

Time Series Analysis on real world datasets

Supervisor : Christian Tomani

Presenters : Carlos Garcia Briones

Jonas Glotz

Thana Guetet

Table of contents

1. Motivation
2. Goals
3. Datasets
4. Workflow
 - a. Data Pipelines
 - b. Model implementation
 - c. Trained models
5. Next Steps
6. Potential Challenges

1. Motivation

- Time dependent data analysis



- Uncertainty of Neural Networks



2. Goals

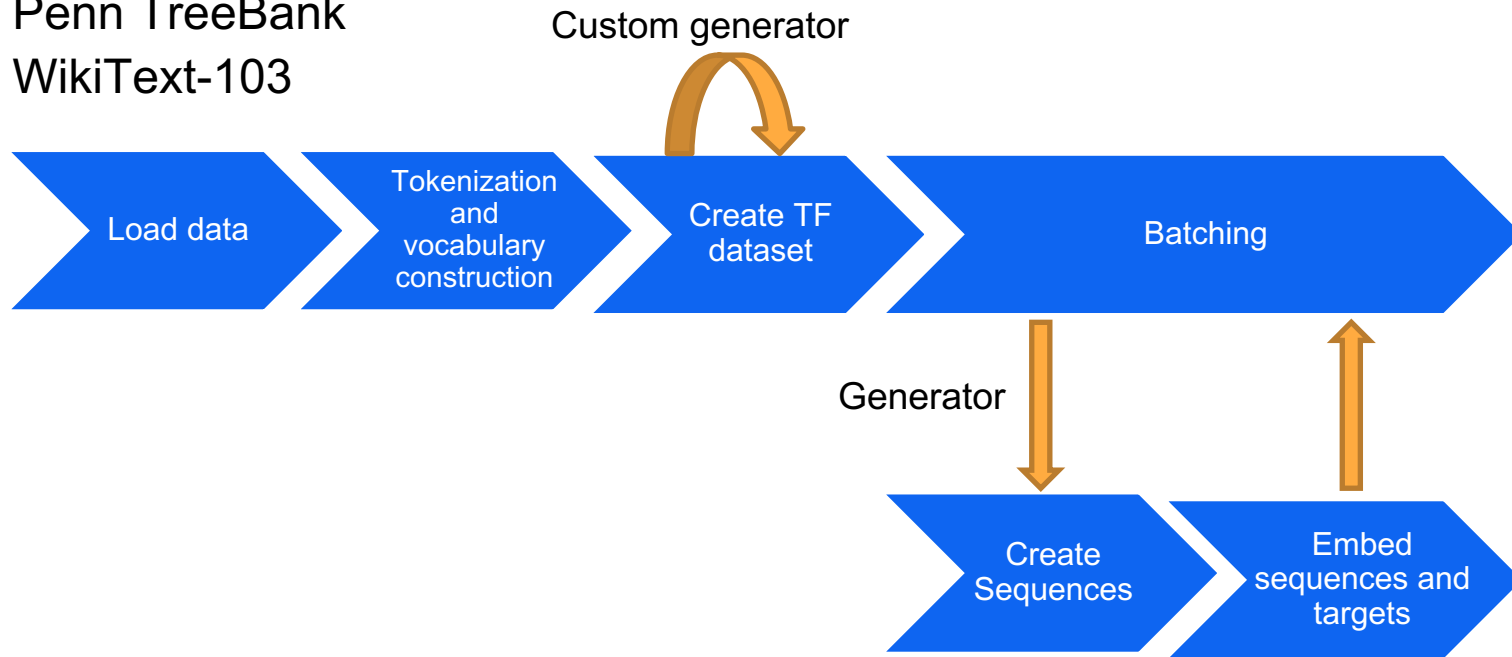
- Implement models with a relatively good performance.
 - Models: RNN, LSTM, GRU, IndRNN, Transformer-XL
- Create models able to deal with perturbed datasets.
- Understand and implement methods for increasing uncertainty awareness.

