

인 공 지 능

[심층학습 최적화 I]

본 자료는 해당 수업의 교육 목적으로만 활용될 수 있음.
일부 내용은 다른 교재와 논문으로부터 인용되었으며, 모든 저작권은 원 교재와 논문에 있음.

미리보기

■ 과학 혹은 공학에서 최적화

- 예) 우주선의 최적궤도, 운영체제의 작업 할당 계획 등

■ 기계학습의 최적화도 매우 복잡함

- 훈련집합으로 학습을 마친 후, 현장에서 발생하는 새로운^{unknown} 샘플을 잘 예측해야 함
 - 즉, 일반화^{generalization} 능력이 좋아야 함
 - 훈련집합은 전체 데이터 (실제, 알 수 없음) 대리자 역할
 - 검증집합은 테스트집합 대리자 역할
 - MSE, log-likelihood 등의 손실함수는 주어진 과업의 학습 성능(=판단 기준) 대리자 역할

■ 기계 학습의 최적화가 어려운 이유

- 대리자 관계
- 매개탐색 공간에서 목적함수의 비볼록^{non-convex} 성질, 고차원 특징 공간, 데이터의 희소성 등
- 긴 훈련 시간

■ 해당 장은 깊은 신경망 최적화에 효과적인 방안을 제시함

5.1 목적함수: 교차 엔트로피와 로그우도

- 5.1.1 평균제곱 오차를 다시 생각하기
- 5.1.2 교차 엔트로피 목적함수
- 5.1.3 소프트맥스^{softmax} 함수와 로그우도 목적함수

5.1 목적함수: 교차 엔트로피와 로그우도

■ 학습 과정의 성능 판정의 중요성

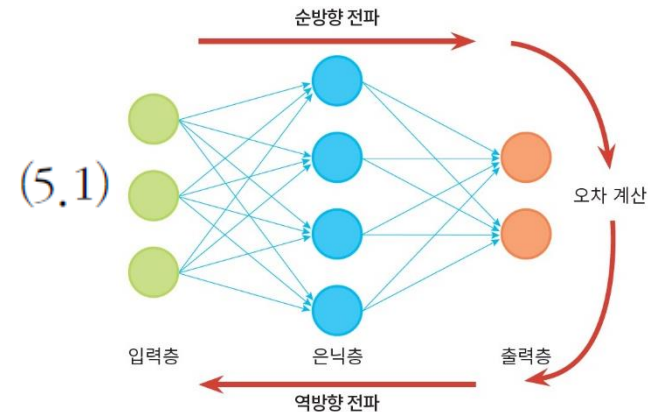
시험에서는 틀린 만큼 합당한 벌점을 받는 것이 중요하다. 그래야 다음 시험에서 심기일전으로 공부하여 틀리는 개수를 줄일 가능성이 크기 때문이다. 틀린 개수에 상관없이 비슷한 벌점을 받는다면 나태해져 성적을 올리는 데 지연이 발생할 것이다. 이러한 원리가 기계 학습에도 적용될까?

