



Help! They dumped a dataset on our doorstep...

... and want us to predict something

By Richard Berendsen



is powered by



My world is:

Nou moe?! - Guust Flater

All our worlds > www.luminis.eu

About me

Richard Berendsen

- Search and data engineer at Luminis
- BSc & MSc Artificial Intelligence
- PhD Information retrieval (aka search)



@rwberendsen



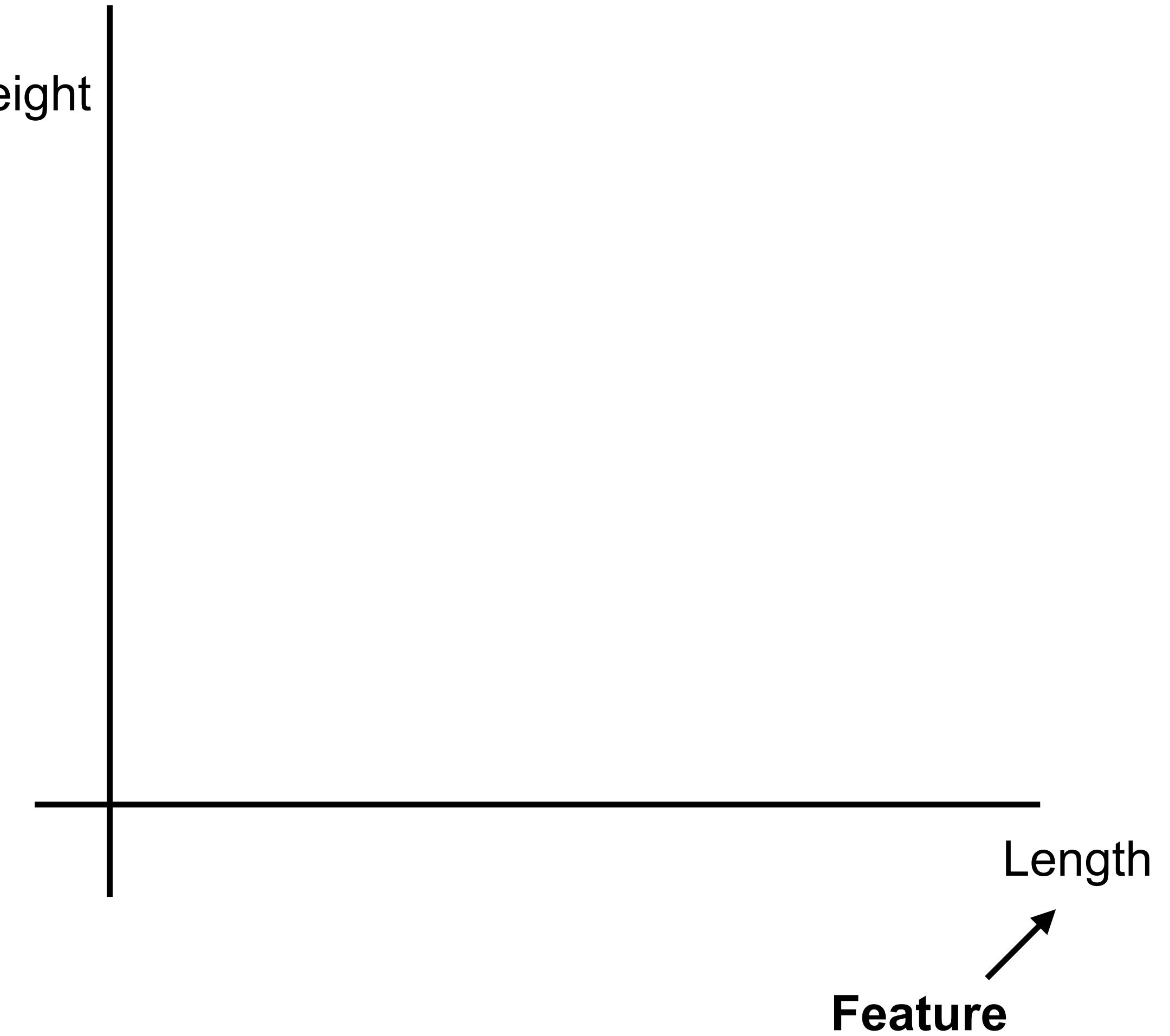
Emphasis

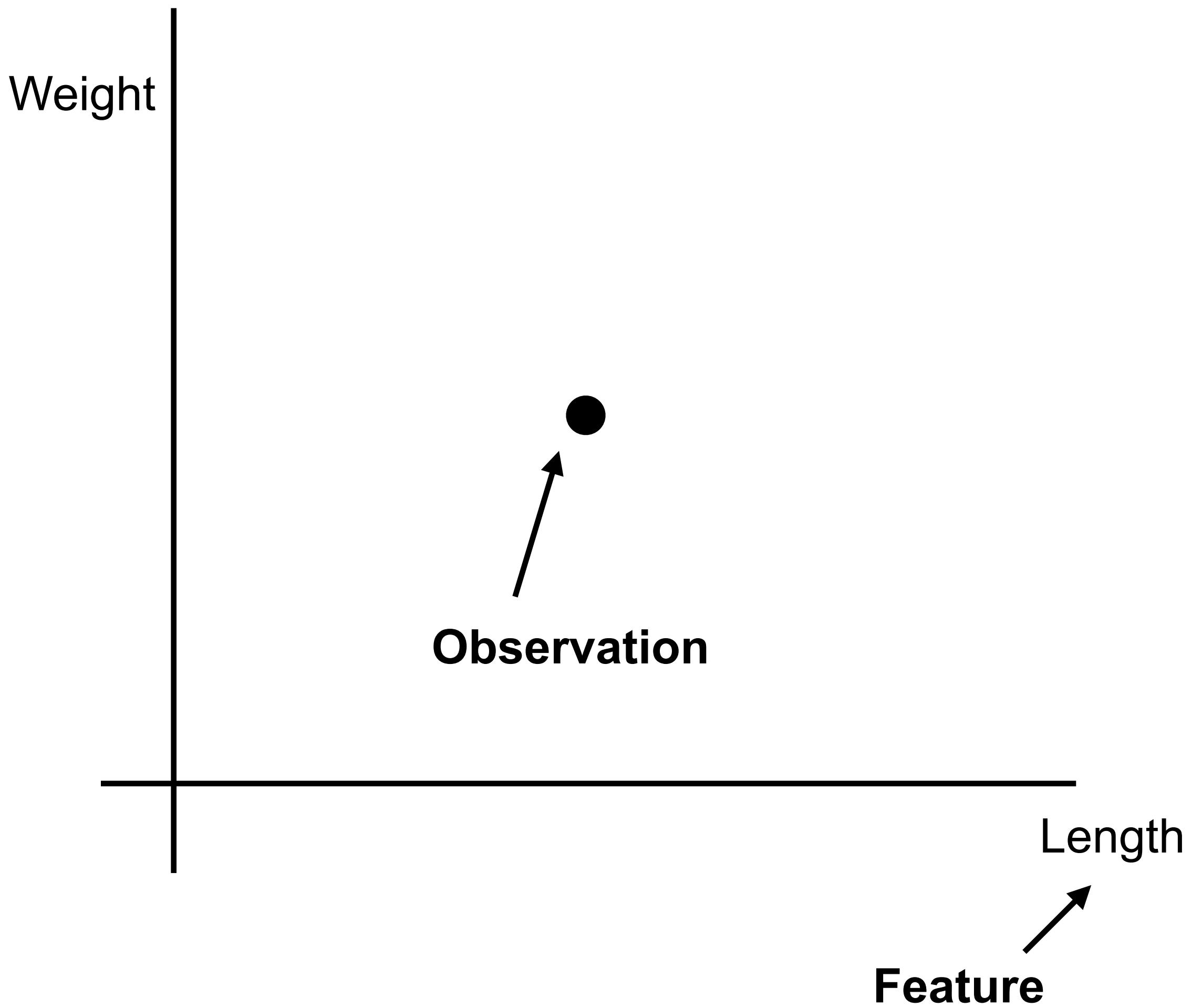
- Data science process
- Choices
- Questions
- Step by step demo

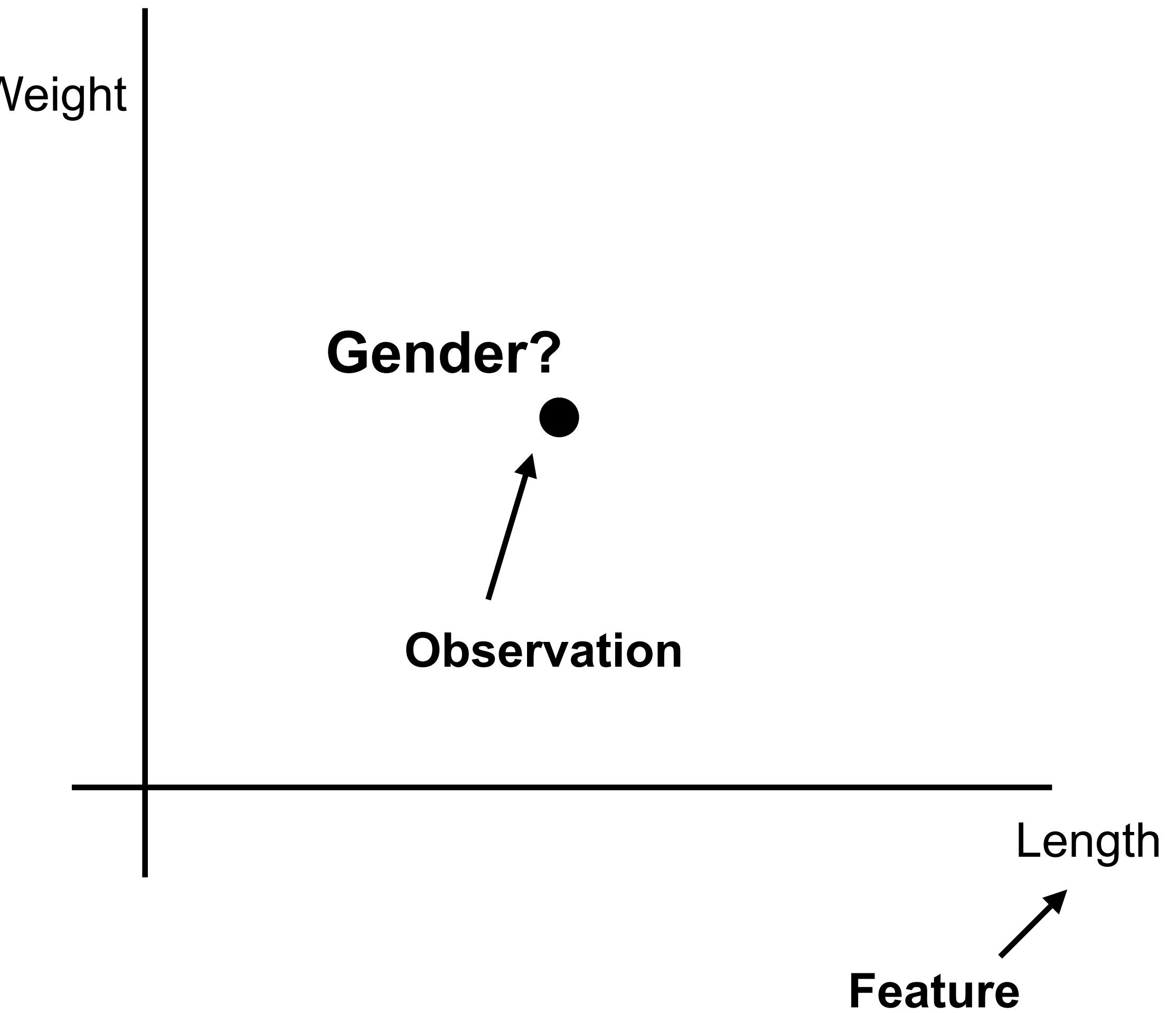
But first...

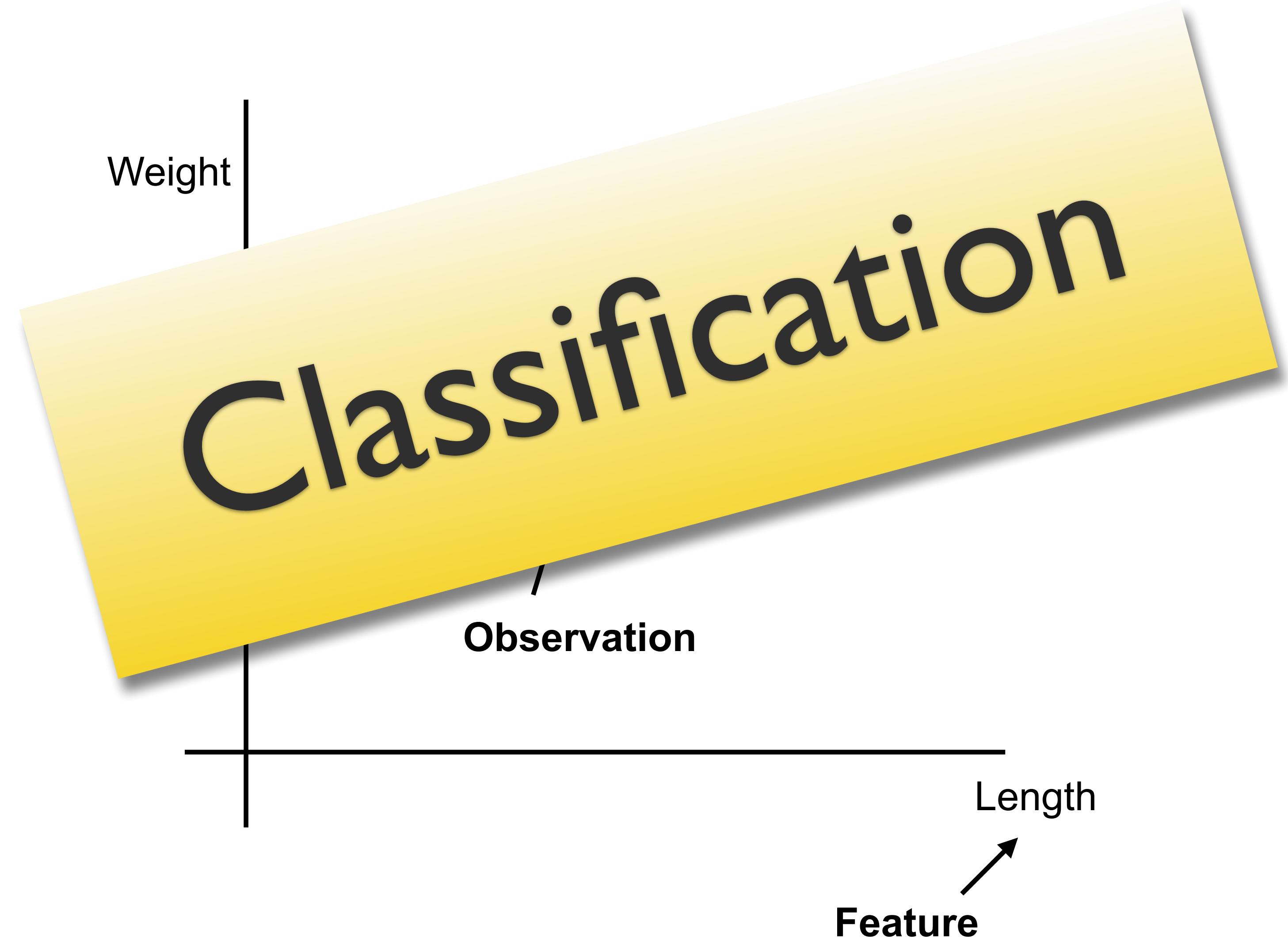
- A little data science intro
- A quick look at the stack

