

1 Deep Discrete Latent Variable Models - Exercises

2 Wilker Aziz

3 September 16, 2020

4 1 Analysing Twitter Data

5 We have a stream of tweets spanning a certain period of time. The data are
6 organised in daily batches of tweets, each batch containing exactly N tweets,
7 each tweet made of at most C characters.

8 Quite a bit of pre-processing has already taken place:

- 9 • Images have been automatically replaced by a token `image:class` where
10 *class* indicates one out of a finite number of categories meant to capture
11 the most salient property of an image;
- 12 • The tweets in the collection are all predominantly expressed in English (as
13 far as an automatic language detection software can tell us), but they do
14 contain hashtags, the codes for images, short URLs, emojis, slang, and all
15 sort of stuff you should expect from social media data, possibly including
16 foreign words.
- 17 • We have one version of the data where we applied tokenization, stripped
18 punctuation, and normalised URLs using some scheme that retains infor-
19 mation about the nature of the content (e.g., *news agency*, *government*
20 *agency*, *science*, etc.), we also applied a number of strategies to reduce
21 vocabulary size down to some manageable constant V_1 .
- 22 • A second version of the dataset has been segmented using an algorithm such
23 as BPE encoding or sentencepiece, and no other form of text normalisation
24 was employed. The vocabulary is some small enough constant V_2 .

25 As there are too many days in the data set, if necessary, it's okay to imagine
26 that days are grouped into overlapping sequences of size T . And, though these
27 sequences are not i.i.d., do pretend they are.

28 1.1 Part I

29 **Context** We want to analyse *trends* in terms of topics that people tweet about.
30 The focus of the analysis is to gain insight about the data we already have and
31 for now we do not have any future predictive task in mind. Initially we don't

32 really know what people talk about, so we would be happy enough to annotate
33 each daily batch with 1 of K latent topics. To gain insight about the topics, we
34 would be happy to see examples of tweets or some other means to ‘label’ the
35 topics for their semantic content.

36 **Task** Design a *tractable temporal model* that can be used to study this dataset.
37 Propose a factorisation of the model as well as its parameterisation, and do
38 employ NN architectures. Be explicit about your parametric choices and mind
39 the domain of their parameters. Present the objective for parameter estimation
40 via gradient-based methods, and discuss how to obtain key quantities such as
41 gradients or gradient estimates.

42 **Analysis** We are looking for insights into

- 43 • can we represent a topic by a list of keywords?
- 44 • how long do topics stay in the platform before going unmentioned?
- 45 • do topics re-emerge after going unmentioned?
- 46 • for topics that re-emerge, how long do they stay unmentioned?
- 47 • the second time a topic emerges, does it trend for a shorter period of time?

48 How can we answer these questions with a trained model?

49 1.2 Part II

50 **Context** We used the model of Part I to study the dataset. After that we
51 selected a sample of daily batches for annotation by humans. Those days have
52 now been annotated with a subset of D possible topics which are known to
53 be relevant throughout the period covered by the dataset. These topics are
54 not mutually exclusive, and the annotation is not associated with any tweet in
55 particular, but with a daily batch.

56 **Task** Adjust the temporal model to account for topics that are not mutually
57 exclusive in a day and to learn from all of the data, that is, the labelled
58 and the unlabelled part. Propose a factorisation of the model as well as its
59 parameterisation, and do employ NN architectures. If you need approximate
60 inference, first explain why that is the case, then go with variational inference and
61 propose a factorisation of the inference model as well as its parameterisation. Be
62 explicit about your parametric choices and mind the domain of their parameters.
63 Present the objective for parameter estimation via gradient-based methods, and
64 discuss how to obtain key quantities such as gradients or gradient estimates.
65 Make sure to account for all data points (observed and unobserved).

66 **Analysis** We are looking for insights into

67 • what is the distribution of number of topics per day?

68 • do certain topics co-occur above chance level?

69 • are there topics whose opposite trends correlate?

70 How can we answer these questions with a trained model?