

Using Conditional Probabilities and Vowel Collocations of a Corpus of Orthographic Representations to Evaluate Nonce Words

Özkan Kılıç

Department of Cognitive Science, Informatics Institute, Middle East Technical University, Ankara, Turkey
e-mail: okilic@ii.metu.edu.tr

Abstract. Nonce words are widely used in linguistic research to evaluate areas such as the acquisition of vowel harmony and consonant voicing, naturalness judgment of loanwords, and children's acquisition of morphemes. Researchers usually create lists of nonce words intuitively by considering the phonotactic features of the target languages. In this study, a corpus of Turkish orthographic representations is used to propose a measure for the nonce word appropriateness for linearly concatenative languages. The conditional probabilities of orthographic co-occurrences and pairwise vowel collocations within the same word boundaries are used to evaluate a list of nonce words in terms of whether they would be rejected, moderately accepted or fully accepted. A group of 50 Turkish native speakers were asked to evaluate the same list of nonce words. Both the method and the participants displayed similar results.

Keywords: Nonce words, Orthographic representations, Conditional probabilities.

1 Introduction

Nonce words are frequently employed in linguistic studies to evaluate areas such as well-formedness [1], morphological productivity [2] and development [3], judgment of semantic similarity [4], and vowel harmony [5]. Nonce words are also used to understand the process of adopting loan words. The majority of loaned words undergo certain phonetic changes to more resemble the lexical entries of the language into which they will be adopted [6]. For example, *television* in Turkish becomes *televizyon* /*televizjon*/ because /*jon*/ is more frequent than /*zm*/ in Turkish¹. Similarly, *train* is adopted as *tren* /*tren*/ because, similar to diphthongs, vowel-to-vowel co-occurrences are not usually allowed in Turkish non-compound words. This phenomenon shows that the speakers of a language are aware of the possible sound frequencies and collocations of their native languages, and they can make judgements on the naturalness of loan

¹In the METU-Turkish Corpus, there are 181 occurrences with the segment /*zm*/ of which only 30 are at the terminating word boundaries. On the other hand, there are 5,945 occurrences with the segment /*jon*/ of which 3,190 are at the terminating word boundaries, excluding the word *televizyon*.

words, recently invented words and nonce words by using their knowledge of the existing Turkish lexis. Thus, the acceptability of nonce words is a logical decision based on known-word statistics.

The acceptability of nonce words can be investigated by experimental investigations through phonotactic properties or factor-based analysis [7]. In the experimental investigations, it is observed that the participants accepted or rejected nonce words according to probable combinations of sounds [1, 8]. In factor-based analysis, the acceptability of nonce words is evaluated through the co-occurrences of syllables or consonant clusters locally [9] or non-locally [10–12] or through nucleus-coda combination probabilities [13].

In this study, the acceptability of nonce words was assessed using the conditional probabilities of the bigram co-occurrences of the orthographic representations locally and the pairwise collocations of the vowels within the same word boundaries. Similar methods within the context of phonotactic modeling had been used for Finnish vowel harmony [14]. Yet in this study, the local bigram phonotactic modeling was used to evaluate Turkish nonce words. Two threshold values were set for the decision to reject, moderately accept and fully accept. The threshold values were computed according to the length of each input string. For the evaluation of the conditional and collocation probabilities, the METU-Turkish Corpus containing about two million words was employed [15]. The list of nonce words was created intuitively. The same list of nonce words evaluated by the method was also given to 50 Turkish native speakers to judge the level of acceptability of each word. The 25 male and 25 female Turkish native speakers, had an average age is 31.26 ($s = 4.11$). The results from the native speakers were very similar to the results provided by the statistical method. In this paper, brief information about Turkish language and plausibility of conditional probabilities will be given then details of the method and the results will be presented.

2 Turkish Language and Conditional Probability

Turkish has 8 vowels and 21 consonants, and it is agglutinative with a considerably complex morphology [16, 17]. While communicating, the word internal structure in Turkish is required to be segmented because Turkish morphosyntax plays a central role in semantic analysis. For example, although Turkish is considered as an SOV language, the sentences are usually in a free order. Thus, the subject and object of a verb can only be determined by the morphological markers as in (1) rather than the word order.

