

# داده‌کاوی

## *Data Mining*

M.M. Pedram  
[pedram@khu.ac.ir](mailto:pedram@khu.ac.ir)  
Kharazmi University  
(Fall 2009)

1

### اهداف درس

- ❖ آموزش مفاهیم یادگیری ماشین و داده کاوی
- ❖ آرایه چشم‌اندازی از زمینه های تحقیقاتی فعلی و آینده
- ❖ آشنایی با کاربردهای موفق
- ❖ آشنایی با اصول انتخاب روش مناسب برای یک مساله خاص

2

## ارزیابی

- ❖ پایان نیم سال + امتحان های کوچک (کوئیز)
- ❖ تکالیف عادی + تکالیف برنامه نویسی
- ❖ ۲ مورد تحقیق + گزارش + ارائه
- ❖ زمان ارائه تحقیق اول: ۲۰ آبان
- ❖ زمان ارائه تحقیق دوم: ۲۰ آذر

► 3- introduction

KHU, M.M.Pedram, [pedram@khu.ac.ir](mailto:pedram@khu.ac.ir)

Data Mining, Fall 2009

3

## سطوح علم



\* معرفت (wisdom)

\* دانش (knowledge)

\* اطلاعات (information)

\* داده (data)

\* جهل (ignorance)

***“Data is not information; information is not knowledge; knowledge is not wisdom.”*** Gary Flake, Principal Scientist & Head of Yahoo! Research Labs, July 2004.

► 4- introduction

KHU, M.M.Pedram, [pedram@khu.ac.ir](mailto:pedram@khu.ac.ir)

Data Mining, Fall 2009

4

## Data Mining Definition

❖ **داده کاوی:** استخراج الگوها یا دانش از حجم زیادی داده.

❖ **DATA MINING:** exploration & analysis

- by **automatic means actionable**
- of **large quantities of data**
- to discover **actionable** patterns & rules

❖ Data mining *a way to utilize massive quantities of data that businesses generate*

## Data Mining Definition

The search for interesting patterns and models,  
in large data collections,  
using statistical and machine learning methods,  
and high-performance computational infrastructure.

**Key point:** applications are

- data-driven and
- compute-intensive

## Characteristics of Data Mining Applications

### ❖ Data

- ▶ Lots of data, numerous sources
- ▶ Noisy: missing values, outliers, interference
- ▶ Heterogeneous: mixed types, mixed media
- ▶ Complex: scale, resolution, temporal, spatial dimensions

▶ 7- introduction

KHU, M.M.Pedram, [pedram@khu.ac.ir](mailto:pedram@khu.ac.ir)

Data Mining, Fall 2009

7

## همگرایی سه تکنولوژی



- ❖ DBMS
- ❖ AI, Machine Learning, Pattern Recognition
- ❖ Data Visualization

▶ 8- introduction

KHU, M.M.Pedram, [pedram@khu.ac.ir](mailto:pedram@khu.ac.ir)

Data Mining, Fall 2009

8

## ❖ اسامی دیگر

- knowledge discovery (mining) in databases (KDD) (*wrong!*)
- knowledge extraction
- data/pattern analysis
- data archeology
- data dredging
- information harvesting
- business intelligence
- ...

## ❖ آیا هر مطالعه ای داده کاوی است؟

- query processing.
- expert systems or small ML/statistical programs

## DM Tasks (Goals)

## ❖ There are two categories of goals (or high level tasks) in DM:

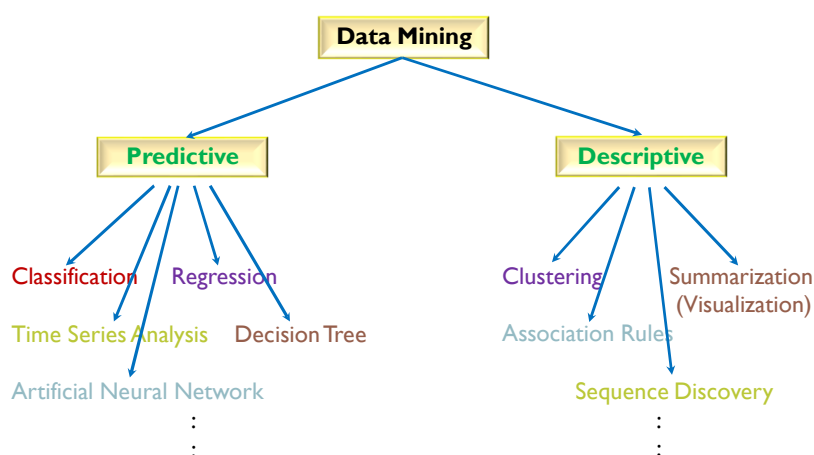
- **Description:** models are constructed to describe particular patterns or relationships in the data
- **Prediction:** models are constructed using historical cases to predict outcomes for new cases

