

Homework 2 – Updated!

Joshua Lumpkin

Vocabulary:

Saved as a JSON, vocab size of 3528

```
▼ root:
  ▼ itos:
    0: "<PAD>"
    1: "<BOS>"
    2: "<EOS>"
    3: "<UNK>"
    4: "a"
    5: "horse"
    6: "on"
    7: "woman"
    8: "the"
    9: "her"
    10: "head"
    11: "head."
    12: "under"
    13: "and"
    14: "she"
    15: "between"
    16: "legs"
    17: "gets"
    18: "goes"
    19: "is"
    20: "horse."
    21: "pooped"
```

Example of preprocessing, tokenizing. I do append BOS and EOS to the caption as well

Original: A woman goes under a horse.

Tokenized: ['a', 'woman', 'goes', 'under', 'a', 'horse.']

Numericalized: [4, 7, 18, 12, 4, 20]

Max caption length: 42

Batch video features shape: torch.Size([32, 80, 4096])

Batch captions shape: torch.Size([32, 42])

Training setup:

- Epochs = 30
- Learning rate = .001
- Batch size = 32

```
S2VTModel(
  (encoder_lstm): LSTM(4096, 500, num_layers=2, batch_first=True, dropout=0.5)
  (decoder_lstm): LSTM(500, 500, num_layers=2, batch_first=True, dropout=0.5)
  (fc): Linear(in_features=500, out_features=3529, bias=True)
)
```

Results:

I am continually working on improving the output. I am displeased with the results. I could not get it to not be fixated on singular terms. It seemed to always converge on a similar sentence structure, like what is shown in the picture. I attempted 5-6 model variations but could not get past this. Some variations I tried: model with attention, singular feature to caption, extensive training with all captions per video, simple models, stepLR. I did not have time to implement beam search.

[illegible]