

V12 Klassifikation (Machine Learning 001)

- Entscheidungsbäume (decision tree): ID3-Algorithmus
- Support Vector-Maschinen (SVM)
- Ein Beispiel zum Einsatz von aktuellen Deep Learning-Methoden

Grobe Unterteilung:

- Methoden für unsupervised classification von nicht-gelabelten Daten -> vor allem Clustering, siehe V3 – Folie 27, V8 – Folie 30 und k-means Clustering
- Methode für supervised classification von gelabelten Daten (HEUTE)

Ereignismenge: Beispiel

Gegeben sei eine Menge an Beobachtungen, z.B. ob ich in der Vergangenheit bei bestimmten äußeren Bedingungen zum Tennisspielen oder Golfspielen gegangen bin.

Wetter	Temperatur	Luftfeuchtigkeit	Wind	Gespielt?
Sonnig	Heiß	Hoch	Schwach	Nein
Sonnig	Heiß	Hoch	Stark	Nein
Bewölkt	Heiß	Hoch	Schwach	Ja
Regen	Mild	Hoch	Schwach	Ja
Regen	Kühl	Normal	Schwach	Ja
Regen	Kühl	Normal	Stark	Nein
Bewölkt	Kühl	Normal	Stark	Ja
Sonnig	Mild	Hoch	Schwach	Nein
Sonnig	Kühl	Normal	Schwach	Ja
Regen	Mild	Normal	Schwach	Ja
Sonnig	Mild	Normal	Stark	Ja
Bewölkt	Mild	Hoch	Stark	Ja
Bewölkt	Heiß	Normal	Schwach	Ja
Regen	Mild	Hoch	Stark	Nein

Angelehnt an:

<http://www.coli.uni-saarland.de/courses/mathe3/SS12/Vorlesungen/decisiontree.pdf>

Entscheidungsbäume: Prinzip

Wir möchten nun einen Entscheidungsbaum konstruieren, mit dem man diese Entscheidungen formalisieren kann.

D.h. wenn es heiß und feucht ist, dann spiele ich nicht

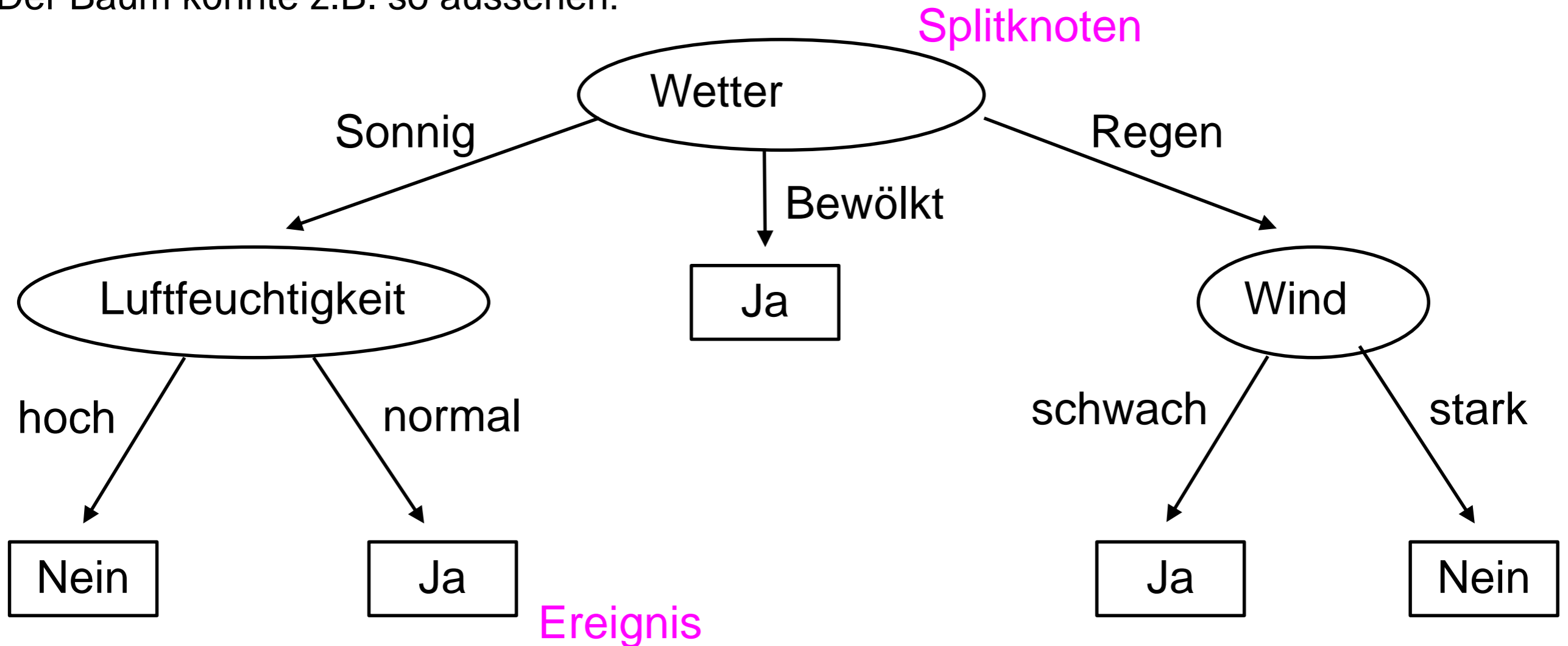
Wenn es kühl ist und starker Wind weht, dann spiele ich nicht.

Wenn es mild ist, die Sonne scheint und der Wind schwach ist, dann spiele ich etc.

Wie baut man den besten Baum?

Entscheidungsbäume

Der Baum könnte z.B. so aussehen:



- Beschreibt dieser Baum korrekt die existierenden Beobachtungen?
- Gibt es vielleicht einen besseren bzw. einfacheren Baum? Zum Beispiel haben wir die Temperatur hier ja gar nicht berücksichtigt.

Konstruiere Split-Knoten

Der von Russ Quinlan entwickelte **ID3-Algorithmus**

ist ein sogenannter *greedy* (gieriger) Algorithmus.

Quinlan, J. R. 1986. Induction of Decision Trees.

Mach. Learn. 1, 1 (Mar. 1986), 81–106



Russ Quinlan
<https://www.rulequest.com/Personal/>

So ein Algorithmus trifft jeweils die lokal bestmögliche Entscheidung.

Wenn eine Entscheidung getroffen wurde, wird diese später nicht mehr revidiert.

Man packt alle Beobachtungen in einen Wurzelknoten und sucht nun das beste Kriterium (d.h. den obersten Splitknoten), die Beobachtungen aufzuteilen.

Informationsgewinn

Ziel: Wir möchten die Entropie in den verbleibenden Klassen möglichst reduzieren.

Für jedes Attribut betrachten wir die Differenz zwischen der vorhandenen Entropie und der verbleibenden Entropie nach Einführung eines Splitknotens bzgl. dieses Attributs:

$$\text{Zugewinn}(S, A) = H(S) - \sum_{v=1}^k \frac{|S_v|}{|S|} H(S_v)$$

S: Datensatz

A: Attribut

H(S) Entropie im Datensatz oder Untermenge davon

S_v Untermenge von S, für die A den Wert v hat

|S| Mächtigkeit von S

$|S_v|$ Mächtigkeit von S_v

Entropie im Datensatz

Zur Berechnung der Entropie benutzen wir die übliche Shannon-Entropie

$$H(S) = - \sum_{c=1}^{|C|} p_c \log_2 p_c$$

S: Datensatz

|C| Anzahl an Kategorien (in diesem Beispiel: 2 oder 3)

p_c Anteil der Instanzen, die Kategorie c angehören

Die Entropie ist am höchsten, wenn die Häufigkeit p_c in allen Kategorien gleich ist.

Die Entropie ist gleich Null, wenn die Häufigkeit in einer Kategorie gleich 1 ist. Warum?

