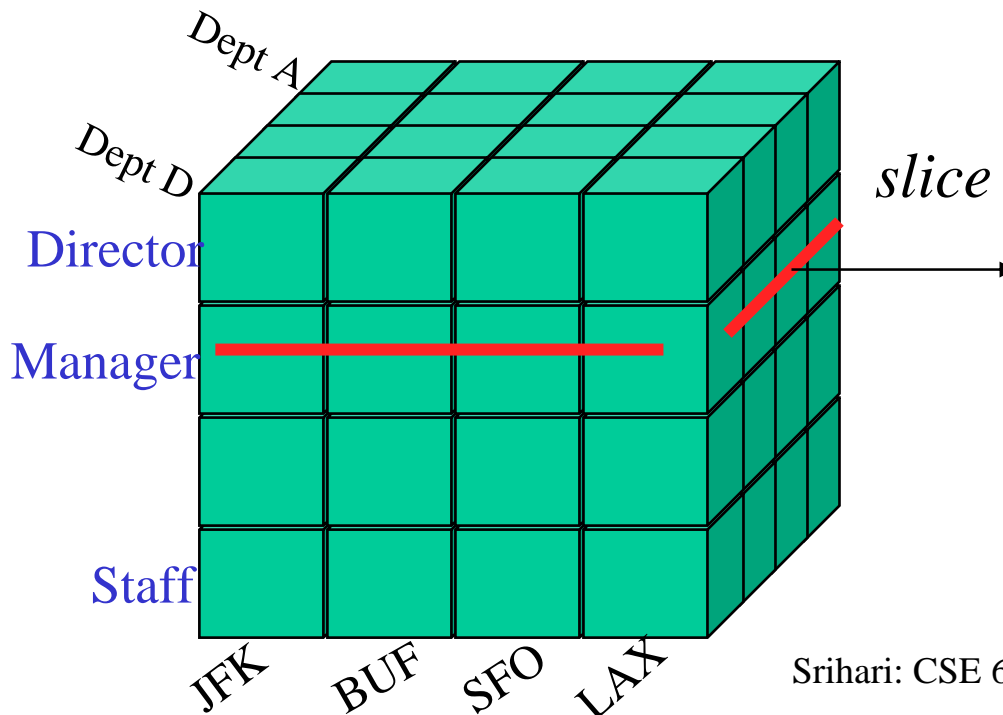


Retrieval by Content

Database Retrieval

- In a Database Context
 - Query is well-defined
 - Operation returns a set of records (or entities) that exactly match required specifications
 - Example query
 - [level = MANAGER] AND [age < 30]
 - Returns list of young employees with significant responsibility



*Drill down to
Records for each
Department, location.
Look up age field*

*Roll-up by East Coast
is another operation*

Retrieval by Content

- More general, less precise queries than Database Retrieval
- Example of Medical Context:
 - Query is a patient record containing
 - Demographic information (age, sex,..)
 - Test Results (blood Tests, physical tests, biomedical time series, X-rays)
 - Search database for similar cases in hospital database
 - To determine diagnoses, treatments, outcomes
 - Exact match is not relevant since it is unlikely there is any other patient that matches exactly
 - Need to determine **similarity** among patients based on different data types (multivariate, time series, image data)

Retrieval Task

- Find the k objects in the database that are most similar to either a specific query or a specific object
- Examples:
 - Searching historical records of Dow Jones index for past occurrences of a particular time series pattern
 - Searching a database of satellite images for evidence of volcano eruptions
 - Searching internet for reviews of restaurants in Buffalo

Retrieval by Content is *Interactive Data Mining*

- User is directly involved in exploring data set by
 - Specifying a query
 - Interpreting results of matching process
- Role of human judgement is not prominent in *predictive* and *descriptive* forms of data mining
- If database is pre-indexed by content then task reduces to standard database indexing
- Instead we have a *query pattern Q*
 - Goal is to infer which other objects are most similar to Q
 - In *Text Retrieval* Q is a short list of query words matched with large sets of documents

