

Overview

Understanding different types of ML algorithms and use cases

Working with numerical and categorical data

Standardization of numerical data input into an ML model

Working with text, representing text data in numerical form

Representing pixel intensities and extracting features from images

Prerequisites and Course Outline

Beginner course on building ML
models using scikit-learn

Comfortable with Python
programming



Software and Skills

Be very comfortable programming in Python (Python 3)

Be comfortable working with Jupyter notebooks

Understand some basics of machine learning



Course Outline

Processing data

- Data preparation, representing text as numbers, representing images as matrices

Building specialized regression models

- Lasso and Ridge regression, Support Vector Regression

Building SVM and gradient boosting models

- Support Vector Machines for text and image classification, Gradient Boosting for regression

Clustering and dimensionality reduction

- Mean-shift clustering, Principal Components Analysis

Understanding Machine Learning

A machine learning algorithm is an algorithm that is able to learn from data

Machine Learning



**Work with a huge
maze of data**

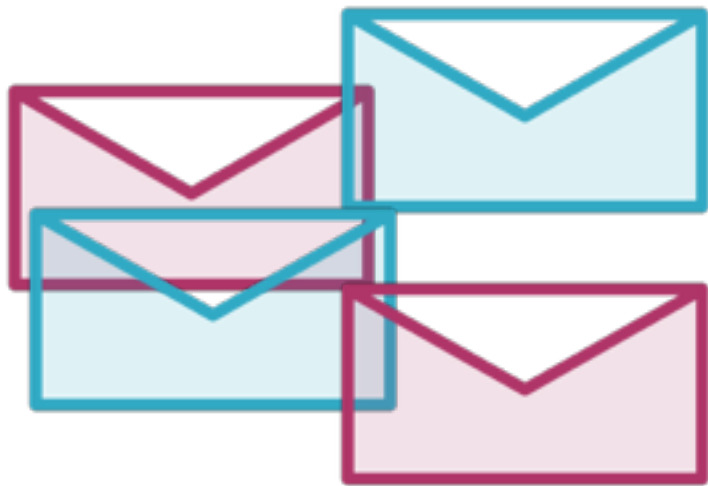


Find patterns



**Make intelligent
decisions**

Machine Learning



Emails on a server

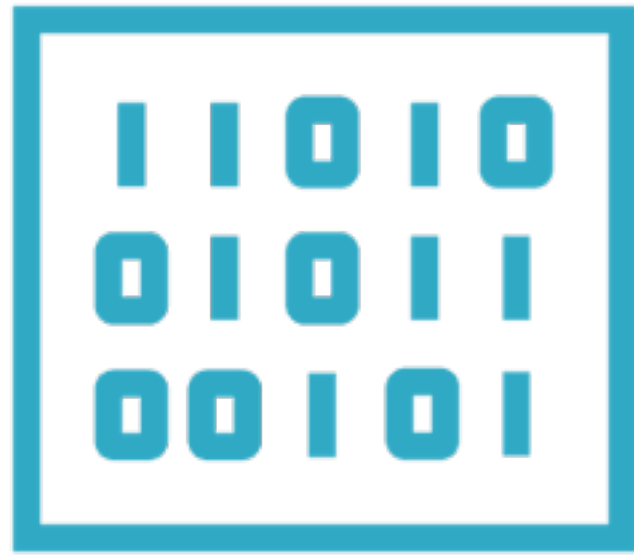


Spam or Ham?



Trash or Inbox

Machine Learning



Images represented
as pixels



Identify edges,
colors, shapes



A photo of a
little girl

Types of Machine Learning Problems



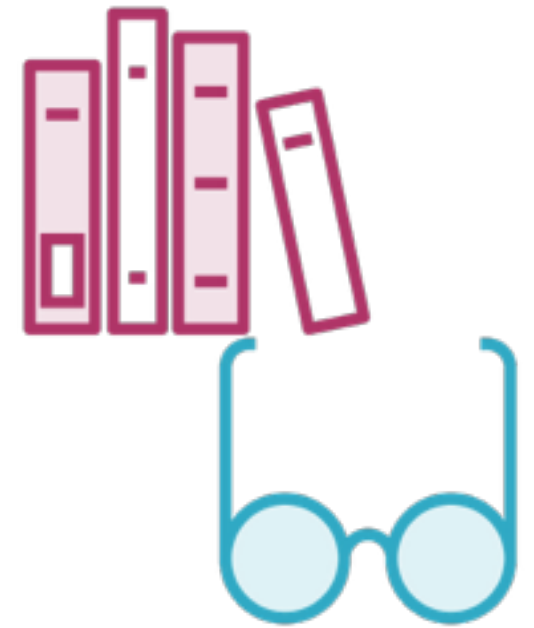
Classification



Regression



Clustering

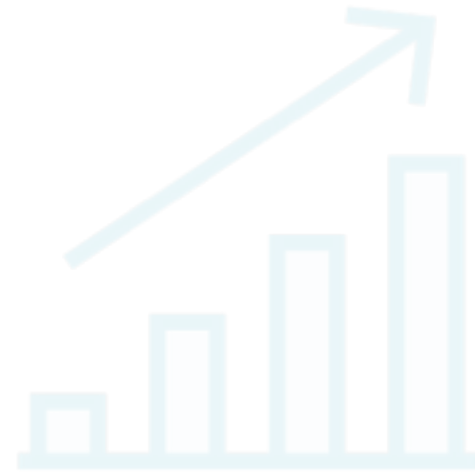


Rule-extraction

Types of Machine Learning Problems



Classification



Regression

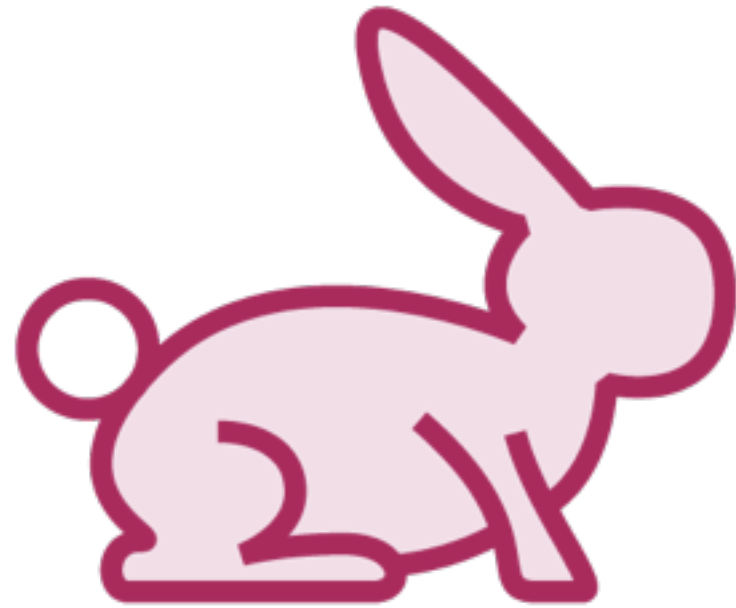


Clustering



Rule-extraction

Whales: Fish or Mammals?



