

---

# Homework

---

- Homework 1: Last chance to ask questions in person!
- Homework 2: Prepare pilot results for your final project.
  - Means you need a final project topic.
  - If you don't have one, come talk to us.
  - If you do have one, be sure to one of us approves.

---

## Relevance feedback (RF)

---

- Basic procedure:
  - User issues query.
  - System returns initial results.
  - User marks documents as relevant or irrelevant.
  - System improves representation of information need based on user's feedback.
  - System returns revised set of results.
- One of more iterations of this process.

---













## Motivation for RF

---













- Hard to formulate a query that retrieves exactly the documents you want.
  - Especially true if you don't know the collection well or are an inexperienced user of the system.
- Easy to look at documents and decide whether or not they match your query.
  - Even if you can't articulate what you're looking for, you'll know it when you see it.
- Also: Can help address some of the issues associated with different models of information seeking.

# Example RF: Before

[Browse](#) [Search](#) [Prev](#) [Next](#) [Random](#)

					
(144473, 16458)	(144457, 252140)	(144456, 262857)	(144456, 262863)	(144457, 252134)	(144483, 265154)
0.0	0.0	0.0	0.0	0.0	0.0
0.0	0.0	0.0	0.0	0.0	0.0
0.0	0.0	0.0	0.0	0.0	0.0
					
(144483, 264644)	(144483, 265153)	(144518, 257752)	(144538, 525937)	(144456, 249611)	(144456, 250064)
0.0	0.0	0.0	0.0	0.0	0.0
0.0	0.0	0.0	0.0	0.0	0.0
0.0	0.0	0.0	0.0	0.0	0.0

# Example RF: After

<a href="#">Browse</a> <a href="#">Search</a> <a href="#">Prev</a> <a href="#">Next</a> <a href="#">Random</a>					
					
(144538, 523493) 0.54182 0.231944 0.309876	(144538, 523835) 0.56319296 0.267304 0.295889	(144538, 523529) 0.584279 0.280881 0.303398	(144456, 253569) 0.64501 0.351395 0.293615	(144456, 253568) 0.650275 0.411745 0.23853	(144538, 523799) 0.66709197 0.358033 0.309059
					
(144473, 16249) 0.6721 0.393922 0.278178	(144456, 249634) 0.675018 0.4639 0.211118	(144456, 253693) 0.676901 0.47645 0.200451	(144473, 16328) 0.700339 0.309002 0.391337	(144483, 265264) 0.70170796 0.36176 0.339948	(144478, 512410) 0.70297 0.469111 0.233859



---

## Rocchio algorithm: Underlying theory

---

- Incorporates relevance feedback into vector space model.
- Find the *optimal query vector*, i.e., the one that maximizes similarity to relevant documents while minimizing similarity to non-relevant documents.

