

# **Data Engineering Project**

## **Module 3**

### **Data Quality Assessment and Data Exploration**

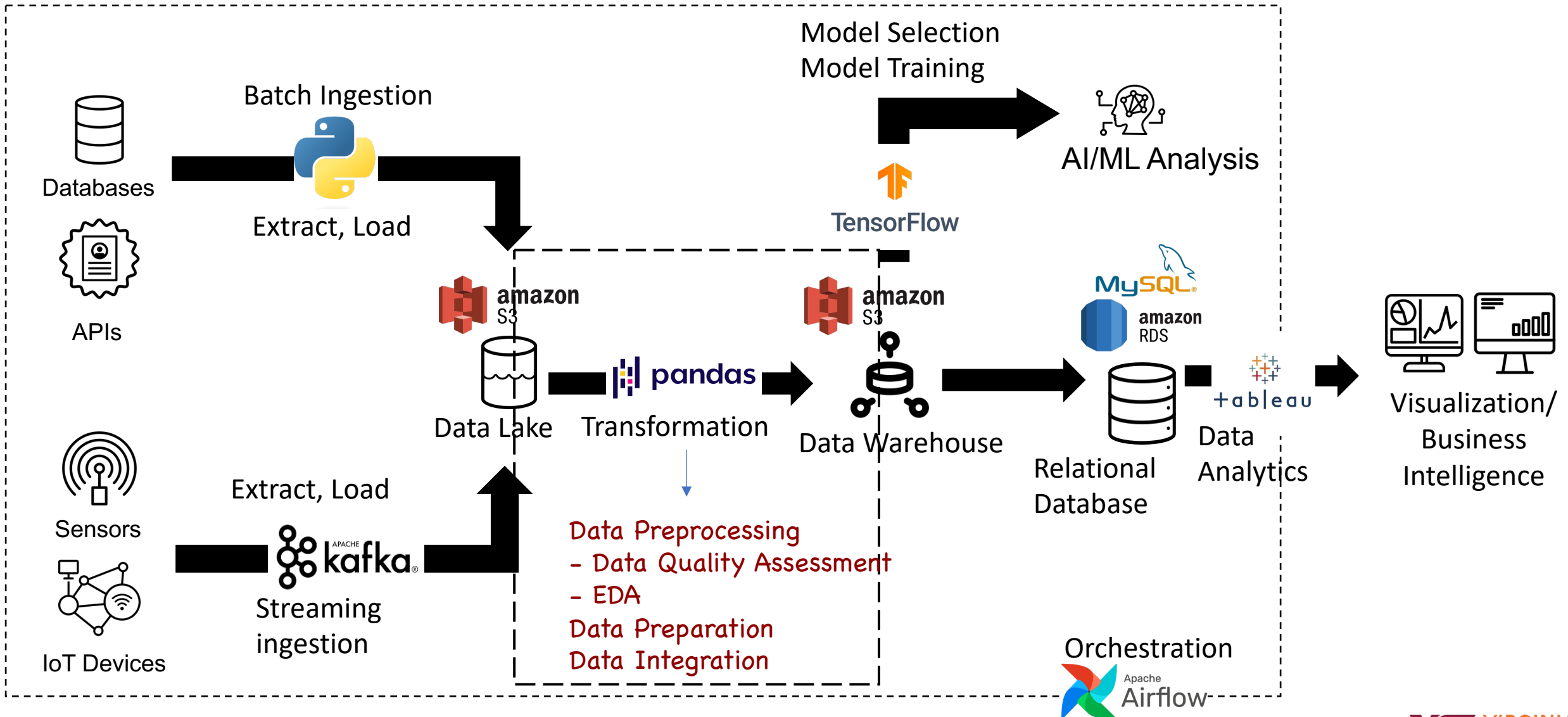
Nektaria Tryfona, PhD

Electrical and Computer Engineering  
Virginia Tech

# Objectives

- Data Quality Assessment
- Exploratory Data Analysis

# Custom Data Engineering Pipeline



# Data Quality Assessment

# The Role of Data Quality in Data-intensive Tasks

Real-world data is noisy, incomplete, inconsistent:

- **Noisy:** errors/ outliers  
Erroneous values: e.g., salary= -10K  
Unexpected values: e.g., salary= 100K when the rest dataset lies in [30K-50K]
- **Incomplete:** missing data  
Missing values: e.g., occupation=""  
Missing attributes of interest: e.g., no information on occupation
- **Inconsistent:** discrepancies in the data  
Example: student grade ranges differ across countries, in USA [A-F] and in GR [0-10]

“Garbage in, garbage out”



“Dirty” data → poor results

# Data Quality Dimensions

## Accuracy

The ages in the dataset match the actual ages of the individuals

## Completeness

The dataset includes unique data on all the employees, include their names, ages, and addresses

## Consistency

Data about a particular customer is the same across all the records in the dataset

## Timeliness

The data in the dataset reflects the most recent sales information

## Relevance

The data in the dataset is relevant to the products the user is interested in

## Validity

The correct format for phone numbers, and does not contain any errors or invalid entries

# Data Quality Dimensions (mostly in ML)

## Amount of Training Data

Does the dataset include enough data for the model at hand?

## Feature Relevance

Analyze relative importance of each feature with respect to the target variable or with other features

## Bias

Is the dataset biased?

## Outliers

How many outliers? The presence of outliers in data increases the misclassification

# **Exploratory Data Analysis (EDA)**

**getting to know our data**



# Exploratory Data Analysis (EDA)

Goal:

- analyze and investigate data sets
- summarize their main characteristics
- determine how best to manipulate data sources to get the meaningful answers
- discover patterns, spot anomalies, test a hypothesis, or check assumptions

# EDA with Pandas DataFrames

**Pandas DataFrame** is a two-dimensional data structure

- i.e., data is aligned in a tabular fashion in rows and columns

Pandas DataFrame consists of three principal components: **data**, **rows**, and **columns**

	<i>Name</i>	<i>Team</i>	<i>Number</i>	<i>Position</i>	<i>Age</i>
0	Avery Bradley	Boston Celtics	0.0	PG	25.0
1	John Holland	Boston Celtics	30.0	SG	27.0
2	Jonas Jerebko	Boston Celtics	8.0	PF	29.0
3	Jordan Mickey	Boston Celtics	NaN	PF	21.0
4	Terry Rozier	Boston Celtics	12.0	PG	22.0
5	Jared Sullinger	Boston Celtics	7.0	C	NaN
6	Evan Turner	Boston Celtics	11.0	SG	27.0

# Basic EDA functions

Function: `df.info()`

Prints a concise summary of a DataFrame including the **dtype** and **columns**, **non-null values**, and **memory usage**

Example:

```
<class 'pandas.core.frame.DataFrame'>
RangeIndex: 458 entries, 0 to 457
Data columns (total 9 columns):
Name           457 non-null object
Team           457 non-null object
Number         457 non-null float64
Position       457 non-null object
Age            457 non-null float64
Height         457 non-null object
Weight         457 non-null float64
College        373 non-null object
Salary         446 non-null float64
dtypes: float64(4), object(5)
memory usage: 32.3+ KB
```

