

# 1\_Measures

## Herleitung

Um die Güte des Clusters zu bewerten benötigt man ein Measure. Intuitiv wird die Bewertung über eine 3D Visualisierung vorgenommen, was jedoch sehr zeitaufwendig ist.

## Recherche

Menge an möglichen Measures:

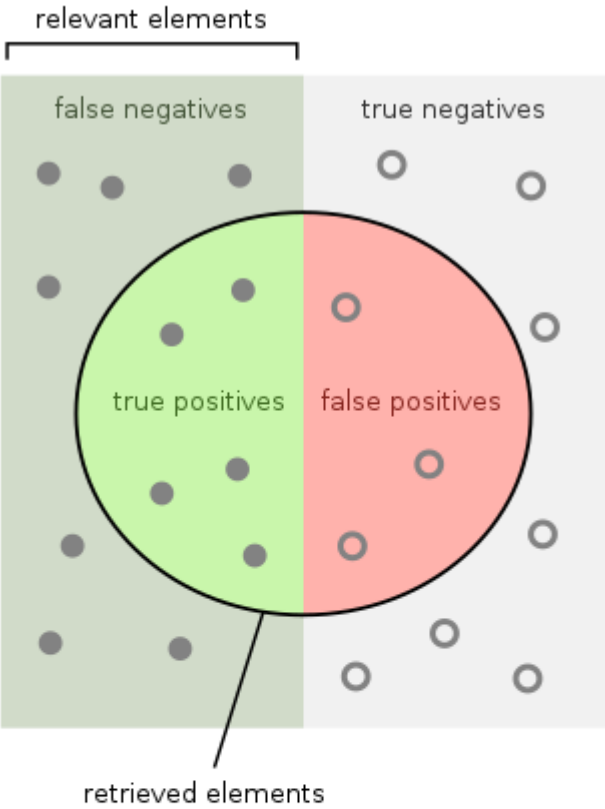
- True Positives
- True Negatives
- False Positives
- False Negatives
- Accuracy
- Precision
- Recall
- F1-Score
- Adjusted Rand Index (ARI)
- Adjusted Mutual Information (AMI)
- Fowlkes-Mallows Index (FMI)
- ggf. weitere (siehe Data Science Vault)

## Hypothese

Es kann ein Measure identifiziert werden, was am besten geeignet ist für die Gütemessung

## Ergebnisse

Folgende Grafik dient zur Argumentation



How many retrieved items are relevant?

Precision =  $\frac{\text{true positives}}{\text{true positives} + \text{false positives}}$

How many relevant items are retrieved?

Recall =  $\frac{\text{true positives}}{\text{true positives} + \text{false negatives}}$

Sources: [4][5][6][7][8][9][10][11][12] view · talk · edit

		Predicted condition			
		Predicted Positive (PP)	Predicted Negative (PN)	Informedness, bookmaker informedness (BM) $= \text{TPR} + \text{TNR} - 1$	Prevalence threshold (PT) $= \frac{\sqrt{\text{TPR} \times \text{FPR}} - \text{FPR}}{\text{TPR} - \text{FPR}}$
Actual condition	Total population $= P + N$				
	Positive (P) <sup>[a]</sup>	True positive (TP), hit <sup>[b]</sup>	False negative (FN), miss, underestimation	True positive rate (TPR), recall, sensitivity (SEN), probability of detection, hit rate, power $= \frac{\text{TP}}{P} = 1 - \text{FNR}$	False negative rate (FNR), miss rate type II error <sup>[c]</sup> $= \frac{\text{FN}}{P} = 1 - \text{TPR}$
	Negative (N) <sup>[d]</sup>	False positive (FP), false alarm, overestimation	True negative (TN), correct rejection <sup>[e]</sup>	False positive rate (FPR), probability of false alarm, fall-out type I error <sup>[f]</sup> $= \frac{\text{FP}}{N} = 1 - \text{TNR}$	True negative rate (TNR), specificity (SPC), selectivity $= \frac{\text{TN}}{N} = 1 - \text{FPR}$
		Positive predictive value (PPV), precision $= \frac{\text{TP}}{\text{PP}} = 1 - \text{FDR}$	False omission rate (FOR) $= \frac{\text{FN}}{\text{PN}} = 1 - \text{NPV}$	Positive likelihood ratio (LR+) $= \frac{\text{TPR}}{\text{FPR}}$	Negative likelihood ratio (LR-) $= \frac{\text{FNR}}{\text{TNR}}$
		Accuracy (ACC) $= \frac{\text{TP} + \text{TN}}{P + N}$	False discovery rate (FDR) $= \frac{\text{FP}}{\text{PP}} = 1 - \text{PPV}$	Negative predictive value (NPV) $= \frac{\text{TN}}{\text{PN}} = 1 - \text{FOR}$	Markedness (MK), deltaP ( $\Delta p$ ) $= \text{PPV} + \text{NPV} - 1$
		Balanced accuracy (BA) $= \frac{\text{TPR} + \text{TNR}}{2}$	F <sub>1</sub> score $= \frac{2 \text{PPV} \times \text{TPR}}{\text{PPV} + \text{TPR}} = \frac{2 \text{TP}}{2 \text{TP} + \text{FP} + \text{FN}}$	Matthews correlation coefficient (MCC) $= \frac{\sqrt{\text{TPR} \times \text{TNR} \times \text{PPV} \times \text{NPV}}}{\sqrt{\text{FNR} \times \text{FPR} \times \text{FOR} \times \text{FDR}}}$	Diagnostic odds ratio (DOR) $= \frac{\text{LR+}}{\text{LR-}}$
			Fowlkes–Mallows index (FM) $= \sqrt{\text{PPV} \times \text{TPR}}$		Threat score (TS), critical success index (CSI), Jaccard index $= \frac{\text{TP}}{\text{TP} + \text{FN} + \text{FP}}$

