

Semantic Science and Machine-Accessible Scientific Theories

David Poole* Clinton Smyth† Rita Sharma†

*Department of Computer Science, University of British Columbia

†Georeference Online

March 2008

Outline

- 1 Semantic Science Vision
 - Ontologies
 - Data
 - Theories
- 2 Theoretical Foundations
 - Probabilistic Prediction
 - Probabilities with Ontologies
 - Existence and Identity Uncertainty
- 3 Fielded Systems

Example: medical diagnosis

Example: people give symptoms and want to know what is wrong with them.

Current Practice (Google)	Vision
<ul style="list-style-type: none">— describe symptoms using keywords— results ranked by popularity (pagerank)— text results	

Example: medical diagnosis

Example: people give symptoms and want to know what is wrong with them.

Current Practice (Google)	Vision
<ul style="list-style-type: none">— describe symptoms using keywords— results ranked by popularity (pagerank)— text results	<ul style="list-style-type: none">— use ontologies

Example: medical diagnosis

Example: people give symptoms and want to know what is wrong with them.

Current Practice (Google)	Vision
<ul style="list-style-type: none">— describe symptoms using keywords— results ranked by popularity (pagerank)— text results	<ul style="list-style-type: none">— use ontologies— theories ranked by relevance and fit to data

Example: medical diagnosis

Example: people give symptoms and want to know what is wrong with them.

Current Practice (Google)	Vision
<ul style="list-style-type: none">— describe symptoms using keywords— results ranked by popularity (pagerank)— text results	<ul style="list-style-type: none">— use ontologies— theories ranked by relevance and fit to data— probabilistic predictions with references to sources

Example: finding a location that contains gold

Given a model of where gold can be found and 25000 location descriptions:

Current Practice	Vision
<ul style="list-style-type: none">— keyword database look-up— results (if any) unranked or ranked by popularity— text— repeat for more and less general terms	<ul style="list-style-type: none">— describe model using ontology— results ranked by fit to model— probabilistic prediction

Example: finding minerals at a location

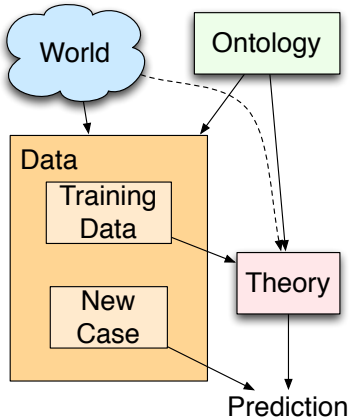
Given one location and 100 models of where minerals can be found:

Current Practice ????	Vision
<ul style="list-style-type: none">— keyword database look-up— results (if any) unranked or ranked by popularity— text	<ul style="list-style-type: none">— describe location and models using ontology— results ranked by relevance and fit to data— probabilistic prediction with references

Notational Minefield

- Theory / hypothesis / model / law (Science)
- Variable (probability and logic and programming languages)
- Model (science, probability and logic)
- Parameter (mathematics and statistics)
- Domain (science and logic and probability and mathematics)
- Object/class (object-oriented programming and ontologies)
- = (probability and logic)
- First-order (logic and dynamical systems)

Our Semantic Science Vision



- Ontologies represent the meaning of symbols.
- Data that adheres to an ontology is published.
- Theories that make (probabilistic) predictions on data are published.
- Data can be used to evaluate theories.
- Theories make predictions on new cases.

AI Traditions

- Expert Systems of 70's and 80's (e.g., Prospector '74-83)
 - Probabilistic models and machine learning.
Bayesian networks, Bayesian X...
 - Ontologies and Knowledge Representations.
Description logic, X logic...

Science in Broadest Sense

We mean *science* in the broadest sense:

- where and when landslides occur
- where to find gold
- what errors students make
- disease symptoms, prognosis and treatment
- what companies will be good to invest in
- what house Mary would like

