



Πανεπιστήμιο Κύπρου - Τεχνητή Νοημοσύνη

MAI612 - ΜΗΧΑΝΙΚΗ ΜΑΘΗΣΗ

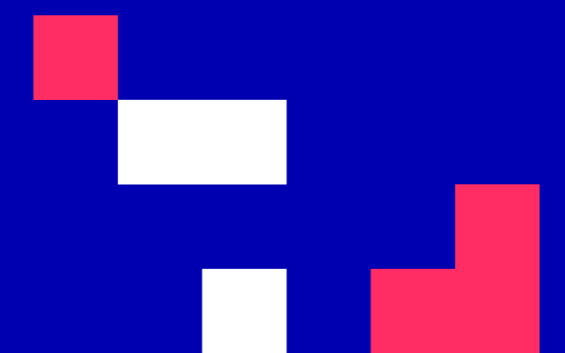
Διάλεξη 2: Προετοιμασία δεδομένων

Βασίλης Βασιλειάδης, PhD

Χειμερινό Εξάμηνο 2022/23



CYENS
CENTRE OF EXCELLENCE





Διάλεξη 2: Προετοιμασία δεδομένων

Μαθησιακά αποτελέσματα

Θα καταλάβετε:

1. τη διαδικασία συλλογής δεδομένων και προετοιμασίας συνόλων δεδομένων για την ανάπτυξη λύσεων μηχανικής μάθησης
2. τη σημασία της προετοιμασίας των δεδομένων
3. τα στάδια της προεπεξεργασίας δεδομένων και του μετασχηματισμού δεδομένων
4. τη σημασία της απεικόνισης των δεδομένων και της διερευνητικής ανάλυσης δεδομένων
5. γιατί χωρίζουμε ένα σύνολο δεδομένων σε σύνολο εκπαίδευσης και σύνολο δοκιμών





Ταξινομητής εικόνων φρούτων με χρήση Μηχανικής Μάθησης (MM)

Είστε μηχανικός MM στην εταιρεία XYZ και σας ανατέθηκε να δημιουργήσετε μια εφαρμογή για κινητά που ταξινομεί διαφορετικές εικόνες φρούτων, στις αντίστοιχες κατηγορίες τους (μήλο, πορτοκαλί κ.λπ.).

Θέλετε να χρησιμοποιήσετε τη MM

Ποιο είναι το πρώτο βήμα της διαδικασίας;

➤ Συλλογή δεδομένων

Πώς;

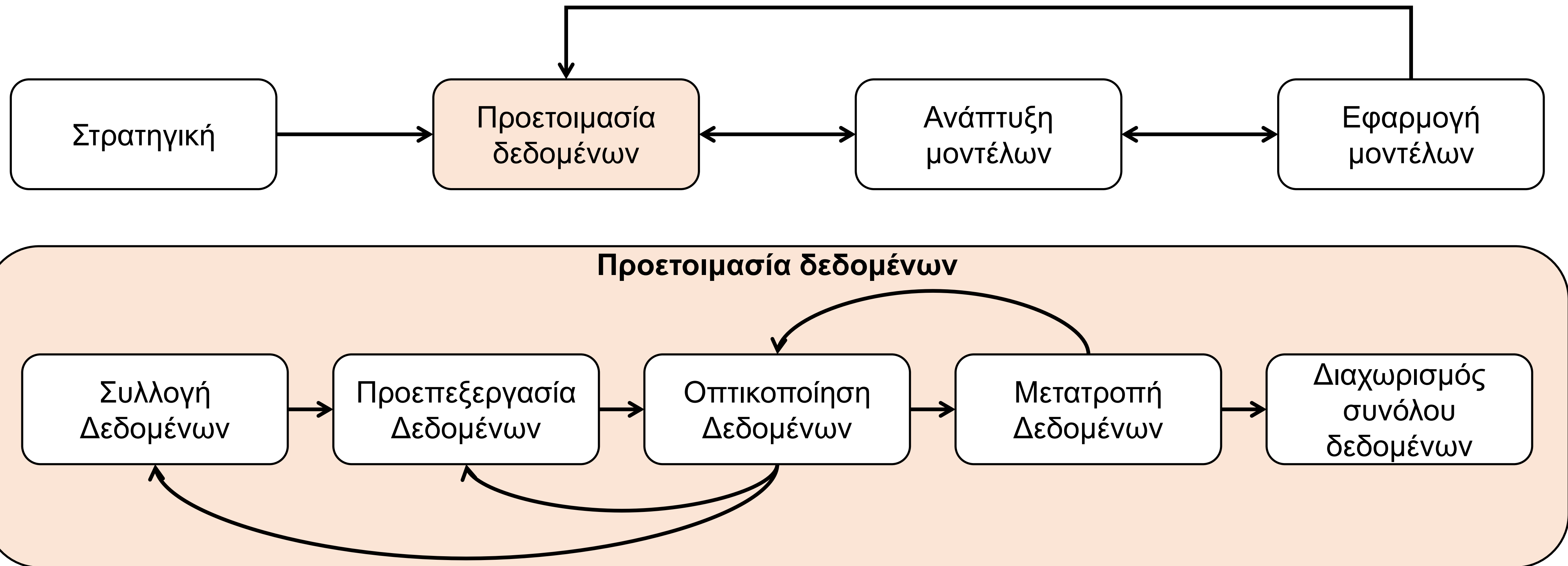
➤ Στο διαδίκτυο

➤ Τραβήξτε φωτογραφίες και συνδέστε τις με ετικέτες



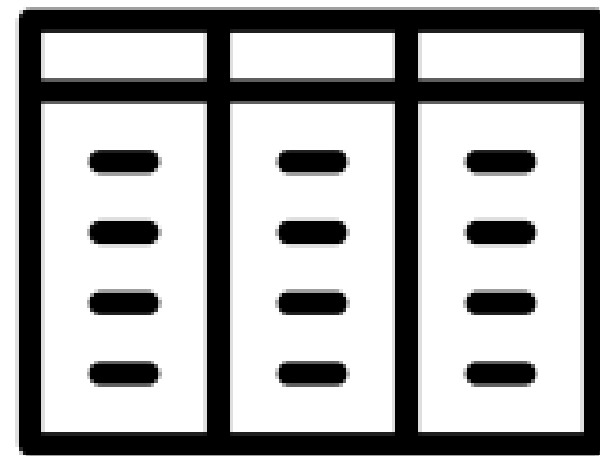


Προετοιμασία δεδομένων





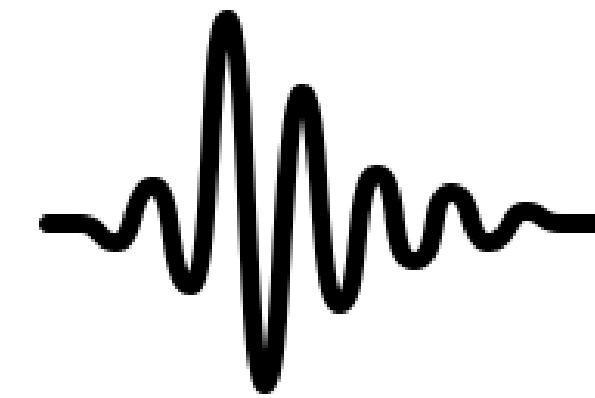
Δεδομένα



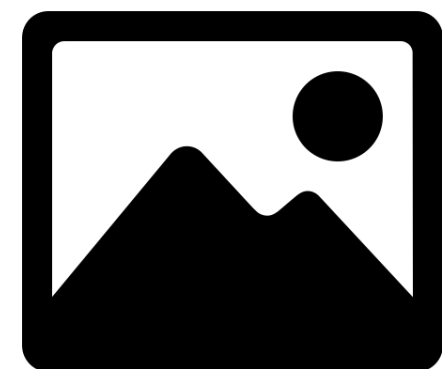
πίνακας



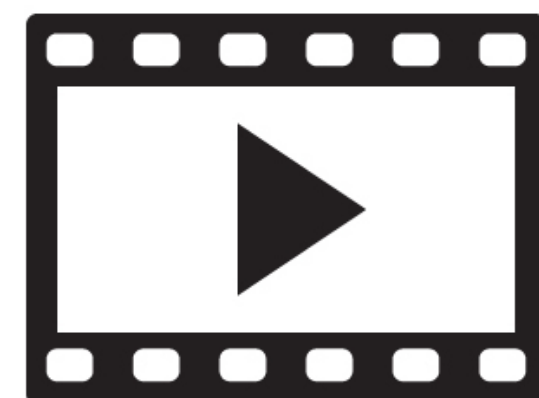
κείμενο



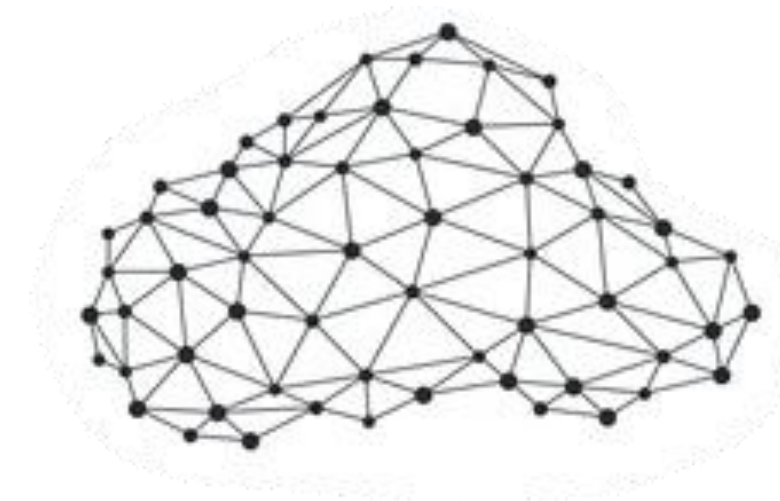
σήματα



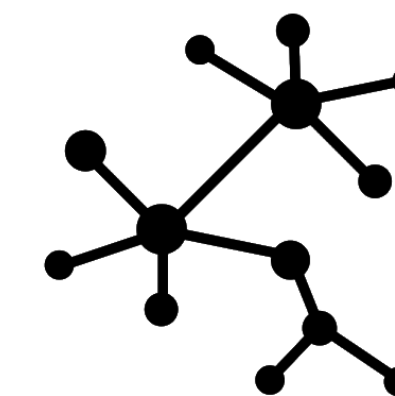
εικόνες



βίντεο



point clouds



γράφοι



Συλλογή δεδομένων

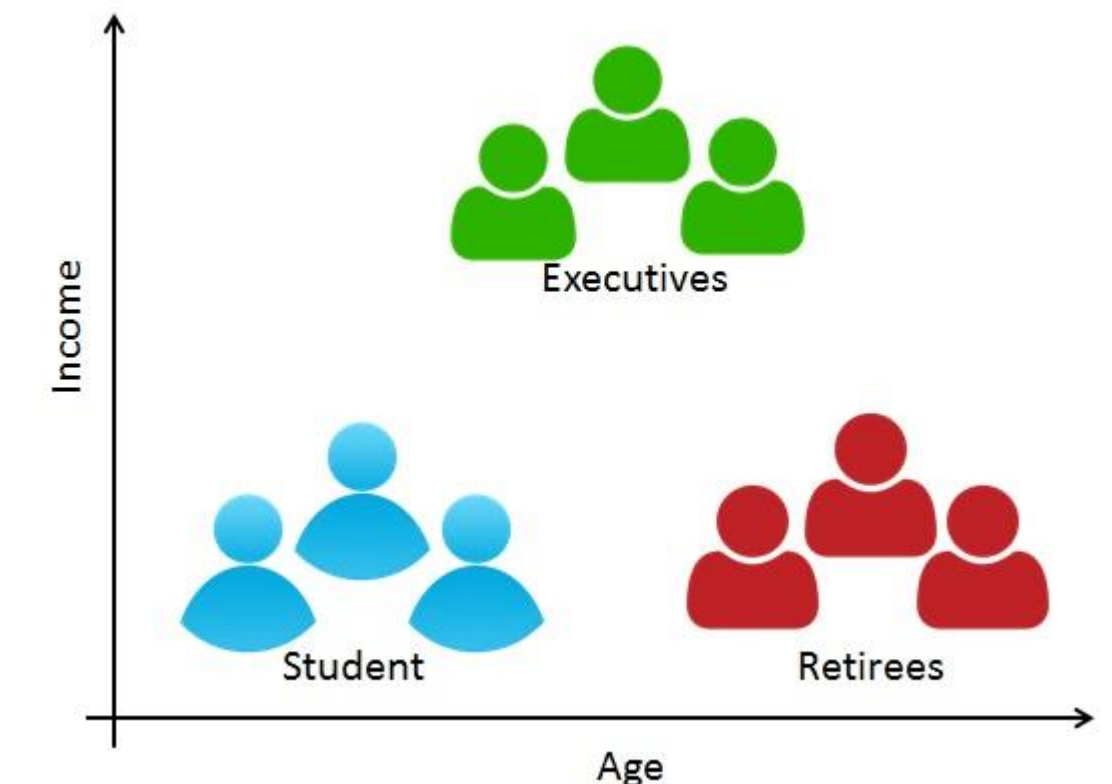




Συλλογή δεδομένων

Ενδέχεται να χρειαστεί να **αποκτηθούν** δεδομένα για τη συγκεκριμένη εργασία

- εικόνες φρούτων για ταξινομητή εικόνων φρούτων
- ερωτηματολόγιο με δημογραφικά στοιχεία για την κατάτμηση των πελατών
- ένας ρομποτικός βραχίονας μαθαίνει να συλλαμβάνει διαφορετικά αντικείμενα



