



# Recherche d'information par le contenu II

Master2 - Partage de données à grande échelle - GMIN307

Alexis JOLY - [alexis.joly@inria.fr](mailto:alexis.joly@inria.fr)

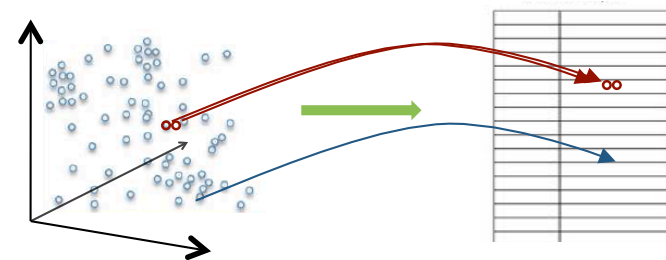
# Etude détaillée de LSH

# LSH: D

LSH (Localité Sensitive Hashing), méthode de hachage introduite par Gionis, Indyk et Motwani en 1998 pour la recherche approximative dans les espaces de grande dimensions

Sensible à la localité = la probabilité de collision des hash code est forte si les vecteurs sont proches et faible si les vecteurs sont éloignés

Formellement, selon Gionis et al. en 1998



$\mathcal{F}$  est une famille de fonctions  $h : \mathbb{R}^d \rightarrow S$  satisfaisant les conditions suivantes pour deux points quelconques  $p, q \in \mathbb{R}^d$  et une fonction  $h$  choisie aléatoirement parmi la famille  $\mathcal{F}$ :

- si  $d(p, q) \leq R$ , alors  $Pr_{h \in H}[h(p) = h(q)] \geq P_1$
- si  $d(p, q) \geq cR$ , alors  $Pr_{h \in H}[h(p) = h(q)] \leq P_2$

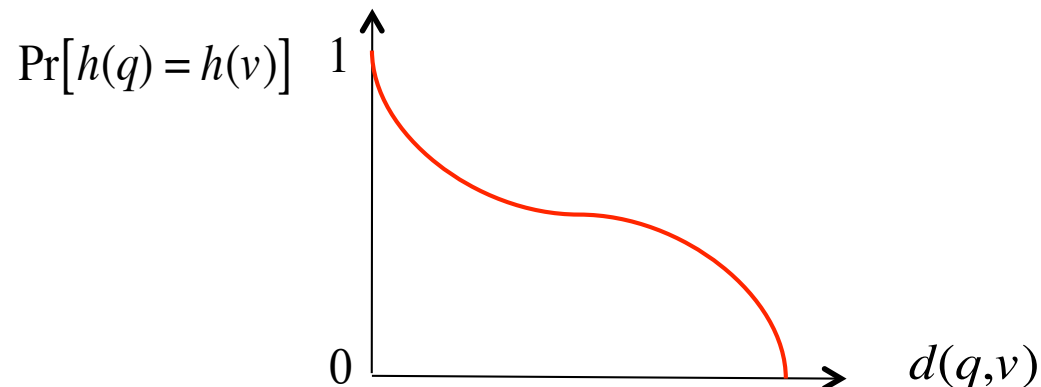
# LSH: Intro

LSH (Localité Sensitive Hashing), méthode de hachage introduite par Gionis, Indyk et Motwani en 1998 pour la recherche approximative dans les espaces de grande dimensions

Pour certaines familles LSH, la probabilité de collision peut même s'exprimer comme une fonction décroissant avec la distance

$$\Pr[h_{\theta}(q) = h_{\theta}(v)]_{p_{\theta}} = f(d(q,v))$$

- $\theta$  is i.i.d drawn from a known **data independent distribution**
- $f(d)$  is the sensitivity function (monotonically increasing from 0 to 1)



# LSH sensible au produit scalaire

Exemple: famille LSH sensible au produit scalaire

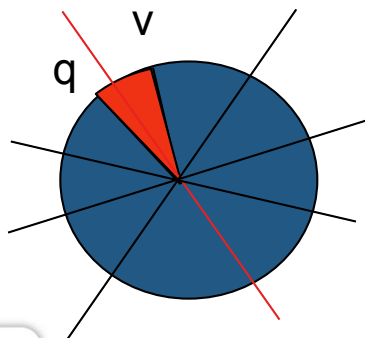
$$h(x) = \text{sgn}(w^T x) \quad p_w = N(0, I)$$

C'est le cas que l'on a déjà rencontré = projections aléatoires + quantification scalaire binaire

On peut montrer que la fonction de sensibilité vaut

$$\text{Pour tous } q, v \in \mathbb{R}^d \quad \Pr[h_w(q) = h_w(v)]_w = 1 - \frac{1}{\pi} \cos^{-1} \left( \frac{q^T v}{\|q\| \|v\|} \right)$$

Interprétation dans le cas normé:



$$\begin{aligned} \Pr[h_w(q) \neq h_w(v)]_w &= \text{angle}(q, v) / \pi = \frac{1}{\pi} \cos^{-1}(q^T v) \\ &= \frac{1}{\pi} \cos^{-1} \left( 1 - \frac{(d(q, v))^2}{2} \right) \end{aligned}$$

