

# Entscheidungsbäume

Masterseminar Klassifikation und Clustering

Tim Sockel, Tatiana Enina

Dozent: PD Dr. Stefan Langer

12.12.2022

# Gliederung

1. Definition
2. Der Algorithmus
3. Vor- und Nachteile des Algorithmus
4. Implementierung
5. Literatur

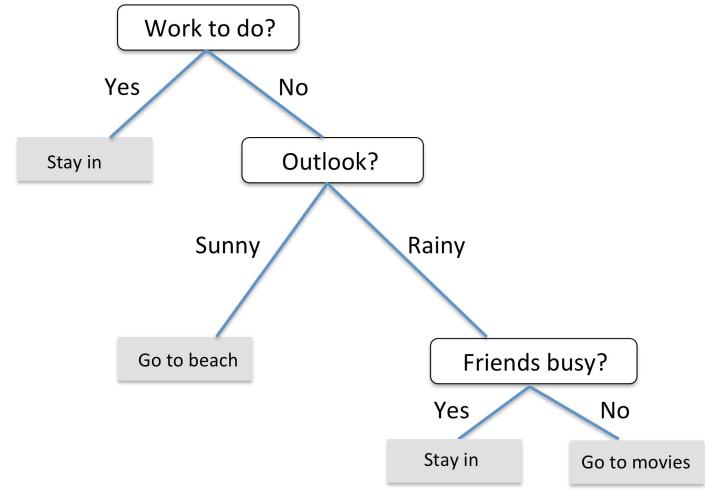
# Definition

**Entscheidungsbaum** – ein Lernalgorithmus, bei dem Klassifikations- oder Regressionsmodelle in Form einer Baumstruktur aufgebaut werden.

⇒ Ein Ansatz zur *mehrstufigen Entscheidungsfindung*: eine komplexe Entscheidung wird in eine Vereinigung mehrerer einfacherer Entscheidungen zerlegt.

## Geschichte:

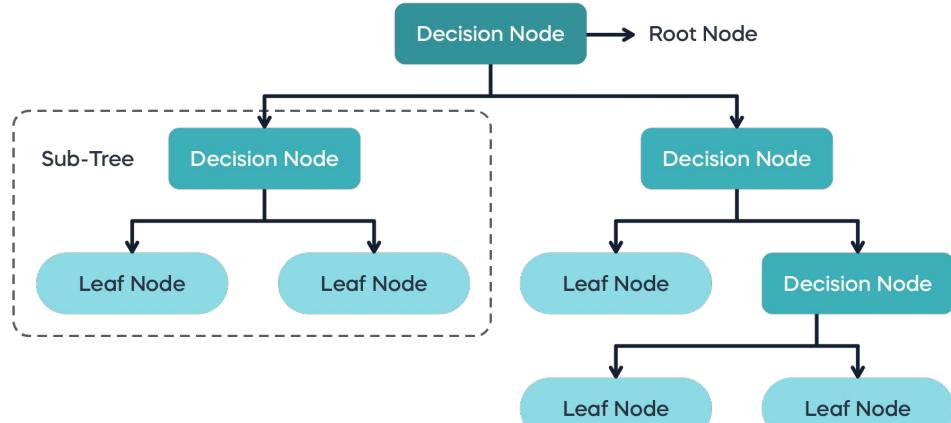
- 1963 – AID project (Morgan, Sonquist): erster Algorithmus mit *Regression*
- 1972 – THAID project (Messenger, Mandell): erster Algorithmus zur *Klassifikation*



Ein Entscheidungsbaum: Beispiel  
(<https://tjmachinelearning.com/lectures/1718/dct/>,  
Zugriff: 12.12.2022)

# Definition

- einfache hierarchische, baumartige Struktur;
- **Bestandteile:**
  - ein Wurzelknoten (eng. *root node*),
  - endliche Menge von internen Knoten (eng. *internal* oder *decision nodes*),
  - endliche Menge von Blattknoten (eng. *leafs* oder *terminal nodes*).
- Prozess der Konstruierung eines Baums
  - Bauminduktion, Baumaufbau (eng. *induction*, *tree building* oder **tree growing**);
- meistens wird die **Top-Down**-Methode verwendet.

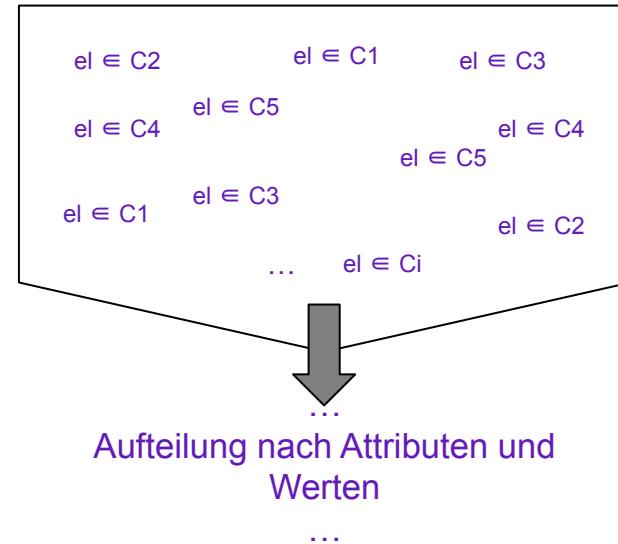


Struktur eines Entscheidungsbaums  
(<https://365datascience.com/tutorials/machine-learning-tutorials/decision-trees/>, Zugriff: 10.12.2022)

