

MLRF Lecture 04

J. Chazalon, LRE/EPITA, 2025

IR evaluation

Lecture 04 part 03

How to evaluate a retrieval system?

We need a set of queries for which we know the expected results
“Ground truth”, aka “targets”, “gold standard”...

To compare 2 methods, we need to use the same database and the same queries.

Many measures / indicators.

Core criterion: is a result relevant (binary classification)?

Precision and Recall

Used to measure the balance between

- Returning many results, hence a lot of the relevant results present in the database, but also a lot of noise
- Returning very few results, leading to less noise, but also less relevant results

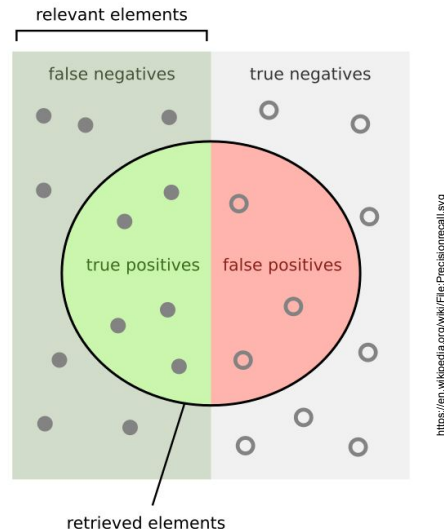
Precision and Recall

Precision (P) is the fraction of retrieved documents that are relevant

$$\text{Precision} = \frac{\#(\text{relevant items retrieved})}{\#(\text{retrieved items})} = P(\text{relevant}|\text{retrieved})$$

Recall (R) is the fraction of relevant documents that are retrieved

$$\text{Recall} = \frac{\#(\text{relevant items retrieved})}{\#(\text{relevant items})} = P(\text{retrieved}|\text{relevant})$$

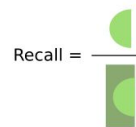


How many retrieved items are relevant?



Precision =

How many relevant items are retrieved?



Recall =

	Relevant	Nonrelevant
Retrieved	true positives (tp)	false positives (fp)
Not retrieved	false negatives (fn)	true negatives (tn)

$$P = tp / (tp + fp)$$

$$R = tp / (tp + fn)$$

F-measure

F measure is the weighted harmonic mean of precision and recall

$$F = \frac{1}{\alpha \frac{1}{P} + (1 - \alpha) \frac{1}{R}} = \frac{(\beta^2 + 1)PR}{\beta^2 P + R} \quad \text{where} \quad \beta^2 = \frac{1 - \alpha}{\alpha}$$

where $\alpha \in [0, 1]$ and thus $\beta^2 \in [0, \infty]$

The default value is $\beta = 1$, leading to:

$$F_{\beta=1} = \frac{2PR}{P + R}$$

or

$$F_1 = \frac{2}{\text{recall}^{-1} + \text{precision}^{-1}} = 2 \frac{\text{precision} \cdot \text{recall}}{\text{precision} + \text{recall}} = \frac{2TP}{2TP + FP + FN}$$

How to evaluate a ranked retrieval system?

When results are ordered, more measures are available.

Common useful measures are:

- The precision-recall graph and the mean average precision
- The ROC graph and the area under it (AUC)

