frac{e^{10} e^9 e^9}{e^{10} + e^{-2} + e^3} \\ [10, -100, -100] &\rightarrow \frac{e^{10} e^{-100} e^{-100}}{e^{10} + e^{-2} + e^3} \end{aligned}$$

Q2. sum of log probabilities  $\sum_j s_j$  vs. loss  $\sum_j -s_j \log p_j$   
 (soft label  $a_2 \times 2$ )

$$\begin{aligned} \text{assume scores:} \\ [10, -2, 3] \\ [10, 9, 9] \\ [10, -100, -100] \\ \text{and } y_i = 0 \end{aligned}$$

SUM Loss = 0  $\Rightarrow$  0

Softmax  $\approx 0$  (sensitive);  $\approx 1$  정상 확률

# More on loss functions: Regression

Intro.

Ex 1. loss function은 모델 예측과 실제의 차이를 정량화하는 것이다.

Ex 2. MLP의 opt loss function은 최소화된다.

Ex 3. output은 distribution이거나 가정된다.

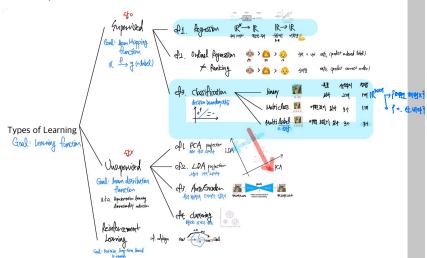
Parameter  $w$ 는  $\hat{w}_M$ 을 찾는다. ...  $\rightarrow$  모델에 맞는 최적의  $w$ 를 찾는다.

$$\begin{aligned}\hat{w}_M &= \underset{w}{\operatorname{argmax}} P_{\text{model}}(y|x, w) \\ &= \underset{w}{\operatorname{argmax}} \prod_{i=1}^N P_{\text{model}}(y_i|x_i, w) \\ &= \underset{w}{\operatorname{argmax}} \sum_{i=1}^N \log(P_{\text{model}}(y_i|x_i, w)) \\ &\quad \text{log-likelihood}\end{aligned}$$

independently and identically distributed

log-likelihood, 주어진 세 개의 예측을 기반으로 확률  
 $w$ 가 어떤  $w$ 를 선택하는지 찾는다.

Recap



Gaussian distribution:

L2 Loss

$$\text{def. } p(y|x, w) = \frac{1}{\sqrt{2\pi\sigma^2}} \cdot \exp\left(-\frac{(y - f_w(x))^2}{2\sigma^2}\right)$$

mean:  $\mu$ , std:  $\sigma$ , thin tails  
↳ penalize outliers  $\rightarrow$

L2 loss  $\rightarrow$

$$\begin{aligned}\hat{w}_M &= \underset{w}{\operatorname{argmax}} \sum_{i=1}^N \log P_{\text{model}}(y_i|x_i, w) \\ &= \underset{w}{\operatorname{argmax}} \sum_{i=1}^N \log\left(\frac{1}{\sqrt{2\pi\sigma^2}} e^{-\frac{(y_i - f_w(x_i))^2}{2\sigma^2}}\right) \\ &= \underset{w}{\operatorname{argmax}} \left(\sum_{i=1}^N \frac{1}{2} (\log \sigma^2) + \sum_{i=1}^N -\frac{(y_i - f_w(x_i))^2}{2\sigma^2}\right)\end{aligned}$$

$$= \underset{w}{\operatorname{argmax}} -\sum_{i=1}^N (y_i - f_w(x_i))^2$$

$w$ 를 찾는다.

$\hat{w}_M$ 은 무엇인가?

등수 최적화법:

$$= \underset{w}{\operatorname{argmin}} \sum_{i=1}^N (f_w(x_i) - y_i)^2$$

L2 loss  $\rightarrow$  square loss

Laplacian distribution:

L1 Loss

$$\text{def. } p(y|x, w) = \frac{1}{2\sigma} \exp\left(-\frac{|y - f_w(x)|}{\sigma}\right)$$

$\mu$ : location,  $\sigma$ : scale, heavy tails  
↳ outliers penalized.

