

Ejercicio 1 Entrenamiento de Word2Vec desde cero

Preguntas:

- Qué representa un vector de palabras en Word2Vec?
- Cuál es la diferencia entre el enfoque CBOW y Skip-Gram?
- Qué significa que dos palabras tengan vectores cercanos?
- Cómo influye el parámetro window en el entrenamiento?
- Por qué es necesario hacer tokenización antes de entrenar?

Respuestas Ejercicio 1:

Un vector numérico multidimensional que captura el significado semántico de la palabra basado en los contextos donde aparece. Cada dimensión representa una característica semántica latente.

CBOW Continuous Bag of Words predice la palabra central basándose en las palabras del contexto circundante. Es más rápido para corpus grandes. Skip-Gram predice las palabras del contexto basándose en la palabra central. Es mejor para corpus pequeños y palabras raras.

Significa que tienen significados similares o aparecen en contextos parecidos. La cercanía se mide con similitud coseno donde valores cercanos a 1.0 indican alta similitud semántica.

Define el tamaño de la ventana de contexto cuántas palabras a cada lado se consideran. Una ventana de 5 significa que se consideran 10 palabras totales alrededor de la palabra objetivo. Ventanas más grandes capturan relaciones semánticas más amplias.

Convierte el texto continuo en unidades discretas tokens o palabras que el modelo puede procesar matemáticamente. Sin tokenización el modelo no puede identificar límites de palabras ni crear el vocabulario.

```
(nlp_env) tuneek@tunek:~/Proyectos/Universidad/InteligenciaArtificial/ProcesamientoDeLenguajeNaturalModerno/Ejercicio1$ python3 Ejercicio1.py
Descargando recursos de NLTK (si es necesario)...
Preparando corpus...
Procesando 500 archivos...
Procesado: 0/500
Procesado: 100/500
Procesado: 200/500
Procesado: 300/500
Procesado: 400/500
Entrenando modelo Word2Vec...
Modelo entrenado exitosamente!
Vocabulario: 11598 palabras

Palabras similares a 'good':
funny: 0.906
enough: 0.889
bad: 0.871
well: 0.870
sure: 0.863

Tiempo de entrenamiento optimizado usando 8 workers
(nlp_env) tuneek@tunek:~/Proyectos/Universidad/InteligenciaArtificial/ProcesamientoDeLenguajeNaturalModerno/Ejercicio1$ |
```

Ejercicio 2 Uso de GloVe y similitud semántica

Preguntas:

- Cuál es la diferencia entre GloVe y Word2Vec en cuanto a su forma de entrenamiento?
- Por qué usamos KeyedVectors en este ejercicio?
- Qué resultados obtuviste al comparar king y queen? Qué interpretas?
- Puedes mencionar un caso donde el análisis semántico con GloVe sería útil en la industria?
- Qué limitaciones tienen los embeddings estáticos como GloVe?

Respuestas Ejercicio 2:

GloVe utiliza estadísticas globales de co-ocurrencia de todo el corpus. Factoriza una matriz de co-ocurrencia palabra-palabra. Word2Vec utiliza ventanas locales de contexto. Entrena con pares palabra-contexto secuenciales.

KeyedVectors es una interfaz optimizada para cargar y usar embeddings preentrenados sin necesidad de reentrenar el modelo. Permite búsquedas eficientes de similitud y analogías.

Similitud de 0.75 alta indicando que el modelo captura correctamente la relación semántica entre conceptos de realidad. Los embeddings entienden roles sociales similares.

E-commerce búsqueda semántica de productos zapatillas deportivas encuentra tennis running shoes. Marketing análisis de sentimientos en redes sociales y agrupación de menciones similares. Recursos Humanos matching automático de CVs con ofertas de trabajo.

