

10/20

# pruning

<1>

When should we prune the model, and by how much?

Early works - 'Loss-based pruning'

① generalization

② speed-up ③ interpretability  
= Less compute → Faster speed.

less decision tree

$\min_{\tilde{\theta} \in T_K(\theta)} L(\tilde{\theta})$   $L(\tilde{\theta})$ : training loss  
the set

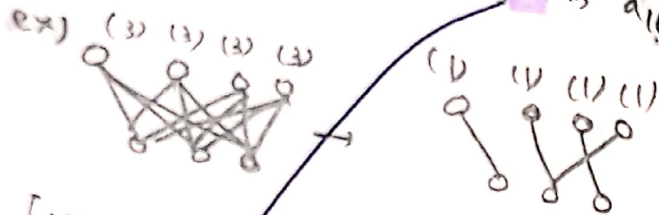
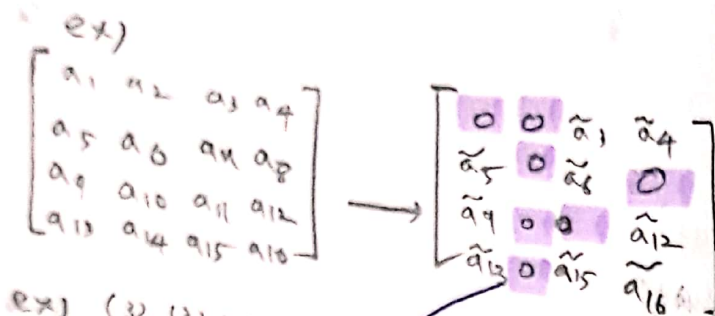
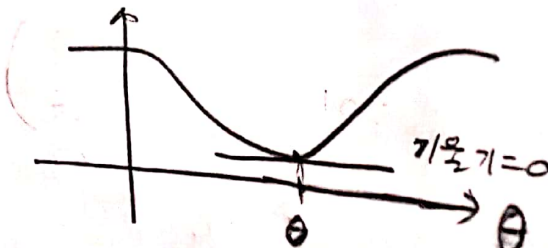
pruning weight = training risk is

minimized right after removing them

$$a \odot b = \begin{bmatrix} a_1 \cdot b_1 \\ a_2 \cdot b_2 \\ \vdots \\ a_n \cdot b_n \end{bmatrix}, a = \begin{bmatrix} a_1 \\ a_2 \\ \vdots \\ a_n \end{bmatrix}, b = \begin{bmatrix} b_1 \\ b_2 \\ \vdots \\ b_n \end{bmatrix}$$

$$10^{12} \binom{10^9}{10^9} = \frac{10^{12}!}{10^9! (10^{12}-10^9)!} \approx (10^{12})^{10^9}$$

$$L(\tilde{\theta}) \approx L(\theta) + (\tilde{\theta} - \theta)^T \cdot \nabla L(\theta) + \frac{1}{2} (\tilde{\theta} - \theta)^T H_{\theta} (\tilde{\theta} - \theta)$$



pruning = Model compression technique  
Zero-ing out the weights

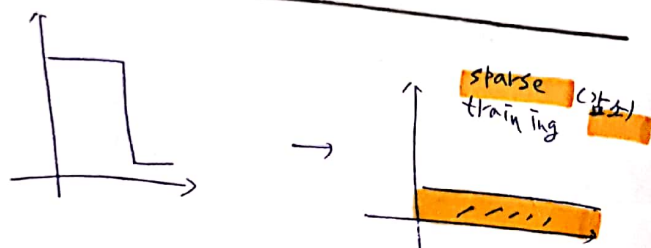
$$\begin{bmatrix} a_1 & 0 \\ 0 & a_4 \end{bmatrix} \begin{bmatrix} b_1 & b_2 \\ b_3 & b_4 \end{bmatrix} = \begin{bmatrix} a_1 \cdot b_1 & a_1 \cdot b_2 \\ a_4 \cdot b_3 & a_4 \cdot b_4 \end{bmatrix}$$

4 multiplication  
0 addition

만약 pruning 안 했으면,  
8 multiplication 1 addition  
4 addition  
ex)  $a_1 \odot b_1 + a_2 \odot b_3$   
2 multiplication

a) How should we select the weight to prune? 'pruning 기준'

- ① Magnitude-based pruning
- ② Loss-based pruning
- ③ Training-based pruning



Remove all weights, except for the top-k elements in terms of  $|\text{weight}|$ .

