



Data Science for Smart Cities

CE88

Prof: Alexei Pozdnukhov

Data Science for Smart Cities



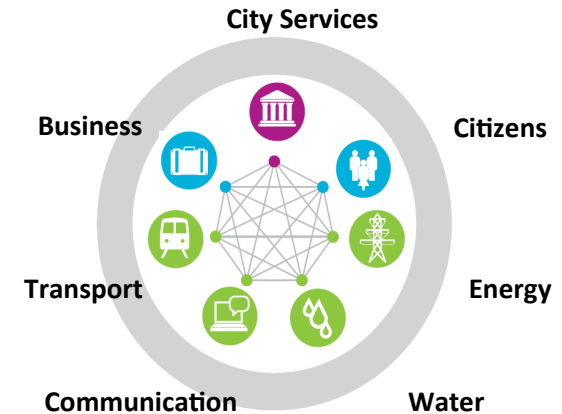
Weeks 1-2. Introduction and motivation: cities as complex systems.

Lecture 1. Introduction to urban systems. Inter-dependent infrastructures with human in the loop.

Lecture 2. Modeling principles. Causality.

Lecture 3. Spatio-temporal nature of urban data.

Lecture 4. Data flows in cities



Weeks 3-5. Urban data collection, handling and processing.

Lecture 5. Data acquisition: measurement and crowd-sensing.

Lecture 6. Community surveys, population census, APIs.

Lecture 7. Demand and supply data exploration.

Lecture 8. Impact of urbanization



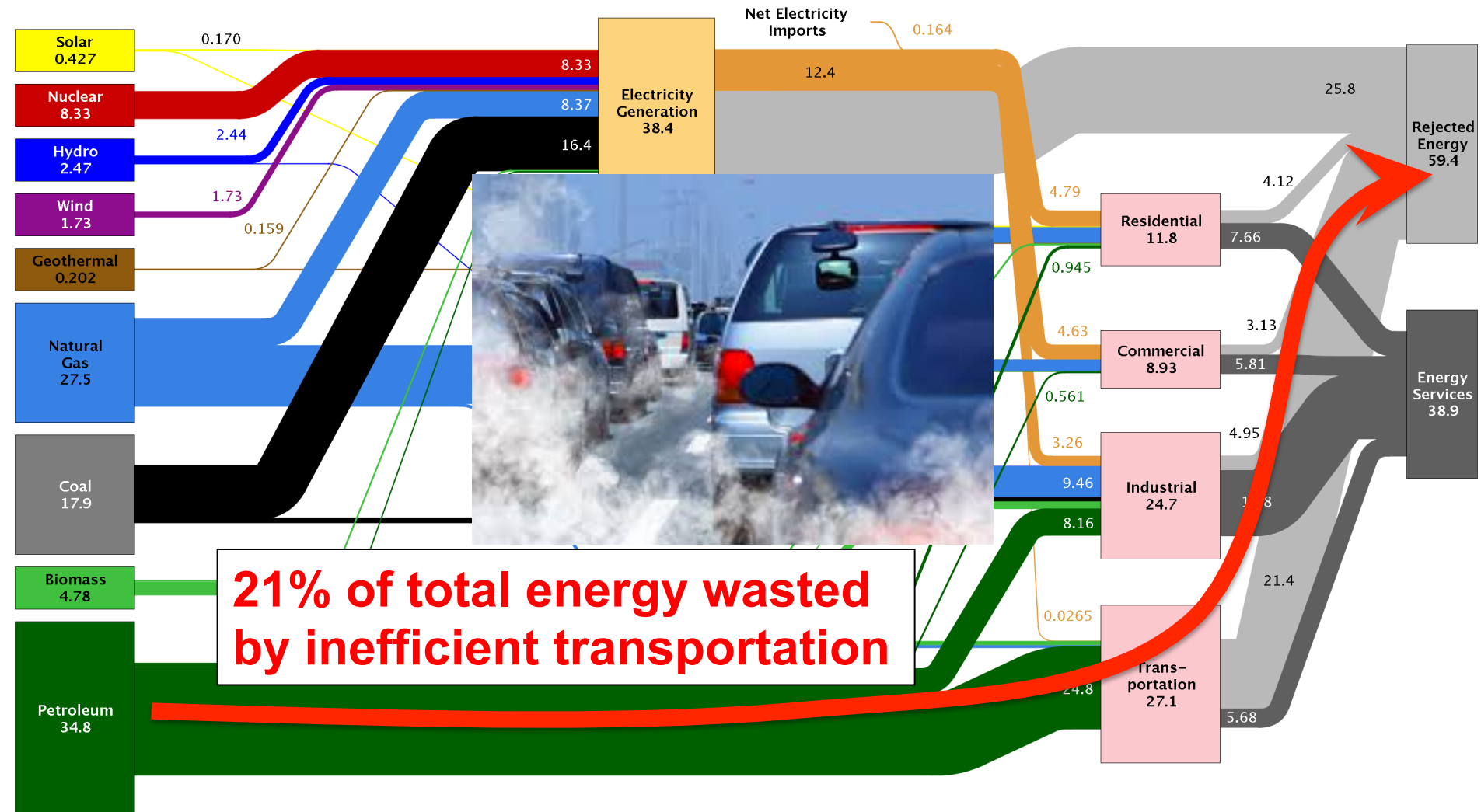
Data exploration, visualization and modeling

