



Neologisms and the Lexicon

Neologisms are novel word forms (New-Word)

Definition

E.g., "soccer mom", "neocon", "gastropub", "affluenza", "chicken hawk"

Neologisms enter the lexicon gradually

Bubbling Under

Can be used without gaining widespread recognition for years

Wikipedia

Reflects Cultural Change

Updated more frequently than print dictionaries, captures the "Zeitgeist"

· ZeitGeist

Neologism Harvester

Trawls Wikipedia looking for neologisms (w.r.t WordNet) and their meanings

Lexicographic Approaches to Neologism Analysis

Explanatory Lexicography

Post-Hoc Analysis

Seeks to dissect and explain neologisms after they have been identified in text

• Predictive Lexicography

Pre-Hoc Creativity

Seeks to predict and verify neologisms using principles of word-formation

• Explanatory Lexicography at Work

Example #1

"Affluenza" = "Affluent" + "Influenza" (overlap: "...fluen...")

Predictive Lexicography at Work

Example # 2

"Chrono-" (= time) + "-onaut" (=traveller) = "chrononaut" (a "time traveller")

Trawling Wikipedia

Uni-Directional Links

Notation: $A \longrightarrow B$

Headwords accompanied by text article that links to relevant other articles

• Bidirectional Reciprocated Links

Notation: $A \leftarrow \longrightarrow B$

Two headwords cross-reference each other, indicating strong mutual association

Contiguous Links

Notation: $A \longrightarrow B$; C

The Wikipedia article for A contains a link to B followed directly by a link to C

Text-Free Analysis (No Parsing Required)

Topological Approach

The text of each article is ignored; only topology of headword xrefs is used

Compound Terms

• Simplest: ADJ-Noun and Noun-Noun collocations Prevalent in English/ PWN

Yokes two branches of an ontology together: e.g., Religious-Music, Chinese-Cuisine

General Structure

Endocentricity

<Modifier : Head>, where Head denotes a hypernym of the compound meaning

Awkward Exceptions

Exocentricity

<Modifier : Head> does not denote a hyponym of {Head}, e.g., "hammer head"

• Inter-Compound Relationships

Specialization of Parts

Applied-Science :: Applied-Physics but Religious-Music :: Christian-Music

Compound Schema I: Head Specialization

Endocentric form: α β

$$\frac{\alpha_{\beta} \rightarrow \alpha_{\gamma} \quad \wedge \quad \beta \text{ is a } \gamma_{\gamma}}{\alpha_{\beta} \text{ is a } \alpha_{\gamma}}$$

PWN senses of β and γ given by this relationship

Examples

Fantasy_Sport -> "Fantasy_Football and Football is a Sport (in PWN)

Vetinary_Science → Vetinary_Medicine and Medicine isa Science (in PWN)

Compound Schema II: Head Anchoring

Anchored in non-compound

$$\frac{\alpha_{\beta} \rightarrow \gamma \quad \wedge \quad \beta \text{ is a } \gamma}{\alpha_{\beta} \text{ is a } \gamma}$$

PWN senses of β and γ given by this relationship

Examples

Applied_Statistics -> Science and Statistics is a Science (in PWN)

British_Rail -> Railway and Rail isa (shorthand for) Railway (in PWN)

Compound Schema III: Compound Expansion

Synonym Creation

$$\frac{\alpha_{\beta} \rightarrow \alpha \quad \wedge \quad \alpha \text{ is a } \beta}{\alpha_{\beta} \text{ synonym of } \alpha}$$

PWN senses of β and γ given by this relationship

Examples

Oscar_award → Oscar and Oscar isa award (in PWN)

Semitic_language -> Semitic and Semitic is a Language (in PWN)

Compound Schema IV: Compound Conflation

Schema linking

$$\alpha_{\beta} \rightarrow \beta_{\gamma} \wedge \beta \text{ is a } \gamma$$

$$\alpha_{\beta} \text{ is a } \beta_{\gamma} \leftarrow \cdots$$

β_γ may be a new compound, previously learned

Examples

Touch_rugby → Rugby_football and Rugby isa football (in PWN)