Retrieval by Content depends on notion of *Similarity*

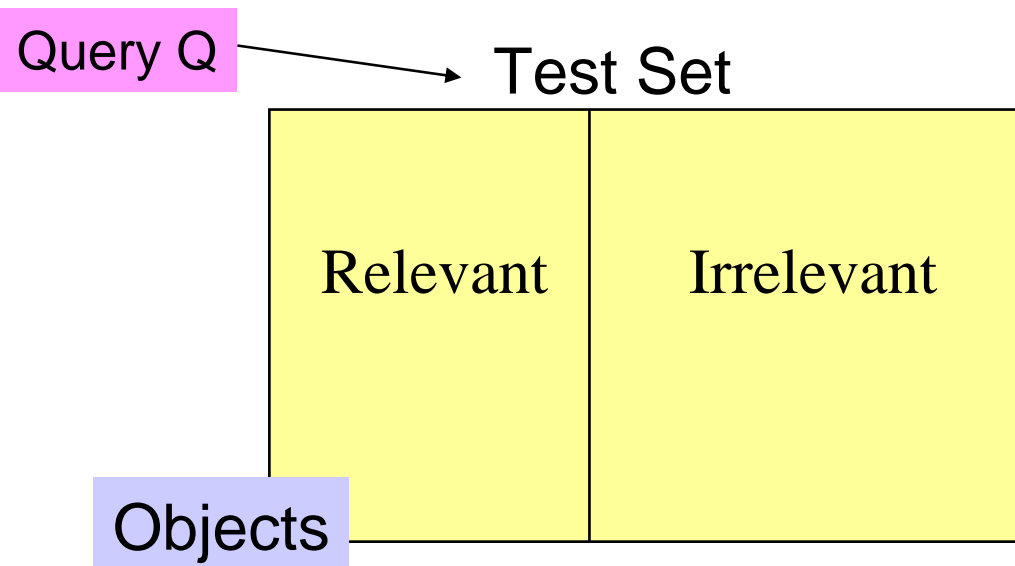
- Either Similarity or Distance is used
- Maximize similarity or minimize distance
- Common to reduce measurements to a standard fixed-length vector and use geometric measures (Euclidean, weighted Euclidean, Manhattan, etc)

Retrieval Performance

- In classification and regression
 - There is an *objective* measure of accuracy of model on unseen test data
 - Comparison of different algorithms and models is straightforward
- In retrieval
 - Performance is *subjective*: relative to a query
 - Ultimate measure is usefulness to user
 - Performance evaluation is difficult
 - Objects in data set need to be labelled as relevant to query

Evaluation of a Retrieval Algorithm

- In response to a specific query Q
- Independent test data set
 - Test data has not been tuned to given query Q
- Objects of the test data set have been pre-classified (truthed) as being relevant or irrelevant to query Q
 - Algorithm is not aware of class labels
 - Who determines whether object is relevant?



