

Business Analytics

Buzzwords





Περιεχόμενα

- Στατιστική & Έλεγχος Υποθέσεων
- Μηχανική Μάθηση (Machine Learning)
 - Απλή Γραμμική Παλινδρόμηση, Πολλαπλή Παλινδρόμηση (Regression)
 - Λογιστική Παλινδρόμηση (Logistic Regression)
 - Κ-Κοντινότεροι Γείτονες (k-nearest neighbors)
 - Νευρωνικά Δίκτυα (Neural Networks)
 - Βαθιά Μάθηση (Deep Learning)
- Βάσεις δεδομένων και SQL (Structured Query Language)



ΣΤΑΤΙΣΤΙΚΗ

ΜΑΘΗΜΑΤΙΚΕΣ ΚΑΙ
ΤΕΧΝΙΚΕΣ ΕΝΝΟΙΕΣ ΜΕ ΤΙΣ
ΟΠΟΙΕΣ ΚΑΤΑΝΟΟΥΜΕ ΤΑ
ΔΕΔΟΜΕΝΑ

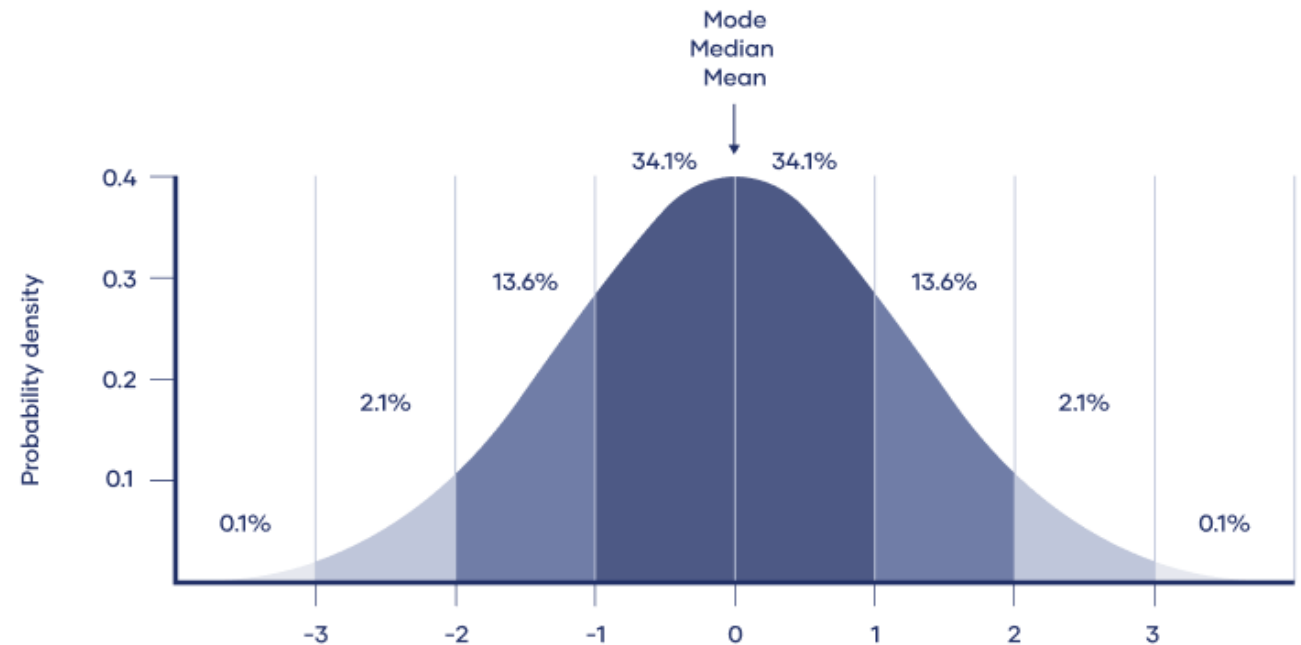




Κεντρικές Τάσεις

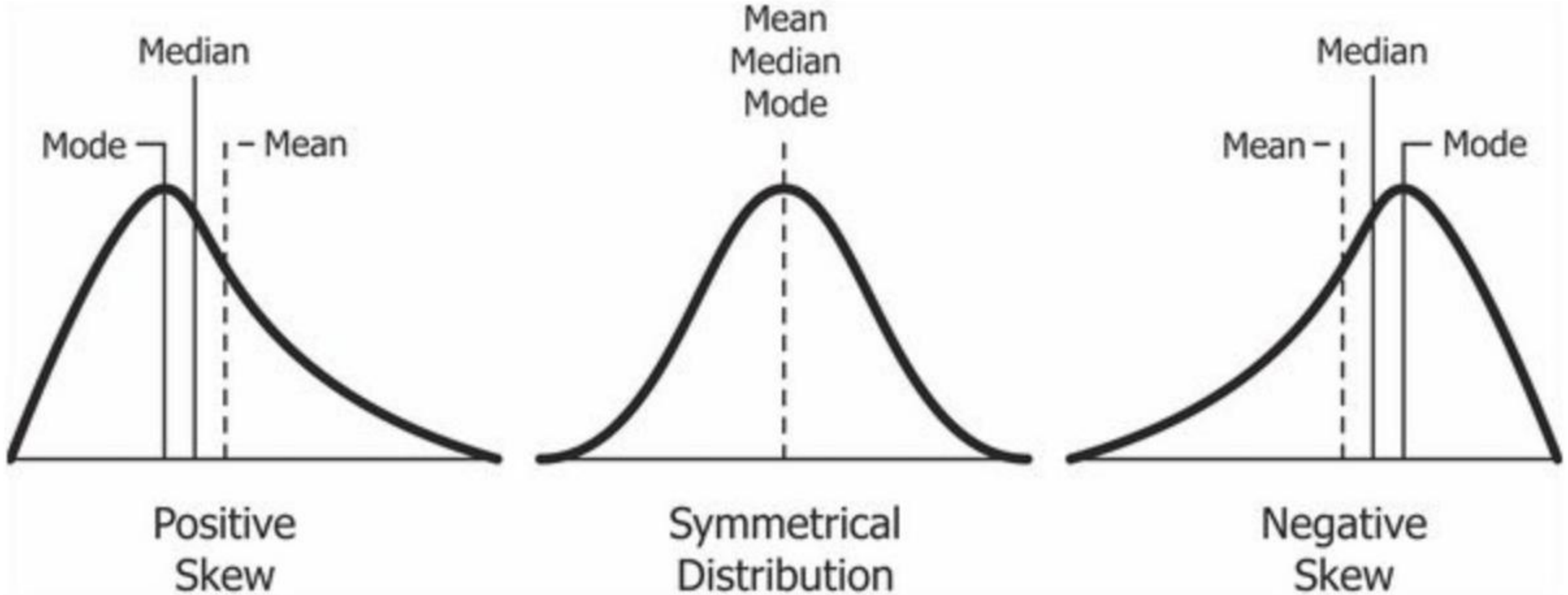
- Συχνά χρειαζόμαστε μια έννοια που να εκφράζει το **κέντρο** των δεδομένων μας. Για αυτό πολλές φορές αναφερόμαστε **στον μέσο όρο** του δείγματος μας. Το άθροισμα δηλαδή δια το πλήθος.
- Βέβαια πολλές φορές είναι **ευαίσθητος** στην ύπαρξη **ακραίων σημείων**, για αυτό κάποιες φορές μπορεί να μας δώσει λάθος συμπεράσματα.

Standard normal distribution





Ασυμμετρία

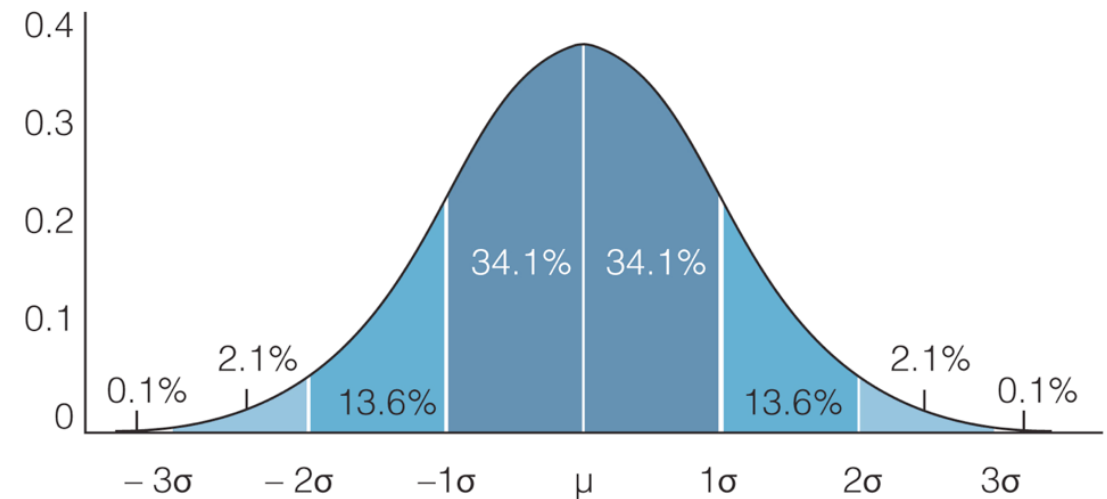




Διασπορά (1/2)

- Ο όρος διασπορά αναφέρεται στο πόσο διασκορπισμένα είναι τα δεδομένα μας. Τα δεδομένα κοντά μηδέν σημαίνουν **καθόλου διασπορά** ενώ μεγάλες τιμές δηλώνουν πως τα δεδομένα είναι **διάσπαρτα**.
- Το range είναι ένας τρόπος έκφρασης της διασποράς αφαιρώντας τον μικρότερο αριθμό από το μεγαλύτερο.
- Ένα άλλο πιο σύνθετο μέτρο διασποράς είναι η **διακύμανση** (var) που στην επόμενη διαφάνεια αποτυπώνουμε πως υπολογίζεται.

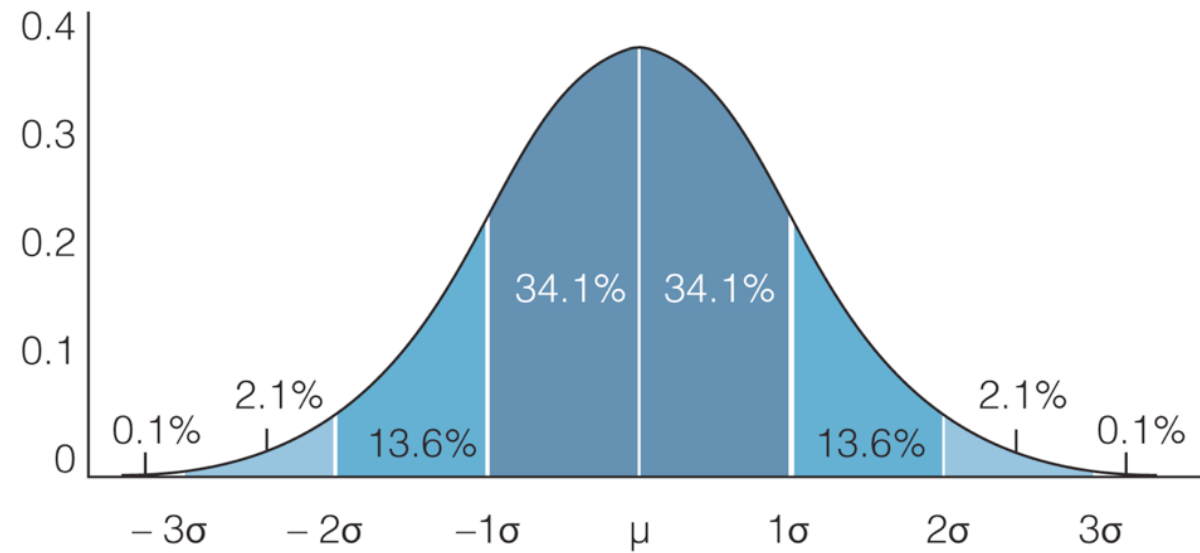
Distribution of Variance





Διασπορά (2/2)

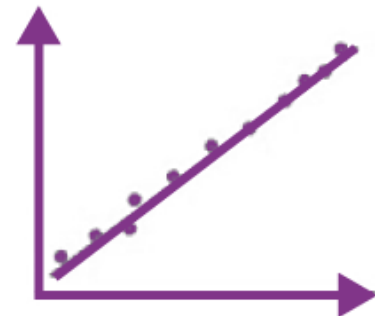
Distribution of Variance



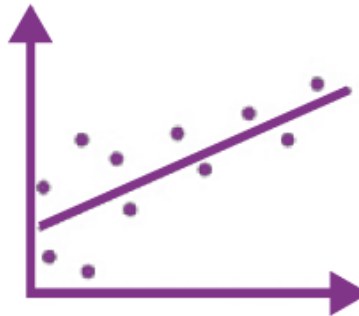


Συσχέτιση

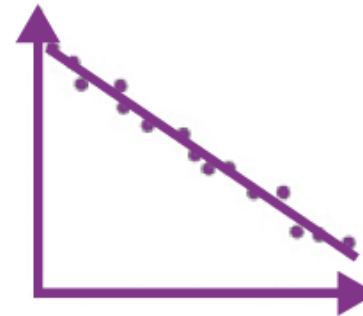
- Η συσχέτιση δεν έχει μονάδες και βρίσκεται μεταξύ του -1 και του 1.



