# LSH One

Yassine Abou Hadid, Clément Préaut, Benjamin Sykes
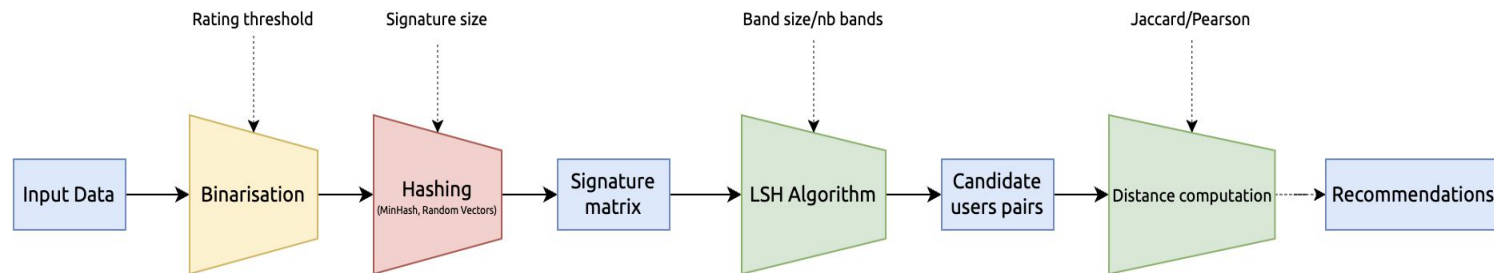
# Recap on LSH : LSH workflow for recommendation
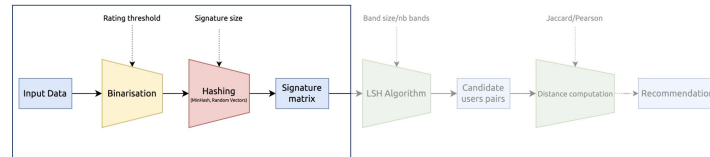
$$R_{user=i,movie=j} = \sum_{k=1}^{N} \text{UserSim}(i,k).R_{user=k,movie=j}$$

**Goal**: only relevant similarities should be taken into consideration

$$R_{user=i,movie=j} = \sum_{k/\exists h_\alpha, h_\alpha(i)=h_\alpha(k)} \text{UserSim}(i,k).R_{user=k,movie=j}$$
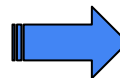
# Signature Matrix Generation



## Binarization

(threshold=3)

|  | movie 1 | movie 2 | movie 3 |
|--------|---------|---------|---------|
| user 1 | 3 | 2 | 4 |
| user 2 | 1 | 2 | 5 |
| user 3 | 4 | 3 | 2 |

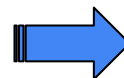|  | movie 1 | movie 2 | movie 3 |
|--------|---------|---------|---------|
| user 1 | 1 | 0 | 1 |
| user 2 | 0 | 0 | 1 |
| user 3 | 1 | 1 | 0 |

| V1 | -3 | 4 | 2 |
|----|----|---|---|

| V2 | 0 | -1 | 2 |
|----|---|----|---|

| V3 | -1 | 2 | -2 |
|----|----|---|----|

$$H_{V_i}(C_j) = \begin{cases} 1 & if \ V_i . C_j^T > 0 \\ 0 & otherwise \end{cases}$$

## Random

## Hyperplans

(3 vectors)

