

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
April 7th, 2015

©Emily Fox 2015

1

Motivating AdaGrad ^{adaptive gradient} (Duchi, Hazan, Singer 2011)

- Assuming $\mathbf{w} \in \mathbb{R}^d$, standard stochastic (sub)gradient descent updates are of the form:

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

^{"step size"}
^{"learning rate"}

- Should all features share the same learning rate?

^{maybe instead:} $\eta_{t,i}$ specific to feature i

- Often have high-dimensional feature spaces
 - Many features are irrelevant \rightarrow ^{small learning rate}
 - Rare features are often very informative
- Adagrad provides a feature-specific adaptive learning rate by incorporating knowledge of the geometry of past observations

©Emily Fox 2015

2

AdaGrad Algorithm

- At time t , estimate optimal (sub)gradient modification A by

estimate of A at time t up to time t past gradients

$$A_t = \left(\sum_{\tau=1}^t g_\tau g_\tau^T \right)^{\frac{1}{2}}$$

in d dims, matrix $\sqrt{\cdot}$ is $O(d^3)$

- For d large, A_t is computationally intensive to compute. Instead,

$\text{diag}(A_t)$ $\Rightarrow A_t = \begin{pmatrix} A_{11} & 0 \\ 0 & A_{dd} \end{pmatrix}$ $A_{t,ii} = \sqrt{\sum_{\tau=1}^t g_{\tau,i}^2}$

- Then, algorithm is a simple modification of normal updates:

$$\mathbf{w}^{(t+1)} = \arg \min_{\mathbf{w} \in \mathcal{W}} \|\mathbf{w} - (\mathbf{w}^{(t)} - \eta \text{diag}(A_t)^{-1} g_t)\|^2_{\text{diag}(A_t)}$$

$\eta_t \rightarrow \eta A^{-1} \rightarrow \eta \text{diag}(A_t)^{-1}$

weigh dimensions by sqrt of sum of past grad. in that dim.

©Emily Fox 2015

AdaGrad Theoretical Guarantees

- AdaGrad regret bound:

$R_\infty := \max_t \|\mathbf{w}^{(t)} - \mathbf{w}^*\|_\infty$ radius of space

$R(T) = \sum_{t=1}^T \ell_t(\mathbf{w}^{(t)}) - \ell_t(\mathbf{w}^*) \leq 2R_\infty \sum_{i=1}^d \|g_{1:T,i}\|_2$ upper bound

Jensen's ineq. selected $\mathbf{w}^{(t)}$'s vs. best in retrospect

– In stochastic setting: $\ell(\mathbf{w}) = \mathbb{E}_x[\ell(\mathbf{w}, x)]$ and $\ell_t(\mathbf{w}) = \ell(\mathbf{w}, x^t)$ Then,

$$\frac{\mathbb{E}[R(T)]}{T} = \mathbb{E} \left[\ell \left(\frac{1}{T} \sum_{t=1}^T \mathbf{w}^{(t)} \right) \right] - \ell(\mathbf{w}^*) \leq \frac{2R_\infty}{T} \sum_{i=1}^d \mathbb{E}[\|g_{1:T,i}\|_2]$$

$\frac{1}{T} \mathbb{E}_{x^{1:T}} \left[\sum_{t=1}^T \ell(\mathbf{w}^*, x^t) \right]$

$= \frac{1}{T} \sum_{t=1}^T \mathbb{E}_{x^t} [\ell(\mathbf{w}^*, x^t)]$

$= \ell(\mathbf{w}^*)$

$\mathbb{E} \left[\frac{1}{T} \sum_{t=1}^T \ell(\mathbf{w}^{(t)}, x^t) \right] \geq \mathbb{E} \left[\ell \left(\frac{1}{T} \sum_{t=1}^T \mathbf{w}^{(t)}, x^t \right) \right]$

$\ell \left(\frac{x+x'}{2}, \frac{x+x'}{2} \right)$

- This really is used in practice!
- Many cool examples. Let's just examine one...

©Emily Fox 2015

4

AdaGrad Theoretical Example

- Expect to out-perform when gradient vectors are *sparse*

- SVM hinge loss example:

$$\ell_t(\mathbf{w}) = [1 - y^t \langle \mathbf{x}^t, \mathbf{w} \rangle]_+ \\ \mathbf{x}^t \in \{-1, 0, 1\}^d \text{ "z"}$$

- If $x_j^t \neq 0$ with probability $\propto j^{-\alpha}$, $\alpha > 1$

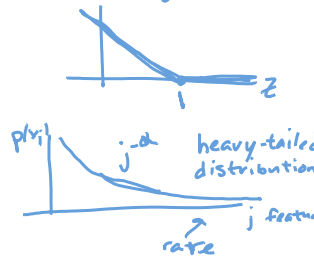
$$\mathbb{E} \left[\ell \left(\frac{1}{T} \sum_{t=1}^T \mathbf{w}^{(t)} \right) \right] - \ell(\mathbf{w}^*) = \mathcal{O} \left(\frac{\|\mathbf{w}^*\|_\infty}{\sqrt{T}} \cdot \max\{\log d, d^{1-\alpha/2}\} \right)$$

