# Language model

- Goal: determine $P(s = w_1 \ldots w_k)$ in some domain of interest

$$P(s) = \prod_{i=1}^{k} P(w_i \mid w_1 \ldots w_{i-1})$$

e.g., $P(w_1 w_2 w_3) = P(w_1) \, P(w_2 \mid w_1) \, P(w_3 \mid w_1 w_2)$

- Traditional n-gram language model assumption:
  "the probability of a word depends only on **context** of $n-1$ previous words"

$$\Rightarrow \widehat{P}(s) = \prod_{i=1}^{k} P(w_i \mid w_{i-n+1} \ldots w_{i-1})$$

- Typical ML-smoothing learning process (e.g., Katz 1987):
  1. compute $\widehat{P}(w_i \mid w_{i-n+1} \ldots w_{i-1}) = \dfrac{\#w_{i-n+1} \ldots w_{i-1} w_i}{\#w_{i-n+1} \ldots w_{i-1}}$ on training corpus
  2. smooth to avoid zero probabilities

# Traditional n-gram language model
## *Limitation 1): curse of dimensionality*

- Example
- train a 10-gram LM on a corpus of 100.000 unique words
- space: 10-dimensional hypercube where each dimension has 100.000 slots
- model training $\leftrightarrow$ assigning a probability to each of the $100.000^{10}$ slots
- **probability mass vanishes** $\rightarrow$ more data is needed to fill the huge space
- the more data, the more unique words! $\rightarrow$ vicious circle
- what about corpuses of $10^6$ unique words?

- $\rightarrow$ in practice, contexts are typically limited to size 2 (trigram model)
  e.g., famous Katz (1987) smoothed trigram model

- $\rightarrow$ such short context length is a limitation: a lot of information is not captured

# Traditional n-gram language model
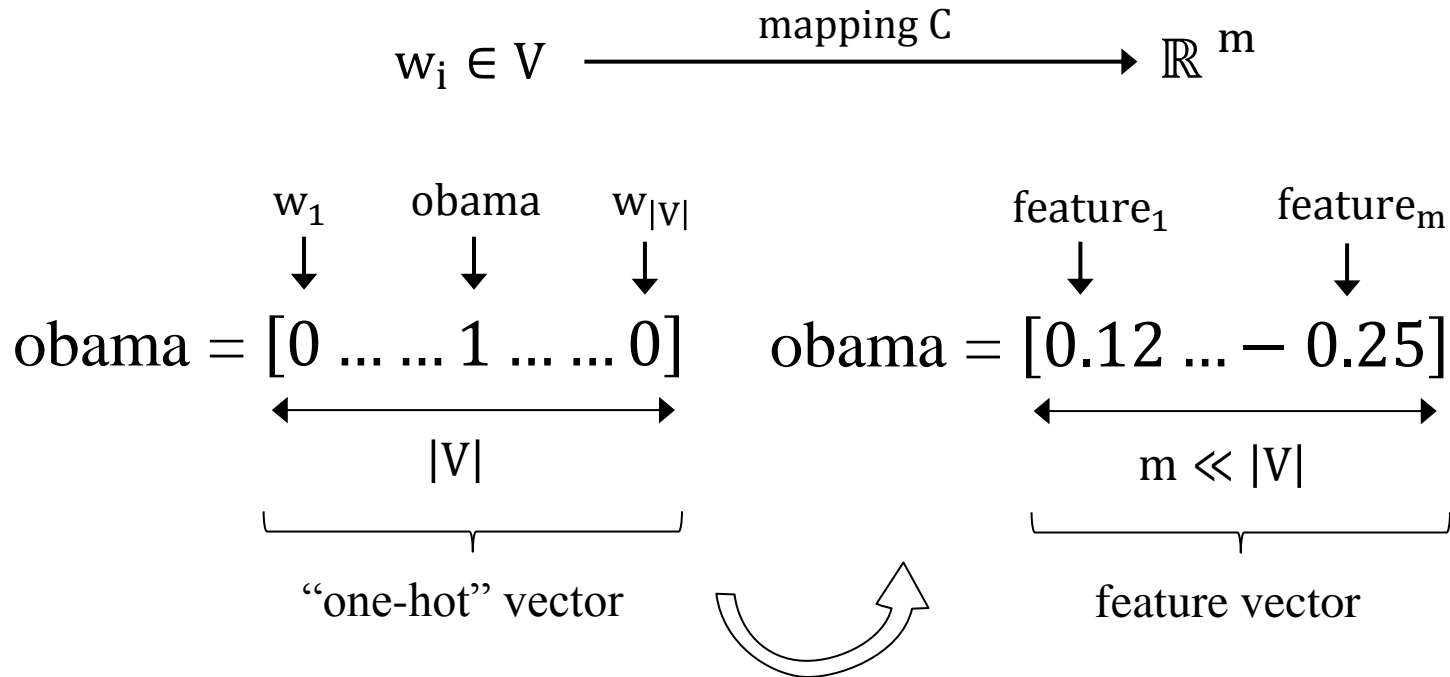## *Limitation 2): word similarity ignorance*

