

---

---

# Wellness Program *For employees*

By Azin Faghihi

---

---

## Problem statement:

How should we tailor the wellness package to each employee based on questionnaires /surveys collected from the employees.

*Possible Role:* Consultant at [SomeCaringStartup.com](https://www.somecaringstartup.com)

*Data:* reddit data r/[Meditation](https://www.reddit.com/r/Meditation), r/[yoga](https://www.reddit.com/r/yoga)

---

---

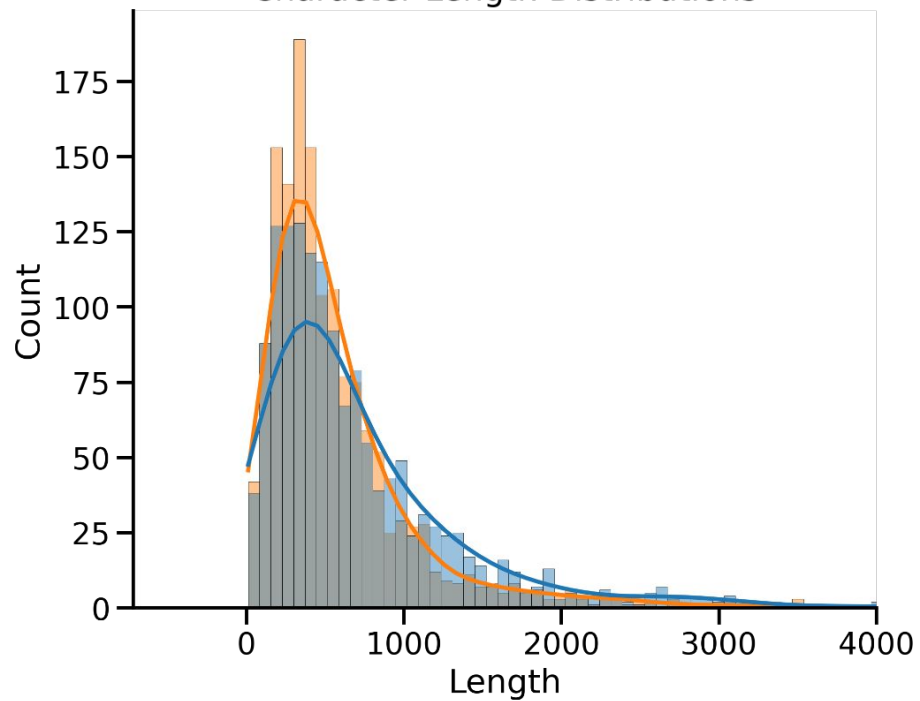
