



Removing the blindfold

Visualising
statistical
models

Hadley Wickham

Assistant Professor
Dobelman Family Junior Chair
Department of Statistics
Rice University

Why?

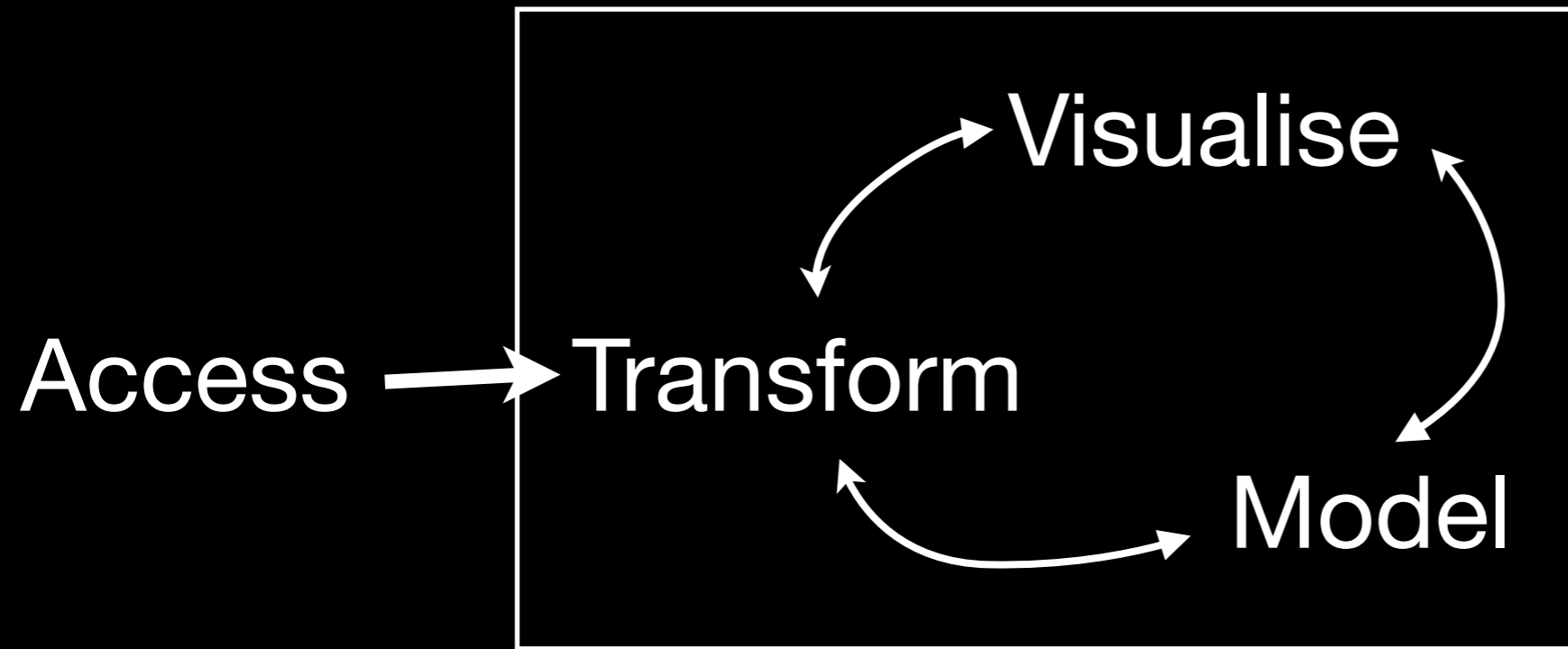
Access

Understand

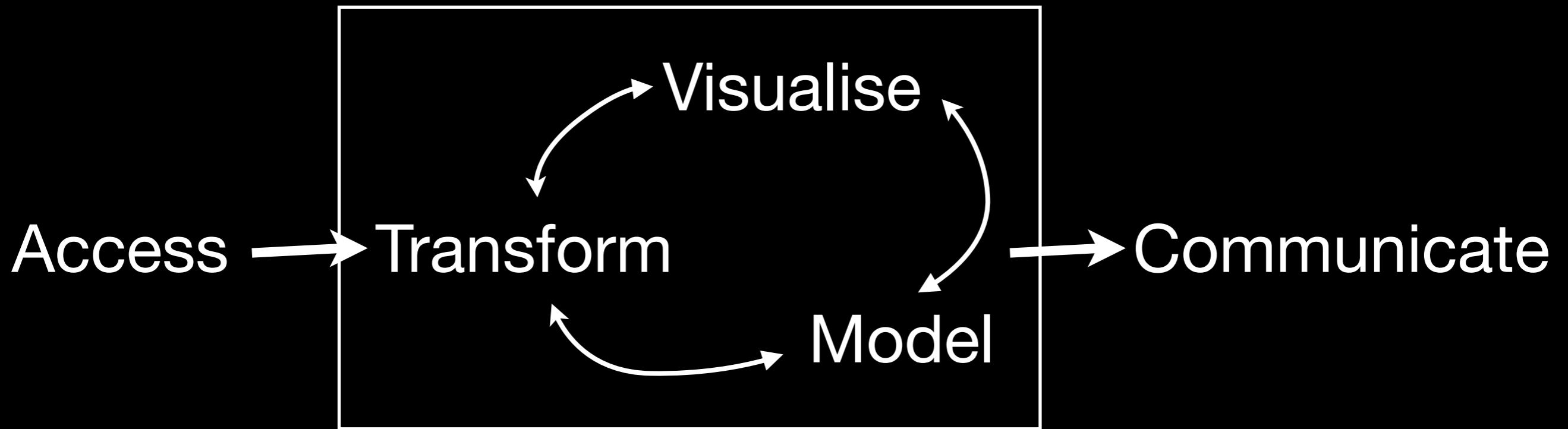
Access →



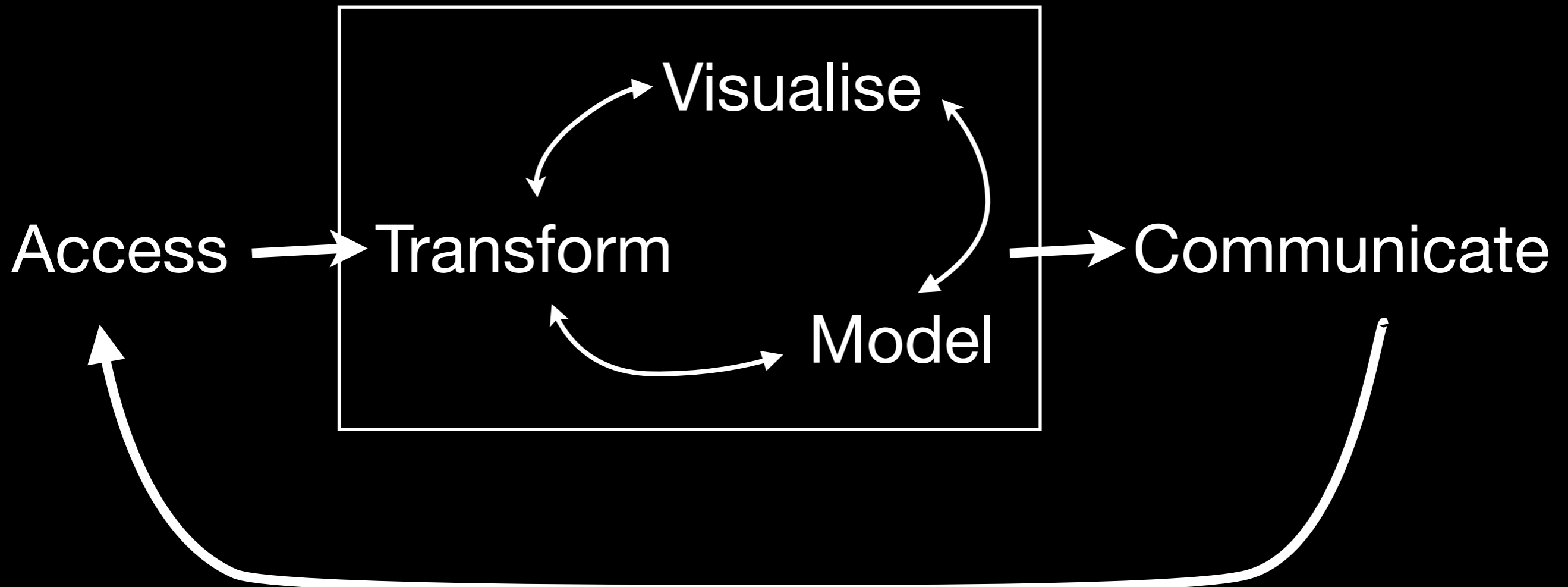
Understand



Understand

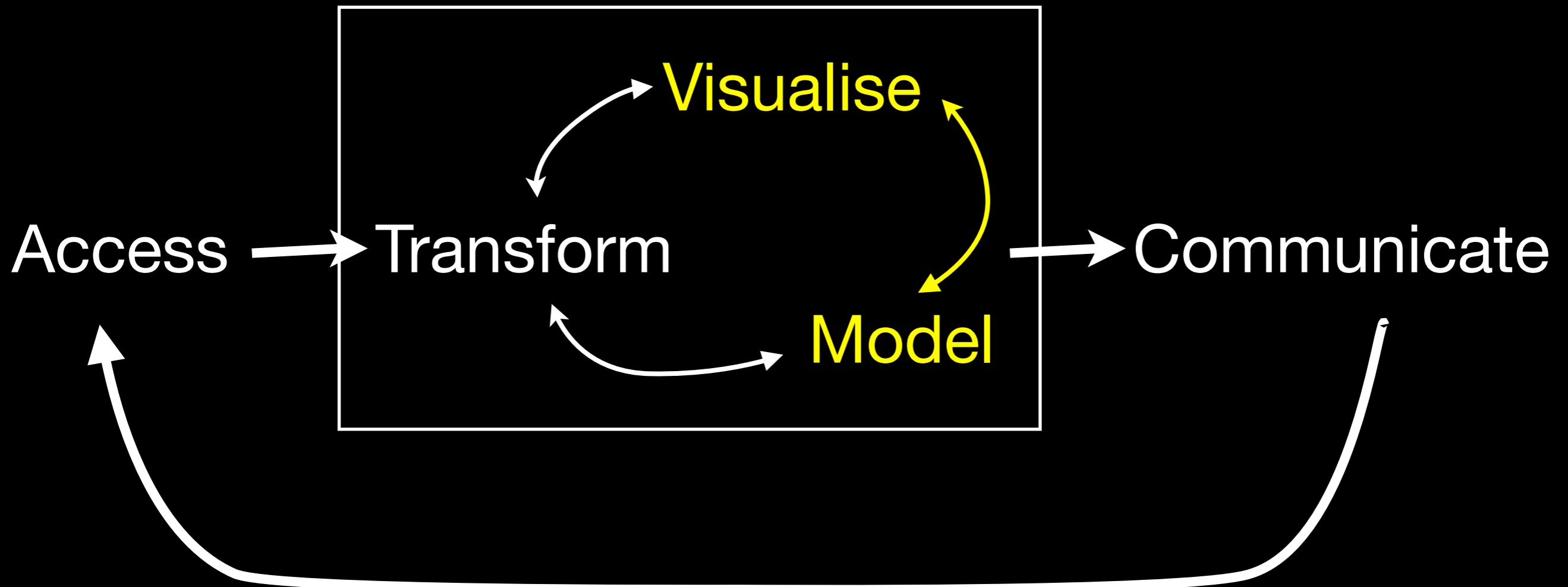


Understand



Visualisation	Model
<ul style="list-style-type: none">+ Uncovers the unexpected- Slow- Cognitive biases	<ul style="list-style-type: none">+ Mathematically well founded+ Fast- Only discovers what we anticipate