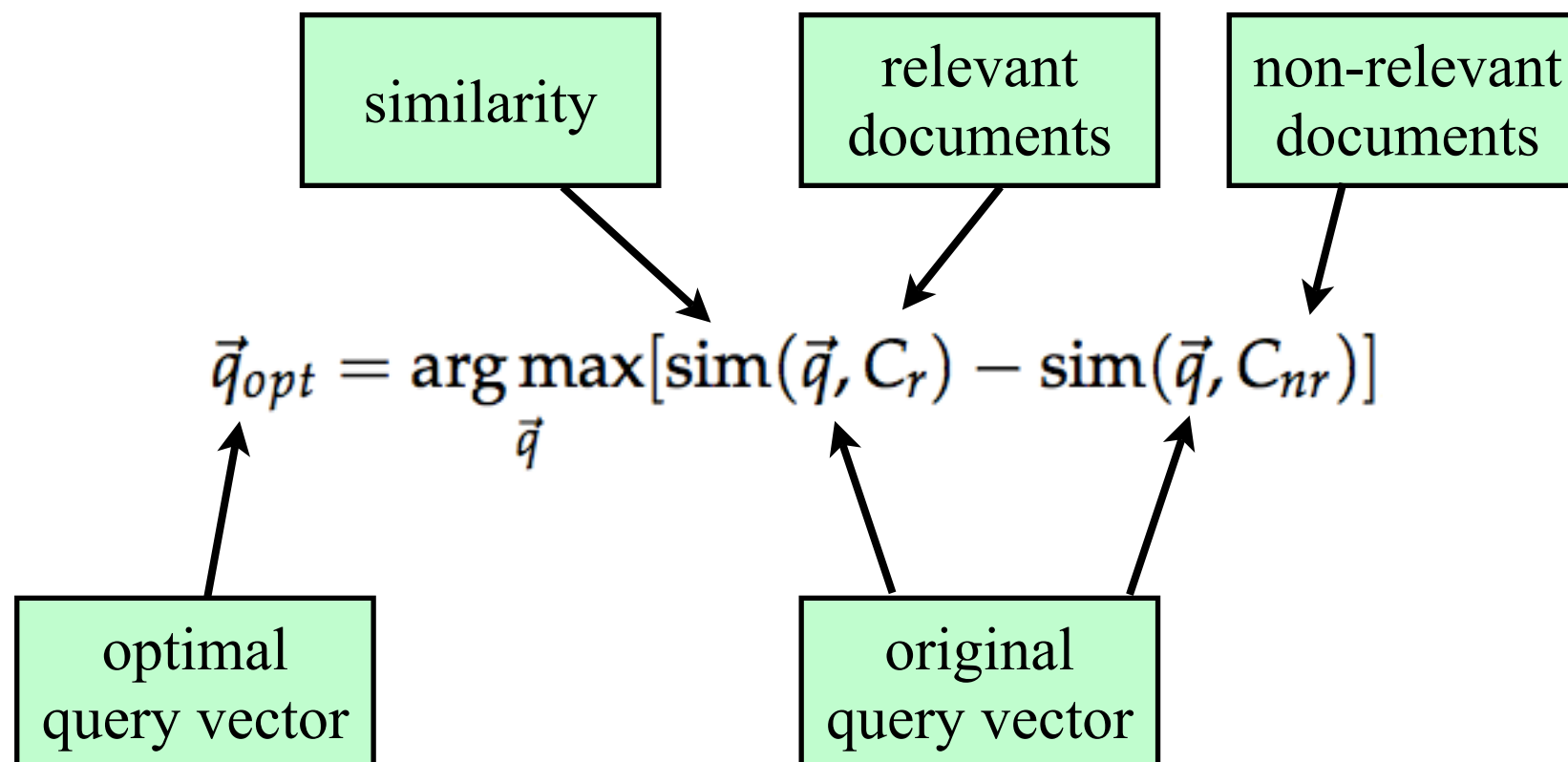
$$\vec{q}_{opt} = \arg \max_{\vec{q}} [\text{sim}(\vec{q}, C_r) - \text{sim}(\vec{q}, C_{nr})]$$

---

## Rocchio algorithm: Underlying theory

---

- Incorporates relevance feedback into vector space model.
- Find the *optimal query vector*, i.e., the one that maximizes similarity to relevant documents while minimizing similarity to non-relevant documents.



---

## Rocchio algorithm: Underlying theory

---

$$\vec{q}_{opt} = \arg \max_{\vec{q}} [\text{sim}(\vec{q}, C_r) - \text{sim}(\vec{q}, C_{nr})]$$

$$\vec{q}_{opt} = \frac{1}{|C_r|} \sum_{\vec{d}_j \in C_r} \vec{d}_j - \frac{1}{|C_{nr}|} \sum_{\vec{d}_j \in C_{nr}} \vec{d}_j$$



---

## Rocchio algorithm: Underlying theory

---

$$\vec{q}_{opt} = \arg \max_{\vec{q}} [\text{sim}(\vec{q}, C_r) - \text{sim}(\vec{q}, C_{nr})]$$

$$\vec{q}_{opt} = \frac{1}{|C_r|} \sum_{\vec{d}_j \in C_r} \vec{d}_j - \frac{1}{|C_{nr}|} \sum_{\vec{d}_j \in C_{nr}} \vec{d}_j$$

