

Capstone Project



Team 5 : Company Classification



Team Members

Arunav Goswami

Nayanjyoti Sharma

Mohammed Saad Pasha

Content

- ❖ Problem statement
- ❖ Data summary
- ❖ Covering insights
- ❖ Challenges
- ❖ Conclusions

Problem Statement

We are given with web scraped data of various businesses and companies.

We need to somehow categorize these businesses and companies across a standard taxonomy (consists of term names and labels that are specific to an organization's information and unique to how that business operates). So that, business can leverage this information and target potential companies.

Data Summary

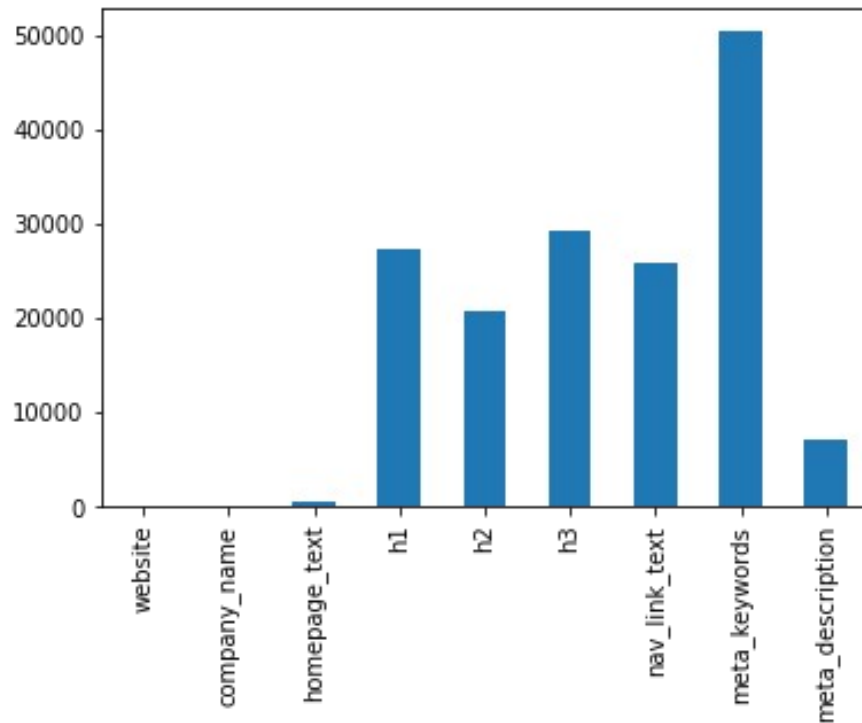
- We have records of 73974 companies with 9 columns/features.
- Each feature is a text data.
- We will be left with only 4K records, if we remove all the null values at the first place.

Data Overview

- **Website:** The website of the company/business
- **Company Name:** The company/business name
- **Homepage Text :** Visible homepage text
- **H1:** The heading 1 tags from the html of the home page
- **H2:** The heading 2 tags from the html of the home page
- **H3:** The heading 3 tags from the html of the home page
- **Navlink text:** The visible titles of navigation links on the homepage (Ex: Home, Services, Product, About Us, Contact Us)
- **Meta keywords:** The meta keywords in the header of the page html for SEO
- **Meta description:** The meta description in the header of the page html for SEO

Techniques Implemented

In the given data there were more null values equally spread across all the major columns. To avoid data loss we have combined all the columns as we are dealing with text. We also analysed individual columns to gain insights.



Data Preprocessing

- ❖ **Handled null values**
- ❖ **Removed punctuations**
- ❖ **Removed digits**
- ❖ **Removed stop words**
- ❖ **Removed special characters**
- ❖ **Stemming/Lemmatization**
- ❖ **Removed short words**



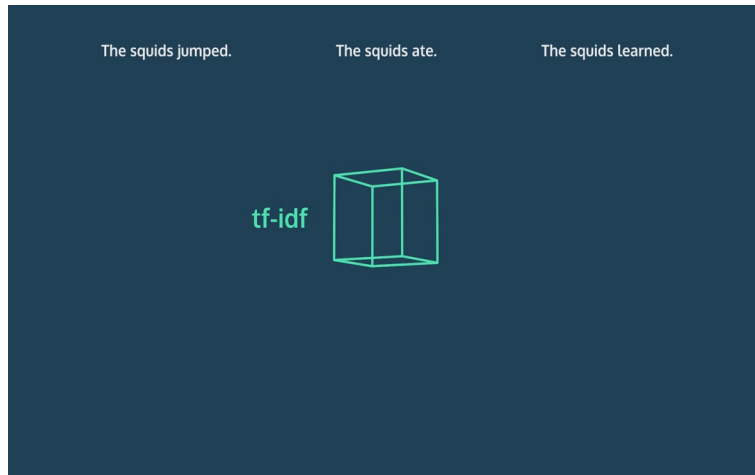
Word Embedding

Word embedding is a term used for the representation of words for text analysis

Few popular techniques are:

- Countvectorizer
- TF-IDF Vectorizer
- Word2Vec

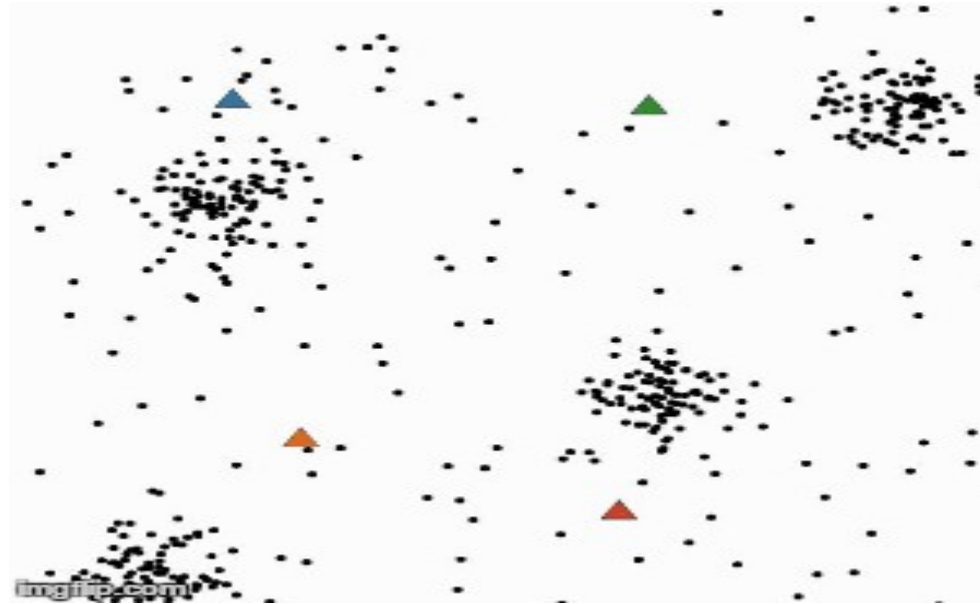
We experimented with different techniques and out of which TF-IDF gave us better results.



