

Aplicación de Procesamiento de Lenguaje Natural en la Ciencia

Gustavo Vazquez

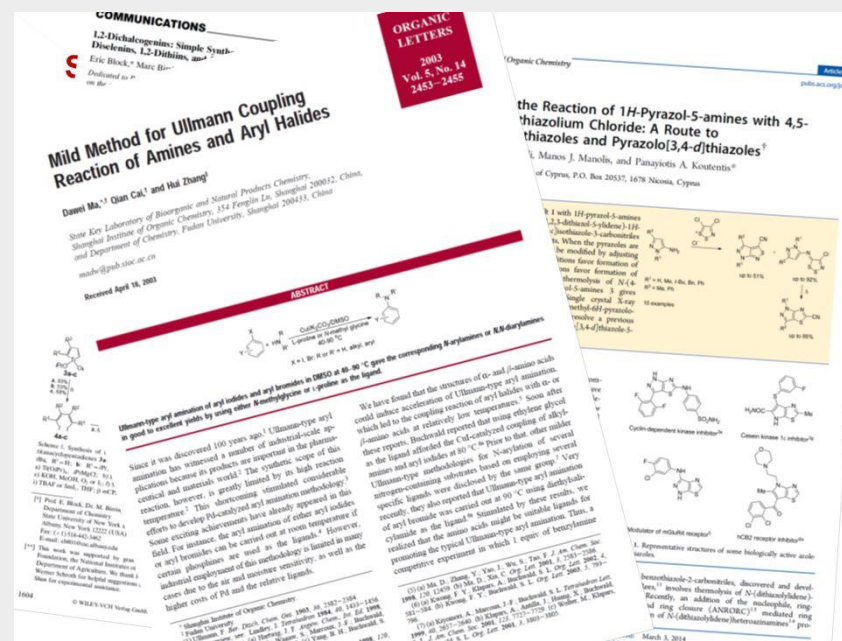
gustavo.vazquez@ucu.edu.uy tw: [@VazquezGust](https://twitter.com/VazquezGust)

Universidad Católica del Uruguay

Depto. de Ciencias de la Computación - Facultad de Ingeniería y Tecnologías

- La ciencia genera un gran volumen de información (típicamente no estructurada)

- Algunos ejemplos:
 - Diseño de Materiales/Diseño de Drogas (Cheminformatics)
 - Biología (Bioinformática)



- One-Hot encoding

{ 'a', 'aaron' .. 'apple' ... 'batman' ... 'joker' ... 'shinchan' ... 'zulu' }



[1, 0, 0, 0, 0...0]



[0, 0, 0, 0, 0...1]

- Problemas de esta representación
 - Sparse
 - No preserva aspectos semánticos

- Word Vectors: qué sucedería si...

```
word2vec('Batman') = [0.9, 0.8, 0.2]
```

```
word2vec('Joker') = [0.8, 0.3, 0.1]
```

```
word2vec('Spiderman') = [0.2, .9, 0.8]
```

```
word2vec('Thanos') = [0.3, 0.1, 0.9]
```

- Ventajas

- Los vectores pueden codificar características / información semántica
- Reducción de la dimensión

¿Cómo obtenemos esta representación?

- Ejemplo: el modelo skip-gram
 - Objetivo: dada una palabra, indicar la probabilidad de que otra aparezca en su contexto

Source Text

The quick brown fox jumps over the lazy dog.

The quick brown fox jumps over the lazy dog.

The quick brown fox jumps over the lazy dog.

The quick brown fox jumps over the lazy dog.

