

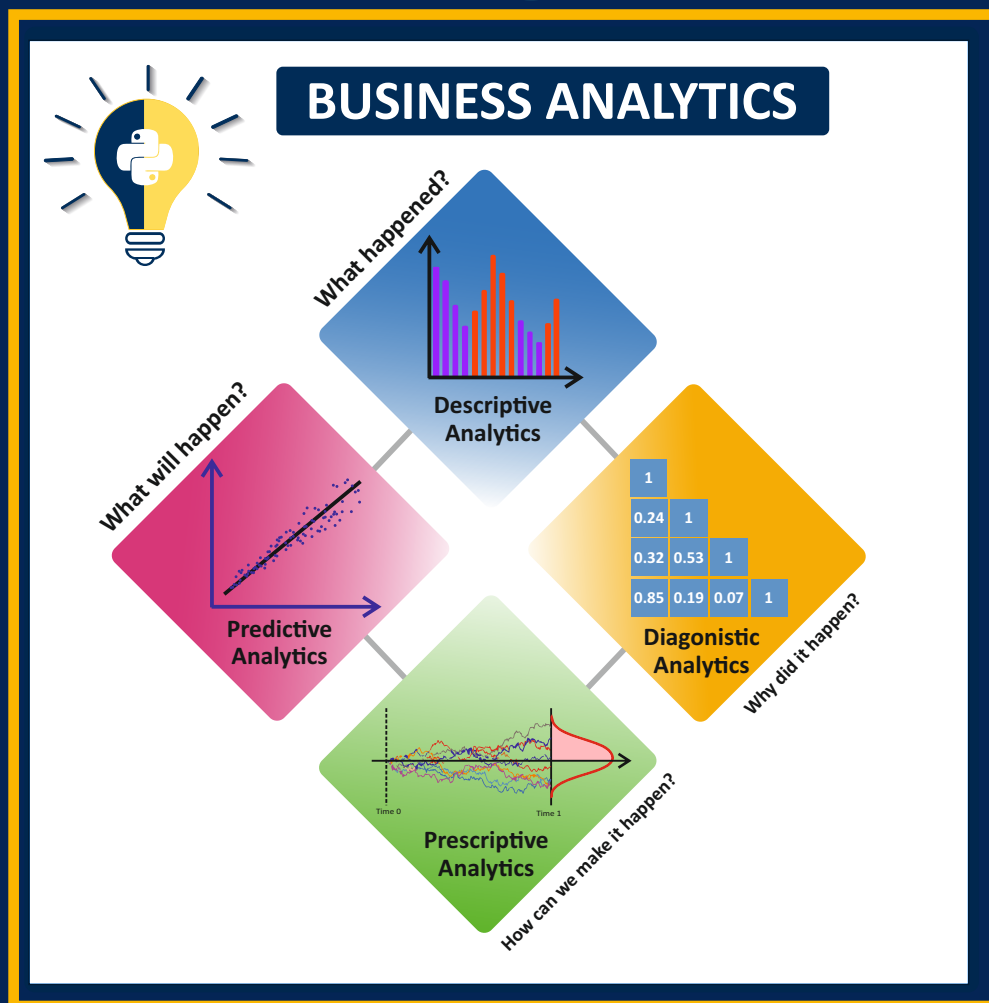
# BUSINESS ANALYTICS

## using Python

50+ hours

Case Study and Project- driven Methodology

Blended Learning Methodology



PEAKS<sup>2</sup>TAILS



## DETAILED CURRICULUM

### MODULE 1 -FEATURE ENGINEERING

#### STEPS OF DATA CLEANING & PROCESSING

- Identifying Data type
- Exploratory data analysis - Data Summarisation and visualisation
- Missing value treatment for Categorical Variables
- Missing value treatment for Continuous Variables
- Outliers treatment for Categorical Variables
- Outliers treatment for Countinuous Variables
- Balancing of data
- Covariaties Creation (Enrichment)
- Dimensions Reductions
- Dummy Coding for Categorical Variables
- Scaling for continuous Variables
- Discretisation or Weight of Evidence
- Data Partitioning

