

# Similarity Learning via Boosting

---

Thiago Rodrigo Ramos

05 de Julho de 2024



I completed my PhD at IMPA, at the Centro PI (Center for Projects and Innovation)



My thesis is result of three applications of Machine Learning to real industrial problems

- Stone Pagamentos (2020): Credit Scoring
  - ExactBoost: directly boosting the margin in combinatorial and non-decomposable metrics
- Dasa (2021): Uncertainty Quantification
  - Split conformal prediction for dependent data
- Rede Globo (2022): Record Linkage
  - Similarity Learning via Boosting

My thesis is result of three applications of Machine Learning to real industrial problems

- Stone Pagamentos (2020): Credit Scoring
  - ExactBoost: directly boosting the margin in combinatorial and non-decomposable metrics
- Dasa (2021): Uncertainty Quantification
  - Split conformal prediction for dependent data
- Rede Globo (2022): Record Linkage
  - Similarity Learning via Boosting

My thesis is result of three applications of Machine Learning to real industrial problems

- Stone Pagamentos (2020): Credit Scoring
  - ExactBoost: directly boosting the margin in combinatorial and non-decomposable metrics
- Dasa (2021): Uncertainty Quantification
  - Split conformal prediction for dependent data
- Rede Globo (2022): Record Linkage
  - Similarity Learning via Boosting

My thesis is result of three applications of Machine Learning to real industrial problems

- Stone Pagamentos (2020): Credit Scoring
  - ExactBoost: directly boosting the margin in combinatorial and non-decomposable metrics
- Dasa (2021): Uncertainty Quantification
  - Split conformal prediction for dependent data
- Rede Globo (2022): Record Linkage
  - Similarity Learning via Boosting



Figure: This scene from “Tropa de Elite” was filmed on IMPA



Figure: Filme “Ricos de amor” da Netflix



Multiple very large datasets with movies information  
(e.g. IMDB, TMDb, Rotten Tomatoes)



# The problem

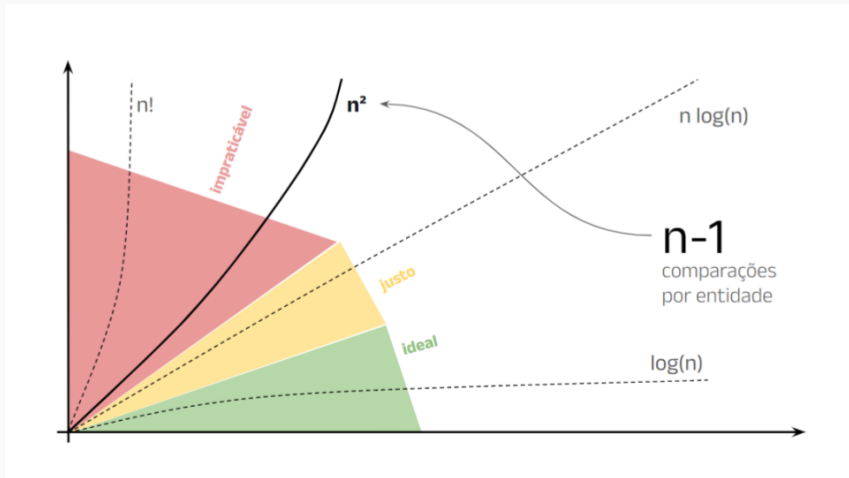
How can we match similar entries from potentially large datasets to create a richer dataset?