Pop_punk -> Punk_Rock and Punk isa Rock (in PWN)

Compound Schema V: Compound Opposition

Antonym Creation

$$\alpha_{\beta} \rightarrow \alpha_{\gamma} \wedge \gamma \text{ antonym } \beta$$

$$\alpha_{\beta} \text{ antonym of } \alpha_{\gamma} \sim PWN \text{ senses of } \beta$$

PWN senses of β and γ given by this relationship

Examples

Second_language -> First_language and Second antonym of First (in PWN)

Synthetic_geometry -> Analytic_geometry and Synthetic opposes Analytic

Compound Schema VI: Head Expansion

See also schema III

$$\alpha_{\beta} \rightarrow \gamma_{\beta} \wedge \gamma_{\beta}$$
 synonym of $\beta \in ...$

PWN senses of β and γ_{β} given by this

Examples

Escort_Carrier -> Aircraft_carrier and Aircraft_carrier syn. of Carrier

Simple_majority → Absolute_Majority and Absolute_Majority syn. Majority

Compound Schema VII: Modifier Specialization

+ vice versa

 $\alpha_{\beta} \rightarrow \gamma_{\beta} \wedge \alpha \mod -isa \gamma$ $\alpha_{\beta} isa \gamma_{\beta}$

PWN senses of α and γ given by this relationship

Examples

Truck_racing -> Auto_racing and Truck isa <u>auto</u>motive_Vehicle (in PWN)

Hindu_music → Religious_Music and Hindu isa Religious_person in (PWN)

Compound Schema VIII: Catch-all Coordination

Weakest schema

$$\frac{\alpha_{\beta} \rightarrow \gamma_{\beta} \wedge \gamma_{\beta} \text{ is a } \beta}{\alpha_{\beta} \text{ coordinate of } \gamma_{\beta} + \beta}$$

PWN senses of β

and γ_β given

by this

Examples

Financial_mathematics → applied_mathematics (isa mathematics)

Constructed_language → natural_language (isa language)

Endocentric Compound Evaluation: Results

For 10,899 Wikipedia headwords matching one or more compound schema:

<i>E.g.</i> ,				
	Schema	# Headwords	# Error	Precision
Dutch_language isa German_language	I	14%	0	1.0
	II	12%	0	1.0
	III	13%	0	1.0
	IV	4%	0	1.0
E.g., Supreme_Court isa state_court	V	1%	0	1.0
	VI	3%	0	1.0
	→ VII	15%	7%	.93
	VIII	70%	0	1.0

In-depth Analysis: Portmanteau Words

Portmanteau (double-pocket) words

Lewis Carroll

A Textual blend of two different words, e.g., "Bollywood", "Infomercial"

General Structure

Prefix + Suffix

One words contributes a prefix, the other a suffix (but, e.g., "Modem")

True Portmanteaux

Double-Scope Blends

Neither component word is present in its entirety, (e.g., "metrosexual")

• Impure Portmanteau

Broad Exceptions

E.g., "Wikipedia", (but not "Wiktionary"), "Gastropub", "Feminazi", ...

Taxonomic Connectives

Precise Taxonomic Placement

Pure ISA

E.g., "Superhero" = "Super-" + "Hero" ⇒ Superhero ISA Hero

• Approximate Taxonomic Placement

Hedging

E.g., "Spintronics" = "Spin" + "Electronics" \Rightarrow Spintronics hedges Electronics

Disambiguation

Sense Priming

E.g., Which sense of "hero" does "Superhero" extend? (not a sandwich obviously)

• Hedging Supports Figurative Portmanteaux

Metaphor

"Affluenza" = "Affluence" + "Influenza" but Affluenza NOT-ISA Influenza

General Approach: Two-Pass Harvesting

A Textual, String-Matching Approach

String-Matching

Let $\alpha\beta$ represent the general form of a headword; analyse with schemata

