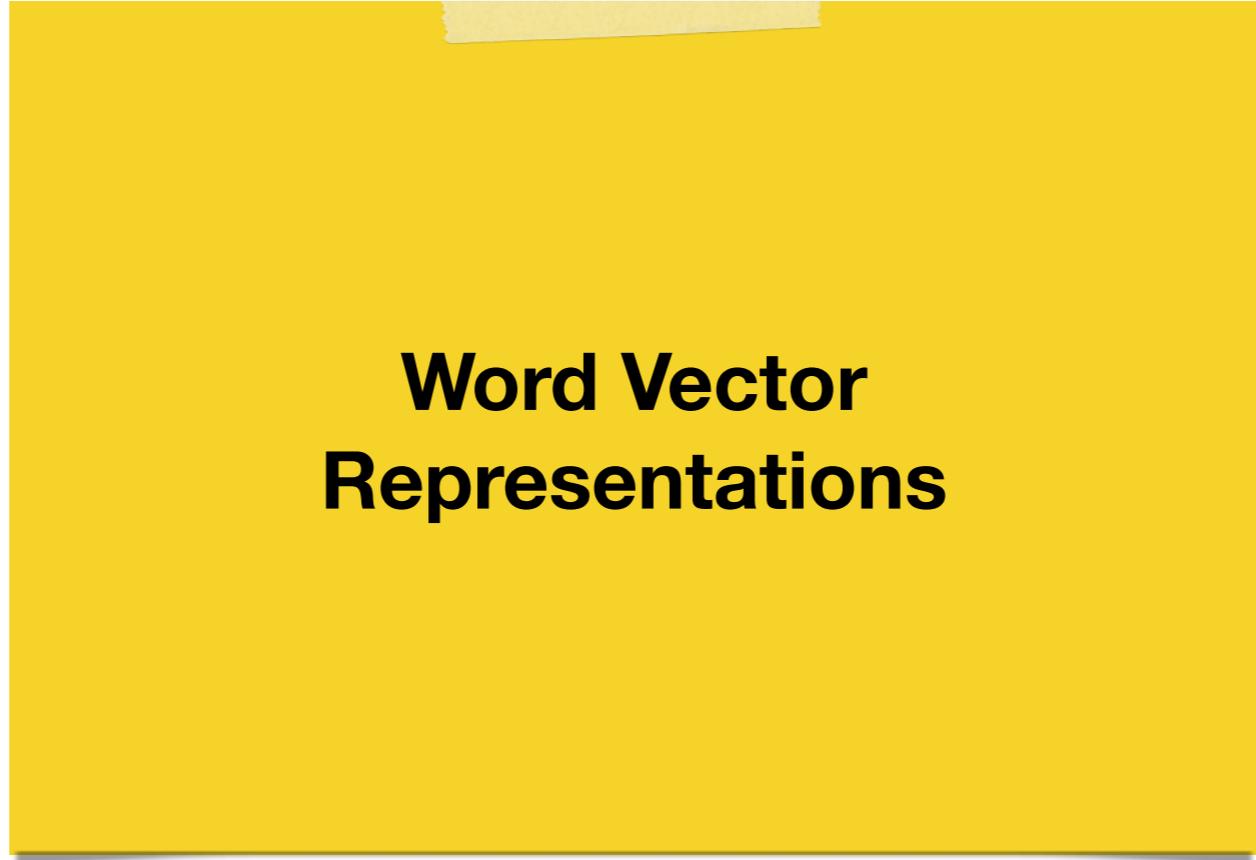


Hacking Human Language

Hendrik Heuer



PyData
London



Word Vector Representations

“You shall know a word
by the company it keeps”

–J. R. Firth 1957

Firth, John R. 1957. A synopsis of linguistic theory 1930–1955. In Studies in
linguistic analysis, 1–32. Oxford: Blackwell.

