**EX**: To estimate the size of set $S$ of labeled elements, sample with replacement from $S$ from a uniform distribution until the first duplicate. If $k$ elements have been drawn, then an unbiased estimator is $|S| = 2k^2 / \pi$.

- - - - - - - - - - - - - - - - - - - - - - - - - - - - - - - - - - - - - - - - - - - - - - - - - - - - - - -

**EX**: Counting distinct words  Flajolet & Martin, "Probabilistic Counting Algorithms for Database Applications", *J. Comp. Syst. Sciences*, 1985 Assume there are $n$ (not necessarily distinct) words in a text, and you want to estimate the # of distinct words.  Let there be a hashing function $h : \{words\} \to \{0,1\}^{5+\lg n}$.  If $s$ is a string of bits, let $\pi(x,b)$, $b \in \{0,1\}$, denote the index of the leftmost bit of $x$ equal to $b$

$signature = s_1 s_2 ... s_{5+\lg n} = 00...0$

   **for** each word $x$ **do** $signature_{\pi(h(x),1)} \leftarrow 1$

$$\mathbf{return}\left( \frac{2^{\pi(signature,0)}}{1.574} \right)$$

If $\pi(signature,0)=4$, then leftmost bits of $signature = 1110$.  If there are 16 distinct words on the tape, $\Pr\{\text{none of hash codings begins } 0001\} = \left( \dfrac{15}{16} \right)^{16} = .356$

\*\*\*\*\*\*\*\*\*\*\*\*\*\*\*\*\*\*\*\*\*\*\*\*\*\*\*\*\*\*\*\*\*\*\*\*\*\*\*\*\*\*\*\*\*\*\*\*\*\*\*\*\*\*\*\*\*\*\*\*\*\*\*\*\*\*\*\*\*\*\*\*\*\*\*\*\*\*\*\*\*\*\*\*

Hash Coding - Chapter II Assume $n$ balls uniformly independently distributed among $m \geq n$ bins. What is probability of no collisions?

$$\prod_{1 \leq j \leq n-1}\left( 1 - \frac{j}{m} \right) \approx \prod_{1 \leq j \leq n-1} e^{-\frac{j}{m}} = e^{-\sum_{1 \leq j \leq n-1}\frac{j}{m}} = e^{-\frac{n(n-1)}{2m}} \approx e^{-\frac{n^2}{2m}}$$

where first approximation justified by $e^{-j/m} = \left( -\dfrac{j}{m} \right)^0 + \dfrac{(-j/m)^1}{1!} + \dfrac{(-j/m)^2}{2!} + ... \approx 1 - \dfrac{j}{m}$

for $j \ll m$. For Birthday Paradox, $e^{-\frac{n^2}{2m}} \approx \dfrac{1}{2}$ implies $\dfrac{n^2}{2m} \approx \ln 2 \Rightarrow n = \sqrt{2m \ln 2}$ and $m=365 \Rightarrow n = 22.49$

A *hash function* $h : U \to [0...m-1]$ where $|U| \gg m$. A *good* $h$ distributes sampled subset $S = \{a_1,...,a_n\}$ of $U$ uniformly among $m$ *buckets* (bins).

**EX**: (Probabilistic set membership)  Want to test if $p \in S = \{a_1,...,a_n\}$  $h : U \to [2^b]$

      Need $\Pr\{\text{false negative}\}=0$ and low $\Pr\{\text{false positive}\}$.

Preprocess $\{h(a_1),...,h(a_n)\}$.

         **for** $i \leftarrow 1$ **to** $2^b$ **do** $A[i] \leftarrow 0$

         **for** $i \leftarrow 1$ **to** $n$ **do** $A[h(a_i)] \leftarrow 1$

Given $p$, test if $\exists x \in S\, h(x) = h(p)$.

         **return** $A[h(x)]$

ANALYSIS:  Fix $x \in S$  $\Pr\{h(x) = h(p)\} = 2^{-b}$.  $\Pr\{h(x) \neq h(p)\} = 1 - 2^{-b}$.