5.1.1 평균제곱 오차 다시 생각하기

■ 평균제곱 오차(MSE) 목적함수

$$e = \frac{1}{2} \|\mathbf{y} - \mathbf{o}\|_2^2$$

- 오차가 클수록 e 값이 크므로 벌점(정량적 성능)으로 활용됨



■ 하지만, 큰 **허점**이 존재

- 왼쪽 상황은 $e = 0.2815$,
- 오른쪽 상황은 $e = 0.4971$ 이므로 오른쪽이 더 큰 벌점을 받아야 마땅함

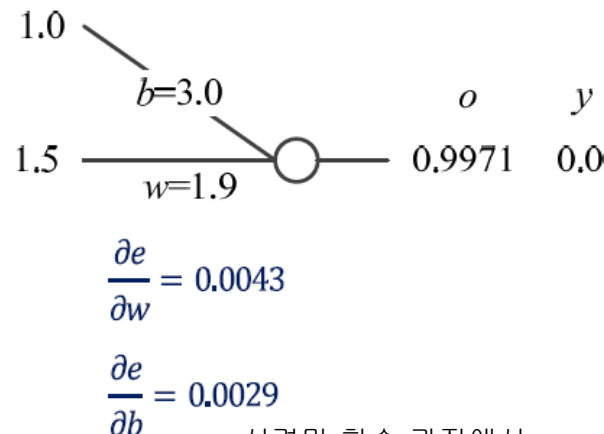
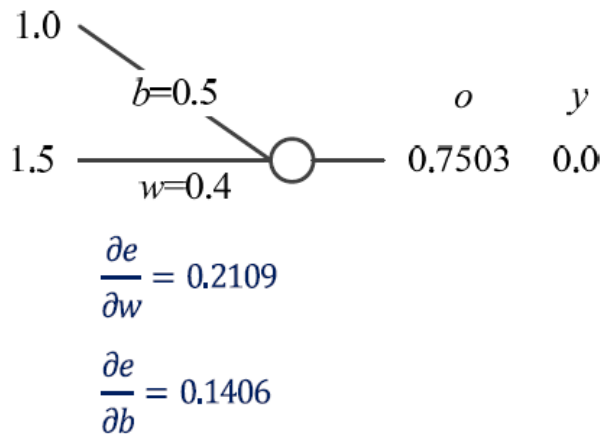


그림 5-1 MSE가 목적함수로서 부적절한 상황

신경망 학습 과정에서
 학습은 오류를 줄이는 방향으로 가중치와 편향을 교정
 ← 큰 교정이 필요함에도 작은 경사도로 작게 갱신됨

5.1.1 평균제곱 오차 다시 생각하기

■ 큰 허점

- 식 (5.3)의 경사도 gradient가 벌점에 해당

$$e = \frac{1}{2}(y - o)^2 = \frac{1}{2}(y - \sigma(wx + b))^2 \quad (5.2)$$

$$\left. \begin{aligned} \frac{\partial e}{\partial w} &= -(y - o)x\sigma'(wx + b) \\ \frac{\partial e}{\partial b} &= -(y - o)\sigma'(wx + b) \end{aligned} \right\} \quad (5.3)$$

- 경사도를 계산해보면 왼쪽 상황의 경사도가 더 큼
→ 더 많은 오류를 범한 상황이 더 낮은 벌점을 받은 꼴
→ 학습이 더딘 부정적 효과

■ 이유

- $w x + b$ (그림의 가로축에 해당)가 커지면 경사도가 작아짐

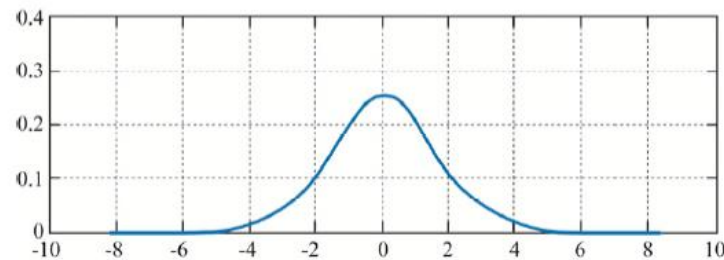
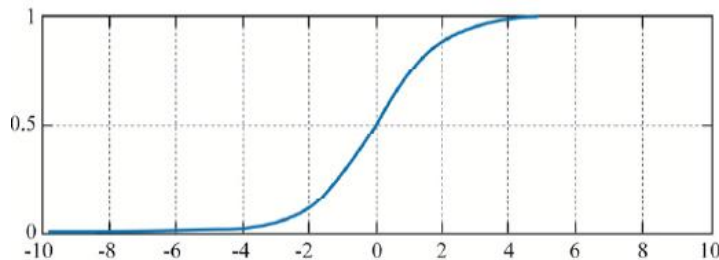


그림 5-2 로지스틱 시그모이드함수와 도함수

5.1.2 교차 엔트로피 목적함수

■ 교차 엔트로피 cross entropy

- 정답^{label}에 해당하는 y 가 확률변수 (부류가 2개라고 가정하면 $y \in \{0,1\}$)
- 확률 분포: P 는 정답, Q 는 신경망 (예측) 출력

