

# Computational Lexical Semantics: Methods and Applications

**Alexander Panchenko**

Language Technology Group, TU Darmstadt, Germany  
panchenko@lt.informatik.tu-darmstadt.de

September 12, 2015

# Plan

- 1 Computational Lexical Semantics
- 2 Semantic Similarity
- 3 Pattern-Based Semantic Similarity Measure
- 4 Hybrid Semantic Similarity Measure
- 5 Word Embeddings
- 6 Applications of Semantic Similarity Measures

# Plan

- 1 Computational Lexical Semantics
- 2 Semantic Similarity
- 3 Pattern-Based Semantic Similarity Measure
- 4 Hybrid Semantic Similarity Measure
- 5 Word Embeddings
- 6 Applications of Semantic Similarity Measures
  - Lexico-Semantic Search
  - Filename Categorization

# Motivation

Я обычно вечером пью **Ягуар** с пацанами на районе.

**Отвертку**

**Водку**

**Пиво**

**Коньяк**



Я наслаждаюсь вождением моего **Ягуара** 68-го  
года выпуска по воскресеньям.



**Бентли**

**Мерседес**

**кабриолет**

**олд-таймер**

# Two Worlds of Computational Lexical Semantics

## Human-Driven Approach



BabelNet



Wiktionary  
[ˈwɪkʃənəri] n.,  
a wiki-based Open  
Content dictionary



UBY

?

## Data-Driven Approach

**JoBimText**  
Linking Language to Knowledge  
with Distributional Semantics

word2vec

Tool for computing continuous distributed representations of words.

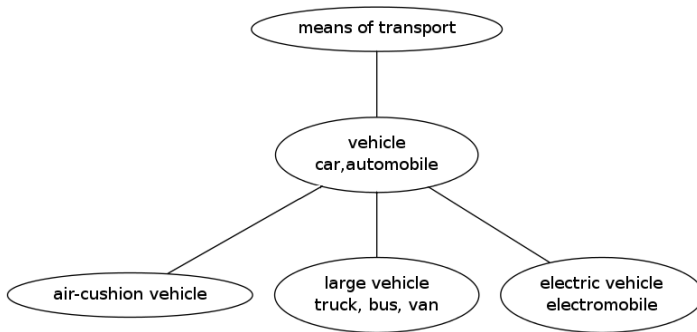


**Semantic Vectors**

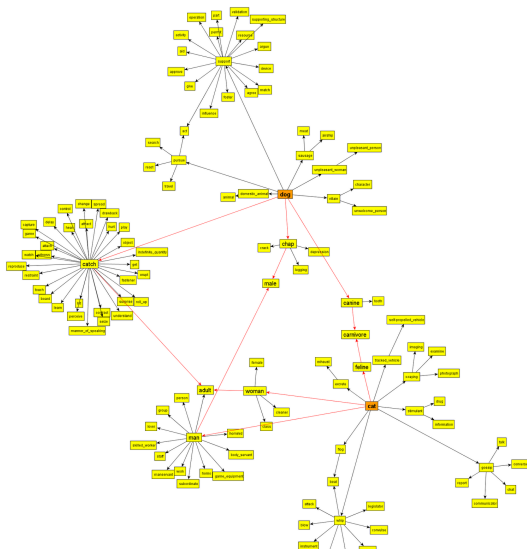
# Semantic Resources

A **semantic resource** is an graph  $(C, R)$ :

- nodes  $C$  represent **terms**;
- edges  $R$  represent untyped **semantic relations**.



## Semantic Resources: WordNet



# A multilingual WordNet: BabelNet.org

[LOG IN](#) [REGISTER](#)

python

ENGLISH

TRANSLATE INTO...

SEARCH

[PREFERENCES](#)

All

Concepts

Named Entities

17 results

■ Noun

Noun



python

Large Old World boar

ID: 00065461n | Concept



python

A soothsaying spirit or a person who is possessed by such a spirit

ID: 00065462n | Concept



Python

(Greek mythology) dragon killed by Apollo at Delphi

ID: 00065463n | Named Entity



Python (programming language)

Python is a widely used general-purpose, high-level programming language.

ID: 01713224n | Named Entity



Monty Python, Python (Monty)

Monty Python was a British surreal comedy group that created Monty Python's Flying Circus, a British television comedy sketch show that first aired on the BBC on 5 October 1969.

ID: 01157670n | Named Entity



Pythonidae, family Pythonidae, Python (snake)

In some classifications a family separate from Boidae comprising



# A multilingual WordNet: BabelNet.org



BabelNet

**Dictionary**

- Images
- Translations
- Sources
- Categories
- External links

[LOG IN](#) [REGISTER](#)

python

ENGLISH

TRANSLATE INTO

SEARCH

INTO...

