

Count-min sketch: extension to negative counters

Check the analysis seen in class, and discuss how to allow $F[i]$ to change by arbitrary values read in the stream. Namely, the stream is a sequence of pairs of elements, where the first element indicates the item i whose counter is to be changed, and the second element is the amount v of that change (v can vary in each pair). In this way, the operation on the counter becomes $F[i] = F[i] + v$, where the increment and decrement can be now seen as $(i, 1)$ and $(i, -1)$.

SOLUTION

The frequency of the item i is represented by $\hat{F}[i] \rightarrow T[j, h_j(i)] = F[i] + X_{ji}$, where X_{ji} represent the garbage introduced by the other counters. If we just increment the counter X_{ji} is going to be positive, and then we can take the $\min_j T[j, h_j(i)]$ to approximate $F[i]$. Instead, if we have also decrement, it could happen that $X_{ji} < 0$ therefore the method for the min is not going to work. In this case, we consider the absolute value of X_{ji} (i.e. $|X_{ji}|$) and we use $\text{median}_j T[j, h_j(i)]$ to approximate $F[i]$. Now let's prove that, with probability $1 - \delta^{1/4}$ holds:

$$F[i] - 3\epsilon\|F\| \leq \hat{F}[i] \leq F[i] + 3\epsilon\|F\|$$

First, let's do some consideration. The value of $\hat{F}[i] = \text{median}_j T[j, h_j(i)]$, and $T[j, h_j(i)] = F[i] + |X_{ji}|$. The first inequality (i.e. $F[i] \leq \hat{F}[i]$) holds because we took the absolute value of X_{ji} . Now we shall prove that $\Pr[|X_{ji}| > 3\epsilon\|F\|] \leq 1/8$. Taken j such that $\hat{F}[i] = \text{median}_j T[j, h_j(i)] = F[i] + |X_{ji}|$, then we have:

$$\begin{aligned} \Pr[F[i] - 3\epsilon\|F\| \leq \hat{F}[i] \leq F[i] + 3\epsilon\|F\|] \\ \Pr[F[i] - 3\epsilon\|F\| \leq F[i] + |X_{ji}| \leq F[i] + 3\epsilon\|F\|] \\ \Pr[-3\epsilon\|F\| \leq |X_{ji}| \leq 3\epsilon\|F\|] \\ \Pr[|X_{ji}| \leq 3\epsilon\|F\|] \end{aligned}$$

From what we have seen in class and by the property of the absolute value, we have $E[|X_{ji}|] \leq E[X_{ji}] = \frac{\epsilon}{e}\|F\|$. Then we can apply the Markov inequality and since universal hash functions are pairwise independent, we have:

$$\Pr[|X_{ji}| > 3\epsilon\|F\|] < \frac{E[|X_{ji}|]}{3\epsilon\|F\|} \leq \frac{\frac{\epsilon}{e}\|F\|}{3\epsilon\|F\|} = \frac{1}{3e} < \frac{1}{8}$$

Let's now define the indicator variable $Y = \sum_{j=0}^r Y_j$ which tell us the number of elements i (columns of the sketch) that have a garbage $|X_{ji}|$ greater than $3\epsilon\|F\|$.

$$Y_j = \begin{cases} 1 & \text{IF } |X_{ji}| > 3\epsilon\|F\| \text{ WITH } p < \frac{1}{8} \\ 0 & \text{OTHERWISE} \end{cases}$$

The median of $\hat{F}[i]$ is going to be a good approximation if we haven't got more than $\frac{r}{2}$ rows such that $|X_{ji}| > 3\epsilon\|F\|$ (that is we want $Y < \frac{r}{2}$).

Therefore, to calculate the probability of error, we calculate the probability of $\Pr[Y \geq \frac{r}{2}]$. Here, we can use the Chernoff's Bound With: $(1 + \lambda)\mu = \frac{r}{2}$, $\mu = E[Y] = rp$.

$$\Pr[Y \geq (1 + \lambda)\mu] < \left[\frac{e^\lambda}{(1 + \lambda)^{1+\lambda}} \right]^\mu = \left[\frac{e}{e(1 + \lambda)^{1+\lambda}} \right]^\mu = \frac{1}{e^\mu} \left[\frac{e}{(1 + \lambda)} \right]^{(1+\lambda)\mu} = \frac{1}{e^{rp}} \left[\frac{1}{(1 + \lambda)} e \right]^{\frac{r}{2}} = \frac{1}{e^{rp}} [2pe]^{\frac{r}{2}}$$

Now we need to prove that $\frac{1}{e^{rp}} [2pe]^{\frac{r}{2}} \leq \delta^{\frac{1}{4}} = \frac{1}{2^{\frac{r}{4}}}$. If we use the reciprocal we have:

$$\begin{aligned} 2^{\frac{r}{4}} &\leq \frac{e^{rp}}{[2pe]^{\frac{r}{2}}} \leq \frac{1}{[2pe]^{\frac{r}{2}}} \\ 2^{\frac{1}{4}} &\leq \frac{1}{\sqrt{2pe}} \end{aligned}$$

Since $e^{rp} \geq 1$. Then, we take $\frac{1}{2pe} > \sqrt{2}$, that is possible just if $p < \frac{1}{2\sqrt{2}e}$, indeed $p = \frac{1}{8}$.