Strong positive correlation



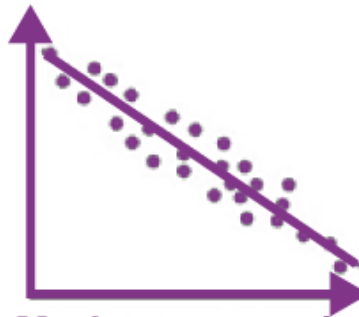
Weak positive correlation



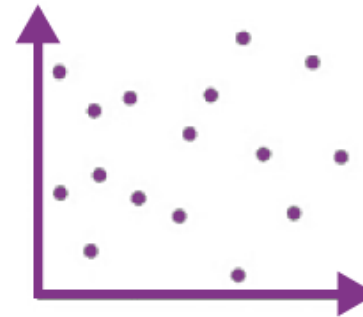
Strong negative correlation



Weak negative correlation



Moderate negative correlation

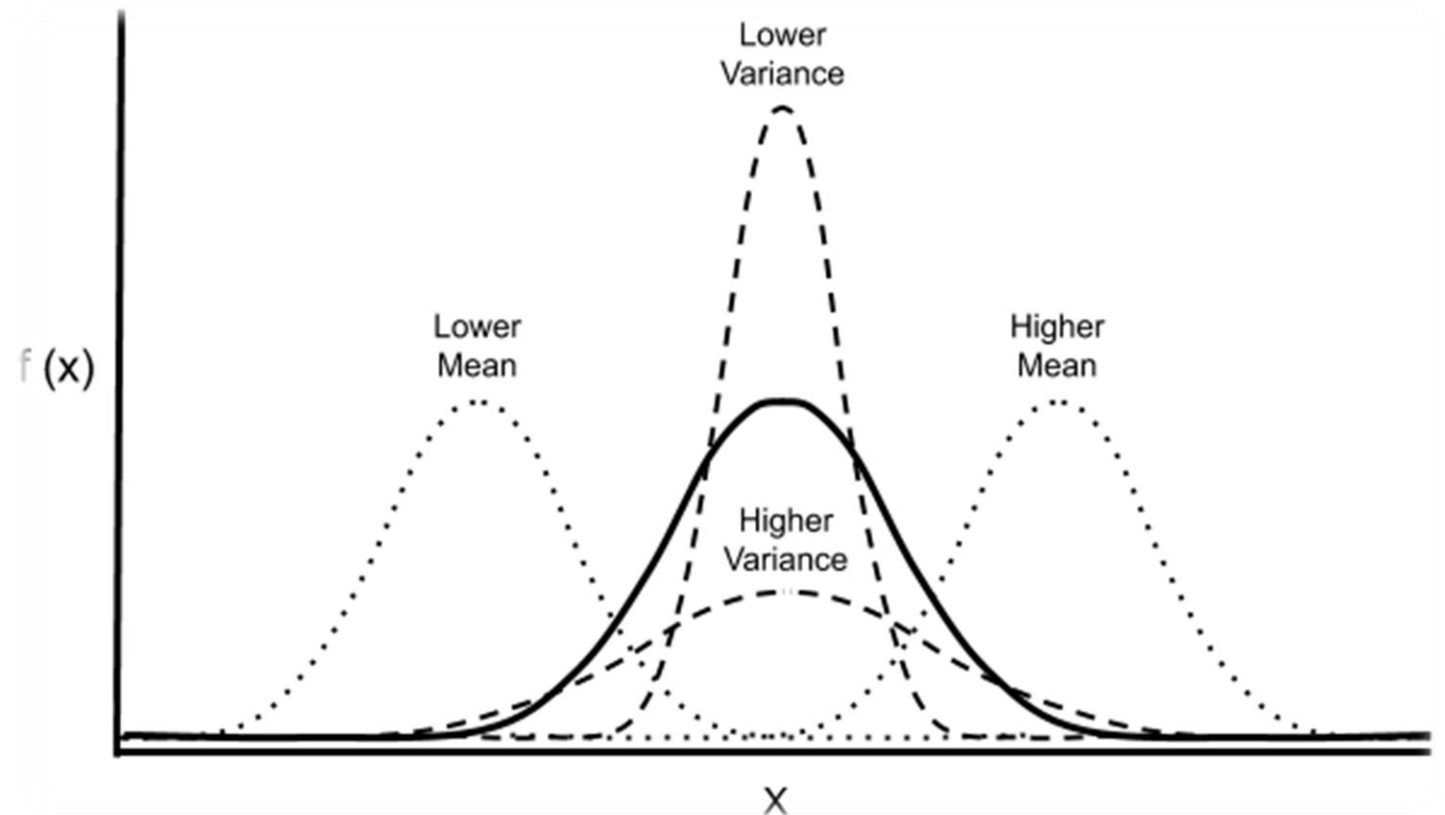


No correlation

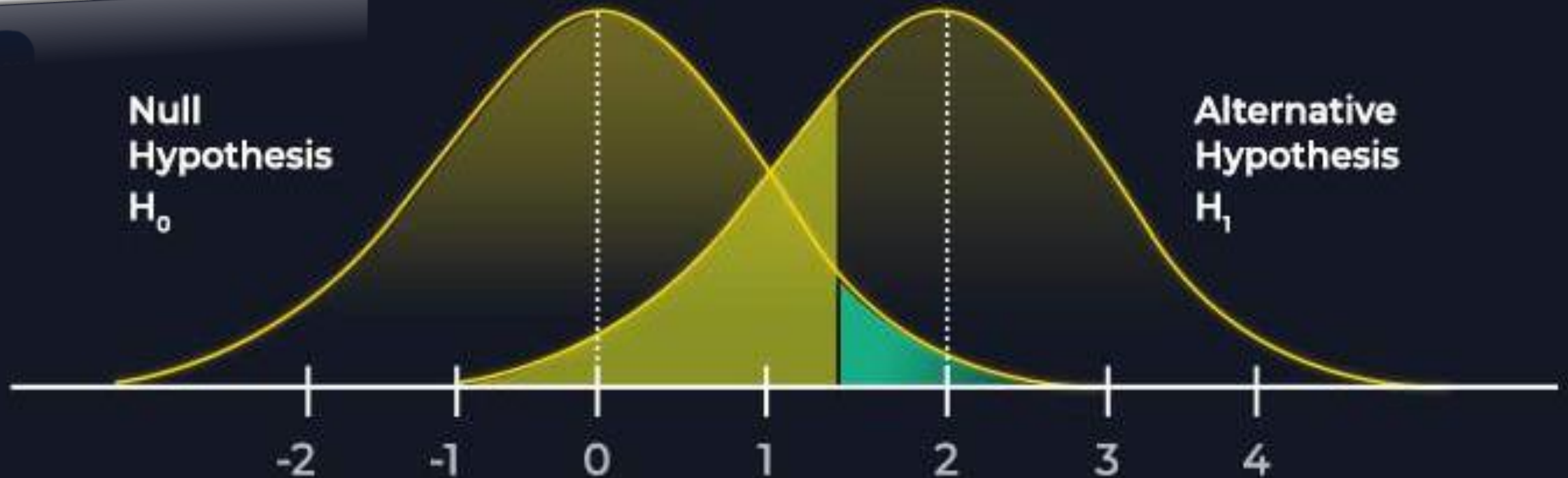


Κανονική Κατανομή

- Η κανονική κατανομή **εξαρτάται** από δύο παραμέτρους το μ (**mean**) και την **τυπική απόκλιση** (σ).
- Όταν $\mu=0$ και $\sigma=1$ η κατανομή λέγεται τυπική.
- Η κανονική κατανομή είναι **συμμετρική** ως προς τον **μέσο όρο**, δείχνοντας πως τα δεδομένα που βρίσκονται κοντά στο μέσο όρο εμφανίζονται πιο συχνά.



HYPOTHESIS TESTING



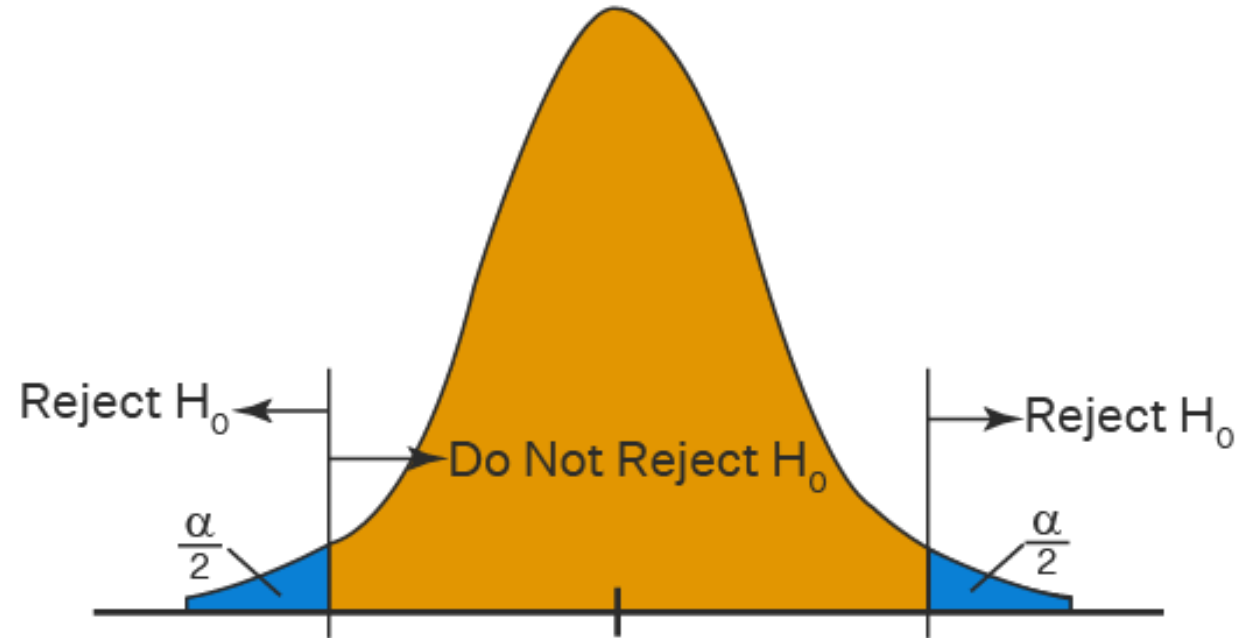
ΕΛΕΓΧΟΣ ΥΠΟΘΕΣΕΩΝ

ΧΡΗΣΙΜΟΠΟΙΩΝΤΑΣ
ΔΕΔΟΜΕΝΑ ΓΙΑ ΤΗΝ
ΕΞΑΓΩΓΗ
ΣΥΜΠΕΡΑΣΜΑΤΩΝ



Έλεγχος Υποθέσεων (1/2)

