

An Introduction to Optimisation

Alasdair Paren - alasdair.paren@eng.ox.ac.uk

What is Optimisation?



The area of maths behind the selection of a best element from some set of candidates, with respect to some criterion or objective



Why is Optimisation Important de Engineering To Machine Learning?

- Most State-Of-The-Art (SOTA) AI systems employ large neural networks trained on large data sets
- Training neural networks requires optimising a very high dimensional, nonconvex optimisation problem, computed over a very large training set

Visualization of a Neural Network Loss Function



Loss Functions



We want to optimise (typically minimise) scalar valued "loss functions" that quantify the error of the model

Desirable properties:

- Bounded below (typically by zero)
- Continuous
- Smooth
- Convex
- Cheap to calculate
- Finite-sum Structure

Types of Optimisation



- Convex / Non-Convex
- Constrained / Non-Constrained
- Discrete / Continuous / Mixed Integer Programming
- Stochastic / Deterministic
- Gradient Free / First Order / Second Order
- Reinforcement Learning

Quiztime (1/4)



If
$$f(w) = w^2 - 4w + 16$$
, what is $\operatorname{argmin}_{w \in \mathbb{R}} f(w)$?
1. $w = -4$
2. $w = -2$
3. $w = 2$

4.
$$w = 4$$

Quiztime(2/4)



If
$$f(w) = w^2 - 4w + 16$$
, whats is min $f(w)$?
1. $f(w) = -8$
2. $f(w) = 0$
3. $f(w) = 8$
4. $f(w) = 12$

Quiztime (3/4)



If $f(w) = \max\{-w, w, \frac{1}{2}(w+3)\}$ (point-wise maximum), what is $\operatorname{argmin}_{w \in \mathbb{R}} f(w)$? 1. -2 2. -1 3. 1 4. 2

Quiztime (4/4)



If $\mathbf{x}^{\top} A \mathbf{x} \ge 0$ for all none zero \mathbf{x} , and a square matrix $A \in \mathbb{R}^{d \times d}$. We call A:

- 1. Full Rank
- 2. A Negative Definite Matrix
- 3. A Negative Semi-Definite Matrix
- 4. A Positive Semi-Definite Matrix

Positive Definiteness



2.1 Positive Definiteness

A square matrix $A \in \mathbb{R}^{d \times d}$ is positive definite if for all none zero $\mathbf{x}, \mathbf{x}^{\top} A \mathbf{x}$ is positive. Formally:

 $\forall \mathbf{x} \in \mathbb{R}^d \backslash 0^d, \quad \mathbf{x}^\top A \mathbf{x} > 0.$

If the above inequality hold in equality A is known as positive semi-definite

Multivariate Functions



For a scalar output multivariate function with *d* inputs what is the Hessian?

Multivariate Functions



2.1 Scalar Output Multivariate Functions

Let us consider a multivariate function $f(\mathbf{w})$:

 $f: \mathbb{R}^d \to \mathbb{R}.$

$$\nabla f \triangleq \begin{bmatrix} \frac{\partial f}{\partial w_1} \\ \frac{\partial f}{\partial w_2} \\ \vdots \\ \frac{\partial f}{\partial w_d} \end{bmatrix} \qquad \qquad H_f \triangleq \begin{bmatrix} \frac{\partial^2 f}{\partial w_1}^2 & \frac{\partial^2 f}{\partial w_1 \partial w_2} & \cdots & \frac{\partial^2 f}{\partial w_1 \partial w_d} \\ \frac{\partial^2 f}{\partial w_2 w_1} & \frac{\partial^2 f}{\partial w_2 \partial w_2} & \cdots & \frac{\partial^2 f}{\partial w_2 \partial w_d} \\ \vdots & \vdots & \ddots & \vdots \\ \frac{\partial^2 f}{\partial w_d w_1} & \frac{\partial^2 f}{\partial w_d \partial w_2} & \cdots & \frac{\partial^2 f}{\partial w_d^2} \end{bmatrix}$$



Suppose we have a scalar function $f = \mathbf{x}^{\top} \mathbf{y}$, where $\mathbf{x}, \mathbf{y} \in \mathbb{R}^d$ then:

 $\nabla_{\mathbf{x}} f(\mathbf{x}) = \mathbf{y}.$

Now suppose we have a scalar function $f = \mathbf{x}^{\top} \mathbf{x}$, where $\mathbf{x} \in \mathbb{R}^d$ then:

 $\nabla_{\mathbf{x}} f(\mathbf{x}) = 2\mathbf{x}$

If in doubt calculate the gradient for one element and then construct the vector of partial derivatives:

$$f(\mathbf{x}) = \sum_{i} x_i^2$$

$$\frac{\partial f(\mathbf{x})}{\partial x_i} = 2x_i$$

$$abla_{\mathbf{x}}f(\mathbf{x}) = egin{bmatrix} rac{\partial f}{\partial x_1} \ rac{\partial f}{\partial x_2} \ dots \ rac{\partial f}{\partial x_d} \end{bmatrix} = egin{bmatrix} 2x_1 \ 2x_2 \ dots \ 2x_d \end{bmatrix} = 2\mathbf{x}$$



If $\mathbf{x} \in \mathbb{R}^d, H \in \mathbb{R}^{d \times d}$ and symmetrical, What do you think the derivative of the following is?

 $f(\mathbf{x}) = \mathbf{x}^\top H \mathbf{x}, \qquad \nabla_x f(\mathbf{x}) = ?$

https://www.math.uwaterloo.ca/~hwolkowi/matrixcookbook.pdf



$$f(\mathbf{x}) = \mathbf{x}^{\top} H \mathbf{x} = \begin{bmatrix} x_1 & \cdots & x_d \end{bmatrix}^T \begin{bmatrix} h_{1,1} & \cdots & h_{1,d} \\ \vdots & \ddots & \vdots \\ h_{d,1} & \cdots & h_{d,d} \end{bmatrix} \begin{bmatrix} x_1 \\ \vdots \\ x_d \end{bmatrix} = \begin{bmatrix} \cdots & \sum_i x_i h_{i,j} & \cdots \end{bmatrix} \begin{bmatrix} x_1 \\ \vdots \\ x_d \end{bmatrix} = \sum_i \sum_j x_j x_i h_{i,j}$$
$$\sum_i \sum_j x_j x_i h_{i,j} = x_1 x_1 h_{1,1} + x_1 x_2 h_{1,2} + x_1 x_3 h_{1,3} + \cdots + x_1 x_d h_{1,d},$$
$$+ x_2 x_1 h_{2,1} + x_2 x_2 h_{2,2} + x_2 x_3 h_{2,3} \cdots + x_2 x_d h_{2,d},$$
$$+ x_3 x_1 h_{3,1} + x_3 x_2 h_{3,2} + x_3 x_3 h_{3,3} \cdots + x_3 x_d h_{3,d},$$
$$+ \cdots,$$
$$+ x_d x_1 h_{d,1} + x_d x_2 h_{d,2} + x_d x_3 h_{d,3} + \cdots + x_d x_d h_{d,d}.$$

https://www.math.uwaterloo.ca/~hwolkowi/matrixcookbook.pdf

$$\frac{\partial}{\partial x_i} \left(\sum_i \sum_j x_j x_i h_{i,j} \right) = 2x_i h_{i,i} + x_2 h_{i,2} + x_3 h_{i,3} + \dots + x_d h_{i,d},$$
$$+ x_2 h_{2,i} + 0 + 0 \dots + 0,$$
$$+ x_3 h_{3,i} + 0 + 0 \dots + 0,$$
$$\dots$$
$$+ x_d h_{d,i} + 0 + 0 + \dots + 0.$$

As H is symmetrical:

$$\frac{\partial}{\partial x_i} \left(\sum_i \sum_j x_j x_i h_{i,j} \right) = 2x_i h_{i,i} + 2x_2 h_{i,2} + 2x_3 h_{i,3} + \dots + 2x_d h_{i,d}$$
$$= 2H_{i,i} \mathbf{x}$$

Thus considering the derivative $w.r.t \mathbf{x}$ rather than x_i :

$$rac{\partial}{\partial \mathbf{x}} \left(\sum_i \sum_j x_j x_i h_{i,j}
ight) = egin{bmatrix} 2H_{1,:}\mathbf{x} \ 2H_{2,:}\mathbf{x} \ dots \ 2H_{d,:}\mathbf{x} \end{bmatrix} = 2H\mathbf{x}$$

https://www.math.uwaterloo.ca/~hwolkowi/matrixcookbook.pdf

Taylor Expansions



2.3 Taylor Expansions

A function f at a point \mathbf{w}_t can be approximated by its **first** order Taylor expansion around this point:

$$f(\mathbf{w}) \approx f(\mathbf{w}_t) + \nabla f(\mathbf{w}_t)^{\top} (\mathbf{w} - \mathbf{w}_t)$$

