



Data Science Bootcamp

Hyunsuk (hyssh@microsoft.com)

Microsoft Technology Center in Seattle

Microsoft Technology Centers

Bootcamp

Goal

The bootcamp will walk you through the concept of data science and enable you to start the data science project quickly



Agenda

- Data Science
 - Five questions
 - Ask the right question
 - How to make a prediction
- Data Science Process
 - Business understanding
 - Data acquisition and understanding
 - Modeling
 - Deployment
 - Customer Acceptance
- Data Science Tool
 - Access and process data
 - Train and deploy model
 - Monitor models and data

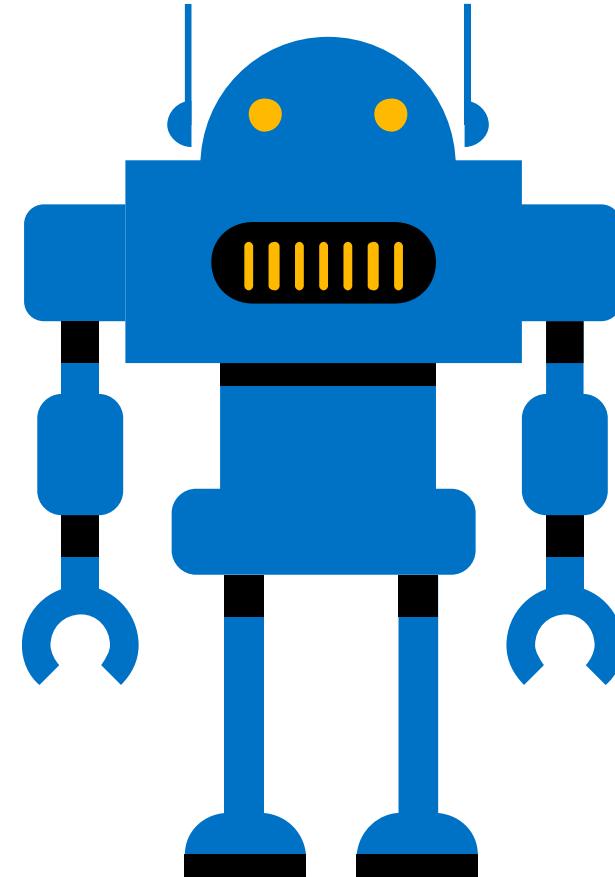
What is Data Science

The Formal one:

"A computer program is said to learn from experience **E** with respect to some class of tasks **T** and performance measure **P** if its performance at tasks in **T**, as measured by **P**, improves with experience **E**."

A Practical Example:

Look at data. Do the thing. Better? No? Look at the data. Do something different. Better? Yes? *Do that again.* (Repeat)



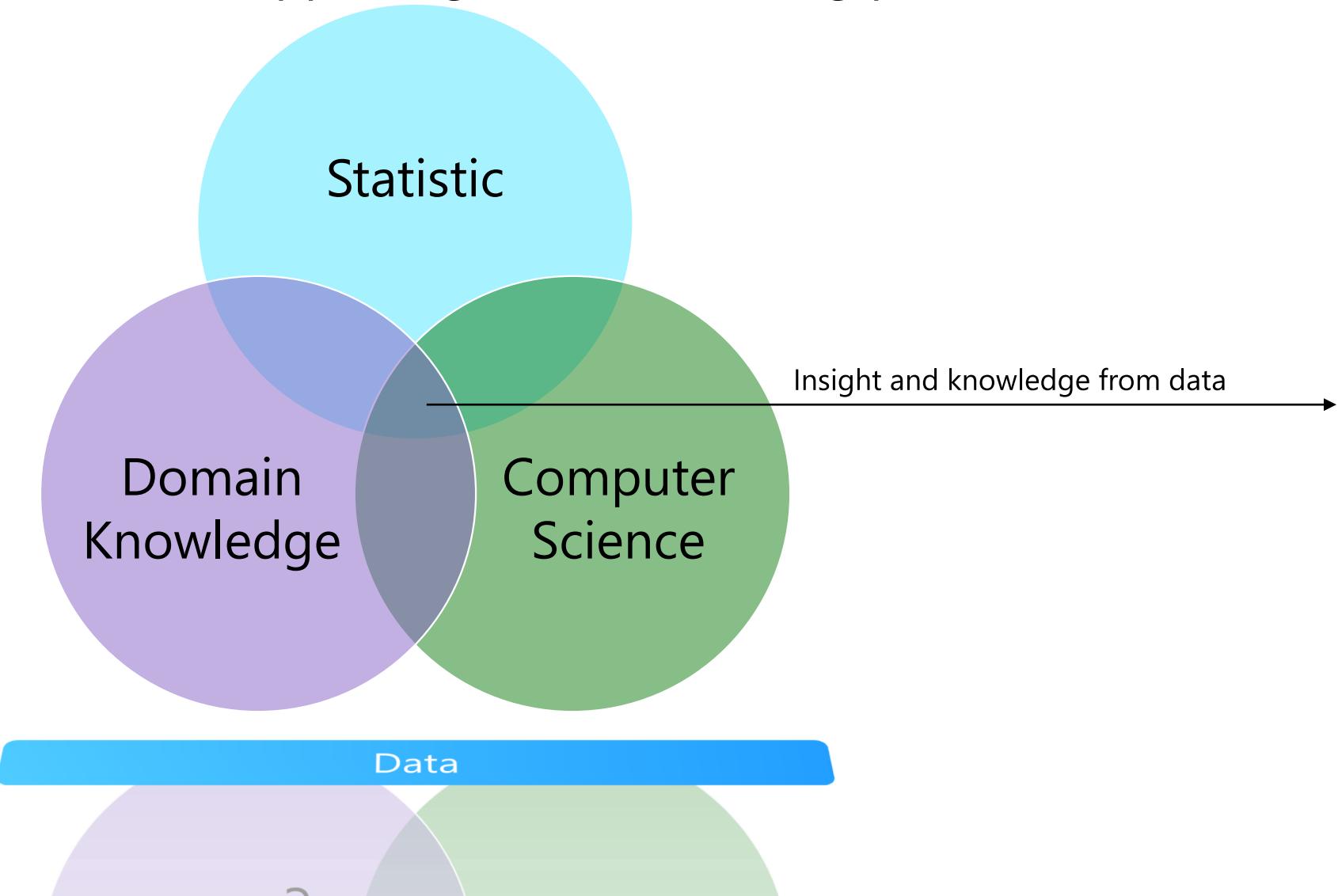
What is Data Science

HS tell me what is data science

 Data Science is an interdisciplinary field that combines statistical methods, computer science techniques, and domain knowledge to extract insights and knowledge from data. It involves using mathematical and computational methods to extract information and insights from structured and unstructured data. The goal of data science is to create value from data by discovering patterns, building predictive models, and supporting decision-making processes. Data Scientists use a variety of tools and techniques to process and analyze large datasets, such as statistical analysis, machine learning, data visualization, and database management. The field of data science is growing rapidly, as the amount of data being generated in various domains continues to increase.

Why Data Science

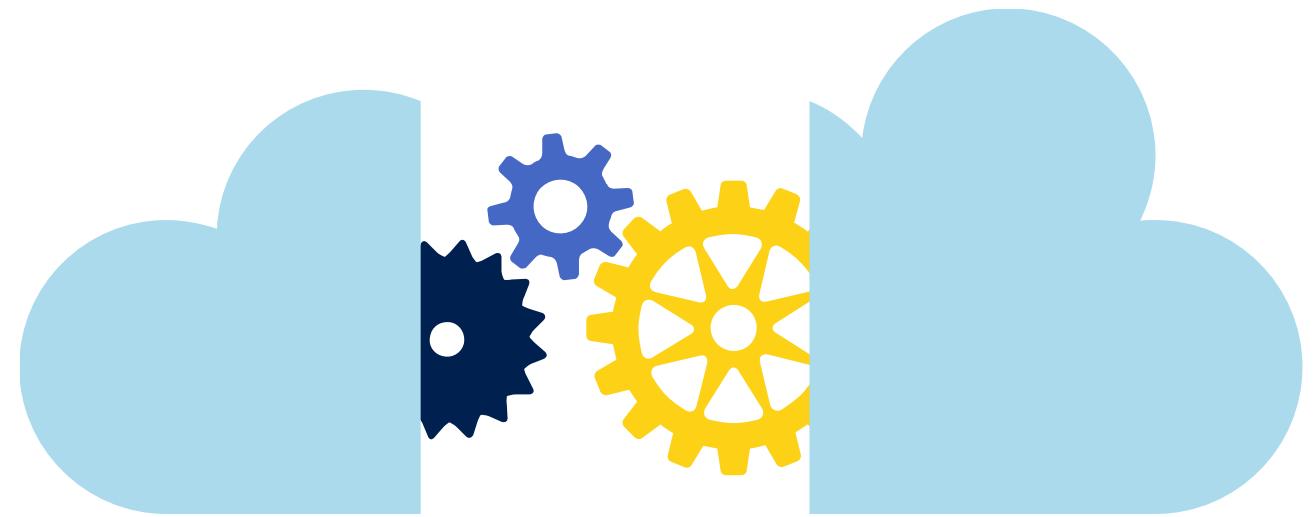
Data science creates value from data by discovering patterns, building predictive models, and supporting decision-making processes.



Value

Agenda

- **Data Science**
- Data Science process – TDSP (Team Data Science Process)
- Data Science Tools

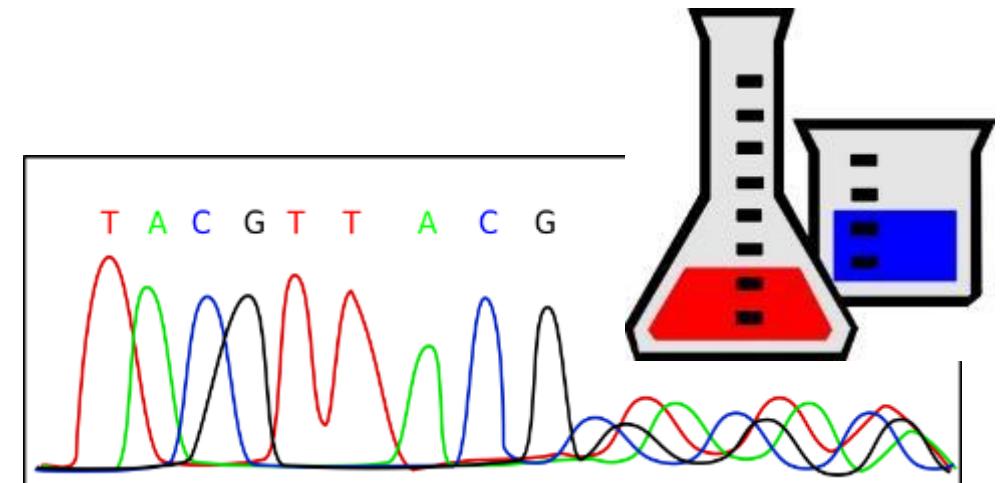


Microsoft Technology Centers

Data Science

Machine Learning

- 1. Five questions**
- 2. Is your data ready?**
- 3. Ask the right question**
- 4. Predict an answer**
- 5. Copy other people's work**



1. Five questions

Data Science

Machine Learning can answer

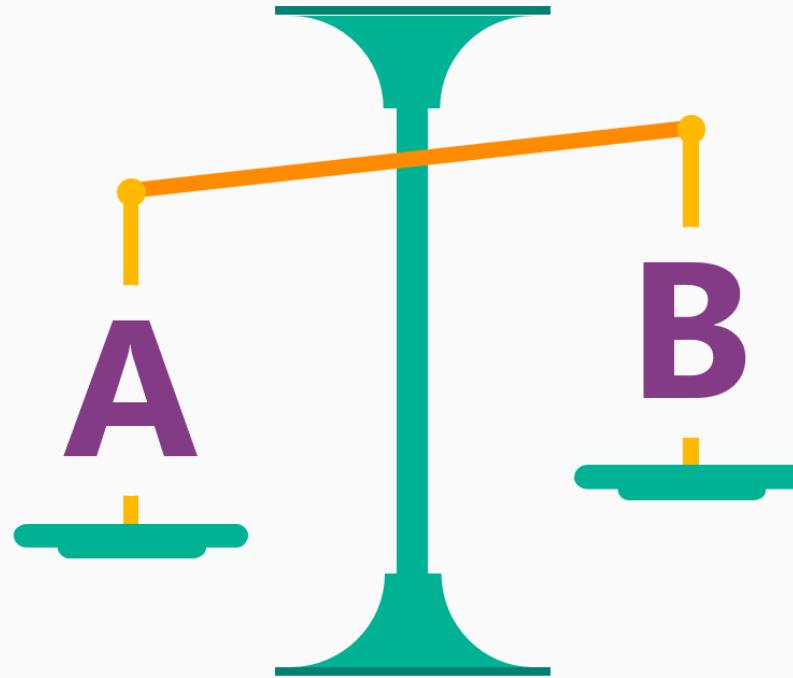
- Is this A or B?
- Is this weird?
- How much? How many?
- How is this organized?
- What should I do?

1. Five questions

Data Science

Is this A or B?

Classification algorithms



1. Five questions

Data Science

Is this weird?

Anomaly detection algorithms

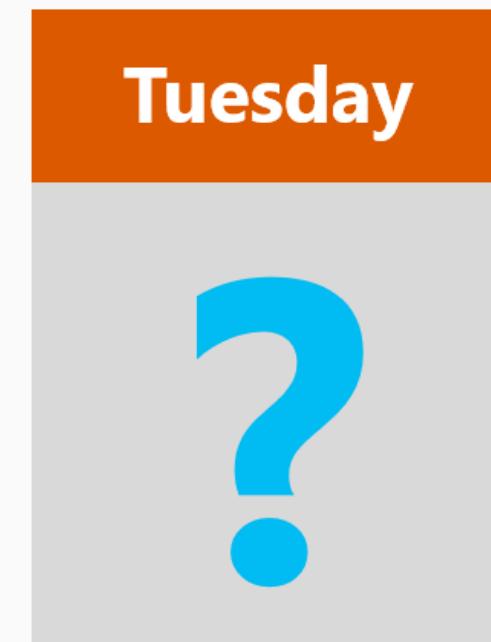


1. Five questions

Data Science

How much? How many?

Regression algorithms

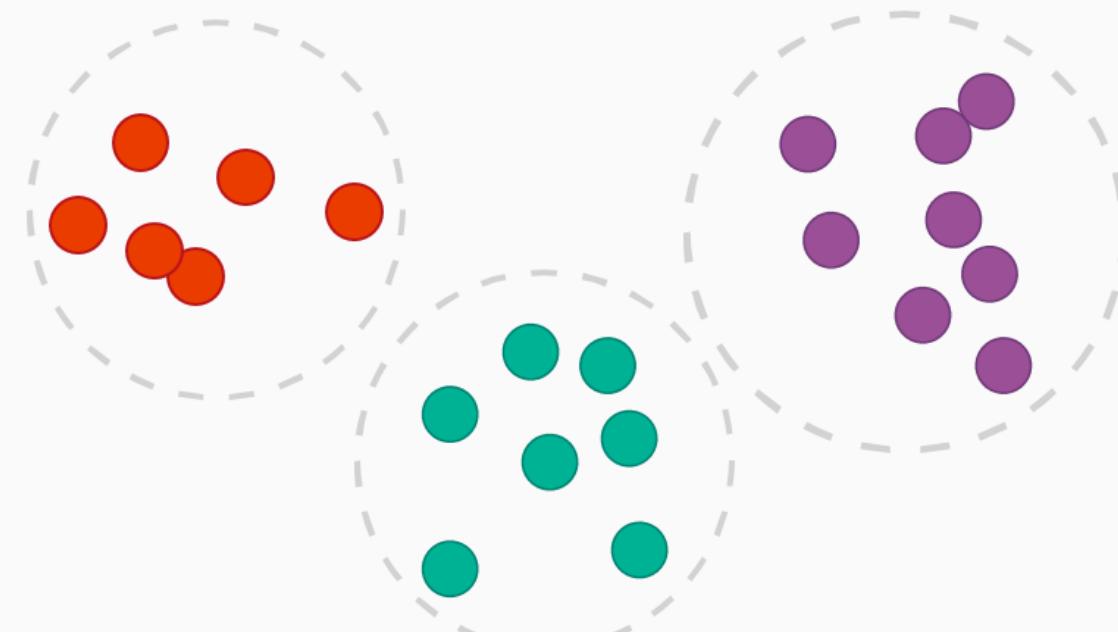


1. Five questions

Data Science

How is this organized?

Clustering Algorithms

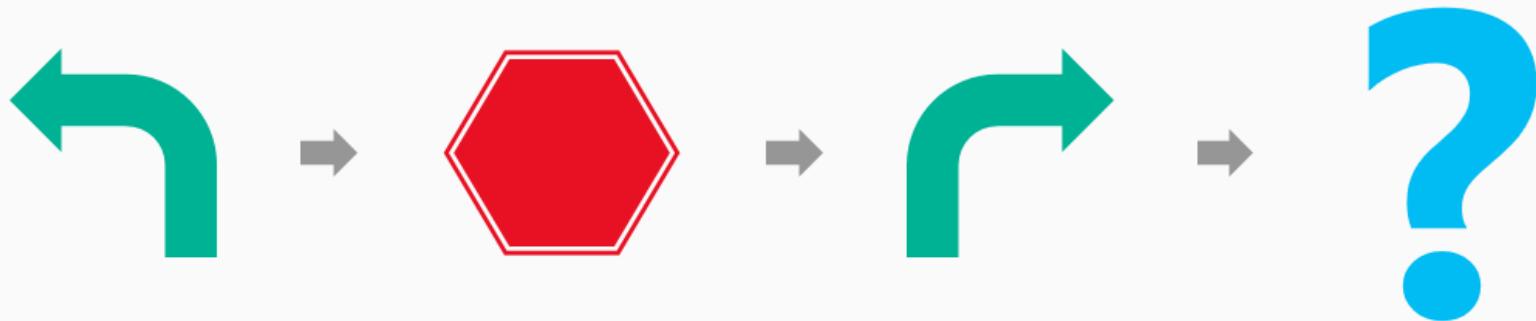


1. Five questions

Data Science

What should I do now?

Reinforcement Learning Algorithms



2. Is your data ready?

- Relevant
- Connected
- Accurate
- Enough to work with

Irrelevant Data

Price of milk (\$/gal)	Red Sox batting avg.	Blood alcohol content (%)
3.79	.304	.03
3.45	.320	.09
4.06	.259	.01
3.89	.298	.05
4.12	.332	.13
3.92	.270	.06
3.23	.294	.10

Relevant Data

Body mass (kg)	Margaritas	Blood alcohol content (%)
103	3	.03
67	5	.09
87	1	.01
52	2	.05
73	5	.13
79	3	.06
110	7	.10

2. Is your data ready?

Data Science

- Relevant
- Connected
- Accurate
- Enough to work with

Disconnected Data

Grill temp. (Fahrenheit)	Weight of beef patty (lb)	Burger rating (out of 10)
	.33	8.2
	.24	5.6
550		7.8
725	.45	9.4
600		8.2
625		6.8
	.49	4.2

Connected Data

Grill temp. (Fahrenheit)	Weight of beef patty (lb)	Burger rating (out of 10)
575	.33	8.2
550	.24	5.6
550	.69	7.8
725	.45	9.4
600	.57	8.2
625	.36	6.8
550	.49	4.2

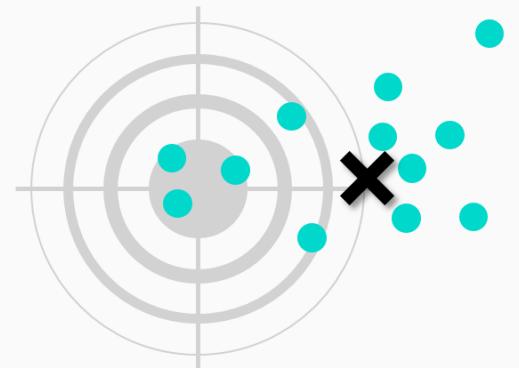
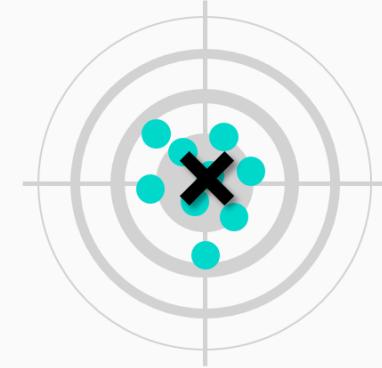
2. Is your data ready?

