



UNIVERSITY OF  
TEXAS  
ARLINGTON

**INSY-5337 WEB AND SOCIAL ANALYTICS  
PROJECT REPORT**

# TABLE OF CONTENTS

<b>INTRODUCTION:</b> .....	<b>3</b>
<b>APPROACH:</b> .....	<b>3</b>
<b>DATA DESCRIPTION &amp; SOURCE:</b> .....	<b>3</b>
DATA SOURCE: .....	3
DATA DESCRIPTION: .....	5
RESEARCH QUESTIONS: .....	6
<b>METHODOLOGY</b> .....	<b>7</b>
DATA CLEANING: .....	7
CALCULATED FIELDS: .....	7
EXPLORATORY ANALYSIS: .....	7
GOOGLE TRENDS: .....	16
MODEL BUILDING .....	18
<b>INTERPRETATION</b> .....	<b>23</b>
<b>RESULTS &amp; CONCLUSION:</b> .....	<b>23</b>

## **INTRODUCTION:**

We are a brand-new company planning to launch a new energy drink. **The Business goal is to make the brand name reach a large population as a first step.** Many celebrities do paid partnerships with companies to promote products on social media platforms such as Facebook, Instagram, and Twitter. We are in search of such a celebrity **to launch and promote the product.**

## **APPROACH:**

In the present world, it is easy to introduce a product, but the major problem lies in getting the product to customers' hands and making people all over the world made aware of the product. So, to reach a wide variety of audiences, we will be taking an athlete with a huge fan following as well as a good reach for their posts on Facebook.

We tried mining the required data from Facebook. But, when we attempted to mine data it gave us a permission error, even after providing permissions for all the things which were mentioned on Facebook. Turns out, we can only access data on our pages. Data extraction from all other pages has been restricted from version 2.8 and the API doesn't allow us to use prior versions.

We overcame the above problem with the help of an online analytics tool called '**popsters**'. In short, popsters is a tool that helps to get compare stats and measure the efficiency of posts on any social media page that we are interested in.

## **DATA DESCRIPTION & SOURCE:**

### **DATA SOURCE:**

The below dataset was created by our **team manually**. We have filtered the top 50 cricket players who were all **active on social media**. And for each player, using their profile link and with the help of the popsters tool, we filled the values for all this column. We have 22 columns of data. Out of which, 3 are Qualitative variables and the remaining 19 are Quantitative variable

We have fixed a Period of **6 Months**, starting from **October 2021 to April 2022**. Though most of the columns were directly imported from the online tools, a few other columns were calculated by us. For examples fields such as Promotional post, Promotional Post Percent, and Average Post Frequency. A detailed description of every column with its uses and formula was given below the dataset in a separate table.

Players	Node ID	Profile Link	Followers	Posts	Likes	Shares	Comments	Avg_Likes	Avg_Shares	Avg_Comments	Likes Ratio	Shares ratio
DAVID MILLER	100044343763446	<a href="https://www.facebook.com/100044343763446">https://www.facebook.com/100044343763446</a>	5493502	6	369803	317	2389	61633	52	398	1.12190%	5.77045E-05
YUVENDRA SINGH CHAHAL	0.03751E-05	<a href="https://www.facebook.com/0.03751E-05">https://www.facebook.com/0.03751E-05</a>	1	2	2	0	0	1	0	0	0.00000%	0.00000E+00
MOAHAMMED SIRAJ	100044179445847	<a href="https://www.facebook.com/100044179445847">https://www.facebook.com/100044179445847</a>	1358079	1	36731	202	1090	36731	202	1090	2.7040%	0.00014871
PRITHVI SHAW	0.0004439566094	<a href="https://www.facebook.com/0.0004439566094">https://www.facebook.com/0.0004439566094</a>	2242883	2	172777	2097	371	88638	1049	185	7.90398%	0.000934957
BHUVENSHWAR KUMAR	10004431735294	<a href="https://www.facebook.com/10004431735294">https://www.facebook.com/10004431735294</a>	10376635	2	32981	103	2536	16490	51	1268	0.1599%	9.92815E-06
YUVRAJ SINGH	10004428688682	<a href="https://www.facebook.com/10004428688682">https://www.facebook.com/10004428688682</a>	19054947	5	81059	3124	1787	16211	624	353	0.42517%	0.000163861
SAM CURRAN	100057260179768	<a href="https://www.facebook.com/100057260179768">https://www.facebook.com/100057260179768</a>	2547182	2	117275	250	2195	58637	125	1097	4.60411%	9.81477E-05
AJINKYA RAHANE	100044289256775	<a href="https://www.facebook.com/100044289256775">https://www.facebook.com/100044289256775</a>	8476878	4	13591	43	378	3397	10	94	0.04010%	5.07262E-06
DAVID WARNER	100044542354906	<a href="https://www.facebook.com/100044542354906">https://www.facebook.com/100044542354906</a>	7695287	2	4681	33	501	2340	16	250	0.06083%	4.28834E-06
RAHUL TRIPATHI	10005772972933	<a href="https://www.facebook.com/10005772972933">https://www.facebook.com/10005772972933</a>	553257	8	204413	1546	4004	25551	193	500	36.94721%	0.002794361
RAHUL TEWATIA	100046826631970	<a href="https://www.facebook.com/100046826631970">https://www.facebook.com/100046826631970</a>	1221964	16	129640	93	1155	8102	5	72	10.60915%	7.77437E-05
SANJU SAMSON	100044281486786	<a href="https://www.facebook.com/100044281486786">https://www.facebook.com/100044281486786</a>	2935248	29	813791	1565	1553	28061	53	536	27.72478%	0.000533175
RAVICHANDRAN ASWIN	100044281486786	<a href="https://www.facebook.com/100044281486786">https://www.facebook.com/100044281486786</a>	4074849	78	203638	2441	4131	2611	31	52	4.99763%	0.000964511
UMESH YADAV	1000445094765801	<a href="https://www.facebook.com/1000445094765801">https://www.facebook.com/1000445094765801</a>	6211195	8	9365	10	1273	4655	8	199	0.00000%	1.17071E-05
ROBIN UTHAPPA	100044339531916	<a href="https://www.facebook.com/100044339531916">https://www.facebook.com/100044339531916</a>	8886016	86	99735	1476	5711	11596	17	66	0.40140%	0.000910901
JASPRIT BUMRAH	100044446160940	<a href="https://www.facebook.com/100044446160940">https://www.facebook.com/100044446160940</a>	8396089	37	175746	5870	20068	47490	168	542	18.70281%	0.000634795
KULDEEP YADAV	100044197371116	<a href="https://www.facebook.com/100044197371116">https://www.facebook.com/100044197371116</a>	1947454	23	153799	345	3736	6686	15	162	7.89744%	0.000177154
ROHIT SHARMA	1000442402905060	<a href="https://www.facebook.com/1000442402905060">https://www.facebook.com/1000442402905060</a>	20931904	20	1332768	5011	40543	66838	250	2027	6.36716%	0.000239395
AB De	100044338481815	<a href="https://www.facebook.com/AbDeVilles17">https://www.facebook.com/AbDeVilles17</a>	13216242	43	4768558	53539	105897	110998	1245	2462	36.08104%	0.004051
MS DHONI	10004481937493	<a href="https://www.facebook.com/10004481937493">https://www.facebook.com/10004481937493</a>	27214451	20	647915	2018	24248	32395	1040	1212	2.38078%	0.000764888
RASHID KHAN	100044281426914	<a href="https://www.facebook.com/100044281426914">https://www.facebook.com/100044281426914</a>	2729178	25	2036184	10649	146942	82527	425	5877	75.59727%	0.003901907
CHRIS GAYLE	100044317545336	<a href="https://www.facebook.com/100044317545336">https://www.facebook.com/100044317545336</a>	13851349	10	93404	380	3544	9340	38	354	0.67433%	2.74342E-05
HARDIK PANDYA	100044386001022	<a href="https://www.facebook.com/100044386001022">https://www.facebook.com/100044386001022</a>	11907722	23	1241491	4456	23454	53977	193	1019	10.42593%	0.000374211
Sachin Tendulkar	1000442420919512	<a href="https://www.facebook.com/1000442420919512">https://www.facebook.com/1000442420919512</a>	37575565	111	6520209	177044	170306	59027	1594	1534	17.43692%	0.004711679
RSIHABH PANT	100057277159152	<a href="https://www.facebook.com/100057277159152">https://www.facebook.com/100057277159152</a>	5582288	42	100768	3096	18924	23884	73	450	17.77786%	0.000646111
SHOBH DHAWAN	100044281486786	<a href="https://www.facebook.com/100044281486786">https://www.facebook.com/100044281486786</a>	1744774	54	942642	1369	35332	3335	253	654	0.00000%	0.00079623
KRUNAL PANDYA	100044260937684	<a href="https://www.facebook.com/100044260937684">https://www.facebook.com/100044260937684</a>	2262069	19	181828	128	2115	8517	22	111	7.15388%	0.000189207
AKSHAR PATEL	100044547304948	<a href="https://www.facebook.com/100044547304948">https://www.facebook.com/100044547304948</a>	2360426	35	103415	155	2213	2926	4	63	4.33884%	6.66661E-05
RAHUL CHAHAR	100018555806290	<a href="https://www.facebook.com/100018555806290">https://www.facebook.com/100018555806290</a>	625181	38	849986	878	6916	22368	23	182	161.84820%	0.001671805
SHREYAS GOPAL	10004447798813	<a href="https://www.facebook.com/10004447798813">https://www.facebook.com/10004447798813</a>	224754	21	148088	157	1069	7051	7	50	3.1740%	0.000698542
KL RAHUL	100044235379214	<a href="https://www.facebook.com/100044235379214">https://www.facebook.com/100044235379214</a>	10284388	54	3027838	11552	71043	56071	213	1315	29.44111%	0.001123256
Virat Kohli	1000447308573	<a href="https://www.facebook.com/1000447308573">https://www.facebook.com/1000447308573</a>	49145016	97	24339754	166883	981214	250925	1720	10115	49.52639%	0.003395726
Surendra Raina	100044298046085	<a href="https://www.facebook.com/100044298046085">https://www.facebook.com/100044298046085</a>	6043506	86	234764	10220	32244	27298	118	374	38.84576%	0.001691071
RAVIN德拉 JADEJA	1000445445196213	<a href="https://www.facebook.com/1000445445196213">https://www.facebook.com/1000445445196213</a>	10484998	30	167153	9211	33529	55717	307	1117	15.41297%	0.000849332
MAYANK AGARWAL	100044343719233	<a href="https://www.facebook.com/100044343719233">https://www.facebook.com/100044343719233</a>	238770	67	469282	600	5037	7004	8	75	19.67829%	0.000251597
SHREYAS IYER	100044290546403	<a href="https://www.facebook.com/100044290546403">https://www.facebook.com/100044290546403</a>	2688432	27	52847	2345	14048	21216	86	520	21.30785%	0.000872256
ISHAN KISHAN	100044148800557	<a href="https://www.facebook.com/100044148800557">https://www.facebook.com/100044148800557</a>	3727674	21	63249	1752	15440	30118	83	735	16.96744%	0.000469998
LAITH MALINGA	10004392598321	<a href="https://www.facebook.com/10004392598321">https://www.facebook.com/10004392598321</a>	3592181	77	61835	11017	55279	8037	143	717	0.22370%	0.003069639
HARBHAIAN SINGH	10004414374060	<a href="https://www.facebook.com/10004414374060">https://www.facebook.com/10004414374060</a>	1024052	92	635480	4876	35922	9081	53	390	0.08680%	0.000476915
WASHINGTON SUNDAR	100051414284787	<a href="https://www.facebook.com/100051414284787">https://www.facebook.com/100051414284787</a>	509723	72	797479	1920	923	11038	26	128	155.91782%	0.003768752