Impact of urbanization: energy use



Estimated U.S. Energy Use in 2014: ~98.3 Quads

Lawrence Livermore
National Laboratory



Source: LLNL 2015. Data is based on DOE/EIA-0035(2015-03), March, 2014. If this information or a reproduction of it is used, credit must be given to the Lawrence Livermore National Laboratory and the Department of Energy, under whose auspices the work was performed. Distributed electricity represents only retail electricity sales and does not include self-generation. EIA reports consumption of renewable resources (i.e., hydro, wind, geothermal and solar) for electricity in BTU-equivalent values by assuming a typical fossil fuel plant "heat rate." The efficiency of electricity production is calculated as the total retail electricity delivered divided by the primary energy input into electricity generation. End use efficiency is estimated as 65% for the residential and commercial sectors 80% for the industrial sector, and 21% for the transportation sector. Totals may not equal sum of components due to independent rounding. LLNL-MI-410527



Sustainable TRANSPORTATION

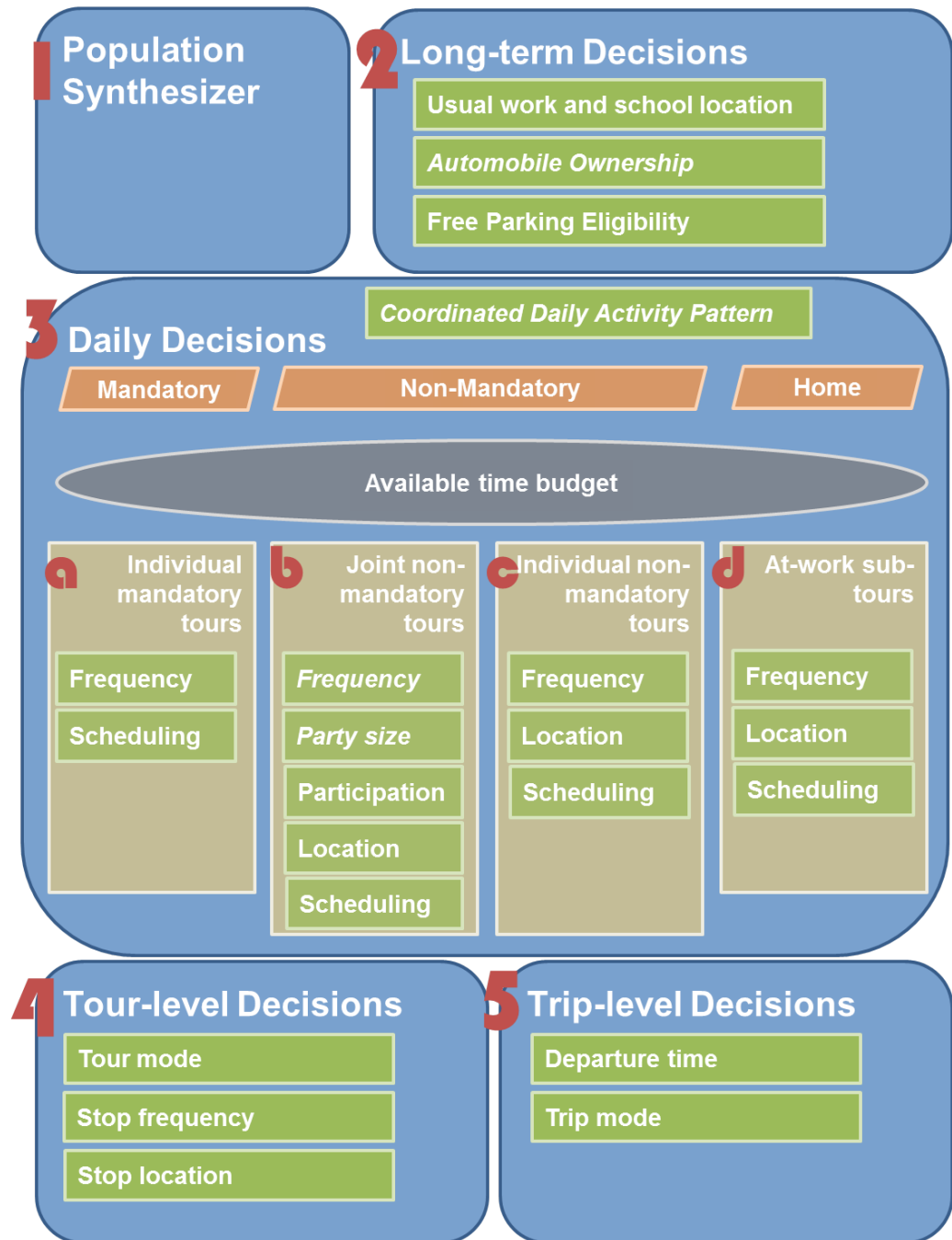
U.S. DEPARTMENT OF
ENERGY

Energy Efficiency &
Renewable Energy

Focus Area	Future New Technologies/Models/Knowledge	Performance Metrics
Decision Science	<ul style="list-style-type: none"> New <u>knowledge and applications</u> of socio-behavioral science to collect and analyze real-world data on transportation decision making, EV and AFV market drivers and barriers, as well as new mobility options. 	For all new transportation as a system studies and models:
Connectivity and Automation	<ul style="list-style-type: none"> An <u>increased understanding</u> of the impact of connected and automated vehicles and their implications on transportation and vehicle technologies, such as electrification and overall mobility. 	<ul style="list-style-type: none"> Survey existing resources Complete gap analysis Propose synthesis and expansion of state-of-the-art analysis
Multi-Modal	<ul style="list-style-type: none"> Dynamic passenger/freight modal and energy-intensity <u>modeling</u> with explicit consideration of consumer/market preferences and energy implications. 	<ul style="list-style-type: none"> Sub-topical deep dives that can inform future technology deployments
Urban Science	<ul style="list-style-type: none"> Integrated <u>city-scale models</u> that explicitly consider energy impacts of urbanization by collecting real-world data and collaborating with local governments 	<ul style="list-style-type: none"> Deliver new cutting-edge transportation system models
Vehicles and Infrastructure	<ul style="list-style-type: none"> Integrated vehicle-fuel <u>models</u> to explore value propositions (consumer and provider), business models and opportunities for increased sustainable transportation deployment. 	<ul style="list-style-type: none"> Apply priority scenario illustrative examples to inform discrete conclusions

