

Monday, 09.10.2023

KLR - 335 Machine Learning

Jupyter Notebook
PyCharm



Wednesday, 10.10.2023

Machine Learning ← still slide 1

Projektarbeit...

Box Filter / OpenAI Hide & Seek ← youtube

Emergent Tool Use from Multi-Agent Interaction

Vorlesungsteil

Testdaten

{

New daten

: neu trainieren

Test Menge

20%

{ Training Menge

80%

? Pareto

Test Menge

Trainingsmenge

Numpy, pandas, matplotlib

Scipy, scikit-learn

Test-set, Train-set

- Stratified Sampling

train-test-split (..)

stratify = housinggen[["income_cat", ... "median
price_cat"]]

--



hast to be category, not continue.

Correlation!

Aussagekraft

Antisipasi data yang hilang → Mengandung None atau NaN

↳ scikit learn

↳ `from sklearn.impute import SimpleImputer`

`imputer = SimpleImputer(strategy = "median")`

`imputer.fit(housing_num, y=None)`

pd.DataFrame

Sparse matrix
in
memory

Compressed Sparse Row format
(CSR)

Freitag, 13.10.2023

Pipeline

↳ hard to debug

`from sklearn.pipeline import Pipeline`

has to have

`fit & transform ← class/function in Pipeline`

sklearn metrics

① fit & transform

↳ Transform...

① fit & predict

← model &
predict

Wahrscheinlichkeitsverhältn. (Probability Distribution)

f: Ereignis \rightarrow "Häufig"

Häufigkeit

- diskret

- kontinuierlich

Galton Board / Galton Brett

~ Binomial Distribution

~ Bernoulli Variable

Log-normal distribution

Bayes ..

Daten x , Ziely

Bedingte Wahrscheinlichkeit ~ 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 cum. data independ

L, harusnya coba korelasi dulu & independent test

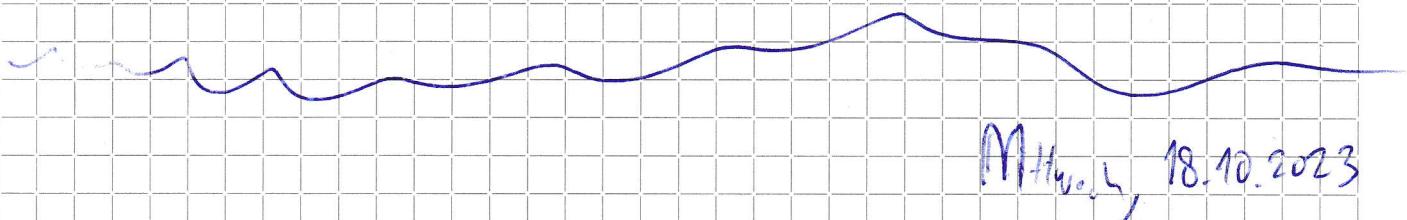
Kategorial verteilt

ca 9

Encoding → Hot
Ordon

Bayes - Klassifikator für nominale Merkmale (Kapitel 9.3.2) } Bsp.
Seite 88 } Fröhle

$$P(i|x) = \frac{\prod_{k=1}^m P(x^{(k)}|i) \cdot P(i)}{\sum_{j=1}^n P(j) \prod_{k=1}^m P(x^{(k)}|j)}$$



M.H., 18.10.2023

Linear Regression
fit in stufen

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

$$y = b + mx$$

Newton Gauß Verfahren

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

Nachenieur ohne Pipeline

Skalierung / Scaling

Linear Regression:

- ~~Skalieren~~ → wenn nur das predict interessiert
Nicht Skalieren
- Skalieren → wenn die Wichtigkeit der Spalten interessiert

Lasso, Ridge

↳ Skalieren, weil es sonst nicht konvergiert.