Using K Means

- **Kmeans is an unsupervised machine learning algorithm used for classification.**
- **To solve this problem statement, we went ahead with K Means algorithm to categorize the companies into multiple clusters.**

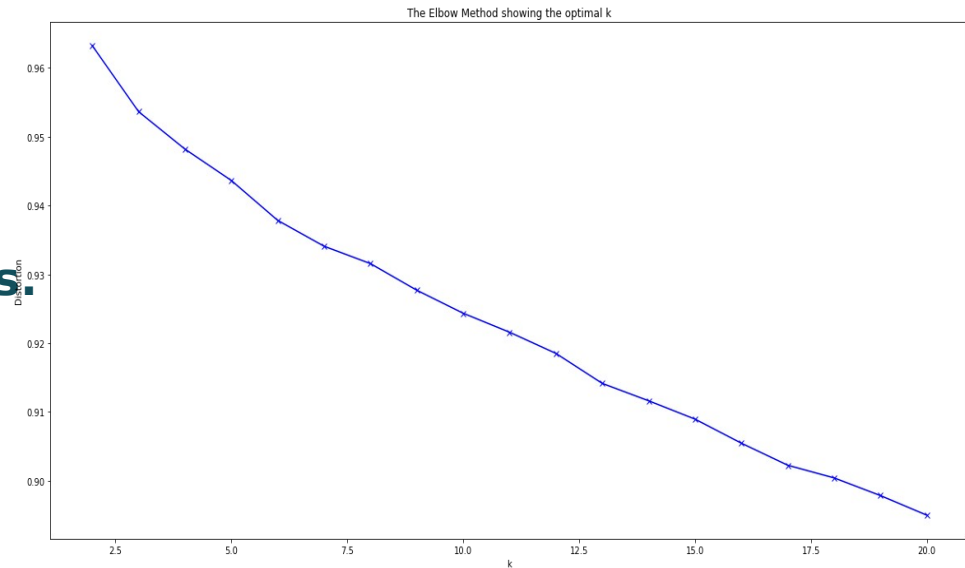
Clusters Insights



INSIGHTS FROM OVERALL PERSPECTIVE..

Optimal number of clusters..

By using Elbow method and Silhouette score, we came up with 12 as an optimal of clusters.



Elbow method

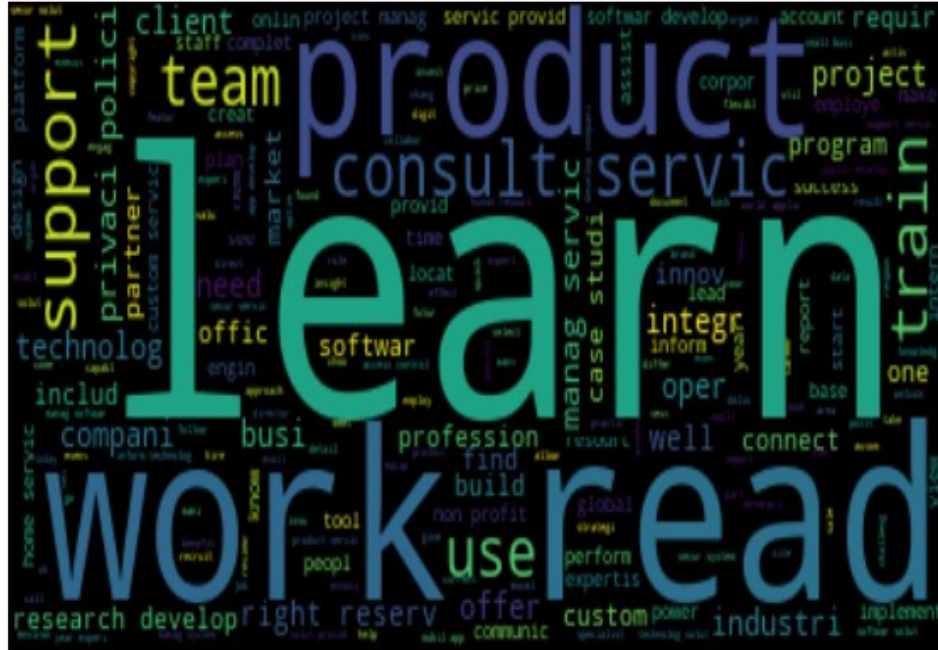
Cluster 1 : PET HEALTHCARE



Example sentences

Company Name	Text
Pleasant Valley Animal Hospital	adopting a pet bereavement loss of a pet..cancer chemotherapy cardiology cat only..
Annabessacook Veterinary Clinic	home meet our team large Animal Services Small Animal Services KDM Farm Resources & Externships..
Columbia Pike Animal Hospital	Request an Appointment Prevent Animal Cruelty Shop Online Access Pet Records..
Doylestown Animal Medical Clinic	proud to provide a full spectrum of routine and specialized medical services designed to care for pets..
Borchard Veterinary Clinic	PET WELLNESS & PREVENTATIVE CARE PET SURGERIES..

