

RUB

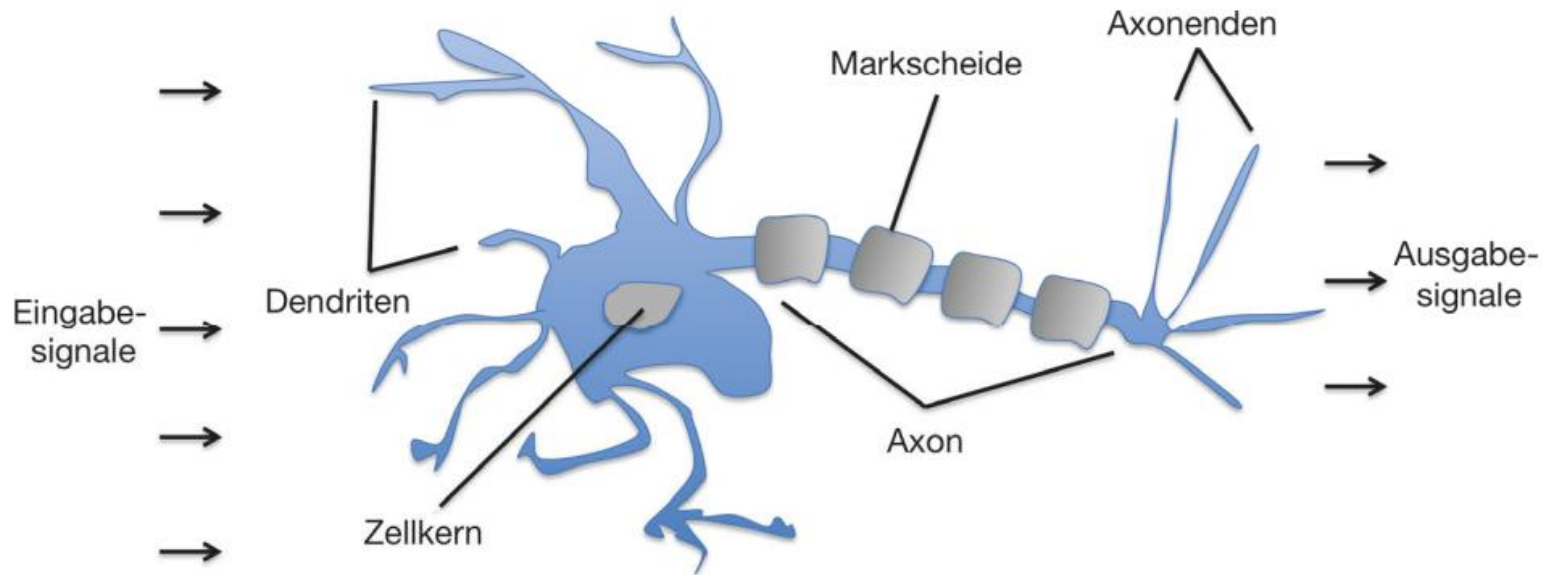
RUHR-UNIVERSITÄT BOCHUM

LERNALGORITHMEN

Training eines Klassifizierungsalgorithmus + Implementieren in Python

Agenda

- Biologischer Hintergrund
- Das **Perzeptron**
 - Grundbegriffe
 - Perzeptron in Python
- **Adaline** (lineare Neuronen)
 - **Gradientenabstiegsverfahren (GD)**
 - Adaline in Python mit GD
 - Problem mit GD
 - **Stochastisches** **Gradientenabstiegsverfahren (SGD)**



Dies ist der ungefähre Aufbau eines **Neurons** (also einer Zelle unseres Gehirns)