- Relevant
- Connected
- Accurate
- Enough to work with

Inaccurate Data



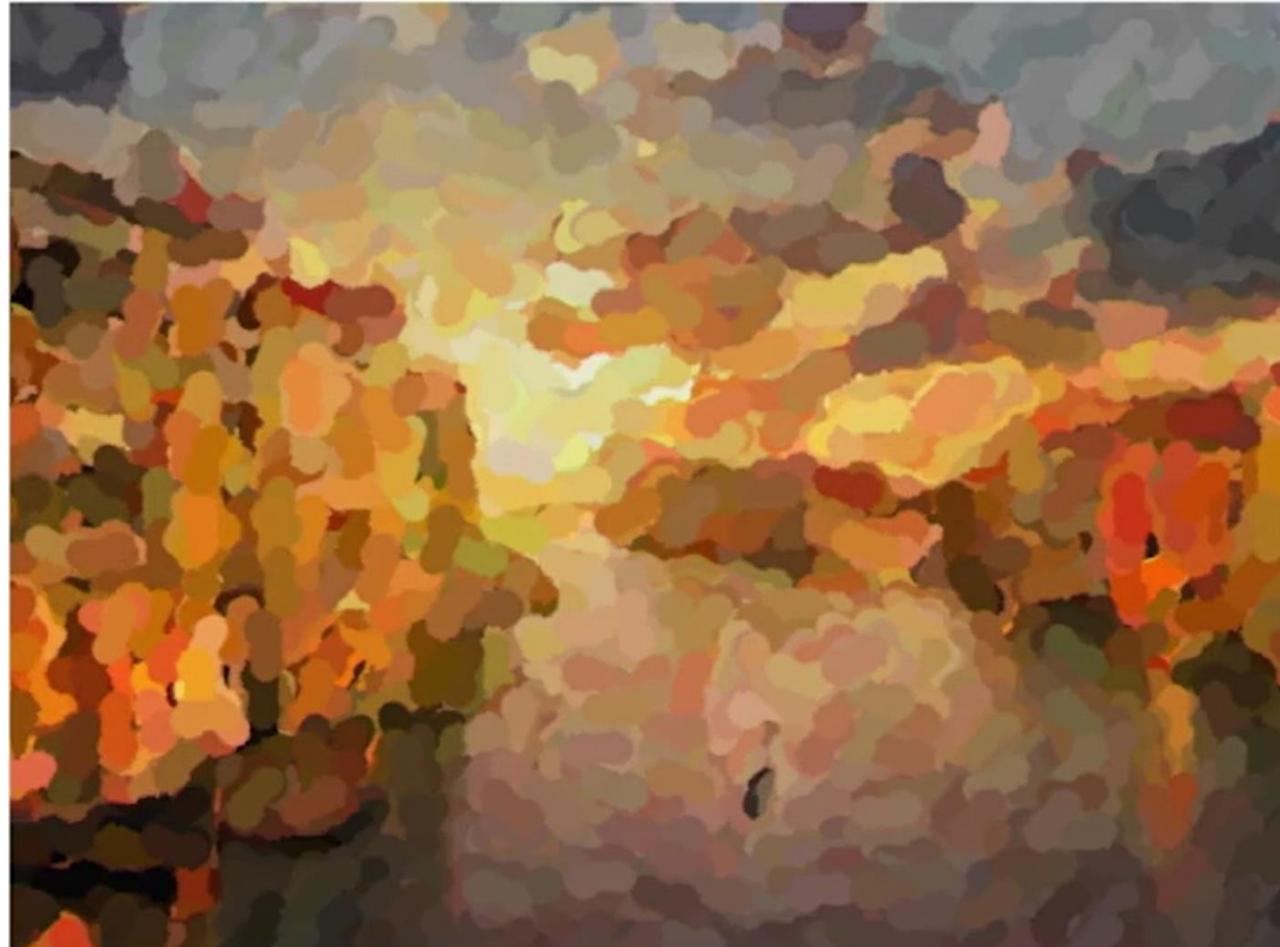
Accurate Data



2. Is your data ready?

Data Science

- Relevant
- Connected
- Accurate
- Enough to work with



2. Is your data ready?

Data Science

- Relevant
- Connected
- Accurate
- Enough to work with



2. Is your data ready?

Data Science

- Relevant
- Connected
- Accurate
- Enough to work with



3. Ask the right question

Data Science

Sharp question



3. Ask the right question

Data Science

Sharp question

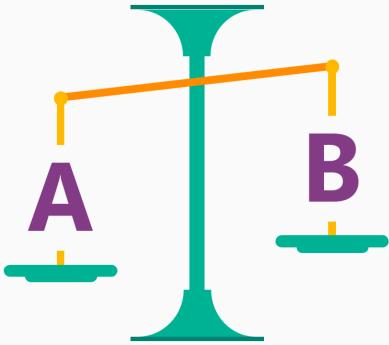
- Will this tire fail in the next 1,000 miles: Yes or no?
- Which brings in more customers: a \$5 coupon or a 25% discount?
- If you have a car with pressure gauges, you might want to know: Is this pressure gauge reading normal?
- If you're monitoring the internet, you'd want to know: Is this message from the internet typical?
- What will my stock's sale price be next week?
- What will the temperature be next Tuesday?
- What will my fourth quarter sales be?
- Which printer models fail the same way?

3. Ask the right question

Reformulate your question

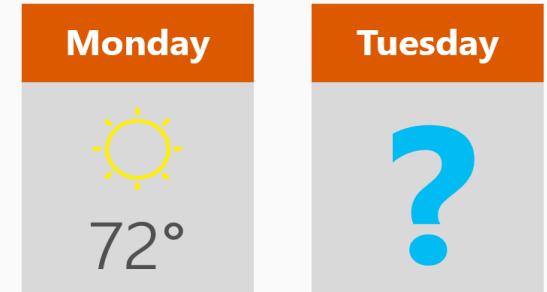
Is this A or B?

Classification algorithms



How much? How many?

Regression algorithms



What will my stock's sale price be next week?

What will the temperature be next Tuesday?

What will my fourth quarter sales be?

Which printer models fail the same way?

4. Predict an answer

How much 1.35 carat diamond cost these days?

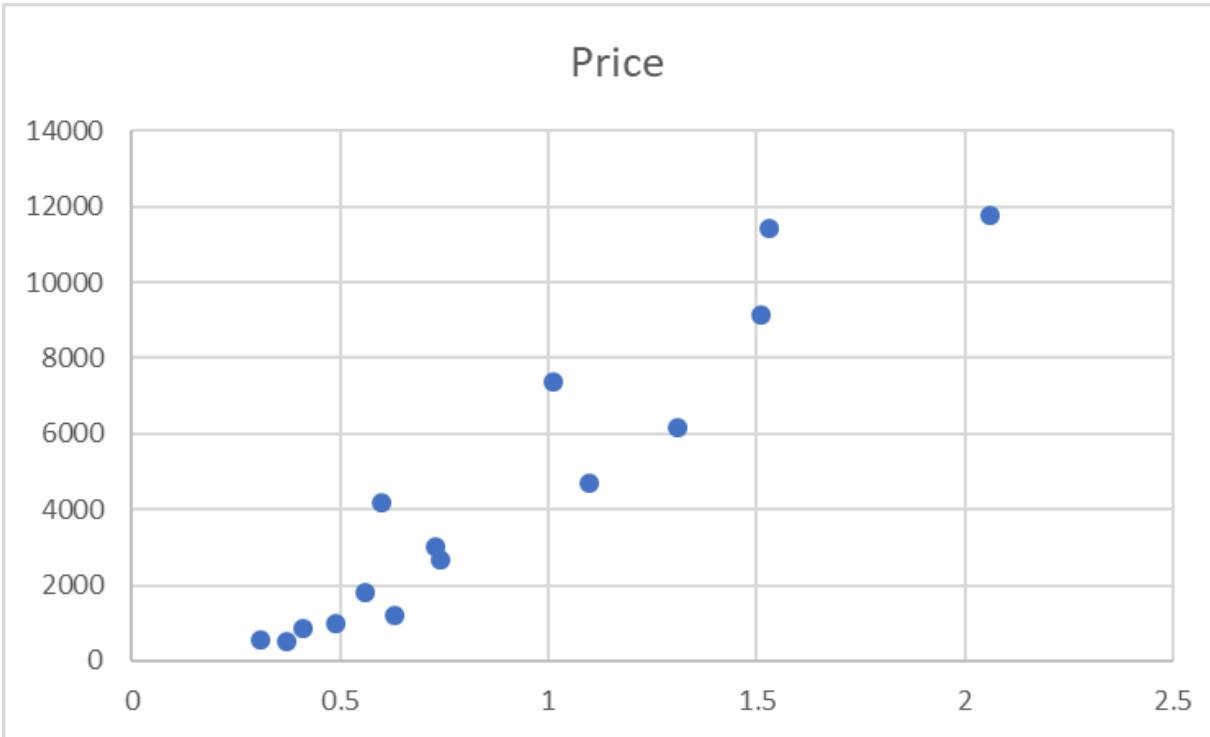
Carats	Price
1.01	7366
0.49	985
0.31	544
1.51	9140
0.37	493
0.73	3011
1.53	11413
0.56	1814
0.41	876
0.74	2690
0.63	1190
0.6	4172
2.06	11764
1.1	4682
1.31	6171

4. Predict an answer

Understand data

- Plot the existing data

Carats	Price
1.01	7366
0.49	985
0.31	544
1.51	9140
0.37	493
0.73	3011
1.53	11413
0.56	1814
0.41	876
0.74	2690
0.63	1190
0.6	4172
2.06	11764
1.1	4682
1.31	6171

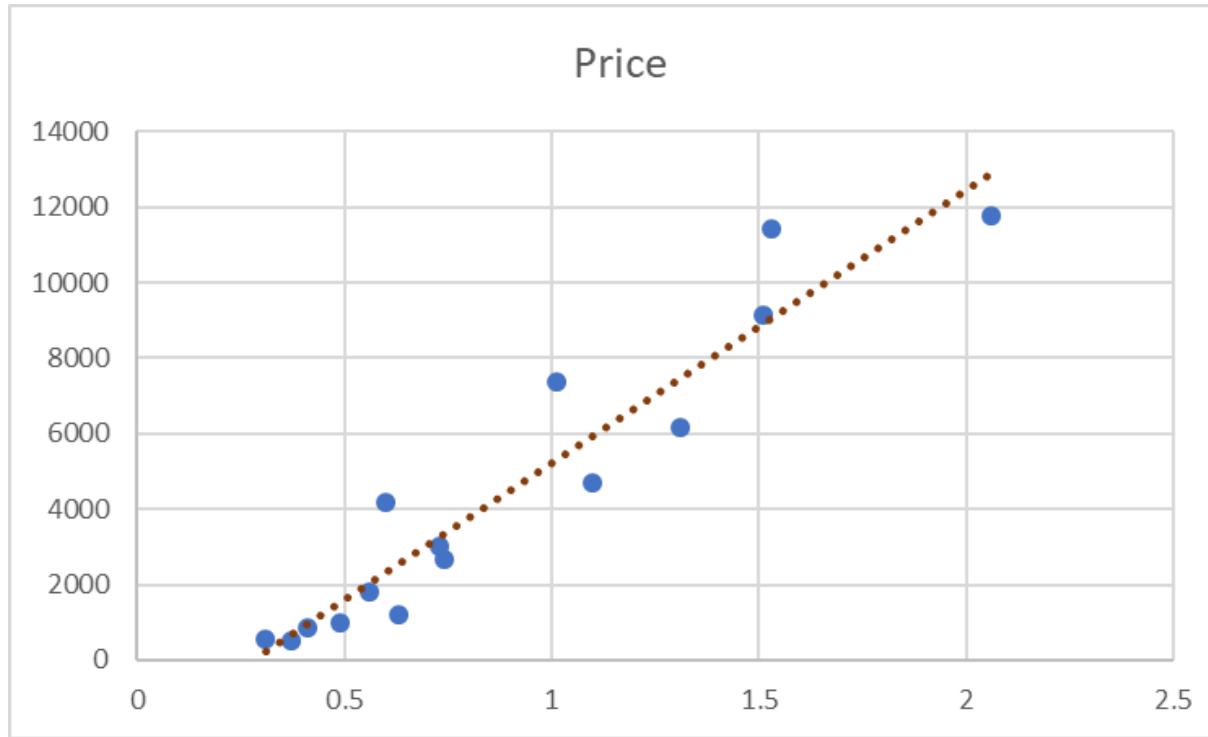


4. Predict an answer

Create a model

- Draw the model through the data points

Carats	Price
1.01	7366
0.49	985
0.31	544
1.51	9140
0.37	493
0.73	3011
1.53	11413
0.56	1814
0.41	876
0.74	2690
0.63	1190
0.6	4172
2.06	11764
1.1	4682
1.31	6171

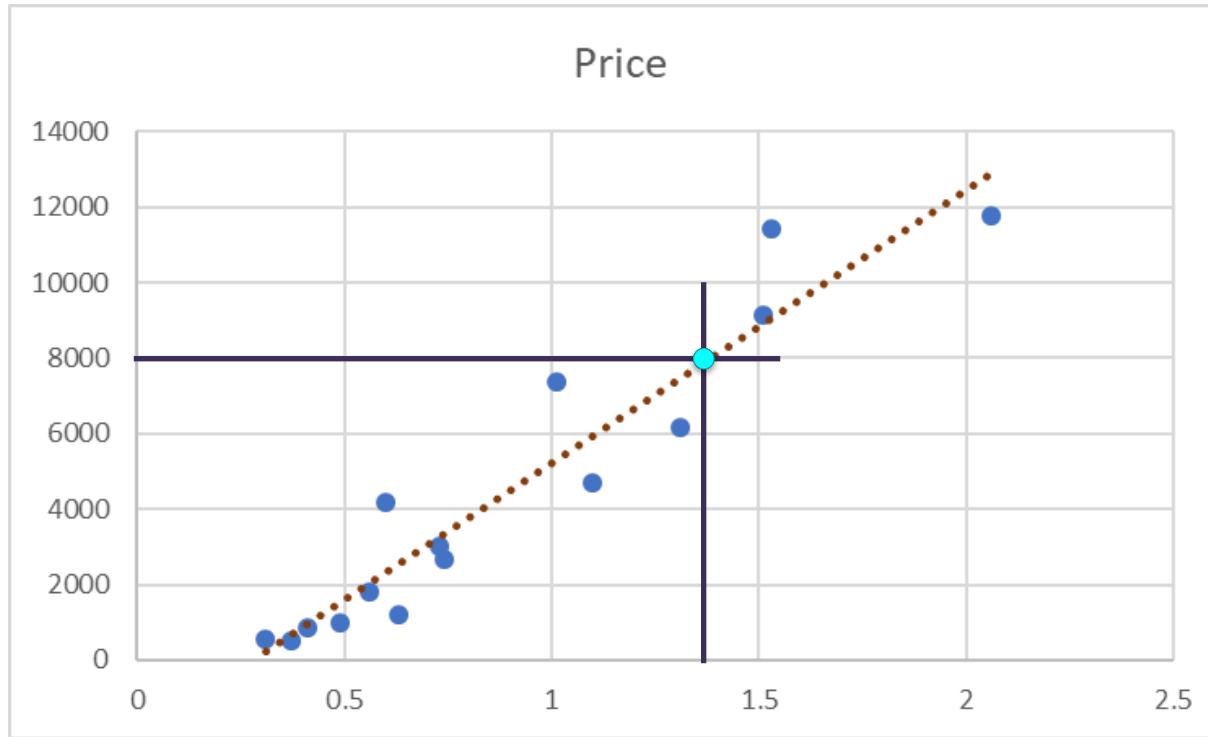


4. Predict an answer

Use the model to find the answer

- How much 1.35 carat diamond cost these days?

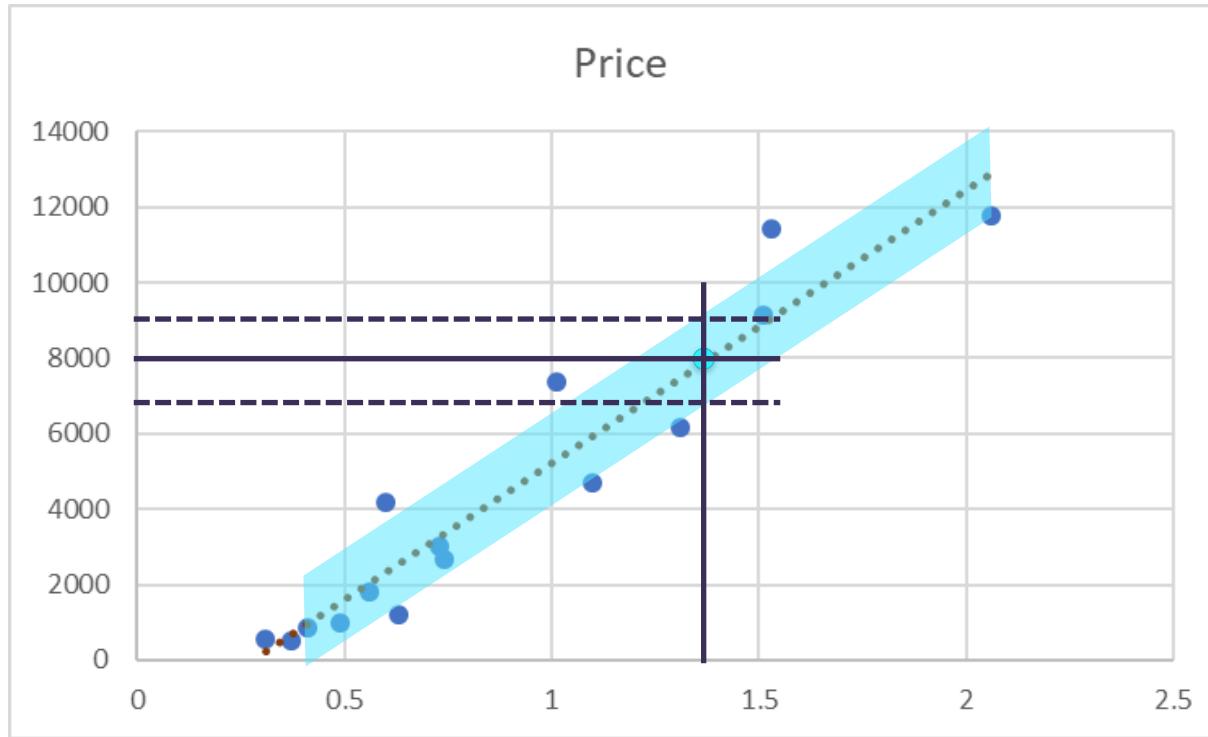
Carats	Price
1.01	7366
0.49	985
0.31	544
1.51	9140
0.37	493
0.73	3011
1.53	11413
0.56	1814
0.41	876
0.74	2690
0.63	1190
0.6	4172
2.06	11764
1.1	4682
1.31	6171