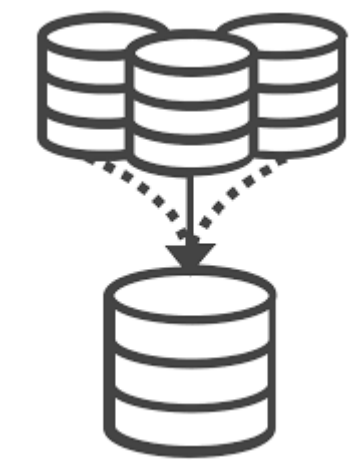


Συλλογή δεδομένων

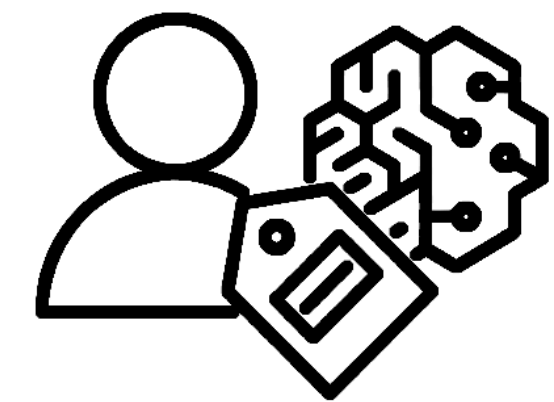
Τα δεδομένα υπάρχουν ήδη:

- μπορεί να χρειαστεί να **ενσωματωθούν** από πολλαπλές πηγές

π.χ. στον εντοπισμό απάτης: χρονοδιάγραμμα μεταξύ δύο διαδοχικών συναλλαγών και απόσταση μεταξύ των τόπων όπου πραγματοποιήθηκαν οι συναλλαγές



- μπορεί να χρειαστεί να **επισημανθούν**
 - όλες οι εποπτευόμενες ρυθμίσεις μάθησης
 - απαιτεί ανθρώπινη προσπάθεια (συχνά ειδικοί τομέα)
 - μπορούν να ανατεθούν σε εξωτερικούς συνεργάτες (υπηρεσίες επισήμανσης δεδομένων, π.χ., Amazon Mechanical Turk)
 - μεταφορά της μάθησης: χρησιμοποιήστε ένα ήδη εκπαιδευμένο μοντέλο με τις **ΕΤΙΚΕΤΕΣ ΤΟΥ**





Συλλογή δεδομένων

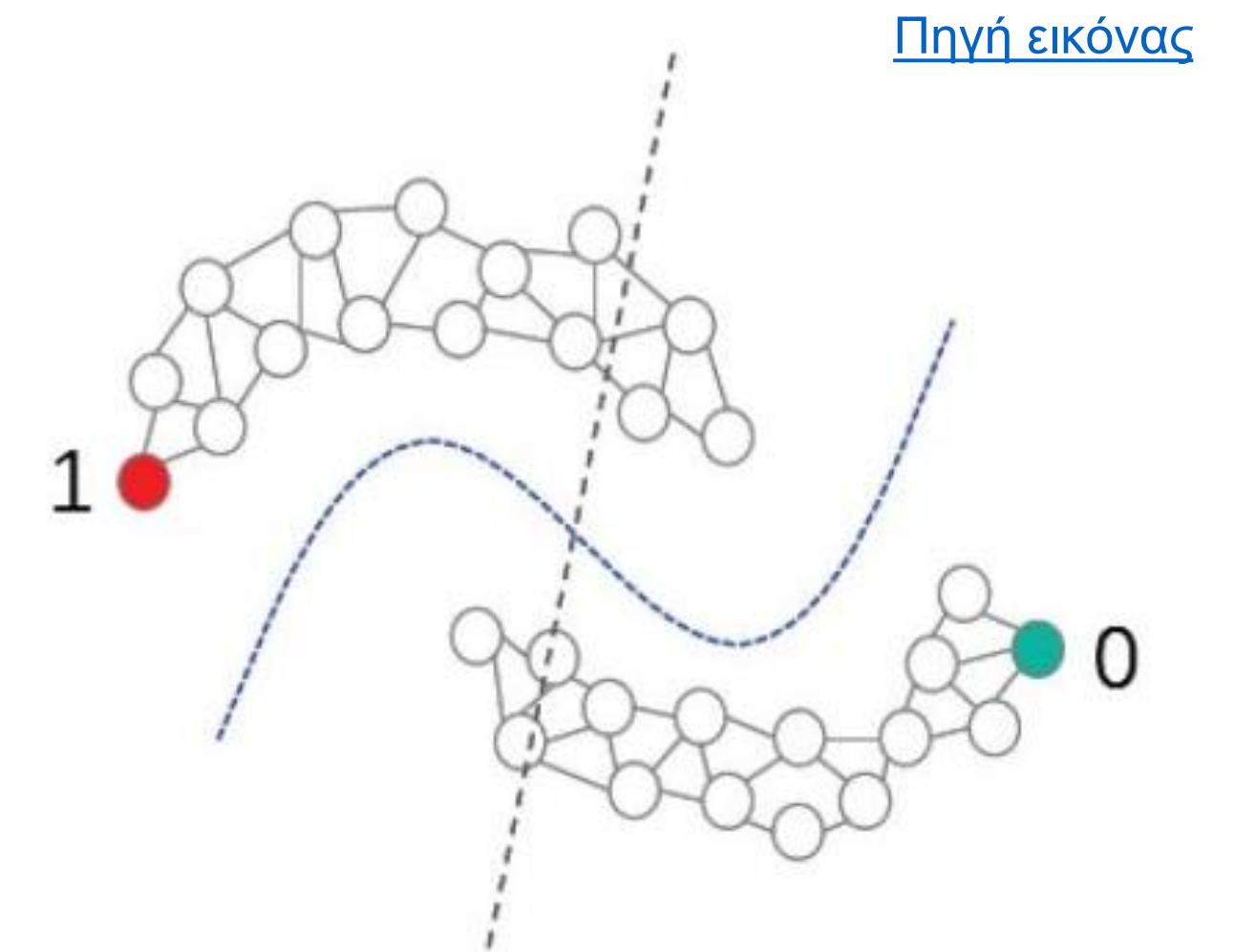
Πιο προηγμένοι τύποι MM που σχετίζονται με τη συλλογή δεδομένων:

Ημι-εποπτευόμενη μάθηση

- Μεγάλος όγκος μη εποπτευόμενων δεδομένων, μικρή ποσότητα δεδομένων με επισήμανση
- Μεταξύ μη εποπτευόμενης και εποπτευόμενης μάθησης

Ενεργός Μάθηση

- Χρησιμοποιείται όταν η επισήμανση είναι δαπανηρή
- Επιλέξτε έξυπνα το επόμενο σημείο δεδομένων για να επισημάνετε με τρόπο που ελαχιστοποιεί την προσπάθεια επισήμανσης



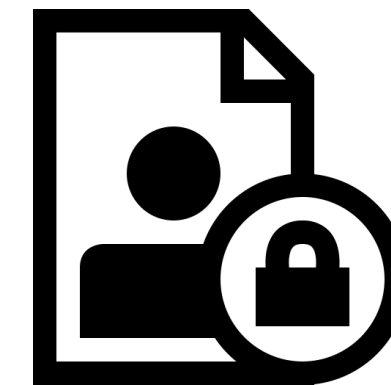


Συλλογή δεδομένων

Άλλα θέματα:

Προστασία προσωπικών δεδομένων

- Ανάγκη να εργαστείτε με ανωνυμοποιημένα δεδομένα
- Το στάδιο της προεπεξεργασίας αφορά ευαίσθητα δεδομένα



Ανισοροπία κλάσεων

- Πρέπει να έχουν περίπου τον ίδιο αριθμό περιπτώσεων για όλες τις κλάσεις σε εργασίες ταξινόμησης



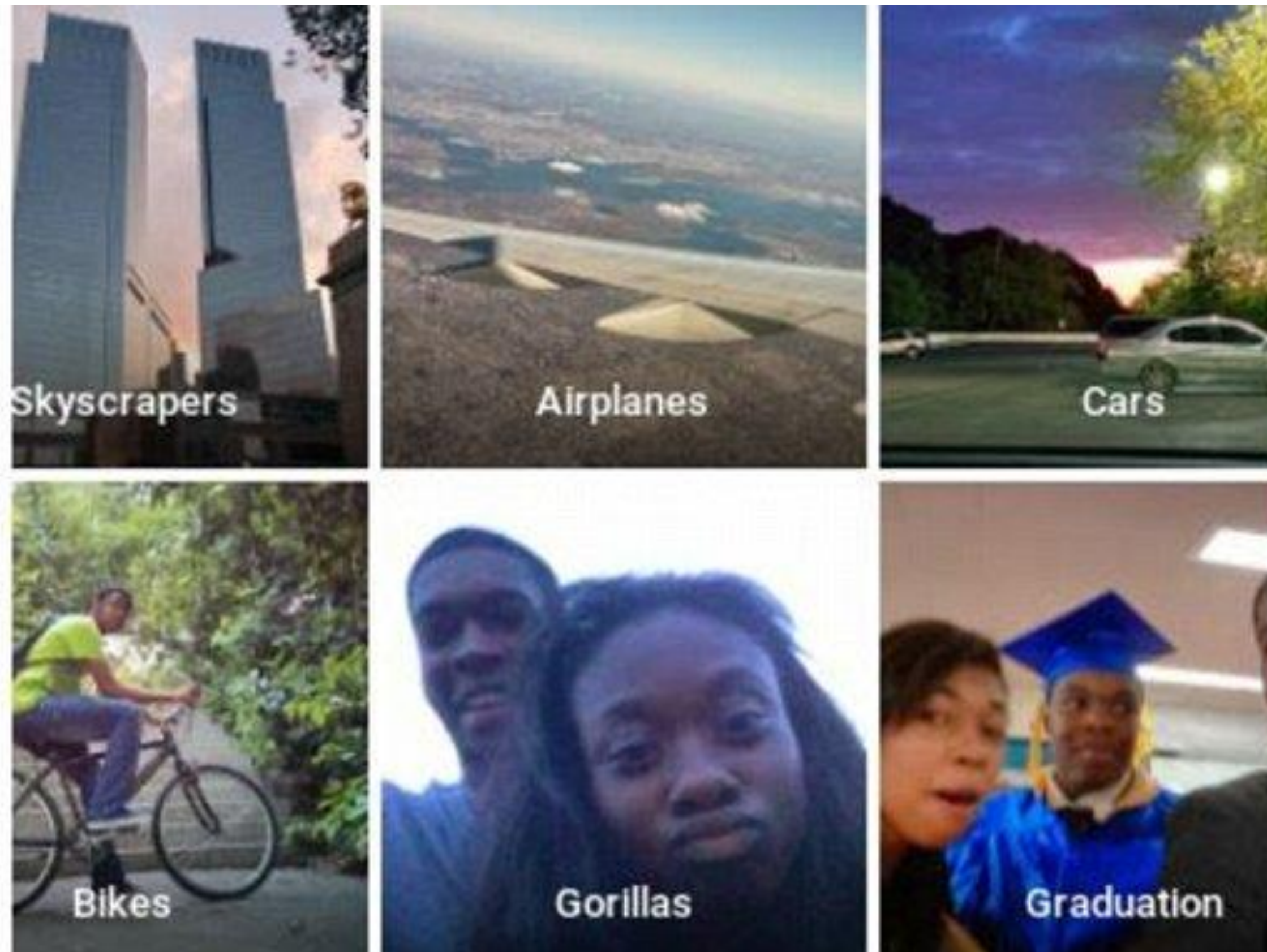
Προκατάληψη & Δικαιοσύνη

- Είναι πολύ σημαντικό να έχουμε ποικιλομορφία στα δεδομένα μας
- Δείτε τις επόμενες διαφάνειες





Προκατάληψη & δικαιοσύνη: Το περιστατικό των «Gorillas»



