

000
001
002
003
004
005
006
007
008
009
010
011
012
013
014
015
016
017
018
019
020
021
022
023
024
025
026
027
028
029
030
031
032
033
034
035
036
037
038
039
040
041
042
043
044
045
046
047
048
049
050
051
052
053

Analysis on Complementary Count Min Sketch

Neeraj Sharma, Kanchi Masalia
Department of Computer Science Engineering
University of Massachusetts Amherst
{neerajsharma, kmasalia}@umass.edu

Abstract

Count Sketch and Count Min Sketch are commonly used for frequency estimation in a streaming setting. This data structures can also be used for approximating data by accumulating it. One interesting application of this data structure is observed in field of Machine Learning for feature selection. In this paper, we present analysis of our proposed data structure called Complementary Count Min Sketch, which is aimed to use less space and running time without much loss in accuracy compared to count sketch. We analyzed it for feature selection as well as frequency estimation wherein the elements can be removed as well from the stream. We analyzed the performance over RCV1 - a high dimensional dataset as well as a synthetic dataset which gave us flexibility to define dimensions of the dataset, nature of dataset and tweak it's sparsity.

1 Introduction

There has been an immense rise in generated data which has led to the era of Big Data. These have heavily influence how we think, build, and maintain applications. For streams with data arriving at high rate, algorithms are needed which use as little processing time and space in order to analyze and provide query response in real time. Count Min Sketch and Count Sketch data structures are used in such scenarios for frequency estimation and finding most frequently occurring items of the stream.

In machine learning and statistics, feature selection is a process in which one chooses features which contribute most to the prediction variable or output. At times feature selection is confused with dimensionality reduction. It is true that both of these help in reducing the features in a dataset, but the difference lies in how they approach this problem. Dimensionality reduction reduces the number of features by creating new features as combinations of existing ones. So all the features are still present in a way, but the total number of features is reduced. But in feature selection, we either retain a feature or remove it completely from the dataset. When data is present in low dimensions, there are many different algorithms available for feature selection. However, when data is present in high dimensions, the training time for the model increases with the dimensions exponentially. Feature selection of high dimensional data requires large amount of memory and time. Usually high dimensional vectors are sparse i.e very few features actually have non-zero values. It is not difficult to load such sparse high dimensional data in memory as we can ignore the zero elements and use data structure like dictionary to have the position of feature and its value. The problem arises with storing and performing operations on dense high dimensional vector. We propose a structures for feature selection of sparse high dimensional vectors using Complementary Count Min Sketches

054 (CCMS) along with maintaining heap of the most important features to preserve the interpretability
055 of features. Previously [1] have used Count Sketch (CS) for feature selection. In this work, we are
056 analyzing the complementary count min sketch.
057

058 2 Related Work

059 Count Sketch was introduced by [2] to find the most frequently occurring item in the data stream.
060 Count Min Sketch was introduced by [3] to find the approximate count of items occurring in the
061 stream of data. Both the data structures have a similar design. However, they have different error
062 guarantees. Count Sketch uses pairwise independent hash functions and sign hash functions for
063 hashing the features into the sketch (Array like data structure). Count Min Sketch uses only hash
064 functions. So the total number of hash functions computed while adding an element to count sketch
065 is twice as compared when the element is added to count min sketch. The estimate of frequency
066 of element provided by count min sketch is an upper bound of the actual frequency of the element
067 whereas the estimate of frequency provided by count sketch could be lower or higher than the actual
068 occurrence of the element. The error estimate for Count Min Sketch is L1 norm of the frequency
069 vector (approximation sketch) and for Count Median Sketch is L2 norm of frequency vector. By
070 Cauchy-Schwarz inequality we know that L1 norm is bounded by $\sqrt{n} * L2norm$ where n is the
071 length of the frequency vector (approximation sketch). However, Count-Min sketch algorithm gives
072 better average error than the Count Sketches when using constant space. Comparative analysis
073 of count sketch vs count min sketch motivated us to go ahead and propose variants of sketch for
074 frequency estimation and its other applications.
075

076 We found an interesting application of this sketch for feature selection technique wherein the sketch
077 is used for compression of feature weights.
078

079 Feature selection for high dimensional data is very important as the training time increases expo-
080 nentially with the dimensions. Due to curse of dimensionality, high dimensional data can easily
081 overfit regression model and thus requires careful hyperparameter tuning. One of the solution to this
082 problem is feature hashing [4] which makes working with high dimensional data computationally
083 feasible, but at the cost of losing the interpretability of features. Consider a 3-gram string “abc”.
084 With feature hashing, one uses a lossy, random hash function $h : strings \rightarrow \{1, 2, \dots, R\}$ to map
085 “abc” to a feature number $h(abc)$ in the range $\{1, 2, \dots, R\}$. This is extremely convenient because
086 it enables one to avoid creating a large look-up dictionary. Furthermore, this serves as a dimen-
087 sionality reduction technique. Unfortunately, this convenience comes at a cost, we lose the identity
088 of the original features. This is not a viable option if one cares about both feature selection and
089 interpretability.
090

091 Another popular approach by [5] is to use greedy thresholding methods combined with stochastic
092 gradient descent to prevent the feature vector from becoming too dense and blowing up in memory.
093 In these methods, the intermediate iterates are regularized at each step, and a full gradient update
094 is never stored nor computed (since this is memory and computation intensive). However, it is well
095 known that greedy thresholding can be myopic and can result in poor convergence.
096

097 [6] have introduced a new sub-linear space sketch: the Weight-Median Sketch. This is for learning
098 compressed linear classifiers over data streams while supporting the efficient recovery of large-
099 magnitude weights in the model. This enables memory-limited execution of several statistical anal-
100 yses over streams, including online feature selection, streaming data explanation.
101

102 [2] has developed a data structure to capture the features that are most discriminative of one stream
103 (or class). The Weight-Median Sketch is built on top of the data structure Count-Sketch, but, instead
104 of sketching counts, it captures sketched gradient updates to the model parameters. The core idea
105 for performing feature selection is figuring out the most discriminative features from the set of
106 features. This memory efficient data structure can be used to accumulate the gradients of the high
107 dimensional feature vector when the model is learning without much loss in approximation. This
can be viewed as dimensionality reduction via random projection. The issue with this approach is

108 we loose the interpretability of the features. Once we accumulate the gradients in count sketch for
109 high dimensional data, it is not possible to decipher the discriminating feature as different features
110 would hash at same index in the sketch and will loose the interpretability.