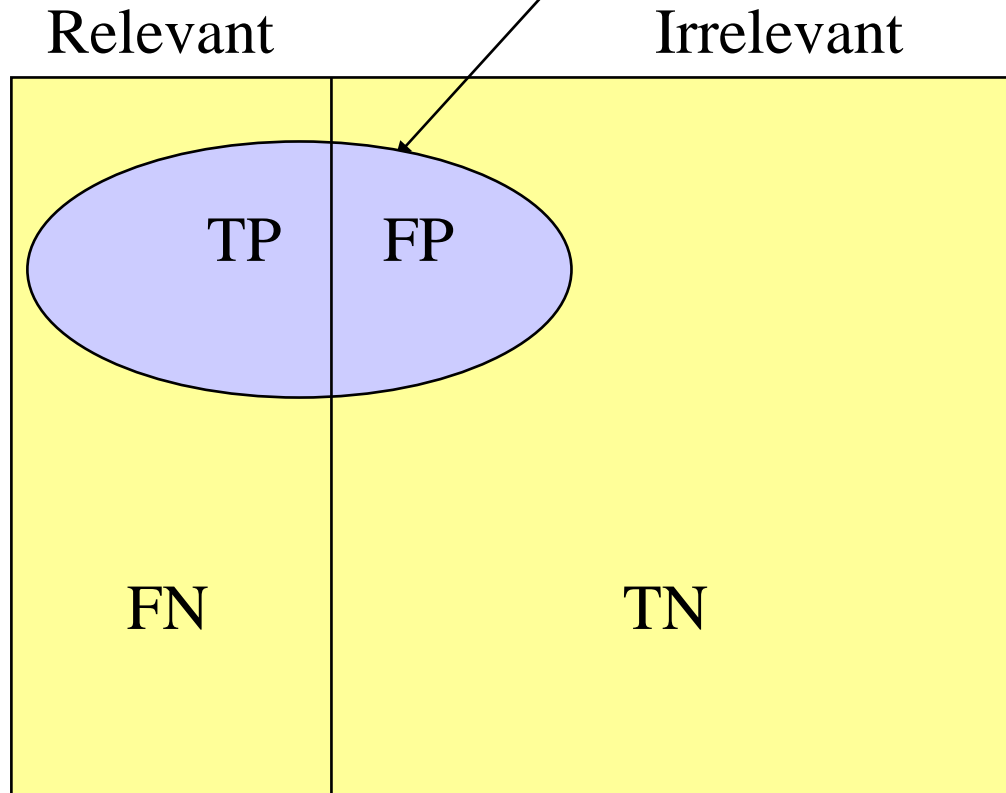
Confusion Matrix

	Truth: Relevant	Truth: Not- Relevant
Algorithm: Relevant	TP	FP
Algorithm: Not Relevant	FN	TN

