

# Recovering unknown governing equations from measurement data

---

Rui Carvalho<sup>1,2,3</sup>

September 29, 2019

<sup>1</sup>Department of Engineering, Durham University, Lower Mountjoy, South Road, Durham, DH1 3LE, UK

<sup>2</sup>Durham Energy Institute, Durham University, South Road, Durham, DH1 3LE, UK

<sup>3</sup>Institute for Data Science, Durham University, South Road, Durham, DH1 3LE, UK

# Linear regression in a nutshell

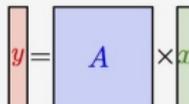
---

# Linear systems in a nutshell

Solving  $y \approx Ax \in \mathbb{R}^m$       $A \in \mathbb{R}^{m \times n}$

Determined ( $m = n$ ):

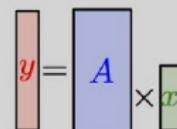
$$x = A^{-1}y$$



Over-determined ( $m > n$ ):

$$\min_x \|Ax - y\|^2$$

$$x = (A^T A)^{-1} A^T y \stackrel{\text{def.}}{=} A^+ y$$

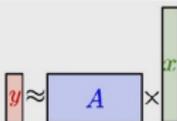


least squares

Under-determined ( $m < n$ ):

$$\min_x \{\|x\| ; Ax = y\}$$

$$x = A^T (AA^T)^{-1} y \stackrel{\text{def.}}{=} A^+ y$$



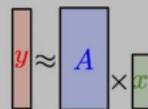
least norm problem

A ill-posed and/or noise:

$$\min_x \|Ax - y\|^2 + \lambda \|x\|^2$$

$$x = (A^T A + \lambda \text{Id}_n)^{-1} A^T y \xrightarrow{\lambda \rightarrow 0} A^+ y$$

$$= A^T (AA^T + \lambda \text{Id}_m)^{-1} y \quad (\text{Woodbury identity})$$



modified from Gabriel Peyré

## Least squares problem

Least squares problem: choose  $x$  to minimise  $f(x) = \|Ax - b\|_2^2$  where  $A \in \mathbf{R}^{m \times n}$  with  $m \geq n$ , and  $b \in \mathbf{R}^m$  are problem data.

- $m \times n$  matrix  $A$  is tall, so  $Ax = b$  is over-determined.
- For most choices of  $b$ , there is no  $x$  that satisfies  $Ax = b$ .
- *Residual*:  $r = Ax - b$ .
- Idea: make residual as small as possible, if not 0.
- Assume that the columns of  $A$  are independent (the Gram matrix  $A^T A$  is invertible), the least-squares approximation problem has the unique solution:

$$x = (A^T A)^{-1} A^T b. \quad (1)$$

- Compare with the solution of the square invertible system  $Ax = b$ :

$$x = A^{-1} b \quad (2)$$

## Regularised approximation

how well the data agrees with the model  $Ax=b$

$$\text{minimise } (\|Ax - b\|, \|x\|) \quad (3)$$

how large are your model parameters

- $A \in \mathbf{R}^{m \times n}$  is a matrix of  $n$  predictors;
- $x \in \mathbf{R}^n$  are the parameters;
- $b \in \mathbf{R}^m$  is a vector of responses.

### Idea:

- We want a good fit of  $Ax = b$ , but we want to do it efficiently, *i.e.*, with small  $\|x\|$ , so we add to the objective a term that penalises large  $x$ .
- Regularisation avoids large  $x$ .

# The Lasso

**Lasso** (*least absolute shrinkage and selection operator*)

$$\text{minimise } \|Ax - b\|_2^2 + \lambda \|x\|_1 \quad (4)$$

- Standardise  $A$ , so that each column has zero mean and unit variance.
- Solution for  $\lambda > 0$  traces out optimal trade-off curve (sweep  $\lambda$  from 0 to  $\infty$ ).
- Convex problem, so we know how to solve it efficiently.

