




POLITECNICO DI MILANO




 **Managing Information Overload**

Elisa Quintarelli

Elisa Quintarelli – 22th April 2015

 **The Challenges for Modern Information Systems** 2

- **Organizations grow and generate more information:** they capture billions of bytes of information about their customers, suppliers and operations
- **The pervasiveness of digital technologies has changed the way individuals interact with the external world (sensor technologies) and with one another (social media),** generating a huge mass of content

 **Data has become as a torrent** that flows through all possible digital channels.

POLITECNICO DI MILANO

Information Overload 3

- The term was used by Alvin Toffler in his book *Future Shock*, already back in 1970
- It refers to the *difficulty of understanding and making decisions when too much information is available*
- This is the main challenge presented by “Big Data”

POLITECNICO DI MILANO

Data Management: What does it mean today? 4

What do users want from us?

- Massive data integration/exchange
- Heterogeneity/mobility
- Incompleteness/uncertainty
- Interaction with the real-world
- Knowledge representation and reasoning
- **Massive data analysis and data mining**

Summarized answers to user queries
Personalization and context awareness

↓

- **Making sense of all this data: extract useful knowledge**


POLITECNICO DI MILANO




 **Context-aware data personalization**

Antonio Miele, Elisa Quintarelli, Emanuele Rabosio, Letizia Tanca


Elisa Quintarelli – 22th April 2015



 **Introduction** 6


- Technological limitations of widespread small mobile devices (e.g.: pda, smart phones, ...) constrain the amount of data which can be loaded on it
- There is a quest for reduction and personalization of the information each user accesses to
- Two main issues:
 - Reduce information noise
 - Satisfy device memory limitations

POLITECNICO DI MILANO


 **Background and Motivations** 7

- The Context-ADDICT methodology tackles the presented problem, performing **context-based data tailoring**
- Current limitations:
 - The customization does not take into account **user preferences**
 - The context is the only driver
 - Data reduction to fit the memory of the device is not supported
 - No memory occupation model is considered
 - The approach is coarse-grain

POLITECNICO DI MILANO

 **The proposed approach** 8

- Preferences express interests on data as **numerical scores or explicit ordering** relations
- Data scoring is commonly used to rank information in several of today's data management applications and search engines



- Fine-grained personalization of data can be performed by means of preferences

POLITECNICO DI MILANO

↘
Contextual preferences – Examples
9

- A user may prefer spicy dishes **at dinner**, but not spicy ones **at lunch**
- A user may prefer comedy movies when **he/she is alone**, but thrillers when **with his/her friends**

POLITECNICO DI MILANO

↘
Running Example
10

- An application for an integrated service of meal order and delivery is considered; it involves several restaurants and meal delivery taxi companies

CUISINES(cuisine_id, description)

DISHES(dish_id, description, isVegetarian, isSpicy, isMildSpicy, wasFrozen, category_id)

RESERVATIONS(reservation_id, customer_id, restaurant_id, date, time)

RESTAURANTS(restaurant_id, name, address, zipcode, city, state, zone_id, rnumber, phone, fax, email, website, openinghourslunch, openinghoursdinner, closingday, capacity, parking, minimumorder, rating)

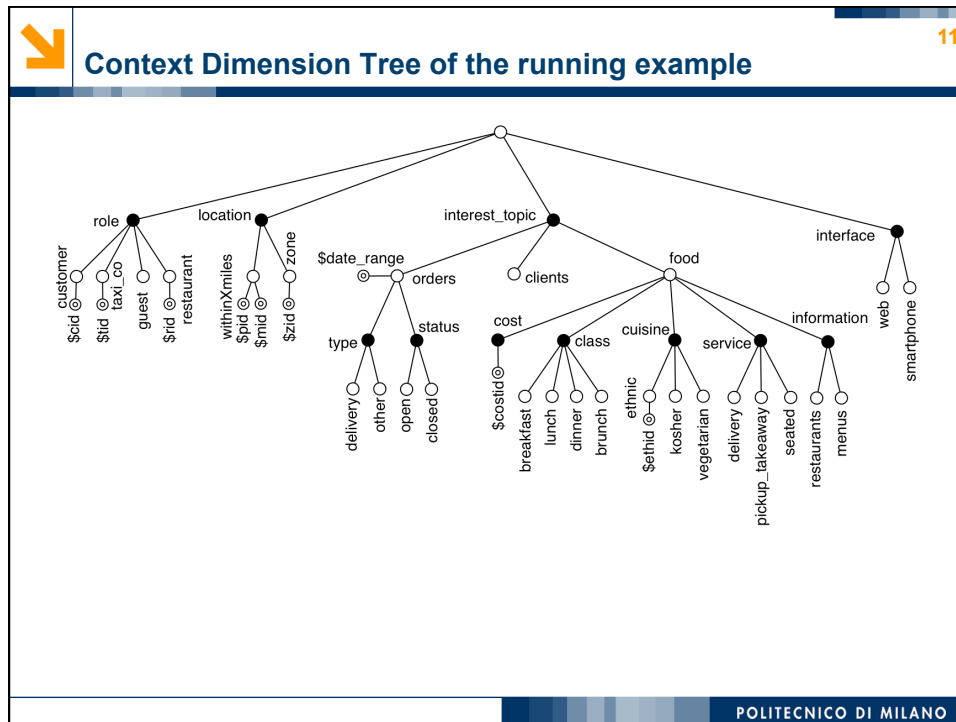
RESTAURANT_DISH(restaurant_id, dish_id)

