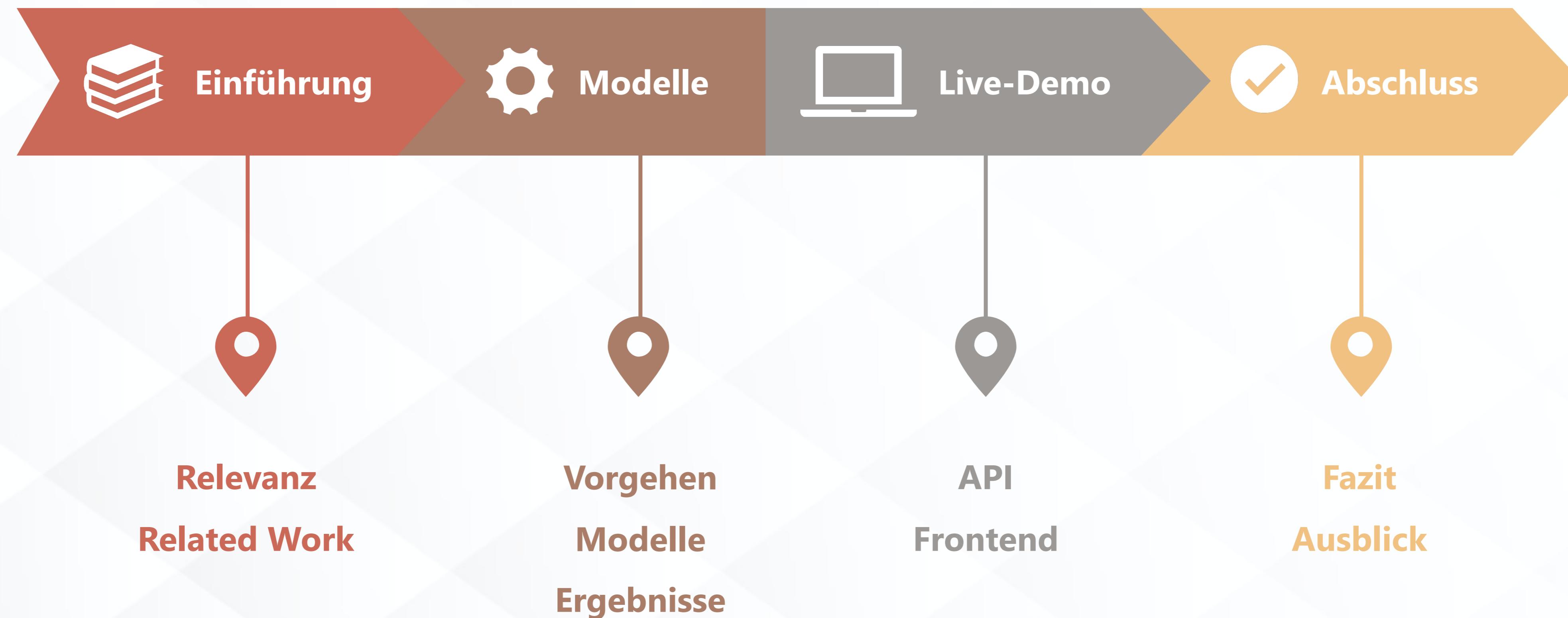




Agenda



Einführung

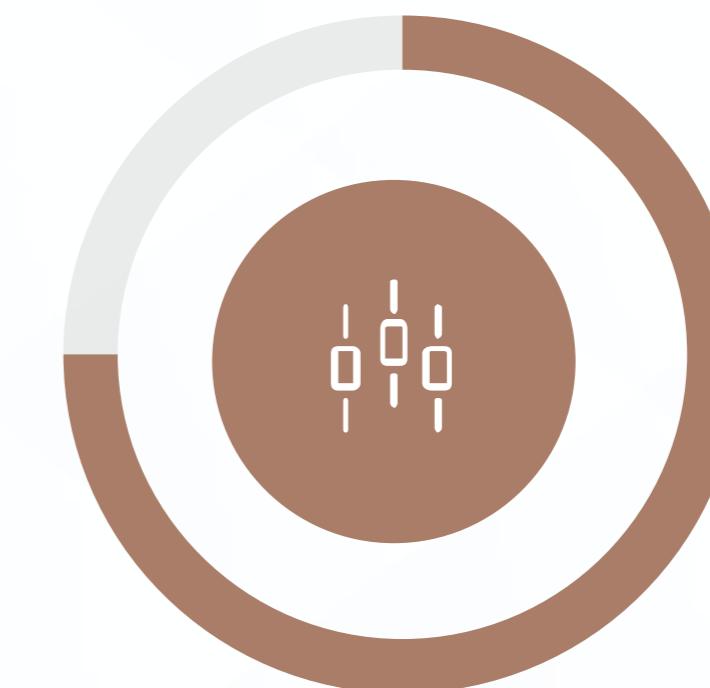


Stilometrie

Untersuchung von Sprachstilen mithilfe statistischer Mittel



Jeder Autor besitzt einen
(unbewussten)
individuellen Schreibstil



Schreibstile lassen sich durch
die konsistente Verwendung
bestimmter Muster statistisch
abgrenzen

