# INTRODUCTION TO DATA SCIENCE
## SESSION # 6: DATA SCIENCE TEAMS

SANKARA NAYAKI K
sankaranayaki@wilp.bits-pilani.ac.in

**BITS** Pilani
Pilani | Dubai | Goa | Hyderabad

The instructor is gratefully acknowledging
the authors who made their course
materials freely available online.

References:

- Introducing Data Science by Cielen, Meysman and Ali

- Storytelling with Data by Cole Nussbaumer Knaflic; Wiley

- Introduction to Data Mining by Tan, Steinbach and Vipin Kumar

- The Art of Data Science by Roger D Peng and Elizabeth Matsui

- Python Data Science Handbook: Essential tools for working with data by Jake VanderPlas

# Table of Contents

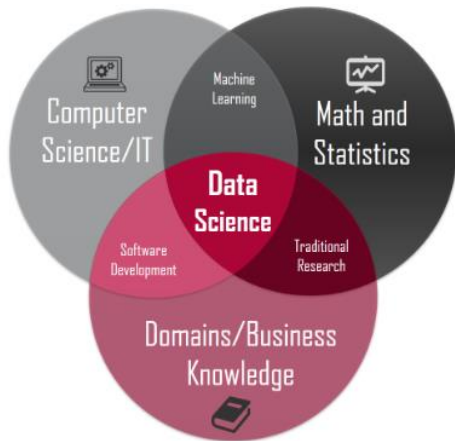1 [Course Handout](#)

2 [Data Science Teams](#)

# Course Handout

M1 Introduction to Data Science

M2 Data Analytics

M3 Data Science Process

**M4 Data Science Teams**

M5 Data and Data Models

M6 Data wrangling and Feature Engineering

M7 Data visualization

M8 Storytelling with Data

M9 Ethics for Data Science

# Table of Contents

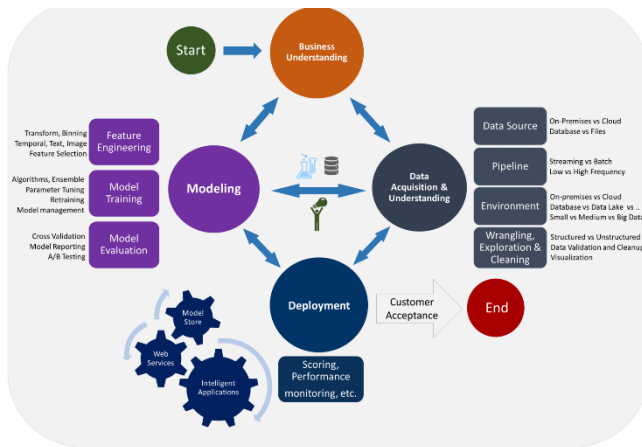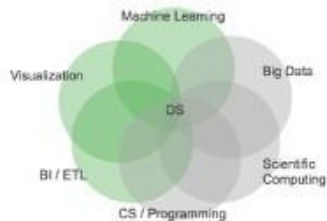1 [Course Handout](#)

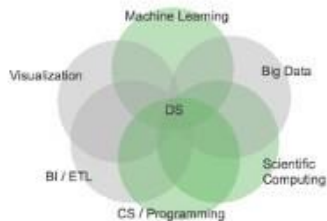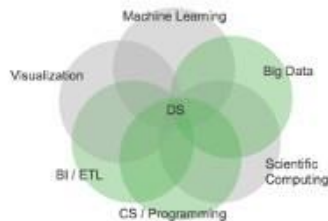2 [Data Science Teams](#)

