# Kserve 開源貢獻

Group 1

# kserve是甚麼
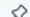
## 🤖 KServe 是什麼？

● KServe 是 在 Kubernetes 上的 AI 模型部署框架

● 幫助我們把訓練好的 AI 模型 ➜ 做成一個 API, 能「即時預測」

● 支援多種模型格式：TensorFlow / PyTorch / MLflow / ONNX...

# kserve是甚麼

# 開源貢獻是甚麼

💡 **什麼是開源（Open Source）？**

● 原始碼公開，任何人都能閱讀、使用、修改、貢獻

● 常見平台 : GitHub (kserve、openCV)

🔧 **可以怎麼貢獻？**

● 改進文件

● 回報錯誤（開 issue）

● 寫程式碼（修 bug、新功能）

● 參與討論（Pull Request 審核、留言）

# 開源貢獻流程

1.尋找適合的問題(ISSUE)

2.在本地端重現問題

3.用各種方式解決問題並跟rewiewers討論

4.提出Pull Request (你自己用甚麼方式解決這個問題)

5.跟rewiewers討論pull request是否能成功合併或需再修改

6.成功被merge(貢獻)

尋找適合問題
並且成功重現

# 尋找適合問題

# GIT指令

- GIT CLONE

- GIT CHECKOUT <BRANCH>

- GIT STATUS

- GIT ADD

- GIT COMMIT

- GIT PUSH

# 問題重現

ImportError: cannot import name 'BentoService' from 'bentoml'
(/usr/local/lib/python3.9/dist-packages/bentoml/__init__.py) #3778

New issue

⊙ Open

LOADBC opened on Jul 8, 2024 · edited by LOADBC        Edits ▾  ⋯

/kind bug

I was trying out some samples on Predict on an InferenceService using BentoML
https://github.com/kserve/kserve/tree/master/docs/samples/bentoml

i installed :-
-pip3 install bentoml scikit-learn
i created :-
-Create the BentoML Service file (iris_classifier.py)

```
from bentoml import env, artifacts, api, BentoService
from bentoml.handlers import DataframeHandler
from bentoml.artifact import SklearnModelArtifact

@env(auto_pip_dependencies=True)
@artifacts([SklearnModelArtifact('model')])
class IrisClassifier(BentoService):

    @api(DataframeHandler)
    def predict(self, df):
        return self.artifacts.model.predict(df)
```

**Assignees**
No one assigned

**Labels**
kind/bug

**Type**
No type

**Projects**
No projects

**Milestone**
No milestone

**Relationships**
None yet

**Development**

9

# 問題重現

-Create the model training and service saving script main.py

```python
from sklearn import svm
from sklearn import datasets

from iris_classifier import IrisClassifier

if __name__ == "__main__":
    # Load training data
    iris = datasets.load_iris()
    X, y = iris.data, iris.target

    # Model Training
    clf = svm.SVC(gamma='scale')
    clf.fit(X, y)

    # Create a iris classifier service instance
    iris_classifier_service = IrisClassifier()

    # Pack the newly trained model artifact
    iris_classifier_service.pack('model', clf)

    # Save the prediction service to disk for model serving
    saved_path = iris_classifier_service.save()
```

-then i ran python3 main.py

-but i got this error:-

Traceback (most recent call last):
File "/root/main.py", line 4, in
from iris_classifier import IrisClassifier
File "/root/iris_classifier.py", line 1, in
from bentoml import env, artifacts, api, BentoService
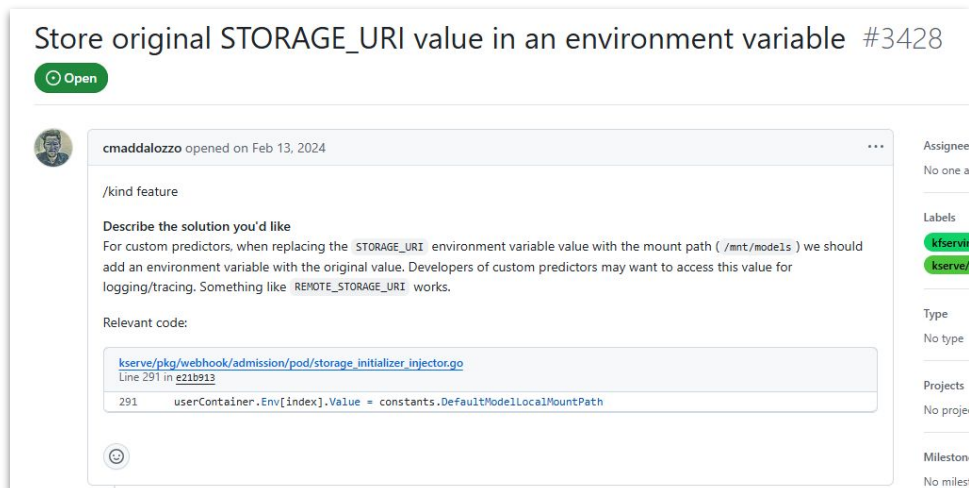ImportError: cannot import name 'env' from 'bentoml' (/usr/local/lib/python3.9/dist-packages/bentoml/**init**.py)