Mammals

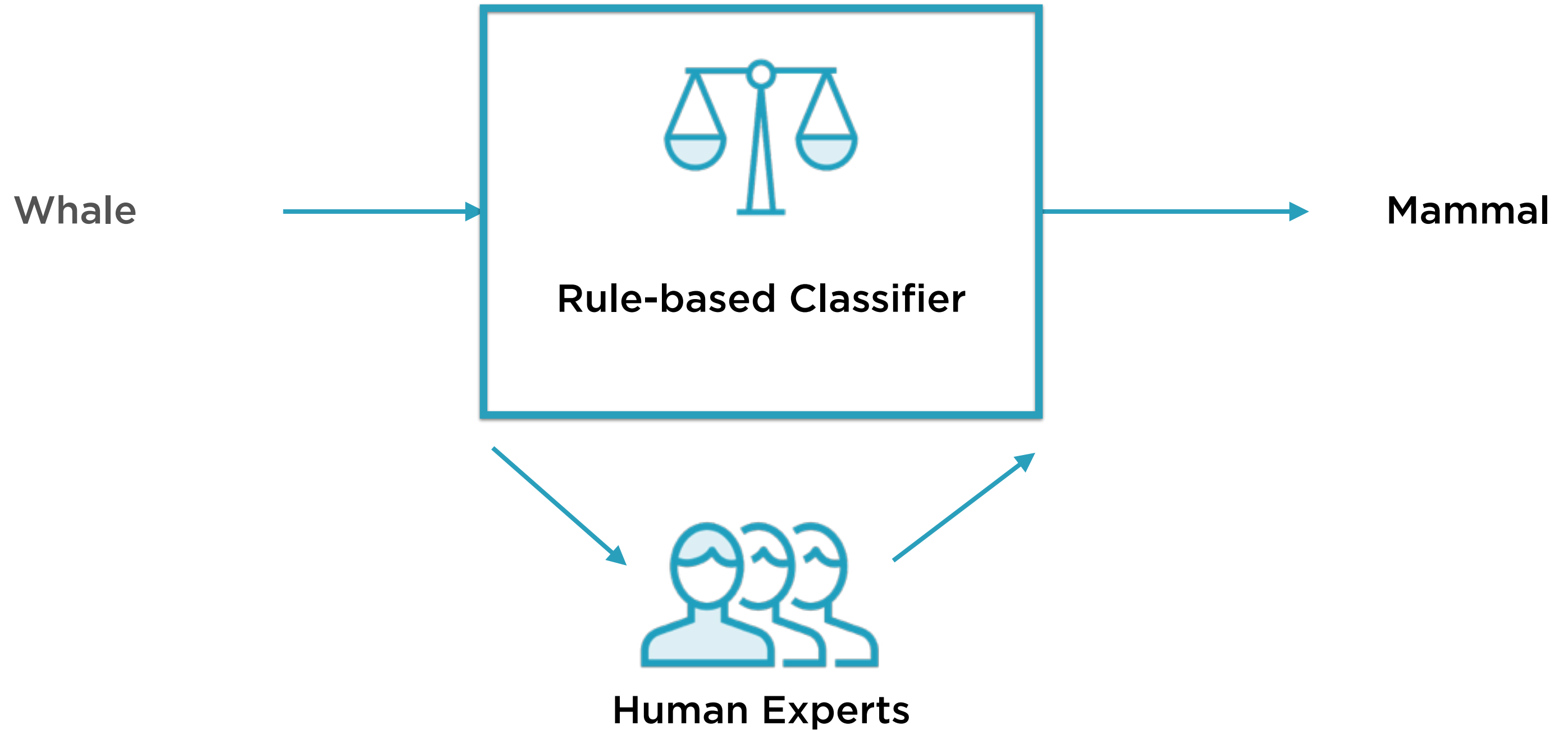
Members of the infraorder
Cetacea



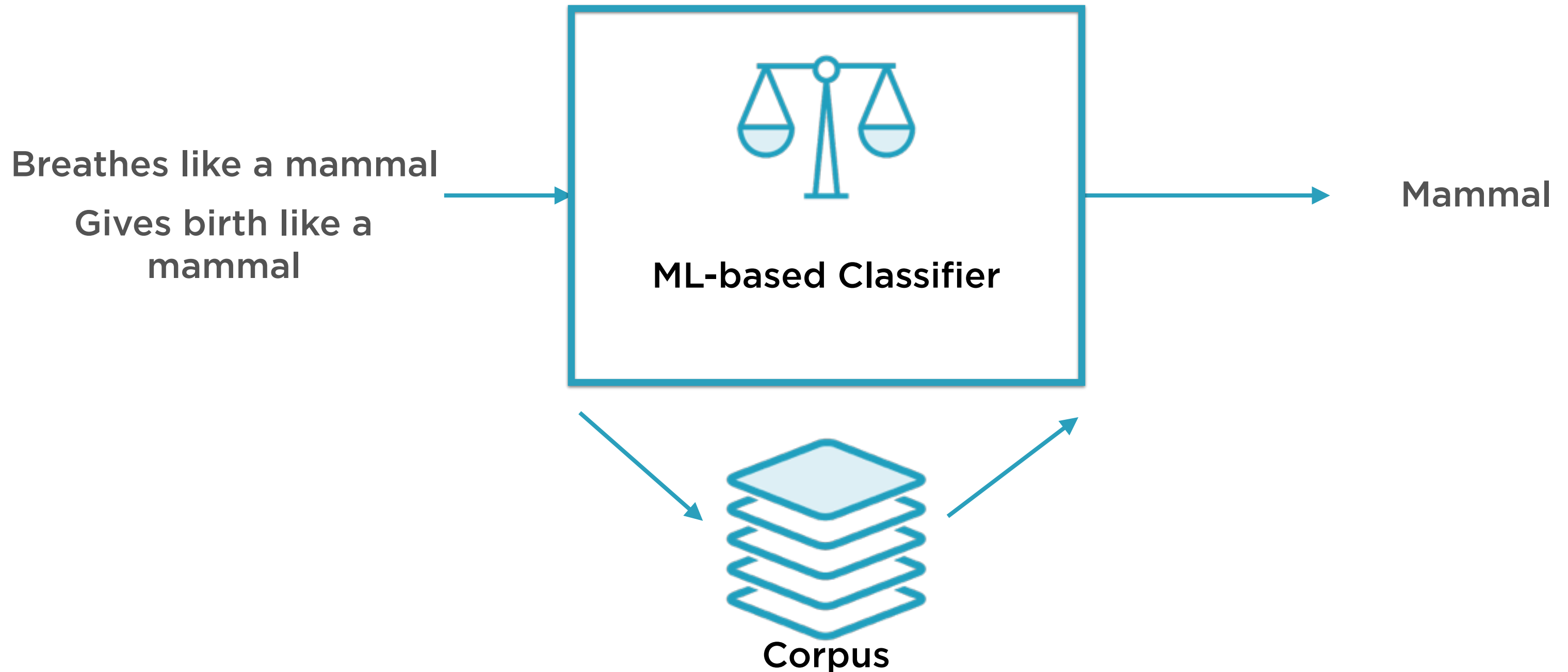
Fish

Look like fish, swim like fish,
move with fish

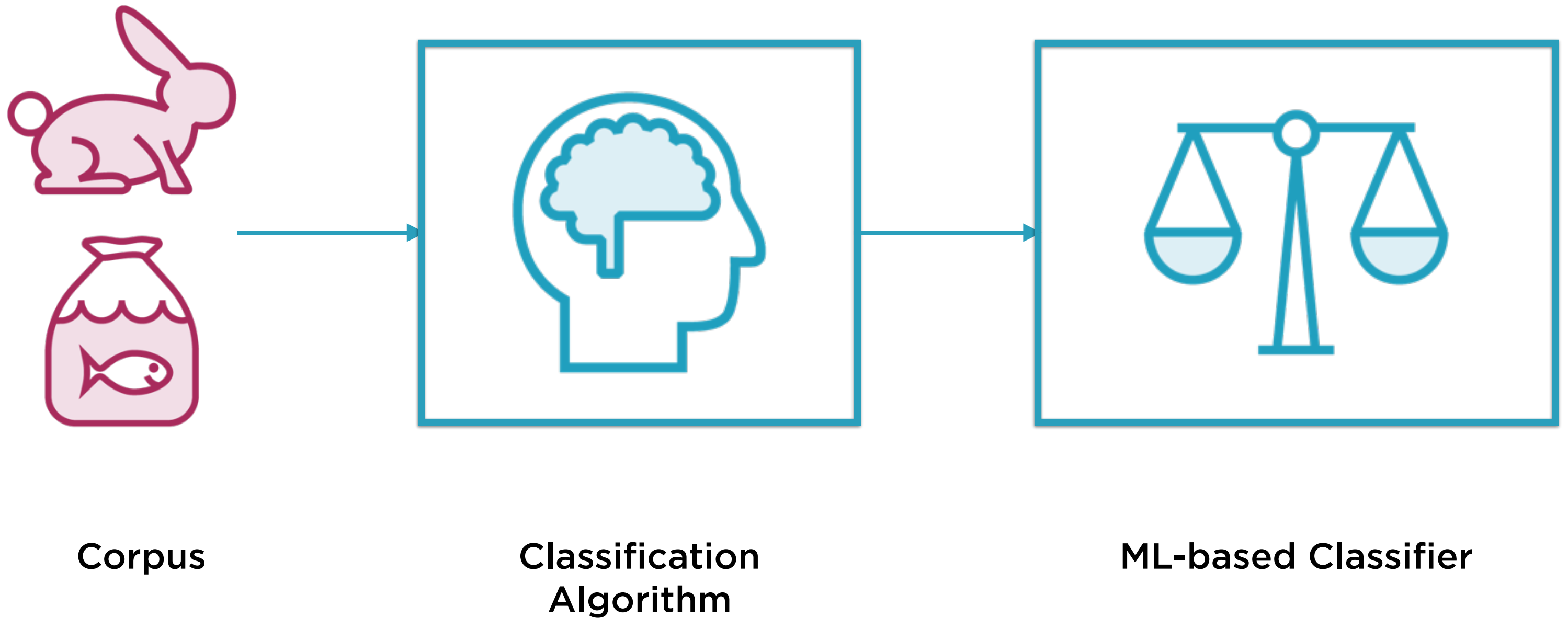
Rule-based Binary Classifier



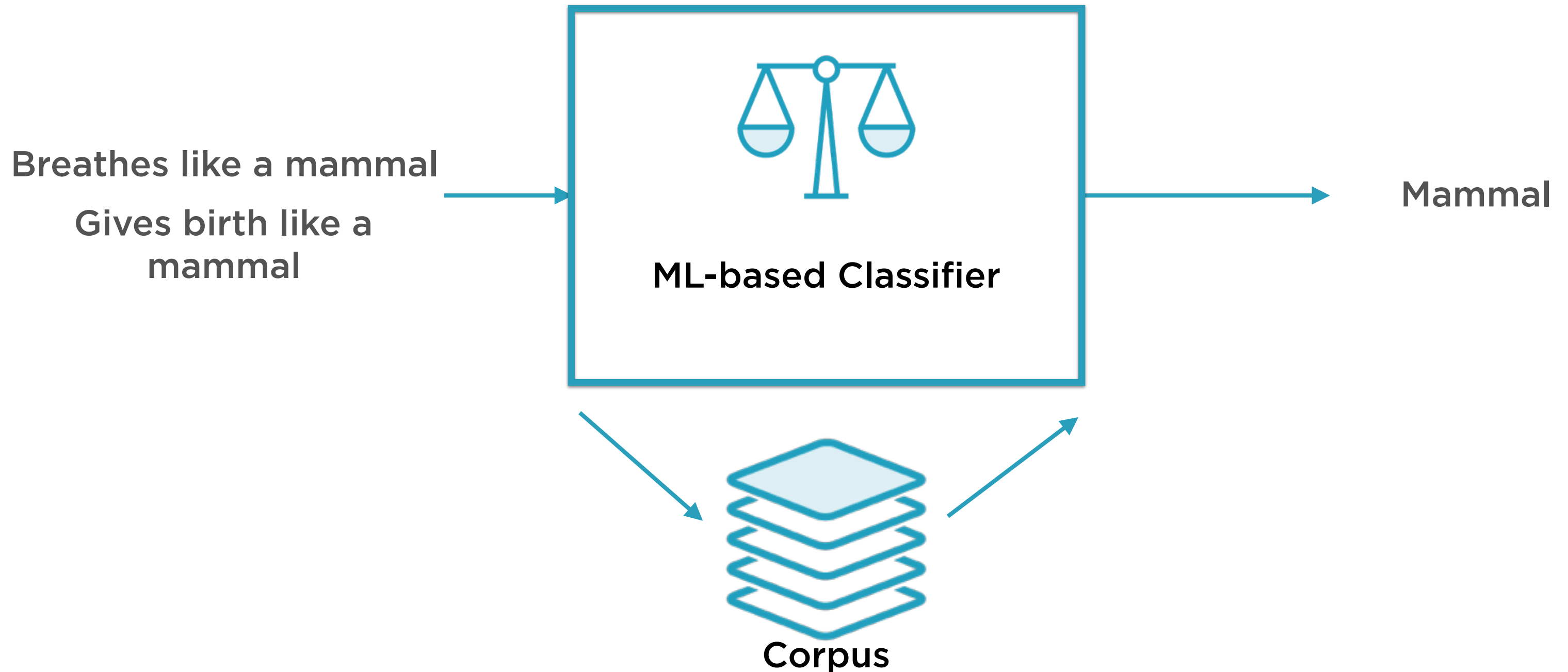
ML-based Binary Classifier



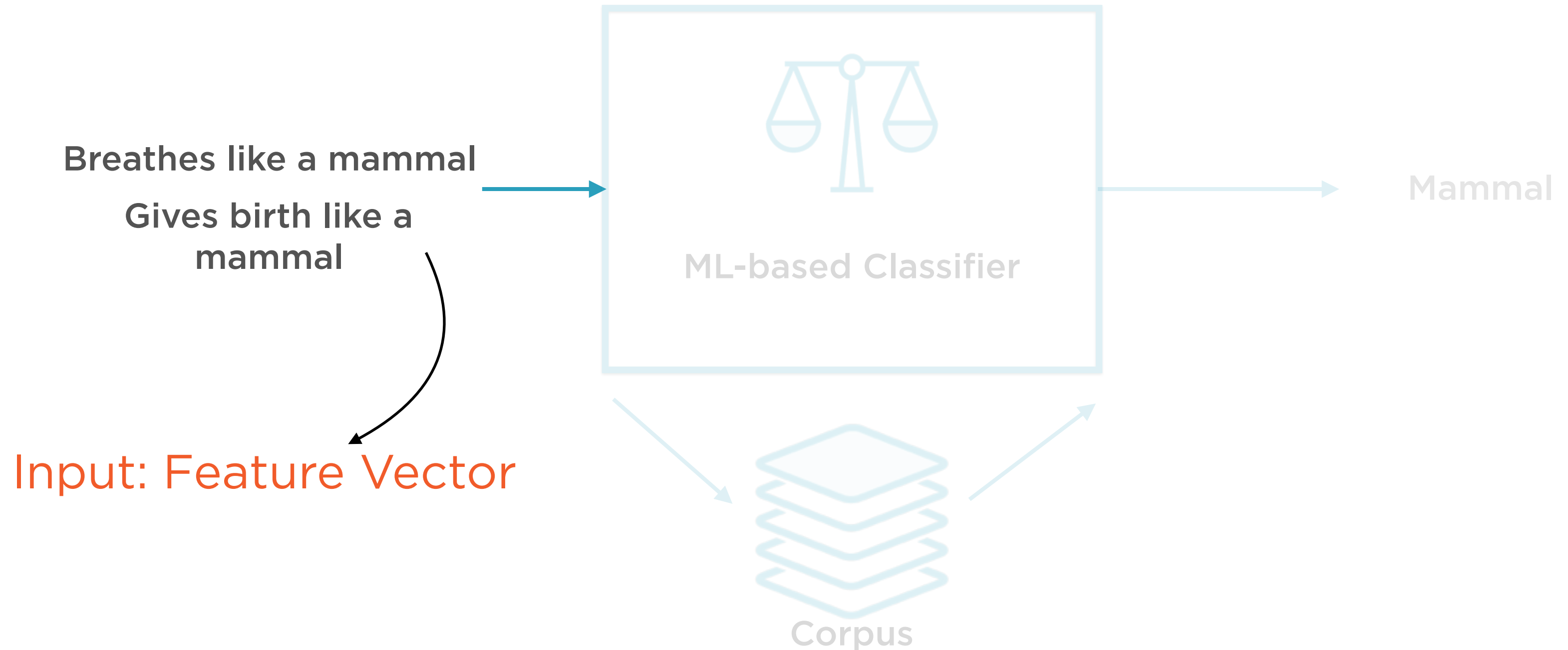
ML-based Binary Classifier



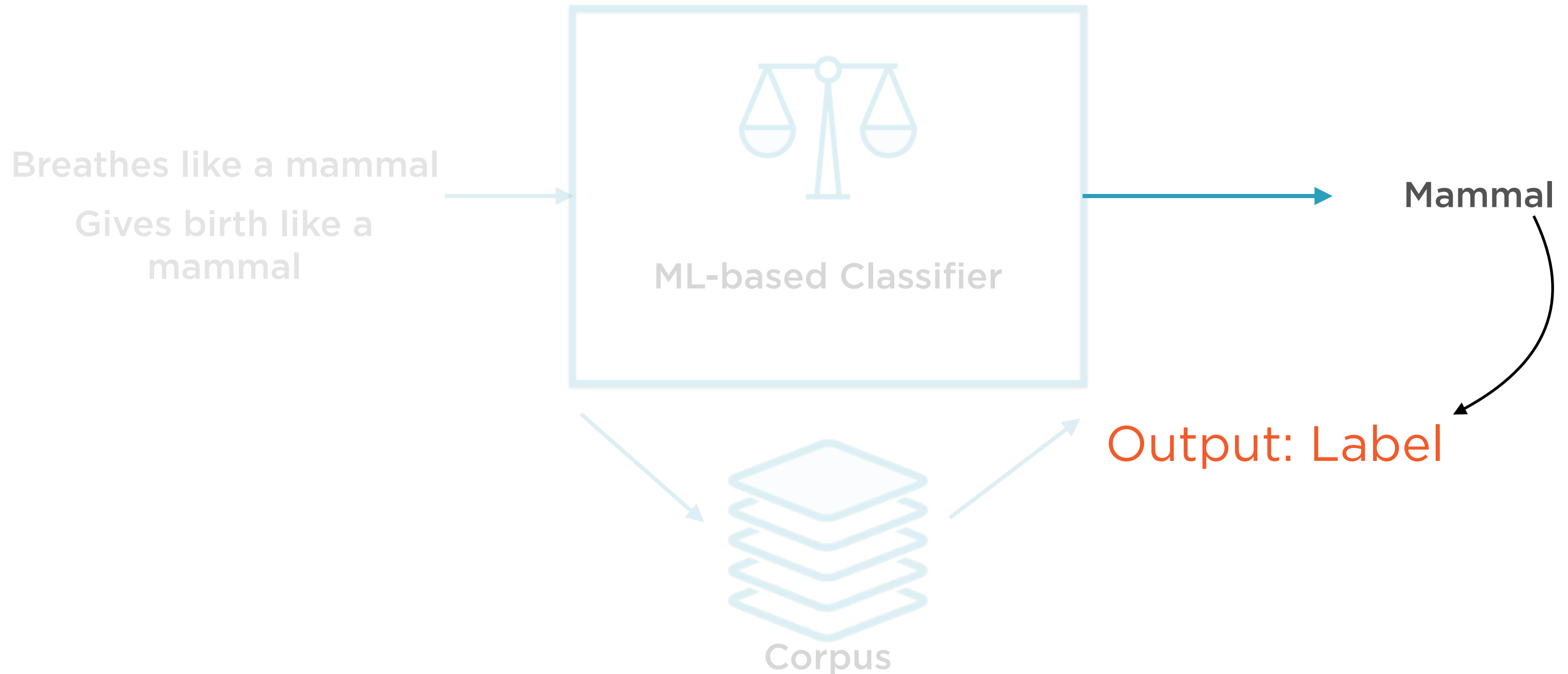
ML-based Binary Classifier



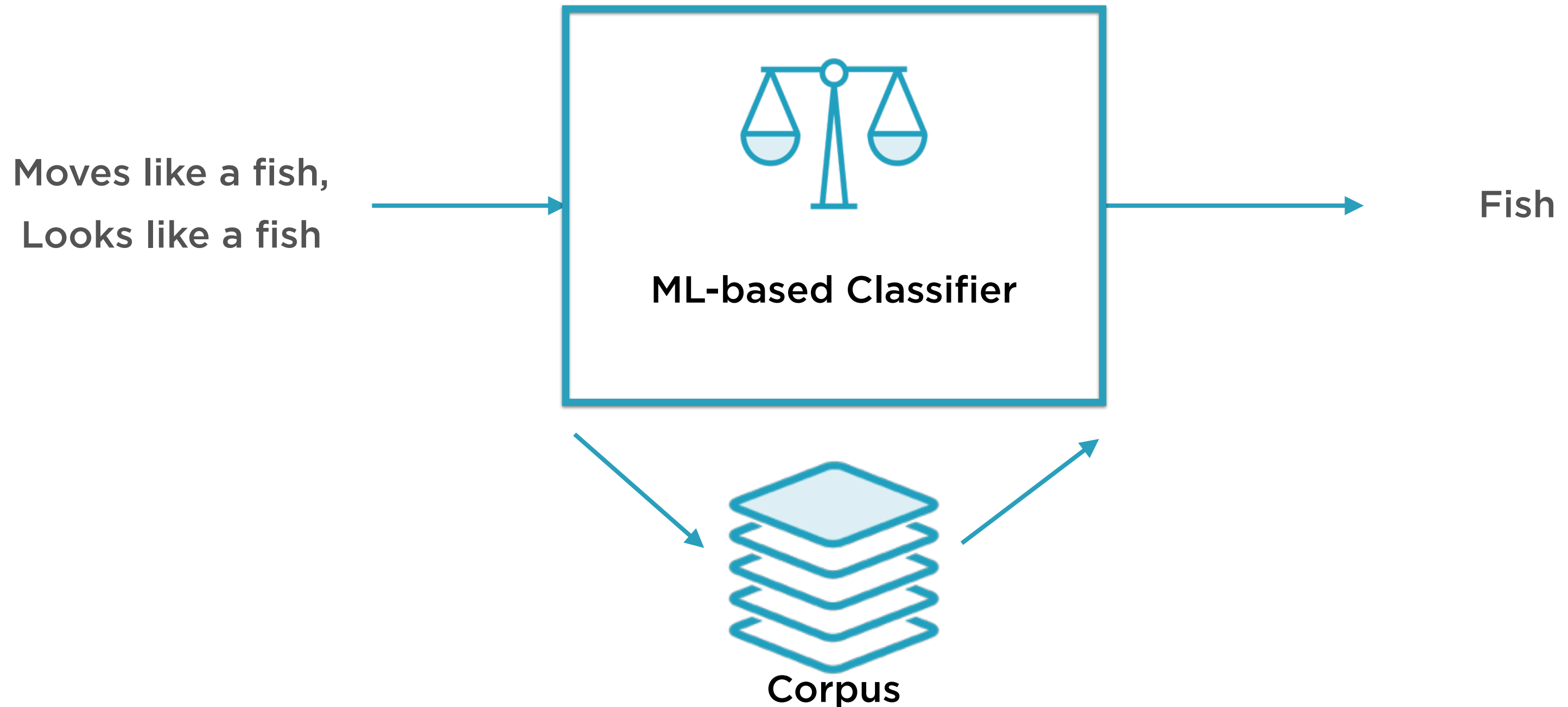
ML-based Binary Classifier



ML-based Binary Classifier



ML-based Binary Classifier



ML-based Binary Classifier

Moves like a fish,
Looks like a fish

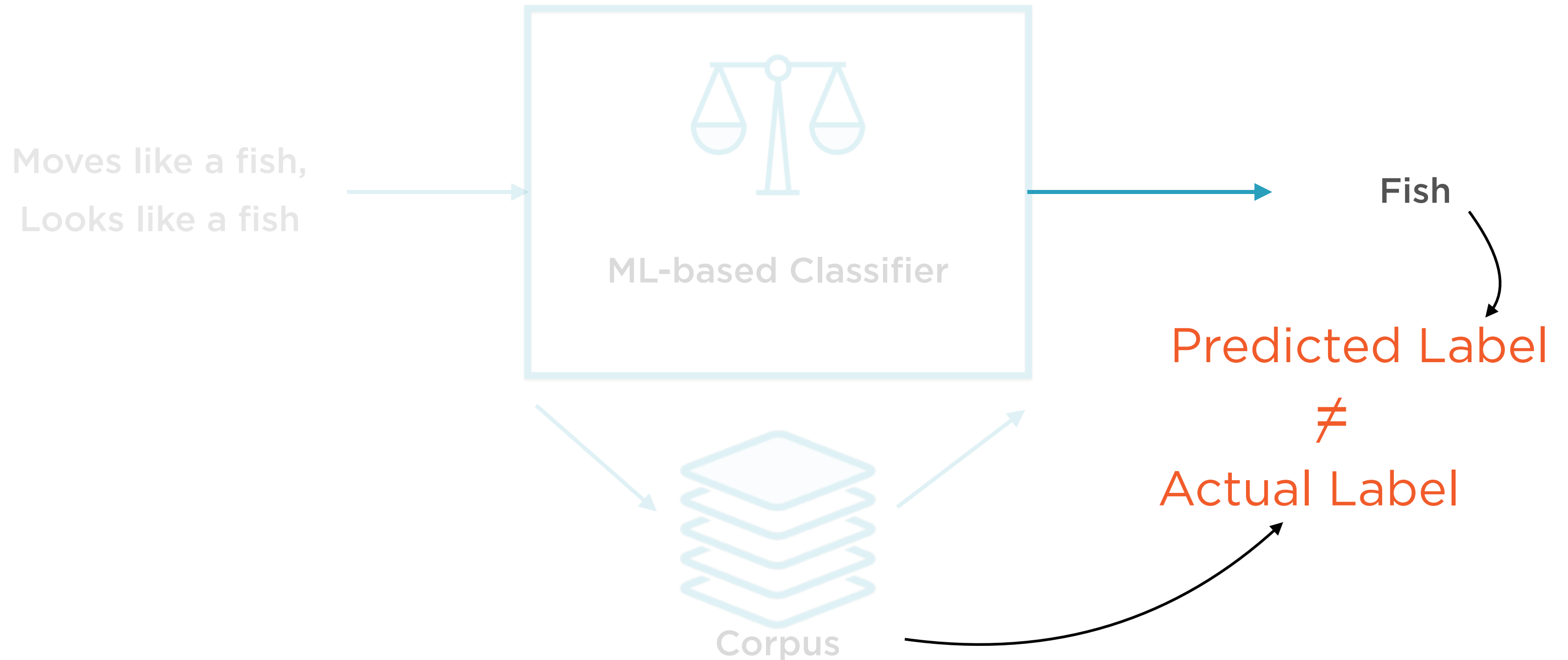
Input: Feature Vector



Fish



ML-based Binary Classifier



“Traditional” ML-based systems still
rely on experts to decide what
features to pay attention to

Traditional ML Models

Regression models: Linear, Lasso,
Ridge, SVR

Classification models: Naive Bayes,
SVMs, Decision trees

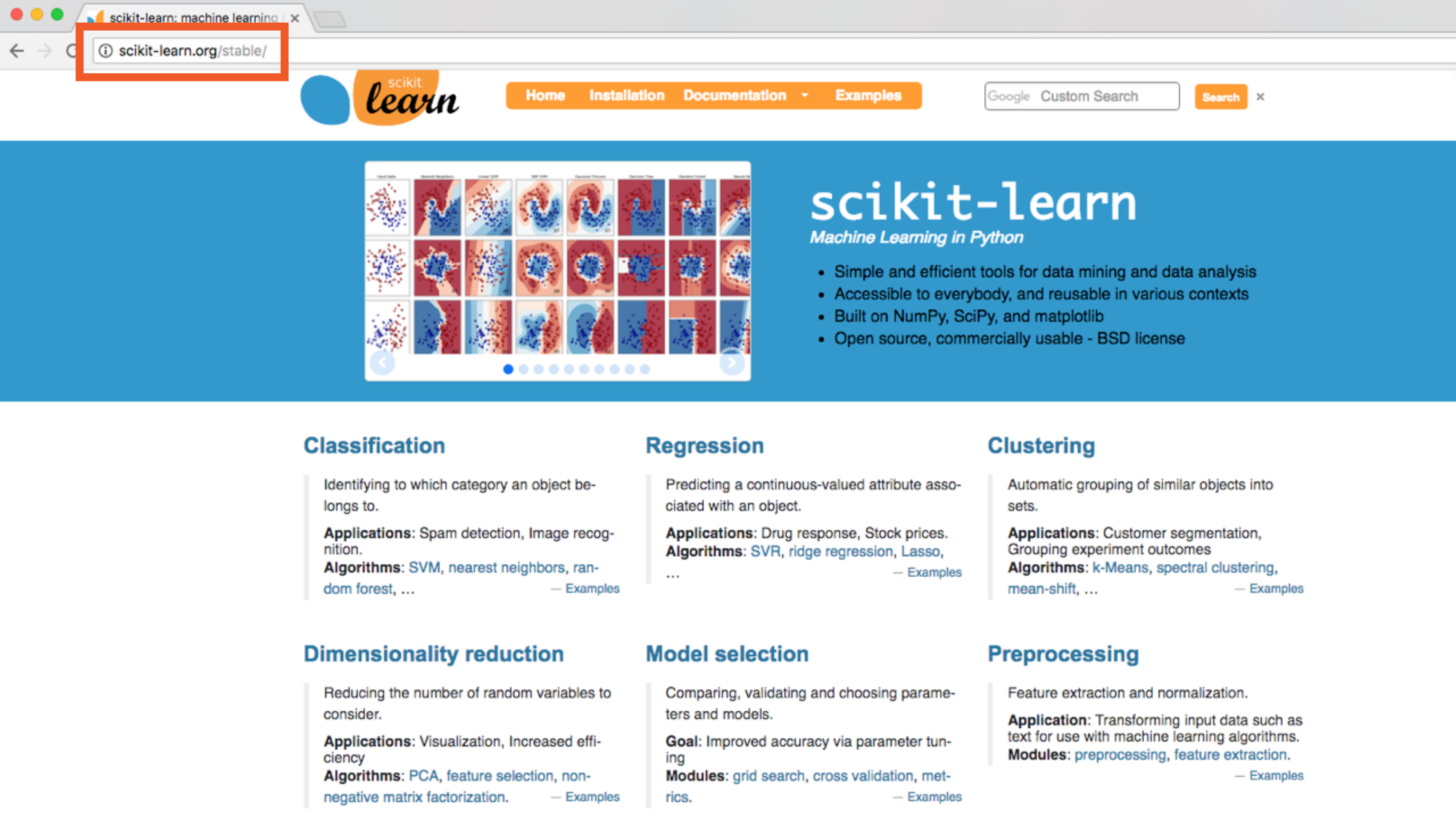
“Representation” ML-based
systems figure out by themselves
what features to pay attention to

Representation
