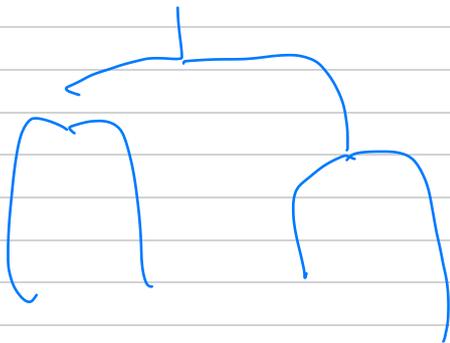


8-5-2025

Regression Trees

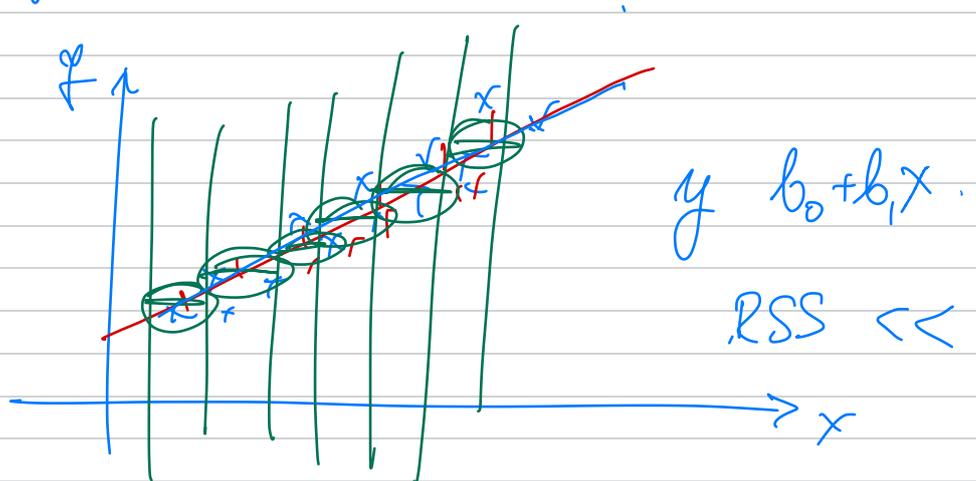
Decision Tree



Linear Regression

$$f(x) = b_0 + b_1 x_1 + \dots + b_p x_p$$

$$1) \quad f(x) = \sum_{i=1}^m c_i \cdot \mathbb{1}(x \in R_i)$$



2) Τα δέντρα είναι πιο εύκολη εγγραφή της προβλεπτικής ανάλυσης

3) Αν υπάρχουν πολλές μεταβλητές x_1, \dots, x_p

με regression: $y = b_0 + b_1 x_1 + \dots + b_k x_k \quad (k \leq p)$

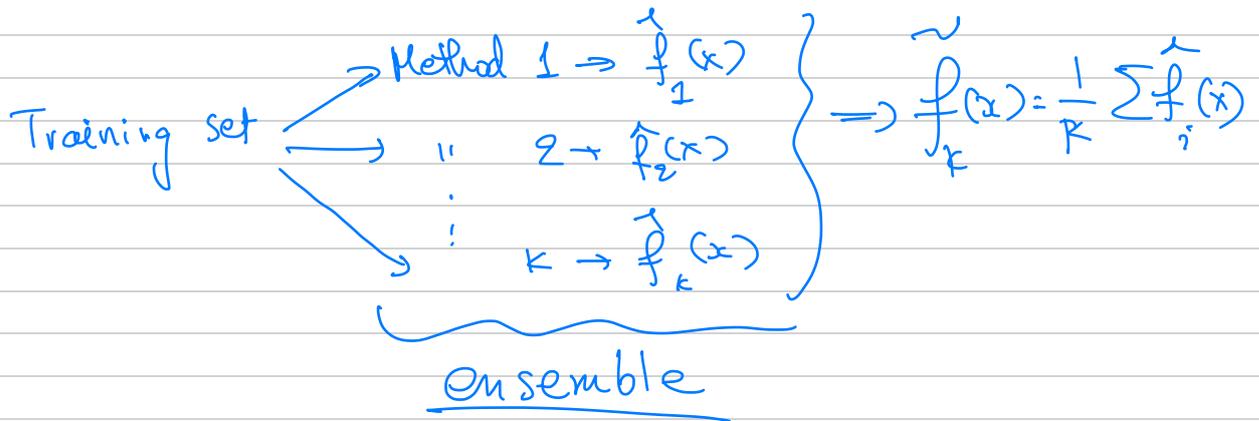
\exists τρόποι να αξιοποιηθεί η πληροφορία μιας μεταβλητής

