

FIT5147 Visualisation Project Report

Nicholas Pennell, 30338913

Introduction

The idea behind my visualization is to create a dashboard to show three findings from my exploration report;

- the change of what pitch types are thrown from game to game and year to year,
- the location of where each pitch is thrown,
- And Where and what happens to each pitch (the outcome).

As well adding some extra information, and visualization features threw the use of filtering the data and having links to videos.

The visualization was design for people with a knowledge of baseball looking to find more visual information about their favorite player all in one place.

Design

To start out the design procedure all ideas for plots and features of each plot were written down, (idea sheet). This yield four main sections.

The time series plot, a plot where the y-axis shows the percentage of what pitches were thrown on the corresponding date, the x-axis. The idea was to have a button to switch between how the data is group, by year or game. As well as add a way to show key events and important dates.

The next was the location plots. A way to show where each pitch is thrown to the batter. Ideas include a scatter plot, heat map, or both with a way to switch. Any of these options would be overlaid with a strike zone and a home plate to help show the ground. The plot would also have a way to link to a video of the pitch, and the data be filtered to look a certain scenario.

The plot to show the outcome of a pitch yielded two ideas. A bar chart to show the count of outcomes or a scatter plot that overlays a baseball field to show where the outcome occurs. Either option would share the data filter options of the location plot.

The last idea was to add a bio page, or a way to share information about the player and/or the visualization.

The first design (sheet 2) was a simple dashboard, it included the time series plot from the idea sheet with a date slider below. A bar chart to count how many pitches were thrown and what type in the same time frame as the time series. With the location and outcome plots below that. The filter options would be on the left side of the page and effect the location and outcome plot.

The second design (sheet 3), was like first but upgraded as a shiny site, adding the bio page while splitting each plot into its own tab to give more space and detail to each plot.

The third design (sheet 4) was a d3 article/slide show, like the second design plots were broken up into their own page. As it would be an article style design each page would have more text and have any linked videos embedded into the page rather than a link to an external source.

The finale design (sheet 5), is a combination of the other sheets. It takes the dashboard idea and uses that as one tab, takes the bio page shiny site idea as uses that as another tab. This idea was selected as I felt it displayed the data in a way that would be easy to follow and use, not overloading the user with information and be achievable to accomplish in the given time frame.

Implementation

For the implementation stage more data was sourced from several places:

- A list of player id's numbers and names, from [baseball prospectus](#),
- Shape information of baseball field's to plot, from [GitHub user bdilley](#),
- Json data was manually scraped from [baseballsavant.mlb.com](#) to get video links of pitches,
- A career information table was scraped from [baseballsavant.mlb.com](#) as well. It was saved to keep load times of the app shorter.
- I then made an event data file with key events, from information gathered from [Baseball Reference](#) and found matching videos on [YouTube](#)

The choice was made to use shiny to create the app, due to having more knowledge and a higher comfort level with it rather than other options. The plots are made using the [plotly](#) package due to its added interactive features over ggplot2. A full list of used packages can be found in the references.

The implementation of the bio page was straight forward, the bio section on the left side of the screen is just a head shot, every MLB player has one and the url is the same just with a different id number, (which was in the original data frame). The table on the bio page was scraped as said above, this table was selected as it has all the key simple stats. The table is then printed using [kableExtra](#) package due, to its html settings as well as it visually appealing look. The Tm column has its elements coloured to match the team they represent, to draw a user's attention to that attribute. The top 3 values are in each other column is also highlighted for the same reason, this time with lighter colour to be less grabbing. The * symbol was also used to show all start years. The foot note at the bottom of the table is added to explain this to the user.

The dashboard page was a much more time-consuming process with many ideas not coming to fruition, and many compromises. The time series plot at the top of the page, starts in the game-by-game view, with each line representing a pitch type. A point can be hover on to show more information about said point. The dots above the lines represent career events, and have more information when hovered on, the brighter pink dots can be clicked on, if done a new window opens to video of the event.

