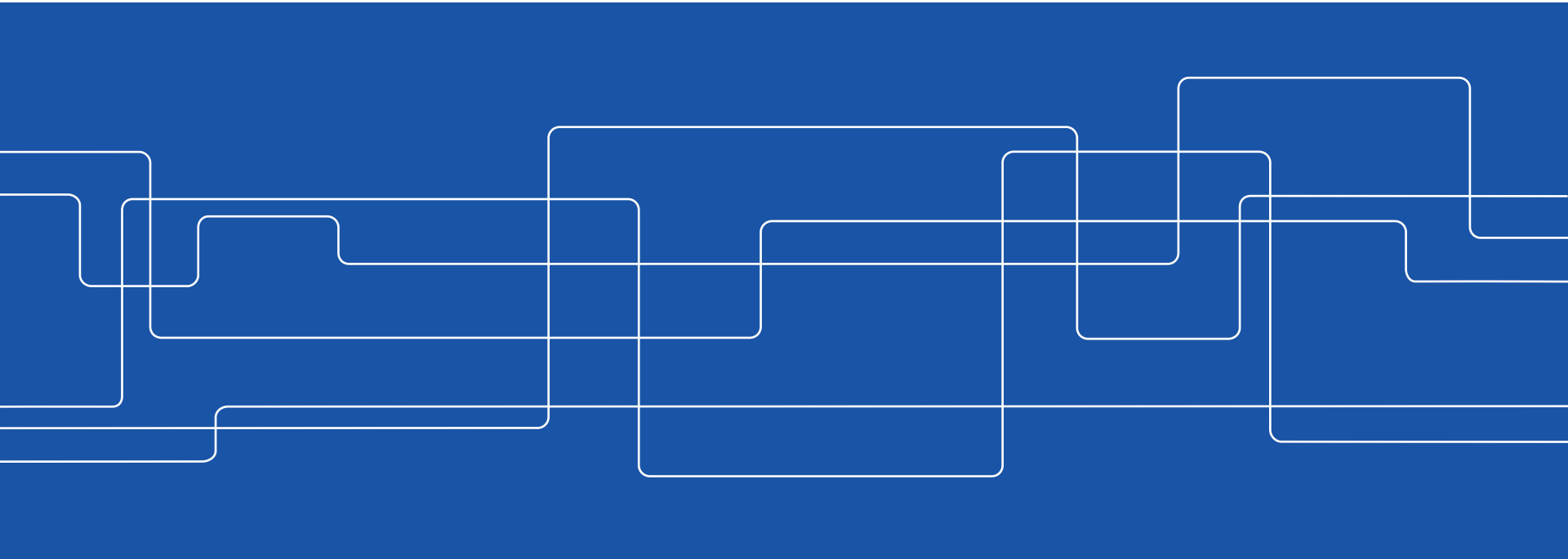




Lecture 12

Introduction to Machine Learning





On the Shoulder of Giants

Much of the material in this slide set is based upon:

"Automated Learning techniques in Power Systems"
by L. Wehenkel, Université Liege

"Probability based learning" Josephine Sullivan, KTH

"Entropy and Information Gain" by F.Aiulli,
University of Padova



Contents

Machine Learning vs Systems Theory

Some definitions

An illustrative example

Information content – entropy

Decision trees

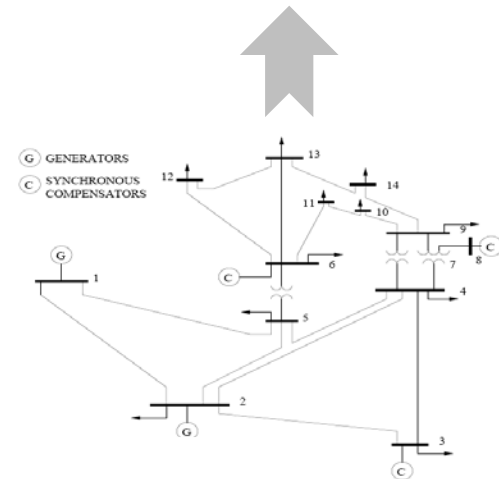
Power Systems Analysis - traditionally

Power system analysis, control and operation is dependent on models

Using the models, analytical and numerical analysis provides decision support for e.g.