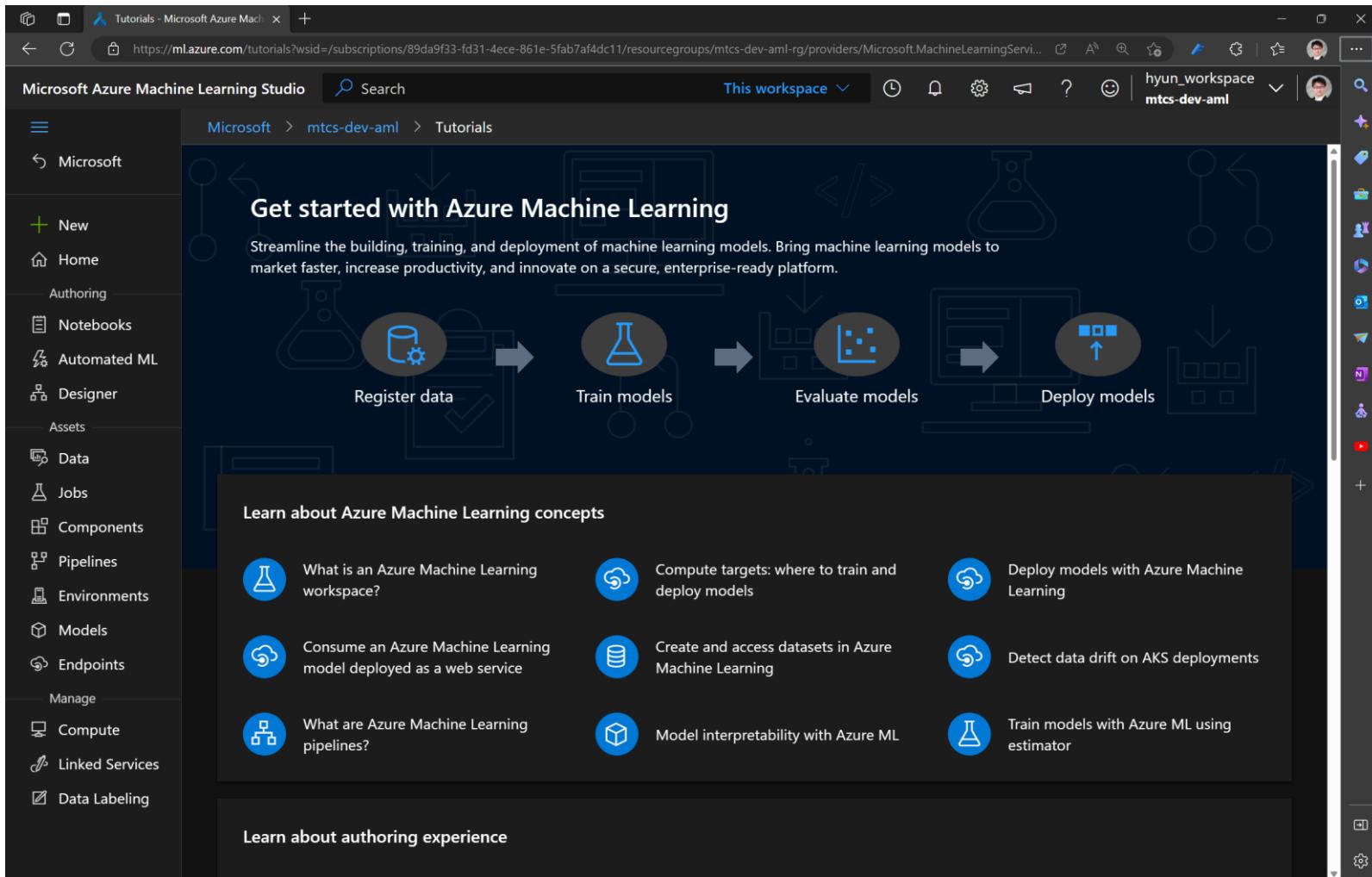
4. Predict an answer

Create a confidence interval

Carats	Price
1.01	7366
0.49	985
0.31	544
1.51	9140
0.37	493
0.73	3011
1.53	11413
0.56	1814
0.41	876
0.74	2690
0.63	1190
0.6	4172
2.06	11764
1.1	4682
1.31	6171



5. Copy other people's work



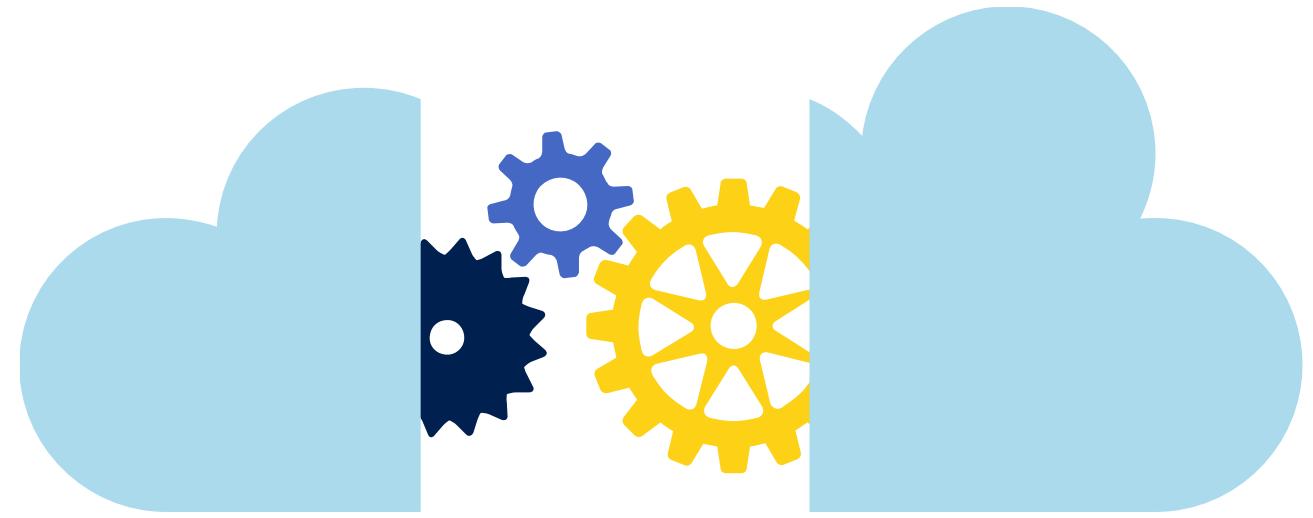
Coffee break

10 minutes



Agenda

- Data Science
- Data Science process – TDSP (Team Data Science Process)
- Data Science Tools

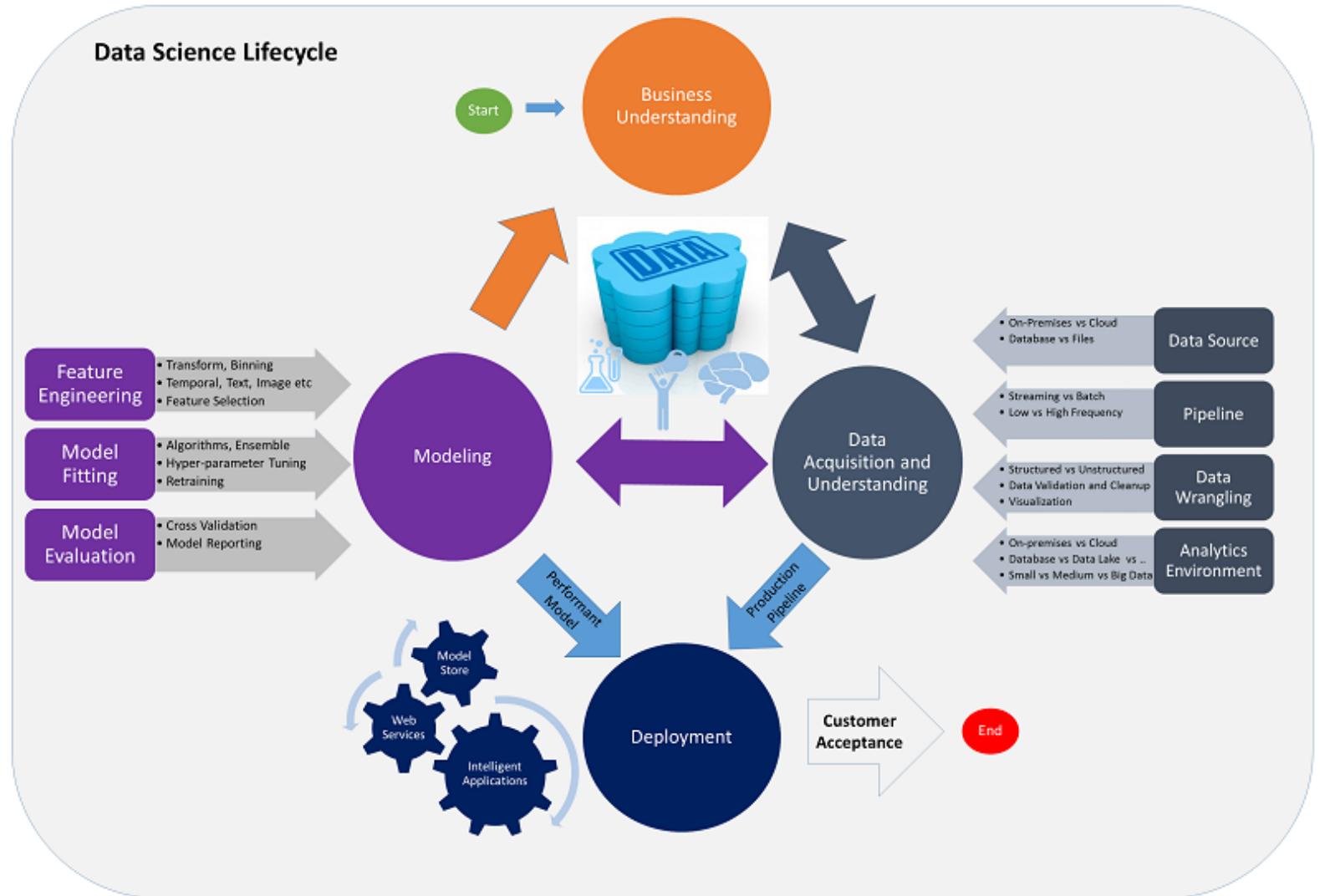


Microsoft Technology Centers

Data Science process – TDSP (Team Data Science Process)

Data Science Lifecycle

1. Business understanding
2. Data acquisition and understanding
3. Modeling
4. Deployment
5. Customer Acceptance



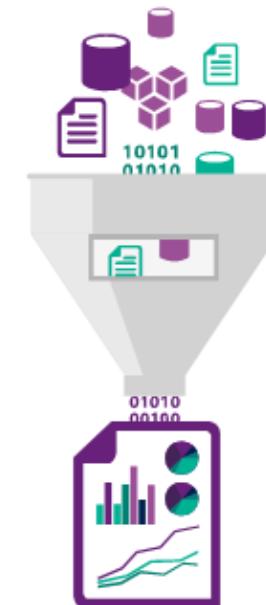
Microsoft Technology Centers

1. Business Understanding

Data Science process – TDSP

Analyze business needs

- Define objectives
 - Identify key business variables
 - Project goal
 - Question
 - Success metrics
- Identify data sources
 - Data sources that contain known examples of answers
 - Data characteristics
 - Data quality
 - Tools and language
 - Security
 - Encryption, Audit, Access control
- Artifacts
 - Documents
 - Data sources
 - Data



1. Business Understanding

Data Science process – TDSP

Business sample case

Wide World Importers is a company that imports and distributes products in multiple countries around the globe.

With several thousand employees, Information Technology is at the heart of our business operations, and has a significant cost.

Since we handle materials in multiple countries, we have a lot of private data, financial information, and other targets which have a high security profile. We are concerned with both external and internal attacks. In addition, many of our employees work in remote locations, some on ships and other challenging environments.

All of our IT systems have been modernized, and we're taking in a significant amount of semi-structured data from computing devices – most of it real-time. After talking with our IT leadership, we need a way to determine anomalies within the data streams we get, and have a way to observe the anomalies in a dashboard so that we can respond to outages, threats, and changes quickly.

1. Business Understanding

Data Science process – TDSP

Design statements

Wide World Importers is a company that imports and distributes products in **multiple countries** around the globe.

With **several thousand employees**, Information Technology is at the heart of our business operations, and has a **significant cost**.

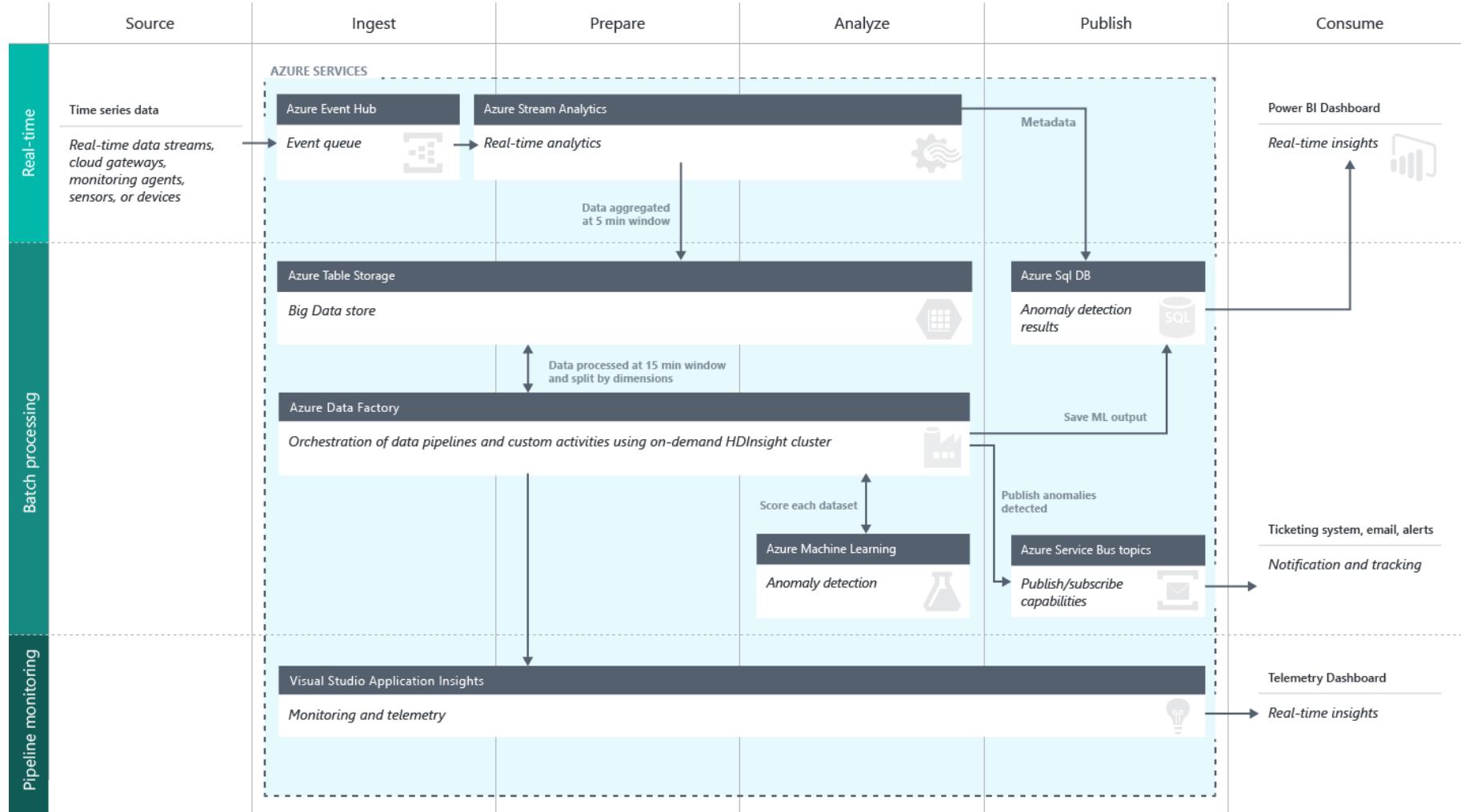
Since we handle materials in multiple countries, we have a **lot of private data, financial information**, and other **targets** which have a **high security profile**. We are concerned with both **external and internal** attacks. In addition, many of our employees work in remote locations, some on ships and other **challenging environments**.

All of our IT systems have been modernized, and we're **taking in a significant amount of semi-structured data** from computing devices – most of it **real-time**. After talking with our IT leadership, we need a way to **determine anomalies** within the **data streams** we get, and have a way to **observe the anomalies** in a **dashboard** so that we can respond to outages, threats, and changes quickly.

1. Business Understanding

Data Science process – TDSP

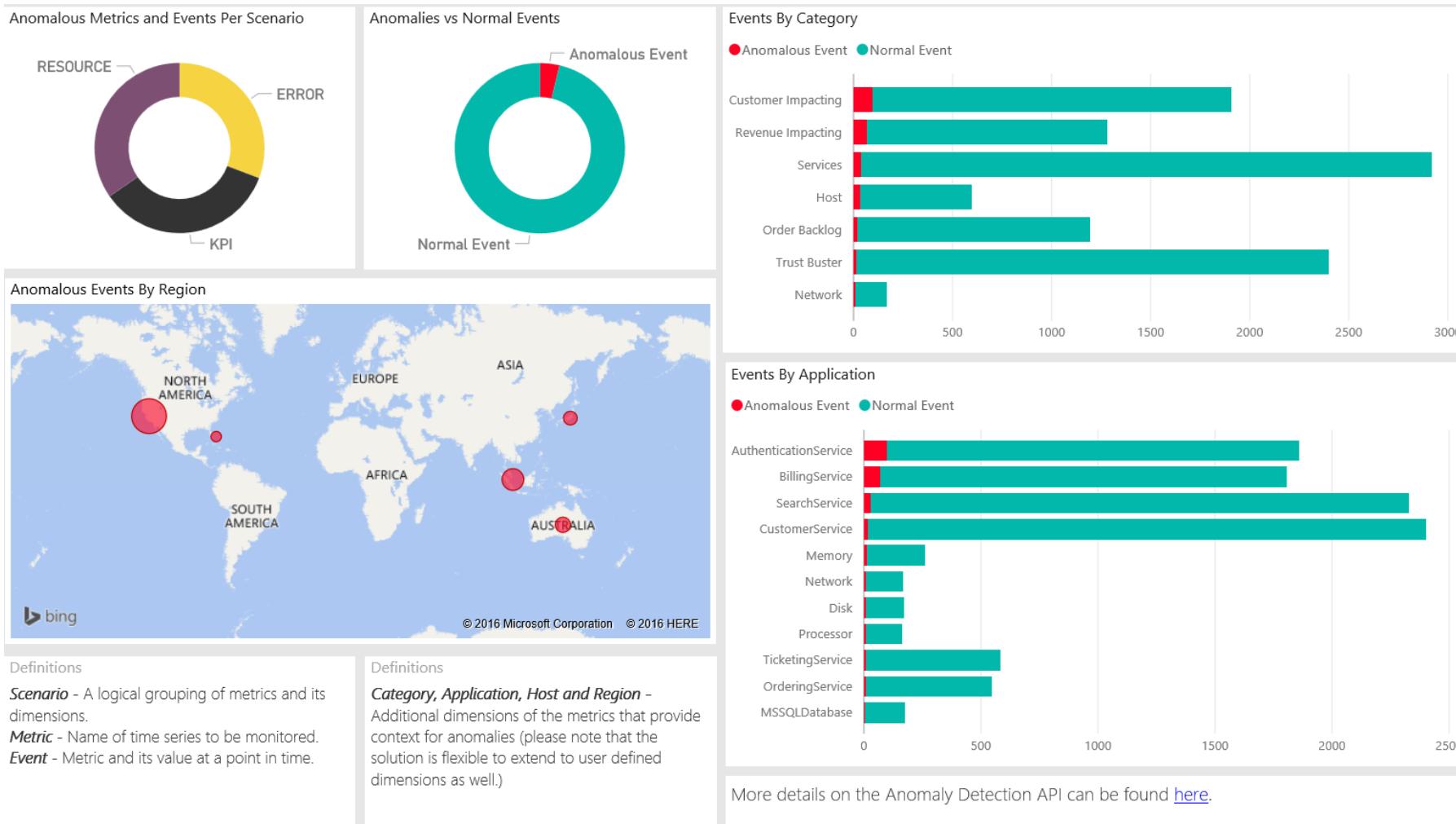
Solution Diagram



1. Business Understanding

Data Science process – TDSP

Business Insights



2. Data Acquisition and Understanding

Data Science process – TDSP

Acquire and understand data

- Ingest the data
 - Move the data from source to target location
- Explore the data
 - Data cleaning
 - Incomplete
 - Noisy
 - Inconsistent
 - Review the quality of data
 - Have better understand the patterns
 - Iterative works
- Set up a data pipeline
 - Score new data
 - Refresh the data regularly
 - Pipeline may be batch-based or a streaming or hybrid
- Artifacts
 - Data quality report
 - Solution architecture
 - Checkpoint decision