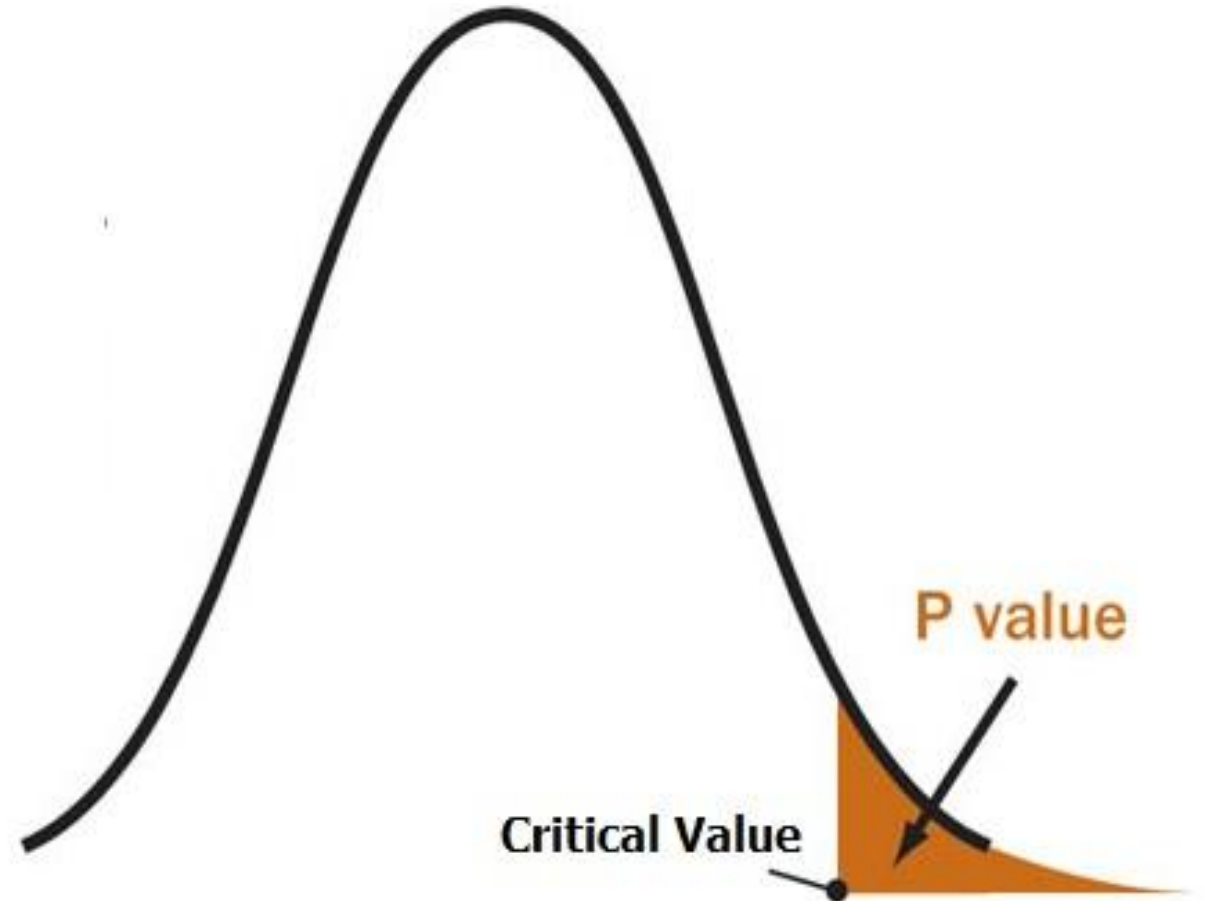
- Στο κλασσικό σενάριο έχουμε μια **μηδενική υπόθεση** την H_0 και μια **εναλλακτική** την H_1 . Χρησιμοποιούμε την στατιστική για να αποφασίσουμε εάν απορρίπτουμε την H_0 ως λανθασμένη ή όχι. Για παράδειγμα μπορούμε να προσδιορίσουμε εάν ο μέσος όρος είναι πιθανό να είναι αληθής.
- Αν ο δειγματικός μέσος είναι **κοντά** στον δηλωμένο μέσο του πληθυσμού, η **μηδενική υπόθεση δεν απορρίπτεται**.
- Αν ο δειγματικός μέσος είναι **μακριά** από τον δηλωμένο μέσο του πληθυσμού, η **μηδενική υπόθεση απορρίπτεται**.





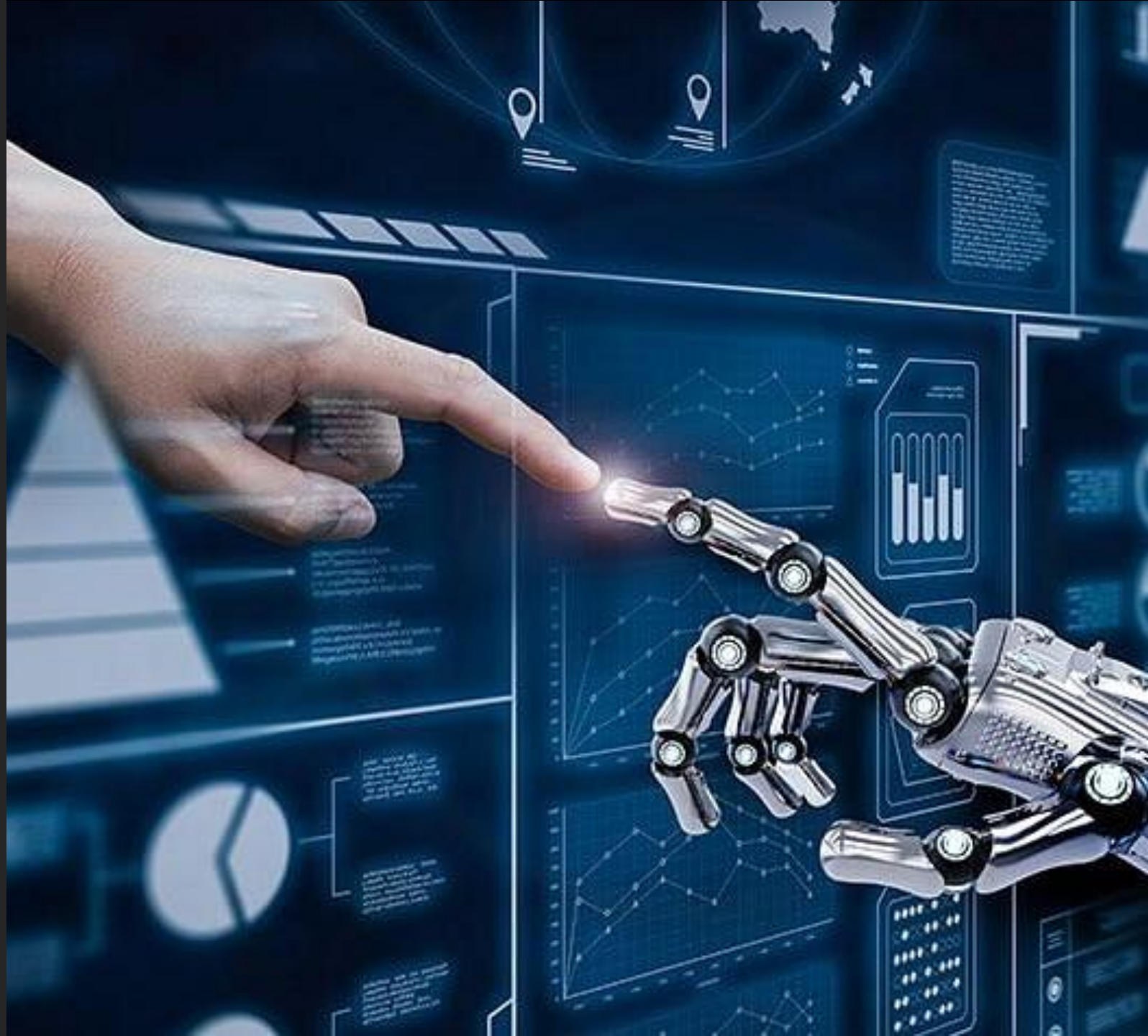
Έλεγχος Υποθέσεων (2/2)

- Ένας **εναλλακτικός τρόπος** για να δεχτούμε ή να απορρίψουμε υποθέσεις είναι η **τιμές p** . Αντί δηλαδή να επιλέγουμε όρια βασισμένοι σε κάποια «βάση» πιθανοτήτων, **υπολογίζουμε την πιθανότητα, υποθέτοντας πως η H_0 είναι αληθής**.
- Έστω δείκτης σπουδαιότητας $\alpha=5\%$, εάν $p\text{-value} > \alpha$ τότε δεν απορρίπτουμε την μηδενική υπόθεση.
- Όσο **μικρότερη** είναι η τιμή p , τόσο **ισχυρότερη** είναι η απόδειξη ότι πρέπει να **απορρίψουμε την μηδενική υπόθεση**.



ΜΗΧΑΝΙΚΗ ΜΑΘΗΣΗ

ΧΡΗΣΗ ΠΡΟΥΠΑΡΧΟΝΤΩΝ
ΔΕΔΟΜΕΝΩΝ ΓΙΑ ΤΗΝ
ΑΝΑΠΤΥΞΗ ΝΕΩΝ
ΜΟΝΤΕΛΩΝ ΠΡΟΒΛΕΨΗΣ





ΠΑΝΕΠΙΣΤΗΜΙΟ ΠΕΙΡΑΙΩΣ

UNIVERSITY OF PIRAEUS

Μηχανική Μάθηση (1/3)

- Η μηχανική μάθηση διερευνά τη μελέτη και την κατασκευή αλγορίθμων που μπορούν να μαθαίνουν από τα δεδομένα και να κάνουν προβλέψεις σχετικά με αυτά. Τέτοιοι αλγόριθμοι λειτουργούν κατασκευάζοντας μοντέλα από πειραματικά δεδομένα, προκειμένου να κάνουν προβλέψεις βασιζόμενες στα δεδομένα ή να εξάγουν αποφάσεις που εκφράζονται ως το αποτέλεσμα



ΠΑΝΕΠΙΣΤΗΜΙΟ ΠΕΙΡΑΙΩΣ

UNIVERSITY OF PIRAEUS

Μηχανική Μάθηση (2/3)

- Υπερπροσαρμογή

Δημιουργία ενός μοντέλου που αποδίδει καλά στα δεδομένα μας έχοντας μέτρια απόδοση σε οποιαδήποτε νέα.

- Υποπροσαρμογή

Δημιουργία ενός μοντέλου που δεν αποδίδει καλά, με αποτέλεσμα να μην έχει νόημα η ανάλυση του.

- Κάποιες φορές **πολύπλοκα μοντέλα** οδηγούν σε **υπερπροσαρμογή**, επομένως ο απλούστερος τρόπος για να μην είναι το μοντέλο μας πολύπλοκο είναι να **διαχωρίσουμε τα δεδομένα μας**.



ΠΑΝΕΠΙΣΤΗΜΙΟ ΠΕΙΡΑΙΩΣ

UNIVERSITY OF PIRAEUS

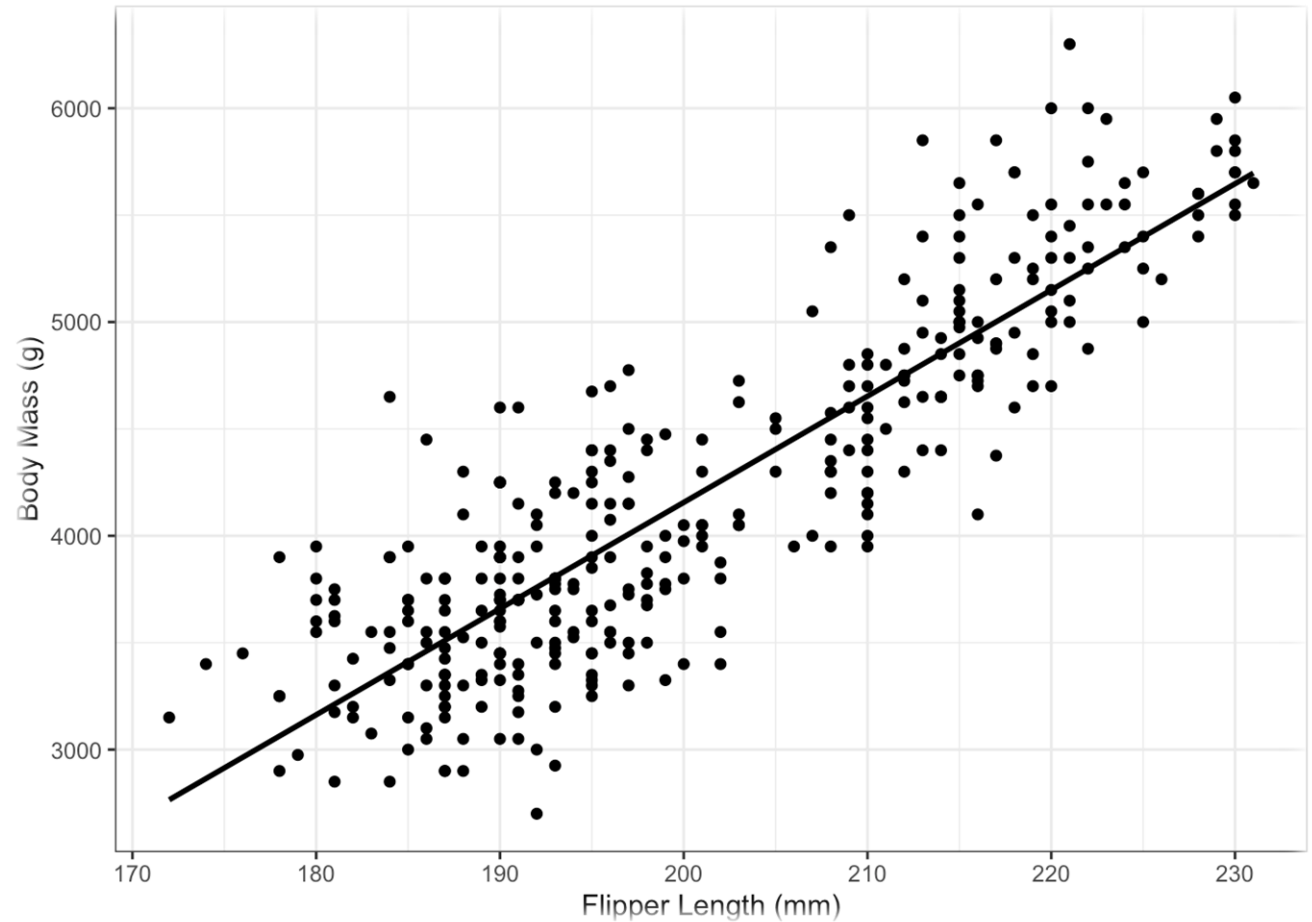
Μηχανική Μάθηση (3/3)