Una sola representación por palabra no maneja polisemia banco financiero versus banco asiento. No se actualizan con nuevos contextos. Sesgo inherente del corpus de entrenamiento. No capturan contexto dinámico de la oración.

```
(nlp_env) tuneek@tunek:~/Proyectos/Universidad/InteligenciaArtificial/ProcesamientoDeLenguajeNaturalModerno/Ejercicio2$ python3 Ejercicio2.py
Cargando modelo GloVe... (esto puede tardar unos momentos)
Modelo cargado en 6.54 segundos
Vocabulario: 100000 palabras

--- Análisis de Similitudes ---
Similitud king-queen: 0.7508
Similitud cat-banana: 0.2738
Similitud computer-technology: 0.7642
Similitud happy-joy: 0.5189
Similitud car-vehicle: 0.8631

--- Analogías ---
king - man + woman = queen (score: 0.7699)
(nlp_env) tuneek@tunek:~/Proyectos/Universidad/InteligenciaArtificial/ProcesamientoDeLenguajeNaturalModerno/Ejercicio2$ |
```

Ejercicio 3 Embeddings personalizados desde corpus local

Preguntas:

- Por qué podrías preferir entrenar tus propios embeddings en vez de usar GloVe?
- Qué características del texto pueden afectar la calidad de los embeddings?
- Cómo se refleja el dominio del texto en los vectores obtenidos?
- Qué cambios harías para mejorar la calidad de tus embeddings?
- Qué usos prácticos tendría este modelo dentro de una empresa?

Respuestas Ejercicio 3:

Vocabulario específico captura jerga y terminología particular de tu dominio. Contexto relevante relaciones semánticas específicas de tu industria aplicación. Datos actualizados refleja tendencias y términos actuales no datos históricos. Control total puedes ajustar parámetros según tus necesidades específicas.

Limpieza textos con ruido HTML caracteres especiales degradan calidad. Tamaño del corpus muy pequeño produce embeddings poco robustos. Diversidad falta de variedad contextual limita las representaciones. Frecuencia palabras muy raras o muy comunes pueden ser problemáticas.

Los términos técnicos del dominio IA ML muestran alta similitud entre sí. Intelligence y artificial tienen 0.89 de similitud reflejando su co-ocurrencia frecuente en el corpus especializado.

Más datos incluir más artículos del dominio. Mejor limpieza preprocesamiento más sofisticado lematización eliminación de stopwords. Ajustar parámetros experimentar con dimensiones del vector tamaño de ventana. Filtrado eliminar palabras muy raras o muy comunes.

Búsqueda documental encontrar documentos técnicos similares. Clustering agrupar tickets de soporte por tema. Recomendaciones sugerir artículos relacionados a empleados. Análisis de tendencias identificar temas emergentes en feedback.

```

(nlp_env) tuneek@tunek:~/Proyectos/Universidad/InteligenciaArtificial/ProcesamientoDeLenguajeNaturalModerno/Ejercicio3$ python3 Ejercicio3.py
Verificando recursos de NLTK...
Cargando dataset...
Dataset cargado: 279577 artículos
Procesando 500 artículos...
Procesado: 0/500
Procesado: 100/500
Procesado: 200/500
Procesado: 300/500
Procesado: 400/500
Tokenizando corpus...
Corpus preparado: 100 oraciones
Entrenando modelo con 8 workers...
Modelo entrenado exitosamente!
Vocabulario: 4220 palabras

Palabras similares a 'intelligence':
  artificial: 0.896
  electricity: 0.859
  mastery: 0.857

Palabras similares a 'artificial':
  continous: 0.898
  intelligence: 0.896
  electricity: 0.867

Palabras similares a 'algorithm':
  proprietary: 0.795
  monitored: 0.792
  vicinity: 0.776

Palabras similares a 'machine':
  validation: 0.862
  intuition: 0.858
  abstractions: 0.853

Palabras similares a 'learning':
  machine: 0.806
  predictive: 0.783
  darpa: 0.776

Palabras similares a 'data':
  motivation: 0.647
  ...

Palabras similares a 'data':
  motivation: 0.647
  sas: 0.646
  server: 0.641

Muestra del vocabulario entrenado:
['the', 'and', 'that', 'this', 'for', 'will', 'with', 'are', 'can', 'from']
(nlp_env) tuneek@tunek:~/Proyectos/Universidad/InteligenciaArtificial/ProcesamientoDeLenguajeNaturalModerno/Ejercicio3$