# Der Algorithmus

## Trainingsdaten:

- Menge  $S$  aus  $n$  Elementen; jedes Element gehört zu einer Klasse  $(C_1, \dots, C_i)$
- Attribute:  $A_1, \dots, A_j$
- Werte:  $(a_1, \dots, a_p)$



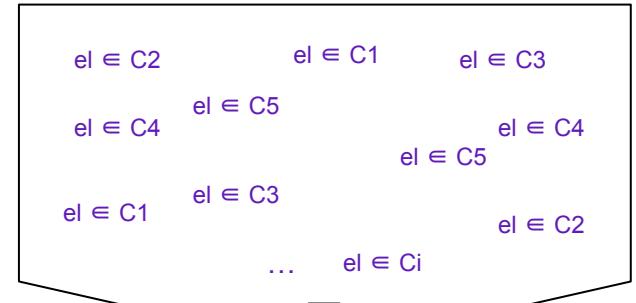
# Der Algorithmus

## Trainingsdaten:

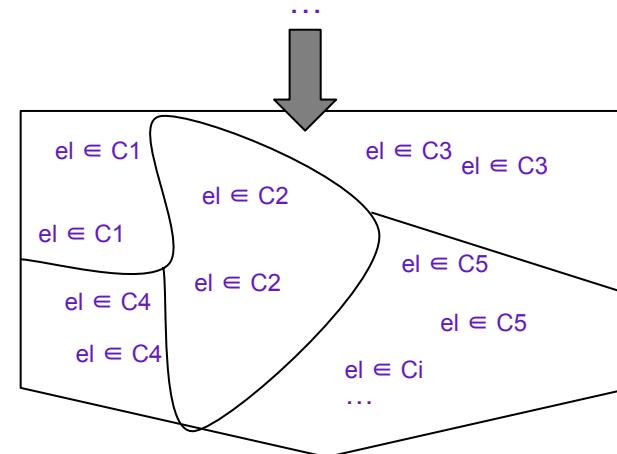
- Menge  $S$  aus  $n$  Elementen; jedes Element gehört zu einer Klasse  $(C_1, \dots, C_i)$
- Attribute:  $A_1, \dots, A_j$
- Werte:  $(a_1, \dots, a_p)$

## Ziel:

Aufgeteilte Menge  $S$ , bei der Elemente jeder Teilmenge eine (möglichst) reine Klasse repräsentieren



Aufteilung nach Attributen und Werten



# Der Algorithmus

## Probleme:

- Generieren (eng. *growing*): Auswahl der Attribute
- Auswahl des Lernstoppkriteriums: Stutzen (eng. *pruning*)

Algorithmischer Top-Down Framework für die Induktion von Entscheidungsbäumen  
(Rokach, Lior & Maimon, 2005: 169)

```
TreeGrowing (S,A,y)
Where:
S - Training Set
A - Input Feature Set
y - Target Feature
Create a new tree T with a single root node.
IF One of the Stopping Criteria is fulfilled THEN
    Mark the root node in T as a leaf with the most
    common value of y in S as a label.
ELSE
    Find a discrete function f(A) of the input
    attributes values such that splitting S
    according to f(A)'s outcomes ( $v_1, \dots, v_n$ ) gains
    the best splitting metric.
    IF best splitting metric > threshold THEN
        Label t with f(A)
        FOR each outcome  $v_i$  of f(A):
            Set Subtree $_i$  = TreeGrowing ( $\sigma_{f(A)=v_i} S, A, y$ ).
            Connect the root node of  $t_r$  to Subtree $_i$  with
            an edge that is labelled as  $v_i$ .
        END FOR
    ELSE
        Mark the root node in T as a leaf with the most
        common value of y in S as a label.
    END IF
END IF
RETURN T

TreePruning (S,T,y)
Where:
S - Training Set
y - Target Feature
T - The tree to be pruned
DO
    Select a node t in T such that pruning it
    maximally improve some evaluation criteria
    IF  $t \neq \emptyset$  THEN T=pruned(T,t)
UNTIL  $t=\emptyset$ 
RETURN T
```

# Der Algorithmus: Generierung eines Baums

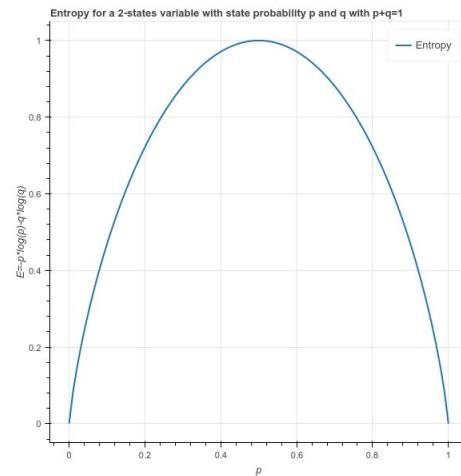
## Entropie:

$$E(S) = - \sum_i p_i \log p_i = - \sum_{i=1}^n \frac{N_i}{N} \log \left( \frac{N_i}{N} \right)$$

$n$  – die Anzahl der Klassen,  $N_i$  – die Anzahl der Elemente einer Klasse  $i$ ,  $N$  – die Gesamtzahl der Elemente.

