arXiv:2004.00088v1 [cs.CL] 31 Mar 2020

# *A Clustering Framework for Lexical Normalization of Roman Urdu*

### ABDUL RAFAE KHAN
*Stevens Institute of Technology, Hoboken, NJ, 07030, USA*
*The Graduate Center Computer Science Department, City University of New York, 365 5th Ave, New York, NY, 10016, USA*

### ASIM KARIM
*Lahore University of Management Sciences, D.H.A, Lahore Cantt., 54792, Lahore, Pakistan*

### HASSAN SAJJAD
*Qatar Computing Research Institute, Hamad Bin Khalifa University, Doha, Qatar*

### FAISAL KAMIRAN
*Information Technology University, Arfa Software Technology Park, Ferozepur Road, Lahore, Pakistan*

### and JIA XU
*Stevens Institute of Technology, Hoboken, NJ, 07030, USA*
*The Graduate Center Computer Science Department, City University of New York, 365 5th Ave, New York, NY, 10016, USA*

## Abstract

Roman Urdu is an informal form of the Urdu language written in Roman script, which is widely used in South Asia for online textual content. It lacks standard spelling and hence poses several normalization challenges during automatic language processing. In this article, we present a feature-based clustering framework for the lexical normalization of Roman Urdu corpora, which includes a phonetic algorithm *UrduPhone*, a string matching component, a feature-based similarity function, and a clustering algorithm *Lex-Var*. UrduPhone encodes Roman Urdu strings to their pronunciation-based representations. The string matching component handles character-level variations that occur when writing Urdu using Roman script. The similarity function incorporates various phonetic-based, string-based, and contextual features of words. The Lex-Var algorithm is a variant of the k-medoids clustering algorithm that groups lexical variations of words. It contains a similarity threshold to balance the number of clusters and their maximum similarity. The framework allows feature learning and optimization in addition to the use of pre-defined features and weights. We evaluate our framework extensively on four real-world datasets and show an F-measure gain of up to 15 percent from baseline methods. We also demonstrate the superiority of UrduPhone and Lex-Var in comparison to respective alternate algorithms in our clustering framework for the lexical normalization of Roman Urdu.

## 1 Introduction

Urdu, the national language of Pakistan, and Hindi, the national language of India, jointly rank as the fourth most widely spoken language in the world (Lewis, 2009). Urdu and Hindi are closely related in morphology and phonology, but use different scripts: Urdu is written in Perso-Arabic script and Hindi is written in Devanagari script. Interestingly, for social media and short messaging service (SMS) texts, a large number of Urdu and Hindi speakers use an informal form of these languages written in Roman script, *Roman Urdu.*

Since Roman Urdu does not have standardized spellings and is mostly used in informal communication, there exist many spelling variations for a word. For example, the Urdu word زندگی [life] is written as *zindagi, zindagee, zindagy, zaindagee* and *zndagi.* The lack of standard spellings inflates the vocabulary of the language and causes sparsity problems. This results in poor performance of natural language processing (NLP) and text mining tasks, such as word segmentation (Durrani and Hussain, 2010), part of speech tagging (Sajjad and Schmid, 2009), spell checking (Naseem and Hussain, 2007), machine translation (Durrani et al., 2010), and sentiment analysis (Paltoglou and Thelwall, 2012). For example, neural machine translation models are generally trained on a limited vocabulary. Non-standard spellings would result in a large number of words unknown to the model, which would result in poor translation quality.

Our goal is to perform *lexical normalization*, which maps all spelling variations of a word to a unique form that corresponds to a single lexical entry. This reduces data sparseness and improves the performance of NLP and text mining applications.

One challenge of Roman Urdu normalization is lexical variations, which emerge through a variety of reasons such as informal writing, inconsistent phonetic mapping, and non-unified transliteration. Compared to the lexical normalization of languages with a similar script like English, the problem is more complex than writing a language informally in the original script. For example, in English, the word *thanks* can be written colloquially as *thanx* or *thx*, where the shortening of words and sounds into fewer characters is done in the same script. During Urdu to Roman Urdu conversion, two processes happen at the same time. (1) Various Urdu characters phonetically map to one or more Latin characters. (2) The Perso-Arabic script is transliterated to Roman script. Since transliteration is a non-deterministic process, it also introduces spelling variations. Fig. 1 shows an example of an Urdu word لڑکے [boys] that can be transliterated into Roman Urdu in three different ways (*larke, ladkay,* or *larkae*) depending on the user's preference. Lexical normalization of Roman Urdu aims to map transliteration variations of a word to one standard form.

Another challenge is that Roman Urdu lacks a standard lexicon or labeled corpus for text normalization to use. Lexical normalization has been addressed for standardized or resource-rich languages like English, e.g., (Jin, 2015; Han et al., 2013; Gouws et al., 2011). For such languages, the correct or the standard spelling of words is known, given the standard existence of the lexicon. Therefore, lexical normalization typically involves finding the best lexical entry for a given word
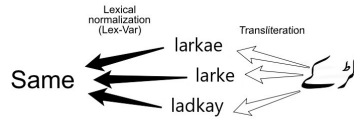
Fig. 1. The lexicon can be varied due to informal writing, non-unified definition of transliteration, phonetic mapping etc.

that does not exist in the standard lexicon. Thus, the proposed approaches aim to find the best set of standard words for a given non-standard word. On the other hand, Roman Urdu is an under-resourced language that does not have a standard lexicon. Therefore, it is not possible to distinguish between an in-lexicon and an out-of-lexicon word, and each word can potentially be a lexical variation of another. Lexical normalization of such languages is computationally more challenging than that of resource-rich languages.

Since we do not have a standard lexicon or labeled corpus for Roman Urdu lexical normalization, we cannot apply a supervised method. Therefore, we introduce an unsupervised clustering framework to capture lexical variations of words. In contrast to the English text normalization by Rangarajan Sridhar (2015); Sproat and Jaitly (2017), our approach does not require prior knowledge on the number of lexical groups or group labels (standard spellings). Our method significantly outperforms the state-of-the-art Roman Urdu lexical normalization using rule-based transliteration (Ahmed, 2009).

In this work, we give a detailed presentation of our framework (Rafae et al., 2015) with additional evaluation datasets, extended experimental evaluation, and analysis of errors. We develop an unsupervised feature-based clustering algorithm, *Lex-Var*, that discovers groups of words that are lexical variations of one another. Lex-Var ensures that each word has at least a specified minimum similarity with the cluster's centroidal word. Our proposed framework incorporates phonetic, string, and contextual features of words in a similarity function that quantifies the relatedness among words. We develop knowledge-based and machine-learned features for this purpose. The knowledge-based features include UrduPhone for phonetic encoding, an edit distance variant for string similarity, and a sequence-based matching algorithm for contextual similarity. We also evaluate various learning strategies for string and contextual similarities such as weighted edit distance and word embeddings. For phonetic information, we develop *UrduPhone*, an encoding scheme for Roman Urdu derived from Soundex. Compared to other available techniques that are limited to English sounds, UrduPhone is tailored for Roman Urdu pronunciations. For string-based similarity features, we define a function based on a combination of the longest common subsequence and edit distance metric. For contextual information, we consider top-k frequently occurring previous and next words or word groups. Finally, we evaluate our framework extensively on four Roman Urdu datasets: two group-chat SMS datasets, one Web blog dataset, and one service-feedback SMS dataset and measure performance against a manually devel-

oped database of Roman Urdu variations. Our framework gives an F-measure gain of up to 15% as compared to baseline methods.

We make the following key *contributions* in this paper:

- We present a general framework for normalizing words in an under-resourced language that allows user-defined and machine-learned features for phonetic, string, and contextual similarity.
- We propose two different clustering frameworks including a k-medoids based clustering (Lex-Var) and an agglomerative clustering (Hierarchical Lex-Var)
- We present the first detailed study of Roman Urdu normalization.
- We introduce UrduPhone for the phonetic encoding of Roman Urdu words.
- We perform an error analysis of the results, highlighting the challenges of normalizing an under-resourced and non-standard language.
- We have provided the source code for our lexical normalization framework.[1]

The remainder of this article is organized as follows. In Section 2, we present our problem statement for the lexical normalization of an under-resourced language. In Section 3, we describe our clustering framework for the lexical normalization of an under-resourced language, including UrduPhone and Lex-Var. In Section 4, we describe the evaluation criterion for the lexical normalization of Roman Urdu, describe the research experiments, and present the results and the error analysis. Section 5 discusses the related work in the lexical normalization of informal language, and Section 6 concludes the paper.

## 2  Task Definition

Roman Urdu is a transliterated form of the Urdu language written in Roman script. It does not have a standardized lexicon. That is, there is no standard spelling for words. Therefore, each word observed in a corpus can potentially be a variant of one or more of the other words appearing in the corpus. The goal of lexical normalization is to identify all spelling variations of a word in a given corpus. This challenging task involves normalizations associated with the following three issues: (1) different spellings for a given word (e.g., *kaun* and *kon* for the word [who]); (2) identically spelled words that are lexically different (e.g., *bahar* can be used for both [outside] and [spring]); and (3) spellings that match words in English (e.g., *had* [limit] for the English word *had*). The last issue arises because of code-switching between Roman Urdu and English, which is a common phenomenon in informal Urdu writing. People often write English phrases and sentences in Urdu conversations or switch language mid-sentence, e.g., *Hi everyone. Kese ha aap log?* [Hi everyone. How are you people?]. In our work, we focus on finding common spelling variations of words (issue (1)), as this is the predominant issue in the lexical normalization of Roman Urdu and do not address issues (2) or (3) explicitly.

Regarding issue (1), we note that while Urdu speakers generally transliterate

---

[1] `https://github.com/abdulrafae/normalization`

Urdu script into Roman script, they also will often move away from the transliteration in favor of a phonetically closer alternative. A commonly observed example is the replacement of one or more vowels with another set of the vowels that has a similar pronunciation (e.g., *janeaey* [to know] can also be written as *janeey*). Here, the final characters 'aey' and 'ey' give the same pronunciation. Another variation of the previous word is *janiey*. Now the character 'i' is replacing the character 'e'. In some cases, users will omit a vowel if it does not impact pronunciation, e.g., *mehnga* [expensive] becomes *mhnga* and similarly *bohut* [very] becomes *bht*. Another common example of this type of omission occurs with nasalized vowels. For example, the Roman Urdu word *kuton* [dogs] is the transliteration of the Urdu word کتوں. But often, the final nasalized Urdu character ں is omitted during conversion, and the Roman Urdu word becomes [kuto]. A similar case is found for words like *larko* [boys], *daikho* [see], *nahi* [no] with final 'n' omitted. We incorporate some of these characteristics in our encoding scheme UrduPhone (See Section 3.3.1 and Table 2 for more details on UrduPhone, its rules, and for complete steps to generate encoding).

We define the identification of lexical variations in an under-resourced language like Roman Urdu as follows: Given words $w_i$ $(i = 1, \ldots, N)$ in a corpus, find the lexical groups $\ell_j$ $(j = 1, \ldots, K)$ to which they belong. Each lexical group can contain one or more words corresponding to a single lexical entry and may represent different spelling variations of that entry in the corpus. In general, for a given corpus, the number of lexical groups $K$ is not known since no standardized lexicon is available. Therefore, we estimate it using normalization.

Clustering is expensive in the specific case of Roman Urdu normalization. Considering an efficient algorithm like k-means clustering, the computational complexity of lexical normalization is $O(NKT)$, where $T$ is the number of iterations required for clustering. By comparison, for languages like English with standardized lexicons, each out-of-vocabulary (OOV or not in the dictionary) word can be a variant of one or more in-vocabulary (IV) words. The computational complexity of lexical normalization in English (given by $O(K(N - K))$ where $K$ and $(N - K)$ are the numbers of IV and OOV words, respectively) is computationally less expensive than the lexical normalization of Roman Urdu.

## 3 Method

In this section, we describe different components of our clustering framework. Section 3.1 formalizes our clustering framework including the algorithm developed. Section 3.2 defines a similarity function used in our clustering algorithm. In Section 3.3 we describe the features used in our system.

