

Introduction to
Assignment Project Exam Help
Informa |
<https://eduassistpro.github.io/>
Add WeChat edu_assist_pro
Lecture 8: E

This lecture

- How do we know if our results are any good?
 - Evaluating a search engine
 - Benchmarks
 - Precision a

Assignment Project Exam Help

<https://eduassistpro.github.io/>

Add WeChat edu_assist_pro

Assignment Project Exam Help

<https://eduassistpro.github.io/>

Add WeChat edu_assist_pro

EVALUATING SEARCH ENGINES

Measures for a search engine

- How fast does it index
 - Number of documents/hour
 - (Average document size)
- How fast does it respond
 - Latency as a function of index size
- Expressiveness of query language
 - Ability to express complex information needs
 - Speed on complex queries
- Uncluttered UI
- Is it free?

Assignment Project Exam Help

<https://eduassistpro.github.io/>

Add WeChat edu_assist_pro

Measures for a search engine

- All of the preceding criteria are *measurable*: we can quantify speed/size
 - we can make expressiveness precise
- The key measures
 - What is this?
 - Speed of response/size of index
 - But blindingly fast, useless answers won't make a user happy
- Need a way of quantifying user happiness

Assignment Project Exam Help

<https://eduassistpro.github.io/>

Add WeChat edu_assist_pro

Measuring user happiness

- Issue: who is the user we are trying to make happy?
 - Depends on the setting
- Web engine:
 - User finds what they want to the engine
 - Can measure
 - User completes their task — means not end
 - See Russell <http://dmrussell.s.com/JCDL-talk-June-2007-short.pdf>
- eCommerce site: user finds what they want and buy
 - Is it the end-user, or the eCommerce site, whose happiness we measure?
 - Measure time to purchase, or fraction of searchers who become buyers?

Measuring user happiness

- Enterprise (company/govt/academic): Care about “user productivity”
 - How much time do my users save when looking for information?
 - Many other readth of access, secure access, etc.

Assignment Project Exam Help

<https://eduassistpro.github.io/>

Add WeChat edu_assist_pro

Happiness: elusive to measure

- Most common proxy: *relevance* of search results
- But how do you measure relevance?
- We will determine, then examine its issues
- Relevance measurement elements:
 1. A benchmark document collection
 2. A benchmark suite of queries
 3. A usually binary assessment of either Relevant or Nonrelevant for each query and each document
 - Some work on more-than-binary, but not the standard

Evaluating an IR system

- Note: the **information need** is translated into a **query**
- Relevance is assessed relative to the **information need** *not* the **Assignment Project Exam Help**
<https://eduassistpro.github.io/>
- E.g., Information need: *I'm information on whether drinking red wine i effective at reducing your risk of heart attacks than white wine.*
Add WeChat edu_assist_pro
- Query: **wine red white heart attack effective**
- You evaluate whether the doc addresses the information need, not whether it has these words

Standard relevance benchmarks

- TREC - National Institute of Standards and Technology (NIST) has run a large IR test bed for many years
- Reuters and other collections used
- “Retrieval tasks” specified
 - sometimes as queries
- Human experts mark, for each query and for each doc, Relevant or Nonrelevant
 - or at least for subset of docs that some system returned for that query

Unranked retrieval evaluation: Precision and Recall

- **Precision:** fraction of retrieved docs that are relevant
= $P(\text{relevant} | \text{retrieved})$
- **Recall:** fraction of relevant docs that are retrieved =
 $P(\text{retrieved} | \text{relevant})$

Add WeChat edu_assist_pro

	Relevant	Not Relevant
Retrieved	tp	fp
Not Retrieved	fn	tn

- Precision $P = \text{tp} / (\text{tp} + \text{fp})$
- Recall $R = \text{tp} / (\text{tp} + \text{fn})$

Should we instead use the accuracy measure for evaluation?

- Given a query, an engine classifies each doc as “Relevant” or “Nonrelevant”
- The **accuracy** of an engine: the fraction of these classifications
$$(tp + tn) / (tp + fp + fn + tn)$$
- **Accuracy** is a commonly used measure in machine learning classification work
- Why is this not a very useful evaluation measure in IR?

