

[ML_08] Multiclass ClassificationOne vs Rest Method

- 이진분류기 C개를 독립적으로 사용하여 class k에 나머지 C-1개 class를 분류
- Class k에 대한 이진분류기를 h_k 라 하면, $h_k(x)$ 가 참인 값을 찾는 것

→ 각 이진분류기에 대해 훈련셋의 불균형은 없음

One vs One Method

- 이진분류기 C(C-2)개를 독립적으로 사용하여 class k에 class l을 분류
- 가장 많은 이진분류기가 선택한 class를 최종 결과로 결정

→ 훈련셋의 불균형은 알리지 않음,

사용하는 이진분류기의 개수가 C에 비례함: 높은 복잡도를 가짐.

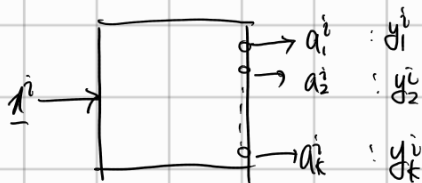
Softmax

- Softmax function

$$\sigma: \mathbb{R}^k \rightarrow (0,1)^k$$

$$\sigma(\underline{z})_i = \frac{e^{z_i}}{\sum_{j=1}^k e^{z_j}} \rightarrow \underline{z} = \begin{bmatrix} z_1 \\ z_2 \\ \vdots \\ z_k \end{bmatrix} \quad \sigma(\underline{z}) = \begin{bmatrix} e^{z_1} \\ e^{z_2} \\ \vdots \\ e^{z_k} \end{bmatrix} / (e^{z_1} + e^{z_2} + \dots + e^{z_k})$$

- Cross Entropy Loss



$$\Rightarrow \frac{1}{N} \sum_{n=1}^N \sum_{k=1}^K -y_k^i \log(a_k^i)$$

CE Loss at output layer

Cross-Entropy Loss

$$J(\omega) = \frac{1}{N} \sum_{n=1}^N \sum_{i=1}^K -y_i^n \log a_i^n$$

$$= \frac{1}{N} \sum_{n=1}^N \sum_{i=1}^K -y_i^n \log \frac{\exp(\omega_i^T \underline{x}^n)}{\sum_{j=1}^K \exp(\omega_j^T \underline{x}^n)}, \quad N = \text{1 샘플 개수}$$

$k = \text{인출층의 개수}$

Gradient of CE Loss

$$\frac{\partial J(\omega)}{\partial \omega_k} = \frac{1}{N} \sum_{n=1}^N \sum_{i=1}^K -y_i^n \frac{\partial}{\partial \omega_k} \log a_i^n$$

→ k번째 노드 출력
가중치의 변화

$$\rightarrow \frac{\partial}{\partial \omega_k} \log a_i^n = \frac{\partial \log a_i^n}{\partial z_k^n} \times \frac{\partial z_k^n}{\partial \omega_k}, \quad z_k^n = \sum_{i=1}^d \omega_{ki} x_i^n$$

$$= \frac{\partial \log a_i^n}{\partial z_k^n} \times \underline{x}^n = \underline{\omega}_k^T \underline{x}^n$$

$$\rightarrow \frac{\partial \log a_i^n}{\partial z_k^n} = \frac{\partial}{\partial z_k^n} \log \frac{\exp(\omega_i^T \underline{x}^n)}{\sum_{j=1}^K \exp(\omega_j^T \underline{x}^n)}$$

$$= \frac{\partial}{\partial z_k^n} \log \frac{\exp(z_i^n)}{\sum_{j=1}^K \exp(z_j^n)}$$

$$= \frac{\partial}{\partial z_k^n} (\log \exp(z_i^n) - \log \sum_{j=1}^K \exp(z_j^n))$$

$$= \{i=k\} - \frac{\exp(z_k^n)}{\sum_{j=1}^K \exp(z_j^n)} = \{i=k\} - a_k^n$$

$$\therefore \frac{\partial J(\omega)}{\partial \omega_k} = \frac{1}{N} \sum_{n=1}^N \sum_{i=1}^K -y_i^n \frac{\partial}{\partial \omega_k} \log a_i^n$$

$$= \frac{1}{N} \sum_{n=1}^N \sum_{i=1}^K -y_i^n (\{i=k\} - a_k^n) \underline{x}^n$$

$$= \frac{1}{N} \sum_{n=1}^N (-y_k^n + \sum_{i=1}^K y_i^n a_k^n) \underline{x}^n$$

$$= \frac{1}{N} \sum_{n=1}^N (a_k^n - y_k^n) \underline{x}^n$$

$\sum_{i=1}^K y_i^n a_k^n = a_k^n$

[ML-09] Overfitting & Regularization

Overfitting & Underfitting

underfitting (과소적합)

- 모델의 용량이 너무 작아 오차가 클 수밖에 없는 상황

underfitting 방지하는법

◦ 비선형모델 등과 같이 용량이 더 큰모델을 사용한다.