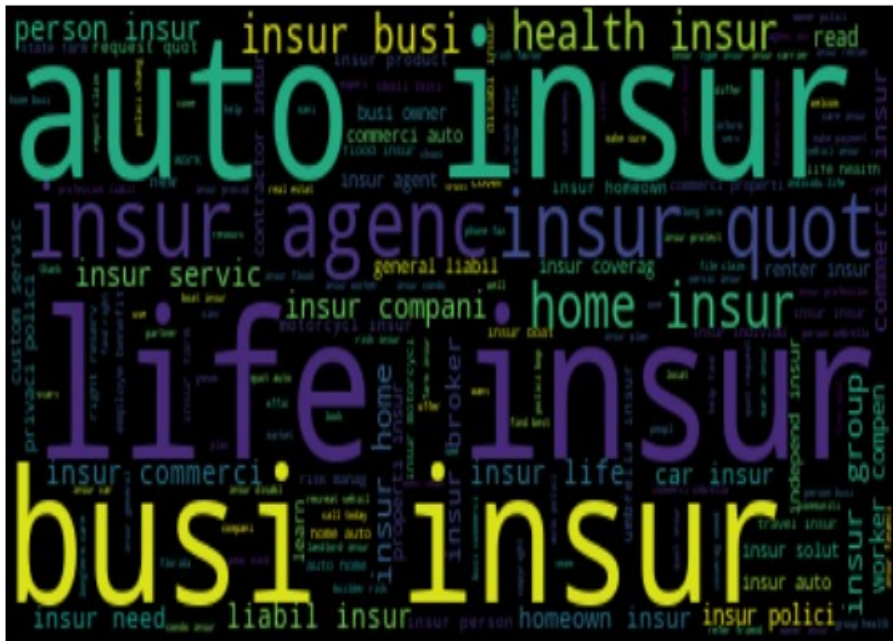
Cluster 2 :IT SERVICES



Example sentences

Company Name	Text
Center for Computer resources	Services & Solutions Microsoft Azure Microsoft Teams Modern Desktop Microsoft 365 Business..IT Consulting Webroot showcase Virtualization Projects..
Appmetric software limited	We believe that software consulting is more of a tactic than a company..reckoning entity in the competitive field of software services and business process outsourcing..
Fortune business transfers & acquisitions	Services Business Acquisition Consultants Preparing to Sell Business Help Selling Business..
Construe solutions	Wireless Solutions Services & Consulting Business consultant API Integrations Mobile Application..

Cluster 3 : INSURANCE



Example sentences

Company Name	Text
Stratford insurance group	We provide personalized support and full-service insurance to individuals and businesses to ensure you are protected in all your adventures..
Petley-hare limited insurance brokers	You're Unique. Your Insurance Should Be Too!
Private client insurance group	We are a global specialty insurance distributor with expertise in property, casualty, professional lines and specialty group benefit products..
W.C. Burgess insurance services ltd	For all your insurance needs, whether it is for your car, home, boat, snowmobile, ATV, business, life or disability insurance we have the answers you need and the prices you want..

Cluster 4 :DIGITAL MARKETING

Example sentences



Company Name	Text
The Digital Marketing People Inc.	Marketing is much more than just a great website..Work with a Digital Marketing Expert to use the right mix of digital marketing services..
ITI marketing	We provide analytical reports that permit clients to better tailor their off-line corporate advertising and marketing strategies to those demographics most responsive to their products..
Allegra Pittsburgh Design & Marketing	We're a full-service marketing and print communications company based in Pittsburgh, PA..
Renegade Advertising Inc	Renegade Advertising is a full service internet advertising agency..Renegade Advertising focuses on online and new media marketing and advertising..

Cluster 5 : LOGISTICS



Example sentences

Company Name	Text
Alpha Global Transport Limited	runs a tight ship with Transport Captain Keith Bonner..experience in both the air freight and general haulage industry..
Clark logistic services	Our flexible warehousing solutions and distribution services are focused on delivering a complete turn-key solution for your business-to-business order fulfillment needs..
Auction Transport services inc.	Leaders in Vehicle Transportation..
Eal & aw Transport group	our fleet consists of over 30 specially equipped trucks and specially equipped trailers to cater for the delivery of windows and doors in both the domestic and commercial market..

Company Name	Text
Medwest Eyecare	Caring, Friendly, Professional Eye Care..Fitting of Glasses Adjustments of Glasses Without Appointment..
Uhhs bedford medical ctr	full service, acute-care community hospital, offering adult and senior emergency services, a state-of-the-art outpatient surgery center, comprehensive imaging facilities..
Wilson pharmacy	Your Health Is Our Business..We also do specialized travel immunizations..
Alexo therapeutics	Oncology is a clinical stage biotechnology company developing innovative immuno-oncology therapies for cancer..

Cluster 7 : FINANCIALS



Example sentences

Company Name	Text
Weitzel Financial Services Inc	(WFS) has been to help our clients reach their financial goals by providing the support and resources needed to make sound, individualized financial decisions..
Greenlink Financial	Leader in consumer unsecured personal loans, helping thousands of people resolve their consumer debt with a personal loan..We've helped thousands of people resolve their debt and find financial freedom. Start Living Your Life Today!
Northwest financial & tax solutions, inc.	we offer a complimentary second opinion and comprehensive portfolio analysis..
Trinity tax service	We offer a broad range of services for business owners, executives and independent professionals..Our services encompass nearly every aspect of financial life.

Cluster 8 : REAL ESTATE & LEGAL AFFAIRS

Example sentences



Company Name	Text
Morris Law group	Morris Law Group is a boutique estate planning and asset protection law.. Through effective planning, our experienced attorneys are dedicated to helping individuals and families preserve and protect their wealth for future generations..
Premier (uk) property management ltd.	portfolio of listed and off-market properties ranging from exclusive domestic and commercial premises to hotels, land and development projects..
White cleland lawyers & consultants	Our Legal team is already highly skilled, but are each equally committed to ongoing practical learning..When a relationship breaks down you need support and assistance in working through the myriad of potential legal issues
Blackstone williams properties, llc	Blackstone Williams provides first-class rentals, sales & property management in Boston's Back Bay & beyond!



