

Procesamiento del Lenguaje Natural (NLP). Word Embeddings y Tokenization

Inteligencia Artificial e Ingeniería del Conocimiento

Constantino Antonio García Martínez

Universidad San Pablo Ceu

- Russell, Stuart J., and Peter Norvig. Artificial intelligence: a modern approach. Pearson, 2016.

¿Qué es NLP?

¿Qué es el Procesamiento del Lenguaje Natural?

Objetivo: Permitir que las máquinas comprendan, interpreten y generen lenguaje humano

Tareas típicas:

- Traducción automática
- Análisis de sentimientos
- Chatbots y asistentes virtuales
- Resumen automático
- Reconocimiento de entidades

1. Ambigüedad

"Vi al hombre con el telescopio"

¿Quién tiene el telescopio?

1. Ambigüedad

"Vi al hombre con el telescopio"

¿Quién tiene el telescopio?

2. Dependencia del Contexto

"banco"

¿Institución financiera o asiento?

1. Ambigüedad

“Vi al hombre con el telescopio”

¿Quién tiene el telescopio?

2. Dependencia del Contexto

“banco”

¿Institución financiera o asiento?

3. Variabilidad

- Mismo significado, diferentes formas: *“no está mal”* == *“está bien”*
- Idiomas, dialectos, registros formales/informales: *“Está en el colo”*
- Evolución constante: neologismos, argot: *“el profe de IA da lache”*

Observación Clave

El lenguaje humano es fundamentalmente diferente de:

- Datos numéricos estructurados
- Imágenes (píxeles en cuadrícula)
- Datos tabulares

Consecuencia: Necesitamos representaciones y algoritmos especializados

Representación de Palabras

Representación Naive: One-Hot Encoding

Idea: Representar cada palabra como un vector con un 1 y el resto 0s

Example: Vocabulario pequeño

$V = \{\text{gato, gata, perro, rey, reina, hombre, mujer}\}$

gato = [1, 0, 0, 0, 0, 0, 0]

gata = [0, 1, 0, 0, 0, 0, 0]

perro = [0, 0, 1, 0, 0, 0, 0]

rey = [0, 0, 0, 1, 0, 0, 0]

...

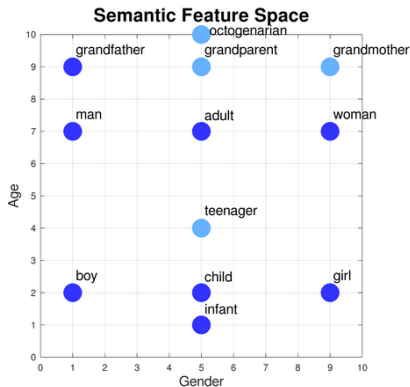
Problemas

1. **Dimensionalidad:** Tamaño del vector = $|V|$ (típicamente 50k-1M palabras)
2. **No captura similitud:** “gato” y “perro” son tan diferentes como “gato” y “democracia”
3. **No captura relaciones semánticas**

Word Embeddings: Representación Densa

Idea: Representar palabras en espacio continuo de dimensión baja (típicamente 50-300 (GloVe) hasta 1024-4096+ en arquitecturas modernas)

Example: Embeddings en dimensión 2



Word Coordinates		
	Gender	Age
grandmother	[9,	9]
grandparent	[5,	9]
octogenarian	[5,	10]
teenager	[5,	4]

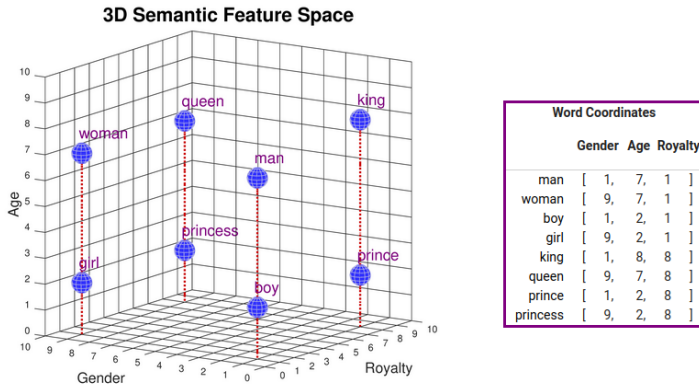
1

¹Fuente: <https://www.cs.cmu.edu/~dst/WordEmbeddingDemo/tutorial.html>

Ventajas

1. **Dimensión fija y pequeña** (no crece con vocabulario)
2. **Captura similitud:** palabras similares \rightarrow vectores cercanos
3. **Captura relaciones:** rey - hombre + mujer \approx reina

Example: Embeddings en dimensión 3



2

²Fuente: <https://www.cs.cmu.edu/~dst/WordEmbeddingDemo/tutorial.html>

Firth (1957)

“Conocerás una palabra por la compañía que tiene”

Idea: Palabras que aparecen en contextos similares tienen significados similares

Firth (1957)

“Conocerás una palabra por la compañía que tiene”

Idea: Palabras que aparecen en contextos similares tienen significados similares

Example: Mismo contexto, palabras similares

- “El _____ doméstico es un animal muy independiente”
- Puede ser: “gato”, “perro”, “conejo”
- NO puede ser: “teorema”, “democracia”, “azul”

Si “gato” y “perro” pueden aparecer en los mismos contextos

→ Sus embeddings deberían ser similares

Firth (1957)

“Conocerás una palabra por la compañía que tiene”

Idea: Palabras que aparecen en contextos similares tienen significados similares

Example: Mismo contexto, palabras similares

- “El _____ doméstico es un animal muy independiente”
- Puede ser: “gato”, “perro”, “conejo”
- NO puede ser: “teorema”, “democracia”, “azul”

Si “gato” y “perro” pueden aparecer en los mismos contextos
→ Sus embeddings deberían ser similares

Fundamento de: Word2Vec, GloVe, FastText, etc.