### 3.1 Clustering Framework: Lex-Var

We develop a new clustering algorithm, Lex-Var, for discovering lexical variations in informal texts. This algorithm is a modified version of the k-medoids algorithm

---

Algorithm 1: Lex-Var

---

**Input:** $\mathcal{L}^* = \{\ell_1^*, \ell_2^*, \ldots, \ell_{K*}^*\}$ (initial clusters; see Table 4), $\mathcal{W} = \{w_1, w_2, \ldots, w_N\}$ (words), $t$ (similarity threshold)

**Output:** $\mathcal{L} = \{\ell_1, \ell_2, \ldots, \ell_K\}$ (predicted clusters)

```
 1  𝓛 = 𝓛*;
 2  repeat
        /* Find cluster centroidal word                                    */
 3      𝓒 = ∅;
 4      for ∀ ℓᵢ ∈ 𝓛 do
 5          𝓡 = ∅;
 6          for ∀ wⱼ ∈ ℓᵢ do
 7              rⱼ = 0;
 8              for ∀ wₖ ∈ ℓᵢ do
 9                  rⱼ = rⱼ + S(wⱼ, wₖ);
10              end
11              𝓡 = 𝓡 ∪ {rⱼ} ;
12          end
13          m = arg maxⱼ(rⱼ ∈ 𝓡);
14          cᵢ = wₘ;
15          𝓒 = 𝓒 ∪ {cᵢ};
16          ℓᵢ = ∅;
17      end
        /* Assign word to clusters                                         */
18      for ∀ wᵢ ∈ 𝓦 do
19          closest = null;
20          maxSim = 0;
21          for ∀ cⱼ ∈ 𝓒 do
22              if S(wᵢ, cⱼ) > t and S(wᵢ, cⱼ) > maxSim then
23                  maxSim = S(wᵢ, cⱼ);
24                  closest = cⱼ;
25              end
26          end
27          if closest != null then                  // Move word wᵢ to cluster ℓⱼ
28              ℓⱼ = ℓⱼ ∪ {wᵢ} | closest ∈ ℓⱼ;
29          else                                     // Move word wᵢ to new cluster ℓ_{|𝓛|+1}
30              ℓ_{|𝓛|+1} = {wᵢ};
31              𝓛 = 𝓛 ∪ {ℓ_{|𝓛|+1}};
32          end
33      end
34  until stop condition Satisfied;
```

---

(Han, 2005) and incorporates an assignment similarity threshold, $t > 0$, for controlling the number of clusters and their similarity spread. In particular, it ensures that all words in a group have a similarity greater than or equal to some threshold, $t$. It is important to note that the k-means algorithm cannot be used here because it requires that the means of numeric features describe the clustered objects. The standard k-medoids algorithm, on the other hand, uses the most centrally located object as a cluster's representative.

Algorithm 1 gives the pseudo-code for Lex-Var. Lex-Var takes as input words ($\mathcal{W}$) and outputs lexical groups ($\mathcal{L}$) for these words. UrduPhone segmentation of the words gives the initial clusters. Lex-Var iterates over two steps until it achieves convergence. The first step finds the centroidal word $c_i$ for cluster $\ell_i$ as the word for which the sum of similarities of all other words in the cluster is maximal. In the second step, each non-centroidal word $w_i$ is assigned to cluster $\ell_j$ if $S(w_i, c_j)$ (see Section 3.2) is maximal among all clusters and $S(w_i, c_j) > t$. If the latter condition is not satisfied (i.e., $S(w_i, c_j) \leq t$), then instead of assigning word $w_i$ to cluster $\ell_j$, it starts a new cluster. We repeat these two steps until a stop condition is satisfied (e.g., a fraction of words that change groups becomes less than a specified
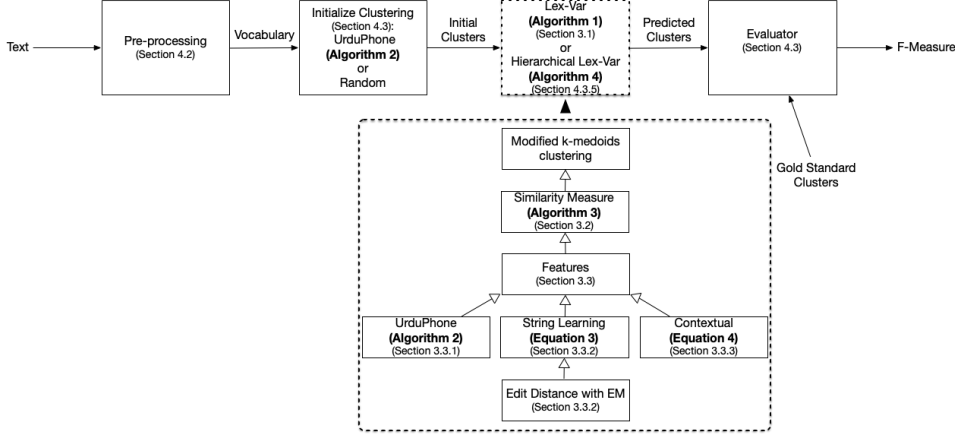
Fig. 2. Flow Diagram for Lex-Var

threshold). The computational complexity of Lex-Var is $O((n^2 + N)KT)$, where $n$ is the maximum number of words in a cluster, which is typically less than $N$.

Fig. 2 shows the details of our clustering framework. The first row of boxes shows the workflow of the system, and the area in the dotted square includes the modules used in our clustering method. The filled arrows indicate the outputs of the algorithms, and the unfilled arrows show modules that apply sub-modules.

After pre-processing the text, we normalize each word in the vocabulary. First, we initialize the clustering using random clustering or UrduPhone clusters. Then, based on the initial clusters, we apply (Hierarchical) Lex-Var algorithm to predict clusters. Finally, we compute the F- Measure based on the gold standard clusters to evaluate our prediction.

The Lex-Var algorithm applies a modified version of the k-medoids clustering, which uses a similarity measure that further consists of different features, including UrduPhone, String Learning, and Contextual feature. The edit distance is a sub-module of the string learning. We learn the substitution cost with various methods such as EM.

### 3.2 Similarity Measure

We compute the similarity between two words $w_i$ and $w_j$ using the following similarity function:

$$S(w_i, w_j) = \frac{\sum_{f=1}^{F} \alpha^{(f)} \times \sigma_{ij}^{(f)}}{\sum_{f=1}^{F} \alpha^{(f)}} \tag{1}$$

Here, $\sigma_{ij}^{(f)} \in [0, 1]$ is the similarity contribution made by feature $f$. $F$ is the total number of features. We will describe each feature in Section 3.3 in detail. $\alpha^{(f)} > 0$ is the weight of feature $f$. These weights are set to one by default and are automat-

Table 1. UrduPhone vs Soundex Encodings

| Word | Soundex Encoding | UrduPhone Encoding |
|------|------------------|--------------------|
| *mustaqbil* [future] | M_2_3_2 | M_1_2_7_9_17 |
| *mustaqil* [constant] | M_2_3_2 | M_1_2_7_17_0 |
| *khirki* [window] | K_6_2_0 | K_19_14_7_0_0 |
| *kursi* [chair] | K_6_2_0 | K_14_1_0_0_0 |
| *ronak* [brightness] | R_5_2_0 | R_11_7_0_0_0 |
| *rung* [color] | R_5_2_0 | R_11_13_0_0_0 |
| *dimaagh* [brain] | D_5_2_0 | D_12_13_19_0_0 |
| *dimaag* [brain] | D_5_2_0 | D_12_13_0_0_0 |
| *please* | P_4_2_0 | P_17_1_0_0_0 |
| *plx* | P_4_2_0 | P_17_3_0_0_0 |

---

**Algorithm 2: UrduPhone**

---

**Input:** $w = \{w_1, \cdots, w_n\}$, a word of length $n$
**Output:** $e = \{e_1, \cdots, e_6\}$, an encoding of length 6
1  $e[0] = uppercase(w[0])$;
2  $j = 1$;
3  **for** $i = 1 \rightarrow n$ **do**
       // Discard duplicates
4      **if** $i + 1 \leq n$ && $w[i] == w[i + 1]$ **then**
5         | continue;
6      **end**
       // Discard Roman Urdu vowels (a,e,i,o,u,y)
7      **if** $w[i] == 'a' \,||\, w[i] == 'e' \,||\, w[i] == 'i' \,||\, w[i] == 'o' \,||\, w[i] == 'u' \,||\, w[i] == 'y'$ **then**
8         | continue;
9      **end**
       // Encode character based on Table 2
10     $e[j] = get\_encoding(w[i])$;
11     $j + +$;
12 **end**
   // Add 0s if encoding length less than 6
13 **while** $j \leq 6$ **do**
14     $e[j] = 0$;
15     $j + +$;
16 **end**

---

ically optimized in Section 3.4 and 4.3.4. The similarity function returns a value in the interval [0, 1] with higher values signifying greater similarity.

### 3.3  Features

The similarity function in Eq. 1 is instantiated with features representing each word. In this work, we use three features: phonetic, string, and contextual, which are computed based on rules or based on learning.

### 3.3.1  UrduPhone

We propose a new phonetic encoding scheme, *UrduPhone*, tailored for Roman Urdu. Derived from Soundex (Knuth, 1973; Hall and Dowling, 1980), UrduPhone encodes

---

[2] https://en.wikipedia.org/wiki/Urdu_alphabet

Table 2. UrduPhone homophone mappings in Roman Urdu

| Characters | Urdu Alphabets | IPA² | Example |
|---|---|---|---|
| q,k | ک , ق | [q], [k] | *qainchi* [scissors], *kitab* [book] |
| c,sh,s | ث , ص , ش , س | [s], [ʃ], [s], [s] | *shadi* [wedding], *sadi* [simple] |
| z,x | ض , ظ , ذ , ز | [z], [z], [z], [z] | *zameen* [earth], *xar* [gold] |
| zh | ژ | [ʒ] | *zhalabari* [hail] |
| kh | خ | [x] | *zakhmi* [injured] |
| d | ڈ , د | [ḍ], [d] | *dahi* [yogurt], *doob* [sink] |
| t | ط , ٹ , ت | [ṭ], [t], [ṱ] | *tareef* [praise], *timatar* [tomato] |
| m | م | [m] | *maut* [death] |
| j | ج | [d͡ʒ] | *jism* [body] |
| g | گ | [g] | *gol* [circular] |
| f | ف | [f] | *fauj* [army] |
| b | ب | [b] | *bjli* [lightening] |
| p | پ | [p] | *pyaz* [onion] |
| l | ل | [l] | *lafz* [word] |
| ch | چ | [t͡ʃ] | *chehra* [face] |
| h | ھ , ہ , ح | [h, ɦ ], [h, ɦ, ø] , [ʰ, ɦ] | *haal* [present], *bahar* [spring], *phal* [fruit] |
| n | ں , ن | [n, ɲ, ɳ, ŋ], [ ˜ ] | *nazar* [sight], *larkioun* [girls] |
| r | ڑ , ر | [r], [ɽ] | *risala* [magazine], *guriya* [doll] |
| w,v | ع , و | [ʋ, uː, oː, ɔː], [aː, oː, eː, ʔ, ʕ, ø] | *waqt* [time], *vada* [promise] |
| bh | بھ | [bʰ] | *bhaag* [run] |
| ph | پھ | [pʰ] | *phool* [flower] |
| jh | جھ | [d͡ʒʰ] | *bojh* [weight], *boj* [weight] |
| th | ٹھ , تھ | [ṭʰ], [tʰ] | *thapki* [pat], *thokar* [stumble] |
| dh | ڈھ , دھ | [ḍ]ʰ, [dʰ] | *udhar* [loan], *dhool* [drum] |
| rh | ڑھ , رھ | [ rʰ], [ɽʰ] | *rhnuma* [guide], *barhna* [to grow] |
| gh | غ | [ɣ] | *ghalat* [wrong] |
| a,i,e,o,u,y | أ , ع ,و , ے , ی , أ | [aː, ʔ, ø], [j, iː, aː], [ɛː,eː], [ʋ, uː, oː, ɔː], [aː, oː, eː, ʔ, ʕ, ø], [ʔ, ø] | *aam* [mango], *ilm* [knowledge], *ullu* [owl] |

consonants by using similar sounds in Urdu and English. UrduPhone differs from Soundex in two ways:

1) UrduPhone's encoding of words contains six characters as opposed to four in Soundex. An increase in encoding length reduces the possibility of mapping semantically different words to one form. Soundex maps different words to a single encoding, which, due to the limited encoding length, can cause errors when trying to find correct lexical variations. See Table 1 for some examples of the differences. For example, *mustaqbil* [future] and *mustaqil* [constant] encode to one form, *MSTQ*, in Soundex but to two different forms using UrduPhone encoding. In a limited number of cases, UrduPhone increases ambiguity by mapping lexical variations of the same

word into different encodings, as in the case of *please* and *plx*. Since these words share a similar context, these variations will map to one cluster with the addition of contextual information. This is also shown during our experiments.

