

# Case Study 1: Estimating Click Probabilities

## Tackling an Unknown Number of Features with Sketching

Machine Learning for Big Data  
CSE547/STAT548, University of Washington

Emily Fox

January 14<sup>th</sup>, 2014

©Emily Fox 2014

1

# Ad Placement Strategies

RECAP

- Companies bid on ad prices

$C_1 \rightarrow \$10$   
 $C_2 \rightarrow \$20$   
 $C_3 \rightarrow \$100$

- Which ad wins? (many simplifications here)

Naively:  $C_3 \rightarrow \$100$

But: paid on clicks

Instead:

$$\begin{aligned}
 \text{e.g. } & \left\{ \begin{aligned} P(\text{click} | C_3) &= 0.01 & E[\$3] &= 0.01 \times 100 = \$1 \\ P(\text{click} | C_2) &= 0.5 & E[\$2] &= 0.5 \times 20 = \$10 \end{aligned} \right.
 \end{aligned}$$

The screenshot shows a Google search for "big data" with approximately 1,300,000,000 results. Several search results are visible, including articles from SAS, Dell, Oracle, and IBM. A red arrow on the right side of the page points to the search results, indicating the context for the ad placement strategies discussed in the text.

©Emily Fox 2014

2

# Learning Problem for Click Prediction

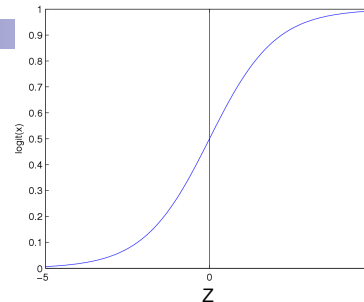
- Prediction task:  $X \rightarrow [0,1]$      $P(\text{click}=1|X)$  ← want to model this
- Features:  $X = (\text{feats of page}, \text{feats ad}, \text{feats user})$  ← Features
- Data:  $(X, Y) \rightarrow (\text{webpage 1}, \text{ad 7}, \text{user 25}, \text{time 12}) \rightarrow \text{click} = 1$ 
  - Batch: fixed dataset  $(X^1, Y^1) \dots (X^N, Y^N)$
  - Online: data as a stream  
 = user arrives at a page  $\rightarrow X^t \rightarrow$  predict  $\hat{Y}^t$ , click?  $\rightarrow$  observe  $Y^t$
- Many approaches (e.g., logistic regression, SVMs, naïve Bayes, decision trees, boosting....)
  - Focus on logistic regression, captures main concepts, ideas generalize to other approaches

©Emily Fox 2014

3

# Logistic Regression

Logistic function (or Sigmoid):  $\frac{1}{1 + \exp(-z)}$



- Learn  $P(Y|X)$  directly
  - Assume a particular functional form
  - Sigmoid applied to a linear function of the data:

$$P(Y = 0|X, W) = \frac{1}{1 + \exp(w_0 + \sum_i w_i X_i)}$$

linear function of features

log odds ← log linear model

**Features can be discrete or continuous!**

©Emily Fox 2014

4

## Standard v. Regularized Updates

- Maximum conditional likelihood estimate

$$\mathbf{w}^* = \arg \max_{\mathbf{w}} \ln \left[ \prod_{j=1}^N P(y^j | \mathbf{x}^j, \mathbf{w}) \right]$$

$$w_i^{(t+1)} \leftarrow w_i^{(t)} + \eta \sum_j x_i^j [y^j - \hat{P}(Y^j = 1 | \mathbf{x}^j, \mathbf{w}^{(t)})]$$

gradient ascent

- Regularized maximum conditional likelihood estimate

$$\mathbf{w}^* = \arg \max_{\mathbf{w}} \ln \left[ \prod_j P(y^j | \mathbf{x}^j, \mathbf{w}) \right] - \frac{\lambda}{2} \sum_{i>0} w_i^2$$