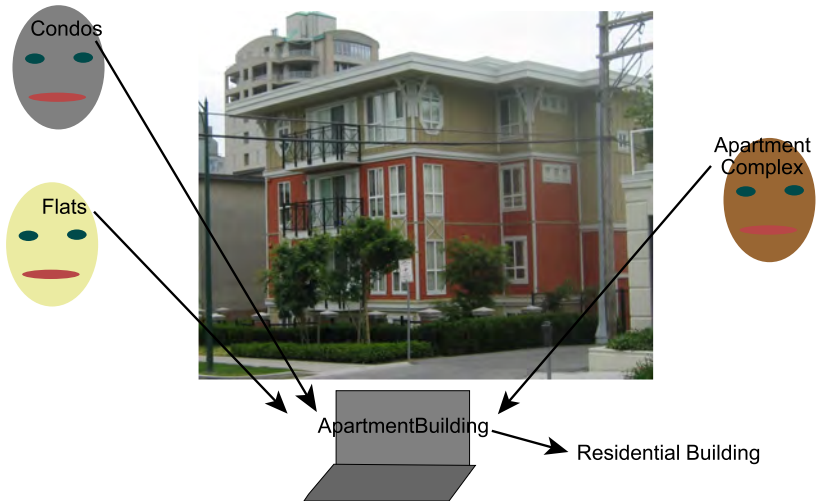
Outline

- 1 Semantic Science Vision
 - Ontologies
 - Data
 - Theories
- 2 Theoretical Foundations
 - Probabilistic Prediction
 - Probabilities with Ontologies
 - Existence and Identity Uncertainty
- 3 Fielded Systems

Ontologies

- In philosophy, **ontology** the study of existence.
- In CS, an **ontology** is a (formal) specification of the meaning of the vocabulary used in an information system.
- Ontologies are needed so that information sources can inter-operate at a semantic level.

Ontologies



Main Components of an Ontology

- **Individuals:** the objects in the world (not usually specified as part of the ontology)
- **Classes:** sets of (potential) individuals
- **Properties:** between individuals and their values

Aristotelian definitions

Aristotle [350 B.C.] suggested the definition of a class C in terms of:

- **Genus**: the super-class
- **Differentia**: the attributes that make members of the class C different from other members of the super-class

"If genera are different and co-ordinate, their differentiae are themselves different in kind. Take as an instance the genus 'animal' and the genus 'knowledge'. 'With feet', 'two-footed', 'winged', 'aquatic', are differentiae of 'animal'; the species of knowledge are not distinguished by the same differentiae. One species of knowledge does not differ from another in being 'two-footed'."

Aristotle, *Categories*, 350 B.C.

An Aristotelian definition

- An **apartment building** is a **residential building** with **multiple units** and **units are rented**.

$$\begin{aligned} ApartmentBuilding &\equiv ResidentialBuilding \& \\ &NumUnits = many \& \\ &Ownership = rental \end{aligned}$$

NumUnits is a property with domain *ResidentialBuilding* and range {*one*, *two*, *many*}

Ownership is a property with domain *Building* and range {*owned*, *rental*, *coop*}.

- All classes can be defined in terms of properties.

Outline

- 1 Semantic Science Vision
 - Ontologies
 - Data
 - Theories
- 2 Theoretical Foundations
 - Probabilistic Prediction
 - Probabilities with Ontologies
 - Existence and Identity Uncertainty
- 3 Fielded Systems

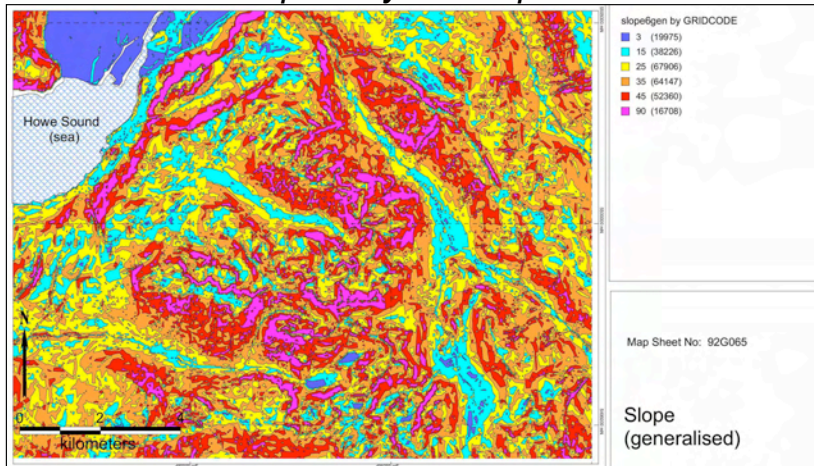
Data

Real data is messy!