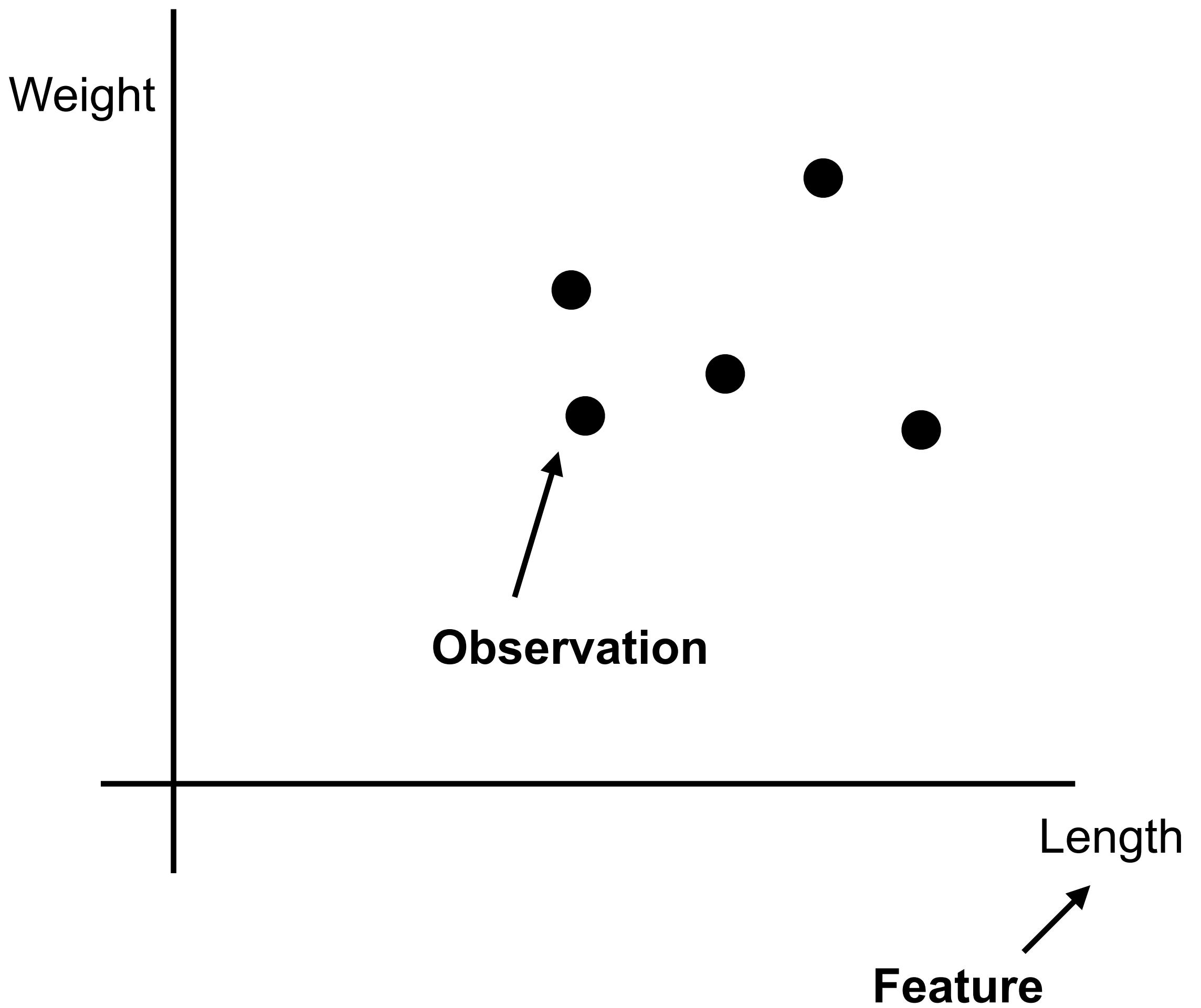


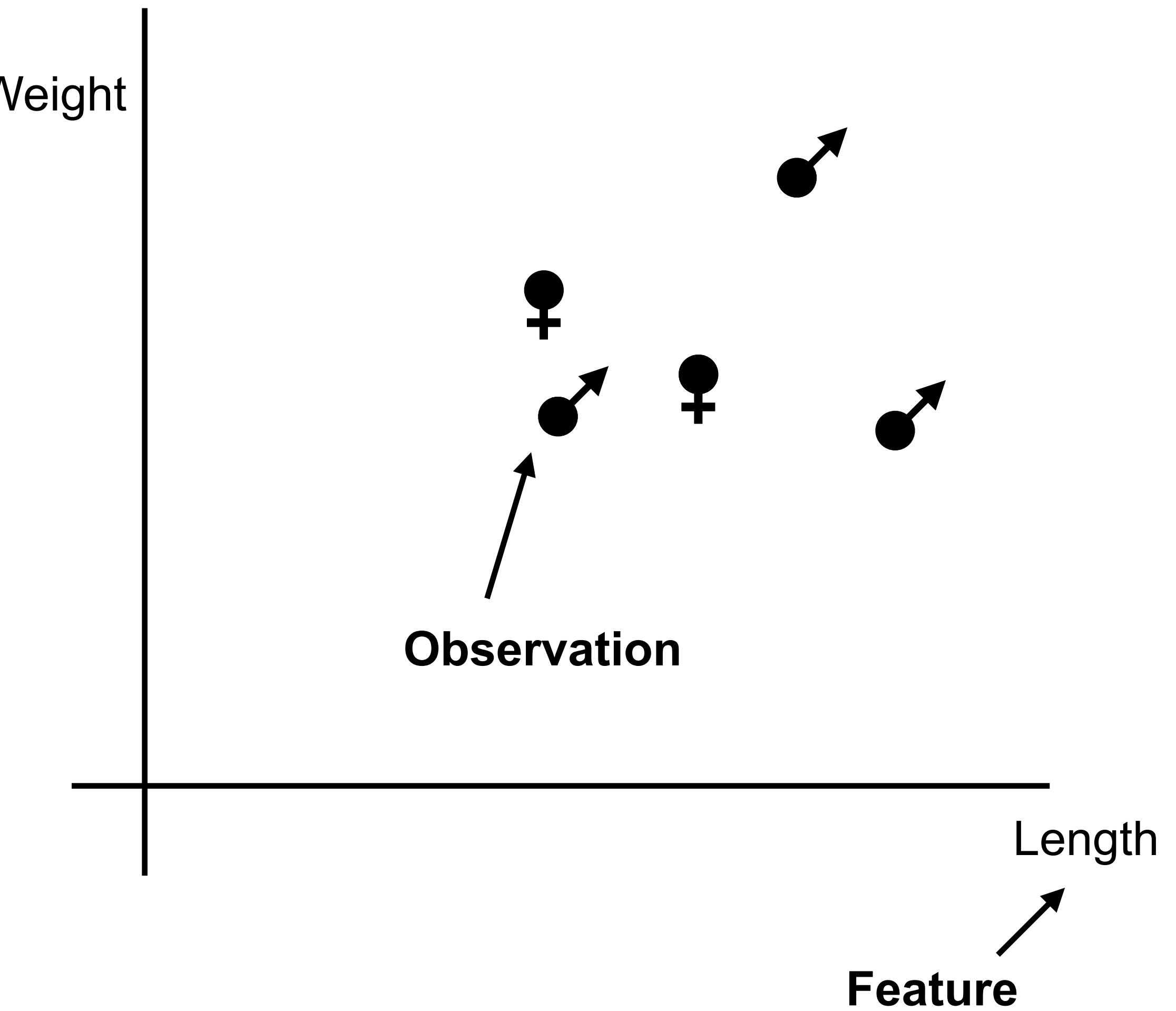


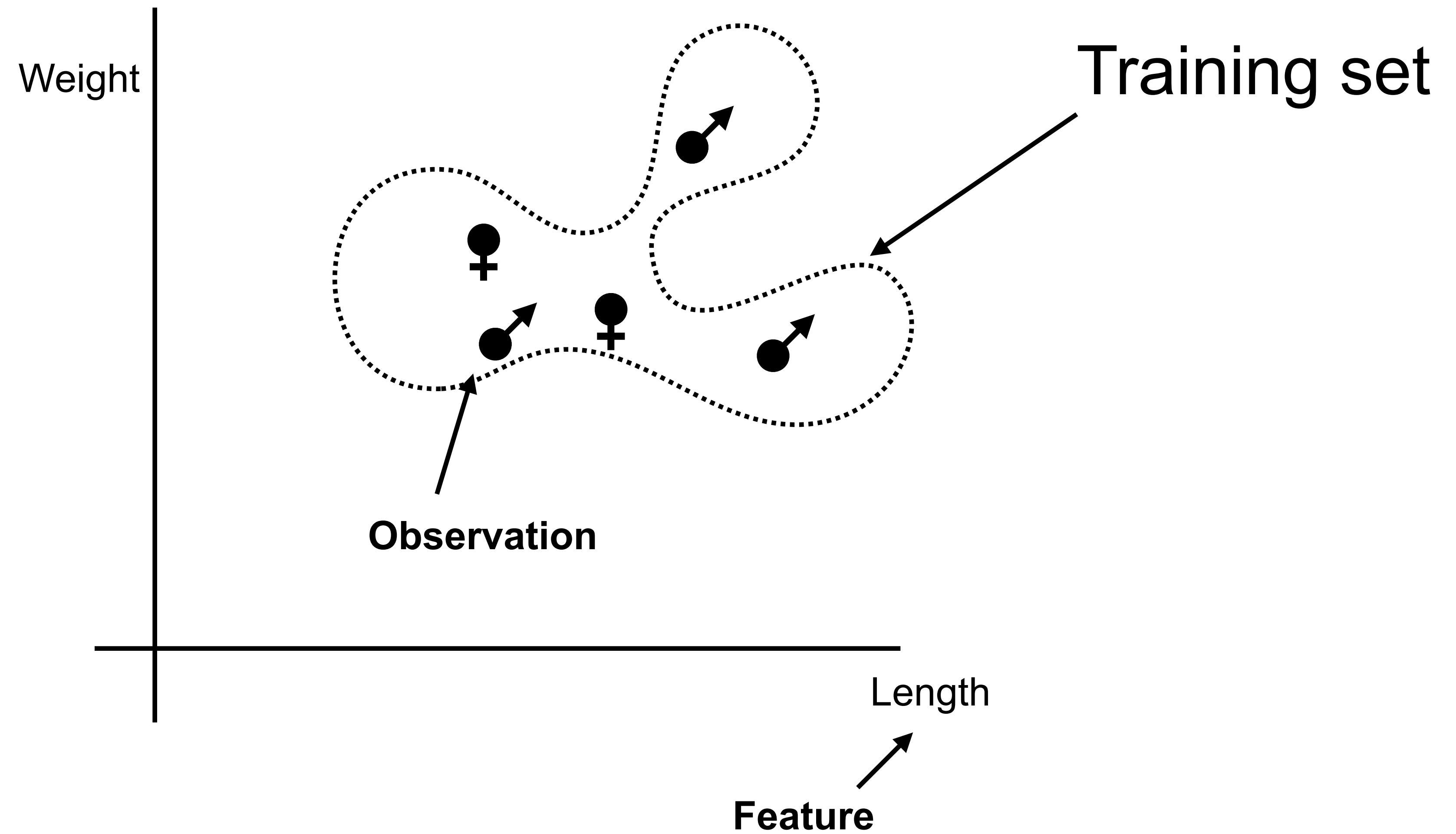


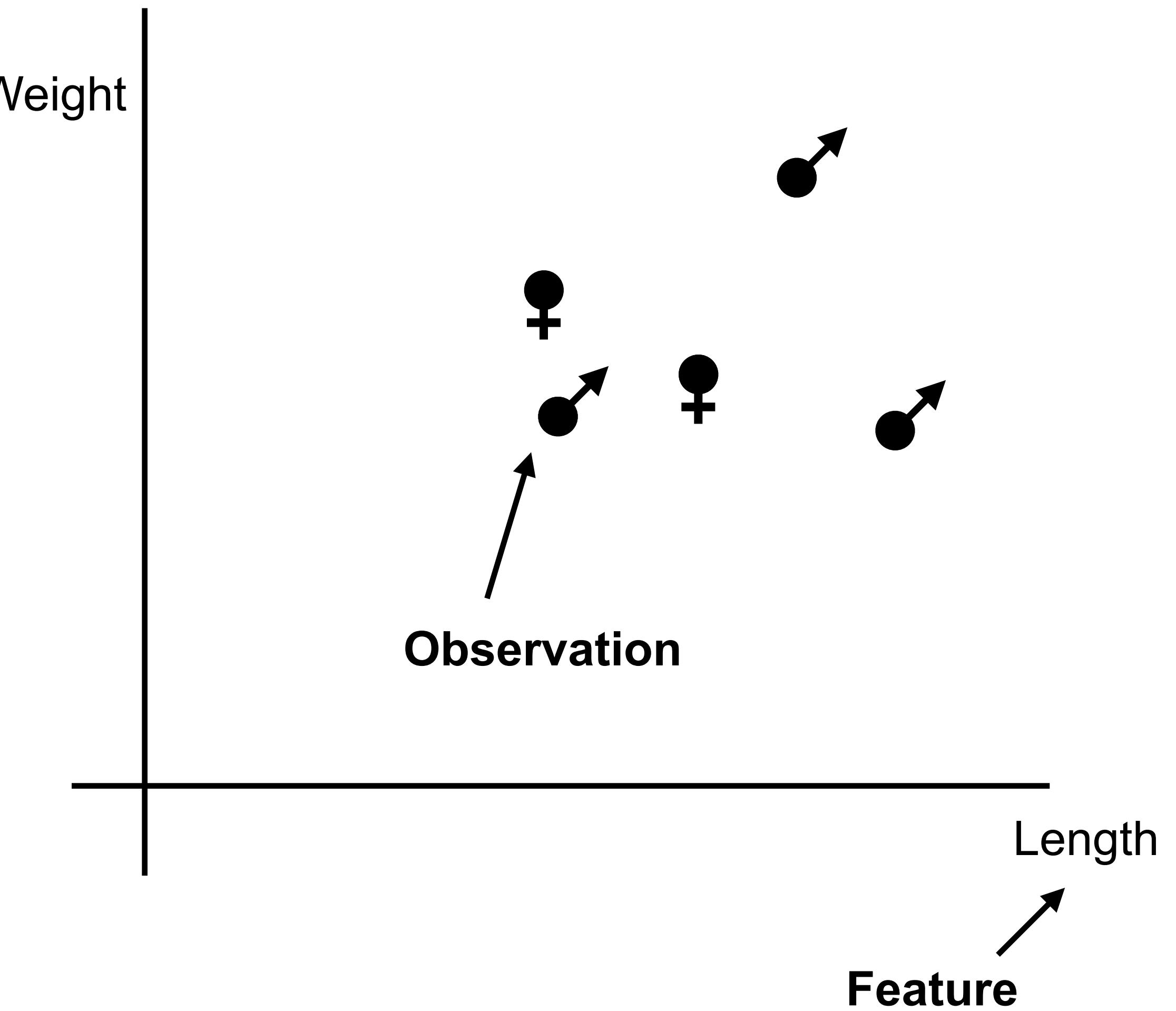


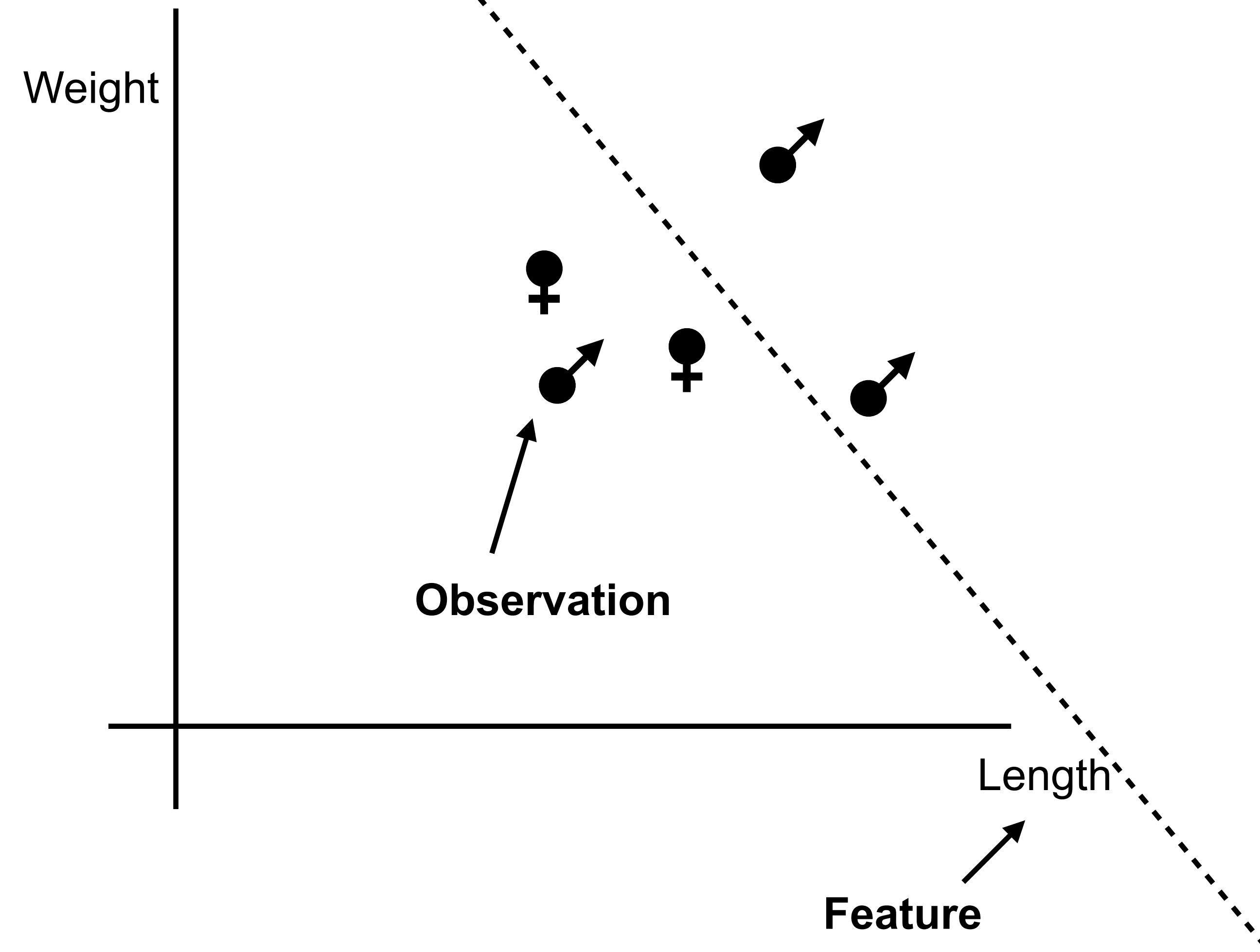


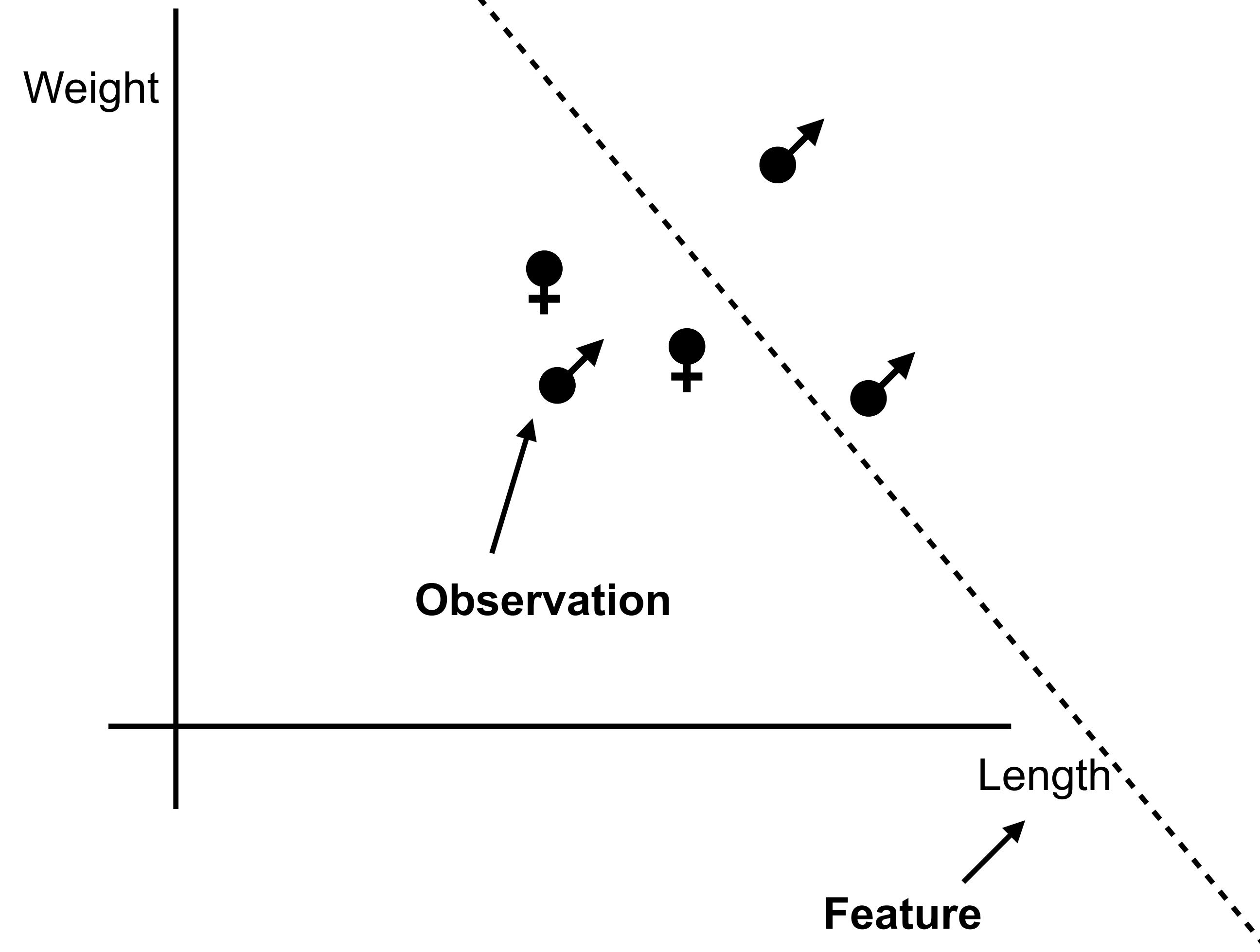


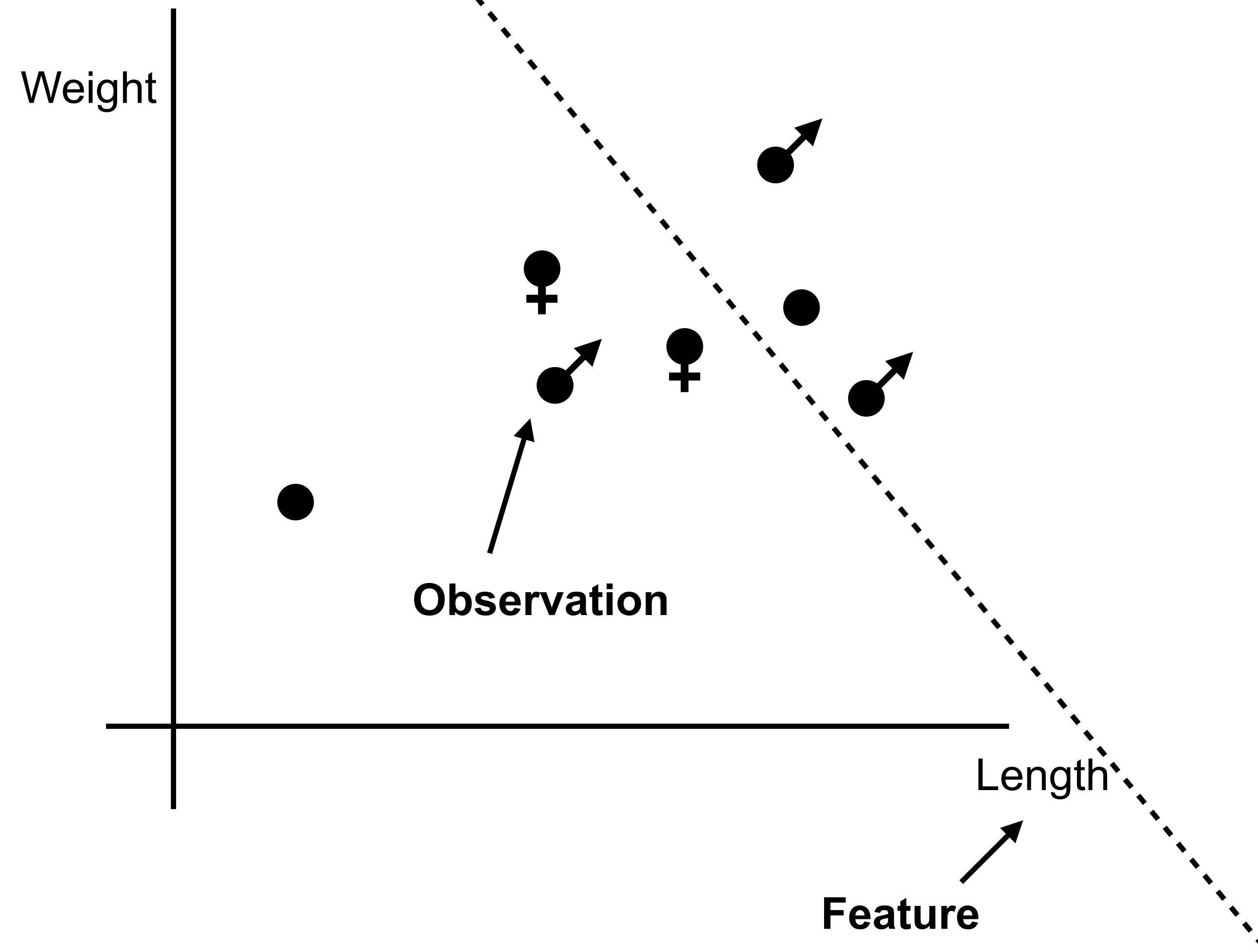


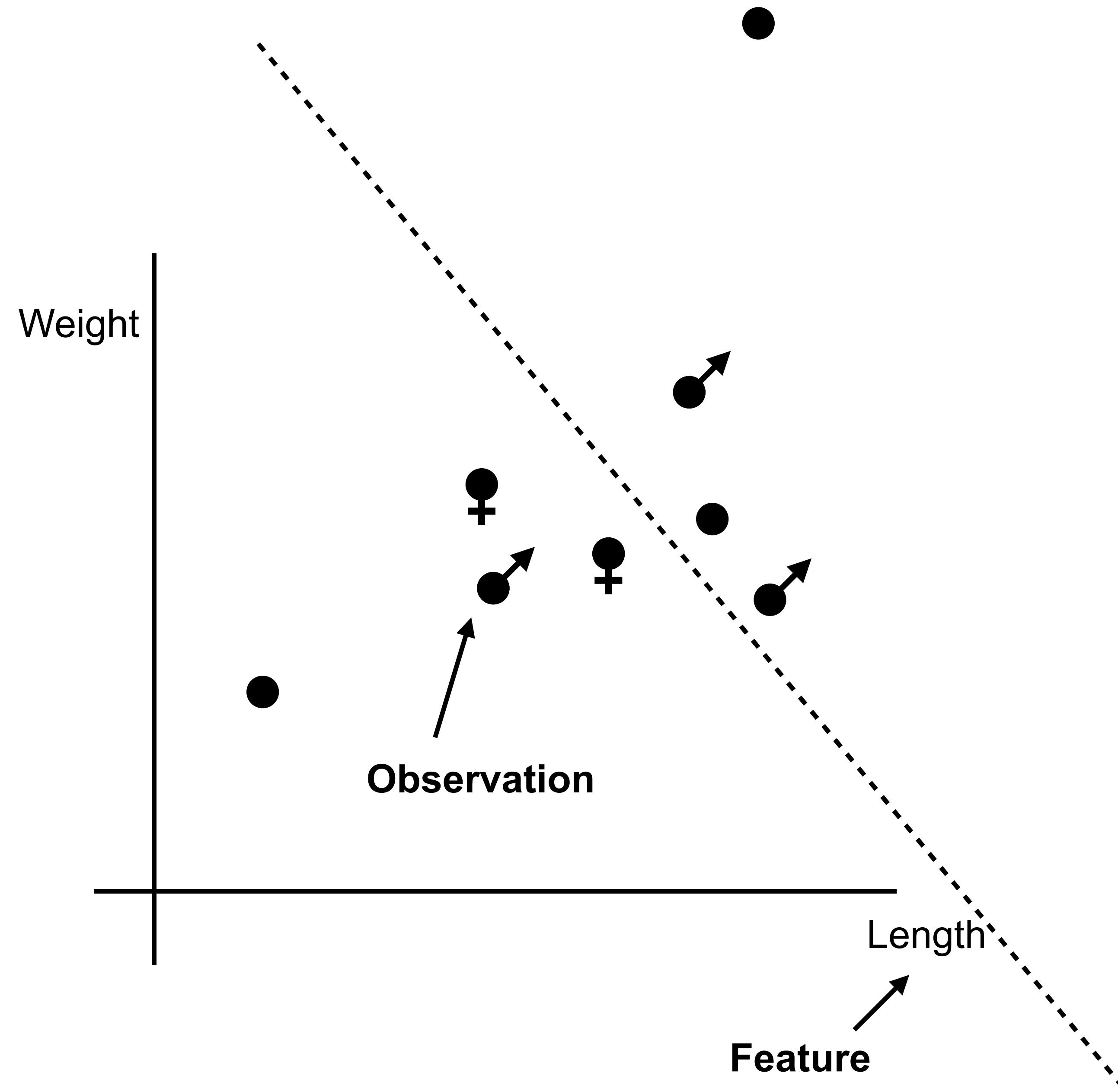


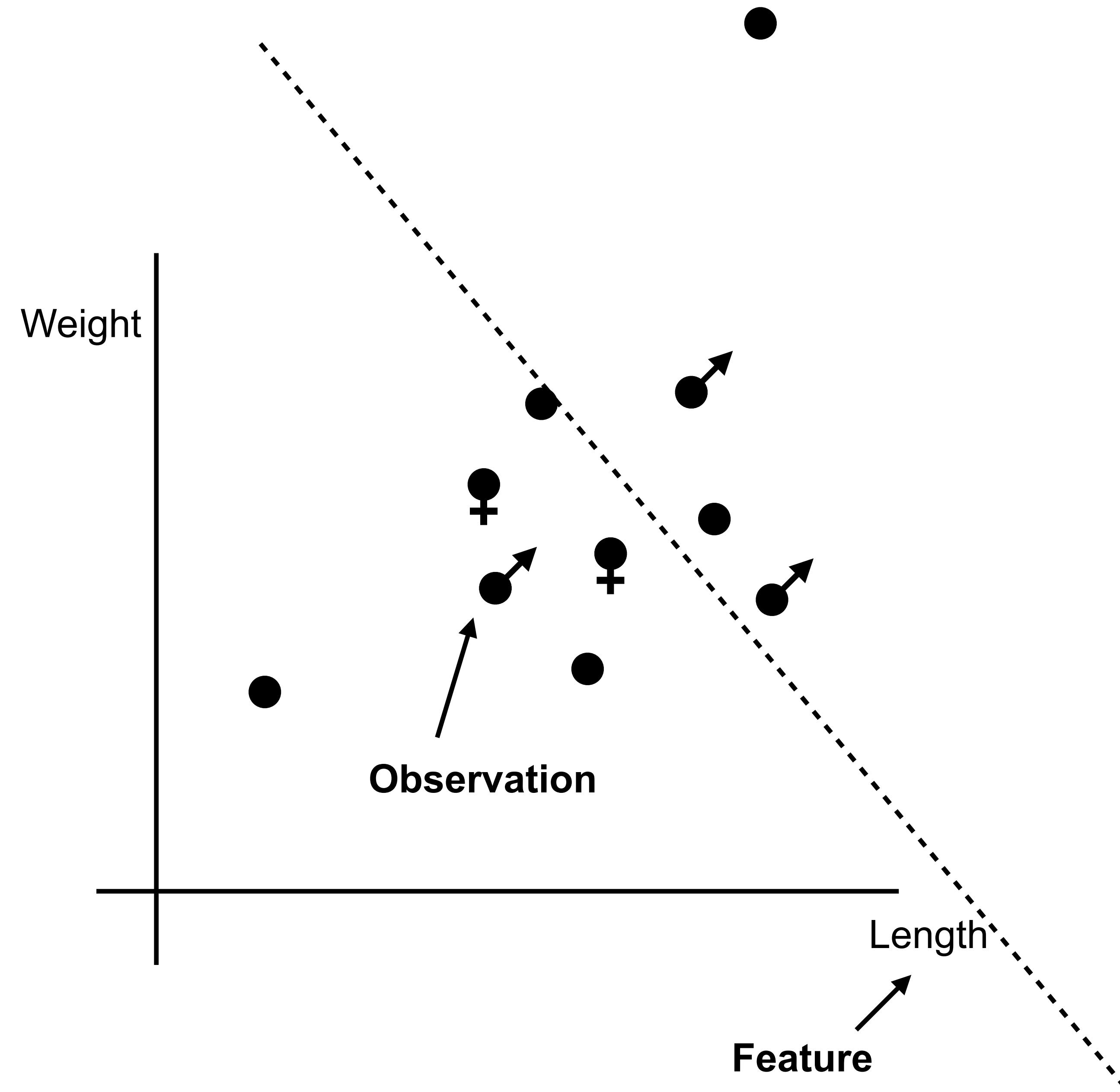


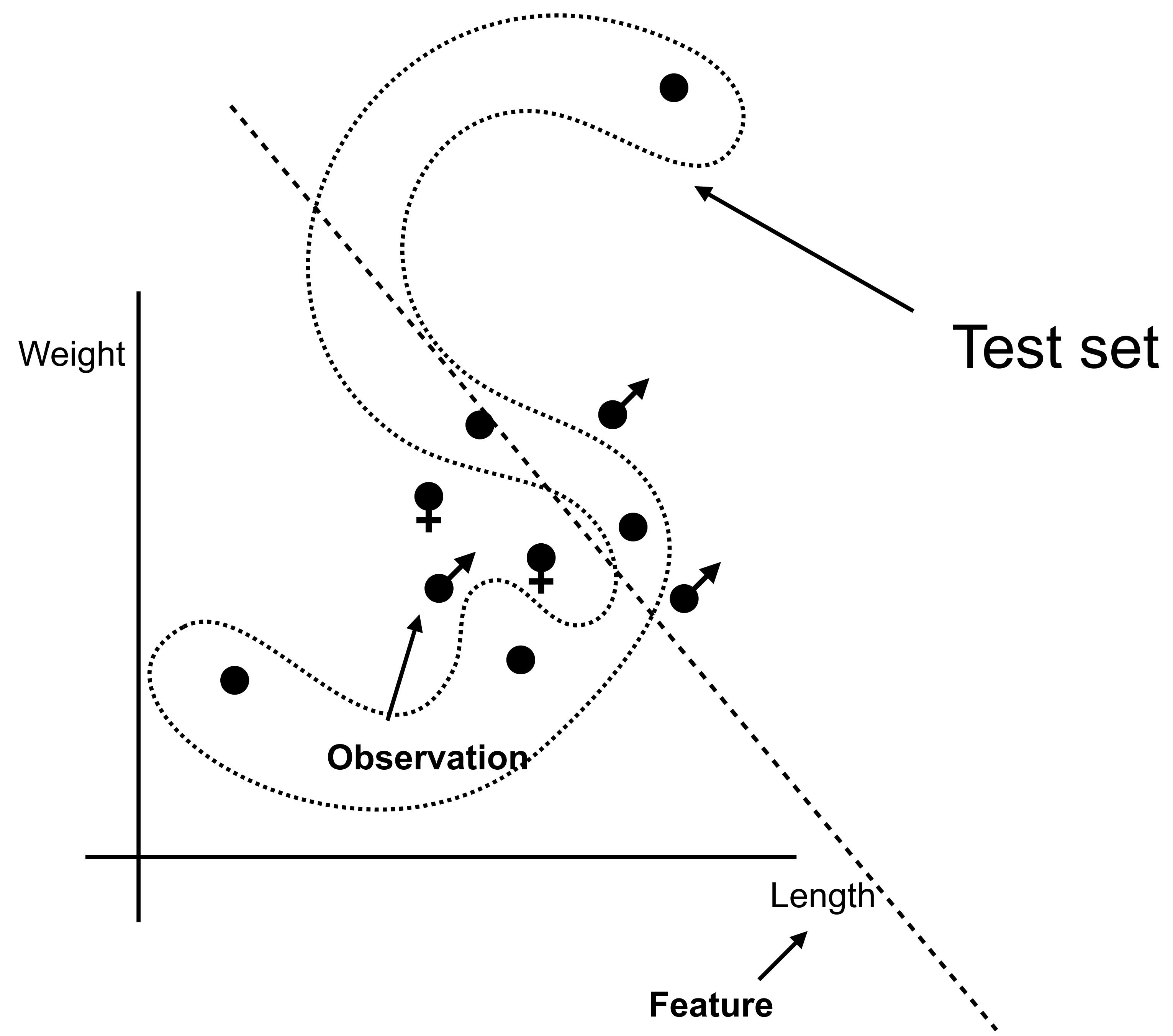


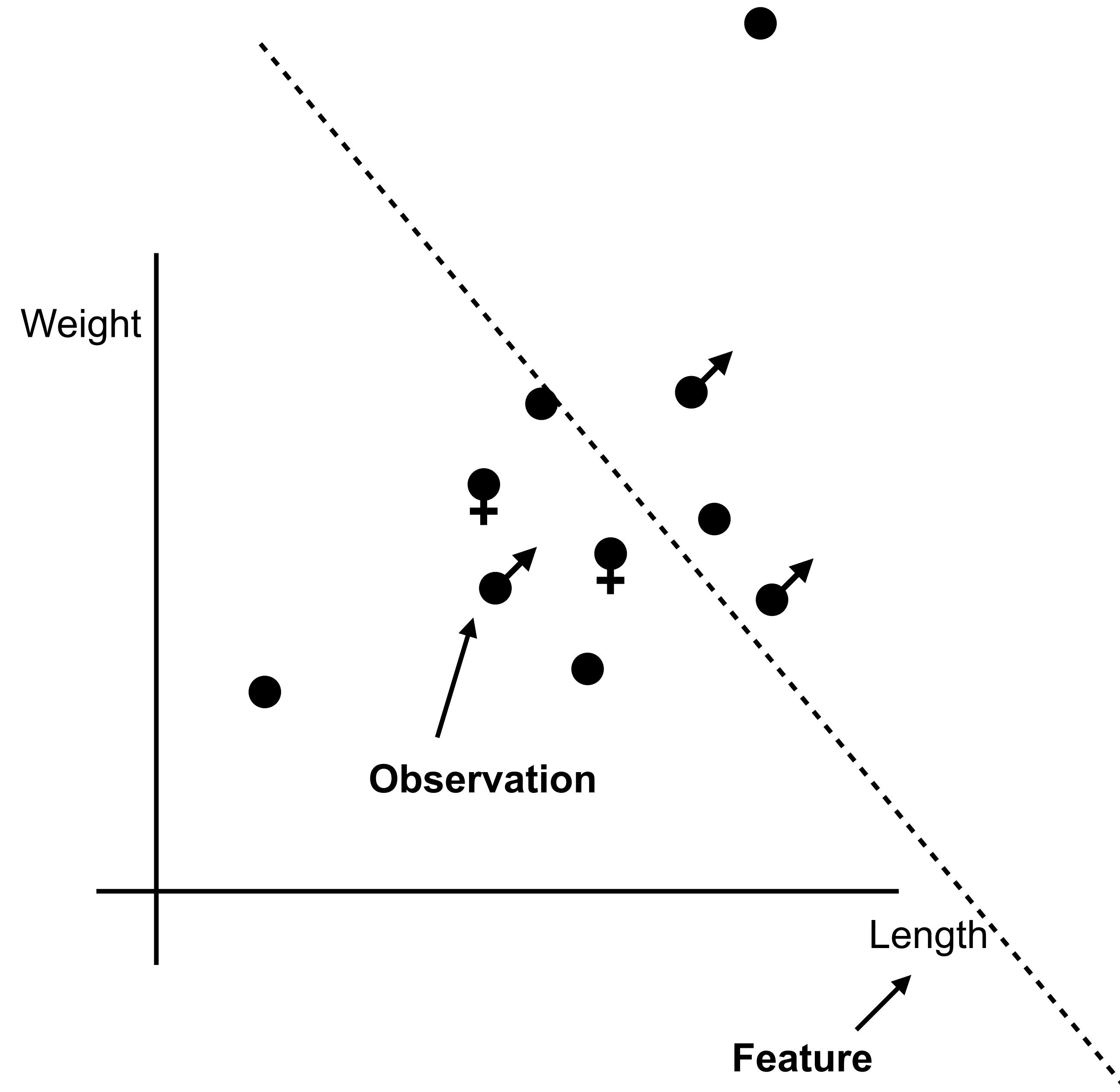