- We should assign similar probabilities to `Obama speaks to the media in Illinois` **and** `the President addresses the press in Chicago`

- This does not happen because of the "one-hot" vector space representation:

$$\text{obama} = [0\ 0\ 0\ 0\ ...\ 0\ 1\ 0\ 0]$$
$$\text{president} = [0\ 0\ 0\ 1\ ...\ 0\ 0\ 0\ 0]$$
$$\left.\vphantom{\begin{matrix}a\\b\end{matrix}}\right\} \quad \overrightarrow{\text{obama}}.\overrightarrow{\text{president}} = \vec{0}$$

$$\text{speaks} = [0\ 0\ 1\ 0\ ...\ 0\ 0\ 0\ 0]$$
$$\text{addresses} = [0\ 0\ 0\ 0\ ...\ 0\ 0\ 1\ 0]$$
$$\left.\vphantom{\begin{matrix}a\\b\end{matrix}}\right\} \quad \overrightarrow{\text{speaks}}.\overrightarrow{\text{addresses}} = \vec{0}$$

$$\text{illinois} = [1\ 0\ 0\ 0\ ...\ 0\ 0\ 0\ 0]$$
$$\text{chicago} = [0\ 1\ 0\ 0\ ...\ 0\ 0\ 0\ 0]$$
$$\left.\vphantom{\begin{matrix}a\\b\end{matrix}}\right\} \quad \overrightarrow{\text{illinois}}.\overrightarrow{\text{chicago}} = \vec{0}$$

- In each case, word pairs share no similarity

- This is obviously wrong

- We need to encode **word similarity** to be able to **generalize**

# Word embeddings: distributed representation of words

- Each unique word is mapped to a point in a real continuous m-dimensional space
- Typically, $|V| > 10^6$, $100 < m < 500$

$$w_i \in V \xrightarrow{\text{mapping C}} \mathbb{R}^m$$

$w_1$     obama     $w_{|V|}$        $\text{feature}_1$     $\text{feature}_m$

$$\text{obama} = [0 \ldots \ldots 1 \ldots \ldots 0] \qquad \text{obama} = [0.12 \ldots -0.25]$$

$|V|$            $m \ll |V|$

"one-hot" vector            feature vector

- Fighting the curse of dimensionality with:
  - **compression** *(dimensionality reduction)*
  - **smoothing** *(discrete to continuous)*
  - **densification** *(sparse to dense)*

- Similar words end up close to each other in the feature space

# Google's word2vec (Mikolov et al. 2013a)

- Key idea of word2vec: achieve better performance not by using a more complex model (i.e., with more layers), but by allowing a **simpler (shallower) model** to be trained on **much larger amounts of data**

- Two algorithms for learning words vectors:

  - **CBOW**: from context predict target (focus of what follows)
  - **Skip-gram**: from target predict context

- Compared to Bengio et al.'s (2003) NNLM:
  - no hidden layer (leads to 1000X speedup)
  - projection layer is shared (not just the weight matrix)
  - context: words from both **history & future**:
    "You shall know a word by the company it keeps" (John R. Firth 1957:11):

```
              …Pelé has called Neymar an excellent player…
…At the age of just 22 years, Neymar had scored 40 goals in 58 internationals…
…occasionally as an attacking midfielder, Neymar was called a true phenomenon…
```

