

Part-of-speech tagging for a Southern Min Corpus

Ching Chu Sun (June Sun)

Department of Linguistics
University of Alberta

American Association for Corpus Linguistics AACL 2009

Southern Min

A Chinese language



Chinese Corpora (McEnery & Xiao)

Corpus	POS	Channel	Variety
LCMC	Yes	Written	China
Sinica	Yes	Mixed	Taiwan
PH	No	Written	China
PFR	Yes	Written	China
LIVAV	No	Written	Mixed
SCCSD	No	Spoken	China
TREC	No	Written	China
Gigaword	No	Written	China
Callhome	No	Spoken	Mixed

No (tagged) corpus of Southern Min is publicly available.

Data sources

The data used in the current study:

Written: <http://taigu.fhl.net/index.html>

Spoken: Part of the face-to face conversations between
3 native speakers of Southern Min

HunPos

(Halácsy, Kornai, & Oravecz 2007)

- HunPos: Finding the most probable sequence of tags for each sentence
- Formula: Calculating the probability of a given sequence of tags

$$\underset{t_1 \dots t_T}{\operatorname{argmax}} P(t_{T+1} | t_T) \prod_{i=1}^T P(t_i | t_{i-1}, t_{i-2}) P(w_j | t_{i-1}, t_i)$$

HunPos

(Halácsy, Kornai, & Oravecz 2007)

The procedure for finding probability of a tag sequence for a sentence:

(1) Product of (A) and (B)

(A) $p(\text{this tag} | \text{prev. 2 tags})$

(B) $p(\text{this word} | \text{this tag \& prev.tag})$

(2) Multiply $p(\text{end tag} | \text{last tag of sequence})$

By using HunPos for tagging POS, I'd like to find out...

- Will there be any significant difference in accuracy for tagging POS in romanized (phonetic) tokens and in Chinese character tokens?
- Will the accuracy of tagging increase with more training data?
- How accurate is the tagger with unseen words (new type words)?
- What tags is the tagger most and least accurate with?

Written

Romanization

0 NUM
 hai2to2 N
 tai5oan5 PW
 e5 AM
 tek8sek4 N
 tai5oan5 PW
 si7 COP
 chit4 NUM
 e5 CF
 hai2to2 N
 , PUNC
 i1 PN
 e5 PM
 tek8cheng4 N
 toh5 ADV
 si7 COP
 hai2hoan7 N
 chin4 INT
 tng5 SV
 , PUNC

Characters

0 NUM
 海島 N
 台灣 PW
 特色 AM
 台灣 N
 台灣 PW
 是 COP
 一 NUM
 海島 CF
 海島 N
 , PUNC
 伊 PN
 特徵 PM
 特徵 N
 就是 ADV
 是 COP
 海岸線 N
 真 INT
 長 SV
 , PUNC

Spoken

Romanization

0 NUM
 tiam3 PREP
 tai5oan5 PW
 long2 ADV
 si7 COP
 kong2 V
 kok4gi3 N
 ma7 PMOD
 khah4 INT
 be7 AV
 a2 FR
 m7ku2 CCL
 chhiuN4 V
 goa2 PNPERS
 hit8 DEM
 kang4 TW
 lok8 V
 nng7 NUM
 e5 CF
 u5bo5 V
 hoN1 QM

Characters

0 NUM
 佇 PREP
 台灣 PW
 1ong2 ADV
 是 COP
 講 V
 國語 N
 0 NUM
 ma7 PMOD
 khah4 INT
 欲 AV
 a2 FR
 m7ku2 CCL
 像 V
 我 PNPERS
 彼 DEM
 kang4 TW
 錄 V
 兩 NUM
 e5 CF
 有無 V
 hoN1 QM

Methodology

	Written (1462)		Spoken (1425)	
	Romanization	Character	Romanization	Character
set 1	278	278	244	244
set 2	141	141	117	117
set 3	146	146	134	134
set 4	149	149	123	123
set 5	94	94	126	126
set 6	110	110	120	120
set 7	106	106	115	115
set 8	106	106	134	134
set 9	156	156	134	134
set 10	176	176	178	178

total tokens

tagged 1184 (wr, wh)

1181 (sr, sh)

