

1. Le problème

Comment accéder aux données de santé?

• Ces données sont-elles suffisantes pour faire la prévention, le soin, et la recherche en santé ?

Ces données sont-elles suffisantes pour faire la prévention, le soin, et la recherche en santé ?

1. Le problème (bis)

Comment accéder aux données *requises pour la prévention, le soin et la recherche* en santé?

- La signification d'une donnée est-elle toujours univoque?
- Est-ce que la donnée est suffisante en elle-même?

épartement d'informatique, Faculté des sciences, Université de Sherbrooke, Québec

1. Le problème (ter)

Comment accéder aux données requises pour la prévention, le soin et la recherche en santé *et les interpréter (correctement)*?

- Ce problème est-il susceptible d'avoir une solution consensuelle, indépendamment des cultures, des sociétés, des états, des ordres professionnels et des intérêts privés?
- Ce problème est-il unique ou ne faudrait-il pas considérer une famille de problèmes, é une famille de solution?

-07-11 Département d'informatique, Faculté des sciences, Université de Sherbecoke, Québec



Il y aura plusieurs solutions partielles aux problèmes d'accès aux données requises pour la prévention, le soin et la recherche en santé.

La couverture de ces solutions demeurera partielle et, pour cette raison, il serait souhaitable qu'elle soit interopérable.

L'interprétation de ces données demeurera plurielle et, pour cette raison, il serait souhaitable qu'elle se fasse sur la base de processus et de modèles documentés, ouverts et traçables.



Puisqu'il y aura plusieurs modèles, il importe de se doter d'une démarche pour les élaborer.

Puisqu'il y aura plusieurs modèles, il faut déterminer un méta-modèle facilitant

- le raisonnement
- l'interopérabilité
- la documentation
- l'ouverture
- la traçabilité

2. Une démarche

- oCaractériser le problème
- •Caractériser les solutions
- oConcevoir, décrire et vérifier les solutions
- oChoisir certaines solutions
- oExpérimenter et valider celles-ci
- •Tirer les conclusions

formatique, Paculté des sciences, Université de Sherbrooke, Québec



Explicitons maintenant notre démarche tout en l'appliquant à notre premier problème: trouver LE méta-modèle! Département d'informatique, Faculté des sciences, Université de Sherbrooke, Québec

3. Un cas

- 3.1 Caractérisation du problème
- 3.2 Caractérisation des solutions
- 3.3. Description d'une première solution
- 3.4. Vérification de la solution
- 3.5. Une première expérimentation
- 3.6. Les premières conclusions

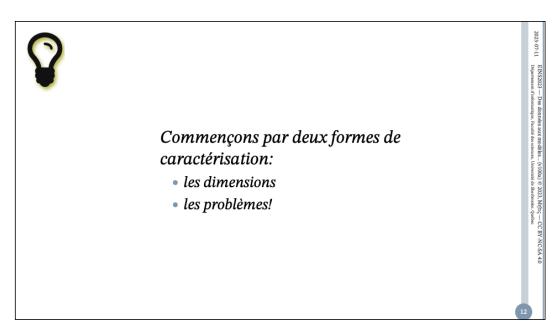
nées aux modèles... (v100a) © 2023, Μήτις — CC BY-NC-SA 4.0 zs, Faculté des sciences, Université de Sherbrooke, Québec

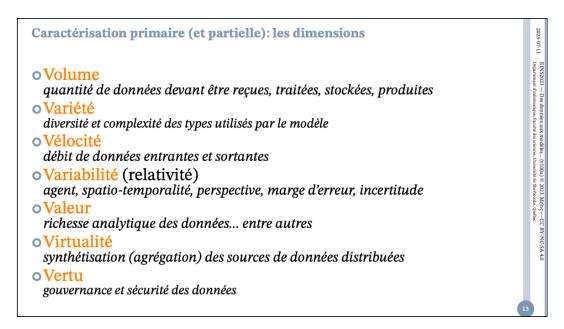
3. Un cas

3.1 Caractérisation du problème

- •Cerner un problème afin d'y trouver une solution, passe par la modélisation et l'analyse dudit problème et donc de sa caractérisation.
- •Comparer un problème avec d'autres est riche en enseignement, voire en solutions existantes ou adaptables.
- Choisir une solution, parmi d'autres, passe par la comparaison de celles-ci aux autres solutions.

natique, Paculté des sciences, Université de Sherbrooke, Québec





- * La mesure de ces dimensions n'est pas aussi triviale qu'on le voudrait.
- * Les procédés et les techniques ne font pas (encore) consensus.



complexe, sensible, débattue

Caractérisation secondaire (et partielle): les besoins

obesoin

chose considérée nécessaire à la définition d'un procédé ou au déroulement d'un processus (d'une activité, d'une tâche, d'une action).