## Naive solution: Check all pairs

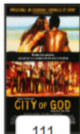


## Naive solution: Check all pairs

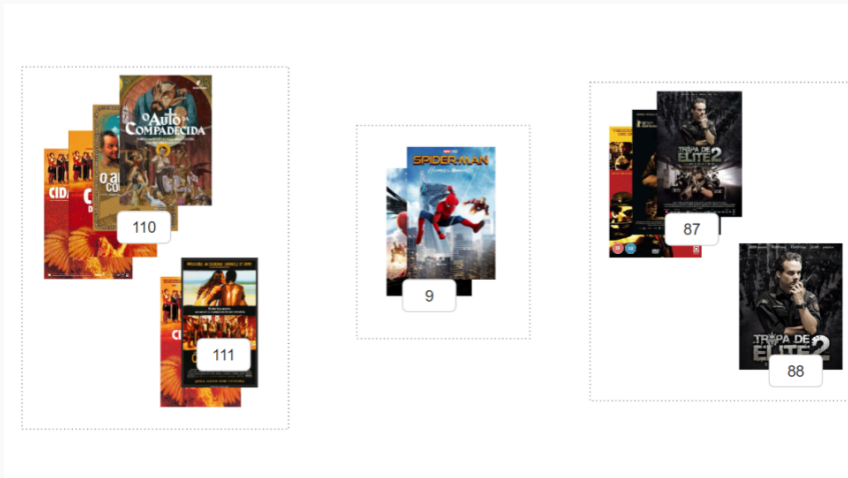




One possible solution is to define a hash code for each entry



One possible solution is to define a hash code for each entry and then block similar movies together



Given datasets  $\mathcal{A}, \mathcal{B}$  for  $A \in \mathcal{A}$  we want to find similar items  $B \in \mathcal{B}$  while doing as few pairwise comparisons as possible.

$$\text{Recall} := \frac{1}{|\mathcal{M}|} \sum_{(\ell, r) \in \mathcal{M}} \mathbf{1}_{[A_\ell \text{ and } B_r \text{ share a block}]};$$

$$\text{RR} := 1 - \frac{1}{|\mathcal{N}|} \sum_{(\ell, r) \in \mathcal{N}} \mathbf{1}_{[A_\ell \text{ and } B_r \text{ share a block}]}$$

$$\text{H} := 2 \frac{\text{Recall} \cdot \text{RR}}{\text{Recall} + \text{RR}}$$

where  $\mathcal{N} := [N_{\mathcal{A}}] \times [N_{\mathcal{B}}]$  denotes all possible pairs and  $\mathcal{M}$  denotes the set of matching pairs:

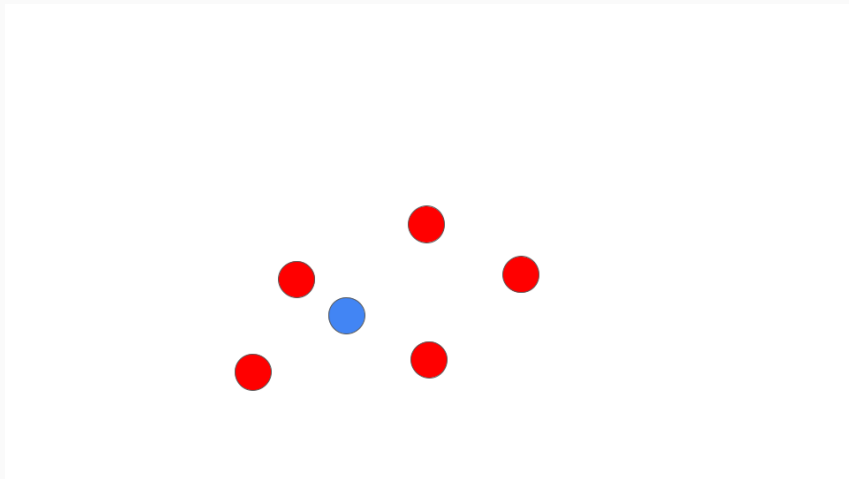
$$\mathcal{M} := \{(\ell, r) \in \mathcal{N}, A_\ell \sim_R B_r, (A_\ell, B_r) \in \mathcal{A} \times \mathcal{B}\}.$$