A function f at a point \mathbf{w}_t can be approximated by its **second** order Taylor expansion around this point:

$$f(\mathbf{w}) \approx f(\mathbf{w}_t) + \nabla f(\mathbf{w}_t)^{\top} (\mathbf{w} - \mathbf{w}_t) + \frac{1}{2} (\mathbf{w} - \mathbf{w}_t)^T \boldsymbol{H}_f (\mathbf{w} - \mathbf{w}_t)$$

Convex Functions





Convex Functions



2.4 Convex Function

A function f is convex if for any \mathbf{x}, \mathbf{y} :

$$f(\mathbf{x}) \ge f(\mathbf{y}) + \nabla f(\mathbf{y})^{\top} (\mathbf{x} - \mathbf{y}),$$

A function f is strictly convex if for any \mathbf{x}, \mathbf{y} :

$$f(\mathbf{x}) > f(\mathbf{y}) + \nabla f(\mathbf{y})^{\top} (\mathbf{x} - \mathbf{y}),$$

A twice differentiable convex function will have a positive semi-definite Hessian. A twice differentiable strictly convex function will have a positive definite Hessian.

Convex Sets





Convex Sets



3.6 Convex Set

A set Ω is convex if for any $\mathbf{x}, \mathbf{y} \in \Omega$ and any $\lambda \in [0, 1]$:

 $\lambda \mathbf{x} + (1 - \lambda) \mathbf{y} \in \Omega$

Alternatively:

A set Ω is convex if there exists a convex function f that $\forall \mathbf{x} \in \Omega \ f(\mathbf{x}) \leq k$ and for all $\mathbf{x} \notin \Omega \ f(\mathbf{x}) > k$ where k is some constant.

Convex Sets



Euclidean Projection let us define the euclidean Projection operation onto a set Ω as:

$$\Pi_{\Omega}(\mathbf{w}_t) = \min_{\mathbf{w}\in\Omega} \|\mathbf{w} - \mathbf{w}_t\|^2.$$

Euclidean Projection onto a convex set decreases the distance to every point in the set, thus it always reduces the distance to an constrain optimum. Formally if we define:

$$\mathbf{w}_{\star} = \operatorname*{argmin}_{\mathbf{w} \in \omega} f(\mathbf{w})$$

Then we have for all \mathbf{w} in \mathbb{R}^d :

$$|\mathbf{w}_{\star} - \mathbf{w}||_2^2 \ge ||\mathbf{w}_{\star} - \Pi_{\Omega}(\mathbf{w})||_2^2$$

Lipschitz Continuity



A function f is C-Lipschitz over a set Ω with respect to a norm $||\cdot||$ if for any $\mathbf{x}, \mathbf{y} \in \Omega$:

 $||f(\mathbf{x}) - f(\mathbf{y})|| \le C||\mathbf{y} - \mathbf{x}||.$

Most commonly started with reference to the ℓ_1 norm, or:

 $|f(\mathbf{x}) - f(\mathbf{y})| \le C|\mathbf{y} - \mathbf{x}|.$

Alternatively:

 $\nabla f(\mathbf{x}) \le C, \ \forall \ \mathbf{x} \in \Omega$

Lipschitz Continuity





Why is Lipschitz Important? (DEPARTMENT OF ENGINEERING SCIENCE UNIVERSITY OF

Q



Smoothness



A function f is β -Smooth over a set Ω with respect to a norm $||\cdot||$ if for any $\mathbf{x}, \mathbf{y} \in \Omega$:

$$||\nabla f(\mathbf{x}) - \nabla f(\mathbf{y})|| \le \beta ||\mathbf{y} - \mathbf{x}||.$$

This is normally defined in terms of the ℓ_1 norm:

 $|\nabla f(\mathbf{x}) - \nabla f(\mathbf{y})| \le \beta |\mathbf{y} - \mathbf{x}|.$

Alternatively:

$$\forall \mathbf{x}, \mathbf{y} \in \Omega, \ |f(\mathbf{x}) - f(\mathbf{y}) - \nabla f(\mathbf{y})^{\top} (\mathbf{x} - \mathbf{y})| \leq \frac{\beta}{2} \|\mathbf{x} - \mathbf{y}\|^2.$$

Finally if f is twice differentiable, then f is β -smooth if and only if for any $\mathbf{x}, \mathbf{y} \in \mathbb{R}^d$:

 $H\mathbf{x}^{\top}\mathbf{y} \le \beta ||\mathbf{x}||||\mathbf{y}||.$

Smoothness



Do you think the below function is smooth?





