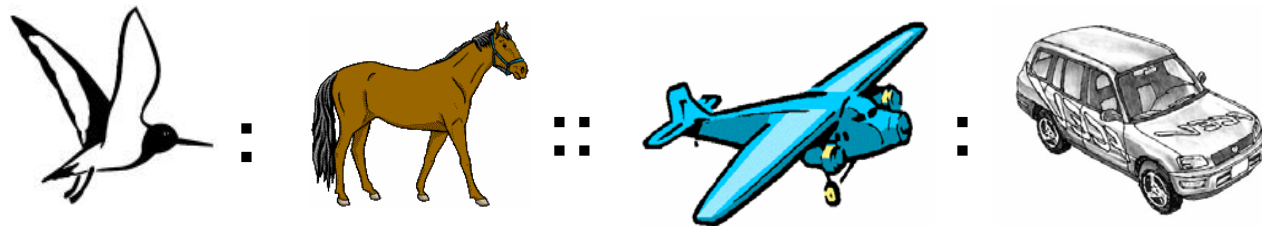


# *Corpus-Based Learning of Analogies and Semantic Relations*



**Peter Turney**

**National Research Council of Canada**  
**(work done with Michael Littman, Rutgers)**

**February 2004**

# ***Outline***



- **introduction**
- **motivation and applications**
  - ubiquity of metaphor
  - classifying semantic relations
- **related work**
- **solving analogy problems**
  - Vector Space Model
  - experiments
- **noun-modifier semantic relations**
  - 30 classes of semantic relations
  - experiments
- **future work**
- **conclusion**

# ***Introduction***



## ***Introduction***

- **verbal analogy has form  $A:B::C:D$** 
  - A is to B as C is to D
- **example**
  - mason:stone::carpenter:wood
  - mason is to stone as carpenter is to wood
- **analogies have been studied at least since Aristotle**
  - Nicomachean Ethics, Book V, Section 3
  - but still not well understood
- **idea**
  - use SAT verbal analogy questions to guide research on computational approaches to analogies

## ***SAT Analogy Question***

---

**Stem:** mason:stone

- 
- Choices:**
- (a) teacher:chalk
  - (b) carpenter:wood
  - (c) soldier:gun
  - (d) photograph:camera
  - (e) book:word

---

**Solution:** (b) carpenter:wood

---

## ***SAT Analogy Question***

- answering SAT analogy questions requires knowledge of semantic relations between words
  - semantic relations are often implicit
- noun-modifier semantic relations always implicit
  - “sleeping dog”
    - noun: **dog**, modifier: **sleeping**
    - the dog is *in a state of* sleeping
  - “sleeping pill”
    - the pill *causes* sleep
  - “sleeping area”
    - an area *for* sleeping in
- want to automatically identify semantic relations

## ***SAT Analogy Question***

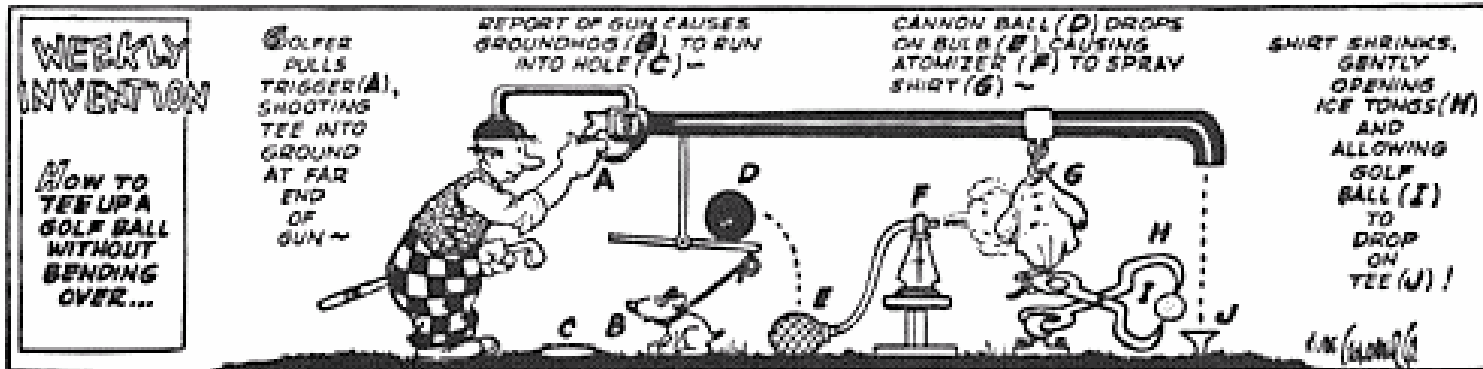
	<b>word pair</b>	<b>semantic relation</b>
<b>Stem:</b>	<b>mason:stone</b>	<b>builds with</b>
<b>Choices:</b>	(a) <b>teacher:chalk</b>	<b>writes with</b>
	(b) <b>carpenter:wood</b>	<b>builds with</b>
	(c) <b>soldier:gun</b>	<b>fights with a</b>
	(d) <b>photograph:camera</b>	<b>is produced by a</b>
	(e) <b>book:word</b>	<b>contains a</b>
<b>Solution:</b>	(b) <b>carpenter:wood</b>	<b>builds with</b>



## ***Our Approach to Verbal Analogies***

- **corpus-based learning**
  - use Web as very large corpus of text
- **vector of phrase frequencies to characterize semantic relation**
  - statistical “signature” of relationship
- **measure similarity of semantic relations by cosine of angle between vectors**
  - Vector Space Model (VSM) of Information Retrieval
  - used by all major search engines to rank hits

## Motivation and Applications



RUBE GOLDBERG (TM) RGI 134

## ***Ubiquity of Metaphor***

- **metaphorical language is very common**
- **metaphors can be understood as verbal analogies**
  - **He *shot down* all of my *arguments*.**
    - **aircraft:shoot\_down::argument:criticize**
  - **You need to *budget* your *time*.**
    - **money:budget::time:schedule**
  - **I *gave* you that *idea*.**
    - **object:give::idea:communicate**
  - ***Life* has *cheated* me.**
    - **charlatan:cheat::life:disappoint**
  - **The Michelson-Morely *experiment* *gave birth* to a new physical theory.**
    - **mother:give\_birth::experiment:initiate**

## ***Evolution of Language***

- language evolution is often metaphorical
- etymology can often be understood as verbal analogy
  - ***Bias***: a partiality that prevents objective consideration of an issue or situation. From the French *biais*, a slant, slope; hence, inclination to one side.
    - bias:person::slant:line
  - ***Disseminate***: cause to become widely known. From the Latin *disseminare*, to scatter seed.
    - disseminate:information::scatter:seed
  - ***Insult***: treat, mention, or speak to rudely. From the Latin *insultare*, to leap upon.
    - insult:character::leap\_upon:body

## ***Noun-Modifier Semantic Relations***

- **algorithm for SAT verbal analogies could classify noun-modifier semantic relations**
- **nearest neighbour supervised learning**
  - **given set of noun-modifier pairs, hand-labeled with semantic relations**
    - **training data**
  - **given new noun-modifier pair, unknown semantic relation**
    - **testing data**
  - **classify by looking for *most analogous* noun-modifier pair in training set**
    - **most analogous = nearest neighbour**