Methodology

- Manually POS tagging principle

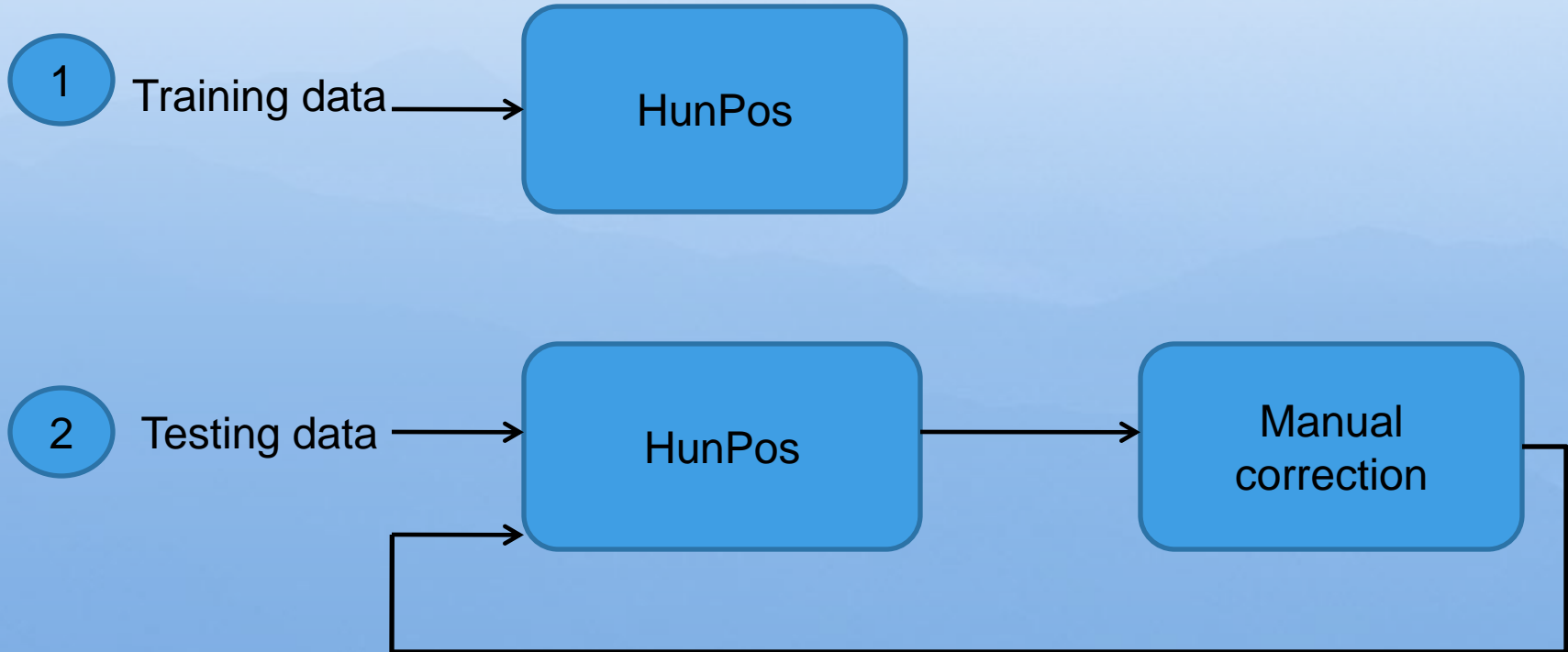
(1) Using “A Dictionary of Southern Min” by Embree (1984) as a starting point