## DM Tasks (Goals)

Another definition:

- ▶ **Descriptive data mining:** Search massive data sets and discover the locations of unexpected structures or relationships, patterns, trends, clusters, and outliers in the data.
- ▶ **Predictive data mining:** Build models and procedures for regression, classification, pattern recognition, or machine learning tasks, and assess the predictive accuracy of those models and procedures when applied to fresh data.

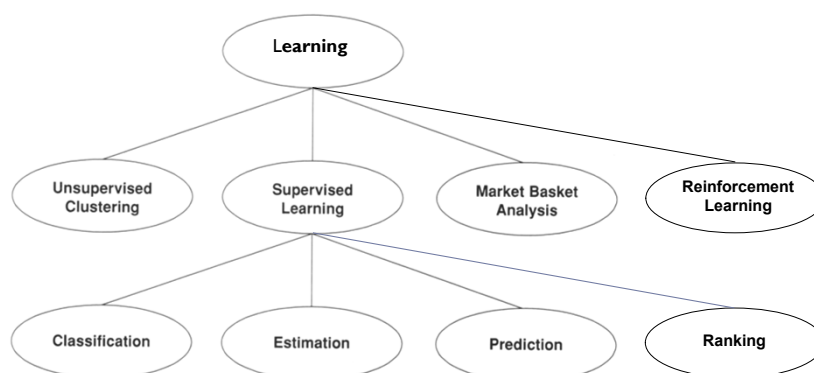
## Data Mining Tech

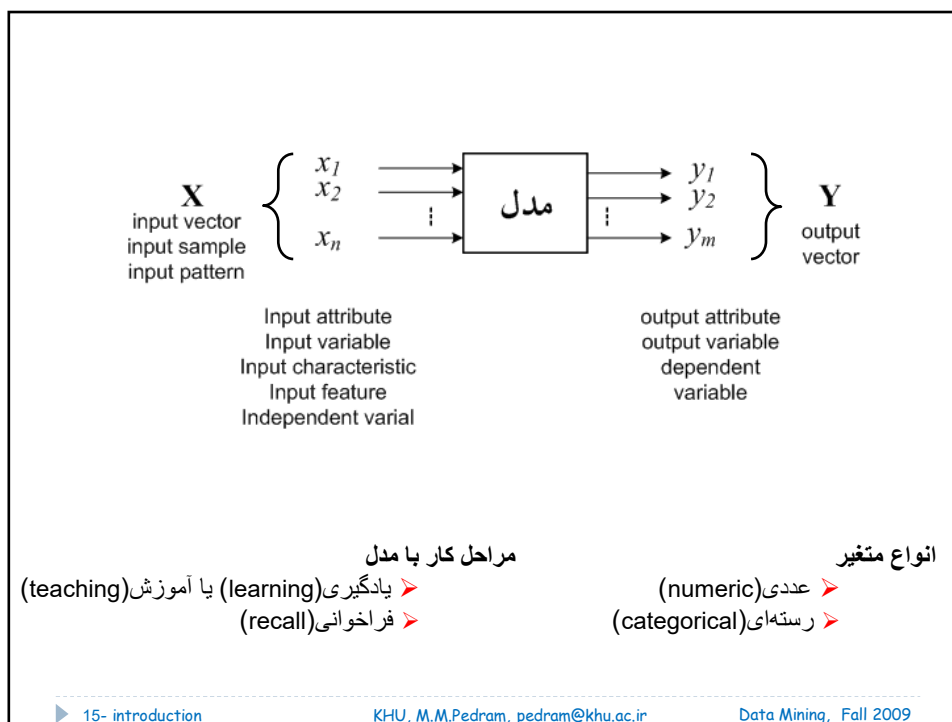


## Other names

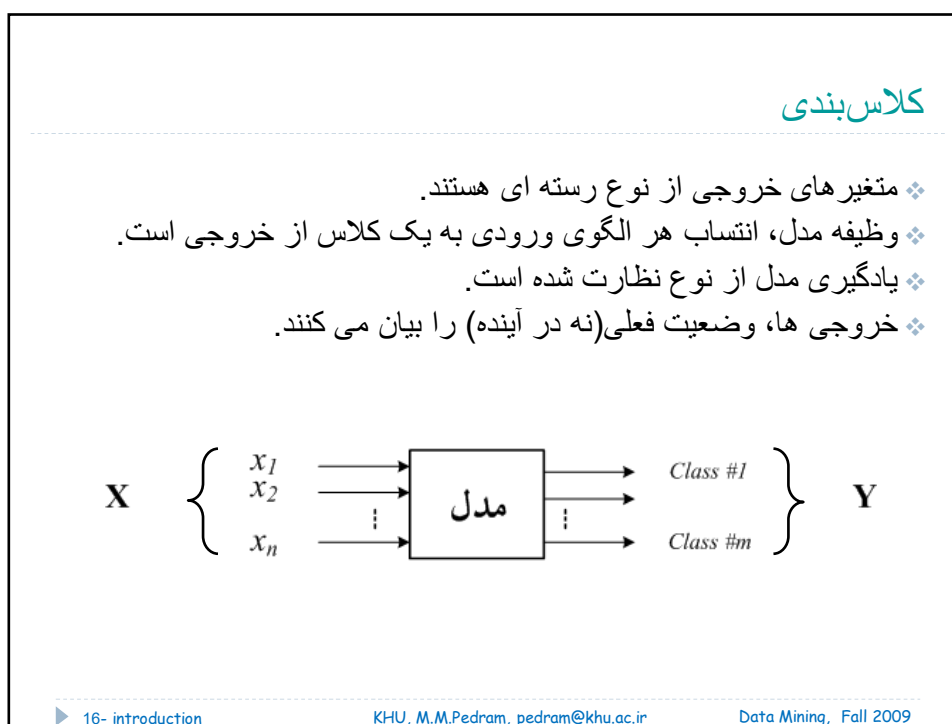
- ❖ **Clustering:** Segmentation, Database Segmentation
- ❖ **Link Analysis:** Discovering Association Rules, Sequential Patterns, and time Sequences
- ❖ **Deviation Detection:** Visualization, Statistics