“You shall know a word
by the company it keeps”

One of the most successful ideas of modern statistical NLP

government debt problems turning into banking crises as has happened in

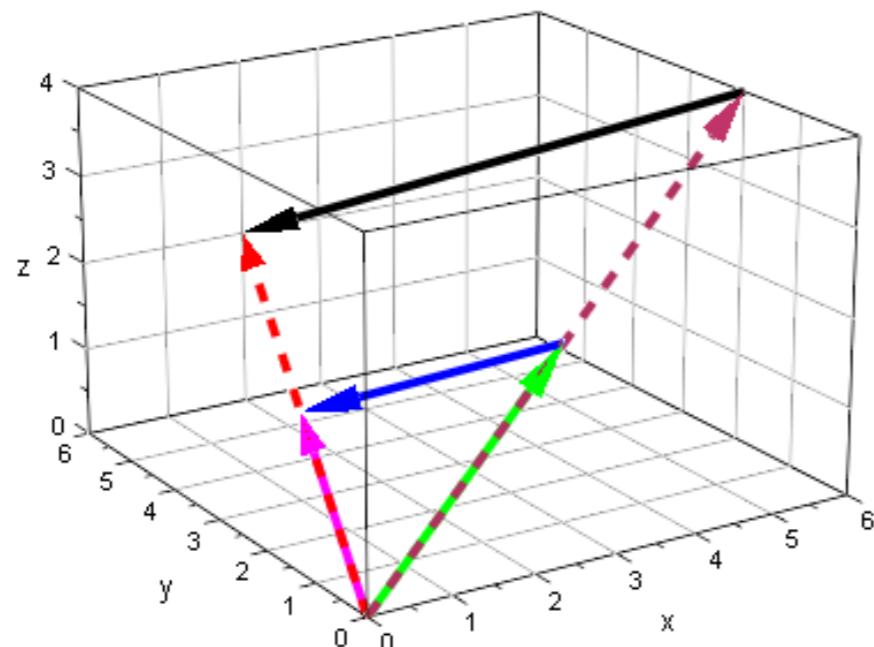
saying that Europe needs unified banking regulation to replace the hodgepodge