These words will represent **Neymar**

# word2vec's Continuous Bag-of-Words (CBOW)

For each training sequence:    input = (context, target) pair: $(w_{t-\frac{n}{2}} \ldots w_{t-1} w_{t+1} \ldots w_{t+\frac{n}{2}}, w_t)$
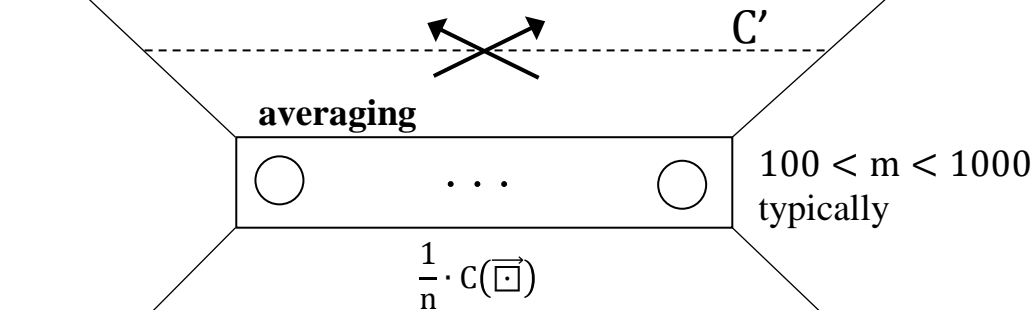
objective: minimize $E = -\log \widehat{P}(w_t \mid w_{t-n/2} \ldots w_{t-1} w_{t+1} \ldots w_{t+n/2})$

**hierarchical softmax.**        $t^{th}$ output $= P(w_i = w_t \mid w_{t-n/2} \ldots w_{t-1} w_{t+1} \ldots w_{t+n/2})$

*OUTPUT LAYER*

○        • • •        ○        |V| probabilities that sum to 1

C'

**averaging**

*PROJECTION LAYER*
*linear*

○        . . .        ○        $100 < m < 1000$ typically

$\frac{1}{n} \cdot C(\overrightarrow{\boxdot})$

table lookup in shared $C_{|V|,m}$

*INPUT LAYER*        $\overrightarrow{\boxdot} =$        1 0 0 0 1 0 0 0 0 0 0 . . . . . . 1 0 0 1 0 0 0 0 0 0 1 0        |V|

0000...0010    ···    0000...0010        0000...0010    ···    0000...0010        $n \cong 8$ typically

input context:        n/2 history words: $w_{t-\frac{n}{2}} \ldots w_{t-1}$    n/2 future words: $w_{t+1} + \cdots + w_{t+\frac{n}{2}}$