- Security
- Stability
- Optimal power flow
- Contingency analysis
- Expansion planning
- Market clearing

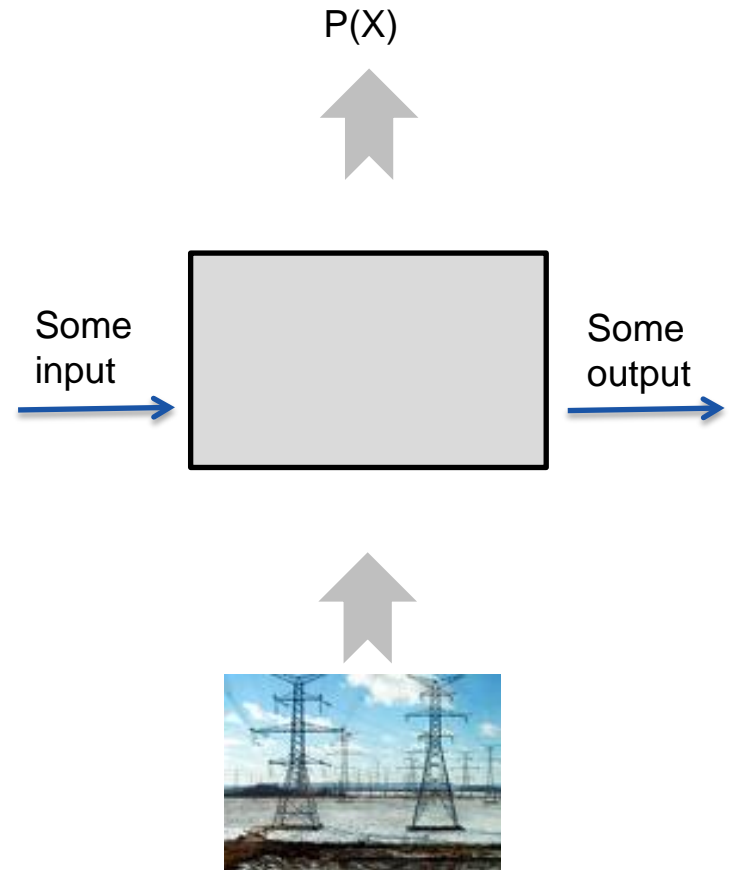
$$0 = -P_i + \sum_{k=1}^N |V_i||V_k|(G_{ik}\cos\theta_{ik} + B_{ik}\sin\theta_{ik})$$
$$0 = -Q_i + \sum_{k=1}^N |V_i||V_k|(G_{ik}\sin\theta_{ik} - B_{ik}\cos\theta_{ik})$$



Power Systems Analysis – An automated learning approach

Understanding states in the power system is established through observation of inputs and outputs without regard to the physical electrotechnical relations between the states.

Adding knowledge about the electrotechnical rules means adding heuristics to the learning.



Given a set of examples (the learning set (LS)) of associated input/output pairs, derive a general rule representing the underlying input/output relationship, which may be used to explain the observed pairs and/or predict output values for any new unseen input.



Why automated (or machine) learning

Historically it developed as computers became more powerful (1990s and on). Reached a peak late 1990s.

Provides complementary perspectives to systems approach:

- Computational efficiency.
Creating models of 10000+ bus systems including load and weather forecasts is challenging using system theoretical approach
- Interpretability
The correlation of variables can more directly be addressed enabling understanding of phenomena across subject fields
- Management of uncertainties
Errors in sensors and parameters can be disregarded



What is the "learning" in learning

Machine learning depends on having data about the system (power system) to be analysed/modeled. Such data should suitably be information about some state of the system.

Typical example of data is time-series of measurements
A trivial example is in weather forecasting, where time-series of temperature, wind speed and humidity from various locations across the country are obviously useful to create forecasts for the weather the next day.

Based on the input "*Learning*" data set, models are created and thereafter the models are fed with real data, and their ability to forecast results is evaluated.



Classes of methods for learning

In **Supervised** learning a set of input data and output data is provided, and with the help of these two datasets the model of the system is created.

For this introductory course, our focus is here. With a short look at unsupervised learning.

In **Unsupervised** learning, no ideal model is anticipated, but instead the analysis of the states is done in order to identify possible correlations between datapoints.

In **Reinforced** learning, the model in the system can be gradually refined through means of a utility function, that tells the system that a certain output is more suitable than another.



Classification vs Regression

Two forms of Supervised learning

Classification: The input data is number of switch operations a circuitbreaker has performed and the output is a notification whether the switch needs maintenance or not. "Boolean"

Regression: Given the wind speed in an incoming weather front, the output is the anticipated production in a set of wind turbines. "Floating point"



Supervised learning - a preview

In the scope of this course, we will be studying three forms of supervised learning.

- Decision Trees

Overview and practical work on exercise session.

- Artificial Neural Networks

Overview only, no practical work.

- Statistical methods – k-Nearest Neighbour

Overview and practical work on exercise session. Also included in Project Assignment

kNN algorithm can also be used for unsupervised clustering.



Steps of developing a learning model

Representation

Which are the inputs and outputs of relevance to the model – what is it that we want the model to be able to tell us?

Attribute selection

Reduction to the minimal set of useful attributes that will be used in creation of the model, making sure to remove non-relevant parameters.

Model selection

Create, or find, the suitable model of a certain structure (Decision tree, ANN, kNN, etc.) depending on the chosen set of input and output parameters.

Interpretation and validation of the model to make sure that the developed model provides output consistent with prior knowledge and physical limitations.