```

Ejercicio 4 Clasificación de texto con BERT

Preguntas:

Cuál es el propósito del proceso de fine-tuning en BERT?

Qué diferencias encontraste entre entrenar con un subconjunto pequeño vs. el dataset completo?

Qué hace el tokenizer en el pipeline de Hugging Face?

Qué métrica usarías para evaluar este modelo?

Por qué es más eficaz BERT que un modelo tradicional como Naive Bayes para clasificación de texto?

Respuestas Ejercicio 4:

Adaptar un modelo preentrenado a una tarea específica clasificación de sentimientos usando datos del dominio objetivo.

Con 100 muestras se logró 100 por ciento accuracy pero con riesgo de overfitting. Un dataset completo daría resultados más generalizables.

Convierte texto en tokens numéricos que BERT puede procesar incluyendo padding y truncación a longitud fija.

Accuracy como se usó pero también F1-score precision y recall para evaluar balance entre clases.

BERT entiende contexto bidireccional y relaciones complejas entre palabras mientras Naive Bayes solo considera frecuencias independientes.

Ejercicio 5 Resumen automático de texto

Cómo se diferencia el resumen extractivo del resumen abstractivo?

¿Qué limitaciones encontraste en los resúmenes generados?

En qué aplicaciones reales sería útil esta técnica?

Ejercicio 6 Análisis de sentimientos con DistilBERT

1 1

Qué cambios podrías hacer para adaptarlo a un nuevo idioma?

Respuestas Ejercicio 6:

Tamaño 60 por ciento más pequeño 66M versus 110M parámetros. Velocidad 60 por ciento más rápido en inferencia. Rendimiento mantiene 97 por ciento del rendimiento de BERT. Memoria menor consumo de GPU CPU.

E-commerce análisis automático de reviews de productos. Redes sociales monitoreo de marca y reputación online. Atención al cliente priorización automática de tickets negativos. Investigación de mercado análisis de opiniones sobre productos servicios.

Confianza muy alta mayor a 99 por ciento para frases con sentimiento claro. Frases neutras como The movie was okay mostraron menor confianza pero correcta clasificación.

Sarcasmo e ironía Qué excelente servicio sarcástico. Contexto cultural referencias locales o culturales específicas. Emociones mixtas opiniones que combinan aspectos positivos y negativos. Negación compleja dobles negaciones o negaciones sutiles.

Modelos multilingües usar mBERT o XLM-RoBERTa. Fine-tuning entrenar con datos específicos del idioma. Traducción traducir texto al inglés analizar y mapear resultado. Modelos nativos usar modelos preentrenados específicos del idioma.

Ejercicio 7 Fine-tuning para QA con transformers

Preguntas:

Qué hace el modelo para identificar la respuesta dentro del contexto?

Por qué es útil tener un modelo preentrenado en SQuAD?

Qué tan preciso fue el modelo en tus pruebas?

Qué desafíos enfrentarías si quisieras entrenar tu propio modelo de QA?