Relevanz

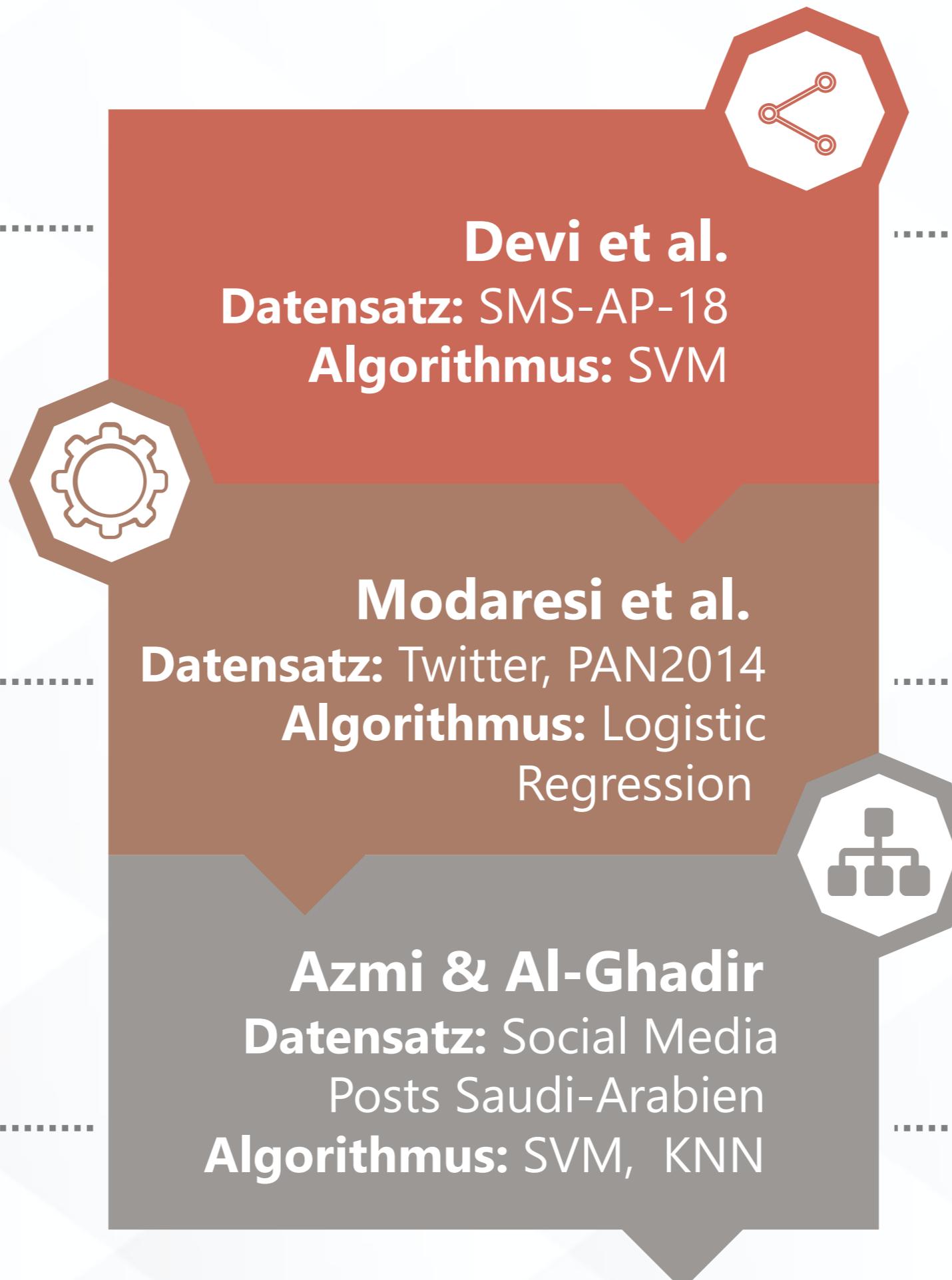


Related Work

Geschlecht:
85,7%
Accuracy

Geschlecht:
75,6%
Accuracy

Geschlecht SVM:
87,3%
Accuracy

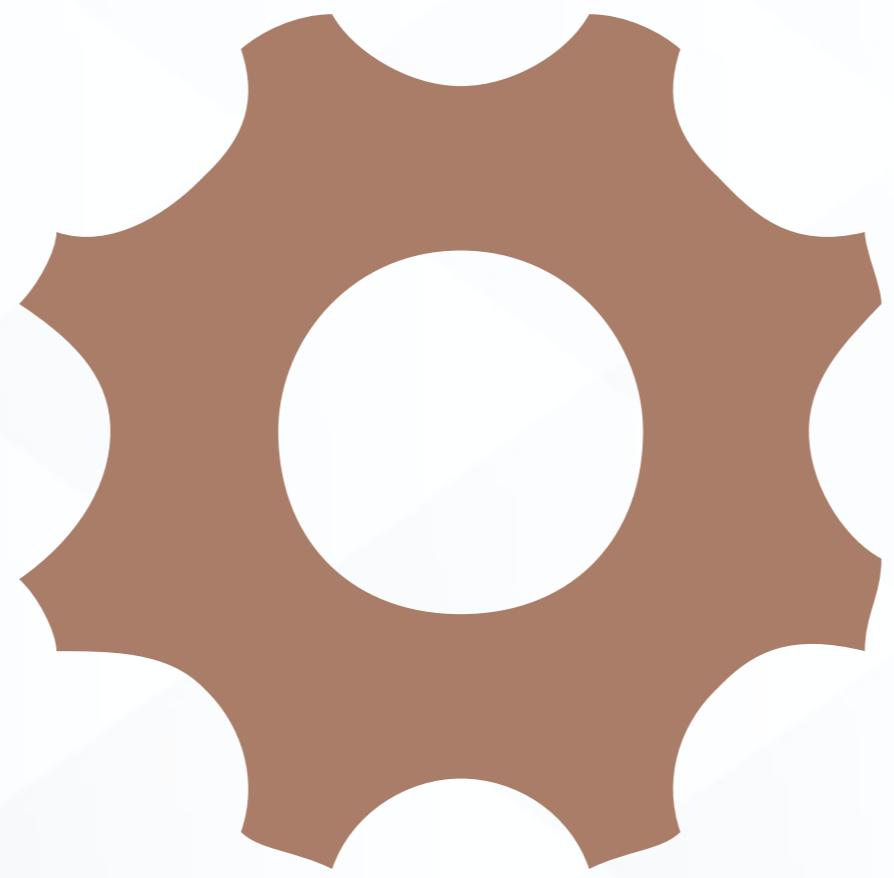


Alter:
64,3%
Accuracy

Alter:
51,7%
Accuracy

Geschlecht KNN:
93,1%
Accuracy

Modelle



Data Cleaning

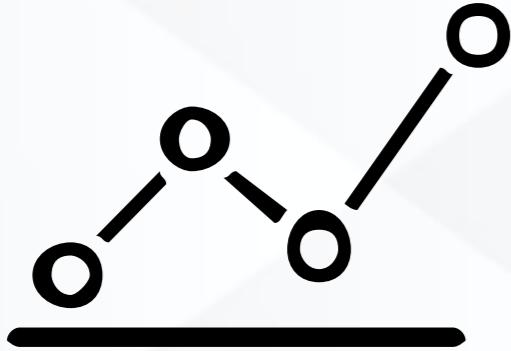
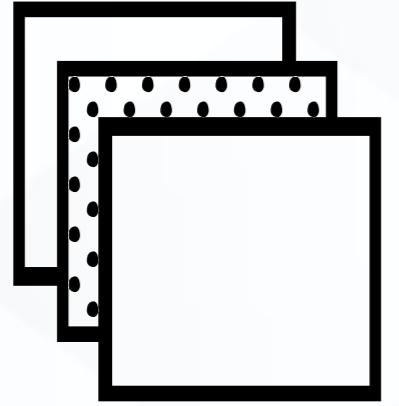


