

Scaling Optimization



Lecture 4 Recap

Neural Network



Source: <u>http://cs231n.github.io/neural-networks-1/</u>

Neural Network



I2DL: Prof. Dai

Compute Graphs → Neural Networks Output layer Input layer $* W_0$ x_0 _oss/ x^*x $\max(0, x)$ + $-y_0$ x_0 cost \hat{y}_0 * W₁ x_1 y_0 Rel U Activation x_1 |2| OSS Weights Input (not arguing this is the (unknowns!) right choice here) We want to compute gradients w.r.t. all weights \boldsymbol{W} e.g., class label/ regression target



Compute Graphs → Neural Networks



Goal: We want to compute gradients of the loss function *L* w.r.t. all weights *w*

 $L = \sum_{i} L_{i}$

L: sum over loss per sample, e.g. L2 loss \rightarrow simply sum up squares: $L_i = (\hat{y}_i - y_i)^2$

 \rightarrow use chain rule to compute partials

$$\frac{\partial L}{\partial w_{i,k}} = \frac{\partial L}{\partial \hat{y}_i} \cdot \frac{\partial \hat{y}_i}{\partial w_{i,k}}$$

We want to compute gradients w.r.t. all weights W AND all biases b

Summary

- We have
 - (Directional) compute graph
 - Structure graph into layers
 - Compute partial derivatives w.r.t. weights (unknowns)



• Next

- Find weights based on gradients

Gradient step: $W' = W - \alpha \nabla_W f_{\{x,y\}}(W)$



Optimization





Gradient Descent

• From derivative to gradient

df(x)

 $\nabla_{r} f(x)$ greatest increase of the function

Direction of

• Gradient steps in direction of negative gradient



Gradient Descent

• From derivative to gradient

 $\frac{\mathrm{d}f(x)}{\mathrm{d}x} \longrightarrow \nabla_{\!x}f(x)$

Direction of greatest increase of the function

• Gradient steps in direction of negative gradient



Gradient Descent

• From derivative to gradient

 $\frac{\mathrm{d}f(x)}{\mathrm{d}x} \longrightarrow \nabla_{\!x} f(x)$

Direction of greatest increase of the function

• Gradient steps in direction of negative gradient





Convergence of Gradient Descent

• Convex function: all local minima are global minima



Source: https://en.wikipedia.org/wiki/Convex_function#/media/File:ConvexFunction.svg

If line/plane segment between any two points lies above or on the graph

Convergence of Gradient Descent

- Neural networks are non-convex
 - many (different) local minima
 - no (practical) way to say which is globally optimal



Source: Li, Qi. (2006). Challenging Registration of Geologic Image Data

Convergence of Gradient Descent



Source: <u>https://builtin.com/data-science/gradient-</u> descent



Gradient Descent: Multiple Dimensions



Source: builtin.com/data-science/gradient-descent

Various ways to visualize...

Gradient Descent: Multiple Dimensions



Source: http://blog.datumbox.com/wp-content/uploads/2013/10/gradient-descent.png

Gradient Descent for Neural Networks



Gradient Descent: Single Training Sample

- Given a loss function L and a single training sample {x_i, y_i}
- Find best model parameters $\theta = \{W, b\}$
- Cost $L_i(\boldsymbol{\theta}, \boldsymbol{x}_i, \boldsymbol{y}_i)$
 - $\boldsymbol{\theta} = \arg\min L_i(\boldsymbol{x}_i, \boldsymbol{y}_i)$
- Gradient Descent:
 - Initialize θ^1 with 'random' values (more on that later)

$$- \boldsymbol{\theta}^{k+1} = \boldsymbol{\theta}^k - \alpha \nabla_{\boldsymbol{\theta}} L_i(\boldsymbol{\theta}^k, \boldsymbol{x}_i, \boldsymbol{y}_i)$$

- Iterate until convergence: $|\theta^{k+1} - \theta^k| < \epsilon$

Gradient Descent: Single Training Sample



(current model)

