

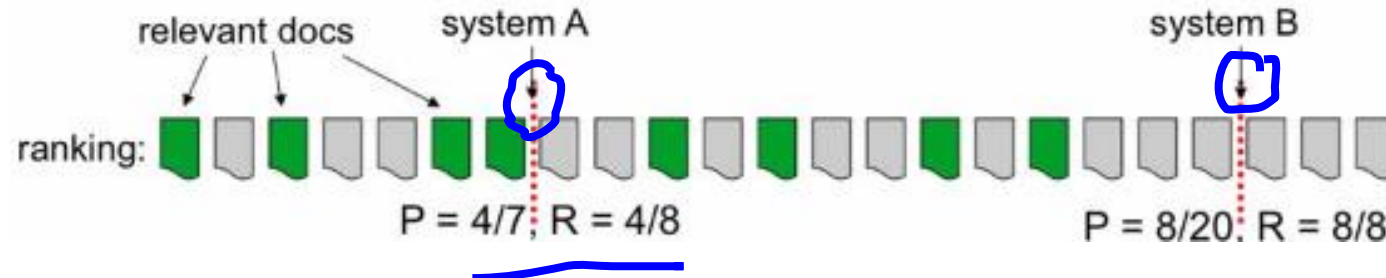
IR evaluation

Ack: Victor Lavrenko

Why not accuracy?

Comparing recall / precision

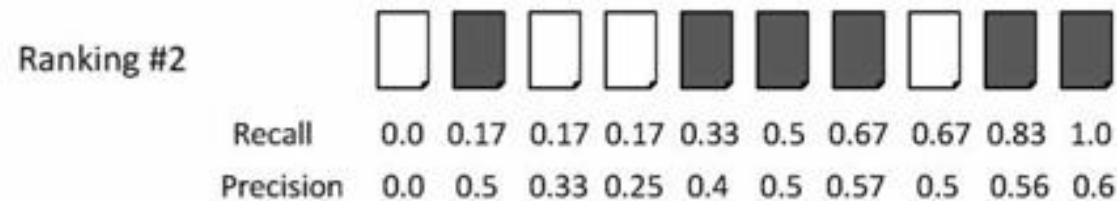
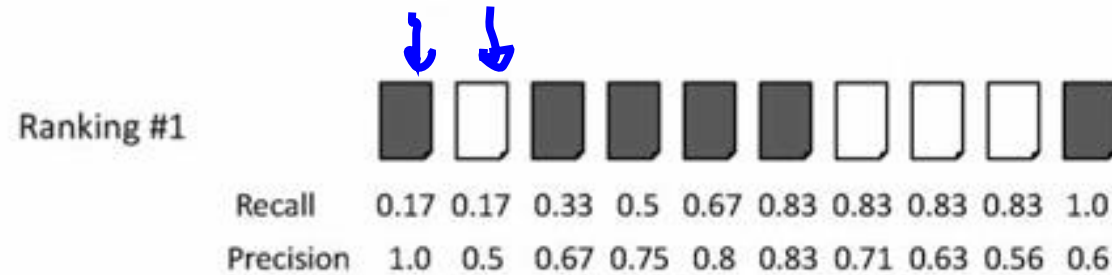
- Which of the following is a better system?
 - system A: recall = 50%, precision = 57%, $F_1=53\%$
 - system B: recall = 100%, precision = 40%, $F_1=57\%$
- Could be the same exact system
 - using different threshold settings
 - R/P , F_1 comparisons often meaningless
 - more informative to compare ranking against ranking



