#### Gradientenverfahren

Wie kann man Minima einer differenzierbaren Abbildung  $f: \mathbb{R}^n \to \mathbb{R}$  finden?

#### Gradientenverfahren

- An jedem Punkt  $x_k \in \mathbb{R}^n$  zeigt der negative Gradient  $d_k := -\nabla f(x_k)$  in die steilste Abstiegsrichtung.
- Für hinreichend kleines  $\alpha_k$  folgt mit Satz über die lokale Linearisierung:

$$f(x_{k+1}) = f(x_k + \alpha_k d_k) = f(x_k) + \alpha_k df(x_k) d_k + R(\alpha_k dk)$$

- Setze  $x_{k+1} = x_k + \alpha_k d_k$
- Es gilt  $f(x_{k+1}) \le f(x_k)$ , falls  $\nabla f(x_k) \ne 0$
- Falls die folge  $f(x_k)$  beschränkt ist, so ist dieser Fixpunkt  $x^*$  ein Minimum, da  $\nabla f(x^*) = 0$  gelten muss.

# Angewandte Mathematik Gradientenverfahren

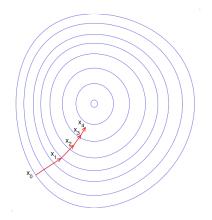


Figure: Quelle: Wikipedia

### Angewandte Mathematik

#### Höhenlinien

Sei  $f: \mathbb{R}^n \to \mathbb{R}$  eine differenzierbare Funktion. Eine Kurve  $\gamma: I \to \mathbb{R}^n$ , auf der f konstant ist, also  $f(\gamma(t)) = c$  für ein festes  $c \in \mathbb{R}$  gilt, heißt Höhenlinie.

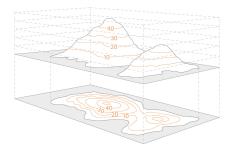


Figure: Quelle:

https://getoutside.ordnancesurvey.co.uk/guides/understanding-map-contour-lines-for-beginners/

### Angewandte Mathematik

#### Höhenlinien

Der Gradient steht senkrecht auf Höhenlinien. Dies bedeutet, dass

$$\langle \nabla f(\gamma(t)), \gamma'(t) \rangle = 0$$

gilt.

#### **Beweis**

Aus  $f(\gamma(t)) = c$  folgt  $\frac{d}{dt}f(\gamma(t)) = 0$ . Mit der Kettenregel folgt  $\frac{d}{dt}f(\gamma(t)) = df(\gamma(t)) \cdot \gamma'(t) = 0$  und damit  $\langle \nabla f(\gamma(t)), \gamma'(t) \rangle = 0$ .

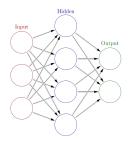
#### Backpropagation

Das Gradientenverfahren angewendet auf eine Lossfunktion eines neuronalen Netzes wird als Backpropagation bezeichnet. Gegeben ist ein neuronales Netz  $f: \Omega \times \mathbb{R}^n \to \mathbb{R}^m$ , und ein Datensatz  $D:=\{(x_i,y_i)\}$  mit  $x_i\in \mathbb{R}^n,y_i\in \mathbb{R}^m$ . Finde Gewichte Omega, so dass Lossfunktion

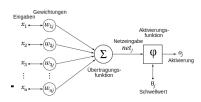
$$L_D:\Omega\subset\mathbb{R}^n\to\mathbb{R}$$

minimal wird. Zum Beispiel

$$L_D(\omega) := \sum_{(x_i, y_i) \in D} (f(\omega, x_i) - y_i)^2$$



Figure



**Figure** 

#### Backpropagation

• Initialisiere k := 0 und zufällige Gewichte  $w_0$ .

- Initialisiere k := 0 und zufällige Gewichte  $w_0$ .
- ullet Initialisiere Genauigkeit  $\epsilon>0$

- Initialisiere k := 0 und zufällige Gewichte  $w_0$ .
- ullet Initialisiere Genauigkeit  $\epsilon>0$
- While  $||\nabla L_D(\omega)|| > \epsilon$

- Initialisiere k := 0 und zufällige Gewichte  $w_0$ .
- Initialisiere Genauigkeit  $\epsilon > 0$
- While  $||\nabla L_D(\omega)|| > \epsilon$
- Bestimme  $\alpha_k$  mit  $L_D(\omega_k + \alpha d_k) = L_D(\omega_k) + \alpha_k dL_D(\omega_k) d_k + R(\alpha_k dk)$

- Initialisiere k := 0 und zufällige Gewichte  $w_0$ .
- Initialisiere Genauigkeit  $\epsilon>0$
- While  $||\nabla L_D(\omega)|| > \epsilon$
- Bestimme  $\alpha_k$  mit  $L_D(\omega_k + \alpha d_k) = L_D(\omega_k) + \alpha_k dL_D(\omega_k) d_k + R(\alpha_k dk)$
- Setze  $\omega_{k+1} := \omega_k + \alpha_k d_k$ .

- Initialisiere k := 0 und zufällige Gewichte  $w_0$ .
- Initialisiere Genauigkeit  $\epsilon>0$
- While  $||\nabla L_D(\omega)|| > \epsilon$
- Bestimme  $\alpha_k$  mit  $L_D(\omega_k + \alpha d_k) = L_D(\omega_k) + \alpha_k dL_D(\omega_k) d_k + R(\alpha_k dk)$
- Setze  $\omega_{k+1} := \omega_k + \alpha_k d_k$ .
- $k \leftarrow k + 1$

#### Mini Batch

• Datensatz D sehr groß (Big Data)

#### Mini Batch

- Datensatz D sehr groß (Big Data)
- Berechnung des Gradienten der Lossfunktion entsprechend aufwendig.