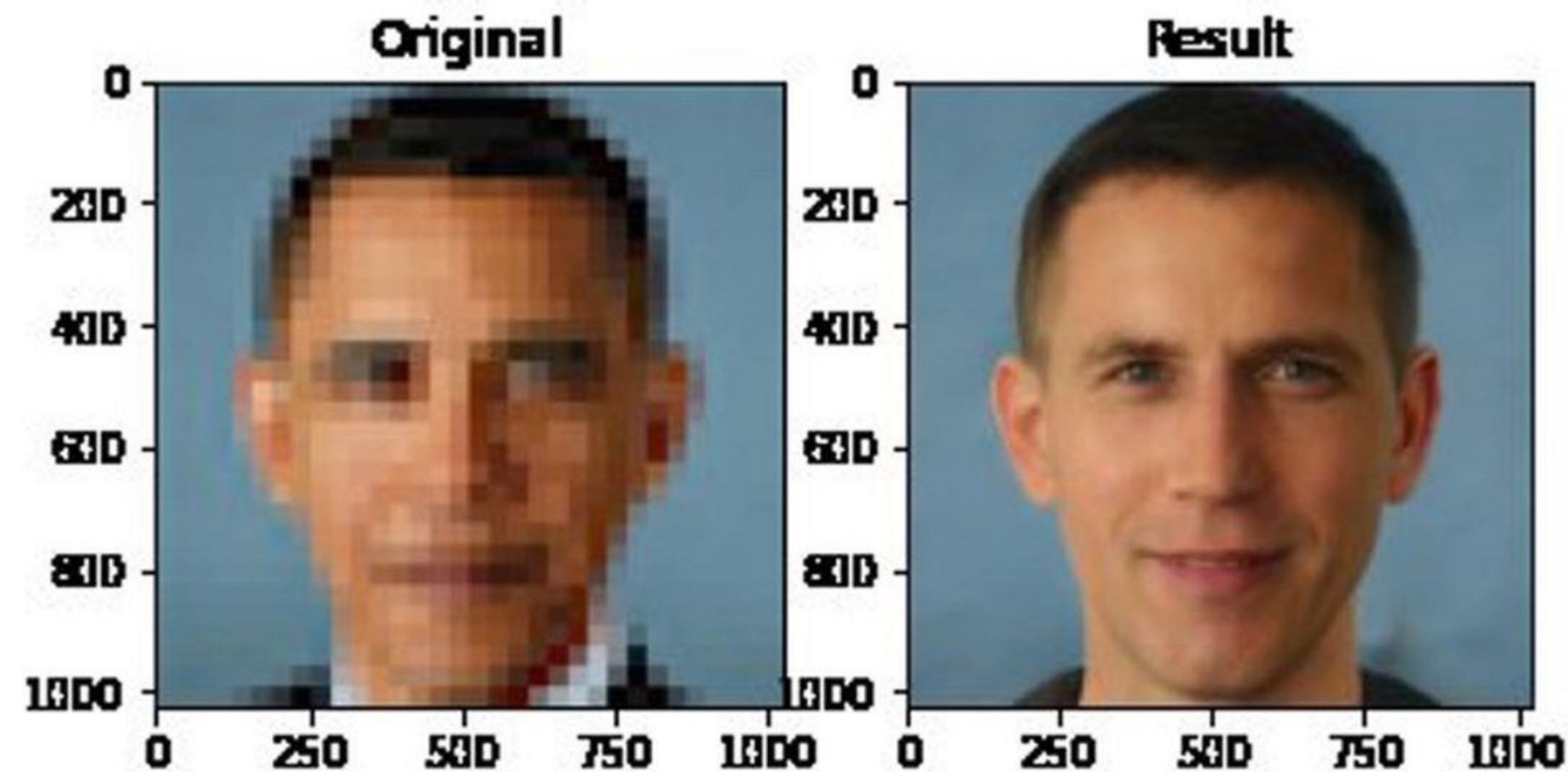
**Μία από αυτές τις
προβλέψεις
είναι πολύ λάθος!**

**Γιατί νομίζετε ότι συνέβη
αυτό;**





Προκατάληψη & δικαιοσύνη: Ο Ομπάμα ανακατασκευάστηκε ως λευκός άνδρας

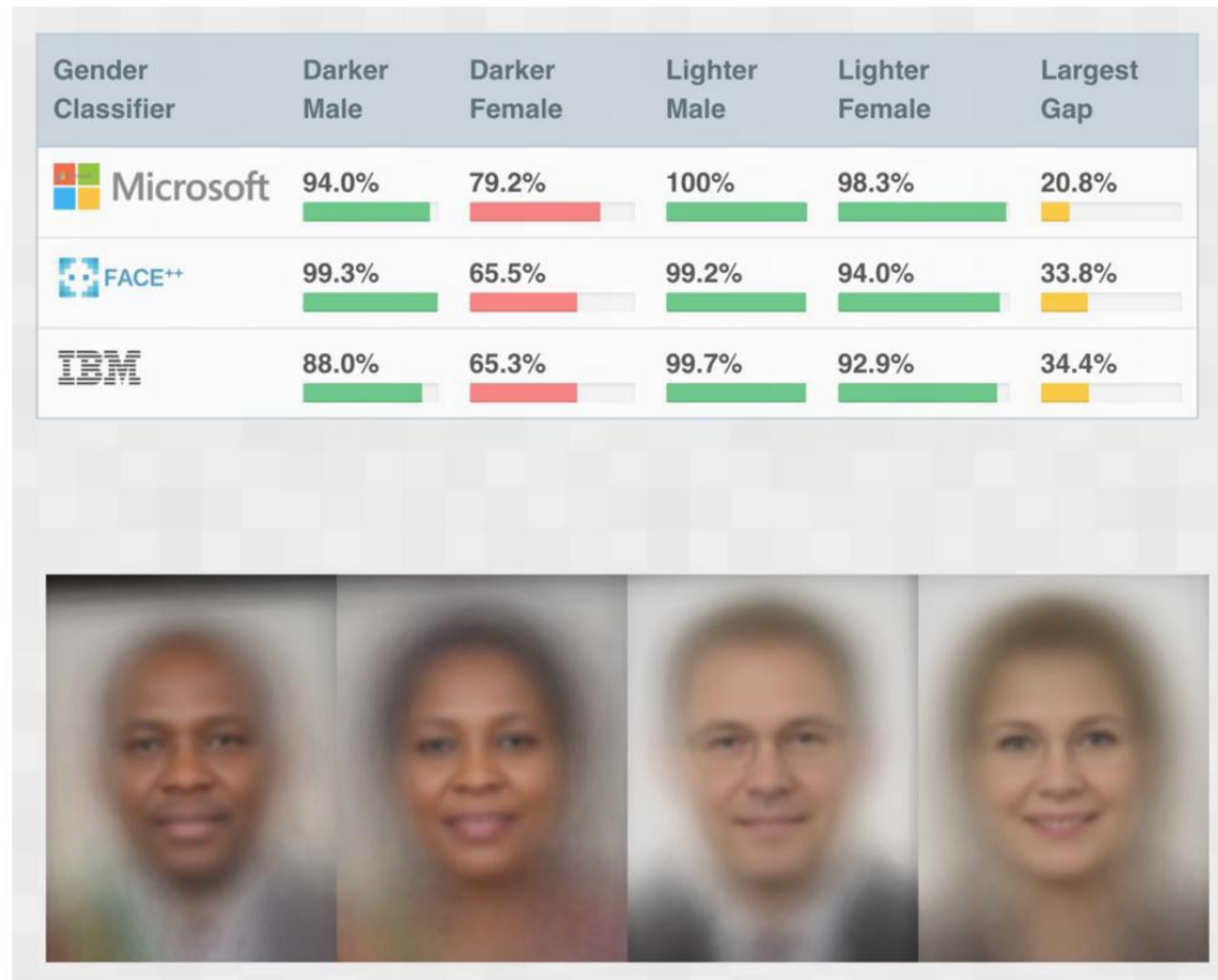


Το σύστημα αναδειγματοληψίας προσώπου κάνει όλους να φαίνονται λευκοί, επειδή εκπαιδεύτηκε σε ένα σύνολο δεδομένων που περιέχει κυρίως εικόνες λευκών ανθρώπων.





Προκατάληψη & δικαιοσύνη: Ποσοστό σφάλματος ανά φύλο και φυλή



Όλα τα μοντέλα έχουν περισσότερο πρόβλημα να ταξινομήσουν σωστά πιο έγχρωμες γυναίκες από λευκούς άντρες.

- Οι ταξινομητές δεν εκπαιδεύτηκαν με επαρκή ποικιλία
- Σημασία της ύπαρξης ισορροπημένων συνόλων δεδομένων

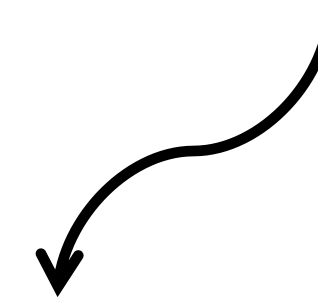
Το μοντέλο σας είναι τόσο καλό όσο τα δεδομένα σας!





Δεδομένα πίνακα

Μεταβλητή/Στοιχείο/Χαρακτηριστικό



Στήλη 1	Στήλη 2	Στήλη 3	Στήλη 4

Παράδειγμα/περίπτωση



Μεταβλητές εισόδου:
στήλες που παρέχονται
στο μοντέλο για την
πρόβλεψη

Μεταβλητή εξόδου: στήλη
που προβλέπεται από το
μοντέλο





Μεταβλητοί τύποι δεδομένων

Ποιοτικός/Κατηγορικός

Φτιαγμένη από **λέξεις**

1. Ονομαστική ονομασία (ονομαζόμενες κατηγορίες)
 - π.χ. φύλο
2. Ordinal (Θέματα Τάξης)
 - π.χ. διάθεση

Ποσοτικός/Αριθμητικός

Φτιαγμένη από **αριθμούς**

1. Διακριτό (integers)
 - π.χ. μέγεθος παπουτσιών
2. Συνεχής (οποιοσδήποτε αριθμός)
 - π.χ. βάρος





Στοιχεία μεταβλητοί τύποι κουίζ

Κατηγορηματικό (ονομαστικό, προφορικό) ή **αριθμητικό** (διακριτό, συνεχές);

1. Εθνικότητα

Ονομαστικό

2. Ηλικία

Συνεχής ή διακριτή

3. Χρώμα

Ονομαστικό

4. Θερμοκρασία

Συνεχής

5. Αριθμός παιδιών

Διακριτή

6. Βαθμολογία ικανοποίησης (δυστυχής/ουδέτερη/ευτυχής)

Ordinal

7. Ημερομηνία

String: Ανάγκη μετασχηματισμού
χαρακτηριστικών





Προεπξεργασία δεδομένων





Προεπεξεργασία δεδομένων

Τι είναι η προεπεξεργασία δεδομένων;

- Διαδικασία εισαγωγής **ακατέργαστων** δεδομένων σε μορφή κατάλληλη για μοντελοποίηση
- Τυπικά περιλαμβάνει: **Καθαρισμός δεδομένων, κωδικοποίηση δεδομένων**

Γιατί είναι σημαντικό;

Οι αλγόριθμοι MM αναμένουν αριθμούς (π.χ. ορισμένοι δεν μπορούν να ασχοληθούν με κατηγορηματικά χαρακτηριστικά)

- Οι αλγόριθμοι MM έχουν απαιτήσεις (π.χ. δεν μπορούν να αντιμετωπίσουν τις ελλείπουσες τιμές)





Προεπεξεργασία δεδομένων

Καθαρισμός δεδομένων

- Διορθώστε σφάλματα: τυπογραφικά λάθη, εσφαλμένη κεφαλαιοποίηση, ασυνέπειες κ.λπ. (π.χ. «Α/Α» έναντι «Άνευ αντικειμένου», «Άνδρας» έναντι «αρσενικού»,...)
- Ελλείποντα δεδομένα (αφαίρεση, αντικατάσταση με 0 ή παρεμβολή)
- Αφαίρεση διπλών γραμμών/στήλων
- Αφαίρεση άσχετων δεδομένων (π.χ. κατά την ανάλυση δεδομένων για πελάτες χιλιετηρίδων, αφαίρεση παλαιότερων γενεών)
- Αφαίρεση ακραίων τιμών (μπορεί να οφείλεται σε σφάλμα στη συλλογή δεδομένων)

Κωδικοποίηση δεδομένων

- Κατηγορηματικές μεταβλητές
- Ανωνυμοποίηση δεδομένων





Κωδικοποίηση δεδομένων

Κωδικοποιήστε τις **κατηγορηματικές** μεταβλητές χρησιμοποιώντας μια αριθμητική αναπαράσταση. Τυπικά:

- **Nominal variables**: κωδικοποίηση «one-hot»: όσες νέες δυαδικές μεταβλητές έχουν ονομαστεί κατηγορίες

Φύλο = {αρσενικό, θηλυκό} -> GenderMale = {0,1}, GenderFemale = {0,1}

Χρώμα = {κόκκινο, πράσινο, μπλε} -> ColorRed = {0,1}, ColorGreen = {0,1}, ColorBlue = {0,1}

- **Τακτικές μεταβλητές**: τακτικός κωδικοποίηση: μετατροπή σε ακέραιες τιμές

Ικανοποίηση = {sad,neutral,happy} -> {0,1,2}





Κωδικοποίηση στοιχείων: One-hot παράδειγμα κωδικοποίησης

```
In [1]: # example of a one-hot encoding
...: from numpy import asarray
...: from sklearn.preprocessing import OneHotEncoder
...:
...: # define data
...: data = asarray([[ 'red' ], [ 'green' ], [ 'blue' ]])
...: print(data)
...:
...: # define one hot encoding
...: encoder = OneHotEncoder(sparse=False)
...:
...: # transform data
...: onehot = encoder.fit_transform(data)
...: print(onehot)
[ 'red' ]
[ 'green' ]
[ 'blue' ]
[[0. 0. 1.]
 [0. 1. 0.]
 [1. 0. 0.]
```





