

An Entropy-Based Approach for preserving Diversity in Evolutionary Topical Search

CECILIA BAGGIO – ROCÍO L. CECCHINI CARLOS M. LORENZETTI – ANA MAGUITMAN

Objetivo del trabajo

Diseñar nuevas técnicas capaces de refinar consultas de manera automática para acumular recursos relevantes y diversos respecto a un contexto temático.

- Diseñar nuevas técnicas capaces de refinar consultas de manera automática para acumular recursos relevantes y diversos respecto a un contexto temático.
- Problema común en enfoques anteriores: uso de estrategias evolutivas que resultan en un Recall pobre debido a la pérdida de diversidad genética.

- Diseñar nuevas técnicas capaces de refinar consultas de manera automática para acumular recursos relevantes y diversos respecto a un contexto temático.
- Problema común en enfoques anteriores: uso de estrategias evolutivas que resultan en un Recall pobre debido a la pérdida de diversidad genética.
- Proponemos una nueva estrategia inspirada en la noción de entropía basada en teoría de la información para favorecer la diversidad de la población con el objetivo de alcanzar buen recall global.

Utilidad de generar buenas consultas

Búsqueda basada en la tarea del usuario

Búsqueda basada en la tarea del usuario

Utilidad de generar buenas consultas

Recolección de recursos web para portales temáticos

Utilidad de generar buenas consultas

- Búsqueda basada en la tarea del usuario
- Recolección de recursos web para portales temáticos
- Búsqueda en la Deep-web

Utilidad de generar buenas consultas

- Búsqueda basada en la tarea del usuario
- Recolección de recursos web para portales temáticos
- Búsqueda en la Deep-web
- Soporte para gestión de conocimiento