#### Mini Batch

- Datensatz D sehr groß (Big Data)
- Berechnung des Gradienten der Lossfunktion entsprechend aufwendig.
- Wende Backpropagation auf Teilräume  $D' \subset D$  an (Minibatch).

#### Mini Batch

- Datensatz D sehr groß (Big Data)
- Berechnung des Gradienten der Lossfunktion entsprechend aufwendig.
- Wende Backpropagation auf Teilräume  $D' \subset D$  an (Minibatch).
- #D' = 1 stochastischer Gradientenabstieg.

## Angewandte Mathematik

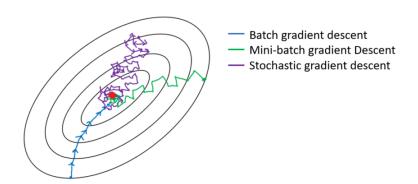


Figure: Quelle: https://towardsdatascience.com/batch-mini-batch-stochastic-gradient-descent-7a62ecba642a



#### Backpropagation

• Initialisiere k := 0 und zufällige Gewichte  $w_0$ .

- Initialisiere k := 0 und zufällige Gewichte  $w_0$ .
- ullet Initialisiere Genauigkeit  $\epsilon>0$

- Initialisiere k := 0 und zufällige Gewichte  $w_0$ .
- Initialisiere Genauigkeit  $\epsilon > 0$
- Wähle Teilmenge  $D_0' \subset D$

- Initialisiere k := 0 und zufällige Gewichte  $w_0$ .
- Initialisiere Genauigkeit  $\epsilon > 0$
- Wähle Teilmenge  $D_0' \subset D$
- While  $||\nabla L_{D'_{k}}(\omega)|| > \epsilon$

- Initialisiere k := 0 und zufällige Gewichte  $w_0$ .
- Initialisiere Genauigkeit  $\epsilon > 0$
- Wähle Teilmenge  $D_0' \subset D$
- While  $||\nabla L_{D'_{k}}(\omega)|| > \epsilon$
- ullet Bestimme  $\alpha_k$  mit

$$L_{D'_k}(\omega_k + \alpha d_k) = L_{D'_k}(\omega_k) + \alpha_k dL_{D'_k}(\omega_k) d_k + R(\alpha_k dk)$$

- Initialisiere k := 0 und zufällige Gewichte  $w_0$ .
- Initialisiere Genauigkeit  $\epsilon > 0$
- Wähle Teilmenge  $D_0' \subset D$
- While  $||\nabla L_{D'_{\iota}}(\omega)|| > \epsilon$
- Bestimme  $\alpha_k$  mit  $L_{D'_k}(\omega_k + \alpha d_k) = L_{D'_k}(\omega_k) + \alpha_k dL_{D'_k}(\omega_k) d_k + R(\alpha_k dk)$
- Setze  $\omega_{k+1} := \omega_k + \alpha_k d_k$ .

- Initialisiere k := 0 und zufällige Gewichte  $w_0$ .
- Initialisiere Genauigkeit  $\epsilon > 0$
- Wähle Teilmenge  $D_0' \subset D$
- While  $||\nabla L_{D'_{\iota}}(\omega)|| > \epsilon$
- Bestimme  $\alpha_k$  mit  $L_{D'_k}(\omega_k + \alpha d_k) = L_{D'_k}(\omega_k) + \alpha_k dL_{D'_k}(\omega_k) d_k + R(\alpha_k dk)$
- Setze  $\omega_{k+1} := \omega_k + \alpha_k d_k$ .
- Wähle neue Teilmenge  $D'_{k+1} \subset D$ .

- Initialisiere k := 0 und zufällige Gewichte  $w_0$ .
- Initialisiere Genauigkeit  $\epsilon > 0$
- Wähle Teilmenge  $D_0' \subset D$
- While  $||\nabla L_{D'_{\iota}}(\omega)|| > \epsilon$
- Bestimme  $\alpha_k$  mit  $L_{D'_k}(\omega_k + \alpha d_k) = L_{D'_k}(\omega_k) + \alpha_k dL_{D'_k}(\omega_k) d_k + R(\alpha_k dk)$
- Setze  $\omega_{k+1} := \omega_k + \alpha_k d_k$ .
- Wähle neue Teilmenge  $D'_{k+1} \subset D$ .
- $k \leftarrow k + 1$



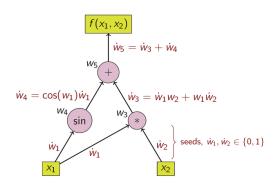


Figure: Quelle: Wikipedia

Automatisches Ableiten in Pytorch Automatisches Ableiten in JAX



#### Konditionszahl

Ist  $f:U\to\mathbb{R}^m$  differnzierbar, so lässt dich mit dem Mittelwertsatz die Konditionszahlen berechnen durch

$$\kappa = ||f'(x)||$$

$$\kappa_{rel} = \frac{||x||}{||f(x)||} ||f'(x)||$$

mit der Operatornorm  $||h|| := \sup_{||x||=1} ||h(x)||$ 

#### Konditionszahl Addition/Subtraction

Für  $f:(x,y):=x\pm y$  ist  $f'(x,y)=(1,\pm 1)$  und dem Betrag als Norm auf  $\mathbb R$  und der Norm ||(x,y)||:=|x|+|y| ist

$$\kappa=1$$
 
$$\kappa_{\it rel}=rac{|{\it a}|+|{\it b}|}{|{\it a}\pm{\it b}|}$$

Für die Addition ist  $\kappa_{rel}=1$  und damit gut konditioniert. Für die Subtraktion zweier fast gleich großer Zahlen ist |a-b|<<|a|+|b| und damit ist in diesem Fall  $\kappa_{rel}>>1$  und damit schlecht konditioniert.