

ภาคผนวก 1

Math Tutorial

1. Summation

$$s = \sum_{i=0}^{n-1} i$$

2. Program of that summation

```
public class Sumup {  
    int n=10;  
    int s = 0;  
    int i = 0;  
    while (i <= (n-1)) {  
        s = s+i;  
        i = i + 1;  
    }  
}
```

Or

```
public class multup2 {  
    int n=10;  
    int s = 0;  
    int i;  
    int a[] = {0,1,2,3,4,5,6,7,8,9,10,11};  
    for (i = 1; i <= n; i++)  
        s = s + (a[i] * a[i+1]);  
}
```

3. Summation

$$S = \sum_{i=1}^n a_i \cdot a_{i+1}$$

```
public class multup {  
    int n=10;  
    int s = 0;  
    int i = 1;  
    int a[] = {0,1,2,3,4,5,6,7,8,9,10,11};  
    while (i <= n) {  
        s = s + (a[i] * a[i+1]);  
        i = i + 1;  
    }  
}
```

or...

```
public class multup2 {  
    int n=10;  
    int s = 0;  
    int i;  
    int a[] = {0,1,2,3,4,5,6,7,8,9,10,11};  
    for (i = 1; i <= n; i++)  
        s = s + (a[i] * a[i+1]);  
}
```

4. Simple tf * idf

$$w_{ik} = tf_{ik} * \log(N / n_k)$$

T_k = term k in document D_i

tf_{ik} = frequency of term T_k in document D_i

idf_k = inverse document frequency of term T_k in C

N = total number of documents in the collection C

n_k = the number of documents in C that contain T_k

$$idf_k = \log\left(\frac{N}{n_k}\right)$$

Inverse Document Frequency

For a collection	$\log\left(\frac{10000}{10000}\right) = 0$
of 10000 documents	$\log\left(\frac{10000}{5000}\right) = 0.301$
($N = 10000$)	$\log\left(\frac{10000}{20}\right) = 2.698$
	$\log\left(\frac{10000}{1}\right) = 4$

5. Similarity Measures

Simple matching (coordination level match)	$ Q \cap D $
Dice's Coefficient	$2 \frac{ Q \cap D }{ Q + D }$
Jaccard's Coefficient	$\frac{ Q \cap D }{ Q \cup D }$
Cosine Coefficient	$\frac{ Q \cap D }{ Q ^{1/2} \times D ^{1/2}}$
Overlap Coefficient	$\frac{ Q \cap D }{\min(Q , D)}$

6. tf*idf Normalization

$$w_{ik} = \frac{(tf_{ik} * \log(N / n_k))}{\sqrt{\left(\sum_{k=1}^t ((tf_{ik})^2 * ([\log(N / n_k)]^2))\right)}}$$

7. Vector Space Similarity

$$sim(D_i, D_j) = \sum_{k=1}^t w_{ik} * w_{jk}$$

8. Vector Space Similarity Measure

$$D_i = w_{d_{i1}}, w_{d_{i2}}, \dots, w_{d_{it}}$$

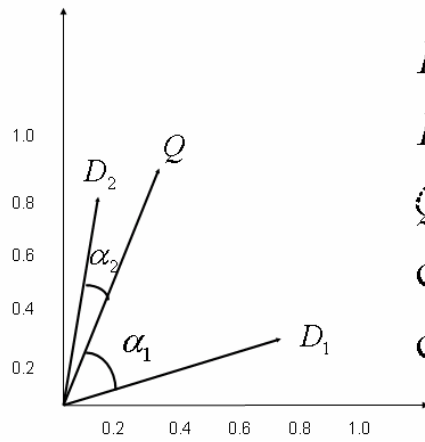
$$Q = w_{q1}, w_{q2}, \dots, w_{qt} \quad w = 0 \text{ if a term is absent}$$

if term weights are normalized: $sim(Q, D_i) = \sum_{j=1}^t w_{qj} * w_{d_{ij}}$

otherwise normalization and the similarity comparison are combined:

$$sim(Q, D_i) = \frac{\sum_{j=1}^t w_{qj} * w_{d_{ij}}}{\sqrt{\sum_{j=1}^t (w_{qj})^2 * \sum_{j=1}^t (w_{d_{ij}})^2}}$$