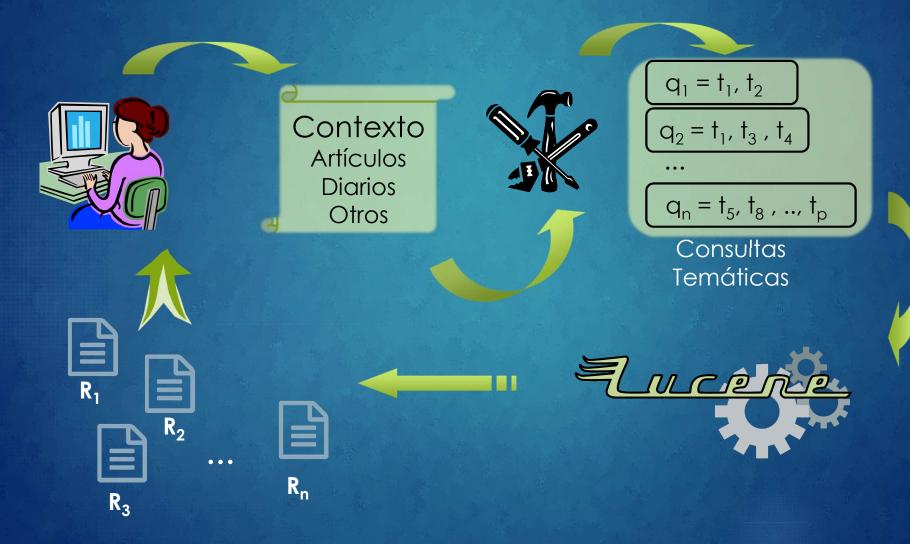


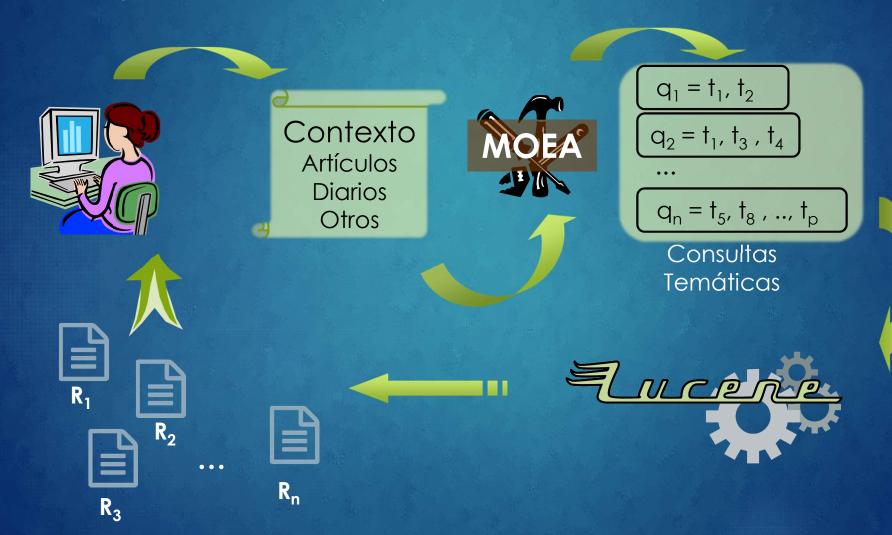


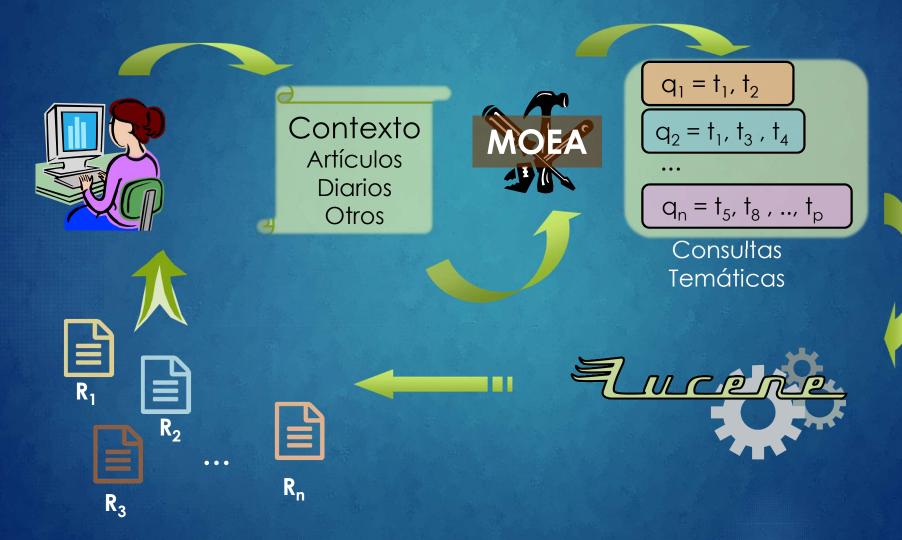












Espacio multidimensional de búsqueda.

- Espacio multidimensional de búsqueda.
- Soluciones subóptimas pueden considerarse efectivas.

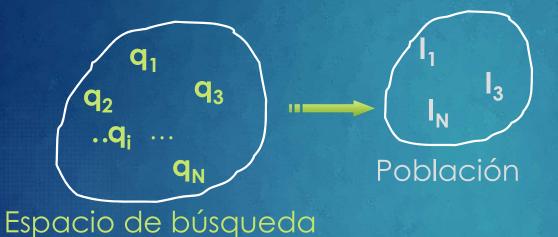
- Espacio multidimensional de búsqueda.
- Soluciones subóptimas pueden considerarse efectivas.
- Múltiples soluciones. Podemos tener:
 - Resultados satisfactorios en distintos individuos.
 - Interés en encontrar más de una consulta.

- Espacio multidimensional de búsqueda.
- Soluciones subóptimas pueden considerarse efectivas.
- Múltiples soluciones. Podemos tener:
 - Resultados satisfactorios en distintos individuos.
 - Interés en encontrar más de una consulta.
- Exploración y Explotación.

- Espacio multidimensional de búsqueda.
- Soluciones subóptimas pueden considerarse efectivas.
- Múltiples soluciones. Podemos tener:
 - Resultados satisfactorios en distintos individuos.
 - Interés en encontrar más de una consulta.
- Exploración y Explotación.
- Múltiples objetivos y posiblemente conflictivos entre sí.