•Ne pas confondre avec

- désir
- attente
- exigence

procédé::
...
processus::
...
activité::
ensemble de...

tâche::
action portant sur une chose
chose::
objet matériel ou information



Nous avons déjà identifié les besoins suivants

- soutenir au moins une démarche de modélisation
- soutenir le raisonnement
- soutenir l'interopérabilité
- soutenir la documentation
- soutenir l'ouverture
- soutenir la traçabilité, voire l'explicabilité

st Voir également article Khnaisser et coll. pour un approfondissement.

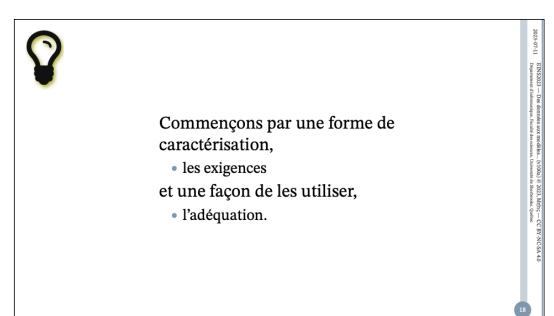
3. Un cas

3.2 Caractérisation des solutions

•En présence d'une solution, il faut s'assurer qu'elle en est bien une.

•En présence de plusieurs solutions, il faut pouvoir les comparer.

tement d'informatique, Paculté des sciences, Université de Sherbrooke, Québec



Les exigences

Les exigences ont pour but de déterminer

- les conditions nécessaires et suffisantes pour qu'une solution soit acceptable;
- les caractéristiques mesurables et pertinentes permettant de comparer deux solutions.

Plusieurs procédés ont été proposés afin de permettre la détermination des exigences. Ces procédés se distinguent par l'organisation des activités suivantes et les méthodes préconisées pour chacune.

- * exploration
- * analyse
- * spécification
- * vérification
- * validation



Exigences pratiques

- Capacité de soutenir un processus de diagnostic puis de choix de traitement (fondée sur des données probantes en regard de pratiques médicales reconnues);
- o corolaire 1: capacité de décrire l'état du patient,
- o corolaire 2: capacité de formuler le diagnostic,
- o corolaire 3: capacité de décrire un traitement,

o ...

• ..

-07-11 Département d'informatique, Faculté des sciences, Université de Sher

es sciences, Université de Sherbrooke, Québec



Exigences induites:

- Capacité de transposer les théories et les modèles scientifiques dans un cadre unifié
 - o conséquence: le méta-modèle doit fournir un formalisme permettant de écrire les théories et les modèles scientifiques.
- Capacité de formuler des hypothèses et de les vérifier
 - conséquence 1: le méta-modèle doit permettre le raisonnement et l'utilisation de données probantes en regard de théories et de modèles scientifiques décrits à l'aide du formalisme.
 - o conséquence 2: le méta-modèle doit permettre la vérification automatisée de prédicats en regard de données.
- Capacité de comparer deux modèles
- Capacité d'évaluer l'adéquation d'un modèle

Remarquons que les deux capacités induites sont applicables tant au méta-modèle qu'aux modèles qu'il permet de décrire (modéliser).

Le processus d'exploration n'est pas terminé :

- * une clarification est certainement possible ;
- * d'autres exigences découlent des besoins exposés ;

Un processus de spécification devrait s'ensuivre afin

- * de s'assurer que chaque exigence soit
 - claire (c'est-à-dire lisible et compréhensible),
 - exacte (c'est-à-dire précise et sans erreurs),
- complète (c'est-à-dire comprenant tous les éléments requis et tous les éléments nécessaires) et
 - concise (c'est-à-dire sans éléments superflus);
- * d'associer à chaque exigence un critère objectif permettant de vérifier qu'elle est satisfaite.

