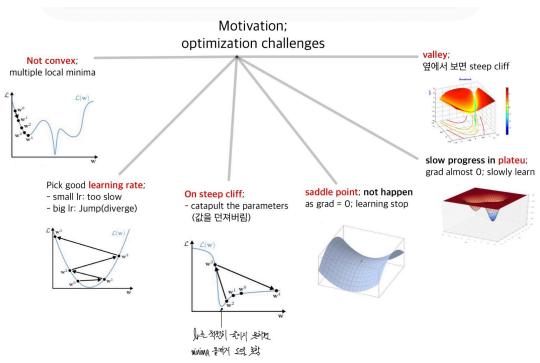


Lec5. optimization

Debugging Tips (Review when you make model.)



Motivation

Optimization algorithms

Debugging Tips

Emergency First Response

- Start with a simple model that is known to work for this type of data (e.g. linear regression, use a standard loss if possible)
- Turn off all bells and whistles, e.g. regularizers, data augmentation
- If fine-tuning a model, double-check the preprocessing: for it to be the same as the one the model's training.
- Verify your input data
- Start with a really small dataset (10 samples). Overfit on it and gradually add more data.
- Start gradually adding back all the pieces that were omitted: augmentation/regulation, custom loss functions, try more complex models.

Visualizing training process

- Understand what happens and update your intuition. For example, the magnitude of the gradients tells you how much the loss function is changing.
- Decompose the visualization theory by themselves and explain it to yourself.
- When in doubt, ask for help!
- After some time, it becomes clear which approach (gradient descent, other optimizers, etc.) is better. These insights will gradually lead to a better understanding of what's going on.

Monitoring stats

- Monitor training metrics to ensure you're getting the best fit. If the validation loss is not decreasing, it may be time to re-examine your model or experiment with different hyperparameters.
- Remember that even though you're not seeing much improvement in training loss, the validation loss may still be decreasing.

Dataset Issues

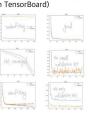
- Check your input data
- Try random input
- Check the data loader
- Make sure input is connected to output
- Verify size of the dataset
- Shuffle the data
- Reduce class imbalance
- Verify number of training examples
- Make sure your batches don't contain a single label
- Use a standard dataset

Implementation issues

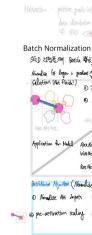
- Try solving a simpler version of the problem
- Make sure your training procedure is correct
- Check your loss function
- Verify loss input
- Adjust loss weights
- Monitor other metrics
- Test any custom layers
- Check for "frozen" layers or variables
- Increase network size
- Check for hidden dimension errors
- Exploding/Vanishing gradients
- Increase/Decrease learning rate
- Try solving a simpler version of the problem
- Make sure your training procedure is correct
- Check your loss function

Optimization strategies

(1) Monitoring the training process (on TensorBoard)



(2) Data Preprocessing

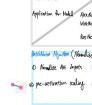


Tip: ① Basic clean: remove outliers > 3 std.
② Batch Normalization: mean <= 0, std >= 1
③ Data Augmentation: 3D rotation, crop, flip, etc.
④ Application in Model: Gradient Descent, Adam, RMSprop, etc.

(3) Data Augmentation

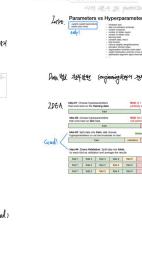


(4) Weight Initialization



Optimization strategies

(5) Learning rate scheduling



(6) Setting hyperparameters



Derivatives

on scalar argument (1D)

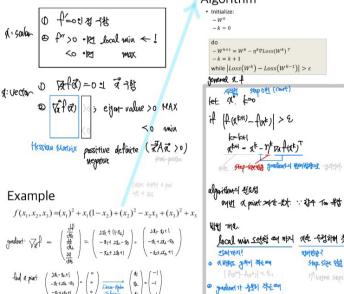
$$\frac{dy}{dx} = \frac{\partial y}{\partial x}$$

Matrix

Hessian

$$H(x) = \nabla^2 f(x) = \left(\begin{array}{ccc} \frac{\partial^2 f}{\partial x_1^2} & \dots & \frac{\partial^2 f}{\partial x_1 \partial x_n} \\ \vdots & \ddots & \vdots \\ \frac{\partial^2 f}{\partial x_n \partial x_1} & \dots & \frac{\partial^2 f}{\partial x_n^2} \end{array} \right)$$