ML Models

Deep learning models such as neural
networks

scikit-learn - a popular, open
source, Python library

Classification, regression, clustering,
dimensionality reduction algorithms



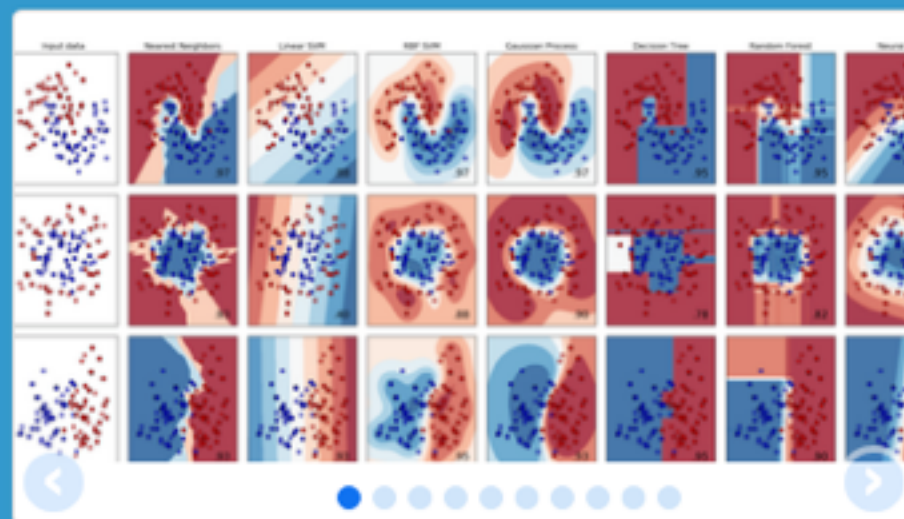
scikit-learn: machine learning | x
scikit-learn.org/stable/



Home Installation Documentation ▾ Examples

Google Custom Search

Search x



scikit-learn

Machine Learning in Python

- Simple and efficient tools for data mining and data analysis
- Accessible to everybody, and reusable in various contexts
- Built on NumPy, SciPy, and matplotlib
- Open source, commercially usable - BSD license

Classification

Identifying to which category an object belongs to.

Applications: Spam detection, Image recognition.

Algorithms: SVM, nearest neighbors, random forest, ... — Examples

Regression

Predicting a continuous-valued attribute associated with an object.

Applications: Drug response, Stock prices.

Algorithms: SVR, ridge regression, Lasso, ... — Examples

Clustering

Automatic grouping of similar objects into sets.

Applications: Customer segmentation, Grouping experiment outcomes

