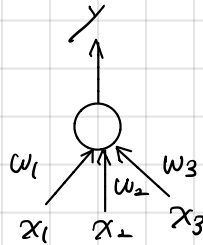


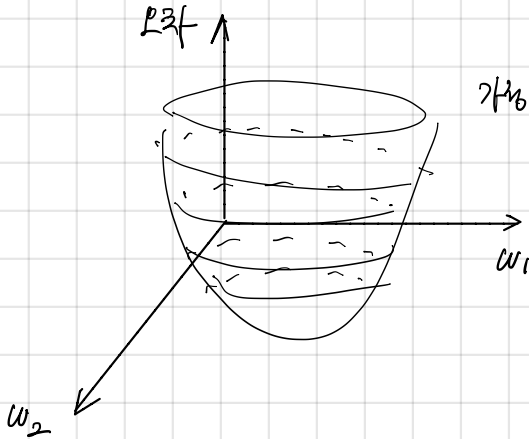
2.1 오차함수 (error function)

$$E = \frac{1}{2} \sum_i (t^{(i)} - y^{(i)})^2$$



$$y^{(i)} = w_1 x_1^{(i)} + w_2 x_2^{(i)} + w_3 x_3^{(i)}$$

이 뉴런의 출력



가중치 ① $\vec{w} = (w_1, w_2)$, ② 편향은 편

$$E = \frac{1}{2} \sum_{i=1}^n [t^{(i)} - (w_1 x_1 + w_2 x_2 + w_0)]^2$$

n : 해당 층의 뉴런 개수



이 뉴런의 출력을 다음 뉴런에

$$ax^2 + by^2 = z$$

→ 이게 뉴런이지...

xy 항은 뉴런끼리
서로 연결하는 거(-)

2.3 델타 규칙과 학습률 (delta rule, learning rate)

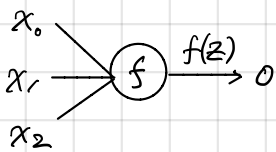
가중치 ← 현재가 + 학습률 (×, learning rate)

$$\Delta w_k = \frac{\partial E(\vec{w})}{\partial w_k} : \text{가중치에 대한 오차함수의 편미분값}$$

$$\begin{aligned} \Delta w_k &= - \epsilon \frac{\partial E}{\partial w_k} = - \epsilon \frac{\partial}{\partial w_k} \left[\frac{1}{2} \sum_{i=1}^n (t^{(i)} - y^{(i)})^2 \right] \\ &= \epsilon \sum_{i=1}^n \{ t^{(i)} - y^{(i)} \} \cdot \frac{\partial y^{(i)}}{\partial w_k} = \sum_{i=1}^n \epsilon x_k (t^{(i)} - y^{(i)}) \end{aligned}$$

2.4 시그모이드 뉴런의 경사 하강법

비선형 vs 선형



$$z = \sum_{k=1}^n w_k x_k \quad (n: \text{입력 뉴런 개수}) : \text{넷 입력값}$$

$$f(z) = \frac{1}{1 + e^{-z}} : \text{시그모이드 함수}$$

$$\Delta w_k = - \epsilon \frac{\partial E}{\partial w_k} = - \epsilon \frac{\partial}{\partial w_k} \left[\frac{1}{2} \sum_{i=1}^n (t^{(i)} - f(x)^{(i)})^2 \right] \quad (n: \text{한 층에 존재하는 뉴런의 개수})$$

$$= \epsilon \sum_{i=1}^n (t^{(i)} - f(x)) \cdot \frac{\partial f(x)}{\partial w_k} = \sum_{i=1}^n \epsilon x_k^{(i)} y^{(i)} (1 - y^{(i)}) (t^{(i)} - y^{(i)})$$

$$\frac{\partial f}{\partial w_k} = \left(\sum_{i=1}^n \frac{\partial f}{\partial z} \frac{\partial z}{\partial w_k} \right) = x_k \frac{e^{-z}}{(1 + e^{-z})^2} = x_k y (1 - y)$$

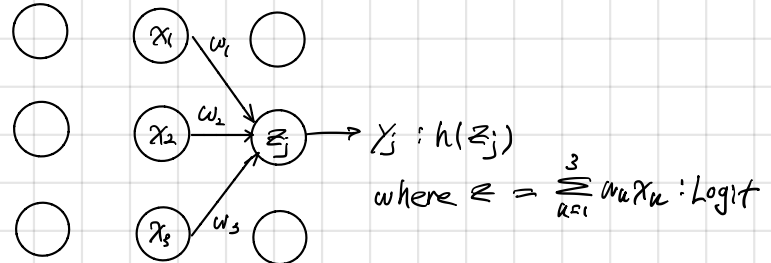
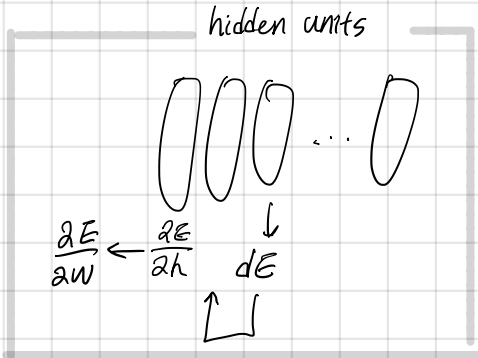


2.5 역전파 알고리즘 (backpropagation)

- hidden unit에서 활성화 함수를 바꿀 때 얼마나 빨리 오차가 변하는지 계산할 수 있다

▶ 개별 가중치를 바꿀 때 오차가 얼마나 빠르게 변하는지 알 수 있다

Error Derivative



출력층에 대한 오차 도함수 $E = \frac{1}{2} \sum_{j \in \text{output}} (t_j - y_j)^2$

$$\frac{\partial E}{\partial y_j} = -(t_j - y_j)$$

계산법 j번째 층에 대한 오차 도함수를 가림,

j-1(k)번째 층에 대한 오차 도함수를 계산해내자.