Población

Representación



Población

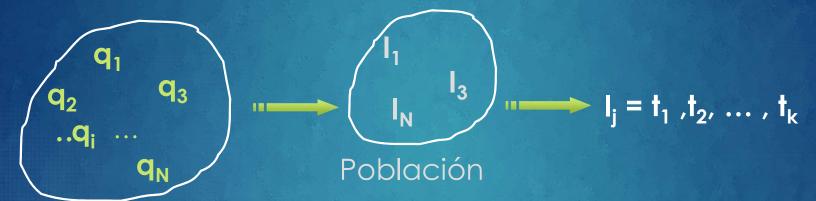
Representación



Espacio de búsqueda

Población

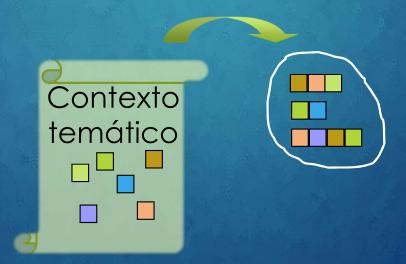
Representación



Espacio de búsqueda

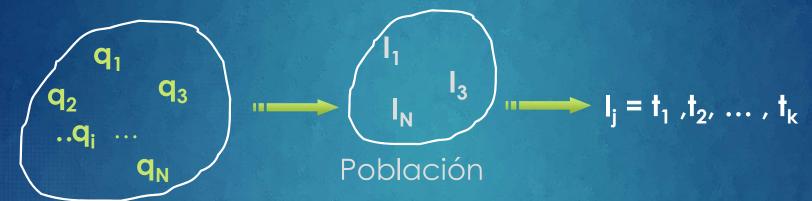
Población

Inicialización



Población

Representación



Espacio de búsqueda

Población

Inicialización



N° de términos?

Función de Fitness

Función de Fitness

Sea Q un espacio de búsqueda y T un conjunto de tópicos. Definimos:

Precision@10(q,t) =
$$\frac{|A_q^{10} \cap D_t|}{|A_q^{10}|}$$

$$Recall(q,t) = \frac{|A_q \cap D_t|}{|D_t|}$$

- D_t : conjunto que contiene a todos los documentos asociados al tópico t.
- \blacktriangleright A_q : conjunto de resultados devuelto por el motor de búsqueda cuando ${f q}$ se usa como consulta.
- ▶ A_q^{10} : es el conjunto de los primeros 10 documentos del ranking A_q .

Función de Fitness

Sea Q un espacio de búsqueda, T un conjunto de tópicos y P una población de consultas:

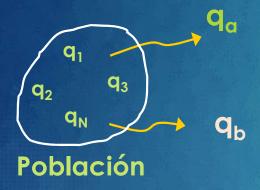
$$Recall_Entr\'opico(q,t,P) = \frac{\sum_{d_i \in A_q \cap D_t} (\mathrm{IQF}(d_i,P)/\log(|P|+1))}{|D_t|}$$

$$IQF(d_i, P) = log((|P| + 1)/n_i)$$

Es la frecuencia inversa en la población de consultas P del documento d_i . Representa la especificidad del documento d_i , donde n_i es el número de consultas que recuperan al documento d_i

Operadores Genéticos

Recombinación

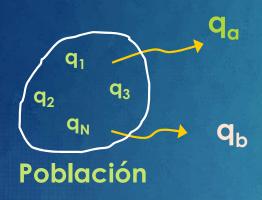


$$q_a = t_1, t_2, \dots, t_k$$

$$q_b = t_1, t_2, t_3, ..., t_m$$

Operadores Genéticos

Recombinación



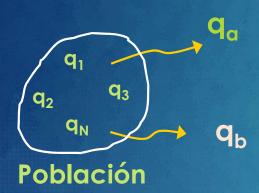
$$q_{\alpha} = t_1, t_2, ..., t_k$$
 $q'_{\alpha} = t_1, t_2, t_3, ..., t_m$

$$q_b = t_1, t_2, t_3, ..., t_m$$

$$q'_b = t_1, t_2, ..., t_k$$

Operadores Genéticos

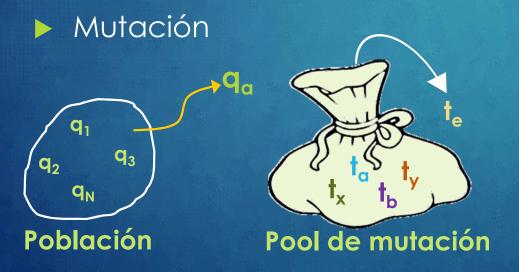
Recombinación



$$q_a = t_1, t_2, ..., t_k$$
 $q'_a = t_1, t_2, t_3, ..., t_m$

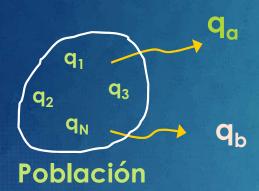
$$q_b = t_1, t_2, t_3, ..., t_m$$

$$q'_b = t_1, t_2, ..., t_k$$



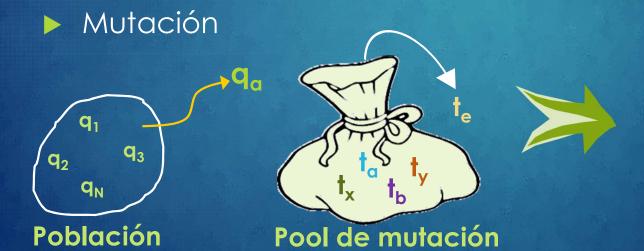
Operadores Genéticos

Recombinación



$$q_a = t_1, t_2, ..., t_k$$
 $q_b = t_1, t_2, t_3, ..., t_m$
 $q'_a = t_1, t_2, t_3, ..., t_m$
 $q'_b = t_1, t_2, ..., t_k$

