

The Textile Plot

A Key to Interaction Through Data

Kumasaka Natsuhiko
Fundamental Sci & Tech
Keio Univ

Ritei Shibata
Dept of Math
Keio Univ

What is Textile Plot?

Possible enhancements of Textile Plot

Interactive environment for understanding data through Textile Plot

Textile Plot

Kumasaka and Shibata 2006, 2007

High dimensional data visualisation
according to parallel coordinate
system

Necessary and sufficient
background information

Any kind of data vector (variable)
including missing values

Automobile Data

Consumer Reports (USA) 1990, S-PLUS

ID	price	country	reliability	type	weight	disp	hp	Manufacturer	mileage
<i>Toyota Corolla</i>	8748	Japan/USA	much better	Small	2390	97	102	Toyota	29
<i>Peugeot 405</i>	15930	France	NA	Compact	2575	116	120	Peugeot	24
<i>Buick Le Sabre V6</i>	16145	USA	average	Large	3325	231	165	Buick	23
<i>Nissan Maxima V6</i>	17899	Japan	much better	Medium	3200	180	160	Nissan	22
<i>Mitsubishi Wagon</i>	14929	Japan	NA	Van	3415	143	107	Mitsubishi	20
<i>Plymouth Laser</i>	10855	USA	NA	Sporty	2840	107	92	Plymouth	26

⋮

<i>Volvo 240</i>	18450	Sweden	average	Compact	2985	141	114	Volvo	23
<i>Volkswagen Jetta</i>	9995	Germany	average	Small	2330	109	100	Volkswagen	26

Records

60 cars

Data Vectors (9 dimensional)

Price

Country

Reliability

Type

Weight

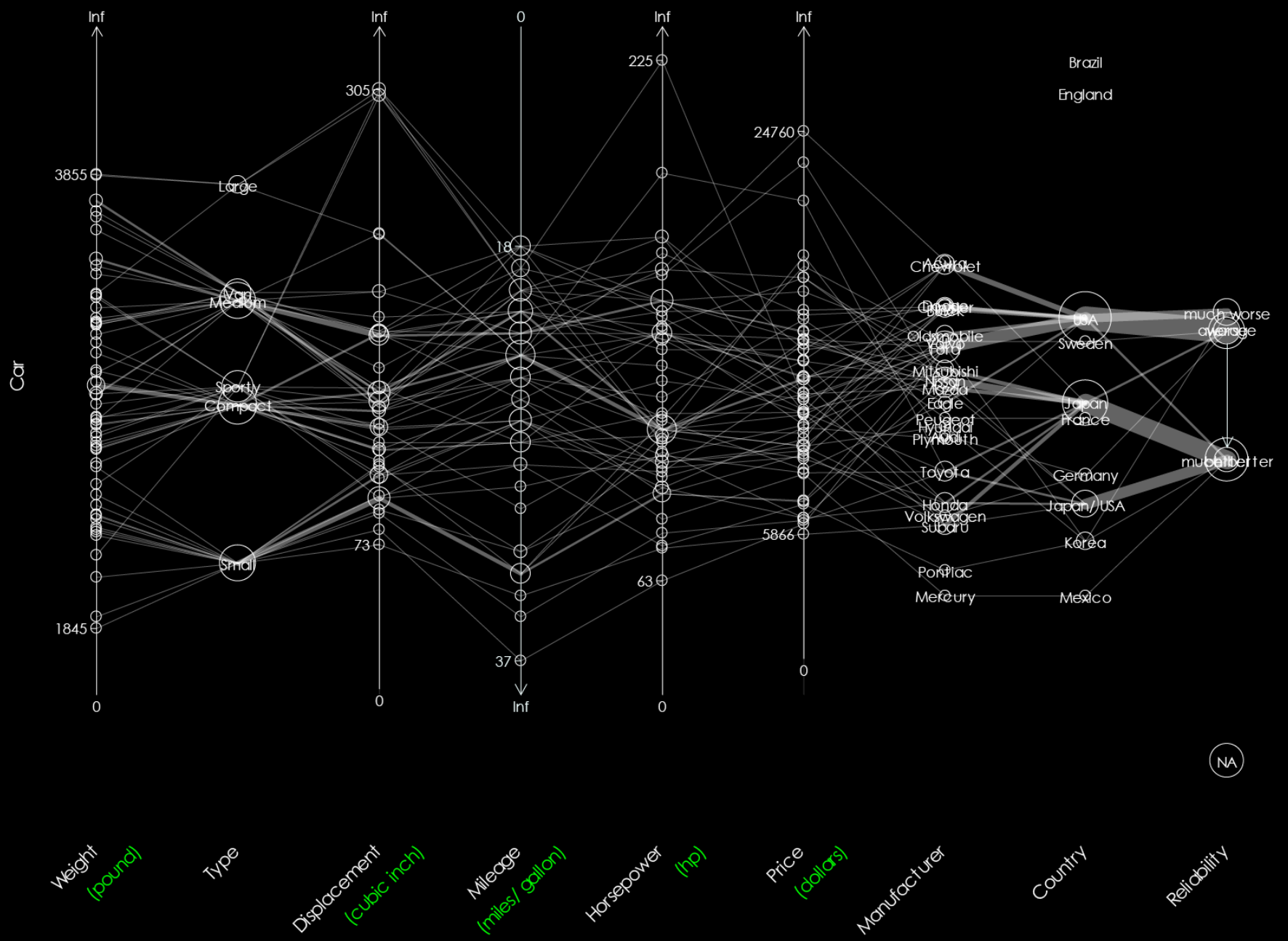
Displacement

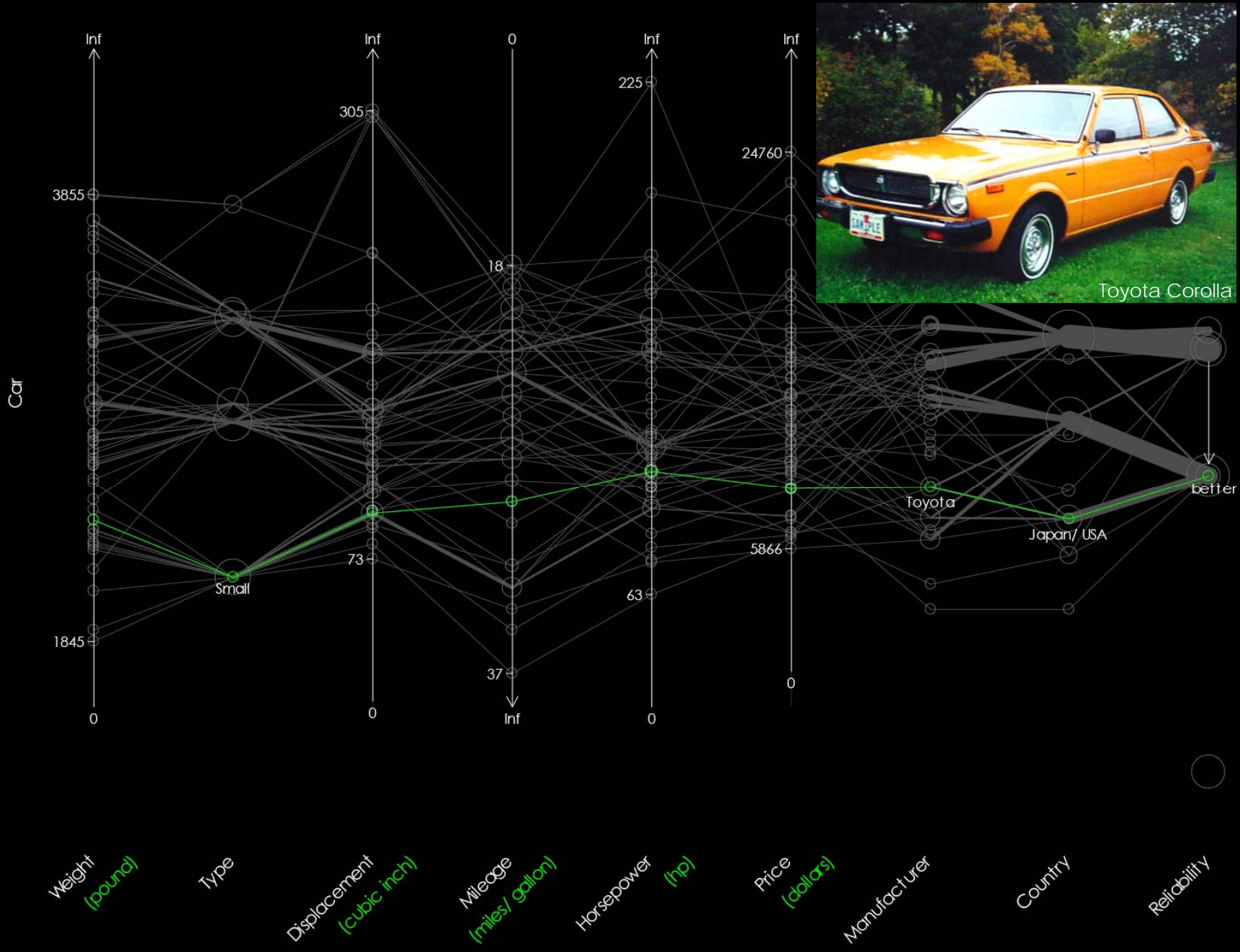
Horsepower

Manufacturer

Mileage

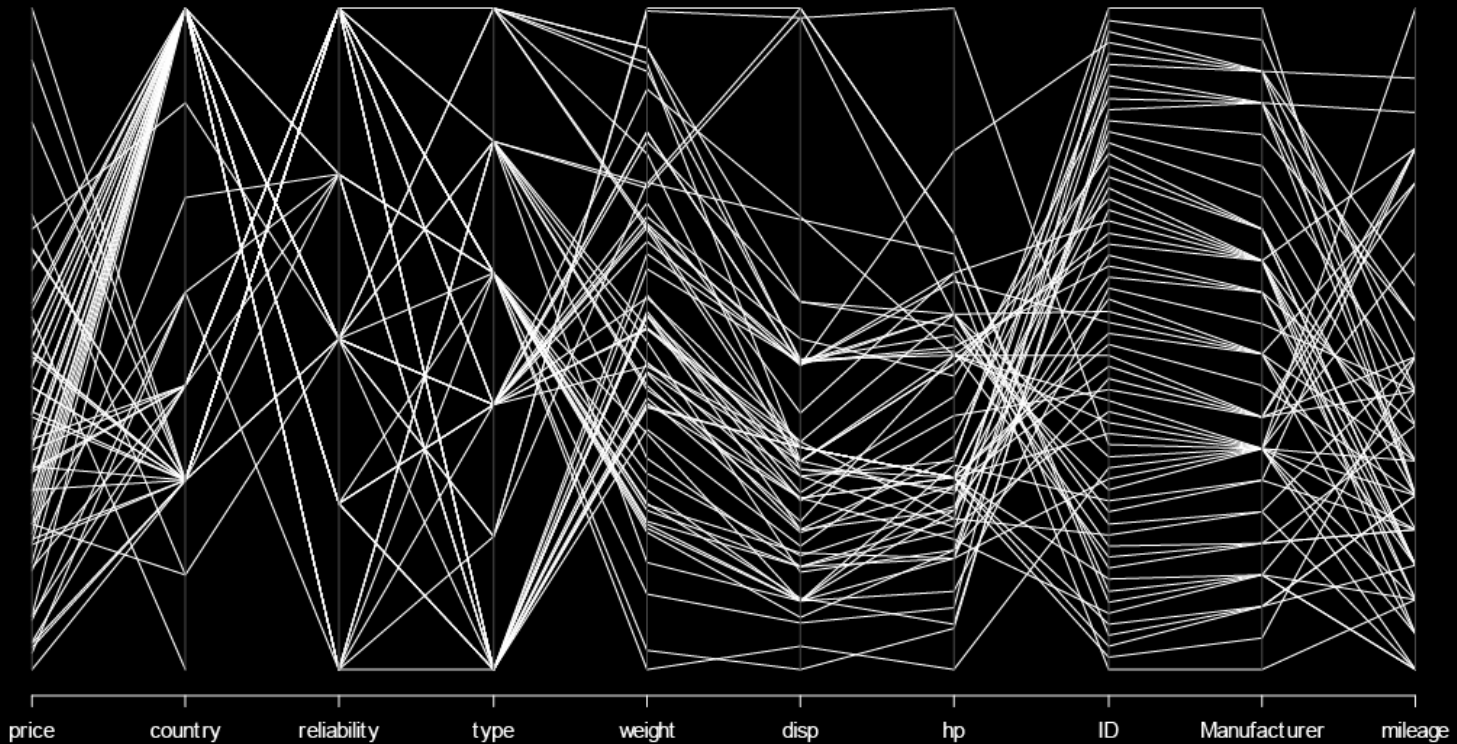






Parallel Coordinate Plots

Inselberg 1985, Wegman 1990



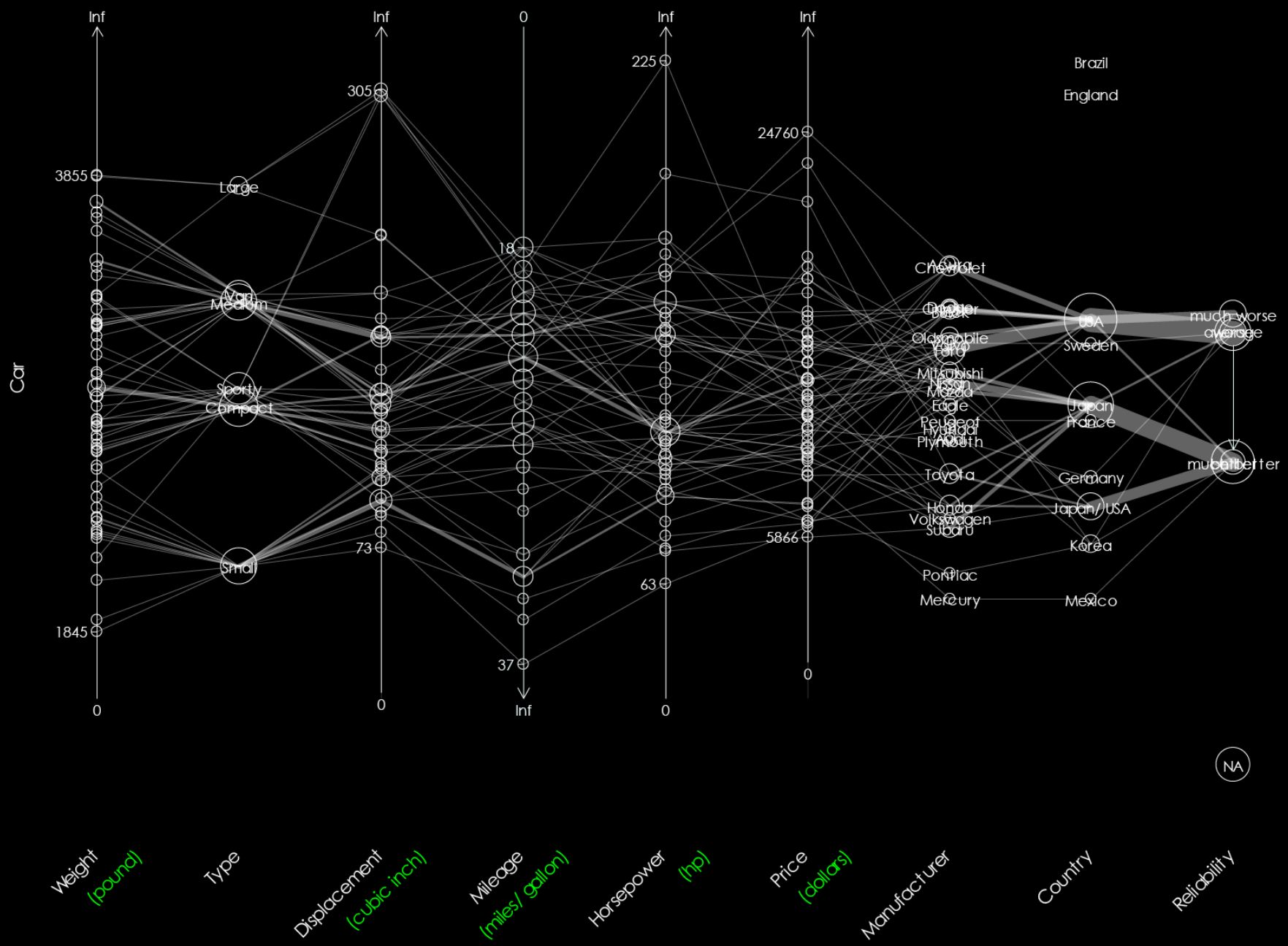
Location and scale of each axis

Independently chosen

Points fill up the range

Levels of categorical data

Equispaced on an axis



Brazil
England

NA

Weight
(pound)

