

# Classement, classification supervisée

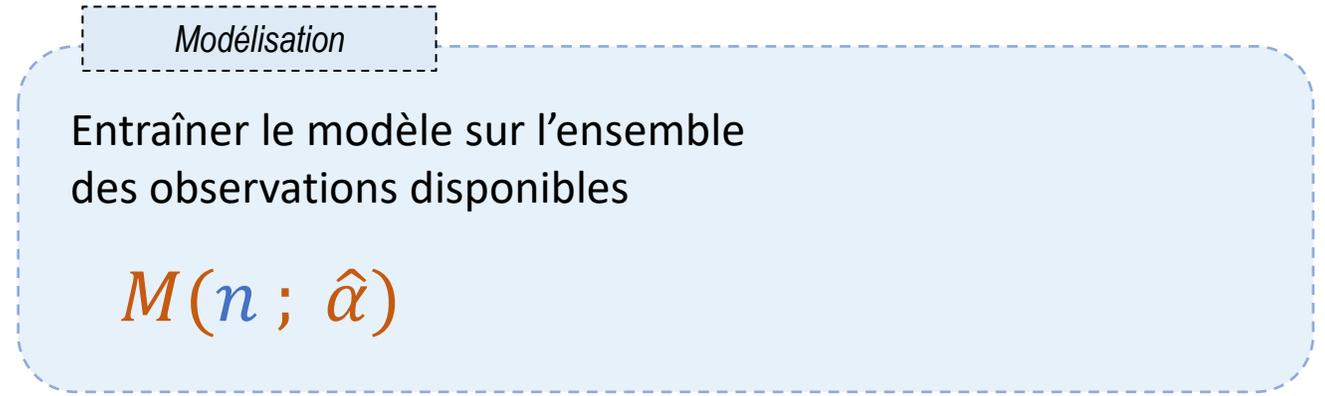
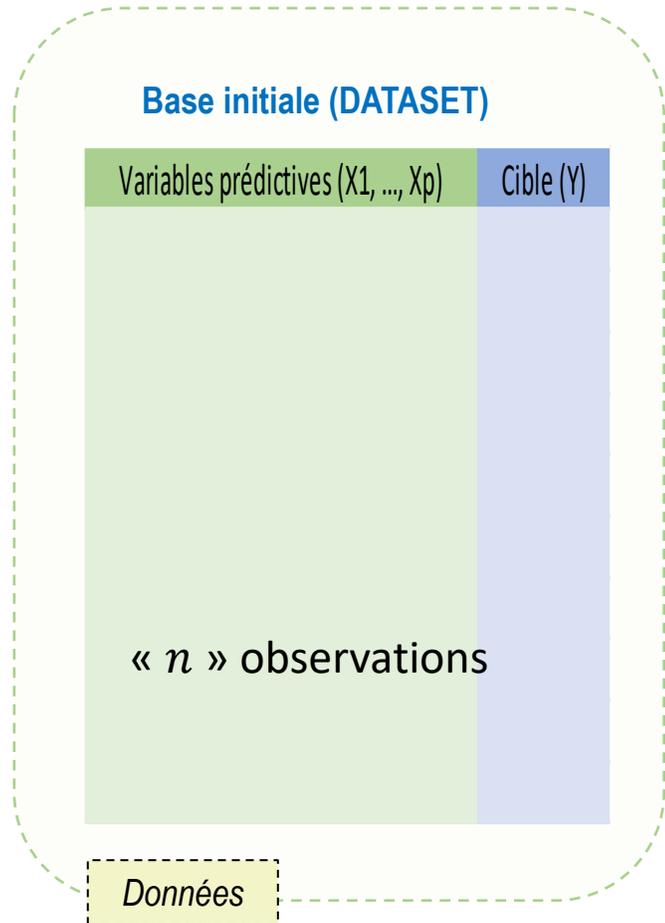
## Base initiale (DATASET)