Feature Engineering



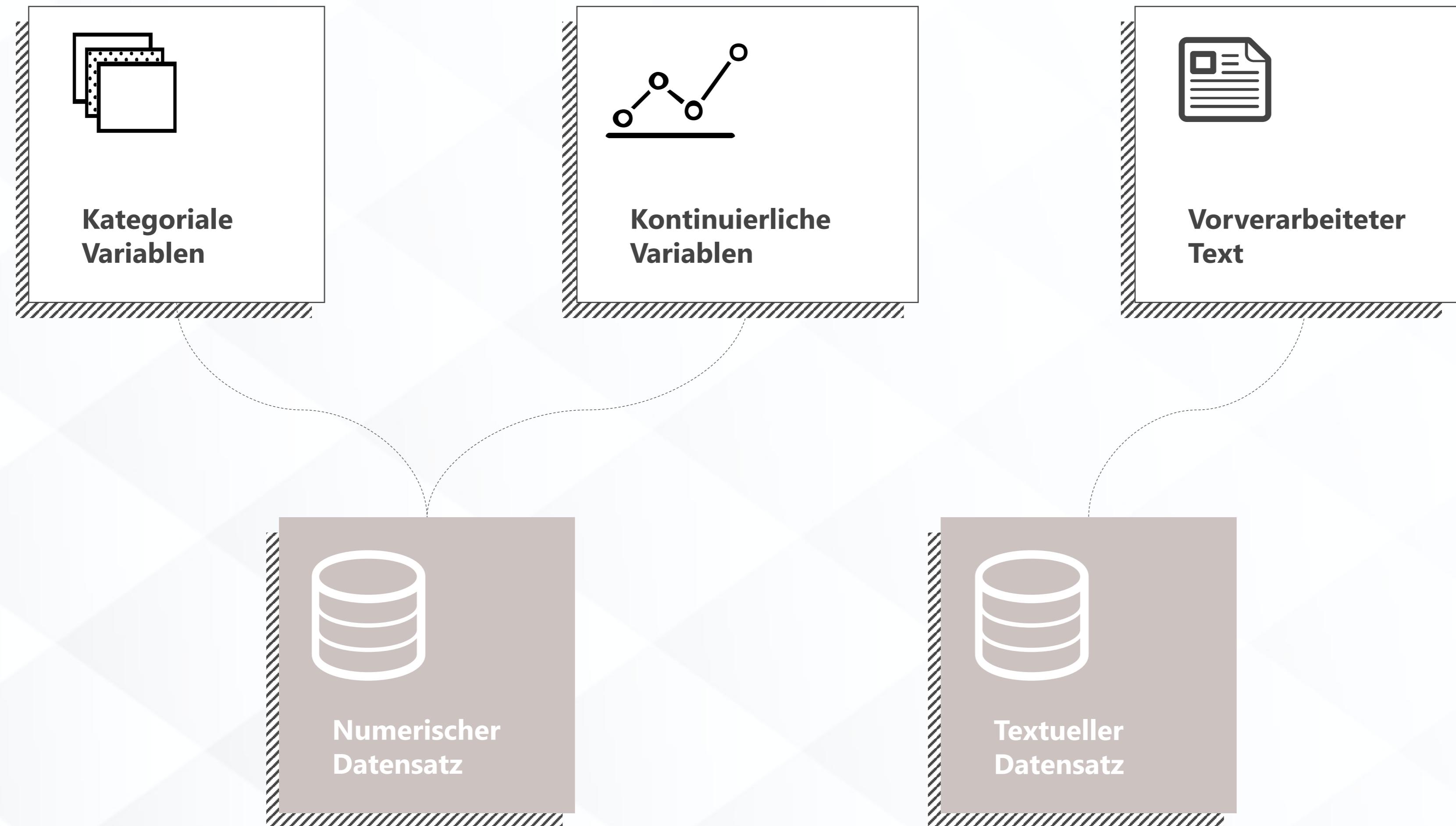
Arten von Features

Vorhandene Arten von Features nach dem Feature Engineering



Aufteilen der Features

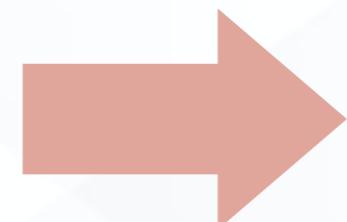
Trennung des Textes von verbleibenden numerischen Features



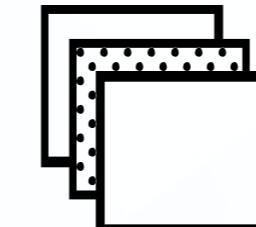
Numerische Transformation

Verarbeitung der nicht-text Features

Gender		
female	1	0
male	0	1
female	1	0
male	0	1

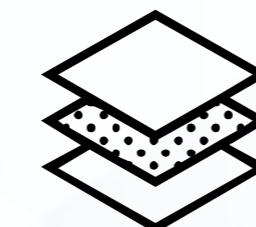


$$Z(x) = \frac{x - \mu}{\sigma}$$



One-Hot Encoding

Für jede Ausprägung einer kategorialen Variable im Datensatz wird eine Variable erstellt. Diese trägt entweder die Ausprägung 1 oder 0, wobei 1 indiziert, ob die Ausprägung der initialen Variable der Ausprägung der neuen Variable entspricht.



Z-Transformation