- $\nabla_{\theta} L_i(\theta^k, x_i, y_i)$ computed via backpropagation
- Typically: dim $(\nabla_{\theta} L_i(\theta^k, x_i, y_i)) = \dim(\theta) \gg 1$ million

Gradient Descent: Multiple Training Samples

- Given a loss function L and multiple (n) training samples {x_i, y_i}
- Find best model parameters $\theta = \{W, b\}$

• Cost
$$L = \frac{1}{n} \sum_{i=1}^{n} L_i(\boldsymbol{\theta}, \boldsymbol{x}_i, \boldsymbol{y}_i)$$

- $\boldsymbol{\theta} = \arg \min L$

Gradient Descent: Multiple Training Samples

• Update step for multiple samples

$$\boldsymbol{\theta}^{k+1} = \boldsymbol{\theta}^k - \alpha \nabla_{\boldsymbol{\theta}} L(\boldsymbol{\theta}^k, \boldsymbol{x}_{\{1..n\}}, \boldsymbol{y}_{\{1..n\}})$$

• Gradient is average / sum over residuals

$$\nabla_{\boldsymbol{\theta}} L(\boldsymbol{\theta}^{k}, \boldsymbol{x}_{\{1..n\}}, \boldsymbol{y}_{\{1..n\}}) = \frac{1}{n} \sum_{i=1}^{n} \nabla_{\boldsymbol{\theta}} L_{i}(\boldsymbol{\theta}^{k}, \boldsymbol{x}_{i}, \boldsymbol{y}_{i})$$

Reminder: this comes from backprop.

• Often people are lazy and just write: $\nabla L = \sum_{i=1}^{n} \nabla_{\theta} L_i$ – omitting $\frac{1}{n}$ is not 'wrong', it just means rescaling the learning rate

Side Note: Optimal Learning Rate

Can compute optimal learning rate α using Line Search (optimal for a given set)

- 1. Compute gradient: $\nabla_{\theta} L = \frac{1}{n} \sum_{i=1}^{n} \nabla_{\theta} L_i$
- 2. Optimize for optimal step α :

$$arg\min_{\alpha} L(\boldsymbol{\theta}^{k} - \alpha \nabla_{\boldsymbol{\theta}} L)$$
$$\boldsymbol{\theta}^{k+1} = \boldsymbol{\theta}^{k} - \alpha \nabla_{\boldsymbol{\theta}} L$$
No

Not that practical for DL since we need to solve huge system every step...

3

Gradient Descent on Train Set

- Given large train set with n training samples $\{x_i, y_i\}$
 - Let's say 1 million labeled images
 - Let's say our network has 500k parameters
- Gradient has 500k dimensions
- n = 1 million
- \rightarrow Extremely expensive to compute

• If we have *n* training samples we need to compute the gradient for all of them which is *O*(*n*)

• If we consider the problem as empirical risk minimization, we can express the total loss over the training data as the expectation of all the samples

$$\frac{1}{n} \left(\sum_{i=1}^{n} L_i(\boldsymbol{\theta}, \boldsymbol{x}_i, \boldsymbol{y}_i) \right) = \mathbb{E}_{i \sim [1, \dots, n]} [L_i(\boldsymbol{\theta}, \boldsymbol{x}_i, \boldsymbol{y}_i)]$$

• The expectation can be approximated with a small subset of the data

$$\mathbb{E}_{i \sim [1, \dots, n]}[L_i(\boldsymbol{\theta}, \boldsymbol{x_i}, \boldsymbol{y_i})] \approx \frac{1}{|S|} \sum_{j \in S} \left(L_j(\boldsymbol{\theta}, \boldsymbol{x_j}, \boldsymbol{y_j}) \right) \text{ with } S \subseteq \{1, \dots, n\}$$

