# MARY ROSE LEGASPI
# DEVELOPMENT PROJECT

MR. MONIR (SUPERVISOR)

# BUILDING SCALABLE SOLUTION FOR PREDICTING HEART DISEASE USING APACHE SPARK MLLIB IN STANDALONE CLUSTER MODE WITH MONGODB DATABASE
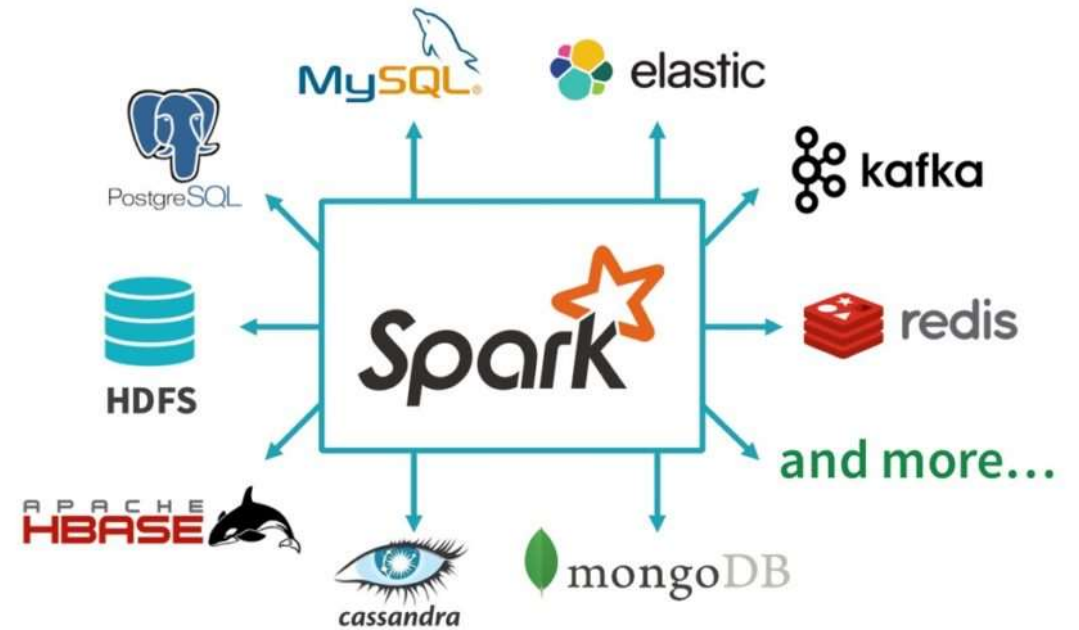
# BACKGROUND AND CONTEXT

# WHAT IS SPARK?

# What is Apache Spark?

- Apache open sourced project originally developed at AMPLab (UC Berkeley)

- Unified general data processing engine that operates varied data workloads and platforms

- Built on top of Hadoop Map Reduce and it extends the MapReduce model to efficiently use more types of computations

# Spark features

100x faster than for large scale data processing

Simple programming layer provides powerful caching and disk persistence capabilities

