

Hands on with the Switchboard Corpus

Ye Tian, Université Paris Diderot

Day 1 PM

The Switchboard Corpus

- The Switchboard-1 Telephone Speech Corpus (LDC97S62) was originally collected by Texas Instruments in 1990.
- A collection of about 2,400 two-sided telephone conversations among 543 speakers (302 male, 241 female) from all areas of the United States.
- A computer-driven robot operator system handled the calls, giving the caller appropriate recorded prompts, selecting and dialing another person (the callee) to take part in a conversation, introducing a topic for discussion and recording the speech from the two subjects into separate channels until the conversation was finished.
- About 70 topics were provided (1) no two speakers would converse together more than once and (2) no one spoke more than once on a given topic.

SWB dialog act corpus

- The Switchboard Dialog Act Corpus (SwDA) extends the Switchboard-1 Telephone Speech Corpus, Release 2, with turn/utterance-level dialog-act tags.
- The dialogue act (speech act type of each utterance) tagging allow the modelling of shallow discourse structure.
- The corpus is also annotated for disfluency and nonverbal elements such as laughter.

Getting the corpus

- A good tutorial:
 - <http://compprag.christopherpotts.net/swda.html>
- The SDA transcripts are a free download:
http://www.stanford.edu/~jurafsky/swb1_dialogact_annot.tar.gz
- The recommended method by Jurafsky and Potts is to use python NLTK package (Natural Language Toolkit). For instructions and information, follow <http://compprag.christopherpotts.net/swda.html> and <http://www.nltk.org/>.
- In this tutorial we use R.
 - I have done some minor processing on the original data file
 - Download act tag list.csv and alltrans_data.csv. The rest can be obtained from the link above.

Let's look at a data file

- Without dialogue tags
 - sw2001A-ms98-a-trans
 - sw2001A-ms98-a-word
- With dialogue act tags and other annotations:
 - sw_0001_4325.utt
- metadata

Annotations – dialogue acts

- See coder's manual for detailed informations

<http://web.stanford.edu/~jurafsky/ws97/manual.august1.html>

List of top dialogue act tags

Speech act	SWBD	Example	Cnt	%
Statement-non-opinion	sd	Me, I'm in the legal department.	72,824	36%
Acknowledge (Backchannel)	b	Uh-huh.	37,096	19%
Statement-opinion	sv	I think it's great	25,197	13%
Agree/Accept	aa	That's exactly it.	10,820	5%
Abandoned or Turn-Exit	% -	So, -	10,569	5%
Appreciation	ba	I can imagine.	4,633	2%
Yes-No-Question	qy	Do you have to have any special training?	4,624	2%
Non-verbal	x	[Laughter], [Throat_clearing]	3,548	2%
Yes answers	ny	Yes.	2,934	1%
Conventional-closing	fc	Well, it's been nice talking to you.	2,486	1%
Wh-Question	qw	Well, how old are you?	1,911	1%
No answers	nn	No.	1,340	1%
Response Acknowledgement	bk	Oh, okay.	1,277	1%
Hedge	h	I don't know if I'm making any sense or not.	1,182	1%
Declarative Yes-No-Question	qy^d	So you can afford to get a house?	1,174	1%

Annotations – PoS tags

- The Penn Treebank Part-of-Speech Tagset

Tag	Description	Example	Tag	Description	Example
CC	coordin. conjunction	<i>and, but, or</i>	SYM	symbol	<i>+, %, &</i>
CD	cardinal number	<i>one, two</i>	TO	“to”	<i>to</i>
DT	determiner	<i>a, the</i>	UH	interjection	<i>ah, oops</i>
EX	existential ‘there’	<i>there</i>	VB	verb base form	<i>eat</i>
FW	foreign word	<i>mea culpa</i>	VBD	verb past tense	<i>ate</i>
IN	preposition/sub-conj	<i>of, in, by</i>	VBG	verb gerund	<i>eating</i>
JJ	adjective	<i>yellow</i>	VCN	verb past participle	<i>eaten</i>
JJR	adj., comparative	<i>bigger</i>	VBP	verb non-3sg pres	<i>eat</i>
JJS	adj., superlative	<i>wildest</i>	VBZ	verb 3sg pres	<i>eats</i>
LS	list item marker	<i>1, 2, One</i>	WDT	wh-determiner	<i>which, that</i>
MD	modal	<i>can, should</i>	WP	wh-pronoun	<i>what, who</i>
NN	noun, sing. or mass	<i>llama</i>	WP\$	possessive wh-	<i>whose</i>
NNS	noun, plural	<i>llamas</i>	WRB	wh-adverb	<i>how, where</i>
NNP	proper noun, sing.	<i>IBM</i>	\$	dollar sign	<i>\$</i>
NNPS	proper noun, plural	<i>Carolinas</i>	#	pound sign	<i>#</i>
PDT	predeterminer	<i>all, both</i>	“	left quote	<i>‘ or “</i>
POS	possessive ending	<i>’s</i>	”	right quote	<i>’ or ”</i>
PRP	personal pronoun	<i>I, you, he</i>	(left parenthesis	<i>[, (, {, <</i>
PRP\$	possessive pronoun	<i>your, one’s</i>)	right parenthesis	<i>],), }, ></i>
RB	adverb	<i>quickly, never</i>	,	comma	<i>,</i>
RBR	adverb, comparative	<i>faster</i>	.	sentence-final punc	<i>. ! ?</i>
RBS	adverb, superlative	<i>fastest</i>	:	mid-sentence punc	<i>: ; ... - -</i>
RP	particle	<i>up, off</i>			