The SciPy Stack (parts we use)

Python

The SciPy Stack



numpy



Python

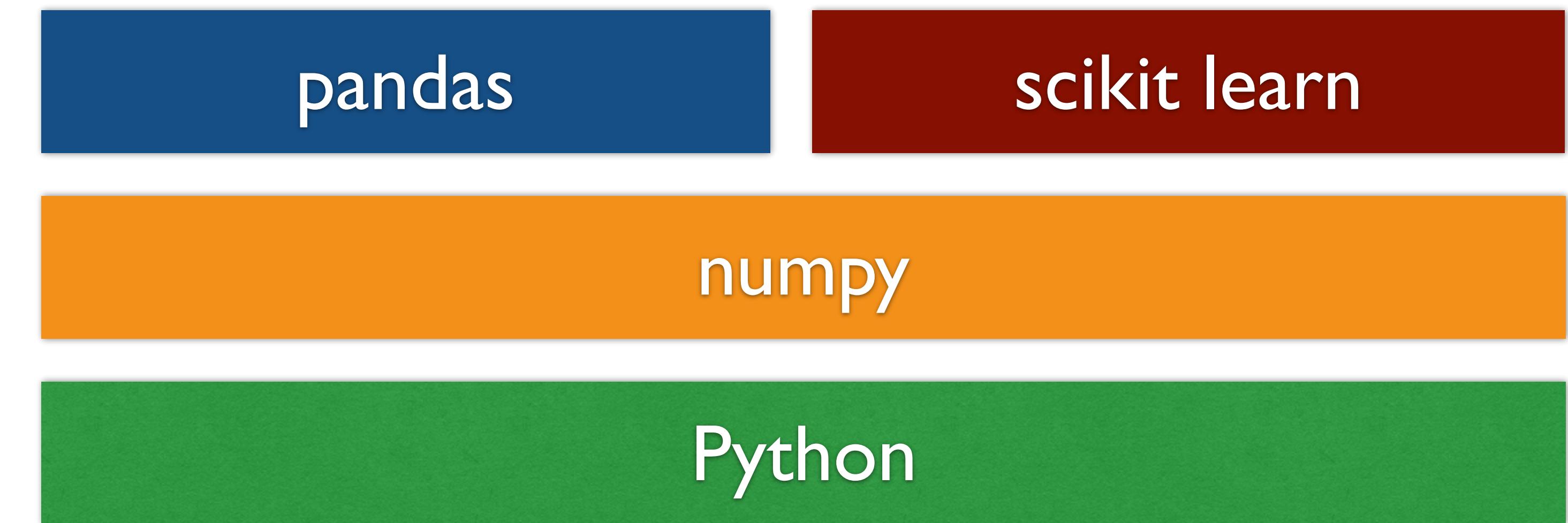
The SciPy Stack

pandas

numpy

Python

The SciPy Stack



The SciPy Stack

Jupyter Notebook

pandas

scikit learn

numpy

Python

Installation requirements

- What I did, anyway.
- Install the latest Anaconda Python 3 distribution: <https://www.continuum.io>
- Install pandas: `conda install pandas`
- Install Jupyter: `conda install jupyter`
- You are now ready to go, just `jupyter notebook`

The data science process

Data science process

- Business model
- Data cleaning
- Exploration
- Experimental setup
- Feature engineering
- Model fitting (with decision trees)

Business model



Business model

- Understand it!
- What to learn / predict?
- Use case?
- Data?

Our aim

- Predict who will order again in the next period
- Predict how much they will spend



Data cleaning

Data frame

Index	Email	...	Paid amount	Refunded amount
Order 1				
Order 2				
Order 3				
Order 4				



