

reproducing “ner and pos when nothing is capitalized”

Andreas Kuster
kustera@student.ethz

Jakub Filipek
balbok@uw

Viswa Virinchi Muppirala
virinchi@uw

Abstract

Capitalization is an important feature in many NLP tasks such as Named Entity Recognition (NER) or Part of Speech Tagging (POS). We are trying to reproduce results of paper which shows how to mitigate a significant performance drop when casing is mismatched between training and testing data. In particular we show that mixing cased and uncased dataset provides the best performance, matching the claims of the original paper. We also show that we got slightly lower performance in almost all experiments we have tried to reproduce, suggesting, that the original paper did not fully disclose all of the experimental setup.

1 Introduction

Previous works have shown that there is a significant performance drop when applying models trained on cased data to uncased data and vice-versa (Wang et al., 2006). Since capitalization is not always available due to real world constraints, there have been some methods trying to use casing prediction (called *truecasing*) to battle this trade-off.

The work we reproduce tries to battle this issue in two popular NLP tasks by mixing both cased and uncased datasets.

2 Contributions

This paper effectively shows how well work from (Mayhew et al., 2019) can be reproduced and how well it applies to a few other settings.

Original paper does show how casing issue in NLP can be effectively solved through a method which requires close to none overhead in terms of development time, and no additional overhead in runtime, especially when compared to methods such as truecasing.

It also serves as a reproducing work to show that

truecasing can be reproduced. We show that while truecasing experiment can be reproduced on the same experiment, but claims on applicability to other datasets fail to reproduce.

2.1 Hypotheses from original paper

Original paper proposes following hypotheses:

- Truecasing fits strongly to data it is presented leading to performance drops on different datasets.
- Mixing cased and uncased data provides the best performance in NER task on CoNLL 2003 dataset.
- Such technique generalizes well on noisy datasets such as Twitter data.
- Mixing cased and uncased data provides the best performance in POS task on Penn Treebank dataset.

2.2 Hypotheses addressed in this work

In addition to hypotheses tested in Section 2.1 we also tested:

- POS:
 - Mixing both cased and uncased data leads to the best performance regardless of word **embedding** used.
 - Mixing both cased and uncased data leads to the best performance regardless of word **dataset** used.
 - Using ELMo model with CRF layer in POS task to outperform the one without, regardless of data casing technique. This is based on (Huang et al., 2015).
- NER:
 - NER task is transferable to other, similar dataset, such as Groningen Meaning Bank (GMB) in this case.

2.3 Experiments

We conducted two main experiments from the paper, and tried to confirm the reproducibility claim that the original paper made.

For the truecasing we are training a simple Bidirectional LSTM with a logical layer on top, as described in both (Mayhew et al., 2019) and (Susanto et al., 2016). Since the former paper, does a great job of mentioning hyper-parameters used in the network, the hyperparameter search is not required.

Part of Speech tagging is the most in depth experiment, which lead a lot of additional hypotheses which got tested. Firstly, we had to find optimal hyper-parameter setting. The search was not too large, since a few of the parameters (such as number of layers) were known, but it still took a significant amount of time. After finding such setting we compared its results with ones reported the original paper. Then we investigated whether claims from this paper are applicable to other datasets and encodings. Lastly, we looked at the CRF layer, and what its impact on performance on the whole model is.

Named Entity Recognition **TO BE DONE**

Here every experiment is listed briefly. Every hypothesis in Section 2.2 should be listed and supported by at least one experiment, plot, table, or other type of data. Every piece of data should be listed here, with the hypothesis it supports.

3 Code

There are no public repos which do mention the project in their *READMEs*.

However, in case of truecasing, primary author of (Mayhew et al., 2019) has two repositories: *truecaser* and *python-truecaser*. The former one refers to the original implementation from (Susanto et al., 2016), while the python one is a port to python. The (probably) primary author of (Susanto et al., 2016) also has a *github repository* possibly related to truecaser we are trying to reproduce. However, to validate results and test ease of application we decided to not use either of these resources and focus on custom implementation. This also removes dependence on Andrej Karpathy's *char-rnn* all of the above mentioned repositories are forks of.