## Fundamental Types of Learning





15



16

## ❖ روش های کلاس بندی

- ❑ درخت تصمیم گیری
- ❑ قواعد
- ❑ شبکه های عصبی
- ❑ گونه های فازی

► 17- introduction

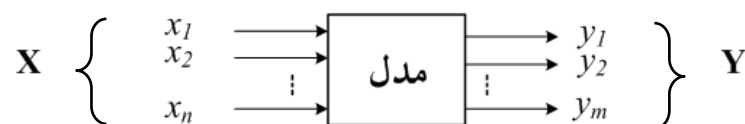
KHU, M.M.Pedram, [pedram@khu.ac.ir](mailto:pedram@khu.ac.ir)

Data Mining, Fall 2009

17

## تخمین

- ❖ متغیرهای خروجی از نوع عددی هستند.
- ❖ وظیفه مدل، انتساب یک عدد در خروجی به هر الگوی ورودی است.
- ❖ یادگیری مدل از نوع نظارت شده است.
- ❖ خروجی ها، وضعیت فعلی (نه در آینده) را بیان می کنند.



► 18- introduction

KHU, M.M.Pedram, [pedram@khu.ac.ir](mailto:pedram@khu.ac.ir)

Data Mining, Fall 2009

18

## ❖ روش‌های تخمین

- ❑ قواعد
- ❑ شبکه‌های عصبی
- ❑ رگرسیون آماری
- ❑ گونه‌های فازی
- ❑ سری‌های زمانی

► 19- introduction

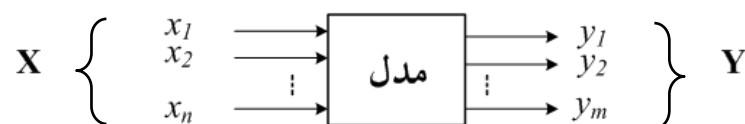
KHU, M.M.Pedram, [pedram@khu.ac.ir](mailto:pedram@khu.ac.ir)