2) We introduce homophone-based groups, which are mapped differently in Soundex. There are several Urdu characters, which map to the same Roman form. For example, *samar* [reward], *sabar* [patience], and *saib* [apple], all start with different Urdu characters that have an identical Roman representation: *s*. We group together homophones such as *w, v* as in *taweez, taveez* [amulet] and *z, x* as in *lolz, lolxx* [laughter] or *zara, xara* [a bit]. One common characteristic with transliteration from Urdu to Roman script is the omission of the Roman character 'h'. For example, the same Urdu word maps to both the Roman words *samajh* & *samaj* [to understand]. This is especially true in the case of digraphs representing Urdu aspirates such as *dh, ph, th, rh, bh, jh, gh, zh, ch*, and *kh*. A problem arises when the longest common subsequence in words (if 'h' is omitted) causes overlaps such as (*khabar* [news], *kabar* [grave]) and (*gari* [car], *ghari* [watch]). Also, when *sh* comes at the end of a word, as in *khawhish, khawhis* [wish]; when 'h' is omitted, the sound is mapped to the character *s*. Similarly, if there is a transcription error, such as *dushman* [enemy] becomes *dusman*, the UrduPhone encoding is identical. Here, the omission of 'h' causes an overlap of the characters س and ش .

The second column of Table 1 shows a few examples of Soundex encodings of Roman Urdu words. In some cases, Soundex maps two semantically different words to one code, which is undesirable in the task of lexical normalization. Table 2 shows a complete list of homophone-based mapping introduced in UrduPhone, and Algorithm 2 shows the process to encode a word into an UrduPhone encoding. Then, we compute the phonetic similarity of words $w_i$ and $w_j$ using Eq. 2.

$$\sigma_{ij}^P = \begin{cases} 1 & \text{if } UrduPhone(w_i) == UrduPhone(w_j) \\ 0 & \text{otherwise} \end{cases} \qquad (2)$$

### 3.3.2 Learning String-similarity

The lexical variations of a word may have a number of overlapping sub-word units, e.g., spelling variations of *zindagi* [life] include *zindagee, zindagy, zaindagee* and *zndagi* with many overlapping sub-word units. To benefit from this overlap, we define a string similarity function as follows:

$$\sigma_{ij}^S = \frac{lcs(w_i, w_j)}{\min[len(w_i), len(w_j)] + edist(w_i, w_j)} \qquad (3)$$

Here, $lcs(w_i, w_j)$ is the length of the longest common subsequence in words $w_i$ and $w_j$, $len(w_i)$ is the length of word $w_i$, and $edist(w_i, w_j)$ is the edit distance between words $w_i$ and $w_j$.

*Edit Distance:* The edit distance allows insertion, deletion and substitution operations. We obtain the cost of edit distance operations in two ways:

**Manually Defined** – In a naive approach, we consider the cost of every operation to be equal and set them to 1. We refer to this edit distance cost as edist$_{man}$. This technique has a downside of considering all operations equally necessary, which is an erroneous assumption. For example, the substitution cost of a Roman character 'a' to 'e' should be less than the cost of 'a' to 'z' because both 'a' and 'e' have related sounds in some contexts. It is possible to use these characters alternatively when transliterating from Perso-Arabic script to Roman Script.

**Automatically Learning Edit Distance Cost** – In this approach, we automatically learn the edit distance cost from the data. Consider a list of word pairs where one word is a lexical variation of another word. One can automatically learn the character alignments between them using an EM algorithm. The inverse character alignment probability serves as the cost for the edit distance operations.

In our case, we do not have a cleaned list of word pairs to learn character alignments automatically. Instead, we try to learn these character alignments from the noisy training data. To do this, we build a list of candidate word pairs by aligning every word to every other word in the corpus as a possible lexical variation. We split the words into characters and run the word-aligner GIZA++ (Och and Ney, 2003). Here, the word-aligner considers every character as a word and every word as a sentence. We use the learned character alignments with one minus their probability as the cost for the edit distance function. We refer to this edit distance cost as edist$_{giza}$.

Since the model learns the cost from the noisy data, likely, it is not a good representative of the accurate edit distance cost that would be learned from the cleaned data. In our alternative method, we automatically refine the list of candidate pairs and learn character alignments from it. In this approach, we consider the problem of lexical variations as a **transliteration mining** problem (Sajjad et al., 2011), where, given a list of candidate word pairs, the algorithm automatically extracts word pairs that are transliterations of each other. For this purpose, we use the unsupervised transliteration mining model of Sajjad et al. (2017), who define the model[3] as a mixture of a transliteration sub-model and a non-transliteration sub-model. The transliteration sub-model generates the source and target character sequences jointly and can model the dependencies between them. The non-transliteration model consists of two monolingual character sequence models that generate source and target strings independently of each other. The parameters of the transliteration sub-model are uniformly initialized and then learned during EM training of the complete interpolated model. During the training process, the model penalizes character alignments that are less likely to be part of a transliteration pair and favors character alignments that are likely to be part of a transliteration pair.

We train the unsupervised transliteration miner on our candidate list of word pairs, similar to the GIZA++ training. Then, we learn the character alignments.

---

[3] `https://github.com/hsajjad/transliteration_mining`

We then use these character alignments with one minus their probability as the cost for the edit distance metric. We refer to this cost as edist$_{miner}$.

### 3.3.3 Context Information

We observe that non-standard variants of a standard word have similar contexts. For example, *truck* and *truk* will be used in similar contexts, which might be very different from *cat*. We used this idea to define a contextual similarity measure between two words. We compare the top-k frequently occurring preceding (previous) and following (next) words' features of the two words in the corpus. The previous and next word's features can be each word's ID, UrduPhone ID, or cluster/group ID (based on initial clustering of the words).

Let $a_1^i, a_2^i, \ldots, a_5^i$ and $a_1^j, a_2^j, \ldots, a_5^j$ be the features (word IDs, UrduPhone IDs, or cluster IDs) for the top-5 frequently occurring words preceding word $w_i$ and $w_j$, respectively. We use the similarity between the two words based on this context as defined by Hassan et al. (2009):

$$\sigma_{ij}^C = \frac{\sum_{k=1}^5 \rho_k}{\sum_{k=1}^5 k} \tag{4}$$

Here, $\rho_k$ is zero for any $a_k^i$ (i.e., the $k$th word in the context of $w_i$) when there exists no match in $a_*^j$ (i.e., in the context of word $w_j$). Otherwise, $\rho_k = 5 - \max[k, l] - 1$ where $a_k^i = a_l^j$ and $l$ is the highest rank (smallest integer) at which a previous match has not occurred. In other words, this measure is the normalized sum of rank-based weights for matches in the two sequences, with more importance given to those occurring in higher ranks. Note that contextual similarity can be computed even if the context sizes of the two words are different, an essential step as a word may not have 5 distinct words preceding it in the corpus.

We combine all the features using our similarity measure from Eq. 1. The code for combining a set of features is in Algorithm 3.

### 3.4 Parameter Optimization

**The feature weights** $\alpha^{(f)}$ used to measure word similarity in Eq. 1 can be tuned to optimize prediction accuracy. For example, by changing the weights in our clustering framework (see Eq. 1), we can make contextual similarity more prominent (by increasing the weight $\alpha^C$) so that words with the same UrduPhone encoding but different contexts are placed in separate clusters (see discussion in Section 4.4). But, we also test with other weight combinations and features, including using both word IDs and UrduPhone IDs to represent the top-5 most frequently occurring previous and next words (rather than just one representation as used in other experiments). We identify corresponding weights for contexts based on word IDs and UrduPhone IDs as $\alpha^{C_1}$ and $\alpha^{C_2}$, respectively. The weights for the phonetic and the string features are $\alpha^P$ and $\alpha^S$, respectively.

We also optimize $n$ variables to maximize an objective function using the Nelder-Mead method (Nelder and Mead, 1965). We use the Nelder-Mead method to max-

---

Algorithm 3: Similarity measure with weighted feature combination

---

**Input:** $w_i, w_j$ (input words), $F$ (set of features used), $\alpha$ (set of feature weights)
**Output:** $sim$ (similarity between $w_i$ & $w_j$)

1   $total\_weight = 0$;
2   $sum = 0$;
3   **if** $phonetic \in F$ **then**
4      $encoded_i = urduphone(w_i)$;
5      $encoded_j = urduphone(w_j)$;
6      $\sigma_{ij}^P = phonetic\_sim(encoded_i == encoded_j)$; // see Eq. 2
7      $sum = sum + \alpha^P \times \sigma_{ij}^P$;
8      $total\_weight = total\_weight + \alpha^P$;
9   **end**
10   **if** $string \in F$ **then**
11      $\sigma_{ij}^S = string\_sim(w_i, w_j)$; // see Eq. 3
12      $sum = sum + \alpha^S \times \sigma_{ij}^S$;
13      $total\_weight = total\_weight + \alpha^S$;
14   **end**
15   **if** $context \in F$ **then**
16      $prev_i = top5prev(w_i)$;
17      $prev_j = top5prev(w_j)$;
18      $\sigma_{ij}^{C1} = context\_sim(prev_i, prev_j)$; // see Eq. 4
19      $sum = sum + \alpha^{C1} \times \sigma_{ij}^{C1}$;
20      $total\_weight = total\_weight + \alpha^{C1}$;
21      $next_i = top5next(w_i)$;
22      $next_j = top5next(w_j)$;
23      $\sigma_{ij}^{C2} = context\_sim(next_i, next_j)$; // see Eq. 4
24      $sum = sum + \alpha^{C2} \times \sigma_{ij}^{C2}$;
25      $total\_weight = total\_weight + \alpha^{C2}$;
26   **end**
27   **if** $word2vec \in F$ **then**
28      $vec_i = word\_vector(w_i)$; // see Section 4.3.5
29      $vec_j = word\_vector(w_j)$;
30      $\sigma_{ij}^W = cosine(vec_i, vec_j)$;
31      $sum = sum + \alpha^W \times \sigma_{ij}^W$;
32      $total\_weight = total\_weight + \alpha^W$;
33   **end**
34   **if** $2skip1gram \in F$ **then**
35      $sigma_{ij}^G = 2skip1gram(w_i, w_j)$; // see Algorithm 5
36      $sum = sum + \alpha^G \times \sigma_{ij}^G$;
37      $total\_weight = total\_weight + \alpha^G$;
38   **end**
39   $sim = \frac{sum}{total\_weight}$;

---

imize the F-measure by optimizing the feature weights of our Similarity function in Eq. 1, as well as **the hyperparameter, threshold $t$,** in Line 21 of Algorithm 1. We apply 10-fold cross-validation on the SMS (small) dataset (Table 9). We will describe the results in Section 4.3.4.

## 4 Experiments

In this section, we first describe our evaluation setup and the datasets used for the experiments. Later, we present the results.

### *4.1 Evaluation Criteria*

Since the lexical normalization of Roman Urdu is equivalent to a clustering task, we can adopt measures for evaluating clustering performance. We need a gold standard database defining the correct groupings of words for evaluation. This database contains groups of words such that all words in a given group are considered lexical variations of a lexical entry. In clustering terminology, words within a cluster are more similar than words across clusters. On the other hand, we typically use the accuracy (i.e., the proportion of OOV words that correctly match IV words) to evaluate the lexical normalization of a standardized language like English. This measure is appropriate because we know the IV words and can be compared to every OOV word.

Bagga and Baldwin (1998) discussed measures for evaluating clustering performance and recommend the use of BCubed precision, recall, and F-measure. These measures possess all four desirable characteristics for clustering evaluation (homogeneity, completeness, rag bag, and cluster size vs. the number of clusters – see Vilain et al. (1995) for details). In the context of the lexical normalization of nonstandard languages, they provide the additional benefit that they are computed for each word separately and then averaged for all words. For example, if a cluster contains all variants of a word and nothing else, then it is considered homogeneous and complete, and this is reflected in its performance measures. These measures are robust in the sense that incorporating small impurities in an otherwise pure cluster impacts the measures significantly (rag bag characteristic), and the tradeoff between cluster size and the number of clusters is reflected appropriately. Other clustering evaluation measures do not possess all these characteristics and, in particular, commonly-used measures like entropy and purity are not based on individual words.

Let $\mathcal{L} = \{\ell_1, \cdots, \ell_K\}$ be the set of output clusters and $\mathcal{L}' = \{\ell'_1, \cdots, \ell'_K\}$ be the set of actual or correct clusters in the gold standard. We define correctness for word pair $w_i$ and $w_j$ as

$$C(w_i, w_j) = \begin{cases} 1 & \text{iff } (w_i \in \ell_m, w_j \in \ell_m) \text{ and } (w_i \in \ell'_n, w_j \in \ell'_n) \\ 0 & \text{otherwise} \end{cases} \quad (5)$$

In other words, $C(w_i, w_j) = 1$ when words $w_i$ and $w_j$ appear in the same cluster $(\ell_m)$ of the clustering and the same cluster $(\ell'_n)$ of the gold standard; otherwise, $C(w_i, w_j) = 0$. By definition, $C(w_i, w_i) = 1$.

