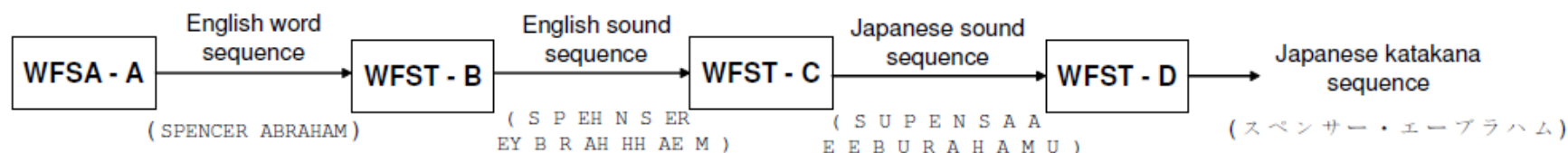


Non-Stochastic Learning of Cross-Language Transliteration Rules From Small Dataset

Варвара Логачева, Эдуард Клышинский
Институт прикладной математики
им. М.В.Келдыша
г. Москва

Существующие методы транскрипции (Knight and Graehl)

K. Knight and J. Graehl. 1998. Machine transliteration. Computational Linguistics, 24(4):599–612.



WFSA-A – разделение входа на слова

WFST-B – генерация фонемного состава английского языка

WFST-C – преобразование в фонемный японского языка

WFST-D – преобразование в азбуку катакана

На каждом этапе используется вероятностная информация

Статистический метод – обучение конечного автомата

Задача обратной транслитерации (японский → английский)

Генерация правил

adria;адрия
adrian;адриан
adriana;адриана
adriane;адриан
adrianna;адрианна
adrienne;адрианн
adrien;адриан
adriene;адриен
adrienn;адрианн
adrienna;адрианна
aldema;альдема
aldena;альдена
aldenaide;альденед
aldenaïse;альденез



$A \rightarrow A$

$AE \rightarrow A$

$\{<\}AI \rightarrow \text{Э}$

$\{*\}AI\{*\} \rightarrow E$

$\{<\}AY \rightarrow \text{Э}$

$\{*\}AY \rightarrow E$

$AIL\{>\} \rightarrow \text{АЙ}$

$AILL\{*\} \rightarrow \text{АЙ}$

$AILLE\{>\} \rightarrow \text{АЙ}$

Первичные правила

$$\text{isVowel} (l) = \begin{cases} \text{true, если } l \text{ – гласная} \\ \text{false, если } l \text{ – согласная} \end{cases}$$

Слово $w = l_1 l_2 \dots l_n$

Граница группы – между l_i и l_{i+1} такими, что

$$\text{isVowel} (l_i) \neq \text{isVowel} (l_{i+1})$$

R		u		gg		ie		r		o
↓		↓		↓		↓		↓		↓
R		u		dzh		e		r		o

M		a		cch		i
↓		↓		↓		↓
M		a		kk		i

Учет контекста

l → л

anjela – анжела

cella – селла

l {a, o, u, i, e, l} → л

l → ль

almelda – альмельда

avital – авиталь

l {>, t, d, p, m} → ль

Деление на слоги

A|d r i a|n a

A|д р и а|н а

A|l d e|n a

A|л ь д е|н а

A|l m e|l i|n e

A|л ь м е|л и н

Порождение сложных правил

Каждую пару слогов можно представить как

$\langle p_{i1}, \dots, p_{ik}, \mathbf{p}_x, p_{ik+1}, \dots, p_{in} \rangle \rightarrow$

$\langle c_{i1}, \dots, c_{ik}, \mathbf{c}_x, c_{ik+1}, \dots, c_{im} \rangle,$

где $\mathbf{p}_x \rightarrow \mathbf{c}_x$ – подстрока, не удовлетворяющая ни одному из существующих правил.

Можно выделить три случая несоответствия \mathbf{p}_x правилам.

Порождение сложных правил

- $рх = \emptyset, сх \neq \emptyset$ $lde \rightarrow \text{льде}$ $l \rightarrow \text{л}$
- $рх \neq \emptyset, сх = \emptyset$ $de \rightarrow \text{ьде}$ $e \rightarrow \text{е}$
- $рх \neq \emptyset, сх \neq \emptyset$ $d \rightarrow \text{ьд}$ $d \rightarrow \text{д}$
 $_ \rightarrow \text{ь}$

l →
ль

Результаты: качество обучения

Исходный язык	СТ	УСТ	ATV
Японский	7005 (99%)	4778 (68%)	1,38
Китайский	4468 (95%)	4173 (89%)	1,06
Немецкий	3484 (82%)	3247 (77%)	1,07
Арабский	2102 (99%)	1793 (85%)	1,19
Шведский	1576 (88%)	905 (50%)	1,61
Польский	1424 (99%)	1174 (81%)	1,2
Испанский	1025 (98%)	777 (74%)	1,33
Французский	678 (89%)	227 (29%)	2,66
Румынский	565 (97%)	295 (52%)	1,78
Словенский	502 (99%)	354 (70%)	1,3
Тагальский	257 (87%)	225 (76%)	1,14
Монгольский	231 (100%)	227 (98%)	1,02
Хинди	160 (99%)	152 (94%)	1,05

Результаты: качество транскрипции

Язык	Тестовая выборка	СТ	УСТ	ANL	AE
Японский	701	664 (94%)	496 (70%)	0,043	1,125
Китайский	468	436 (93%)	406 (86%)	0,026	1,23
Немецкий	421	339 (80%)	325 (77%)	0,032	1,176
Арабский	211	169 (80%)	160 (75%)	0,054	1,5
Шведский	178	149 (83%)	82 (46%)	0,131	1,33
Польский	144	138 (95%)	123 (85%)	0,027	1,5
Испанский	105	99 (94%)	76 (72%)	0,038	1,125



Конец