

# HTRやってみたver4

7-22 個別ゼミ 工藤滉青

# アウトライン

- コードの理解を進める
- いろんな人が書いた筆記スタイルで精度評価
- 生成したテキストが正しいかどうかの処理をしているのか(入っていたらオフにする)

# 1.コードの理解を進める

実行コマンド：python -m src train

```
1  from .cli import cli
2
3  main = cli
4
5  if __name__ == "__main__":
6      main()
7
```

\_\_main\_\_.py

```
import typer
from pathlib import Path

from .main import TrocrPredictor, main_train, main_validate

cli = typer.Typer(name="TrOCR")

@cli.command()
def train(local_model: bool = False):
    main_train(local_model)

@cli.command()
def validate(local_model: bool = True):
    main_validate(local_model)

@cli.command()
def predict(image_paths: list[str], local_model: bool = True):
    # コマンドライン引数から画像パスを受け取り、予測を行う
    predictions = TrocrPredictor(local_model).predict_for_image_paths(image_paths)
    for path, (prediction, confidence) in zip(image_paths, predictions):
        print(f"Path:\t\t{path}\nPrediction:\t{prediction}\nConfidence:\t{confidence}\n")
```

cli.py

## 1.コードの理解を進める

```
def main_train(use_local_model: bool = False):  
    processor = load_processor()  
    train_dataset = HCRDataset(paths.train_dir, processor)  
    train_dataloader = DataLoader(train_dataset, constants.batch_size, shuffle=True, num_workers=constants.num_workers)  
  
    val_dataset = HCRDataset(paths.val_dir, processor)  
    val_dataloader = DataLoader(val_dataset, constants.batch_size, num_workers=constants.num_workers)  
  
    model = load_model(use_local_model)  
    init_model_for_training(model, processor)  
  
    context = Context(model, processor, train_dataset, train_dataloader, val_dataset, val_dataloader)  
    train(context, constants.train_epochs)  
    debug_print(f"Saving model to {paths.model_path}...")  
    model.save_pretrained(paths.model_path)
```

main.py(一部抜粋)

## 1.コードの理解を進める

processor中身

```

1  ✓ [
2  ✓   ViTImageProcessor {
3     "do_normalize": true,
4     "do_rescale": true,
5     "do_resize": true,
6     "image_mean": [0.5, 0.5, 0.5],
7     "image_processor_type": "ViTImageProcessor",
8     "image_std": [0.5, 0.5, 0.5],
9     "resample": 2,
10    "rescale_factor": 0.00392156862745098,
11  ✓    "size": {
12      "height": 384,
13      "width": 384
14    }
15  },
16  ✓  RobertaTokenizerFast {
17    "name_or_path": "microsoft/trocr-base-handwritten",
18    "vocab_size": 50265,
19    "model_max_length": 512,
20    "is_fast": true,
21    "padding_side": "right",
22    "truncation_side": "right",
23  ✓    "special_tokens": {
24      "bos_token": "<s>",
25      "eos_token": "</s>",
26      "unk_token": "<unk>",
27      "sep_token": "</s>",
28      "pad_token": "<pad>",
29      "cls_token": "<s>",
30      "mask_token": "<mask>"
31    },
32    "clean_up_tokenization_spaces": true
33  }
34  ]
35

```

[illegible]