Beispiel für den Wetterdatensatz

Wir berechnen jetzt, welches Attribut an der Wurzel (d.h. dem obersten Splitknoten) in unserem Beispielproblem den größten Informationsgewinn bringen würde:

$$\text{Zugewinn}(S, A) = H(S) - \sum_{v=1}^k \frac{|S_v|}{|S|} H(S_v)$$

Die Gesamtentropie $H(S)$ ist – mit 14 Instanzen (Ereignissen), davon 5 x Nein, 9 x Ja –

$$H(S) = -\frac{9}{14} \log_2 \frac{9}{14} - \frac{5}{14} \log_2 \frac{5}{14} = 0.940286$$

Für das Attribut Wind = „schwach“ gab es 8 Instanzen - 6 mal wurde gespielt, 2 mal nicht

$$H(S_{\text{schwach}}) = -\frac{6}{8} \log_2 \frac{6}{8} - \frac{2}{8} \log_2 \frac{2}{8} = 0.8112781$$

Für das Attribut Wind = „stark“ gab es 6 Instanzen - 3 mal wurde gespielt, 3 mal nicht

$$H(S_{\text{stark}}) = -\frac{3}{6} \log_2 \frac{3}{6} - \frac{3}{6} \log_2 \frac{3}{6} = 1$$

Beispiel für den Wetterdatensatz

$$\text{Gain}(S, A) = H(S) - \sum_{v=1}^k \frac{|S_v|}{|S|} H(S_v)$$

Aus der Gesamtentropie und der Entropie für die beiden Ereignisse Wind = schwach / stark ergibt sich

$$\begin{aligned} \text{Gain}(S, \text{Wind}) &= H(S) - \frac{|S_v|}{|S|} H(S_v) - \frac{|S_v|}{|S|} H(S_v) = 0.940286 - \frac{8}{14} \cdot 0.8112781 - \frac{6}{14} \cdot 1 \\ &= 0.04812709 \end{aligned}$$

Auf ähnliche Weise bekommt man

$$\text{Gain}(S, \text{Wetter}) = 0.247$$

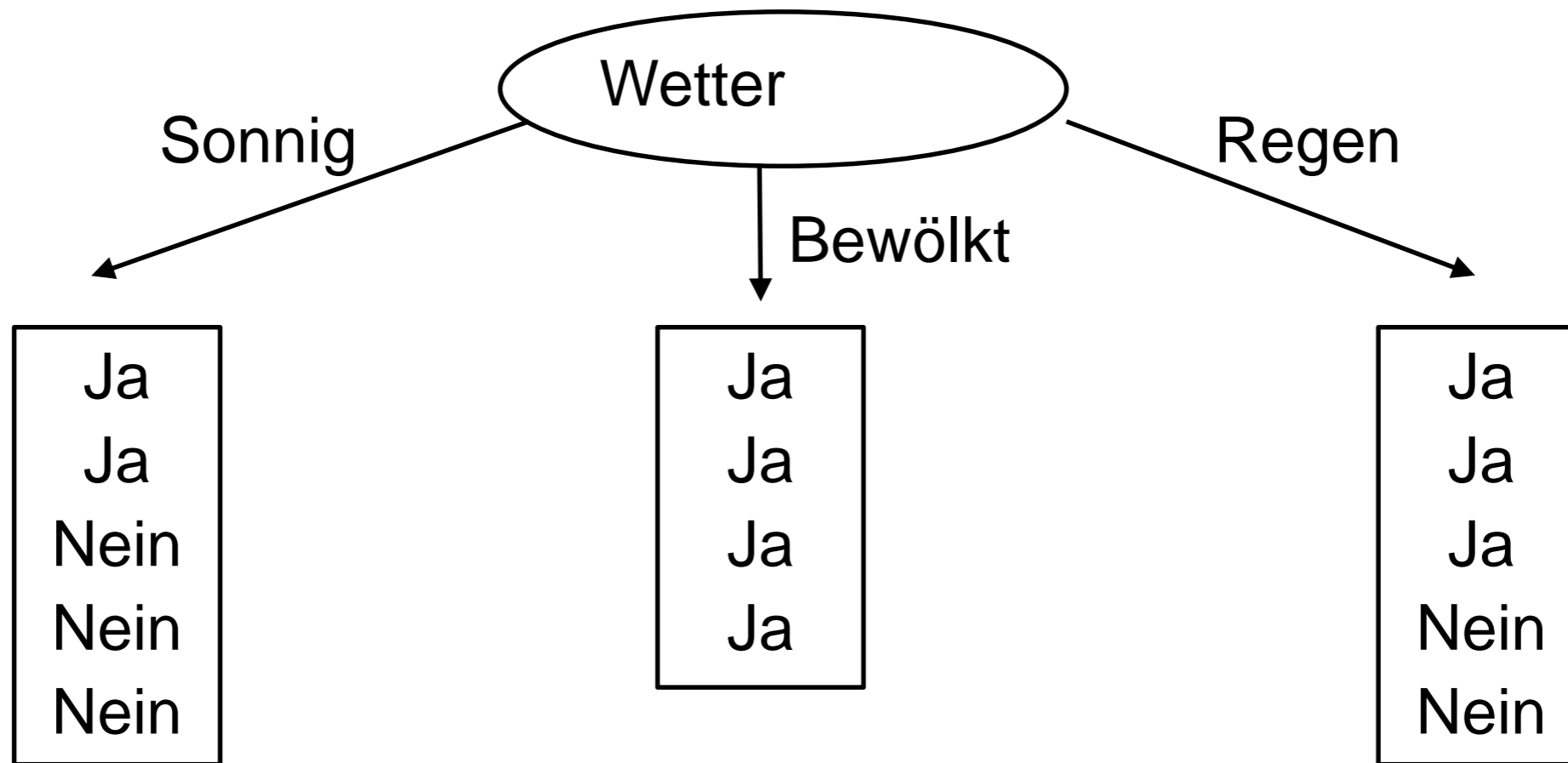
$$\text{Gain}(S, \text{Temperatur}) = 0.029$$

$$\text{Gain}(S, \text{Luftfeuchtigkeit}) = 0.152$$

„Wetter“ bringt also den größten Informationsgewinn und wird als Baumwurzel gewählt.

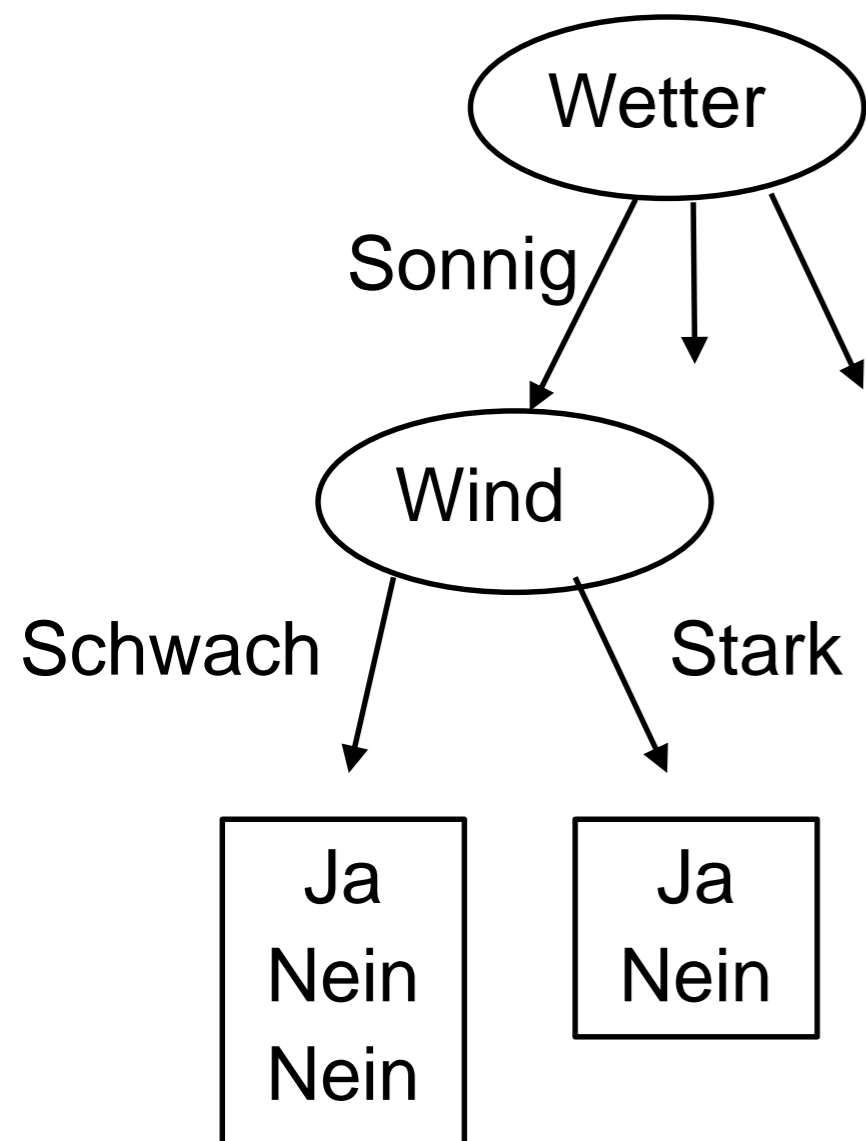
Entscheidungsbäume

Der Baum nach Einführung des Splitknotens „Wetter“ sieht folgendermaßen aus:

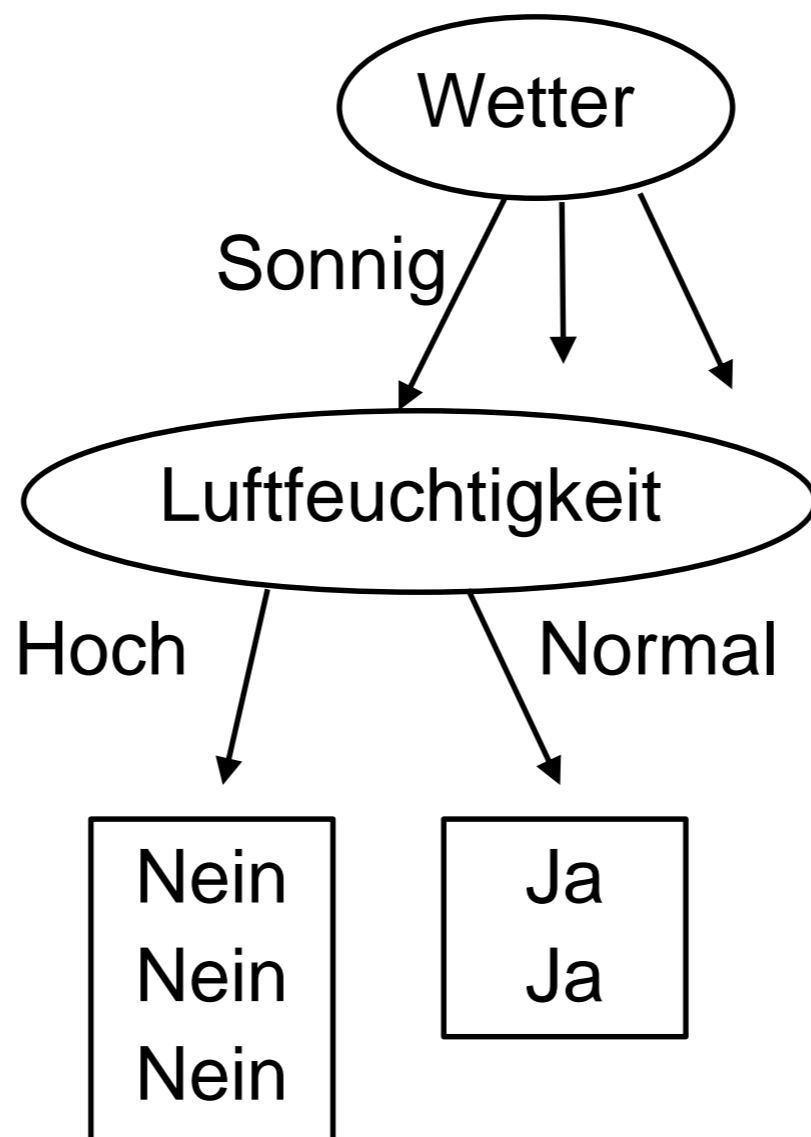


Im nächsten Schritt muss man für jeden der 3 Attributwerte, z.B. „sonnig“, dasselbe Verfahren für die Untermenge an Ereignissen rekursiv noch einmal anwenden.

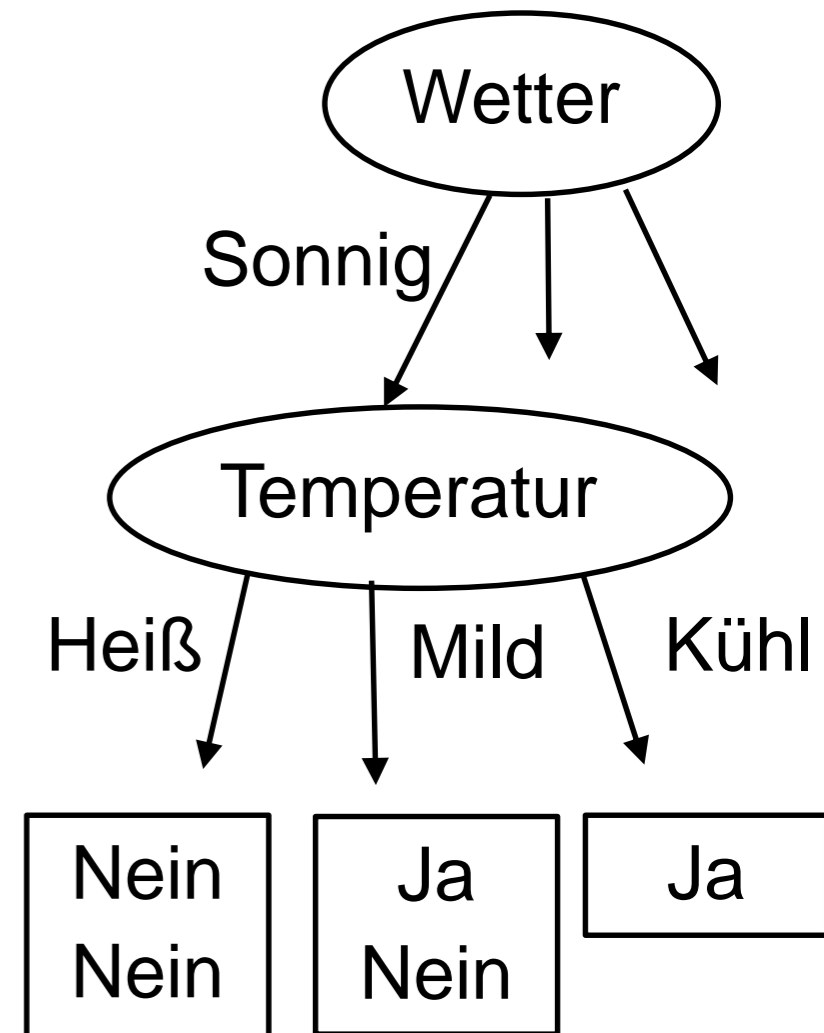
Entscheidungsbaum Wetter - Rekursion



$$\text{Gain}(S, \text{Wind}) = 0.020$$



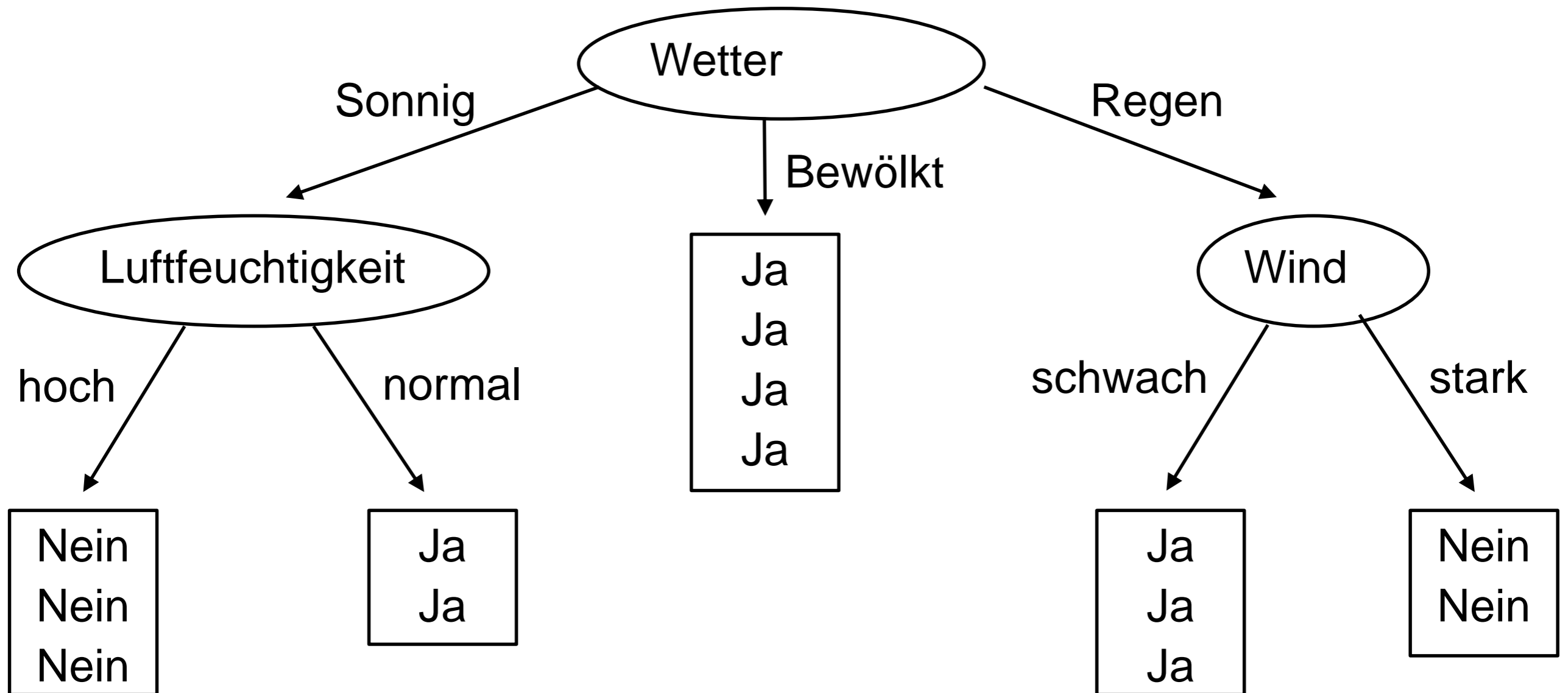
$$\text{Gain}(S, \text{Luftfeuchtigkeit}) = 0.971$$



$$\text{Gain}(S, \text{Temperatur}) = 0.571$$

Endergebnis – ID3-Algorithmus

Damit ergibt sich folgender Baum:



- In diesem Baum werden alle Instanzen unter einem Blatt jeweils gleich klassifiziert (Entropie = 0)
- Das Attribut „Temperatur“ wurde gar nicht benutzt.

