



UNIVERSITY OF  
BIRMINGHAM

# Week 3. Classification and Logistic Regression

**Dr. Shuo Wang**



# Overview

- Classification: parametric and non-parametric
- Logistic regression algorithm
- K-Nearest Neighbours (kNN)



# Some concepts

Recall:

- Classification: predict categorical labels, e.g. spam detection.
- Regression: predict real values, e.g. stock price prediction.

Some learning algorithms:

- Linear regression: a linear and **parametric** model for regression problems.
- Logistic regression: a linear and **parametric** model for **classification** problems (contrary to its name!)
- K-Nearest Neighbours (KNN) : a **non-parametric** model that can be used for both classification and regression problems.

# Parametric and Non-parametric Models

## Parametric models:

- A model that summarizes data with a finite set of parameters.
- Make assumptions on data distributions.
- E.g. linear/logistic regression, neural networks

## Non-parametric models:

- A model that cannot be characterized by a bounded set of parameters.
- No assumptions on data distributions.
- E.g. instance-based learning that generate hypotheses using training examples, including kNN, SVM, decision trees, etc.

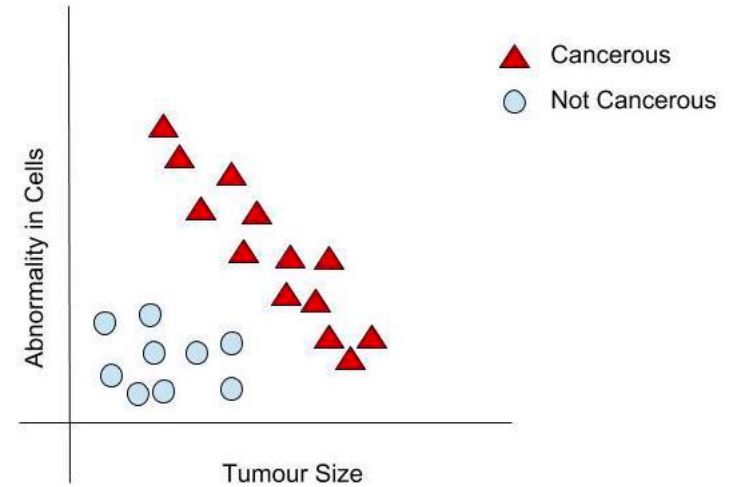
# Logistic regression

Similar to linear regression,

- 1) Model formulation
- 2) Cost function
- 3) Learning algorithm by gradient descent

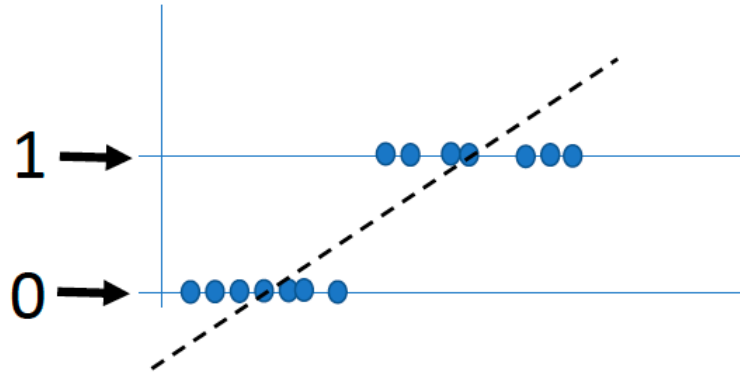
# Model formulation

- We want to put a boundary between 2 classes.



Decision boundary

# Can we use linear regression to classify them?



# Model

- We change the linear model slightly by passing it through a nonlinearity.
- If  $x$  has 1 attribute, we have:

$$h(x; \mathbf{w}) = \sigma(w_0 + w_1x) = \frac{1}{1 + e^{-(w_0 + w_1x)}}$$

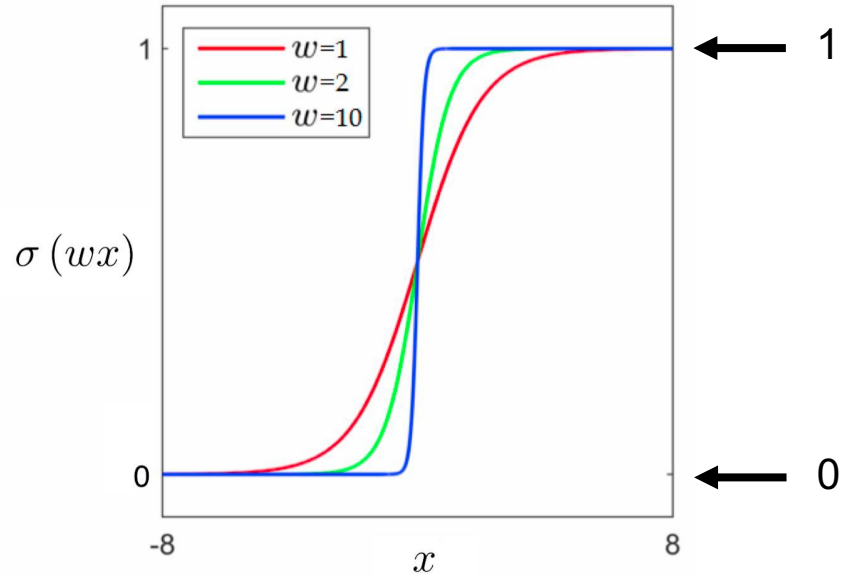
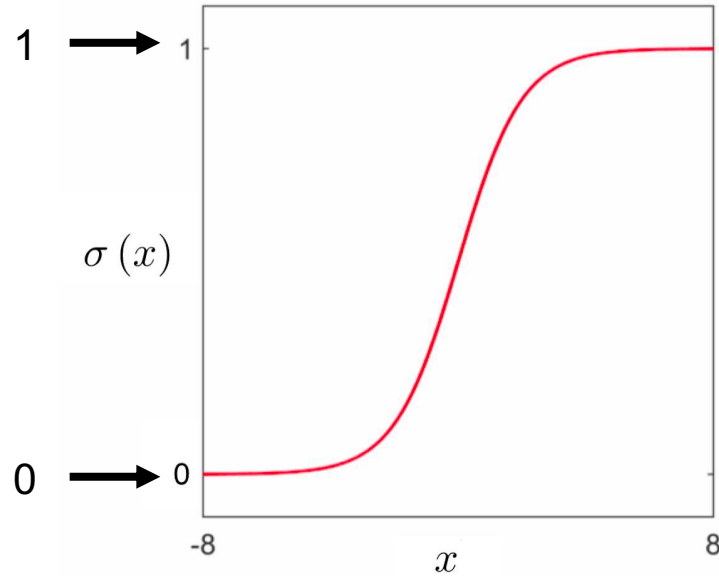
- The function  $\sigma(u) = \frac{1}{1+e^{-u}}$  is called **sigmoid** function or **logistic** function.

Go to <https://www.desmos.com/calculator>

Type:  $y = 1/(1+\exp(-(w_0+w_1x)))$



# Sigmoid function



- It is a smoothed version of a step function (e.g.  $-1$  for negative numbers and  $+1$  for positive numbers).
- Also seen in neural networks.

# Model

- If  $\mathbf{x}$  has  $d$  attributes, we have:

$$h(\mathbf{x}; \mathbf{w}) = \sigma(w_0 + w_1x_1 + \cdots + w_dx_d) = \frac{1}{1+e^{-(\mathbf{w}^T\mathbf{x})}}, \text{ where}$$

All components of  $\mathbf{w}$  are free parameters.

$$\mathbf{w} = \begin{pmatrix} w_0 \\ w_1 \\ w_2 \\ \vdots \\ w_d \end{pmatrix}, \mathbf{x} = \begin{pmatrix} 1 \\ x_1 \\ x_2 \\ \vdots \\ x_d \end{pmatrix} \in R^d$$

# Meaning of the sigmoid function

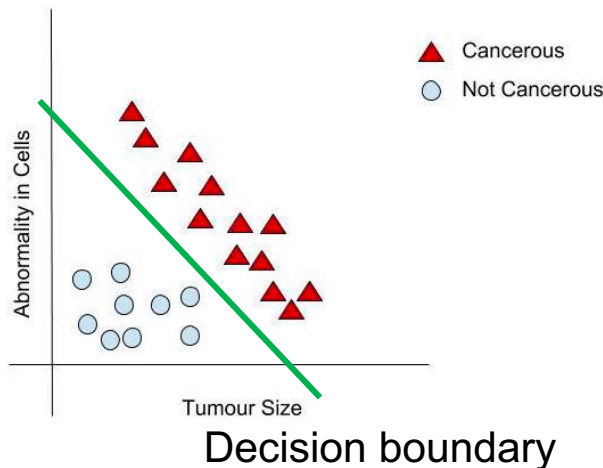
- The sigmoid function takes a single argument (note,  $\mathbf{w}^T \mathbf{x}$  is one number)
- It always returns a value between 0 and 1. The meaning of this value is **the probability that the label is 1**.