Precision and Recall Definitions

Obtained from Confusion Matrix

Objects returned for query Q

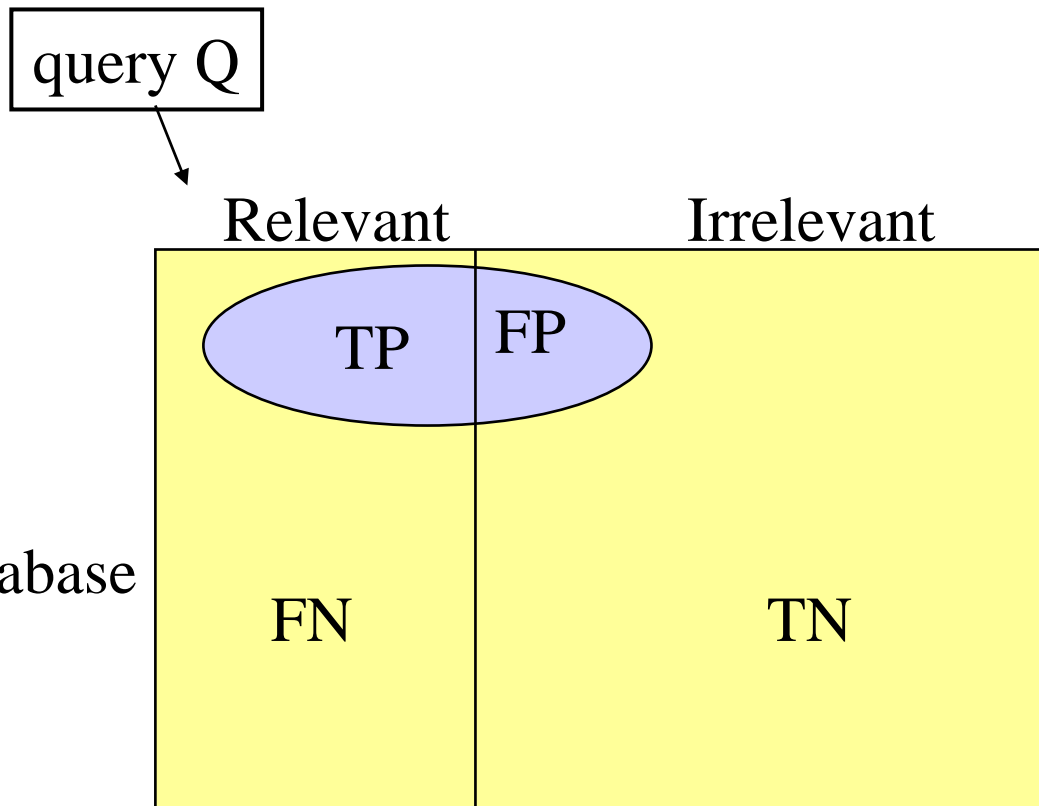


$$\text{Recall} = \frac{\text{TP}}{\text{TP} + \text{FN}} \times 100\%$$

$$\text{Precision} = \frac{\text{TP}}{\text{TP} + \text{FP}} \times 100\%$$

Observations about Precision and Recall

1. Numerator is same for precision and recall: no of correct returned
2. Denominator for **precision** is all that is returned
3. Denominator for **recall** is all that is relevant



