

# 1 Begriffe

## 2 Definition Formel & Algorithmen

### 2.1 Specification of learning tasks

Realworld  $\rightarrow$  Modelworld

#### 1. Reale Welt:

- $O$  - Menge an Objekten
- $C$  - Menge an Klassen
- $\gamma : O \rightarrow C$  - idealer Classifier für  $O$

Aufgabe-Klassifizierung:

- bestimme von  $o \in O$  die Klasse  $\gamma(o) \in C$

#### 2. Model Welt

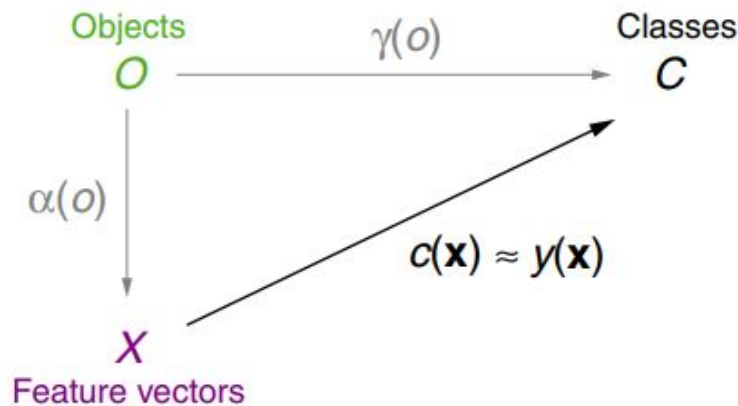
- $X$  - Menge von feature Vektoren
- $C$  - Menge von Klassen
- $c : X \rightarrow C$  - idealer Classifier für  $X$  ( $c$  ist unbekannt)
- $D = \{(\mathbf{x}_1, c(\mathbf{x}_1)), \dots, (\mathbf{x}_n, c(\mathbf{x}_n))\}$  - Menge von Beispielen (*bereits klassifiziert*)

Todo: Schätze  $c(\mathbf{x})$ , welche implizit durch  $D$  gegeben sind, durch Model-Funktion  $y(\mathbf{x})$

### 2.2 Machine Learning:

1. Collect real-world examples of the form  $(o, \gamma(o)), o \in O$
2. abstract the objects towards feature vectors  $\mathbf{x} \in X$ , where  $\mathbf{x} = \alpha(o)$
3. Formuliere Model-Funktion:  $y : X \rightarrow C, \mathbf{x} \mapsto y(\mathbf{x})$
4. Nutze Statistik, Theorie und Algorithmen aus ML um den fit zwischen  $c(\mathbf{x})$  und  $y(\mathbf{x})$  zu Maximieren, sodass  $y(\mathbf{x}) \approx c(\mathbf{x}) \approx \gamma(o)$

$O$	is a set of objects	Mushrooms
$C$	is a set of classes	essbar oder nicht
$X$	Set/ Menge der Feature Vektoren	Vektoren mit den Eigenschaften der Pilze
$\mathbf{X}$	Feature Vektor; Feature Space; Feature Domain	
$\mathbf{x}$	one feature vector	ein Pilz
$x$	single feature	ein eigenschaft von einem Pilz (z.B. braun)
$\mathbf{w}$	Weight-Vektor	
$D = (x_1, c(x_1)), \dots, ((x_n, c(x_n)))$	classification knowledge alias the example set	schon bestimmte Pilze
$D = \{(x_1, y_1), \dots, ((x_n, y_n))\}$	same but for regression	Mietpreise von Wohnungen
$\alpha(o) = x$	model formation function	Funktion zur Bestimmung der Pilz-Eigenschaft
$y(x)$ (Ziel $y(x) = c(x)$ )	model function/ classifier	Bestimmung ob essbar anhand der Eigenschaft
$\gamma(o)$ mit $O \rightarrow C$	ideal target function / ideal classifier für O	Pilzexperte der Pilze bestimmt
$c(x)$ mit $c: X \rightarrow C$	ideal target function/ ideal classifier for X ; indirekt gegeben durch D	Theoretische Ideale Bestimmung aller Pilze
$h(x), h: X \rightarrow \{0, 1\}$	hypothesis/ model function to approximate $c(x)$	wenn es sonnig ist mache ich Sport (Regeln/ H
$y_i$	ground truth for $x_i \in X$	der Mietpreis für eine Wohnung
$H$	hypothesis space	Menge alle möglichen Hypothesen/Regeln
$V_{H,D} = \{h(x)   h(x) \in H \wedge (\forall (x, c(x)) \in D : h(x) = c(x))\}$	version space :	Wenn es
$\mathcal{X} \in X / \mathcal{C} \in C$	random variables	ein random Pilz bzw. eine random Klasse
$p(x, c) = P(X = x, C = c)$	the probability of the joint event that $x \in X$ & $x$ belongs to class $c \in C$	die Wahrscheinlichkeit das ein Pilz zu einer bes



## 2.3 LMS: Least Mean Squared

Ziel: Fitting  $y(x)$ ; Anpassung der weights, sodass Klassifizierungsfehler möglichst gering sind.

Input:  $D$  - Trainingsdaten  $(\mathbf{x}, c(\mathbf{x}))$  mit  $\mathbf{x} \in \mathbf{R}^p$  und Zielklasse  $c(\mathbf{x})$   $\eta$   
 Learning rate, kleine positive Konstante

IGD( $D, \eta$ )

```

1. initialize_random_weights(w), t = 0
2. REPEAT
3.   t = t + 1
4.   FOREACH (x, c(x)) ∈ D DO
5.     δ = c(x) - y(x) // y(x) (*) = wTx, δ ∈ R.
6.     Δw (*) = η · δ · x // δ · x is the derivative of ℓ2(c(x), y(x)).
7.     w = w + Δw
8.   ENDDO
9. UNTIL(convergence(c(D), y(D)) OR t > tmax)
10. return(w)

```

## 2.4 Lineare Regression

Grundformel für lineare Gerade, mit:

$y(x)$  - abhängige Variabel  $x$  - unabhängige Variable  $w_1$  - Anstieg der Geraden  $w_0$  - Schnittpunkt der y-Achse

$$y(x) = w_0 + w_1 \cdot x \quad (1)$$

Wobei das minimale  $w_0$  und  $w_1$  sich ergeben aus:

$$w_1 = \frac{\sum_{i=1}^n (x_i - \bar{x} \cdot (y_i - \bar{y}))}{\sum_{i=1}^n (x_i - \bar{x})^2} \quad (2)$$

$$w_0 = \bar{y} - w_1 \cdot \bar{x} \quad (3)$$

**Goodness of Modelfit, Regressionerror** Residual sum of Squares(RSS).

Der Residue wird aus der Differenz zwischen (beobachteten) Realwelt-Wert  $y_i$  und geschätzten/modelierten Wert  $y(\mathbf{x}_i)$

$$RSS(\mathbf{w}) = \sum_{i=1}^n (y_i - y(\mathbf{x}_i))^2 \quad (4)$$

**Higher-Dimensional Feature Space** ML:ll-16

## 2.5 Concept Learning

Setting:

$X$  - Menge an Feature Vektoren

$C$  - Ist eine Menge mit zwei Klassen *Beispiel*:  $\{0,1\}, \{ja, nein\}$

$c : X \rightarrow C$  -idealer Klssifizierer für  $X$

$D = \{(\mathbf{x}_1, c(\mathbf{x}_1)), \dots, (\mathbf{x}_n, c(\mathbf{x}_n))\} \subseteq X \times C$  - Menge mit Beispielen

Todo:

Schätze  $c(x)$ , was implizit durch  $D$  mit feature-Value-Muster

## 3 Definitionen Text

### 3.1 Supervised learning

Eine Funktion mithilfe von input-output-Daten lernen;  
automatisierte Klassifikation mit von Menschen bereits klassifizierten  
Daten als Grundlage

*Beispiel: optical character recognition*

### 3.2 Unsupervised learning

identifiziert/findet selbstständig Muster und Strukturen in Daten;

- automatisierte Kategorisierung durch Cluster Analysis
- Parameter Optimierung durch Expectation Maximation
- Feature Extrahierung durch Factor Analysis

*Beispiel: intrusion detection in a network data stream*

### 3.3 Reinforcement learning

”Learn, adapt, or optimize a behavior strategy in order to maximize the  
ownbenefit by interpreting feedback that is provided by the environment.”

*Beispiel: program to play tetris*

### 3.4 Feature Vektor

Ein Feature Vektor ist ein Vektor in dem jede Dimension eine  
Eigenschaft(Feature) des beschriebenen Objektes enthält.

$$\mathbf{x} = \begin{bmatrix} x_1 \\ x_2 \\ \cdot \\ \cdot \\ \cdot \\ x_n \end{bmatrix} \quad \text{Beispiel: Gitarre} = \begin{bmatrix} \text{Farbe : blau} \\ \text{Baujahr : 1997} \\ \text{Material : Holz} \\ \text{Elektrisch : ja} \end{bmatrix}$$

### 3.5 Ground Truth

Überprüfung der Klassifizierung eines Lernprozesses auf Richtigkeit für  
gewolltes Model.

*Beispiel: Überprüfung eines Spamfilter nach falsch kategorisierten Mails*

## 4 Linear Models

### 4.1 Overfitting

- fitting: der Prozess die Parameter einer Modelfunktion  $y$  so anzupassen das sie der der Beispieldaten  $D$  am besten passen
- overfitting: "Fitting the data more than is warranted"
- alias besser passen als berechtigt?
- Gründe:
  - zu komplizierte Modelfunktion (zu viele Features)
  - zu wenig Daten in  $D$
  - zu viel Datenrauschen
  - $D$  ist zu biased alias nicht repräsentativ
- Folgen:
  - kleiner Error auf  $D_{tr}$  anber großer Error auf  $D_{test}$  und IRL
  - loss of inductive Bias
  - increase of variance as a result of sensitivity to noise
- Overfitting finden:
  - Visuell untersuchen für Fälle mit Dimensionen  $< 3$  sonst embedding oder projizieren in kleinere Dimensionen
  - Validieren: wenn  $Err_{fit} = Err_{val}(y) - Err_{tr}(y)$  zu groß ist
- Overfitting vermeiden:
  - Early stopping through model selection: nach  $m$  schritten überprüfen ob sich  $Err_{fit}$  noch verkleinert und stoppen wenn er sich vergrößert
  - Qualität (schlechte Beispiele raus) und / oder Quantität (mehr Daten gleichen Rauschen aus) von  $D$  verbessern
  - Manually enforcing a higher bias by using a less complex hypothesis space alias Removing Features: In this approach, irrelevant features are removed from the dataset. This enhances the algorithm's ability to generalize
  - Regularization (WUHU!)

#### 4.1.1 Well- and Ill-posed problems

A mathematical problem is called well-posed if

1. a solution exists,
2. the solution is unique,
3. the solution's behavior changes continuously with the initial conditions.

Otherwise, the problem is called ill-posed.

#### 4.2 Regularization

Automatic adjustment of the loss function to penalize model complexity. Let  $L(\mathbf{w})$  denote a loss function used to optimize the parameters  $\mathbf{w}$  of a model function  $y(\mathbf{x})$ . Regularization introduces a trade-off between model complexity and inductive bias:

$$\mathcal{L}(\mathbf{w}) = L(\mathbf{w}) + \lambda * R(\mathbf{w})$$

where  $\lambda \geq 0$  controls the impact of the regularization term  $R(\mathbf{w}) \geq 0$ .  $\mathcal{L}$  is called "objective function".