Data Mining, Fall 2009

19

## پیش‌بینی

- ❖ متغیرهای خروجی از نوع عددی ویا رشته‌ای هستند.
- ❖ وظیفه مدل، انتساب یک عدد یا کلاس در هر خروجی به هر الگوی ورودی است.
- ❖ یادگیری مدل از نوع نظارت شده است.
- ❖ خروجی‌ها، وضعیت در آینده را بیان می‌کنند.



► 20- introduction

KHU, M.M.Pedram, [pedram@khu.ac.ir](mailto:pedram@khu.ac.ir)

Data Mining, Fall 2009

20

## ❖ روش های پیش بینی

- ❑ درخت تصمیم گیری
- ❑ قواعد
- ❑ شبکه های عصبی
- ❑ رگرسیون آماری
- ❑ گونه های فازی
- ❑ سری های زمانی

## خوشه بندی

- ❖ متغیرهای خروجی از نوع رسته ای هستند.
- ❖ وظیفه مدل، کشف و استخراج ساختارهایی در داده های ورودی.
- ❖ یادگیری مدل از نوع بدون نظارت است.

## کاربردها

- تعیین روابط معنادار در داده ها.
- ارزیابی عملکرد یک مدل که یادگیری نظارت شده، دارد.
- تعیین بهترین مجموعه از متغیرهای ورودی برای یادگیری با نظارت.
- آشکارسازی برون هشته ها (outliers).
- برون هشته: داده ای که به طور عادی با سایر نمونه ها، گروه بندی نمی شود.

❖ روش‌های خوشه‌بندی

- آماری (k-means و c-means)
- شبکه‌های عصبی خودسازمان‌ده
- گونه‌های فازی

## تحلیل سبد بازار (Market Basket Analysis)

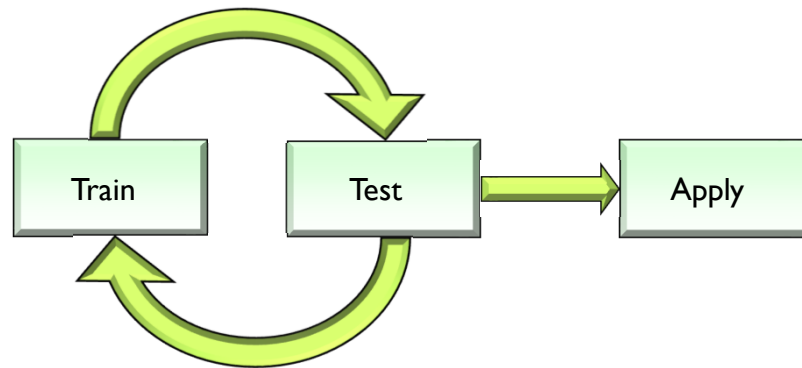
یافتن ارتباطات جالب بین محصولات خریده‌فروشی.

❖ روش‌های تحلیل سبد بازار

- قواعد هم‌باشی (association rules)

## Model Development

Training and Testing (Model validation)



► 25- introduction

KHU, M.M.Pedram, [pedram@khu.ac.ir](mailto:pedram@khu.ac.ir)

Data Mining, Fall 2009

25

## Evaluation of classifiers

❖ For any data set that is used to test a classifier (model), a *confusion matrix* can be built.

❖ Classification

- Binary class
- Multi-class

► 26

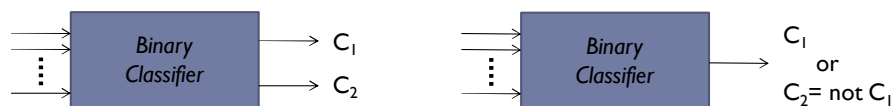
KHU, M.M.Pedram, [pedram@khu.ac.ir](mailto:pedram@khu.ac.ir)

Data Mining, Fall 2009

26

## Evaluation of classifiers

### ❖ Binary classification



		Classifier Output	
		Positive	Negative
Actual	Positive	124	15
	Negative	8	84

*agreement* (circled around the 124 and 84 values)