Finalement, la vérification et la validation devraient s'ensuivre afin de

minimiser les risques de devoir recommencer la modélisation qui en découlera.

Département d'informatique, Faculté des sciences, Université de Sherbrooke, Québ

erbrooke, Québec

L'adéquation

L'adéquation détermine à la fois la façon d'utiliser les exigences, de les concilier, mais aussi de les compléter au moment d'une évaluation globale d'une solution.

Plusieurs procédés ont été proposés afin de permettre la détermination des exigences. Ces procédés se distinguent par l'organisation des activités suivantes et les méthodes préconisées pour chacune.

- * exploration
- * analyse
- * spécification
- * vérification
- * validation

L'adéquation (suite)

Critères absolus

- Validité (conformité au modèle)
- Efficacité (conformité aux exigences)
- Cohérence (non-contradiction interne)

Critères relatifs

- Complétude (couverture «suffisante» du problème)
- Efficience (consommation «acceptable» de ressources)
- Évolutivité (adpatation «aisée» aux changements)

Méta-critères

- Réfutabilité (permet d'exprimer et d'évaluer les «falsificateurs potentiels»)
- Acceptabilité (permet d'exprimer et d'évaluer les «critères éthiques»)

L'adéquation (fin)

Les critères relatifs sont souvent décomposés en sous-critères auxquels sont associés un protocole de mesure et un poids de façon à permettre une évaluation globale pondérée. D'autres méthodes d'évaluation globale sont possibles.

La réfutabilité est considérée en regard d'une épistémologie donnée; en conséquence, une solution doit préciser son cadre épistémologique.

L'acceptabilité est considérée en regard d'une éthique donnée; en conséquence, une solution doit préciser cadre éthique.

Le cas du méta-modèle.

Au fin de l'exercice, posons

* pour l'épistémologie, celle de Karl Popper avec la réfutabilité discrète

(de préférence binaire);

* pour l'éthique, celle de Karl Popper avec le principe de la nécessaire protection

de la liberté par l'État… protection qui passe par sa limitation nécessaire

et suffisante!

Nous y reviendrons après avoir décrit le méta-modèle, au moment de l'évaluer.

Département d'informatique, Faculté des sciences, Université de Sherbrooke, Québec

3. Un cas

3.3 Description d'une (première) solution

Cheminement

- Raisonnement
 - o la capacité de soutenir le raisonnement, voire de l'automatiser, est au coeur des capacités recherchées et des exigences à satisfaire.
- Logique
 - o la logique décrit les lois du raisonnement gouvernant les prédicats.
- Relation
 - o les relations sont une représentation privilégiée des prédicats qui en facilitent le calcul grâce à une algèbre appropriée [adéquate].
- Typage
 - o le typage, fondé sur la théorie des ensembles, réduit les risques d'ambigüité, de paradoxes et d'indécidabilité dans la formulation des prédicats et le calcul des prédicats.

Ce cheminement et la théorie relationnelle qui en découle sont le fruit de

travaux que Edgard F. Codd entrepris au cours des années 1960 et mena jusqu'au début

des années 1990. La publication séminale date de 1969.

ovaleur, représentation, type
otype de base et sous-type
otype scalaire et type non scalaire
otype prédéfini et constructeur de type

Modèle relationnel Concepts • valeur • variable • fonction • état • procédure • automatisme Structure • attribut • tuple • relation • base

- oUn attribut est un couple formé d'un identifiant a et d'un type D, noté a:D.
- •Par abus de langage, lorsque le contexte le permet, il est usuel de désigner l'attribut par son seul identifiant; ainsi écrit-on l'attribut a.

Un tuple correspond à une observation, un «fait».

Un attribut correspond à une dimension de l'observation, à un «aspect» d'un fait.

La relation est un ensemble d'observations de même type, de même «nature».

Deux tuples sont de même type si leurs entêtes sont les mêmes.

Transposé dans le domaine de la logique

- un tuple est une proposition (un énoncé vrai sur le monde);
- une relation est un prédicat (...).

Un prédicat peut être défini

- par énumération (l'ensemble de tous les énoncés vrais et eux seuls);
- par compréhension (la caractérisation nécessaire et suffisante des relations entre les variables).