- Ένας άλλος τρόπος για την αντιμετώπιση της υπερπροσαρμογής είναι ο συμβιβασμός της Μεροληψίας και της διακύμανσης.
- Υψηλή μεροληψία και χαμηλή διακύμανση αντιστοιχούν σε υποπροσαρμογή.
- Χαμηλή μεροληψία και υψηλή διακύμανση αντιστοιχούν σε υπερπροσαρμογή.
- Η υψηλή μεροληψία αντιμετωπίζεται προσθέτοντας περισσότερα χαρακτηριστικά.
- Η υψηλή διακύμανση αντιμετωπίζεται αφαιρώντας χαρακτηριστικά.

Όσο περισσότερα δείγματα έχουμε τόσο δυσκολότερο είναι να έχουμε υπερπροσαρμογή, χωρίς αυτό να σημαίνει απαραίτητα πως θα βοηθήσουν στην μεροληψία του δείγματος.

ΑΠΛΗ ΓΡΑΜΜΙΚΗ ΠΑΛΙΝΔΡΟΜΗΣΗ

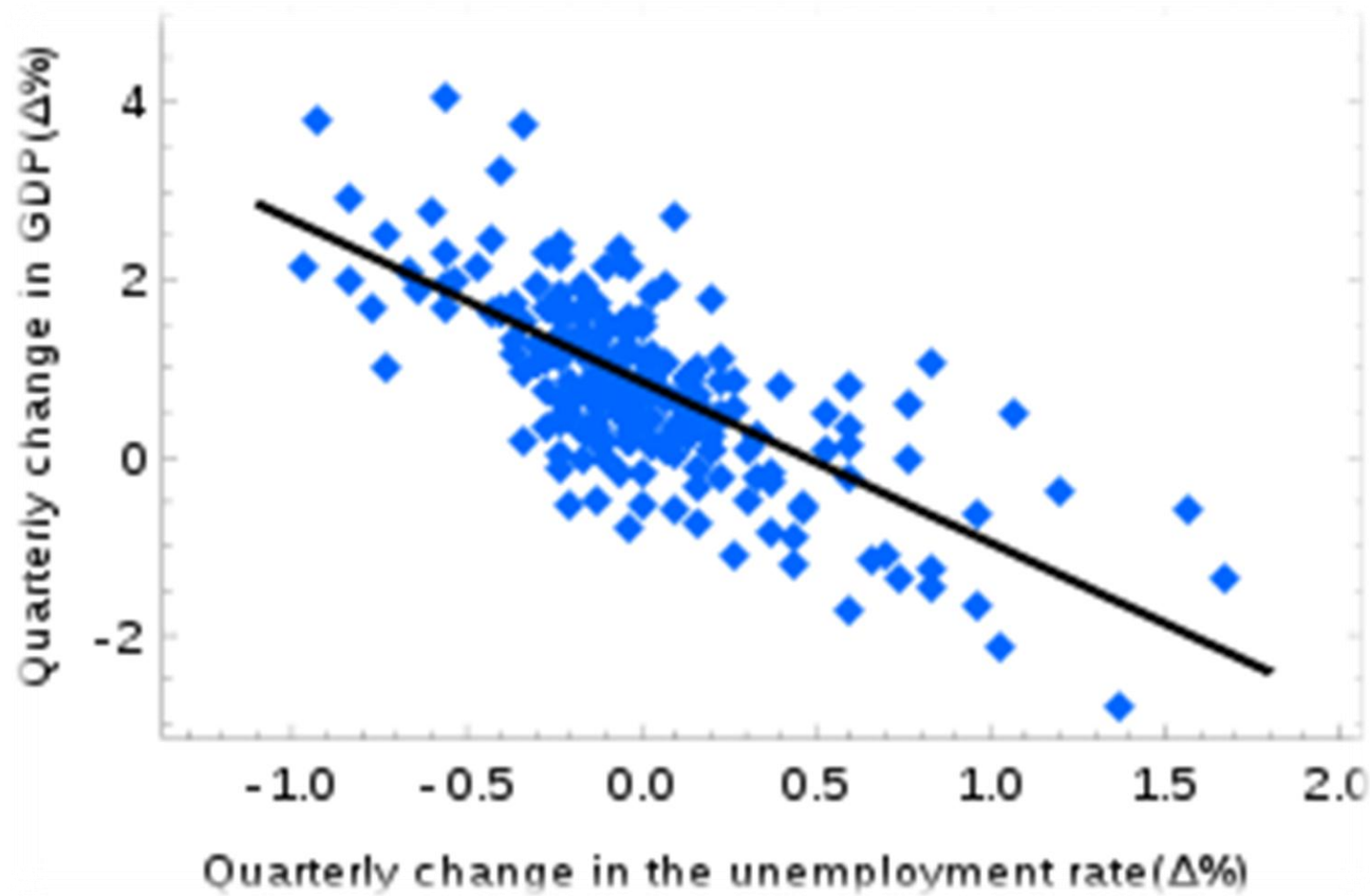
$$Y = A + BX + E$$





Απλή Γραμμική Παλινδρόμηση(1/2)

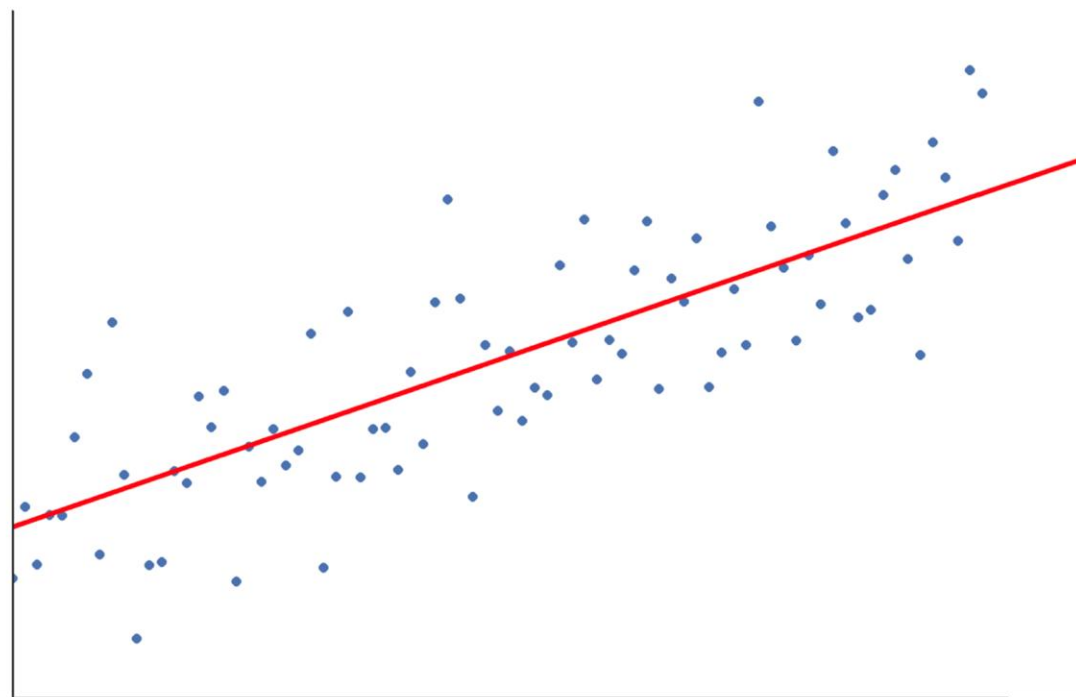
- Στην παραπάνω εξίσωση ο όρος **E** είναι ο όρος **σφάλματος**. Για κάθε ζεύγος A και B χρειάζεται να υπολογίσουμε το σφάλμα. Βέβαια αυτό που θα θέλαμε να γνωρίζαμε είναι το συνολικό σφάλμα σε ολόκληρο το σύνολο δεδομένων.
- Στην ουσία η απλή γραμμική παλινδρόμηση είναι η **μοντελοποίηση μιας εξαρτημένης μεταβλητής και μιας ανεξάρτητης μεταβλητής**.
- Στο παράδειγμα της πολλαπλής παλινδρόμησης όμως έχουμε ένα μεγαλύτερο πλήθος ανεξάρτητων μεταβλητών.





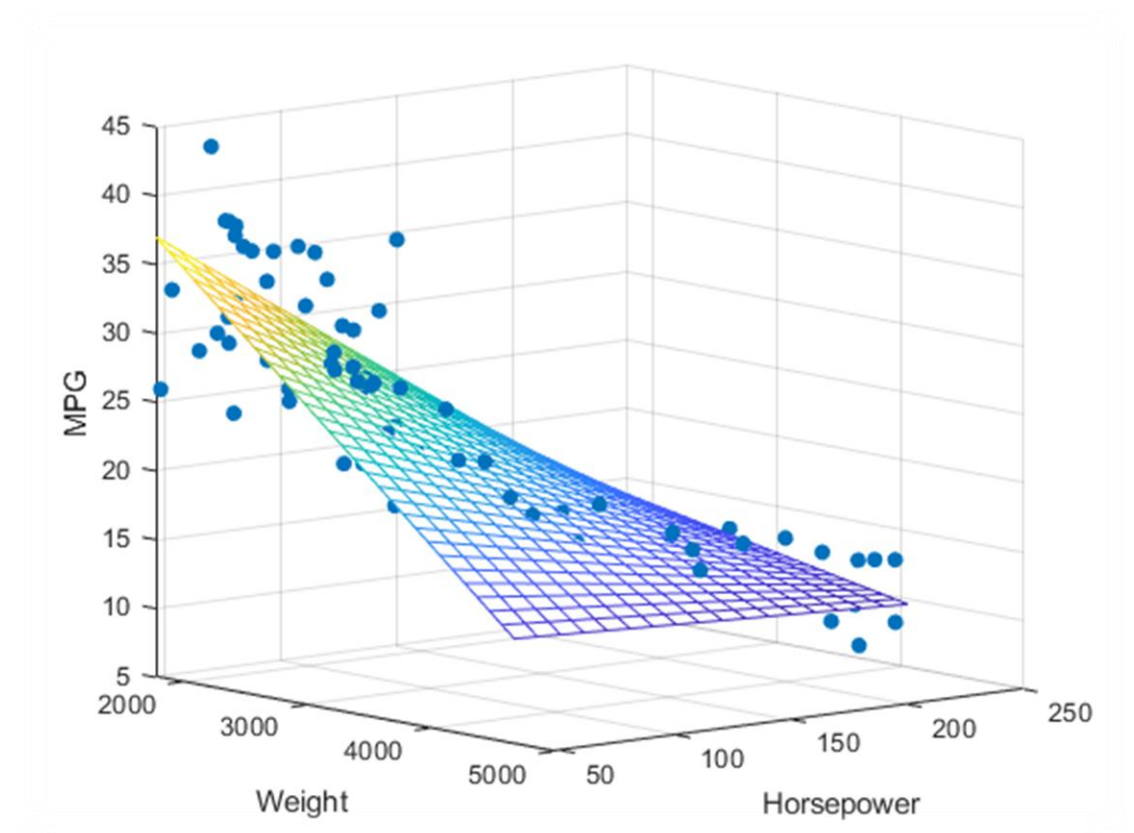
Απλή Γραμμική Παλινδρόμηση (2/2)

- Η Μέθοδος Ελαχίστων Τετραγώνων μας δίνει τη δυνατότητα να επιλέξουμε συγκεκριμένες τιμές για A και B ελαχιστοποιώντας το σφάλμα.
- Βέβαια δεν αρκεί μόνο αυτό για να βεβαιωθούμε πως έχουμε επιλέξει το σωστό για το μοντέλο μας.
- Ο συντελεστής R^2 είναι ένας συντελεστής που μετρά το ποσοστό της συνολικής απόκλισης της εξαρτημένης μεταβλητής.