Precision-recall graph

Plotting the points

For a given query

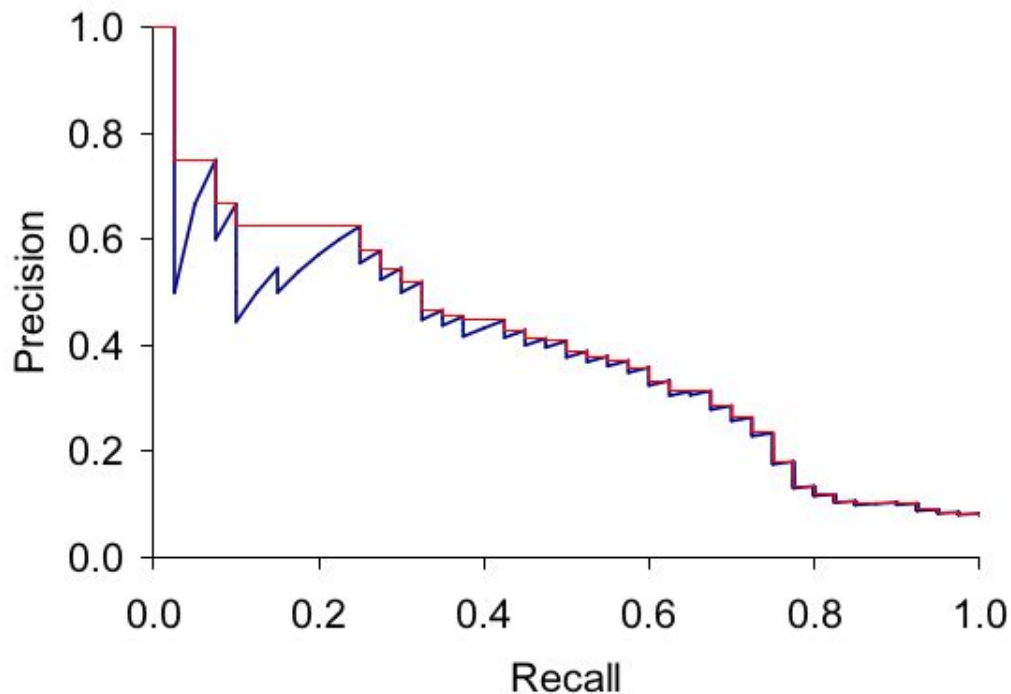
For each result

if the result is relevant

set $x = \#tp / \#expected$

set $y = \#tp / \#returned$

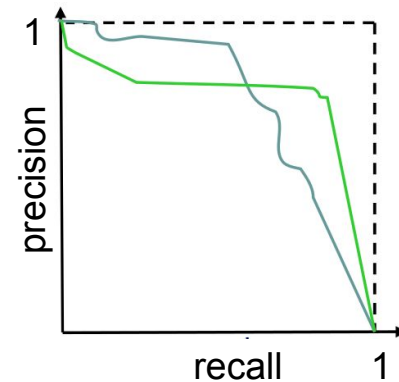
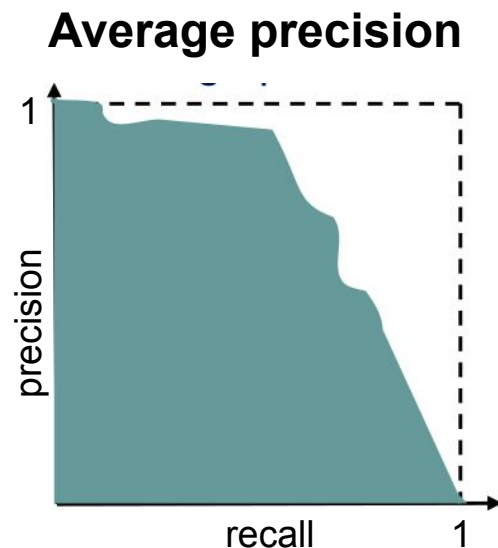
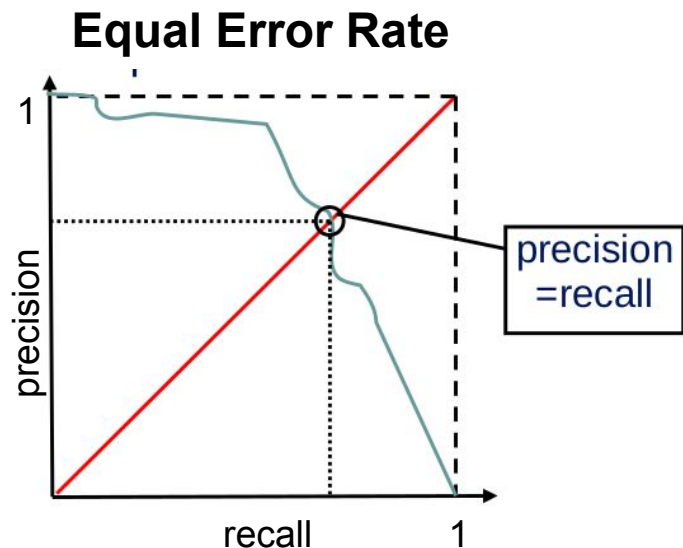
The recall always increases while we scan the result list.



Equal Error Rate and Average Precision

Note: the PR graph does not provide a total order
⇒ need more indicators

Which one is the best?



Mean average precision at k — mAP (@k)

Mean of the average precision of several queries,
when considering **k results for each query**

⇒ makes evaluation tractable with very large databases

Computed for each query using the [trapezoid technique](#)

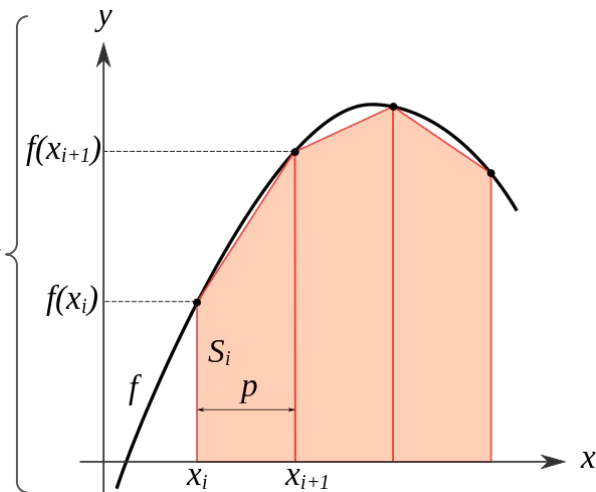
General algorithm:

For each query q_i in the test set with expected results e_i :

Retrieve the list **ret_i** of **k** best results

Compute the AP **ap_i** given e_i and **ret_i**

Compute the mean AP over all **ap_i**



Example: Compute the AP for a given query

For this query and the following results, plot the precision/recall graph and compute the average precision.