#### 4.2.1 Regularized Linear Regression

$$\mathcal{L}(\mathbf{w}) = \sum_{i=1}^n (y_i - \mathbf{w}^T \mathbf{x}_i)^2 + \lambda \cdot \vec{w}^T \vec{w}$$

Estimate  $\mathbf{w}$  by minimizing the residual sum of squares:

$$\hat{w} = \underset{\mathbf{w} \in \mathbf{R}^{p+1}}{\operatorname{argmin}} \mathcal{L}(\mathbf{w})$$

$$\rightsquigarrow RSS(\mathbf{w}) = (\mathbf{y} - \mathbf{X}\mathbf{w})^T (\mathbf{y} - \mathbf{X}\mathbf{w}) + \lambda \mathbf{w}^T \mathbf{w}$$

Ableitung bilden um  $RSS(\mathbf{w})$  zu minimieren und man kommt auf:

$$\mathbf{w} = (X^T X + \operatorname{diag}(0, \lambda, \dots, \lambda))^{-1} X^T \mathbf{y}$$

$$\operatorname{diag}(0, \lambda, \dots, \lambda) = \begin{bmatrix} 0 & 0 & 0 & \dots & 0 \\ 0 & \lambda & 0 & \dots & 0 \\ 0 & 0 & \lambda & \dots & 0 \\ \vdots & \vdots & \vdots & \ddots & \vdots \\ 0 & 0 & 0 & \dots & \lambda \end{bmatrix}$$

$$\hat{y}(\mathbf{x}_i) = \hat{\mathbf{w}}^T \mathbf{x}_i$$

Um so höher  $\lambda$  um so einfacher ist die Funktion - Regularization archived!

## 5 Neural Networks

### 5.1 Perception Learning

Idee: Lass mal ein Gehirn programmieren!

Typisches Beispiel: Schrifterkennung

$$y(\mathbf{x}) = 1 \Leftrightarrow \left( \sum_{j=0}^p w_j x_j \right) \geq 0$$

sonst ist  $y(\mathbf{x}) = 0$

- wenn  $w_0 = -\theta$  und  $x_0 = 1$  (canonical form)
- sonst  $y(\mathbf{x}) = 1 \Leftrightarrow \left( \sum_{j=1}^p w_j x_j - \theta \right) \geq 0$
- $0 = w_0 + w_1 * x_1 + w_2 * x_2$  für 2D Fälle

#### 5.1.1 PT Algorithm

$PT(D, \eta)$

1. *initialize\_random\_weights*( $\mathbf{w}$ ),  $t = 0$
2. **REPEAT**
3.    $t = t + 1$
4.    $(\mathbf{x}, c(\mathbf{x})) = \text{random\_select}(D)$
5.    $\delta = c(\mathbf{x}) - y(\mathbf{x})$     //  $y(\mathbf{x}) \stackrel{(*)}{=} \text{heaviside}(\mathbf{w}^T \mathbf{x}) \in \{0, 1\}$ ,  $c(\mathbf{x}) \in \{0, 1\} \rightsquigarrow \delta \in \{0, 1, -1\}$
6.    $\Delta \mathbf{w} \stackrel{(*)}{=} \eta \cdot \delta \cdot \mathbf{x}$
7.    $\mathbf{w} = \mathbf{w} + \Delta \mathbf{w}$
8. **UNTIL** (*convergence*( $c(D), y(D)$ ) **OR**  $t > t_{\max}$ )
9. *return*( $\mathbf{w}$ )

- If a separating hyperplane between  $X_0$  and  $X_1$  exists, the PT algorithm will converge. If no such hyperplane exists, convergence cannot be guaranteed.
- A separating hyperplane can be found in polynomial time with linear programming. The PT Algorithm, however, may require an exponential number of iterations.
- Classification problems with noise are problematic

### 5.2 Gradient Descent

- Finde den kürzesten Weg in ein Min/Max über partielle Ableitungen
- The gradient of a function is the direction of steepest ascent or descent.
- in der VL ist ein Beweis den ich nicht abtippe weil irrelevant



### 5.2.1 Linear Regression + Squared Loss

$$L_2(\mathbf{w}) = \frac{1}{2} \cdot \sum_{(\mathbf{x}, c(\mathbf{x})) \in D} (c(\mathbf{x}) - y(\mathbf{x}))^2$$

Jetzt müssen wir für jedes  $w_i$  aus  $\mathbf{w}$  eine partielle Ableitung machen um den weight vector zu updaten ( $\mathbf{w} = \mathbf{w} + \Delta \mathbf{w}$ )

$$\Delta \mathbf{w} = \frac{\delta}{\delta w_i} L_2(\mathbf{w}) = \eta \cdot \sum_D (c(\mathbf{x}) - \mathbf{w}^T \mathbf{x}) \cdot \mathbf{x}$$

$\eta$  = learning rate - a small positiv constant - legen wir selbst fest

### 5.2.2 The Batch Gradient Descent (BGD) Algorithm

BGD( $D, \eta$ )

```
1. initialize_random_weights( $\mathbf{w}$ ),  $t = 0$ 
2. REPEAT
3.    $t = t + 1$ 
4.    $\Delta \mathbf{w} = 0$ 
5.   FOREACH  $(\mathbf{x}, c(\mathbf{x})) \in D$  DO
6.      $\delta = c(\mathbf{x}) - y(\mathbf{x})$  //  $y(\mathbf{x}) \stackrel{(*)}{=} \mathbf{w}^T \mathbf{x}$ ,  $\delta \in \mathbb{R}$ .
7.      $\Delta \mathbf{w} \stackrel{(*)}{=} \Delta \mathbf{w} + \eta \cdot \delta \cdot \mathbf{x}$  //  $\delta \cdot \mathbf{x}$  is the derivative of  $\ell_2(c(\mathbf{x}), y(\mathbf{x}))$ .
8.   ENDDO
9.    $\mathbf{w} = \mathbf{w} + \Delta \mathbf{w}$  //  $\Delta \mathbf{w}$  is  $-\eta \cdot \nabla L_2(\mathbf{w})$  here.
10. UNTIL (convergence( $c(D), y(D)$ ) OR  $t > t_{\max}$ )
11. return( $\mathbf{w}$ )
```

- wichtig: immer wenn irgendwo  $\mathbf{w}^T \mathbf{x}$  steht haben wir  $x_0 = 1$  zu  $\mathbf{x}$  hinzugefügt
- funktionsweise BGD. wir berechnen über die Ableitung in welche Richtung wir müssen und gehen dann einen Schritt der große  $\eta$
- die "convergence" schaut ob der Squared Loss noch größer als ein  $\varepsilon$  ist (das wir auch festlegen)
- BGD ist nicht der schnellste (bestenfalls linear) aber sehr einfach (Newton-Raphson algorithm, BFGS algorithm sind z.B. schneller)
- BGD nimmt den global loss: loss of all examples in  $D$  ("batch gradient descent") (Schritt in Richtung die für alle Punkte am besten ist)
- man kann auch den (squared) loss in Bezug auf einzelne Beispiele nehmen (pointwise loss) (dann gehts halt im Zickzack runter)  
berechnet sich dann  $\ell_2(c(\mathbf{x}), y(\mathbf{x})) = \frac{1}{2}(c(\mathbf{x}) - \mathbf{w}^T \mathbf{x})^2$

