

Final Project: Character-Level Language Modeling for Text Generation via Deep Markov Models

Armaan Kohli - ECE467 Natural Language Processing
Spring 2020

Remarks

We attempted to use a deep markov model (DMM) to make a character-level language model. This work is based on recent developments in the understanding of discrete time series, such as MIDI, as well as natural language processing. Using a DMM, we were able to generate text that yielded quantitative performance approaching state of the art for character-based language models. However, training remains unstable and more research into DMMs is likely required for their performance for language models to improve.

Deep Markov Models

Traditional markov models are a method representing complex temporal dependencies in observed data. A markov model has a chain of latent variables, with each latent (or hidden) variable in the chain is conditioned on the previous latent variable. This is a useful approach, but if we want to represent complex data with complex dynamics, such as text, we would like to be able to model dynamics that are potentially highly non-linear.

This brings forth the idea of a deep markov model, wherein we allow the transition probabilities governing the dynamics of the latent variables as well as the emission probabilities that govern how the observations are generated by the latent dynamics to be parametrized by (non-linear) neural networks. DMMs were first used in the setting of polyphonic music generation. Using a MIDI representation of musical notes, Krishnan et. al were able to generate high-quality songs and learn a representation of electronic health record data [1].

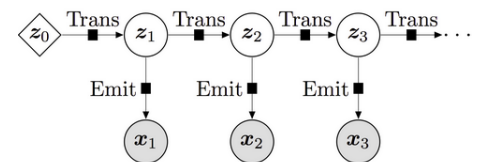


Figure 1: An illustration of a DMM. Each of the black squares represent an RNNs that determine the probability of emission or transmission. Image replicated from Pyro documentation. [ADD CITATION]

Even though this method was originally designed for music generation, character-level language models can be thought of in a similar way. At each time step, music can be represented by an 88-dimensional binary vector. Similarly, characters in a phrase can be represented by a one-hot vector with a dimension given by the size of the learned dictionary. Research by the Harvard Intelligent Probabilistic Systems (HIPS) group takes a similar approach, using the a neural network for both polyphonic music generation and character-level language modeling, the only change being the distribution from which the data is drawn from, the observation likelihood (Bernoulli vs categorical) [2]. HIPS uses a generative flow model for character-level language modelling as opposed to a DMM, however. The inference strategy we're going to use called variational inference (VI), which requires specifying a parametrized family of distributions that can be used to approximate the posterior distribution over the latent random variables. Due to the complex temporal relations we seek to model, we can expect the posterior distribution to be highly non-trivial, necessitating a probabilistic approach. Thus, we use PyTorch as

our choice of deep learning framework, as well as Pyro, a probabilistic programming language integrated into PyTorch to effectively sample and perform VI on our model.

Implementation Details

We use a single-layer RNN for our emission and transmission probabilities. Our objective function is the ELBO (evidence-based lower bound) with a KL-annealing term β , inspired by [3].

$$\mathcal{F}(\theta, \phi, \beta; \mathbf{x}, \mathbf{z}) \geq \mathcal{L}(\theta, \phi; \mathbf{x}, \mathbf{z}) = \mathbb{E}_{q_{\phi}(\mathbf{z}|\mathbf{x})}[\log_{p_{\theta}}(\mathbf{x}|\mathbf{z})] - \beta D_{KL}(q_{\phi}(\mathbf{z}|\mathbf{x})||p(\mathbf{z})) \quad (1)$$

We use Monte Carlo estimates of the KL divergence term.

We train the language model using the Penn Treebank (PTB) corpus. We perform treat every line in the corpus as a distinct sequence, or sentence, and tokenize each character in each sentence, adding the <unk> token for low-frequency or unknown words, and <eos> to demarcate the end of a sentence. The size of the dictionary was 52 We opt for a batch size of 16. In order to generate a character embedding, we simply encoded our character dictionary as a one-hot 52-dimensional vector. This was an appropriate choice due to the small dictionary size. For full details see github.com/armaank/textDMM for the full codebase and the parameters used to train the network.