The following expressions give the BCubed precision $P(w_i)$ and recall $R(w_i)$ for a word $w_i$:

$$P(w_i) = \frac{\sum_{j=1}^{N} C(w_i, w_j)}{|\ell_m|} \quad (6)$$

$$R(w_i) = \frac{\sum_{j=1}^{N} C(w_i, w_j)}{|\ell'_m|} \quad (7)$$

Here, $\ell_m$ and $\ell'_m$ identify the cluster in the clustering and gold standard, respectively, that contain word $w_i$. The summation for Eq. 6 & Eq. 7 is over all the words $j$. Finally, we define the BCubed F-measure $F(w_i)$ of word $w_i$ in the usual manner as:

$$F(w_i) = 2 \times \frac{P(w_i) \times R(w_i)}{P(w_i) + R(w_i)} \tag{8}$$

We compute the overall BCubed precision, recall, and F-measure of the clustering as the average of the respective values for each word. For example, we calculate the F-measure of the clustering as $\frac{\sum_{i=1}^{N} F(w_i)}{N}$.

### 4.2 Datasets

We utilize four datasets in our experimental evaluation. The first and second datasets, SMS (small) and SMS (large), are obtained from Chopaal, an internet-based group SMS service.[4] These two versions are from two different time periods and do not overlap. The third dataset, Citizen Feedback Monitoring Program (CFMP) dataset, is a collection of SMS messages sent by citizens as feedback on the quality of government services (e.g., healthcare facilities, property registration).[5] The fourth dataset, Web dataset, is scraped from Roman Urdu websites on news,[6] poetry,[7] SMS,[8] and blogs.[9] Unless mentioned otherwise, the SMS (small) dataset is used for the experiment. All four datasets are pre-processed with the following steps: (1) Remove single-word sentences; (2) add tags to URLs, email addresses, time, year, and numbers with at least four digits; (3) Collapse more than two repeating groups to only two (e.g., *hahahaha* to *haha*); (4) Replace punctuations with space; (5) Replace multiple spaces with single space. For the SMS (small) and SMS (large) datasets, we carry out an additional step of removing group messaging commands.

We evaluate the performance of our framework against a manually annotated database of Roman Urdu variations developed by Khan and Karim (2012). This database, which we refer to as the 'gold standard', is developed from a sample of the SMS (small) dataset. It maps each word to a unique ID representing its standard or normal form. There are 61,000 distinct variations in the database, which map onto 22,700 unique IDs. The number of variations differs widely for different unique IDs. For example, *mahabbat* [love] has over 70 variations such as *muhabaat*, *muhabbat*, and *mhbt*. The gold standard database also includes variations of English language words that are present in the dataset.

Table 3 shows statistics of the datasets in comparison with the evaluation gold

---

[4] `http://chopaal.org`
[5] `http://cfmp.punjab.gov.pk/`
[6] `http://www.shashca.com`, `http://stepforwardpak.com/`
[7] `https://hadi763.wordpress.com/`
[8] `http://www.replysms.com/`
[9] `http://roman.urdu.co/`

Table 3. Datasets and gold standard database statistics

| Dataset | SMS (small) | SMS (large) | CFMP | Web |
|---|---|---|---|---|
| Message Count | 159,158 | 1,994,136 | 183,083 | 5,423 |
| Unique words | 89,692 | 366,583 | 101,395 | 21,800 |
| Overlap with Gold Standard (OGS) | 57,699 | 51,477 | 23,112 | 12,634 |
| OGS and context information $\geq 1$ | 51,133 | 49,272 | 18,516 | 9,773 |
| UrduPhone IDs Previous Case | 11,146 | 9,738 | 4,683 | 6,171 |
| OGS and context information $\geq 5$ | 12,852 | 30,856 | 1,414 | 2,479 |
| UrduPhone IDs for Previous Case | 4,218 | 6,681 | 1,305 | 2,175 |

standard database. The "Overlap with Gold Standard" means the number of words in the vocabulary of a dataset that also appear in the gold standard lexicon Khan and Karim (2012). The table also gives the number of words that appear in the gold standard *and* have at least (1) one preceding and at least one following word (context size $\geq 1$), and (2) five distinct preceding and following words in the dataset (context size $\geq 5$). We evaluate these numbers of words for the respective datasets. The UrduPhone IDs of a dataset gives the number of distinct encodings of the evaluation words in the dataset (corresponding to the number of initial clusters).

### 4.3 Experimental Results and Analysis

We conduct different experiments to evaluate the performance of our clustering framework for lexical normalization of Roman Urdu. We test different combinations of features (UrduPhone, string, and/or, context) and different representations of contextual information (UrduPhone IDs or word IDs). We also establish two baseline methods for comparisons.

Table 4 gives the details of each experiment's setting. Exp. 1 and 2 are baselines corresponding to segmentation using UrduPhone encoding and string similarity-based clustering (with initial random clusters equal to the number of UrduPhone segments), respectively. The remaining experiments utilize different combinations of features (string, phonetic, and context) in our clustering framework. Here, for string-based features, we used manually defined edit distance rules.[10] The initial clustering in these experiments is given by segmentation via UrduPhone encoding. In Exp. 3 no contextual information is utilized, while in Exp. 4 and Exp. 5 the context is defined by the top-5 most frequently occurring previous and next words (context size $\geq$

---

[10] Section 3.3.2 presents a comparison of using automatically learned edit distance rules with manually defined rules.

5) represented by their UrduPhone IDs and word IDs, respectively. In Exp. 2 to 5, we select the similarity threshold $t$ such that the number of discovered clusters is as close as possible to the number of actual clusters in the gold standard for each dataset. This is done to make the results comparable across different settings. During our experiments, we observed that a threshold within a range of $0.25 - 0.3$ was optimal for smaller datasets, including Web & CFMP, and $0.4 - 0.45$ gave the best performance for larger datasets, including SMS (small) & SMS (large). However, we also tried to find the optimum threshold value using the Nelder-Mead method (see Table 9), which maximizes the F-Measure.

Figures 4a, 4b, 4c, and 4d show performance results on SMS (small), SMS (large), CFMP, and Web datasets, respectively. The x-axes in these figures show the experiment IDs from Table 4, while the left y-axes give the BCubed precision, recall, and F-measure, and the right y-axes describe the difference between the number of predicted and actual clusters.

The baseline experiment of segmentation via UrduPhone encoding (Exp. 1) produces a high recall and a low precision value. This is because UrduPhone tends to group more words in a single cluster, which decreases the total number of clusters and results in an overall low F-measure. The second baseline of string-based clustering (Exp. 2) gives similar values for precision and recall since the average number of clusters is closer to that of the gold standard. Although the F-measure increases over Exp. 1, string-based similarity alone does not result in sound clustering.

Combining the string and phonetic features in our clustering framework (Exp. 3) results in an increase in precision and recall values as well as a marked increase in F-measure from the baselines (e.g., there is an increase of 9% for the SMS (small) dataset, see Fig. 4a). When contextual information is added (via UrduPhone IDs in Exp. 4 and word IDs in Exp. 5), precision, recall, and F-measure values increase further. For example, for the SMS (small) dataset, the F-measure increases from 77.4% to 79.7% (2% gain) and from 77.4% to 80.3% (3% gain) from Exp. 3 to Exp. 4 and Exp. 5, respectively.
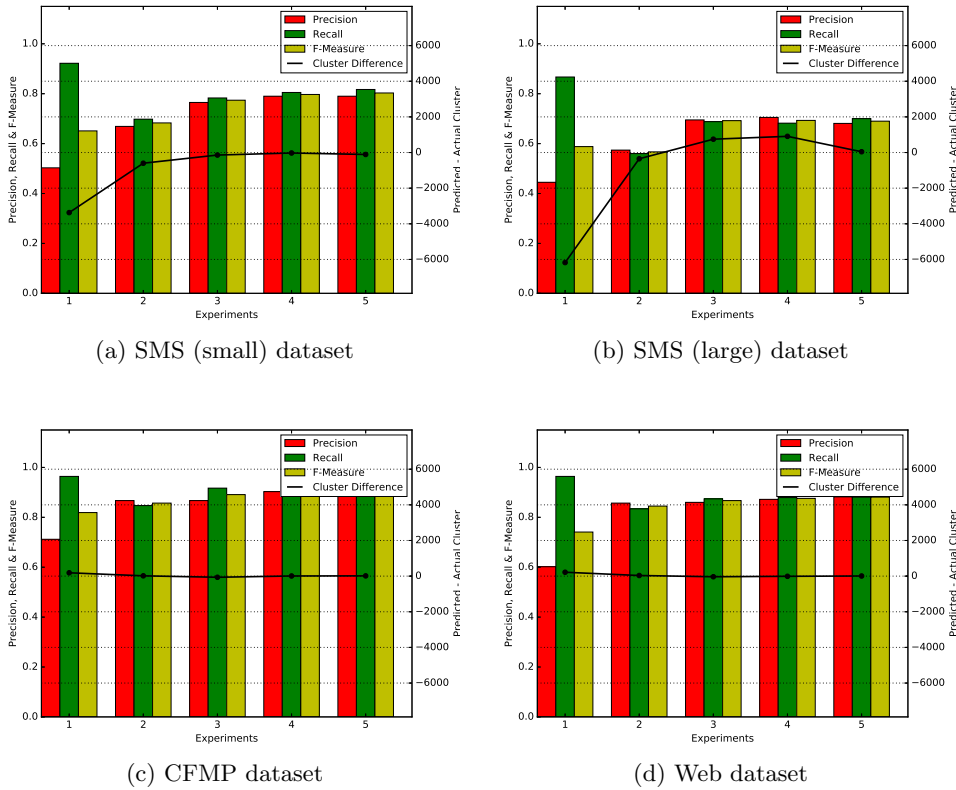
The higher performance values obtained for the CFMP and Web datasets (Fig. 4c and Fig. 4d) are due to fewer variations in these datasets, as evidenced by their fewer numbers of unique words in comparison to the SMS datasets.

Overall, our clustering framework using string, phonetic, and contextual features shows a significant F-measure gain when compared to baselines Exp. 1 and Exp. 2. We obtain the best performances when we use UrduPhone and string similarity, and when the context is defined using Word IDs (Exp. 5).

Table 4. Details of experiments' settings

| Exp. | Initial clusters | String | Phonetic | Context |
|------|------------------|--------|----------|---------|
| 1 | UrduPhone | ✗ | ✗ | – |
| 2 | Random | ✓ | ✗ | – |
| 3 | UrduPhone | ✓ | ✓ | – |
| 4 | UrduPhone | ✓ | ✓ | UrduPhone ID |
| 5 | UrduPhone | ✓ | ✓ | Word ID |

Fig. 3. Performance results for experiments in Table 4



(a) SMS (small) dataset

(b) SMS (large) dataset

(c) CFMP dataset

(d) Web dataset

### 4.3.1 UrduPhone

We compare UrduPhone with Soundex and its variants[11] for lexical normalization of Roman Urdu. All the phonetic encoding algorithms are used to group/segment

---

[11] We use Apache Commons Codec for DoubleMetaphone (`https://commons.apache.org/proper/commons-codec/apidocs/org/apache/commons/codec/language/DoubleMetaphone.html`) & NLTK-Trainer's phonetic library (`https://github.com/japerk/nltk-trainer/blob/master/nltk_trainer/featx/phonetics.py`) for the remaining

Table 5. Comparison of UrduPhone with other algorithms on the SMS (small) dataset. Single clusters are clusters with one word only. Actual clusters = 7,589

| Algorithm | Precision | Recall | F-measure | Clusters | Single Clusters |
|---|---|---|---|---|---|
| Soundex | 0.216 | 0.960 | 0.353 | 1,647 | 525 |
| Metaphone | 0.468 | 0.871 | 0.601 | 3,906 | 2,061 |
| Double Metaphone Primary Encoding | 0.295 | 0.931 | 0.448 | 2388 | 1008 |
| Double Metaphone Alternative Encoding | 0.280 | 0.927 | 0.430 | 2291 | 964 |
| Caverphone | 0.286 | 0.885 | 0.433 | 2,498 | 1,315 |
| NYSIIS | 0.584 | 0.668 | 0.623 | 6,550 | 4,376 |
| UrduPhone | 0.508 | 0.923 | **0.655** | 4,272 | 2,399 |

words based on their encoding and then evaluated against the gold standard. Table 5 shows the results of this experiment on the SMS (small) dataset.

We observe that UrduPhone outperforms Soundex, Caverphone, and Metaphone while NYSIIS's F-measure is comparable to that of UrduPhone. NYSIIS produces a large number of single-word clusters (4,376 have only one word out of 6,550 groups), which negatively impacts its recall. UrduPhone produces fewer clusters (and fewer one-word clusters), giving high recall. This property of UrduPhone is desirable for initial clustering in our clustering framework, as Lex-Var can split them but cannot collapse them.