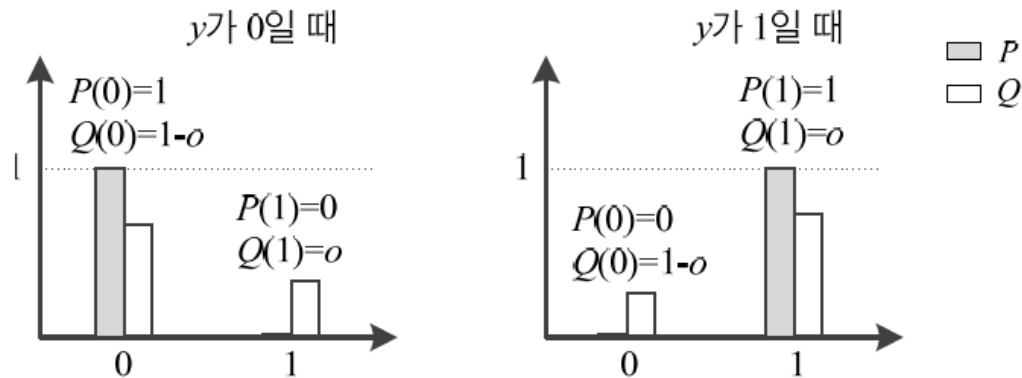


그림 5-3 레이블 y 가 0일 때와 1일 때의 P 와 Q 의 확률분포

- 확률분포를 통일된 수식으로 쓰면,

$$\begin{aligned} P(0) &= 1 - y & Q(0) &= 1 - o & \text{신경망 출력으로 표기하면,} \\ P(1) &= y & Q(1) &= o & o = \sigma(z) \text{이고 } z = wx + b \end{aligned}$$

- 식 (2.47) 교차 엔트로피 $H(P, Q) = -\sum_x P(x) \log_2 Q(x) = -\sum_{i=1, \dots, k} P(e_i) \log_2 Q(e_i)$ 을 적용
 $\rightarrow H(P, Q) = -\sum_{y \in \{0,1\}} P(y) \log_2 Q(y)$

5.1.2 교차 엔트로피 목적함수

■ 간단한 예

- $P(x)$ 와 $Q(x)$ 의 값을 극단적인 경우로 가정 1 혹은 0으로 단순화 가정
- 교차 엔트로피 $H(P, Q) = -\sum P(x)\log Q(x)$ 를 구하면,

잘못된 학습으로 오분류된 손실은

$$-[1 \ 0] \begin{bmatrix} \log 0 \\ \log 1 \end{bmatrix} = -(-\infty + 0) = \infty$$

잘된 학습으로 제대로 분류된 손실은

$$-[1 \ 0] \begin{bmatrix} \log 1 \\ \log 0 \end{bmatrix} = -(0 + 0) = 0$$

5.1.2 교차 엔트로피 목적함수

■ 교차 엔트로피 목적함수

$$e = -(y \log_2 o + (1 - y) \log_2(1 - o)), \quad \text{이때, } o = \sigma(z) \text{이고 } z = wx + b \quad (5.4)$$

■ 역할을 잘 수행하는지 확인

- y 가 1, o 가 0.98일 때 (예측이 잘된 경우)
 - 오류 $e = -(1 \log_2 0.98 + (1 - 1) \log_2(1 - 0.98)) = 0.0291$ 로서 낮은 값
- y 가 1, o 가 0.0001일 때 (예측이 잘못된 경우, 혹은 오분류된 경우)
 - 오류 $e = -(1 \log_2 0.0001 + (1 - 1) \log_2(1 - 0.0001)) = 13.2877$ 로서 높은 값

5.1.2 교차 엔트로피 목적함수

■ 공정한 벌점을 부여하는지 확인 (MSE의 느린 학습 문제를 해결 확인)