3. Datasets

- NLP Datasets
 - Penn TreeBank
 - WikiText-103
 - Lambada
- Music Datasets
 - JSB Chorales
 - Nottingham
- Sensor Dataset
 - InsectWingBeat

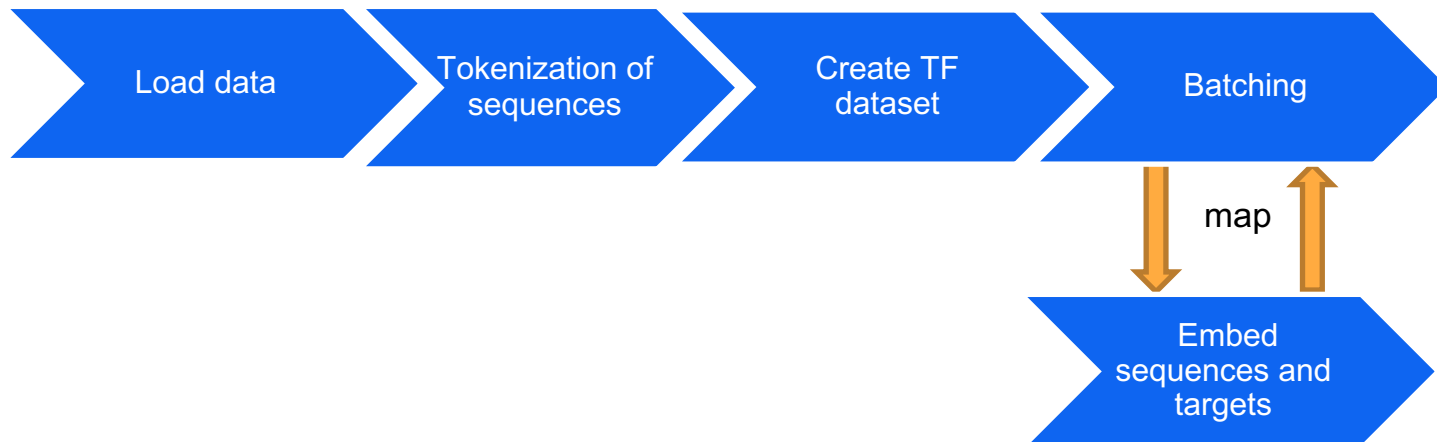
4. Workflow – Data Pipelines

- Penn TreeBank
- WikiText-103



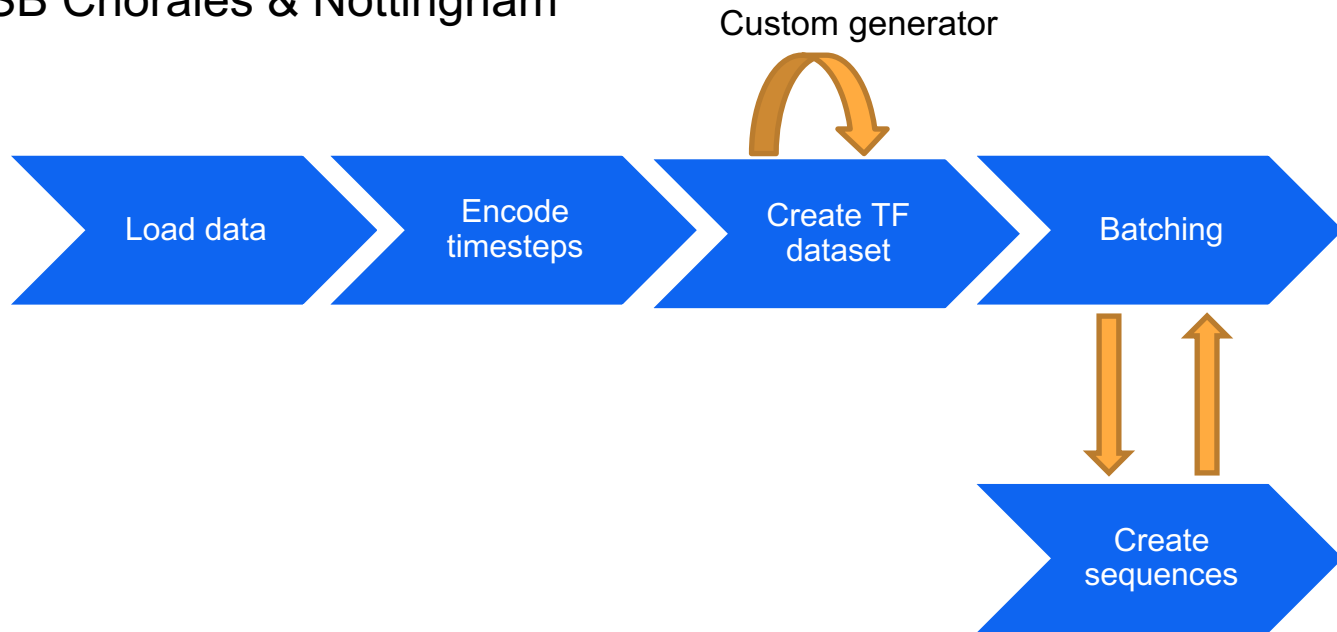
4. Workflow – Data Pipelines

- LAMBADA



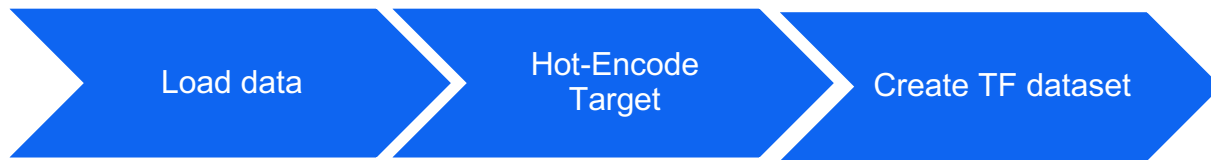
4. Workflow – Data Pipelines

- JSB Chorales & Nottingham



4. Workflow – Data Pipelines

- InsectWingBeat



4. Workflow – Model Implementation

- Using Keras Sequential API to implement baseline Models
