

---

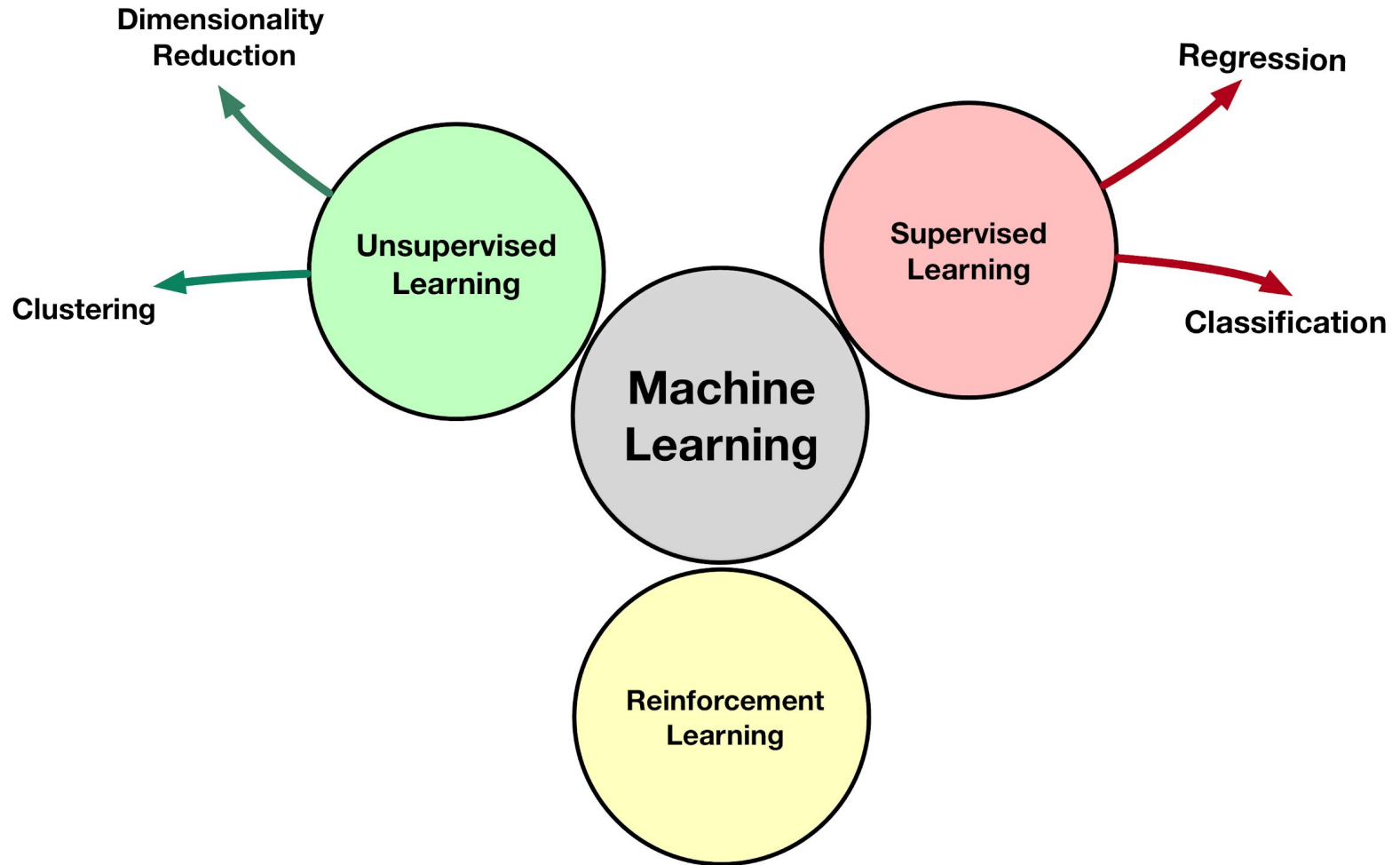
---

# Supervised learning

Ali Madani  
Farnoosh Khodakarami

---

---



# Supervised vs Unsupervised Learning

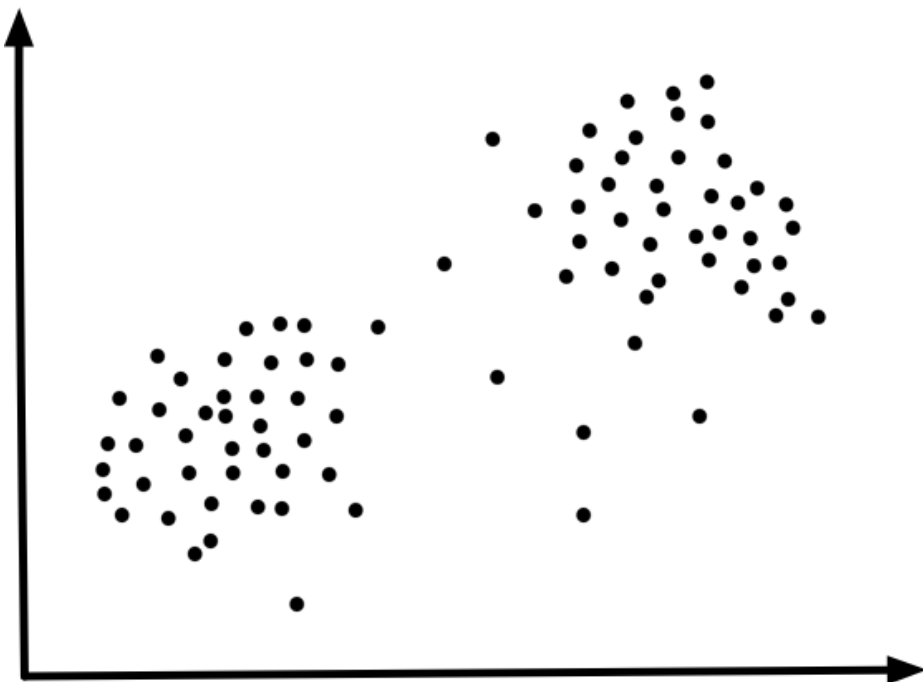
## Unsupervised Learning

- **No Knowledge** of output
- data is **unlabeled**
- Self guided learning
- **Goal:** determine data patterns/grouping

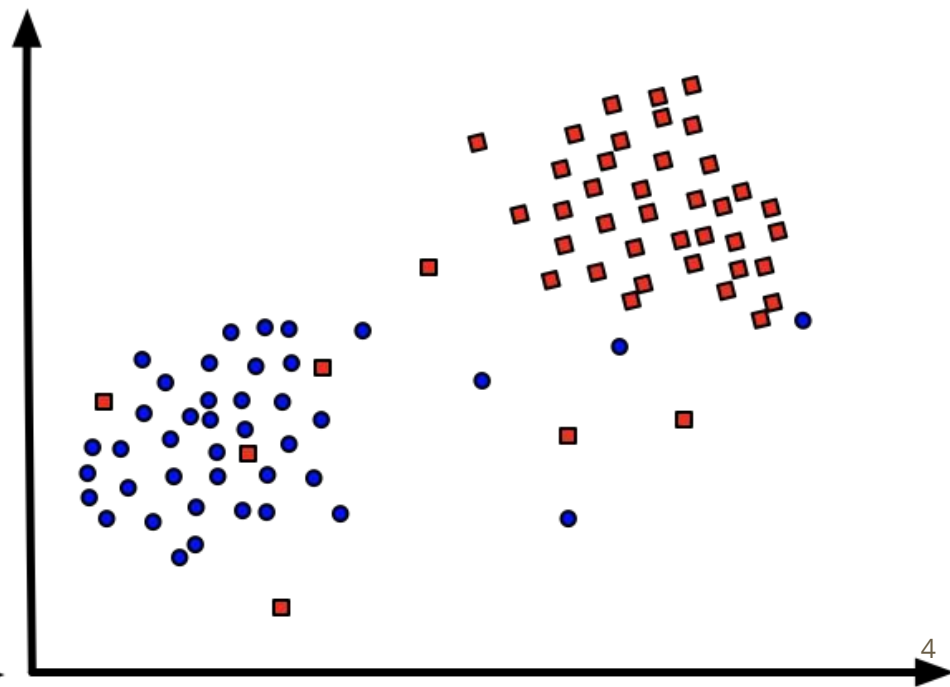
## Supervised Learning

- **Knowledge** of output
- data is **labeled** with class or value
- **Goal:** predict value label or class label

