

Introducción y Problema

En los últimos años, los sistemas de recomendación han evolucionado hacia arquitecturas complejas basadas en aprendizaje profundo. Sin embargo, modelos más simples como los autoencoders lineales (LAEs) han mostrado resultados competitivos debido a su eficiencia computacional y menor riesgo de sobreajuste. A pesar de sus ventajas, los LAEs enfrentan dos limitaciones críticas: el sesgo de popularidad y el sesgo de vecinos cercanos, lo que impide modelar desemejanzas o correlaciones negativas. Estos sesgos afectan la diversidad y calidad de las recomendaciones generadas. Por otro lado, los sistemas de recomendación *memory-based* son altamente interpretables, pero típicamente presentan métricas de similitud únicamente positivas, lo que impide modelar desemejanzas o correlaciones negativas.

Motivación

A pesar del auge de modelos complejos basados en aprendizaje profundo, los enfoques lineales y memory-based siguen siendo competitivos debido a su eficiencia y facilidad de interpretación. Sin embargo, la mayoría de las investigaciones se enfocan en compararlos por separado. Nuestra motivación radica en explorar cómo estos paradigmas pueden colaborar: ¿Es posible combinar la normalización adaptativa de DAN con la capacidad de Sapling para capturar disimilaridades? ¿Qué sucede si las salidas de un modelo se usan como entradas del otro?

Propuesta

Este proyecto propone combinar ambas metodologías de distintas maneras, utilizando *pipelines* pasando la matriz que genera Sapling Similarity a DAN, y viceversa, generando un modelo ponderado entre ambos métodos y replicando ambos métodos en su forma original. Esto se aplicará en distintos datasets para revisasr si es posible obtener las ventajas que ofrecen ambos métodos o alguno prevalece sobre otro, esto mediante el ajuste de los hiperparámetros de cada algoritmo.

Algoritmos

Sapling Similarity, técnica propuesta por Albora et al. (2023), permite modelar tanto similitudes como disimilaridades:

- Bean (base): muestra cuántos ítems conecta el usuario i
- Hojas: separan los ítems según si el usuario j está o no conectado
- Si saber que j está conectado aumenta la probabilidad de que i también lo esté → similitud positiva
- Si la reduce → similitud negativa