# Introduction

# DATA SCIENCE LIFECYCLE
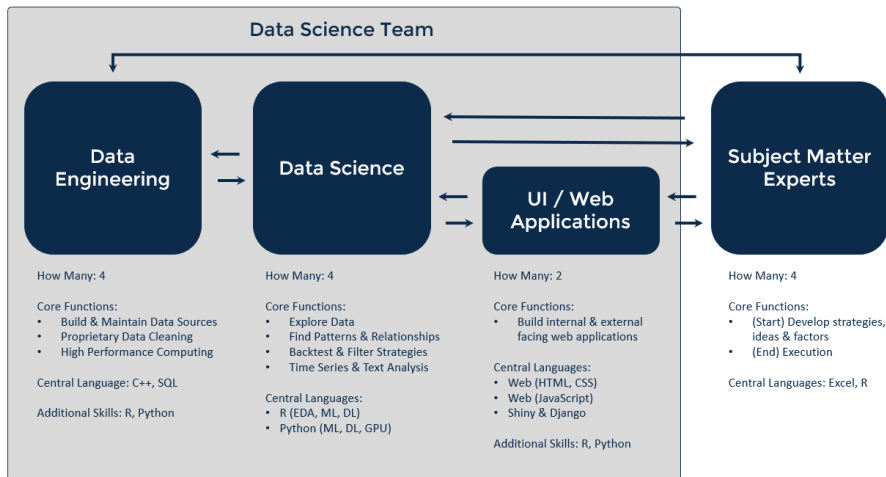
# DATA TEAM

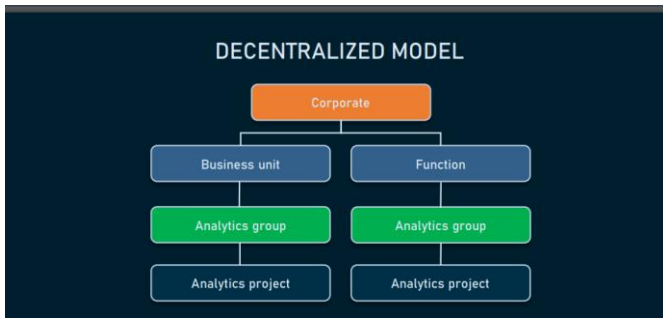| Statistician / Analyst | Research / Computational Scientist | Developer / Engineer |

# DATA SCIENCE TEAM



## Data Science Team

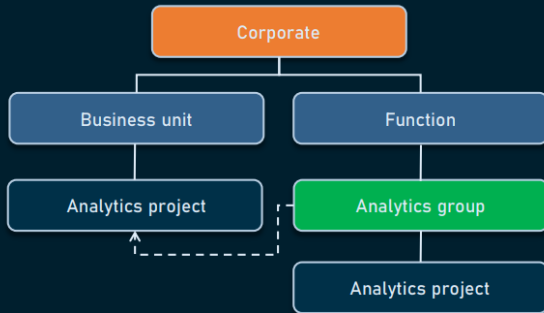**Data Engineering**

How Many: 4

Core Functions:
- Build & Maintain Data Sources
- Proprietary Data Cleaning
- High Performance Computing

Central Language: C++, SQL

Additional Skills: R, Python

**Data Science**

How Many: 4

Core Functions:
- Explore Data
- Find Patterns & Relationships
- Backtest & Filter Strategies
- Time Series & Text Analysis

Central Languages:
- R (EDA, ML, DL)
- Python (ML, DL, GPU)