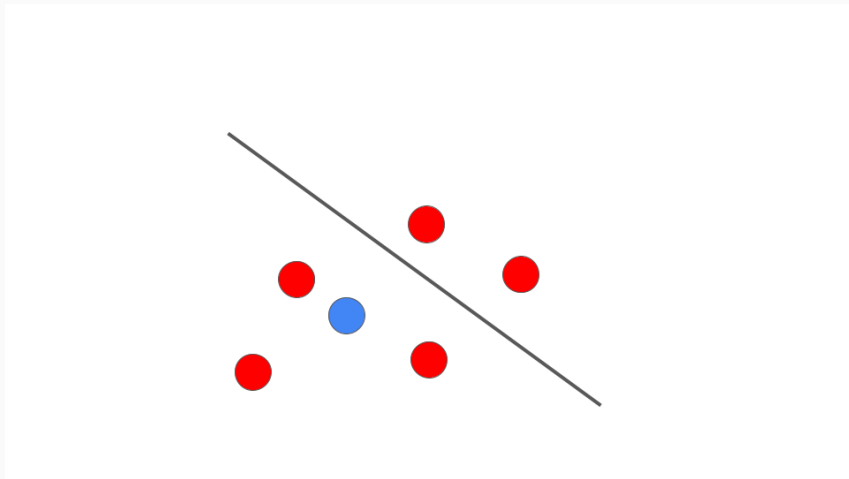


**A possible solution is to use Locality Sensitive Hashing (LSH)**

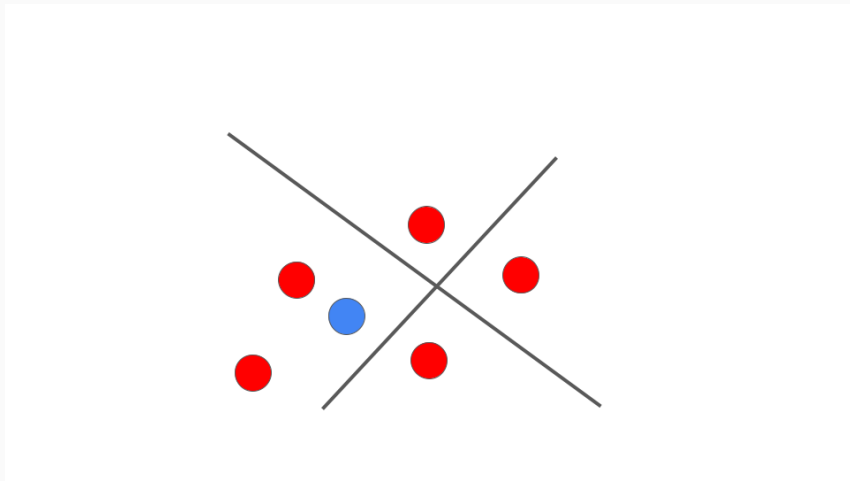
How to find out which red point is closest to the blue point?



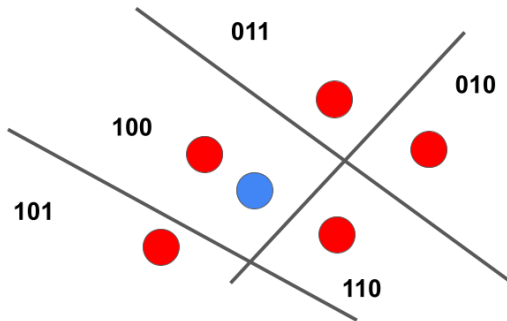
Select a random hyperplane...



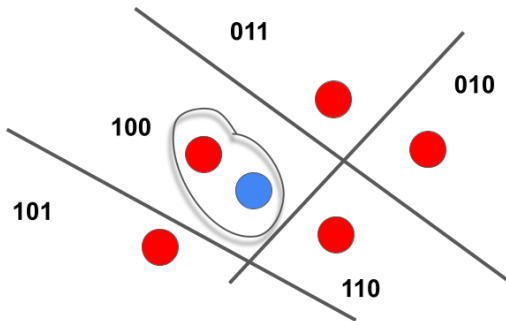
and another...