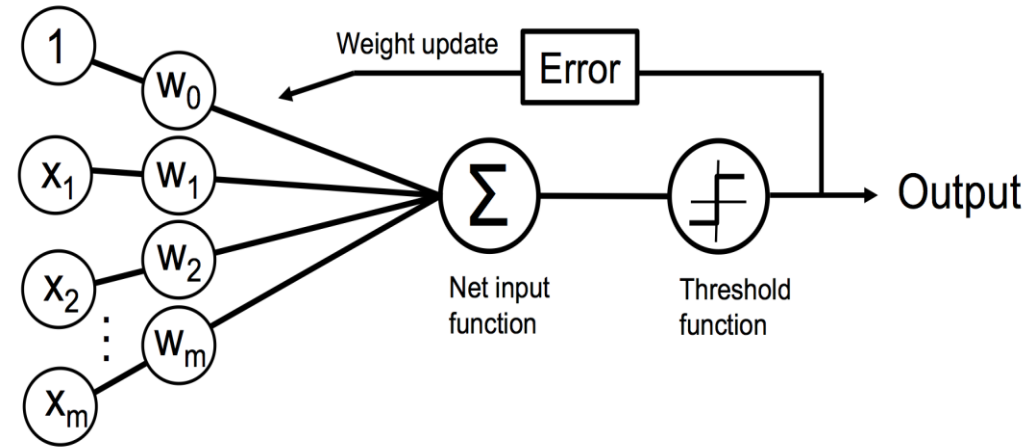
Neuron kann sortieren und selektieren

Wichtige Begriffe **Schwellenwert** und **Aktivieren**

Idee: Programmieren eines ähnlichen Modells, welches wie das Gehirn Sortierungen durchführen kann

Das Perzeptron

Das Perzeptron



Neuronales Netz mit nur einem Layer

Wir brauchen ein paar Definitionen:

- **Eingabevektor: x**
- **Gewichtungsvektor: w**
- **Nettoeingabe: z**
- **Aktivierungsfunktion: $\Phi(z)$**
- **Ausgabewert: \hat{y}**
- **Lernrate: η**

Grundbegriffe für das Perzeptron

Unser Perzeptron soll eine binäre Ausgabe haben

Wir brauchen **Eingabevektor** und **Gewichtungsvektor**:

- Eingabevektor: $\mathbf{x} = (x_1, x_2, \dots, x_n)$
- Gewichtungsvektor: $\mathbf{w} = (w_1, w_2, \dots, w_n)$

Damit berechnen wir die **Nettoeingabe**.

- Nettoeingabe: $\mathbf{z} = \mathbf{w}^T \mathbf{x} = w_1 * x_1 + w_2 * x_2 + \dots + w_n * x_n$

Grundbegriffe für das Perzeptron

Aktivierungsfunktion: $\Phi(z) = \begin{cases} 1, & \text{wenn } z \geq \theta \\ -1, & \text{andernfalls} \end{cases}$

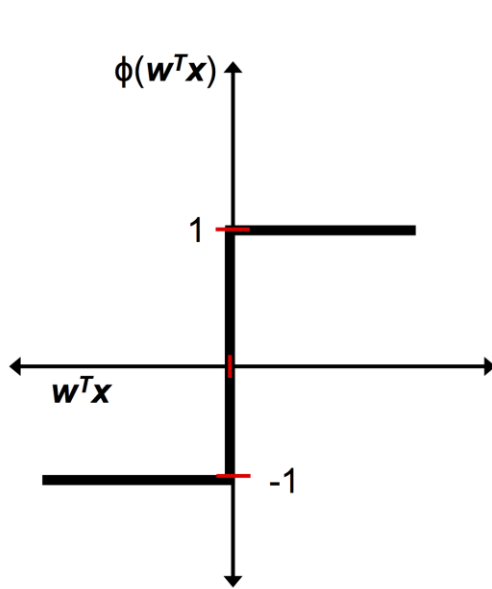
Nullgewichtung dadurch führen wir $w_0 = -\theta$ ein, sowie $x_0 = 1$