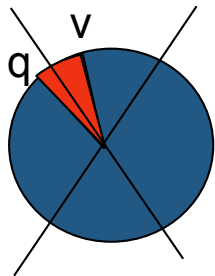
# LSH sensible au produit scalaire

Exemple: famille LSH sensible au produit scalaire

$$h(x) = \text{sgn}(w^T x) \quad p_w = N(0, I) \quad \Pr[h_w(q) = h_w(v)]_w = 1 - \frac{1}{\pi} \cos^{-1} \left( \frac{q^T v}{\|q\| \|v\|} \right)$$

On forme un hash code de D bits en utilisant D fonctions de hachage de la famille  $\mathcal{F}$

$$z(x) = [h_1(x), \dots, h_D(x)] = \text{sign}(Wx)$$



$$z(q) = 1 \text{ hash code} = [0 \ 1 \ 0 \ 1 \ 0 \ 1 \ 1 \ 0 \ \dots \ 0 \ 0 \ 1 \ 0]$$



$$z(v) = 1 \text{ hash code} = [1 \ 0 \ 1 \ 1 \ 1 \ 1 \ 0 \ 1 \ \dots \ 0 \ 1 \ 0 \ 0]$$



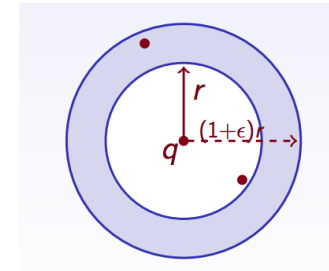
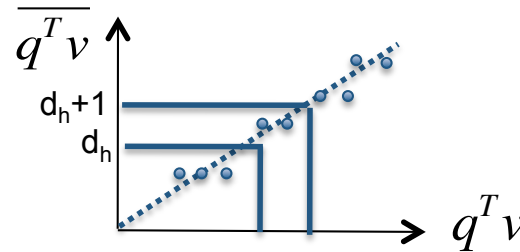
$$\approx \frac{D}{\pi} \cos^{-1} \left( \frac{q^T v}{\|q\| \|v\|} \right)$$

# LSH sensible au produit scalaire

La distance de Hamming entre deux hash code est un estimateur de la probabilité de collision et donc du produit scalaire entre  $q$  et  $v$

$$\overline{q^T v} = \|q\| \|v\| \cos\left(\frac{\pi}{D} d_h(z(q), z(v))\right)$$

Variance



La variance  $\sigma(k, q^T v)$  de l'estimateur décroît avec le nombre de bits et converge vers le produit scalaire exact

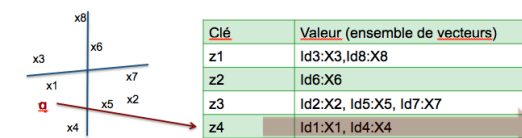
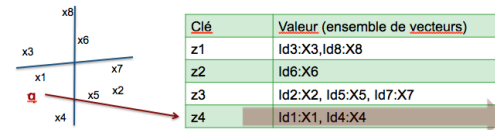
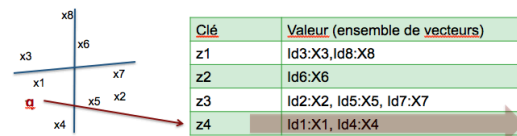
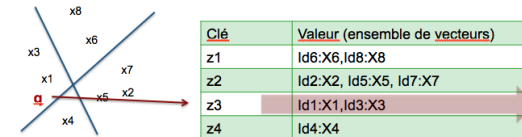
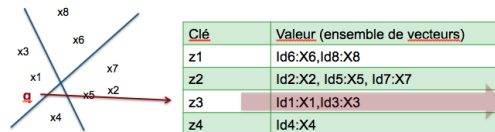
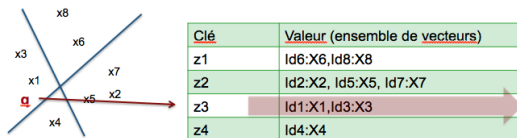
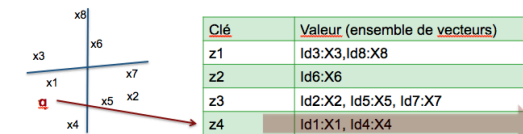
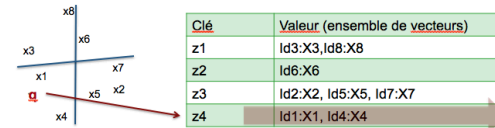
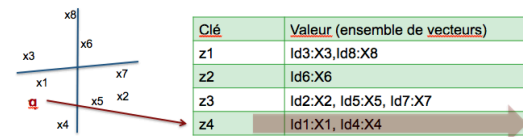
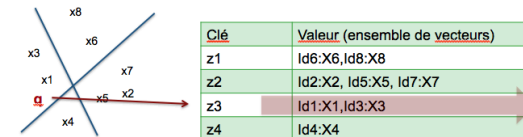
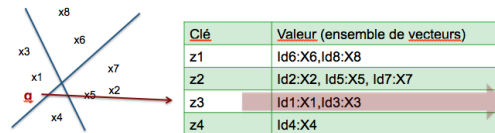
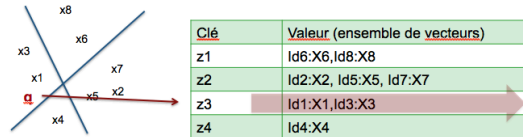
$$\lim_{k \rightarrow \infty} \overline{q^T v} = q^T v$$