(2) Customized tags

- Without default tags

- A new tag was added when necessary

Procedure of training the tagger

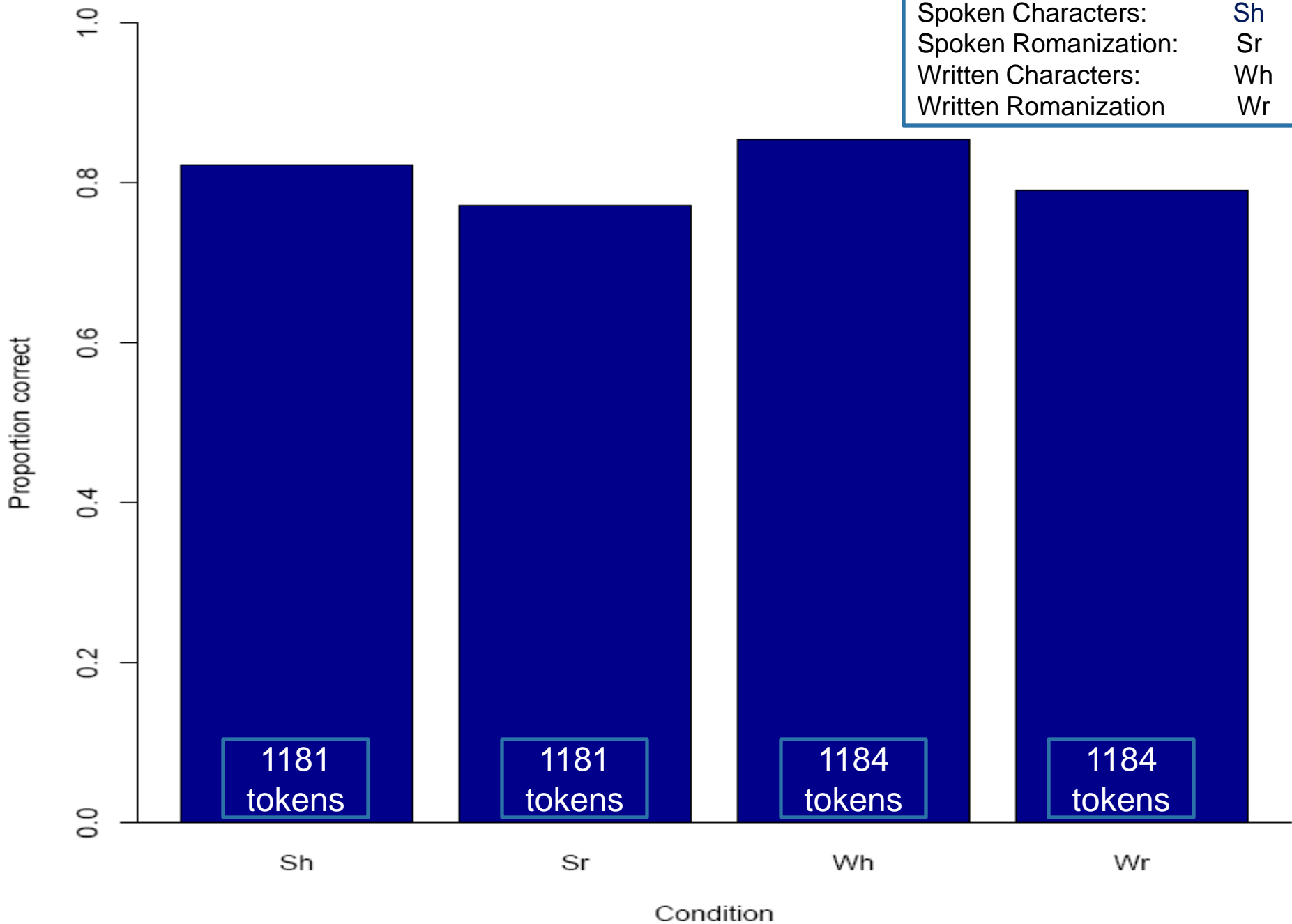


Question 1:

Will there be any significant difference in accuracy for tagging POS in romanized tokens and in Chinese character tokens?

Overall accuracy by condition

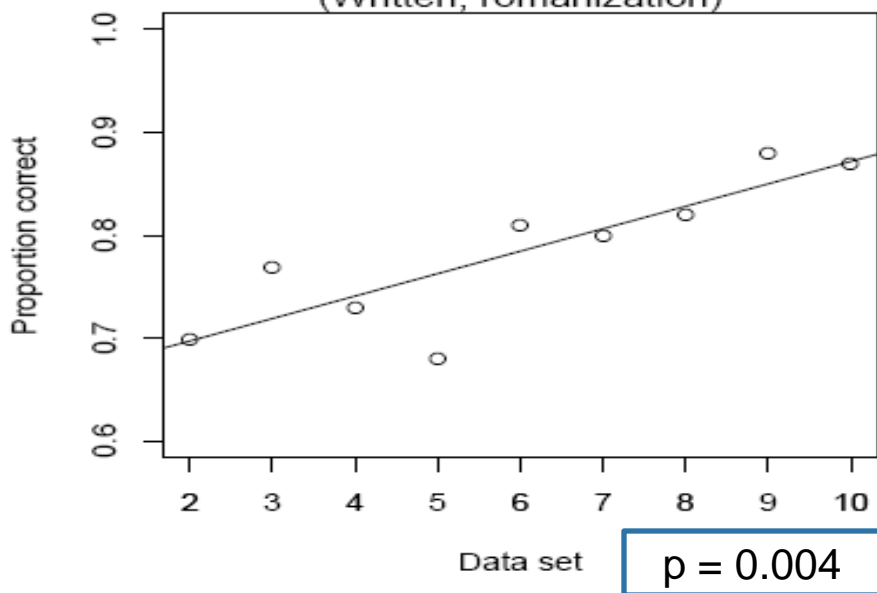
Spoken Characters:	Sh
Spoken Romanization:	Sr
Written Characters:	Wh
Written Romanization	Wr