RESTAURANT_CUISINE(restaurant_id, cuisine_id)

RESTAURANT_SERVICE(restaurant_id, service_id)

SERVICES(service_id, name, description)


POLITECNICO DI MILANO



Preference Model – Overview 12

- **Quantitative preferences**
 - Numerical scores expressed in the range [0;1]
- Two different types of preferences:
 - σ -preferences
 - Expressed on data tuples
 - Acting in horizontal on each table
 - π -preferences
 - Expressed on relation attributes
 - Acting in vertical on each table

POLITECNICO DI MILANO




Preference Model – Preferences on data tuples

13

- **σ -preference** $P_\sigma(R)$ on the relation $R(X)$ is defined as

$$\langle C, SQ_\sigma, S_\sigma \rangle$$
- C : the **context configuration**
- SQ_σ : the **selection query** composed by a selection on $r(X)$ possibly semi-joined on selections of related tables (only on foreign keys)
- S_σ : the **numerical score**

POLITECNICO DI MILANO



Preference Model – Preference on data tuples Example


14

```

CP1 = ⟨C1 = role = client("Smith") ∧ location = zone("CentralSt."),
      SQσ1 = σisVegetarian=1(dishes),
      Sσ1 = 0.1)
CP2 = ⟨C2 = role = client("Smith") ∧ location = zone("CentralSt."),
      SQσ2 = σisSpicy=1(dishes),
      Sσ2 = 1)
CP3 = ⟨C3 = role = client("Smith") ∧ location = zone("CentralSt."),
      SQσ3 = restaurant ⋈ restaurant_cuisine ⋈ σcuisine.description="Chinese" cuisine,
      Sσ3 = 1)
CP4 = ⟨C4 = role = client("Smith") ∧ location = zone("CentralSt."),
      SQσ4 = restaurant ⋈ restaurant_cuisine ⋈ σcuisine.description="Indian" cuisine,
      Sσ4 = 0.3)

```

POLITECNICO DI MILANO




Preference Model – Preferences on relation attributes

15

- **π -preference** $P_{\pi}(R)$ on the relation $R(X)$ is defined as

$$\langle C, A_{\pi}, S_{\pi} \rangle$$
- C : the **context configuration**
- A_{π} : the **list of attributes** on the relation $R(X)$
- S_{π} : the **numerical score**

