<u>Generalization of Perception</u> <u>convergence Algorithm (Generalized</u> <u>PCA)</u>

- use of preference
- multi-level relevance

Let documents and queries be represented as n-dimensional vectors. D is the set of all document representations. Let <- represent a user preference relation over D.

Defn. 1: A preference relation, \leq is <u>linear</u> if there <u>exists</u> a query vector <u>q</u> such that for any <u>d</u>, <u>d</u>¹ \in D, <u>d</u> $\leq \cdot \underline{d}^1 \Leftrightarrow \underline{dq}^T \leq \underline{d}^1 \underline{q}^T$ If an algorithm is available to find such a <u>q</u>, then we can use <u>q</u> to achieve <u>perfect</u> <u>ranking</u>.

Perfect ranking is difficult to achieve and <u>not necessary</u>.

Defn. 2: A preference relation, $<\cdot$ is <u>weakly linear</u> if there <u>exists</u> a query vector <u>q</u> such that for any <u>d</u>, <u>d</u>¹ \in D, <u>d</u> $<\cdot$ <u>d</u>¹ => <u>dq</u>^T < <u>d</u>¹<u>q</u>^T

-- This condition implies $\underline{dq}^{T} < \underline{d}^{1}\underline{q}^{T} => not (\underline{d} < \cdot \underline{d}^{1})$

TWO CLASS RELEVANCE ASSUMED IN THE FIGURES BELOW





--That is, less preferable documents will not be ranked higher or equal to more preferable documents. But equally preferred documents may have different RSVs.

--If an algorithm is available to find such a \underline{q} , then \underline{q} can be used to achieve <u>acceptable ranking</u>.

APPROACH

Suppose $\underline{d} < \cdot \underline{d}^1$ holds for \underline{d} , \underline{d}^1 in the training set. Let training set be such that user preference defined by $< \cdot$ is weakly linear. By definition, $\underline{d} < \cdot \underline{d}^1 => \underline{d} \ \underline{q}^T < \underline{d}^1 \underline{q}^T$

OR

$$\underline{d} < \cdot \underline{d}^1 = \ge \underline{b} \underline{q}^T > 0,$$

where $\underline{b} = \underline{d}^1 - \underline{d}$

Thus, we need to find \underline{q} such that $\underline{b}_{\alpha} \underline{q}^{T} > 0$ for all α , such that \underline{b}_{α} is difference between two vectors in the training set and first of the two is more preferred.

<u>Generalized Perception Criterion</u> $J(\underline{q}) = \sum -\underline{b}q^{T},$ $\underline{b} \in y(\underline{q}),$ Where $y(\underline{q}) = \{\underline{b} = \underline{d}^{1} - \underline{d} \mid \underline{d} < \cdot \underline{d}^{1} \text{ and } \underline{b}q^{T} \le 0\}$

<u>Generalized Perception Conv.</u> <u>Algorithm</u>

1. choose \underline{q}^0 , let k=0 2. let q^k be the query vector in iteration k. compute $y(\underline{q}^k)$ If $y(\underline{q}^k) = \phi$, then Output \underline{q}^k and terminate 3. Let $q^{k+1} = \underline{q}^k + \underline{\rho}_k$ Σ $b \in y(q^k)$ 4. Let k=k+1, go to step 2

$$\underline{\text{Example}} \\
D = \{\underline{d}_1, \underline{d}_2, \underline{d}_3, \underline{d}_4\} \\
\underline{d}_1 = (1, 1, 0, 1) \\
\underline{d}_2 = (1, 0, 1, 0) \\
\underline{d}_3 = (0, 1, 1, 0) \\
\underline{d}_4 = (0, 1, 0, 1) \\
\text{Let preference relation be,} \\
\underline{d}_1 < \underline{d}_2, \ \underline{d}_4 < \underline{d}_2, \ \underline{d}_1 < \underline{d}_3, \ \underline{d}_2 < \underline{d}_3, \\
\underline{d}_4 < \underline{d}_3,$$



$$g^{0} = (0,0,0,0)$$

y (g⁰) = y
g^{1}=q^{0} + \sum \underline{b}_{i}
b_i \in y(q⁰)
= (-1,-1,4,-4)
y (q¹) = {b_{4}}
q^{2} = q^{1} + \underline{b}_{4}
= (-2,0,4,-4)
y (q²) = ϕ , terminate

Theorem: Generalized PCA. Will terminate if preference relation for documents in the training set is <u>weakly linear</u>.

