

Précision

$$\text{precision} = \frac{|\{\text{relevant documents}\} \cap \{\text{retrieved documents}\}|}{|\{\text{retrieved documents}\}|}$$

Rappel :

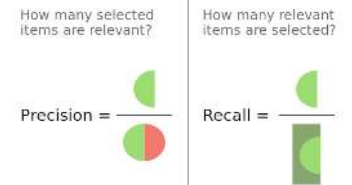
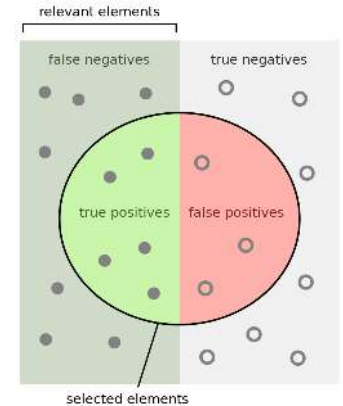
$$\text{recall} = \frac{|\{\text{relevant documents}\} \cap \{\text{retrieved documents}\}|}{|\{\text{relevant documents}\}|}$$

Matrice de confusion :

		classe réelle	
		chat	chien
classe prédite	chat	5	2
	chien	3	3

		classe réelle	
		Pos	Neg
classe prédite	Pos	TP	FP
	Neg	FN	TN

$$\text{Precision} = \frac{TP}{TP + FP} \quad \text{Recall} = \frac{TP}{TP + FN}$$



- **sensitivity, recall, hit rate, or true positive rate (TPR) :**

$$\text{TPR} = \frac{TP}{P} = \frac{TP}{TP + FN} = 1 - \text{FNR}$$

- **specificity, selectivity or true negative rate (TNR) :**

$$\text{TNR} = \frac{TN}{N} = \frac{TN}{TN + FP} = 1 - \text{FPR}$$

- **precision or positive predictive value (PPV) :**

$$\text{PPV} = \frac{TP}{TP + FP} = 1 - \text{FDR}$$

- **negative predictive value (NPV) :** $\text{NPV} = \frac{TN}{TN + FN} = 1 - \text{FOR}$

- **miss rate or false negative rate (FNR) :**

$$\text{FNR} = \frac{FN}{P} = \frac{FN}{FN + TP} = 1 - \text{TPR}$$

- **fall-out or false positive rate (FPR) :** $\text{FPR} = \frac{FP}{N} = \frac{FP}{FP + TN} = 1 - \text{TNR}$

- **false discovery rate (FDR) :** $\text{FDR} = \frac{FP}{FP + TP} = 1 - \text{PPV}$

- **false omission rate (FOR) :** $\text{FOR} = \frac{FN}{FN + TN} = 1 - \text{NPV}$

Sources : Fawcett (2006),[4] Powers (2011),[1] Ting (2011),[5] and CAWCR[6]

- **Threat score (TS) or Critical Success Index (CSI)** $\text{TS} = \frac{TP}{TP + FN + FP}$

- **accuracy (ACC) :** $\text{ACC} = \frac{TP + TN}{P + N} = \frac{TP + TN}{TP + TN + FP + FN}$

- **Balanced Accuracy** = $\frac{\text{TPR} + \text{TNR}}{2}$

- **F1 score** is the harmonic mean $\left(\frac{1 + \frac{1}{b}}{2}\right)^{-1}$ of precision and sensitivity :

$$F_1 = 2 \cdot \frac{\text{PPV} \cdot \text{TPR}}{\text{PPV} + \text{TPR}} = \frac{2TP}{2TP + FP + FN}$$

- **Matthews correlation coefficient [1975] (MCC, value $\in [-1, 1]$)** = Pearson's phi coefficient [1946], Yule phi coefficient [1912] :

$$\text{MCC} = \frac{TP \times TN - FP \times FN}{\sqrt{(TP + FP)(TP + FN)(TN + FP)(TN + FN)}}$$

- **Informedness or Bookmaker Informedness (BM) :** $\text{BM} = \text{TPR} + \text{TNR} - 1$

- **Markedness (MK) :** $\text{MK} = \text{PPV} + \text{NPV} - 1$

Sources : Fawcett (2006),[4] Powers (2011),[1] Ting (2011),[5] and CAWCR[6]