Exploratory Data Analysis

Demo

3. Modeling

Data Science process – TDSP

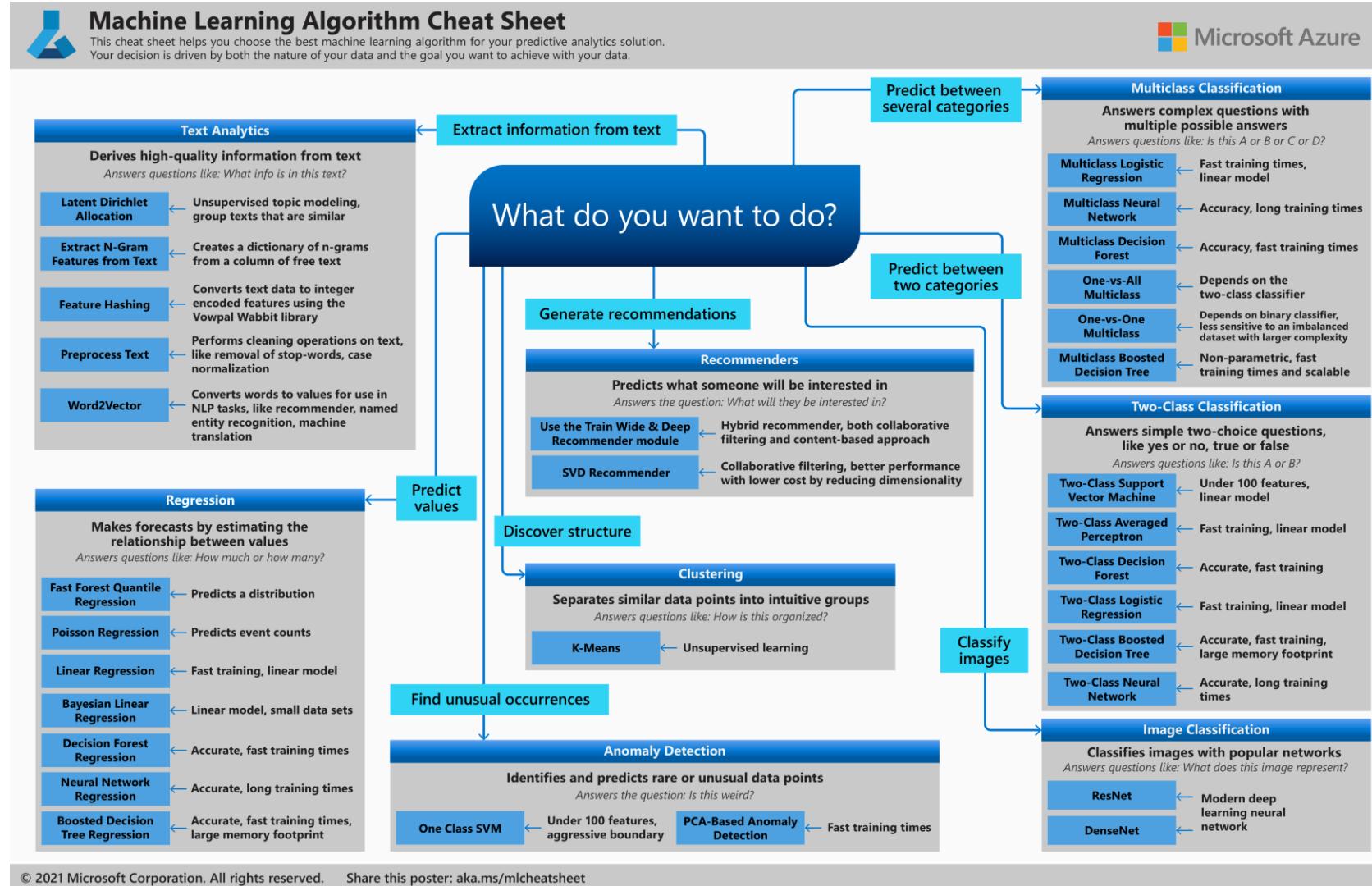
Train models

- Goals
 - Optimal data features for the machine learning model
 - Build informative machine learning model that predicts the target most accurately
 - Experiment for machine learning model that is suitable for production
- How to do it
 - Feature engineering
 - Aggregate and transform the raw variables to create the features
 - Requires a creative combination of domain expertise and insights
 - Model training
 - How to choose algorithms
 - Split, Build, Evaluate, Determine the best solution
- Artifacts
 - Feature sets
 - Modeling report
 - Checkpoint decision

3. Modeling

How to choose algorithms

Data Science process – TDSP



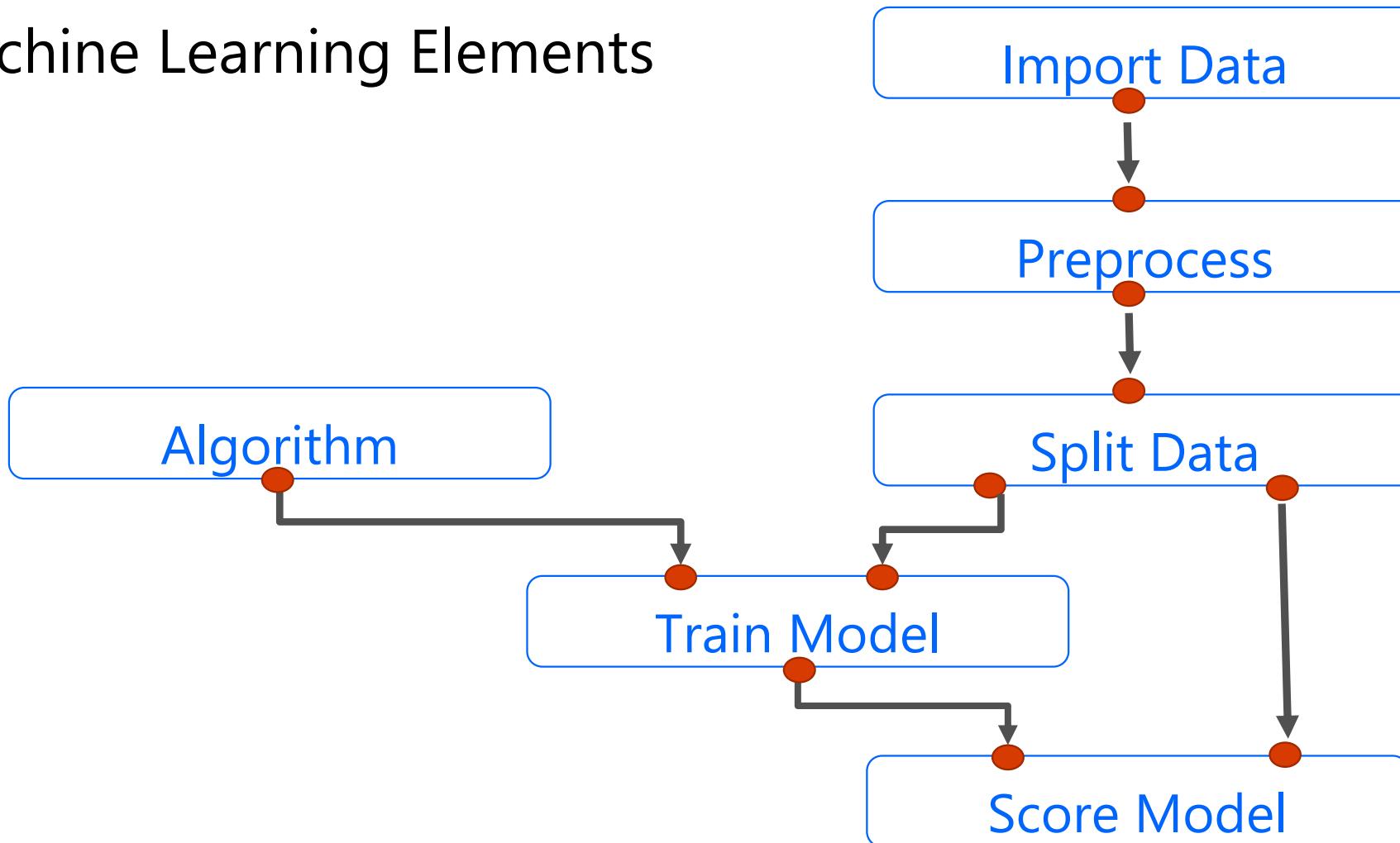
Microsoft Technology Centers

3. Modeling

Data Science process – TDSP

Model Training

- Basic Machine Learning Elements



4. Deployment

Data Science process – TDSP

Operationalize models

- Operationalize the model
 - Apply trained model to application for real-time or batch basis predictions
 - Expose the model as API interface to consume it easy
- Artifacts
 - Status dashboard of system health and key metrics
 - Final modeling report with deployment details
 - Final solution architecture document

5. Customer Acceptance

Data Science process – TDSP

Finalize the project deliverables

- System validation
 - Confirm the deployed model and pipeline are meeting customer needs
 - Monitor systems
- Project hand-off
- Artifacts
 - Project final report

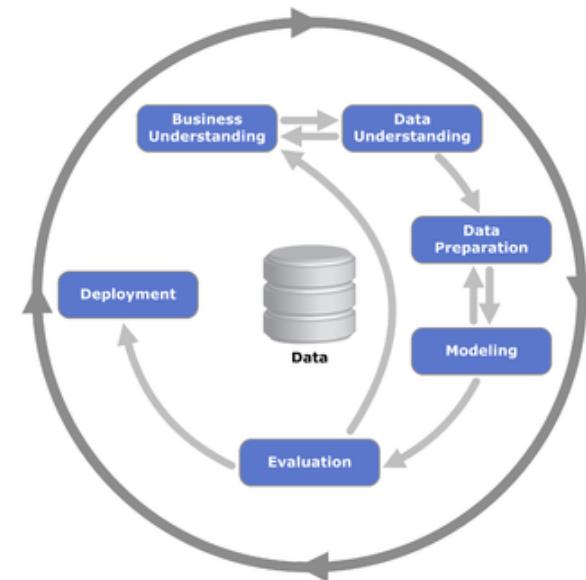
Data Science process – TDSP

Data Science process – TDSP

Conclusion

Team Data Science Process Lifecycle

- Increase the chance of a successful completion of a complex data science project
- Continue to move a data science project forward toward a clear engagement end point
- Modeled as a sequence of iterated steps



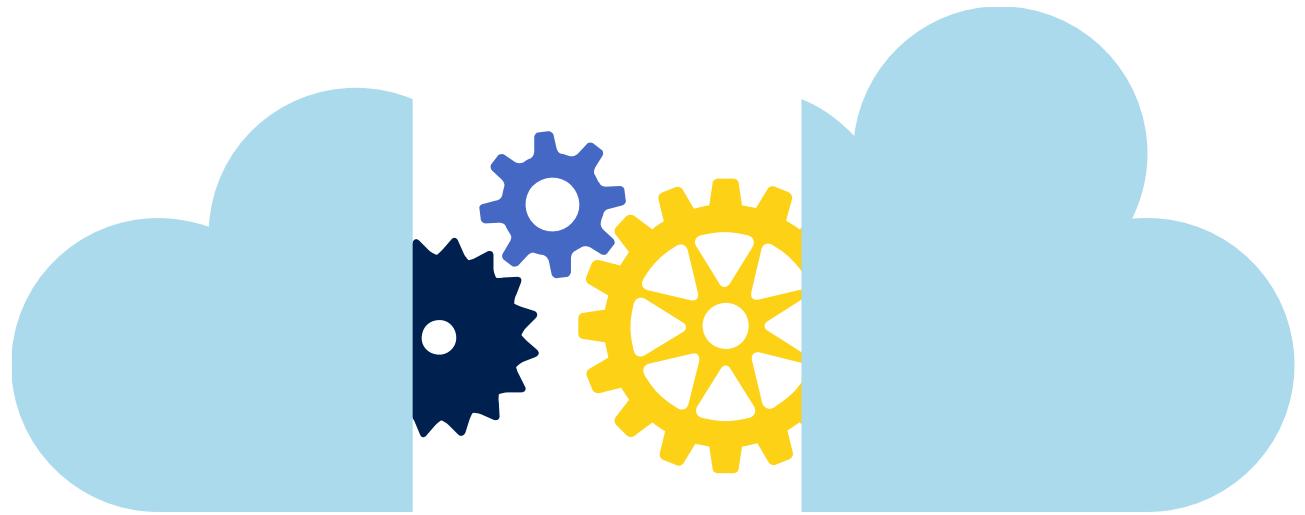
Coffee break

10 minutes



Agenda

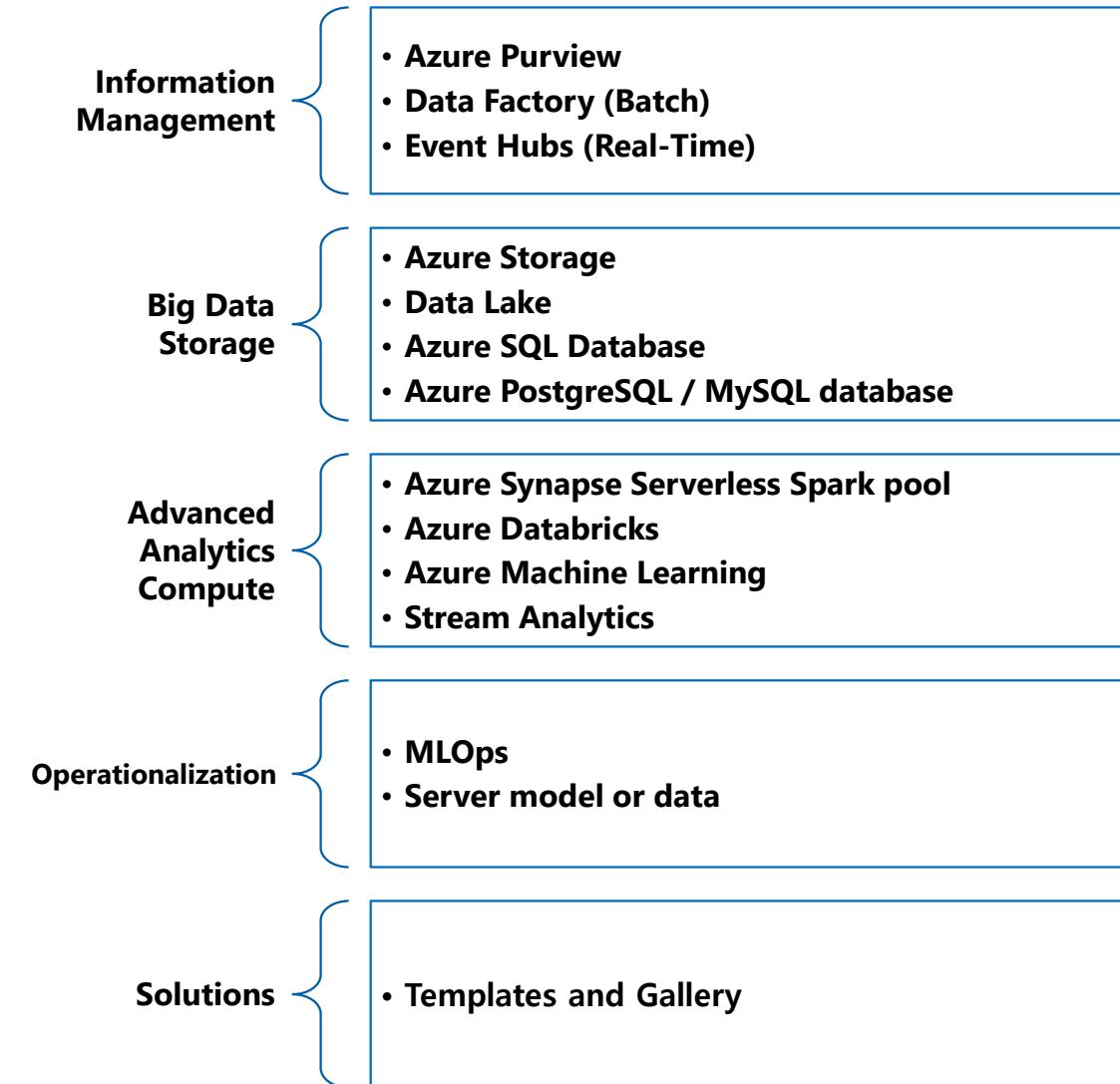
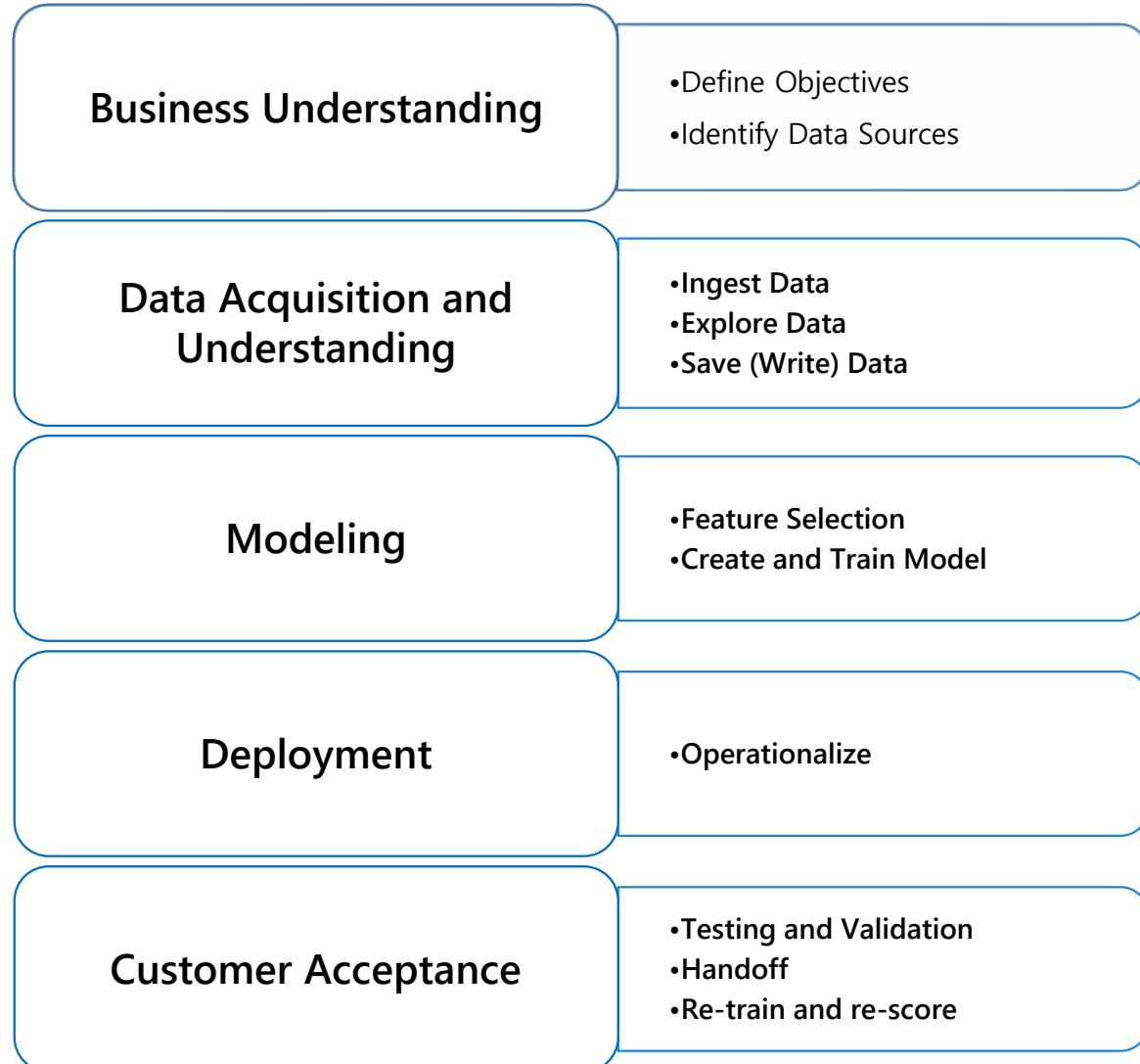
- **Data Science**
- **Data Science process – TDSP (Team Data Science Process)**
- **Data Science Tools**