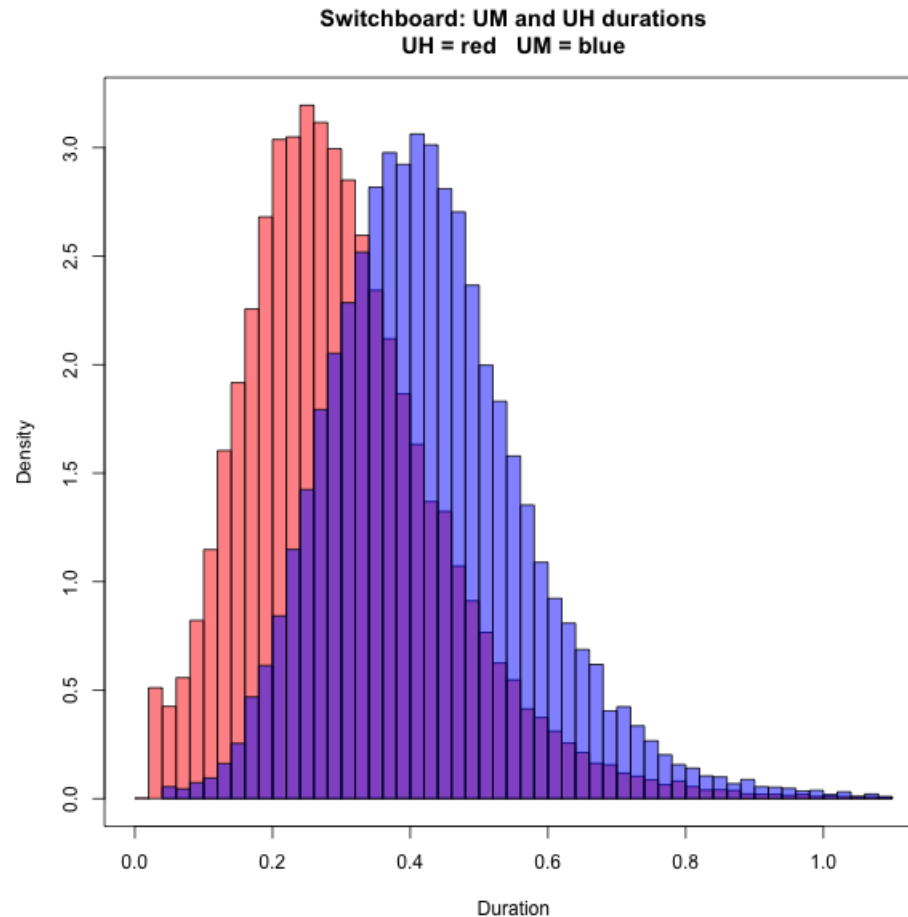
Figure 9.1 Penn Treebank part-of-speech tags (including punctuation).

Annotations – parsed trees

```
>>> [tree.pprint() for tree in utt.trees]
["(S
  (CC And)
  (NP-SBJ (PRP it))
  (VP
    (BES 's)
    (NP-PRD
      (NP (DT a) (JJ small) (NN office))
      (SBAR
        (WHNP-1 (WDT that))
        (S
          (NP-SBJ (PRP she))
          (VP (VBZ works) (PP-LOC (RB in) (NP (-NONE- *T*-1))))))
      (-DFL- E_S))"]
>>> utt.tree_lemmas(wn_lemmatize=True)
[('And', 'CC'), ('it', 'PRP'), (''s', 'BES'), ('a', 'DT'), ('small', 'JJ'), \
('office', 'NN'), ('that', 'WDT'), ('she', 'PRP'), ('works', 'VBZ'), ('in', 'RB'), \
('*T*-1', '-NONE-'), ('E_S', '-DFL-')]
```

Example studies

- <http://languagelog.idc.upenn.edu/nll/?p=14991>



Let's do something

- First a bit of regular expressions
 - Search for a word/phrase
 - Search for dialogue acts
 - Search for words with specific PoS
-
- Files used:
 - Switchboard search engine.R
 - swda tagged_treebank search engine.R

Regular expressions

- Regular expressions are the standard notation for characterizing text sequences.
- First developed by Kleene (1956)
- We use regular expressions to search for patterns in a corpus.
- The simplest kind of regular expression is a sequence of characters. E.g. /linguistics/.
- Regular expressions are case sensitive.
- Let's work on some examples
- https://en.wikibooks.org/wiki/R_Programming/Text_Processing

References

- Kleene, Stephen C. (1956). Shannon, Claude E.; McCarthy, John, eds. Representation of Events in Nerve Nets and Finite Automata. Automata Studies. Princeton University Press. pp. 3–42.