Puedes imaginar una aplicación de esta técnica en tu entorno profesional?

```
(nlp_env) tune@tune:~/Proyectos/Universidad/InteligenciaArtificial/ProcesamientoDeLenguajeNaturalModerno/Ejercicio6$ python3 Ejercicio6.py
Cargando modelo de análisis de sentimientos en GPU...
No model was supplied, defaulted to distilbert/distilbert-base-uncased-finetuned-sst-2-english and revision 714eb0f (https://huggingface.co/distilbert/distilbert-base-uncased-finetuned-sst-2-english).
Using a pipeline without specifying a model name and revision in production is not recommended.
Device set to use cuda:0
Modelo cargado exitosamente!

--- Análisis de Sentimientos en Lote ---
"This new laptop is amazing!"
Sentimiento: POSITIVE (confianza: 1.000)

"This was the worst customer service ever."
Sentimiento: NEGATIVE (confianza: 1.000)

"The movie was okay, nothing special."
Sentimiento: NEGATIVE (confianza: 0.993)

"I absolutely love this product!"
Sentimiento: POSITIVE (confianza: 1.000)

"The weather is nice today."
Sentimiento: POSITIVE (confianza: 1.000)

"I'm feeling sad about the news."
Sentimiento: NEGATIVE (confianza: 0.999)

"This restaurant has excellent food!"
Sentimiento: POSITIVE (confianza: 1.000)

"The service was slow and disappointing."
Sentimiento: NEGATIVE (confianza: 1.000)

(nlp_env) tune@tune:~/Proyectos/Universidad/InteligenciaArtificial/ProcesamientoDeLenguajeNaturalModerno/Ejercicio6$ |
```

Respuestas Ejercicio 7:

El modelo identifica spans segmentos de texto que probablemente contienen la respuesta. Predice tokens de inicio y fin calculando probabilidades para cada posición posible en el contexto.

SQuAD contiene 100000 más preguntas-respuestas de alta calidad anotadas por humanos. El preentrenamiento proporciona comprensión de patrones pregunta-respuesta. Transfer learning para nuevos dominios. Base sólida para fine-tuning específico.

Precisión variable según complejidad. Respuestas directas fechas nombres alta precisión. Conceptos abstractos precisión moderada. Preguntas ambiguas menor confianza pero respuestas relevantes.

Datos anotados necesitas miles de pares pregunta-respuesta anotados manualmente. Costo computacional entrenamiento requiere GPU potentes durante días semanas. Expertise conocimiento profundo de arquitecturas transformer. Evaluación métricas complejas F1 EM para validar calidad.

Soporte técnico FAQ automático que responde preguntas de usuarios. Documentación sistema de búsqueda inteligente en manuales técnicos. Legal búsqueda de precedentes y respuestas en documentos legales. Educación asistente virtual para estudiantes con dudas académicas.

Ejercicio 8 Chatbot con Hugging Face más Gradio

Preguntas:

Qué diferencia a un chatbot basado en reglas de uno basado en modelos generativos como DialoGPT?

Cómo maneja el modelo el historial de la conversación?

Qué problemas encuentras en la coherencia de las respuestas?

Qué harías para mejorar la fluidez y precisión del chatbot?

Qué otros modelos podrías probar en lugar de DialoGPT?

```
(nlp_env) tune@tune:~/Proyectos/Universidad/InteligenciaArtificial/ProcesamientoDeLenguajeNaturalModerno/Ejercicio7$ python3 Ejercicio7.py
Cargando modelo de Question Answering en GPU...
Device set to use cuda:0
Modelo cargado exitosamente!

--- Sistema de Preguntas y Respuestas ---

1. Pregunta: ¿Cuándo fue fundada la Universidad Técnica de Oruro?
/home/tune/Proyectos/Universidad/InteligenciaArtificial/ProcesamientoDeLenguajeNaturalModerno/nlp_env/lib/python3.12/site-packages/transformers/pipelines/question_answering.py:390: FutureWarning: Passing a list of SQuAD examples to the pipeline is deprecated and will be removed in v5. Inputs should be passed using the 'question' and 'context' keyword arguments instead.
  warnings.warn(
  Respuesta: La Universidad Técnica de Oruro fue fundada en 1892.
  Confianza: 0.1859
  Posición en texto: 0-52
  Contexto: **La Universidad Técnica de Oruro fue fundada en 1892.** Es una de las universidades más antiguas de...

2. Pregunta: ¿En qué se especializa la universidad?
  Respuesta: ingeniería y tecnología
  Confianza: 0.2469
  Posición en texto: 125-148
  Contexto: La Universidad Técnica de Oruro fue fundada en 1892. Es una de las universidades más antiguas de Bol...

3. Pregunta: ¿Qué es la inteligencia artificial?
  Respuesta: una rama de la informática
  Confianza: 0.6299
  Posición en texto: 30-56
  Contexto: La inteligencia artificial es **una rama de la informática** que busca crear sistemas capaces de rea...

4. Pregunta: ¿Quién creó Python y cuándo?
  Respuesta: 1991
  Confianza: 0.1195
  Posición en texto: 83-87
  Contexto: Python es un lenguaje de programación de alto nivel creado por Guido van Rossum en **1991**. Es cono...

Memoria GPU utilizada: 0.13 GB
(nlp_env) tune@tune:~/Proyectos/Universidad/InteligenciaArtificial/ProcesamientoDeLenguajeNaturalModerno/Ejercicio7$ |
```