### MODULE 2 -MODEL BUILDING

#### STEPWISE REGRESSION

- Removing problems of Multicollinearity,
- Selecting Important Variables
- Correcting Problems of autocorrelation
- Checking problems of Non-Linearity
- Then correcting the Problems of Heteroskedasticity





## DETAILED CURRICULUM

### TYPES OF REGRESSION

- Multiple Regression
- Dummy Independent Variable - Dummy Regression
- Dummy Dependant Variable - Logistic Regression
- Penalised Regression - Ridge & Lasso Regression
- Forward Selection, Forward stagewise and least angle Regression
- Bayesian regression with Spike & Slab Selection
- Support Vector Regression
- Principal Component Regression

### LOGISTIC REGRESSION

- General Linear Modelling in Excel using 4 main link functions Normal, Poisson, Binomial, Gamma.
- Multinomial Logit, Ordered Logit Model

### DECISION TREE

- Classification and Decision Tree
- Classification & Regression Tree
- CHAID
- Boosting, Bagging & Random Forest

### SUPPORT VECTOR MACHINE

- Primal & Dual Formulation
- Linear SVM, Non Linear SVM using Slack Variables
- Kernel Trick and Radial Basis Function

### LINEAR DISCRIMINANT ANALYSIS

- Maximum Likelihood
- Fisher's Discriminant
- Bayesian Discrimination

### K -NEAREST NEIGHBOUR

- 1.K-Nearest neighbour
- 2.K-means prototype



## DETAILED CURRICULUM

### NAIVE BAYES CLASSIFIER

- Bayes Theorem
- Naïve Bayes Classifier

### CLUSTERING

- Bottom up Clustering a.k.a K means clustering
- Bottom up Clustering a.k.a Hierarchical clustering
- Top down Clustering a.k.a Minimal Spanning Tree
- Clustering using Expectation Maximization
- Mixed Variables Clustering - K means Prototype

## MODULE 3 - MODELLING VALIDATION

### VALIDATION METRICS

- Gini/AR
- AUROC/ CAP
- Unconditional Entropy
- Conditional Entropy
- Kullback-Leibler Divergence
- Kolmogorov-Smirnov (KS)
- Information Value

## MODULE 4 - PROJECTS

### PROJECTS

- Application to Credit data
- Application to HR data
- Application to Fraud detection
- Application to Marketing data

## BACKGROUND

### BACKGROUND

The amount of data that is being generated on a daily basis has increased multiple times, the technological advancement in terms of storing this data has improved considerably. Therefore, in turn the ability to analyse the data and generate insights for data driven decision making becomes primarily importance However the availability of manpower with data science and machine learning skills is limited. The objective of the course is to introduce the concepts of data science and machine learning to the participants using python.



## OBJECTIVE

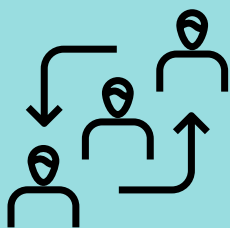
### OBJECTIVE

Develop predictive models using various statistical and machine learning techniques, Interpret and evaluate various models and its generalization, Hands on experience on the usage of open notebooks in Python like Jupiter.



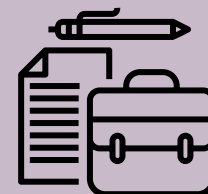
### WHO CAN ATTEND

Beginner candidates from various quantitative backgrounds, like Engineering, Finance, Maths, Business Management who are looking for Business Analytics training to start their career in the field of Analytics and Data Science.



### PEDADOGY

Professionals working in analytics field or students interested to make a career in analytics