Algorithms: k-Means, spectral clustering, mean-shift, ... — Examples

Dimensionality reduction

Reducing the number of random variables to consider.

Applications: Visualization, Increased efficiency

Algorithms: PCA, feature selection, non-negative matrix factorization. — Examples

Model selection

Comparing, validating and choosing parameters and models.

Goal: Improved accuracy via parameter tuning

Modules: grid search, cross validation, metrics. — Examples

Preprocessing

Feature extraction and normalization.

Application: Transforming input data such as text for use with machine learning algorithms.

Modules: preprocessing, feature extraction. — Examples

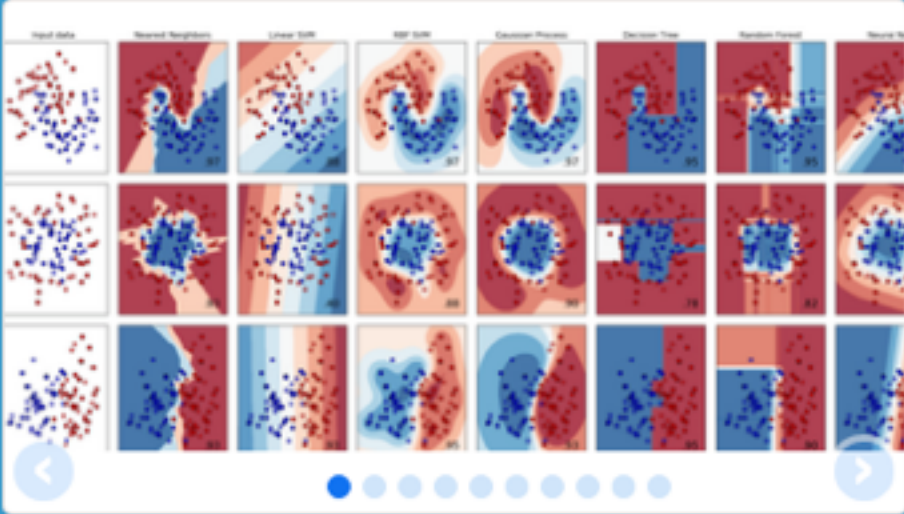
scikit-learn: machine learning

scikit-learn.org/stable/

scikit-learn

HomeInstallationDocumentationExamples

Google Custom SearchSearch



scikit-learn

Machine Learning in Python

- Simple and efficient tools for data mining and data analysis
- Accessible to everybody, and reusable in various contexts
- Built on NumPy, SciPy, and matplotlib
- Open source, commercially usable - BSD license

Classification

Identifying to which category an object belongs to.