1



2



3



4



5



6



7



8



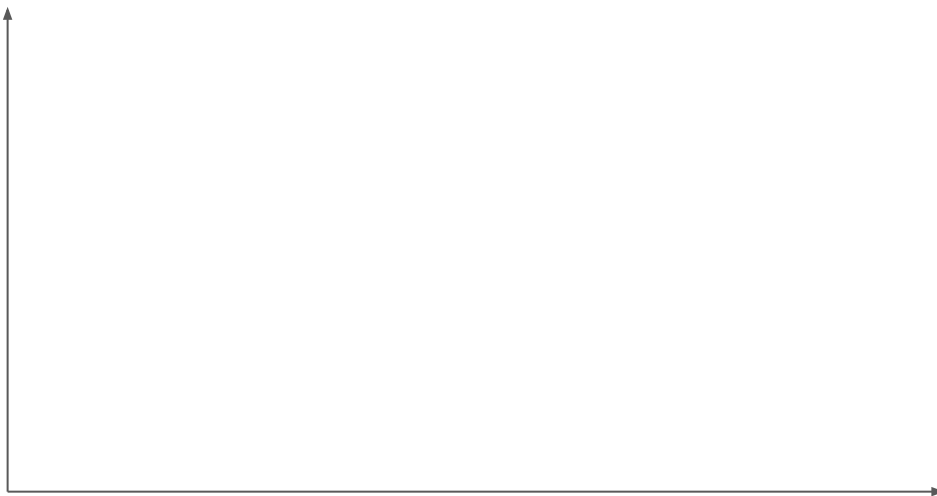
9



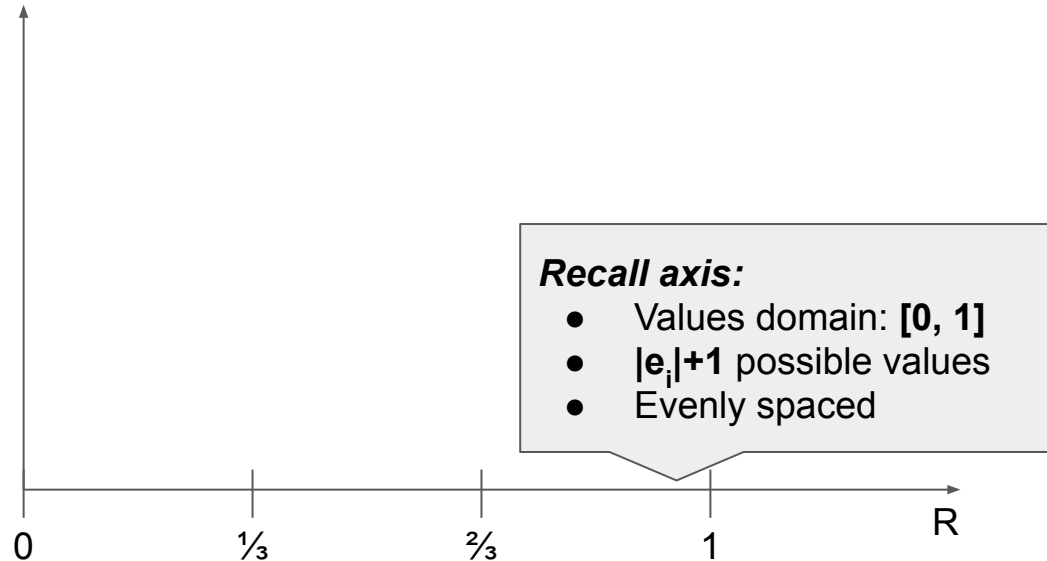
10



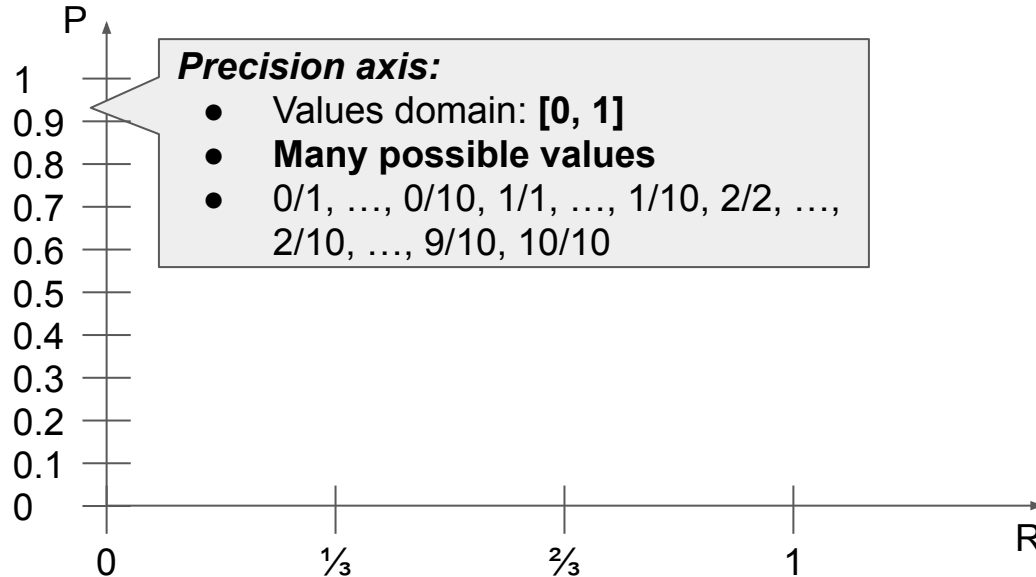
Case 1: assume $|e_i| = 3$



Case 1: assume $|e_i| = 3$

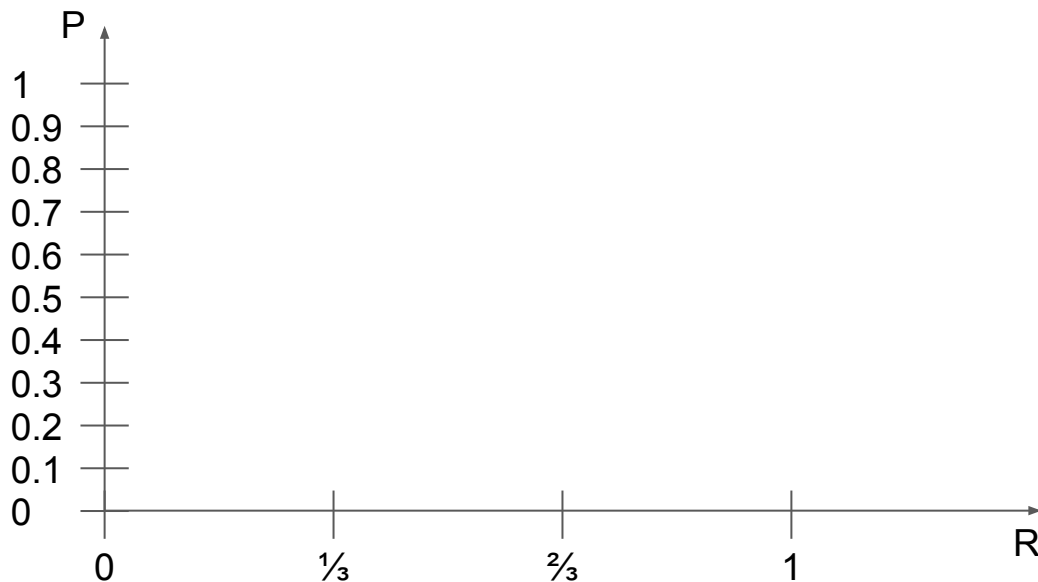


Case 1: assume $|e_i| = 3$



Case 1: assume $|e_i| = 3$

Check the **first** result:
Is it **relevant**?



Case 1: assume $|e_i| = 3$

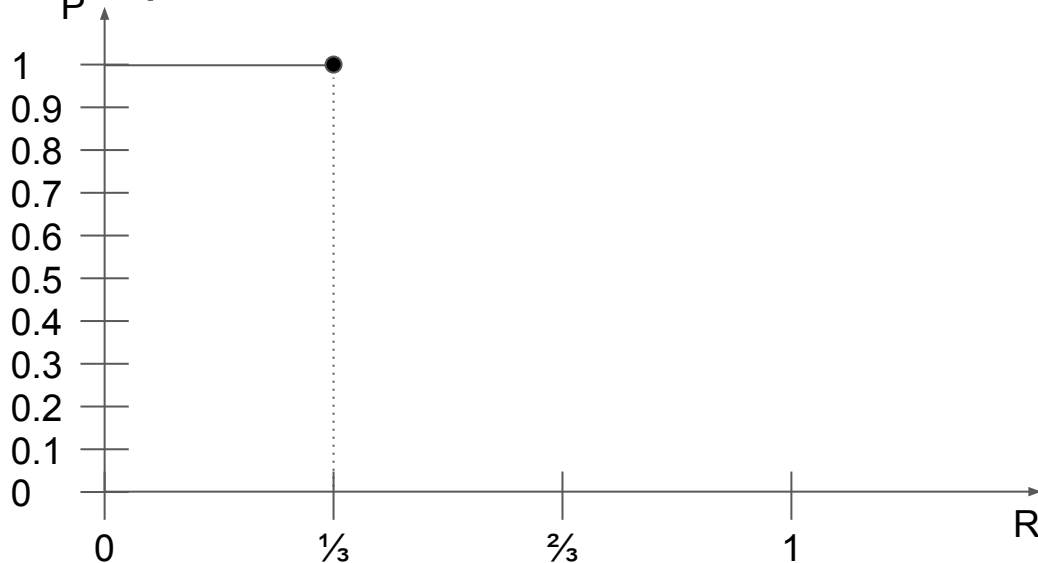
Check the **first** result:

It it **relevant**? **YES**

⇒ **Compute current precision:**

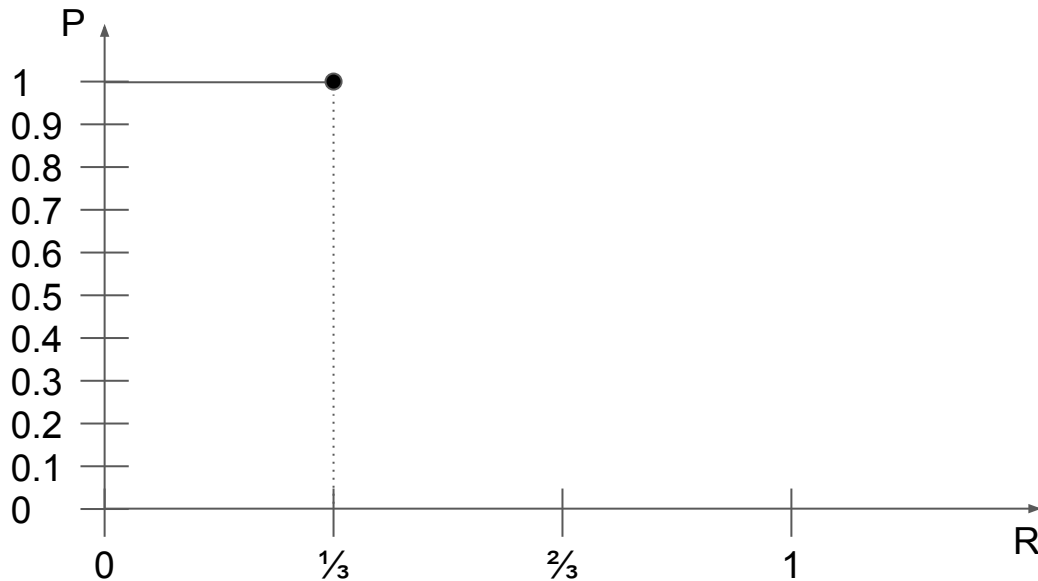
$$1 \text{ relevant} / 1 \text{ retrieved} = 1$$