## True Positives

- Führt dazu, dass eine Adjazenzmatrix gefüllt nur mit 1, die beste Alternative ist.
- Alle Wörter sind im selbem Cluster
- **ungeeignet**

## True Negatives

- Führt dazu, dass eine Adjazenzmatrix gefüllt nur mit 0, die beste Alternative ist.
- Alle Wörter sind eigene Cluster
- **ungeeignet**

## False Positives

- Führt dazu, dass eine Adjazenzmatrix gefüllt nur mit 0, die beste Alternative ist.
- Alle Wörter sind eigene Cluster
- **ungeeignet**

## False Negatives

- Führt dazu, dass eine Adjazenzmatrix gefüllt nur mit 1, die beste Alternative ist.
- Alle Wörter sind im selbem Cluster
- **ungeeignet**

# Accuracy

Da die Adjazenzmatrix eine sparse-Matrix ist, kann die Accuracy ebenfalls nicht genutzt werden. Die große Anzahl von True Negatives, die erreicht werden, wenn keine Verbindungen hergestellt werden, überwiegt bei weitem die wenigen Fälle, in denen True Positivs einen positiven Effekt hätten.

ungeeignet

# Precision

Schon geeigneter, bezieht aber keine negativen Fälle mit ein. Kann aber genutzt werden.

ungeeignet

# Hyperparameter Optimierung

## Visuelle Beurteilung

# Recall

Nicht geeignet, da auch hier ein guter Recall erreicht wird, wenn ich alle Verbindungen zulasse.

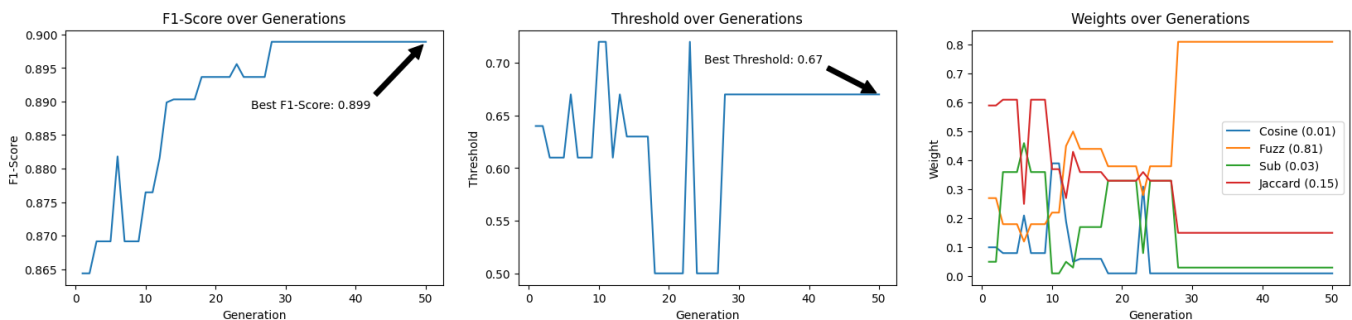
ungeeignet

# F1-Score

Schon geeigneter, bezieht aber keine negativen Fälle mit ein. Kann aber genutzt werden.