- | | |
|---------------------------------|-----------------------------|
| (1) <i>Köpek adam-ı ısırdı.</i> | <i>Köpeğ-i adam ısırdı.</i> |
| Dog man-Acc bit | Dog-ACC man bit |
| The dog bit the man. | The man bit the dog. |

The description of Turkish word structure depends heavily on morphophonological constraints and morphotactics. In Turkish morphotactics, the continuation of a morpheme is determined by the preceding morpheme or by the stem as in (2).

- | | |
|----------------------|------------------|
| (2) <i>ev-de-ki</i> | <i>*ev-ki-de</i> |
| house-Loc-Rel | |
| The one in the house | |

These morphotactic constraints in Turkish are captured by statistical models based on conditional probabilities [18, 19]. In addition to morphotactics, the morphophonology of Turkish needs a brief explanation because nonce words have to mimic this morphophonology.

Vowel harmony is dominantly effective in Turkish morphophonology in order to preserve the roundedness and the frontness of vowels within the same word boundaries. While a morpheme with a vowel is concatenated to a string, its vowel is modified with respect to the roundedness and frontness properties of the most recent vowel in the string as in (3).

(3) <i>ev-ler</i>	<i>oda-lar</i>	<i>bil-di</i>	<i>duy-du</i>
house - Plu	room - Plu	know - Past	hear - Past
houses	rooms	knew	heard

Another important phenomenon in Turkish morphophonology is voicing. If some of the strings terminating with the voiceless consonant, ‘*p, t, k, ç*’, are followed by the suffixes starting with vowels, then the consonants are voiced as ‘*b, d, ğ, c*’ as in (4).

(4) <i>sonuç</i>	<i>sonuc-um</i>	<i>kanat</i>	<i>kanad-ı</i>
result	result -1S.Poss	wing	wing - Acc
	my result		he wing

Consonant assimilation is also important in Turkish morphophonology. The initial consonants of some morphemes undergo an assimilation operation if they are attached to the strings terminating in the voiceless consonants, ‘*p, t, k, ç, f, s, ş, h, g*’, as in the surface forms of the Turkish past tense *-DI* in (5).

(5) <i>at-tı</i>	<i>konuş-tu</i>
throw - Past	speak - Past
threw	spoke

The final Turkish morphophonological phenomena that need to be briefly mentioned are deletion and epenthesis occurring as in (6).

(6) <i>hak</i>	<i>hakk-ım</i>	<i>isim</i>	<i>ism-im</i>
right	right - 1S.Poss	name	name - 1S.Poss
	my right		my name

The Turkish morphophonological phenomena described above occur in the co-occurrences of the orthographic representations in the concatenating positions except in vowel harmony and the deletion. This results in high conditional probabilities evaluated using the frequencies of the pairs of consecutive orthographic representations. Since the vowel harmony and deletion take place after or before the concatenation positions, their pairwise collocations within the same word boundaries are also required to be utilized in the statistical model.

The transition probability between *A* and *B* is simply based on the conditional probability statistics as in (7).

$$(7) P(B|A) = (\text{frequency of } AB) / (\text{frequency of } A)$$

Infants are reported to successfully discriminate speech segments using transitional probabilities of syllable pairs [20, 21]. Adults also make use of transitional probabilities between word classes to acquire syntactic rules [22]. Similarly, transition probabilities are dominantly used in unsupervised morphological segmentation and disambiguation [18, 19], [23–25].

Statistical approaches to linguistics support the empiricist view; and they provide an explanatory account of linguistic phenomena such as the decrease in performance errors and language variations. Considering the properties of the Turkish language, using the conditional probabilities of orthographic representations and the collocations of vowels within the same word boundaries is a plausible method to decide whether nonce words or loan words will be *rejected*, *moderately accepted* or *accepted*.