$E(S)=0$ , falls  $p_i=1$  für eine Klasse  $i$  in  $S$

$E(S)=1$  für 2 Klassen, falls  $p_1, p_2 = 1/2$



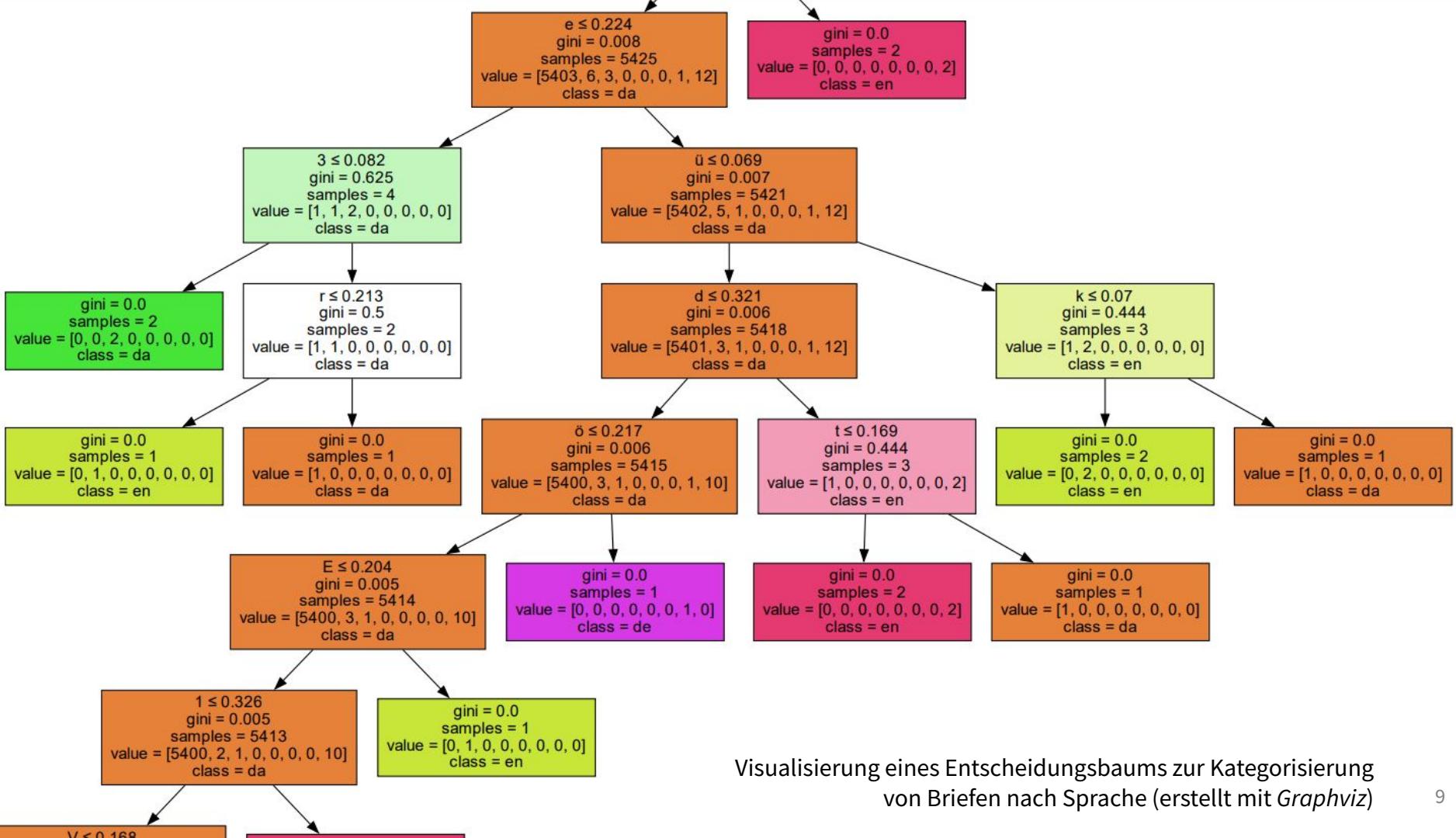
## Gini-Index:

$$\text{Gini}(S) = 1 - \sum_{i=1}^n p_i^2$$

$n$  - die Anzahl der Klassen,  
 $p_i$  - die Wahrscheinlichkeit einer Klasse i

kleiner Gini-Index  $\Leftrightarrow$  geringe Unreinheit  
großer Gini-Index  $\Leftrightarrow$  hohe Unreinheit

Grafik für  $E = p \log(p) - q \log(q)$ ,  $q+p=1$   
(<https://bricaud.github.io/personal-blog/entropy-in-decision-trees/>, Zugriff: 8.12.22)



# Der Algorithmus: Beispiele von Klassifikatoren

- ID3: Entropie, Informationsgewinn (eng. *Information gain*);
- C4.5: Informationsgewinn;
- CART (Classification and Regression Tree): Gini Index.