We were not able to find publicly available code for either NER or POS parts of the original paper.

Hence we needed to reimplement the code from

scratch. All our work is in *public github repository*. We tried to make all results as easily accessible as possible, which means that there is significant overlap between this report and the *README* on that repository.

Truecasing and NER were implemented using PyTorch (Paszke et al., 2019), while POS was implemented in Keras (Chollet et al., 2015), on top of TensorFlow (Abadi et al., 2015).

4 Experimental setup and results

4.1 Datasets

Datasets were a common resource about all three parts of experiments. Hence we will describe them separate here, and for each experiment specify which exact dataset was used:

- CoNLL2003 - From (Tjong Kim Sang and De Meulder, 2003). Not publicly available. Shared as part of class.
- Peen Tree Bank (PTB) - From (Marcus et al., 1993). Not publicly available. Shared as part of class. We used a loader from *nlTK* (after appending nonpublic data to it). In particular:
 - Sections 0-18 are used for training
 - Sections 19-21 are used for development, i.e. choosing hyperparameters
 - Sections 22-24 are used for test, i.e. reporting accuracy
- Twitter - From (Derczynski et al., 2016). Specifically we found *authors github repository* with data available.
- Wikipedia - From *github repository*, which we suspect is authored by primary author of (Susanto et al., 2016). This data was used for truecasing experiment.

Additionally, in case of NER and POS each of these datasets occurs in 5 different experimental flavors:

- Cased (C) - Standard dataset, as downloaded.
- Uncased (U) - Uncased dataset. It's a standard dataset, with an additional step of converting everything to lowercase.
- Cased + Uncased (C + U) - Combination of both cased and uncased flavors. Hence it is twice the size of either cased or uncased.

- Half Mixed (C + U 50) - Combination of cased and uncased flavors. However, only 50% of data is lowercase, and remaining half is as-is.
- Truecase Test (TT) - Truecased Test dataset. Training data is the same as cased one. Test is converted to uncased version, and then a trucasert is run on it.
- Truecase All (TA) - Truecased Train and Test dataset. The transformation described for test in TT is applied to both train and test.

4.2 Padding

To remove repetition in description of experiments we also want to discuss padding and how results should be understood.

All experiments described below used padding to the maximum sentence size for training. However, for results on both validation (development) and test such padding was either removed or ignored, and hence it is not affecting the results reported in this paper.

4.3 Truecasing

4.3.1 Model description

We used the exact same model as described in (Susanto et al., 2016), a 2 layer, bidirectional LSTM. Since encoding was not specified, nor any character level encoding was mentioned in class we used a PyTorch (Paszke et al., 2019) Encoder layer, which learned weights from bag-of-words to a feature-sized word representation. Note that in this model size of word embeddings is equal to hidden state size (both 300 dimensional). Then on top of these two layer a binary linear layer was applied.

Hence output of this model is 2 dimensional, one specifying that character should be upper cased (true), other that character should be left as is.

Implementation of this model took 2 hours, mostly due to transposes required by LSTM layers in PyTorch framework. It also helped that assignment 2 required us to implement LSTM model, and we could reuse parts of it.

4.3.2 Hyperparameters

We used a 2 layer, bidirectional LSTM with hidden size of 300. On top of that a linear layer, with 2 classes was applied, resulting in binary output. We used Adam (Kingma and Ba, 2014) optimizer

with default settings (learning rate of 0.001, betas of 0.9 and 0.999), and batch size of 100, as suggested in the (Susanto et al., 2016). Since these were the hyperparameters specified, we did not perform hyperparameter search.