Type

Displacement
(cubic inch)

Mileage
(miles/gallon)

Horsepower
(hp)

Price
(dollars)

Manufacturer

Country

Reliability

Car

Inf
3855
1845
0

Inf
305
0

0
18
37
Inf

Inf
225
63
0

Inf
24760
5866
0

Chevrolet
Buick
Oldsmobile
Ford
Mitsubishi
Mazda
Eagle
Peugeot
Plymouth
Toyota
Honda
Volkswagen
Subaru
Pontiac
Mercury

Brazil
England
USA
Sweden
Japan
France
Germany
Japan/USA
Korea
Mexico

much worse average
much better

Large

Medium

Sporty
Compact

Small

305

18

225

24760

3855

1845

73

37

63

5866

Chevrolet

Buick

Oldsmobile

Ford

Mitsubishi

Mazda

Eagle

Peugeot

Plymouth

Toyota

Honda

Volkswagen

Subaru

Pontiac

Mercury

Brazil

England

USA

Sweden

Japan

France

Germany

Japan/USA

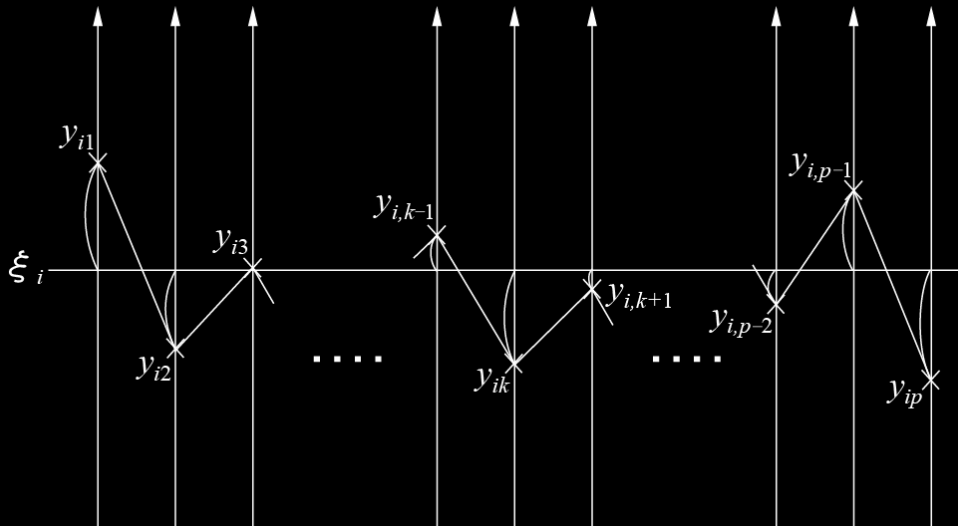
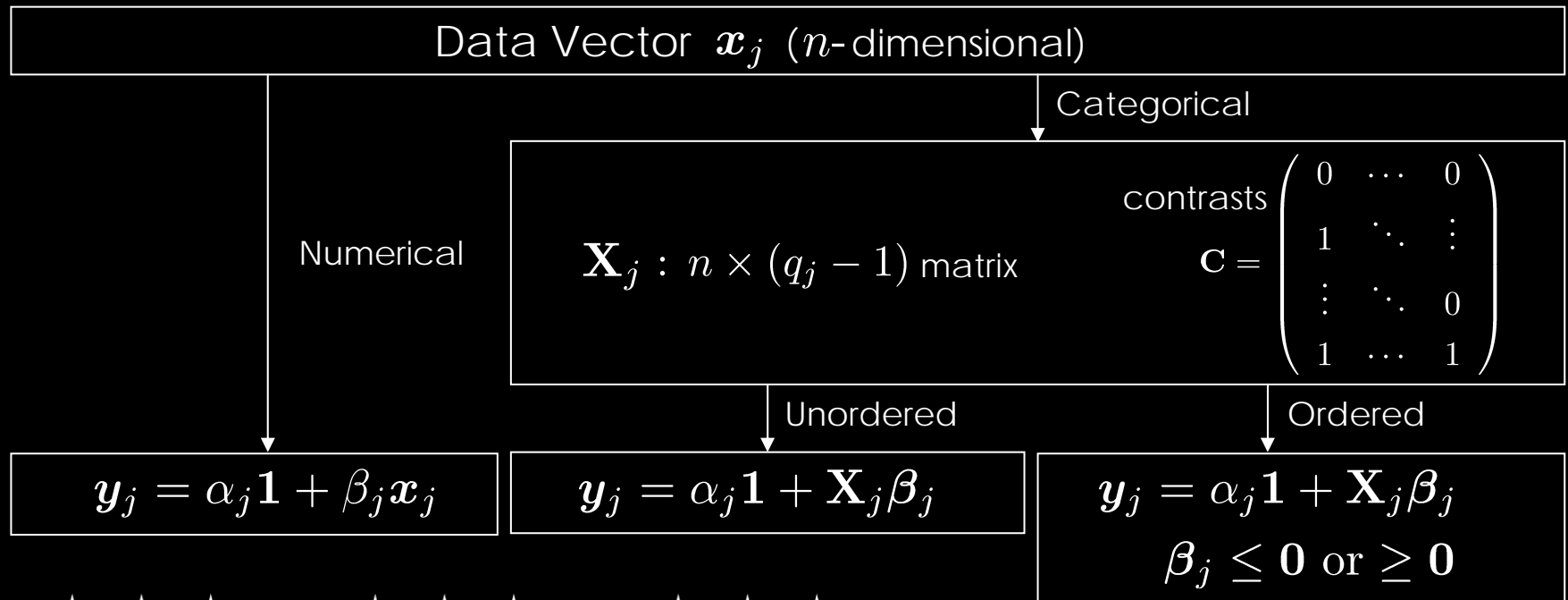
Korea

Mexico

much worse average

much better

Choice of Locations and Scales



Deviations for the i th observation

The sum of squares between **coordinate vectors** and an **ideal coordinate vector** is minimised.

$$\sum_{j=1}^p \|\mathbf{y}_j - \boldsymbol{\xi}\|^2 \rightarrow_{\{\alpha_j, \beta_j\}_{j=1}^p, \boldsymbol{\xi}} \min$$

SNP Data on Calpain-10

Iwasaki, et al. 2005

Data Vectors (14-dimensional)

Phenotype (affected or non-affected)

12 Genotypes on SNP loci

Genotype on INDEL locus

Joint work with

Dr Kamitsuji at Stagen (<http://www.stagen.co.jp/>)

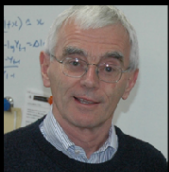
Goal

To know the danger of a disease from genotypes

SNP / INDEL Locus

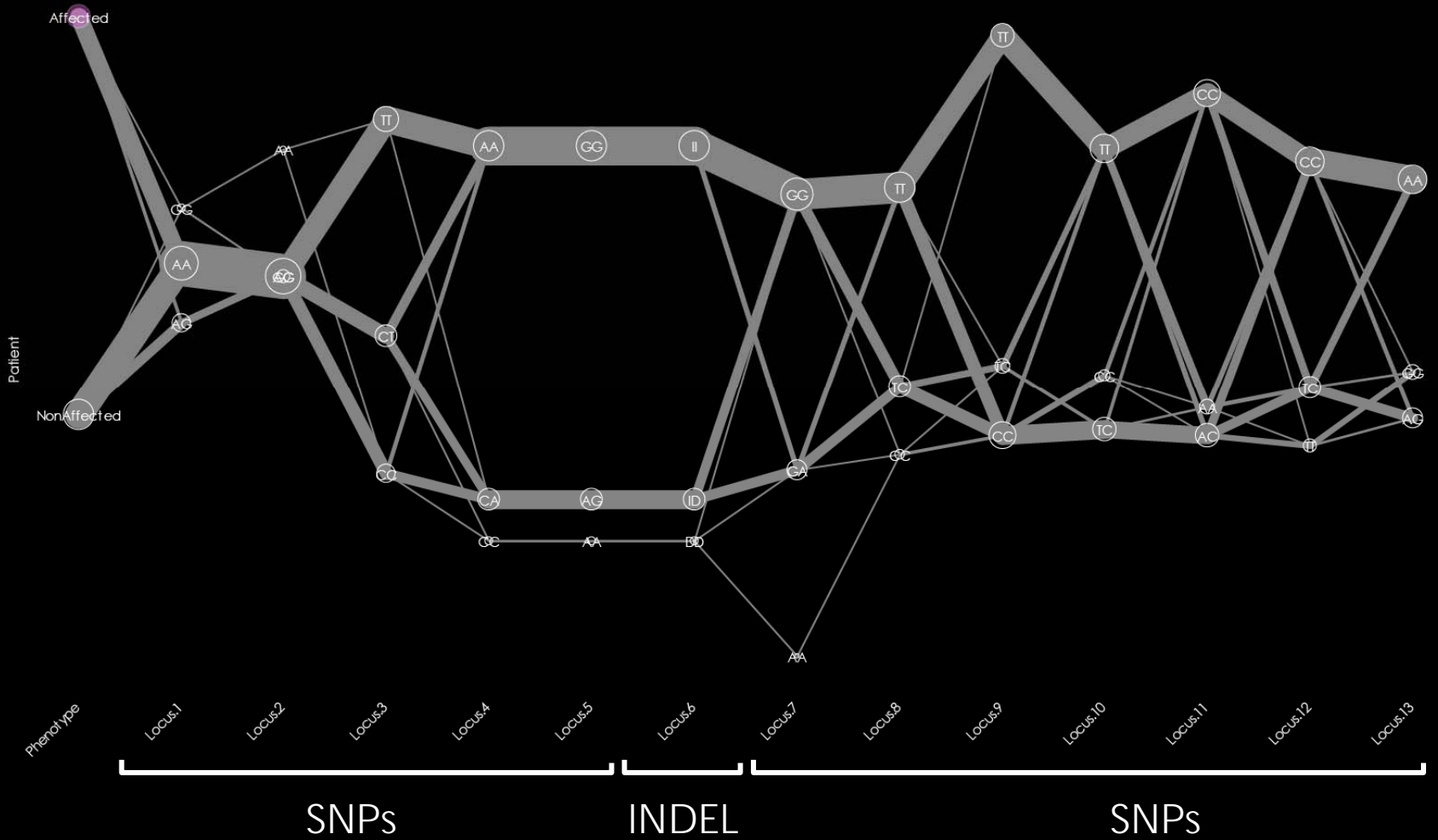
3 Genotype Sequences

⋮	⋮	⋮	
AA	AA	AA	
TT	TT	TT	
GG	GG	GG	
CC	GC	GG	← SNP Locus
AA	AA	AA	← SNP Locus
AA	AA	TT	
CC	CC	CC	
CC	CC	CC	
TT		T	← INDEL Locus
GG		G	
CC		C	
AA	AA	AA	
TT	TT	TT	
CG	CC	CC	← SNP Locus
CC	CC	CC	
AA	AA	AA	
⋮	⋮	⋮	