Optimisation Problems

The formulation you will find in the literature



min $f(\mathbf{w})$ s.t. $g_i(\mathbf{w}) = 0$ for i = 1, ..., ns.t. $h_j(\mathbf{w}) \ge 0$ for i = 1, ..., m

Equality constraints Inequality constraints

Optimisation Standard Set Up



$\begin{array}{ll} \min & f(\mathbf{w}),\\ \text{s.t.} & \mathbf{w} \in \Omega \end{array}$

 $\Omega \subseteq \mathbb{R}^d,$ $f: \Omega \to \mathbb{R}.$



Convex Optimisation

 $\begin{array}{ll} \min & f(\mathbf{w}) \\ \text{s.t.} & \mathbf{w} \in \Omega \end{array}$

$$\begin{split} \Omega &\subseteq \mathbb{R}^d, \\ \Omega \text{ is a convex set,} \\ f &: \Omega \to \mathbb{R}, \\ f \text{ is a convex function.} \end{split}$$



Linear Programming

 $\begin{array}{ll} \min \quad \boldsymbol{c}^{\top} \boldsymbol{w}, \\ \text{s.t.} \quad A \boldsymbol{w} \leq \boldsymbol{b}, \\ \boldsymbol{w} \geq \boldsymbol{0}, \\ \text{and} \quad \boldsymbol{w} \in \mathbb{R}^{d}. \end{array}$

$$f(\mathbf{w}) = \mathbf{c}^{\top}\mathbf{w}, \quad \Omega = \{\mathbf{w} \mid A\mathbf{w} \in \mathbb{R}^d, \ \mathbf{w} \leq \mathbf{b}, \ ext{and} \ \mathbf{w} \geq 0\}$$



Discrete Optimisation

 $\begin{array}{ll} \min & f(\mathbf{w}) \\ \text{s.t.} & \mathbf{w} \in \Omega \end{array}$

$$\begin{split} &\Omega \subseteq \mathbb{R}^d,\\ &\text{Where }\Omega \text{ is a discrete set. For example }\Omega = \{0,1\},\\ &f:\Omega \to \mathbb{R}. \end{split}$$



Integer Programming

$$\begin{array}{ll} \min \quad \boldsymbol{c}^{\top} \mathbf{w} \\ \text{s.t.} \quad A \mathbf{w} + \mathbf{s} \leq \boldsymbol{b}, \\ \mathbf{w} \geq 0, \\ \mathbf{s} \geq 0, \\ \text{and} \quad \mathbf{w} \in \mathbb{Z}^d. \end{array}$$

$$f(\mathbf{w}) = \mathbf{c}^{\top} \mathbf{w}, \quad \Omega = \{\mathbf{w} \mid A\mathbf{w} + \mathbf{s} \in \mathbb{Z}^d, \ \mathbf{w} \le \mathbf{b}, \text{ and } \mathbf{w} \ge 0, \mathbf{s} \ge 0\},$$
$$\mathbb{Z} \triangleq \{\dots, -3, -2, -1, 0, 1, 2, 3, \dots\}$$

Focus of this Lecture:



Unconstrained Optimisation

min $f(\mathbf{w})$, s.t. $\mathbf{w} \in \mathbb{R}^d$, $f : \mathbb{R}^d \to \mathbb{R}$.

Finite Sum Loss Functions (#)

DEPARTMENT OF



Finite Sum loss functions 5.11

$$f(\mathbf{w}) \triangleq rac{1}{N} \sum_{z=1}^{N} \ell_z(\mathbf{w}) pprox \mathbb{E}_{z \in \mathcal{Z}}[\ell_z(\mathbf{w})]$$

Can you think of a type of loss that would not naturally exhibit a finite sum structure?

