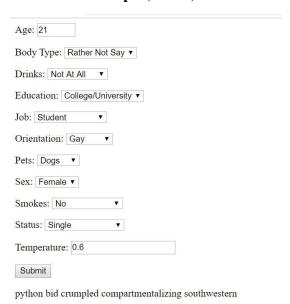
Generating Online Dating Profiles

Pushpak Raj Gautam

Description

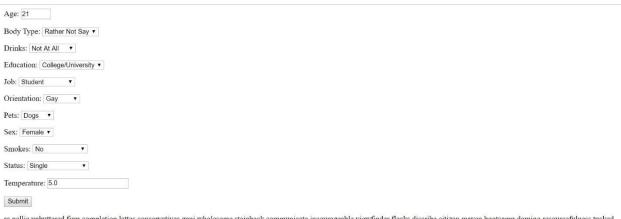
DateGen is a tool that generates an introduction bio for a person based on their attributes like age, sex, social habits, job, education, et cetera.

A successful example (sort of)



The sentence above doesn't make much sense but is a bit coherent. This example generates a lot of words and clearly gets to <End of Sentence> (word limit is 50). This is an LSTM model and for each word generation, the set of attributes along with the previous word is fed in as input. The model contains a hidden layer and a linear layer.

A failed example



rc nellie unbuttered firm completion lattes conservatives ganj wholesome steinbeck communicate incourageable viewfinder flecks discribe citizen maven bootcamp domino resourcefulness tucked depart scaredy fountain restful pupusas bckward virgin remembering contrasts preteen incandescent chesszy planetarium astute morcheba droles abstain sinner tortuously settlers carport billy skin unfalteringly openness latinos tube bind except

The example above is generated using a temperature of 5.0 (in the softmax layer). Thus, the words are way more random and the output hits the limit of 50 words.

Creativity

The paper that I based my work on used character-level output generation. When I tried this, my output wasn't legible (due to the dataset's small size). So, I changed this to a word level output generation model. By using lower cases and removing punctuation marks, I got to a vocabulary size of around 28k words. The other attributes, when one-hot encoded, gave a 102-dimensional vector. The interesting thing is that I was able to one-hot encode everything (words and attributes) on the GPU and still get a reasonable memory consumption. In the output below (that was generated while training), one can see the shape which corresponds to a word limit of 51 words, a batch of 64 points and an input vector of 28623 elements.

```
Total Loss: tensor(9.5215, device='cuda:0')
Batch: 20
(51, 64, 28623)
Total Loss: tensor(9.4785, device='cuda:0')
Batch: 21
(51, 64, 28623)
Total Loss: tensor(9.4066, device='cuda:0')
Batch: 22
(51, 64, 28623)
Total Loss: tensor(9.2980, device='cuda:0')
Batch: 23
(51, 64, 28623)
Total Loss: tensor(9.3027, device='cuda:0')
Batch: 24
(51, 64, 28623)
Total Loss: tensor(9.3486, device='cuda:0')
Batch: 25
(51, 64, 28623)
Total Loss: tensor(8.9415, device='cuda:0')
Batch: 26
(51, 64, 28623)
Total Loss: tensor(8.9328, device='cuda:0')
Batch: 27
(51, 64, 28623)
Total Loss: tensor(8.7440, device='cuda:0')
```

Secondly, I created a web interface for my model so that it could be easier to display the generated output. I did this using Django. Here, I also added a temperature option so that I could show how the model behaves with different degrees of randomness.