Minimization



on vector argument (2D)

figure

Jacobian

$$J(f) = \begin{bmatrix} \frac{\partial f}{\partial x_1} & \dots & \frac{\partial f}{\partial x_n} \\ \vdots & \ddots & \vdots \\ \frac{\partial f}{\partial x_n} & \dots & \frac{\partial f}{\partial x_1} \end{bmatrix}$$

gradient = ∇f T means of all derivatives

gradient = $\begin{bmatrix} \frac{\partial f}{\partial x_1} \\ \vdots \\ \frac{\partial f}{\partial x_n} \end{bmatrix}$

gradient = ∇f T means of all derivatives

gradient = $\begin{bmatrix} \frac{\partial f}{\partial x_1} & \dots & \frac{\partial f}{\partial x_n} \end{bmatrix}$

gradient = $\begin{bmatrix} \frac{\partial f}{\partial x_1} \\ \vdots \\ \frac{\partial f}{\partial x_n} \end{bmatrix}$

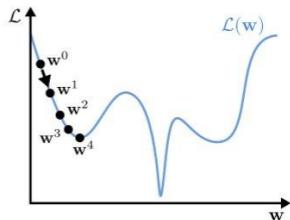
gradient = $\begin{bmatrix} \frac{\partial f}{\partial x_1} & \dots & \frac{\partial f}{\partial x_n} \end{bmatrix}$

Example

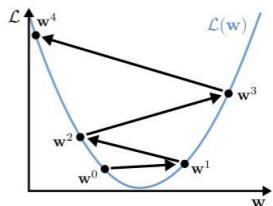
$$f(x_1, x_2, x_3) = x_1^2 + x_1(1-x_2) + (x_2-1)^2 - x_2x_3 + x_3^2 + x_1$$
$$\text{gradient: } \nabla f = \begin{bmatrix} \frac{\partial f}{\partial x_1} \\ \vdots \\ \frac{\partial f}{\partial x_n} \end{bmatrix} = \begin{bmatrix} 2x_1 + 1 - x_2 \\ -x_1 + 2x_2 - x_3 \\ -x_2 \end{bmatrix} = \begin{bmatrix} 2x_1 + 1 - x_2 \\ -x_1 + 2x_2 - x_3 \\ -x_2 \end{bmatrix}$$
$$\text{Jacobian: } J(f) = \begin{bmatrix} \frac{\partial^2 f}{\partial x_1^2} & \dots & \frac{\partial^2 f}{\partial x_1 \partial x_n} \\ \vdots & \ddots & \vdots \\ \frac{\partial^2 f}{\partial x_n \partial x_1} & \dots & \frac{\partial^2 f}{\partial x_n^2} \end{bmatrix} = \begin{bmatrix} 2 & -1 & 0 \\ -1 & 2 & -1 \\ 0 & -1 & 2 \end{bmatrix}$$
$$\text{grad grad: } \nabla(\nabla f) = \begin{bmatrix} \frac{\partial^2 f}{\partial x_1^2} & \dots & \frac{\partial^2 f}{\partial x_1 \partial x_n} \\ \vdots & \ddots & \vdots \\ \frac{\partial^2 f}{\partial x_n \partial x_1} & \dots & \frac{\partial^2 f}{\partial x_n^2} \end{bmatrix} = \begin{bmatrix} 2 & -1 & 0 \\ -1 & 2 & -1 \\ 0 & -1 & 2 \end{bmatrix}$$

Motivation; optimization challenges

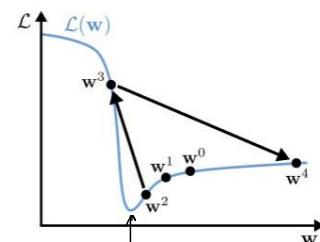
Not convex;
multiple local minima



Pick good **learning rate**;
 - small lr: too slow
 - big lr: Jump(diverge)

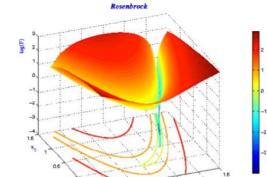


On steep cliff;
 - catapult the parameters
(값을 던져버림)

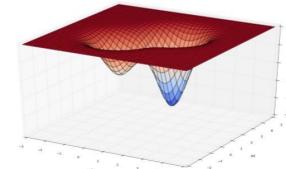


값을
던져버리자
Local
minima
를
넘어가자

valley;
옆에서 보면 steep cliff



slow progress in plateau;
grad almost 0; slowly learn



Optimization algorithms

Gradient Descent

↓ pick mini-batch ($B=2^n$ data) 매번 다 계산할까?