Une base de données est un ensemble de variable de relation définies par leurs valeurs (donc par énumération) et leurs contraintes (donc par compréhension).

En conclusion, une base de données est la représentation d'un système logique.

Dans la pratique, il est souvent difficile d'établir un ensemble de contraintes nécessaires et suffisantes.

Mais on tente de s'en approcher le plus possible.

Pour une théorie des types plus complète, voir (entre autres) IFT 232, IFT 339 et IGE 487.

```
• Soit a_i des identifiants distincts et D_j des types, un tuple t est défini comme suit:

• t \triangleq (\{a_1:D_1, a_2:D_2, ..., a_n:D_n\}; \{(a_1,v_1), (a_2,v_2), ..., (a_n,v_n)\})

• avec \forall i: 1 \leq i \leq deg(t) \Rightarrow val(t, a_i) \in def(t, a_i)

• def(t) = \{a_1:D_1, a_2:D_2, ..., a_n:D_n\} entête de t

• def(t, a_i) = D_i type de l'attribut a_i de t

• val(t) = \{(a_1,v_1), (a_2,v_2), ..., (a_n,v_n)\} valeur de t

• val(t, a_i) = v_i valeur de de l'attribut a_i de t

• deg(t) = n degré de t

• deg(t) = \{a_1, a_2, ..., a_n\} les identifiants d'attributs de t
```

Un tuple correspond à une observation, un «fait».

Un attribut correspond à une dimension de l'observation, à un «aspect» d'un fait.

La relation est un ensemble d'observations de même type, de même «nature».

Deux tuples sont de même type si leurs entêtes sont les mêmes.

Transposé dans le domaine de la logique

- un tuple est une proposition (un énoncé vrai sur le monde);
- une relation est un prédicat (...).

Un prédicat peut être défini

- par énumération (l'ensemble de tous les énoncés vrais et eux seuls);
- par compréhension (la caractérisation nécessaire et suffisante des relations entre les variables).

Une base de données est un ensemble de relations définies par leurs valeurs (donc par énumération) et leurs contraintes (donc par compréhension). En conclusion, une base de données est la représentation d'un système logique.

Dans la pratique, il est souvent difficile d'établir un ensemble de contraintes nécessaires et suffisantes.

Le cas deg(t) = 0, est important.

Il existe un seul tuple possible (pourquoi?):

• t0 = ({}, {})

Notation simplifiée fondée sur l'ordre d'énumération des attributs (les identifiants d'attributs et leurs types étant déterminés par ailleurs):

• $t = \langle v_1, v_2, ..., v_n \rangle$

Finalement, la notation t.ai désigne l'attribut ai:Di dans def(t).

```
• Soit a_i des identifiants distincts, D_j des types et t_k des tuples, une relation R est définie comme suit:

• R \triangleq (\{a_1:D_1, a_2:D_2, ..., a_n:D_n\}; \{t_1, t_2, ..., t_m\})
• avec \forall i: 1 \leq i \leq card(R) \Rightarrow def(R) = def(t_i)
• Où

• def(R) = \{a_1:D_1, a_2:D_2, ..., a_n:D_n\} entête de R
• def(R, a_i) = D_i type de a_i de R
• val(R) = \{t_1, t_2, ..., t_m\} valeur de R
• deg(R) = n degré de R
• card(R) = m cardinalité de R
• id(R) = \{a_1, a_2, ..., a_n\} identifiants d'attributs de R
```

Un prédicat peut être défini

- par énumération (l'ensemble de tous les énoncés vrais et eux seuls);
- par compréhension (la caractérisation nécessaire et suffisante des relations entre les variables).

Le cas deg(R) = 0 est important.

Il existe deux relations possibles (pourquoi?):

- $R0 = (\{\}, \{\})$
- $R1 = (\{\}, \{t0\})$

et elles sont très importantes, comme le zéro et le un pour les entiers!