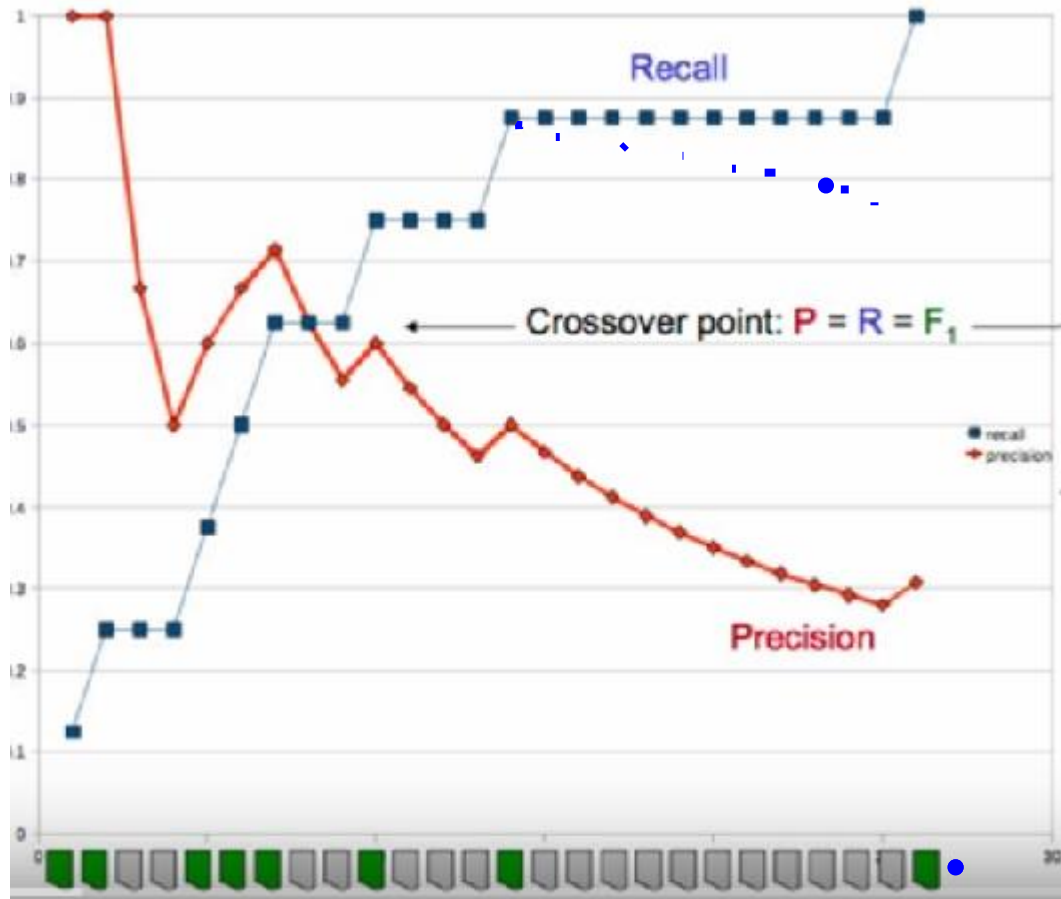
Recall / Precision and ranking

- Search engine produces a ranking, not a set
 - can compute recall, precision at every rank

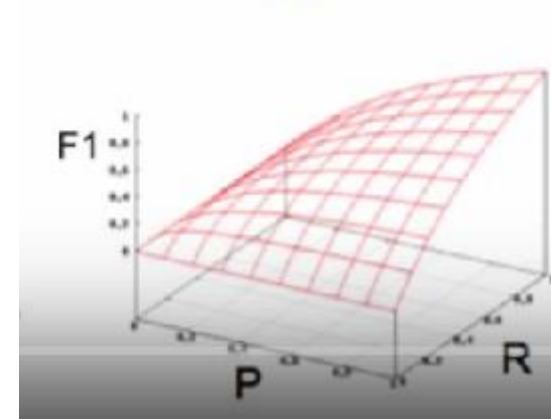
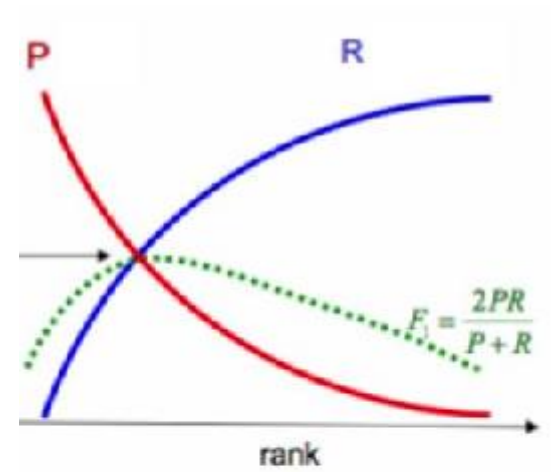
 = the relevant documents