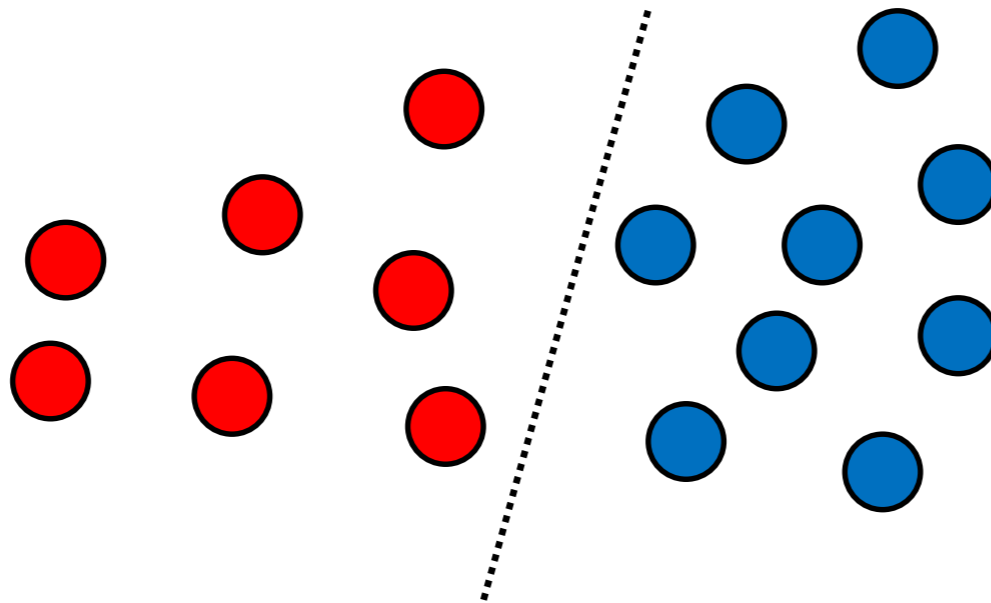
Eigenschaften des ID3-Algorithmus

Der Algorithmus ist „greedy“, d.h. unvollständig. Es werden nicht alle möglichen Bäume betrachtet.

Es ist im Prinzip möglich, dass ein „besserer“ (d.h. kleinerer) Baum nicht gefunden wird. Dieser könnte als erste Wurzel ein Attribut benutzen, das anfangs nicht den höchsten Informationsgewinn bietet.

Support Vektor Maschinen

Gegeben sei wiederum eine gelabelte Datenwolke (rot/blau)



Wir möchten gerne eine Trennlinie konstruieren, die wir verwenden können, um weitere Datenpunkte in rot/blau zu klassifizieren.

Wir könnten „einfach“ eine Gerade dazwischen legen, etwa die gestrichelte Linie.

Präsentation angelehnt an https://www.tu-chemnitz.de/urz/itime/documents/Vortrag_Jens_Poenisch_SVM.pdf

SVM – Konstruktion der Trennebene

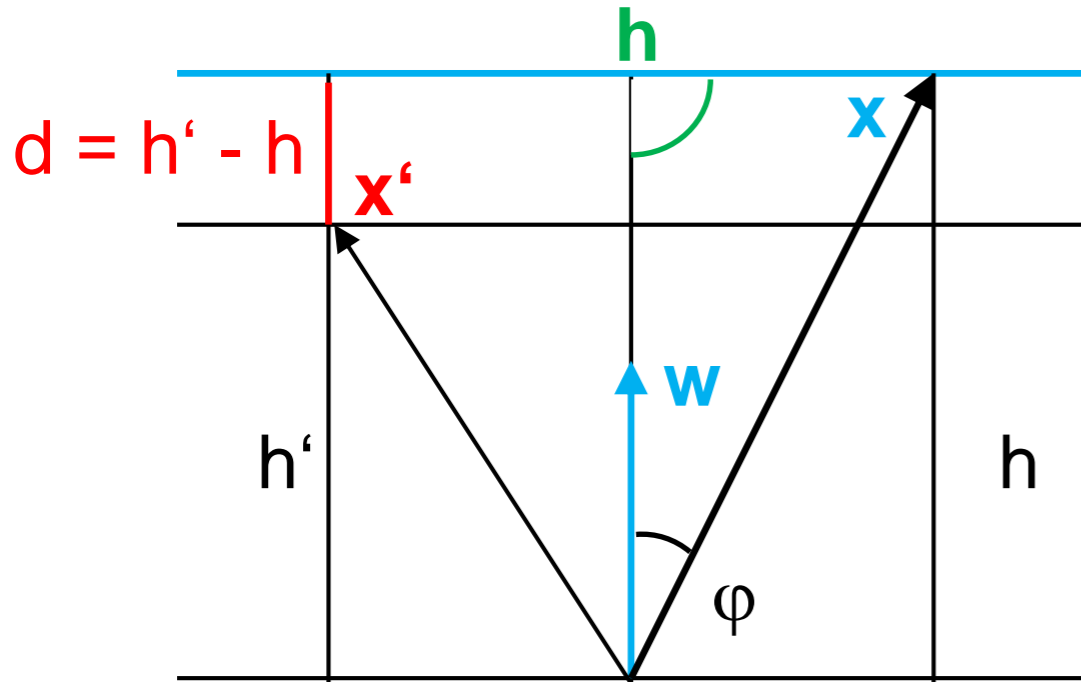
Bezeichnungen:

$\mathbf{x} = (x_1, \dots, x_n)^T$ Vektor mit n Komponenten

$\langle \mathbf{x}, \mathbf{y} \rangle = x_1 y_1 + \dots + x_n y_n$ Skalarprodukt der Vektoren \mathbf{x} und \mathbf{y}

$\|\mathbf{x}\| = \sqrt{x_1^2 + \dots + x_n^2}$ Euklidische Norm (= Länge) des Vektors \mathbf{x}

Abstand eines Punktes von der Ebene



\mathbf{w} ist ein Vektor senkrecht zur blauen Ebene,

\mathbf{x} ist der Vektor zu einem beliebigen Punkt auf der blauen Ebene,

φ ist der Winkel zwischen \mathbf{w} und \mathbf{x} .

Abstand der Ebene vom Ursprung $h = \|\mathbf{x}\| \cos \varphi$

Skalarprodukt $\langle \mathbf{w}, \mathbf{x} \rangle = \|\mathbf{w}\| \cdot \|\mathbf{x}\| \cdot \cos \varphi$

Damit gilt $h = \frac{\langle \mathbf{w}, \mathbf{x} \rangle}{\|\mathbf{w}\|}$

Wir setzen $b := -h \cdot \|\mathbf{w}\|$

$h = \frac{\langle \mathbf{w}, \mathbf{x} \rangle}{\|\mathbf{w}\|}$ formen wir um in

$\frac{\langle \mathbf{w}, \mathbf{x} \rangle}{\|\mathbf{w}\|} - h = 0$ bzw. in die

Ebenengleichung $\frac{\langle \mathbf{w}, \mathbf{x} \rangle + b}{\|\mathbf{w}\|} = 0$

Der Punkt \mathbf{x}' liegt auf einer Parallelebene mit Abstand

$h' = \frac{\langle \mathbf{w}, \mathbf{x}' \rangle}{\|\mathbf{w}\|}$ vom Ursprung, also

gilt für den Abstand von der Ebene

$$d = h' - h = h' + \frac{b}{\|\mathbf{w}\|} = \frac{\langle \mathbf{w}, \mathbf{x}' \rangle + b}{\|\mathbf{w}\|}$$