Players	Comments Ratio	Promotional Posts	Promotional Post percent	Average Post Frequency	Engagement Rate / post	Engagement Rate / day	Highest Active Day	Highest day active percent	Time of the day
DAVID MILLER	0.0010434877	3/20	10	0.032396703	1.12%	0.00%	Sunday	40.00%	03:00
YUVENDRA SINGH CHAHAL	0.03751E-05	2/2	100	0.010988011	0.54%	0.00%	Saturday	85.15	08:00
MOAHAMMED SIRAJ	0.000802604	0/30	0	0.0054945505	2.80%	2.80%	Tuesday	100%	03:00
PRITHVI SHAW	0.0001654512	15/20	75	0.010988011	4.01%	1.00%	Wednesday	100%	02:00
BHUVENSHWAR KUMAR	0.000244395	1/20	5	0.010988011	0.07%	0.00%	Wednesday	98.02%	10:00
YUVRAJ SINGH	0.26832E-05	8/20	40	0.027472527	0.69%	0.04%	Friday	72.13%	23:00
SAM CURRAN	0.0008617137	0/2	0	0.010988011	2.35%	1.57%	Thursday	64.36%	04:00
AJINKYA RAHANE	0.45919E-05	2/20	10	0.021970022	0.04%	0.00%	Wednesday	50.99%	06:00
DAVID WARNER	6.51048E-05	2/2	100	0.010988011	0.03%	0.02%	Thursday	40.00%	13:00
MAHENDRA PANT	0.000196448	4/48	50	0.043956044	0.88%	0.22%	Friday	37.55%	14:00
ROHIL CHAHAR	0.0001969448	3/10	30	0.126573628	0.46%	0.06%	Friday	26.49%	21:00
ROBIN UTHAPPA	0.001977988	2/20	10	0.0472527473	0.40%	0.21%	Thursday	34.31%	14:00
JASPRIT BUMRAH	0.002135797	1/20	5	0.203296703	0.51%	0.11%	Tuesday	32.64%	05:00
KULDEEP YADAV	0.001919482	2/20	10	0.126373626	0.35%	0.30%	Tuesday	31.57%	09:00
ROHIT SHARMA	0.00169369	8/20	40	0.10988011	0.33%	0.04%	Sunday	30.73%	05:00
AB De	0.008012641	8/10	80	0.236263736	0.87%	0.22%	Thursday	30.72%	13:00
MS DHONI	0.000809097	14/15	93	0.133333333	0.10%	0.02%	wednesday	28.64%	05:00
RASHID KHAN	0.05384112	3/15	20	0.137362637	3.25%	0.52%	sunday	27.87%	13:00
CHRIS GAYLE	0.000196448	1/15	50	0.05384112	0.70%	0.10%	Wednesday	27.50%	07:00
HARDIK PANDYA	0.001969448	3/10	30	0.126573628	0.46%	0.06%	Friday	26.49%	21:00
Sachin Tendulkar	0.004532316	4/20	20	0.06088011	0.17%	0.10%	Thursday	25.38%	09:00
RSIHABH PANT	0.003380098	10/20	50	0.230769231	0.44%	0.11%	Sunday	25.21%	03:00
SHIKHAR DHAWAN	8/20	40	0.029670329	0.20%	0.08%	Sunday	24.94%	09:00	
KRUNAL PANDYA	5/20	25	0.104395604	0.38%	0.04%	Tuesday	24.84%	02:00	
AKSHAR PATEL	0.0009347543	7/20	35	0.192307692	0.13%	0.15%	Saturday	24.12%	06:00
RAHUL CHAHAR	0.013168793	8/20	40	0.208791209	4.30%	0.00%	Monday	23.64%	03:00
SHREYAS GOPAL	0.004736511	3/20	15	0.115384615	3.16%	0.38%	Thursday	23.52%	01:00
KL RAHUL	0.000907849	6/20	30	0.296703297	0.56%	0.17%	Friday</td		

## DATA DESCRIPTION:

Column Name	Description
Players	Name of the Cricket Players
Node ID	Unique Node ID which represents the Player
Profile Link	The Players Facebook Profile Link
Followers	The Total Number of followers that the player has on their Facebook Profile as of April 2022
Posts	The Total Number of posts published on the Player's profile as of April 2022
Likes	The Total Count of like for all the posts combined as of April 2022. Likes Include all the reactions from the Audience
Shares	The Total Count of Shares for all the posts published as of April 2022. Shares of Shares will also come under this category.
Comments	The Total Count of Comments for all the posts combined as of April 2022.
Average Likes	Average Likes are calculated by the formula Total number of likes / Total number of posts
Average Shares	Average Likes are calculated by the formula Total number of Shares / Total number of posts
Average Comments	Average Likes are calculated by the formula Total number of Comments / Total number of posts
Likes Ratio	The likes Ratio is calculated by the formula Total Number of Likes / Total Number of Followers
Shares Ratio	The shares Ratio is calculated by the formula Total Number of Likes / Total Number of Followers
Comments Ratio	Comments Ratio is calculated by the formula Total Number of Likes / Total Number of Followers
Promotional Post	The number of Promotional Posts to the Number of Total Posts. Total Post include both personal posts and promotional posts
Promotional Post Percentage	The Promotional Post is converted into the form of a Percentage
Average Post Frequency	To Find How Frequently player posts on their social media handle. Calculated by the formula Total Number of Post/ Number of Days
Engagement Rate / Post	Engagement Rate is a metric that shows how many followers were engaged with the published posts. $ERpost = (\text{likes} + \text{shares} + \text{comments}) / \text{count of followers}$
Engagement Rate / Day	$ERday = (\text{total likes} + \text{total shares} + \text{total comments for the day}) / \text{count of followers}$

Highest Active Day	The Day on which the Player is highly Active. It helps in finding out the pattern.
Highest Day Active Percent	The Percentage of Activity of the most active day of the week
Time of the Day	The Time of the Day when there is a sudden burst of users in the account

## RESEARCH QUESTIONS:

To ensure that the athlete meets our requirements, we have framed the following set of questions. The questions are framed in such a way that it will not only help us in finding the exact match but also will help us in finding the outliers.

- What are the average Likes, Comments, and shares per post for the last 6 months period?
- What is the Number of views to several follower's ratio for all the players?
- What are the Like Ratio, Comment Ratio, and Share Ratio per post?
- What is the Number of Promotional Posts that the players have posted?
- Identifying the type of products promoted and frequency of posts for that product
- What is the total number of posts saved by the users?
- Which type of products reach a wider audience? (Clothing, accessories, food, etc.)
- At what point in time there is an increase in users on the player's Facebook Profile?
- At what day there is an increase in users on the player's Facebook Profile?
- How frequently do the players posts?