The time series plot at the top of the page, starts in the game-by-game view, with each line representing a pitch type. A point can be hover on to show more information about said point. The dots above the lines represent career events, and have more information when hovered on, the brighter pink dots can be clicked on, if done a new window opens to video of the event.

The location and hit spray chart are located below the time series. They share a row, so the user can view all plots at one time. The location plot shares colours based on the pitch type with the time series and is also faceted by pitch type to allow the user to a better view. The strike zone (the square) is the average size of an MLB strike zone and is added to show reference as well as the home plate is added to show where the ground is. The option to switch between a heat map and a scatter plot was scraped once I had to program the facet feature for Plotly myself. Some points (2019 only) have point that are clickable and will open a window to a video and some information on baseballsavant.mlb.com.

The spray chart has the same click feature as the location plot, but it over lays a shape of the baseball field using the data from Git Hub user bdilday. A major downside with this plot is that the legend was removed due to some scenarios lead to the legend being to large and blocking the plot itself. The hover information does say what the outcome is though. The text around the baseball field shape is the distance in feet, all measurements are in imperial as that's the standard in the for baseball and the MLB, to the edge of the field that is plotted. The filter setting on the left side has the option I felt where the most important there are many more that could be added, the list was cut short to keep it all in one page. The date option effects all three options, while the other options just effect the location and it spray chart. The starting filter setting where set as such to show off all features but to also keep load time down.

User guide

The bio page is straight forward and has new user interaction. To get to the dashboard a user must just click on the dashboard tab.

