

SoSe 2024

NLP-gestützte Data Science

Übung 1

Manuel Stoeckel Kevin Bönisch
Prof. Dr. Alexander Mehler

Frist: 06. Juni 2024

Übung 1: Word2Vec

50 P

In dieser Übung werden wir uns mit dem word2vec-Modell (Mikolov, Chen et al., 2013; Mikolov, Sutskever et al., 2013) beschäftigen, welches in PyTorch zu implementieren ist. Lesen Sie insb. das Paper *Distributed Representations of Words and Phrases and their Compositionality*, Mikolov, Sutskever et al. (2013)!

Definitionen

Das word2vec Skip-Gram-Modell für *SoftMax* und *Negative Sampling* (NS) sind gegeben als:

$$\frac{1}{T} \sum_{t=1}^T \sum_{c=t-C}^{t+C} F_{\text{obj}}(w_c, w_t) \quad | \quad 0 < c \leq T, c \neq t \quad (1)$$

$$F_{\text{SoftMax}}(w_o, w_l) = \log \frac{\exp(v'_{w_o} v_{w_l})}{\sum_{w=1}^W \exp(v'_w v_{w_l})} \quad (2)$$

$$F_{\text{NS}}(w_o, w_l) = \log \sigma(v'_{w_o} v_{w_l}) + \sum_{n=1}^k \log \sigma(v'_{w_n} v_{w_l}) \quad (3)$$

Hier ist T die Satzlänge, C die Kontextfenstergröße, v und v' sind jeweils die Input- und Output-Embeddings und (2) bzw. (3) sind die *objective functions* der *SoftMax* bzw. *Negative Sampling* Varianten. Die *negative samples* w_n in (3) werden zufällig aus der *unigram distribution* $w_n \sim P_V(w)$ gezogen, wobei $P(w_n)$ für ein bestimmtes Wort w_n mit Frequenz $f(w_n)$ im Trainingskorpus und Vokabulargröße V gegeben ist als:

$$P(w_n) = \frac{f(w_n)^{3/4}}{\sum_{j=0}^V f(w_j)^{3/4}} \quad (4)$$

Kontextwörter w_c werden dabei im Rahmen des *Subsampling of Frequent Words* (Mikolov, Sutskever et al., 2013, §2.3) mit Wahrscheinlichkeit $S(w_c)$ **nicht** aus einem Kontext entfernt, wobei $S(w_c)$ abweichend vom Paper wie folgt definiert ist:

$$S(w_c) = \left(\sqrt{\frac{f(w_c)}{t}} + 1 \right) \cdot \frac{t}{f(w_c)} \quad | \quad S(w_c) \in [0, 1] \quad (5)$$

Dabei ist t ein Hyper-Parameter, üblicherweise 0.001, und $f(w_c)$ die Frequenz wie zuvor.

1.1 Text Processing

20 P

- › Implementieren Sie die Klassen:
 - ›› **TokenizedSentence**, **PreTokenizer** und **Tokenizer**,
 - ›› das **Dataset**, in welchem die Samples für das aggregiert werden, sowie
 - ›› den **NegativeSampler** und die Dataset-Variante **DatasetNegativeSampling**.

Hinweise

- › Im Gegensatz zu einer Mindestanzahl an Vorkommen im Trainings-Korpus für den Tokenizer, wie im originalen word2vec, sollen Sie einen Tokenizer mit einer maximalen Größe implementieren. Dabei werden die häufigsten Token gespeichert und weniger häufige Token verworfen.
- › Der **PreTokenizer** soll in der Vorverarbeitung **alle Satzzeichen**¹ von anderen Wörtern **einzel**n abtrennen.
- › Achten Sie darauf auch das Subsampling von häufigen Wörtern zu implementieren.

1.2 Word2Vec

25 P

- › Implementieren Sie die Klassen **SkipGramSoftMax** und **SkipGramNegativeSampling**.

1.3 Ergebnisse

5 P

Verwenden Sie Ihren Code und dokumentieren Sie Ihre Ergebnisse (in einer begleitenden PDF).