Once at this stage, the model can be used to predict new outputs based on alternative inputs – all within the constraints of its applicability as determined from the above process.



Contents

Machine Learning vs Systems Theory

Some definitions and notation

An illustrative example

Information content – entropy

Decision trees



Some further notation

The universe \mathbf{U} is the set of all possible objects \mathbf{O}

Objects have attributes denoted with $\mathbf{a(o)}$

With $\mathbf{a(X)}$, where \mathbf{X} is a subset of \mathbf{U} , we mean all possible values that $\mathbf{a(.)}$ can take in \mathbf{X}

For any subset \mathbf{V} of $\mathbf{a(U)}$ we denote by $\mathbf{a^{-1}(V)}$ the set of objects that exist in \mathbf{U} and whose attribute values are in \mathbf{V}



A database (yet a definition of the term)

A database is a subset of **objects** described by a certain number of **attributes** providing information about each object.

In a power system application, a database is for instance a number of simulated scenarios of varying load and generation profiles.

- Candidate Attributes are those that are used to develop (learn) the rule – also called Learning Set (LS)
- Selected Attributes are those that are used in the rule to validate it (also called Test Set)

Classification problems

Classification problems are such that we want to assign each object in the database to a specific class \mathbf{C} .

Each object in U can be assigned to one of mutually excluding classes \mathbf{C} , where $\mathbf{c}(\mathbf{o})$ denotes the class of \mathbf{o} .

$$\mathbf{C} \triangleq \{c_1, \dots, c_m\}$$

This can formally be specified as:

$$\mathbf{C} = \{\mathbf{C}_1, \dots, \mathbf{C}_m\} \quad : \quad \mathbf{C}_i \triangleq \{o \in U \mid c(o) = c_i\}.$$



Learning Sets and Test sets

In a classification problem the Learning set is

$$LS \triangleq \{(v^1, c^1), (v^2, c^2), \dots, (v^N, c^N)\},$$

where

$$v^k = (v_1^k, v_2^k, \dots, v_n^k)^T = a(o_k)$$

And

$$c^k = c(o_k)$$



Decision rules for classification

A Decision function, is a function $d(.)$ on a object such that

$$d(a_1(o), \dots, a_n(o)) \text{ or simply } d(o) : U \mapsto C.$$

A decision rule will create partitions of the universe U
(separate all the objects in different subsets D_i)

$$D_i \triangleq d^{-1}(c_i) = \{o \in U \mid d(o) = c_i\} \quad (i = 1, \dots, m).$$

The Hypothesis space is the space of all possible decision rules d that can operate on objects o in U



Regression problems

Regression problems takes as input objects O and provide as output a real number $y(U)$ (not a class).

$$r(a_1(o), \dots, a_n(o)) \text{ or simply } r(o) : U \mapsto y(U).$$



Probabilities – some recap and intro

A random variable X denotes a quantity that is uncertain, for example:

- Rolling a die, flipping a coin
- Measuring temperature

The probability distribution $\mathbf{P}(\mathbf{x})$ of a random variable describes that the variable will have different values when observed, and also which values we are more likely to get.

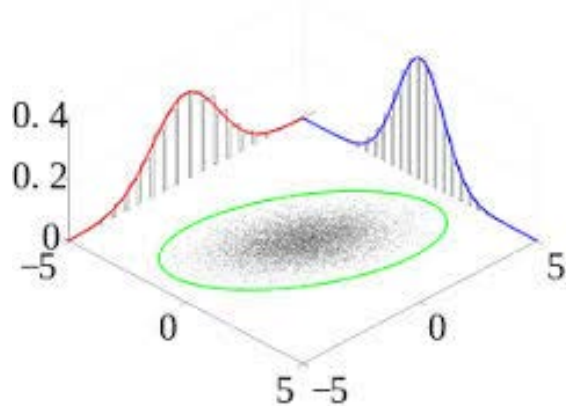
Probability Distribution of X

| Outcome | Probability |
|---------|-------------|
| 2 | 0.03 |
| 3 | 0.05 |
| 4 | 0.08 |
| 5 | 0.11 |
| 6 | 0.14 |
| 7 | 0.17 |
| 8 | 0.14 |
| 9 | 0.11 |
| 10 | 0.08 |
| 11 | 0.05 |
| 12 | 0.03 |

Joint Probability Distributions

Assume that you have two random variables X and y

If you observe multiple paired instances of X and y , this is described in the joint probability density function $P(x,y)$.



By summing (integrating) of one of the variables (e.g. y) you get the PDF for the other (e.g. X)

$$P(x) = \int_y P(x, y) dy$$



Conditional Probability

For a joint probability distribution $P(x,y)$ the conditional probability denoted $P(X|Y)$ is the probability of getting the value X , when we know $y=y^*$

$$P(x | y = y^*) = \frac{P(x, y = y^*)}{P(y = y^*)} = \frac{P(x, y = y^*)}{\int_x P(x, y = y^*) dx}$$



Classification probability

To simplify notation, we say that

$$P^i(\mathbf{X}) \triangleq P(\mathbf{X} | C_i).$$

Denotes the probability P_i of finding \mathbf{X} in a given class C_i

Remember, Classes are usually few (1-3), like safe – unsafe.