Considering all the above questions, we have decided to have the columns presented in the dataset. All the columns such as likes, comments, shares, followers, or engagement focuses on **measuring the reach** of a particular post by a player. Analysis and modelling are done based on the collected data to determine the appropriate brand ambassador for our product – Energy Drink.

## METHODOLOGY

### DATA CLEANING:

The dataset which we used is **created by our team manually** so, we made sure that there are **no missing values** and while collecting data of 3 players we observed that they have not posted any promotional related posts almost for the past one year. Therefore, we have removed those players as they will not provide any value to the product promotion. We have collected enough information for each player, so we analyze various factors contributing to the reach of a post.

### CALCULATED FIELDS:

In the dataset, a few columns were collected directly from the tool. Excel computations are used for the following columns:

- **Likes / Comments / Share ratio** - computed with their respective field and follower's count.
- **promotional post percentage** - promotional posts /total posts
- **Average post frequency** - Average of Posts for the stipulated time.

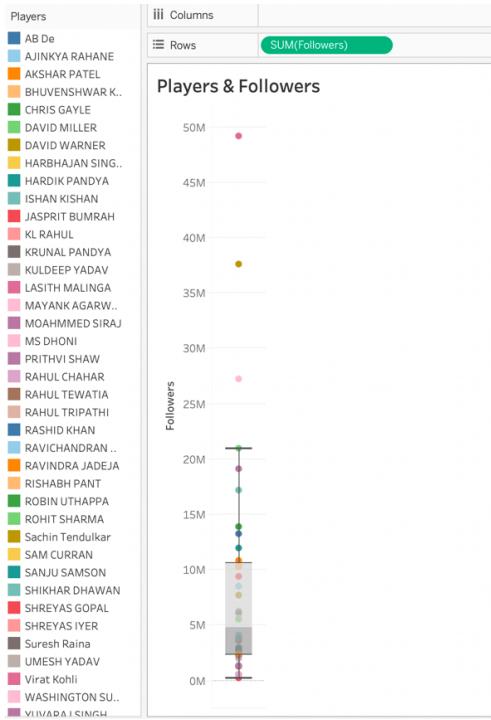
### EXPLORATORY ANALYSIS:

We have performed many types of exploratory analysis **for better understanding of the data**. Based on our analysis we were able to get an idea of who might be our **potential brand ambassadors**. The following depicts the different visualization used.

### FOLLOWER DISTRIBUTION:

We use box plot to determine the relationship between the players and the follower distribution. The box plot also known as a box-and-whisker plots shows the distribution of

values along an axis. The boxes indicate the middle 50% percent of the data which is the middle quartiles of the data's distribution. It also portrays the distribution of our data, outliers and median. The box plot here is mainly used to determine the players and follower distribution. We distribute various players representing different countries based on the follower count.



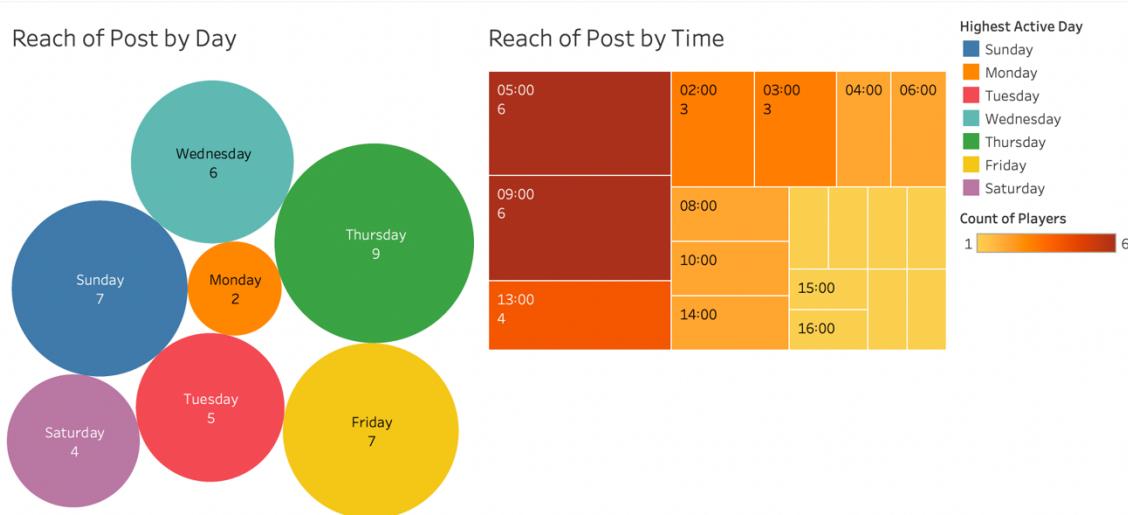
In the image we can find the players with huge follower count viz **Virat Kohli, Sachin Tendulkar and MS Dhoni would be depicted as outliers** because the primary goal of this project is to pick a celebrity who can be cost efficient and someone who can reach a wider spectrum of audiences. It is difficult to pick someone from the mentioned players as they would charge more than our intended budget as they are the superstars of their respective fields. It is evident that we can't chose someone from the outliers.

Hence, we decide to into account the upper 25% band from the box plot. The players we concentrate mainly here are Rohit Sharma, AB de Villiers, Jadeja, Chris Gayle, Yuvraj Singh etc. We must consider that there are no negative outliers in the box plot because we take the followers count into consideration.

## POST REACH DISTINGUISHED BY DAY & TIME:

The purpose of bubble chart is to display data as a cluster of circles. Each of the value in the dimension field represents a circle whereas the value of measure represents the size of those circles.

Tree map is used to visualize data that can be hierarchical for comparative analysis. Tree map plays an important part to analyze the anomalies in the data set.



We chose the bubble graph and the tree map to show the distribution of engagement of posts on social media and the maximum reach it can have in terms of days and time.

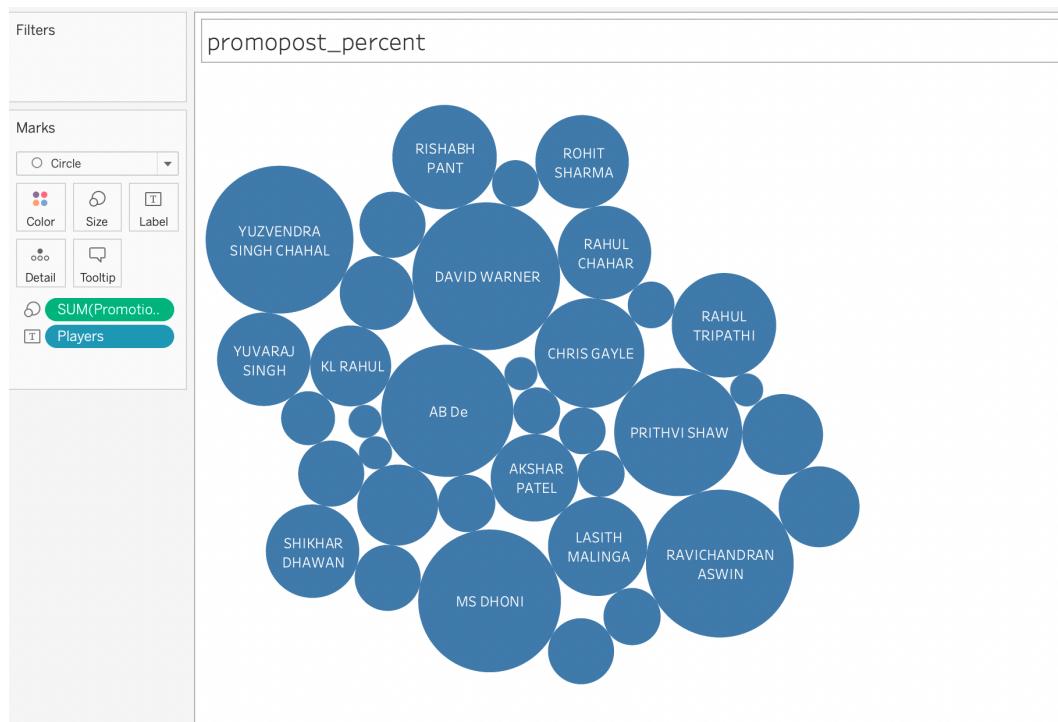
The bubble graph here displays seven days of the week, we use bubble graph here to find out which day of the week has the highest reach of posts on social media.

The graph indicates that **Thursday, Friday, and Sunday** the reach for the posts on social media are high compared to other days of the week. It is understood that many people use social media during weekends more than may be the reason why reach is pretty high.

The tree map implies that the reach ratio is higher between **05:00 AM and 09:00 AM**. It is because people wake up generally at time. As the current generation access social media, the first thing they wake up, that is the reason why the reach of social posts is comparatively high during that time frame.