$$\Rightarrow z = w_0 * x_0 + w_1 * x_1 + \dots + w_n * x_n$$

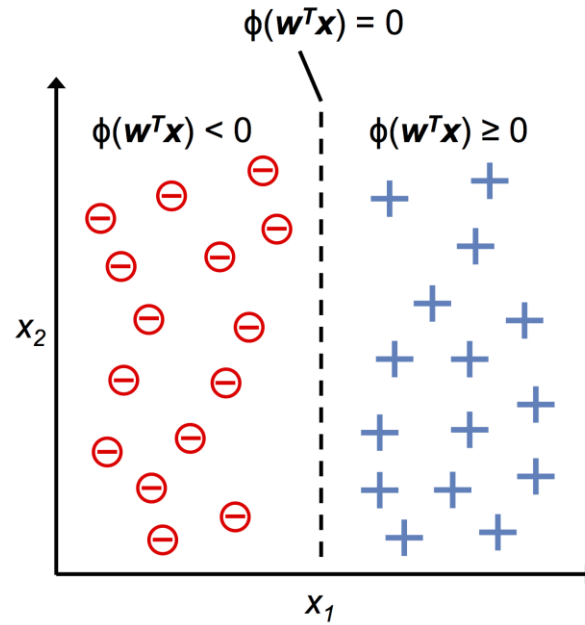
$$\Rightarrow \Phi(z) = \begin{cases} 1, & \text{wenn } z \geq 0 \\ -1, & \text{andernfalls} \end{cases}$$

Später brauchen wir noch den **Ausgabewert** $\Phi(z) = \hat{y}$,

sowie $y = \text{Realwert}$



Links:
Aktivierungsfunktion $\Phi(w^T x)$



Rechts:
Trennung von zwei Klassen, durch Anwendung von $\Phi(w^T x)$

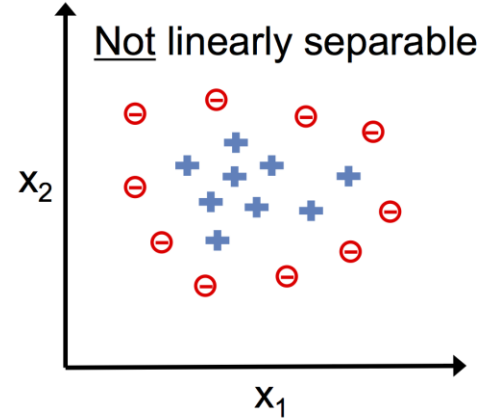
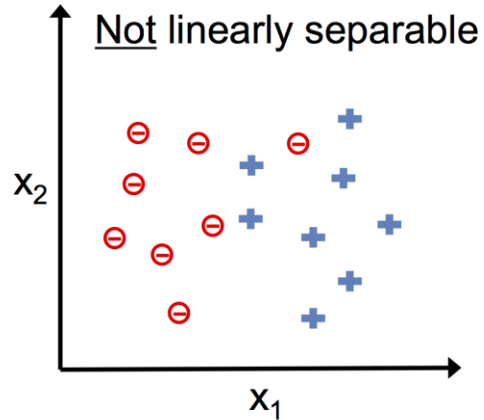
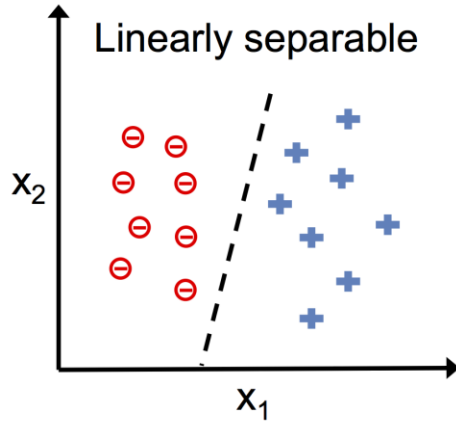
Training des Perzeptrons

1. Die Gewichtungen werden mit 0 oder kleinen Werten initialisiert (später mit 0)
2. Für alle Trainingsobjekten x^i passiert folgendes
 1. Berechnung des Ausgabewerts \hat{y}
 2. Aktualisierung der Gewichtungen
 1. $w_j := w_j + \Delta w_j$
 2. $\Delta w_j = \eta(y^{(i)} - \hat{y}^{(i)})x_j^{(i)}$

Hierbei ist es wichtig zu erkennen, dass $\hat{y}^{(i)} > y^{(i)} \Rightarrow w_j(\text{neu}) < w_j(\text{alt})$

Andersrum genauso.