FIT5147 Visualisation Project

Chris Sale	Player Bio	Dash Board																																																																																																																																																																																																																							
	Carrer MLB Pitching Statistics <table border="1"> <thead> <tr> <th>Season</th> <th>Tm</th> <th>LG</th> <th>BF</th> <th>W</th> <th>L</th> <th>ERA</th> <th>G</th> <th>GS</th> <th>SV</th> <th>IP</th> <th>H</th> <th>R</th> <th>ER</th> <th>HR</th> <th>BB</th> <th>SO</th> <th>WHIP</th> </tr> </thead> <tbody> <tr><td>2010</td><td>CWS</td><td>AL</td><td>92</td><td>2</td><td>1</td><td>1.93</td><td>21</td><td>0</td><td>4</td><td>23.1</td><td>15</td><td>5</td><td>5</td><td>2</td><td>10</td><td>32</td><td>1.07</td></tr> <tr><td>2011</td><td>CWS</td><td>AL</td><td>288</td><td>2</td><td>2</td><td>2.79</td><td>58</td><td>0</td><td>8</td><td>71</td><td>52</td><td>22</td><td>22</td><td>6</td><td>27</td><td>79</td><td>1.11</td></tr> <tr><td>2012*</td><td>CWS</td><td>AL</td><td>772</td><td>17</td><td>8</td><td>3.05</td><td>30</td><td>29</td><td>0</td><td>192</td><td>167</td><td>66</td><td>65</td><td>19</td><td>51</td><td>192</td><td>1.14</td></tr> <tr><td>2013*</td><td>CWS</td><td>AL</td><td>866</td><td>11</td><td>14</td><td>3.07</td><td>30</td><td>30</td><td>0</td><td>214.1</td><td>184</td><td>81</td><td>73</td><td>23</td><td>46</td><td>226</td><td>1.07</td></tr> <tr><td>2014*</td><td>CWS</td><td>AL</td><td>685</td><td>12</td><td>4</td><td>2.17</td><td>26</td><td>26</td><td>0</td><td>174</td><td>129</td><td>48</td><td>42</td><td>13</td><td>39</td><td>208</td><td>0.97</td></tr> <tr><td>2015*</td><td>CWS</td><td>AL</td><td>854</td><td>13</td><td>11</td><td>3.41</td><td>31</td><td>31</td><td>0</td><td>208.2</td><td>185</td><td>88</td><td>79</td><td>23</td><td>42</td><td>274</td><td>1.09</td></tr> <tr><td>2016*</td><td>CWS</td><td>AL</td><td>907</td><td>17</td><td>10</td><td>3.34</td><td>32</td><td>32</td><td>0</td><td>226.2</td><td>190</td><td>88</td><td>84</td><td>27</td><td>45</td><td>233</td><td>1.04</td></tr> <tr><td>2017*</td><td>BOS</td><td>AL</td><td>851</td><td>17</td><td>8</td><td>2.9</td><td>32</td><td>32</td><td>0</td><td>214.1</td><td>165</td><td>73</td><td>69</td><td>24</td><td>43</td><td>308</td><td>0.97</td></tr> <tr><td>2018*</td><td>BOS</td><td>AL</td><td>617</td><td>12</td><td>4</td><td>2.11</td><td>27</td><td>27</td><td>0</td><td>158</td><td>102</td><td>39</td><td>37</td><td>11</td><td>34</td><td>237</td><td>0.86</td></tr> <tr><td>2019</td><td>BOS</td><td>AL</td><td>286</td><td>1</td><td>7</td><td>4.35</td><td>12</td><td>12</td><td>0</td><td>68.1</td><td>55</td><td>37</td><td>33</td><td>11</td><td>19</td><td>98</td><td>1.08</td></tr> <tr><td>10 Seasons</td><td></td><td></td><td>6218</td><td>104</td><td>69</td><td>2.95</td><td>299</td><td>219</td><td>12</td><td>1550.2</td><td>1244</td><td>547</td><td>509</td><td>159</td><td>356</td><td>1887</td><td>1.03</td></tr> </tbody> </table> <p><i>Note:</i> highlight top 3 for each year, and Tm(Team) * All Star Year</p>	Season	Tm	LG	BF	W	L	ERA	G	GS	SV	IP	H	R	ER	HR	BB	SO	WHIP	2010	CWS	AL	92	2	1	1.93	21	0	4	23.1	15	5	5	2	10	32	1.07	2011	CWS	AL	288	2	2	2.79	58	0	8	71	52	22	22	6	27	79	1.11	2012*	CWS	AL	772	17	8	3.05	30	29	0	192	167	66	65	19	51	192	1.14	2013*	CWS	AL	866	11	14	3.07	30	30	0	214.1	184	81	73	23	46	226	1.07	2014*	CWS	AL	685	12	4	2.17	26	26	0	174	129	48	42	13	39	208	0.97	2015*	CWS	AL	854	13	11	3.41	31	31	0	208.2	185	88	79	23	42	274	1.09	2016*	CWS	AL	907	17	10	3.34	32	32	0	226.2	190	88	84	27	45	233	1.04	2017*	BOS	AL	851	17	8	2.9	32	32	0	214.1	165	73	69	24	43	308	0.97	2018*	BOS	AL	617	12	4	2.11	27	27	0	158	102	39	37	11	34	237	0.86	2019	BOS	AL	286	1	7	4.35	12	12	0	68.1	55	37	33	11	19	98	1.08	10 Seasons			6218	104	69	2.95	299	219	12	1550.2	1244	547	509	159	356	1887	1.03
Season	Tm	LG	BF	W	L	ERA	G	GS	SV	IP	H	R	ER	HR	BB	SO	WHIP																																																																																																																																																																																																								
2010	CWS	AL	92	2	1	1.93	21	0	4	23.1	15	5	5	2	10	32	1.07																																																																																																																																																																																																								
2011	CWS	AL	288	2	2	2.79	58	0	8	71	52	22	22	6	27	79	1.11																																																																																																																																																																																																								
2012*	CWS	AL	772	17	8	3.05	30	29	0	192	167	66	65	19	51	192	1.14																																																																																																																																																																																																								
2013*	CWS	AL	866	11	14	3.07	30	30	0	214.1	184	81	73	23	46	226	1.07																																																																																																																																																																																																								
2014*	CWS	AL	685	12	4	2.17	26	26	0	174	129	48	42	13	39	208	0.97																																																																																																																																																																																																								
2015*	CWS	AL	854	13	11	3.41	31	31	0	208.2	185	88	79	23	42	274	1.09																																																																																																																																																																																																								
2016*	CWS	AL	907	17	10	3.34	32	32	0	226.2	190	88	84	27	45	233	1.04																																																																																																																																																																																																								
2017*	BOS	AL	851	17	8	2.9	32	32	0	214.1	165	73	69	24	43	308	0.97																																																																																																																																																																																																								
2018*	BOS	AL	617	12	4	2.11	27	27	0	158	102	39	37	11	34	237	0.86																																																																																																																																																																																																								
2019	BOS	AL	286	1	7	4.35	12	12	0	68.1	55	37	33	11	19	98	1.08																																																																																																																																																																																																								
10 Seasons			6218	104	69	2.95	299	219	12	1550.2	1244	547	509	159	356	1887	1.03																																																																																																																																																																																																								