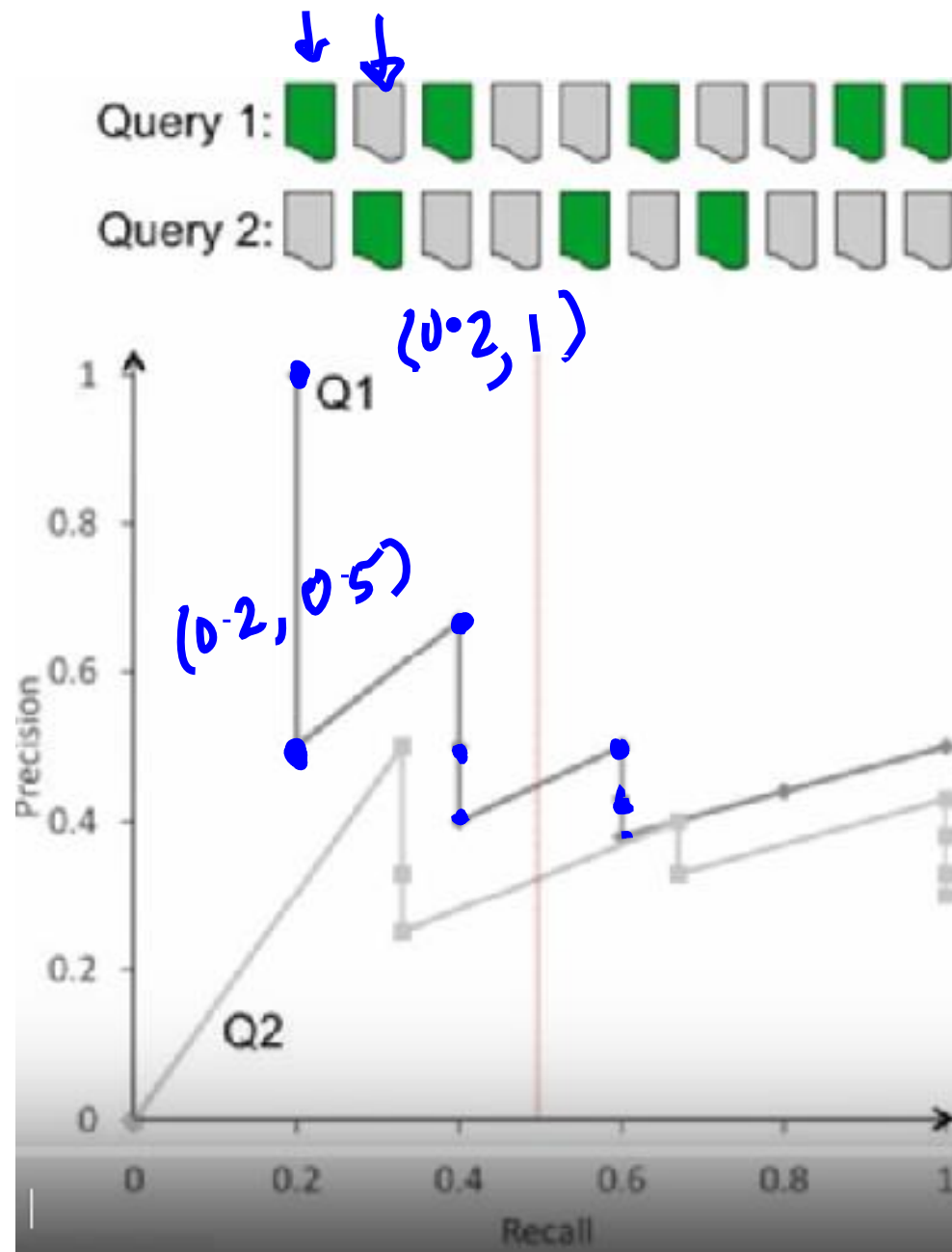


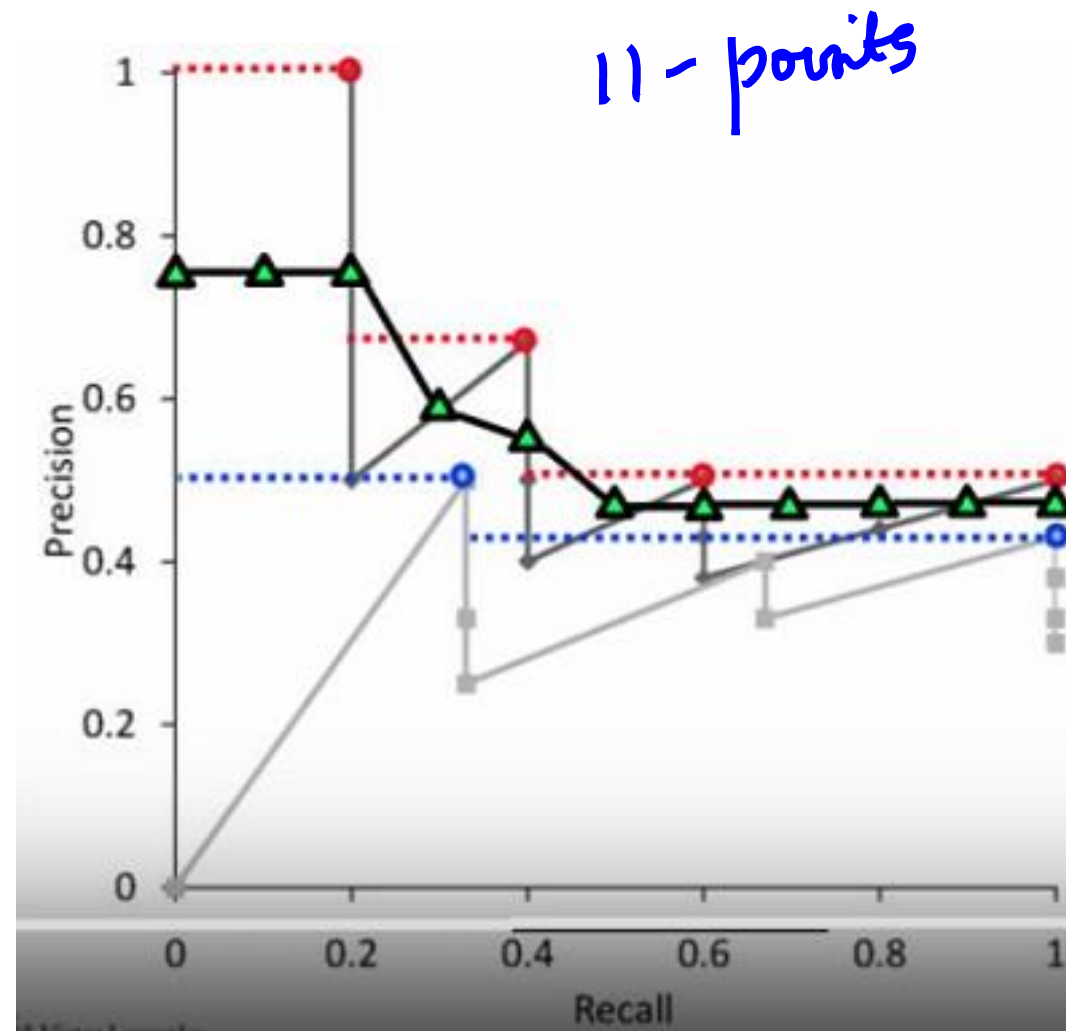