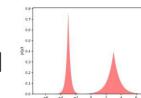
L1 loss  $\rightarrow$

$$\begin{aligned}\hat{w}_M &= \underset{w}{\operatorname{argmax}} \sum_{i=1}^N \log\left(\frac{1}{2\sigma} e^{-\frac{|y_i - f_w(x_i)|}{\sigma}}\right) \\ &= \underset{w}{\operatorname{argmax}} \left(\sum_{i=1}^N -\log \sigma + \sum_{i=1}^N \frac{|y_i - f_w(x_i)|}{\sigma}\right) \\ &= \underset{w}{\operatorname{argmax}} \sum_{i=1}^N |f_w(x_i) - y_i| \\ &= \underset{w}{\operatorname{argmin}} \sum_{i=1}^N |f_w(x_i) - y_i|\end{aligned}$$

L1 loss  
absolute loss  
more robust than L2

Multimodal distribution:

a.k.a. laplace distribution \* N



$$\text{def. } p_{\text{model}}(y|x, w) = \sum_{m=1}^M \tau_w^{(m)}(x) \frac{1}{2g_w^{(m)}(x)} \exp\left(-\frac{|y - f_w^{(m)}(x)|}{2g_w^{(m)}(x)}\right)$$

$\rightarrow$  X.

# More on loss functions: classification

Intro.

predict 1 value sigmoid, 1 value softmax?

C1 BCE Binary Cross Entropy loss

C2 CE Cross Entropy loss

Background: 2 classes

Bernoulli Distribution

$$p(y) = \mu^y (1-\mu)^{1-y}$$

Reg label  
 $\begin{cases} 0 & p(0) = 1-\mu \\ 1 & p(1) = \mu \end{cases}$

$$\hat{w}_{\text{HL}} = \underset{\omega}{\operatorname{argmax}} \sum_{i=1}^n \log P_{\text{Model}}(y_i | x_i, \omega)$$

$y_i \in \{0, 1\}$

$$\begin{aligned} &= \underset{\omega}{\operatorname{argmax}} \sum_{i=1}^n \log (\hat{f}_\omega(x_i)^{y_i} \cdot (1 - \hat{f}_\omega(x_i))^{1-y_i}) \\ &= \underset{\omega}{\operatorname{argmax}} \sum_{i=1}^n \log \hat{f}_\omega(x_i)^{y_i} + \sum_{i=1}^n \log (1 - \hat{f}_\omega(x_i))^{1-y_i} \\ &= \underset{\omega}{\operatorname{argmax}} \sum_{i=1}^n [y_i \cdot \log \hat{f}_\omega(x_i) + (1-y_i) \cdot \log (1 - \hat{f}_\omega(x_i))] \\ &= \underset{\omega}{\operatorname{argmin}} \sum_{i=1}^n [-y_i \cdot \log \hat{f}_\omega(x_i) - (1-y_i) \cdot \log (1 - \hat{f}_\omega(x_i))] \end{aligned}$$

BCE loss  
 $\leftarrow$  1가지 예측 가능  
 Loss 최소화 시키는 그때 가짜 예측  
 argmin으로 바꿔야 함

도수에 대해?  
 $(0, 1)$

Logistic Regression

각 주제를 때마다 같은 확률 분포로 바꿔야 함  
 $(0, 1)$  평가 by Sigmoid

$$p(y|x) = \begin{cases} h(x) = \sigma(\omega^T x + b) & \text{if } y=1 \\ 1-h(x) & \text{if } y=0 \end{cases}$$

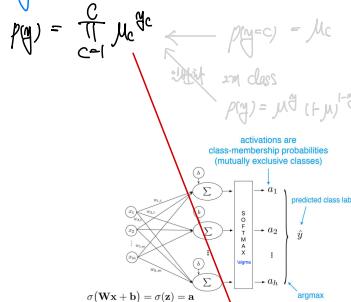
이제 각 주제에 맞게 바꿔야 함

$$= \frac{1}{1+e^{-\omega^T x - b}}$$

Min. von  $P(y_0, y_1 | x_0, x_1) = \prod_{i=1}^n p(y_i | x_i)$   $\rightarrow$  log loss