3 The Method

Let s be a string such that $s = u_1u_2 \dots u_n$, where u_i is a letter in the Turkish alphabet. The string s is unified with the empty strings σ and ε such that $s = \sigma u_1u_2 \dots u_n \varepsilon$, where σ denotes the initial word boundary and ε denotes the terminal word boundary. The overall transition probability of the string s is evaluated from the METU-Turkish Corpus using Formula 1.

$$P_t(s) = \prod_1^{n+1} P(u_i|u_{i-1}) \quad (1)$$

For example, using the Formula 1, $P(a|\sigma)$ gives the probability of the strings starting with the letter a , and $P(b|a)$ estimates the probability of the substring ab in the corpus. Now let v be a subset of the string s such that $v = u_{i,1}u_{j,2} \dots u_{k,m}$ where $u_{k,m}$ is the m^{th} vowel in the k^{th} location of the string s . The overall vowel collocations of the string s are estimated from the substring of vowels v using Formula 2.

$$P_c(v) = \prod_2^m \frac{g(v_{i-1}v_i)}{f(v_{i-1})} \quad \text{if } |v| > 1$$

$$P_c(v) = \frac{f(v_i)}{\text{CorpusSize}} \quad \text{if } |v| = 1 \quad (2)$$

In the Formula 2, the function $f(v_i)$ gives the frequency of the words that contain the vowel v_i as a substring in the corpus. The function $g(v_{i-1}v_i)$ gives the frequency of words in which the vowels v_{i-1} and v_i are collocating not necessarily in immediately consecutive positions but within the same word boundaries. The acceptability probability of the string s is calculated by $P_a(s) = P_t(s)P_c(v)$. The acceptability decision of the string s in the method is made by using the Formula 3.

$$\begin{aligned} \text{Accept} & \quad \text{if} & \quad P_a(s) \geq 10^{-(t+v)} \\ \text{Moderately accept} & \quad \text{if} & \quad 10^{-(t+v+1)} \leq P_a(s) < 10^{-(t+v)} \\ \text{Reject} & \quad \text{if} & \quad 10^{-(t+v+1)} > P_a(s) \end{aligned} \quad (3)$$

where t is the number of transitions (which is *the length of the string* + 1) and v is the number of the vowel collocations (which is *the number of the vowels* - 1) in the string. If the string s has only one vowel, then $v = 1$.

The method was applied to the list of nonce words given in the following section. The same list was also given to the 50 Turkish native speakers to evaluate the acceptability of each item. The comparison of the results from the method and the native speakers is given below.

4 Results

The nonce word *talar* is evaluated as in (8)

(8)

$$\begin{aligned} P_a(\textit{talar}) &= P_t(\sigma\textit{talar}\varepsilon) x P_c(aa) \\ &= P(t|\sigma)P(a|t)P(l|a)P(a|l)P(r|a)P(\varepsilon|r) x P_c(aa) \\ &= 7.66e - 06 x P_c(aa) = 7.66e - 06 * 4.75e - 01 = 3.63e - 06 \end{aligned}$$

Since $P_a(\textit{talar}) \geq 10^{-(6+1)}$, in which 6 conditional probability estimations and 1 vowel collocation are evaluated, the nonce word *talar* is accepted. The word list was evaluated by the 50 selected Turkish speakers. The distribution of the native speaker responses and the results of the method are given in Table 1.

For 82% of the words the Turkish native speaker's responses are in agreement with the results from the method. The method failed to simulate the responses from the participants in 18% of the results.

5 Discussions and Conclusion

The acceptability of loan words and nonce words is mainly determined by the phonological properties of the target language and the current approaches are syllable-based [7–13]. Since there are no lexical entries for nonce words, the method in this study tries to estimate the acceptability of the words using the bigram conditional probabilities and collocations of the orthographic representations within the word boundaries, which is a simplified way of inducing Turkish morphophonology.

The nonce word *ülü* was rejected by the method but accepted by the participants. A possible reason might be that the nonce word *ülü* sounds similar to an existing Turkish word *ölü* 'death'. Similarly, the responses for the nonce word *nort* were in disagreement. This nonce word has a similar pronunciation to an English word *north* and the most of the participants also knew English as a foreign language. Therefore, the participants might also make use of their foreign language knowledge to evaluate nonce words.

