

Randomized Algorithms

William Cohen

Outline

- Randomized methods: today
 - SGD with the hash trick (recap)
 - Bloom filters
- Later:
 - count-min sketches
 - locality sensitive hashing

THE HASH TRICK: A REVIEW

Hash Trick - Insights

- Save memory: don't store hash keys
- Allow collisions
 - even though it distorts your data some
- Let the learner (downstream) take up the slack

Learning as optimization for regularized logistic regression

- Algorithm: $w^j = w^j + \lambda(y - p)x^j - \lambda 2\mu w^j$
- Initialize arrays W, A of size R and set $k=0$
- For each iteration $t=1, \dots, T$
 - For each example (\mathbf{x}_i, y_i)
 - Let V be hash table so that $V[h] = \sum_{j: \text{hash}(j) \% R == h} x_i^j$
 - $p_i = \dots$; $k++$
 - For each hash value $h: V[h] > 0$:
 - » $W[h] *= (1 - \lambda 2\mu)^{k-A[j]}$
 - » $W[h] = W[h] + \lambda(y_i - p^i)V[h]$
 - » $A[h] = k$

Learning as optimization for regularized logistic regression

- Initialize arrays W, A of size R and set $k=0$
- For each iteration $t=1, \dots, T$
 - For each example (\mathbf{x}_i, y_i)
 - $k++$; let V be a new array of size R ; let $tmp=0$
 - For each $j: \mathbf{x}_i^j > 0: V[\text{hash}(j) \% R] += \mathbf{x}_i^j$
 - Let $ip=0$
 - For each $h: V[h] > 0$:
 - $W[h] *= (1 - \lambda 2\mu)^{k-A[h]}$
 - $ip += V[h] * W[h]$
 - $A[h] = k$
 - $p = 1/(1 + \exp(-ip))$
 - For each $h: V[h] > 0$:
 - $W[h] = W[h] + \lambda(y_i - p^i)V[h]$

$$V[h] = \sum_{j: \text{hash}(j) \% R == h} x_i^j$$

regularize $W[h]$'s

$$p \equiv \frac{1}{1 + e^{-\mathbf{v} \cdot \mathbf{w}}}$$

$$w^j = w^j + \lambda(y - p)x^j - \lambda 2\mu w^j$$

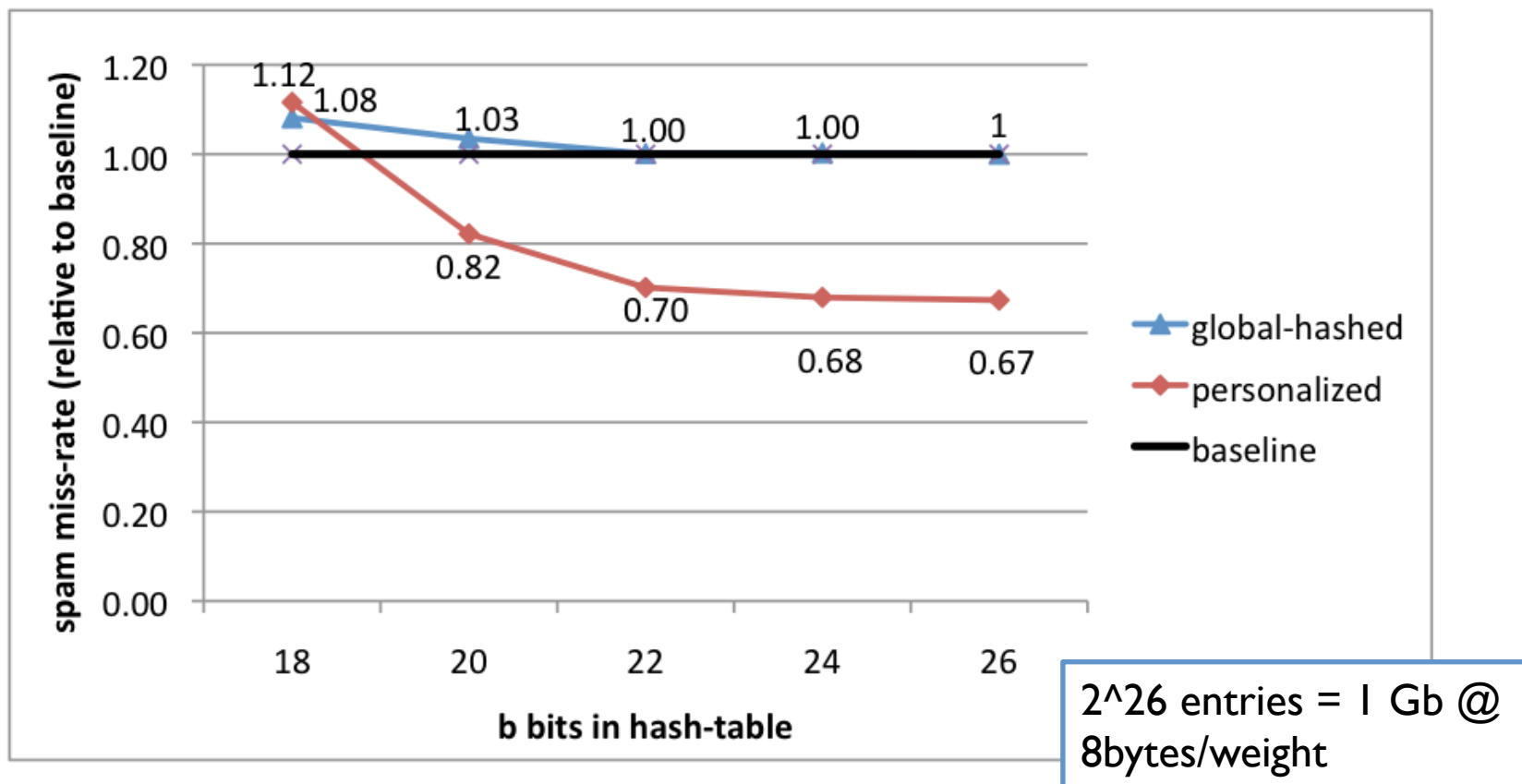


Figure 2. The decrease of uncaught spam over the baseline classifier averaged over all users. The classification threshold was chosen to keep the not-spam misclassification fixed at 1%. The hashed global classifier (*global-hashed*) converges relatively soon, showing that the distortion error ϵ_d vanishes. The personalized classifier results in an average improvement of up to 30%.

Data Sets	#Train	#Test	#Labels
RCV1	781,265	23,149	2
Dmoz L2	4,466,703	138,146	575
Dmoz L3	4,460,273	137,924	7,100

Table 1: Text data sets. #X denotes the number of observations in X.

	HLF (2^{28})		HLF (2^{24})		HF		no hash	U base	P base
	error	mem	error	mem	error	mem	mem	error	error
L2	30.12	2G	30.71	0.125G	31.28	2.25G (2^{19})	7.85G	99.83	85.05
L3	52.10	2G	53.36	0.125G	51.47	1.73G (2^{15})	96.95G	99.99	86.83

Table 5: Misclassification and memory footprint of hashing and baseline methods on DMOZ. HLF: joint hashing of labels and features. HF: hash features only. no hash: direct model (not implemented as too large, hence only memory estimates—we have 1,832,704 unique words). U base: baseline of uniform classifier. P base: baseline of majority vote. mem: memory used for the model. Note: the memory footprint in HLF is essentially independent of the number of classes used.

MOTIVATING BLOOM FILTERS

A variant of feature hashing

- Hash each feature *multiple times* with different hash functions
- Now, each w has k chances to *not* collide with another useful w'
- An easy way to get multiple hash functions
 - Generate some random strings s_1, \dots, s_L
 - Let the k -th hash function for w be the ordinary hash of concatenation $w \bullet s_k$

$$V[h] = \sum_k \sum_{j: \text{hash}(j \cdot s_k) \% R = h} x_i^j$$

A variant of feature hashing

$\mathbf{a} \neq \mathbf{b}$ are binary vectors

$V(\mathbf{a})$	1	0	1	0
$V(\mathbf{b})$	1	0	1	0

$V(\mathbf{a})$	1	0	1	0
	0	1	1	0
$V(\mathbf{b})$	1	0	1	0
	1	1	0	0

1	1	2	0
---	---	---	---

1	1	1	0
---	---	---	---

- An easy way to get multiple hash functions
 - Generate some random strings s_1, \dots, s_L
 - Let the k -th hash function for w be the ordinary hash of concatenation $w \bullet s_k$

$$V[h] = \sum_k \sum_{j: \text{hash}(j \cdot s_k) \% R = h} x_i^j$$

A variant of feature hashing

- Why would this work?

$$V[h] = \sum_k \sum_{j: \text{hash}(j \cdot s_k) \% R = h} x_i^j$$

- Claim: with 100,000 features and 100,000,000 buckets:
 - $k=1 \rightarrow \Pr(\text{any feature duplication}) \approx 1$
 - $k=2 \rightarrow \Pr(\text{any feature duplication}) \approx 0.4$
 - $k=3 \rightarrow \Pr(\text{any feature duplication}) \approx 0.01$

Hash Trick - Insights

- Save memory: don't store hash keys
- Allow collisions
 - even though it distorts your data some
- Let the learner (downstream) take up the slack
- Here's another famous trick that exploits these insights....