# Der Algorithmus: Abbruch des Lernprozesses

## ***Pre-Pruning:***

- Baumtiefenbegrenzung (`max_depth`)
- Beste Mindestanzahl der Elemente bei internen (Split-)Knoten (`min_samples_split`)
- Beste Mindestanzahl der Elemente bei Endknoten (`min_samples_leaf`)

## ***Post-Pruning (z. B. ‘Minimal Cost-Complexity Pruning’):***

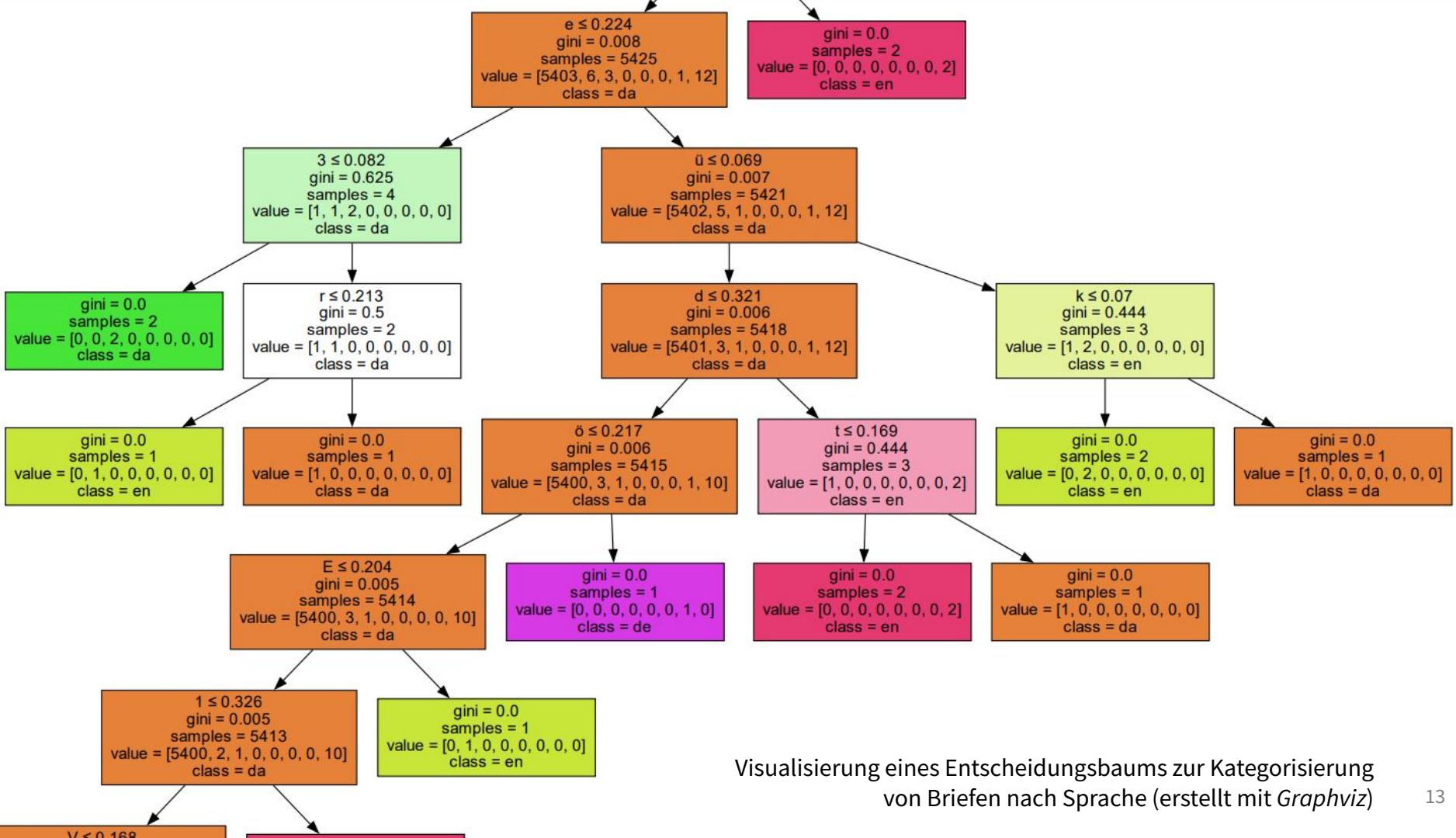
- Entscheidungsbaum wird zunächst mit voller Baumtiefe generiert
- `ccp_alpha`-Wert kann darauf basierend bspw. mittels einer Funktion berechnet oder geschätzt werden
- Entscheidungsbaum kann mittels verschiedenen `ccp_alpha`-Werten trainiert werden und die Evaluationsergebnisse können verglichen werden

# Vorteile

- Weniger Datenvorbereitung während der Vorverarbeitung;
- Merkmalsskalierung nicht erforderlich;
- schneller Trainingsprozess;
- können automatisch mit fehlenden Werten umgehen;
- leicht zu verstehen und visualisieren.

# Nachteile

- Hinzufügen neuer Datenpunkte führt zur Neugenerierung des Gesamtbaums (Knoten müssen neu berechnet werden);
- werden durch Noise beeinflusst und instabil;
- das Ergebnis ist teils von der Reihenfolge der eingegebenen Texte abhängig;
- für große Datasets nicht geeignet, Gefahr von Overfitting.  
→ Lösung: Stutzen (*pruning*).