**UI / Web Applications**

How Many: 2

Core Functions:
- Build internal & external facing web applications

Central Languages:
- Web (HTML, CSS)
- Web (JavaScript)
- Shiny & Django

Additional Skills: R, Python

**Subject Matter Experts**

How Many: 4

Core Functions:
- (Start) Develop strategies, ideas & factors
- (End) Execution

Central Languages: Excel, R

# DATA SCIENCE – SKILL SET

## NECESSARY AND PREFERRED DATA SCIENCE SKILLS
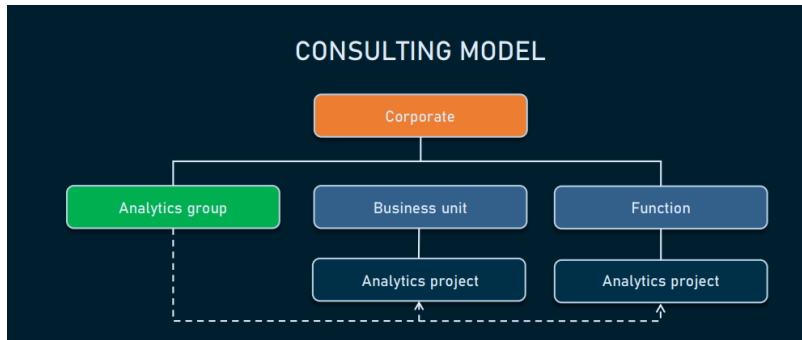
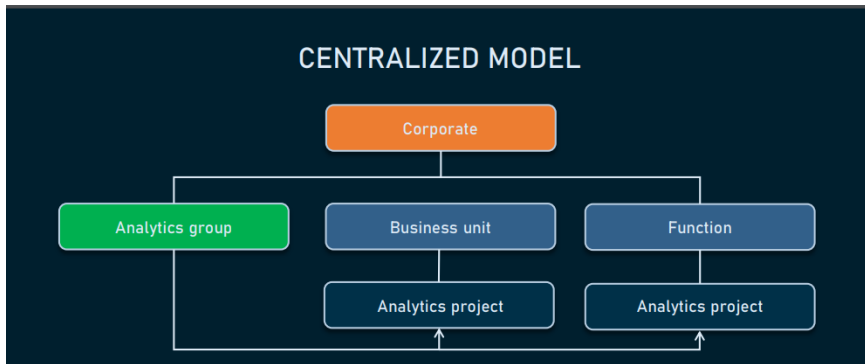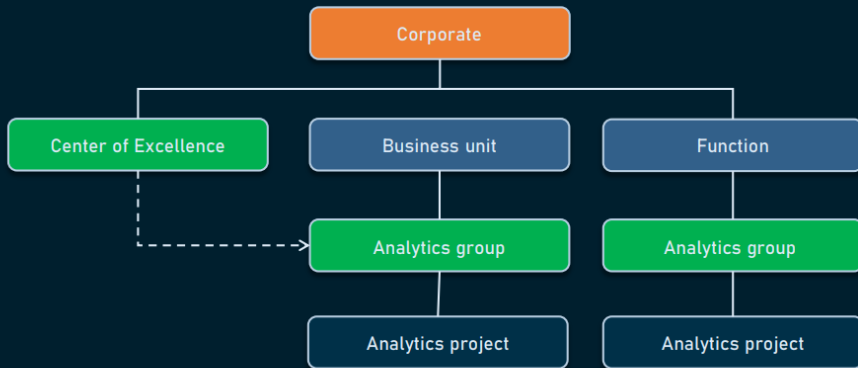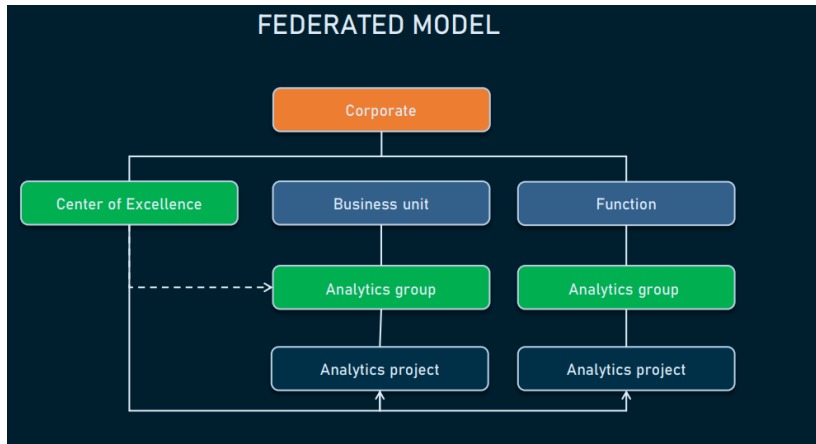| | | |
|---|---|---|
| Analytics | R/SAS | necessary |
| Coding | R, Python, Java, C/C++ | necessary |
| Databases | SQL, NoSQL (MongoDB, CouchDB, Cassandra, MemcacheDB, etc.) | necessary |
| Big Data Processing | Hadoop, Spark, Flink | preferred |
| Algorithms and Models | Regression models, Hidden Markov models, Support Vector Machines, Dimensionality Reduction algorithms, Ensemble algorithms, Decision Trees, Clustering | necessary |
| Frameworks and Libraries | TensorFlow, Theano, CNTK, scikit-learn, Caffe, Spark MLlib, etc. | preferred |
| Domain knowledge | Understanding of company goals, industry fundamentals, business problems, finding new ways to leverage data | preferred |
| Other | Intellectual curiosity, communication and presentation skills | preferred |