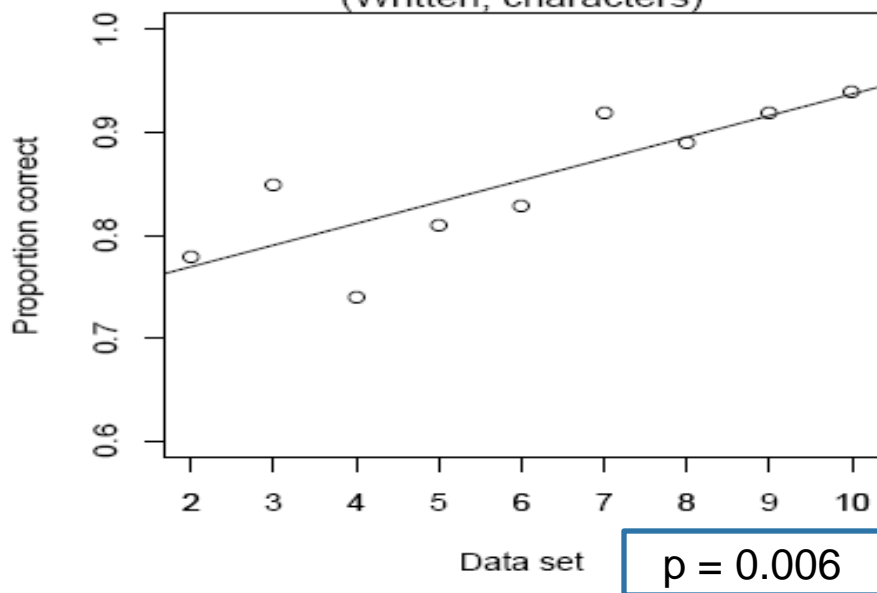
Question 2:

Will the accuracy of tagging increase with more training data?

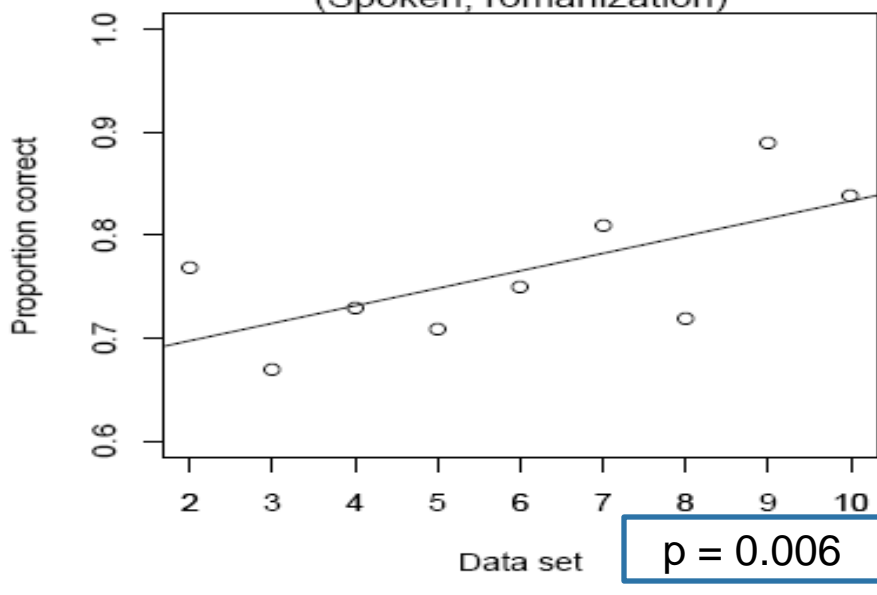
Accuracy
(Written, romanization)



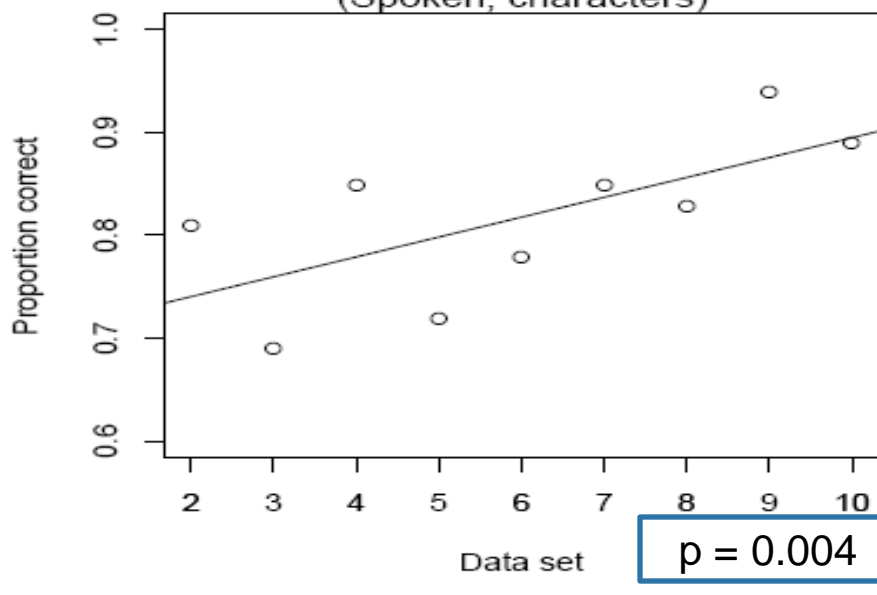
Accuracy
(Written, characters)



