→

YARN
|

```
Resource
   manages
Application
   manages
containes
```

7/07/2023

* Rprogramming

• Statistical Measures

→ Descriptive         : Inferential

1. central tendancy or moments (eg. Mean, mode,
                                         median)
2. Measures of positions → Quantiles (Fractions)
3.     "        " Distribution → Skewners and kurtosis
4.     "        " Dispersion (Spread) → variance, Std.deviation,
                                      covariance, corelation
                                                      Coeff

→ Inferential

   • Probability
       eg: population, Sample

          → P-test
          → t-test
          → Z-test
          → chi-square test

# * installing packages in R

install.packages("moments")
library(moments)

```
> x = c(4, 1, 8, 9, 10)
> min(x)
> max(x)
> sort(x)
> range(x)
> var(x)
> sd(x)
> cov(x, y)
> cor(x, y)
> skewness(x)
> kurtosis(x)
> mean(x)
> median(x)
```

Std. devi $\Rightarrow \sigma = \sqrt{\frac{\epsilon(x-\bar{x})^2}{n}}$

var $= \sigma^2 = \frac{\epsilon(x-\bar{x})^2}{n}$

population $\Rightarrow /n$

Sample $\Rightarrow /n-1$

⇒ mode has no predefined function in 'R'

```
> quantile(x, probs = 0.10)
Q1 = quantile(x, probs = 0.25) → quartile
     quantile(x, probs = 0.50) → median
Q3 = quantile(x, probs = 0.75) → quartile
> IQR (Q3 - Q1)
> data()        // display all predefined data sets
> data(iris)    // iris dataset
> str(iris)     // gives attributes
> summary(iris) // gives summary
```

Quantile $= \dfrac{fraction per}{100}$ (n+1)

> iris [1:150,]    // entire dataset

> iris [1:10,]    // 10 rows

> iris [1:150, -5]    // displays only 4 columns

> iris [1:150, c(-4,-5)]    // removes last two col,
                            // more than one value sthd
                                give in vector

# Programs

## 1. Linear Regression

> x = c (1, 2, 3, 4)

> y = c (2, 4, 6, 8)

> z = data·frame (x, y)

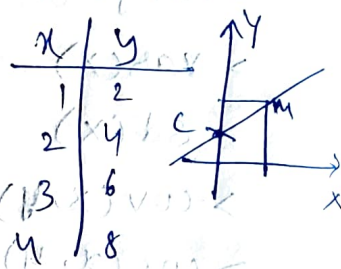> lr = lm (y ~ x),        // linear model is predefined func
        data = z)

                         // multiple => lm (y~(x₁+x₂)
> print (lr)    // gives summary        , data=)

> new = data·frame (x = 5)

> p = predict (lr, new)

> print (p)

## 2. Multiple Regression

# 3. Decision Tree Model (DTM)

> recursive partition regression trees

```r
> install.packages (rpart)
> library (of class)
> library ( rpart.plot)

> x = iris [1:140, ]
> y = iris [141:150, ]
> dtm = rpart (Species ~ (sepal.length + sepal. width +
                        petal.length + petal.width))
            (or)                      data = x , method --"class"

  dtm = rpart (Species ~., data = x , method = "class")


> plot (dtm)
> text (dtm)    // gives summary based on branch

> P = predict (dtm, y[,-5] , type = " class")
> print (P)
> table (y[,5], P)
```

|     | S | VC | V |
|-----|---|----|---|
| S   | . |    | 0 |
| VC  | 0 | 0  | 0 |
| V   | 0 | 0  | 10 |

• Random forest : Iterative of multiple DTM's

# Scala

⇒ Scalable Language

| Hadoop 1.0 | Hadoop 2.0 | Hadoop 3.0 |
|---|---|---|
| Hdfs + MapReduce | Hdfs + MapReduce + YARN | Hdfs + MapReduce + Spark → load balancing |

Modules developed by
R/Python

Modules developed by:
Scala

## Scala programs

### Syntax

```
object mca
{

    def main (args: Array [String])
    {
        //

    }

}

class classname
{
    def f1()
    {
    }

    def f2()
    {
    }
}
```

→ Val : value which is constant
→ Var : variable which is not-fixed

Syntax :   Val x: Int = 95;  → Immutable

Var x : Int = 100;  → mutable

→ object declaration

Syntax :  Var obj = new classname();

* object mca
{
    def main (args : Array [String])
    {
        var x: Int = 50 ;
        Vas y: Double = 49.99 ;
        Vas z: String = "Ram";
        Print (x + " " + y + " " + s) ;
    }
}

* Reading &/p directly from keyboard

var x: Int = Scala. io. Std In . read Line. to Int ;