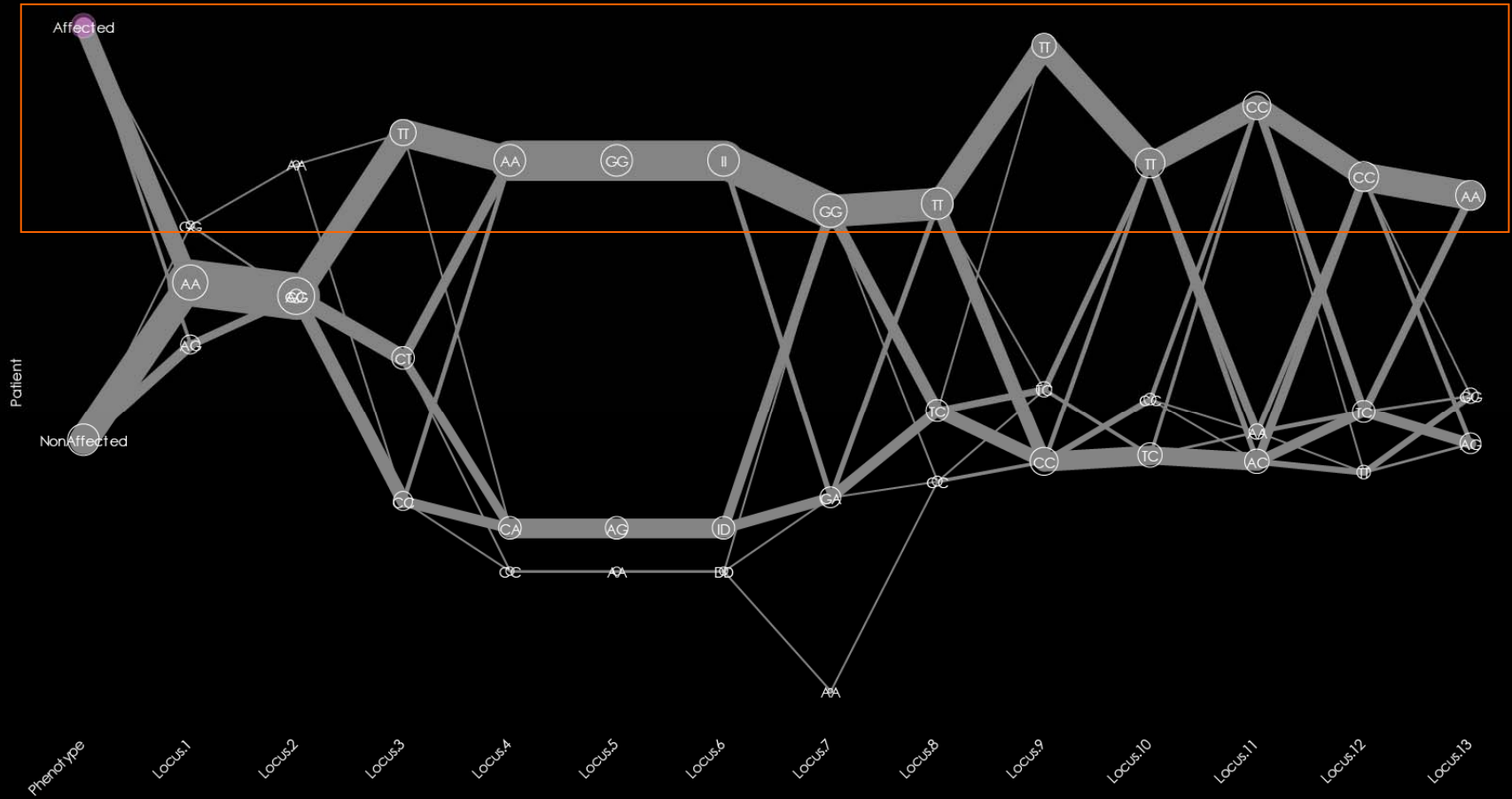


SNP = Single Nucleotide Polymorphism
 INDEL = INsertion and DEletion

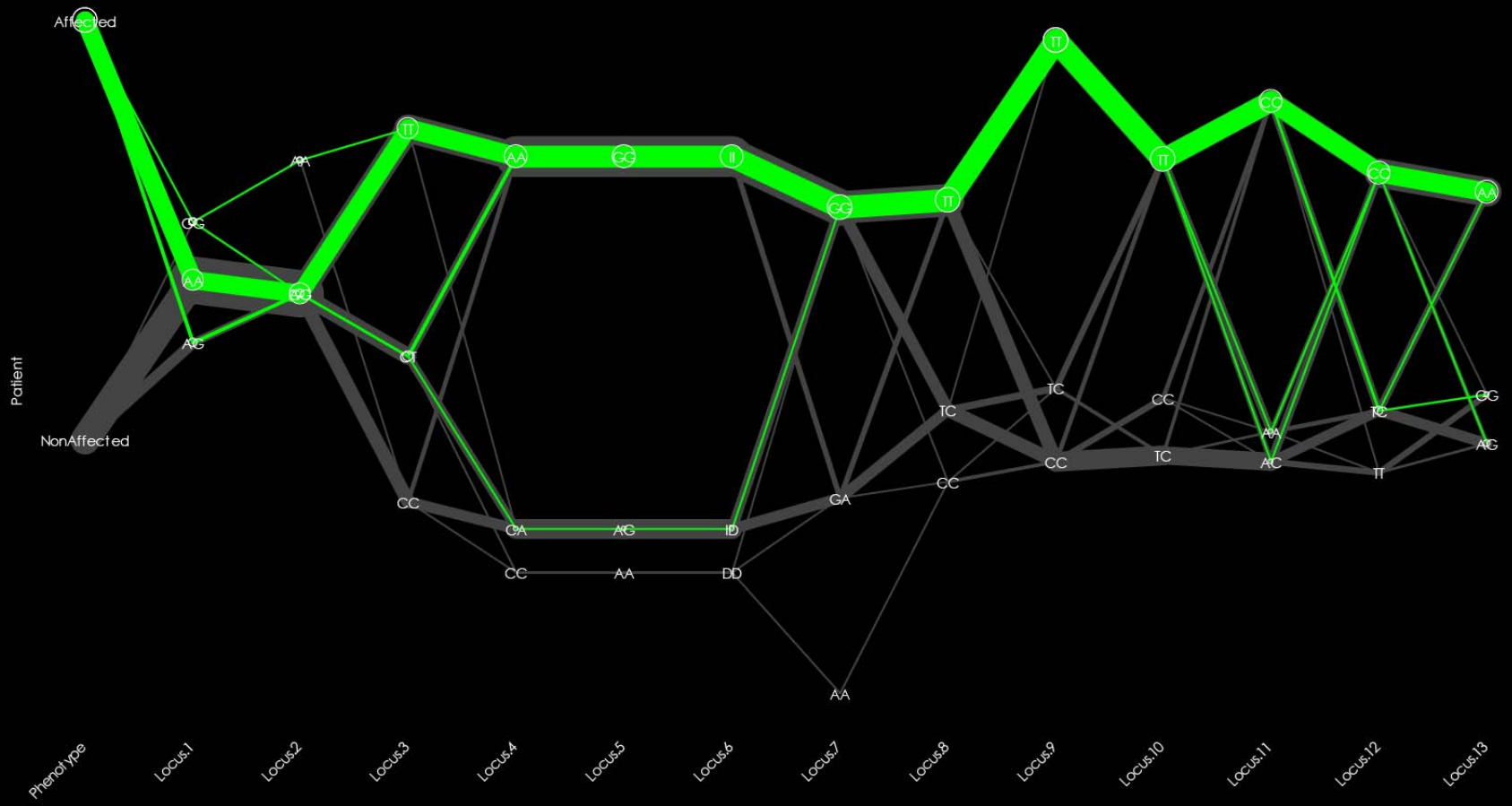
Textile Plot of SNP Data



Recessive Inheritance



Enhancements of Textile Plot



Enhancements of Textile Plot

Transform

- single data vector (log, inverse, etc)
- multiple data vectors (normalisation, etc)

Substitute Data Point

Extract/Remove

- records
- data vectors

Re-weight

- outliers

Order Data Vectors

- a priori
- distance to ideal coordinates
- distance between coordinate vectors

Highlight Connecting Lines

Select Graphical Parameters

- color, transparency, size, etc.

Enhancements of Textile Plot

Transform

- single data vector (log, inverse, etc)
- multiple data vectors (normalisation, etc)

Substitute Data Point

Extract/Remove

- records
- data vectors

Re-weight

- outliers

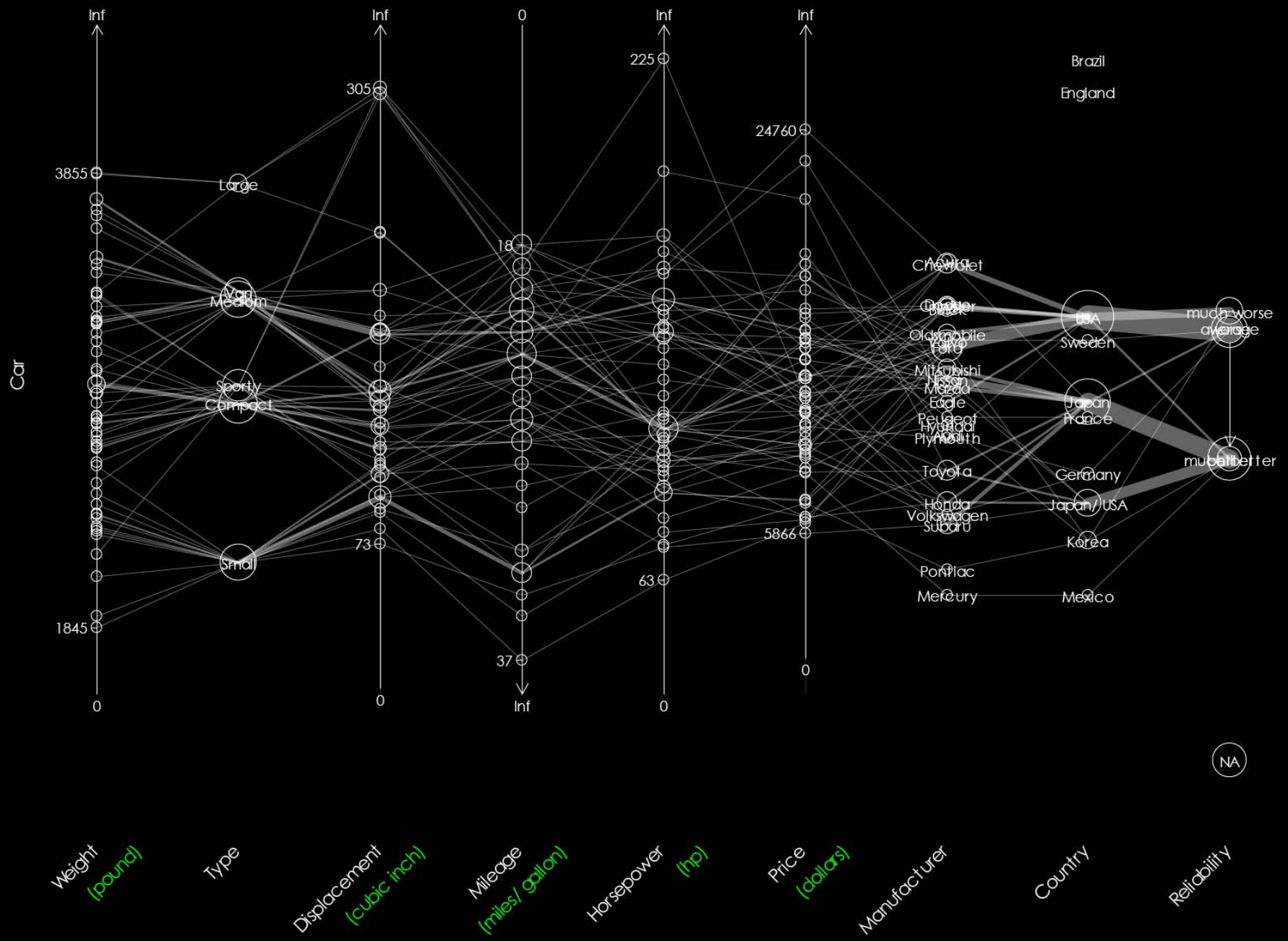
Order Data Vectors

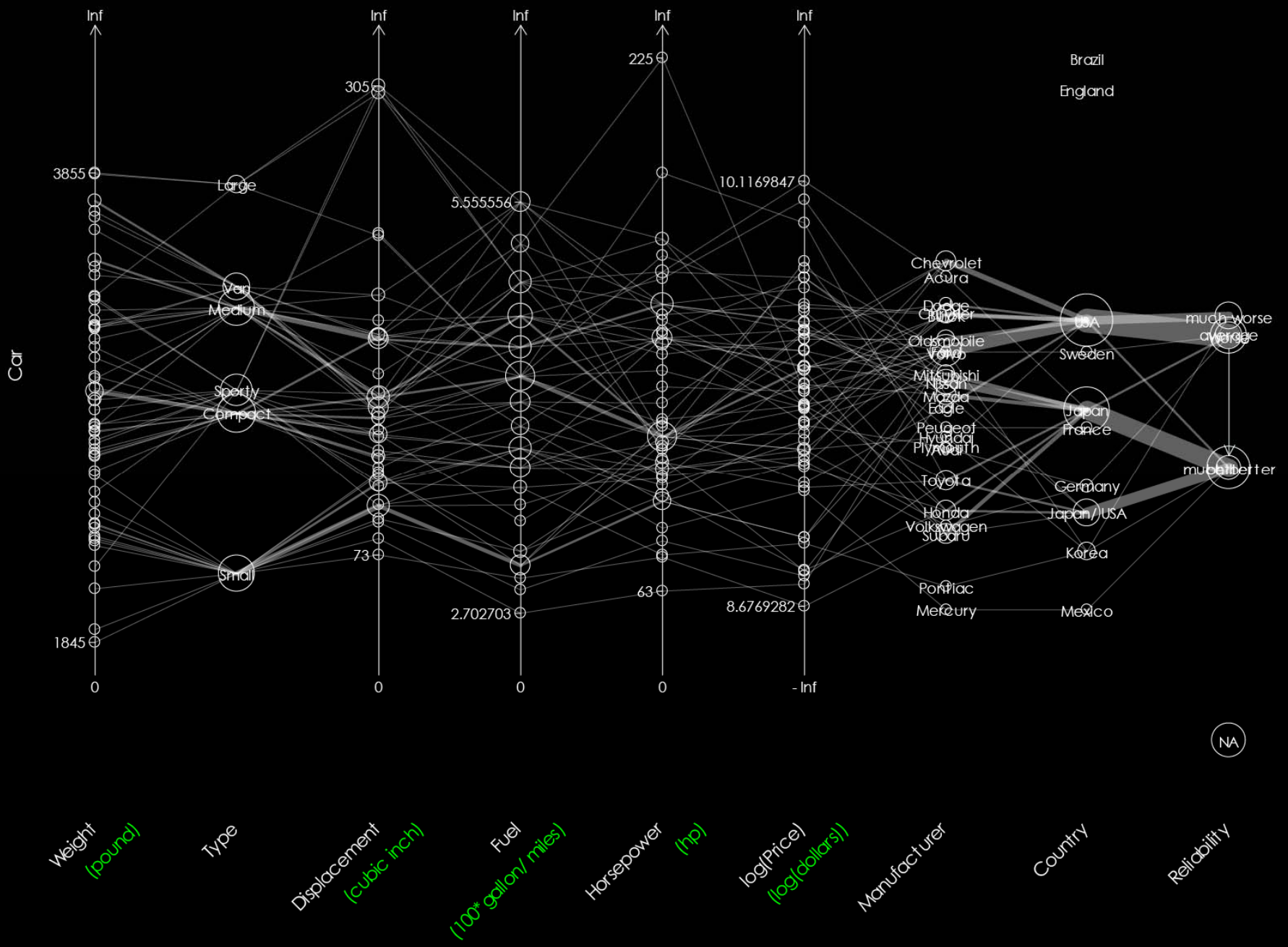
- a priori
- distance to ideal coordinates
- distance between coordinate vectors

Highlight Connecting Lines

Select Graphical Parameters

- color, transparency, size, etc.





Enhancements of Textile Plot

Transform

- single data vector (log, inverse, etc)
- multiple data vectors (normalisation, etc)