[PREFERENCES](#)

English

Arabic

Chinese

French

German

Greek

Hebrew

[+ all preferred languages](#)

bn:01713224n • NOUN • Named Entity • Categories: Class-based programming languages, Cross-platform free software, Dutch inventions, Dynamically typed programming languages...

🇬🇧 **Python (programming language)** • Python philosophy • Python programming • Python 3K • Python programming language • Pythonista • .pyc • PEP8 • Python prog • Python computer language • Py3K • PythonLanguage • Python3 • .pyo • Python2 • Pyston • Python code • .py (file extension) • The Zen of Python • Python3000 • Python implementations • Python 2 • Python scripting language • Pythonistas • Python language • Python Programming Language language • Pythonic • Python Enhancement Proposal

[More](#)

Python is a widely used general-purpose, high-level programming language. [definitions](#)

IS-A: [high-level language](#) • [Dynamic programming language](#) • [object-oriented programming language](#)

**EXPLORE NETWORK****Translations**

Python, Python语言, Python程序设计语言, Python编程语言

پایتھن, Python, پایتون, پایتون, پایتون

🇬🇧 Python, Python philosophy, Python programming, Python 3K, Python programming language, Pythonista, .pyc, PEP8, Python prog, Python computer language, Py3K, PythonLanguage, Python3, .pyo, Python2, Pyston, Python code, .py, The Zen of Python, Python3000, Python implementations, Python 2, Python scripting language.

# Plan

- 1 Computational Lexical Semantics
- 2 Semantic Similarity**
- 3 Pattern-Based Semantic Similarity Measure
- 4 Hybrid Semantic Similarity Measure
- 5 Word Embeddings
- 6 Applications of Semantic Similarity Measures
  - Lexico-Semantic Search
  - Filename Categorization

# Semantic Similarity Measures

## Definition

A semantic similarity measure quantifies semantic relatedness input terms  $c_i, c_j$  with the similarity score  $s_{ij} = \text{sim}(c_i, c_j)$ :

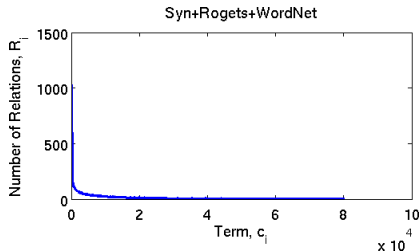
$$s_{ij} = \begin{cases} 1 & \text{if } \langle c_i, c_j \rangle \text{ is a pair of } \textit{syn}, \textit{hyper}, \textit{cohyponym} \\ 0 & \text{otherwise} \end{cases}$$

## Properties

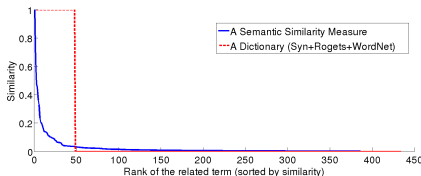
- Nonnegativity:  $0 \leq s_{ij} \leq 1$ ;
- Reflexivity:  $s_{ij} = 1 \Leftrightarrow c_i = c_j$ ;
- Symmetry:  $s_{ij} = s_{ji}$ ;

# Semantic Similarity Measures

- Many dissimilar pairs, few similar pairs:  $s_{ij} \sim \exp(\lambda)$ :



- Similarity distribution of the term “doctor”:



# Plan

- 1 Computational Lexical Semantics
- 2 Semantic Similarity
- 3 Pattern-Based Semantic Similarity Measure**
- 4 Hybrid Semantic Similarity Measure
- 5 Word Embeddings
- 6 Applications of Semantic Similarity Measures
  - Lexico-Semantic Search
  - Filename Categorization

# A live demo

■ <http://serelex.org>

ferrari

Results count: 367

1. [porsche](#)
2. [maserati](#)
3. [lamborghini](#)
4. [marque](#)
5. [bmw](#)
6. [audi](#)
7. [mercedes](#)
8. [jaguar](#)
9. [mclaren](#)
10. [dodge](#)
11. [bugatti](#)
12. [lancia](#)
13. [lotus](#)
14. [isuzu](#)
15. [tvr](#)
16. [vw](#)
17. [nissan](#)
18. [rally driving](#)
19. [fiat](#)
20. [mazda](#)

[Show all results...](#)

computational linguistics

Results count: 88

1. [psycholinguistics](#)
2. [machine learning](#)
3. [computer science](#)
4. [knowledge representation](#)
5. [cognitive science](#)
6. [artificial intelligence](#)
7. [information retrieval](#)
8. [neuoinformatics](#)
9. [natural language](#)
10. [graduate student](#)
11. [library science](#)
12. [distributed computing](#)
13. [research community](#)
14. [information processing](#)
15. [subfield](#)
16. [language acquisition](#)
17. [computer vision](#)
18. [soas](#)
19. [sociolinguistics](#)
20. [data mining](#)