Microsoft Technology Centers

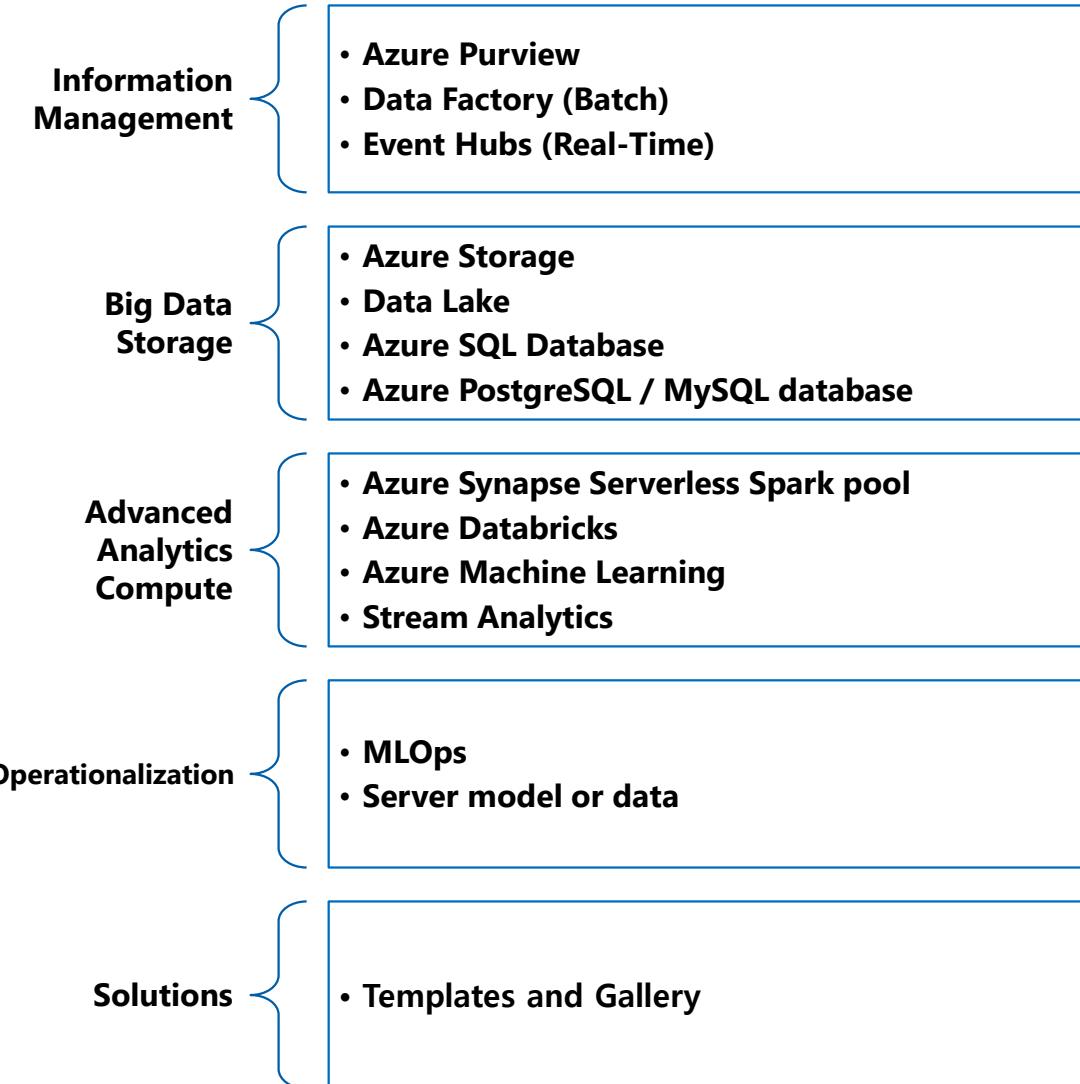
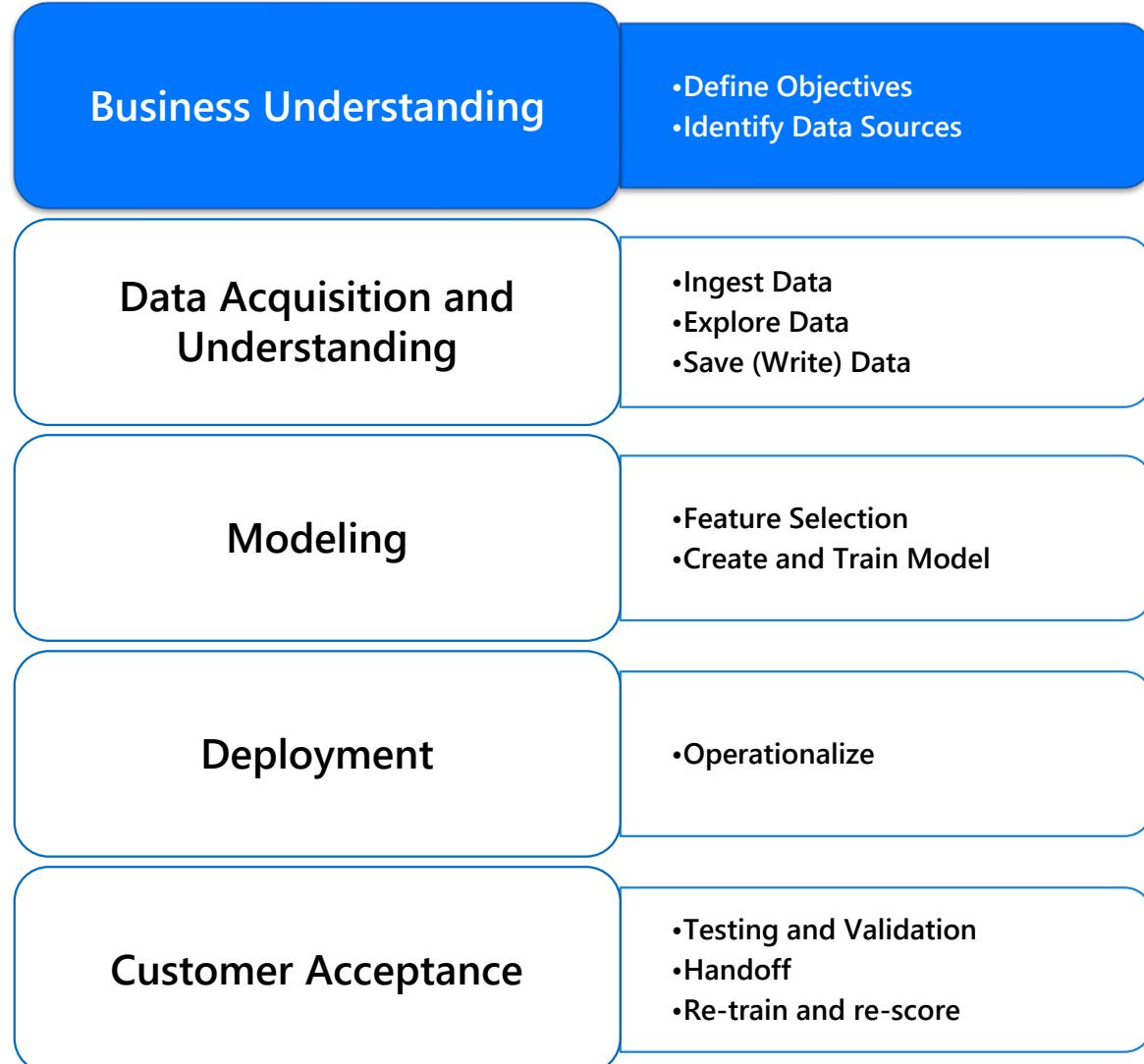
Data Science Tools

Data Science Tools



Data Science Tools

Data Science Tools



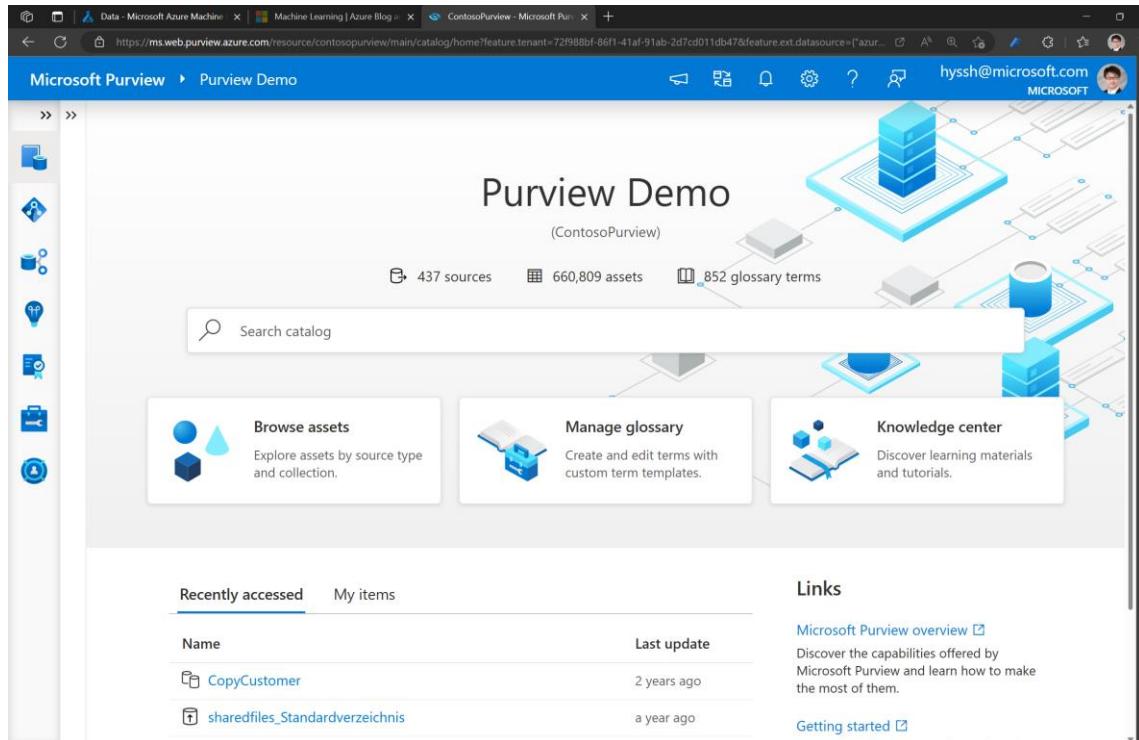
Microsoft Technology Centers

1. Business Understanding

Data Science Tools

Discover data

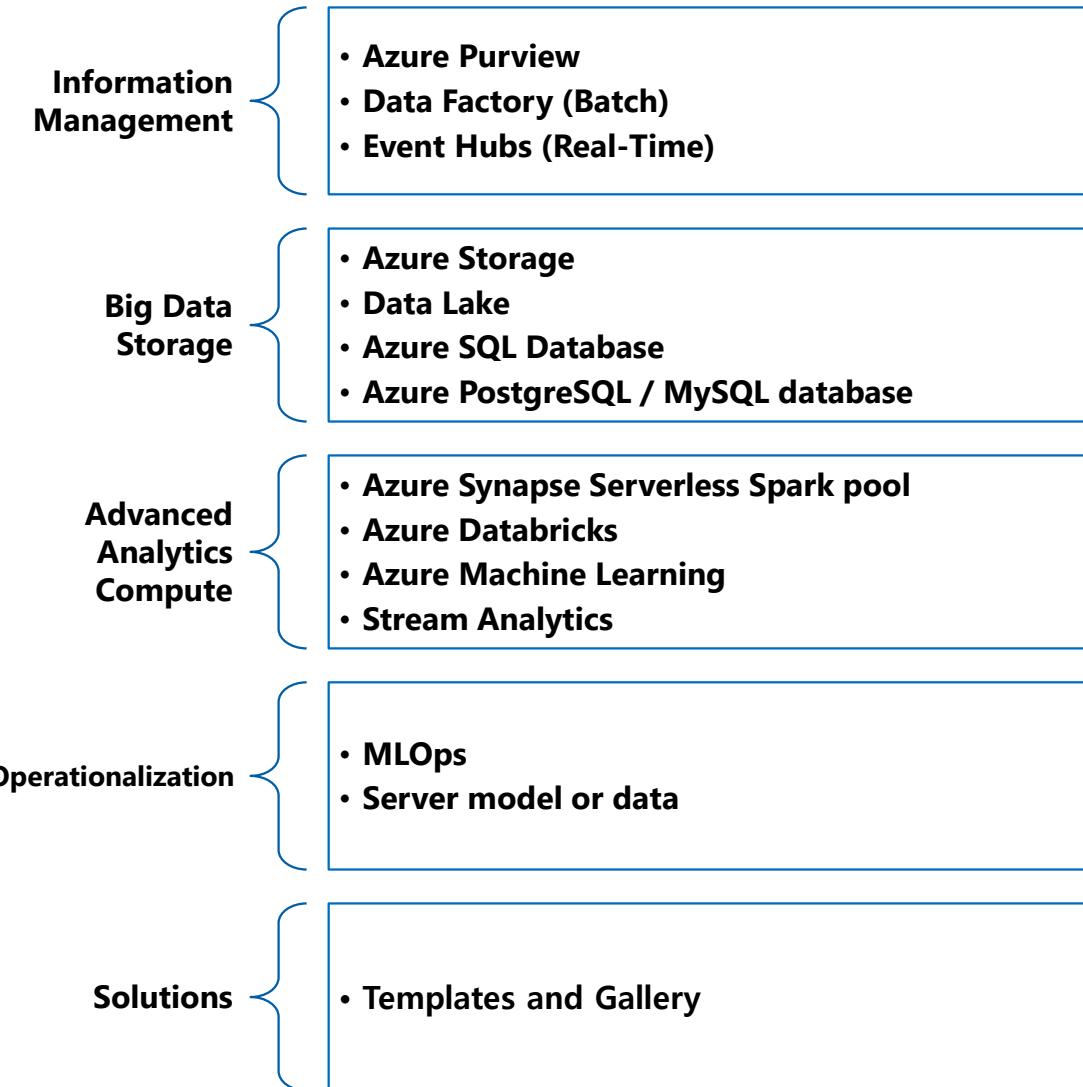
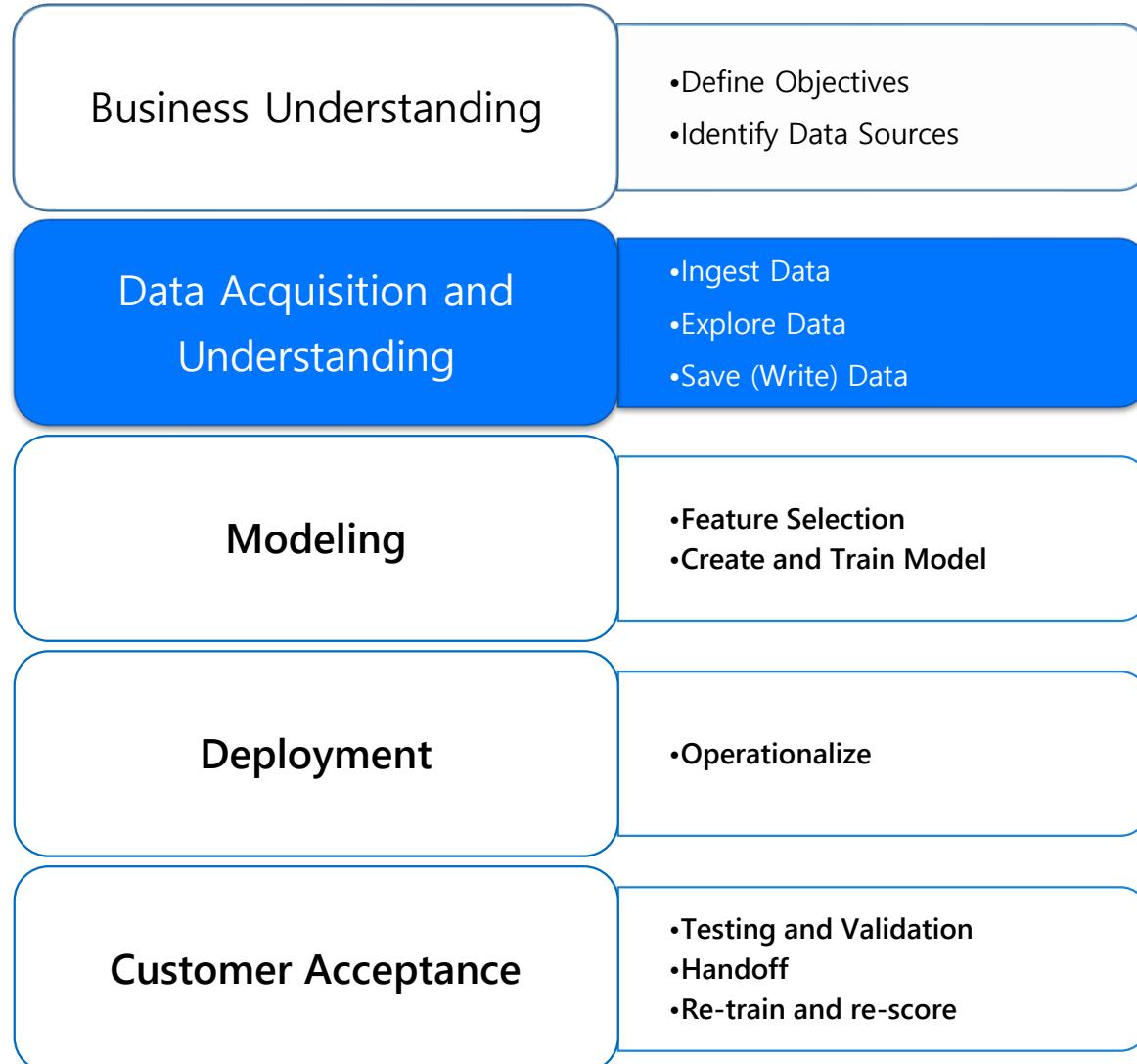
- Find data sources using Azure Purview



- Spend less time looking for data, and more time getting value from it
- Register enterprise data sources, discover data assets and unlock their potential, and capture tribal knowledge to make data understandable
- Bridge the gap between IT and the business, allowing everyone to contribute their insights, tags, and descriptions
- Intuitive search and filtering to understand the data sources and their purpose
- Let your data live where you want; connect using tools you choose
- Integrate into existing tools and processes with open REST APIs

Data Science Tools

Data Science Tools

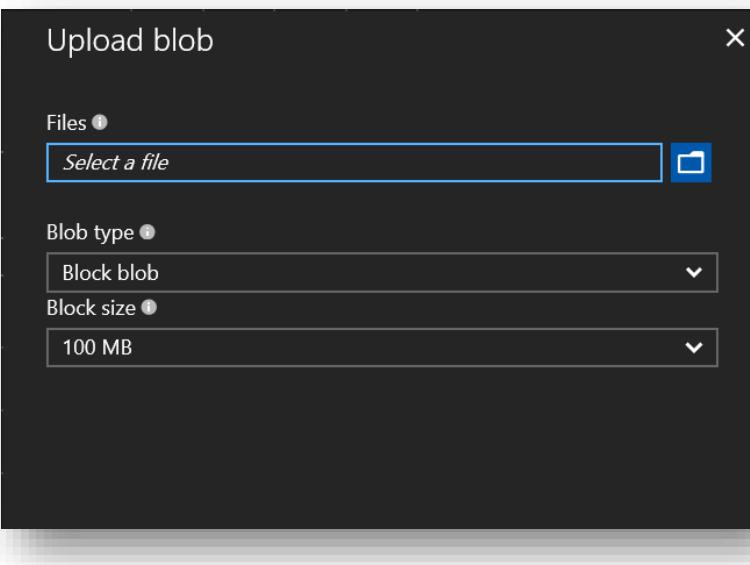


2. Data Acquisition and Understanding

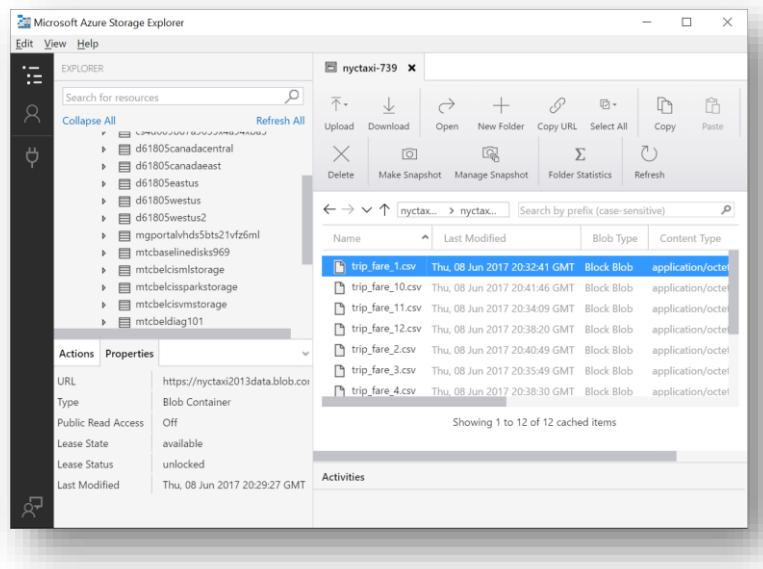
Data Science Tools

Acquire and understand data

- Azure portal
- PowerShell
- Azure Data Factory
- Azure Event Hubs
- Azure storage SDKs (.NET, Node.js, python, C++, etc.)
- AzCopy (blob, file, and table only)
- Import/Export service



[Azure Portal]



[Azure Storage Explorer]

The screenshot shows a Windows Command Prompt window titled 'cmd.exe' with the path 'C:\windows\system32\cmd.exe'. The command entered is 'C:\Program Files\Microsoft AzCopy>AZCOPY C:\myvhds\ http://hkstorage.blob.core.windows.net/myvhds /destKey:SivqRonJU/L9fhboXx1MnwCKyhlrae2o/Zt5HTyrsu1gMfuHo5h08uxVBN219xA8cIILxEKpyqx/x7pFo7dpn4g== /S /V /blobtype:page'. The output shows 'Finished transfer: demovhd.vhd'. Below that is a 'Transfer summary:' section with statistics: 'Total files transferred: 1', 'Transfer successfully: 1', and 'Transfer failed: 0'.

```
C:\Program Files\Microsoft AzCopy>AZCOPY C:\myvhds\ http://hkstorage.blob.core.windows.net/myvhds /destKey:SivqRonJU/L9fhboXx1MnwCKyhlrae2o/Zt5HTyrsu1gMfuHo5h08uxVBN219xA8cIILxEKpyqx/x7pFo7dpn4g== /S /V /blobtype:page
Finished transfer: demovhd.vhd

Transfer summary:
-----
Total files transferred: 1
Transfer successfully: 1
Transfer failed: 0
```