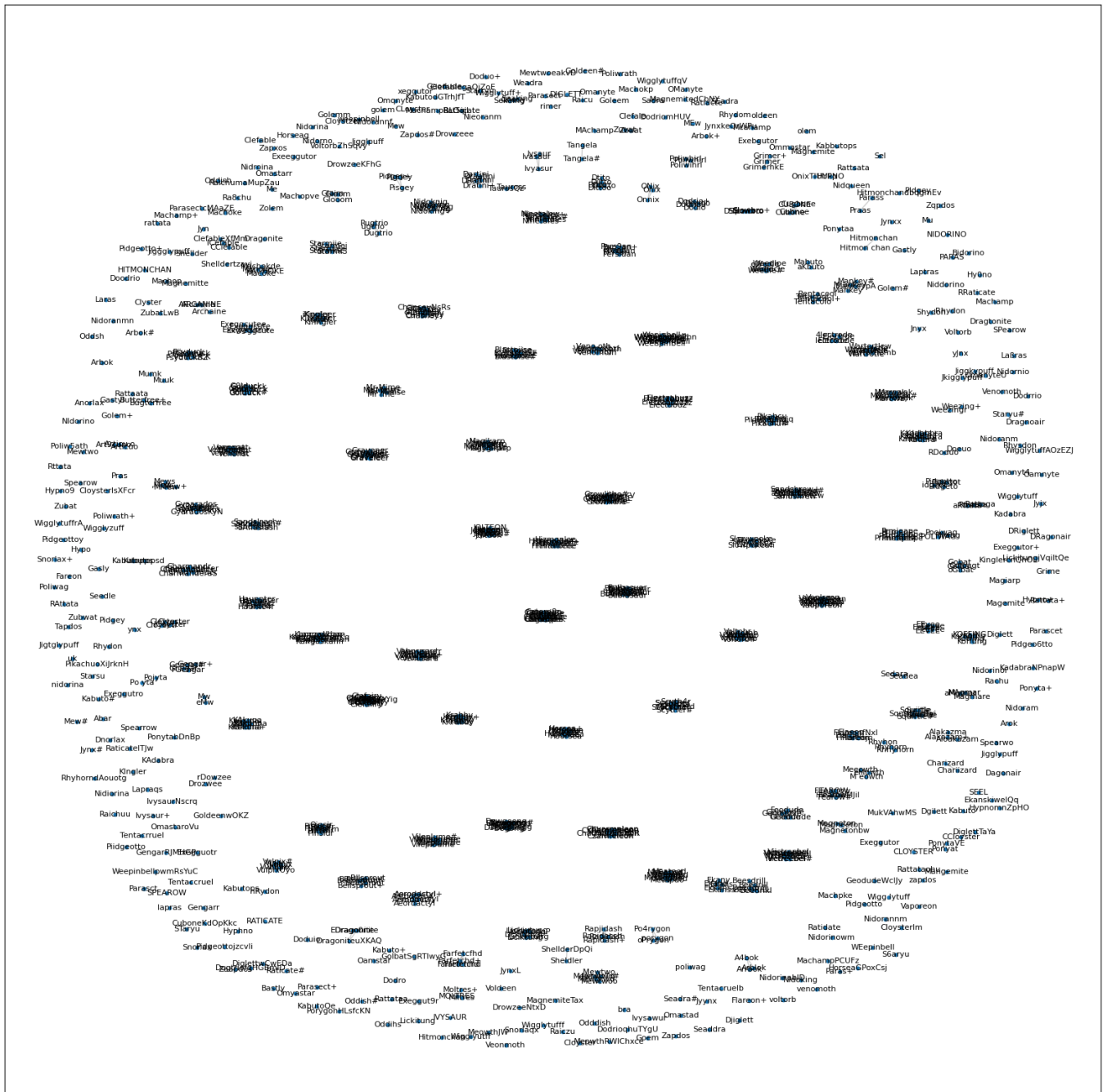
geeignet

# Hyperparameter Optimierung



Merkwürdigerweise verläuft die Fitness nicht monoton steigend.

## Visuelle Beurteilung



Größere Cluster und weniger einzelne Wörter

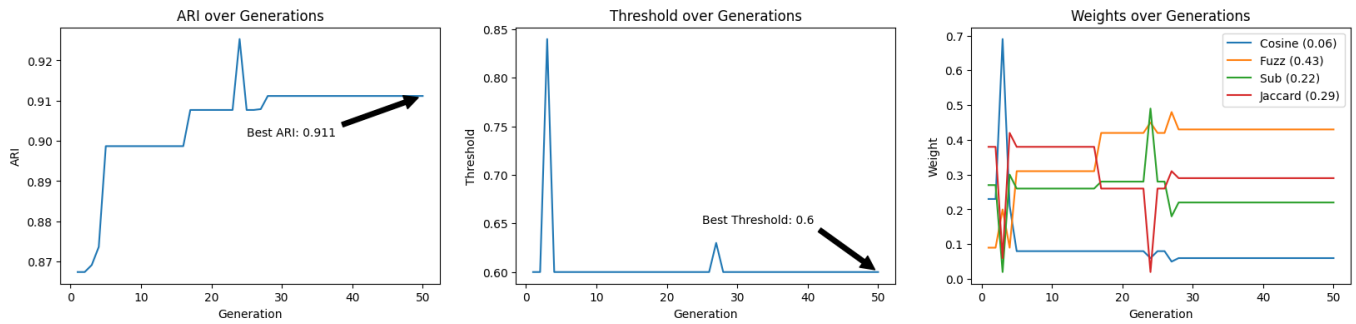
[F1\\_Score.html](#)

## ARI

Geeignet für Clusteranalysen, bei denen die Form und Größe der Cluster variieren können und die Bewertung der Ähnlichkeit zwischen zwei Clusterzuordnungen erforderlich ist.