Example model structure

MTC Travel Demand model



Today

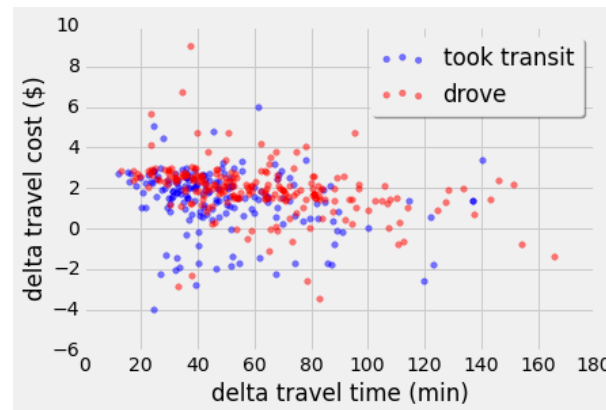


Reminder on the goals of Exploratory Data Analysis:

- suggest a hypothesis about the causes of the phenomena
- state the problem of data analysis
- assess the validity of the assumptions
- select an algorithm to approach the stated problem

Taxonomy of algorithms

Approaches to model specification

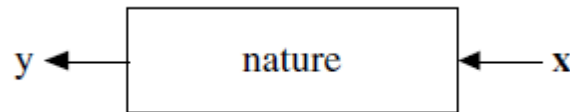


Mini Lab 7: specify and apply a simple predictor algorithm

“Statistics starts with data”

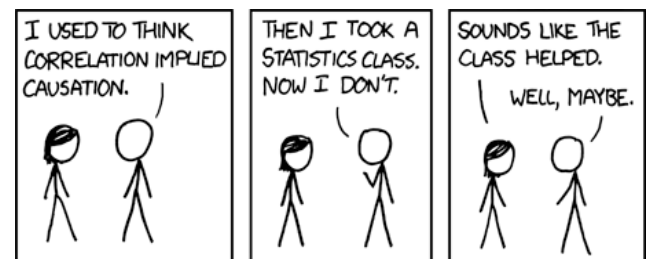


Think of the data as being generated by a black box in which input variables x (predictor or inputs, independent or explanatory variables) go in one side, and on the other side the response variables y come out. Inside the black box, nature functions to associate the predictor variables with the response variables, so the picture is like this:

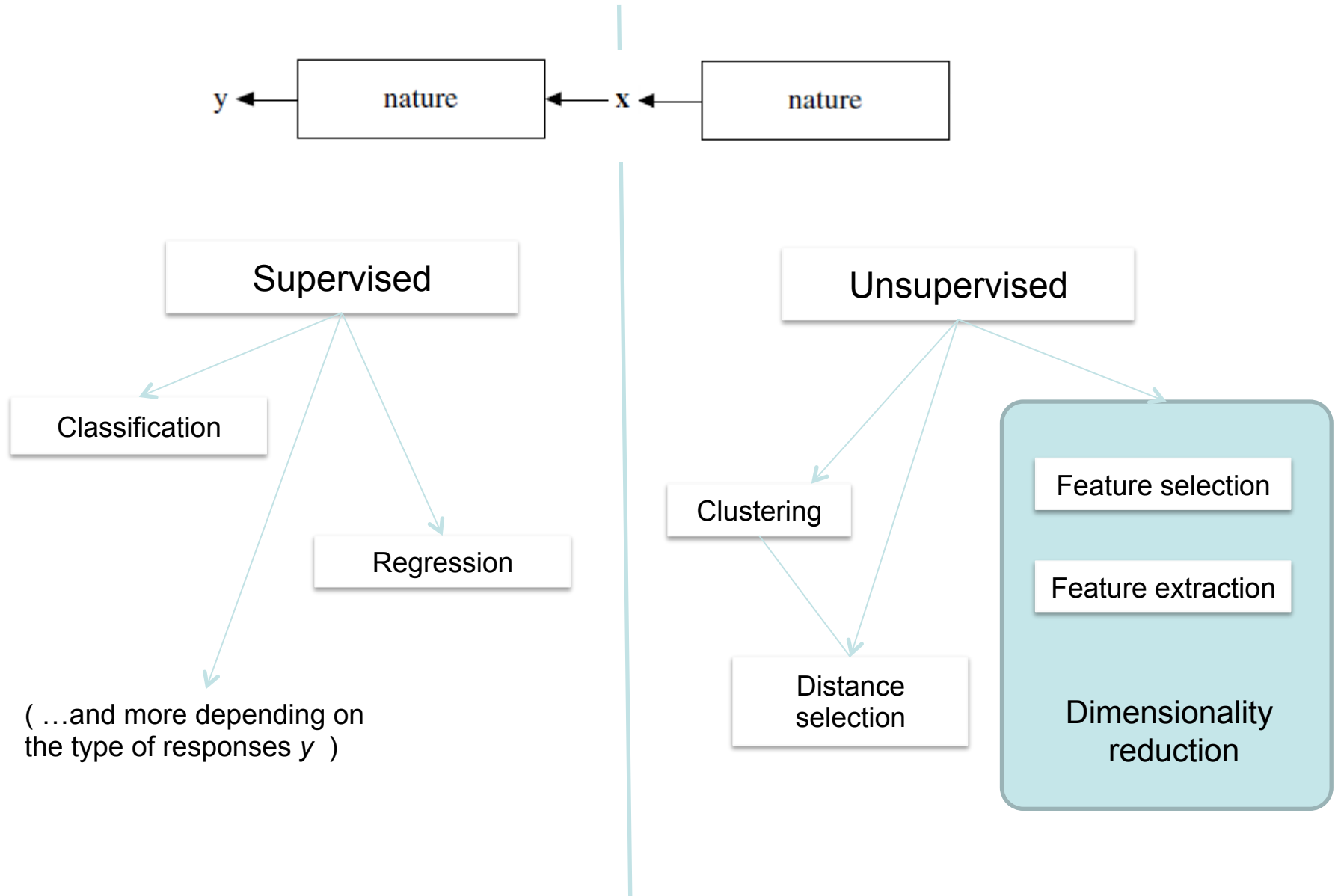


There are at least **two** goals in analyzing the data:

- To be able to predict what the responses are going to be to future input variables;
- To extract some information about *how* nature is associating the response variables to the input variables.



Taxonomy of problems and algorithms



Supervised modeling



Nowhere it is set in stone how nature black box should be modelled when the only thing you got is empirical data. It is **our** choice to specify an appropriate model. How the translation from subject-matter problem to statistical model is done is often the most critical part of an analysis.

Best thing one can do is to specify a model in a way that achieves both goals we just stated.

One can assume a stochastic data model for the inside of the black box, i.e. assume that data are generated by independent draws from:

response variables = f (predictor variables, random noise, parameters).

Then the values of the parameters can be estimated from the **sample** and the model then used for information (inference) and/or prediction for the **population**.

Statistical modelling



If one starts with “assume that the data are generated by the following model: ... “, this implies that by imagination and by looking at the data, one can invent a reasonably good parametric class of models for a complex mechanism devised by nature.

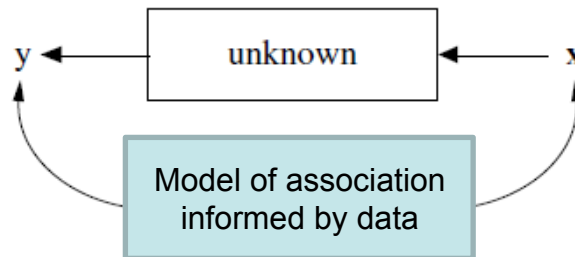
The conclusions made with such approach could be true about the model's mechanism, but not necessarily about the nature's mechanism. If the model is a poor emulation of nature, the conclusions maybe wrong.