# 1.コードの理解を進める

data中身

train\_datasetはHCRDatasetのインスタンス化で  
ディレクトリtrain/,val,内の画像読み込みetc

train:9288,val:1030

```
Loaded 9288 samples from C:\Users\user\ai\trocr\train
9288
dict_keys(['idx', 'input', 'label'])
torch.Size([3, 384, 384])
torch.Size([71])
9
Loaded 1030 samples from C:\Users\user\ai\trocr\val
torch.utils.data.dataloader.DataLoader object at 0x00000207F5AA6070>
Loaded 9288 training samples and 1030 validation samples.
['idx': 0, 'input': tensor([[[[0.9765, 0.9765, 0.9529, ..., 0.9922, 0.9922, 0.9922],
[0.9765, 0.9765, 0.9529, ..., 0.9922, 0.9922, 0.9922],
[0.9765, 0.9765, 0.9529, ..., 0.9922, 0.9922, 0.9922],
...,
[0.9608, 0.9765, 0.9843, ..., 0.9843, 0.9608, 0.9765],
[0.9608, 0.9765, 0.9843, ..., 0.9843, 0.9608, 0.9765],
[0.9608, 0.9765, 0.9843, ..., 0.9843, 0.9608, 0.9765]],
[[[0.9765, 0.9765, 0.9529, ..., 0.9922, 0.9922, 0.9922],
[0.9765, 0.9765, 0.9529, ..., 0.9922, 0.9922, 0.9922],
[0.9765, 0.9765, 0.9529, ..., 0.9922, 0.9922, 0.9922],
...,
[0.9608, 0.9765, 0.9843, ..., 0.9843, 0.9608, 0.9765],
[0.9608, 0.9765, 0.9843, ..., 0.9843, 0.9608, 0.9765],
[0.9608, 0.9765, 0.9843, ..., 0.9843, 0.9608, 0.9765]],
[[[0.9765, 0.9765, 0.9529, ..., 0.9922, 0.9922, 0.9922],
[0.9765, 0.9765, 0.9529, ..., 0.9922, 0.9922, 0.9922],
[0.9765, 0.9765, 0.9529, ..., 0.9922, 0.9922, 0.9922],
...,
[0.9608, 0.9765, 0.9843, ..., 0.9843, 0.9608, 0.9765],
[0.9608, 0.9765, 0.9843, ..., 0.9843, 0.9608, 0.9765],
[0.9608, 0.9765, 0.9843, ..., 0.9843, 0.9608, 0.9765]],
[[[0.9765, 0.9765, 0.9529, ..., 0.9922, 0.9922, 0.9922],
[0.9765, 0.9765, 0.9529, ..., 0.9922, 0.9922, 0.9922],
[0.9765, 0.9765, 0.9529, ..., 0.9922, 0.9922, 0.9922],
...,
[0.9608, 0.9765, 0.9843, ..., 0.9843, 0.9608, 0.9765],
[0.9608, 0.9765, 0.9843, ..., 0.9843, 0.9608, 0.9765],
[0.9608, 0.9765, 0.9843, ..., 0.9843, 0.9608, 0.9765]]]), 'label': tensor([ 0, 35507, 76
35, 6, 56, 47, 3033, 10, 367, 11505,
536, 4, 2, 1, 1, 1, 1, 1, 1, 1,
1, 1, 1, 1, 1, 1, 1, 1, 1, 1,
1, 1, 1, 1, 1, 1, 1, 1, 1, 1,
1, 1, 1, 1, 1, 1, 1, 1, 1, 1,
1, 1, 1, 1, 1, 1, 1, 1, 1, 1,
1, 1, 1, 1, 1, 1, 1, 1, 1, 1,
1])}]
```

python

📋 コピーする ✎ 編集する

```
train_dataloader = DataLoader(
    train_dataset,                # → HCRDataset ( __getitem__ あり)
    constants.batch_size,        # → 1回に取り出すサンプル数 (例: 16, 32)
    shuffle=True,                # → 各エポックでサンプルをランダムに並び替える
    num_workers=constants.num_workers # → 並列でデータを読み込むワーカープロセス数
)
```

