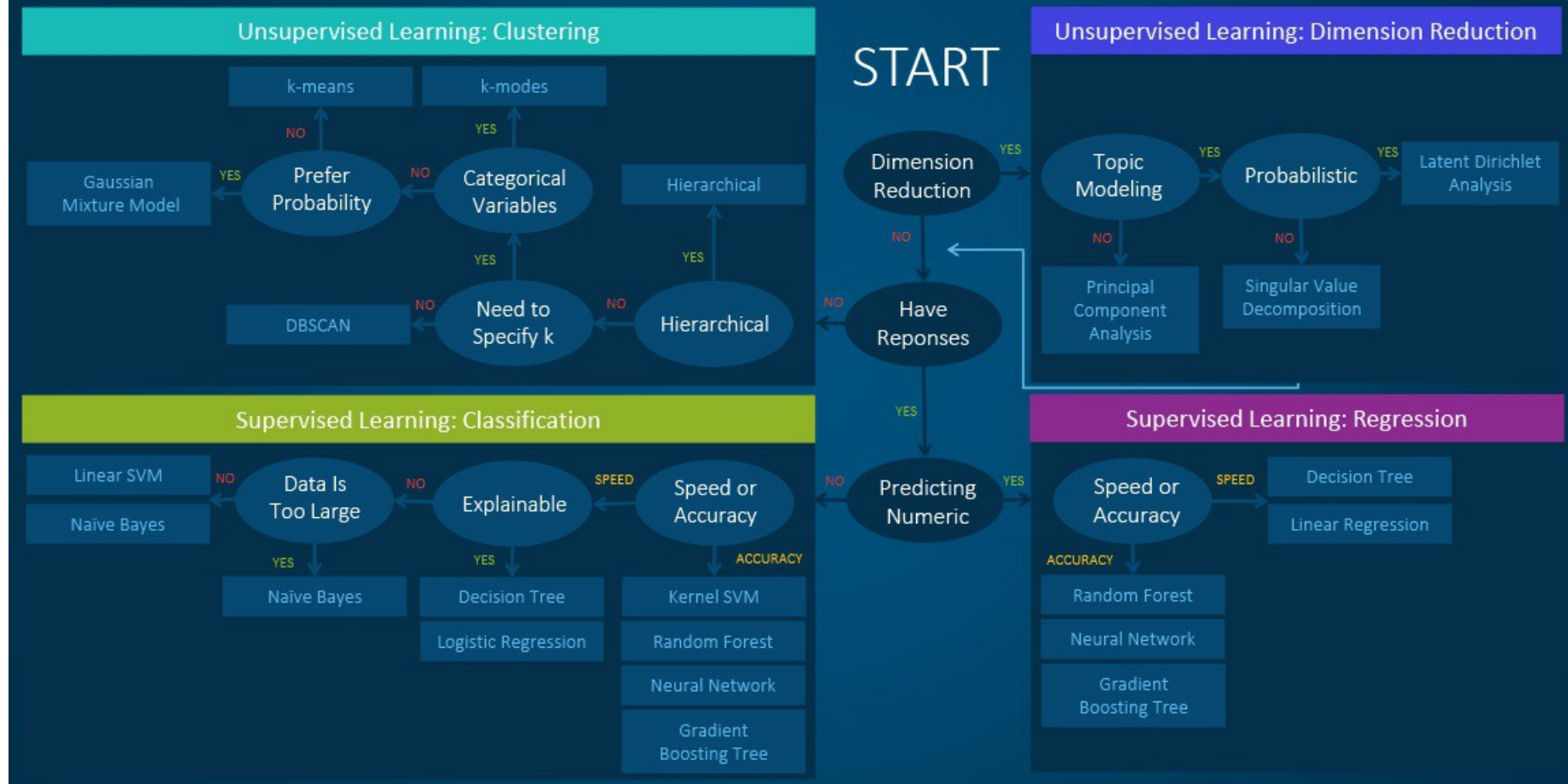


Machine Learning Algorithmen

Machine Learning Algorithms Cheat Sheet



DATA SCIENCE

MACHINE LEARNING		
Tag	Beschreibung	Beispiel
<ul style="list-style-type: none"> CRISP-DM (immer im Kreislauf von den Arbeitsschritten von 1.-8.) cross-industry-standard-process-for-data-mining 	<ul style="list-style-type: none"> <u>Struktur von Data Science Projekten:</u> 	<ol style="list-style-type: none"> 1. Fragestellung erkennen 2. Business Understanding (Recherche in den Bereichen/Branchen): einzelne Problemstellungen 3. Data Understanding (Daten beziehen von...) 4. Austausch: Business Understanding - Data Understanding 5. Data Preparation (Daten vorbereiten, bei Defiziten z.B. Mittelwerte) 6. Modelling (auch in Rücksprache mit Data Preparation) 7. Evaluation (Daten werden ausgewertet, hohe Wahrscheinlichkeit für die Aussagekraft des Modells, z.B. 99%) auch in Rücksprache mit Business Understanding 8. Deployment (Einsatz/Ausführen des Modells) auch in Rücksprache mit Business Understanding
<ul style="list-style-type: none"> Datensatz mit: Features(Spalten, 1-3 oder mehr) 1 Label (z.B. Spezies bei Tieren) viele Samples (Zeilen, z.B. 1000) 	<ul style="list-style-type: none"> Wie Algorithmen lernen Bsp.: Welche Spezies ist der Pinguin? Wichtig: gute Verteilung von Daten, z.B. gleichviele Daten von den verschiedenen 	<ul style="list-style-type: none"> Man geht von einem Datensatz aus (.csv Datei) ca. 70 – 98% davon wird als Trainingsdatensatz -Training Data- verwendet: --->Struktur finden, trainiertes Modell entwickelt aus den Features (Train_Test_Split (Datensatz splitten),Fitting(Datensatz um zu trainieren)) ca. 30 – 2% Testdatensatz (=Rest vom Trainingsdatensatz) -Validation Data-: --->Struktur testen (Validierungs-Datensatz) nach und nach Samples (Features ohne Labels) aus dem