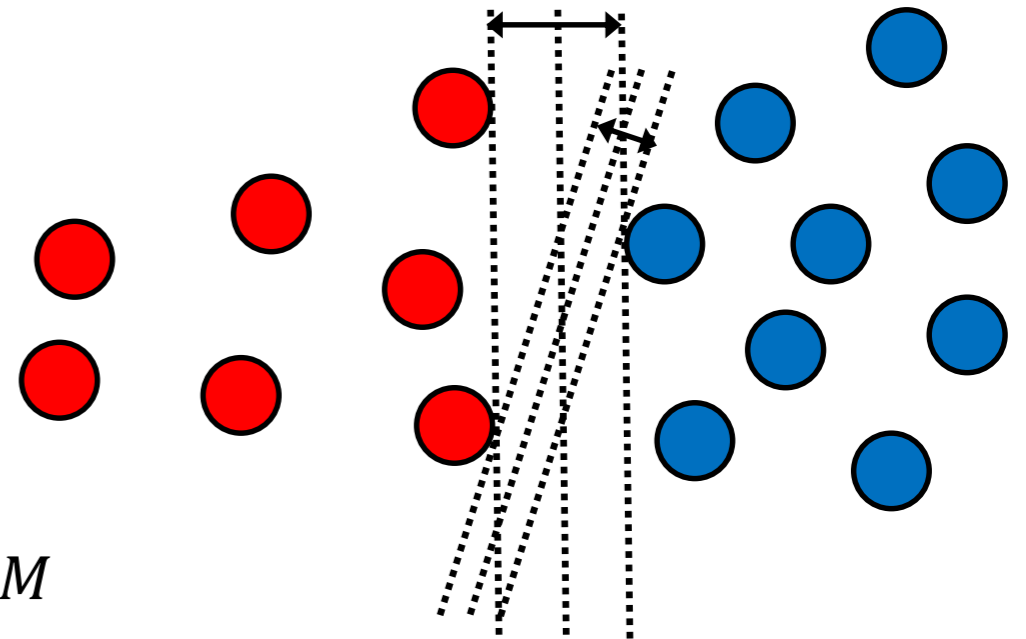
SVM – Grundidee

Ziel: lege die Hyperebene (in zwei Dimensionen ist dies einfach eine Gerade) so, dass der kleinste Abstand eines Punktes zur Ebene möglichst groß ist und auf der „richtigen“ Seite/Klasse liegt.

$$t \left(\frac{\langle \mathbf{w}, \mathbf{x}' \rangle + b}{\|\mathbf{w}\|} \right) \rightarrow \max, t \in \{-1, +1\}$$

Man kann dies umformen in

$$\arg \min_{\mathbf{w}, b} \frac{1}{2} \|\mathbf{w}\|^2, \text{ wobei } t_m (\langle \mathbf{w}, \mathbf{x}_m \rangle + b) \geq 1, m = 1, \dots, M$$



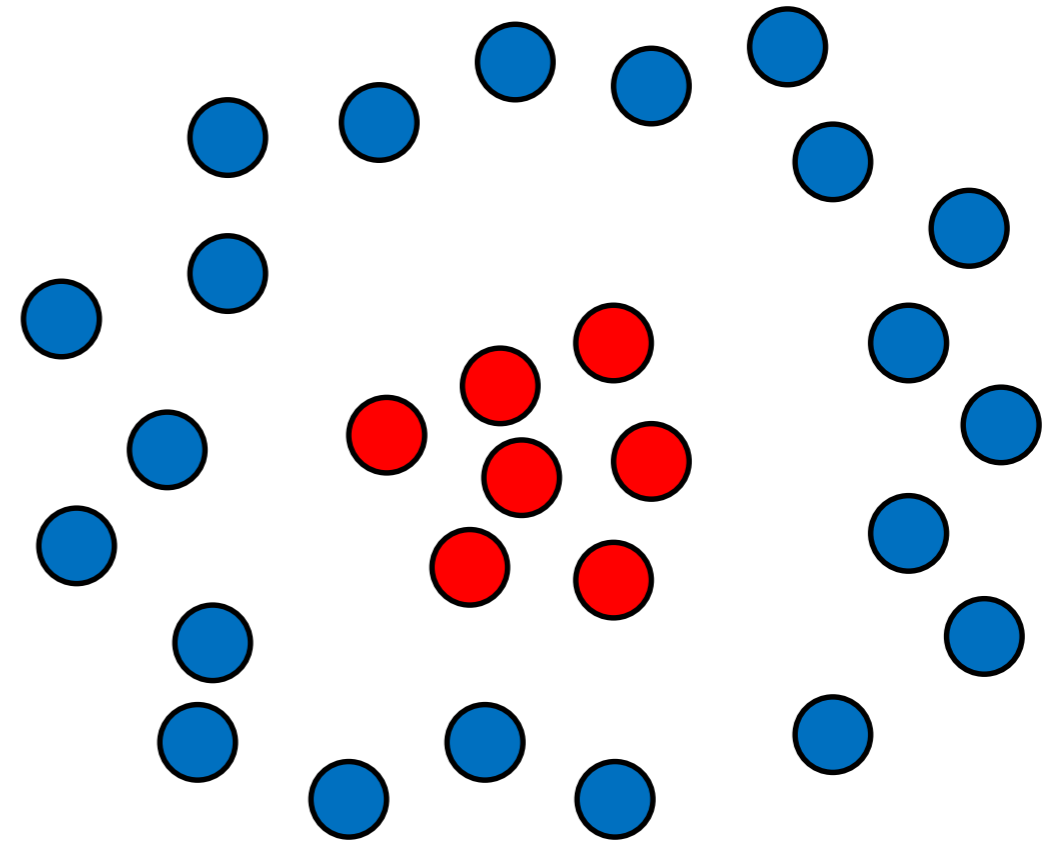
- Man möchte einen möglichst breiten «Streifen» unterhalb und oberhalb der Ebene erhalten, der frei von Trainingspunkten ist.
- Die Punkte auf dem «Rand» heißen *Supportvektoren*.
- Neue Punkte werden durch diese Trennebene hoffentlich «richtig» klassifiziert.
- Wähle die Ebene so, dass der kleinste Abstand eines Punktes zur ihr maximiert wird.

Problem 1: keine lineare Trennung möglich

Es könnten jedoch auch Fälle auftreten, wie der rechts gezeigte, bei dem sich die beiden Gruppen von Datenpunkten offensichtlich nicht durch eine Gerade trennen lassen.

In diesem Fall wäre ein Kreis um die roten Punkte am besten geeignet.

Allgemein überprüft man meist, ob sich die Datenpunkte nach der Transformation mit einer Kernelfunktion (typisch: radialer Gauss-Kernel, polynomisch, sigmoidal) besser in 2 Klassen separieren lassen.



Gauß-Kernel

$$k(\mathbf{x}_i, \mathbf{x}_k) = e^{-\gamma \|\mathbf{x}_i - \mathbf{x}_k\|^2}$$

Problem 2: Behandlung von Ausreißern

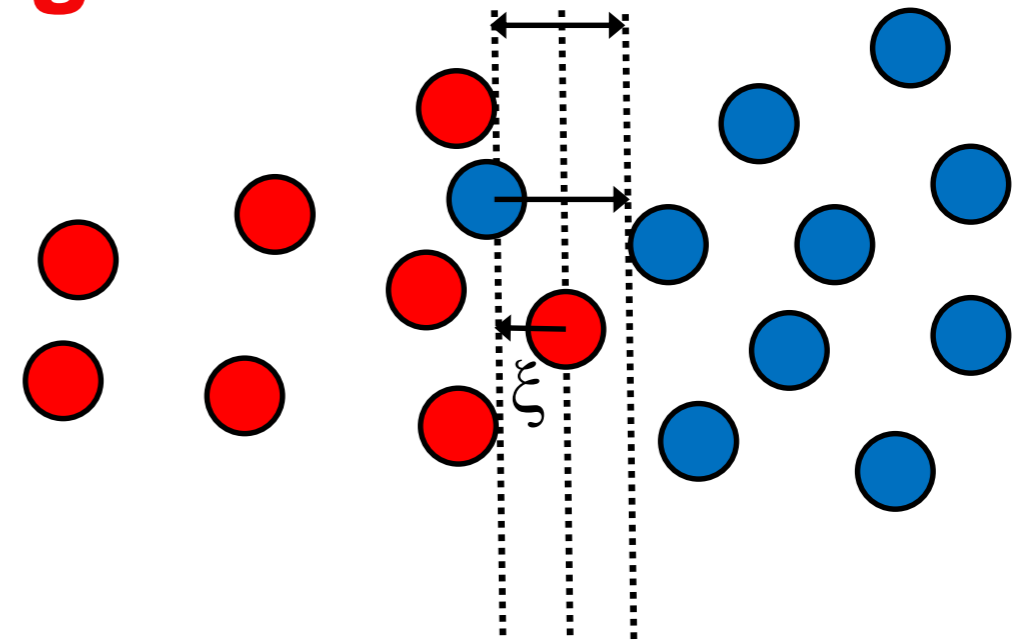
Durch «Ausreißer» liegen einzelne Datenpunkte in der «falschen» Klasse, es ist keine saubere Trennung möglich.

