# Evaluating LLM from AI Quick Actions Catalog

## Referenced Documentation

https://docs.oracle.com/en-us/iaas/data-science/using/ai-quick-actions-evaluations.htm

## Description

With deployed models, you can create a model evaluation to evaluate its performance. You can choose a dataset from Object Storage or upload one from the storage of the notebook you're working in. BERTScore and ROUGE are the evaluation metrics available for measuring model performance. You can save the model evaluation result in Object Storage. You can set the model evaluation parameters. Under advanced options, you can choose the compute instance shape for the evaluation and optionally enter the Stop sequence.
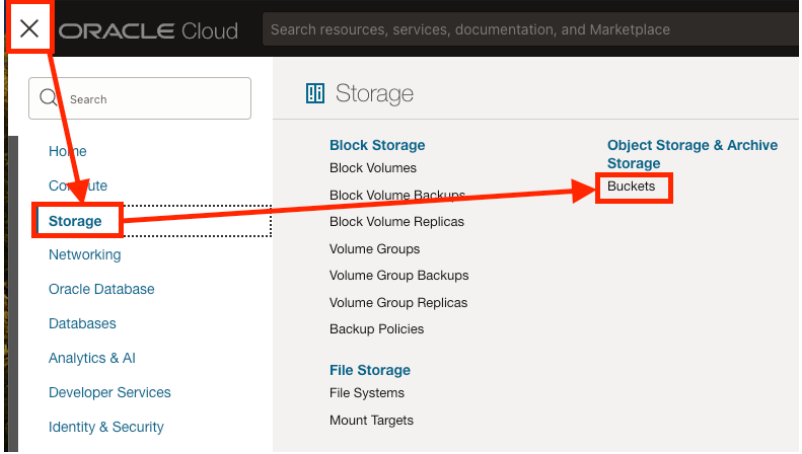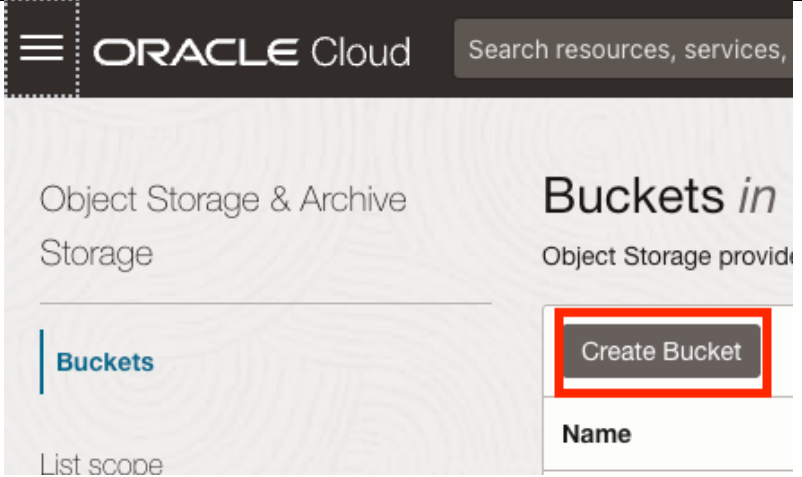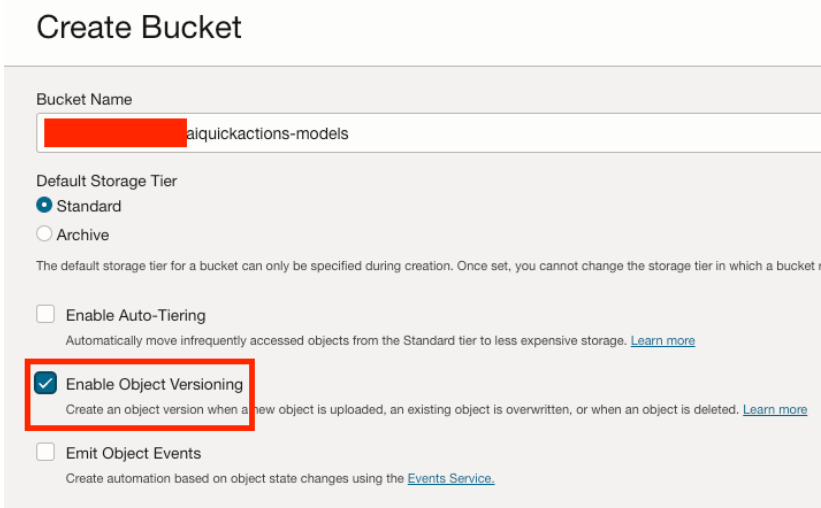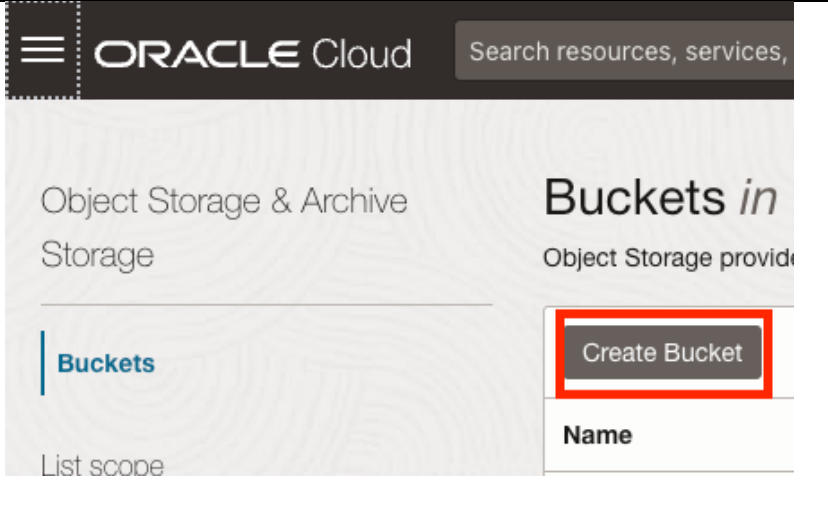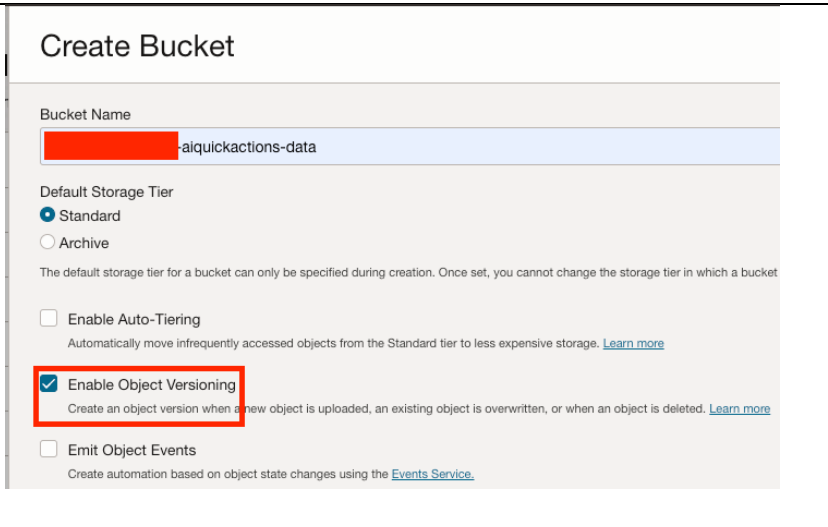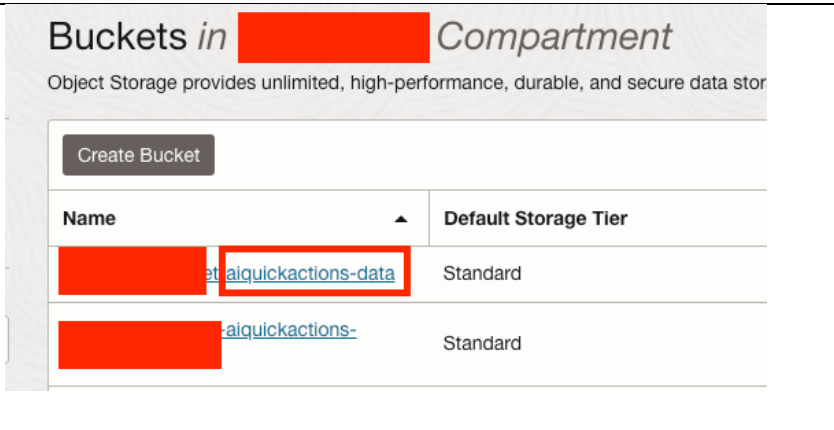
## Pre-Requisites