- bzw. die weight adaptation:  $\Delta \mathbf{w} = \eta \cdot (c(\mathbf{x}) - \mathbf{w}^T \mathbf{x}) \cdot \mathbf{x}$
- für  $BGD_\sigma$  wird Zeile 9 zu  
 $\mathbf{w} = \mathbf{w} + \Delta \mathbf{w} + \eta \cdot 2\lambda \cdot \begin{pmatrix} 0 \\ \vec{\mathbf{w}} \end{pmatrix}$

### 5.2.3 The Incremental Gradient Descent IGD Algorithm

IGD( $D, \eta$ )

```

1. initialize_random_weights( $\mathbf{w}$ ),  $t = 0$ 
2. REPEAT
3.    $t = t + 1$ 
4.   FOREACH  $(\mathbf{x}, c(\mathbf{x})) \in D$  DO
5.      $\delta = c(\mathbf{x}) - y(\mathbf{x})$  //  $y(\mathbf{x}) \stackrel{(*)}{=} \mathbf{w}^T \mathbf{x}$ ,  $\delta \in \mathbb{R}$ .
6.      $\Delta \mathbf{w} \stackrel{(*)}{=} \eta \cdot \delta \cdot \mathbf{x}$  //  $\delta \cdot \mathbf{x}$  is the derivative of  $\ell_2(c(\mathbf{x}), y(\mathbf{x}))$ .
7.      $\mathbf{w} = \mathbf{w} + \Delta \mathbf{w}$ 
8.   ENDDO
9. UNTIL (convergence( $c(D), y(D)$ )) OR  $t > t_{\max}$ 
10. return( $\mathbf{w}$ )

```

- kleinere Schritte als BGD
- can better avoid getting stuck in a local minimum of the loss function then BGD

### 5.2.4 Linear Regression + Squared Loss

$$L_{0/1}(\mathbf{w}) = \sum_D \frac{1}{2} \cdot (c(\mathbf{x}) - \text{sign}(\mathbf{w}^T \mathbf{x}))$$

$L_{0/1}(\mathbf{w})$  cannot be expressed as a differentiable function alias es kann nicht abgeleitet werden damit ist gradient descent nicht möglich

### 5.2.5 Logistic Regression + Logistic Loss + Regularization

Wie oben nur mit neuer Formel für  $\Delta \mathbf{w}$ :

$$\Delta \mathbf{w} = -\eta \cdot \nabla \mathcal{L}_\sigma(\mathbf{w}) = \eta \cdot \sum_D (c(\mathbf{x}) - \sigma(\mathbf{w}^T \mathbf{x})) \cdot \mathbf{x} - \eta \cdot 2\lambda \cdot \begin{pmatrix} 0 \\ \vec{\mathbf{w}} \end{pmatrix}$$

logistic loss Formel:

$$\mathcal{L}_\sigma(\mathbf{w}) = \sum_D -c(\mathbf{x}) \cdot \log(y(\mathbf{x})) - (1 - c(\mathbf{x})) \cdot \log(1 - y(\mathbf{x})) + \lambda \cdot \vec{\mathbf{w}}^T \vec{\mathbf{w}}$$

## 5.3 Multilayer Perceptron

### 5.3.1 Linear Separability

2 Klassen sind teilbar wenn ich da eine gerade Linie / Ebene / Hyperplane dazwischen packen kann.... oder:

Two sets of feature vectors,  $X_0, X_1$ , sampled from a  $p$ -dimensional feature space  $\mathbf{X}$ , are called linearly separable if  $p+1$  real numbers,  $w_0, w_1, \dots, w_p$ , exist such that the following conditions holds:

1.  $\forall \mathbf{x} \in X_0 : \sum_{j=0}^p w_j x_j < 0$
2.  $\forall \mathbf{x} \in X_1 : \sum_{j=0}^p w_j x_j \geq 0$

Problem: viele Probleme sind nicht linear separierbar Lösung: wir zeichnen mehrere Linien! (nehmen multilayer perceptron)

- The first, second, and third layer of the shown multilayer perceptron are called input, hidden, and output layer respectively
- input units perform no computation but only distribute the values to the next layer
- Compared to a single perceptron, the multilayer perceptron poses a significantly more challenging training (= learning) problem, which requires continuous (and non-linear) threshold functions along with sophisticated learning strategies.
- a continuous and non-linear threshold function:

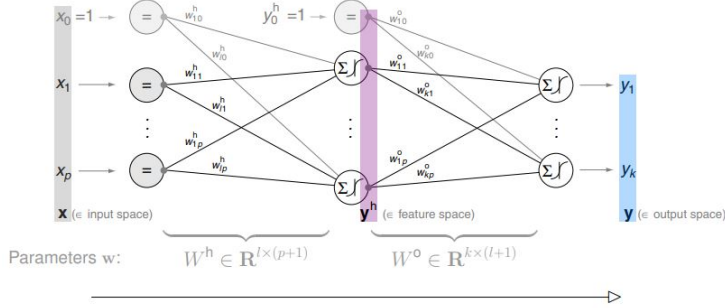
$$\sigma(z) = \frac{1}{1 + e^{-z}} \Rightarrow \frac{\delta\sigma(z)}{\delta z} = \sigma(z) \cdot (1 - \sigma(z))$$

- und damit:  $y(\mathbf{x}) = \sigma(\mathbf{w}^T \mathbf{x}) = \frac{1}{1 + e^{-\mathbf{w}^T \mathbf{x}}}$
- eine Alternative zu  $\sigma$  ist:

$$\tanh(z) = \frac{e^z - e^{-z}}{e^z + e^{-z}} = \frac{e^{2z} - 1}{e^{2z} + 1}$$

- A “multilayer” perceptron with linear threshold functions can be expressed as a single linear function and hence is equivalent to the power of a single perceptron only  $\Rightarrow$  Employing a nonlinear is necessary
- Multilayer perceptrons are also called multilayer networks or (artificial) neural network
-