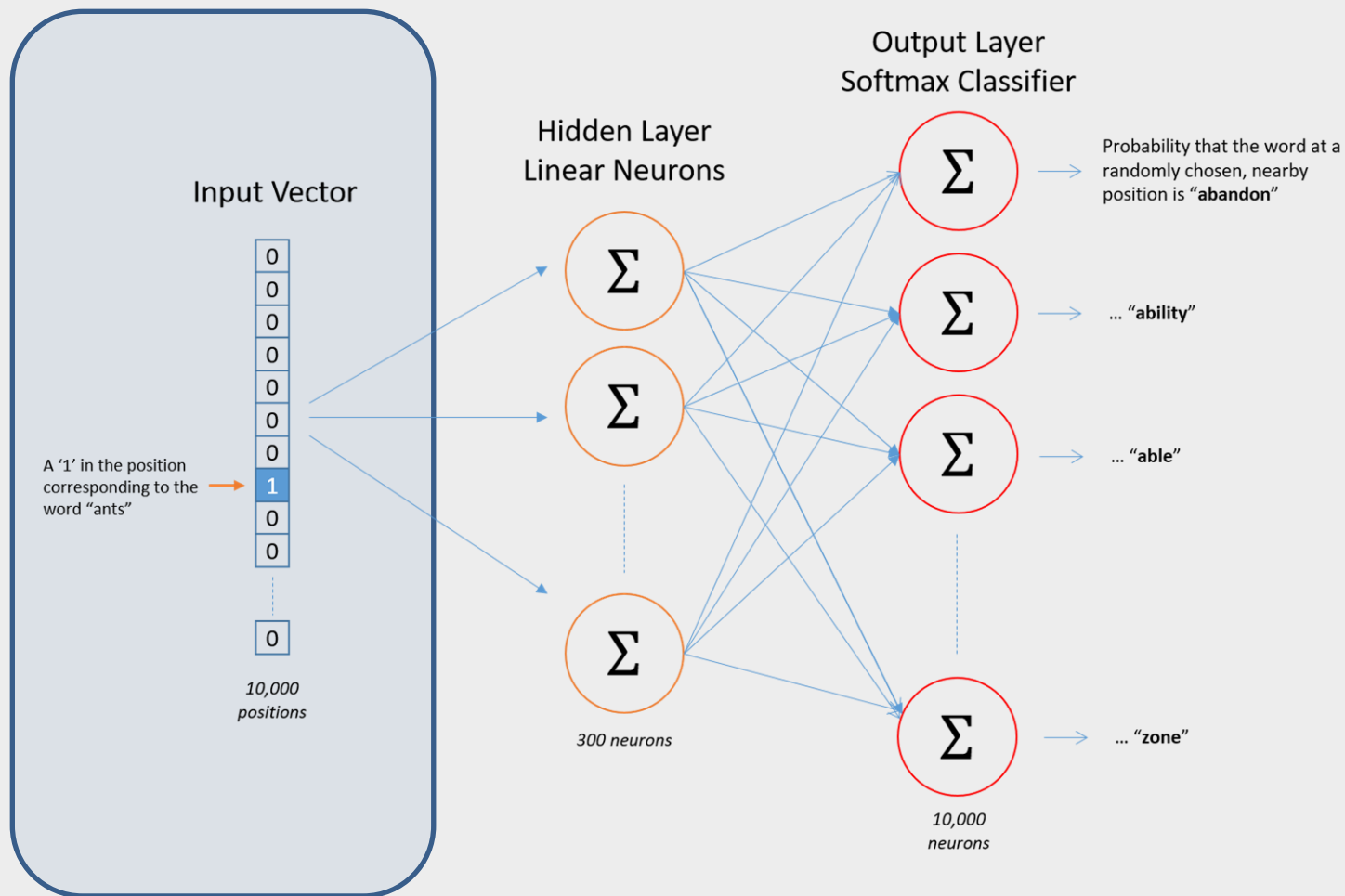
- Contexto: tamaño de ventana 2

Source Text

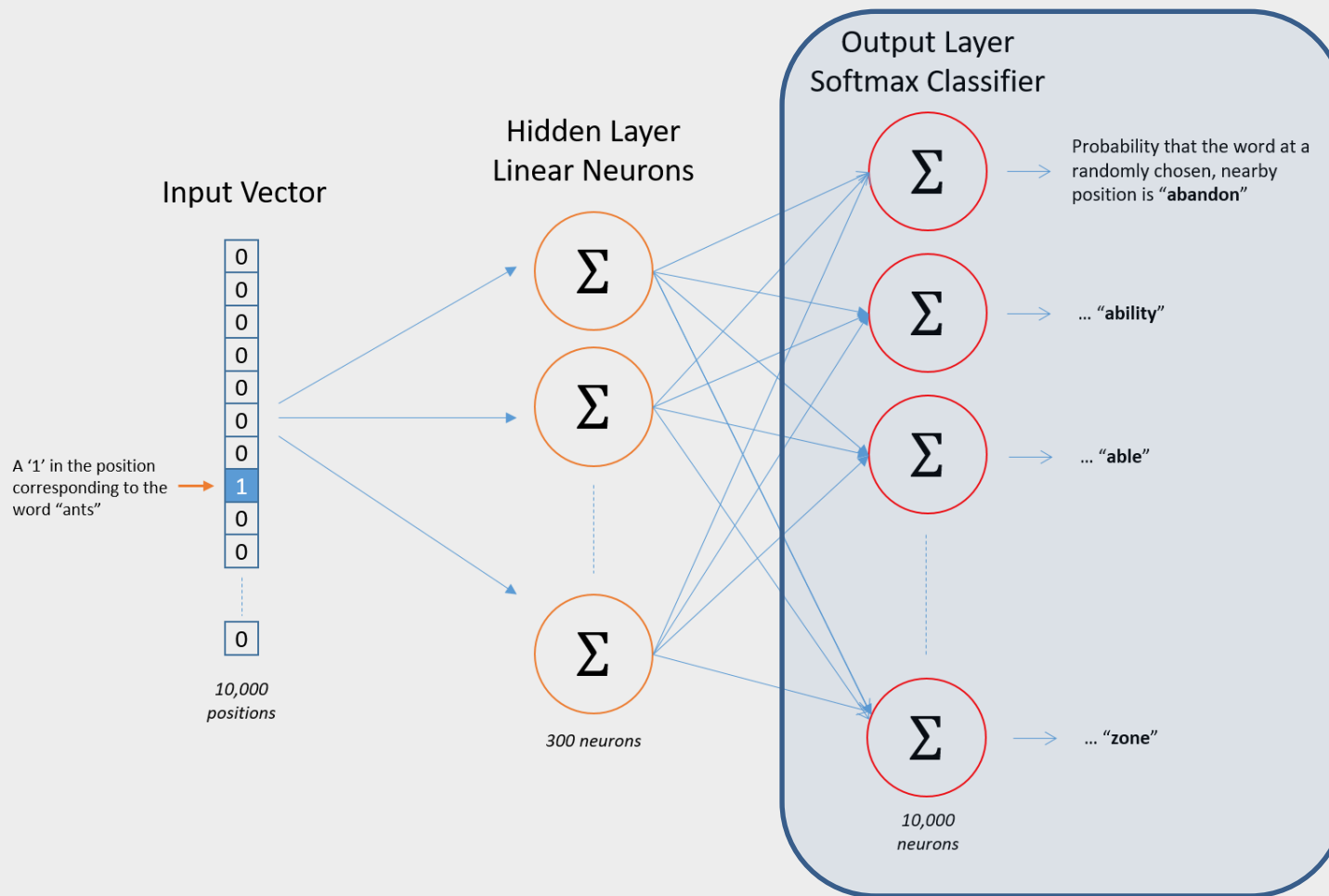
Training
Samples

<div style="display: inline-block; border: 1px solid black; padding: 2px;"> The quick brown fox jumps over the lazy dog. ➔ </div>	<p>(the, quick) (the, brown)</p>
<div style="display: inline-block; border: 1px solid black; padding: 2px;"> The quick brown fox jumps over the lazy dog. ➔ </div>	<p>(quick, the) (quick, brown) (quick, fox)</p>
<div style="display: inline-block; border: 1px solid black; padding: 2px;"> The quick brown fox jumps over the lazy dog. ➔ </div>	<p>(brown, the) (brown, quick) (brown, fox) (brown, jumps)</p>
<div style="display: inline-block; border: 1px solid black; padding: 2px;"> The quick brown fox jumps over the lazy dog. ➔ </div>	<p>(fox, quick) (fox, brown) (fox, jumps) (fox, over)</p>

- Representación de la información
 - Entrada es un vector de n dimensiones (una por cada palabra) - One-hot vector