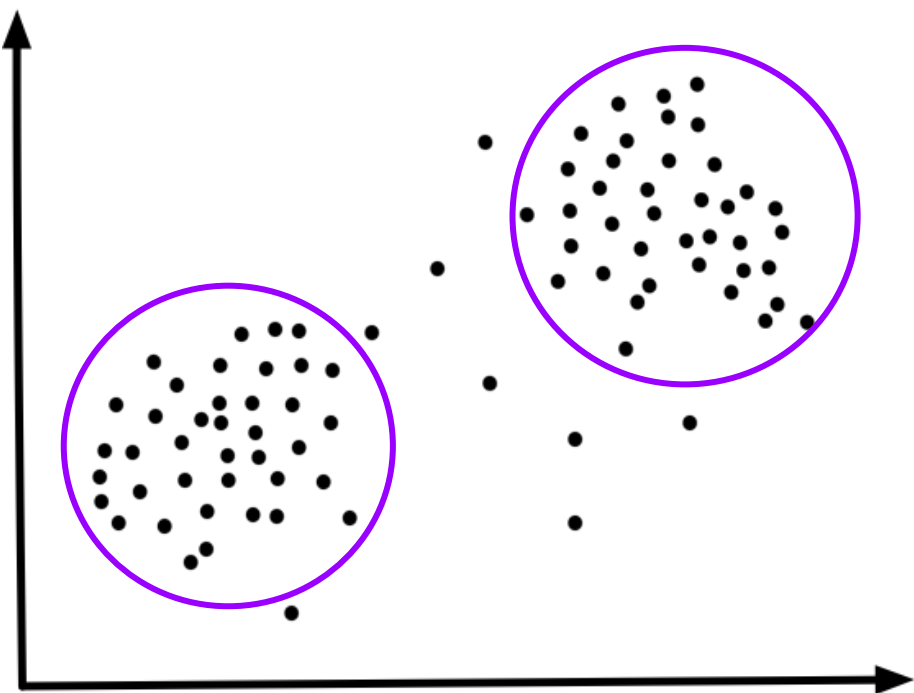
## Unsupervised Learning



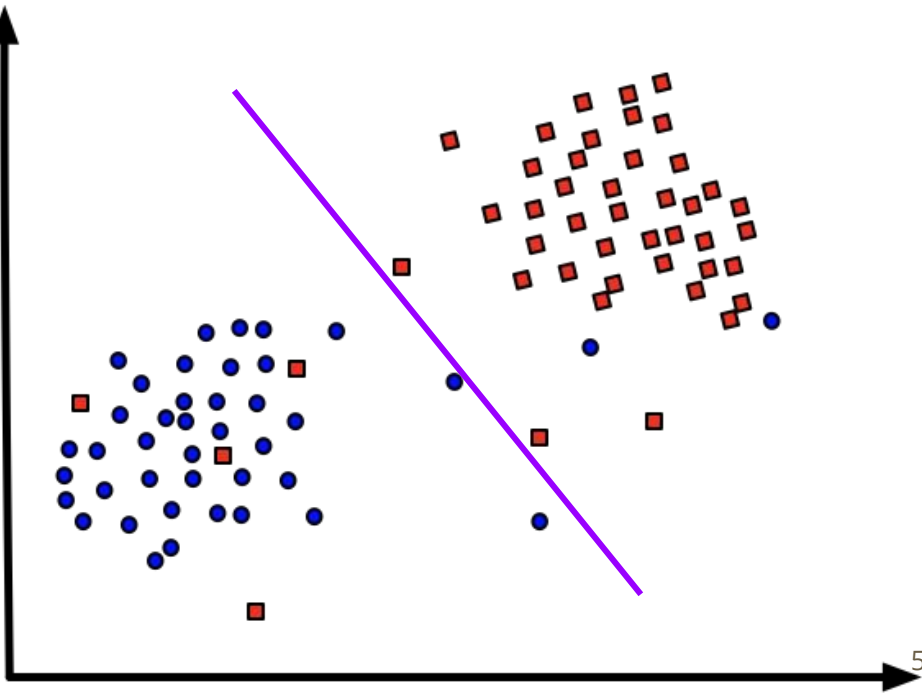
## Supervised Learning

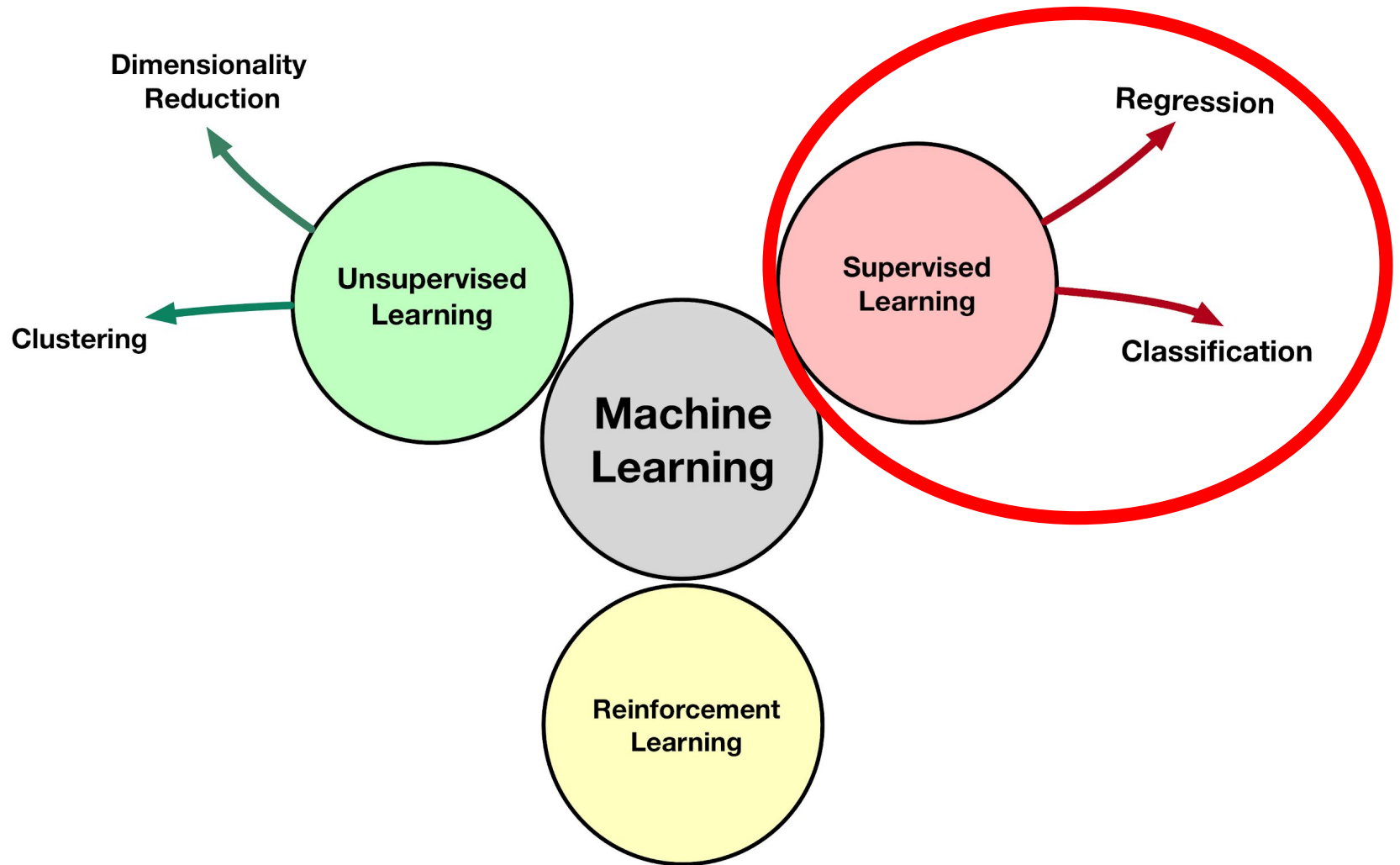


## Unsupervised Learning



## Supervised Learning

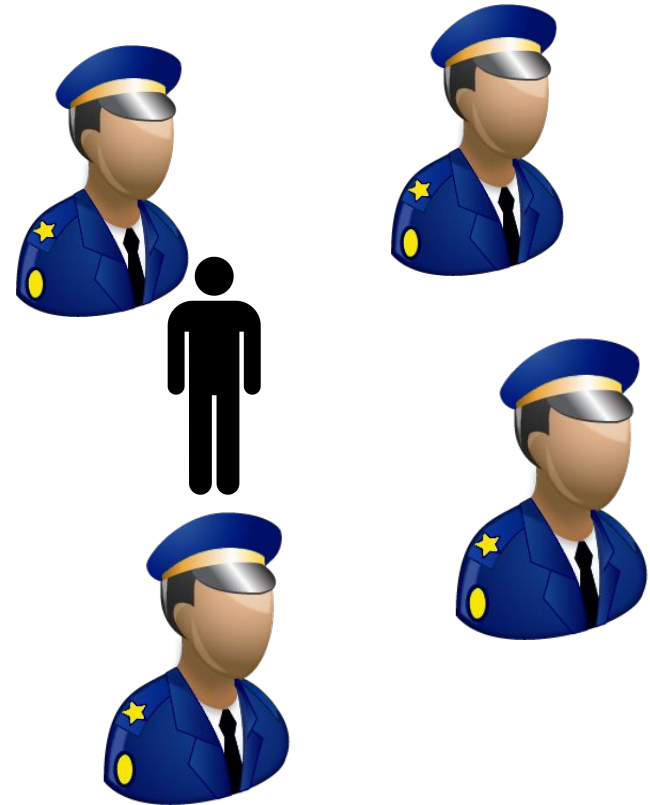
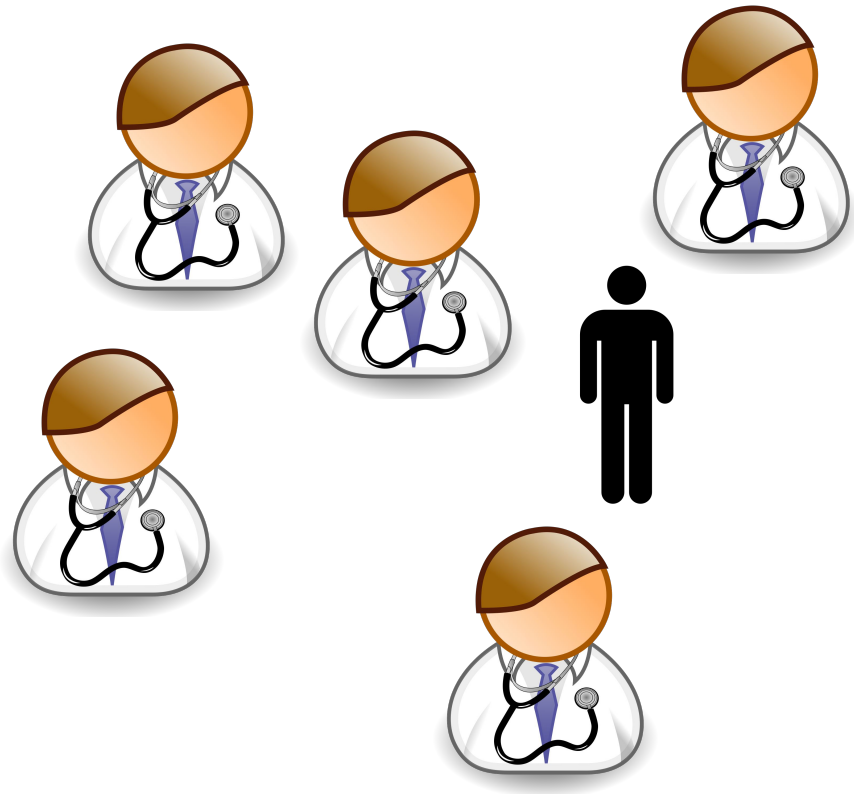






