

SCAM: A Copy Detection Mechanism for Digital Documents

Narayanan Shivakumar, Hector Garcia-Molina
Department of Computer Science
Stanford University
Stanford, CA

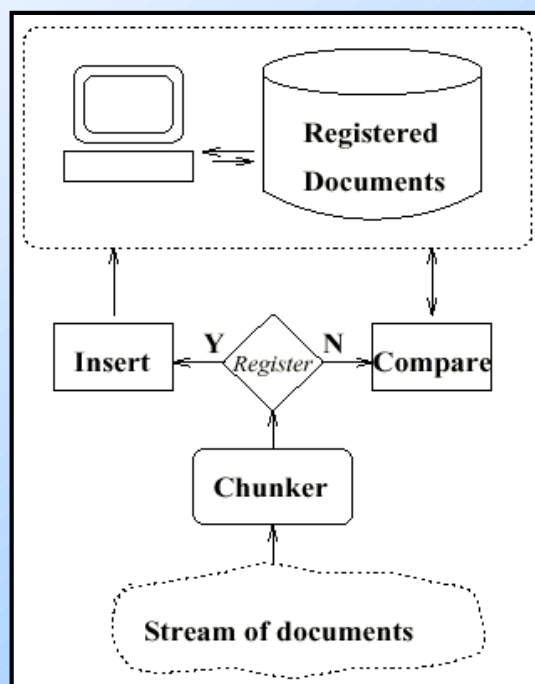
Presented by Dennis Peper

Introduction

- Digital libraries may become popular in the future.
- We don't know how to prevent illegal copying effectively.
- We can detect illegal copying. Methods are:
 - Signature based (e.g. Watermarking.)
 - Registration based (SCAM's approach)
 - Original document is registered with a server.
 - Subsequent documents are **compared** with pre-registered documents.

Copy Detection Overview

- Document is broken up into chunks.
 - A chunk could be a word, sentence or sequence of sentences, depending on the detection method.
 - SCAM is based on word chunking.
- The chunks of a registered document are stored in a database.
- A new document is broken up and compared to the chunks in the database.
- If **overlap** exceeds a certain threshold (depending on detection system used), then the author is notified.



Measuring Overlap

- Vocabulary $V = \{a,b,c,d,e,f,g,h\}$ are the unique chunks in our library.
- Given a registered document $R = \langle a,b,c \rangle$. What should our system return, when comparing R to other documents D_1, D_2, \dots ?

<u>Document</u>	<u>Comparison with R</u>
$D_1 = \langle a,b \rangle$	Quite similar
$D_2 = \langle a,b,c \rangle$	Exact replica
$D_3 = \langle a^k \rangle \quad k \geq 1$	Somewhat similar at low k, not very similar for high k. e.g. when $k=2$, $D_3 = \langle a,a \rangle$
$D_4 = \langle a,b,c,d,e,f,g,h \rangle$	Significant overlap

Vector Space Model

- We need a method to compare registered document R to incoming documents.
- F(d) is the frequency vector. (e.g. $F(R)=\langle 1,1,1,0,\dots \rangle$)
- Dot product of normalized frequency vectors.
 - $R=\langle a,b,c \rangle$ $V(R)=\langle 1/3,1/3,1/3,0,\dots \rangle$
 - $D_1=\langle a,b \rangle$ $V(D_1)=\langle 1/2,1/2,0,0,\dots \rangle$
 - $V(R)\cdot V(D_1)=1/3 \times 1/2 + 1/3 \times 1/2 = 1/3$

There is 1 'a' chunk, 1 'b' chunk, etc...

<u>Document</u>	<u>Dot Product</u>
$D_2=\{a,b,c\}$	$V(R)\cdot V(D_2)=1/3$
$D_3=\{a^k\} \quad k \geq 1$	$V(R)\cdot V(D_3)=1/3$
$D_4=\{a,b,c,d,e,f,g,h\}$	$V(R)\cdot V(D_4)=1/8$

- The overlap of R with D_2 and D_4 is more significant than that with D_1 or D_3 .
- The dot product does not report this.

