

Improving Accuracy of Relevance Assessment for Math Search using Rendered Expressions

by

Matthias Reichenbach

Project submitted in Partial Fulfillment of the Requirements for the
degree of Master of Science in Human-Computer Interaction

Rochester Institute of Technology

B. Thomas Golisano College

of

Computing and Information Sciences

Department of Information Science and Technology

May 2013

Rochester Institute of Technology
B. Thomas Golisano College
of
Computing and Information Sciences

Master of Science in Human-Computer Interaction

Project Approval Form

Student Name: Matthias Reichenbach

Project Title: Improving Accuracy of Relevance Assessment for Math Search
using Rendered Expressions

Project Committee

Name

Signature

Date

Dr. Richard Zanibbi

Chair

Dr. Evelyn Rozanski

Committee Member

Dr. Michael Yacci

Committee Member

Acknowledgments

I am grateful to several people that helped me in different ways in my path to complete this project. I want to thank my advisors for their dedication to the project, and specially Dr. Richard Zanibbi for providing invaluable support and guidance, Dr. Anurag Agarwal for his help in defining first the problem and then the search tasks and Katherine Zanibbi for informing the initial experimental design.

Finally I am want to thank the U.S. Department of State and, more specifically, the Fulbright Scholarship Program for allowing me to pursue my desire to learn more about the Human-Computer Interaction field.

Abstract

Improving Accuracy of Relevance Assessment for Math Search using Rendered Expressions

Matthias Reichenbach

Supervising Professor: Dr. Richard Zanibbi

Finding ways to help users assess relevance when they search using math expressions is critical for making Math Information Retrieval (MIR) systems easier to use. We designed a study where participants completed artificial search tasks involving mathematical expressions, in one of two different summary styles and across two different information needs, and measured response time and relevance assessment. One summary style was based on Google's result summaries with the math expressions in it rendered as math and one with the original linearized text form of the expressions used as control. We found that participants with the rendered summary style performed significantly better. On average, they had an assessment accuracy 17.18% higher and reported having fewer problems reading the results than participants in the control summary style. This suggests that search engines would benefit from properly rendering math expressions in their summaries and opens the possibility of showing more non text-based extracts of relevant information as a means to increase accuracy in relevance assessments.

Contents

Acknowledgments	iii
Abstract	iv
1 Introduction	1
1.1 Problem	1
1.2 Related Work	3
1.3 Existing MIR Solutions	5
2 Methodology	8
2.1 Experimental Design	8
2.2 Participants	9
2.3 Independent Variable Levels	10
2.3.1 Information Needs	10
2.3.2 Summary Styles	11
2.4 Procedure	13
2.4.1 Tasks	15
2.4.2 Final Questionnaire	15
2.5 Materials	17
3 Results	19
3.1 Demographics	19
3.2 Experiment	20
3.3 Exit Questionnaire	24
3.4 Summary	26
4 Discussion	27
4.1 Learning Effects	27
4.2 Summary Style Effect	28
4.3 Information Need Effect	29

4.4 Summary	30
5 Conclusions	31
Bibliography	33
A Screening Questionnaire	35
B Post-test Questionnaire	39
C Recruiting Poster	41
D Recruiting Email	43
E Pre-generated Hit Results	44
E.1 Familiarization Task	45
E.2 Informational Need Task	46
E.3 Resource Need Task	48

List of Tables

2.1	Experimental group arrangement	9
2.2	Scenarios used to prompt participants to each information need.	10
2.3	Hit Summary Generation Procedure	13
2.4	Scenario and query terms used in the familiarization task.	14
2.5	Each of the counterbalanced orders used to present the tasks' search hits. . .	16
3.1	Mean and standard deviation for each summary style group and information need	23

List of Figures

1.1	Summary styles compared by Aula [1].	3
1.2	Sample results from searches with a math expression.	5
2.1	Example of the two summary styles used. From top to bottom: summary style used as control obtained from Google search (SS1) and summary style that shows the same information but with the math expressions properly formatted (SS2).	12
2.2	Interface of experiment website.	17
3.1	Distribution of the highest level of education achieved among the participants.	20
3.2	Distribution of participant's reported frequencies for the need to express mathematical notation in a computer.	20
3.3	Mean response time by hit position in log scales for SS2 with a trend line. .	21
3.4	Mean accuracy by log of response time. Mean accuracy values aggregate the accuracy of all responses that fall between integer values of time in the log scale.	21
3.5	Average accuracy achieved by all participants for each of the hits. Hits 1-10 are from the informational need and hits 11-20 are for the resource need. . .	22
3.6	Average time taken to assess relevance by all participants for each of the hits. Hits 1-10 are from the informational need and hits 11-20 are for the resource need.	22
3.7	Profile graph showing Accuracy by Information Need and Summary Style .	23
3.8	Profile graph showing Time by Information Need and Summary Style . . .	23
3.9	Participants response to the statement "I'm familiar with the math involved in these tasks".	25
3.10	Participants response to the statement "I have had information needs similar to the tasks I just completed"	25
3.11	Participants response to the statement "I had no problems reading the results presented"	25
3.12	Profile graph showing Time by Information Need and Summary Style . . .	25

Chapter 1

Introduction

In her authoritative book, Hearst [7] defines the objective of search user interfaces is to “aid users in the expression of their information needs, in the formulation of their queries, in the understanding of their search results, and in keeping track of the progress of their information seeking efforts.”

In the context of Math Information Retrieval (MIR) systems, previous research has been done in developing systems to help formulate queries that include mathematical expressions and search engines that understand them. Zanibbi and Blostein [14] provide a comprehensive overview of current techniques for recognition and retrieval of mathematical expressions. These have been used by, for example, Sasarak et al. [11] who designed a web based interface where users can draw the math expressions they want to search without requiring them to know math coding languages such as LaTeX or MathML.

Less research has been done in helping users understand the search results they get from MIR systems. This is an important area of research because it can improve users success rate and satisfaction with the systems. In section 1.2 we provide a summary of related work about this area for text-based search systems and its relation to MIR systems. In the following section we describe the specific problem our study addresses.

1.1 Problem

Previous research has described the different information needs from users regarding search of math expressions [15]. If a Math Information Retrieval (MIR) system is to address all of

these different needs effectively, it should present a specialized interface for each. Specifically, the summary of each hit may need to be formatted or styled differently depending on the information need in order to better help the user assess the relevance of each result.

For example, users searching for a theorem proof may benefit more from seeing the title of the section where the expression appears than users searching for example uses of expressions. The same idea can be extended to all the different types of users' information needs. Zaho et.al. [15] separates them into informational needs (names/alias, definition, derivation, etc.) and resource needs (papers, tutorials, slides, books, etc.). They additionally provide categories (definition, example, proof, etc.) into which sections of math related documents can be classified and suggest that MIR systems should display only the categories the user is searching for. Pattern recognition systems would be able to determine the information need and the desired result category from the query expression. However, these are not clear-cut categories and some queries will fall into more than one.

Different ways of presenting results (summary styles) has been shown to have an effect on the ability of users to assess relevance [1]. In addition, specific styles of result summaries work best for different information needs. For example, when users want to find some piece of information — called an information need in [3] — the search is better served with longer result summaries [6]. On the other hand, when users want to find a specific website or resource — called a navigational need in [3] — it is better to show short summaries [6].

To our knowledge, no previous research has been carried out to help define the most appropriate summary styles for different information needs when using MIRs. The ability to adapt the summary style of the results depending on the information need can be critical in a MIR system that wants to maximize the users ability to distinguish relevant from irrelevant results. We tried to shed light into this problem by comparing two different summary styles and two different information needs. Specifically, we wanted to find if there is a difference in relevance assessment when the matched math expressions in the summary are displayed linearized as text and when they are correctly formatted as math (see Figure 2.1).

List:

[HCI 2004 Design for Life: Upcoming Deadline: 7th May 2004...](#)

- Annual Conference is taking place at **Leeds** Metropolitan University...
 - May 7th **2004** is the deadline for industry...
 - ...concerns are traditional ones for **HCI**; others are...
- www.chiplace.org/modules.php?op=modload&name=News&file=article&sid=237

Normal-bolded:

[HCI 2004 Design for Life: Upcoming Deadline: 7th May 2004 ...](#)

... Annual Conference is taking place at **Leeds** Metropolitan University ... May 7th **2004** is the deadline for industry ...
concerns are traditional ones for **HCI**; others are ...

www.chiplace.org/modules.php?op=modload&name=News&file=article&sid=237

Normal-plain:

[HCI 2004 Design for Life: Upcoming Deadline: 7th May 2004 ...](#)

... Annual Conference is taking place at Leeds Metropolitan University ... May 7th 2004 is the deadline for industry ...
concerns are traditional ones for HCI; others are ...

www.chiplace.org/modules.php?op=modload&name=News&file=article&sid=237

Figure 1.1: Summary styles compared by Aula [1].

1.2 Related Work

Previous research on summary styles has been focused on text search engines. These usually consist of informal and formal usability tests where users are given artificial tasks. For example Aula [1] compared different result modalities (see Figure 1.1) by asking users to select the link with the best answer in the shortest amount of time from a list of predefined results with only one having the answer. No significant differences were discovered in error rates, but there were differences in response time. Presenting the results summaries as a list of sentences was shown to significantly improve performance over presenting the sentences as a paragraph. However, making the query terms bold in the summary was shown to have no effect and in some cases even hindered time performance.

Kickmeier and Albert [8] found that changing the density of salient words in a result summary had an effect on participant's ability to determine relevancy. In their experiment, 550 participants performed two tasks where they had to assess relevancy of one result summary. They modified the density, defined as the proportion of salient words in the summary, by making random words in the paragraph bold and measured response time and accuracy of the answer. Results showed that the best performance was obtained when the

density of the bolded word was between 10% to 25%, without regard to what words are the ones being bolded.

Previous research has shown that the order in which the results are displayed also affects how likely they are to be selected as relevant [5]. Participants were asked to select the most relevant result from a list of 10 results for 12 search tasks. The most relevant result was moved between six positions (1, 2, 4, 5, 7 and 8) and the response time and accuracy were measured. They found an effect for the position for both accuracy and time, with accuracy dropping from 83% when the most relevant result was in position 1 to around 11% when the most relevant result was in position 8.

Cutrell and Guan [6] found that different information needs are better served with summaries of different lengths. Participants performed search tasks with varying difficulty and with navigational or informational needs. The time to complete the tasks was measured and compared across conditions. They found participants performed navigational tasks better with short summaries, but performed better with long summaries in the informational tasks.

In their seminal research, Tombros and Sanderson [12] showed the advantages of displaying the matched query terms with their document context in the summaries — called query-biased result summaries. The query biased summary was constructed by combining the top ranking sentences based on their query term frequency. Participants rated a list of 50 documents as relevant or not for a specific query with a time limit of 5 minutes. Participants in the query-biased condition rated more documents on average in the time limit, had to consult the linked documents less and had a higher precision rate than the participants in the pre-defined summaries condition.

These results, while important, cannot be directly applied to MIR systems. Youssef [13] provides a possible implementation of hit content summarization for MIR systems. His system fragments each document into small units (e.g.: equations, sentences, tables and graphs) and indexes each separately. The document summaries are then constructed with the top matching fragments with respect to the query. He proposes that a user study should be conducted to test the system's usability. This gap is what this project address,

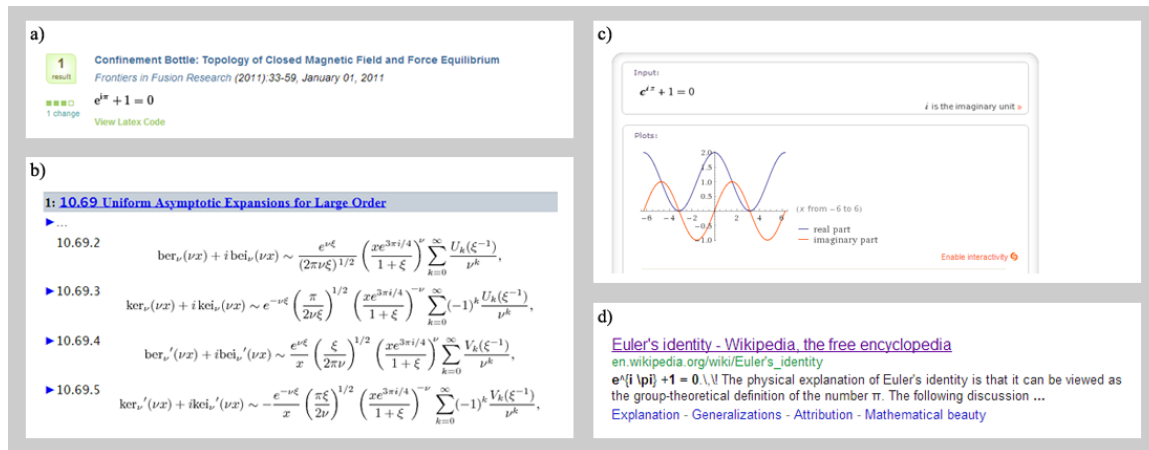


Figure 1.2: First result returned by the following search engines when searching for $e^{i\pi} + 1 = 0$ as of 01/22/2013: a) latexsearch.com, b) dlmf.nist.gov, c) wolframalpha.com and d) google.com. LaTeX syntax was used to specify the query: $e^{i\pi} + 1 = 0$

trying to shed light onto defining the appropriate summaries styles for MIR systems based on the user's information need.

1.3 Existing MIR Solutions

The following is a short discussion of existing MIR solutions and other related systems that can be used to search for expressions. These systems differ, among other things, in their interpretation of math expressions in the queries and the amount of context surrounding the matched terms in the query-biased summaries. Query-biased summaries are summaries that are modified depending on the query terms. They usually contain the query terms in the context in which they appear in the document. [12]

LaTeX Search. This search service is provided by Springer with the intention of allowing researchers to search through their scientific publications for mathematical expressions. Each result consists of the title and other general information of the matched document, the matched expression(s) in the document and a ranking of similarity (see Figure 1.2.a).

The result's summary style can be classified as a query-biased with expression level context. No surrounding or additional content from the document is shown together with

the matched expression.

Digital Library of Mathematical Functions. This project was started to revise and digitize Abramowitz and Stegun's Handbook of Mathematical Functions with Formulas, Graphs, and Mathematical Tables[2]. When searching, it displays the matching section's title, along with the up-to 5 top matching expressions in it [9]. If query terms other than math expressions are used, they are also used in matching and relevance ranking.

The result's summary style can be classified as query-biased with sentence level context. Sentences are only shown if the matching search terms (words or expressions) are in a sentence (see Figure 1.2.b).

Wolfram Alpha. Wolfram Alpha is a question answering system not a search engine. Instead of searching through a database of indexed documents it tries to interpret the query and produce an answer based on its knowledge base. It is able to synthesize different data sources and combine them in different ways, such as plots and tables.

However, it can still be used to gain general information about mathematical expressions. The results page shows an interpretation of the query followed by information derived from Wolfram Alpha's knowledge base. For mathematical expressions it may include graphs, integrals, derivations and more (see Figure 1.2.c).

Google. Google's search engine is not designed to handle math expressions. It interprets terms as text query terms and as such understands no mathematical rules. It will produce no matches unless the actual string representation of the math expression is the same in the query as in the document. Still, it can be used to search for specific math expressions given that the textual representation of the expression (e.g. in LaTeX or MathML) in the query matches what is indexed from the documents.

The results consist of the matched page's title and a query-biased summary with two or three lines of total context extracted from the document. The context can consist of text surrounding the matched query terms or other parts of the document (see Figure 1.2.d).

Each of the systems presented employs a different strategy for generating and presenting their hit summaries. Since their indexed documents differ greatly, a direct comparison

of user preference and performance is hard. However, a study that uses artificial tasks and pre-generated results can modify the summary style to test for their effect.

This is the study we describe in the following sections. We measured relevance assessment accuracy and response time from participants across different conditions. We compared two different summary styles: one with rendered math expressions and one without. Similar to what Aula [1] did, we created pre-defined search tasks and search results, formatted in the two different summary styles.

To test for differences similar to what Guan and Cutrell [6] found for text search, we also compared two different information needs, based on the distinctions by Zhao, et al. [15] of mathematical information needs. To avoid the ordering effects found by Guan and Cutrell [5], and similar to what Kickmeier and Albert [8] did, we presented the search results one by one in a counter-balanced order across participants. The following chapter describes the methodology of our experiment in detail.

Chapter 2

Methodology

The goal of this study was to define summary styles that help users assess relevance of searches that include mathematical expressions. Two different summary styles were considered: one similar to a traditional Google search result and one based on it but with the math expressions in it properly formatted. It is hypothesized that the properly formatted expressions should help with readability and thus allow users to assess relevance faster, similar to what Aula [1] found for text search.

In addition, we wanted to test if participants' ability to assess relevance was affected by their information need. We considered two different types of information needs based on the categories defined by Zhao et al. [15]. One of the information needs, the "informational need", gives more importance to the content of the search while the other, the "resource need" is more focused on the form. We hypothesized that better readability should have a larger effect on the informational need.

2.1 Experimental Design

Participants were divided into two groups based on each of the summary styles in our test. All participants performed three tasks: one familiarization task and one for each type of information need, presented in a counterbalanced order (see Table 2.1). By exposing each participant to only one summary style, this design allowed us to avoid showing the same participants the same hit result more than once, without having to design multiple tasks of comparable difficulty.

	SS1	SS2
IN1 - IN2	A ₁	B ₁
IN2 - IN1	A ₂	B ₂

Table 2.1: Distribution of the tasks between groups A and B. Each group is divided into two subgroups to counterbalance for presentation order of the information needs. INx represents the information needs and SSx represents the summary styles.

The experiment design conforms to a mixed factorial design where the summary style condition was between subjects and the information need condition was within subjects. Both of these conditions represent the independent variables of our study and they are explained in detail in section 2.3.

The measured variables, or dependent variables (DV), were the response time the participants took to assess if each hit is relevant to the task's information need and the accuracy of their assessment. These variables allowed us to measure if there was any difference in performance and were used to represent users' experience in a real MIR system.

2.2 Participants

Participants were recruited to be graduate and undergraduate RIT students both male and female. They completed a pre-screening questionnaire to assess if they met the required level of math proficiency — defined as having completed two or more math college courses — and experience with computer systems and search engines. Participants were also required to have normal or corrected to normal vision and hearing. Each participant was offered \$10.00 for their participation in the study for 30 minute sessions.

They were recruited through posters and emails as well as from Calculus 1 to 3 and Linear Algebra classes taught by Prof. Agarwal from the Department of Mathematics. Posters were added in the appropriate boards around campus and emails were sent to RIT's mailing list of graduates and undergraduates in the College of Science and the Golisano College of Computing and Information Sciences. Participants were scheduled to perform the experiment one by one during 30 minute time slots mostly between 2:00 PM and 6:00

Informational Need (IN1)	Resource Need (IN2)
Search for a proof [15].	Search for a tutorial [15].
<i>You have just finished attending a Linear Algebra class. Today's topic involved finding the inverse matrices through their adjoint matrix, but the professor did not explain how the formula $A^{-1} = \frac{1}{\det A} \cdot \text{adj}A$ was derived and you want to find that out. You go to a math search engine and search for '$A^{-1} = \frac{1}{\det A} \cdot \text{adj}A$ proof'.</i>	<i>Your friend is having trouble understanding derivatives of polynomials and you have agreed to help him. You need to be prepared to explain that to him so you want to find tutorials showing $\frac{d}{dx}ax^b = abx^{b-1}$. You go to a math search engine and search for '$\frac{d}{dx}ax^b = abx^{b-1}$ tutorial'.</i>

Table 2.2: Scenarios used to prompt participants to each information need.

PM in the Usability Testing Lab or the Eye-Tracking Lab at RIT.

Appendix C contains the poster used to recruit participants and appendix D the email. Appendix A shows the screening questionnaire. Participants that answered the question "please indicate how many courses have you taken about mathematics" with the option "0-1" were not selected for the study. The question "How frequently do you need to express mathematical notation when using a computer" was used to balance participants' skill at writing and understanding math code (e.g. LaTeX) across groups.

2.3 Independent Variable Levels

2.3.1 Information Needs

We used two different information needs based on the distinction made by Zhao et al. [15] between informational needs and resource needs. For each level of this IV an associated task that relates to a specific information need was created and participants were prompted to the information need by means of a short scenario. (See Table 2.2 for the paragraphs used.)

The first information need (IN1) was the informational need. These needs arise from users that need additional information related to some knowledge they lack. Examples of these needs are searches for names/aliases, definitions, derivations, explanations, examples,

problems/solutions, graphs/charts, algorithms, applications and related entities [15]. The task related to this need in our experiment was to search for the proof of an expression (see Table 2.2). The phrasing in IN1 emphasize the importance of finding a proof while leaving in what resource type unspecified.

The second information need (IN2) was the resource need. These needs arise from users that want to find a specific type of resource to support their need. Examples are searches for papers, tutorials, slides, course websites, books, code, toolkits and data [15]. The task related to this need in our experiment was to search for a tutorial (see Table 2.2). The phrasing of IN2 makes it clear that the main objective of the search is to find a tutorial and leave the specifics of its content as secondary.

The search results for the information needs were selected from the results page of Google after searching with the task's query terms. The expressions in the query terms were converted to their LaTeX representation and stripped of special characters to make it suitable for Google search. A hit was considered relevant if, and only if, it contained at least some portion of the query expression and the query term. Five hits matching this criteria were selected from the search results. Non relevant hits were selected from search hits that did not contain the query expression but did contain some other expression. In some cases, additional searches were made to generate summaries that met this criteria.

Both tasks were designed to require a level of math understanding that could be met by most second year students that have had two or more college level math courses. Participants were expected to have background knowledge of the general topics of the task to help them in assessing relevance.

2.3.2 Summary Styles

Two summary styles were used corresponding to the two levels of our summary styles independent variable. The first level (SS1) was used as a control. The hit results were styled based on how they appeared in the Google's results page, effectively using it as the "gold standard" (see Figure 2.1). Removing the result's URL and any other links besides

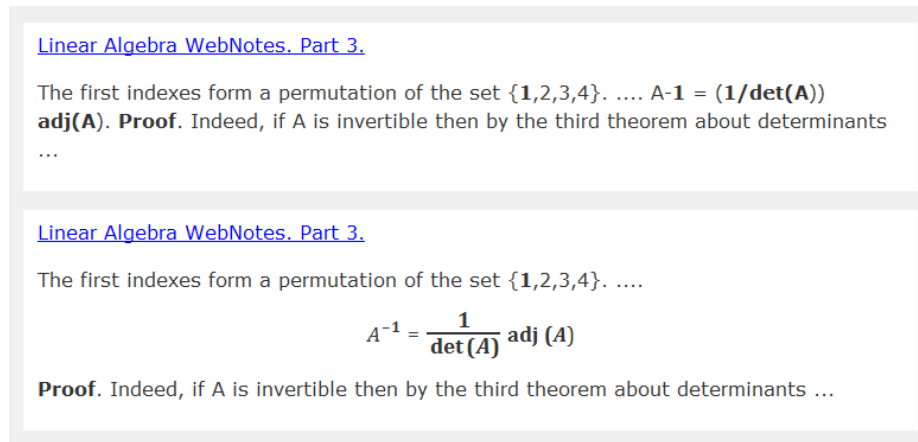


Figure 2.1: Example of the two summary styles used. From top to bottom: summary style used as control obtained from Google search (SS1) and summary style that shows the same information but with the math expressions properly formatted (SS2).

the title were the only modifications to the original summaries. Table 2.3 describes the specific procedure used for creating this summary style.

The second level (SS2) was our experimental condition. SS1 was used as the base for SS2, but with every math expression in it properly formatted (see Figure 2.1). Expressions in the result summaries were converted from their original code (e.g. LaTeX) when available, or visually when not, to MathML — a W3C standard for describing mathematical notation in XML — using MathJax¹. The converted code was then rendered in our experiment website by Mozilla Firefox’s native MathML rendering engine.

Aula [1] found that modifying the presentation style of summaries had an effect on the participants’ ability to determine relevancy for text search. Our expectation is that the modified summary styles in our experiment will similarly produce a difference in relevance assessment. Figure 2.1 shows examples of result summaries in each of the summary style levels and table 2.3 has the detailed description of how each one is created.

Previous research showed that making any part of the result summary text bold can affect relevancy evaluations [8] and that in some cases making matching terms bold has no effect when compared to not making them bold [1]. For this reason this initial study did

¹<http://www.mathjax.org/>

Hit Summary Generation Procedure	
Summary style 1	Summary style 2
<ol style="list-style-type: none"> 1. Obtain the result summary for the document as generated by Google search 2. If the document is relevant but does not contain any expressions in the summary, force Google to generate a summary with the expression using quotes around the expression and the "site:" operator. 3. Remove the URL 4. Remove any smart links² or any other extra links besides the document title 5. Keep the same words bolded 	<ol style="list-style-type: none"> 1. Start with the SS1 2. Render the math expressions <ul style="list-style-type: none"> • If the summary contains LaTeX code or other, use that to convert to MathML • If not, search for the original expression in the document and use that • If the expression is not in LaTeX or other standard encoding, translate to MathML visually replicating the expression's structure in the original document. 3. Keep only the exact portion of the expressions shown in SS1

Table 2.3: Procedure to build summary styles one and two for a given search result

not test for the effect of making the query terms bold in the summaries and instead kept the same parts bold from the original Google summary in both summary styles.

2.4 Procedure

Participants selected for the experiment were scheduled to meet one-on-one with the experimenter during predetermined time slots, mostly between 2:00 PM and 6:00PM. The meetings took place during two weeks, with 8 sessions scheduled per day

The experiment was performed in the Usability Testing Lab and the Eye Tracking Lab in the Golisano building in RIT. Once there, participants were instructed to take the seat in front of the computer fitted for the experiment. A consent form was handed to them in the

Familiarization Task
<p>Your classmate has heard of Pascal’s triangle but doesn’t understand how it relates to math. You want to find one or more resources to help explain to your classmate how the equation $(x + y)^2 = x^2 + 2xy + y^2$ relates to Pascal’s Triangle. You go to a search engine and search using the following keywords ‘Pascal triangle $(x + y)^2 = x^2 + 2xy + y^2$’.</p>

Table 2.4: Scenario and query terms used in the familiarization task.

entrance, with five minutes scheduled for them to read it and provide consent. During this time the experimenter answered any questions the participants had making sure that they understood the experiment.

Once the previous period was over the experimenter described one more time the experiment and the tasks the participants had to perform. He restated that the evaluation is of the system they were testing, not of them. Finally the participants were told to read the familiarization task and follow the instructions on the screen. The familiarization task had the same structure as the experimental tasks (see Section 2.4.1) but with only four result hits. (See the Table 2.4 for the scenario used for this task.) Similar to Aula [1], they were verbally told to “respond as quickly as possible, but take your time to make sure that you carefully consider whether a search result is relevant before you click Yes or No, even if it takes you longer than it usually does when you search, that is fine.”

The experiment’s website then guided the participants through the tasks. It started with the familiarization task — with summaries in their group’s summary style — to expose the participants to the system and get them to see how it works, followed by a time to ask questions. After this time the experimenter stated that he won’t be able to answer any more questions because the tasks are timed. After the familiarization task, the website showed them each of the experiment’s tasks in their group’s order. Participants were again asked to respond as quickly as possible, but take their time to make sure that they carefully

consider whether a search result is relevant before you click Yes or No, both verbally and with written instructions on-screen.

After finishing the tasks, participants were taken to an online questionnaire. It was designed to measure subjective responses to the system, the summary styles and the tasks. Before leaving, participants were given \$10.00 as compensation for their time.

2.4.1 Tasks

Each task in the experiment corresponded to one of the levels of the information need independent variable or the familiarization task. They started with a short description of the information need (see Table 2.2) and the pre-defined query used to meet the information need and generate the results. Participants were asked to read the tasks and, when ready, click the start button. At this moment the system started measuring response times and relevance assessments. Each of the hit results related to the search were displayed one at a time.

The presentation order was counterbalanced among participants to minimize ordering effects. Each hit for the tasks was presented in one of the 10 positions at least once across all participants and was preceded and followed by a different hit each time. Table 2.5 shows each of the 10 different orders used.

Participants were asked to determine if each hit is relevant or not to the information need, similar to [8]. Time was measured from the moment they saw the hit to the moment they made their relevance assessment by clicking the respective button. After each decision the participant was shown a new hit until they rated a total of ten hits for the experimental tasks and four for the familiarization task.

2.4.2 Final Questionnaire

The post-test questionnaire was administered through Google Forms immediately after the completion of the tasks. It consisted of several questions that asked the participants to rate, on Likert scales, their subjective impressions of the tasks. Importantly, it asked them to

Hit presentation order									
1	2	10	3	9	4	8	5	7	6
2	3	1	4	10	5	9	6	8	7
3	4	2	5	1	6	10	7	9	8
4	5	3	6	2	7	1	8	10	9
5	6	4	7	3	8	2	9	1	10
6	7	5	8	4	9	3	10	2	1
7	8	6	9	5	10	4	1	3	2
8	9	7	10	6	1	5	2	4	3
9	10	8	1	7	2	6	3	5	4
10	1	9	2	8	3	7	4	6	5

Table 2.5: Each of the counterbalanced orders used to present the tasks' search hits.

compare both tasks directly and to choose which one they found easier to complete (see Appendix B).

On the first question in the questionnaire, participants were asked to state how strongly they agreed or disagreed with several statements. Statements 1, 3 and 4 tested if there was a difference in how participants perceived the level of difficulty of the tasks. Statements 2 and 5 tested the participants' familiarity with the math used and the information needs expressed in the tasks.

A difference in the answer pattern is expected in the first group of statements if the difference in summary style has an effect in the ability to assess relevance. The second group of statements should have the same answer pattern, regardless of an effect from summary style, if the two groups have a similar composition.

The second question asked participants to reflect on how they assessed relevance. Of particular importance was what they looked at in the hits. With this question we wanted to gain some insight into how the summary style might change the participants' strategy. For instance, would participants in SS2 look more at the expressions and participants in SS1 look more for the text?

Information need task

You have just finished attending a Linear Algebra class. Today's topic involved finding the inverse matrices through their adjoint matrix, but the professor did not explain how the formula $A^{-1} = \frac{1}{\det A} \text{adj } A$ was derived and you want to find that out.

You go to a search engine and search using the following keywords

$A^{-1} = \frac{1}{\det A} \text{adj } A$ proof
 Search

The search engine returns 10 results. Below you will see each of them one by one. You should decide whether each link is relevant to your search or not.

Please respond as quickly as possible, but take your time to make sure that you carefully consider whether a search result is relevant before you click Yes or No.

[Chapter 3 Determinants](#)

$$\left(\frac{1}{\det(A)} \text{Adj}(A)\right)A = I_3.$$

So, $A^{-1} = \frac{1}{\det(A)} \text{Adj}(A)$. So, the **proof** is complete when A is a 3×3 matrix. **Proof** in the general case: This means, A is an $n \times n$ matrix and ...

Is this link relevant?

Yes
 No

Figure 2.2: Interface of experiment website.

2.5 Materials

The experiment and data collection were performed by a custom made online system. It guided the participants, collected their responses to the tasks and showed them the final questionnaire. The description of the two experimental information needs and the familiarization task were included in the system, as well as each of the ten hit results formatted in each of the two summary styles.

Figure 2.2 shows a screen shot of the interface. Each of the three tasks had a "card" like the one in the figure that slid into view from the right for the participant to complete it. This was done so that participants could clearly see that the task had changed and that they needed to read a new information need.

Each of these cards had two sections. The top section (white) described the information need and showed the query terms in a mock-up search bar. This section was visible throughout the completion of the task so participants could refer to it if they needed to see the query terms or remember something about the information need.

The bottom section (gray) showed each of the hit results for the task. The buttons at the bottom were color coded so that participants could quickly identify them. Once they selected one of them, the current hit result slid out of view towards the left of the screen and a new hit slid into view from the right. This allowed participants to clearly see that the hit had changed and that they needed to make a new assessment.

The system was run on a server with Apache, PHP and MySQL. Additionally, client computers, used by the participants, had access to the server and were running Windows 7 and the Firefox browser, and had a standard keyboard and mouse.

Materials used outside the system were the consent form and a sign-off sheet to track the payments. The pre-screening and final questionnaires were setup using Google Forms. Posters and email for recruiting were also used (see Appendix C and D).

Chapter 3

Results

In the following sections we present the results obtained from the experiment. We start with a summary of the demographic data and the screening questionnaire. Next we report the data collected from the experiment with the corresponding statistical tests for learning effects and mean differences. We conclude with the results from the exit questionnaire with statistical test for differences in distributions.

3.1 Demographics

A total of 38 participants completed the experiment. All participants reported having normal, or corrected to normal, vision and hearing. Additionally, all participants indicated not having any problems, such as dyslexia, when reading from a computer screen. 73.7% (n=28) of participants were male and 26.3% (n=10) were female.

92.1% (n=35) of participants reported being between the ages of 18 and 24 with the rest reporting being between 25 and 34. 76.3% (n=29) reported their highest level of education as some college. See Figure 3.1 for more details on the level of education for participants.

Figure 3.2 shows the frequency for which they need to express mathematical notation in a computer. Each frequency level was balanced across both summary style groups when possible and a random assignment was used when the count for the frequency level was odd.

In the following sections a p-value of 0.05 or less will be considered statistically significant. P-values are reported at three different levels of significance (0.05, 0.01 and 0.005)

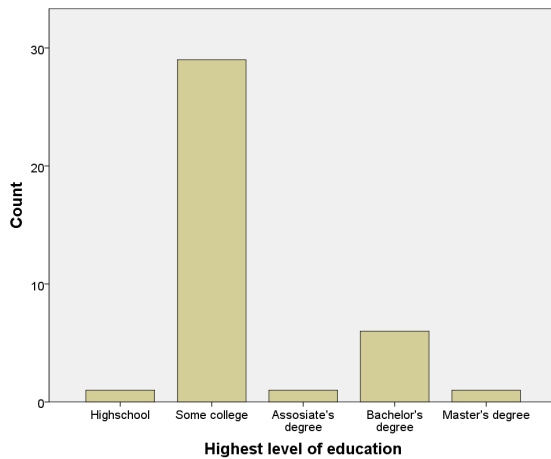


Figure 3.1: Distribution of the highest level of education achieved among the participants.

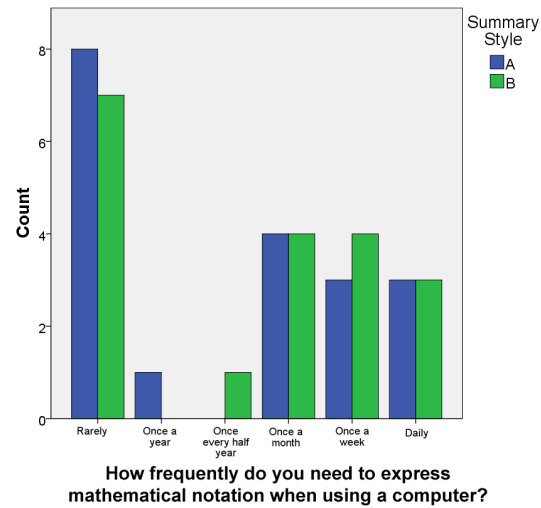


Figure 3.2: Distribution of participant's reported frequencies for the need to express mathematical notation in a computer.

based on where the calculated value falls.

3.2 Experiment

The mean response time taken by all participants to assess relevancy for each hit was 12.94 seconds ($sd = 5.77, n = 757$) and overall accuracy was 75.56%. A Pearson Correlation test was performed to test for learning effect across both summary styles. A small correlation between hit position and time was found for SS2 ($r = -0.143, p < 0.01$) but not for SS1 ($p > 0.05$) (see Figure 3.3). No correlation was found between accuracy and hit position for both summary styles ($p > 0.05$). A small negative correlation between time and accuracy was found for SS1 ($r = -0.114, p < 0.05$) but not for SS2 (see Figure 3.4).

Figures 3.5 and 3.6 show the detailed results for each individual hit across all participants. It can be observed that, except for hit 4, time taken to assess relevance is fairly constant. Hit 4 was relevant and has the longest expression in its task when measured as the horizontal space taken to display it.

On the other hand, accuracy varies more with Hits 2, 13 and 19 going below 50% in the control style group. Hit 13 was relevant and one of only two hits in its task that displayed

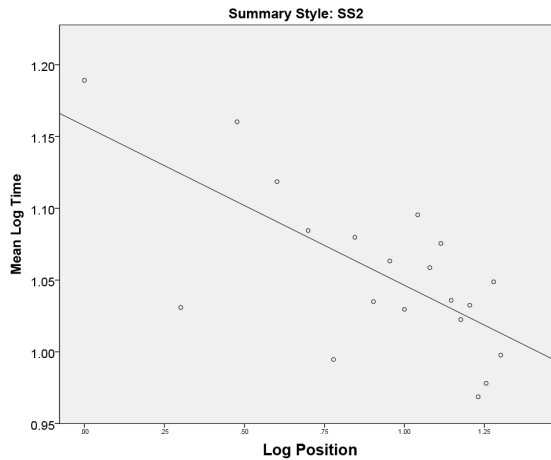


Figure 3.3: Mean response time by hit position in log scales for SS2 with a trend line.

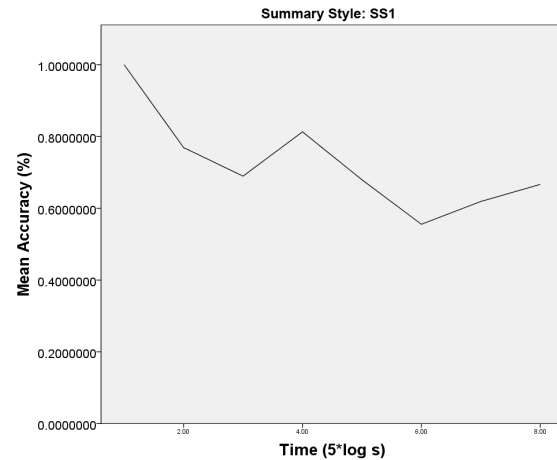


Figure 3.4: Mean accuracy by log of response time. Mean accuracy values aggregate the accuracy of all responses that fall between integer values of time in the log scale.

a date, this one having "Nov 6, 1998" and the other one being the not-relevant Hit 20 with "Oct 25, 2010". Hit 19 is not relevant and is one of two hits in its task that contains LaTeX code, the other one is the relevant Hit 12. There is nothing particularly different for the relevant Hit 2.

The detailed data collected from the experiment was summarized by participant and task. An accuracy score was calculated as the percentage of correct assessments and the time was calculated as the average time to make a relevance assessment for the hits in the task. The mean time to decide was 12.93 seconds ($sd = 4.66, n = 76$) and the mean accuracy was 75.57% ($sd = 16.60\%, n = 76$). Table 3.1 shows the mean and standard deviation for each combination of summary style and information need. Figures 3.7 and 3.8 show a profile view of the Accuracy and Time by Information Need and Summary Style.

Mixed-effects model ANOVA was used to test for differences in means. It allowed us to test for differences among multiple groups without increasing our chance of committing Type I errors. In addition, ANOVAs are robust to departures from normality specially when

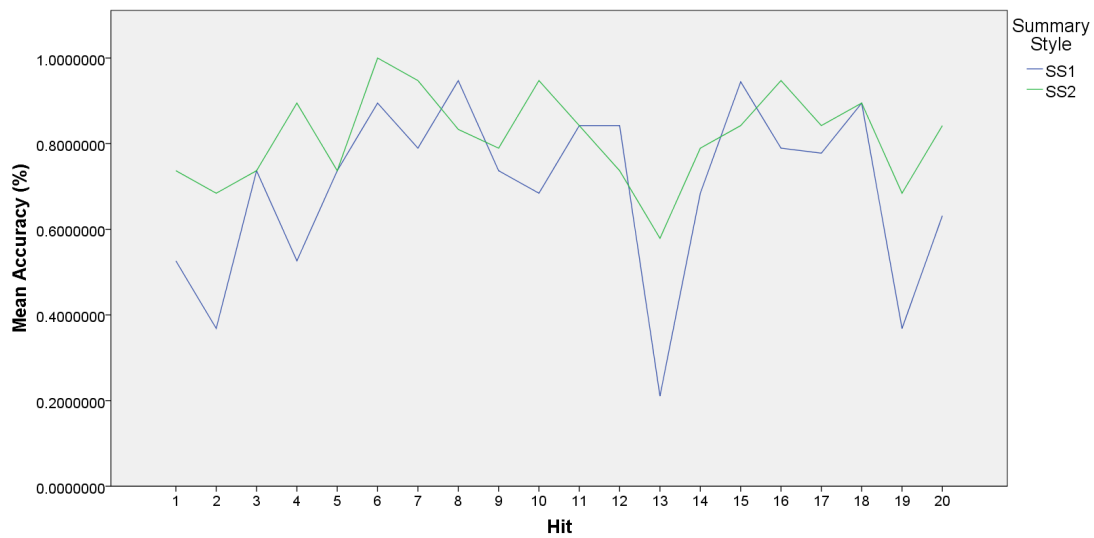


Figure 3.5: Average accuracy achieved by all participants for each of the hits. Hits 1-10 are from the informational need and hits 11-20 are for the resource need.

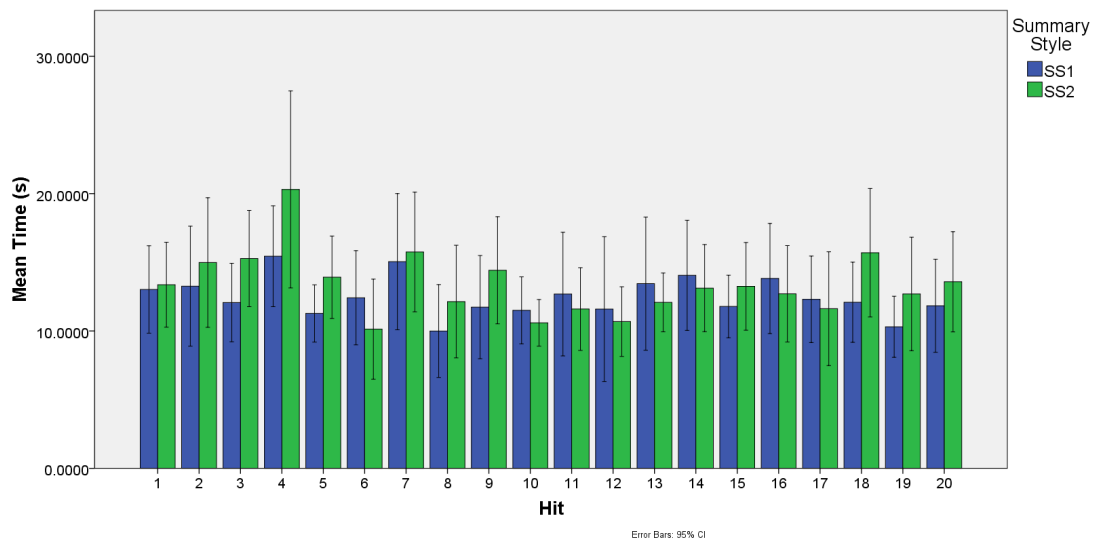


Figure 3.6: Average time taken to assess relevance by all participants for each of the hits. Hits 1-10 are from the informational need and hits 11-20 are for the resource need.

Summary		Accuracy (%)		Response Time (s)	
Style	Need	Mean	SD	Mean	SD
Control (n=19)	Informational	69.47	13.11	12.58	4.55
	Resource	69.71	20.78	12.39	4.79
Rendered (n=19)	Informational	83.10	12.01	14.06	5.11
	Resource	80.00	15.63	12.70	4.35
Total (n=38)		75.57	16.60	12.93	4.66

Table 3.1: Mean and standard deviation for each summary style group and information need

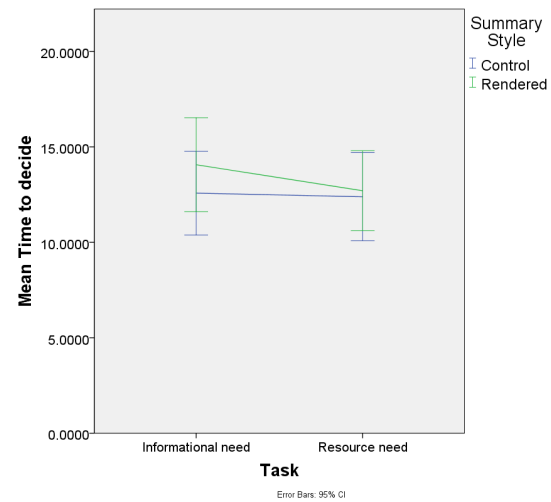
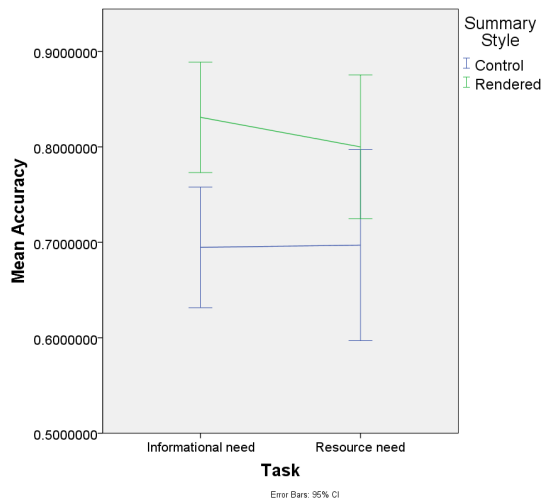


Figure 3.7: Profile graph showing Accuracy by Information Need and Summary Style

Figure 3.8: Profile graph showing Time by Information Need and Summary Style

the group sizes are the same, as in our case.

A 2 (Information Needs) x 2 (Summary Style) mixed-effects factorial ANOVA was performed on the response time. Time was not found to change by information need ($F(1, 36) = 1.407, p > 0.05$) or by summary style ($F(1, 36) = 0.427, p > 0.05$). In addition, no interaction effect was found either ($F(1, 36) = 0.802, p > 0.05$).

Another 2 (Information Needs) x 2 (Summary Style) mixed-effects factorial ANOVA was performed on accuracy scores. Accuracy scores were found to not change by information need ($F(1, 36) = 0.211, p > 0.05$) and no interaction effect was shown ($F(1, 36) =$

0.286, $p > 0.05$).

However, accuracy scores did change based on summary style ($F(1, 36) = 8.730, p < 0.01$). The mean accuracy for the rendered condition was higher than for the control condition by 11.96 percentage points. This translates to an increase in accuracy of 17.18% on average for participants in the rendered condition.

3.3 Exit Questionnaire

This section summarizes the responses participants gave for the exit questionnaire. Since a statistically significant difference was found between summary styles, graphics are displayed with both groups separated. The non-parametric Mann-Whitney test for independent samples was used to compare the answer distributions. This test is more efficient than t-tests for non normal distributions and works better on the ordinal data from our questions.

Figure 3.9 shows participants response to the statement "I'm familiar with the math involved in these tasks". Possible answers range from "Strongly agree" to "Strongly disagree". A Mann-Whitney Independent Samples test was run. No significant difference was shown between the two groups ($p > 0.05$). Similar behavior is shown in Figure 3.10 for the answers to the statement "I have had information needs similar to the tasks I just completed". A Mann-Whitney Independent Samples test again shows no statistical difference between both groups ($p > 0.05$).

Figure 3.11 shows the response of participants to the statement "I had no problems reading the results presented". A Mann-Whitney Independent Samples test was run between summary style groups. Responses were shown to be statistically different among both groups ($p < 0.005$).

Figure 3.12 shows participants response to the question "Which of the two tasks did you find the easiest?". A Mann-Whitney Independent Samples test was run. No significant difference was found between both summary style groups ($p > 0.05$). Of the people that chose the Resource Need as being easier, 71.43% ($n = 20$) explained their choice based on more familiarity with the math used (calculus vs. linear algebra) while of the people that

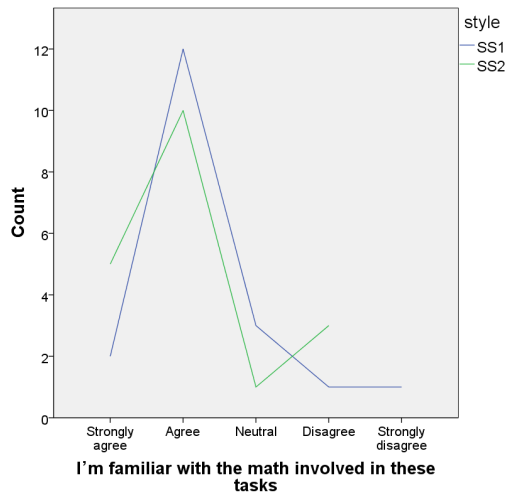


Figure 3.9: Participants response to the statement "I'm familiar with the math involved in these tasks".

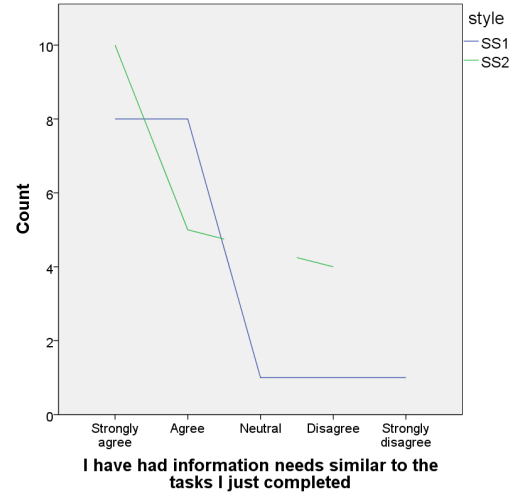


Figure 3.10: Participants response to the statement "I have had information needs similar to the tasks I just completed"

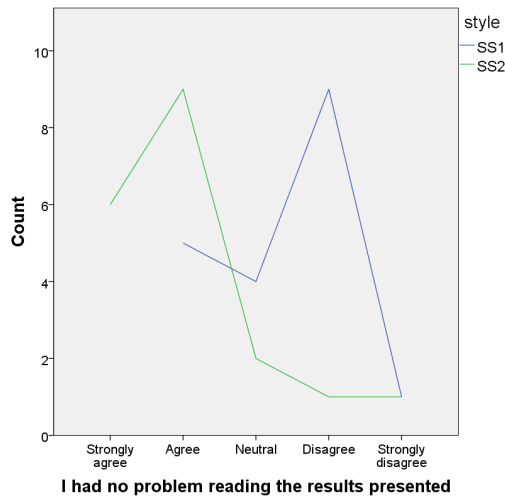


Figure 3.11: Participants response to the statement "I had no problems reading the results presented"

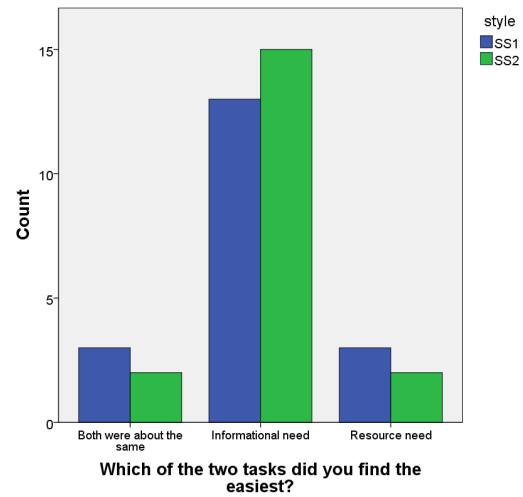


Figure 3.12: Profile graph showing Time by Information Need and Summary Style

chose the Informational Need 60.00% ($n = 3$) did the same.

3.4 Summary

Participants on the rendered summary (SS2) had statistically higher accuracy rates than participants in the control summary style (SS1). They showed an increase in performance of 17.18%. Additionally, participants in SS2 tended to agree more with the statement "I had no problems reading the results presented" than participants in SS1.

No statistical difference was found between the two information needs for either response time or accuracy. The same was shown for the two summary styles for response time. Additionally, no interaction effects between information need and summary style was shown for both response time and accuracy.

Both groups agreed similarly to the statements "I'm familiar with the math involved in these tasks" and "I have had information needs similar to the tasks I just completed". There was no difference in the answer distributions of both groups to the question "Which of the two tasks did you find the easiest?", however most participants selected the informational need as easier on both groups because of more familiarity with the math involved.

Chapter 4

Discussion

Great care was taken to make both summary style groups balanced. Figure 3.2 shows that both groups are very similar with respect to their experience using computers to create math expressions. Balancing this skill across the two groups was important as many of the hit results in the control summary style (SS1) contain LaTeX and other code used to input math expressions into a computer.

The similarity of both groups can be seen more clearly when looking at the answer patterns to the statements "I'm familiar with the math involved in these tasks" and "I have had information needs similar to the tasks I just completed" (Figures 3.9 and 3.10 respectively). Both questions probed participants about their past experience with the math involved in the tasks. Both summary style groups were shown to have statistically equal distributions for their answers, supporting our claim that there are not significant differences in the composition of both groups.

4.1 Learning Effects

Participants tended to assess relevance faster towards the end of the hit results, but only for the rendered summary style (SS2). Assessing relevance in SS2 was a new skill for participants and as such the response time seems to follow the Power Rule of Practice for response times [10]. The Power Rule of Practice states that reaction time for a task decreases linearly with the logarithm of the practice task. Tasks that follow this rule show an approximately straight line when the logarithm of the reaction time is plotted against the

logarithm of the practice task. Figure 3.3 shows this plot and the roughly linear correlation in our data. In our case, the reaction time is the time taken to assess relevance and the practice task is to assess the hit in a specific position.

We believe that this behavior is not seen in SS1 because assessing relevance with that style is already a learned and practiced skill in common search engines. What is not clear from our data is whether with enough practice the improvements in response time in SS2 will take its value statistically below response times in SS1.

4.2 Summary Style Effect

Our results support the hypothesis that properly formatting math expressions in the summaries improves users' ability to assess relevance. Evidence for this conclusion come from three sources: differences in accuracy, differences in perceived difficulty reading the summaries and a violation to the usual speed-accuracy trade-off.

Participants in the SS2 condition had a higher accuracy when assessing relevance of search results. On average, they were 17.18% more accurate than their counterparts in SS1. However, as opposed to what Aula [1] showed for text search, we found the difference in the accuracy and not in the time taken. Even though participants were told to take as much time as they needed to accurately assess relevance, the average times across conditions were remarkably similar. One possible explanation for this behavior is that it was not easy for our participants to assess relevance.

Not only were participants more accurate in SS2, they also reported having less overall difficulty reading the result summaries. In their responses to the statement "I had no problem reading the results presented" responses from the SS2 were mostly in agreement while in SS1 the answers were mostly in disagreement. This suggests that participant in SS1 had a higher cognitive load when analyzing the hits. This could lead to users having an improved user experience and higher satisfaction when using systems that implement SS2.

Finally, a negative correlation was shown between response time and accuracy for SS1.

This is opposed to what is usually expected as the speed-accuracy trade-off. According to an explanation by Busemeyer [4]:

Violations of the speed-accuracy trade-off can happen when the following two conditions are present: (a) the preference state is initially biased in the direction of the correct alternative and (b) the discriminability between the correct and incorrect alternative is low.

If this violation indeed happens in SS1, then our task should meet these two criteria. SS2, however, does not present this correlation between response time and accuracy. In both summary styles we can assume that the first condition is met in a similar fashion. It follows logically then that the reason the effect is not seen in SS2 is that the second condition is not met, which supports our claim that SS2 improves participants ability to assess relevancy when compared to SS1.

4.3 Information Need Effect

There were no effects shown in the data with respect to the change in information need. Accuracy and response time remained statistically constant across both conditions. However, these results are greatly confounded by the fact that the majority of participants chose the Resource Need as the easiest task on account of them being more familiar with the math used. Ideally, responses would have been distributed evenly in favor of both tasks or, if uneven, had most people explain their choice with respect to some intrinsic value of the information need.

However, it is interesting to point out that although the majority of participants chose the Resource Need as an easier task, participants did not take less time to assess relevance nor were they more accurate.

4.4 Summary

Assessing relevance under the rendered summary style (SS2) shows a learning effect that seems to follow the Power Rule of Practice [10]. It is unclear from our data if with enough practice response times in SS2 would be faster than in the control summary style (SS1). SS2 improved participants' ability to assess relevance. Participants in SS2 had higher accuracy and reported having less difficulty reading the result summaries. Additionally, participants in SS1 showed a violation to the speed-accuracy trade-off which happens, in part, due to low discriminability between correct and incorrect choices [4]. This effect was not shown for participants in SS2 suggesting that discriminability was higher in this summary style. Finally, there was no effect found for information need, but this result is greatly confounded since most participants reported being more familiar with one of the tasks based on the familiarity with the math used.

Chapter 5

Conclusions

Finding ways to help users assess relevance when they search using math expressions is critical for making MIR systems easier to use. We designed a study to test for effects on response time and relevance assessment between two different summary styles and across two different information needs. We used a summary style based on Google's result summaries with the math expressions in it rendered as math and one with the original linearized text form of the expressions used as control.

Our results provide evidence to support our hypothesis that the users' ability to assess relevance improves when properly rendering math expressions in the search results. Participants in the rendered condition of our study had 17.18% better accuracy and reported having less problem reading the results. Only participants in this condition showed a learning effect that, if extrapolated, could mean shorter response times, compared to the control summary style, to assess relevance once users are more familiar with it.

Additionally, participants in the non-rendered condition showed a violation of the speed-accuracy trade-off that can happen when the discriminability between correct and incorrect alternatives is low. Since this effect was not found in the rendered condition group, this suggests, indirectly, that the discriminability was higher in the rendered summary style.

Users are used to search results that contain mostly text and links, but our results suggest that current search systems, such as Google, should make a concerted effort to properly render the mathematical expressions shown in their summaries. Of course this is not an easy task as extracting the formatting information from the matched documents and rendering the results page in a way that does not break it is not trivial. Additionally, our results

hold for hit results analyzed one by one and different results could be obtained if hits are shown in a more traditional results page with multiple hits.

Our results do not support our hypothesis that any effect from summary style would be greater for the informational need. Unfortunately, given the confound raised by participants' higher familiarity with the math in one of the tasks, we cannot conclude if there is or not a difference in effect based on the information need. This is a clear avenue for future research as different needs may be better served using different summary styles.

Our results also suggest an interesting question: should non text-based extracts of relevant information in the matched documents be shown in the result summaries? We provided evidence that suggests that the answer is yes for mathematical expressions. Future research could focus on testing whether other types of multimedia (e.g. audio, figures, tables, etc.) also help in relevance assessment.

Continuing with our line of research, we want to test different summary styles. Now that we have evidence that Google search summaries with properly rendered matched expressions is better, what other modifications would help users? Possibilities include modifying the summaries to: increase the amount of document context that surround the matched expression in the document, show math expressions that surround the matched expression or vary the proportion of expressions to text.

References

- [1] A. Aula. Enhancing the readability of search result summaries. In *Proceedings of the Conference HCI 2004: Design for Life*, pages 1–4, 2004.
- [2] R. F. Boisvert and D. W. Lozier. *A Century of Excellence in Measurements Standards and Technology*, chapter Handbook of Mathematical Functions, pages 135–139. 2001.
- [3] Andrei Broder. A taxonomy of web search. *SIGIR Forum*, 36(2):3–10, September 2002.
- [4] JeromeR. Busemeyer. Violations of the speed-accuracy tradeoff relation. In Ola Svenson and A.John Maule, editors, *Time Pressure and Stress in Human Judgment and Decision Making*, pages 181–193. Springer US, 1993.
- [5] Zhiwei Guan and Edward Cutrell. An eye tracking study of the effect of target rank on web search. In *Proceedings of the SIGCHI Conference on Human Factors in Computing Systems*, CHI '07, pages 417–420, New York, NY, USA, 2007.
- [6] Zhiwei Guan and Edward Cutrell. What are you looking for? an eye-tracking study of information usage in web search. pages 407–416, 2007.
- [7] M. Hearst. *Search User Interfaces*. Search User Interfaces. Cambridge University Press, 2009.
- [8] M.D. Kickmeier and D. Albert. The effects of scanability on information search: An online experiment. In *Proc. of HCI*, pages 33–36, 2003.
- [9] Bruce R. Miller and Abdou Youssef. Technical aspects of the digital library of mathematical functions. *Annals of Mathematics and Artificial Intelligence*, 38(1-3):121–136, May 2003.

- [10] A. Newell and P.S. Rosenbloom. *Mechanisms of Skill Acquisition and the Law of Practice*. CMU-CS-80-145. Carnegie-Mellon University, Department of Computer Science, 1980.
- [11] Christopher Sasarak, Kevin Hart, Richard Pospesel, David Stalnaker, Lei Hu, Robert LiVolsi, Siyu Zhu, and Richard Zanibbi. min: A multimodal web interface for math search. *Symp. Human-Computer Interaction and Information Retrieval*, pages online, 4pp, 2012.
- [12] Anastasios Tombros and Mark Sanderson. Advantages of query biased summaries in information retrieval. In *Proceedings of the 21st annual international ACM SIGIR conference on Research and development in information retrieval*, SIGIR '98, pages 2–10, New York, NY, USA, 1998.
- [13] Abdou S. Youssef. Methods of relevance ranking and hit-content generation in math search. In *Calculus 07 / MKM 07: Proceedings of the 14th symposium on Towards Mechanized Mathematical Assistants*, pages 393–406, Berlin, Heidelberg, 2007.
- [14] Richard Zanibbi and Dorothea Blostein. Recognition and retrieval of mathematical expressions. *International Journal on Document Analysis and Recognition (IJDAR)*, 15(4):331–357, 2012.
- [15] Jin Zhao, Min-Yen Kan, and Yin Leng Theng. Math information retrieval: user requirements and prototype implementation. In *Proceedings of the 8th ACM/IEEE-CS joint conference on Digital libraries*, JCDL '08, pages 187–196, New York, NY, USA, 2008.

Appendix A

Screening Questionnaire

Hit Summary Relevance Assessment Study

* Required

Gender *

☐ Male

☐ Female

☐ Other:

Age range *

Minors are not eligible to participate

☐ 18-24

☐ 25-34

☐ 35-44

☐ 45-54

☐ 55-65

☐ 65-74

☐ 75+

What is the highest level of education that you have completed *

☐ Some highschool

☐ Highschool

☐ Some college

☐ Associate's degree

☐ Bachelor's degree

☐ Master's degree

☐ Professional degree (e.g. MD, LLD)

☐ PhD

If you studied at college, please indicate your major discipline(s) of study

e.g. Electrical Engineering, Chemistry, Music, Computer Science

If you studied at college, please indicate how many courses have you taken about mathematics *

☐ 0-1

☐ 2-3

☐ 4-7

☐ 8+

Please indicate the name of the mathematics courses you have taken *

How frequently do you need to look up mathematical information? Examples of mathematical information include function definitions *

e.g. trigonometric and statistical functions), definitions for mathematical symbols, function plots, mathematical models (e.g. environmental or physical models), theorems, and proofs

- ☐ Rarely
- ☐ Once a year
- ☐ Once every half year
- ☐ Once a month
- ☐ Once a week
- ☐ Daily

How frequently do you need to express mathematical notation when using a computer, such as for writing technical documents or in using computer programs such as Matlab, Mathematica, or Maple? *

- ☐ Rarely
- ☐ Once a year
- ☐ Once every half year
- ☐ Once a month
- ☐ Once a week
- ☐ Daily

Do you have normal or corrected to normal vision?

- ☐ Yes
- ☐ No

Do you have normal or corrected to normal hearing?

- ☐ Yes
- ☐ No

Do you have any impediment when reading text of a computer screen? *

e.g. some forms of dyslexia

- ☐ Yes
- ☐ No

Please provide your email so we can contact you in case you qualify for the study

We will not share your email with anyone not involved in the study and will only use it to contact you

in case you qualify

Submit

Never submit passwords through Google Forms.

Powered by


This content is neither created nor endorsed by Google.

[Report Abuse](#) - [Terms of Service](#) - [Additional Terms](#)

Appendix B

Post-test Questionnaire

1) Please state how much you agree or disagree with the following statements

	Strongly dis- agree	Disagree	Neutral	Agree	Strongly agree
Deciding if the results were relevant or not was hard					
I have had information needs similar to the tasks I just completed					
The summaries presented were clear and understandable					
I had no problem reading the results presented					
I am familiar with the math involved in these tasks					

2) How did you assess the relevancy of the result summaries? i.e. what did you look at, what did you think about, etc.

3) Which of the two tasks did you find the easiest?

- Task 1
- Task 2
- Both were about the same

4) Why?

5) Please provide any additional comments that you have below.

Appendix D

Recruiting Email

To: XXXXX

Subject: Seeking Participants for Hit Summary Relevance Assessment Experiment

The Document and Pattern Recognition Lab (DPRL) at RIT is looking for participants in an experiment studying relevance assessment of search engine results from mathematical queries. Knowing the best way to generate and present hit summaries will improve users performance when searching using mathematical expressions.

The study is expected to last 30 minutes. Participants will be paid \$10 for their time.

If you would like to participate in the project, please go to <http://bit.ly/108pkCO> and complete the questionnaire. If you have any questions regarding the study please contact Matthias Reichenbach (msr5919@rit.edu, (585) 214-9785). Any questions about your rights as a participant may be directed to Heather Foti (Associate Director, Human Subjects Research Office, RIT: hmfsrcs@rit.edu (585) 475-7673) or Associate Prof. Richard Zanibbi (Principal Investigator, rlaz@cs.rit.edu, (585) 475-5023).

Sincerely,

Matthias Reichenbach

Graduate Student, M.S. HCI Program

Appendix E

Pre-generated Hit Results

Following is the list of all the hits used in our experiment. They are organized by task, with the rendered version displayed on the left and the control on the right. For each of the tasks, the first half of the hits are relevant to the information need and the second half is not.

E.1 Familiarization Task

[Pascal's triangle](#)

Although named after Blaise Pascal, who studied it, this arithmetic triangle has ...

$$(x + y)^2 = 1x^2 + 2xy + 1y^2$$

... **Pascal's triangle** also has the following properties: ...

[Pascal's triangle](#)

Although named after Blaise Pascal, who studied it, this arithmetic triangle has ... $(x + y)^2 = 1x^2 + 2xy + 1y^2$... **Pascal's triangle** also has the following properties: ...

[How to factorize the polynomial \$\(2x - y\)^6 - \(2x + y\)^6\$ - WyzArt Tutoring](#)

The expression $(2x - y)^6 - (2x + y)^6$ is a difference of squares since it is $[(2x - y)^3]^2 - [(2x + y)^3]^2$ is greater than 3 are complicated and it is best to use **Pascal's triangle**. The x exponent in the first term will be 6 and then descend from there as the ...

[How to factorize the polynomial \$\(2x - y\)^6 - \(2x + y\)^6\$ - WyzArt Tutoring](#)

The expression $(2x - y)^6 - (2x + y)^6$ is a difference of squares since it is $[(2x - y)^3]^2 - [(2x + y)^3]^2$ is greater than 3 are complicated and it is best to use **Pascal's triangle**. The x exponent in the first term will be 6 and then descend from there as the ...

[Pascal's Triangle by Shelby M on Prezi](#)

-Row Sums Row 0 Row 1 Row 2 Works Cited "Binomial theorem. ... in **Pascal's Triangle**
2 1, 9, 36, 84, 126, 126, 84, 36, 9, 1 1, 2, 1

$$x^2 + 2xy + y^2$$

$$x + y = 2$$

$$x + y = 9$$

...

[Pascal's Triangle by Shelby M on Prezi](#)

-Row Sums Row 0 Row 1 Row 2 Works Cited "Binomial theorem. ... in **Pascal's Triangle**
2 1, 9, 36, 84, 126, 126, 84, 36, 9, 1 1, 2, 1 $x + 2xy + y^2$ $(x + y)^2 = 2(x + y) = 9$...

[Algebraic expressions, expanding and factoring algebraic expressions](#)

$$... ab - b^2 = a^2 - b^2, ...$$

$$c) (x + y) \cdot (x^2 - xy + y^2) = x^3 - x^2y + xy^2 + x^2y - xy^2 + y^3. ...$$

$$c) (2x + y)^3 = (2x)^3 + 3 \cdot (2x)^2 \cdot y + 3 \cdot (2x) \cdot y^2 + y^3 = 8x^3 + 12x^2y + 6xy^2 + y^3$$

... The binomial coefficients can also be obtained by using **Pascal's triangle**.

[Algebraic expressions, expanding and factoring algebraic expressions](#)

... $ab - b^2 = a^2 - b^2, c) (x + y) \cdot (x^2 - xy + y^2) = x^3 - x^2y + xy^2 + x^2y - xy^2 + y^3 = x^3 + y^3. ... c) (2x + y)^3 = (2x)^3 + 3 \cdot (2x)^2 \cdot y + 3 \cdot (2x) \cdot y^2 + y^3 = 8x^3 + 12x^2y + 6xy^2 + y^3. ...$ The binomial coefficients can also be obtained by using **Pascal's triangle**.

E.2 Informational Need Task

[lecture notes on determinants - PlanetMath](#)

Then, $A \operatorname{adj}(A) = \det(A)I$. Furthermore, if A is invertible, then $A^{-1} = \frac{1}{\det(A)} \operatorname{adj}(A)$.
Let's consider the **proof** for the 3×3 case. We aim to show ...

[lecture notes on determinants - PlanetMath](#)

Then, $A \operatorname{adj} A = \det(A) I$. Furthermore, if A is invertible, then $A^{-1} = \frac{1}{\det(A)} \operatorname{adj} A$. Let's consider the **proof** for the 3×3 case. We aim to show ...

[Math Refresher: Adjoint of a Matrix](#)

Then $A(\operatorname{adj} A) = (\det A)I$ and $(\operatorname{adj} A)A = (\det A)I$ **Proof:** (1) Let $D = A(\operatorname{adj} A)$ (2) Then, $d_{i,k} = \sum_{j=1}^n a_{i,j} \operatorname{cof}_{k,j}(A)$. [See Definition 1, here for ...

[Math Refresher: Adjoint of a Matrix](#)

Then $A(\operatorname{adj} A) = (\det A)I$ and $(\operatorname{adj} A)A = (\det A)I$ **Proof:** (1) Let $D = A(\operatorname{adj} A)$ (2) Then, $d_{i,k} = \sum_{j=1}^n a_{i,j} \operatorname{cof}_{k,j}(A)$. [See Definition 1, here for ...

[3.4. The determinant of a matrix](#)

Proof. 2) If $A \neq 0$, then $\frac{1}{\det A} \operatorname{adj} A \cdot A = I \Rightarrow A^{-1} = \frac{1}{\det A} \operatorname{adj} A$. 1) If A has an inverse, then $I = AA^{-1} \Rightarrow 1 = \det I = \det(AA^{-1}) = \det A \cdot \det A^{-1} \Rightarrow \det A \neq 0$...

[3.4. The determinant of a matrix](#)

Proof. 2) If $\det A \neq 0$, then $(1 / \det A) \operatorname{adj} A \cdot A = I \Rightarrow A^{-1} = (1 / \det A) \operatorname{adj} A$. 1) If A has an inverse, then $I = AA^{-1} \Rightarrow 1 = \det I = \det(AA^{-1}) = \det A \cdot \det A^{-1} \Rightarrow \det A \neq 0$...

[Linear Algebra WebNotes. Part 3.](#)

The first indexes form a permutation of the set $\{1,2,3,4\}$

$$A^{-1} = \frac{1}{\det(A)} \operatorname{adj}(A)$$

Proof. Indeed, if A is invertible then by the third theorem about determinants ...

[Linear Algebra WebNotes. Part 3.](#)

The first indexes form a permutation of the set $\{1,2,3,4\}$ $A^{-1} = (1/\det(A)) \operatorname{adj}(A)$. **Proof.** Indeed, if A is invertible then by the third theorem about determinants ...

[Linear Algebra WebNotes. Part 3.](#)

The first indexes form a permutation of the set $\{1,2,3,4\}$

$$A^{-1} = \frac{1}{\det(A)} \operatorname{adj}(A)$$

Proof. Indeed, if A is invertible then by the third theorem about determinants ...

[Linear Algebra WebNotes. Part 3.](#)

The first indexes form a permutation of the set $\{1,2,3,4\}$ $A^{-1} = (1/\det(A)) \operatorname{adj}(A)$. **Proof.** Indeed, if A is invertible then by the third theorem about determinants ...

[linear algebra - Eigenvalues of adjoint of non-singular matrix ...](#)

How could we **prove** that ... then eigenvalues of **adj** A are

$$\frac{\det A}{\lambda_1}, \frac{\det A}{\lambda_2}, \frac{\det A}{\lambda_3}, \dots$$

[linear algebra - Eigenvalues of adjoint of non-singular matrix ...](#)

How could we **prove** that ... then eigenvalues of **adj** A are $\frac{\det A}{\lambda_1}, \frac{\det A}{\lambda_2}, \frac{\det A}{\lambda_3}, \dots$

[Lecture 30](#)

Thus **adj** $(A) = [A_{ij}]^T$. ADJOINT THEOREM. For any $n \times n$ matrix A ,

$$A \cdot \text{adj}(A) = \text{adj}(A) \cdot A = \det(A) \cdot I_n.$$

$$\text{If } \det(A) \neq 0, \text{ then } A^{-1} = \frac{1}{\det(A)} \text{adj}(A).$$

PROOF. We **prove** that $A \dots$

[Lecture 30](#)

Thus **adj** $(A) = [A_{ij}]^T$. ADJOINT THEOREM. For any $n \times n$ matrix A , $A \text{adj}(A) = \text{adj}(A) A = \det(A) I_n$. If $\det(A) \neq 0$, then $A^{-1} = \frac{1}{\det(A)} \text{adj}(A)$. **PROOF.** We **prove** that $A \dots$

Proof: $\text{adj}(\text{adj}(A)) = (\det(A))^{n-2} \cdot A$

Proof: $\text{adj}(\text{adj}(A)) = (\det(A))^{n-2} \cdot A$ for $A \dots A) = (\det(A))^{n-1} \cdot I_n$ and using (2) $\text{adj}(A) \cdot \dots$

Proof: $\text{adj}(\text{adj}(A)) = (\det(A))^{n-2} \cdot A$

Proof: $\text{adj}(\text{adj}(A)) = (\det(A))^{n-2} \cdot A$ for $A \dots A) = (\det(A))^{n-1} \cdot I_n$ and using (2) $\text{adj}(A) \cdot \dots$

[Chapter 3 Determinants](#)

$$\left(\frac{1}{\det(A)} \text{Adj}(A)\right)A = I_3.$$

So, $A^{-1} = \frac{1}{\det(A)} \text{Adj}(A)$. So, the **proof** is complete when A is a 3×3 matrix. **Proof** in the general case: This means, A is an $n \times n$ matrix and ...

[Chapter 3 Determinants](#)

$\text{adj}(A)A = I_3$. So, $A^{-1} = \frac{1}{\det(A)} \text{adj}(A)$. So, the **proof** is complete when A is a 3×3 matrix. **Proof** in the general case: This means, A is an $n \times n$ matrix and ...

[linear algebra - How do I prove that \$\det A = \det A^t\$ - Mathematics](#)

Can anyone please comment whether my **proof** is correct or not? Attempted solution: If ... What is $\frac{\det(A+ti)}{\det(B+ti)}$ as $t \rightarrow 0$? ... How to **prove** that if $\det(A) = 0$ then $\det(\text{adj}(A)) = 0$? a question ...

[linear algebra - How do I prove that \$\det A = \det A^t\$ - Mathematics](#)

Can anyone please comment whether my **proof** is correct or not? Attempted solution: If ... What is $\frac{\det(A+ti)}{\det(B+ti)}$ as $t \rightarrow 0$? ... How to **prove** that if $\det(A) = 0$ then $\det(\text{adj}(A)) = 0$? a question ...

E.3 Resource Need Task

[Calculus Made Easier: A Calculus Tutorial](#)

$(x^n)' = nx^{n-1}$ So, if $f(x) = x^4$, then $f'(x) = (4)x^3 = 4x^3$. And, for any polynomial term with a constant factor, the general derivative formula is $(cx^n)' = (c)(n)x^{n-1}$ where c ...

[Calculus Made Easier: A Calculus Tutorial](#)

$(xn)' = nxn-1$ So, if $f(x) = x^4$, then $f'(x) = (4)x^3 = 4x^3$. And, for any polynomial term with a constant factor, the general derivative formula is $(c xn)' = (c)(n)xn-1$ where c ...

[Quotient Rule for Derivatives - HMC Calculus Tutorial](#)

back to the math **tutorial** index ...

$$f'(x) = \frac{(x-3)\frac{d}{dx}[2x+1] - (2x+1)\frac{d}{dx}[x-3]}{(x-3)^2} = \frac{(x-3)(2) - (2x+1)(1)}{(x-3)^2} = -\frac{7}{(x-3)^2}$$

[Quotient Rule for Derivatives - HMC Calculus Tutorial](#)

back to the math **tutorial** index ... $f'(x) = \frac{\frac{d}{dx}\{(x-3)\} \frac{d}{dx}\{2x+1\} - (2x+1) \frac{d}{dx}\{x-3\}}{(x-3)^2} = \frac{(x-3)(2) - (2x+1)(1)}{(x-3)^2} = -\frac{7}{(x-3)^2}$.

[Differentiation - The University of Akron](#)

Nov 6, 1998 - ... throughout the rest of these **tutorials** and the rest of your Calculus course.

$$\frac{dx^r}{dx} = rx^{r-1}$$

It is this formula that will be our working formula.

[Differentiation - The University of Akron](#)

Nov 6, 1998 - ... throughout the rest of these **tutorials** and the rest of your Calculus course. $dx^r dx = rx^{r-1}$. It is this formula that will be our working formula.

[Tutorial for Derivatives of Powers, Sums and Constant Multiples](#)

The notation $\frac{d}{dx}$ means the derivative with respect to x. Thus, for instance,

$$\frac{d}{dx}[x^3] = 3x^2 \quad \frac{d}{dx}[1] = 0 \quad \frac{d}{dx}\left[\frac{1}{x}\right] = \frac{-1}{x^2}$$

... We can find the derivative of more complicated ...

[Tutorial for Derivatives of Powers, Sums and Constant Multiples](#)

The notation $d dx$ means the derivative with respect to x. Thus, for instance, $d dx x^3 = 3x^2 d dx [1] = 0 d dx 1/x = -1/x^2$... We can find the derivative of more complicated ...

[Tutorial for Derivatives of Powers, Sums and Constant Multiples](#)

The notation $\frac{d}{dx}$ means the derivative with respect to x. Thus, for instance,

$$\frac{d}{dx}[x^3] = 3x^2 \quad \frac{d}{dx}[1] = 0 \quad \frac{d}{dx}\left[\frac{1}{x}\right] = \frac{-1}{x^2}$$

... We can find the derivative of more complicated ...

[Tutorial for Derivatives of Powers, Sums and Constant Multiples](#)

The notation $d dx$ means the derivative with respect to x. Thus, for instance, $d dx x^3 = 3x^2 d dx [1] = 0 d dx 1/x = -1/x^2$... We can find the derivative of more complicated ...

[Tutorial: The Indefinite Integral](#)

$\int 2x dx = x^2 + C$, The indefinite integral of $2x$ with respect to x is $x^2 + C$ $\int 4x^3 dx = x^4 + C$,
The indefinite ... Function, Antiderivative, Formula. x^n ($n \neq -1$) ...

[Tutorial: The Indefinite Integral](#)

$\int 2x \, dx = x^2 + C$, The indefinite integral of $2x$ with respect to x is $x^2 + C$. $\int 4x^3 \, dx = x^4 + C$, The indefinite ... Function, Antiderivative, Formula. x^n ($n \neq -1$) ...

[Differential Calculus: Techniques of Differentiation](#)

This **tutorial** is going to show you the *easy* side of differentiating. In the examples from the ...

$$\frac{d}{dx} x^n = n * x^{(n-1)}, \text{ where } n \text{ is a constant.}$$

So, to evaluate the ...

[Differential Calculus: Techniques of Differentiation](#)

This **tutorial** is going to show you the *easy* side of differentiating. In the examples from the ... $d \rightarrow xn = n * x(n-1)$, where n is a constant dx . So, to evaluate the ...

[Derivative of the natural log - Straight Dope Message Board](#)

If you know that $\frac{d}{dx} e^{f(x)} = f'(x) \cdot e^{f(x)}$, where $f(x)$ is some ... $\frac{d}{dx}$ of $a^3 + b^2 + c$ is **1 ? No**.
The derivative of that is 0, because all the ...

[Derivative of the natural log - Straight Dope Message Board](#)

If you know that d/dx ($\exp(f(x)) = f'(x) * \exp(f(x))$), where $f(x)$ is some ... d / dx of $a^3 + b^2 + c$ is **1 ? No**. The derivative of that is 0, because all the ...

[derivative of \$x^n\$ - PlanetMath](#)

... the casual rule

$$\frac{d}{dx} (x^n) = nx^{n-1}$$

for positive integer values of n can ... First notice that

$$(x-a)(x^{n-1} + x^{n-2}a + \dots + xa^{n-2} + a^{n-1}) = x^n \dots$$

[derivative of \$x^n\$ - PlanetMath](#)

... the casual rule $\frac{d}{dx} (x^n) = nx^{n-1}$ for positive integer values of n can ... First notice that $(x-a)(x^{n-1} + x^{n-2}a + \dots + xa^{n-2} + a^{n-1}) = x^n \dots$

[Differentiate the function \$y = \sin^4\(\sqrt{u}\)\$ - Math - Questions & Answers](#)

Oct 25, 2010 - $y = \sin^4(\sqrt{u})$. $y' = 4 \cos^3(\sqrt{u}) - \sin(\sqrt{u}) \frac{1}{2\sqrt{u}}$... Also we use $\frac{d}{dx} x^n = nx^{n-1}$ to differentiate the given expression. $y = \sin \dots$

[Differentiate the function \$y = \sin^4\(\sqrt{u}\)\$ - Math - Questions & Answers](#)

Oct 25, 2010 - $y = \sin^4(\sqrt{u})$. $y' = 4\cos^3(\sqrt{u}) * \sin(\sqrt{u}) * 1/(2*\sqrt{u})$... Also we use $d/dx(x^n) = nx^{n-1}$ to differentiate the given expression. $y = \sin \dots$