---

**General look at the data**

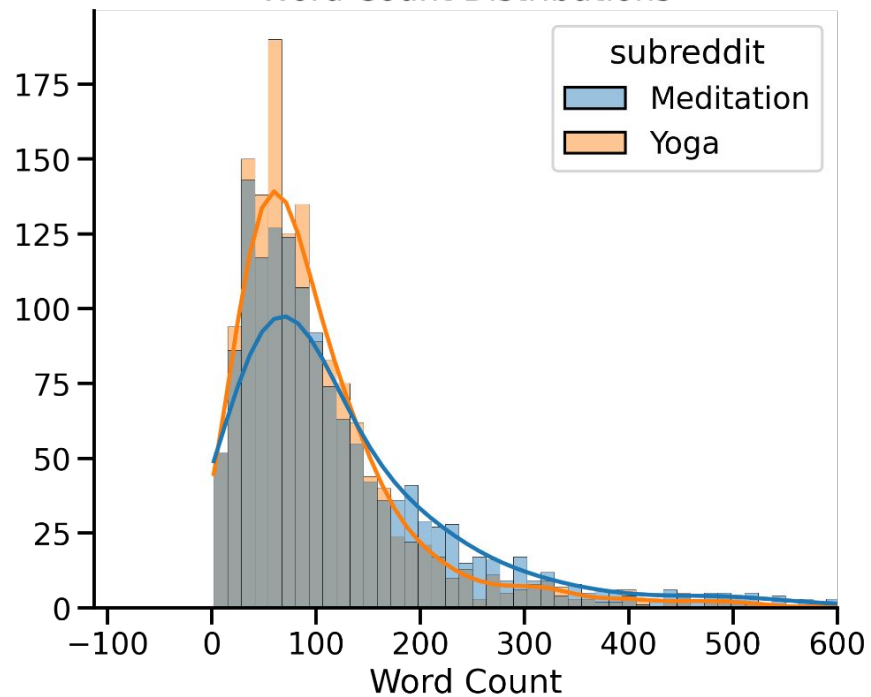
---

## Title+Content

### Character Length Distributions



### Word Count Distributions

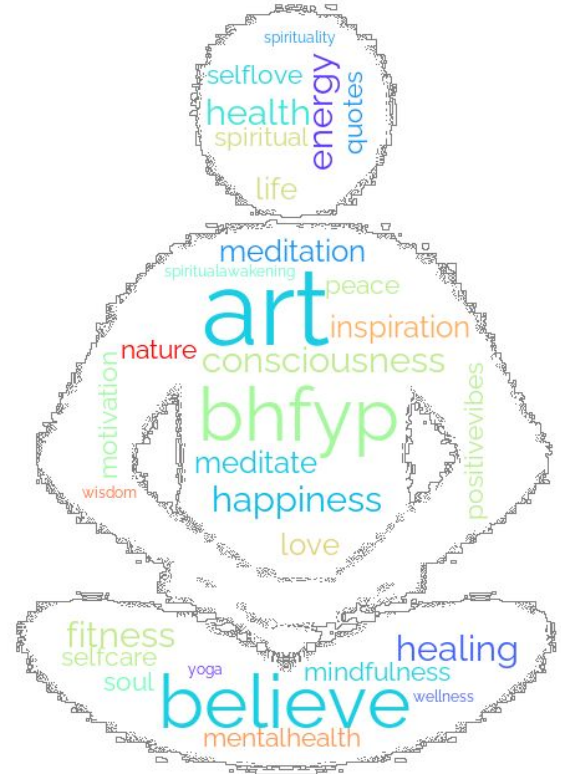
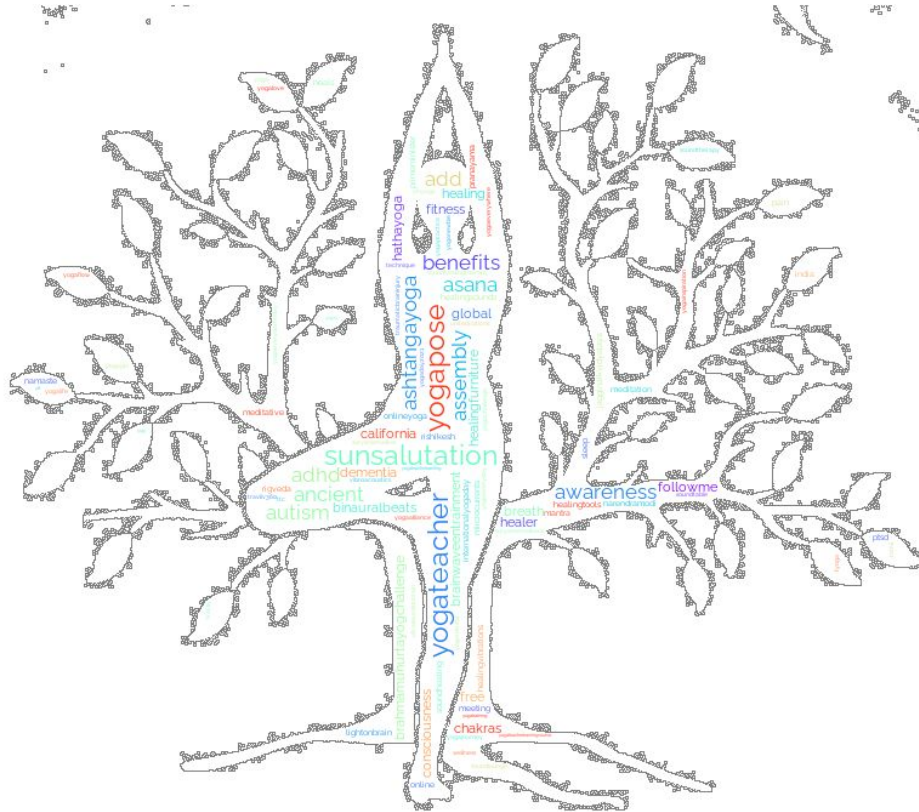


**# rows   # hashtags**

**subreddit**

<b>Meditation</b>	1480	33
<b>Yoga</b>	1480	99

# hashtags



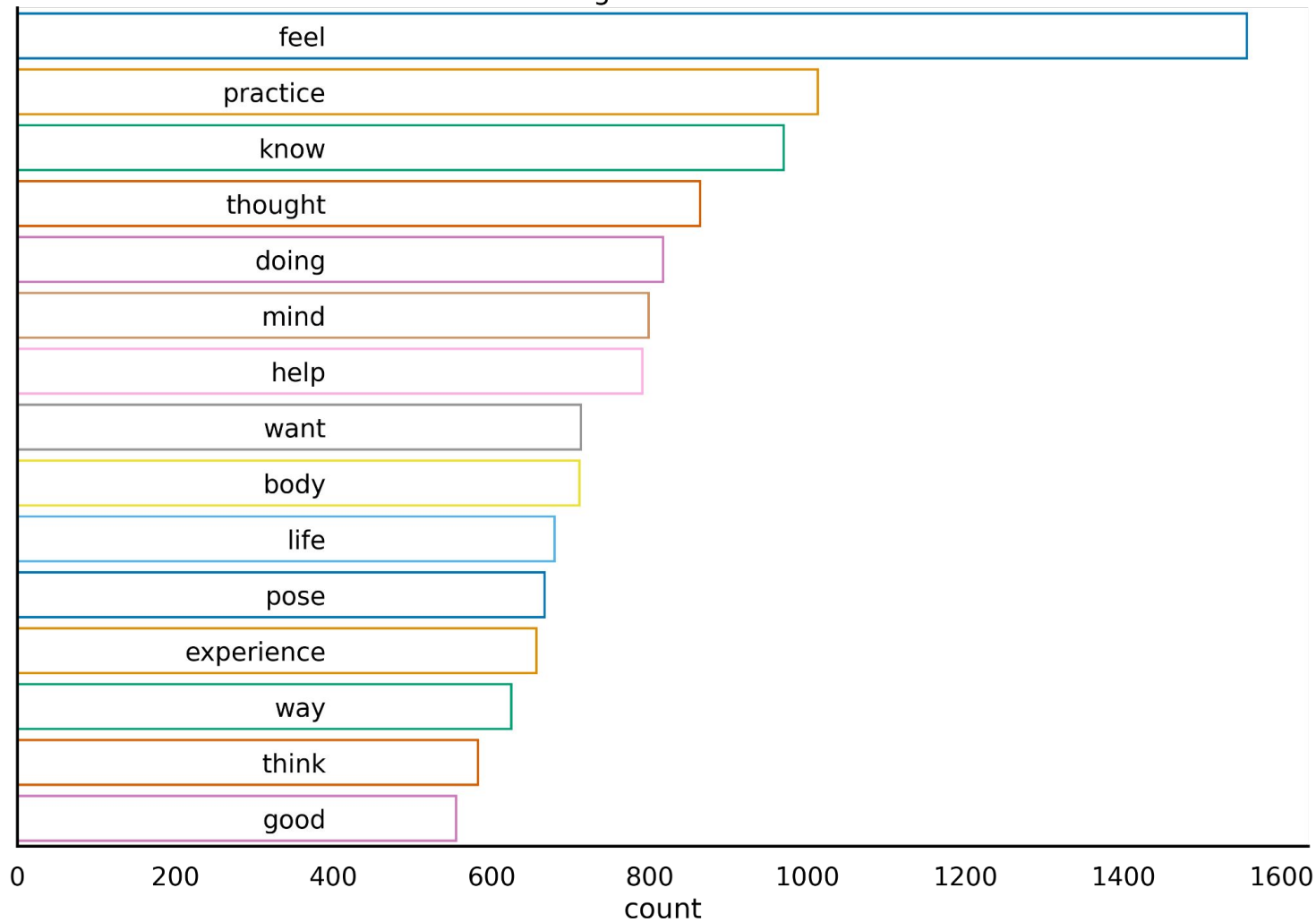
---

---

Which *words* appear commonly  
in these two subreddits?