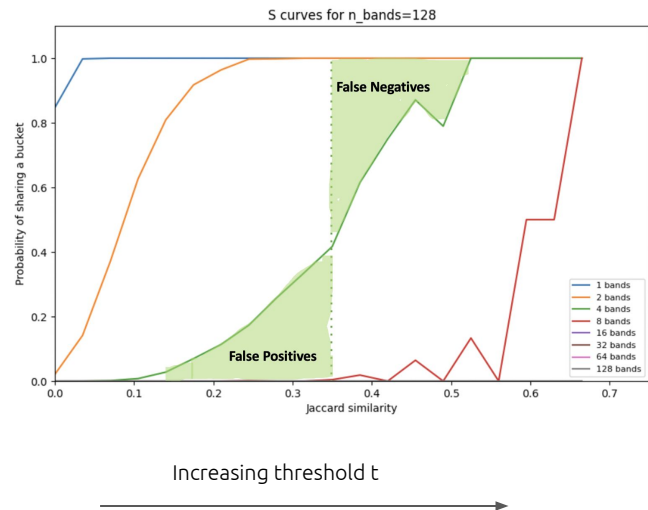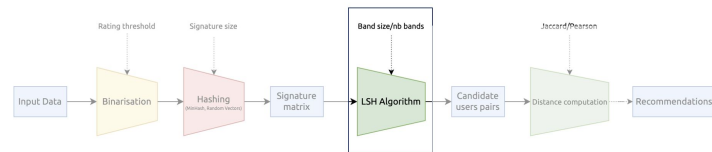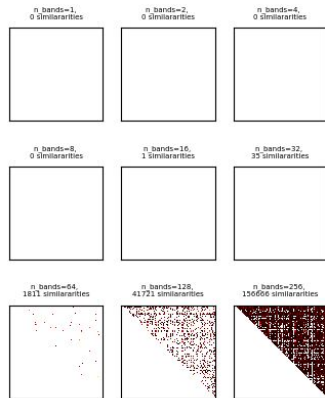| 1 | 0 | 1 |
|---|---|---|
| 0 | 0 | 1 |
| 1 | 1 | 0 |

| 0 | 1 | 0 |
|---|---|---|
| 1 | 1 | 0 |
| 1 | 0 | 1 |

**Signature Matrix**

3

# Approximating the s-curve



r=2

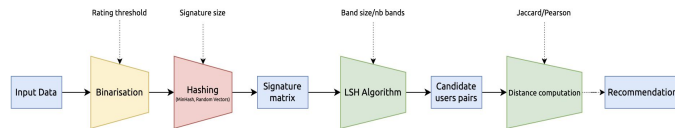| user 1 | 4 | 1 | 2 | 0 | 2 | 1 | 3 | 2 |
| user 2 | 2 | 4 | 3 | 1 | 2 | 1 | 3 | 2 |
| user 3 | 4 | 1 | 1 | 0 | 3 | 4 | 2 | 3 |

r=1     r=4

Evolution of computed similarities on full matrix

n_bands=1,
0 similararities

n_bands=2,
0 similararities

n_bands=4,
0 similararities

n_bands=8,
0 similararities

n_bands=16,
1 similararities

n_bands=32,
35 similararities

n_bands=64,
1811 similararities

n_bands=128,
41721 similararities

n_bands=256,
156666 similararities

Input Data → Binarisation → Hashing (Minhash, Random vectors) → Signature matrix → LSH Algorithm → Candidate users pairs → Distance computation → Recommendations

Rating threshold     Signature size     Band size/nb bands     Jaccard/Pearson

S curves for n_bands=128

Probability of sharing a bucket

False Negatives

False Positives

Jaccard similarity

1 bands
2 bands
4 bands
8 bands
16 bands
32 bands
64 bands
128 bands

Increasing threshold t

$$threshold = (1/b)^{1/r}$$

Trade-off **speed/false negatives** vs **false positives**

4

# Optimisation → scaling



**LSH is made for large data**

$$\mathcal{O}(n_{\mathrm{users}}.n_{\mathrm{movies}}.8bits) \xrightarrow{\text{hash}} \mathcal{O}(n_{\mathrm{users}}.s_{\mathrm{signature}}.8bits)$$

**Precomputation of the signature matrix**

*All the computations for **similarity evaluation** can be then made **offline***

**Computation of the sparse similarity matrix (numpy)**

$$R_{i,m} = \boldsymbol{S}_{i,.}\mathrm{R}_{,m}$$

| 4 |
|---|
| 2 |
| 4 |
| 1 |
| 0 |

Ratings for movie **m**

| 0.2 | 0.7 | 0.4 | 0.3 | 0.5 |
|-----|-----|-----|-----|-----|

Similarities for user **i**
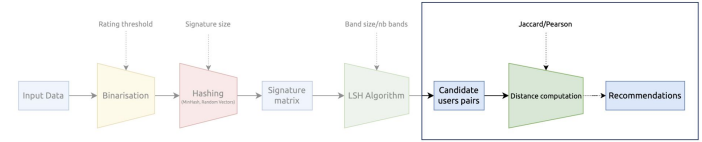
**R** = 2.05

# Curve interpretation

# Thank you for your attention!