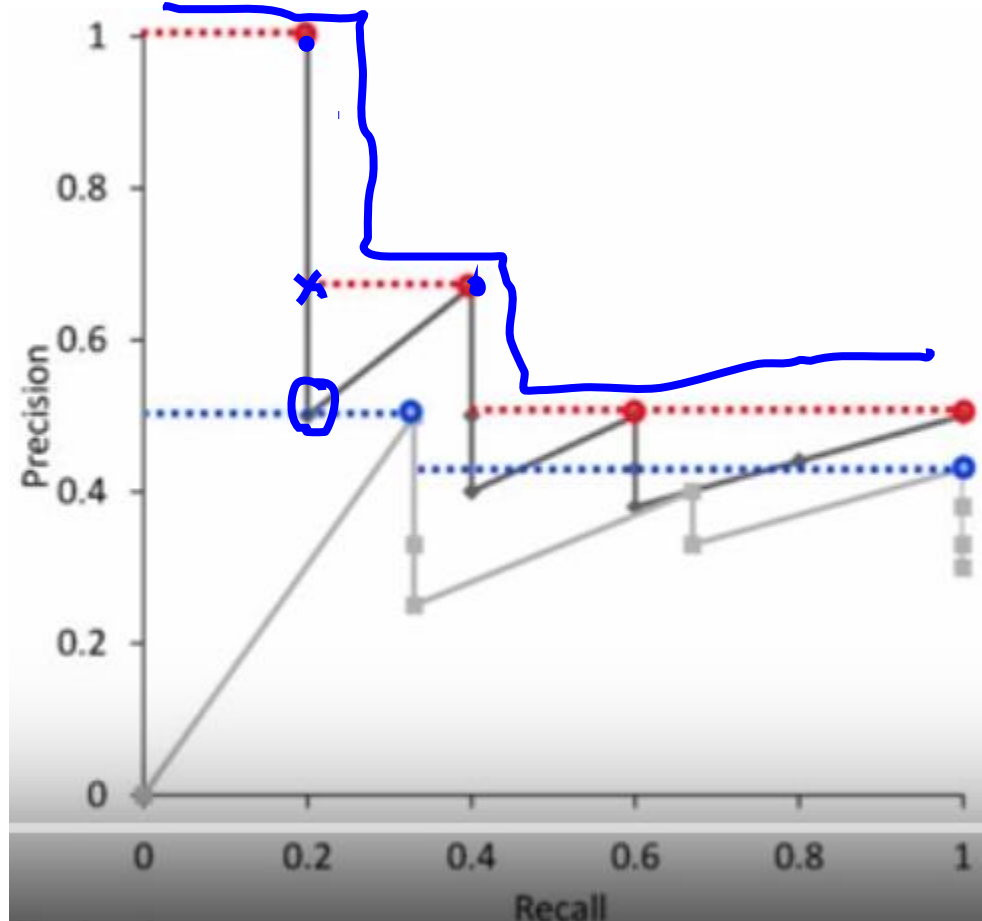
$$\underline{F = P = R}$$

