Vector. Representation of Docs.

Cosine Distance

Latent Semantic Analysis. (LSA) $\Rightarrow$
- Discover Latent topics
- Dimension Reduction.
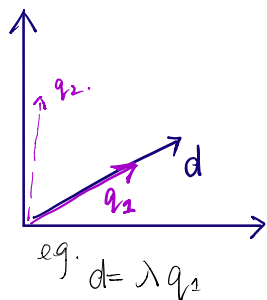- Represent words in terms of topics

---

Vector Notation.

vec(d) of d : V dimensional Vector    $W_{i(m),d} = f_{i(m)} \times IDF(i(m))$

**Document Vector :**

$(0, \cdots 0, W_{i(1)d}, 0 \cdots 0, W_{i(2)d}, \cdots \cdots W_{i(m),d}, 0, \cdots 0)$

- $D = \{d_1, \cdots d_N\}$   set of docs contains N documents

- No. of different words in. $D$ : **V** (vocab size).

- one document d : M terms $t_{i(1)}, \cdots \underline{t_{i(m)}}$
  $\hookrightarrow f_{i(m)}$ is the frequency

**if $d_1, d_2$ are docs**

- $vec(d_1) = vec(d_2) \Leftrightarrow d_1 = d_2$ ?

- $vec(d_1) = \underline{\lambda} vec(d_2)$
  same proportion.
  same words.

<u>**Document Length**</u>   $Len = \sqrt{\Sigma w^2}$



eg. $d = \lambda q_1$

$q_1 . d$ : Same words. Same proportion.

$q_2 . d$ : different words

$\Rightarrow$ greater angle : $vec(d)$ and $vec(q)$
less similar between q and d.

## Cosine Similarity

- Doc $d$ and query $q$.

$$CSim(q,d) = \cos\theta = \frac{vec(q) \cdot vec(d)}{\|q\| \cdot \|d\|} = \frac{\sum_{t \in q \cap d} W_{tq} \cdot W_{td}}{\|q\| \cdot \|d\|}$$

dimension $= N$

- Doc $d_1$ and $d_2$

$$CSim(d_1, d_2) = \cos\theta$$

$$= Sim(q,d)$$