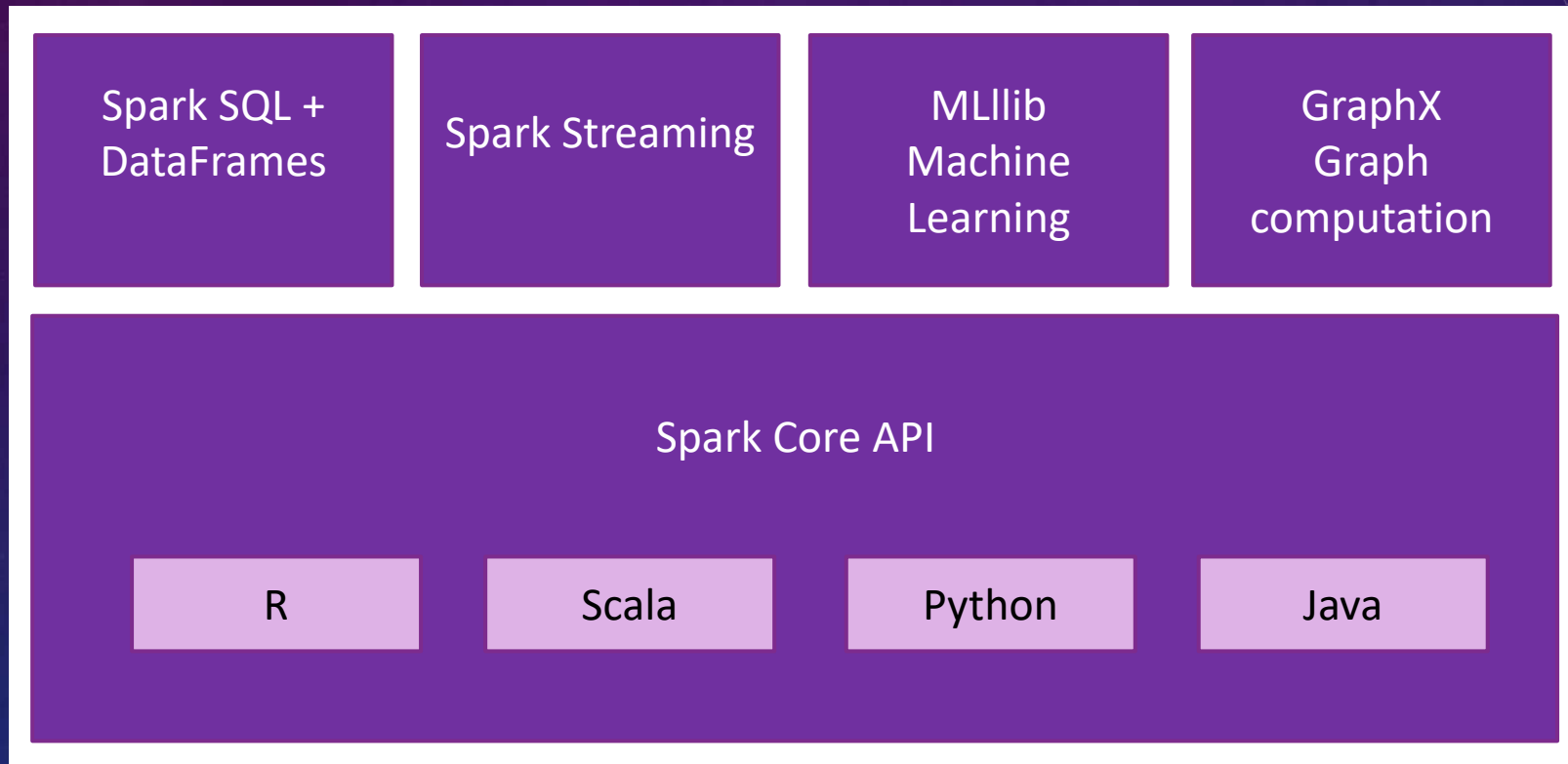Can be programmed in Scala, Java, Python, and R.

