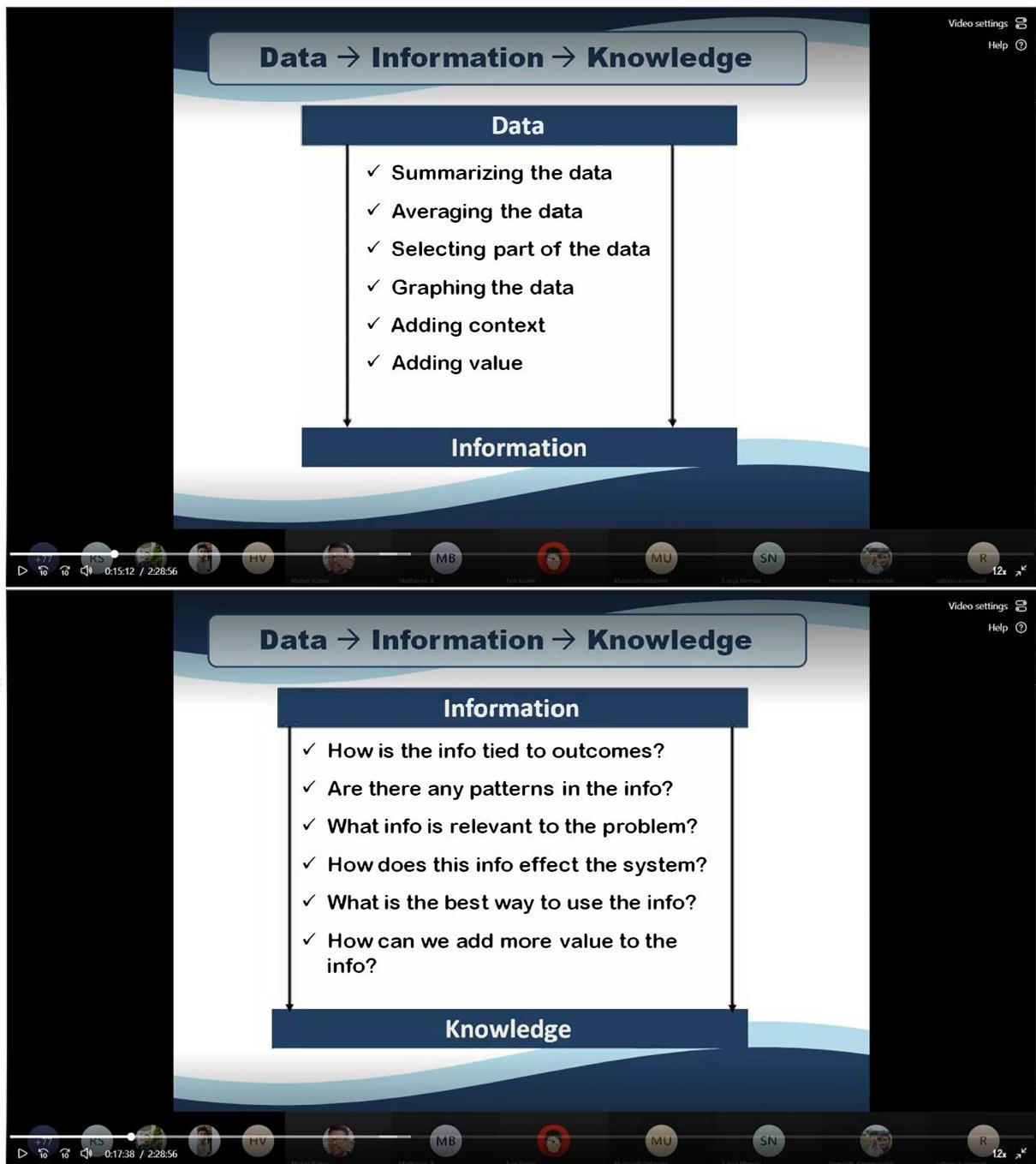


Semantic Information Extraction Approaches for Industries

Dr. K. Rajbabu
Manager
Controls & Instrumentation / FB
BHEL, Trichy

Overview of the Discussion

- About Information
- Extraction vs Retrieval
- Information Extraction (IE) Approaches
- Information Extraction for Business Intelligence
- Knowledge Requirement
- Semantic Knowledge
- Ontology Model
- OBIE Approaches
- Recent Developments in Information Extraction



Information Extraction

- **The process of obtaining pertinent information (facts) from documents.**

- Examples: *The forest area* in India extended to about *75 million hectares*, which in terms of **geographical area** is approximately 22 percent of the total land.

- **What is the area of forest in India?**



Extraction Vs Retrieval

- **Retrieval:** Obtaining **relevant** information from a collection of information resources. Searches can be based on metadata or on full-text indexing.
 - gets sets of relevant documents
 - Ex. Google
- **Extraction:** automatically extracting **structured information** from unstructured and/or semi-structured machine-readable documents.
 - Extraction gets facts out of documents
 - Google result of currency rate, gold rate etc
 - NLP and Machine Learning solution



Information Retrieval

- **DATABASES**
 - Structured schema in advance
 - Semantic correlation between queries and data is clear.
 - Strong theoretical foundation
- **IR**
 - No schema (unstructured natural language text)
 - Unclear semantic correlation between queries and data
 - We get inexact, approximated answers
 - Theory not well understood (especially Natural Language Processing)

Information Extraction (IE) Example

Bridgestone Sports Co. said Friday it had set up a joint venture in Taiwan with a local concern and a Japanese trading house to produce golf clubs to be supplied to Japan.

The joint venture, Bridgestone Sports Taiwan Co., capitalized at 20 million new Taiwan dollars, will start production in January 1990 with production of 20,000 iron and "metal wood" clubs a month.