Applications: Spam detection, Image recognition.

Algorithms: SVM, nearest neighbors, random forest, ... — Examples

Regression

Predicting a continuous-valued attribute associated with an object.

Applications: Drug response, Stock prices.

Algorithms: SVR, ridge regression, Lasso, ... — Examples

Clustering

Automatic grouping of similar objects into sets.

Applications: Customer segmentation, Grouping experiment outcomes

Algorithms: k-Means, spectral clustering, mean-shift, ... — Examples

Dimensionality reduction

Reducing the number of random variables to consider.

Applications: Visualization, Increased efficiency

Algorithms: PCA, feature selection, non-negative matrix factorization. — Examples

Model selection

Comparing, validating and choosing parameters and models.

Goal: Improved accuracy via parameter tuning

Modules: grid search, cross validation, metrics. — Examples

Preprocessing

Feature extraction and normalization.

Application: Transforming input data such as text for use with machine learning algorithms.

Modules: preprocessing, feature extraction. — Examples

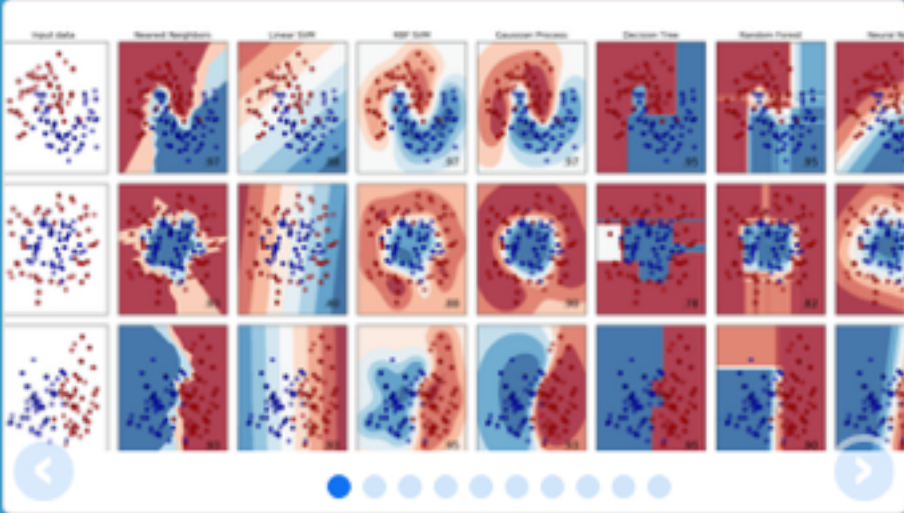
scikit-learn: machine learning

scikit-learn.org/stable/

scikit-learn

HomeInstallationDocumentationExamples

Google Custom SearchSearch



scikit-learn

Machine Learning in Python

- Simple and efficient tools for data mining and data analysis
- Accessible to everybody, and reusable in various contexts
- Built on NumPy, SciPy, and matplotlib
- Open source, commercially usable - BSD license

Classification

Identifying to which category an object belongs to.

Applications: Spam detection, Image recognition.

Algorithms: SVM, nearest neighbors, random forest, ... — Examples

Regression

Predicting a continuous-valued attribute associated with an object.

Applications: Drug response, Stock prices.

Algorithms: SVR, ridge regression, Lasso, ... — Examples

Clustering

Automatic grouping of similar objects into sets.

Applications: Customer segmentation, Grouping experiment outcomes

Algorithms: k-Means, spectral clustering, mean-shift, ... — Examples

Dimensionality reduction

Reducing the number of random variables to consider.

Applications: Visualization, Increased efficiency

Algorithms: PCA, feature selection, non-negative matrix factorization. — Examples

Model selection

Comparing, validating and choosing parameters and models.

Goal: Improved accuracy via parameter tuning

Modules: grid search, cross validation, metrics. — Examples

Preprocessing

Feature extraction and normalization.

Application: Transforming input data such as text for use with machine learning algorithms.

Modules: preprocessing, feature extraction. — Examples

Supervised and Unsupervised Learning

Types of ML Algorithms



Supervised

Labels associated with the training data is used to correct the algorithm



Unsupervised

The model has to be set up right to learn structure in the data

Types of ML Algorithms



Supervised

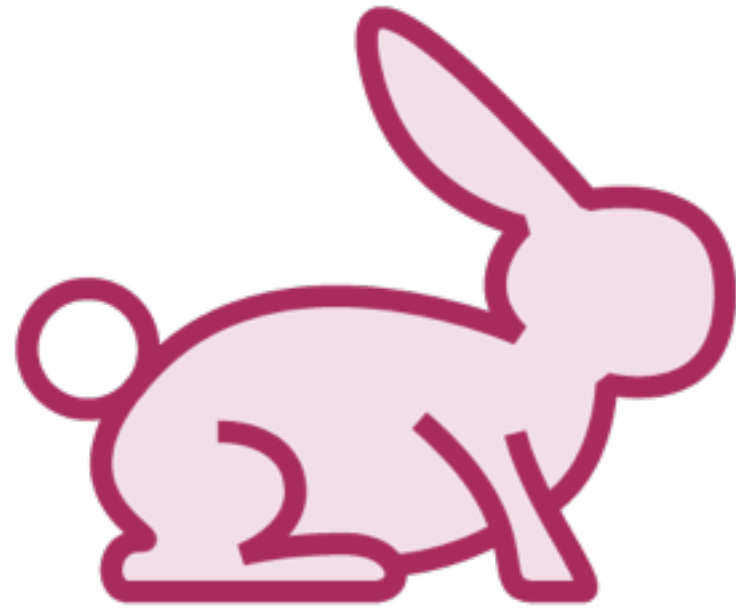
Labels associated with the training data is used to correct the algorithm



Unsupervised

The model has to be set up right to learn structure in the data

Whales: Fish or Mammals?



Mammals

Members of the infraorder
Cetacea



Fish

Look like fish, swim like fish,
move with fish

Whales: Fish or Mammals?