sum of all the  
relevant document  
vectors

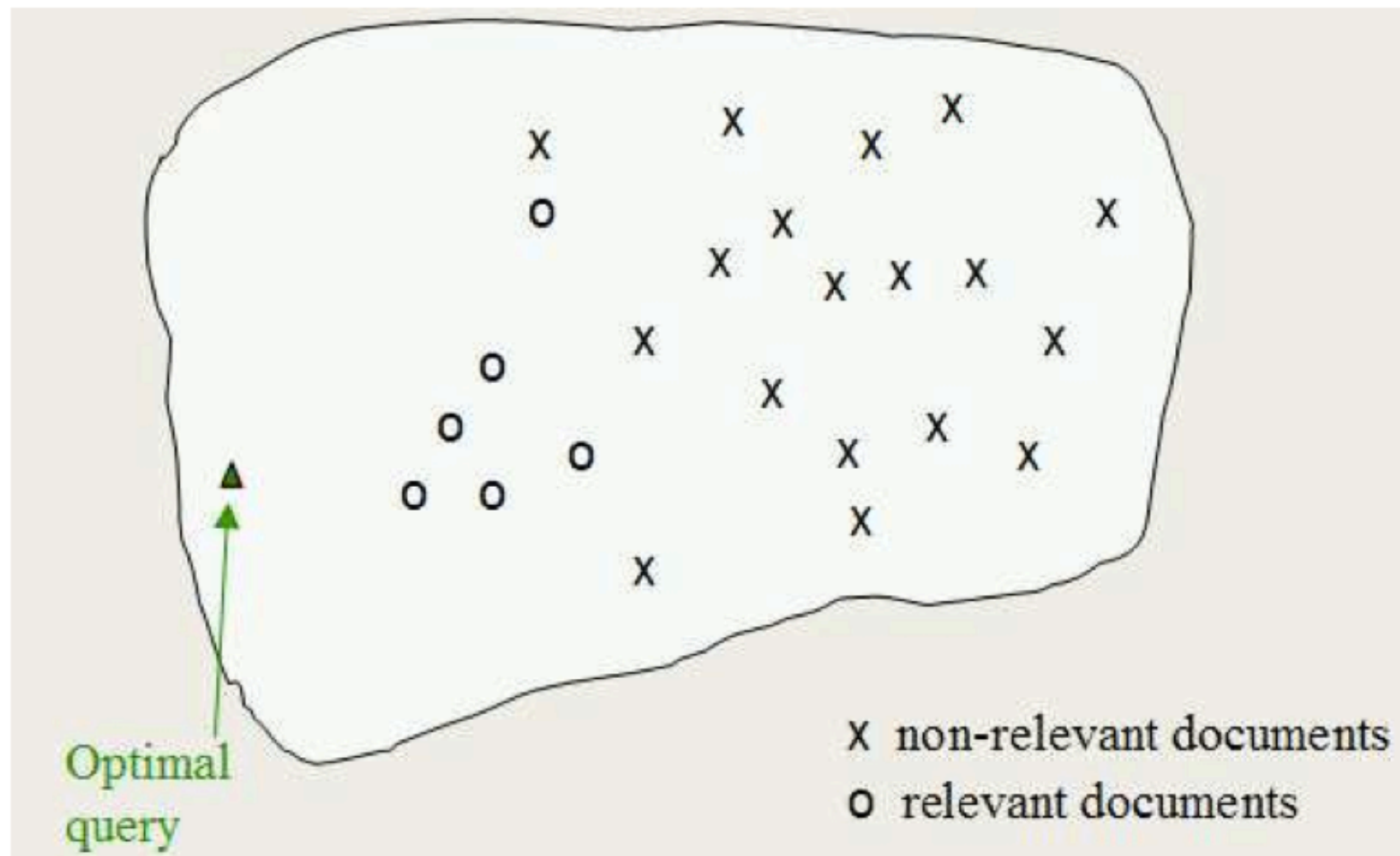
sum of all the non-  
relevant document  
vectors

- ➔ Optimal query is vector difference between the centroids of the relevant and non-relevant documents.

---

# Rocchio algorithm: Underlying theory

---



➡ Problem: We don't know the *full* set of relevant documents!