Squared Loss

Squared Loss :
$$\ell_z(\mathbf{y}_z, \mathbf{y}_z^*) = ||\mathbf{y} - \mathbf{y}^*||_2^2$$

Cross Entropy

Cross Entropy Loss :
$$\ell_z(\mathbf{y}_z, \mathbf{y}_z^*) = -\sum_{c \in C} y_c^* \log(y_c)$$

For the cross entropy loss both y_z and y_z^* should be vectors denoting a probability distribution $(\sum_i y_{z,i} =$ 1, $y_{z,i} \ge 0, \forall i$). How can we ensure we have this property? One choice is the softmax.

$$\operatorname{Softmax}(\mathbf{x}) = \frac{\exp \mathbf{x}}{\sum_{i} \exp x_{i}}$$

Regularisation



5.12 Regularisation Functions

 $F(\mathbf{w}) = f(\mathbf{w}) + \lambda R(\mathbf{w})$

L2 Regularisation - Weight Decay

L2 Regularisation : $R_{\ell_1}(\mathbf{w}) = ||\mathbf{w}||_2^2$

L1 Regularisation - Lasso

L1 Regularisation : $R_{\ell_2}(\mathbf{w}) = ||\mathbf{w}||_1$



Optimisation Algorithms

Types of Optimiser



- Gradient Free Only uses function value information
- First Order Methods use function value and gradient information
- Second order methods use function value, gradient and Hessian information
- Stochastic Optimisers only use approximate function information which has been evaluated on a subset of the training data set

+ Many others

Gradient Descent (GD)





$$\mathbf{w}_{t+1} = \operatorname*{argmin}_{\mathbf{w} \in \mathbb{R}^d} \left\{ f(\mathbf{w}_t) +
abla f(\mathbf{w}_t)^{ op} (\mathbf{w} - \mathbf{w}_t)
ight\}.$$

Gradient Descent (GD)



$$\mathbf{w}_{t+1} = \operatorname*{argmin}_{\mathbf{w}\in\Omega} \left\{ \frac{1}{2\eta_t} \|\mathbf{w} - \mathbf{w}_t\|^2 + f(\mathbf{w}_t) + \nabla f(\mathbf{w}_t)^\top (\mathbf{w} - \mathbf{w}_t) \right\}.$$



Wt



Gradient Descent(GD)



$$\mathbf{w}_{t+1} = \underset{\mathbf{w} \in \Omega}{\operatorname{argmin}} \left\{ \frac{1}{2\eta_t} \| \mathbf{w} - \mathbf{w}_t \|^2 + f(\mathbf{w}_t) + \nabla f(\mathbf{w}_t)^\top (\mathbf{w} - \mathbf{w}_t) \right\}$$

W

f(w)

Gradient Descent



5.1 Gradient Descent

$$\mathbf{w}_{t+1} = \operatorname*{argmin}_{\mathbf{w}\in\Omega} \left\{ \frac{1}{2\eta_t} \|\mathbf{w} - \mathbf{w}_t\|^2 + f(\mathbf{w}_t) + \nabla f(\mathbf{w}_t)^\top (\mathbf{w} - \mathbf{w}_t) \right\}.$$

$$\begin{split} &\frac{\partial}{\partial \mathbf{w}} \left(\frac{1}{2\eta_t} \| \mathbf{w} - \mathbf{w}_t \|^2 + f(\mathbf{w}_t) + \nabla f(\mathbf{w}_t)^\top (\mathbf{w} - \mathbf{w}_t) \right), \\ &= \frac{\partial}{\partial \mathbf{w}} \left(\frac{1}{2\eta_t} \left(\| \mathbf{w} \|^2 - 2\mathbf{w}_t^\top \mathbf{w} + \| \mathbf{w}_t \|^2 \right) + f(\mathbf{w}_t) + \nabla f(\mathbf{w}_t)^\top (\mathbf{w} - \mathbf{w}_t) \right), \\ &= \frac{1}{2\eta_t} \left(2\mathbf{w} - 2\mathbf{w}_t + 0 \right) + 0 + \nabla f(\mathbf{w}_t), \\ &= \frac{1}{\eta_t} \left(\mathbf{w} - \mathbf{w}_t \right) + \nabla f(\mathbf{w}_t). \end{split}$$

Setting the gradient to zero and rearranging give the desired output.

$$\mathbf{w}_{t+1} = \mathbf{w}_t - \eta_t \nabla f(\mathbf{w}_t).$$

Gradient Descent



$$\mathbf{w}_{t+1} = \mathbf{w}_t - \eta_t \nabla f(\mathbf{w}_t).$$



Gradient Descent



$$\mathbf{w}_{t+1} = \mathbf{w}_t - \eta_t \nabla f(\mathbf{w}_t).$$



Gradient Descent (GD)





f(w)

Gradient Descent(GD)





f(w)

Gradient Descent(GD)





f(w)

What went wrong here?