We also used an Out-Of-Vocabulary (OOV) rate of 0.5%. Due to initial mistake we first defined out-of-vocabulary rate in following manner. When each token was drawn it had a 0.5% chance to be masked with OOV token.

However, after clarification from TA we switched OOV to be done at the initialization of tokens. In particular sorted tokens in order of increasing occurrence, and masked all least occurring tokens (such that their sum contained at least 0.5% of total amount of tokens) with OOV token. This skewed our OOV distribution towards uncommon tokens. This is a desired behavior, due to the fact that the differences between datasets will be in these rare-occurring tokens.

For wikipedia dataset, with first approach we got a dictionary of size 64, and with second approach it was 41.

4.3.3 Results

We trained these models for 30 epochs, with crossentropy loss. We recorded both training and validation loss for both, and chose the models with lowest validation loss for each of OOV settings. In both cases, such point was reached after around 7-8 epochs.

The performance of the model matched one claimed both by original paper, and original truecasing paper on Wikipedia dataset. However, there was a significant mismatch in performance on other datasets provided in the original paper. This can be seen in Table 1. A drop-off between datasets is expected, because as explained in original paper, various sources will use various casing schemes specific to the domain. For reference Table 2 contains results from investigated papers. We can see that there is not a high difference in performance between results for Wikipedia dataset, but there is a 10% difference in both CoNLL and PTB between our results and original ones. The possible reason for it might be the fact that original paper used fork of, before mentioned, [char-rnn](#). We looked at it, searching for differences between our intuition and actual implementation, but could not find any. This leaves us in conclusion that, in terms of performance, this experiment can be only partially reproduced, and it is possible that additional

Dataset	Initial OOV	Proper OOV
Wikipedia	92.65	92.71
CoNLL Train	66.03	65.32
CoNLL Test	63.49	63.28
PTB 01-18	78.53	78.73
PTB 22-24	78.47	78.69

Table 1: F1 scores of best performing model (chosen on validation set of Wikipedia dataset) for various datasets. Initial OOV stands for at-the-read technique of creating OOV tokens, while proper is the former one of two described above.

Dataset	Original
Wikipedia (Susanto)	93.19
Wikipedia	93.01
CoNLL Train	78.85
CoNLL Test	77.35
PTB 01-18	86.91
PTB 22-24	86.22

Table 2: F1 scores from reference papers. All entries but “Wikipedia (Susanto)” are from (Mayhew et al., 2019).

techniques (such as dropout, etc.) were applied during training of the model. However, in terms of between-datasets trends the paper describes, we have confirmed them, and **were able to (partially) reproduce the first hypothesis of the paper.**

4.4 Part of Speech Tagging

4.4.1 Model description

Model for POS task is not well described in the original paper, and rather references other papers for implementation details (Ma and Hovy, 2016). This left us with only a rough sketch of a model, and a hyperparameter search for most of variables.

Model used here is a single BiLSTM layer, followed by a Time Distributed Dense layer, followed by an output CRF layer. We used ELMo as the encoding for the input to BiLSTM layer.

4.4.2 Hyperparameters

As mentioned before, the task required a hyperparameter search. It was performed on standard (Cased) version of PTB dataset. Due to time constraints we opted for grid search with following settings:

One of the surprises of optimal setting is that both of dropout’s are zero. However, as seen on Fig. 1, we can see that all settings converge to

Hyperparameter	Values	Best
LSTM Hidden Units	$2^{\{0,1,2,3,5,7,9\}}$	512
LSTM Drop.	0.0, 0.2, 0.4	0.0
LSTM Recurr. Drop.	0.0, 0.2, 0.4	0.0
learning rate	0.1, 0.001	0.001

Table 3: Hyperparamter Search for POS experiment. Values investigated are in the middle column, while optimal combination of hyperparameters is in the right one.

roughly the same performance on validation set, with 0.0 having the fastest convergence.

