

# مقایسه حفاری داده‌ها با تحلیل داده‌ها

Data Analysis and Data Mining Comparison

عادل محمدپور

دانشگاه صنعتی امیرکبیر (پلی‌تکنیک تهران)

AdelM.ir

adel@aut.ac.ir

# مفهوم داده‌کاوی

عبارت داده‌کاوی مترادف با یکی از عبارت‌های استخراج دانش، برداشت اطلاعات، و ارسای داده‌ها

و حتی لایروبی کردن داده‌ها است که در حقیقت کشف دانش در پایگاه داده‌ها را توصیف می‌کند

- Knowledge Discovery of Database (KDD)

# کشف دانش در پایگاه داده‌ها

**Knowledge Discovery in Data** is the *non-trivial* process of identifying

- *valid*
- *novel*
- *potentially useful*
- and ultimately *understandable patterns* in data.

From *Advances in Knowledge Discovery and Data Mining*, Fayyad, Piatetsky-Shapiro, Smyth, and Uthurusamy, (Chapter 1), AAAI/MIT Press 1996

# تعاريف داده‌کاوي

- داده‌کاوي در حقيقت کشف ساختارهاي جالب توجه، غيرمنتظره و با ارزش از داخل مجموعه وسيعي از داده‌ها مي‌باشد و فعاليتي است که اساساً با آمار و تحليل دقيق داده‌ها منطبق است
- داده‌کاوي فرايند کشف رابطه‌ها، الگوها و روندهاي جديد معني‌داري است که به بررسي حجم وسيعي از اطلاعات ذخيره شده در انبارها با فناوري‌هاي تشخيص الگو (مانند رياضي و آمار) مي‌پردازد
- داده‌کاوي يا به تعبير ديگر کشف دانش در پايگاه داده‌ها، استخراج غير بديهي اطلاعات بالقوه مفيد از روي داده‌هايي است که قبلاً ناشناخته مانده‌اند. اين مطلب برخي از روش‌هاي فني مانند خوشه‌بندي، خلاصه‌سازي داده‌ها، فراگيري قاعده‌هاي رده‌بندي، يافتن ارتباط شبکه‌ها، تحليل تغييرات و کشف بي‌قاعدهگي‌ها را شامل مي‌شود

چرا حفاري داده‌ها بجاي تحليل داده‌ها

# داده‌های بیشتری تولید می‌شوند

Bank, telecom, other business transactions ...

Scientific Data: astronomy, biology, etc

Web, text, and e-commerce

# داده‌های بیشتری ذخیره می‌شوند

Storage technology faster and cheaper

DBMS capable of handling bigger DB

# مثال

Europe's Very Long Baseline Interferometry (VLBI) has 16 telescopes, each of which produces **1 Gigabit/second** of astronomical data over a 25-day observation session  
storage and analysis a big problem