Provable Progress

DEPARTMENT OF ENGINEERING SCIENCE

Theorem 5.1. Let us assume f is convex and smooth with constant β then if we select $\eta \leq \frac{2}{\beta}$ we will have monotonic decrease in function value after each step.

Proof. From our assumption that f is convex and smooth with constant β then by the definition of the smoothness we have:

$$orall \mathbf{x}, \mathbf{y}, \ |f(\mathbf{x}) - f(\mathbf{y}) -
abla f(\mathbf{y})^{ op} (\mathbf{x} - \mathbf{y})| \leq rac{eta}{2} \|\mathbf{x} - \mathbf{y}\|^2.$$

Rearranging gives:

$$f(\mathbf{x}) \leq f(\mathbf{y}) + \nabla f(\mathbf{y})^{\top} (\mathbf{x} - \mathbf{y}) + \frac{\beta}{2} \|\mathbf{x} - \mathbf{y}\|^2.$$

Let $\mathbf{x} = \mathbf{w}_{t+1} = \mathbf{w}_t - \eta \nabla f(\mathbf{w}_t)$, $\mathbf{y} = \mathbf{w}_t$ and hence $\mathbf{x} - \mathbf{y} = -\eta \nabla f(\mathbf{w}_t)$. Plugging this in gives:

$$f(\mathbf{w}_{t+1}) \leq f(\mathbf{w}_t) + \eta \left(-\|\nabla f(\mathbf{w}_t)\|^2 + \frac{\beta\eta}{2} \|\nabla f(\mathbf{w}_t)\|^2 \right),$$

$$f(\mathbf{w}_{t+1}) \leq f(\mathbf{w}_t) - \eta \left(1 - \frac{\beta\eta}{2} \right) \|\nabla f(\mathbf{w}_t)\|^2.$$

Hence if $1 - \frac{\beta \eta}{2} \ge 0$ we will have a decrease. Rearranging this condition gives the desired result.



5.3 What if we don't know what the smoothness constant is?

Simply try a bunch of different values, and keep the one that works best. This process is know as cross validation or hyperparameter tuning.

5.4 What if we don't have a smooth function?

What if our function is only Lipschitz continuous and not smooth? well we can still prove GD asymptotically converges to the optimum for convex function but only with a decreasing step size such as $\eta_t = \frac{1}{\sqrt{t}}$.







5.5 Line Search Methods

Line search methods contain two key components:

- 1. A method for proposing points, typically backtracking procedure.
- 2. A condition to determine where a point is accepted.



Backtracking Line Search At time t a point we first select a point $\mathbf{w}_t - \eta_{t,0} \nabla f(\mathbf{w}_t)$ check the point meets the acceptance conditions, if not $\eta_{t,k+1} = \gamma \cdot \eta_{t,k}$ where $\gamma \in (0,1)$.

Armijo-Goldstein Conditions

$$f(\mathbf{w}_t - \eta_{t_k} \nabla f(\mathbf{w}_t)) \le f(\mathbf{w}_t) - c\eta_{t_k} \|\nabla f(\mathbf{w}_t)\|^2, \quad 0 < c.$$

Wolfe Conditions

$$f(\mathbf{w}_t - \eta_{t_k} \nabla f(\mathbf{w}_t)) - f(\mathbf{w}_t) \le c \eta_{t_k} \|\nabla f(\mathbf{w}_t)\|^2,$$

$$\nabla f(\mathbf{w}_t - \eta_{t_k} \nabla f(\mathbf{w}_t))^\top \nabla f(\mathbf{w}_t) \ge c \|\nabla f(\mathbf{w}_t)\|^2 \quad 0 < c_1 < c_2 < 1$$

Question: Why might one prefer the Armijo-Goldstein Conditions?











