

CISC 372

Advanced Data Analytics

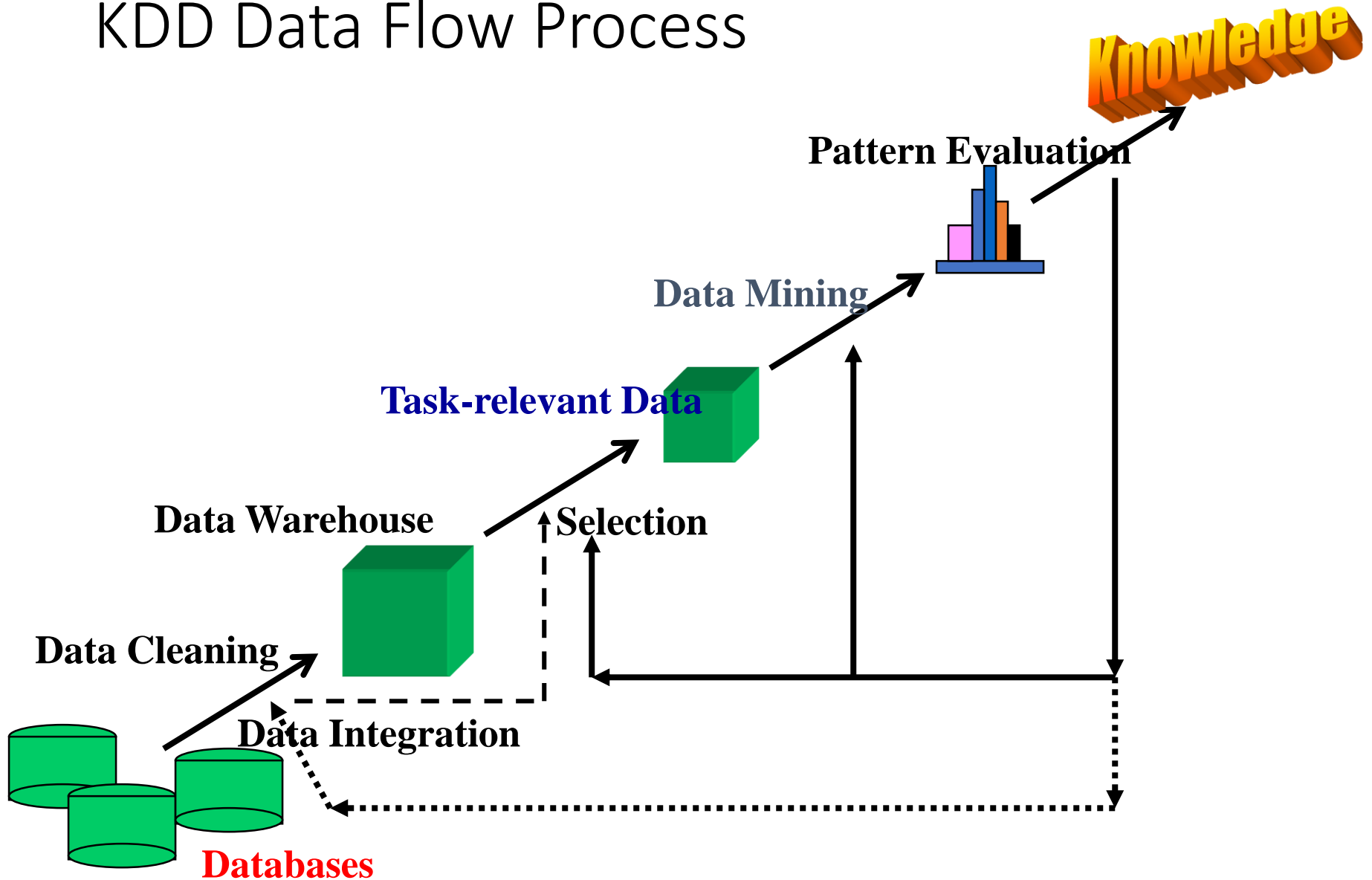
L5 – Infrastructure & Preprocessing

<https://l1nna.com/372>

Today

- KDD Process
- Data Attributes
- Data Characteristics
- Data Preprocessing

KDD Data Flow Process



Data Objects

- Data sets are made up of data objects.
- A **data object** represents an entity.
- Examples:
 - sales database: customers, store items, sales
 - medical database: patients, treatments
 - university database: students, professors, courses
- Also called *samples* , *examples*, *instances*, *data points*, *objects*, *tuples*, and *records*.
- Data objects are described by **attributes**.
- Database rows -> data objects; columns -> attributes.