BLOOM FILTERS

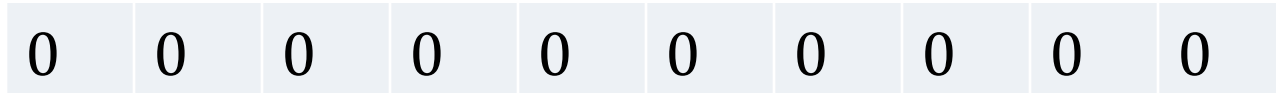
Bloom filters

- Interface to a Bloom filter
 - `BloomFilter(int maxSize, double p);`
 - `void bf.add(String s);` // insert `s`
 - `bool bd.contains(String s);`
 - // If `s` was added return true;
 - // else with probability at least $1-p$ return false;
 - // else with probability at most p return true;
 - I.e., a noisy “set” where you can test membership (and that’s it)

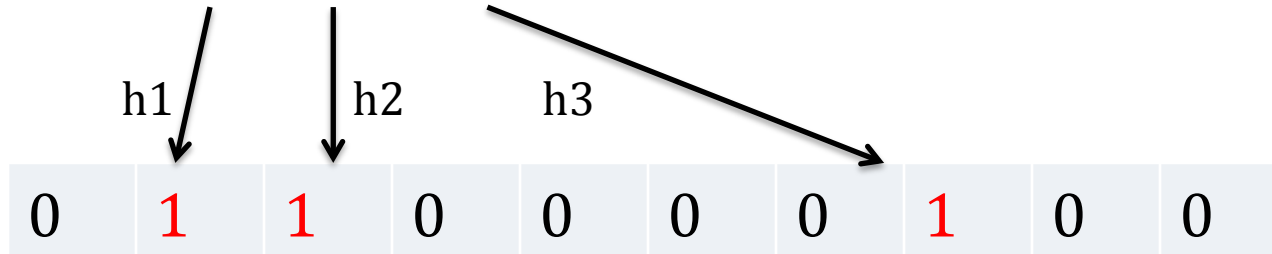
Bloom filters

- An implementation
 - Allocate M bits, $\text{bit}[0] \dots, \text{bit}[1-M]$
 - Pick K hash functions $\text{hash}(1,s), \text{hash}(2,s), \dots$
 - E.g: $\text{hash}(i,s) = \text{hash}(s + \text{randomString}[i])$
 - To add string s:
 - For $i=1$ to k , set $\text{bit}[\text{hash}(i,s)] = 1$
 - To check contains(s):
 - For $i=1$ to k , test $\text{bit}[\text{hash}(i,s)]$
 - Return “true” if they’re all set; otherwise, return “false”
 - We’ll discuss how to set M and K soon, but for now:
 - Let $M = 1.5 * \text{maxSize}$ *// less than two bits per item!*
 - Let $K = 2 * \log(1/p)$ *// about right with this M*

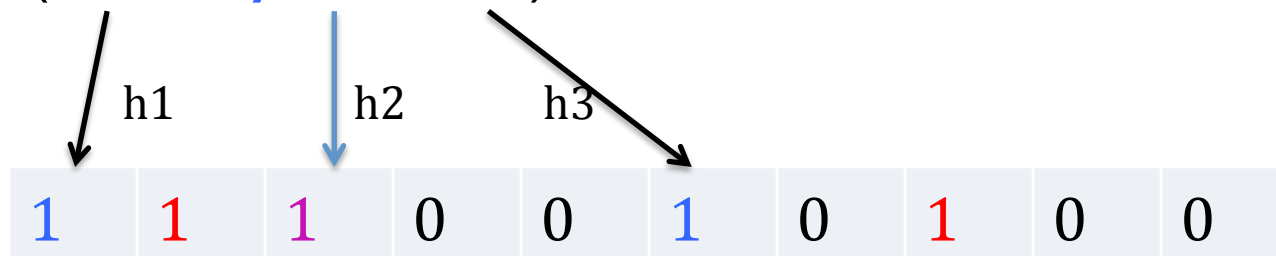
Bloom filters



bf.add("fred flintstone"):



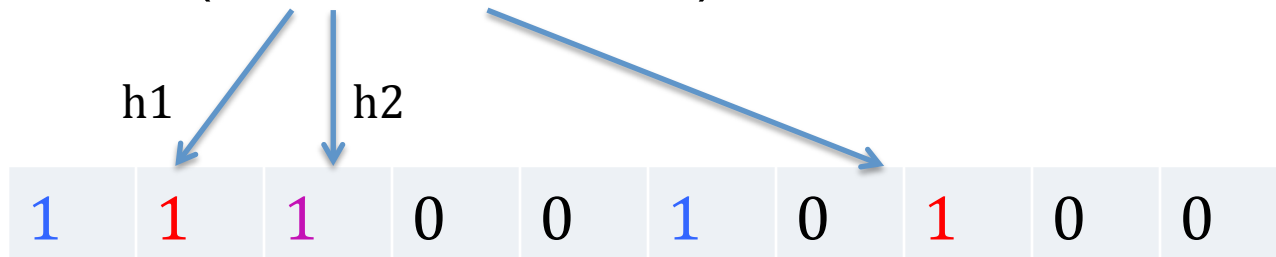
bf.add("barney rubble"):



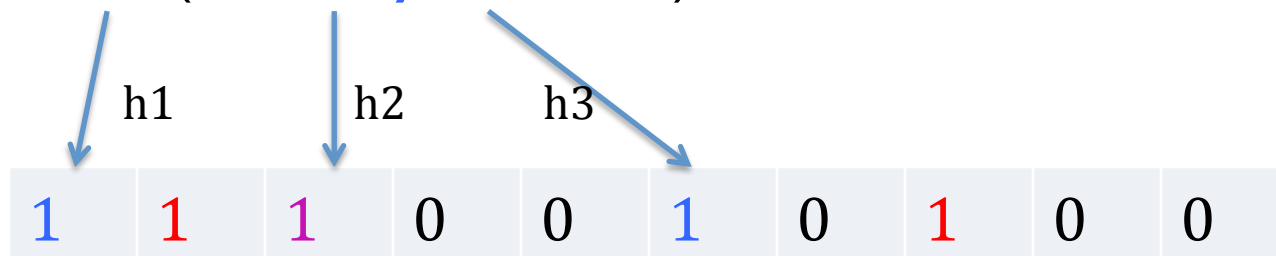
Bloom filters



bf.contains (“fred flintstone”):



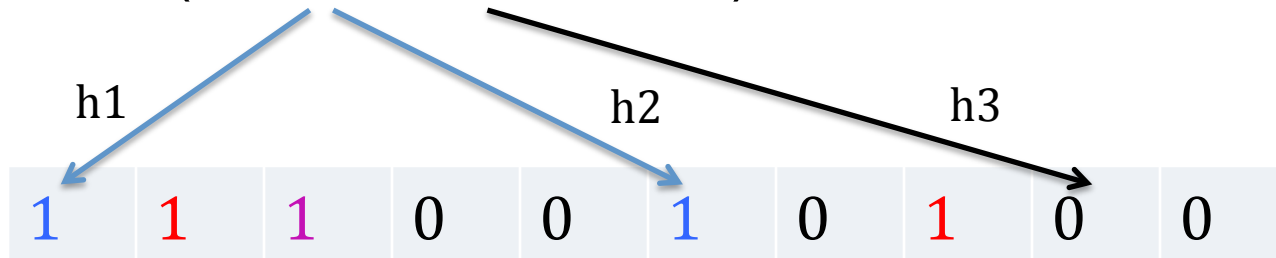
bf.contains (“barney rubble”):



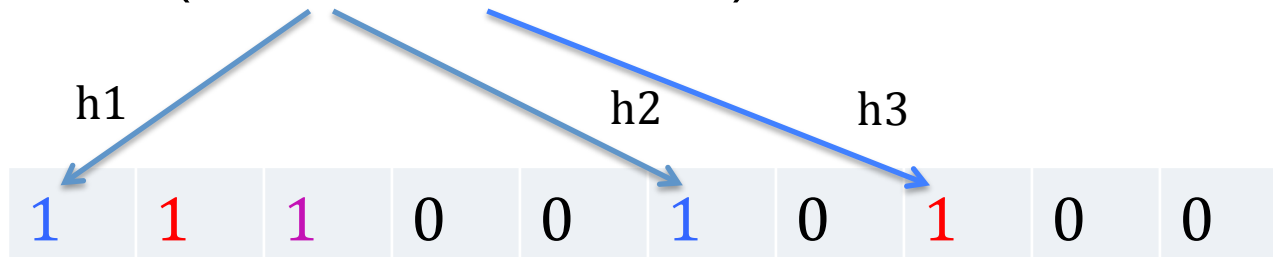
Bloom filters



bf.contains(“**wilma flintstone**”):



bf.contains(“**wilma flintstone**”):



Bloom filters: analysis

- Analysis (m bits, k hashers):

- Assume $\text{hash}(i,s)$ is a random function

- Look at $\Pr(\text{bit } j \text{ is unset after } n \text{ add's})$: $\left(1 - \frac{1}{m}\right)^{kn}$

- ... and $\Pr(\text{collision}) = \Pr(\text{not all } k \text{ bits set})$

$$f(m,n,k) = \left(1 - \left[1 - \frac{1}{m}\right]^{kn}\right)^k \approx \left(1 - e^{-kn/m}\right)^k$$

- fix m and n and minimize k :

$$k = \frac{m}{n} \ln 2 \approx 0.7 \frac{m}{n}$$

Bloom filters

$$\left(1 - \left[1 - \frac{1}{m}\right]^{kn}\right)^k \approx \left(1 - e^{-kn/m}\right)^k$$

- Analysis:
 - Plug optimal $k=m/n \cdot \ln(2)$ back into $\Pr(\text{collision})$:

$$f(m,n) = p = \left(1 - e^{-(m/n \ln 2)n/m}\right)^{(m/n \ln 2)}$$

- Now we can fix any two of p, n, m and solve for the 3rd:
E.g., the value for m in terms of n and p :

$$m = -\frac{n \ln p}{(\ln 2)^2}.$$

Bloom filters

- Interface to a Bloom filter
 - `BloomFilter(int maxSize /* n */, double p);`
 - `void bf.add(String s); // insert s`
 - `bool bf.contains(String s);`
 - `// If s was added return true;`
 - `// else with probability at least $1-p$ return false;`
 - `// else with probability at most p return true;`
 - I.e., a noisy “set” where you can test membership (and that’s it)