POLITECNICO DI MILANO



Preference Model – Preferences on relation attributes

Example

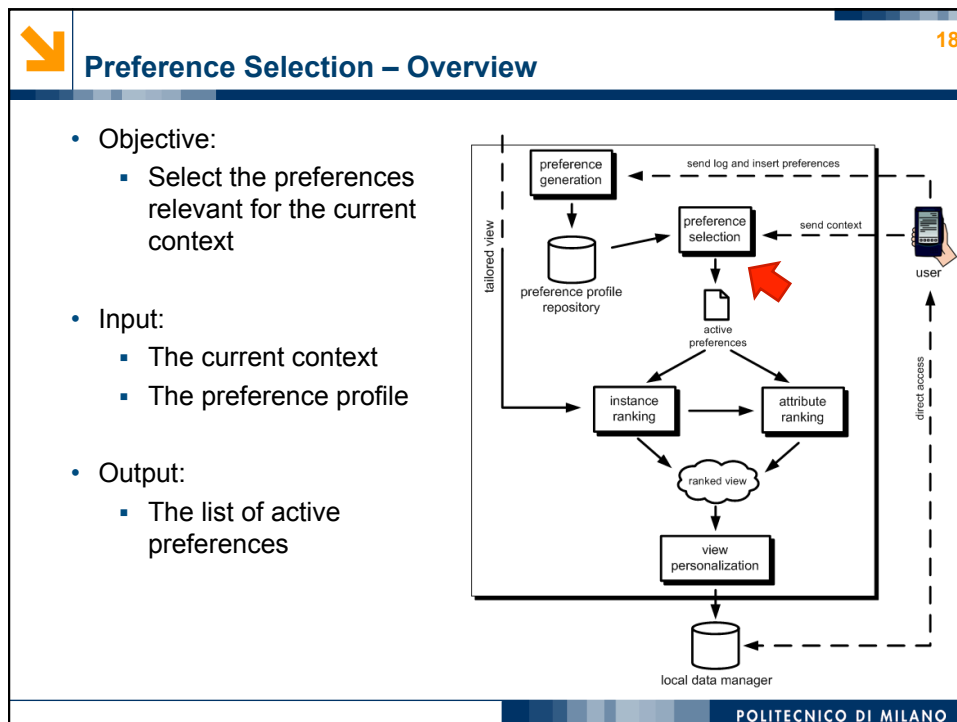
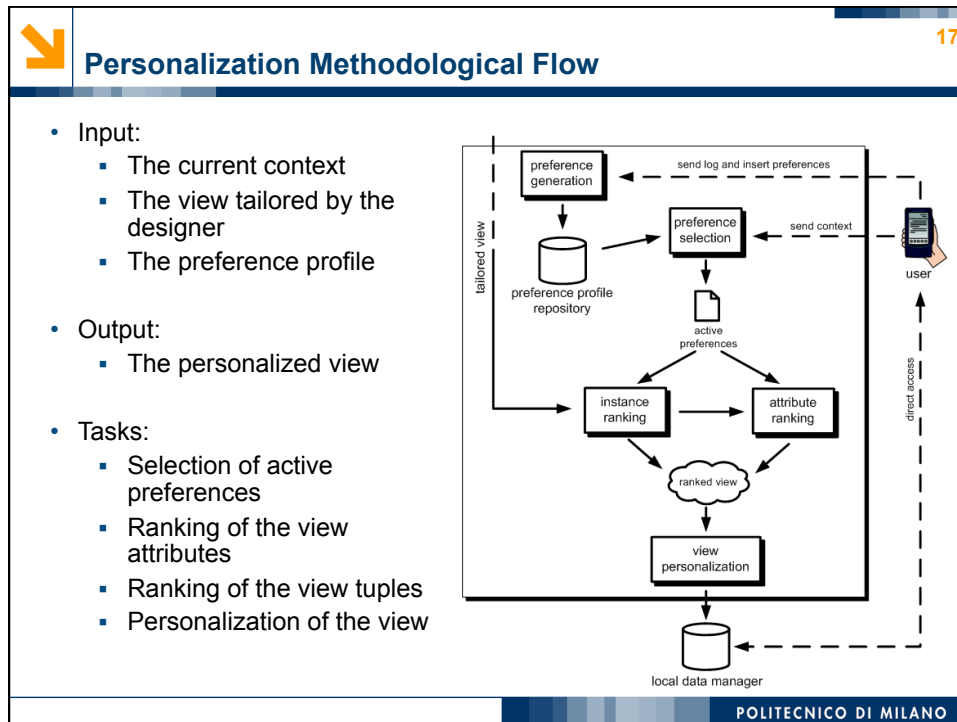
16

```

CP5 = ⟨C5   = role = client("Smith") ^ location = zone("CentralSt.")
        ^ interface = smartphone,
        Aπ5  = {name, phone, zipcode},
        Sπ5  = 1⟩
CP6 = ⟨C6   = role = client("Smith") ^ location = zone("CentralSt.")
        ^ interface = smartphone,
        Aπ6  = {address, city, rnumber, fax, email, website},
        Sπ6  = 0.2⟩

```

POLITECNICO DI MILANO



Preference Selection – Relevance relation 19

- Relevance:**
 - A partial order relationship (abstractness) is defined among context configurations (CDT hierarchical structure)
 - A preference is **relevant for the current context** if its context configuration is equal or more abstract than the current context
- Distance index:**

$$dist(C_1, C_2) = |asc_dim(C_1) - asc_dim(C_2)|$$
 - $asc_dim(C_i)$ = number of dimension nodes that are ancestor of the instantiated context value nodes of C_i
 - Represents the number of dimension nodes present on the CDT between the two context configurations


POLITECNICO DI MILANO

Preference selection – Relevance relation Example 20

$C_1 = \langle \text{role} = \text{client}(\text{"Smith"}) \wedge \text{location} = \text{zone}(\text{"CentralSt."}) \rangle$
 $C_2 = \langle \text{role} = \text{client}(\text{"Smith"}) \wedge \text{location} = \text{zone}(\text{"CentralSt."}) \wedge \text{cuisine} = \text{vegetarian} \wedge \text{class} = \text{lunch} \rangle$