This creates a partition



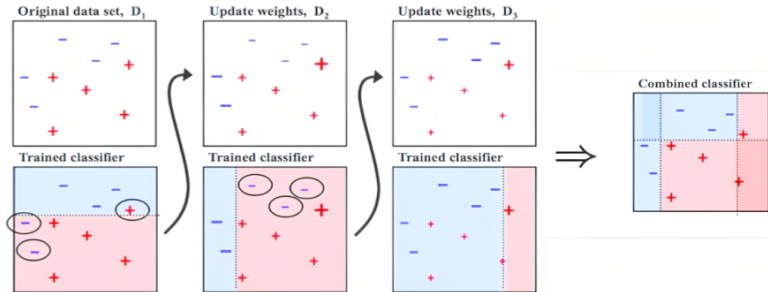
Compare only points in the same partition!





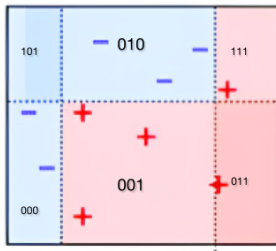
**Boost can be used to learn these hyperplanes effectively using a data-driven approach!**

# Similarity hashing: Boosting





Combined classifier



# Similarity hashing: Boosting

**Step 1:** The model learns simple similarity rules from the data

Movies			
name	year	genre	country
cidade de deus	2002	action	brazil
spider-man	2002	adventure	usa
robocop	2014	action	usa
cidade de deus	2002	drama	brazil

⋮

# Similarity hashing: Boosting

**Step 1:** The model learns simple similarity rules from the data

Movies			
name	year	genre	country
cidade de deus	2002	action	brazil
spider-man	2021	adventure	usa
robocop	2014	action	usa
cidade de deus	2002	drama	brazil

- Rule<sub>1</sub>: Is it from 2002?
- Rule<sub>2</sub>: Is it from Brazil?
- Rule<sub>3</sub>: Name starts with “c”?
- ⋮
- Rule<sub>T</sub>: The second letter in its name is an “i”?

# Similarity hashing: Boosting

**Step 1:** The model learns simple similarity rules from the data

Movies			
name	year	genre	country
cidade de deus	2002	action	brazil
spider-man	2021	adventure	usa
robocop	2014	action	usa
cidade de deus	2002	drama	brazil

- **Rule<sub>1</sub>:** Is it from 2002?
- **Rule<sub>2</sub>:** Is it from Brazil?
- **Rule<sub>3</sub>:** Name starts with “c”?
- **⋮**
- **Rule<sub>T</sub>:** The second letter in its name is an “i”?

# Similarity hashing: Boosting

**Step 1:** The model learns simple similarity rules from the data

Movies			
name	year	genre	country
cidade de deus	2002	action	brazil
spider-man	2021	adventure	usa
robocop	2014	action	usa
cidade de deus	2002	drama	brazil

- **Rule<sub>1</sub>:** Is it from 2002?
- **Rule<sub>2</sub>:** Is it from Brazil?
- **Rule<sub>3</sub>:** Name starts with “c”?
- **⋮**
- **Rule<sub>T</sub>:** The second letter in its name is an “i”?

# Similarity hashing: Boosting

**Step 1:** The model learns simple similarity rules from the data

Movies			
name	year	genre	country
cidade de deus	2002	action	brazil
spider-man	2021	adventure	usa
robocop	2014	action	usa
cidade de deus	2002	drama	brazil

- **Rule<sub>1</sub>:** Is it from 2002?
- **Rule<sub>2</sub>:** Is it from Brazil?
- **Rule<sub>3</sub>:** Name starts with “c”?
- $\vdots$
- **Rule<sub>T</sub>:** The second letter in its name is an “i”?

# Similarity hashing: Boosting

**Step 1:** The model learns simple similarity rules from the data

Movies			
name	year	genre	country
cidade de deus	2002	action	brazil
spider-man	2021	adventure	usa
robocop	2014	action	usa
cidade de deus	2002	drama	brazil

- **Rule<sub>1</sub>:** Is it from 2002?
- **Rule<sub>2</sub>:** Is it from Brazil?
- **Rule<sub>3</sub>:** Name starts with “c”?
- $\vdots$
- **Rule<sub>T</sub>:** The second letter in its name is an “i”?