Bloom filters: demo

Bloom filters

- An example application
 - Finding items in “sharded” data
 - Easy if you know the sharding rule
 - Harder if you don’t (like Google n-grams)

```
furter:google_ngram wcohen$ ls -alh *2gram* | tail
-rw-rw-rw- 1 13527 _lpoperator 264M Sep 17 2011 googlebooks-eng-all-2gram-20090715-90.csv.zip
-rw-rw-rw- 1 13527 _lpoperator 264M Sep 17 2011 googlebooks-eng-all-2gram-20090715-91.csv.zip
-rw-rw-rw- 1 13527 _lpoperator 264M Sep 17 2011 googlebooks-eng-all-2gram-20090715-92.csv.zip
-rw-rw-rw- 1 13527 _lpoperator 264M Sep 17 2011 googlebooks-eng-all-2gram-20090715-93.csv.zip
-rw-rw-rw- 1 13527 _lpoperator 264M Sep 17 2011 googlebooks-eng-all-2gram-20090715-94.csv.zip
-rw-rw-rw- 1 13527 _lpoperator 263M Sep 17 2011 googlebooks-eng-all-2gram-20090715-95.csv.zip
-rw-rw-rw- 1 13527 _lpoperator 264M Sep 17 2011 googlebooks-eng-all-2gram-20090715-96.csv.zip
-rw-rw-rw- 1 13527 _lpoperator 264M Sep 17 2011 googlebooks-eng-all-2gram-20090715-97.csv.zip
-rw-rw-rw- 1 13527 _lpoperator 264M Sep 17 2011 googlebooks-eng-all-2gram-20090715-98.csv.zip
-rw-rw-rw- 1 13527 _lpoperator 264M Sep 17 2011 googlebooks-eng-all-2gram-20090715-99.csv.zip
```


Bloom filters

- An example application
 - Finding items in “sharded” data
 - Easy if you know the sharding rule
 - Harder if you don’t (like Google n-grams)
- Simple idea:
 - Build a BF of the contents of each shard
 - To look for *key*, load in the BF’s one by one, and search only the shards that probably contain *key*
 - Analysis: you won’t miss anything, you might look in some extra shards
 - You’ll hit $O(1)$ extra shards if you set $p=1/\text{\#shards}$

Bloom filters

- An example application
 - discarding singleton features from a classifier
- Scan through data once and check each w :
 - if `bf1.contains(w)`: `bf2.add(w)`
 - else `bf1.add(w)`
- Now:
 - `bf1.contains(w)` \Leftrightarrow w appears \geq once
 - `bf2.contains(w)` \Leftrightarrow w appears ≥ 2 x
- Then train, ignoring words not in `bf2`

Bloom filters

- An example application
 - discarding rare features from a classifier
 - seldom hurts much, can speed up experiments
- Scan through data once and check each w :
 - if `bf1.contains(w)`:
 - if `bf2.contains(w)`: `bf3.add(w)`
 - else `bf2.add(w)`
 - else `bf1.add(w)`
- Now:
 - `bf2.contains(w)` \Leftrightarrow w appears $\geq 2x$
 - `bf3.contains(w)` \Leftrightarrow w appears $\geq 3x$
- Then train, ignoring words not in `bf3`

THE COUNT-MIN SKETCH

A variant of feature hashing

- Hash each feature *multiple times* with different hash functions
- Now, each w has k chances to *not* collide with another useful w'
- Get multiple hash functions as in Bloom filters
- *Part Bloom filter, part hash kernel*
 - *but predates either, called “count-min sketch” -- Cormode and Muthukrishnan*

Bloom filters

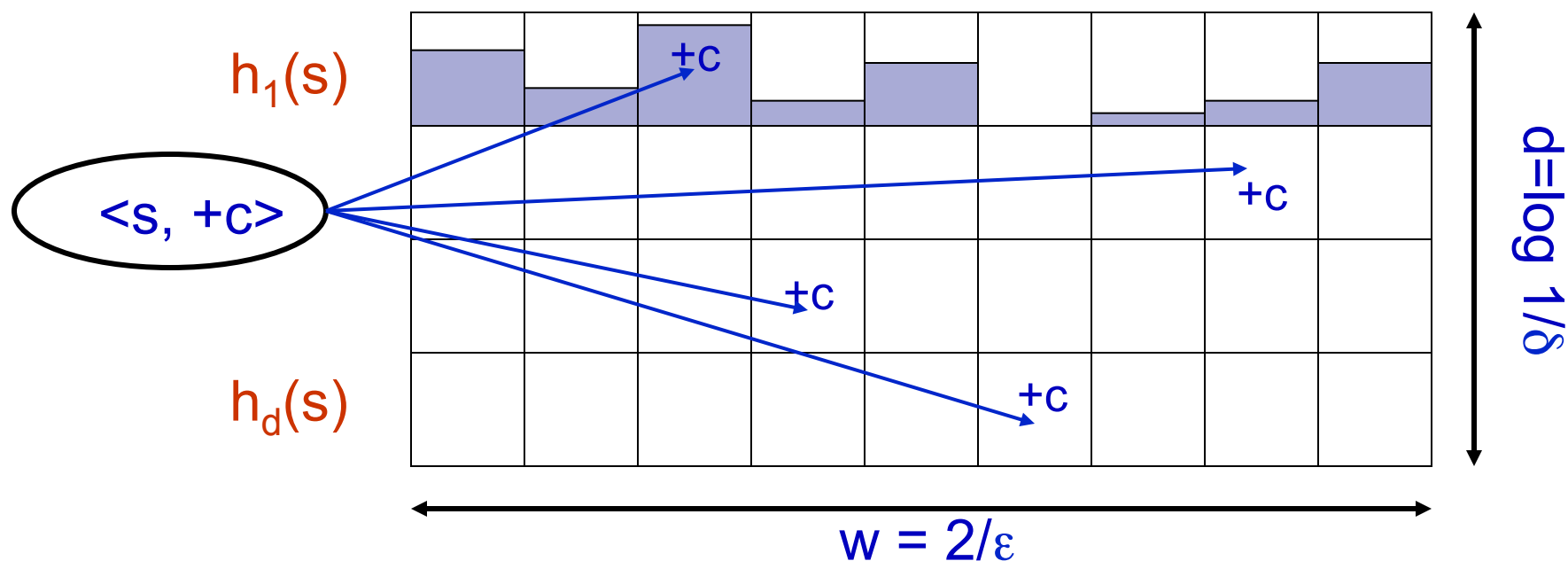
- An implementation
 - Allocate M bits, $\text{bit}[0] \dots, \text{bit}[1-M]$
 - Pick K hash functions $\text{hash}(1,s), \text{hash}(2,s), \dots$
 - E.g: $\text{hash}(i,s) = \text{hash}(s + \text{randomString}[i])$
 - To add string s:
 - For $i=1$ to k , set $\text{bit}[\text{hash}(i,s)] = 1$
 - To check contains(s):
 - For $i=1$ to k , test $\text{bit}[\text{hash}(i,s)]$
 - Return “true” if they’re all set; otherwise, return “false”
 - We’ll discuss how to set M and K soon, but for now:
 - Let $M = 1.5 * \text{maxSize}$ *// less than two bits per item!*
 - Let $K = 2 * \log(1/p)$ *// about right with this M*

Bloom Filter → Count-min sketch

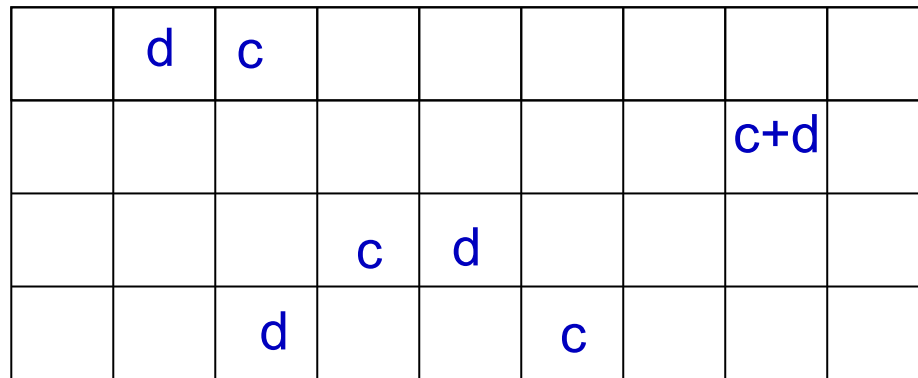
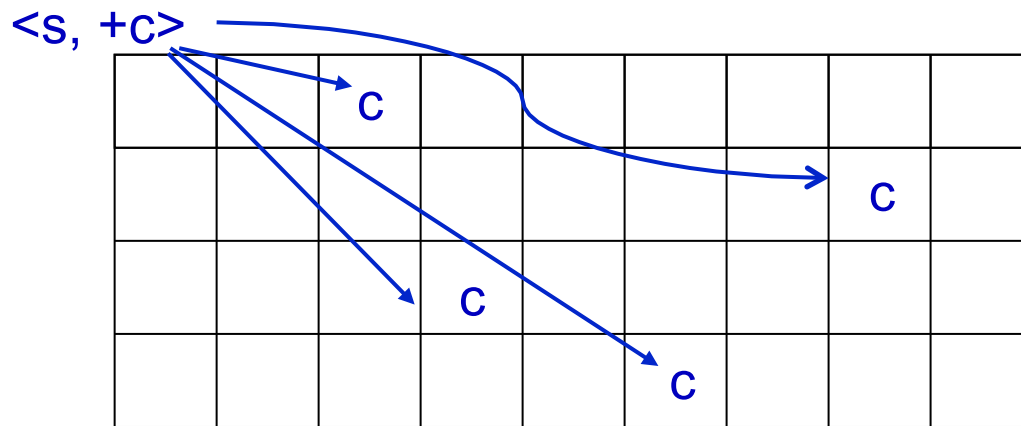
- An implementation
 - Allocate a matrix CM with d rows, w columns
 - Pick d hash functions $h_1(s), h_2(s), \dots$
 - To increment counter $A[s]$ for s by c
 - For $i=1$ to d , set $CM[i, hash(i,s)] += c$
 - To retrieve value of $A[s]$:
 - For $i=1$ to d , retrieve $M[i, hash(i,s)]$
 - Return **minimum** of these values
 - Similar idea as Bloom filter:
 - if there are d collisions, you return a value that's too large; otherwise, you return the correct value.

