

Lab-09-1 ReLU

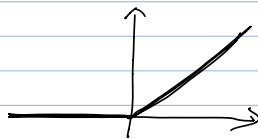
sigmoid는 gradient를 구할 때

이 부분은 0과 아주 가까운 값이나 1과 아주 가까운 값이 나옴

back propagation 구할 때 layer가 많으면 앞쪽 값을 구할 때 gradient가 거의 0이 됨 \Rightarrow vanishing gradient

ReLU: $f(x) = \max(0, x)$

$x = \text{torch.nn.relu}(x)$



drop_last: 마지막 batch를 버리지 말 것

`torch.nn.init.normal_ (linear#. weight)` 으로 초기화

마지막 linear는 relu하지 않고 `CrossEntropyLoss()` 이용
optim. Adam 사용

Lab-09-2 Weight initialization

0으로 초기화하는 것은 안좋은 방법임. back propagation 계산에 가장 좋음.

* Restricted Boltzmann Machine

- layer 사이에 connection 있음.

- 입력 x 가 들어가면 y 를 만드는 forward가 있고

- y 가 들어가면 x 를 복원하는 x' 이 생성됨. \Rightarrow encoding

* Pre-training

RBM: 특이한 layer를 학습하고 그 다음 layer를 하나씩
추가해서 학습한다.

Fine-tuning: RBM 마지막 layer까지 back propagation으로
최종 결과 계산.

* Xavier / He initialization

- Xavier Normal initialization

$$W \sim N(0, \text{Var}(W))$$

$$\text{Var}(W) = \sqrt{\frac{2}{n_{in} + n_{out}}}$$

- Xavier Uniform initialization

$$W \sim U\left(-\sqrt{\frac{6}{n_{in} + n_{out}}}, \sqrt{\frac{6}{n_{in} + n_{out}}}\right)$$

- He normal initialization

$$W \sim N(0, \text{Var}(W))$$

$$\text{Var}(W) = \sqrt{\frac{2}{n_{in}}}$$

- He Uniform initialization

$$W \sim U\left(-\sqrt{\frac{6}{n_{in}}}, \sqrt{\frac{6}{n_{in}}}\right)$$

$$\text{std} = \text{gain} * \text{math.sqrt}(2.0 / (\text{fan}_{in} + \text{fan}_{out}))$$

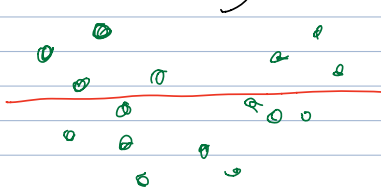
$$a = \text{math.sqrt}(3.0) * \text{std}$$

with torch.no_grad():

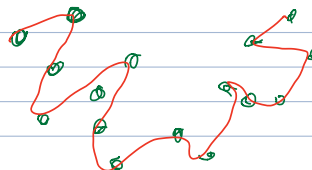
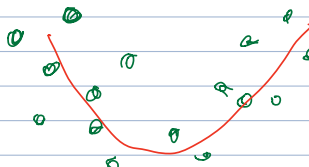
return torch.Uniform(-a, a)

Lab-09-3 Dropout

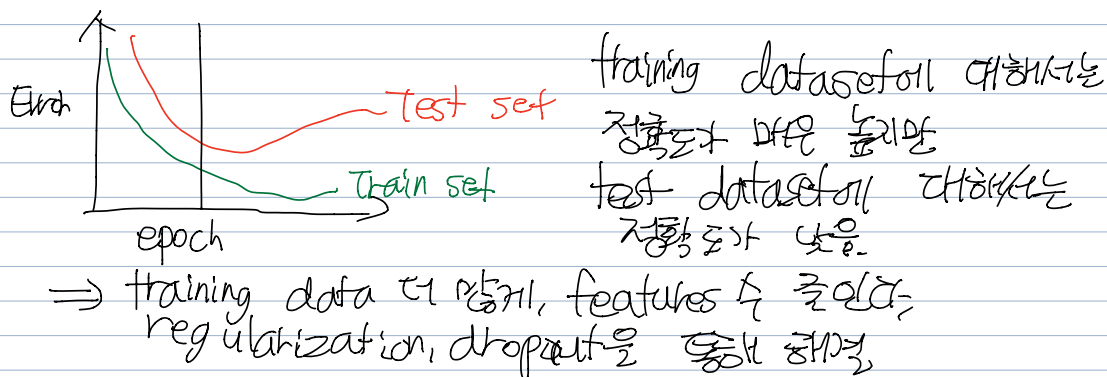
* Overfitting



Underfitting



Overfitting



* Dropout

dropout을 이용하여 무작위로 layer를 선택하여 전파함.

Overfitting 방지, 성능 향상

`dropout = torch.nn.Dropout(p=drop_prob)`

* Train & Eval mode

`model.train()` 버릴 노드와 안버릴 노드를 2대2대 가르킨다.

`model.eval()` 모든 노드를 이용하여 test 할 것이다.

꼭 정의하라!

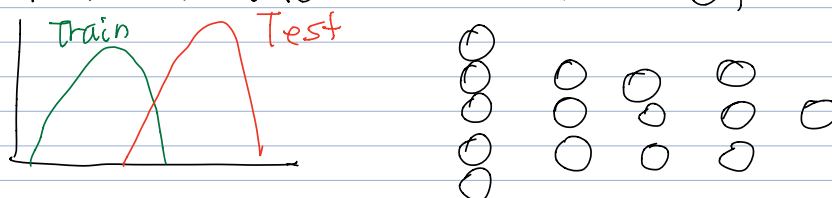
Lab-09-4 Batch Normalization

Gradient Exploding: back propagation 계산에서 너무 큰 값이
나오거나 NaN이 뜨는 현상 \leftrightarrow Gradient Vanishing

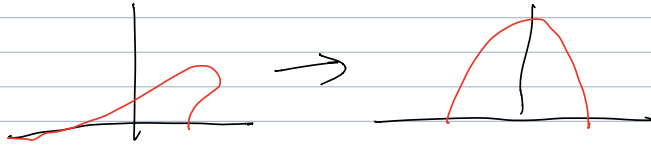
⇒ activation function 바꾸기,

initialization 조심하기, 학습률 (learning rate)

* Internal Covariate Shift: Gradient Exploding/Vanishing 판독.



* Batch normalization



$$\mu_B = \frac{1}{m} \sum_{i=1}^m x_i \quad \sigma_B^2 = \frac{1}{m} \sum_{i=1}^m (x_i - \mu_B)^2 \quad \hat{x}_i = \frac{x_i - \mu_B}{\sqrt{\sigma_B^2 + \epsilon}}$$

$$y_i = \gamma \hat{x}_i + \beta \equiv \text{BN}_{\gamma, \beta}(x_i)$$

$$\text{bn_prediction} = \text{bn_model}(X)$$

$$\text{bn_loss} = \text{Chitotion}(\text{bn_prediction}, Y)$$

batch normalization을 쓰면 더 적은 loss,
더 높은 정확도를 얻음.