	Name	Team	Number	Position	Age	Height	Weight	College	Salary
0	Avery Bradley	Boston Celtics	0.0	PG	25.0	6-2	180.0	Texas	7730337.0
1	Jae Crowder	Boston Celtics	99.0	SF	25.0	6-6	235.0	Marquette	6796117.0
2	John Holland	Boston Celtics	30.0	SG	27.0	6-5	205.0	Boston University	NaN
3	R.J. Hunter	Boston Celtics	28.0	SG	22.0	6-5	185.0	Georgia State	1148640.0
4	Jonas Jerebko	Boston Celtics	8.0	PF	29.0	6-10	231.0	NaN	5000000.0
5	Amir Johnson	Boston Celtics	90.0	PF	29.0	6-9	240.0	NaN	12000000.0
6	Jordan Mickey	Boston Celtics	55.0	PF	21.0	6-8	235.0	LSU	1170960.0
7	Kelly Olynyk	Boston Celtics	41.0	C	25.0	7-0	238.0	Gonzaga	2165160.0
8	Terry Rozier	Boston Celtics	12.0	PG	22.0	6-2	190.0	Louisville	1824360.0
9	Marcus Smart	Boston Celtics	36.0	PG	22.0	6-4	220.0	Oklahoma State	3431040.0
10	Jared Sullinger	Boston Celtics	7.0	C	24.0	6-9	260.0	Ohio State	2569260.0
11	Isaiah Thomas	Boston Celtics	4.0	PG	27.0	5-9	185.0	Washington	6912869.0
12	Evan Turner	Boston Celtics	11.0	SG	27.0	6-7	220.0	Ohio State	3425510.0
13	James Young	Boston Celtics	13.0	SG	20.0	6-6	215.0	Kentucky	1749840.0
14	Tyler Zeller	Boston Celtics	44.0	C	26.0	7-0	253.0	North Carolina	2616975.0
15	Bojan Bogdanovic	Brooklyn Nets	44.0	SG	27.0	6-8	216.0	NaN	3425510.0

DataFrame Example

# Basic EDA functions

**Function:** `df.describe()`

Generates descriptive statistics including those that summarize the **central tendency**, **dispersion** and **shape of a dataset's distribution**, excluding **NaN (Not a Number)** values

Example:

	Number	Age	Weight	Salary
<b>count</b>	457.000000	457.000000	457.000000	4.460000e+02
<b>mean</b>	17.678337	26.938731	221.522976	4.842684e+06
<b>std</b>	15.966090	4.404016	26.368343	5.229238e+06
<b>min</b>	0.000000	19.000000	161.000000	3.088800e+04
<b>25%</b>	5.000000	24.000000	200.000000	1.044792e+06
<b>50%</b>	13.000000	26.000000	220.000000	2.839073e+06
<b>75%</b>	25.000000	30.000000	240.000000	6.500000e+06
<b>max</b>	99.000000	40.000000	307.000000	2.500000e+07

	Name	Team	Number	Position	Age	Height	Weight	College	Salary
0	Avery Bradley	Boston Celtics	0.0	PG	25.0	6-2	180.0	Texas	7730337.0
1	Jae Crowder	Boston Celtics	99.0	SF	25.0	6-6	235.0	Marquette	6796117.0
2	John Holland	Boston Celtics	30.0	SG	27.0	6-5	205.0	Boston University	NaN
3	R.J. Hunter	Boston Celtics	28.0	SG	22.0	6-5	185.0	Georgia State	1148640.0
4	Jonas Jerebko	Boston Celtics	8.0	PF	29.0	6-10	231.0	NaN	5000000.0
5	Amir Johnson	Boston Celtics	90.0	PF	29.0	6-9	240.0	NaN	12000000.0
6	Jordan Mickey	Boston Celtics	55.0	PF	21.0	6-8	235.0	LSU	1170960.0
7	Kelly Olynyk	Boston Celtics	41.0	C	25.0	7-0	238.0	Gonzaga	2165160.0
8	Terry Rozier	Boston Celtics	12.0	PG	22.0	6-2	190.0	Louisville	1824360.0
9	Marcus Smart	Boston Celtics	36.0	PG	22.0	6-4	220.0	Oklahoma State	3431040.0
10	Jared Sullinger	Boston Celtics	7.0	C	24.0	6-9	260.0	Ohio State	2569260.0
11	Isaiah Thomas	Boston Celtics	4.0	PG	27.0	5-9	185.0	Washington	6912869.0
12	Evan Turner	Boston Celtics	11.0	SG	27.0	6-7	220.0	Ohio State	3425510.0
13	James Young	Boston Celtics	13.0	SG	20.0	6-6	215.0	Kentucky	1749840.0
14	Tyler Zeller	Boston Celtics	44.0	C	26.0	7-0	253.0	North Carolina	2616975.0
15	Bojan Bogdanovic	Brooklyn Nets	44.0	SG	27.0	6-8	216.0	NaN	3425510.0

