# Optimization

Let $f: \mathbb{R}^d \to \mathbb{R}$ be a scalar-valued function. ($f(x_1,..,x_d$

**Gradient**:

$$\nabla_x f(x) = \begin{bmatrix} \frac{\partial f}{\partial x_1} \\ \frac{\partial f}{\partial x_2} \\ \vdots \\ \frac{\partial f}{\partial x_d} \end{bmatrix}$$

$d \times 1$

**Hessian**:

$$\nabla_x^2 f(x) = \begin{bmatrix} \frac{\partial^2 f}{\partial x_1^2} & \frac{\partial^2 f}{\partial x_1 \partial x_2} & \cdots & \frac{\partial^2 f}{\partial x_1 \partial x_d} \\ \frac{\partial^2 f}{\partial x_2 \partial x_1} & \frac{\partial^2 f}{\partial x_2^2} & \cdots & \frac{\partial^2 f}{\partial x_2 \partial x_d} \\ \frac{\partial^2 f}{\partial x_d \partial x_1} & & \cdots & \frac{\partial^2 f}{\partial x_d^2} \end{bmatrix}$$

$d \times d$

⊛ In the case of

$f: \mathbb{R}^d \to \mathbb{R}^m$ being a vector-valued function, i.e.:

$$f(x) = \left( f_1(x_1,...,x_d), f_2(x_1,...,x_d), ---, f_m(x_1,x_2,...,x_d) \right)$$

**Jacobian**:

$$\nabla_x f(x) = \begin{bmatrix} \frac{\partial f_1}{\partial x_1} & \cdots & \frac{\partial f_m}{\partial x_1} \\ \vdots & & \\ \frac{\partial f_1}{\partial x_d} & \cdots & \frac{\partial f_m}{\partial x_d} \end{bmatrix}$$

$d \times m$

---

**Setup**: Given a model with $\overset{\text{parameters}}{\nabla} \vartheta \in \mathbb{R}^d$, i.e. $\vartheta = (\vartheta_1, \vartheta_2, ..., \vartheta_d)$

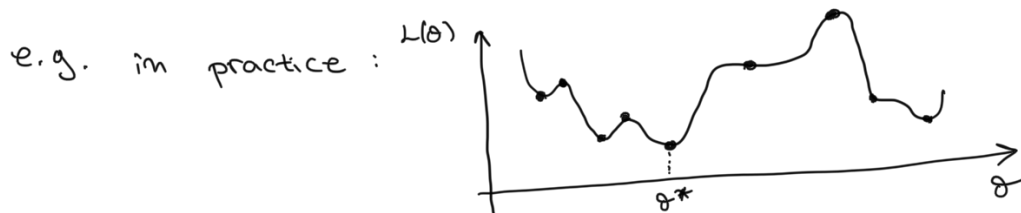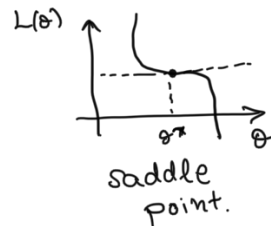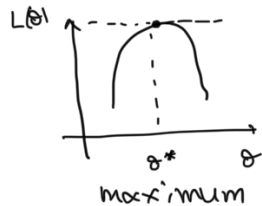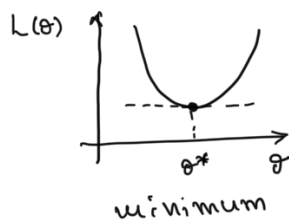and a loss/likelihood $L(\vartheta)$, then our goal is to identify a set of parameters $\vartheta^*$ such that:

$$\vartheta^* = \arg\min_\vartheta L(\vartheta)$$

We need to identify the critical points for which :

$$\nabla_\theta L(\theta) = 0$$

$L : \mathbb{R}^d \to \mathbb{R}$

This condition is true for : (i) minima, (ii) maxima, (iii) saddle poin

minimum
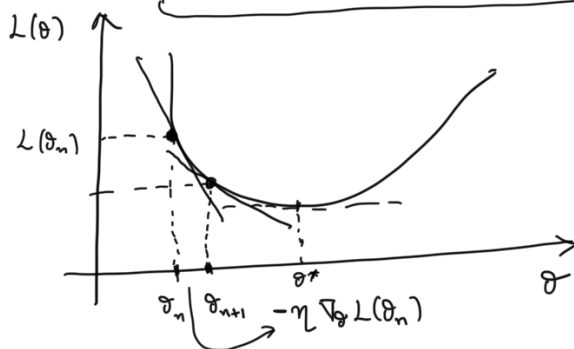
maximum

saddle point.