11
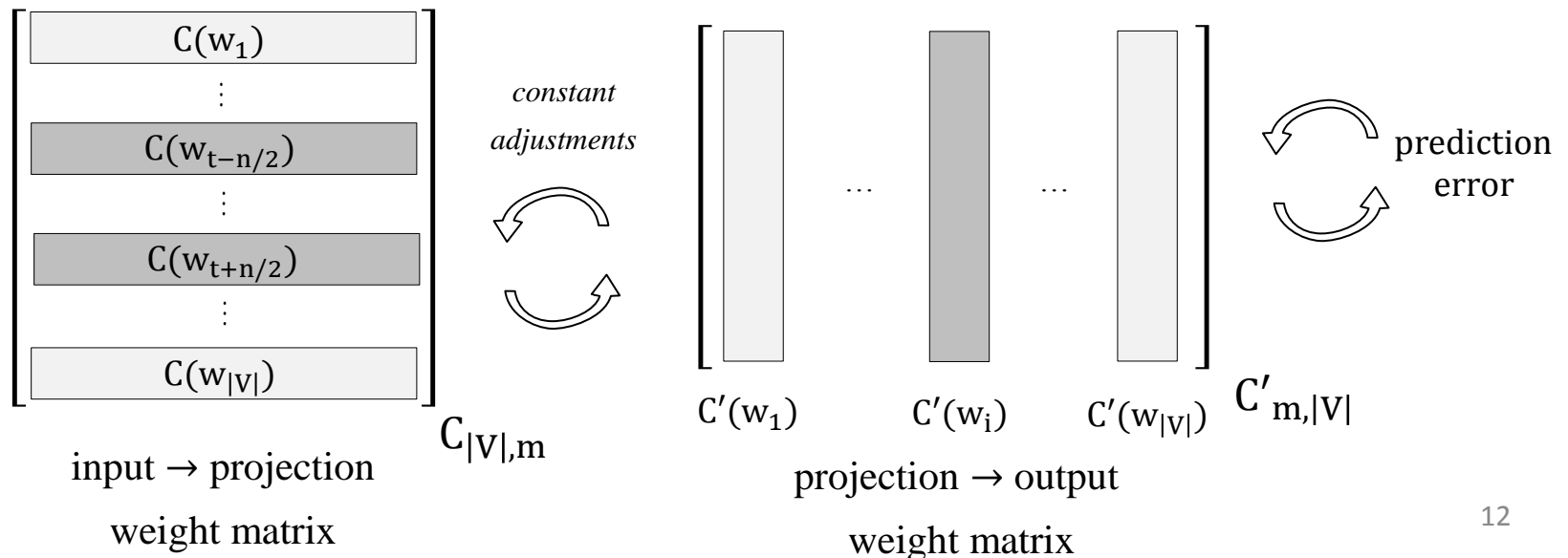
# Weight updating intuition

- For each (context, target=$w_t$) pair, only the word vectors from matrix C corresponding to the context words are updated
- Recall that we compute P ($w_i = w_t$ I context) $\forall \, w_i \in V$ . We compare this distribution to the true probability distribution (1 for $w_t$, 0 elsewhere)
- If P ($w_i = w_t$ I context) is **overestimated** (i.e., $> 0$, happens in potentially $|V| - 1$ cases), some portion of C'($w_i$) is **subtracted** from the context word vectors in C, proportionally to the magnitude of the error
- Reversely, if P ($w_i = w_t$ I context) is **underestimated** ($< 1$, happens in potentially 1 case), some portion of C'($w_i$) is **added** to the context word vectors in C
  $\rightarrow$ at each step the words move away or get closer to each other in the feature space $\rightarrow$ clustering
  $\rightarrow$ analogy with a **spring force** layout. See online demo with Chrome



