





## Ad-hoc IR: Terminologies

Terminologies:

- Query
  - Representative data of user's information need: text (default) and other media
- Document
  - Data candidate to satisfy user's information need: text (default) and other media
- Database|Collection|Corpus
  - A set of documents
- Corpora
  - A set of databases
  - Valuable corpora from TREC (Text Retrieval Evaluation Conference)

























Retrieval Models: Latent Semantic Indexing Query: Machine Learning Protein									
[		C1	C2	C3	C4	B1	B2	B3	
	information	1	1	0	0	0	0	0	
	retrieval	1	1	0	0	0	0	0	
	machine	1	1	1	1	0	0	0	
	learning	0	1	1	1	0	0	0	
	system	1	0	1	0	0	0	0	
	protein	0	0	0	0	0	1	1	
	gene	0	0	1	0	1	1	0	
	mutation	0	0	0	0	0	1	1	
	expression	0	0	0	0	1	0	1	
R	epresenta [0 0 1	tion of 1 1 0 1 0	he que <mark>0 0]<sup>⊤</sup></mark>	ry in the	e term v	vector s	pace:		









































































![](_page_26_Figure_0.jpeg)

![](_page_26_Figure_1.jpeg)

![](_page_27_Figure_0.jpeg)

![](_page_27_Figure_1.jpeg)

![](_page_28_Figure_0.jpeg)

![](_page_28_Figure_1.jpeg)

![](_page_29_Figure_0.jpeg)

![](_page_29_Figure_1.jpeg)

![](_page_30_Figure_0.jpeg)

![](_page_30_Picture_1.jpeg)

![](_page_31_Figure_0.jpeg)

![](_page_31_Figure_1.jpeg)

![](_page_32_Figure_0.jpeg)

![](_page_32_Figure_1.jpeg)

![](_page_33_Figure_0.jpeg)

Tex	t Catego	rization:	Evaluatio	on a					
Contingency Table Per Category (for all docs)									
	Truth: True	Truth: False							
Predicted Positive	а	b	a+b						
Predicted Negative	С	d	c+d						
	a+c	b+d	n=a+b+c+d						
a: number of th c: number of fa n: total numbe	ruly positive docs alse negative docs r of test documen	b: number of fal s d: number of tru ts	se-positive docs Ily-negative docs						

![](_page_34_Figure_0.jpeg)

![](_page_34_Figure_1.jpeg)

![](_page_35_Figure_0.jpeg)

Choices of Similarity FunctionsEuclidean distance
$$d(\vec{x}_1, \vec{x}_2) = \sqrt{\sum_v (x_{1v} - x_{2v})^2}$$
Kullback Leibler  
distance $d(\vec{x}_1, \vec{x}_2) = \sum_v x_{1v} \log \frac{x_{1v}}{x_{2v}}$ Dot product $\vec{x}_1 * \vec{x}_2 = \sum_v x_{1v} * x_{2v}$ Cosine Similarity $\cos(\vec{x}_1, \vec{x}_2) = \frac{\sum_v x_{1v} * x_{2v}}{\sqrt{\sum_v x_{1v}^2} \sqrt{\sum_v x_{2v}^2}}$ Kernel functions $k(\vec{x}_1, \vec{x}_2) = e^{-d(\vec{x}_1, \vec{x}_2)/2\sigma^2}$  (Gaussian Kernel)Automatic learning of the metrics

![](_page_36_Figure_0.jpeg)

![](_page_36_Picture_1.jpeg)

![](_page_37_Figure_0.jpeg)

![](_page_37_Figure_1.jpeg)

![](_page_38_Figure_0.jpeg)

![](_page_38_Figure_1.jpeg)

![](_page_39_Figure_0.jpeg)

![](_page_39_Figure_1.jpeg)

![](_page_40_Figure_0.jpeg)

![](_page_40_Figure_1.jpeg)

![](_page_41_Figure_0.jpeg)

![](_page_41_Figure_1.jpeg)

![](_page_42_Picture_0.jpeg)

## Linear SVM

•Set the derivative of the Lagrangian to be zero and calculate W by  $a_i$ , plug new form of w into the Lagrangian, the optimization problem can be written in terms of  $a_i$  (the dual problem)

$$\mathbf{w} = \sum_{i=1}^{n} \alpha_i y_i \mathbf{x}_i$$

Plug new form of w into the Lagrangian, the optimization problem can be written in terms of  $a_i$  (the dual problem)

max. 
$$W(\alpha) = \sum_{i=1}^{n} \alpha_i - \frac{1}{2} \sum_{\substack{i=1,j=1 \\ \mu}}^{n} \alpha_i \alpha_j y_i y_j \mathbf{x}_i^T \mathbf{x}_j$$
  
subject to  $\alpha_i \ge 0, \sum_{i=1}^{n} \alpha_i y_i = 0$ 

The above optimization problem is a quadratic program problem, which means there is a global maximum of  $a_i$  can always be found

## 

![](_page_43_Figure_0.jpeg)

![](_page_43_Figure_1.jpeg)

![](_page_44_Figure_0.jpeg)

![](_page_44_Figure_1.jpeg)

![](_page_45_Figure_0.jpeg)

![](_page_45_Figure_1.jpeg)

![](_page_46_Figure_0.jpeg)

![](_page_46_Figure_1.jpeg)

![](_page_47_Figure_0.jpeg)

![](_page_47_Figure_1.jpeg)

![](_page_48_Figure_0.jpeg)

![](_page_48_Figure_1.jpeg)

![](_page_49_Figure_0.jpeg)

![](_page_49_Figure_1.jpeg)

![](_page_50_Figure_0.jpeg)

![](_page_50_Figure_1.jpeg)

![](_page_51_Figure_0.jpeg)

![](_page_51_Figure_1.jpeg)

![](_page_52_Figure_0.jpeg)

![](_page_52_Picture_1.jpeg)

![](_page_53_Figure_0.jpeg)

![](_page_53_Figure_1.jpeg)

## Complete Link Agglomerative Clustering

• Use minimum similarity of pairs:

$$sim(c_i,c_j) = \min_{x \in c_i, y \in c_j} sim(x, y)$$

- Makes "tighter," spherical clusters that are typically preferable.
- After merging  $c_i$  and  $c_j$ , the similarity of the resulting cluster to another cluster,  $c_k$ , is:

 $sim((c_i \cup c_j), c_k) = min(sim(c_i, c_k), sim(c_j, c_k))$ 

![](_page_54_Figure_6.jpeg)

![](_page_55_Figure_0.jpeg)

![](_page_55_Figure_1.jpeg)

![](_page_56_Figure_0.jpeg)

![](_page_56_Picture_1.jpeg)

![](_page_57_Figure_0.jpeg)

![](_page_57_Figure_1.jpeg)

![](_page_58_Figure_0.jpeg)

![](_page_58_Figure_1.jpeg)

![](_page_59_Figure_0.jpeg)

![](_page_59_Figure_1.jpeg)