Although the method does not assume to utilize any property of Turkish phonology and it does not implement any phonologic filtering mechanism, it is able to mimic, in a remarkable way, a large number of the responses from the participants. Indeed, this study does not propose that acceptability is based on raw orthographical representations rather than syllables and phonemes. Instead, it underlines that simple pairwise conditional properties and vowel collocations from a corpus can give an estimation of

Table 1. The results of the method and the results of the participants (Bold text indicates a strong similarity of the results)

Nonce Words	Results of the Method	Responses of the Participants		
		Reject	Moderately Accept	Accept
<i>öğtar</i>	Reject	96%	4%	
<i>söykül</i>	Reject	96%	4%	
<i>talar</i>	Accept			100%
<i>telüti</i>	Reject	64%	28%	8%
<i>prelüs</i>	Reject	84%	14%	2%
<i>katutak</i>	ModeratelyAccept	8%	50%	42%
<i>par</i>	Accept		14%	86%
<i>öçgöş</i>	Reject	100%		
<i>jeklürt</i>	Reject	100%		
<i>böşems</i>	Reject	88%	12%	
<i>trüğat</i>	Reject	96%	4%	
<i>cakeyas</i>	Reject	92%	8%	
<i>çörottu</i>	Reject	74%	16%	10%
<i>döyyal</i>	Reject	78%	22%	
<i>efföl</i>	Reject	92%	8%	
<i>aznı</i>	Reject	32%	60%	8%
<i>fretanit</i>	Reject	64%	30%	6%
<i>erttiçe</i>	ModeratelyAccept	36%	64%	
<i>goytar</i>	Reject	38%	52%	10%
<i>hekkürük</i>	Reject	41%	47%	12%
<i>henatiya</i>	ModeratelyAccept	36%	64%	
<i>taberarul</i>	Reject	84%	16%	
<i>gövük</i>	Reject	30%	44%	26%
<i>sör</i>	ModeratelyAccept		78%	22%
<i>perolus</i>	Reject	84%	16%	
<i>kletird</i>	Reject	98%	2%	
<i>ojuçı</i>	Reject	100%		
<i>ürtanıg</i>	Reject	94%	6%	
<i>lezğaji</i>	Reject	100%		
<i>lamafi</i>	ModeratelyAccept		64%	36%
<i>nort</i>	Reject	38%	42%	20%
<i>netik</i>	Accept		18%	82%
<i>meşipir</i>	ModeratelyAccept		24%	76%
<i>oblan</i>	ModeratelyAccept		58%	42%
<i>öftik</i>	Reject	62%	34%	4%
<i>özola</i>	ModeratelyAccept	32%	60%	8%
<i>ayora</i>	Accept		72%	28%
<i>sengri</i>	ModeratelyAccept	32%	68%	
<i>sakkütan</i>	Reject	58%	34%	8%
<i>şepilt</i>	Reject	78%	22%	
<i>şür</i>	ModeratelyAccept		78%	22%
<i>puhaptı</i>	ModeratelyAccept	38%	44%	18%
<i>upapık</i>	Reject	54%	28%	18%
<i>ülü</i>	Reject	28%	52%	20%
<i>yukta</i>	ModeratelyAccept		74%	26%
<i>zerafip</i>	Reject	54%	34%	12%
<i>upgur</i>	Reject	70%	16%	14%
<i>kujmat</i>	Reject	90%	10%	
<i>lertic</i>	Reject	94%	6%	
<i>düleri</i>	Accept		64%	36%

the acceptability of a list of nonce words. This can be used by researchers that need an evaluation for the nonce words for their studies when no phonologically annotated corpus with syllables exists.

6 Limitations and development

The method needs to be tested with larger word lists. The method is successful because there is a close correspondence between phonotactics and orthotactics in Turkish. It requires improvements in terms of the morphophonological properties of target languages. The method uses exact orthographic representations. Thus, it requires an additional phonological similarity measure for the representations to increase the success rate.

The threshold values for the acceptability decisions depend on word lengths. They also need to be improved with respect to the target languages. The method also needs to be tested and adapted for the languages with ablaut or umlaut phenomena such as English and German, and the templatic languages such as Arabic and Hebrew.