Contents

Machine Learning vs Systems Theory

Some definitions

An illustrative example

Information content – entropy

Decision Trees



OMIB – further information

In the OMIB system the following parameters influence security

- Amount of active and reactive power of the generator (
- Amount of load nearby the generator (PI)
- Voltage magnitudes at the load bus and at the infinite bus
Short-circuit reactance X_{inf} , representing the effect of variable topology in the large system represented by the infinite bus.

In the example, Voltages at generator and Infinite bus are assumed similar and constant for simplicity

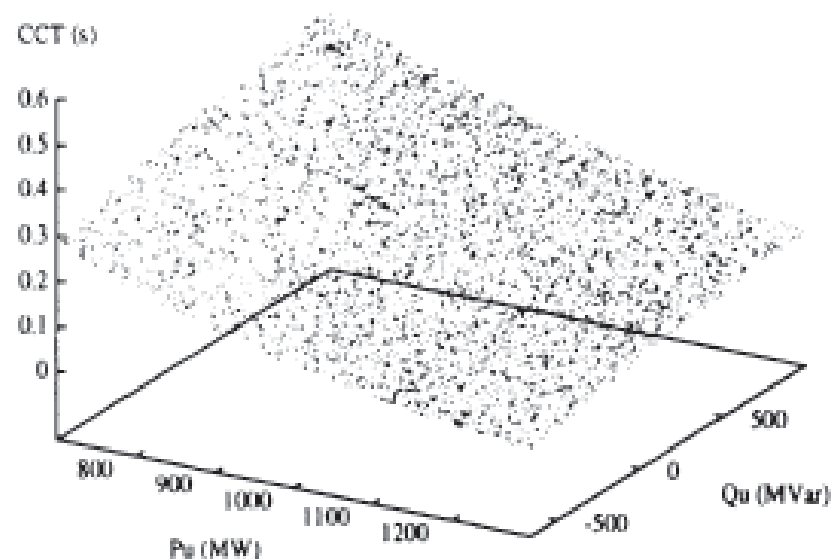
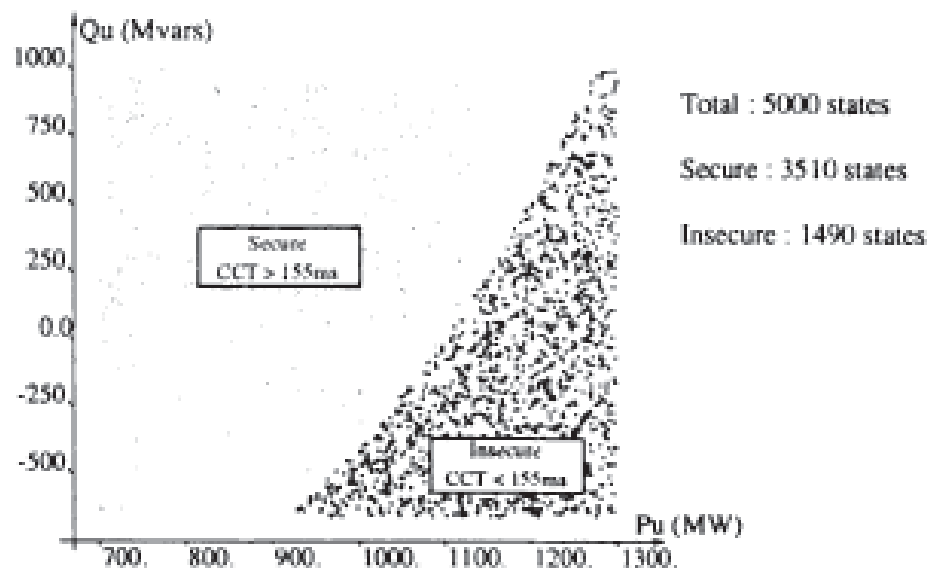
Our database of objects with attributes

In a simulator, we randomly sample values for P_u and Q_u creating a database with 5000 samples (objects) and for each object we have a set of attributes (P_u , Q_u , V_1 , P_1 , V_{inf} , X_{inf} , CCT) as per below.