Respuestas Ejercicio 8:

Basado en reglas respuestas predefinidas limitado a patrones específicos predecible pero rígido. Generativo DialoGPT crea respuestas nuevas más natural y flexible pero puede ser impredecible o generar contenido inapropiado.

Concatena la conversación previa como contexto para la siguiente respuesta. Mantiene memoria de intercambios anteriores hasta el límite de tokens del modelo típicamente 3 a 5 intercambios.

Pérdida de contexto después de varias interacciones. Respuestas repetitivas o genéricas.

Dificultad para mantener personalidad consistente. Puede generar información contradictoria.

Fine-tuning entrenar con datos específicos del dominio. Filtros de seguridad evitar respuestas inapropiadas. Memoria extendida sistemas de memoria a largo plazo. Retrieval-augmented combinar con base de conocimientos externa.

GPT-3.5 4 más capaces pero requieren API. LLaMA modelo open-source más reciente. Claude excelente para conversaciones naturales. Blenderbot especializado en conversaciones casuales.

Ejercicio 9 Proyecto integrador Clasificador más Sentimiento más Resumen

Preguntas:

Qué tarea resultó más precisa el resumen la clasificación o el análisis de sentimiento?

Qué tan bien se adaptaron los modelos preentrenados a tu texto personalizado?
Cómo integrarías este pipeline en una aplicación web real?
Qué parte del pipeline automatizarías o optimizarías con otra herramienta?
Qué mejoras podrías hacer si el texto estuviera en otro idioma o jerga regional?

```
(nlp_env) tune@tunek:~/Proyectos/Universidad/InteligenciaArtificial/ProcesamientoDeLenguajeNaturalModerno/Ejercicio8$ python3 Ejercicio8.py
Cargando modelo DialoGPT en GPU...
Usando modelo: microsoft/DialoGPT-medium
Modelo cargado exitosamente!

=== Chatbot Listo ===
Escribe 'quit' para salir

Tú: hola
The attention mask is not set and cannot be inferred from input because pad token is same as eos token. As a consequence, you may observe unexpected behavior. Please pass your input's 'attention_mask' to obtain reliable results.
Bot: The one true god, you mean.

Tú: como estas?
Bot: Porque no los dos?

Tú: como?
Bot: Porque?

Tú: que?
Bot: Que?

Tú: jaja
Bot: Jaja. No?

Tú: que dices?
Bot: Todo, si.

Tú: wtf
Bot: Amen brother

Tú: |
```

Respuestas Ejercicio 9:

Clasificación temática fue la más precisa 0.88 a 0.92 confianza seguida por análisis de sentimientos 0.67 a 0.91. El resumen siendo generativo es más difícil de evaluar objetivamente.