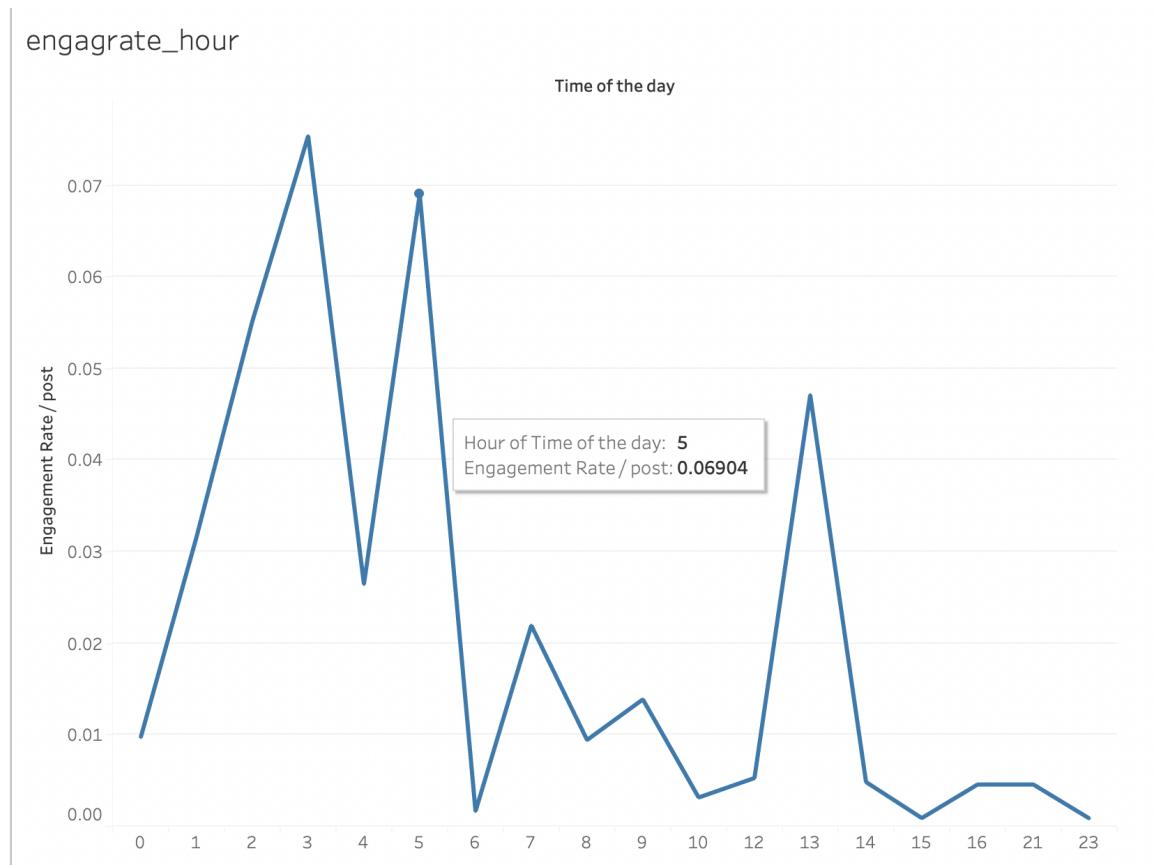
## PLAYERS BY PROMOTION:

We use bubble chart again to analyze which players does more in comparison to others. This to get a basic idea of how the **individual players fare in endorsements**. By the bubble chart below we can see that David Warner, Chahal, Dhoni and Ashwin endorse more compared to their contemporaries.



## ENGAGEMENT RATE:

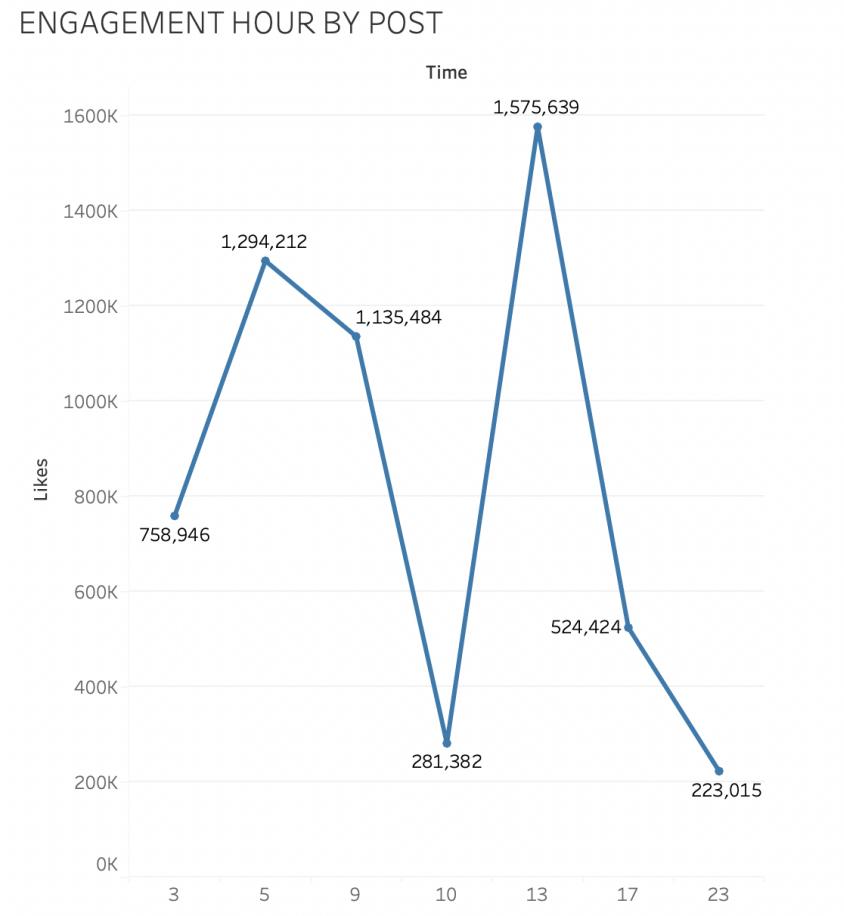
When comparing the engagement rate with hours of the day, we can find that the maximum reach of posts peak at **3:00 AM in the morning**. Intuitively, we can say that 3:00 AM posts are seen by the **Gen Z generation**. The posts at morning **5:00 AM are accessed by the millennials**. The rest of people access social media between 5:00 AM and 13:00 PM. The thirteenth hour marks the hour after lunch time or the lunch time.



## PROOF:

To confirm whether our hunch is right, we took a post of AB de Villiers for a period of 5 days and cumulated the count of likes for the timings - 3:00,5:00, 10:00,13:00,17:00,23:00.

**This analysis confirmed that the likes are highest at 13:00 AM and the second highest being 5:00AM as shown in the graph.**



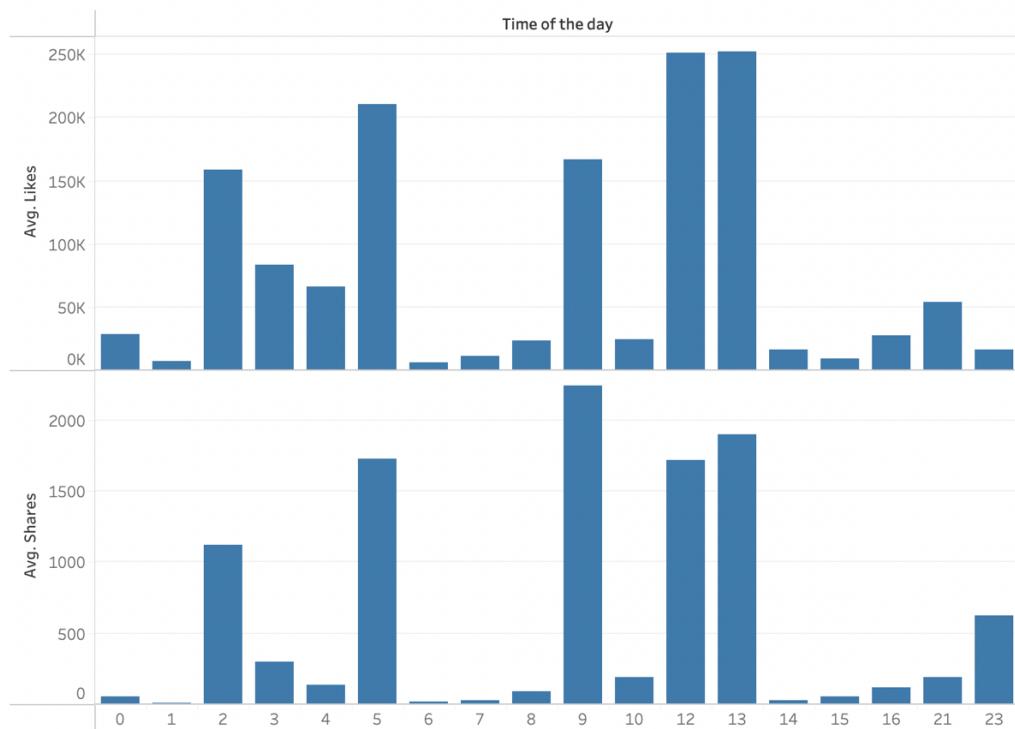
The graph also confirms that the highest engagement rates are 3:00, 5:00, 9:00 and 13:00, the pattern is same as analyzed in our whole dataset.