## ***Noun-Modifier Semantic Relations***

- **applications for noun-modifier classification**
  - machine translation
    - translate “electron microscope” to another language
    - is semantic relation *purpose* or *instrument*?
  - information extraction
    - extract parties involved from news about wars
    - in “cigarette war”, relation is *topic*, not *agent*
  - word sense disambiguation
    - “plant” might be industrial plant or living plant
    - helpful to know that relation in “plant food” is *beneficiary*, not *source*

## ***Related Work***



## ***Metaphor and Analogy***

- **French (2002)**
  - general survey of literature on analogy and metaphor
  - all work in survey involved hand-built knowledge-bases
  - no prior work in machine learning cited
- **Dolan (1995)**
  - extracting knowledge automatically from a dictionary
  - discovered “conventional” metaphors
  - no systematic evaluation
- **Marx et al. (2002)**
  - clustering algorithm; could discover analogies between clusters of words, but not between individual words
- **Veale (2003)**
  - extracts analogies of form adjective:noun::adjective:noun from WordNet



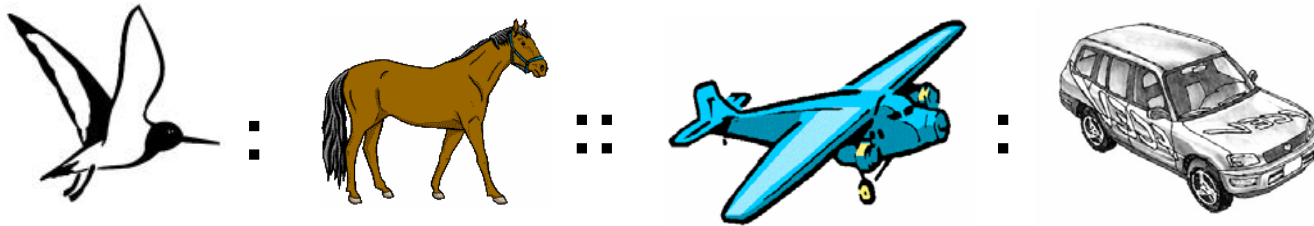
## ***Vector Space Model***

- **VSM first developed in Information Retrieval**
  - similarity of query to document measure by cosine of angle between query vector and document vector
  - Salton and McGill (1983), Salton (1989)
- **cosine also used to measure word similarity**
  - Lesk (1969), Ruge (1992), Pantel and Lin (2002)
- **our use of cosine for analogies is new**
  - we measure similarity of word pairs (similarity semantic relations), not individual words (similarity of concepts)

## ***Noun-Modifier Semantic Relations***

- **Nastase and Szpakowicz (2003)**
  - use supervised learning to classify 600 noun-modifier pairs
  - same data as we use here
  - algorithm uses features from WordNet, rather than corpus-based features
  - still in “exploratory” phase of research
- **Rosario and Hearst (2001) and Rosario et al. (2002)**
  - semantic relations in medical text
    - domain-specific
  - use features from MeSH (Medical Subject Headings) and UMLS (Unified Medical Language System)
- **no prior corpus-based approach**

## ***Solving Analogy Problems***



## ***Solving Analogy Problems***

- **assign scores to candidate analogies  $A:B::C:D$** 
  - **for multiple-choice questions, guess highest scoring choice**
- **quality of analogy depends on degree of similarity between semantic relation  $R_1$  of  $A:B$  and semantic relation  $R_2$  of  $C:D$** 
  - **difficulty is that  $R_1$  and  $R_2$  are implicit**
- **attempt to learn  $R_1$  and  $R_2$  using unsupervised learning from a very large corpus**

## Vector Space Model

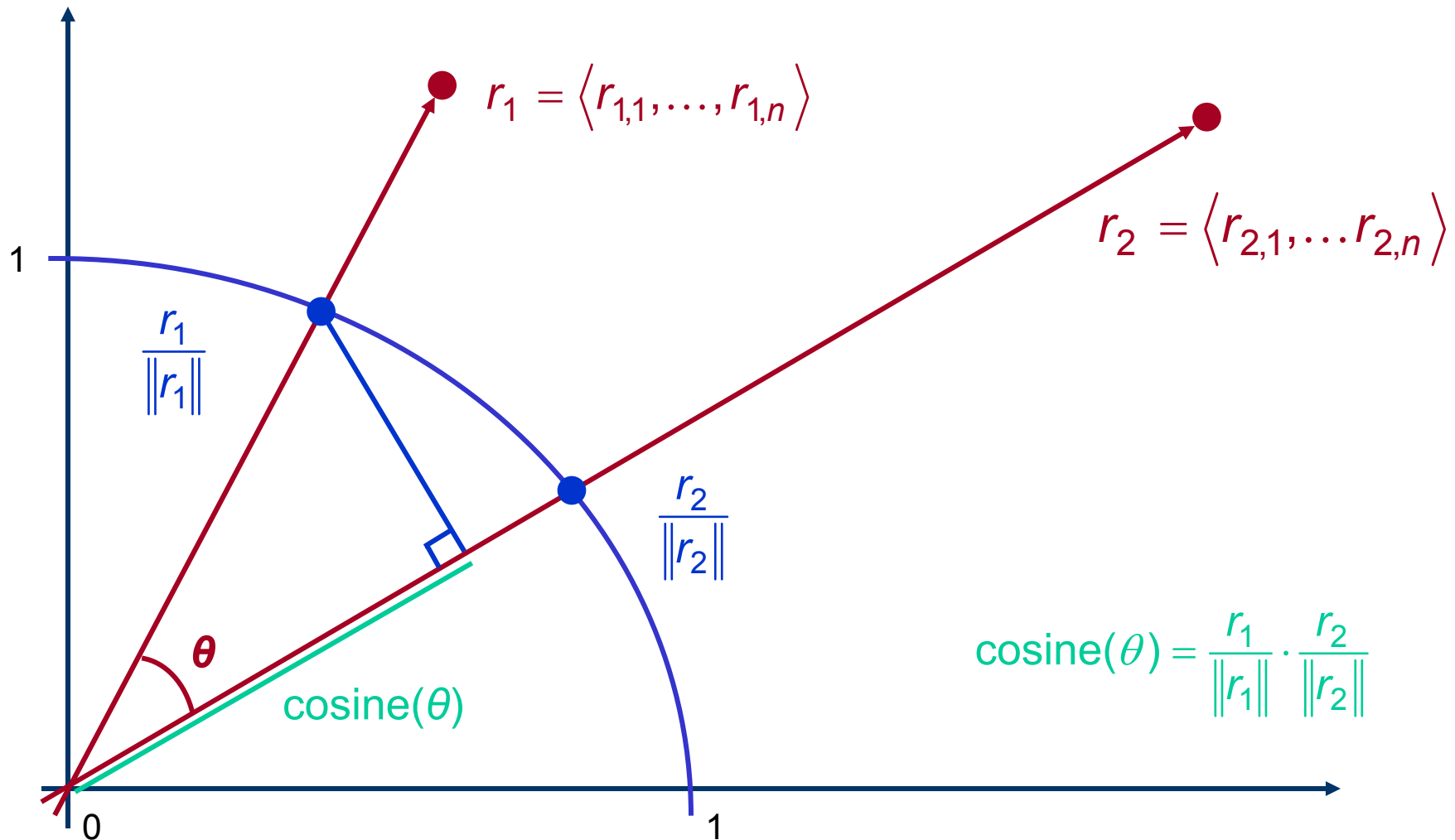
- create vectors,  $r_1$  and  $r_2$ , that represent features of  $R_1$  and  $R_2$

$$r_1 = \langle r_{1,1}, \dots, r_{1,n} \rangle \quad r_2 = \langle r_{2,1}, \dots, r_{2,n} \rangle$$