TIE-UP-1 Relationship: TIE-UP Entities: "Bridgestone Sport Co." "a local concern" "a Japanese trading house" Joint Venture Company: "Bridgestone Sports Taiwan Co." Activity: ACTIVITY-1 Amount: NT\$200000000	ACTIVITY-1 Activity: PRODUCTION Company: "Bridgestone Sports Taiwan Co." Product: "iron and 'metal wood' clubs" Start Date: DURING: January 1990
--	--

Information Extraction Steps

- Pre-processing
- Noise Removal
- Concept Identification
- Relating the Concepts
- Unifying Results

▷ ⏴ 18 🔍 KS 🔍 0:30:57 / 2:28:56 HV SS MB MU SN R 12x ↗

Challenges in Information Extraction

- Junks
- Relevancy
- Prioritization
- Context relativity

▷ ⏴ 18 🔍 KS 🔍 0:31:09 / 2:28:56 HV SS MB MU SN R 12x ↗

The image consists of two vertically stacked screenshots of a video player interface, likely from a video conference or presentation.

Top Screenshot:

- Section Header:** IE Approaches
- List:**
 - Unsupervised
 - Wrapper based Extraction
 - Natural Language Processing
 - Stochastic Model (statistical or probabilistic)
 - Supervised
 - Machine Learning Techniques

Bottom Screenshot:

- Section Header:** IE Approaches
- List:**
 - Heuristic Learning based Approach
 - Rule based Approach
 - Classification based Approach
 - Knowledge based Approach

Common Interface Elements:

- Video settings icon
- Help icon
- Progress bar: 0:32:28 / 2:28:56
- Participants list: HV, SS, MB, MU, SN, R
- Speaker icon: HV
- Microphone icon: MB
- Video icon: MU
- Audio icon: SN
- Indicator icon: R
- Volume icon: 12x

Video settings Help

Taxonomy

```

graph TD
    DM[Data Mining] --> S[Supervised]
    DM --> SS[Semi-supervised]
    DM --> U[Unsupervised]
    S --> C[Classification]
    S --> R[Regression]
    SS --> H[Hybrid]
    U --> CL[Clustering]
    U --> AS[Association]
    C --> NBC[Naive Bayes Classifier]
    C --> NN[Neural Network]
    C --> DT[Decision Tree]
    C --> SVM[Support Vector Machines (SVMs)]
    C --> KNN[K-Nearest Neighbors]
    R --> LR[Linear Regression]
    R --> LRg[Logistic Regression]
    H --> S_H[Supervised]
    H --> U_H[Unsupervised]
    CL --> KM[K-mean]
    CL --> KMd[K-medoid]
    CL --> ML[Multilevel]
    AS --> AP[Apriori]
  
```

0:35:39 / 2:28:56

Surendra SS MB MU SN R 12x

Video settings Help

Vector Space Model

- Algebraic model of representing text documents
- Term- Document Matrix
- Similarity Measures: cosine similarity

```

graph TD
    D([Documents]) --> DP[Data Preprocessing: remove punctuations, stop words, to lower]
    SQ([Search Queries]) --> DP
    DP --> TDM[Term Document Matrix]
    TDM --> DTf[Document TfIdf]
    TDM --> QTf[Queries TfIdf]
    DTf --> CCS[Calculate Cosine similarity between each document and each query]
    QTf --> CCS
    CCS --> SRTK[Sort and select Top K relevant documents with large cosine similarity values]
  
```

0:35:41 / 2:28:56

Surendra SS MB MU SN R 12x

Video settings Help

Probabilistic Model

Ex.: Mission to Mars, Playing Chess

Assume data $\mathbf{X} = \{x_n\}_{n=1}^N$ generated from a **probability distribution** $p(x|\theta)$, in an i.i.d. (independent and identically distributed) fashion

$$x_1, \dots, x_N \sim p(x|\theta)$$

The form of $p(x|\theta)$ (also called **likelihood**) depends on the type of the data

Assumptions about parameter θ can be encoded via a **prior distribution** $p(\theta)$

- Also corresponds to imposing a regularizer over θ (helps in generalization)

Goal: To **estimate parameter θ** , given data \mathbf{X}

Variations of this general view subsume most machine learning problems

- Regression, classification, clustering, dimensionality reduction, etc.

Video settings Help

Probabilistic Model

<u>Distribution</u>	<u>Domain</u>	<u>Picture</u>	<u>Parametric Form</u>
Binomial	Binary		$Bin(x N,\theta) \propto \theta^x (1-\theta)^{N-x}$
Multinomial	K classes		$Mult(x \theta) \propto \prod \theta_k^{x_k}$
Beta	[0,1]		$Beta(\theta \alpha,\beta) \propto \theta^{\alpha-1} (1-\theta)^{\beta-1}$
Gamma	[0,∞)		$Gam(x \alpha,b) \propto x^{\alpha-1} \exp(-bx)$
Dirichlet	Simplex		$Dir(\theta \alpha) \propto \prod \theta_k^{\alpha_k-1}$
Gaussian	Reals		$Nor(x \mu,\sigma^2) \propto \exp(-(x-\mu)^2/2\sigma^2)$

Video settings Help

Probabilistic Model

- Naive Bayes classifier is based on the Naive Bayes algorithm $p(y|x) = \frac{p(y)p(x|y)}{p(x)}$
- Logistic Regression classifier is based on the Logistic function $f(x) = \frac{1}{1 + e^{-x}} = \frac{e^x}{1 + e^x}$
- Conditional Random Fields (CRF) $P(y, X, \lambda) = \frac{1}{Z(X)} \exp\left\{\sum_{i=1}^n \sum_j \lambda_j f_i(X, i, y_{i-1}, y_i)\right\}$
Where: $Z(x) = \sum_{y \in \mathcal{Y}} \sum_{i=1}^n \sum_j \lambda_j f_i(X, i, y'_{i-1}, y'_i)$

Support Vector Machine(SVM)

- A Discriminative classifier which takes training data, algorithm and outputs a optimal hyperplane which categorizes new examples

A line drawn between these data points classify the black dots and blue squares.

Support Vector Machine(SVM)

The first scatter plot shows two distinct clusters of points (purple and blue) separated by a vertical line, labeled "Linearly separable data". The second scatter plot shows points of various colors (purple, yellow, green) overlapping significantly, labeled "Non linearly separable data".