The Topological Context

Wikipedia

The set of Wikipedia cross-references for a given headword lphaeta

Pass I: Learning from Easy Cases

Easy

Harvest obvious examples (in rich contexts) first, and learn from these cases

Pass II: Applying learnt patterns to Hard Cases

Hard

When topological context is insufficient, use experience of easy cases as a guide

Portmanteau Schema I: Explicit Extension

Easy

$$\frac{\alpha\beta \rightarrow \beta \quad \wedge \quad \alpha\beta \rightarrow \alpha\gamma}{\alpha\beta \quad isa \quad \beta}$$

Examples

Gastropub ("Gastropub" \rightarrow "pub" and "Gastropub" \rightarrow "Gastronomy")

Feminazi ("Feminazi" → "Nazi" and "Feminazi" → "Feminism")

Portmanteau Schema II: Suffix Alternation

Easy

$$\frac{\alpha\beta \rightarrow \alpha\gamma \quad \wedge \quad \beta \leftrightarrow \gamma}{\alpha\beta \quad hedges \quad \alpha\gamma}$$

Examples

"man" → "boy" "woman" → "girl" "bit" → "byte" "toxin" → "bacteria"

Fangirl ("Fangirl" \rightarrow "Fanboy" and "boy" \rightarrow "girl")

Portmanteau Schema III: Partial Suffix

Easy

$$\frac{\alpha\beta \rightarrow \gamma\beta \wedge (\alpha\beta \rightarrow \alpha \vee \alpha\beta \rightarrow \delta \rightarrow \alpha)}{\alpha\beta \text{ hedges } \gamma\beta}$$

Examples

"Metrosexual" -> "Heterosexual" \Lambda "Metrosexual" -> "Metro"

"Pomosexual" \rightarrow "Homosexual" \wedge "Pomosexual" \rightarrow "Postmodernism" \rightarrow "pomo"

Portmanteau Schema IV: Consecutive Blending

Easy

$$\alpha\beta \rightarrow \alpha\gamma$$
; $\delta\beta$

αβ hedges δβ

Examples

"sharpedo" → "shark"; "torpedo" hedges "torpedo"

"Spanglish" -> "Spanish"; "English" hedges "English"

Portmanteau Schema IVa: Partial Suffix

Easy

$$\alpha\beta \rightarrow \alpha\gamma$$
; $\delta\beta \wedge \alpha\beta \rightarrow portmanteau$

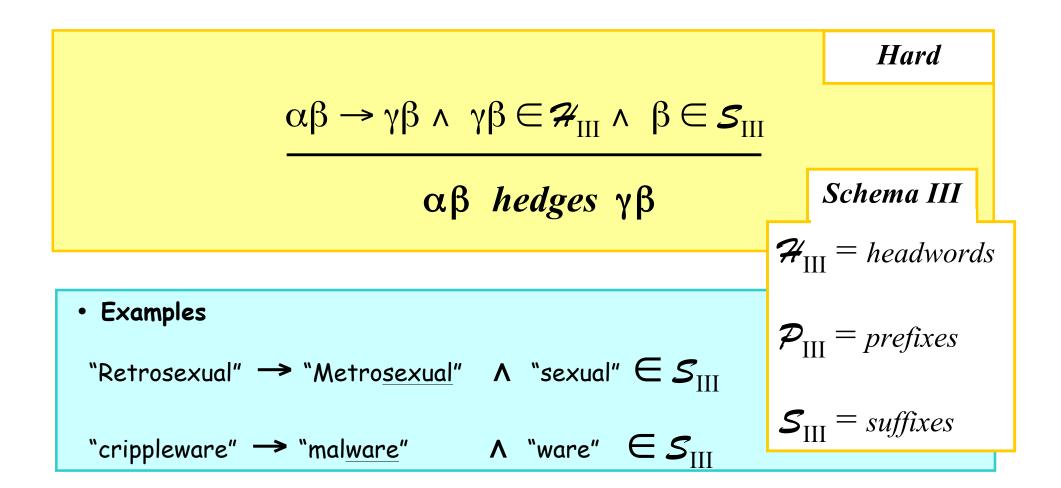
αβ hedges γβ

Examples

"Spork" \rightarrow "Spoon" Λ "Fork" hedges "English" (2 characters)