- Multiple levels of abstraction
- Multiple levels of detail
- Uses the vocabulary from many ontologies: rocks, minerals, top-level ontology, . . .
- Rich meta-data:
 - Who collected each datum? (identity and credentials)
 - Who transcribed the information?
 - What was the protocol used to collect the data?
(Chosen at random or chosen because interesting?)
 - What were the controls — what was manipulated, when?
 - What sensors were used? What is their reliability and operating range?

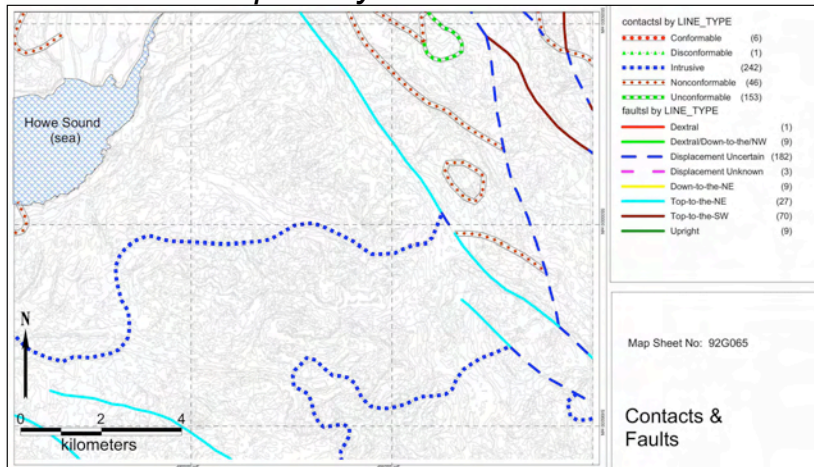
Example Data in Geology (I)

Input Layer: Slope



Example Data in Geology (II)

Input Layer: Structure



Data is theory-laden

- Sapir-Whorf Hypothesis [Sapir 1929, Whorf 1940]: people's perception and thought are determined by what can be described in their language. (Controversial in linguistics!)
- A stronger version for information systems:

What is stored and communicated by an information system is constrained by the representation and the ontology used by the information system.

- Ontologies must come logically prior to the data.
- Data can't make distinctions that can't be expressed in the ontology.
- Different ontologies result in different data.

Outline

- 1 Semantic Science Vision
 - Ontologies
 - Data
 - Theories
- 2 Theoretical Foundations
 - Probabilistic Prediction
 - Probabilities with Ontologies
 - Existence and Identity Uncertainty
- 3 Fielded Systems

Theories make predictions on data

A **theory** is a procedure that makes a prediction on data.

- Theories can make whatever predictions they like about data:
 - definitive predictions
 - point probabilities
 - probability ranges
 - ranges with confidence intervals
 - qualitative predictions
- For each prediction type, we need ways to judge predictions on data
- Users can use whatever criteria they like to evaluate theories (e.g., taking into account simplicity and elegance)

Theory Ensembles

- How can we compare theories that differ in their generality?
- Theory A makes predictions about all cancers. Theory B makes predictions about lung cancers. Should the comparison between A and B take into account A 's predictions on non-lung cancer?

Theory Ensembles

- How can we compare theories that differ in their generality?
- Theory A makes predictions about all cancers. Theory B makes predictions about lung cancers. Should the comparison between A and B take into account A 's predictions on non-lung cancer?
- What about theory C : *if lung cancer, use B 's prediction, else use A 's prediction?*

Theory Ensembles

- How can we compare theories that differ in their generality?
- Theory A makes predictions about all cancers. Theory B makes predictions about lung cancers. Should the comparison between A and B take into account A 's predictions on non-lung cancer?
- What about theory C : *if lung cancer, use B 's prediction, else use A 's prediction?*
- Proposal: make **theory ensembles** the norm.
 - Judge theories by how well they fit into ensembles.
 - Ensembles can be judged by simplicity.
 - Theory designers don't need to game the system by manipulating the generality of theories

Dynamics of Semantic Science

- Anyone can design their own ontologies.
 - People vote with their feet what ontology they use.
 - Need for semantic interoperability leads to ontologies with mappings between them.
- Ontologies evolve with theories:
 - A theory hypothesizes unobserved features or useful distinctions
 - add these to an ontology
 - other researchers can refer to them
 - reinterpretation of data
- Ontologies can be judged by the predictions of the theories that use them
 - the role of the vocabulary is to describe useful distinctions.