Understand



Neural networks

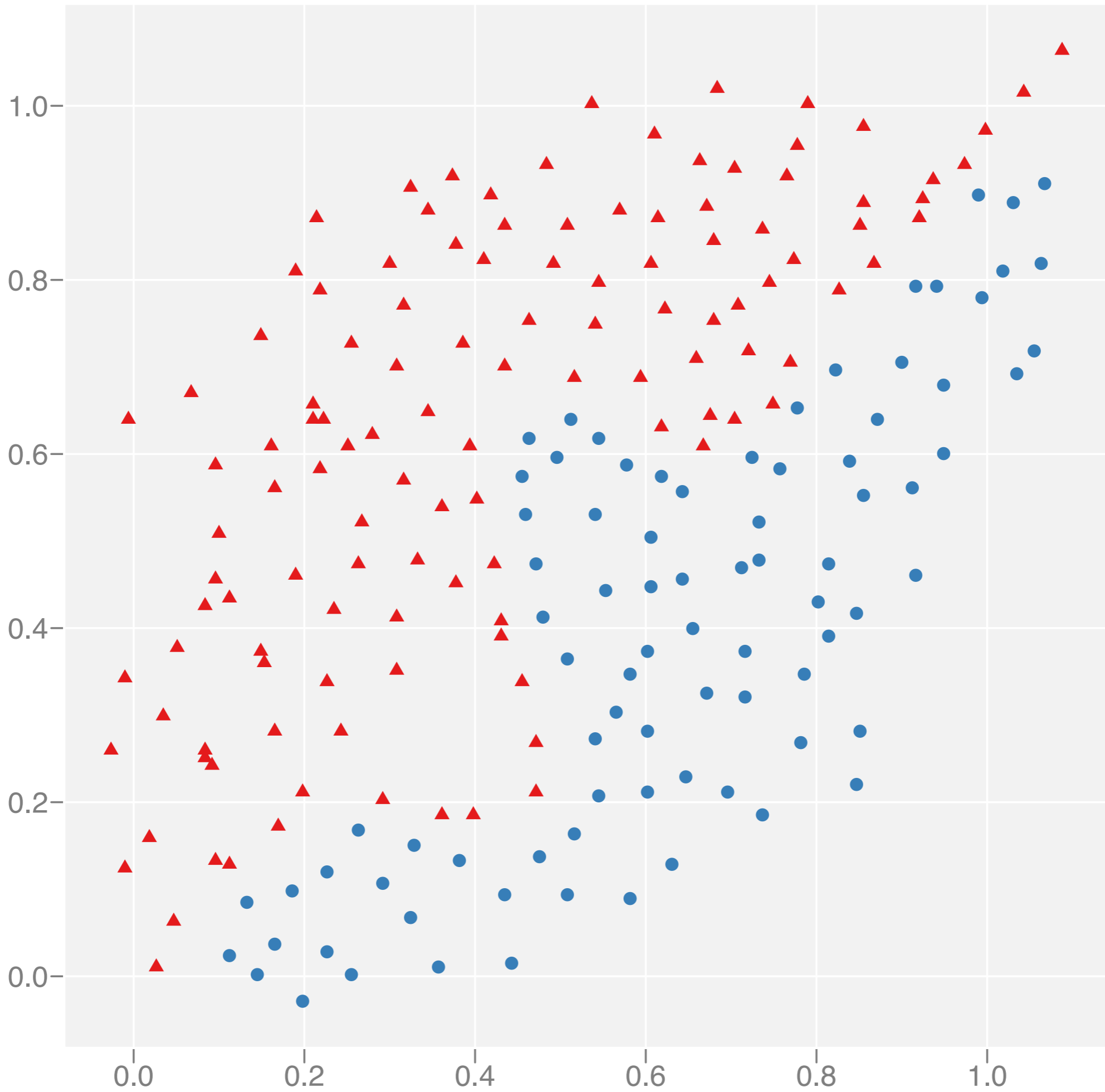
Display the model
in the data space

Look at many
members of a collection

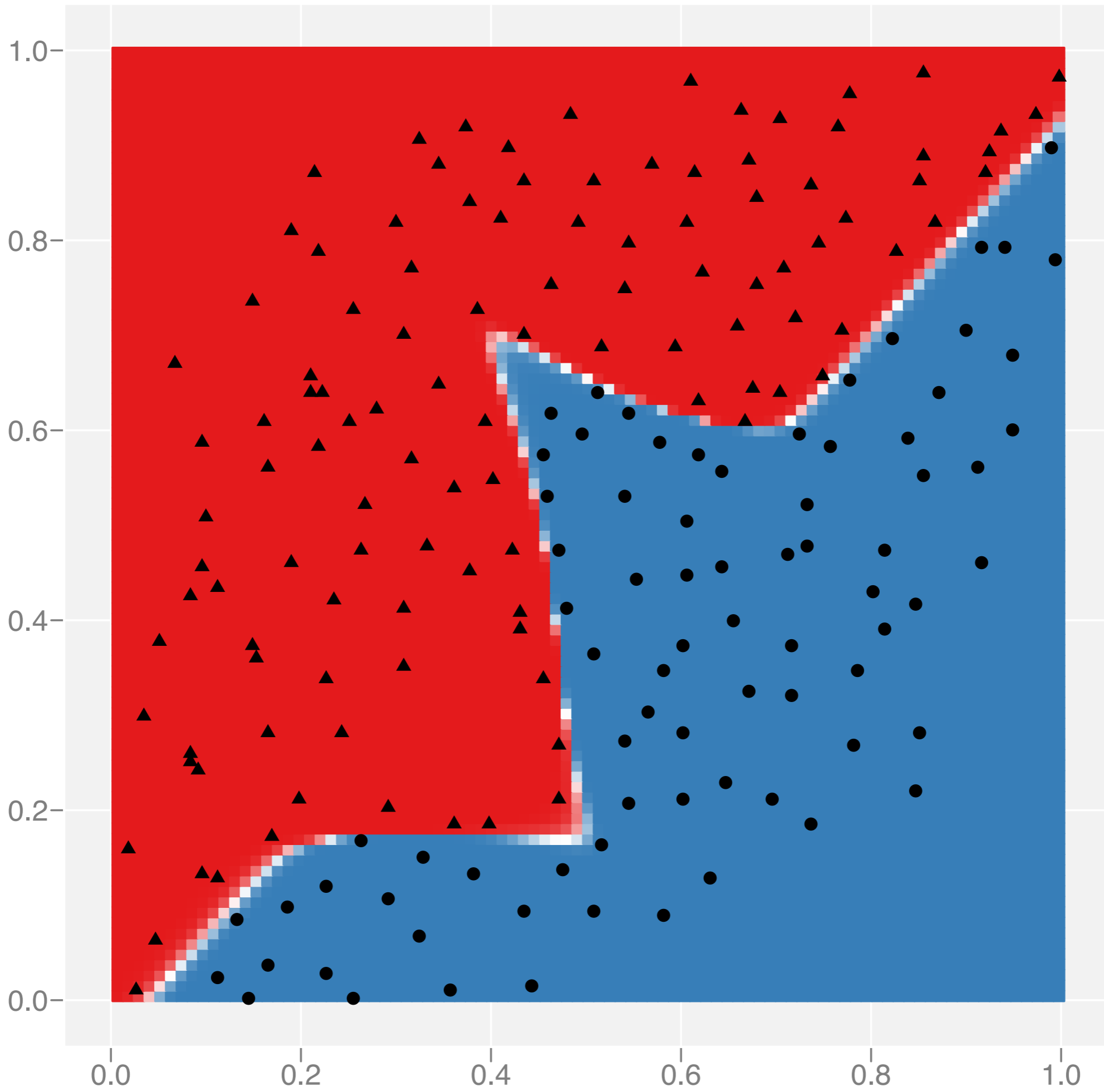
Explore the process of
fitting, not just the end result

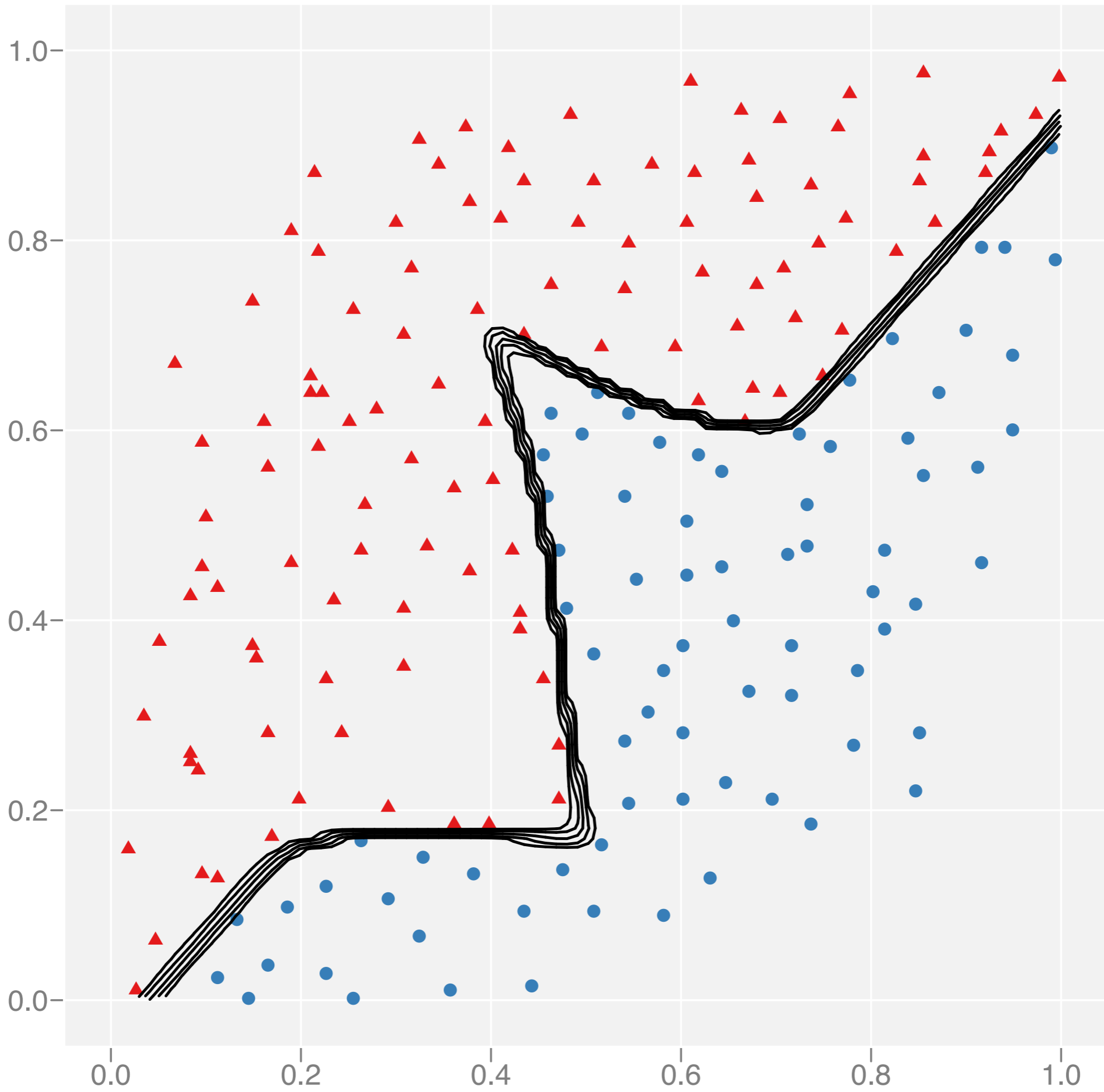
Neural networks

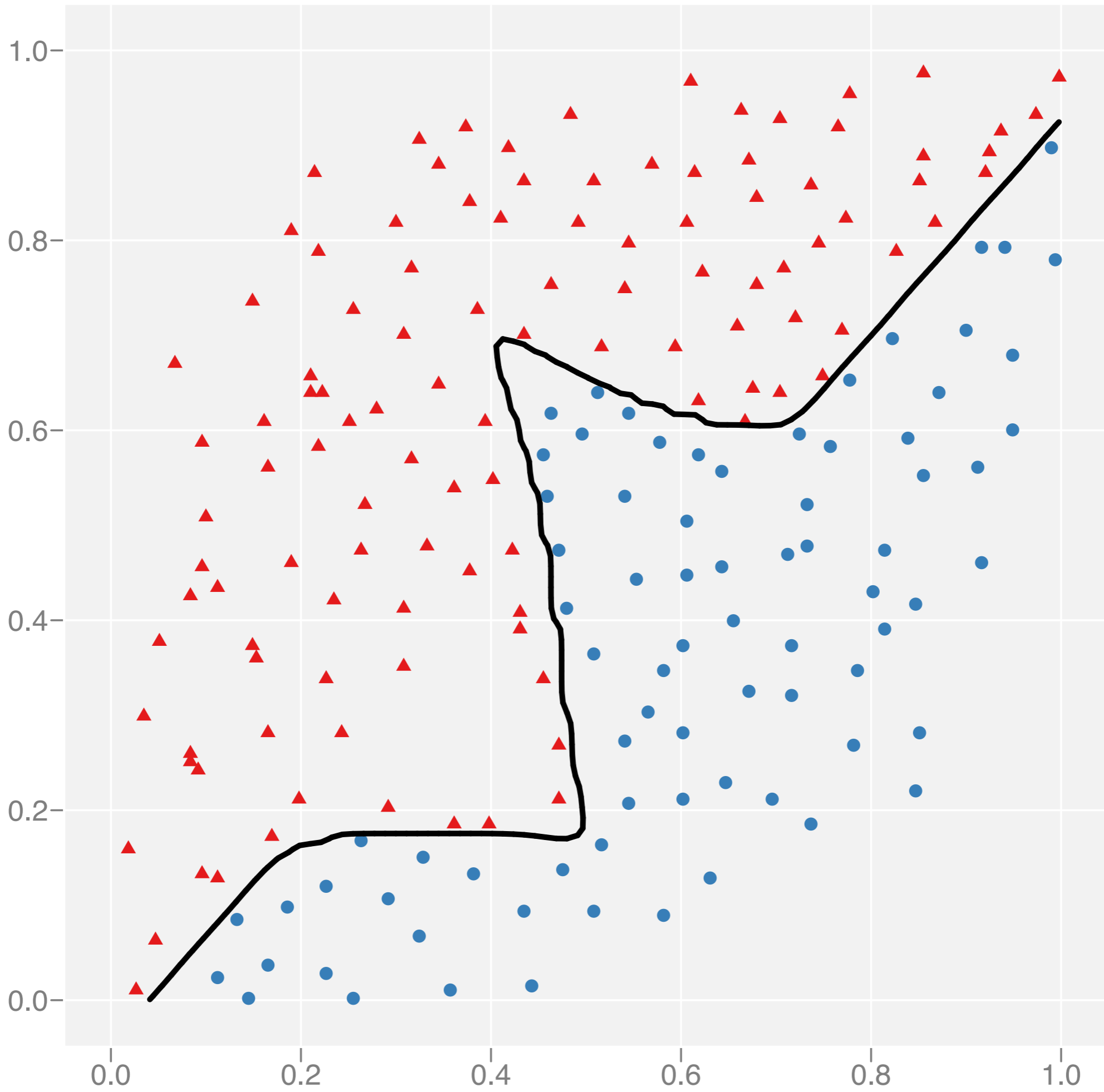
- Modelled on the way that brains work
- Normally treated as a black box. Can we gain more insight into how they work?
- Single hidden-layer neural network:
nnet R package