# Machine Learning Algorithms

# K Nearest Neighbor





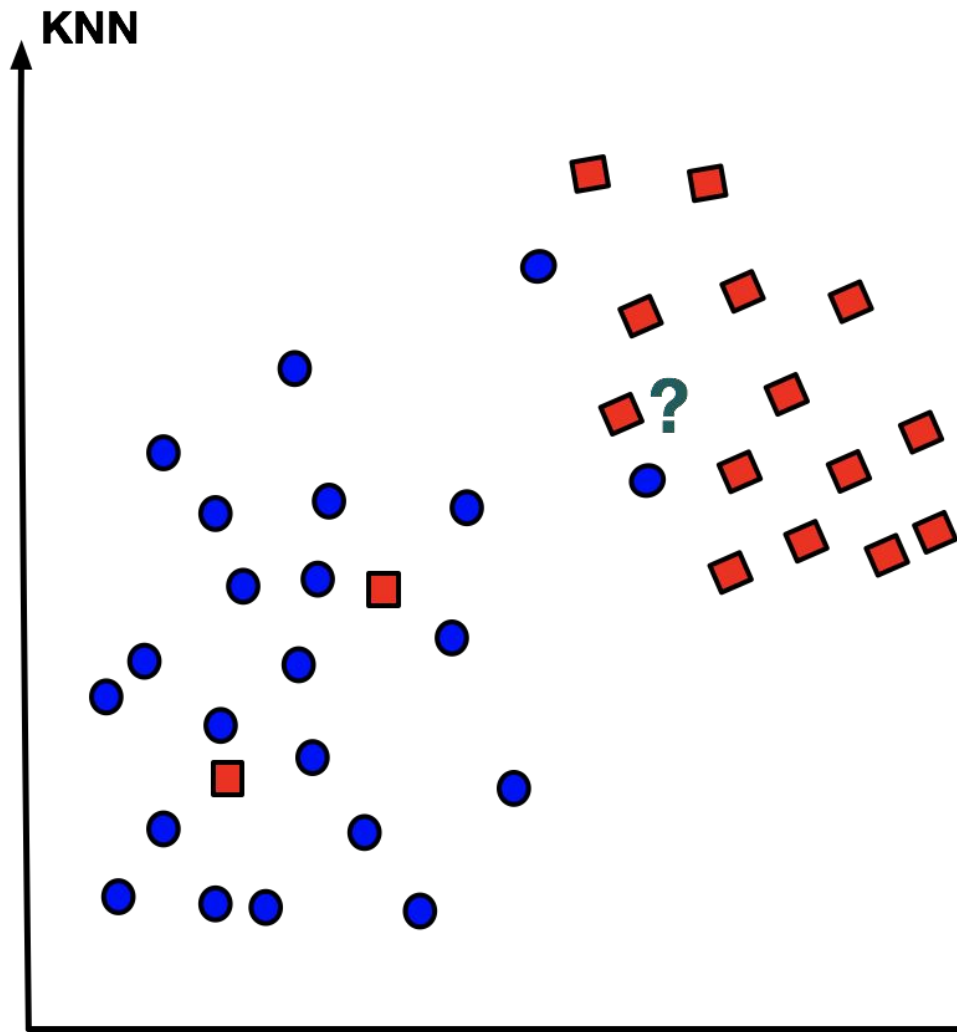
# K Nearest Neighbor

- Can be used both for classification and regression.
- Uses **feature similarity** to predict values of any new data points.
- The output based on the majority vote (for classification)
- or mean (or median, for regression)