Substitute Data Point

Extract/Remove

- records
- data vectors

Re-weight

- outliers

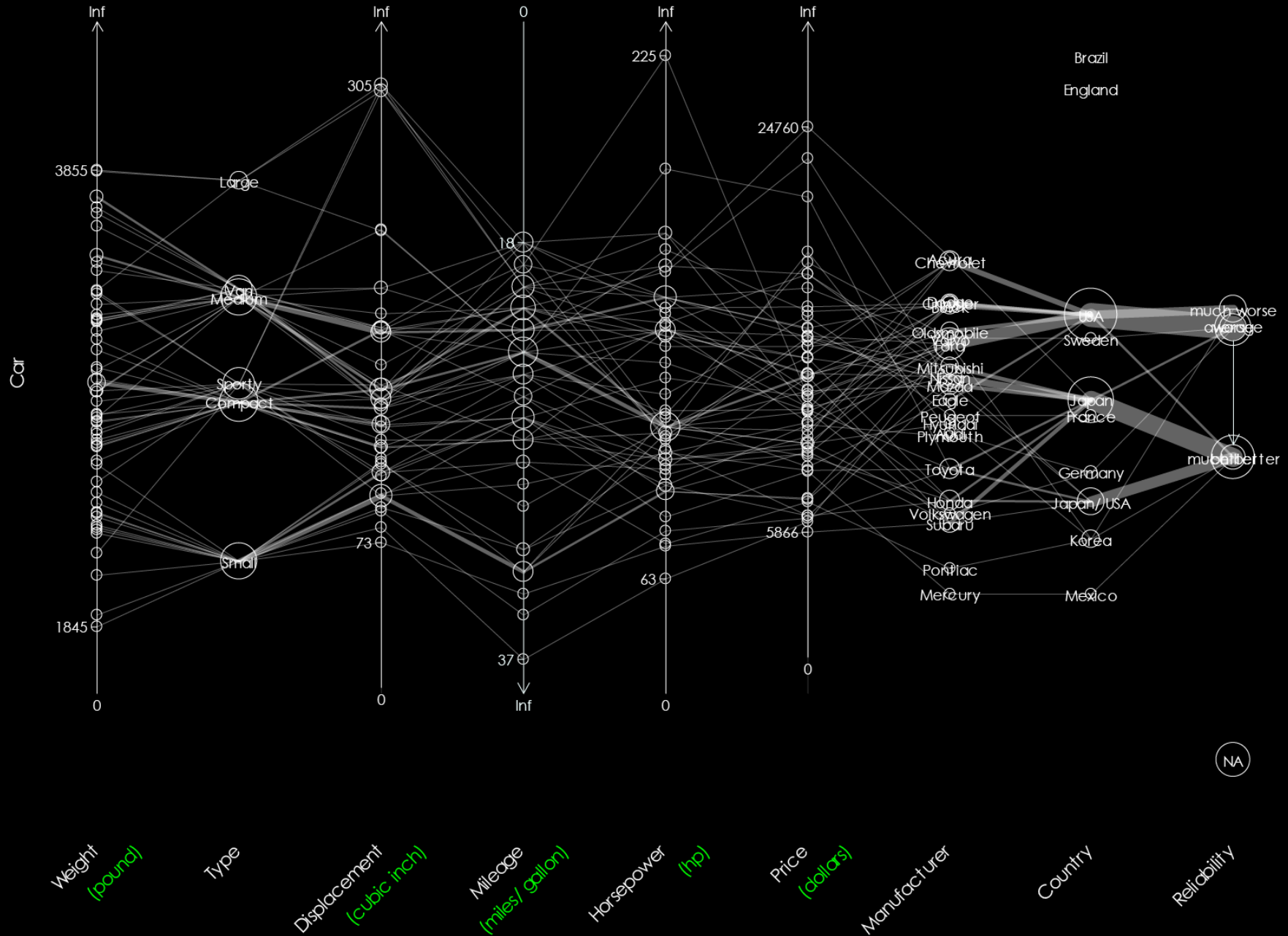
Order Data Vectors

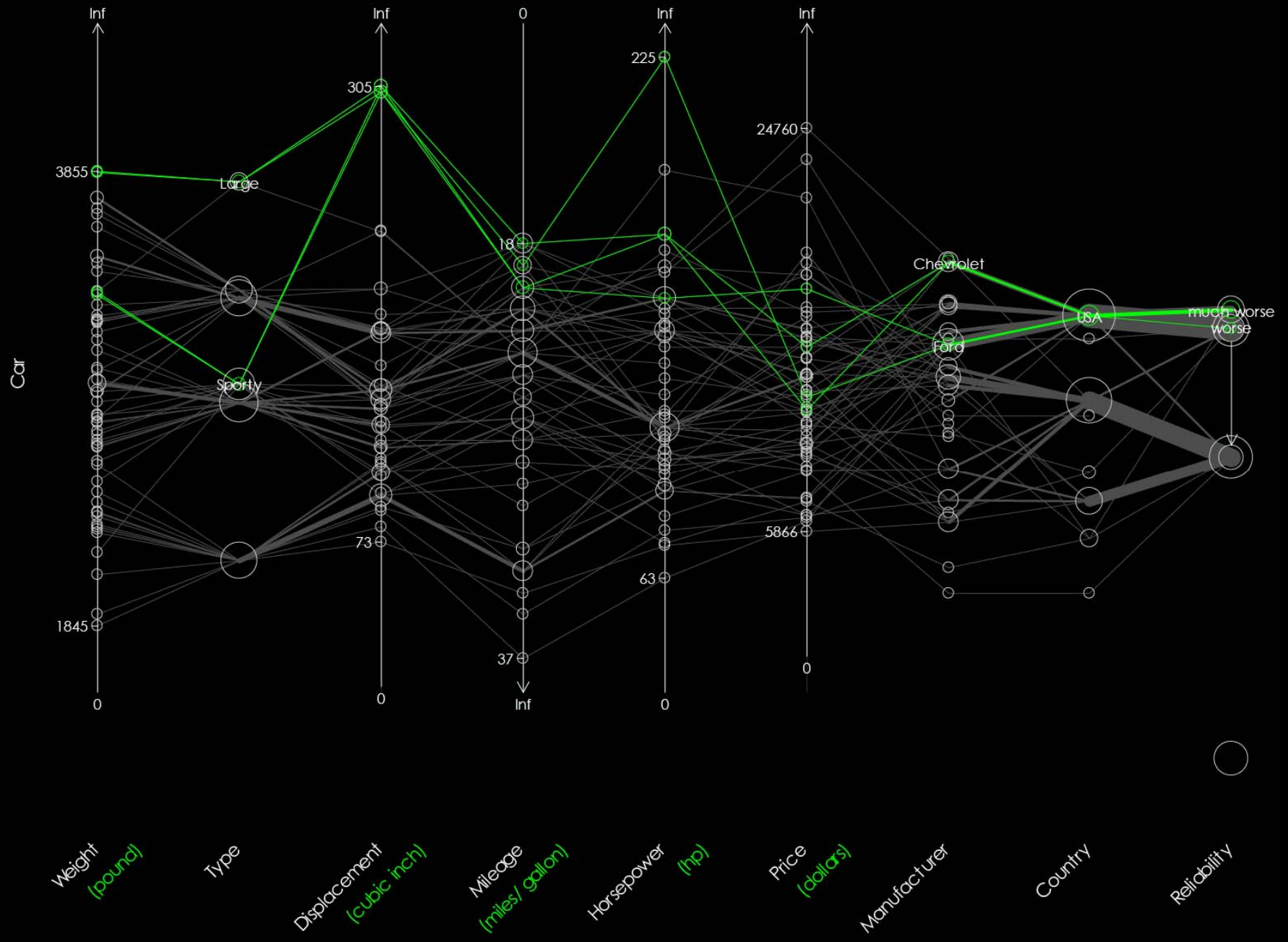
- a priori
- distance to ideal coordinates
- distance between coordinate vectors

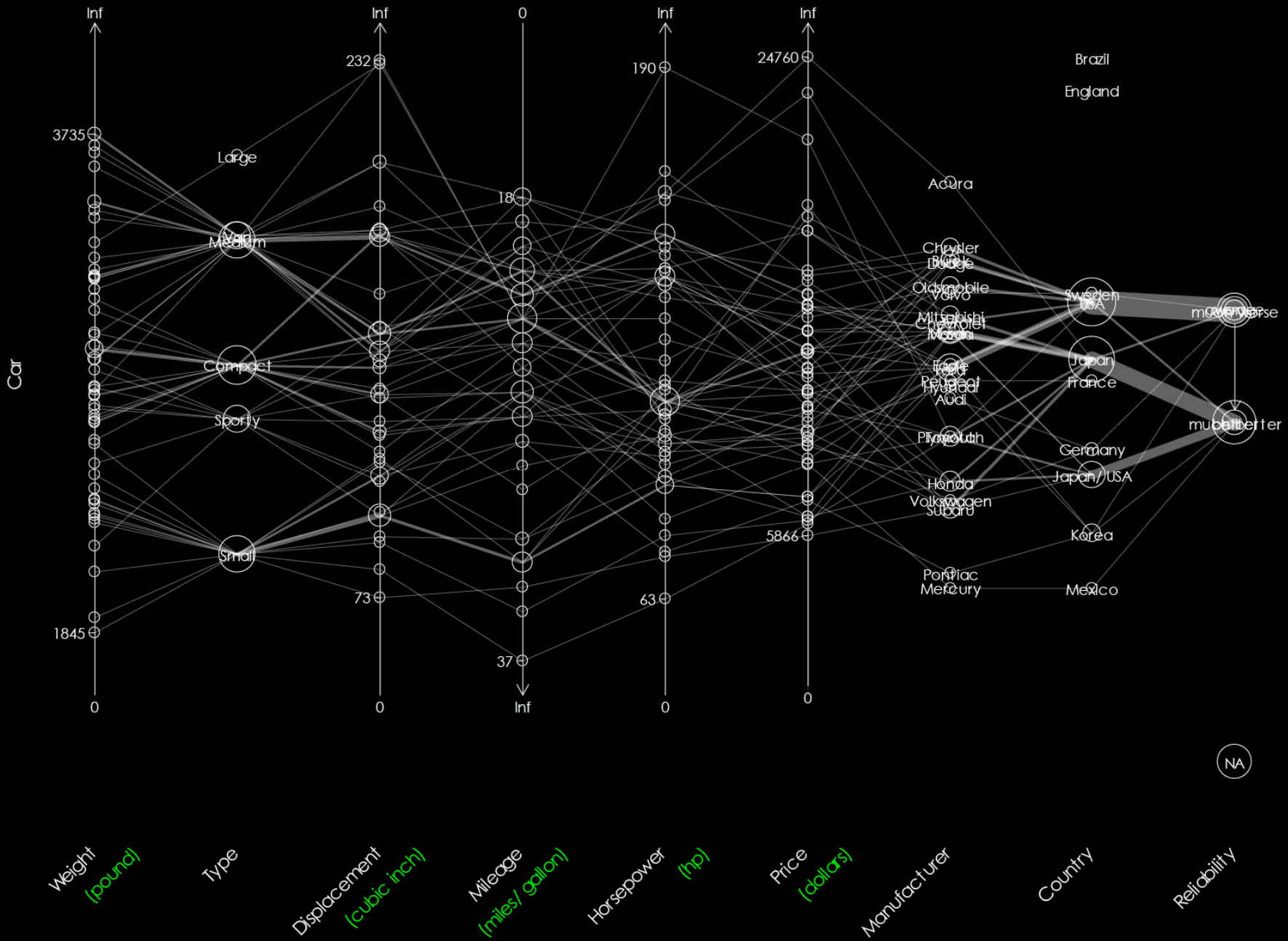
Highlight Connecting Lines

Select Graphical Parameters

- color, transparency, size, etc.







Enhancements of Textile Plot

Transform

- single data vector (log, inverse, etc)
- multiple data vectors (normalisation, etc)

Substitute Data Point

Extract/Remove

- records
- data vectors

Re-weight

- outliers

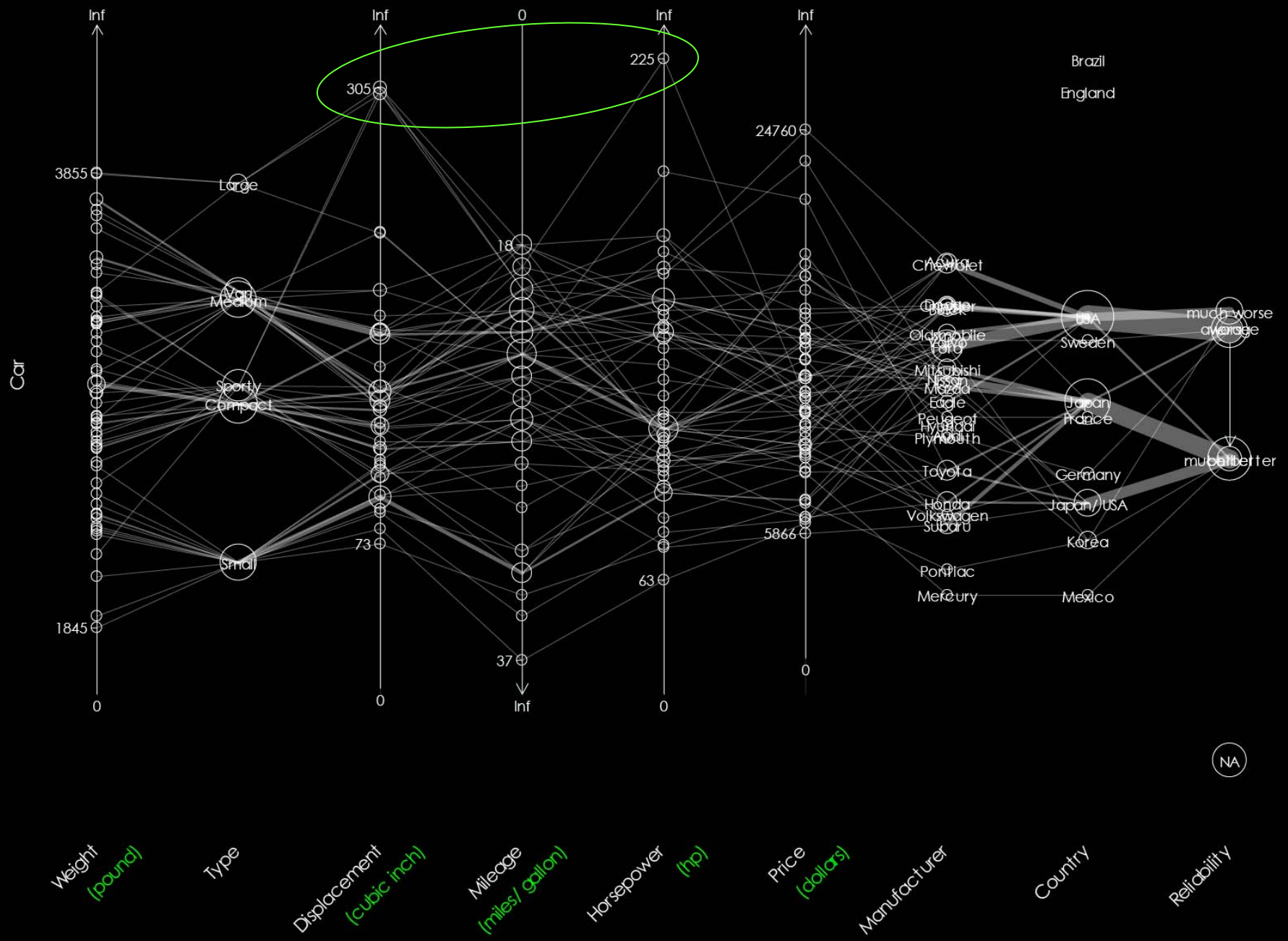
Order Data Vectors

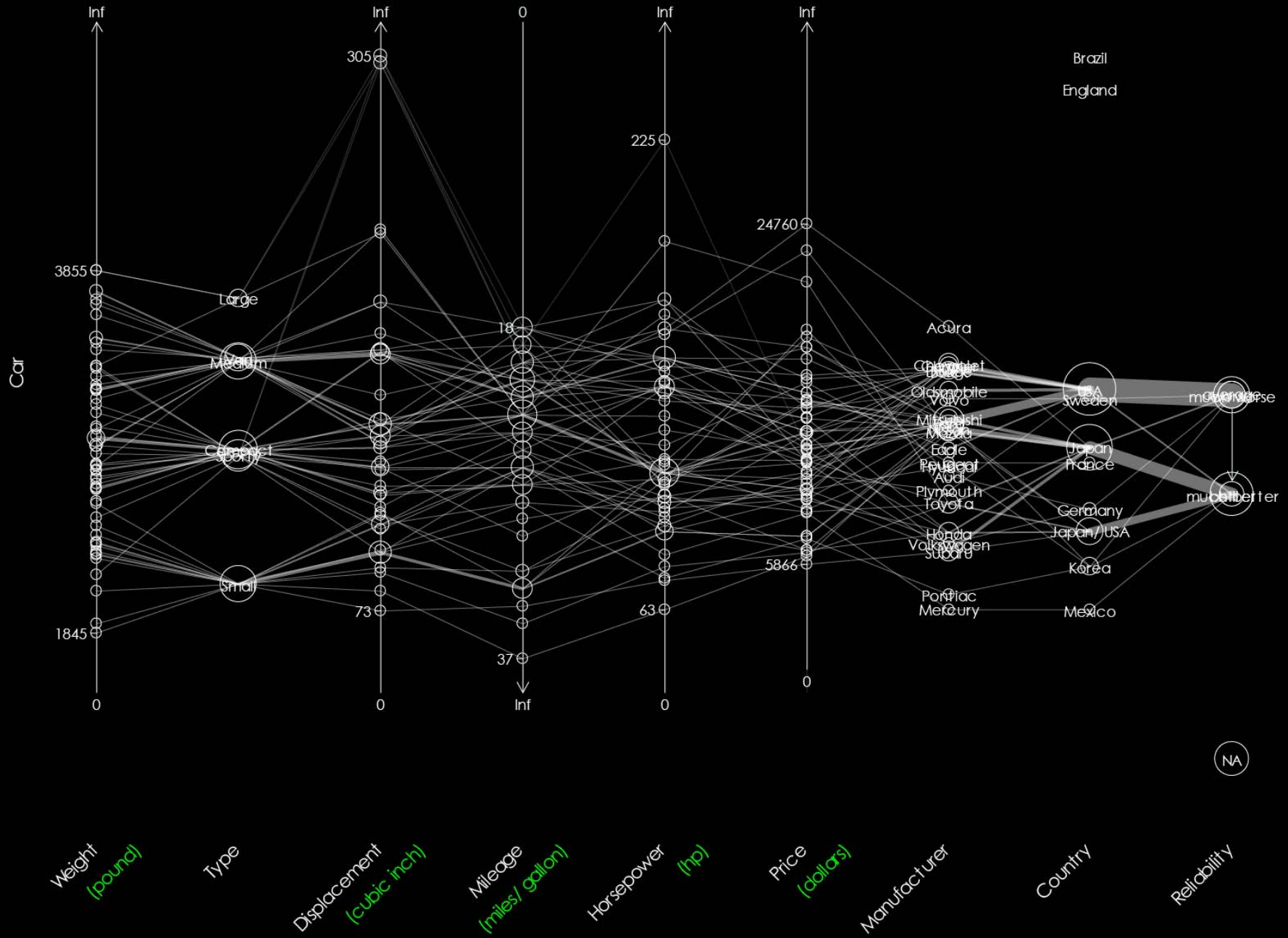
- a priori
- distance to ideal coordinates
- distance between coordinate vectors

Highlight Connecting Lines

Select Graphical Parameters

- color, transparency, size, etc.