111 [1] have implemented a method to maintain the interpretability of the features using the Weight-
112 Median sketch and maintaining a heap of most discriminating features. This method accurately and
113 efficiently performs feature selection on real-world, large-scale datasets with billions of dimensions.
114

115 We will be leveraging this methodology to maintain the interpretability of features. We propose
116 using our variants of sketch to maintain weights which would use less space compared to Weight
117 median sketch built on count sketch.
118

119 3 Methodology

120 Here we are presenting the methods and algorithms which we used for our study.

121 3.1 Zipfian Distribution and Power Law Distribution

122 Many real data distributions such as sizes of cities, word frequencies, citations of papers, web page
123 access frequencies, and file transfer size and duration are often characterized by the Zipfian, Pareto,
124 or Power-law distributions which only differ by the choice of parameters. The zipfian distribution
125 with parameter $\alpha > 0$ is a discrete distribution stating that the k^{th} largest frequency f_k has a fre-
126 quency proportional to $k^{-\alpha}$. We can see, $\alpha = 0$ generates a uniform distribution whereas the larger
127 α the more skewed the distribution gets. Zipfian law is cumulative form of Power law distribu-
128 tion. Zipfian and Power law distributions are especially interesting with regards to the heavy hitters
129 problem since this problem looks for frequencies which are significantly larger than the rest of the
130 data. For increasingly skewed zipfian distributions the elements with such frequencies become more
131 frequent, since only a few of the overall frequencies account for more of the total frequency.
132

133 3.2 Count Min Sketch

134 The Count Min Sketch (CMS) is a randomized method closely related to bloom filters. The count
135 min sketch data structure can only be used when Δ is positive. Count Min Sketch has d random
136 pairwise independent hash functions $h_j, j \in \{1, 2, \dots, d\}$ to map the vector's components to bins w .
137 $h_j : \{1, 2, \dots, p\} \rightarrow \{1, 2, 3, \dots, w\}$. Every component i is hashed into bin $S(j, h_j(i))$. The count-
138 min-sketch supports two operations: **UPDATE**(item i , increment Δ). The update operation updates
139 the sketch with any observed increment. For an increment Δ to an item i , the sketch is updated by
140 adding Δ to the cell $S(j, h_j(i)) \forall j \in \{1, 2, \dots, d\}$. The **QUERY** operation returns the estimate for
141 component i , the min of all the d different associated counters.
142

143 **Algorithm 1:** Count Min Sketch

144 v universal hash functions h_j

145 Initialize count-min-sketch matrix $S \in R^{v,w} = 0$

146 **Update**(item i , increment: Δ)

147 Update component i with update Δ

148 $S(j, h_j(i)) = S(j, h_j(i)) + \Delta \quad \forall j \in \{1 \dots d\}$

149 **Query** (item i):

150 Query Sketch for an estimate for item i

151 **return** $\text{Min}(S_{j, h_j(i)})$

152 for any item i , when we query its value from Count Min Sketch, it will always give an overestimate
153 of its value as other item also can map to the same position in count min sketch.
154

162
163
164
165
166
167
168
169
170
171
172
173
174
175
176
177
178
179
180
181
182
183
184
185
186
187
188
189
190
191
192
193
194
195
196
197
198
199
200
201
202
203
204
205
206
207
208
209
210
211
212
213
214
215

$$S_j[h_j(i)] = f_i(x) + \sum_{y \neq i: h_j(y) = h_j(i)} f(y)$$

Where $f_i(x)$ is actual value of item i and $\sum_{y \neq i: h_j(y) = h_j(i)} f(y)$ is error in estimating value for item i . Expected Error:

$$E\left[\sum_{y \neq i: h_j(y) = h_j(i)} f(y)\right] = \sum_{y \neq i} Pr(h_j(i) = h_j(y))f(y)$$

$$\implies \sum_{y \neq i} \frac{1}{w} f(y) \leq \frac{N}{w}$$

where N is all items in stream. Using Markov inequality we can say:

$$f(i) \leq S_j[h_j(i)] \leq f(i) + \frac{\epsilon * N}{K}$$

where K is Heavy Hitters Count and ϵ is error rate. So in this way we can say that count min sketch can be used to approximate frequency of every item i in a stream up to error $\frac{\epsilon n}{k}$ with probability $\geq 1 - \delta$ in $O(\log(\frac{1}{\delta}) \frac{k}{\epsilon})$ space and time $O(\log \delta^{-1})$ where δ is failure probability.

For frequency estimation of heavy hitters under stream settings, we can use $\Delta = 1$ in algorithm 1.

3.3 Count Median Sketch (Count Sketch)

The algorithm uses d random hash functions similar to count-min sketch but it also uses d random sign functions as well to map the components of vectors randomly to $\{+1, -1\}$ i.e. $s_i : \{1, 2, \dots, n\} \rightarrow \{+1, -1\}$. The count-sketch (CS) also similarly supports two operations, **UPDATE**(item i , update δ) and **QUERY**(item i). The **UPDATE** operation updates the sketch with any observed update. It may be increment as well as decrement. For an update Δ to an item i , the sketch is updated by adding $s_j(i)\Delta$ to the cell $S(j, h_j(i)) \forall j \in \{1, 2, \dots, d\}$. The **QUERY** operation returns an estimate for component i , the median of all the d different associated counters.