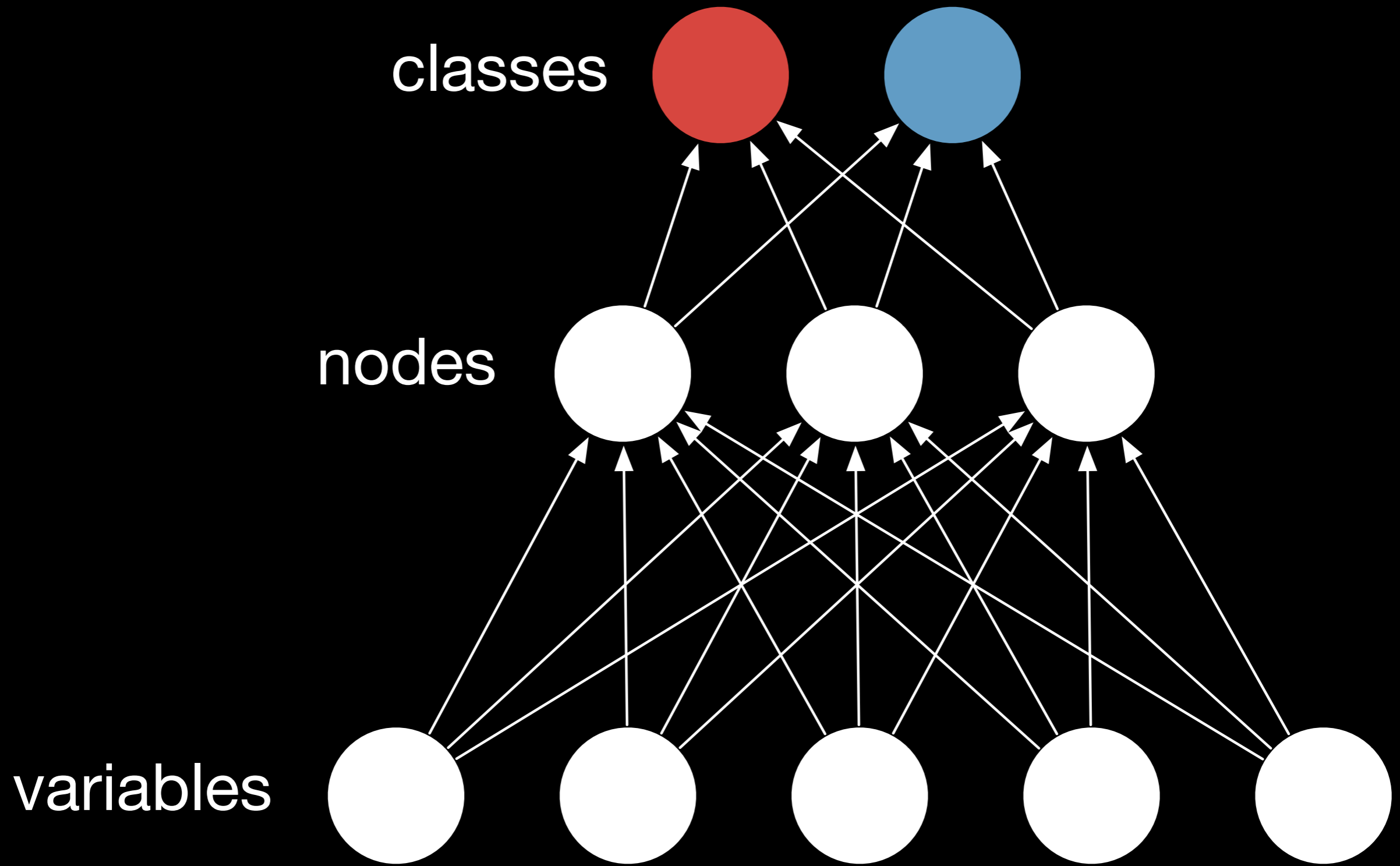
Display the model in
data space



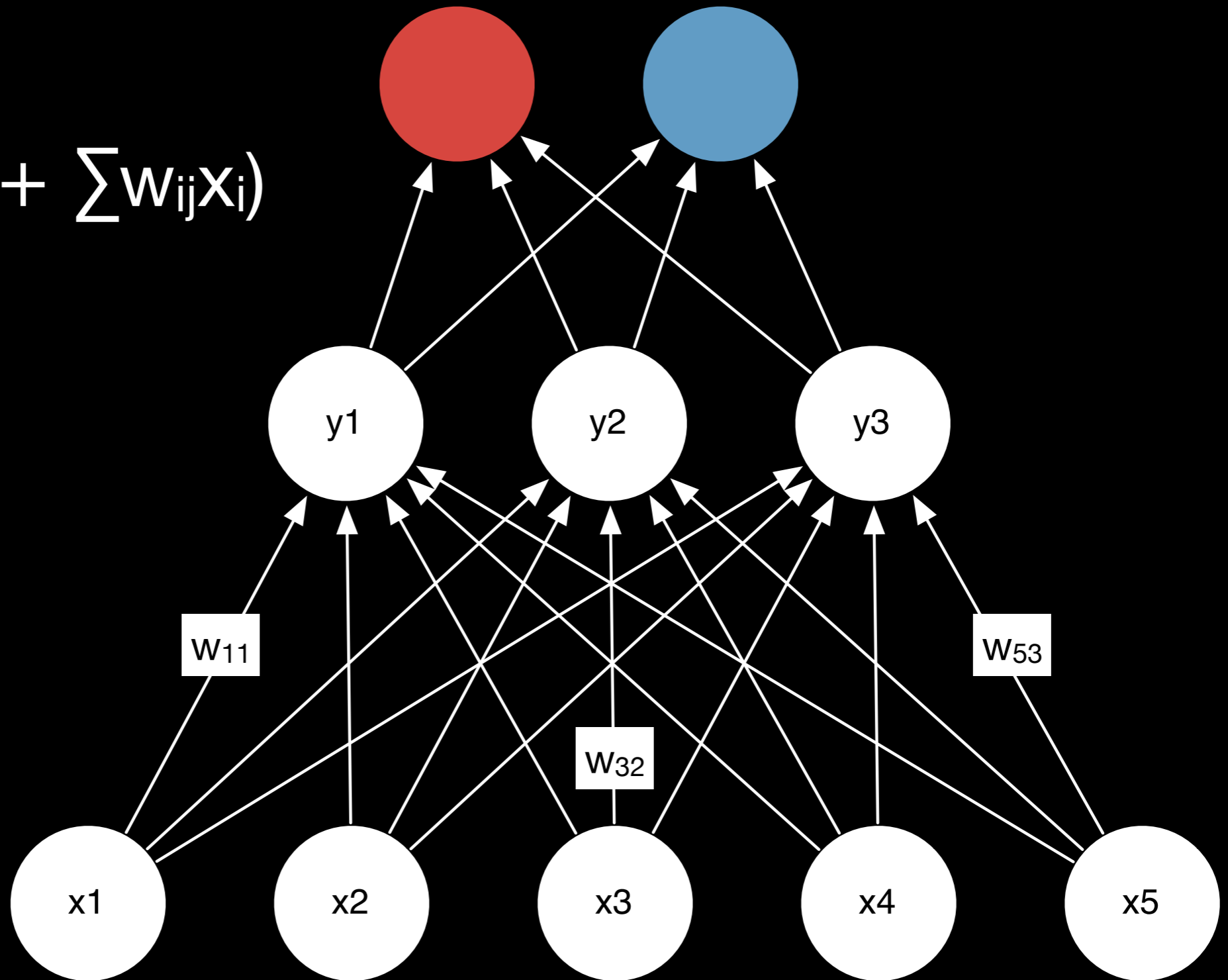




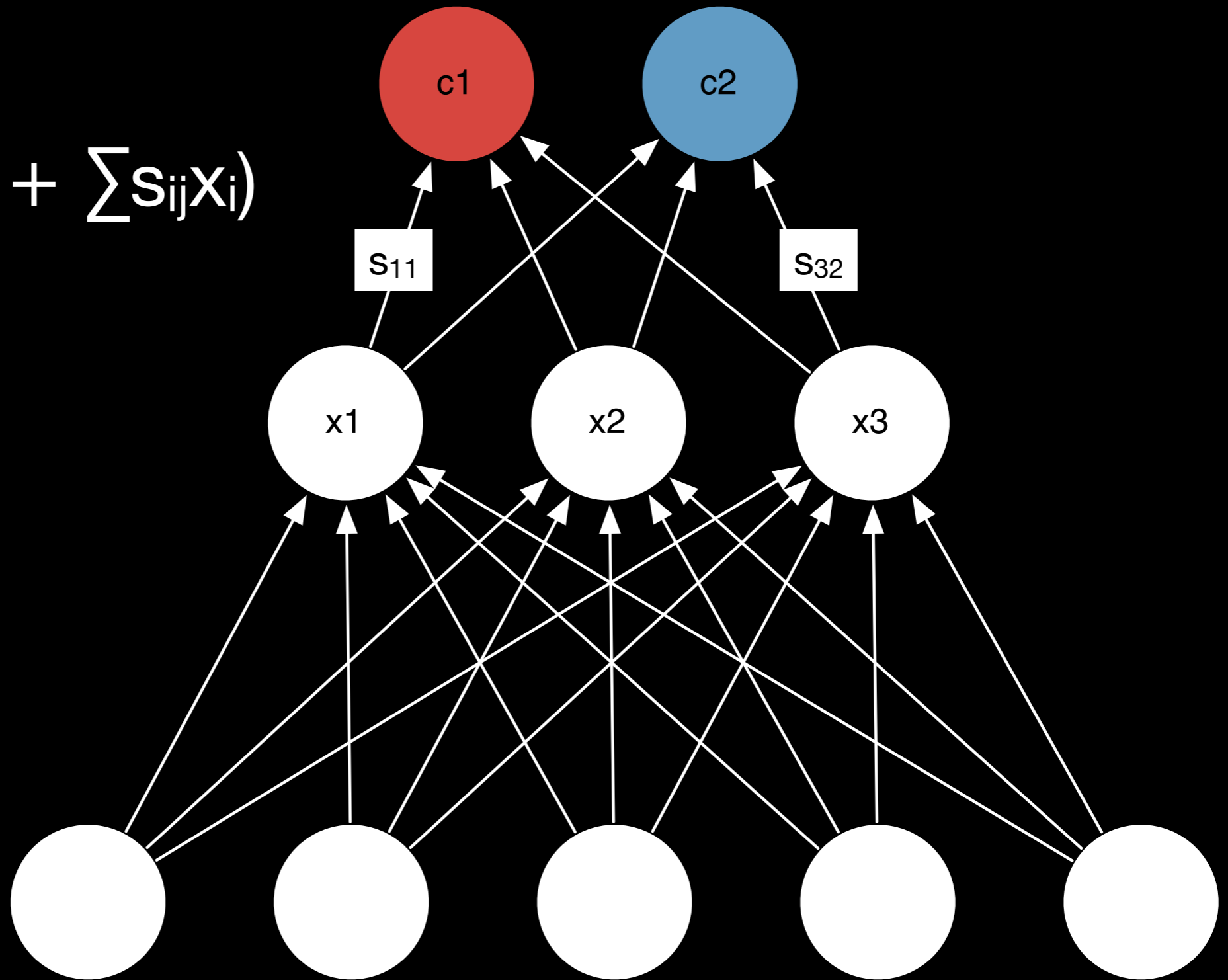
How do neural
networks work?

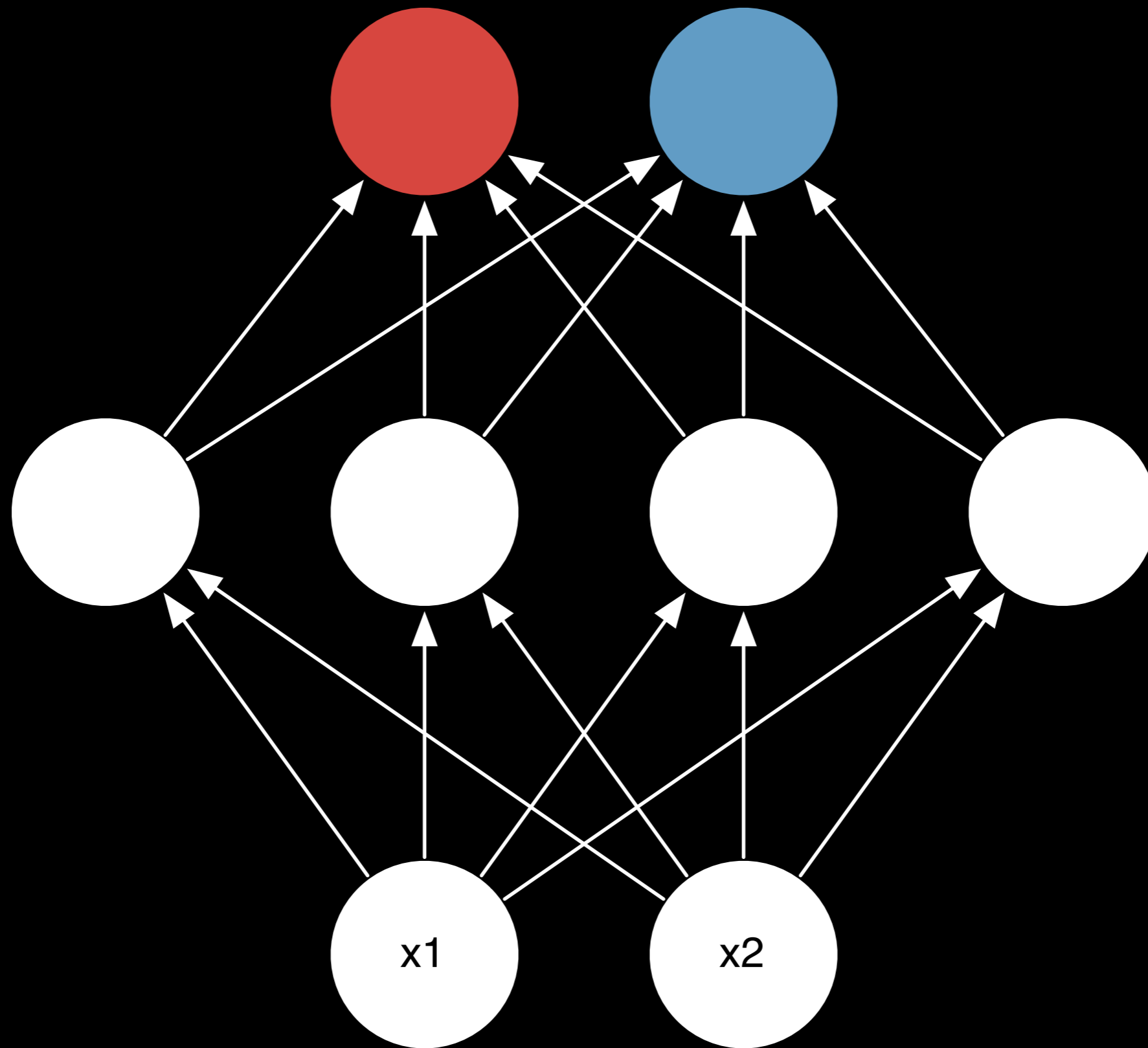


$$y_j = \text{logit}(\alpha_j + \sum w_{ij}x_i)$$

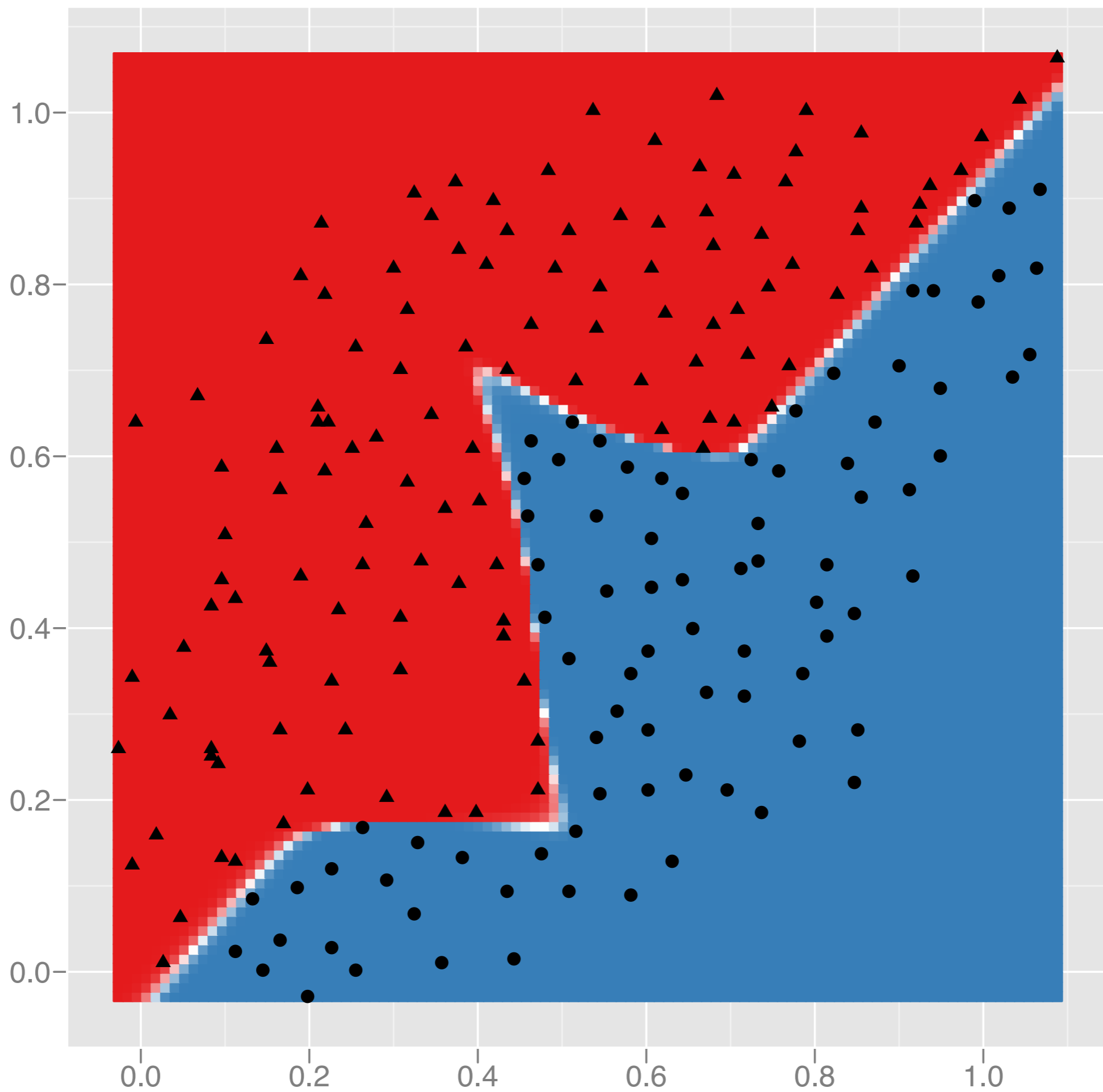


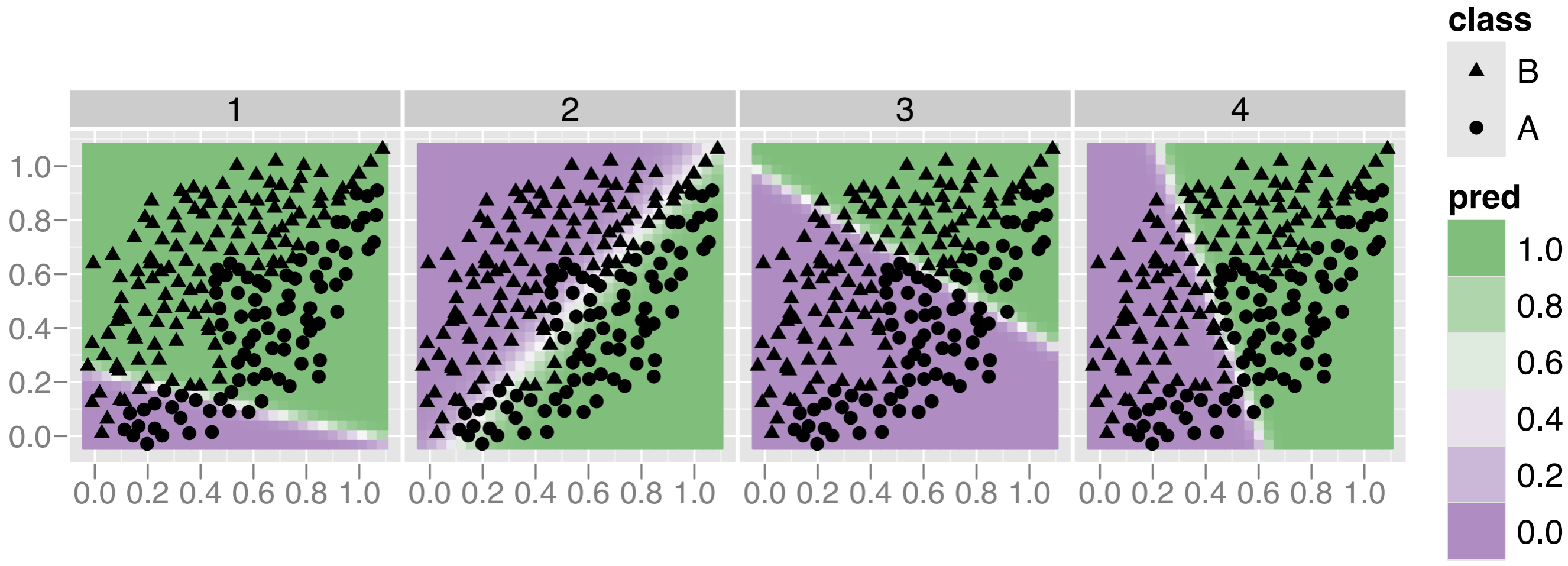
$$c_j = \text{logit}(\alpha_j + \sum s_{ij}x_i)$$