- Implement the required policies - https://docs.oracle.com/en-us/iaas/data-science/using/ai-quick-actions-set-up.htm
- Ensure you have your OCI Data Science GPU service limits raised for the GPU Shapes you plan to use. This can be done from OCI Console
- Provisioned OCI Data Science Project and Notebook Session (Must be deactivated and reactivated if created before the policies where implemented).
- Deployed AI Quick Actions Model.
- OCI Log Group & Log Created (Optional)

## Guide
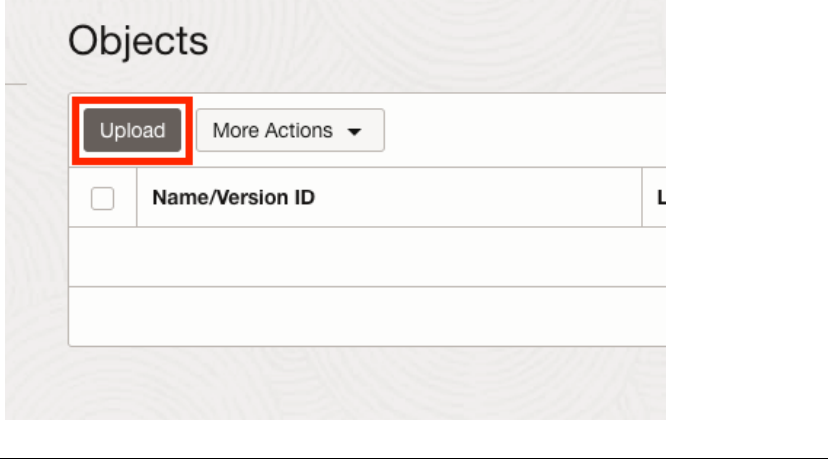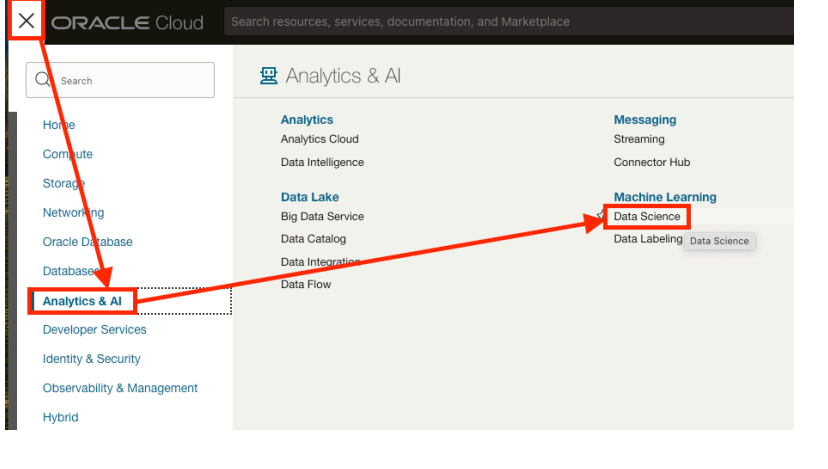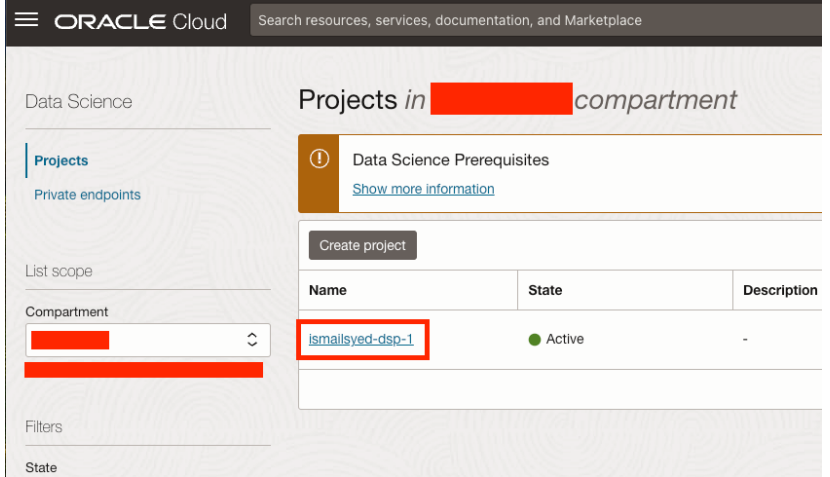
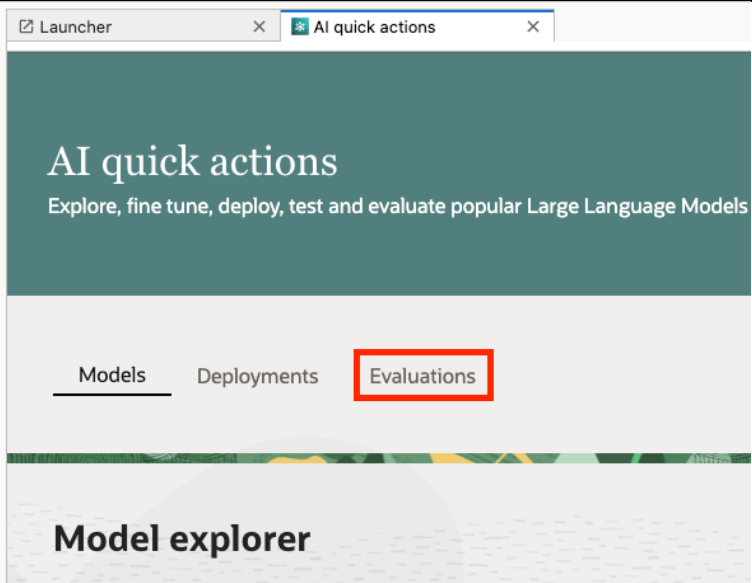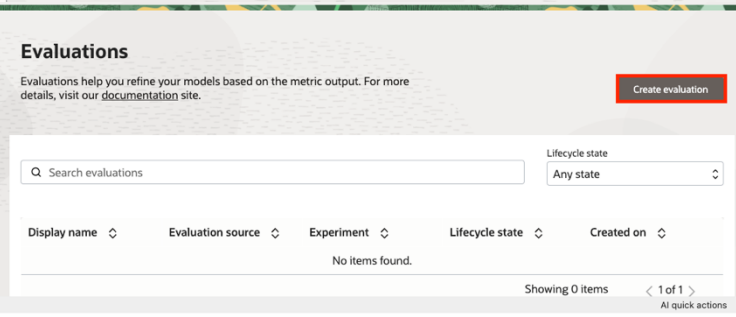| Step | Screenshot |
|---|---|
| Login to the Cloud Console.<br><br>***cloud.oracle.com*** |  |