Algorithm 2: Count Median Sketch

v universal hash functions h_j
 v random sign functions s_j
 Initialize count-sketch tensor $S \in R^{v,w} = 0$

Update(item i , update: Δ)

Update component i with update Δ
 $S(j, h_j(i)) = S(j, h_j(i)) + s_j(i)\Delta \quad \forall j \in \{1 \dots d\}$

Query(item i):

Query Sketch for an estimate for item i
return Median($S_{j, h_j(i)} s_j(i)$)

Let us analyze the Count Median Sketch: Let us assume the possible different items coming into the stream will be M . Let us define an indicator function

$$Y_j = \begin{cases} 1, & \text{if } h(j) = h(j^*) \\ 0, & \text{otherwise} \end{cases}$$

Approximate value of item j :

$$\hat{f}_{j^*} = s(j) * S[h(j)]$$

216
217
218
219
220
221
222
223
224
225
226
227
228
229
230
231
232
233
234
235
236
237
238
239
240
241
242
243
244
245
246
247
248
249
250
251
252
253
254
255
256
257
258
259
260
261
262
263
264
265
266
267
268
269

$$\begin{aligned} \hat{f}_{j^*} &= s(j^*) \sum_{j=1}^M f_j s(j) Y_j \\ \implies & s(j^*)^2 f_j y_j^* + \sum_{j \neq j^*} f_j s(j^*) s(j) Y_j \\ \implies & f_j + \sum_{j \neq j^*} f_j s(j^*) s(j) Y_j \end{aligned}$$

Since $s(j^*)^2 = 1$ and $Y_{j^*} = 1$, By Linearity of Expectation:

$$E(\hat{f}_{j^*}) = f_{j^*} + \sum_{j \neq j^*} f_j E[s(j^*) s(j) Y_j]$$

Since s is pairwise independent and is independent of Y_j which is solely function of h , for $j \neq j^*$, we have:

$$\begin{aligned} E[s(j^*) s(j) Y_j] &= 0 \\ \implies E[\hat{f}_{j^*}] &= f_{j^*} \end{aligned}$$

As \hat{f}_{j^*} is an unbiased estimator of f_{j^*} , let us compute its variance:

$$\begin{aligned} \text{Var}(\hat{f}_{j^*}) &= E[\hat{f}_{j^*} - f_{j^*}]^2 \\ \implies E[\sum_{i \neq j^*} \sum_{j \neq j^*} f_i f_j s(i) s(j) Y_i Y_j] \end{aligned}$$

Since $s(j^*)^2 = 1$

$$\implies \sum_{i \neq j^*} \sum_{j \neq j^*} f_i f_j E[s(i) s(j) Y_i Y_j]$$

As for $i \neq j$, since s is pairwise independent and independent of h , $E[s(i) s(j) Y_i Y_j] = 0$. Therefore the only terms in the variance that survive are when $i = j$.

$$\implies \text{Var}(\hat{f}_{j^*}) = \sum_{j \neq j^*} f_j^2 E[Y_j^2]$$

As $E[Y_j^2] = E[Y_j]$ and Y_j is indicator function so $E[Y_j] = \text{Pr}(h(j) = h(j^*)) = \frac{1}{w}$. Therefore,

$$\text{Var}(\hat{f}_{j^*}) = \sum_{j \neq j^*} \frac{f_j^2}{w} = \frac{\|f\|_2^2 - f_{j^*}^2}{w} = \frac{\|f_{-j^*}\|_2^2}{w}$$

By Chebyshev Inequality:

$$\text{Pr}(\|\hat{f}_{j^*} - f_{j^*}\| \geq \epsilon \|f_{-j^*}\|) \leq \frac{\text{Var}(\hat{f}_{j^*})}{\epsilon^2 \|f_{-j^*}\|_2^2}$$

The query has a running time proportional to the depth of the sketch. As for the update procedure of Count-Median Sketch, two invocations of a hash function, a multiplication, and an addition is required at each row. All these operations are constant and add up to $O(d)$ time. The median of the d estimates must be found, which can be done in linear $O(d)$ time, yielding a total of $O(d) = O(\ln \delta^{-1})$ running time for the point query overall.

3.4 Comparison of Count Min and Count Sketch

The Count-Median Sketch provides a better guarantee, since it guarantees that \hat{v}_i is within an additive factor of $\epsilon \|v\|_2$ of the true frequency v_i with probability $1 - \delta$. This guarantee is significantly stronger in most cases as $\|v\|_2 \leq \|v\|_1$. The cost of this guarantee is significantly larger, since the Count-Median Sketch requires $O(\log(\frac{1}{\delta}) \frac{k}{\epsilon^2})$ in comparison to count min sketch which takes $O(\log(\frac{1}{\delta}) \frac{k}{\epsilon})$ space to support updates and queries in the same time as Count-Min Sketch.

If we compare the precision of the two sketches, it is observed that for data with a near uniform distribution, the Count-Median Sketch provides smaller errors, whereas when the data becomes more and more skewed, the Count-Min Sketch provides the smallest error. This is not surprising since the Count-Median Sketch still has a relationship with the L2-norm, which increases when the data becomes more skewed implying that the error increases as well.

Under equal space constraint, The precision of the Count-Median Sketch is also expected to change drastically, due to the decrease in space usage. In fact it is expected that the decrease in space implies that the error guarantee now is bounded according to the L1-norm instead of the L2-norm according to [7]. So Count-Median Sketch can in fact be shown to provide an error guarantee according to the L1-norm, by changing the width to be equal to the width of a Count-Min Sketch.

Summing it up, we can say that the Count-Min Sketch and the Count-Median Sketch with width $w = O(\epsilon^{-1})$ is indeed comparable, and that comparing the sketches according to space, precision and running times gives very similar results, where different data distributions determines which sketch performs the best. The only notable difference is in the running time of the query algorithms where the Count-Min Sketch in general seems to be faster than the Count-Median Sketch. It is likely due to calculating the minimum compared to calculating the median is faster and easier operation. This is where the motivation came for us to come up with a data structure inspired from count min sketch for accumulation of gradients.

3.5 Complementary Count Min Sketch

Count min sketch doesn't support reducing frequency or removing an item in a stream of updates. As if we decrement the frequency the error bound will not hold which states that the frequency estimated by count min would be equal or would be an overestimate of the true frequency. However, Count Sketch supports reducing frequency or removing an item and also holds the error bounds as it gives the estimate as an median. As we discussed in 3.4 that count min sketch is faster in comparison to count sketch and with equal space bound count sketch and count min sketch both provides error within L1-norm bound. To support negative updates in count min sketch we proposed this novel data structure - Complementary Count Min Sketch (CCMS). In Complementary (positive and negative) Count Min Sketch, we will have two sketches S^{pos} and S^{neg} which will accumulate the positive and negative updates respectively. So for **UPDATE**(item i , update Δ), we will first check if Δ is positive or negative. If Δ is positive then will update it in S^{pos} , the same way we did in Count Min Sketch's update method. If Δ is negative then we will update the absolute value of Δ in S^{neg} . For **QUERY**(item i), we will query the item in S^{pos} and S^{neg} and will return the $min_j S^{pos}[j, h_j(i)] - min_j S^{neg}[j, h_j(i)]$.