$$q_b = t_1, t_2, t_3, ..., t_m$$
 $q'_b = t_1, t_2, ..., t_k$



Punto de mutación

$$q_a = t_1, t_2, t_3, t_4, ..., t_p$$
 $q'_a = t_1, t_2, t_e, t_4, ..., t_p$

Operadores Genéticos

Selección

Operadores Genéticos

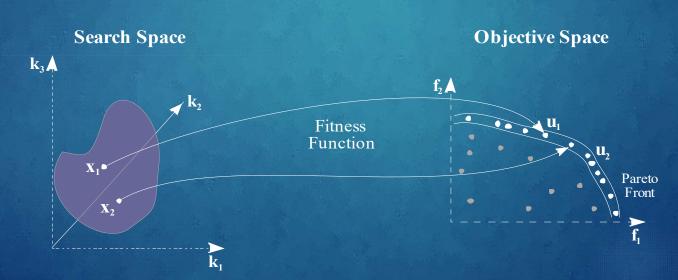
Selección

Método **NSGA-II** (Non-dominated Sorting Genetic Algorithm) para evolucionar consultas guiado por los los objetivos dados en las combinaciones:

Co1: Precision@10 y Recall

Co2: Precision@10 y Entropic-Recall

Basado en el concepto de dominancia de Pareto



An Entropy-Based Approach for preserving Diversity in Evolutionary Topical Search

JAIIO 2016 - Buenos Aires

Análisis de la performance sobre conjuntos de entrenamiento y prueba

Métricas utilizadas

Dado un tópico t, una población de consultas P, una consulta $q \in P$ y el conjunto $A_q = \{a_1, ..., a_n\}$ de recursos recuperados para q:

Métricas utilizadas

Dado un tópico t, una población de consultas P, una consulta $q \in P$ y el conjunto $A_q = \{a_1, ..., a_n\}$ de recursos recuperados para q:

$$\overline{Precision@10}(P,t) = \frac{\sum_{q \in P} Precision@10(q,t)}{|P|}$$

Métricas utilizadas

Dado un tópico t, una población de consultas P, una consulta $q \in P$ y el conjunto $A_q = \{a_1, ..., a_n\}$ de recursos recuperados para q:

$$Precision@10(P,t) = \frac{\sum_{q \in P} Precision@10(q,t)}{|P|}$$

► Global_Recall(P,t) =
$$\frac{|A(P) \cap D_t|}{|D_t|}$$
donde $A(P) = \bigcup_{q_i \in P} A_{q_i}$

Métricas utilizadas

Métricas utilizadas

- $\overline{Jaccard_Similarity_Index}(P) = \\ \underline{\sum_{q_i,qj\in P,\,i\neq j} Jaccard_Similarity_Index(A_{q_i}^*,A_{q_j}^*)} \\ |P|.(|P|-1)$
- Donde:

$$Jaccard_Similarity_Index(A_{q_i}^*, A_{q_j}^*) = \frac{\left|A_{q_i}^* \cap A_{q_j}^*\right|}{\left|A_{q_i}^* \cup A_{q_j}^*\right|}$$

$$\forall A_{q_i}^* = A_{q_i} \cap D_t$$

La similitud promedio de Jaccard se calcula promediando la similitud de Jaccard entre todos los pares de conjuntos de documentos recuperados, restringido solo a los documentos relevantes.

JAIIO 2016 - Buenos Aires

Corpus de Documentos y Diseño de Experimentos

- 448 tópicos del tercer nivel de DMOZ.
 - Reportamos resultados para los tópicos Bodypainting y Aquaculture.
- Más de 100 URLs por tópico.
- Más de 350.000 páginas.
- Lucene fue usado para indexar los documentos y para crear un motor de búsqueda.
- Cada tópico fue dividido en 2/3 para entrenamiento y 1/3 para testing.
- Para cada uno de las pruebas se realizaron 20 corridas con los siguientes parámetros genéticos:
 - \rightarrow n = 100, g = 150
 - Probabilidad de cruzamiento = 0,7
 - Probabilidad de mutación = 0,03
 - ► Longitud máxima inicial de las consultas = 32 palabras