Linearly separable data Non linearly separable data

Support Vector Machine(SVM)

The left diagram, labeled "Raw Data", shows two classes of points (black dots and blue squares) in a 2D space defined by red axes. The right diagram, labeled "Line as Hyperplane", shows the same data points with a red line drawn through them, representing the decision boundary. A blue speech bubble labeled "Outliers" points to a purple point located on the wrong side of the line.

Raw Data Line as Hyperplane

Support Vector Machine(SVM)

- Tuning Parameters
 - Kernel
 - Regularization
 - Gamma
 - Margin

Maximum margin decision hyperplane

Support vectors

Margin is maximized

Video settings Help

0:41:44 / 2:28:56

SS MB MU SN R

Support Vector Machine(SVM)

- Kernel – Mathematical function for transforming data using linear algebra.
- Different SVM Algorithm use different kernel functions
 - Linear
 - NonLinear
 - Radial Basis Functions (RBF)
 - Sigmoid
 - Polynomial
 - Exponential

Video settings Help

0:41:58 / 2:28:56

SS MB MU SN R

Video settings Help

Decision Tree

- Data is break down by making decisions using series of conditions.
- Classification and Regression algorithms
- Business Management
- CRM
- Fraudulent Statement
- Energy Consumption

The diagram illustrates a decision tree structure. It starts with a single blue square at the top labeled "root". This root branches into two blue squares, which further branch into four blue squares, and so on. Some of these blue squares are labeled "node". At the bottom level, there are several light blue squares labeled "leaf". Arrows point from the labels "root", "node", and "leaf" to their respective parts of the tree.

Video settings Help

Neural Networks

The page displays four types of neural network architectures:

- Single Layer Feedforward Network:** Shows an input layer with three green nodes and an output layer with three pink nodes. Every node in the input layer is connected to every node in the output layer.
- Multi-Layer Feedforward Network:** Shows an input layer with three green nodes, two hidden layers (Hidden Layer 1 with three orange nodes and Hidden Layer 2 with three pink nodes), and an output layer with two blue nodes. Every node in one layer is connected to every node in the next layer.
- Recurrent Network:** Shows a circular structure where each node (green, orange, blue, pink) is connected to both the node above it and the node below it, forming a loop.
- Lattice Network:** Shows a grid of nodes. The input layer has nodes $x_1, x_2, x_3, \dots, x_n$. The hidden layer has nodes $\mu_{11}, \mu_{12}, \mu_{13}, \dots, \mu_{1n}$. The output layer has nodes $y_1, y_2, y_3, \dots, y_n$. Connections are shown between adjacent nodes in both horizontal and vertical directions.

Video settings Help

Neural Networks - Application

- **Classification:** Assigning each object to a known specific class
- **Clustering:** Grouping together objects similar to each other
- **Pattern Association:** Presenting of an input sample triggers the generation of specific output pattern
- **Function approximation:** Constructing a function generating almost the same outputs from input data as the modeled process
- **Optimization:** Optimizing function values subject to constraints
- **Forecasting:** Predicting future events on the basis of past history
- **Control:** Determining values for input variables to achieve desired values for output variables

Deep Learning

Machine Learning

```

graph LR
    Input[Input] --> Feature[Feature extraction]
    Feature --> Classification[Classification]
    Classification --> Output[Output]
    
```

Deep Learning

```

graph LR
    Input[Input] --> Feature[Feature extraction + Classification]
    Feature --> Output[Output]
    
```

The image consists of two vertically stacked screenshots of a video player interface, likely from a platform like YouTube or a similar video hosting service. Both screenshots show a presentation slide with the title "Natural Language Processing" at the top.

Screenshot 1 (Top):

- Sentence Segmentation
- Tokenization
- Part of Speech (POS) tagging
- Named Entity Recognition (NER)
- Chunk parsing
- Relation Detection
- Leading open-source tool: GATE/ANNIE, MonkeyLearn, Aylien, IBM Watson, Google Cloud, Amazon Comprehend
- Python Tools : NLTK, Spacy, TextBlob, Textacy, PyTorch-NLP
- JAVA tools: OpenNLP, StanfordNLP, CogCompNLP

Screenshot 2 (Bottom):

- Pattern matching & rule learning (regular expressions, FSAs)
- Statistical learning (HMMs, MRFs, etc.)
- Lexicon lookups (name dictionaries, geo gazetteers, etc.)
- Text mining in general