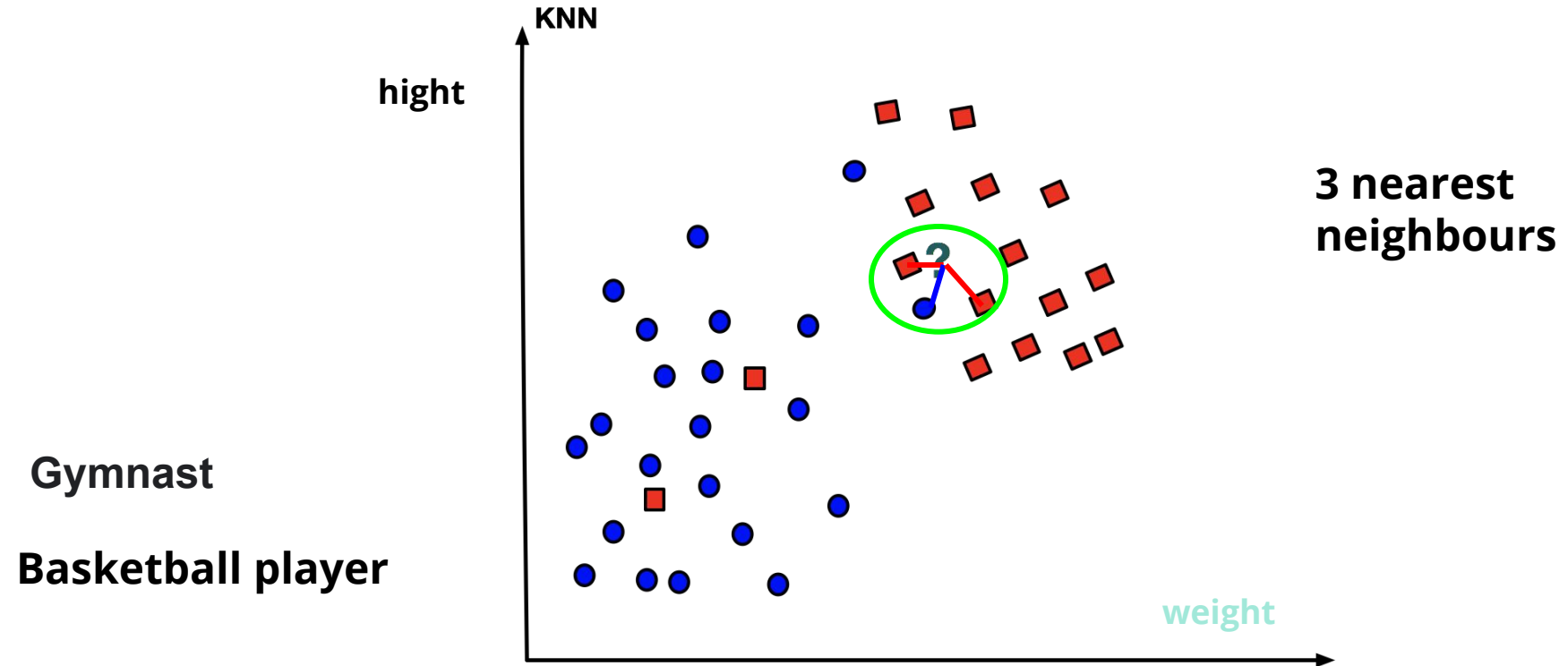
## K Nearest Neighbor

Pick a value  $k$

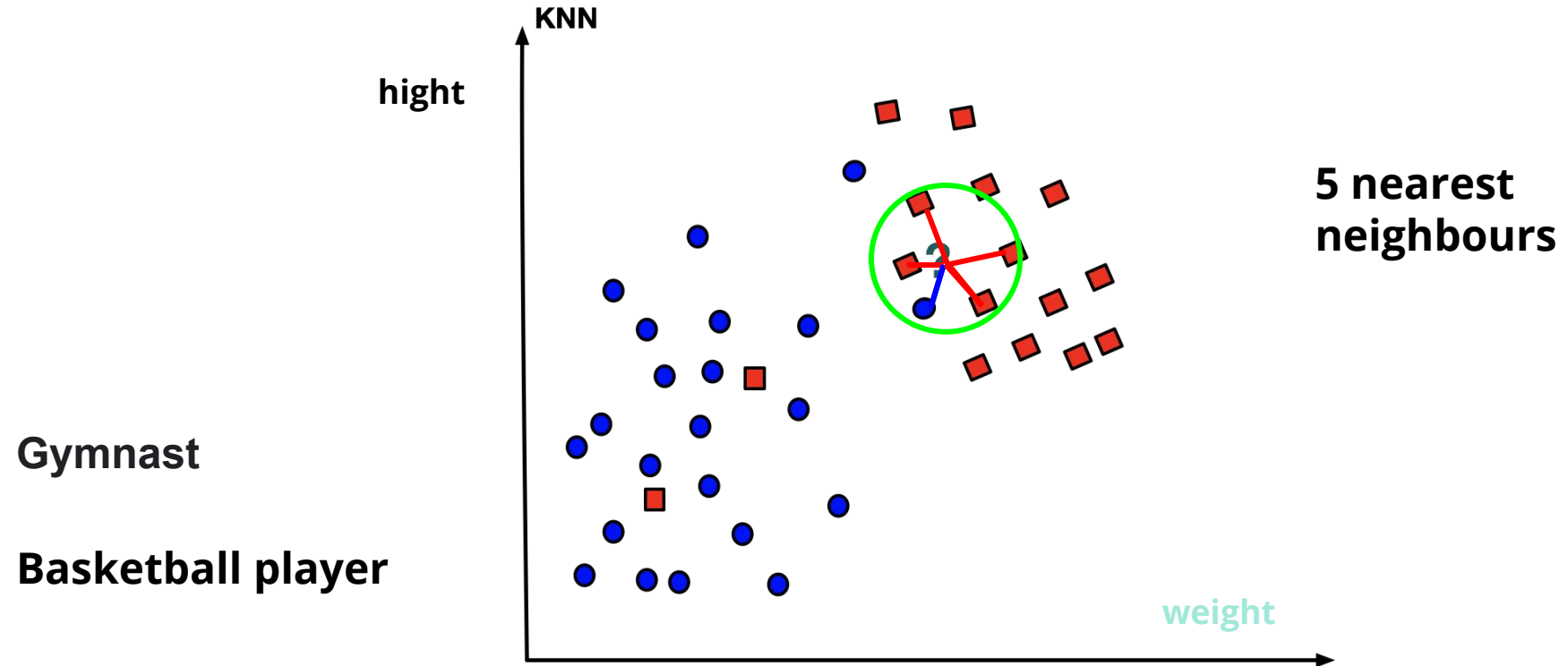
Use  $x$ 's K-Nearest  
Neighbors to vote  
on what  $x$ 's label  
should be.



