

Video 1

Overview of Information Integration

It refers to the problem of using many different info. sources to accomplish a task.

- 10.00 → Explain the data integration problem
→ Define integrated views and schema mapping
- 11.00 → Describe the impact of increasing the no. of data source
→ Appreciate the need to use data compression.
- 12.00 → Describe Record Linking, Data Exchange and Data Fusion
task.
Z one vs multiple
data source
warehouse

Record linking

3.00 The "Big data" problem

- many sources -

- ① Hundreds of tables
② schema mapping problem
is a combinational challenge.

- 6.00 → pay-as-you-go model

- 6.00 → pay-as-you-go model
- ① Only integrate sources that are needed when needed
- 7.00 → probabilistic schema mapping

Design mediated Schema

Attribute Grouping

Data Integration Scenarios

→ 4 data sources each with one relation.

10.00 **Data Exchange** :- Given a source database with a finite no. of relations, a set of schema mappings, and a set of constraints that the target schema must satisfy, the data exchange problem is to find a finite target database such that both the schema mappings and the target constraints are satisfied.

2.00 **Using Codebooks** :- Logical observations identifies raw and codes is a database.

3.00 **Using Compressed data**

4.00 **Compression** :- Encoded representation of data so that it uses less space.

5.00 → **Dictionary encoding**

6.00 **Ontological data** :-

Ontology queries are graph queries.

7.00 **Ontology** :- A set of terms of a domain
- Relationship bⁿ the terms.

Snomed :- A medical ontology used for clinical data.

NOTES

Data Fusion :-

- Data sources

- Data Items

 - A product

 - A part of product

 - A feature of a product

 - ...

} value

- Using data from a subset of sources find the true value or a true value distribution of a data item.
- Assemble all such values for the real-world entity represented by the data items.

Too many sources :-

- Too many sources = too many values.
- Voting to select the "right" value.

- simple voting can be problematic -

Veracity Problem

- source Reliability
- Copy Detection
- Statistical techniques to estimate
 - Trustworthiness of sources.
 - Bias introduced by copies.
 - True distribution of values for data items

Source Selection :-

The problem :- 1) choose only useful sources.

2) Adding sources first improves integration accuracy then reduces it.

The solⁿ :- Order candidate source based on a measure of "goodness"

→ Add sources until the marginal benefit is less than the marginal cost

→ Current techniques scale well.

Quiz 3 :- Information Integration

1. What is the main problem with big data information integration?
→ Many sources.
2. Two possible solⁿs associated with "big data" info. integration as mentioned in lect?
→ • Probabilistic Schema mapping
• Pay-as-you-go Model.
3. Mediated schemas?
→ Schemas created from integrating 2 or more schemas.
4. In attribute grouping, how would one evaluate if two attributes should go together?
→ • probability of Two attributes Co-occurring
• similarity of Attributes.
5. What is data item?
→ Data that represents an aspect of a real-world entity.
6. What is data fusion?
→ Extracting the true value of a data item

7. What is a potential problem of having too many data sources as mentioned in text?

→ Too many data value.

8. What do we mean by the true value of a data item?

→ Extrapolated data from a data item that represent the worth of that item.

9. potential method to deal with too many data source

→ compare and weigh each source by their trustworthiness.