Question: what does this look like if $d=1$?

CM Sketch Structure



- Each string is mapped to one bucket per row
- Estimate $A[j]$ by taking $\min_k \{ CM[k, h_k(j)] \}$
- Errors are always **over-estimates** i.e. with prob $> 1-\delta$
- Analysis: $d = \log 1/\delta$, $w = 2/\epsilon \rightarrow$ error is usually less than $\epsilon \|A\|_1$



- You can find the *sum* of two sketches by doing element-wise summation
- Also, you can compute a weighted sum of MC sketches
- Same result as adding $\langle s, +c \rangle$ and then $\langle t, +d \rangle$ to an empty sketch

CM Sketch Guarantees

- *[Cormode, Muthukrishnan '04]* CM sketch guarantees approximation error on point queries less than $\epsilon \|A\|_1$ in space $O(1/\epsilon \log 1/\delta)$
 - Probability of more error is less than $1-\delta$
- This is sometimes enough:
 - Estimating a multinomial: if $A[s] = \Pr(s|\dots)$ then $\|A\|_1 = 1$
 - Multiclass classification: if $A_x[s] = \Pr(x \text{ in class } s)$ then $\|A_x\|_1$ is probably small, since most x 's will be in only a few classes

CM Sketch Guarantees

- *[Cormode, Muthukrishnan '04]* CM sketch guarantees approximation error on point queries less than $\epsilon \|A\|_1$ in space $O(1/\epsilon \log 1/\delta)$
- CM sketches are also accurate for *skewed* values---i.e., only a few entries s with large $A[s]$

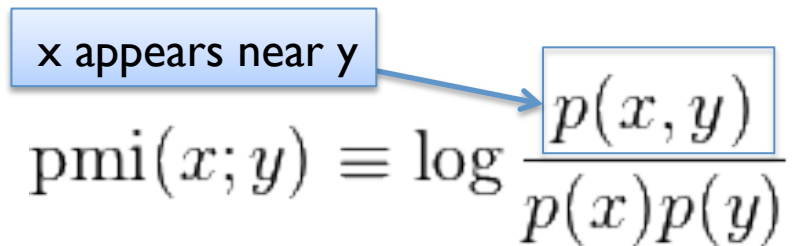
Lemma 1 (Cormode and Muthukrishnan [6], Eqn 5.1) *Let y be an vector, and let \tilde{y}_i be the estimate given by a count-min sketch of width w and depth d for y_i . Let the k largest components of y be $y_{\sigma_1}, \dots, y_{\sigma_k}$, and let $t_k = \sum_{k' > k} y_{\sigma_{k'}}$ be the weight of the “tail” of y . If $w \geq \frac{1}{3k}$, $w > \frac{\epsilon}{\eta}$ and $d \geq \ln \frac{3}{2} \ln \frac{1}{\delta}$, then $\tilde{y}_i \leq y_i + \eta t_k$ with probability at least $1-\delta$.*

Theorem 3 (Cormode and Muthukrishnan [6], Theorem 5.1) *Let y represent a Zipf-like distribution with parameter z . Then with probability at least $1-\delta$, y can be approximated to within error η by a count-min sketch of width $O(\eta^{-\min(1, 1/z)})$ and depth $O(\ln \frac{1}{\delta})$.*

An Application of a Count-Min Sketch

- Problem: find the semantic orientation of a work (positive or negative) using a large corpus.
- Idea:
 - positive words co-occur more frequently than expected near positive words; likewise for negative words
 - so pick a few pos/neg seeds and compute

x appears near y


$$\text{pmi}(x; y) \equiv \log \frac{p(x, y)}{p(x)p(y)}$$

$$\text{SO}(w) = \sum_{p \in \text{Pos}} \text{PMI}(p, w) - \sum_{n \in \text{Neg}} \text{PMI}(n, w)$$

An Application of a Count-Min Sketch

x appears near y

$$\text{pmi}(x; y) \equiv \log \frac{p(x, y)}{p(x)p(y)}$$

$$\text{SO}(w) = \sum_{p \in \text{Pos}} \text{PMI}(p, w) - \sum_{n \in \text{Neg}} \text{PMI}(n, w)$$

Example: Turney, 2002 used two seeds, “excellent” and “poor”

$$\text{SO}(\text{phrase}) = \log_2 \left(\frac{\text{hits}(\text{phrase NEAR 'excellent'}) \text{hits}('excellent')}{\text{hits}(\text{phrase NEAR 'poor'}) \text{hits}('excellent')} \right)$$

In general, $\text{SO}(w)$ can be written in terms of logs of products of counters for w , with and without seeds

An Application of a Count-Min Sketch


- Use 2B counters, 5 hash functions, “near” means a 7-word window, GigaWord (10 Gb) and GigaWord + Web news 50 Gb)

Data	Exact	CM-CU	CMM-CU	LCU-WS
GW	<i>74.2</i>	<i>74.0</i>	65.3	<i>72.9</i>
GWB50	81.2	80.9	74.9	78.3

Table 2: Evaluating Semantic Orientation on accuracy metric using several sketches of 2 billion counters against exact. Bold and italic numbers denote no statistically significant difference.

An Application of a Count-Min Sketch

CM-CU: CM with “conservative update” - for $\langle j, +c \rangle$ increment counters just enough to make the new estimate for j grow by c



Data	Exact	CM-CU	CMM-CU	LCU-WS
GW	74.2	<i>74.0</i>	65.3	72.9
GWB50	81.2	80.9	74.9	78.3

Table 2: Evaluating Semantic Orientation on accuracy metric using several sketches of 2 billion counters against exact. Bold and italic numbers denote no statistically significant difference.

LOCALITY SENSITIVE HASHING (LSH)

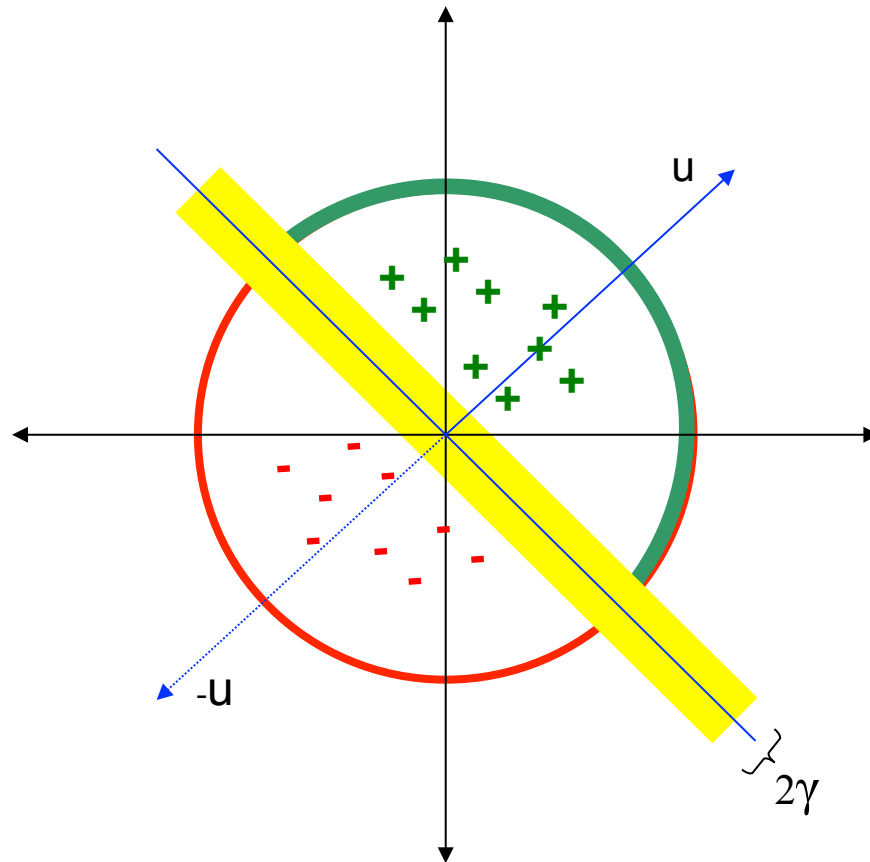
LSH: key ideas

- Goal:
 - map feature vector \mathbf{x} to bit vector \mathbf{bx}
 - ensure that \mathbf{bx} preserves “similarity”