Problem mit dem Perzeptron

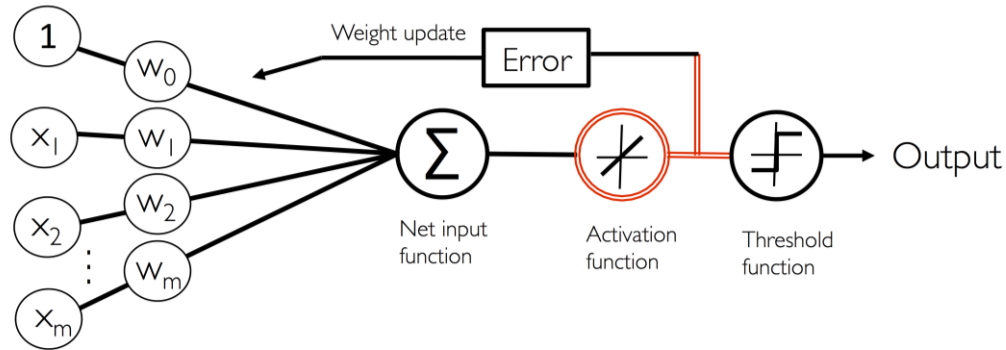
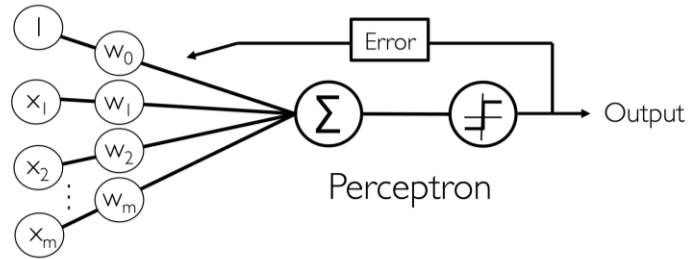


Das Perzeptron in Python

Adaline

Und das

Gradientenabstiegsverfahren

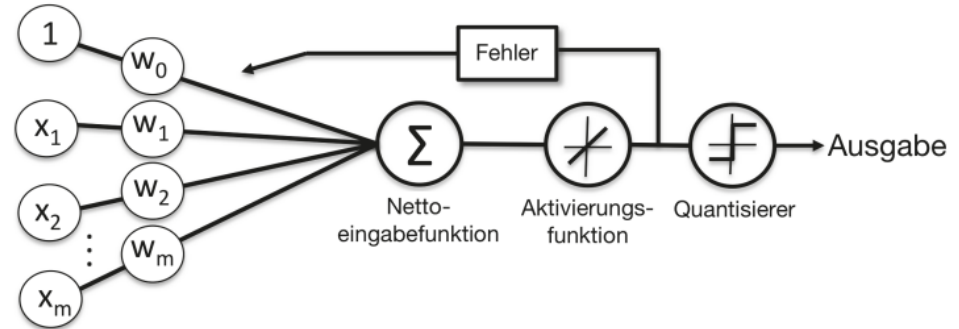


Adaptive Linear Neuron (Adaline)

Perzeptron im Vergleich mit Adaline

Adaline

- Aktivierungsfunktion:
lineare Funktion
- Aktivierungsfunktion beginnt als
Identitäts Abbildung
 $\Phi(w^t x) = w^t x$
- Gewichts Anpassung durch
Minimierung der Errorfunktion
- Quantisierer sorgt wieder für
Binären Output



Errorfunktion mit dem GD minimieren

Bestimmung einer Zielfunktion, diese wird oft als Errorfunktion (Straffunktion) bezeichnet.