Accuracy
(Spoken, romanization)



Accuracy
(Spoken, characters)

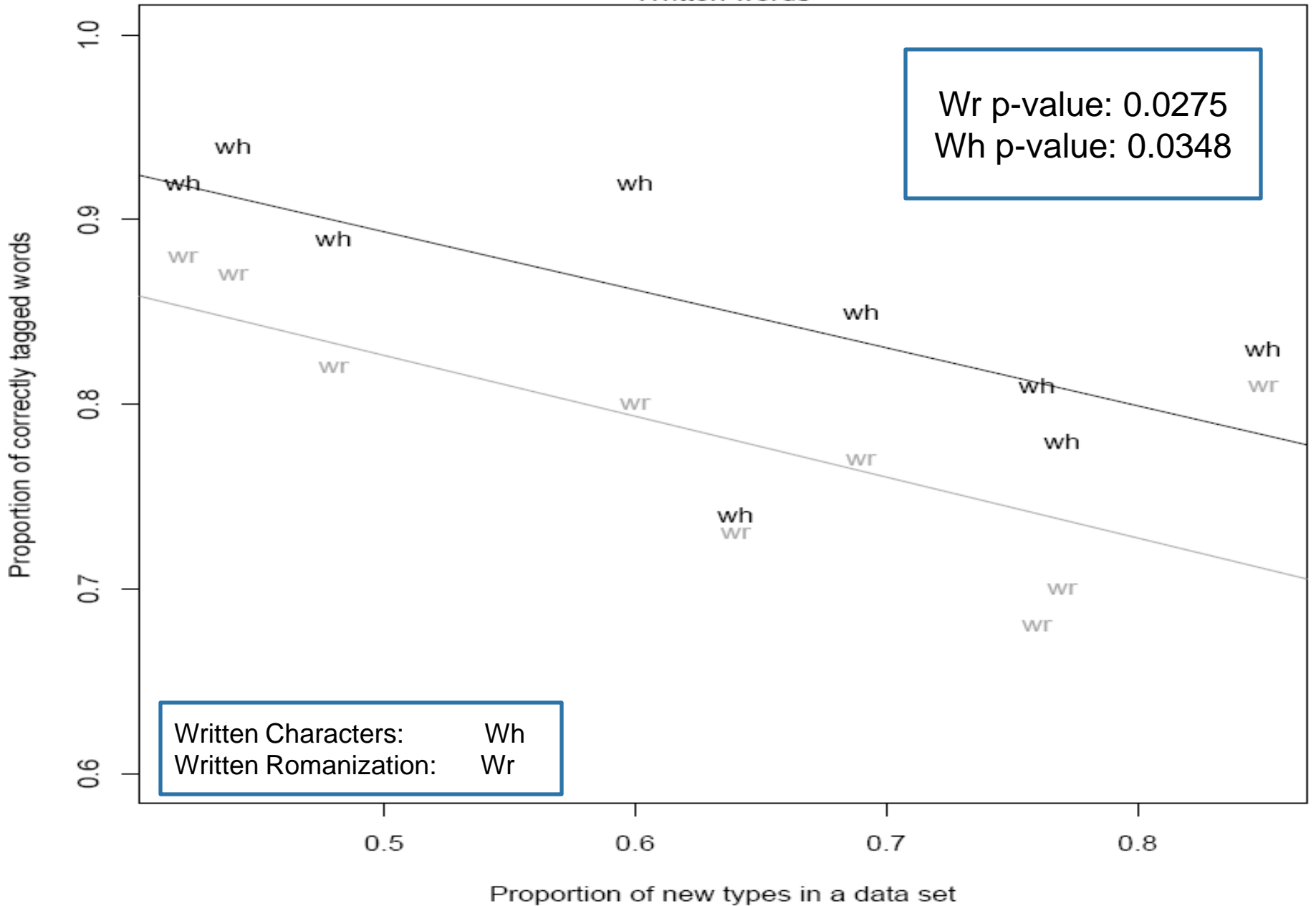


Question 3:

How accurate is the tagger with unseen words?