Κωδικοποίηση στοιχείων: Ordinal παράδειγμα κωδικοποίησης

```
In [1]: # example of an ordinal encoding
...: from numpy import asarray
...: from sklearn.preprocessing import OrdinalEncoder
...:
...: # define data
...: satisfaction_rating = asarray(['sad', 'neutral', 'happy'])
...: print(satisfaction_rating)
...:
...: # define ordinal encoding
...: encoder = OrdinalEncoder()
...:
...: # transform to numbers
...: result = encoder.fit_transform(satisfaction_rating)
...: print(result)
[['sad']
 ['neutral']
 ['happy']]
[[2.]
 [1.]
 [0.]]
```





ΟΠΤΙΚΟΠΟΙΗΣΗ ΔΕΔΟΜΕΝΩΝ





Οπτικοποίηση δεδομένων

Ανάλυση διερευνητικών δεδομένων

- υπάρχουν εμφανή μοτίβα (μπορεί να δώσει συμβουλές σχετικά με τη μέθοδο για να επιλέξετε);
- υπάρχουν προφανή προβλήματα (θόρυβος ετικέτας ή υπερβολικές τιμές);
- περιλαμβάνει βασικά στατιστικά στοιχεία (π.χ. μέσος όρος, διακύμανση, συντελεστής συσχέτισης) και γραφήματα

- Δεδομένα πίνακα με μικρό αριθμό χαρακτηριστικών: χρησιμοποιήστε **pair plot**

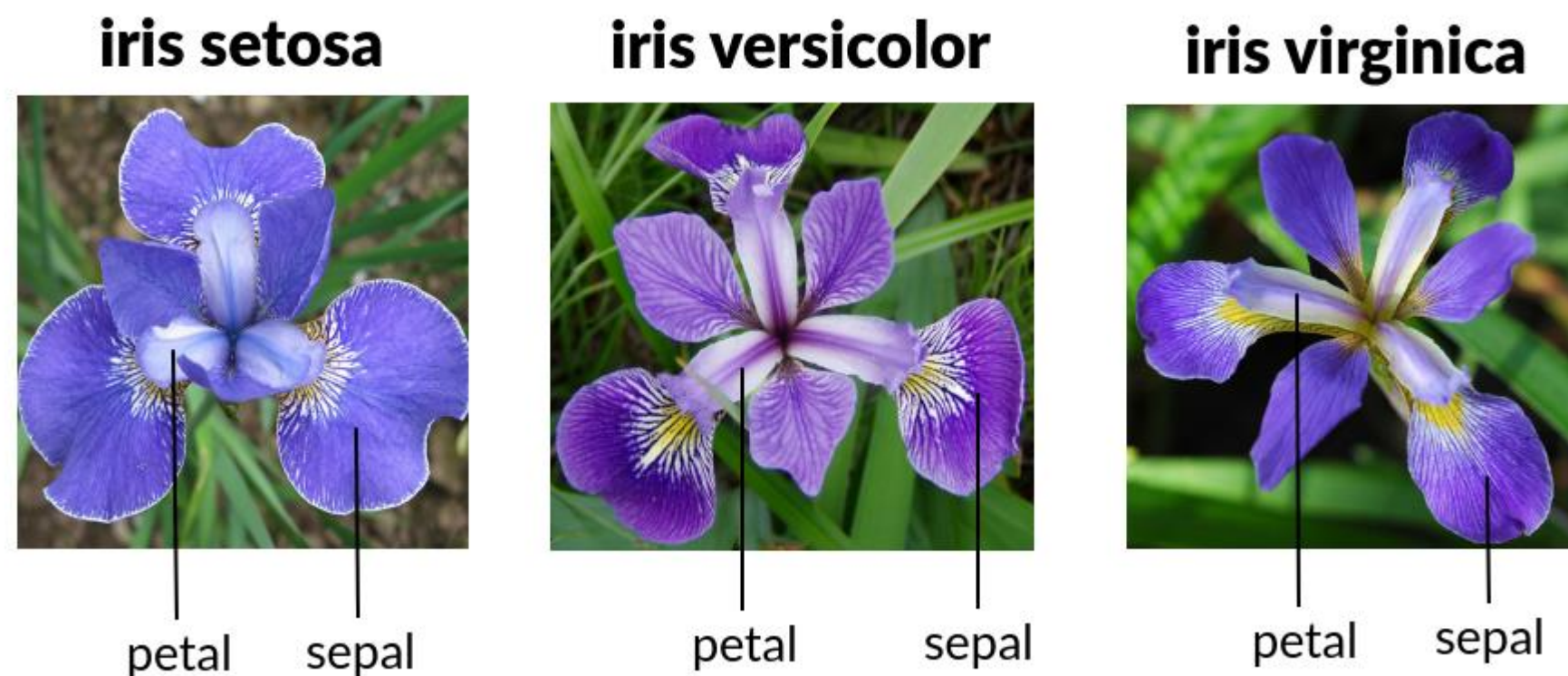
- Υψηλής διάστασης στοιχεία: χρήση της **μείωσης της διάστασης**
 - π.χ., ανάλυση βασικών στοιχείων
 - θα μελετήσει σε μεταγενέστερες διαλέξεις





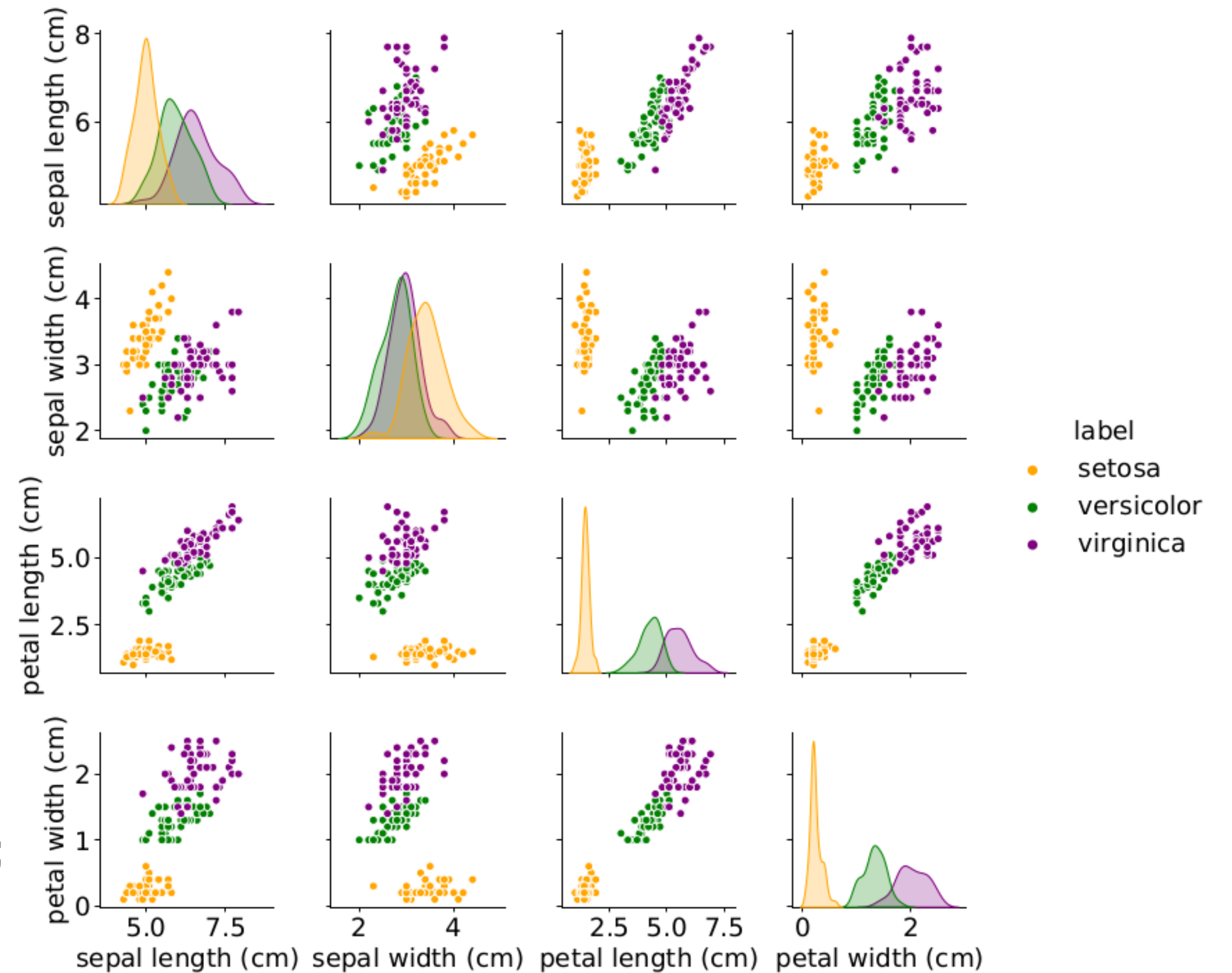
Σύνολο δεδομένων Iris

- 3 κλάσεις των 50 περιπτώσεων η κάθε μία
- 4 χαρακτηριστικά (μήκος/πλάτος πέτου/σεπάλ σε cm)



Ποια τάξη είναι γραμμικά διαχωρίσιμη από τις άλλες;

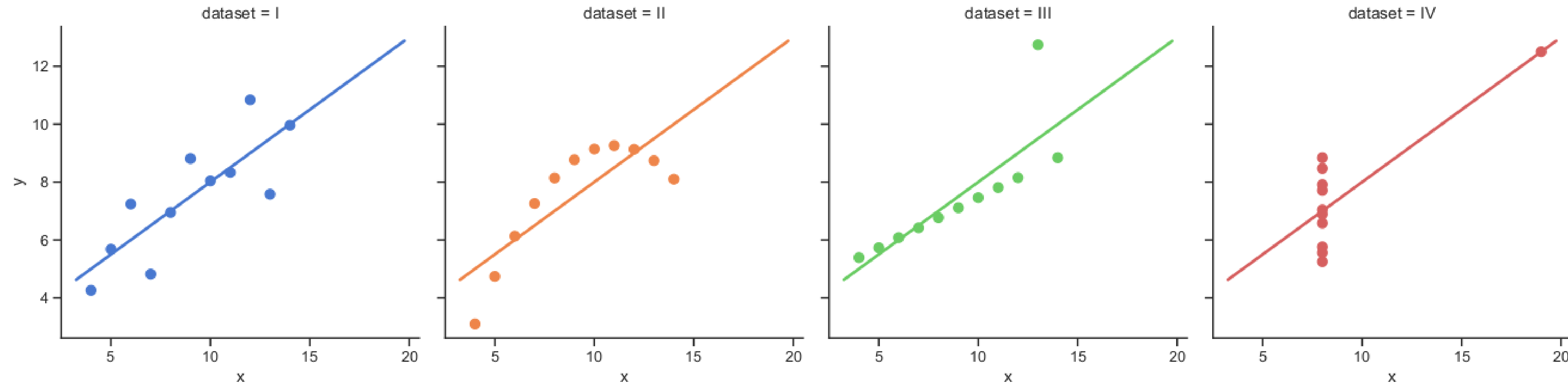
Πηγή: Μέρφι, Κ. (2022). Πιθανοτική μηχανική μάθηση — μια εισαγωγή. MIT Press





Σημασία της απεικόνισης των δεδομένων σε σχέση με τις συνοπτικές στατιστικές

Πηγή: Murphy, K. (2022). Probabilistic Machine Learning - An Introduction. MIT Press