Bei Adeline betrachten wir die Errorfunktion J:

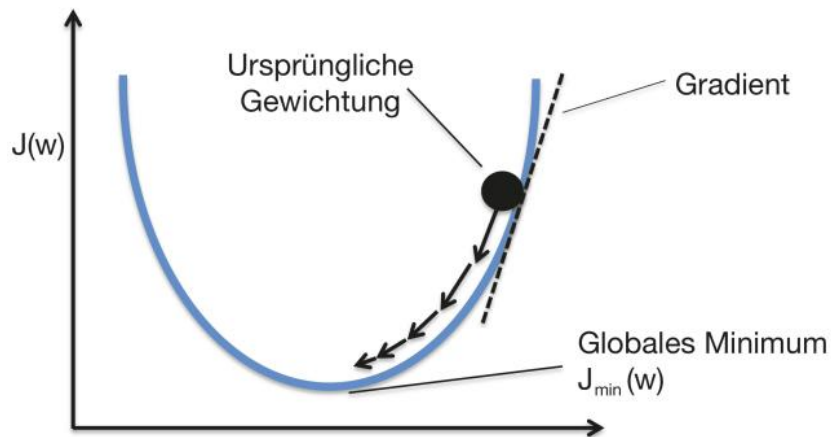
$$J(w) = \frac{1}{2} \sum_i \left(y^{(i)} - \Phi(z^{(i)}) \right)^2$$

Also die quadrierte Abweichung (MSE)

Diese Zielfunktion hat den praktischen Vorteil konvex zu sein

⇒ hat Minimum

Gradientenabstiegsverfahren



Wir aktualisieren nun die Gewichtungen so, dass wir uns einen Schritt vom Gradienten entfernen, also in Richtung Minimum gehen

Die Gewichtungen werden nun so aktualisiert:

$$w := w + \Delta w$$

$$\Delta w := -\eta \nabla J(w)$$

Um ∇J zu berechnen brauchen wir die partiellen Ableitungen der Errorfunktion

$$\frac{\partial J}{\partial w_j} = - \sum_i (y^{(i)} - \Phi(z^{(i)})) x_j^{(i)}$$

$$\Rightarrow \Delta w_j = -\eta \frac{\partial J}{\partial w_j} = \eta \sum_i (y^{(i)} - \Phi(z^{(i)})) x_j^{(i)}$$

Perzeptron Adeline Unterschiede

- Perzeptron hat ganzzahlige Ausgabe von $\Phi(z)$
- Perzeptron Gewichtaktualisierung nach jedem Objekt

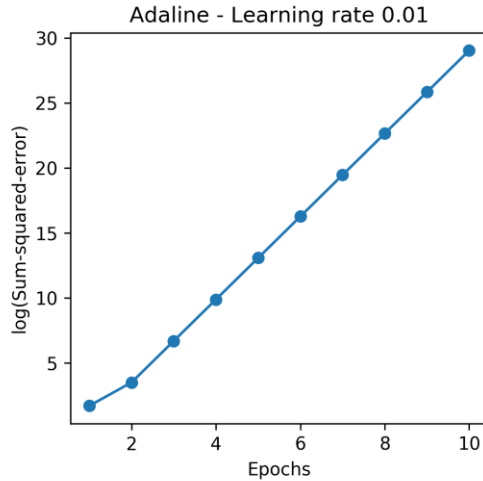
$$\begin{aligned}w_j &:= w_j + \Delta w_j \\w_0 &= \eta(y^{(i)} - \hat{y}^{(i)}) \\ \Delta w_{j \in [1,n]} &= \eta(y^{(i)} - \hat{y}^{(i)})x_j^{(i)}\end{aligned}$$

- Adeline hat eine reelle Zahl als Ausgabe von $\Phi(w^t x)$
- Adeline Gewichtaktualisierung beruht auf allen Objekten

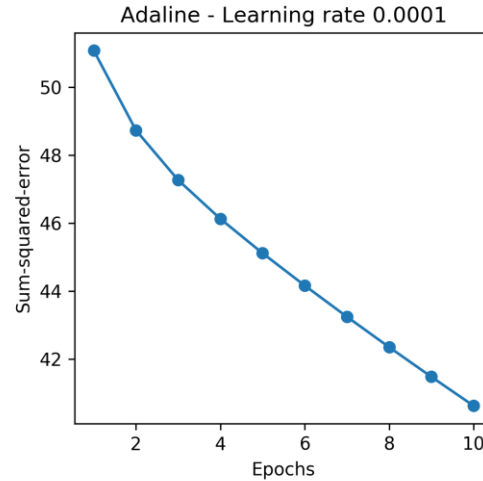
$$\begin{aligned}w_0 &= \eta * \sum errors, \text{ wobei } errors = (y - X^T * w) \\ w_{i \in [1,n]} &= \eta * X^T * (y^{(i)} - X^T * w^{(i)})\end{aligned}$$