Bayes

↳ nicht skalieren

 $O(n)$ → Big O ← Aufwand von Algorithmus.Vector n Komponenten → addieren: ~~O(n²)~~ $n \sim O_n$

Matrix n Zeilen, 2 Spalten, Summen der Spalten

 $\tilde{A}r + b \sim O(n)$ Matrix n Zeilen f ~~für~~ f Spalten
 $O(n \cdot f)$,Bayes: $f \cdot O(n \cdot f)$
produkt (af)Linear regression: $O(n \cdot f^2)$

Nächste Nachbarn (h NN)

↳ Latel

↳ Eigenschaften

↳ Gewichtung der Eigenschaften

↳ Gruppierung nach Ähnlichkeit

↳ totale Übereinstimmung

↳ relative kurzer Abstand

KD Tree

+

Recursion

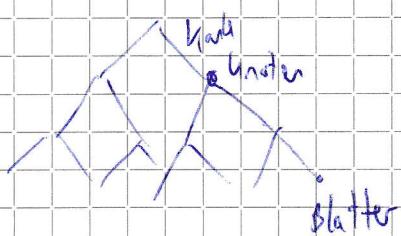
Ball Tree

~~Decision Tree~~ hNN (nearest neighbor)

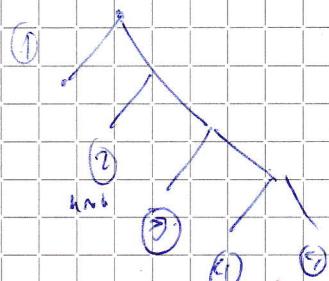
- Distanz:
 - Euclidean (L_2)
 - Haversine
 - Manhattan, cityblock (L_1)
 - Cosine

Decision Tree:

$O(\log(n))$



Pfadlängen:



- Anzahl Blätter: wenig ist besser

↳ Anzahl der Regeln

- Pfadlänge: kurz ist besser

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

- Pfadlängenvar.: weniger ist besser

↳ durchschnittl. Regelgröße

Entropy

Decision Tree

- ↳ minimizing standard deviation \leftarrow continuous data
- ↳ famous of overfitting.

Cross validation

~~gini~~ Decision Tree Classifier

Criterion:

- gini \rightarrow Gini impurity
- $-\log_{10} \{ \}$ Shannon Information Gain
- entropy

Splitter:

- best
- random

Warum benötzen? \leftarrow erklären / Begründen

- ↳ Impul...
↳ Ordinal ...
↳ scalar ...
↳ Algorithmo ...
- } Prozess ...?

SVM = Support Vector Machine

- ↳ takentime in sklearn.

```
from sklearn.datasets import make_blobs
```

← New
version

Steigungsfaktor = gradient

Achsenabschnitt = fit.b.intercept_

Randbereich = margin

SVM = Support Vector Machine

↳ Standard Scaler \rightarrow for SVM \leftarrow so scaling may be needed

$$k(x, x') = \exp\left(-\frac{\|x - x'\|^2}{2\sigma^2}\right)$$

$$\gamma = \frac{1}{2\sigma^2} \quad \Leftarrow \text{for GridSearch}$$

$$k(x, x') = \exp\left(-\gamma \|x - x'\|^2\right)$$

\uparrow
Gamma in GridSearch

γ = gamma in SVC in sklearn

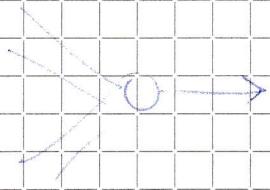
gs = GridSearchCV (svc, parameter)

gs.fit (X, y)

best = gs.best_estimator_

best.predict (Xnew)

Neuronale Netzwerke: \rightarrow Kapitel 7 Früchte



$$e = \hat{y} - y$$

def: residual
outflow

$$\hat{y}_i = \sum_j w_{ij} x_j$$

$$\frac{\partial e}{\partial y_i} = -1$$

$$\frac{\partial e}{\partial w_j} = \cancel{\sum_i} \cancel{\hat{y}_i} x_j$$

$$\frac{\partial e}{\partial w} = x_j$$

$$\frac{\partial e^2}{\partial w} \rightarrow \frac{\partial e}{\partial w} \cdot 2e = 2e x_j$$

Projekt

With score many

- CV - score \leftarrow cross validation

- gs - score \leftarrow grid search score

def __init__(self, scalar=True, scaler=MinMaxScaler()):

 self.sc = scaler

 self.scalar = scalar

def fit(self, X, y)

 if ~~scalar~~ self.scalar:

 X_shal = self.sc.fit_transform(X)

 else

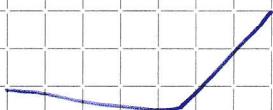
 X_shal = X

predict(\tilde{X})

\hookrightarrow scale or not scale...