Finalement, la notation R.ai désigne l'attribut ai:Di dans def(R).

oSoit v_i des identifiants distincts, D_j des types de relation et r_k des (valeurs de) relations, une base (de données) B est définie comme suit:

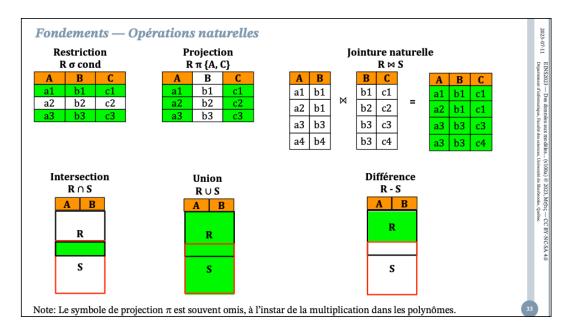
```
• B \triangleq (\{\mathbf{v}_1: \mathbf{D}_1, \mathbf{v}_2: \mathbf{D}_2, ..., \mathbf{v}_n: \mathbf{D}_n\}; \{\mathbf{r}_1, \mathbf{r}_2, ..., \mathbf{r}_m\})
• avec \forall i: 1 \leq i \leq card(B) \Rightarrow def(B, \mathbf{v}_i)=def(\mathbf{r}_i)
• Où
• def(B) = \{\mathbf{v}_1: \mathbf{D}_1, \mathbf{v}_2: \mathbf{D}_2, ..., \mathbf{v}_n: \mathbf{D}_n\} entête de B
• def(B, \mathbf{a}_i) = \mathbf{D}_i type de \mathbf{a}_i de B
• val(B) = \{\mathbf{r}_1, \mathbf{r}_2, ..., \mathbf{r}_m\} valeur de B
• deg(B) = n degré de B
• id(B) = \{\mathbf{v}_1, \mathbf{v}_2, ..., \mathbf{v}_n\} ensemble des identifiants de
```

variables de relation de B

nee's aux modenes... (v 100a) % 2025, Wijtis — CC BT-NC-SA 4-0 n, Faculté des sciences, Université de Sherbrooke, Québec

L'algèbre relationnelle

ounion, intersection, différence
orenommage
oprojection, restriction, jointure



Opérations relationnels propres: 3

Opérations ensemblistes: 3 (avec le produit cartésien 4)

Opérations structurelles: 1 (renommage)

IGE 487

Le produit est-il vraiment nécessaire?

Et le renommage?

Quel est l'ensemble de base minimal (nécessaire et suffisant)?

Cet ensemble est-il unique?

Voir

DATE, C. J.; DARWEN, H.

Databases, types, and the relational model: The third manifesto.

3rd ed., Addison-Wesley Inc., 2008.

ISBN 0-321-39942-0

Fondements — Opération de renommage

 La présence de l'entête dans chacun des tuples et chacune des relations permet de définir une opération structurelle, le renommage. Renommage R ρ A:C

Α	В	ρ A:C =	C	В
a1	b1		a1	b1
a2	b2		a2	b2
a3	b3		a3	b3

- L'entête d'une relation est conservé dans le catalogue du SGBDR.
- •Le catalogue est la description des modèles relationnels du SGBDR sous la forme d'une BD dont le modèle relationnel est lui-même dans le catalogue, comme tous les autres modèles relationnels de toutes les autres BD du SGBDR.

34

Pourquoi, et surtout comment, les opérateurs structurels (aussi appelés meta-opérateurs, tels que le renommage) peuvent-ils être considérés redondants?

Astuce: prendre le catalogue en considération.

Esquisser un modèle relationnel rudimentaire du catalogue et illustrer votre réponse.

3. Un cas
3.4 Vérification de la solution

La théorie relationnelle est une solution

• acceptable?
• adéquate?

Acceptabilité (satisfaction des exigences)

•La démonstration de la satisfaction des quatre exigences spécifiées plus tôt dépasse la portée et le temps imparti.

oLes écrits de nombreux scientifiques (et parmi les plus grands) le démontrent; écrits qui furent revus, publiés et rendus accessibles.