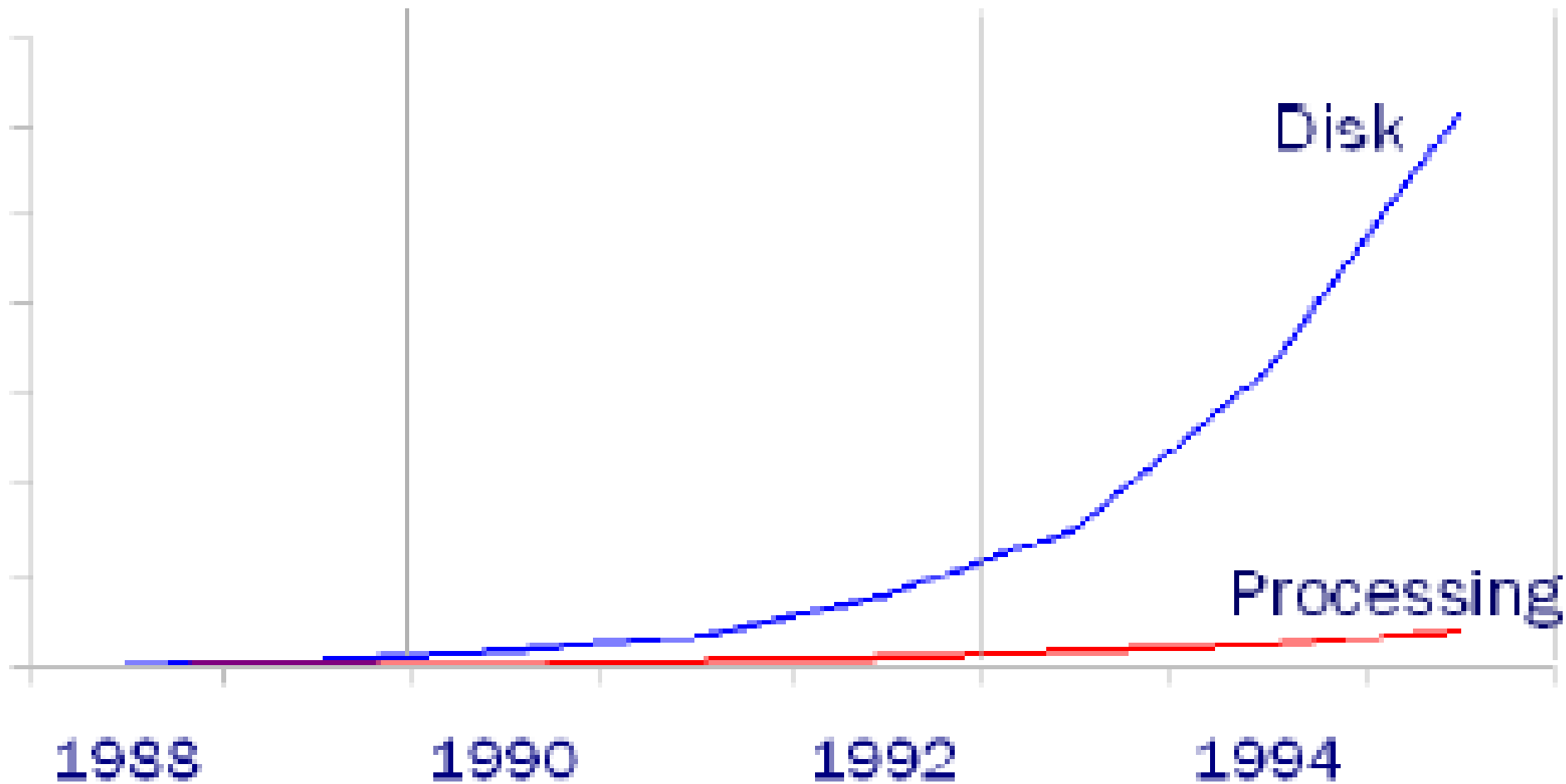
Walmart reported to have 24 Tera-byte DB

AT&T handles billions of calls per day

data cannot be stored -- analysis is done on the fly



# چرا حفاري داده‌ها بجاي تحليل داده‌ها



Knowledge Discovery in Data is the *non-trivial* process of identifying

- *valid*

- *novel*

- potentially *useful*

- and ultimately *understandable patterns* in data.

# ارتباط داده‌کاوی با علوم مختلف

**Machine Learning**

**Visualization**

**Data Mining and Knowledge Discovery**

**Statistics**

**Databases**

- **هوش مصنوعی** در زمینه تحلیل آماری، هوش مصنوعی عبارت از شیوه‌های خودکار برای به کارگیری روش‌های آماری در تحلیل داده‌ها است

- الگوریتم **یادگیری ماشین**، ابزاری است که می‌تواند آزمایش‌ها (مثال‌های مشاهده شده) را با در نظر گرفتن رده بعضی از توابع و معیارهای اندازه‌گیری، یاد بگیرد