[Show all results...](#)

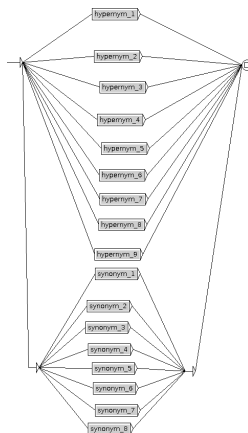
# Lexico-syntactic patterns

## ■ 18 patterns that extract **hypernyms**, **co-hyponyms** and **synonyms**

- |  |   |
|--|---|
| <input type="checkbox"/> <i>such NP as NP, NP[,]</i> and/or NP;      | <input type="checkbox"/> <i>NP, for example, NP, NP[,]</i> and/or NP; |
| <input type="checkbox"/> <i>NP such as NP, NP[,]</i> and/or NP;      | <input type="checkbox"/> <i>NP, i. e. [,]</i> NP;                     |
| <input type="checkbox"/> <i>NP, NP [,]</i> or other NP;              | <input type="checkbox"/> <i>NP (or NP);</i>                           |
| <input type="checkbox"/> <i>NP, NP [,]</i> and other NP;             | <input type="checkbox"/> <i>NP means the same as NP;</i>              |
| <input type="checkbox"/> <i>NP, including NP, NP [,]</i> and/or NP;  | <input type="checkbox"/> <i>NP, in other words[,]</i> NP;             |
| <input type="checkbox"/> <i>NP, especially NP, NP [,]</i> and/or NP; | <input type="checkbox"/> <i>NP, also known as NP;</i>                 |
| <input type="checkbox"/> <i>NP: NP, [NP,]</i> and/or NP;             | <input type="checkbox"/> <i>NP, also called NP;</i>                   |
| <input type="checkbox"/> <i>NP is DET ADJ.Superl NP;</i>             | <input type="checkbox"/> <i>NP alias NP;</i>                          |
| <input type="checkbox"/> <i>NP, e. g., NP, NP[,]</i> and/or NP;      | <input type="checkbox"/> <i>NP aka NP.</i>                            |

# Patterns are encoded as FSTs

- Finite State Transducers (FSTs)
- Open source corpus processing tool Unitex:  
<http://igm.univ-mlv.fr/~unitex/>





## A pattern encoded as an FST



- Take into account linguistic variation
- Unlike string-based patterns (Bollegala et al., 2007)

# Patterns extract concordances

- such diverse {[occupations]} as {[doctors]}, {[engineers]} and {[scientists]}[PATTERN=1]
- such {non-alcoholic [sodas]} as {[root beer]} and {[cream soda]}[PATTERN=1]
- {traditional[food]}, such as {[sandwich]}, {[burger]}, and {[fry]}[PATTERN=2]

# Corpus

Corpus Wikipedia+ukWaC:  $2.9 \cdot 10^{12}$  tokens

Name	# Documents	# Tokens	# Lemmas	Size
WaCky	2,694,815	$2,026 \cdot 10^9$	3,368,147	5.88 Gb
ukWaC	2,694,643	$0.889 \cdot 10^9$	5,469,313	11.76 Gb
WaCky + ukWaC	5,387,431	$2.915 \cdot 10^9$	7,585,989	17.64 Gb

Table 2.5: Corpora used by the PatternSim measure.

## Extracted concordances

- Wikipedia – 1.196.468
- ukWaC – 2.227.025
- WaCyclopedia+ukWaC – 3.423.493

# PatternSim Semantic Similarity

$$s_{ij} = \sqrt{p_{ij}} \cdot \frac{2 \cdot \mu_b}{b_{i*} + b_{*j}} \cdot \frac{P(c_i, c_j)}{P(c_i)P(c_j)}.$$

- $P(c_i, c_j) = \frac{e_{ij}}{\sum_{ij} e_{ij}}$  – extraction probability of the pair  $\langle c_i, c_j \rangle$ ,  $e_{ij}$  – frequency of co-occurrence of  $c_i$  and  $c_j$  in concordances  $K$
- $P(c_i) = \frac{f_i}{\sum_i f_i}$  – probability of the term  $c_i$ ,  $f_i$  – frequency of  $c_i$
- $b_{i*} = \sum_{j: e_{ij} \geq \beta} 1$  – the number of extractions for term  $c_i$  with the frequency  $\geq \beta$ ,  $\mu_b = \frac{1}{|C|} \sum_{i=1}^{|C|} b_{i*}$  – the average number of extractions per term
- $p_{ij} \in [1; 18]$  – number of distinct patterns which extracted the relation  $\langle c_i, c_j \rangle$

# Semantic Relation Ranking

- Precision is **comparable or better** w.r.t. the baselines;
- Recall is **lower** w.r.t. the baselines.

