Explaining A Deep Learning Based Model for Spelling Correction of "de/da" Clitics in Turkish through Understanding Their Dependence to Various Regions of A Sentence



Advisors: Suzan Üsküdarlı – Onur Güngör

Department of Computer Engineering, Boğaziçi University



Introduction and Motivation

- A deep learning-based spellchecker is built to detect the spelling errors of Turkish "de/da" clitics with a SOTA score of 88.26% F1-Score on this task.
- The spellchecker is treated as a black-box which is used to spellcheck perturbed sentences to provide an explanation in terms of the most significant features related to the prediction.
- A Game with a Purpose (GWAP) type web application called "-de/-da Takıntısı" was developed to gather information from Turkish speakers regarding how they predict the correct spelling of perturbed sentences containing "de/da". These results are being analyzed and compared with the findings of the model explanation.

Perturbation Based Explanation Approach

- There has been efforts to build explanation techniques for black-box NLP models [1]. In this study, a modelagnostic explanation approach is proposed to interpret sequence labeling models.
- Sentences are perturbed with a certain strategy and their effect on the model is observed.

Tenis topu evde kaldı. — Tenis evde kaldı.

Perturbation example 1: Remove the two previous word.

Arabayı ben de sürdüm. — Arabayı ben de düm.

Perturbation example 2: Remove a common prefix in the following word.

Algorithm 1 Perturbation based explanation algorithm > Set of sentences which contain "de/da ⊳ Set of perturbation functions Perturbations> "-de/-da" Spellchecker function \triangleright set of $\overline{\Delta P}$ for each perturbation 4: Output: Perturbation_Impact 5: for $pr \in Perturbations$ do $w_d \leftarrow get_deda_word(s)$ $\langle w_d, p_0 \rangle \leftarrow spellchecker(s, w_d)$ $s_{pr} \leftarrow pr(s)$ $\langle w_d, p_{pr} \rangle \leftarrow spellchecker(s_{pr}, w_d)$ $\Delta P_{pr} = p_{pr} - p_0$ $\Delta P_{pr_{\tau}} \leftarrow \Delta P_{pr_{\tau}} + \Delta P_{pr}$ $Perturbation_Impact[i].append(\langle pr, \overline{\Delta_{Ppr}} \rangle)$ 7: return Perturbation_Impact

Perturbation based explanation algorithm \mathbf{w}_d is a word that contains a "de/da" and **spellchecker** is the model trained to check the spelling of "de/da" words and it returns the labeling probability of a given word in a sentence.

The average Δ_{P_i} 's for each perturbation strategy indicates the effect of each perturbation to the prediction of the model. If the Δ_{P_i} for a particular strategy is relatively high, then this means that that perturbation is important for determining the "de/da" usage in a sentence.