Enhancements of Textile Plot

Transform

- single data vector (log, inverse, etc)
- multiple data vectors (normalisation, etc)

Substitute Data Point

Extract/Remove

- records
- data vectors

Re-weight

- outliers

Order Data Vectors

- a priori
- distance to ideal coordinates
- distance between coordinate vectors

Highlight Connecting Lines

Select Graphical Parameters

- color, transparency, size, etc.

Order of Data Vectors

Data Vectors $\mathbf{x}_1, \dots, \mathbf{x}_p$

A priori order

Distance to Ideal coord vec

The further left coordinate vector is closer to the ideal coordinate vector

$$\|\mathbf{y}_j - \hat{\boldsymbol{\xi}}\|^2 \leq \|\mathbf{y}_{j+1} - \hat{\boldsymbol{\xi}}\|^2$$

Ideal coordinate vector gives a set of ideal coordinates for each records

Distance between data vecs

Sum of absolute values of slopes

$$d(\mathbf{x}_j, \mathbf{x}_k) = \sum_{i=1}^n |y_{ij} - y_{ik}|$$

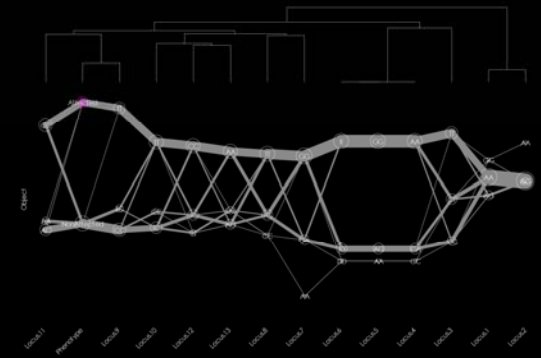
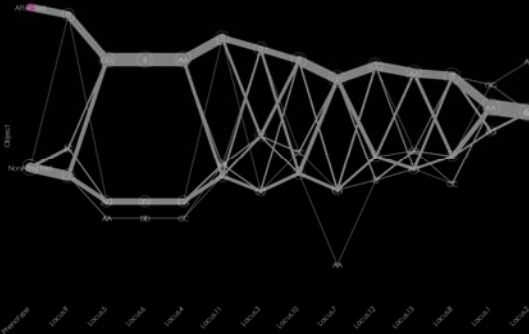
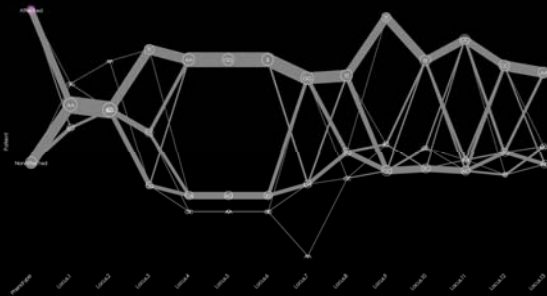
Linkage method

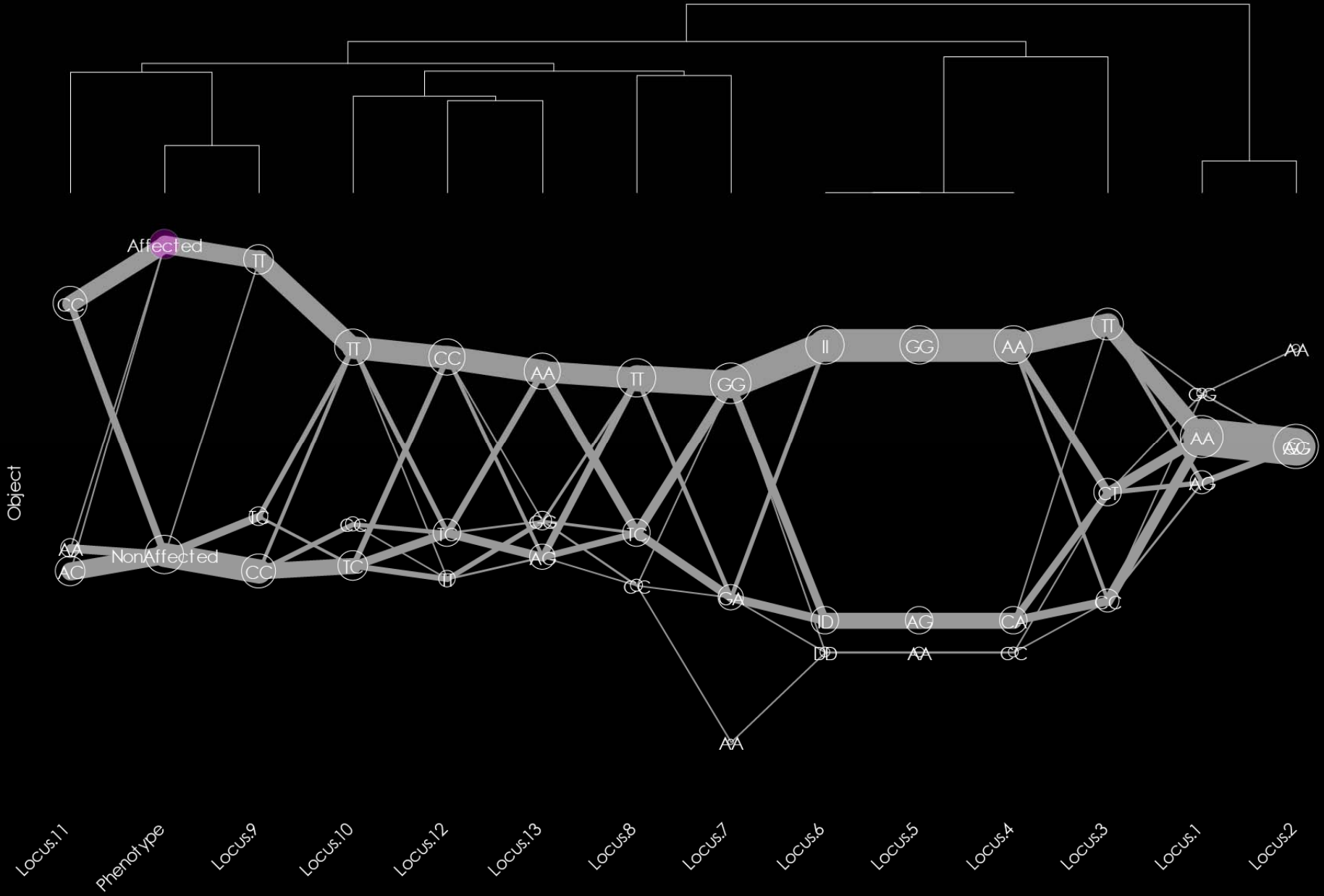
Ordered single end-linkage (Hurley 2004)

As it is

Classification of records

Classification of variables





Enhancements of Textile Plot

Transform

- single data vector (log, inverse, etc)
- multiple data vectors (normalisation, etc)

Substitute Data Point

Extract/Remove

- records
- data vectors

Re-weight

- outliers

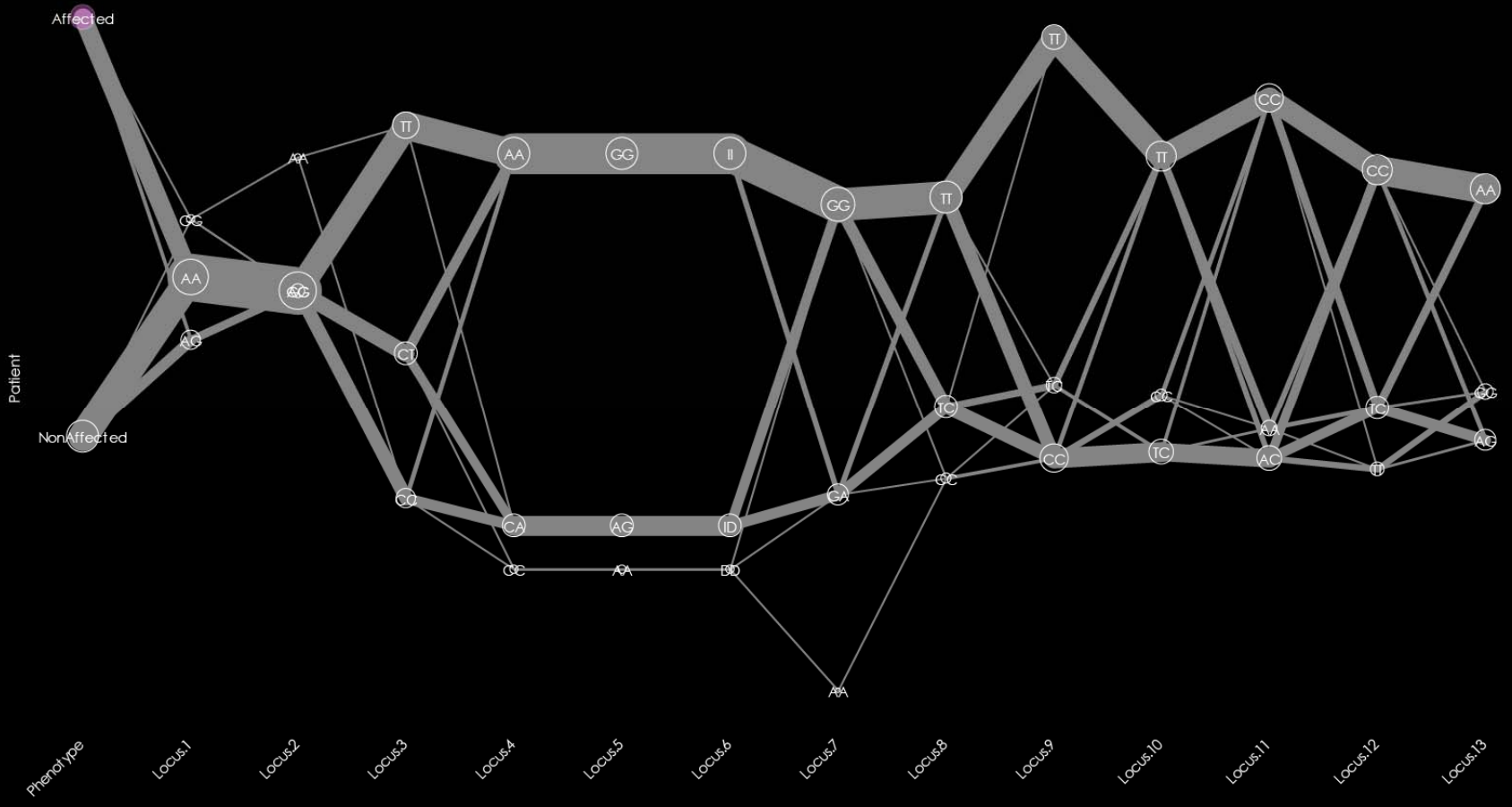
Order Data Vectors

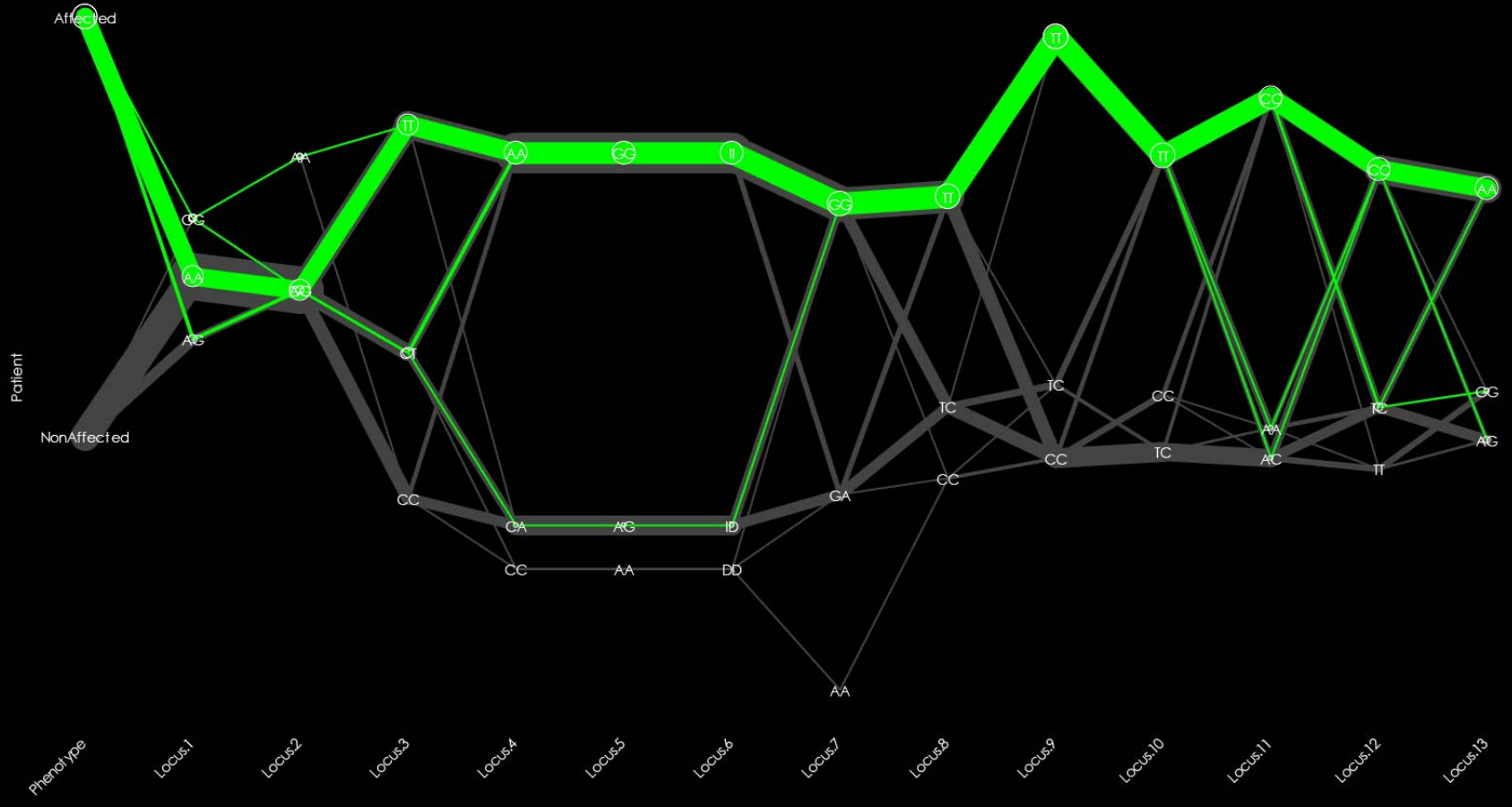
- a priori
- distance to ideal coordinates
- distance between coordinate vectors

Highlight Connecting Lines

Select Graphical Parameters

- color, transparency, size, etc.