⇒ **Recall = 1 relev. / 3 expected = $\frac{1}{3}$**



Case 1: assume $|e_i| = 3$

Check the **next** result:
It it **relevant**?



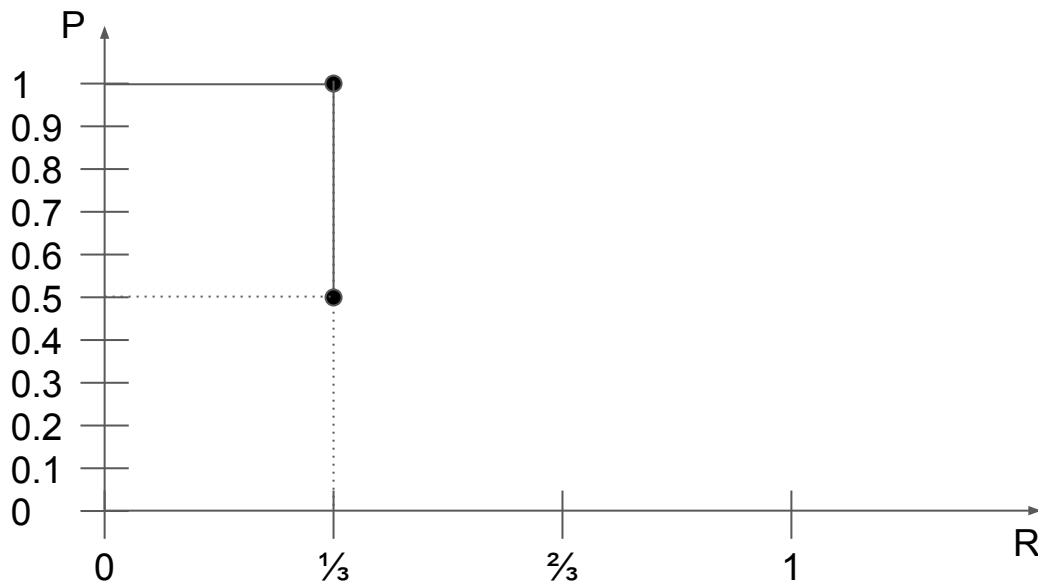
Case 1: assume $|e_i| = 3$

Check the **next** result:

Is it **relevant**? **NO**

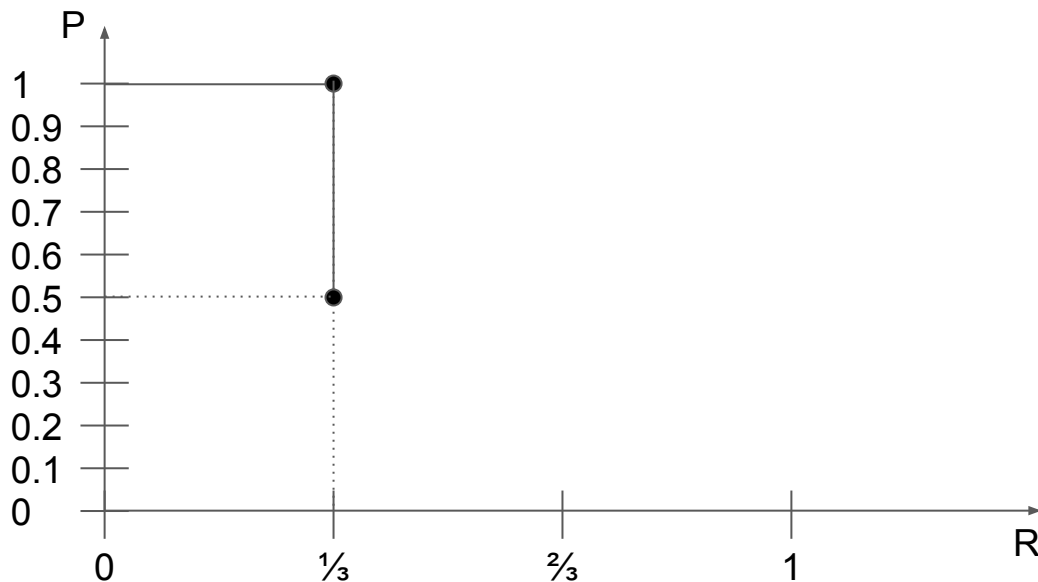
⇒ $P@2 = 1$ relevant / 2 retrieved = $\frac{1}{2}$

⇒ $R@2$ is unchanged



Case 1: assume $|e_i| = 3$

Check the **next** result:
It it **relevant**?



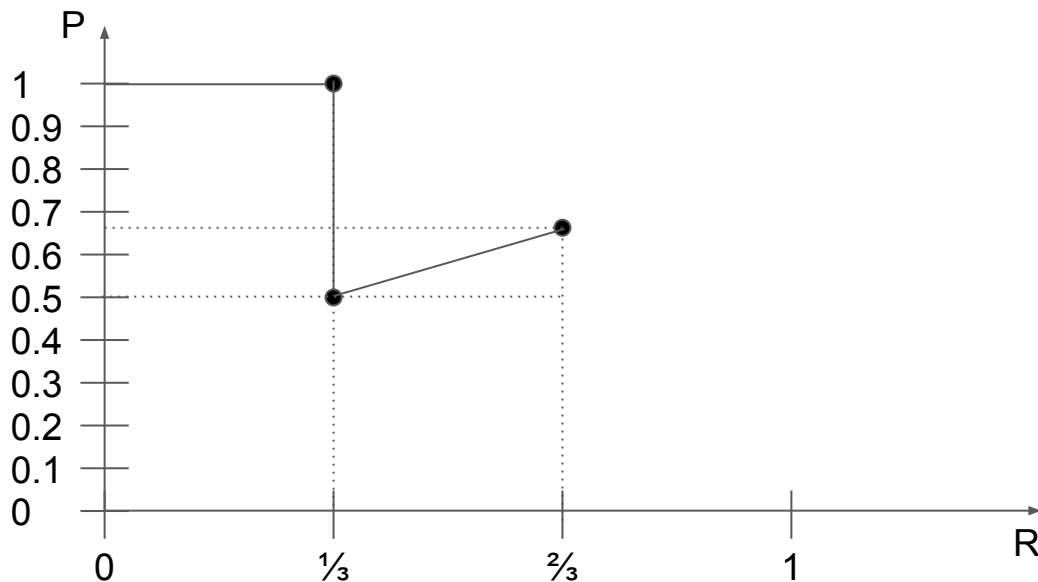
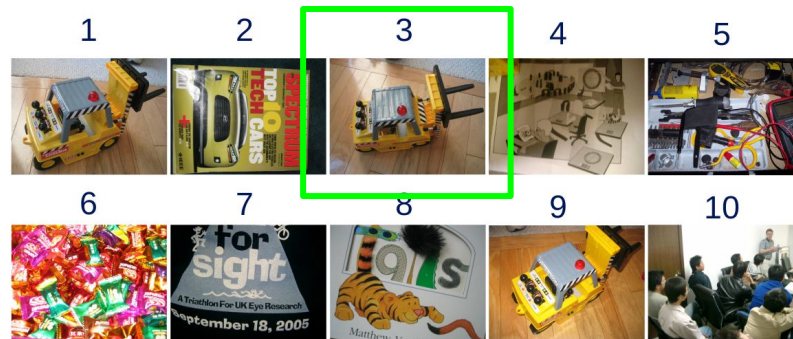
Case 1: assume $|e_i| = 3$