ΠΟΛΛΑΠΛΗ ΠΑΛΙΝΔΡΟΜΗΣΗ

$$Y = A + B_1X_1 + B_2X_2 + B_3X_3 + E$$

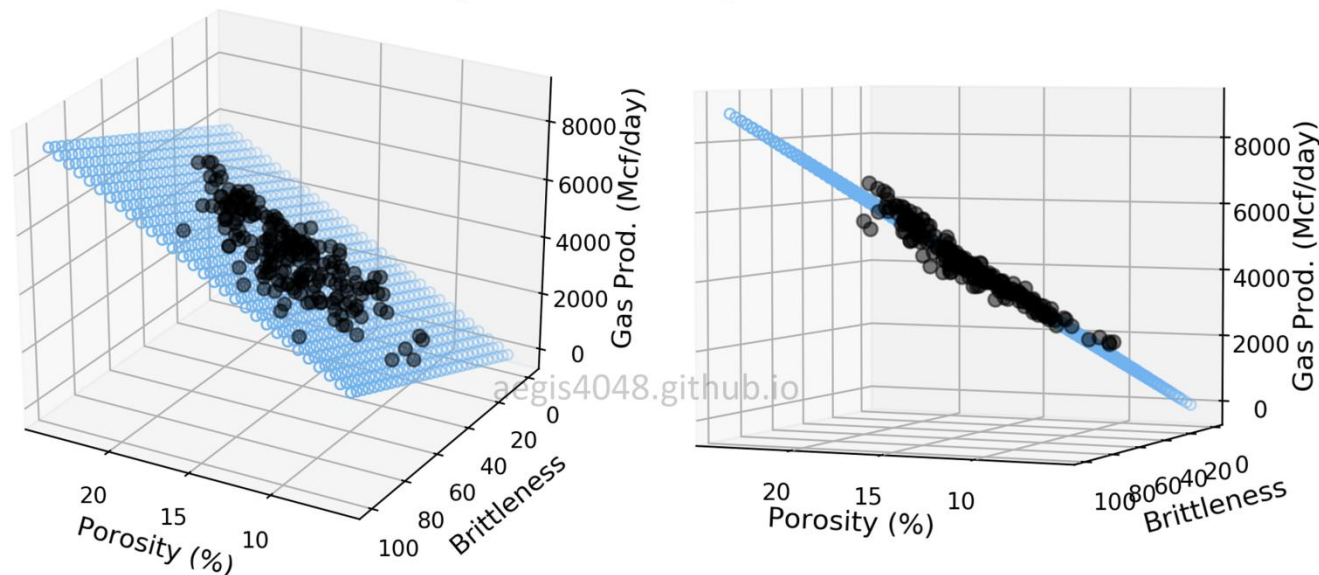




Πολλαπλή Παλινδρόμηση(1/2)

- Ένα γραμμικό μοντέλο με περισσότερες ανεξάρτητες μεταβλητές. Όπως κάναμε και στο απλό γραμμικό μοντέλο θα επιλέξουμε “beta” ώστε να ελαχιστοποιήσουμε το άθροισμα των τετραγωνικών σφαλμάτων με την διαφορά πως σε αυτήν την περίπτωση το beta θα δέχεται ένα διάνυσμα αυθαίρετου μήκους.
- Γενικά η προσθήκη νέων μεταβλητών σε μια παλινδρόμηση θα αυξάνουν συνεχώς το $R \text{ sqrt.}$ Για αυτό το λόγο θα χρειαστεί να εξετάζουμε και τα τυπικά σφάλματα.

3D multiple linear regression model

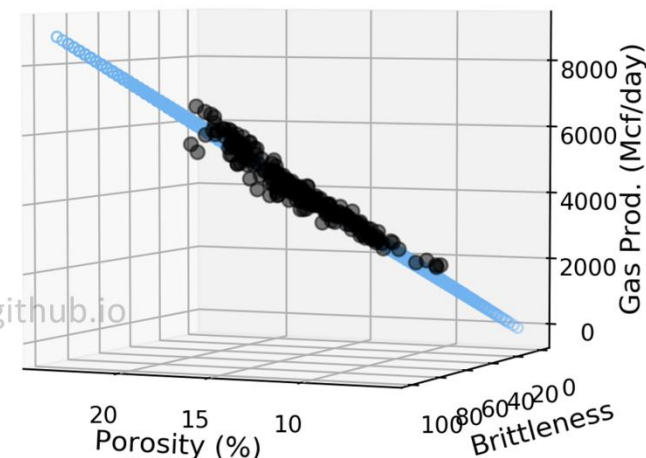
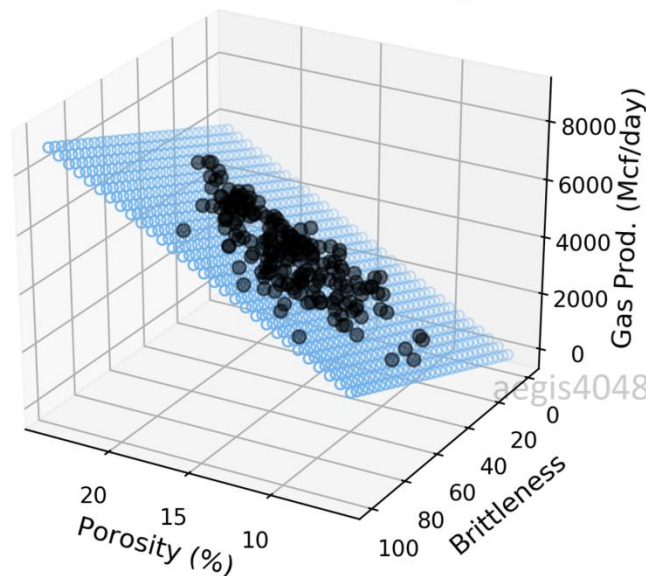




Πολλαπλή Παλινδρόμηση(2/2)

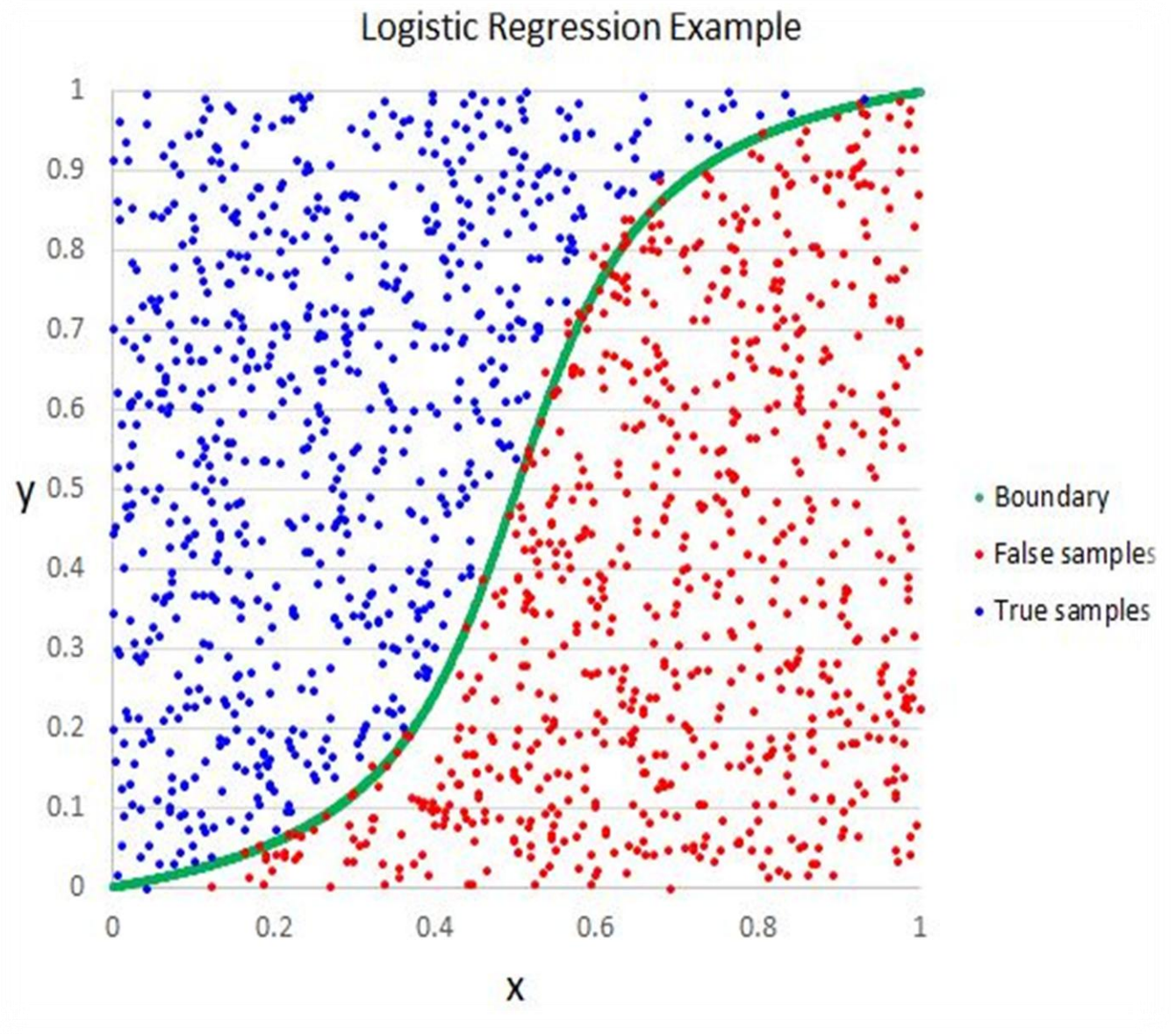
- Όσο **περισσότερες μεταβλητές** χρησιμοποιούμε τόσο πιο πιθανό είναι να «**υπερπροσαρμόσουμε**» το μοντέλο μας.
- Δεύτερον όσο περισσότερους **μη μηδενικούς** συντελεστές έχουμε στο μοντέλο μας τόσο πιο **δύσκολο** είναι να βγάλουμε νόημα μέσα από αυτούς.