We also test our clustering framework by replacing UrduPhone with NYSIIS as the phonetic algorithm. In Exp. 5 on the SMS (small) dataset, we find that the F-measure increases by only 5% over the NYSIIS baseline (Table 5), which is lower than the F-measure achieved with UrduPhone (Fig. 4a).

In another experiment, we analyze the effect of encoding length on the performance of the algorithm. We use the SMS (small) dataset to generate UrduPhone encodings of different sizes and cluster the words accordingly. Fig. 4 summarizes the results. We see an increase in F-measure with an increase in encoding length until length seven and eight, where we achieve similar performance.

Table 2 defines the UrduPhone rules based on well-known techniques used for phonetic encoding schemes (dropping vowels) and on common knowledge of how people write Roman Urdu. As an additional experiment, we try to learn these rules from some datasets and use them to define our encoding scheme. We call this approach UrduPhone$_{prob}$. Jiampojamarn et al. (2007) propose an alignment tool[12] based on the initial work of Ristad and Yianilos (1998). Instead of mapping each grapheme to a single phoneme, their method creates a many-to-many mapping. We use an Urdu script, and Roman Urdu transliteration parallel corpus scraped from

---

[12] https://github.com/letter-to-phoneme/m2m-aligner

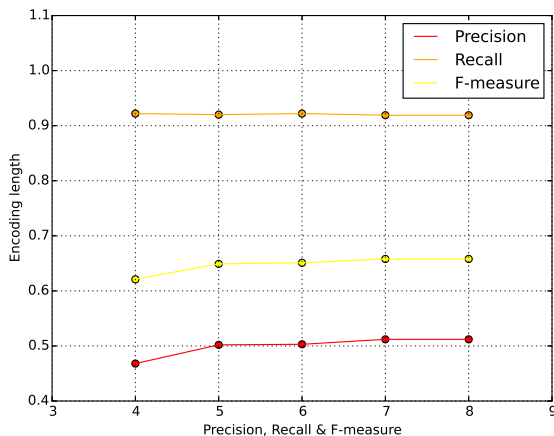Fig. 4. Effect of varying UrduPhone encoding length on SMS (small) dataset (Exp 5)



Table 6. Experiments using UrduPhone, learning rules from Urdu-Roman Urdu transliteration corpus

| Features | Precision | Recall | F-measure |
|---|---|---|---|
| UrduPhone (Exp. 1) | 0.508 | 0.923 | 0.655 |
| UrduPhone + String + Context (Exp. 5) | 0.790 | 0.817 | 0.803 |
| UrduPhone$_{prob}$ | 0.503 | 0.922 | 0.651 |
| UrduPhone$_{prob}$ + String + Context | 0.512 | 0.919 | 0.658 |

the internet.[13] Unlike the Roman Urdu words in our experiment dataset, these have more standardized spellings. We use a maximum length of two as a parameter for training the model. Our output is probabilities of Roman Urdu characters mapping to Urdu script characters or to null.

We use the maximum probability mapping rules to define our UrduPhone$_{prob}$ encodings. We experimented with using UrduPhone$_{prob}$ as the feature in our system and also in combination with other string and context features. Table 6 shows the results.

### 4.3.2 String-similarity

In section 3.3.2, using the SMS (small) dataset, we compare the performance of three methods used to calculate edit distance cost – manually defined (edist$_{man}$),

---

[13] http://www.ijunoon.com/transliteration/

Table 7. Varying edit distance cost for SMS (small) dataset. Learning character pair alignment probabilities

| String feature | Precision | Recall | F-measure |
|---|---|---|---|
| $\text{edist}_{man}$ (Exp. 5) | 0.790 | 0.817 | 0.803 |
| $\text{edist}_{giza}$ | 0.786 | 0.817 | 0.802 |
| $\text{edist}_{miner}$ | 0.794 | 0.813 | 0.803 |

automatically learned using GIZA++ ($\text{edist}_{giza}$), and automatically learned using unsupervised transliteration mining ($\text{edist}_{miner}$).[14]

For each word in our vocabulary, we found the 100 closest pairs, where closeness is defined by our similarity function as described in Eq. 1 using UrduPhone, $\text{edist}_{man}$ for the string similarity, and context of previous and next Word IDs as the feature set. We created a list of candidate word pairs by pairing every word with every other word in the cluster of 100 closest words. We take each Roman Urdu word as a sequence of Roman characters and its original Urdu script as a sequence of Urdu characters. We learn the alignment between the above two character-sequences in two different ways. First, we apply GIZA++ and learn the alignment with the Expectation-Maximization (EM) algorithm. Second, we implemented an unsupervised transliteration mining tool, details see Sajjad et al. (2017). Here, GIZA++ considers every word pair in the list of candidate pairs as a correct word pair to learn character alignments, whereas the transliteration mining tool penalizes the pairs that are less likely to be transliterations of each other during the training process. Since our list of candidate pairs is a mix of correct and incorrect pairs, the character alignments learned by the transliteration miner are likely to be better. The edit distance cost for each pair of characters can be computed from character alignments as $cost(char_i, char_j) = |1 - P(char_i, char_j)|$. Our string similarity function uses these edit distance costs instead of manually defined costs. Table 7 reports the results for both of these experiments using the SMS (small) dataset. The F-measure of the cost learned by the miner and GIZA++ is competitive with the manually defined cost. $\text{edist}_{giza}$ is affected by the noise in the data, which can be seen in its low precision compared to other methods. $\text{edist}_{miner}$ achieved the highest precision, though it has the lowest recall.

### 4.3.3 Context Size

The experiments presented in the previous section used a context of top-5 frequently occurring previous and next words. Here, we study the effect of varying context size on the performance of our clustering framework. Table 8 shows the F-measure for all experiments with two different context sizes on the SMS (small) dataset. Decreasing the minimum context list size to one increases the number of words to evaluate;

---

[14] The experiment reported in previous sections used the manually defined edit distance cost, which associates cost of 1 for each insertion, deletion, and substitution operation.

Table 8. Performance (F-measure) with two different context sizes. Details of the
experiments are given in Table 4.

| Exp. | SMS (small) | SMS (large) | CFMP | Web |
|---|---|---|---|---|
| | | Context Size = 5 | | |
| 1 | 0.651 | 0.588 | 0.852 | 0.831 |
| 2 | 0.683 | 0.567 | 0.857 | 0.845 |
| 3 | 0.774 | 0.692 | 0.891 | 0.867 |
| 4 | 0.797 | **0.693** | 0.900 | 0.876 |
| 5 | **0.803** | 0.690 | **0.917** | **0.881** |
| | | Context Size = 1 to 5 | | |
| 1 | 0.593 | 0.576 | 0.616 | 0.641 |
| 2 | 0.542 | 0.537 | 0.598 | 0.756 |
| 3 | 0.658 | 0.645 | 0.712 | 0.785 |
| 4 | 0.617 | 0.642 | 0.692 | 0.778 |
| 5 | 0.637 | 0.640 | 0.695 | 0.794 |

therefore, results are reported for all experiments with context size between 1 and
5, even though Exp. 1 to 3 do not use contextual information. Decreasing the
minimum context list size to one also explains the lower performance values for
these experiments as compared to those with context size of at least 5.

We see that context size of 1 to 5 (including words with contexts defined by at
least 1 to 5 top previous/next words) is less effective in lexical normalization and
sometimes even negatively impacts performance. For example, for the SMS (small)
and CFMP datasets, Exp. 3 (no contextual information) performs better than Exp.
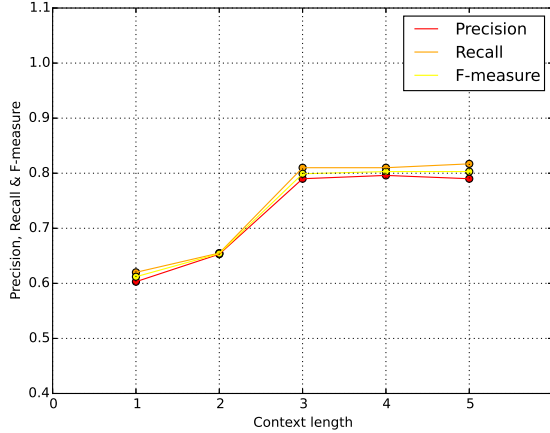4 and Exp. 5 due to the noisy nature of shorter contexts.

For further analysis, we carried out experiments where we changed the context
length from 1 to 5; an approach that differs from the previous experiments in which
we used context size = 5 & $\geq$ 1. Fig. 5 describes the results of the tests carried
out on the SMS (small) dataset. We see a significant increase in performance when
context size changes from 2 to 3. After 3, there is a slight performance increase.
The best F-measure is from a context size of 4 and 5.

### 4.3.4 Parameters: Feature Weights and Clustering Threshold

*Feature Weights* As discussed in Section 3.4, we test the impact of changing the
weights in our clustering framework (see Eq. 1). We assumed that all features have
equal weights in experiments presented in Section 4.3. Then, we change the feature
weights to emphasize different features. The increased weights caused words to
break their initial UrduPhone clusters in favor of better contextual similarity, but
the overall performance did not change. We tried several combinations, including
using both the contexts (i.e., word IDs and UrduPhone IDs).

Table 9 shows the performance of our clustering framework on the SMS (small)

Fig. 5. Effect of varying context size on SMS (small) dataset (Exp 5)



dataset with different feature weight combinations. As a comparison, we show results for Exp. 5 (context represented by word IDs only) and have the following observations with respect to F-Measure. (1) F-measure does not improve when using both word IDs and UrduPhone IDs to represent the context. (2) F-measure degrades when removing the phonetic similarity feature. (3) F-measure achieves the highest value when we set a higher weight to phonetic and contextual similarity than to string similarity.

We also use the Nelder-Mead method to maximize the F-measure by optimizing the feature weights of our similarity function in Eq. 1, as well as the threshold mentioned in line 21 of Algorithm 1 on cross-validation set (see Section 3.4). The average F-measure is slightly better than what we observed with manual selection of weights in Exp. 5 (described in Table 4).

*Clustering Threshold* We analyze the performance of Exp. 5 (best setting) for the SMS (small) dataset with varying threshold $t$ (Fig. 6). The value of $t$ controls the number of clusters smoothly, and precision increases with this number while F-measure reaches a peak when the number of predicted groups is close to that of the gold standard.
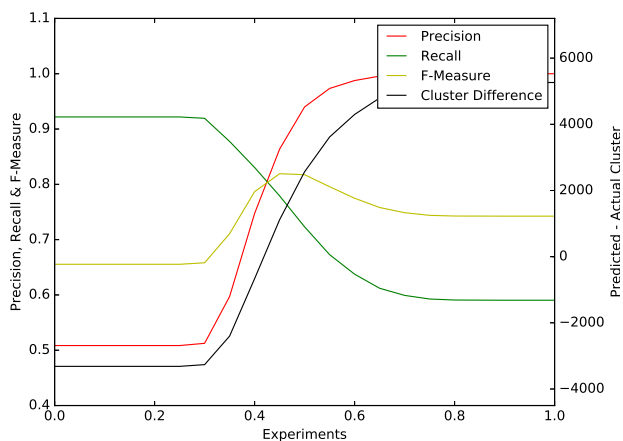
### 4.3.5 Comparison with Other Clustering Methods and Variations

In addition to our k-medoids based Lex-Var clustering method, we propose using agglomerative hierarchical clustering (Hierarchical Lex-Var) as our clustering framework for lexical normalization. To reduce the search complexity at each merge decision, we form (once) and search within the ten most similar words for each word (neighborhood). At each merge decision, we merge the two most similar words and/or groups (if either word is part of a group) in their respective neighborhoods.