altexsoft
software and engineering

# MODELS



DECENTRALIZED MODEL

# MODELS



FUNCTIONAL MODEL

# MODELS

# MODELS

**CENTRALIZED MODEL**

# MODELS

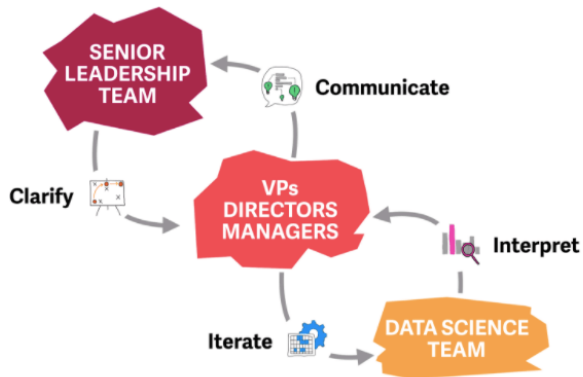# MODELS
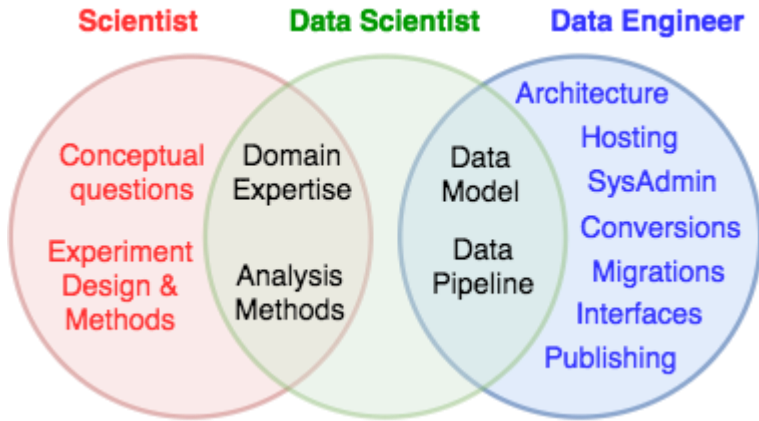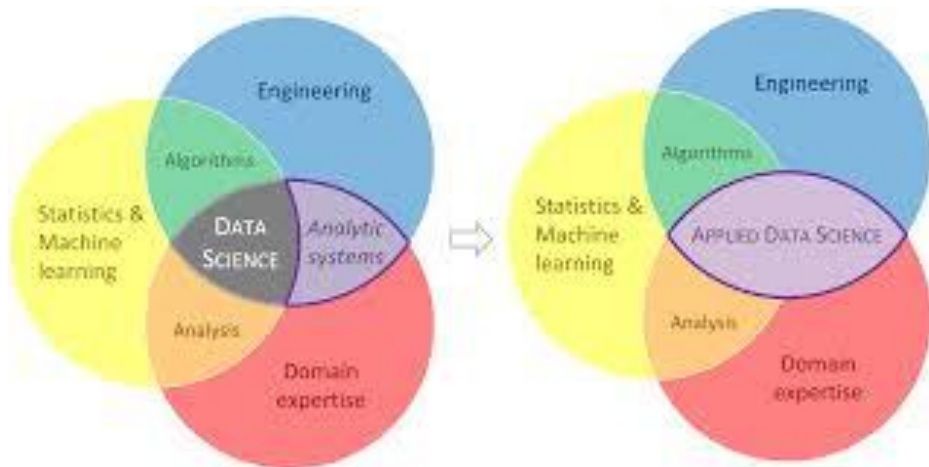