e.g. in practice :

## Gradient descent :

We want to minimize $L(\theta)$, i.e. $\theta^* = \arg\min_\theta L(\theta)$

Pick an initial guess $\theta_0$ and use update rule :

Update Rule :

$$\theta_{n+1} = \theta_n - \eta \nabla_\theta L(\theta_n)$$

$d \times 1 \quad d \times 1 \quad 1 \times 1 \quad d \times 1$

discretize version $\longrightarrow$

$\xleftarrow{\text{Forward Euler}}$

$$\frac{d\theta}{dt} = -\nabla_\theta L(\theta)$$

gradient flow system

- This is a first-order method as it relies on a linear approximation of $L(\theta)$ around some point $\theta_n$.

- It is guaranteed to converge to a critical point (e.g. local min / max) of $L(\theta)$, assuming that $\eta$ is properly chose

$\eta$ : step-size / learning rate (user need to adjust)

## Newton's algorithm :

Let's use a Taylor expansion of $L(\theta)$ around $\theta_n$ :

$$\mathcal{L}(\vartheta) \simeq \underset{1\times 1}{\mathcal{L}(\vartheta_n)} + \underset{1\times d}{g_n^T} \underset{d\times 1}{(\vartheta - \vartheta_n)} + \frac{1}{2} \underset{1\times d}{(\vartheta - \vartheta_n)^T} \underset{d\times d}{H_n} \underset{d\times 1}{(\vartheta - \vartheta_n)} \quad \Bigg| \quad g_n := \nabla_\vartheta \mathcal{L}(\vartheta_n$$

$$H_n := \nabla_\vartheta^2 \mathcal{L}(\vartheta$$

$$\simeq \mathcal{L}(\vartheta_n) + g_n^T (\vartheta - \vartheta_n) + \frac{1}{2}\left[ \vartheta^T H_n \vartheta - 2\vartheta^T H_n \vartheta_n + \vartheta_n^T H_n \vartheta_n \right]$$

Let us now find critical points:

$$\nabla_\vartheta \mathcal{L}(\vartheta) = 0 \implies 0 + g_n^T + H_n \vartheta - H_n \vartheta_n = 0$$

$$\implies \underset{\text{critical point}}{\underline{\vartheta}} = \vartheta_n - H_n^{-1} g_n^T$$

critical point of the $\underbrace{\text{Taylor approximation}}_{\text{quadratic}}$ to $\mathcal{L}(\vartheta)$

> Newton update rule: $\vartheta_{n+1} = \vartheta_n - H_n^{-1} g_n^T$



- No need for tuning parameters, but extra cost for computing and inverting the Hessian.

- $2^{nd}$-order algorithm, faster convergence vs gradient descent

- utilizes the underlying geometry by exploiting curvature information

### Example: Linear regression

Recall the loss function for MLE estimation:

$$\mathcal{L}(\vartheta) := - \log p(y \mid \underline{w}^T x, \underline{\sigma^2}), \quad \vartheta := \underset{d\times 1}{\{w}, \underset{1\times 1}{\sigma^2\}}$$

Gradient: $\nabla_w \mathcal{L}(w) = -X^T y + X^T X w$

Hessian: $\nabla_w^2 \mathcal{L}(w) = X^T X$

Gradient descent: $W_{n+1} = W_n - \eta \left[ -X^T y + X^T X w_n \right]$

Newton: $W_{n+1} = W_n - (X^T X)^{-1} \left[ X^T X w - X^T y \right]$

⊛ **Remark** : Similar updates can be formulated for $\sigma^2$ !

## Limitations :

- Gradient descent converges slowly. Choosing $\eta$ is an art.

- ⊙ Scalability to $\underbrace{big \ data.}_{\substack{evaluating \ the \ loss \\ over \ a \ big \ data-set \\ is \ expensive.}}$    $\vartheta_{n+1} = \vartheta_n - \eta \nabla_\vartheta \underbrace{L(\vartheta_n)}_{}$

  typically: $\propto \sum_{i=1}^{n} \underbrace{(y_i - w_i^T}$

- Exact Hessians are often very hard to compute $\xrightarrow{appr}$ Quasi–Newton methods.