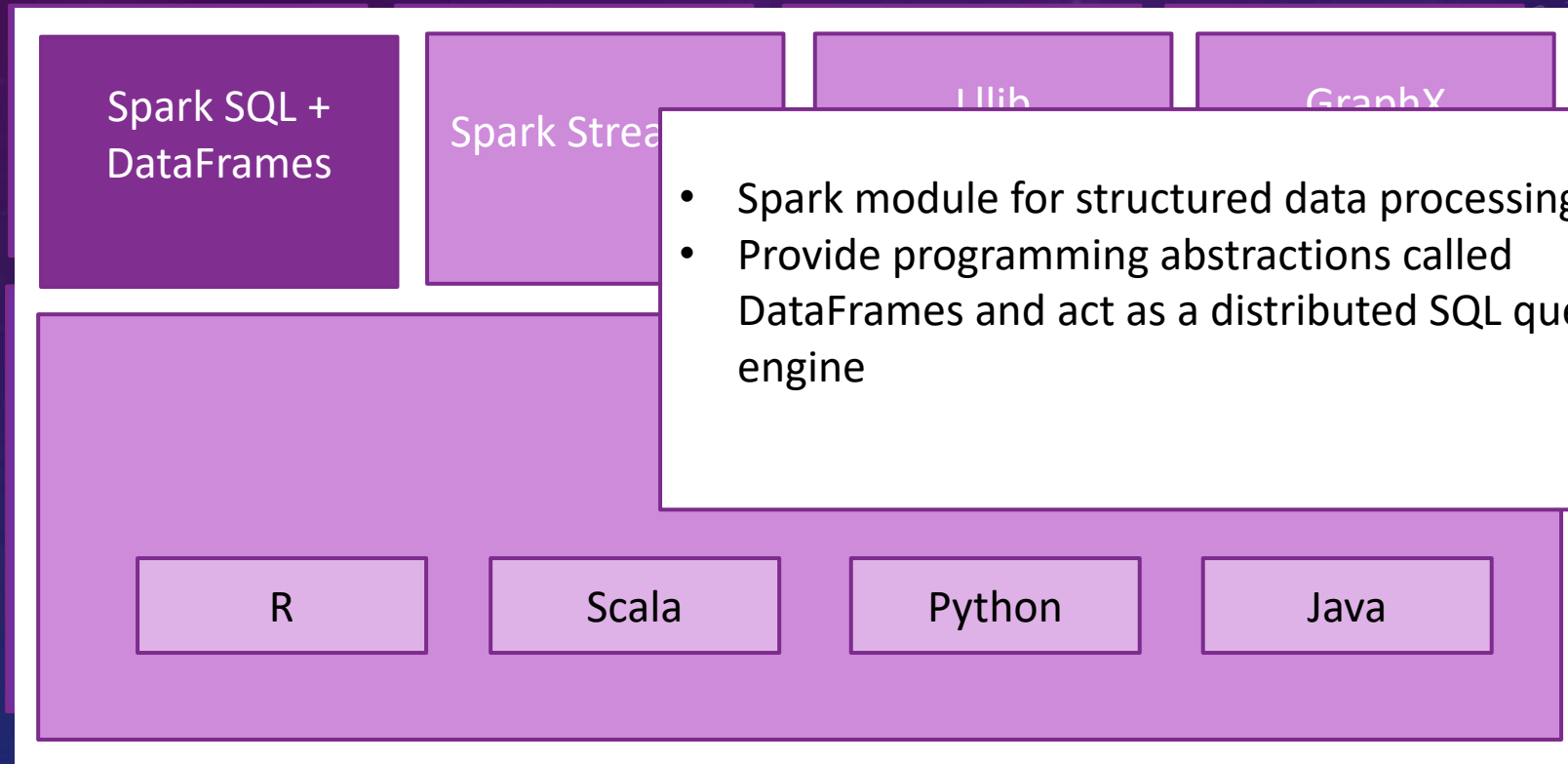Can be deployed through Mesos, Yarn, EC2 or Sparks