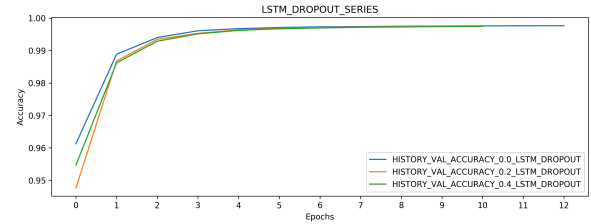


Figure 1: Performance of various settings for different dropout values in LSTM layer.

Sizes of Time Distributed Dense layer, and CRF layer are the same, and equal to number of classes given dataset. Additionally size of ELMo encoding (and hence size of input to BiLSTM layer) is 1024.

Due to time constraints we also used early stopping. Training was stopped either after 40 epochs or if validation accuracy improvement was smaller than 0.001 over 4 epochs.

4.4.3 Results

We tested the optimal model on 5 different setups, as described in list of flavors in Sec. 4.1. The dataset used in this case was PTB, with splits as described before.

Exp.	Test (C)	Test (U)	Avg
C	97.30	88.29	92.78
U	96.51	96.51	96.51
C + U	97.51	96.59	97.05
C + U 50	97.12	96.19	96.66
TT	95.04	95.04	95.04
TA	96.61	96.61	96.61

Table 4: Accuracies for the POS task, in our setup. Averages of two best performing flavors are emboldened for readability.

We can see that we agree with the hypothesis that mixing cased and uncased provides the best

Exp.	Our	Original
C	92.78	93.26
U	96.51	97.45
C + U	97.05	97.57
C + U 50	96.66	97.61
TT	95.04	95.21
TA	96.61	97.38

Table 5: Comparison of average accuracies for the POS task. Results on the left are using our code, while ones on the right are reported in the original paper. For both sources 2 best results are emboldened.

performance on the Penn Tree Bank as can be seen in Table 4. This confirms the hypothesis.

However, as in case of Truecasing, we conclude that in terms of absolute performance we did get different results as reference paper. This can be seen in Table 5, where there is about 1% difference in accuracy between our implementation and the original one. In addition to this the gain due to mixing cased and uncased data is sometimes smaller than the difference between implementations, such as in case of Uncased flavor, where reference has higher accuracy than our C + U flavor.

Overall, **we confirm the hypothesis**, but with a grain of salt due to lack of transparency in implementation in the original paper.

4.5 Named Entity Recognition

4.5.1 Model description

See syllabus. If you had to implement the model, record how long each part took.

4.5.2 Hyperparameters

See syllabus. Describe what you can.

4.5.3 Results

Each experiment should have: a clear explanation of how it was run, the high-level takeaway, a pointer to the hypothesis it supports, and it should say whether it reproduces the experiments in the original paper or not.

5 Experiments beyond the original paper

5.1 Datasets

In addition to the datasets used in original paper, we additionally used:

- Brown - From [Brown University](#). We used a loader from `nlTK`.

- CoNLL2000 - From ([Tjong Kim Sang and Buchholz, 2000](#)). We used a loader from `nlTK`.
- PTB Reduced (PTB R) - Same as Penn Tree Bank described in Section 4.1, but:
 - Train consists of sections 0-4
 - Validation consists of sections 5-6
 - Test consists of sections 7-8
- **ADD Groningen Meaning Bank (GMB) dataset**

5.2 Part of Speech Tagging

The aim of the additional experiments is to find out if the hypothesis from the paper is more generally applicable. We run the same experiments on LSTM models with different word embeddings, with or without the CRF layer, and on different datasets.

5.2.1 Hyperparameters

Due to results in Section 4.4.2 the same model was used as in the original experiment.

One important note for encodings is that they change number of inputs to BiLSTM layer, thus slightly modifying its size. In particular both GloVe and word2vec have size of 300, which is much smaller than 1024 in ELMo.

5.2.2 Results

Since there were 3 experiments roughly separate experiments let us present them separately.

