

Wahrscheinlichkeitsverteilung (Probability Distribution)

f: Ereignis \rightarrow "Häufig"
Häufigkeit

- diskret
- kontinuierlich

Galton Board / Galton Brett

\sim Binomial Distribution

\sim Bernoulli Verteilung

Log-normal distribution

Bayes..

Daten x, Ziely

Bedingte Wahrscheinlichkeit \sim Conditional probability

$$P(A|B) = \frac{P(A \cap B)}{P(B)}$$

A given B

A oder B, d
durch B

$$P(A|B) * P(B) = P(A \cap B)$$

$$P(A|B) * P(B) = P(B|A) * P(A) \quad \text{Durchschnitt}$$

naive Bayes \rightarrow Bayes dengan asumsi data independen

\rightarrow harus cek korelasi dulu & independen test

--

Kategorial verteilung

6199

Dienstag, 17.10.2023

Encoding \rightarrow Hot
 Ordinal

Bayesian - Klassifikator für nominale Merkmale (Kapitel 4.3.2) } Buch
Seite 88 } Fräcke

$$P(i, x) = \frac{\prod_{h=1}^n P(x^{(h)} | i) \cdot P(i)}{\sum_{j=1}^M P(j) \prod_{h=1}^n P(x^{(h)} | j)}$$

Mittwoch, 18.10.2023

Linear Regression
fit in sklearn

$$y = \text{intercept}_- + \text{coef}_- * x$$

$$y = b + mx$$

Newton Gauss Verfahren

$$x_{n+1} = x_n - \frac{f(x_n)}{f'(x_n)}$$

Nacheinander ohne Ableite

Skalierung / Scaling

Linear Regression:

- ~~Skalieren~~ Nicht skalieren \rightarrow wenn nur das predict interessiert
- Skalieren \rightarrow wenn die Wichtigkeit der Spalten interessiert

Lasso, Ridge

\rightarrow Skalieren, weil es sonst nicht konvergiert.

Bayes

\rightarrow nicht skalieren

$O(n) \rightarrow$ Big O \leftarrow Aufwand von Algorithmus

Vector n Komponenten \rightarrow addieren: ~~$O(n)$~~ $n \sim O(n)$

Matrix n Zahlen, 2 Spalten, Summierung der Spalte

$O(n)$

Matrix n Zeilen ~~n~~ f Spalten

$O(n \cdot f)$

Bayes: $f \cdot O(n \cdot f)$
produkt $O(f)$

Linear Algebra: $O(n \cdot f^2)$

Nächste Nachbarn (kNN)

↳ Label

↳ Eigenschaften

↳ Gewichtung der Eigenschaften

↳ Gruppierung nach Ähnlichkeit

↳ totale Übereinstimmung

↳ relative kurzer Abstand

KD Tree

|

Rechts

Ball Tree

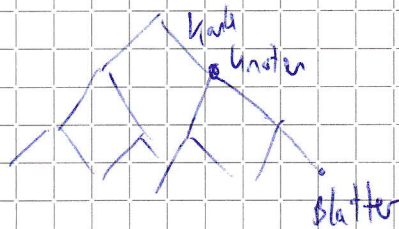
~~Decision Tree~~ kNN (nearest neighbor)

Distance:

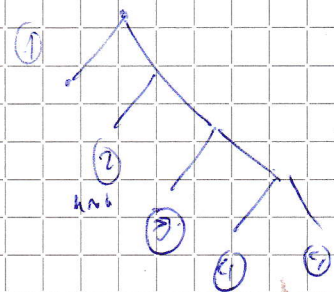
- Euklidisch (L2)
- Haversine
- Manhattan, cityblock (L1)
- Cosine

Decision Tree:

$O(\log(n))$



Pfadlänge



- Anzahl Blätter: weniger ist besser

↳ Anzahl der Regeln

- Pfadlänge: kurz ist besser

↳ Länge der ~~längsten~~ längsten Regeln (worst case)

- Pfadlängenscore: weniger ist besser

↳ durchschnittl. Regellänge

Entropy