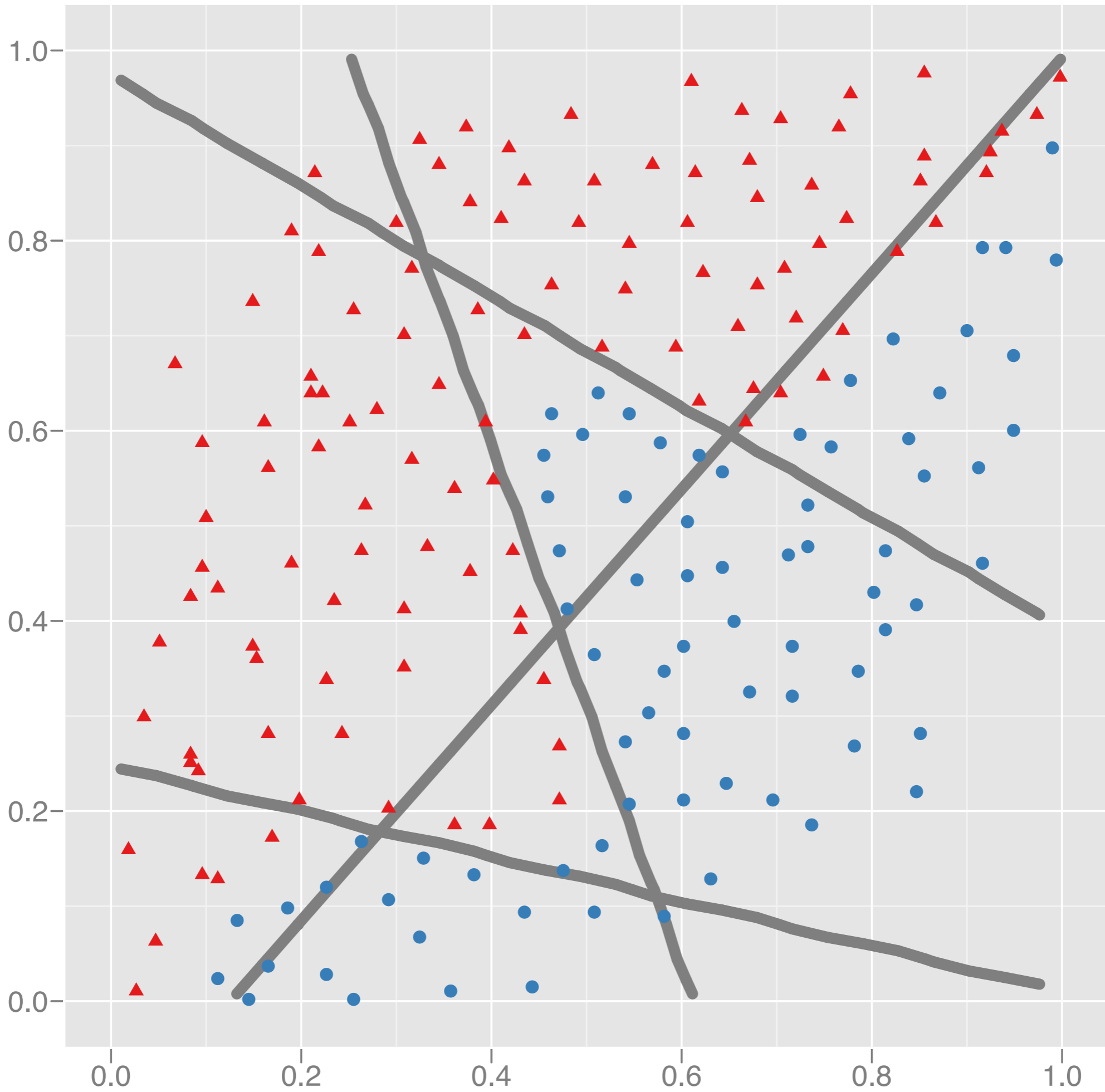


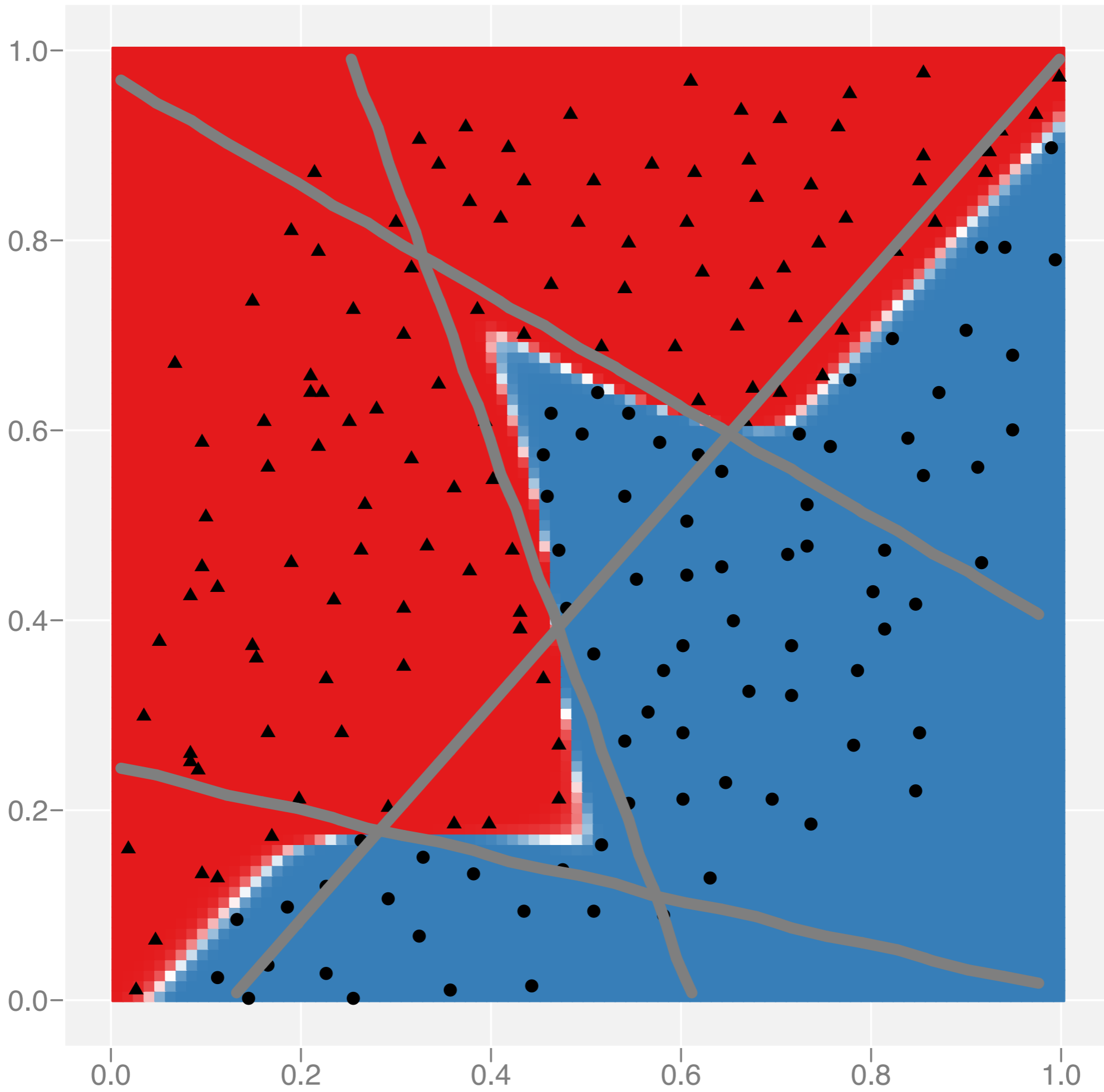


Look at all members
of the collection

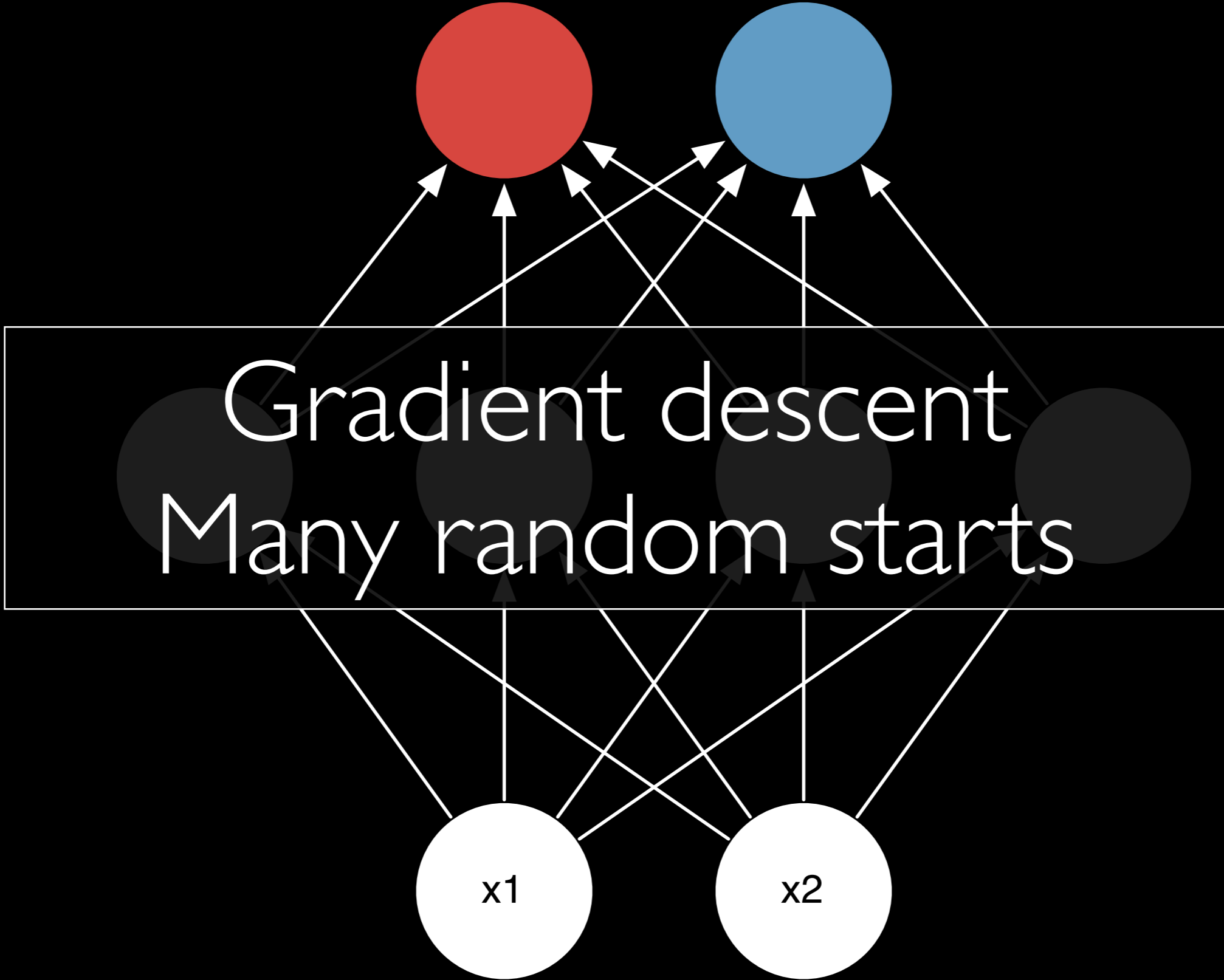


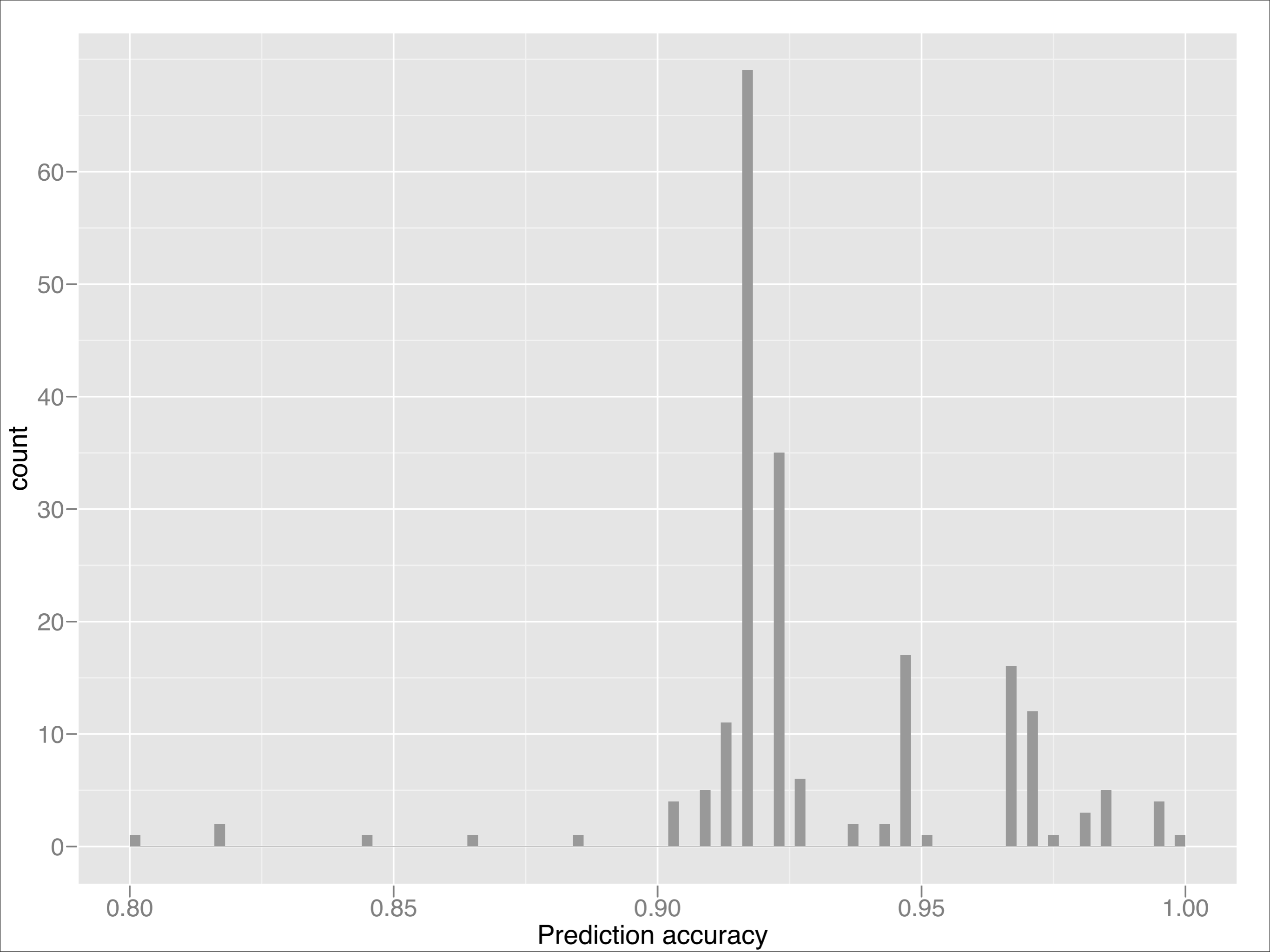


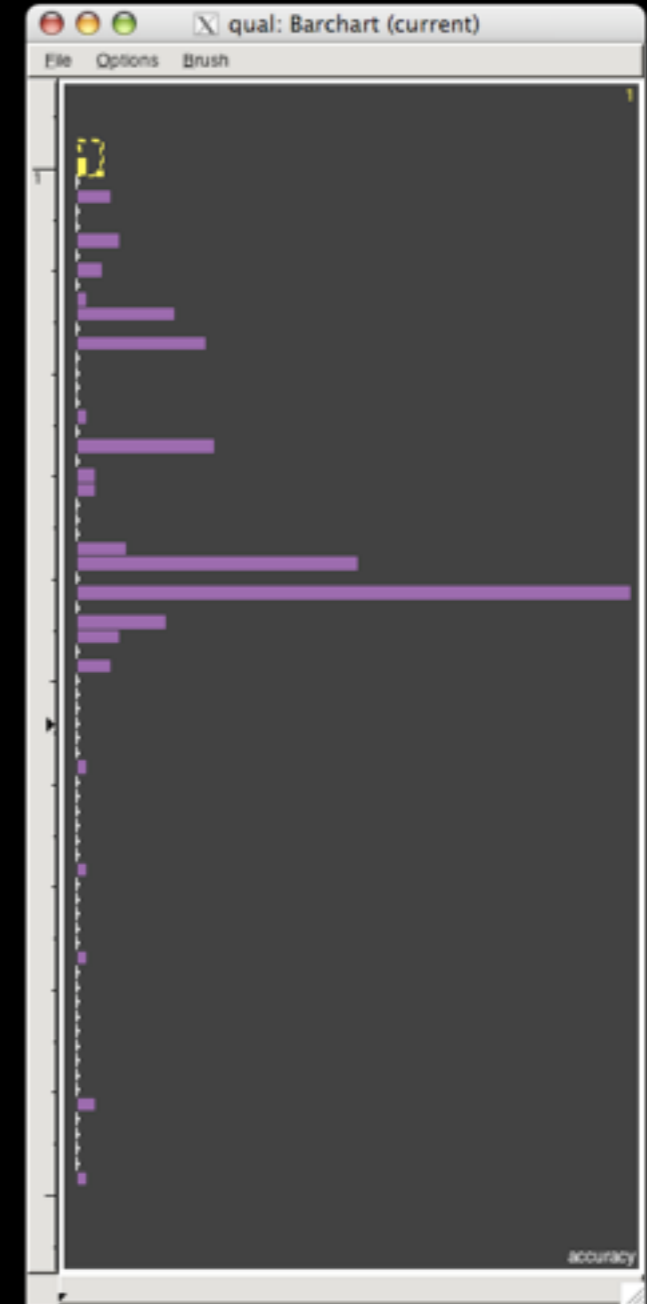
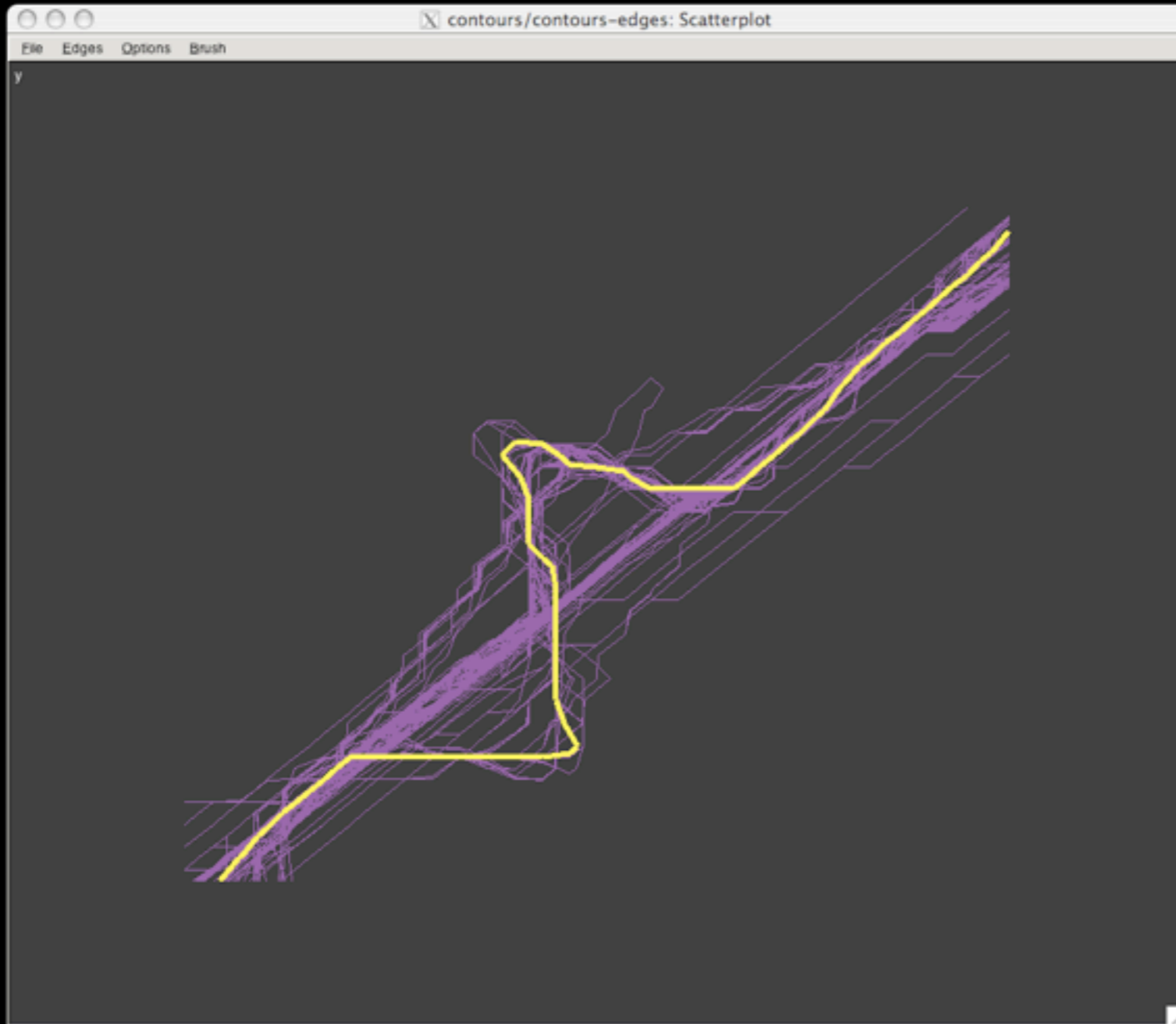




How did I find
that model?



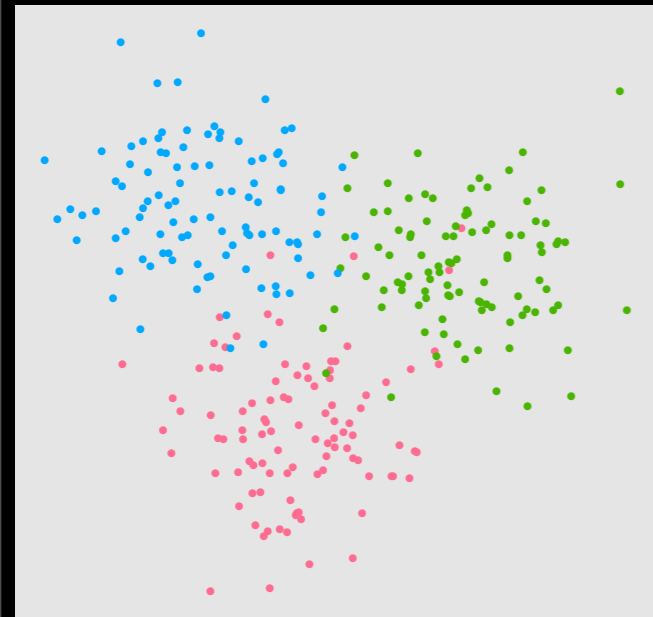




Classification algorithm