A multilayer perceptron  $y(x)$  with one hidden layer and  $k$ -dimensional output layer:



Forward pass computation (aka. forward propagation) :

$$y(x) = \sigma(W^o y^h(x)) = \sigma\left(W^o \left(\sigma(W^h x)\right)\right)$$

### 5.3.2 Forward propagation

The input data is fed in the forward direction through the network. Each hidden layer accepts the input data, processes it as per the activation function and passes to the successive layer

(a) Propagate  $x$  from input to hidden layer: (IGD<sub>MLP</sub> algorithm, Line 5)

$$W^h \in \mathbb{R}^{l \times (p+1)} \quad x \in \mathbb{R}^{p+1}$$

$$\sigma \left( \begin{bmatrix} w_{10}^h & \dots & w_{1p}^h \\ \vdots & & \vdots \\ w_{l0}^h & \dots & w_{lp}^h \end{bmatrix} \begin{bmatrix} 1 \\ x_1 \\ \vdots \\ x_p \end{bmatrix} \right) = \begin{bmatrix} y_1^h \\ \vdots \\ y_l^h \end{bmatrix}$$

(b) Propagate  $y^h$  from hidden to output layer: (IGD<sub>MLP</sub> algorithm, Line 6)

$$W^o \in \mathbb{R}^{k \times (l+1)} \quad y^h \in \mathbb{R}^{l+1} \quad y \in \mathbb{R}^k$$

$$\sigma \left( \begin{bmatrix} w_{10}^o & \dots & w_{1l}^o \\ \vdots & & \vdots \\ w_{k0}^o & \dots & w_{kl}^o \end{bmatrix} \begin{bmatrix} 1 \\ y_1^h \\ \vdots \\ y_l^h \end{bmatrix} \right) = \begin{bmatrix} y_1 \\ \vdots \\ y_k \end{bmatrix}$$

propagation: Batch Mode

Forward

(a) Propagate  $\mathbf{x}$  from input to hidden layer: (IGD<sub>MLP</sub> algorithm, Line 5)

$$W^h \in \mathbf{R}^{l \times (p+1)} \quad D \subset \mathbf{R}^{p+1}$$

$$\sigma \left( \begin{bmatrix} w_{10}^h & \dots & w_{1p}^h \\ \vdots & & \vdots \\ w_{l0}^h & \dots & w_{lp}^h \end{bmatrix} \begin{bmatrix} 1 & \dots & 1 \\ x_{11} & \dots & x_{1n} \\ \vdots & & \vdots \\ x_{p1} & \dots & x_{pn} \end{bmatrix} \right) = \begin{bmatrix} y_{11}^h & \dots & y_{1n}^h \\ \vdots & & \vdots \\ y_{l1}^h & \dots & y_{ln}^h \end{bmatrix}$$

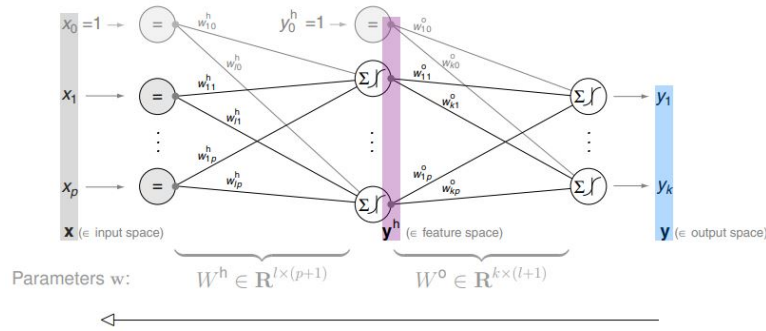
(b) Propagate  $\mathbf{y}^h$  from hidden to output layer: (IGD<sub>MLP</sub> algorithm, Line 6)

$$W^o \in \mathbf{R}^{k \times (l+1)}$$

$$\sigma \left( \begin{bmatrix} w_{10}^o & \dots & w_{1l}^o \\ \vdots & & \vdots \\ w_{k0}^o & \dots & w_{kl}^o \end{bmatrix} \begin{bmatrix} 1 & \dots & 1 \\ y_{11}^h & \dots & y_{1n}^h \\ \vdots & & \vdots \\ y_{l1}^h & \dots & y_{ln}^h \end{bmatrix} \right) = \begin{bmatrix} y_{11} & \dots & y_{1n} \\ \vdots & & \vdots \\ y_{k1} & \dots & y_{kn} \end{bmatrix}$$

### 5.3.3 Backwards Propagation

The considered multilayer perceptron  $\mathbf{y}(\mathbf{x})$ :



Weight update (aka. backward propagation) wrt. the global squared loss:

$$L_2(\mathbf{w}) = \frac{1}{2} \cdot \text{RSS}(\mathbf{w}) = \frac{1}{2} \cdot \sum_{(\mathbf{x}, \mathbf{c}(\mathbf{x})) \in D} \sum_{o=1}^k (c_o(\mathbf{x}) - y_o)^2$$

- $L_2(\mathbf{w})$  usually contains various local minima

$$\begin{bmatrix} \frac{\delta L_2(\mathbf{w})}{\delta w_{10}^o} & \dots & \frac{\delta L_2(\mathbf{w})}{\delta w_{1l}^o} \\ \vdots & & \vdots \\ \frac{\delta L_2(\mathbf{w})}{\delta w_{k0}^o} & \dots & \frac{\delta L_2(\mathbf{w})}{\delta w_{kl}^o} \end{bmatrix} \equiv \nabla^o L_2(\mathbf{w})$$

die selbe Formel gibt es nochmal nur mit "h" statt "o"

Update of weight matrix  $W^o$ :  $W^o = W^o + \delta W^o$  mit:

$$\Delta W^o = -\eta \cdot \nabla^o L_2(\mathbf{w}) = \eta \cdot \sum_D [(W^{oT} \delta^o) \odot y^h(\mathbf{x}) \odot (1 - y^h(\mathbf{x}))]_{1, \dots, l} \otimes \mathbf{x}$$

weiß irgendwer was diese lange Formel sagen will? Nein? okay.