Variables prédictives ( $X_1, \dots, X_p$ )	Cible ( $Y$ )
« $n$ » observations	

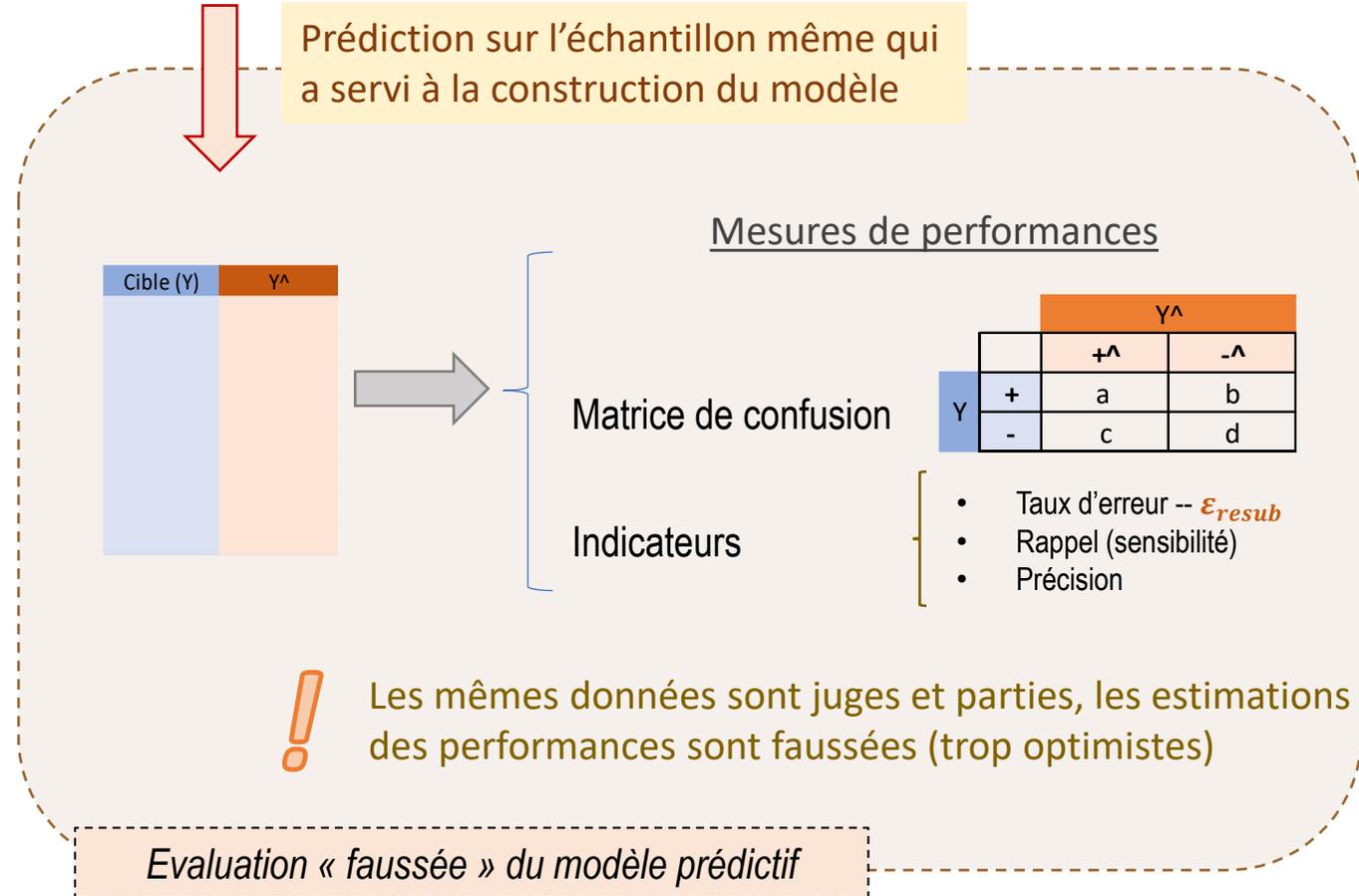
Comment construire un modèle prédictif  $M(n ; \alpha)$   
(ex. régression logistique, arbre de décision, SVM, etc.)

Comment évaluer les performances (ex. taux d'erreur)

# Evaluation en « resubstitution »



Prédiction sur l'échantillon même qui a servi à la construction du modèle



# Schéma « holdout » (apprentissage – test)

## Base initiale (DATASET)

Variables prédictives ( $X_1, \dots, X_p$ )	Cible (Y)
(0)	



## Base d'apprentissage (TRAIN SET)

Variables prédictives ( $X_1, \dots, X_p$ )	Cible (Y)
(1) $n_{train}$	

Variables prédictives ( $X_1, \dots, X_p$ )	Cible (Y)
(3) $n_{test}$	

## Base de test (TEST SET)

Modélisation (2)

$M(n_{train}; \hat{\alpha})$  (2)

Prédiction (4)

Application des règles prédictives sur l'échantillon test. Construction de la prédiction  $\hat{Y}$

Cible (Y)	$\hat{Y}$
(4)	