Exp.	word2vec	GloVe	ELMo
C	83.71	91.01	92.78
U	80.97	94.67	96.51
C + U	86.15	96.36	97.05
C + U 50	85.01	95.35	96.66
TT	85.74	93.82	95.04
TA	86.64	95.20	96.61

Table 6: Average accuracies for different encodings. Note that ELMo encodings are the exact copy of results from 4. Similarly to other tables 2 best scores are emboldened for readability.

Table 6 shows results for various encodings. We see that we consistently mixed dataset appears to be in the top 2 results. We can also see that encoding has a strong effect on absolute performance of the model, especially in case of word2vec which has a 10% drop in accuracy relative to GloVe and ELMo. This is an expected behavior as ELMo is

known to outperform both word2vec and GloVe in many cases.

There is an interesting notion that both of Truecasing flavors tend to perform better than Half Mixed one. However, due to significant drop in word2vec with respect to other encodings we believe that it rather an artifact of the encoding itself rather than a theme. Additionally, performance of all three scenarios is rather close further pointing to encoding specific problem.

Overall, we believe that the first additional hypothesis is **supported by the above results**.

Exp.	PTB R	Brown	CoNLL2000
C	92.36	89.50	92.86
U	95.48	92.91	96.83
C + U	96.22	96.47	99.23
C + U 50	95.71	93.92	97.16
TT	94.62	92.11	95.40
TA	95.35	92.62	96.79

Table 7: Average accuracies for different datasets.

Table 7 shows results for different dataset. Here results are clear, mixing cased and uncased datasets provides the best result regardless of a dataset. In fact for Brown CoNLL2000 gaps between this technique and others is even larger than on the original PTB dataset. Hence, these results **strongly support our second additional hypothesis**.

Exp.	No CRF	CRF
C	92.88	92.78
U	96.52	96.51
C + U	97.02	97.05
C + U 50	96.83	96.66
TT	94.78	95.04
TA	96.56	96.61

Table 8: Average accuracies for the original model with CRF and without CRF.

Table 8 compares performance of the original model with or without the last layer. Considering that most of results are within a few hundredths between two models, with TT being 0.26% away from the original implementation, we conclude that CRF can be eliminated from the model without major differences in performance. This is very beneficial on the implementation and runtime side, since CRF is a non-standard library, making model are complicated, and does not have multi-gpu sup-

port, causing it to run much longer. These results, **strongly support our third additional hypothesis**, and we recommend potential future users of BiLSTM to not include CRF layer (at least for POS task).

6 Computational requirements

Due to relative lack of overlap between three sections each of us used different computational resources:

- Truecasing:
 - GPU: NVIDIA RTX 2080 Super
 - CPU: AMD Ryzen 7 3700x
 - Runtime: 10h (Original Paper + Debugging)
- POS:
 - GPU: 2x NVIDIA V100
 - CPU: AMD EPYC 7501
 - Runtime: 15h (Original Paper + Additional Experiments + Debugging)
- NER:
 - GPU: **TO BE FILLED**
 - CPU: **TO BE FILLED**
 - Runtime: **TO BE FILLED**

This is required for every report. Include every item listed in the syllabus, plus any other relevant statistics (e.g. if you have multiple model sizes, report info for each). Include the information listed in the syllabus, including the **total** number of GPU hours used for all experiments, and the number of GPU hours for **each** experiment.

7 Discussion and recommendations

Conclude the report, and summarize which hypotheses from the original paper were reproducible. Include suggestions for future researchers – what was hard? What worked easily?