FEDERATED MODEL

# MODELS



DEMOCRATIC MODEL

# DEFINING DATA SCIENCE TEAM



The data science team

# DATA SCIENTIST – JOB ROLE

# UNDERSTANDING DATA SCIENCE

# DATA SCIENCE TEAM EFFORT

| | ⚙ Data Engineers | 🧠 Data Scientists | 🏢 Software Engineers | 📊 Data Storyteller/Translators |
|---|---|---|---|---|
| **What They Do** | • Create Data pipelines.<br>• Evaluate Databases<br>• Design Schemas<br>• Perform ETL | • Apply statistical/Machine learning techniques to solve business problems<br>• Perform R&D<br>• Innovate new solutions<br>• Develop Data science products | • Help design UI (front end coding)<br>• Do backend coding<br>• Help deploy data science solution in production<br>• Automate the entire process | • Communicate Data Science solutions in Business friendly/ non technical terms<br>• Understand business requirements and translate them to Data science problems<br>• Design persuasive Data visualizations |
| **Skill Set** | • Knowledge of Databases<br>• Scripting skills (Linux commands)<br>• Knowledge of Cloud technologies<br>• SQL commands | • Knowledge of statistical and mathematical concepts<br>• Knowledge of various statistical/ML algorithms<br>• Scripting skills (R/Python)<br>• SQL commands | • Knowledge of Programming concepts<br>• Programming languages<br>• Knowledge of Databases<br>• Knowledge of Restful APIs<br>• Scripting skills (Linux commands) | • High level understanding of statistics and ML concepts<br>• Business acumen<br>• Good soft skills<br>• Creativity<br>• Persuasion and articulation |
| **Tools Used** | Hadoop  amazon webservices<br>TERADATA  ORACLE DATABASE | python  R<br>IBM  §sas<br>SPSS | python  Java<br>django  C# | +tableau  Office |

R Venkat Raman

# MATURITY LEVELS WITH DATA

| | Data Engineering | | | Data Science | | | Data as 'Culture' |
|---|---|---|---|---|---|---|---|
| **Phases** | Data Engineering | | | Data Science | | | Data as 'Culture' |
| **Maturity** | Data Collection | Data Storage | Data Transformation | Reporting | Insights | Consumption | Decisions |
| **Activitie** | Int/External <br><br> Logs, IOT <br><br> Stage/Stream | Un/Structured <br><br> SQL, Spark.. <br><br> Data lake.. | ETL <br><br> Cleaning <br><br> Preparation | Metrics/KPI <br><br> Aggregates <br><br> Reports | EDA <br><br> ML <br><br> AI | Narrative <br><br> Info Design <br><br> Data Stories | Change Mgmt <br><br> Workflows <br><br> Actions |

https://techcrunch.com/2019/12/13/when-and-how-to-build-out-your-data-science-team/

BITS Pilani, Deemed to be University under Section 3 of UGC Act, 1956

# MATURITY LEVELS WITH DATA



| Maturity Phases | Data Engineering | | | Data Science | | | Data as 'Culture' |
|---|---|---|---|---|---|---|---|
| | Data Collection | Data Storage | Data Transformation | Reporting | Insights | Consumption | Decisions |
| **Activities** | Int/External<br><br>Logs, IOT<br><br>Stage/Stream | Un/Structured<br><br>SQL, Spark..<br><br>Data lake.. | ETL<br><br>Cleaning<br><br>Preparation | Metrics/KPI<br><br>Aggregates<br><br>Reports | EDA<br><br>ML<br><br>AI | Narrative<br><br>Info Design<br><br>Data Stories | Change Mgmt<br><br>Workflows<br><br>Actions |

https://towardsdatascience.com/whats-the-secret-sauce-to-transforming-into-a-unicorn-in-data-science-94082b01c39d

# SKILLS

Knowledge of machine learning

Understand multiple analytical functions

**What skills make a DATA SCIENTIST?**

Strong knowledge of Python, SAS, R, Scala

Ability to work with unstructured data from various sources like video and social media

Hands-on experience in SQL database coding

# 5 ROLES & SKILLS IN DATA SCIENCE

| | 1. Data Translator | 2. Data Scientist |
|---|---|---|
| **Responsibilities** | - Own from inception to adoption<br>- Translate across domain & data<br>- Act as a glue in the team | - Devise analytics approach<br>- Analyze data & identify insights<br>- Build ML models |
| **Skills** | - Domain expertise<br>- Business analysis & solutioning<br>- Interpersonal & mentoring skills | - Statistics and machine learning<br>- Identify & interpret insights<br>- Scripting skills |
| **Closest role** | Business analyst, Domain experts | Statistician, ML experts |