### **ANALYSIS BY LIKES AND SHARES:**

When endorsing a brand or product, we realized that just the promotional posts alone doesn't suffice to reach a large chunk of people. The reach is best measured by the shares

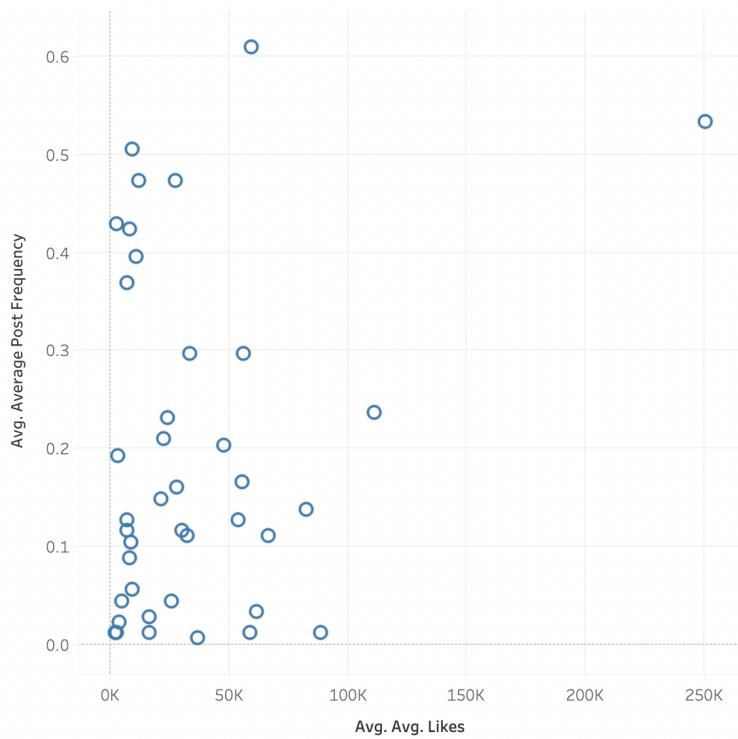
made on social media. Each Post shared would compound the reach made to a large group of people indicating an exponential growth of a product or a brand.

As we have seen in the engagement rate, we noticed that there are a higher number of likes and shares at 3:00 AM, 5:00 AM, 9:00 AM, 12:00 AM and 13:00 PM respectively. **Therefore, as a part of our strategy, we believe that we should promote our product during this time frame.**

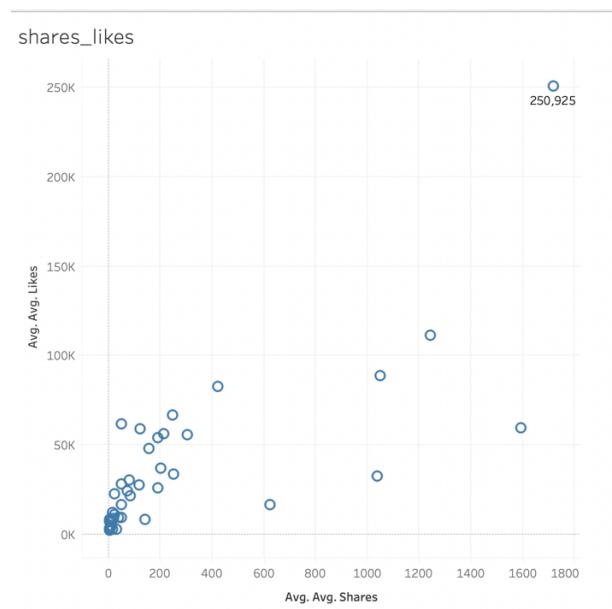


We use scatter plot to visualize the relationship between numerical variables. You create a scatter plot by placing at least one measure on the columns shelf and one measure on the rows shelf. We use the scatter plot her to determine the relationship between the likes to the post frequency.

In the scatter plot below, we can find that **Virat Kohli** dominates the entire players in all the major categories of **likes/post frequency and has the highest number of reaches**. AB De Villiers, from the scatter plot, has a moderate number of posts and has a better entanglement of posts with the reach compared to the other contemporary players.



When comparing the reach of endorsement by considering the number of likes and shares with respect to the players. We identify that Virat Kohli and Sachin Tendulkar has the highest ratio. We also realized that apart from AB DE Villiers; MS Dhoni, Prithivi Shaw and Yuvraj Singh have a great reach for their posts. We concluded that the rest of the players kind of have the same effect in terms of reach.



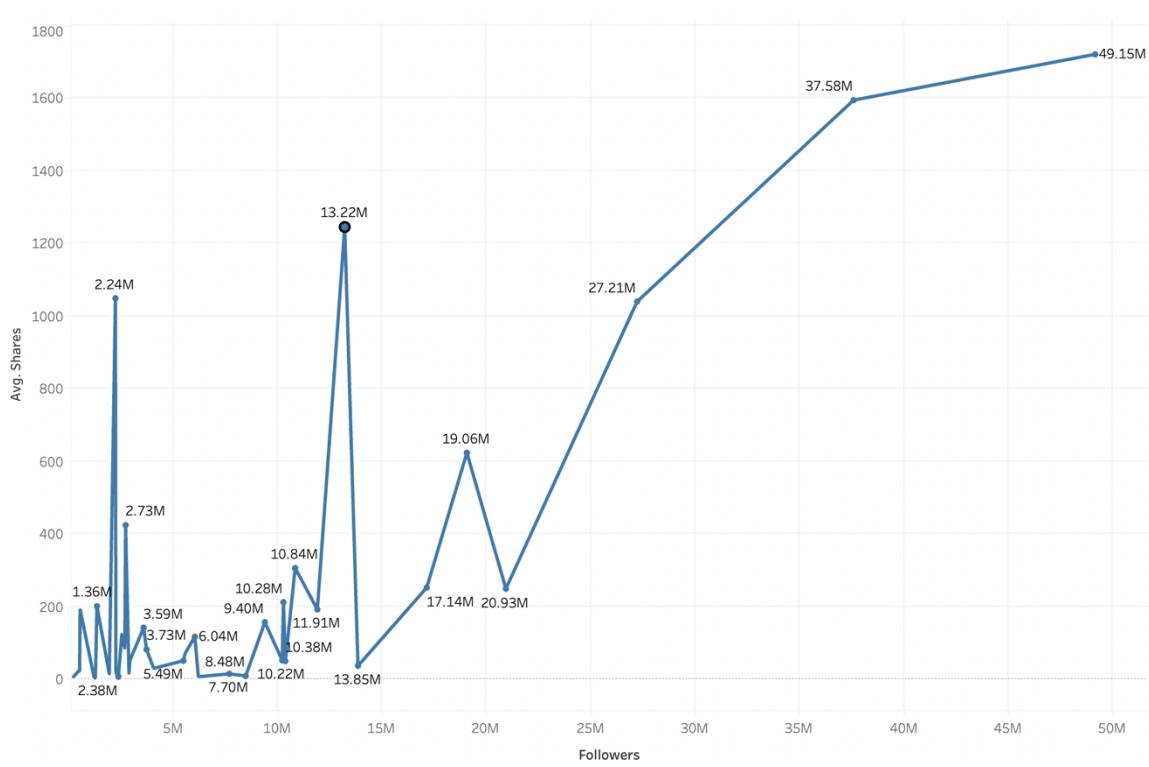
Having a huge number of followers doesn't necessarily infer that it can guarantee reach.

**Share is considered as a good form of measurement to understand a post's reach.**

Therefore, followers and posts are analyzed together on how far a post about a product can reach to different section of audiences.

We found out that apart from the obvious highest,

- AB DE Villiers(13.2 Million) has a great reach for his posts even though the followers are comparatively less than the top few players.
- Prithivi Shaw (2.24 Million) has a decent number of reach for his promotional posts on social media with a relatively less number of followers.

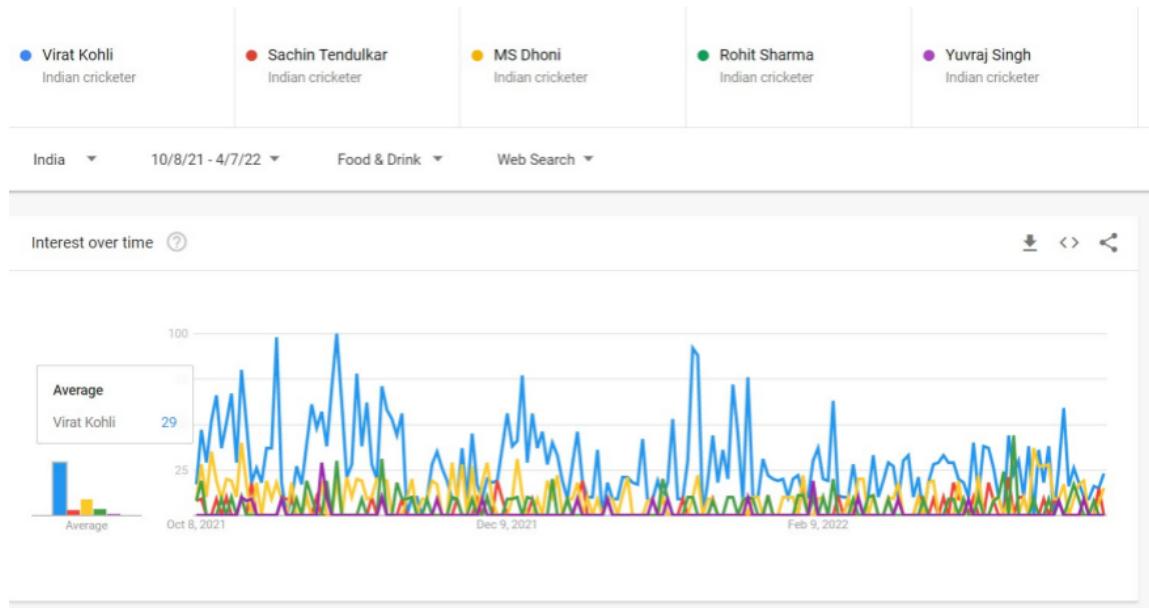


## GOOGLE TRENDS:

To extend the exploratory analysis we have used google trends and through this website, we will be able to analyze, the popularity of the top search queries in google search across various regions and languages and it uses graphs to compare the search volume of different queries over time.