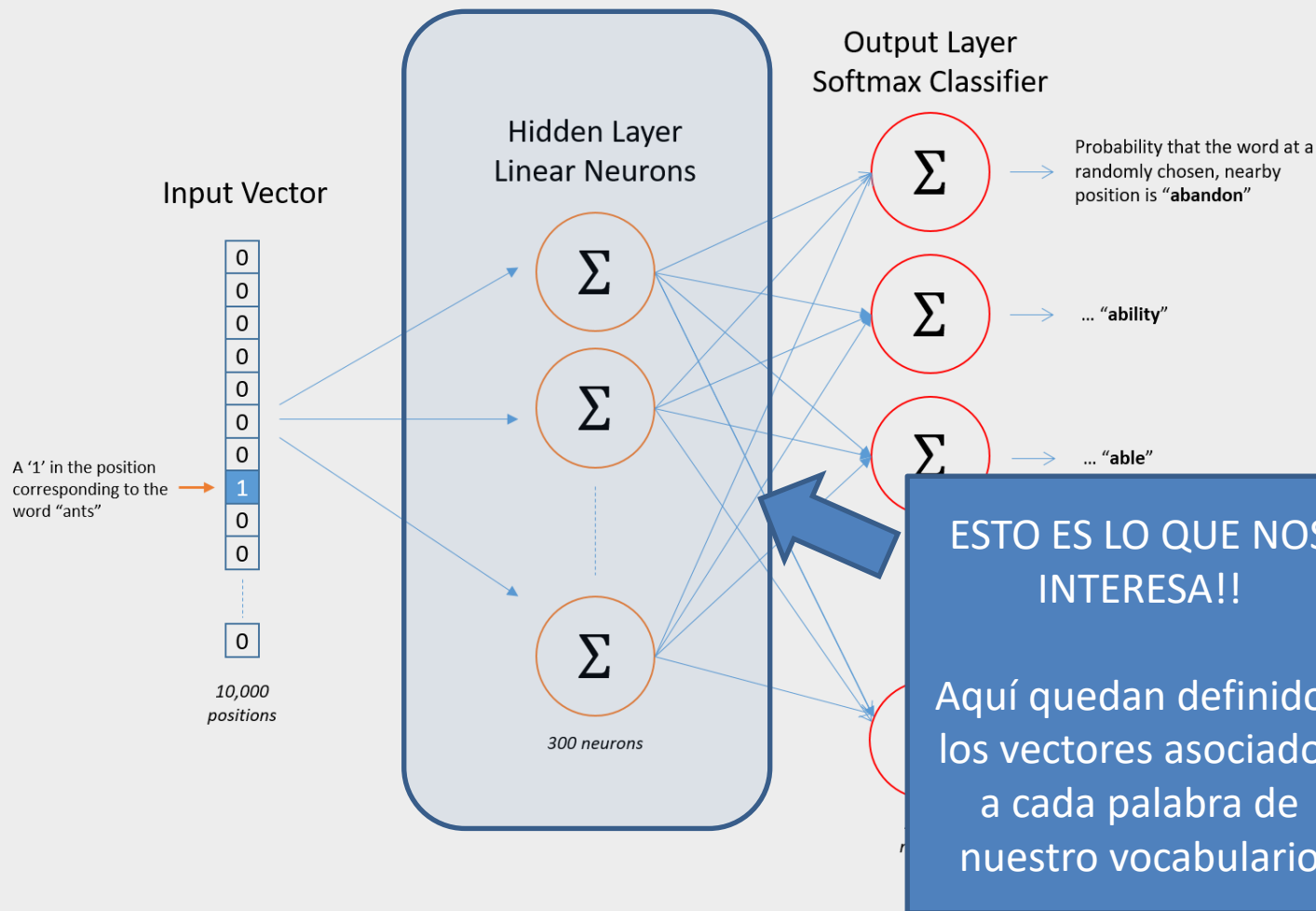
- Representación de la información
 - Entrada es un vector de n dimensiones (una por cada palabra) - One-hot vector
 - Salida es un vector de n dimensiones (cada posición es la probabilidad de que esa palabra esté en la cercanía de la palabra de la entrada)

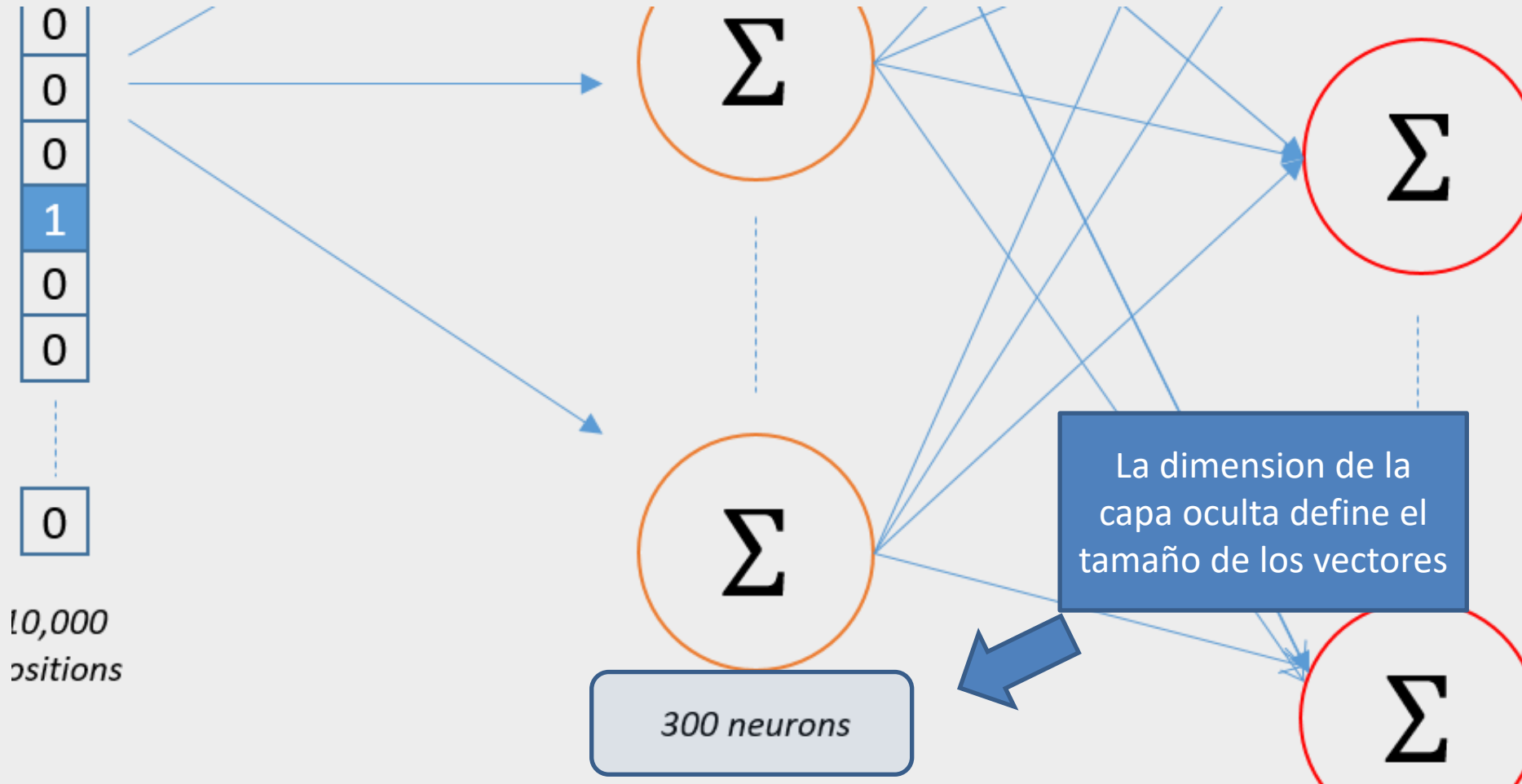


Entrenamos la red, pero la
tiramos a la basura...



**PLOT
TWIST**





Si dos palabras se usan en contextos similares,
entonces la salida de la red debería ser similar...

Si dos palabras se usan en contextos similares,
entonces la salida de la red debería ser similar...

... y sus vectores asociados
(word vectors) también

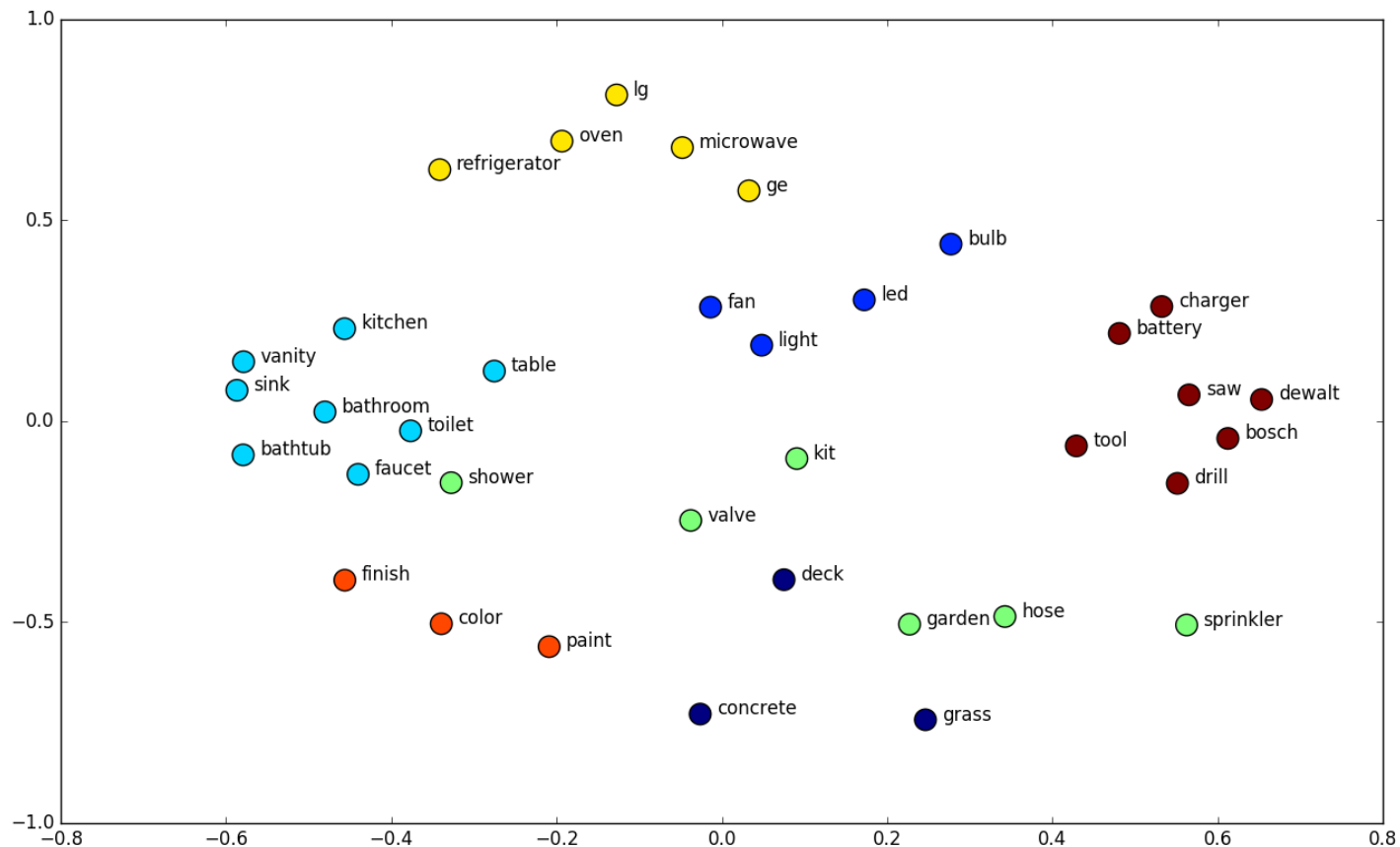
- Algunas consideraciones
 - Si bien el entrenamiento corresponde a un problema supervisado, no exige ningún tipo de proceso de etiquetado
 - Ídem con las relaciones semánticas (son capturadas por los word vectors a partir de los contextos)
 - El entrenamiento se realiza utilizando cualquier corpus disponible: Wikipedia, Google News, un dominio específico, etc.