SGD



Algorithm

- ① initialize $w^0 \in \mathbb{R}^B$
- ② Draw random mini batch $\{(x_1, y_1), \dots, (x_B, y_B)\}$
- ③ for i^{th} element (x_i, y_i)
 - forward: $y_b \rightarrow \hat{y}_b$
 - backward: $\nabla_w L_b(w^t) = \nabla_w L_b(\hat{y}_b, y_b, w)$
- ④ Update gradient $w^{t+1} = w^t - \eta \frac{1}{B} \sum \nabla_w L_b(w^t)$ Loss

problem of SGD : how to find good η ?

- ① Jittering in shallow dimension
- ② Stuck in local minima & saddle point
- ③ noisy gradient
: gradient이 있고, minibatch이 gradient

TMI: prof: SGD works with mathematics later!

TMI:

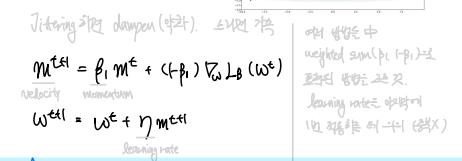
대각선 흥미로운
Momentum의 관점
(堆疊!)

$$\begin{aligned} \text{Vel} &= p_{\text{wt}} - \alpha \nabla f_{\text{wt}} \quad \text{Vel} = p_{\text{wt}} + \nabla f_{\text{wt}} \\ \alpha_{\text{wt}} &= \alpha_{\text{wt}} + \text{Vel}_{\text{wt}} \quad \alpha_{\text{wt}} = \alpha_{\text{wt}} - \eta \text{Vel}_{\text{wt}} \\ \text{Vel}_{\text{wt}} &= \alpha_{\text{wt}} + (\alpha_{\text{wt}} - \eta \nabla f_{\text{wt}}) \quad = \alpha_{\text{wt}} - (\alpha_{\text{wt}} - \eta \nabla f_{\text{wt}}) \nabla f_{\text{wt}} \end{aligned}$$

Solution of SGD

Momentum

SGD + Momentum



Scaling

AdaGrad Adaptive Gradient Estimation

```
grad_squared = 0
while True:
    dx = compute_gradient(x)
    grad_squared += dx * dx
    x -= learning_rate * dx / (np.sqrt(grad_squared) + 1e-7)
```

grad_squared += dx * dx
→ 00

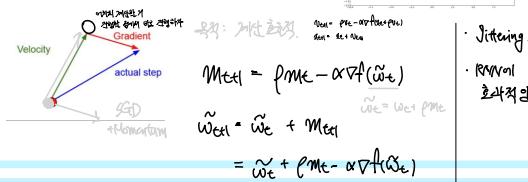
scaling: dimension 별로 각 gradient 평균화 (parameter)

lr: adaptive lr는 정지

Adam(good) Momentum + RMSprop

```
first_moment = 0
second_moment = 0
for t in range(num_iterations):
    dx = compute_gradients(x)
    first_moment = beta1 * first_moment + (1 - beta1) * dx
    second_moment = beta2 * second_moment + (1 - beta2) * dx * dx
    first_unbias = first_moment / (1 - beta1 ** t)
    second_unbias = second_moment / (1 - beta2 ** t)
    x -= learning_rate * first_unbias / (np.sqrt(second_unbias) + 1e-7))
```

SGD + Nesterov Momentum



RMSProp

```
grad_squared = 0
while True:
    dx = compute_gradients(x)
    grad_squared = decay_rate * grad_squared + (1 - decay_rate) * dx * dx
    x -= learning_rate * dx / (np.sqrt(grad_squared) + 1e-7)
```

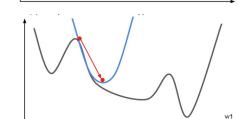
decay-rate ≠ 1, weighted sum 더해
이후 초기 설정으로, learning OK (잘 adaptive 됨)

Understanding with Math later!

Momentum

Bias correction

AdaGrad / RMSProp



TMI: Second-order optimization

RNN은 X에 따라 수렴하는 경향

수렴하기 어렵고, 초기 반복에서 초기화가 필요

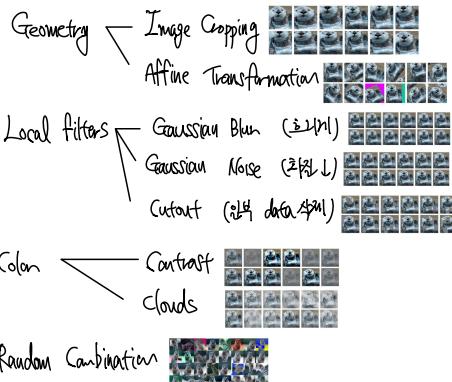
Optimization strategies

[1] Monitoring the training process (on TensorBoard)