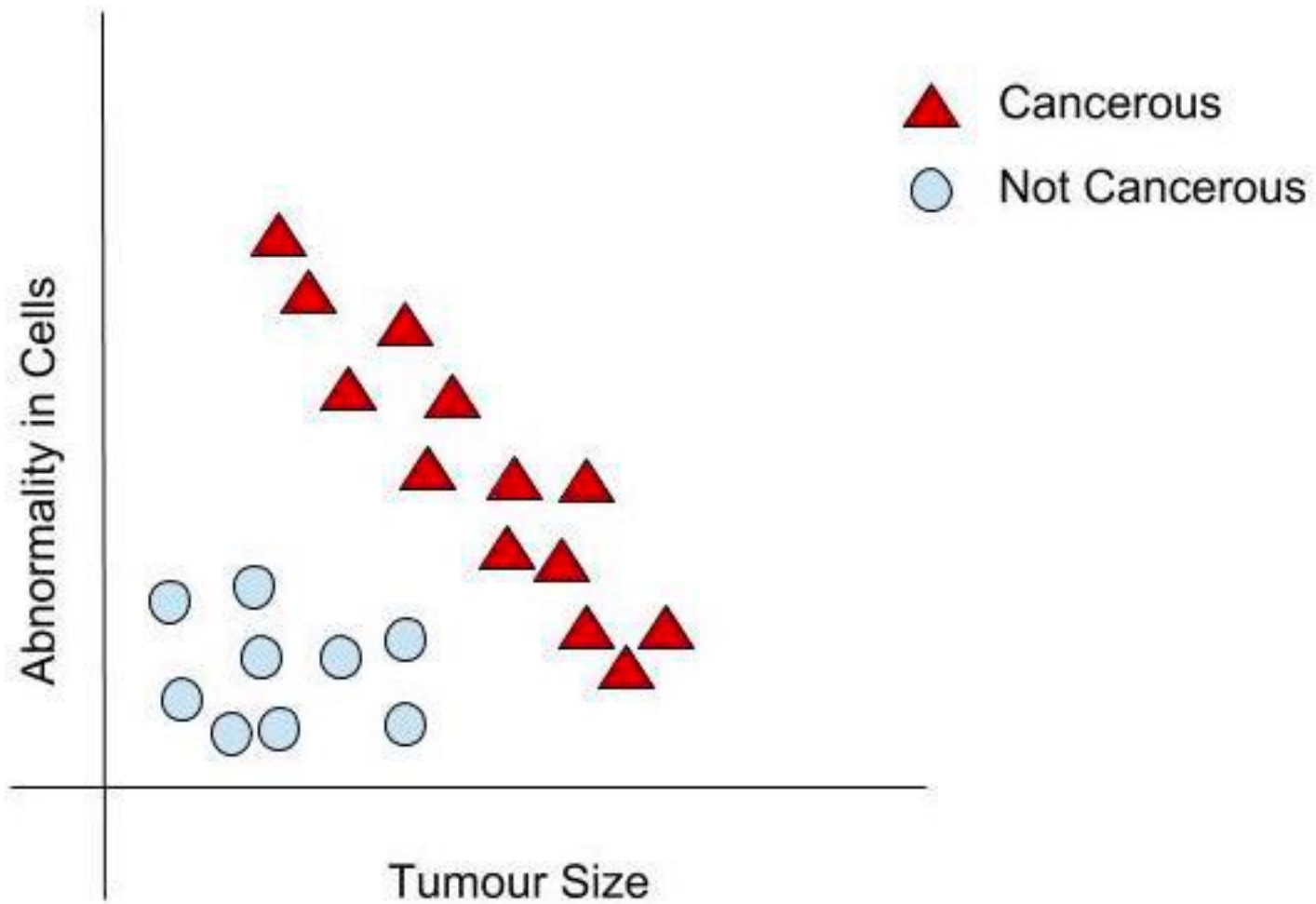
$$\sigma(\mathbf{w}^T \mathbf{x}) = P(y = 1 | \mathbf{x}; \mathbf{w})$$

If it is smaller than 0.5, then we predict label 0.