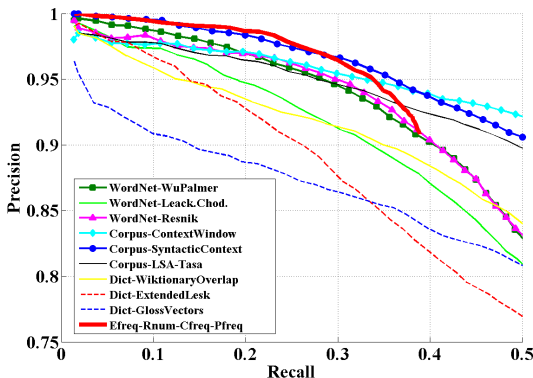
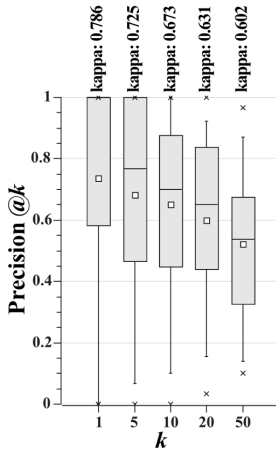


Figure: Precision-Recall graphs (the BLESS dataset).

# Semantic Relation Extraction



- $Precision@1 \approx 0.80$ ;
- “Good” coverage:

computational linguistics Search  
System finds semantically related words.  
For example, cottage cheese

Results count: 88

- 1 [psycholinguistics](#)
- 2 [machine learning](#)
- 3 [computer science](#)
- 4 [knowledge representation](#)
- 5 [cognitive science](#)
- 6 [artificial intelligence](#)
- 7 [information retrieval](#)
- 8 [neuroinformatics](#)
- 9 [natural language](#)
- 10 [graduate student](#)
- 11 [library science](#)
- 22 [distributed annotation](#)

Word to search for:

Display Options:

Key: "S:" = Show Synset (semantic) relations, "W:" = Show W

Display options for sense: (gloss) "an example sentence"

**Noun**

- **S: (n) computational linguistics** (the use of computers applications)
  - **direct hyponym / full hyponym**
    - **S: (n) machine translation, MT** (the use of co one language to another)
  - **direct hypernym / inherited hypernym / sister term**
    - **S: (n) linguistics** (the scientific study of langu

# Plan

- 1 Computational Lexical Semantics
- 2 Semantic Similarity
- 3 Pattern-Based Semantic Similarity Measure
- 4 Hybrid Semantic Similarity Measure**
- 5 Word Embeddings
- 6 Applications of Semantic Similarity Measures
  - Lexico-Semantic Search
  - Filename Categorization

# Hybrid vs Single Measures

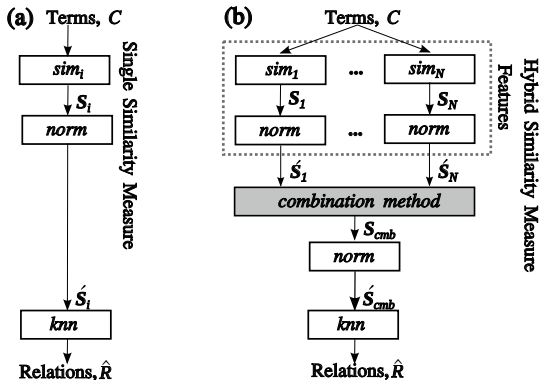


Figure: Semantic relation extractor based on:

- (a) a **single** similarity measure;
- (b) a **hybrid** similarity measure.



# 16 Features = 16 Single Similarity Measures

- 5 **network-based** measures :
  - 1 WuPalmer;
  - 2 Leacock and Chodorow;
  - 3 Resnik;
  - 4 Jiang and Conrath;
  - 5 Lin.
- 3 **web-based** measures (NGD-Yahoo/Bing/Google);
- 5 **corpus-based** measures:
  - 2 distributional (BDA, SDA)
  - 1 lexico-syntactic patterns (PatternSim)
  - 2 other co-occurrence based (LSA, NGD-Factiva)
- 3 **definition-based** measures
  - 1 ExtendedLesk;
  - 2 GlossVectors;
  - 3 DefVectors-WktWiki.