Example sentences

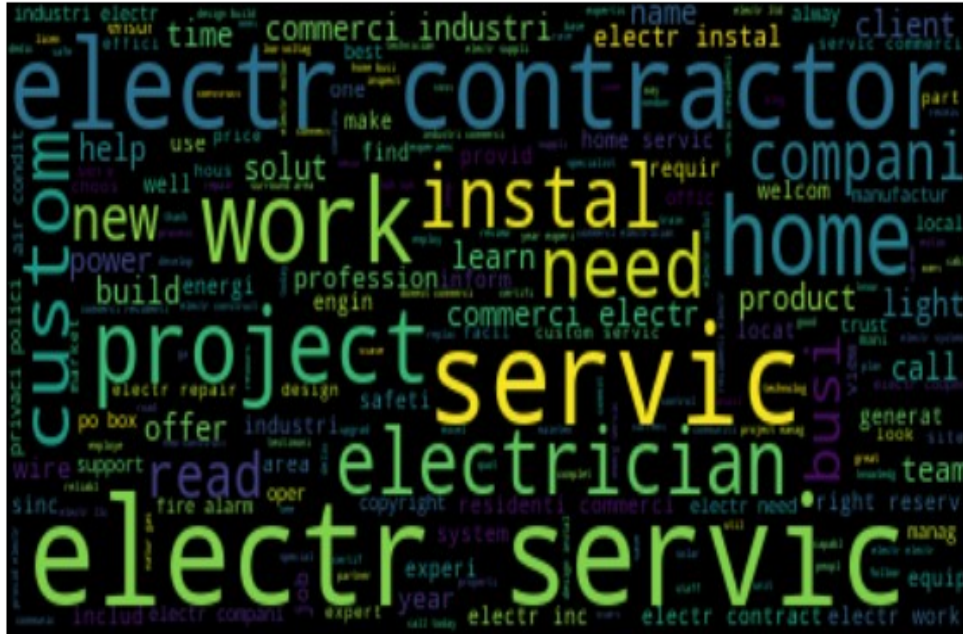
Company Name	Text
hawaiian hotels & resorts	breathtaking beaches and incredible natural wonders. Explore Now the best of Hawaii Escape to The Big Island for an unforgettable vacation or romantic honeymoon in the heart of
refinery hotel	imaginative reinterpretation of the Colony Arcade, Āia former hat factory and Prohibition-era tea room. More than 100 years of Garment District history are on display at our luxury hotel near
the wellington apartment hotel,	Accommodation in Kangaroo Point near the Gabba, East Brisbane. -†Perfect Self Contained Apartment Accommodation for holiday and corporate travellers, groups and conferences near Brisbane CBD. Come and experience our Luxury Central Brisbane
lower plenty hotel,	Award-winning and welcoming Friends, family or work colleagues. No matter who you, Āre dining with, our Bistro is the perfect place to entertain your tastebuds. Named the Best



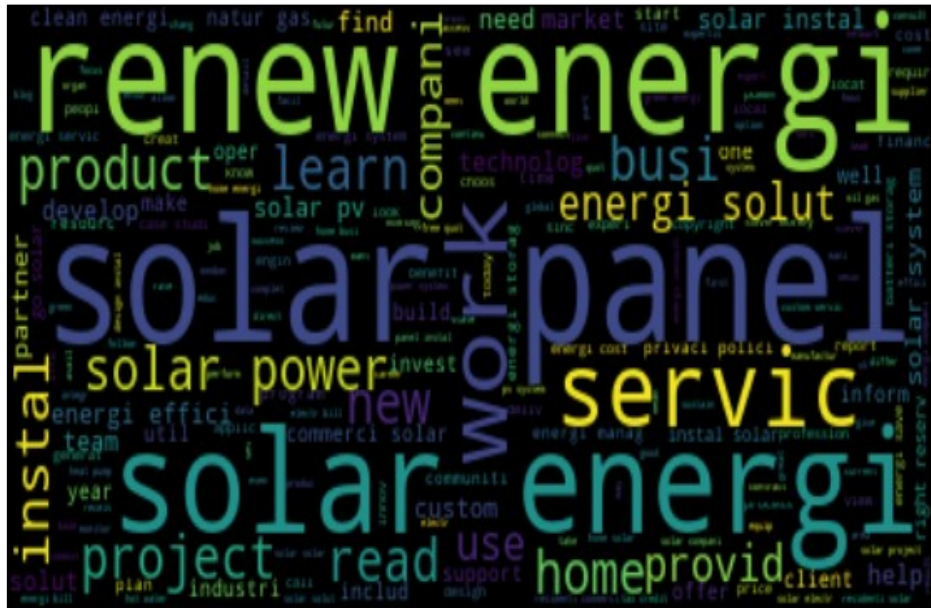
Company Name	Text
PrestaShop	freemium, open source e-commerce platform..You can launch your online store right now and start to sell online..
WooCommerce Pvt Ltd	e-commerce plugin for WordPress. Designed for small to large-sized online merchants using WordPress..
OXID eShop	ecommerce solution built using object oriented programming and PHP..
Currys	British electrical retailer operating in the United Kingdom..Free Delivery or Order & Collect..