 - RNN
 - LSTM
 - GRU
- Json configuration files
 - Add flexibility, documentation and effectiveness on training models with different datasets and architectures.
- Building an automated experiment documentation

4. Workflow – Model Implementation

- Example of documented experiments

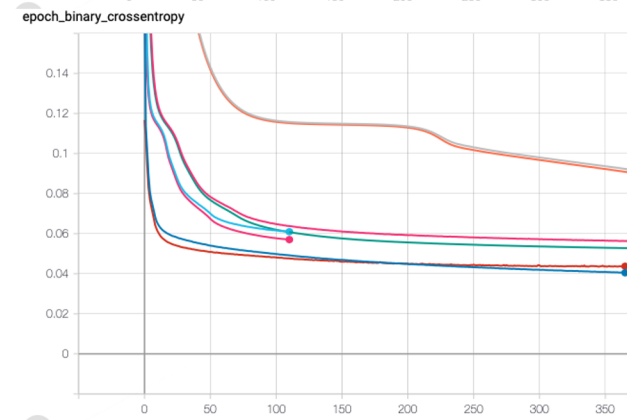
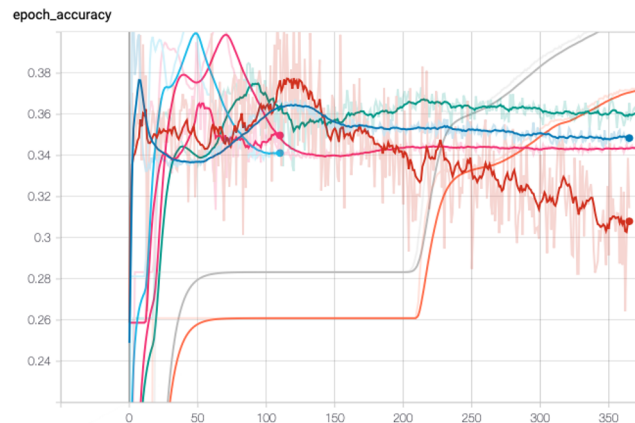
```
ptb_lstm_2020_05_24_08_39_35
IWB_lstm_3_2020_05_23_14_15_20
IWB_lstm_3_2020_05_23_14_12_32
nh_lstm_3_256_2020_05_23_13_19_30
Chorales_lstm_3_256_2020_05_23_13_1...
```