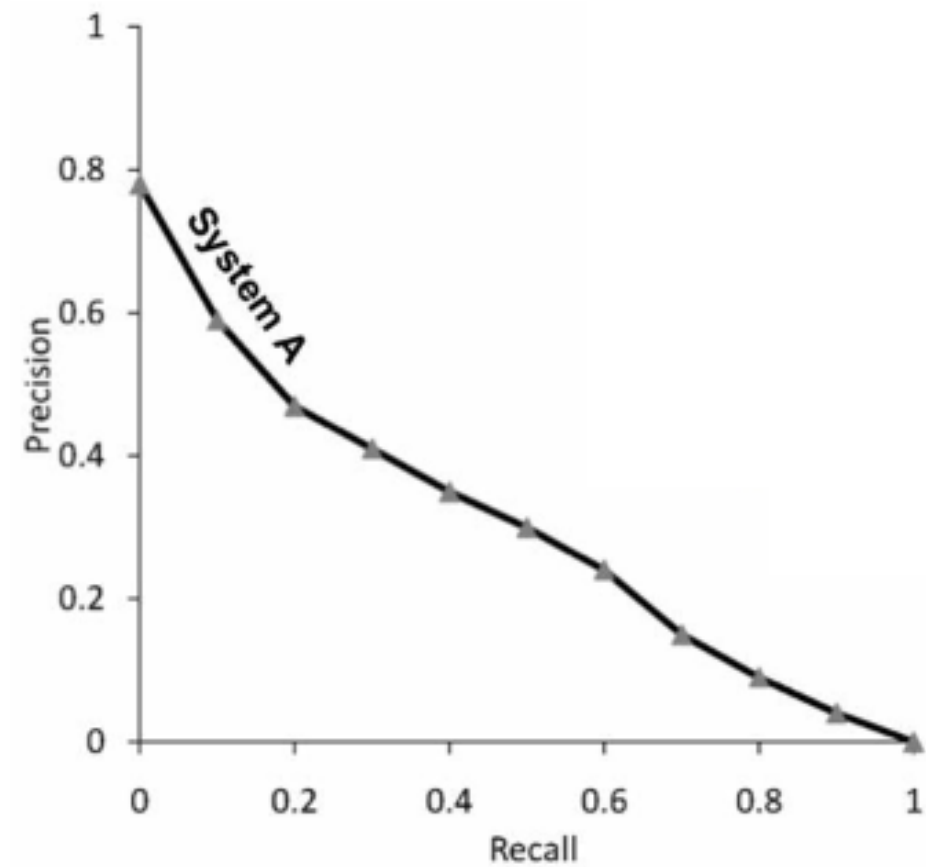
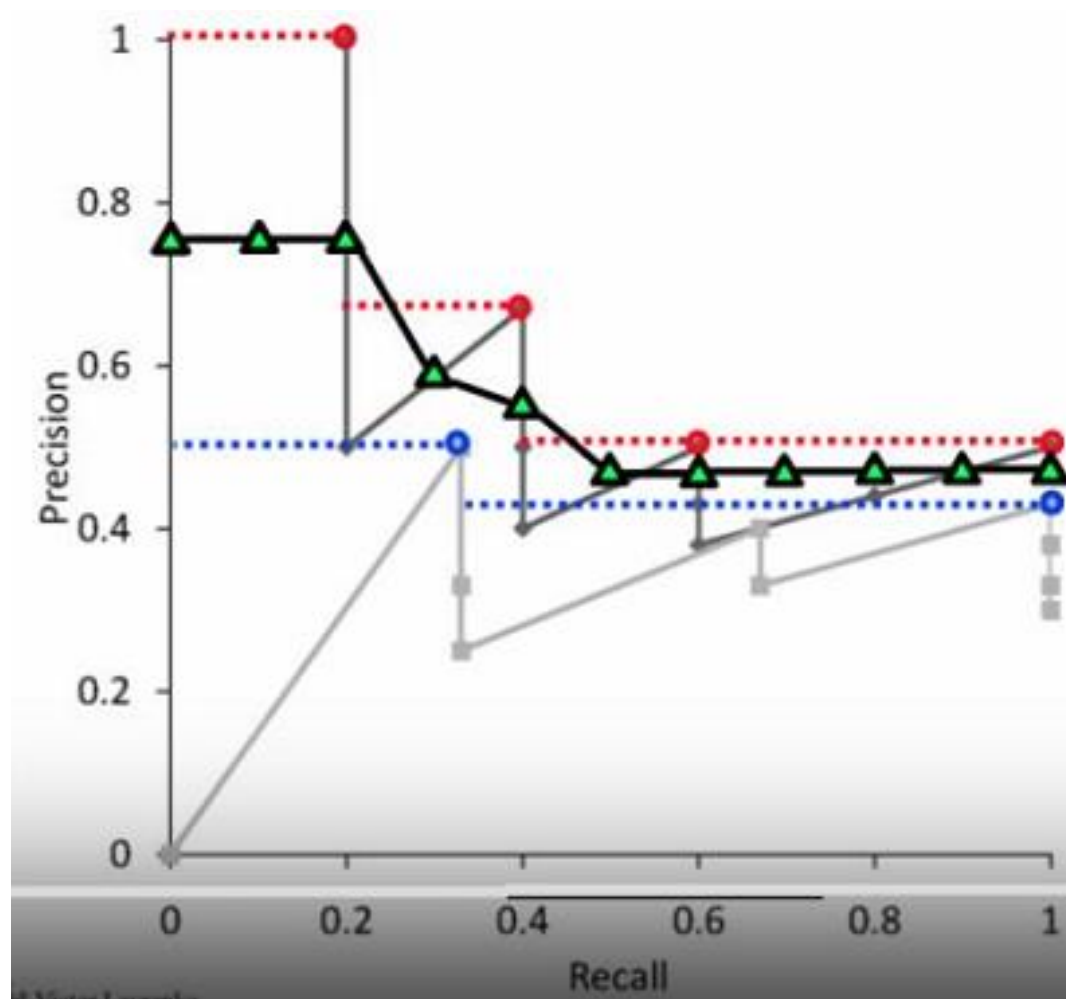


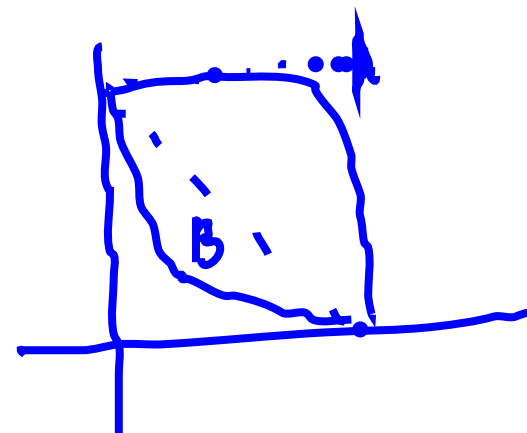
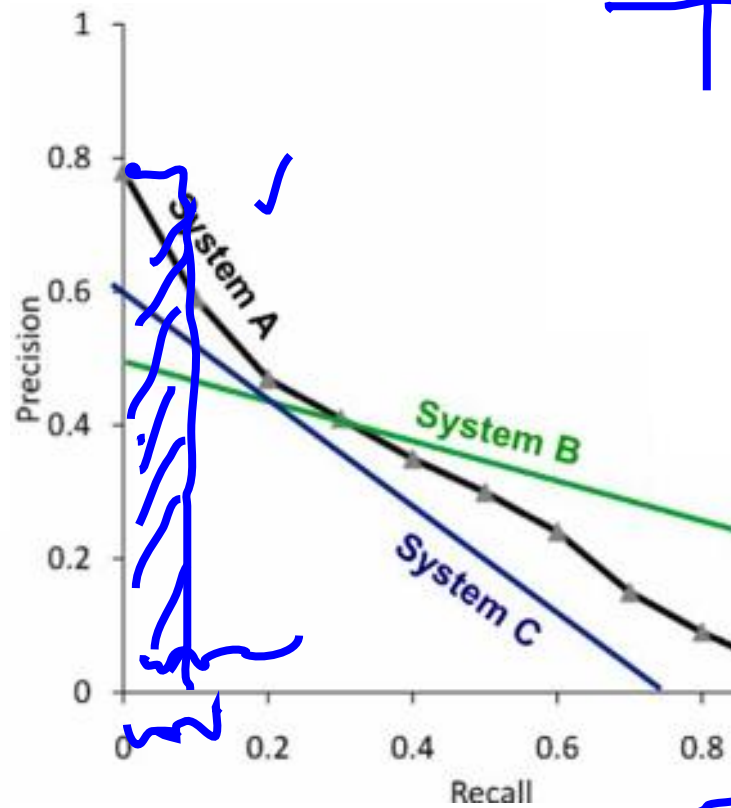
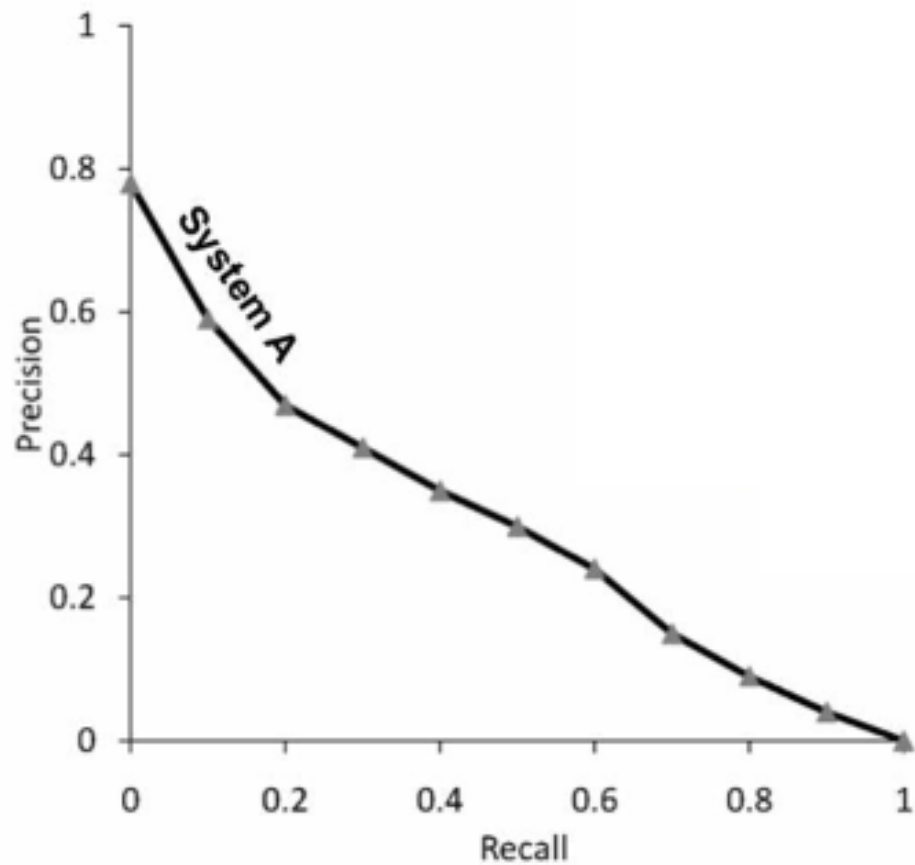
→ Ranking











→ Mean Reciprocal Rank
 → Mean Avg. Precision
 → πDCG