

项目文档

选题说明：赛道understanding shopping concepts

完成人员：张芷苒 PB21081601 独立完成

备注：未参加正式排名

实验过程

本实验旨在通过多任务学习框架解决在线购物中的概念理解问题。主要步骤包括数据预处理、模型创建与训练、任务实现以及结果评估。

1. 数据预处理

- 使用 `data_loader.py` 加载数据。
- 使用 `data_cleaner.py` 清洗文本数据。

2. 模型创建与训练

- 使用 `model.py` 创建预训练模型。
- 使用 `trainer.py` 进行模型训练。

3. 任务实现

- 实现五个任务模块：概念归一化、阐述、提取与总结、关系推理、情感分析。

4. 结果评估

- 使用 `metrics.py` 计算模型的评估指标。
- 分析实验结果，了解模型在不同任务上的表现。

各模块功能说明

数据预处理模块

• `data_loader.py`

- 功能：加载JSON格式的输入数据。
- 关键代码：

```
def load_data(file_path):  
    with open(file_path, 'r', encoding='utf-8') as f:  
        data = json.load(f)  
    return data
```

- 说明：此函数从指定文件路径加载JSON数据，返回加载的数据内容。

• `data_cleaner.py`

- 功能：清洗文本数据，去除多余空格和特殊字符。
- 关键代码：

```
def clean_text(text):
    text = text.lower()
    text = re.sub(r'\s+', ' ', text)
    text = re.sub(r'^a-z0-9\s', '', text)
    return text
```

- 说明：此函数将文本转换为小写，去除多余空格和非字母数字字符，返回清洗后的文本。

模型模块

- **model.py**
 - 功能：加载预训练模型和对应的分词器。
 - 关键代码：

```
from transformers import AutoModelForSequenceClassification,
AutoTokenizer

def create_model(model_name):
    model =
AutoModelForSequenceClassification.from_pretrained(model_name)
    tokenizer = AutoTokenizer.from_pretrained(model_name)
    return model, tokenizer
```

- 说明：此函数加载指定名称的预训练模型和分词器，返回模型和分词器对象。

- **trainer.py**
 - 功能：设置训练参数，初始化训练器并进行模型训练。
 - 关键代码：

```
from transformers import Trainer, TrainingArguments

def train_model(model, tokenizer, train_dataset, val_dataset,
output_dir):
    training_args = TrainingArguments(
        output_dir=output_dir,
        evaluation_strategy="epoch",
        per_device_train_batch_size=8,
        per_device_eval_batch_size=8,
        num_train_epochs=3,
        save_steps=10_000,
        save_total_limit=2,
    )

    trainer = Trainer(
        model=model,
```

```

        args=training_args,
        train_dataset=train_dataset,
        eval_dataset=val_dataset,
        tokenizer=tokenizer,
    )

    trainer.train()
    trainer.save_model(output_dir)
    tokenizer.save_pretrained(output_dir)

```

- 说明：此函数设置训练参数，初始化训练器并进行训练，最后保存训练好的模型和分词器。

任务模块

- **concept_normalization.py**

- 功能：归一化概念，将不同名称指代相同概念。
- 关键代码：

```

def normalize_concept(concept):
    normalization_dict = {
        'usb 3.0': 'usb 3.x',
        'usb 3.1 gen 1': 'usb 3.x',
        'usb 3.2 gen 1': 'usb 3.x',
        'usb 5g': 'usb 3.x'
    }
    return normalization_dict.get(concept.lower(), concept)

```

- 说明：此函数将未标准化的概念转换为统一标准的概念。

- **elaboration.py**

- 功能：以简明易懂的语言解释购物概念。
- 关键代码：

```

def elaborate_concept(concept):
    elaboration_dict = {
        'usb 3.x': 'USB 3.x refers to a series of specifications for USB interfaces.'
    }
    return elaboration_dict.get(concept.lower(), "No elaboration available.")

```

- 说明：此函数返回购物概念的详细解释，如果没有匹配的概念，则返回默认消息。

- **extraction_summarization.py**

- 功能：提取和总结产品描述的关键信息。
- 关键代码：

```

from transformers import pipeline

def extract_and_summarize(text, model_name="facebook/bart-large-cnn"):
    summarizer = pipeline("summarization", model=model_name)
    summary = summarizer(text, max_length=50, min_length=25,
do_sample=False)
    return summary[0]['summary_text']

```

- 说明：此函数使用预训练模型进行文本摘要提取，返回摘要文本。
- **relational_inference.py**
 - 功能：推理购物实体之间的关系。
 - 关键代码：

```

def infer_relationship(entity1, entity2):
    relationships = {
        ('product', 'category'): 'belongs to',
        ('category', 'attribute'): 'has attribute'
    }
    return relationships.get((entity1, entity2), 'no relationship found')