---

# Rocchio algorithm

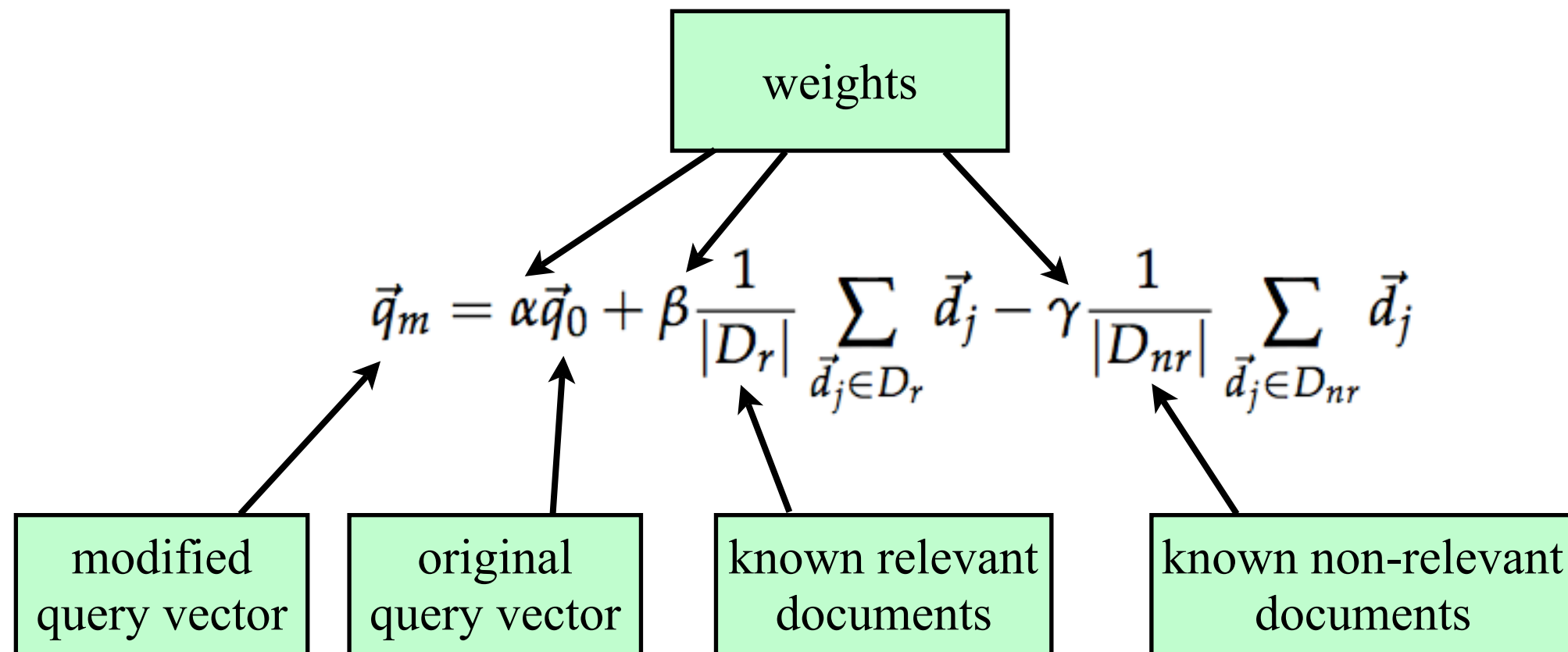
---

$$\vec{q}_m = \alpha \vec{q}_0 + \beta \frac{1}{|D_r|} \sum_{\vec{d}_j \in D_r} \vec{d}_j - \gamma \frac{1}{|D_{nr}|} \sum_{\vec{d}_j \in D_{nr}} \vec{d}_j$$