standalone cluster manager
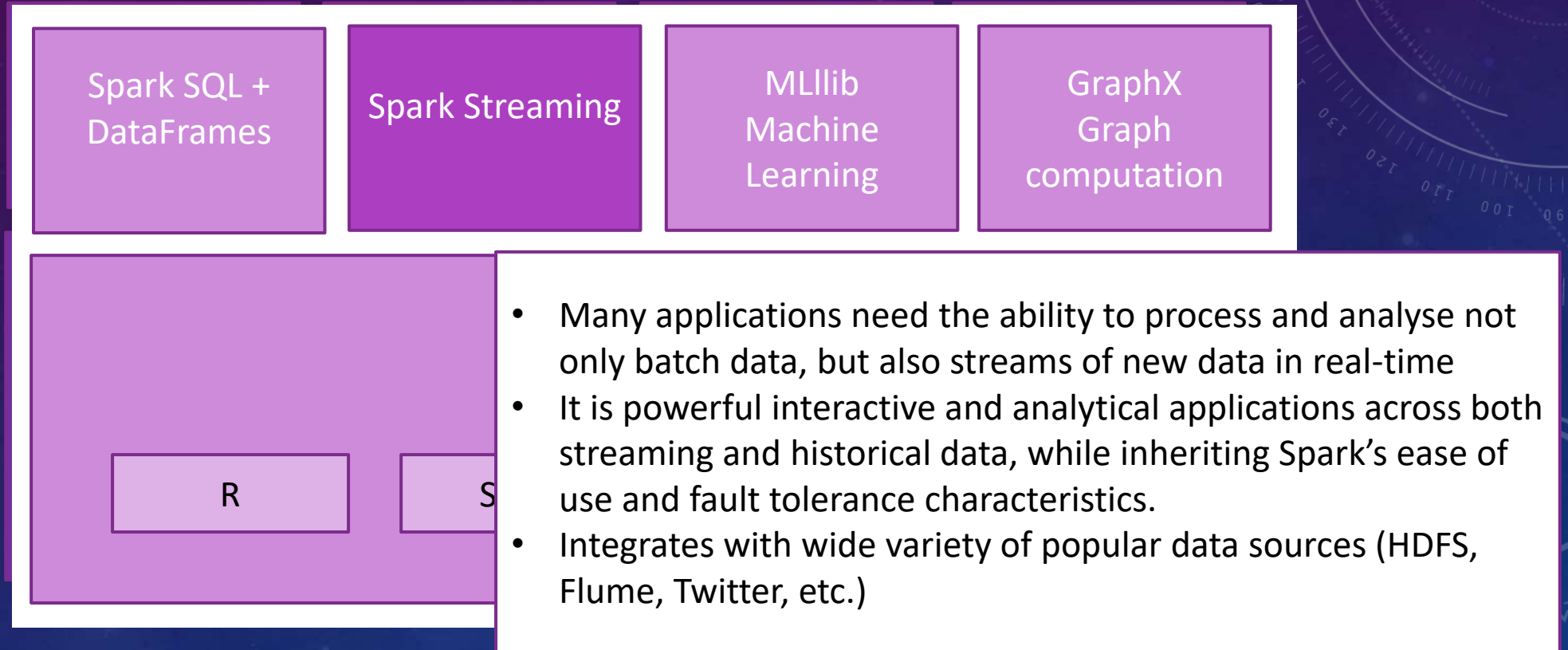
# Apache Spark Domain scenarios