$$w_i^{(t+1)} \leftarrow w_i^{(t)} + \eta \left\{ -\lambda w_i^{(t)} + \sum_j x_i^j [y^j - \hat{P}(Y^j = 1 | \mathbf{x}^j, \mathbf{w}^{(t)})] \right\}$$

reg derivative  $\leftarrow$  move towards zero

©Emily Fox 2014

5

## Stochastic Gradient Ascent for Logistic Regression

- Logistic loss as a stochastic function:

$$E_{\mathbf{x}} [\ell(\mathbf{w}, \mathbf{x})] = E_{\mathbf{x}} \left[ \ln P(y | \mathbf{x}, \mathbf{w}) - \frac{\lambda}{2} \|\mathbf{w}\|_2^2 \right]$$

- Batch gradient ascent updates:

$$w_i^{(t+1)} \leftarrow w_i^{(t)} + \eta \left\{ -\lambda w_i^{(t)} + \frac{1}{N} \sum_{j=1}^N x_i^{(j)} [y^{(j)} - P(Y = 1 | \mathbf{x}^{(j)}, \mathbf{w}^{(t)})] \right\}$$

$O(Nd)$

problem for large or streaming

- Stochastic gradient ascent updates:

- Online setting:

$$w_i^{(t+1)} \leftarrow w_i^{(t)} + \eta_t \left\{ -\lambda w_i^{(t)} + x_i^{(t)} [y^{(t)} - P(Y = 1 | \mathbf{x}^{(t)}, \mathbf{w}^{(t)})] \right\}$$

1 data point at a time

©Emily Fox 2014

6

# AdaGrad in Euclidean Space

$$x_t = (0 \ 0 \ 0 \ 1 \ 0 \ 0 \ 0 \ 1 \ 0 \ 0 \ 0 \ 0)$$

- For  $\mathcal{W} = \mathbb{R}^d$ ,

no constraints on  $w$

$$w^{(t+1)} \leftarrow w^{(t)} - \eta \text{diag}(A_t)^{-1} g_t$$

- For each feature dimension,

$$w_i^{(t+1)} \leftarrow w_i^{(t)} - \eta_{t,i} g_{t,i}$$

where

$$\eta_{t,i} = \eta / A_{t,ii}$$

- That is,

$$w_i^{(t+1)} \leftarrow w_i^{(t)} - \frac{\eta}{\sqrt{\sum_{\tau=1}^t g_{\tau,i}^2}} g_{t,i}$$

- Each feature dimension has its own learning rate!

- Adapts with  $t$
- Takes geometry of the past observations into account
- Primary role of  $\eta$  is determining rate the first time a feature is encountered

©Emily Fox 2014

7

# Problem 1: Complexity of Update Rules for Logistic Regression

- Logistic regression update:

stoch. grad. asc.

$$w_i^{(t+1)} \leftarrow w_i^{(t)} + \eta_t \left\{ -\lambda w_i^{(t)} + x_i^{(t)} [y^{(t)} - P(Y = 1 | \mathbf{x}^{(t)}, \mathbf{w}^{(t)})] \right\}$$

- Complexity of updates:

- Constant in number of data points ✓
- In number of features?  $O(d)$ 
  - Problem both in terms of computational complexity and sample complexity

what if we have 1B features??

- What can we do with very high dimensional feature spaces?

- Kernels not always appropriate, or scalable ← "kernel trick"
- What else?

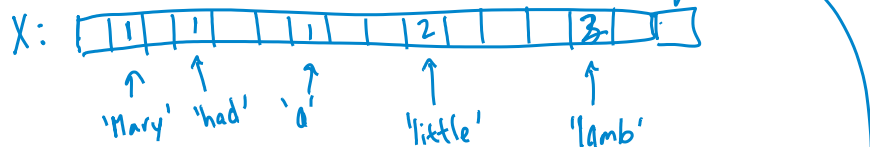
©Emily Fox 2014

8

## Problem 2: Unknown Number of Features

- For example, bag-of-words features for text data:

- “Mary had a little lamb, little lamb...”