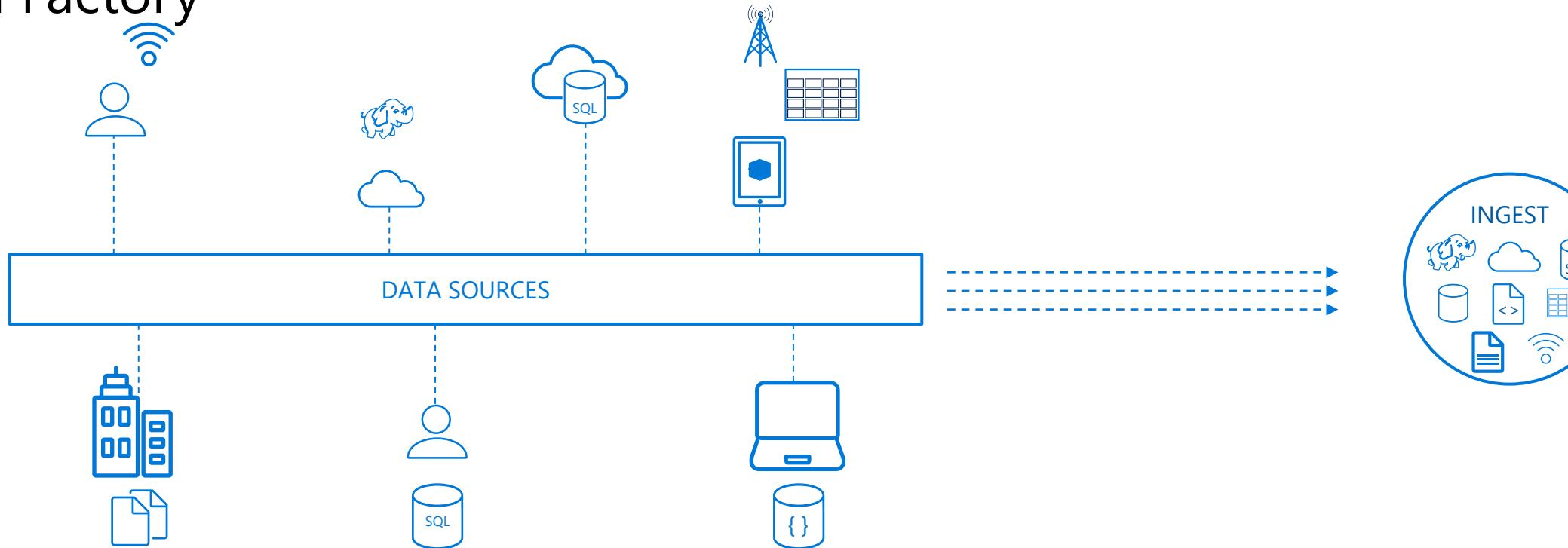
[AZcopy]

2. Data Acquisition and Understanding

Data Science Tools

Ingest Data

- Azure Data Factory



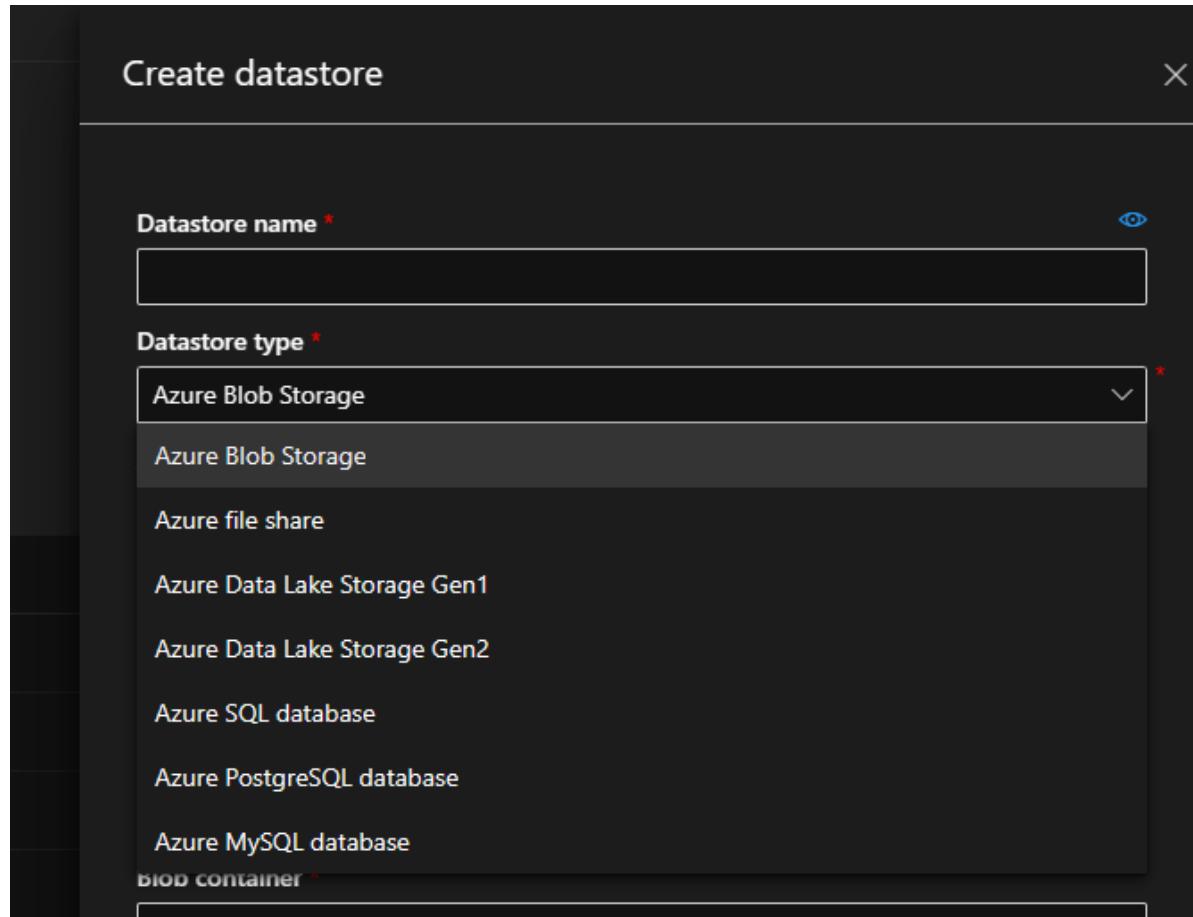
- Create, schedule, orchestrate, and manage data pipelines
- Visualize data lineage
- Connect to on-premises and cloud data sources
- Monitor data pipeline health
- Automate cloud resource management
- Move relational data for Hadoop processing
- Transform with Hive, Pig, or custom code

2. Data Acquisition and Understanding

Data Science Tools

Link data to Azure Machine Learning Studio

- Datastore

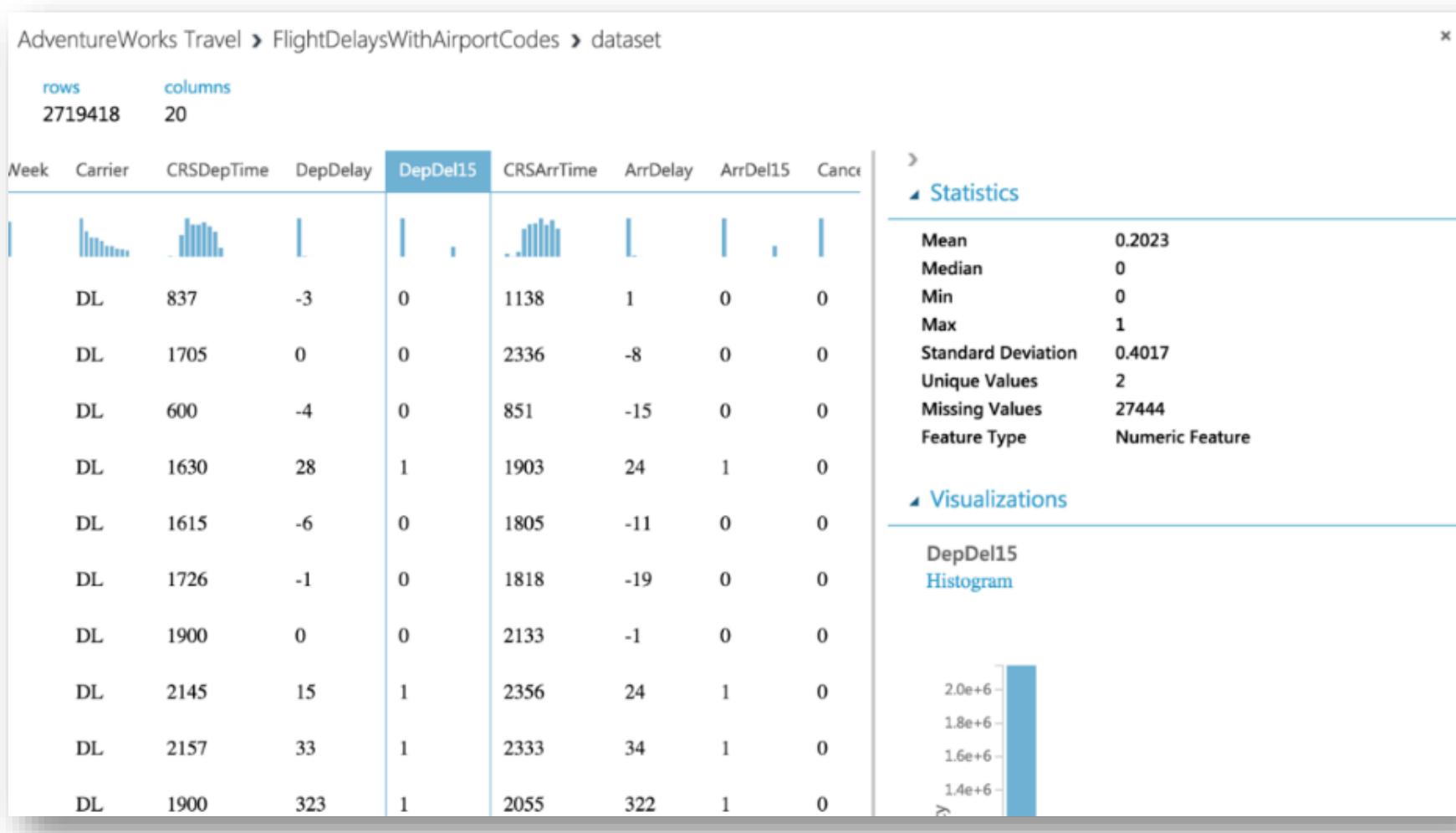


2. Data Acquisition and Understanding

Data Science Tools

Tools to understand data

- Data Asset in Azure Machine Learning Studio

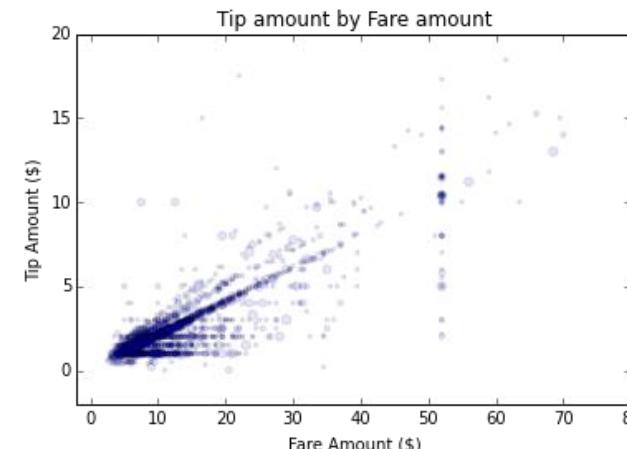
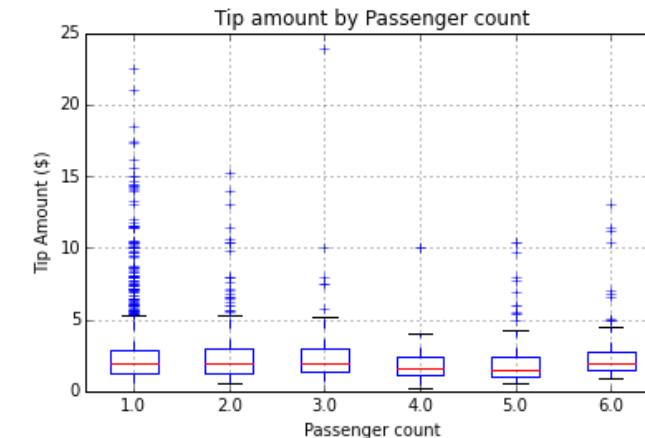
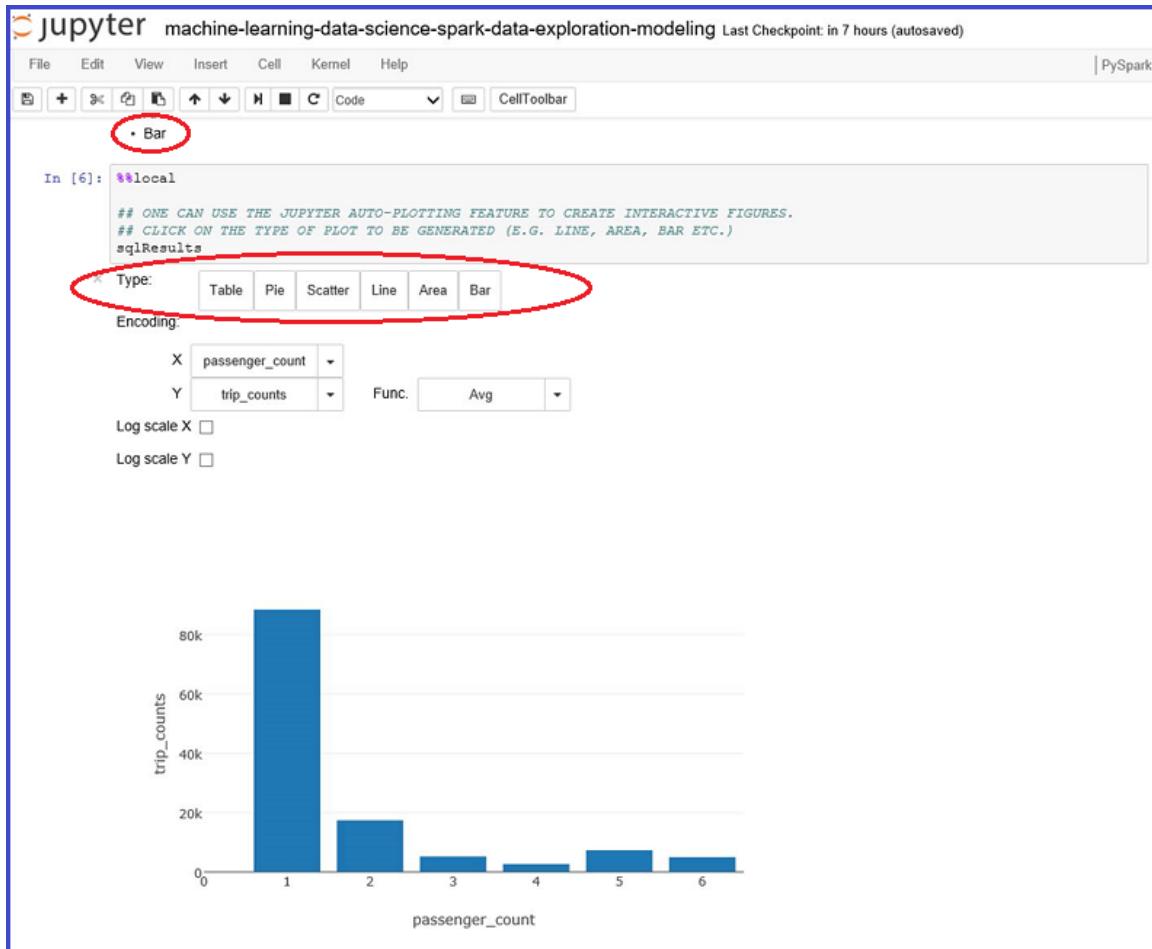


2. Data Acquisition and Understanding

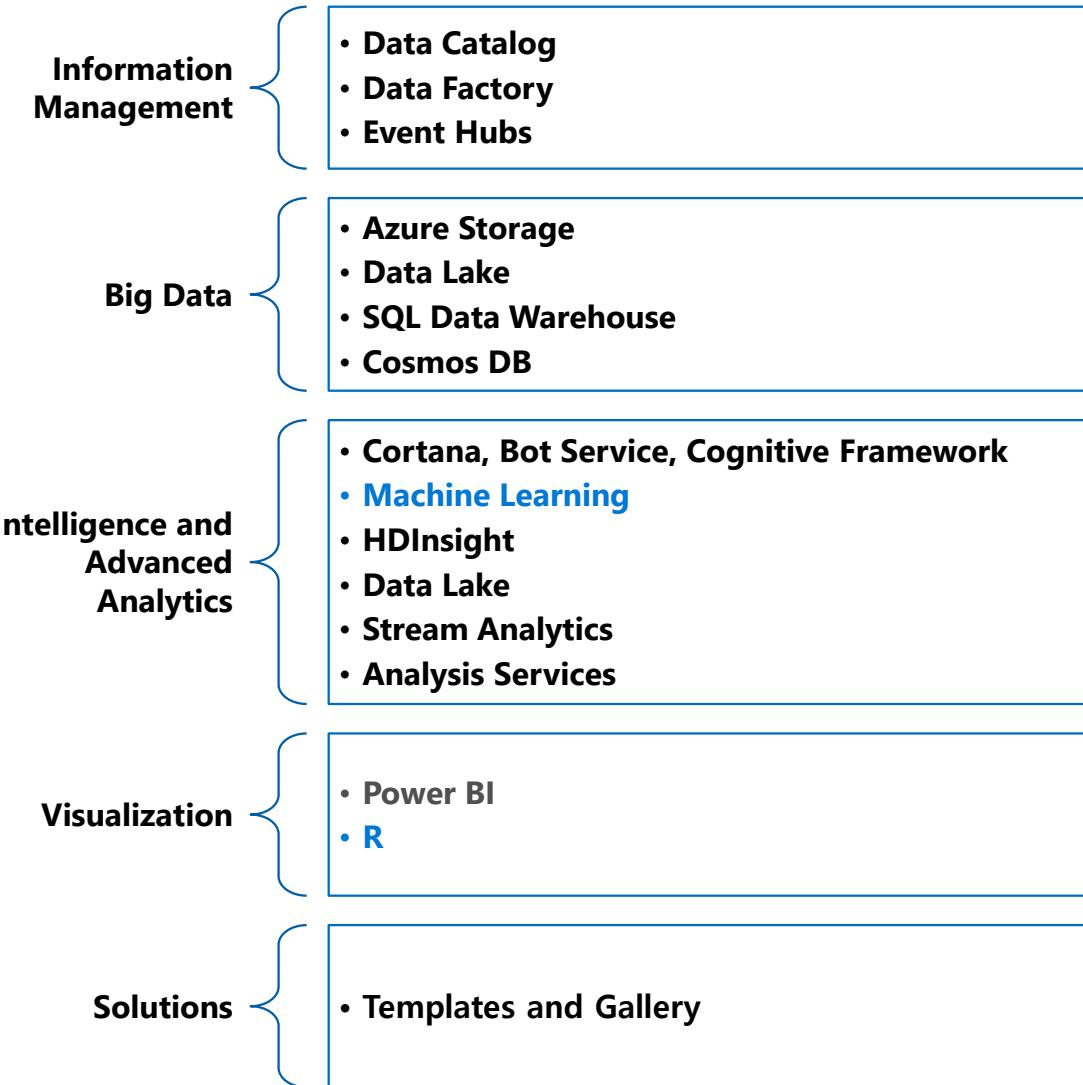
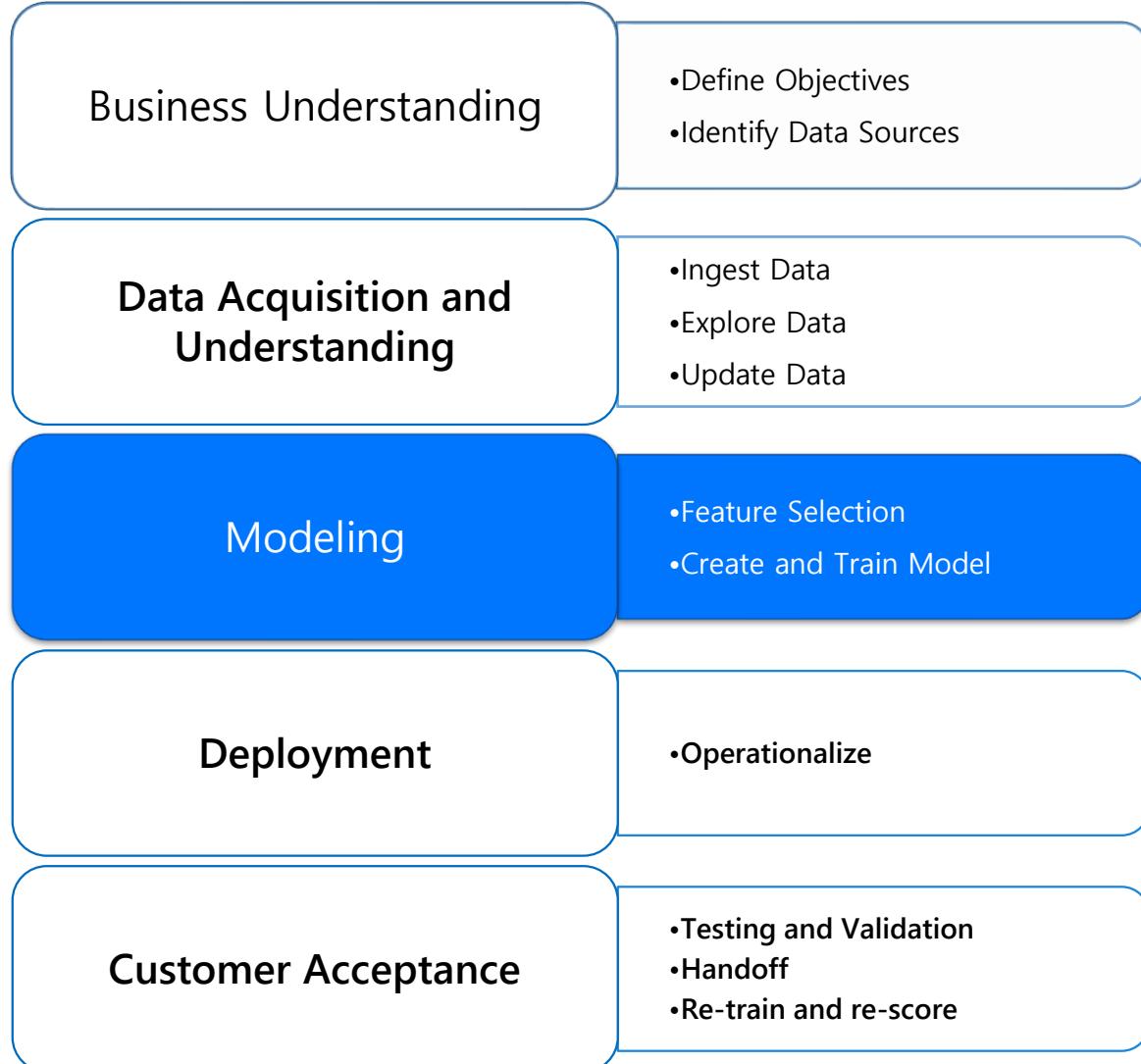
Data Science Tools