Les  $k$  plus proches voisins dans l'espace de Hamming convergent vers les  $k$  plus proches voisins exactes

# LSH: algorithme de base pour l'indexation et la recherche

Tables multiples: création de L tables dont les clé sont générées par k fonctions de hachage (hash code de taille k pour chaque table). Taille index  $O(Ln)$

Recherche = accès simple dans chaque table  $O(L)$





# LSH: algorithme de base

Probabilité de collision dans une table de k bits:

$$\Pr[h_w(q) = h_w(v)]_w = f(d(q,v)) \longrightarrow \Pr[z(q) = z(v)]_w = (f(d(q,v)))^k$$

Fonction de sensibilité

Probabilité de non collision dans une table:

$$\Pr[z(q) \neq z(v)]_w = 1 - (f(d(q,v)))^k$$

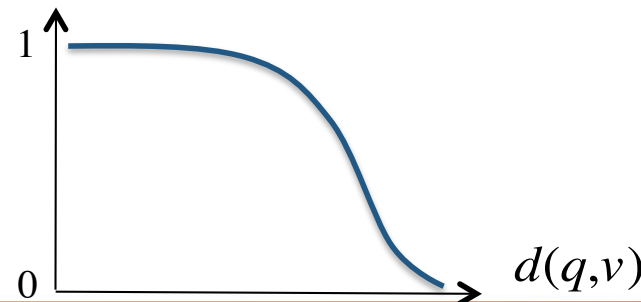
↓  
Très faible

Probabilité de non collision dans L tables:

$$(1 - (f(d(q,v)))^k)^L$$

Probabilité de collision dans au moins une table:

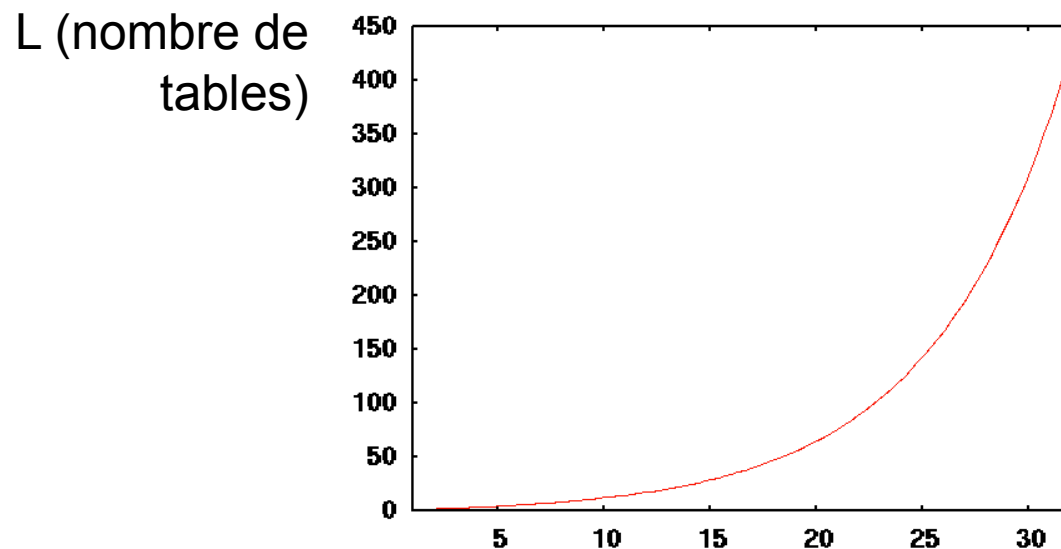
$$1 - (1 - (f(d(q,v)))^k)^L$$



# LSH: algorithme de base

Si on veut retrouver  $\alpha$  (e.g 95%) des  $v$  tels que  $d(q,v) < r$ , il faut que la probabilité de collision dans au moins une table soit supérieure à  $\alpha$

$$1 - \left(1 - (f(r))^k\right)^L > \alpha \quad \longrightarrow \quad L > \frac{\ln(1 - \alpha)}{\ln(1 - f(r)^k)}$$



A qualité  $\alpha$  constante et rayon  $r$  constant, le nombre de tables doit croître pour compenser l'augmentation de  $k$

k (nombre de bits par table)