3D multiple linear regression model



ΛΟΓΙΣΤΙΚΗ ΠΑΛΙΝΔΡΟΜΗΣΗ

LOGISTIC REGRESSION



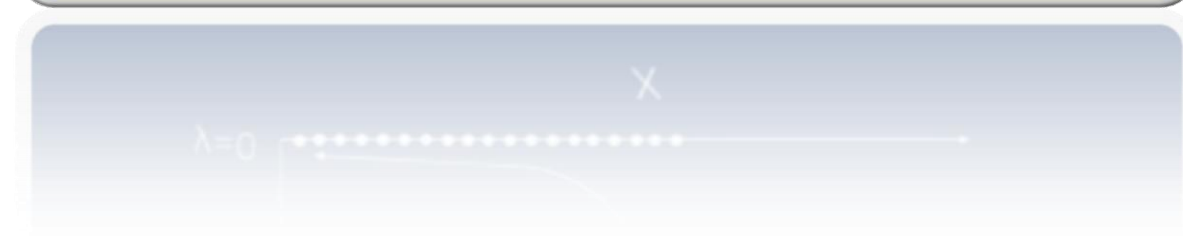
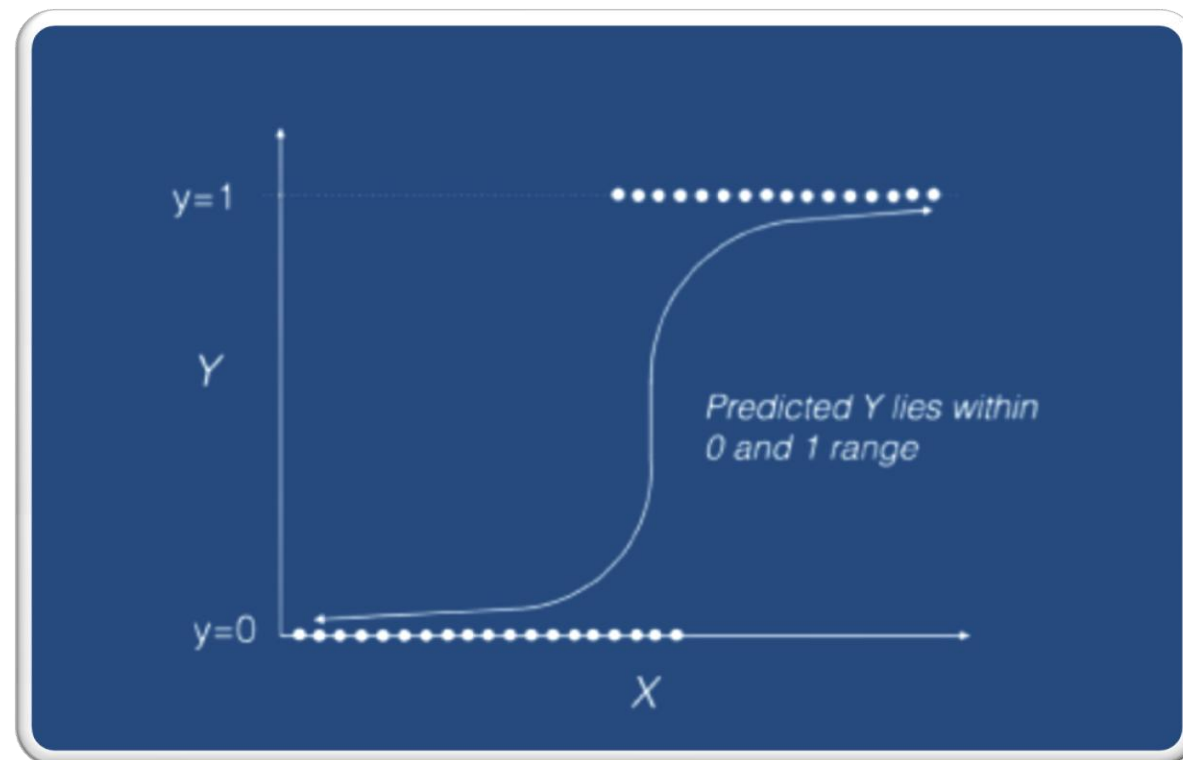


Λογιστική Παλινδρόμηση(1/2)

- Η λογιστική παλινδρόμηση είναι μια S-shaped συνάρτηση. Όσο η είσοδος της παίρνει υψηλότερες θετικές τιμές τόσο πιο πολύ πλησιάζει στο 1.

- Η συνάρτηση είναι της μορφής :

$Y_i = f(X_i\beta) + \epsilon_i$, όπου f λογιστική συνάρτηση.



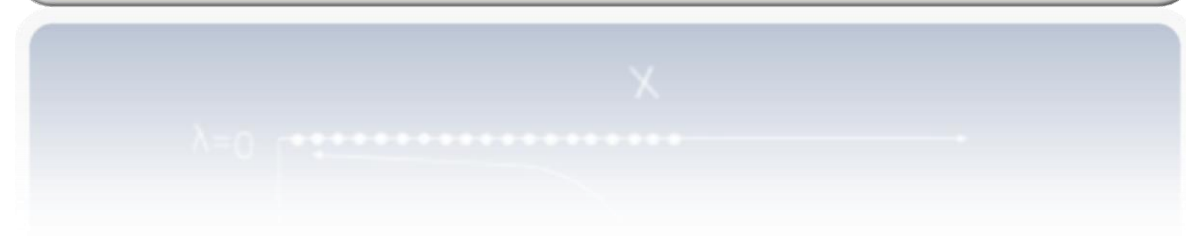
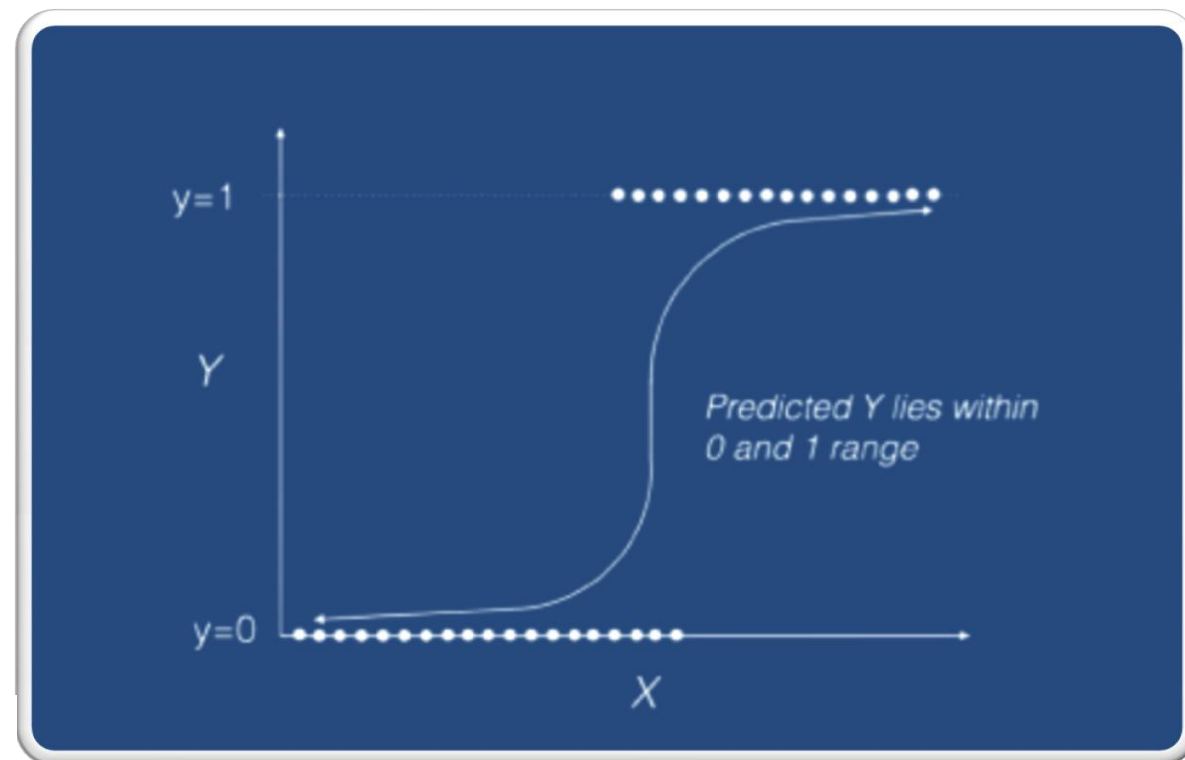


Λογιστική Παλινδρόμηση(2/2)

- Μια **βασική διαφορά** είναι πως σε αντίθεση με την απλή γραμμική παλινδρόμηση που χρησιμοποιεί την ελαχιστοποίηση τετραγώνων, η **λογιστική παλινδρόμηση εφαρμόζει την τεχνική της μέγιστης πιθανοφάνειας**.

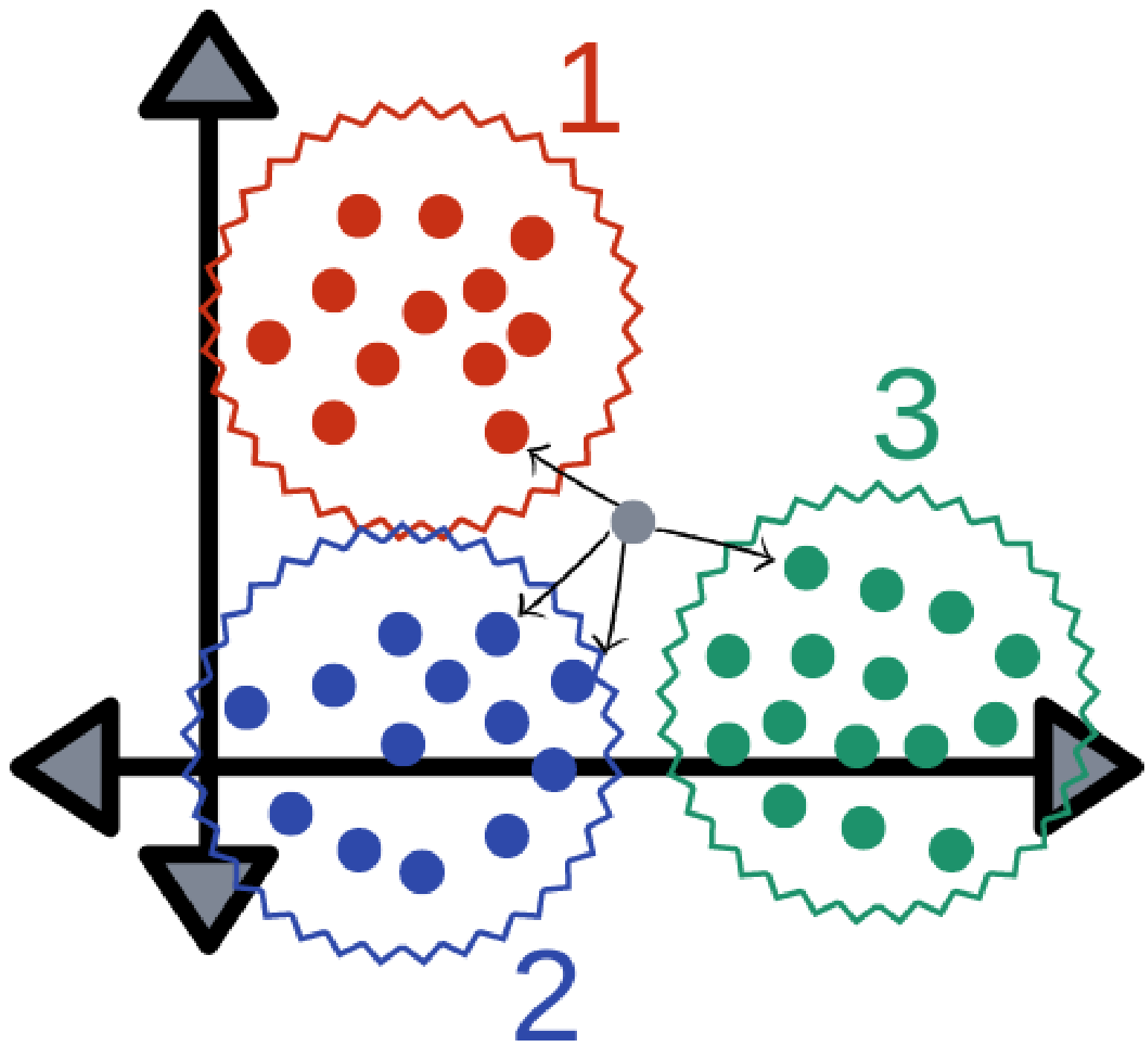
Αν X_1, \dots, X_n ανεξάρτητες, τότε $L(\underline{\vartheta}|\underline{x}) = \prod_{i=1}^n f(x_i; \underline{\vartheta})$.

Ερμηνεία: Για δεδομένο δείγμα παρατηρήσεων \underline{x} , η συνάρτηση πιθανοφάνειας $L(\underline{\vartheta}|\underline{x})$ εκφράζει το πόσο πιθανό είναι να έχουμε παρατηρήσει τα δεδομένα \underline{x} που παρατηρήσαμε για τις διάφορες τιμές της άγνωστης παραμέτρου $\underline{\vartheta}$.