Results & Discussion

We were able to

These aforementioned issues might be resolved by using a different method for KL annealing, which can improve stability during training. Furthermore, we use a Monte Carlo estimates of the KL divergence, leading to higher variance gradient estimates of the ELBO loss, which can also destabilize performance during early training periods. We might also trying using an LSTM architecture to parametrize our transmission and emission probabilities over the so-called ‘vanilla’ RNN. On a related note, one possibility is that exploding gradients are caused by lengthy input sequences, so one way to resolve this issue would be to only train on shorter sequences of characters.

Conclusion

In conclusion, we were able to successfully train a DMM as a character-level language model and achieve performance close to that of traditional character-level language models using purely RNNs/LSTMs. However, though more research is needed to improve DMMs for NLP tasks. The power of DMMs and other probabilistic models is their flexibility, in that the same model can generate MIDI music, missing EHR data and text with only the changing distribution governing the observation likelihood.

References

- [1] R. G. Krishnan, U. Shalit, and D. Sontag, “Structured inference networks for nonlinear state space models,” 2016.

- [2] Z. M. Ziegler and A. M. Rush, "Latent normalizing flows for discrete sequences," 2019.
- [3] I. Higgins, L. Matthey, A. Pal, C. Burgess, X. Glorot, M. M. Botvinick, S. Mohamed, and A. Lerchner, "beta-vae: Learning basic visual concepts with a constrained variational framework," in *ICLR*, 2017.

Appendix A: Code

The code below is dmm.py, the main model code.

```

1  """dmm
2  """
3  import argparse
4  import os
5
6  import numpy as np
7  import torch
8  import torchtext
9  import pyro
10
11 import torch.nn as nn
12 import pyro.distributions as dist
13 import pyro.poutine as poutine
14
15 from torch.autograd import Variable
16 from pyro.distributions import TransformedDistribution
17
18 import utils
19
20
21 class Emitter(nn.Module):
22     """
23     parameterizes the categorical observation likelihood  $p(x_t|z_t)$ 
24     """
25
26     def __init__(self, input_dim, z_dim, emission_dim):
27         super().__init__()
28         """
29         initilize the fcns used in the network
30         """
31         self.lin_z_to_hidden = nn.Linear(z_dim, emission_dim)
32         self.lin_hidden_to_hidden = nn.Linear(emission_dim, emission_dim)
33         self.lin_hidden_to_input = nn.Linear(emission_dim, input_dim)
34         self.relu = nn.ReLU()
35
36         pass
37
38     def forward(self, z_t):
39         """
40         given  $z_t$ , compute the probabilities that parameterizes the categorical distribution  $p(x_t|z_t)$ 
41         """
42         h1 = self.relu(self.lin_z_to_hidden(z_t))
43         h2 = self.relu(self.lin_hidden_to_hidden(h1))
44         probs = torch.sigmoid(
45             self.lin_hidden_to_input(h2)
46         ) # might need to change to argmax, max?, softmax?
47