28

https://techhq.com/2019/12/a-complete-data-science-team-requires-more-than-just-data-scientists/

BITS Pilani, Deemed to be University under Section 3 of UGC Act, 1956

# 5 ROLES & SKILLS IN DATA SCIENCE

## 3. Information Designer  4. ML Engineer

| | 3. Information Designer | 4. ML Engineer |
|---|---|---|
| **Responsibilities** | - Ensure consumption of insights<br>- Design information architecture<br>- Understand user, drive adoption | - Package data science solution<br>- Productionizing, DevOps<br>- Data pipelines/integration |
| **Skills** | - Information design<br>- User centered design<br>- Aspects of interface/visual design | - Software engineering<br>- Data handling<br>- Front-end / Back-end coding |
| **Closest role** | UX Designer, Interaction designer | Software engineer, Data architect |

29

https://techhq.com/2019/12/a-complete-data-science-team-requires-more-than-just-data-scientists/

BITS Pilani, Deemed to be University under Section 3 of UGC Act, 1956

## 5. Data Science Manager

**Responsibilities**
- Identify roadmap & scale maturity
- Ensure biz value from data science
- Drive a culture of data

30

**Skills**
- Project management
- Business analysis, solutioning
- Team handling

**Closest role**

Project manager, Business analyst

https://techhq.com/2019/12/a-complete-data-science-team-requires-more-than-just-data-scientists/

BITS Pilani, Deemed to be University under Section 3 of UGC Act, 1956

# 5 CORE ROLES IN DATA SCIENCE



Data Science Manager

Data Translator

Data Scientist

Information Designer

ML Engineer

# ROLES – EXTENDED WITH KNOWLEDGE

| | Domain Expertise | Technical Knowledge | Quantitative Skills |
|---|:---:|:---:|:---:|
| Data Scientist | Some | Some | Significant |
| Data Engineer | Minimal | Significant | Some |
| Data Science Architect | Minimal | Significant | Minimal |
| Data Science Developer | Minimal | Significant | Some |
| Product Owner | Some | Minimal | Minimal |
| Data/Business Analyst | Some | Some | Some |
| Process Master | Minimal | Minimal | Minimal |
| Subject Matter Expert | Significant | Minimal | Minimal |

Significant Expertise: ●    Some Expertise: ◔    Minimal Expertise: ○

# DATA SCIENCE & ANALYTICS

# DE vs DA vs DJ



Data Journalist
- Experience as working journalist
- Data visualization
- Graphics creation
- Cartography tools
- Adobe Creative Suite
- CSS, D3, HTML5

Data Analyst
- Experience of data analysis / working as a data analyst
- Data mining
- Data management

- Domain knowledge
- Statistical analysis
- Data sourcing
- Collaborative team skills
- Python
- SQL
- R, Excel

- Quantified experience

- Data design & modelling
- SAS

- Experience in data engineering /processing/warehousing
- Experience with large IT /amounts of raw data
- Hadoop/MapReduce/Hive/Pig
- ETL

Data Engineer

# DATA SCIENTIST



**Languages**
R, SAS, Python, Matlab, SQL, Hive, Pig, Spark

**Skills & Talents**
✓ Distributed computing
✓ Predictive modeling
✓ Story-telling and visualizing
✓ Math, Stats, Machine Learning

## DATA SCIENTIST
'AS RARE AS UNICORNS'

**Role**
Cleans, massages and organizes (big) data

**Mindset**
Curious data wizard

HIRED BY

Google  Microsoft  Adobe

# DATA ANALYST



DATA ANALYST
DATA DETECTIVE

**Role**
Collects, processes and performs statistical data analyses

**Mindset**
Intuitive data junkie with high "figure-it-out" quotient

**Languages**
R, Python, HTML, Javascript, C/C++, SQL

**Skills & Talents**
✓ Spreadsheet tools (e.g. Excel)
✓ Database systems (SQL and NO SQL based)
✓ Communication & visualization
✓ Math, Stats, Machine Learning

HIRED BY

IBM (hp) DHL

# BUSIESS ANALYST

# THREE EMERGING ROLES

# 1. DATA STORYTELLER

## Senator Voting Patterns

When senator 'X' votes a 'Yea' or 'Nay' what are the chances that senator 'Y' would do the same? This tool allows you to find out the similarity in voting patterns of senators of the 115th Congress.