- ▶ **Error** = "proportion of incorrect classification"  
 $= (8 + 15) / (124 + 84 + 8 + 15) = 23/231$
- ▶ **Accuracy** = 1 - Error = "proportion of correct predictions" =  $208/231$
- ▶ **Precision** = "proportion of predicted positive cases that were correct"  
 $= 124 / (124 + 8)$

▶ 27

KHU, M.M.Pedram, [pedram@khu.ac.ir](mailto:pedram@khu.ac.ir)

Data Mining, Fall 2009

27

## Evaluation of classifiers

- ▶ Entries in the confusion matrix for binary classification have names:

		Classifier Output	
		Positive	Negative
Actual	Positive	124 → TP	15 → FN
	Negative	8 → FP	84 → TN

**TP = hit** = "number of positive cases correctly identified" = 124

**FN = a miss** = "number of positive cases incorrectly classified as negative" = 15

**FP = false alarm** = "number of negative cases incorrectly classified as positive" = 8

**TN = correct rejection** = "number of negative cases correctly identified" = 84

▶ 28

KHU, M.M.Pedram, [pedram@khu.ac.ir](mailto:pedram@khu.ac.ir)

Data Mining, Fall 2009

28

## Evaluation of classifiers

- Entries in the confusion matrix for binary classification have names:

		Classifier Output	
		Positive	Negative
Actual	Positive	124 → TP	15 → FN
	Negative	8 → FP	84 → TN

**Recall = Sensitivity = SENS = TP rate** = “proportion of positive cases correctly identified” =  $TP/(TP+FN) = 124 / (124 + 15)$

**FN rate** = “proportion of positive cases incorrectly classified as negative”  
 $= 15 / (124 + 15)$

**FP rate** = “proportion of negative cases incorrectly classified as positive” =  $8 / (84 + 8)$

**Specificity = SPEC = TN rate** = “proportion of negative cases correctly identified”  
 $= 84 / (84 + 8)$

▶ 29

KHU, M.M.Pedram, [pedram@khu.ac.ir](mailto:pedram@khu.ac.ir)

Data Mining, Fall 2009

29

## Evaluation of classifiers

- Entries in the confusion matrix for binary classification have names:

		Classifier Output		
		Positive	Negative	
Actual	Positive	124 → TP	15 → FN	Recall or SENS
	Negative	8 → FP	84 → TN	SPEC
		Precision	NPV	

**Recall = Sensitivity = SENS = TP rate** =  $TP/(TP + FN) = 124 / (124 + 15)$

**Precision = PPT = Positive Predictive value** = “proportion of correctly classified ones as positive case to the all ones classifies as positive”  
 $= TP/(TP + FP) = 124 / (124 + 8)$

**NPV = Negative Predictive value** = “proportion of correctly classified ones as negative case to the all ones classifies as negative” =  $84 / (84 + 15)$

**Specificity = SPEC = TN rate** =  $TN/(TN + FP) = 84 / (84 + 8)$

▶ 30

KHU, M.M.Pedram, [pedram@khu.ac.ir](mailto:pedram@khu.ac.ir)

Data Mining, Fall 2009

30

## Evaluation of classifiers

### ❖ Multi-Class classification



		Classifier Output		
		$C_1$	$C_2$	$C_3$
Actual	$C_1$	140	20	22
	$C_2$	17	54	8
	$C_3$	12	4	76

agreement

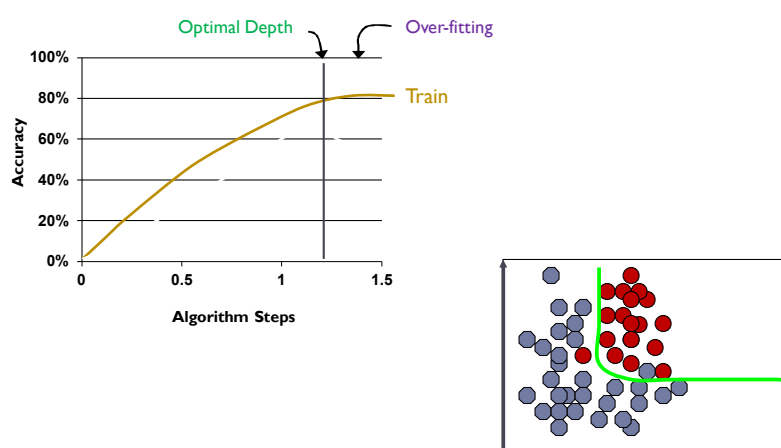
▶ 31

KHU, M.M.Pedram, [pedram@khu.ac.ir](mailto:pedram@khu.ac.ir)

Data Mining, Fall 2009

31

## Over-fitting



❖ Running too many epochs can result in over-fitting.