↖ These words will represent *banking* ↗

Quoted after Socher

word2vec

Representing a word with a vector



linguistics =

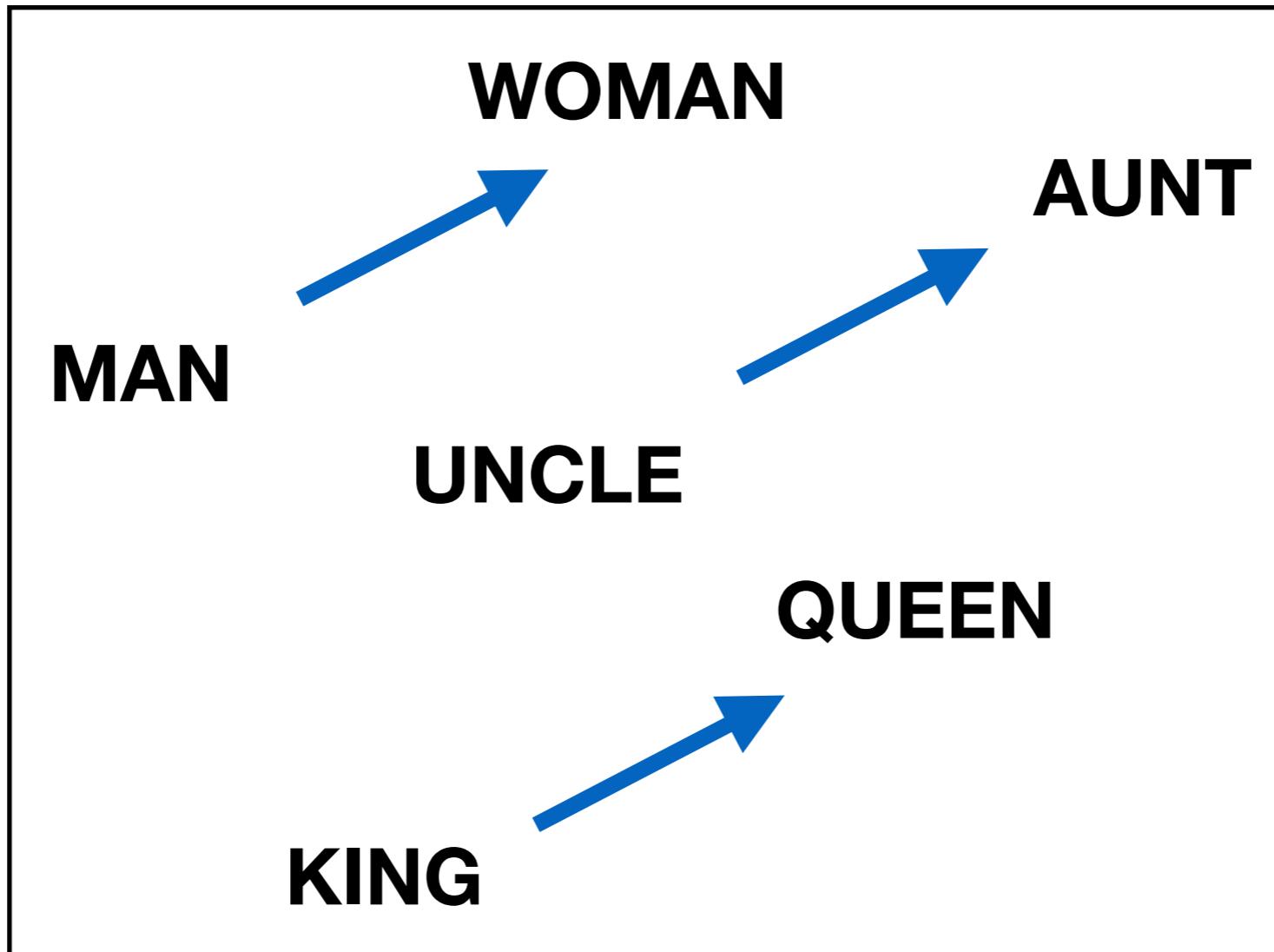
$$\begin{bmatrix} 0.286 \\ 0.792 \\ -0.177 \\ -0.107 \\ 0.109 \\ -0.542 \\ 0.349 \\ 0.271 \end{bmatrix}$$