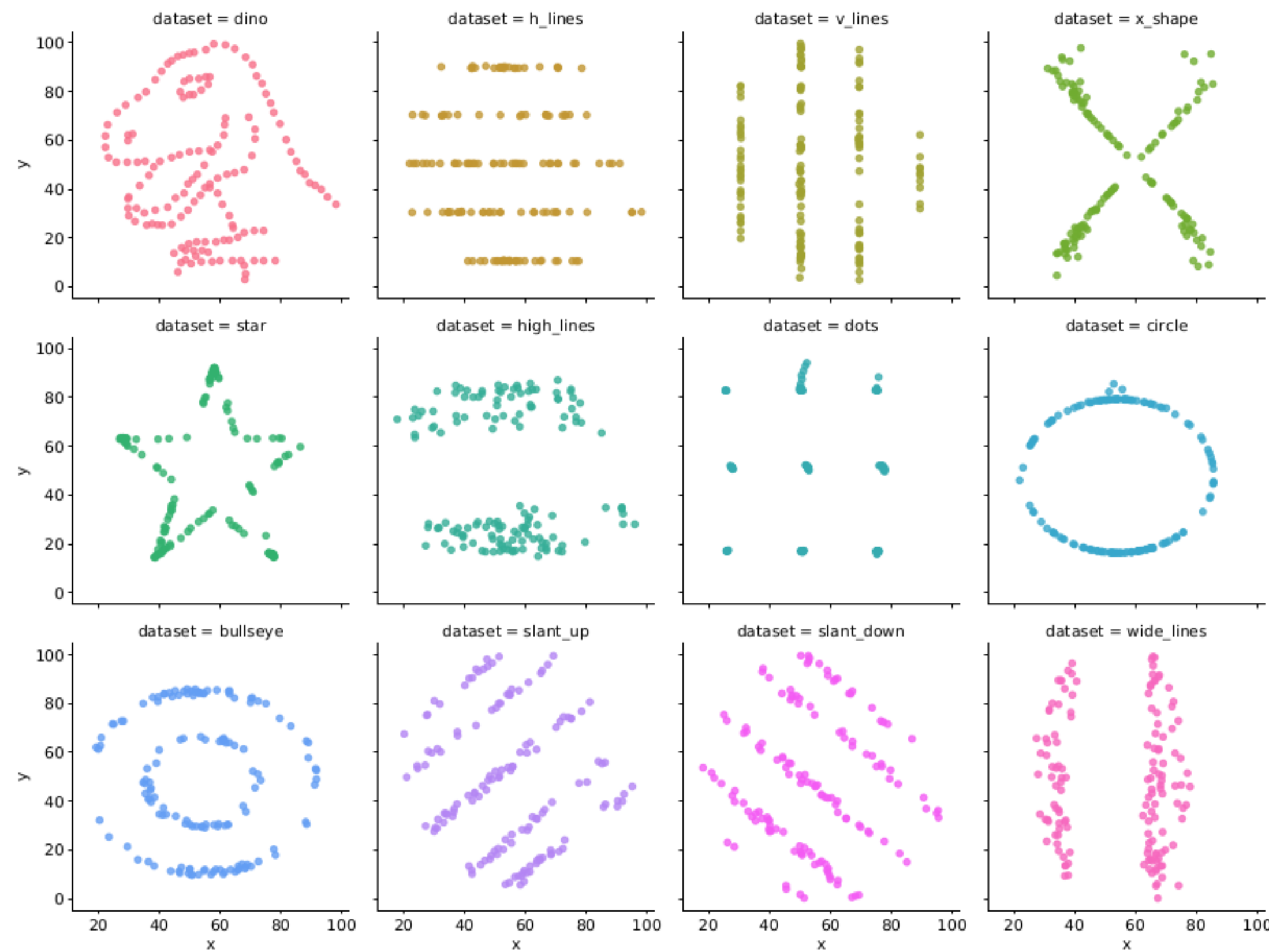


Το κουαρτέτο του Anscombe: Όλα αυτά τα σύνολα δεδομένων έχουν τον ίδιο μέσο όρο, συντελεστή διακύμανσης και συσχέτισης.





Σημασία της απεικόνισης των δεδομένων σε σχέση με τις συνοπτικές στατιστικές



Πηγή: Murphy, K. (2022). Probabilistic Machine Learning - An Introduction. MIT Press

Ντουζίνα Datasaurus: Όλα αυτά τα σύνολα δεδομένων έχουν τον ίδιο μέσο όρο, συντελεστή διακύμανσης και συσχέτισης.





Μετασχηματισμός δεδομένων





Μετασχηματισμός δεδομένων

Τι είναι ο μετασχηματισμός των δεδομένων;

- Διαδικασία μετασχηματισμού δεδομένων σε μορφή **καταλληλότερη** για μοντελοποίηση
- Τυπικά περιλαμβάνει: **Κλιμάκωση Χαρακτηριστικών, Επιλογή Χαρακτηριστικών, Εξαγωγή Χαρακτηριστικών και Κατασκευή Χαρακτηριστικών (ή Μηχανική)**
- Μερικές φορές: **Αύξηση δεδομένων, Δειγματοληψία δεδομένων**

Γιατί είναι σημαντικό;

- Οι αλγόριθμοι MM έχουν απαιτήσεις (π.χ., η πολυπλοκότητα ορισμένων αλγορίθμων κλιμάκων με αριθμό σημείων δεδομένων)
- Η απόδοση του μοντέλου εξαρτάται από τα δεδομένα (π.χ., διαφορετικές περιοχές χαρακτηριστικών επηρεάζουν πολλούς αλγορίθμους)





Χαρακτηριστικό γνώρισμα Κλιμάκωση

Ελάχιστη-μέγιστη κανονικοποίηση (ή γραμμική κλιμάκωση):

- χρησιμοποιημένος όταν το χαρακτηριστικό γνώρισμα κατανέμεται ομοιόμορφα σε μια σταθερή περιοχή
- συνήθως κλιμακώνεται στο εύρος [0,1]
- μπορεί να επηρεαστεί σε μεγάλο βαθμό από τις ακραίες τιμές

$$x_{scaled} = \frac{x - x_{min}}{x_{max} - x_{min}}$$

Τυποποίηση (ή εξομάλυνση z-score):

- χρησιμοποιείται όταν δεν γνωρίζουμε τη σειρά χαρακτηριστικών
- οι κλιμακούμενες κατανομές χαρακτηριστικών έχουν μέσο = 0 και sd=1
- λιγότερο επηρεαζόμενες από ακραίες τιμές

$$x_{scaled} = \frac{x - mean}{sd}$$





Επιλογή χαρακτηριστικών

Διαδικασία που επιλέγει ένα υποσύνολο χαρακτηριστικών M από το αρχικό σύνολο χαρακτηριστικών N ($M < N$), έτσι ώστε ο χώρος χαρακτηριστικών να μειώνεται βέλτιστα σύμφωνα με ένα συγκεκριμένο κριτήριο.

Ρόλος:

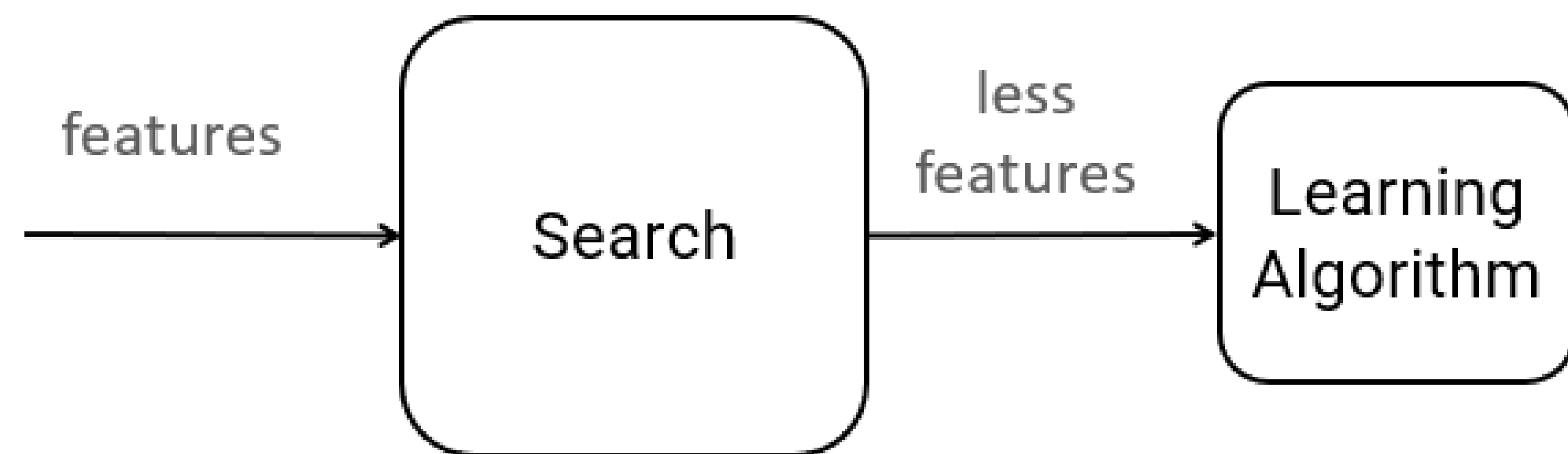
1. μειώστε τη διάσταση του χώρου χαρακτηριστικών
2. επιταχύνετε έναν αλγόριθμο εκμάθησης
3. βελτίωση της προγνωστικής ακρίβειας ενός μοντέλου
4. βελτίωση της κατανόησης των αποτελεσμάτων