Αυτό είναι όμοιο στο decision tree

4) Γενικά ένα "μεγάλο" decision tree έχει μεγάλο variance. //

ΕΠΕΡΙΧΕΙΡΗΣΕΙΣ

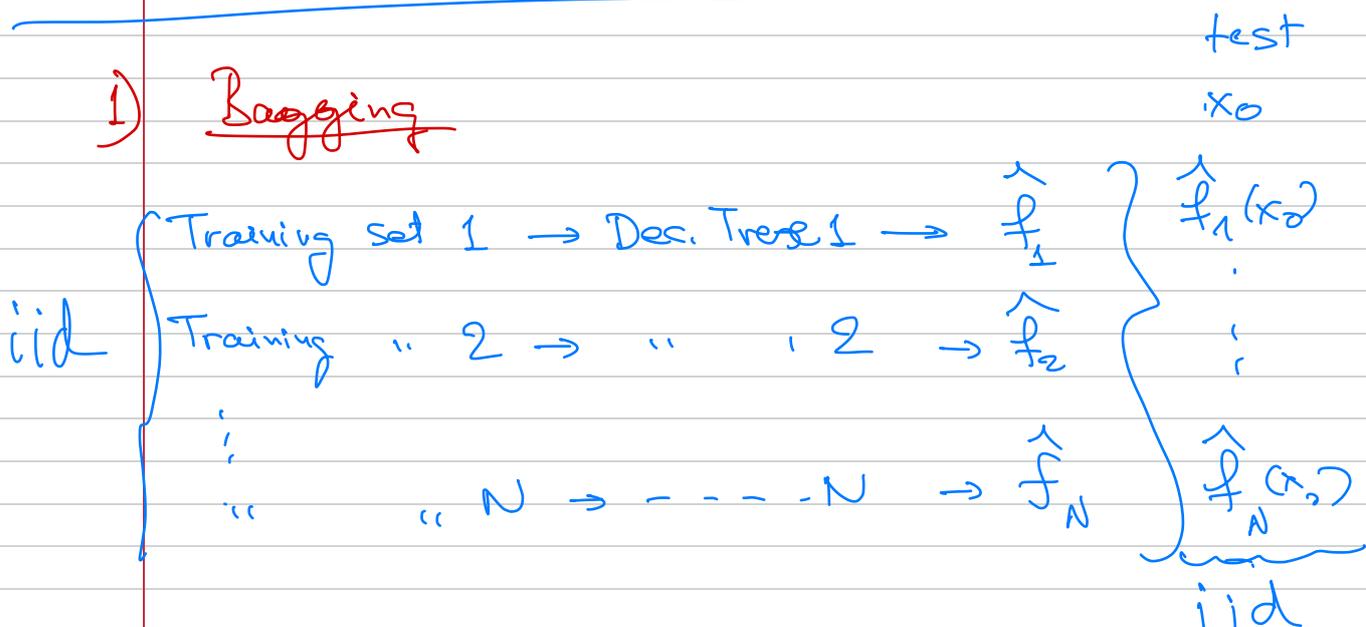
Τεχνική ιδέα : model ensemble.



Η εφαρμογή εδώ

- 1) Bagging
- 2) Random Forests
- 3) Boosting

1) Bagging



Εδώ $V(x_0) = \text{Var} \hat{f}_1(x_0)$

$$\text{Av } \bar{f}_N(x_0) = \frac{1}{N} \sum_{i=1}^N \hat{f}_i(x_0) \Rightarrow E[\bar{f}_N(x_0)] = E[f(x_0)]$$

$$\text{Var}\left[\bar{f}_N(x_0)\right] = \frac{1}{N} V(x_0)$$

Av έχουμε μόνο ένα training set
μεγεθός n ("training n ")

κ' λάβουμε N bootstrap δείγματα
από αυτό

- 1) κάθε bootstrap set iid από εμπειρική κατανομή των δειγμάτων
- 2) Η εμπειρική κατανομή των δειγμάτων \approx
πραγματική κατανομή των (X, Y)
όταν n μεγάλο

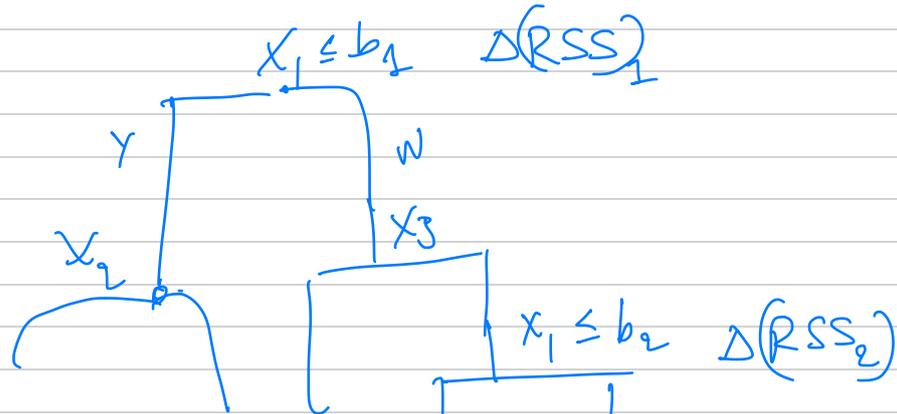
Bagging Algorithm

- 1) Bootstrap samples $b = 1, \dots, B$
από training set
- 2) Dec. tree \in κάθε $b \rightarrow \hat{f}_b(x)$
- 3)
$$\hat{f}(x) = \frac{1}{B} \sum_{b=1}^B \hat{f}_b(x)$$

Variance importance factors (VIF) Συντελεστής X_1

Σε κάθε δείγμα $b = 1, \dots, B$

υπολογίζουμε τα συνολικά μέσων των RSS που προκύπτει από αυτή τη μεταβλητή X_1 .



$(VIF)_{X_1}$ = συνολικά μέσων $\Delta(RSS)$ εφ' αυτών των X_1
(μέσος όρος σε $b = 1, \dots, B$)

2) Τυχαία Δείγμα

Εάν σε n X_i είναι από ομογενή
Τότε θα υπάρξει στα αποτελέσματα bootstrap samples
και τα $\hat{f}_1, \dots, \hat{f}_b$ θα έχουν
δυνατά ομοιότητα

Εάν X_1, X_2 i.i.d. αλλά $\text{Cov}(X_1, X_2) \neq 0$

$$\text{Var}\left(\frac{X_1 + X_2}{2}\right) \neq \frac{1}{2} \text{Var}(X_1)$$

$$\text{Var}(X_1 + X_2) = 2 \text{Var}(X_1) + 2 \text{Cov}(X_1, X_2)$$

$$\text{Var}\left(\frac{X_1 + X_2}{2}\right) = \frac{1}{2} \text{Var}(X_1) + \frac{1}{2} \text{Cov}(X_1, X_2)$$