# 1.コードの理解を進める

model中身

```
VisionEncoderDecoderModel(  
  (encoder): ViTModel(  
    (embeddings): ViTEmbeddings(  
      (patch_embeddings): ViTPatchEmbeddings(  
        (projection): Conv2d(3, 768, kernel_size=(16, 16), stride=(16, 16))  
      )  
    )  
    (dropout): Dropout(p=0.0, inplace=False)  
  )  
  (encoder): ViTEncoder(  
    (layer): ModuleList(  
      (0): ViTLayer(  
        (attention): ViTAttention(  
          (self_attn): ViTMultiheadAttention(  
            (key_proj): Linear(in_features=768, out_features=768, bias=True)  
            (value_proj): Linear(in_features=768, out_features=768, bias=True)  
            (q_proj): Linear(in_features=768, out_features=768, bias=True)  
            (v_proj): Linear(in_features=768, out_features=768, bias=True)  
            (attn_dropout): Dropout(p=0.0, inplace=False)  
            (proj): Linear(in_features=768, out_features=768, bias=True)  
            (dropout): Dropout(p=0.0, inplace=False)  
            (activation_dropout): Dropout(p=0.0, inplace=False)  
            (activation): Gelu()  
            (layer_norm1): LayerNorm((768,), eps=1e-12, elementwise_affine=True)  
            (layer_norm2): LayerNorm((768,), eps=1e-12, elementwise_affine=True)  
            (dropout2): Dropout(p=0.0, inplace=False)  
            (fc1): Linear(in_features=768, out_features=768, bias=True)  
            (fc2): Linear(in_features=768, out_features=768, bias=True)  
            (dropout3): Dropout(p=0.0, inplace=False)  
            (activation2): Gelu()  
            (proj2): Linear(in_features=768, out_features=768, bias=True)  
            (dropout4): Dropout(p=0.0, inplace=False)  
          )  
        )  
      )  
    )  
  )  
  (pooler): ViTPooler(  
    (dense): Linear(in_features=768, out_features=768, bias=True)  
    (activation): Tanh()  
  )  
)
```

• • •

```
    )  
    (layernorm_before): LayerNorm((768,), eps=1e-12, elementwise_affine=True)  
    (layernorm_after): LayerNorm((768,), eps=1e-12, elementwise_affine=True)  
  )  
)  
)  
)  
(layernorm): LayerNorm((768,), eps=1e-12, elementwise_affine=True)  
(pooler): ViTPooler(  
  (dense): Linear(in_features=768, out_features=768, bias=True)  
  (activation): Tanh()  
)  
)
```

• • •

```
(decoder): TrOCRForCausalLM(  
  (model): TrOCRDecoderWrapper(  
    (decoder): TrOCRDecoder(  
      (embed_tokens): Embedding(50265, 1024, padding_idx=1)  
      (embed_positions): TrOCRLearnedPositionalEmbedding(514, 1024)  
      (layernorm_embedding): LayerNorm((1024,), eps=1e-05, elementwise_affine=True)  
      (layers): ModuleList(  
        (0): TrOCRDecoderLayer(  
          (self_attn): TrOCRAttention(  
            (k_proj): Linear(in_features=1024, out_features=1024, bias=True)  
            (v_proj): Linear(in_features=1024, out_features=1024, bias=True)  
            (q_proj): Linear(in_features=1024, out_features=1024, bias=True)  
            (out_proj): Linear(in_features=1024, out_features=1024, bias=True)  
          )  
        )  
      )  
    )  
  )  
)
```

• • •

```
    (final_layer_norm): LayerNorm((1024,), eps=1e-05, elementwise_affine=True)  
  )  
)  
)  
)  
)  
(output_projection): Linear(in_features=1024, out_features=50265, bias=False)  
)  
)
```