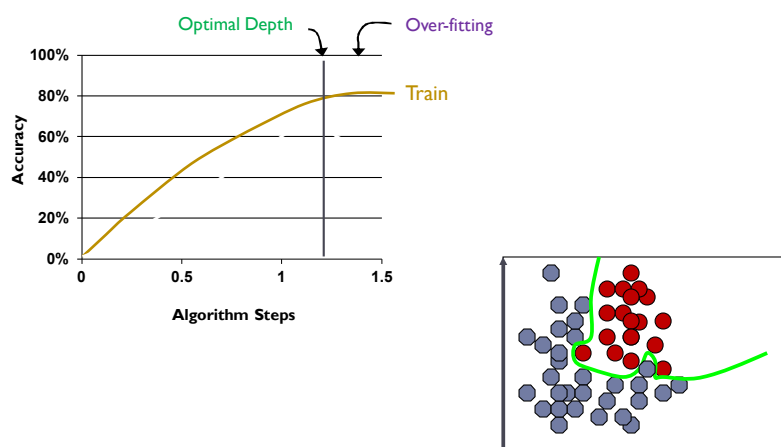
▶ 32

KHU, M.M.Pedram, [pedram@khu.ac.ir](mailto:pedram@khu.ac.ir)

Data Mining, Fall 2009

32

## Over-fitting



❖ Running too many epochs can result in over-fitting.

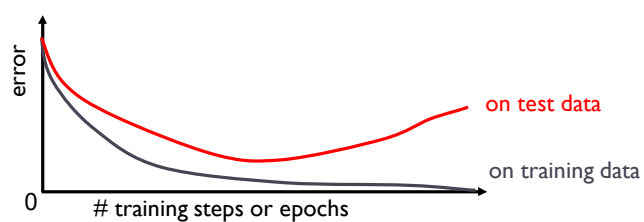
► 33

KHU, M.M.Pedram, [pedram@khu.ac.ir](mailto:pedram@khu.ac.ir)

Data Mining, Fall 2009

33

## Over-fitting



► 34

KHU, M.M.Pedram, [pedram@khu.ac.ir](mailto:pedram@khu.ac.ir)

Data Mining, Fall 2009

34

## Scientific vs. Commercial Data Mining

### Goals:

- ▶ **Science – Theories:** Need for insight and theory-based models, interpretable model structures, generate domain rules or causal structures, support for theory development
- ▶ **Commercial – Profits:** black boxes OK

### Types of data:

- ▶ **Science** – Images, sensors, simulations
- ▶ **Commercial** - Transaction data
- ▶ **Both** - Spatial and temporal dimensions, heterogeneous

### IT (information technology) tools fit both enterprises

- ▶ **Database systems** (Oracle, DB2, etc),
- ▶ **integration tools** (Information Integrator),
- ▶ **web services** (Blue Titan, .NET)

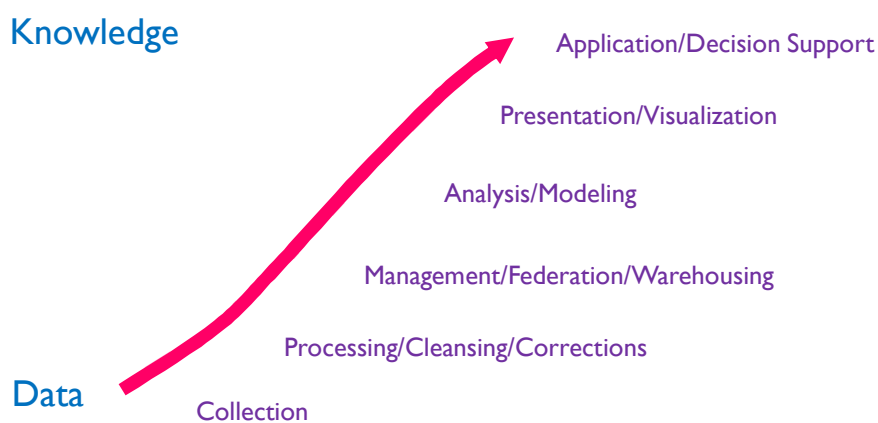
▶ 35- introduction

KHU, M.M.Pedram, [pedram@khu.ac.ir](mailto:pedram@khu.ac.ir)