Επιλογή χαρακτηριστικών

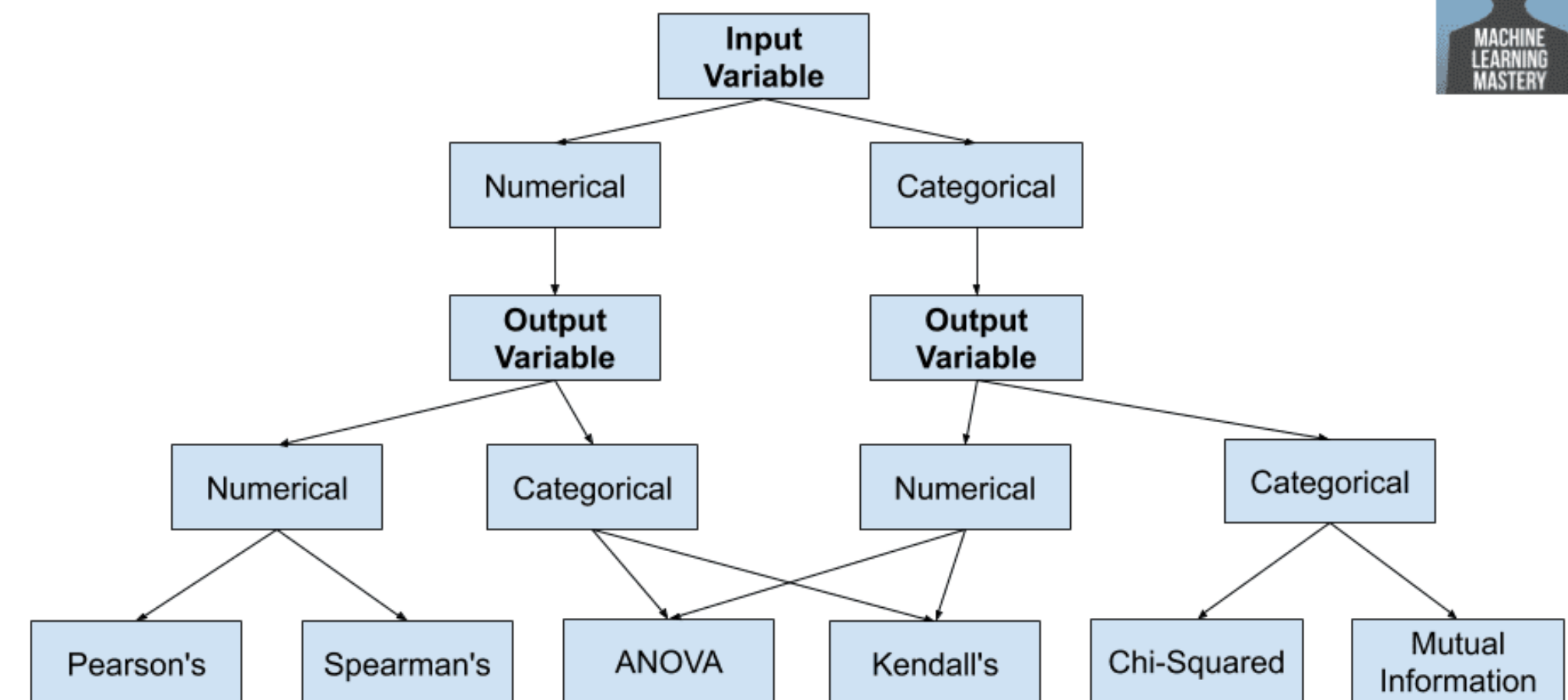
Μέθοδοι φίλτρου



Επιλέξτε υποσύνολα των χαρακτηριστικών με βάση τη **σχέση τους με το στόχο**

Συνήθως με τη χρήση στατιστικών τεχνικών

How to Choose a Feature Selection Method



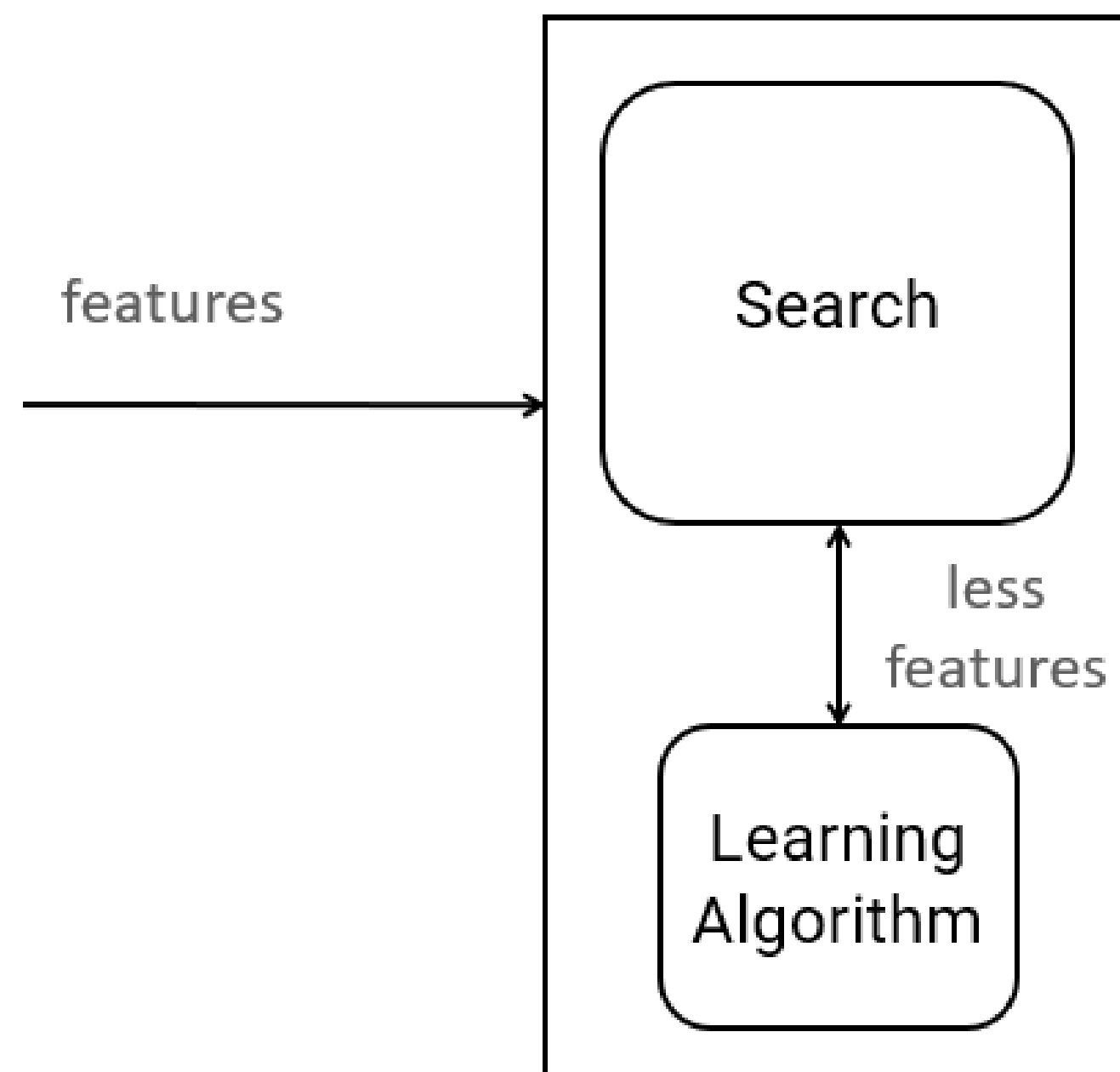
Copyright © MachineLearningMastery.com





Επιλογή χαρακτηριστικών

Μέθοδοι περιτυλίγματος



Αναζήτηση για **καλά αποδοτικά** υποσύνολα χαρακτηριστικών

Παραδείγματα:

- Αναδρομική Εξάλειψη Χαρακτηριστικών
- Ευρετικοί αλγόριθμοι αναζήτησης (π.χ. γενετικοί αλγόριθμοι, βελτιστοποίηση σμήνους σωματιδίων)

Υπολογιστικό κόστος;

Για κάθε υποσύνολο, ένα νέο μοντέλο εκπαιδεύεται και δοκιμάζεται για να επιτευχθεί η ακρίβεια





Επιλογή χαρακτηριστικών

Εγγενείς/έμμεσες/ενσωματωμένες μέθοδοι

Αυτόματη επιλογή χαρακτηριστικών κατά τη διάρκεια της εκπαίδευσης

Παραδείγματα:

- Μοντέλα που βασίζονται σε κανόνες
- Δενδρικά μοντέλα (θα μελετηθούν σε μεταγενέστερη διάλεξη)
- L1 τακτοποίηση (θα μελετήσει σε μια μεταγενέστερη διάλεξη)





Εξαγωγή χαρακτηριστικών (μείωση διάστασης)

Διαδικασία που εξάγει ένα σύνολο **νέων χαρακτηριστικών** M από τα αρχικά χαρακτηριστικά N ($M < N$) μέσω κάποιας λειτουργικής χαρτογράφησης.

Στόχος: αναζητήστε ένα ελάχιστο σύνολο **νέων χαρακτηριστικών** μέσω κάποιου μετασχηματισμού σύμφωνα με κάποιο **μέτρο απόδοσης**.

Προσέγγιση:

- Ανάλυση κύριων συστατικών (και άλλες μέθοδοι μείωσης της διάστασης)
- Νευρωνικά δίκτυα
- Θα μελετηθούν σε μεταγενέστερες διαλέξεις





Κατασκευή χαρακτηριστικών (μηχανική χαρακτηριστικών)

Διαδικασία που ανακαλύπτει **ελλείπουσες πληροφορίες** σχετικά με τις σχέσεις μεταξύ των χαρακτηριστικών και **αυξάνει** το χώρο των δυνατοτήτων με τη συναγωγή ή τη δημιουργία **πρόσθετων** χαρακτηριστικών.

Μπορεί να γίνει με τη χρήση:

- Αυτοματοποιημένες μέθοδοι - παραδείγματα:
 - Αριθμητικά χαρακτηριστικά: πολυωνυμική επέκταση
 - Ονομαστικά χαρακτηριστικά γνωρίσματα: σύνδεση, αποσύνδεση, άρνηση
- Γνώση τομέα (π.χ., $\text{SurfaceArea} = \text{Ύψος} \times \text{Πλάτος}$)

Στόχος: αυξήστε την εκφραστική δύναμη των αρχικών χαρακτηριστικών





Παράδειγμα: Πολυωνυμικά χαρακτηριστικά

```
In [1]: import numpy as np
...: from sklearn.preprocessing import PolynomialFeatures
...: X = np.arange(6).reshape(3, 2)
...: print(X)
...:
...: poly = PolynomialFeatures(2)
...: poly.fit_transform(X)
[[0 1]
 [2 3]
 [4 5]]
Out[1]:
array([[ 1.,  0.,  1.,  0.,  0.,  1.],
       [ 1.,  2.,  3.,  4.,  6.,  9.],
       [ 1.,  4.,  5., 16., 20., 25.]])
```

Πολυώνυμο 2ου βαθμού:

Είσοδος: $[a, b]$

Έξοδος: $[1, a, b, a^2, ab, b^2]$.





Κουίζ

- Η επιλογή χαρακτηριστικών συνήθως μειώνει τον αριθμό των χαρακτηριστικών

ΣΩΣΤΟ

- Η εξαγωγή χαρακτηριστικών συνήθως αυξάνει τον αριθμό των χαρακτηριστικών

ΛΑΘΟΣ

- Χαρακτηριστικό κατασκευής αυξάνει συνήθως τον αριθμό των χαρακτηριστικών

ΛΑΘΟΣ





Αύξηση δεδομένων (ή υπερδειγματοληψία δεδομένων)

Τεχνική που εισάγει πρόσθετα σημεία δεδομένων (περιστάσεις)

Χρησιμοποιείται όταν έχουμε:

- ένα **μικρό σύνολο δεδομένων**
- **σύνολο δεδομένων με ανισορροπία** (π.χ. 100 δείγματα κατηγορίας 0, 10 δείγματα της κατηγορίας 1)

Προσεγγίσεις:

- απλή τυχαία υπερδειγματοληψία (επαναλαμβανόμενα δεδομένα)
- προσθέστε μικρό θόρυβο Gaussian (για αριθμητικά χαρακτηριστικά)
- Smote [1] και παραλλαγές (για αριθμητικά χαρακτηριστικά)

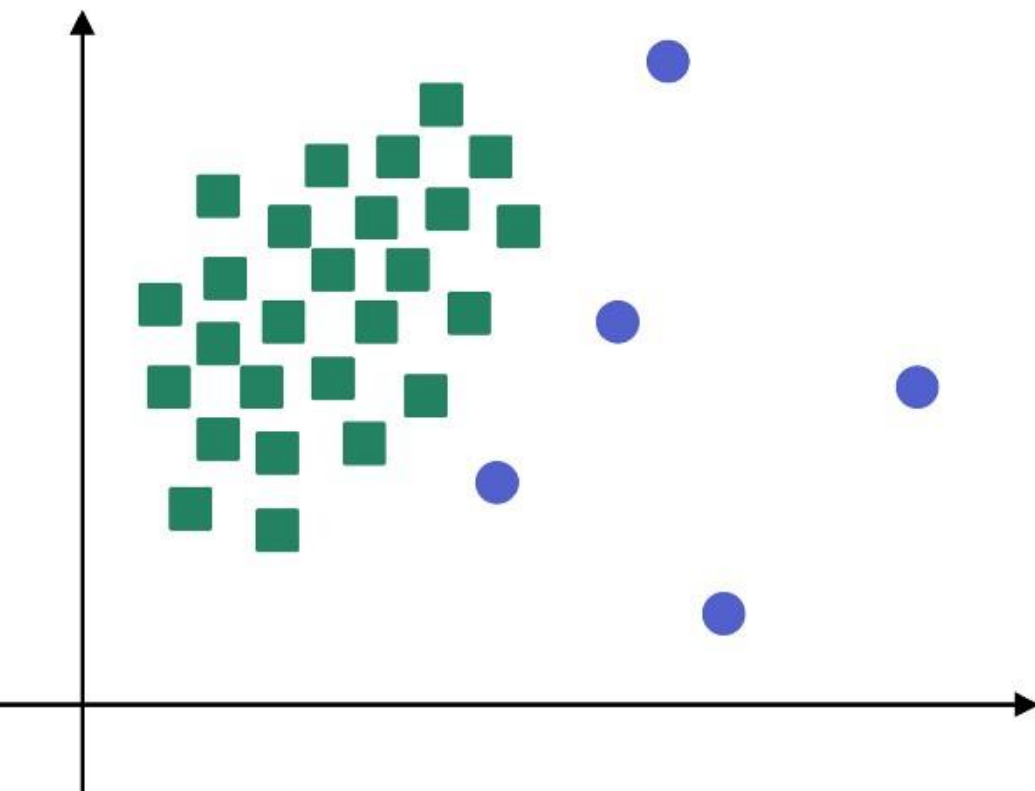
Δείτε [imbalanced-learn library](#)

[1] Nitesh et al. Smote: synthetic minority over-sampling technique. Journal of artificial intelligence research, 16:321–357, 2002.