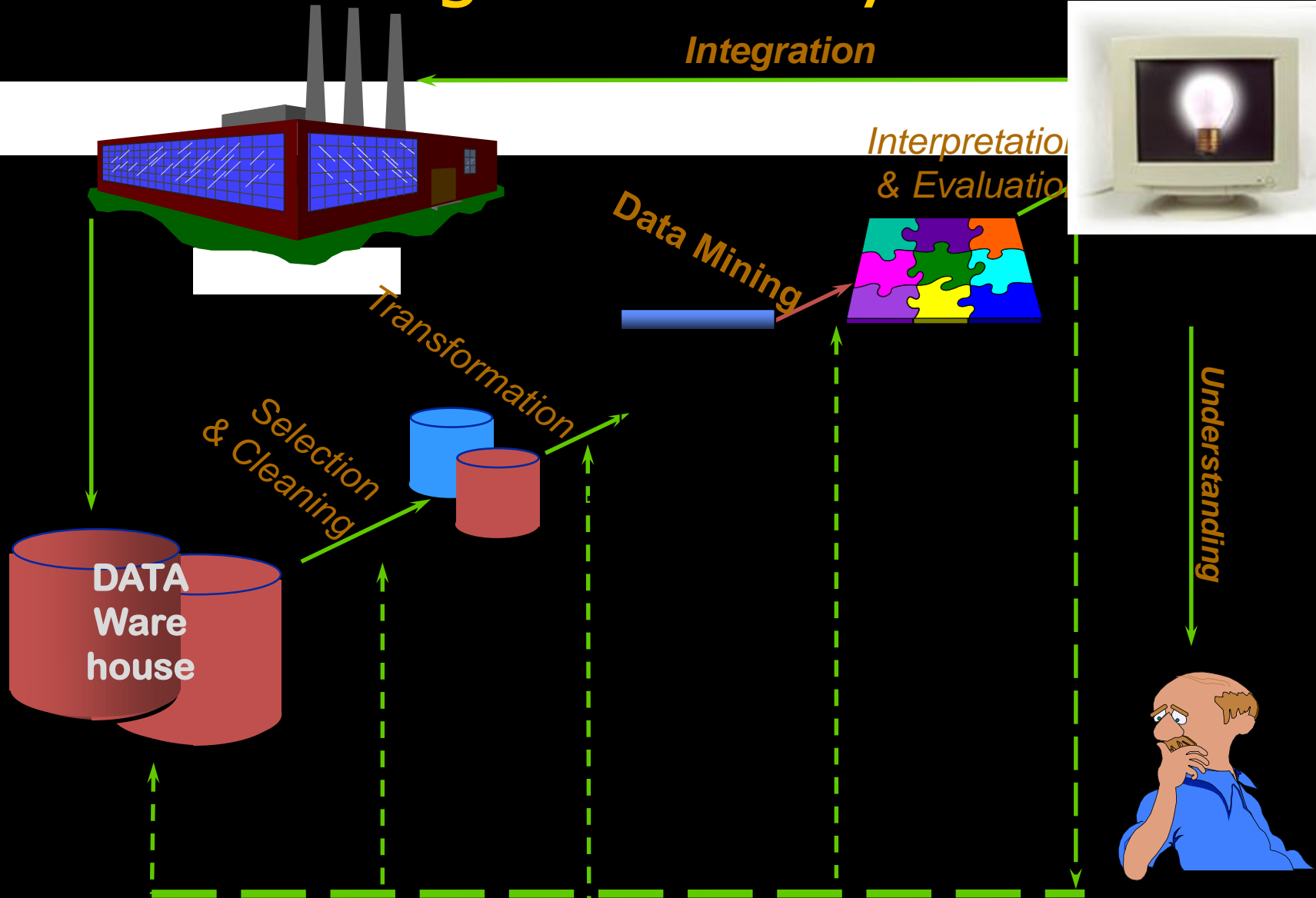
# گام‌های داده‌کاوی

- (1) درک قلمرو
- (2) آماده کردن مجموعه داده‌ها
- (3) کشف الگوها (داده‌کاوی)
- (4) پردازش بعد از کشف الگو
- (5) استفاده از نتایج