# Similarity hashing: Boosting

**Step 1:** For each Rule, the model learns positive weights associated to its relevance and an error value

Movies			
name	year	genre	country
cidade de deus	2002	action	brazil
spider-man	2002	adventure	usa
robocop	2014	action	usa
cidade de deus	2002	drama	brazil

- Rule<sub>1</sub> has relevance  $\alpha_1 = 0.31$
- Rule<sub>2</sub> has relevance  $\alpha_2 = 0.29$
- Rule<sub>3</sub> has relevance  $\alpha_3 = 0.25$
- $\vdots$
- Rule<sub>T</sub> has relevance  $\alpha_T = 0.015$



# Similarity hashing: Boosting

**Step 1:** For each Rule, the model learns positive weights associated to its relevance and an error value

Movies			
name	year	genre	country
cidade de deus	2002	action	brazil
spider-man	2021	adventure	usa
robocop	2014	action	usa
cidade de deus	2002	drama	brazil

- **Rule<sub>1</sub>** has relevance  $\alpha_1 = 0.31$
- **Rule<sub>2</sub>** has relevance  $\alpha_2 = 0.29$
- **Rule<sub>3</sub>** has relevance  $\alpha_3 = 0.25$
- $\vdots$
- **Rule<sub>T</sub>** has relevance  $\alpha_T = 0.015$

# Similarity hashing: Boosting

- Using such rules and its weights, we can construct a similarity function between items  $A$  and  $B$  given by:

$$f^*(A, B) = \sum_{i=1}^T \alpha_i \text{Rule}_i(A) \text{Rule}_i(B)$$

where  $\text{Rule}_i(x) = 1$  if  $x$  satisfies  $\text{Rule}_i$  and  $-1$  otherwise

- Intuitively, for we want for some  $\theta \in (0, 1)$ :
  - if  $A$  and  $B$  are match then  $f^*(A, B) \geq \theta \approx 1$
  - if  $A$  and  $B$  are not match then  $f^*(A, B) \leq -\theta \approx -1$ ,

# Similarity hashing: Boosting

- Using such rules and its weights, we can construct a similarity function between items  $A$  and  $B$  given by:

$$f^*(A, B) = \sum_{i=1}^T \alpha_i \text{Rule}_i(A) \text{Rule}_i(B)$$

where  $\text{Rule}_i(x) = 1$  if  $x$  satisfies  $\text{Rule}_i$  and  $-1$  otherwise

- Intuitively, for we want for some  $\theta \in (0, 1)$ :
  - if  $A$  and  $B$  are match then  $f^*(A, B) \geq \theta \approx 1$
  - if  $A$  and  $B$  are not match then  $f^*(A, B) \leq -\theta \approx -1$ ,

## Theorem

If  $\theta > 0$ , then  $f^*$  satisfies the previous condition with probability at least  $1 - \varepsilon$ , where:

$$\varepsilon := 2^T \prod_{t=1}^T \text{error}_t^{1/2-\theta} (1 - \text{error}_t)^{\theta-1/2} \quad (1)$$

$$+ \frac{8}{\theta} (\mathcal{R}_{\mathcal{S}_{A,n}}(\mathcal{K}) + \mathcal{R}_{\mathcal{S}_{B,n}}(\mathcal{K})). \quad (2)$$

Furthermore, if there exists  $\gamma > 0$  such that for all  $t \in [T]$ ,  $\gamma \leq (1/2 - \text{error}_t)$  and  $\theta \leq 2\gamma$ , then the term in (1) decreases exponentially with  $T$ .