Erlaube einzelne Ausreißer durch eine „Slackvariable“ ξ . Man addiert einen Korrekturwert auf die eigentlich verletzte Randbedingung, um diese zu erfüllen.

Die Summe der Korrekturen soll möglichst klein sein.

Beim Trainieren der SVM variiert man den „Regularisierungsparameter“ C , also den Einfluss von ξ .

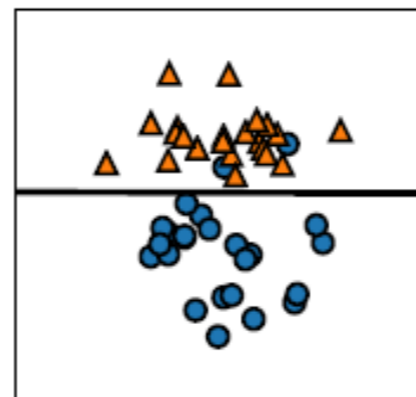
Je größer C , desto komplexer das Modell (höhere «Bestrafung» der Ausreißer).



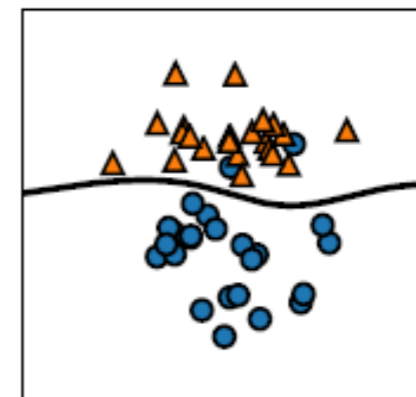
$$\arg \min_{\mathbf{w}, b, \xi_1, \dots, \xi_M} \left\{ \frac{1}{2} \|\mathbf{w}\|^2 + C \sum_{m=1}^M \xi_m \right\}$$

$$\text{wobei } t_m (\langle \mathbf{w}, \mathbf{x}_m \rangle + b) + \xi_m \geq 1, \xi_m \geq 0, m = 1, \dots, M$$

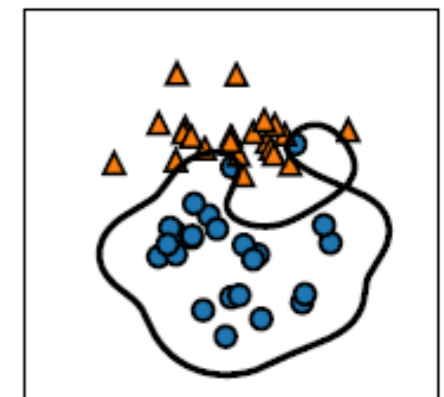
$\gamma = 0.005, C = 0.5$



$\gamma = 0.1, C = 100$



$\gamma = 1.0, C = 1000$



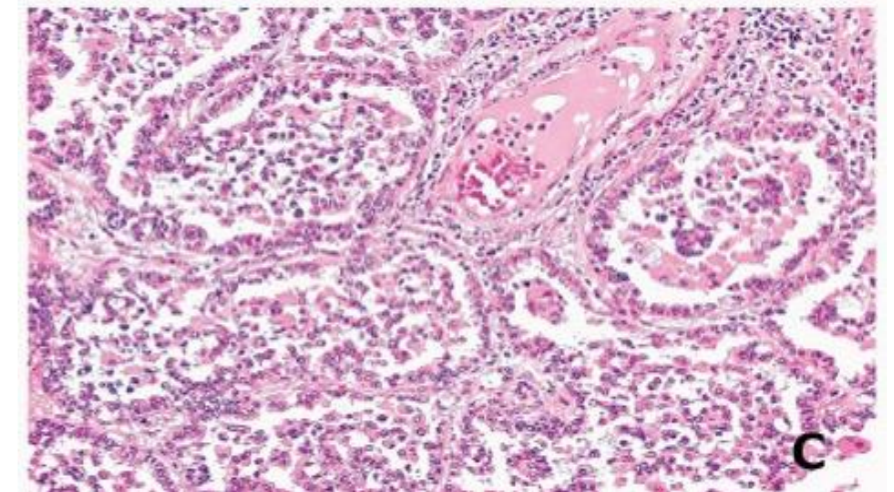
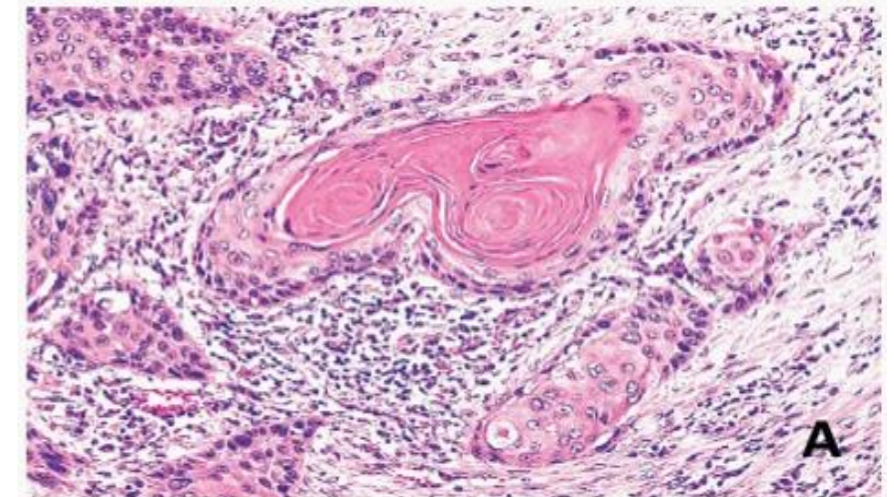
Fallstudie: Klassifizierung von Lungenkrebs anhand von histologischen Aufnahmen

Classification and mutation prediction from non-small cell lung cancer histopathology images using deep learning

Nicolas Coudray^{1,2,9}, Paolo Santiago Ocampo^{3,9}, Theodore Sakellaropoulos⁴, Navneet Narula³, Matija Snuderl³, David Fenyö^{5,6}, Andre L. Moreira^{3,7}, Narges Razavian^{8*} and Aristotelis Tsirigos^{1,3*}

- 2 häufigste Formen von Lungentumoren:
- LUSC – lung squamous cell carcinoma (SCC):
- SCCs treten in Plattenepithelzellen auf.
- LUAD – lung adenocarcinoma - Adenocarcinome entstehen in Drüsen an verschiedenen Stellen im Körper.
- Beides sind **nicht-kleinzellige Lungentumore**

SCC



AD - BAC

Coudray et al. Nature Medicine 24, 1559–1567 (2018)
<http://www2.keelpno.gr/blog/?p=1391>

Behandlung von LUAC / LUSC

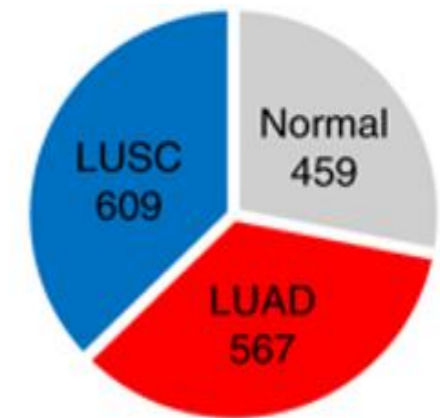
- Stadium I – Operation oder Bestrahlung
- Stadium II – Operation und Chemotherapie oder Bestrahlung
- Stadium III – sequenzielle oder gleichzeitige Chemotherapie und Bestrahlung
- **Stadium IV – Patienten-Genom wird wichtig**
 - Zytotoxische Kombinationschemotherapie
 - Kombinationschemotherapie mit monoklonalen Antikörpern
 - Erhaltungstherapie nach primärer Chemotherapie
 - EGFR Tyrosinkinase-Inhibitoren
 - ALK Inhibitoren (für Patienten mit ALK-Translokationen)
 - ROS1 Inhibitoren (für Patienten mit ROS1-Umlagerung)
 - BRAFV600E und MEK Inhibitoren (für Patienten mit BRAFV600E Mutationen)
 - Immun Checkpoint-Inhibitoren mit bzw. ohne Chemotherapie

https://www.cancer.gov/types/lung/hp/non-small-cell-lung-treatment-pdq#section/_48406

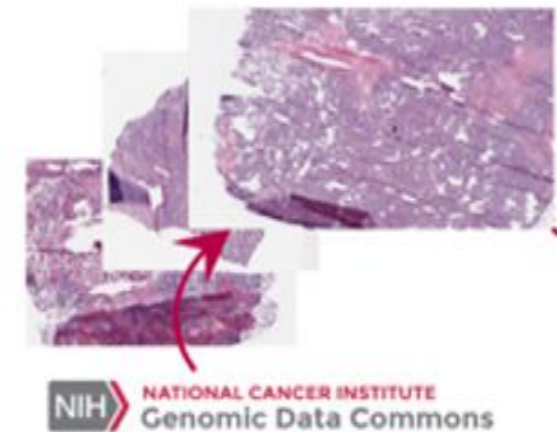
Klassifizierung von Tumorgewebe

Q: Können Deep Learning-Methoden LUAD / LUSC / normale (gesunde) Proben mit ähnlicher Genauigkeit wie ein medizinischer Experte (Pathologe) klassifizieren?