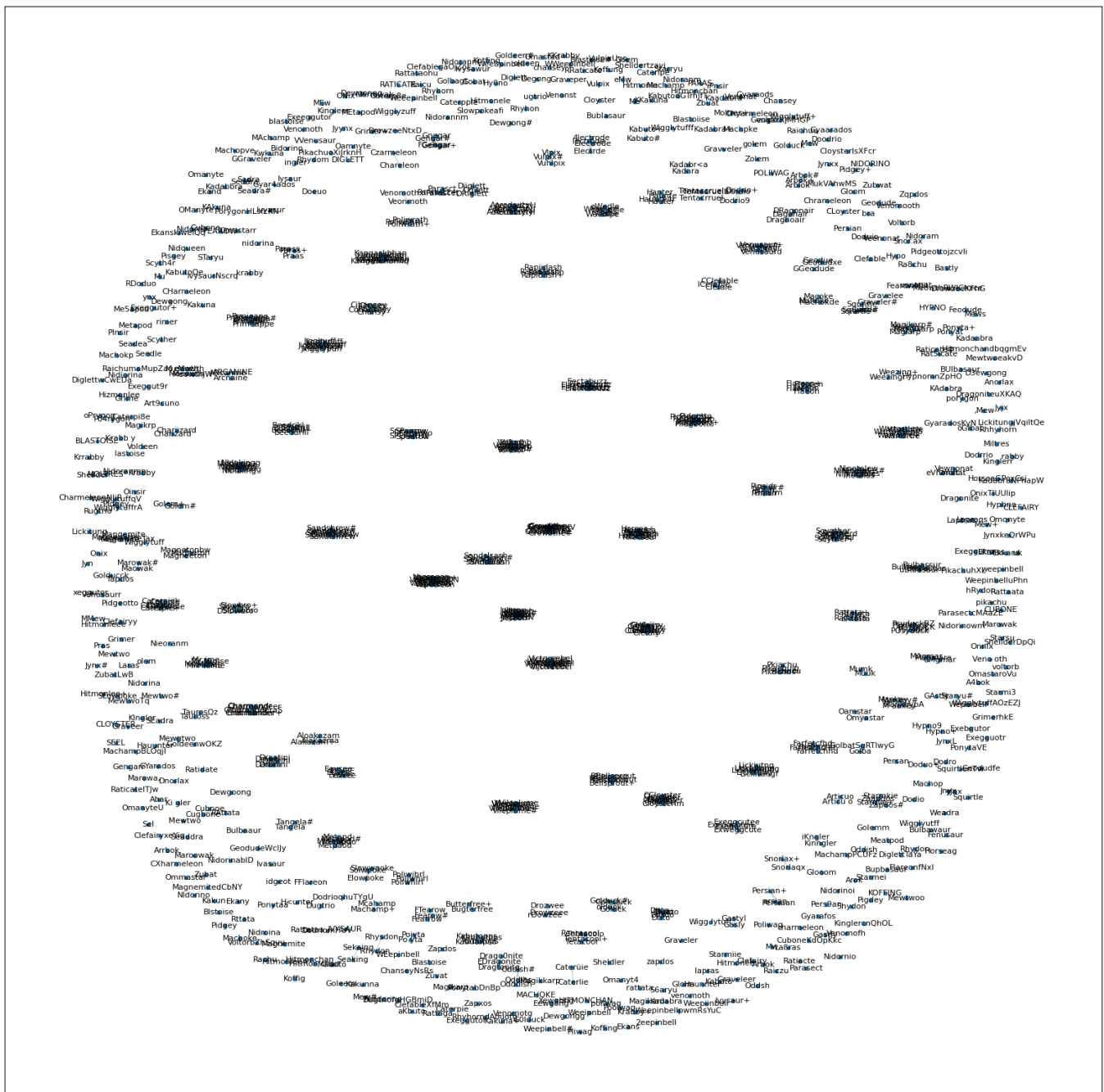
geeignet

## Hyperparameter Optimierung



Merkwürdigerweise verläuft die Fitness nicht monoton steigend. Gewichte sind ähnlich zu AMI. Threshold ist wieder halbwegs gleich.

## Visuelle Beurteilung



Wieder sehr viele einzelnen Wörter ohne Cluster.