constant adjustments

prediction error

$C(w_1)$

$\vdots$

$C(w_{t-n/2})$

$\vdots$

$C(w_{t+n/2})$

$\vdots$

$C(w_{|V|})$

$C_{|V|,m}$

input $\rightarrow$ projection weight matrix

$C'(w_1)$    $C'(w_i)$    $C'(w_{|V|})$    $C'_{m,|V|}$

projection $\rightarrow$ output weight matrix

12

# word2vec facts

- Complexity is $n * m + m * \log|\mathbf{V}|$ (Mikolov et al. 2013a)
- On Google news 6B words training corpus, with $|\mathbf{V}| \sim 10^6$:
  - CBOW with m = 1000 took **2 days** to train on **140 cores**
  - Skip-gram with m = 1000 took **2.5 days** on **125 cores**
  - NNLM (Bengio et al. 2003) took **14 days** on **180 cores**, for m = 100 only!
    (note that m = 1000 was not reasonably feasible on such a large training set)
- word2vec training speed $\cong$ 100K-5M words/s
- Quality of the word vectors:
  - ↗ significantly with **amount of training data** and **dimension of the word vectors** (m),
    with diminishing relative improvements
  - measured in terms of accuracy on 20K semantic and syntactic association tasks.
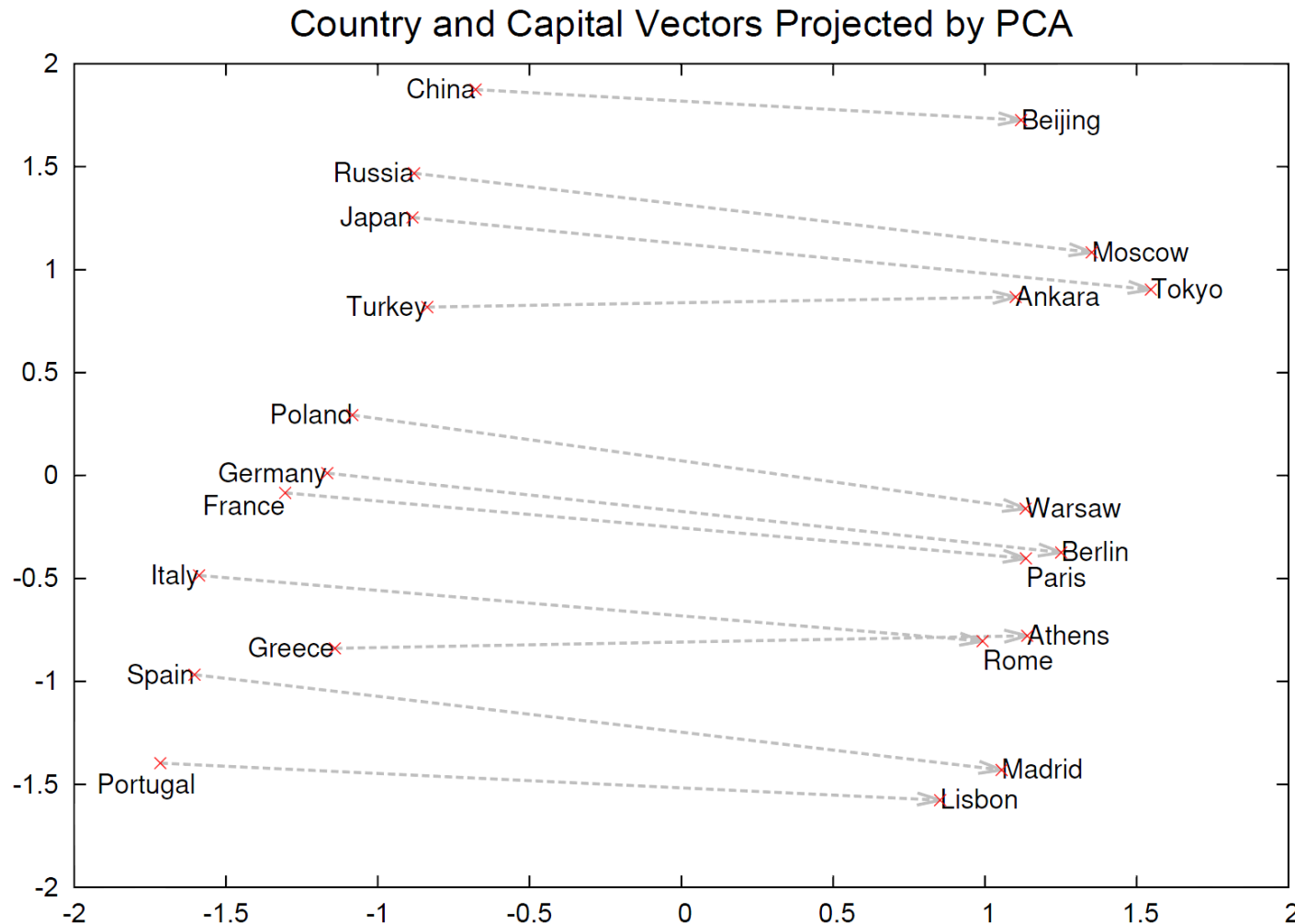    e.g., words in **bold** have to be returned:

| Capital-Country | Past tense | Superlative | Male-Female | Opposite |
|---|---|---|---|---|
| Athens: **Greece** | walking: **walked** | easy: **easiest** | brother: **sister** | ethical: **unethical** |

Adapted from Mikolov et al. (2013a)