Overfitting (과잉적합)

- 용량이 너무커서 학습데이터에서 잡음까지 수용

→ 즉, 훈련집합에만 적합하게 되어 훈련집합에는 높은 성능을 보여주지만, 다른 새로운 데이터가 들어오게 되면 성능이 떨어짐.

Overfitting의 원인

- 1. 충분하지 않은 training dataset.
- 2. 너무 용량이 큰 parameter

Overfitting 해결법

- 1. 충분히 많은 training data를 사용한다.
→ 만약 training data가 적다면 데이터 증대 방법을 통해
같은 training data 증대 가능

2. 검증집합을 이용한 모델선택

- 훈련집합과 테스트 집합과 다른 별도의 검증집합을 준비한다.
- 검증집합에 대해 최의 성능을 보인 모델을 선택한다.

3. 규제(Regularization)

- 용량이 충분히 큰모델 + 다양한 feature 고려
+ 다양한 규제 기법을 적용

Regularization (규제)

$$J_{reg}(\theta) = \underbrace{J(\theta)}_{\text{목적함수}} + \lambda \underbrace{R(\theta)}_{\text{규제항}}$$

→ 규제항 R(θ)로 무엇을 사용할 것인가? → 가중치 감쇠

L2 Norm 사용 : $R(\theta) = \|\theta\|_2^2$

L1 Norm 사용 : $R(\theta) = \|\theta\|_1$

⇒ 최종계를 원할 때까지 당기는 효과

L2 Norm Regularization (Ridge)

$$J_{reg}(\theta; \mathbf{X}, \mathbf{y}) = J(\theta; \mathbf{X}, \mathbf{y}) + \lambda \|\theta\|_2^2$$

$$\nabla J_{reg}(\theta; \mathbf{X}, \mathbf{y}) = \nabla J(\theta; \mathbf{X}, \mathbf{y}) + 2\lambda\theta$$

- Parameter Update

$$\begin{aligned} \theta &= \theta - \rho \nabla J_{reg}(\theta; \mathbf{X}, \mathbf{y}) && \rightarrow L_2 \text{ 규제는 } \theta \text{를} \\ &= \theta - \rho (\nabla J(\theta; \mathbf{X}, \mathbf{y}) + 2\lambda\theta) && 2\lambda \text{의 배수로 } \theta \text{를} \\ &= (1 - 2\rho\lambda)\theta - \rho \nabla J(\theta; \mathbf{X}, \mathbf{y}) && \text{압축시키려는 성} \end{aligned}$$

L1 Norm Regularization (Lasso)

$$J_{reg}(\theta; \mathbf{X}, \mathbf{y}) = J(\theta; \mathbf{X}, \mathbf{y}) + \lambda \|\theta\|_1$$

$$\nabla J_{reg}(\theta; \mathbf{X}, \mathbf{y}) = \nabla J(\theta; \mathbf{X}, \mathbf{y}) + \lambda \text{sgn}(\theta)$$

- Parameter Update

$$\begin{aligned} \theta &= \theta - \rho \nabla J_{reg}(\theta; \mathbf{X}, \mathbf{y}) \\ &= \theta - \rho (\nabla J(\theta; \mathbf{X}, \mathbf{y}) + \lambda \text{sgn}(\theta)) \\ &= \theta - \rho \nabla J(\theta; \mathbf{X}, \mathbf{y}) - \underbrace{\rho \lambda \text{sgn}(\theta)}_{\text{penalty가 적용됨}} \end{aligned}$$

→ 0이 커져서
sgn이 0이
즉, 0으로 수렴하려는
경향을 가진다

⇒ L1 규제의 효과 : 0이 되는 가중치가 많이 발생

→ 즉 학습에 불필요한 도움이 되지 않는 가중치의 경우 0이 됨

→ 특징 선택 효과가 있음