Check the **next** result:

It it **relevant**? **YES**

⇒ $P@3 = 2 \text{ relevant} / 3 \text{ retrieved} = \frac{2}{3}$

⇒ Add a point at next recall value ($\frac{2}{3}$)

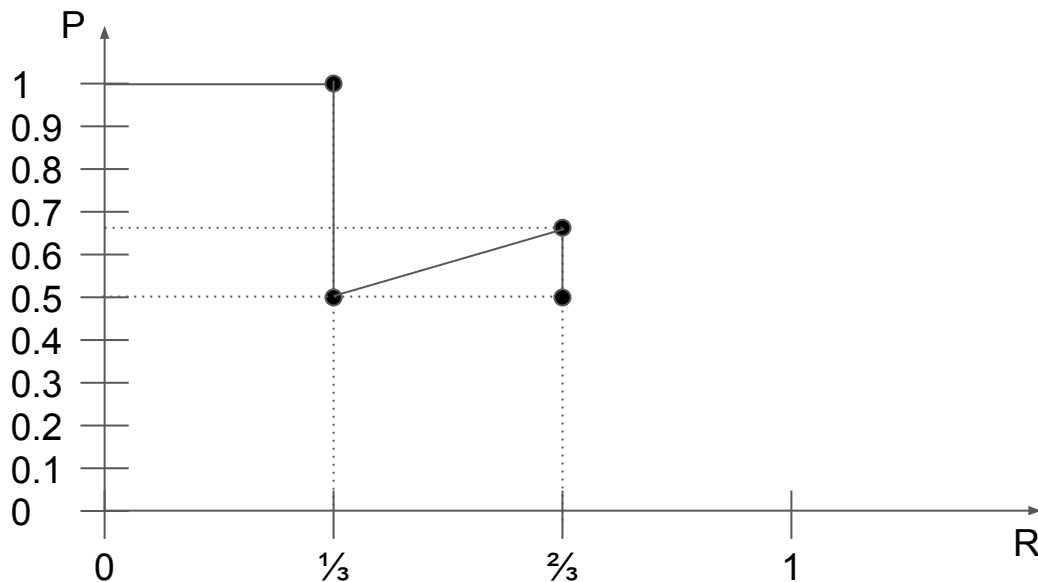


Case 1: assume $|e_i| = 3$

And we keep going...

$$P@4 = 2/4 = 1/2$$

$R@4$ = unchanged

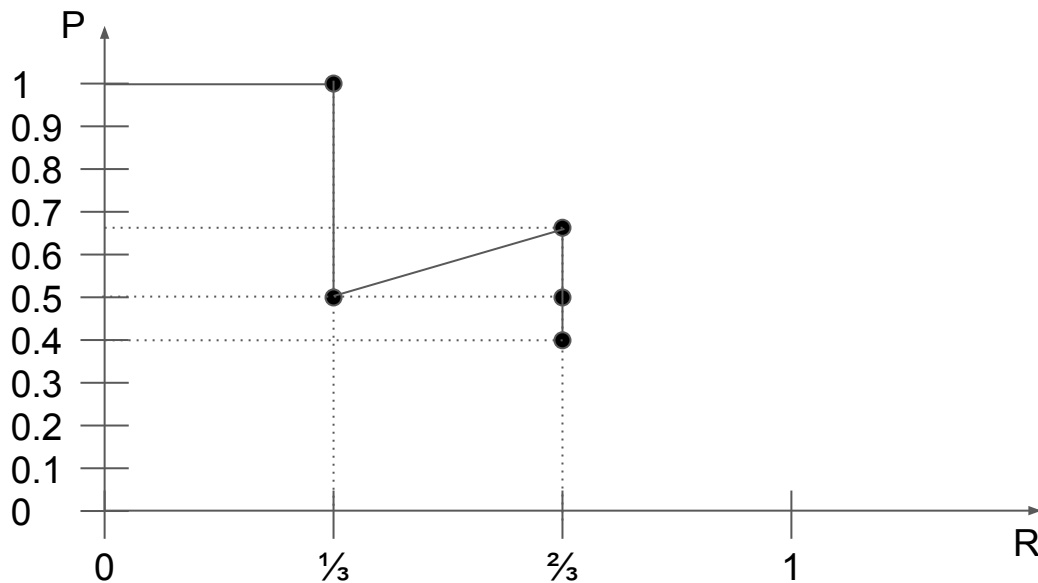


Case 1: assume $|e_i| = 3$

And we keep going...

$$P@5 = 2/5 = 0.4$$

$R@5$ = unchanged

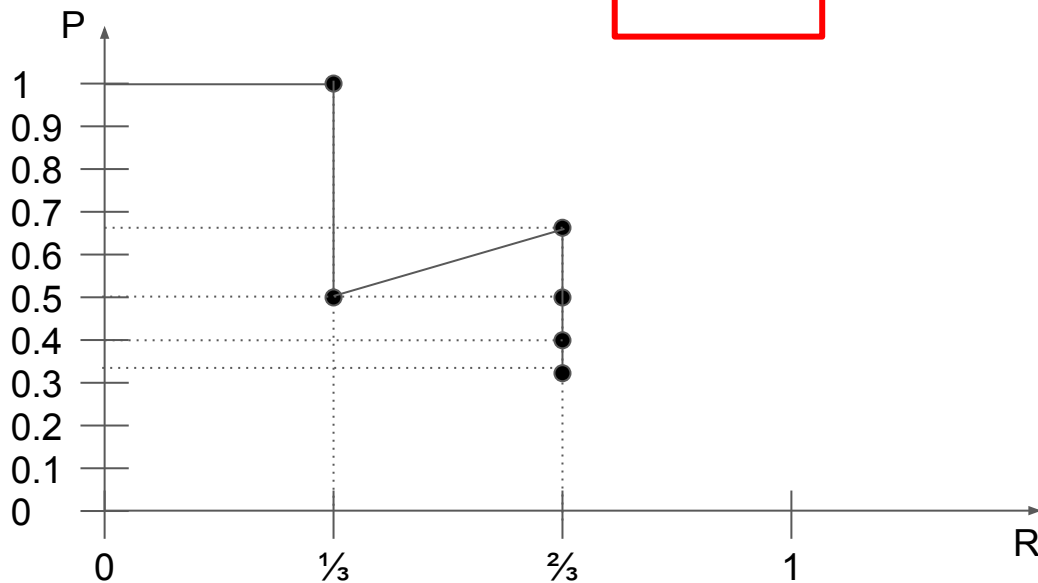


Case 1: assume $|e_i| = 3$

And we keep going...