Can also be written as:

$$\underbrace{\sum_{i=1}^n \left( y_i - \overbrace{\beta_0}^{\text{intercept}} - \sum_{j=1}^p \beta_j x_{ij} \right)^2}_{\text{Residual Sum of Squares}} + \lambda \sum_{j=1}^p |\beta_j| \quad (5)$$

# **Understanding the energy consumption of electric vehicle charging stations**

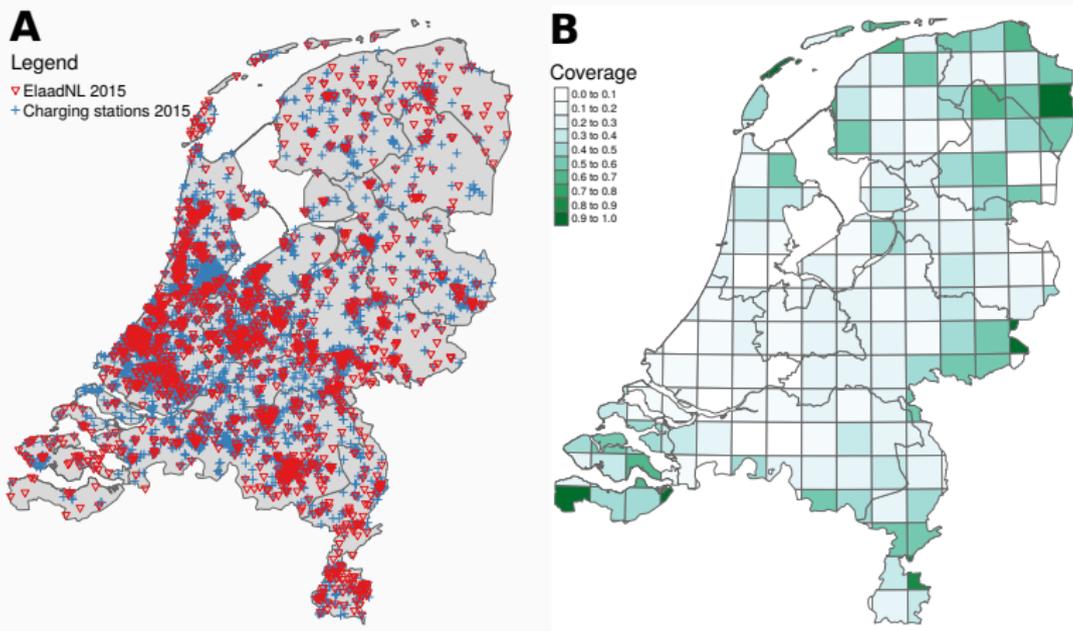
---

- Electric Vehicle (EV) charging
  - ELAAD
- Socio-economic, demographic, built environment and land use
  - Population cores.
  - Ambient population grid  $1000m \times 1000m$ , Landscan 2012.
  - Corine Land Use and Land Cover.
  - Neighbourhoods data.
  - Energy Atlas
  - Traffic flows
  - OSM amenities

## ELAAD raw dataset

- ElaadNL: Dutch research organisation involved in the development and deployment of EV charging technologies.
- 1,747 georeferenced charging stations.
- 54,000 users, each identified by a unique *id*.
- 1,060,763 charging events.
- Data collected between January 2012 and March 2016.