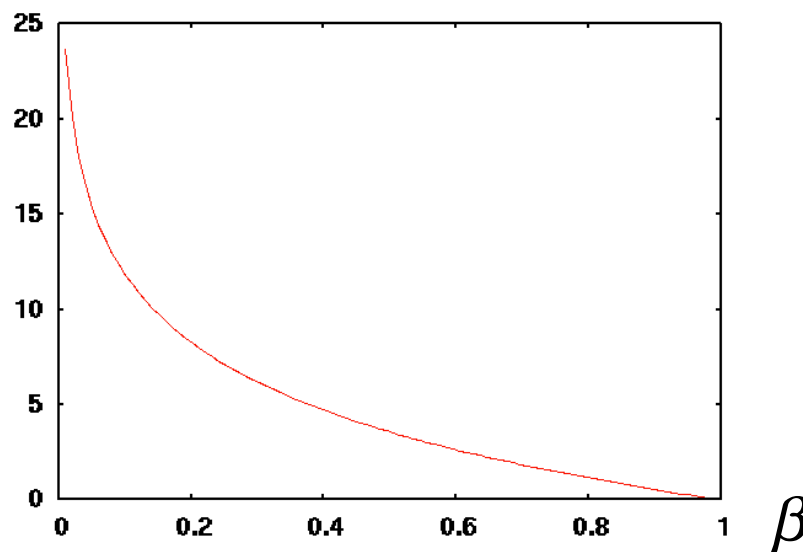
# LSH: algorithme de base

D'un autre côté, l'augmentation de  $k$  permet de réduire la probabilité de fausses collisions dans chaque table. Si on veut que la probabilité de collision soit faible lorsque  $d(q,v) > r + \epsilon$

Proba de collision dans une table

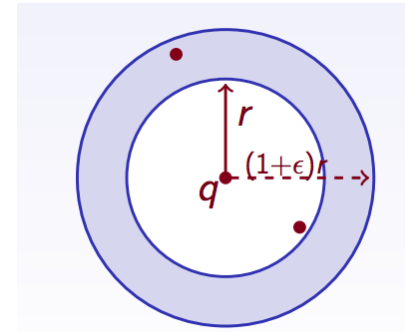
$$n(f((1 + \epsilon)r))^k < \beta \quad \longrightarrow \quad k > \frac{\ln(\beta)}{\ln(f((1 + \epsilon)r))}$$

$k$  (nombre de bits par table)



# Propriété fondamentale de LSH dans le cas général

Requête à un rayon près  $(1+\epsilon)$ -approximatives



- Given: radius  $R$ , error parameter  $\epsilon$  and query point  $q$ :
  - if there exists data point  $p$  s.t.  $q \in B(p, R)$ , return YES and a point (or all points)  $p'$  s.t.  $q \in B(p', (1 + \epsilon)R)$ ,
  - if  $q \notin B(p, R)$  for all data points  $p$ , return NO,
  - if closest data point to  $q$  is at distance between  $R$  and  $R(1 + \epsilon)$  then return YES or NO

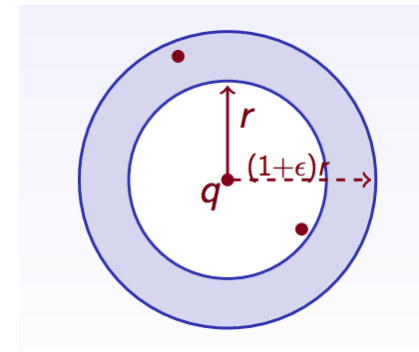
# Propriété fondamentale de LSH dans le cas général

Famille de fonctions de hachage locality-sensitive  $\mathcal{H}$  avec paramètres  $(\epsilon, r, P_1, P_2)$ :

- Si  $\|p - q\| \leq r$  alors  $\Pr_{\mathcal{H}}[h(p) = h(q)] \geq P_1$
- Si  $\|p - q\| \geq (1 + \epsilon)r$  alors  $\Pr_{\mathcal{H}}[h(p) = h(q)] \leq P_2$

## Theorem

Toute famille de fonctions de hachage  $(\epsilon, r, P_1, P_2)$ -locality-sensitive conduit à un algorithme pour la recherche à un rayon  $r$  près  $(1 + \epsilon)$ -approximative avec  $\sim n^\rho$  de temps de recherche et  $n^{1+\rho}$  de taille d'index



$$\rho = \frac{\log(1/P_1)}{\log(1/P_2)} < 1$$

# Propriété fondamentale de LSH dans le cas général

Pour avoir une probabilité d'erreur au plus égale à  $\delta$  il faut prendre  $k, L$  tels que

$$P_2^k n \approx 1 \quad L \approx (1/P_1)^k \log(1/\delta)$$

Avec ces contraintes on obtient en effet:

$$k = \frac{\log n}{\log(1/P_2)}$$

$$L = (1/P_1)^{\frac{\log n}{\log(1/P_2)}} \log(1/\delta) = n^{\frac{\log(1/P_1)}{\log(1/P_2)}} \log(1/\delta) = n^\rho \log(1/\delta)$$