Here, we have taken the **top 10 players** based on their follower count, we have searched the trend of their name for the period of 6 months i.e., from 8<sup>th</sup> October 2021 to 7<sup>th</sup> April 2022 which is the same data as our dataset, and we can select a specific category of search as our product is an energy drink, **we chose food & drinks category for analysis**.

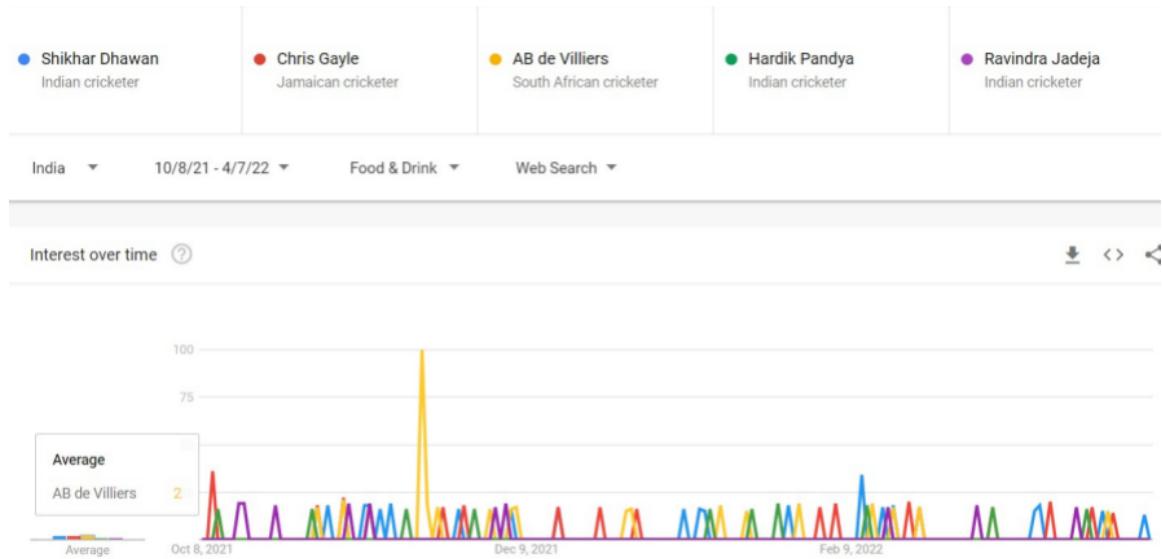
We have started our trend analysis for the first top 5 players they are Virat Kohli, Sachin Tendulkar, MS Dhoni, Rohit Sharma, and Yuvraj Singh & we observed that Virat Kohli is dominating as well among those players.



Screenshot of comparison graph of first top 5 players from google trends:

Here the colour blue indicates Virat Kohli and we can observe that he has high spikes overall when compared with other players and an **average of 29** which is the highest and that we can be able to analyze from the bar graph on left.

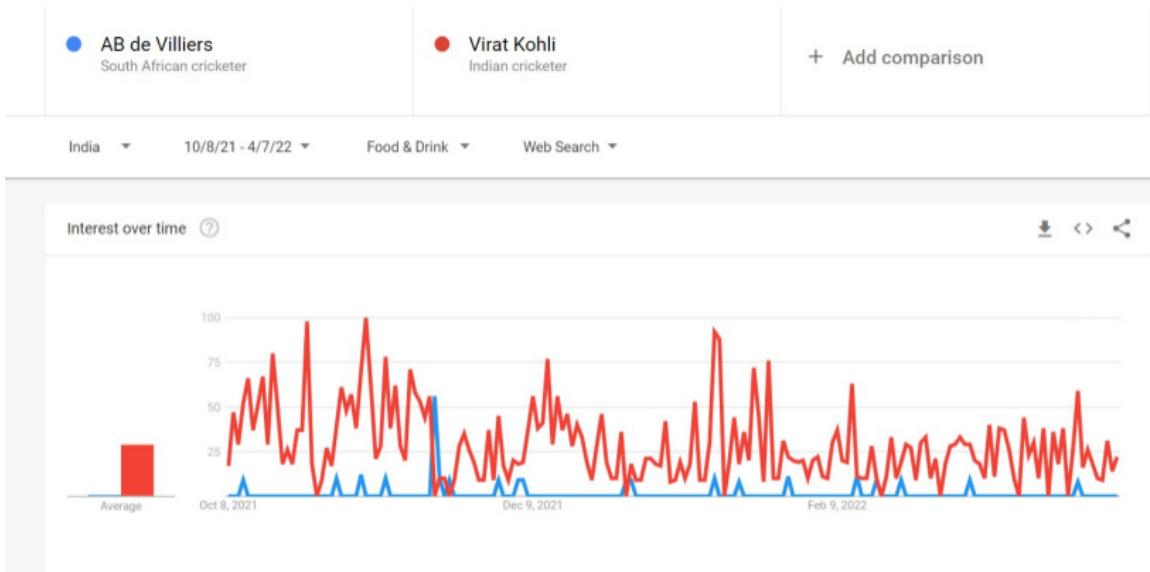
Next, we continued trend analysis with the next top 5 players' comparison they are Shikhar Dhawan, Chris Gayle, AB de Villiers, Hardik Pandya, and Ravindra Jadeja & we observed that AB de Villiers has an edge over other players.



#### Screenshot of comparison graph of next top 5 players from google trends:

Here, the yellow colour represents AB de Villiers and we can see that he has little high spikes than other players with an average of 2 which is slightly higher than and that we can figure out from the bar graph on left and we have searched google what caused his high spike and we figured that he announced his retirement during that time.

As now we have compared the top 10 players, we continued with the top 2 players from both comparisons and analyzed them for better decision making and those players are Virat Kohli and AB de Villiers



Screenshot of comparison graph of Virat Kohli and AB de Villiers:

Here, the red represents Virat Kohli and blue AB de Villiers, and we can observe that a single spike of AB de Villiers is matching with Virat Kohli.

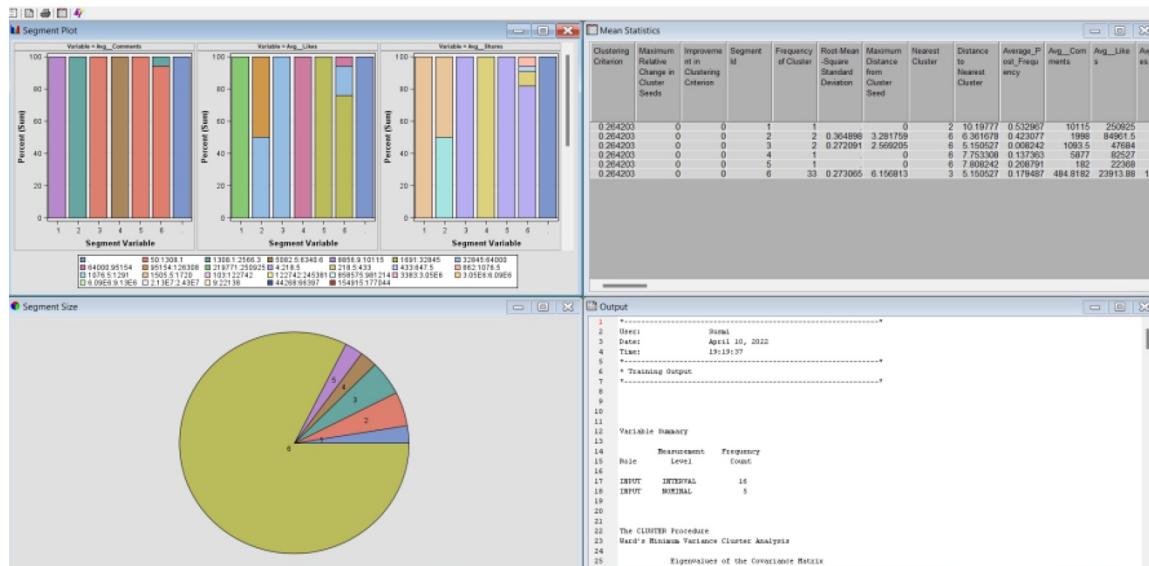
## MODEL BUILDING

The major objective of our business problem is to find a cricket player who would be effective in promoting a product to the public. This is done based on the collected dataset of the top cricket players and by considering the various factors affecting the reach of a promotional post made by them.