- What's the dimensionality of  $x$ ? *size of vocabulary ... millions*
- What if we see new word that was not in our vocabulary?
  - Obamacare
  - Theoretically, just keep going in your learning, and initialize  $w_{\text{Obamacare}} = 0$
  - In practice, need to re-allocate memory, fix indices, ... A big problem for Big Data

©Emily Fox 2014

9

## What Next?

- Hashing & Sketching!
  - Addresses both dimensionality issues and new features in one approach!

- Let's start with a much simpler problem: Is a string in our vocabulary?

- Membership query

- How do we keep track?

- Explicit list of strings

- Very slow

*scan the list*

{ 'Mary', 'had', 'a', 'little', 'lamb', 'Obamacare' }

- Fancy Trees and Tries

- Hard to implement and maintain

- Hash tables?

$h(\text{'Mary'}) = 7$

$h(\text{'Obamacare'}) = 7$

→ { 'Mary', 'Obamacare' }

©Emily Fox 2014

10

# Hash Functions and Hash Tables

- Hash functions map **keys** to integers (bins):

$$h: X \rightarrow \{1, \dots, m\}$$

- Keys can be integers, strings, objects,...



- Simple example: **mod**

- $h(i) = (a \cdot i + b) \% m$

$$a=7 \quad b=11 \quad m=32$$

$$i=4 \Rightarrow h(i) = 39 \% 32 = 7$$

- Random choice of (a,b) (usually primes)  $\rightarrow$  random hash fcn
- If inputs are uniform, bins are uniformly used
- From two results can recover (a,b), so not pairwise independent  $\rightarrow$  Typically use fancier hash functions  $\{h(i), h(j)\}$

- Hash table:

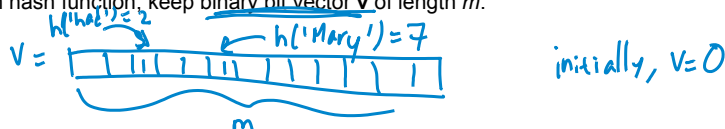
- Store list of objects in each bin
- Exact, but storage still linear in size of object ids, which can be very long
  - E.g., hashing very long strings, entire documents

*need this for our theory*

# Hash Bit-Vector Table-based Membership Query

- Approximate queries with one-sided error: Accept false positives only
  - If we say no, element is not in set
  - If we say yes, element is very to be likely in set

- Given hash function, keep binary bit vector **v** of length **m**:



- Query Q(i): Element i in set?

- $v(h(i)) = 0 \Rightarrow Q(i) = \text{no!}$
- $v(h(i)) = 1 \Rightarrow Q(i) = \text{probably yes}$

- Collisions:

$$h(\text{'Obamacare'}) = 7 \Rightarrow v(h(\text{'Obamacare'})) = v(7) = 1$$

*but Obamacare not in set*

- Guarantee: One-sided errors, but may make many mistakes

- How can we improve probability of correct answer?

*probability of collision =  $\frac{1}{m}$*   
*bc we saw 'Mary'*

# Bloom Filter: Multiple Hash Tables

- Single hash table -> Many false positives
- Multiple hash tables with independent hash functions
  - Apply  $h_1(i), \dots, h_p(i)$ , set all bits to 1



- Query  $Q(i)$ ?
  - if  $\forall j, h_j(i) = 1$   
 $Q(i) = \text{very probably yes}$
  - else  $Q(i) = \text{no}$
- Significantly decrease probability of false positives

'Mary' + 'Obamacare'  
collide in  $h_1$  but not  $h_2$

©Emily Fox 2014

13

# Analysis of Bloom Filter

- Want to keep track of  $n$  elements with false positive probability of  $\delta > 0 \dots$  how large  $m$  &  $p$ ?
  - $p$  ← # of hashes