Results of the Perturbation

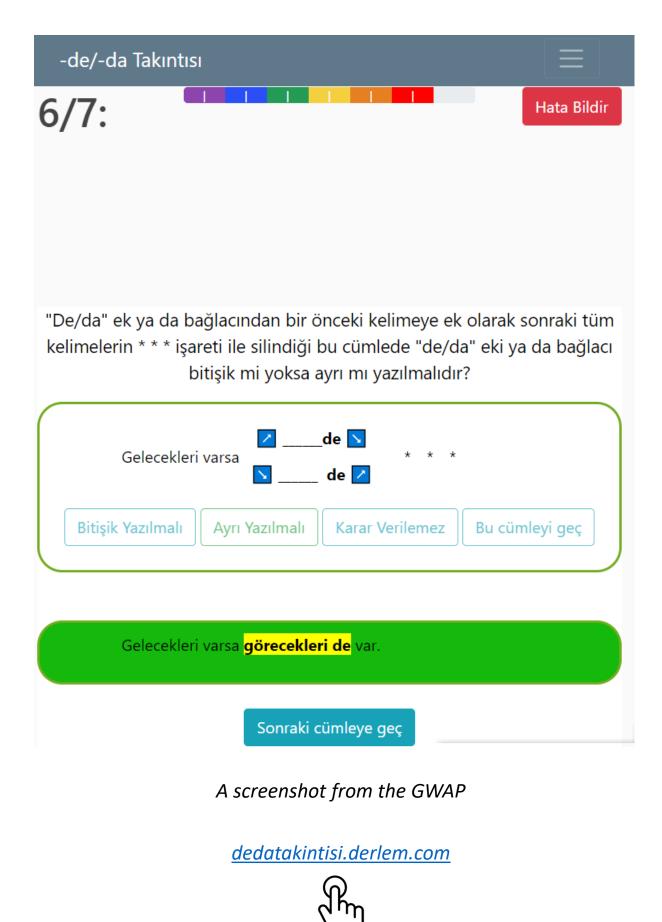
Perturbation Description	Δ_{P}
Remove the first letter of the previous word	0.0021
Remove the last letter of the previous word	0.0033
Remove the previous word	0.33
Remove the two previous word	0.32
Remove the three previous word	0.32
Remove the following word	0.0002
Remove the two following word	0.0026
Remove the three following word	0.0009
Remove a common prefix in the previous word	0.0016
Remove a common suffix in the previous word	0.03
Remove a common prefix in the following word	0.013
Remove a common suffix in the following word	0.012

Results of different perturbation techniques

These results show that the most important elements in a sentence determining the decisions of our model are the words which comes before the "de/da" clitic.

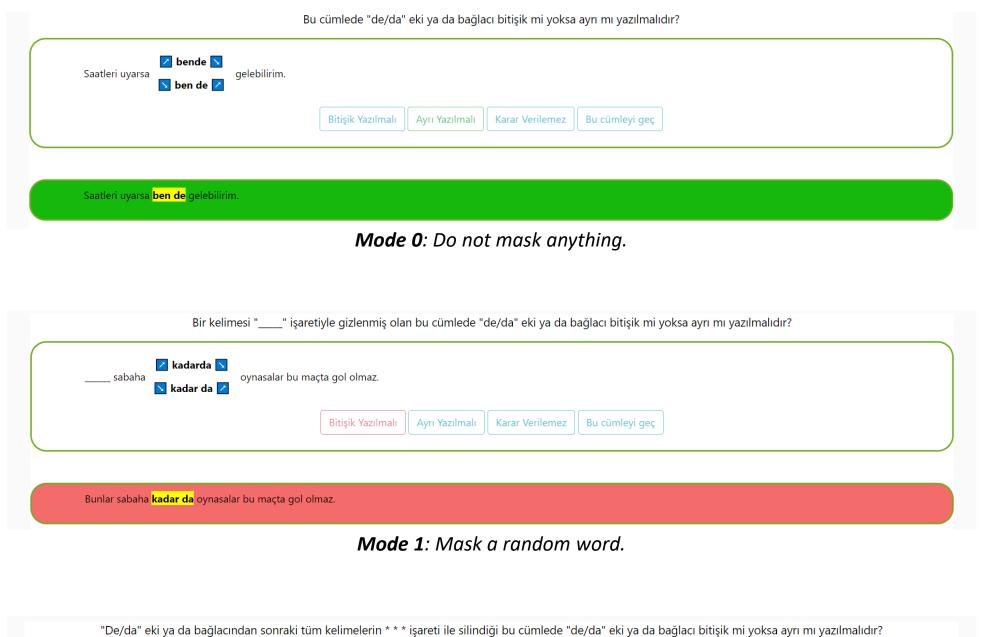
Human Based Computation Game

- We developed a GWAP to validate the findings of the perturbation-based explanation approach.
- In this game some words are masked in a sentence and the user is asked to predict the correct spelling of the "de/da" clitic.



Game Modes

There are 7 modes in the game, each of which measures the significance of a particular region of a sentence.





