

Setup:

- $N = 1,000,000$ docs in corpus
- $R = 8$ relevant docs
- System retrieves 9 rel. docs in the top 10 results
- What is performance of set retrieved?
- Standard ML measure: accuracy

acc = frac. of instances correctly predicted

error = $1 - \text{acc} = \dots$ incorrect

$$\text{acc} = \frac{\# \checkmark}{N} = \frac{4 + 999986}{1,000,000} = 99.999\%$$

$$\text{em} = \frac{\# X}{N} = \frac{6 + 4}{1,000,000} = \frac{10}{1,000,000} = 0.001\%$$

LIST

<u>Rank</u>	<u>Rel.</u>
1	R
2	N
3	N
4	R
5	N
6	R
7	N
8	N
9	R
10	N

<u>Correct?</u>	
✓	
X	
X	
✓	4 ✓
X	
✓	
X	6 X
X	
✓	
X	
✓	
X	
X	4 X
✓	
X	999986 ✓
✓	
X	
✓	
X	
X	

Search engine returns nothing!

$$\text{em} = \frac{8}{1,000,000} = 0.0008\%$$

How to measure performance w/ massive data imbalance?

- ⇒ ① Better set-level metrics
② Ranking based metrics

Set-level metrics: precision, recall, F1

precision = frac. of ret. docs that are rel.

$$= \frac{|rel. \cap ret.}|}{|ret.}| = \frac{4}{10} = 0.4$$

recall = frac. of rel. docs that are ret.

$$= \frac{|rel. \cap ret.}|}{|rel.}| = \frac{4}{8} = 0.5$$

How to combine prec. & recall? maybe just average = $\frac{0.4 + 0.5}{2} = 0.45$

But easy to game in IR

① return just the top rel. doc

$$prec = \frac{1}{1} = 1$$

$$rec = \frac{1}{8} = 0.125$$

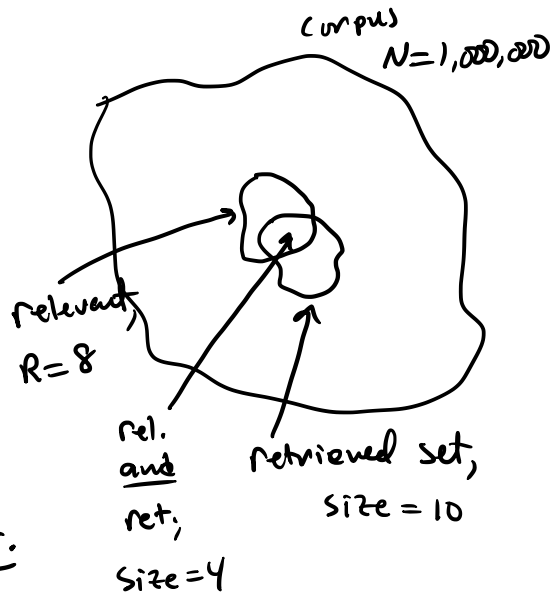
$$\frac{p+r}{2} = \frac{1+0.125}{2} = 0.5625$$

② return everything

$$prec = \frac{8}{1,000,000} = 0.000008 \approx 0$$

$$rec = \frac{8}{8} = 1$$

$$\frac{p+r}{2} \approx \frac{0+1}{2} = 0.5$$



ML: confusion matrix
 (binary classif)
~~true~~ / ground truth
 P / N

Predicted	P	TP	FP
	N	FN	TN

relevant set $\Rightarrow P = TP + FN$

retrieved set = predicted positives

$P = TP + FP$

TP precision = fraction of retrieved

$$\frac{|\text{retrieved} \cap \text{relevant}|}{|\text{retrieved}|} = \frac{TP}{TP + FP}$$

TP recall = fraction of all good

$$\frac{|\text{retrieved} \cap \text{relevant}|}{|\text{relevant}|} = \frac{TP}{TP + FN}$$



$$\text{ML accuracy} = \frac{\text{TP} + \text{TN}}{\text{TP} + \text{FP} + \text{TN} + \text{FN}}$$

accuracy not appropriate when $\text{TN} + \text{FP}$ is very large
(data skewed to negative side)

How to combine #'s when you want to penalize for any # small?

