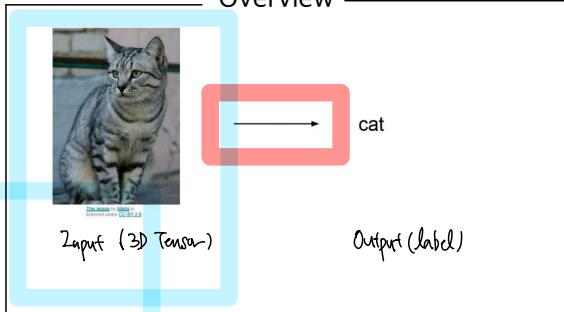


Example Problem: Image Classification

Overview



Problem: semantic gap

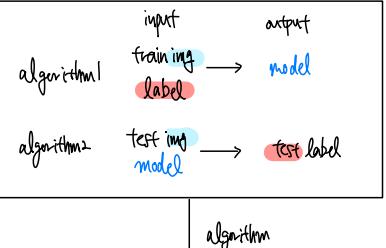
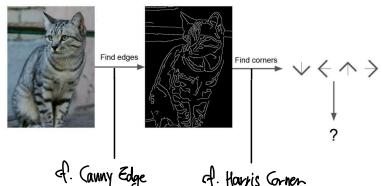
3D Tensor (오자). 실제적 차원 감소

Challenge

Viewpoint variation	
illumination	
background clutter	
occlusion	
deformation	
intraclass variation	
context	

Solutions(Past)

= hard coding



Solutions(ML)

데이터 정리 → train → evaluate
ML

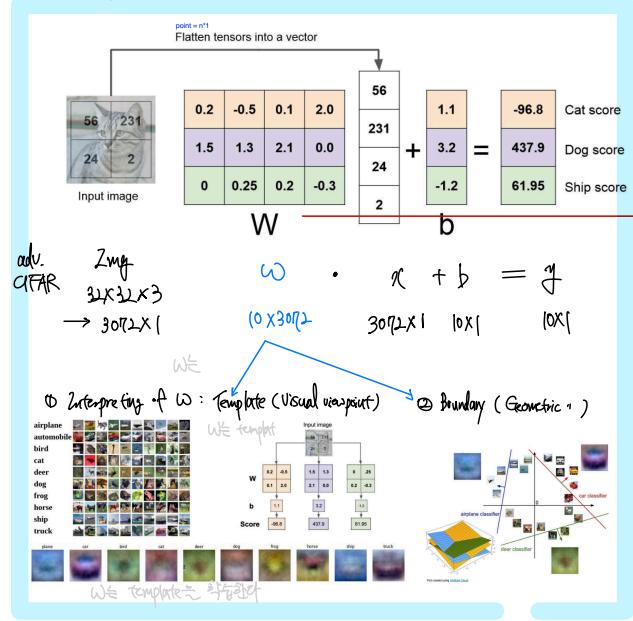
CIFAR10



Linear Classifier

Parametric approach:

Linear Classifier



How about non-linear cases?

Class 1: 10 blue squares
 Class 2: 10 red circles
 Class 3: 10 green dots

clf.

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



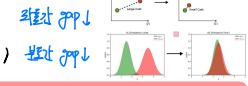
Ans: 3D tensor: 3 (3차원은 정수인 integer label)

Do: $(\hat{y}) \rightarrow \hat{y}_i = \hat{y}$

※: with 액스터치가 되어서도 되어줄

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

② Target = $\begin{cases} 1 & (\hat{y}_i \geq \hat{y}_j) \\ 0 & \text{otherwise} \end{cases}$



Multiclass SVM Loss

$$L_i = \sum_{j \neq y_i} \max(0, \hat{y}_j - \hat{y}_{y_i} + 1)$$

$\therefore \hat{y}_i \geq \hat{y}_j \forall j \neq i \Rightarrow L_i = 0$

exercise



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



$$\text{cat: } L_i = \max(0, 5.1 - 3.2 + 1) + \max(0, 2.5 - 3.2 + 1) = 2.9$$

$$\text{car: } L_i = 0 \quad \because \text{no other image is the target label.}\quad \text{frog: } L_i = (5.1 + 1) + (2.5 + 1) = 12.9$$