# Propriété fondamentale de LSH dans le cas général

L'espérance du nombre d'objets rencontrés à moins de  $(1 + \epsilon)r$  de la requête est  $P_2^k L n \approx L$

Pour les vrais voisins à  $r$ -près la probabilité d'être hachés dans une même case que  $q$  est au moins

$$1 - (1 - (1/P_1)^k)^L \geq 1 - (1/e)^{\frac{L}{(1/P_1)^k}} \geq 1 - \delta$$

Taille de l'index  $\mathcal{O}(Ln) \approx n^{1+\rho+o(1)}$

Recherche  $\mathcal{O}(L) \approx n^{\rho+o(1)}$

# LSH sensible aux distances $L_p$

## $p$ -Stable distributions

- **$p$ -stable distribution** ( $p \geq 0$ ): A distribution  $\mathcal{D}$  over  $\mathfrak{R}$  s.t.
  - $n$  real numbers  $v_1 \dots v_n$ ,
  - i.i.d. variables  $X_1 \dots X_n$  with distribution  $\mathcal{D}$ ,
  - r.v.  $\sum_i v_i X_i$  has the same distribution as the variable  $(\sum_i |v_i|^p)^{1/p} X = l_p(\mathbf{v})X$ , where  $X$  is a r.v. with distribution  $\mathcal{D}$
- E.g.  $p$ -Stable distr for  $p = 1$  is Cauchy distr, for  $p = 2$  is Gaussian distr

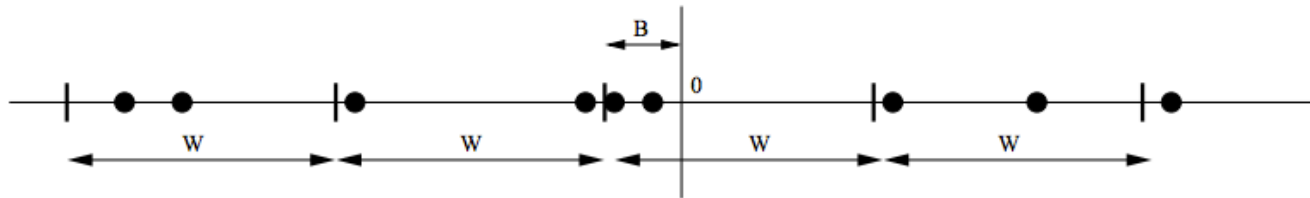


# LSH sensible aux distances $L_p$

## How are $p$ -Stable distributions useful?

- Consider a vector  $\mathbf{X} = \{X_1, X_2, \dots, X_N\}$ , where each  $X_i$  is drawn from a  $p$ -Stable distr
- For any pair of vectors  $\mathbf{a}, \mathbf{b}$   $\mathbf{a} \cdot \mathbf{X} - \mathbf{b} \cdot \mathbf{X} = (\mathbf{a} - \mathbf{b}) \cdot \mathbf{X}$  (by linearity)
- Thus  $\mathbf{a} \cdot \mathbf{X} - \mathbf{b} \cdot \mathbf{X}$  is distributed as  $(l_p(\mathbf{a} - \mathbf{b}))X'$  where  $X'$  is a  $p$ -Stable distr r.v.
- Using multiple independent  $\mathbf{X}$ 's we can use  $\mathbf{a} \cdot \mathbf{X} - \mathbf{b} \cdot \mathbf{X}$  to estimate  $l_p(\mathbf{a} - \mathbf{b})$  [Ind'01]

# LSH sensible aux distances $L_p$



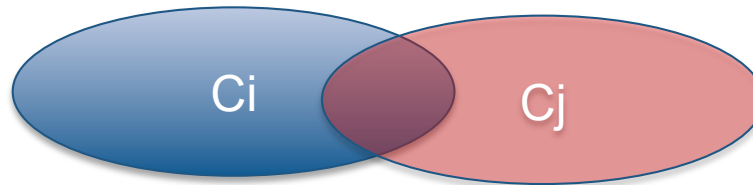
- Consider  $h_{\mathbf{a},b} \in \mathcal{H}^w$ ,  $h_{\mathbf{a},b}(\mathbf{v}) : \mathcal{R}^d \rightarrow \mathcal{N}$
- $\mathbf{a}$  is a  $d$  dimensional random vector whose each entry is drawn from a  $p$ -stable distr
- $b$  is a random real number chosen uniformly from  $[0, w]$  (random shift)
- $h_{\mathbf{a},b}(\mathbf{v}) = \lfloor \frac{\mathbf{a} \cdot \mathbf{v} + b}{w} \rfloor$

# LSH sensible à la distance de Jaccard

MinHash: the min-wise independent permutations locality sensitive hashing scheme

Jaccard distance = mesure de similarité entre des ensembles d'objets

$$\text{sim}_J(C_i, C_j) = \frac{|C_i \cap C_j|}{|C_i \cup C_j|}$$



- View sets as columns of a matrix; one row for each element in the universe.  $a_{ij} = 1$  indicates presence of item  $i$  in set  $j$
- Example