ML-based Classifier

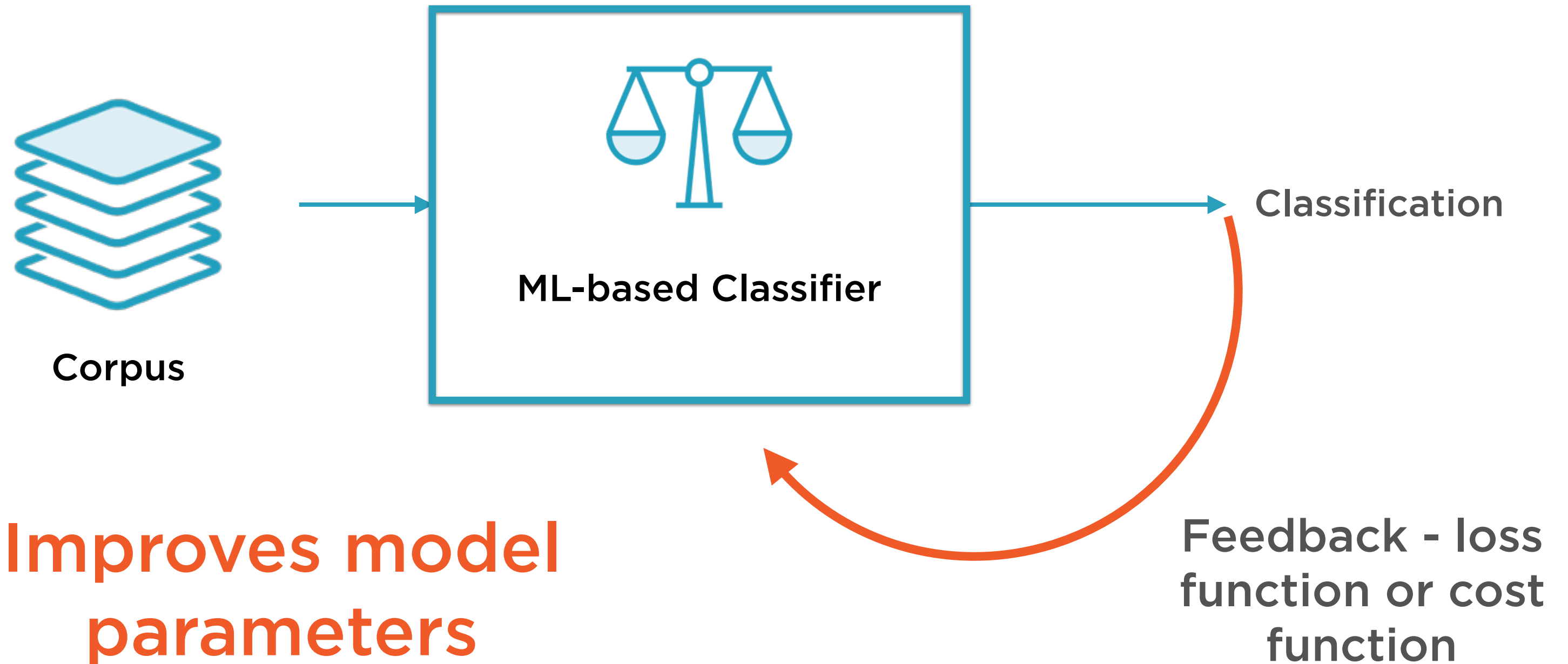
Training

Feed in a large corpus of data
classified correctly

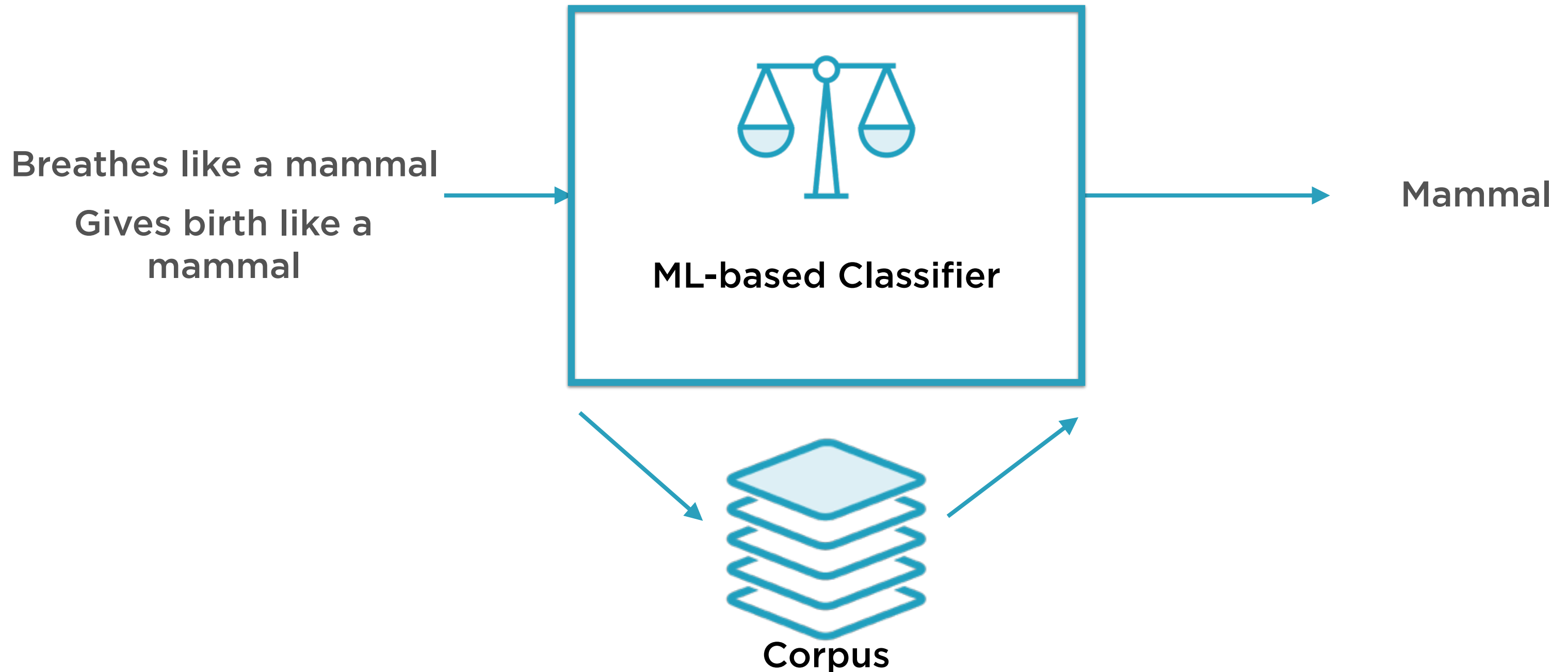
Prediction

Use it to classify new instances
which it has not seen before

Training the ML-based Classifier



ML-based Binary Classifier



x Variables

The attributes that the ML algorithm focuses on are called **features**

Each data point is a list - or **vector** - of such features

Thus, the input into an ML algorithm is a **feature vector**

Feature vectors are usually called the x variables

y Variables

The attributes that the ML algorithm tries to predict are called **labels**

Types of labels

- categorical (classification)
- continuous (regression)

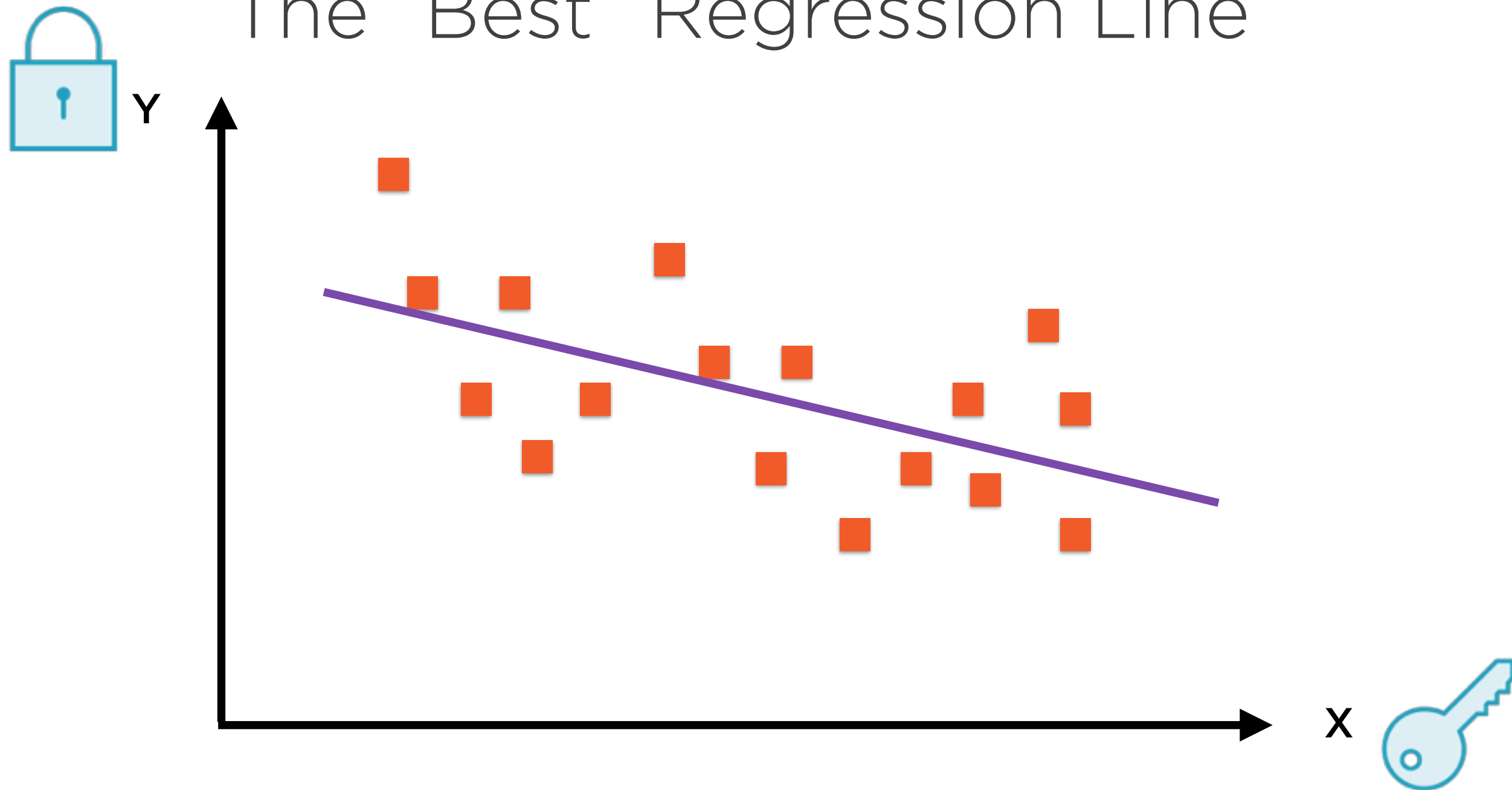
Labels are usually called the y variables

$$y = f(x)$$

Supervised Machine Learning

Most machine learning algorithms seek to “learn” the function f that links the features and the labels

The “Best” Regression Line



Linear Regression involves finding the “best fit” line via a training process

$$y = Wx + b$$

$$f(x) = Wx + b$$

Linear regression specifies, up-front, that the function f is linear

```
def doSomethingReallyComplicated(x1, x2...):  
    ...  
    ...  
    ...  
    return complicatedResult
```

$f(x) = \text{doSomethingReallyComplicated}(x)$

ML algorithms such as neural network can “learn” (reverse-engineer) pretty much anything given the right training data

Types of ML Algorithms



Supervised

Labels associated with the training data is used to correct the algorithm



Unsupervised

The model has to be set up right to learn structure in the data

Unsupervised Learning does not have:

- y variables
- a labeled corpus



Supervised Learning

Input variable x and output variable y

Learn the mapping function $y = f(x)$

Approximate the mapping function so
for new values of x we can predict y

Use existing dataset to **correct** our
mapping function approximation



Unsupervised Learning

Only have input data **x** - no output data

Model the underlying structure to learn more about data

Algorithms **self discover** the patterns and structure in the data

Unsupervised ML Algorithms

Clustering

Identify patterns in data items e.g.
K-means clustering

Dimensionality reduction

Identify significant factors that drive
data e.g. PCA

Continuous and Categorical Data

Continuous and Categorical Variables

Continuous

Can take an infinite set of values
(height, weight, income...)

Categorical

Can take a finite set of values (Male/
Female, Day of week...)

Categorical variables that can take just two
values are called **binary variables**

Continuous and Categorical Variables

Continuous

Can take an infinite set of values
(height, weight, income...)

Categorical

Can take a finite set of values (Male/
Female, Day of week...)

Standardizing Data: Mean and Variance

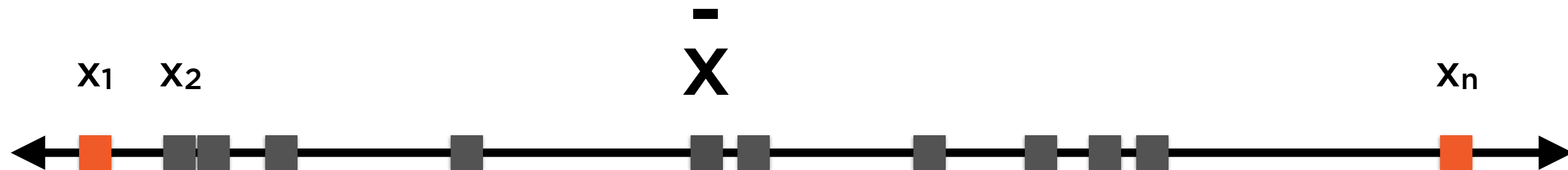
Mean as Headline



The mean, or average, is the one number that best represents all of these data points