"Sporgery" \rightarrow "Spam" Λ "Forgery" hedges "Forgery" (2 characters)

Portmanteau Schema V: Suffix Completion



Portmanteau Schema VI: Separable Suffix

Hard

$$\frac{\alpha\beta \to \beta \land \alpha \in (\mathcal{P}_{I} \cup \mathcal{P}_{II} \cup \mathcal{P}_{III})}{}$$

αβ isa β

Examples

"Gastroshop" → "shop"

\Lambda "gastro" $\in \mathcal{P}_{\mathrm{III}}$

"antiprism" → "prism"

 \wedge "anti" $\in \mathcal{P}_{ ext{III}}$

Portmanteau Schema VII: Prefix Completion

Hard

$$\alpha \gamma \rightarrow \alpha \land \langle \gamma, \beta \rangle \in 7_{I}$$

Schema I

αβ isa β

Suffix replacement pairings for schema I

Examples

"Logicism" \rightarrow "logic" \land <"ism", "Nazi"> $\in 7_I$ so logicnazi ISA Nazi

"Psychology" o "psycho" Λ <"ology", "technology"> $\in \mathcal{7}_{\mathrm{I}}$

Portmanteau Schema VIII: Recombination

Hard

$$\alpha\beta \rightarrow \alpha\gamma \wedge \alpha\beta \rightarrow \delta\beta \wedge \alpha \in \mathcal{P}_{_{\text{III}}} \wedge \beta \in \mathcal{S}_{_{\text{III}}}$$

αβ hedges δβ

Examples

"geonym"
$$\rightarrow$$
 "geography" Λ "geonym" \rightarrow "toponym" (hedges geonym)

"dubtitle"
$$\rightarrow$$
 "dubbed" Λ "dubtitle" \rightarrow "subtitle" (hedges subtitle)

Evaluation: Set-Up

• # of Atomic Headwords in Wikipedia (June 2005)

Wikipedia

Wikipedia contains 152,060 atomic headwords at this time of download

WordNet Version: 1.6

WordNet

Few differences were achieved using WordNet 2.1 instead

• # of Wikipedia Entries Matching a ZeitGeist schema

Initial Selection

4676 headwords match at least one Zeitgeist schema

Metaphor is a useful too for ontological development

Pre-Filtering

1385 already in WordNet; 1083 analyses yield non-PWN parent or hedge

Portmanteau Evaluation: Results

For 2048 Wikipedia headwords matching one or more Portmanteau schema:

	•	•	•		
E.g.,					
	Schema	# Headwords		# Error	Precision
Rubbergate from Watergate	I	710	29%	11	.985
	II	144	5%	0	1.0
	III	330	13%	5	.985
	IV	82	3%	2	.975
E.g., Retrosexual from Metrosexual	V	161	6%	0	1.0
	VI	321	13%	16	.95
	VII	340	14%	32	.90
	VIII	320	13%	11	.965

Conclusions

• A Linguistics-Lite Approach to Neologisms is Feasible

Lightweight

No text parsing or morphological analysis; all relevant morphemes are learned

Taxonomic Hedging is required

Uncertainty

Word-forms are not deterministic w.r.t. taxonomic placement, approx. needed

• Link Topology offers context-specific insights

Grounding

E.g., Microsurgery \rightarrow <u>Microscopy</u> + surgery \Rightarrow "surgery done with a microscope"

Not biased towards English

Multilingual

Linguistics-lite means no language-bias - applicable to other languages / wikis?