

## Trusted Zone

### 1) Combine All data

• IMDb → monthly, disjoint

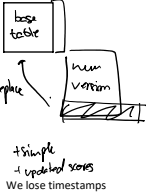
base table	movies actors directors
17/24	
08/24	

• Netflix users → eventful, disjoint

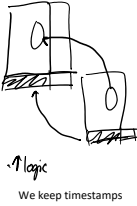
base table	
18/24	
09/10/24	
..	

• Netflix users → daily, incremental

option 1:



option 2:



option 3:

movie, score, timestamp
movie, score, timestamp
movie, score, timestamp
movie, score, timestamp

↓ the table can get very big...  
↑ we have all information

### 2) Check Data Quality: generic steps to avoid problems in data.

• Profiling: extract min, value, IR range... → quantitative  
with... → qualitative } Analyse the data with which we are working. Detect outliers!

extreme  
less extreme  
outliers that  
could be  
mistakes

• Removing duplicates: Some might have appeared in step ①, combining data sources.

• Consistent format: remove special characters...

• Data normalization:  $m_1, [a_1, a_2, a_3] \rightarrow \begin{matrix} m_1 & a_1 \\ m_2 & a_2 \\ m_3 & a_3 \end{matrix}$   
(categorical columns)

• Misspellings: typos!

• Difficult to detect...

① Easy way → use a spellchecker

② More advanced way → count times that each category appears + Levenshtein distance

③ Even ↑ tricky → Lemmatization/Stemming

• Outliers  
• Missing values  
• Scaling } We don't handle this categories since the applied methods can add bias before training the model. We will handle them in the exploitation zone.

EXPLOITATION ZONE → Expose data to the analysts

Guidelines: how to do the exploitation zone! Some tips

- Explain in detail why and what you do!
- Separate entities (static info) vs actions (dynamic info)
- Reduce redundancy and increase query efficiency
- Prepare for more data, don't "overfit" to the data we have now.
- Build Ground truth (KPIs) → We want to give to the analyst the more information, so we should generate KPIs for them to analyse

Actors = IMDb actors

opt 1

m1	a1
m1	a2
m2	a3

Option 2: the good one (do this)

One table for actors: a1 \_

a2 \_

a3 \_

Ona table for acts: a1 m1

a2 m1

a3 m2

Movies = IMDb movies + Netflix movies

Join operation, big table containing all the information

Possible issues:

- Not easy join bc of different granularities
- Overlapping information (repeated) in the attributes → Decide which one you take!

WE DON'T INCLUDE THE SCORES IN THE TABLE! we will set them in another table, and do the same as in the actors configuration

Scores table will be:

Movie score	IMDb score	Net timestamp	Avg. Score

The avg score could be, for example, an indication of a KPI

User views (3rd data source)

User	KPIs

User movies

User	movie	timestamp