$$\bar{x} = \frac{x_1 + x_2 + \dots + x_n}{n}$$

Variation Is Important Too

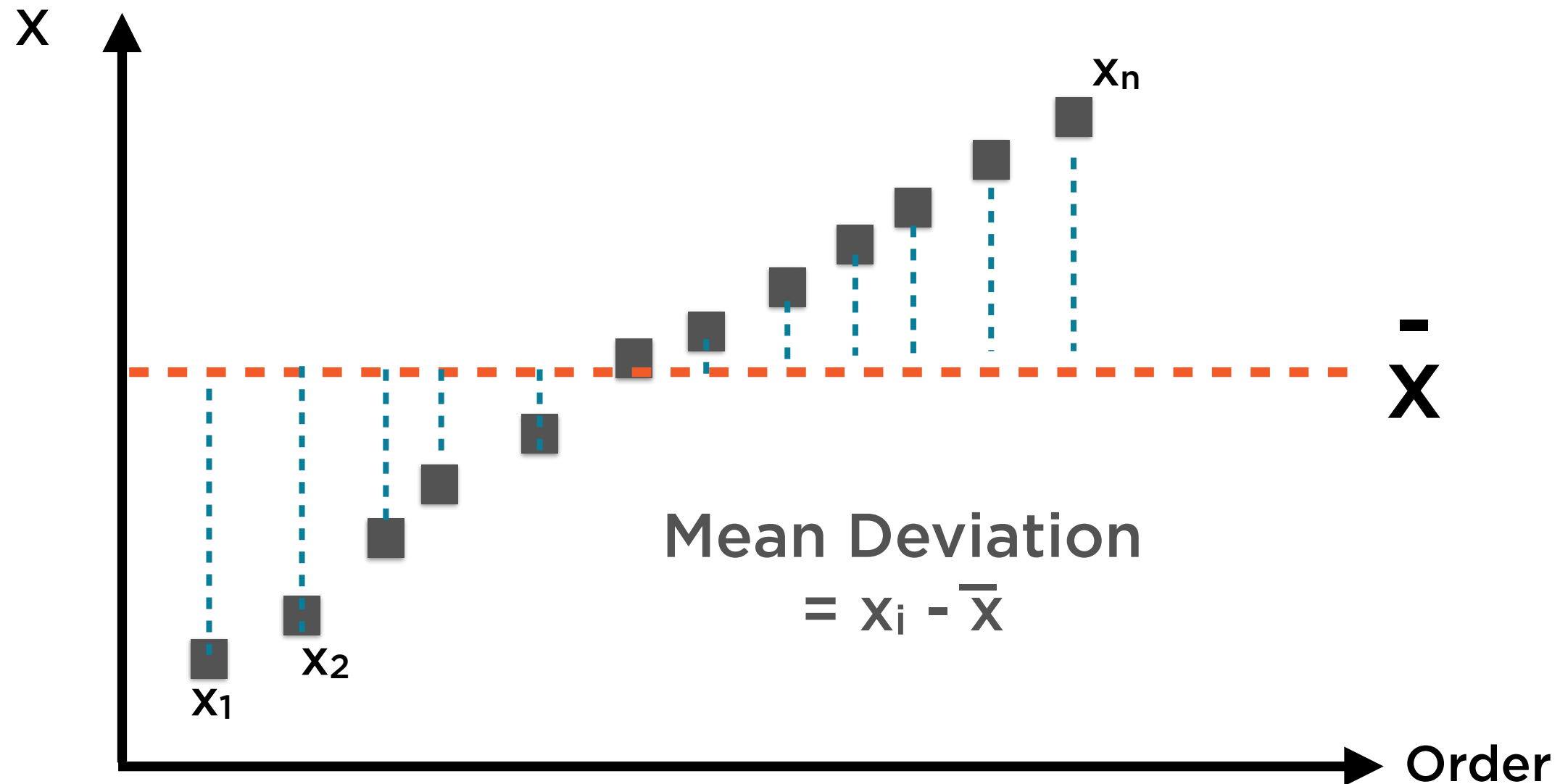


“Do the numbers jump around?”

$$\text{Range} = X_{\max} - X_{\min}$$

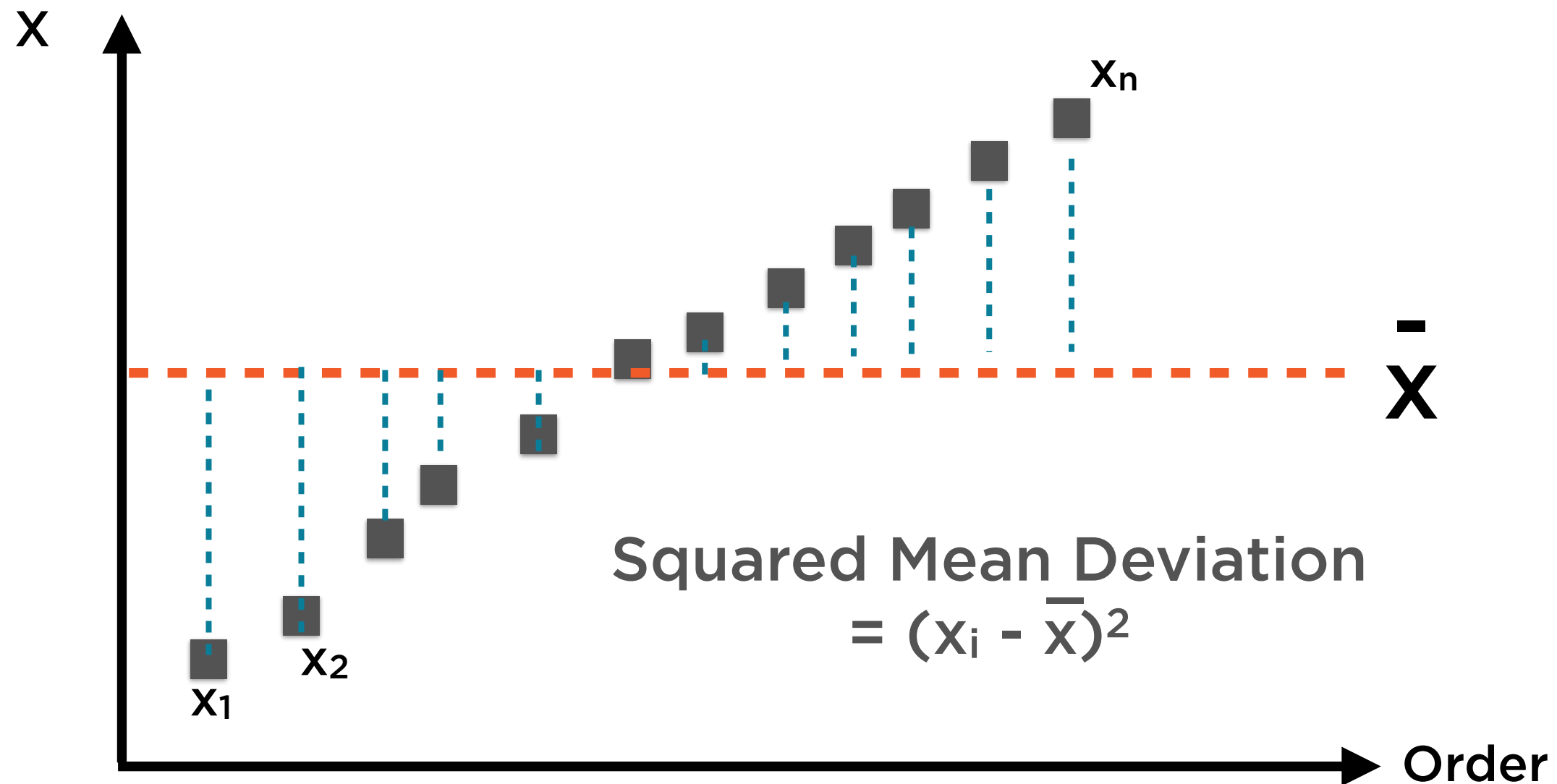
The range ignores the mean, and is swayed by outliers - that's where variance comes in

Variance as Asterisk



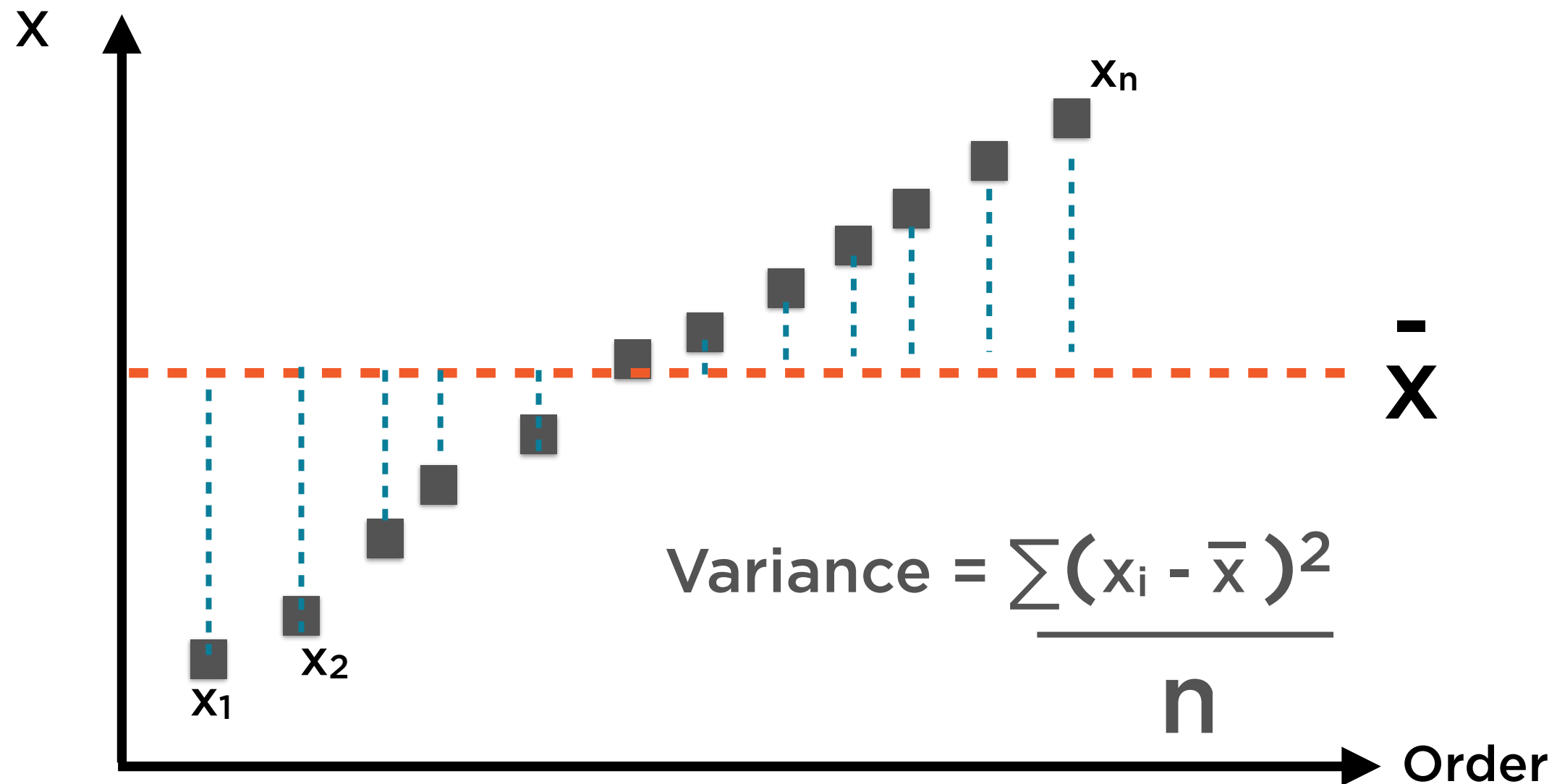
Variance is the second-most important number to summarize this set of data points