DataFrame Example

# Basic EDA functions

Function: `df.isnull()`

Detects **missing** values. Returns a boolean same-sized object indicating if the values are empty (or missing)

For example, the following command displays the number of **empty** values for each column in descending order

```
df.isnull().sum().sort_values(ascending = False)
```

```
Example:      College      85
            Salary      12
            Name        1
            Team        1
            Number      1
            Position    1
            Age         1
            Height      1
            Weight      1
            dtype: int64
```

	Name	Team	Number	Position	Age	Height	Weight	College	Salary
0	Avery Bradley	Boston Celtics	0.0	PG	25.0	6-2	180.0	Texas	7730337.0
1	Jae Crowder	Boston Celtics	99.0	SF	25.0	6-6	235.0	Marquette	6796117.0
2	John Holland	Boston Celtics	30.0	SG	27.0	6-5	205.0	Boston University	NaN
3	R.J. Hunter	Boston Celtics	28.0	SG	22.0	6-5	185.0	Georgia State	1148640.0
4	Jonas Jerebko	Boston Celtics	8.0	PF	29.0	6-10	231.0	NaN	5000000.0
5	Amir Johnson	Boston Celtics	90.0	PF	29.0	6-9	240.0	NaN	12000000.0
6	Jordan Mickey	Boston Celtics	55.0	PF	21.0	6-8	235.0	LSU	1170960.0
7	Kelly Olynyk	Boston Celtics	41.0	C	25.0	7-0	238.0	Gonzaga	2165160.0
8	Terry Rozier	Boston Celtics	12.0	PG	22.0	6-2	190.0	Louisville	1824360.0
9	Marcus Smart	Boston Celtics	36.0	PG	22.0	6-4	220.0	Oklahoma State	3431040.0
10	Jared Sullinger	Boston Celtics	7.0	C	24.0	6-9	260.0	Ohio State	2569260.0
11	Isaiah Thomas	Boston Celtics	4.0	PG	27.0	5-9	185.0	Washington	6912869.0
12	Evan Turner	Boston Celtics	11.0	SG	27.0	6-7	220.0	Ohio State	3425510.0
13	James Young	Boston Celtics	13.0	SG	20.0	6-6	215.0	Kentucky	1749840.0
14	Tyler Zeller	Boston Celtics	44.0	C	26.0	7-0	253.0	North Carolina	2616975.0
15	Bojan Bogdanovic	Brooklyn Nets	44.0	SG	27.0	6-8	216.0	NaN	3425510.0

DataFrame Example

# Basic EDA functions

**Function:** `df.apply(pd.Series.nunique)`

Counts the number of **unique values** in each column of a DataFrame

Example:

```
Name      457
Team       30
Number     53
Position    5
Age        22
Height     18
Weight     87
College    118
Salary     309
dtype: int64
```