Let us do the Mathematical Analysis of Complementary Count Min Sketch: From Count Min Sketch analysis we know: For Single Count Min Sketch

$$S_j[h_j(i)] = f_i(x) + \sum_{y \neq i: h_j(y) = h_j(i)} f(y)$$

Where $f_i(x)$ is actual value of item i and $\sum_{y \neq i: h_j(y) = h_j(i)} f(y)$ is error in estimating value for item i . Similarly:

$$S_j^{pos}[h_j(i)] = f_i(x) + \sum_{y \neq i: h_j(y) = h_j(i)} f(y)$$

324
325
326
327
328
329
330
331
332
333
334
335
336
337
338
339
340
341
342
343
344
345
346
347
348
349
350
351
352
353
354
355
356
357
358
359
360
361
362
363
364
365
366
367
368
369
370
371
372
373
374
375
376
377

Algorithm 3: Complementary Count Min Sketch

v universal hash functions h_j

Initialize positive count-min-sketch matrix $S^{pos} \in R^{v,w} = 0$, negative count-min-sketch matrix $S^{neg} \in R^{v,w} = 0$

Update (item i , update: Δ)

Update component i with update Δ

if $\Delta > 0$ **then**

 | $S^{pos}(j, h_j(i)) = S^{pos}(j, h_j(i)) + \Delta \quad \forall j \in \{1 \dots d\}$

else

 | $S^{neg}(j, h_j(i)) = S^{neg}(j, h_j(i)) + abs(\Delta) \quad \forall j \in \{1 \dots d\}$

end

Query (item i):

Query Sketch for an estimate for item i

return $Min(S_{j,h_j(i)}^{pos}) - Min(S_{j,h_j(i)}^{neg})$

Where $f_i(x)$ is actual value of item i and $\sum_{y \neq i: h_j(y)=h_j(i)} f(y)$ is error in estimating value for item i .

$$S_j^{neg}[h_j(i)] = f_i(x) + \sum_{y \neq i: h_j(y)=h_j(i)} f(y)$$

Where $f_i(x)$ is actual value of item i and $\sum_{y \neq i: h_j(y)=h_j(i)} f(y)$ is error in estimating value for item i from negative items of stream.

For Complementary Count Min Sketch:

$$E[g(i)] = E[f^{pos}(i) + \sum_{y \neq i: h_j(y)=h_j(i)} f^{pos}(y)] - E[f^{neg}(i) + \sum_{y' \neq i: h_j(y')=h_j(i)} f^{neg}(y')]$$

Here y and y' are hash functions for positive and negative sketch. $g(x)$ is true estimate of particular item.

$$\implies E[f^{pos}(i) - f^{neg}(i)] - E[\sum_{y \neq i: h_j(y)=h_j(i)} f^{pos}(y) - \sum_{y' \neq i: h_j(y')=h_j(i)} f^{neg}(y')]$$

As

$$\implies E[f^{pos}(i) - f^{neg}(i)] = E[f(i)] = f(i)$$

So Error from CCMS:

$$Error = E[\sum_{y \neq i: h_j(y)=h_j(i)} f^{pos}(y) - \sum_{y' \neq i: h_j(y')=h_j(i)} f^{neg}(y')]$$

$$Error = I(y = y') \cdot E[\sum_{y \neq i: h_j(y)=h_j(i)} f^{pos}(y) - f^{neg}(y)] + I(y \neq y') \cdot E[\sum_{y' \neq y \neq i: h_j(y')=h_j(i)} f^{pos}(y) - f^{neg}(y')]$$

$$\implies \frac{1}{m^2} E[\sum_{y \neq i: h_j(y)=h_j(i)} f^{pos}(y) - f^{neg}(y)] + \frac{1}{m(m-1)} E[\sum_{y' \neq y \neq i: h_j(y')=h_j(i)} f^{pos}(y) - f^{neg}(y')]$$

378 So upper bound for the error will be:
 379

$$381 \implies UpperBound = \frac{1}{m^2} E\left[\sum_{y \neq i: h_j(y) = h_j(i)} f^{pos}(y) \right] + \frac{1}{m(m-1)} E\left[\sum_{y \neq y' \neq i: h_j(y) = h_j(i)} f^{pos}(y') \right]$$

384
 385 So lower bound for the error will be:
 386

$$387 \implies LowerBound = \frac{1}{m^2} E\left[\sum_{y \neq i: h_j(y) = h_j(i)} f^{neg}(y) \right] + \frac{1}{m(m-1)} E\left[\sum_{y \neq y' \neq i: h_j(y) = h_j(i)} f^{neg}(y') \right]$$

392 We observe the error is bounded between error from negative count min sketch to positive count min
 393 sketch. Thus the estimated frequency can be lower as well as greater than true frequency. However,
 394 If there are proportionate amount of negative and positive updates these error terms will cancel out
 395 each other and would be close to the true estimate.
 396

397 3.5.1 Variants of Complementary Count Min Sketch

398 **Complementary Count Min Sketch with Different Hash Functions:**

399 In the previously proposed complementary count min sketch, we used same hash functions for positive
 400 as well as negative count min sketch meaning positive and negative updates will have same
 401 positions in sketches. Thus, if hash of two particular items are colliding with each other then they
 402 will collide in both the sketches. We tried to analyze if we have different hash functions for positive
 403 and negative count min sketches, would it help us to minimize the error. Algorithm is described here
 404 4.
 405

407 **Algorithm 4:** Complementary Count Min Sketch using different Hash Functions

