Image Similarity and OCR-Based Analysis

Overview

This project aims to perform two main tasks:

- 1. **Image Similarity Calculation**: For sets 1 to 7, the project calculates the similarity between a main image and two test images using a pre-trained VGG16 model.
- 2. **OCR-Based Value Extraction**: For sets 8 and 9, the project extracts specific numerical values ("TOTAL WIN" and "BET") from two test images using OCR (Optical Character Recognition) with PaddleOCR.

Intuition Behind the Approach

1. Image Similarity Calculation:

- o **Model**: We use the VGG16 model pre-trained on ImageNet for feature extraction.
- o **Process**: The images are pre-processed and fed into the VGG16 model to extract feature vectors. Cosine similarity is then calculated between the feature vectors of the main image and the test images to determine their similarity percentages.

2. OCR-Based Value Extraction:

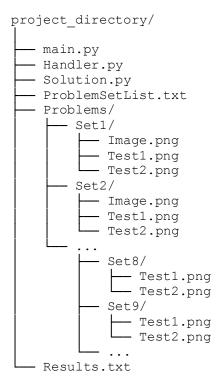
- o **Model**: PaddleOCR is used for extracting text from images.
- Process: The OCR model detects text regions in the images, and specific target words ("TOTAL WIN" and "BET") are identified. The nearest numerical value to these target words is extracted and cleaned (commas removed) for further use.

0

Step-by-Step Guide to Run the Project

1. Project Structure

Ensure your project directory has the following structure:



2. Setting Up the Environment

Ensure you have the necessary libraries installed:

```
pip install certifi paddleocr tensorflow scikit-learn pillow
```

3. Running the Project

Follow these steps to execute the project:

- 1. Initialization:
 - o Ensure the Results.txt file is empty or non-existent.
 - o The ProblemSetList.txt file should list all the sets to be processed.
- 2. Execution:
 - o Run the main.py script to start the process.

```
python main.py
```

4. Output

- The Results.txt file will be populated with the results.
 - o For sets 1 to 7: Image similarity percentages between the main image and test images.
 - o For sets 8 and 9: Extracted numerical values for "TOTAL WIN" and "BET".

Code Files

```
main.py
```

```
from Handler
import start, execute

def run():
    start()
    execute()

if __name__ == '__main__':
    run()

Handler.py
```

```
if len(result) == 2:
    save(','.join(map(str, result)))
```

Solution.py

```
import certifi
from paddleocr import PaddleOCR
import os
import re
import numpy as np
from PIL import Image
import tensorflow as tf
from tensorflow.keras.applications.vgg16 import VGG16, preprocess input
from tensorflow.keras.preprocessing import image
from tensorflow.keras.models import Model
from sklearn.metrics.pairwise import cosine similarity
os.environ['SSL CERT FILE'] = certifi.where()
class Solution:
    def init (self):
        self.base model = VGG16(weights='imagenet', include top=False, pooling='avg')
        self.model = Model(inputs=self.base model.input,
outputs=self.base model.output)
        self.ocr = PaddleOCR(use angle cls=True, lang='en')
    def preprocess image(self, img path):
        img = image.load img(img path, target size=(224, 224))
        img_array = image.img_to_array(img)
        img array = np.expand dims(img array, axis=0)
        img array = preprocess input(img array)
        return img_array
    def extract features(self, img path):
        img array = self.preprocess image(img path)
        features = self.model.predict(img array)
        return features
    def calculate similarity(self, main image path, test image path1,
test image path2):
        main features = self.extract features(main image path)
        test1_features = self.extract_features(main_image_path)
        test2 features = self.extract features(test image path2)
        similarity1 = cosine_similarity(main_features, test1_features)[0][0]
        similarity2 = cosine similarity(main features, test2 features)[0][0]
        percentage1 = similarity1 * 100
        percentage2 = similarity2 * 100
        return [percentage1, percentage2]
    def process image sets(self, base path, set number):
        main image path = f'{base path}/Set{set number}/Image.png'
        test image path1 = f'{base_path}/Set{set_number}/Test1.png'
        test image path2 = f'{base path}/Set{set number}/Test2.png'
        return self.calculate similarity(main image path, test image path1,
test image path2)
    def extract numerical value(self, text):
        pattern = r'[-+]?\d{1,3}(?:,\d{3})*(?:\.\d+)?'
        match = re.search(pattern, text)
            return match.group().replace(',', '') # Remove commas from the extracted
number
        else:
            return None
    def find nearest numerical value (self, target word, ocr result, max distance):
        target box = None
```

```
for (box, text) in zip(ocr result['boxes'], ocr result['texts']):
            if target word in text:
                target box = box
                break
        if target box is None:
            return ""
        target_x, target_y = target box[0][0], target box[0][1]
        target center = np.array([target x + (target box[1][0] - target x) / 2,
                                  target y + (target box[2][1] - target y) / 2])
        closest value = None
        closest distance = float('inf')
        for box, text in zip(ocr result['boxes'], ocr result['texts']):
            numerical value = self.extract numerical value(text)
            if numerical_value is not None:
                box x, box y = box[0][0], box[0][1]
                box center = np.array([box x + (box[1][0] - box x) / 2,
                                       box_y + (box[2][1] - box_y) / 2])
                distance = np.linalg.norm(target center - box center)
                if distance < closest distance and distance <= max distance:
                    closest distance = distance
                    closest value = numerical value
        return closest_value if closest_value is not None else ""
    def find nearest values for targets (self, target words, ocr result, max distance):
        results = {}
        for word in target words:
            results[word] = self.find nearest numerical value(word, ocr result,
max distance)
        return results
    def process images for sets(self, base path, set number, target words,
max distance):
        results = {word: [] for word in target words}
        image files = ['Test1.png', 'Test2.png']
        for image file in image files:
            img path = os.path.join(base path, f'Set{set number}', image file)
            if not os.path.exists(img path):
                print(f"File {img_path}) not found. Skipping.")
                continue
            try:
                img = Image.open(img path)
            except Exception as e:
                print(f"Error loading image {img path}: {e}")
                continue
            ocr result = self.ocr.ocr(img_path, cls=True)
            parsed result = {'boxes': [], 'texts': [], 'scores': []}
            for line in ocr result:
                for word info in line:
                    parsed result['boxes'].append(word info[0])
                    parsed result['texts'].append(word info[1][0])
                    parsed result['scores'].append(word info[1][1])
            nearest values = self.find nearest values for targets(target words,
parsed result, max distance)
            for word in target words:
                results [word] .append (nearest values [word])
        return results
    def get answer(self, problem):
        base path = './Problems'
        if problem == 'Set8':
            target words = ["TOTAL WIN"]
            return self.process_images_for_sets(base_path, 8, target words, 150)["TOTAL
WIN"]
        elif problem == 'Set9':
            target words = ["BET"]
            return self.process_images_for_sets(base_path, 9, target_words, 150)["BET"]
        else:
```

```
set_number = int(problem.split('Set')[1])
return self.process_image_sets(base_path, set_number)
```

ProblemSetList.txt

Set1 Set2 Set3 Set4 Set5 Set6 Set7 Set8 Set9

Key Points

- Image Similarity: Uses VGG16 for feature extraction and cosine similarity for comparison.
- OCR: PaddleOCR for text detection and extraction, targeting specific words and nearest numerical values.
- Flexibility: The solution is adaptable for different sets based on the problem requirements.

By following this guide, you should be able to run the project and understand the intuition behind the approach. The provided code processes images according to their respective requirements, and the results are saved in a structured manner.