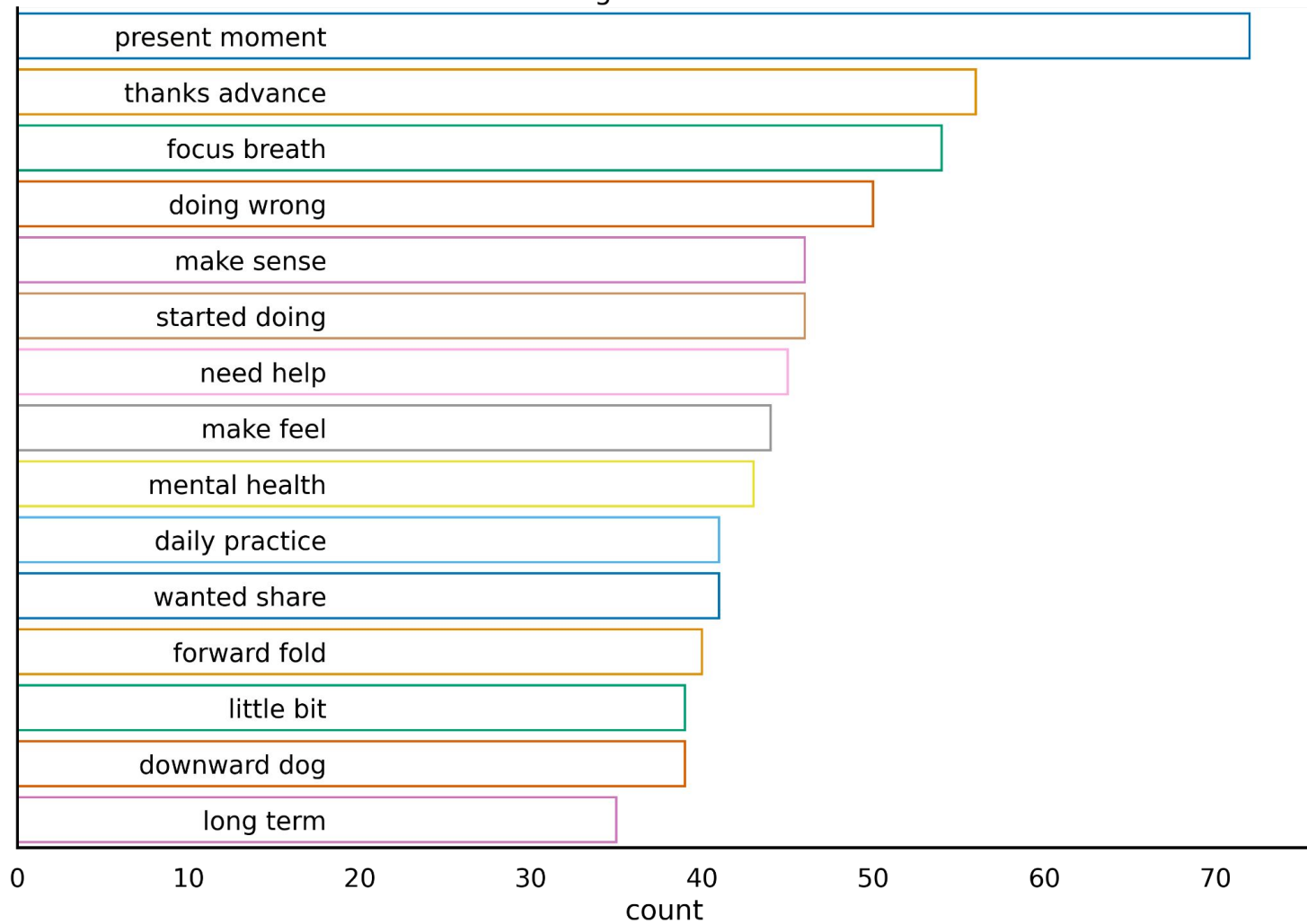
---

1-gram Words

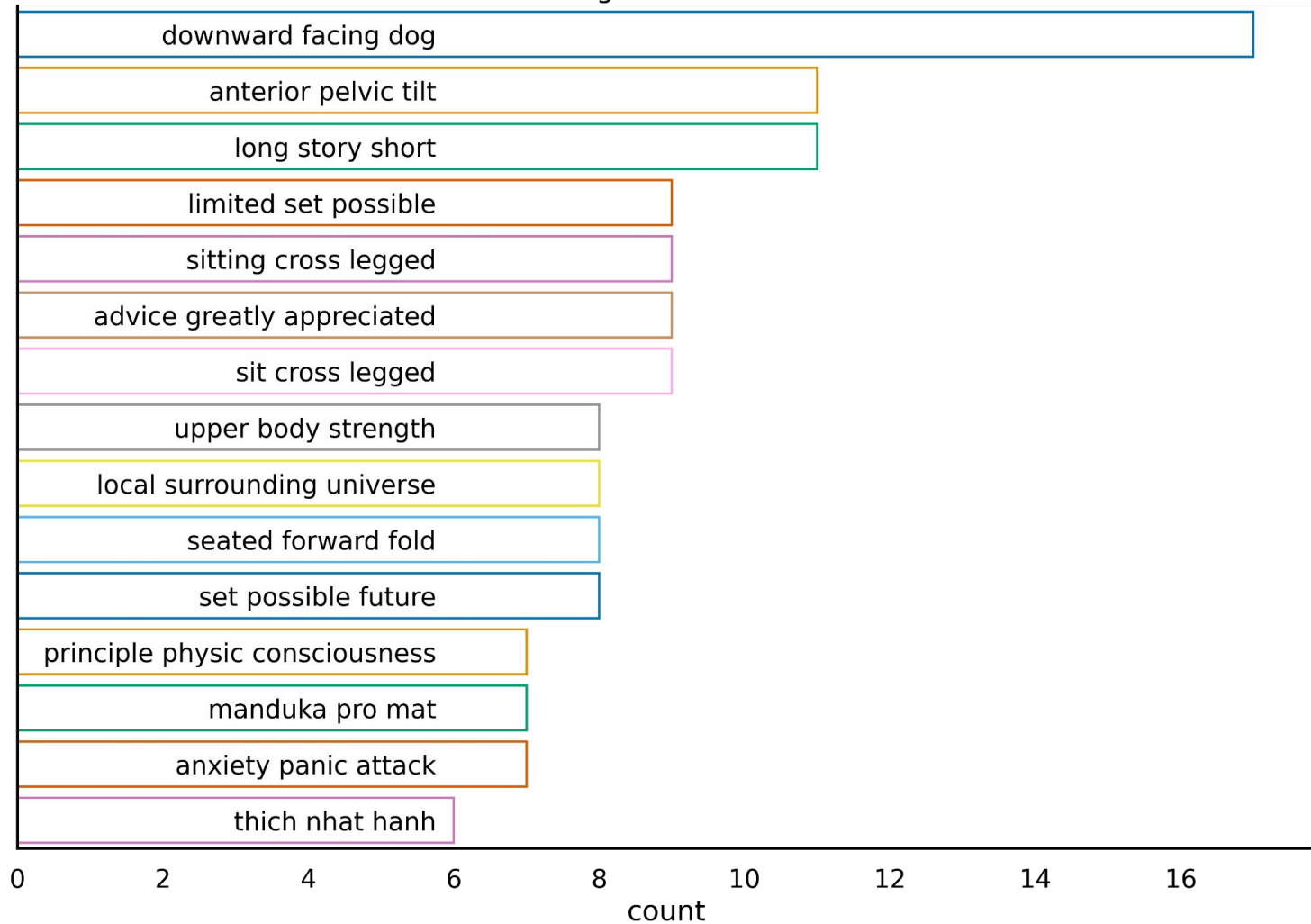