2.いろいろな人が書いた筆記スタイルで精度評価

I ate an apple. pen by和住さん

I have a banana. hoge byきよや

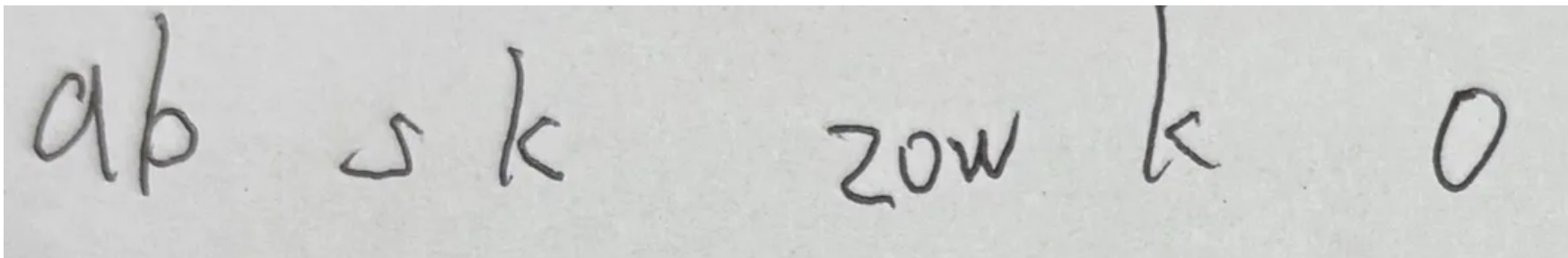
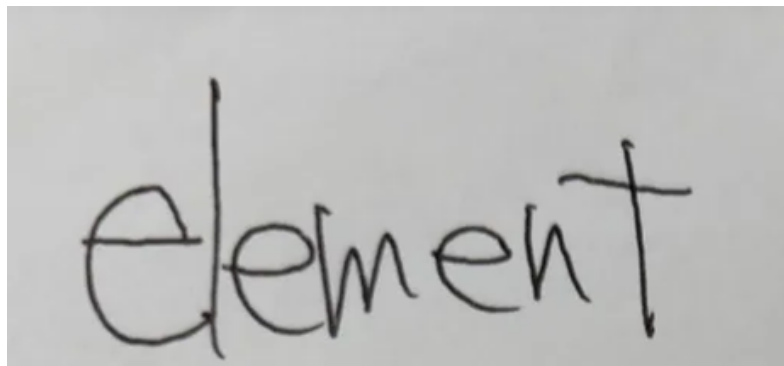
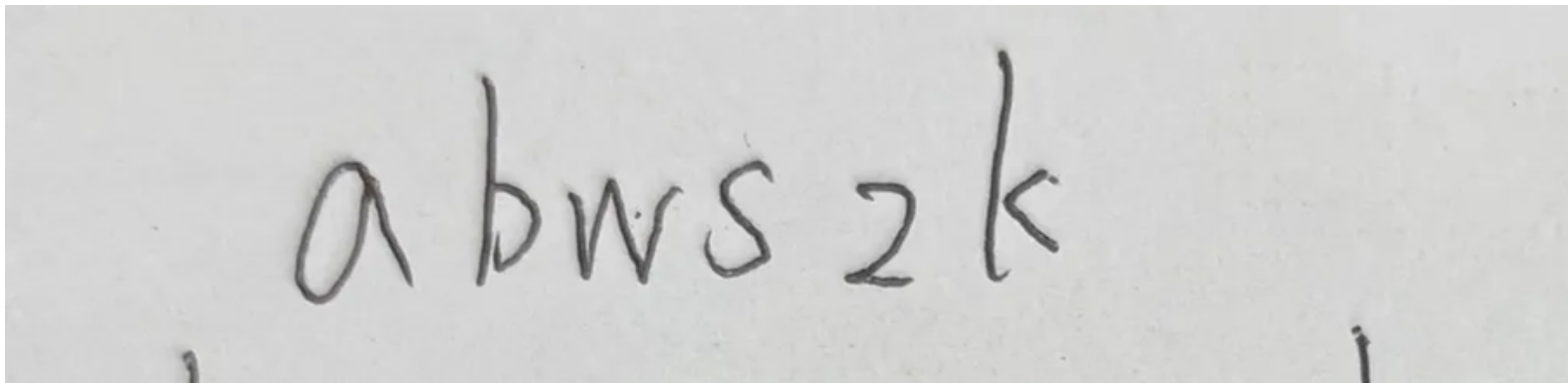
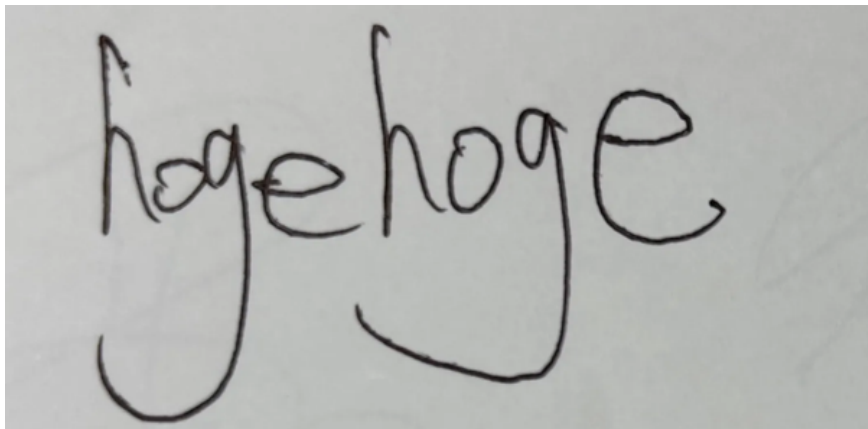
eri prosess by岡本さん