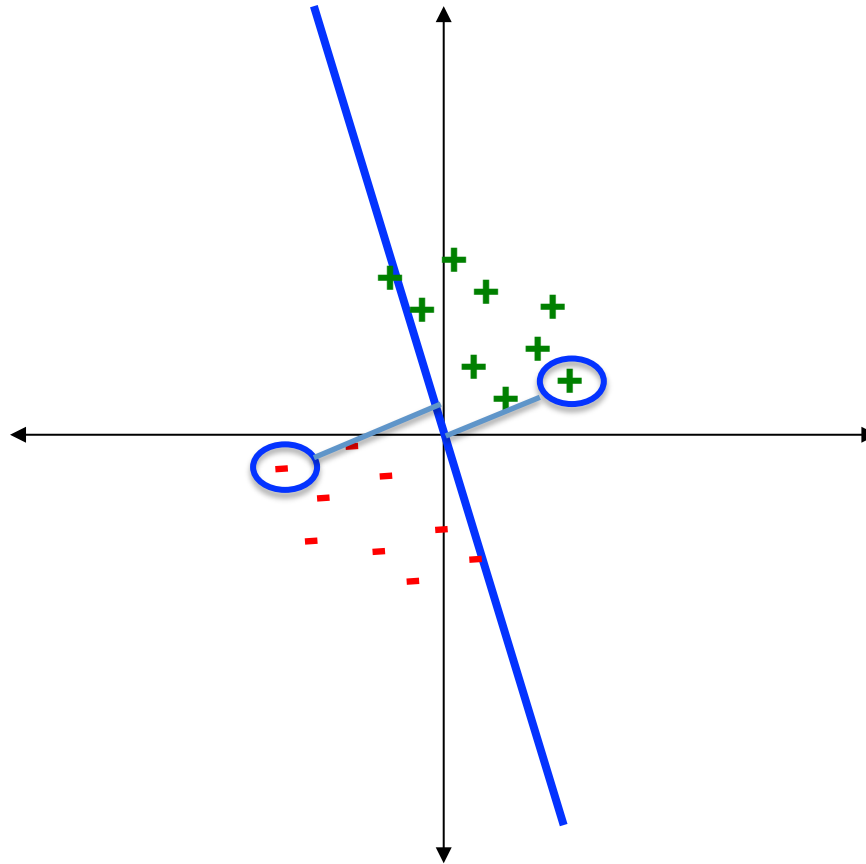
Random Projections



Random projections



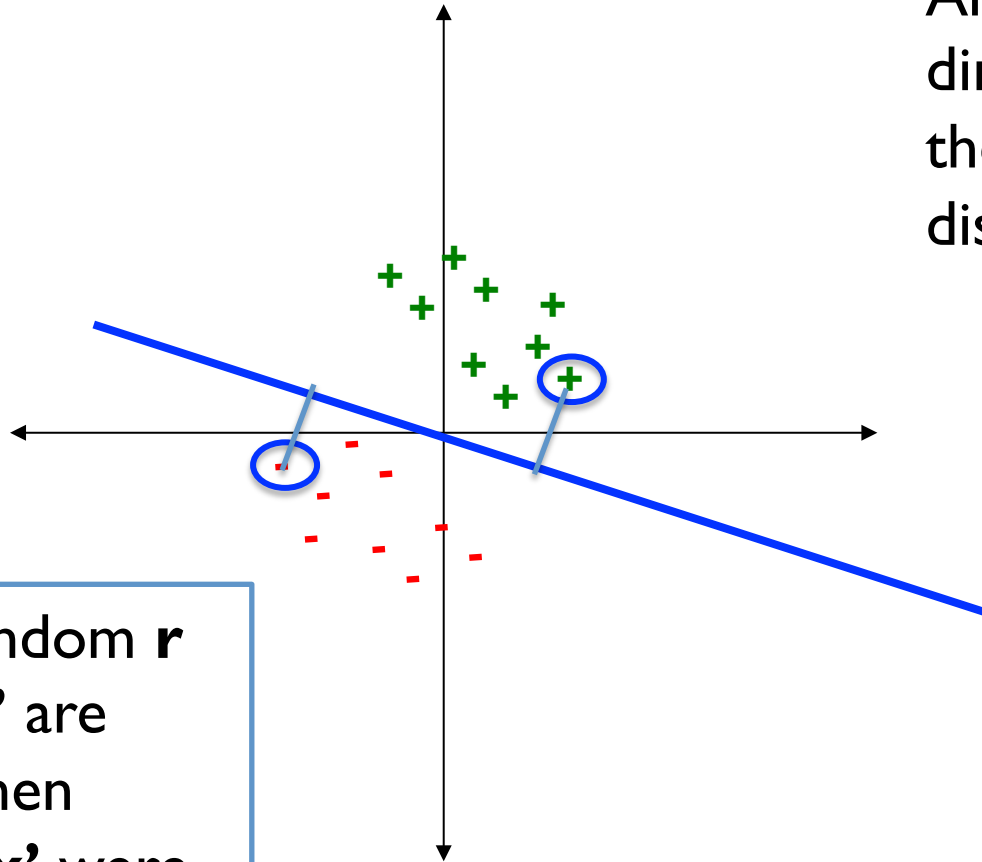
Random projections



To make those points “close” we need to project to a direction orthogonal to the line between them

Random projections

Any other direction will keep the distant points distant.



So if I pick a random \mathbf{r} and $\mathbf{r} \cdot \mathbf{x}$ and $\mathbf{r} \cdot \mathbf{x}'$ are closer than γ then *probably* \mathbf{x} and \mathbf{x}' were close to start with.

LSH: key ideas

- Goal:
 - map feature vector \mathbf{x} to bit vector \mathbf{bx}
 - ensure that \mathbf{bx} preserves “similarity”
- Basic idea: use *random projections* of \mathbf{x}
 - Repeat many times:
 - Pick a random hyperplane \mathbf{r} by picking random weights for each feature (say from a Gaussian)
 - Compute the inner product of \mathbf{r} with \mathbf{x}
 - Record if \mathbf{x} is “close to” \mathbf{r} ($\mathbf{r} \cdot \mathbf{x} \geq 0$)
 - the next bit in \mathbf{bx}
 - Theory says that if \mathbf{x}' and \mathbf{x} have small cosine distance then \mathbf{bx} and \mathbf{bx}' will have small Hamming distance