Cluster 11 : HOME SERVICES

Example sentences



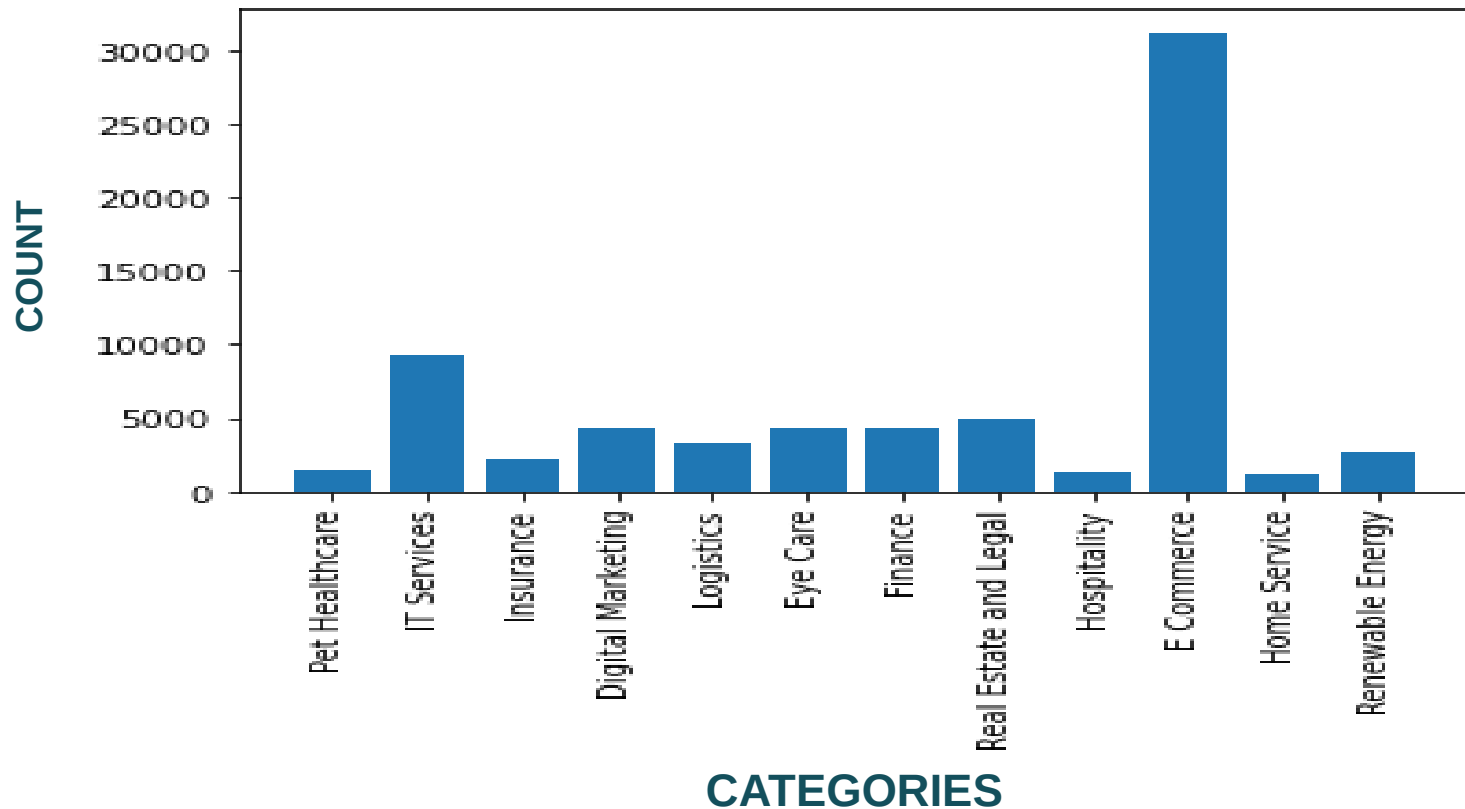
Company Name	Text
auction transport services inc.	HOME ABOUT SERVICES CONTACT EVENTS JOBS FranVbais AUCTION TRANSPORT SERVICES,customer service driven organization, where safety and consistently clean deliveries-are always a priority,
custom electric sd	Custom Electric SD. We handle all of your electrical contracting needs. We specialize in recessed lighting, and both residential and commercial wiring
d-tech electrical contractors ltd	by our NICEIC Approved Contractor, Construction Line and CHAS accreditations. All of our directly employed electricians are full qualified and highly experienced, which enable us to maintain our
electrical installations inc,	is what makes Electrical Installations shine. Our expertise as an electrical contractor adds an additional dimension of experience and capability that other system integrators cannot offer. Our



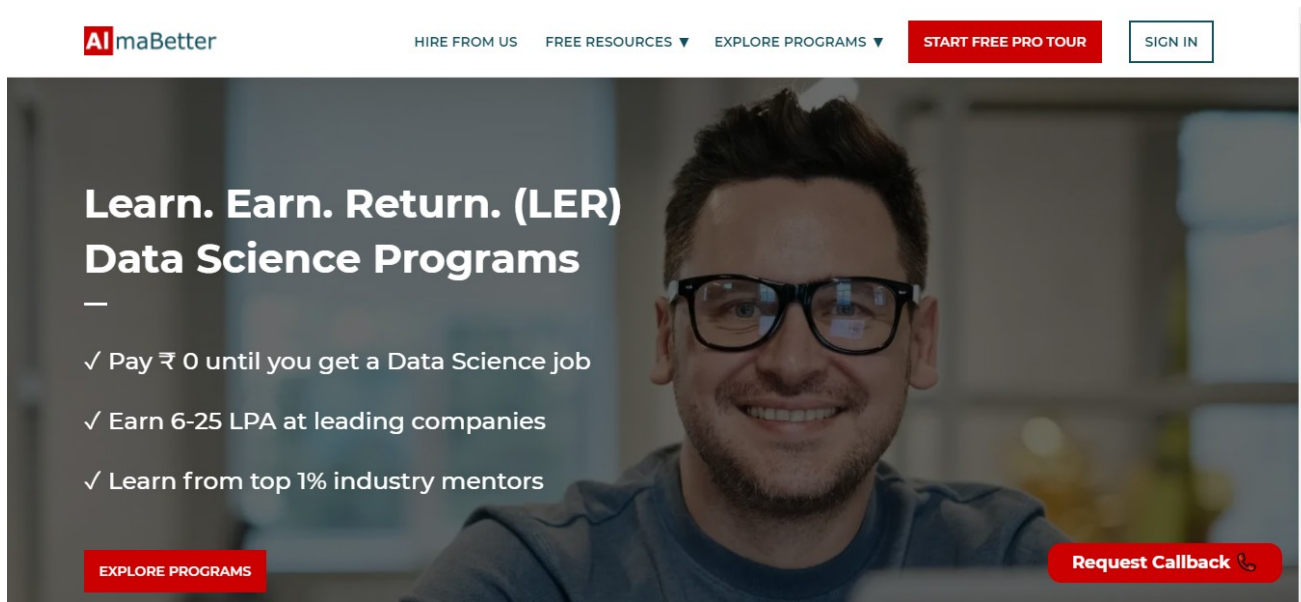
Example sentences

Company Name	Text
Commercial & Industrial Microturbine,	Industrial-Scale Alternative Energy Developer & Financier, reduce their energy supply costs, and—become more sustainable and resilient
solarpowergb.com	Welcome to Solar Power GB, the premier provider of solar PV, solar thermal, solar panel cleaning and energy saving products in Essex. ",Renewable Energy
appleblossomenergy.com	Insulation Duct Insulation Reflective Barriers Wall Insulation Floor Insulation Insulation for Builders Cellulose Insulation
sd-windenergy.com	No Shut Down Wind Turbines#sep#Share this page#sep#News & Updates#sep#A Unique Delta Rotor Design",Renewable Energy

Distribution of Categories



INSIGHTS FROM HOMEPAGE..