- Wie bei Perzeptron Sammlung der Fehler um Konvergenz zu erkennen
 - Diesmal Cost genannt in `self.cost_` gespeichert

Gradientenabstiegsverfahren



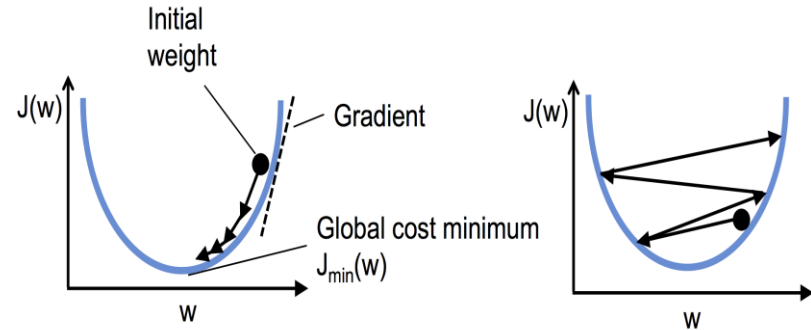
Links Adeline Abweichungen pro Epoche bei einer konstanten Lernrate von 0,01



Rechts Adeline Abweichungen pro Epoche bei einer konstanten Lernrate von 0,0001

Beide Lernraten bringen keine ersichtliche Konvergenz.

- Rechts Lernrate zu groß
- Links Lernrate zu klein



Adeline in Python mit Gradientenabstiegsverfahren

Problem beim Gradientenabstiegsverfahren

Wir haben nun gesehen wie das Gradientenabstiegsverfahren funktioniert

Bei großen Datenmengen \Rightarrow großer Rechenaufwand!

$$\Delta w_j = -\eta \frac{\partial J}{\partial w_j} = \eta \sum_i \left(y^{(i)} - \Phi(z^{(i)}) \right) x_j^{(i)}$$