# Implementation of the baseline measures

- **Semantic Vectors:**  
<https://code.google.com/p/semanticvectors/>
- **S-Space Package:**  
<https://code.google.com/p/airhead-research/>
- **WordNet::Similarity:**  
<http://wn-similarity.sourceforge.net>
- **NLTK:** <http://nltk.googlecode.com/svn/trunk/doc/howto/wordnet.html>
- **WikiRelate!**
- **PatternSim:** <http://serelex.org>
- **Web-based metrics:**  
<http://cwl-projects.cogsci.rpi.edu/msr>
- **LSA:** <http://lsa.colorado.edu>

# Supervised Combination of Measures

## 8 Logistic Regression

- A binary **logistic regression**;
- **Positive examples** – synonyms, hyponyms, co-hyponyms;
- **Negative examples** – random relations;
- A relation  $\langle c_i, t, c_j \rangle \in R$  is represented with a **vector of pairwise similarities**:  $\mathbf{x} = (s_{ij}^1, \dots, s_{ij}^N)$ ,  $N = \overline{2, 16}$ ;
- Category  $y_{ij}$ :

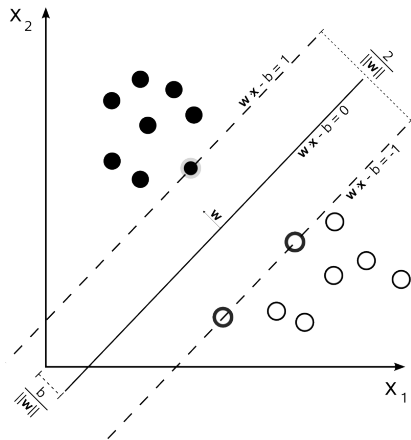
$$y_{ij} = \begin{cases} 0 & \text{if } \langle c_i, t, c_j \rangle \text{ is a random relation} \\ 1 & \text{otherwise} \end{cases}$$

- Using the model  $(w_1, \dots, w_K)$  for combination:

$$s_{ij}^{cmb} = \frac{1}{1 + e^{-z}}, z = \sum_{k=1}^K w_k s_{ij}^k + w_0.$$

# Supervised Combination Methods

## 9 SVM.



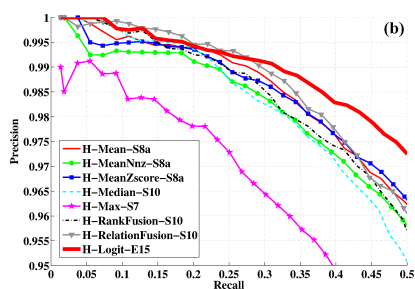
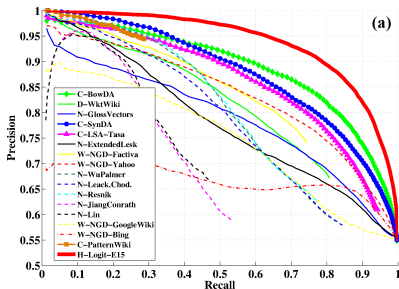
- The weights  $\mathbf{w}$  and the support vectors  $SV$ :

$$\mathbf{w} = \sum_{x_i \in SV} \alpha_i y_i \mathbf{x}_i.$$

- Using the model

$$s_{ij}^{cmb} = \mathbf{w}^T \mathbf{x} + b = \sum_{k=1}^K w_i s_{ij}^k + b.$$

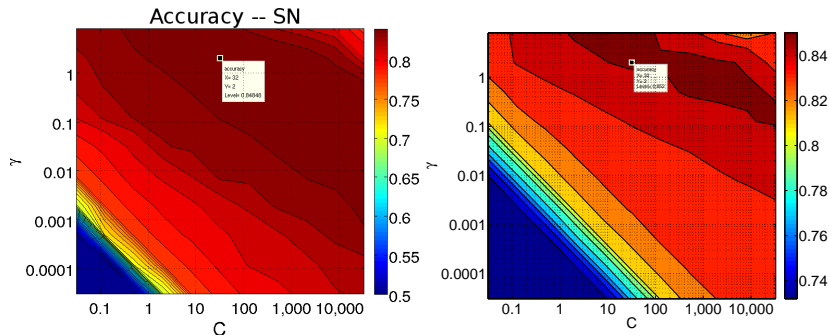
# Hybrid Similarity Measures



Precision-Recall graphs calculated on the BLESS dataset:

- (a) 16 single measures and the best hybrid measure Logit-E15;
- (b) 8 hybrid measures.

# Supervised Hybrid Similarity Measures



**Figure:** Meta-parameter optimization with the grid search of the C-SVM-radial-E15 measure.

# Plan

- 1 Computational Lexical Semantics
- 2 Semantic Similarity
- 3 Pattern-Based Semantic Similarity Measure
- 4 Hybrid Semantic Similarity Measure
- 5 Word Embeddings**
- 6 Applications of Semantic Similarity Measures
  - Lexico-Semantic Search
  - Filename Categorization

# Key facts

**Word embedding** is a dense vector representing a word obtained during a training of a neural network on a big corpus.

- Semantic similarity is cosine of word embeddings.
- The training process is known as representation learning.
- Learning methods rely on distributional hypothesis of Harris (1954), similarly to classical distributional models.
- Most popular representation learning methods:
  - Continuous Bag of Words (CBOW)
  - Skip-Gram Model
  - Global Vectors of for Word Representation (GloVe)



