

Estimating Word's Predictability on Lexical Processing

using Latent Semantic Analysis – Verification from Eye Movement Data

Hsueh-Cheng Wang¹, Minglei Chen¹, Hwawei Ko¹, and Walter Kintsch²

¹ National Central University & ² Colorado University at Boulder

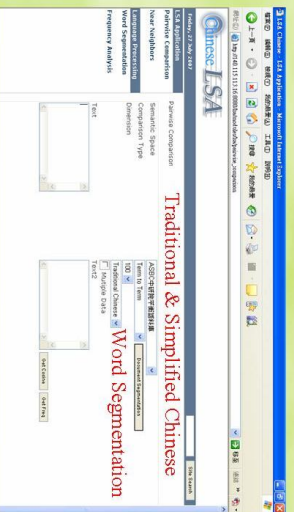
Word Predictability & Lexical Processing

- E-Z Reader Model (Reichle, et al., 2003)
 - Early Stage (L1)
 - Orthographic Form
 - Familiar Check
 - $t(L_1) = [\beta_1 r - \beta_2 \ln(\text{frequency})](1 - \theta \text{predictability})$
 - Late Stage (L2)
 - Phonological / Semantic Form
 - Completion of Lexical Access
 - $t(L_2) = \Delta(\beta_1 r - \beta_2 \ln(\text{frequency}))(1 - \text{predictability})$
- Predictability Effect is stronger on L2 than on L1
- SWIFT Model (Engbert, et al., 2002)
 - $L_0 = (1 - p_0) L_0$ $L_0 = \alpha - \beta \log(f_j)$

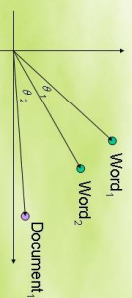
- High-predictable words are read faster than low-predictable words.
 - First Fixation Duration
 - Gaze Duration
 - Total Time
- High-predictable words are skipped more and re-fixated less than low-predictable words.
 - Skip Rate
 - Regression
- Chinese reader exploit target word predictability during reading. (Rayner, et al. 2005)

Latent Semantic Analysis (LSA)

- From ASBC (5M Chinese Words Corpus)
 - Term-to-Documents Co-occurrence Matrix
 - 49021 Terms (Words)
 - 40463 Documents (Paragraphs)
 - Local & Global Weighting
 - Singular Value Decomposition (SVD)
 - High-Dimension Semantic Space
 - Mean: 0.067 Std: 0.115



Pairwise Comparison Application



Estimates of Predictability by LSA, Cloze Task and Transitional Probability

LSA vs. Cloze Task

- Re-analysis of Materials for P/U targets in Rayner, et al., 2004.
- Semantic Space of General English Reading (<http://lsa.colorado.edu/>)

Context_Pre	Context_Post	Target Word	Predictability	Freq	LSA
Most cowboys know how to ride a	if necessary.	horse	P	H	0.62
Most cowboys know how to ride a	if necessary.	canal	U	L	0.1
In the desert, many Arabe ride a	to get around.	horse	U	H	0.29
In the desert, many Arabe ride a	to get around.	canal	P	L	0.64

LSA CAN distinguish High- / Low- Predictability

Eye Movement Data vs. Estimated Predictability

Correlations Among Predictors

Variable	Infreq	Wordlength	Strokes	LSA	TP
Infreq	—	-.639	-.570	.180	.306
Wordlength		—	.044	-.128	-.318
Strokes			—	-.116	-.278
LSA				—	.301
TP					—

Unstandardized Regression Coefficients

Variable	Mean	Std D	FFD	Mean	Std D	GD	Mean	Std D	TotalTime
W/Freq	-5.2**	3.3		-6.6**	3.9		-10.9**	9.7	
Wordlength	24.5**	8.7		35.1**	13.6		113.8**	31.5	
Strokes	1.1**	0.7		1.5**	0.9		3.7**	1.7	
LSA	-5.2	14.4		-14.1	29.2		-76.2**	55.5	
Predictability	-56.1**	46.0		-43.5**	47.0		36.6	174.5	
TP									

descript	ride	candidate	explanation	TP	LSA (沙漠)
骑 马	马	[mao] horse	horse	0.231	0.35
骑 马	马	[mao] horse	horse	0.139	0.17
骑 马	马	[mao] horse	horse	0.134	0.07
骑 马	马	[mao] horse	horse	0.134	0.07
骑 马	马	[mao] horse	horse	0.134	0.07
骑 马	马	[mao] horse	horse	0.134	0.07
骑 马	马	[mao] horse	horse	0.134	0.07
骑 马	马	[mao] horse	horse	0.134	0.07
骑 马	马	[mao] horse	horse	0.134	0.07
骑 马	马	[mao] horse	horse	0.134	0.07

- LSA might be adopted to calculate contextual constraint and represent word predictability on lexical processing.

Poster and Oral Presentation in 14th European Conference of Eye Movements



ECCEM 2007

14th European Conference on Eye Movements
August 19-23, Potsdam, Germany