Κ-ΠΛΗΣΙΕΣΤΕΡΟΙ ΓΕΙΤΟΝΕΣ

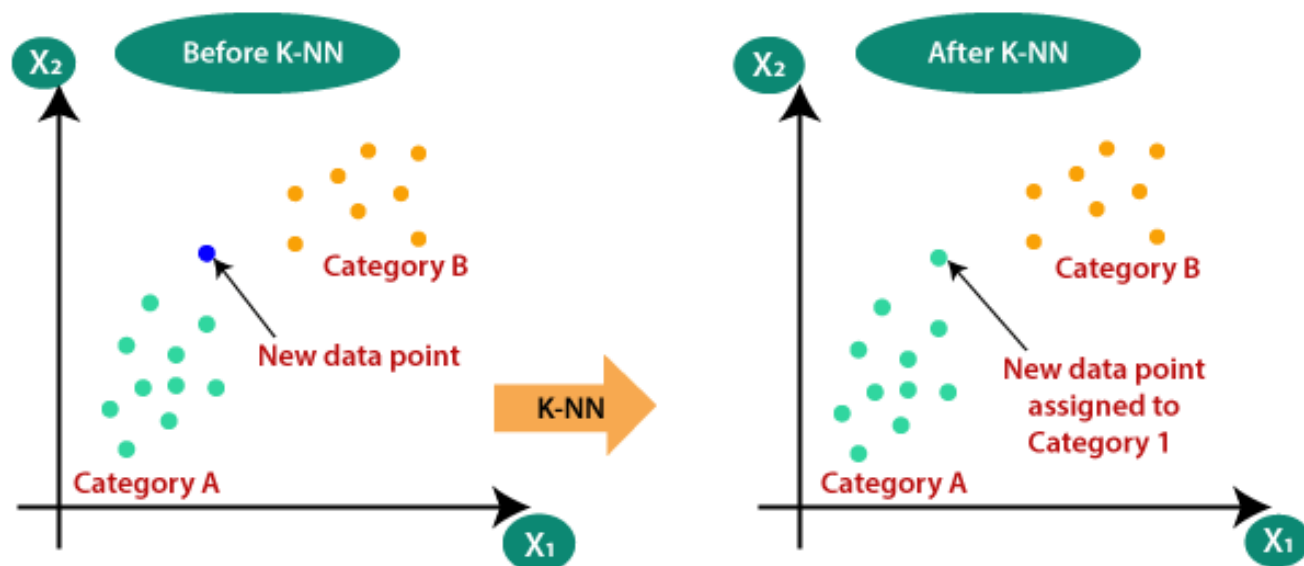
K-NEAREST
NEIGHBORS
(KNN)





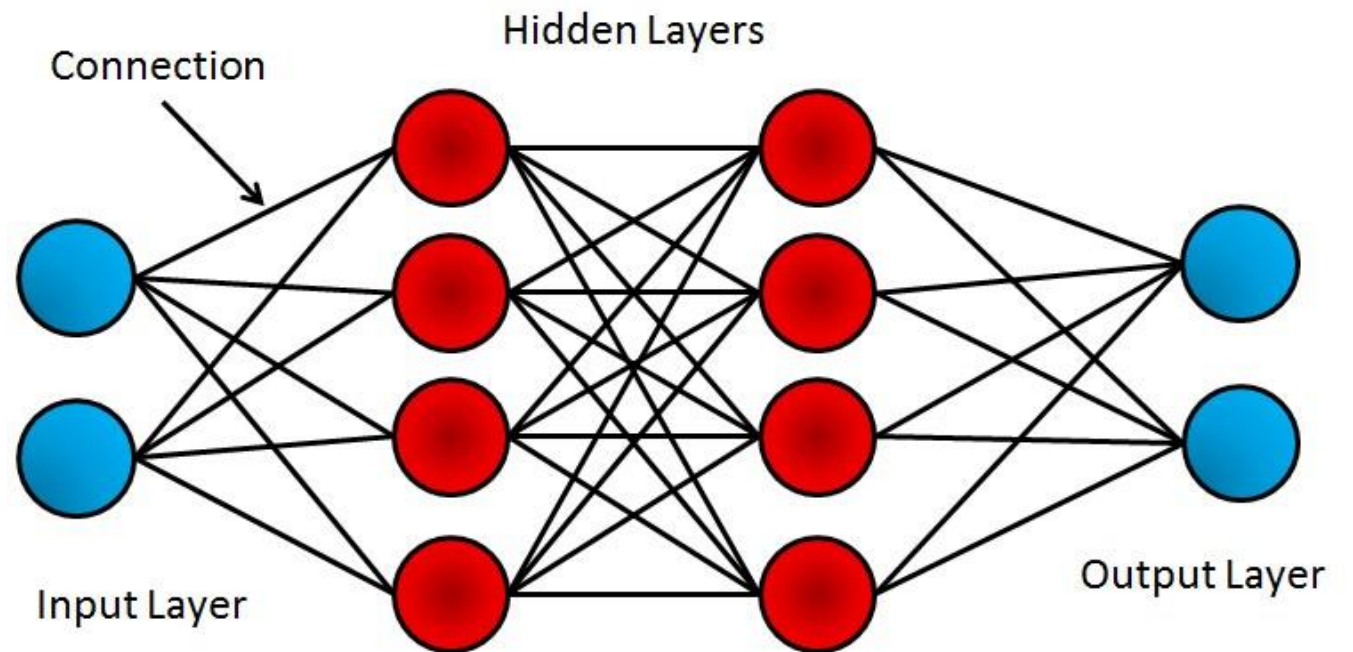
k-nearest neighbors (KNN)

- Λειτουργεί βρίσκοντας **αποστάσεις** μεταξύ ενός ερωτήματος μέσω όλων των παραδειγμάτων στα δεδομένα, επιλέγοντας αριθμούς K πιο κοντά στο ερώτημα. Έπειτα τα συλλέγει σε μία λίστα βρίσκοντας την πιο **συχνή παρατήρηση** ή υπολογίζοντας τον **μέσο όρο**.
- Ένας σύνηθες και εύκολος **τρόπος** είναι να λάβουμε υπόψη το **άθροισμα των τετραγωνικών σφαλμάτων** κάθε σημείου σαν συνάρτηση του K .
- Βέβαια το μοντέλο αυτό συναντά προβλήματα σε υψηλότερες διαστάσεις. **Καθώς το πλήθος των διαστάσεων αυξάνεται, αυξάνεται και η μέση απόσταση τους.**



ΝΕΥΡΩΝΙΚΑ ΔΙΚΤΥΑ

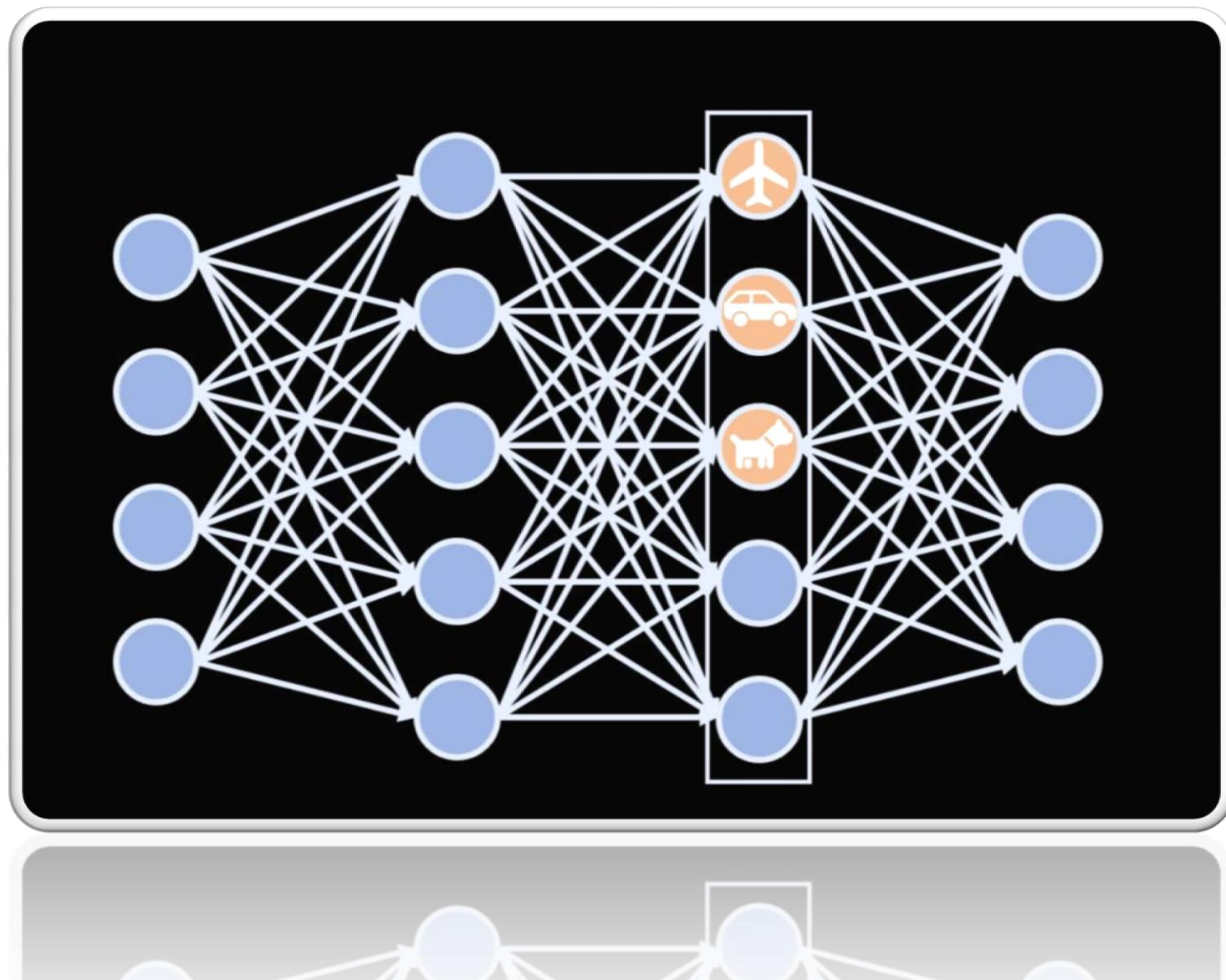
ΕΞΕΤΑΖΟΝΤΑΣ ΤΙΣ ΕΞΟΔΟΥΣ
ΤΩΝ ΑΛΛΩΝ ΝΕΥΡΩΝΩΝ ΠΟΥ
ΤΟΝ ΤΡΟΦΟΔΟΤΟΥΝ





Νευρωνικά Δίκτυα

- Τα νευρωνικά δίκτυα αποτελούνται από κόμβους που συνδέονται με άλλους κόμβους.
- **Κάθε σύνδεση** μεταξύ των κόμβων αποτελείται από αριθμούς που **εκφράζουν την κλήση μιας ευθείας** η οποία εκφράζει μια απλή παλινδρόμηση πάνω στα δεδομένα μας.
- Οι κόμβοι που βρίσκονται ανάμεσα στους κόμβους εισόδου και εξόδου ονομάζονται **Hidden layers**. Και είναι ιδιαίτερα βασικά διότι στην αρχή κατασκευής ενός νευρωνικού δικτύου είναι απαραίτητος ο καθορισμός της ποσότητας των Hidden layers.



The diagram illustrates a deep neural network architecture. It consists of five layers of nodes: an input layer with 5 blue nodes, two hidden layers each with 4 blue nodes, an output layer with 3 orange nodes (each containing a different icon: a plane, a car, and a dog), and a final output layer with 5 blue nodes. The nodes are fully connected to the nodes in the adjacent layers by white lines. A red rectangular box highlights the output layer (the 3 orange nodes), and a green rectangular box highlights the input layer (the 5 blue nodes).