# 問題重現

**Environment:**

- KServe Version: 0.11
- Kubeflow version: no
- Kind version: 0.23.0

解決問題
並且提出PR

# Issue #3428 : Store original STORAGE_URI value in an environment variable

Store original STORAGE_URI value in an environment variable #3428

⊙ Open

cmaddalozzo opened on Feb 13, 2024

/kind feature

**Describe the solution you'd like**
For custom predictors, when replacing the `STORAGE_URI` environment variable value with the mount path ( `/mnt/models` ) we should add an environment variable with the original value. Developers of custom predictors may want to access this value for logging/tracing. Something like `REMOTE_STORAGE_URI` works.

Relevant code:

kserve/pkg/webhook/admission/pod/storage_initializer_injector.go
Line 291 in e21b913

```
291        userContainer.Env[index].Value = constants.DefaultModelLocalMountPath
```

Assignee
No one a

Labels
kfservi
kserve/

Type
No type

Projects
No proje

Milestone
No miles

Issue 介紹：

- STORAGE_URI (環境變數):
  告訴 KServe 從哪裡下載模型
- 模型被下載後, KServe 會把
  這個變數的值換掉, 換成一
  個本地的路徑
- 希望在 KServe加上一個新的
  環境變數會保留「原本的路
  徑值」

# Issue #3428 : Store original STORAGE_URI value in an environment variable

storage_initializer_injector.go

- **兩個地方可控制功能**
  - ○ 平台預設開關：整個平台、所有pod預設行為
  - ○ 個別複寫開關：單一pod可自行決定是否開啟

```go
func (mi *StorageInitializerInjector) InjectStorageInitializer(pod *corev1.Pod) error {
    // Only inject if the required annotations are set
    srcURI, ok := pod.ObjectMeta.Annotations[constants.StorageInitializerSourceUriInternalAnnotationKey]
    if !ok {
        return nil
    }

    enabledRemoteEnv := mi.config.EnableRemoteStorageEnv
    if ann, ok := pod.Annotations[EnableRemoteStorageEnvAnnotation]; ok &&
        strings.ToLower(ann) == "true" {
        enabledRemoteEnv = true
    }
}
```

# Issue #3428 : Store original STORAGE_URI value in an environment variable

storage_initializer_injector.go

- **注入環境變數 ( 在找到container後新增 )**
  - 只在 enabledRemoteEnv == true 時執行。
  - 保證不重覆注入（利用hasEnv 檢查）。

```
// Find the kserve-container (this is the model inference server) and transformer container and the worker-container
userContainer := getContainerWithName(pod, constants.InferenceServiceContainerName)
transformerContainer := getContainerWithName(pod, constants.TransformerContainerName)
workerContainer := getContainerWithName(pod, constants.WorkerContainerName)

if enabledRemoteEnv && userContainer != nil && !hasEnv(userContainer, RemoteStorageEnvVarName) {
        addOrReplaceEnv(userContainer, RemoteStorageEnvVarName, srcURI)
}
```

# Issue #3428 : Store original STORAGE_URI value in an environment variable

- make test 測試已通過

# Issue #3428 : Store original STORAGE_URI value in an environment variable

- 提出PR

feat: support remote storage URI injection for serving runtimes #4492

New issue

⟂ Open  de0725 wants to merge 2 commits into kserve:master from de0725:feature/storage-uri-revamp

💬 Conversation 0   ⊙ Commits 2   ☑ Checks 1   🗎 Files changed 3    +217 −0

**de0725** commented yesterday

**What this PR does / why we need it:**

1. Injects REMOTE_STORAGE_URI= into the main user container.
2. Respects both storageSpec-based and PVC-based source URIs.
3. Adds corresponding unit tests to ensure correctness and backward compatibility.
   REMOTE_STORAGE_URI.Usage.and.Testing.Guide (1).docx