Proportion of correctly tagged words as a function of proportion of new types

Written words

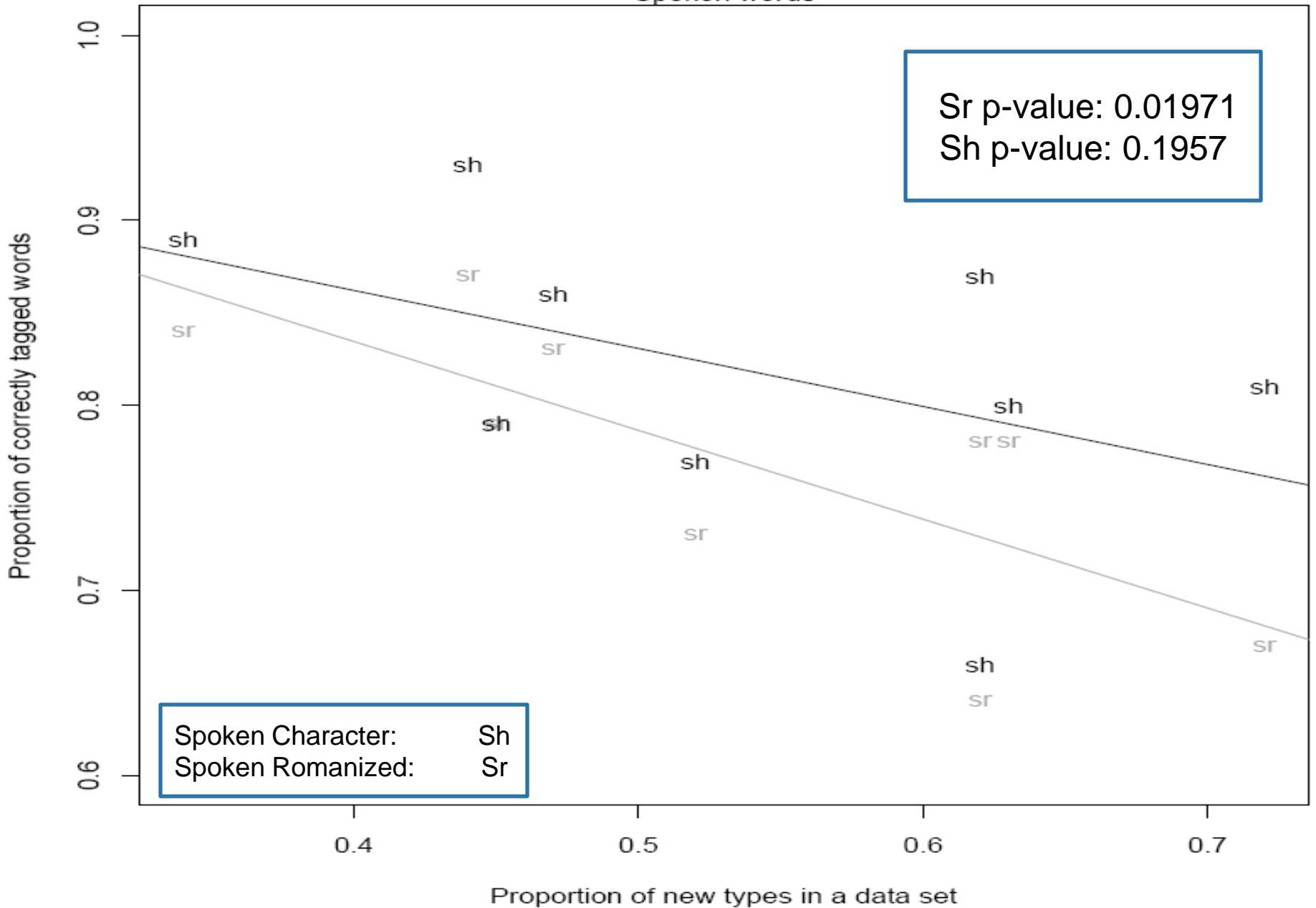


Wr p-value: 0.0275
Wh p-value: 0.0348

Written Characters: Wh
Written Romanization: Wr

Proportion of correctly tagged words as a function of proportion of new types

Spoken words



Question 4:

What tags is the tagger most and least accurate with?