Optimal Brain Surgeon

Massib & Stark (1998)

known to perform better but requires optimized value much more compute.

approximate Hessian with a Fisher matrix of the model

Hessian

$\nabla f \rightarrow$  vector  $(f_x, f_y)$

$f(x, y)$ 에 대한 second derivative (2차 미분) 한 것에 대한 성분은 총 4개.

$$f''(x, y) = f_{xx}, f_{xy}, f_{yx}, f_{yy}$$

이것을 행렬로 나타내는 것 = Hessian

$$\begin{bmatrix} f_{xx} & f_{xy} \\ f_{yx} & f_{yy} \end{bmatrix}$$

$$D(a, b) = \text{discriminant} = f_{xx} \cdot f_{yy} - f_{xy} \cdot f_{yx}$$

= a, b에 대한 discriminant 값.

critical point 찾기  $\Rightarrow$   $\frac{\partial f}{\partial x} = 0$  일 때,  $f(x)$ 의 비변을 해서  $f'(x) = 0$  이 되는  $x$  값을 찾는다.

$f'(x) = 0$  은 만족시키는 모든  $x$  값이 극대 와 극소의  $x$  값이 아니다.

나, 항상 극대, 극소인지를 확인하기 위해  
이제 함수  $f''(x)$  는 구하여  $f''(x) = 0$  은  
만족시키  $x$  값의 전후의 부호를 확인하여  
극소 값, 극대 값을 찾을 수 있다.  
이제 함수가 어떤 함수에서 Hessian  
matrix = discriminant.

## Training-based pruning

rather, they modify the loss function, gradient update to encourage having many zeros.

Monte Carlo approximation to update  $\theta$ ,  $\phi$  simultaneously.

$$\text{Gaussian } \tilde{\theta} = M\theta, M_i \sim N(1, \sigma)$$

Magnitude-based pruning.

At the

same time, easier to use

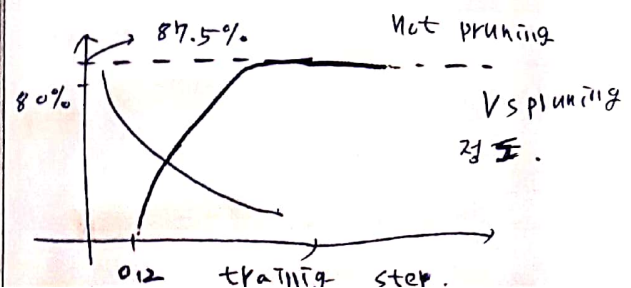
① No extra compute/memory

② No data

pruning scheduling

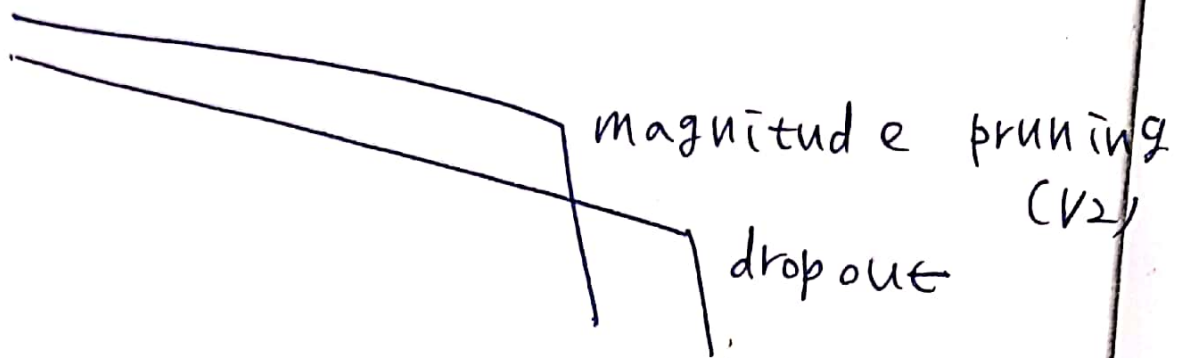
= training / retraining steps, fraction to

remove at each step best to have gradual scheduling.



cubic scheduling.

pruning scheduling 활동작 여부  
= ? ?



★ Magnitude-based = <sup>①</sup> easiest and cheapest to use. ② requires gradual pruning and retraining. ③ requires well-tuned hyperparameters.

Loss-based = more compute / memory - intensive.

training-based = very compute / memory-intensive.

known to be unstable for large-scale

tasks.