

Lecture 11: Optimization

Goals for Today:

Fundamental Properties:

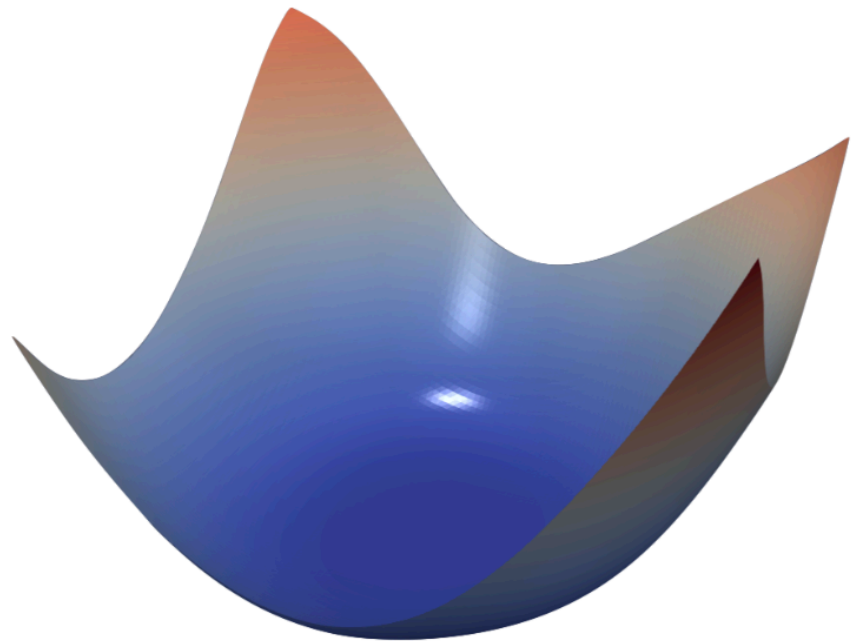
- Convexity
- Smoothness

Gradient Descent:

- Basic implementation
- Obvious drawbacks
- Learning Schedules

Types of Acceleration:

- Stochastic Gradient
- Momentum
- Adaptive Learning



What is optimization?

Find a set of parameters, β , that minimizes $f(\mathbf{x}, \mathbf{y}, \beta)$: $\arg \min_{\beta} f(\mathbf{x}, \mathbf{y}, \beta)$

For example: in ordinary least squares, $f(\mathbf{x}, \mathbf{y}, \beta)$ would be the mean squared error, \mathbf{x} would be the set of predictor variables, \mathbf{y} would be the known values, and β would be the set of coefficients and intercept that minimize $f(\mathbf{x}, \mathbf{y}, \beta)$.

For a narrow class of $f(\mathbf{x}, \mathbf{y}, \beta)$ we can analytically solve for β (ordinary linear least squares and the common variants that we discussed last time), but for non-linear $f(\mathbf{x}, \mathbf{y}, \beta)$ we need to use numerical optimization to find β .

Optimizers: these are the algorithms that we will use for finding the β that minimize our objective function. The most common algorithms involve following gradients, but this is still an active area of research with new specialized algorithms regularly appearing.

Objective function: this is the $f(\mathbf{x}, \mathbf{y}, \beta)$ that we are trying to minimize. When this function refers to an error of some kind (e.g., mean squared error) it is sometimes also called a **loss function** and we will use the notation $L(\mathbf{x}, \mathbf{y}, \beta)$.

What is optimization?

An important distinction: when discussing optimization, β will be reserved for the parameters that we are trying to determine and \mathbf{x} will usually be considered fixed by the problem (e.g., the sensor values in the Dow dataset). If we were trying to find the minimum of a mathematical function (e.g., $f(\mathbf{x}) = \mathbf{x}^2 - e^{\mathbf{x}}$) then the “ \mathbf{x} ” value that minimizes $f(\mathbf{x})$ would play the role of β .

Parameterization
(Loss function minimization) :

$$\arg \min_{\beta} f(\mathbf{x}, \mathbf{y}, \beta)$$

VS.