97-11 DEM SOMES — LESS MONIMES AUX HOUGERS... (V. 1004) & SOMES, MIJIN, — CE DE VE CON-4
Département d'informatique, l'éculié des sciences, Université de Sherbeoke, Québec

Adéquation

Critères absolus

Validité: oui Efficacité: oui Cohérence: oui

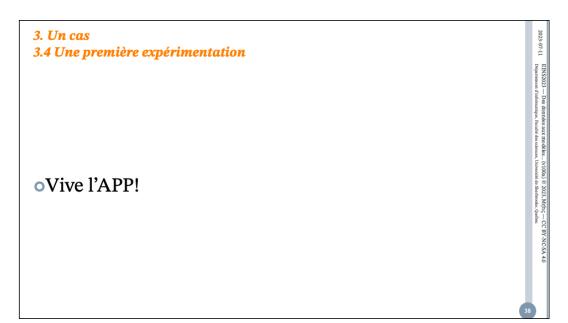
Critères relatifs

- Complétude: suffisante (complétude au sens de Turing démontrée)
- Efficience: la plus efficiente connue parmi les solutions générales
- Évolutivité: pas meilleure que les autres!

Méta-critères

- Réfutabilité...
- Acceptabilité (éthique)...
- À vous de voir, en fonction de l'épistémologie et de l'éthique choisie!

atique, Paculté des sciences, Université de Sherbrooke, Québoc



3. Un cas

3.4 Les premières conclusions

- oEn fait, il s'agit plutôt des conclusions anticipées.
- •Autrement dit, les hypothèses à vérifier lors de l'expérimentation.
- •Hypothèses découlant de l'analyse du méta-modèle et de celle de la littérature scientifique contemporaine.

יאטערעפי – Des données aux moderes... (עוניים) איי עריק – איי איי איי איי עריק. – איי איי איי איי איי איי איי urtement d'informatique, Paculté des sciences, Université de Sherbrooke, Québec

Acquis et défis

- oPourquoi présenter acquis et défis conjointement?
- •Parce que de nombreux acquis théoriques tardent encore à être mis à disposition, sinon utilisés, en pratique!

ди *soutes — ис*в мящеев вид приевев... (у 1004) % 2023, иг]ну — сс Département d'informatique, Faculté des sciences, Université de Sherbrooke, Québec

o Ambigüité:

• Une certaine amélioration grâce à l'algèbre relationnelle et aux raisonneurs, mais le problème demeure en théorie (lire notre ami Gödel) et en pratique (complexité algorithmique des algorithmes de raisonnement).

o Complétude:

• La finitude résout théoriquement le problème de complétude... mais persiste toujours en pratique (voir complexité algorithmique et volume de données).

o Données manquantes:

 Plusieurs solutions ont été proposées dont deux seulement sont solides (Codd:logique_4V et Date:décomosition); malheureusement, aucun langage utilisé couramment ne permet de les utiliser commodément. nt d'informatique, Faculté des sciences, Université de Sherbrooke, Québec

Acquis et défis (lot 2)

O Agents:

• Problème résolu en théorie — en pratique: complexité et expressivité.

• * Axes spatiaux:

 Problème résolu en théorie — en pratique: outils disponibles, mais la complexité demeure et l'expressivité ainsi que l'efficience demande encore des avancées significatives.

• Axes temporels:

• Plusieurs approches intéressantes, aucune exhaustive, peu d'outils en pratique... mais qu'est-ce que le temps?

o Cohérence:

• Que se passe-t-il si on ajoute l'exigence suivante «Capacité de transposer les règles de pratiques médicales» dans le modèle?

-11 Département d'informatique, Faculté des sciences, Université de Sherbrooke, Québec

Pistes

Utiliser les ontologies appliquées pour décrire les prédicats:

- S'appuyer sur les connaissances.
- Réduire (encore plus) les sources d'ambigüité.

Université de Sherbrooke, Québec

Trois réflexions (en guise de conclusion)

- «Du problème à la solution, en passant par les besoins, les modèles, les exigences et l'adéquation.»
- «Des données aux modèles en passant par la connaissance et les relations.»
- o«Du raisonnement aux relations en passant par la logique»

```
Merci à
Pythagore (vers -580 à -495), Socrate (vers -470 à -399),
Platon (vers -428 à -347), Aristote (vers -384 à -322),
Eulcide (vers -300)...
Descartes (1596 à 1650), Pascal (1623 à 1662), Russel (1872 à 1970),
Wittgenstein (1889 à 1951), Gödel (1906 à 1978), Floyd (1936-2001),
Codd (1923 à 2003), Chomsky (1928...), Hoare (1934...) et Date (1941...).
```