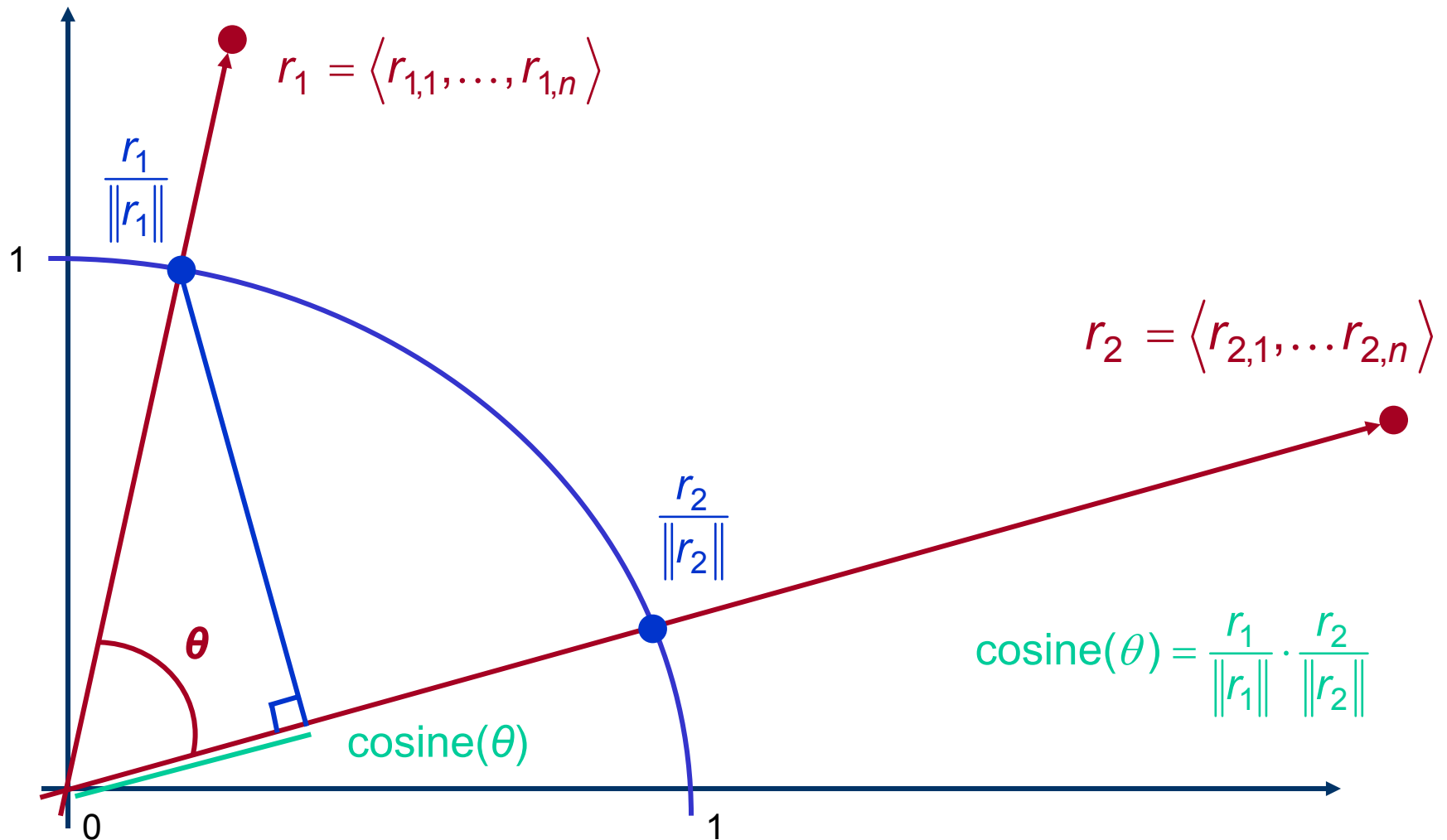
- measure the similarity of  $R_1$  and  $R_2$  by the cosine of the angle  $\theta$  between  $r_1$  and  $r_2$

$$\text{cosine}(\theta) = \frac{\sum_{i=1}^n r_{1,i} r_{2,i}}{\sqrt{\sum_{i=1}^n (r_{1,i})^2} \cdot \sqrt{\sum_{i=1}^n (r_{2,i})^2}} = \frac{r_1 \cdot r_2}{\sqrt{r_1 \cdot r_1} \cdot \sqrt{r_2 \cdot r_2}} = \frac{r_1 \cdot r_2}{\|r_1\| \cdot \|r_2\|}$$