Once the dashboard is selected the user can change setting on the left side of the page. The date filter range takes dates from the start of 2010 to the start of 2020. The drop down box filters based on the outcome of a pitch. Any number of the tick boxes can be selected at the one time.



A user may click on a bright pink dot on the time series to open a video of the event, or a point on either the location plot or the spray chart if the point has an accoupling link. Points on any of the plots can be hovered over to show more information.

Conclusion

For me this is I would have liked to spend more time on and make many changes to. There are a few things that I'm proud of, the chunk of code that facets the location plot for example, or the extra json data used to link videos. Many ideas where scraped dew to a lack of time to get them working are ways to fit it all in. A few changes or things I would like to add give more time are, more filter options for the data, a way to better link all the plots and have them interact with each other. Several extra java script functions, better use of space on the page and the ability to embed the linked videos in the page as well as an option to select a view more players. Some of this could be achieved quickly in r, but to achieve the results I would like more time needs to be spent learning and implementing with D3 and JavaScript.

references

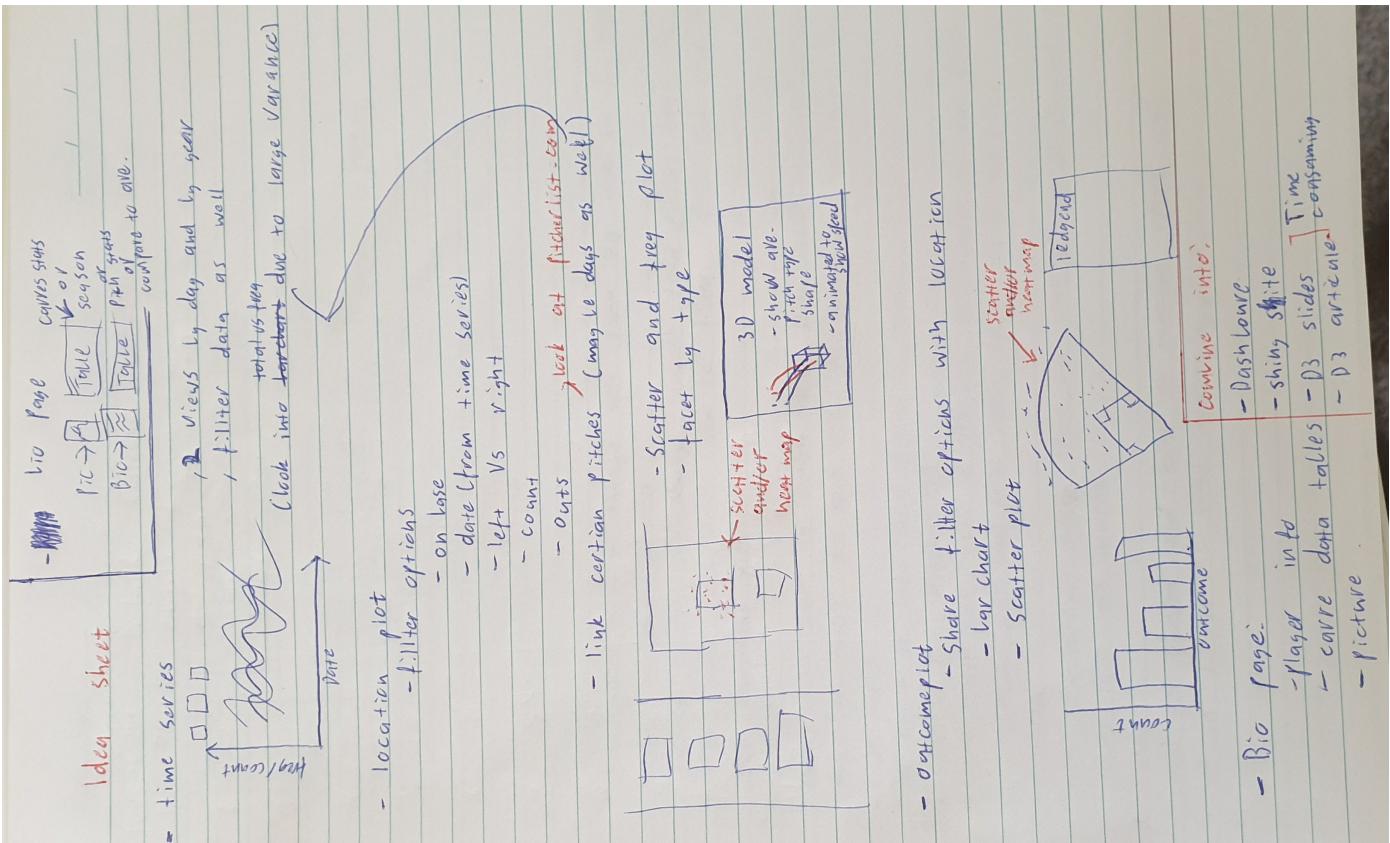
- [baseball prospectus](#)
- [baseballsavant.mlb.com](#)
- [Baseball Reference](#)
- [YouTube](#)
- [mlb.com](#)

Packge list:

- shiny
- htmlwidgets, html and javascript
- shinythemes, shiny themes
- readr, read in files
- jsonlite, read json
- dplyr, data wrangling
- tidyr, data wrangling
- data.table, data wrangling
- lubridate, data wrangling
- XML, scraping
- RCurl, scraping
- rlist, scraping
- plotly, plot
- knitr, html table
- kableExtra, html table
- RColorBrewer, colour pallet

Appendix

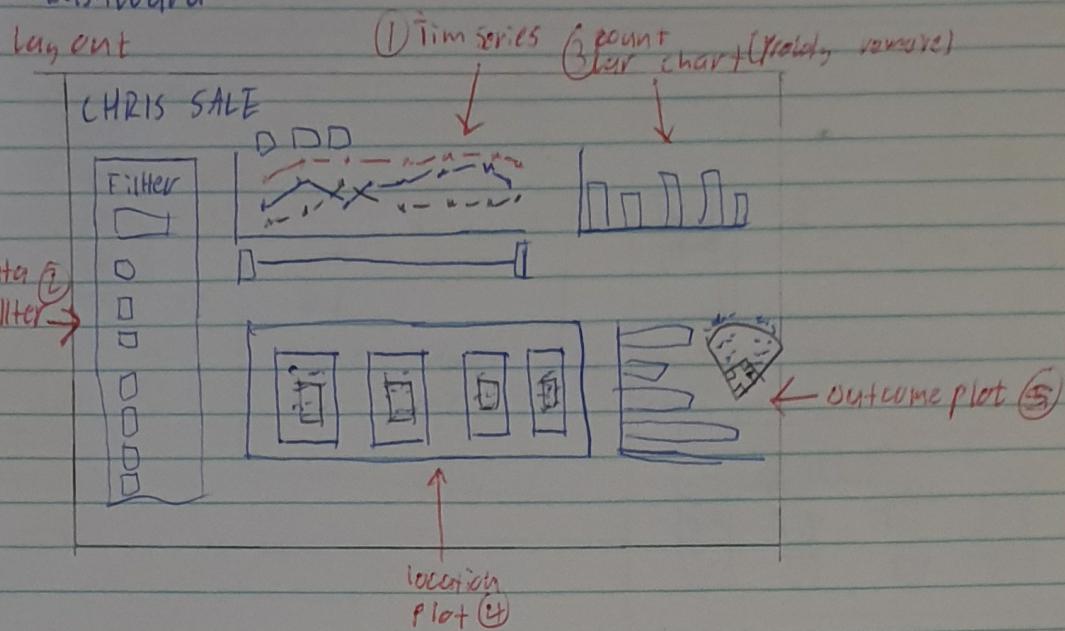
note: my five sheet method was not done 100% like the examples, as i found this way best for me to explore my ideas.