Digression: Other kinds of means?

arithmetic mean: $\frac{x_1 + x_2}{2}$ or $\frac{x_1 + x_2 + \dots + x_n}{n}$ - straight average

geometric mean: $\sqrt{x_1 \cdot x_2}$ or $\sqrt[n]{x_1 \cdot x_2 \cdot x_3 \cdot \dots \cdot x_n}$

harmonic mean: $\frac{1}{(\frac{1}{x_1} + \frac{1}{x_2})/2} = \frac{2}{\frac{1}{x_1} + \frac{1}{x_2}} = \frac{2x_1x_2}{x_1 + x_2}$ or $\frac{n}{\sum_{i=1}^n \frac{1}{x_i}}$

Then: h.m. \leq g.m. \leq a.m. w/equality iff all #'s same

In I.R., we use the h.m., which we call F1

$$F1 = \text{ham. mean}(\text{prec}, \text{rec}) = \frac{2}{\frac{1}{\text{prec}} + \frac{1}{\text{rec}}} = \frac{2 \cdot \text{prec} \cdot \text{rec}}{\text{prec} + \text{rec}}$$

	p	r	a.m.	F1
<u>Example</u>	0.4	0.5	0.45	0.444
	1	0.125	0.5625	0.222
	0.000008	1	0.500004	0.000016

weighted F_β

- What about **ranked retrieval evaluation**?
- ROC curves → used widely in ML & data analysis
 - ↳ receiver operator curve

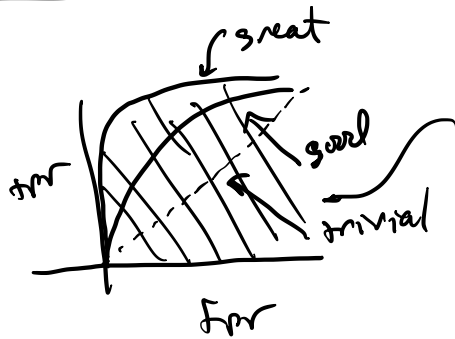
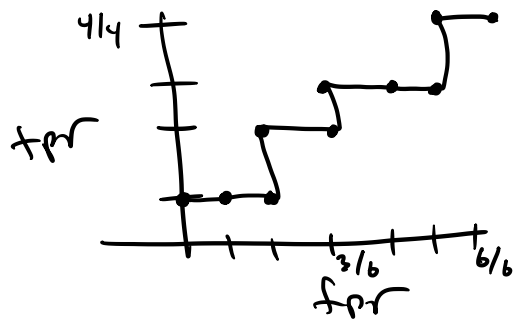
original setup
R=8 N=10

- true positive rate
- false positive rate

universe of just 10 items

		<u>tp</u>	<u>fp</u>
1	R	1/4	0/6
2	N	1/4	1/6
3	N	1/4	2/6
4	R	2/4	2/6
5	N	2/4	3/6
6	R	3/4	3/6
7	N	3/4	4/6
8	N	3/4	5/6
9	R	4/4	5/6
10	N	4/4	6/6

<u>tp</u>	<u>fp</u>
1/8	0/999,992
1/8	1/999,992
1/8	2/999,992
1/4	⋮
⋮	⋮
⋮	⋮
⋮	⋮
⋮	⋮
⋮	⋮
⋮	⋮
⋮	⋮



area under curve
→ AUC

↑
trivial
ROC
curve

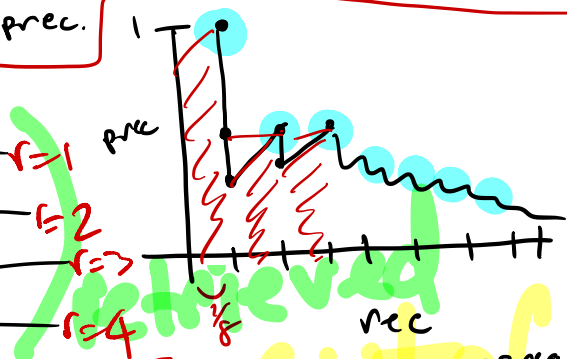
IR: Ranked Retrieval Metrics

- prec-rec curves
- avg. prec.
- R-prec.

$N=1,000,000$
 $R=8$

Rank	Q _{rel}	rel.	prec	rec
1	R	1/1	1/1	1/8
2	N	1/2	1/2	1/8
3	N	1/3	1/3	1/8
4	R	2/4	2/4	2/8
5	N	2/5	2/5	2/8
6	R	3/6	3/6	3/8
7	N	3/7	3/7	3/8
8	N	3/8	3/8	3/8
9	R	4/9	4/9	4/8
10	N	4/10	4/10	4/8

⋮	⋮	⋮	⋮	⋮
⋮	R	~0	~0	5/8
⋮	⋮	⋮	⋮	⋮
⋮	R	~0	~0	6/8
⋮	⋮	⋮	⋮	⋮
⋮	R	~0	~0	7/8
⋮	⋮	⋮	⋮	⋮
⋮	R	~0	~0	8/8



1-query
 $\text{avg. prec.} = \text{avg. of prec. at each rel. doc}$

$$= \frac{1/1 + 2/4 + 3/6 + 4/9 + \dots + 0}{8} = 0.3524$$

1-query
 $R\text{-prec} = (\text{prec} = \text{rec}) = 3/8 = 0.375$

multiple queries HW1: 25 queries

MAP = mean-avg-precision

$$\frac{1}{|Q|} \sum_{q \in Q} AP(q)$$