# Distributional Hypothesis: "You shall know the word by the company it keeps" (Firth, 1957)

he curtains open and the **stars** **shining** in on the barely  
ars and the **cold** , close **stars** " . And neither of the w  
rough the **night** with the **stars** **shining** so **brightly** , it  
made in the **light** of the **stars** . It all boils down , wr  
surely under the **bright** **stars** , thrilled by ice-white  
sun , the **seasons** of the **stars** ? Home , alone , Jay pla  
m is dazzling snow , the **stars** have **risen** **full** and **cold**  
un and the **temple** of the **stars** , driving out of the hug  
in the **dark** and now the **stars** **rise** , **full** and amber a  
bird on the **shape** of the **stars** over the **trees** in front  
But I could n't see the **stars** or the **moon** , only the  
they love the **sun** , the **stars** and the **stars** . None of  
r the **light** of the **shiny** **stars** . The splash of flowing w  
man 's first **look** at the **stars** ; various exhibits , aer  
rief information on both **stars** and **constellations**, inc

## Construct vector representations

	shining	bright	trees	dark	look
<b>stars</b>	38	45	2	27	12

# Skip-Gram Model

- Probability that word  $w$  appears in some context  $c$ :

$$P(D = 1|w, c; \theta) = \frac{1}{1 + \exp^{-V_c W_w}}$$

- Assign high probability to  $(c, w)$  pairs which appear in texts (*corp*) and low probability to the ones which cannot (*rand*):

$$\theta^* = \underset{\theta}{\operatorname{argmax}} \prod_{(c,w) \in \text{corp}} P(D = 1|w, c; \theta) \prod_{(c,w) \in \text{rand}} (1 - P(D = 1|w, c; \theta))$$

- $\theta = (V, W)$  are two matrices with columns  $V_c$  and  $W_w$  contain vectors of *dim* dimensions of the context  $c$  and the word  $w$ .
- The "embedding" of the word  $w$  is the vector  $W_w$ .

# Shared Task on Russian Semantic Similarity



- [www.dialog-21.ru/en/evaluation/2015/semantic](http://www.dialog-21.ru/en/evaluation/2015/semantic)
- **Relatedness track:** Synonyms, Hypernyms, Human Judgements about Semantic Similarity
- **Associative track:** Free associations

# Evaluation

word1	word2	sim
петух (cock)	петушок (cockerel)	0.952
побережье (coast)	берег (shore)	0.905
тип (type)	вид (kind)	0.852
миля (mile)	километр (kilometre)	0.792
чашка (cup)	посуда (tableware)	0.762
птица (bird)	петух (cock)	0.714
война (war)	войска (troops)	0.667
улица (street)	квартал (block)	0.667
...	...	...
доброволец (volunteer)	девиз (motto)	0.091
аккорд (chord)	улыбка (smile)	0.088
энергия (energy)	кризис (crisis)	0.083
бедствие (disaster)	площадь (area)	0.048
производство (production)	экипаж (crew)	0.048
мальчик (boy)	мудрец (sage)	0.042
прибыль (profit)	предупреждение (warning)	0.042
напиток (drink)	машина (car)	0.000
сахар (sugar)	подход (approach)	0.000
лес (forest)	погост (graveyard)	0.000
практика (practice)	учреждение (institution)	0.000

# Best Systems

Model ID	HJ	RT-AVEP	AE-AVEP	AE2-AVEP	Method Description
5-ae-3	0.7071	0.9185	0.9550	0.9835	Word2vec (skip-gram, window size 10, 300d vectors) on ruwac + lib.ru + ru-wiki, bigrams on the same corpus, synonym database, prefix dictionary, orthographic similarity
5-rt-3	<b>0.7625</b>	<b>0.9228</b>	0.8887	<b>0.9749</b>	Word2vec (skip-gram, window size 10, 300d vectors) on ruwac + lib.ru + ru-wiki, synonym database, prefix dictionary, orthographic similarity
1-ae-1	0.6378	<b>0.9201</b>	<b>0.9277</b>	<b>0.9849</b>	Decision trees based on n-grams (Wikipedia titles and search queries), morphological features and Word2Vec
15-rt-2	0.6537	0.9034	<b>0.9123</b>	0.9646	Word2vec trained on 150G of texts from lib.rus.ec (skip-gram, 500d vectors, window size 5, 3 iteration, min cnt 5)
16-ae-1	0.6395	0.8536	<b>0.9493</b>	0.9565	GloVe (100d vectors) on RuWac (lemmatized, normalized)
9-ae-9	<b>0.7187</b>	0.8839	0.8342	0.9517	Word2vec CBOW with window size 5 on Russian National Corpus, augmented with skip-gram model with context window size 20 on news corpus
17-rt-1	<b>0.7029</b>	0.8146	0.8945	0.9490	Distributional vector-based model, window size 5, trained on RUWAC and NRC, plmi-weighting
9-ae-6	<b>0.7044</b>	0.8625	0.8268	0.9649	Word2vec CBOW model with context window size 10 trained on web corpus
15-rt-1	0.6213	0.8472	<b>0.9120</b>	<b>0.9669</b>	Word2vec trained on 150G of texts from lib.rus.ec (skip-gram, 100d vectors, window size 10, 1 iteration, min cnt 100)
1-rt-3	0.4939	<b>0.9209</b>	0.8500	<b>0.9723</b>	Logistic regression trained on synonyms, hyponyms and hypernyms on word2vec features with AUC maximization
12-rt-3	0.4710	<b>0.9589</b>	0.5651	0.7756	Applying knowledge extracted from Wikipedia and Wiktionary for computing semantic relatedness