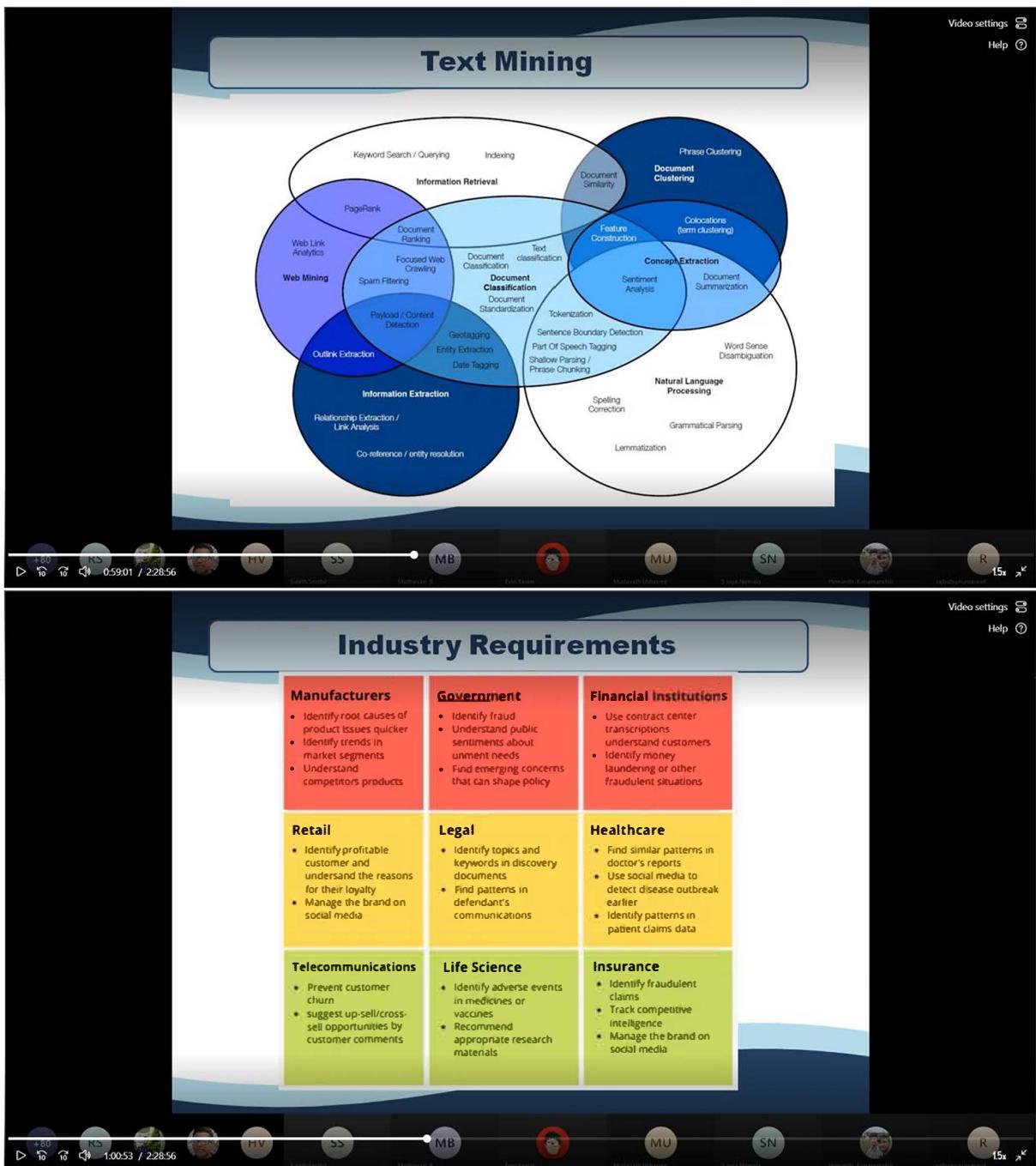
Both screenshots include a standard video player control bar at the bottom, showing icons for play, pause, volume, and other media controls. The video progress bar indicates the current time as 0:45:27 / 2:28:56. The title "Natural Language Processing" is displayed in a blue header bar on both slides.

Discussion

- Supervised learning outperforms compared to unsupervised learning
- Domain specific approach
- Hybrid solutions
- Issue based approach based on inputs and expected outcome

Text Mining

The diagram consists of several overlapping blue circles. The main central circle is labeled 'Text Mining'. Other circles include 'Data Mining' (top left), 'AI and Machine Learning' (top right), 'Computational Linguistics' (bottom right), 'Statistics' (top right), 'Library and Information Sciences' (bottom left), and 'Databases' (middle left). Various sub-fields are listed within the intersections: 'Information Extraction' (between Data Mining and AI), 'Natural Language Processing' (between AI and Computational Linguistics), 'Concept Extraction' (between Computational Linguistics and Text Mining), 'Web Mining' (between Text Mining and Databases), 'Information Retrieval' (between Databases and Library Sciences), 'Document Clustering' (between Data Mining and Text Mining), and 'Document Classification' (between Data Mining and AI).



Industrial Needs for Information Extraction

Process of IE Solution

- Identify business problem
- Analyze data
- Find samples
- Identify / Standardize formats
- Get together domain experts, technical staff, NLP engineers and potential users
- Transform the problem in to IE task
- Clear the annotation types
- Clear the annotation guidelines
- Apply suitable algorithms
- Iterate, evaluate and improvise
- Ensure continuous enhancement through training

Need of Information Extraction in Industry

- Increasing thrust towards digitization
- Cut-Throat Competitiveness
- Showcase technical capability in Business
- Industry 4.0

Video settings  Help 

Digitization

- **Digitization:** conversion of text, pictures or sound into a digital form that can be processed by a computer. E.g. scanning
- In India, 53% of industries are using data analytics
- More than 90% of the data can be used for decision making
- 33% of the industries plan to invest more than 8% of their annual revenue in digitization

Courtesy: *Digitization in Industrial sector in India to grow to 65% in next five years Jul 21, 2016* <http://pwc.in>

Video settings  Help 

Unstructured Text

Unstructured Text:
Text-oriented; complex data structures which are hidden, collapsed & highly interpretable in nature e.g. e-mails, reports, business application documents etc.

Industry	Structured Data (%)	Unstructured Data (%)	Semi-Structured Data (%)
Retail	67	17	16
Travel/Hospitality/Airlines	60	21	19
Energy & Resources	60	25	15
Insurance	58	21	21
Consumer Goods	55	18	27
Banking/ Financial Services	55	26	20
Life Sciences	52	29	19
Manufacturing	52	28	20
Total	51	27	21
High Tech	45	30	25
Telecommunications	45	28	27
Utilities	39	40	21
Media and Entertainment	33	58	9

*80% of business-relevant information originates in unstructured form, primarily text
Courtesy: Big Data Study <http://sites.tcs.com/big-data-study/industries-unstructured-data/> 2013

Information Extraction

Method	Percentage
Text Based IE	49%
Image Based IE	27%
Video Based IE	15%
Audio Based IE	9%

■ Text Based IE ■ Image Based IE ■ Audio Based IE ■ Video Based IE

Unstructured Text

Characteristics of unstructured data

- Multiple formats (text, images, audio, video, blogs, and websites, etc.)
- Schema-less due to non-standardization
- Diverse sources (e.g. social media, clouds, sensors, etc.)

Courtesy: Digitization in Industrial sector in India to grow to 65% in next five years Jul 21, 2016 <http://pwc.in>

Jupyter NLTK_Mining Last Checkpoint: 13 hours ago (autosaved)

```
In [6]: 1: # Importing necessary Library
2: import pandas as pd
3: import numpy as np
4: import nltk
5: import os
6: nltk.download('punkt')
7: import nltk.corpus.sample_text for performing tokenization
8: text = "In Brazil they drive on the right-hand side of the road. Brazil has a large coastline on the eastern side of South America"
9: from nltk.tokenize import word_tokenize # Passing the string text into word tokenize for breaking the sentences
10: token = word_tokenize(text)
11: token
[12: ]
[nltk_data] Downloading package punkt to
[nltk_data]   C:\Users\rajbabu\AppData\Roaming\nltk_data...
[nltk_data] Unzipping tokenizers\punkt.zip.
```

Out[6]: ['In', 'Brazil', 'they', 'drive', 'on', 'the', 'right-hand', 'side', 'of', 'the', 'road']

The image shows two screenshots of a Jupyter Notebook interface, likely from a video recording, demonstrating NLTK text mining. Both screenshots show the same notebook environment with a Python 3 kernel.