Sheet 1

Dashboard

Layout



Operations

1. Time series

- has buttons above to change view, from by year to by date, (maybe by month/quarter)
- possible option to change between total pitches and %
- slider under time series use to filter data for location and outcome plots

2. Data filter

- mostly check boxes in sections, sections being, on(0,1,2), count(0-0,0-1,1-0,-1)
- drop down select for hand, left, right, both

3. Count bar chart

- bar chart to show total pitches by year

4. Location plot, outcome plot

- same scatter plot over a frequency plot(law, alpha) to show pitch location, with some point linking to a gif of the pitch

5. Outcome plot

- same outcome plot from report, but with a baseball diamond shading where hits are located.

Discussion

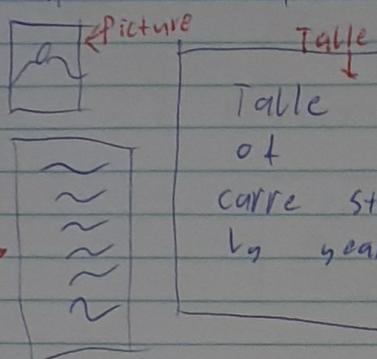
- could get cluttered and hard to manage space, -ve
- all in R, so less time consuming, +ve
- hard to explain to a user how to use

Shiny site

layout tab1 Bio

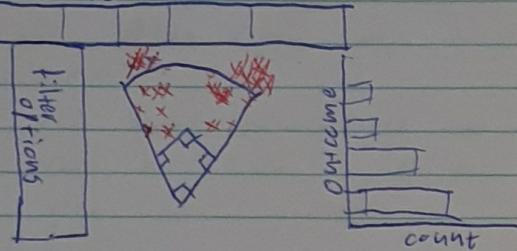
Chris Sale

tab1



layout tab4 outcome

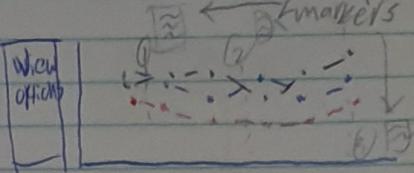
Chris Sale



layout tab2 time series

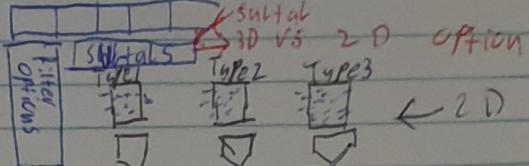
Chris Sale

tab2



layout tab3 location

Chris Sale



3D option



← 3D

- tab1 (Bio page)

- picture of chris sale
- a little bio about him, full name, age, college, Draft ~~date~~ ^{No.}, etc
- Stat summary table, ~~stats~~ ^{highlight} highlight standout stats
- maybe add radar chart to show some stats compared to ave.