$$\text{Recall} = \frac{\text{TP}}{\text{TP} + \text{FN}} \times 100\%$$

Recall=1 means the whole truth

$$\text{Precision} = \frac{\text{TP}}{\text{TP} + \text{FP}} \times 100\%$$

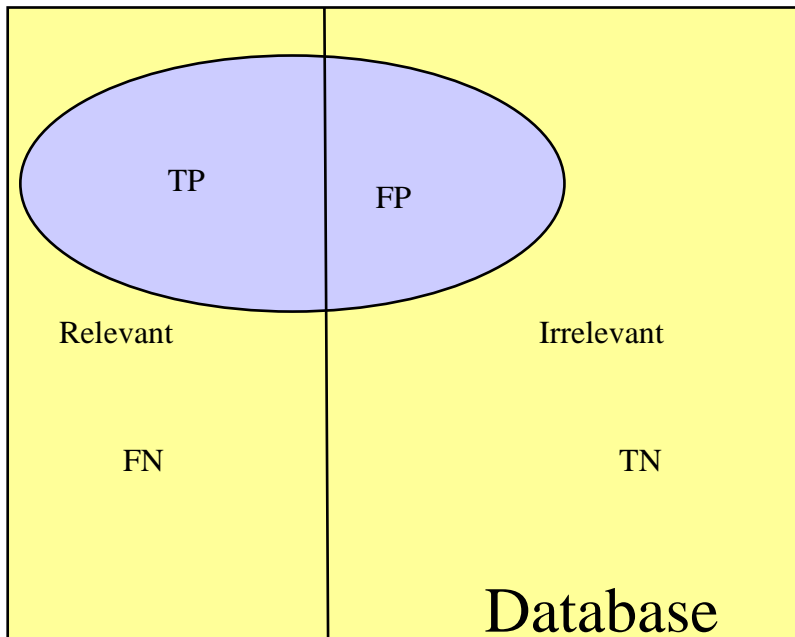
Precision=1 means nothing but the truth

Precision versus Recall

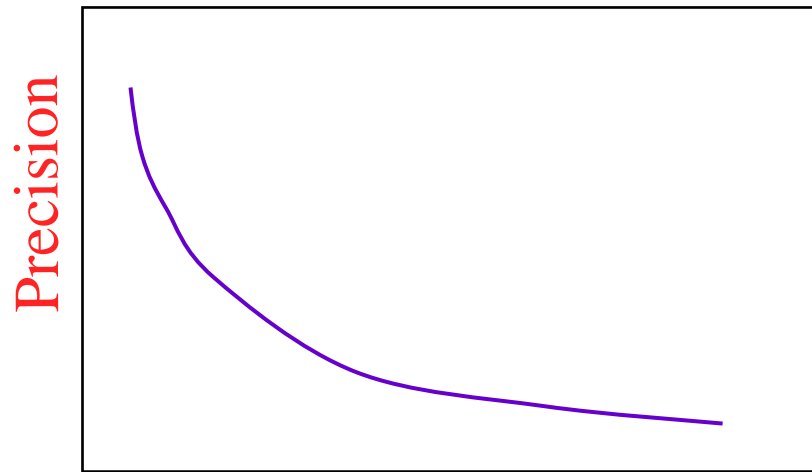
- Assume that the results of retrieval have been pre-classified as relevant or irrelevant w.r.t query Q
- If algorithm uses a distance measure to rank objects, then a threshold T is used
 - then K_T objects are returned as closer than threshold T to query object Q
- If we run the retrieval algorithm with a set of values of T we get different pairs of (recall, precision) values– giving recall-precision characterization
 - Relative to query Q , particular data set, labeling of the data

Precision-Recall Relationship

Precision-Recall
are evaluated
w.r.t. a set of queries



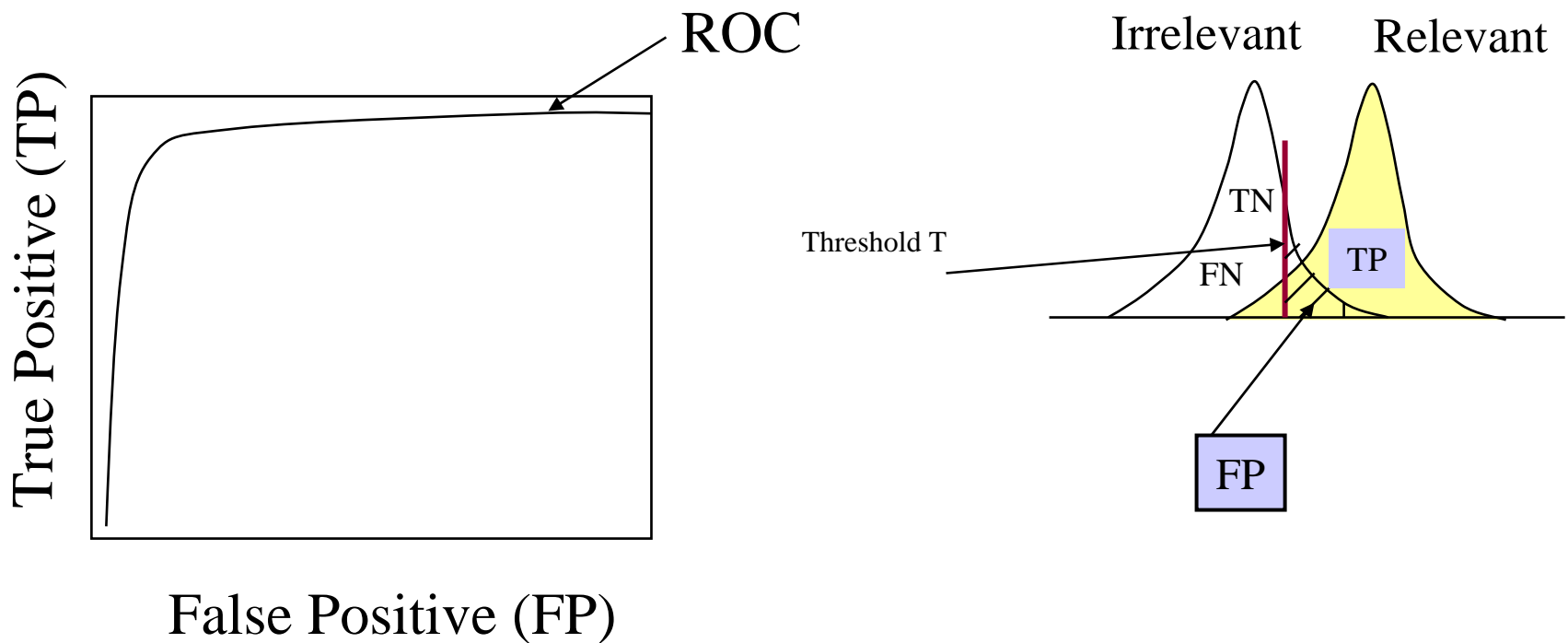
Typically an inverse relationship:
as FP is decreased (to increase precision),
TP also decreases
and FN increases (decreasing recall)