Table 9. Performance with different weights for features (Exp. 5 on SMS (small) dataset). $\alpha^P$ = Weight of phonetic feature, $\alpha^S$ = Weight of string feature, $\alpha^{C_1}$ = Weight of context using Word ID , $\alpha^{C_2}$ = Weight of context using UrduPhone ID.

| Experiment | Precision | Recall | F-measure |
|---|---|---|---|
| Exp. 5 | 0.790 | 0.817 | 0.803 |
| Nelder-Mead method | 0.797 | 0.843 | 0.819 |
| $\alpha^P = 1.0, \alpha^S = 1.0, \alpha^{C_1} = 1.0, \alpha^{C_2} = 1.0$ | 0.777 | 0.814 | 0.795 |
| $\alpha^P = 1.0, \alpha^S = 1.0, \alpha^{C_1} = 2.0, \alpha^{C_2} = 0.0$ | 0.784 | 0.810 | 0.797 |
| $\alpha^P = 1.0, \alpha^S = 1.5, \alpha^{C_1} = 2.0, \alpha^{C_2} = 0.0$ | 0.801 | 0.812 | 0.807 |
| $\alpha^P = 1.0, \alpha^S = 1.0, \alpha^{C_1} = 1.5, \alpha^{C_2} = 0.0$ | 0.801 | 0.811 | 0.806 |
| $\alpha^P = 1.5, \alpha^S = 1.0, \alpha^{C_1} = 2.0, \alpha^{C_2} = 0.0$ | 0.701 | 0.819 | 0.805 |
| $\alpha^P = 1.0, \alpha^S = 1.0, \alpha^{C_1} = 0.0, \alpha^{C_2} = 2.0$ | 0.768 | 0.781 | 0.774 |
| $\alpha^P = 1.0, \alpha^S = 1.0, \alpha^{C_1} = 2.0, \alpha^{C_2} = 1.5$ | 0.736 | 0.763 | 0.749 |
| $\alpha^P = 1.0, \alpha^S = 1.0, \alpha^{C_1} = 1.5, \alpha^{C_2} = 0.5$ | 0.793 | 0.809 | 0.801 |
| $\alpha^P = 0.0, \alpha^S = 1.0, \alpha^{C_1} = 1.0, \alpha^{C_2} = 0.0$ | 0.754 | 0.758 | 0.756 |
| $\alpha^P = 0.0, \alpha^S = 1.0, \alpha^{C_1} = 1.5, \alpha^{C_2} = 0.0$ | 0.710 | 0.726 | 0.717 |
| $\alpha^P = 1.5, \alpha^S = 1.0, \alpha^{C_1} = 1.0, \alpha^{C_2} = 0.0$ | 0.802 | 0.811 | 0.807 |
| $\alpha^P = 1.5, \alpha^S = 1.0, \alpha^{C_1} = 1.5, \alpha^{C_2} = 0.0$ | 0.804 | 0.813 | 0.808 |
| $\alpha^P = 2.0, \alpha^S = 1.0, \alpha^{C_1} = 2.0, \alpha^{C_2} = 0.0$ | 0.813 | 0.809 | 0.811 |
| $\alpha^P = 2.0, \alpha^S = 1.5, \alpha^{C_1} = 2.0, \alpha^{C_2} = 0.0$ | 0.791 | 0.815 | 0.803 |

Fig. 6. Effect of varying threshold $t$ on SMS (small) dataset (Exp 5)



Algorithm 4 describes the Hierarchical Lex-Var Clustering algorithm. We tested with a neighborhood size of 10 and 100. The results are mentioned in Table 10.

Hierarchical Lex-Var, when used instead of Lex-Var, results in slightly better performance. However, it is significantly slower than Lex-Var. Even with our neighborhood-based optimization, hierarchical clustering takes hours to converge,

Table 10. Performance of Hierarchical Lex-Var on SMS (small) dataset.

| Experiment | Precision | Recall | F-measure |
|---|---|---|---|
| Exp. 5 | 0.790 | 0.817 | 0.803 |
| Nelder-Mead method | 0.797 | 0.843 | 0.819 |
| Neighborhood= 10 | 0.793 | 0.837 | 0.815 |
| Neighborhood= 100 | 0.771 | 0.849 | 0.808 |

while our Lex-Var algorithm converges in minutes when processing the SMS (small) dataset.

---

**Algorithm 4: Hierarchical Lex-Var**

---

**Input:** $\mathcal{W} = \{w_1, w_2, \ldots, w_N\}$ (words), $t$ (similarity threshold), $K$ (neighborhood size)
**Output:** $\mathcal{L} = \{\ell_1, \ell_2, \ldots, \ell_K\}$ (predicted clusters)
1 $\mathcal{L} = \mathcal{W}$;
  /* Create a similarity matrix                                              */
2 **for** $\forall\, w_i\, \in\, \mathcal{W}$ **do**
3    **for** $\forall\, w_j\, \in\, \mathcal{W} \setminus w_i$ **do**
4       $Sim\_matrix_{i,j} = S(w_i, w_j)$;
5    **end**
6 **end**
  /* Create $K$ sized neighborhoods for each word                           */
7 $\mathcal{N} = \{\}$;
8 **for** $\forall\, w_i\, \in\, \mathcal{W}$ **do**
9    $n_i = get\_max\_k(Sim\_matrix_{i,}, K)$;                       // Get $K$ most similar words to $w_i$
10    $\mathcal{N} = \mathcal{N} \cup \{n_i\}$;
11 **end**
12 **repeat**
     /* Assign word to clusters                                          */
13    **for** $\forall\, w_i\, \in\, \mathcal{W}$ **do**
14       $closest = $ null;
15       $maxSim = 0$;
16       **for** $\forall\, w_j\, \in\, \mathcal{N}_i$ **do**
17          **if** $S(w_i, n_j) > t$ and $S(w_i, n_j) > maxSim$ **then**
18             $maxSim = S(w_i, n_j)$;
19             $closest = n_j$;
20          **end**
21       **end**
22       **if** $closest\, != $ null **then**              // Move word $w_i$ to cluster $\ell_j$
23          $\ell_j = \ell_j \cup \{w_i\} \mid closest \in \ell_j$;
24       **end**
25    **end**
26 **until** *stop condition Satisfied*;

---

Additionally, we compare our clustering framework with other clustering methods as independent approaches. We also test with variations in similarity features of our clustering framework. We report the following experiments:

1. Rule-based transliteration: Each word in the vocabulary was transliterated based on the method by Ahmed (2009). The final words were mapped to an Urdu word dictionary of around 150,000 words.[15] Each Urdu word acted as a cluster label.

---

[15] `https://raw.githubusercontent.com/urduhack/urdu-words/master/words.txt`

2. Brown clustering: Brown clustering is a hierarchical clustering method for grouping words based on their contextual usage in a corpus (Brown et al., 1992). We use this as an independent approach for the lexical normalization of Roman Urdu.

3. Word2Vec clustering: Word2Vec represents words appearing in a corpus by fixed-length vectors that capture their contextual usage in the corpus (Mikolov et al., 2013). The Word2Vec model is generated using the gensim[16] python package to learn vectors for each Roman word. For learning the word vectors, we used the minimum count of 5, dimension size of 100, and 10 iterations. Words are clustered using K-Means clustering on word vectors, and we report the performance for lexical normalization of Roman Urdu.

4. 2-skip-1-grams: In our clustering framework for lexical normalization, we use the 2-skip-1-gram approach with Jaccard coefficient (Jin, 2015) to compute string similarity (rather than our string similarity function (Eq. 3)). Algorithm 5 shows the 2-skip-1-gram algorithm.

5. 2-skip-1-gram + string feature: We use both 2-skip-1-gram and our string similarity functions for computing string similarity in our clustering framework for lexical normalization.

6. 'h' omitted UrduPhone: We use a modified version of UrduPhone in our clustering framework for lexical normalization. The modified version discards aspirated characters in the encoding. For example, encoding for *mujhay* [me] becomes identical to that for *mujay* [me] to handle 'h' omission.

7. Word2Vec Vectors (50): We generate Word2Vec vecctors of size 50. We use the cosine similarity of these vectors instead of the contextual similarity described in Equation 4.

8. Word2Vec Vectors (100): We increase the size of Word2Vec vectors to 100.

9. Word2Vec Words: Word2Vec vectors are used to find the ten most similar words for each word. These neighboring words define the context of each word, and contextual similarity is computed using Eq. 4. We use our clustering framework for lexical normalization.

10. Word IDs + Word2Vec Words: We use two contextual features: top-5 frequently occurring previous/next words represented by word IDs (like in Exp. 5) and top-10 most-similar words according to Word2Vec (as above).

Table 11 summarizes the results. Experiment 1 is a rule-based lexical normalization method. Experiments 2 and 3 are independent clustering methods for lexical normalization. We also modify string features (experiments 4 and 5), phonetic features (experiment 6), and contextual features (experiments 7, 8, 9, and 10), respectively, in our clustering framework.

We can make the following observations from these experiments. (1) Rule-based transliteration performs slightly lower than our clustering method (2) Brown clustering and Word2Vec clustering are unsuitable for lexical normalization as evidenced by their poor performance. (3) Word2Vec-based context (either Word2Vec

---

[16] https://github.com/RaRe-Technologies/gensim

---

Algorithm 5: 2-skip-1-gram

---

**Input:** $w_i, w_j$ (pair of words)
**Output:** $\sigma_{ij}^{G}$ (2-skip-1-gram similarity)

1  $m = length(word_i)$;
2  $A = \phi$;
3  **for** $k \in 1, m-2$ **do**
4  $\quad$ $X = \{word_i[k]\}, \{word_i[k+2]\}$;
5  $\quad$ $A = A \cup X$;
6  **end**
7  $n = length(word_j)$;
8  $B = \phi$;
9  **for** $l \in \{1 \cdots n-2\}$ **do**
10 $\quad$ $Y = \{word_j[l]\}, \{word_j[l+2]\}$;
11 $\quad$ $B = B \cup Y$;
12 **end**
13 $\sigma_{ij}^{G} = \frac{|A \cap B|}{|A \cup B|}$;

---

Table 11. Performance of other clustering methods and variations in our framework on SMS (small) dataset.

|  | Experiment | Precision | Recall | F-measure |
|---|---|---|---|---|
| Other methods | Rule-based (Ahmed (2009)) | 0.833 | 0.765 | 0.797 |
|  | Brown clustering | 0.024 | 0.447 | 0.046 |
|  | Word2Vec clustering | 0.350 | 0.221 | 0.271 |
| Additional features | 2-skip-1-gram | 0.782 | 0.810 | 0.796 |
|  | 2-skip-1-gram + String feature | 0.791 | 0.799 | 0.795 |
|  | 'h' omitted UrduPhone | 0.796 | 0.808 | 0.802 |
|  | Word2Vec Vectors (50) | 0.782 | 0.802 | 0.792 |
|  | Word2Vec Vectors (100) | 0.795 | 0.803 | 0.799 |
|  | Word2Vec Words | 0.777 | 0.779 | 0.778 |
|  | Word IDs + Word2Vec Words | 0.780 | 0.808 | 0.793 |

vectors or similar words) and 2-skip-1-gram-based string features do not outperform our context and string features. One possible reason for the low performance of Brown clustering and Word2Vec could be the small size of the training data. These algorithms require a huge amount of data to learn.

### 4.3.6 Lexical Normalization of English Text

To test the robustness of our dataset for other languages, we experimented with an English dataset provided by Derczynski et al. (2013) and used in the W-NUT 2015 task.[17] The gold standard we used is the lexical normalization dictionary provided by the University of Melbourne.[18] The dataset has more than 160,000 messages containing 60,000 unique words. After pre-processing (the same pre-processing steps as

---

[17] `https://noisy-text.github.io/norm-shared-task.html`
[18] Available on the W-NUT 2015 website

Table 12. Performance of Lex-Var on English dataset. We used Soundex &
UrduPhone encodings as phonetic features

| Language | Phonetic Encoding | Precision | Recall | F-measure |
|---|---|---|---|---|
| Roman Urdu (Exp. 5) | UrduPhone | 0.790 | 0.817 | 0.803 |
| English | Soundex | 0.950 | 0.948 | 0.949 |
| English | UrduPhone | 0.967 | 0.961 | 0.965 |

for the Roman Urdu datasets), we get a 2,700 word-overlap with the gold standard.
For the phonetic encoding, we tested with Soundex and UrduPhone.

Table 12 summarizes the results along with the best results for the Roman Urdu
dataset from Table 4. We observe an F-measure of more than 90% with both en-
coding schemes, with UrduPhone performing better than Soundex. This difference
in performance is presumably due to the extended encoding size in UrduPhone,
which makes it possible to keep more information about the original word.

### *4.4 Error Analysis*

To gain a better understanding of our clustering framework, we analyze the output
of different experiments with examples of correct and incorrect lexical normaliza-
tion. While lexical normalization based on UrduPhone mappings (Exp. 1) is a good
starting point for finding word variations, it produces some erroneous groupings.
We summarize these groupings as follows:

1. Words that differ only in their vowels are in the same cluster:
   - *takiya* [pillow], *tikka* [grilled meat], *take*
   - *khalish* [pain], *khuloos* [sincerity]
   - *baatain* [conversations], *button*
   - *doosra* [another], *desire*
   - *separate*, *spirit*, *support*

2. Same words having different consonants map to different groups:
   - *mujhse*, *mujse* meaning [from me]
   - *kuto*, *kuton* meaning [dogs]
   - *whose*, *whoze*
   - *skool*, *school*

3. Words whose abbreviations or short forms do not have the same UrduPhone
   mapping:
   - *government*, *govt*
   - *private*, *pvt*
   - *because*, *coz*
   - *forward*, *fwd*

Exp. 4 and Exp. 5 can separate words initially clustered incorrectly (group 1.)
(e.g., *baatain* [conversations] and *button*, *spirit* and *support*) due to contextual

information and similarity differentiating the variations. Despite using phonetic variations in combination with contextual feature we see incorrect clusterings in the two experiments. We can divide these inaccuracies into several groups.