- Example of a Configuration file

```
1 {
2   "dataset": "InsectWingBeat",
3   "dataset_args": {...},
4   "model": "RNN",
5   "model_args": {
6     "recurrent_units": [256, 128, 128],
7     "linear_hidden_units": [
8       {"units": 128, "activation": "relu"},
9       {"activation": "softmax", "name": "output"}
10    ]
11  },
12  "learning_rate": 1e-3,
13  "compile": {
14    "loss": "categorical_crossentropy",
15    "optimizer": "adam",
16    "metrics": ["accuracy"]
17  },
18  "logdir": "log",
19  "experiment_name": "some_experiment",
20  "fit": {"epochs": 10, "verbose": 2}
21 }
```

4. Workflow – Trained Models

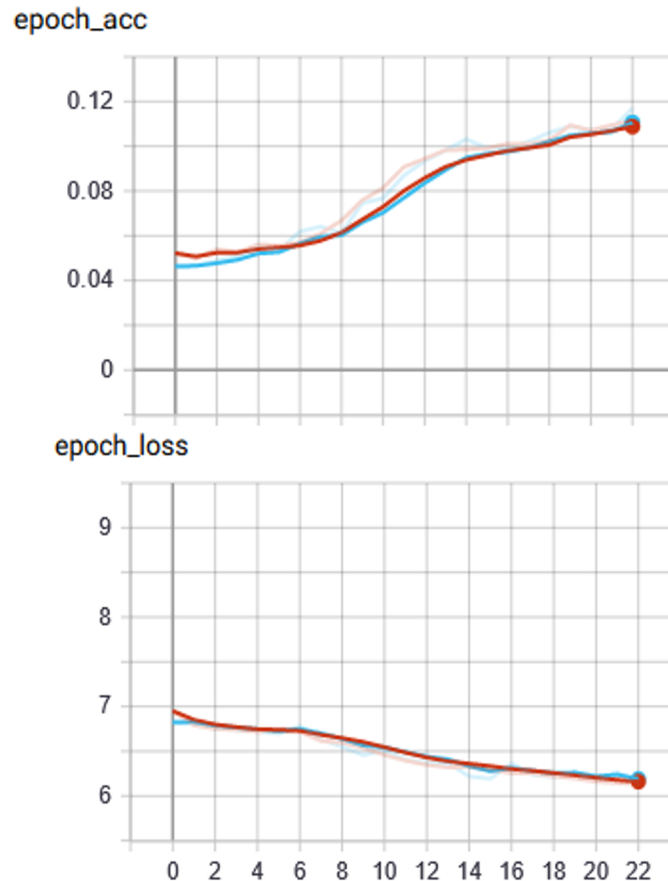
- Nottingham dataset
- Lstm model
- Hyperparameters:
 - Layers: 1,2
 - Hidden units: {256, 500}
 - Learning rate: {1e-4, 1e-5}



4. Workflow – Trained Models

- Penn TreeBank Dataset
- Validation perplexity of ~ 67.7
(= validation loss of 6.0804)

- $\text{perp.} = 2^{H(x)}$
- H is the (cross-)entropy,
i.e. the loss



5. Next steps

1. Implement the Transformer-XL and Independent Neural Network models.
2. Implement models that accept sequences with varying lengths to use with the Lambada dataset.
3. Use an informed grid search to determine better suitable hyperparameters and train the baseline models.
4. Train the newly implemented models on the chosen datasets.
5. Expand training script:
 - a. Use tensorboard to log other metrics (e.g. F1 Score)

6. Potential challenges

- Train models as big as Transformer-XL on the provided resources with limited time.
- The NLP datasets take a long time (>30 min per epoch) to train.
 - Specially hard to optimize with limited resources
- The Lambada dataset is fairly small in comparison to the other NLP datasets, while being the most challenging one

Thank you for listening!