## 2-gram Words



### 3-gram Words



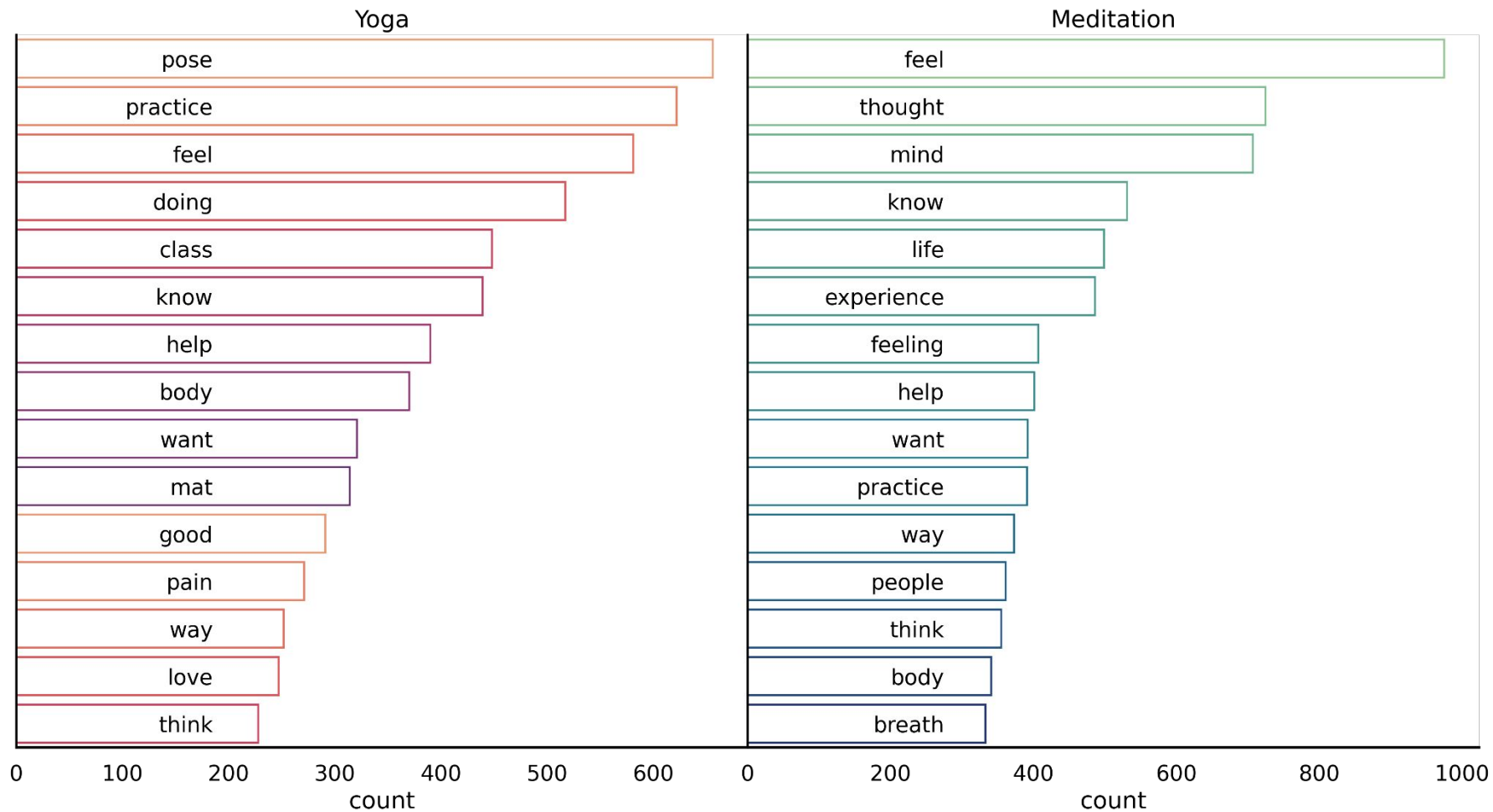
---

---

Which *words* appear in *Yoga & Meditation* each?