$$f: \mathbb{R}^p \rightarrow \{1, 2, \dots, k\}$$

Input

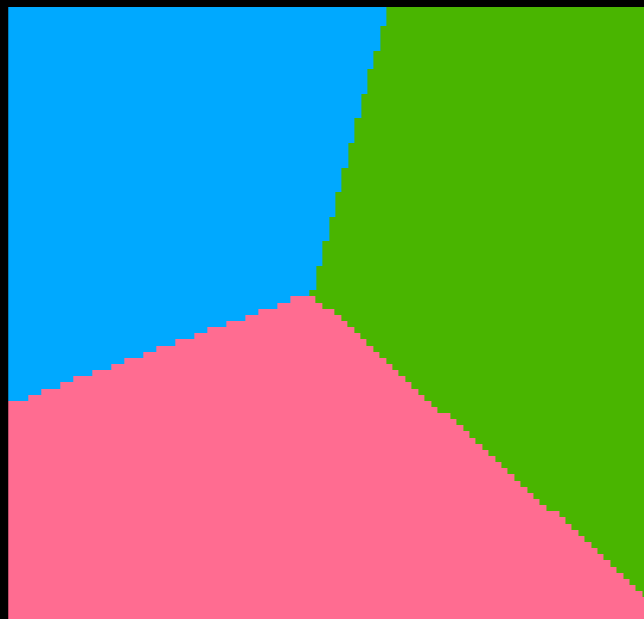
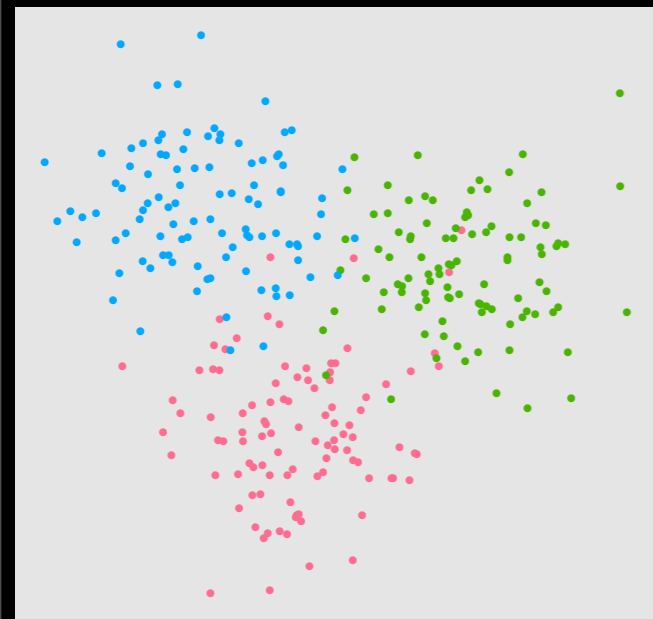


Classification algorithm

$$f: \mathbb{R}^p \rightarrow \{1, 2, \dots, k\}$$

Input

Prediction



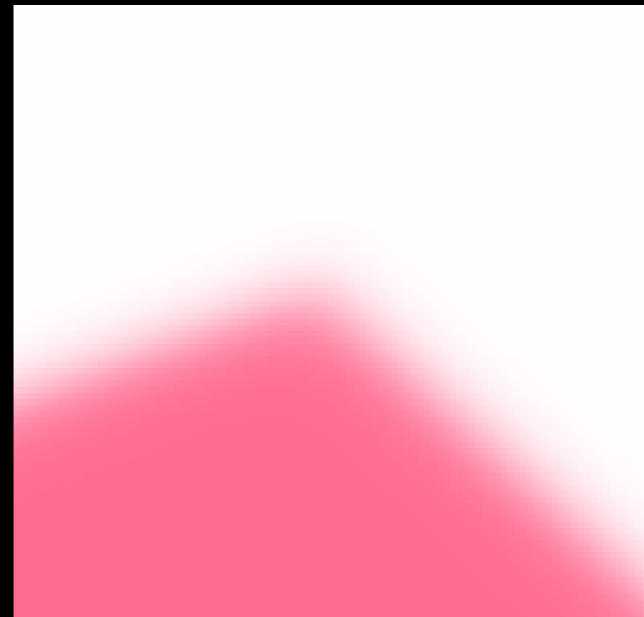
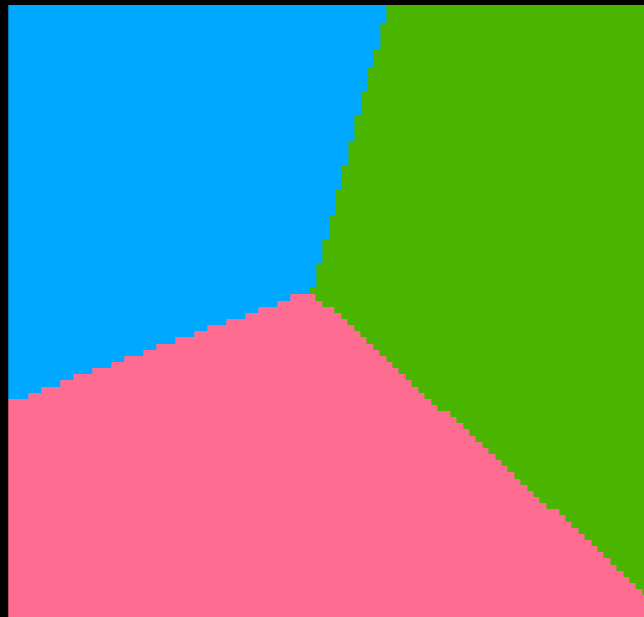
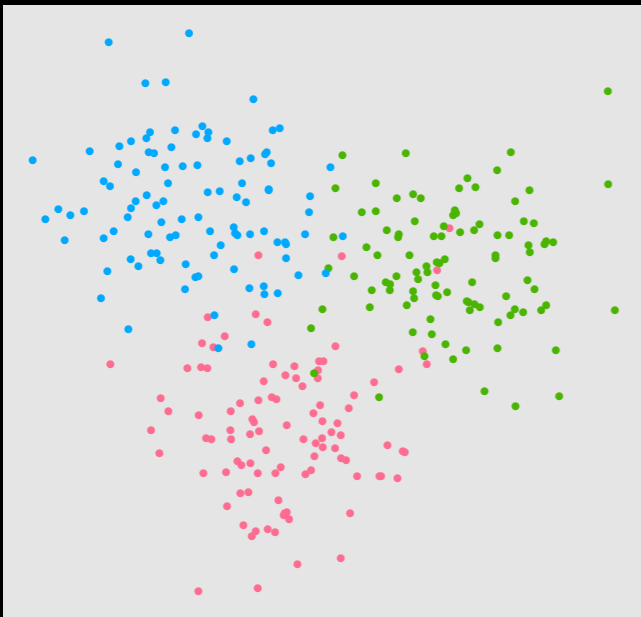
Classification algorithm