Table 1.1. Sample of OMIB operating states

| State Nb | P_u (MW) | Q_u (MVar) | V_1 (p.u.) | P_1 (MW) | V_{inf} (p.u.) | X_{inf} (Ω) | CCT (s) |
|----------|------------|--------------|--------------|------------|------------------|------------------------|---------|
| 1 | 876.0 | -193.7 | 1.05 | -100 | 1.05 | 60 | 0.236 |
| 2 | 1110.9 | -423.2 | 1.05 | -100 | 1.05 | 60 | 0.112 |
| 3 | 980.1 | 79.7 | 1.05 | -100 | 1.05 | 60 | 0.210 |
| 4 | 974.1 | 217.1 | 1.05 | -100 | 1.05 | 60 | 0.224 |
| 5 | 927.2 | -618.5 | 1.05 | -100 | 1.05 | 60 | 0.158 |
| ... | ... | ... | ... | ... | ... | ... | ... |
| 2276 | 1090.4 | -31.3 | 1.05 | -100 | 1.05 | 60 | 0.157 |
| ... | ... | ... | ... | ... | ... | ... | ... |
| 4984 | 1090.2 | -20.0 | 1.05 | -100 | 1.05 | 60 | 0.158 |
| ... | ... | ... | ... | ... | ... | ... | ... |

Plot of database content





Our task

By using the sampled datapoints (5000 total) our task is now to determine a model that maps the input parameters to a classification of the system being Secure or Unsecure.

To start with, we select from the samples

1. A Learning set (LS) which are the objects (with attributes) that we will use to create the model
2. A Test set (TS) which are the objects (with attributes) we will use to validate the model



Contents

Machine Learning vs Systems Theory

Some definitions

An illustrative example

Information content – entropy

Decision Trees



How to measure information content

Entropy **H** is a measure of *Unpredictability*.

Defined as:

$$H = - \sum p_i \log p_i$$

Where

p_i is the probability of event i



Some examples

1. Flipping a coin
2. Rolling a 6 sided die
3. Rolling a loaded 6 sided die



Entropy in a Dataset

$$H_C(\mathbf{X}) \triangleq - \sum_{i=1, \dots, m} P(C_i | \mathbf{X}) \log P(C_i | \mathbf{X}),$$

The classification Entropy, is the entropy related to the probability of a value \mathbf{X} belonging to a class C_i

Or simply put, how difficult is it to guess which partition of U , i.e. Class C_i that an object o belongs to.

An example of classification entropy

| Color | Size | Shape | Eadible? |
|--------|-------|-----------|----------|
| Yellow | Small | Round | Yes |
| Yellow | Small | Round | No |
| Green | Small | Irregular | Yes |
| Green | Large | Irregular | No |
| Yellow | Large | Round | Yes |
| Yellow | Small | Round | Yes |
| Yellow | Small | Round | Yes |
| Yellow | Small | Round | Yes |
| Green | Small | Round | No |
| Yellow | Large | Round | No |
| Yellow | Large | Round | Yes |
| Yellow | Large | Round | No |
| Yellow | Large | Round | No |
| Yellow | Large | Round | No |
| Yellow | Small | Irregular | Yes |
| Yellow | Large | Irregular | Yes |



Entropy example

Entropy for the example data set is calculated as:

$$I(all_data) = -\left[\left(\frac{9}{16}\right)\log_2\left(\frac{9}{16}\right) + \left(\frac{7}{16}\right)\log_2\left(\frac{7}{16}\right)\right]$$