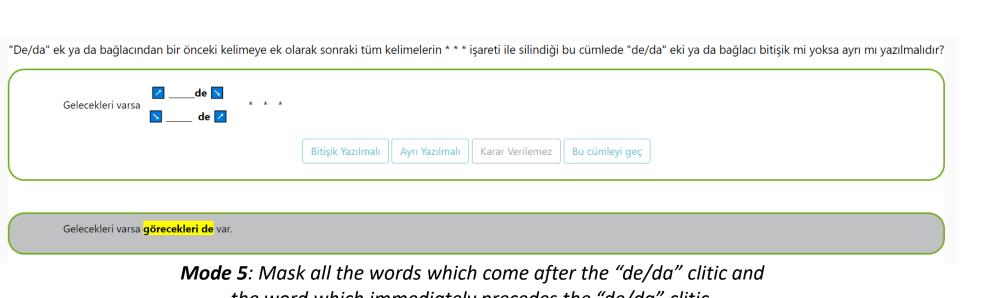
Mode 2: Mask all the words which come after the "de/da" clitic

× * * ■ ülke	senin gibi çok ins	na ihtiyaç var.			
		Bitişik Yazılmalı Ayrı Ya	zılmalı Karar Verilemez	Bu cümleyi geç	

excluding the one which immediately precedes the "de/da" clitic.

Bitişik Yazılmalı Ayrı Yazılmalı Karar Verilemez Bu cümleyi geç

Mode 4: Mask all the words which come before the "de/da" clitic



the word which immediately precedes the "de/da" clitic

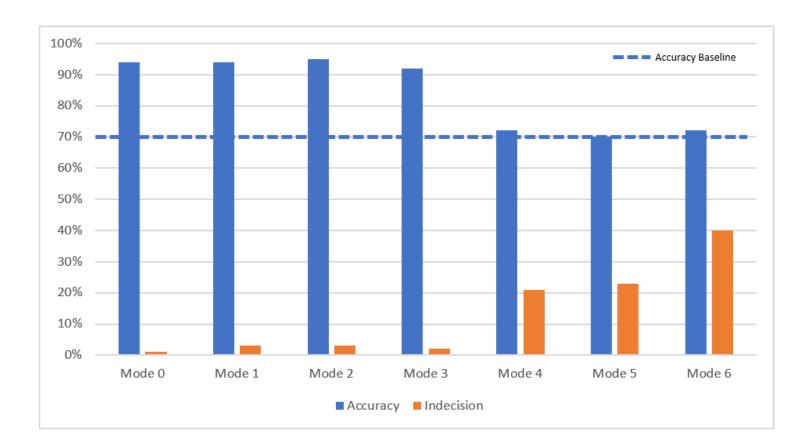


The "de/da" Spellchecker



Preliminary Results of the GWAP Experiment

447 people participated to the experiment and 13536 questions were solved in total.



- According to results of Mode 0, our participants had 94% success in determining of correct spelling of "de/da" clitic.
- Accuracy of the Mode 2 and 4 are 95% and 72% respectively. This shows that words which come **before** the "de/da" clitic are more important than the words which come after.
- Accuracy of the Mode 2 and 5 are 95% and 70% respectively. This shows that the word which immediately precedes the "de/da" clitic has a crucial role in the correct spelling.

Discussion and Future Work

- The findings of the explanation approach can be used to improve the spellchecker.
- This model has the potential to be integrated into a general-purpose spellchecker.

Acknowledgements

We would like to thank Suzan Üsküdarlı and Onur Güngör for their exceptional support throughout this year. We are also grateful to everyone who contributed to the experiment.

References

[1]: Ribeiro, M. T., Singh, S., & Guestrin, C. (2016, August). "Why should i trust you?" Explaining the predictions of any classifier. In Proceedings of the 22nd ACM SIGKDD international conference on knowledge discovery and data mining (pp. 1135-1144).