	$C_1$	$C_2$	
	0	1	$\text{sim}_J(C_1, C_2) = 2/5 = 0.4$
	1	0	
	1	1	
	0	0	
	1	1	
	0	1	
	0	1	

Hash function=

- Randomly **permute** rows
- Hash  $h(C_i) =$  index of first row with 1 in column  $C_i$
- **Suprising Property**

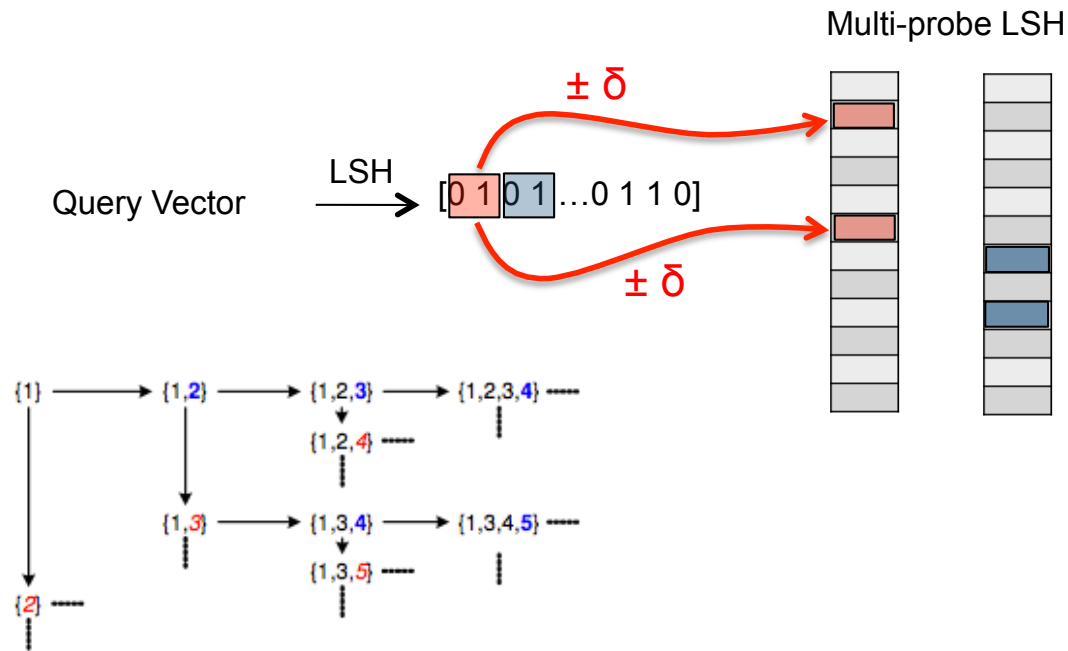
$$P[h(C_i) = h(C_j)] = \text{sim}_J(C_i, C_j)$$

# Limitations de LSH et variantes

# Problème de LSH: espace mémoire

LSH peut nécessiter un grand nombre de tables pour atteindre des qualités suffisantes ce qui rend la méthode problématique lorsque les données sont grandes

→ Utilisation d'accès multiples pour réduire le nombre de tables



# Exemple: LSH pour vecteurs binaires de grande dimension

Vecteurs d'entrée déjà binarisés (souvent le cas en audio fingerprinting)

- Cube de Hamming:  $\{0, 1\}^d$  avec distance de Hamming

Fonction de hachage = une réduction de  $\{0, 1\}^d$  dans  $\{0, 1\}^k$ . Le choix des composantes se fait selon une loi binomiale:

- Choisir aléatoirement un ensemble de composantes selon une loi binomiale de paramètre  $\frac{1}{2l}$ , (espérance du nombre de composantes  $\frac{d}{2l}$ )
- Affecter aléatoirement 0 ou 1 à toutes ces positions, affecter les autres à 0. Soit  $r$  le vecteur résultat.
- $h_r(p) = r \cdot p$

= Projection dans des sous espaces aléatoires de dimension avec  $k \ll d$  typiquement:  
 **$k = \log(d)$**

# Exemple: LSH pour vecteurs binaires de grande dimension

**Recherche à  $(\delta_1, \delta_2)$ -près** = conserver les vecteurs à moins de  $d \cdot \delta_1$  bits de la requête, rejeter ceux à plus de  $d \cdot \delta_2$  bits de la requête

**Indexation/pré-traitement:**

- Projeter tous les  $p$  dans  $A(p) = h_{r_1}(p) \dots h_{r_k}(p)$

- Pour tout  $v \in \{0, 1\}^k$  précalculer tous les  $(\frac{\delta_1 + \delta_2}{2})k$ -voisins

=

Buckets voisines à un rayon près

**Recherche:**

- Calculer  $A(q) = h_{r_1}(q) \dots h_{r_k}(q)$

- Retrouver et vérifier explicitement les  $(\frac{\delta_1 + \delta_2}{2})k$ -voisins de  $A(q)$

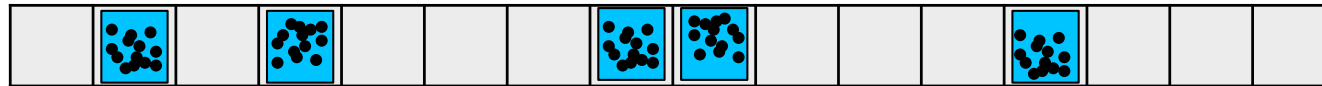
=

Visite des buckets voisines de la requête dans chaque table.