$$\therefore L = \frac{1}{3} (L_i + L_c + L_f) = \frac{1}{3} (9.9) \approx 9.9$$

trivial question

Q1. car loss of -0.5199
 Loss는 대체로 높아야 한다

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

Q3. ≈ 0 일 때 Loss

Q4. $y_i = y_j$ case에서 정답이 예상한 것과 같음

Q5. 5 칸 평균으로 예상한 것과 같음

Q6. What if we used

$$L_i = \sum_{j \neq y_i} \max(0, \hat{y}_j - \hat{y}_{y_i} + 1)$$

A1. 대체로 X
 $\max(0, \cdot) = 0 \quad \because \text{정답은 정답인 정답 label}$

A2. min = 0, max = infinity

A3. C-1
 $\stackrel{\text{0이면}}{=} \stackrel{\text{C-1인}}{=}$

A4. Just adding 1 $\because \sum \max(0, \hat{y}_j - \hat{y}_{y_i} + 1) = \max(0, 1) = 1$

A5. Just Rescaling

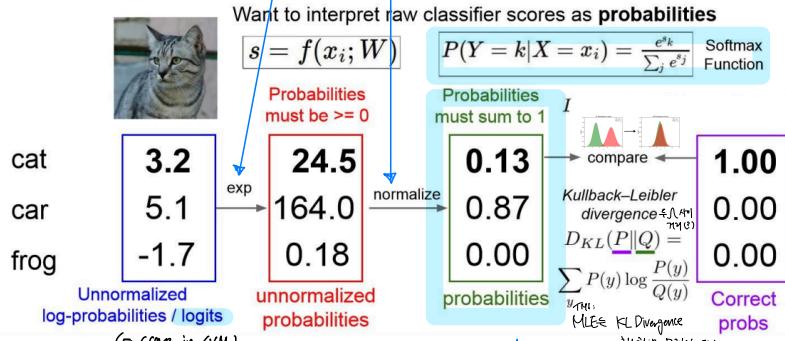
A6. 5 칸 loss + 5 칸 squared hinge loss

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

Motivation

장수는 확률로 표기

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



Score는 확률로

비교가능

loss를 정의해야지! : $-\log$ 은 확률

Putting it all together:

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

*(-log은 확률이 아님) ① 정상 확률 확정 → loss 최소화

-log은 확률 gradient (loss 정의가 있음)
-然是 sensitive하게 되어, 확률이 0에 가까울 때.

Def. Putting it all together:

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

trivial questions.

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

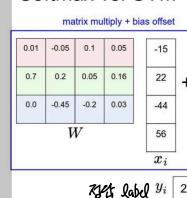
Q2. if s_i is equal, C classes of L_i are same

A1. 0/∞

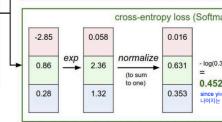
A2. $-\log \left(\frac{1}{C} \right) = \log C$

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?

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

$$\text{Softmax } L = \frac{1}{N} \sum_i L_i = \frac{1}{3} \left(-\log \left(\frac{e^{10}}{e^{10} + e^{-2} + e^3} \right), -\log \left(\frac{e^9}{e^{10} + e^{-2} + e^3} \right), -\log \left(\frac{e^9}{e^{10} + e^{-2} + e^3} \right) \right) \approx 0$$

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

Q2. sum of log 확률 vs. loss의 정의에 대한

(right label 0x2)

assume scores:

[10, -2, 3]

[10, 9, 9]

[10, -100, -100]

and $y_i = 0$

SUM Loss = 0 정의

Softmax 확률 (sensitive); ≈ 1 정상 확률

More on loss functions: Regression

Intro.

Ex 1. loss function은 미분 가능한 함수.

Ex 2. MLP의 opt loss function 구조화.

Ex 3. output은 distribution이거나 가정된다.

Parameter w is set by ...