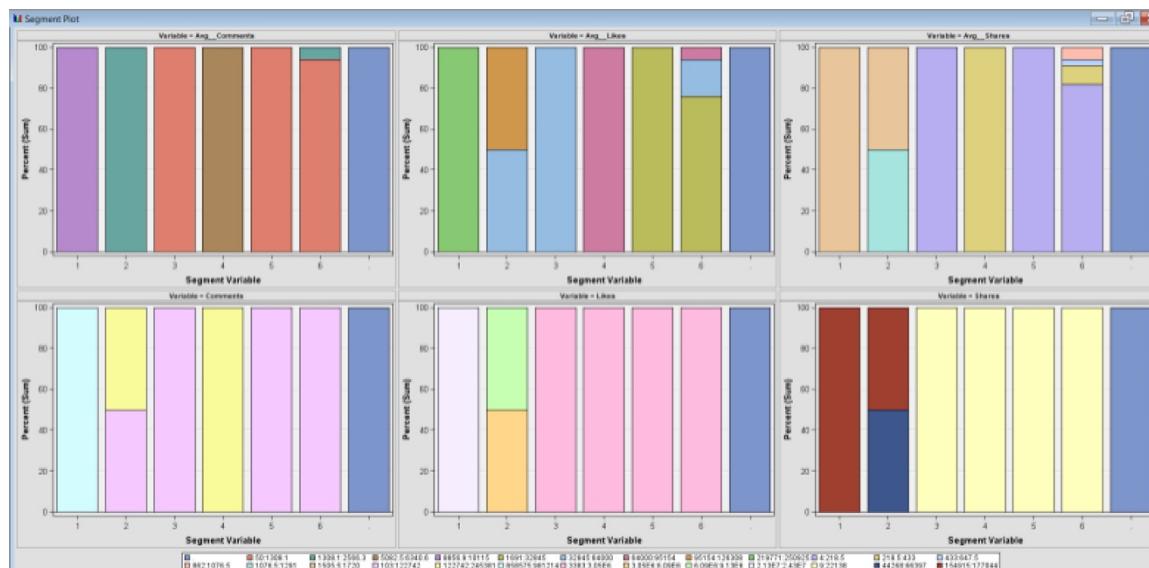
Though exploratory analysis is sufficient to find the best player with high promotional posts, model building is done for our future purposes to bring out intuitive results. Here, as a part of model building, K-means clustering, and K-Nearest Neighbor are performed for the required analysis using SAS enterprise miner. K-means clustering is performed to group the players with similar characteristics, and K-NN is attempted to check the categorization of the players.

Regression is not performed using the data since our goal is not to predict using the data.

## K-MEANS CLUSTERING:



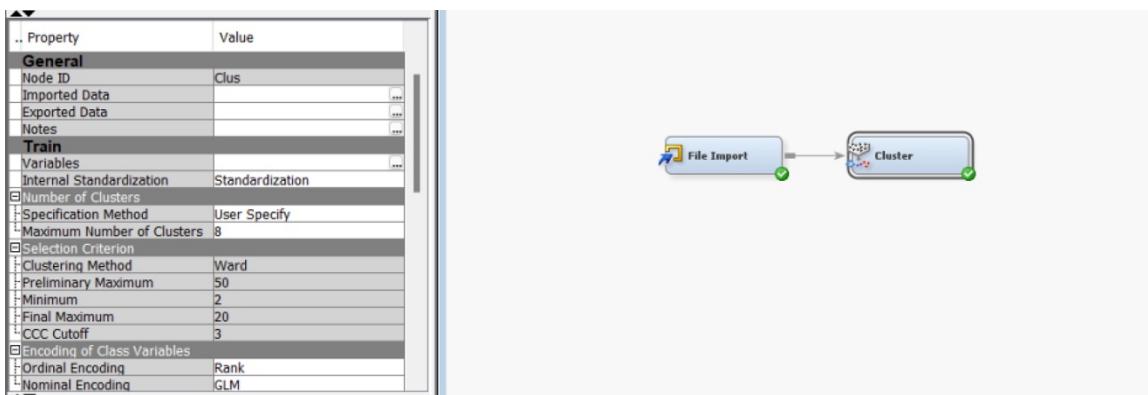
The Excel data file was used to perform clustering, where players are clustered into groups with similar characteristics. The default clustering resulted in six clusters at the beginning, with cluster 2 containing only one player, Virat Kohli, because of his remarkable reach of posts and following on the social media that made him stand out among all the other players and belong in a separate cluster.



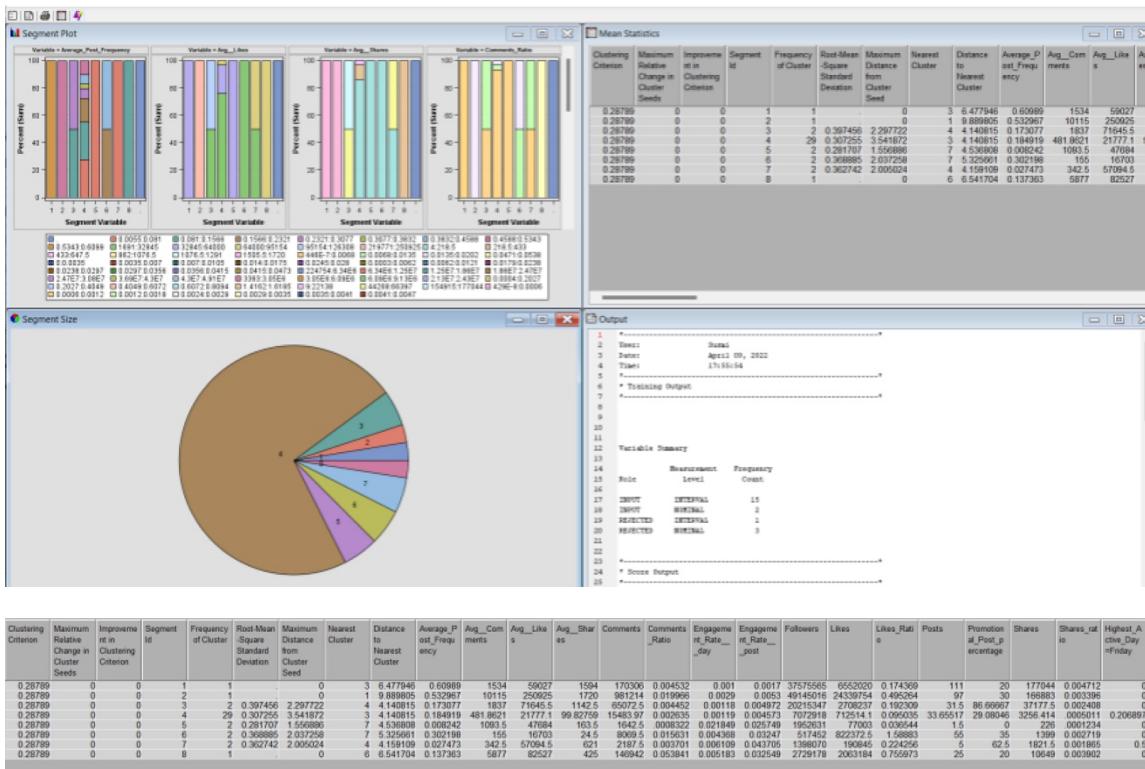
The above figure represents the segmentation of the cluster concerning the various variables. With this as a reference, new data could be categorized into any of these clusters with similar characteristics correspondingly. **The variable Importance shows the variables that are been considered to form the cluster breakdown.**

#### Variable Importance

Variable Name	Label	Number of	Number of	Importance
		Splitting Rules	Surrogate Rules	
Avg_Comments		1	0	1.00000
Comments	Comments	0	1	0.98742
Avg_Likes		0	1	0.97468
Likes	Likes	0	1	0.97468
Avg_Shares		0	1	0.97468
Shares	Shares	0	1	0.97468



In the case of User Specified Clustering, when the number of clusters has been set to **8 Clusters, most of the players converge towards one as like in the default clustering.** While the other players have some metrics that sets them apart from the group. This User specified clustering helped in analyzing the similarities among the players that lie in between the extremes.



## Variable Importance

Variable Name	Label	Number of Splitting	Number of Surrogate	Importance
		Rules	Rules	
Followers	Followers	0	2	1.00000
Engagement_Rate__post		1	0	0.77247
Engagement_Rate__day		0	1	0.73863
Likes_Ratio		0	1	0.72699
Comments_Ratio		0	1	0.72699
Shares_ratio		0	1	0.71517
Avg_Shares		1	0	0.69163
Shares	Shares	0	1	0.68293
Likes	Likes	0	1	0.67411
Avg_Likes		0	1	0.67411
Average_Post_Frequency		0	1	0.66518

This data using the user-defined clustering with 8 clusters is like the box plot distribution analyzed during the exploratory study earlier. This result is in accordance with the output of the exploratory analysis of the data. Here, the major clusters considered are:

- Cluster 1 – Virat Kohli
- Cluster 2 – Sachin Tendulkar
- Cluster 3 – AB de Villiers and Prithvi Shaw

Apart from considering only the basic common factors like, comments, and shares, the other factors like followers, engagement rate, and Post Frequency were also considered, as it was because of these variables, unique cluster characteristics were obtained, and the **clusters are found to be more realistic and intuitive than the default clustering.**