[2] Data Augmentation

By data augmentation → data augmentation

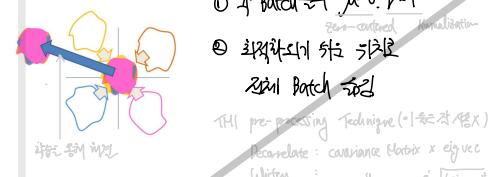


[3] Data Preprocessing

Motivation : positive grads (cf. sigmoid) \rightarrow 편향된 최적화
data distribution of each layer is different
cf. SGD

Batch Normalization

SGD 즐겁지 않음 Batch size, but Batch size distribution \neq
Normalize 1st layer ; gradient \rightarrow $\frac{\partial L}{\partial \mu} = \frac{1}{N} \sum_{i=1}^N \frac{\partial L}{\partial \mu_i}$ $\frac{\partial L}{\partial \sigma^2} = \frac{1}{N} \sum_{i=1}^N \frac{\partial L}{\partial \sigma_i^2}$
Solution (Not finish!) \rightarrow Batch size, $\mu=1$



Application for Model

AlexNet: Each 2image -> 2image
VGGNet: " -> Each channel μ
ResNet: " -> Each channel μ
Each channel σ

BatchNorm Algorithm (Normalize hidden layers)

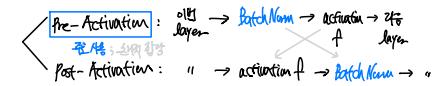
$$\text{① Normalize Net Inputs } \hat{x}_j^{(i)} = \frac{x_j^{(i)} - \mu_i}{\sqrt{\sigma_i^2 + \epsilon}}$$

$$\text{② pre-activation scaling } a_j^{(i)} = \gamma_j \cdot \hat{x}_j^{(i)} + \beta_j$$

Note

- ① Batch Norm's training algorithm \neq test
 - μ, σ \rightarrow sample size
 - ? (2번)
 - sample size \rightarrow μ, σ \rightarrow 1번
 - batch size = (μ, σ)
 - batch size \neq 1번
- ② Batch Norm's covariance shift \rightarrow $\mu \rightarrow \mu X$
작간 수치는 그대로 ; 다른 계산에 영향을 주거나 다른 계산에 영향을 주거나

③ BatchNorm 적용 방식



④ BatchSize 개진으로 성능 안정화.

⑤ BatchSize 1으로 훈련할 경우의 x 연산 처리 대비

⑥ BatchNorm의 input layer는 optimizer가 아님 이거 카운트 백포그 \rightarrow computation graph에 포함되지 않음

⑦ Batch Normalization for ConvNet



⑧ 각 채널 Normalization \rightarrow Layer Normalization



Optimization strategies

[4] Weight Initialization

Bad initialization

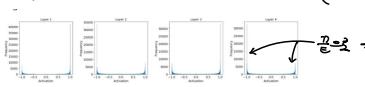
① Constant Initialization

weight = $\frac{1}{\text{size}}$ 일정 $\rightarrow \text{grad} = 0$

② Small random number $\rightarrow \text{grad vanishing}$ ($b=0.01$)



③ Large " " \rightarrow ($b=0.2$)

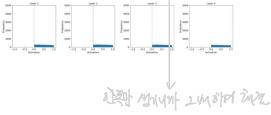


Xavier Initialization ($b^2 = 1/D_{in}$)

$$\text{Var}(y) \approx \text{Var}(w^T x) = D_{in} \text{Var}(x_i w_i) \\ = D_{in} (\mathbb{E}[x_i^2 w_i^2] - \mathbb{E}[x_i w_i]^2) \\ = D_{in} (\mathbb{E}[x_i^2] \mathbb{E}[w_i^2] - \mathbb{E}[x_i]^2 \mathbb{E}[w_i]^2) \\ = D_{in} \mathbb{E}[x_i^2] \mathbb{E}[w_i^2] \\ = D_{in} \text{Var}(x_i) \text{Var}(w_i)$$

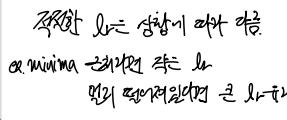
ReLU only
positive weight

He Initialization ($b^2 = 2/D_{in}$)