# Spark Components

# Spark Components
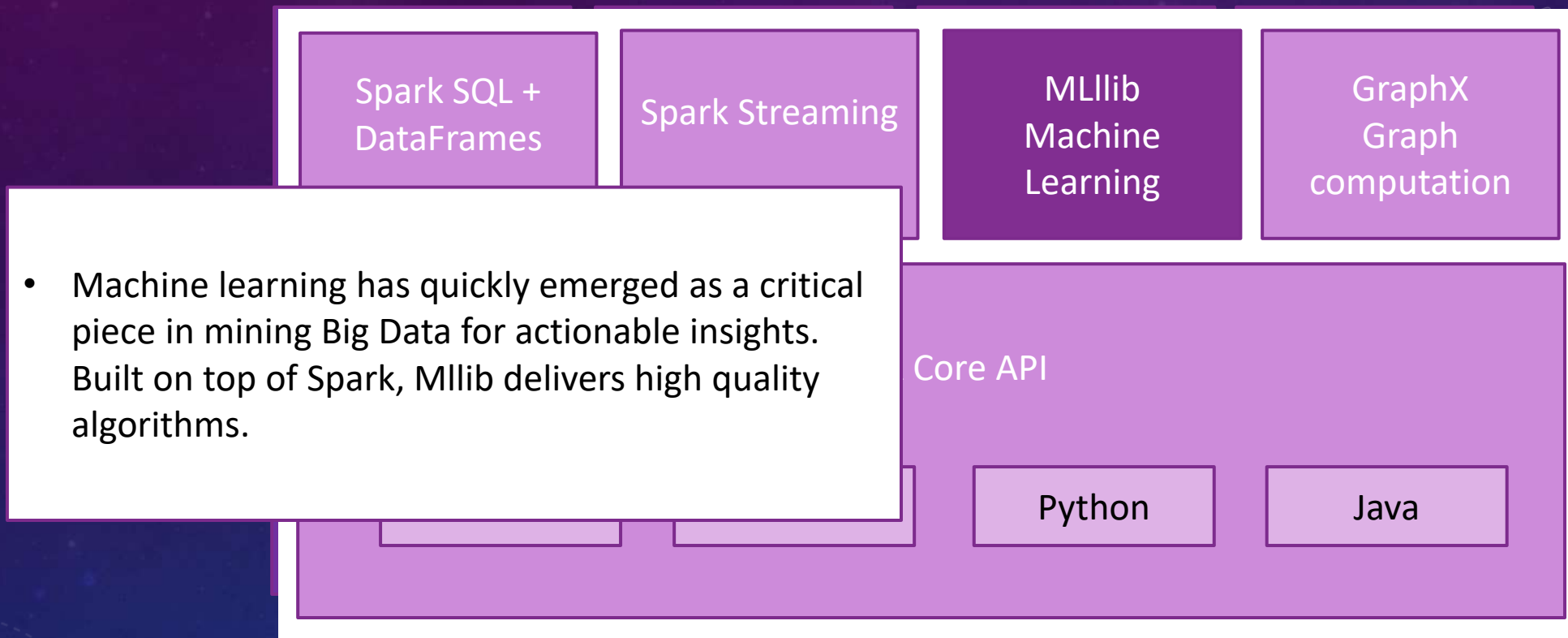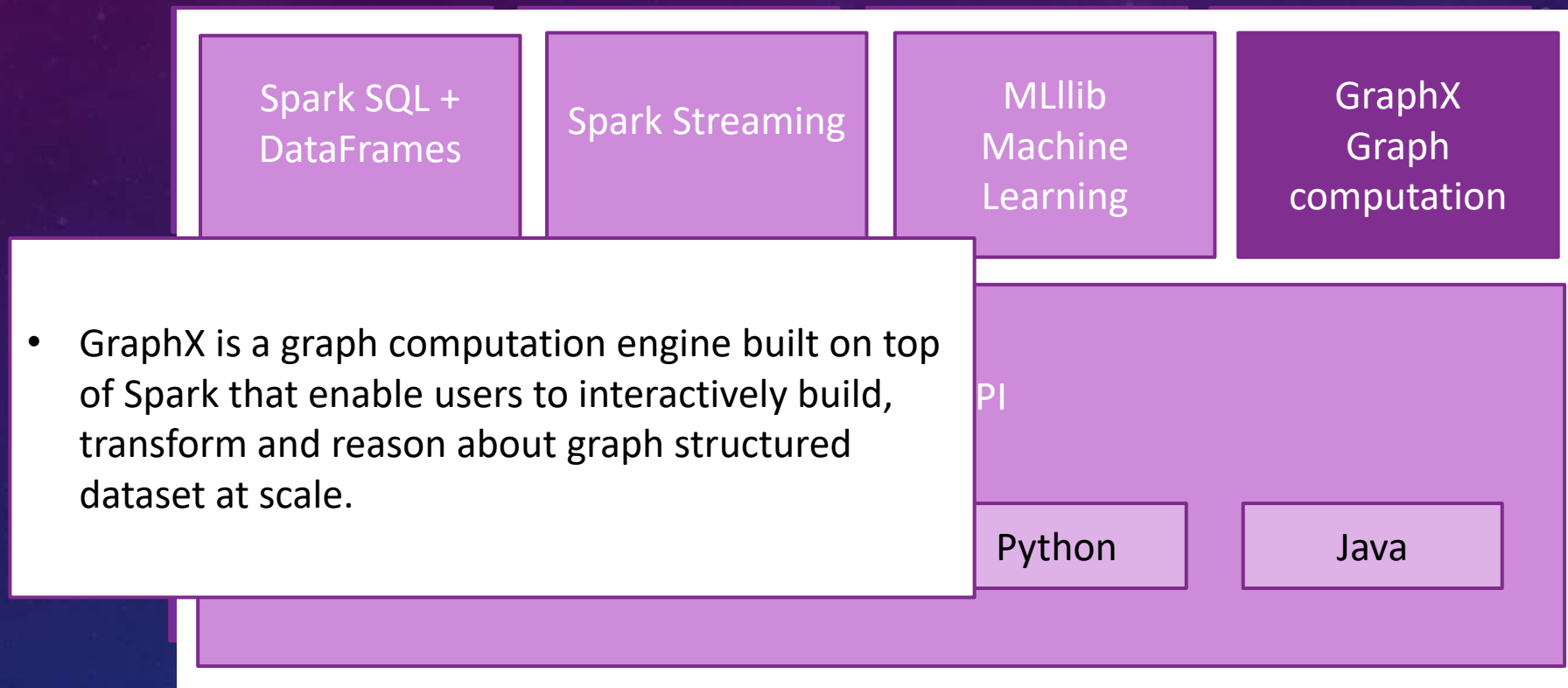
Spark SQL + DataFrames