AI maBetter

[HIRE FROM US](#) [FREE RESOURCES ▼](#) [EXPLORE PROGRAMS ▼](#) [START FREE PRO TOUR](#) [SIGN IN](#)

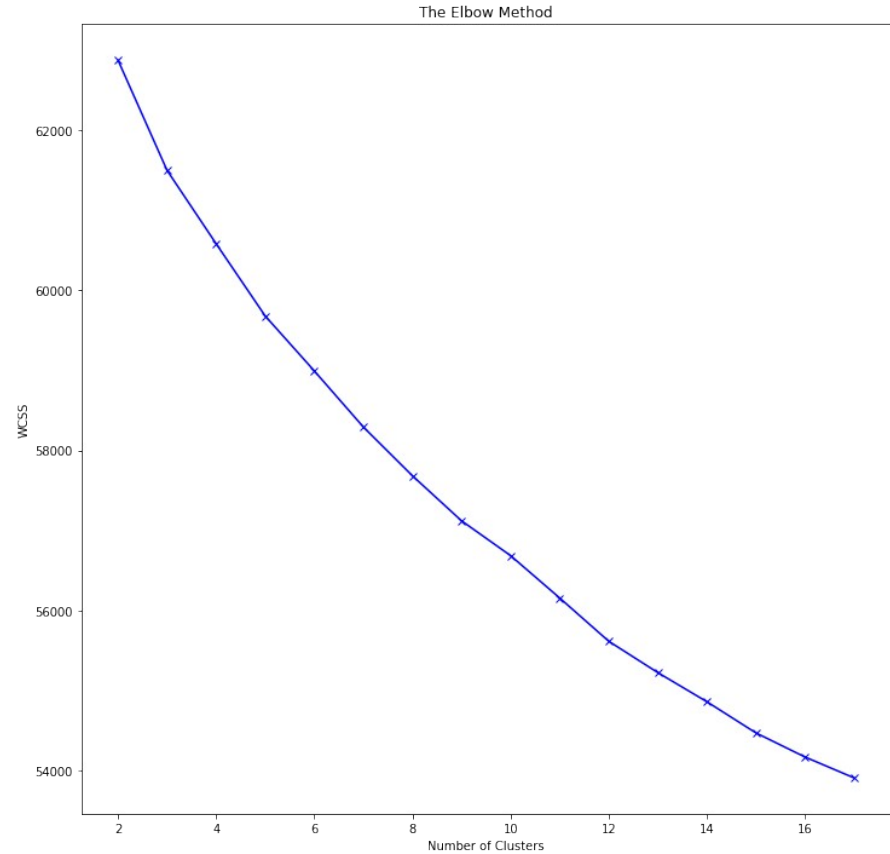
Learn. Earn. Return. (LER) Data Science Programs

- ✓ Pay ₹ 0 until you get a Data Science job
- ✓ Earn 6-25 LPA at leading companies
- ✓ Learn from top 1% industry mentors

[EXPLORE PROGRAMS](#) [Request Callback ☎](#)

Optimal Clusters

By using Elbow method and Silhouette score, we came up with 9 as an optimal number of clusters.