- 도함수를 구하면,

$$\frac{\partial e}{\partial w} = -\left(\frac{y}{o} - \frac{1-y}{1-o}\right) \frac{\partial o}{\partial w}$$

$$= -\left(\frac{y}{o} - \frac{1-y}{1-o}\right) x \sigma'(z)$$

$$= -x \left(\frac{y}{o} - \frac{1-y}{1-o}\right) o(1-o)$$

$$= x(o - y)$$

식 (5.4) 대입, 연쇄법칙 적용

$(\log f(x))' = f'(x)/f(x)$

$$e = -(y \log_2 o + (1-y) \log_2 (1-o))$$

$$o = \sigma(z) \text{이고 } z = wx + b$$

$$\sigma'(z) = \sigma(z)(1 - \sigma(z)) = o(1 - o)$$

$$\left. \begin{aligned} \frac{\partial e}{\partial w} &= x(o - y) \\ \frac{\partial e}{\partial b} &= (o - y) \end{aligned} \right\} \quad (5.5)$$

- 경사도를 계산해 보면, 오류가 더 큰 오른쪽에 더 큰 벌점 (경사도) 부과

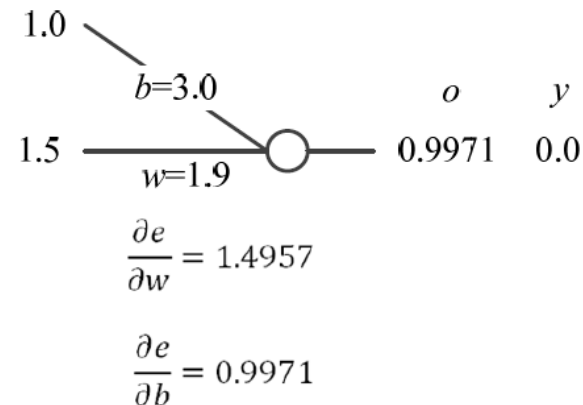
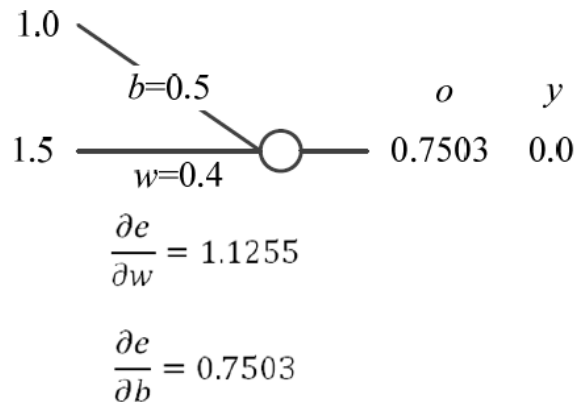


그림 5-4 교차 엔트로피를 목적함수로 사용하여 느린 학습 문제를 해결

5.1.2 교차 엔트로피 목적함수

■ 식 (5.4)를 c 개의 출력 노드를 가진 경우로 확장

- 출력 벡터 $\mathbf{o} = (o_1, o_2, \dots, o_c)^T$ 인 상황으로 확장 ([그림 4-3]의 DMLP)