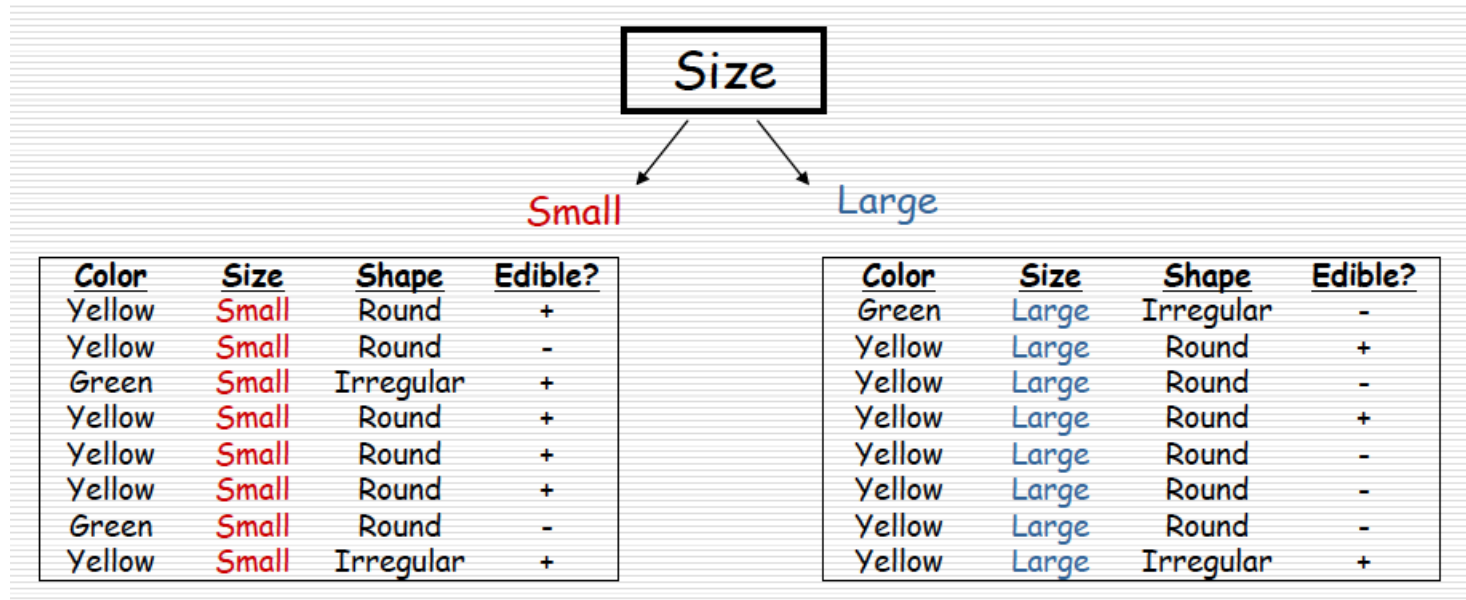
Giving: 0,9836

Is this reasonable?

Information Gain

The reduction in Entropy achieved by partitioning the dataset differently.

Lets separate for instance per the attribute Size.





Information Gain calculation

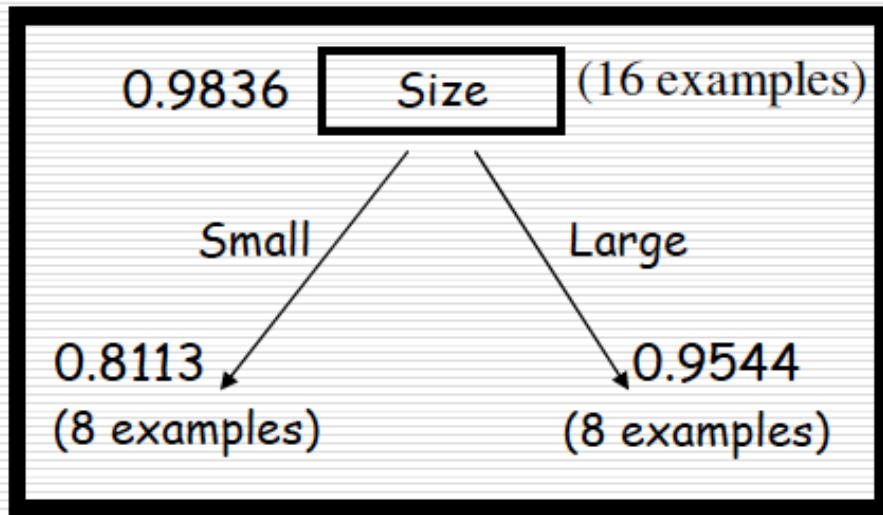
The two partitions has their own entropy value.

We can calculate for each possible attribute its expected entropy. This is the degree to which the entropy would change if partitioned based on this attribute.

To determine resulting entropy, you add the entropies of the two partitions, weighted by the proportion of examples from the parent node that ended up in that partition.

$$G(S, A) = I(S) - \sum_{v \in \text{Values}(A)} \frac{|S_v|}{|S|} I(S_v)$$

Hence for the example



Entropy of left child is 0.8113
 $I(\text{size}=\text{small}) = 0.8113$

Entropy of right child is 0.9544
 $I(\text{size}=\text{large}) = 0.9544$

$$I(S_{\text{size}}) = (8/16) * .8113 + (8/16) * .9544 = .8828$$



Contents

Machine Learning vs Systems Theory

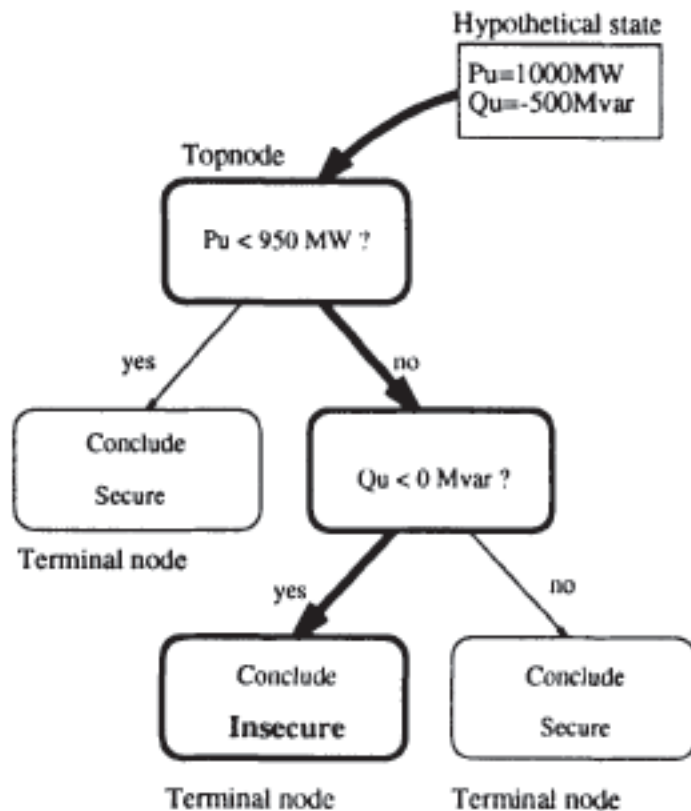
Some definitions

An illustrative example

Information content – entropy

Decision Trees

Back to our Power System example



Perhaps we can partition our dataset according to some attribute?