■ Training

Testing

Evaluación de la performance durante el entrenamiento

14

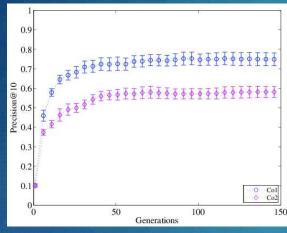
Evaluación de la performance durante el entrenamiento

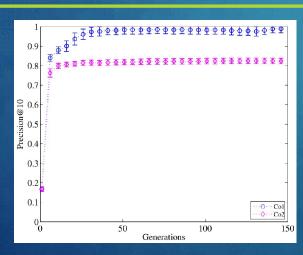
Evolución

Precision@10

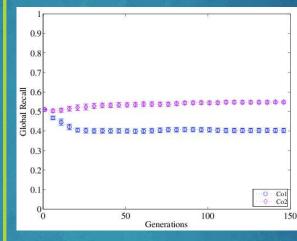
Tópico: Body Painting

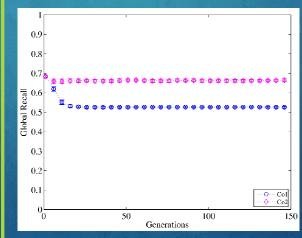
> Tópico: Aquaculture





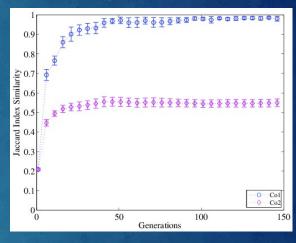


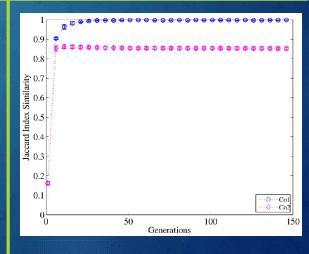




Evolución







Evaluación sobre un nuevo corpus

Evaluación sobre un nuevo corpus

Cálculo de los I.C (al 95%) de las métricas de Precision@10, Global Recall y Mean Jaccard Index para cada tópico sobre el conjunto de test, usando las poblaciones de consultas evolucionadas en comparación con las consultas obtenidas a partir de la descripción del tópico.

TESTING C.I.				
Topics	Metrics	Co1 - Co2 First Generation	Co1 Last Generation	Co2 Last Generation
GlobalRecall	[0.482, 0.482]	[0.362, 0.388]	[0.480, 0.525]	
Jaccard-Similarity-Index	[0.237, 0.272]	[0.975, 0.999]	[0.580, 0.662]	
AQUACULTURE	Precision@10	[0.158, 0.186]	[0.773, 0.845]	[0.735, 0.775]
	GlobalRecall	[0.696, 0.708]	[0.511, 0.529]	[0.672, 0.686]
	Jaccard-Similarity-Index	[0.191, 0.210]	[0.985, 0.997]	[0.822, 0.835]

Hemos propuesto una estrategia para evolucionar consultas temáticas con énfasis en preservar la diversidad.

- Hemos propuesto una estrategia para evolucionar consultas temáticas con énfasis en preservar la diversidad.
- Formulamos una función de aptitud inspirada en la noción de entropía basada en teoría de la información que mejora la cobertura a nivel poblacional.

- Hemos propuesto una estrategia para evolucionar consultas temáticas con énfasis en preservar la diversidad.
- Formulamos una función de aptitud inspirada en la noción de entropía basada en teoría de la información que mejora la cobertura a nivel poblacional.
- Esta función de aptitud es utilizada en combinación con precisión a 10, dando lugar a una estrategia MOEA que alcanza buena cobertura global con pérdida moderada de la precisión.

- Hemos propuesto una estrategia para evolucionar consultas temáticas con énfasis en preservar la diversidad.
- Formulamos una función de aptitud inspirada en la noción de entropía basada en teoría de la información que mejora la cobertura a nivel poblacional.
- Esta función de aptitud es utilizada en combinación con precisión a 10, dando lugar a una estrategia MOEA que alcanza buena cobertura global con pérdida moderada de la precisión.
- Trabajo futuro:
 - Experimentos con ≠ valores para los parámetros del AE.
 - Método de selección.
 - Aplicación de Programación Genética.
 - Función de aptitud

¡Muchas Gracias!

CECILIA BAGGIO – ROCÍO L. CECCHINI CARLOS M. LORENZETTI – ANA MAGUITMAN

{cb, rlc, cml, agm}@cs.uns.edu.ar







¡Muchas Gracias!



CECILIA BAGGIO – ROCÍO L. CECCHINI CARLOS M. LORENZETTI – ANA MAGUITMAN

{cb, rlc, cml, agm}@cs.uns.edu.ar







Arquitectura propuesta