Muy bien para contenido general. Los modelos preentrenados en datasets diversos manejan correctamente textos sobre tecnología economía y salud. Menor precisión esperada para jerga muy específica.

API REST endpoints separados para cada tarea summarize sentiment classify. Microservicios cada modelo en contenedor independiente. Cache Redis para cachear resultados de textos frecuentes. Load balancing distribuir carga entre múltiples instancias GPU.

Preprocesamiento Apache Kafka para streaming de datos. Batch processing Apache Spark para procesar grandes volúmenes. Monitoring MLflow para tracking de métricas y versiones. Orquestación Kubernetes para scaling automático.

Modelos multilingües mBERT XLM-R mT5 para múltiples idiomas. Detección de idioma identificar automáticamente el idioma del texto. Fine-tuning local entrenar con datos específicos de la región jerga. Traducción pipeline de traducción automática más procesamiento en inglés.


```
(nlp_env) tuneek@tunek:~/Proyectos/Universidad/InteligenciaArtificial/ProcesamientoDeLenguajeNaturalModerno/Ejercicio9$ python3 Ejercicio9.py
=== Pipeline NLP Integrado en GPU ===
Cargando modelos...
Device set to use cuda:0
No model was supplied, defaulted to distilbert/distilbert-base-uncased-finetuned-sst-2-english and revision 714eb0f (https://huggingface.co/distilbert/distilbert-base-uncased-finetuned-sst-2-english).
Using a pipeline without specifying a model name and revision in production is not recommended.
Device set to use cuda:0
No model was supplied, defaulted to facebook/bart-large-mnli and revision d7645e1 (https://huggingface.co/facebook/bart-large-mnli).
Using a pipeline without specifying a model name and revision in production is not recommended.
Device set to use cuda:0
Todos los modelos cargados exitosamente!

=====
TEXTO 1 - Analisis Completo
=====
Texto original (445 caracteres):
El avance de la inteligencia artificial está cambiando el mundo. Los modelos de machine learning est...

RESUMEN:
    Los modelos de machine learning están revolucionando industrias enteras. Sin embargo, también plantean desaf

SENTIMIENTO:
    NEGATIVE (confianza: 0.967)

CLASIFICACIÓN TEMÁTICA:
    Tema principal: tecnologia (confianza: 0.558)
    Temas secundarios:
        - salud: 0.183
        - medio ambiente: 0.109

Tiempo de procesamiento: 3.78 segundos

=====
TEXTO 2 - Analisis Completo
=====
Texto original (462 caracteres):
La situación económica mundial muestra signos de recuperación tras la crisis. Los indicadores financ...

RESUMEN:
    La inflación sigue siendo una preocupación para muchos países. Los bancos central

SENTIMIENTO:
    NEGATIVE (confianza: 0.974)

CLASIFICACIÓN TEMÁTICA:
    Tema principal: economia (confianza: 0.415)
    Temas secundarios:
        - politica: 0.188
        - medio ambiente: 0.179

Tiempo de procesamiento: 3.53 segundos

=====
TEXTO 3 - Analisis Completo
=====
Texto original (448 caracteres):
Los nuevos descubrimientos médicos ofrecen esperanza para el tratamiento de enfermedades raras. Los ...

RESUMEN:
    Los nuevos descubrimientos médicos ofrecen esperanza para el trat

SENTIMIENTO:
    NEGATIVE (confianza: 0.930)

CLASIFICACIÓN TEMÁTICA:
    Tema principal: salud (confianza: 0.413)
    Temas secundarios:
        - tecnologia: 0.179
        - medio ambiente: 0.129

Tiempo de procesamiento: 3.50 segundos

Memoria GPU utilizada: 1.66 GB

Pipeline NLP completo ejecutado exitosamente!
(nlp_env) tuneek@tunek:~/Proyectos/Universidad/InteligenciaArtificial/ProcesamientoDeLenguajeNaturalModerno/Ejercicio9$ |
```