$dist(C_1, C_2) = |asc_dim(C_1) - asc_dim(C_2)| = |2 - 5| = 3$

POLITECNICO DI MILANO



Preference selection - Steps


21

- Select from the context profile preferences relevant for the current context
- Assign to each relevant preference cp a **relevance index**:

$$relevance(cp) = \frac{dist(C_{curr_context}, C_{root_context}) - dist(C_{cp}, C_{curr_context})}{dist(C_{curr_context}, C_{root_context})}$$

- It is a relevance percentage w.r.t. the current context
- **Relevance = 1** for preferences expressed on the current context
- **Relevance = 0** for preferences expressed on the root context

POLITECNICO DI MILANO



Preference selection – Example

22

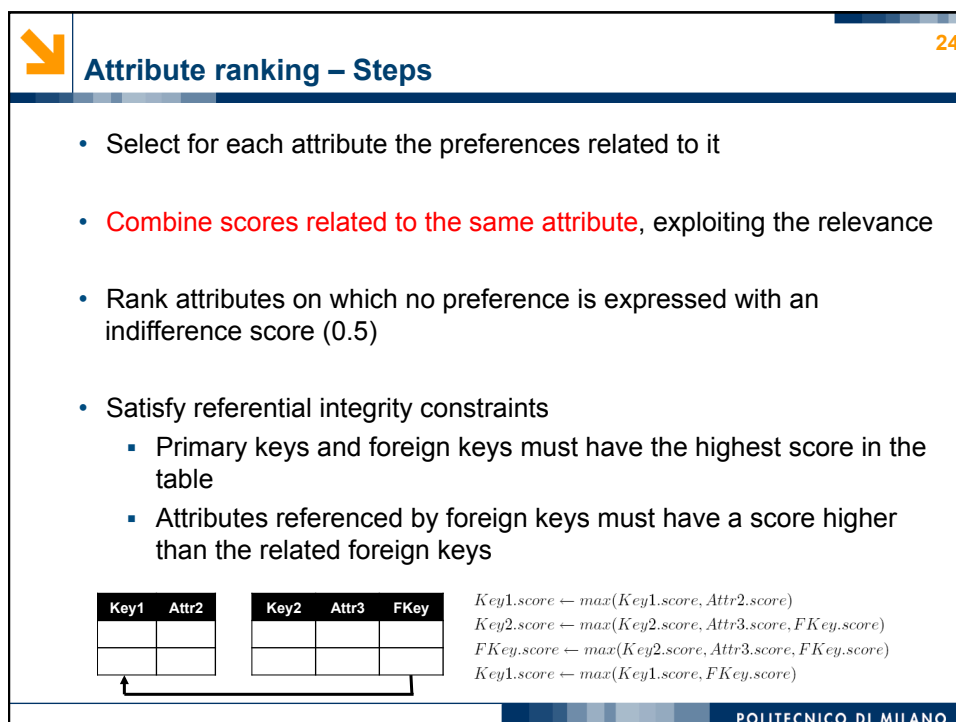
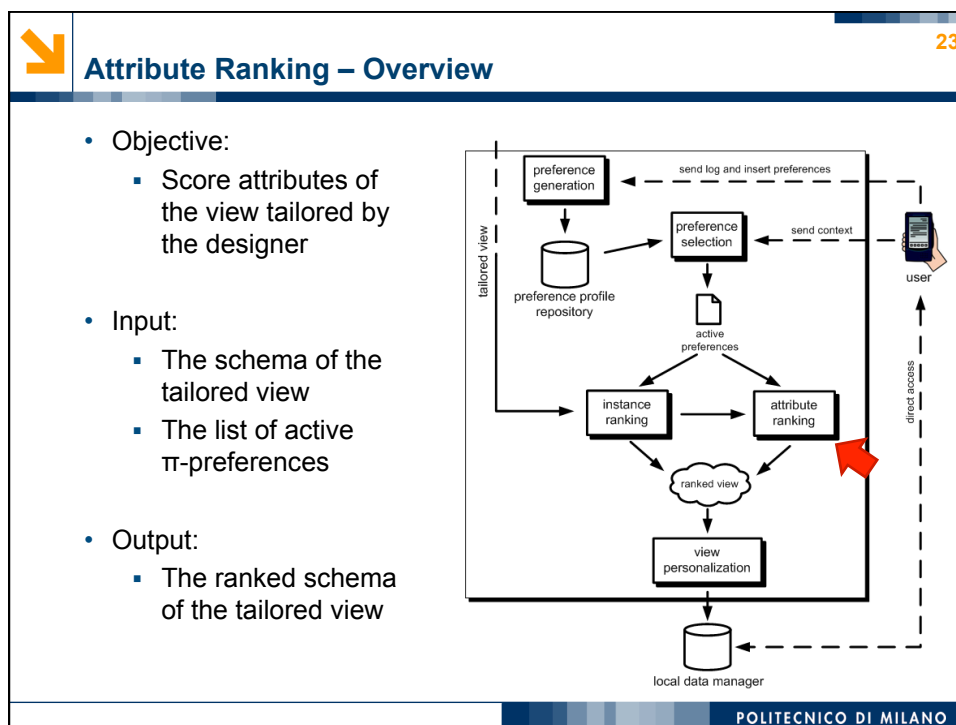
Current Context: $dist(C_{curr}, C_{root}) = 4$