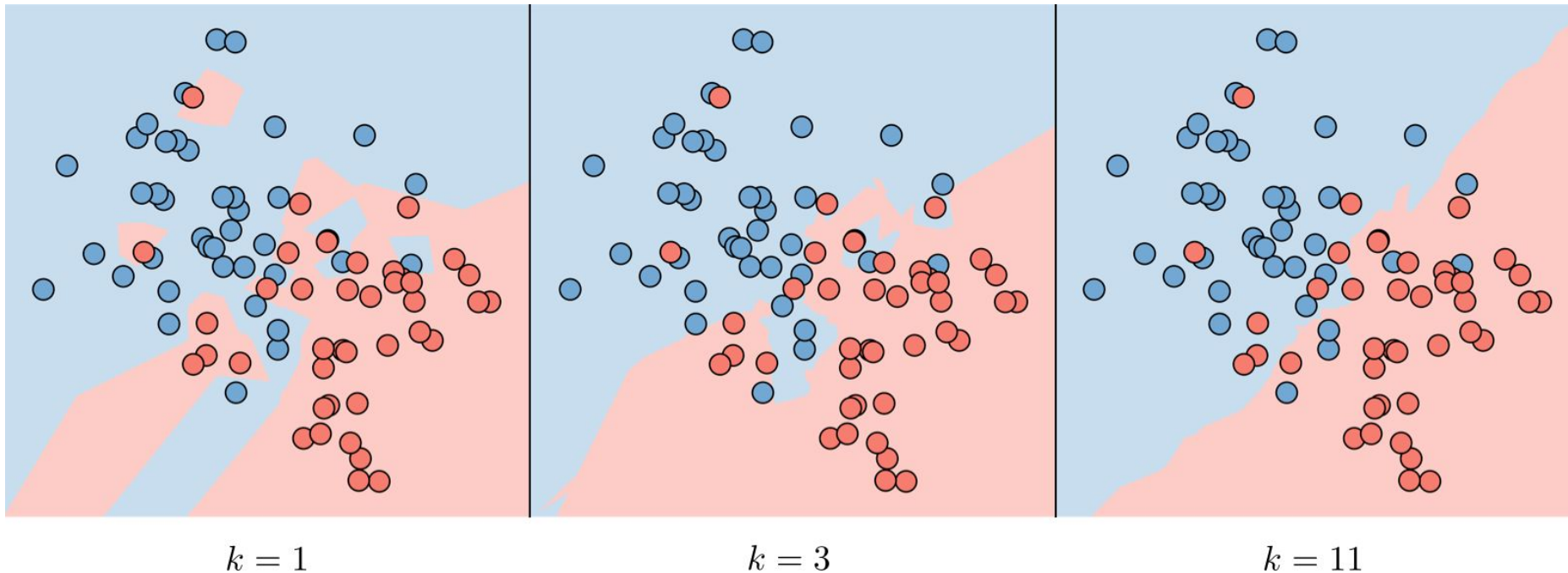
# K Nearest Neighbor



# K Nearest Neighbor



# K Nearest Neighbor



# Iris DataSet



*Iris virginica*



*Iris setosa*



*Iris versicolor*

Petal

Sepal

# Linear Regression



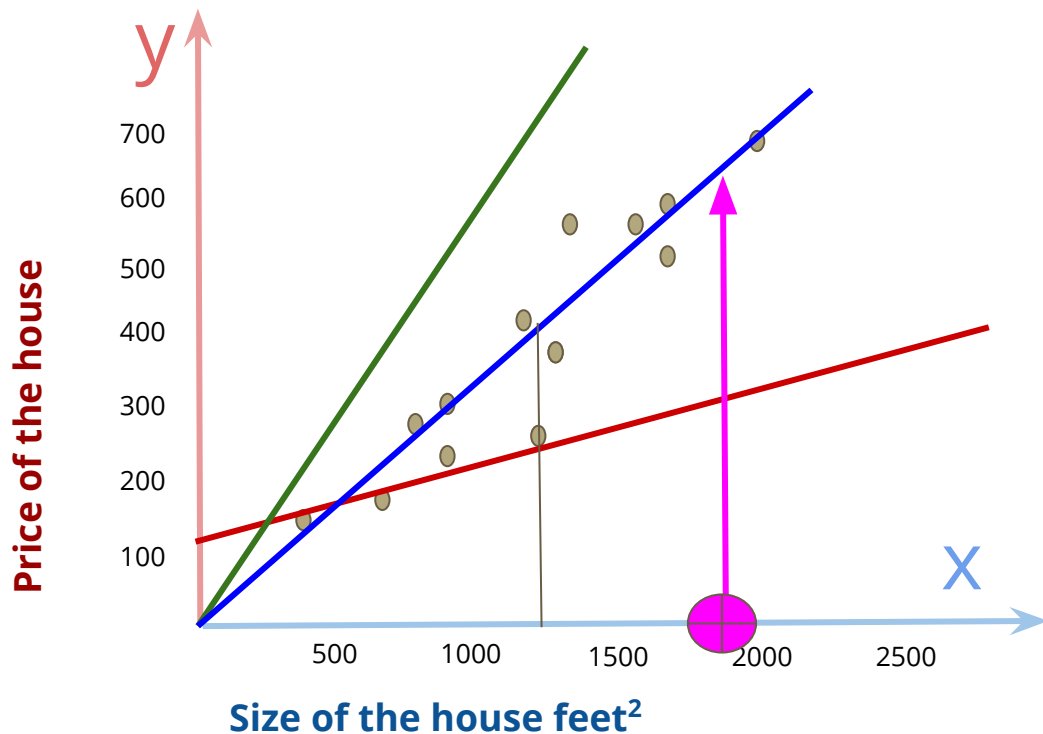


# Linear Regression

Linear regression is the simplest and most widely used statistical technique

A linear model expresses the target output value in terms of a sum of weighted input variables.

# Linear Regression



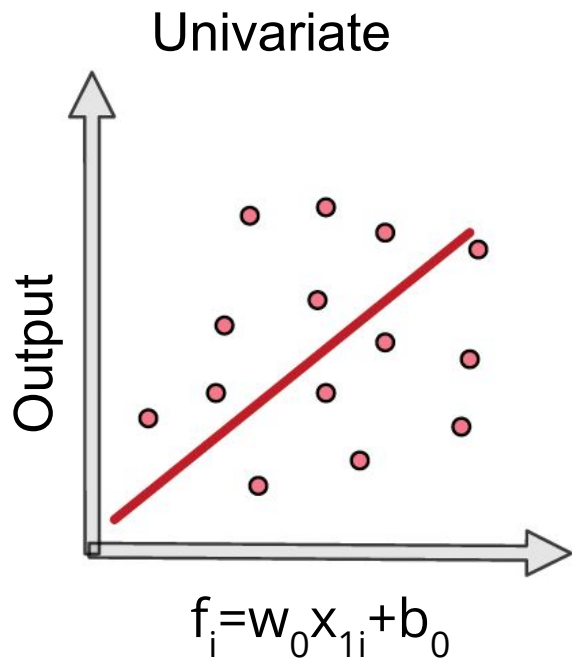
$$f_i = w_0 x_i + b_0$$

Mean squared error

$$MSE = \frac{1}{N} \sum_{i=1}^N (f_i - y_i)^2$$

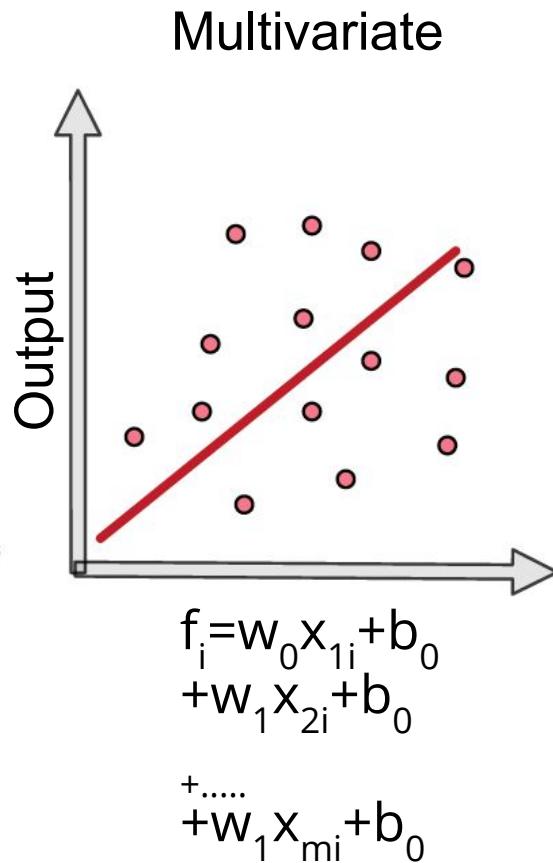
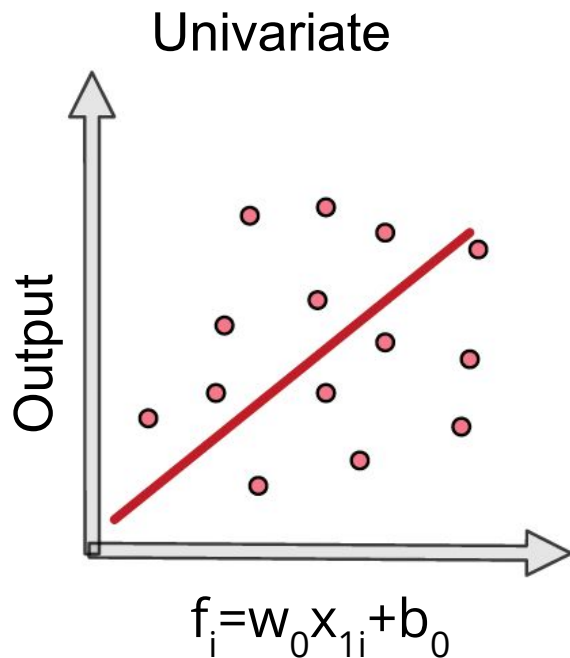
where  $N$  is the number of data points,  
 $f_i$  the value returned by the model and  
 $y_i$  the actual value for data point  $i$ .