Function minimization :

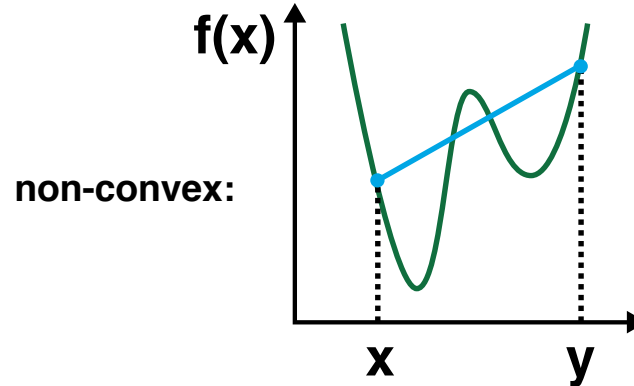
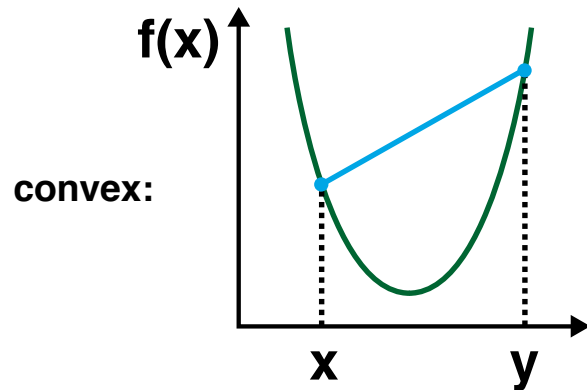
$$\arg \min_{\mathbf{x}} f(\mathbf{x})$$

There is no fundamental distinction in the way these two problems are solved, other than that in loss function minimization, only a subset of the parameters are being optimized (β).

Convex vs. Non-Convex

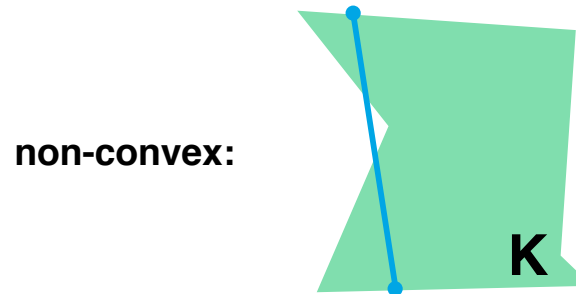
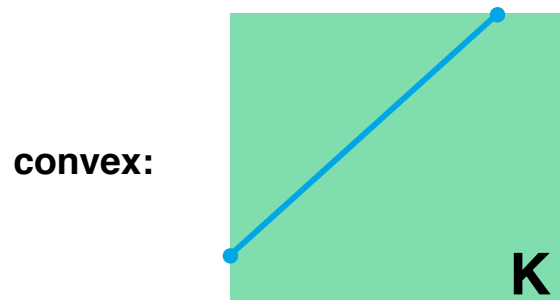
A function, f , is convex if for any \mathbf{x} and \mathbf{y} in \mathbf{R}^n (Jensen's Inequality):

$$\forall \lambda \in [0, 1], f(\lambda \mathbf{x} + (1 - \lambda)\mathbf{y}) \leq \lambda f(\mathbf{x}) + (1 - \lambda)f(\mathbf{y})$$



A set, K , is convex if for any \mathbf{x} and \mathbf{y} in K :

$$\forall \lambda \in [0, 1], \lambda \mathbf{x} + (1 - \lambda)\mathbf{y} \in K$$



A function, f , of a random variable \mathbf{x} is convex if:

$$f(E[X]) \leq E(f[X])$$

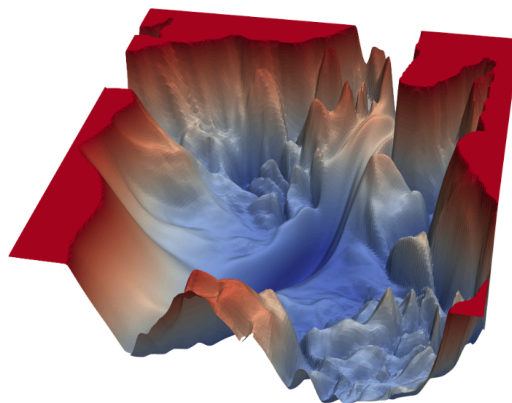
Complex models are typically non-convex

Which would you want to optimize?