If it is larger than 0.5, then we predict label 1.

- There is a slim chance that the sigmoid outputs exactly 0.5. The set of all possible inputs for which this happens is called the **decision boundary**.



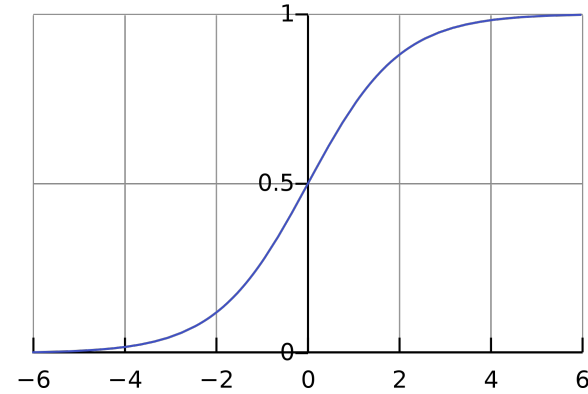


# Example

- Suppose we have 2 input attributes, so our model is

$$h(\mathbf{x}; \mathbf{w}) = \sigma(w_0 + w_1x_1 + w_2x_2)$$

- Suppose we know that  $w_0 = -1$ ,  $w_1 = 1$ ,  $w_2 = 1$ .
- Q1: When do we predict 1? What is the decision boundary?
- Q2: Is the decision boundary of logistic regression always linear?



# Logistic regression

Similar to linear regression,

- 1) Model formulation
- 2) Cost function
- 3) Learning algorithm by gradient descent

# Cost function - Recall

- The loss expresses an error, so it must be always non-negative.
- Absolute value loss (L1 loss):

$$L1 = |f(x) - y|$$

- Mean squared error loss (L2 loss): *used in linear regression*

$$L2 = (f(x) - y)^2$$

$$g(w_0, w_1) = \frac{1}{N} \sum_{n=1}^N (f(x^{(n)}; w_0, w_1) - y^{(n)})^2$$

*Empirical loss  
used by LR*

Loss for the n-th training example

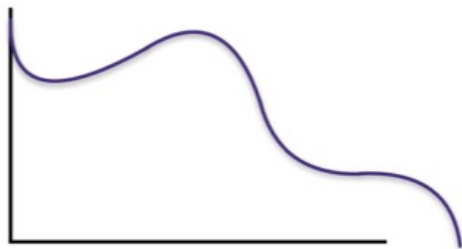
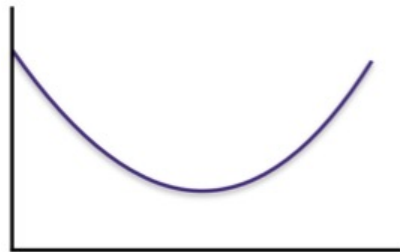
- 0/1 loss:

$$L_{0/1} = 0 \text{ if } f(x) = y, \text{ else } 1$$

# Cost function

- MSE using sigmoid function does not work.
  - The MSE function becomes concave(not convex) – too wriggly (due to discrete output labels and bounded sigmoid output between (0,1)).
  - Gradient descent does not work well on non-convex functions. (local minimum)
  - Go to <https://www.desmos.com/calculator>, and observe:

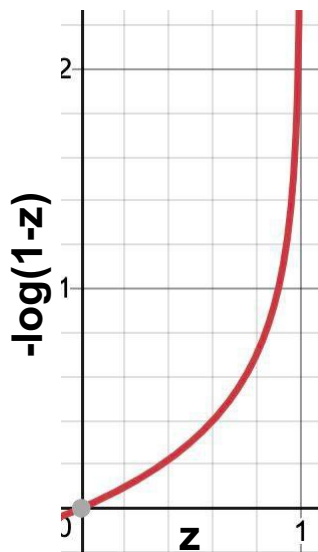
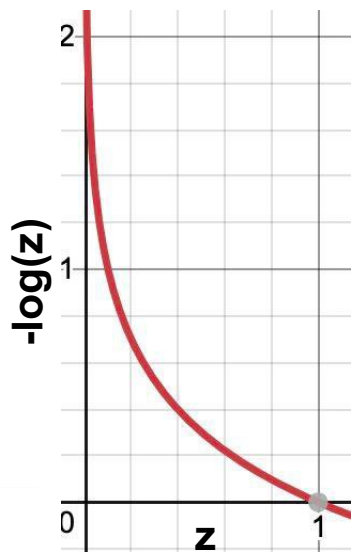
$$y = \left( 1 - \frac{1}{1 + e^{-(w_0 + w_1 x)}} \right)^2$$





# Logistic cost function

For each  $(x, y)$  pair,  $cost(h(x; \mathbf{w}), y) = \begin{cases} -\log(\overbrace{h(x; \mathbf{w})}^{\mathbf{z}}), & \text{if } y = 1 \\ -\log(1 - \underbrace{h(x; \mathbf{w})}_{\mathbf{z}}), & \text{if } y = 0 \end{cases}$



# Write the cost function in a single line

- $g(\mathbf{w}) = \frac{1}{N} \sum_{n=1}^N \text{cost}(h(\mathbf{x}^{(n)}; \mathbf{w}), y^{(n)})$ , where
$$\text{cost}(h(\mathbf{x}; \mathbf{w}), y) = \begin{cases} -\log(h(\mathbf{x}; \mathbf{w})), & \text{if } y = 1 \\ -\log(1 - h(\mathbf{x}; \mathbf{w})), & \text{if } y = 0 \end{cases}$$
- $g(\mathbf{w}) = -\frac{1}{N} \sum_{n=1}^N (y^{(n)} \log h(\mathbf{x}^{(n)}; \mathbf{w}) + (1 - y^{(n)}) \log(1 - h(\mathbf{x}^{(n)}; \mathbf{w})))$

This logistic loss is also called **cross-entropy** loss.

# Logistic regression – what we want to do

- Given training data

$$(x^{(1)}, y^{(1)}), (x^{(2)}, y^{(2)}), \dots, (x^{(N)}, y^{(N)}), \text{ where } y \in \{0, 1\}$$

- Fit the model

$$y = h(\mathbf{x}; \mathbf{w}) = \sigma(\mathbf{w}^T \mathbf{x})$$

- By minimizing the cross-entropy cost function

$$g(\mathbf{w}) = -\frac{1}{N} \sum_{n=1}^N (y^{(n)} \log h(\mathbf{x}^{(n)}; \mathbf{w}) + (1 - y^{(n)}) \log(1 - h(\mathbf{x}^{(n)}; \mathbf{w})))$$

# Logistic regression

Similar to linear regression,

- 1) Model formulation
- 2) Cost function
- 3) Learning algorithm by gradient descent

# Training by gradient descent

- We use gradient descent (again like linear regression) to minimize the cost function, i.e. to find the best weight values  $\mathbf{w}$ .
- The gradient vector is\*:

$$\nabla g(\mathbf{w}) = -\left(y^{(n)} - h(\mathbf{x}^{(n)}; \mathbf{w})\right) \mathbf{x}^{(n)}, \text{ where } \mathbf{w} = \begin{pmatrix} w_0 \\ w_1 \\ w_2 \\ \vdots \\ w_d \end{pmatrix}, \mathbf{x} = \begin{pmatrix} 1 \\ x_1 \\ x_2 \\ \vdots \\ x_d \end{pmatrix} \in R^d$$

- We plug this into the general gradient descent algorithm given last week.

\* This follows after differentiating the cost function w.r.t.w –we omit the lengthy math!

# Training by gradient descent

While not converged

for  $n = 1, 2 \dots N$  //each example in the training set

$$\mathbf{w} := \mathbf{w} + \alpha(y^{(n)} - h(\mathbf{x}^{(n)}; \mathbf{w}))\mathbf{x}^{(n)}$$

Return  $\mathbf{w}$ .