## 5.4 IGD for Multilayer Perceptrons

```

1. initialize_random_weights( $W^h, W^o$ ),  $t = 0$ 
2. REPEAT
3.    $t = t + 1$ 
4.   FOREACH  $(\mathbf{x}, \mathbf{c}(\mathbf{x})) \in D$  DO
5.      $\mathbf{y}^h = (\sigma_{(W^h \mathbf{x})}^1)$  // forward propagation,  $\mathbf{x}$  is extended by  $x_0 = 1$ 
6.      $\mathbf{y} = \sigma(W^o \mathbf{y}^h)$ 
7.      $\delta^o = (\mathbf{c}(\mathbf{x}) - \mathbf{y}) \odot \mathbf{y} \odot (\mathbf{1} - \mathbf{y})$  // backward propagation
8.      $\Delta W^o = \eta \cdot (\delta^o \otimes \mathbf{y}^h|_{1,...,d})$ 
9.      $\delta^h = [(W^{oT} \delta^o) \odot \mathbf{y}^h \odot (\mathbf{1} - \mathbf{y}^h)]_{1,...,d}$ 
10.     $\Delta W^h = \eta \cdot (\delta^h \otimes \mathbf{x})$ 
11.     $W^o = W^o + \Delta W^o$  // weight update
12.     $W^h = W^h + \Delta W^h$ 
13.  ENDDO
14. UNTIL(convergence( $\mathbf{c}(D), \mathbf{y}(D)$ )) OR  $t > t_{\max}$ 
15. return( $W^h, W^o$ )

```

- $\odot$  = Hadamard product
- $\otimes$  =  $\mathbf{vw}^T$  dyadic product

another formel for  $\nabla^o L_2(\mathbf{w})$ :

$$\frac{\delta}{\delta w_{ij}^o} = - \sum_D (c_i(\mathbf{x}) - y_i(\mathbf{x})) \cdot y_i(\mathbf{x}) \cdot (1 - y_i(\mathbf{x})) \cdot y_j^h(\mathbf{x})$$

## 6 Decision Trees

### 6.1 Splitting

Das splitting von  $X$  ist die Teilung in mutually exclusive Teilmengen  $X_1, \dots, X_s$  alias  $X = X_1 \cup \dots \cup X_s$  mit  $X_j \neq \emptyset$  and  $X_j \cap X_{j'} \neq \emptyset$ , where  $j, j' \in \{1, \dots, s\}, j \neq j'$ . Eine Teilung  $X_1, \dots, X_s$  von  $X$  induces eine Teilung  $D_1, \dots, D_s$  von  $D$ , bei der  $D_j, j = 1, \dots, s$ , als  $\{(\mathbf{x}, c(\mathbf{x})) \in D | x \in X_j\}$  definiert wird.

Oder in Beispielen: Wir teilen unsere Pilze in 2 (oder mehr) Haufen anhand einer Eigenschaft, z.B. Farbe. Statt also einem Haufen beschrifteter Pilze (unser Beispieldatenset  $D$ ) haben wir 3 Haufen: weiße Pilze, braune Pilze und rote Pilze. (die wiederum weiter aufgeteilt werden bis wir für jeden Haufen klar sagen können: essbar oder nicht essbar) Die Teilung hängt von der Art der Eigenschaften/ Features ab. Diese können

- numerisch/ quantitative sein: Zahlen und Ratios, z.B. die Temperatur in °C, das Alter in Jahren oder Geld in Währung (Rechnen mit + & - und/oder \* & / ergibt Sinn)

- kategorisch/ Qualitative sein: Namen, Farben, IDs, Postleitzahlen, Hausnummern, etc. (nur Vergleichen [=, ≠ bzw. manchmal noch < & >] ergibt Sinn)

Daraus folgen verschiedene Teilmöglichkeiten:

- $m$ -ary splitting induced by a (nominal) feature  $A$  with finite domain (braune, weiße und rote Pilze):

$$A = \{a_1, \dots, a_m\} : X = \{\mathbf{x} \in X : \mathbf{x}|_a = a_1\} \cup \dots \cup \{\mathbf{x} \in X : \mathbf{x}|_a = a_m\}$$

- Binary splitting induced by a (nominal) feature  $A$  (Punkte vs keine Punkte):

$$A' \subset A : X = \{\mathbf{x} \in X : \mathbf{x}|_a \in A'\} \cup \{\mathbf{x} \in X : \mathbf{x}|_a \notin A'\}$$

- Binary splitting induced by an ordinal feature  $A$  (Temp <25 °C und Temp ≥25°C):

$$v \in \text{dom}(A) : X = \{\mathbf{x} \in X : \mathbf{x}|_a \succeq v\} \cup \{\mathbf{x} \in X : \mathbf{x}|_a \prec v\}$$

- $\mathbf{x}|_{a_i}$  sind die  $\mathbf{x}$  die eine eigenschaft  $a_i$  erfüllen (z.B. braun)

### 6.1.1 Possible criteria for splitting of $X(t)$

1. **Size of  $D(t)$ :**  $D(t)$  is not split if  $|D(t)|$  is below a threshold.
2. **Purity of  $D(t)$ :**  $D(t)$  is not split if all examples in  $D(t)$  are members of the same class
3. **Impurity reduction of  $D(t)$ :**  $D(t)$  is not split if its impurity reduction,  $\Delta_i$ , is below a threshold (wenn es sich nicht lohnt weil es sich nicht signifikant verbessert)

## 6.2 Decision Trees Definiton

Let  $X$  be a set of features and  $C$  a set of classes. A decision tree  $T$  for  $X$  and  $C$  is a finite tree with a distinguished root node. A non-leaf node  $t$  of  $T$  has assigned

- 1 a set  $X(t) \subset X$ ,
- 2 a splitting of  $X(t)$ , and
- 3 a one-to-one mapping of the subsets of the splitting to its successors.

$X(t) = X$  iff  $t$  is root node. A leaf node of  $T$  has assigned a class from  $C$ .