Screenshot 1:

- Code Cell [7]:**

```
# finding the frequency distinct in the tokens
# Importing FreqDist Library from nltk and passing token into FreqDist
from nltk.probability import FreqDist
fdist = FreqDist(token)
fdist
```
- Output [7]:**

```
FreqDist({'the': 3, 'Brazil': 2, 'on': 2, 'side': 2, 'of': 2, 'In': 1, 'they': 1, 'drive': 1, 'right-hand': 1, 'road': 1, ...})
```
- Code Cell [8]:**

```
# To find the frequency of top 10 words
```

Screenshot 2:

- Code Cell [7]:**

```
# finding the frequency distinct in the tokens
# Importing FreqDist library from nltk and passing token into FreqDist
from nltk.probability import FreqDist
fdist = FreqDist(token)
fdist
```
- Output [7]:**

```
FreqDist({'the': 3, 'Brazil': 2, 'on': 2, 'side': 2, 'of': 2, 'In': 1, 'they': 1, 'drive': 1, 'right-hand': 1, 'road': 1, ...})
```
- Code Cell [8]:**

```
# To find the frequency of top 10 words
fdist1 = fdist.most_common(10)
```
- Output [8]:**

```
[('the', 3),
 ('Brazil', 2),
 ('on', 2),
 ('side', 2),
 ('of', 2),
 ('In', 1),
 ('they', 1),
 ('drive', 1),
 ('right-hand', 1)]
```

Limitations

- **Unstructured big data issues** - Noise, Data quality, Data diversity, Data dimensionality, Data sparsity, Data modeling
- **Unstructured data usability** - Variation in perspective, Data ambiguities, Semantic understanding, Context understanding, Relevance from user's perspective
- **Language and domain issues** - Lack of multilingual system, Poor morphological languages, Language ambiguity, Language modeling, Language knowledge, Language understanding, Domain-specific data
- **Capability issues** - Volume of unstructured big data, Optimal feature extraction and selection, Automatic and semantic labeling, Lack of large labeled corpus, Limitations of ML and rule-based techniques, Selection of technique

Proposed Methodology

- Advanced data pre-processing
- IE should meet pragmatics with semantics
- Advance cross and multilingual systems
- Advanced hybrid IE techniques.

Business Requirements

- Text Summarization
- Content Recommendation
- Document classification
- Topic Extraction
- Document Search and Retrieval
- Question answering
- Sentiment Analysis

Video settings Help

Why Supervised Approach?

- Need for accuracy rather than approx.
- Uses prevailing knowledge for decision

SO, HOW TO REPRESENT AND USE KNOWLEDGE?

Video settings Help

The image shows a video player interface with two slides displayed vertically.

Common Knowledge Representations

- **Lists** - linked lists are used to represent hierarchical knowledge
- **Trees** - graphs which represent hierarchical knowledge.
- **Schemas** - used to represent commonsense or stereotyped knowledge.
- **Rule-based representations**
- **Logic-based representations**
- **Semantic networks** - nodes and links - stored as propositions.

50

Video settings Transcript Help

1:47:27 / 2:28:56

Semantic Requirement

"I have a PhD"
and
"I am a doctor,"

- These two **semantically different** entities, represent the **same concept**.
- Database retrieval based on **Concepts** will extract only a subset of **Knowledge**.
- One needs retrieval based on **Semantics**, to extract all the necessary **Knowledge!**

1:47:44 / 2:28:56

Video settings Transcript Help

Video settings Transcript Help

Semantic Network

- Semantic networks or concept maps reflect cognition.
- Visual models of concepts of some specific domain connected by some type of relationship (link/arc).
- Can be easily extended and are easy to learn.

```

graph TD
    Taste -- "has attribute" --> variety
    Shape -- "has attribute" --> variety
    Color -- "has attribute" --> variety
    variety -- "has attribute" --> variety
    Mango -- "a kind of" --> variety
    Grape -- "a kind of" --> variety
    Apple -- "a kind of" --> variety
  
```

Video settings Transcript Help

Semantic Network

```

graph TD
    AS[Agricultural Science] -- "is branch of" --> PP[Plant Pathology]
    PP -- "also known as" --> PP2[phytopathology]
    PP -- "deals with" --> PD[Plant Disease]
    PP -- "deals with" --> PDM[Plant Disease Management]
    PD -- "is caused by" --> P[Pathogen]
    P -- "a kind of" --> BP[Biotic Pathogen]
    P -- "a kind of" --> AP[Abiotic Pathogen]
    BP -- "a kind of" --> F[Fungi]
    BP -- "a kind of" --> B[Bacteria]
    BP -- "a kind of" --> N[Nematode]
    AP -- "a kind of" --> F2[Frost]
    AP -- "a kind of" --> AP2[Air Pollutant]
    AP -- "a kind of" --> T[Toxicant]
  
```

Video settings Transcript Help

Semantic Web Languages

- **RDF (Resource Description Framework)**

- Triples: <subject> <property> <object>
- RDF is a data model for objects ("resources") and relations between them. These data models can be represented in an XML syntax.



Semantic Web Languages

- **RDFS (RDF Schema)**

- A vocabulary for describing properties (subclass, subproperty, domain, range) and classes of RDF resources, with a semantics for generalization-hierarchies of such properties and classes.



The image shows two consecutive frames from a video player, both displaying a slide titled "Semantic Web Languages".