References

Martín Abadi, Ashish Agarwal, Paul Barham, Eugene Brevdo, Zhifeng Chen, Craig Citro, Greg S. Corrado, Andy Davis, Jeffrey Dean, Matthieu Devin, Sanjay Ghemawat, Ian Goodfellow, Andrew Harp, Geoffrey Irving, Michael Isard, Yangqing Jia, Rafal Jozefowicz, Lukasz Kaiser, Manjunath Kudlur, Josh Levenberg, Dandelion Mané, Rajat Monga, Sherry

- Moore, Derek Murray, Chris Olah, Mike Schuster, Jonathon Shlens, Benoit Steiner, Ilya Sutskever, Kunal Talwar, Paul Tucker, Vincent Vanhoucke, Vijay Vasudevan, Fernanda Viégas, Oriol Vinyals, Pete Warden, Martin Wattenberg, Martin Wicke, Yuan Yu, and Xiaoqiang Zheng. 2015. [TensorFlow: Large-scale machine learning on heterogeneous systems](#). Software available from tensorflow.org.
- François Chollet et al. 2015. Keras. <https://keras.io>.
- Leon Derczynski, Kalina Bontcheva, and Ian Roberts. 2016. [Broad twitter corpus: A diverse named entity recognition resource](#). In *Proceedings of COLING 2016, the 26th International Conference on Computational Linguistics: Technical Papers*, pages 1169–1179, Osaka, Japan. The COLING 2016 Organizing Committee.
- Zhiheng Huang, Wei Xu, and Kai Yu. 2015. [Bi-directional LSTM-CRF models for sequence tagging](#). *CoRR*, abs/1508.01991.
- Diederik P. Kingma and Jimmy Ba. 2014. [Adam: A method for stochastic optimization](#).
- Xuezhe Ma and Eduard Hovy. 2016. [End-to-end sequence labeling via bi-directional LSTM-CNNs-CRF](#). In *Proceedings of the 54th Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pages 1064–1074, Berlin, Germany. Association for Computational Linguistics.
- Mitchell P. Marcus, Beatrice Santorini, and Mary Ann Marcinkiewicz. 1993. [Building a large annotated corpus of English: The Penn Treebank](#). *Computational Linguistics*, 19(2):313–330.
- Stephen Mayhew, Tatiana Tsygankova, and Dan Roth. 2019. [ner and pos when nothing is capitalized](#). *CoRR*, abs/1903.11222.
- Adam Paszke, Sam Gross, Francisco Massa, Adam Lerer, James Bradbury, Gregory Chanan, Trevor Killeen, Zeming Lin, Natalia Gimelshein, Luca Antiga, Alban Desmaison, Andreas Kopf, Edward Yang, Zachary DeVito, Martin Raison, Alykhan Tejani, Sasank Chilamkurthy, Benoit Steiner, Lu Fang, Junjie Bai, and Soumith Chintala. 2019. [Pytorch: An imperative style, high-performance deep learning library](#). In *Advances in Neural Information Processing Systems 32*, pages 8024–8035. Curran Associates, Inc.
- Raymond Hendy Susanto, Hai Leong Chieu, and Wei Lu. 2016. [Learning to capitalize with character-level recurrent neural networks: An empirical study](#). In *Proceedings of the 2016 Conference on Empirical Methods in Natural Language Processing*, pages 2090–2095, Austin, Texas. Association for Computational Linguistics.
- Erik F. Tjong Kim Sang and Sabine Buchholz. 2000. [Introduction to the conll-2000 shared task: Chunking](#). In *Proceedings of the 2nd Workshop on Learning Language in Logic and the 4th Conference on Computational Natural Language Learning - Volume 7*, ConLL 00, page 127132, USA. Association for Computational Linguistics.
- Erik F. Tjong Kim Sang and Fien De Meulder. 2003. [Introduction to the CoNLL-2003 shared task: Language-independent named entity recognition](#). In *Proceedings of the Seventh Conference on Natural Language Learning at HLT-NAACL 2003*, pages 142–147.
- Wei Wang, Kevin Knight, and Daniel Marcu. 2006. [Capitalizing machine translation](#). In *Proceedings of the Human Language Technology Conference of the NAACL, Main Conference*, pages 1–8, New York City, USA. Association for Computational Linguistics.