Vector Space Model (Cont.)

2. Cosine Similarity measure. If R and Q are documents:

$$sim(R, Q) = \frac{\sum_{i=1}^N \alpha_i^2 * F_i(R) * F_i(Q)}{\sqrt{\sum_{i=1}^N \alpha_i^2 * F_i^2(R) * \sum_{i=1}^N \alpha_i^2 * F_i^2(Q)}}$$

- The higher the frequency of the word, the less the word contributes to the similarity result.
- $R=\{a,b,c\}$ $F(R)=\langle 1,1,1,0,\dots \rangle$
- $D_1=\{a,b\}$ $F(D_1)=\langle 1,1,0,0,\dots \rangle$
- $sim(R, D_1) = 0.2$

α_i is a weighting associated with the i^{th} chunk, set to 1 in the example.

HIDDEN SLIDE

<u>Document</u>	<u>Cosine Similarity Measure</u>
$D_2=\{a,b,c\}$	$sim(R, D_2)=1$
$D_3=\{a^k\} \quad k \geq 1$	$sim(R, D_3)=0.58 \longrightarrow F(D_3)=\langle k,0,0,\dots \rangle$
$D_4=\{a,b,c,d,e,f,g,h\}$	$sim(R, D_4)=0.61$

- Similarity for D_3 is independent of k. Ideally similarity decreases as k increases.
- Similarity for D_4 is low considering the entirety of R contained within D_4 .

Relative Frequency Model - SCAM(?)

- For documents R and S define a closeness set $c(R,S)$ to contain all words w_i with similar frequencies $F_i(R)$ and $F_i(S)$.
 - “Similar” is defined by $s_i = F_i(R)/F_i(S) + F_i(S)/F_i(R) < \epsilon$.
 - ϵ must be chosen in the range 2.0001 to infinity.
 - If either $F_i(S)$ or $F_i(R) = 0$, then don't include s_i in $c(R,S)$
- Note: if a word occurs the same number of times in both documents, it is added.
- If $F_i(R)=3$ and $F_i(S)=2$, then $c(R,S)$ contains w_i if $\epsilon < 2.17$.
- We can see that ϵ is a tolerance factor.

Relative Frequency Model (Cont.)

- Subset measure. How much of a subset is R wrt S? (my understanding.)
- $\text{subset}(R,S) = \left(\sum_{w_i \in c(R,S)} \alpha_i^2 F_i(R)F_i(S) \right) / \left(\sum_i \alpha_i^2 F_i^2(R) \right)$
- α_i is a weighting associated with the i^{th} chunk, this allows adding an importance factor.
- $\text{sim}(R,S) = \max \{ \text{subset}(R,S), \text{subset}(S,R) \}$
- If $\text{sim}(R,S) \geq 1$, set the result to 1, no extra information is gained as both documents are very similar anyway.
- $R = \{a,b,c\}$ $D_1 = \{a,b\}$
- $\text{sim}(R,D_1) = \max \{ (1 \times 1 + 1 \times 1) / 3, (1 \times 1 + 1 \times 1) / 2 \} = 1$ ($\epsilon = 2.0001$ or 3)

Document	sim result, $\epsilon = 2.0001$	sim result, $\epsilon = 3$
$D_2 = \{a,b,c\}$	$\text{sim}(R,D_2) = 1$	$\text{sim}(R,D_2) = 1$
$D_3 = \{a^k\} \quad k \geq 1$	$\text{sim}(R,D_3) = 1, k=1$ $\text{sim}(R,D_3) = 0, k > 1$	$\text{sim}(R,D_3) = 1, k=1$ $\text{sim}(R,D_3) = 2/3, k=2$ $\text{sim}(R,D_3) = 0, k \geq 3$
$D_4 = \{a,b,c,d,e,f,g,h\}$	$\text{sim}(R,D_4) = 1$	$\text{sim}(R,D_4) = 1$

Conclusion

- What about detecting copying from multiple sources?
- Compared against another method, COPS, and results favored SCAM.
- Has higher false positive rate.

Thank you, any questions?