Outline

- 1 Semantic Science Vision
 - Ontologies
 - Data
 - Theories
- 2 Theoretical Foundations
 - Probabilistic Prediction
 - Probabilities with Ontologies
 - Existence and Identity Uncertainty
- 3 Fielded Systems

Why Probabilistic Prediction?

- Probabilities are what you get from data.
(Most suggested measures of prediction accuracy are optimized by probabilistic prediction!)
- There is a well defined procedure for combining background knowledge with data (conditioning).
- Probabilities are what is needed (with utilities) to make decisions.

Probabilistic Prediction

- The role of models in prediction:
Given a description of a new case,

$$P(\text{prediction}|\text{description}) \\ = \sum_{m \in \text{Models}} \left(\frac{P(\text{prediction}|m \& \text{description}) \times P(m|\text{description})}{P(m|\text{description})} \right)$$

Models is a set of mutually exclusive and covering set of hypotheses.

Probabilistic Prediction

- The role of models in prediction:
Given a description of a new case,

$$P(\text{prediction}|\text{description}) \\ = \sum_{m \in \text{Models}} \left(\frac{P(\text{prediction}|m\&\text{description}) \times P(m|\text{description})}{P(m|\text{description})} \right)$$

Models is a set of mutually exclusive and covering set of hypotheses.

- What features of the description are predictive?
- How do the features interact?
- What are the appropriate probabilities? (How can these be learned with limited data?)

Outline

- 1 Semantic Science Vision
 - Ontologies
 - Data
 - Theories
- 2 Theoretical Foundations
 - Probabilistic Prediction
 - **Probabilities with Ontologies**
 - Existence and Identity Uncertainty
- 3 Fielded Systems

Random Variables and Triples

- Reconcile:
 - random variables of probability theory
 - individuals, classes, properties of modern ontologies

Random Variables and Triples

- Reconcile:
 - random variables of probability theory
 - individuals, classes, properties of modern ontologies
- For functional properties:
random variable for each $\langle individual, property \rangle$ pair,
where the domain of the random variable is the range of
the property.
- For non-functional properties:
Boolean random variable for each
 $\langle individual, property, value \rangle$ triple.

First-order probabilistic models

- Individuals are not known until run time.
- Therefore the random variables are not known until run time (and they change for each situation).
- We want to build the models before we know the random variables.

First-order probabilistic models

- Individuals are not known until run time.
- Therefore the random variables are not known until run time (and they change for each situation).
- We want to build the models before we know the random variables.

→ First-order probabilistic models

- Idea: if you are a Bayesian, you need to treat every individual that you have the same knowledge about the same (exchangability).
- Probabilities are specified for all individuals.

Probabilities and Aristotelian Definitions

Aristotelian definition

$$\begin{aligned} \textit{ApartmentBuilding} &\equiv \textit{ResidentialBuilding} \& \\ &\textit{NumUnits} = \textit{many} \& \\ &\textit{Ownership} = \textit{rental} \end{aligned}$$

leads to probability over property values

$$\begin{aligned} &P(\langle A, \textit{type}, \textit{ApartmentBuilding} \rangle) \\ &= P(\langle A, \textit{type}, \textit{ResidentialBuilding} \rangle) \times \\ &\quad P(\langle A, \textit{NumUnits}, \textit{many} \rangle \mid \langle A, \textit{type}, \textit{ResidentialBuilding} \rangle) \times \\ &\quad P(\langle A, \textit{Ownership}, \textit{rental} \rangle \mid \langle A, \textit{NumUnits}, \textit{many} \rangle, \\ &\quad \quad \langle A, \textit{type}, \textit{ResidentialBuilding} \rangle) \end{aligned}$$

Type uncertainty \longrightarrow uncertainty over property values.