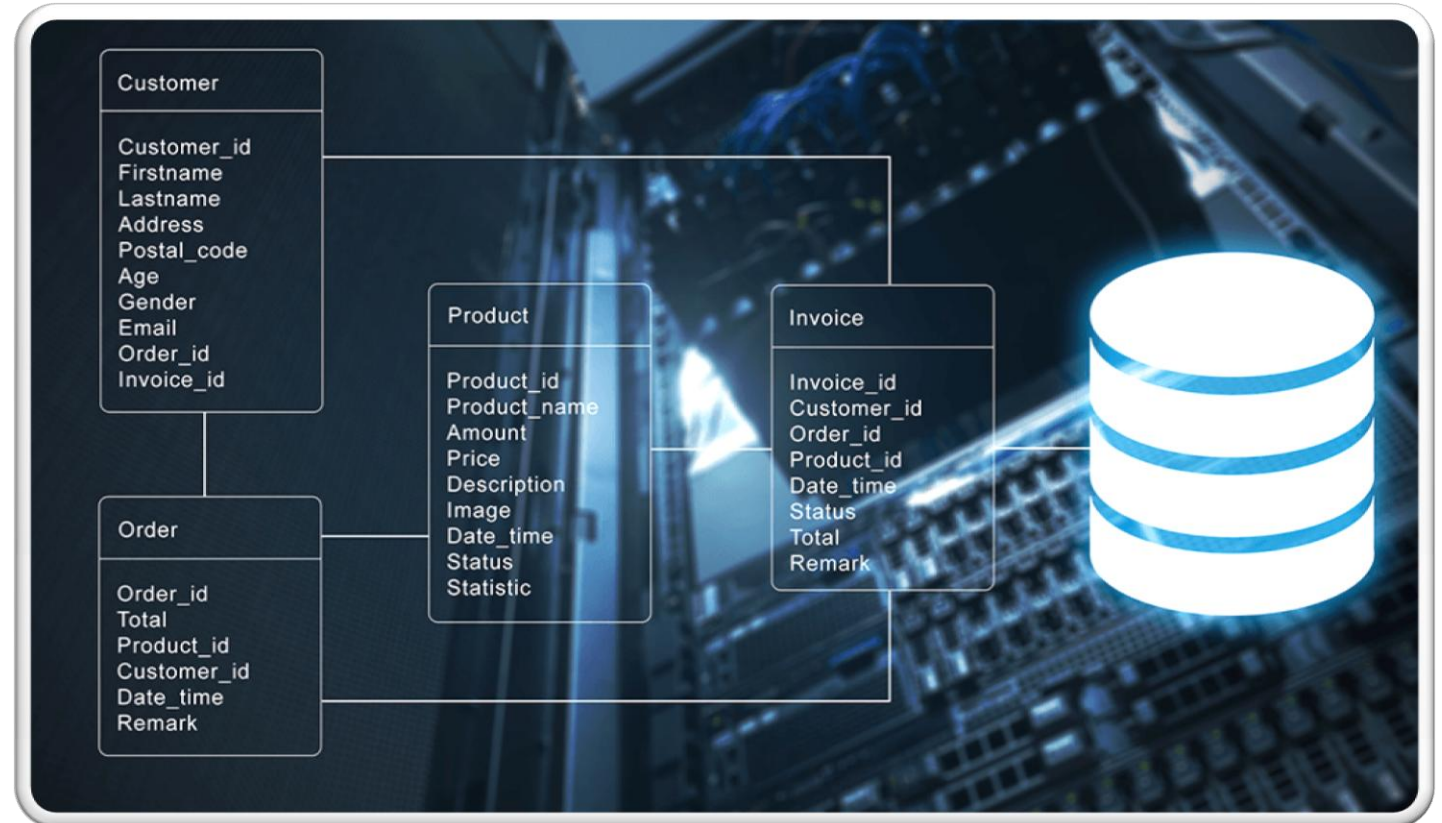


ΒΑΘΙΑ ΜΑΘΗΣΗ DEEP LEARNING

ΤΑ ΒΑΘΥΤΕΡΑ
ΝΕΥΡΩΝΙΚΑ ΔΙΚΤΥΑ

Βάσεις Δεδομένων και SQL - Structured Query Language

ΣΥΣΤΗΜΑΤΑ ΣΧΕΔΙΑΣΜΕΝΑ ΓΙΑ
ΤΗΝ ΑΠΟΘΗΚΕΥΣΗ ΚΑΙ
ΑΝΑΖΗΤΗΣΗ ΤΩΝ ΔΕΔΟΜΕΝΩΝ



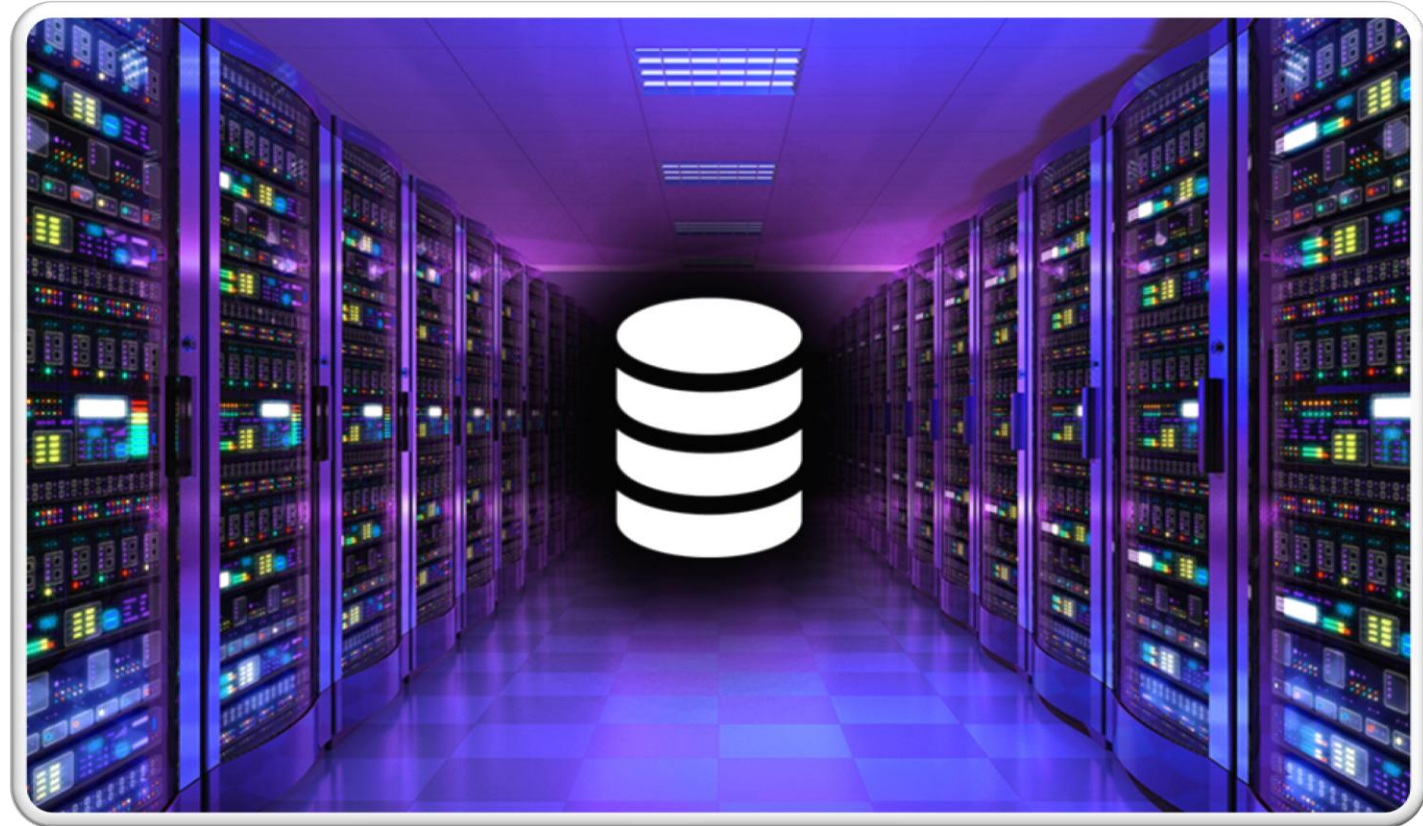


ΠΑΝΕΠΙΣΤΗΜΙΟ ΠΕΙΡΑΙΩΣ

UNIVERSITY OF PIRAEUS

Βάσεις Δεδομένων και SQL (1/4)

- Ο κύριος όγκος των βάσεων δεδομένων είναι relational (σχεσιακές), όπως η PostgreSQL, MySQL και ο SQL server, οι οποίες αποθηκεύουν τα δεδομένα σε πίνακες.





Βάσεις Δεδομένων και SQL (2/4)

Employee table

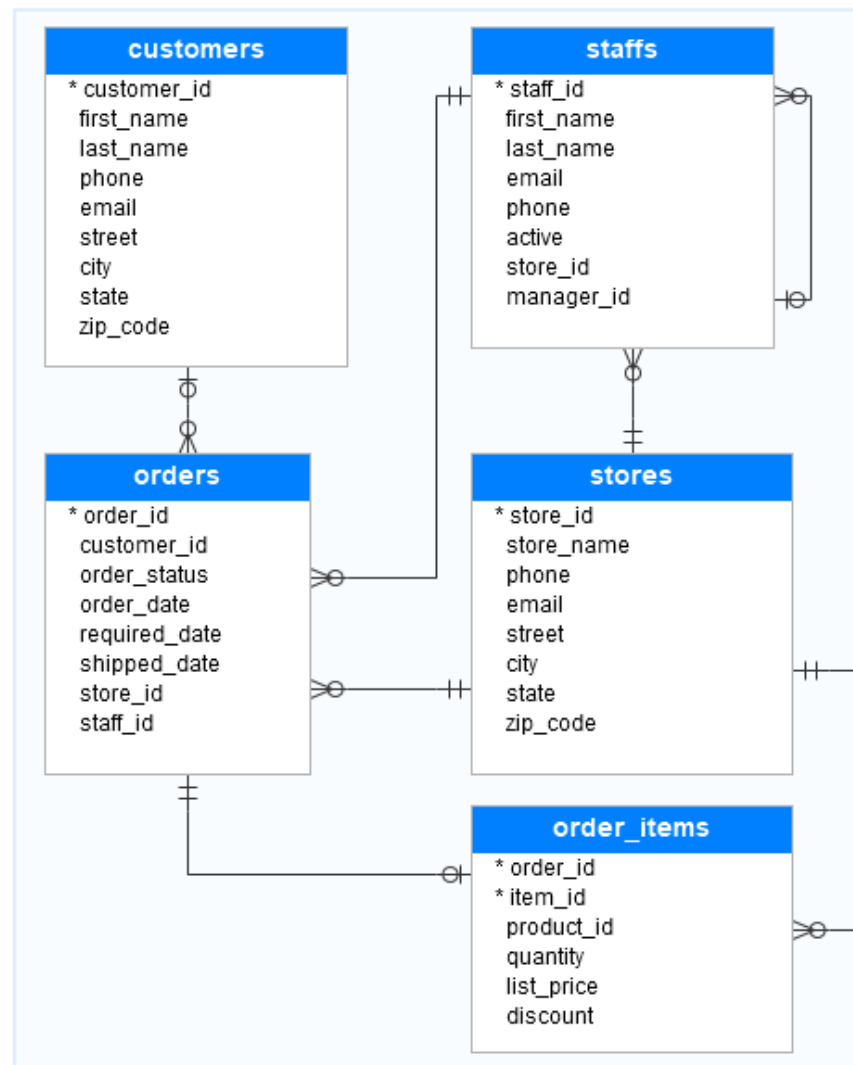
Empno (PK)	Ename	Job	Deptno (FK)
101	A	Salesman	10
102	B	Manager	10
103	c	Manager	20

Department table

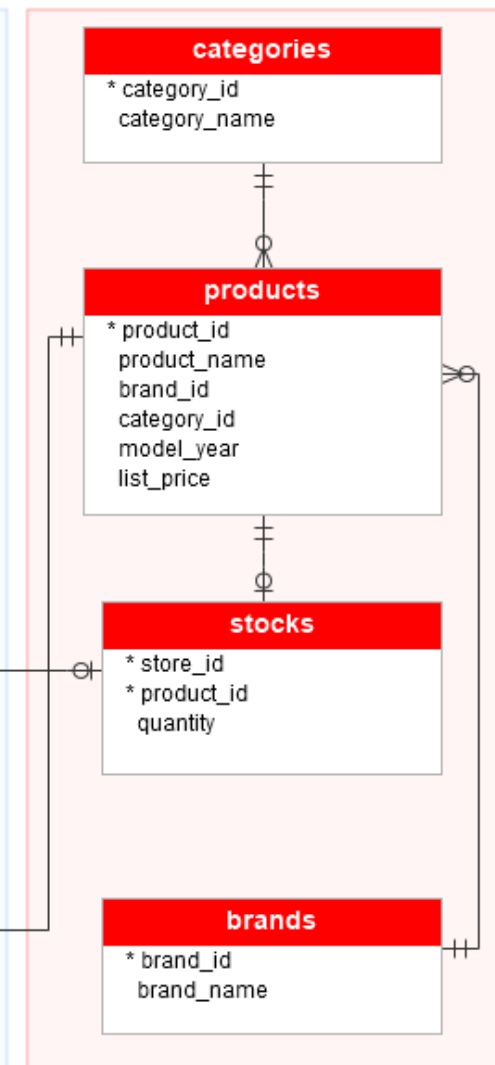
Deptno (PK)	dname	loc
10	Sales	Chicago
20	Sales	Chicago
30	Finance	New York



Sales



Production





Βάσεις Δεδομένων και SQL (3/4)

Basic SQL query structure

```
SELECT Attributes  
FROM relations  
WHERE condition
```

For example:

```
SELECT sid,sname  
FROM students  
WHERE sid=1122
```




Βάσεις Δεδομένων και SQL (4/4)

```
-- get employees who joined company in 2000
```

SELECT clause	{	SELECT		Comment
		first_name		
FROM clause	{	FROM		
		employees		
WHERE clause	{	WHERE		
		YEAR(hire_date) = 2000		
			Predicate	