- decision trees (DT) haben kein Problem mit Rauschen/ Labelnoise (falsch klassifizierte Beispiele) da wenn diese in der Minderheit sind nicht berücksichtigt werden (wenn ich 5 mal bei sonnigem, warmen, windstillen Wetter Sport mache und 1 mal nicht) haben wir ein geteilten Knoten (siehe Bild) und entscheiden und meistens für die Mehrheit (Sport machen).
- manchmal (essbar/ nicht essbar) ergibt jedoch auch eine andere Entscheidung Sinn
- DT können nur klassifizieren nicht Regression

### 6.2.1 classification of some $x \in X$ given a decision tree $T$

1. Find the root node  $t$  of  $T$
2. if  $t$  is a non-leaf node, find among its successors that node  $t'$  whose subset of the splitting of  $X(t)$  contains  $\mathbf{x}$ . Repeat this step with  $t = t'$
3. If  $t$  is a leaf node, label  $x$  with the respective class

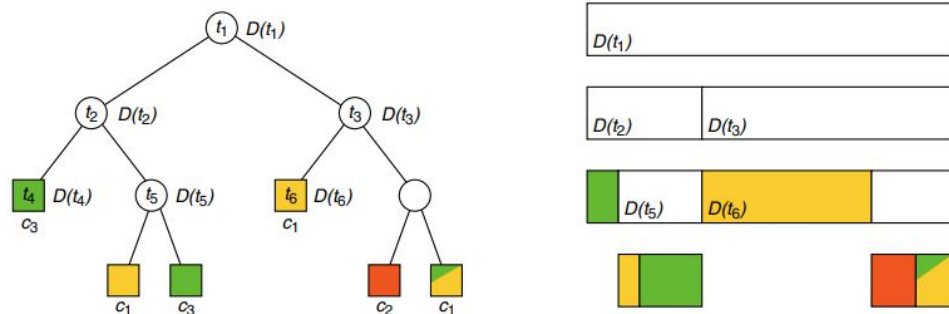
The set of possible decision trees forms the hypothesis space  $H$

- The classification of an  $x \in X$  determines a unique path from the root node of  $T$  to some leaf node of  $T$
- At each non-leaf node a particular feature of  $x$  is evaluated in order to find the next node along with a possible next feature to be analyzed

### 6.2.2 Notations

Let  $T$  be a decision tree for  $X$  and  $C$ , let  $D$  be a set of examples [setting], and let  $t$  be a node of  $T$ . Then we agree on the following notation:

- $X(t)$  denotes the subset of  $X$  that is represented by  $t$ .
- $D(t)$  denotes the subset of the example set  $D$  that is represented by  $t$ , where  $D(t) = (\mathbf{x}, c(\mathbf{x})) \in D | \mathbf{x} \in X(t)$ . (see the splitting definition)





- The set  $X(t)$  is comprised of those members  $x$  of  $X$  that are filtered by a path from the root node of  $T$  to the node  $t$ .
- $leaves(T)$  denotes the set of all leaf nodes of  $T$
- ab hier:  $t, T : X \rightarrow C$  statt  $y_t, y_T : X \rightarrow C$

### 6.2.3 DT-construct

*DT-construct*( $D$ )

```

1.  $t = newNode()$ 
    $label(t) = representativeClass(D)$ 

2. IF  $impure(D)$ 
   THEN  $criterion = splitCriterion(D)$ 
   ELSE  $return(t)$ 

3.  $\{D_1, \dots, D_s\} = decompose(D, criterion)$ 

4. FOREACH  $D'$  IN  $\{D_1, \dots, D_s\}$  DO
    $addSuccessor(t, DT-construct(D'))$ 
ENDDO

5.  $return(t)$ 

```

- Since *DT-construct* assigns to each node of a decision tree  $T$  a class, each subtree of  $T$  (as well as each pruned version of a subtree of  $T$ ) represents a valid decision tree on its own
- $representativeClass(D)$  : Returns a representative class for the example set  $D$ . Note that, due to pruning, each node may become a leaf node.
- $impure(D)$  Evaluates the (im)purity of a set  $D$  of examples
- $splitCriterion(D)$  Returns a split criterion for  $X(t)$  based on the examples in  $D(t)$ .
- $decompose(D, criterion)$  Returns a splitting of  $D$  according to criterion

- *addSuccessor*( $t, t'$ ) Inserts the successor  $t'$  for node  $t$

#### 6.2.4 DT-classify

*DT-classify*( $x, t$ )

```
1. IF isLeafNode( $t$ )
   THEN return( $label(t)$ )
   ELSE return(DT-classify( $x, splitSuccessor(t, x)$ ))
```

- *isLeafNode*( $t$ ) Tests whether  $t$  is a leaf node
- *splitSuccessor*( $t, x$ ) Returns the (unique) successor  $t'$  of  $t$  for which  $x \in X(t')$  holds

### 6.3 Evaluation of DT

Um noch Bias zu haben (und somit für nicht gesehene Beispiele abstrahieren zu können) müssen wir den Baum jetzt noch kürzen.  
(Overfitting vermeiden)

1. Size: Among all decision trees of minimum classification error we choose the one of smallest size
2. Classification Error: Quantifies the rigor according to which a class label is assigned to  $x$  in a leaf node of  $T$ , based on the examples in  $D$ . If all leaf nodes of a decision tree  $T$  represent a single example of  $D$ , the classification error of  $T$  with respect to  $D$  is zero. (Wollen wir meistens nicht)

#### 6.3.1 measuring Size

- Leaf node number (wie viele Endknoten)
- Tree height (Längster Weg im Baum von Wurzel bis Blatt)
- External path length (Summe aller Wege von Wurzel bis Blätter = der Platz den es braucht alle in einem Baum gespeicherten Regeln zu speichern)
- Weighted external path length (wie external path length aber wir schauen für jeden Weg noch wie viele Beispiele aus  $D$  klassifiziert werden) → Ein kurzer Weg der viel klassifiziert + ein sehr langer weg der sehr wenig klassifiziert kann besser sein als 2 Halblange Wege (die beliebig aufteilend klassifizieren)
- Weighted external path length ist meistens bevorzugt

### 6.3.2 Classification Error

$$label(t) = \operatorname{argmax}_{c \in C} \frac{|\{\mathbf{x}, c(\mathbf{x}) \in D(t) : c(\mathbf{x}) = c\}|}{|D(t)|}$$

**Missclassification rate**/ Error of node classifier t based on that:

$$Err(t, D(t)) = \frac{|\{\mathbf{x}, c(\mathbf{x}) \in D(t) : c(\mathbf{x}) \neq label(t)\}|}{|D(t)|} = 1 - \max_{c \in C} \frac{|\{\mathbf{x}, c(\mathbf{x}) \in D(t) : c(\mathbf{x}) = c\}|}{|D(t)|}$$