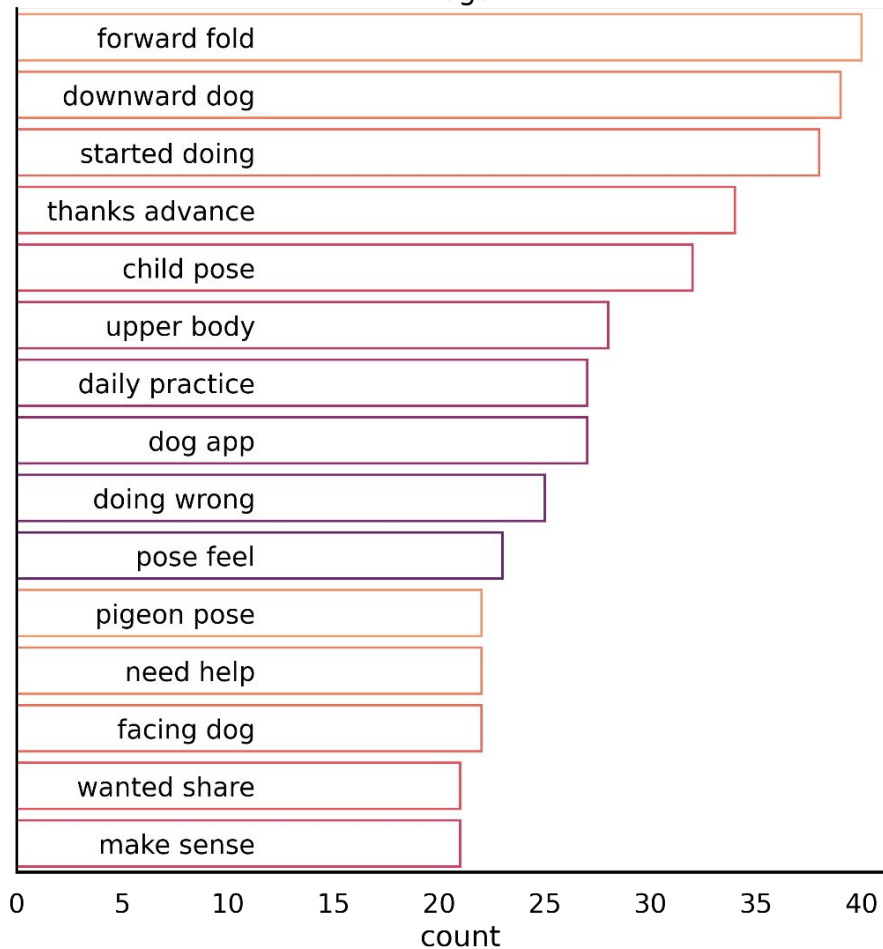
---

# 1-gram Words

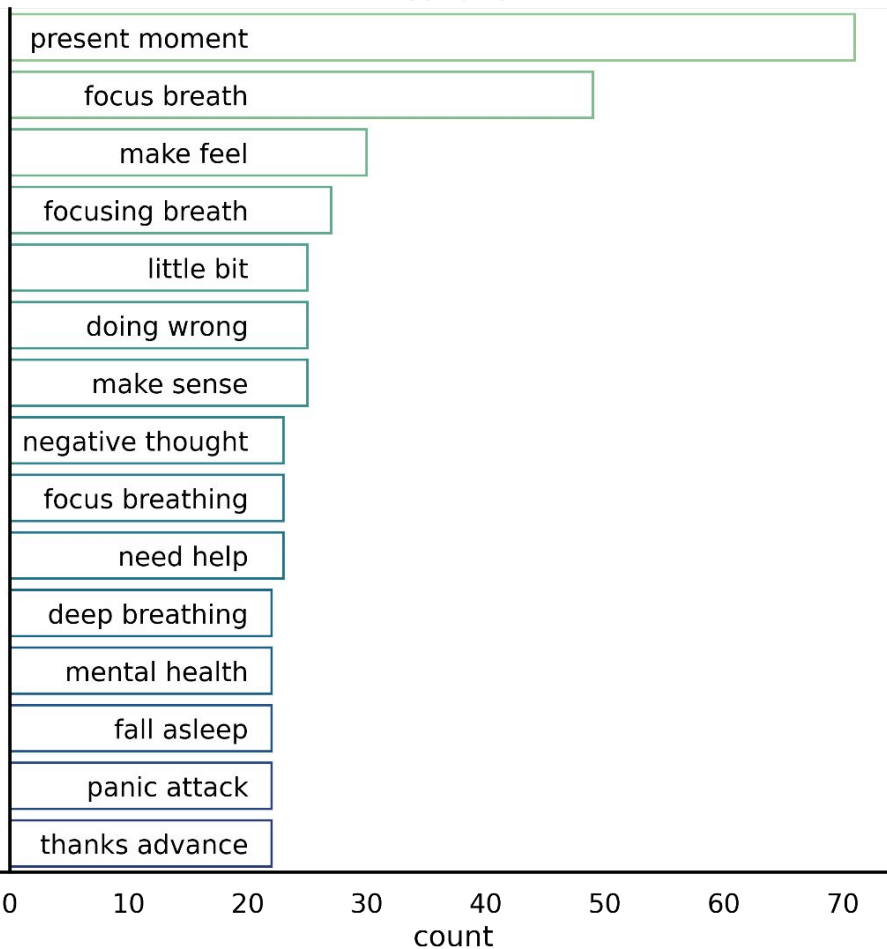


# 2-gram Words

Yoga



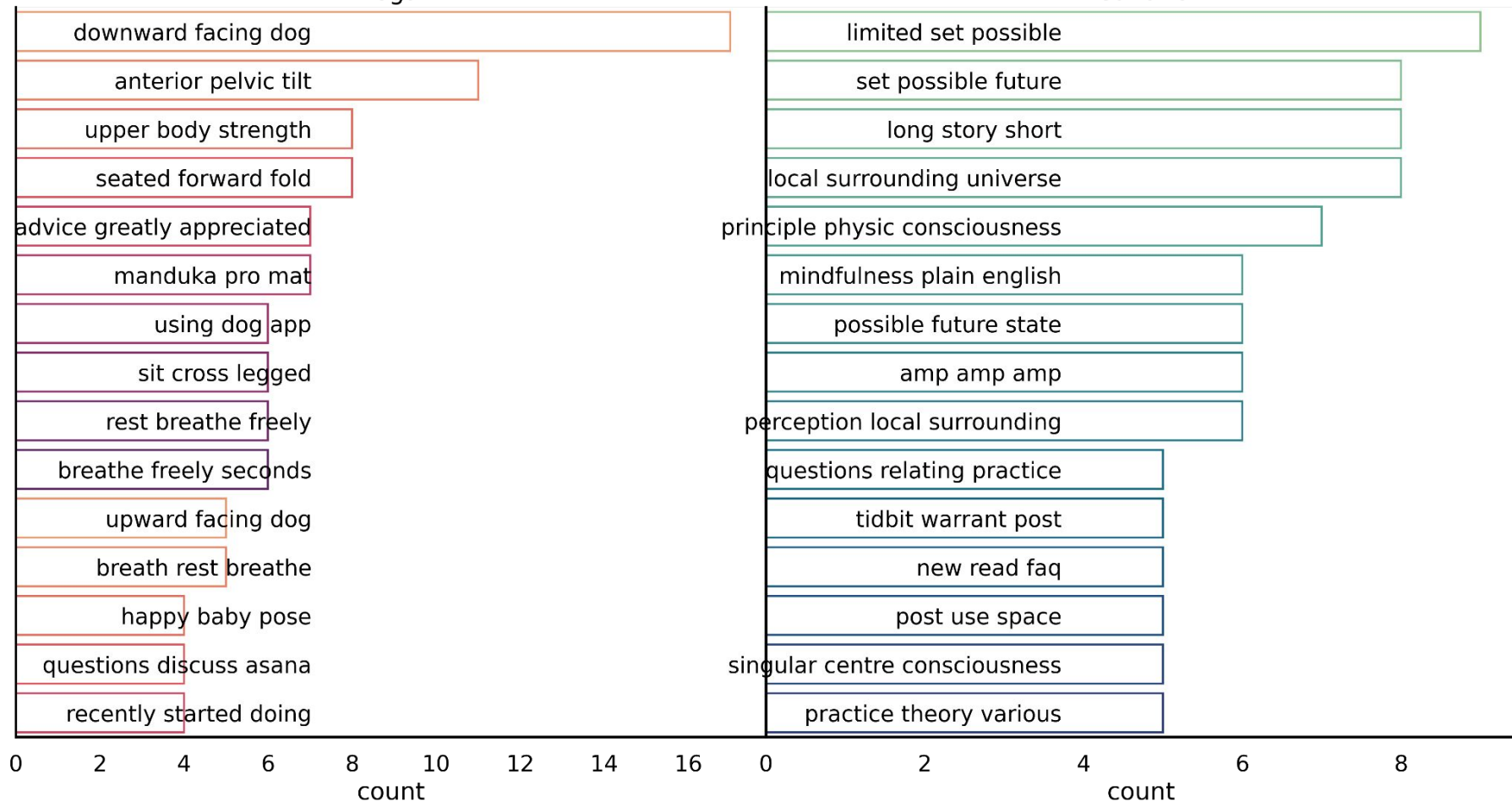
Meditation



# 3-gram Words

Yoga

Meditation



---

---

# Modeling

---

		Validation Score	Train Score	Test Score	Best Score (GS)	Train Score (GS)	Test Score (GS)	F1 - score	Recall (Sensitivity)	Specificity (True Negative Rate)	Precision	Accuracy
Estimator	Transformer											
MultinomialNB	TfidfVectorizer	0.881	0.991	0.884	0.881	0.941	0.883	0.885	0.905	0.860	0.866	0.883
	CountVectorizer	0.890	0.994	0.880	0.883	0.924	0.875	0.875	0.872	0.878	0.878	0.875