# Univariate versus Multivariate Modeling



$$MSE = \frac{1}{N} \sum_{i=1}^N (f_i - y_i)^2$$

# Univariate versus Multivariate Modeling



$$MSE = \frac{1}{N} \sum_{i=1}^N (f_i - y_i)^2$$

# Diabetes

## Ten baseline variables:

age, sex, body mass index, average blood pressure, and six blood serum measurements

n = 442 diabetes patients

## Target value:

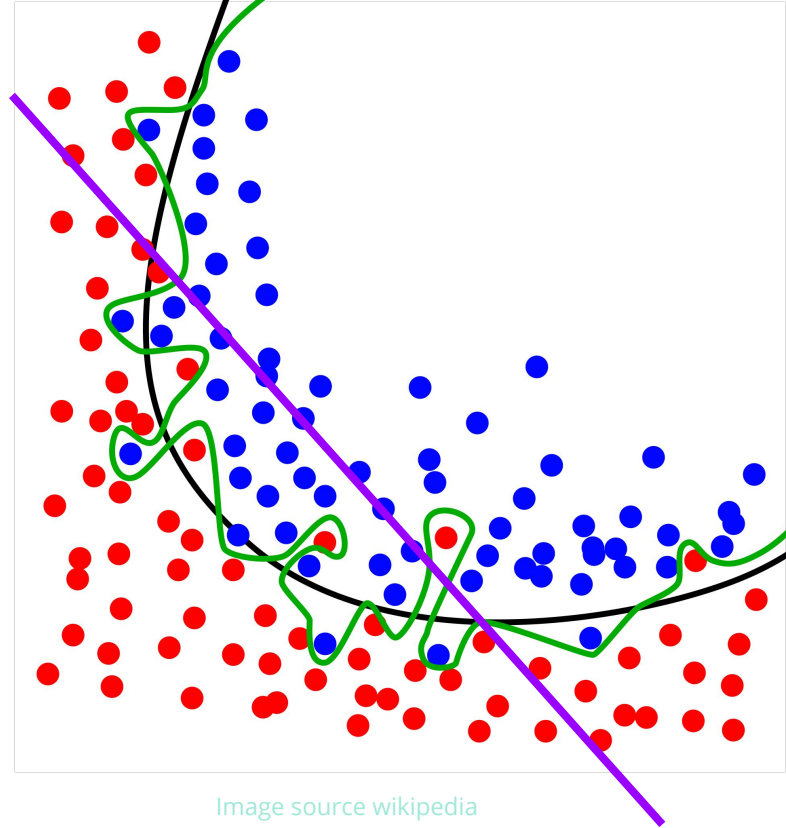
A quantitative measure of disease progression one year after baseline.



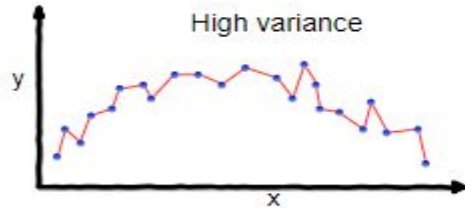
# Overfitting

**Overfitting:** Good performance on the training data, poor generalization to other data.

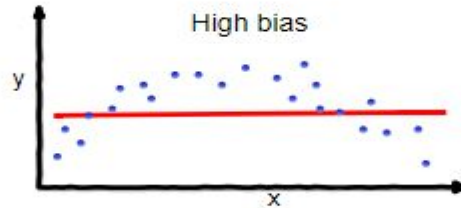
**Underfitting:** Poor performance on the training data and poor generalization to other data