jupyter

Exploration

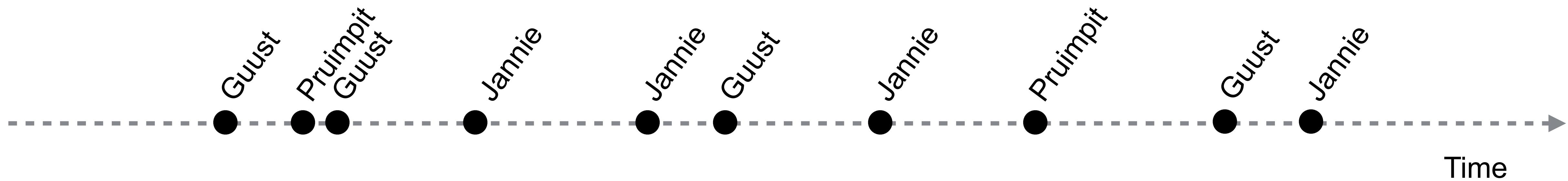


Our retail data

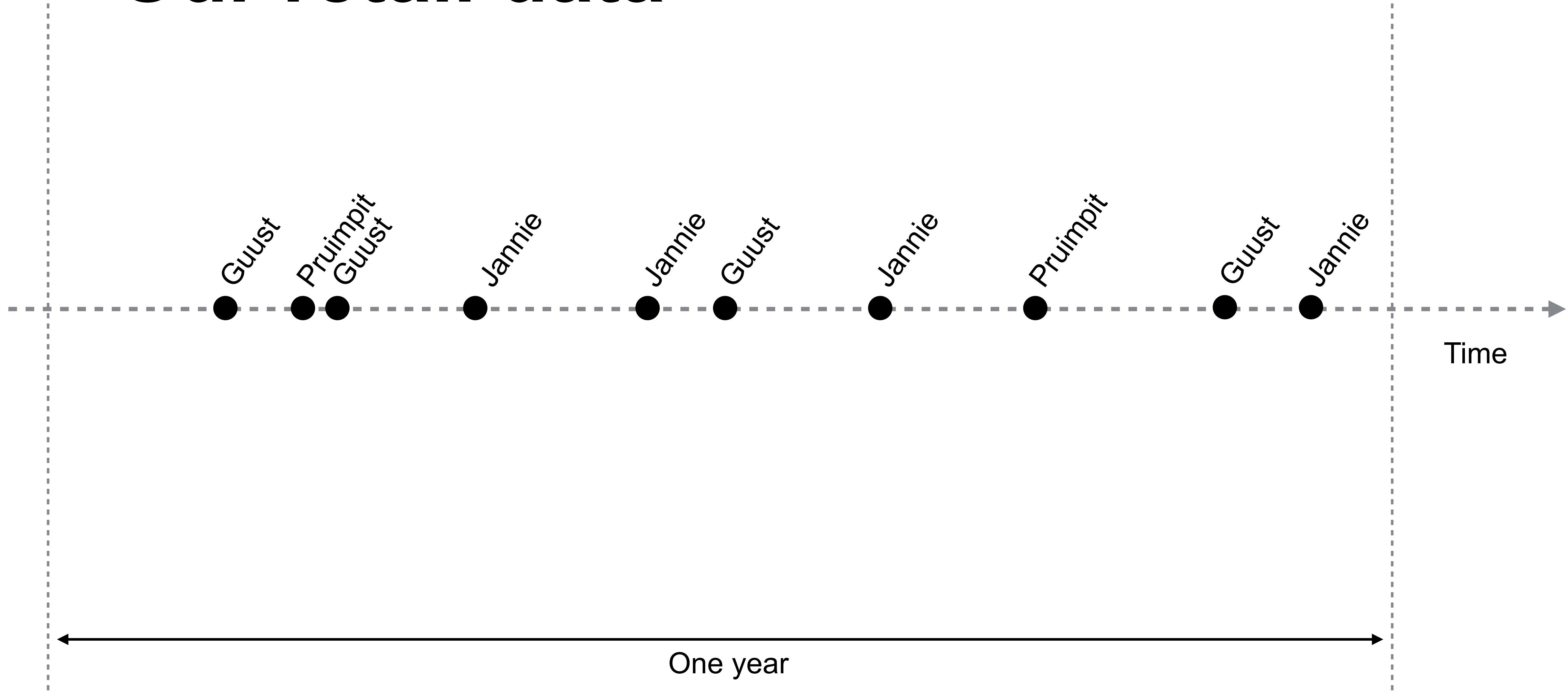


Time

Our retail data



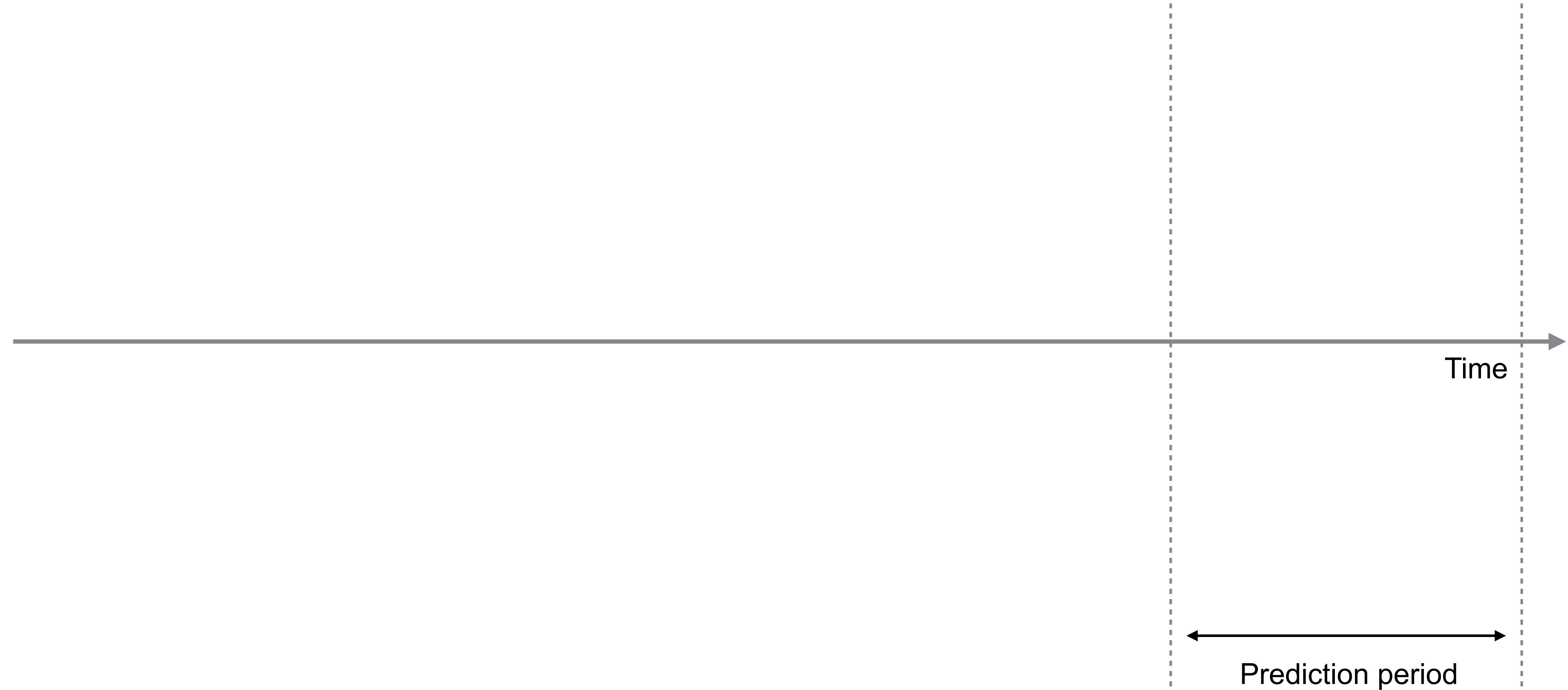
Our retail data

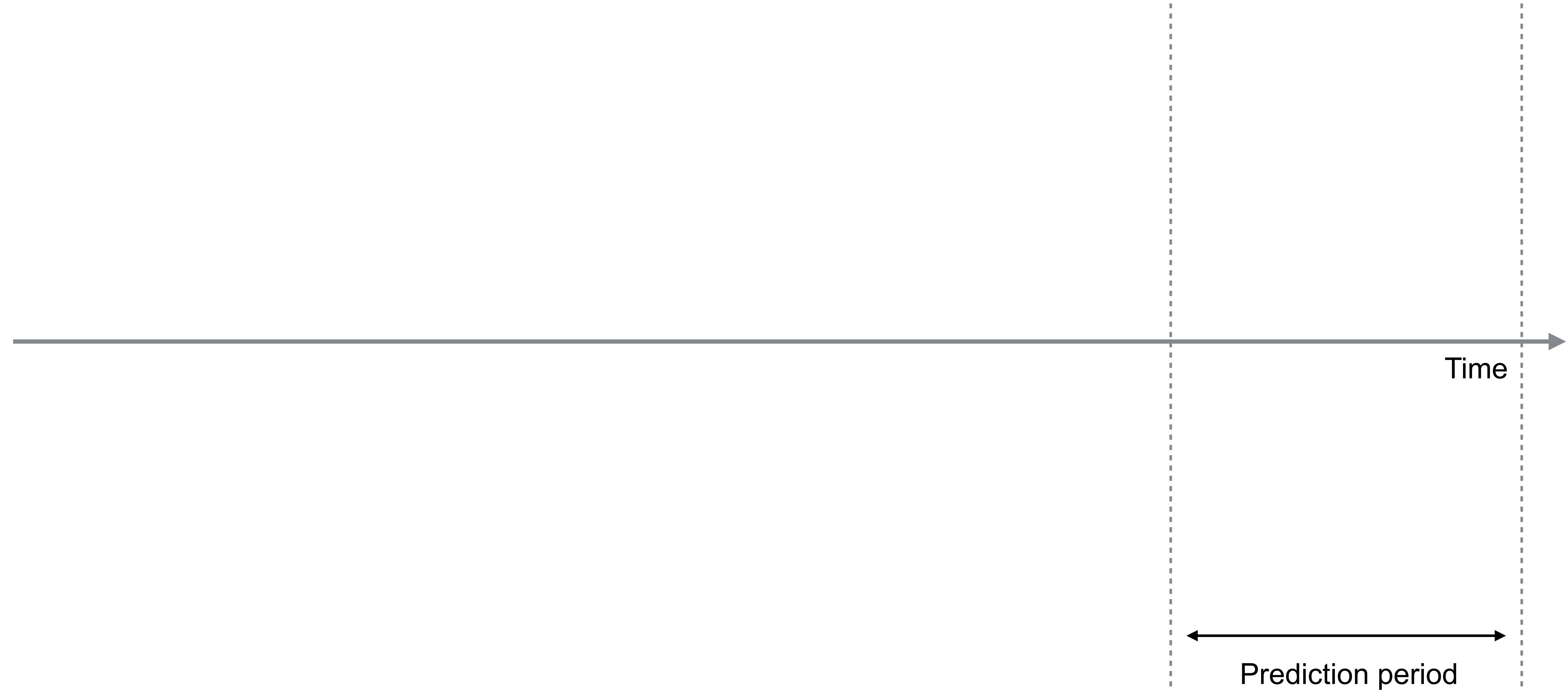


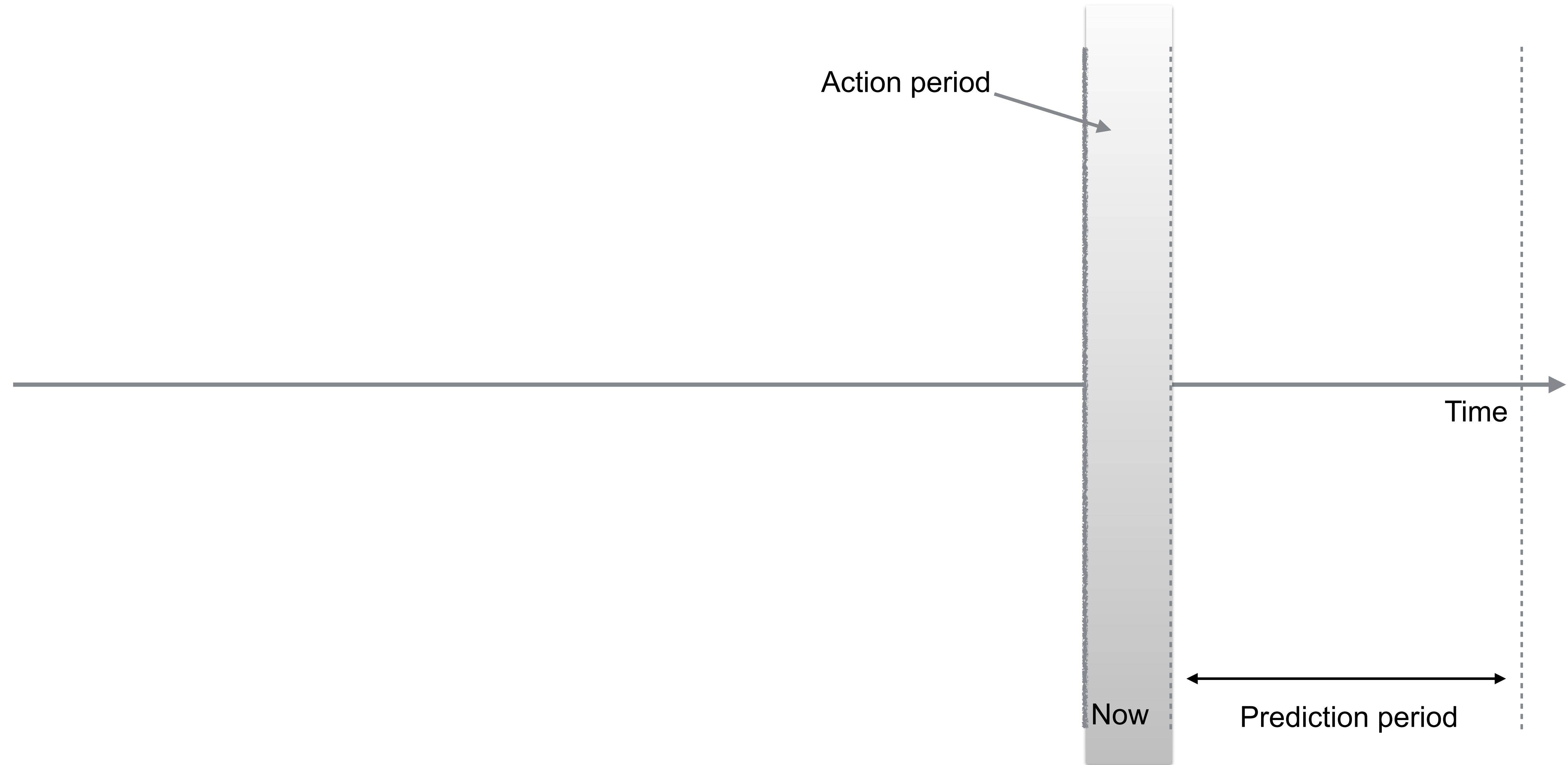
The Jupyter logo features the word "jupyter" in a dark grey sans-serif font, centered within a white circle. The circle is partially overlaid by two thick, orange curved bands forming a smile-like shape. There are also three small grey circles positioned at the top and bottom intersections of the circle and the orange bands.