**Which issue(s) this PR fixes** (optional, in `fixes #<issue number>(, fixes #<issue_number>, ...)` format, will close the issue(s) when PR gets merged):
Fixes #3428

**Type of changes**
Please delete options that are not relevant.

☐ Bug fix (non-breaking change which fixes an issue)
☑ New feature (non-breaking change which adds functionality)
☐ Breaking change (fix or feature that would cause existing functionality to not work as expected)
☐ This change requires a documentation update

**Feature/Issue validation/testing:**
Added unit test: TestRemoteStorageUriEnvInjection
Verified behavior under both ConfigMap and Annotation-based configurations
Ran make test, all unit tests passed

Please describe the tests that you ran to verify your changes and relevant result summary. Provide instructions so it can be reproduced.
Please also list any relevant details for your test configuration.

Reviewers
No reviews

Assignees
No one assigned

Labels
None yet
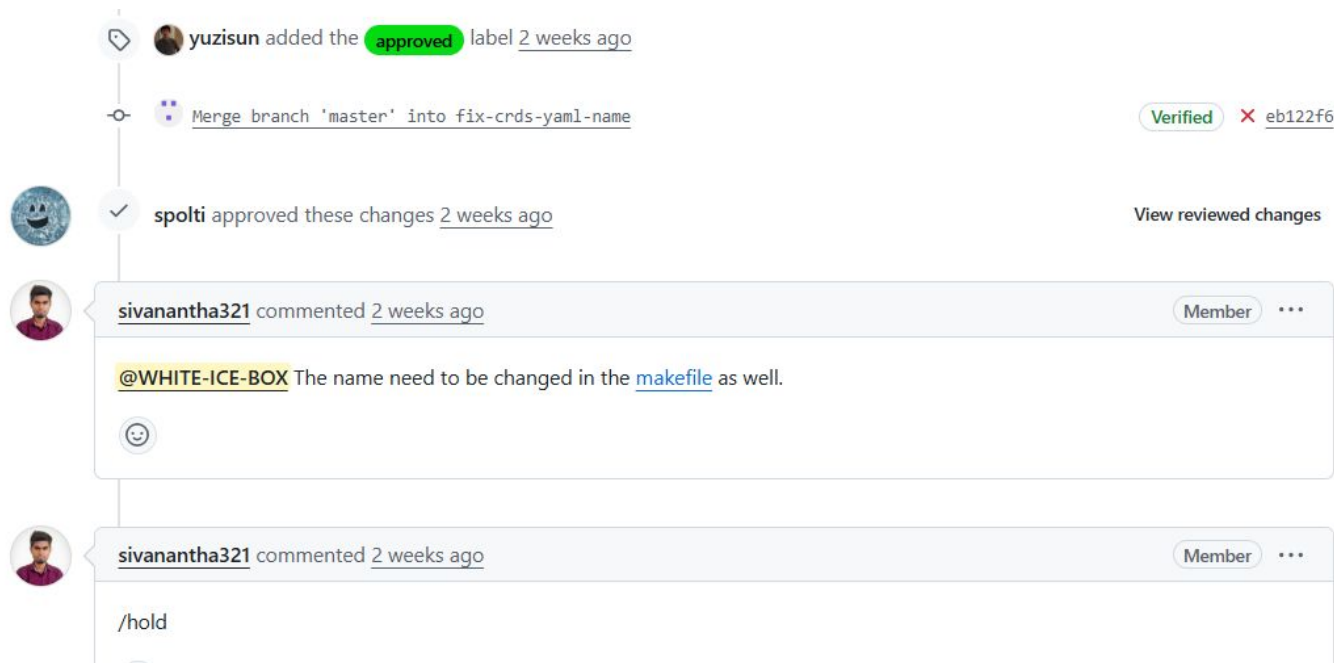
Projects
None yet

Milestone
No milestone

Development
Successfully merging this pull request may close these issues.
⊙ Store original STORAGE_URI value in an environ...

1 participant

# 討論PR可行性
# 並且成功被MERGE

# Rename CRD file to reflect all KServe CRDs (Fixes #4396) #4494

# Rename CRD file to reflect all KServe CRDs (Fixes #4396) #4494

# Rename CRD file to reflect all KServe CRDs (Fixes #4396) #4494

# QNA

🧠 **為什麼要貢獻？**

● 學技術, 練 Git

● 建立個人作品集

● 跟世界各地的高手合作