---

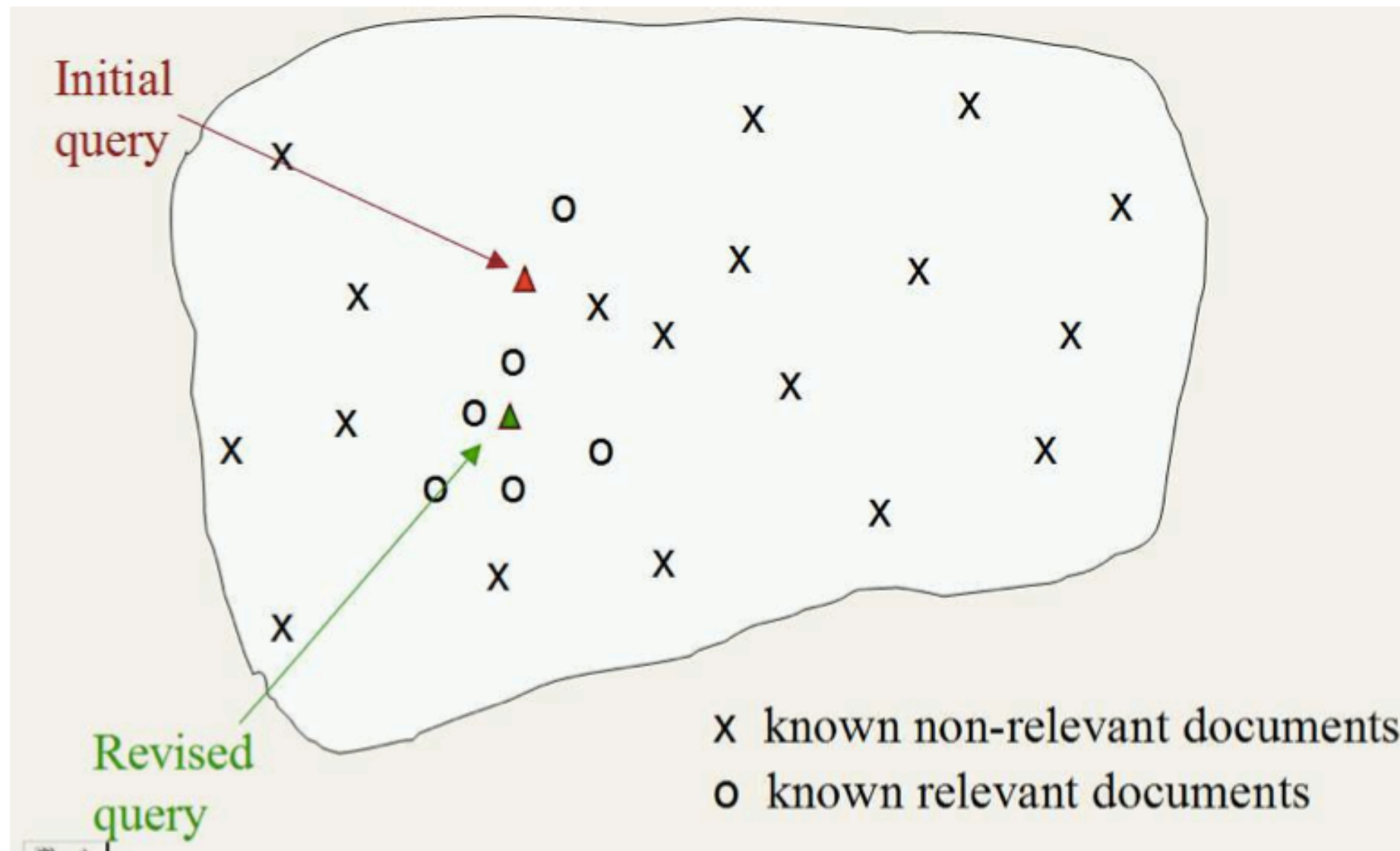
## Rocchio algorithm

---



- Weights control balance between trusting the original query and trusting the judged document set.
  - Fewer judged documents, more weight for alpha.
  - More judged documents, more weight for beta and gamma.

# Rocchio algorithm: Picture example



---

## Rocchio algorithm: Number example

---

new query vector =  $\alpha \cdot$  original query vector +  
 $\beta \cdot$  relevant document vectors -  
 $\gamma \cdot$  non-relevant document vectors

0	4	0	8	0	0
---	---	---	---	---	---

2	4	8	0	0	2
---	---	---	---	---	---

8	0	4	4	0	16
---	---	---	---	---	----



---

## Rocchio algorithm: Number example

---

new query vector =  $\alpha \cdot$  original query vector +  
 $\beta \cdot$  relevant document vectors -  
 $\gamma \cdot$  non-relevant document vectors

0	4	0	8	0	0
---	---	---	---	---	---

 $\alpha = 1$ 

2	4	8	0	0	2
---	---	---	---	---	---

 $\beta = 0.5$ 

8	0	4	4	0	16
---	---	---	---	---	----

 $\gamma = 0.25$ 

Typically  $\beta > \gamma$ ,  
since positive  
feedback is more  
meaningful.