Vectors are directions in space

Quoted after Socher

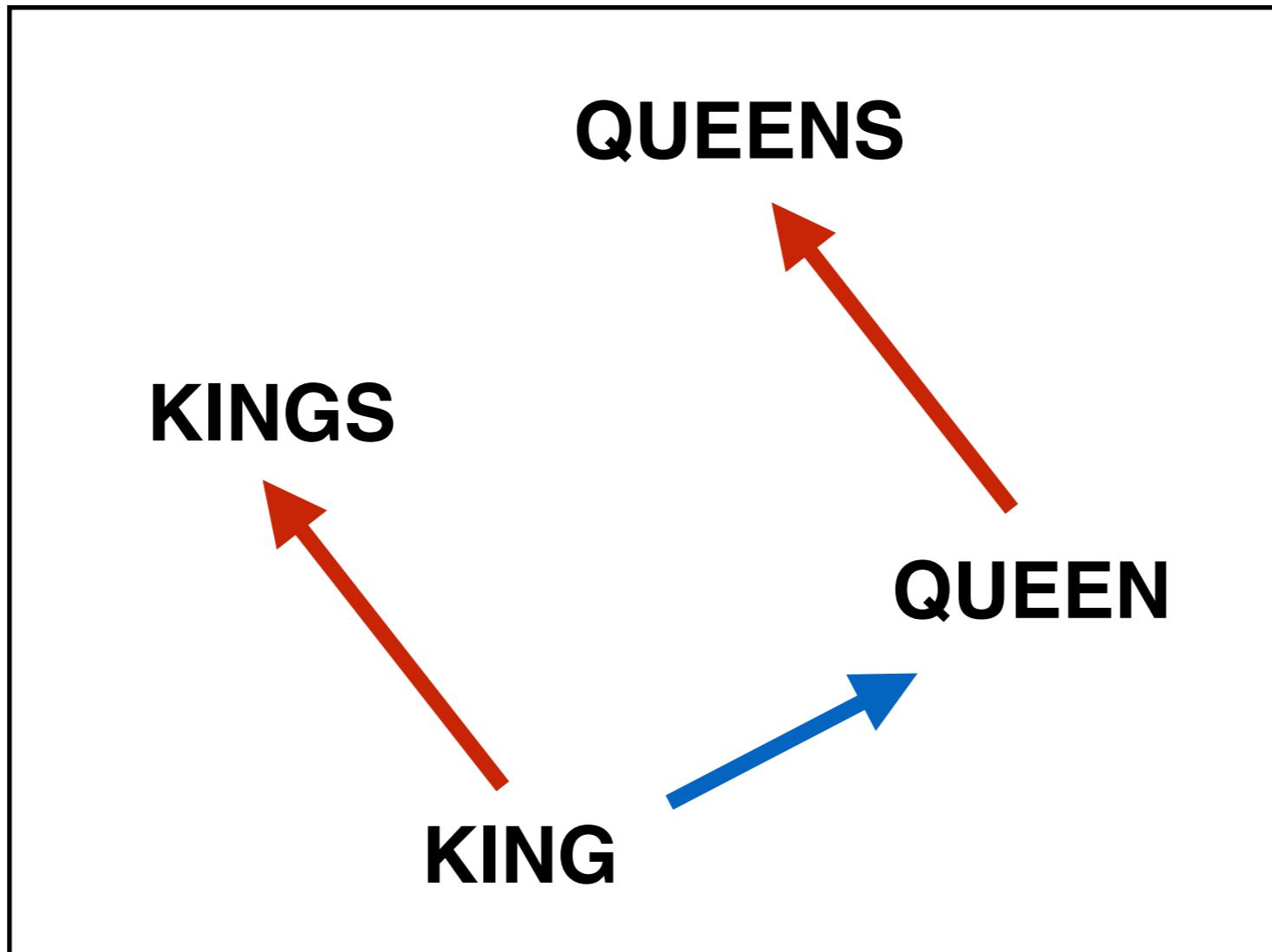
word2vec

Vectors can encode relationships



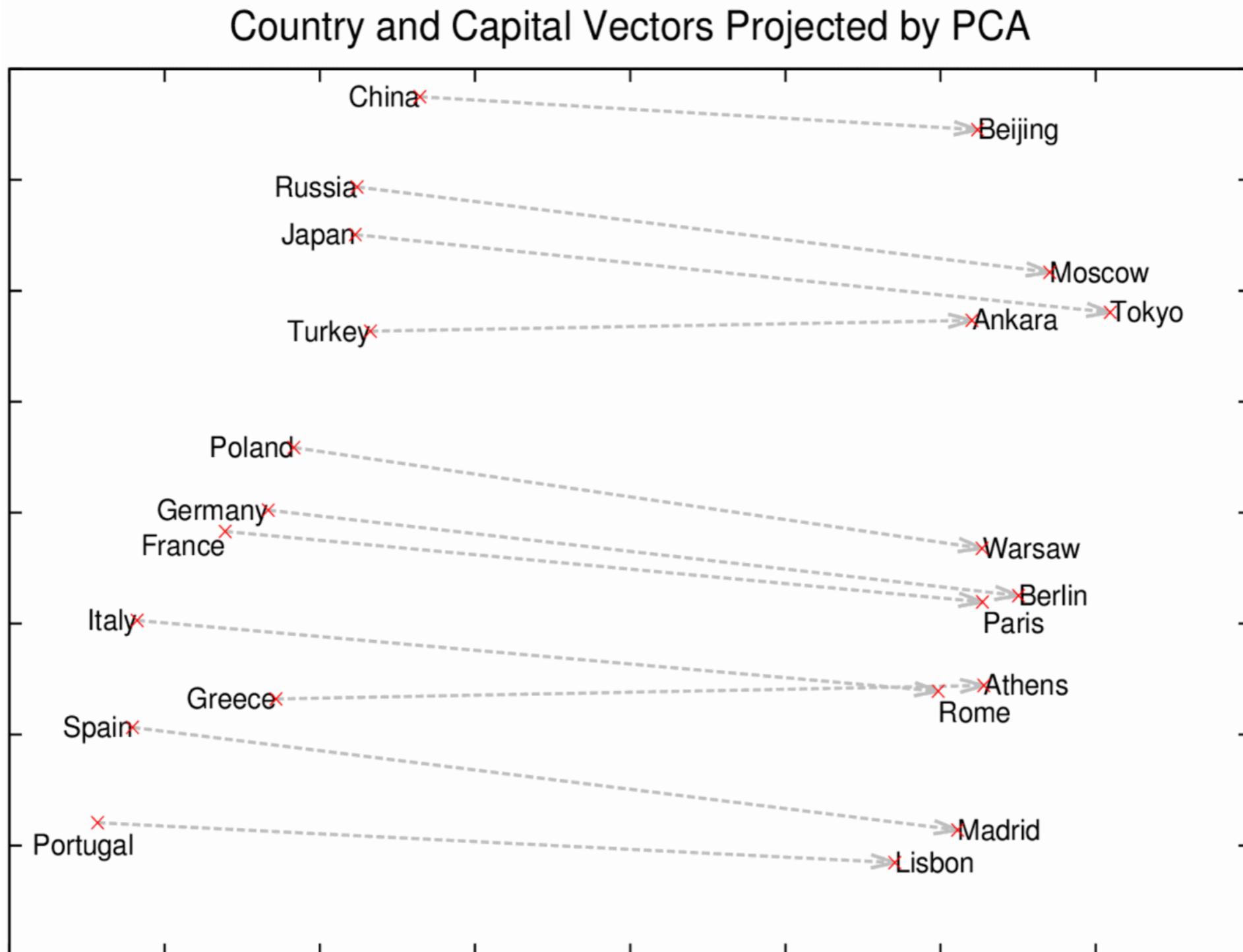
word2vec

Vectors can encode relationships



word2vec

Vectors can encode relationships



T. Mikolov, K. Chen, G. Corrado, and J. Dean, 'Efficient Estimation of Word Representations in Vector Space', CoRR, vol. abs/1301.3781, 2013 [Online].
Available: <http://arxiv.org/abs/1301.3781>

word2vec

Analogy puzzles

England is to **London** as
Germany is to ?

598.7ms **[["Berlin",0.563393235206604],["Dusseldorf",0.5625754594802856],
["Munich",0.5460122227668762],["Budapest",0.5285829901695251],
["Düsseldorf",0.5266501903533936]]**

England is to **Cameron** as
Germany is to ?

556.8ms **[["Merkel",0.5016422867774963],["Schroeder",
0.49941977858543396],["Klaus",0.4981233477592468],["Schröder",
0.4947296977043152],["Peer_Nils",0.492642343044281]]**

word2vec

Analogy puzzles

fast is to **fastest** as
slow is to ?

806.2ms `[["slowest",0.7025301456451416],["slower",0.6236234307289124],