$\Pr\{\text{false positive}\} \geq 1 - \left(1 - 2^{-b}\right)^n \geq 1 - e^{-n/2^b}$ If we want $\Pr\{\text{false positive}\} \leq c \implies$

$1 - e^{-n/2^b} \geq 1 - c \implies b \geq \lg \dfrac{n}{\ln\left(1/(1-c)\right)}$ We need $\Omega(\lg n)$ bits.

**Ex**: Bloom Filters Boolean $A[0..m-1]$ boolean, initialized to 0

$k$ independent hash functions $h_1, \ldots, h_k : U \to [0, \ldots, m-1]$

> **for** $i \leftarrow 1$ **to** $n$ **do**
> > **for** $j \leftarrow 1$ **to** $k$ **do**
> > > $A\left[h_j(a_i)\right] \leftarrow 1$
>
> MEMBER$?(p, A)$
> > **for** $j \leftarrow 1$ **to** $k$ **do**
> > > **if** $A\left[h_j(a_i)\right] = 0$ **then return** $p \notin S$
> >
> > **return** $p \in S$

$\Pr\{\text{false negative}\} = 0$. $\Pr\{\text{false positive}\} - \forall j \ \Pr\{A[j] = 0\} = \left(1 - \dfrac{1}{m}\right)^{kn} \approx e^{-kn/m}$

So $\Pr\{\text{false positive}\} = \left(1 - \left(1 - \dfrac{1}{m}\right)^{kn}\right)^k \approx \left(1 - e^{-kn/m}\right)^k$.


Choose larger or smaller $k$?

> · Larger $k \implies$ For $p \notin S$ more chances to find a 0 in $A$
> · Larger $k \implies$ more 0's in $A$

Want to choose $k$ to minimize $\left(1 - e^{-kn/m}\right)^k$? Taking derivative, global minimum is

$k = (\ln 2)(m/n) \implies$ So $\Pr\{\text{false positive}\} = 2^{-k} \approx 0.6185^{m/n}$, which falls exponentially with $m/n$.

\*\*\*\*\*\*\*\*\*\*\*\*\*\*\*\*\*\*\*\*\*\*\*\*\*\*\*\*\*\*\*\*\*\*\*\*\*\*\*\*\*\*\*\*\*\*\*\*\*\*\*\*\*\*\*\*\*\*\*\*\*\*\*\*\*\*\*\*\*\*\*\*\*\*\*\*

**FINDING AVERAGE GRADE** $n > 2$ friends want to compute average grade while conveying no knowledge about their grades, $s_1, \ldots, s_n$. No player cheats; no trusted $(n+1)^{\underline{st}}$ party.

> Choose $m >$ largest possible $\displaystyle\sum_{n \geq i \geq 1} s_i$
>
> Each player $i$ chooses $x_{i,1}, \ldots, x_{i,n-1} \sim U[0, m-1]$
>
> Each player $i$ chooses $x_{i,n}$ to satisfy $\displaystyle\sum_{n \geq k \geq 1} x_{i,k} = s_i \bmod m$
>
> Each player $i$ distributes one of $\left\{x_{i,1}, \ldots, x_{i,n}\right\}$ (without replacement) to every

player

> Each player $i$ computes and announces sum of numbers held (mod $m$) (say $S_i$)
>
> Average grade $= \displaystyle\sum_{n \geq i \geq 1} s_i / n = \left.\displaystyle\sum_{n \geq i \geq 1} S_i \ (\text{mod m})\right/ n$

<u>Claim</u>: For any set $T$ of friends, $|Y| < n-1$, the only information they can infer about other grades is $\sum_{i \notin T} S_i$

**DEF**: Events $A$, $B$ *independent* if $\Pr[A] = \Pr[A|B]$. *k-wise independent.*

(Return to) <u>Claim</u>: $(\forall i)\left(\{x_{i1}, \ldots x_{in}\} \text{ is } (n-1) - \text{wise independent}\right)$