$$f: \mathbb{R}^p \rightarrow \{1, 2, \dots, k\}$$

Input

Prediction

Probabilities



Most also provide class membership probabilities

$$f: \mathbb{R}^p \rightarrow [0, 1]^k$$

Classification algorithm

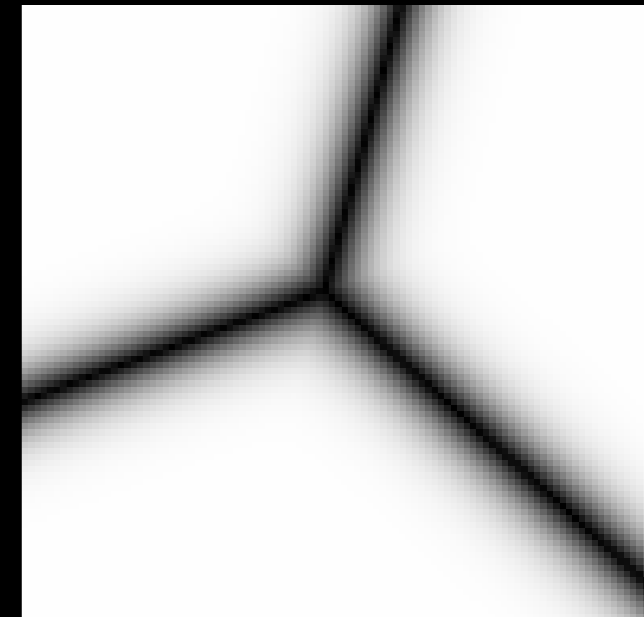
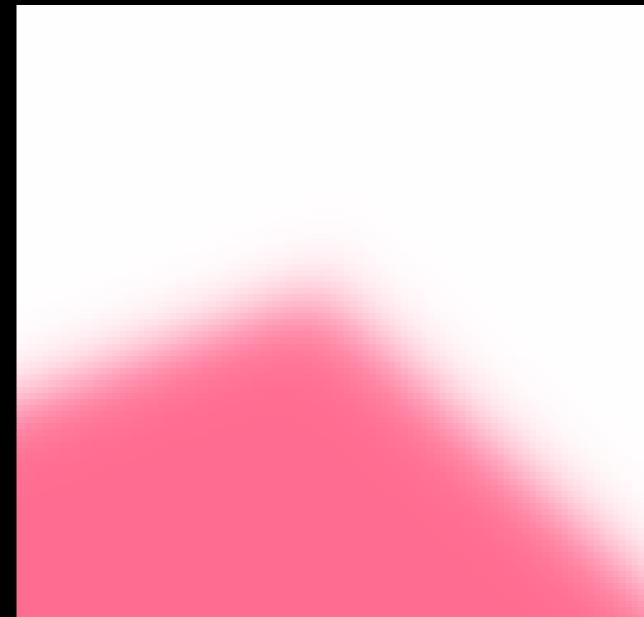
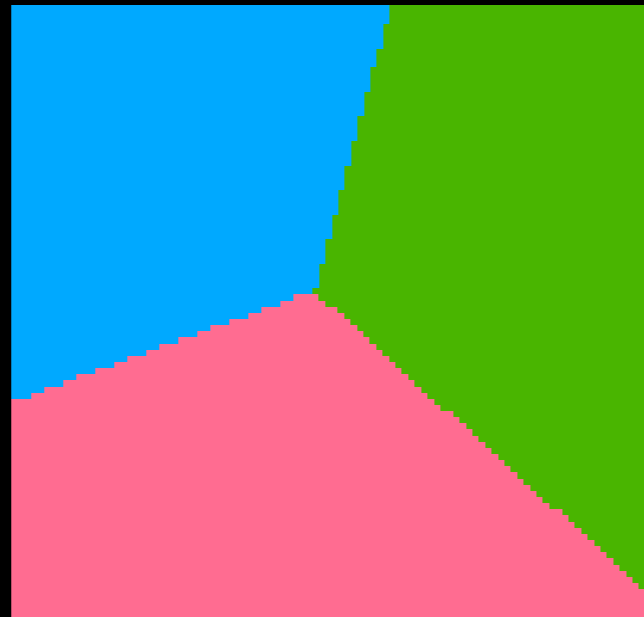
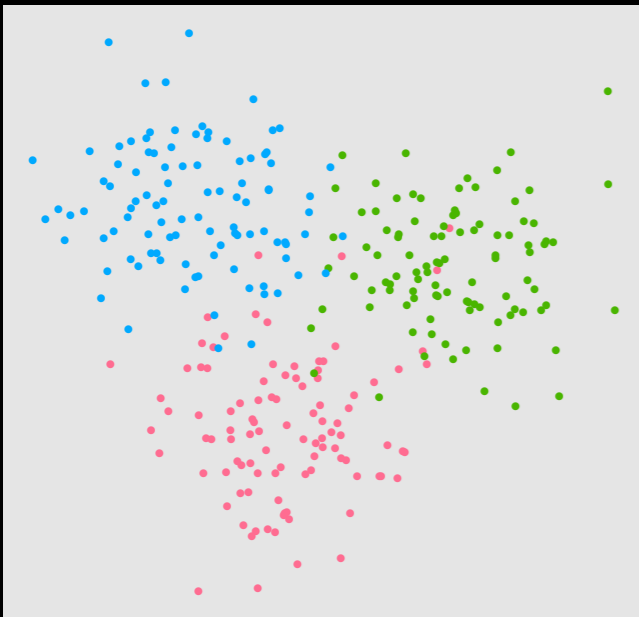
$$f: \mathbb{R}^p \rightarrow \{1, 2, \dots, k\}$$

Input

Prediction

Probabilities

Advantage



Most also provide class membership probabilities

$$f: \mathbb{R}^p \rightarrow [0, 1]^k$$

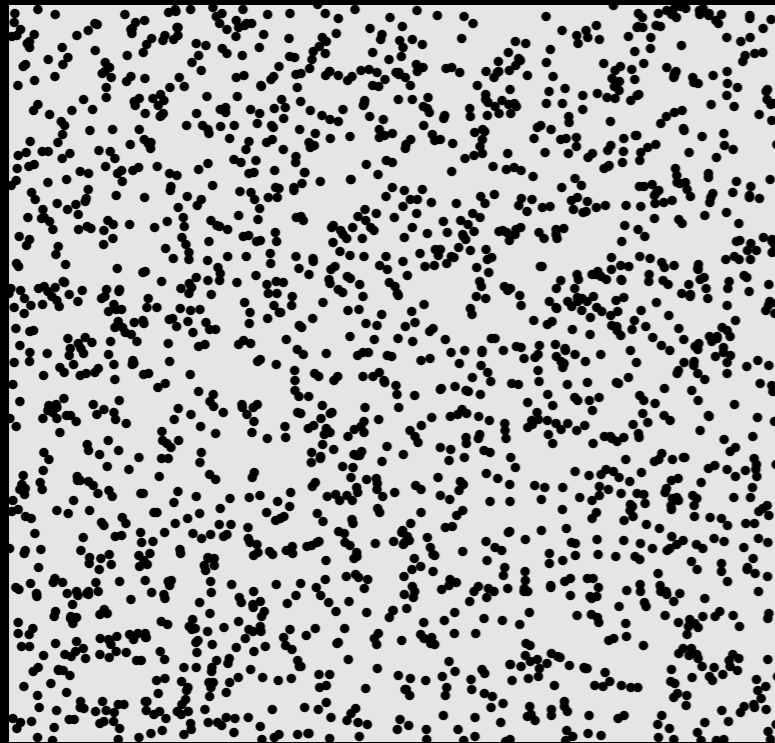


$P(\text{best}) -$
 $P(\text{second best})$

How to find the boundaries?

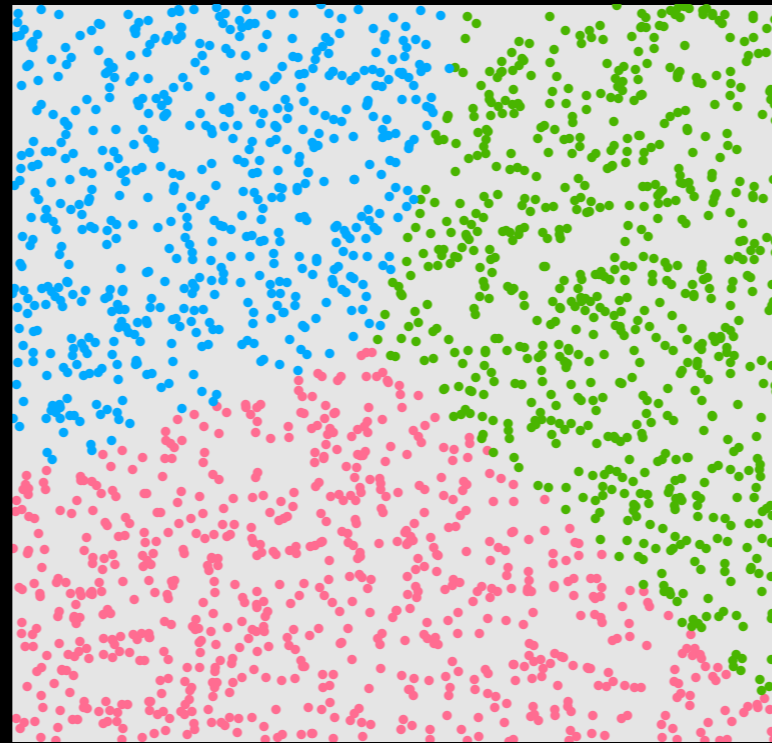
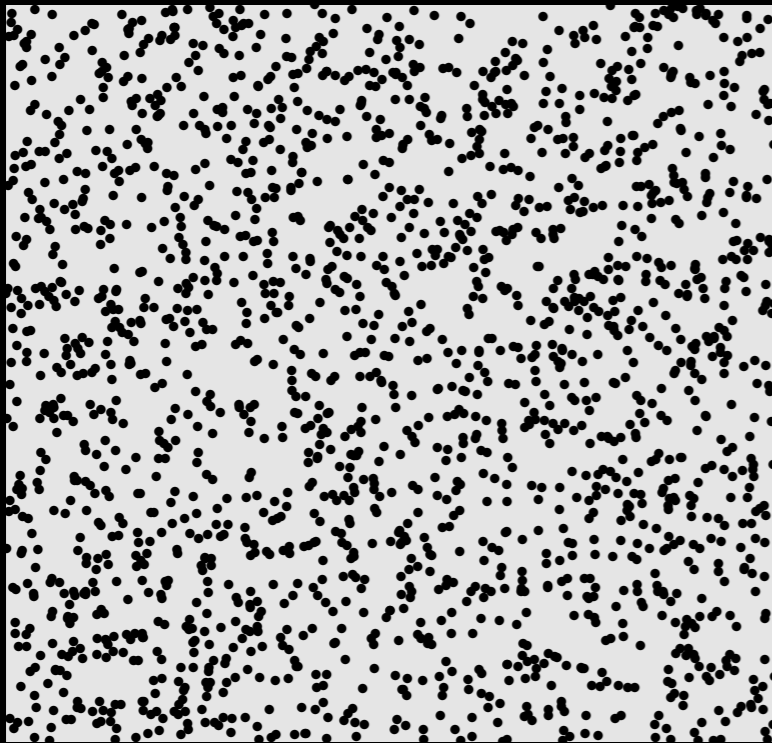
How to find the boundaries?

Random sample



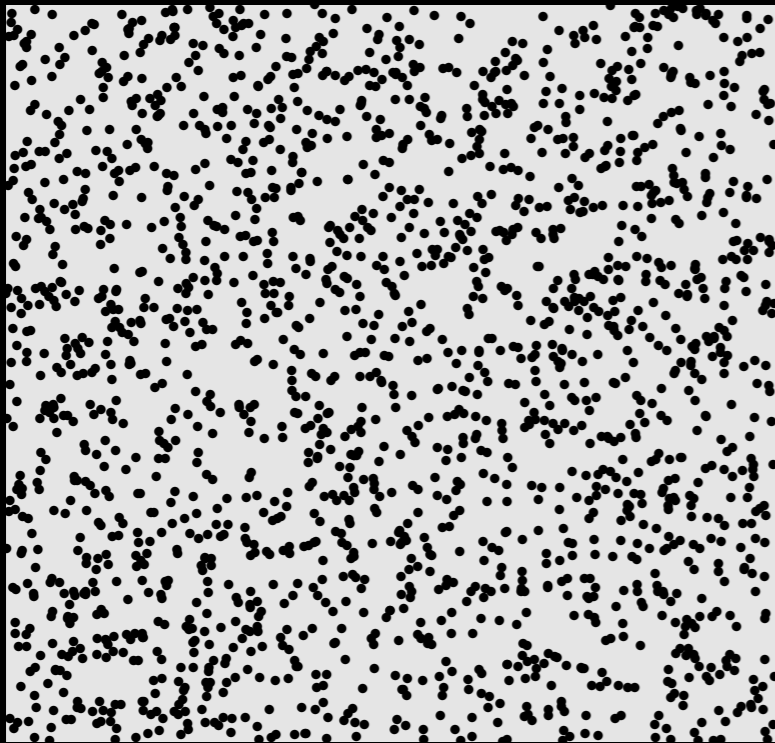
How to find the boundaries?

Random sample Classify

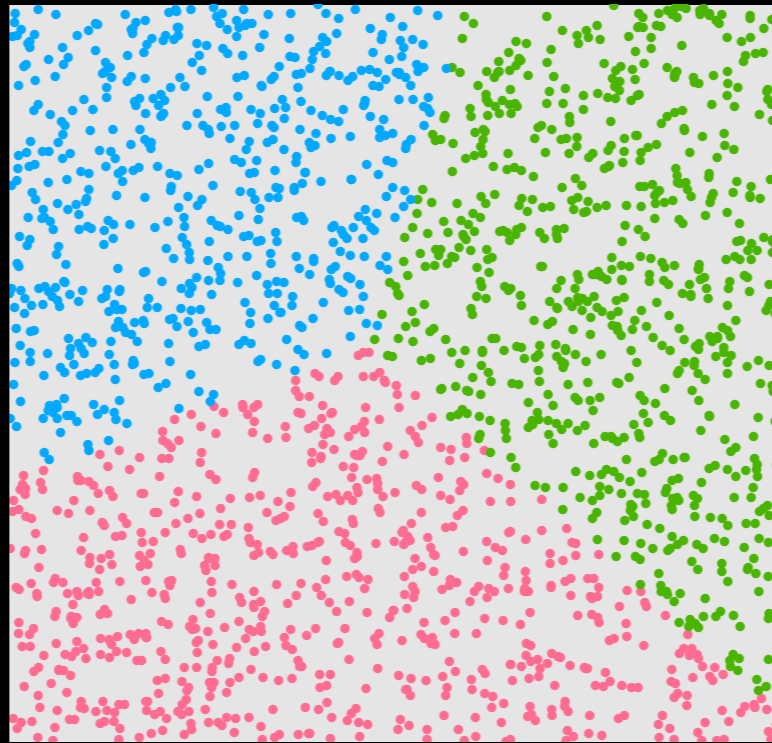


How to find the boundaries?