Online Generation of Locality Sensitive Hash Signatures

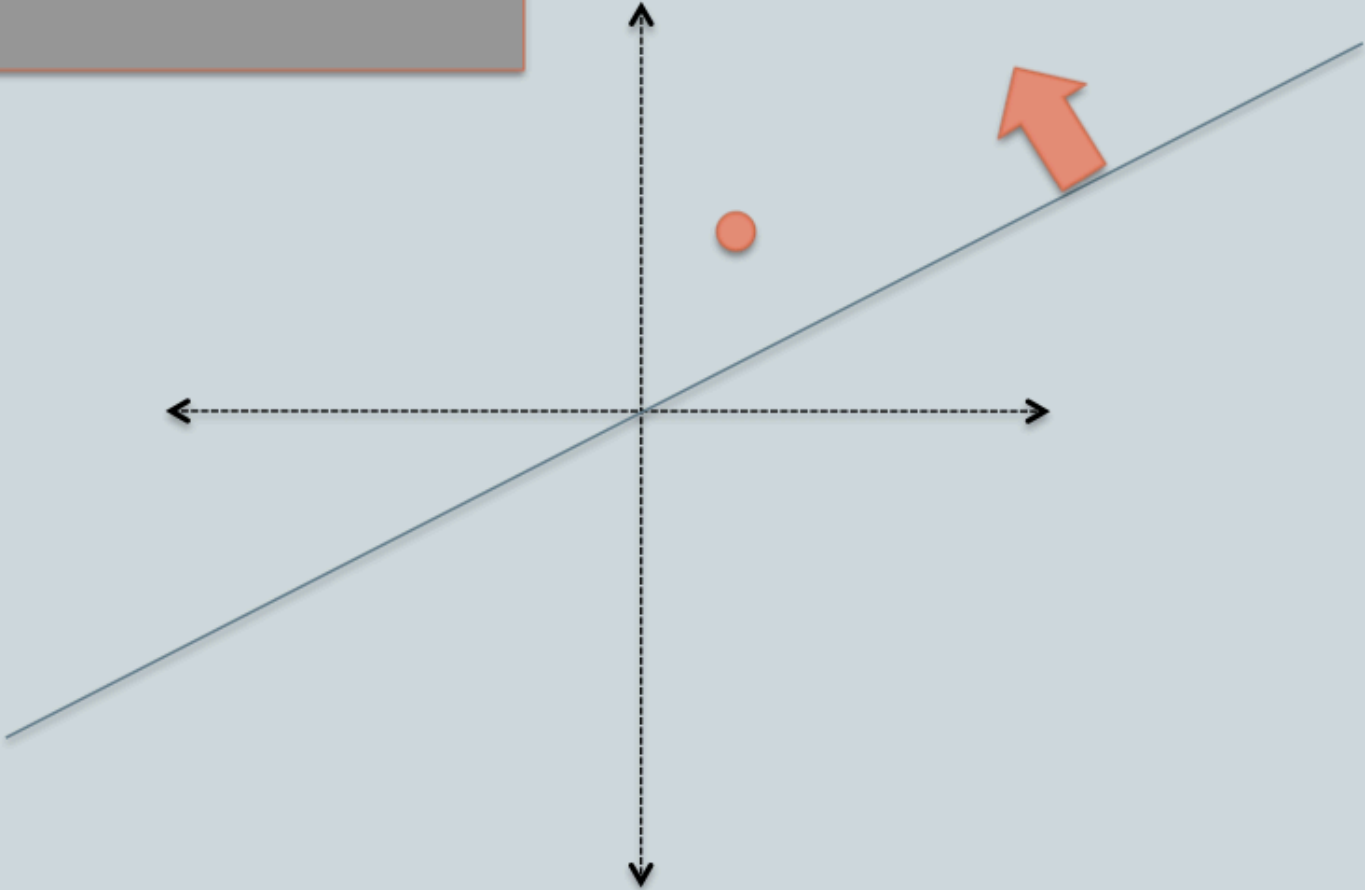
Benjamin Van Durme and Ashwin Lall

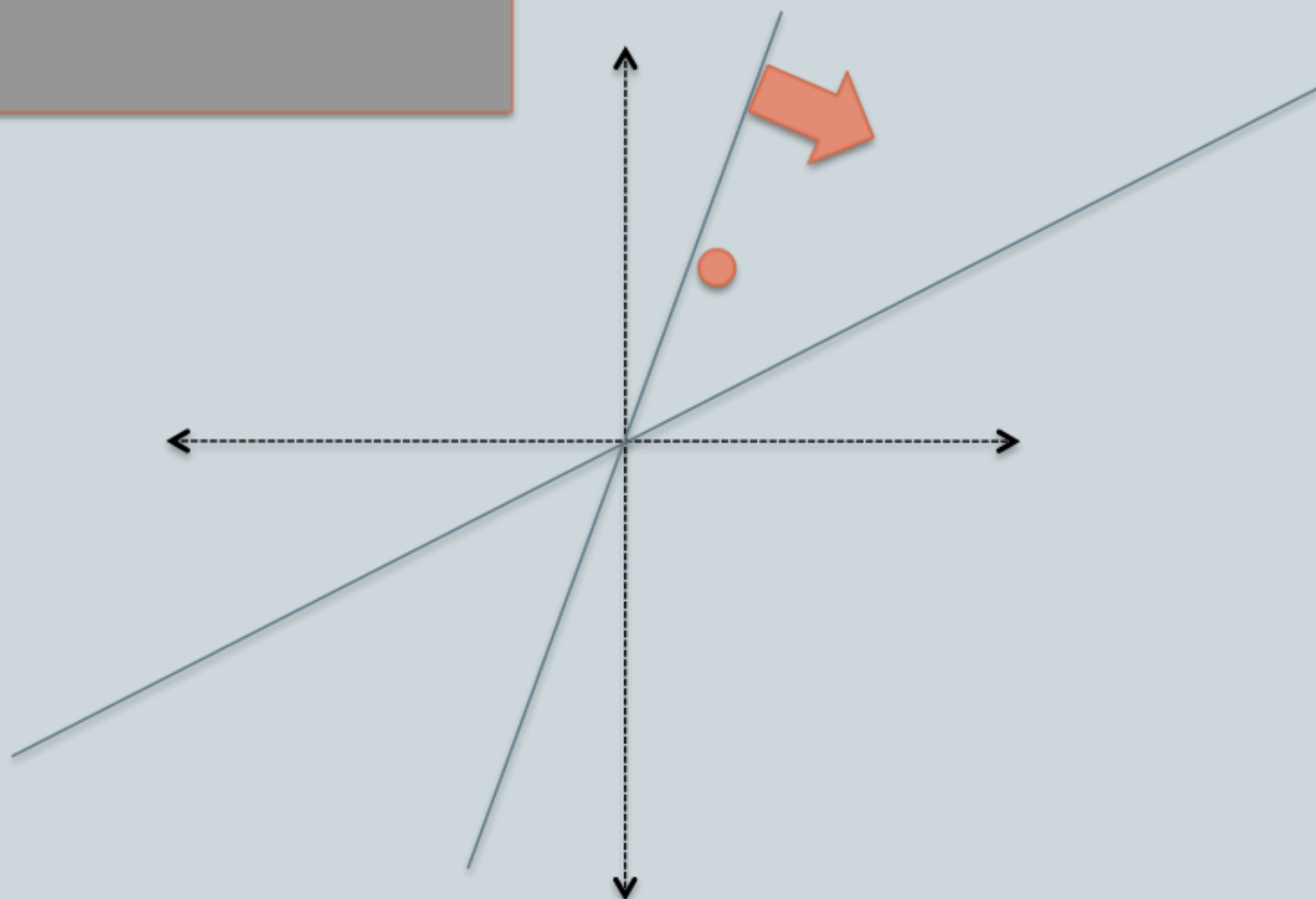


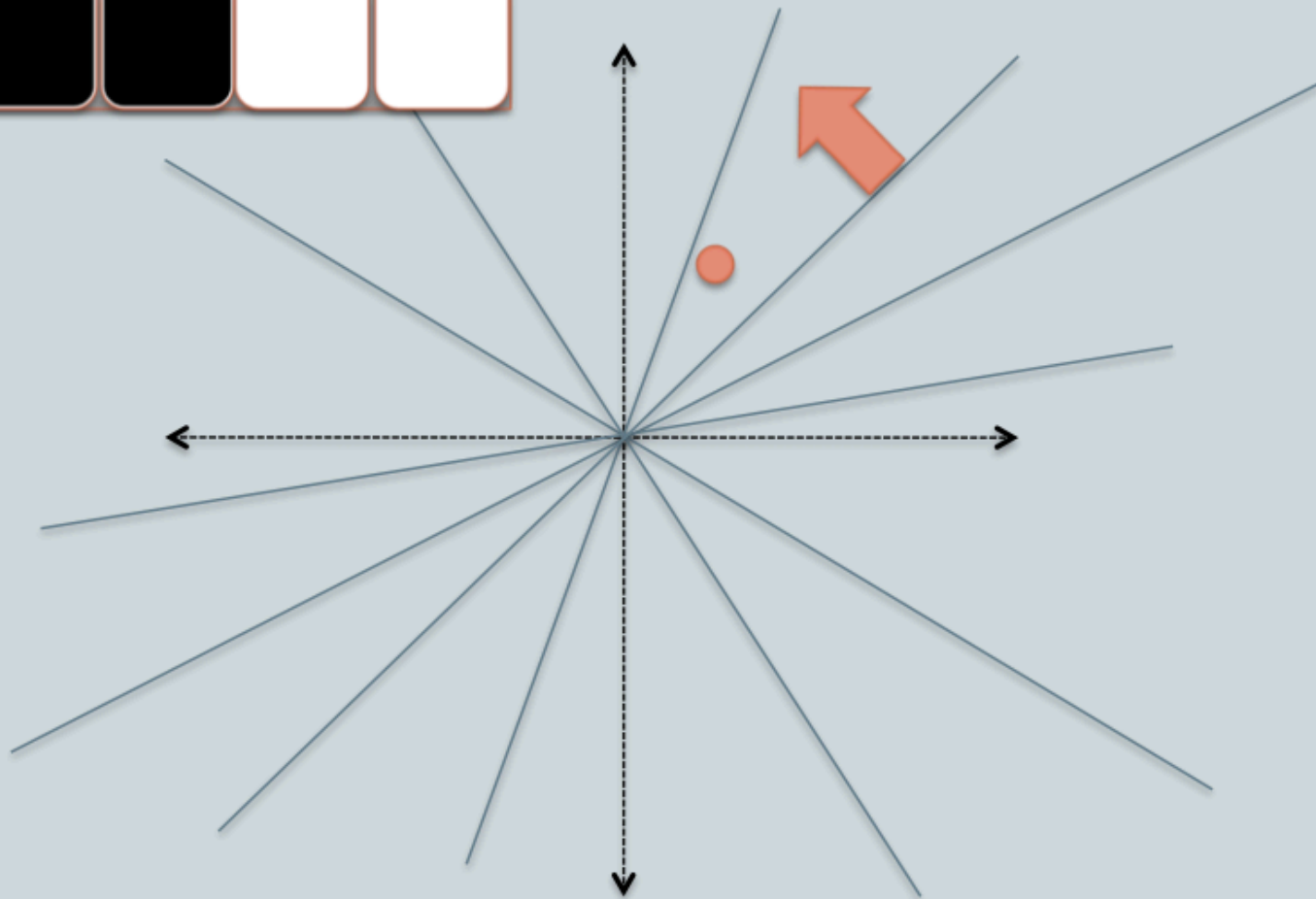
human language technology
center of excellence