$C_{curr} = \langle role = client("Smith") \wedge location = zone("CentralSt.") \wedge information = restaurants \rangle$

Preference Profile:

$dist(C_{curr}, C_1) = 0$ relevance = 1	✓	$CP_1 = \langle C_1 = role = client("Smith") \wedge location = zone("CentralSt.") \wedge information = restaurant, SQ_{\sigma 1} = restaurant \times restaurant_cuisine \times \sigma_{cuisine.description="Chinese"} cuisine, S_{\sigma 1} = 0.8 \rangle$
$dist(C_{curr}, C_2) = 1$ relevance = 0.75	✓	$CP_2 = \langle C_2 = role = client("Smith") \wedge information = restaurant, SQ_{\sigma 2} = restaurant \times restaurant_cuisine \times \sigma_{cuisine.description="Chinese"} cuisine, S_{\sigma 2} = 0.5 \rangle$
$dist(C_{curr}, C_3) = 3$ relevance = 0.25	✓	$CP_3 = \langle C_3 = role = client("Smith"), SQ_{\sigma 3} = \sigma_{isVegetarian=1}(dishes), S_{\sigma 3} = 0.3 \rangle$
$dist(C_{curr}, C_4) = 3$ relevance = 0.25	✓	$CP_4 = \langle C_4 = role = client("Smith"), SQ_{\sigma 4} = \sigma_{isMildSpicy=1}(dishes), S_{\sigma 4} = 0.8 \rangle$
	✗	$CP_5 = \langle C_5 = role = client("Smith") \wedge location = zone("CentralSt.") \wedge interface = smartphone, A_{\pi 5} = \{name, zipcode, phone\}, S_{\pi 5} = 0.8 \rangle$
	✗	$CP_6 = \langle C_6 = role = client("Smith") \wedge location = zone("CentralSt.") \wedge interface = smartphone, A_{\pi 6} = \{address.city\}, S_{\pi 6} = 0.2 \rangle$

POLITECNICO DI MILANO



Attribute ranking – Example
25

RESTAURANTS(restaurant_id, name, address, zipcode, city, phone, fax, email, website, openinghourslunch, openinghoursdinner, closingday, capacity, parking)
 RESTAURANT_CUISINE(restaurant_id, cuisine_id)
 CUISINES(cuisine_id, description)

↓

$$P_{\pi_1} = \langle \{name, cuisine.description, phone, closingday\}, 1 \rangle, R = 1$$

$$P_{\pi_2} = \langle \{address, city, state, phone\}, 0.1 \rangle, R = 0.2$$

$$P_{\pi_3} = \langle \{fax, email, website\}, 0.1 \rangle, R = 0.2$$

RESTAURANTS(restaurant_id:1, name:1, address:0.1, zipcode:0.5, city:0.1, phone:1, fax:0.1, email:0.1, website:0.1, openinghourslunch:0.5, openinghoursdinner:0.5, closingday:1, capacity:0.5, parking:0.5)
 RESTAURANT_CUISINE(restaurant_id:0.5, cuisine_id:0.5)
 CUISINES(cuisine_id:1, description:1)