Top Frame:

- Title:** Semantic Web Languages
- Section:** • OWL (Web Ontology Language)
- List:**
 - OWL adds more vocabulary for describing properties and classes
 - Relation between classes (e.g. disjointness)
 - cardinality (e.g. "exactly one"), equality, richer typing of properties
 - characteristics of properties (e.g. symmetry)
 - enumerated classes.
 - There are constraints on classes and the types of relationships permitted between them. These provide semantics by allowing systems (reasoners) to infer additional information and provide classification based on the data explicitly provided.

Video controls at the bottom include: play/pause, volume, full screen, and a progress bar showing 1:52:35 / 2:28:56. Other controls include: Video settings, Transcript, Help, and a zoom icon (1x).

Bottom Frame:

 - Title:** Semantic Web Languages
 - Section:** • OWL (Web Ontology Language)
 - List:**
 - There are constraints on classes and the types of relationships permitted between them. These provide semantics by allowing systems (reasoners) to infer additional information and provide classification based on the data explicitly provided.

Video controls at the bottom include: play/pause, volume, full screen, and a progress bar showing 1:53:36 / 2:28:56. Other controls include: Video settings, Transcript, Help, and a zoom icon (1x).

The image shows a video player interface with two slides displayed sequentially.

Top Slide: OWL Types

Section Headers:

- **OWL Full**
- **OWL DL (Description logic)**
- **OWL Lite (even more restricted)**

Bottom Slide: Sample OWL

Code Content:

```
<?xml version="1.0"?>
<rdf:RDF
  xmlns:temporal="http://swrl.stanford.edu/ontologies/built-ins/3.3/temporal.owl#"
  xmlns:rdfs="http://www.w3.org/2000/01/rdf-schema#"
  (...many more namespaces...)
  xml:base="http://www.owl-ontologies.com/unnamed.owl">
<owl:Ontology rdf:about="">
<owl:imports rdf:resource="http://swrl.stanford.edu/ontologies/built-ins/3.3/query.owl"/>
<owl:imports rdf:resource="http://www.w3.org/2003/11/swrl"/>
  (...many more imports...)
</owl:Ontology>
<owl:Class rdf:ID="Rotational_displacement">
  <rdfs:label rdf:datatype="http://www.w3.org/2001/XMLSchema#string">
    Rotational displacement
  </rdfs:label>
  <rdfs:subClassOf>
    <owl:Class rdf:ID="Solid_displacement"/>
  </rdfs:subClassOf>
  <protege:subclassesDisjoint rdf:datatype="http://www.w3.org/2001/XMLSchema#boolean">
    true
  </protege:subclassesDisjoint>
  <owl:disjointWith>
    <owl:Class rdf:ID="Bending_displacement"/>
  </owl:disjointWith>
```

The video player interface includes standard controls (play/pause, volume, progress bar), a transcript button, help button, and video settings.

Ontology - Definition

- The study of being or existence
- Describes the basic categories and relationships of being to define entities and types of entities
- Study of the conceptions of reality
- A form of knowledge representation about the world or some part of it
- A set of definitions of formal vocabulary
- Includes an agreement to use a vocabulary in ways that are consistent (but not complete) with respect to the theory specified by the ontology

Video settings Transcript Help

Ontology – Why?

- To share common understanding of the structure of information among people or software agents
- To enable reuse of domain knowledge
- To make domain assumptions explicit
- To separate domain knowledge from the operational knowledge
- To analyze domain knowledge

Video settings Transcript Help

Ontology in e-Commerce

- Taxonomies provide :
 - Controlled shared vocabulary (search engines, authors, users, databases, programs/agents all speak same language)
 - Site Organization and Navigation Support
 - Expectation setting (left side of many web pages)
 - “Umbrella” Upper Level Structures (for extension)
 - Browsing support (tagged structures such as Yahoo!)
 - Sense disambiguation

Ontology - Evolution

- Simple ontologies can be built by non-experts
- Verity's Topic Editor, Collaborative Topic Builder, GFP, Chimaeras, Protégé, OIL-ED, etc.
- Ontologies can be semi-automatically generated
 - from crawls of site such as yahoo!, Amazon, excite, etc.
 - Semi-structured sites can provide starting points
- Ontologies are exploding(business pull instead of technology push)
 - most e-commerce sites are using - Amazon, Yahoo! Shopping
 - Expanding Business interest
 - Markup Languages growing XML, RDF, DAML, RuleML, xxML
 - “Real” ontologies are becoming more central to applications

Ontology - Representation

Ontologies may be expressed

Informally using natural language (e.g., in philosophy and sometimes biology)

Formally using a mathematical language, e.g., first-order logic (**or a fragment**)

Video settings Transcript Help

Formal Representation

- Every relation should be given/defined with a formal first order logical definition during its implementation

Video settings Transcript Help

Informal Representation

- OWL Model using Natural Language representation

Video settings Transcript Help

Ontology - Elements

- Instances (Individuals)
- Classes
- Attributes (properties)
- Relationships
- Hierarchical structure

Video settings Transcript Help

Instances

- Also called Individuals
- The most basic components of an ontology
- The actual, concrete objects (e.g., animals, bones, cars, etc.)
- An ontology does not require the inclusion of instances, but a main purpose of an ontology is to provide a means of classifying individuals, even if those instances are not explicitly part of the ontology



Examples of Instances

- Instance: *Sanjay*
 - an instance of *Student*
- Instance: *CSE Conference Hall*
 - an instance of a Conference Hall
- Instance: *Maruti Swift*
 - An instance of a hatchback car



Classes

- Also called Concepts
- Abstract groups, sets, or collections of objects.
- May contain individuals, other classes, or a combination of both

Attributes

- Also called Properties
- Features, characteristics, or parameters that objects can have and share
- Objects in the ontology are described by assigning attributes to them

Examples of Attributes

- Class: *Hatchback Car*
 - Attribute: *Number_of_doors*: Value: *four*
 - Attribute: *Engine Displacement* : Value: *1200cc*
- Instance: *Maruti Swift*
 - Attribute: *Name*: Value: *Maruti Swift*
 - Attribute: *Color*: Value: *Silver*
 - Attribute: *Registration_state*: Value: *Tamilnadu*



Relationships

- Ways that objects interact with one another
- An attribute whose value is another object in the ontology
- Most important type is subsumption: *is_a*
- Vertical Relation
- Horizontal Relation