Enhancements of Textile Plot

Transform

- single data vector (log, inverse, etc)
- multiple data vectors (normalisation, etc)

Substitute Data Point

Extract/Remove

- records
- data vectors

Re-weight

- outliers

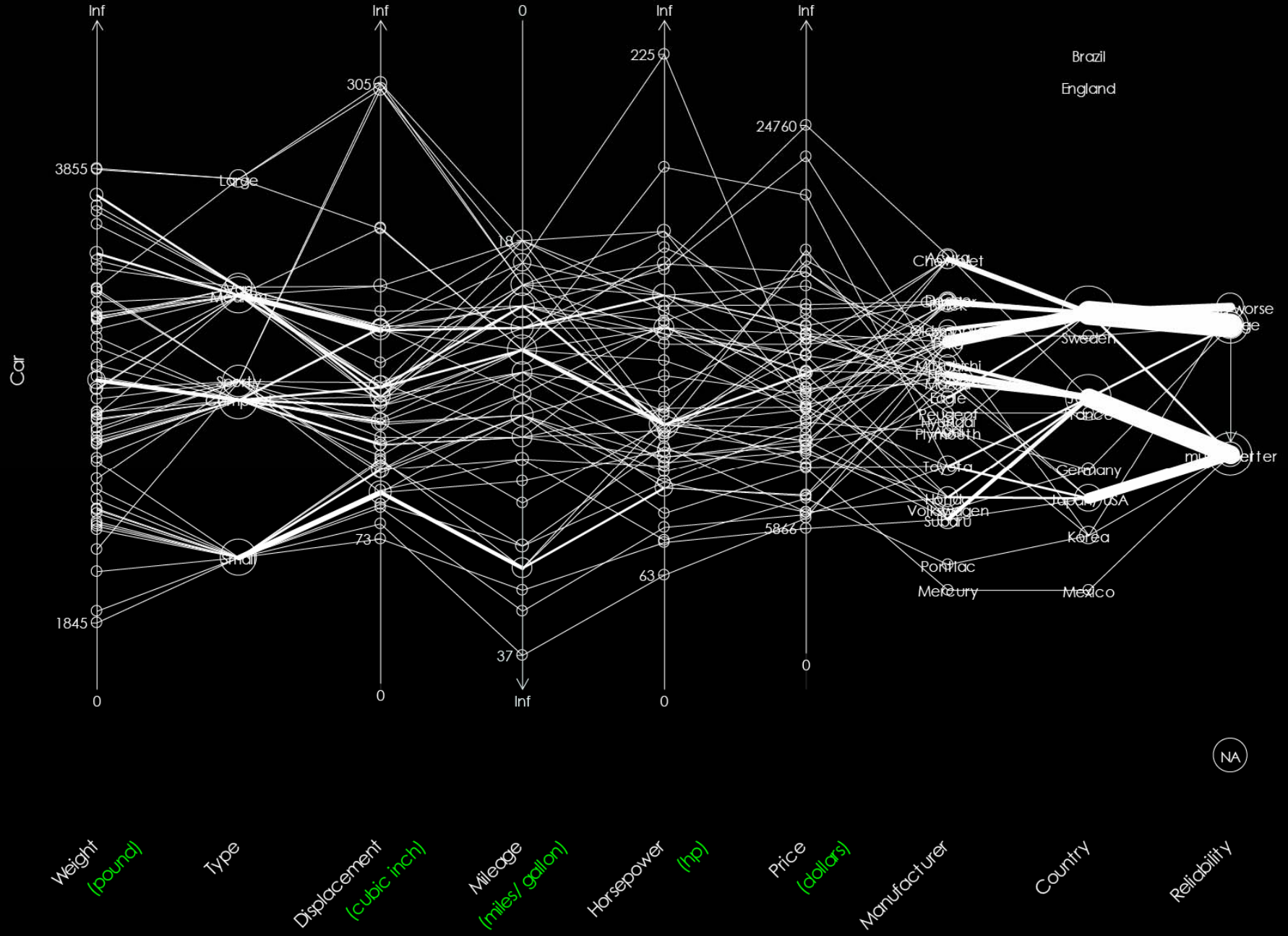
Order Data Vectors

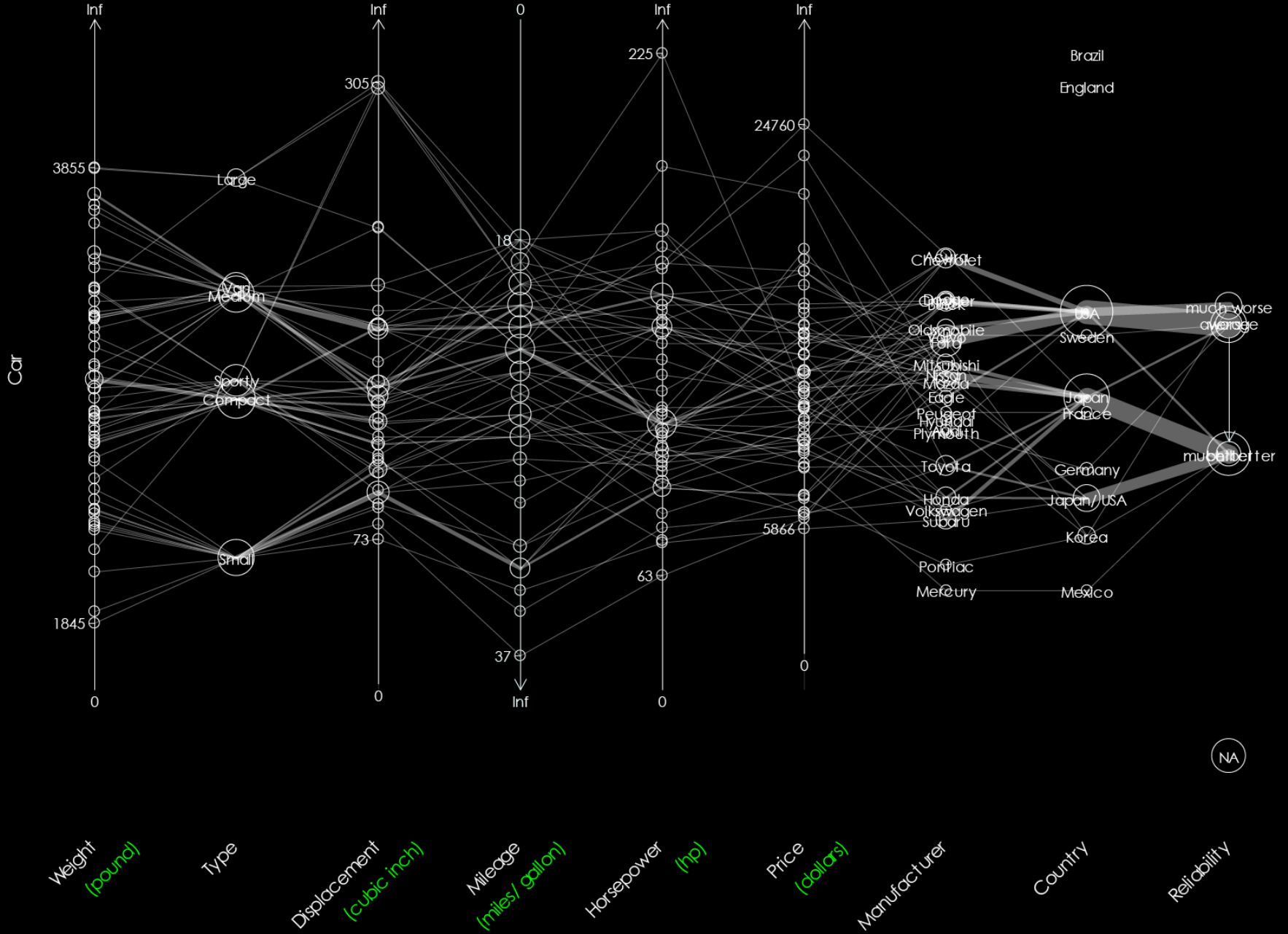
- a priori
- distance to ideal coordinates
- distance between coordinate vectors

Highlight Connecting Lines

Select Graphical Parameters

- color, transparency, size, etc.





Enhancements of Textile Plot

Transform

- single data vector (log, inverse, etc)
- multiple data vectors (normalisation, etc)

Substitute Data Point

Extract/Remove

- records
- data vectors

Re-weight

- outliers

Order Data Vectors

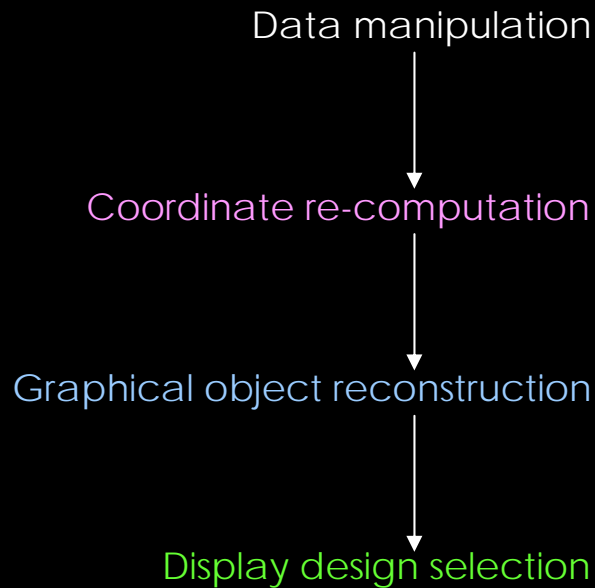
- a priori
- distance to ideal coordinates
- distance between coordinate vectors

Highlight Connecting Lines

Select Graphical Parameters

- color, transparency, size, etc.

Enhancements of Textile Plot



Transform

- single data vector (log, inverse, etc)
- multiple data vectors (normalisation, etc)

Substitute Data Point

Extract/Remove

- records
- data vectors

Re-weight

- outliers

Order Data Vectors

- a priori
- distance to ideal coordinates
- distance between coordinate vectors

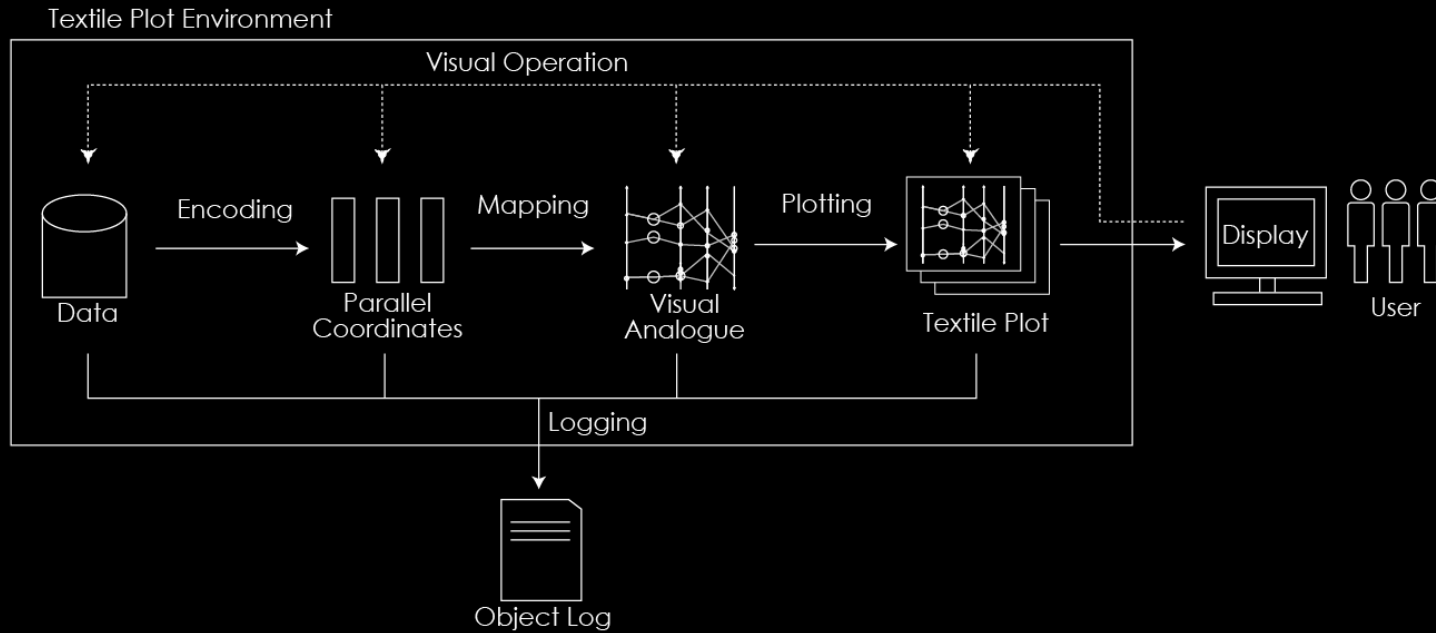
Highlight Connecting Lines

Select Graphical Parameters

- color, transparency, size, etc.

