# [Evaluation] - Metric Config and Management

| Research | ☐ Confident AI Blog - Resources to help teams stay confident in AI 🔴 Blog - Langfuse |
| --- | --- |
| | 🎯 Evaluating AI Agents - DeepLearning.AI |
| Design | 🟦 Papago AI-Evaluation |
| Relevant PRDs | 🔲 [Evaluation] - Create and manage Test Case for Agent Evaluation |

## Business Context

As an Agent Building Platform, our mission is to help users (agent builders) maximize their ROI from the effort invested in building and refining agents. Similar to traditional software development, having clear indicators that signal improvements after each update is crucial for product growth and success. In this phase of agent optimization, users seek solutions to construct a systematic evaluation workflow that provides consistent and reliable insights through metrics that matter most to them.
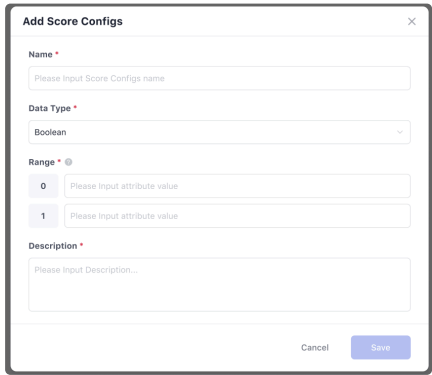
## Requirement

### 1. Metrics  🔖 FE-2346: [Eval Config] list of Metrics  `LAUNCHED`

| Requirement | Wireframe / Design |
| --- | --- |
| Add Eval Config:<br><br>1. Metric Config<br><br>2. Evaluator<br><br>Metric tab show all metrics |  |
| Metrics are searchable, filterable by data type, evaluators, creators | |

| | |
|---|---|
| Metrics is also support pagination and adjust display number. Inherit behavior from existing listing | |
| Columns are immutable. Shown columns are:<br><br>1. Name: text<br>2. Type: selected type<br>3. Criteria: text, each criteria separate by comma `,`<br>4. Description: long text<br>5. Creator: Workflow member<br>6. Reference in (Run): Number of reference to Evaluator | |
| Switchable column: Status: Active, Draft | |
| Action include edit, remove | |

## 2. Create a new metric

| Requirement | Wireframe / Design |
|---|---|
| User can create a new metric with:<br><br>1. Name<br>2. Type<br>3. Criteria (Range)<br>4. Description (Explain why the metric matter and how to measure the metric) |  |

1. Boolean type will only have two value **True, False**. Non editable. Allow editing this would create confusion in the future.
2. **Numeric** will have default 0-100 given this range is very common.
3. **Categorical** unchanged
4. Add a new type Range Qualities to allow users to config cases that have quantitative score to be map with its qualitative label:
   a. **0-100%:** 0% error, 25% bad, 50% normal, ... 100% <quality label>
   b. **1-7:** 1 label A, 2 label B, ....

**Add Metric Configs** ×

Name *
Accuracy

Type *
Categorical

Criteria * ⓘ

| 0 | Correct,100% | 🗑 |
| 1 | Partially correct | 🗑 |
| 2 | Error,100% | 🗑 |

⊕ Add Category

Description *
Only one option can be selected

Cancel  **Save**

---

**Add Metric Configs** ×

Name *
Accuracy

Criteria *
Boolean

Criteria * ⓘ

| 0 | True |
| 1 | False |

Description *
Only one option can be selected

Cancel  **Save**

---

**Add Metric Configs** ×

Name *
Accuracy

Type *
Range Qualities

Criteria * ⓘ

| | From (inc) | To (exc) | Attribute value | |
| 0 | 0 | — Input Range | Please Input attribute value | |
| 1 | Input Range | — Input Range | Please Input attribute value | 🗑 |
| 2 | Input Range | — 10 | Please Input attribute value | |

⊕ Add Range

Description *
Select value within range

Cancel  **Save**

## 3. Manage metrics

| Requirement | Wireframe / Design |
| --- | --- |

| | |
|---|---|
| Metric can be switch from active to draft |  |
| Metric can be deleted. There will be a confirmation pop-up when deleting a metric<br><br>**Deleting or changing a metric will not affect the result of experiment runs or production audit run** | |

## 4. Metric Usage

Metrics could be use in Evaluator configuration and will be in the scope of Evaluator PRD