$$e = - \sum_{i=1,c} (y_i \log_2 o_i + (1 - y_i) \log_2 (1 - o_i)) \quad (5.6)$$

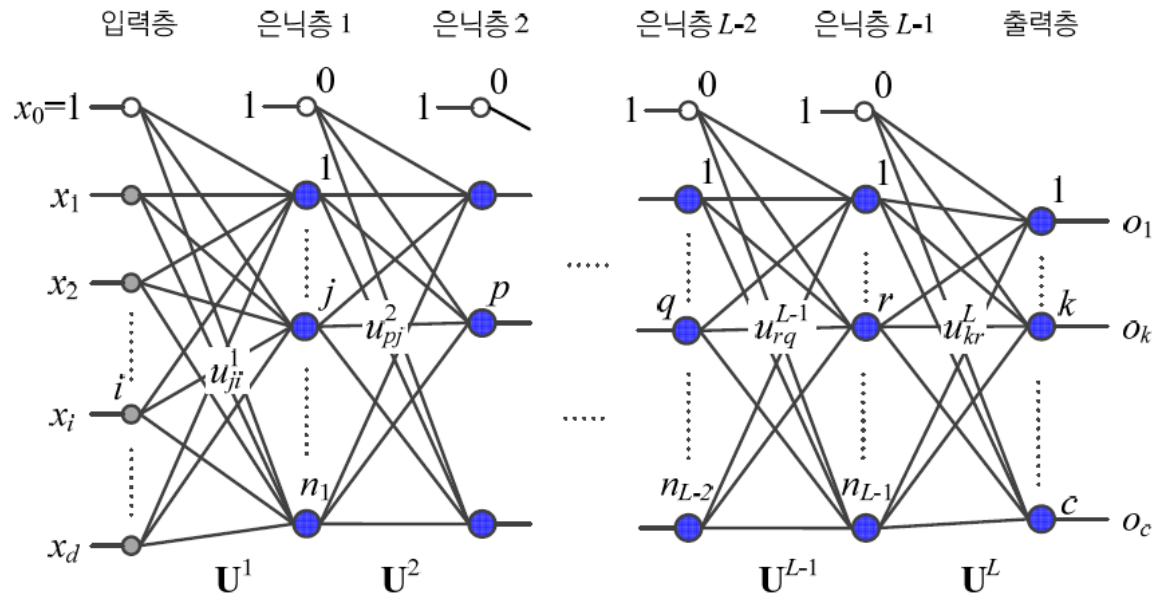


그림 4-3 깊은 MLP(DMLP)의 구조

- 예) $y \in \{\text{개, 고양이, 사람}\}$ 경우, $P(y) = [0, 0, 1]$ 와 $Q(y) = [0.2, 0.3, 0.5]$ 일 때, 교차 엔트로피는 $-\log(0.5)$

5.1.3 소프트맥스 활성함수와 로그우도 목적함수