Reference Model

Kumasaka and Shibata 2007



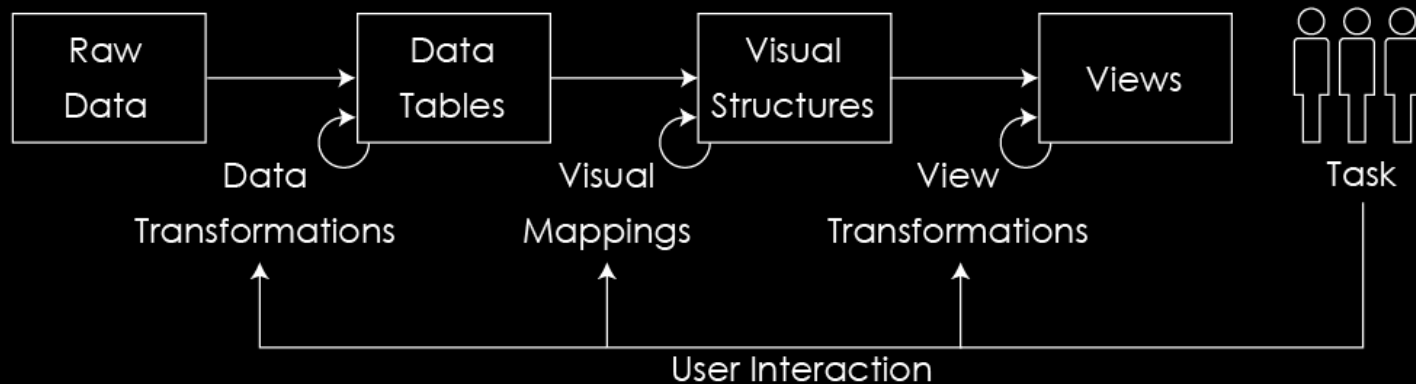
Sequence of objects

Visual operation

Object log

Reference Model for Information Visualization

Card, et al. 1999

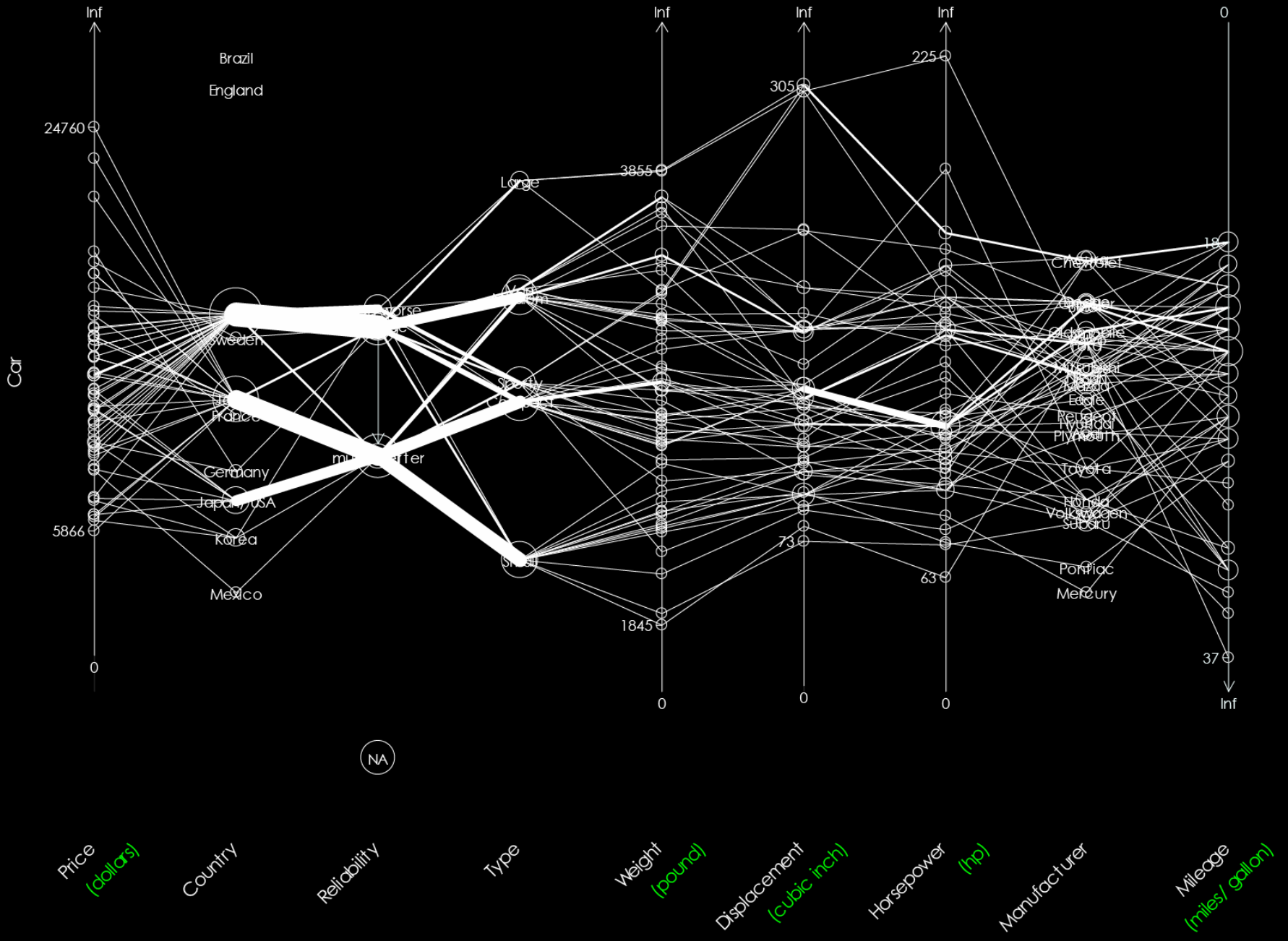


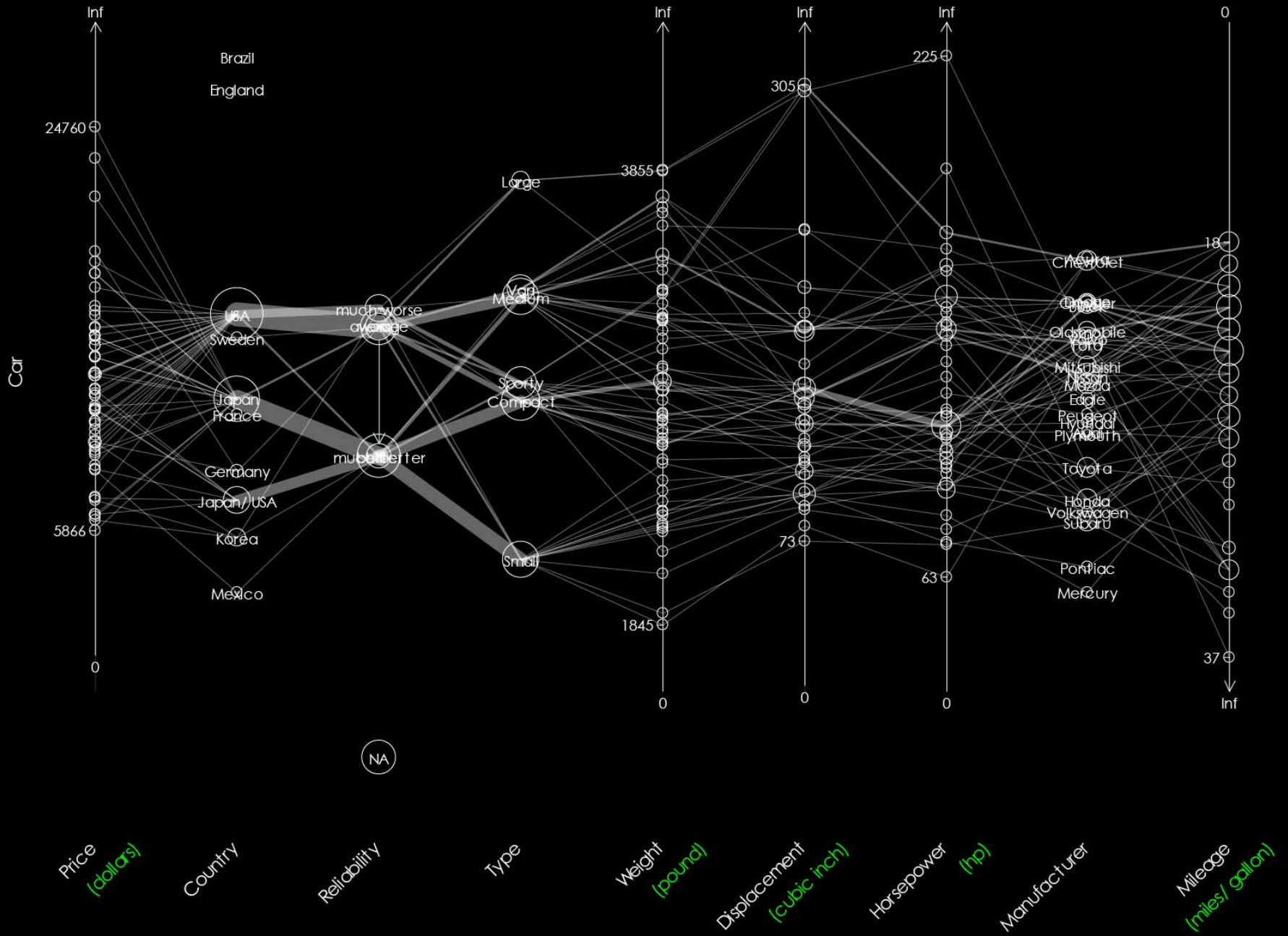
Raw Data idiosyncratic formats

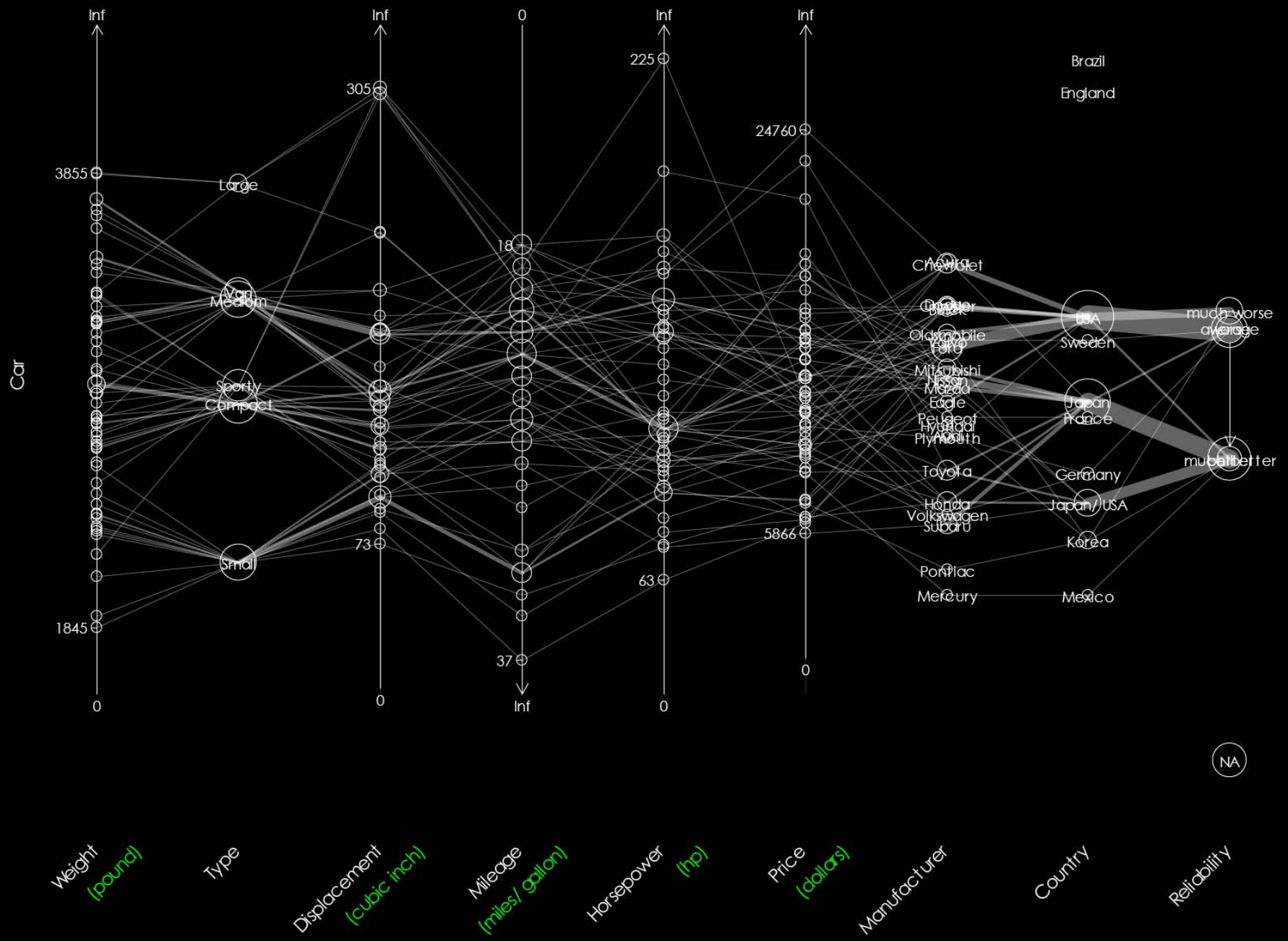
Data Tables relations (cases by variables) + metadata

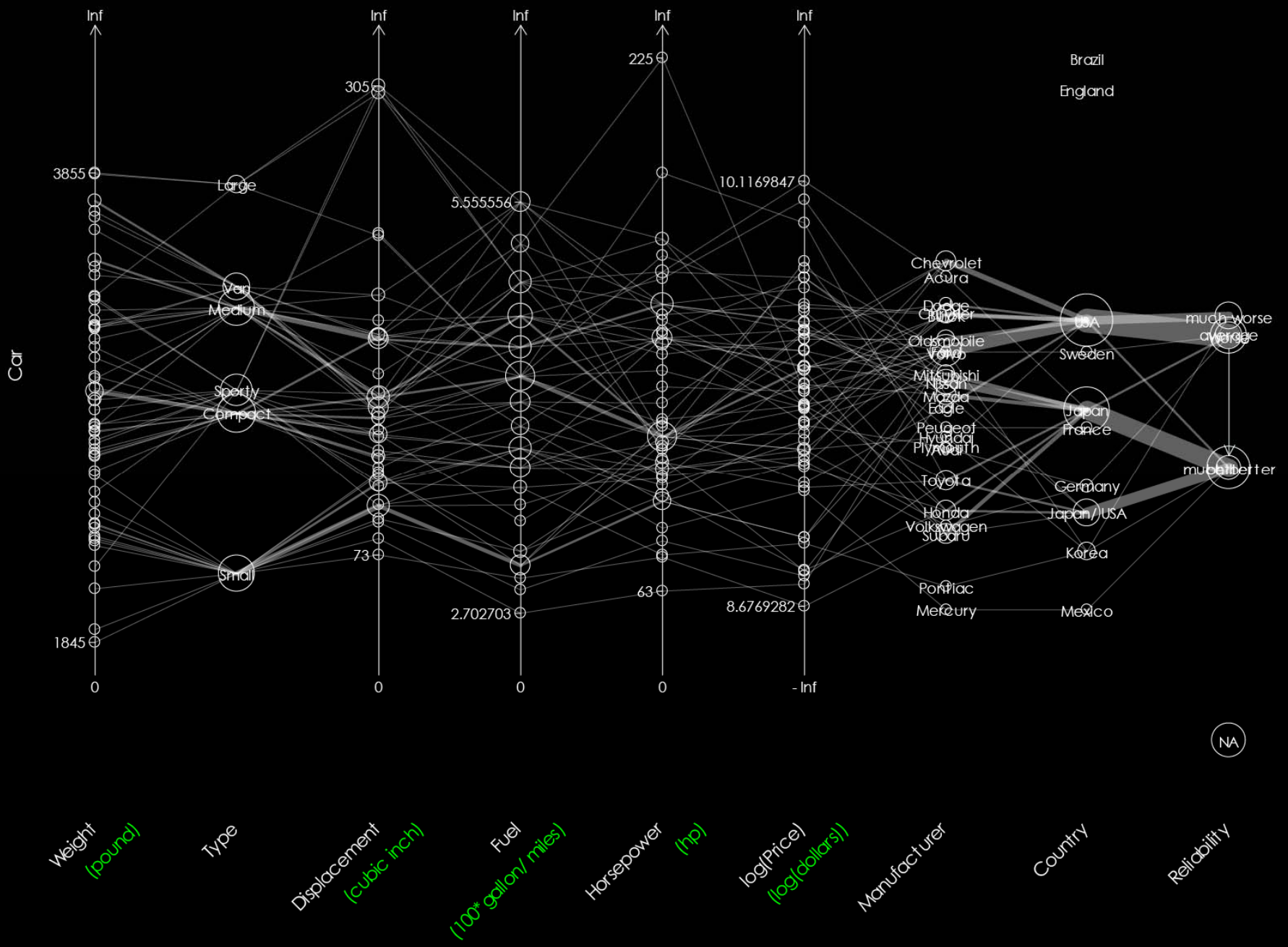
Visual Structures spatial substrates + marks + graphical properties