- tab2 (time series)

- time series, same options as in the dashboard, with some date markers for events or notes

- tab3 (location)

Subtab 1 - 2D: location plot same as in dashboard, with filter options on the left

Subtab 2 - 3D: a 3D version of the location plot

- tab4 (outcome)

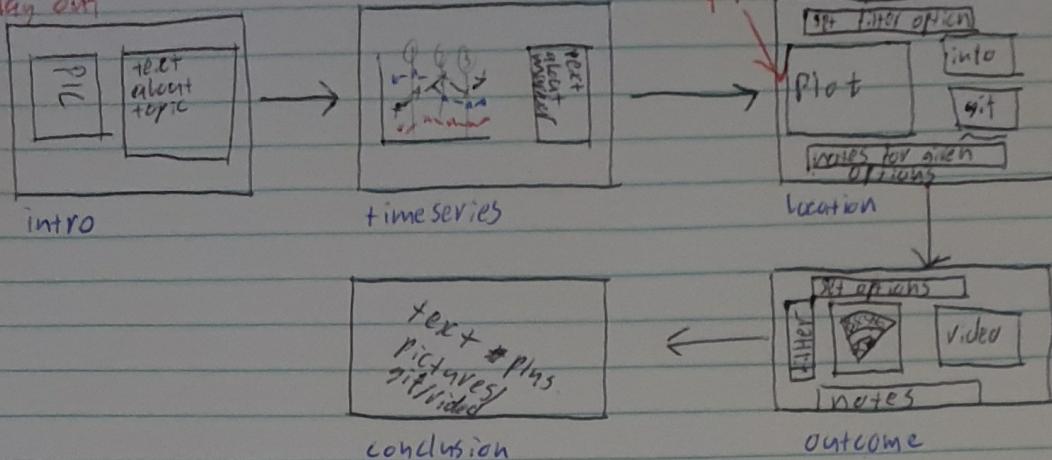
- outcome plot, same as Dashboard, filter options on the left, like location

Discussion

- again all R, good for time, +ve
- hard to show one story, -ve
- filter options, might not be shared

D3 article/slide show

lay out



Intro

- would show and introduce who chris sale is
- Intro @ the topic abit about each slide and a quick link to each timeseries
- the same timeseries as the others
- when an event is clicked on explain the event and why it's important

location

- same location plot
- when filter settings to the left, also the same
- when a point is clicked on display info about it as well as a gif of the pitch/pitch type
- Some predefined filter options, with notes on there importants

Outcome

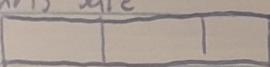
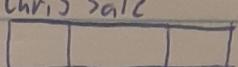
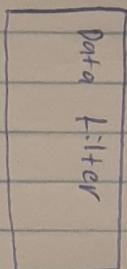
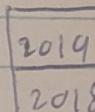
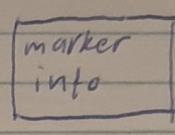
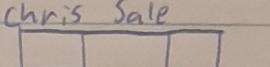
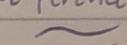
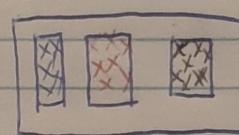
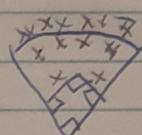
- Similar set up to location, with filter, set options ~~and~~ or video of
- q - video of hit [probably just home runs]
- might put info if extra data can be sourced

Conclusion

- summarise findings
- references
- glossary

Discussion

- uses D3 so takes more time / -ve
- more info, +ve
- stronger narrative, +ve
- better looking

Chris Sale   Bio	Chris Sale   Data Filter
 2019  2018 Career summary	 marker info
tab 1  Glossary  References 	tab 2  

- Use `plotly` in R to create the time series and scatter plot, due to its slider, buttons to change view and link points to a url.
 - outcome plot `ggplot` plus a package to map the shape of a baseball field, `geomMLBStadiums`
 - use `shiny` to host/reactive
 - scrape summary table from `MLB.com`
 - glossary terms from, Baseball Reference
 - Radar plot with `ggplot`
 - makers info from `MLB.com` articles
 - source ang gif's from `pitcherlist`
 - Viedos from `MLB.com`