# DEMO MODELS

DURING THE PROGRAM YOU WILL LEARN TO  
CREATE EXCEL MODELS LIKE SHOWN BELOW

```
In [28]: 1 edibles = ["ham", "spam", "eggs", "nuts"]
2 for food in edibles:
3     if food == "spam":
4         print("No more spam please!")
5         break
6     print("Great, delicious " + food)
7 else:
8     print("I am so glad: No spam!")
9     print("Finally, I finished stuffing myself")

Great, delicious ham
No more spam please!
Finally, I finished stuffing myself

In [45]: 1 for num in range(2, 10):
2     if num % 2 == 0:
3         print("Found an even number", num)
4         continue
5     print("Found a number", num)

Found an even number 2
Found a number 3
Found an even number 4
Found a number 5
Found an even number 6
Found a number 7
Found an even number 8
Found a number 9

In [46]: 1 for x in 'abcd':
2     for y in 'abcd':
3         print(x, y)
4         print('----')

a a
a b
a c
a d
----
b a
b b
```

```
In [74]: 1 plt.scatter(y=sales,x=newspaper)
Out[74]: <matplotlib.collections.PathCollection at 0x5bcf9910>

In [75]: 1 plt.scatter(y=sales,x=radio)
Out[75]: <matplotlib.collections.PathCollection at 0x5b87130>
```

```
Important Variable Selection

n [206]: 1 # use lprice as the dependent variable and all the other variables as independent
2 # lprice = b1*crime + b2*rooms + dist + radial + stratio + lowstat + lproptax
3 # calculate pvalue
4 # if pvalue is greater than 5% we will delete the independent variable with the

n [207]: 1 # model = ols('lprice ~ crime+rooms+dist+radial+stratio+lowstat+lproptax', data)
2
3 def equation(dependent, independent):
4     for i in dependent:
5         z = i+'-'
6         for j in independent:
7             z = z + j + '+'
8         z = z[:-1]
9     return z

n [208]: 1 eq = equation(y, x)

Autocorrelation

n [232]: 1 # we will store the error terms in a variable - 'error'
2 # we will store the lag1 term of the error - errorlag1
3 # 'error' - dependent variable and rooms, stratio, lowstat, lproptax and error
4 # find p-value
5 # if pvalue or error lag1 is less than 5%, we will say that problem of autocorre

n [233]: 1 x.head()
Out[233]:
```

	rooms	stratio	lowstat	lproptax
0	0.57	15.3	4.98	5.690300
1	0.42	17.8	0.14	5.488938
2	7.18	17.8	4.03	5.488938
3	7.00	18.7	2.94	5.402878
4	7.15	18.7	5.33	5.402878

```
Pseudo R-squared

In [24]: 1 prsq = (model.llnull-model.llf)/model.llnull
2 prsq
Out[24]: 0.18224598528527883

Confusion Matrix

In [360]: 1 pred = model.predict()
2 model.pred_table()
Out[360]: array([[495., 204.],
 [ 75., 225.]])

ROC curve

In [361]: 1 from sklearn import metrics
2
3 fpr, tpr, _ = metrics.roc_curve(y, pred)
4
5 plt.plot([0,1],[0,1],linestyle='--')
6 plt.plot(fpr,tpr,marker='.')
7 plt.show()
Out[361]: <function matplotlib.pyplot.show(*args, **kw)>
```

## FREQUENTLY ASKED QUESTIONS

### PREREQUISITE



Knowledge of Basic Excel

### CERTIFICATE



Silver Certificate on successful completion of projects .  
Gold Certification on passing a 2 hours MCQ based exam.

### FEES



Rs.20000

### DURATION



50+ hours

## ABOUT THE TRAINER



Karan Aggarwal is one of India's leading trainers in Financial Modelling, Risk Modelling, Data Analytics and academic programs like Financial Risk Manager (FRM) & Actuarial Science. He has spearheaded several solution accelerators and spreadsheet-based prototypes in Risk and Analytics space. Karan has also authored a number of books on Advanced Excel, Statistical Modelling, Risk Modelling & Machine Learning. He is widely regarded for his problem-solving, thought leadership and intrapreneurship skills. His analytical mindset, solid fundamentals & the thirst to keep learning set him apart as a true authority in this field. Karan has also been awarded the Young Indian Entrepreneur Award by the Confederation Of Indian Industries in the year 2017.



# OUR TRAINEES WORK IN



# OUR SERVICES

1



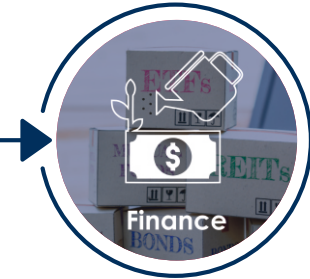
2



3



4



98 74 98 74 98