# ELAAD: spatial data

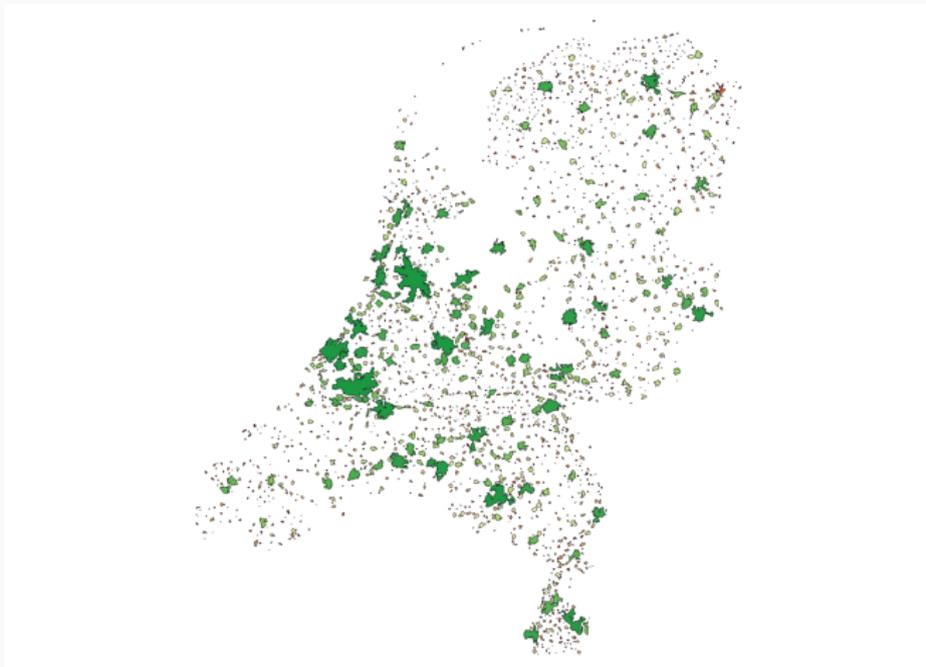


**Figure 1:** Public charging stations in the ElaadNL data set in 2015 (triangles) shown together with the Charging stations 2015 dataset (crosses). In the year 2015, 17 786 publicly available connectors for slow charging were operational in the NL. We identified 8 400 unique positions of charging stations, i.e. considering the distribution of connectors at charging stations observed in the ElaadNL dataset, this data covers 78.3% of all stations. In panel (B), we estimated the spatial representativeness of the ElaadNL data sets by calculating the ratio between the number of station in ELaadNL and in Charging stations 2015 located in squared cells of a regular grid. In the largest cities, Amsterdam and Rotterdam, the data contains a small percentage of all charging stations.

# Predictors: GIS data

- **Vector data**
  - **Polygon data**
    - population cores,
    - neighbourhoods data,
    - energy atlas,
    - liveability,
    - land use and land cover (urban atlas, CBS land cover).
  - **Polyline data**
    - traffic flow data.
  - **Point data**
    - OSM amenities,
    - OpenChargeMap.
- **Raster data**
  - LandScan - ambient population.

## Vector polygon data: Population cores (2168 cores)



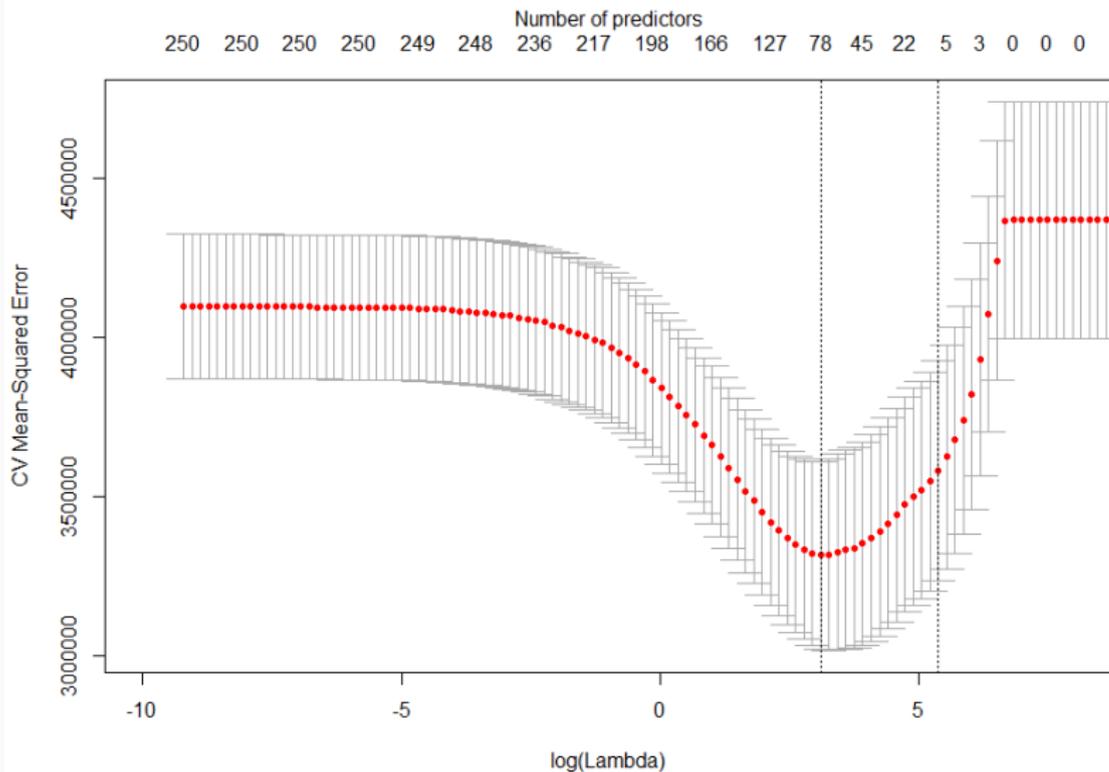
**Figure 2:** Population cores are continuous spatial units with at least 25 homes or 50 inhabitants (102 predictors). Source: Statistics Netherlands <https://opendata.cbs.nl/>

