

Case Study 1: Estimating Click Probabilities

Database Issues

Tackling an Unknown Number of Features with Sketching

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

Sham Kakade

April 7th, 2015

©Sham Kakade 2016

1

Problem 1: Complexity of LR Updates

- Logistic regression update:

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

- Complexity of updates:

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

1B features

- What can we do with very high dimensional feature spaces?
 - Kernels not always appropriate, or scalable
 - What else?

©Sham Kakade 2016

2

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
 - Fancy Trees and Tries
 - Hard to implement and maintain
 - Hash tables?

{'mary', 'had', 'a', ... 'lamb'}

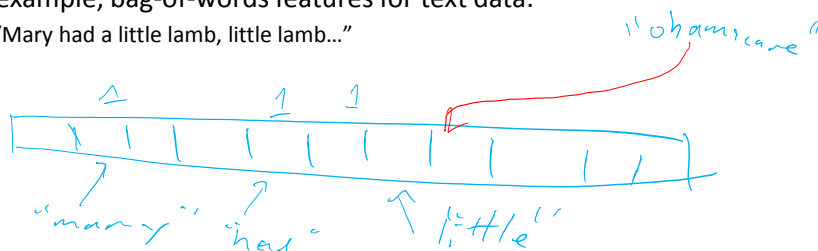
$h('mary')$
 $h('obamacare')$

©Sham Kakade 2016

3

Problem 2: Unknown Number of Features

- For example, bag-of-words features for text data:
 - “Mary had a little lamb, little lamb...”



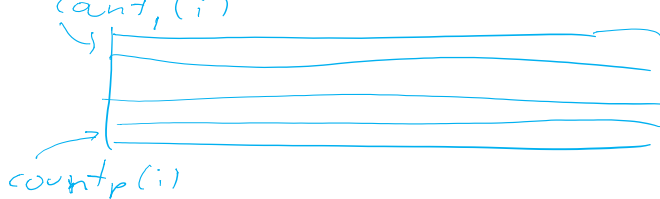
- What's the dimensionality of \mathbf{x} ? *← size of vocab.*
- What if we see new word that was not in our vocabulary?
 - Obamacare
 - Theoretically, just keep going in your learning, and initialize $\mathbf{w}_{\text{Obamacare}} = 0$
 - In practice, need to re-allocate memory, fix indices,... A big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

©Sham Kakade 2016

4

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 :

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

row index.
input

©Sham Kakade 2016

5

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 $\text{Count}[j, k]$?

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

- Thus:

highest upper bound.

- Return:

$$a_i \leq \min_j \text{count}(j, h_j(i)) \geq a_i$$

©Sham Kakade 2016

6

Analysis of Count-Min Sketch

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

- Set: $O\left(\frac{1}{\epsilon}\right)$

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

- Then, after seeing n elements:

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

} want +

- With probability at least $1-\delta$

©Sham Kakade 2016

7

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 :

$$O(\epsilon)$$

- Bounding expected value:

$$E[I_{i,j,k}] = \Pr(h_j(i) = h_j(k)) = \frac{1}{m}$$

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

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

- Bounding colliding mass:

$$E[X_{i,j}] \leq O(n\epsilon)$$

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

©Sham Kakade 2016

8

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

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

Handwritten notes:
 \downarrow (2 ^{positive} r.v.)
 $\leq E[Z]$
 $\frac{1}{\alpha}$

- Applying to the Count-Min sketch:

union bound

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 Δ_i to element i :

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

- Each $\text{Count}[j, h_j(i)]$ no longer an upper bound on a_i
- How do we make a prediction?

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

- Bound: $|\hat{a}_i - a_i| \leq 3\epsilon \|\mathbf{a}\|_1$
 - With probability at least $1 - \delta^{1/4}$, where $\|\mathbf{a}\| = \sum_i |a_i|$

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:
- Making a prediction:
- Scales to huge problems, great practical implications...