POLITECNICO DI MILANO

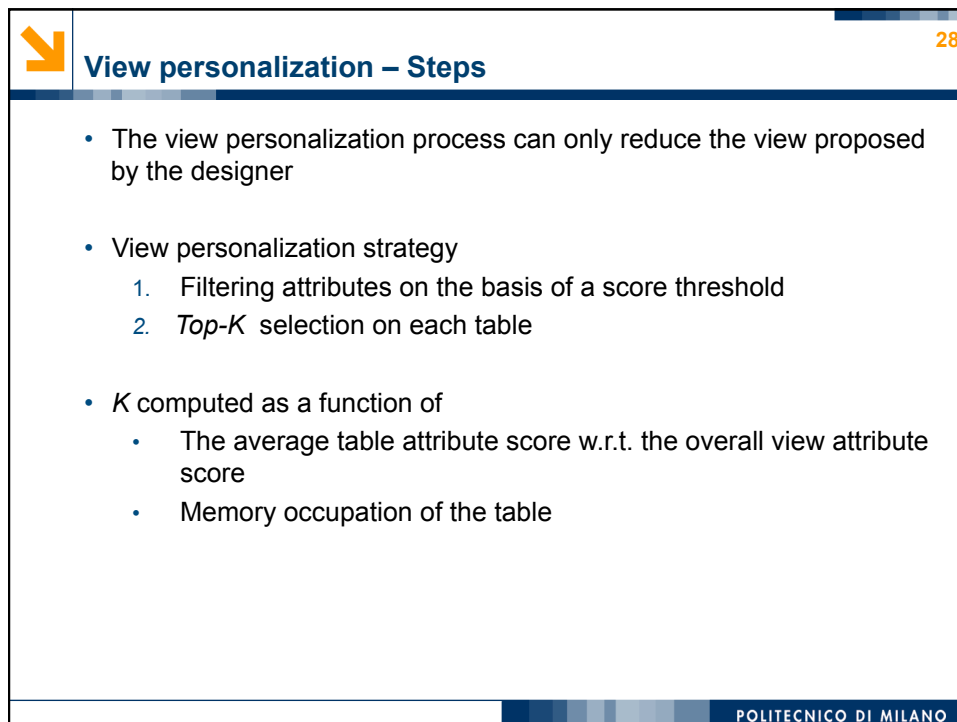
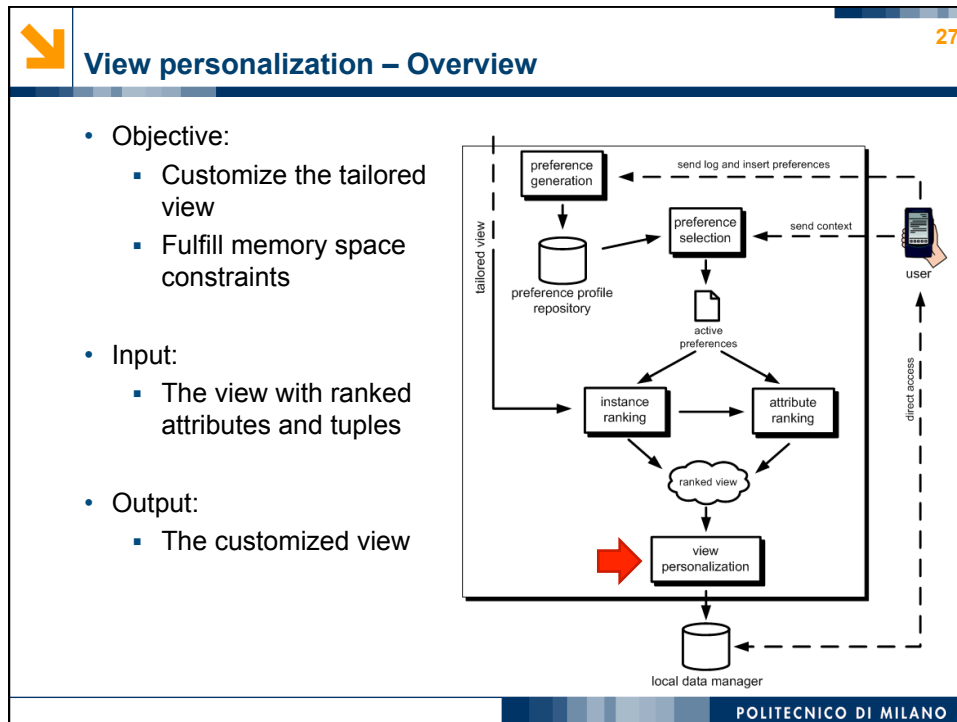
Tuple ranking – Overview
26


- Objective:
 - Score tuples of the view tailored by the designer
- Input:
 - The tailored view
 - The list of active σ -preferences
- Output:
 - The tailored view with ranked data tuples

```


            graph TD
            subgraph Tailored_View [tailored view]
            PG[preference generation] --> PR[preference profile repository]
            PR --> PS[preference selection]
            PS --> AP[active preferences]
            AP --> IR[instance ranking]
            AP --> AR[attribute ranking]
            IR --> RV[ranked view]
            AR --> RV
            RV --> VP[view personalization]
            VP --> LDM[(local data manager)]
            end
            LDM --> TV[tailored view]
            TV --> PG
            U[user] --> PS
            PS --> U
            U --> LDM
            LDM -.-> U
            
```

