

How Learning differs from pure optimization

$$J(\theta) = \mathbb{E}_{(x,y) \sim P_{\text{data}}} L(f(x; \theta), y)$$

- Empirical Risk minimization

$$\hat{J}(x, y) \sim P_{\text{data}} L(f(x; \theta), y) = \frac{1}{m} \sum_{i=1}^m L(f(x_i; \theta), y_i)$$

- Surrogate Loss functions and early stopping

Pure optimization: Minimize/zero out the gradient.

Learning: Aim to minimize the generalization/validation error.

- Batch and minibatch algorithms.

- Compute overall training data sometimes hard to accomplish, and some training data may be related.

- Stochastic methods: Randomly sample batches.

Online methods: Data come one-by-one sequentially from on-line.

- When we get endless data source, how we arrange and utilize is more important than solving over-fitting.

Challenges in Neural Network Optimization

- Ill Conditioning

- Hessian matrix may be ill-conditioned/poorly-conditioned so that the error surface has more traps and too complicated.

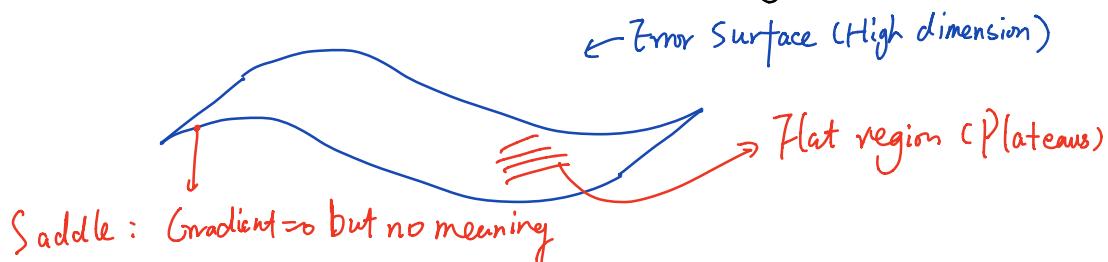
$$\frac{1}{2} \epsilon^2 g^T H g - \epsilon g^T g : \text{when } g^T H g \text{ exceeds } g^T g$$

- Ill-conditioning \approx Gradient norm $\uparrow \neq$ Training failed

- Local minima \Rightarrow we can plot the gradient over time.

- Model Identifiability Problem \Rightarrow Sometimes caused by weight space symmetry

- Local minimum + High cost \Rightarrow Problem
- Plateaus, Saddle Points and other flat region



Learning algorithms should escape or avoid points of these kinds.

- Cliffs and Exploding Gradients
- Long terms Dependencies

$$w^t = (\nabla \text{diag}(w) w^t)^t = \nabla \text{diag}(w^t) w^t$$

we get vanishing and exploding gradients problem.

- Inexact Gradients
- Poor Correspondence between local and global Structure
 - During learning process, the direction of local optimization may be useless or opposite to global optimal.
 - Initial weights thus are important.
- Theoretical Limits for Optimization

Back to the conclusion that we really need focus on validation error.

Basic Algorithms

- Stochastic (batch selected) gradient descent.

When we talking about SGD, we are talking on:

$$\hat{g} \leftarrow \frac{1}{m} \nabla_{\theta} \sum_i \|f(x_i; \theta), y_i\| \quad m \text{ is the size of a minibatch}$$

$$\theta \leftarrow \theta - \epsilon_k \hat{g} \quad G_k = (I - \epsilon_k \hat{g}) G_{k-1} + \epsilon_k \hat{g}, \text{ which means that we first go quick then go gently.}$$

Though Stochastic minibatch may sometimes mis-leading, but it has many other benefits, say rapid progress when beginning

- Momentum

Solving two aspects of problem:

① Poor-Conditioning Hessian.

② Variance in stochastic gradient.

$$v \leftarrow \alpha v - \eta \nabla_{\theta} \left(\frac{1}{m} \sum_i L(f(x_i; \theta), y_i) \right)$$

$$\theta \leftarrow \theta + v$$

Update rules with momentum is actually considering a sequence of gradients. Similar with a physical system with two force affecting a particle: Loss $\boxed{\nabla \ell}$ and viscous drag $\boxed{-\eta v_t}$

- Nesterov Momentum

Calculating gradients with current θ adding by current momentum.

Parameter Initializing Strategies

- Initialization often breaks symmetry.
- Bigger value can keep stability but always failed with vanishing or exploding problem.
- We prefer small weights cause it gives more optimization space.
- Many initialization tricks: Sparse initialization, try with a mini-batch. etc.