Relationships – Broad Types

- Compositional Relation
 - Component of, element of, material of
- Spatial Relation
 - Has direct contact to, has non-direct contact to etc
- Role Relation
 - Operend, operator, resource, input, output
- Dependency Relation
 - Base of, depend on, aim at etc
- Influence Relation
 - Influencing, is opposing, is supporting
- Temporal Relation
 - After, before, co-occur, follows, proceeds
- General Relation
 - Alternative, criteria, represents, delivers etc



Examples of Relationships

➤ Instance: Maglia's car

- Relationship: *is_a*: Class: *Prius*
- Relationship: *is_same_color_as*: Instance: *Leopold's car*
- Relationship: *is_slower_than*: Instance: *Trueb's car*

➤ Class: *Prius*

- Relationship: *is_a*: Class: *Toyota*
- Relationship: *is_smaller_than*: Class: *4-Runner*



Hierarchical structure

- Inherent in classification system
- Defined by relationships among classes
- Attributes of superclasses are inherited by subclasses
- Most commonly used are *is_a* and *part_of*

is-a Classification

Class: *Vehicle*

Class: *Truck*

Class: *Car*

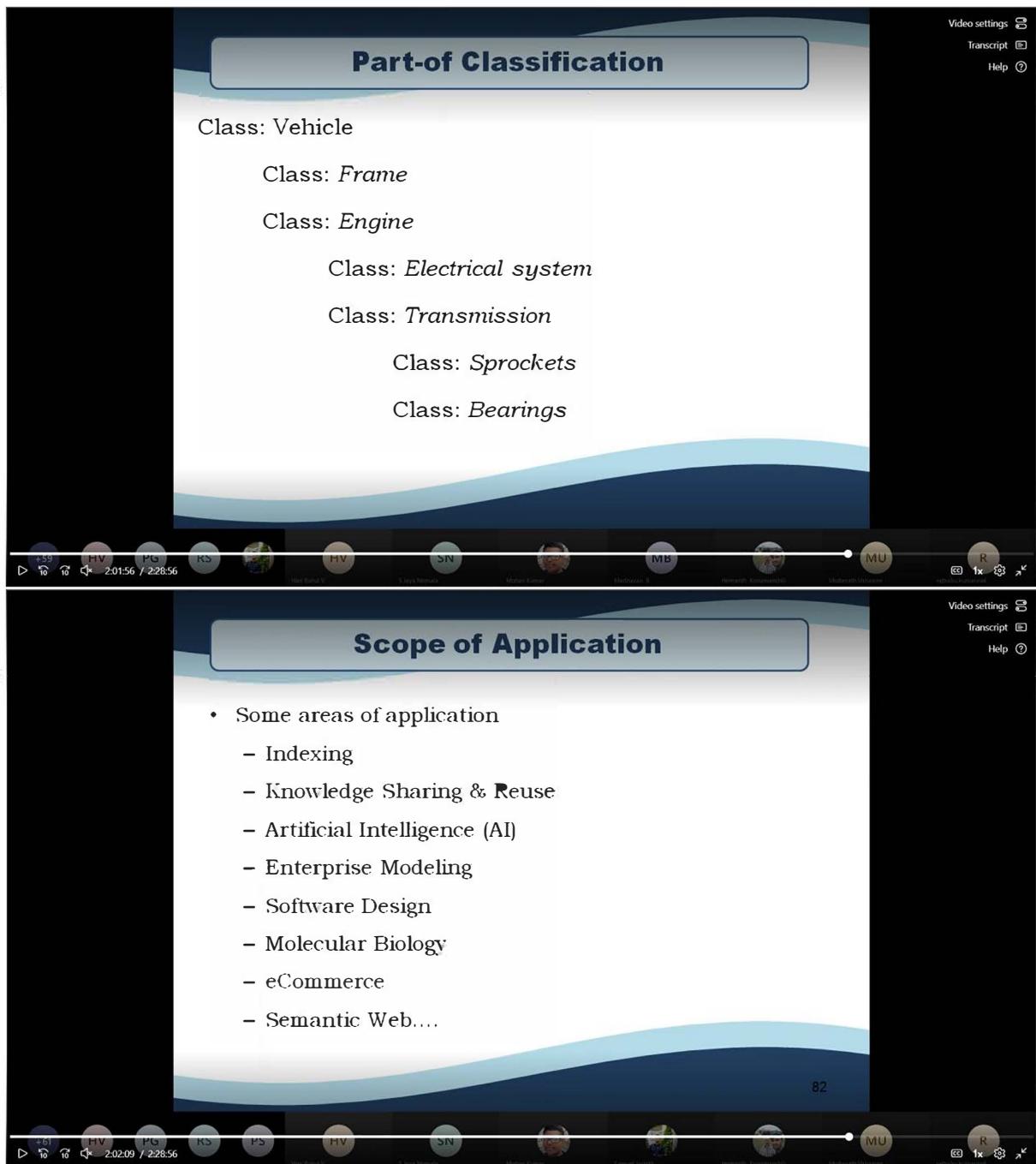
Class: *Honda*

Class: *Toyota*

Class: *Gasoline*

Class: *Hybrid*

Class: *Prius*



Part-of Classification

Class: Vehicle

Class: Frame

Class: Engine

Class: Electrical system

Class: Transmission

Class: Sprockets

Class: Bearings

Scope of Application

- Some areas of application
 - Indexing
 - Knowledge Sharing & Reuse
 - Artificial Intelligence (AI)
 - Enterprise Modeling
 - Software Design
 - Molecular Biology
 - eCommerce
 - Semantic Web....

82

Design Guidelines

- Need not contain all possible information about the domain.
- Need not generalize or specialize more than what your application needs.
- Need not contain all possible properties of and distinctions among classes in the hierarchy.