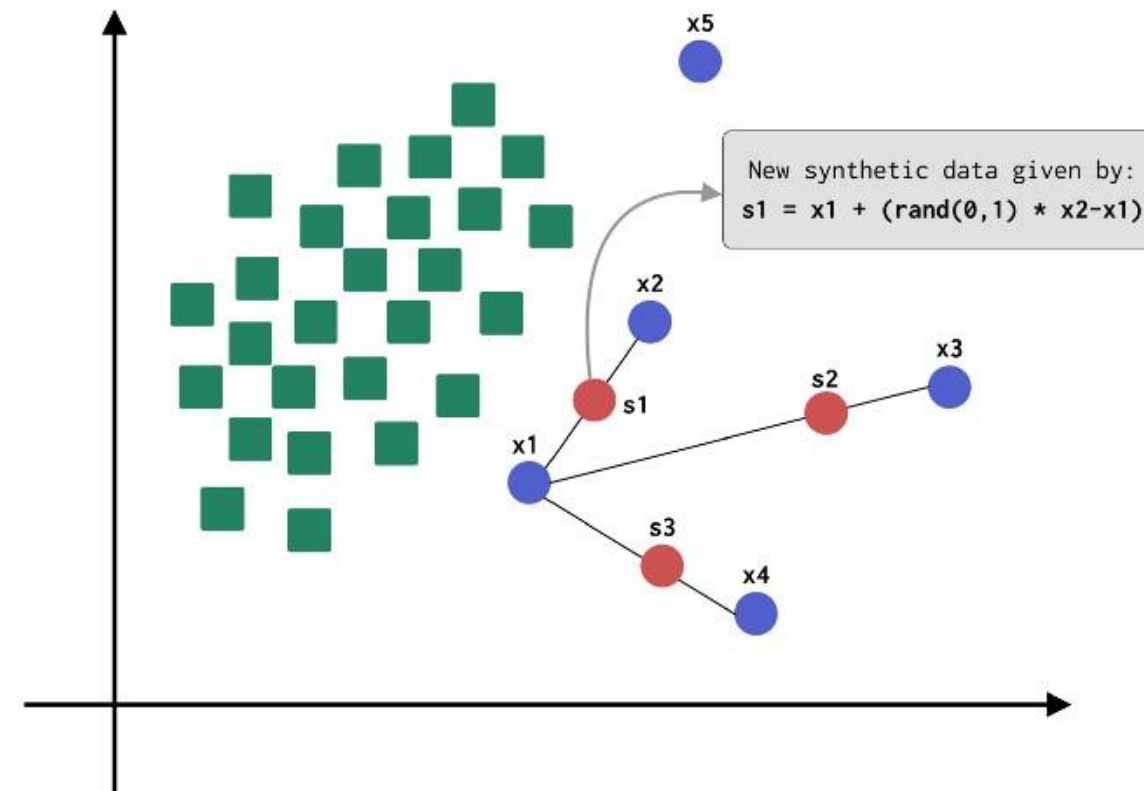




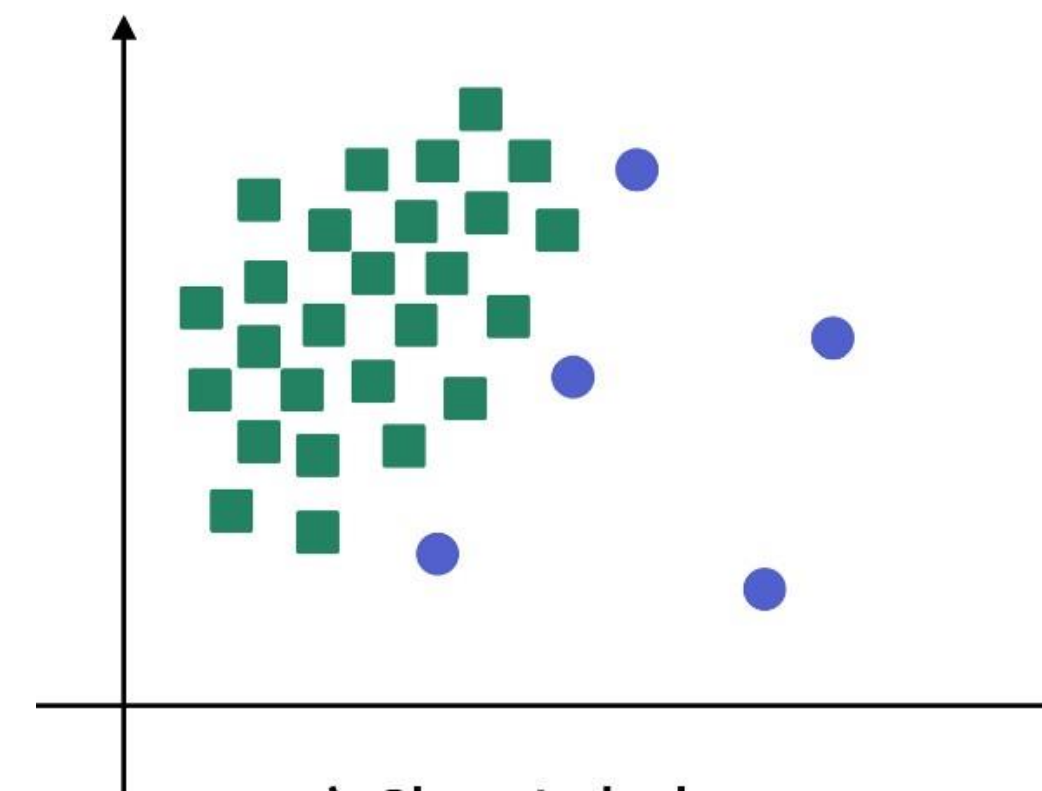
Αύξηση δεδομένων (ή υπερδειγματοληψία δεδομένων)



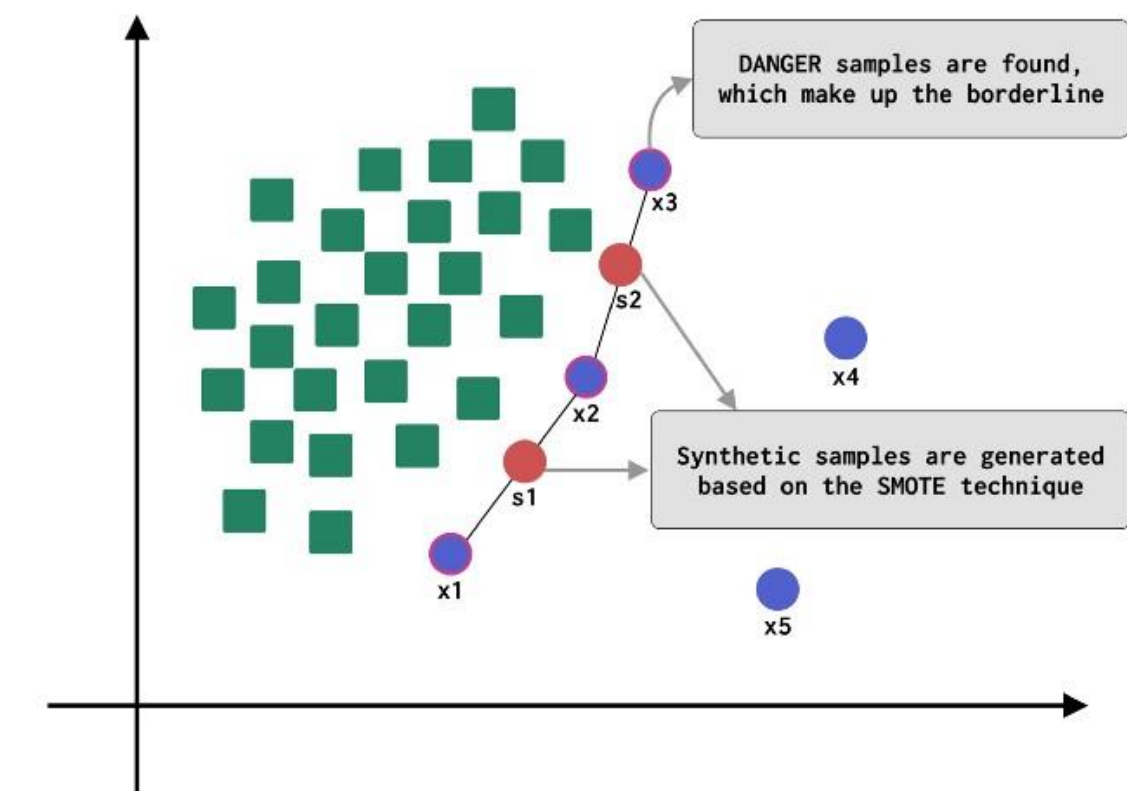
a) Class Imbalance



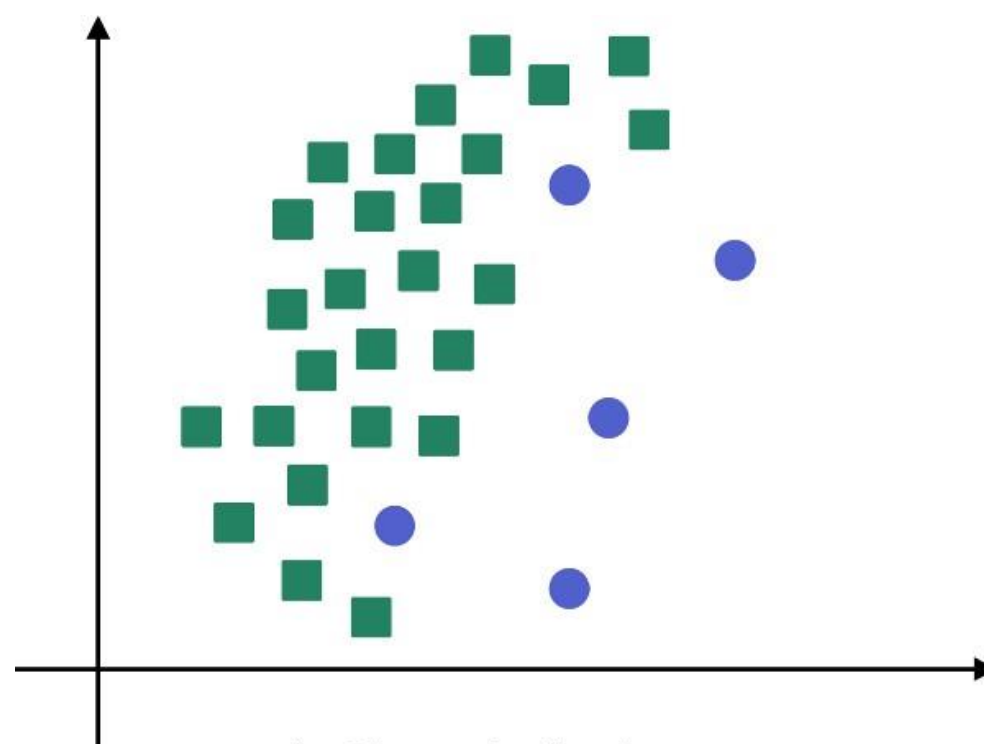
b) SMOTE



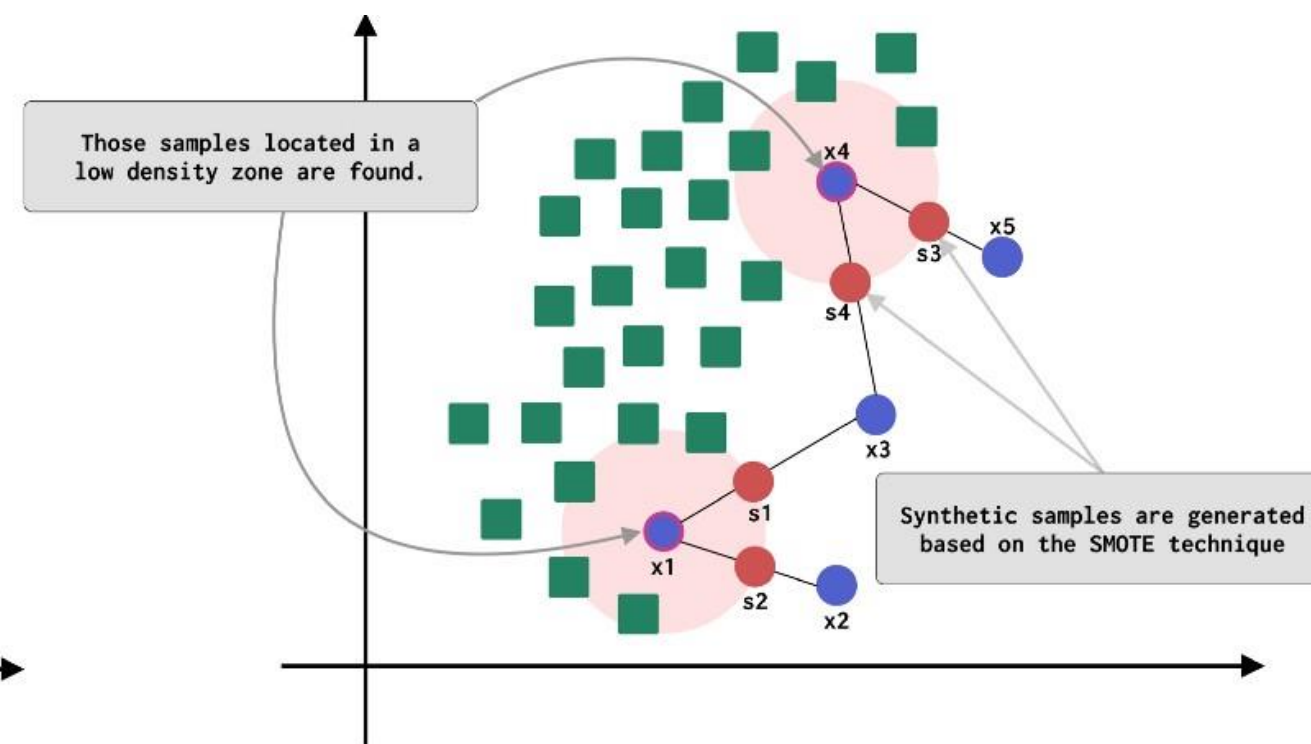
a) Class Imbalance



b) Borderline-SMOTE



a) Class Imbalance



b) ADASYN

Πηγή [ΕΙΚόνων](#)





Δειγματοληψία δεδομένων (ή μη δειγματοληψία δεδομένων)

Τεχνική που αφαιρεί τα σημεία δεδομένων (instances)

Χρησιμοποιείται όταν έχουμε:

- ένα **τεράστιο σύνολο δεδομένων** (που συνήθως δεν μπορεί να χωρέσει στη μνήμη)
- **σύνολο δεδομένων με ανισορροπία** (π.χ. 1000 δείγματα κατηγορίας 0, 100 δείγματα της κατηγορίας 1)

Προσεγγίσεις:

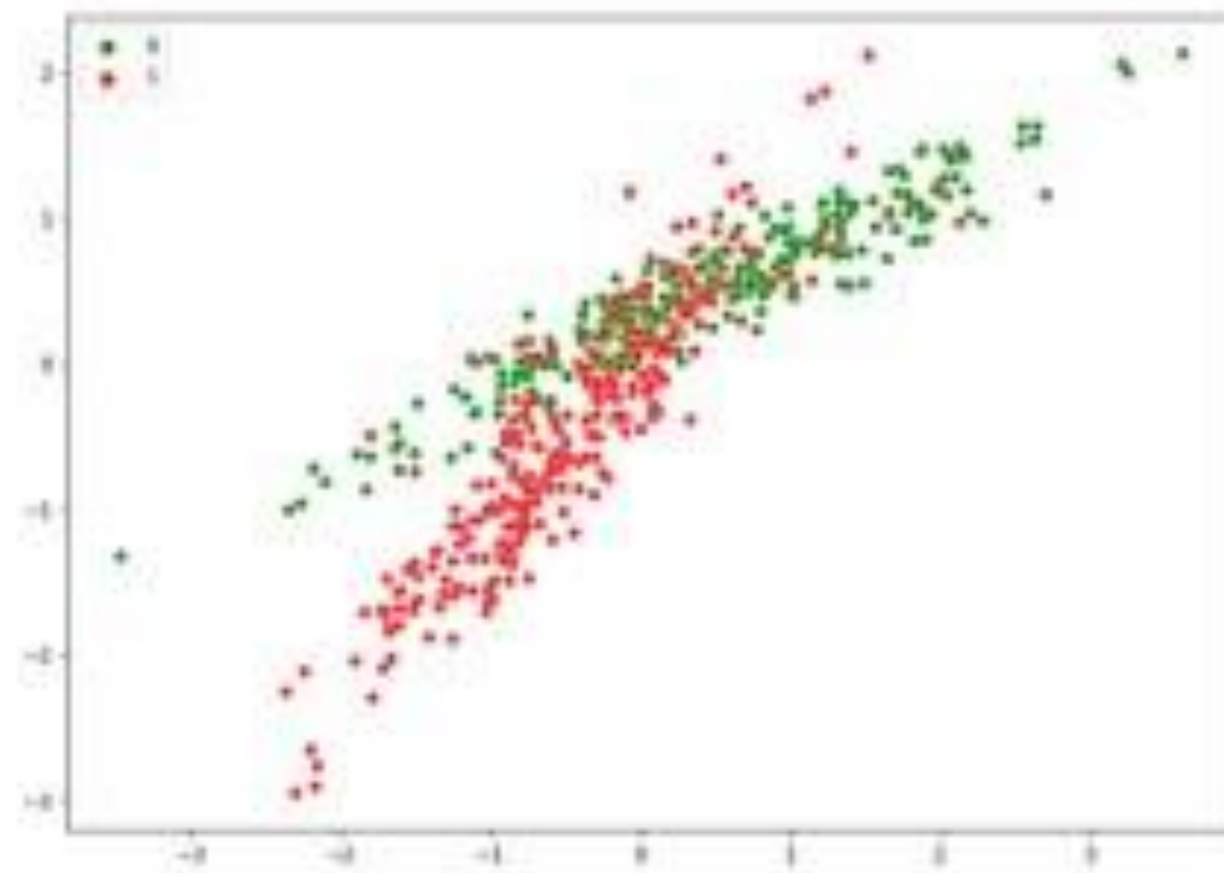
- απλή τυχαία δειγματοληψία
- δειγματοληψία συστάδων