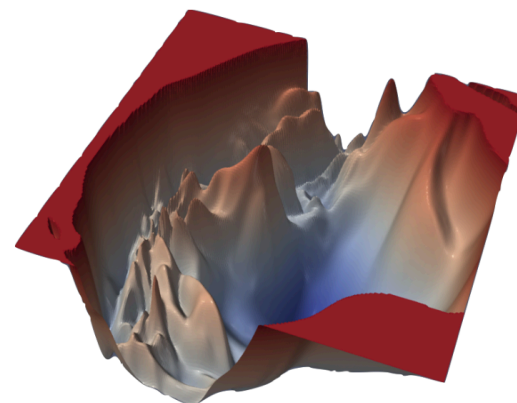
Which performs the best?

A lot of deep learning research circulates around the fact that the two questions don't necessarily have the same answer.

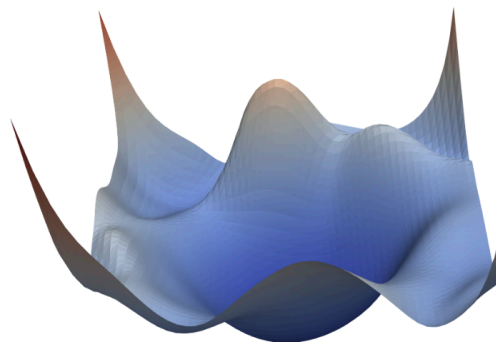
VGG-56



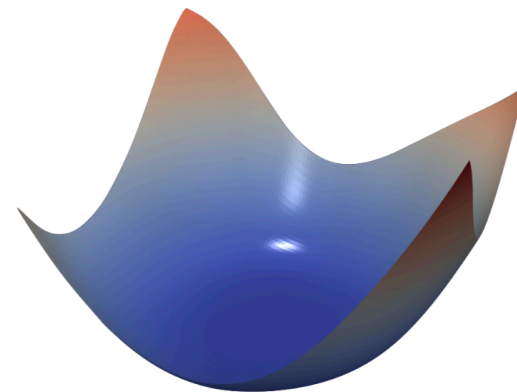
VGG-110



Renset-56



Densenet-121



3D projection of a high dimensional loss surface associated with 4 neural networks for image recognition tasks. The lower two are designed to have relatively better convexity and smoothness. Easier to optimize but not necessarily better. <https://www.cs.umd.edu/~tomg/projects/landscapes/>

Smooth vs. Non-Smooth

- For **smooth** functions, all differentials are continuous.

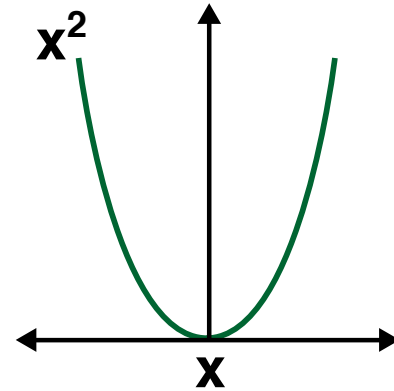
- For **non-smooth** functions at least one differential is discontinuous.

- **Lipschitz continuity** generalizes the degree of smoothness to continuous functions.

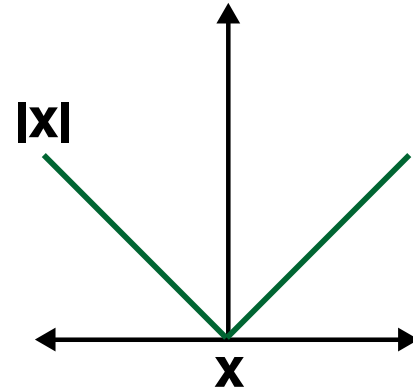
$$\frac{||f(\mathbf{x}_1) - f(\mathbf{x}_2)||}{||\mathbf{x}_1 - \mathbf{x}_2||} \leq K$$

- Functions that are non-smooth and non-convex are formally hopeless to optimize.

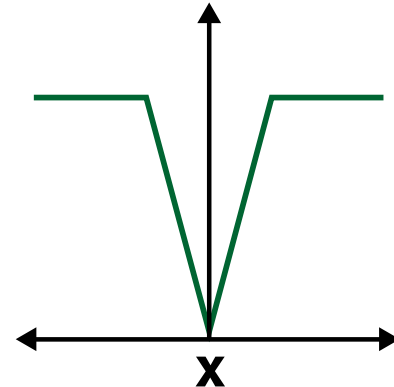
smooth
and
convex



non-smooth
and
convex



non-smooth
and
non-convex



Types of Objective Functions

Convexity and smoothness are determined by both the model and loss function that we are trying to optimize.