```

```

48     return probs
49
50
51 class GatedTransition(nn.Module):
52     """
53     parameterizes the gaussian latent transition probability  $p(z_t | z_{t-1})$ 
54     """
55
56     def __init__(self, z_dim, transition_dim):
57         super().__init__()
58         """
59         initilize the fcns used in the network
60         """
61         self.lin_gate_z_to_hidden = nn.Linear(z_dim, transition_dim)
62         self.lin_gate_hidden_to_z = nn.Linear(transition_dim, z_dim)
63         self.lin_proposed_mean_z_to_hidden = nn.Linear(z_dim, transition_dim)
64         self.lin_proposed_mean_hidden_to_z = nn.Linear(transition_dim, z_dim)
65         self.lin_sig = nn.Linear(z_dim, z_dim)
66         self.lin_z_to_loc = nn.Linear(z_dim, z_dim)
67
68         self.lin_z_to_loc.weight.data = torch.eye(z_dim)
69         self.lin_z_to_loc.bias.data = torch.zeros(z_dim)
70
71         self.relu = nn.ReLU()
72         self.softplus = nn.Softplus()
73
74     pass
75
76     def forward(self, z_t_1):
77         """
78         Given the latent  $z_{t-1}$  we return the mean and scale vectors that parameterize the
79         (diagonal) gaussian distribution  $p(z_t | z_{t-1})$ 
80         """
81         # compute the gating function
82         _gate = self.relu(self.lin_gate_z_to_hidden(z_t_1))
83         gate = torch.sigmoid(self.lin_gate_hidden_to_z(_gate))
84         # compute the 'proposed mean'
85         _proposed_mean = self.relu(self.lin_proposed_mean_z_to_hidden(z_t_1))
86         proposed_mean = self.lin_proposed_mean_hidden_to_z(_proposed_mean)
87         # assemble the actual mean used to sample  $z_t$ , which mixes a linear transformation
88         # of  $z_{t-1}$  with the proposed mean modulated by the gating function
89         loc = (1 - gate) * self.lin_z_to_loc(z_t_1) + gate * proposed_mean
90         # compute the scale used to sample  $z_t$ , using the proposed mean from
91         # above as input the softplus ensures that scale is positive
92         scale = self.softplus(self.lin_sig(self.relu(proposed_mean)))
93         # return loc, scale which can be fed into Normal
94         return loc, scale
95
96
97 class Combiner(nn.Module):
98     """

```

parameterizes $q(z_t | z_{t-1}, x_{t:T})$, which is the basic building block of the guide (i.e. the variational distribution). The dependence on $x_{t:T}$ is through the hidden state of the RNN

"""

def __init__(self, z_dim, rnn_dim):

super().__init__()

 """

 initilize the fcns used in the network

 """

 self.lin_z_to_hidden = nn.Linear(z_dim, rnn_dim)

 self.lin_hidden_to_loc = nn.Linear(rnn_dim, z_dim)

 self.lin_hidden_to_scale = nn.Linear(rnn_dim, z_dim)

 self.tanh = nn.Tanh()

 self.softplus = nn.Softplus()

pass

def forward(self, z_t_1, h_rnn):

 """

 Given the latent z_{t-1} at a particular time as well as the hidden state of the RNN $h(x_{t:T})$ we return the mean and scale vectors that parameterize the (diagonal) gaussian distribution $q(z_t | z_{t-1}, x_{t:T})$

 """

 # combine the rnn hidden state with a transformed version of z_{t-1}

 h_combined = 0.5 * (self.tanh(self.lin_z_to_hidden(z_t_1)) + h_rnn)

 # use the combined hidden state to compute the mean used to sample z_t

 loc = self.lin_hidden_to_loc(h_combined)

 # use the combined hidden state to compute the scale used to sample z_t

 scale = self.softplus(self.lin_hidden_to_scale(h_combined))

 # return loc, scale which can be fed into Normal

return loc, scale

class DMM(nn.Module):

 """

 module for the model and the guide (variational distribution) for the DMM

 """

def __init__(

 self,

 input_dim=52,

 z_dim=100,

 emissions_dim=100,

 transition_dim=200,

 rnn_dim=600,

 num_layers=1,

 dropout=0.0,

):

super().__init__()