Example: <u>Learning by sample (Gen-</u> PCA) K=0 $\underline{q}^0 = (0,0,0,0)$ b_1 arrives $\rightarrow b_1$ is misclassified $q^1 = b_1$ <u>b</u>₂ arrives \rightarrow check <u>b</u>₂q^{1T} Not misclassified $q^2 = \underline{q}^1$ <u>b</u>₃ arrives \rightarrow check b₃q^{2T} Not misclassified $q^{3} = q^{2}$

Rocchio defined: an ideal (our definition for the same is "acceptable ranking) query q:

 $\rho(\underline{q}, \underline{d}) < \rho(\underline{q}, \underline{d}')$ for all $\underline{d} \in \text{nrel}$ $\underline{d}' \in \text{rel}$

Didn't know how to get acceptable ranking, So suggested an approximate solution. He finds the optimal query vector by maximizing the following criterion function with respect to q.

$$C = \frac{1}{n_0} \sum_{\underline{d'} \in rel} \rho(\underline{q}, \underline{d'}) - \frac{1}{n_1} \sum_{\underline{d} \in nrel} \rho(\underline{q}, \underline{d})$$

where n_0 . is: # of rel doc. in the set of rel.

n₁. is: # of nrel doc. in the set of nrel. By the definition of ρ (q, d) $C = \frac{q}{\|\underline{q}\|} [\frac{1}{n_0} \sum_{\underline{d'} \in rel} \frac{\underline{d'}}{\|\underline{d'}\|} - \frac{1}{n_1} \sum_{\underline{d} \in nrel} \frac{\underline{d}}{\|\underline{d}\|}]$

C is maximized if q is chosen to be: $\underline{q}^* = k \left[\frac{1}{n_0} \sum_{\underline{d'} \in rel} \frac{\underline{d'}}{\|\underline{d'}\|} - \frac{1}{n_1} \sum_{\underline{d} \in nrel} \frac{\underline{d}}{\|\underline{d}\|} \right]$

because C= $\underline{q}^* \bullet A = kA \bullet A$, where $A = \frac{1}{n_0} \sum_{\underline{d'} \in rel} \frac{\underline{d'}}{\|\underline{d'}\|} - \frac{1}{n_1} \sum_{\underline{d} \in nrel} \frac{\underline{d}}{\|\underline{d}\|}$ and k is a

positive constant. (i.e., the direction of optimal query coincides with direction of A).

The query q* is referred to as THE optimal query by Rocchio

Problem

- The optimal query is not necessarily a solution vector that gives acceptable ranking even if the preference relation <• is weakly linear.
- 2. The criterion function C has a maximum only if the query vector q is normalized.

3. This method is only defined for two-level relevance.

Ex. Same example we used for PCA. $\{d_1, d_4\}$ is rel and $\{d_2, d_3\}$ is nrel.

k≠0

$$q^* = k \left[\frac{1}{2} \left(\frac{\underline{d}_2}{\|\underline{d}_2\|} + \frac{\underline{d}_3}{\|\underline{d}_3\|} \right) - \frac{1}{2} \left(\frac{\underline{d}_1}{\|\underline{d}_1\|} + \frac{\underline{d}_4}{\|\underline{d}_4\|} \right) \right]$$
$$= \frac{1}{2} \left(\frac{(1,1,0,1)}{\sqrt{3}} + \frac{(0,1,0,1)}{\sqrt{2}} \right) - \frac{1}{2} \left(\frac{(1,0,1,0)}{\sqrt{2}} + \frac{(0,1,1,0)}{\sqrt{2}} \right)$$
$$= \frac{1}{2} \left(\frac{1}{\sqrt{3}} - \frac{1}{\sqrt{2}}, \frac{1}{\sqrt{3}} + \frac{1}{\sqrt{2}} - \frac{1}{\sqrt{2}}, \frac{-2}{\sqrt{2}}, \frac{1}{\sqrt{3}} + \frac{1}{\sqrt{2}} \right)$$
$$= \frac{1}{2\sqrt{6}} \left(\sqrt{2} - \sqrt{3}, \sqrt{2}, -2\sqrt{3}, \sqrt{2} + \sqrt{3} \right)$$

From the Generalized PCA:

Choose $q^0 = 0$, at the first iteration,

$$y(\underline{q}^{0}) = \{ \underline{b} = (\frac{\underline{d}'}{\|\underline{d}'\|} + \frac{\underline{d}}{\|\underline{d}\|} | \underline{d} \in nrel, \underline{d}' \in rel) \}$$

$$q^{1} = q^{0} + \sum_{\underline{b} \in y(\underline{q}^{0})} \underline{b} = \sum_{\underline{d}' \in rel} \sum_{\underline{d} \in rel} (\frac{\underline{d}'}{\|\underline{d}'\|} - \frac{\underline{d}}{\|\underline{d}\|}) = n_{0}n_{1} [\frac{1}{n_{0}} \sum_{\underline{d}' \in rel} \frac{\underline{d}'}{\|\underline{d}'\|} - \frac{1}{n_{1}} \sum_{\underline{d} \in nrel} \frac{\underline{d}}{\|\underline{d}\|}] \approx q^{*}$$

Rocchio's method is the result obtained after first iteration of generalized PCA. (Here we assume: (i) training vectors are normalized (ii) the user is giving feedback using 2-level relevance)

Cycle Vs. Iteration (of Relevance Feedback) Training Sample. I:

