Sketching Count Statistics for Approximate Bayesian Inference Diana Cai and Zoe Ashwood

Modern data analysis problems

- Massive in size, possible infinite data stream
- High dimensional
- Complex models are needed for modeling

Probabilistic modeling and Bayesian inference

- Powerful for identifying interpretable, latent structure in data
- Provides a framework for making predictions about future observations
- Incorporate prior knowledge, share statistical strength across model

This project: we explore approximating the sufficient statistics in probabilistic models using a hashing-based probabilistic data structure (count-min sketch) perform Bayesian inference in two ubiquitous probabilistic models

[Cormode and Muthukrishnan, 2003]

Data stream $x_1, ..., x_M$ of tokens in $\{1,...,K\}$

Data stream $x_1, ..., x_M$ of tokens in $\{1,...,K\}$

N hash functions randomly generated from a *pairwise-independent* hash family

 $h_n: [K] \rightarrow [J], J \le K, n=1,...,N$

Data stream $x_1, ..., x_M$ of tokens in $\{1,...,K\}$

For a random hash h in the family, Pr(h(x)=h(y)) at most 1/J

N hash functions randomly generated from a pairwise-independent hash family

 $h_n: [K] \rightarrow [J], J \le K, n=1,...,N$

Data stream $x_1, ..., x_M$ of tokens in $\{1,...,K\}$

For a random hash h in the family, Pr(h(x)=h(y)) at most 1/J

N hash functions randomly generated from a pairwise-independent hash family

 $h_n: [K] \rightarrow [J], J \le K, n=1,...,N$

Data structure: N x J counter array: C[n,j]

h ₁			
h ₂			
h ₃			

Data stream $x_1, ..., x_M$ of tokens in $\{1,...,K\}$

For a random hash h in the family, Pr(h(x)=h(y)) at most 1/J

N hash functions randomly generated from a pairwise-independent hash family

 $h_n: [K] \rightarrow [J], J \le K, n=1,...,N$

Data structure: N x J counter array: C[n,j]

Update(C, k): hash k and update the counts

 $C[n,h_n(k)] += 1$, for all n=1,...,N

h ₁			
h ₂			
h ₃			

Data stream $x_1, ..., x_M$ of tokens in $\{1,...,K\}$

For a random hash h in the family, Pr(h(x)=h(y)) at most 1/J

N hash functions randomly generated from a pairwise-independent hash family

 $h_n: [K] \rightarrow [J], J \le K, n=1,...,N$

Data structure: N x J counter array: C[n,j]

Update(C, k): hash k and update the counts

 $C[n,h_n(k)] += 1$, for all n=1,...,N

h ₁	+1		
h ₂			
h ₃			

 $h_1(k) = 1$

Data stream $x_1, ..., x_M$ of tokens in $\{1,...,K\}$

For a random hash h in the family, Pr(h(x)=h(y)) at most 1/J

N hash functions randomly generated from a pairwise-independent hash family

 $h_n: [K] \rightarrow [J], J \le K, n=1,...,N$

Data structure: N x J counter array: C[n,j]

Update(C, k): hash k and update the counts

 $C[n,h_n(k)] += 1$, for all n=1,...,N

h ₁	+1			
h ₂			+1	
h ₃				

 $h_1(k) = 1$ $h_2(k) = 4$

Data stream $x_1, ..., x_M$ of tokens in $\{1,...,K\}$

For a random hash h in the family, Pr(h(x)=h(y)) at most 1/J

N hash functions randomly generated from a pairwise-independent hash family

 $h_n: [K] \rightarrow [J], J \le K, n=1,...,N$

Data structure: N x J counter array: C[n,j]

Update(C, k): hash k and update the counts

 $C[n,h_n(k)] += 1$, for all n=1,...,N

h_1	+1			
h ₂			+1	
h ₃		+1		

$$h_1(k) = 1$$

 $h_2(k) = 4$
 $h_3(k) = 2$

Data stream $x_1, ..., x_M$ of tokens in $\{1,...,K\}$

For a random hash h in the family, Pr(h(x)=h(y)) at most 1/J

N hash functions randomly generated from a pairwise-independent hash family

 $h_n: [K] \rightarrow [J], J \le K, n=1,...,N$

Data structure: N x J counter array: C[n,j]