Activation function

RELU



Batch Learning & Incremental Learning

Voting Classifier \rightarrow for project ???

\hookrightarrow find best Classifier.

} vs. Grid Search

Random Forest

↳ has many trees

Bagging:

↳ Bagging Classifier & Bagging Regressor.

Bootstrap ↲ in Random Forest

bob-score ↲ out-of-bag ↳ only if Bootstrap is true

↳ RandomForestClassifier

Bagging vs Boosting

Dimensionality Reduction } Feature Engineering
 - feature selection
 - PCA

Composite Estimator

PCA (Hauptkomponentenanalyse)

- Kovarianz-Matrix

Zentriert: $(\mathbf{x} - \bar{\mathbf{x}})$ ← mean will be zero

$\mathbf{x} = \mathbf{x} - \mathbf{x}_{\text{mean}}$

np.linalg.eigh (M)

↓
eigenvekt.

eigh → eigen von hermitian
 eig → eigen

Hermitian Matrix \Rightarrow symmetrisch, diagonale Spiegel
 diagonal mirror

pca. components_ ← ~~all~~ eigen vectors

pca. singular-values ← eigen values

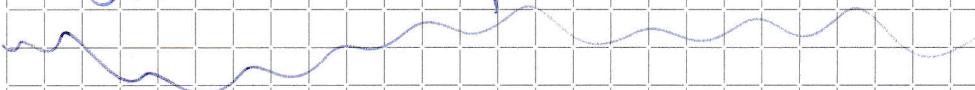
→ sorted from large to small

pca. explained-variance ← sorting pca. components based on it

pca. noise-variance

pca. explained-variance-ratio.

TSNE ~ Manifold



Projekt:

↳ Ziel von Projekt ① Richtigkeit
↓ was ist minimiert/maximiert

Recall = $\frac{\text{gefundene in Vorrau}}{\text{alle ech.}}$

Precision = $\frac{\text{richtig in Vorrau}}{\text{alle Vorburau}}$

Bericht → PDF
↳ Wind

Freitag, 27.10.2023

Clustering

↳ DB SCAN

↳ Hierarchical Clustering

↳ K-Means

DB SCAN ?

↳ Kernpunkt

↳ Randpunkt

↳ Rauschpunkt

Kann es nicht random

when is cluster good?

KD-Baume ..

Montag, 30.10.2023

Red Wine & White Wine → CSV Datei.

output = kategorisch.

- Bayesian? score? →
- kNN?
- Neural network? X
- PCA?

Duplicate?

Ausreißer

df1.loc[df1.duplicated(), :]

Woj für Sorte (Metric)
↳ precision
↳ recall

↳ classification_report.

↳ confusion matrix

Nur Weiß Wein ↲

↳ outlier outlier ...

Dienstag, 31.10.2023

3, 4 schlich

5, 6 Θ^h

7, 8, 9 gJ

SVM ov

PCA + Logist. Regress.

Precision 80%. Gute Weissweine

Oder

Recall

~~fit~~ solver : liblinear -> penalty: l2 ↗ over
solver : liblinear -> penalty: l1 ↗ multiclass auto.

solver -> SGD -> l2 ↗ over
l1 ↗ m. theory

solver -> SGD -> l1 ↗ over
l2 ↗ over
elastic ↗ over
ridge ↗ over

~~Mittwoch~~, 02.11.2023
Donnerstag,

Klassen

- cross-val. score(...)

①

{ 2 Klassen ?

Abgabe: Freitag mittag

[Aufgabe behalten] ???

pickle Datei laden

Gaussian Naive Bayes ???

, , , +



Log Reg.