Tools to understand data

- Use Seaborn, Matplotlib in Python



Data Science Tools

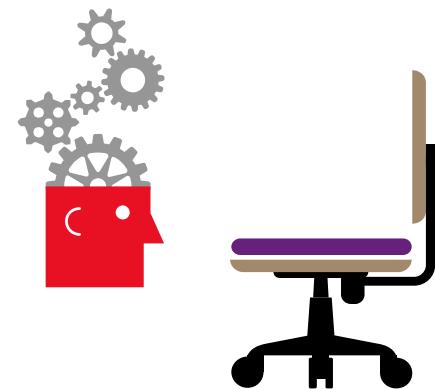
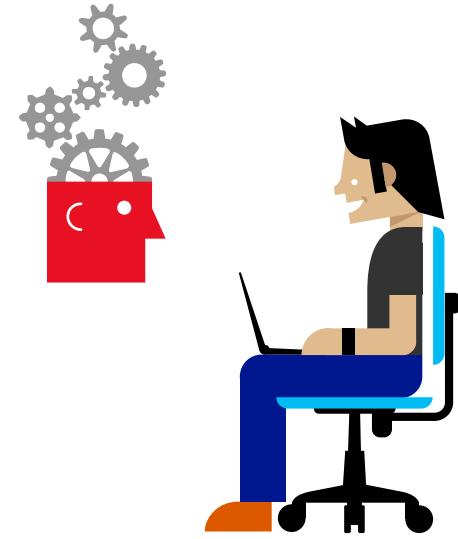


3. Modeling

Data Science Tools

Machine Learning Algorithms

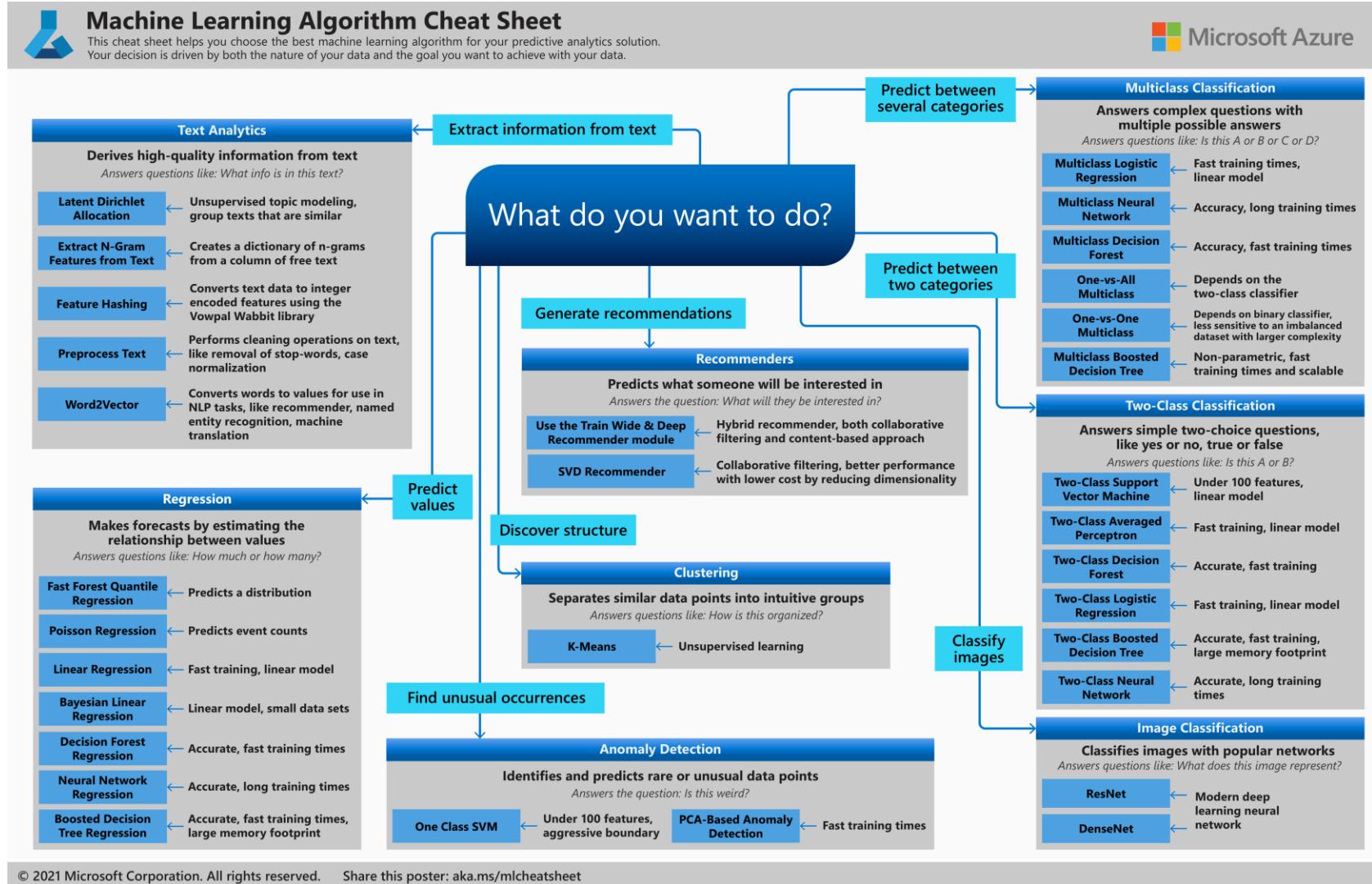
- Split into two main categories
 - Supervised learning
 - Predicting the future
 - Learn from known past examples to predict future
 - Labels provided
 - Unsupervised learning
 - Making sense of data
 - Understanding the past
 - Learning the structure of data
 - Labels no provided



3. Modeling

Machine Learning Algorithms

Data Science Tools

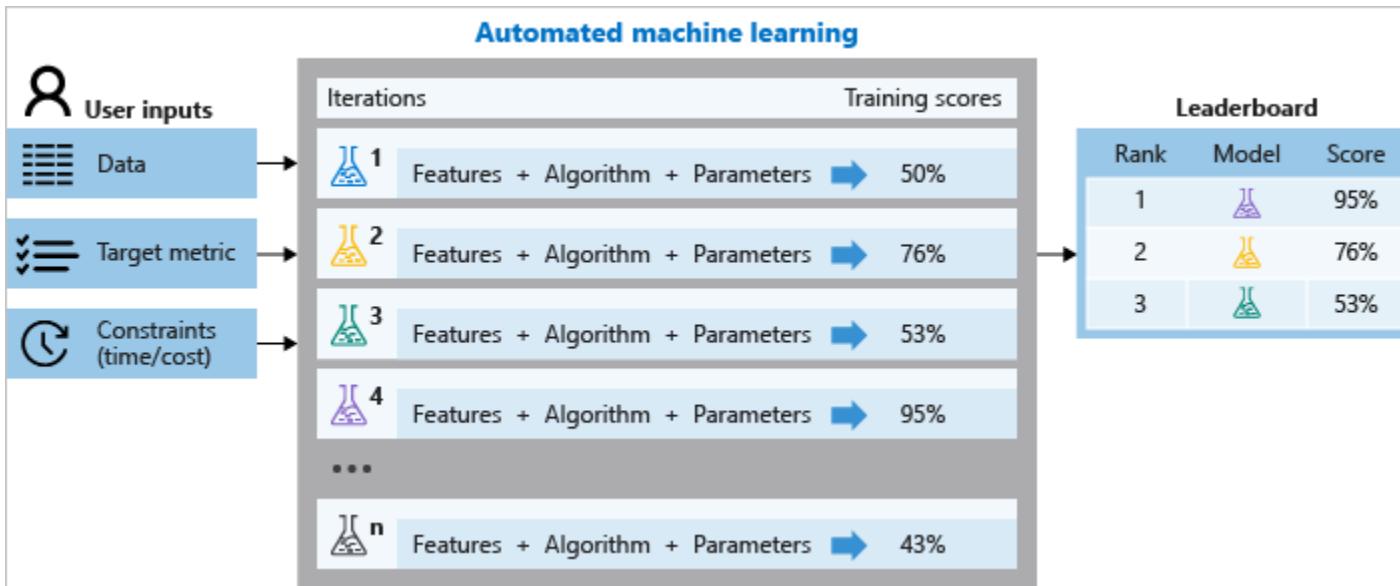


3. Modeling

Data Science Tools

Model training

- AutoML



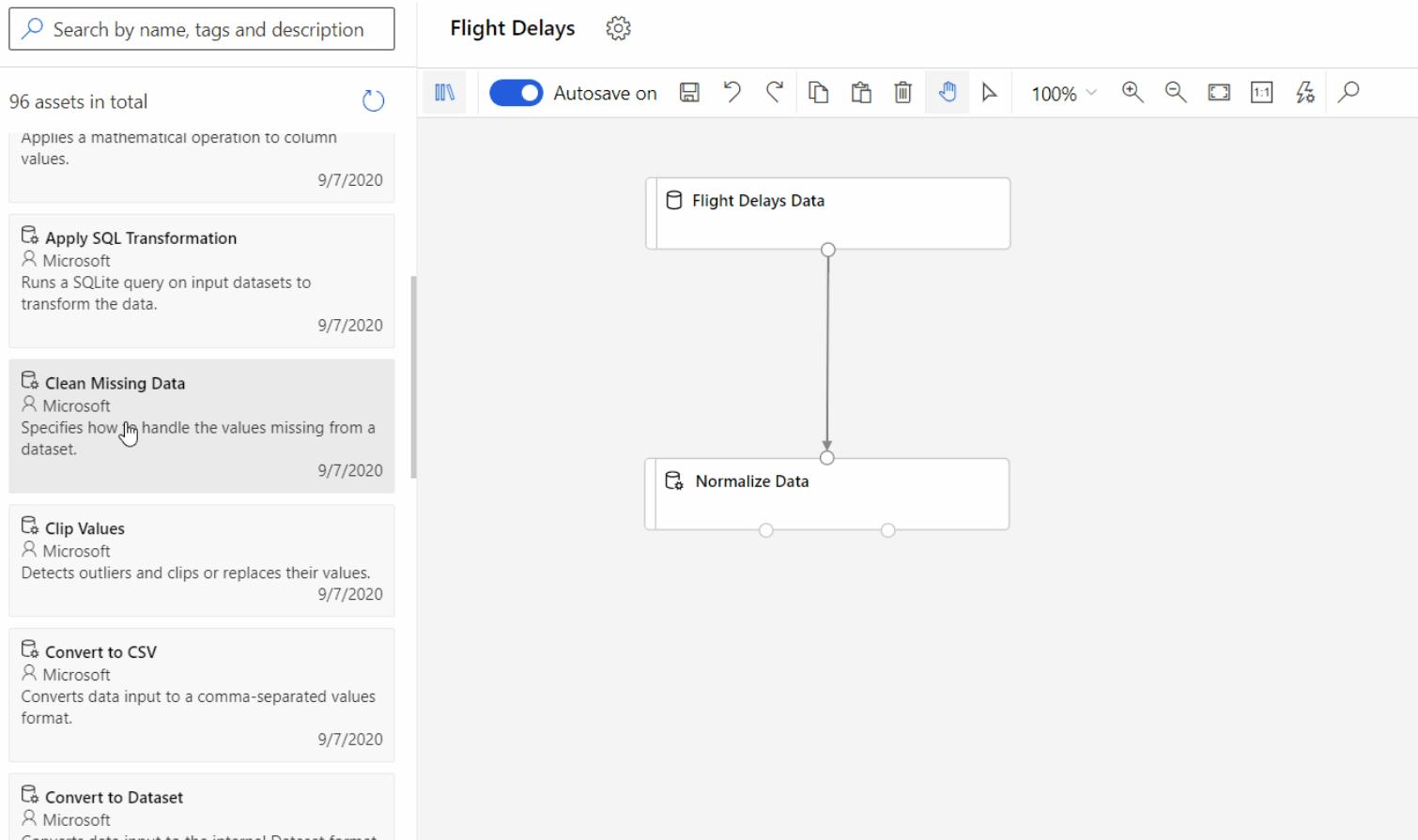
Automated machine learning, also referred to as automated ML or AutoML, is the process of automating the time-consuming, iterative tasks of machine learning model development. It allows data scientists, analysts, and developers to build ML models with high scale, efficiency, and productivity all while sustaining model quality. Automated ML in Azure Machine Learning is based on a breakthrough from our [Microsoft Research](#) division.

3. Modeling

Model training

- Designer

Data Science Tools



3. Modeling

Data Science Tools

Model training

- Notebook (Code First)

Microsoft Azure Machine Learning Studio

Search within your workspace

Microsoft > sdg-ws > Notebooks

Microsoft

New

Home

Author

Notebooks

Automated ML

Designer

Assets

Notebooks

Microsoft > sdg-ws > Notebooks

Notebooks

Files Samples

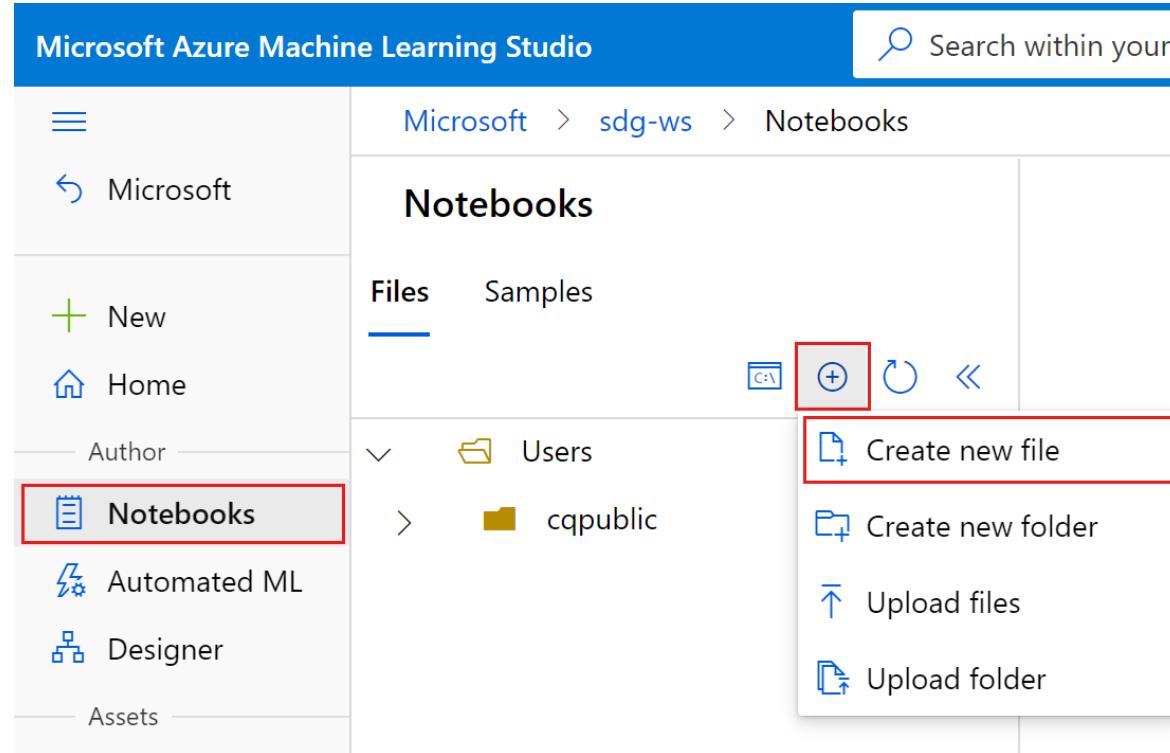
Users cpublic

Create new file

Create new folder

Upload files

Upload folder



Testing a new notebook

Use markdown cells to add nicely formatted content to the notebook.

[2]

```
1 print("Hello, world!")  
✓ <1 sec
```

Hello, world!