$$P@6 = 2/6 = 1/3$$

$R@6$ = unchanged

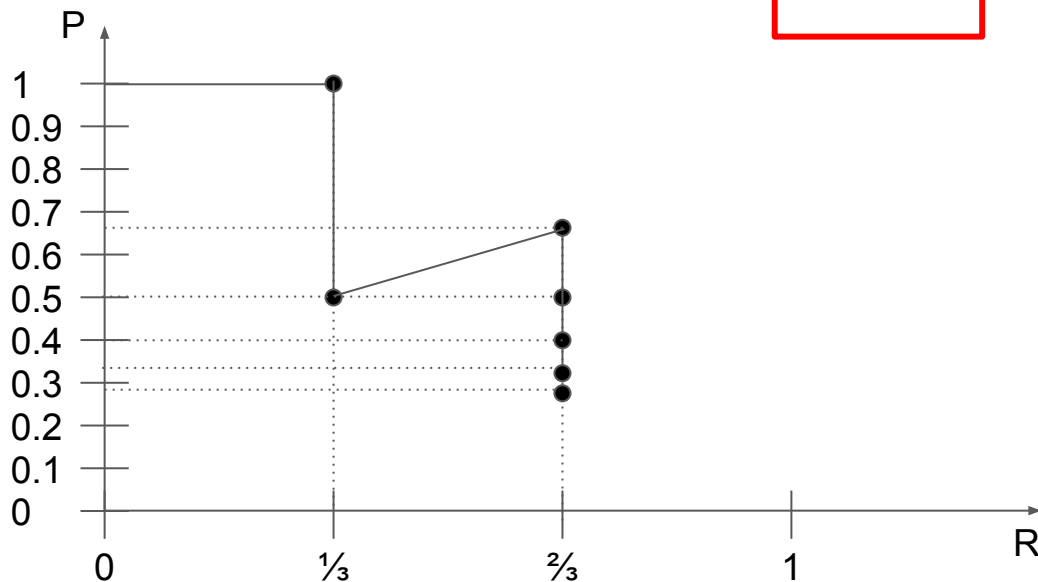


Case 1: assume $|e_i| = 3$

And we keep going...

$P@7 = 2/7 = 0.285...$

$R@7 = \text{unchanged}$

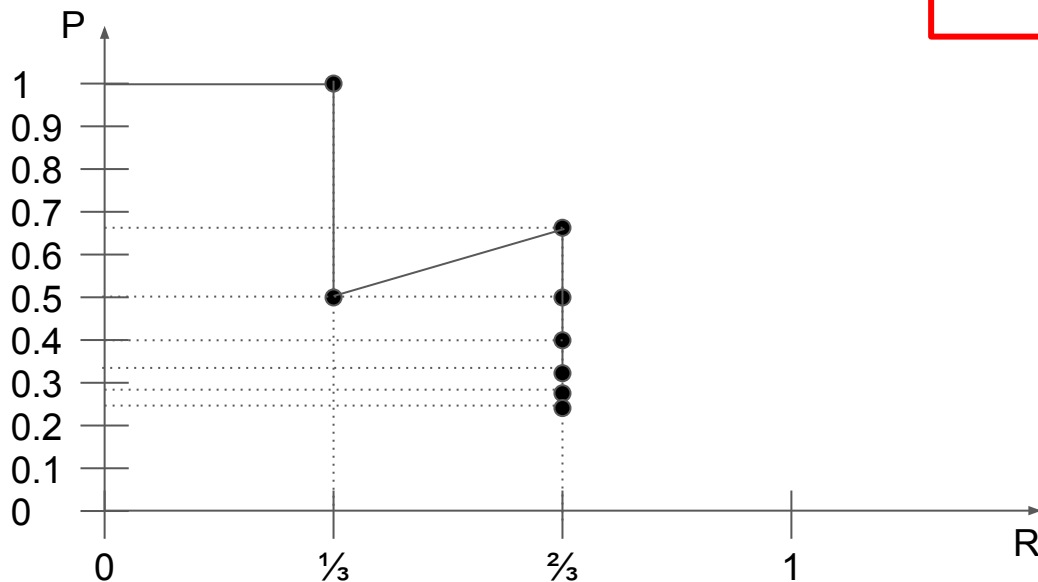


Case 1: assume $|e_i| = 3$

And we keep going...

$P@8 = 2/8 = 1/4$

$R@8 = \text{unchanged}$

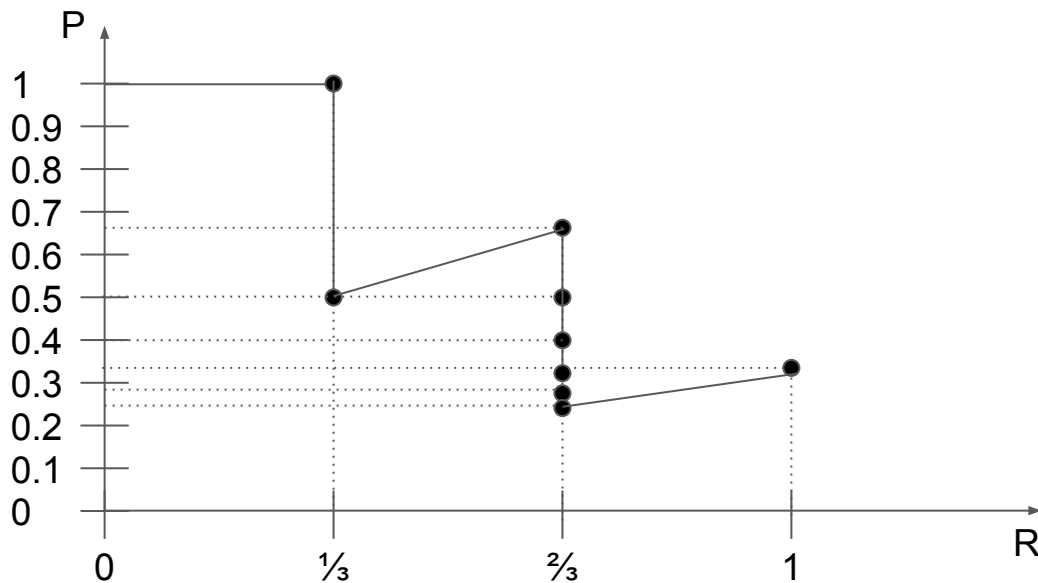


Case 1: assume $|e_i| = 3$

And we keep going...