Mesure des performances (5)

Confrontation entre valeurs observées de Y et prédictions

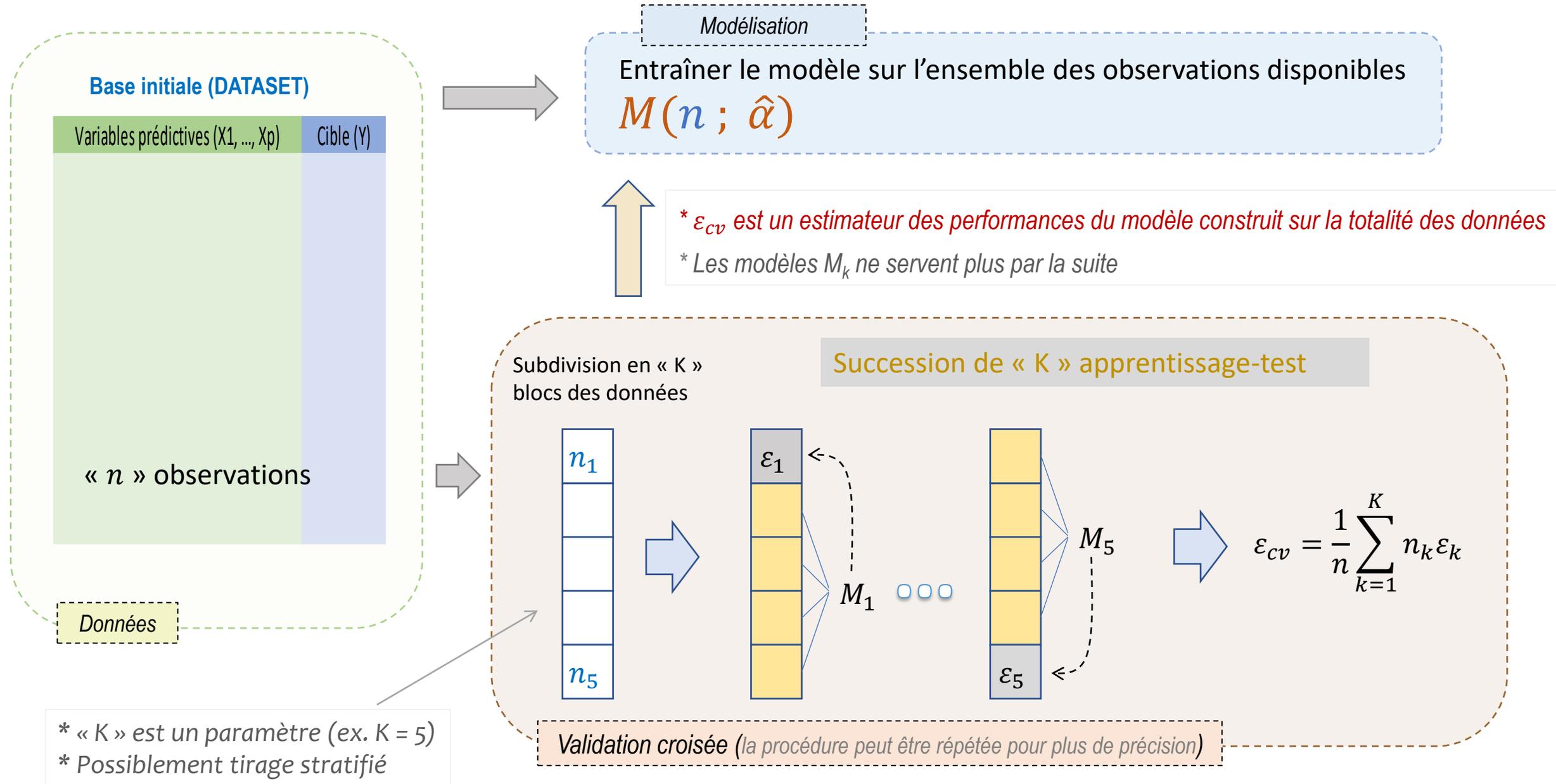
Matrice de confusion + Indicateurs, dont  $\epsilon_{test}$

Evaluation du modèle prédictif (4 et 5)

Subdivision des données (0) en échantillons d'apprentissage (1) et de test (3).

Possiblement tirage stratifié, surtout si classes déséquilibrées (cf. `train_test_split` de « scikit\_learn »)

# Validation croisée (K-Fold cross-validation)



# Leave-one-out (« K = n » cross-validation)