# Skip-Gram on 150Gb of lib.rus.ec corpus

1	бренди	лафрейг шерри шивас коньяка мадеры лафрозгй содовой бакарди коньяк коньяком кофе пунша солодового бокал виски ликера полбокала кй
2	физика	физике эйнштейна электродинамику генетика радиопизика геометрия физики релятивистская математика начертательная наука астробиолог
3	линия	вертикальная полоска точка линиями инстербургская линию прочерченная конфигурация мовиуса полукружная позиция лотербуря прольче оч
4	марс	окололунной луну мимас тригоне нептун явин сатурн плутона марса секстиле марсе венере венера планет сатурна меркурий залбатану мерку
5	омар	халим джеззар хафиз аль халид абдаллах абд муса ахмад хусейн амир абулаббас омару ибн умар исмаил хасан мохаммед хамид мухаммед а
6	книга	страница книжища трилогия брошюра переиздавалась рукопись книге издана изданная многотомная сборник брошюрка пьеса книги тетрадьчи
7	кредит	агропромбанке ссуда кредиту кредитов кредитом русфинанс кредиты безаналу заемщик заем кредитам беспроцентный беспроцентным банком
8	выгода	убыток профит выгод профиту максимизируется услуга маржа инвестиция взаимовыгодность выгодой невыгоды возмездная сумма барыш св
9	пол	линолеум ненагерты цементный паркет линолеумовый паркетный ковер наземь дионосопулос линолиум пола стукнувшуюся неоструганные по
10	еда	едой провизия закуска трапеза вкусная невкусна поесть жратва жарочка легкоусвояемая невкусной водка питье выпивка шамовка еде вкусно ж
11	авария	атомэнергетические энергосетях авиакатастрофа ошибка аварией неприятность диверсия денебойлой москваленстрое авариях протечка ходы
12	страна	супердержавя столица нация страной грузия родина мильдендо полуколонияльная цивилизация империя россия греция стране планета сверх
13	юрист	консультант журналист бреслауэр гитлин присяжный консультирующий правовед прокурор геронтолог юристом политэкономист юристов пове
14	имущество	отказанное собственность отчуждаемое награбленное выморочное секвестрованное недвижимого залогодателя недвижимостью собственность
15	попкорн	сперминт батончики мюслями гамбургеры пиццу гранолы гранолу киткат пепси арахисовым твикс арахисовое сэндвичи фрозист бургеры корн
16	владение	крещеную неотчуждаемое собственность бенефициарное владением югерами владеют пронию владения владению владею владении наслед
17	вещь	вещичку херь безделушка вещью штуку подробность ценная самая вещь бесполезная носильная своекта неизносима вещьца драгоценс
18	природа	бездременна сущность доразвилась дуальна нетварная человеческая недействительна техносфера несотворенная расклассифицированная при
19	рождаемость	демографического брачности прироста деторождаемости недетедобыча покупательная сверхсмертность репродуктивность смертности рождае
20	кровопролитие	грабежи сражение побоище братоубийственная пролитие усобной братоубийственное кровопролития предотвратить бесчинство междуусобиц
21	епископ	епископа поппон лагрийский архидиакон митрополит монтейский священник хинкомар прелат настоятель бовезский нарбонский бовезский архие
22	доктор	шатонней глайстер доктора лагенторп стифен врача чакк нейропсихиатр мистер сайболт доктором невролог бичма доктору мейленштейн обер
23	изображение	увеличенное объемное изображением телеизображение изображения изображению голографическое очертание спроецировалось видеоизоб
24	огурец	цуккини пошинкуйте салат огурцом репчатая накрошенный стручковый баклажан лимон полпучка помидоры шинкуется помидорчик огурца арбу
25	стекло	запотевающее витринное стеклом смотровое разбитое ветрозашитное бемское оконное полуопущенное стекле стекла незатонированное бэм
26	улика	кастете компра зацепка шульцевскую килбэрна зацепочка заказуха мэллинса дактилограф сгитальному куплетовские платежка улике уликах
27	отмывание	отмычки обналичка отмывка отмывке букмекерство наркобизнес невозврат офшорные зарабатывание мошенничество крышевание отмывани
28	страх	снедавший ненависть панический обуял побарывает закрывшийся атавистичное овладевавший трепет оценяющий всепоглощающий обесси
29	среда	социосистема абиотическую высококонкурентная средой суперблагоприятная сложноорганизованная реагенная макроокружение биосистема
30	сезон	фестиваль дект ташитру сезоне муссонов муссонных июль дождей муссонных триместр сноубордный осень маямар лета летом нессон