$\begin{array}{l} {\rm Minibatch} \\ {\rm choose \ subset \ of \ trainset \ } m \ll n \end{array}$

$$B_i = \{\{x_1, y_1\}, \{x_2, y_2\}, \dots, \{x_m, y_m\}\}\$$

$$\{B_1, B_2, \dots, B_{n/m}\}$$

- Minibatch size is hyperparameter
 - Typically power of 2 → 8, 16, 32, 64, 128...
 - Smaller batch size means greater variance in the gradients
 - → noisy updates
 - Mostly limited by GPU memory (in backward pass)
 - E.g.,
 - Train set has $n = 2^{20}$ (about 1 million) images
 - With batch size m = 64: $B_{1 \dots n/m} = B_{1 \dots 16,384}$ minibatches

(Epoch = complete pass through training set)

$$\boldsymbol{\theta}^{k+1} = \boldsymbol{\theta}^k - \alpha \nabla_{\boldsymbol{\theta}} L(\boldsymbol{\theta}^k, \boldsymbol{x}_{\{1..m\}}, \boldsymbol{y}_{\{1..m\}})$$

k now refers to k-th iteration

$$\nabla_{\boldsymbol{\theta}} L = \frac{1}{m} \sum_{i=1}^{m} \nabla_{\boldsymbol{\theta}} L_i$$

 $\sim m$ training samples in the current minibatch

Gradient for the *k*-th minibatch

Note the terminology: iteration vs epoch

Convergence of SGD

Suppose we want to minimize the function $F(\theta)$ with the stochastic approximation

$$\theta^{k+1} = \theta^k - \alpha_k H(\theta^k, X)$$

where $\alpha_1, \alpha_2 \dots \alpha_n$ is a sequence of positive step-sizes and $H(\theta^k, X)$ is the unbiased estimate of $\nabla F(\theta^k)$, i.e.

$$\mathbb{E}\big[H\big(\theta^k,X\big)\big] = \nabla F\big(\theta^k\big)$$

Robbins, H. and Monro, S. "A Stochastic Approximation Method" 1951.

I2DL: Prof. Dai

Convergence of SGD

$$\theta^{k+1} = \theta^k - \alpha_k H(\theta^k, X)$$

converges to a local (global) minimum if the following conditions are met:

1)
$$\alpha_n \ge 0, \forall n \ge 0$$

2) $\sum_{n=1}^{\infty} \alpha_n = \infty$
3) $\sum_{n=1}^{\infty} \alpha_n^2 < \infty$
4) $F(\theta)$ is strictly convex

The proposed sequence by Robbins and Monro is $\alpha_n \propto \frac{\alpha}{n}$, for n > 0

Problems of SGD

- Gradient is scaled equally across all dimensions
 - \rightarrow i.e., cannot independently scale directions
 - → need to have conservative min learning rate to avoid divergence
 - → Slower than 'necessary'
- Finding good learning rate is an art by itself
 → More next lecture

Gradient Descent with Momentum



Source: A. Ng

We're making many steps back and forth along this dimension. Would love to track that this is averaging out over time.

Would love to go faster here... I.e., accumulated gradients over time



Exponentially-weighted average of gradient Important: velocity \boldsymbol{v}^{k} is vector-valued!

[Sutskever et al., ICML'13] On the importance of initialization and momentum in deep learning I2DL: Prof. Dai

Gradient Descent with Momentum



Step will be largest when a sequence of gradients all point to the same direction

Hyperparameters are α, β β is often set to 0.9

$$\boldsymbol{\theta}^{k+1} = \boldsymbol{\theta}^k + \boldsymbol{v}^{k+1}$$

Gradient Descent with Momentum

• Can it overcome local minima?