408 v universal hash functions h_j^{pos}, h_j^{neg} for positive and negative count min sketches
 409 Initialize positive count-min-sketch matrix $S^{pos} \in R^{v,w} = 0$
 410 Initialize negative count-min-sketch matrix $S^{neg} \in R^{v,w} = 0$
 411

412 **Update (item i, update: Δ)**

413 Update component i with update Δ

414 **if** $\Delta > 0$ **then**

415 | $S^{pos}(j, h_j^{pos}(i)) = S^{pos}(j, h_j^{pos}(i)) + \Delta \quad \forall j \in \{1 \dots d\}$

416 **else**

417 | $S^{neg}(j, h_j^{neg}(i)) = S^{neg}(j, h_j^{neg}(i)) + abs(\Delta) \quad \forall j \in \{1 \dots d\}$

418 **end**

420 **Query (item i):**

421 Query Sketch for an estimate for item i

422 **return** $Min(S_{j, h_j^{pos}(i)}^{pos}) - Min(S_{j, h_j^{neg}(i)}^{neg})$
 423

425 **Complementary Count Min Sketch with Single Hash Function:**

426 We realized that the previous approach is quite similar to having a single count min sketch where in
 427 we use different hash functions for positive and negative updates. This is synonymous to using one
 428 count min sketch where we update the positive one as usual and for the negative we get the hash of
 429 negative number and update at that position. Algorithm is described here 5.
 430

431 **Conservative Complementary Count Min Sketch:**

As we studied in our last semester, count min sketch with conservative updates gives more accurate

432 **Algorithm 5:** Complementary Count Min Sketch using Single Count Min Sketch

433 v universal hash functions h_j
 434 Initialize count-min-sketch matrix $S \in R^{v,w} = 0$
 435
 436 **Update (item i, update: Δ)**
 437 Update component i with update Δ
 438 **if** $\Delta > 0$ **then**
 439 | $S(j, h_j(i)) = S(j, h_j(i)) + \Delta \quad \forall j \in \{1\dots d\}$
 440 **else**
 441 | $S(j, h_j(-i)) = S(j, h_j(-i)) + \text{abs}(\Delta) \quad \forall j \in \{1\dots d\}$
 442 **end**

443
 444 **Query (item i):**
 445 Query Sketch for an estimate for item i
 446 **return** $\text{Min}(S(j, h_j(i)) - \text{Min}(S(j, h_j(-i))) \quad \forall j \in \{1\dots d\}$

448
 449
 450 and lower bounds on error. In standard count min sketch, the update operation updates the sketch
 451 with any observed increment. For an increment Δ to an item i, the sketch is updated by adding Δ
 452 to the cell $S(j, h_j(i)) \forall j \in \{1, 2, \dots, d\}$. In this variant, instead of incrementing each counter, we first
 453 compute $M = \min_{j \in \{1, \dots, d\}} S(j, h_j(i))$. Then we only increment $S(j, h_j(i))$ if $S(j, h_j(i)) = M$.
 454 We follow this technique for both positive and negative count min sketch. Query method remains
 455 same to the standard count min sketch method. Algorithm is described here 6.

457 **Algorithm 6:** Conservative Complementary Count Min Sketch

458 v universal hash functions h_j
 459 Initialize positive count-min-sketch matrix $S^{pos} \in R^{v,w} = 0$
 460 Initialize negative count-min-sketch matrix $S^{neg} \in R^{v,w} = 0$
 461

462 **Update (item i, update: Δ)**
 463 Update component i with update Δ
 464 Compute $M = \min_{j \in \{1, \dots, d\}} S(j, h_j(i))$
 465 **if** $\Delta > 0$ **then**
 466 | $S^{pos}(j, h_j(i)) = S^{pos}(j, h_j(i)) + \Delta \quad \forall j \in \{1\dots d\} \text{ s.t. } S(j, h_j(i)) = M$
 467 **else**
 468 | $S^{neg}(j, h_j(i)) = S^{neg}(j, h_j(i)) + \text{abs}(\Delta) \quad \forall j \in \{1\dots d\} \text{ s.t. } S(j, h_j(i)) = M$
 469 **end**

470
 471 **Query (item i):**
 472 Query Sketch for an estimate for item i
 473 **return** $\text{Min}(S_{j, h_j(i)}^{pos}) - \text{Min}(S_{j, h_j(i)}^{neg})$

474

475

476

477

478

479

480

481

482

483

484

485

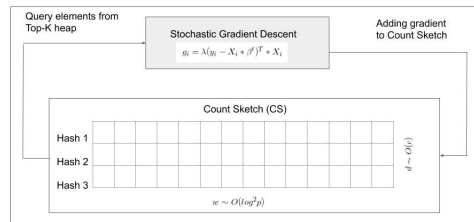
3.6 Feature selection techniques

3.6.1 Feature selection Using Sketches

Consider the feature selection problem in the high dimensional setting where we are given a dataset (X_i, y_i) for $i \in [n]$. Each data point $X_i \in R^p$ and label $y_i \in R$. We are interested in finding the k-sparse feature vector $\beta \in R^p$ from below optimization problem which solves our feature selection task where k non zero elements are the selected features.

$$\min_{\|\beta\|_0=k} \|y - X\beta\|_2$$

486 Where $X = [X_1; X_2; \dots; X_n]$ and $y = [y_1; y_2; \dots; y_n]$ denote the data matrix and label vector and l_o
487 norm $\|\beta\|_0$ denotes the number of non zero entries in $\|\beta\|_0$.
488 We are interested in solving the feature selection problem for high-dimensional datasets where the
489 number of features p is so large that a dense vector (or matrix) of size p cannot be stored explicitly in
490 memory. Sketch data structures allow us to accumulate the gradients updates over several iterations
491 because of linear aggregation. We will be using the same algorithm as described in [1]. Our novel
492 contribution is to use above described sketch variants instead of Count Sketch. First, we initialize
493 the Sketch S and the feature vector $\beta^{t=0}$ with zeros entries. The sketch hashes a p -dimensional
494 vector into $O(\log^2 p)$ buckets as fig(1). At iteration t , this algorithm selects a random row X_i from
495 the data matrix X and computes the stochastic gradient update term using the learning rate λ . For
496 logistic regression, the gradient of softmax function can be defined as $g_i = \lambda * (y_i - X_i \beta^t) X_i$.
497 As the data vector X_i and the corresponding stochastic gradient term are sparse, we only add the
498 non-zero entries of the stochastic gradient term $\{g_{ij} : \forall j, g_{ij} > 0\}$ to the Count-Sketch S . After
499 adding the non-zero entry to the sketch, we perform query operation for those items and insert it
500 into the heap structure - top k if the absolute value from query operation output is greater than
501 the minimum absolute value in heap along with feature position. In gist, we are maintaining the
502 most discriminating features (ones with high absolute weights) in Top-K. This structure is solely
503 responsible for maintaining the interpretability of the features. After performing this stochastic
504 gradient update step for all the examples and performing the described operation. Finally, we query
505 the Top-K values of the sketch as the final output. This is detailed in Algorithm 7
506