■ 소프트맥스softmax 함수

$$o_j = \frac{e^{s_j}}{\sum_{i=1,c} e^{s_i}} \quad (5.7)$$

■ 동작 예

- 소프트맥스는 최대^{max}를 모방

← 출력 노드의 중간 계산 결과 s_i^L 의 최댓값을 더욱 활성화하고 다른 작은 값들은 억제

- 모두 더하면 1이 되어 확률 모방

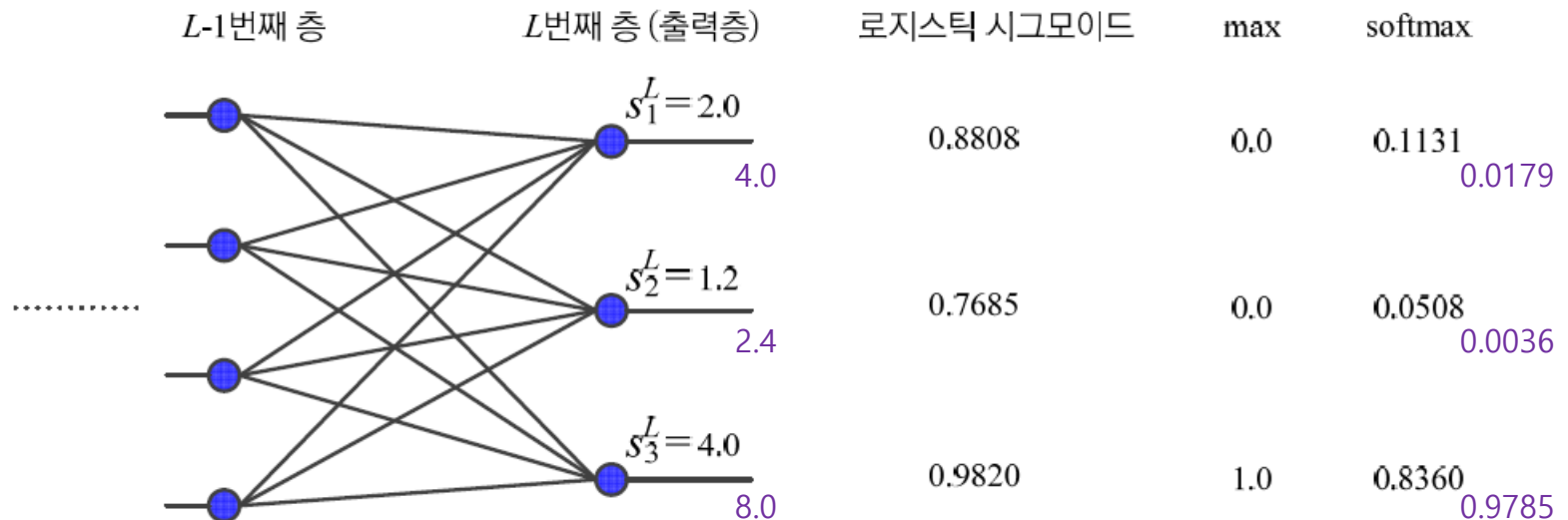
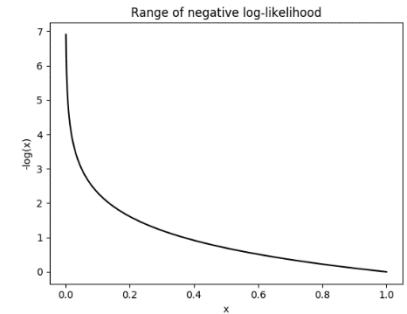