# گام‌های داده‌کاوی

- (1) بیان مسئله
- (2) انتخاب داده
- (3) آماده‌سازی داده‌ها
- (4) مدل‌سازی
- (5) ارزیابی مدل
- (6) توسعه مدل

# Knowledge Discovery Process



# آمار و داده‌کاوی

- نوع داده‌ها در آمار کمی است ولی داده‌کاوی داده‌های کمی، کیفی و متنی را پوشش می‌دهد.
- نوع متغیرهای ورودی در آمار عددی ولی در داده‌کاوی متغیرها از نوع عددی، طبقه‌ای و متنی می‌باشند.
- نوع تمرکز در آمار روی مدل ولی در داده‌کاوی تمرکز روی الگو است.
- زیربخش‌های اصلی آمار برآورد، توزیع‌های احتمال، آزمون فرضیه، امتیازبندی مدل و پیشگویی است، ولی زیربخش‌های داده‌کاوی مدل‌بندی پیشگویی‌ها، بخش‌بندی پایگاه داده‌ها و ساخت فرضیه است.



# آمار و داده‌کاوی

- در آمار جستجو برای دستیابی به نتایج محدود به جستجوهای جهت‌دار بوده و با نتایج نیز آشنا هستیم. ولی در داده‌کاوی جستجوها اکثراً از طریق کاربر تعیین شده و ممکن است جهت‌دار باشند، اما اصولاً روش خودکار است و نوع نتایج نامعلوم است.
- در مقایسه با آمار، داده‌کاوی توجه کمتری به ویژگیهای جانبی استنباط‌های بزرگ نمونه‌ای دارد و فلسفه کلی یادگیری، بیشتر شامل ملاحظه پیچیدگی مدل‌ها و محاسباتی که آنها نیاز دارند، است.

# آمار و داده‌کاوی

- هدف اصلی از به کارگیری روش‌ها در آمار استفاده از برآوردها و توزیع‌ها برای ادغام اطلاعات، ولی در داده‌کاوی هدف اصلی کشف دانش مورد علاقه است.
- شکل نتایج در آمار به صورت مدل‌های کلی برآورد می‌شود ولی در داده‌کاوی مدل‌های موضعی محاسبه می‌شوند.
- ارزش اطلاعاتی نتایج در آمار معلوم و محدود است در حالی که در داده‌کاوی نامعلوم و نامحدود می‌باشد.

# آمار و داده‌کاوی

- با يك فرضيه شروع مي‌شود
- به فرضيه احتياجي ندارد

- آمارشناسان بايد رابطه‌هايي را ايجاد کنند که به فرضيه مربوط شود
- الگوريتم‌هاي داده‌کاوي خودکار روابط را ايجاد مي‌کند

- داده‌ها عددي هستند
- از هر نوع داده‌اي مي‌توانند استفاده کنند

# آمار و داده‌کاوی

- داده‌های نابجا و نادرست ضمن تحلیل تشخیص داده می‌شوند
- داده‌کاوی به داده‌های صحیح و درست طبقه‌بندی شده بستگی دارد
- نتایج کار آمارشناسان برای مدیران نیاز به تفسیر ندارد
- تفسیر نتایج داده‌کاوی آسان نیست و برای تحلیل و بیان آن‌ها به متخصصان آمار نیاز است

# آمار و داده‌کاوی

- نمونه‌گیری
- کم کردن هزینه
- فروکاهی داده‌ها
- مخرب بودن
- اندازه زیاد داده‌ها
- زمان بر بودن سرشماری
- ضروری نبودن ذخیره آن‌ها
- هزینه محاسباتی
- محدودیت سخت افزار و نرم افزار
- زمان بر بودن تحلیل