

# ניהול פרויקטים וניתוח מערכות בעולם ה BIG DATA

קורס 2425 – 40 שעות

## אודות הקורס

ארגונים רבים שוקלים לעבור לעולם ה-Big-Data הם נתקלים במושגים חדשים רבים אינספור פלטפורמות, ספקי שירותים, וטכנולוגיות חדשות. פערי הידע של הארגונים אינם מתנקזים לכדי נושא אחד אלא מתפרשים על פני כמה תחומים שכל אחד מהם הוא עולם בפני עצמו. לשם כך על מנהלי הפרויקט/ מנתחי המערכות וראשי הצוותים נדרשים לדעת ולהכיר עולם זה, בשלב הראשון יש צורך להבין את היתרונות העסקיים הטמונים בעולם זה, להכיר ברמה גבוהה את המרכיבים של עולם זה. מידע זה נדרש כיום לכל מי שעוסק בתחום המחשוב הארגוני

האתגרים העיקריים הניצבים כיום בפני ארגונים השוקלים או נמצאים בתהליך של מעבר לעולם ה-Big-Data:

- הבנת הצורך והערך העסקי בעולם הביג דטה
- הבנה ממעוף הציפור: מה הן מושגים והמונחים המרכיבים את עולם ה Big-Data
- אפיון הרכיבים המשותפים לעולם זה לעומת אילו היכולים להיעזר בסביבות המסורתיות
- איתור ניקוי ואיסוף המידע
- אחסון, המידע הרב והלא מובנה בבסיסי נתונים מהסוג החדש
- ניתוח ועיבוד מידע לא-מובנה ורב (non-structured data)
- הצגתו בכלים פשוטים ובעלי משמעות אסטרטגית לארגון
- היכולת לבצע ניתוחים סטטיסטיים, כריית נתונים וביצוע למידה ממוחשבת ( Machine learning)

## מבנה הקורס :

1. הבנת הקשר וההבדל בין עולם ה BI לעולם ה BIG DATA
2. הקמת פרויקט ביג דטה המורכב מ 5 שלבים:
  - Data exploration
  - Data cleansing
  - Data preparation
  - Data analyze
  - Data visualization

במהלך הקורס נלמד את עולם הביג דטה דרך שלבי ניהול אפיון ומימוש הפרוייקט פרט לפרוייקט הקבוצתי, במהלך כל הקורס יתרגלו החניכים את הנושאים הנלמדים בתרגול מעשי.

## תוצרים לחניכי הקורס

- יכולת איתור מקורות המידע אפשריים, כולל מתודולוגיה והכלים לאיסוף הנתונים ממקורות מידע אלו, טיפול בניקיון הנתונים.
- שיקולים לבחירת בסיס נתונים (NoSql) ומידול מודל הנתונים
- הבנה היכולת ניתוח ועיבוד מידע לא-מובנה ורב (non-structured data)
- היכרות וכלים פשוטים להצגת הנתונים ובעלי משמעות אסטרטגית לארגון
- יכולת לבצע ניתוחים סטטיסטיים, כריית נתונים וביצוע למידה ממוחשבת (Machine learning)

## כלים:

- תבניות אפיון וניהול ככלי עזר למנהל הפרוייקט ומנתח המערכות
- כלי ניהול משימות בעולם נתוני עתק כדוגמת TRELLO

## מטרות הקורס

הבנת העולם החדש, הכרת מונחים ומושגים ובהמשך נלמד לנתח ולנהל פרוייקטי ביג דטה. נלמד כיצד לנצל מערכות פשוטות וזולות לאגירת נתונים גדולים, כיצד לבחור, להגדיר ולנצל את היתרונות של סוגי מערכות חדשים ( NoSQL ) וכיצד להתחבר למקורות מידע מתוך הארגון, מהעולם הנייד ומתוך רשתות חברתיות

## קהל יעד

אנשי מערכות מידע כדוגמת מפתחים, מנתחי מערכות, ראשי צוותים, מנהלי פרויקטים, אנשי QA, מומחי יישום ורפרנטים, בעלי ניסיון מעשי בעולם מערכות המידע ומעוניינים להכיר, להבין להתמחות ולהתמקצע בתחום מערכות big data. הקורס מיועד הן לאנשי IT ולמנהלים המבקשים להיכנס לתחום ה- big data וללמוד כיצד ליישם פרויקטים מעולמות התוכן הרלוונטיים.