## Polygon data: Population cores (2168 cores, 84 attributes)

- number of persons in private household,
- number of persons in private households, 0 to 15 years,
- number of persons in private households, 15 to 25 years,
- number of persons in private households 25 to 45 years,
- number of persons in private households, 45 to 65 years,
- number of persons in private households 65 years or older,
- number of persons in one-person households,
- number of people in multi-person household with children,
- number of people in multi-person household without children,
- percentage of working population, 15-24 years old,
- number of households of two persons,
- number of households of three persons,
- the number of residential units,
- ...



# Lasso



## Lasso fit

- Lasso does variable selection on 240 predictors.
- At the optimal  $\lambda$ , we reduce the number of predictors to 79 (about 1/3 of the original predictors).
- $R^2 = 0.362$  at optimum value ( $R^2(\text{adjusted}) = 0.316$ ).
- The  $\lambda_{min}$  is the one which minimises the error in cross-validation.  
The  $\lambda_{1\sigma}$  is the  $\lambda$  value within 1 standard error of  $\lambda_{min}$ .

Coefficient	Meaning
-27.8627	The percentage of working population working in Mining, Manufacturing and Construction.
20.3313	The percentage of working population employed in commercial services
18.169	The percentage working population engaged in agriculture, forestry and fisheries, industry, commercial and non-commercial services
-0.1884	The percentage of the number of multi-person households without children
4.1384	Number of business Services
14.5717	Property unknown (no link between the addresses of the Key Registers Addresses and the housing register Cadaster).
57.8468	Average income per inhabitant
-0.9093	Average distance of all residents in an area to the nearest shops for groceries.

# **Data-driven discovery of dynamical systems**

---

# Data-driven discovery of dynamical systems

Goal of computationally-oriented scientists:

Inferring a (typically nonlinear) model from observations that both correctly identifies the underlying dynamics *and* generalises qualitatively and quantitatively to unmeasured parts of the phase, parameter, or application space.

- ODE or PDE system described by

$$u_t = N(u, x, t; \overline{\mu}) \quad (7)$$

- Our objective is to discover  $N(\cdot)$  given only time-series measurements of the system.
- A key assumption (prior) is that the true  $N(\cdot)$  is comprised of only a few terms, making the model sparse in the space of all possible combinations of functions.
- For example, Burgers' equation

$$N = -uu_x + \mu u_{xx} \quad (8)$$

and the harmonic oscillator

$$N = -i\mu x^2 u - i\hbar u_{xx}/2 \quad (9)$$

each have only two terms.



