

Q) How to determine how efficient Bloom filters are?

→ Say there are n sets whose membership is to be represented using a bloom filter of m bits, with a false +ve rate of ϵ

→ "P" of a single bit not set any time: $(1 - \frac{1}{m})$

> Assuming the hash fns are very random

→ "P" of none of the m bits are set: $(1 - \frac{1}{m}) \times (1 - \frac{1}{m}) \times \dots$
 $= (1 - \frac{1}{m})^m$

→ when there are k hash fns: $(1 - \frac{1}{m})^{nk}$

→ when we test for an element that doesn't exist k hashes

"P" of all of those k bits to be set: "P" is false +ve

$$P = \left(1 - \left(1 - \frac{1}{m}\right)^{nk}\right)^k \quad \text{--- (1)}$$

→ let's approximate m being very large to find true efficiency of bloom filters

$$\lim_{m \rightarrow \infty} \left(1 - \frac{1}{m}\right)^k$$

$$= \lim_{m \rightarrow \infty} \left(1 + \frac{(-1)}{m}\right)^k \quad \text{--- (2)}$$

for real world, the approximation of (1) = (2) is close enough:

$$\Rightarrow \lim_{m \rightarrow \infty} \left(1 + \frac{(-1)}{m}\right)^k = e^{-1}$$

Using:

$$\left[\lim_{n \rightarrow \infty} \left(1 + \frac{x}{n}\right)^n = e^x \right]$$

$$\left[\lim_{m \rightarrow \infty} \left(1 - \frac{1}{m}\right)^m = e^{-1} \right]$$

(1)

$$\left(1 - \frac{1}{m}\right)^k = \left(\left(1 - \frac{1}{m}\right)^m\right)^{k/m} = e^{-k/m} \approx \left(1 - \frac{1}{m}\right)^k \quad \textcircled{3}$$

Sub ③ in ①

$$P = \left(1 - e^{-\frac{kn}{m}}\right)^k$$

$$P(k) = e^{\ln\left(1 - e^{-\frac{kn}{m}}\right)^k} \Rightarrow P(k) = e^{k \ln\left(1 - e^{-\frac{kn}{m}}\right)}$$

> let's differentiate the power of e . to make it easier to differentiate against k as the minima remains same. (they have the same x intercept any ways)

$$\frac{d}{dk} \left[k \ln\left(1 - e^{-\frac{kn}{m}}\right) \right] = \ln\left(1 - e^{-\frac{kn}{m}}\right) + k \frac{\frac{n}{m} \times e^{-\frac{kn}{m}}}{1 - e^{-\frac{kn}{m}}}$$

$$\text{let } e^{-\frac{kn}{m}} = q \Rightarrow \ln(q) = -\frac{kn}{m}$$

$$\textcircled{4} \Rightarrow \ln(1-q) - \ln(q) \frac{q}{1-q} = 0 \quad (\text{eq to 0 to find min})$$

$$\Rightarrow q = 1/2 \Rightarrow e^{-kn/m} = 1/2$$

$$-\frac{kn}{m} = \ln(1/2)$$

$$k = \frac{\ln(2)m}{n} \quad \textcircled{5}$$

$$\text{false rate} = \epsilon = \left(1 - e^{-\frac{kn}{m}}\right)^k \quad (\text{from eq } \textcircled{0})$$

LHS

$$= \left(1 - e^{-\ln(2) \frac{m}{n} \times \frac{n}{m}}\right)^{\ln(2) \frac{m}{n}}$$

$$= \left(1 - e^{-\ln(2)}\right)^{\frac{\ln(2)m}{n}} = \left(1 - e^{\ln(1/2)}\right)^{\frac{\ln(2)m}{n}}$$

$$= \left(1 - 1/2\right)^{\frac{\ln(2)m}{n}} = \left(1/2\right)^{\frac{\ln(2)m}{n}}$$

RHS

$$\epsilon = \frac{1}{2} \log_{1/2}(\epsilon) \Rightarrow$$

$$\frac{1}{2} \log_{1/2}(\epsilon) = \frac{1}{2} \frac{\ln(2)m}{n}$$

(LHS) (RHS)

$$\log_{1/2}(\epsilon) = \frac{\ln(2)m}{n} \Rightarrow \frac{m}{n} = \frac{\log_{1/2}(\epsilon)}{\ln(2)}$$

$$\log_{1/2}(\epsilon) = \ln(2) \frac{m}{n}$$

$$\frac{\log_2(\epsilon)}{\log_2(1/2)} = \ln(2) m/n$$

$$(-1) \log_2(\epsilon) = \ln(2) m/n$$

$$\textcircled{f2} \quad \leftarrow \quad \boxed{m/n = -1.44 \log_2(\epsilon)}$$

So, to decipher what $\textcircled{f2}$ is saying, say we want to rep n sets using m bits:

$$m = n \times 1.44 \times \log_2(\epsilon)$$

You would need $1.44 \times \log_2(\epsilon)$ # of bits for each member of each set of n , where the ideal minima would be $1 \times \log_2(\epsilon)$ for each member of n , where ϵ is the false +ve rate.

IDEALLY: $\log_2(\text{rate})$ would be the # of required bits.