The proof relies on concentration of measure, margin theory and Rademacher complexity properties

# Similarity hashing: Boosting

**Step 2:** Given rules  $\text{Rule}_t$ , weights  $\alpha_t > 0$  and the similarity function  $f^*$  from the previous step

- To construct a single-bit hash we draw a random rule  $R$  following the distribution  $\mathbb{P}(R = \text{Rule}_i) = \alpha_i$
- We say that two items  $A$  and  $B$  have the same single-bit hash if  $R(A) = R(B)$
- It is easy to show that

$$\mathbb{P}[R(A) = R(B)] = \frac{1 + f^*(A, B)}{2}.$$

- We concatenate several single-bit hashes to construct a final hash function

# Similarity hashing: Boosting

**Step 2:** Given rules  $\text{Rule}_t$ , weights  $\alpha_t > 0$  and the similarity function  $f^*$  from the previous step

- To construct a single-bit hash we draw a random rule  $R$  following the distribution  $\mathbb{P}(R = \text{Rule}_i) = \alpha_i$
- We say that two items  $A$  and  $B$  have the same single-bit hash if  $R(A) = R(B)$
- It is easy to show that

$$\mathbb{P}[R(A) = R(B)] = \frac{1 + f^*(A, B)}{2}.$$

- We concatenate several single-bit hashes to construct a final hash function

# Similarity hashing: Boosting

**Step 2:** Given rules  $\text{Rule}_t$ , weights  $\alpha_t > 0$  and the similarity function  $f^*$  from the previous step

- To construct a single-bit hash we draw a random rule  $R$  following the distribution  $\mathbb{P}(R = \text{Rule}_i) = \alpha_i$
- We say that two items  $A$  and  $B$  have the same single-bit hash if  $R(A) = R(B)$
- It is easy to show that

$$\mathbb{P}[R(A) = R(B)] = \frac{1 + f^*(A, B)}{2}.$$

- We concatenate several single-bit hashes to construct a final hash function

# Similarity hashing: Boosting

**Step 2:** Given rules  $\text{Rule}_t$ , weights  $\alpha_t > 0$  and the similarity function  $f^*$  from the previous step

- To construct a single-bit hash we draw a random rule  $R$  following the distribution  $\mathbb{P}(R = \text{Rule}_i) = \alpha_i$
- We say that two items  $A$  and  $B$  have the same single-bit hash if  $R(A) = R(B)$
- It is easy to show that

$$\mathbb{P}[R(A) = R(B)] = \frac{1 + f^*(A, B)}{2}.$$

- We concatenate several single-bit hashes to construct a final hash function



# Similarity hashing: Boosting

**Step 2:** Given rules  $\text{Rule}_t$ , weights  $\alpha_t > 0$  and the similarity function  $f^*$  from the previous step

- To construct a single-bit hash we draw a random rule  $R$  following the distribution  $\mathbb{P}(R = \text{Rule}_i) = \alpha_i$
- We say that two items  $A$  and  $B$  have the same single-bit hash if  $R(A) = R(B)$
- It is easy to show that

$$\mathbb{P}[R(A) = R(B)] = \frac{1 + f^*(A, B)}{2}.$$

- We concatenate several single-bit hashes to construct a final hash function

# Similarity hashing: Boosting

We combine several single-bit hashes via the following algorithm:

---

**Algorithm** Algorithm to construct the hash codes

---

**Require:**  $k, L \in \mathbb{N}$ , convex weights  $(\alpha_t)_{t=1}^T$ , Rules  $(\text{Rule}_t)_{t=1}^T$

```

1: for  $i \leftarrow 1$  to  $L$  do
2:   for  $j \leftarrow 1$  to  $k$  do
3:      $g_{i,j} \leftarrow \text{Rule}_t$  with probability  $\alpha_t$ 
4:   end for
5:    $g_i \leftarrow (g_{i,1}, \dots, g_{i,k})$ 
6: end for
7:  $g \leftarrow (g_1, \dots, g_L)$ 
8: return  $g$ 

```

---

Two elements  $A$  and  $B$  are tested for similarity if  $g_i(A) = g_i(B)$  for some  $i = 1, \dots, L$