File Home Share View Manage application Tools

This PC > Local Disk (E) > softwares > Protege-5.5.0

Name	Date modified	Type	Size
app	3/14/2019 3:35 PM	File folder	
bin	3/14/2019 3:35 PM	File folder	
bundles	3/14/2019 3:35 PM	File folder	
conf	3/14/2019 3:35 PM	File folder	
jre	3/18/2019 6:50 AM	File folder	
plugins	3/14/2019 3:35 PM	File folder	
pizza.owl.xml	3/19/2022 9:36 AM	XML File	160 KB
Protege.exe	3/14/2019 2:53 PM	Application	130 KB
Protege4j.ini	3/14/2019 2:54 PM	Configuration sett...	1 KB
run.bat	3/14/2019 2:54 PM	Windows Batch File	1 KB

Top Panel (Annotations View):

Active ontology: pizza (http://www.co-ode.org/ontologies/pizza/2.0.0) : [E:\softwares\Protege-5.5.0\pizza.owl.xml]

Annotations tab selected.

Class hierarchy: AmericanHot

Individual: AmericanHot

Annotations for AmericanHot:

- rdfs:label [language: en] AmericanHot
- rdfs:label [language: pt] AmericanHot
- skos:prefLabel [language: en]
- Description: AmericanHot
- Equivalent To: None
- SubClass Of: None
- General class axioms: None

No Reasoner set. Select a reasoner from the Reasoner menu Show Inferences

Bottom Panel (OntoGraf View):

Active ontology: pizza (http://www.co-ode.org/ontologies/pizza/2.0.0) : [E:\softwares\Protege-5.5.0\pizza.owl.xml]

OntoGraf tab selected.

Class hierarchy: Mild

Individual: Mild

OntoGraf diagram showing relationships between toppings:

```

graph TD
    Mild --> ParmesanTopping
    Mild --> FetaTopping
    Mild --> RosemaryTopping
    Mild --> HerbSpiceTopping
    Mild --> CajunSpiceTopping
    Mild --> LeekTopping
    Mild --> SeafoodTopping
    Mild --> NutTopping
    Mild --> MeatTopping
    Mild --> CheeseTopping
    Mild --> SpicyTopping
    Mild --> ArtichokeTopping
    Mild --> MushroomTopping
    Mild --> GorgonzolaTopping
    Mild --> HotSauceTopping
    Mild --> Spiciness
    Mild --> Food
    Mild --> PizzaTopping
    ParmesanTopping --> Food
    FetaTopping --> Food
    RosemaryTopping --> Food
    HerbSpiceTopping --> Food
    CajunSpiceTopping --> Food
    LeekTopping --> Food
    SeafoodTopping --> Food
    NutTopping --> Food
    MeatTopping --> Food
    CheeseTopping --> Food
    SpicyTopping --> Food
    ArtichokeTopping --> Food
    MushroomTopping --> Food
    GorgonzolaTopping --> Food
    HotSauceTopping --> Food
    Spiciness --> Food
    Food --> PizzaTopping
  
```

No Reasoner set. Select a reasoner from the Reasoner menu Show Inferences

Ontology Design Approach

- There is no one correct way to model a domain.
- There are always viable alternatives
- Best solution depends on Applications and Extensions
- Iterative Process
- Concepts in Ontology close to objects (physical/logical) and relationships in the domain of interest
 - Objects are generally nouns
 - Relationships are generally verbs in a sentence

Video settings Transcript Help

Design Guidelines

- What is the domain that it will cover?
- For what are we going to use the ontology?
- For what types of questions should the information in the ontology provide answers? (i.e., *competency questions*)
- Who will use it?

Video settings Transcript Help

Design Guidelines

- Need not contain all possible information about the domain.
- Need not generalize or specialize more than what your application needs.
- Need not contain all possible properties of and distinctions among classes in the hierarchy.

Design Steps

- Define classes
- Arrange in Taxonomic hierarchy
- Sub-class/super-class model
- Define slots and Describe allowed values for these slots
- Fill values for slots for instances

Re-Use Existing Ontology

- If they exist, Sure..
- Problems in merging Ontologies ?
 - Format Conflicts
 - Same concept, different representation

Approaches

- Define Classes & Hierarchy
- Top-Down Approach
- Bottom-Up Approach
- Mixed
- Object Oriented Programming Analogy
- What do we get ?
 - Hierarchical arrangement of concepts
 - If class P is a super-class of class Q, every instance of Q is an instance of P.
 - Implication: class Q represents a “kind-of” P.

The image displays two identical screenshots of a video player interface, likely from a presentation or lecture video. The top slide is titled "Consistency Checks" and contains the following content:

➤ Ensure Class hierarchy correctness

- All siblings at a same level in the tree should have same level of generality
- Synonyms of classes are NOT different classes
- Check the “is-a”, “kind-of” relationships
- How many is too many & how few is too few?
- Multiple Inheritance

The bottom slide is also titled "Consistency Checks" and contains the following content:

- When do you introduce a new class ?
 - Subclass of a class usually
 - Have additional properties that the super-class does not have.
 - Have restrictions different from the super-class
 - A new class or a new property-value ?
 - Class “Black Bike” or simply property of class “Bike” that takes value “Black” ?
 - An Instance or a Class ?
 - Individual Instances are the most specific concepts represented in a Knowledge base.

Both slides include standard video player controls at the bottom, such as play/pause, volume, and a progress bar indicating the video is at 2:25:08 / 2:28:56. The interface also features a navigation bar with icons for search, refresh, and other functions, along with user profile pictures and names like "Samuel Hirsch", "Hari Balaji V", "S. Siva Narayana", "Mohan Kumar", "Ponnuthurai Karuppiah", "Madhavikuttyan", and "Video settings", "Transcript", and "Help".

Ontology Tools

Ontology-development becomes more accessible
Protégé

Developed at Stanford Medical Informatics
Is an extensible and customizable toolset for
constructing knowledge bases
developing applications that use these
knowledge bases

<http://protege.stanford.edu>

What is Protégé?

- An ontology editor
- A knowledge-acquisition tool
- A platform for knowledge-based applications

Protégé – Ontology Editor

What makes Protégé Different?

- Easy-to-use graphical interface
- Scalability
 - currently can handle up to 5 million concepts
- Plugin architecture
 - active international community of plugin developers
- It's a platform for other applications
 - Integration with Eclipse (Mayo Clinic)
 - A server and a client for (Semantic) Web Services
- Open source