9. Computing Similarity Scores



$$D_1 = (0.8, 0.3)$$

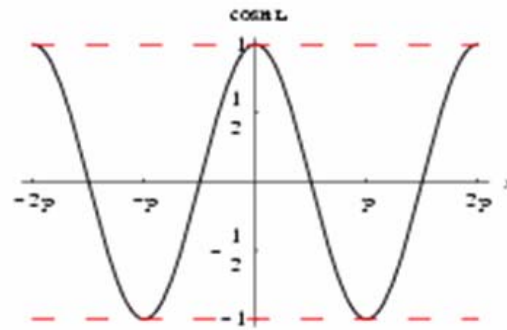
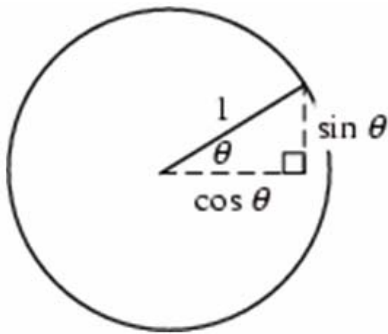
$$D_2 = (0.2, 0.7)$$

$$Q = (0.4, 0.8)$$

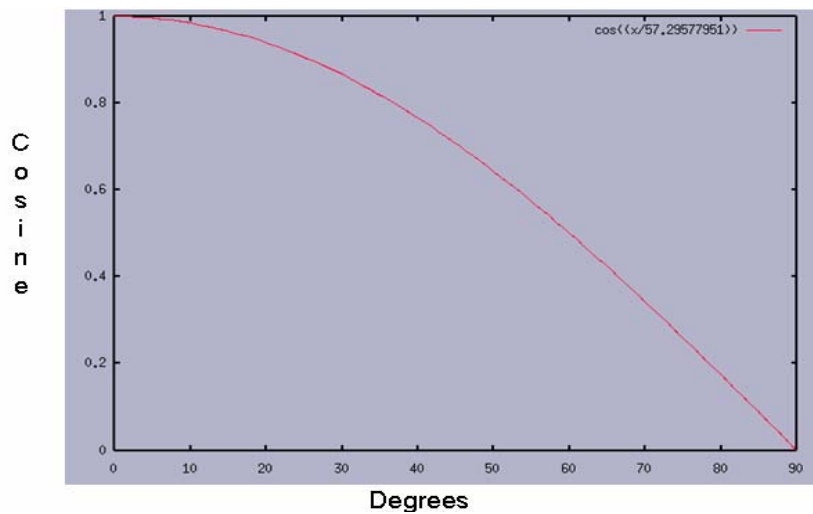
$$\cos \alpha_1 = 0.74$$

$$\cos \alpha_2 = 0.98$$

What's Cosine Anyway?



Cosine vs. Degrees



Computing a Similarity Score

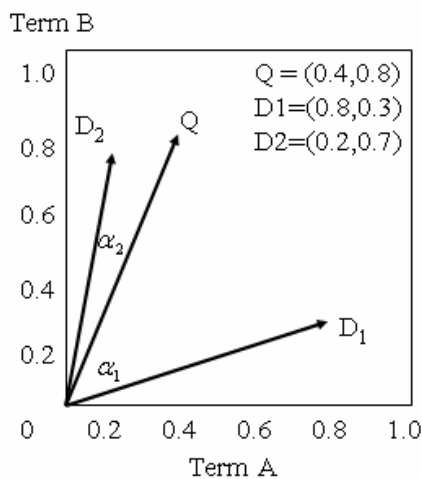
Say we have query vector $Q = (0.4, 0.8)$

Also, document $D_2 = (0.2, 0.7)$

What does their similarity comparison yield?

$$\begin{aligned} \text{sim}(Q, D_2) &= \frac{(0.4 \cdot 0.2) + (0.8 \cdot 0.7)}{\sqrt{[(0.4)^2 + (0.8)^2] \cdot [(0.2)^2 + (0.7)^2]}} \\ &= \frac{0.64}{\sqrt{0.42}} = 0.98 \end{aligned}$$

Vector Space Matching



$$\begin{aligned} D_1 &= (d_{i1}, w_{di1}; d_{i2}, w_{di2}; \dots; d_{it}, w_{dit}) \\ Q &= (q_{i1}, w_{qi1}; q_{i2}, w_{qi2}; \dots; q_{it}, w_{qit}) \end{aligned}$$

$$\text{sim}(Q, D_1) = \frac{\sum_{j=1}^t w_{qj} w_{dj}}{\sqrt{\sum_{j=1}^t (w_{qj})^2 \sum_{j=1}^t (w_{dj})^2}}$$

$$\begin{aligned} \text{sim}(Q, D_2) &= \frac{(0.4 \cdot 0.2) + (0.8 \cdot 0.7)}{\sqrt{[(0.4)^2 + (0.8)^2] \cdot [(0.2)^2 + (0.7)^2]}} \\ &= \frac{0.64}{\sqrt{0.42}} = 0.98 \end{aligned}$$

$$\text{sim}(Q, D_1) = \frac{.56}{\sqrt{0.58}} = 0.74$$