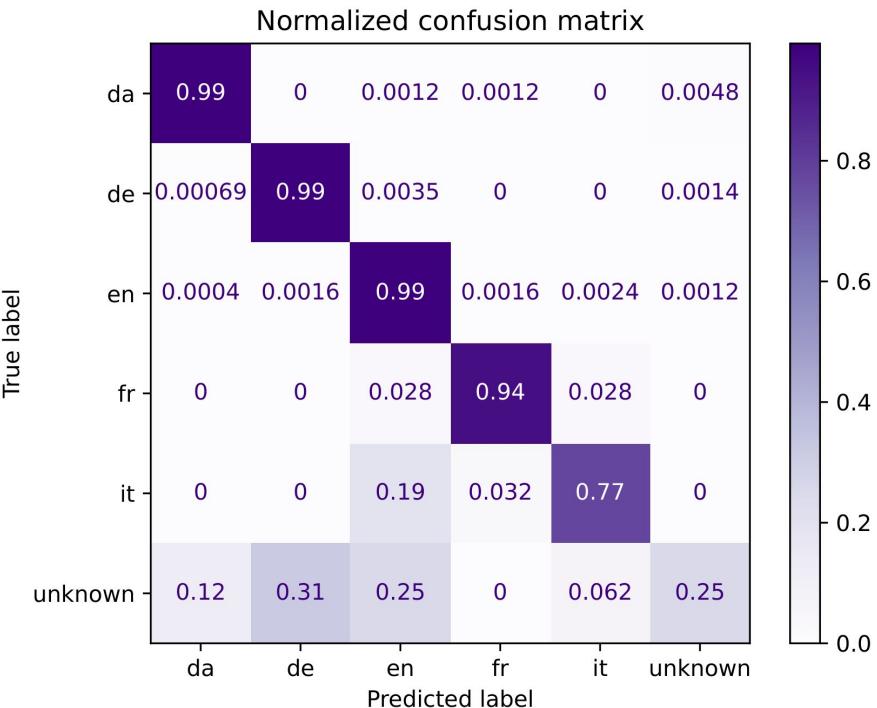


# Implementierung

Briefsammlung – Kategorisierung  
nach Sprache

Methode: Zeichenbasierte TF-IDF-  
Matrix; ohne sog. *Lowercasing*

	Precision	Recall	F1	Support
da	1.00	0.99	0.99	838
de	0.99	0.99	0.99	1443
en	0.99	0.99	0.99	2517
fr	0.85	0.94	0.89	36
it	0.75	0.77	0.76	31
unk.	0.31	0.25	0.28	16
acc.			0.99	4881
macro	0.81	0.82	0.82	4881
weight.	0.99	0.99	0.99	4881

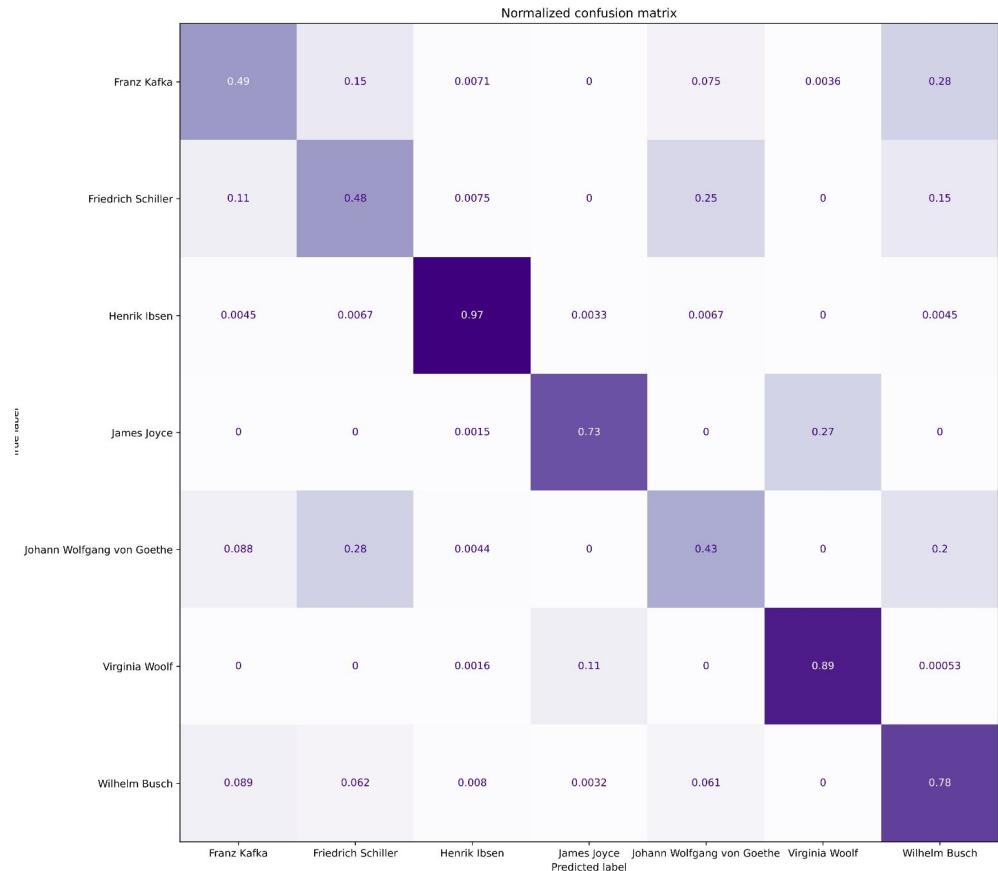


# Implementierung

Briefsammlung – Kategorisierung  
nach Autorinnen und Autoren

Methode: Wortbasierte TF-IDF-Matrix

	Prec.	Recall	F1	Support
Kafka	0.56	0.49	0.52	280
Schiller	0.46	0.48	0.47	266
Ibsen	0.98	0.97	0.98	897
Joyce	0.69	0.73	0.71	682
Goethe	0.43	0.43	0.43	228
Woolf	0.90	0.89	0.89	1901
Busch	0.75	0.78	0.76	627
acc.			0.80	4881
macro	0.68	0.68	0.68	4881
weight.	0.80	0.80	0.80	4881