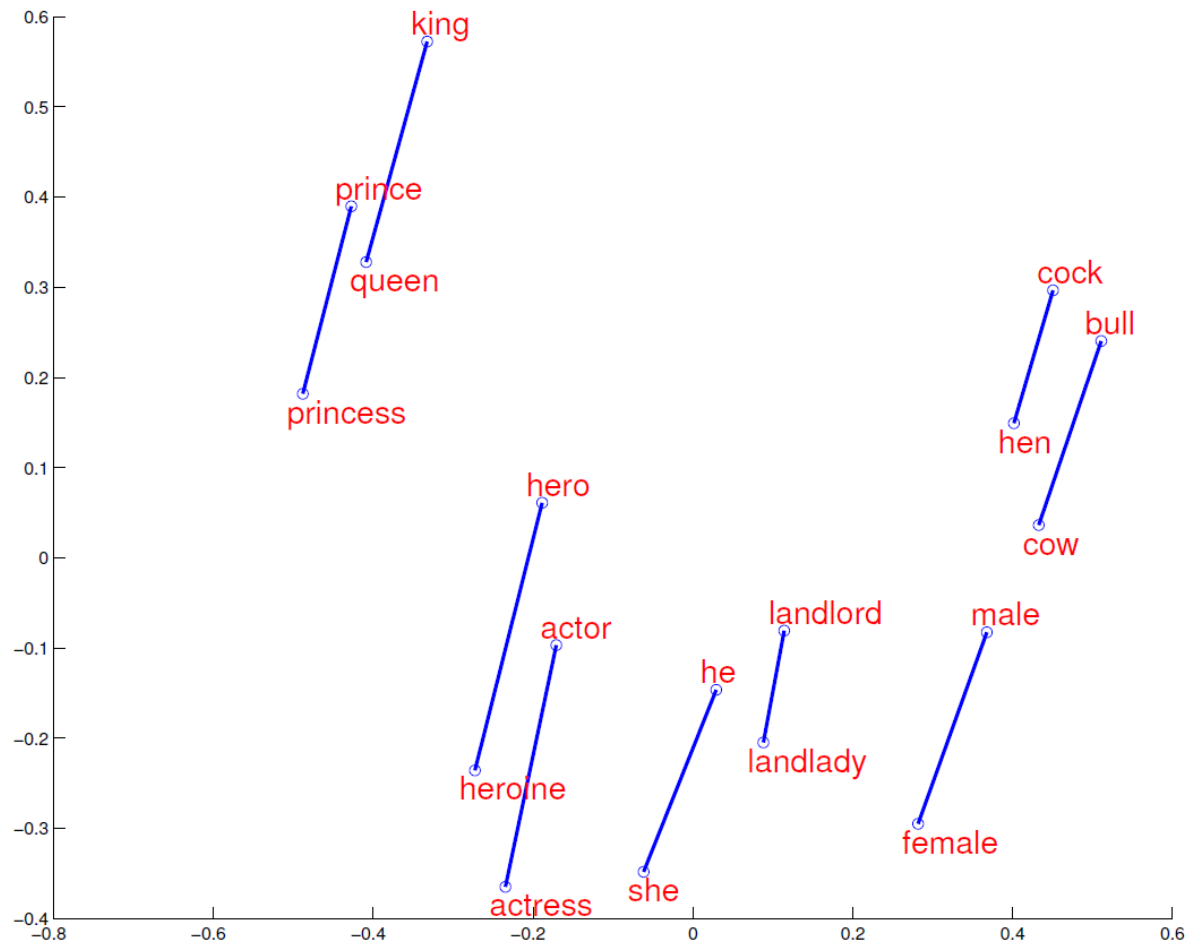- Best NNLM: 12.3% overall accuracy. Word2vec (with Skip-gram): 53.3%

References: http://www.scribd.com/doc/285890694/NIPS-DeepLearningWorkshop-NNforText#scribd
https://code.google.com/p/word2vec/

# Remarkable properties of word2vec's word vectors

## Country and Capital Vectors Projected by PCA
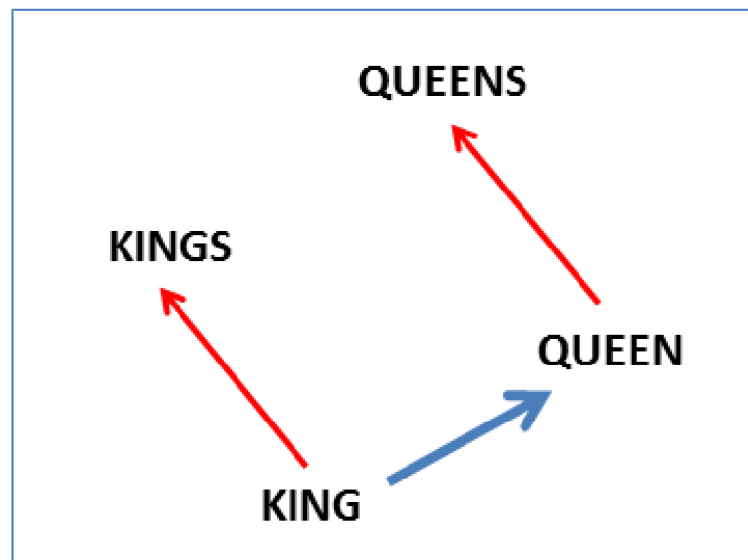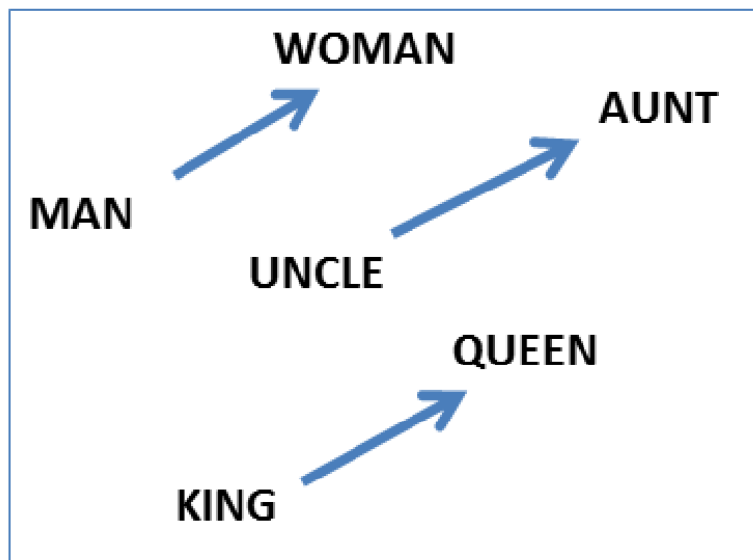