그림 5-5 출력층의 활성함수로 로지스틱 시그모이드와 softmax 비교

5.1.3 소프트맥스 활성화 함수와 로그우도 목적 함수

■ 음의 로그우도 negative log-likelihood 목적 함수

$$e = -\log_2 o_y \quad (5.8)$$



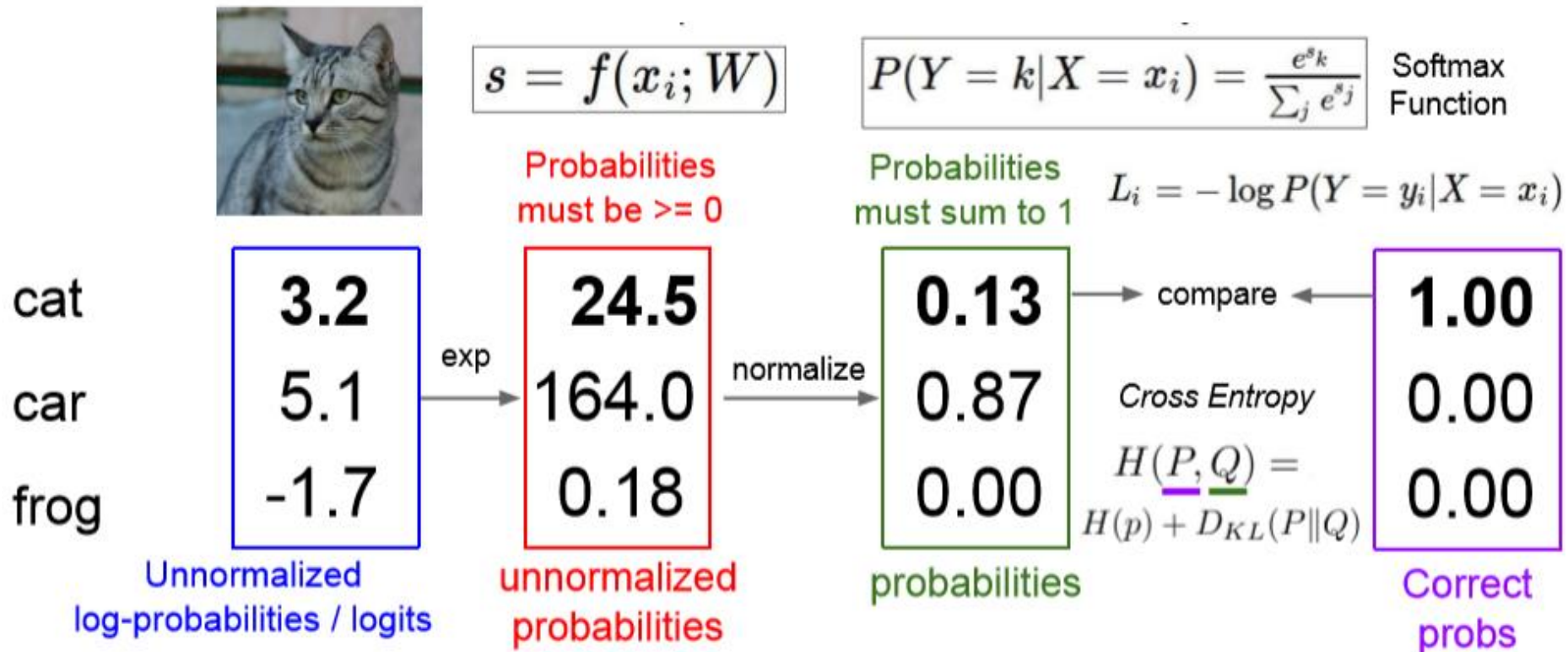
- 모든 출력 노드값을 사용하는 MSE나 교차 엔트로피와 달리 o_y 라는 하나의 노드만 적용
- o_y 는 샘플의 정답에 해당하는 노드의 출력값
 - 동작 예시1) [그림 5-5]에서 현재 샘플이 두 번째 부류라면 o_y 는 o_2
 $e = -\log_2 0.0508 = 4.2990$. 잘못 분류한 셈이므로 목적함수값이 큼
 - 동작 예시2) [그림 5-5]에서 현재 샘플이 세 번째 부류라면 o_y 는 o_3
 $e = -\log_2 0.8360 = 0.2584$. 제대로 분류한 셈이므로 목적함수값이 작음