 ARI.html

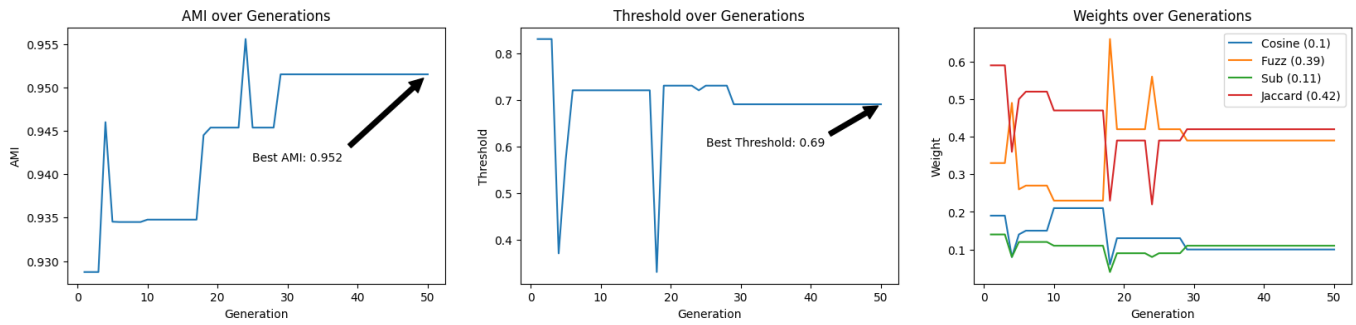
# AMI

Geeignet für Clusteranalysen, bei denen die Form und Größe der Cluster variieren können und die Bewertung der Ähnlichkeit zwischen zwei Clusterzuordnungen erforderlich ist.