- Temps de recherche logarithmique, taille de l'index polynomiale

## Autre problème de LSH: mauvais balancement

Sur certaines données réelles, LSH peut conduire à des partitions loin de l'uniformité. Certaines bucket peuvent concentrer presque tous les points:



- Augmentation du nombre de fausses collisions
- Augmentation du temps de raffinement
- Problématique dans les contextes distribués (load balancing)

Solution: utiliser des fonctions de hachage dépendant des données. De nombreuses fonctions « data dependent » ont été proposées récemment dans la littérature (spectral-hashing, KLSH, RMMH, etc.)

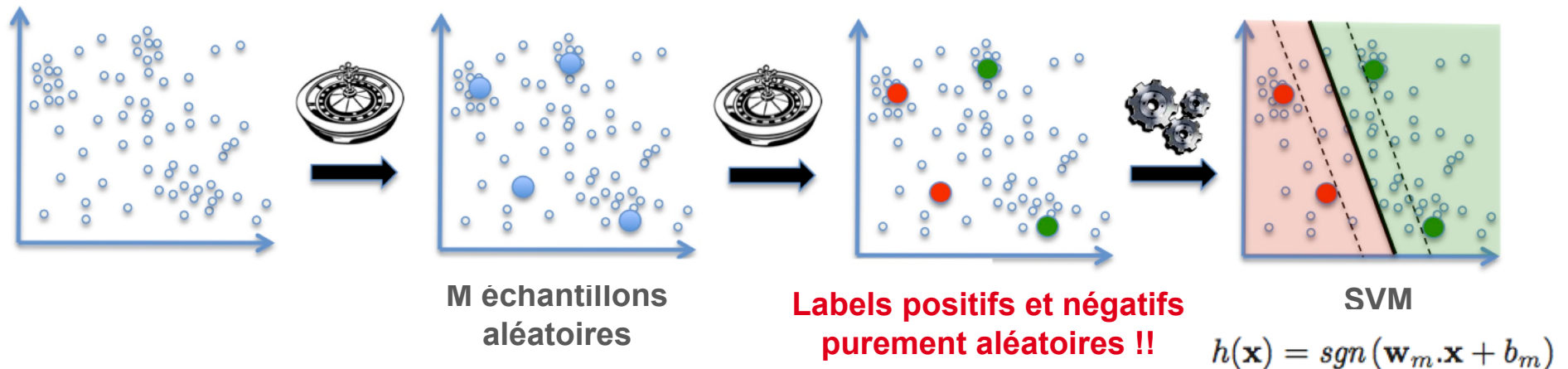


# Random Maximum Margin Hashing

RMMH = une famille de fonctions de hachage basée sur le **partitionnement aléatoire des données**

Avantage = fonctions indépendantes + adaptabilité aux données

Basé sur **un concept théorique nouveau** = l'apprentissage aléatoire de **classifieurs** (1 par fonction de hachage = 1 par bit)