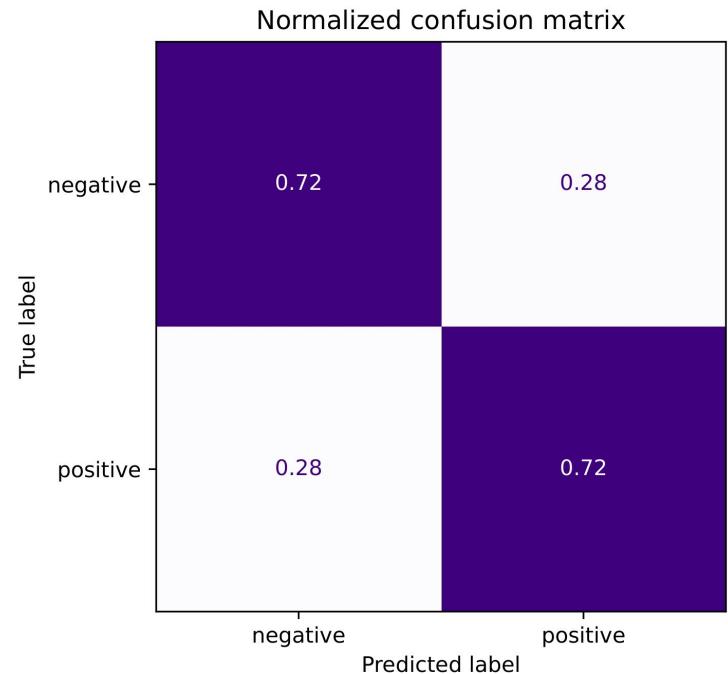


# Implementierung

IMDb-Filmrezensionen – Sentiment-Analyse

Methode: TF-IDF-Matrix aus Uni- und Bigrammen;  
45.000 häufigste Merkmale

	Precision	Recall	F1	Support
Negativ	0.72	0.72	0.72	4985
Positiv	0.72	0.72	0.72	5015
acc.			0.72	10000
macro	0.72	0.72	0.72	10000
weight.	0.72	0.72	0.72	10000