	Labels, um eine gute Vorhersage zu generieren	<p>Testdatensatz werden im trainierten Modell eingesetzt ---> Prediction (Vorhersage), also eines der Labels</p> <ul style="list-style-type: none"> nach und nach Samples (Labels) aus dem Testdatensatz wird im trainierten Modell -Test Data- eingesetzt----> korrektes Label <ul style="list-style-type: none"> ---> Fehler/Treffer Tabelle: <ul style="list-style-type: none"> Prediction = <u>korrektes</u> Label (Ja/Nein)
<ul style="list-style-type: none"> Accuracy Confusion-Matrix 	<ul style="list-style-type: none"> Evaluieren 	<ul style="list-style-type: none"> Trefferquote über richtige Vorhersage trefferanzahl / Gesamtanzahl (=accuracy) Abgleichen mit realen Werten (actual values) <ul style="list-style-type: none"> Möglichkeiten: richtig (true) als Wert/Kategorie (Value/Klasse) und falsch(false) als Wert/Kategorie
<ul style="list-style-type: none"> Zyklus zwischen Daten-Produkt-Kunden-Daten 	<ul style="list-style-type: none"> Selbst lernendes Produkt wenn alle 3 Elemente steigen (Produkt-Kunden-Daten) profitieren alle untereinander 	<ul style="list-style-type: none"> Auto: selbstfahrend Daten (wenn nicht vorhanden, dann manuelle verwenden oder z.B captcha[als Nutzer die Ampel in 9 Bildern finden,Text erkennen]) Kunden (bezahlen für deren Daten zu sammeln) Produkt (wenn nicht vorhanden, dann als gegeben betrachten, aber manuell bedienen)
<ul style="list-style-type: none"> Fundamentale Arbeit in der Data Science 	<ul style="list-style-type: none"> Fragestellungen 	<ul style="list-style-type: none"> 3 Typen: <ol style="list-style-type: none"> Descriptive (klare, beschreibende Fragen) Exploratory (offene, untersuchende Fragen) Prediction (Vorhersagen, klar, aber auch offen)
<ul style="list-style-type: none"> Algorithmen werden mit Labels trainiert 	<ul style="list-style-type: none"> Supervised Learning <p>Voraussetzung: historische Daten mit hohen Warscheinlichkeiten für Zukunftsvorhersage</p>	<ul style="list-style-type: none"> Input, bei dem der Output klar ist (Output ist das Label, z.B. Hund oder Katze als Vorhersage, Spam: Ja/Nein) Algorithmus lernt, indem er das Ergebnis mit dem Label vergleicht ---->Fehler wird minimiert, Modell wird modifiziert

	<ul style="list-style-type: none"> • <u>Regression</u> • <u>Klassifizierung</u> 	<ul style="list-style-type: none"> • Vorhersage von kontinuierlichen Werten • Aussage JA / NEIN (Treffer oder nicht Treffer)
<ul style="list-style-type: none"> • 	<ul style="list-style-type: none"> • Unsupervised Learning 	<ul style="list-style-type: none"> • Keine vorhandenen Labels <ul style="list-style-type: none"> ◦ <u>exploratory data analysis</u> z.B. Geschäfte einer Stadt ohne Kenntnis über die Kategorie - Bar,Cafe,Friseur,etc. ◦ <u>Markt segmentation</u> (innerhalb eines business,z.B. Versicherungsbereiche) ◦ <u>social network</u> ◦ <u>clustering businesses in a given area</u> •
<ul style="list-style-type: none"> • <u>Overfitting</u>:(zu viele passende Trainingsdaten vorhanden) • <u>Underfitting</u> (zu wenig passende Trainingsdaten vorhanden) 	<ul style="list-style-type: none"> • Typische Probleme bei Machine Learning 	<ul style="list-style-type: none"> ◦ Nahezu perfekte Vorhersage (100%), Fehler bei Trainingsdaten nach validierung immer gering und bei einsatz in Testdaten groß ◦ ----> für Vorhersage eher z.B. eher über eine Gerade im einem Diagramm die Trainingsdaten interpretieren als über jeden einzelnen vorhandenen Datenpunkt (Sample)