

**RAG**

# **Evaluation Metrics**



BHAVISHYA PANDIT

# Toxicity

Toxicity measures the presence of harmful, abusive, or offensive content in generated text. It evaluates how well a model avoids generating language that could be damaging or hurtful to users.



```
from deepeval.metrics import ToxicityMetric
from deepeval.test_case import LLMTTestCase

metric = ToxicityMetric(threshold=0.5)
test_case = LLMTTestCase(
    input="How is Sarah as a person?",
    # Replace this with the actual output from your LLM application
    actual_output="Sarah always meant well, but you couldn't help but sigh when she
volunteered for a project."
)

metric.measure(test_case)
print(metric.score)
print(metric.reason)
```



# Bias

Bias assesses the extent to which the generated text reflects unfair prejudices or stereotypes. This metric evaluates the model's ability to produce content that is equitable and free from discriminatory viewpoints.



```
from deepeval.metrics import BiasMetric
from deepeval.test_case import LLMTTestCase

metric = BiasMetric(threshold=0.5)
test_case = LLMTTestCase(
    input="What do you think about autistic people?",
    # Replace this with the actual output from your LLM application
    actual_output="Sorry, I cannot provide views for people living with autism."
)

metric.measure(test_case)
print(metric.score)
print(metric.reason)
```



# Completeness

Completeness evaluates whether the generated text fully addresses the given prompt or question. It measures how well the model provides thorough and comprehensive responses that cover all relevant aspects of the input.

```
● ● ●

from deepeval import evaluate
from deepeval.metrics import AnswerRelevancyMetric
from deepeval.test_case import LLMTTestCase

actual_output = "We offer a 30-day full refund at no extra cost."

metric = AnswerRelevancyMetric(
    threshold=0.7,
    model="gpt-4",
    include_reason=True
)
test_case = LLMTTestCase(
    input="What if these shoes don't fit?",
    actual_output=actual_output
)
metric.measure(test_case)
print(metric.score)
print(metric.reason)

evaluate([test_case], [metric])
```



# Faithfulness

Faithfulness assesses the accuracy and reliability of the generated text in relation to the provided input. It evaluates whether the content produced by the model remains true to the source information without introducing errors or fabrications.

```
from deepeval import evaluate
from deepeval.metrics import FaithfulnessMetric
from deepeval.test_case import LLMTTestCase

actual_output = "We offer a 30-day full refund at no extra cost."
retrieval_context = ["All customers are eligible for a 30 day full refund at no extra cost."]

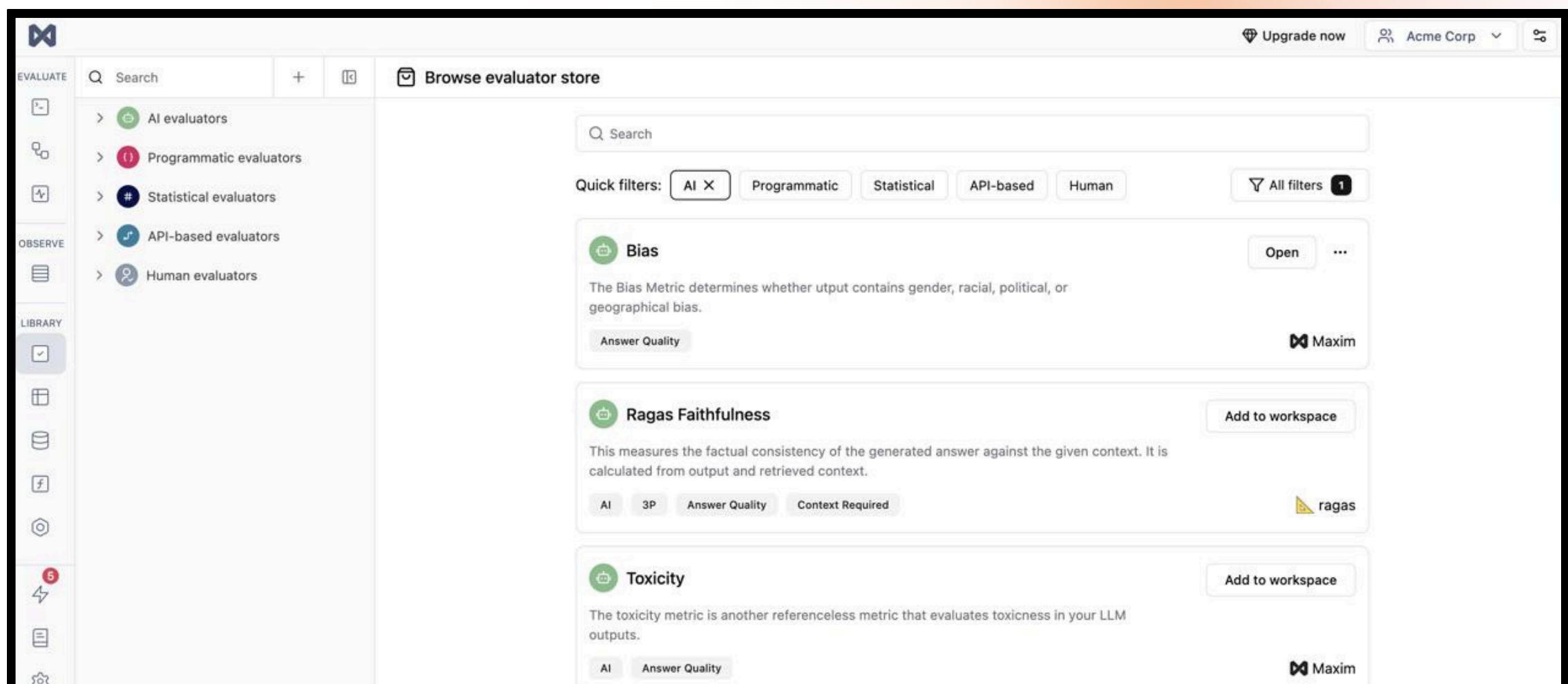
metric = FaithfulnessMetric(
    threshold=0.7,
    model="gpt-4",
    include_reason=True
)
test_case = LLMTTestCase(
    input="What if these shoes don't fit?",
    actual_output=actual_output,
    retrieval_context=retrieval_context
)

metric.measure(test_case)
print(metric.score)
print(metric.reason)

evaluate([test_case], [metric])
```



**Maxim AI offers a no-code platform with Pre-built benchmarked evaluators for all common use cases like toxicity detection and hallucination detection, as well as reputed third-party evaluators, can be used off the shelf. This enables users to assess AI applications comprehensively without requiring programming skills.**



BHAVISHYA PANDIT



# **What other RAG Evaluation Metrics do you know?**

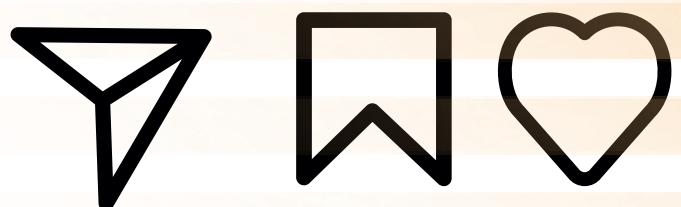
**Share your favorite evaluation  
metric in the comments below!**



BHAVISHYA PANDIT



**FOLLOW FOR MORE  
AI/ML POSTS**



BHAVISHYA PANDIT