Spark Strea...

Mllib

GraphX

- Spark module for structured data processing.
- Provide programming abstractions called DataFrames and act as a distributed SQL query engine

| R | Scala | Python | Java |

# Spark Components

| Spark SQL +<br>DataFrames | Spark Streaming | MLlib<br>Machine<br>Learning | GraphX<br>Graph<br>computation |
|---|---|---|---|

| R | | S |
|---|---|---|

- Many applications need the ability to process and analyse not only batch data, but also streams of new data in real-time
- It is powerful interactive and analytical applications across both streaming and historical data, while inheriting Spark's ease of use and fault tolerance characteristics.
- Integrates with wide variety of popular data sources (HDFS, Flume, Twitter, etc.)

# Spark Components

| Spark SQL + DataFrames | Spark Streaming | MLlib Machine Learning | GraphX Graph computation |

- Machine learning has quickly emerged as a critical piece in mining Big Data for actionable insights. Built on top of Spark, Mllib delivers high quality algorithms.

Core API

| Python | Java |

# Spark Components

| Spark SQL +<br>DataFrames | Spark Streaming | MLlib<br>Machine<br>Learning | GraphX<br>Graph<br>computation |

- GraphX is a graph computation engine built on top of Spark that enable users to interactively build, transform and reason about graph structured dataset at scale.

PI

| Python | Java |

# Spark Components

- Spark core is the underlying general execution engine for the Spark platform that all other functionality is built on top of.

- Provides in-memory computing capabilities to deliver speed, a generalized execution model to support a wide variety of applications, and Java, Python, Scala and R API for ease of development

| MLlib Machine Learning | GraphX Graph computation |
|---|---|

**Spark Core API**

| R | Scala | Python | Java |
|---|---|---|---|

# CONTEXT

# Context



**+**



- Big Data analytics
- Faster

- Health care data increasing rapidly
- Problems in health care
  - High cost
  - High waste
  - Low quality

# Context



+



Health care data volume is increasing and is expected to skyrocket as new ways of collecting health data in various forms are continuously emerging (i.e. patient centered e-health record, and wearables)

- Health care data increasing rapidly
- Problems in health care
  - High cost
  - High waste
  - Low quality

According to AbuKhousa and Campbell, the healthcare system has a massive wealth of information but knowledge poor, " there is a lack of effectual analysis tools to discover knowledge contained in the databases of these systems "

# Context



**+**



Overall spending ($170 billion spend on health in 2015-16 )

- Health care data increasing rapidly
- Problems in health care
  - High cost
  - High waste
  - Low quality

# Context



**+**

Poor quality of health (increasing mortality and morbidity rate on preventable diseases)

- Health care data increasing rapidly
- Problems in health care
  - High cost
  - High waste
  - Low quality

# PURPOSE

- Understand big data processing in health care

- Learn factors that contribute to heart disease

- Build machine learning model to predict heart disease.

- Use the Spark framework to implement analysis.

- Find significant risk factors of coronary heart disease

# DELIVERABLES

- Build a predictive analytics platform to predict heart disease using Spark's Machine Learning library module in Standalone cluster mode with Mongo DB as the database.

- Determine which features or feature subsets contribute to risk of heart disease.