Attributes

- **Attribute (or dimensions, features, variables)**: a data field, representing a characteristic or feature of a data object.
 - *E.g., customer_ID, name, address*
- Types:
 - Nominal (aka **categorical**)
 - Binary (aka **binominal**)
 - Ordinal
 - Continuous

Attribute Types

- **Nominal:** categories, states, or “names of things”
 - *Hair_color* = {auburn, black, blond, brown, grey, red, white}
 - marital status, occupation
- **Binary** (aka binominal)
 - Nominal attribute with only 2 states (0 and 1)
 - Symmetric binary: both outcomes equally important
 - e.g., gender
 - Asymmetric binary: outcomes not equally important.
 - e.g., medical test (positive vs. negative)
 - Convention: assign 1 to most important outcome (e.g., HIV positive)

Attribute Types

- **Ordinal**

- Values have a meaningful order (ranking) but magnitude between successive values is not known.
- *Size = {small, medium, large}*, grades, army rankings

- **Continuous Attribute**

- Has real numbers as attribute values
 - E.g., temperature, height, or weight
- Practically, real values can only be measured and represented using a finite number of digits

Types of Data Sets

- Record
 - Relational records
 - Data matrix, e.g., numerical matrix, crosstabs
- Transaction data
- Document data
- Graph and network
 - World Wide Web
 - Social or information networks
 - Molecular Structures
- Ordered
 - Video data: sequence of images
 - Temporal data: time-series
 - Sequential Data: transaction sequences
 - Genetic sequence data
- Spatial, image and multimedia:
 - Spatial data: maps
 - Image data
 - Video data

Diagram illustrating an EMPLOYEE table with annotations:

Employee ID	First name	Last name	Title	Date of birth	Address ID
1008	Jane	Smith	Sales	2/14/1969	302
1009	Joe	Diner	Clerk	4/16/1968	884
1010	Ed	Smithers	Sales	9/22/1964	992
1011	Tom	Massee	Mgr	4/5/1947	42
1012	Julie	Vahnne	Clerk	6/11/1972	223
1013	John	Smith	Clerk	3/12/1970	302
1014	Donald	Winter	Clerk	5/18/1969	55

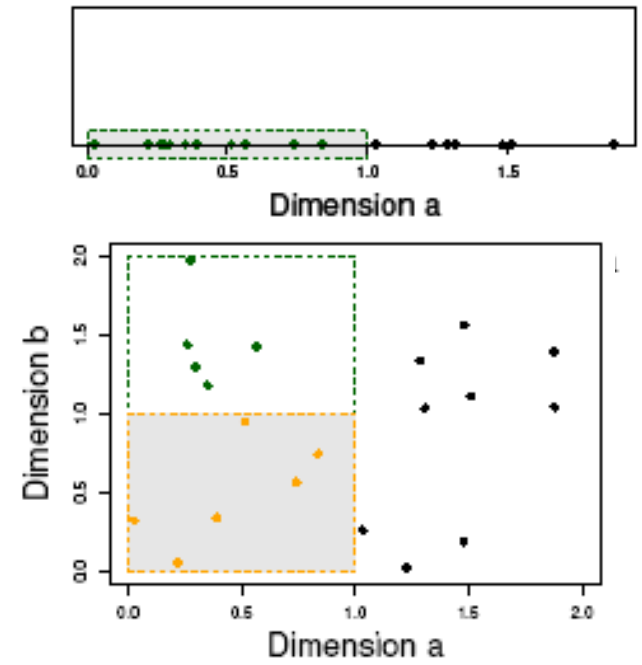
Annotations:

- Record:** Points to the entire row for Employee ID 1010.
- Field:** Points to individual cells (e.g., 'Ed', 'Smithers', 'Sales', '9/22/1964', '992').
- EMPLOYEE table:** Points to the entire table structure.

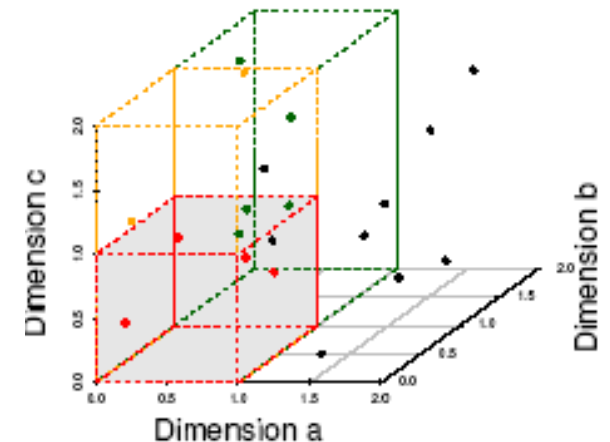
Month	REGION 1	REGION 2	REGION 3	REGION 4	REGION 5	TOTAL
April	13	33	76	2	47	171
May	17	55	209	1	143	425
June	8	63	221	1	127	420
July	13	104	240	6	123	486
August	18	121	274	9	111	533
September	25	160	239	2	88	514
October	9	88	295	2	127	521
November	2	86	292	2	120	502
December	1	128	232	6	155	522
TOTAL	106	838	2078	31	1041	4094