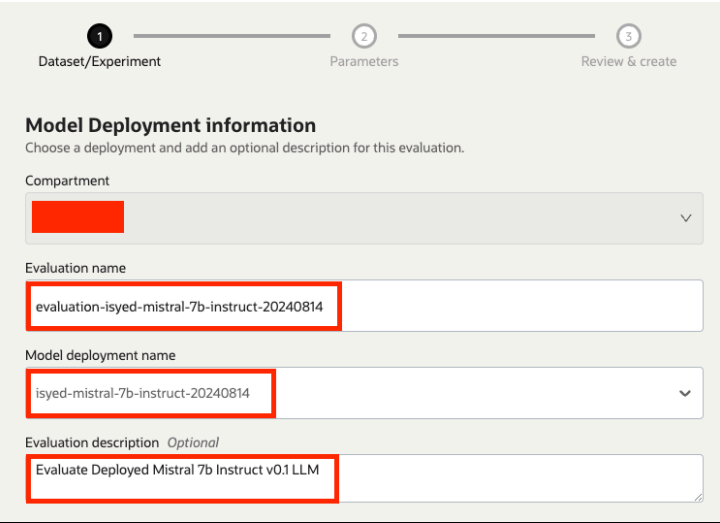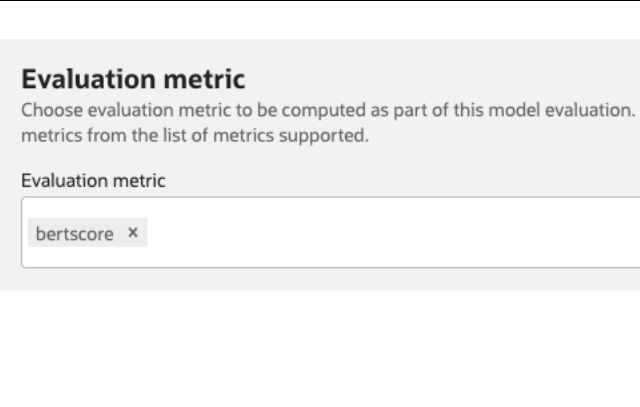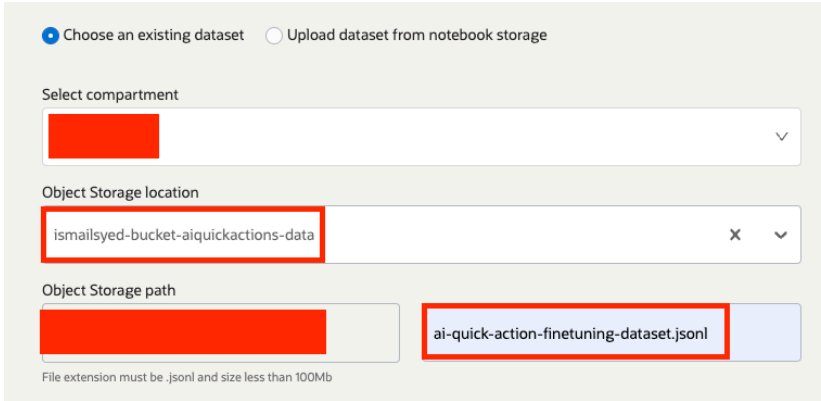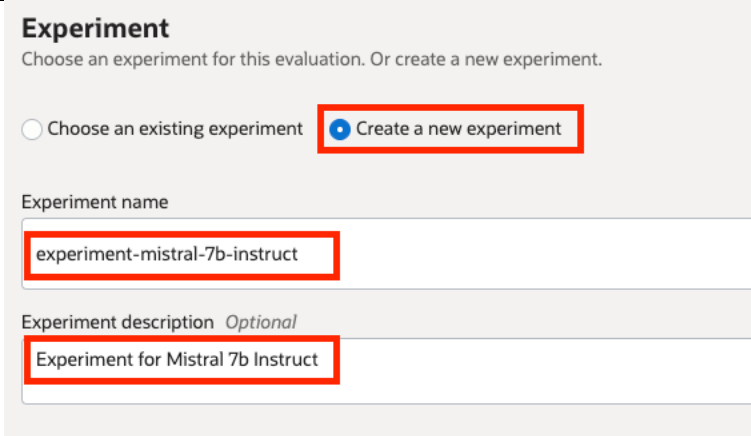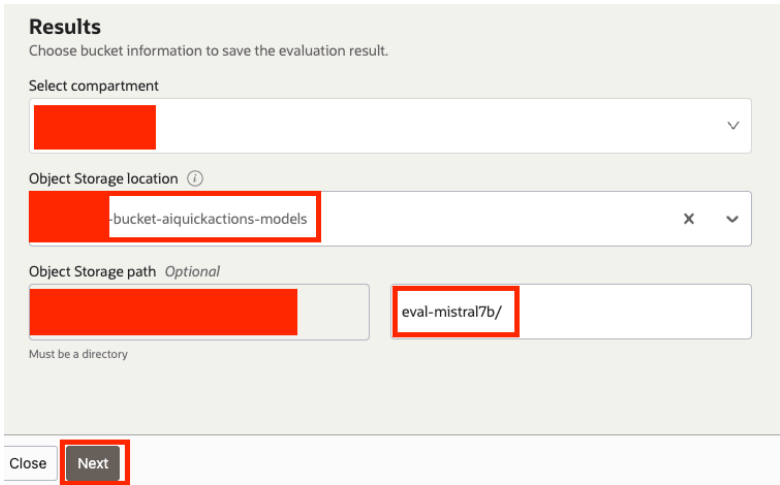| | |
|---|---|
| First, we will create two buckets to store our evaluation data and our evaluation results.<br><br>**Navigate to OCI Menu > Storage > Buckets.** |  |
| **Click Create Bucket.** |  |
| The first bucket we create will be for our evaluation results.<br><br>**Enter a Name.**<br><br>**Enable Object Versioning.**<br><br>**Click Create.** |  |