Update(C, k): hash k and update the counts

 $C[n,h_n(k)] += 1$, for all n=1,...,N

Query(C, k): returns estimate for the count of k

For all n, return the minimum count C[n,h_n(k)]

h ₁	+1			
h ₂			+1	
h ₃		+1		

 $h_1(k) = 1$ $h_2(k) = 4$ $h_3(k) = 2$

[Cormode and Muthukrishnan, 2003]

For a random hash h in the family,

Data stream $x_1, ..., x_M$ of tokens in $\{1,...,K\}$

N hash functions randomly generated from a *pairwise-independent* hash family

 $h_n: [K] \rightarrow [J], J \le K, n=1,...,N$

Data structure: N x J counter array: C[n,j]

Update(C, k): hash k and update the counts

 $C[n,h_n(k)] += 1$, for all n=1,...,N

Query(C, k): returns estimate for the count of k

For all n, return the minimum count C[n,h_n(k)]

sketch after M observations:

Pr(h(x)=h(y)) at most 1/J



 $h_1(k) = 1$ $h_2(k) = 4$ $h_3(k) = 2$

[Cormode and Muthukrishnan, 2003]

For a random hash h in the family,

Data stream $x_1, ..., x_M$ of tokens in $\{1,...,K\}$

N hash functions randomly generated from a *pairwise-independent* hash family

$$h_n: [K] \rightarrow [J], J \le K, n=1,...,N$$

Data structure: N x J counter array: C[n,j]

Update(C, k): hash k and update the counts

 $C[n,h_n(k)] += 1$, for all n=1,...,N

Query(C, k): returns estimate for the count of k

For all n, return the minimum count C[n,h_n(k)]

sketch after M observations:

Pr(h(x)=h(y)) at most 1/J



[Cormode and Muthukrishnan, 2003]

For a random hash h in the family,

Data stream $x_1, ..., x_M$ of tokens in $\{1,...,K\}$

N hash functions randomly generated from a *pairwise-independent* hash family

$$h_n: [K] \rightarrow [J], J \le K, n=1,...,N$$

Data structure: N x J counter array: C[n,j]

Update(C, k): hash k and update the counts

 $C[n,h_n(k)] += 1$, for all n=1,...,N

Query(C, k): returns estimate for the count of k

For all n, return the minimum count C[n,h_n(k)]

sketch after M observations:

Pr(h(x)=h(y)) at most 1/J



 $h_1(k) = 1$ $h_2(k) = 4$ $h_3(k) = 2$

[Cormode and Muthukrishnan, 2003]

1. The <u>estimated</u> count ≥ <u>true</u> count

observed count = true count + counts of colliding tokens

1. The <u>estimated</u> count ≥ <u>true</u> count

observed count = true count + counts of colliding tokens

 \Rightarrow estimated count := minimum observed count \ge true count



1. The <u>estimated</u> count ≥ <u>true</u> count

observed count = true count + counts of colliding toker

N: # of hash functions J: range of hash function M: total # of tokens

 \Rightarrow estimated count := minimum observed count \ge true count

1. With probability 1 - exp(-N),

estimated count ≤ true count + (e/J) * M

1. The <u>estimated</u> count ≥ <u>true</u> count

```
observed count = true count + counts of colliding toker
```

N: # of hash functions J: range of hash function M: total # of tokens

 \Rightarrow estimated count := minimum observed count \ge true count

1. With probability 1 - exp(-N),

estimated count ≤ true count + (e/J) * M [Homework 1, Q2!]