Algorithm 2 Back Tracking Line Search

Require: γ : learning rate **Require:** α : scaling factor **Require:** f: objective function **Require:** \mathbf{w}_0 : initial parameter vector **Require:** *c*: Armijo-Goldstein Hyper-parameter **Require:** k_{max} : max number of inner iterations for $t = 1, \dots, T$ do $\eta_{t_0} = \gamma$ for k in range k_{max} do if $f(\mathbf{w}_t - \eta_{t_k} \nabla f(\mathbf{w}_t)) \leq f(\mathbf{w}_t) - c\eta_{t_k} \|\nabla f(\mathbf{w}_t)\|^2$ then $\mathbf{w}_{t+1} = \mathbf{w}_t - \eta_{t_k} \nabla f(\mathbf{w}_t)$ Break else $\eta_{t_k} = \alpha \cdot \eta_{t_{k-1}}$ end if end for end for return \mathbf{w}_T

Secord Order Methods



Second Order Methods

$$\mathbf{w}_{t+1} = \operatorname*{argmin}_{\mathbf{w}\in\Omega} \left\{ f(\mathbf{w}_t) + \nabla f(\mathbf{w}_t)^\top (\mathbf{w} - \mathbf{w}_t) + \frac{1}{2} (\mathbf{w} - \mathbf{w}_t)^T \boldsymbol{H}_f(\mathbf{w} - \mathbf{w}_t) \right\}.$$

$$\begin{split} &\frac{\partial}{\partial \mathbf{w}} \left(f(\mathbf{w}_t) + \nabla f(\mathbf{w}_t)^\top (\mathbf{w} - \mathbf{w}_t) + \frac{1}{2} (\mathbf{w} - \mathbf{w}_t)^T \boldsymbol{H}_f(\mathbf{w} - \mathbf{w}_t) \right), \\ &= \frac{\partial}{\partial \mathbf{w}} \left(f(\mathbf{w}_t) + \nabla f(\mathbf{w}_t)^\top (\mathbf{w} - \mathbf{w}_t) + \frac{1}{2} \left(\mathbf{w}^T \boldsymbol{H}_f \mathbf{w} - 2\mathbf{w}^\top \boldsymbol{H}_f \mathbf{w}_t + \mathbf{w}_t^T \boldsymbol{H}_f \mathbf{w}_t \right) \right), \\ &= \left(0 + \nabla f(\mathbf{w}_t) + \frac{1}{2} \left(2\boldsymbol{H}_f \mathbf{w} - 2\boldsymbol{H}_f \mathbf{w}_t + 0 \right) \right) \\ &= \nabla f(\mathbf{w}_t) + \boldsymbol{H}_f \mathbf{w} - \boldsymbol{H}_f \mathbf{w}_t \end{split}$$

Secord Order Methods



Setting the gradient equal to zero:

$$0 = \nabla f(\mathbf{w}_t) + \boldsymbol{H}_f \mathbf{w} - \boldsymbol{H}_f \mathbf{w}_t,$$
$$-\nabla f(\mathbf{w}_t) = \boldsymbol{H}_f \mathbf{w} - \boldsymbol{H}_f \mathbf{w}_t,$$
$$-\boldsymbol{H}_f^{-1} \nabla f(\mathbf{w}_t) = \mathbf{w} - \mathbf{w}_t,$$
$$\mathbf{w} = \mathbf{w}_t - \boldsymbol{H}_f^{-1} \nabla f(\mathbf{w}_t).$$

What do you notice about this update?

Stochastic Gradient Descent (SGD)

ę

$$f(\mathbf{w}) \triangleq \frac{1}{N} \sum_{z=1}^{N} \ell_z(\mathbf{w}) \approx \mathbb{E}_{z \in \mathcal{Z}}[\ell_z(\mathbf{w})]$$

Gradient Descent

$$\mathbf{w}_{t+1} = \operatorname*{argmin}_{\mathbf{w}\in\Omega} \left\{ \frac{1}{2\eta_t} \|\mathbf{w} - \mathbf{w}_t\|^2 + f(\mathbf{w}_t) + \nabla f(\mathbf{w}_t)^\top (\mathbf{w} - \mathbf{w}_t) \right\}.$$

$$\mathbf{w}_{t+1} = \mathbf{w}_t - \eta_t \nabla f(\mathbf{w}_t).$$

Stochastic Gradient Descent

$$\mathbf{w}_{t+1} = \operatorname*{argmin}_{\mathbf{w}\in\Omega} \left\{ \frac{1}{2\eta_t} \|\mathbf{w} - \mathbf{w}_t\|^2 + \ell_{z_t}(\mathbf{w}_t) + \nabla \ell_{z_t}(\mathbf{w}_t)^\top (\mathbf{w} - \mathbf{w}_t) \right\}.$$