Summe sehr groß für viele $i \Rightarrow$ eine Iteration dauert sehr lang

Und da kommt das stochastische Gradientenabstiegsverfahren zur Rettung

Large-scale Machine Learning

stochastisches

Gradientenabstiegsverfahren

Stochastisches Gradientenabstiegsverfahren

Anstatt Gewichtung mit allen Testobjekten

$$\Delta w_j = -\eta \sum_i \left(y^{(i)} - \Phi(z^{(i)}) \right) x^{(i)}$$

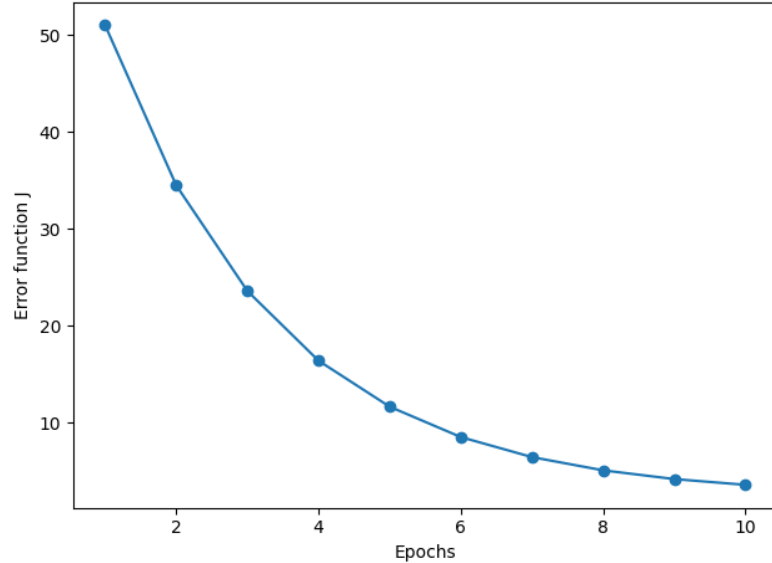
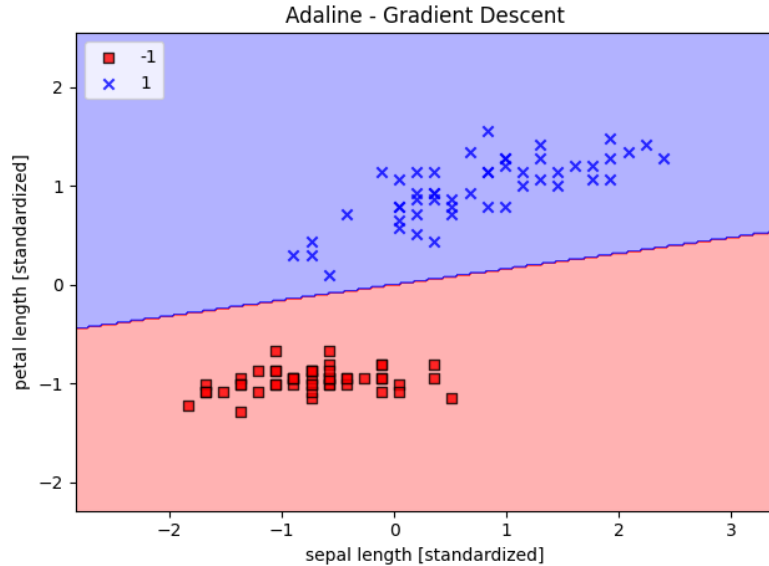
Benutzen für die Aktualisierung nur die Abweichung eines zufälligen Trainingsobjekt x^i

$$\eta \left(y^{(i)} - \Phi(z^{(i)}) \right) x^{(i)}$$

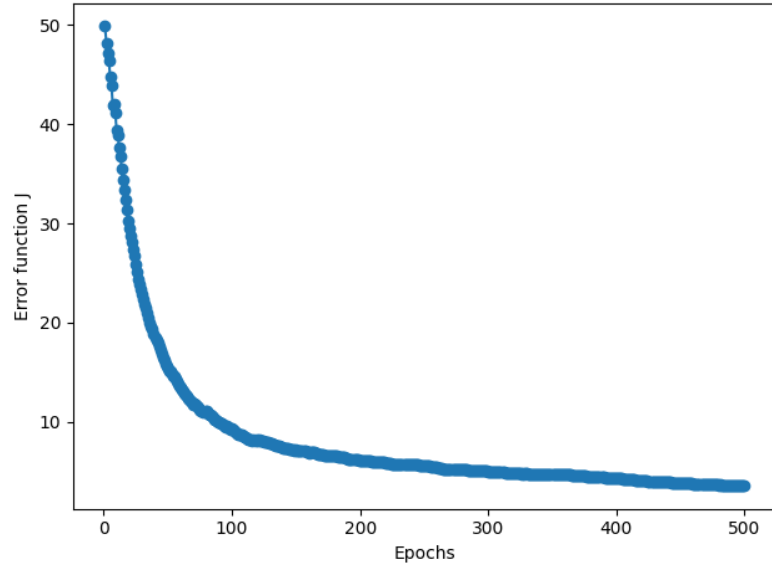
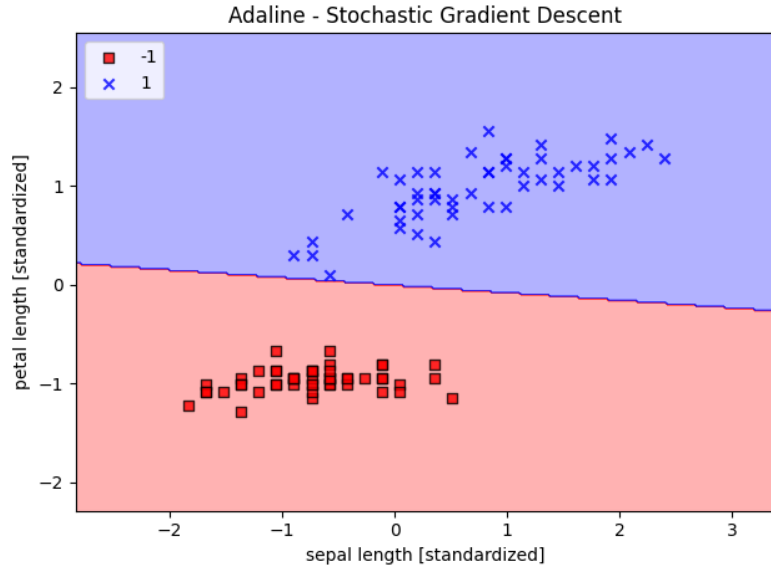
Deutlich ungenauer als GD, aber auch deutlich schnellere Iterationen