$$P@9 = 3/9 = 1/3$$

$$R@9 = 3/3 = 1$$

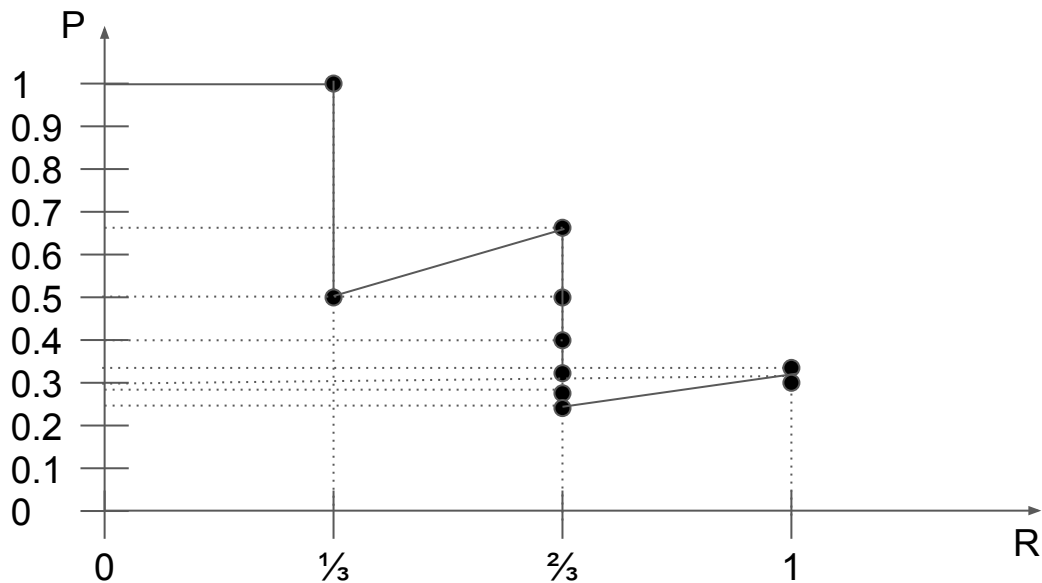


Case 1: assume $|e_i| = 3$

It does not change the AP here...

$$P@10 = 3/10$$

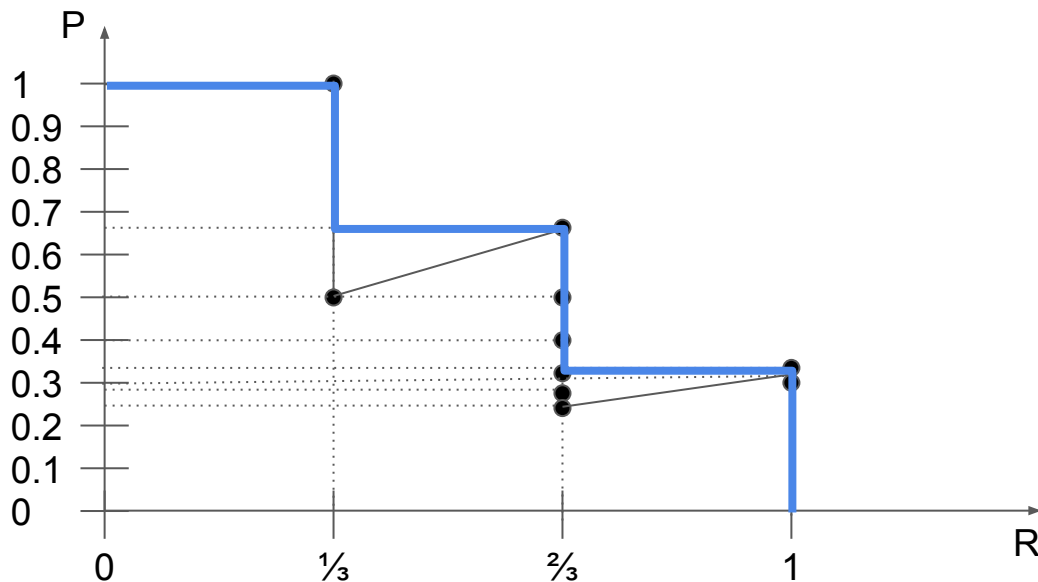
$$R@10 = 3/3 = 1$$



Case 1: assume $|e_i| = 3$

And we are done!

A common approximation is to take
only the **upper envelope** of the curve...



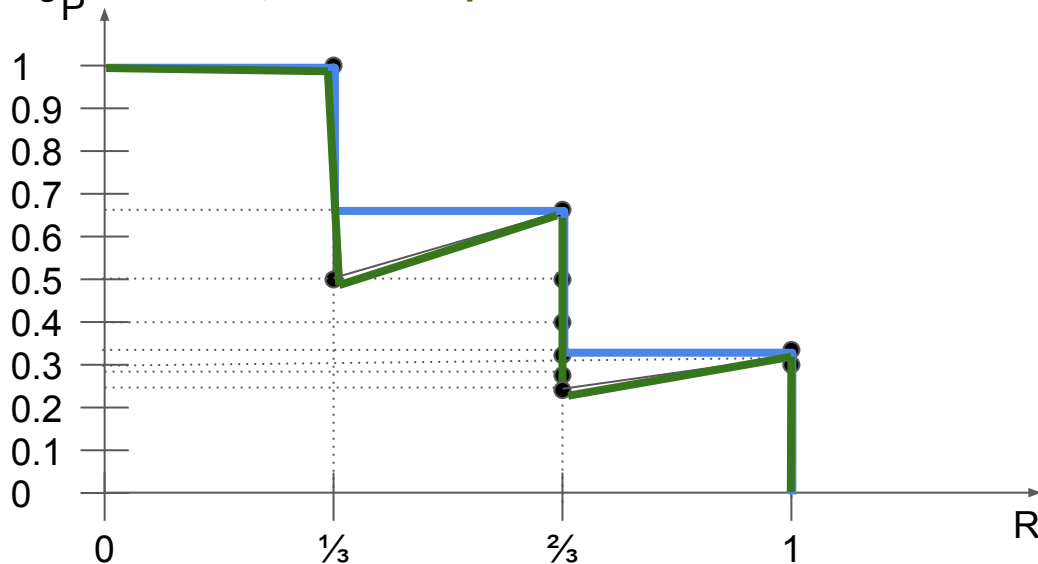
Case 1: assume $|e_i| = 3$

And we are done!