- Simple analysis yields:
  - dim of each

$$m = \frac{n \log_2 \frac{1}{\delta}}{\ln 2} \approx 1.5n \log_2 \frac{1}{\delta}$$

$$p = \log_2 \frac{1}{\delta}$$

prob. of mistakes exp. decreasing w/ # hash fens

versus  $\frac{1}{m}$  by making hash table longer

©Emily Fox 2014

14

# Sketching Counts

want sketch to keep track of  $w^{(t)}$

- Bloom Filter is super cool, but not what we need...
  - We don't just care about whether a feature existed before, but to keep track of counts of occurrences of features! (assuming  $x_i$  integer)

- Recall the LR update:  $(1 - \eta_t \lambda) w_i^{(t)} + x_i^{(t)} \eta_t (y^{(t)} - P(\cdot))$

$$w_i^{(t+1)} \leftarrow w_i^{(t)} + \eta_t \left\{ -\lambda w_i^{(t)} + x_i^{(t)} [y^{(t)} - P(Y = 1 | \mathbf{x}^{(t)}, \mathbf{w}^{(t)})] \right\}$$

- Must keep track of (weighted) counts of each feature:
  - E.g., with sparse data, for each non-zero dimension  $i$  in  $\mathbf{x}^{(t)}$ :

For all entries of hash  
- multiply current  $w^{(t)}$  by  $(1 - \eta_t \lambda)$

For all  $x_i^{(t)} \neq 0$   
-  $w_i^{(t+1)} \neq x_i^{(t)} \cdot \text{const}$  ←  $\eta_t (y^{(t)} - P(\cdot))$

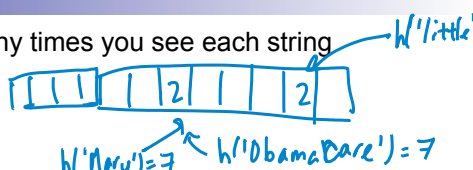
- Can we generalize the Bloom Filter?

©Emily Fox 2014

15

# Count-Min Sketch: single vector

- Simpler problem: Count how many times you see each string
- Single hash function:  $h$ 
  - Keep Count vector of length  $m$
  - every time see string  $i$ :



$$\text{Count}[h(i)] \leftarrow \text{Count}[h(i)] + 1$$

See 'Mary'  $\Rightarrow \text{Count}[7] = 2$       $Q('Mary') = \text{Count}[7] = 2 > 1$   
 'Obamacare'  $\Rightarrow \text{Count}[7] = 2$

- Again, collisions could be a problem:

■  $a_i$  is the count of element  $i$ : true counts

$$\text{Count}[j] = \sum_{i: h(i)=j} a_i$$

$Q(i) \rightarrow \text{return } \hat{a}_i = \text{Count}[h(i)] \geq a_i$  over-est. true count of 'Mary'

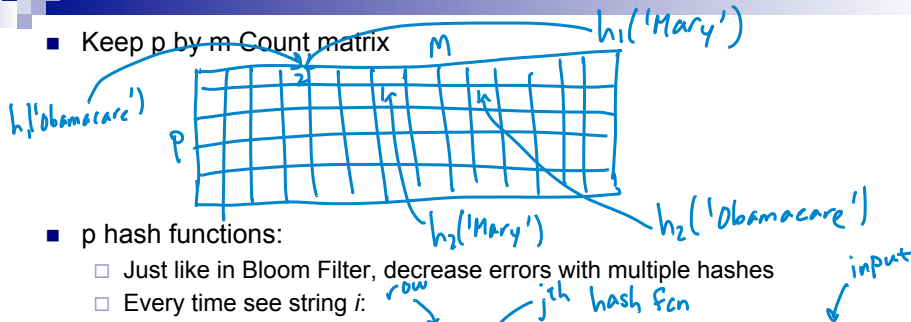
©Emily Fox 2014

16



# Count-Min Sketch: general case

- Keep  $p$  by  $m$  Count matrix



- $p$  hash functions:

- Just like in Bloom Filter, decrease errors with multiple hashes
- Every time see string  $i$ :

$$\forall j \in \{1, \dots, p\} : \text{Count}[j, h_j(i)] \leftarrow \text{Count}[j, h_j(i)] + 1$$

©Emily Fox 2014

17

# Querying the Count-Min Sketch

$$\forall j \in \{1, \dots, p\} : \text{Count}[j, h_j(i)] \leftarrow \text{Count}[j, h_j(i)] + 1$$

- Query  $Q(i)$ ?