Assignment Project Exam Help

<https://eduassistpro.github.io/>

Add WeChat edu_assist_pro

Why not just use accuracy?

- How to build a 99.9999% accurate search engine on a low budget....

Assignment Project Exam Help

<https://eduassistpro.github.io/>

Add WeChat edu_assist_pro

0 matching results found.

- People doing information retrieval *want to find something* and have a certain tolerance for junk.

Precision/Recall

- You can get high recall (but low precision) by retrieving all docs for all queries!
- Recall is a ratio of the number of docs retrieved to the number of relevant docs

Assignment Project Exam Help

<https://eduassistpro.github.io/>
Add WeChat edu_assist_pro

- In a good system, precision decreases as either the number of docs retrieved or recall increases
 - This is not a theorem, but a result with strong empirical confirmation

Difficulties in using precision/recall

- Should average over large document collection/
query ensembles
- Need human relevance assessments
 - People aren't <https://eduassistpro.github.io/>
- Assessments have to be bin
 - Nuanced assessments?
- Heavily skewed by collection/authorship
 - Results may not translate from one domain to another

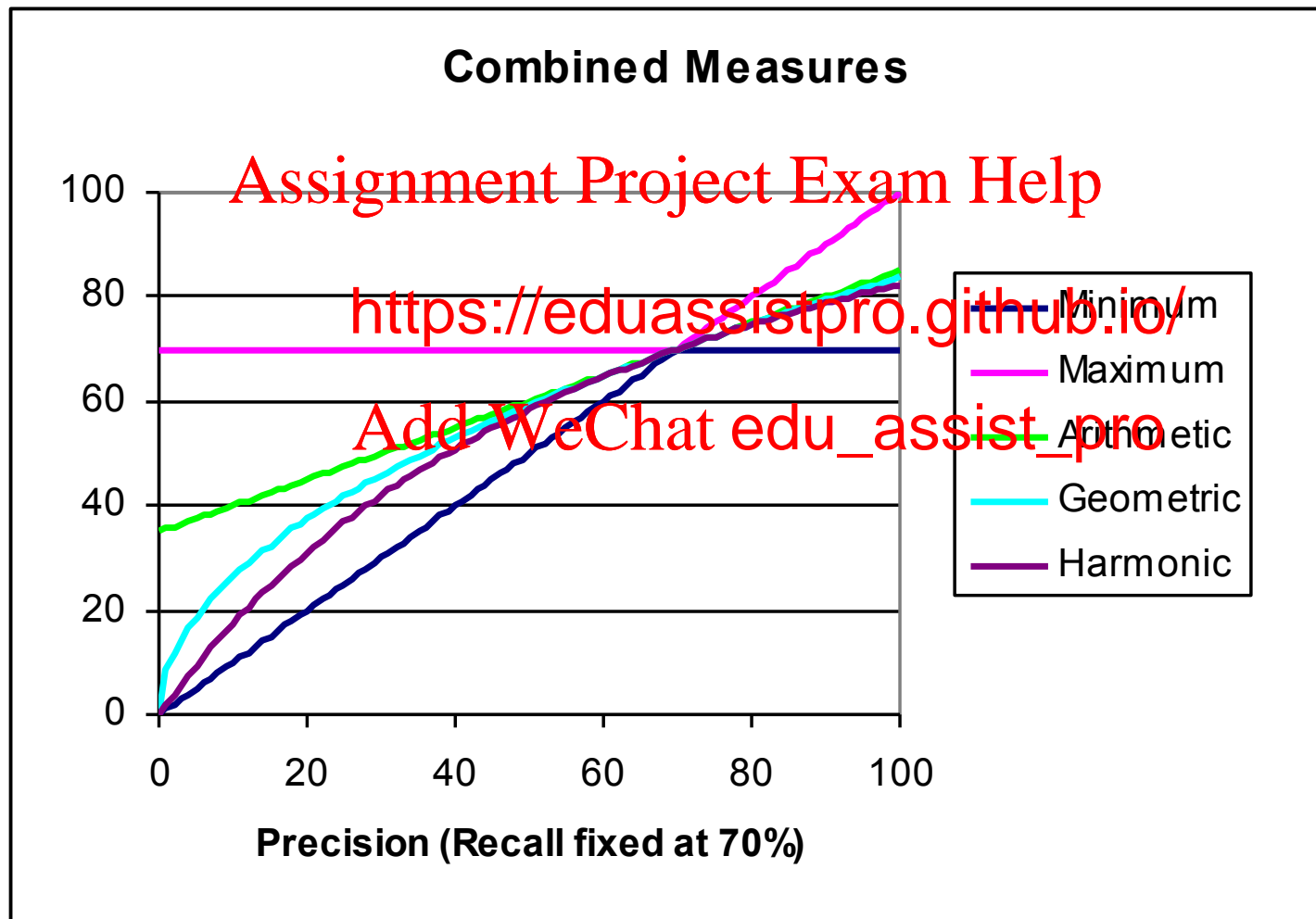
A combined measure: F

- Combined measure that assesses precision/recall tradeoff is **F measure** (weighted harmonic mean):

$$F = \frac{2\alpha R + (1-\alpha)P}{2}$$

- People usually use balanced F_1 measure
 - i.e., with $\beta = 1$ or $\alpha = \frac{1}{2}$
- Harmonic mean is a conservative average
 - See CJ van Rijsbergen, *Information Retrieval*

F_1 and other averages



Evaluating ranked results

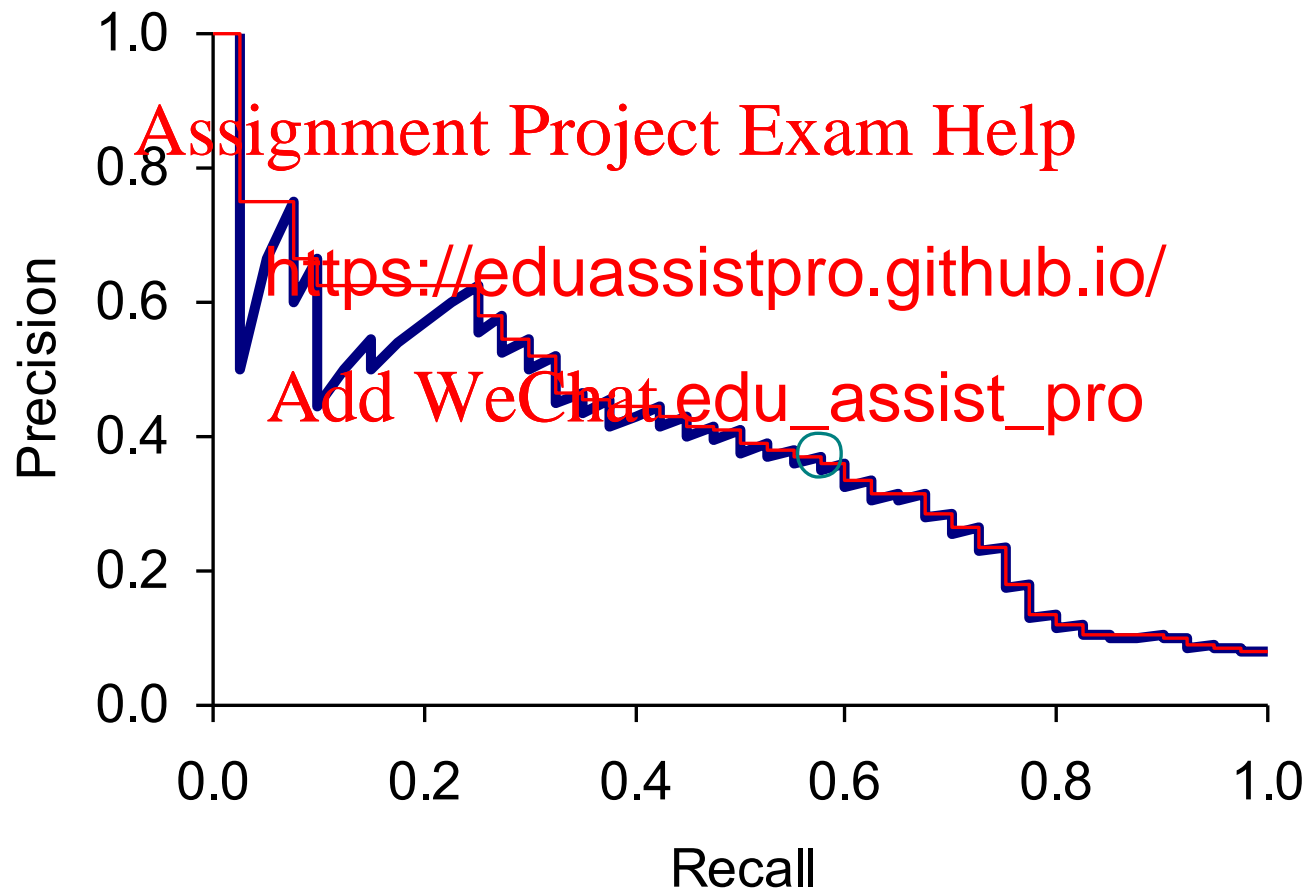
- Evaluation of ranked results:
 - The system can return any number of results
 - By taking various numbers of the top returned documents (levels of recall) produce a *precision-recall curve*