RandomForestClassifier	TfidfVectorizer	0.841	0.999	0.855	0.866	0.999	0.855	0.861	0.896	0.813	0.827	0.855
DecisionTreeClassifier	CountVectorizer	0.774	0.999	0.787	0.778	0.842	0.785	0.803	0.876	0.694	0.741	0.785
	TfidfVectorizer	0.752	0.999	0.757	0.760	0.820	0.762	0.788	0.883	0.642	0.711	0.762

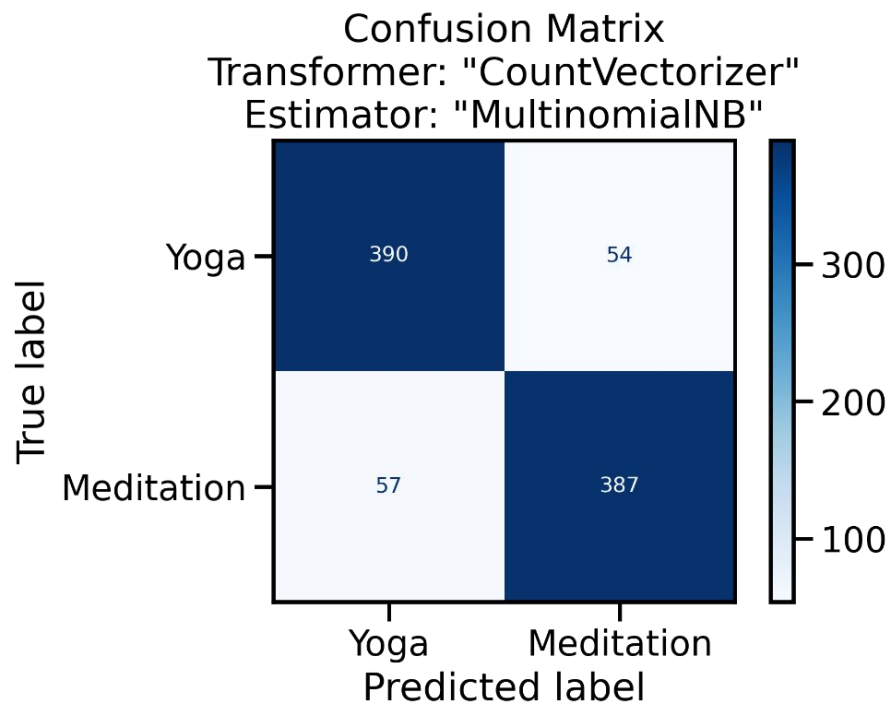


---

---

Modeling: *Multinomial Naive  
Bayes*

---



Accuracy: 88%

Precision: 88%

F1 - score: 87%

Recall (Sensitivity): 87%

Specificity (True Negative Rate): 88%

max\_df: 0.9

max\_features: 3000

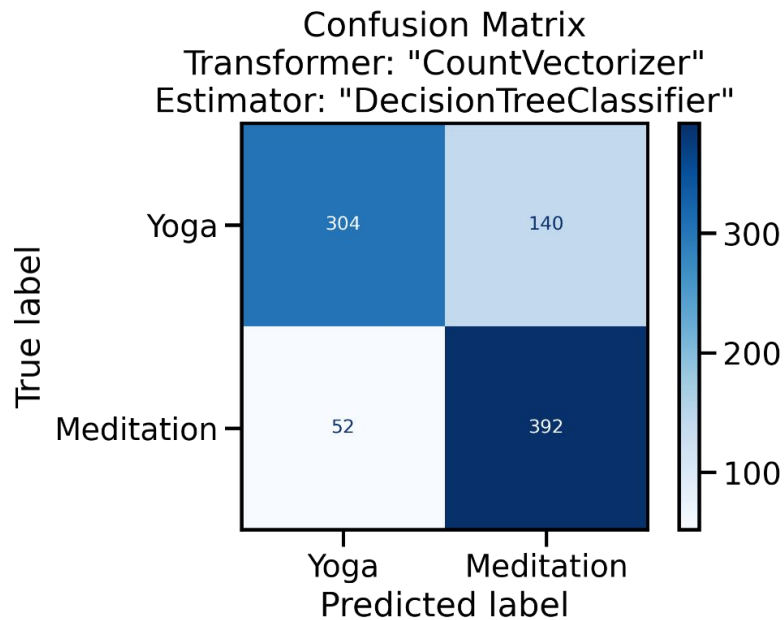
min\_df: 2

ngram\_range: (1, 1)

