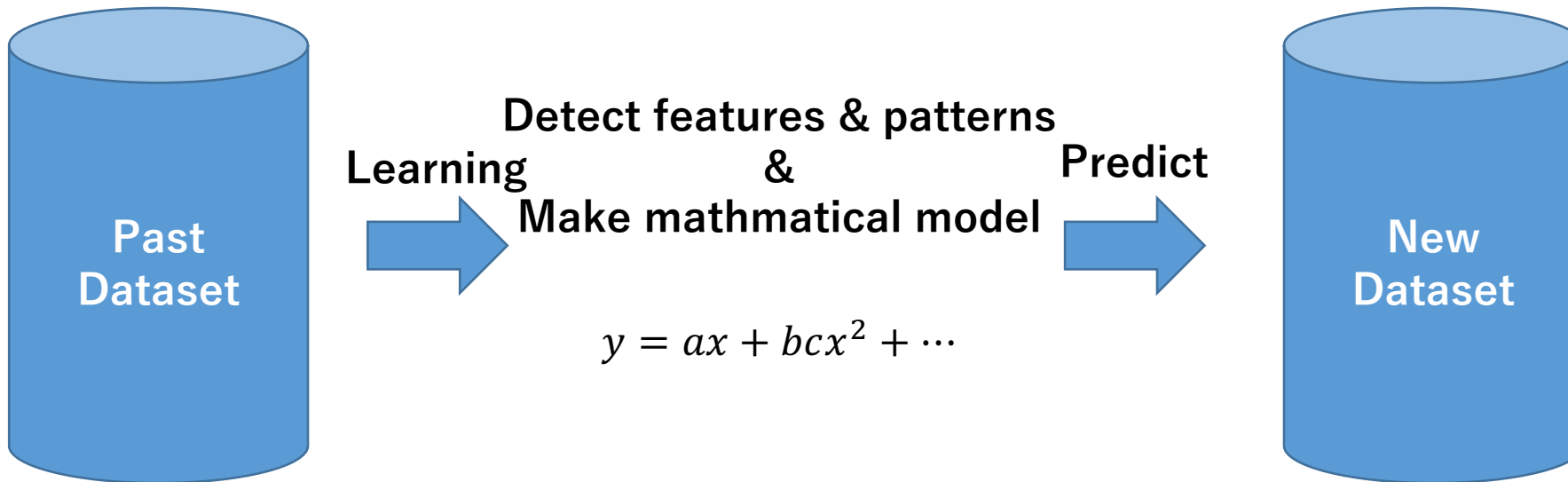# Difference between machine learning and statistics

# &

# Application to genome analysis

# What is Machine Learning？

Machine learning (ML) is the study of algorithms and mathematical models
that computer systems use to progressively improve their performance on a specific task…(Wikipedia)

**Past Dataset**

**Learning**

**Detect features & patterns
&
Make mathmatical model**

$$y = ax + bcx^2 + \cdots$$

**Predict**

**New Dataset**

# Example



*Digitaria ciliaris*

*Digitaria violascens*

# Example



= *Digitaria ciliaris*

= *Digitaria ciliaris*

= *Digitaria violascens*

= *Digitaria violascens*

⋮          ⋮

**Learning**
**From past data**

**Detect patterns**

(ex. *Digitaria violascens*
= Hairless leaf sheath)

**Prediction**

**New Data**

= ？？？

**97%:*Digitaria violascens***
**3%: *Digitaria ciliaris***

# Example

**Past Data**

| Sequence | Crop yield |
|----------|------------|
| ATTGAC⋯ | 2,0kg |
| GAGGTA⋯ | 3,6kg |
| TGCCGC⋯ | 1,1kg |
| ATCGAA⋯ | 2.1kg |
| … | … |

**Learning**

**Detect features & patterns
&
Make mathmatical model**
$$CropYield = ax + by + cz \dots$$

**Prediction**

**New Dataset**

| Sequence | Crop yield |
|----------|------------|
| GAAAAC⋯ | ??? |
| TTAGGG⋯ | ??? |
| … | … |

| Sequence | Predicted Crop yield |
|----------|------------|
| GAAAAC⋯ | 1.82kg |
| TTAGGG⋯ | 1.11kg |
| … | … |

# Difference between machine learning and statistics

…method & algorithm are common

- ## Statistics

Description statistics ⋯ Make data easier for people
                                to understand using statistics or visualize
        ex) mean, variance

Estimated statistics ⋯Consideration on data

        ex) Estimate mean value of population
            Estimate the position of related genes at SNP-index

- ## Machine learning
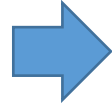
Make highly accurate prediction models from past data

# Difference between machine learning and statistics

··· Different purposes / strategies

- ## Statistics ··· Selection of models and methods and validity are important

    ex) $CropYield = ax + by + cz \ldots$

    ···**Good Accuracy**

    Which variable are important ?
    Why this model get good accuracy ?

    ex) $a$ is big $\rightarrow x$ has large effect

- ## Machine learning ··· How to achieve high prediction accuracy is important

$\rightarrow$ Focus on prediction accuracy even in various methods, black box model

ex) $CropYield = \sqrt[3]{a^8 x^{27} z^{55}} + \lim_{n \to \infty} \left( \frac{\int b \sin \alpha n}{n} \right)^n + \cdots$

···**Very high accuracy**

Sometimes it's difficult to understand model

ex) $a$ is big $\rightarrow$ ？？？

# Application to genome analysis

ex1) **To know gene function** → It is important how easy it is to understand the model

$$CropYield = \beta_1 gene_1 + \beta_2 gene_2 + \beta_3 gene_3 \dots$$

Gene Effects $\beta$



It shows that
Genes in this position affect Crop Yield

ex2) **Apply Genomic Breeding** → The prediction accuracy of new dataset is also important

CropYield +



Model

New Dataset



Prediction Accuracy

## Summary

- Statistics and machine learning are trying to solve problems by using data, but the goal and strategy is different.

- Especially genome analysis

    Gene Data → To find gene function, evolutionary mechanism…etc

    ・What does data mean ?

    ・What can we detect from models ?