We'll use the notation $\mathbf{L}()$ for a loss function, $f()$ for our model, \mathbf{x} for our features, $\boldsymbol{\beta}$ for our parameters, and \mathbf{y} for our training data vector.

$$\mathbf{L} [f(\mathbf{x})]$$

e.g., function minimization

$$\mathbf{L} [f(\mathbf{x}), \mathbf{x}]$$

e.g., function minimization with constraints

$$\mathbf{L} [f(\mathbf{x}, \boldsymbol{\beta}), \mathbf{y}]$$

e.g., least-squares optimization

Least squares with L_2 regularization example:

$$\mathbf{L} [f(\mathbf{x}, \boldsymbol{\beta}), \boldsymbol{\beta}, \mathbf{y}] = \frac{1}{N} \sum_i^N [f(\mathbf{x}_i, \boldsymbol{\beta}) - \mathbf{y}_i]^2 + \alpha \|\boldsymbol{\beta}\|_2$$

We seek $\boldsymbol{\beta}^*$ such that:

$$\mathbf{L} [f(\mathbf{x}, \boldsymbol{\beta}^*), \mathbf{y}] = \min (\mathbf{L} [f(\mathbf{x}, \boldsymbol{\beta}), \mathbf{y}])$$

How do we know if we have found a minimum?

$$\|\nabla \mathbf{L} [f(\mathbf{x}, \boldsymbol{\beta}^*), \mathbf{y}]\| = 0 \quad \|\nabla^2 \mathbf{L} [f(\mathbf{x}, \boldsymbol{\beta}^*), \mathbf{y}]\| > 0$$

Numerically,

$$\|\nabla \mathbf{L} [f(\mathbf{x}, \boldsymbol{\beta}^*), \mathbf{y}]\| \leq \epsilon_1 \quad \|\nabla^2 \mathbf{L} [f(\mathbf{x}, \boldsymbol{\beta}^*), \mathbf{y}]\| \geq \epsilon_2$$

For non-convex problems, $\boldsymbol{\beta}^*$ is probably not a global minimum

Formulating an Optimization Problem

- Let's say we have a process with sensor readings for n variables, \mathbf{x} , and binary labelled historical data, \mathbf{y} , indicating when the process yielded a product that passed inspection (i.e., 1 or -1 for pass/fail). We've collected t samples (each containing \mathbf{x}, \mathbf{y}) that we want to use to train a model.
- Our goal is to find a linear separator, $\boldsymbol{\beta} \in \mathbb{R}^n$ (a hyperplane with $\boldsymbol{\beta}$ defining the normal) that divides conditions \mathbf{x}_{pass} that yield 1 and \mathbf{x}_{fail} that yield -1. If no clean linear separator exists, then we want to minimize the number of misclassifications. **What should we use for our loss function?**

We might use the dot product of the normal and \mathbf{x} :

$$\mathbf{L} [f(\mathbf{x}, \boldsymbol{\beta}), \mathbf{y}] = \frac{1}{t} \sum_i^t \delta (\text{sign} [\boldsymbol{\beta}^T \mathbf{x}] \neq y_i)$$

This loss captures our definition, but it is **non-smooth** and **non-convex**.

An alternative, (mostly) smooth and convex loss function would be:

$$\mathbf{L} [f(\mathbf{x}, \boldsymbol{\beta}), \mathbf{y}] = \frac{1}{t} \sum_i^t \max \left\{ 0, 1 - \underbrace{y_i \boldsymbol{\beta}^T \mathbf{x}}_{\text{Correct classification yields (+)}} \right\}$$

Hinge Loss

Gradient Descent

The vanilla gradient descent algorithm is the parent of the more complicated algorithms that are used in this course:

$$\beta_{i+1} = \beta_i - \alpha \nabla_{\beta} \mathbf{L}$$

where, β is a vector holding the parameters of in the model

α is a scalar called the “learning rate” or “step-size”

\mathbf{L} is the loss function

$\nabla_{\beta} \mathbf{L}$ is gradient of loss function with respect to β

For a typical supervised example:

$$\nabla_{\beta} \mathbf{L} [f(\beta, \mathbf{x}), \mathbf{y}] = \frac{1}{N} \sum_i^N \nabla_{\beta} \mathbf{L} [f(\beta, \mathbf{x}_i), \mathbf{y}_i]$$

Loss compares predictions against specific samples.