Proof sketch:

 X_{kn} := counts tokens colliding with k on h_n = observed - true $E[X_{kn}] \le (1/J) * M$ (expected colliding mass) $Pr(estimate > true + (e/J) * M) = Pr(\forall n, \{X_{kn} > e * E[X_{kn}]\}) \le exp(-N)$

Dirichlet-multinomial (DM) model

<u>Multinomial occupancy scheme</u>: K bins, M balls thrown independently in each bin Probability θ_k of landing into the kth bin

Let $x_1,...,x_M \sim$ **Multinomial(\theta)**. (M tokens in [K])

Under the DM model, we assume that the probabilities are random: $\theta \sim \text{Dirichlet}_{\kappa}(\alpha)$

Goal is to compute the posterior distribution using Bayes' rule: $\theta \in \Delta^{K}$

 $p(\theta \mid x_1,...,x_M) = p(\theta) p(x_1,...,x_M \mid \theta) = \text{Dirichlet}(\theta \mid \alpha + c), \text{ where } c = (c_1,...,c_k)$

Predictive distribution:
$$p(x_{M+1} = k | x_1,...,x_M) = \int p(x_{M+1}) p(\theta | x_1,...,x_M) d\theta$$

Sketching multinomial counts

For the DM model, we just have to keep track of the posterior through the sufficient statistic vector c, i.e., the counts of the tokens.

In a streaming setting, we can directly apply the CMsketch for inference:

for t = 1, 2, ...

- 1. Observe x_t
- 2. Update(C, x_t): (hash x_t and update all the counters)
- 3. Evaluate predictive distribution

 $p(x_{M+1} = k \mid x_1, ..., x_M) = (\alpha + \hat{c}_k) (\sum_j \alpha_j + \hat{c}_j)^{-1}$

Generated 10K size dataset from the Dirichlet-multinomial model



Applied the CMsketch (N hash functions with range J) to approximate the counts

M=10,000; K=1000 N hash functions, J range

 $epsilon(k) := (\hat{c}_{k} - c_{k}) / M$

plotted mean over all k

1 hash function: errors are much worse, more variance



M=10,000; K=10,000 N hash functions, J range

 $epsilon(k) := (\hat{c}_{k} - c_{k}) / M$

plotted mean over all k



M=10,000; K=1000; N hash functions, J range

 $MAE(\hat{c},c) = sum(\hat{c}_k - c_k) / K$

For 20 trials: (mean over trials, sd computed over trials)



M=10,000; K=10,000; N hash functions, J range

 $MAE(\hat{c},c) = sum(\hat{c}_{k}-c_{k}) / K$

For 20 trials: (mean over trials, sd computed over trials)



M=10,000; K=10,000; N hash functions, J=1000

Posterior approximation: 1-step ahead predictive probabilities (streaming)

M=10,000; K=10,000; N hash functions

Posterior approximation: 1-step ahead predictive probabilities (streaming)

M=10,000; K=10,000; N hash functions, J range

Posterior approximation: MAP estimate (the mode of posterior distribution)

Topic modeling with Latent Dirichlet Allocation (LDA)

- Document corpus: D documents, over vocabulary set [N]
- Topic: probability distribution over words in the vocabulary set
- Goal: uncover K topics in corpus

Topic modeling with Latent Dirichlet Allocation (LDA)

- K topics, D documents (each of length n words)
- Each topic is a distribution over the words in the vocabulary set [N]

Draw topic $\beta_k \sim \text{Dirichlet}_N(.)$, k = 1,...,K, $(\beta_k \in \Delta^N)$

• Each document contains some proportion of the topics

Draw $\theta_d \sim \text{Dirichlet}_{K}(.), \quad d = 1, ..., D, \quad (\theta_d \in \Delta^K)$

• For the jth word in document d, assign that word to a topic

Draw $z_{dj} \sim \text{Multinomial}(\theta_d), \quad j = 1, ..., n, \quad (z_{dj} \in [K])$

• Each word in a document is then drawn from that topic

Draw $w_{dj} \sim \text{Multinomial}(\beta_i)$, where $i = z_{dj}$, j = 1,...n, $(w_{dj} \in [N])$

Gibbs sampling with approximate count statistics

• We want to know the posterior:

$$p(z|w) = \frac{p(z,w)}{\sum_{z} p(z,w)}$$

- We can write analytic form using model assumptions
- But cannot evaluate on any real corpus: sum of Kⁿ terms in denominator is intractable (n is number of words in entire corpus)
- Must use an approximation method: can use Gibbs sampling
- Iterative algorithm that is guaranteed to converge to the true posterior