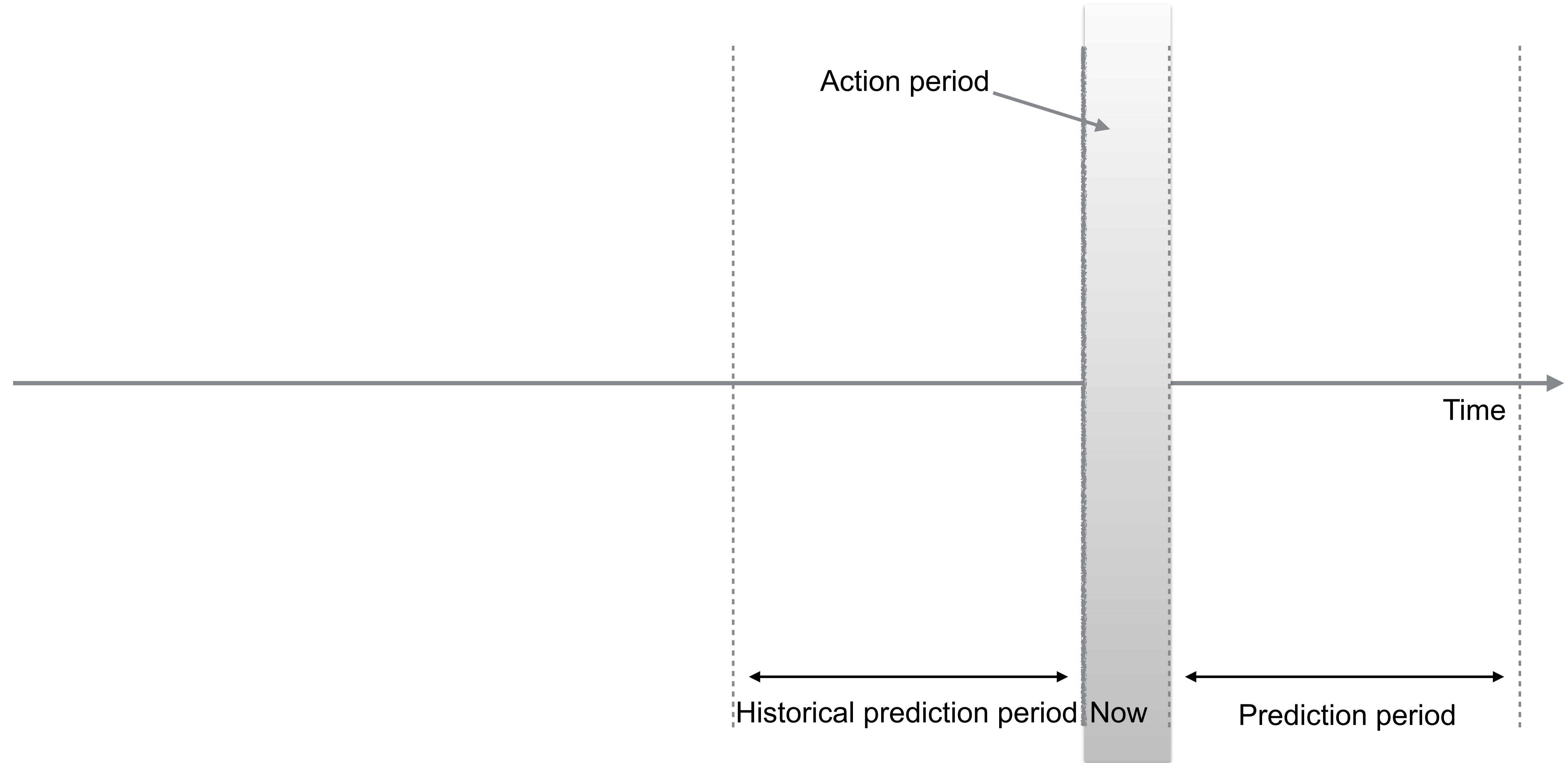
jupyter

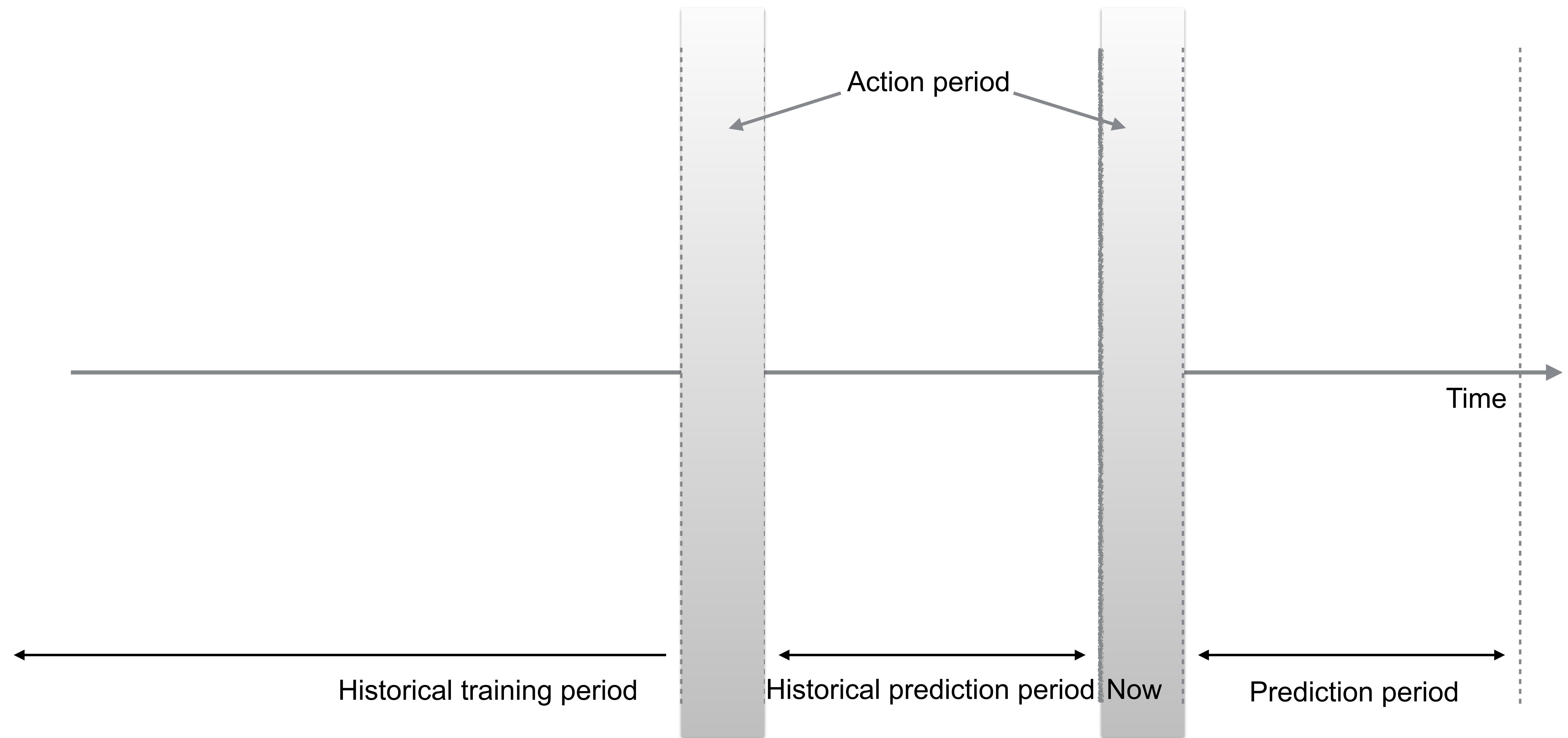
Experimental setup

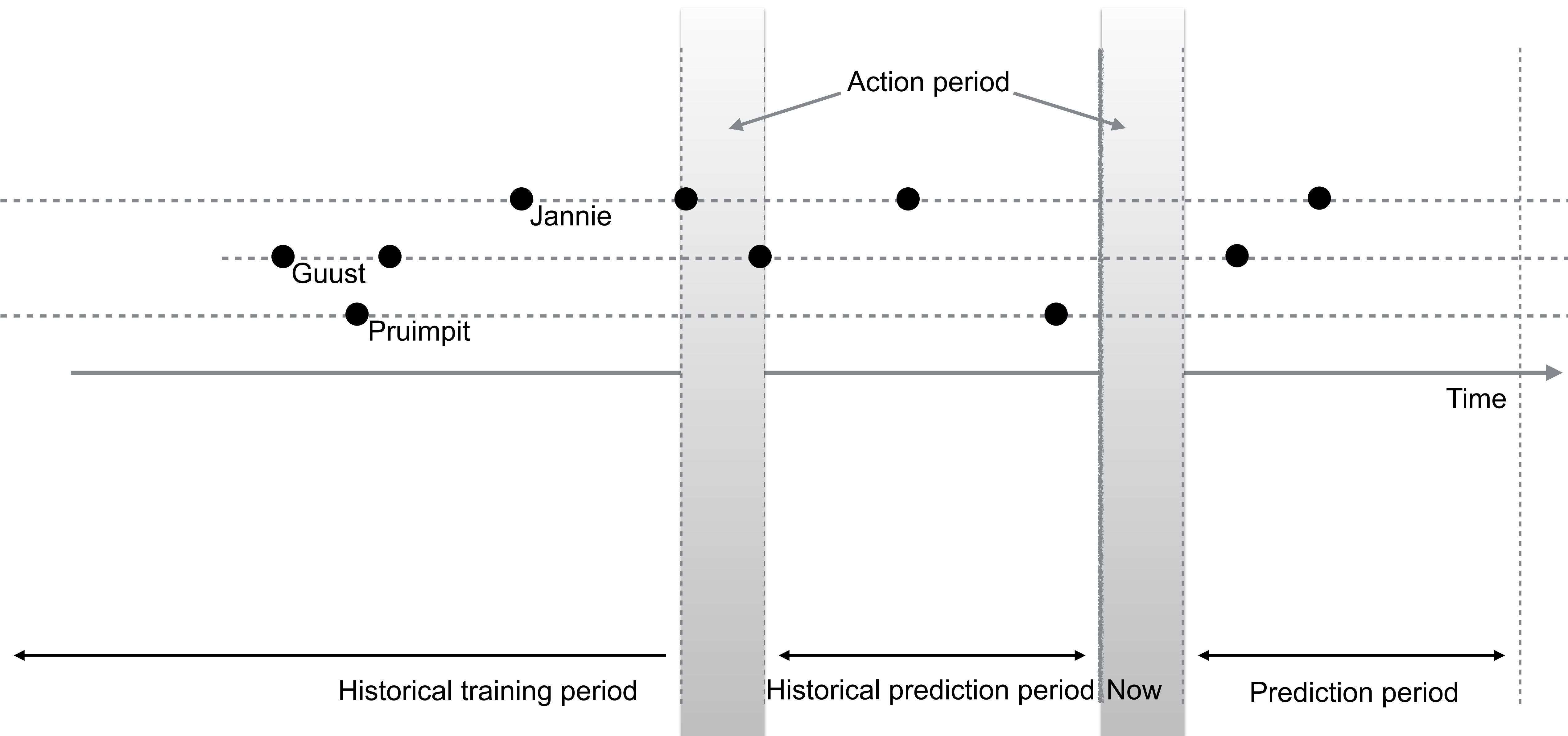


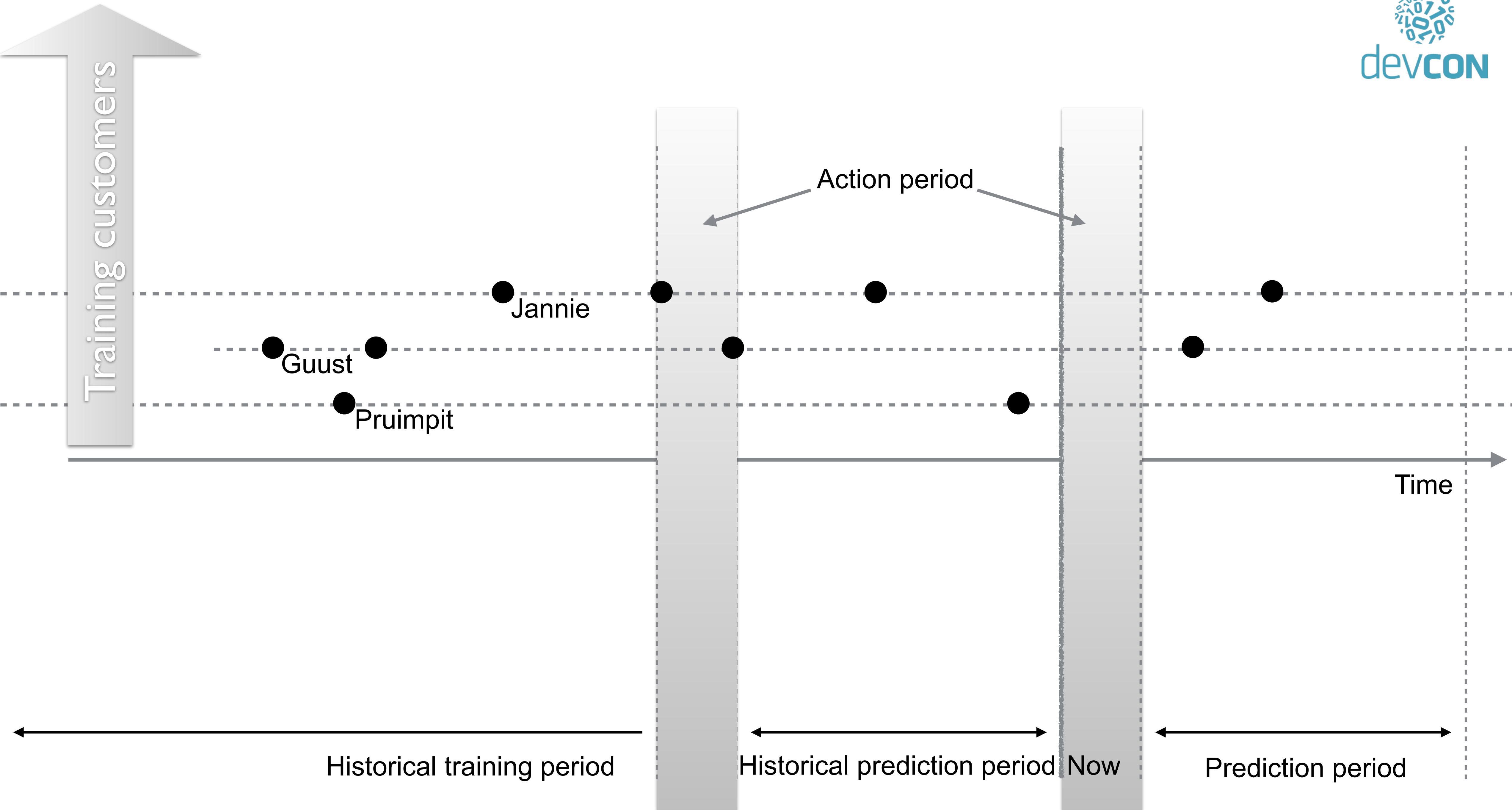


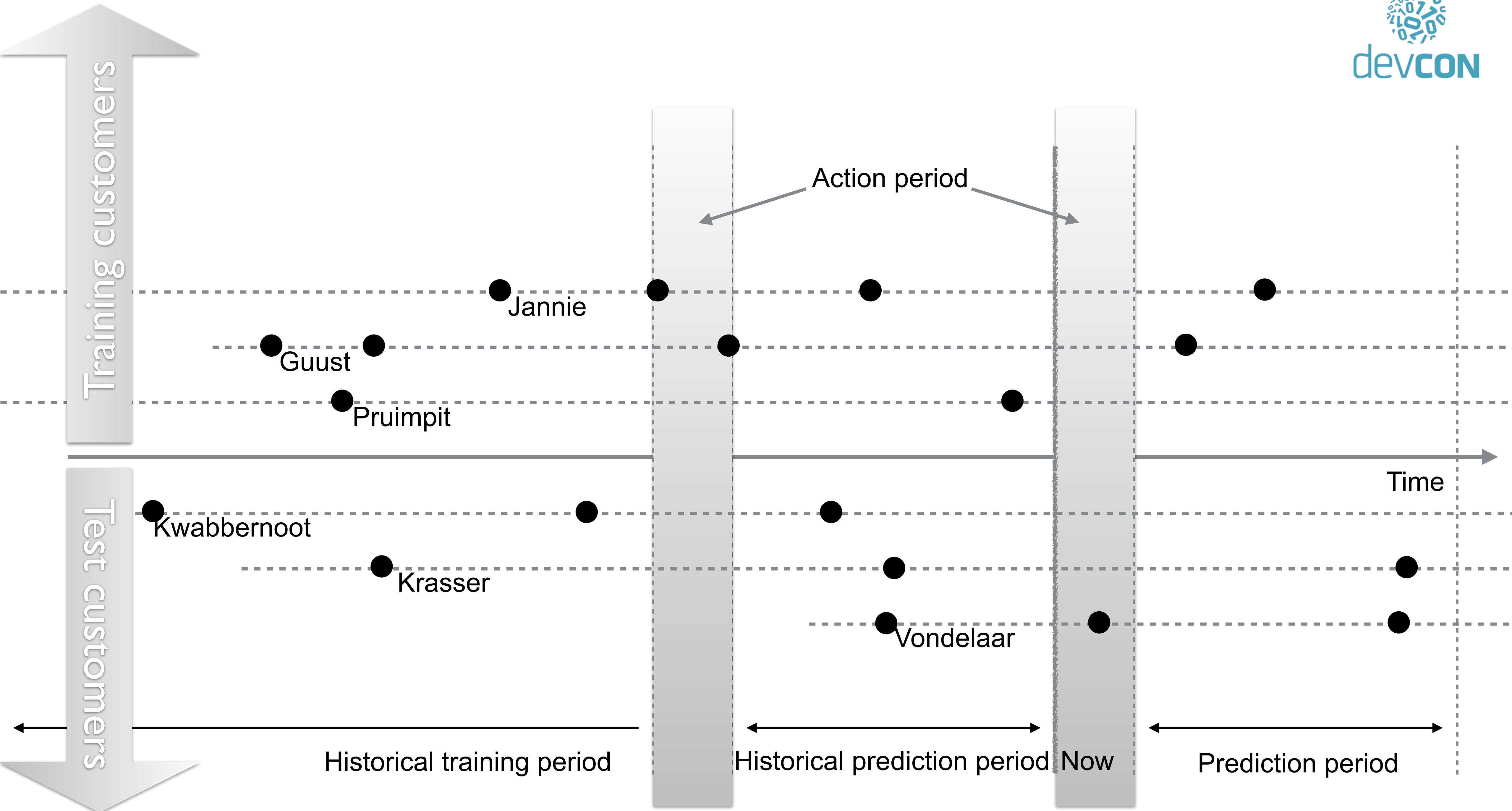


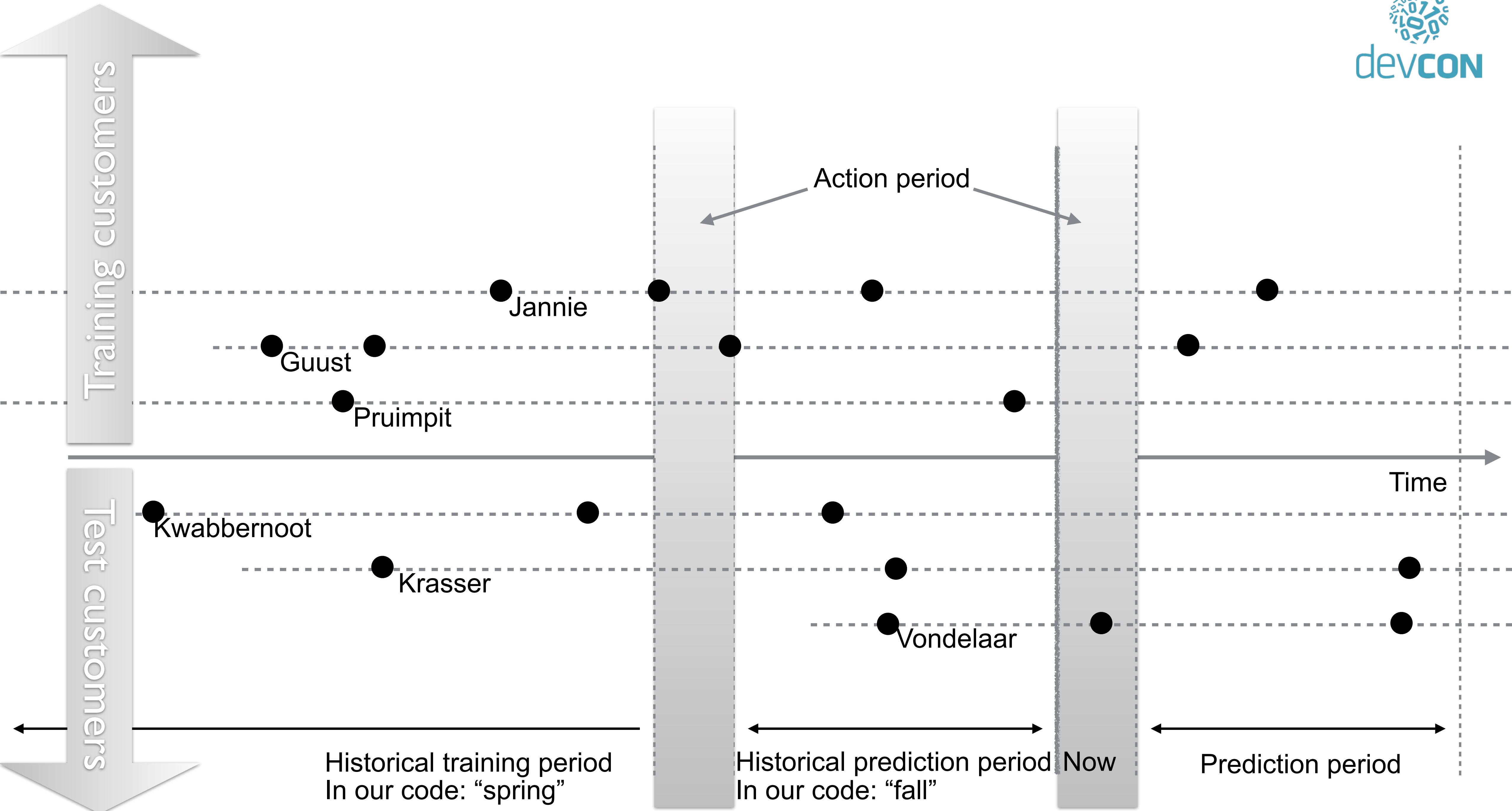












Evaluation

- How many customers that came back did we predict correctly?
- How many that left?
- How close was our prediction of what they spent?