For function minimization:

$$\nabla_{\beta} \mathbf{L} [f(\beta, \mathbf{x}), \mathbf{y}] \rightarrow \nabla_{\mathbf{x}} f(\mathbf{x})$$

For root finding/function minimization, there are no samples, $f(\mathbf{x})$ is itself the loss function (same distinction as we made earlier).

Problem 1: the gradient evaluation can be expensive for lots of samples

Problem 2: small α leads to slow convergence.

Problem 3: large α means oscillations/instability

Problem 4: neglects obvious sources of info (previous iterations, second derivative)

Stochastic Gradient Descent

Stochastic gradient descent speeds up gradient descent by approximating the gradient evaluation using a **single** randomly chosen sample out of N .

$$\underbrace{\beta_{i+1} = \beta_i - \frac{\alpha}{N} \sum_j^N \nabla_{\beta} \mathbf{L} [\beta_i, \mathbf{x}_j, \mathbf{y}_j]}_{\text{Gradient Descent}} \rightarrow \underbrace{\beta_{i+1} = \beta_i - \alpha \nabla_{\beta} \mathbf{L} [\beta_i, \mathbf{x}_j, \mathbf{y}_j]}_{\text{Stochastic Gradient Descent}}$$

- We want a cheap **unbiased estimator** for the gradient. We get the algorithm from the fact that $E(\nabla_{\beta} \mathbf{L} [f(\beta, \mathbf{x}_j), \mathbf{y}_j]) = \nabla_{\beta} \mathbf{L} [f(\beta, \mathbf{x}), \mathbf{y}]$ when j is chosen from a random uniform distribution across N samples.
- Further generalizing:

$$\underbrace{\beta_{i+1} = \beta_i - \frac{\alpha}{N} \sum_j^N \nabla_{\beta} \mathbf{L} [\beta_i, \mathbf{x}_j, \mathbf{y}_j]}_{\text{Gradient Descent}} \rightarrow \underbrace{\beta_{i+1} = \beta_i - \frac{\alpha}{M} \sum_j^{M < N} \nabla_{\beta} \mathbf{L} [\beta_i, \mathbf{x}_j, \mathbf{y}_j]}_{\substack{\text{Mini-Batch Gradient Descent} \\ (m=1, \text{ we recover SGD})}}$$

- SGD and mini-batch GD largely solve the issue of evaluating the gradient with respect to large numbers of samples and are still the default in many deep learning applications.

Problem 1: Still haven't solved the issue with choosing α .

Problem 2: Still neglects potential info (previous iterations, second derivative)

Learning Schedules

Take a moment and ask yourself what you might want α to do in a generic problem

Start high and end low?

This idea is commonly referred to as “Learning Rate Annealing” and has many flavors

1) 1/t Decay:

$$\alpha_t = \frac{\alpha_0}{1 + rt}$$

2) Exponential Decay:

$$\alpha_t = \alpha_0 e^{-rt}$$

3) Step Decay:
(0-indexing)

$$\alpha_t = \alpha_0 r^{\text{floor}\left(\frac{1+t}{d}\right)}$$

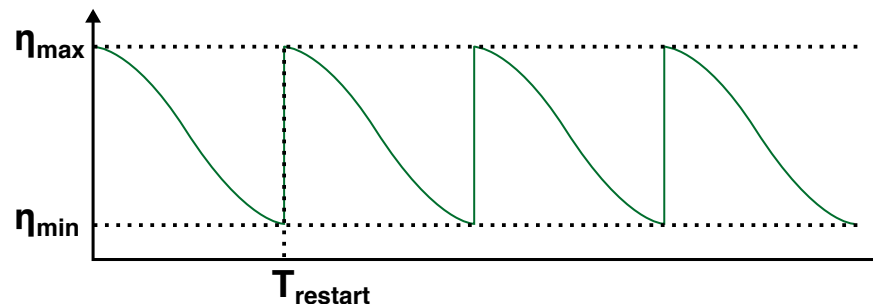
where, α_0 is an initial learning rate
 r is a decay rate
 d is decay interval in steps

All achieve the qualitative high \rightarrow low behavior necessary to heuristically increase convergence while permitting discovery of “deep” minima.