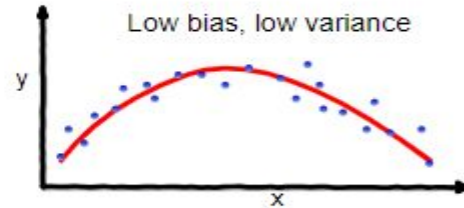
# Bias-Variance Tradeoff



**overfitting**

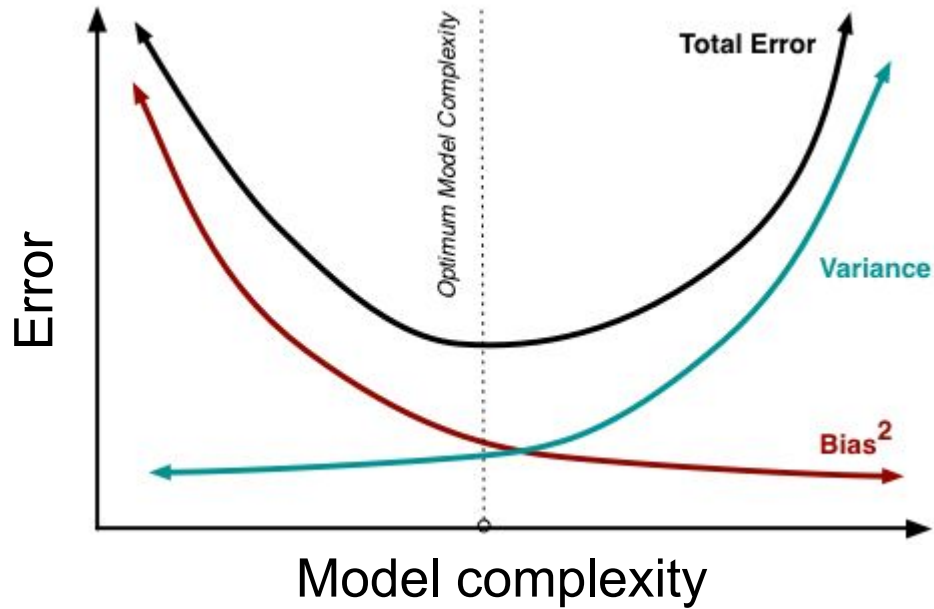


**underfitting**



**Good balance**

# Bias–Variance Tradeoff

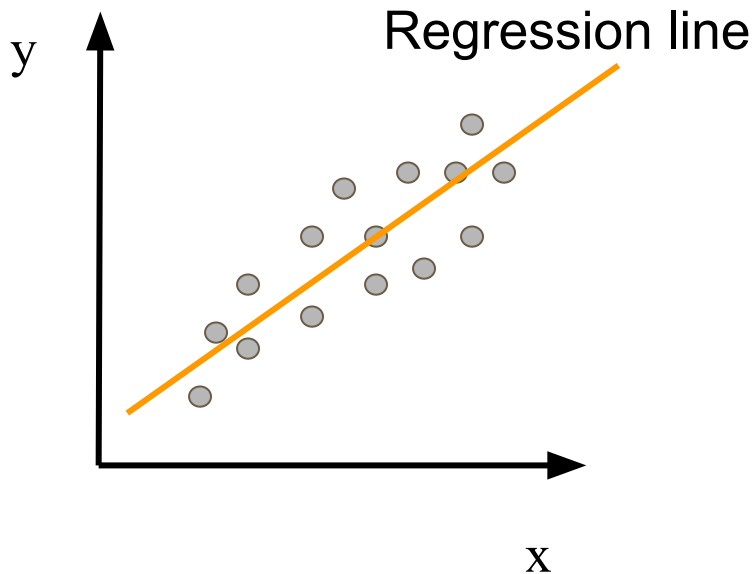




# Logistic Regression

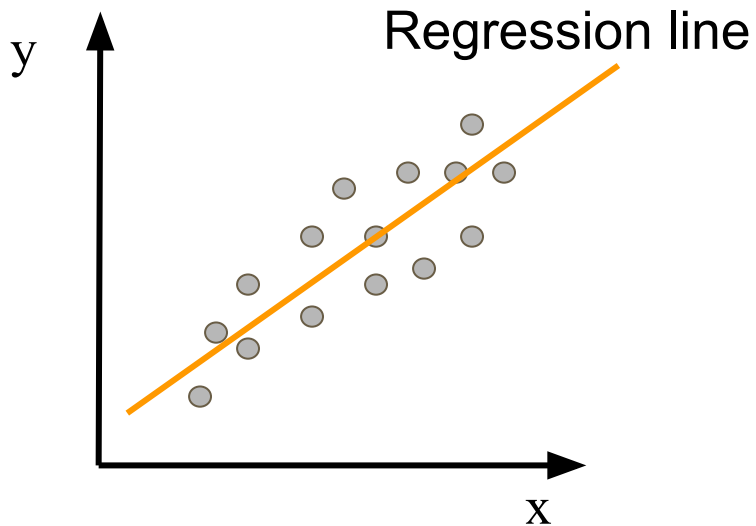
# Linear versus Logistic Regression

## Regression (linear regression)

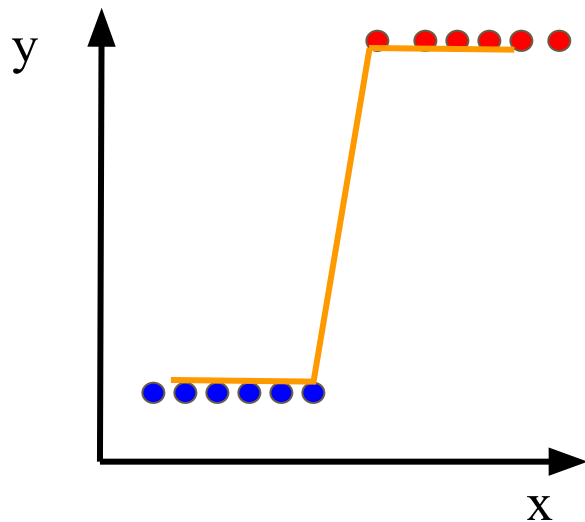


# Linear versus Logistic Regression

**Regression (linear regression)**



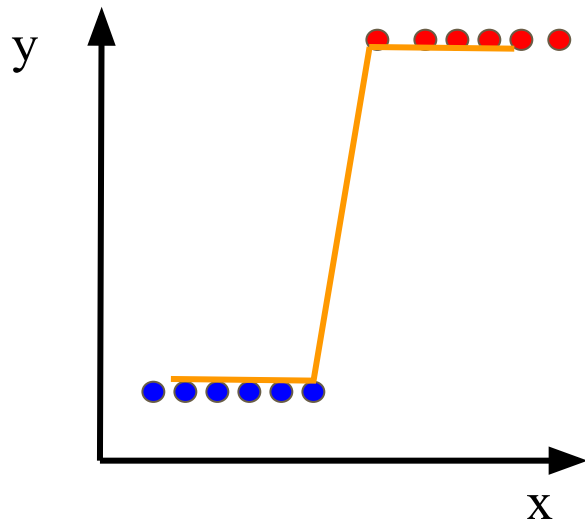
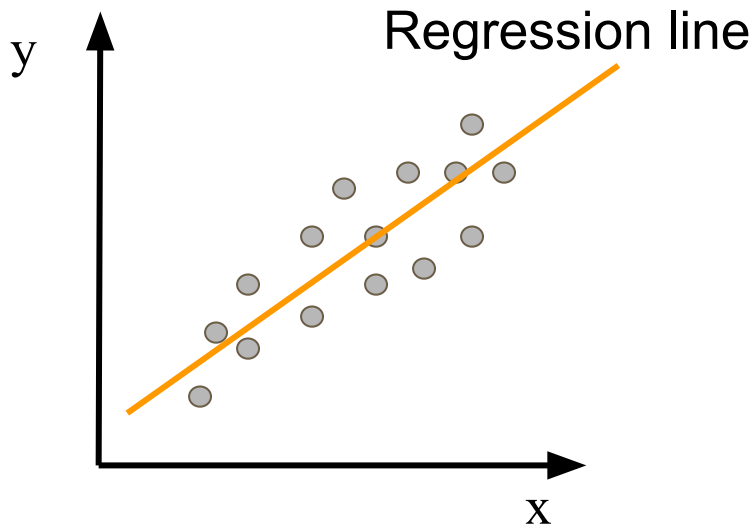
**Classification (logistic regression)**



# Linear versus Logistic Regression

**Regression (linear regression)**

**Classification (logistic regression)**



We need a smooth function that gives us this trend (Sigmoid Function)

# Linear versus Logistic Regression

Linear regression

$$f_i = \sum_i w_i x_i + b_0$$

Logistic regression

$$f_i = \frac{1}{1 + e^{\sum_i w_i x_i + b_0}}$$

W will be identified to minimize cost

$$\text{Cost}(w) = \text{function}(w, f_i, y_i)$$

# Bayes rule and Naive Bayes classifier

# What we know when training a model

$$p(X_1 = x_1 \mid \text{Class} = 1)$$

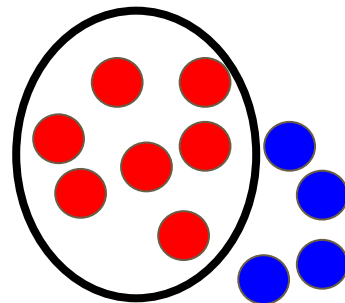
$$p(X_2 = x_2 \mid \text{Class} = 1)$$



$$p(X_m = x_m \mid \text{Class} = 1)$$

● Class=1

● Class=2

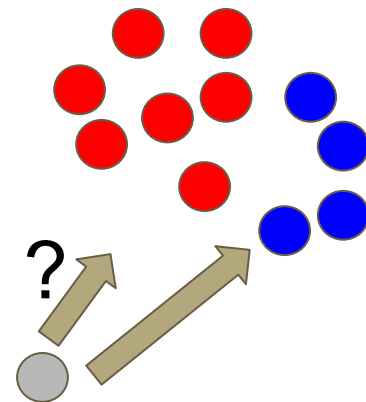


# What do we care about?

● Class=1

● Class=2

$$p(\text{Class}=\textcolor{red}{1} \mid \mathbf{X}_1=\mathbf{x}_1, \mathbf{X}_2=\mathbf{x}_2, \dots, \mathbf{X}_m=\mathbf{x}_m) = ?$$





# Bayes rule is useful to figure out the relationship

$$p(A|B)p(B)=p(B|A)p(A)$$

$$p(\text{Class}=\textcolor{red}{1} | \mathbf{X}_1=\mathbf{x}_1, \mathbf{X}_2=\mathbf{x}_2, \dots, \mathbf{X}_m=\mathbf{x}_m) *$$

$$p(\mathbf{X}_1=\mathbf{x}_1, \mathbf{X}_2=\mathbf{x}_2, \dots, \mathbf{X}_m=\mathbf{x}_m) =$$

$$p(\mathbf{X}_1=\mathbf{x}_1, \mathbf{X}_2=\mathbf{x}_2, \dots, \mathbf{X}_m=\mathbf{x}_m | \text{Class}=\textcolor{red}{1}) p(\text{Class}=\textcolor{red}{1})$$

# The relationship looks complicated

$$\text{WWW} * p(\mathbf{X}_1=\mathbf{x}_1, \mathbf{X}_2=\mathbf{x}_2, \dots, \mathbf{X}_m=\mathbf{x}_m) = \\ p(\mathbf{X}_1=\mathbf{x}_1, \mathbf{X}_2=\mathbf{x}_2, \dots, \mathbf{X}_m=\mathbf{x}_m \mid \text{Class}=\mathbf{1}) p(\text{Class}=\mathbf{1})$$

**WWW**: What We Want

$$p(\text{Class}=\mathbf{1}) : \text{easy to calculate} \quad p(\text{Class} = i) = \frac{N_i}{\sum_i^C N_i}$$

# Naive Bayes

Naive Bayes classifier is called **Naive** as it assumes each feature will independently contribute in prediction of a class for each data point

$$p(X_1=x_1, X_2=x_2, \dots, X_m=x_m) = p(X_1=x_1) p(X_2=x_2) \dots p(X_m=x_m)$$

$$p(X_1=x_1, X_2=x_2, \dots, X_m=x_m | \text{Class}=\textcolor{red}{1}) =$$

$$p(X_1=x_1 | \text{Class}=\textcolor{red}{1}) p(X_2=x_2 | \text{Class}=\textcolor{red}{1}) \dots p(X_m=x_m | \text{Class}=\textcolor{red}{1})$$

# Bayesian approach for problem solving

# Bayesian versus frequentist

## Frequentist:

- Variation of data and derived quantities in terms of fixed model parameters

## Bayesian:

- Variation of beliefs about parameters in terms of fixed observed data

**Note.** The difference become clear in complicated problems like *Interpretation of uncertainty*.

# Bayesian versus frequentist

## Frequentist:

- If an experiment is repeated many times, in 95% of these cases the computed confidence interval will contain the true  $\theta$ .

**Note.** In general, a frequentist 95% confidence interval is not 95% likely to contain the true value. This very common mistake is a Bayesian interpretation of a frequentist construct.

# Bayesian versus frequentist

## Frequentist:

- If an experiment is repeated many times, in 95% of these cases the computed confidence interval will contain the true  $\theta$ .

**Note.** In general, a frequentist 95% confidence interval is not 95% likely to contain the true value. This very common mistake is a Bayesian interpretation of a frequentist construct.