	Name	Team	Number	Position	Age	Height	Weight	College	Salary
0	Avery Bradley	Boston Celtics	0.0	PG	25.0	6-2	180.0	Texas	7730337.0
1	Jae Crowder	Boston Celtics	99.0	SF	25.0	6-6	235.0	Marquette	6796117.0
2	John Holland	Boston Celtics	30.0	SG	27.0	6-5	205.0	Boston University	NaN
3	R.J. Hunter	Boston Celtics	28.0	SG	22.0	6-5	185.0	Georgia State	1148640.0
4	Jonas Jerebko	Boston Celtics	8.0	PF	29.0	6-10	231.0	NaN	5000000.0
5	Amir Johnson	Boston Celtics	90.0	PF	29.0	6-9	240.0	NaN	12000000.0
6	Jordan Mickey	Boston Celtics	55.0	PF	21.0	6-8	235.0	LSU	1170960.0
7	Kelly Olynyk	Boston Celtics	41.0	C	25.0	7-0	238.0	Gonzaga	2165160.0
8	Terry Rozier	Boston Celtics	12.0	PG	22.0	6-2	190.0	Louisville	1824360.0
9	Marcus Smart	Boston Celtics	36.0	PG	22.0	6-4	220.0	Oklahoma State	3431040.0
10	Jared Sullinger	Boston Celtics	7.0	C	24.0	6-9	260.0	Ohio State	2569260.0
11	Isaiah Thomas	Boston Celtics	4.0	PG	27.0	5-9	185.0	Washington	6912869.0
12	Evan Turner	Boston Celtics	11.0	SG	27.0	6-7	220.0	Ohio State	3425510.0
13	James Young	Boston Celtics	13.0	SG	20.0	6-6	215.0	Kentucky	1749840.0
14	Tyler Zeller	Boston Celtics	44.0	C	26.0	7-0	253.0	North Carolina	2616975.0
15	Bojan Bogdanovic	Brooklyn Nets	44.0	SG	27.0	6-8	216.0	NaN	3425510.0

DataFrame Example



# Basic EDA functions

**Function:** `df.duplicated()`

Returns Boolean Series denoting duplicate rows

Example:

```
0      False
1      False
2      False
3      False
4      False
...
453    False
454    False
455    False
456    False
457    False
Length: 458, dtype: bool
```

	Name	Team	Number	Position	Age	Height	Weight	College	Salary
0	Avery Bradley	Boston Celtics	0.0	PG	25.0	6-2	180.0	Texas	7730337.0
1	Jae Crowder	Boston Celtics	99.0	SF	25.0	6-6	235.0	Marquette	6796117.0
2	John Holland	Boston Celtics	30.0	SG	27.0	6-5	205.0	Boston University	NaN
3	R.J. Hunter	Boston Celtics	28.0	SG	22.0	6-5	185.0	Georgia State	1148640.0
4	Jonas Jerebko	Boston Celtics	8.0	PF	29.0	6-10	231.0	NaN	5000000.0
5	Amir Johnson	Boston Celtics	90.0	PF	29.0	6-9	240.0	NaN	12000000.0
6	Jordan Mickey	Boston Celtics	55.0	PF	21.0	6-8	235.0	LSU	1170960.0
7	Kelly Olynyk	Boston Celtics	41.0	C	25.0	7-0	238.0	Gonzaga	2165160.0
8	Terry Rozier	Boston Celtics	12.0	PG	22.0	6-2	190.0	Louisville	1824360.0
9	Marcus Smart	Boston Celtics	36.0	PG	22.0	6-4	220.0	Oklahoma State	3431040.0
10	Jared Sullinger	Boston Celtics	7.0	C	24.0	6-9	260.0	Ohio State	2569260.0
11	Isaiah Thomas	Boston Celtics	4.0	PG	27.0	5-9	185.0	Washington	6912869.0
12	Evan Turner	Boston Celtics	11.0	SG	27.0	6-7	220.0	Ohio State	3425510.0
13	James Young	Boston Celtics	13.0	SG	20.0	6-6	215.0	Kentucky	1749840.0
14	Tyler Zeller	Boston Celtics	44.0	C	26.0	7-0	253.0	North Carolina	2616975.0
15	Bojan Bogdanovic	Brooklyn Nets	44.0	SG	27.0	6-8	216.0	NaN	3425510.0

DataFrame Example

<https://pandas.pydata.org/docs/>

# Summary

- Data Quality Assessment
- Exploratory Data Analysis



# **Data Engineering Project**

## **Module 3**

### **Data Quality Assessment and Data Exploration**

Nektaria Tryfona, PhD

Electrical and Computer Engineering  
Virginia Tech