Gibbs sampling with approximate count statistics

- T iterations; each iteration involves sampling a new topic assignment for each word in the corpus
- The distribution that the topic assignment is sampled from is constantly being updated
- 5 data structures involved: one is n_{k, w}, a matrix in R^{KxN}, containing the number of times word w appears in topic k. Problem: too large for large vocabularies, N.
- Core idea: can we approximate $n_{k,w}$ with K count-min sketches? How is the posterior distribution, $p(z|w) = \frac{p(z,w)}{\sum_{z} p(z,w)}$ affected? Can we still recover the topics in a corpus?

Topic-word probabilities and sketching

• With probability $1 - \frac{1}{e^{H}}$, counts are in range:

• Now: we want to estimate topic-word probabilities, $\beta_w^{k, ext{true}}$

$$\equiv \frac{n_{k,w}^{\mathrm{true}} + \eta}{\sum_w n_{k,w}^{\mathrm{true}} + N\eta}$$

• With same high probability, these are in range:

$$\frac{1}{(1+\frac{N}{J})}(\beta_w^{k,\mathrm{true}}+\frac{1}{J}) \leq \beta_w^{k,\mathrm{sketch}} \leq \beta_w^{k,\mathrm{true}}$$

 $n_{k,w}^{ ext{true}} \leq n_{k,w}^{ ext{sketch}} \leq n_{k,w}^{ ext{true}} + rac{1}{J} \sum n_{k,w}^{ ext{true}}$

Topic-word probabilities and sketching

LDA Toy Corpus: 100,000 words; N = 2500

LDA Toy Corpus: 100,000 words; N = 2500

Results: Conditional Loglikelihood of Corpus

(Larger values are better)

Conditional loglikelihood of corpus, log(P(w|z))

Iteration

Recall Topics:

Results: Recovered Topics - without sketch

Results: Recovered Topics - with sketch

(memory use $\frac{4}{5}$ of vocabulary: H = 2, J = 1000)

Conclusions and Future Directions

- Demonstrated incorporation of count-min sketch into inference algorithms for two ubiquitous probabilistic models.
- Useful when number of unique tokens (DM case) or vocabulary size (LDA) is too large to be stored in memory
- Were able to recover posterior probabilities in both cases:
 - $\circ~$ DM-case: stream of M tokens. Were able to recover probability, θ_k , that k^{th} unique token appears in stream
 - LDA: were able to recover latent topics in toy corpus
- In future:
 - Applying LDA sketch-based inference on different document corpuses: corpuses with denser topics and on real-world corpuses
 - Investigating additional sketches: we are interested in proportions, not counts. Are there sketches that are better at preserving proportions?
 - Applying LDA sketch-based inference in streaming setting

Thanks!

Dirichlet-multinomial (DM) model

6

<u>Multinomial occupancy scheme</u>: K bins, M balls thrown independently in each bin Probability θ_k of landing into the kth bin

Let $x_1,...,x_M$ be M tokens in [K]. Likelihood has a **Categorical(\theta)** distribution:

 $p(x_1,...,x_M \mid \theta) = \prod_n \prod_k \theta_k^{1(x_m=k)} = \prod_k \theta_k^{c_k}$, where $c_k := #$ of balls in k^{th} bin

Under the DM model, we assume that the probabilities are random: $\theta \sim \text{Dirichlet}_{\kappa}(\alpha)$ Goal is to compute the posterior distribution using Bayes' rule: $\theta \in \Delta^{\kappa}$

$$p(\theta \mid x_1,...,x_M) = p(\theta) p(x_1,...,x_M \mid \theta) = Dirichlet(\theta \mid \alpha + c), \text{ where } c = c_1,...,c_k)$$

Predictive distribution:
$$p(x_{M+1} = k | x_1,...,x_M) = \int p(x_{M+1}) p(\theta | x_1,...,x_M) d\theta$$