---

Modeling: *Decision Tree Classifier*

---



Accuracy: 78%

Precision: 74%

F1 - score: 80%

Recall (Sensitivity): 88%

Specificity (True Negative Rate): 68%

max\_depth: 30

min\_samples\_leaf: 5

min\_samples\_split: 30

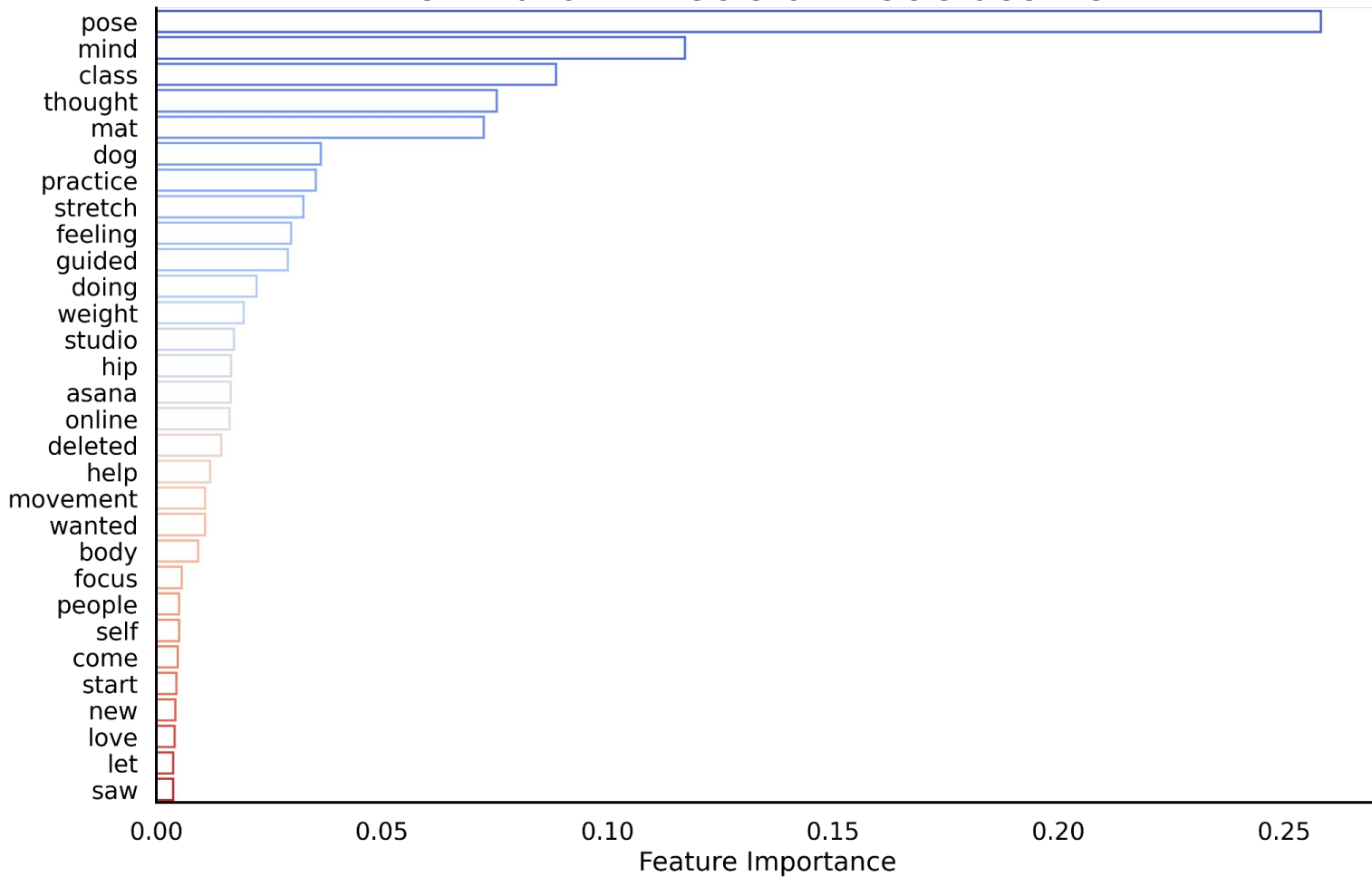
max\_df: 0.9

min\_df: 2

ngram\_range: (1, 2)

# Transformer: "TfidfVectorizer"

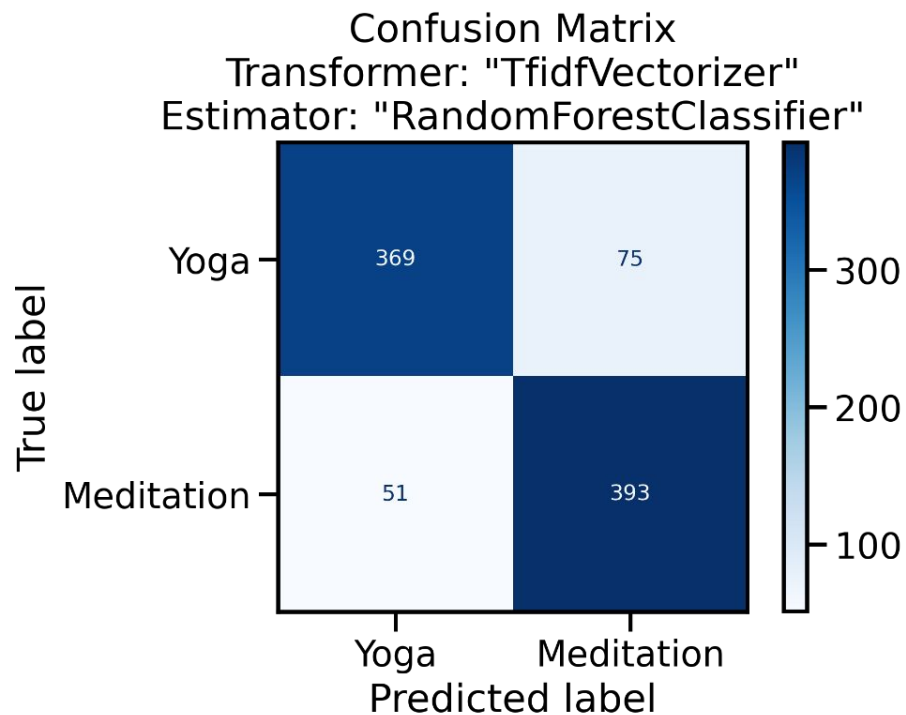
## Estimator: "DecisionTreeClassifier"