4) Warm Restarting:

$$\alpha_t = \alpha_{\min} + \frac{1}{2} (\alpha_{\max} - \alpha_{\min}) \left(1 + \cos \left[\frac{t_{i-0}\pi}{t_i} \right] \right) \quad \begin{array}{l} t_{i-0} \text{ intervals} \\ \text{since last} \\ \text{restart} \end{array}$$

Heuristic for finding deeper and broader minima



Momentum Variants

The intuition behind momentum is that the previous gradient evaluation(s) potentially give us good information about the direction of the minimum. The analogy comes from the physics of an object rolling down the hill. It will start slowly, but if it is consistently downhill in the same direction, then the object “accelerates” toward the minimum. The change in algorithm is tiny but powerful:

$$\mathbf{m}_{i+1} = \eta \mathbf{m}_i - \alpha \nabla_{\beta} \mathbf{L} [\beta_i]^*$$

$$\beta_{i+1} = \beta_i + \mathbf{m}_{i+1}$$

where, η is a new hyperparameter and we recover GD/SGD when $\eta=0$.
(typical values [0.5,0.9])

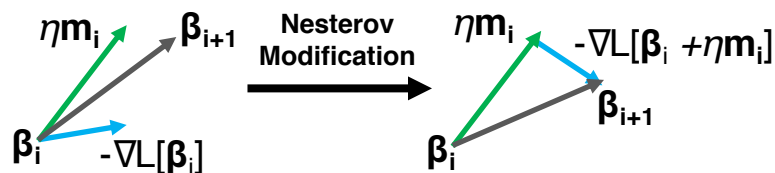
*not showing x and y dependence

An even better variant developed by Nesterov calculates the gradient at $\beta_i + \eta \mathbf{m}_i$ (i.e., *after* applying the momentum displacement). That is, on step $i+1$ we know the update from momentum alone ($\beta \mathbf{m}_i$) before we calculate the gradient—so why not calculate the gradient at that point, $\beta_i + \eta \mathbf{m}_i$, since it is probably closer to β_{i+1} than β_i is?

$$\mathbf{m}_{i+1} = \eta \mathbf{m}_i - \alpha \nabla_{\beta} \mathbf{L} [\beta_i + \eta \mathbf{m}_i]$$

$$\beta_{i+1} = \beta_i + \mathbf{m}_{i+1}$$

Nesterov Momentum



Adaptive Learning Rates

“Adaptive” meaning parameter-specific learning rates that are dynamically updated during the optimization depending on the optimization surface.

AdaGrad (adaptive gradient algorithm, diagonal version)

$$\mathbf{A}_{i+1} = \mathbf{A}_i + \text{diag}(\nabla_{\beta} \mathbf{L})^2$$

Diagonal matrix with cumulative square of the update derivatives

$$\beta_{i+1} = \beta_i - \frac{\alpha}{\sqrt{\mathbf{A}_{i+1}} + \text{diag}(\varepsilon)} \nabla_{\beta} \mathbf{L}$$

Large gradients get damped and small get accelerated

Note 1: Monotonic decrease in the learning rate makes it susceptible to early termination.

Note 2: eps is a small fixed value (e.g., 1E-4) to avoid division by zero.

Note 3: sqrt is actually important. What would be the problem with full matrix AdaGrad?

RMSProp (root mean square propagation)

$$\mathbf{A}_{i+1} = \eta \mathbf{A}_i + (1 - \eta) \text{diag}(\nabla_{\beta} \mathbf{L})^2$$

mixes in a fraction 1- η of the current gradients

$$\beta_{i+1} = \beta_i - \frac{\alpha}{\sqrt{\mathbf{A}_{i+1}} + \text{diag}(\varepsilon)} \nabla_{\beta} \mathbf{L}$$

Note 1: Averaging of the previous gradients addresses early stopping issues.

Note 2: η is a new hyperparameter that is typically fixed (e.g., 0.9, 0.99, 0.999 etc.)

Note 3: In principle, this addresses the early stopping issues of AdaGrad.

Extensions Under Active Research

AdaDelta: Uses an exponential average of past gradients for the update step.

Adam: Adds a sliding window to AdaGrad and a form of momentum

Shampoo: Full matrix AdaGrad is in principle better than the diagonal version, but square roots of non-diagonal matrices are expensive. Shampoo interpolates between full-matrix and diagonal AdaGrad.

AdaFactor: Shampoo + reduced memory requirements (good for really large neural networks)

GGT: Uses recent gradients to approximate full matrix AdaGrad

Extreme Tensoring: Approximates full-matrix AdaGrad using a low-rank tensor

Recommended Reading:

Introductory Lectures on Convex Optimization, **Yurii Nesterov** (2004)