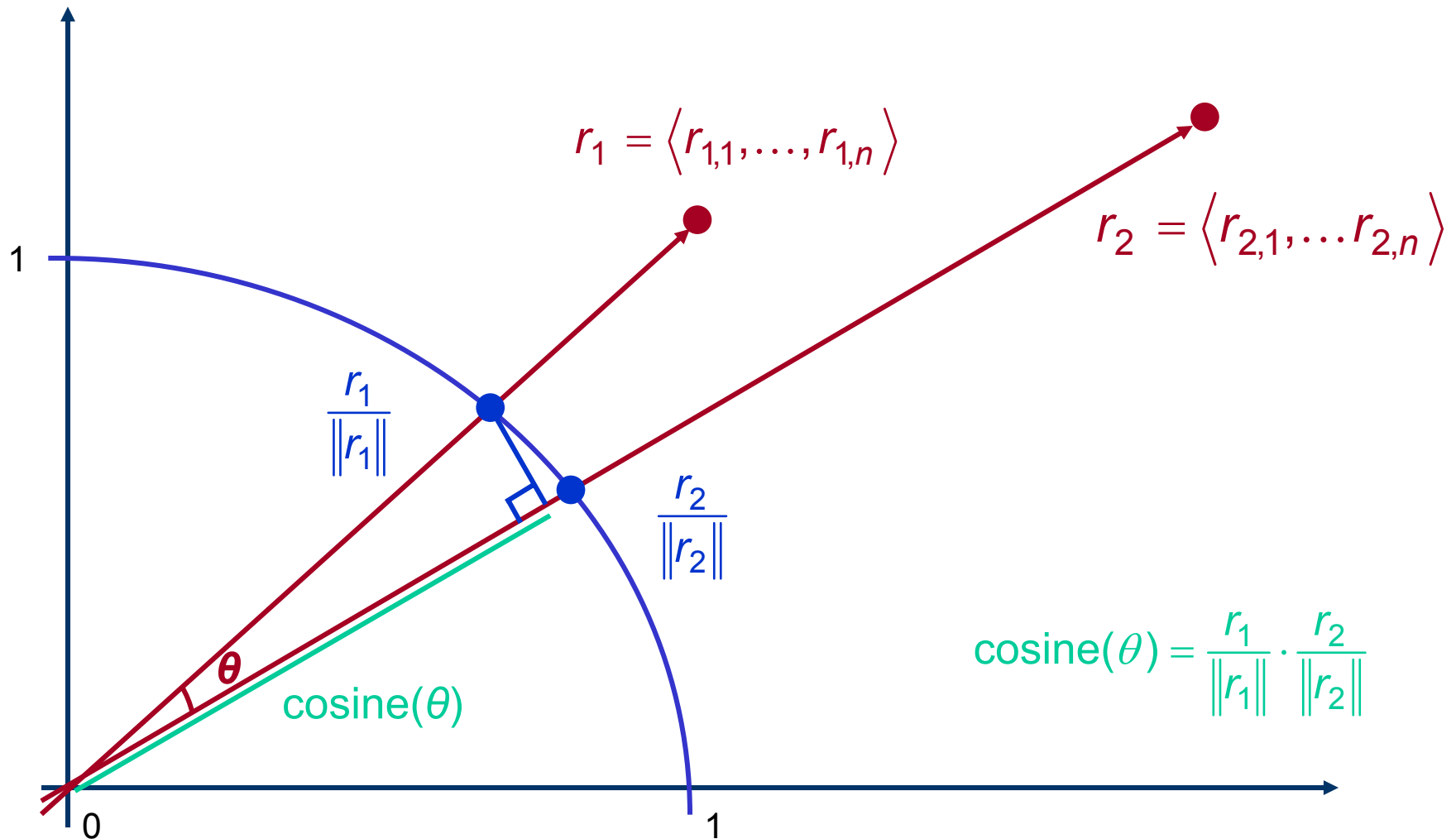
## Vector Space Model



# Vector Space Model



# Vector Space Model





## ***Generating Vectors***

- **given word pair X:Y**
  - use joining term J to make phrases “X J Y” and “Y J X”
  - search web for frequencies of phrases “X J Y” and “Y J X”
  - N joining terms results in vector of 2N numbers
  - take logarithm of frequencies
- **example**
  - word pair: “mason:stone”
  - joining terms: “with”, “to”, “for”, “of”, ...
  - search AltaVista: “mason with stone”, “stone with mason”, ...
  - note number of hits (matching web pages)
  - vector of logs of hits
  - 64 joining terms; 128 elements in vector

## Algorithm

- given candidate analogy A:B::C:D
  - traffic:street::water:riverbed
- generate vector for A:B and vector for C:D
  - $r_1$  for A:B and  $r_2$  for C:D
- calculate cosine of angle between vectors
  - $\text{cosine}(r_1, r_2)$
- cosine is score for candidate analogy
  - $\text{score}(\text{traffic:street::water:riverbed}) = \text{cosine}(r_1, r_2)$
- similar pattern of frequencies implies small angle between vectors, implies large cosine
  - note importance of vector length normalization
  - frequent words result in longer vectors
  - we care about *direction*, not *length*

## ***Algorithm***

- **example**
  - **traffic:street::water:riverbed**

<b>query</b>	<b>traffic:street</b>	<b>water:riverbed</b>
<b>“X in the Y”</b>	<b>615 hits</b>	<b>91 hits</b>
<b>“Y on X”</b>	<b>6 hits</b>	<b>0 hits</b>
<b>“Y with X”</b>	<b>478 hits</b>	<b>11 hits</b>
<b>“X from the Y”</b>	<b>136 hits</b>	<b>14 hits</b>
<b>“X when Y”</b>	<b>2 hits</b>	<b>0 hits</b>
<b>total</b>	<b>1237 hits</b>	<b>116 hits</b>

## ***Algorithm***

- **example**
  - **traffic:street::water:riverbed**
  - **one of the SAT questions**

<b>Stem pair:</b>	<b>traffic:street</b>	<b>Cosine</b>
<b>Choices:</b>	(a) <b>ship:gangplank</b>	<b>0.31874</b>
	(b) <b>crop:harvest</b>	<b>0.57234</b>
	(c) <b>car:garage</b>	<b>0.68757</b>
	(d) <b>pedestrians:feet</b>	<b>0.49725</b>
	(e) <b>water:riverbed</b>	<b>0.69265</b>

## ***Evaluation Metrics***

- **374 SAT analogy questions, 5 choices each**

$$\text{precision} = \frac{\text{number of correct guesses}}{\text{total number of guesses made}}$$

$$\text{recall} = \frac{\text{number of correct guesses}}{\text{maximum possible number correct}}$$

$$F = \frac{2 \times \text{precision} \times \text{recall}}{\text{precision} + \text{recall}}$$