Here are a few senators with interesting voting patterns - Joe Manchin (The Democrat who votes like a Republican), Thad Cochran, Robert Menendez, Heidi Heitkamp, Joe Donnelly, Elizabeth Warren, Claire McCaskill, Kirsten Gillibrand, Angus King, Bernie Sanders, John McCain, Rand Paul & John Isakson.

King (I-ME)

Powered by GRAMENER.COM
Data Courtesy: www.senate.gov

### Voting Patterns of Senators

The dark stroked circle at the center is the selected senator. The distance between the senator and other senators around him/her defines the voting similarity score. Closer to the center greater the similarity in voting pattern and vice versa.

Click on any senator to view the Voting Similarity score.

🔴 Rep    🔵 Dem    🟢 Ind

King (I-ME)    Carper (D-DE)

On what issues do Senator King (I-ME) & Senator Carper (D-DE) agree & disagree? Click on the image of Senator Carper (D-DE) to find out.

**88%**

**Voting Similarity**

**Role Highlights**

- **Dashboards are NOT data stories**
- **Stories=visual+context+narrative**
- **Fields: Journalism, creative arts**

https://gramener.com/playground/senate/similarity

**Role Highlights**

- **Human side of data insights**
- **More practical, 'accurate' results**
- **Fields: Social sciences**

# 3. DATA ETHICIST



A top BFSI player wanted a scientific way to identify peers, for employee feedback.

Was there an alternative to manually screening for peer review?

This visual shows the network of email exchanges between people.

Look for the closest neighbors. The distance is a function of email exchange.

**Role Highlights**

- Ensure trust & fairness

- Act as a collective conscience

- Fields: Law, Humanities

Gramener Email Communication Analysis          Brochure (PDF)   Video demo   Gramener Employee Data ▾

Employee email connections

Pratap_Vardhan

| Recipient | Email | | Count |
|---|---|---|---|
| Hr | hr@gramene. | | 126 |
| S_Anand | s.anand@gr. | | 96 |
| Naveen_Gattu | naveen.gat. | | 72 |
| Gitlab | gitlab@cod. | | 57 |
| Mukul_Taneja | mukul.tane. | | 36 |
| Kathir_Mani | kathir.man. | | 30 |
| Soumya_Ranjan | soumya.ran. | | 30 |
| Talent | talent@gra. | | 30 |
| Raghunandh | raghunandh. | | 24 |
| Siva_Sangubotla | siva.sangu. | | 24 |

Network of top 10 recipients out of 17 for Pratap_Vardhan

Link distance is inversely proportional to the count, thickness is proportional to the count and right click on person to collapse.

● Normal person
● Interacted person
● Selected person

BITS Pilani, Deemed to be University under Section 3 of UGC Act, 1956

# 3 EMERGING SKILLS IN DATA SCIENCE



**1. Data Storyteller**

**2. Behavioral Psychologist**

**3. Data Ethicist**

| Phases | Data Engineering | | | Data Science | | | Data as 'Culture' |
|---|---|---|---|---|---|---|---|
| **Maturity** | Data Collection | Data Storage | Data Transformation | Reporting | Insights | Consumption | Decisions |
| **Activitie** | Int/External  Logs, IOT  Stage/Stream | Un/Structured  SQL, Spark..  Data lake.. | ETL  Cleaning  Preparation | Metrics/KPI  Aggregates  Reports | **Data Scientist**  EDA  ML  AI  **Data Ethicist** | **Behavioral Psychologist**  Info Design  **Info Designer**  Data Stories  **Storyteller**  Packaged App  **ML Engineer** | Workflows  Actions  Change Mgmt |

**Data Science**  **Data Translator**

Gramener

# DATA SCIENCE PROJECT HIERARCHY

# DATA ANALYTICS PROCESS



**Data Analytics Process**

**Monitor Results** 7
- Ongoing testing
- Looking for evidence of a need to refresh the model, or replace with a new one
- Can spawn new ideas to test