Assignment Project Exam Help

<https://eduassistpro.github.io/>

Add WeChat edu_assist_pro

A precision-recall curve



Averaging over queries

- A precision-recall graph for one query isn't a very sensible thing to look at
- You need to average performance over a whole bunch of queries
<https://eduassistpro.github.io/>
- But there's a technical issue
 - Precision-recall calculations points on the graph
 - How do you determine a value (interpolate) between the points?

Interpolated precision

- Idea: If locally precision increases with increasing recall, then you should get to count that...
- So you max of precisions to right of value

Assignment Project Exam Help

<https://eduassistpro.github.io/>

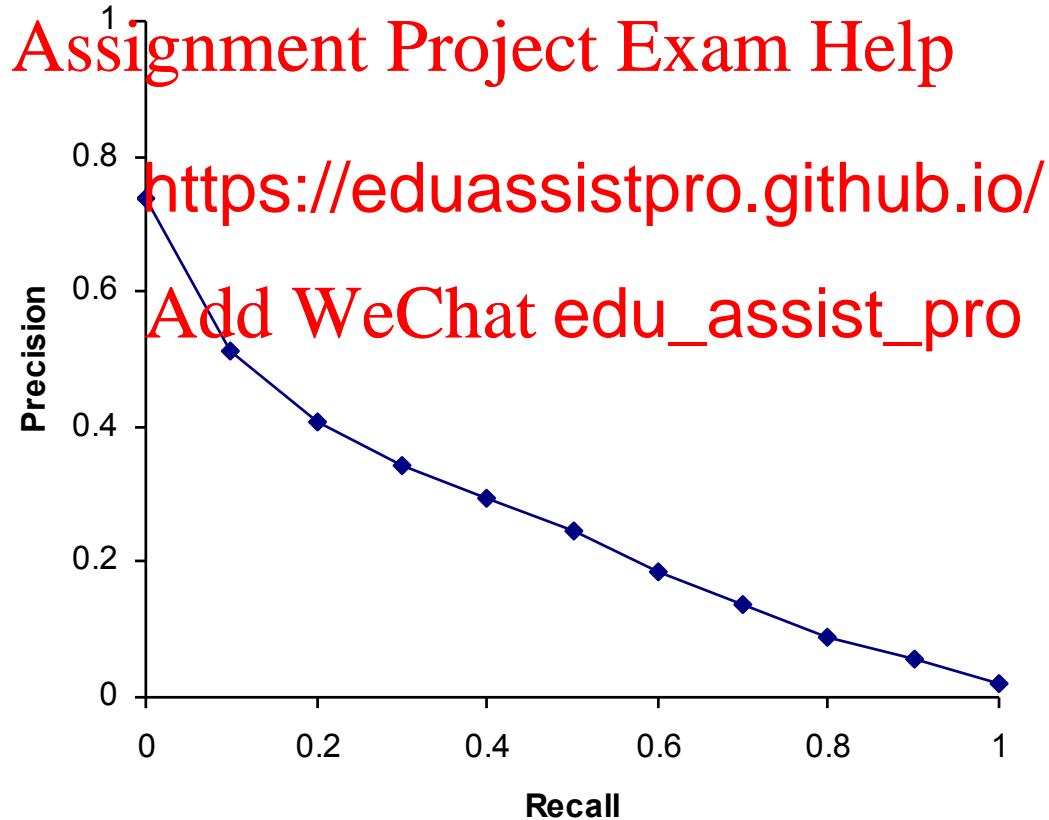
Add WeChat edu_assist_pro

Evaluation

- Graphs are good, but people want summary measures!
 - Precision at fixed retrieval level
 - Precision-at-k: Precision of top k results
 - Perhaps good match: all people want are results pages
 - But: averages badly and has parameter of k
 - 11-point interpolated average
 - The standard measure in the early TREC competitions: you take the precision at 11 levels of recall varying from 0 to 1 by tenths of the documents, using interpolation (the value for 0 is always interpolated!), and average them
 - Evaluates performance at all recall levels

Typical (good) 11 point precisions

- SabIR/Cornell 8A1 11pt precision from TREC 8 (1999)



Yet more evaluation measures...

- Mean average precision (MAP)
 - Average of the precision value obtained for the top k documents, each time a relevant doc is retrieved
 - Avoids interpolation levels
 - MAP for que ave.
 - Macro-averaging each query
- R-precision
 - If have known (though perhaps incomplete) set of relevant documents of size Rel , then calculate precision of top Rel docs returned
 - Perfect system could score 1.0.

Variance

- For a test collection, it is usual that a system does crummily on some information needs (e.g., MAP = 0.1) and excellently on others (e.g., MAP = 0.7)
- Indeed, it is usual that the variance in performance across queries is much greater than the variance between different systems on the same query.
- That is, there are easy information needs and hard ones!

Assignment Project Exam Help

<https://eduassistpro.github.io/>

Add WeChat edu_assist_pro

CREATING TEST COLLECTIONS FOR IR EVALUATION

Test Collections

Assignment Project Exam Help

<https://eduassistpro.github.io/>

Add WeChat edu_assist_pro

From document collections to test collections

- Still need
 - Test queries
 - Relevance assessments
- Test queries <https://eduassistpro.github.io/>
 - Must be germane to docs av
 - Best designed by domain ex
 - Random query terms generally not a good idea
- Relevance assessments
 - Human judges, time-consuming
 - Are human panels perfect?

Unit of Evaluation

- We can compute precision, recall, F, and ROC curve for different units.
- Possible units
 - Documents (<https://eduassistpro.github.io/>)
 - Facts (used in some TREC ev)
 - Entities (e.g., car companies)
- May produce different results. Why?

Assignment Project Exam Help

Add WeChat edu_assist_pro

Kappa measure for inter-judge (dis)agreement

- Kappa measure
 - Agreement measure among judges
 - Designed for
 - Corrects for
- $\text{Kappa} = [P(A) - P(E)] / [1 - P(E)]$
- $P(A)$ – proportion of time judges agree
- $P(E)$ – what agreement would be by chance
- Kappa = 0 for chance agreement, 1 for total agreement.

Kappa Measure: Example

P(A)? P(E)?

	Judge 2: Relevant	Judge 2: Nonrelevant
Judge 1: Relevant	300	20
Judge 1: Nonrelevant	10	70

Total assessment: 400

- $P(A) = 370/400 = 0.9250$
- $P(\text{nonrelevant}) = (10+20+70+70)/800 = 0.2125$
- $P(\text{relevant}) = (10+20+300+300)/800 = 0.7875$
- $P(E) = 0.2125^2 + 0.7875^2 = 0.6653$
- $\text{Kappa} = (0.9250 - 0.6653)/(1 - 0.6653) = 0.7759$

Kappa Example

- $P(A) = 370/400 = 0.9250$
- $P(\text{nonrelevant}) = (10+20+70+70)/800 = 0.2125$
- $P(\text{relevant}) = \frac{7875}{7875}$
- $P(E) = 0.2125$
- $\text{Kappa} = (0.9250 - 0.6653) / (1 - 0.6653) = 0.7759$
- $\text{Kappa} > 0.8 = \text{good agreement}$
- $0.67 < \text{Kappa} < 0.8 \rightarrow \text{"tentative conclusions"} \text{ (Carletta '96)}$
- Depends on purpose of study
- For >2 judges: average pairwise kappas

TREC

- TREC Ad Hoc task from first 8 TRECs is standard IR task
 - 50 detailed information needs a year
 - Human evaluation of pooled results returned
 - More recently , HARD
- A TREC query (TR <https://eduassistpro.github.io/>)
 - <top>
 - <num> Number: 225
 - <desc> Description:
What is the main function of the Federal Emergency Management Agency (FEMA) and the funding level provided to meet emergencies?
Also, what resources are available to FEMA such as people, equipment, facilities?
 - </top>

Standard relevance benchmarks:

Others

- GOV2
 - Another TREC/NIST collection
 - 25 million web pages
 - Largest collection
 - But still 3 orders of magnitude smaller than what Google/Yahoo/MSN in
- NTCIR
 - East Asian language and cross-language information retrieval
- Cross Language Evaluation Forum (CLEF)
 - This evaluation series has concentrated on European languages and cross-language information retrieval.
- Many others

Assignment Project Exam Help

<https://eduassistpro.github.io/>

Add WeChat edu_assist_pro

Interjudge Agreement: TREC 3

Assignment Project Exam Help

<https://eduassistpro.github.io/>

Add WeChat edu_assist_pro

Impact of Inter-judge Agreement

- Impact on **absolute** performance measure can be significant (0.32 vs 0.39)
- Little impact on **relative** performance <https://eduassistpro.github.io/>
- Suppose we want to know if algorithm A is better than algorithm B
- A standard information retrieval experiment will give us a reliable answer to this question.

Critique of pure relevance

- Relevance vs Marginal Relevance
 - A document can be redundant even if it is highly relevant
 - Duplicates
 - The same information from different sources
 - Marginal relevance of utility for the user.
- Using facts/entities as evaluation units more directly measures true relevance.
- But harder to create evaluation set

$$MMR \stackrel{\text{def}}{=} \text{Arg} \max_{D_i \in R \setminus S} \left[\lambda (\text{Sim}_1(D_i, Q)) - (1 - \lambda) \max_{D_j \in S} \text{Sim}_2(D_i, D_j) \right]$$

Can we avoid human judgment?

- No
- Makes experimental work hard
 - Especially on a large scale
- In some very simple cases
 - E.g.: for approximate vector comparison, we can compare the cosine distance of the closest docs to those found by an approximate retrieval algorithm
- But once we have test collections, we can reuse them (so long as we don't overtrain too badly)

Evaluation at large search engines

- Search engines have test collections of queries and hand-ranked results
- Recall is difficult to measure on the web
- Search engines often use $nDCG_k$, e.g., $k = 10$
- ... or measures that compare getting rank 1 right to getting rank 10 right.
 - **NDCG** (Normalized Cumulative Discounted Gain)
- Search engines also use non-relevance-based measures.
 - Clickthrough on first result
 - Not very reliable if you look at a single clickthrough ... but pretty reliable in the aggregate.
 - Studies of user behavior in the lab
 - A/B testing

A/B testing

- Purpose: Test a single innovation
- Prerequisite: You have a large search engine up and running.
- Have most users use old system
- Divert a small proportion of traffic (e.g., 1%) to the new system that includes a new feature
- Evaluate with a controlled experiment: e.g., clickthrough on first result
- Now we can directly see if the innovation improves user happiness.
- Probably the evaluation methodology that large search engines trust most
- In principle less powerful than doing a multivariate regression analysis, but easier to understand

Assignment Project Exam Help

<https://eduassistpro.github.io/>

Add WeChat edu_assist_pro

Assignment Project Exam Help

<https://eduassistpro.github.io/>

Add WeChat edu_assist_pro

RESULTS PRESENTATION

Result Summaries

- Having ranked the documents matching a query, we wish to present a results list
- Most commonly, a list of the document titles plus a short summary

Assignment Project Exam Help

<https://eduassistpro.github.io/>

Add WeChat edu_assist_pro

Resources for this lecture

- IIR 8
- MIR Chapter 3
- MG 4.5
- Carbonell and <https://eduassistpro.github.io/> se of MMR, diversity-based reranking, ng documents and producing summaries.

Assignment Project Exam Help

<https://eduassistpro.github.io/>

Add WeChat edu_assist_pro