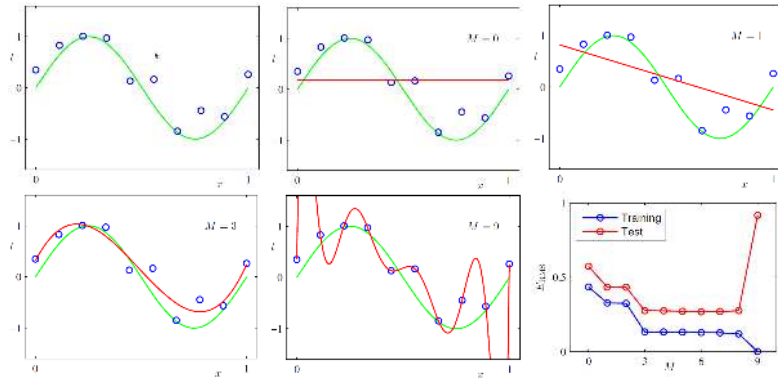




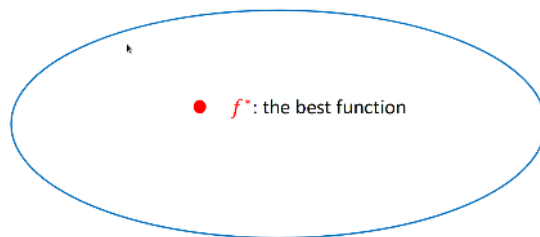




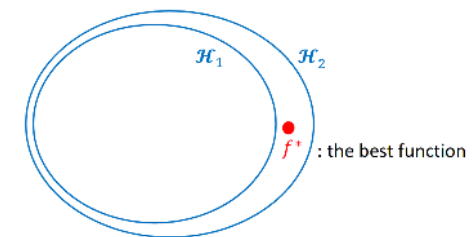
Figure from Machine Learning and Pattern Recognition, Bishop

- Très forte connaissance *a priori* : $\mathcal{H} = \{f^*\}$
⇒ Pas besoin de données!  f^* : the best function
- Très forte connaissance *a priori* : $\mathcal{H} = \{f^*, f_1\}$
⇒ Quelques points suffisent pour détecter f^*  f^* : the best function

- Très forte connaissance *a priori* : $\mathcal{H} = \{f^*\}$
⇒ Pas besoin de données!  f^* : the best function
- Très forte connaissance *a priori* : $\mathcal{H} = \{f^*, f_1\}$
⇒ Quelques points suffisent pour détecter f^*  f^* : the best function
- Ensemble de données gigantesque : infini
La classe des hypothèses \mathcal{H} peut être l'ensemble de toutes les fonctions



- Très forte connaissance *a priori* : $\mathcal{H} = \{f^*\}$
⇒ Pas besoin de données!  f^* : the best function
- Très forte connaissance *a priori* : $\mathcal{H} = \{f^*, f_1\}$
⇒ Quelques points suffisent pour détecter f^*  f^* : the best function
- Ensemble de données gigantesque : infini
La classe des hypothèses \mathcal{H} peut être l'ensemble de toutes les fonctions
- Scénarii classiques :
ensemble fini de données, \mathcal{H} de capacité moyenne, $f^* \notin \mathcal{H}$



- Supposons
 - on a appris un modèle (arbre de décision, SMV, NN...)
 - Les données d'entraînement sont parfaitement apprises
- Est-ce un bon modèle ?
- L'erreur d'entraînement n'est pas suffisante :
 - il faut confronter le modèle à de nouveaux exemples : **exemple de généralisation**
 - Un modèle peut classer parfaitement les données d'entraînement et ne pas savoir classer de nouveaux exemples :
 - car les exemples d'entraînement sont incomplets (échantillon d'une distribution)
 - car les exemples d'entraînement peuvent être bruitées (erreur d'étiquetage, erreur de mesure sur les caractéristiques...)
- Rappel : Formalisation de l'induction
 étant donné une **fonction objectif** l et un échantillon provenant d'une **distribution inconnue** D , le but est de calculer une **fonction** f dont l'erreur a une faible espérance sur D :

$$\mathbb{E}_{(x,y) \sim D} \{l(y, f(x))\} = \sum_{(x,y)} D(x,y) l(y, f(x))$$

- Considérons une hypothèse h et
 - taux d'erreur sur l'ensemble d'entraînement $erreur_{train}(h)$
 - taux d'erreur sur toutes les données $erreur_{true}(h)$
- On dit que h sur-apprend l'ensemble d'apprentissage si