---

## Stochastic gradient descent (SGD)

In many ML applications the loss functions factorize across data points, i.e. they can be written as a summation over individual data points :    $L(\vartheta) \propto \sum_{i=1}^{n} L_i(\vartheta)$    (e.g. linear regression

(this follows from assuming an $\overrightarrow{i.i.d}$ likelihood).

In this case, a standard "full batch" gradient descent approach would take the form :    $\vartheta_{n+1} = \vartheta_n - \eta \underbrace{\nabla_\vartheta L(\vartheta_n)}_{\substack{true \\ gradient}} = \vartheta_n - \eta \sum_{n=1}^{n} \nabla_\vartheta L_i(\vartheta_n$

In **stochastic** gradient descent,
the true gradient is approximated using a single example :

$$\vartheta_{n+1} = \vartheta_n - \eta \underbrace{\nabla_\vartheta L_i(\vartheta_n)}_{\substack{very \ noisy \\ approx. \ to \\ the \ true \ gradient.}} \ , \quad \begin{array}{l} i \ \text{is chosen at rando} \\ \text{at each iteration of} \\ \text{SGD algorithm} \end{array}$$

A compromise between these two extremes is to approximate the true gradient over a "mini-batch" of data :

$$\vartheta \quad = \vartheta - \eta \frac{1}{\sum}^{m} \nabla_\vartheta L(\vartheta_n) \qquad m \ll n$$

$$u_{n+1} = u_n \qquad \frac{1}{m} \sum_{i=1}^{m} \cdots \qquad )$$

At every iteration we randomly pick a mini-batch of data by sub-sampling our full data-set.

- A complete looping cycle over the entire data-set is called an "epoch".

## Modern variants of SGD

- **SGD with momentum :**

$u_n$ : momentum variable

$$\begin{cases} u_{n+1} = \gamma u_n + \eta \nabla_\theta L(\theta_n) \quad , \quad u_0 = 0 \\ \theta_{n+1} = \theta_n - u_{n+1} \end{cases} \quad , \quad \begin{cases} \eta: \text{ learning rate} \\ \gamma \begin{cases} 0 = SGD \\ 0.9 \text{ is a typical} \\ \text{value in prac.} \end{cases} \end{cases}$$

- **Nesterov accelerated gradient method (NAG) :**

$$\begin{cases} u_{n+1} = \gamma u_n + \eta \nabla_\theta L(\theta_n - \gamma u_n) \\ \theta_{n+1} = \theta_n - u^{n+1} \end{cases}$$

## Adaptive learning rate approaches :

**1.) RMS prop :**  $E[g^2]_n :=$ average of the square gradients.

$$\nwarrow E\left[(\nabla_\theta L(\theta_n))^2\right]_{d\times 1}$$

running avg of the sq. gradients

$d\times 1$

$$\begin{cases} E[g^2]_{n+1} = \gamma E[g^2]_n + (1-\gamma) g_n^2 \quad , \quad \text{typically} \begin{cases} \gamma = 0.9 \\ \eta = 10^{-3} \end{cases} \\ \theta_{n+1} = \theta_n - \dfrac{\eta}{\sqrt{E[g^2]_{n+1}} + \varepsilon} g_n \\ \qquad\qquad\qquad \curvearrowright 10^{-8} \end{cases}$$

**2.) Adam (adaptive moment estimation) :**

$$\begin{cases} \rightarrow m_{n+1} = \beta_1 m_n + (1-\beta_1) g_n \quad \rightarrow \quad \text{running average estimator of the expected gradie} \\ \rightarrow v_{n+1} = \beta_2 v_n + (1-\beta_2) g_n^2 \quad \rightarrow \quad -\text{"}- \\ \qquad\qquad\qquad\qquad\qquad\qquad\qquad\qquad \text{of the variance of the } g \end{cases}$$

$m, v$ are typically initialized to be zero

..., ... are typically ... ... ... to be zero, hence the estimates are biased towards zero. To counteract this bias, we can consider a correction:

$$\hat{m}_n = \frac{m_n}{1 - \beta_1^n} \quad , \quad \hat{v}_n = \frac{v_n}{1 - \beta_2^n}$$

Adam update rule $\Big\}$

$$\theta_{n+1} = \theta_n - \underbrace{\frac{\eta}{\sqrt{\hat{v}_{n+1}} + \varepsilon}}_{d \times 1} \hat{m}_{n+1}$$