- What is in  $\text{Count}[j, k]$ ?

$$\text{Count}[j, k] = \sum_{i: h_j(i)=k} a_i$$

- Thus:

$$Q(i)$$

$$\text{each } \text{count}[j, h_j(i)] \geq a_i$$

- Return:

$$\hat{a}_i \triangleq \min_j \text{Count}[j, h_j(i)] \geq a_i$$

↑ tightest upper bound

©Emily Fox 2014

18

# Analysis of Count-Min Sketch

$$\hat{a}_i = \min_j \text{Count}[j, h(i)] \geq a_i$$

- Set:
 
$$m = \left\lceil \frac{e}{\epsilon} \right\rceil$$

↑ length of each hash

$$p = \left\lceil \ln \frac{1}{\delta} \right\rceil$$

↑ # of hashes

false pos. rate

- Then, after seeing n elements:

$$a_i \leq \hat{a}_i \leq a_i + \epsilon n$$

} high prob. statement

- With probability at least  $1 - \delta$

©Emily Fox 2014

19

# Proof of Count-Min for Point Query with Positive Counts: Part 1 – Expected Bound

- $I_{i,j,k}$  = indicator that i & k collide on hash j:

$$(i \neq k) \wedge (h_j(i) = h_j(k))$$

- Bounding expected value:

$$E[I_{i,j,k}] = P(h_j(i) = h_j(k)) = \frac{1}{m} \leq \frac{\epsilon}{e}$$

- $X_{i,j}$  = total colliding mass on estimate of count of i in hash j:

$$X_{i,j} = \sum_{k \neq i} I_{i,j,k} a_k$$

← add their counts

$$\text{Count}[j, h_j(i)] = a_i + X_{i,j}$$

↑ colliding counts

- Bounding colliding mass: sum over words  $\neq$  'Mary'

$$E[X_{i,j}] = \sum_{k \neq i} a_k E[I_{i,j,k}] \leq \frac{n\epsilon}{e}$$

≤ n

- Thus, estimate from each hash function is close in expectation

©Emily Fox 2014

20

## Proof of Count-Min for Point Query with Positive Counts: Part 2 – High Probability Bounds

- What we know:  $Count[j, h_j(i)] = a_i + X_{i,j}$   $E[X_{i,j}] \leq \frac{\epsilon}{e} n$

- Markov inequality: For  $z_1, \dots, z_k$  positive iid random variables

$$P(\forall z_i : z_i > \alpha E[z_i]) < \alpha^{-k}$$

$$P(z_i > a) < \frac{E[z_i]}{a}$$

- Applying to the Count-Min sketch:

$$\begin{aligned} P(\hat{a}_i > a_i + \epsilon n) &= P(\forall j, Count[j, h_j(i)] > a_i + \epsilon n) \\ &= P(\forall j, a_i + X_{i,j} > a_i + \epsilon n) \\ &\leq P(\forall j, X_{i,j} > \epsilon E[X_{i,j}]) < e^{-p} = \delta \end{aligned}$$

"α" = α E[z<sub>i</sub>]  
and use ind. of z<sub>i</sub>
small prob.