Anzahl an Proben aus
TCGA (The Cancer Genome Atlas)



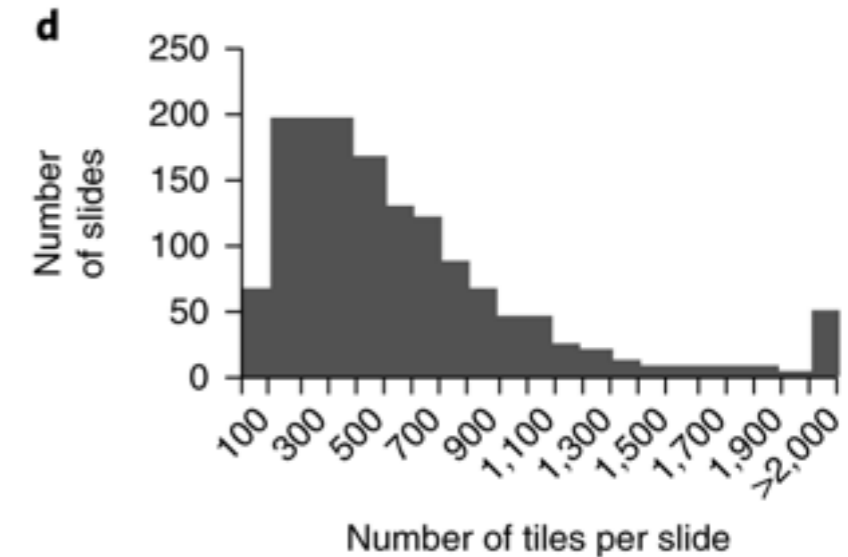
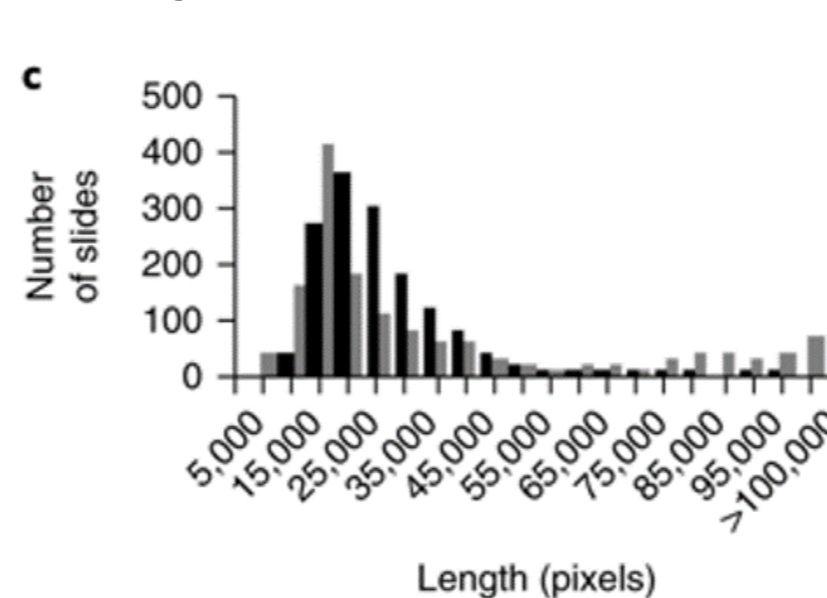
i Download
from GDC
database



Coudray et al. Nature Medicine 24, 1559–1567 (2018)

Klassifizierung von Tumorgewebe

- Die einzelnen Aufnahmen sind „zu groß“ um sie direkt als Input für neuronales Netzwerk zu verwenden



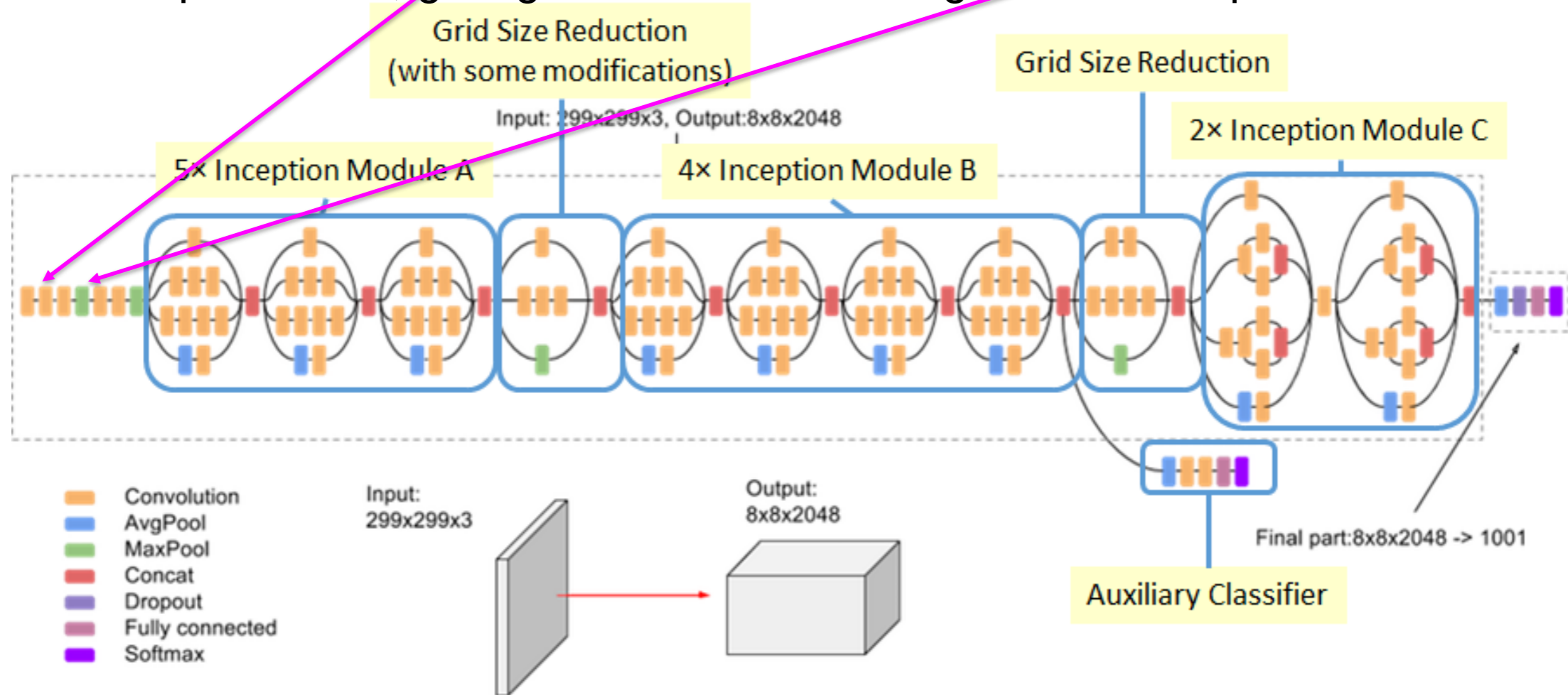
- Idee: spalte jedes Bild in kleine Quadraten mit 512×512 Pixeln auf.