Important Characteristics of Structured Data

- Dimensionality
 - Curse of dimensionality
- Sparsity
 - Only presence counts
- Resolution
 - Patterns depend on the scale
- Distribution
 - Centrality and dispersion



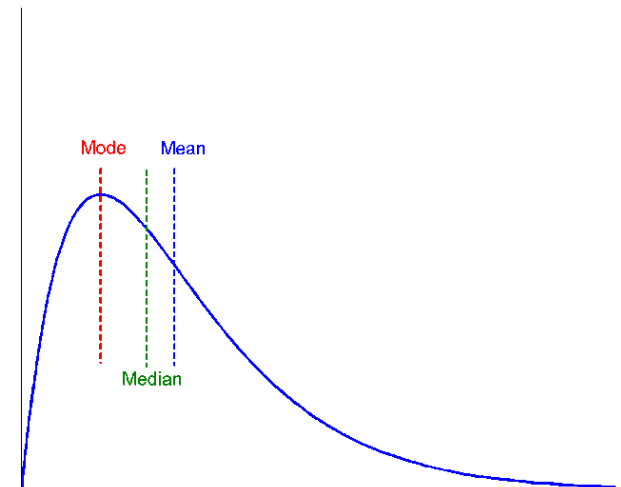
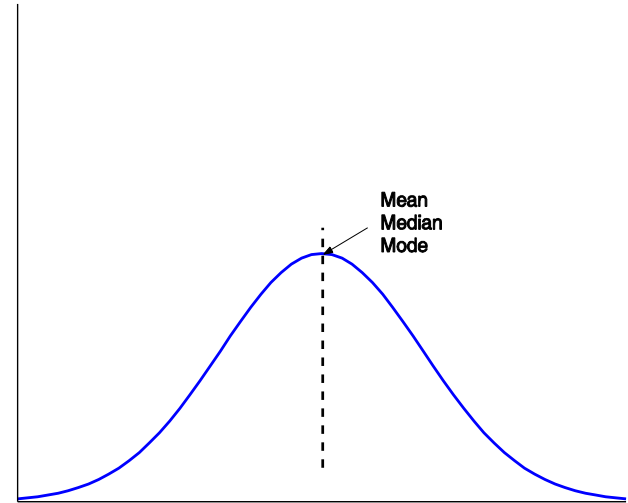
(b) 6 Objects in One Unit Bin



(c) 4 Objects in One Unit Bin

Important Characteristics of Structured Data

- Dimensionality
 - Curse of dimensionality
- Sparsity
 - Only presence counts
- Resolution
 - Patterns depend on the scale
- Distribution
 - Centrality and dispersion



Data Does not Come Prepared (mostly)

- Data cleaning
 - Preprocess data in order to reduce noise and handle missing values
- Relevance analysis (feature selection)
 - Remove the irrelevant or redundant attributes
 - e.g., do not need both age and birth year
- Data transformation
 - Generalize and/or normalize data
 - e.g. do not need to know specific address of a customer, city or postal code is good enough.

Data Does not Come Prepared (mostly)

- Data cleaning
 - Preprocess data in order to reduce noise and handle missing values
- Relevance analysis (feature selection)
 - Remove the irrelevant or redundant attributes
 - e.g., do not need both age and birth year
- Data transformation
 - Generalize and/or normalize data
 - e.g. do not need to know specific address of a customer, city or postal code is good enough.