©Sham Kakade 2016

11

Hash Kernels

- Count-Min sketch not designed for negative updates

- Biased estimates of dot products

$$x \cdot x' \quad \mathbb{E} E[\phi(x) \cdot \phi(x')]$$

- **Hash Kernels:** Very simple, but powerful idea to remove bias

- Pick 2 hash functions:

- h : Just like in Count-Min hashing

$$h: \text{string} \rightarrow \{1, \dots, m\}$$

- ξ : Sign hash function

$$\xi: \text{string} \rightarrow \{-1, 1\}$$

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

- Define a “kernel”, a projection ϕ for \mathbf{x} :

$$\phi = [\quad]$$

$$\phi_i = \sum_{j: h(j)=i} \xi(j) a_j$$

add $\xi(j) \cdot \underbrace{a_j}_{\text{see } j \text{ } c_j \text{ times}}$ to index $h(j)$

©Sham Kakade 2016

12

Hash Kernels Preserve Dot Products

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

a - counts

- Hash kernels provide unbiased estimate of dot-products!

$$\mathbb{E} [\phi(\mathbf{x}) \cdot \phi(\mathbf{y})] = \mathbf{x} \cdot \mathbf{y}$$

r.f. : HW!

- Variance decreases as $O(1/m)$

same as

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

$$m = O\left(\frac{\log \frac{N}{\delta}}{\epsilon^2}\right)$$

Random Proj. arguments!

J.L.

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

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

©Sham Kakade 2016

13

Learning With Hash Kernels

- Given hash kernel of dimension m , specified by h and ξ
 - Learn m dimensional weight vector
- Observe data point \mathbf{x}
 - Dimension does not need to be specified a priori!
- Compute $\phi(\mathbf{x})$:
 - Initialize $\phi(\mathbf{x})$
 - For non-zero entries j of \mathbf{x} :

$$P(Y=1|\phi(\mathbf{x}), \mathbf{w}^t) = \frac{e^{\phi(\mathbf{x}) \cdot \mathbf{w}^t}}{1 + e^{\phi(\mathbf{x}) \cdot \mathbf{w}^t}}$$

- Use normal update as if observation were $\phi(\mathbf{x})$, e.g., for LR using SGD:

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

©Sham Kakade 2016

14

Interesting Application of Hash Kernels: Multi-Task Learning

- Personalized click estimation for many users:

- One global click prediction vector \mathbf{w} :

$$\frac{e^{\mathbf{w} \cdot \mathbf{x}}}{1 + e^{\mathbf{w} \cdot \mathbf{x}}}$$

- But...

- A click prediction vector \mathbf{w}_u per user u :

instead

$$\frac{e^{\mathbf{w}_u \cdot \mathbf{x}}}{1 + e^{\mathbf{w}_u \cdot \mathbf{x}}}$$

- But...

- Multi-task learning: Simultaneously solve multiple learning related problems:

- Use information from one learning problem to inform the others

- In our simple example, learn both a global \mathbf{w} and one \mathbf{w}_u per user:

- Prediction for user u :

$$(\mathbf{w} + \mathbf{w}_u) \cdot \mathbf{x} = \mathbf{w} \cdot \mathbf{x} + \mathbf{w}_u \cdot \mathbf{x}$$

- If we know little about user u :

$$\mathbf{w} \cdot \mathbf{x}$$

- After a lot of data from user u :

$$\mathbf{w}_u \neq 0 \quad \text{get user specific info.}$$

$$\frac{e^{(\mathbf{w} + \mathbf{w}_u) \cdot \mathbf{x}}}{1 + e^{(\mathbf{w} + \mathbf{w}_u) \cdot \mathbf{x}}}$$

©Sham Kakade 2016

15

Problems with Simple Multi-Task Learning

- Dealing with new user is annoying, just like dealing with new words in vocabulary
- Dimensionality of joint parameter space is HUGE, e.g. personalized email spam classification from Weinberger et al.:
 - 3.2M emails
 - 40M unique tokens in vocabulary
 - 430K users
 - 16T parameters needed for personalized classification!