Result from cross validation for individual tags

Spoken romanization

Total tags:31		Correct prop.
copula (41) progressive marker (5)	final particle(27) conjunction (8) post clause modifier (3)	100%
personal pronoun (100/104) specifier (20/22) negation1 (18/19)	clause connector (18/20) attributive marker (18/19) negation 2(11/12)	90%-99%
verb (159/185) numeral (24/29) intensifier (14/17) discourse marker (8/10)	adverb (64/76) filler (38/44) preposition (20/25) question marker(10/12) perfective marker (7/8)	80%-89%
noun (108/136) auxiliary verb (15/21)	predicate modifier (17/23)	70%-79%
place word(25/47)	classifier (19/33) time word (10/19)	50-59%
stative verb (12/32) location direction(2/5)	pronoun (5/11) numeral modifier(0/1)	Below 50%

Result from cross validation for individual tags

Spoken characters

Total tags: 31		Correct prop.
copula (41) post clause modifier (3)	determiner (3) progressive marker (5)	100%
attributive marker (18/19) personal pronoun (93/103)	negation 2(11/12)	90%-99%
verb (151/185) specifier (19/22) question marker (10/12) conjunction(7/8)	adverb (63/76) clause connector (17/20) discourse marker (8/10)	80%-89%
noun (105/136) preposition (18/24) intensifier (12/17)	numeral (22/29) predicate modifier (18/23)	70%-79%
place word (29/47) auxiliary verb (14/21) location direction (3/5)	classifier (20/33) perfective marker (5/8)	60%-69%
time word (11/19)		50-59%
stative verb (12/32) numeral modifier (0/1)	negation 1(9/19) pronoun (4/12)	Below 50%

Result from cross validation for individual tags

Written romanization

Total tags: 29			Correct prop.
sentence punctuation (24) numeral modifier (2) predicate modifier (2)	copula (16) classifier (2)	negation1 (4)	100%
attributive marker (62/64) intensifier (9/10)	place word (31/32)		90%-99%
noun (161/199) conjunction (10/12)	punctuation (124/142) time word (5/6)		80%-89%
adverb (49/63) clause connector (10/13)	preposition (14/18) post verb predicate (6/8)		70%-79%
verb (90/133) pronoun (16/24)	numeral (20/31)		60%-69%
stative verb (24/46) specifier (1/2)	negation 2 (1/2)		50-59%
location direction (5/14) personal pronoun(2/5) noun suffix (0/1)	auxiliary verb (4/12) determiner (0/1) pronoun suffix (0/1)		Below 50%

Result from cross validation for individual tags

Written characters

Total tags: 29			Correct prop.
sentence punctuation (25) numeral modifier (2)	copula (16) predicate modifier (2)	negation 1(4)	100%
attributive marker (63/64) punctuation(130/142)	intensifier (9/10) place word (29/32)		90%-99%
noun (160/199)	time word (5/6)		80%-89%
conjunction (9/12)	preposition (14/18)		70%-79%
verb (84/133) pronoun (15/24)	adverb (42/63) clause connector (8/13)		60%-69%
auxiliary verb (6/12) location direction (8/14) specifier (1/2)	classifier (1/2) negation 2 (1/2)		50-59%
stative verb (18/46) post verb predicate (2/6) determiner (0/1) pronoun suffix (0/1)	numeral (15/31) personal pronoun (2/5) noun suffix (0/1)		Below 50%

Tagging difficulty

	Spoken character	Spoken romanization	Written character	Written romanization
80%-89%	Verb (151/185)	Verb (159/185)	Noun (161/199)	Noun (160/199)
70%-79%	Noun (105/136)	Noun (108/136)		
60%-69%			Verb (90/133)	Verb (84/133)
Below 50%	Stative verb (12/32)	Stative verb (12/32)	Stative verb (24/46)	Stative verb (18/46)

Tagging difficulty

Easy for manual labor but difficult for HunPos

Stative verb

	HunPos	me
goa2	PNPERS	PNPERS
siuN7	N	V
ka7	CONJ	CONJ
he7	N	SP
kang5khuang3	N	SV
lah4	DM	DM

	HunPos	me
goa2	PNPERS	PNPERS
kam2kak4	V	V
chiok4	INT	INT
su1si4	N	SV
e5	AM	AM
kong2	V	V

	HunPos	me
li2	PNPERS	PNPERS
na7	ADV	PMOD
ti7	PREP	PREP
hak1hau7	N	N
khah4	INT	INT
chiap8	V	SV
kong2	V	V
kok4gu2	N	N
hoN1	QM	QM