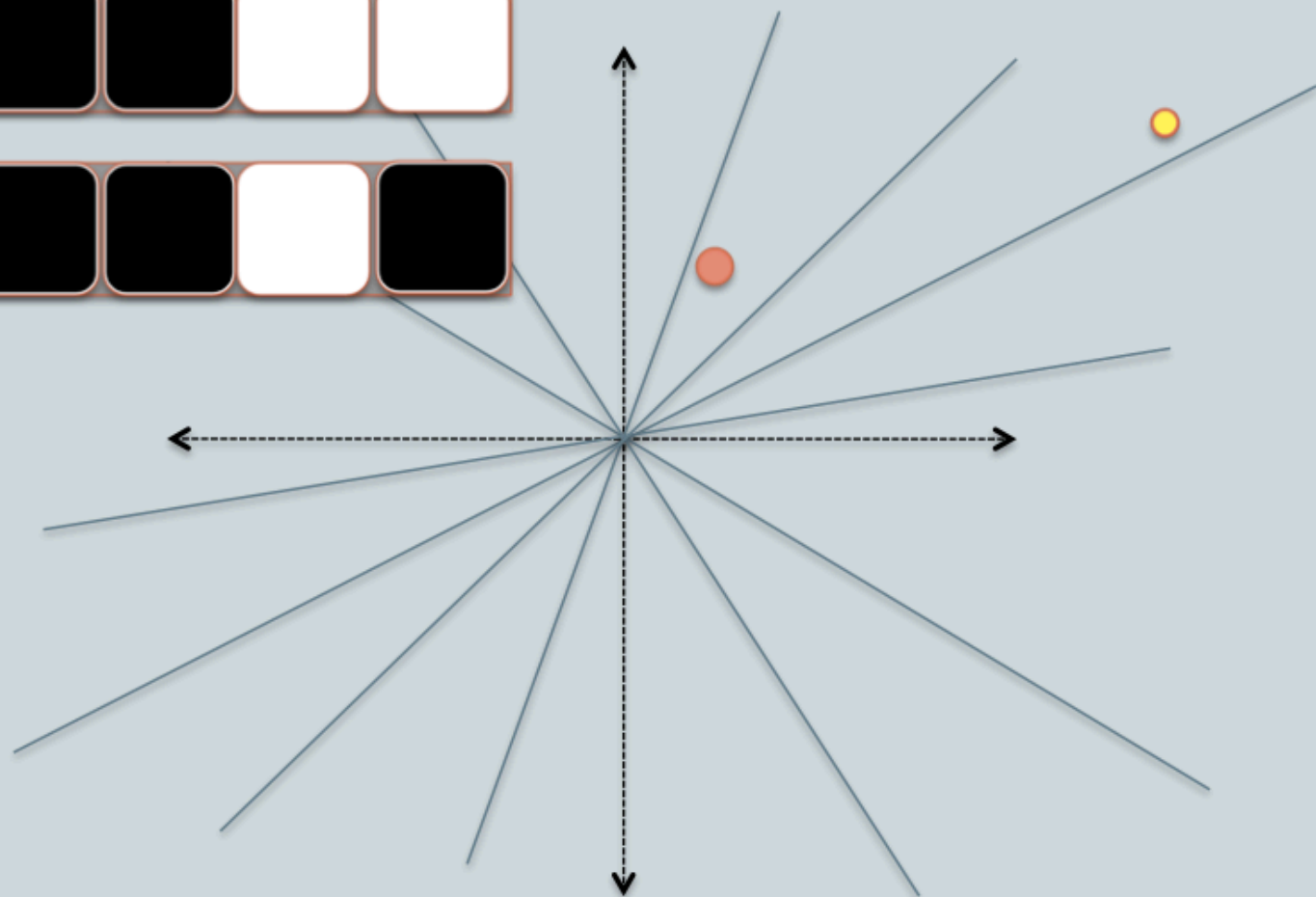
JOHNS HOPKINS
UNIVERSITY

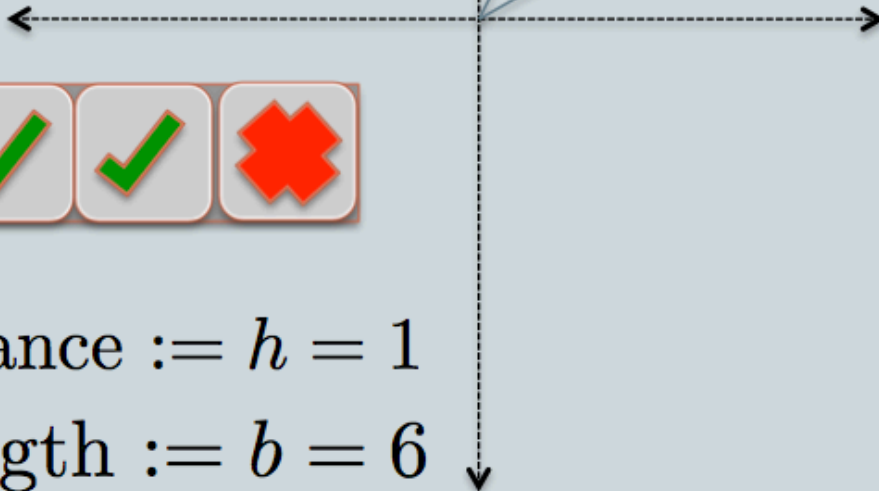
DENISON
UNIVERSITY







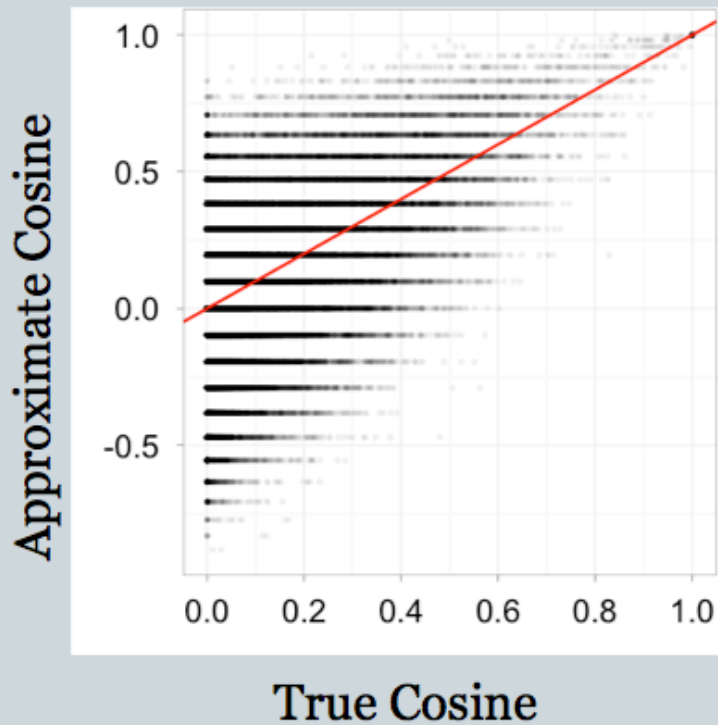




Hamming Distance $:= h = 1$
Signature Length $:= b = 6$

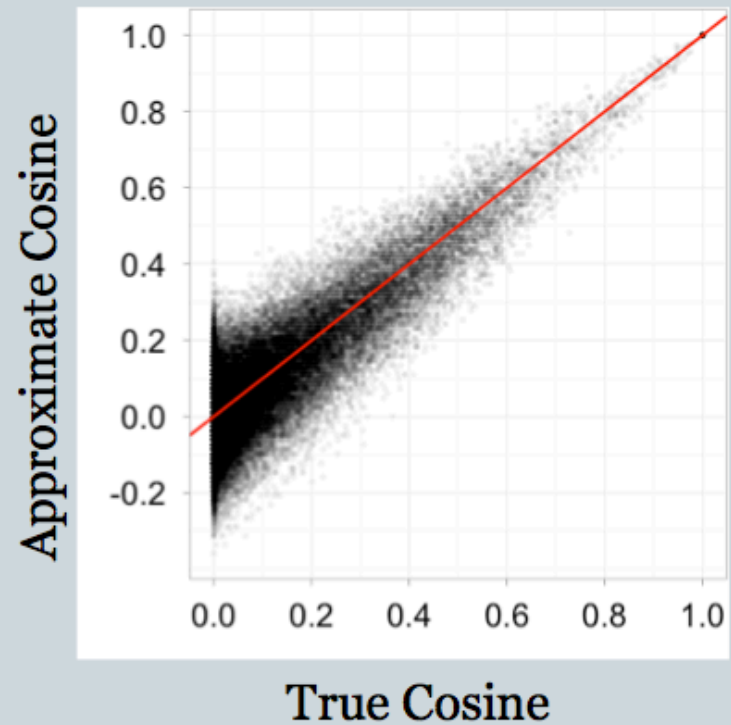
$$\begin{aligned}\cos(\theta) &\approx \cos\left(\frac{h}{b}\pi\right) \\ &= \cos\left(\frac{1}{6}\pi\right)\end{aligned}$$

32 bit signatures



Cheap

256 bit signatures



Accurate

LSH applications

- Compact storage of data
 - and we can still compute similarities
- LSH also gives very fast approximations:
 - approx nearest neighbor method
 - just look at other items with $\mathbf{bx}' = \mathbf{bx}$
 - also very fast nearest-neighbor methods for Hamming distance
 - very fast clustering
 - cluster = all things with same \mathbf{bx} vector

Locality Sensitive Hashing (LSH) and Pooling Random Values

LSH algorithm

- Naïve algorithm:
 - Initialization:
 - For $i=1$ to outputBits:
 - For each feature f :
 - » Draw $r(f,i) \sim \text{Normal}(0,1)$
 - Given an instance \mathbf{x}
 - For $i=1$ to outputBits:
 - LSH[i] =
 $\text{sum}(\mathbf{x}[f] * r[i,f] \text{ for } f \text{ with non-zero weight in } \mathbf{x}) > 0 ?$
1 : 0
 - Return the bit-vector LSH