$$\textcircled{1} \quad \frac{\partial E}{\partial x_j} = \sum_k \frac{\partial E}{\partial z_j} \frac{\partial z_j}{\partial x_j} = \sum_k w_{kj} \frac{\partial E}{\partial z_j}$$

$$\frac{\partial z_j}{\partial x_j} = w_{1j}$$

$$\textcircled{2} \quad \frac{\partial E}{\partial z_j} = \frac{\partial E}{\partial y_j} \frac{dy_j}{dz_j} = y_j(1-y_j) \frac{\partial E}{\partial y_j}$$

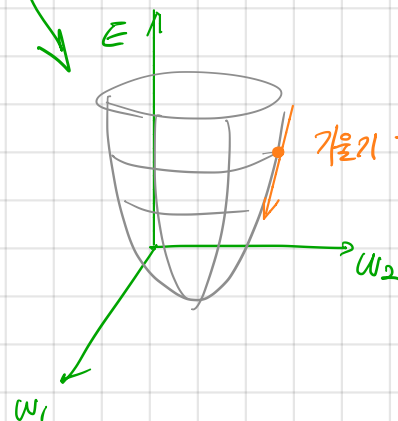
$$\frac{dy_j}{dz_j} = y_j(1-y_j)$$

활성화함수의 시그모이드로 가정

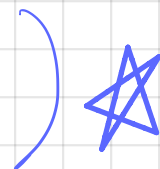
$$\textcircled{1} + \textcircled{2} \Rightarrow \frac{\partial E}{\partial x_j} = \sum_k w_{kj} \frac{\partial E}{\partial z_j} = \sum_k w_{kj} y_j(1-y_j) \frac{\partial E}{\partial y_j}$$

$$\frac{\partial E}{\partial w_{1j}} = \frac{\partial z_j}{\partial w_{1j}} \frac{\partial E}{\partial z_j} = x_1 y_j(1-y_j) \frac{\partial E}{\partial y_j}$$

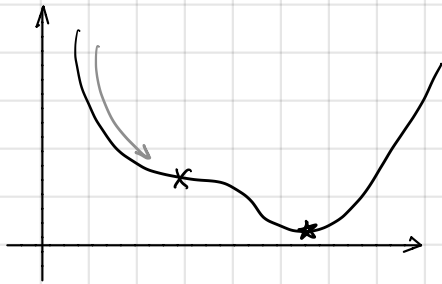
$$\Delta w_{1j} = - \sum_{k \in \text{output}} \epsilon_j^{(k)} y_j^{(k)} (1-y_j^{(k)}) \frac{\partial E}{\partial y_j^{(k)}}$$



가중치 $\frac{\partial E}{\partial w_{1j}}$: 가장 가파른 하강 경로를 따라 다음 가중치를 설정



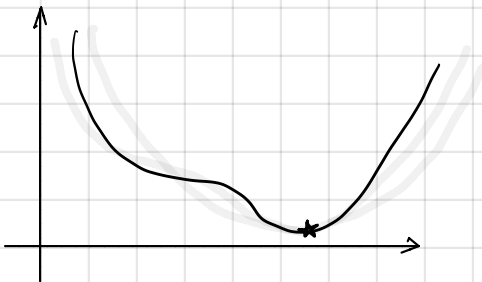
2.6 확률적 경사 하강법(stochastic gradient descent)과 미니배치 경사 하강법 (batch gradient descent)



단일가중치 설정 후 경사하강법을 사용했을 때
발생 가능한 문제점 : local-minimum을 찾아서 고립될 수 있음



해결책 1. 확률적 경사 하강법



단일 정적 오차 곡면X, 동적 오차 곡면O
단점 : 시간 오래걸림

2.7 테스트 데이터와 검증 데이터 그리고 과적합

딥러닝의 딜레마적 문제 : 많은 뉴런 X 수십개의 층 = 과적합 문제 vs 과도하게 일반화

1. 딥러닝은 복잡한 문제로 아주 복잡한 문제를 풀고 과적합 방지를 위해 추가 대책을 적용
2. 학습 데이터와 테스트 데이터를 나눈다.
3. 과적합이 시작되자마자 학습을 멈추는 것이 좋다. ➡ 학습 과정을 epoch으로 나눈다
 - 학습 데이터 크기가 d , 미니배치 경사하강법의 배치크기를 b 라고할 때 d/b 번의 epoch이 수행
 - 각 epoch 끝에 모델이 얼마나 잘 학습되었는지 검증 데이터 세트를 통해 평가한다
4. 하이퍼파라미터 튜닝 : gridsearch

2.8 신경망에서 과적합 막기

1. 정형화 Regulation (규제)

- 중요한 피처의 가중치 값을 늘리는 항을 추가

$$\text{objective} = \text{error} + \lambda f(\theta)$$

하이퍼 파라미터 λ : 규제 X
시행 횟수 필요

1) L2 정형화

$$E_{L2} = E + \frac{1}{2} \lambda w^2$$

큰 값은 2배 크게, 작은 값은 2배 작게...

- 경사하강법이 진행됨에 따라 결국 모든 가중치가 선형적으로 0으로 감소하기 때문에 일반적으로 가중치 감쇠(weight decay)라고 불림

2) L1 정형화

$$E_{L1} = E + \lambda |w|$$

- 최적화 하는 동안 가중치 벡터를 드문드문하게 만듦
- 입력 노이즈에 매우 잘 견딤
- 어떤 피처가 예측 결과에 기여하고 있는지 이해하고자 할 때 유용