Outline

- 1 Semantic Science Vision
 - Ontologies
 - Data
 - Theories
- 2 Theoretical Foundations
 - Probabilistic Prediction
 - Probabilities with Ontologies
 - Existence and Identity Uncertainty
- 3 Fielded Systems

Existence and Identity Uncertainty

Theory about what house Mary would like:

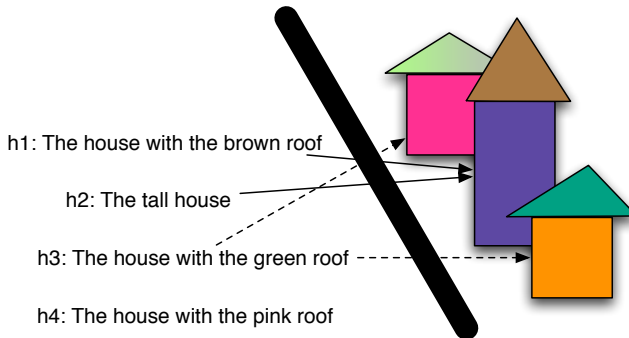
Whether Mary likes an house depends on:

- Whether there is a bedroom for daughter Sam
- Whether Sam's room is green
- Whether there is a bedroom for Mary
- Whether Mary's room is large
- Whether they share

Existence and Identity

Symbols

Individuals



Clarity Principle

Clarity principle: probabilities must be over well-defined propositions.

- What if an individual doesn't exist?
 - $house(h4) \wedge roof_colour(h4, pink) \wedge \neg exists(h4)$

Clarity Principle

Clarity principle: probabilities must be over well-defined propositions.

- What if an individual doesn't exist?
 - $house(h4) \wedge roof_colour(h4, pink) \wedge \neg exists(h4)$ **X**

Want: probability that there exists an object that matches some description. Name the the object that exists.

Clarity Principle

Clarity principle: probabilities must be over well-defined propositions.

- What if an individual doesn't exist?
 - $house(h4) \wedge roof_colour(h4, pink) \wedge \neg exists(h4)$ **X**
Want: probability that there exists an object that matches some description. Name the the object that exists.
- What if more than one individual exists? Which one are we referring to?
 - In a house with three bedrooms, which is the second bedroom?

Clarity Principle

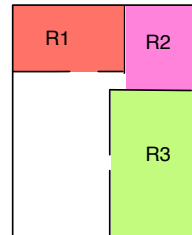
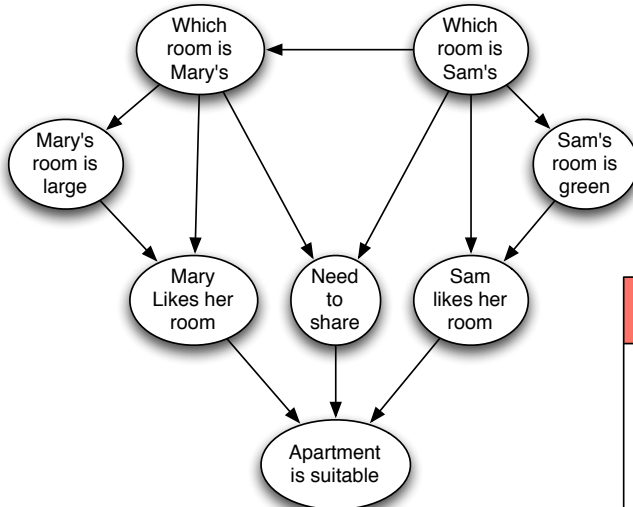
Clarity principle: probabilities must be over well-defined propositions.

- What if an individual doesn't exist?
 - $house(h4) \wedge roof_colour(h4, pink) \wedge \neg exists(h4)$ **x**

Want: probability that there exists an object that matches some description. Name the the object that exists.

- What if more than one individual exists? Which one are we referring to?
 - In a house with three bedrooms, which is the second bedroom?
- **Note:** Reified individuals are special:
 - Non-existence means the relation is false.
 - Well defined what doesn't exist when existence is false.
 - Same description implies the same individual.

Role assignments



Outline

- 1 Semantic Science Vision
 - Ontologies
 - Data
 - Theories
- 2 Theoretical Foundations
 - Probabilistic Prediction
 - Probabilities with Ontologies
 - Existence and Identity Uncertainty
- 3 Fielded Systems

Expert Models

What if the models are provided by the experts in the field?

- not covering — only provide positive models
- not exclusive — they are often refinements of each other
- described at various levels of abstraction and detail
- often the experts don't know the probabilities and there is little data to estimate them

Providing Probabilities

Experts are reluctant to give probabilities:

- No data from which to estimate them
- People who want to make decision use more information than provided in our theories
- Difficult to combine marginal probabilities with new information to make decisions
- It is *not* because decision theory is inappropriate. Decision makers use probabilities and utilities.