$$\text{Precision} = \frac{TP}{TP+FP}$$
$$\text{Recall} = \frac{TP}{TP+FN}$$

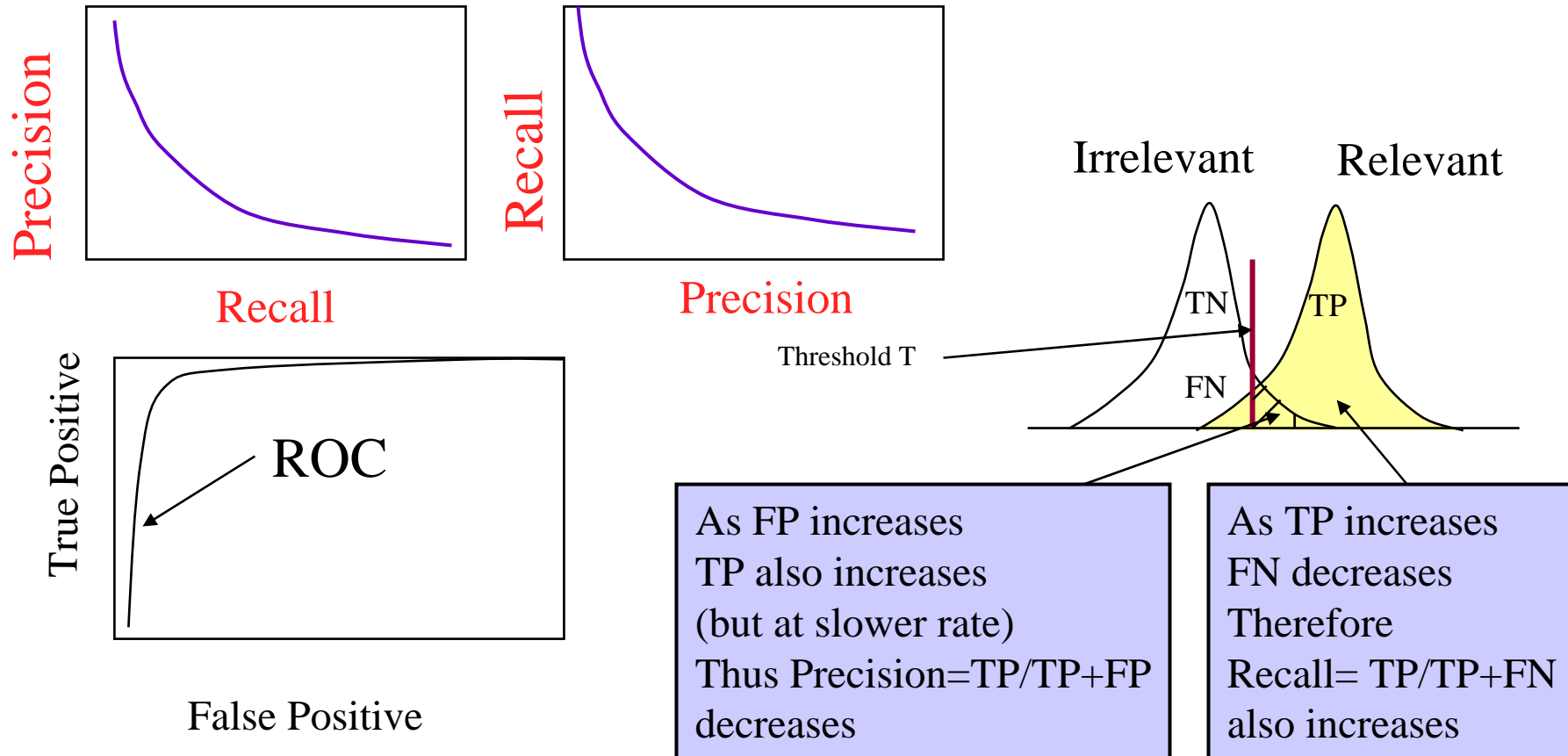
How is Precision-Recall related to ROC?