**Data Acquisition & Management**
- Gather data – often from multiple systems and LOBs 2
- Prepare data – making it ready for use by the analytic model(s) 3

**Model Deployment** 6
- Running on production data at useful intervals or when called by business applications
- Exporting in useable formats

1 Business Objective

**Model Validation** 5
- Independent testing model performance on sample data
- Comparing model predictions to actual results

**Model Creation** 4
- Can be in-house developed or delivered as part of an analytics package
- Become useful only when actionable results are manifest

Source: Celent

**Requirement Flow**

| Data Acquisition | Data exploration | Feature Engineering | Model Selection | Experimentation | Prediction |

**Implementation Flow**

| Data Engineer | Data Scientist | Data Scientist | Data Scientist | Data Scientist | Business User |

# ROLES IN TIMELINE

# Managing Data Science Team

- **Knowledge Management**
- **Attracting Top Talent**
- **Hiring Process**
- **Onboarding**
- **Retention and Management**

Focus –

- On boarding and evaluating the success of team
- Working with other teams

# HOW TO CREATE AN ENVIRONMENT FOR TEAM SUCCESS?



- Promote collaboration within teams

- Align closely with business users

- Measure outcomes through business value

- Adopt process frameworks for consistency

- Up-skill continuously on tech skillsets

- Nudge cross-training across disciplines

# Common Challenges

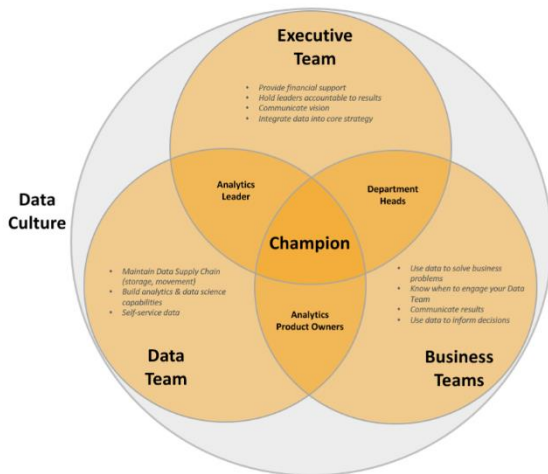1. Hiring a balanced data science team - build a cross-functional data science team that enables your organization to get insights from data and build production ready models.( Data Scientist, machine learning Engineer, Data Architect/Engineer, SW Developer, Research Scientist)
2. Retaining the team and Growing the team
3. Translating the business goals to smaller chunks of tasks, and defining measurable KPIs for the Data Science Team to work on achieving these KPIs.
4. Transforming Data Science team output/deliverables to a business understandable form, with key focus on Data Visualization. Hence try to bridge the gap between the Business Teams who relatively less/non-technical and the very technical Data Science Team
5. Engage and keep team motivated during the failures, and also keep the Senior leadership aligned with the fact that Data Science Projects are not like any SW Engineering project which can very Agile and give results every 7 days.

# DATA CULTURE

BITS Pilani, Deemed to be University under Section 3 of UGC Act, 1956

# Data driven decision making

- **Definition** – When it is data and not instincts that drives the business decisions.
- **Examples** – Fraud detection in Loans, Credit Cards (Cibil scores); Insurance, Six sigma projects to improve efficiency; Target advertising in e-commerce; Product Roadmap planning, Team planning
- 6 Steps to Data Driven decision making-
  1. Strategy – Define clear Business goals
  2. Identifying key data focus areas – Data is everywhere, flowing from multiple sources. Based on domain knowledge define key focus data sources which seem to impact the most, easier to access, reliable and clean
  3. Data Collection & Storage – Defining data architecture to collect, store, archive i.e. manage data. Connect multiple data sources, clean, prepare and organize
  4. Data Analytics – Analyzing the data and derive key insights
  5. Turning insights to Actions – business actions to be taken based on the findings from key insights from data
  6. Operationalize and Deploy – Using IT systems, automate the data collection, storage, analysis and presenting the key highlights

# THANK YOU