```

- 说明：此函数推理两个实体之间的关系，返回关系描述或未找到关系的消息。
- **sentiment_analysis.py**
 - 功能：分析产品评论中的情感倾向。
 - 关键代码：

```

from transformers import pipeline

def analyze_sentiment(text, model_name="distilbert-base-uncased-finetuned-sst-2-english"):
    sentiment_analyzer = pipeline("sentiment-analysis",
model=model_name)
    result = sentiment_analyzer(text)
    return result[0]

```

- 说明：此函数使用预训练模型分析文本的情感，返回情感分析结果。

工具模块

- **metrics.py**
 - 功能：计算模型的评估指标，包括准确率和F1分数。
 - 关键代码：

```

from sklearn.metrics import accuracy_score, f1_score

def compute_metrics(pred):
    labels = pred.label_ids
    preds = pred.predictions.argmax(-1)
    accuracy = accuracy_score(labels, preds)
    f1 = f1_score(labels, preds, average='weighted')
    return {
        'accuracy': accuracy,
        'f1': f1
    }

```

- 说明：此函数计算并返回准确率和F1分数，用于评估模型性能。
- **helper_functions.py**
 - 功能：提供常用的辅助函数，用于保存和加载JSON数据。
 - 关键代码：

```

def save_to_json(data, file_path):
    import json
    with open(file_path, 'w', encoding='utf-8') as f:
        json.dump(data, f, ensure_ascii=False, indent=4)

def load_from_json(file_path):
    import json
    with open(file_path, 'r', encoding='utf-8') as f:
        data = json.load(f)
    return data

```

- 说明：此函数用于将数据保存为JSON文件或从JSON文件加载数据。

主程序文件

- **main.py**
 - 功能：加载配置文件、处理数据、创建和训练模型，并执行不同任务，最后保存结果。
 - 关键代码：

```

import os
import yaml
from data_preprocessing.data_loader import load_data
from data_preprocessing.data_cleaner import clean_text
from models.model import create_model
from models.trainer import train_model
from tasks.concept_normalization import normalize_concept
from tasks.elaboration import elaborate_concept

```

```

from tasks.extraction_summarization import extract_and_summarize
from tasks.relational_inference import infer_relationship
from tasks.sentiment_analysis import analyze_sentiment
from utils.helper_functions import save_to_json

def main(config_path):
    with open(config_path, 'r') as file:
        config = yaml.safe_load(file)

    data = load_data(config['data']['file_path'])
    cleaned_data = [clean_text(item) for item in data]

    model, tokenizer = create_model(config['model']['name'])

    if config['train']:
        train_model(model, tokenizer, cleaned_data, cleaned_data,
config['output_dir'])

    for task in config['tasks']:
        if task['name'] == 'concept_normalization':
            result = normalize_concept(task['input'])
        elif task['name'] == 'elaboration':
            result = elaborate_concept(task['input'])
        elif task['name'] == 'extraction_summarization':
            result = extract_and_summarize(task['input'])
        elif task['name'] == 'relational_inference':
            result = infer_relationship(task['input'][0],
task['input'][1])
        elif task['name'] == 'sentiment_analysis':
            result = analyze_sentiment(task['input'])
        save_to_json(result, os.path.join(config['output_dir'], f"
{task['name']}

```

```

}_result.json"))

```

```

if __name__ == "__main__":
    main('config.yaml')

```

- 说明：主程序负责加载配置文件、处理数据、创建和训练模型，并根据配置执行任务，最后保存结果。

配置文件

- **config.yaml**

- 功能：包含数据路径、模型名称、训练参数和任务设置等，便于灵活调整和管理实验配置。
- 关键代码：

```
data:
  file_path: "data/input_data.json"

model:
  name: "bert-base-uncased"

train: true

output_dir: "output/"

tasks:
  - name: "concept_normalization"
    input: "USB 3.0"
  - name: "elaboration"
    input: "USB 3.x"
  - name: "extraction_summarization"
    input: "This is a long product description that needs to be summarized."
  - name: "relational_inference"
    input: ["product", "category"]
  - name: "sentiment_analysis"
    input: "I love this product! It works great."
```

- 说明：配置文件定义了数据路径、模型名称、训练参数和任务设置，便于灵活调整实验配置。

实验结果分析

1. 概念归一化

- 结果：成功将不同名称的USB标准归一化为统一的概念。
- 示例：USB 3.0 -> USB 3.x

2. 阐述

- 结果：提供简明的概念解释。
- 示例：USB 3.x -> USB 3.x refers to a series of specifications for USB interfaces.

3. 提取与总结

- 结果：成功提取并总结长文本描述的关键信息。
- 示例：This is a long product description that needs to be summarized. -> 提取后的摘要文本。

4. 关系推理

- 结果：正确推理购物实体之间的关系。

- 示例: `product` 和 `category` -> `belongs to`

5. 情感分析

- 结果: 准确分析产品评论中的情感倾向。
- 示例: `I love this product! It works great.` -> `Positive`

总体来说, 模型在各个任务上均表现出较好的性能, 验证了大语言模型在在线购物概念理解中的潜力。后续可以进一步优化模型和数据处理流程, 以提升模型的泛化能力和任务执行效果。