Cluster 1 : ENERGY & UTILITIES

Cluster 0

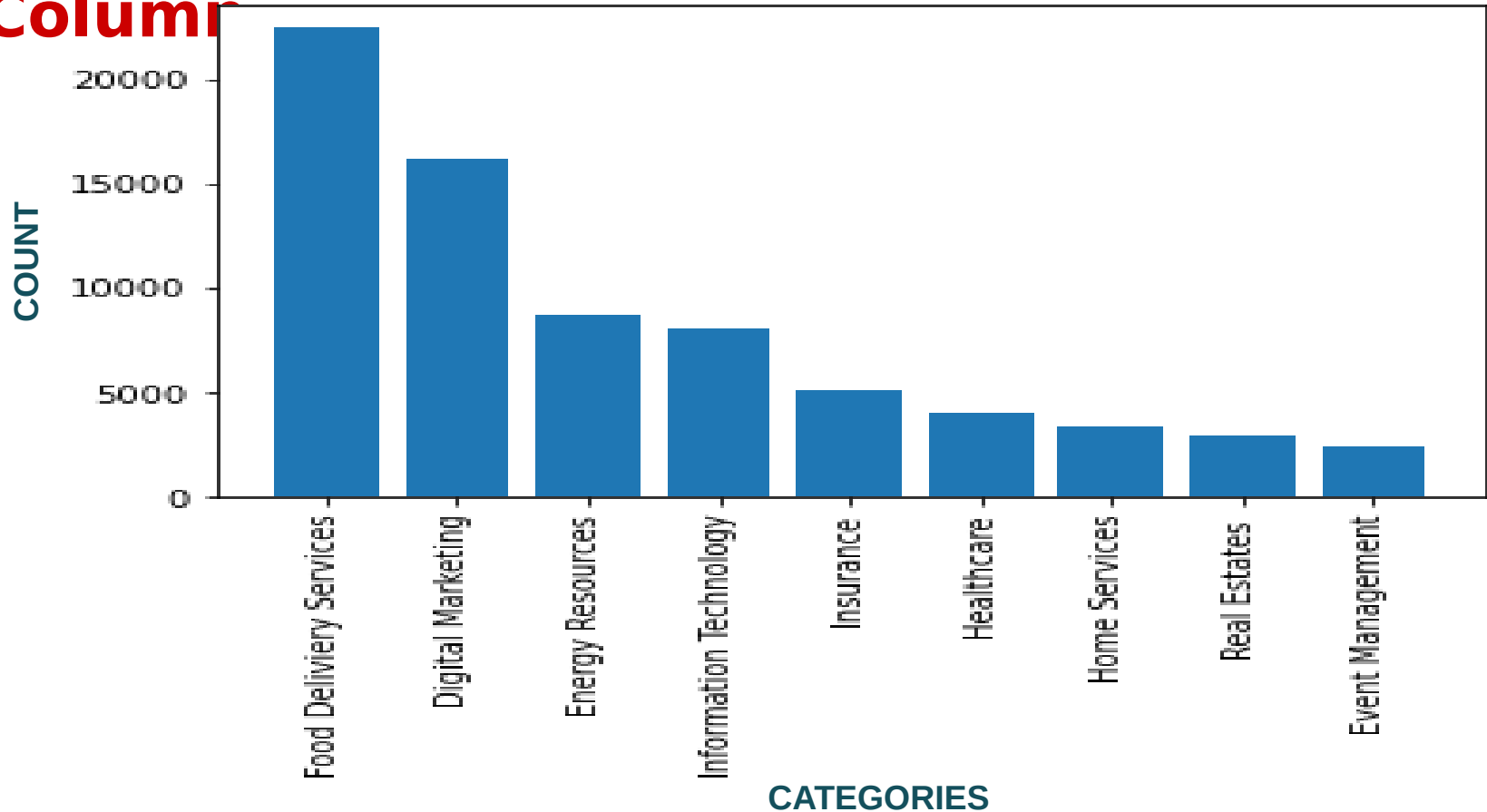


Cluster 9 : HOME SERVICES

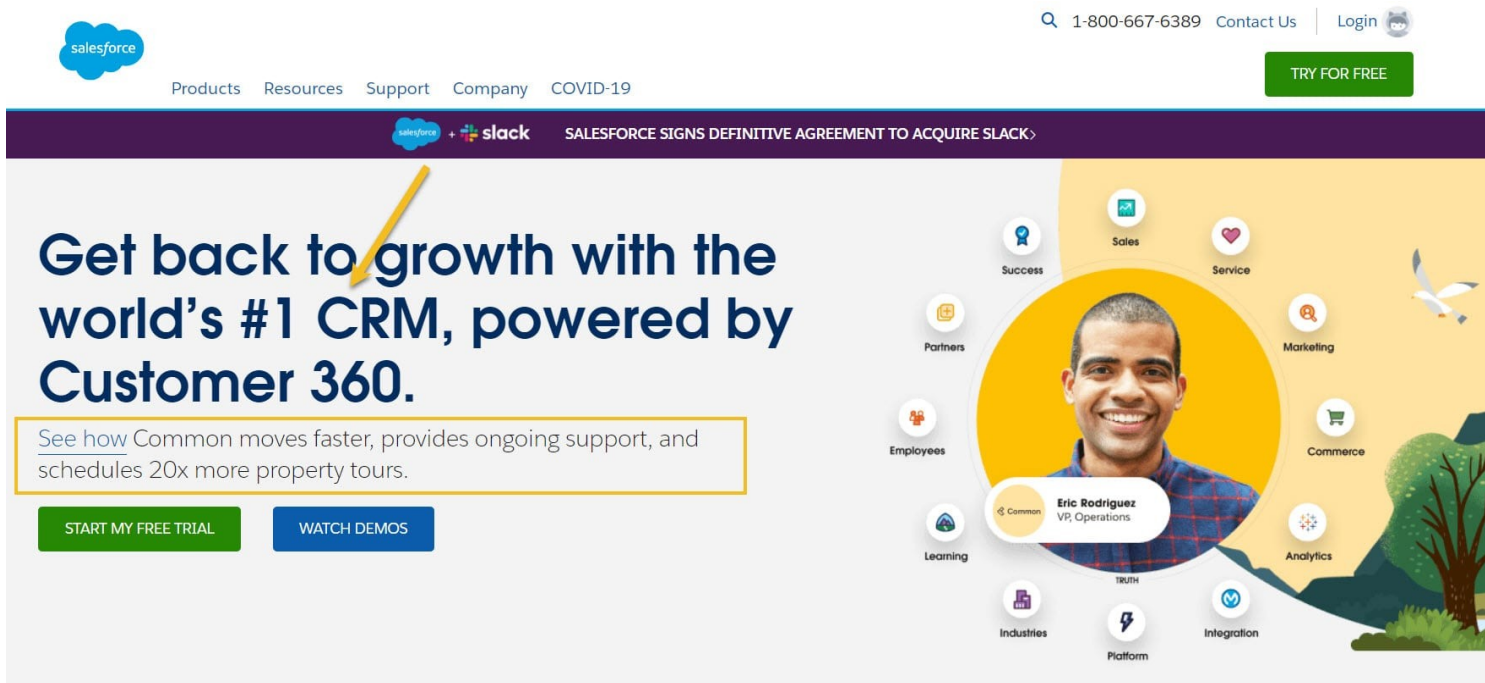
Cluster 8



Distribution based on 'Homepage_txt' Column



INSIGHTS FROM HEADERS..



The image shows the top portion of the Salesforce website. At the top right is a red 'AI' logo. Below it is the Salesforce logo and a navigation bar with links for Products, Resources, Support, Company, and COVID-19. To the right of these links are search and contact options, including a phone number, 'Contact Us', 'Login', and a 'TRY FOR FREE' button. A purple banner below the navigation bar features the Salesforce and Slack logos, followed by the headline 'SALESFORCE SIGNS DEFINITIVE AGREEMENT TO ACQUIRE SLACK'. The main hero section has a large headline 'Get back to growth with the world's #1 CRM, powered by Customer 360.' with a yellow pencil icon pointing to the word 'growth'. Below the headline is a text box stating 'See how Common moves faster, provides ongoing support, and schedules 20x more property tours.' with two buttons: 'START MY FREE TRIAL' and 'WATCH DEMOS'. On the right side of the hero section is a circular graphic with a portrait of Eric Rodriguez, VP of Operations at Common, surrounded by various business function icons like Sales, Service, Marketing, Commerce, Analytics, Integration, Platform, Industries, Learning, Employees, Partners, and Success.

salesforce

Products Resources Support Company COVID-19

1-800-667-6389 Contact Us Login

TRY FOR FREE

salesforce + slack SALESFORCE SIGNS DEFINITIVE AGREEMENT TO ACQUIRE SLACK

Get back to growth with the world's #1 CRM, powered by Customer 360.