- Analyse result by comparing the ground truth to predicted outcome via graph or table.

# OVERVIEW ON THE APPROACH USED

# PROJECT MANAGEMENT APPROACH

# Project Architecture

# Project Approach/Methodology

# Project Approach/Methodology

Literature Review

↓

Installation, set-up and configuration

↓

Understanding the big data framework

↓

Understanding heart disease physiology

→

Finding datasets and managing Mongodb

↓

Connecting application to spark cluster

↓

Creating a predictive modelling pipeline

↓

Result evaluation

# Project Approach/Methodology

# Project Approach/Methodology

Framingham Heart Study datasets

| | age | sex | education | currentSmoker | cigsPerDay | heartRate | BMI | glucose | diabetes | sysBP | diaBP | BPMeds | prevalentHyp | prevalentStroke | totChol | T |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| 0 | 40 | Male | 3 | Yes | 70 | 98 | 31.57 | 80 | No | 132.0 | 86.0 | No | Yes | No | 210.0 | |
| 1 | 56 | Male | 1 | Yes | 60 | 70 | 29.64 | 85 | No | 125.0 | 79.0 | No | No | No | 246.0 | |
| 2 | 59 | Male | 1 | Yes | 60 | 70 | 25.05 | 84 | No | 153.5 | 105.0 | No | Yes | No | 298.0 | |
| 3 | 58 | Male | 2 | Yes | 60 | 75 | 32.00 | 65 | No | 150.0 | 97.0 | No | Yes | No | 250.0 | |
| 4 | 39 | Male | 1 | Yes | 60 | 59 | 23.60 | 78 | No | 112.0 | 65.0 | No | No | No | 215.0 | |

- csv format
- 4241 records

Framingham Heart Study datasets

| Categorical variables | Numeric/continuous variables |
| --- | --- |
| Sex | age |
| Education | cigsPerDay |
| currentSmoker | heartrate |
| Diabetes | BMI |
| BPMeds | Glucose |
| prevalentHyp | sysBP |
| prevalentStroke | diaBP |
| label | totChol |

- Importing datasets

- Managing database using mongoshell

# Project Approach/Methodology

Literature Review

↓

Installation, set-up and configuration

↓

Understanding the big data framework

↓

Understanding heart disease physiology

Finding datasets and managing Mongodb

↓

Connecting application to spark cluster

↓

Creating a predictive modelling pipeline

↓

Result evaluation

# Deployment stage: Starting Mongodb service



Starting Mongodb service

# Deployment stage : starting the spark master cluster

# Deployment stage: starting 4 worker nodes

# Deployment stage: starting 4 worker nodes



Application running

# Deployment stage:



Application User Interface

Application running but
No task executed yet

# Deployment stage:

# Deployment stage:

# Deployment stage: Stopping Spark Session

**Stop Spark Session**

```
In [13]:    1  spark.stop()
```

![Spark logo] 2.3.1 **Spark Master at spark://mary-VirtualBox:7077**

**URL:** spark://mary-VirtualBox:7077
**REST URL:** spark://mary-VirtualBox:6066 *(cluster mode)*
**Alive Workers:** 4
**Cores in use:** 4 Total, 0 Used
**Memory in use:** 8.0 GB Total, 0.0 B Used
**Applications:** 0 Running, 1 Completed
**Drivers:** 0 Running, 0 Completed
**Status:** ALIVE

No running applications

## Workers (4)

| Worker Id | Address | State | Cores | Memory |
|---|---|---|---|---|
| worker-20181012120722-10.0.2.15-36303 | 10.0.2.15:36303 | ALIVE | 1 (0 Used) | 2.0 GB (0.0 B Used) |
| worker-20181012120725-10.0.2.15-46477 | 10.0.2.15:46477 | ALIVE | 1 (0 Used) | 2.0 GB (0.0 B Used) |
| worker-20181012120727-10.0.2.15-42737 | 10.0.2.15:42737 | ALIVE | 1 (0 Used) | 2.0 GB (0.0 B Used) |
| worker-20181012120730-10.0.2.15-45341 | 10.0.2.15:45341 | ALIVE | 1 (0 Used) | 2.0 GB (0.0 B Used) |