**geeignet**

## Hyperparameter Optimierung

## Recht langsame Performance

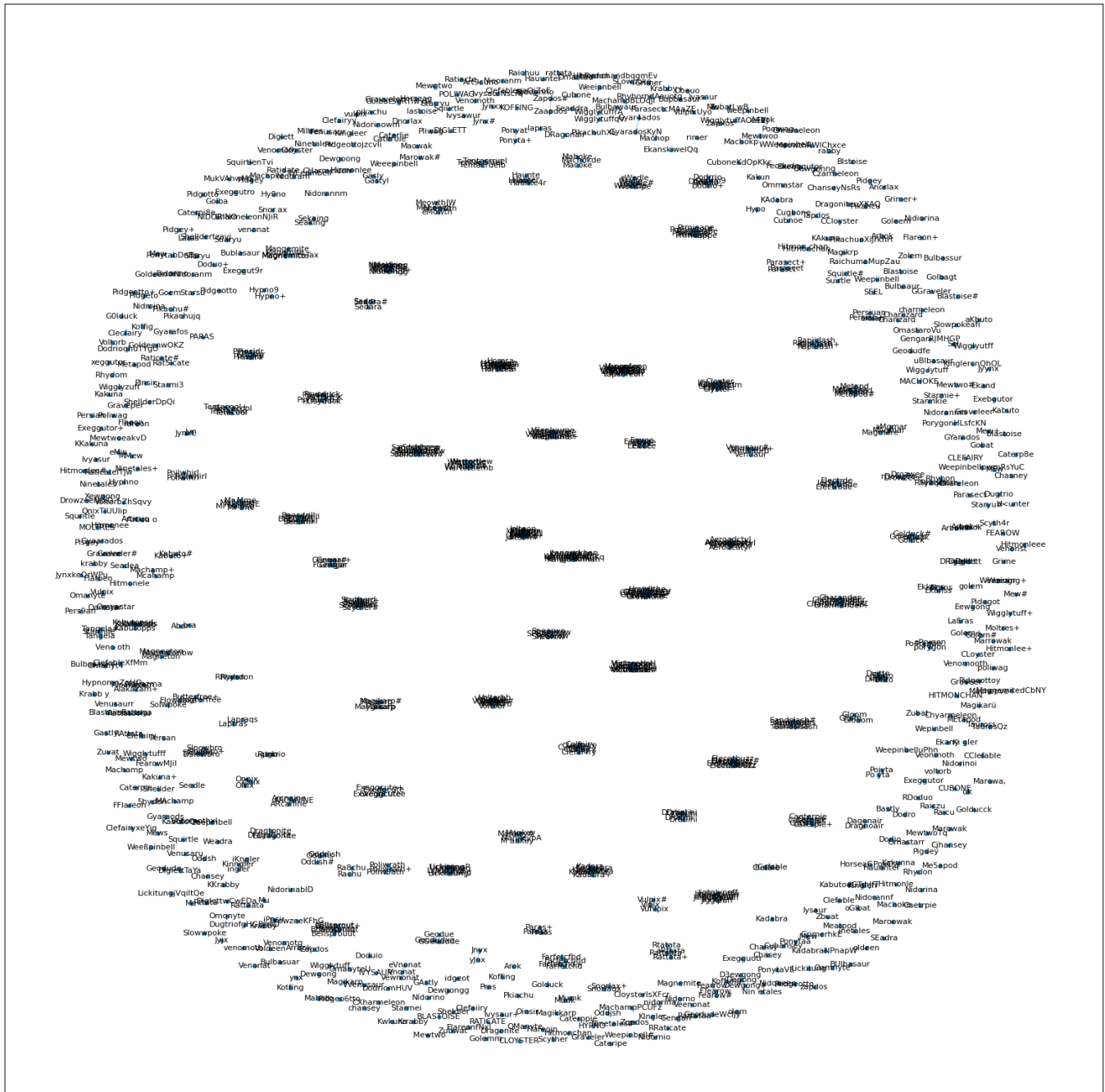


Merkwürdigerweise verläuft die Fitness nicht monoton steigend. Gewichte weichen zu F1-Score und FMI ab. Threshold ist ähnlich.

## Visuelle Beurteilung



## Viele einzelne Wörter ohne Cluster



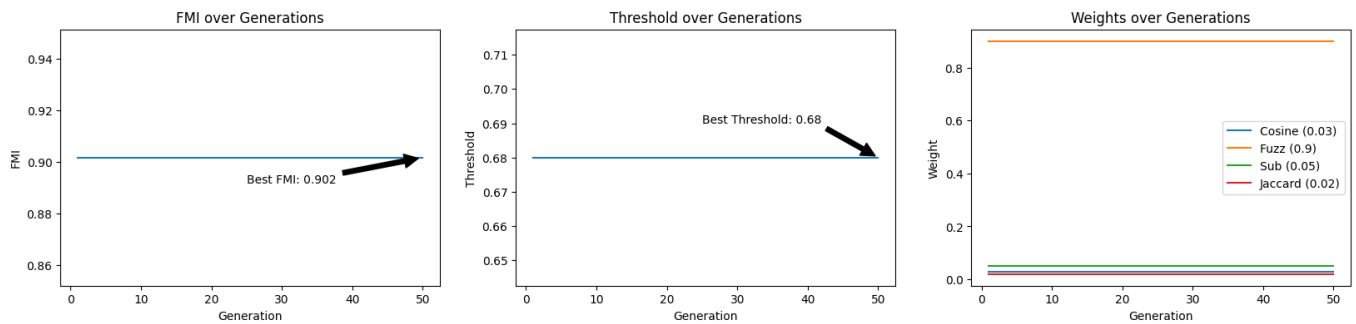
AMI.html

## FMI

Gut geeignet, wenn die Größe der Cluster bekannt ist und die Bewertung der Ähnlichkeit zwischen den Clusterzuordnungen wichtig ist.