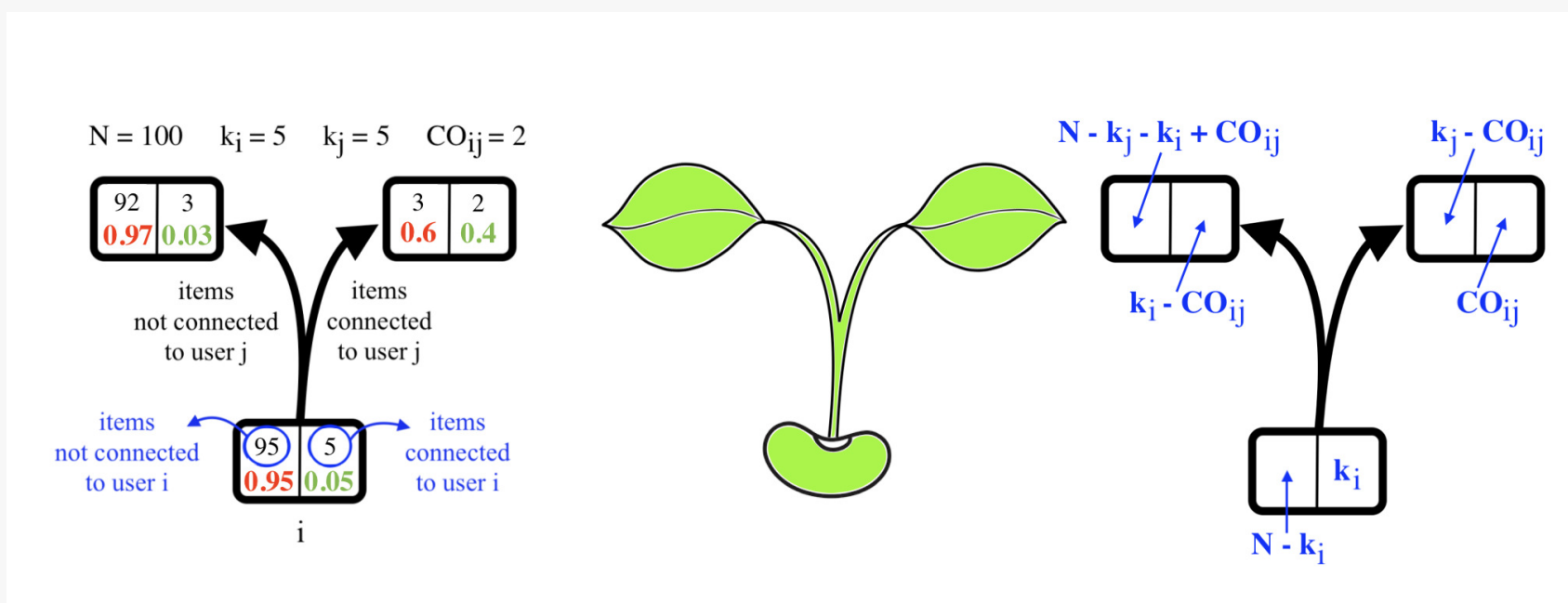


Figura 1: Funcionamiento de Sapling Similarity

Data-Adaptive Normalization (DAN) comienza desde EASE, que resuelve:

$$\mathbf{W} = (\mathbf{X}^T \mathbf{X} + \lambda \mathbf{I})^{-1} \mathbf{X}^T \mathbf{X}, \quad \text{con } \text{diag}(\mathbf{W}) = 0$$

DAN modifica esto:

- Escala las filas de \mathbf{X} según la frecuencia de usuarios (β) y columnas según popularidad de ítems (α).
- Introduce términos de relajación adaptativos.

$$\mathbf{G} = \tilde{\mathbf{X}}^T \tilde{\mathbf{X}} + \lambda + \frac{\text{drop}_p}{1 - \text{drop}_p}$$
$$\mathbf{W} = (\mathbf{I} - \mathbf{P} \cdot \lambda) \cdot \frac{1}{\text{item_scale}}$$

Metodología

Se implementaron 5 enfoques de recomendación:

1. **DAN**: $\mathbf{R}_u = \mathbf{x}_u \cdot \mathbf{W}$
2. **Sapling Similarity**:
$$\mathbf{R} = (1 - \gamma) \cdot \mathbf{S}_U \cdot \mathbf{X} + \gamma \cdot \mathbf{X} \cdot \mathbf{S}_I$$
3. **Pipeline**: usa salida de DAN como entrada a Sapling
4. **Combine**: combinación lineal:

$$\mathbf{R} = (1 - \alpha) \cdot \mathbf{R}_{DAN} + \alpha \cdot \mathbf{R}_{Sapling}$$

5. **Regularization**:

$$\mathbf{G} \leftarrow \mathbf{G} + \eta (\mathbf{I} - \mathbf{S}_{Sapling})$$

Evaluado con **Precision@k**, **Recall@k** y **nDCG@k**.

Resultados

amazon_products			
model	precision@20	recall@20	ndcg@20
combine	0.0001	0.0007	0.0003
dan	0.0001	0.0003	0.0001
pipeline	0.0001	0.0013	0.0003
regularization	0.0002	0.0029	0.0008
sapling	0.0013	0.0234	0.0083

amazon_reviews_books			
model	precision@20	recall@20	ndcg@20
combine	0.0000	0.0000	0.0000
dan	0.0003	0.0000	0.0004
pipeline	0.0005	0.0021	0.0008
regularization	0.0003	0.0034	0.0009
sapling	0.0003	0.0037	0.0012

yelp_reviews			
model	precision@20	recall@20	ndcg@20
combine	0.0000	0.0000	0.0000
dan	0.0000	0.0000	0.0000
pipeline	0.0001	0.0027	0.0006
regularization	0.0035	0.0708	0.0331
sapling	0.0038	0.0763	0.0279