["slowed",0.5842559337615967],["slowing",0.5462259650230408],["quickest",
0.5290436744689941]]`

wake is to **woken** as
be is to ?

929.9ms `[["been",0.41698968410491943],["tobe",0.40402814745903015],
["are",0.3866569399833679],["being",0.3746173679828644],["notbe",
0.36837878823280334]]`

word2vec

Analogy puzzles

communism is to **Karl_Marx** as
capitalism is to ?

544.7ms [["Capitalism",0.5884973406791687],["capitalist",
0.5700926184654236],["**Friedrich_Hayek**",0.5352163314819336],
["**Milton_Friedman**",0.5348755121231079],["**John_Maynard_Keynes**",
0.5335651636123657]]



Sweden

Most similar words

Word	Cosine distance
norway	0.760124
denmark	0.715460
finland	0.620022
switzerland	0.588132
belgium	0.585835
netherlands	0.574631
iceland	0.562368
estonia	0.547621
slovenia	0.531408
floorball_federation	0.529570
luxembourg	0.529477
czech_republic	0.528778
slovakia	0.526340
romania	0.524281
kista	0.522488
helsinki_vantaa	0.519936
swedish	0.519901
balrog_ik	0.514556
portugal	0.502495
russia	0.500196
slovakia_slovenia	0.496051
ukraine	0.495712