---

## Rocchio algorithm: Number example

---

new query vector =  $\alpha \cdot$  original query vector +  
 $\beta \cdot$  relevant document vectors -  
 $\gamma \cdot$  non-relevant document vectors

<table><tr><td>0</td><td>4</td><td>0</td><td>8</td><td>0</td><td>0</td></tr></table>	0	4	0	8	0	0	$\alpha = 1$		<table><tr><td>0</td><td>4</td><td>0</td><td>8</td><td>0</td><td>0</td></tr></table>	0	4	0	8	0	0
0	4	0	8	0	0										
0	4	0	8	0	0										
<table><tr><td>2</td><td>4</td><td>8</td><td>0</td><td>0</td><td>2</td></tr></table>	2	4	8	0	0	2	$\beta = 0.5$	+	<table><tr><td>1</td><td>2</td><td>4</td><td>0</td><td>0</td><td>1</td></tr></table>	1	2	4	0	0	1
2	4	8	0	0	2										
1	2	4	0	0	1										
<table><tr><td>8</td><td>0</td><td>4</td><td>4</td><td>0</td><td>16</td></tr></table>	8	0	4	4	0	16	$\gamma = 0.25$	-	<table><tr><td>2</td><td>0</td><td>1</td><td>1</td><td>0</td><td>4</td></tr></table>	2	0	1	1	0	4
8	0	4	4	0	16										
2	0	1	1	0	4										

---

Typically  $\beta > \gamma$ ,  
since positive  
feedback is more  
meaningful.

---

## Rocchio algorithm: Number example

---

new query vector =  $\alpha \cdot$  original query vector +  
 $\beta \cdot$  relevant document vectors -  
 $\gamma \cdot$  non-relevant document vectors

<table><tr><td>0</td><td>4</td><td>0</td><td>8</td><td>0</td><td>0</td></tr></table>	0	4	0	8	0	0	$\alpha = 1$		<table><tr><td>0</td><td>4</td><td>0</td><td>8</td><td>0</td><td>0</td></tr></table>	0	4	0	8	0	0
0	4	0	8	0	0										
0	4	0	8	0	0										
<table><tr><td>2</td><td>4</td><td>8</td><td>0</td><td>0</td><td>2</td></tr></table>	2	4	8	0	0	2	$\beta = 0.5$	+	<table><tr><td>1</td><td>2</td><td>4</td><td>0</td><td>0</td><td>1</td></tr></table>	1	2	4	0	0	1
2	4	8	0	0	2										
1	2	4	0	0	1										
<table><tr><td>8</td><td>0</td><td>4</td><td>4</td><td>0</td><td>16</td></tr></table>	8	0	4	4	0	16	$\gamma = 0.25$	-	<table><tr><td>2</td><td>0</td><td>1</td><td>1</td><td>0</td><td>4</td></tr></table>	2	0	1	1	0	4
8	0	4	4	0	16										
2	0	1	1	0	4										
			<hr/>												
			<table><tr><td>-1</td><td>6</td><td>3</td><td>7</td><td>0</td><td>-3</td></tr></table>	-1	6	3	7	0	-3						
-1	6	3	7	0	-3										
			<table><tr><td>0</td><td>6</td><td>3</td><td>7</td><td>0</td><td>0</td></tr></table>	0	6	3	7	0	0						
0	6	3	7	0	0										

Typically  $\beta > \gamma$ ,  
since positive  
feedback is more  
meaningful.

Negative term  
weights become 0.