Global Vectors (GloVe): Aprende embeddings basándose en co-ocurrencias

Loss Function de GloVe (simplificada)

$$L = \sum_{i,j} f(X_{ij}) \left(w_i^T w_j + b_i + b_j - \log X_{ij} \right)^2$$

- X_{ij} : número de veces que palabra i y j co-ocurren
- w_i, w_j : embeddings de palabras
- $f(X_{ij})$: función de peso

1. Antes (2013-2018): Word2Vec, GloVe, FastText

- Entrenados por separado del modelo principal
- El modelo principal se entrenaba sobre estos embeddings
- Usualmente, los embeddings se mantenían congelados o se aplicaba fine-tuning

Evolución de las Representaciones de Palabras

1. Antes (2013-2018): Word2Vec, GloVe, FastText

- Entrenados por separado del modelo principal
- El modelo principal se entrenaba sobre estos embeddings
- Usualmente, los embeddings se mantenían congelados o se aplicaba fine-tuning

2. Ahora (2018+): BERT, GPT, LLaMA, etc.

- Los embeddings se aprenden de forma conjunta con el resto de la red neuronal

Evolución de las Representaciones de Palabras

1. Antes (2013-2018): Word2Vec, GloVe, FastText

- Entrenados por separado del modelo principal
- El modelo principal se entrenaba sobre estos embeddings
- Usualmente, los embeddings se mantenían congelados o se aplicaba fine-tuning

2. Ahora (2018+): BERT, GPT, LLaMA, etc.

- Los embeddings se aprenden de forma conjunta con el resto de la red neuronal

Los embeddings iniciales en ambos casos son estáticos (1 palabra = 1 vector). Sin embargo, estudiaremos que tras pasar por varias capas Transformer, **los embeddings se vuelven contextuales**

Example: Embedding contextual

“Me senté en el **banco**” \neq “Fui al **banco**” (diferentes vectores)

Tokenización: De Palabras a Subpalabras

Con word embeddings, cada palabra única necesita su propio vector

Problemas

1. Palabras fuera de vocabulario (OOV):

- Vocabulario de entrenamiento: {gato, gata, perro, ...}
- Palabra nueva: "gatito" → ¿cómo representarla?
- Solución común: token especial <UNK> (pierde información)

Con word embeddings, cada palabra única necesita su propio vector

Problemas

1. Palabras fuera de vocabulario (OOV):

- Vocabulario de entrenamiento: {gato, gata, perro, ...}
- Palabra nueva: "gatito" → ¿cómo representarla?
- Solución común: token especial <UNK> (pierde información)

2. Palabras relacionadas tratadas como independientes:

- "jugar", "jugando", "jugó", "jugador" → 4 vectores diferentes
- No capturan que comparten raíz "jug-"

Con word embeddings, cada palabra única necesita su propio vector

Problemas

1. Palabras fuera de vocabulario (OOV):

- Vocabulario de entrenamiento: {gato, gata, perro, ...}
- Palabra nueva: "gatito" → ¿cómo representarla?
- Solución común: token especial <UNK> (pierde información)

2. Palabras relacionadas tratadas como independientes:

- "jugar", "jugando", "jugó", "jugador" → 4 vectores diferentes
- No capturan que comparten raíz "jug-"

3. Vocabulario enorme:

- Español: 100k palabras comunes
- Cada palabra flexionada = entrada diferente
- "hablar": hablo, hablas, habla, hablamos, habláis, hablan, hablaba, ...

Solución: Tokenización por Subpalabras (BPE)

Byte Pair Encoding (BPE): Dividir palabras en unidades más pequeñas (subpalabras)

Example: Tokenización BPE

Vocabulario de tokens: {gat, o, a, ito, perr, jugar, jug, ando, ador}

- “gato” → [“gat”, “o”]
- “gata” → [“gat”, “a”]
- “gatito” → [“gat”, “ito”] (¡funciona aunque no esté en vocabulario!; ¡muy útil para erratas!)
- “jugando” → [“jug”, “ando”]
- “jugador” → [“jug”, “ador”]

Solución: Tokenización por Subpalabras (BPE)

Byte Pair Encoding (BPE): Dividir palabras en unidades más pequeñas (subpalabras)

Example: Tokenización BPE

Vocabulario de tokens: {gat, o, a, ito, perr, jugar, jug, ando, ador}

- “gato” → [“gat”, “o”]
- “gata” → [“gat”, “a”]
- “gatito” → [“gat”, “ito”] (¡funciona aunque no esté en vocabulario!; ¡muy útil para erratas!)
- “jugando” → [“jug”, “ando”]
- “jugador” → [“jug”, “ador”]

Ventajas de BPE

1. **Sin palabras OOV:** Cualquier palabra se puede representar
2. **Comparte información:** “jug-” tiene mismo embedding en las variantes
3. **Vocabulario más pequeño:** 50k tokens vs 100k+ palabras
4. **Eficiente para múltiples idiomas**

Example: `tokenization.py`