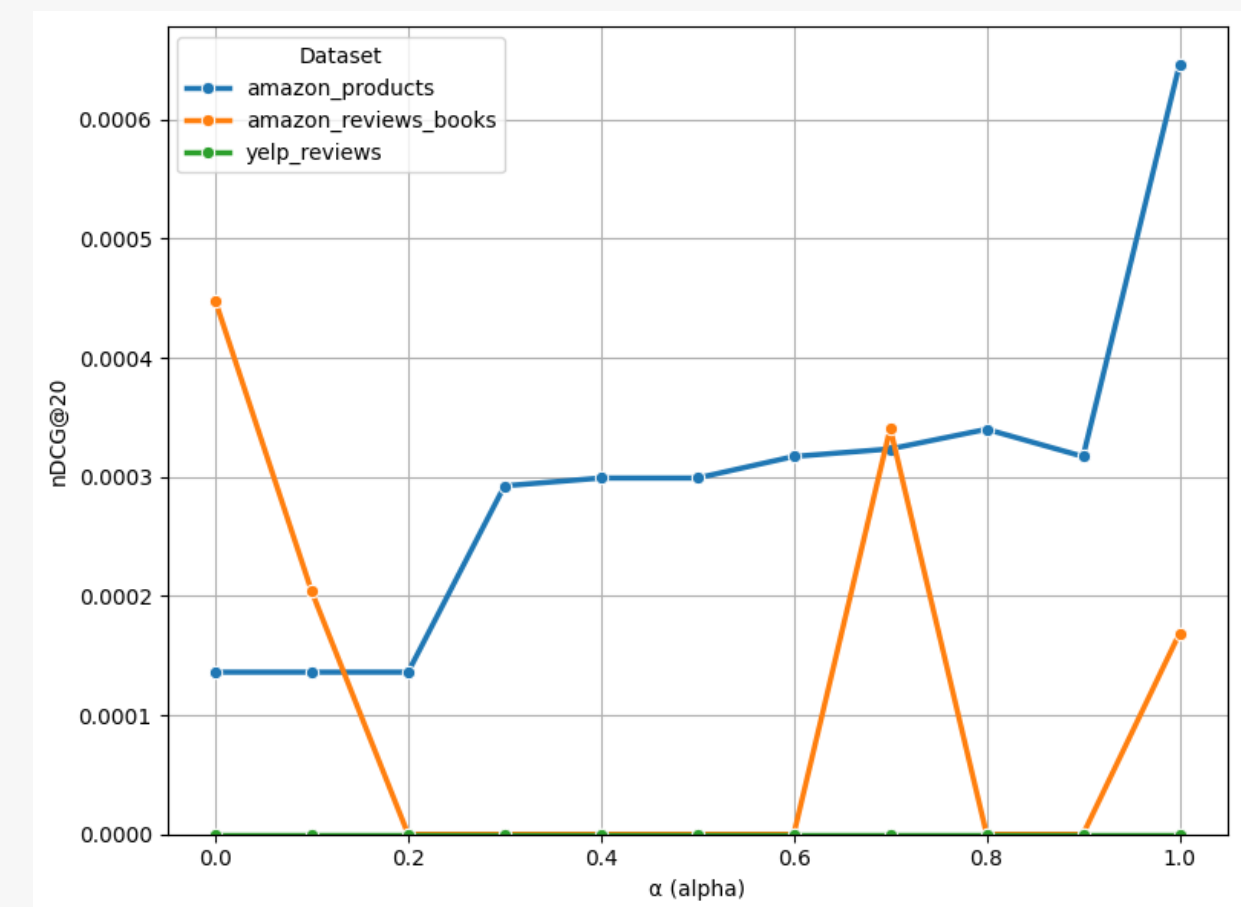


Figura 2: Sensitivity Analysis

Discusión

- **Sapling Similarity** fue el más robusto en los tres datasets.
- **Regularization** destacó en *yelp_reviews*, posiblemente por refuerzo local.
- **Combine** y **Pipeline** no superaron los métodos base.
- **DAN** requiere una estructura más rica para mejorar su rendimiento.

Conclusiones:

- La integración de DAN con Sapling es prometedora si se normalizan bien.
- Los modelos que modelan disimilaridad explícitamente fueron mejores.

Trabajo futuro:

- Probar normalización con GCN.
- Variantes temporales de Sapling.
- Agregar información multimodal a \mathbf{X} .