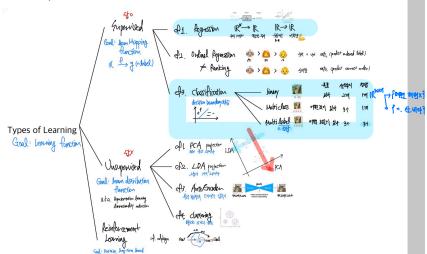
$$\hat{w}_{ML} = \underset{w}{\operatorname{argmax}} P_{\text{model}}(y|x, w)$$

$$= \underset{\substack{i \rightarrow j \\ \text{indep and identically distributed}}}{\operatorname{argmax}} \prod_{i=1}^N P_{\text{model}}(y_i|x_i, w)$$

$$= \underset{\substack{i \rightarrow j \\ \log}}{\operatorname{argmax}} \sum_{i=1}^N \log(P_{\text{model}}(y_i|x_i, w))$$

\log SIRE, 주어진 세 개의 확률을 같은 확률
 w 에 의해 하나의 확률로 고려되게 함

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: μ , std: σ , thin tails
↳ penalize outliers \propto

L2 loss \Leftrightarrow

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

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

w 를 찾는다.

\hat{w}_{ML} 은 무엇인가?

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

L2 loss square loss

L2 loss를 생각 \Leftrightarrow Gaussian 분포를 생각하는 것

Laplacian distribution:



L1 Loss

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

μ : location, b : scale, heavy tails
↳ outlier penalize:

L1 loss \Leftrightarrow

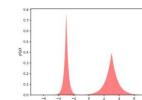
$$\begin{aligned} \hat{w}_{ML} &= \underset{w}{\operatorname{argmax}} \sum_{i=1}^N \log \left(\frac{1}{2b} e^{-\frac{|y_i - f_w(x_i)|}{b}} \right) \\ &= \underset{w}{\operatorname{argmax}} \left(\sum_{i=1}^N -\log b + \sum_{i=1}^N -\frac{|y_i - f_w(x_i)|}{b} \right) \\ &= \underset{w}{\operatorname{argmax}}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

L1 loss
absolute loss
more robust than L2



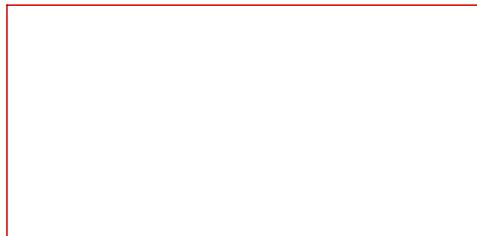
Multimodal distribution:

a.k.a. laplace distribution * N



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

즉 X. 유행이 전·세 번째



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

$$\begin{aligned} \hat{w}_{\text{BL}} &= \underset{\omega}{\operatorname{argmax}} \sum_{i=1}^n \log P_{\text{Model}}(y_i | x_i, \omega) \\ &= \underset{\omega}{\operatorname{argmax}} \sum_{i=1}^n \log (f_w(x_i)^{y_i} \cdot (1-f_w(x_i))^{1-y_i}) \\ &= \underset{\omega}{\operatorname{argmax}} \sum_{i=1}^n \log f_w(x_i)^{y_i} + \sum_{i=1}^n \log (1-f_w(x_i))^{1-y_i} \\ &= \underset{\omega}{\operatorname{argmax}} \sum_{i=1}^n [y_i \log f_w(x_i) + (1-y_i) \log (1-f_w(x_i))] \\ &= \underset{\omega}{\operatorname{argmin}} \sum_{i=1}^n [-y_i \log f_w(x_i) - (1-y_i) \log (1-f_w(x_i))] \end{aligned}$$

BCE loss
 \leftarrow 가능한 틀리지도
 loss 최소화 시킬 때 고려해
 argmin으로 바꿔야 함

Logistic Regression

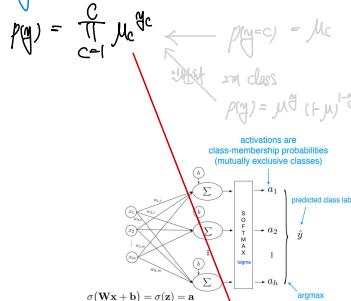
각 주제를 때마다 같은 확률 분포로 바꿔야 함
 (0,1) 범위 by Sigmoid

$$\begin{aligned} \text{2nd var} \quad p(y|x) &= \begin{cases} h(x) = \sigma(w^T x + b) & if y=1 \\ 1-h(x) & if y=0 \end{cases} \quad \text{이제 각 주제에 맞게 바꿔야 함} \\ &= \frac{1}{1+e^{-w^T x - b}} \quad \text{loss} \rightarrow \log \frac{p(y|x)}{1-p(y|x)} \\ \text{Min. von } p(y_0, y_1, y_2, y_3) &= \prod_{i=1}^4 p(y_i|x_i) \end{aligned}$$