Nesterov Momentum

• Look-ahead momentum

$$\widetilde{\boldsymbol{\theta}}^{k+1} = \boldsymbol{\theta}^k + \boldsymbol{\beta} \cdot \boldsymbol{\nu}^k$$

$$\boldsymbol{v}^{k+1} = \boldsymbol{\beta} \cdot \boldsymbol{v}^k - \boldsymbol{\alpha} \cdot \nabla_{\boldsymbol{\theta}} L(\boldsymbol{\widetilde{\theta}}^{k+1})$$

$$\boldsymbol{\theta}^{k+1} = \boldsymbol{\theta}^k + \boldsymbol{v}^{k+1}$$

Nesterov, Yurii E. "A method for solving the convex programming problem with convergence rate O (1/k^ 2)." *Dokl. akad. nauk Sssr.* Vol. 269. 1983.

I2DL: Prof. Dai

Nesterov Momentum

- First make a big jump in the direction of the previous accumulated gradient.
- Then measure the gradient where you end up and make a correction.



blue vectors = standard momentum

Source: G. Hinton

$$\widetilde{\boldsymbol{\theta}}^{k+1} = \boldsymbol{\theta}^k + \beta \cdot \boldsymbol{v}^k$$
$$\boldsymbol{v}^{k+1} = \beta \cdot \boldsymbol{v}^k - \alpha \cdot \nabla_{\boldsymbol{\theta}} L(\widetilde{\boldsymbol{\theta}}^{k+1})$$
$$\boldsymbol{\theta}^{k+1} = \boldsymbol{\theta}^k + \boldsymbol{v}^{k+1}$$

I2DL: Prof. Dai

Root Mean Squared Prop (RMSProp)



Source: Andrew. Ng

• RMSProp divides the learning rate by an exponentially-decaying average of squared gradients.

Hinton et al. "Lecture 6.5-rmsprop: Divide the gradient by a running average of its recent magnitude." COURSERA: Neural networks for machine learning 4.2 (2012): 26-31.

RMSProp

$$\begin{split} \boldsymbol{s}^{k+1} &= \boldsymbol{\beta} \cdot \boldsymbol{s}^{k} + (1-\boldsymbol{\beta}) \left[\nabla_{\boldsymbol{\theta}} L \circ \nabla_{\boldsymbol{\theta}} L \right] \\ \boldsymbol{\theta}^{k+1} &= \boldsymbol{\theta}^{k} - \boldsymbol{\alpha} \cdot \frac{\nabla_{\boldsymbol{\theta}} L}{\sqrt{\boldsymbol{s}^{k+1}} + \boldsymbol{\epsilon}} \end{split} \\ \end{split}$$
Element-wise multiplication



RMSProp



RMSProp

• Dampening the oscillations for high-variance directions

- Can use faster learning rate because it is less likely to diverge
 - → Speed up learning speed
 - → Second moment

Adaptive Moment Estimation (Adam)

Idea : Combine Momentum and RMSProp

$$m^{k+1} = \beta_1 \cdot m^k + (1 - \beta_1) \nabla_{\theta} L(\theta^k) \longleftarrow \qquad \text{First momentum:} \\ mean of gradients \\ v^{k+1} = \beta_2 \cdot v^k + (1 - \beta_2) [\nabla_{\theta} L(\theta^k) \circ \nabla_{\theta} L(\theta^k)]$$

$$\boldsymbol{\theta}^{k+1} = \boldsymbol{\theta}^k - \alpha \cdot \frac{\boldsymbol{m}^{k+1}}{\sqrt{\boldsymbol{v}^{k+1}} + \boldsymbol{\epsilon}}$$

Note: This is not the update rule of Adam

Second momentum: variance of gradients

Q. What happens at k = 0? A. We need bias correction as $m^0 = 0$ and $v^0 = 0$

1 . .