Alle kontinuierlichen Variablen sind mittels Z-Transformation standardisiert worden, so dass ihr Mittelwert 0 und die Standardabweichung 1 entspricht. Dies vermeidet numerische Instabilitäten.

Textuelle Transformation

Transformation der Texte in TF*IDF-Vektoren

TF*IDF

Term Frequency

Relativer Anteil eines Terms in einem Dokument

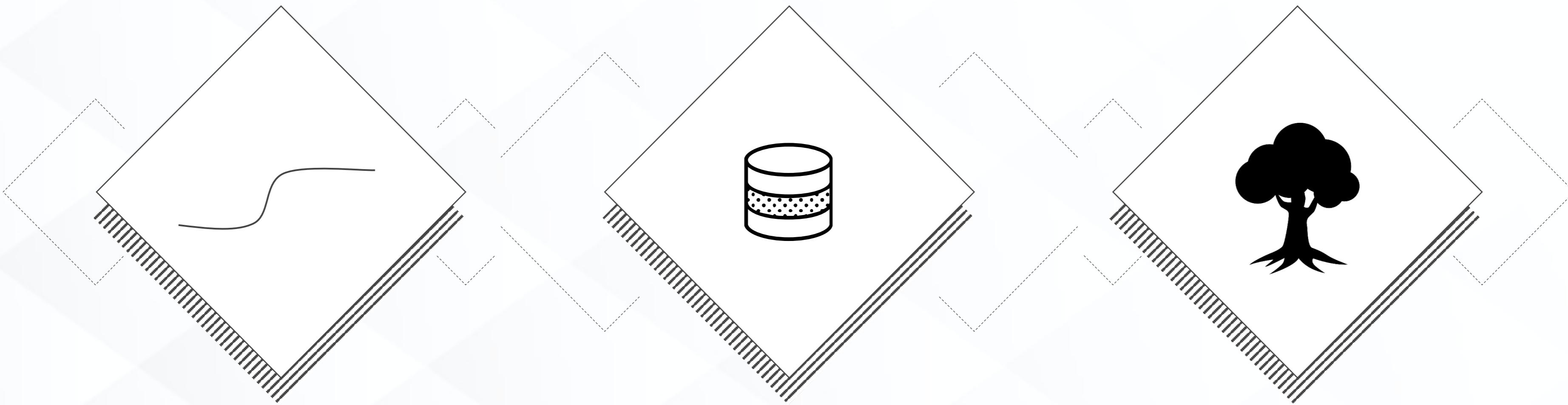
*

Inverse Document Frequency

Inverser Relativer Anteil an Dokumenten, die diesen Term enthalten

Supervised Learning

Modelle die zur Zielwertvorhersage verwendet worden sind



Logistische Regression

Die Logistische Regression kann mittels One-Versus-Rest-Verfahren für univariante und multinominale Klassifikation verwendet werden.

Support Vector Machine

Die Support Vector Machine bestimmt Entscheidungsgrenzen, so dass die Differenz zwischen den Klassen maximiert wird.

XGBoost

Baumbasierter Algorithmus, der auf Gradientenverfahren und Boosting zurückgreift.