©Emily Fox 2014

21

## But updates may be positive or negative

$$w_i^{(t+1)} \leftarrow w_i^{(t)} + \eta_t \left\{ -\lambda w_i^{(t)} + x_i^{(t)} [y^{(t)} - P(Y=1 | \mathbf{x}^{(t)}, \mathbf{w}^{(t)})] \right\}$$

- Count-Min sketch for positive & negative case

- $a_i$  no longer necessarily positive

- Update the same: Observe change  $\Delta_i$  to element  $i$ :

$$\forall j \in \{1, \dots, p\} : Count[j, h_j(i)] \leftarrow Count[j, h_j(i)] + \Delta_i$$

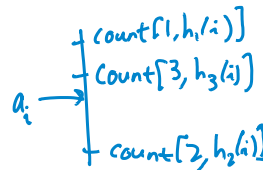
- Each  $Count[j, h_j(i)]$  no longer an upper bound on  $a_i$

- How do we make a prediction?

$$\hat{a}_i = \text{median } Count[j, h_j(i)]$$

- Bound:  $|\hat{a}_i - a_i| \leq 3\epsilon \|\mathbf{a}\|_1$

- With probability at least  $1 - \delta^{1/4}$ , where  $\|\mathbf{a}\|_1 = \sum_i |a_i|$



©Emily Fox 2014

22

# Finally, Sketching for LR

$$w_i^{(t+1)} \leftarrow w_i^{(t)} + \eta_t \left\{ -\lambda w_i^{(t)} + x_i^{(t)} [y^{(t)} - P(Y = 1 | \mathbf{x}^{(t)}, \mathbf{w}^{(t)})] \right\}$$

- Never need to know size of vocabulary!
- At every iteration, update Count-Min matrix:

$$\forall j, k \text{ Count}[j, k] = (1 - \eta_t \lambda) \text{Count}[j, k]$$

$$\forall x_i^{(t)} \neq 0$$

$$\forall j \text{ Count}[j, h_j(i)] += x_i^{(t)} \cdot \text{const} \quad \leftarrow \eta_t (y^{(t)} - P(Y=1|\dots))$$

- Making a prediction:

Remember our est. of  $w_i^{(t)}$ :  $\text{median}_j \text{Count}[j, h_j(i)]$

Make pred:  
 $-\log \text{odds} = w_0^{(t)} + \sum_{i: x_i \neq 0} \text{median}_j \text{Count}[j, h_j(i)] x_i^{(t)}$

- Scales to huge problems, great practical implications...

©Emily Fox 2014

23

# Hash Kernels

- Count-Min sketch not designed for negative updates
- Biased estimates of dot products
- Hash Kernels:** Very simple, but powerful idea to remove bias
- Pick 2 hash functions:

$h$ : Just like in Count-Min hashing

$$h: X \rightarrow \{1, \dots, m\}$$

$\xi$ : Sign hash function

$$\xi: X \rightarrow \{+1, -1\}$$

- Removes the bias found in Count-Min hashing (see homework)

- Define a "kernel", a projection  $\phi$  for  $\mathbf{x}$ :



$$\phi_i(\mathbf{x}) = \sum_{j: h(j)=i} f(j) X_j$$

can think of as random projection of  $\mathbf{x}$

IF  $X_j = 7$   
 $h(j) = 4$   
 $f(j) = -1$   
 add  $\downarrow$   $-7$   
 to bin 4

©Emily Fox 2014

24

# Hash Kernels Preserve Dot Products

$$\phi_i(\mathbf{x}) = \sum_{j:h(j)=i} \xi(j)\mathbf{x}_j$$

- Hash kernels provide unbiased estimate of dot-products!

$$E_{h,f}[\phi(x) \cdot \phi(y)] = x \cdot y \quad \text{pf: by homework}$$

- Variance decreases as  $O(1/m)$  ← gets better w/ more dims

- Choosing  $m$ ? For  $\epsilon > 0$ , if

$$m = O\left(\frac{\log \frac{N}{\delta}}{\epsilon^2}\right) \quad \text{log in data size}$$

- Under certain conditions...
- Then, with probability at least  $1-\delta$ :

$$(1 - \epsilon)\|\mathbf{x} - \mathbf{x}'\|_2^2 \leq \|\phi(\mathbf{x}) - \phi(\mathbf{x}')\|_2^2 \leq (1 + \epsilon)\|\mathbf{x} - \mathbf{x}'\|_2^2$$

©Emily Fox 2014

25