See how Common moves faster, provides ongoing support, and schedules 20x more property tours.

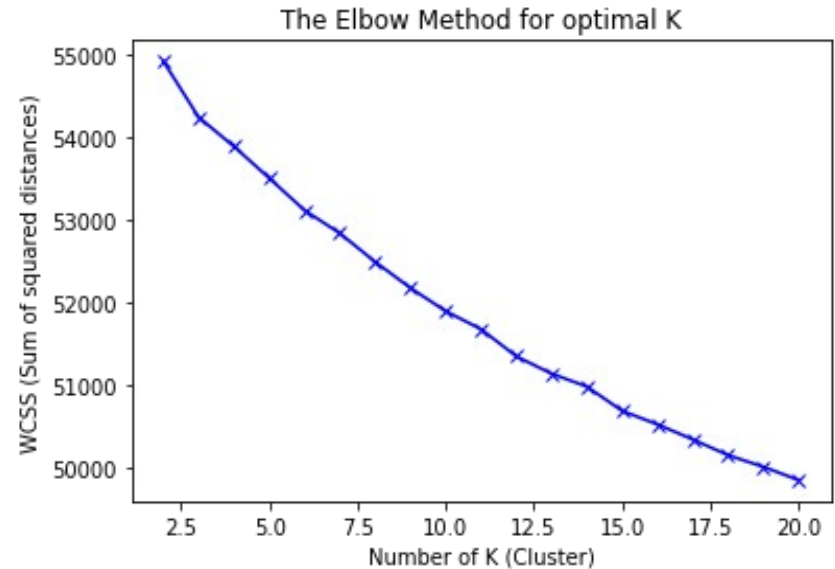
START MY FREE TRIAL WATCH DEMOS

Success Sales Service Marketing Commerce Analytics Integration Platform Industries Learning Employees Partners

Eric Rodriguez
VP, Operations

Optimal Clusters

By using Elbow method and Silhouette score, we came up with 12 as an optimal number of clusters.



Cluster 1 : REAL ESTATE



Cluster 2 : LEGAL AFFAIRS



Cluster 3 : FOOD SERVICES



Cluster 4 : IT SERVICES



Cluster 5 : LOGISTICS



Cluster 6 : INSURANCES



Cluster 7 : HEALTHCARE



Cluster 8 : INDUSTRIALS



Cluster 9 : EATERY INDUSTRY





Cluster 11 : HOSPITALITY



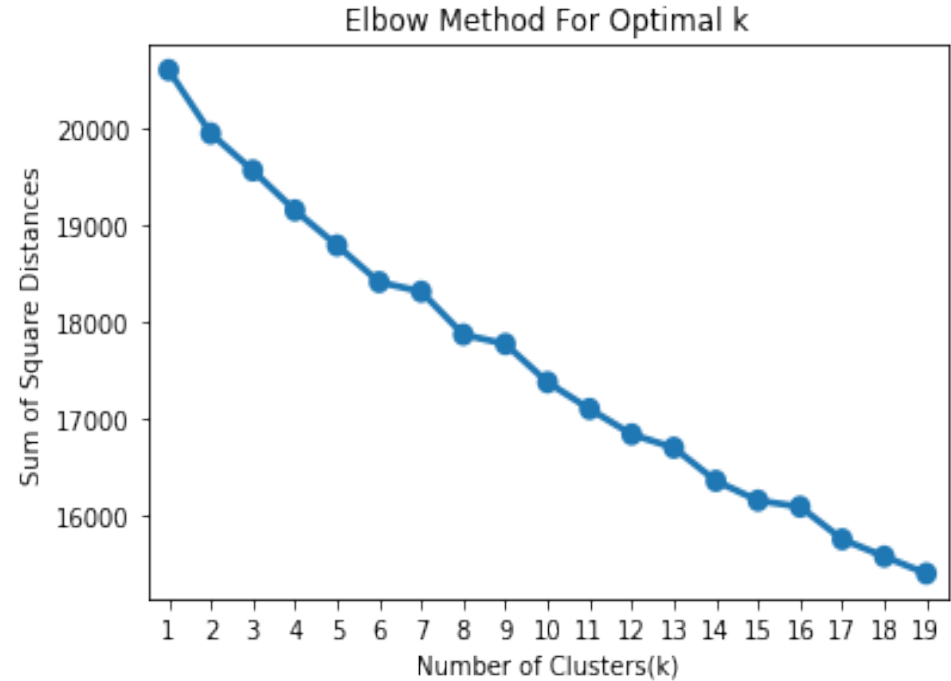
Cluster 12 : ENERGY & UTILITIES



INSIGHTS FROM Meta_keywords & Meta_description..

Optimal Clusters

By using Elbow method and Silhouette score, we came up with 9 as an optimal number of clusters



Cluster 4 : FINANCIALS

Cluster 3



Cluster 7 : HOME SERVICES

Cluster 6



Cluster 8 : INDUSTRIALS

Cluster 7



Challenges

- **Text Preprocessing as there were characters from multiple scripts.**
- **There were lot of null values that we had to handle carefully.**
- **Limited resources : Dataset turned to be large for computation power we have.**
- **Deciding the number of optimal clusters was bit difficult since the dataset is really vast and can have more number of clusters.**

Conclusion

- **As we observed “E-commerce” was the top cluster from overall perspective, we can further dig down and do sub-clustering to gain more insights.**
- **We were able to categorize businesses/companies in 9 to 12 clusters based on different features and on the combination of whole features as well.**
- **Those clusters represent various industries and we can approach those companies of each cluster with relevant business proposals.**

Q & A