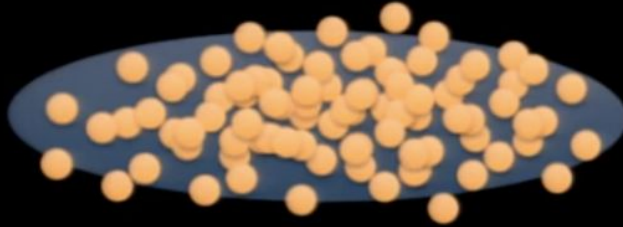
## Bayesian:

- Given our observed data there is a 95% probability that the value of  $\theta$  lies within the credible region

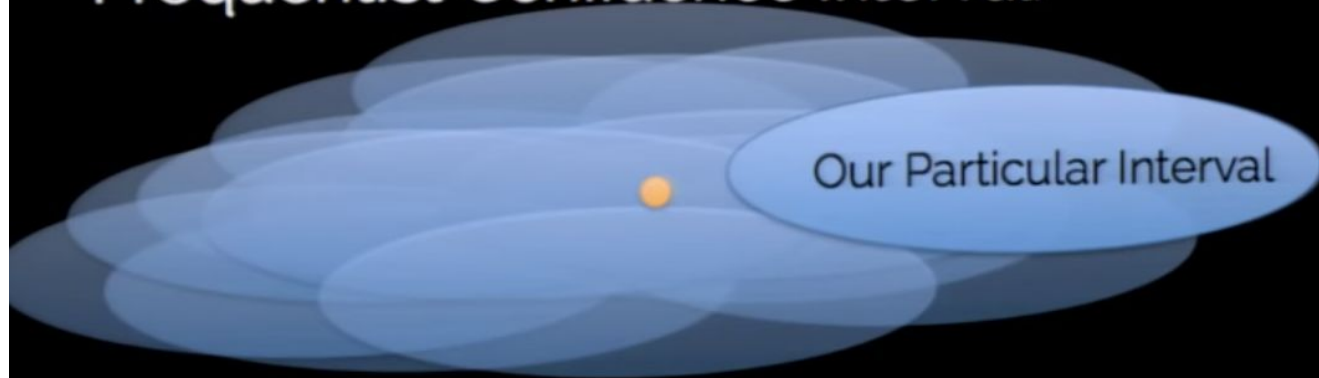
# Confidence vs. Credibility

- - Parameter
- - Interval

Bayesian Credible Region:



Frequentist Confidence Interval:





# Conversation between a (frequentist) statistician and a scientist:

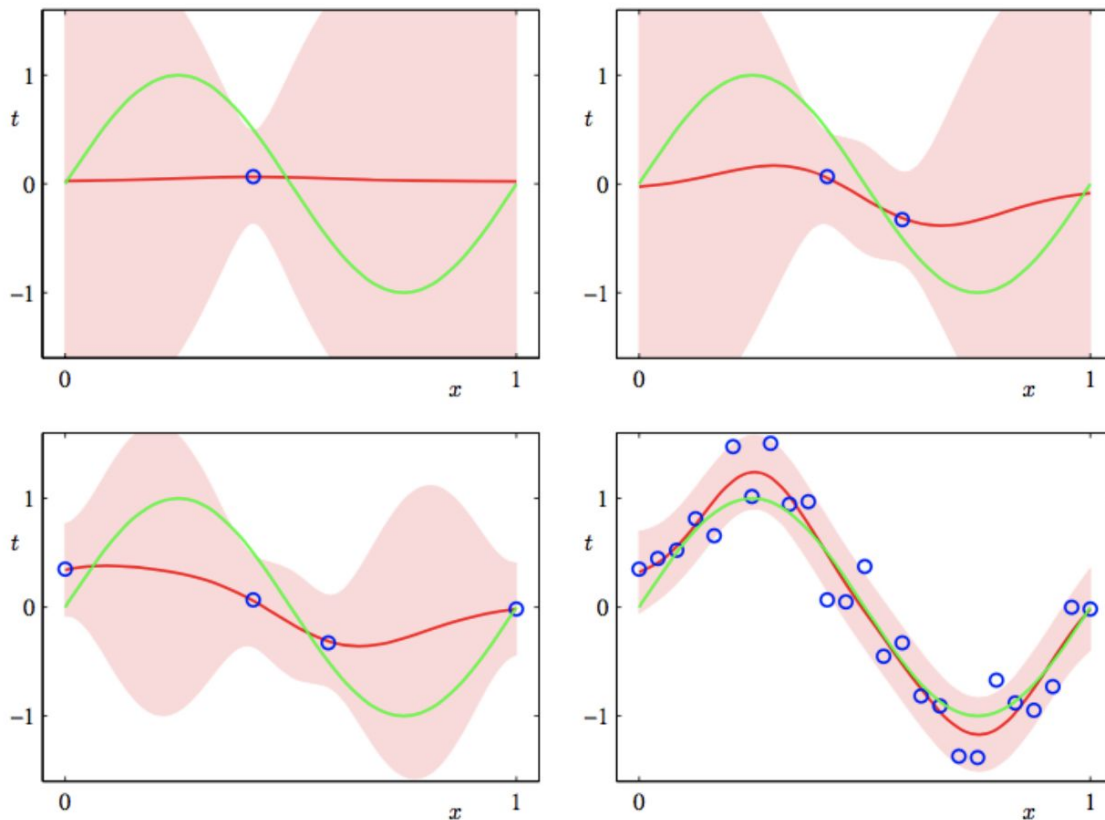
## **Statistician:**

- 95% of such confidence intervals in repeated experiments will contain the true value. (Referring to chance in that context is meaningless. The 95% refers to the interval itself)
  - The long term limiting frequency of the procedure or constructing this interval ensures that 95% of the resulting ensemble of intervals contains the value.

## **Scientist:**

- So there is a 95% chance that the value is in this interval?

# Example in regression



Pattern Recognition  
and machine learning,  
Bishop.





# Iris DataSet



*Iris virginica*



*Iris setosa*



*Iris versicolor*

Petal

Sepal

# Diabetes

## Ten baseline variables:

age, sex, body mass index, average blood pressure, and six blood serum measurements

n = 442 diabetes patients

## Target value:

A quantitative measure of disease progression one year after baseline.



# Breast cancer dataset

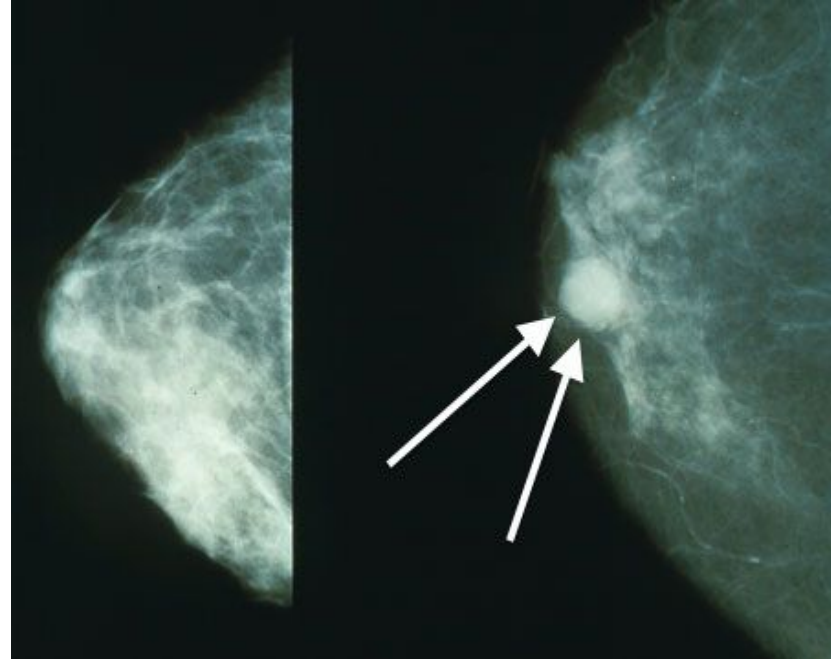
The breast cancer dataset is a classic and very easy binary classification dataset.

## Features :

Computed from a digitized image of a fine needle aspirate (FNA) of a breast mass.

## Target values:

Benign /Malignant



# Extra useful information



# Useful links

## Installation instructions

- [scikit-learn](#)
- [IPython](#)

## Data Sets

- [scikit-learn DataSet](#)

## scikit-learn: machine learning in Python :

- <https://scikit-learn.org/stable/>

## Useful cheat sheets:

- <https://www.analyticsvidhya.com/blog/2017/02/top-28-cheat-sheets-for-machine-learning-data-science-probability-sql-big-data/>