Bob I like industry

14	ata\test\image3.png	I ate an apple.
15	ata\test\image4.png	I. have a banana,
16	ata\test\image5.png	er I
17	ata\test\image6.png	Boe'_br_br_br_br_gr_gr_gr_
18	ata\test\image7.png	pen
19	ata\test\image8.png	I was a student.
20	ata\test\image9.png	hope
21	ata\test\imagee.png	I
22	ata\test\imageeee.png	I like industry



2.いろいろな人が書いた筆記スタイルで精度評価



by俺

byきよや

	column 1	column 2	column 3
1	image_path	prediction	confidence
2	data\test\text_htr\imagee1.png	lugehope	0.514011561870575
3	data\test\text_htr\imagee2.png	element	0.910437822341919
4	data\test\text_htr\imagee3.png	, abwszk,,,	0.5754343867301941
5	data\test\text_htr\imagee4.png	Ab.s k-zow k-O	0.26034966111183167

### 3.生成したテキストが正しいかどうかの処理をしているのか(入っていたらオフにする)

```
def predict(
    processor: TrOCRProcessor, model: VisionEncoderDecoderModel, dataloader: DataLoader
) -> tuple[list[tuple[int, str]], list[float]]:
    output: list[tuple[int, str]] = []
    confidence_scores: list[tuple[int, float]] = []

    with torch.no_grad():
        model.eval()
        for i, batch in enumerate(dataloader):
            debug_print(f"Predicting batch {i+1}")
            inputs: torch.Tensor = batch["input"].to(constants.device)

            generated_ids = model.generate(inputs, return_dict_in_generate=True, output_scores = True)
            generated_text = processor.batch_decode(generated_ids.sequences, skip_special_tokens=True)

            ids = [t.item() for t in batch["idx"]]
            output.extend(zip(ids, generated_text))

            # Compute confidence scores
            batch_confidence_scores = get_confidence_scores(generated_ids)
            confidence_scores.extend(zip(ids, batch_confidence_scores))

    return output, confidence_scores

def get_confidence_scores(generated_ids) -> list[float]:
    # Get raw logits, with shape (examples,tokens,token_vals)
    logits = generated_ids.scores
    logits = torch.stack(list(logits),dim=1)

    # Transform logits to softmax and keep only the highest (chosen) p for each token
    logit_probs = F.softmax(logits, dim=2)
    char_probs = logit_probs.max(dim=2)[0]

    # Only tokens of val>2 should influence the confidence. Thus, set probabilities to 1 for tokens 0-2
    mask = generated_ids.sequences[:, :-1] > 2
    char_probs[mask] = 1

    # Confidence of each example is cumulative product of token probs
    batch_confidence_scores = char_probs.cumprod(dim=1)[: , -1]
    return [v.item() for v in batch_confidence_scores]
```

- 文法的正確性
- 単語の妥当性
- 文脈の一貫性

上記のような処理はしてなさそう

単純にモデルの出力ロジットを単語に置き換えてるっぽい

## 4. 今後の展望

OCRなどのファインチューニング用に、日本語文が入ったテキストファイル（例：1行1文 × 10,000行）を入れてみて日本語htrモデルにしてみる？？できるのか

画像が中央に来ていないときに精度悪いのかも？？  
もっと回りきたない場合は？？

二行以降もできるようにした方がよき？？

# trocrの論文読んでみた

TrOCRモデルでは、エンコーダの初期化に事前学習済みのDeiT（Touvronら、2021）および  
BEiT  
（Bao、Dong、Wei、2021）モデルが使用