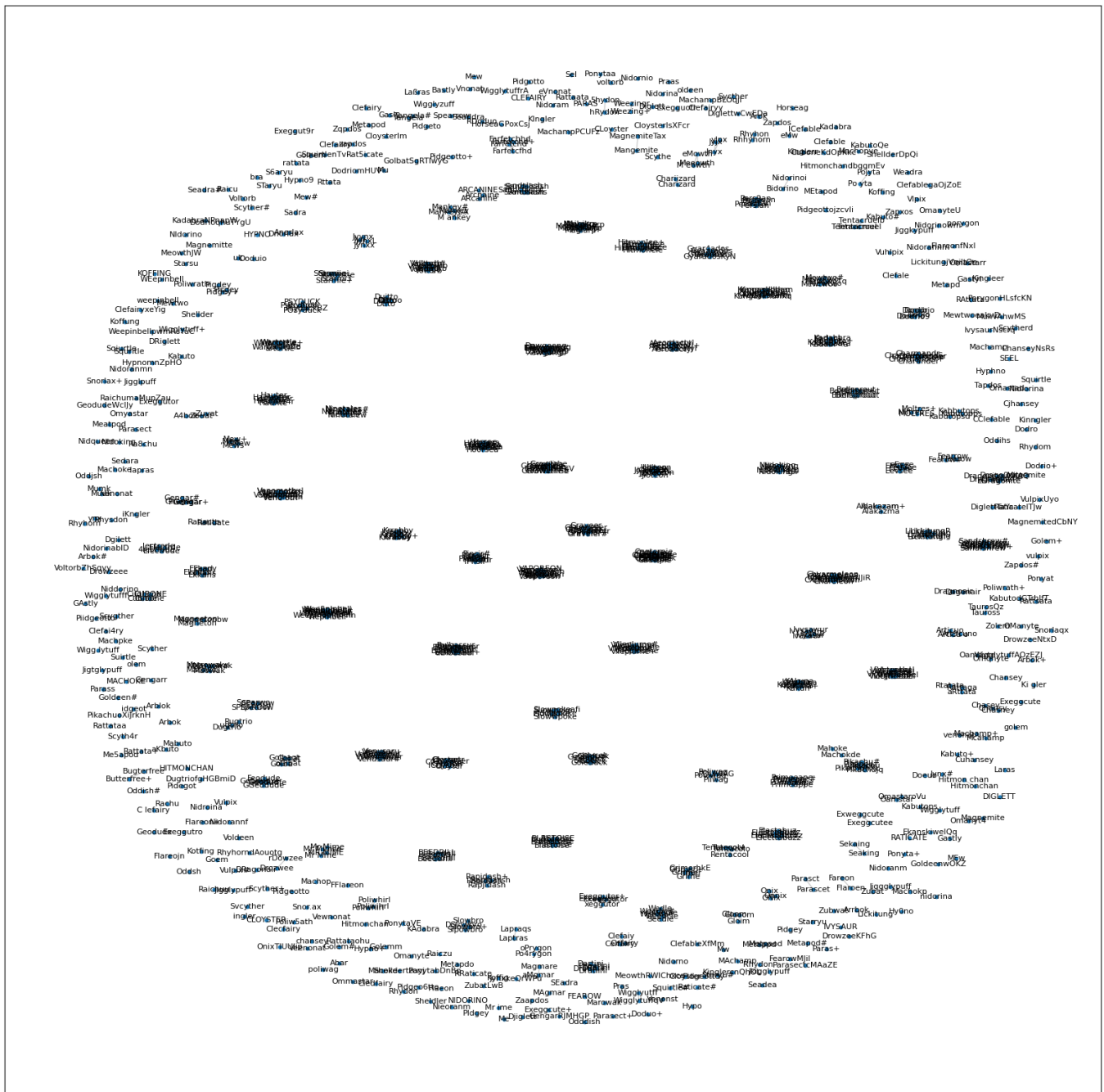
bedingt geeignet

## Hyperparameter Optimierung



Funktioniert recht gut. Gewichte sind ähnlich zu F1-Score. Sehr schnelle Performance

## Visuelle Beurteilung



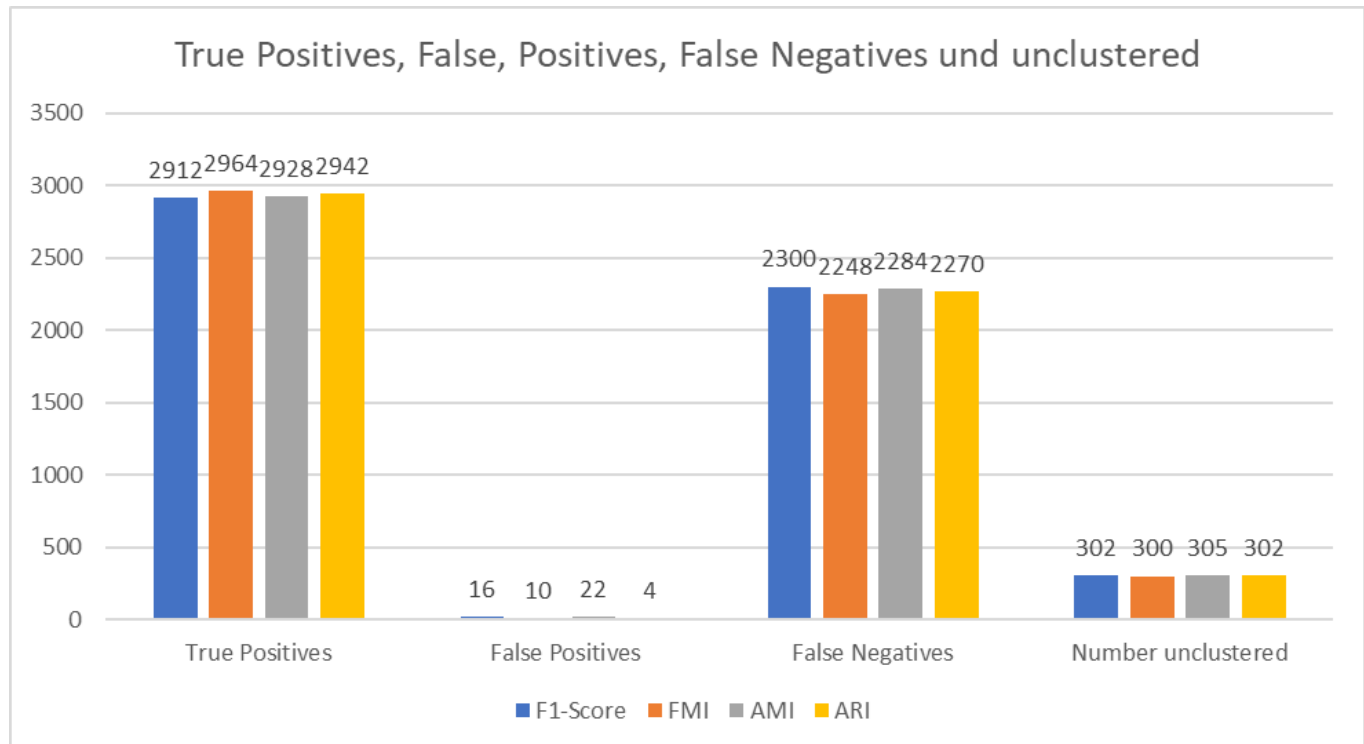
Bei optischer Bewertung, fallen einige nicht-geclusterte Wörter auf, die augenscheinlich beim