Lets try $P_u < 950\text{MW}$

Equivalent If-Then rules :

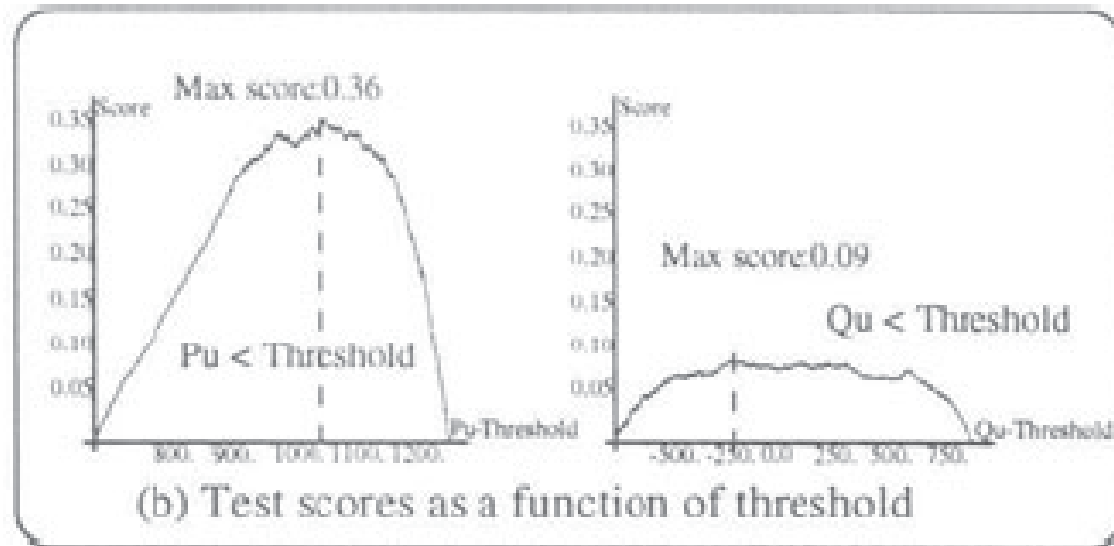
Rule 1 : If $(P_u < 950\text{MW})$ then Conclude Secure

Rule 2 : If $(P_u > 950\text{MW})$ and $(Q_u < 0\text{Mvar})$ then Conclude Insecure

Rule 3 : If $(P_u > 950\text{MW})$ and $(Q_u > 0\text{Mvar})$ then Conclude Secure

Finding best partition.

Starting with the candidate attributes (P_u and Q_u) in our case
We check which of the values for P_u and Q_u that create the most valuable partition in terms of information gain.



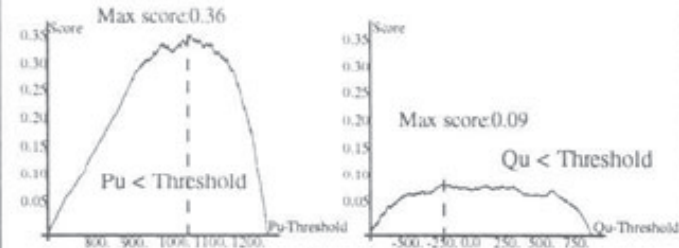
$P_u > 1096,2$ MW is the best partition

Gradual expansion of the Decision Tree

Tree at step 0

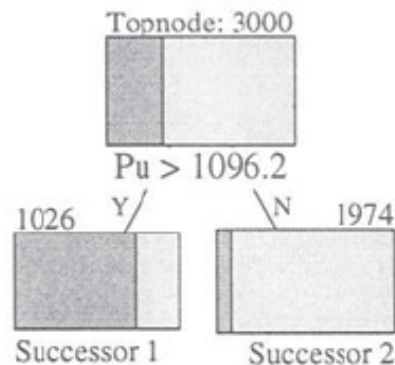


(a) Top node of the tree



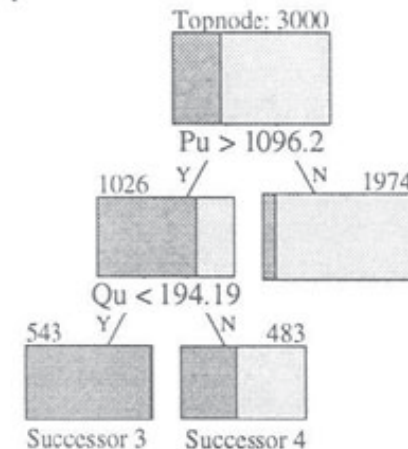
(b) Test scores as a function of threshold