 $\mathbf{w}_{t+1} = \mathbf{w}_t - \eta_t \nabla \ell_{z_t}(\mathbf{w}_t).$















https://production-media.paperswithcode.com/methods/Screen_Shot_2020-05-28_at_3.25.40_PM_Y687HvA.png

SGD with momentum

$$\begin{split} \mathbf{m}_0 &= 0, \\ \mathbf{m}_{t+1} &= \mu \mathbf{m}_t - \eta_t \nabla \ell_{z_t}(\mathbf{w}_t), \\ \mathbf{w}_{t+1} &= \mathbf{w}_t + \mathbf{m}_t. \end{split}$$

Nestarov Momentum

$$\begin{split} \mathbf{m}_0 &= 0, \\ \mathbf{m}_{t+1} &= \mu \mathbf{m}_t - \eta_t \nabla \ell_{z_t}(\mathbf{w}_t), \\ \mathbf{w}_{t+1} &= \mathbf{w}_t - \eta_t \nabla \ell_{z_t}(\mathbf{w}_t) + \mu \mathbf{m}_t. \end{split}$$





SGD with momentum

$$\begin{split} \mathbf{m}_0 &= 0, \\ \mathbf{m}_{t+1} &= \mu \mathbf{m}_t - \eta_t \nabla \ell_{z_t}(\mathbf{w}_t), \\ \mathbf{w}_{t+1} &= \mathbf{w}_t + \mathbf{m}_t. \end{split}$$

What happens to the update if the gradient is constant?

What is causing this zig-zagging?



• Poorly scaled dimension – high condition number

What is causing this zig-zagging?



Dimension with very different scales:

RMSprop

$$oldsymbol{v}_k \leftarrow 0.9 oldsymbol{v}_{k-1} + 0.1
abla \ell_{z_t} (\mathbf{w}_t)^2 \ \mathbf{w}_{k+1} \leftarrow \mathbf{w}_k - rac{\eta}{\sqrt{oldsymbol{v}_k + \epsilon}} \ell_{z_t} (\mathbf{w}_t)$$



Adam



\mathbf{Adam}

$$egin{aligned} \hat{m{m}}_k &= rac{m{m}_k}{1-eta_1^t}, \ m{m}_k \leftarrow eta_1m{m}_{k-1} + (1-eta_1)
abla \ell_{z_t}(m{w}_t) \ \hat{m{v}}_k &= rac{m{v}_k}{1-eta_2^t}, \ m{v}_k \leftarrow eta_2m{v}_{k-1} + (1-eta_2)
abla \ell_{z_t}(m{w}_t)^2 \ m{w}_{k+1} \leftarrow m{w}_k - rac{\eta}{\sqrt{\hat{m{v}}_k} + \epsilon} \hat{m{m}}_k \end{aligned}$$

Gradient Free Optimisation Greening



- What if you don't have access to the gradient?
- One must use optimisers that rely on function values
- Many gradient free optimisation algorithms exist
- In this class we use the "3 point" method as an illustrative example

Gradient Free Optimisation



UNIVERSITY OF

Algorithm 1 Non-Stochastic Three Points Method

Require: η : learning rate **Require:** *f*: objective function **Require:** \mathcal{D} : distribution over trial direction that spans \mathbb{R}^d for example $P(\mathcal{D} = \mathbf{e}_i) = 1/d$ **Require:** \mathbf{w}_0 : initial parameter vector for $t = 1, \dots, T$ do $m{p}_{t} \sim \mathcal{D}$ $\mathcal{W}_t \triangleq \{\mathbf{w}_t - \eta \boldsymbol{p}_t, \mathbf{w}_t, \mathbf{w}_t + \eta \boldsymbol{p}_t\}$ $\mathbf{w}_{t+1} \leftarrow \operatorname{argmin}_{\mathbf{w} \in \mathcal{W}_t} f(\mathcal{W}_t)$ {try all three values and pick the best} end for return \mathbf{w}_T

What do you think might be the disadvantages of this approach?



Time for some exercises