Data frame layout

Index	Independent 1	Independent 2	Returning Customer?	Revenue
Customer 1				
Customer 2				
Customer 3				
Customer 4				

Feature engineering



Customer journey

- So let's take a look at a typical one







A large, ornate spider with long legs and a patterned body sits on its intricate web. A white rectangular sign with the words "Physical store" written in black is positioned diagonally across the spider's web. The background is a blurred, colorful bokeh effect of green and blue lights.

Physical store



A large, ornate spider with long legs and a patterned body sits on a complex web against a background of blurred green foliage. Two white, semi-transparent rectangular cards are positioned near the spider's head. The card on the left contains the text "Physical store" and the card on the right contains the text "Website".

Physical store

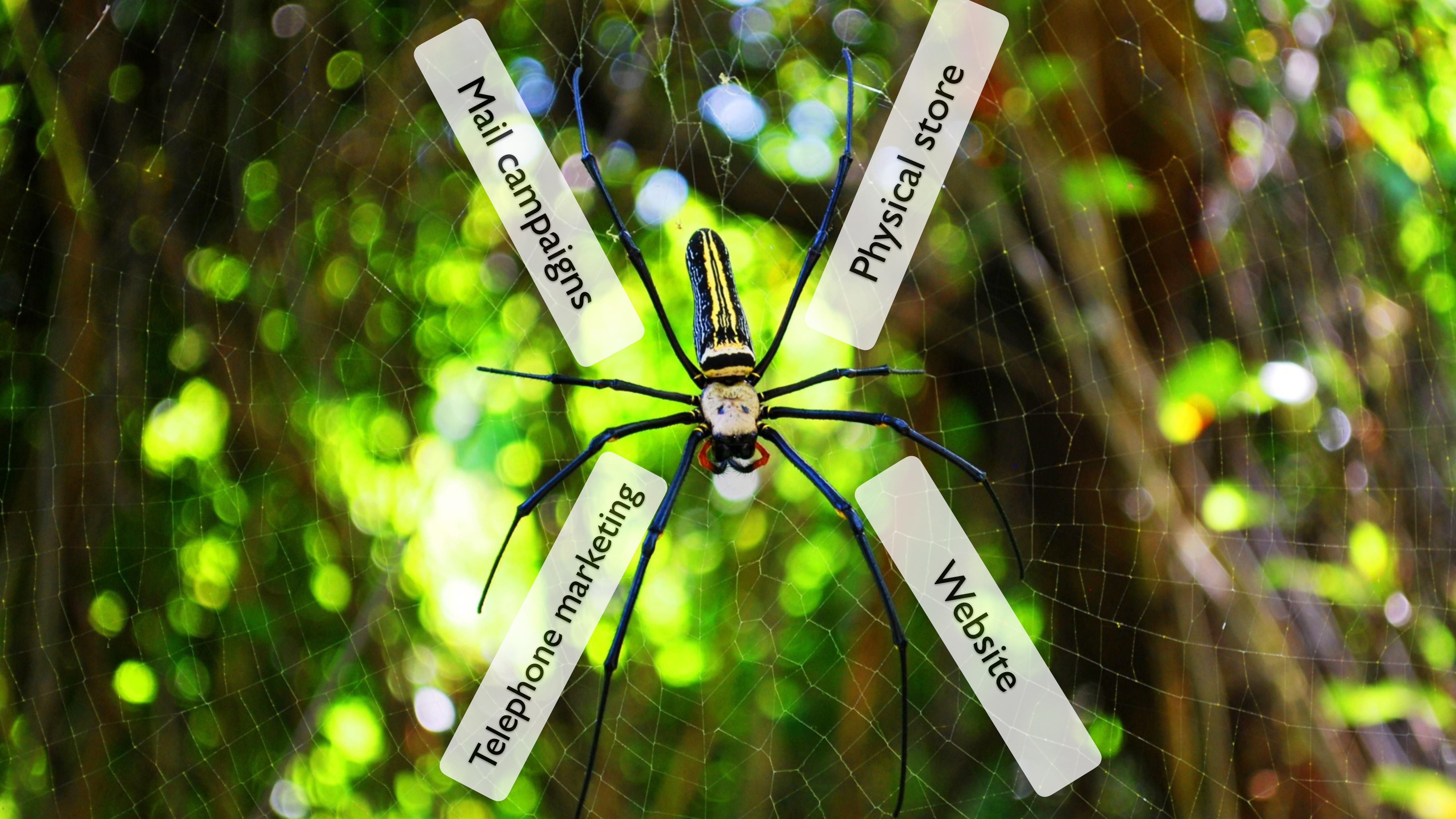
Website



Mail campaigns

Physical store

Website

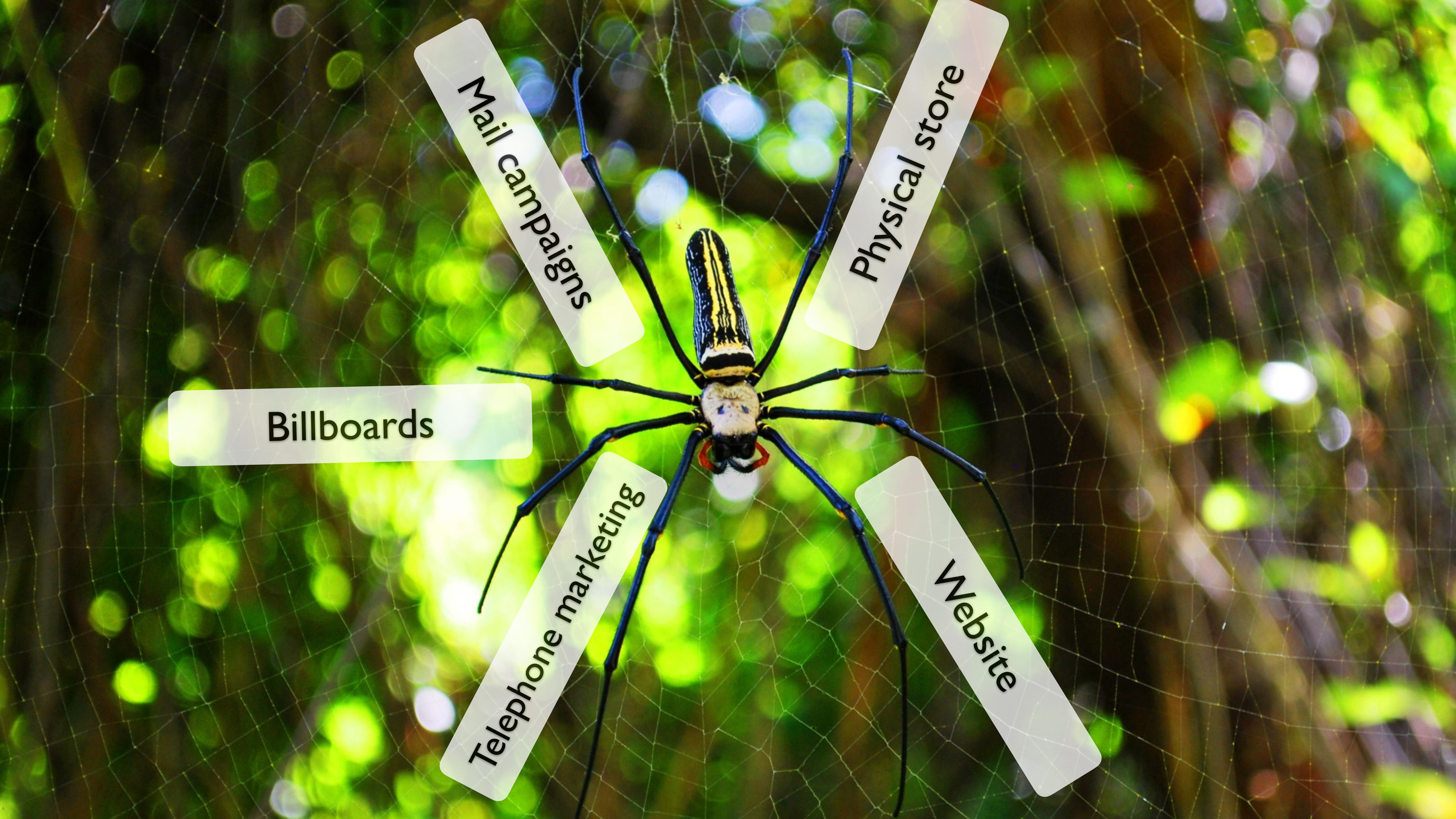


Mail campaigns

Physical store

Telephone marketing

Website



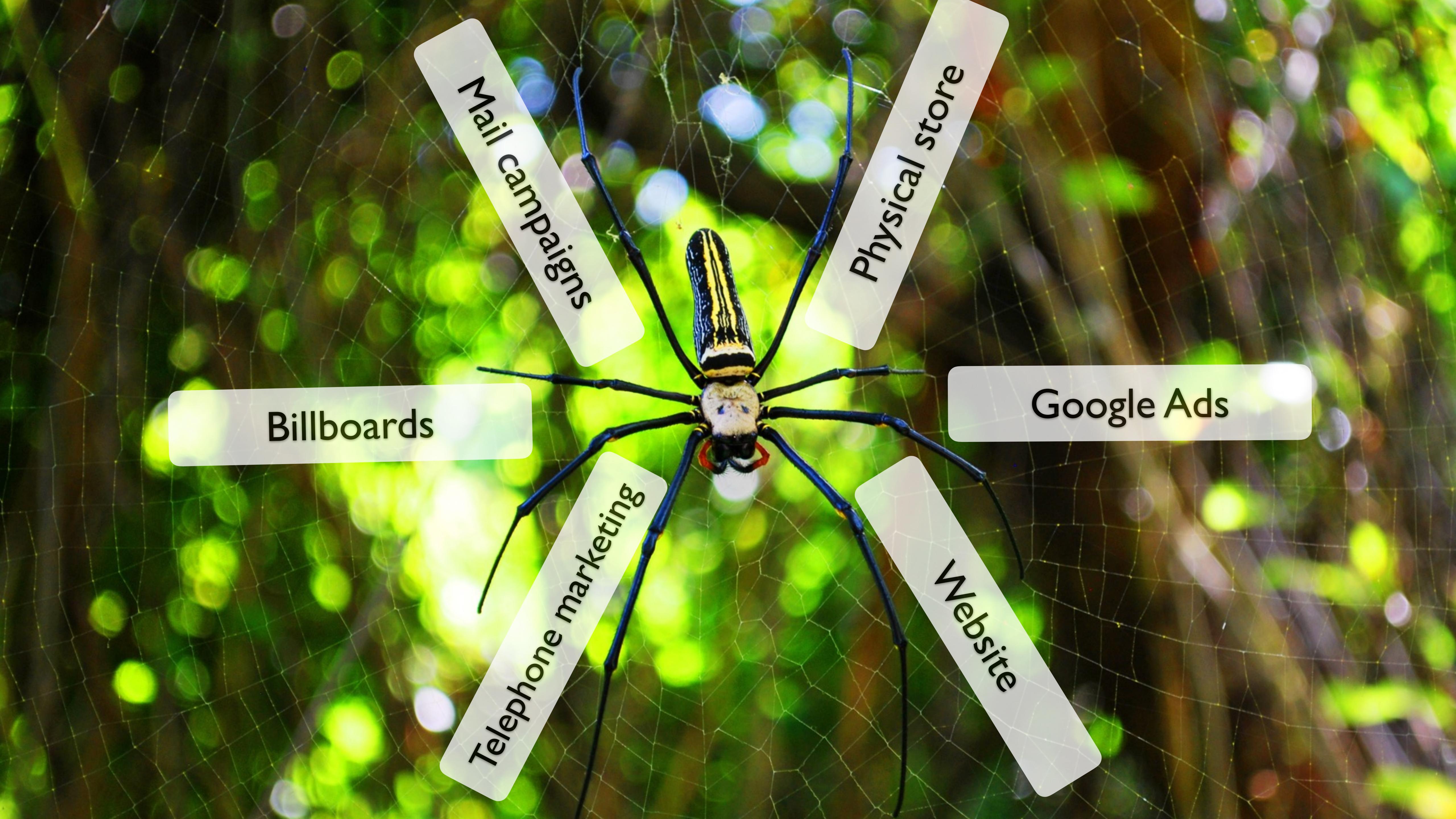
Billboards

Telephone marketing

Mail campaigns

Physical store

Website



Billboards

Telephone marketing

Mail campaigns

Website

Physical store

Google Ads

Customer buying behaviour

- How many orders
- How much per order
- How much time between orders
- Paying behaviour