| | |
|---|---|
| We will then create the bucket for our evaluation data.<br><br>**Click Create Bucket.** |  |
| **Enter a Name.**<br><br>**Enable Object Versioning.**<br><br>**Click Create.** |  |
| **Click on our Data bucket we just created.** |  |

| | |
|---|---|
| **Click Upload.** |  |
| Upload the **ai-quick-action-finetuning-dataset.jsonl**<br><br>Which should come with this guide.<br><br>The dataset must be in **jsonl** format. |  |
| Now navigate to your Data Science Projects.<br><br>**OCI Menu > Analytics & AI > Data Science.** |  |
| Open up your existing Data Science Project. |  |

| | |
|---|---|
| ***Click on*** your existing Data Science Notebook Session.<br><br>Note – This does not have to be a GPU Shape. |  |
| ***Click on Open.***<br><br>This will open up your Data Science Notebook Session.<br><br>***You will have to reauthenticate.*** |  |
| If the policies within the pre-requisites have been implemented correctly you should be able to open up the AI Quick Actions Extension within the Launcher.<br><br>***Click AI quick actions.*** |  |

| | |
|---|---|
| **Click on the Evaluations Tab.** |  |
| **Click on Create Evaluation.** |  |
| **Select your Compartment.**<br><br>**Enter a name for your evaluation.**<br><br>**Select your existing deployed model.**<br><br>**Enter an Evaluation Description.** |  |
| Select your Evaluation Metric. **I have selected Bert Score.**<br><br>You can select multiple evaluation metrics. |  |

| | |
|---|---|
| Select your Evaluation dataset.<br><br>***Choose an existing dataset.***<br><br>***Select your compartment.***<br><br>***Select your Object Storage Bucket (Data Bucket)***<br><br>***Enter path/name of your evaluation dataset.*** | ○ Choose an existing dataset    ○ Upload dataset from notebook storage<br><br>Select compartment<br>[▅▅▅▅▅]    ⌄<br><br>Object Storage location<br>ismailsyed-bucket-aiquickactions-data    ×    ⌄<br><br>Object Storage path<br>[▅▅▅▅▅▅▅▅▅]    ai-quick-action-finetuning-dataset.jsonl<br>File extension must be .jsonl and size less than 100Mb |
| Define new Experiment.<br><br>***Select Create new experiment.***<br><br>***Enter Experiment Name.***<br><br>***Enter Experiment Description.*** | **Experiment**<br>Choose an experiment for this evaluation. Or create a new experiment.<br><br>○ Choose an existing experiment    ● Create a new experiment<br><br>Experiment name<br>experiment-mistral-7b-instruct<br><br>Experiment description  *Optional*<br>Experiment for Mistral 7b Instruct |
| Define where to save evaluation results.<br><br>***Select Compartment.***<br><br>***Select Object Storage Location (Model Bucket).***<br><br>***Enter path where to save results.***<br><br>***Click Next.*** | **Results**<br>Choose bucket information to save the evaluation result.<br><br>Select compartment<br>[▅▅▅▅▅]    ⌄<br><br>Object Storage location  ⓘ<br>[▅▅▅]-bucket-aiquickactions-models    ×    ⌄<br><br>Object Storage path  *Optional*<br>[▅▅▅▅▅▅▅▅]    eval-mistral7b/<br>Must be a directory<br><br>Close    Next |