Data Mining, Fall 2009

35

## Knowledge Discovery Process

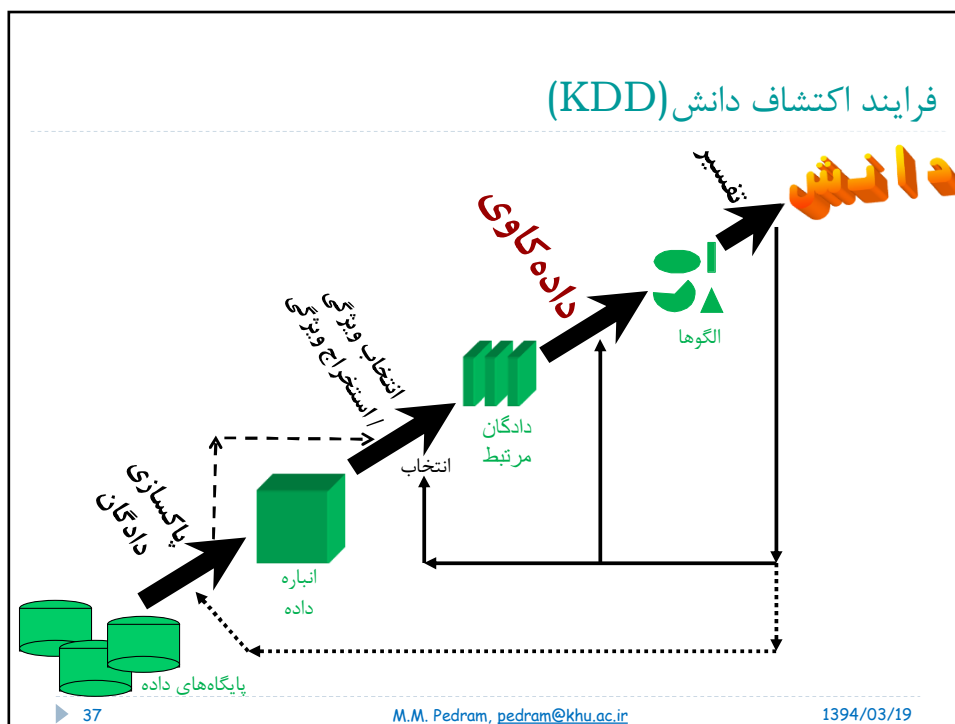


▶ 36- introduction

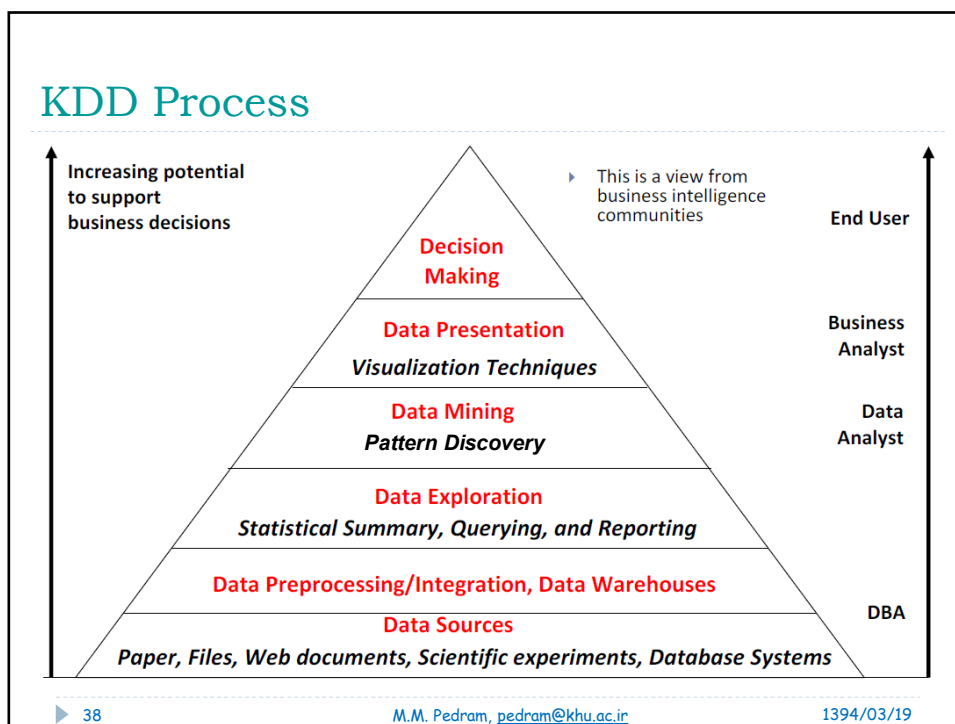
KHU, M.M.Pedram, [pedram@khu.ac.ir](mailto:pedram@khu.ac.ir)