POLITECNICO DI MILANO




 **View personalization - Example** 29

RESTAURANTS(restaurant_id:1, name:1, address:0.1, zipcode:0.5, city:0.1, phone:1, fax:0.1, email:0.1, website:0.1, openinghourslunch:0.5, openinghoursdinner:0.5, closingday:1, capacity:0.5, parking:0.5)
 RESTAURANT_CUISINE(restaurant_id:0.5, cuisine_id:0.5)
 CUISINES(cuisine_id:1, description:1)

 threshold = 0.5

RESTAURANTS(restaurant_id:1, name:1, zipcode:0.5, phone:1, openinghourslunch:0.5, openinghoursdinner:0.5, closingday:1, capacity:0.5, parking:0.5)
 RESTAURANT_CUISINE(restaurant_id:0.5, cuisine_id:0.5)
 CUISINES(cuisine_id:1, description:1)

POLITECNICO DI MILANO

 **View personalization – Example (2)** 30

RESTAURANTS(restaurant_id:1, name:1, zipcode:0.5, phone:1, openinghourslunch:0.5, openinghoursdinner:0.5, closingday:1, capacity:0.5, parking:0.5)
 RESTAURANT_CUISINE(restaurant_id:0.5, cuisine_id:0.5)
 CUISINES(cuisine_id:1, description:1)


 compute average schema score

Table	Average Score
RESTAURANTS	0.72
RESTAURANT_CUISINE	0.5
CUISINES	1
RESTAURANT_SERVICE	0.5
SERVICE	0.6
RESERVATION	0.72

POLITECNICO DI MILANO

View personalization – Example (3) 31

Table	Average Score
RESTAURANTS	0.72
RESTAURANT_CUISINE	0.5
CUISINES	1
RESTAURANT_SERVICE	0.5
SERVICE	0.6
RESERVATION	0.72

↓ order tables and partition available space (2Mb)

Table	Average Score	Memory (Mb)
CUISINES	1	0.50
RESTAURANTS	0.72	0.35
RESERVATION	0.72	0.35
SERVICE	0.6	0.30
RESTAURANT_CUISINE	0.5	0.25
RESTAURANT_SERVICE	0.5	0.25

POLITECNICO DI MILANO

Preference generation – Overview 32

- Objective:
 - Generate preference profiles analyzing log data, extracting knowledge in terms of association rules
- Input:
 - User activity log
 - Old preferences
- Output:
 - σ - and π -preferences
- The sub-tasks are performed independently for σ - and π -preferences

POLITECNICO DI MILANO

Mining σ -preferences Log synchronization

33

- Input: log stored on the client device**
 - Textual recording of SQL queries
- Output: log on the server**
 - A relational table for each table in the database
 - The log associated with a table R(X) contains:
 - An attribute for each black node of the CDT
 - An attribute for each attribute belonging to X
 - An attribute for each attribute of each table reachable from R(X) through foreign key constraints
 - A row for each tuple returned in user query answers

POLITECNICO DI MILANO

Mining σ -preferences Server log – Example

34

```

SELECT DISTINCT dishes.description
FROM dishes, restaurants, restaurant_dish
WHERE restaurants.restaurant_id= restaurant_dish.restaurant_id
AND dishes.dish_id= restaurant_dish.dish_id
AND restaurants.closingday='Monday'
    
```

RESTAURANTS		
restaurant_id	name	closingday
r1	Pizzeria Rita	Monday
r2	Cing Restaurant	Tuesday
r3	Cantina Mariachi	Monday

DISHES	
dish_id	description
p1	Pizza Margherita
p2	Pizza Napoli
p3	Pizza Capricciosa

RESTAURANT_DISH	
restaurant_id	dish_id
r1	p1
r1	p3
r2	p2
r2	p3
r3	p1

Log of the table Restaurant_dish

id	Context dimensions				Restaurant_dish attributes		Restaurants attributes		Dishes attributes		
	role	int-topic	cuisine	...	rd.r_id	rd.d_id	r.r_id	r.name	r.closingday	d.d_id	d.descr
1	client(cli1)	food	veg	...	r1	p1	r1	Pizzeria Rita	Monday	p1	Margh
2	client(cli1)	food	veg	...	r1	p3	r1	Pizzeria Rita	Monday	p3	Capr...
3	client(cli1)	food	veg	...	r3	p1	r3	Cantina Mariachi	Monday	p1	Margh

POLITECNICO DI MILANO

Mining σ -preferences

Association rule mining

35

- **Input:** *server log*
- **Output:** *association rules*

- Knowledge is extracted from the log by means of association rules
- An association rule is an implication in the form

$$A \rightarrow B$$
- Quality indexes for association rules:

$$\text{Support} = \text{num. of data with } A \cup B$$

$$\text{Confidence} = \frac{\text{num. of data with } A \cup B}{\text{num. of data with } A}$$

POLITECNICO DI MILANO

Mining σ -preferences

Association rule mining (2)