**Misclassification rate** of decision tree classifier T:

$$Err(T, D) = \sum_{t \in leaves(T)} \frac{|D(t)|}{|D|} \cdot Err(t, D(t))$$

argmax gibt das  $\mathbf{x}$  zurück für das eine Funktion  $f(\mathbf{x})$  maximal ist, max gibt  $y(\mathbf{x})$  zurück (Beispiel Wohnungspreise anhand der Lage:  $argmax(\mathbf{x}) = \text{"Innenstadt"}$  und  $max(\mathbf{x}) = \text{"75000 Euro"}$ ) (und genauso für argmin vs min)

**Misclassification cost** of node classifier t:

$$Err_{cost}(t, D(t)) = \min_{c' \in C} \sum_{c \in C} \frac{|\{\mathbf{x}, c(\mathbf{x}) \in D(t) : c(\mathbf{x}) = c\}|}{|D(t)|} \cdot cost(c'|c)$$

**Misclassification costs** of decision tree classifier T:

$$Err_{cost}(T, D) = \sum_{t \in leaves(T)} \frac{|D(t)|}{|D|} \cdot Err_{cost}(t, D(t))$$

Wenn  $c' = c$  ist  $cost(c'|c) = 0$ , verschiedene c können verschiedene Kosten haben (es ist schlimmer einen toxischen Pilz als essbar einzustufen als einen essbaren als toxisch)

## 6.4 Impurity Functions

An impurity function  $\iota : [0; 1]^k \rightarrow \mathbb{R}$  (mit  $k \in \mathbb{N}$ ) is a function defined on the standard  $k - 1 - simplex$ , denoted  $\Delta^{k-1}$ , for which the following properties hold:

1.  $\iota$  becomes minimum at points  $(1, 0, \dots, 0), (0, 1, \dots, 0), \dots, (0, \dots, 0, 1)$
2.  $\iota$  is symmetric with regard to its arguments,  $p_1, \dots, p_k$
3.  $\iota$  becomes maximum at point  $(1/k, \dots, 1/k)$

simplex = einfachste Form zwischen Punkten in k- Dim raum (also  $k = 3 \rightarrow \text{Dreieck}$ )?

### 6.4.1 Strict Impurity Function

genauso wie Impurity aber

3 wird zu  $\iota(\lambda p + (1 - \lambda)p') > \lambda\iota(p) + (1 - \lambda)\iota(p'), 0 < \lambda < 1, p \neq p'$

### 6.4.2 Impurity of an Example Set

Impurity of  $D = \iota(D)$  :

$$\iota(D) = \iota\left(\frac{|\{(\mathbf{x}, c(\mathbf{x})) \in D : c(\mathbf{x}) = c_1\}|}{|D|}, \dots, \frac{|\{(\mathbf{x}, c(\mathbf{x})) \in D : c(\mathbf{x}) = c_k\}|}{|D|}\right)$$

eine Funktion ist maximal impure wenn alle ihr zugeordneten Klassen  $c$  gleich wahrscheinlich sind (50% Wahrscheinlichkeit essbar, 50% Wahrscheinlichkeit toxisch).

### 6.4.3 Impurity Reduction

D aufgeteilt in  $D_1, \dots, D_s$ , durch das splitting von X, dann ist die resulting impurity reduction:

$$\Delta\iota(D, \{D_1, \dots, D_s\}) = \iota(D) - \sum_{j=1}^s \frac{|D_j|}{|D|} \cdot \iota(D_j)$$

### 6.4.4 Impurity Functions Based on the misclassification Rate

Definition für 2 classes:

$$\iota_{misclass}(p_1, p_2) = 1 - \max\{p_1, p_2\} = \begin{cases} p_1 & \text{if } 0 \leq p_1 \leq 0.5 \\ 1 - p_1 & \text{otherwise} \end{cases}$$

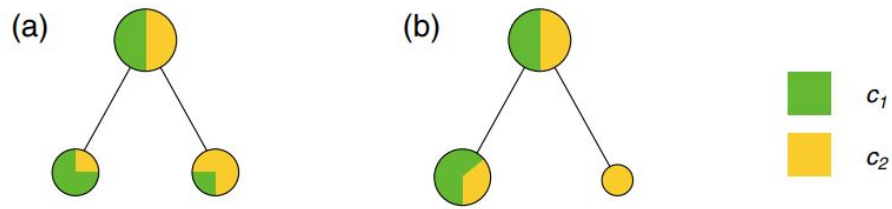
$$\iota_{misclass}(D) = 1 - \max\left\{\frac{|\{(\mathbf{x}, c(\mathbf{x})) \in D : c(\mathbf{x}) = c_1\}|}{|D|}, \frac{|\{(\mathbf{x}, c(\mathbf{x})) \in D : c(\mathbf{x}) = c_2\}|}{|D|}\right\}$$

Definition für k classes:

$$\iota_{misclass}(p_1, \dots, p_k) = 1 - \max_{i=1, \dots, k} p_i$$

$$\iota_{misclass}(D) = 1 - \max_{c \in C} \frac{|\{(\mathbf{x}, c(\mathbf{x})) \in D : c(\mathbf{x}) = c\}|}{|D|}$$

Er hat in dem Video erklärt warum er hier plötzlich  $p_1$  und  $p_2$  benutzt statt  $c_1$  und  $c_2$ , hat allerdings nicht so viel Sinn ergeben.



$$\Delta \iota_{misclass} = \iota_{misclass}(D) - \left( \frac{|D_1|}{|D|} \cdot \iota_{misclass}(D_1) + \frac{|D_2|}{|D|} \cdot \iota_{misclass}(D_2) \right)$$

left splitting:  $\Delta \iota_{misclass} = \frac{1}{2} - \left( \frac{1}{2} \cdot \frac{1}{4} + \frac{1}{2} \cdot \frac{1}{4} \right) = \frac{1}{4}$

right splitting:  $\Delta \iota_{misclass} = \frac{1}{2} - \left( \frac{3}{4} \cdot \frac{1}{3} + \frac{1}{4} \cdot 0 \right) = \frac{1}{4}$

Das Bild erklärt das Rechnen besser als alle Formeln.