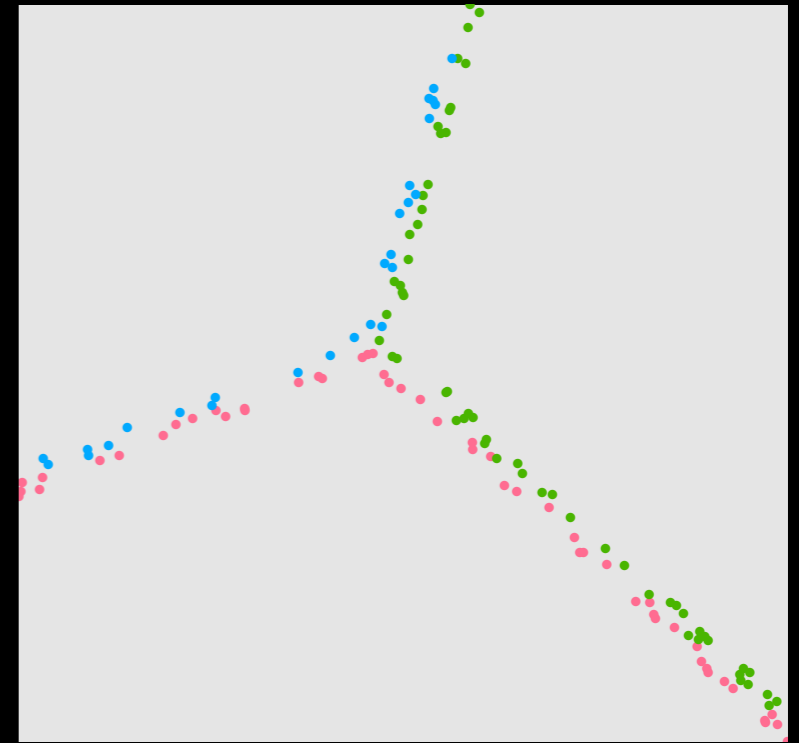
Random sample Classify



Classify

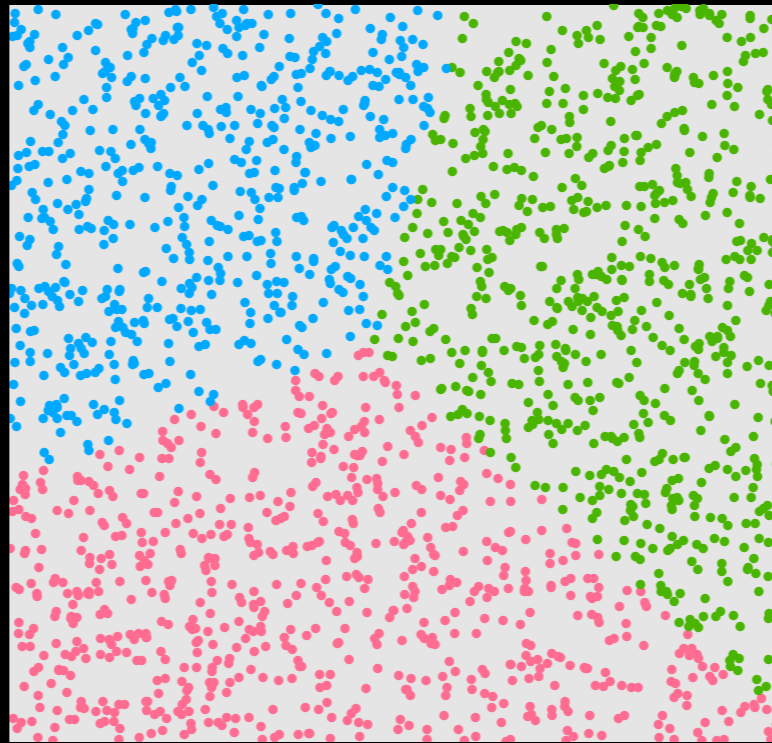
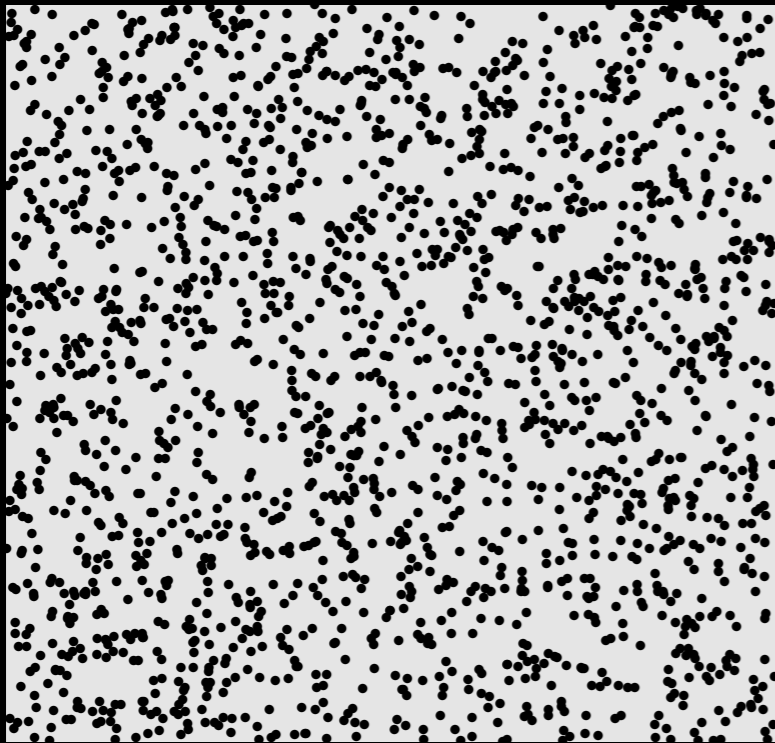


Low advantage

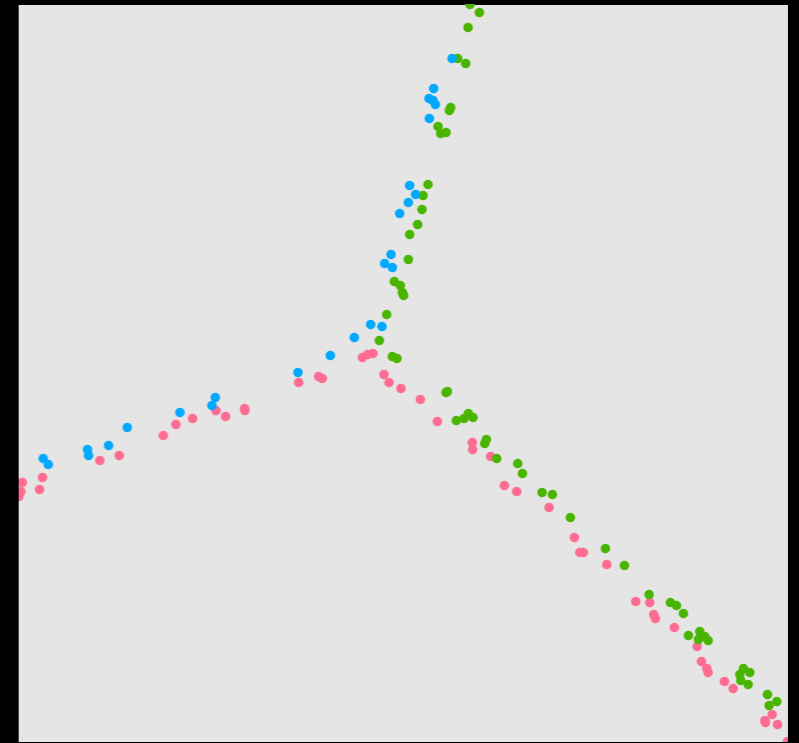


How to find the boundaries?

Random sample Classify



Low advantage



Crude method, but works for
all classification algorithms and
for moderate dimensionality

Ensembles of linear models

Display the model
in the data space

Look at many
members of a collection

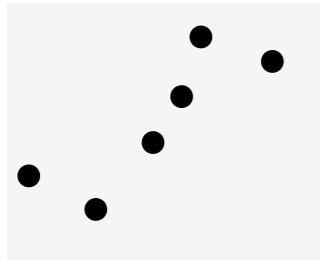
Explore the process of
fitting, not just the end result

Data

- Fertility in French-speaking Swiss provinces in the late 1800's
- Predict fertility based on:
 - proportion of agricultural workers
 - average performance on an army examination
 - amount of higher education
 - proportion of Catholics
 - infant mortality

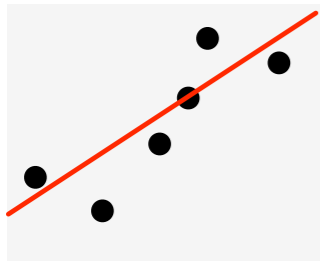
Model

- Linear models with all combinations of covariates (2^p models)
- What can looking at all models tell us that looking at just a few can't?