- › Trainieren Sie einen Tokenizer an dem kleinen Beispiel-Korpus (`data/train_enwiki.txt.gz`). Was sind die häufigsten Token?
- › Trainieren Sie die Modelle auf dem tokenisierten Korpus. Können Sie die Ergebnisse von Mikolov et al. reproduzieren?
 - ›› Verwenden Sie Standard-Hyperparameter zum Training. Sie können die GPUs auf den Rechnern der RBI nutzen, um Ihr Training zu beschleunigen.
 - ›› Falls das Training zu lange dauert, können Sie es auch frühzeitig unterbrechen.
 - ›› Dokumentieren Sie auch Ihre Zwischenergebnisse.
- › Wählen Sie sinnvolle Hyperparameter für den Tokenizer und Ihre word2vec Modelle beim Training!

Extra-Aufgaben

CBOW

5 P

In Aufgabe 1.2 war nur das Skip-Gram-Modell zu implementieren. Das CBOW-Modell wird im ersten word2vec Papier ausführlicher erläutert (Mikolov, Chen et al., 2013). Im Prinzip unterscheiden sich die Modelle aber nur dahingehend, dass beim Skip-Gram-Modell ein *input* Zielwort mit den Embeddings mehrerer *output* Kontextwörter verglichen wird, während beim CBOW-Modell der Durchschnitt der *input* Kontextwörter mit dem *output* Embedding des Zielworts verglichen wird.

- › Implementieren Sie nun die CBOW-Modelle **CbowSoftMax** und **CbowNegativeSampling**.

¹Siehe: <https://docs.python.org/3/library/string.html>

Literatur

- Goldhahn, Dirk, Thomas Eckart und Uwe Quasthoff (2012). „Building Large Monolingual Dictionaries at the Leipzig Corpora Collection: From 100 to 200 Languages“. In: *Proceedings of the Eighth International Conference on Language Resources and Evaluation, LREC 2012, Istanbul, Turkey, May 23-25, 2012*. Hrsg. von Nicoletta Calzolari et al. *Die Trainings- und Testdaten für die Übungsaufgabe wurden dieser Korpusammlung entnommen*. European Language Resources Association (ELRA), S. 759–765. URL: <http://www.lrec-conf.org/proceedings/lrec2012/summaries/327.html>.
- Mikolov, Tomáš, Kai Chen et al. (2013). „Efficient Estimation of Word Representations in Vector Space“. In: *1st International Conference on Learning Representations, ICLR 2013, Scottsdale, Arizona, USA, May 2-4, 2013, Workshop Track Proceedings*. Hrsg. von Yoshua Bengio und Yann LeCun. arXiv: [1301.3781](https://arxiv.org/abs/1301.3781).
- Mikolov, Tomáš, Ilya Sutskever et al. (2013). „Distributed Representations of Words and Phrases and their Compositionality“. In: *Advances in Neural Information Processing Systems 26: 27th Annual Conference on Neural Information Processing Systems 2013. Proceedings of a meeting held December 5-8, 2013, Lake Tahoe, Nevada, United States*. Hrsg. von Christopher J. C. Burges et al., S. 3111–3119. URL: <https://proceedings.neurips.cc/paper/2013/hash/9aa42b31882ec039965f3c4923ce901b-Abstract.html>.

Revisionen

Rev. 3 mit Änderungen vom 29.05.2024, 10:20 Uhr

- › Formel (5) korrigiert.

Rev. 2 mit Änderungen vom 28.05.2024, 12:30 Uhr

- › Einige Fehler im Code-Template gefixt:
 - ›› Ein `ABC` entfernt: `class TokenizedSentenceABC(ABC) → class TokenizedSentence(NamedTuple)`.
 - ›› Type Conversion in Tests `class TestDataset.test_*_samples(): {tuple(map(int, s)) for s in ...}`.

Rev. 1 mit Änderungen vom 27.05.2024, 12:30 Uhr

- › Formeln angepasst.
- › Einige Fehler im Code-Template gefixt:
 - ›› `class DatasetABC.sentence_to_samples → class DatasetABC._sentence_to_samples` umbenannt.
 - ›› `class DatasetABC.window_size` Attribut und entsprechenden `window_size: int` Parameter hinzugefügt.
 - ›› Import: `from altair import Self → from typing import Self`
 - ›› Input zu `.forward` in `test_calculate_objective` angepasst.