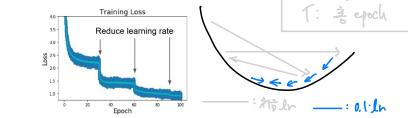


linear activation function

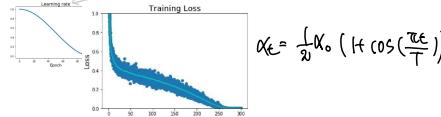
[5] Learning rate scheduling



① 90 epochs of lr = 0.1 막상 (cf. ResNet)



② Cosine- η lr減衰법 사용



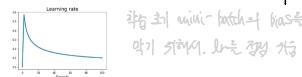
③ Linear- η lr減衰법 사용 (cf. Best)

$$\eta_t = \eta_0 (1 - \frac{t}{T})$$

④ Inverse sqrt η_t (cf. transformer; attention is all you need)

$$\eta_t = \frac{\eta_0}{\sqrt{t}}$$

⑤ lr 증가하기 (a.k.a. linear warm up)



[6] Setting hyperparameters

하이퍼파라미터 설정

Intro.

Parameters vs Hyperparameters

- minibatch size
- data normalization schemes
- number of epochs
- number of layers
- number of hidden units
- learning rates
- random seed, why?
- loss functions
- various weights (embedding dims)
- various model types
- regularization schemes (more later)
- weight initialization schemes (more later)
- optimization algorithm (more later)
- ...

Data 품질 조건 설정 (engineering strategy 설정하기)

ZDEA

Idea #1: Choose hyperparameters that work best on the training data

train

BAD: K = 1 always works perfectly on training data

Idea #2: choose hyperparameters that work best on test data

train test

BAD: No idea how algorithm will perform on new data

Idea #3: Split data into train, val, choose hyperparameters on val and evaluate on test

train validation test

Better!

Idea #4: Cross-Validation: Split data into folds, try each fold as validation and average the results

fold 1 fold 2 fold 3 fold 4 fold 5 test

fold 1 fold 2 fold 3 fold 4 fold 5 test

fold 1 fold 2 fold 3 fold 4 fold 5 test

fold 1 fold 2 fold 3 fold 4 fold 5 test

averaging

평균 사용

Debugging Tips

Emergency First Response

- Start with a simple model that is known to work for this type of data cf. VGG for images, use a standard loss if possible
- Turn off all bells and whistles, e.g. regularization and data augmentation
- If finetuning a model, double check the preprocessing, for it should be the same as the original model's training
- Verify that the input data is correct
- Start with a really small dataset(1-10 samples). Overfit on it and gradually add more data
- Start gradually adding back all the pieces that were omitted: augmentation/regularization, custom loss functions, try more complex models.

Visualize the training process

- ▶ Monitor the activations, weights, and updates of each layer. Make sure their magnitudes match. For example, the magnitude of the updates to the parameters (weights and biases) should be 1×10^{-3} .
- ▶ Consider a visualization library like Tensorboard and Crayon. In a pinch, you can also print weights/biases/activations.
- ▶ Be on the lookout for layer activations with a mean much larger than 0. Try Batch Norm or ELUs.
- ▶ Weight histograms should have an approximately Gaussian (normal) distribution, after some time. For biases, these histograms will generally start at 0, and will usually end up being approximately Gaussian. Keep an eye out for parameters that are diverging or biases that become very large.

Overcoming NaNs

- ▶ Decrease the learning rate, especially if you are getting NaNs in the first 100 iterations.
- ▶ NaNs can arise from division by zero or natural log of zero or negative number.
- ▶ Try evaluating your network layer by layer and see where the NaNs appear.

Dataset Issues

- Check your input data
- Try random input
- Check the data loader
- Make sure input is connected to output
- Verify noise in the dataset
- Shuffle the dataset
- Reduce class imbalance
- Verify number of training examples
- Make sure your batches don't contain a single label
- Use a standard dataset

Training Issues

- Solve for really small dataset
- Check weights initialization
- Change your hyperparameters
- Reduce regularization
- Give it time
- Switch from Train to Test mode
- **Visualize the training process**
- Try a different optimizer
- Exploding/Vanishing gradients
- Increase/Decrease learning rate
- **Overcoming NaNs**

Data Normalization/ Augmentation Issues

- standardize the features
- Check for too much data augmentation
- Check the preprocessing of your pretrained model
- Check the preprocessing for train/validation/test set

Implementation Issues

- Try solving a simpler version of the problem
- Make sure your training procedure is correct
- Check your loss function
- Verify loss input
- Adjust loss weights
- Monitor other metrics
- Test any custom layers
- Check for "frozen" layers or variables
- Increase network size
- Check for hidden dimension errors
- Explore gradient checking
- Try solving a simpler version of the problem
- Make sure your training procedure is correct
- Check your loss function