Sweden

Most similar words

Word	Cosine distance
norway	0.760124
denmark	0.715460
finland	0.620022
switzerland	0.588132
belgium	0.585835
netherlands	0.574631
iceland	0.562368
estonia	0.547621
slovenia	0.531408
floorball_federation	0.529570
luxembourg	0.529477
czech_republic	0.528778
slovakia	0.526340
romania	0.524281
kista	0.522488
helsinki_vantaa	0.519936
swedish	0.519901
balrog_ik	0.514556
portugal	0.502495
russia	0.500196
slovakia_slovenia	0.496051
ukraine	0.495712

Harvard

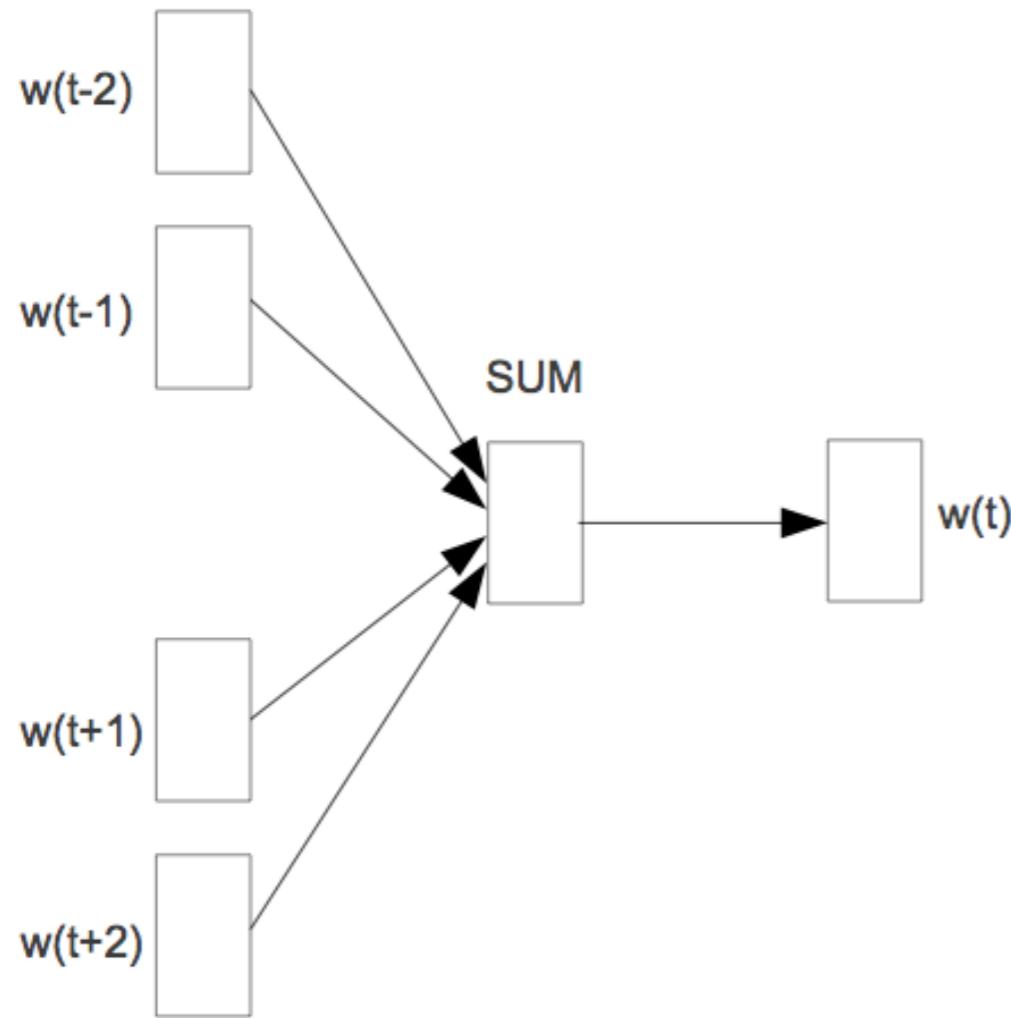
Most similar words

Word	Cosine distance
yale	0.638970
cambridge	0.612665
university	0.597709
faculty	0.588422
harvey_mudd	0.578338
johns_hopkins	0.575645
graduate	0.570294
undergraduate	0.565881
professor	0.563657
mcgill	0.562168
ph_d	0.558665
california_berkeley	0.555539
yale_university	0.550480
harvard_crimson	0.549848
princeton	0.544070
college	0.542838
oxford	0.531948
barnard_college	0.530800
professors	0.529959
princeton_university	0.529763
ucl	0.527395
doctorates	0.526292

word2vec

How it is trained

government debt problems turning into banking crises as has happened in
saying that Europe needs unified banking regulation to replace the hodgepodge