## K NEAREST NEIGHBOR:

Nearest neighbor	Nearest neighbor	Nearest neighbor	Nearest neighbor	Players	Followers	Posts	Likes	Shares	Comments	Node_ID	Profile_Link	Avg_Likes	Avg_Shares
22.0	14.0	6.0	1.0	Sachin Tendulkar	111	6552020	177044	170306	1.0004426919512E14	https://www.facebook.com/100...	59027.0	1594.0	
6.0	14.0	1.0	2.0	Virat Kohli	49145016	97	24339754	166083	981214	1.0004427304957E14	https://www.facebook.com/100...	250925.0	1720.0
39.0	12.0	4.0	3.0	KL RAHUL	10284388	54	3027838	11552	71043	1.00044235379214E14	https://www.facebook.com/100...	56071.0	213.0
3.0	39.0	9.0	4.0	RAVIN德拉 JADEJA	10844998	30	1671536	9211	35529	1.0004451469218E14	https://www.facebook.com/100...	55717.0	907.0
16.0	37.0	19.0	5.0	RISHABH PANT	5562289	42	1003578	3095	18924	1.00044527715691E14	https://www.facebook.com/100...	23094.0	73.0
14.0	28.0	6.0	22.0	ROHIT SHARMA	20931904	20	1332768	5011	40543	1.000442400200614	https://www.facebook.com/100...	66638.0	250.0
10.0	19.0	39.0	7.0	ISHAN KISHAN	3727674	21	632491	1752	15440	1.0004414880557E14	https://www.facebook.com/100...	30118.0	93.0
33.0	4.0	39.0	8.0	HARDIK PANDYA	11907722	23	1241491	4456	23454	1.00044386601022E14	https://www.facebook.com/100...	53977.0	193.0
11.0	34.0	10.0	9.0	RASHID KHAN	2729178	25	2063188	10649	146942	1.0004428146976E14	https://www.facebook.com/100...	82527.0	425.0
13.0	11.0	10.0	34.0	SANJU SAMSON	2905248	29	913791	1965	15953	1.0004428149679E14	https://www.facebook.com/100...	20051.0	53.0
34.0	10.0	13.0	11.0	SHREYAS IYER	26684932	27	572847	2345	14048	1.000442905946403E14	https://www.facebook.com/100...	21216.0	86.0
3.0	4.0	12.0	39.0	JASPRIT BUMRAH	9395089	37	1757148	20056	10004444176940E14	https://www.facebook.com/100...	47490.0	159.0	
17.0	27.0	11.0	13.0	MAYANK AGARWAL	2384770	67	469282	600	5037	1.0004443719233E14	https://www.facebook.com/100...	7004.0	8.0
20.0	22.0	6.0	14.0	MS DHONI	27214451	20	647915	20816	24248	1.0004481937943E14	https://www.facebook.com/100...	32395.0	1040.0
4.0	3.0	9.0	15.0	AB De	13216242	43	4766559	50539	105997	1.0004428940391E14	https://www.facebook.com/Ab...	110996.0	1245.0
19.0	37.0	5.0	16.0	Suresh Raina	6043500	86	234764	10220	32444	1.00044290346095E14	https://www.facebook.com/Su...	27298.0	118.0
20.0	23.0	17.0	27.0	PRITHVI SHAW	2242883	2	177277	2097	371	1.00044395660694E14	https://www.facebook.com/100...	89638.0	1049.0
10.0	7.0	38.0	18.0	RAVICHANDRAN ASWIN	4074691	70	206368	2441	4131	1.00044160351507E14	https://www.facebook.com/100...	2611.0	31.0
32.0	5.0	19.0	37.0	UMESH YADAV	6211138	8	39665	70	1273	1.00044509476307E14	https://www.facebook.com/100...	4958.0	8.0
23.0	17.0	27.0	20.0	KULDEEP YADAV	1947454	23	153799	345	9736	1.0004419737116E14	https://www.facebook.com/100...	6666.0	15.0
33.0	4.0	21.0	8.0	CHRIS GAYLE	13051349	10	93404	300	2544	1.00044317545336E14	https://www.facebook.com/100...	9340.0	38.0
21.0	28.0	6.0	22.0	YUVRAJ SINGH	19054947	5	81059	3124	1767	1.00044286896682E14	https://www.facebook.com/100...	16211.0	624.0
29.0	27.0	17.0	23.0	ANSHAR PATEL	2369426	35	102415	195	2213	1.00044647304968E14	https://www.facebook.com/100...	2926.0	4.0
26.0	35.0	31.0	24.0	RAHUL TEWATIA	1221964	16	129640	95	1195	1.000466263197E14	https://www.facebook.com/100...	8102.0	5.0
36.0	26.0	30.0	25.0	WASHINGTON SUNDAR	509723	72	794749	1920	9223	1.00051414294787E14	https://www.facebook.com/100...	11038.0	26.0
30.0	36.0	25.0	26.0	RAHUL TRIPATHI	553257	8	204413	1546	4004	1.000579727933E14	https://www.facebook.com/100...	25551.0	193.0
29.0	23.0	17.0	27.0	KRUNAL PANDYA	2262069	19	161828	428	2115	1.00044260937694E14	https://www.facebook.com/100...	8517.0	22.0
6.0	22.0	21.0	28.0	SHIKHAR DHawan	17140774	54	1802822	13609	25332	1.0004427427052E14	https://www.facebook.com/100...	23395.0	253.0
17.0	27.0	29.0	23.0	SAM CURRAN	2547182	2	117275	250	2195	1.0005700719768E14	https://www.facebook.com/100...	58637.0	125.0
36.0	26.0	30.0	25.0	RAHUL CHAHAR	525181	38	849996	678	6916	1.000185590509E14	https://www.facebook.com/100...	22368.0	23.0
20.0	24.0	31.0	35.0	YUVENDRA SINGH C.	1281491	2	3303	9	103	1.0004429481371E14	https://www.facebook.com/100...	1691.0	4.0
37.0	40.0	19.0	32.0	DAVID WARNER	7695287	2	4681	33	501	1.00044542354906E14	https://www.facebook.com/100...	2340.0	16.0
40.0	4.0	39.0	33.0	BHUVENESHWAR KUMAR	10376635	2	32961	103	2536	1.00044317545336E14	https://www.facebook.com/100...	16490.0	51.0
13.0	11.0	10.0	34.0	ROBIN UTHAPPA	28890108	86	997335	1476	5711	1.0004439533919E14	https://www.facebook.com/100...	11596.0	17.0
20.0	24.0	31.0	35.0	MOAHMMED SIRAJ	1358079	1	36731	202	1990	1.00044794459476E14	https://www.facebook.com/100...	36731.0	202.0
30.0	26.0	25.0	36.0	SHREYAS GOPAL	224754	21	148008	157	1069	1.000444779981E14	https://www.facebook.com/100...	7051.0	7.0
18.0	5.0	19.0	37.0	DAVID MILLER	5493502	6	369803	317	2389	1.000443763944E14	https://www.facebook.com/100...	61633.0	52.0
10.0	18.0	38.0	7.0	LASTH MALINGA	3592181	77	618985	11017	55279	1.0004392599321E14	https://www.facebook.com/100...	9037.0	143.0
12.0	4.0	39.0	33.0	HARBHAJAN SINGH	10324052	92	835480	4876	35922	1.000441337408E14	https://www.facebook.com/100...	9081.0	53.0
33.0	32.0	12.0	40.0	AJINKYA RAHANE	8476878	4	13591	43	378	1.0004428925775E14	https://www.facebook.com/100...	3997.0	10.0
39.0	32.0	33.0	40.0										

The K Nearest Neighbor is used to determine the players with similar attributes in the dataset. This algorithm groups the data points based on the closeness to the other data points. 4 nearest neighbors are used to see which players go along together in terms of the variables with respect to the promotions. **The value of K is chosen as 4 since it has the lowest misclassification error rate.** In 4NN, as one of the neighbors is itself on calculating

the shortest distance, the other three points are considered the actual neighbors. From the results, it could be inferred that Sachin, Dhoni, Yuvaraj, and Rohit are neighbors with similar traits. **When new data comes in, we can match them with their neighbors and find out their characteristics easily.**

## INTERPRETATION

From both exploratory analysis and modelling, it could be assured that Virat Kohli is the best player for any kind of Facebook Promotion. But Virat Kohli, Sachin Tendulkar, and MS Dhoni are the top cricketers and are considered highly influential among the public on the social media platforms with high number of followers and their charges for a promotional post range between **1 – 5 crore rupees**.

On the other hand, it was also found that **AB De Villiers and Prithvi Shaw** have great reach for their posts, be it promotional or normal, with a good post engagement rate. For a promotional post, **AB De Villiers charges nominally around 30 - 50 Lakhs with 13.2 million engagements**, whereas Prithvi Shaw's post reach is good, but he can reach only 2.24M.

Therefore, Prithvi Shaw would be a good fit for smaller promotions and AB De Villiers would be considered a good cricket player for making effective promotions. These players do not have as many followers as Virat Kohli or MS Dhoni, but their posts have met the desirable reach with the public and attained a good engagement rate.

## RESULTS & CONCLUSION:

With all the analysis and clustering techniques, we can conclude that choosing **AB de Villiers as the brand ambassador would be ideal, as it would be economical cutting costs for the company and promotes the product to a larger audience**. In addition, the post must be published in the analyzed peak timings and days drawn from the exploratory analysis to meet the expected reach among the public.

This would help the brand spend a decent amount of charge to the cricket player for the promotional pots and get the desired reach and exposure to the product as well as the brand among the public.