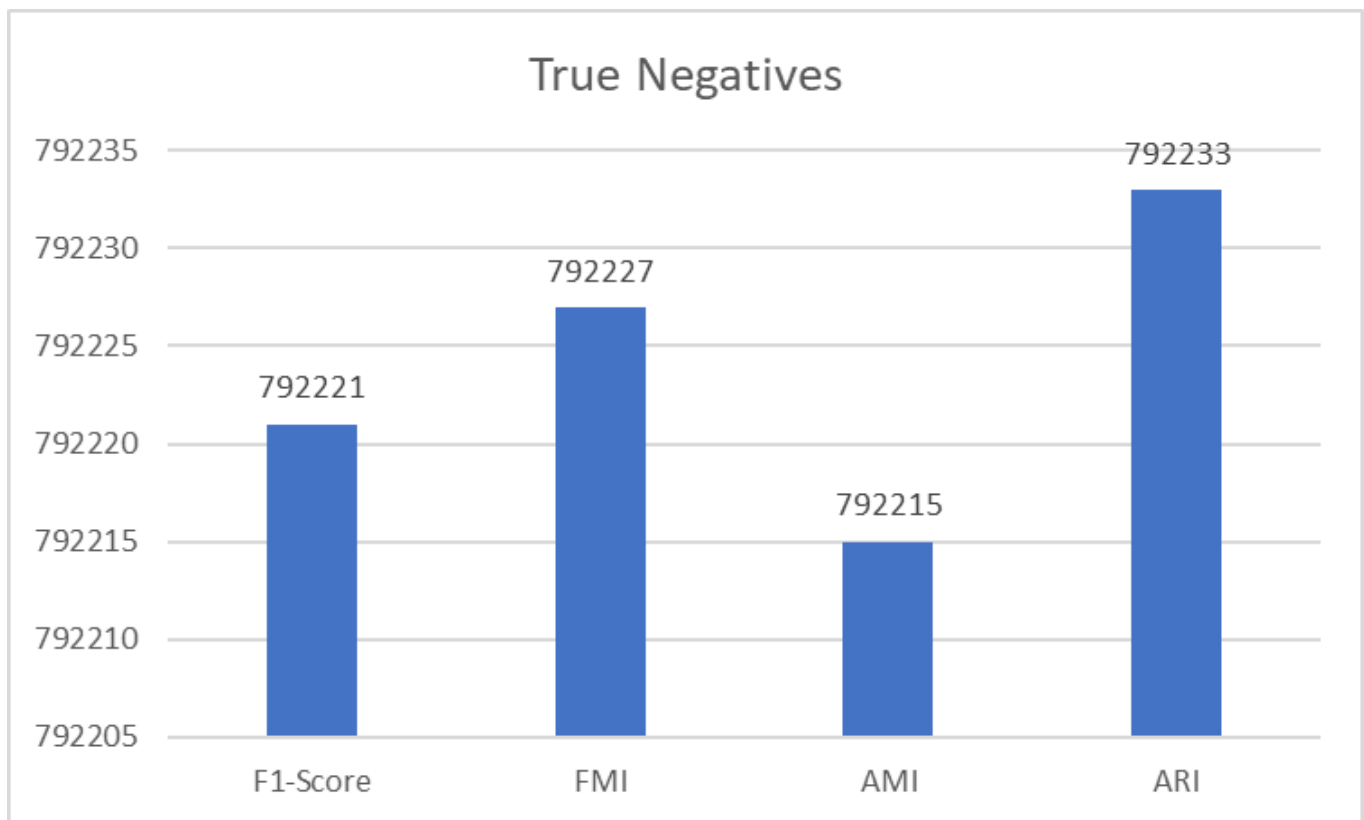
F1-Score geclustert waren (siehe Nidoran-Familie)

 FMI.html

## Vergleich

Betrachten wir die Konfusionsmatrix und die nicht geclusterten Wörter, haben wir eine Entscheidungsgrundlage zur Auswahl der besten Metrik. Die Ergebnisse basieren auf einer Hyperparameteroptimierung mittels Evolutionärem Algorithmus und sind deshalb zufallsbedingt. Neue Durchläufe können zu anderen Ergebnissen führen.





Ziel soll es sein alle True-Werte zu maximieren und alle False-Werte sowie die unclustered Menge zu minimieren. Dies erreicht am besten FMI oder ARI. In drei von fünf Fällen ist aber FMI besser (True Positives, False Negatives, Number Unclustered), weshalb sich für FMI entschieden wird.