





## Example Docs: Austen's Sense and Sensibility, Pride and Prejudice; Bronte's Wuthering Heights SaS PaP WН affection 115 58 20 7 11 jealous 10 0 gossip 2 6 WH SaS PaP affection 0,996 0,993 0,847 jealous 0,087 0,120 0,466 0,017 0,000 0,254 gossip cos(SAS, PAP) = .996 x .993 + .087 x .120 + .017 x 0.0 = 0.999 cos(SAS, WH) = .996 x .847 + .087 x .466 + .017 x .254 = 0.889 F. Aiolli - Sistemi Informativi 14 Dip. di Matematica Pura ed Applicata 2007/2008 Advantages of the VSM

- Flexibility. The most decisive factor in imposing VSM. The same intuitive geometric interpretation has been re-applied, apart from relevance feedback, in different contexts
  - Automatic document categorization
  - Automatic document filtering
  - Document clustering
  - Term-term similarity computation (terms are indexed by documents, dual)





 $rac{d_i q_j}{min(||d_i||^2,||q_j||^2)}$ 

## Conversion from a distance

Minkowsky Distances

$$L_p(x,z) = (\sum_{i=1}^{n} |x_i - z_i|^p)^{\frac{1}{p}}$$

When  $p = \infty$ ,  $L_{\infty} = \max_i(|x_i - z_i|)$ 

A similarity measure taking values in [0,1] can always be defined as

$$s_{p,\lambda}(x,z) = e^{-\lambda L_p(x,z)}$$

Where  $\lambda \in (0, +\infty)$  is a constant parameter

Dip. di Matematica Pura ed Applicata F. Aiolli - Sistemi Informativi 2007/2008 20

## Kernel functions

A kernel function K(x,z) is a (generally non-linear) function which corresponds to an inner product in some expanded feature space,

i.e.  $K(x,z) = \phi(x) \cdot \phi(z)$ 

Example: For 2-dimensional spaces  $x=(x_1,x_2)$ 

$$K(x,z) = (1 + x \cdot y)^2$$

is a kernel where

 $\phi(x) = (1, x_1^2, \sqrt{2}x_1x_2, x_2^2, \sqrt{2}x_1, \sqrt{2}x_2)$ 

Dip. di Matematica Pura ed Applicata F. Aiolli - Sistemi Informativi 2007/2008 21