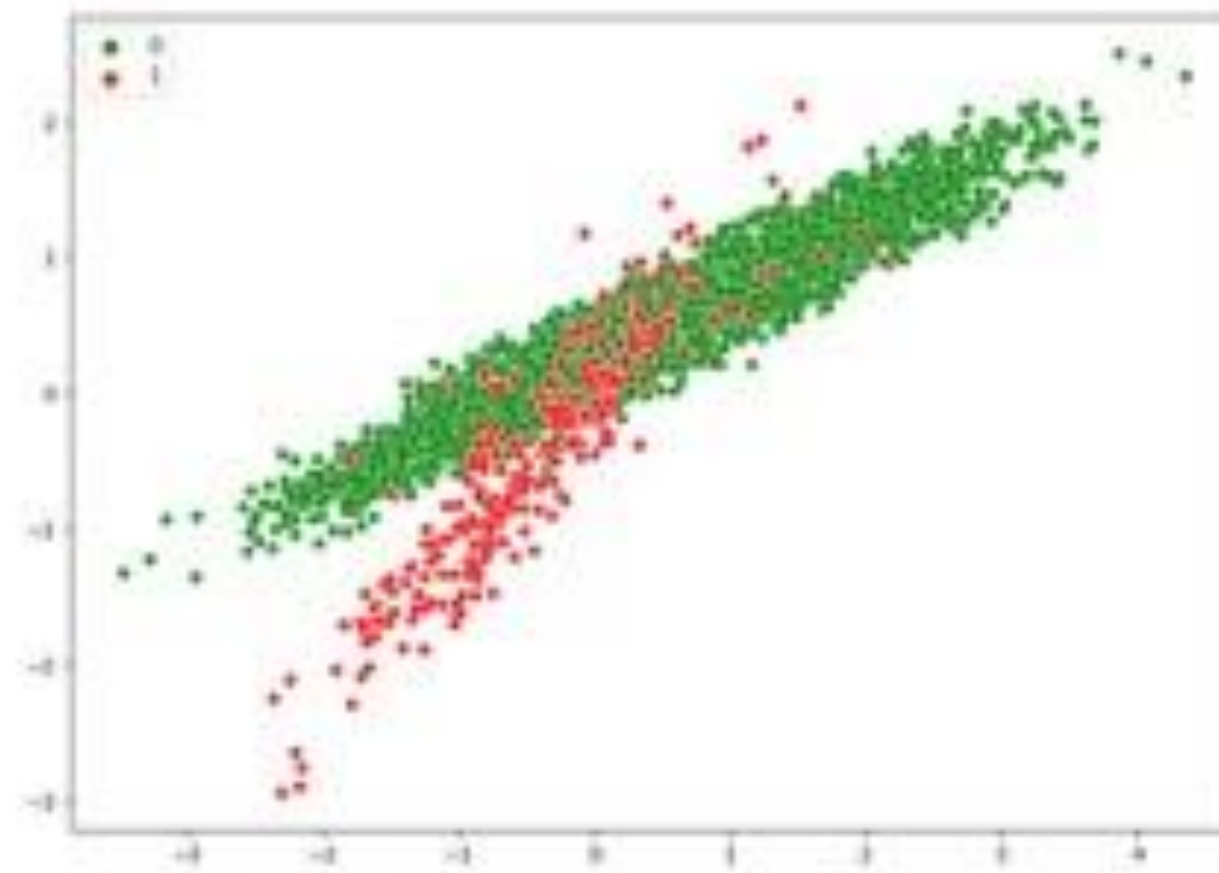
Υποδειγματικότητα και υπερδειγματοληψία δεδομένων

Under-sampling



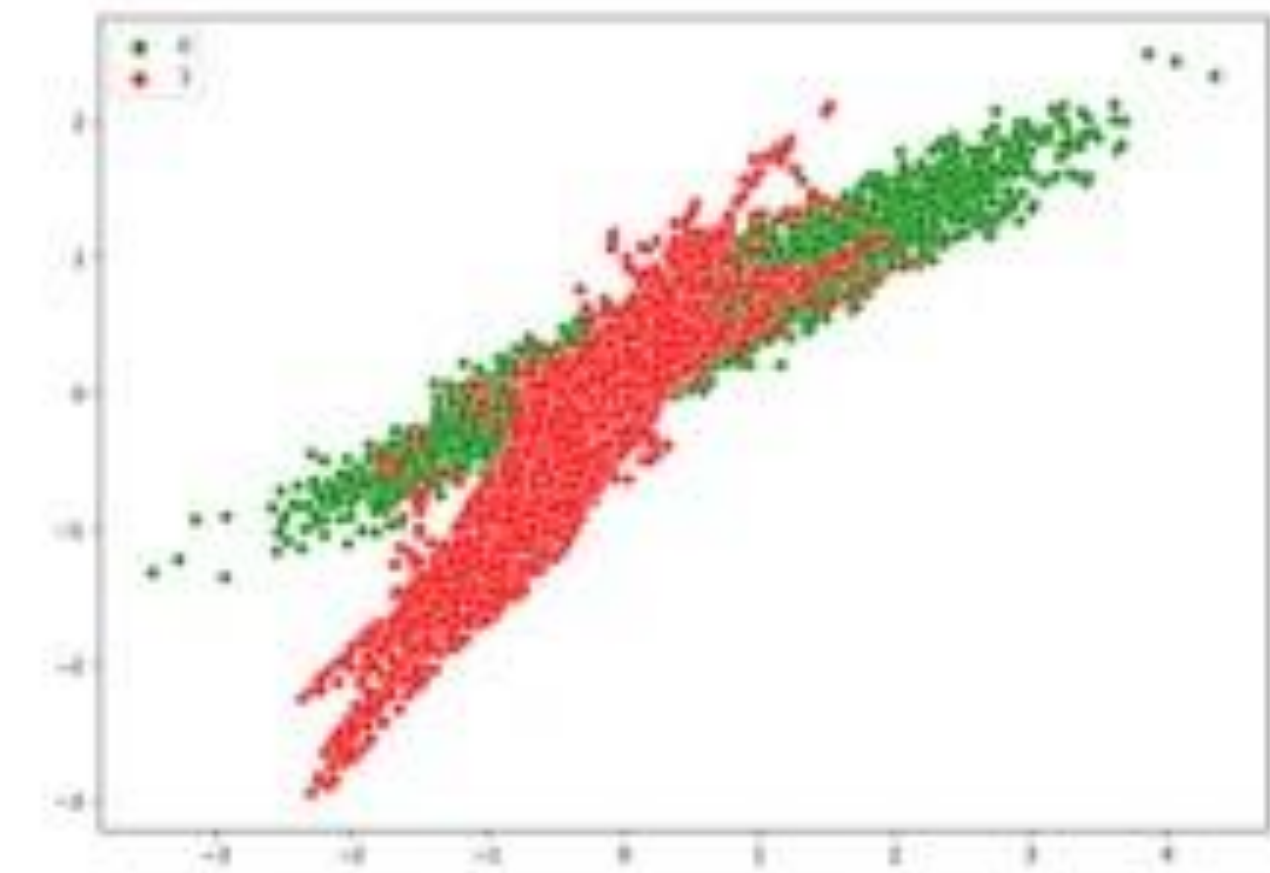
Samples of class 1 = 267
Samples of class 0 = 267

Over-sampling



Samples of class 1 = 4733
Samples of class 0 = 4733

SMOTE



Samples of class 1 = 4733
Samples of class 0 = 4733



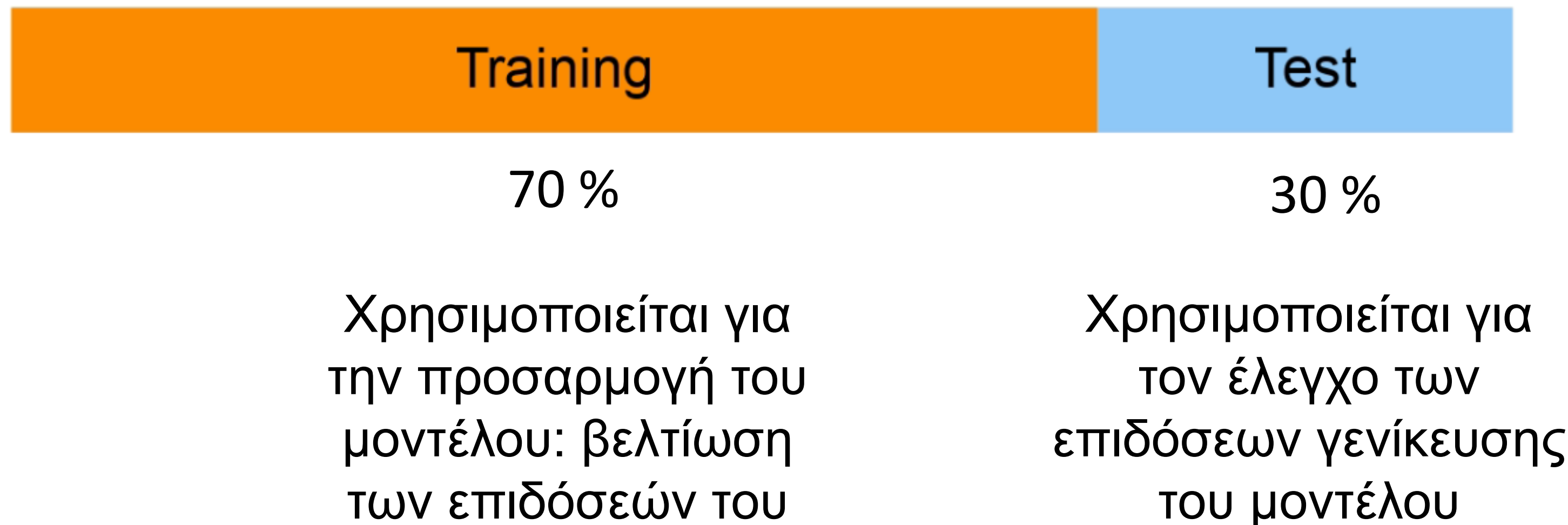


Διαχωρισμός συνόλου δεδομένων





Διαχωρισμός συνόλου δεδομένων



Γιατί το χρειαζόμαστε αυτό;

Το σύνολο δεδομένων είναι ένα δείγμα από την υποκείμενη κατανομή δεδομένων.

Αν χρησιμοποιήσουμε ολόκληρο το σύνολο, διατρέχουμε τον κίνδυνο «υπερπροσαρμογής» του μοντέλου στις αποχρώσεις του δείγματος.

Περισσότερα για αυτό σε μεταγενέστερες διαλέξεις...





Συνήθες πρόβλημα: Διαρροή δεδομένων

Κατά την προετοιμασία των δεδομένων διαρρέεται η γνώση του συνόλου διακρατούμενων (δοκιμών)

Χαρακτηριστικό παράδειγμα: η κλιμάκωση χαρακτηριστικών εξετάζει ολόκληρο το σύνολο δεδομένων αντί μόνο για το σύνολο εκπαίδευσης

Μπορεί να οδηγήσει σε λανθασμένες επιδόσεις του μοντέλου σε νέα δεδομένα

Πώς να το αποφύγετε: τοποθετήστε τη διαδικασία μετατροπής δεδομένων στο σύνολο εκπαίδευσης και αξιολογήστε μόνο το σύνολο δοκιμών

Το σύνολο δοκιμών δεν χρησιμοποιείται μόνο για την αξιολόγηση του εκπαιδευμένου μοντέλου, αλλά και για τη διαδικασία μετασχηματισμού δεδομένων.





Επόμενη Διάλεξη

Μέρος 2: Εποπτευόμενη Μάθηση

- Παλινδρόμηση
- Ταξινόμηση



MAI4CAREU

Master programmes in Artificial
Intelligence 4 Careers in Europe



Σας ευχαριστούμε