Observation

Obs ID
Original data
Model-observation summaries



Model

Model ID
Model fit statistics

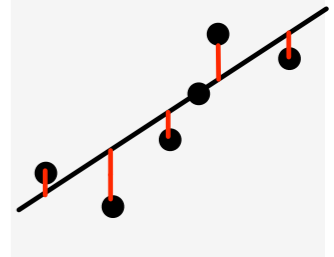
$f(\alpha, \beta)$

Estimate

Estimate ID
n
model-estimate summaries

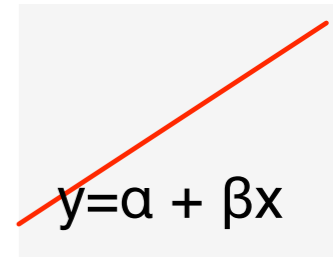
Model-Observation

Obs ID
Model ID
Diagnostics
Fit quality



Model-Estimate

Model ID
Estimate ID
Raw
Standardised
Uncertainty



1

many

1

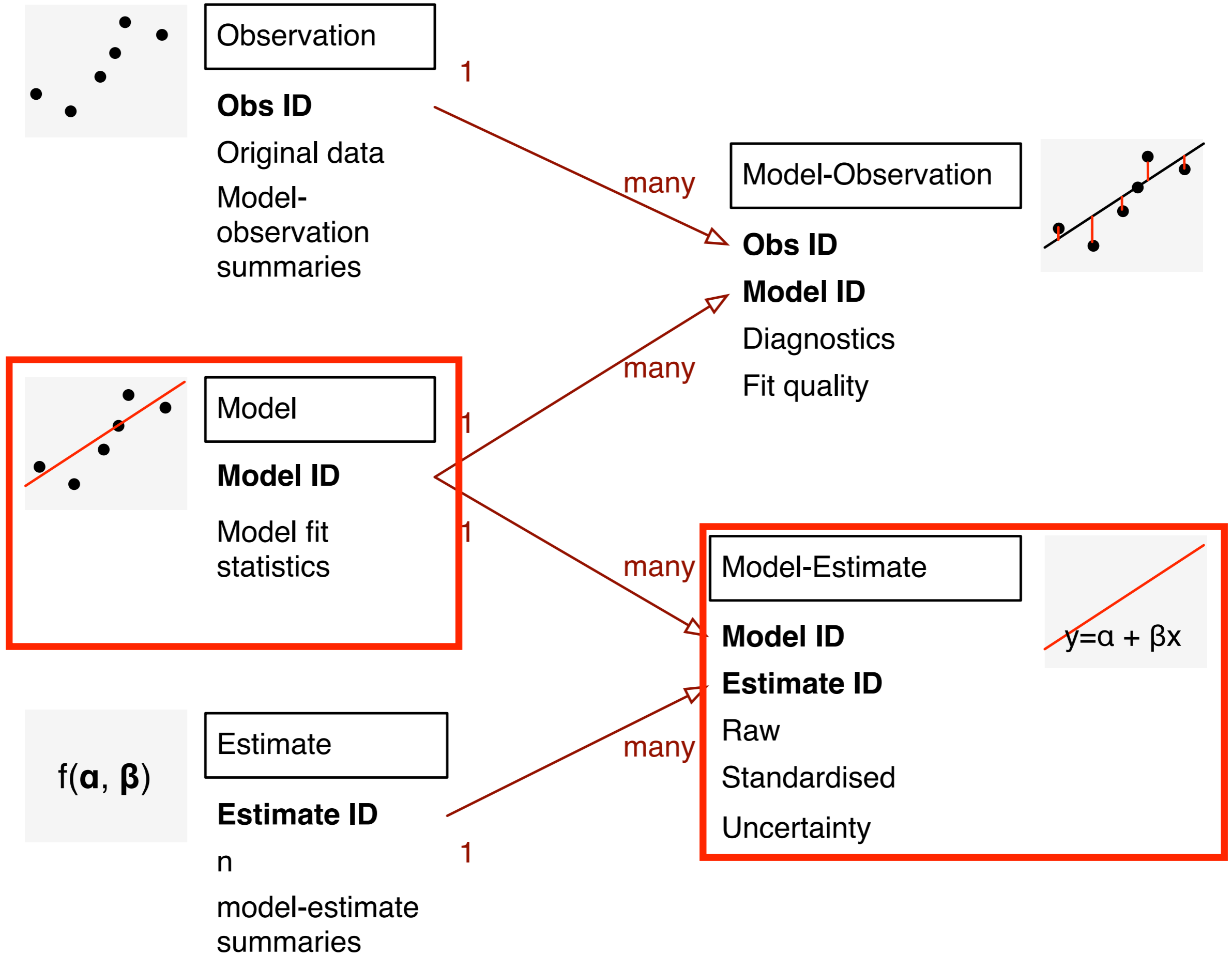
many

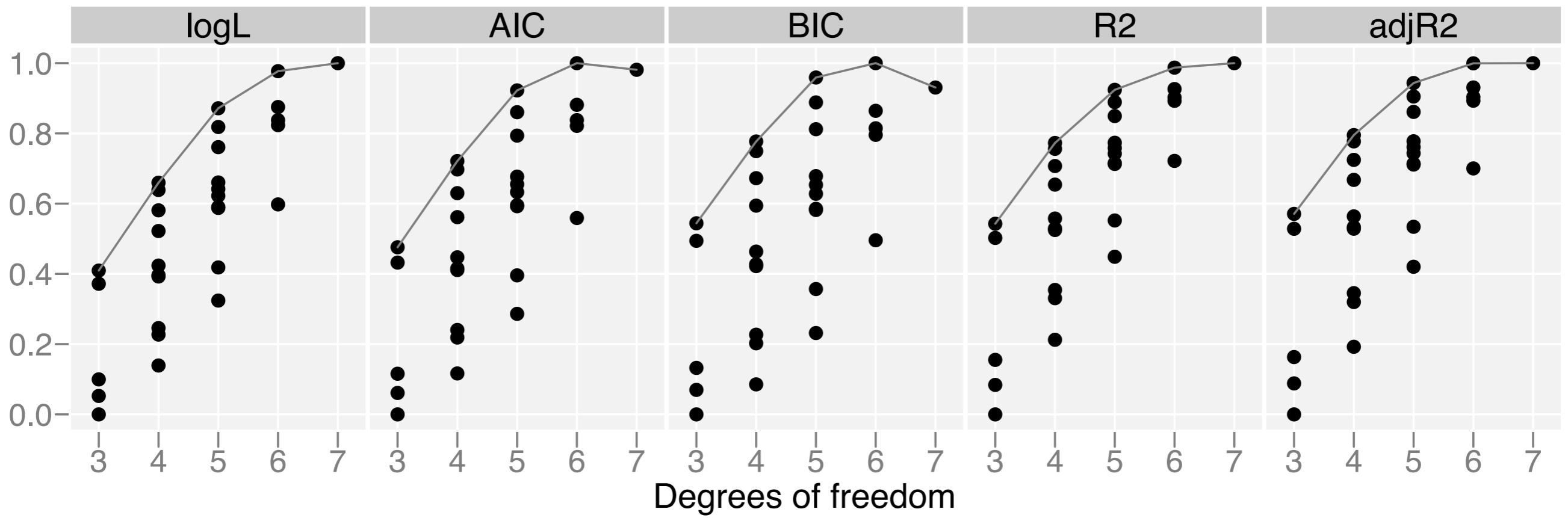
1

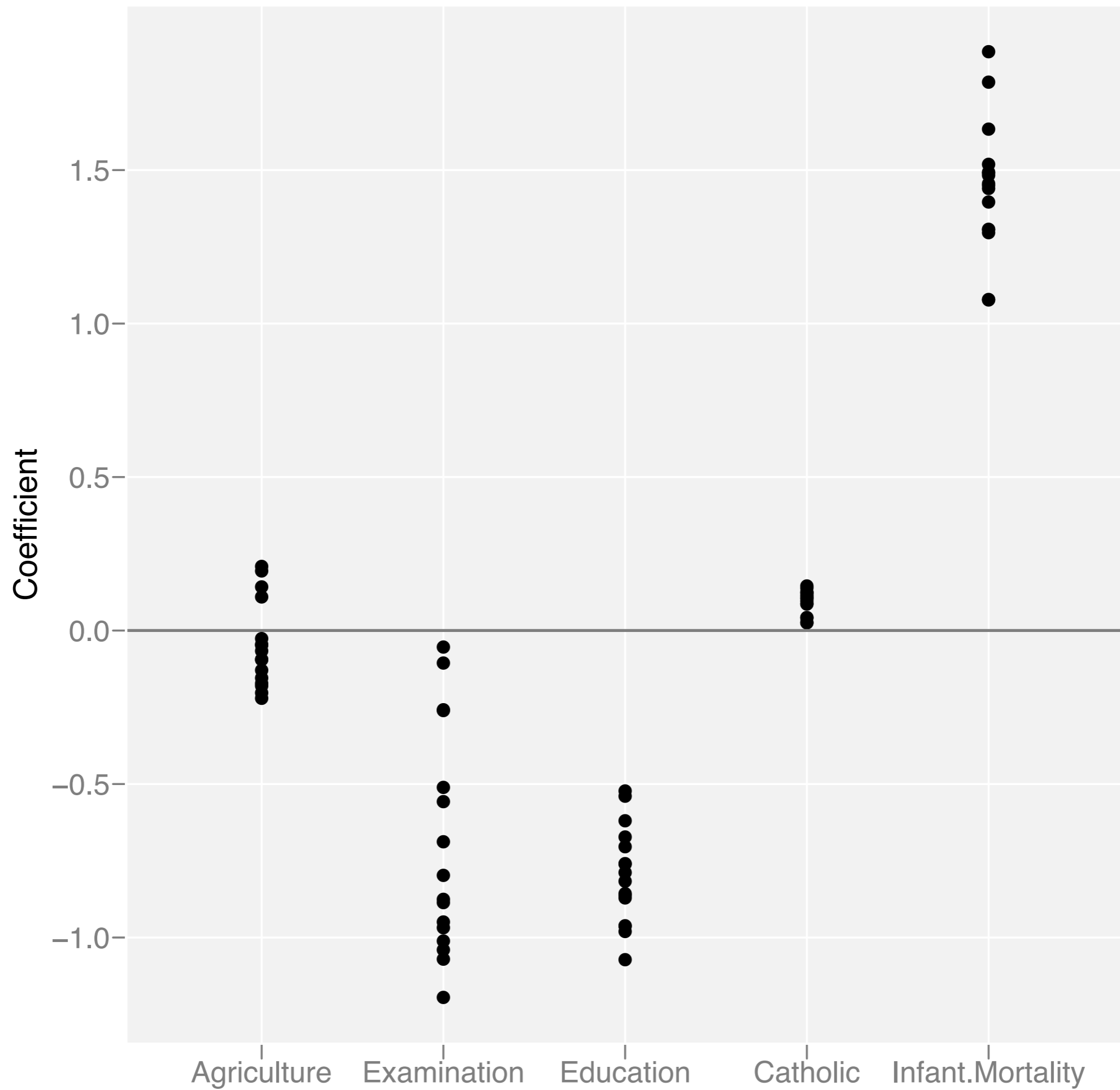
many

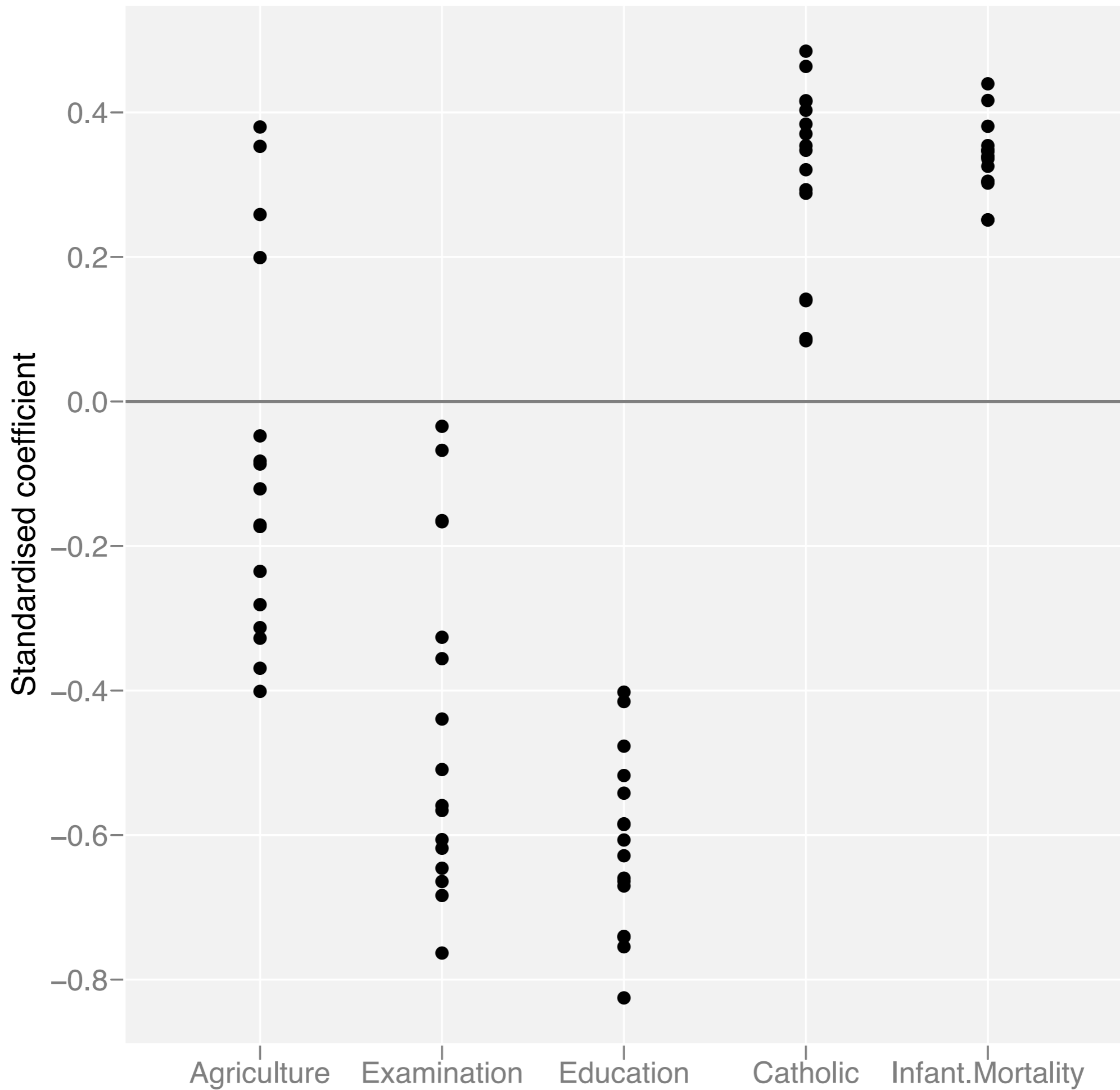
1

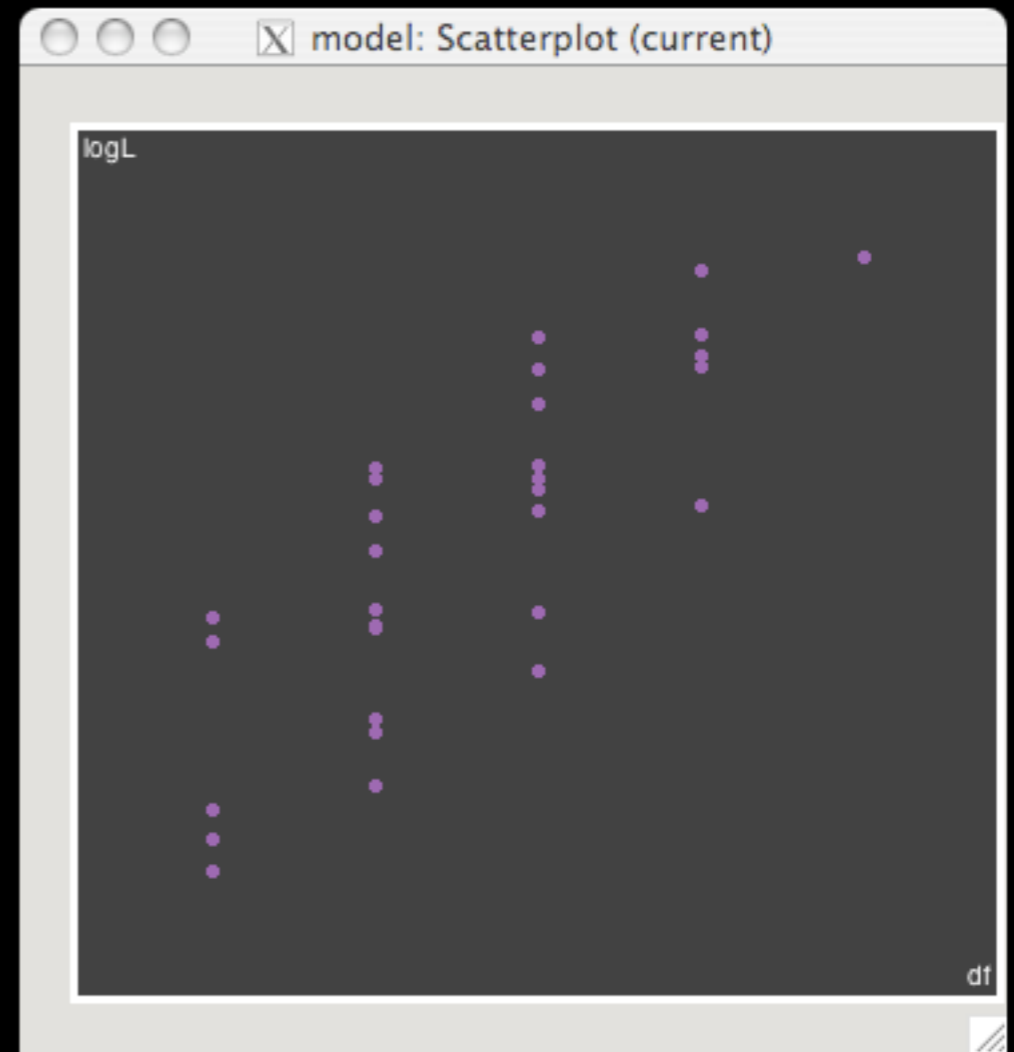
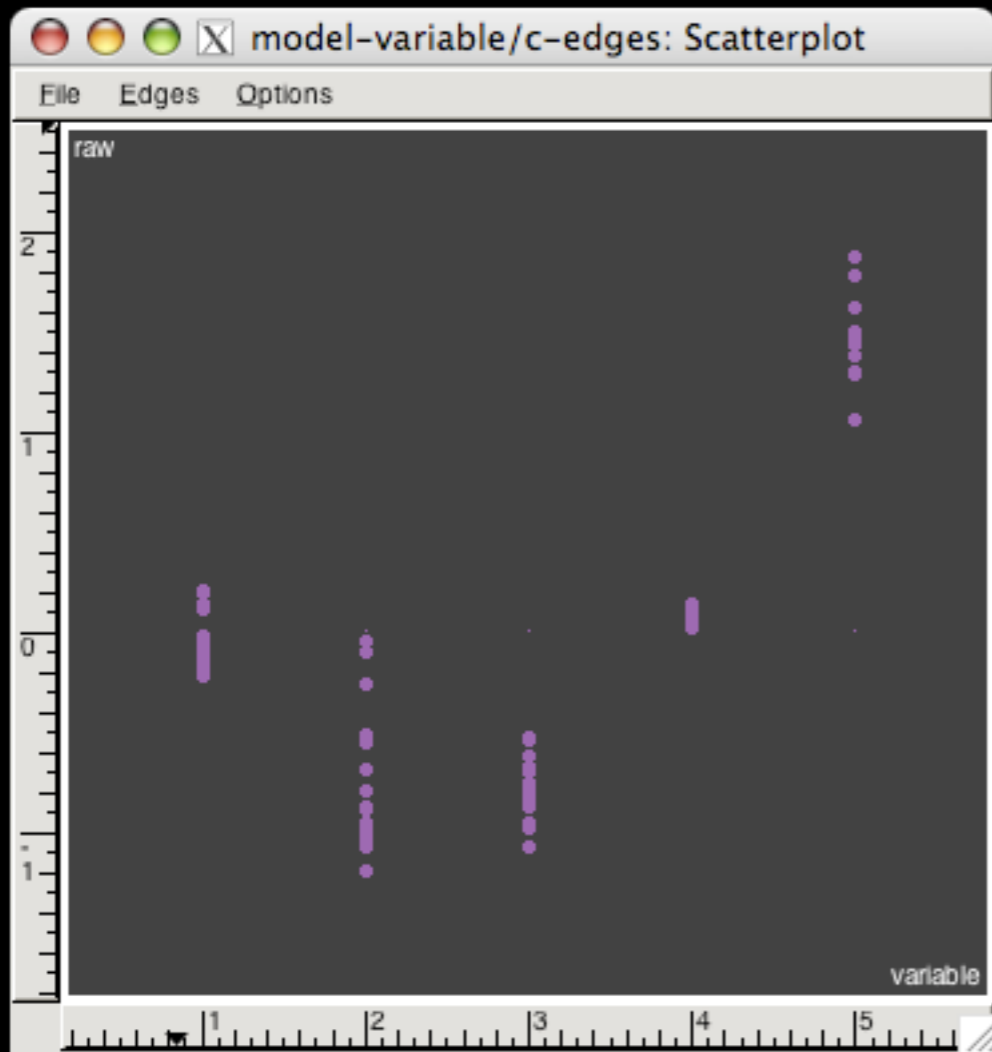
many











Conclusions

Other methods

- MANOVA
- Self-organising maps (clusterfly)
- Hierarchical clustering (clusterfly)
- Classification methods (classifly)
- Projection pursuit (tourr)

The future

- Better iteration between modelling and visualisation
- Foundations to make interactive graphics easy to produce in R