## Theorem

Consider datasets  $\mathcal{A}$  and  $\mathcal{B}$  such that  $|\mathcal{A}| = N_{\mathcal{A}}$  and  $|\mathcal{B}| = N_{\mathcal{B}}$ . Suppose our condition holds for  $\theta > 0$ . Then, given  $\gamma \in (0, 1)$ , if we set:

$$\rho := \frac{\log\left(\frac{2}{1+\theta}\right)}{\log\left(\frac{2}{1-\theta}\right)}, \quad k := \lceil \log_{\frac{2}{1+\theta}} N_{\mathcal{A}} \cdot N_{\mathcal{B}} \rceil \quad \text{and} \quad L := \left\lceil \frac{2(N_{\mathcal{A}} \cdot N_{\mathcal{B}})^{\rho} \log(1/\gamma)}{1 + \theta} \right\rceil,$$

then

$$\mathbb{E} [\text{Recall}] \geq (1 - \gamma)(1 - \varepsilon), \quad \mathbb{E} [\text{RR}] \geq \left(1 - \frac{|\mathcal{M}| + L}{N_{\mathcal{A}} \cdot N_{\mathcal{B}}}\right) (1 - \varepsilon).$$

Both expectations are with respect to the randomness in the hash code.

# Similarity hashing: Boosting

DATASET	BB	CANOPY	KLSH	TLSH	SPECT	AG	CTT	HYBRID
ABT_BUY	<b>0.911</b>	0.761	0.365	0.625	0.263	0.503	0.907	0.822
AMZ_GG	<b>0.877</b>	0.605	0.515	0.281	0.518	0.539	0.810	0.849
DBLP_ACM	0.993	0.850	0.895	0.861	0.662	0.696	0.993	<b>0.998</b>
DBLP_SCH	0.989	0.891	0.691	0.543	0.602	0.670	<b>0.991</b>	0.983
RESTAURANT	0.988	0.785	0.937	0.838	0.519	0.728	<b>0.997</b>	0.997
RLDATA500	<b>0.992</b>	0.829	0.969	0.982	0.691	0.717	0.966	0.966
RLDATA10K	<b>0.999</b>	0.929	0.926	0.987	0.755	0.800	0.957	0.926
MUSICBRAINZ	0.991	0.101	0.944	0.950	0.662	0.737	<b>0.994</b>	0.992
WM_AMZ	<b>0.943</b>	0.017	0.495	0.005	0.577	0.558	<b>0.943</b>	0.942
AVERAGE	<b>0.965</b>	0.641	0.749	0.675	0.583	0.660	0.951	0.942

Table: Harmonic mean for each model and dataset.

# Similarity hashing: Boosting

At each iteration  $t = 1, \dots, T$  of the boosting process, our model assigns a weight  $\alpha_t^*$  to a specific feature of the dataset to block. This weight serves as an indicator of the significance of this feature in matching entities.

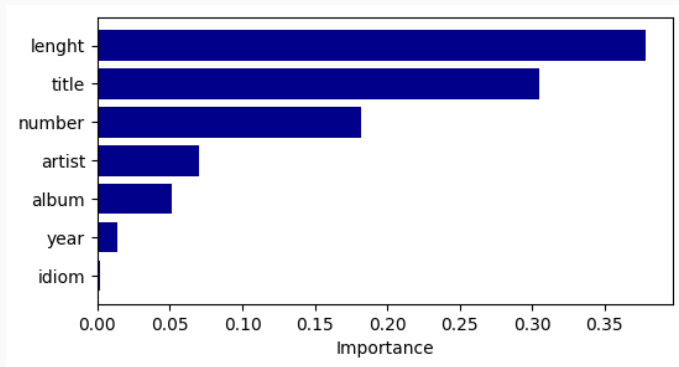


Figure: Feature relevance identified by the model during the boosting step for the `musicbrainz` dataset.



**Thank you!**