| | |
|---|---|
| Here you can define your parameters for the LLM. *I have left as the default.*<br><br>Select your Instance shape to run the evaluation. *I have left as the default. VM.Standard.E3.Flex.*<br><br>*Click Next.* |  |
| Review all the Evaluation configuration.<br><br>*Click Submit.* |  |
| This will kick start an OCI Data Science Job in the background to run the Evaluation.<br><br>The Evaluation Status will be displayed as *In Progress.* |  |

| | |
|---|---|
| Once the Evaluation Job is completed the status will update to **Succeeded**. |  |
| If you scroll down, you will see the evaluation metrics displayed. In our case, the BERT Score. |  |
| Head back to the OCI Console and visit your Object Storage Buckets.<br><br>**Navigate to OCI Menu > Storage > Buckets.** |  |
| **Select your Models Object Storage Bucket** which we used to save the Evaluation Metrics to. |  |

| | |
|---|---|
| Under our Evaluation path, **expand the results for our Evaluation Job Run**. |  |
| **Use the 3 dots to Download the report.html** file which gets automatically generated by the Evaluation Job Run. |  |
| We can then take a look at the Model Evaluation Report. Here we can see: **Description of the Evaluation Metric. Overview of the Evaluation Metrics Calculated.** |  |
| We also get a **Box Plot of the BERT F1 Score broken down by the different Categories defined in our Evaluation Dataset.** We can see it performed better on Math related questions compared to the  NULL Category. |  |

| | |
|---|---|
| **We get a list of all the parameters the Model was invoked with.** | **Parameters**<br><br>▼ Model Parameters<br><br>| Parameter | Value |<br>|---|---|<br>| model | odsc-llm |<br>| temperature | 0.700000 |<br>| top_p | 0.990000 |<br>| n | 1 |<br>| max_tokens | 500 |<br>| seed | None |<br>| stop | [] |<br>| stream | False | |
| **Finally, we get a list of each individual sample in our Evaluation Dataset with the Prompt vs the Completion, vs the Response generated by the Model.** | **Results**<br><br>▼ Results Table<br><br>Copy  Excel  CSV                    Search:<br><br>| prompt | completion | category | response |<br>|---|---|---|---|<br>| <s>[INST] <<SYS>> You are a helpful assistant <</SYS>> A bat and a ball cost $1.10 in total. The bat costs $1.00 more than the ball. How much does the ball cost? [/INST] | The ball costs 5 cents. | _ | To solve this problem, we can use algebra. Let's represent the cost of the ball as "x" dollars. The cost of the bat is then $1.00 more than the cost of the ball, so it costs $1.00 + x dollars. The total cost of the bat and ball is $1.10, so we can write the equation: $1.00 + x = $1.10 Now, we can solve for x by subtracting $1.00 from both sides of the equation: x = $1.10 - $1.00 x = $0.10 So, the ball costs $0.10 dollars. |<br>| <s>[INST] <<SYS>> You are a helpful assistant <</SYS>> A bat and a ball cost $1.10 in total. The bat costs $1.00 more than the ball. How much does the ball cost? [/INST] | The ball costs 5 cents. | _ | The ball costs $0.10. | |