36

- We are interested in **σ -rules**, correlating contexts and data
- A σ -rule on a relation $R(X)$ is a triple:

$$\langle C \rightarrow \text{cond}, \text{sup}, \text{conf} \rangle$$
- C : a **context**
- Cond : a **conjunction of conditions** in the form $A = \text{value}$, where A is an attribute belonging to $R(X)$ or to a relation reachable from $R(X)$ through foreign keys
- sup : the **support** of the association rule $C \rightarrow \text{cond}$
- conf : is the **confidence** of the association rule $C \rightarrow \text{cond}$
- Example:

$$\langle \text{location} = \text{zone}(\text{"Central St."}) \wedge \text{interest} - \text{topic} = \text{food} \rightarrow$$

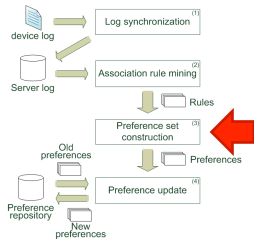
$$\text{isVegetarian} = \text{true} \wedge \text{isSpicy} = \text{false}, 0.2, 0.7 \rangle$$

POLITECNICO DI MILANO

Mining σ -preferences Preference-set construction

37

- **Input:** σ -rules
- **Output:** σ -preferences
- Three sub-tasks:
 - **Removal of redundant rules**, whose satisfaction is implied by other ones
 - **Computation of the score** of the σ -preferences



The flowchart illustrates the process of preference-set construction. It starts with 'device log' and 'Server log' as inputs. 'device log' goes through 'Log synchronization' (1) to produce 'Rules'. 'Server log' goes through 'Association rule mining' (2) to produce 'Rules'. 'Rules' then go through 'Preference set construction' (3) to produce 'Preferences'. 'Preferences' are then used for 'Preference update' (4), which involves 'Old preferences' and 'New preferences' from a 'Preference repository'.

POLITECNICO DI MILANO

Mining σ -preferences Score computation

38

- A σ -rule $r = \langle C \rightarrow cond, sup, conf \rangle$ may determine a σ -preference $p = \langle C, \langle SQ_\sigma, score \rangle \rangle$ such that $p.C = r.c$ and $p.SQ_\sigma = C.cond$
- The mining procedure can discover only preferences indicating user interest in the interval $(0.5, 1]$, and **only if the confidence of the rule is greater than the frequency of the data satisfying it**
- The score is computed as follows:

$$\Delta = r.conf - f$$

$$p.score = \min((1 + \gamma \cdot \Delta) \cdot 0.5, 1), \text{ if } \Delta > 0$$
- f is the frequency of the data satisfying $r.cond$ in the data set accessed by the user

POLITECNICO DI MILANO

Mining π -preferences Log synchronization

39

- **Input:** *log stored on the client device*
 - Textual recording of SQL queries
- **Output:** *log on the server*
 - An attribute for each black node of the CDT
 - An attribute for each attribute of each table
 - A row for each query
 - Non-contextual columns contain '1' if the associated attribute has been accessed in the query
 - A query accesses an attribute if it is contained in the select or where clause

```

SELECT restaurants.address,
restaurants.phone
FROM restaurants
WHERE restaurants.closingday <> 'lunedì'

```

```

SELECT dishes.description
FROM dishes
WHERE dishes.category_id='high'

```

Context dimensions			Attributes of the database relations										
role	interest-topic	cuisine	r.rest_id	r.addr	r.city	r.phone	r.closingday	d.dish_id	d.descr	d.isSpicy	d.isMildSpicy	d.cat_id	
client(Smith)	food	veget		1		1	1						
client(Smith)	food	veget								1		1	

POLITECNICO DI MILANO

Mining π -preferences Association rule mining

40

- **Input:** *server log*
- **Output:** *association rules*

- We are interested in π -rules, correlating contexts and attributes
- A π -rule on a relation $R(X)$ is a triple:

$$\langle C \rightarrow Attr, sup, conf \rangle$$
- C : a **context**
- $Attr$: an **attribute**
- sup : is the **support** of the association rule $C \rightarrow Attr$
- $conf$: is the **confidence** of the association rule $C \rightarrow Attr$

- Example:

$$\langle \text{interest - topic = food} \rightarrow \text{isVegetarian}, 0.1, 0.7 \rangle$$

POLITECNICO DI MILANO

Mining π -preferences Preference-set construction

41

- **Input:** π -rules
- **Output:** π -preferences

• Three sub-tasks:

- **Modification of the confidences of some rules** due to missing attributes
- **Computation of the score** of the π -preferences
- **Introduction of some not-mined preferences** to deal with preference-propagation problems

POLITECNICO DI MILANO

ADaPT: Automatic Data Personalization Based on Contextual Preferences [ICDE 2014]

42

What do we want to do this afternoon?
Why don't we watch a movie?
Let's look on the ADaPT app on my mobile phone

Movie list has changed!

POLITECNICO DI MILANO