- Previously best known method: $\mathbb{E} \left[\ell \left(\frac{1}{T} \sum_{t=1}^T \mathbf{w}^{(t)} \right) \right] - \ell(\mathbf{w}^*) = \mathcal{O} \left(\frac{\|\mathbf{w}^*\|_\infty}{\sqrt{T}} \cdot \sqrt{d} \right)$

©Emily Fox 2015

5

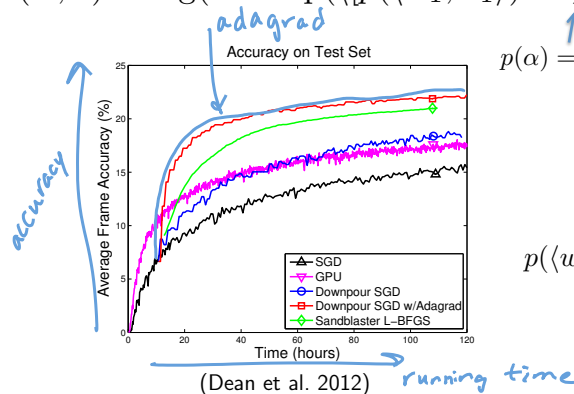
Hinge loss:
 $[1-z]_+$



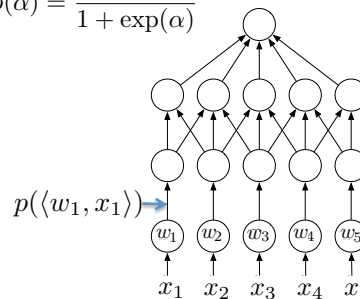
Neural Network Learning

- Very non-convex problem, but use SGD methods anyway

$$\ell(w, x) = \log(1 + \exp(\langle p(\langle w_1, x_1 \rangle) \cdots p(\langle w_k, x_k \rangle), x_0 \rangle))$$



$$p(\alpha) = \frac{1}{1 + \exp(\alpha)}$$



Distributed, $d = 1.7 \cdot 10^9$ parameters. SGD and AdaGrad use 80 machines (1000 cores), L-BFGS uses 800 (10000 cores)

©Emily Fox 2015

Images from Duchi et al. ISMP 2012 slides

6

What you should know about Logistic Regression (LR) and Click Prediction

- Click prediction problem:
 - Estimate probability of clicking
 - Can be modeled as logistic regression
- Logistic regression model: Linear model
- Gradient ascent to optimize conditional likelihood
- Overfitting + regularization
- Regularized optimization
 - Convergence rates and stopping criterion
- Stochastic gradient ascent for large/streaming data *focus*
 - Convergence rates of SGD
- AdaGrad motivation, derivation, and algorithm

©Emily Fox 2015

7

Problem 1: Complexity of LR Updates

- Logistic regression update:

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

stochastic gradient ascent

- Complexity of updates:

- Constant in number of data points ✓
- In number of features? *did*
 - Problem both in terms of computational complexity and sample complexity

what if we have 1B features?

- What can we with very high dimensional feature spaces?

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

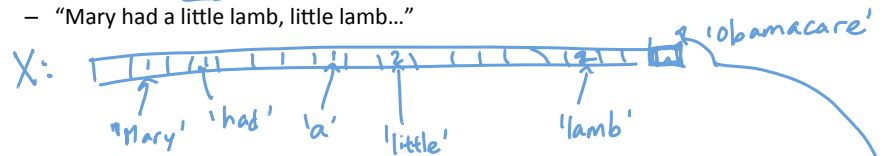
©Emily Fox 2015

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 2015

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

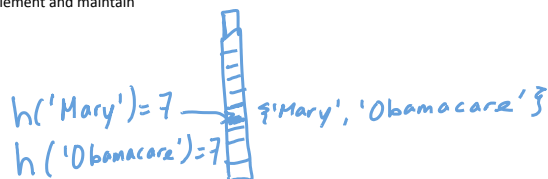
scan the list

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

- Fancy Trees and Tries

- Hard to implement and maintain

- Hash tables?



©Emily Fox 2015

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
- 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

©Emily Fox 2015

11

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 likely to be in set
- Given hash function, keep binary bit vector v of length m :
 - $h('had') = 7$
 - $h('Mary') = 7$
 - $v =$ [bit vector of length m]
 - initially, $v = 0$
- Query $Q(i)$: Element i in set?
 - $v(h(i)) = 0 \Rightarrow Q(i) = \text{no!}$
 - $v(h(i)) = 1 \Rightarrow Q(i) = \text{probably yes}$ (or if m small compared to vocab then "I have no idea")
- Collisions:
 - $h('Obamacare') = 7 \Rightarrow v(h('Obamacare')) = v(7) = 1$
 - but 'Obamacare' not in set
 - bc we saw 'Mary'
- Guarantee: One-sided errors, but may make many mistakes
 - How can we improve probability of correct answer?