Usage behaviour

- Complaints
- Newsletter opt-in / opt-out

Sharing behaviour

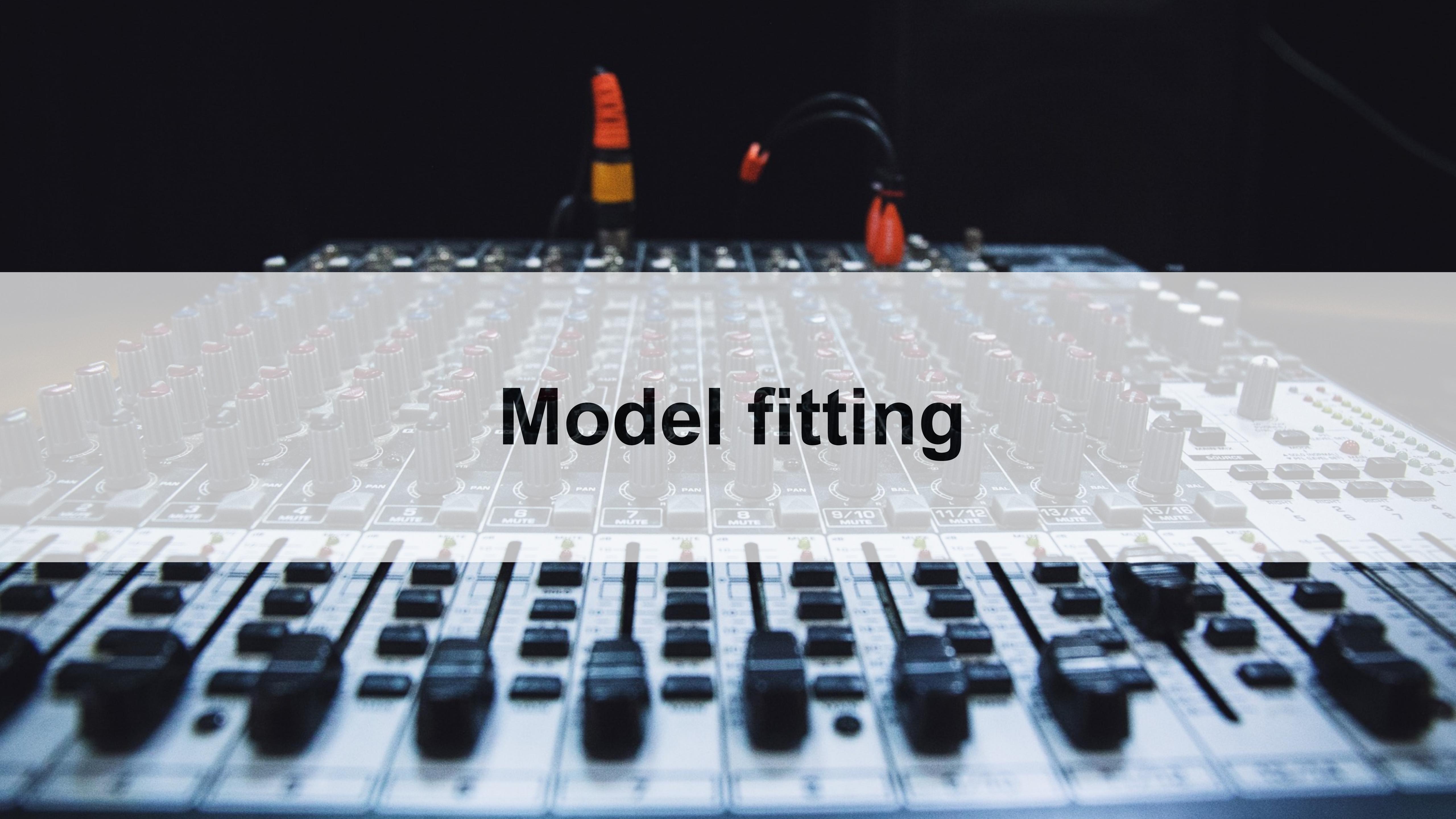
- Social media
- Member get member

Discarding behaviour

- Behaviour observed around ending of customer relationship

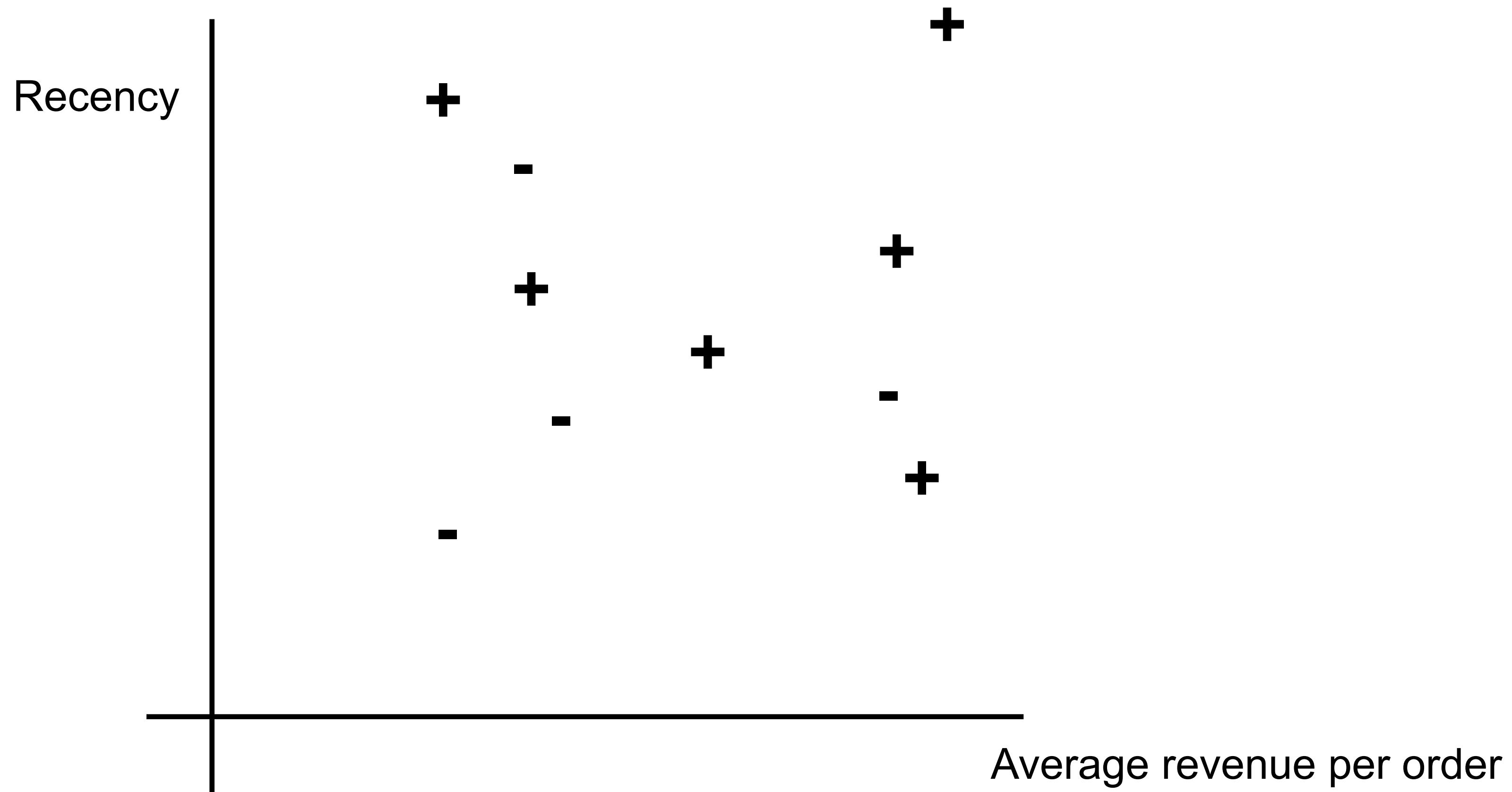


jupyter

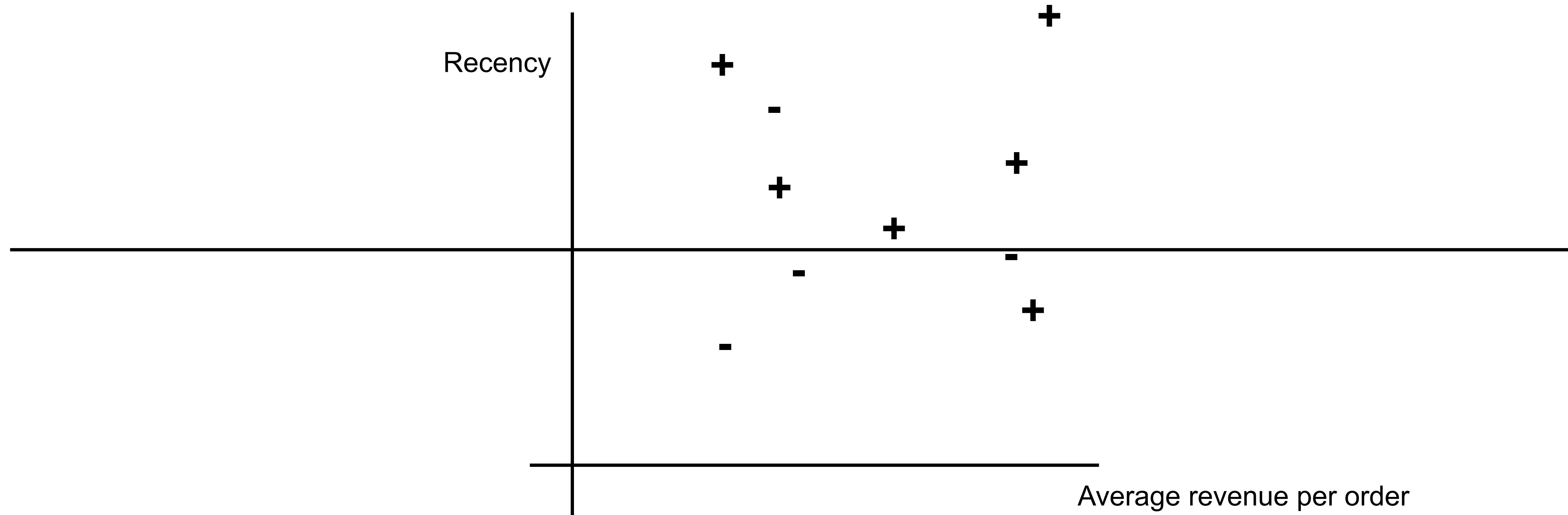


Model fitting

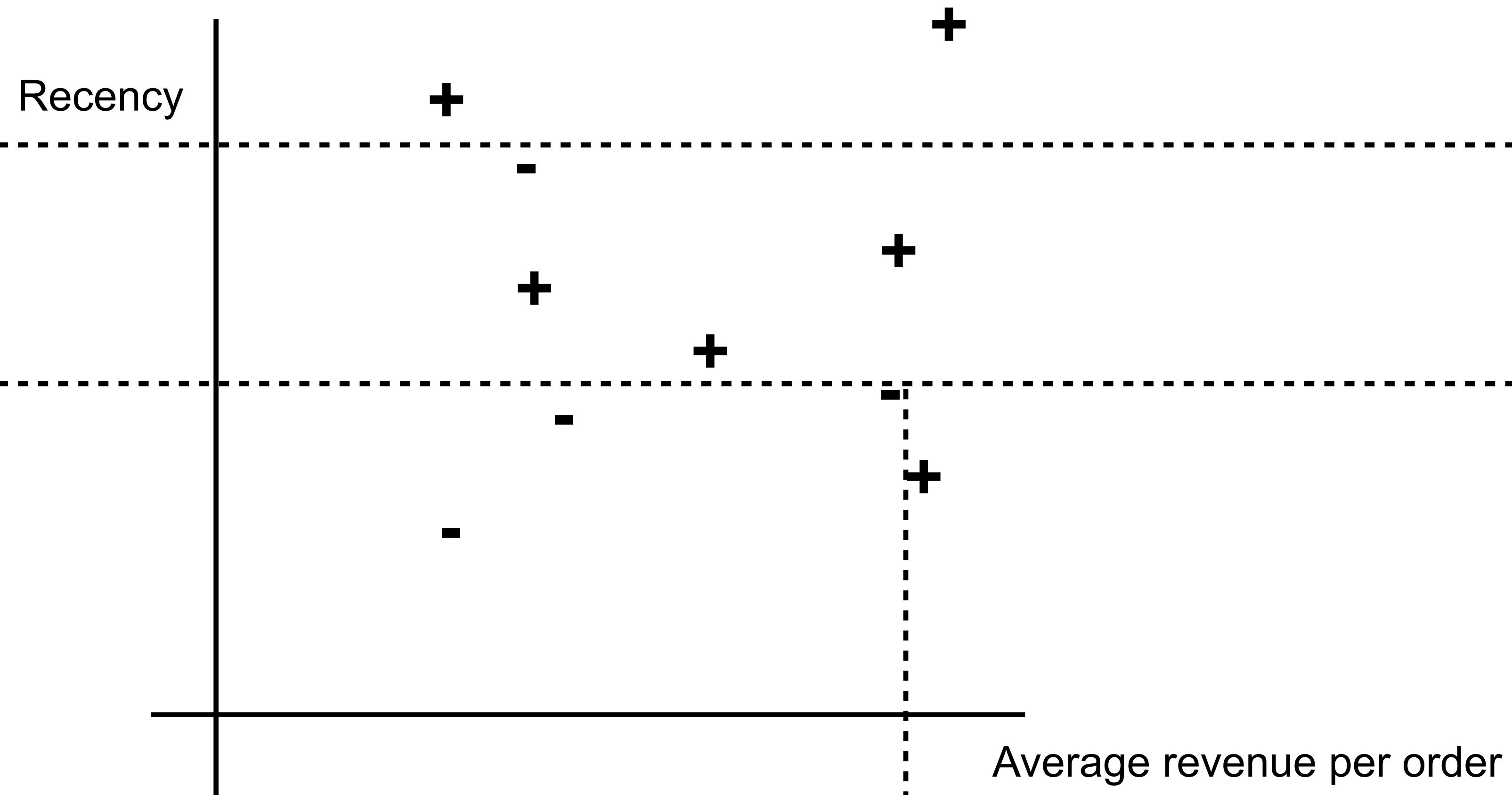
Decision trees



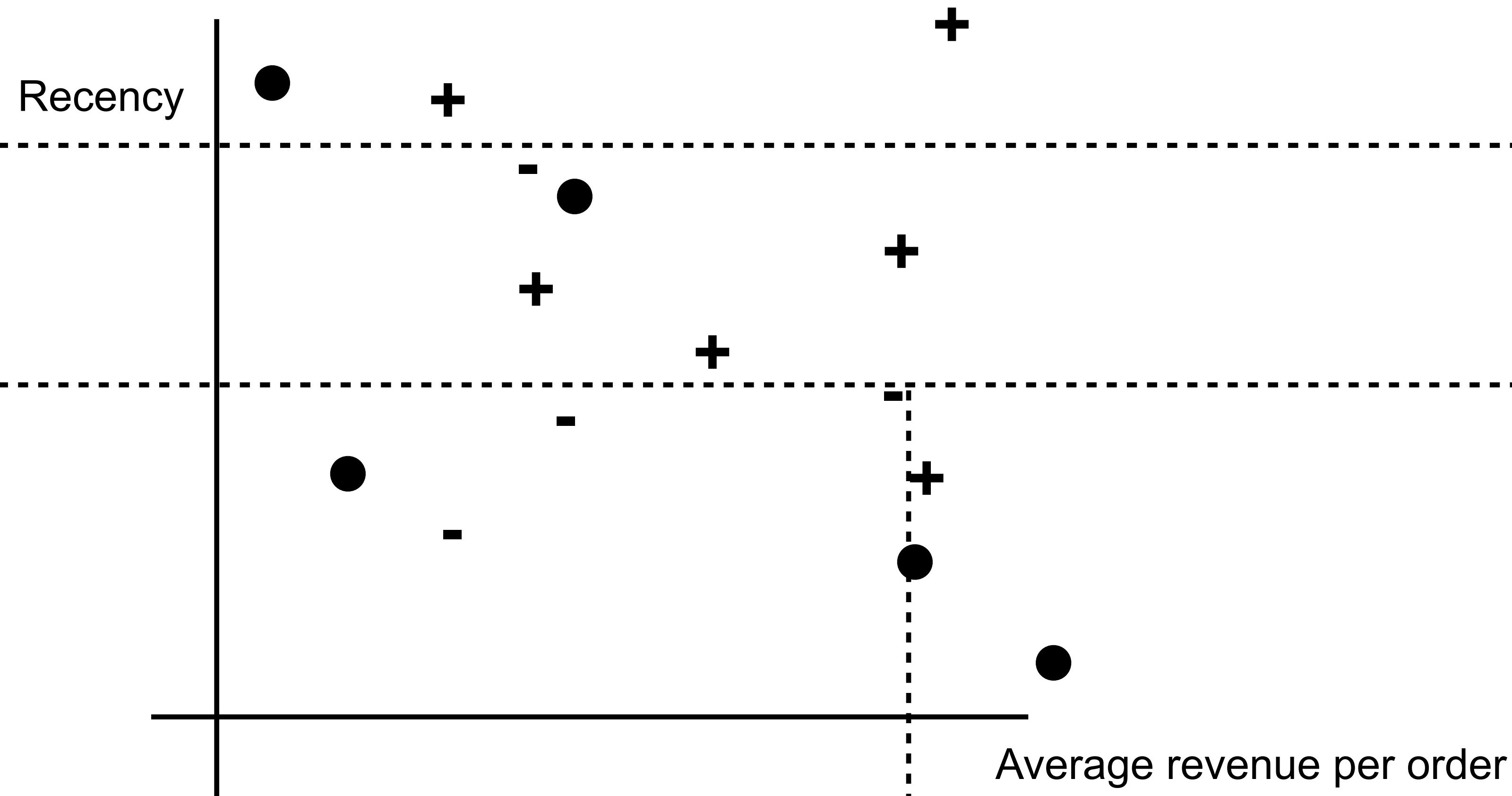
Decision trees



Decision trees



Decision trees



The Jupyter logo features the word "jupyter" in a dark grey sans-serif font, centered within a white circle. The circle is partially obscured by two thick, orange curved bands that intersect at the bottom. Small dark grey circles are positioned at the intersections of the curves.

jupyter

Take away points

- Understand business model
- Careful experimental setup
- Interactive data science:
 - Fun
 - Productive

Thank you

Any

Questions?



is powered by



My world is:

Nou moe?! - Guust Flater

All our worlds > www.luminis.eu

Pearson correlation

$$r = \frac{\sum_i (x_i - \bar{x})(y_i - \bar{y}) / (n - 1)}{s(x)s(y)}.$$

Precision

- $\text{tp} / (\text{fp} + \text{tp})$

Recall

- $\text{tp} / \text{tp} + \text{fn}$

F β -score

- $(1 + \beta^2) (\text{precision} \times \text{recall}) / (\beta^2 \text{precision} + \text{recall})$

Gini impurity in node m

$$H(X_m) = \sum_k p_{mk}(1 - p_{mk})$$

Summing over class labels k

Mean square error

$$\text{MSE} = \frac{1}{n} \sum_{i=1}^n (\hat{Y}_i - Y_i)^2$$

R²: coefficient of determination

$$R^2(y, \hat{y}) = 1 - \frac{\sum_{i=0}^{n_{\text{samples}}-1} (y_i - \hat{y}_i)^2}{\sum_{i=0}^{n_{\text{samples}}-1} (y_i - \bar{y})^2}$$