What we do

- Use qualitative probabilities: $\{ \textit{always}, \textit{usually}, \textit{sometimes}, \textit{rarely}, \textit{never} \}$.
- With thousands of instances and hundreds of models, find the most likely and the rationale.
- Independence assumptions.

Example Model

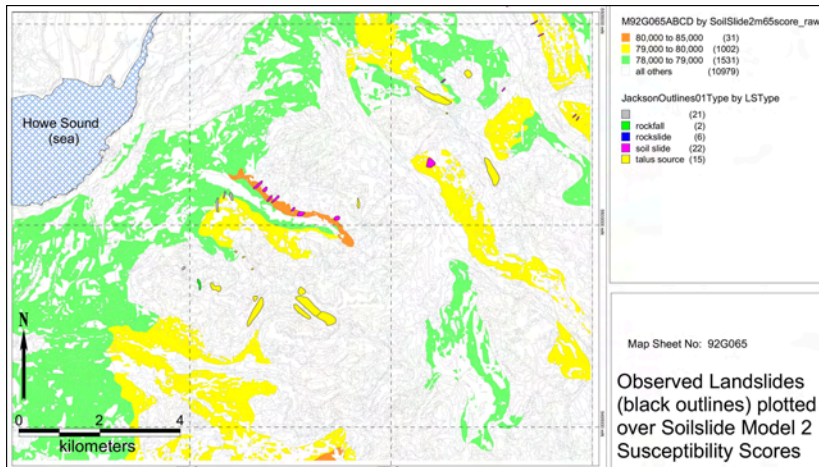
Prototype SoilSlide Model (Jackson, 2007)

	Description	Presence	Comment
Bedrock	SoilSlide01	model	
Terrain	SoilSlide02	model	
Primary	Component - Component1	always	Secondary Primary Terrain unit is USUALLY C if Primary is R (This is the Primary component)
SOMETIMES	Layer - Layer 1	always	Minor terrain unit will ALWAYS be M or C if Major Terrain Unit is R alone
Commercial	SurficialMaterial - Bedrock	always	Minor terrain unit will ALWAYS be M or C if Major Terrain Unit is R alone
areas of	SurficialMaterial - <other values>	never	Minor terrain unit will ALWAYS be M or C if Major Terrain Unit is R alone
Secondary	Component - Component2	always	Secondary Primary Terrain unit is USUALLY C if Primary is R (This is the Secondary component)
USUALLY	Layer - Layer 1	always	
Minor te	SurficialMaterial - Colluvium	always	Minor terrain unit will ALWAYS be M or C if Major Terrain Unit is R alone
C if Maj	Slope - Gentle	never	NEVER on slopes 14 degrees or less
Thus, we	Slope - Plain	never	NEVER on slopes 14 degrees or less
this by sayin	Slope - Moderate	usually	USUALLY on slopes between 20 and 40 degrees
ALWAYS ass	Slope - Moderately Steep	usually	USUALLY on slopes between 20 and 40 degrees
that contain	Slope - Steep	rarely	RARELY on slopes 41 to 60 degrees
whether the	Slope - Very Steep	never	RARELY on slopes 41 to 60 degrees
components	SurficialMaterial - Morainal Material (Till)	always	Minor terrain unit will ALWAYS be M or C if Major Terrain Unit is R alone
	GeomorphProcess - Gully Erosion	sometimes	SOMETIMES associated with V or A
	GeomorphProcess - SnowAvalanches	sometimes	SOMETIMES associated with V or A

19
40
it
ative
ve
will be
slips

Example Prediction from a Model

Test Results: Model SoilSlide02



Conclusion

- Demand from funders, scientists and users.
- Complementary to Semantic web.
- Representing, reasoning and learning complex probabilistic theories is largely unexplored.
- Still lots of work to be done!

To Do

- Fundamental research on complex probabilistic models.
- Build infrastructure to allow publishing and interaction of ontologies, data, theories, theory ensembles, evaluation criteria, meta-data.
- Build inverse semantic science web:
 - Given a theory, find relevant data
 - Given data, find theory ensembles
 - Given a new case, find relevant theory ensembles with explanations