Es ist wichtig, dass die Trainingsobjekte, deren Aktualisierungen wir beobachten zufällig gewählt werden

Gradientenabstiegsverfahren



Stochastisches Gradientenabstiegsverfahren



Stochastisches Gradientenabstiegsverfahren

Um das stochastische Gradientenabstiegsverfahren genauer zu machen wird oft eine **adaptive Lernrate** benutzt

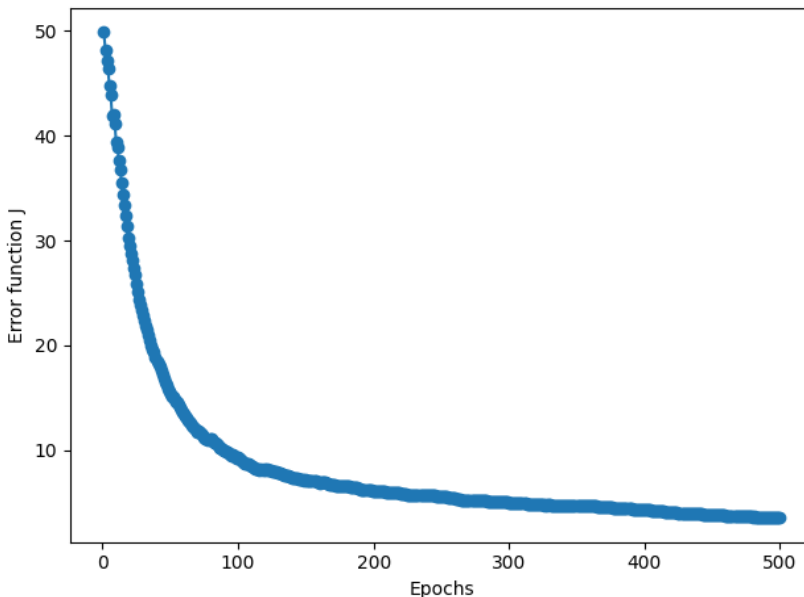
z.B. $\frac{c_1}{[Anzahl\ der\ Iterationen] + c_2}$ wobei c_1 und c_2 Konstanten sind

Ein Beispiel für eine adaptive Lernrate wäre auch **Adam** (Dieser ist eine ausgefeilte Version einer adaptiven Lernrate)

Dass GSD wird auch oft beim *online Learning* verwendet, wo es um große in Echtzeit eintreffende Datenmengen geht (Werbung).

Stochastisches Gradientenabstiegsverfahren

- Schnelle Iterationen sind gut bei großen Datenmengen
 - Daher beliebt bei Online Learning
 - Ungenauer als GD
- Mittelweg zwischen GD und SGD (*Mini-Batch-Learning*)
 - Hierbei wendet man das GD auf kleinere Teilmengen an



Errorfunktion bei stochastischem Gradientenabstiegsverfahren

Vielen Dank fürs Zuhören

Falls es noch Fragen gibt gerne
stellen!

Frohe Weinachten!



Quellen

BildQuellen:

<https://www.publicdomainpictures.net/de/view-image.php?image=28179&picture=frohe-weihnachten-4>

RaschkaSebastian Machine Learning mit Python von 2017

TextQuelle:

RaschkaSebastian Machine Learning mit Python von 2017