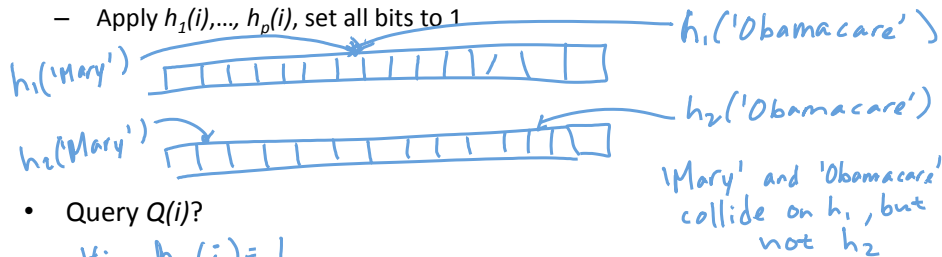
prob. of collision $= 1/m$

12

©Emily Fox 2015

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

©Emily Fox 2015

13

Analysis of Bloom Filter

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

dim of each \nwarrow \nearrow # of hashes

- Simple analysis yields:

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

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

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

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

©Emily Fox 2015

14

Sketching Counts

For LR:
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:

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

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

- 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 $w_i^{(t)}$ by $(1 - \eta_t \lambda)$

For all $x_i^{(t)} \neq 0$
- $w_i^{(t+1)} += x_i^{(t)} \cdot \text{const}$ $\nearrow \eta_t (y^{(t)} - P(Y=1|\dots))$

- Can we generalize the Bloom Filter?

©Emily Fox 2015

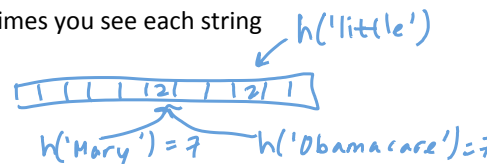
15

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' $\rightarrow \text{count}[7] = 2$
 'Obamacare' $\rightarrow \text{count}[7] = 2$

$Q('Mary')$
 $= \text{Count}[7]$
 $= 2 > 1$
 \uparrow
 true count
 of 'Mary'

- Again, collisions could be a problem:

- a_i is the count of element i :

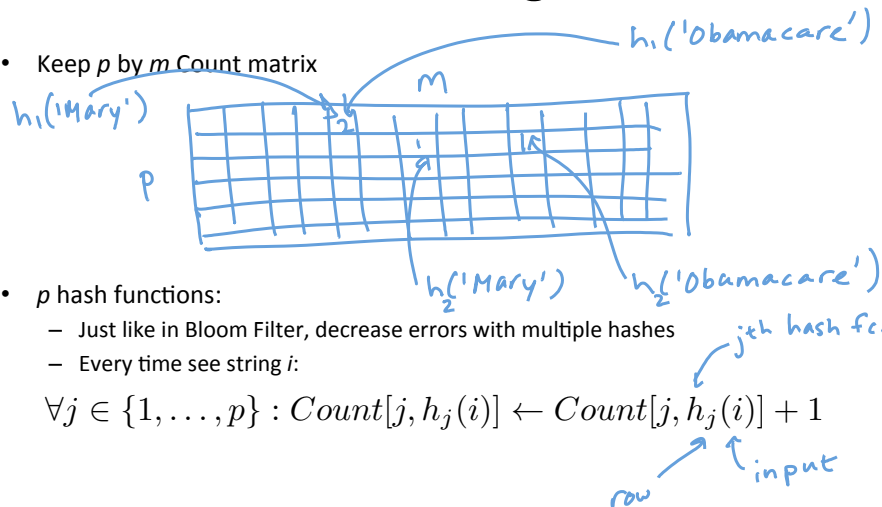
$$\text{Count}[j] = \sum_{i: h(i)=j} a_i$$

$$Q(i) \rightarrow \text{return } \hat{a}_i = \text{Count}[h(i)] \geq a_i$$

©Emily Fox 2015

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 2015

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:

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

\nwarrow tightest upper bound

©Emily Fox 2015

18

Analysis of Count-Min Sketch

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

- Set:

$$m = \left\lceil \frac{e}{\epsilon} \right\rceil \quad p = \left\lceil \ln \frac{1}{\delta} \right\rceil$$

length of each hash # 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 2015

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 :

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

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

- 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 if collide
sum over words \neq 'Mary'

- Bounding colliding mass:

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

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

'Mary' counts

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

©Emily Fox 2015

20

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

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

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

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

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

"a" = $\alpha E[z_i]$
and use ind. of z_i

- Applying to the Count-Min sketch:

$$\begin{aligned} P(\hat{a}_i > a_i + \epsilon n) &= P(\forall j, \text{Count}[j, h(i)] > a_i + \epsilon n) \\ &= P(\forall j, \cancel{a_i} + X_{i,j} > \cancel{a_i} + \epsilon n) \quad \leftarrow P = \left[\frac{1}{\delta} \right] \\ &\leq P(\forall j, X_{i,j} > \epsilon E[X_{i,j}]) < e^{-P} \leq \delta \end{aligned}$$

small prob. \square

©Emily Fox 2015

21