[Kingma et al., ICLR'15] Adam: A method for stochastic optimization

I2DI : Prof. Dai

Adam : Bias Corrected

Combines Momentum and RMSProp

 $\boldsymbol{m}^{k+1} = \beta_1 \cdot \boldsymbol{m}^k + (1 - \beta_1) \nabla_{\boldsymbol{\theta}} L(\boldsymbol{\theta}^k) \qquad \boldsymbol{v}^{k+1} = \beta_2 \cdot \boldsymbol{v}^k + (1 - \beta_2) [\nabla_{\boldsymbol{\theta}} L(\boldsymbol{\theta}^k) \circ \nabla_{\boldsymbol{\theta}} L(\boldsymbol{\theta}^k)]$

- m^k and v^k are initialized with zero
 - → bias towards zero
 - → Need bias-corrected moment updates

Update rule of Adam

$$\widehat{\boldsymbol{m}}^{k+1} = \frac{\boldsymbol{m}^{k+1}}{1 - \beta_1^{k+1}} \qquad \widehat{\boldsymbol{v}}^{k+1} = \frac{\boldsymbol{v}^{k+1}}{1 - \beta_2^{k+1}} \longrightarrow \boldsymbol{\theta}^{k+1} = \boldsymbol{\theta}^k - \alpha \cdot \frac{\widehat{\boldsymbol{m}}^{k+1}}{\sqrt{\widehat{\boldsymbol{v}}^{k+1}} + \epsilon}$$

Adam

• Exponentially-decaying mean and variance of gradients (combines first and second order momentum)



There are a few others...

- 'Vanilla' SGD
- Momentum
- RMSProp
- Adagrad
- Adadelta
- AdaMax
- Nada
- AMSGrad

Adam is mostly method of choice for neural networks!

It's actually fun to play around with SGD updates. It's easy and you get pretty immediate feedback ⓒ

Convergence



Source: <u>http://ruder.io/optimizing-gradient-descent/</u>

Convergence



Source: <u>http://ruder.io/optimizing-gradient-descent/</u>

Convergence



Source: https://github.com/Jaewan-Yun/optimizer-visualization

Jacobian and Hessian

- Derivative $f: \mathbb{R} \to \mathbb{R}$ $\frac{df(x)}{dx}$
- Gradient $f: \mathbb{R}^m \to \mathbb{R}$ $\nabla_x f(x) \qquad \left(\frac{\mathrm{d}f'(x)}{\mathrm{d}x_1}, \frac{\mathrm{d}f'(x)}{\mathrm{d}x_2}\right)$

• Jacobian $f: \mathbb{R}^m \to \mathbb{R}^n$ $J \in \mathbb{R}^{n \times m}$

• Hessian $f: \mathbb{R}^m \to \mathbb{R}$ $\mathbf{H} \in \mathbb{R}^{m \times m}$



• Approximate our function by a second-order Taylor series expansion

$$L(\boldsymbol{\theta}) \approx L(\boldsymbol{\theta}_0) + (\boldsymbol{\theta} - \boldsymbol{\theta}_0)^T \nabla_{\boldsymbol{\theta}} L(\boldsymbol{\theta}_0) + \frac{1}{2} (\boldsymbol{\theta} - \boldsymbol{\theta}_0)^T \mathbf{H} (\boldsymbol{\theta} - \boldsymbol{\theta}_0)$$

First derivative Second derivative (curvature)

• Differentiate and equate to zero

$$\theta^* = \theta_0 - H^{-1} \nabla_{\theta} L(\theta)$$
 Update step
We got rid of the learning rate!

SGD
$$\boldsymbol{\theta}_{k+1} = \boldsymbol{\theta}_k - \alpha \nabla_{\boldsymbol{\theta}} L(\boldsymbol{\theta}_k, \mathbf{x}_i, \mathbf{y}_i)$$

• Differentiate and equate to zero

$$\boldsymbol{\theta}^* = \boldsymbol{\theta}_0 - \mathbf{H}^{-1} \nabla_{\boldsymbol{\theta}} L(\boldsymbol{\theta})$$
 Update step

Parameters of a network (millions)

k

Number of elements in the Hessian

 k^2

Computational complexity of 'inversion' per iteration

$$\mathcal{O}(k^3)$$

• Gradient Descent (green)

 Newton's method exploits the curvature to take a more direct route



Source: https://en.wikipedia.org/wiki/Newton%27s_method_in_optimization

 $J(\boldsymbol{\theta}) = (\mathbf{y} - \mathbf{X}\boldsymbol{\theta})^T (\mathbf{y} - \mathbf{X}\boldsymbol{\theta})$

Can you apply Newton's method for linear regression? What do you get as a result?