1. Words that have different UrduPhone mappings but are in fact the same. These are not clustered in the final outcome.
   - [*mujy*] and [*mujhy*] meaning [me]
   - [*oper*] and [*uper*] meaning [up]
   - [*prob*] and [*problem*]
   - [*mornin*] and [*morng*]
   - [*number*] and [*numbers*]
   - [*please*] and [*plx,plz*]

2. Words that have the same UrduPhone mapping and are lexical variants but are not clustered in the same group:
   - [*tareeka*] and [*tareka*] meaning [way]
   - [*zamaane*] and [*zamany*] meaning [times]
   - [*msg*] and [*message*]
   - [*morng*] and [*morning*]
   - [*cmplete,complet,complete*] and [*cmplt*]

3. Words that are different but have the same UrduPhone mapping and are clustered together:
   - *maalik* [owner], *malika* [queen], *malaika* [angels]
   - *nishaan* [vestige], *nishana* [target]
   - *tareka* [way], *tariq* [a common name meaning 'a night visitor']
   - *what*, *white*
   - *waiter*, *water*

A closer look at the examples reveals that some words that have the same Urdu-Phone mapping and should cluster together are found in separate groups (group 2.). This result is due to low context similarity between the words, which causes them not to group (e.g., *tareeka* and *tareka* meaning [way] have a contextual similarity of 0.23, even though they have the same UrduPhone mapping).

Another prominent issue is that words in separate clusters in UrduPhone remain separated in the output of Exp. 4 and Exp. 5 (groups 2. and 3.). This observation highlights the point that our experiments do not perform well at handling abbreviations (e.g., *prob* and *problem*), plurals (e.g., *number* and *numbers*), and some phonetic substitutes (e.g., *please* and *plx*). Our framework separates Roman Urdu words that can be written with an additional consonant (e.g., *mujy* and *mujhy* meaning [me]). It also maps words that start with a different vowel (e.g., *oper* and *uper* meaning [up]).

To tackle the issue of low contextual similarity not overcoming the difference in UrduPhone mapping, we doubled the weight assigned to the context feature. This adjustment produces almost no change in overall performance when compared to standard (Exp. 4 and Exp. 5). However, this adjustment causes more words with different UrduPhone mappings to be clustered together, usually incorrectly:

- *acha* [okay], *nahaya* [bathe], *sucha* [truthful]
- *maalom* [know], *manzor* [approve]
- *chalang* [jump], *thapar* [slap]
- *darzi* [tailor], *pathar* [stone]
- *azmaya* [to try], *sharminda* [ashamed]

Furthermore, as the same UrduPhone mappings do not restrict the clusters, this variation produces interesting combinations. The words in the groups below, although not lexical variants of each other, have strong contextual similarity and sometimes can even be replaced (for the other) in the sentence.

- *admi* [man], *larkay* [boys], *larki* [girl]
- *kufr* [to unbelieve in God], *shirk* [to associate partners with God]
- *shak* [suspicion], *yaqeen* [certainty]
- *loves*, *likes*
- *private*, *pvt*
- *cud*, *may*, *would*
- *tue*, *tuesday*, *wed*
- *blocked*, *kicked*
- *gov*, *government*

## 5 Previous Work

Normalization of informal text messages and tweets has been a research topic of interest (Sproat et al., 2001; Kaufmann and Kalita, 2010; Clark and Araki, 2011; Wei et al., 2011; Pinto et al., 2012; Ling et al., 2013; Sidarenka et al., 2013; Roy et al., 2013; Chrupała, 2014; Desai and Narvekar, 2015), with the vast majority of the work limited to English and other resource-rich languages. Our work focuses on Roman Urdu, an under-resourced language, that does not have a gold standard corpus with standard word forms. We restrict our task to finding lexical variations in informal text, a challenging problem because every word is a possible variation of every other word in the corpus. Additionally, the spelling variation problem of Roman Urdu inherits inconsistencies that occur due to the transliteration of Urdu words from Perso-Arabic script to Roman script. In our work, we model these inconsistencies separately and in combination with other features.

Researchers have used phonetic, string, and contextual knowledge to find lexical variations in informal text.[19] Pinto et al. (2012); Han et al. (2012); Zhang et al. (2015) used phonetic-based methods to find lexical variations.

Contractor et al. (2010) used string edit distance based on the longest common subsequence ratio and edit distance of Consonant Skeletons (Prochasson et al., 2007) of the IV-OOV words. Gouws et al. (2011) used a sizable English corpus to

---

[19] Spelling correction is also considered as a variant of text normalization (Damerau, 1964; Naseem, 2004; Fossati and Di Eugenio, 2007). Here, we limit ourselves to the previous work on short text normalization.

extract candidate lexical variations and re-score them based on lexical similarity. We also use lexical similarity as a feature in our clustering framework but do not have a reference to a Roman Urdu corpus with standard word forms. Jin (2015) also generated an OOV-IV list by using the Jaccard Index (Levandowsky and Winter, 1971) between $k$-skip-$n$-grams of string $s$ and standard word forms. As we do not have these in Roman Urdu, we consider every word as a possible lexical variation of every other word in the corpus. Similar to Jin (2015), we use k-skip-n-grams in our additional experiments and find that they perform slightly worse than our algorithm. Chrupała (2014) used Conditional Random Field (Lafferty, 2001) to learn the sequence of edits from labeled data.

Han et al. (2012) used word similarity and word context to enhance performance by initially extracting OOV (out-of-vocabulary) – IV (in-vocabulary) pairs using contextual similarity and then re-ranking them based on string and phonetic distances. In contrast, we define a similarity function that considers all three features together to find lexical variations of a word. Unlike previous approaches, we have a small corpus from which to extract contextually similar word pairs. Also, there is no standard Roman Urdu dictionary that can be used to annotate words as either IV or OOV. Li and Liu (2014) defined similarity measure as a combination of the longest common subsequence, term frequency, and inner product of word embeddings. We use the longest common subsequence as part of the string similarity feature. In our additional experiments, we test with a cosine similarity of word embeddings (Table 11). Li and Liu (2014) used a combination of string similarity and vector-based similarity to generate a candidate list, which was re-ranked using a character-level machine translation model (Pennell and Liu, 2011) and Jazzy Spell Checker,[20] etc. Yang and Eisenstein (2013) used an unsupervised approach that learns string edit distance, lexical, and contextual features using a log-linear model and sequential Monte Carlo approximation.

Singh et al. (2018); Costa Bertaglia and Volpe Nunes (2016) used word embeddings to find similar standard and non-standard words for text normalization. Chrupała (2014) used character-level neural text embeddings (Chrupała, 2013) as added information from unlabeled data for better performance. Rangarajan Sridhar et al. (2014) used deep neural networks to learn distributed word representations. We experimented with word embeddings as a feature in our similarity measure in the supplementary experiments Table 11.

Hany Hassan (2013) used a 5-gram language model to create a contextual similarity lattice and applied Markov random walk for lexicon generation. Their approach uses a linear combination of contextual feature and string similarity (longest common subsequence ratio and edit distance), which is very similar to our approach. However, unlike Hany Hassan (2013), we assume that every Roman Urdu word is a noisy word and thus can not separate nodes of the graph into standard and non-standard forms. Sproat and Jaitly (2017) used a recurrent neural network to normalize text. Pennell and Liu (2011); Li and Liu (2014) used a character-level

---

[20] http://jazzy.sourceforge.net/

machine translation system for the normalization task. Lusetti et al. (2018) used an encoder-decoder architecture where different levels of granularity were used for the target-side language model, e.g. characters and words. Wang and Ng (2013) used a beam-search decoder with integrated normalization operations such as missing word recovery and punctuation correction to normalize non-standard words. Our work, however, is limited to grouping the lexical variations of Roman Urdu words. However, we do not have any labeled data or parallel data available to build such a translation system. Our proposed method is robust since it learns from user data, and it groups abbreviations and their complete forms together in one cluster.

Almeida et al. (2016) used a standard English dictionary and an informal English dictionary to normalize words to their root forms. In our case, we do not use a standard dictionary as one does not exist for Roman Urdu words. Ling et al. (2013) automatically learned normalization rules using a parallel corpus of informal text. Irvine et al. (2012) used manually prepared training data to build an automatic normalization system for the Roman Urdu script. Unlike Irvine et al. (2012), we propose an unsupervised approach, which does not require labeled data. Additionally, our approach to the Roman Urdu normalization problem does not require us to have a corresponding Urdu script form for each Roman word.

*Phonetic encoding schemes* There have been several sound-based encoding schemes used in the literature to group similar sounding words together. Here, we summarized a few of the schemes in the context of lexical normalization.

The Soundex algorithm (Knuth, 1973; Hall and Dowling, 1980) encodes the first letter and the following three consonants of a word with consonants having a similar place of articulation sharing the same code. The NYSIIS method (Taft, 1970), designed by the New York Police for American names, employ more sophisticated encoding rules based on multi-character n-grams and relative vowel positioning. The Metaphone algorithm (Philips, 1990), developed in 1990 as a Soundex variant, incorporates English pronunciation rules for phonetic encoding of words. Other, more-recent variations include Caverphone (Wang, 2009) and Double Metaphone;[21] they include complex grammatical rules for phonetic encoding of words. The Double Metaphone algorithm also differs from others in that it generates up to two encodings for each word – one reflects the basic version of the word's pronunciation, and the other reflects an alternative pronunciation based on other languages. This is particularly useful when comparing foreign names with their anglicized versions. For example, the names *Catherine* and *Katrina* have a common code *KTRN*. Previous algorithms like Metaphone and Soundex do not provide such a capability.

Most of these schemes are designed for English and European languages and are not sufficiently expressive, especially for lexical normalization or when applied to another family of languages.

We propose a method to find lexical variations in Roman Urdu that uses string edit distance like Contractor et al. (2010), sound-based encoding like Pinto et al.

---

[21] http://en.wikipedia.org/wiki/Metaphone

(2012), and contextual information like Han et al. (2012) combined in a discriminative framework. In contrast to previous work, our method does not use a resource of standard word forms to find lexical variations.

## 6 Conclusion and future work

Roman Urdu is a transliterated form of the Urdu language written in Roman script, used in informal communication in social media and SMS texts. It does not have a standard lexicon, which results in an extensive use of lexical variations that hamper automatic processing. Our framework for lexical normalization of Roman Urdu is an unsupervised model meant to address this important issue. Our clustering framework incorporates customized phonetic encoding, string-based matching, and contextual similarity. We conducted an extensive evaluation of our framework on four real-world datasets. We used manually generated gold standard containing Roman Urdu lexical variations with their standard forms (Khan and Karim, 2012). We show that our framework effectively discovers lexical variations in Roman Urdu corpora with significant improvement over the baseline methods.

Our work brings us one step closer to automatically generating a normalized Roman Urdu corpus. We can cluster spelling variations of a word and then map them to the most frequent form, and can use this corpus to develop NLP applications. In the future, we would like to extrinsically evaluate our normalization procedure on several NLP tasks, such as POS tagging and machine translation.

## 7 Acknowledgement

## References

Ahmed, T. (2009). Roman to Urdu transliteration using wordlist. In *Proceedings of the Conference on Language and Technology*, Poznan, Poland.

Almeida, T. A., Silva, T. P., Santos, I., and Gómez Hidalgo, J. M. (2016). Text normalization and semantic indexing to enhance instant messaging and SMS spam filtering. *Knowledge-Based Systems*, 108(C):25–32.

Bagga, A. and Baldwin, B. (1998). Algorithms for scoring coreference chains. In *Proceedings of the 1st International Conference on Language Resources and Evaluation Workshop on Linguistics Coreference*, pages 563–566, Granada, Spain.

Brown, P. F., deSouza, P. V., Mercer, R. L., Pietra, V. J. D., and Lai, J. C. (1992). Class-based n-gram models of natural language. *Computational Linguistics*, 18(4):467–479.

Chrupała, G. (2013). Text segmentation with character-level text embeddings. In *Proceedings of the International Conference on Machine Learning: Workshop on Deep Learning for Audio, Speech and Language Processing*, Atlanta, Georgia, USA.

Chrupała, G. (2014). Normalizing tweets with edit scripts and recurrent neural embeddings. In *Proceedings of the 52nd Annual Meeting of the Association for Computational Linguistics: Short Papers*, pages 680–686, Baltimore, Maryland, USA. Association for Computational Linguistics.