507
508
509
510
511
512
513
514
515
516
517 Figure 1: Schematic of Algorithm

520 **Algorithm 7:** Feature Selection: Using Sketch

521 **Result:** The top-k heavy-hitters from the Sketch

522 **Initialize:** $\beta^0 = 0$, S (Sketch), λ (learning rate)

523 **while not stopping criteria do**

524 Find the gradient update $g_i = \lambda(y_i - X_i \beta^t)^T X_i$

525 Add the gradient update to the sketch $g_i \rightarrow S$

526 Get the top-k heavy hitters from the sketch $\beta^{t+1} \leftarrow S$

527 **end**

528
529
530 **3.6.2 Greedy Thresholding**

531
532 In the feature selection algorithms, the class of hard thresholding algorithms have the smallest mem-
533 ory footprint. Hard thresholding algorithms retain only the top-k values and indices of the entire
534 feature vector using $O(k \log(p))$ memory. An algorithm where, after each gradient update, a hard
535 threshold is applied to the features. Only the top-k features are kept active, while the rest are set to
536 zero. This algorithm generates the following iterates for the i^{th} variable in an stochastic gradient
537 descent (SGD) framework.
538
539

$$\beta^{t+1} \leftarrow H_k(\beta^t - 2\lambda(y_i - X_i \beta^t)^T X_i)$$

540 The sparsity of the feature vector B^t , enforced by the hard thresholding operator H^k , alleviates
 541 the need to store a vector of size $O(p)$ in the memory in order to keep track of the changes of
 542 the features over the iterates. As in this algorithm, we only retains the top-k elements of β , the
 543 hard thresholding procedure greedily discards the information of the non top-k coordinates from the
 544 previous iteration. In particular, it clips off coordinates that might add to the support set in later
 545 iterations. This drastically affects the performance of hard thresholding algorithms, especially in
 546 real-world scenarios where by the design matrix X is not random, normalized, or well-conditioned.
 547 In this scenario, the gradient terms corresponding to the true support typically arrive in lagging order
 548 and are prematurely clipped in early iterations by H^k .
 549

550 3.6.3 Logistic Regression with Heap 551

552 In this approach we maintain a heap similar to the the heap as we maintained for feature selection
 553 using sketch technique. After each iteration of SGD, we update the heap in order to maintain the
 554 top-k features with high magnitude where k is user defined. This ensures that after each iteration,
 555 we just have weights for k features.
 556

557 3.7 Top K Recovery 558

559 Our aim for feature selection of high dimensional data set is to find most discriminating features i.e
 560 features which have most impact on the true label. In order to do this, we maintain a heap kind of
 561 data structure which maintains fixed size of the high magnitude feature weights and their indexes.
 562 This structure is updated after each update on count sketch i.e when the weight of any feature is
 563 modified. At the end of training, the Top-K heap is used to recover the K-sparse weight vector. The
 564 key idea of recovery lies in that a suitably high dimensional sparse signal can be inferred from very
 565 few linear observations.
 566

567 4 Dataset 568

569 We have used two datasets in our experiment.
 570
 571

572 4.1 RCV1 573

574 Reuters Corpus Volume I (RCV1) contains over 800,000 manually categorized newswire stories.
 575 This dataset contains the Non-zero values cosine-normalized, log TF-IDF values for each document.
 576 All the features are real and between 0 and 1. RCV1 [8] dataset is categorized into 4 groups to cap-
 577 ture the major subjects of a story. The 4 groups are Economics(ECAT), Corporate/Industrial(CCAT),
 578 Government/Social(GCAT) and Markets(MCAT). This data was processed and converted into binary
 579 classes where in positive class includes CCAT, ECAT and negative class includes GCAT, MCAT. ¹
 580

581 The statistics of these datasets are summarized in table 1.
 582

Dataset	Dimension	Train Size	Test Size
RCV1	47236	20242	677399

583 Table 1: Dataset Statistics
584
585

586 As the data set has large number of features, data is represented in dictionary format where the key
 587 is the position of feature and value is the feature value.
 588

589 For simplicity purpose we chose to implement binary classification using logistic loss. Our proposed
 590 methodologies can be easily extended for multiclass classification.
 591
 592

593 ¹<https://www.csie.ntu.edu.tw/~cjlin/libsvmtools/datasets/binary.htmlrcv1.binary>

594
595
596
597
598
599
600
601
602
603
604
605
606
607
608
609
610
611
612
613
614
615
616
617
618
619
620
621
622
623
624
625
626
627
628
629
630
631
632
633
634
635
636
637
638
639
640
641
642
643
644
645
646
647

4.2 Synthetic Dataset

To analyze the performance of Complementary Count Min Sketch for feature selection and heavy hitters, the results of RCV1 dataset were not easy to interpret due its high dimensions and lack of true model parameters. In order to mitigate this, we came up with the idea of generating a synthetic dataset where we can play with the number of examples, feature dimensions and sparsity of features. We also generated the true feature weights and calculated the true labels based on it. Knowing true model parameters was a major advantage of working with synthetic dataset. This helped us to analyze the behaviour of sketch and top-k heap at any stage easily. We created dataset using Gaussian distribution as well as with power law distribution with by inducing sparsity as per choice.