- Dadurch vergrößert sich die Menge an Trainingsdaten stark.
- Spalte Daten in 70% fürs Training, 15% für Validierung, und 15% für Testen auf.
- -> **etwa 1 Million Bildblöcke**

Coudray et al. Nature Medicine 24, 1559–1567 (2018)

Deep learning Modell

- Die Autoren verwendeten eine convolutional neural network-Architektur namens inception v3 architecture³⁶, die von Google entwickelt wurde.
- Zu Beginn gibt es 5 convolution-Knoten, die mit 2 max pooling Operationen verknüpft werden, gefolgt von 11 aneinandergereihten inception Modulen

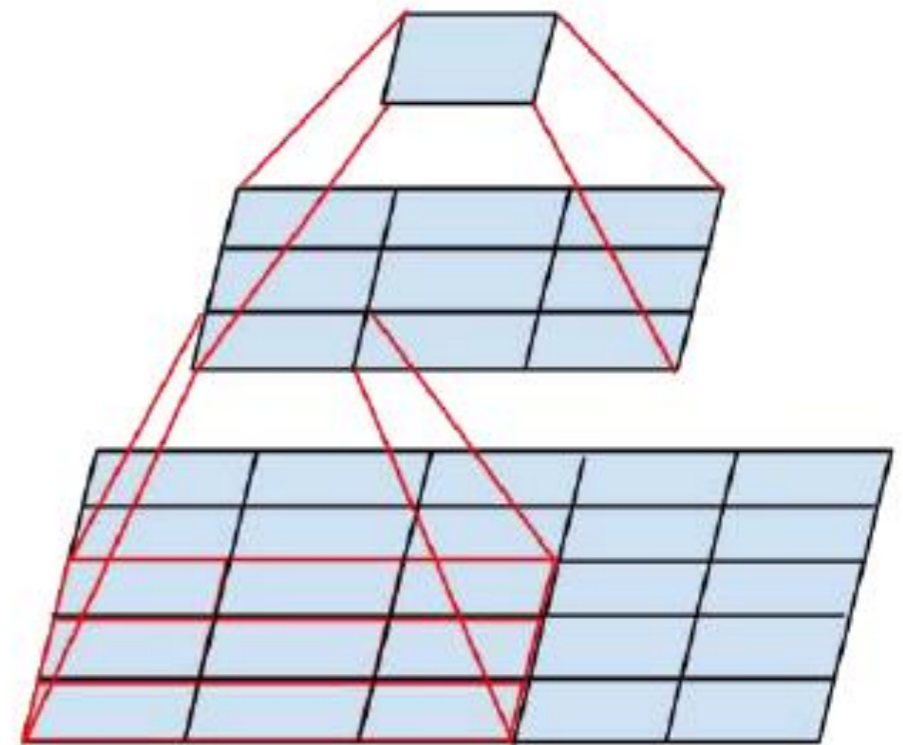


- „convolution“ = engl. für die Berechnung eines Faltungsintegrals = räumliche Summierung

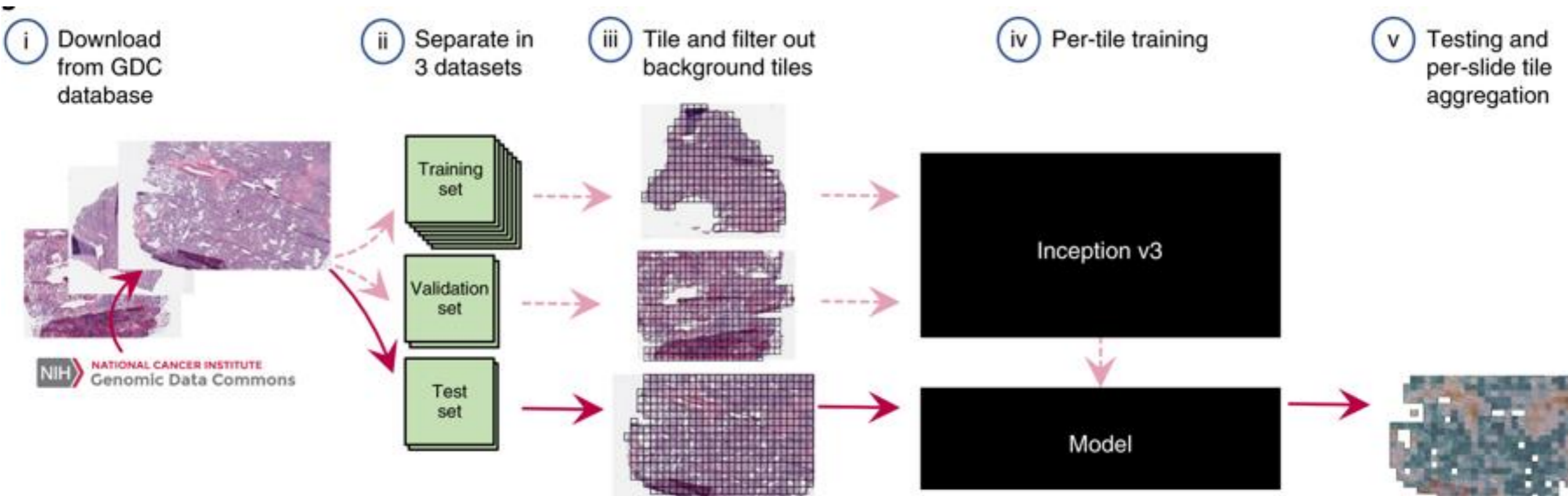
medium.com

inception v3 architecture36

- Die hintereinander angeordneten convolution-Knoten dienen dazu, die Anzahl an Parametern zu reduzieren, die trainiert werden müssen.
- Im Beispiel unten ersetzen zwei 3×3 convolutions eine 5×5 convolution
- Statt einer Schicht mit $5 \times 5 = 25$ Parametern verwendet man daher zwei Schichten mit je 3×3 Parametern, also insgesamt nur 18.



Workflow



Die Klassifizierung von gesundem versus Tumorgewebe klappte sehr gut (~0.99 AUC “area under the curve”).

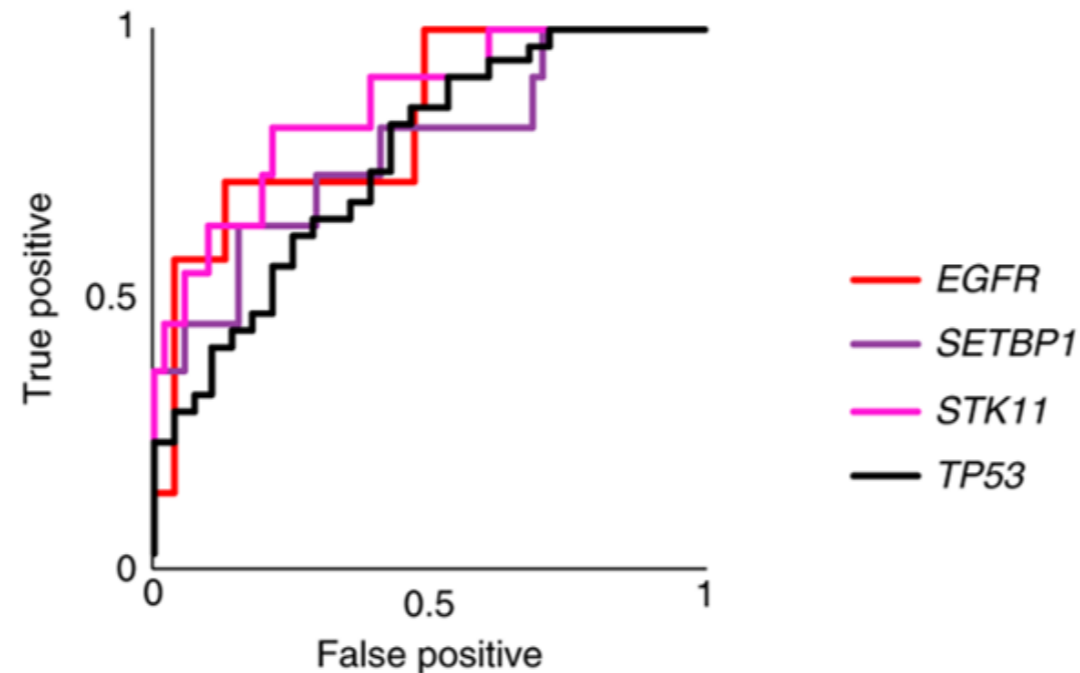
Unterscheidung von LUSC und LUAD ebenfalls mit 0.97 AUC möglich.

Das ist dieselbe Genauigkeit, die man erhält, wenn man 3 ausgebildete Pathologen bittet, dieselben Bilder zu klassifizieren.

Coudray et al. Nature Medicine 24, 1559–1567 (2018)

Klassifizierung von genetischen Varianten

- Können CNNs das Auftreten bestimmter Genmutationen aus den Bildern erkennen?



- Einigermaßen. Die Genauigkeit (AUC = das Flächenintegral unter der Kurve) reicht von 0.64 (LRP1B) bis 0.84 (STK11).
- Falls es mehr Trainingsdaten gäbe, sollten sich die Ergebnisse noch verbessern.
- Tool kann Pathologen bereits heute bei Routine-Tätigkeiten unterstützen.

Coudray et al. Nature Medicine 24, 1559–1567 (2018)

Klausur-relevanter Vorlesungsstoff

Vorlesung	Folien
1	16-25, 29
2	3-45
3	6-21, 32-45
4	21, 23
5	11-14, 19-34, 38-39, 41
6	1-35, 39
7	5-6, 8-12, 16-23
8	9-13, 21-29, 33-42
9	7-8, 14-18, 24-33
10	11-26
11	8-9, 13-15, 25
12	1-19

Relevant sind immer die Foliennummern unten rechts.