- Receiver Operating Characteristics (ROCs) are used to characterize performance of binary classifiers with variable thresholds



Relationship between Precision-Recall and ROC

- Receiver Operating Characteristics (ROCs) are used to characterize performance of binary classifiers with variable thresholds



Thus ROC is inverse of Recall-Precision Plot

Combined Measure of Retrieval

- Harmonic Mean of Precision and Recall

$$\frac{1}{F} = \frac{1}{2} \left(\frac{1}{P} + \frac{1}{R} \right)$$

- Or

$$F = \frac{2P \bullet R}{P + R}$$

- If you travel at 20 mph one way and 40 mph the other way, the average speed is given by the harmonic mean of 26.6 mph
- Harmonic mean is appropriate when the average of a rate is desired

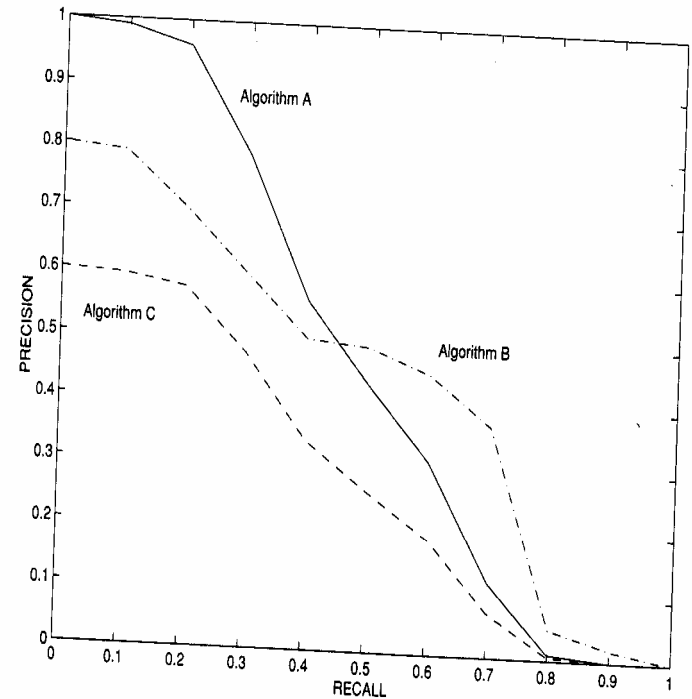
Precision-Recall of several algorithms

Precision-Recall are evaluated
w.r.t. the same data set and a
set of queries

Cannot distinguish between two
algorithms

Except at say:

1. Precision = recall
2. Precision when a certain no
are retrieved
3. Average precision over
multiple recall levels



A simple (synthetic) example of precision-recall curves for three hypothetical query algorithms. Algorithm A has the highest precision for low recall values, while Algorithm B has the highest precision for high recall values. Algorithm C is universally worse than A or B, but we cannot make a clear distinction between A or B unless (for example) we were to operate at a specific recall value. The actual recall-precision numbers shown are fairly typical of current text-retrieval algorithms; e.g., the ballpark figures of 50% precision at 50% recall are not uncommon across different text-retrieval applications.

Precision-Recall Properties

- Should average over large corpus/query ensembles
- Need human assessments
 - People aren't reliable assessors
- Assessments have to be binary
 - Nuanced assessments?