## Running Applications (0)

| Application ID | Name | Cores | Memory per Executor | Submitted Time | User | State | Duration |
|---|---|---|---|---|---|---|---|

## Completed Applications (1)

| Application ID | Name | Cores | Memory per Executor | Submitted Time | User | State | Duration |
|---|---|---|---|---|---|---|---|
| app-20181012120820-0000 | SparkDemo_DecisionTrees_FraminghamDataset_specific | 4 | 1024.0 MB | 2018/10/12 12:08:20 | mary | FINISHED | 1.1 h |

# FINDINGS/OUTCOMES

# Project Approach/Methodology

Rank of the most important feature with GBT

| Accuracy | Precision | Recall | f1 |
|----------|-----------|--------|-----|
| 85.08% | 81.30% | 85.08% | 81.59% |

Age, body mass index, blood sugar level, systolic blood pressure and cholesterol are the most important determining factor to assess risk of having coronary heart disease

# Findings/outcome

| | | Ground Truth | | | |
|---|---|---|---|---|---|
| | Total Population 905 | Condition Positive 135 | Condition Negative 770 | Prevalence 14.92% | |
| Prediction | Predictive Outcome Positive 135 | True Positive 114 | False Positive 21 | Positive Predictive Value 84.44% | False Discovery Rate 15.5% |
| | Predictive Outcome Negative 770 | False Negative 21 | True Negative 749 | False Omission Rate 2.72% | Negative Predictive Value 97.27% |
| | Accuracy 85.08% | True Positive Rate 85.08% | False Positive Rate 2.73% | | |
| | | False Negative Rate 15% | True Negative Rate 97.27% | | |

Digest features and determine which specific subsets or range is likely contributing to coronary heart disease.

```
Here are the results!
A ensemble using GBT accuracy  : 86.17%  precision:74.26%  recall:86.17%  f1:79.77%
###################################################################################
***********************************************************************************
     risk   features
0   0.379652      60s
1   0.244219      50s
2   0.208924      40s
3   0.166460      30s
4   0.000745      70s
###################################################################################
```

# Findings/outcome: by Glucose level



```
Here are the results!
---------------------------------------------------------------------------------
A single decision tree accuracy: 85.41%   precision:73.16%   recall:85.41%   f1:78.82%
---------------------------------------------------------------------------------
A randomForestEnsemble accuracy: 85.55%   precision:73.18%   recall:85.55%   f1:78.88%
---------------------------------------------------------------------------------
A ensemble using GBT accuracy  : 85.41%   precision:73.16%   recall:85.41%   f1:78.82%
##################################################################################
      risk features
0  0.255947  125-174
1  0.182892  100-124
2  0.139361    40-99
3  0.124695  275-324
4  0.113635  325-399
5  0.098622  225-274
6  0.084848  175-224
7  0.000000     400+
##################################################################################
```

```
A randomForestEnsemble accuracy: 86.60%  precision:74.99%  recall:86.60%  f1:80.37%
-----------------------------------------------------------------------------------
A ensemble using GBT accuracy   : 86.60%  precision:74.99%  recall:86.60%  f1:80.37%
###################################################################################
        risk features
0   0.290246       1-9
1   0.210968     40-49
2   0.129764     30-39
3   0.124156     20-29
4   0.110447     50-59
5   0.066960     10-19
6   0.062431     60-69
7   0.005028       70+
###################################################################################
```



Smoking Relevance — bar chart of Risk Percentage vs Number of cigarette (1-9, 40-49, 30-39, 20-29, 50-59, 10-19, 60-69, 70+)

# IMPLICATIONS

# Implications

- Health care authorities may use the findings to focus their health education not only to elderly but also the young individuals as the risk for heart disease increased significantly higher as they grow older.

- Regulatory commissions to study sugar content of consumable goods and reduce use of sugar.

- Promote healthy lifestyle by quitting cigarette smoking as oppose to cutting down.

- Although the dataset is not a big data, the framework could work using the big data. Thus, applicable for scalable projects in predicting heart disease.

THANK YOU FOR LISTENING