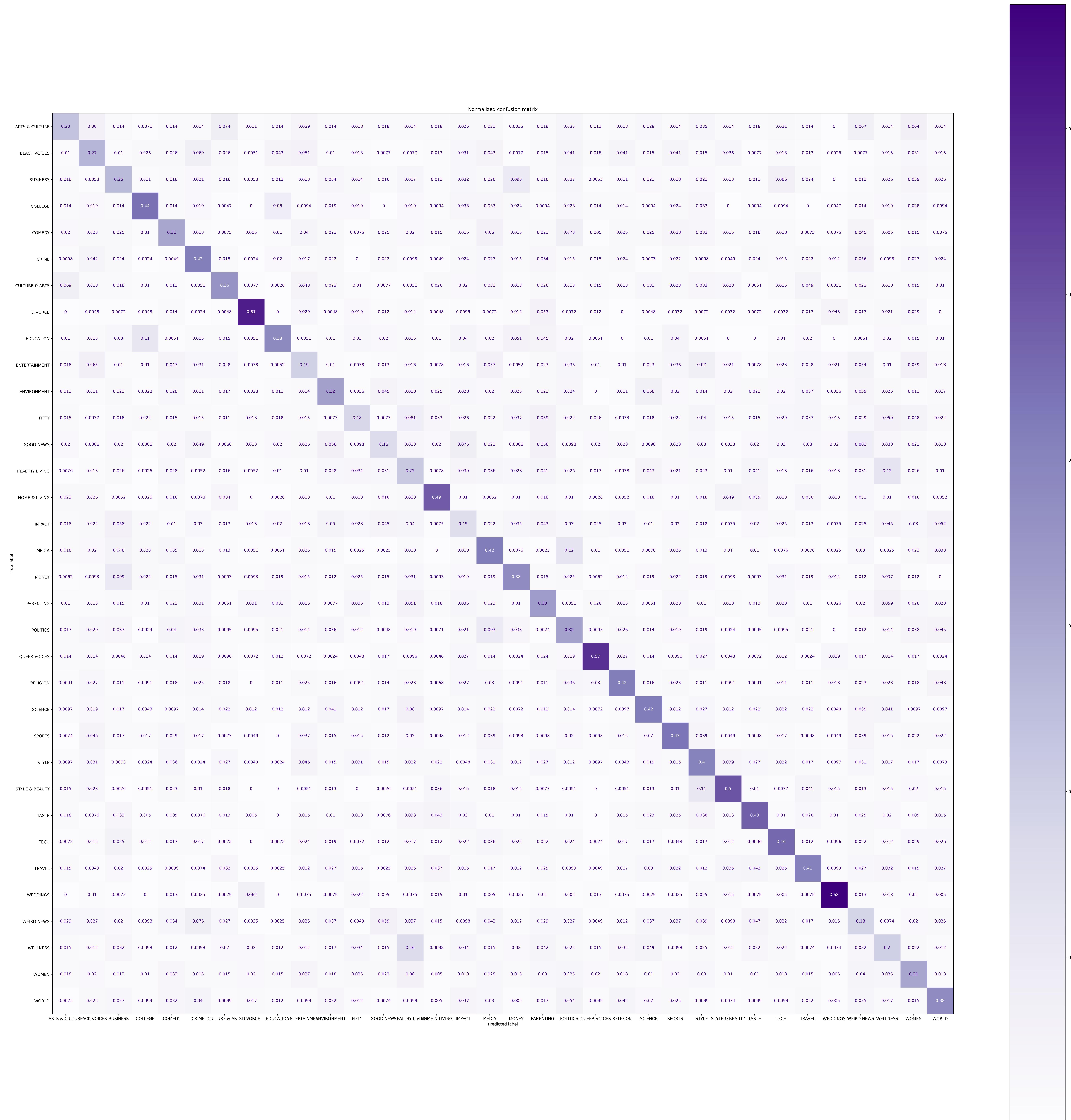
# Implementierung

Huffington Post – Kategorisierung  
nach Ressorts

Methode: Wortbasierte TF-IDF-Matrix;  
ohne englische Stoppwörter

- weitere Optimierung mittels  
*Minimal Cost-Complexity Pruning*  
(`ccp_alpha = 0.00006`)
- ⇒ Accuracy von 36 %  
konnte auf 41 % verbessert  
werden

	Prec.	Recall	F1	Support
Arts & Culture	0.27	0.23	<b>0.25</b>	282
Black Voices	0.29	0.27	<b>0.28</b>	392
Business	0.26	0.26	<b>0.26</b>	380
College	0.40	0.44	<b>0.42</b>	212
Comedy	0.32	0.31	<b>0.32</b>	399
Crime	0.40	0.42	<b>0.41</b>	409
Culture & Arts	0.40	0.36	<b>0.38</b>	391
Divorce	0.68	0.61	<b>0.64</b>	419
Education	0.33	0.38	<b>0.36</b>	198
Entertainment	0.22	0.19	<b>0.21</b>	387
Environment	0.31	0.32	<b>0.32</b>	355
Fifty	0.20	0.18	<b>0.19</b>	273
Good News	0.19	0.16	<b>0.17</b>	305
Healthy Living	0.19	0.22	<b>0.20</b>	386
Home & Living	0.51	0.49	<b>0.50</b>	386
Impact	0.17	0.15	<b>0.16</b>	400
Media	0.32	0.42	<b>0.36</b>	395
Money	0.36	0.38	<b>0.37</b>	324
Parenting	0.30	0.33	<b>0.32</b>	391
Politics	0.29	0.32	<b>0.30</b>	420
Queer Voices	0.63	0.57	<b>0.60</b>	415
Religion	0.48	0.42	<b>0.45</b>	440
Science	0.41	0.42	<b>0.42</b>	414
Sports	0.41	0.43	<b>0.42</b>	410
Style	0.34	0.40	<b>0.37</b>	413
Style & Beauty	0.52	0.50	<b>0.51</b>	391
Taste	0.47	0.48	<b>0.48</b>	397
Tech	0.46	0.46	<b>0.46</b>	416
Travel	0.41	0.41	<b>0.41</b>	405
Weddings	0.68	0.68	<b>0.68</b>	400
Weird News	0.17	0.18	<b>0.18</b>	408
Wellness	0.21	0.20	<b>0.20</b>	407
Women	0.29	0.31	<b>0.30</b>	400
World	0.40	0.38	<b>0.39</b>	404
<b>accuracy</b>			<b>0.36</b>	12824
<b>macro avg</b>	0.36	0.36	<b>0.36</b>	12824
<b>weighted avg</b>	0.37	0.36	<b>0.37</b>	12824