## דרישות קדם

- ניסיון בעולם מערכות המידע

## תכני הקורס

Module Title	Module Description
From BI to Big Data.	<ul style="list-style-type: none"> <li>▪ What is Business Intelligence</li> <li>▪ How Important is BI</li> <li>▪ From data to knowledge</li> <li>▪ BI Life cycle</li> <li>▪ Requirements in BI project</li> <li>▪ Front end / Dashboards / Data visualization</li> <li>▪ Data modeling</li> <li>▪ ETL</li> <li>▪ What if</li> <li>▪ Prediction</li> <li>▪ In memory db</li> <li>▪ Data mining</li> </ul>
Introduction to Big Data	<p><b>Introduction to Big Data</b></p> <ul style="list-style-type: none"> <li>▪ Big Data Characteristics and Use-Cases</li> <li>▪ Big Data from Business Perspective Challenges and limitations of Big Data Implementation</li> </ul> <p><b>The Data in Big Data</b></p> <ul style="list-style-type: none"> <li>▪ Structured vs. Semi-structured vs. Unstructured Data</li> <li>▪ Various Phases of Data Processing</li> <li>▪ Problems with O/R Mapping</li> </ul>

Module Title	Module Description
	<p><b>Introduction to NoSQL</b></p> <ul style="list-style-type: none"> <li>▪ NoSQL Definition</li> <li>▪ Business Drivers for NoSQL</li> <li>▪ New NoSQL Paradigms</li> <li>▪ NoSQL Categories:                             <ul style="list-style-type: none"> <li>○ Key-value Store</li> <li>○ Column Family</li> <li>○ Document DBs</li> <li>○ Graph DBs</li> </ul> </li> </ul> <p><b>Hadoop Overview</b></p> <ul style="list-style-type: none"> <li>▪ Apache Hadoop Architecture: HDFS and MapReduce</li> <li>▪ MapReduce and Big Data – Real Life Examples</li> </ul> <p><b>Introduction to Big Data Analytics</b></p> <ul style="list-style-type: none"> <li>▪ Data and Analytical Complexity</li> <li>▪ Real-Time Analytics</li> <li>▪ Big Data Visualization</li> </ul>
<p><b>Big Data: What it Means to IT Managers on the Front Lines</b></p>	<ul style="list-style-type: none"> <li>▪ Making sense of Big Data</li> <li>▪ Turning Big Data into something useful</li> <li>▪ Storage</li> <li>▪ Security</li> <li>▪ Data reconciliation</li> <li>▪ Information extraction</li> <li>▪ Insight distribution</li> <li>▪ Successfully Navigating Big Data</li> </ul>
<p><b>Project management (project steps) in big data</b></p>	<ol style="list-style-type: none"> <li>1. <b>Data exploration</b></li> <li>2. <b>Data cleansing</b></li> <li>3. <b>Data preparation</b></li> <li>4. <b>Data analyze</b></li> <li>5. <b>Data display</b></li> </ol> <p><b>1+3 data exploration and data preparation</b></p> <p><b>Introduction</b></p> <p>There are no shortcuts for data exploration. If you are in a state of mind, that machine learning can sail you away from every data storm, trust me, it won't. After some point of time, you'll realize that you are struggling at improving model's accuracy. In such situation, data exploration techniques will come to your rescue.</p> <p><b>Steps of Data Exploration and Preparation</b></p> <ul style="list-style-type: none"> <li>▪ Variable Identification</li> </ul>

Module Title	Module Description
	<ul style="list-style-type: none"> <li>▪ Univariate Analysis</li> <li>▪ Bi-variate Analysis</li> <li>▪ Missing values treatment</li> <li>▪ Outlier treatment</li> <li>▪ Variable transformation</li> <li>▪ Variable creation</li> </ul> <p><b>Missing Value Treatment</b></p> <ul style="list-style-type: none"> <li>▪ Why missing value treatment is required ?</li> <li>▪ Why data has missing values?</li> <li>▪ Which are the methods to treat missing value ?</li> </ul> <p><b>Techniques of Outlier Detection and Treatment</b></p> <ul style="list-style-type: none"> <li>▪ What is an outlier?</li> <li>▪ What are the types of outliers ?</li> <li>▪ What are the causes of outliers ?</li> <li>▪ What is the impact of outliers on dataset ?</li> <li>▪ How to detect outlier ?</li> <li>▪ How to remove outlier ?</li> </ul> <p><b>The Art of Feature Engineering</b></p> <ul style="list-style-type: none"> <li>▪ What is Feature Engineering ?</li> <li>▪ What is the process of Feature Engineering ?</li> <li>▪ What is Variable Transformation ?</li> <li>▪ When should we use variable transformation ?</li> <li>▪ What are the common methods of variable transformation ?</li> <li>▪ What is feature variable creation</li> </ul>
<b>Data Cleansing</b>	<ul style="list-style-type: none"> <li>▪ Introduction</li> <li>▪ Motivation</li> <li>▪ Data quality                             <ul style="list-style-type: none"> <li>○ Validity:</li> <li>○ DE cleansing</li> <li>○ Accuracy:</li> <li>○ Completeness:</li> <li>○ Consistency:</li> <li>○ Uniformity:</li> </ul> </li> <li>▪ The process of data cleansing                             <ul style="list-style-type: none"> <li>○ Data auditing:</li> <li>○ Workflow specification:</li> <li>○ Workflow execution:</li> <li>○ Post-processing and controlling:</li> </ul> </li> <li>▪ De-cleanse                             <ul style="list-style-type: none"> <li>○ Parsing:</li> <li>○ Data transformation</li> </ul> </li> </ul>