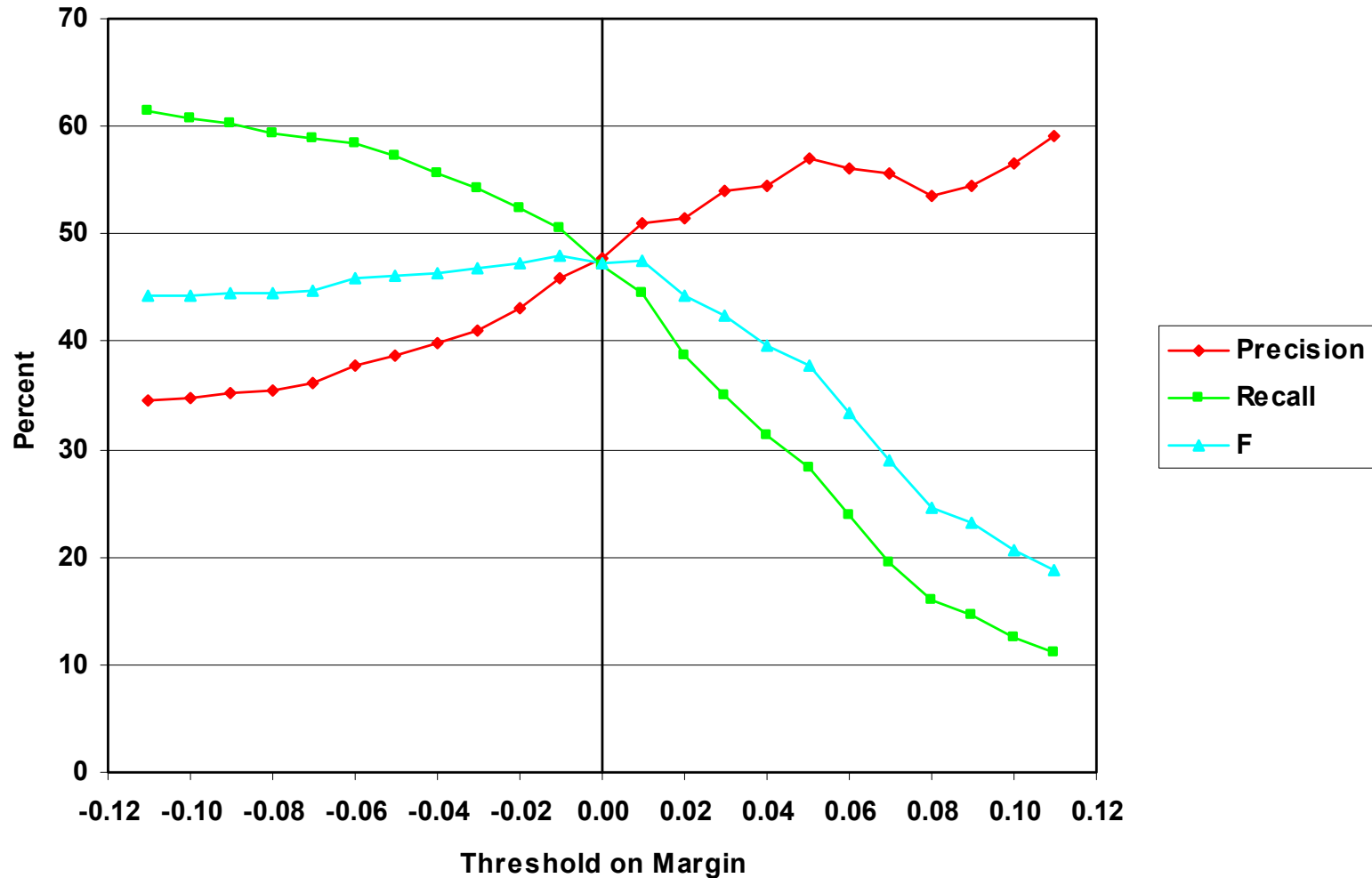
## ***Results on 374 SAT Questions***

	<b>Number</b>	<b>Percent</b>
<b>Correct</b>	<b>176</b>	<b>47.1%</b>
<b>Incorrect</b>	<b>193</b>	<b>51.6%</b>
<b>Skipped</b>	<b>5</b>	<b>1.3%</b>
<b>Total</b>	<b>374</b>	<b>100.0%</b>
<b>Precision</b>	<b>176 / 369</b>	<b>47.7%</b>
<b>Recall</b>	<b>176 / 374</b>	<b>47.1%</b>
<b>F</b>		<b>47.4%</b>

## ***Human Performance on SAT***

<b>Note</b>	<b>Percent correct (no skipping)</b>	<b>SAT I raw score verbal</b>	<b>SAT I scaled score verbal</b>	<b>Percentile rank</b>
	<b>100%</b>	<b>78</b>	<b>800</b>	<b>100.0</b>
	<b>92%</b>	<b>70</b>	<b>740</b>	<b>98.0</b>
	<b>82%</b>	<b>60</b>	<b>645</b>	<b>88.5</b>
	<b>71%</b>	<b>50</b>	<b>580</b>	<b>74.0</b>
<b>College-bound mean</b>	<b>57%</b>	<b>36</b>	<b>504</b>	<b>48.0</b>
<b>VSM algorithm</b>	<b>47%</b>	<b>26</b>	<b>445</b>	<b>29.0</b>
	<b>41%</b>	<b>20</b>	<b>410</b>	<b>18.5</b>
	<b>30%</b>	<b>10</b>	<b>335</b>	<b>5.5</b>
<b>Random guessing</b>	<b>20%</b>	<b>0</b>	<b>225</b>	<b>0.5</b>

# Precision versus Recall





## Generating Analogies

- SAT test is about *recognizing* analogies
- what about *generating* analogies?
- maybe reduce *generation* to *recognition*?
  - randomly create candidate word pairs
  - see which pair is most similar to given stem pair
- step towards generation
  - 374 questions, 5 skipped = 369 not skipped
  - merge all 369 correct answer pairs
  - for each of 369 stem pairs, select correct answer pair from set of 369 choices
  - how often is correct choice among top 10 choices?
  - random guessing:  $10/369 = 2.7\%$

## ***Generating Analogies***

<b>Rank</b>	<b>Matches</b>	<b>Matches</b>	<b>Cumulative</b>	<b>Cumulative</b>
<b>#</b>	<b>#</b>	<b>%</b>	<b>#</b>	<b>%</b>
<b>1</b>	<b>31</b>	<b>8.4%</b>	<b>31</b>	<b>8.4%</b>
<b>2</b>	<b>19</b>	<b>5.1%</b>	<b>50</b>	<b>13.6%</b>
<b>3</b>	<b>13</b>	<b>3.5%</b>	<b>63</b>	<b>17.1%</b>
<b>4</b>	<b>11</b>	<b>3.0%</b>	<b>74</b>	<b>20.1%</b>
<b>5</b>	<b>6</b>	<b>1.6%</b>	<b>80</b>	<b>21.7%</b>
<b>6</b>	<b>7</b>	<b>1.9%</b>	<b>87</b>	<b>23.6%</b>
<b>7</b>	<b>9</b>	<b>2.4%</b>	<b>96</b>	<b>26.0%</b>
<b>8</b>	<b>5</b>	<b>1.4%</b>	<b>101</b>	<b>27.4%</b>
<b>9</b>	<b>5</b>	<b>1.4%</b>	<b>106</b>	<b>28.7%</b>
<b>10</b>	<b>3</b>	<b>0.8%</b>	<b>109</b>	<b>29.5%</b>

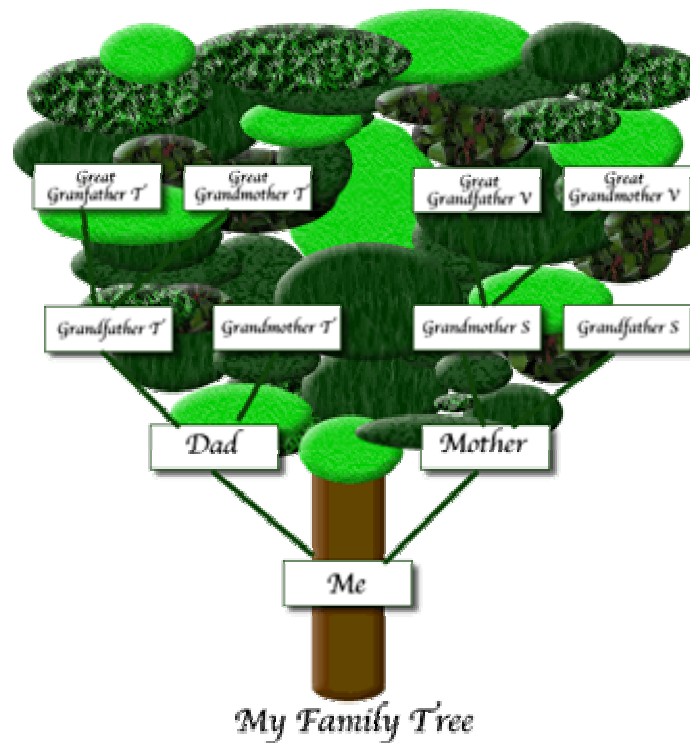
## ***Generating Analogies***

<b>Rank</b>	<b>Word pair</b>	<b>Cosine</b>	<b>Question #</b>
<b>Stem</b>	<b>tourniquet:bleeding</b>		<b>46</b>
<b>1</b>	<b>antidote:poisoning</b>	<b>0.7540</b>	<b>308</b>
<b>2</b>	<b>belligerent:fight</b>	<b>0.7482</b>	<b>84</b>
<b>3</b>	<b>chair:furniture</b>	<b>0.7481</b>	<b>107</b>
<b>4</b>	<b>mural:wall</b>	<b>0.7430</b>	<b>302</b>
<b>5</b>	<b>reciprocate:favor</b>	<b>0.7429</b>	<b>151</b>
<b>6</b>	<b>menu:diner</b>	<b>0.7421</b>	<b>284</b>
<b>7</b>	<b>assurance:uncertainty</b>	<b>0.7287</b>	<b>8</b>
<b>8</b>	<b>beagle:dog</b>	<b>0.7210</b>	<b>19</b>
<b>9</b>	<b>canvas:painting</b>	<b>0.7205</b>	<b>5</b>
<b>10</b>	<b>ewe:sheep</b>	<b>0.7148</b>	<b>261</b>

## ***Execution Time***

- **experiments presented here required 287,232 queries to AltaVista**
  - 374 analogy questions
  - × 6 word pairs per question
  - × 128 queries per word pair
  - = 287,232 queries
- **as courtesy to AltaVista, inserted a five second delay between each query**
  - processing 287,232 queries took about seventeen days

## ***Noun-Modifier Semantic Relations***



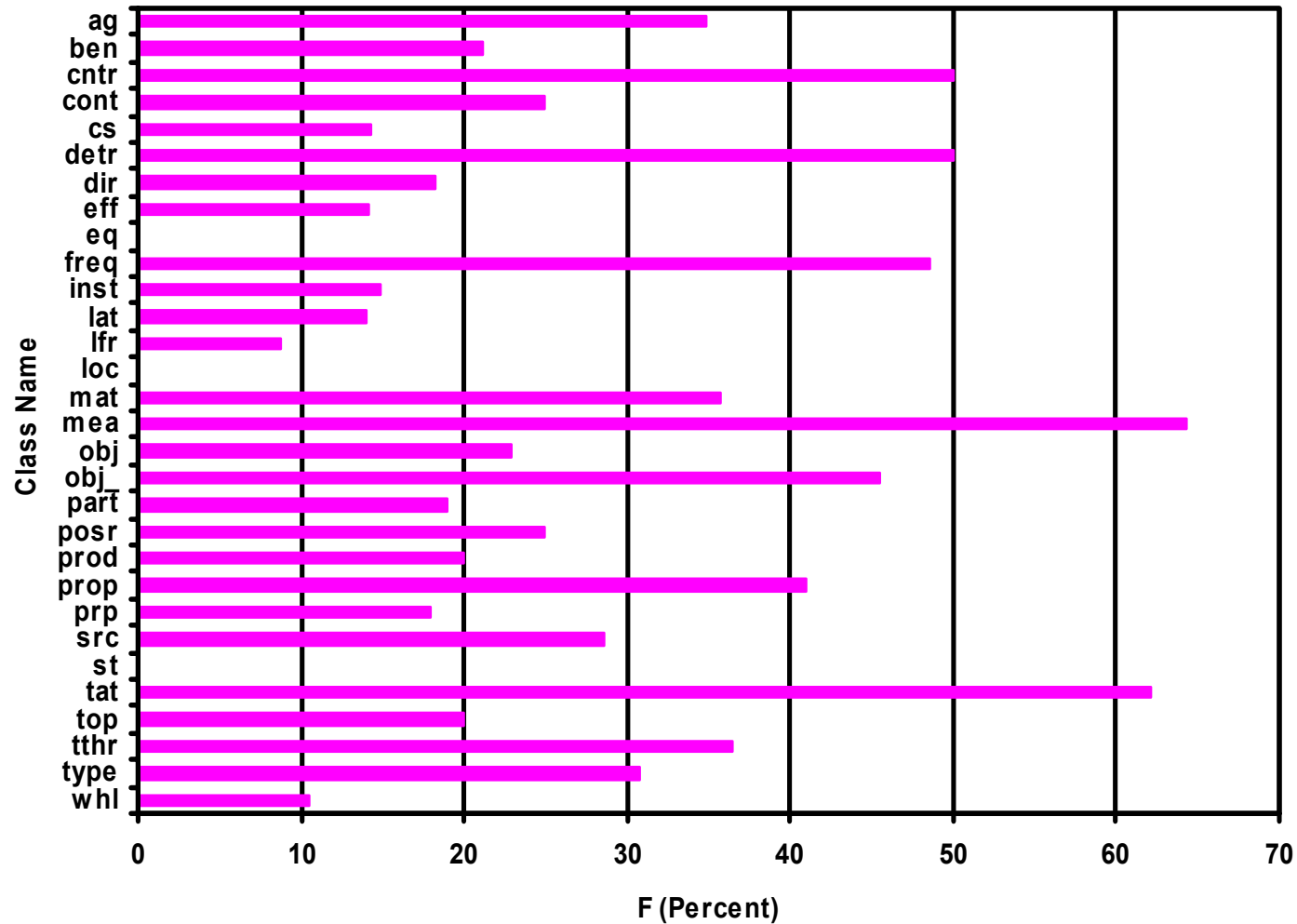
## ***Noun-Modifier Semantic Relations***

- **nearest neighbour supervised learning**
  - given set of noun-modifier pairs, hand-labeled with semantic relations (Nastase and Szpakowicz, 2003)
    - training data
  - given new noun-modifier pair, unknown semantic relation
    - testing data
  - classify by looking for *most analogous* noun-modifier pair in training set
    - most analogous = nearest neighbour
- **nearest neighbour = cosine**
  - cosine(training pair, testing pair)
  - vector of 128 elements, same joining terms as before