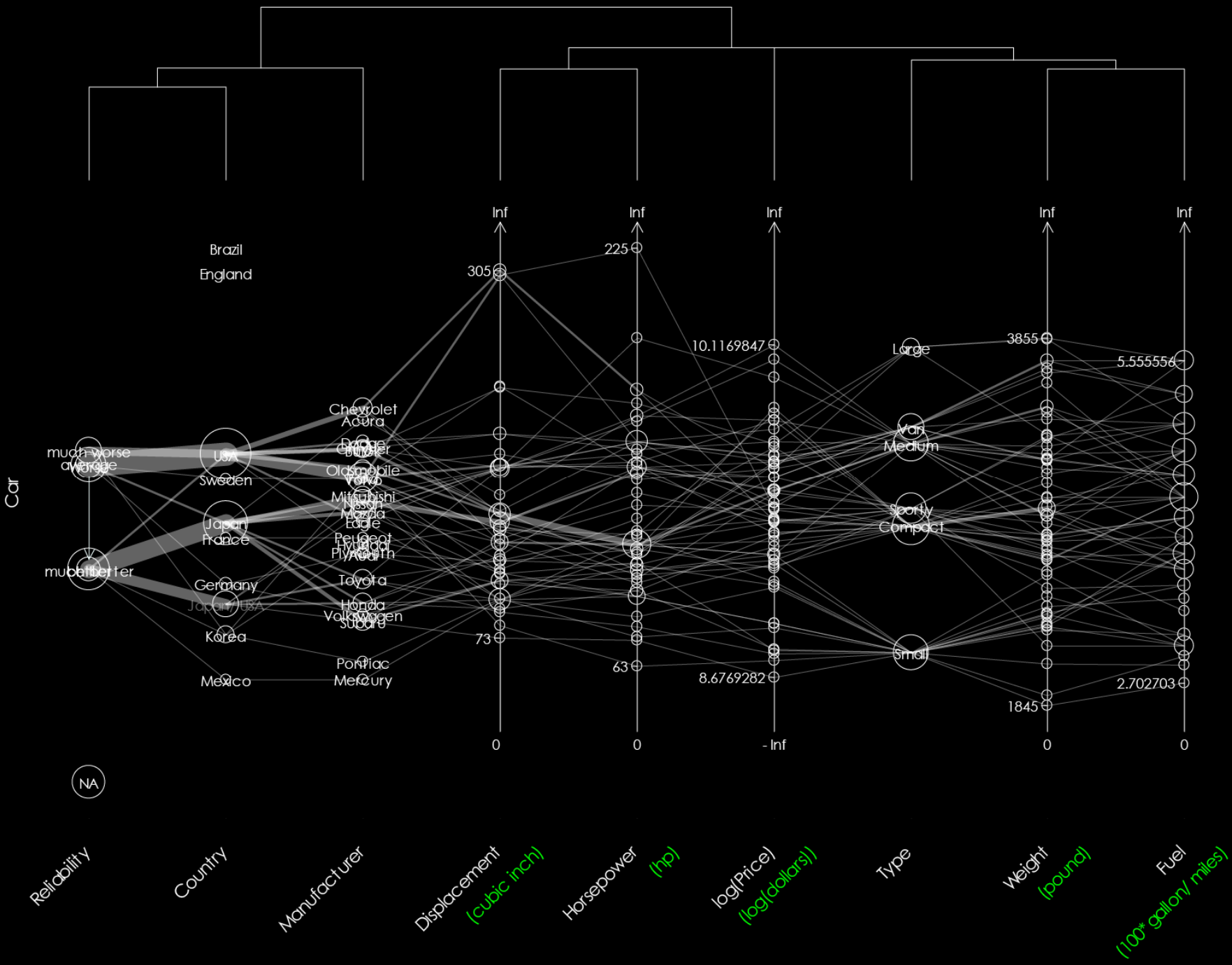
Views graphical parameters (positions, scaling, clipping,...)

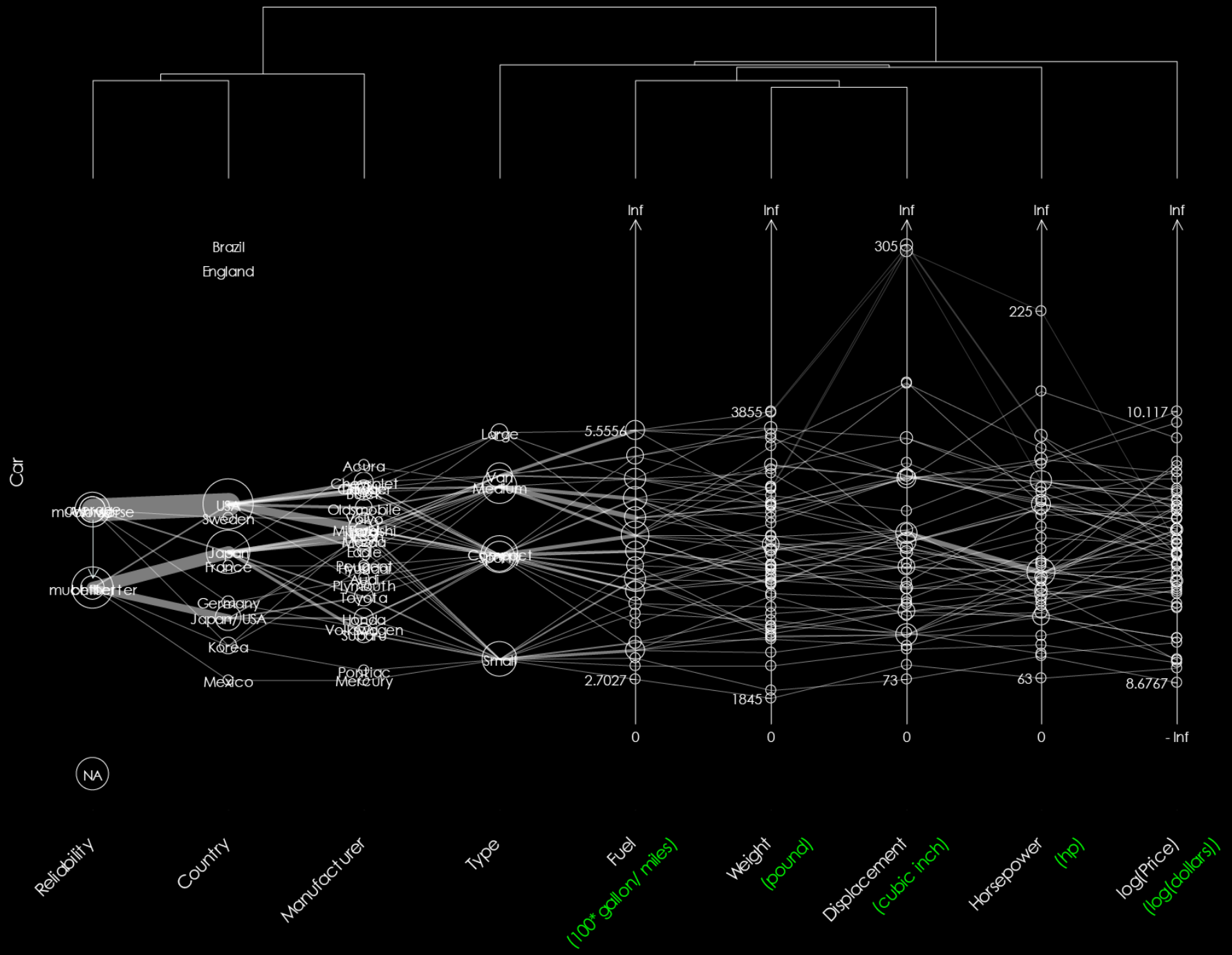


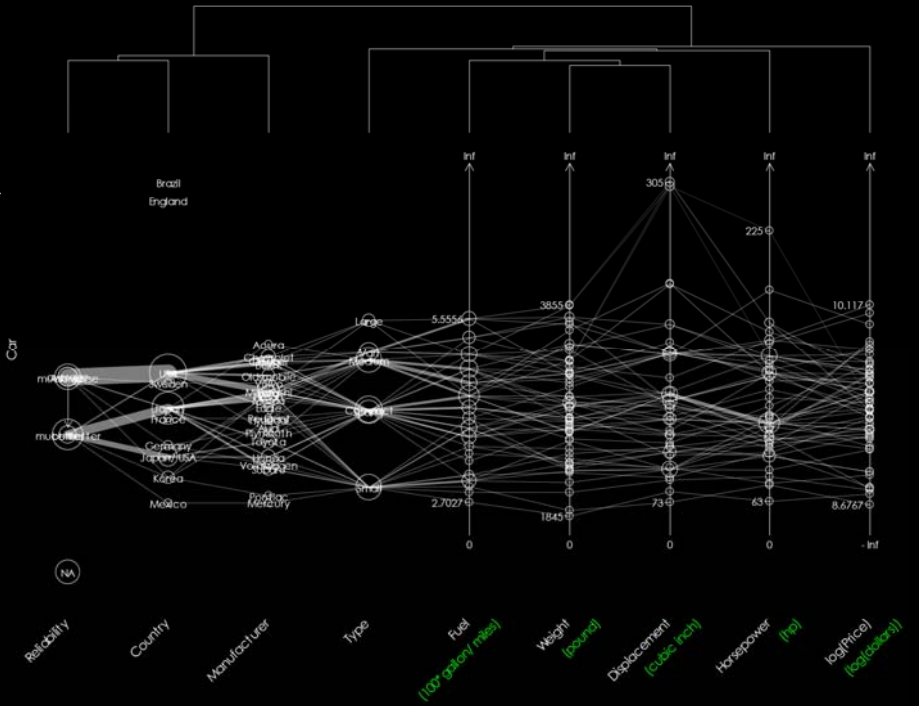
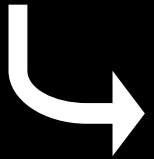
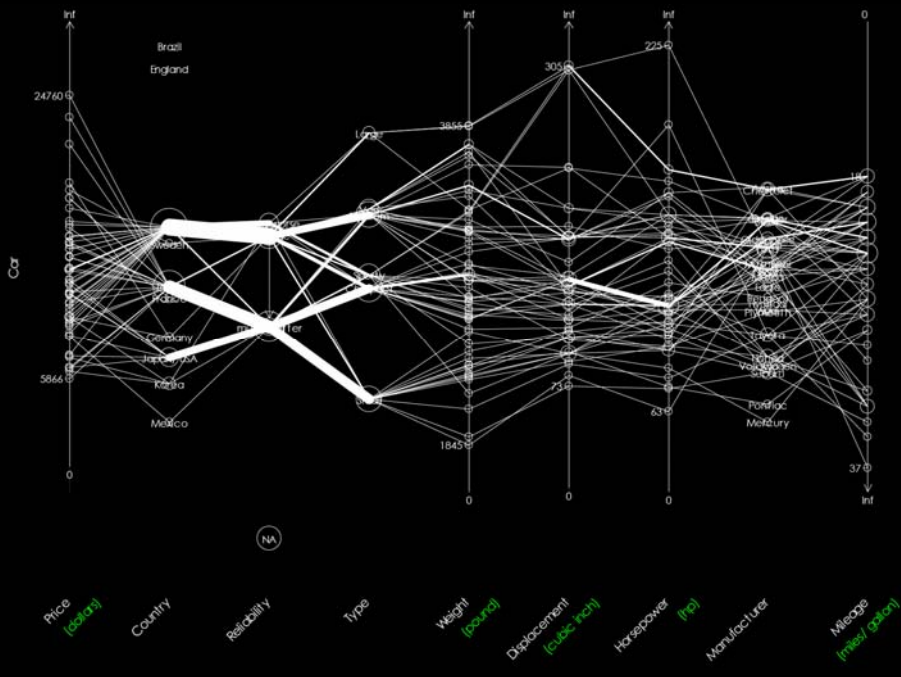




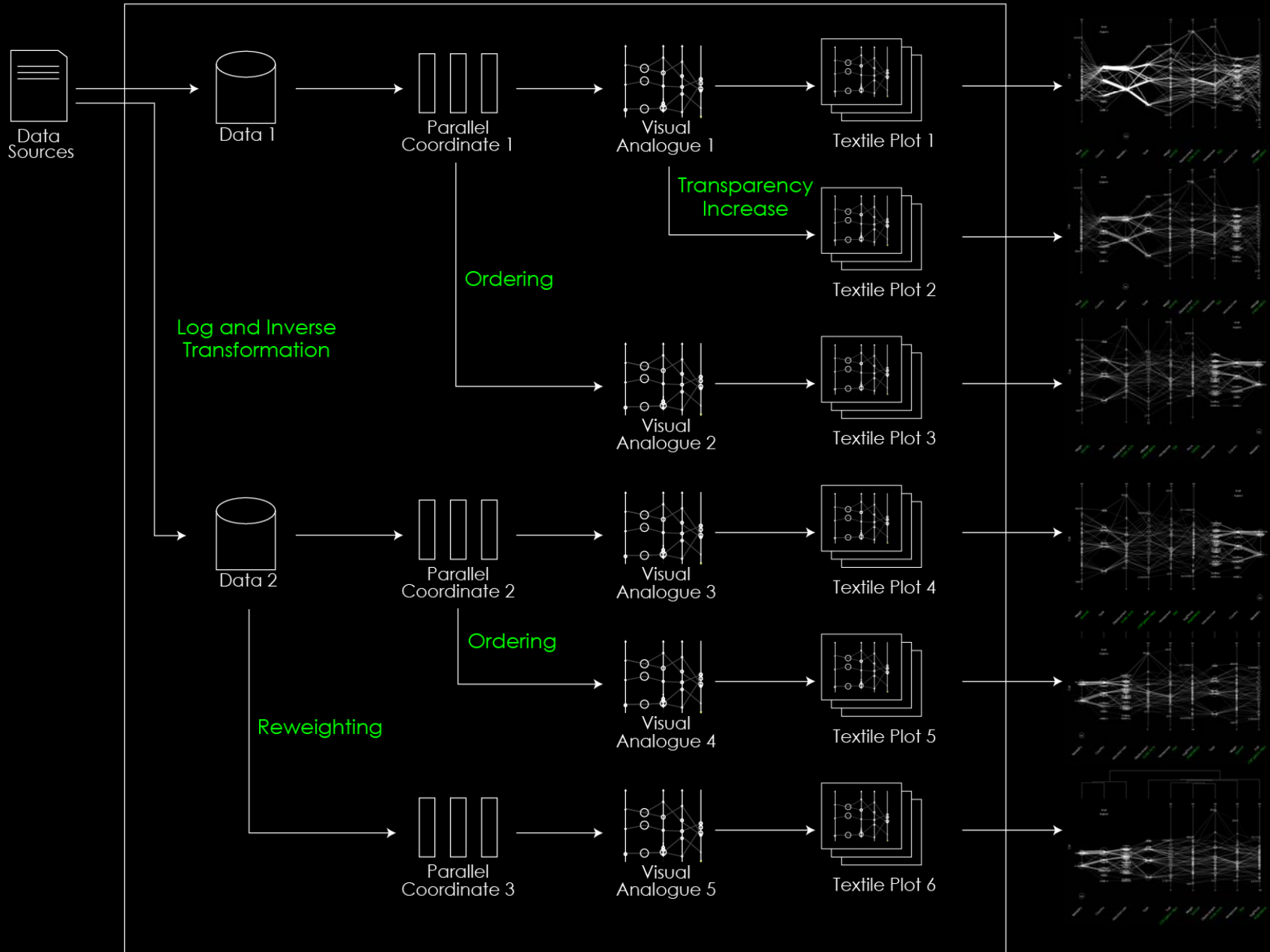








Automobile Data Analysis Log



Bibliography 1

Card, S. K., et al. (1999), *Information Visualization*, Morgan Kaufmann Pub.

Hurley, Catherine B. (2004), Clustering Visualisations of Multidimensional Data, *Journal of Computational and Graphical Statistics*, **13** 788--806.

Inselberg, A. (1985), The plane with parallel coordinates, *The Visual Computer* **1** 69--91.

Iwasaki, N. et al. (2005), Genetic Variants in the calpain-10 gene and the development of type 2 diabetes in the Japanese population. *Journal of Human Genetics* **50** 92--98.

Kumasaka, N. and Shibata, R., High Dimensional Data Visualisation: the Textile Plot, *Computational Statistics & Data Analysis*, submitted.

Bibliography 2

Kumasaka, N. and Shibata, R. (2006), Implementation of Textile Plot, *Proceedings in COMPSTAT 2006* Roma.

Kumasaka, N. and Shibata, R. (2007), The Textile Plot Environment, *統計数理特集号*, in Japanese.

Wegman, E. (1990), Hyperdimensional data analysis using parallel coordinates. *Journal of The American Statistical Association* **85** 664--675.