

Signature files (Faloutsos' 92)

- Quick and dirty, based on hashing
- Underestimate the distances (important!!)

Maria Luisa Sapino (BDM 2018)

Signature files

- Quick and dirty, based on hashing
- Underestimate the distances (important!!)

Kasim	00110101
Selcuk	10010110

(bit-or) 10110111

Maria Luisa Sapino (BDM 2018)

Signature files

- Quick and dirty, based on hashing
- Underestimate the distances (important!!)

Kasim	00110101
Selcuk	10010110

(bit-or) 10110111

- Length of the signature?
- Ratio of the 1s to 0s?

Maria Luisa Sapino (BDM 2018)

Signature files

- l of B bits are randomly set
 - too many 1s would cause false hits
- Minimize false hits

Maria Luisa Sapino (BDM 2018)

Signature files

- l of B bits are randomly set
 - too many 1s would cause false hits
- Minimize false hits
- If there are b keywords, the probability that a given bit is set is

$$1 - \left(1 - \frac{1}{B}\right)^{bl}$$

Maria Luisa Sapino (BDM 2018)

Signature files

- l of B bits are randomly set
 - too many 1s would cause false hits
- Minimize false hits
- If there are b keywords, the probability that a given bit is set is

$$1 - \left(1 - \frac{1}{B}\right)^{bl} \approx 1 - e^{-b\alpha} \quad \alpha = \frac{l}{B}$$

Maria Luisa Sapino (BDM 2018)

Signature files

- l of B bits are randomly set
 - too many 1s would cause false hits
- Minimize false hits
- Probability that l random bits in the query are also set in the signature

$$(1 - e^{-b\alpha})^l = (1 - e^{-b\alpha})^{\alpha B}$$

Maria Luisa Sapino (BDM 2018)

Signature files

- l of B bits are randomly set
 - too many 1s would cause false hits
- Minimize false hits
- Probability that l random bits in the query are also set in the signature

$$(1 - e^{-b\alpha})^l = (1 - e^{-b\alpha})^{\alpha B}$$

Minimized if $\alpha = \frac{\ln(2)}{b}$

Maria Luisa Sapino (BDM 2018)

Signature files

- l of B bits are randomly set
 - too many 1s would cause false hits
- Minimize false hits
- Optimal l is $B \frac{\ln(2)}{b}$

Maria Luisa Sapino (BDM 2018)

Signature files

- l of B bits are randomly set
 - too many 1s would cause false hits
- Minimize false hits

- Optimal l is $B \frac{\ln(2)}{b}$

- False hit rate under the optimal l is $\frac{1}{2^l}$

Maria Luisa Sapino (BDM 2018)

Signature files

- l of B bits are randomly set
 - too many 1s would cause false hits
- Minimize false hits

- Optimal l is $B \frac{\ln(2)}{b}$

There is a tradeoff between the length of the signature and the false hits!

- False hit rate under the optimal l is $\frac{1}{2^l}$

Maria Luisa Sapino (BDM 2018)

Signature files

- l of B bits are randomly set
 - too many 1s would cause false hits
- Minimize false hits

- Optimal l is $B \frac{\ln(2)}{b}$

There is a tradeoff between the length of the signature and the false hits!

- False hit rate under the optimal l is $\frac{1}{2^l}$

Experimentally: a good α is 0.5

Maria Luisa Sapino (BDM 2018)