$$erreur_{train}(h) < erreur_{true}(h)$$

évaluation du sur-apprentissage : $erreur_{true}(h) - erreur_{train}(h)$

- **Problème** : on ne connaît pas $erreur_{true}(h)$
Solution :
 - on met de côté un ensemble de test
 ce sont des exemples qui seront utilisés pour l'évaluation
 - on ne s'en sert pas pour l'entraînement du modèle
 - après l'entraînement, on calcule $erreur_{test}(h)$.
- **Underfitting** : un algo qui n'apprend pas suffisamment de la phase d'apprentissage
Overfitting : un algo qui prend trop en compte les particularités de l'ensemble d'apprentissage. l'hypothèse se généralise mal
- **ce qu'on cherche** : Un modèle qui ne sur-apprend ni ne sous-apprend pas

Le dilemme (ou compromis) biais variance : minimiser simultanément deux sources d'erreurs qui empêchent les algorithmes d'apprentissage supervisé de généraliser au-delà de leur échantillon d'apprentissage :

- **Le biais** : erreur provenant d'hypothèses erronées
 manque de relations pertinentes entre les données en entrée et les sorties prévues (sous-apprentissage).
- **La variance** : erreur due à la sensibilité aux petites fluctuations de l'échantillon d'entraînement. \Rightarrow peut entraîner un surapprentissage.

Décomposition biais-variance de l'erreur quadratique

- ens d'entraînement : x_1, \dots, x_n et y_i .
- $y_i = f(x_i) + \epsilon_i$, où le bruit, ϵ_i a une moyenne nulle et une variance σ^2 .
- Trouver une fonction \hat{f} qui se généralise à des points extérieurs à l'ensemble d'entraînement
- son erreur attendue sur un échantillon test x

$$E \left[(y - \hat{f}(x))^2 \right] = \text{Biais} [\hat{f}(x)]^2 + \text{Var} [\hat{f}(x)] + \sigma^2$$