4.2.1 Synthetic Dataset using Gaussian Distribution

As we saw in section 3.4 that under uniform distribution count sketch has smaller error in estimation of frequency in comparison to count min sketch. We wanted to verify this for our data structure. We created the dataset matrix using standard Gaussian distribution. We randomly selected the important features and generated random weights and based on these data points and feature weight assigned labels to each example. We experimented with different dataset sparsity level as well as feature sparsity levels.

4.2.2 Synthetic Dataset using Power Law Distribution

As we saw in section 3.4 for skewed distributions like Power Law and Zipfian, count min sketch should provide smallest error. To understand the performance of our novel data structure we generated the synthetic data examples under power law distribution. We used multiple values of $\alpha = 2, 3$ and selected minimum value of $k=1$ and maximum value of $k=$ number of features. We experimented with different dataset sparsity levels and feature sparsity levels to understand the performance of our data structure in comparison to other sketch variants for heavy hitters as well as feature selection task.

5 Experiments and Results

5.1 Analyzing Gradient Updates

Complementary Count Min Sketch will give least error when the error in positive and negative count min sketch negate each other. Since our experiments with RCV1 and KDD dataset performed well with CCMS data structure, we believed this was due to negation of the sketch errors. So we decided to analyze the gradient updates of each dimension of the feature vector. We performed comparative analysis of gradient updates for different approaches of feature selection using all the variants (including CCMS, conservative CCMS, CS, logistic regression, and standard logistic regression with top-k) to understand how gradient updates are being store and how top-k is changing over time.

In figure 5 we present few samples for gradient updates. Analyzing these sample gradient updates, we observe that feature 13 and feature 29 are behaving completely opposite. Where in feature 13, number of gradient updates of count sketch are lesser in comparison to standard logistic regression method, in feature 29 number of updates are more for count sketch. If we analyze feature 25, gradients fluctuations have higher magnitudes for count sketches in comparison to logistic regression updates but number of updates in logistic regression are way highers then number of updates in count sketch. To sum up, we can say we didn't find any specific pattern in gradient updates to conclude anything.

648
649
650
651
652
653

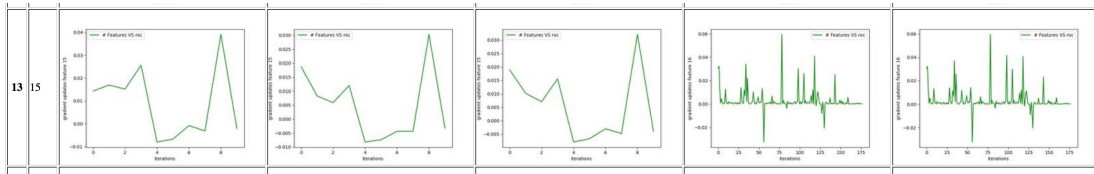


Figure 2: Gradient updates for feature 13

654
655
656
657
658
659
660

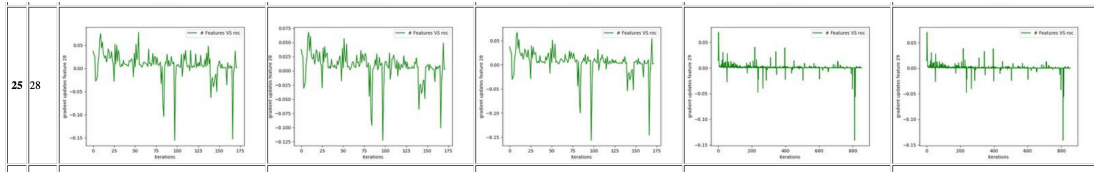


Figure 3: Gradient updates for feature 25

661
662
663
664
665
666
667

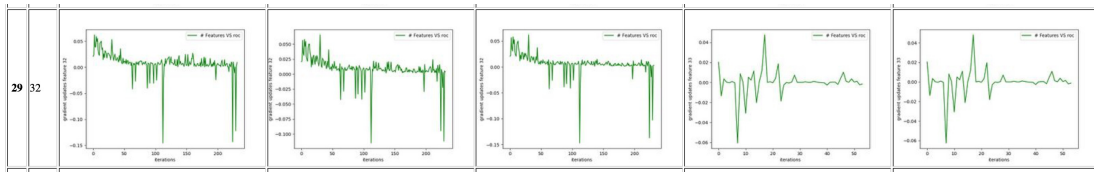


Figure 4: Gradient updates for feature 29

668
669
670
671
672
673
674

Figure 5: In these gradient updates, first column belongs to CCMS, second column belongs to complementary CCMS, third column belongs to CS, fourth column belongs to logistic regression and fifth belongs to logistic regression with heap

675 5.2 Performance of Datasets for feature selection

676
677
678
679

We performed experiments using different feature selection techniques for both RCV1 dataset and synthetic dataset.

680
681
682

RCV1 Dataset follows power law. The results of RCV1 dataset with top $K = 8000$ are summarised in the table below:

683
684
685
686
687
688
689
690
691
692
693
694
695
696

	Num Hash Fun.	Sketch Size	Total Space	Total Time (Sec.)	Accuracy
Logistic Regression	1	47236	47236		97.65%
Logistic Regression with Heap	1	47236	47236		97.64%
CS	3	19000	57000	45.53	95.33%
CCMS	2	14000	56000	34.42	95.36%
Conservative CCMS	2	14000	56000	44.84	94.98%
Greedy Thresholding	1	47236	47236		81.94%

697
698

Table 2: Results for RCV1 Dataset

699
700
701

We have approximately allocated equal space to all the sketches. We observe that accuracy of logistic regression outperforms all of them. All the sketches nearly give same accuracy. The time taken by CCMS is relatively less compared to the other sketches which is as expected.

702
703
704
705
706
707
708
709
710
711
712
713
714
715
716

	Num Hash Func	Sketch Size	Total Space	Accuracy
Logistic Regression	1	10000	10000	91.32%
Logistic Regression with Heap	1	10000	10000	92.45%
CS	3	2000	6000	78.55%
CCMS	2	1500	6000	74.77%
Conservative CCMS	2	1500	6000	74.41%
Greedy Thresholding	1	10000	10000	74.42%

717
718
719

Table 3: Results for Gaussian Synthetic Dataset

720
721
722
723
724

Above are the summarised results for synthetic Gaussian dataset with top $K = 100$. As expected, the standard logistic regression outperforms as we do not compress any gradients. We observed all other approaches nearly perform similarly. We haven't stated number of most important features recovered as all the methods perform poorly when we consider a Gaussian dataset.