■ 소프트맥스와 로그우도

- 소프트맥스는 최댓값이 아닌 값을 억제하여 0에 가깝게 만든다는 의도 내포
- 신경망에 의한 샘플의 정답에 해당하는 노드만 보겠다는 로그우도와 잘 어울림
- 따라서 둘을 결합하여 사용하는 경우가 많음

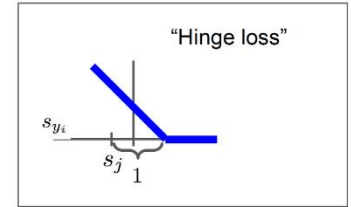
5.1.3 소프트맥스 활성화 함수와 로그우도 목적 함수

- 소프트맥스와 교차 엔트로피 목적 함수
 - 로그우도 손실 함수 ~ 교차엔트로피 최소화

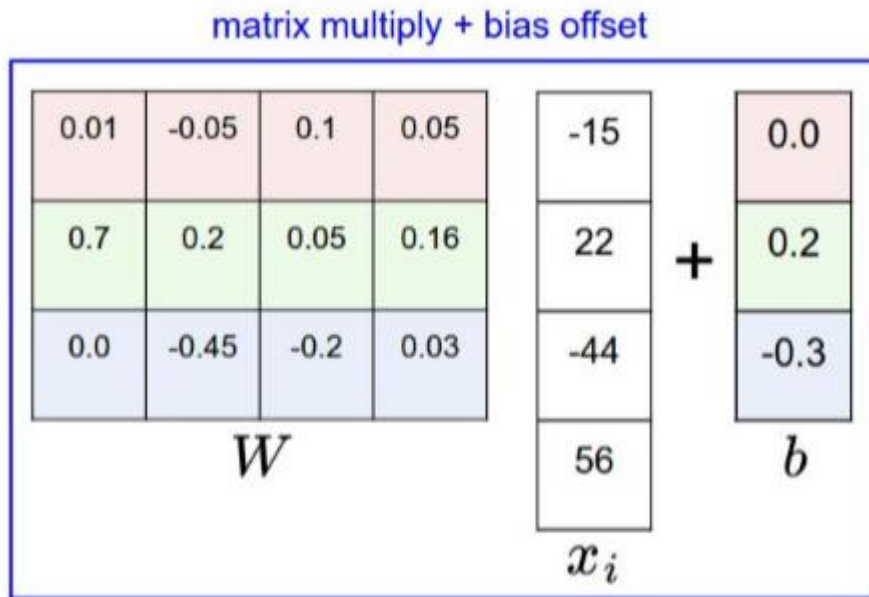


5.1.3 소프트맥스 활성화 함수와 로그우도 목적 함수

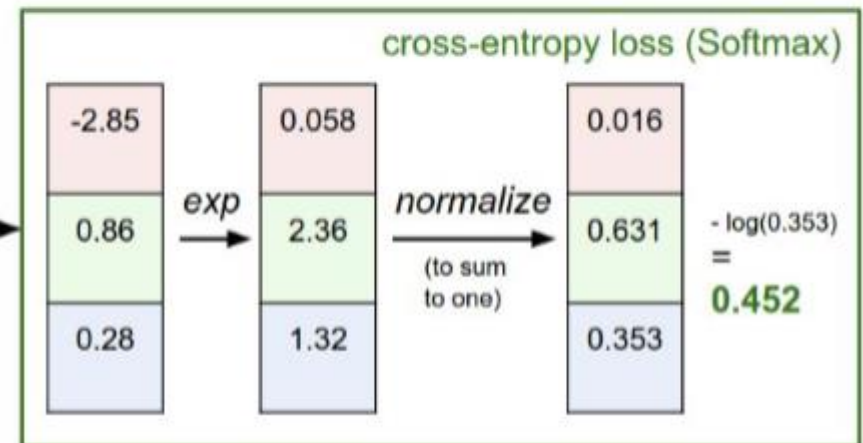
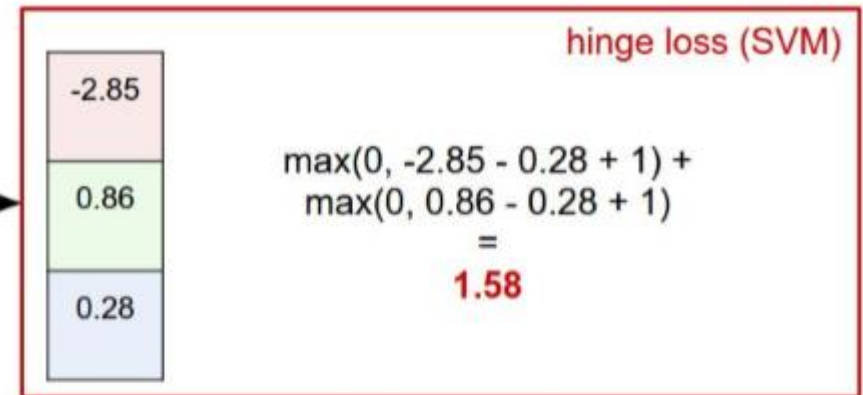
■ 소프트맥스와 힌지로스 hinge loss 출력 비교



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



분류 정답 y_i 2



5.1.3 소프트맥스 활성화 함수와 로그우도 목적 함수

■ 소프트맥스 분류기 | softmax classifier

- 다항 로지스틱 회귀 분석 multinomial logistic regression의 예
- 분류기의 최종 값을 확률로 표현
- 소프트맥스와 로그우도 목적 함수



$$s = f(x_i; W)$$

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

Probabilities
must be ≥ 0

Probabilities
must sum to 1

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

cat
car
frog

3.2
5.1
-1.7

Unnormalized
log-probabilities / logits

exp

24.5
164.0
0.18

unnormalized
probabilities

normalize

0.13
0.87
0.00

probabilities

$$\rightarrow L_i = -\log(0.13) = 2.04$$

Maximum Likelihood Estimation
Choose probabilities to maximize
the likelihood of the observed data



5.1.3 소프트맥스 활성화 함수와 로그우도 목적 함수

■ 미니배치 단위 예

Negative Log Likelihood (NLL) Loss