Vas y: String = Scala. io. StdIn. read Line;

object peuledu
{
    def main (args: Array [String])
    {
        var z: Double = Scala. io, std zn. readLine. to Double;
        Print ( z ) ;
    }
}

# * Static Single dimensional array

Object areyKoushik
{
    def main (args : Array [String])
    {
        var x = Array (3,4,5,8,6);
        var y = Array (Array (1,2,3,4),
                           Array (4,3,2,1)) ;

        // var x: Array [int] = Array (1,2,3,4);

        preinste

        for (i ← 0 to 4)
        {
            print (x(i));
        }

        for (i ← 0 to 1)
        {
            for (j ← 0 to 4)
            {
                print (y(i)(j));
            }
        }
    }
}

# * Dynamic Single D & Two D array

```scala
object nasavnemasta
{
    def main (args: Array [String])
    {
        var x = Array [Int] (5);
        for (i ← o to 4)
        {
            x(i) = Scala. io. StdIn. readLine. toInt;
        }

        for (i ← o to4)
        {
            printth(x(i));
        }

        var y = Array. of Dim [Int] (2,5);
        for (i ← o to 1)
        {
            for (j ← o to 4)
            {
                y(i)(j) = Scala. io. StdIn. read Line. toInt;
            }
        }

        for (i ← o to 1)
        {
            for (j ← o to 4)
            {
                Println (y(i)(j));
            }
        }

        for (i ← o to 1;
             j ← o to 4)
        {
            print (y(i)(j))
        }
    }
}
```

# Exercise

→ Create method overloading

→ constructor overloading

→ simple inheritance

→ method overriding

28|7

→ Real time appl. for method overloading & constr.
overloading

4|08

## * Spark

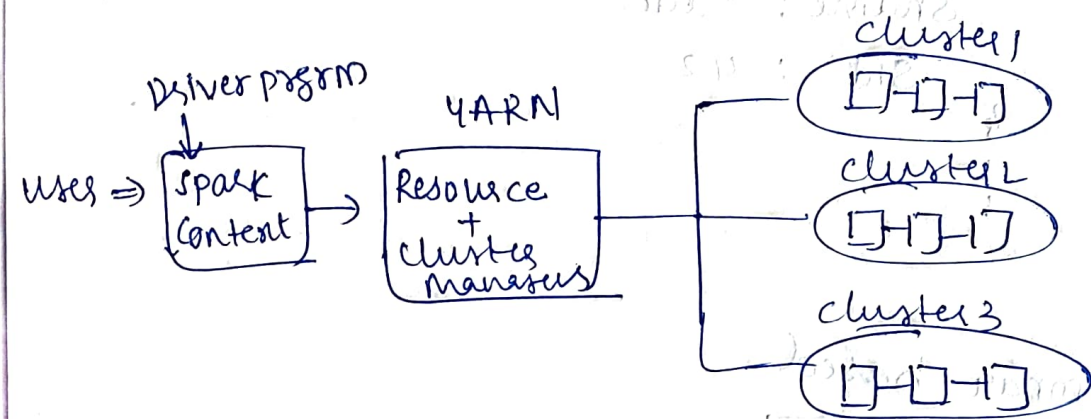→ Hadoop 3.0

• Spark Ecosystem :-

| Data Management | Ozzie | Zookeeper | R, Java Python, Scala |
|---|---|---|---|
| Data Accessing | Spark MLlib | Spark Streaming | Spark SQL |
| Data Processing | Map Reduce | | Spark Core → RDD |
| Data Storage | Rdbms | | Spark HBASE |

- RDD : Re-silent Distributed Databases
  └→ Spark Context

    • Re-silent ⇒ fault-tolerance

→ In-memory processing is done in Spark



## * NOSQL

→ Not Only SQL Databases

  1. Key - value based

  2. columnar - based

  3. Graph - based

  4. Document - based

11/05

→ BASE properties

    BA → Basically Available

    S → Soft state

    E → Eventually consistency

1. Key -value based

   relation : attributekey = " value"

ℰ  Student : "sname" = "xyz"

   Student :  " sid"  = 45

## 2. Document - based

→ json / BSON

Student
```
{
    {
        Sname: "adc"
        Sid  : 42
    }
    ⋮
}
```

BSON
↓
either 1 0r0

## 3. columnar - based

→ Dynamically create column for the row we
needed instead of creating a column for
entire table like in RDBMS

↙
row-based Seggregation

| abc | 48 | 94385 | |
|-----|----|-------|-----|
| xyz | 49 | 43695 | 9.6 |

column -based

| abc | | 40 | | 984 |
|-----|---|----|---|-----|
| xyz | | 49 | | 866 |

- row based takes more traversal time than
column based

- row based → can increase as many columns
as we need