Module Title	Module Description
	<ul style="list-style-type: none"> <li>○ Duplicate elimination:</li> <li>○ Statistical methods:</li> <li>▪ Challenges and problems</li> </ul>
<b>Data exploration with ETL, ELT, and ETLT</b>	<ul style="list-style-type: none"> <li>▪ The ETL Bottleneck in Big Data Analytics</li> <li>▪ Apache Hadoop for Big Data</li> <li>▪ ETL, ELT, and ETLT with Apache Hadoop:</li> <li>▪ Choosing the Physical Infrastructure for ETL with Hadoop:                             <ul style="list-style-type: none"> <li>○ Compute</li> <li>○ Memory</li> <li>○ Storage</li> <li>○ Network</li> <li>○ Software</li> </ul> </li> </ul>
<b>Data collection and integration from outside resources</b>	<ul style="list-style-type: none"> <li>▪ Introduction</li> <li>▪ Hadoop distributors (Cloudera, Hortonworks, MapR)</li> <li>▪ Building blocks of Hadoop (NameNode, DataNode...)</li> <li>▪ Introduction to HDFS</li> <li>▪ Map-reduce pattern</li> <li>▪ Distribute cache</li> <li>▪ Introduction to Hive for ad-hoc queries                             <ul style="list-style-type: none"> <li>○ Hive basics</li> <li>○ Hive data types</li> <li>○ HiveQL</li> </ul> </li> <li>▪ Pig:                             <ul style="list-style-type: none"> <li>○ Introduction to Pig as data flow language</li> <li>○ Pig Latin basic expressions</li> <li>○ Operators for data processing</li> </ul> </li> <li>▪ YARN (Map-Reduce 2)                             <ul style="list-style-type: none"> <li>○ Motivation for YARN.</li> <li>○ Architecture.</li> <li>○ Features.</li> </ul> </li> </ul>
<b>Data modeling in Big data world</b>	<ul style="list-style-type: none"> <li>▪ Saving data methods:                             <ul style="list-style-type: none"> <li>○ Hadoop ,Files, XML, CSV, RDBMS,OLAP, Tabular, No-SQL db</li> </ul> </li> <li>▪ RDBMS challenge in Big-data world</li> <li>▪ No-SQL vs traditional relational data</li> <li>▪ Scaling vs. consistency</li> <li>▪ No-Sql database types:                             <ul style="list-style-type: none"> <li>○ Key Value DB</li> <li>○ Column-Family DB</li> <li>○ Document DB</li> <li>○ graph DB</li> </ul> </li> </ul>

Module Title	Module Description
	<ul style="list-style-type: none"> <li>▪ Transaction in No-SQL</li> <li>▪ Applying map-reduce</li> <li>▪ No-SQL leading implementations</li> </ul>
<b>Date inquiry in Big data</b>	<ul style="list-style-type: none"> <li>▪ Hive                             <ul style="list-style-type: none"> <li>○ Introduction to Hive for ad-hoc queries</li> <li>○ Hive basics</li> </ul> </li> <li>▪ Hbase                             <ul style="list-style-type: none"> <li>○ Introduction to Pig as data flow language</li> <li>○ Introduction to Hbase for processing huge tables</li> <li>○ Hbase data model</li> <li>○ Hbase vs. RDBMS</li> </ul> </li> <li>▪ Data mining ,using R</li> </ul>
<b>Data visualization</b>	<ul style="list-style-type: none"> <li>▪ 1:2:3 Method to data visualization</li> <li>▪ What are the changes that should be made in our methods</li> <li>▪ Alerts &amp; exceptions</li> <li>▪ Links &amp; networks</li> <li>▪ Hue – the new gate to big data</li> </ul>
	<ul style="list-style-type: none"> <li>▪ Presentation/ project</li> </ul>