# Data-driven discovery of dynamical systems: method

Method:

- Construct a library  $\Theta(U)$  of candidate linear, nonlinear, and partial derivative terms for the right-hand side.
- Each column of  $\Theta(U)$  contains the values of a candidate term evaluated using the collected data.
- In this library, one can write the dynamics as

$$U_t = \Theta(U)\xi \quad (10)$$

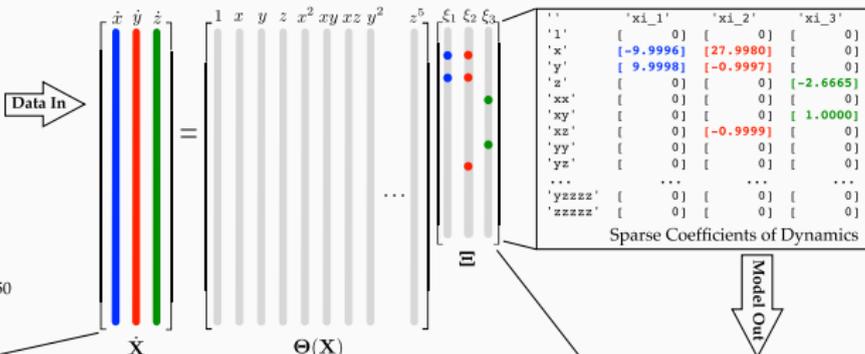
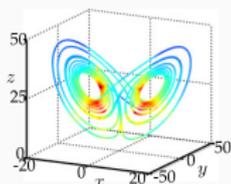
where

- $U_t$  is a vector of time derivatives of the measurement data.
- $\xi$  is a sparse vector, with each nonzero entry corresponding to a functional term to be included in the dynamics.
- Finding the sparsest vector  $\xi$  consistent with the measurement data is now feasible with advanced methods in sparse regression, which makes it possible to find the most parsimonious model while circumventing a combinatorial search.

# Identifying the Lorenz Equations

## I. True Lorenz System

$$\begin{aligned}\dot{x} &= \sigma(y - x) \\ \dot{y} &= x(\rho - z) - y \\ \dot{z} &= xy - \beta z.\end{aligned}$$



## III. Identified System

$$\begin{aligned}\dot{x} &= \Theta(\mathbf{x}^T)\xi_1 \\ \dot{y} &= \Theta(\mathbf{x}^T)\xi_2 \\ \dot{z} &= \Theta(\mathbf{x}^T)\xi_3\end{aligned}$$

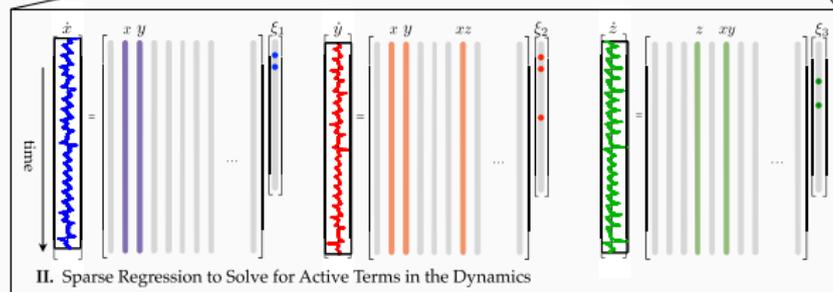
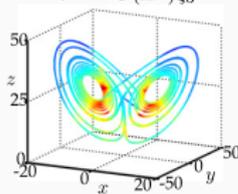


Fig. 1. Schematic of the SINDy algorithm, demonstrated on the Lorenz equations. Data are collected from the system, including a time history of the states  $\mathbf{X}$  and derivatives  $\dot{\mathbf{X}}$ ; the assumption of having  $\mathbf{X}$  is relaxed later. Next, a library of nonlinear functions of the states,  $\Theta(\mathbf{X})$ , is constructed. This nonlinear feature library is used to find the fewest terms needed to satisfy  $\dot{\mathbf{X}} = \Theta(\mathbf{X})\Xi$ . The few entries in the vectors of  $\Xi$ , solved for by sparse regression, denote the relevant terms in the right-hand side of the dynamics. Parameter values are  $\sigma = 10, \beta = 8/3, \rho = 28, (x_0, y_0, z_0)^T = (-8, 7, 27)^T$ . The trajectory on the Lorenz attractor is colored by the adaptive time step required, with red indicating a smaller time step.