This approach is most solid when indeed the nature of the observed association between x and y can be identified, for example, when a known physical process is observed under noise (noise in the process parameters or as additive random factors to the observable).

Statistical modelling



Most often, the inside of the nature box is complex and unknown.



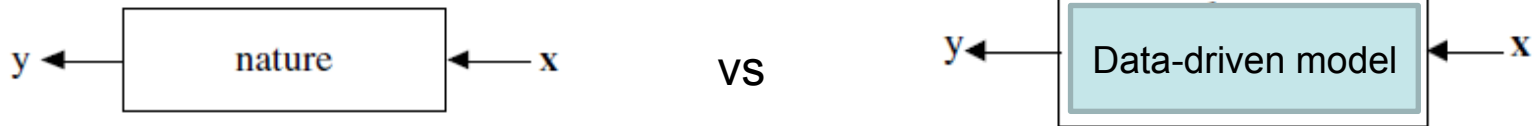
The approach to find an algorithm that operates on x to predict the responses y must be well thought through. Particularly, to make statistical inference about the nature of the x to y association, one has to:

- control for confounding factors
- use randomized controlled trials
- validate and test models on previously unobserved (out-of-sample) data

Analysis of models



When the model is built, we are comparing:



The great advantage of a model is if it produces a simple and understandable picture of the relationship between the input variables and responses.

Examples are:

- decision rule based on thresholds

```
if (travel distance X > 10 miles):  
    person Y travels by car  
else:  
    person Y travels by bike
```

- model based on empirical evidence

```
if (most similar person to person Y  
    travels by car):  
    person Y travels by car
```

$$y = b_0 + \sum_{m=1}^M b_m x_m + \varepsilon$$

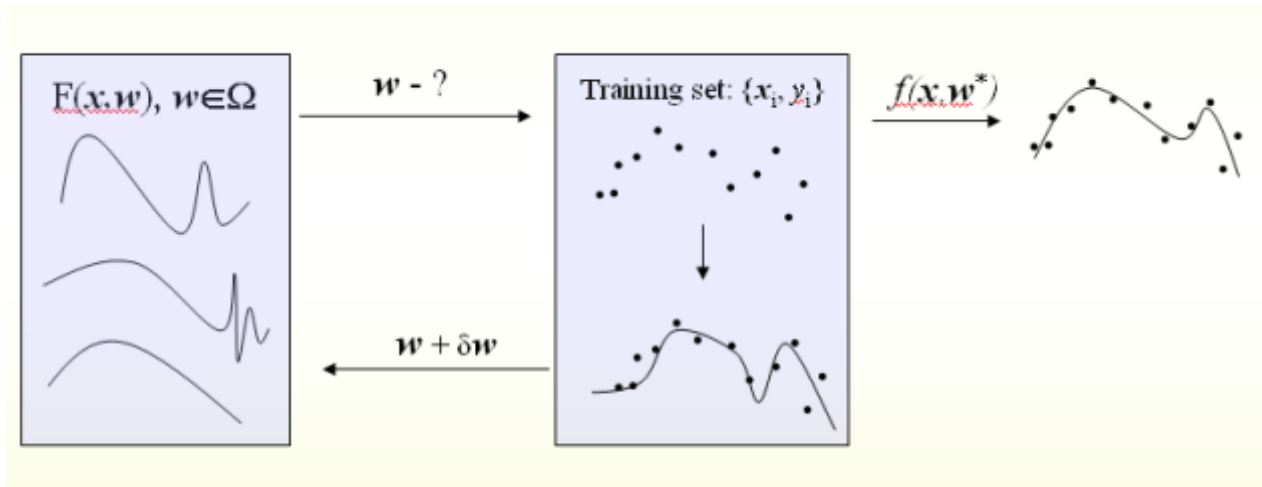
- model linear in parameters

Note: such models may seem oversimplified, and not fitting the data perfectly. The choice between accuracy and interpretability is a difficult one. In a choice between accuracy and interpretability, practitioners often go for interpretability. However, a model does not have to be simple to provide reliable information about the relation between predictor and response variables.