Mikolov et al. (2013b)

regularities between words are encoded in the difference vectors
e.g., there is a constant **country-capital** difference vector

14

# Remarkable properties of word2vec's word vectors



constant **female-male** difference vector

# Remarkable properties of word2vec's word vectors



constant **male-female** difference vector

constant **singular-plural** difference vector

- Vector operations are supported and make intuitive sense:

$$w_{king} - w_{man} + w_{woman} \cong w_{queen}$$

$$w_{einstein} - w_{scientist} + w_{painter} \cong w_{picasso}$$

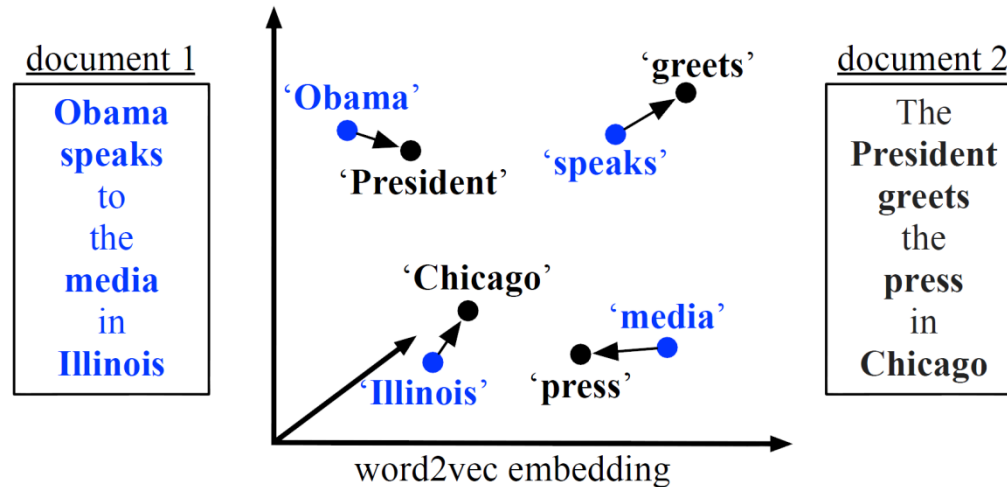$$w_{paris} - w_{france} + w_{italy} \cong w_{rome}$$

$$w_{his} - w_{he} + w_{she} \cong w_{her}$$

$$w_{windows} - w_{microsoft} + w_{google} \cong w_{android}$$

$$w_{cu} - w_{copper} + w_{gold} \cong w_{au}$$

- Online demo (scroll down to end of tutorial)

picture taken from http://www.scribd.com/doc/285890694/NIPS-DeepLearningWorkshop-NNforText#scribd

# Application to document classification



document 1

**Obama**
**speaks**
to
the
**media**
in
**Illinois**

'Obama' 'President' 'speaks' 'greets'

'Chicago' 'media' 'Illinois' 'press'

word2vec embedding

document 2

The
**President**
**greets**
the
**press**
in
**Chicago**

With the BOW representation $D_1$ and $D_2$ are at equal distance from $D_0$. Word embeddings allow to capture the fact that $D_1$ is closer.

$D_1$ **Obama** **speaks** to the **media** in **Illinois.**

$1.07 = 0.45 + 0.24 + 0.20 + 0.18$

$D_0$ The **President** **greets** the **press** in **Chicago.**

$1.63 = 0.49 + 0.42 + 0.44 + 0.28$

$D_2$ The **band** **gave** a **concert** in **Japan.**

Kusner, M. J., Sun, E. Y., Kolkin, E. N. I., & EDU, W. From Word Embeddings To Document Distances. Proceedings of the 32nd International Conference on Machine Learning, Lille, France, 2015. JMLR: W&CP volume 37.