## ***30 Semantic Relations***

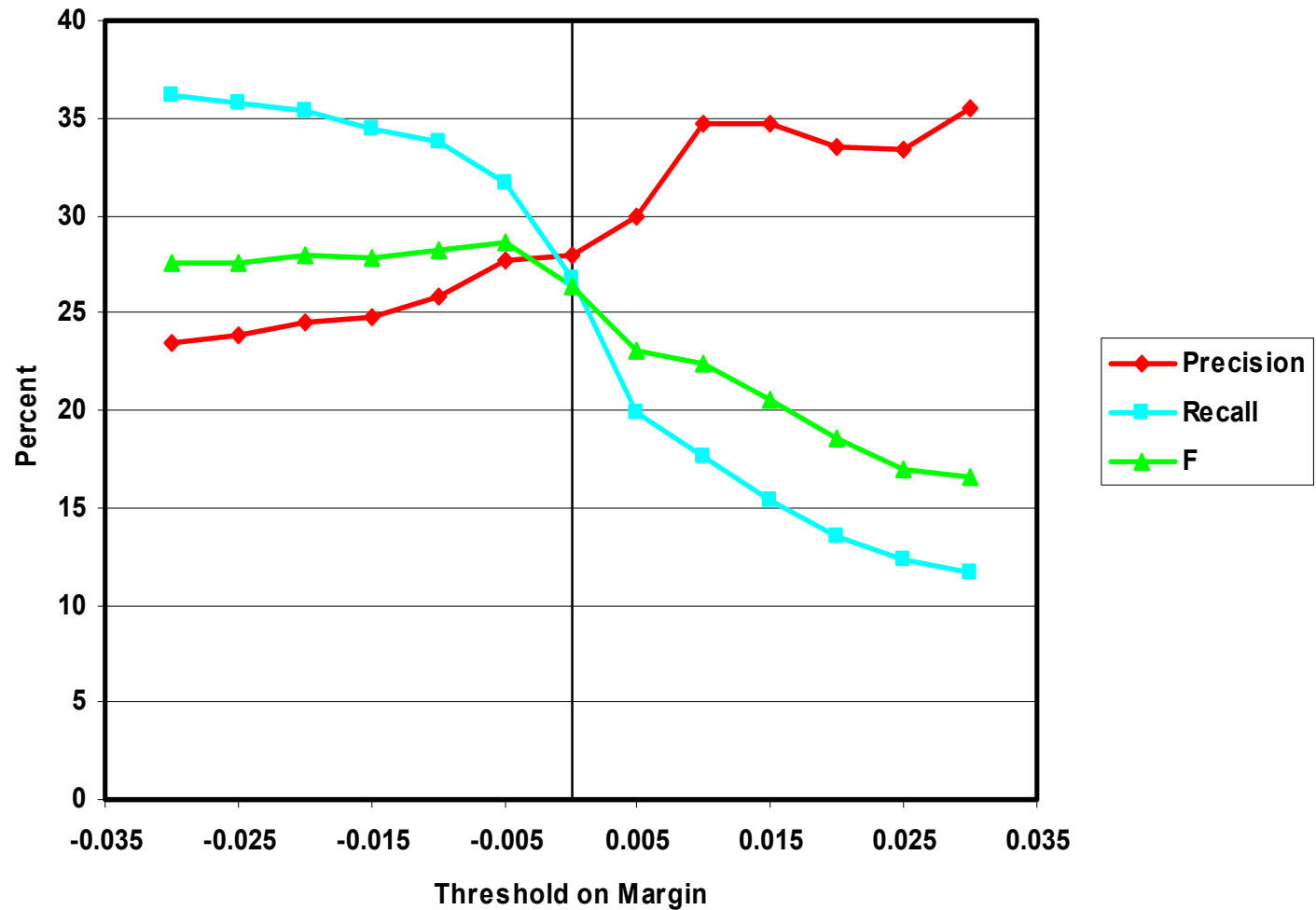
<b>Relation</b>	<b>Example</b>	<b>Relation</b>	<b>Example</b>		
<b>1</b>	<b>cause</b>	<b>flu virus</b>	<b>16</b>	<b>object property</b>	<b>sunken ship</b>
<b>2</b>	<b>effect</b>	<b>exam anxiety</b>	<b>17</b>	<b>part</b>	<b>printer tray</b>
<b>3</b>	<b>purpose</b>	<b>concert hall</b>	<b>18</b>	<b>possessor</b>	<b>national debt</b>
<b>4</b>	<b>detraction</b>	<b>headache pill</b>	<b>19</b>	<b>property</b>	<b>blue book</b>
<b>5</b>	<b>frequency</b>	<b>daily exercise</b>	<b>20</b>	<b>product</b>	<b>plum tree</b>
<b>6</b>	<b>time at</b>	<b>morning exercise</b>	<b>21</b>	<b>source</b>	<b>olive oil</b>
<b>7</b>	<b>time through</b>	<b>six-hour meeting</b>	<b>22</b>	<b>stative</b>	<b>sleeping dog</b>
<b>8</b>	<b>direction</b>	<b>outgoing mail</b>	<b>23</b>	<b>whole</b>	<b>daisy chain</b>
<b>9</b>	<b>location</b>	<b>home town</b>	<b>24</b>	<b>container</b>	<b>film music</b>
<b>10</b>	<b>location at</b>	<b>desert storm</b>	<b>25</b>	<b>content</b>	<b>apple cake</b>
<b>11</b>	<b>location from</b>	<b>foreign capital</b>	<b>26</b>	<b>equative</b>	<b>player coach</b>
<b>12</b>	<b>agent</b>	<b>student protest</b>	<b>27</b>	<b>material</b>	<b>brick house</b>
<b>13</b>	<b>beneficiary</b>	<b>student discount</b>	<b>28</b>	<b>measure</b>	<b>expensive book</b>
<b>14</b>	<b>instrument</b>	<b>laser printer</b>	<b>29</b>	<b>topic</b>	<b>weather report</b>
<b>15</b>	<b>object</b>	<b>metal separator</b>	<b>30</b>	<b>type</b>	<b>oak tree</b>

## *F for the 30 Classes*





# Precision versus Recall



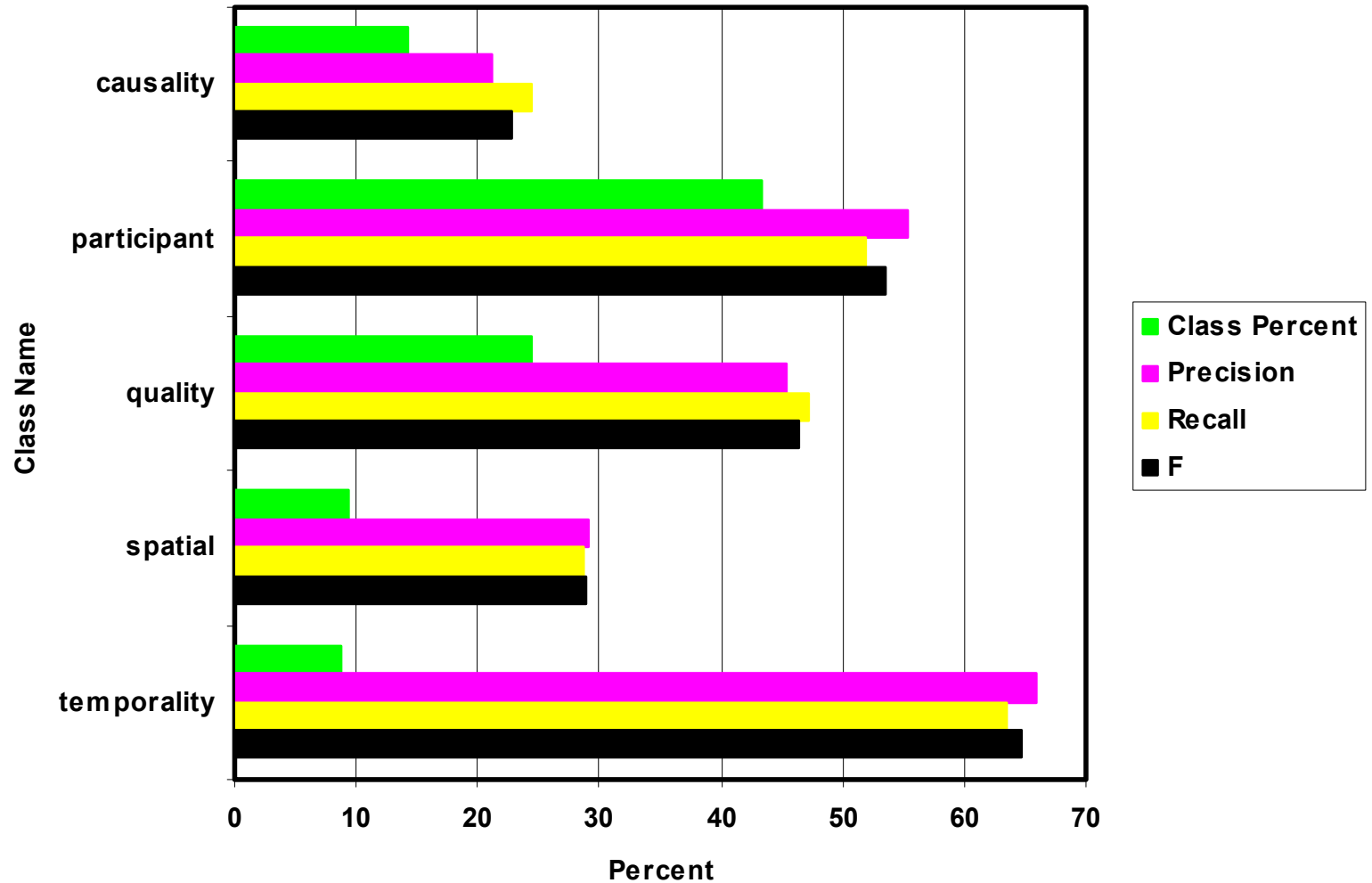
## **30 Semantic Relations**

- **F when precision and recall are balanced**
  - 26.5%
- **F for random guessing**
  - 3.3%
- **much better than random guessing**
  - but still much room for improvement
- **30 classes is hard**
  - too many possibilities for confusing classes
- **try 5 classes instead**
  - group classes together