References

1. Hammond, M. : Gradience, phonotactics, and the lexicon in English phonology. *Int. J. of English Studies* **4** (2004) 1–24
2. Anshen, F., Aronoff, M.: Producing morphologically complex words. *Linguistics* **26** (1988) 641–655
3. Dabrowska, E.: Low-level schemas or general rules? The role of diminutives in the acquisition of Polish case inflections. *Language Sciences* **28** (2006) 120–135
4. MacDonald, S., Ramsar, M. : Testing the distributional hypothesis: The influence of context on judgements of semantic similarity. *Proc. of the 23rd Annual Conference of the Cognitive Science Society, University of Edinburgh* (2001)
5. Pycha, A., Novak, P., Shosted, R., Shin, E.: Phonological rule-learning and its implications for a theory of vowel harmony. *Proc. of WCCFL 22, G. Garding and M. Tsujimura (Eds.)* (2003) 423–435
6. Kawahara, S: OCP is active in loanwords and nonce words: Evidence from naturalness judgment studies. *Lingua* (to appear)
7. Albright, A.: From clusters to words: Grammatical models of nonce word acceptability. Handout of talk presented at 82nd LSA, Chicago, January 3 (2008)
8. Shademan, S.: From clusters to words: Grammatical models of nonce word acceptability. *Grammar and Analogy in Phonotactic Well-formedness Judgments*. Ph. D. thesis, University of California, Los Angeles (2007)
9. Hay, J., Pierrehumbert, J., Beckman, M.: Speech perception, well-formedness and the statistics of the lexicon. In: J. Local, R. Ogden, and R. Temple (Eds.), *Phonetic Interpretation: Papers in Laboratory Phonology VI*. Cambridge: Cambridge University Press (2004)
10. Frisch, S. A., Zawaydeh, B. A.: The psychological reality of OCP-Place in Arabic. *Language* **77** (2001) 91–106
11. Koo, H., Callahan, L.: Tier-adjacency is not a necessary condition for learning phonotactic dependencies. *Language and Cognitive Processes* **77** (2011) 1–8
12. Finley, S.: Testing the limits of long-distance learning: learning beyond a three-segment window. *Cognitive Science* **36** (2012) 740–756

13. Treiman, R., Kessler, B., Knewasser, S., Tincoff, R., and Bowman, M.: English speakers' sensitivity to phonotactic patterns. In: M. B. Broe and J. Pierrehumbert (eds.), *Papers in Laboratory Phonology V: Acquisition and the Lexicon*. Cambridge: Cambridge University Press (2000) 269–282
14. Goldsmith, J., Riggle, J.: Information theoretic approaches to phonological structure: the case of Finnish vowel harmony. *Natural Language & Linguistic Theory* (to appear)
15. Say, B., Zeyrek, D., Ofłazer, K., Özge, U.: Development of a corpus and a treebank for present-day written Turkish. *Proc. of the Eleventh International Conference of Turkish Linguistics* (2002)
16. Göksel, A., Kerslake, C.: *Turkish: A Comprehensive Grammar*. Routledge: London and New York (2005)
17. Lewis, G.: *Turkish Grammar*, Second edition. Oxford: University Press (2000)
18. Kılıç, Ö., Bozşahin, C.: Semi-supervised morpheme segmentation without morphological analysis. *Pro. of the LREC 2012 Workshop on Language Resources and Technologies for Turkic Languages*, İstanbul, Turkey (2012)
19. Yatbaz, M. A., Yuret, D.: Unsupervised morphological disambiguation using statistical language models. *Pro. of the NIPS 2009 Workshop on Grammar Induction, Representation of Language and Language Learning*, Whistler, Canada (2009)
20. Aslin, R.N., Saffran, J.R., Newport, E.L.: Computation of conditional probability statistics by human infants. *Psychological Science* **9** (1998) 321–324
21. Gomez, R. L.: Variability and detection of invariant structure. *Psychological Science* **13** (2002) 431–436
22. Kaschak, M. P., Saffran, J. R.: Idiomatic syntactic constructions and language learning. *Cognitive Science* **30** (2006) 43–63
23. Creutz, M., Lagus, K.: Unsupervised models for morpheme segmentation and morphology learning. *ACM Tran. on Speech and Language Processing* **4(1)** (2007)
24. Bernhard, D.: Unsupervised morphological segmentation based on segment predictability and word segments alignment. *Proc. of 2nd Pascal Challenges Workshop* (2006) 19–24
25. Demberg, V.: A language-independent unsupervised model for morphological segmentation. *Ann. Meet. of Assoc. for Computational Linguistics* **45(1)** (2007) 920–927