	HunPos	me
ma7	PMOD	PMOD
si7	COP	COP
be7	AV	AV
lian2tng2	V	SV

	HunPos	me
gao2	PNPERS	PNPERS
chiok4	ADV	INT
tho2ia3	V	SV
kong2	V	V

	HunPos	me
so2i2	CCL	CCL
li2	PNPERS	PNPERS
su7lim5	V	PW
chian7	TW	INT
sek4	V	SV

Tagging difficulty

Easy for manual labor but difficult for HunPos

Noun

	HunPos	me
akong4	TW	N
ama2	TW	N
the3hiu1	V	V
ti7	PREP	PREP
chhu3	N	N
li2	PNPERS	DL

	HunPos	me
[Winnie]	N	PN
si7	COP	COP
ti7	PREP	PREP
tai5oan5	PW	PW
to2	ADV	ADV
bo5	NEG1	NEG1
kong2	V	V
tai5gi3	V	N
ah4	FP	FP

	HunPos	me
chhiuN1	QM	V
gaon2	PNPERS	PNPERS
ti7ti7	CCL	N
in4	PNPERS	PNPERS
giaN2	V	N

Tagging difficulty

Easy for manual labor but difficult for HunPos

Verb

	HunPos	me
goa2	PNPERS	PNPERS
siuN7	N	V
ka7	CONJ	CONJ
he7	N	SP
kang5khuang3	N	SV
lah4	DM	DM

	HunPos	me
m7	V	NEG2
bat4	ADV	AV
u5	V	V
tui3	PREP	PREP
goa7	SV	SV
hoat4tian2	N	V
,	PUNC	PUNC
a1si7	CCL	CCL

	HunPos	me
,	PUNC	PUNC
ㄝ	PUNC	PUNC
sio2kiat4	ADV	N
ㄍ	V	PUNC
koh4	ADV	ADV
toa7liong5	V	SV
ㄝ	PUNC	PUNC
ko3iong7	ADV	V
ㄍ	V	PUNC
tong4te7	N	N
lang5	N	N
,	PUNC	PUNC

Tagging difficulty

Difficult for manual labor & HunPos

	HunPos	me
hoa1jin5	N	N
e5	AM	AM
long5bin5	N	N
ju2	ADV	SV
lai5	DL	V
ju2	V	SV
che7	N	SV
,	PUNC	PUNC

ju2 lai5 ju2 che7
 ? come ? plenty
 “The more...the ...er”

	HunPos	me
i2	V	V
tit4	N	V
tioh8	ADV	PVMOD
mo7ik8	V	N
siong7	N	PREP
e5	AM	AM
li7sun5	N	N
ui7	N	V
bok8tek8	N	N

i2 tit4 tioh8...
 ? get ?...
 ‘In order to get...’

	HunPos	me
,	PUNC	PUNC
ui7	N	V
tioh8	ADV	PVMOD
beh4	AV	AV
seng4san3	V	V
kam4chia3	N	N
,	PUNC	PUNC

ui7 tioh8...
 for ?...
 ‘in order to...’

Summary

- (1) Overall: Spoken character and written character had higher tagging accuracy; spoken romanization and written romanization had lower tagging accuracy.
- (2) The larger the training data, the higher the accuracy.

Summary

- (3) The percentage of unseen words had an impact on the tagging accuracy.
- (4) The number of model tags or tokens in the training data didn't play a significant role in tagging accuracy. Certain tags got mixed up due to similar syntactic patterns.

References

- Embree, Bernard L.M. 1984. *A Dictionary of Southern Min*. Taipei Language Institute.
- Feng, Zhiwei. 2006. *Evolution and present situation of corpus research in China*. International Journal of Corpus Linguistics, 11:2, pp. 173-207. John Benjamins Publishing Company.
- Halácsy, P., Kornai, A., and Oravecz, Cs. 2007. *HunPos - an open source trigram tagger*. In Proceedings of the 45th Annual Meeting of the Association for Computational Linguistics, Companion Volume, Proceedings of the Demo and Poster Sessions, pp. 209-212. Prague, Czech Republic. 2007. Association for Computational Linguistics.
- McEnery, Anthony and Xiao Zhonghua. 2004. *The Lancaster Corpus of Mandarin Chinese: A Corpus for Monolingual and Contrastive Language Study*. Proceedings of the Fourth International Conference on Language Resources and Evaluation (LREC) 2004, pp.1175-78. Available from:
<http://www.lancs.ac.uk/postgrad/xiaoz/papers/231.pdf>
- Zhang Ruixiong. 2009. *Huang Jia San Dai Mushi* (The three generations of priests).
<http://taigu.fhl.net/index.html>.

Thank you.