[3]

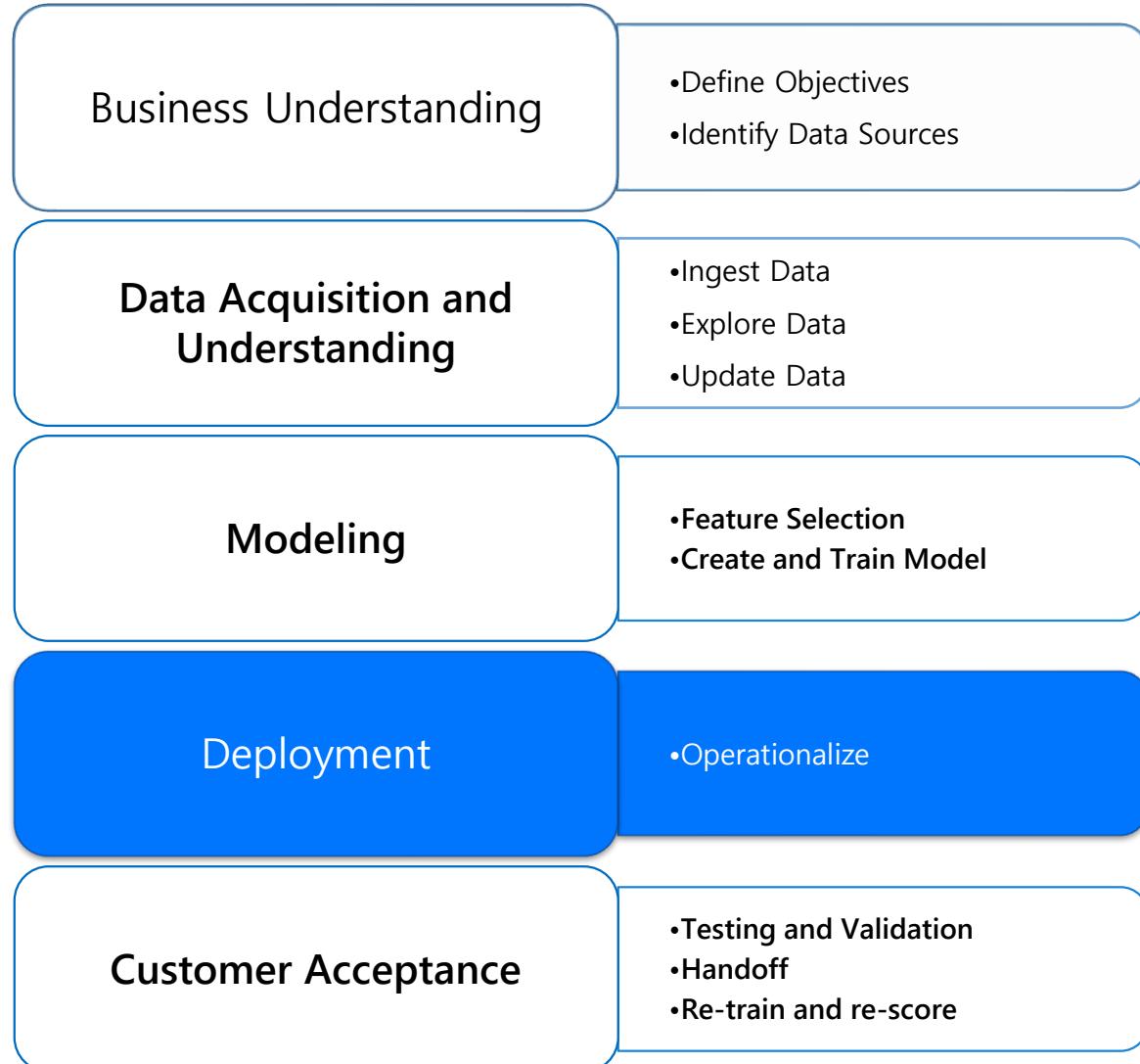
```
1 two = 1 + 1  
2 print("One plus one is ",two)  
✓ <1 sec
```

... One plus one is 2

+ Code + Markdown



Data Science Tools



4. Deployment

Data Science Tools

Operationalize models

- Azure Machine Learning - Publish as a Web Service

Microsoft Azure Machine Learning Studio

Search This workspace

movie-recommender

Details Test Consume Deployment logs

Basic consumption info

REST endpoint
http://465592e1-56d6-4c98-b0ff-0d3c1490c21a.westus2.azurecontainer.io/score

Authentication

Primary key
..... Regenerate

Secondary key
..... Regenerate

Consumption option

Consumption types

C# Python R

```
40
41 body = str.encode(json.dumps(data))
42
43 url = 'http://465592e1-56d6-4c98-b0ff-0d3c1490c21a.westus2.azurecontainer.io/score'
44 # Replace this with the primary/secondary key or AMLToken for the endpoint
45 api_key = ''
46 if not api_key:
47     raise Exception("A key should be provided to invoke the endpoint")
48
49
50 headers = {'Content-Type': 'application/json', 'Authorization':('Bearer ' + api_key)}
51 req = urllib.request.Request(url, body, headers)
52
53 try:
54     response = urllib.request.urlopen(req)
55
56     result = response.read()
57     print(result)
58 except urllib.error.HTTPError as error:
59     print("The request failed with status code: " + str(error.code))
60
61
62 # Print the headers - they include the request ID and the timestamp, which are useful for debugging the
63 print(error.info())
64 print(error.read().decode("utf-8"))
65 print(error.headers)
```

Microsoft Azure Machine Learning Studio

Search This workspace

movie-recommender

Details Test Consume Deployment logs

Input data to test real-time endpoint

Test

Test result

```
"Inputs": {
    "Input1": [
        {
            "UserId": 0,
            "Movie Name": "Star Wars (1977)",
            "Rating": 5
        },
        {
            "UserId": 290,
            "Movie Name": "Star Wars (1977)",
            "Rating": 5
        },
        {
            "UserId": 79,
            "Movie Name": "Star Wars (1977)",
            "Rating": 4
        }
    ],
    "GlobalParameters": {}
}
```

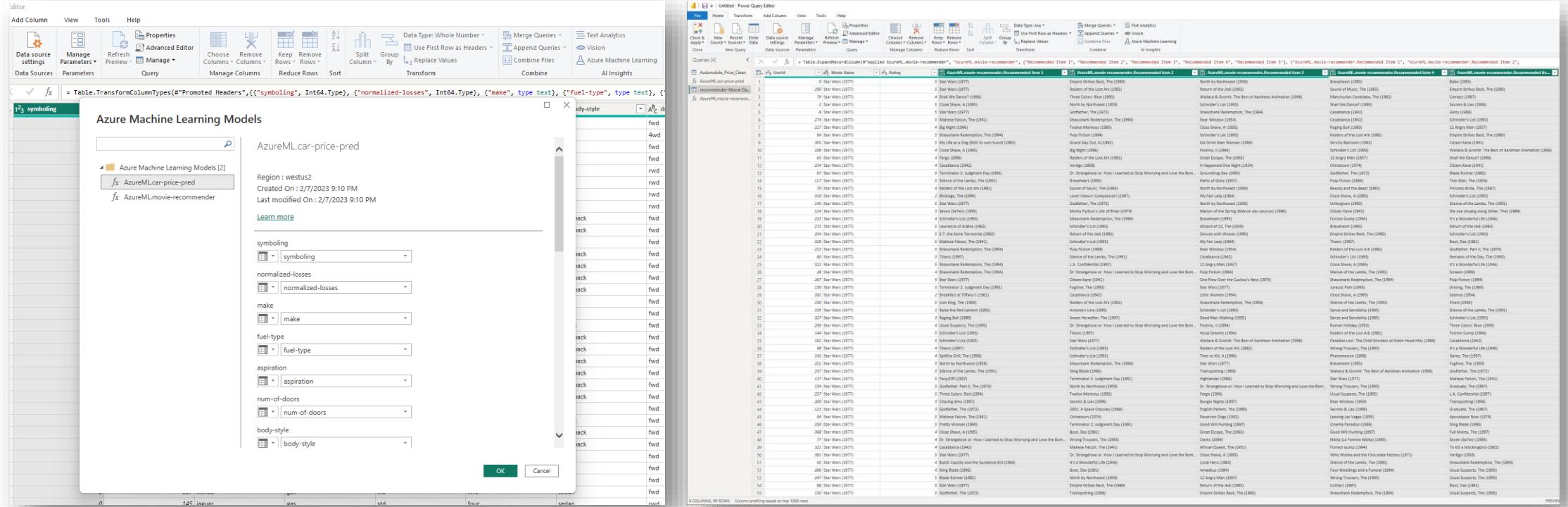
```
{
    "Results": {
        "WebServiceOutput0": [
            {
                "User": "0",
                "Recommended Item 1": "Star Wars (1977)",
                "Recommended Item 2": "Empire Strikes Back, The (1980)",
                "Recommended Item 3": "North by Northwest (1959)",
                "Recommended Item 4": "Braveheart (1995)",
                "Recommended Item 5": "Babe (1995)"
            },
            {
                "User": "290",
                "Recommended Item 1": "Star Wars (1977)",
                "Recommended Item 2": "Raiders of the Lost Ark (1981)",
                "Recommended Item 3": "Return of the Jedi (1983)",
                "Recommended Item 4": "Sound of Music, The (1965)",
                "Recommended Item 5": "Empire Strikes Back, The (1980)"
            },
            {
                "User": "79",
                "Recommended Item 1": "Shall We Dance? (1996)",
                "Recommended Item 2": "Three Colors: Blue (1993)",
                "Recommended Item 3": "Wallace & Gromit: The Best of Aardman Animations",
                "Recommended Item 4": "Manchurian Candidate, The (1962)",
                "Recommended Item 5": "Contact (1997)"
            }
        ]
    }
}
```

4. Deployment

Operationalize models

- Power BI Integration

Data Science Tools

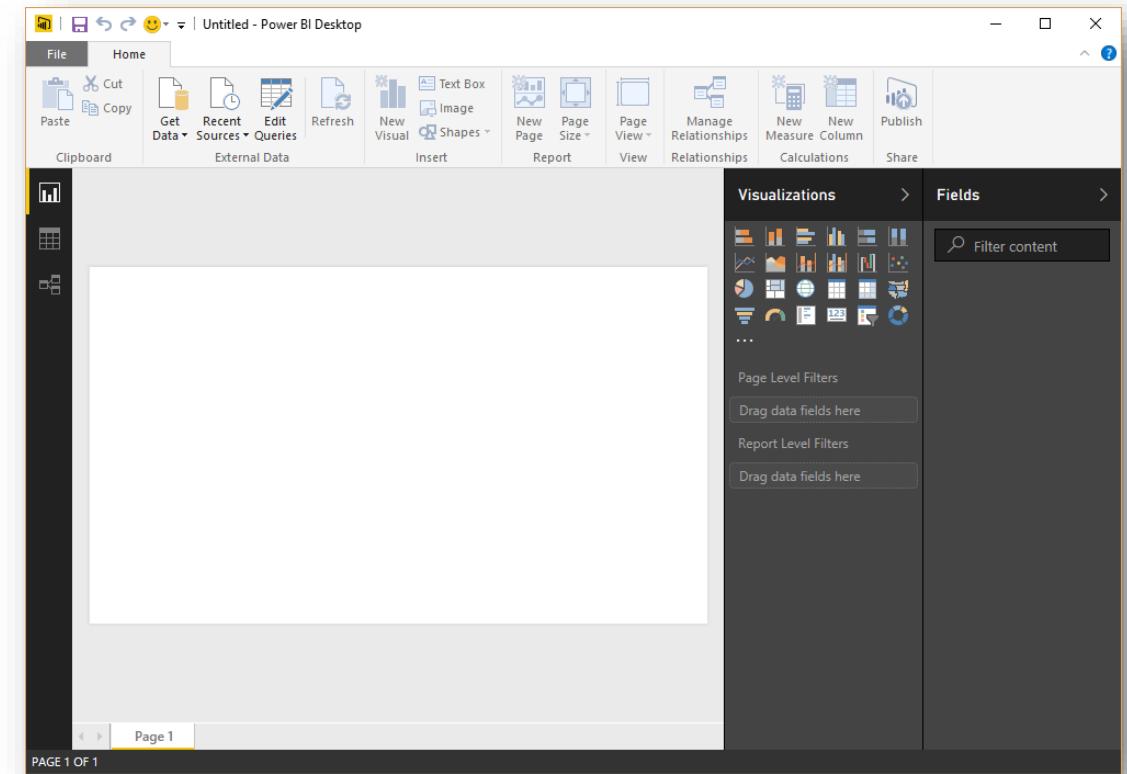


4. Deployment

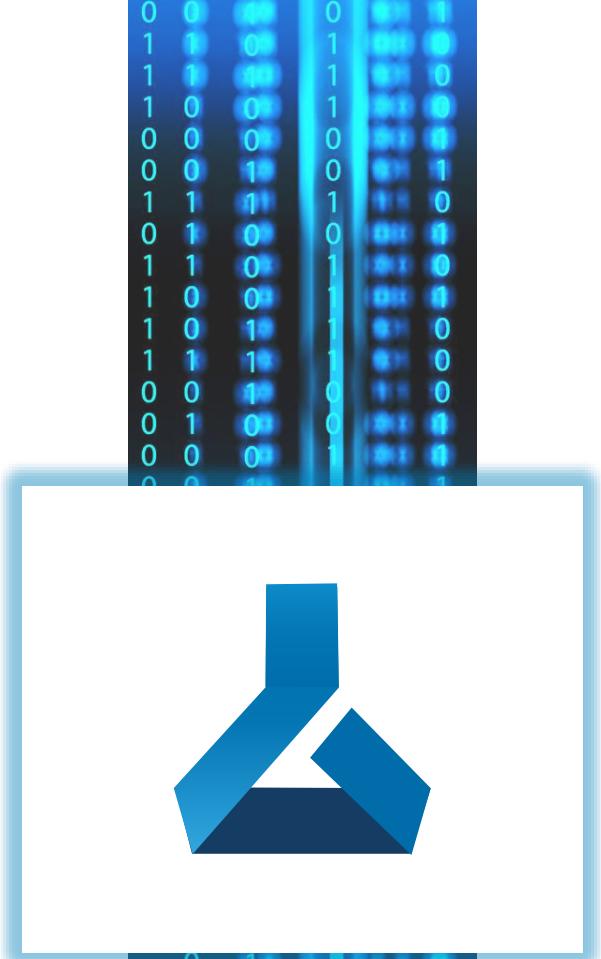
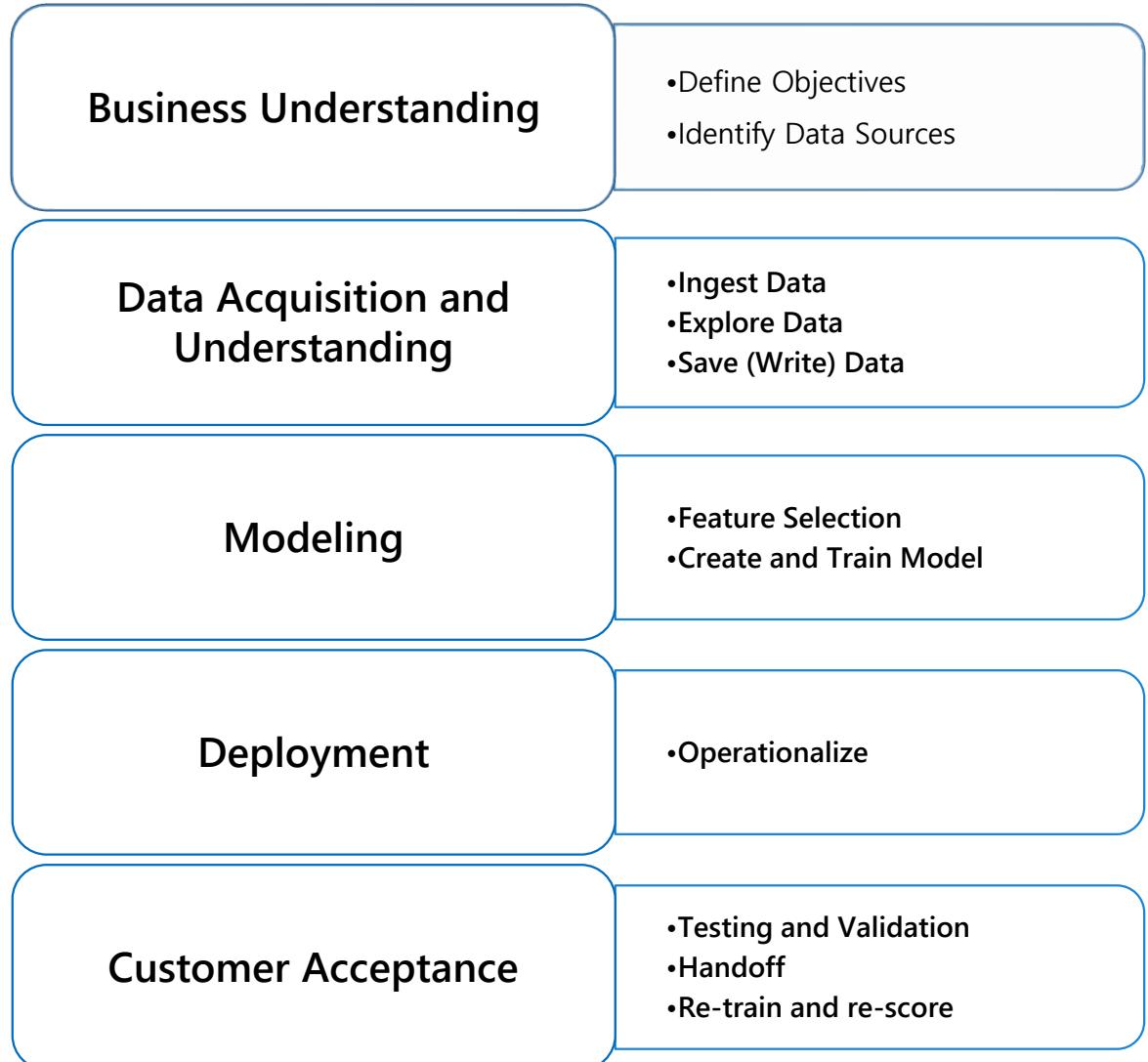
Data Science Tools

Report

- Power BI
 - A Reporting System for Multiple Data Sources
 - Available in:
 - Web Portal
 - Power BI Desktop
 - Microsoft Excel
 - Mobile apps (iOS, Android, Windows)
 - Author
 - Connect to Data
 - Shape the Data
 - Model the Data
 - Report on the Data
 - Publish
 - Local
 - To Service



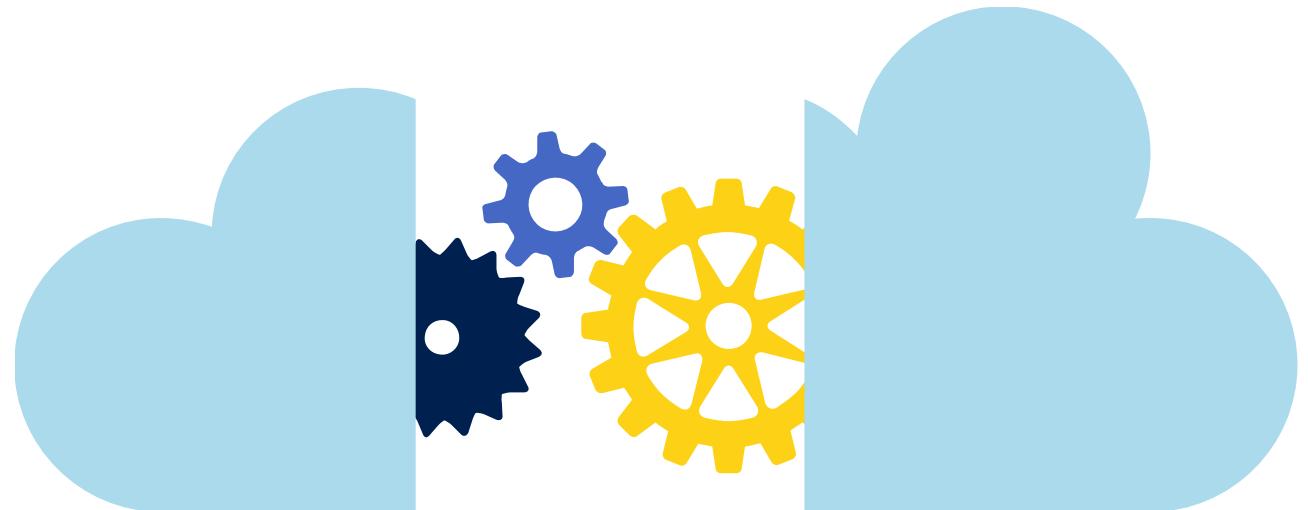
Data Science Tools



Microsoft Technology Centers

Agenda

- **Data Science**
- **Data Science process – TDSP (Team Data Science Process)**
- **Data Science Tools**
- **And next**



Microsoft Technology Centers



How to start Data Science Project at work?



HS

how to start data science project at work?



Starting a data science project at work can be a complex and challenging task, but it can also be very rewarding. Here are some steps to help you get started:

1. Define the problem: Start by identifying a specific problem or business challenge that you aim to solve with your data science project.
2. Gather data: Collect relevant data that can be used to address the problem. This may involve acquiring data from internal sources, purchasing data from external providers, or scraping data from websites.
3. Clean and prepare the data: Clean and pre-process the data to make it ready for analysis. This may involve dealing with missing values, converting data into a format that can be easily analyzed, and normalizing data to remove any biases.
4. Explore the data: Use data visualization and other exploratory data analysis techniques to gain a deeper understanding of the data and identify any patterns or relationships.
5. Develop a model: Based on the insights gained from the exploratory data analysis, develop a model to make predictions or support decision-making.
6. Validate and evaluate the model: Validate the model by comparing its predictions to actual outcomes, and evaluate its performance to determine if it meets the desired accuracy criteria.
7. Deploy the model: Deploy the model into a production environment, where it can be used to make real-time predictions or support decision-making.
8. Monitor and maintain the model: Continuously monitor the model to ensure it is functioning correctly, and make any necessary updates or improvements over time.

These are just general guidelines, and the specific steps involved in starting a data science project will vary depending on the problem being solved, the data being used, and the methods being employed. It's important to have a solid understanding of statistics, machine learning, and computer programming, and to work with a team of skilled data scientists, engineers, and domain experts.