BFGS and L-BFGS

- Broyden-Fletcher-Goldfarb-Shanno algorithm
- Belongs to the family of quasi-Newton methods
- Have an approximation of the inverse of the Hessian

$$\boldsymbol{\theta}^* = \boldsymbol{\theta}_0 - \mathbf{H}^{-1} \nabla_{\boldsymbol{\theta}} L(\boldsymbol{\theta})$$

- BFGS $\mathcal{O}(n^2)$
- Limited memory: L-BFGS $\mathcal{O}(n)$

Gauss-Newton

- $x_{k+1} = x_k H_f(x_k)^{-1} \nabla f(x_k)$
 - 'true' 2nd derivatives are often hard to obtain (e.g., numerics)

$$- H_f \approx 2J_F^T J_F$$

- Gauss-Newton (GN): $x_{k+1} = x_k - [2J_F(x_k)^T J_F(x_k)]^{-1} \nabla f(x_k)$
- Solve linear system (again, inverting a matrix is unstable):

$$2(J_F(x_k)^T J_F(x_k))(x_k - x_{k+1}) = \nabla f(x_k)$$

Solve for delta vector

Levenberg

- Levenberg
 - "damped" version of Gauss-Newton:
 - $(J_F(x_k)^T J_F(x_k) + \lambda \cdot I) \cdot (x_k x_{k+1}) = \nabla f(x_k)$

Tikhonov regularization

- The damping factor λ is adjusted in each iteration ensuring: $f(x_k) > f(x_{k+1})$
 - if the equation is not fulfilled increase λ
 - →Trust region
- \rightarrow "Interpolation" between Gauss-Newton (small λ) and Gradient Descent (large λ)

Levenberg-Marquardt

• Levenberg-Marquardt (LM)

 $(J_F(x_k)^T J_F(x_k) + \lambda \cdot diag(J_F(x_k)^T J_F(x_k))) \cdot (x_k - x_{k+1})$ = $\nabla f(x_k)$

- Instead of a plain Gradient Descent for large λ , scale each component of the gradient according to the curvature.
 - Avoids slow convergence in components with a small gradient

Which, What, and When?

• Standard: Adam

• Fallback option: SGD with momentum

• Newton, L-BFGS, GN, LM only if you can do full batch updates (doesn't work well for minibatches!!)

This practically never happens for DL Theoretically, it would be nice though due to fast convergence

General Optimization

- Linear Systems (Ax = b)
 - LU, QR, Cholesky, Jacobi, Gauss-Seidel, CG, PCG, etc.
- Non-linear (gradient-based)
 - Newton, Gauss-Newton, LM, (L)BFGS ← second order
 - Gradient Descent, SGD

← first order

- Others
 - Genetic algorithms, MCMC, Metropolis-Hastings, etc.
 - Constrained and convex solvers (Langrage, ADMM, Primal-Dual, etc.)

Please Remember!

- Think about your problem and optimization at hand
- SGD is specifically designed for minibatch
- When you can, use 2nd order method \rightarrow it's just faster

• GD or SGD is <u>not</u> a way to solve a linear system!

Next Lecture

- This week:
 - Check exercises
 - Check office hours \bigcirc

- Next lecture
 - Training Neural networks



See you next week 🕑

Some References to SGD Updates

- Goodfellow et al. "Deep Learning" (2016),
 Chapter 8: Optimization
- Bishop "Pattern Recognition and Machine Learning" (2006),
 - Chapter 5.2: Network training (gradient descent)
 - Chapter 5.4: The Hessian Matrix (second order methods)
- https://ruder.io/optimizing-gradient-descent/index.html
- PyTorch Documetation (with further readings)
 - <u>https://pytorch.org/docs/stable/optim.html</u>