Main: Multiple Classes

Categorical Distribution



Example

$$\mathcal{L}^{(1)} = [(-1) \cdot \log(0.3792) + (-0) \cdot \log(0.3104) + (-0) \cdot \log(0.3104)] = 0.969692...$$

$$\mathbf{Y}_{\text{model}} = \begin{bmatrix} 1 & 0 & 0 \\ 0 & 0 & 1 \\ 0 & 0 & 0 \end{bmatrix}, \text{ Softmax outputs} = \begin{bmatrix} 0.3792 & 0.3104 & 0.3104 \\ 0.2698 & 0.4147 & 0.3156 \\ 0.2698 & 0.2978 & 0.4541 \end{bmatrix}$$

$$\begin{aligned} \mathcal{L}^{(2)} &= [(-1) \cdot \log(0.3792) + (-1) \cdot \log(0.3104) + (-1) \cdot \log(0.3104)] \\ \mathcal{L}^{(3)} &= [(-1) \cdot \log(0.2698) + (-1) \cdot \log(0.4147) + (-1) \cdot \log(0.2790)] \\ \mathcal{L}^{(4)} &= [(-1) \cdot \log(0.2698) + (-1) \cdot \log(0.3104) + (-1) \cdot \log(0.3104)] \\ \mathcal{L}^{(5)} &= [(-1) \cdot \log(0.2698) + (-1) \cdot \log(0.3104) + (-1) \cdot \log(0.3104)] \\ &\approx 0.9335 \end{aligned}$$

Softmax 누적 연속성

$$P(Y=k|X=x_i) = \frac{e^{x_k}}{\sum_j e^{x_j}} \text{ Softmax Function}$$

$$\text{softmax}(\lambda) = \left( \frac{e^{\lambda_1}}{e^{\lambda_1} + e^{\lambda_2}}, \dots, \frac{e^{\lambda_n}}{e^{\lambda_1} + e^{\lambda_n}} \right)$$

let  $\lambda_0 = 0$  (to degree of freedom = 1)

$$= \left( \frac{e^{\lambda_1}}{1+e^{\lambda_1}}, \dots, \frac{e^{\lambda_n}}{1+e^{\lambda_1}} \right)$$

$$b(\lambda) = \frac{1}{1+e^{-\lambda}}$$

$$= (b(\lambda_1), 1 - b(\lambda_1))$$

sigmoid function은 multi-class 확장 가능한 것



$$\text{softmax}(s) = \left( \frac{e^{s_1}}{\sum_i e^{s_i}}, \dots, \frac{e^{s_n}}{\sum_i e^{s_i}} \right) \rightarrow \text{softmax}(x) = \text{softmax}(x = \max_{i=1,2,3} s_i)$$

softmax(x) = softmax( $x = \max_{i=1,2,3} s_i$ )

• 문제: 이동평균을 사용하는  
 즉시 미리 확률을 예상할 때 softmax 사용.

# Regularization

Motivation:  
what is good W?  $w_1$  vs  $w_2$

example: SVM Loss

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

cat	3.2	1.3	2.6
car	5.1	4.9	9.8
frog	-1.7	2.0	4.0
Losses:	2.9		

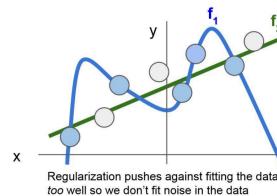
$$\begin{aligned} \textcircled{1} \text{ for } w_1, \\ L_{w_1} &= \max(0, 1.3 - 4.9 + 1) \\ &\quad + \max(0, 2.0 - 4.9 + 1) \\ &= 0 \\ \textcircled{2} \text{ for } w_2, \\ L_{w_2} &= \max(0, 2.6 - 9.8 + 1) \\ &\quad + \max(0, 4.0 - 9.8 + 1) \\ &= 0 \end{aligned}$$

Goal of Regularization:  
prevent the model "doing too well"



Why regularize?

- Express preferences over weights
- Make the model *simple* so it works on test data
- Improve optimization by adding curvature



Understanding: term by term

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

Data loss
Regularization

Regularization Strength  
 $\gamma$  (hyperparameter)

TH1:  $\gamma$  越大  $\hat{w}$  越小  $\Leftrightarrow$  越能 simplify  $w$  by Occam's Razor

Simple Examples

L2 regularization:  $R(W) = \sum_k \sum_l W_{k,l}^2$  spread  $w$  亂子 cf. [0.25, 0.25, 0.25]

L1 regularization:  $R(W) = \sum_k \sum_l |W_{k,l}|$  集中 cf. [1, 0, 0, 0]

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