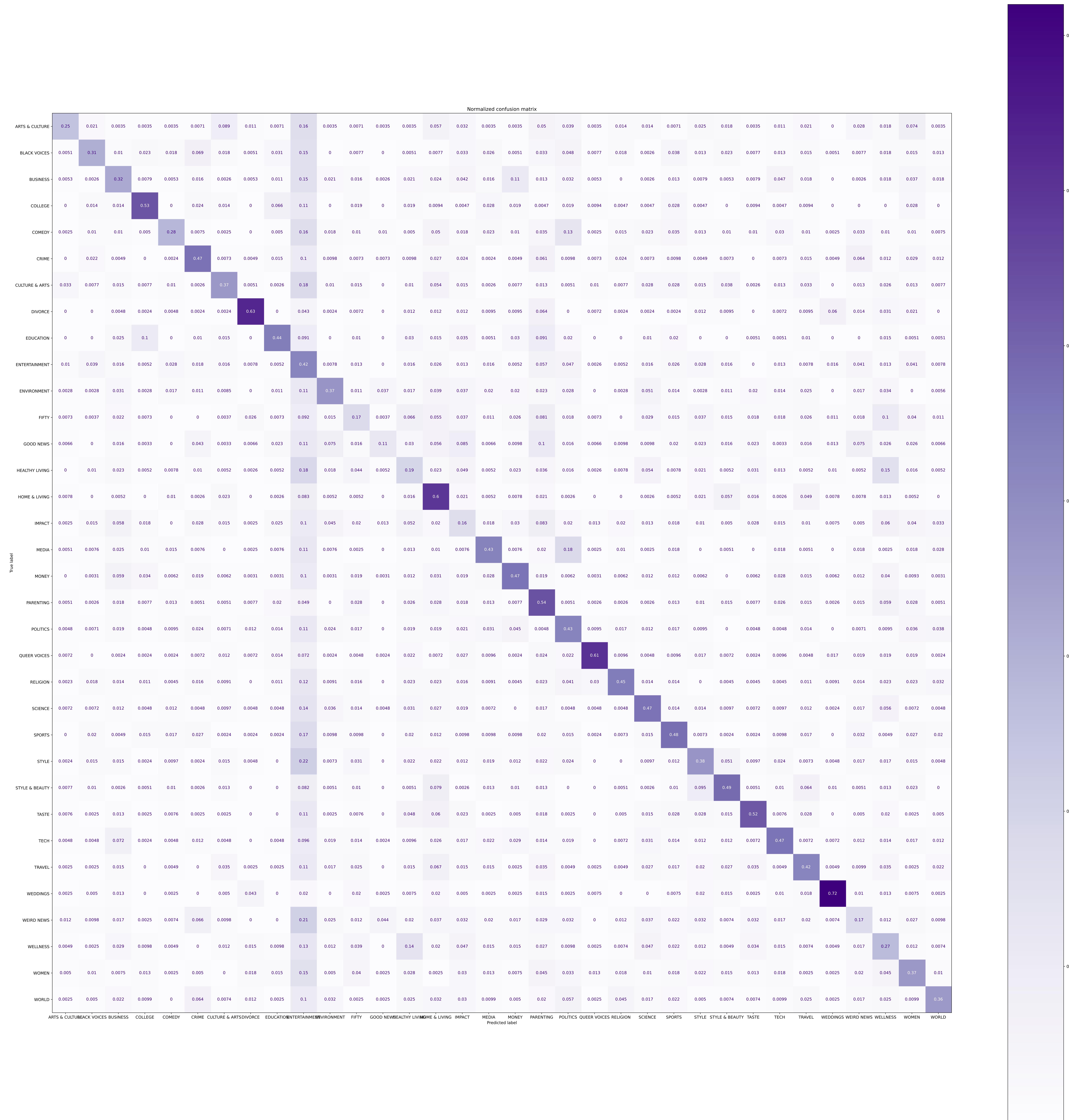
# Implementierung

Huffington Post – Kategorisierung  
nach Ressorts

Methode: Wortbasierte TF-IDF-Matrix;  
ohne englische Stoppwörter

- weitere Optimierung mittels  
*Minimal Cost-Complexity Pruning*  
(`ccp_alpha = 0.00006`)
- ⇒ Accuracy von 36 %  
konnte auf 41 % verbessert  
werden

	Prec.	Recall	F1	Support
Arts & Culture	0.52	0.25	<b>0.34</b>	282
Black Voices	0.53	0.31	<b>0.39</b>	392
Business	0.34	0.32	<b>0.33</b>	380
College	0.51	0.53	<b>0.52</b>	212
Comedy	0.54	0.28	<b>0.37</b>	399
Crime	0.50	0.47	<b>0.48</b>	409
Culture & Arts	0.52	0.37	<b>0.43</b>	391
Divorce	0.77	0.63	<b>0.69</b>	419
Education	0.43	0.44	<b>0.44</b>	198
Entertainment	0.10	0.42	<b>0.16</b>	387
Environment	0.44	0.37	<b>0.40</b>	355
Fifty	0.19	0.17	<b>0.18</b>	273
Good News	0.38	0.11	<b>0.18</b>	305
Healthy Living	0.20	0.19	<b>0.20</b>	386
Home & Living	0.38	0.60	<b>0.47</b>	386
Impact	0.18	0.16	<b>0.17</b>	400
Media	0.51	0.43	<b>0.47</b>	395
Money	0.45	0.47	<b>0.46</b>	324
Parenting	0.34	0.54	<b>0.42</b>	391
Politics	0.34	0.43	<b>0.38</b>	420
Queer Voices	0.80	0.61	<b>0.69</b>	415
Religion	0.63	0.45	<b>0.52</b>	440
Science	0.49	0.47	<b>0.48</b>	414
Sports	0.49	0.48	<b>0.48</b>	410
Style	0.43	0.38	<b>0.40</b>	413
Style & Beauty	0.53	0.49	<b>0.51</b>	391
Taste	0.60	0.52	<b>0.56</b>	397
Tech	0.54	0.47	<b>0.50</b>	416
Travel	0.45	0.42	<b>0.44</b>	405
Weddings	0.77	0.72	<b>0.74</b>	400
Weird News	0.24	0.17	<b>0.20</b>	408
Wellness	0.24	0.27	<b>0.25</b>	407
Women	0.37	0.37	<b>0.37</b>	400
World	0.52	0.36	<b>0.43</b>	404
<b>accuracy</b>			<b>0.41</b>	12824
<b>macro avg</b>	0.45	0.40	<b>0.41</b>	12824
<b>weighted avg</b>	0.45	0.41	<b>0.42</b>	12824



# Literatur

- Loh, W. 2014. Fifty Years of Classification and Regression Trees. *International Statistical Review*. International Statistical Institute (ISI). 82, 3, 329–348.
- Manning, C. and Schütze, H. 2001. *Foundations of Statistical Natural Language Processing*. MIT Press, Cambridge MA.
- Murthy, S. 1997. Automatic construction of decision trees from data: A Multi-disciplinary survey. Kluwer Academic Publishers, Boston.
- Rokach, L. and Maimon, O. 2005. Decision Trees. *The Data Mining and Knowledge Discovery Handbook*, 165-192. Springer, Boston, MA.
- Pedregosa et al. 2011. Scikit-learn: Machine Learning in Python. *JMLR* 12. 2825-2830.
- Ellson J. et al. 2004. *Graphviz and Dynagraph – Static and Dynamic Graph Drawing Tools*.