---

## Probabilistic RF

---

- Alternative to reweighting the query a vector space: build a *classifier* using the user relevance feedback.
- Using a Naive Bayes classifier:

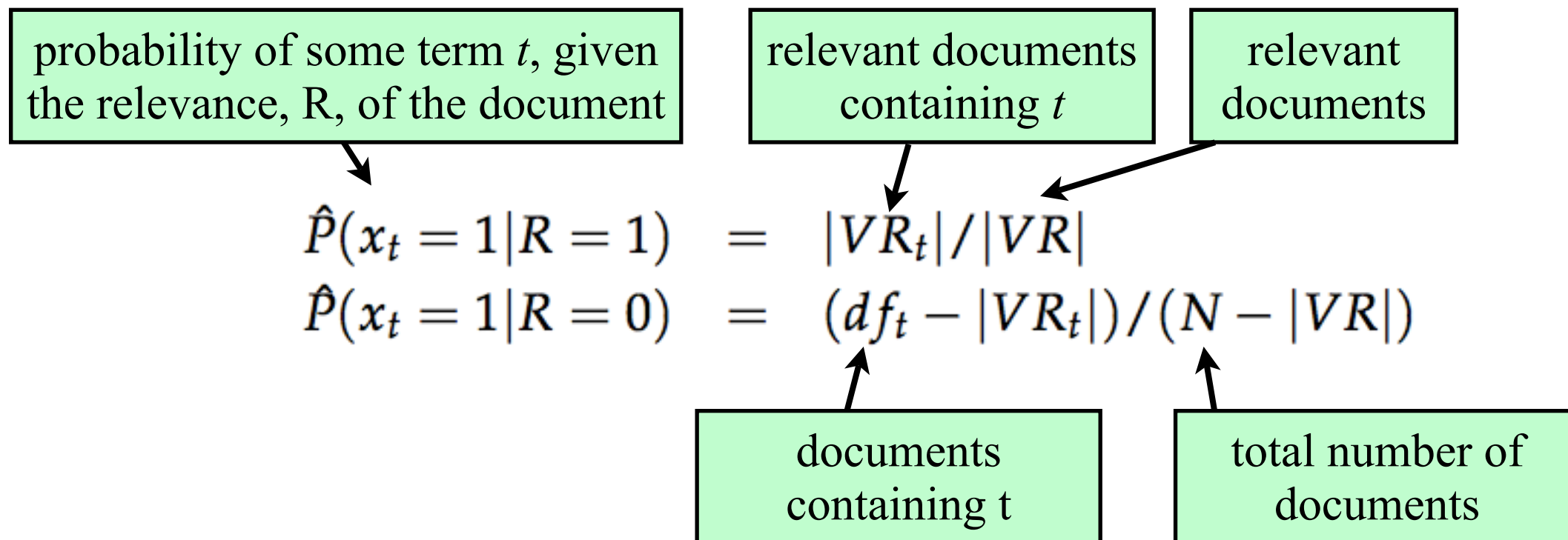
$$\begin{aligned}\hat{P}(x_t = 1 | R = 1) &= |VR_t| / |VR| \\ \hat{P}(x_t = 1 | R = 0) &= (df_t - |VR_t|) / (N - |VR|)\end{aligned}$$

---

## Probabilistic RF

---

- Alternative to reweighting the query a vector space: build a *classifier* using the user relevance feedback.
- Using a Naive Bayes classifier:



- A little more on this later, but do note that all information about the original query is lost in this formulation.

---

## Evaluation

---

- Compare precision and recall of results retrieved using  $q_0$  to those retrieved with  $q_m$ .
  - Over all documents: works great but cheating.
  - Over residual documents (excludes those assessed by user): not cheating but performance degrades, especially when there are few relevant documents.
  - Use two sets of documents: train and test.
- Other evaluation measures:
  - how does RF compare to user query reformulation?
  - how many documents found per time unit?
- Empirically, one round of RF is useful: diminishing returns on multiple iterations.



---

## Problems with RF

---

- RF doesn't work well:
  - when query contains misspellings (e.g., *Mitt Rmoney*)
  - in cross-language IR scenarios
  - in document sets with vocabulary mismatch (e.g., *myocardial infarction* vs. *heart attack*)
  - when documents don't actually cluster very well
- Other issues:
  - High computing cost: queries can end up being huge.
  - Users don't like giving explicit feedback: too much work.
- Alternatives to explicit user RF: coming up next.