y_0, y_1, y_2, y_3 alive

Main: Multiple Classes

Categorical Distribution



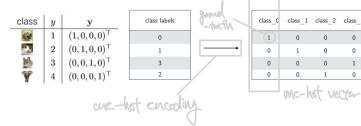
$$\begin{aligned} \hat{w}_{\text{ML}} &= \underset{\omega}{\operatorname{argmax}} \sum_{i=1}^n P_{\text{Model}}(y_i | x_i, \omega) \\ &= \underset{\omega}{\operatorname{argmax}} \sum_{i=1}^n \log \left(\prod_{c=1}^C f_w^{y_i}(x_i)^{y_i c} \right) \\ &= \underset{\omega}{\operatorname{argmin}} \sum_{i=1}^n \sum_{c=1}^C \log f_w^{y_i}(x_i)^{y_i c} \\ &= \underset{\omega}{\operatorname{argmin}} \sum_{i=1}^n \sum_{c=1}^C -y_i c \log f_w^{y_i}(x_i) \end{aligned}$$

CE loss

정답 (ground-truth) by count loss count. 나중에 0

TH1: one-hot vector

정답 class 번호에 대해서만 1을 갖는 vector



Example

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

$$\mathbf{Y}_{\text{model}} = \begin{bmatrix} 1 & 0 & 0 & 0 \\ 0 & 0 & 1 & 0 \\ 0 & 0 & 0 & 1 \end{bmatrix} \quad \text{Softmax outputs} = \begin{bmatrix} 0.3792 & 0.3104 & 0.3104 \\ 0.2698 & 0.2978 & 0.3351 \\ 0.2698 & 0.2978 & 0.3351 \end{bmatrix} \rightarrow \text{by model vector}$$

$$\begin{aligned} \mathcal{L}^{\text{BL}} &= (-1) \cdot \log(0.3792) + (-0) \cdot \log(0.3104) + (-0) \cdot \log(0.3104) \\ &= (-1) \cdot \log(0.3792) + (-1) \cdot \log(0.3104) + (-1) \cdot \log(0.3104) \\ &= 0.969692... \end{aligned}$$

$$\begin{aligned} \mathcal{L}^{\text{CE}} &= (-1) \cdot \log(0.3792) + (-1) \cdot \log(0.3104) + (-1) \cdot \log(0.3104) \\ &= (-1) \cdot \log(0.3792) + (-1) \cdot \log(0.3104) + (-1) \cdot \log(0.3104) \\ &= 0.9335 \end{aligned}$$

Softmax 누적 연습

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

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

let $\lambda_2 = 0$ (to degree of freedom = 1)

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

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

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

sigmoid function이 multi-class 확장될 것

$$\text{softmax}(S) = \left(\frac{e^{S_1}}{\sum_i e^{S_i}}, \dots, \frac{e^{S_c}}{\sum_i e^{S_i}} \right) \rightarrow \text{softmax}(x) = \text{softmax}(x = \max_{i=1,2,\dots,c} S_i)$$

TH1: soft argmax(x)
 정답은 1

• 문제: 이전步 충돌 argmax
 즉, 미지 확률로 예상 X가 softmax(1).

Regularization

Motivation:
what is good W? w_1 vs w_2

example: SVM Loss

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

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

$$f = w_1 x_{1b}$$

$$= w_1 x_{1b} + w_2 x_{2b}$$

① for w_1 ,

$$\begin{aligned} L_{w_1} &= \max(0, 1.3 - 4.9 + 1) \\ &\quad + \max(0, 2.0 - 4.9 + 1) \\ &= 0 \end{aligned}$$

② for w_2 ,

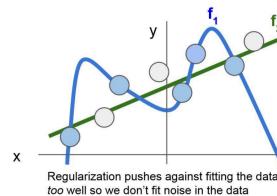
$$\begin{aligned} L_{w_2} &= \max(0, 2.6 - 9.8 + 1) \\ &\quad + \max(0, 4.0 - 9.8 + 1) \\ &= 0 \end{aligned}$$

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
 γ (hyperparameter)

TH1: γ 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, 0.25]

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

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

부록