CBOW

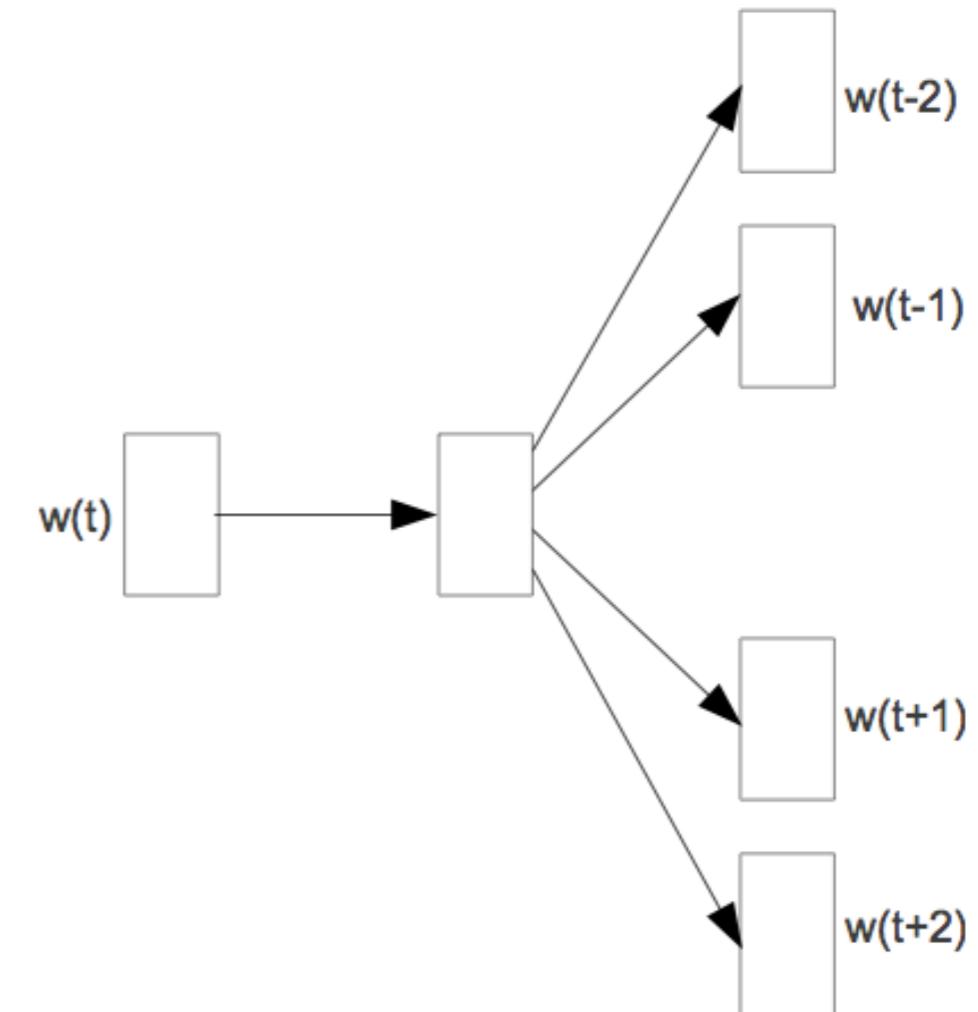
predict the current word

input

$w_{i-2}, w_{i-1}, w_{i+1}, w_{i+2}$

output

w_i



Skip-gram

predict the surrounding words

input

w_i

output

$w_{i-2}, w_{i-1}, w_{i+1}, w_{i+2}$.