Variance as Asterisk



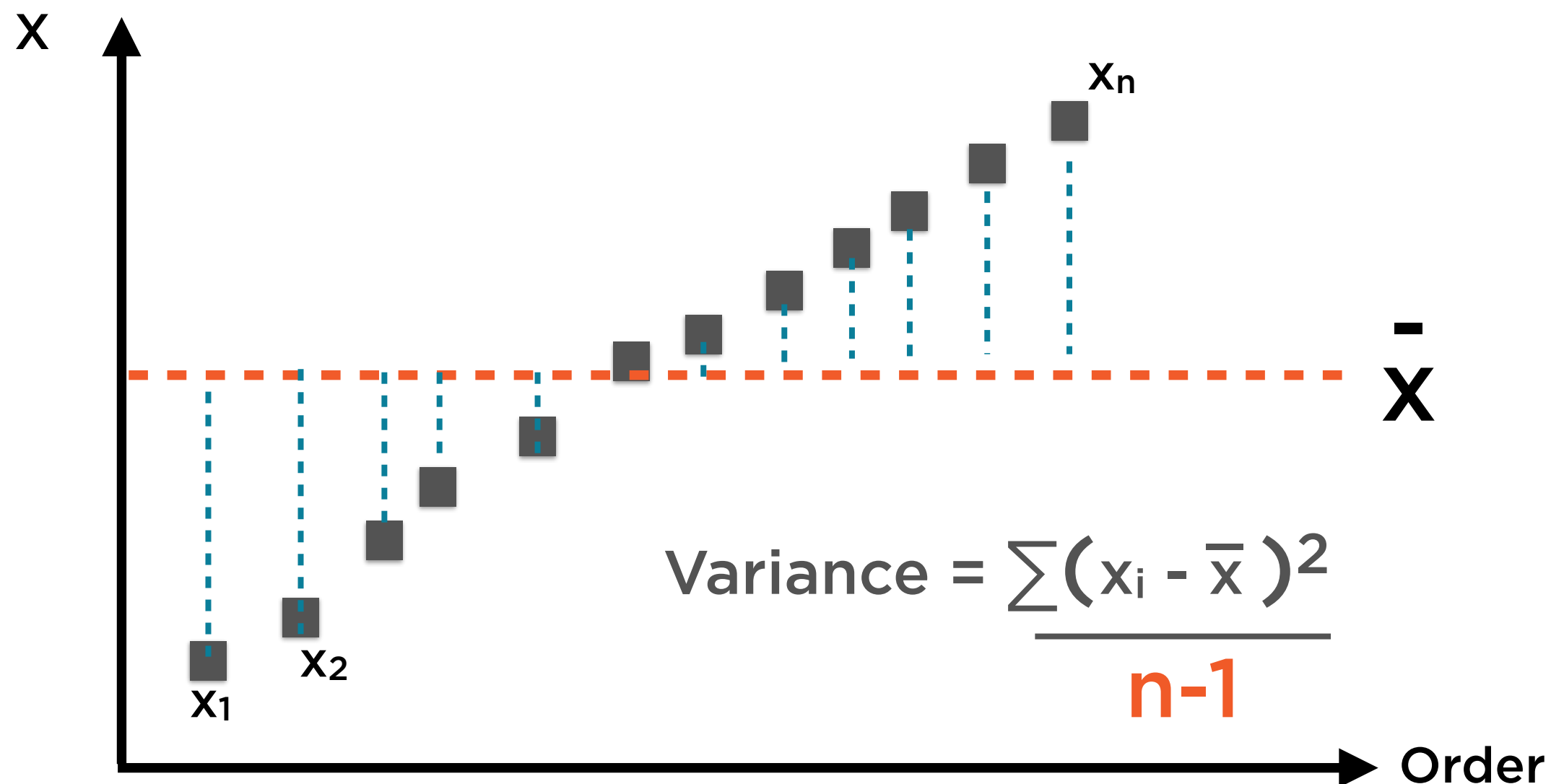
Variance is the second-most important number to summarize this set of data points

Variance as Asterisk



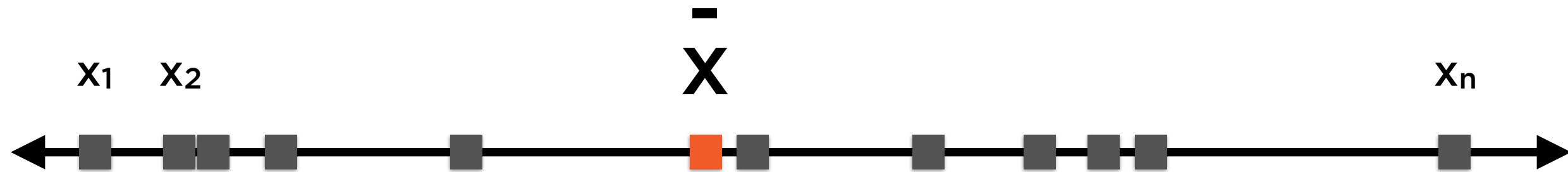
Variance is the second-most important number to summarize this set of data points

Variance as Asterisk



We can improve our estimate of the variance by tweaking the denominator - this is called **Bessel's Correction**

Mean and Variance



Mean and variance succinctly summarize a set of numbers

$$\bar{x} = \frac{x_1 + x_2 + \dots + x_n}{n}$$

$$\text{Variance} = \frac{\sum (x_i - \bar{x})^2}{n-1}$$

Variance and Standard Deviation

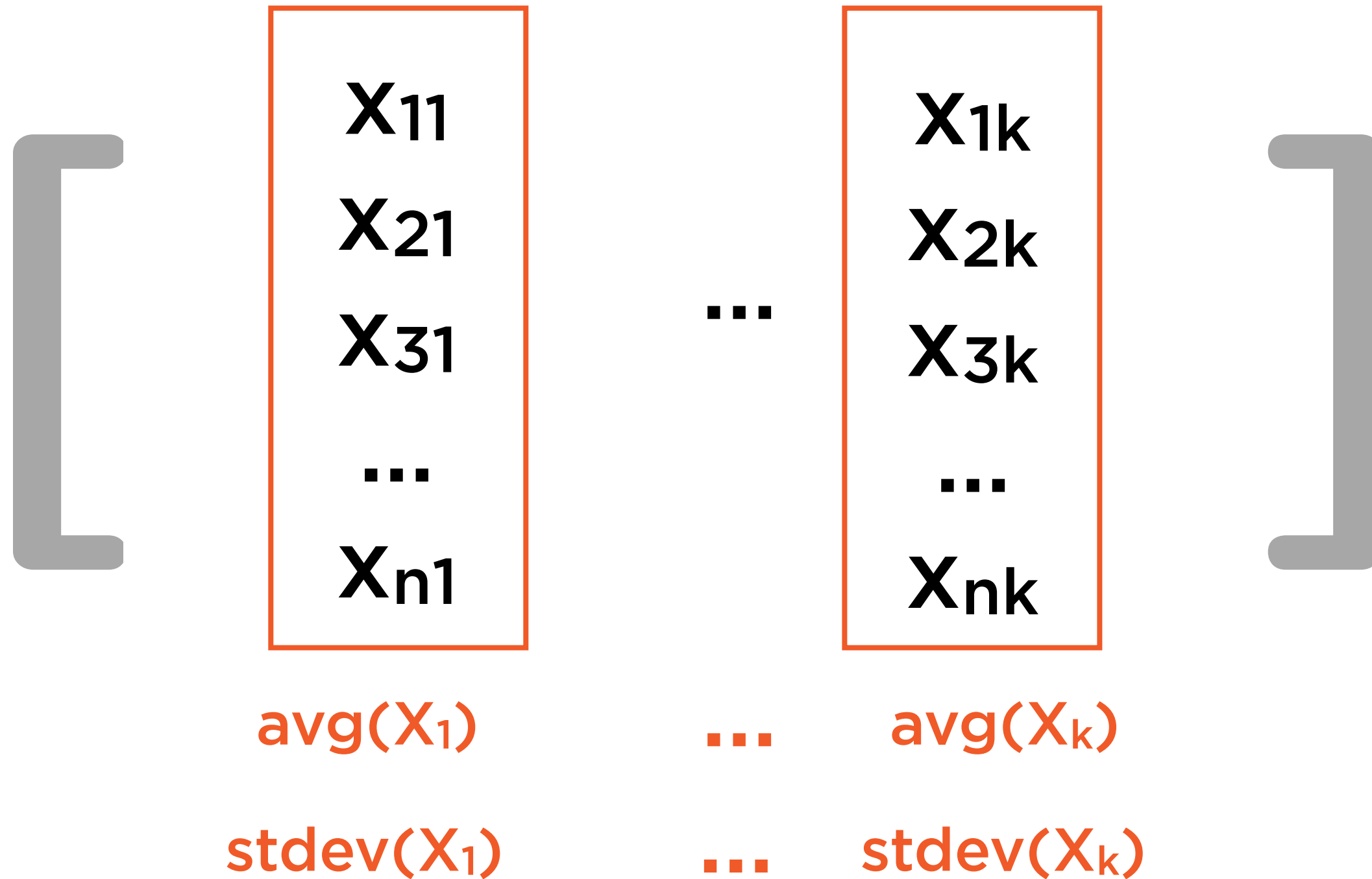


Standard deviation is the square root of variance

$$\text{Variance} = \frac{\sum (x_i - \bar{x})^2}{n-1}$$

$$\text{Std Dev} = \sqrt{\frac{\sum (x_i - \bar{x})^2}{n-1}}$$

Standardizing Data



Standardizing Data

$$\begin{bmatrix} \frac{x_{11} - \text{avg}(X_1)}{\text{stdev}(X_1)} & \frac{x_{1k} - \text{avg}(X_k)}{\text{stdev}(X_k)} & \dots \\ \dots & \dots & \dots \\ \frac{x_{n1} - \text{avg}(X_1)}{\text{stdev}(X_1)} & \frac{x_{nk} - \text{avg}(X_k)}{\text{stdev}(X_k)} & \dots \end{bmatrix}$$

Each column of the standardized data has mean 0 and variance 1



Standardized Data

Many techniques work best on standardized data

Standardization prevents some (high-variance) data series from dominating

Examples:

- Principal Components Analysis
- Lasso/Ridge Regression

Continuous and Categorical Variables

Continuous

Can take an infinite set of values
(height, weight, income...)

Categorical

Can take a finite set of values (Male/
Female, Day of week...)



Categorical Data

**Continuous data can be ordered,
categorical data can not**

ML algorithms only operate on numbers

**Categorical data need to be encoded as
numbers**

**Numerical encodings of categorical data
should never be ordered**



Categorical Data

Continuous data can be ordered,
categorical data can not

ML algorithms only operate on numbers

**Categorical data need to be encoded as
numbers**

Numerical encodings of categorical data
should never be ordered

One-hot Encoding

Sunday

Monday

Tuesday

Wednesday

Thursday

Friday

Saturday

One-hot Encoding

	Sunday	Monday	Tuesday	Wednesday	Thursday	Friday	Saturday
Monday	0	1	0	0	0	0	0
Thursday	0	0	0	0	1	0	0
Saturday	0	0	0	0	0	0	1