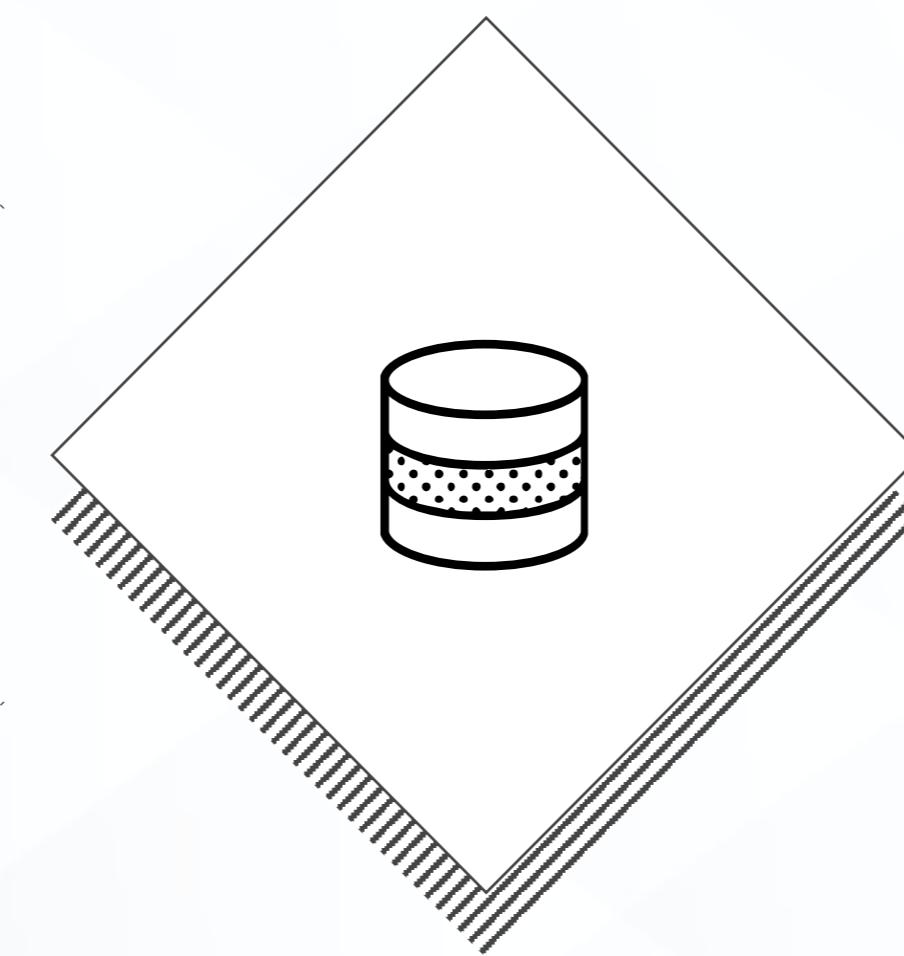
Parametertuning

durch Cross Validation



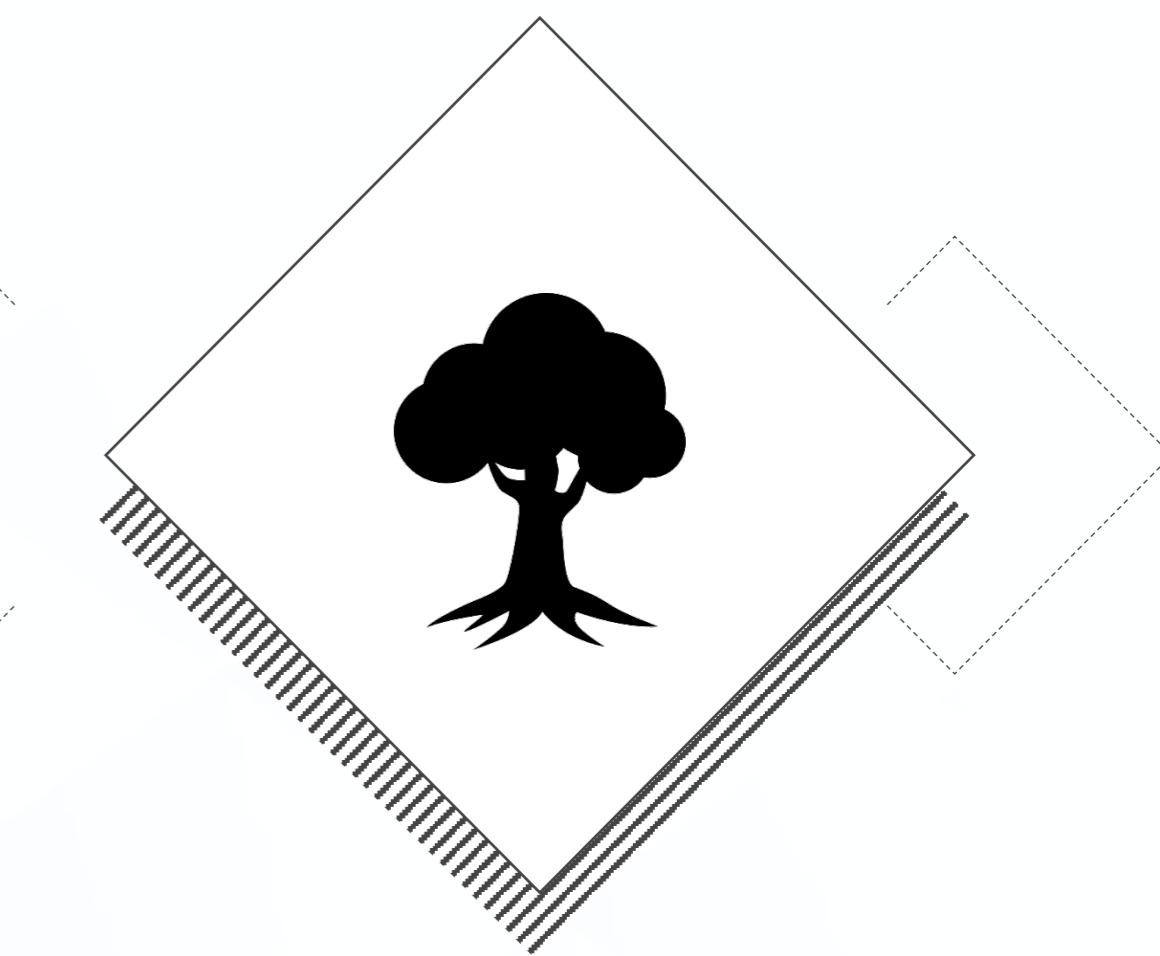
Logistische Regression

Alpha: 0,001; 0,00001
Penalty: elasticnet;
L1-Norm;
L2-Norm



Support Vector Machine

Loss: epsilon insensitive;
squared epsilon insensitive;
hinge;
squared hinge
C: 0,8; 1; 1,2