---

Modeling: *Random Forest Classifier*

---



Accuracy: 86%

Precision: 84%

F1 - score: 86%

Recall (Sensitivity): 89%

Specificity (True Negative Rate): 83%

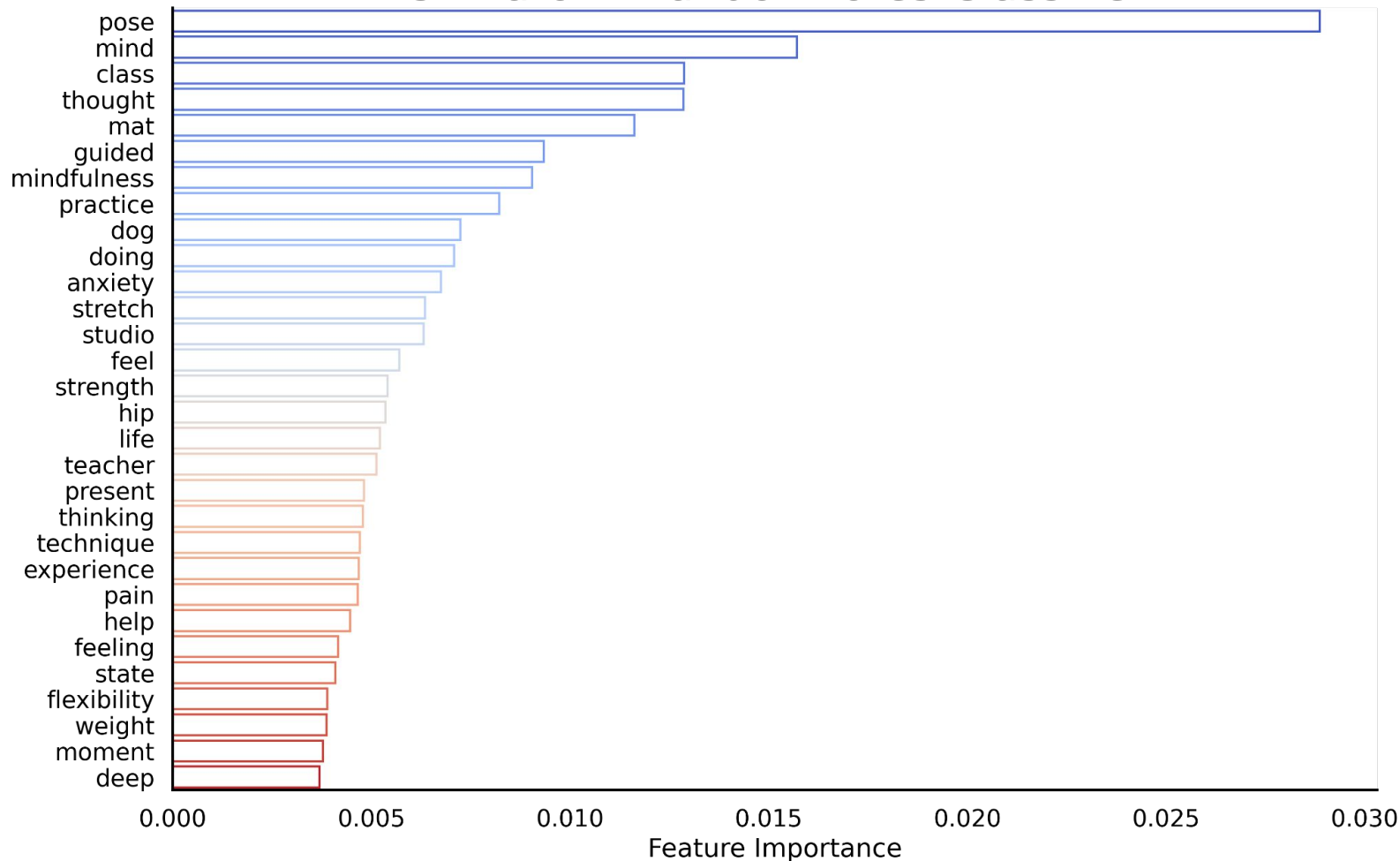
max\_depth: None

n\_estimators: 100

ngram\_range: (1, 1)

# Transformer: "TfidfVectorizer"

## Estimator: "RandomForestClassifier"





---

---

Modeling: *summary*

---

		Validation Score	Train Score	Test Score	Best Score (GS)	Train Score (GS)	Test Score (GS)	F1 - score	Recall (Sensitivity)	Specificity (True Negative Rate)	Precision	Accuracy
Estimator	Transformer											
MultinomialNB	TfidfVectorizer	0.881	0.991	0.884	0.881	0.941	0.883	0.885	0.905	0.860	0.866	0.883
	CountVectorizer	0.890	0.994	0.880	0.883	0.924	0.875	0.875	0.872	0.878	0.878	0.875
RandomForestClassifier	TfidfVectorizer	0.841	0.999	0.855	0.866	0.999	0.855	0.861	0.896	0.813	0.827	0.855
DecisionTreeClassifier	CountVectorizer	0.774	0.999	0.787	0.778	0.842	0.785	0.803	0.876	0.694	0.741	0.785
	TfidfVectorizer	0.752	0.999	0.757	0.760	0.820	0.762	0.788	0.883	0.642	0.711	0.762

---

## Recommendations:

- The corpus of data collected from the employees throughout the first two years of working at the company could be used to determine if Meditation or Yoga would be a better option for the employee.
-

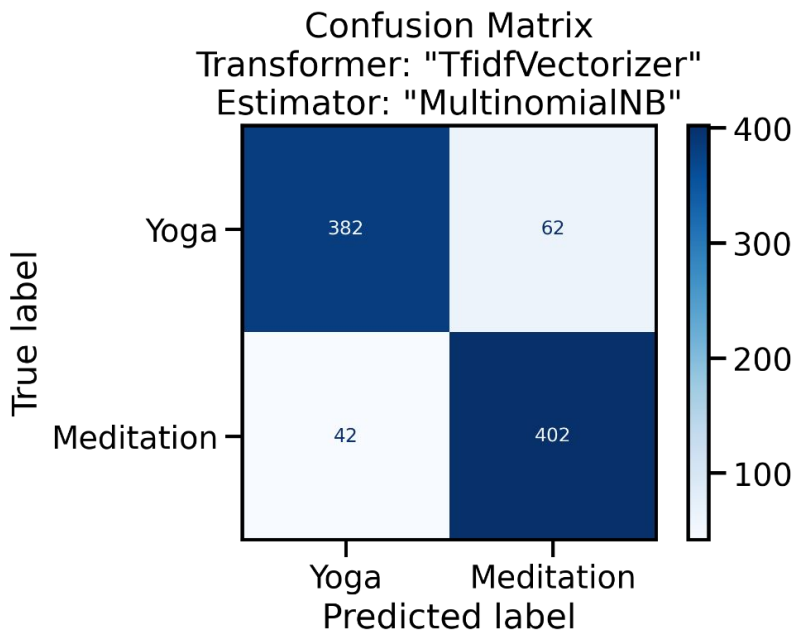
---

---

## Future Work:

- Try other classification models
  - Consider the posts from unique users
-





Accuracy: 88%

Precision: 87%

F1 - score: 89%

Recall (Sensitivity): 91%

Specificity (True Negative Rate): 86%

max\_features: 4000

ngram\_range: (1, 1)