 """

```

150     instantiate modules used in the model and guide
151     """
152     self.emitter = Emitter(input_dim, z_dim, emission_dim)
153     self.transition = GatedTransition(z_dim, transition_dim)
154     self.combiner = Combiner(z_dim, rnn_dim)
155
156     # TODO: alter dropout scheme
157     if num_layers == 1:
158         rnn_dropout = 0.0
159     else:
160         rnn_dropout = dropout
161
162     # TODO: add option for bidirectional rnn?
163     self.rnn = nn.RNN(
164         input_size=input_dim,
165         hidden_size=rnn_dim,
166         nonlinearity="relu",
167         batch_first=True,
168         bidirectional=False,
169         num_layers=num_layers,
170         dropout=rnn_dropout,
171     )
172     """
173     define learned parameters that define the probability distributions  $P(z_{1:T})$  and  $q(z_{1:T})$  and hidden
174     state of rnn
175     """
176     self.z_0 = nn.Parameter(torch.zeros(z_dim))
177     self.z_q_0 = nn.Parameter(torch.zeros(z_dim))
178     self.h_0 = nn.Parameter(torch.zeros(1, 1, rnn_dim))
179
180     pass
181
182     def model(self, batch, reversed_batch, batch_mask, batch_seqlens, kl_anneal=1.0):
183         """
184         the model defines  $p(x_{1:T}|z_{1:T})$  and  $p(z_{1:T})$ 
185         """
186         # maximum duration of batch
187         Tmax = batch.size(1)
188
189         # register torch submodules w/ pyro
190         pyro.module("dmm", self)
191
192         # setup recursive conditioning for  $p(z_t|z_{1:t-1})$ 
193         z_prev = self.z_0.expand(batch.size(0), self.z_0.size(0))
194
195         # sample conditionally independent text across the batch
196         with pyro.plate("z_batch", len(batch)):
197             # sample latent vars z and observed x w/ multiple samples from the guide for each z
198             for t in pyro.markov(range(1, Tmax + 1)):
199                 # compute params of diagonal gaussian  $p(z_t|z_{1:t-1})$ 

```

```

200         z_loc, z_scale = self.trans(z_prev)
201
202         # sample latent variable
203         with poutine.scale(scale=kl_anneal):
204             z_t = pyro.sample(
205                 "z_%d" % t,
206                 dist.Normal(z_loc, z_scale)
207                 .mask(batch_mask[:, t - 1 : t])
208                 .to_event(1),
209             )
210
211         # compute emission probability from latent variable
212         emission_prob = self.emitter(z_t)
213
214         # observe x_t according to the Categorical distribution defined by the emitter
215         probability
216         pyro.sample(
217             "obs_x_%d" % t,
218             dist.OneHotCategorical(emission_prob)
219             .mask(batch_mask[:, t - 1 : t])
220             .to_event(1),
221             obs=batch[:, t - 1, :],
222         )
223
224         # set conditional var for next time step
225         z_prev = z_t
226     pass
227
228 def guide(self, batch, reversed_batch, batch_mask, batch_seqlens, kl_anneal=1.0):
229     """
230     the guide defines the variational distribution  $q(z_{1:T}|x_{1:T})$ 
231     """
232     # maximum duration of batch
233     Tmax = batch.size(1)
234
235     # register torch submodules w/ pyro
236     pyro.module("dmm", self)
237
238     # to parallelize, we broadcast rnn into contiguous gpu memory
239     h_0_contig = self.h_0.expand(
240         1, batch_size(0), self.rnn.hidden_size
241     ).contiguous()
242
243     # push observed sequence through rnn
244     rnn_output, _ = self.rnn(batch_reversed, h_0_contig)
245
246     # reverse and unpack rnn output
247     rnn_output = utils.pad_reverse(rnn_output, batch_seqlens)
248
249     # setup recursive conditioning
250     z_prev = self.z_q_0.expand(batch_size(0), self.z_q_0.size(0))

```



```
250
251 with pyro.plate("z_batch", len(mini_batch)):
252
253     for t in pyro.markov(range(1, Tmax + 1)):
254
255         z_loc, z_scale = self.combiner(z_prev, rnn_output[:, t - 1, :])
256
257         z_dist = dist.Normal(z_loc, z_scale)
258
259         assert z_dist.event_shape == ()
260         assert z_dist.batch_shape[-2:] == len(batch) == self.z_q_0.size(0)
261
262         # sample z_t from distribution z_dist
263         with pyro.poutine.scale(scale=kl_anneal):
264             z_t = pyro.sample(
265                 "z_%d" % t, z_dist.mask(batch[:, t - 1 : t]).to_event(1)
266             )
267
268         # set conditional var for next time step
269         z_prev = z_t
270
271     pass
```