XGBoost

Learning Rate: 0,1; 1; 1,5
Max Depth: 3; 6
N Estimators: 400; 800

Modell Evaluation



Alter

Text: Linear Support Vector Regression, 49,109%

Numerisch: XGBoost, 50,891%

Negativer MSE: -4,648



Geschlecht

Text: Logistische Regression, 45,734%

Numerisch: XGBoost, 54,266%

F1-Score: 80,30%



Thema

Text: Linear Support Vector Classification, 38,950%

Numerisch: XGBoost, 61,050%

F1-Score: 39,307%



Sternzeichen

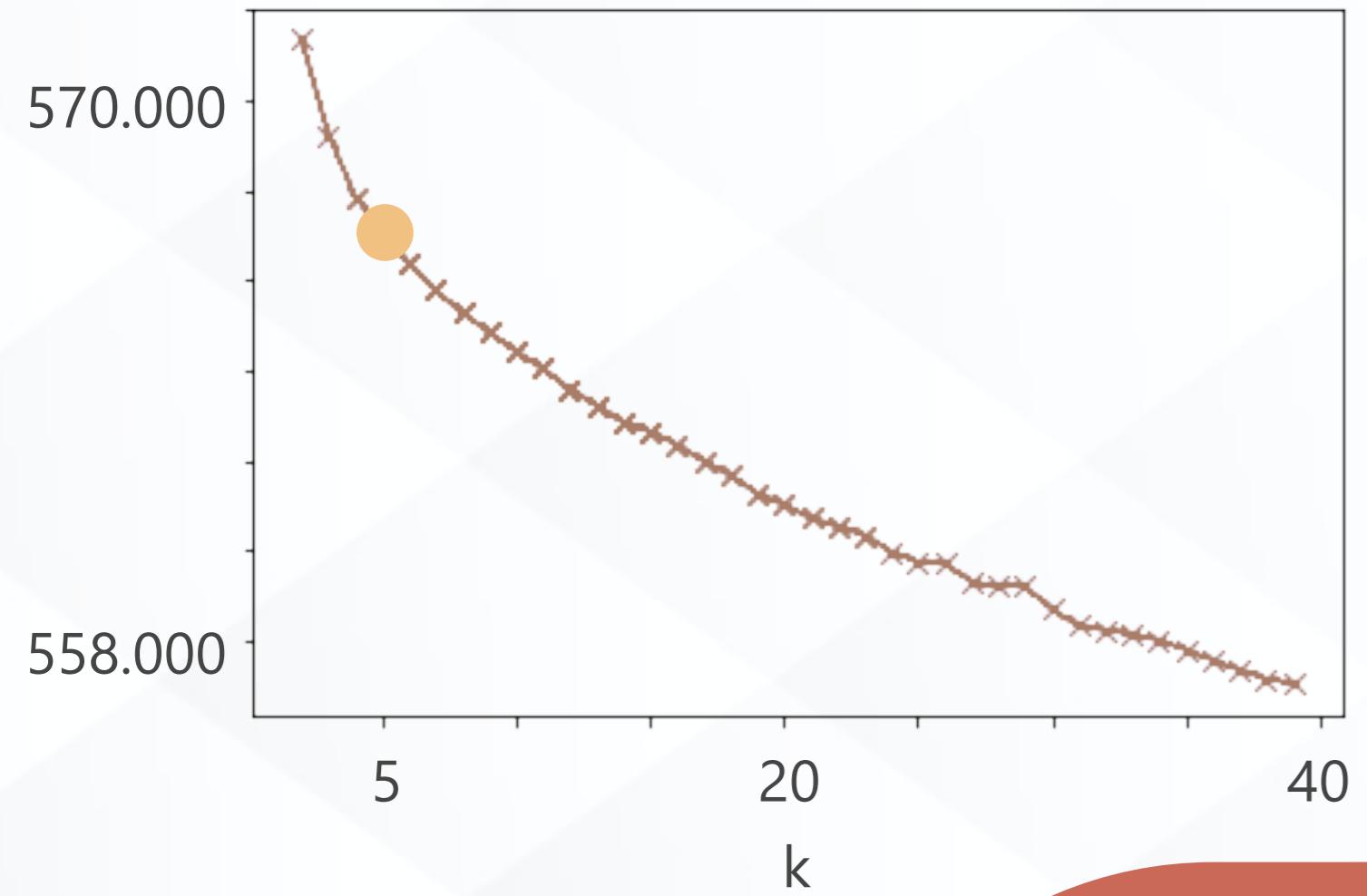
Text: logistische Regression, 46,390%

Numerisch: XGBoost, 53,710%

F1-Score: 42,20%

Unsupervised Learning

k-means



Numerisch

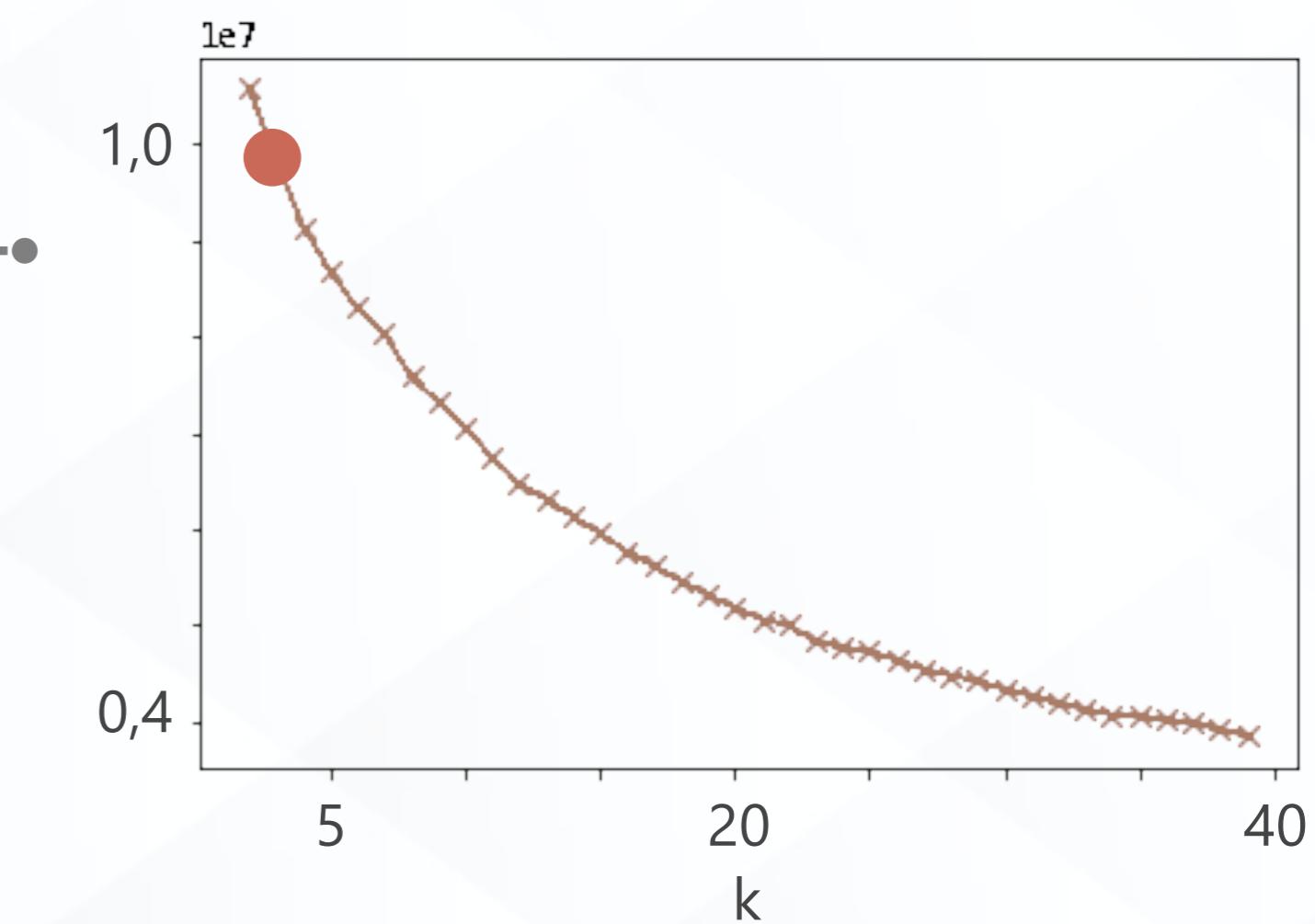
Basierend auf
stilistischen
Eigenheiten

"Hobby Publisher"