Tree at step 1



(c) Tree after the topnode was developed

Tree at step 2



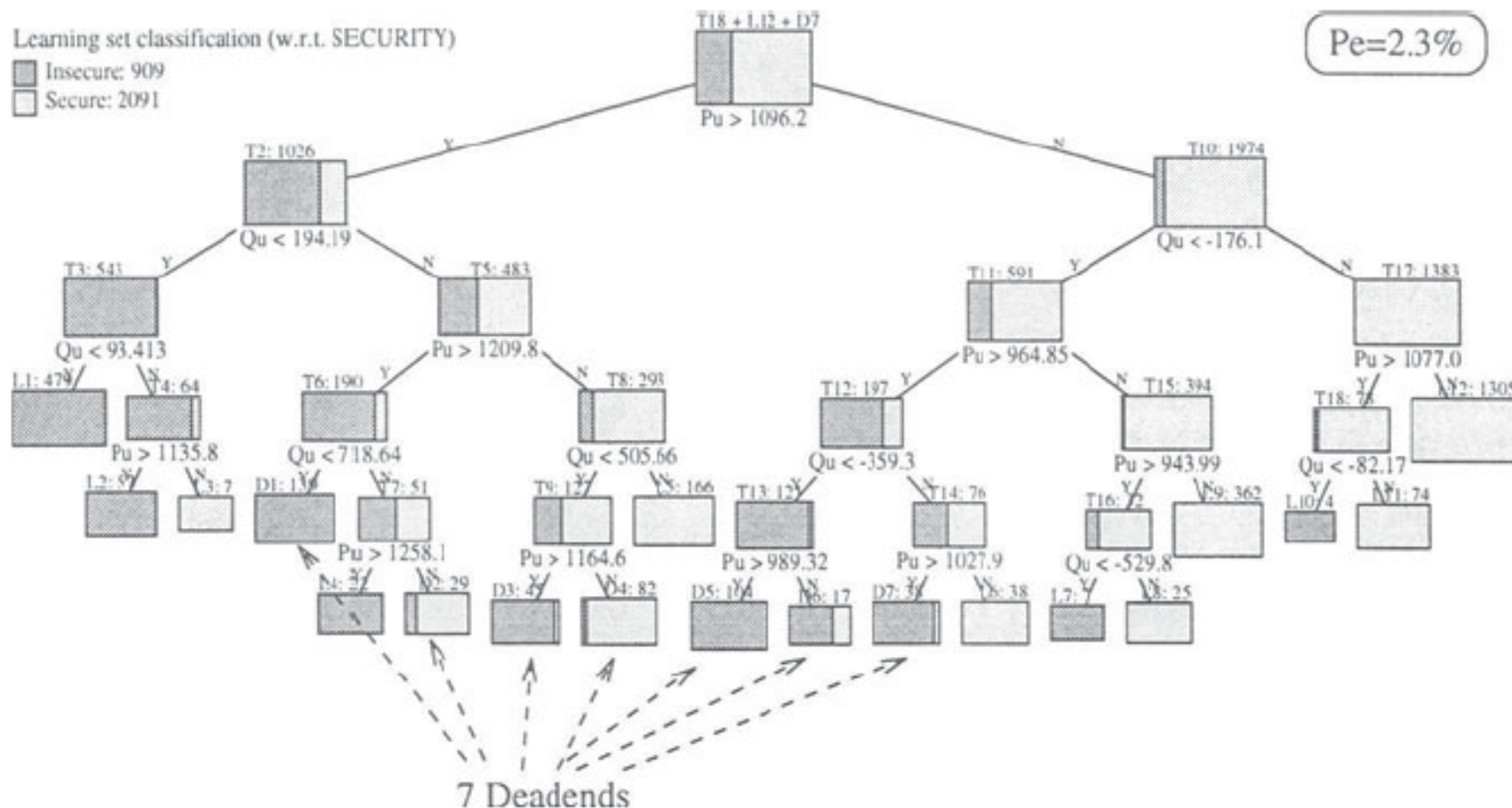
(d) Tree after the first successor was developed

Complete Decision Tree

Learning set classification (w.r.t. SECURITY)

■ Insecure: 909
□ Secure: 2091

Pe=2.3%





How to stop?

The splitting of data sets continues until either:

A perfect partition is reached – i.e. One which perfectly explains the content of the class – a *leaf*

One where no information is gained no matter how the data set is split. – a *deadend*.



Validation of the Decision Tree

By using the Test Set (2000 samples) we can calidate the Decision tree.

By testing for each Object in the Test Set, we determine if the Decision tree provides the right answer for the Object.

In this particular example, the probability o error can be determined to 2,3. I.e. Of the 2000 samples 46 were classified to the wrong class.



Contents

Machine Learning vs Systems Theory

Some definitions

An illustrative example

Information content – entropy

Decision Trees