Clark, E. and Araki, K. (2011). Text normalization in social media: Progress, problems and applications for a pre-processing system of casual english. *Procedia-Social and Behavioral Sciences*, 27:2–11.

Contractor, D., Faruquie, T. A., and Subramaniam, L. V. (2010). Unsupervised cleansing of noisy text. In *Proceedings of the 23rd International Conference on Computational Linguistics: Poster*, pages 189–196, Beijing, China. Association for Computational Linguistics.

Costa Bertaglia, T. F. and Volpe Nunes, M. d. G. (2016). Exploring word embeddings for unsupervised textual user-generated content normalization. In *Proceedings of the 2nd Workshop on Noisy User-generated Text*, pages 112–120, Osaka, Japan.

Damerau, F. J. (1964). A technique for computer detection and correction of spelling errors. *Communications of the Association for Computing Machinery*, 7(3).

Derczynski, L., Ritter, A., Clark, S., and Bontcheva, K. (2013). Twitter part-of-speech tagging for all: Overcoming sparse and noisy data. In *Proceedings of the International Conference on Recent Advances in Natural Language Processing*, Hissar, Bulgaria. Association for Computational Linguistics.

Desai, N. and Narvekar, M. (2015). Normalization of noisy text data. *Procedia Computer Science*, 45:127–132.

Durrani, N. and Hussain, S. (2010). Urdu word segmentation. In *Proceedings of the Annual Conference of the North American Chapter of the Association for Computational Linguistics – Human Language Technologies*, pages 528–536, Los Angeles, California, USA. Association for Computational Linguistics.

Durrani, N., Sajjad, H., Fraser, A., and Schmid, H. (2010). Hindi-to-Urdu machine translation through transliteration. In *Proceedings of the 48th Annual Meeting of the Association for Computational Linguistics*, pages 465–474, Uppsala, Sweden. Association for Computational Linguistics.

Fossati, D. and Di Eugenio, B. (2007). A mixed trigrams approach for context sensitive spell checking. In *Computational Linguistics and Intelligent Text Processing*, pages 623–633. Springer.

Gouws, S., Hovy, D., and Metzler, D. (2011). Unsupervised mining of lexical variants from noisy text. In *Proceedings of the 1st workshop on Unsupervised Learning in Natural Language Processing*, pages 82–90, Edinburgh, Scotland. Association for Computational Linguistics.

Hall, P. A. V. and Dowling, G. R. (1980). Approximate string matching. *Association for Computing Machinery Computing Surveys*, 12(4):381–402.

Han, B., Cook, P., and Baldwin, T. (2012). Automatically constructing a normalisation dictionary for microblogs. In *Proceedings of the Joint Conference on Empirical Methods in Natural Language Processing and Computational Natural*

*Language Learning*, pages 421–432, Jeju Island, Korea. Association for Computational Linguistics.

Han, B., Cook, P., and Baldwin, T. (2013). Lexical normalization for social media text. *Association for Computing Machinery Transactions on Intelligent Systems and Technology*, 4(1):5.

Han, J. (2005). *Data Mining: Concepts and Techniques*. Morgan Kaufmann Publishers Inc., San Francisco, California.

Hany Hassan, Arul Menezes, H. H. A. (2013). Social text normalization using contextual graph random walks. In *Proceedings of the 51st Annual Meeting of the Association for Computational Linguistics*, Sofia, Bulgaria.

Hassan, M. T., Junejo, K. N., and Karim, A. (2009). Learning and predicting key web navigation patterns using bayesian models. In *Computational Science and Its Applications*, pages 877–887. Springer.

Irvine, A., Weese, J., and Callison-Burch, C. (2012). Processing informal, Romanized Pakistani text messages. In *Proceedings of the 2nd Workshop on Language in Social Media*, LSM '12, pages 75–78, Montreal, Canada. Association for Computational Linguistics.

Jiampojamarn, S., Kondrak, G., and Sherif, T. (2007). Applying many-to-many alignments and hidden markov models to letter-to-phoneme conversion. In *Proceedings of the Conference of the North American Chapter of the Association for Computational Linguistics – Human Language Technologies: Main Conference*, pages 372–379, Rochester, New York, USA. Association for Computational Linguistics.

Jin, N. (2015). Ncsu-sas-ning: Candidate generation and feature engineering for supervised lexical normalization. In *Proceedings of the Workshop on Noisy User-generated Text*, pages 87–92, Beijing, China. Association for Computational Linguistics.

Kaufmann, M. and Kalita, J. (2010). Syntactic normalization of twitter messages. In *Proceedings of the International Conference on Natural Language Processing*, Kharagpur, India.

Khan, O. and Karim, A. (2012). A rule-based model for normalization of SMS text. In *Proceedings of Institute of Electrical and Electronics Engineers 24th International Conference on Tools with Artificial Intelligence*, pages 634–641, Athens, Greece.

Knuth, D. E. (1973). *The Art of Computer Programming: Volume 3, Sorting and Searching*. Addison-Wesley.

Lafferty, J. (2001). Conditional random fields: Probabilistic models for segmenting and labeling sequence data. In *Proceedings of the 18th International Conference on Machine Learning*, pages 282–289, San Francisco, California, USA.

Levandowsky, M. and Winter, D. (1971). Distance between sets. *Nature*, 234:34–35.

Lewis, M. P., editor (2009). *Ethnologue: Languages of the World*. SIL International, Dallas, Texas, USA, Sixteenth edition.

Li, C. and Liu, Y. (2014). Improving text normalization via unsupervised model and discriminative reranking. In *Proceedings of the Association for Computational Linguistics: Student Research Workshop*, pages 86–93, Baltimore, Maryland, USA. Association for Computational Linguistics.

Ling, W., Dyer, C., Black, A. W., and Trancoso, I. (2013). Paraphrasing 4 microblog normalization. In *Proceedings of the Conference on Empirical Methods in Natural Language Processing*, pages 73–84, Seattle, Washington, USA.

Lusetti, M., Ruzsics, T., Göhring, A., Samardžić, T., and Stark, E. (2018). Encoder-decoder methods for text normalization. In *Proceedings of the 5th Workshop on Natural Language Processing for Similar Languages, Varieties and Dialects*, pages 18–28. Association for Computational Linguistics.

Mikolov, T., Sutskever, I., Chen, K., Corrado, G., and Dean, J. (2013). Distributed representations of words and phrases and their compositionality. In *Proceedings of the 26th International Conference on Neural Information Processing Systems*, page 31113119, Red Hook, New York, USA. Curran Associates Inc.

Naseem, T. (2004). A hybrid approach for Urdu spell checking. Master's thesis, Master of Science (Computer Science) thesis at the National University of Computer & Emerging Sciences.

Naseem, T. and Hussain, S. (2007). A novel approach for ranking spelling error corrections for Urdu. *Language Resources and Evaluation*, 41(2):117–128.

Nelder, J. A. and Mead, R. (1965). A simplex method for function minimization. *Computer Journal*, 7:308–313.

Och, F. J. and Ney, H. (2003). A systematic comparison of various statistical alignment models. *Computational Linguistics*, 29(1):19–51.

Paltoglou, G. and Thelwall, M. (2012). Twitter, MySpace, Digg: Unsupervised sentiment analysis in social media. *Association for Computing Machinery Transactions on Intelligent Systems and Technology*, 3(4):66.

Pennell, D. and Liu, Y. (2011). A character level machine translation approach for normalization of SMS abbreviations. In *Proceedings of the 5th International Joint Conference on Natural Language Processing*, pages 974–982, Chiang Mai, Thailand.

Philips, L. (1990). Hanging on the metaphone. *Computer Language Magazine*, 7(12):39–44.

Pinto, D., Ayala, D. V., Alemán, Y., Gómez-Adorno, H., Loya, N., and Jiménez-Salazar, H. (2012). The soundex phonetic algorithm revisited for SMS text representation. In *Proceedings of the 15th International Conference on Text, Speech and Dialogue*, pages 47–55, Brno, Czech Republic.

Prochasson, E., Viard-Gaudin, C., and Morin, E. (2007). Language models for handwritten short message services. In *Proceedings of the 9th International Conference on Document Analysis and Recognition*, pages 83–87, Parana, Brazil. Institute of Electrical and Electronics Engineers Computer Society.

Rafae, A., Qayyum, A., Uddin, M. M., Karim, A., Sajjad, H., and Kamiran, F. (2015). An unsupervised method for discovering lexical variations in Roman Urdu informal text. In *Proceedings of the Conference on Empirical Methods in Natural Language Processing*, pages 823–828, Lisbon, Portugal. Association for Computational Linguistics.

Rangarajan Sridhar, V. K. (2015). Unsupervised text normalization using distributed representations of words and phrases. In *Proceedings of the 1st Workshop on Vector Space Modeling for Natural Language Processing*, pages 8–16, Denver, Colorado, USA. Association for Computational Linguistics.

Rangarajan Sridhar, V. K., Chen, J., Bangalore, S., and Shacham, R. (2014). A framework for translating SMS messages. In *Proceedings of the 25th International Conference on Computational Linguistics: Technical Papers*.

Ristad, E. S. and Yianilos, P. N. (1998). Learning string edit distance. *Institute of Electrical and Electronics Engineers Transactions on Pattern Recognition and Machine Intelligence*, 20(5):522–532.

Roy, S., Dhar, S., Bhattacharjee, S., and Das, A. (2013). A lexicon based algorithm for noisy text normalization as pre-processing for sentiment analysis. *International Journal of Research in Engineering and Technology*, 02.

Sajjad, H., Fraser, A., and Schmid, H. (2011). An algorithm for unsupervised transliteration mining with an application to word alignment. In *Proceedings of the 49th Conference of the Association for Computational Linguistics – Human Language Technologies*, Portland, Oregon, USA.

Sajjad, H. and Schmid, H. (2009). Tagging Urdu text with parts of speech: A tagger comparison. In *Proceedings of the 12th Conference of the European Chapter of the Association for Computational Linguistics*, pages 692–700, Athens, Greece.

Sajjad, H., Schmid, H., Fraser, A., and Schütze, H. (2017). Statistical models for unsupervised, semi-supervised and supervised transliteration mining. *Computational Linguistics*, 43(2).

Sidarenka, U., Scheffler, T., and Stede, M. (2013). Rule-based normalization of German twitter messages. In *Proceedings of the Gesellschaft fur Sprachtechnologie und Computerlinguistik Workshop Verarbeitung und Annotation von Sprachdaten aus Genres internetbasierter Kommunikation*, Darmstadt, Germany.

Singh, R., Choudhary, N., and Shrivastava, M. (2018). Automatic normalization of word variations in code-mixed social media text. *Computing Research Repository*, abs/1804.00804.

Sproat, R., Black, A. W., Chen, S. F., Kumar, S., Ostendorf, M., and Richards, C. (2001). Normalization of non-standard words. *Computer Speech & Language*, 15(3):287–333.

Sproat, R. and Jaitly, N. (2017). An rnn model of text normalization. In *Proceedings of Interspeech*, Stockholm, Sweden.

Taft, R. (1970). Name search techniques. *Special report. Bureau of Systems Development, New York State Identification and Intelligence System*.

Vilain, M., Burger, J., Aberdeen, J., Connolly, D., and Hirschman, L. (1995). A model-theoretic coreference scoring scheme. In *Proceedings of the 6th Message Understanding Conference*, pages 45–52, Columbia, Maryland, USA.

Wang, J., editor (2009). *Encyclopedia of Data Warehousing and Mining, Second Edition (4 Volumes)*. IGI Global.

Wang, P. and Ng, H. T. (2013). A beam-search decoder for normalization of social media text with application to machine translation. In *Proceedings of the Conference of the North American Chapter of the Association for Computational Linguistics – Human Language Technologies*, pages 471–481, Atlanta, Georgia, USA. Association for Computational Linguistics.

Wei, Z., Zhou, L., Li, B., Wong, K.-F., Gao, W., and Wong, K.-F. (2011). Exploring tweets normalization and query time sensitivity for twitter search. In *Proceedings of the 20th Text Retrieval Conference*, Gaithersburg, Maryland, USA.

Yang, Y. and Eisenstein, J. (2013). A log-linear model for unsupervised text nor-
   malization. In *Proceedings of the Conference on Empirical Methods in Natural
   Language Processing: Meeting of SIGDAT, a Special Interest Group of the Asso-
   ciation for Computational Linguistics*, pages 61–72, Seattle, Washington, USA.
Zhang, X., Song, J., He, Y., and Fu, G. (2015). Normalization of homophonic
   words in Chinese microblogs. In *Intelligent Computation in Big Data Era*, pages
   177–187. Springer.