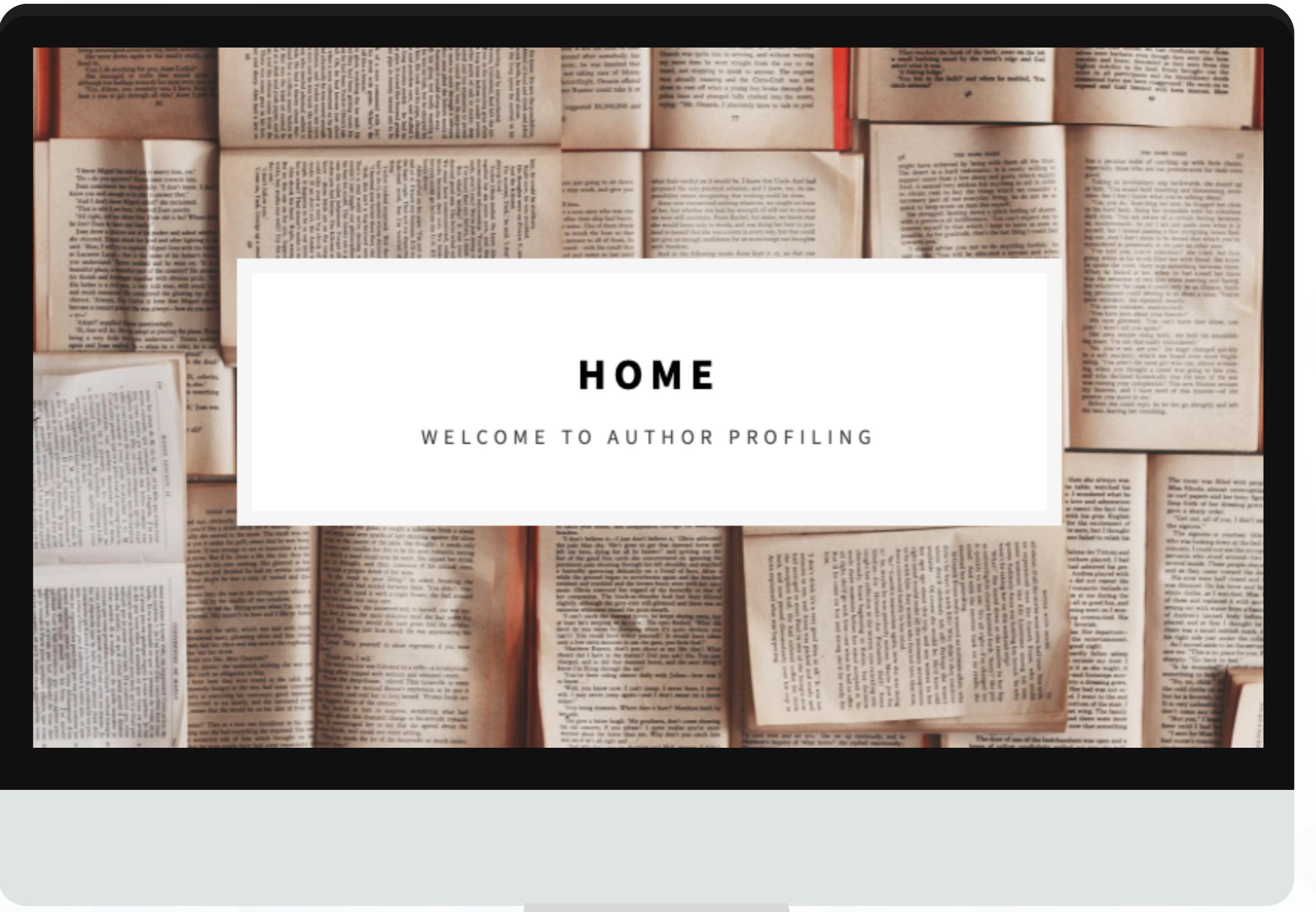


Textuell

Basierend auf
Wortwahl
"Negation-Lover"



Live-Demo



Abschluss



Fazit

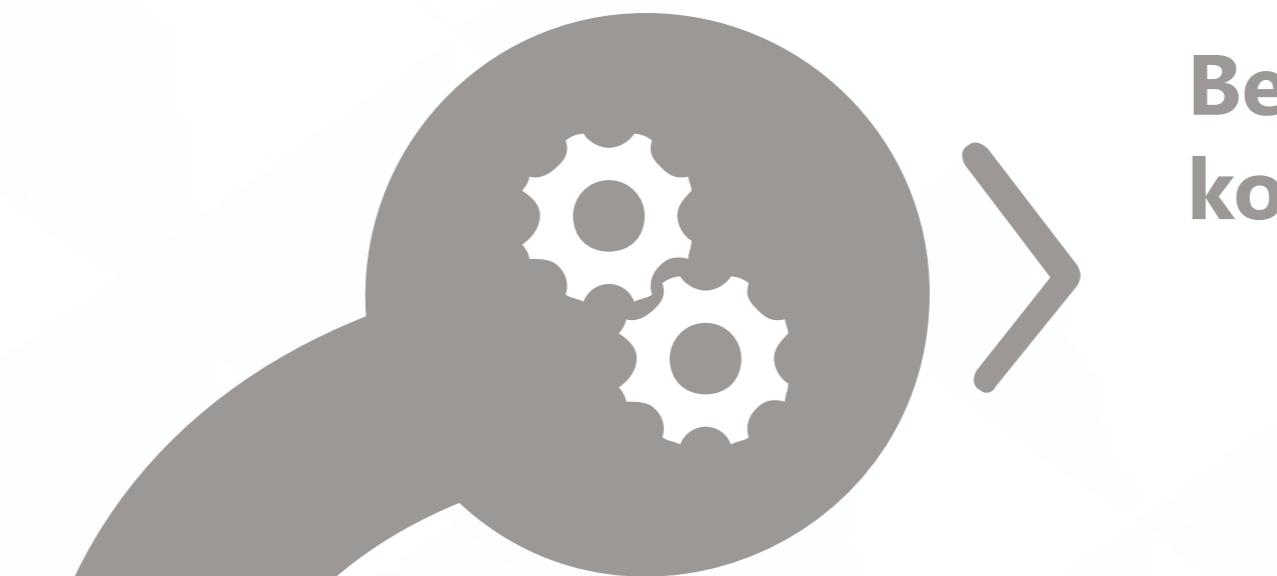
Modelle ähnlich gut wie Related Work



Mehrwert durch breiten Vergleich von Algorithmen und Hyperparametern



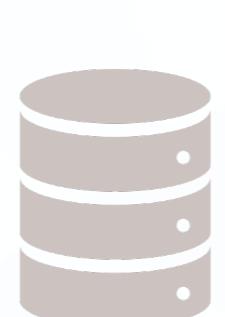
Bessere Ergebnisse durch kombinierte Modelle



Wenige neue Cluster, dafür eindeutige Merkmale



Ausblick



Diverserer und hochwertigerer Datensatz

Weitere Clusteringalgorithmen

Verwendung von Word-Embeddings