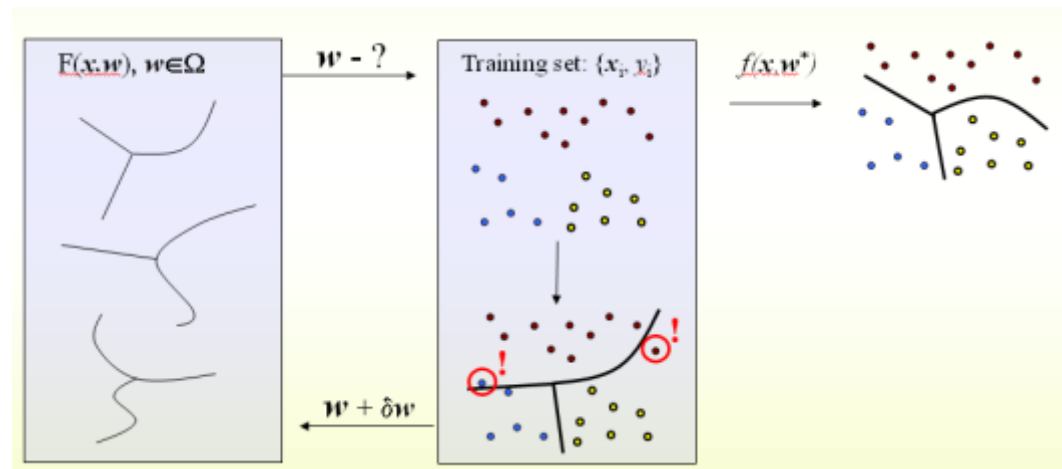
Algorithmic framework



Regression



Classification



Geometry of the input space



A simple prediction method is a ‘nearest neighbor’:

Given an unseen situation \mathbf{x}_0 , provide a likely y_0 that can be associated to it:

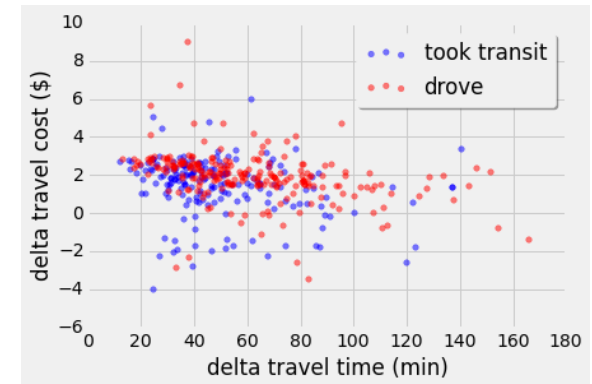
- find an example \mathbf{x} in the available dataset that is most similar to \mathbf{x}_0
- assign an observed response variable y as your best guess for y_0

What do we mean by ‘nearest’ or ‘most similar’?

It has to be an informed decision made by us!

Ideas:

- use Euclidean distances,
- normalize the input variables,
- or
- apply scaling to input variables justified by domain knowledge.



Take away ideas



- There are multiple tasks one can think of when dealing with data, so keep in mind a taxonomy of algorithms and methods
- Methods deal with a sample and are used to get information (infer) and/or predict for the entire population
- Methods can be conceptually represented as an algorithmic model of nature



- Algorithms can have parameters that one can 'tune' to achieve better performance

In the Mini-lab:

- 1) we will learn a simple programming concept to implement algorithms
- 2) we will solve a simple prediction problem