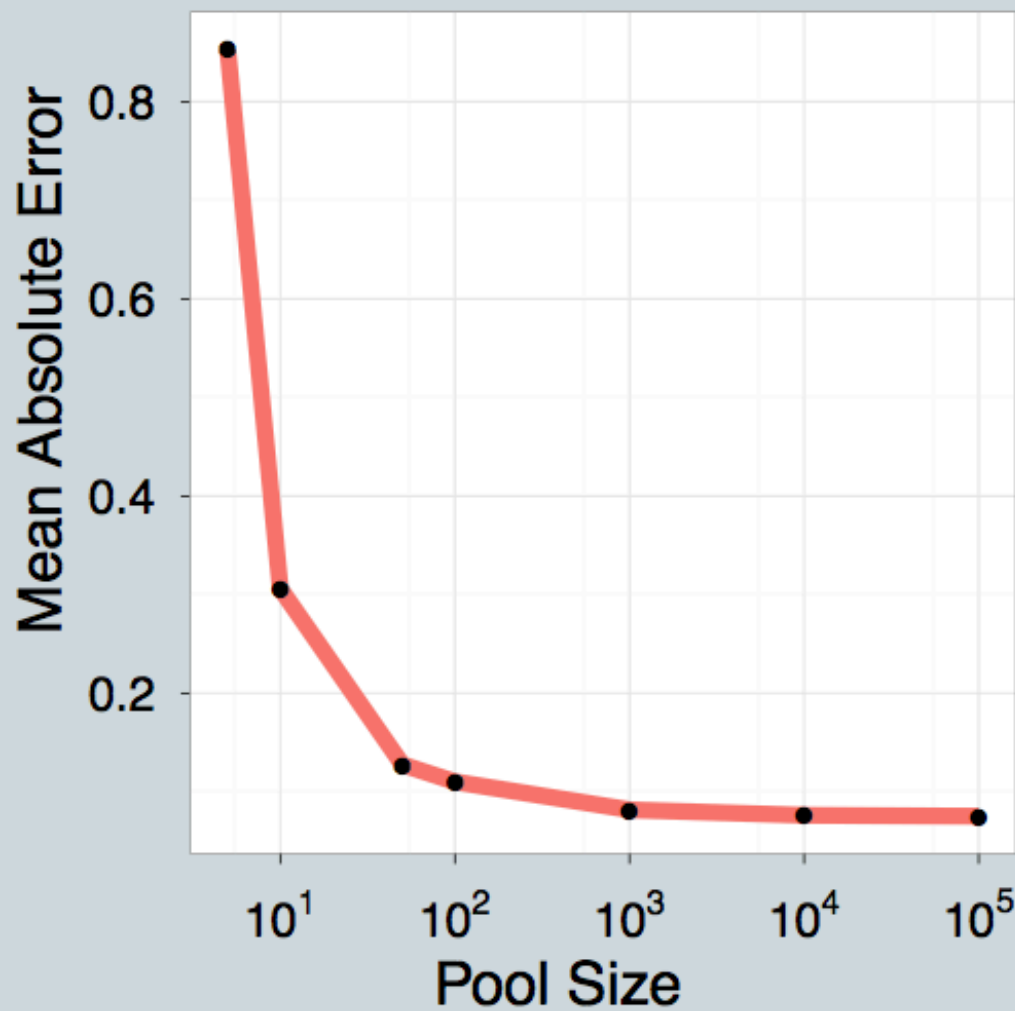
LSH algorithm

- But: storing the k *classifiers* is expensive in high dimensions
 - For each of 256 bits, a dense vector of weights for every feature in the vocabulary
- Storing seeds and random number generators:
 - Possible but somewhat fragile

LSH: “pooling” (van Durme)

- Better algorithm:
 - Initialization:
 - Create a pool:
 - Pick a random seed s
 - For $i=1$ to $poolSize$:
 - » Draw $pool[i] \sim \text{Normal}(0,1)$
 - For $i=1$ to $outputBits$:
 - Devise a random hash function $hash(i,f)$:
 - » E.g.: $hash(i,f) = \text{hashCode}(f) \text{ XOR } \text{randomBitString}[i]$
 - Given an instance \mathbf{x}
 - For $i=1$ to $outputBits$:
 - $LSH[i] = \text{sum}(\mathbf{x}[f] * \text{pool}[hash(i,f) \% \text{poolSize}] \text{ for } f \text{ in } \mathbf{x}) > 0 ? 1 : 0$
 - Return the bit-vector LSH

The Pooling Trick



LSH: key ideas: pooling

- Advantages:
 - with pooling, this is a compact re-encoding of the data
 - you don't need to store the \mathbf{r} 's, just the pool

Locality Sensitive Hashing (LSH) in an On-line Setting

LSH: key ideas: online computation

- Common task: distributional clustering
 - for a word w , $\mathbf{x}(w)$ is sparse vector of words that co-occur with w
 - cluster the w 's

$$\vec{v} \in \mathbb{R}^d$$

$$\vec{r}_i \sim N(0, 1)^d$$

$$h_i(\vec{v}) = \begin{cases} 1 & \text{if } \vec{v} \cdot \vec{r}_i \geq 0, \\ 0 & \text{otherwise.} \end{cases}$$

if $\vec{v} = \sum_j \vec{v}_j$

then $\vec{v} \cdot \vec{r}_i = \sum_j \vec{v}_j \cdot \vec{r}_i$

Break into local products

Online

$$h_{it}(\vec{v}) = \begin{cases} 1 & \text{if } \sum_j^t \vec{v}_j \cdot \vec{r}_i \geq 0, \\ 0 & \text{otherwise.} \end{cases}$$

Algorithm 1 STREAMING LSH ALGORITHM

Parameters:

m : size of pool

d : number of bits (size of resultant signature)

s : a random seed

h_1, \dots, h_d : hash functions mapping $\langle s, f_i \rangle$ to $\{0, \dots, m-1\}$

INITIALIZATION:

- 1: Initialize floating point array $P[0, \dots, m-1]$
- 2: Initialize H , a hashtable mapping words to floating point arrays of size d
- 3: **for** $i := 0 \dots m-1$ **do**
- 4: $P[i] :=$ random sample from $N(0, 1)$, using s as seed

ONLINE:

- 1: **for** each word w in the stream **do**
- 2: **for** each feature f_i associated with w **do**
- 3: **for** $j := 1 \dots d$ **do**
- 4: $H[w][j] := H[w][j] + P[h_j(s, f_i)]$

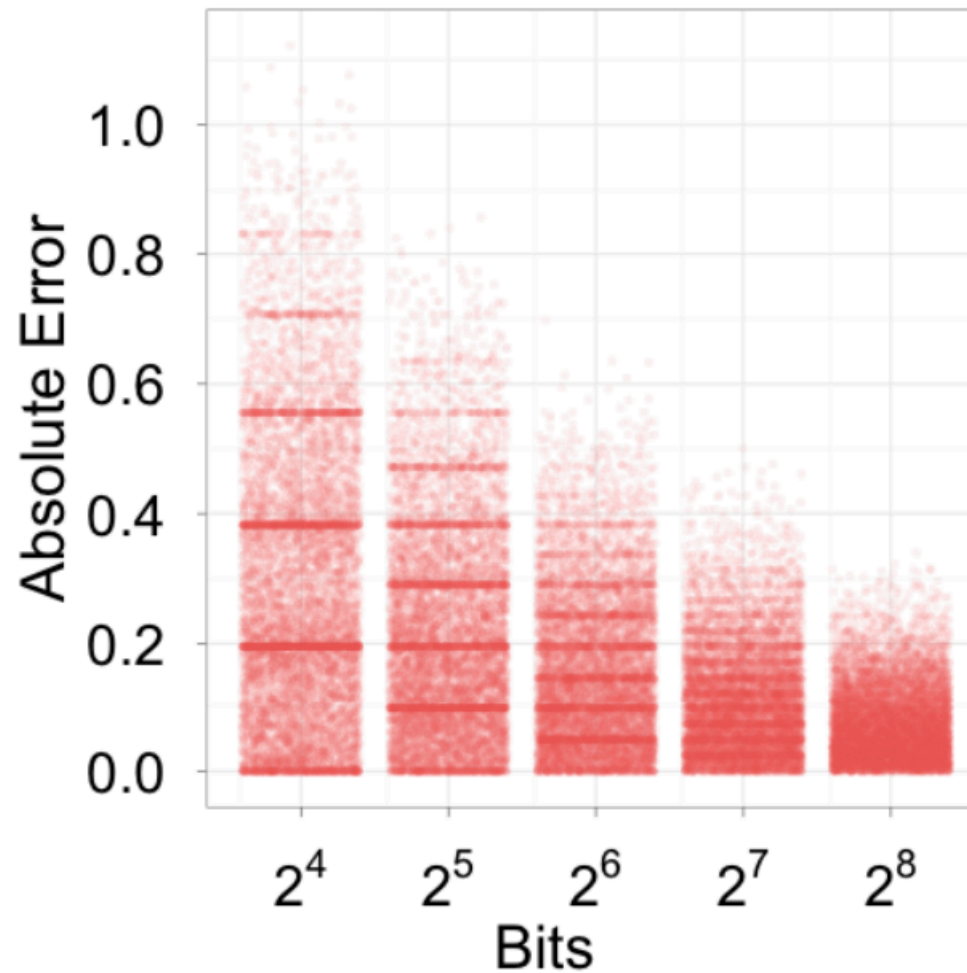
SIGNATURECOMPUTATION:

- 1: **for** each $w \in H$ **do**
 - 2: **for** $i := 1 \dots d$ **do**
 - 3: **if** $H[w][i] > 0$ **then**
 - 4: $S[w][i] := 1$
 - 5: **else**
 - 6: $S[w][i] := 0$
-

Experiment

- Corpus: 700M+ tokens, 1.1M distinct bigrams
- For each, build a feature vector of words that co-occur near it, using on-line LSH
- Check results with 50,000 actual vectors

Experiment



Closest based on true cosine

London

Milan_{.97}, Madrid_{.96}, Stockholm_{.96}, Manila_{.95}, Moscow_{.95}
ASHER₀, Champaign₀, MANS₀, NOBLE₀, come₀
Prague₁, Vienna₁, suburban₁, synchronism₁, Copenhagen₂

London

Milan_{.97}, Madrid_{.96}, Stockholm_{.96}, Manila_{.95}, Moscow_{.95}
ASHER₀, Champaign₀, MANS₀, NOBLE₀, come₀
Prague₁, Vienna₁, suburban₁, synchronism₁, Copenhagen₂
Frankfurt₄, Prague₄, Taszar₅, Brussels₆, Copenhagen₆
Prague₁₂, Stockholm₁₂, Frankfurt₁₄, Madrid₁₄, Manila₁₄
Stockholm₂₀, Milan₂₂, Madrid₂₄, Taipei₂₄, Frankfurt₂₅

Closest based on 32 bit sig.'s

Cheap