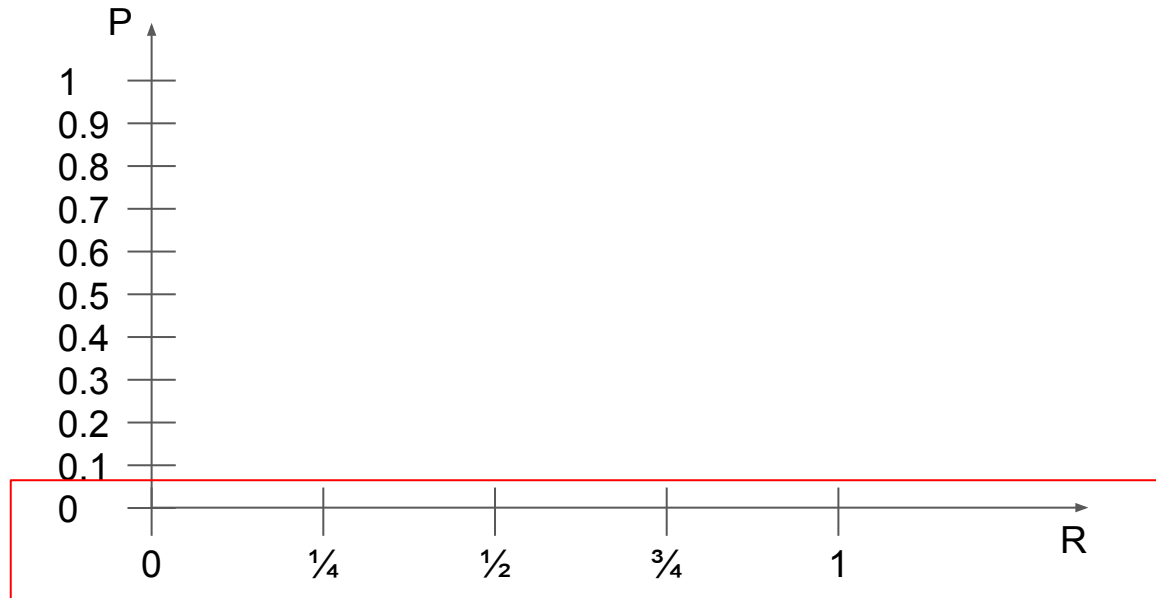
A common approximation is to take
only the upper envelope of the curve...

But **good libraries** go for the **full, exact computation.**



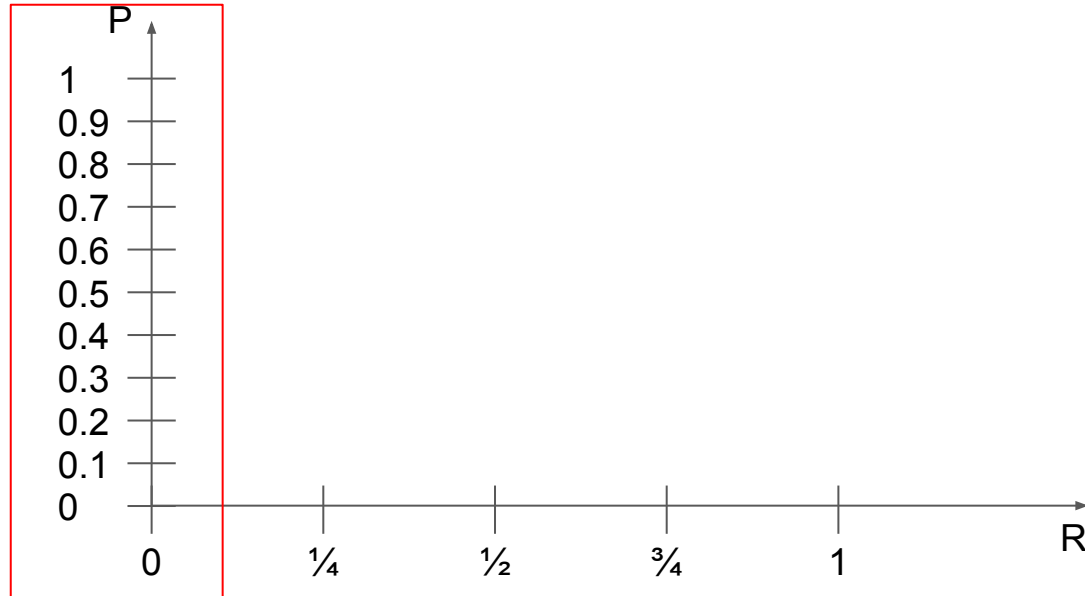
Case 2: what if $|e_i| = 4$?

1. Adjust R values.



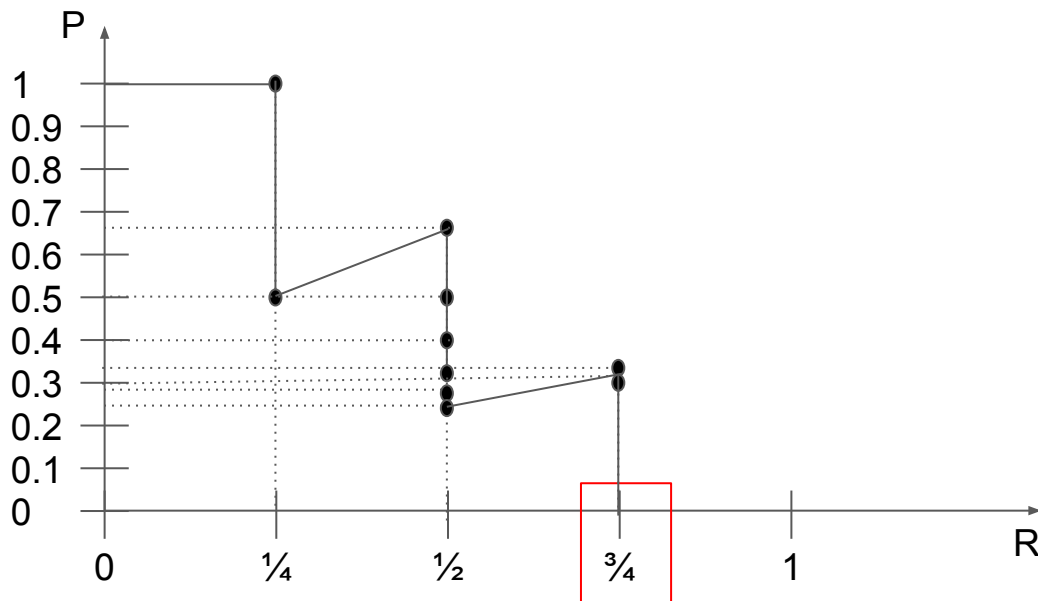
Case 2: what if $|e_i| = 4$?

1. Adjust R values.
2. P values do not change if k does not change.



Case 2: what if $|e_i| = 4$?

1. Adjust R values.
2. P values do not change if k does not change.
3. Here, it would imply that we did not get all relevant results (very common in practice) \Rightarrow we stop the curve before the 1



ROC & others

[next lecture, more useful for classification]

Ground truthing issues

Do we have to annotate all images within a dataset for all our test queries?

No! Use “**distractors**”: samples that you know, for sure, not to be relevant to any query.