725
726
727
728
729
730
731
732
733
734
735
736
737
738
739

	Num Hash Func	Sketch Size	Total Space	Accuracy	Positions Recovered
Logistic Regression	1	10000	10000	87.5%	80
Logistic Regression with Heap	1	10000	10000	86.35%	80
CS	3	2000	6000	87.80%	80
CCMS	2	1500	6000	87%	80
Conservative CCMS	2	1500	6000	88.67%	80
Greedy Thresholding	1	10000	10000	54.74%	80

740
741

Table 4: Results for Power Law Synthetic Dataset

742
743
744
745

From the summarised results for synthetic power law following dataset with top $K = 100$. We observed all approaches nearly perform similarly except greedy thresholding. All approaches are able to require the true important features.

746 5.3 Comparing Top-K Heap for various sketches

747
748
749
750
751
752

For RCV1 dataset, we were not aware of true parameters of the model and were clueless what should be right baseline to compare with. So we decided to compare the most important features obtained from each technique. Here we observed something very strange, the match in most important features varied a lot for each technique used. These results made us pivot towards experiments with synthetic dataset.

753
754
755

Also, we analyzed the magnitude of updates for most impacting features which vary a lot. We notice the range of weights for CCMS is much wider compared to cs. We can say range of weights of CCMS is closer to the range given by logistic regression.

756
757
758
759
760
761
762
763
764
765
766
767
768
769
770
771
772
773
774
775
776
777
778
779
780
781
782
783
784
785
786
787
788
789
790
791
792
793
794
795
796
797
798
799
800
801
802
803
804
805
806
807
808
809

	Logistic Regression	Logistic Regression with Top K	Count Sketch	Complementary Count Sketch	Conservative Complementary Count Sketch
Logistic Regression	100	90.45	52.45	27.78	48.31
Logistic Regression with Top K	90.45	100	55.81	28.63	51.81
Count Sketch	52.45	55.81	100	29.47	55.72
Complementary Count Sketch	27.78	26.83	29.47	100	29.35
Conservative Complementary Count Sketch	48.31	51.81	55.72	29.35	100

Figure 6: Top K Overlap for all feature selection methods over RCV1 Dataset

Minimum weights:

Logistic Regression	Logistic Regression with Top K	Count Sketch	Complementary Count Sketch	Conservative Complementary Count Sketch
-35.90	-36.10	-5.75	-20.27	-6.18

Maximum weights:

Logistic Regression	Logistic Regression with Top K	Count Sketch	Complementary Count Sketch	Conservative Complementary Count Sketch
21.32	21.48	-9.10	11.29	11.48

Figure 7: Maximum and Minimum weights for all feature selection methods over RCV1 Dataset

5.4 Frequency estimation

To estimate frequency we created a synthetic dataset with followed power law distribution with varying α . The dataset included positive and negative numbers where negative number indicated deletion of the particular element from the stream. We observed that the error in frequency estimation for frequent hitters using CCMS was decent when the data set was very skewed. Also, we observe the first proposed CCMS outperforms compared to the CCMS with different hash functions 4 and CCMS using single sketch 5.

alpha values	CS loss	CCMS loss	CCMS using different hash functions	CCMS using Single Sketch
2	3.27	12.79	13.87	100.17
3	0.025	0.27	0.20	1.84

Table 5: Mean Square Loss for top 30 items from variants of count sketches

810
811
812
813
814
815
816
817
818
819
820
821
822
823
824
825
826
827
828
829
830
831
832
833
834
835
836
837
838
839
840
841
842
843
844
845
846
847
848
849
850
851
852
853
854
855
856
857
858
859
860
861
862
863

6 Conclusion

The analysis of complementary count min sketch does not give a better bound than count sketch in terms of error. However, they have promising results in terms of time taken compared to count sketch. From the results we observed under same space constraint count sketch and complementary count min sketch both give same error bound which is $L1$ -norm of true frequency vector when dataset is very skewed.

References

- [1] Amirali Aghazadeh, Ryan Spring, Daniel LeJeune, Gautam Dasarathy, Anshumali Shrivastava, and Richard G. Baraniuk. MISSION: ultra large-scale feature selection using count-sketches. In *Proceedings of the 35th International Conference on Machine Learning, ICML 2018, Stockholm, Sweden, July 10-15, 2018*, pages 80–88, 2018.
- [2] Moses Charikar, Kevin Chen, and Martin Farach-Colton. Finding frequent items in data streams. In *Proceedings of the 29th International Colloquium on Automata, Languages and Programming, ICALP '02*, pages 693–703, Berlin, Heidelberg, 2002. Springer-Verlag.
- [3] Graham Cormode and Marios Hadjieleftheriou. Methods for finding frequent items in data streams. *The VLDB Journal*, 19(1):3–20, February 2010.
- [4] Kilian Weinberger, Anirban Dasgupta, John Langford, Alex Smola, and Josh Attenberg. Feature hashing for large scale multitask learning. In *Proceedings of the 26th Annual International Conference on Machine Learning, ICML '09*, pages 1113–1120, New York, NY, USA, 2009. ACM.
- [5] A. Maleki. Coherence analysis of iterative thresholding algorithms. In *2009 47th Annual Allerton Conference on Communication, Control, and Computing (Allerton)*, pages 236–243, 2009.
- [6] Kai Sheng Tai, Vatsal Sharan, Peter Bailis, and Gregory Valiant. Sketching linear classifiers over data streams. In *Proceedings of the 2018 International Conference on Management of Data, SIGMOD '18*, pages 757–772, New York, NY, USA, 2018. ACM.
- [7] Morten Houmller Nygaard Jonas Nicolai Hovmand. Estimating frequencies and finding heavy hitters, 2006.
- [8] David D. Lewis, Yiming Yang, Tony G. Rose, and Fan Li. Rcv1: A new benchmark collection for text categorization research. *J. Mach. Learn. Res.*, 5:361–397, December 2004.