Data Mining, Fall 2009

36

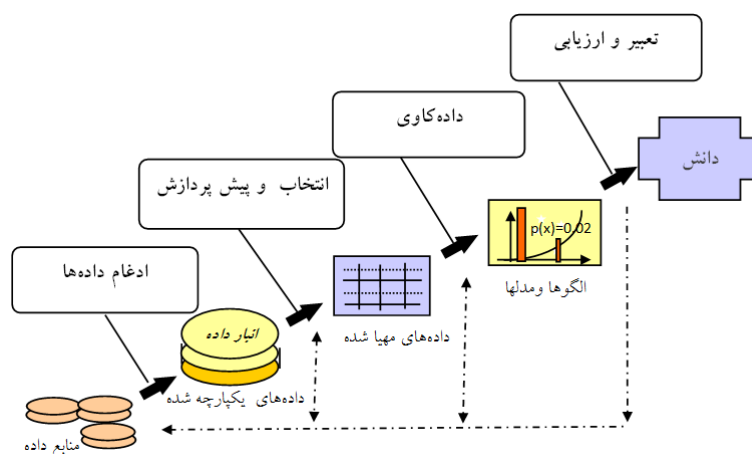


37



38

## اکتشاف دانش در پایگاه داده ها (KDD)



► 39- introduction

KHU, M.M.Pedram, pedram@khu.ac.ir

Data Mining, Fall 2009

39

### 1. تشکیل انبار داده:

- با توجه به عنوان، این مرحله برای تشکیل محیطی پیوسته و یکپارچه جهت انجام مراحل بعدی و داده کاوی در آن، انجام می‌گیرد.
- انبار داده مجموعه پیوسته و طبقه‌بندی شده است که دائماً در حال تغییر بوده و پویا است و برای کاوش آماده می‌شود.

### 2. انتخاب داده‌ها و پیش‌پردازش:

- در این مرحله برای کم کردن هزینه‌های عملیات داده‌کاوی، داده‌هایی از پایگاه داده انتخاب می‌شوند که مورد مطالعه هستند و هدف داده‌کاوی دادن نتایجی در مورد آن‌هاست.
- برای انجام عملیات داده‌کاوی لزوماً باید تبدیلات خاصی روی داده‌ها انجام گیرد ممکن است این تبدیلات خیلی راحت و مختصر مثل تبدیل *byte* به *integer* باشد یا خیلی پیچیده و زمان‌بر و با هزینه‌های بالا مثل تعریف صفات جدید و یا تبدیل و استخراج داده‌ها از مقادیر رشته‌ای و ... باشد.

► 40- introduction

KHU, M.M.Pedram, pedram@khu.ac.ir

Data Mining, Fall 2009

40

### 3. داده‌کاوی:

➤ در این مرحله با استفاده از تکنیک‌های داده‌کاوی داده‌ها مورد کاوش قرار گرفته، دانش نهفته در آنها استخراج شده و الگوسازی صورت می‌گیرد.

### 4. تفسیر نتیجه:

➤ در این مرحله نتایج و الگوهای ارائه شده توسط ابزار داده‌کاوی مورد بررسی قرار گرفته و نتایج مفید معین می‌شود و در صورت عدم دقت، روند تکرار می‌شود.