# 30 Semantic Relations

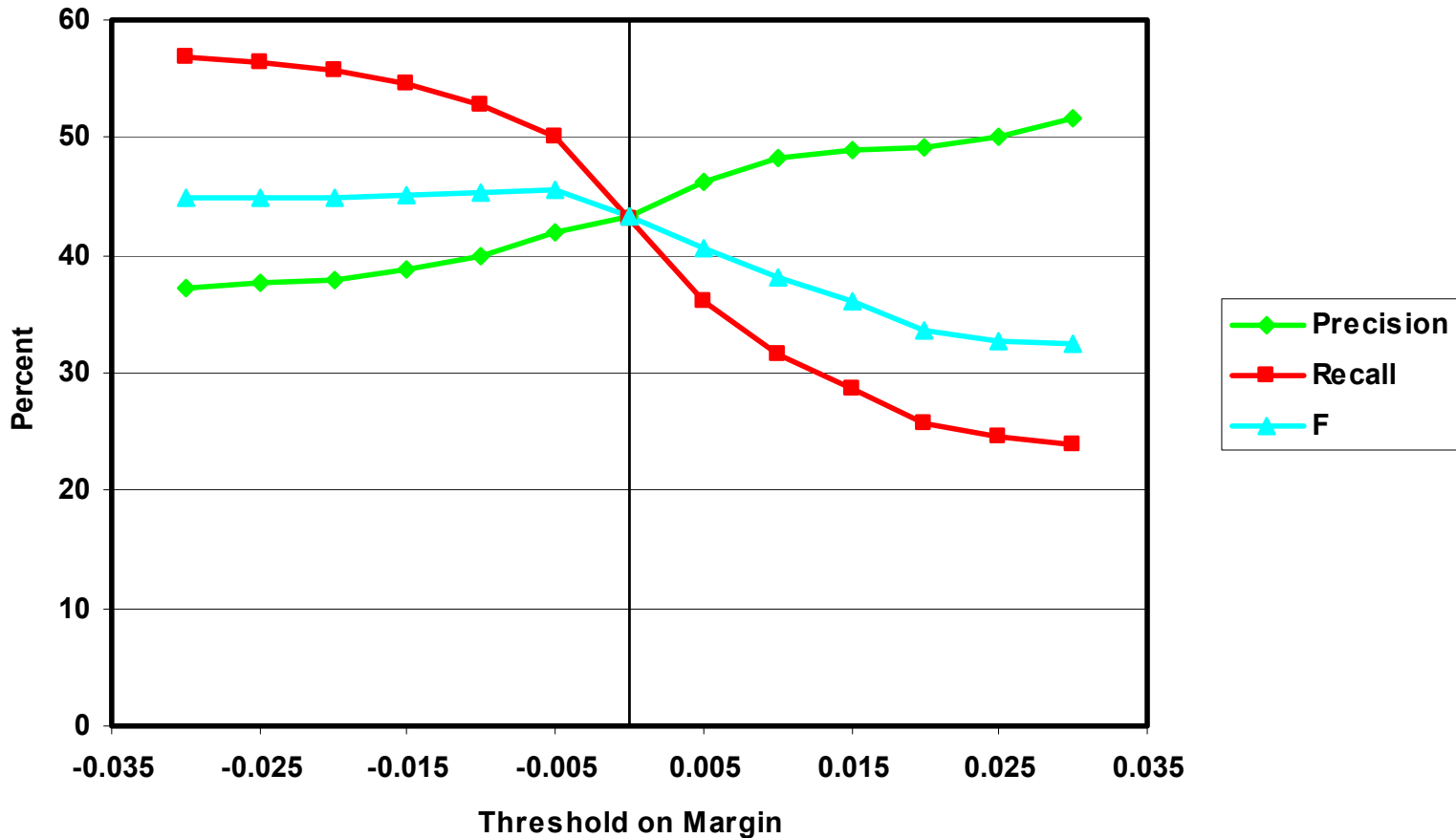
Relation	Example	Relation	Example
1	cause flu virus	16	object property sunken ship
2	effect <b>Causality</b>	17	part printer tray
3	purpose	18	possessor national debt
4	detraction headache pill	19	property blue book
5	frequency <b>Temporality</b>	20	product <b>Participant</b>
6	time at	21	source olive oil
7	time through six-hour meeting	22	stative sleeping dog
8	direction outgoing mail	23	whole daisy chain
9	location <b>Spatial</b>	24	container film music
10	location at m	25	content apple cake
11	location from foreign capital	26	equative ach
12	agent student protest	27	material <b>Quality</b> ise
13	beneficiary <b>Participant</b>	28	measure expensive book
14	instrument	29	topic weather report
15	object metal separator	30	type oak tree

## *F for the 5 Classes*



# Precision versus Recall

Precision and Recall with Varying Thresholds for 5 Classes



## ***5 Semantic Relations***

- **F when precision and recall are balanced**
  - **43.2%**
- **F for random guessing**
  - **20.0%**
- **better than random guessing**
- **better than 30 classes**
  - **26.5%**
  - **but still room for improvement**

## ***Execution Time***

- **experiments presented here required 76,800 queries to AltaVista**
  - 600 word pairs
  - × 128 queries per word pair
  - = 76,800 queries
- **as courtesy to AltaVista, inserted a five second delay between each query**
  - processing 76,800 queries took about five days

## ***Future Work***





## ***Future Work***

- **much room for experimentation in choice of joining terms**
  - but experiments take long time to run
- **progress in hardware will allow searching local database of AltaVista size**
  - recently acquired 16 CPU Beowulf Cluster
  - terabyte corpus from University of Waterloo
- **variations on VSM**
  - LSA, GVSM, term weighting schemes, ...

## ***Conclusion***



## ***Conclusion***

- **analogy and metaphor play a central role in human cognition and language**
  - Lakoff and Johnson (1980), Hofstadter *et al.* (1995), French (2002)
- **SAT-style analogy questions are a simple but powerful and objective tool for investigating these phenomena**
  - can express many metaphors as verbal analogies
- **promising first attempt at corpus-based learning of analogies**
  - first objective evaluation on human-level tests
  - not yet ready for real-world applications, but soon
    - classifying semantic relations in noun-modifier pairs