The same, written component-wise

While not converged

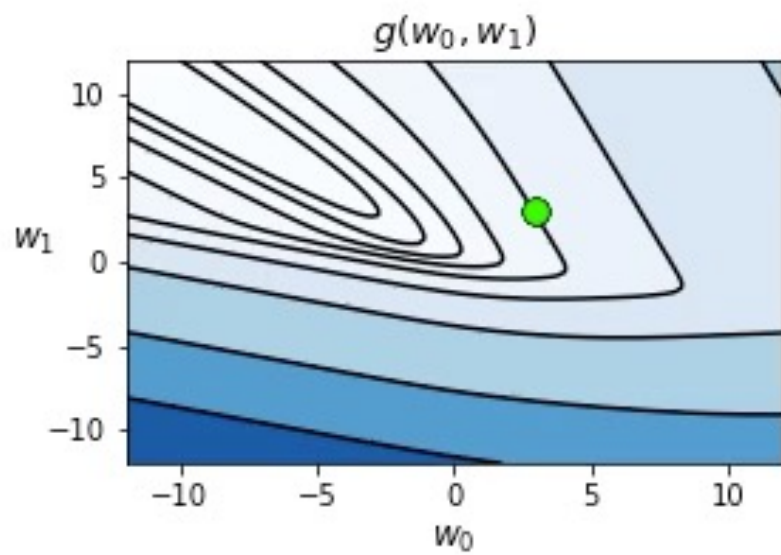
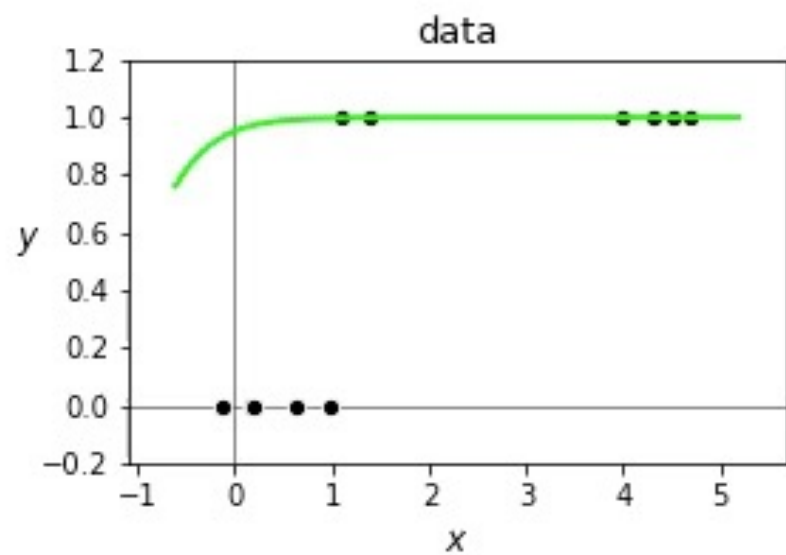
for  $n = 1, 2 \dots N$  // each example in the training set

$$w_0 := w_0 + \alpha(y^{(n)} - h(\mathbf{x}^{(n)}; \mathbf{w}))$$

for  $i = 1, \dots, d$

$$w_i := w_i + \alpha(y^{(n)} - h(\mathbf{x}^{(n)}; \mathbf{w}))x_i^{(n)}$$

Return  $\mathbf{w}$ .



# Extensions

- We studied logistic regression for linear binary classification
- There are extensions, such as:
  - Nonlinear logistic regression: instead of linear function inside the exp in the sigmoid, we can use polynomial functions of the input attributes
  - Multi-class logistic regression: uses a multi-valued version of sigmoid
- Details of these extensions are beyond of our scope in this module.



# Examples of application of logistic regression

- Face detection: classes consist of images that contain a face and images without a face
- Sentiment analysis: classes consist of written product-reviews expressing a positive or a negative opinion
- Automatic diagnosis of medical conditions: classes consist of medical data of patients who either do or do not have a specific disease

# Are we done?

- Logistic regression is not the only classifier.
- There are many others: decision trees, KNN, neural networks, SVM, etc.
- Which one is the best?

# No Free Lunch (NFL)

Simply to say:

- No single machine learning algorithm is universally the best-performing algorithm for all problems.

Theorem(Wolpert; also Hume 200 years ago):

- Given any distribution that generates the  $x$  of  $S$ , and any training set of size  $N$ .  
For any learner  $A$ ,

$$\frac{1}{|F|} \sum_{f \in F} \text{err}(A(S) \text{ on task } f) = \frac{1}{2}$$

Implication:

- If learner  $A1$  is better than learner  $A2$  for a task  $f$ , then there is another task  $g$  for which learner  $A2$  is better than learner  $A1$ .
- So we need to know many learning methods & try them on the task at hand.

$S$ : training set;  $h$ : a classifier;  $H$ : set of all classifiers;  $A$ : learner that chooses  $h$  from  $H$  based on  $S$ ;  $f$ : task that  $A$  tries to learn;  $F$ : set of all possible tasks;  $\text{err}(h \text{ on task } f)$ : generalisation error.