- Ejemplo:

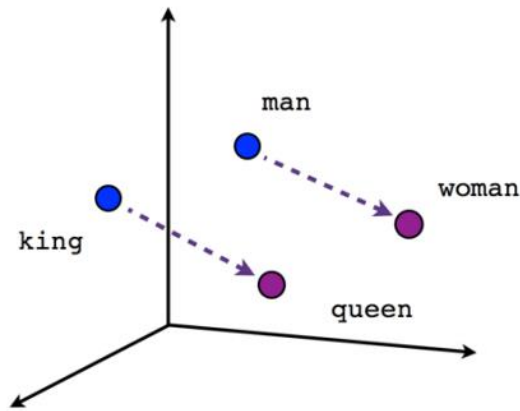
```
# look up top 6 words similar to 'shocked'  
w1 = ["shocked"]  
model.wv.most_similar (positive=w1,topn=6)
```

```
[('horrificed', 0.80775386095047),  
 ('amazed', 0.7797470092773438),  
 ('astonished', 0.7748459577560425),  
 ('dismayed', 0.7680633068084717),  
 ('stunned', 0.7603034973144531),  
 ('appalled', 0.7466776371002197)]
```

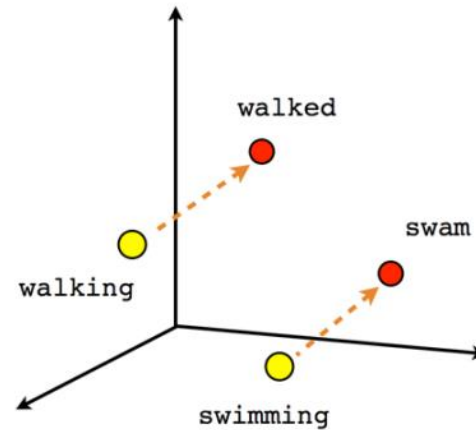
- Preservación de la relación semántica



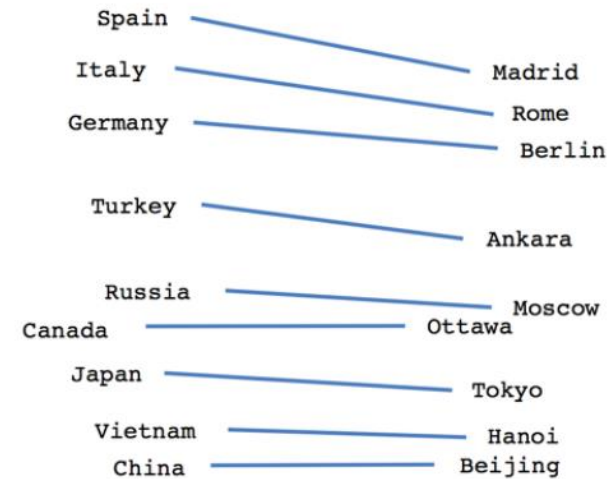
- Preservación de la relación semántica



Male-Female

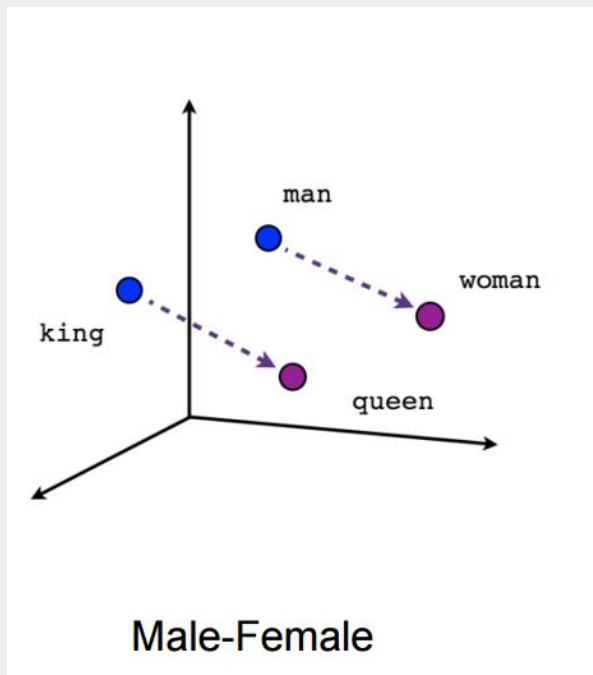


Verb tense



Country-Capital

- Preservación de la relación semántica
 - Word analogy

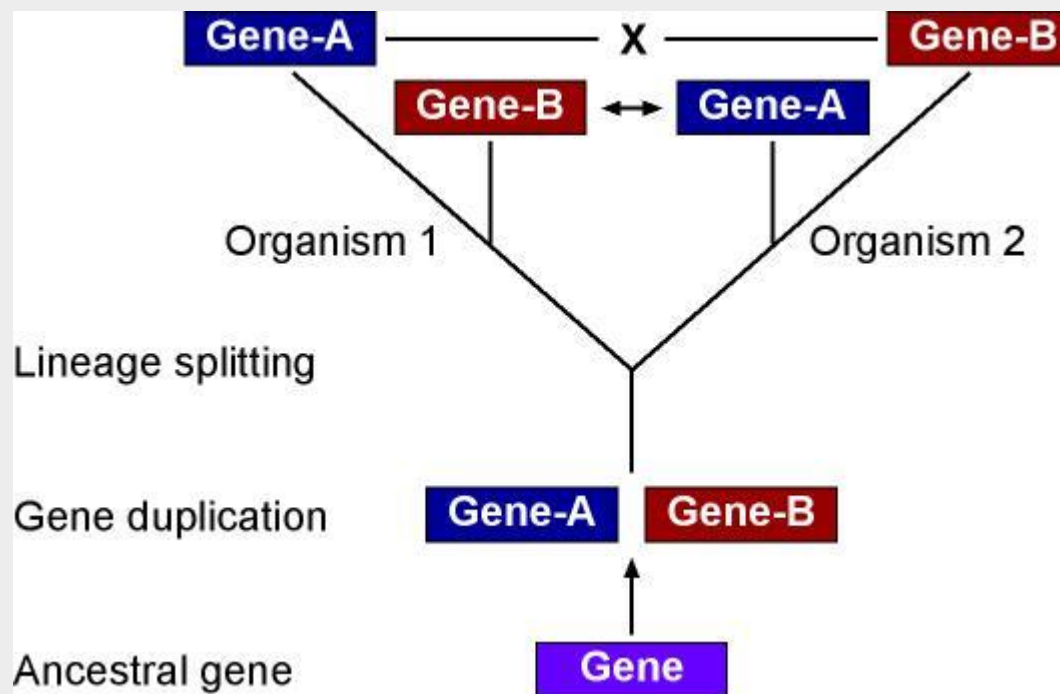


king-man+woman -> ?

Algunas aplicaciones

- Knowledge discovery: Bioinformática
 - Costos de secuenciación cada vez más económicos
 - Estudios sobre nuevos organismos
 - La información es reportada en forma de texto libre (papers)
 - No existen estándares para la sistematización de la información

- Ejemplo: genes ortólogos



GenA-Organismo1+Organismo2 = ?


GenX = ?

- Knowledge discovery: ciencia de materiales