# Plan

- 1 Computational Lexical Semantics
- 2 Semantic Similarity
- 3 Pattern-Based Semantic Similarity Measure
- 4 Hybrid Semantic Similarity Measure
- 5 Word Embeddings
- 6 Applications of Semantic Similarity Measures
  - Lexico-Semantic Search
  - Filename Categorization

# Plan

- 1 Computational Lexical Semantics
- 2 Semantic Similarity
- 3 Pattern-Based Semantic Similarity Measure
- 4 Hybrid Semantic Similarity Measure
- 5 Word Embeddings
- 6 Applications of Semantic Similarity Measures
  - Lexico-Semantic Search
  - Filename Categorization



# Search for Related Words: the List and the Graph

■ <http://serelex.org>

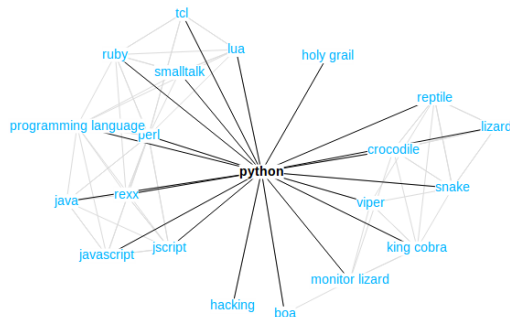


System finds semantically related words.  
For example, [lamborghini](#)

Results count: 412

- 1 [perl](#)
- 2 [ruby](#)
- 3 [programming language](#)
- 4 [tcl](#)
- 5 [monitor lizard](#)
- 6 [lua](#)
- 7 [boa](#)
- 8 [crocodile](#)
- 9 [holy grail](#)
- 10 [hacking](#)
- 11 [java](#)
- 12 [king cobra](#)
- 13 [reptile](#)
- 14 [snake](#)
- 15 [jscript](#)
- 16 [javascript](#)
- 17 [viper](#)
- 18 [smalltalk](#)
- 19 [rexx](#)
- 20 [lizard](#)

[Show all results...](#)

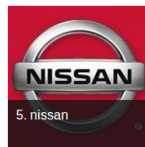


## Lexico-Semantic Search

# Search for Related Words: the Images

citroen

System finds semantically related words.  
For example, [linux](#)



# Plan

- 1 Computational Lexical Semantics
- 2 Semantic Similarity
- 3 Pattern-Based Semantic Similarity Measure
- 4 Hybrid Semantic Similarity Measure
- 5 Word Embeddings
- 6 Applications of Semantic Similarity Measures
  - Lexico-Semantic Search
  - **Filename Categorization**

# Short text classification with Vocabulary Projection

**18XGirls Yulia**



{}



{ekaterina, sonya, daughter}

**HD Widget Android**



{gadget, menu, button}

**Plan9 Unix**



{}



{linux macintosh, solaris, freebsd, BSD, window, platform, novell, sco}

**Sexart 10 04 05 Nedda A Presenting**



{}



{adina, gilda, mimi, juliette, marguerite, heroine, lucia, lui, role}

# Evaluation of the Vocabulary Projection

Training Dataset	Test Dataset	Accuracy	Accuracy (voc. projection)
Gallery (train)	Gallery	96.41	96.83 (+0.42)
PirateBay Title+Desc+Tags	PirateBay Title+Desc+Tags	98.92	98.86 (−0.06)
PirateBay Title+Tags	PirateBay Title+Tags	97.73	97.63 (−0.10)
Gallery	PirateBay Title+Desc+Tags	90.57	91.48 (+0.91)
Gallery	PirateBay Title+Tags	84.23	88.89 (+4.66)
PirateBay Title+Desc+Tags	Gallery	88.83	89.04 (+0.21)
PirateBay Title+Tags	Gallery	91.16	91.30 (+0.14)

**Table:** Performance of an C-SVM linear classifier (10-fold cross validation).