Algorithms with adaptive learning rates

- AdaGrad

Dividing all the historical gradients summation as the learning rate

- RMSprop \Rightarrow Adam algorithms

- RMSprop is a way of combining momentum with Adagrad. Adam does more:

For each minibatch with size M .

$$g \leftarrow \frac{1}{M} \nabla_{\theta} \sum_i l(f(x_i; \theta), y_i) \quad \text{Compute gradient}$$

$$s \leftarrow p_1 s + (1 - p_1) g \quad \text{Corresponding to momentum}$$

$$r \leftarrow p_2 r + (1 - p_2) g \odot g \quad \text{Corresponding to Adagrad.}$$

$$\hat{s} \leftarrow \frac{s}{1 - p_1^t} \quad \begin{matrix} \nearrow \\ \text{Rescale them in case they sometimes being} \end{matrix}$$

$$\hat{r} \leftarrow \frac{r}{1 - p_2^t} \quad \begin{matrix} \nearrow \\ \text{too small.} \end{matrix}$$

$$\theta \leftarrow \theta - \epsilon \frac{\hat{s}}{\sqrt{\hat{r}} + \delta} \quad \begin{matrix} \nearrow \\ \text{momentum} \end{matrix}$$

\nearrow Adagrad.

Approximate Second-Order Methods

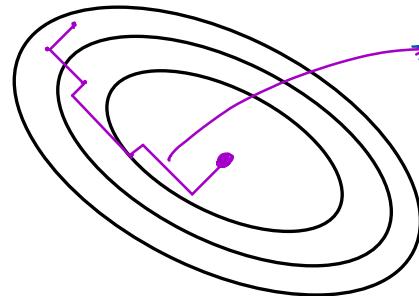
Note that for poorly conditioning Hessian, second orders methods may fail this discussion here introducing Newton method and some other optimization for saving computational loss.

- Newton's method

$$\theta^* = \theta_0 - H^{-1} \nabla_{\theta} J(\theta)$$

- Conjugate Gradients

An optimized version of steepest gradient descent



steepest descent / line search according to
vertical direction of last time

$$d_t \cdot d_{t-1} = 0$$

$$d_t = \nabla_{\theta} J(\theta) + \underbrace{\beta_t d_{t-1}}_{\text{Approximate } H}$$

Conjugate Gradients define a direction of line search, which is more efficient than steepest descent.

- BFGS

Saving computational cost of getting \mathbf{H} using a part informative \mathbf{M} .

Optimization Strategies and Meta-Algorithm

- Batch Normalization

- An important assumption of deep learning: Updating a layer with other layers stay stable.

↳ But we cannot achieve. → Vanishing/Exploding gradients.

- Form of batch normalization

$$\mu = \frac{1}{m} \sum_i H_i \quad (\text{Here } H \text{ is design matrix})$$

$$S = \sqrt{S + \frac{1}{m} \sum_i (H - \mu)^2} \quad (S \text{ is a stable factor})$$

- BN can degenerate linear accumulation of past units but preserve non-linear and other attributes.

- Learnable parameters γ and β :

$$\gamma H' + \beta \rightarrow \hat{H}$$

Here these two parameters can be learned by deep neural network easily and increasing the generalization ability.

- Coordinate Descent : learn partly the parameters each time.

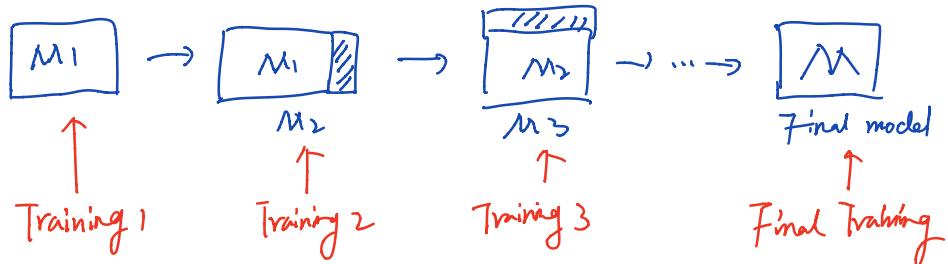
- Polyak Averaging : Average gradients from a period, then update.

- Designing models to aid optimization: Auxiliary link / structure

- Continuation methods and curriculum learning : Learning from J_0 to J .
Simple Hard

- Supervised pretraining:

- Greedy supervised pretraining:



- FitNet:

