

# ANALY482 MEETING MINUTES WITH SUPERVISOR(S)

<b>Date:</b>	18 January 2017
<b>Time:</b>	14:30 – 15:00
<b>Venue:</b>	Meeting Rm 4-1 (School of Information Systems)
<b>Attendees:</b>	Lu Ning, Song Rui, Dina, Prof Kam
<b>Absentees:</b>	
<b>Agenda:</b>	Finding out if our approach meets Prof Kam's requirements for the project

S/N	Things Discussed/Done	Remark
1.	Updates about our domain summary and methodology	<ol style="list-style-type: none"> <li>1. Limitation is not referring to data we have               <ol style="list-style-type: none"> <li>1. show the example why we have this kind of limitation</li> </ol> </li> <li>2. Raw data are not made into data file               <ol style="list-style-type: none"> <li>1. Lu Ning: organize data into separate txt files base on the domain name                   <ol style="list-style-type: none"> <li>I. Prof: need to translate raw data file into a data table/ extract them into different column that capture certain characteristics, URL, time spend, etc. know how to dealt with missing value. Then will be able to do exploratory (summary frequency count, distribution...)</li> <li>II. based on queries, find a way to put them into column</li> </ol> </li> </ol> </li> <li>3. Prof: Capture both behavior (pdf is being viewed online), some people only read pdf online. But some will download directly. Later then see how we can decide on the analysis. (don't eliminate possible behavior first, find first then drop later)</li> <li>4. prof: Review the URL if we can differentiate view online or actual download. (okay)</li> <li>5. we can see the link but cannot re-query. this analysis is very unstable. Point to individual publisher via proxy.               <ol style="list-style-type: none"> <li>1. Prof: The current system may create grey areas for the analysis. (proxy on a proxy)</li> </ol> </li> </ol>

		<p>6. Prof: bench marking is not really a feasible solution. (e.g. same bank on different location cannot reapply the model been built)</p> <p>7. Proxy Approach - Downloads</p> <ol style="list-style-type: none"> <li>1. Prof: URL will lead you back to the publisher. some publisher is expired.</li> <li>2. Prof: two sets of data (pdf vs non-pdf) cannot draw the conclusion on which set to choose from. PDF has 3 sub-groups there. We should consider each group on their detailed behaviors (e.g. The third bar, what are the characteristics? To determine if it is a noise, maybe is a scan image?)</li> <li>3. Prof: Take all collection as the base, then look at the file size and draw a fairer conclusion.</li> <li>4. Prof: There is way to match the URL that (replace the front part of the URL) (via manual sampling to find out the pattern of the)</li> </ol> <p>8. Prof: Do not support (pattern approach), because before we explore our situation how our lib been used, we cannot assume we can use other people approach. We should do a proper understanding our situation first, if ours are identical to theirs, then we can follow their approach. But if our approach is really contradicting, then we must change.</p>
--	--	---

<b>Item Due (Team)/Action(s)</b>
<p>Deadline: 8 February 2017</p> <ol style="list-style-type: none"> <li>1. Inform Aaron on our revised methodology</li> </ol>