nature

Letter | Published: 03 July 2019

Unsupervised word embeddings capture latent knowledge from materials science literature

Vahe Tshitoyan , John Dagdelen, Leigh Weston, Alexander Dunn, Ziqin Rong, Olga Kononova, Kristin A. Persson, Gerbrand Ceder  & Anubhav Jain 

Nature **571**, 95–98(2019) | [Cite this article](#)

42k Accesses | **13** Citations | **1562** Altmetric | [Metrics](#)

Abstract

The overwhelming majority of scientific knowledge is published as text,

ferromagnetic–NiFe+IrMn



antiferromagnetic

- Algunas consideraciones
 - “Entrenamiento” fuertemente dependiente del corpus disponible
 - “dificultad” de acceso a la información (copyright)
 - Sesgo de los datos: estereotipos
 - doctor–man+woman = nurse

- Otras líneas de investigación de nuestra facultad
- Procesamiento de señales / reconocimiento de imágenes y patrones (Alvaro Pardo)
 - Sistemas de recomendación clínica y NLP para procesamiento de historias clínicas (Ernesto Ocampo)
 - Machine Learning en bioinformática – Urugenomes (Lucía Spangenberg)
 - Cheminformatics (Gustavo Vazquez)

Contactos

¡ Muchas Gracias !

<http://www.ucu.edu.uy/>

Gustavo Vazquez
gustavo.vazquez@ucu.edu.uy