

## Wörterliste ML Datenaufbereitung

<b>Instanzen</b>	<ul style="list-style-type: none"> <li>- Einzelne Datenpunkte oder «Fälle», die das System lernen soll – z.B. ein Kunde bei TeleFlow</li> <li>- Herausforderungen: Daten über Instanzen kommen oft aus mehreren Quellen und müssen mittels Joins zusammengeführt werden.</li> </ul>
<b>Klassenattribut</b>	<ul style="list-style-type: none"> <li>- Das Ziel, das vorhergesagt werden soll (z.B. «Kunde kündigt Vertrag: Ja oder Nein»)</li> <li>- Herausforderungen: Muss teilweise manuell erstellt oder berechnet werden <ul style="list-style-type: none"> <li>o Manuell: z.B. bei der Frage, ob eine E-Mail «interessant» ist.</li> <li>o Komplex: z.B. ob ein Kunde kündigt → aus seinem Verhalten abgeleitet</li> </ul> </li> </ul>
<b>Attribute</b>	<ul style="list-style-type: none"> <li>- Merkmale/Beschreibungen einer Instanz (z.B. Wohnort, Anzahl, Tickets, letzte Reparatur)</li> <li>- Typische Aufgaben: <ul style="list-style-type: none"> <li>o Transformation: z.B. Datum in «Saison» umwandeln</li> <li>o Feature Creation: neue Attribute aus bestehenden erstellen.</li> <li>o Aggregation: historische Daten zusammenfassen.</li> <li>o N:m-Problem: Einer Instanz können mehrere Werte eines Attributs zugeordnet sein</li> <li>o Feature Selection: irrelevante Attribute entfernen (z.B: IDs, Datumsangaben)</li> </ul> </li> </ul>
<b>Fehlende Werte</b>	<p>Umgangsmöglichkeiten:</p> <ul style="list-style-type: none"> <li>- Ignorieren: Sonderwert einfügen</li> <li>- Löschen 1: Instanzen mit fehlenden Attributwerten löschen(Zeile)</li> <li>- Löschen 2: Attribut mit fehlenden Werten aus allen Datenobjekten entfernen (Spalte)</li> <li>- Berechnen(«impute») 1: Mittelwert/Median/Modus einfügen</li> <li>- Berechnen(«impute») 2: Wert aus ähnlichster Instanz des Datensatzes kopieren</li> </ul>
<b>Cleaning (Datenbereinigung)</b>	<ul style="list-style-type: none"> <li>- Beseitigung von: <ul style="list-style-type: none"> <li>o Duplikaten</li> <li>o Ausreisern</li> <li>o Inkonsistenzen</li> </ul> </li> </ul>
<b>Typkonversion</b>	<ul style="list-style-type: none"> <li>- Einfache Konversion: z.B. Monat «3» als Kategorie behandeln, nicht als Zahl.</li> </ul>

	<ul style="list-style-type: none"> <li>- Diskretisierung: numerische Werte in Kategorien einteilen (z.B. «jung», «alt»)</li> <li>- One-Hot-Encoding: Kategorien in Zahlenvektoren umwandeln:</li> </ul> <p>Vor One-Hot-Encoding:</p> <table border="1" data-bbox="695 422 1049 676"> <thead> <tr> <th>Nahrungsmittel</th><th>Kalorien</th></tr> </thead> <tbody> <tr> <td>Apple</td><td>95</td></tr> <tr> <td>Chicken</td><td>231</td></tr> <tr> <td>Broccoli</td><td>50</td></tr> </tbody> </table> <p>Nach One-Hot-Encoding:</p> <table border="1" data-bbox="695 759 1208 968"> <thead> <tr> <th>Apple</th><th>Chicken</th><th>Broccoli</th><th>Calories</th></tr> </thead> <tbody> <tr> <td>1</td><td>0</td><td>0</td><td>95</td></tr> <tr> <td>0</td><td>1</td><td>0</td><td>231</td></tr> <tr> <td>0</td><td>0</td><td>1</td><td>50</td></tr> </tbody> </table>	Nahrungsmittel	Kalorien	Apple	95	Chicken	231	Broccoli	50	Apple	Chicken	Broccoli	Calories	1	0	0	95	0	1	0	231	0	0	1	50
Nahrungsmittel	Kalorien																								
Apple	95																								
Chicken	231																								
Broccoli	50																								
Apple	Chicken	Broccoli	Calories																						
1	0	0	95																						
0	1	0	231																						
0	0	1	50																						
<b>Normalisierung</b>	<ul style="list-style-type: none"> <li>- Alle Werte in einen Bereich zwischen 0 und 1 bringen</li> <li>- Formel:</li> </ul> $x_{neu} = \frac{x_{alt} - min}{max - min}$																								
<b>Skalierung</b>	<ul style="list-style-type: none"> <li>- Alle Werte so anpassen, dass sie: <ul style="list-style-type: none"> <li>○ einen Mittelwert von 0 haben.</li> <li>○ und eine Standardabweichung von 1.</li> </ul> </li> <li>- Formel:</li> </ul> $x_{neu} = \frac{x_{alt} - \bar{x}}{\sigma}$																								