©Sham Kakade 2016

16

Hash Kernels for Multi-Task Learning

- Simple, pretty solution with hash kernels:
 - Very multi-task learning as (sparse) learning problem with (huge) joint data point z for point x and user u : $d + (\# \text{ users}) \times d$

$$z_{x,u} = (\underbrace{x_1, \dots, x_d}_{\text{global}}, \underbrace{0, 0, \dots, 0}_d, \dots, \underbrace{x_1, \dots, x_d, 0, \dots, 0}_{\text{user } u})$$

- Estimating click probability as desired:

$$w = (\bar{w}, \bar{w}_1, \dots, \bar{w}_u, \dots, \bar{w}_{\# \text{ users}})$$

\hookrightarrow "weight vect. for z " $z_{x,u} \cdot w = w \cdot x + w_u \cdot x$

- Address huge dimensionality, new words, and new users using hash kernels:

©Sham Kakade 2016

17

Simple Trick for Forming Projection $\phi(x,u)$

don't form $z_{x,u}$

- Observe data point x for user u
 - Dimension does not need to be specified a priori and user can be new!

- Compute $\phi(x,u)$:

- Initialize $\phi(x,u) = 0$
- For non-zero entries j of x_j :

- E.g., $j = \text{'Obamacare'}$
- Need two contributions to ϕ :
 - Global contribution
 - Personalized Contribution

- Simply:

hash bucket

$$\phi = g(\text{'Obamacare'}) \cdot x_j + g(\text{'Obamacare'} \rightarrow i) \cdot x_j$$

user \rightarrow

- Learn as usual using $\phi(x,u)$ instead of $\phi(x)$ in update function

©Sham Kakade 2016

18

Results from Weinberger et al. on Spam Classification: Effect of m

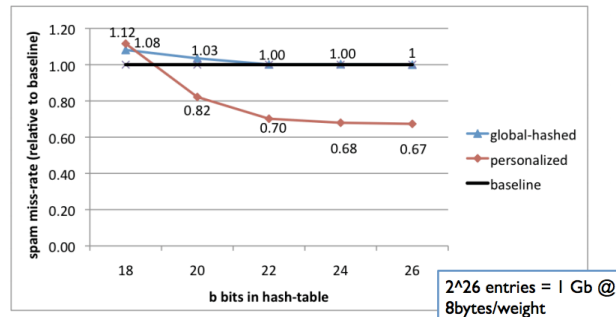


Figure 2. The decrease of uncaught spam over the baseline classifier averaged over all users. The classification threshold was chosen to keep the not-spam misclassification fixed at 1%. The hashed global classifier (*global-hashed*) converges relatively soon, showing that the distortion error ϵ_d vanishes.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

©Sham Kakade 2016

19

Results from Weinberger et al. on Spam Classification: Multi-Task Effect

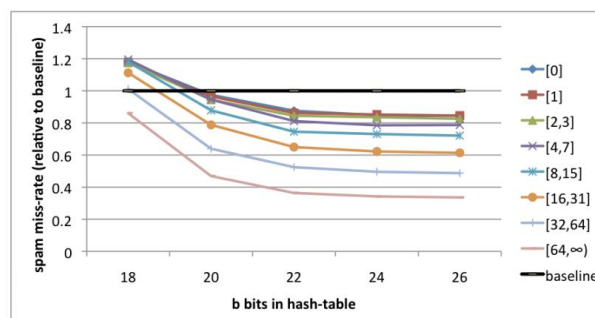


Figure 3. Results for users clustered by training emails. For example, the bucket [8, 15] consists of all users with eight to fifteen training emails. Although users in buckets with large amounts of training data do benefit more from the personalized classifier (up to 65% reduction in spam), even users that did not contribute to the training corpus at all obtain almost 20% spam-reduction.

©Sham Kakade 2016

20

What you need to know

- Hash functions
- Bloom filter
 - Test membership with some false positives, but very small number of bits per element
- Count-Min sketch
 - Positive counts: upper bound with nice rates of convergence
 - General case
- Application to logistic regression
- Hash kernels:
 - Sparse representation for feature vectors
 - Very simple, use two hash function (Can use one hash function...take least significant bit to define ξ)
 - Quickly generate projection $\phi(\mathbf{x})$
 - Learn in projected space
- Multi-task learning:
 - Solve many related learning problems simultaneously
 - Very easy to implement with hash kernels
 - Significantly improve accuracy in some problems (if there is enough data from individual users)

database

learning

next