

# پیش پردازش داده‌ها

جواد وحدت

[vahdatjavad@gmail.com](mailto:vahdatjavad@gmail.com)

# این مراحل جزء مهمترین گام‌های تحلیل داده محسوب می‌شوند

به طور کلی پیش پردازش داده به دو دسته عمده تقسیم می‌شود که عبارتند از:

- انتخاب اشیا داده و ویژگی‌ها (attribute) برای تحلیل
- ایجاد کردن یا تغییر دادن ویژگی‌ها

# ۱-یکپارچه سازی داده ها – Data Integration

شامل ترکیب داده ها از چند منبع غیر همگون، انتقال آن ها به یک انبار منسجم و ایجاد یک دید یکتا به داده ها.

برای ترکیب داده ها از دو روش کمک میگیریم :

- **اتصال محکم** Tight Coupling – از این به عنوان مولفه بازیابی اطلاعات استفاده میکنیم.

- **اتصال سست** Loose Coupling – در این روش رابطی تهیه می گردد که در آن query را از کاربر گرفته و به شیوه ای که پایگاه داده مبدا متوجه آن شود، تبدیل می کند.

## ۲-تجميع داده – Data Aggregation

به تركيب متغيرها باهم و ادغام در يك متغير تجميع داده گفته مي شود.  
به منظور:

- کاهش داده
- تغيير مقياس ( به عنوان مثال: شهرها مي توانند به مناطق و استان ها و يا کشورها تبديل شوند و يا روزها به هفته ها و ماه ها و سال ها تبديل شوند.)
- ثبات بيشتر داده

## ۳-نمونه گیری - Sampling

**نمونه گیری** یک روش معمول برای انتخاب یک **زیر مجموعه** از داده ها برای تحلیل داده است.

- نمونه گیری تصادفی عادی Simple Random Sampling
- نمونه گیری بدون جایگزینی Sampling without replacement
- نمونه گیری با جایگزینی Sampling with replacement
- نمونه گیری طبقه بندی شده Stratified sampling
- نمونه گیری تصاعدی Progressive Sampling

## ۴- کاهش ابعاد – Dimension Reduction

Curse of Dimensionality – خطای classify کردن دیتا افزایش پیدا میکند.

- آنالیز مولفه اصلی – PCA

- تجزیه مقدارهای منفرد – SVD

## ۵-انتخاب زیر مجموعه ای از ویژگی ها – Feature Subset Selection

- ویژگی های غیر مرتبط Irrelevant

- ویژگی های زائد Redundant

سه راه استاندارد برای انتخاب زیرمجموعه ها وجود دارد:

- روش های تعبیه شده Embedded (به عنوان بخشی از الگوریتم داده کاوی اتفاق می افتد ، به طور دقیق خود الگوریتم تصمیم میگرد که از چه ویژگی هایی استفاده کند و چه ویژگی هایی استفاده نکند.)

- روش های مبتنی بر فیلتر Filter (انتخاب ویژگی با استفاده از روش هایی مستقل از عملیات داده کاوی قبل از اجرای الگوریتم های داده کاوی.)

- روش های بسته بندی Wrapper (از الگوریتم هدف به عنوان یک جعبه سیاه برای پیدا کردن بهترین مجموعه از ویژگی ها استفاده می کنند.)

## ۶-خلق ویژگی – Feature Creation

- استخراج ویژگی
- ساخت ویژگی
- نگاشت داده ها به فضای جدید



## ۷- گسسته سازی و دوگانه سازی – Discretization and Binarization

**گسسته سازی** Discretization : گسسته سازی روند تبدیل یک متغیر پیوسته به یک متغیر ترتیبی است.

**دوگانه سازی** Binarization : دوگانه سازی یک ویژگی پیوسته یا طبقه ای categorical را به یک یا چند متغیر دودویی تبدیل می کند.

Categorical Value	Integer Value	$x_1$	$x_2$	$x_3$
<i>awful</i>	0	0	0	0
<i>poor</i>	1	0	0	1
<i>OK</i>	2	0	1	0
<i>good</i>	3	0	1	1
<i>great</i>	4	1	0	0

# ۸-تغییر شکل متغیر ها – Variable Transformation

- تابعی ایست که تمامی مقادیر ویژگی مورد نظر را به مقادیر جایگزینی تبدیل می کند. به نحوی که هر کدام از مقادیر قبلی با یکی از مقادیر جدید مشخص می شود.
- توابع ساده مثل توان، لگاریتم، قدر مطلق
- نرمال سازی – برای تنظیم اختلاف میان ویژگی ها به لحاظ تناوب رخداد، میانگین، واریانس و بازه به کار می روند
- استاندارد سازی – در آمار به اختلاف میانگین ها تقسیم بر میانگین واریانس ها اشاره دارد.