# Random Maximum Margin Hashing

## Pourquoi ça marche ?

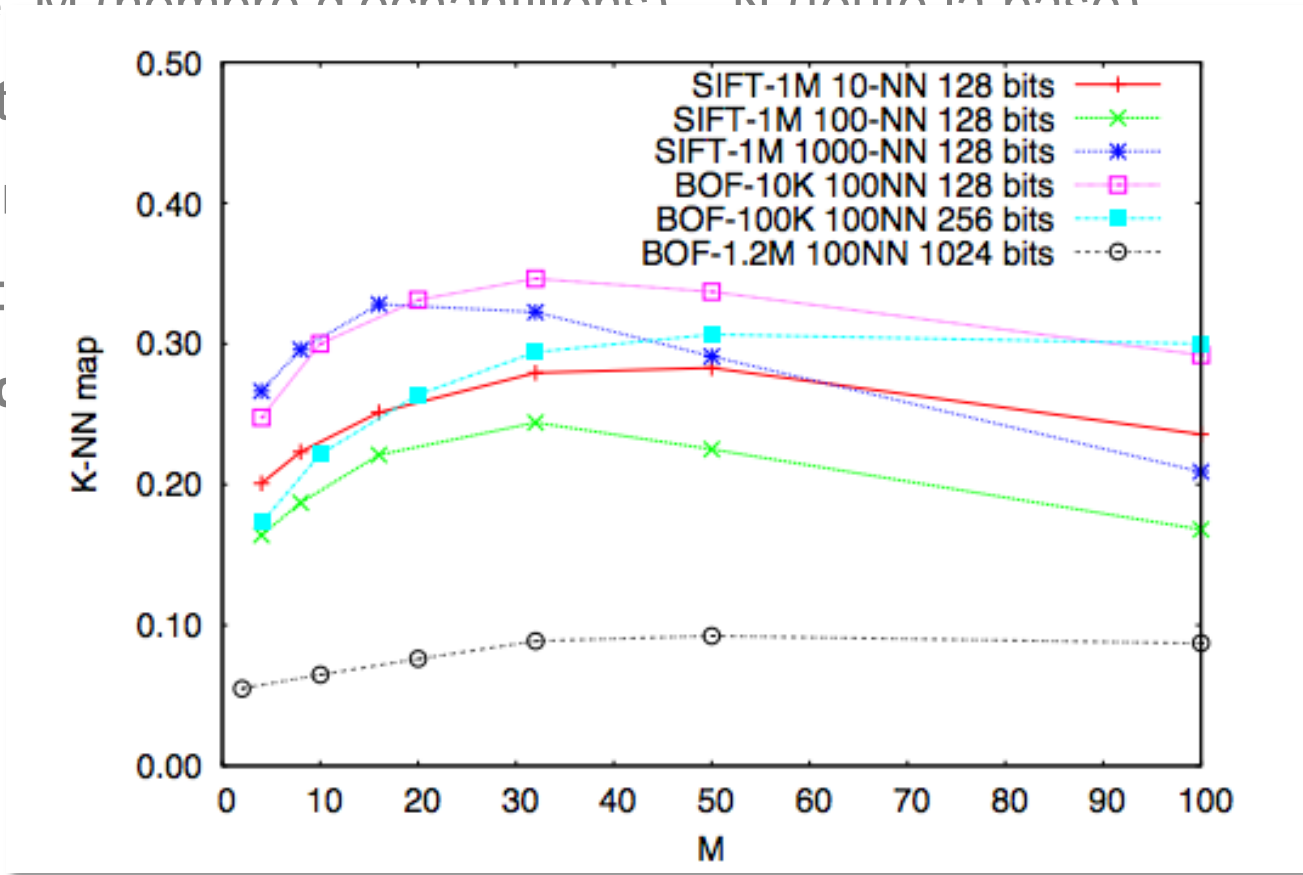
Cas extrême:  $M$  (nombre d'échantillons) =  $N$  (toute la base)

Uniformité

Overfitting

$M$  trop faible:

$M$  optimal =  $\epsilon$



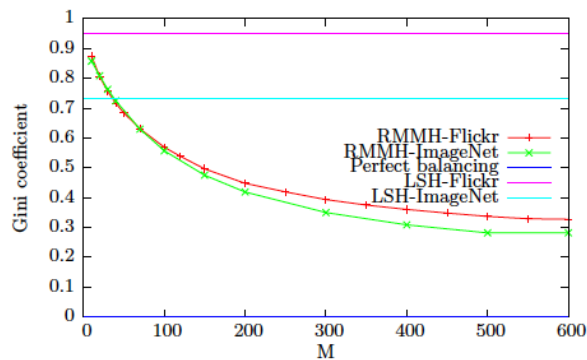
age

# Random Maximum Margin Hashing

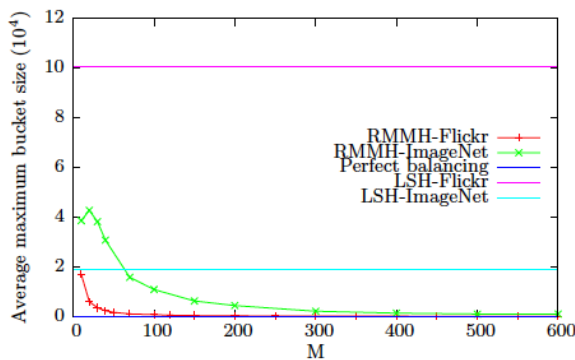
## Pourquoi ça marche ?

Cas extrême:  $M$  (nombre d'échantillons) =  $N$  (toute la base)

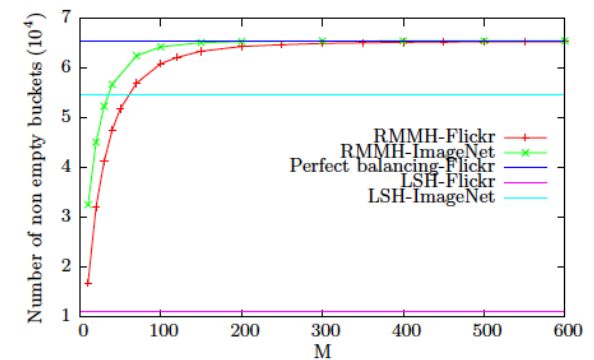
Uniformité parfaite de la partition apprise ( $\approx$ Voronoi)



(a) Gini coefficient



(b) Avg. Max. bucket size



(c) Number of non empty buckets

# Random Maximum Margin Hashing

## Evaluation on ImageNet

1.2M images described by BoVW features (dim 1000)

Supervised classification across 1000 semantic categories, 1000-NN