Normalization

- Scaling individual samples to have unit norm.
- Could be L1 norm or L2 norm
 - Very important for vector space model (think about dot-product)

```
1 from sklearn import preprocessing
2 import numpy as np
3 from numpy.linalg import norm
4
5 X_train = np.array([[ 1., -1.,  2.],
6                     [ 2.,  0.,  0.],
7                     [ 0.,  1., -1.]])
8 X_scaled = preprocessing.normalize(X_train, norm='l2')
9
10 print(X_scaled)
11 print('norm', norm(X_scaled, axis=1))
12 print('std ', X_scaled.std(axis=0))
13 print('mean', X_scaled.mean(axis=0))
```



```
[[ 0.40824829 -0.40824829  0.81649658]
 [ 1.          0.          0.          ]
 [ 0.          0.70710678 -0.70710678]]
norm [1. 1. 1.]
std  [0.41053309 0.46075826 0.62254262]
mean [0.4694161  0.0996195  0.03646327]
```

Standardization

- A common requirement
 - Making the feature a Gaussian **with zero mean** and **unit variance**

```
1 from sklearn import preprocessing
2 import numpy as np
3 X_train = np.array([[ 1., -1.,  2.],
4                     [ 2.,  0.,  0.],
5                     [ 0.,  1., -1.]])
6 X_scaled = preprocessing.scale(X_train)
7
8 print(X_scaled)
9 print('mean', X_scaled.mean(axis=0))
10 print('std ', X_scaled.std(axis=0))
```

```
↳ [[ 0.          -1.22474487  1.33630621]
    [ 1.22474487  0.          -0.26726124]
    [-1.22474487  1.22474487 -1.06904497]]
mean [0. 0. 0.]
std  [1. 1. 1.]
```

Categorical Encoding

- Converting categorical attribute (aka nominal) into numeric features (or a one-hot encoding vector)

```
1 from sklearn import preprocessing
2 import numpy as np
3 from numpy.linalg import norm
4
5 enc = preprocessing.OrdinalEncoder()
6 X = [
7     ['male', 'from US', 'uses Safari'],
8     ['female', 'from Europe', 'uses Firefox']]
9 enc.fit(X)
10
11 enc.transform([
12     ['female', 'from US', 'uses Safari']
13 ])
```

```
array([[0., 1., 1.]])
```

Continuous Values? Discretization

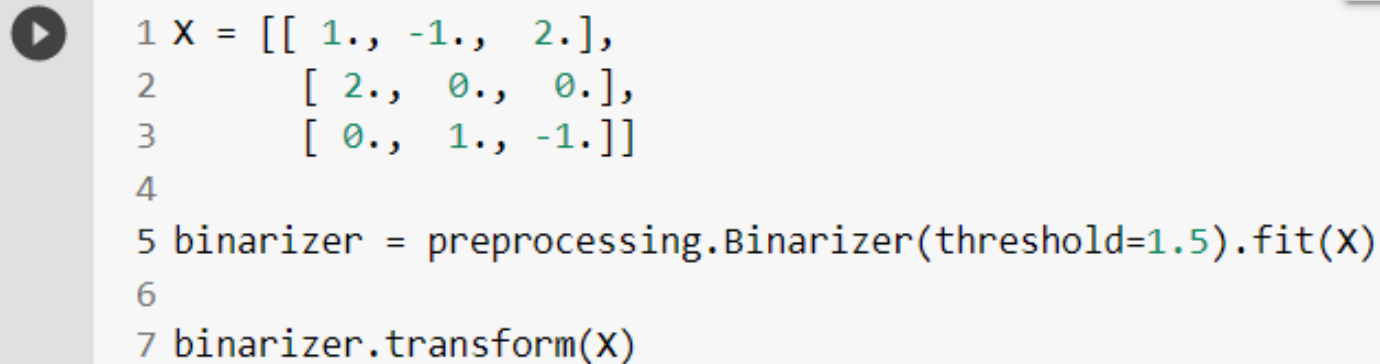
- Bin-based Discretization

```
1 from sklearn import preprocessing
2 import numpy as np
3 from numpy.linalg import norm
4
5 X = np.array([[ -3.,  5., 15 ],
6               [  0.,  6., 14 ],
7               [  6.,  3., 11 ]])
8
9 est = preprocessing.KBinsDiscretizer(
10     n_bins=[3, 2, 2], encode='ordinal'
11 ).fit(X)
12 |
13 est.transform(X)
```

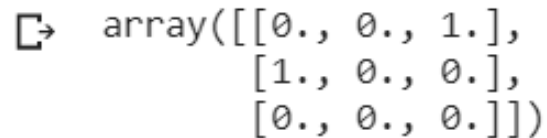
```
➤ array([[0., 1., 1.],
         [1., 1., 1.],
         [2., 0., 0.]])
```


Continuous Values? Discretization

- Thresholding numerical features => binominal



```
1 X = [[ 1., -1.,  2.],  
2      [ 2.,  0.,  0.],  
3      [ 0.,  1., -1.]]  
4  
5 binarizer = preprocessing.Binarizer(threshold=1.5).fit(X)  
6  
7 binarizer.transform(X)
```



```
array([[0., 0., 1.],  
       [1., 0., 0.],  
       [0., 0., 0.]])
```