où

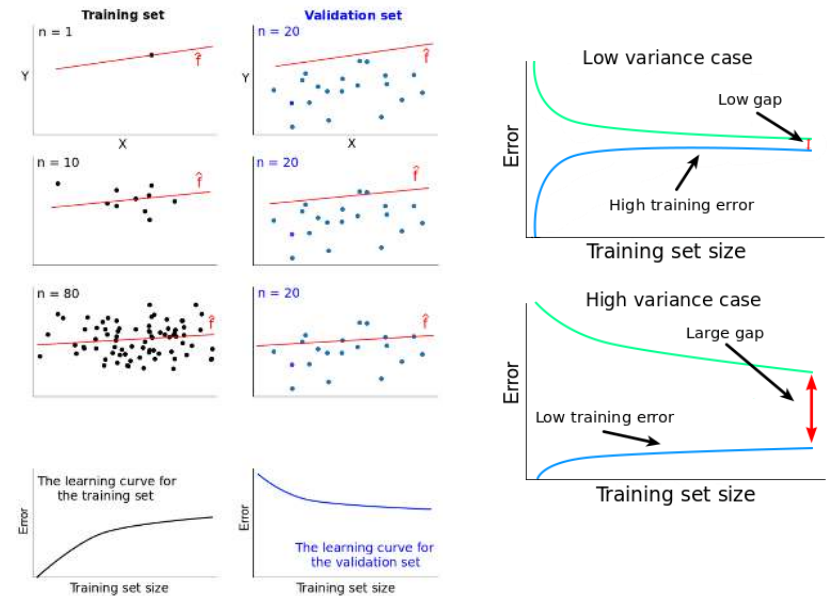
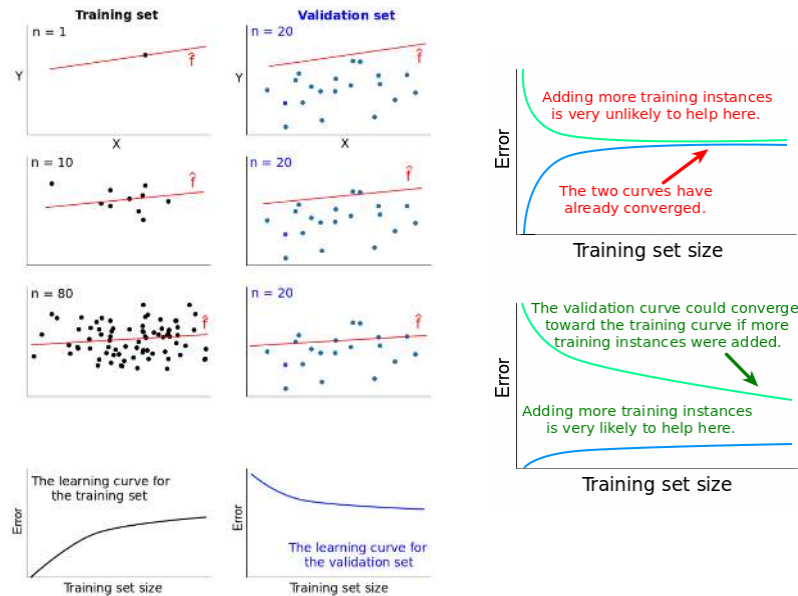
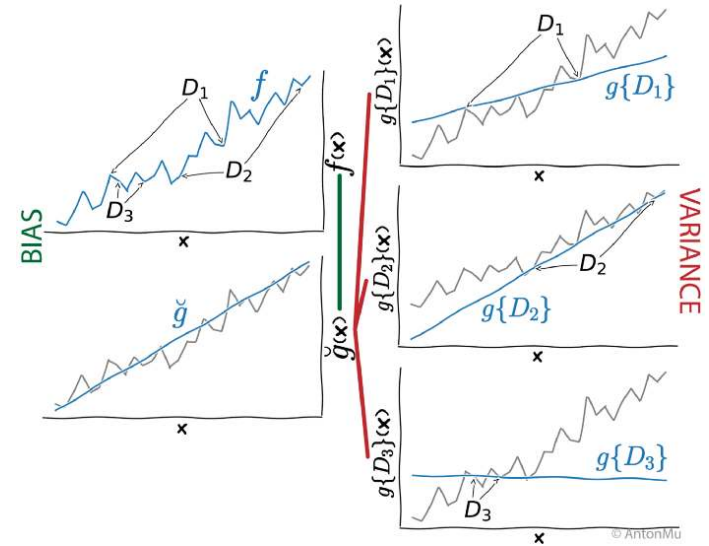
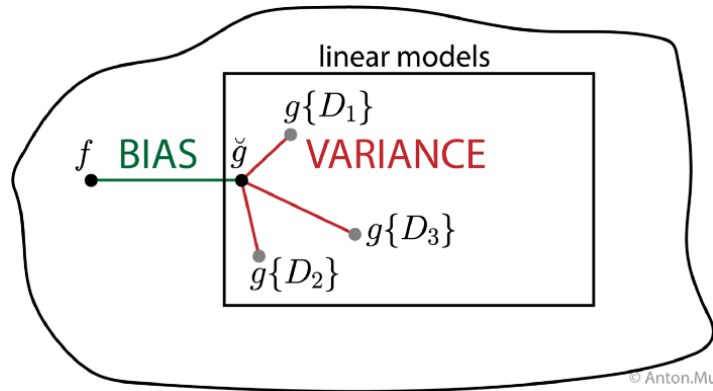
$$\text{Biais} [\hat{f}(x)] = E [\hat{f}(x) - f(x)]$$

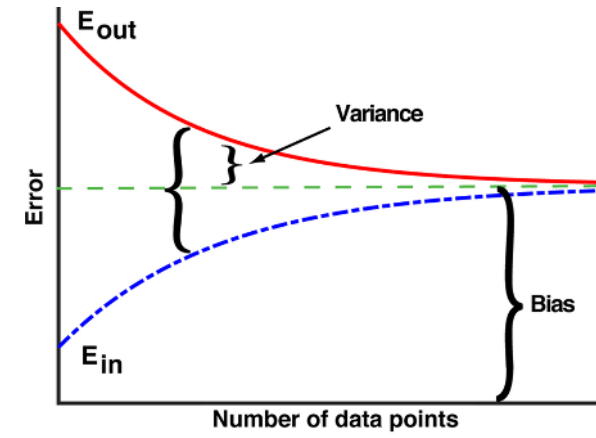
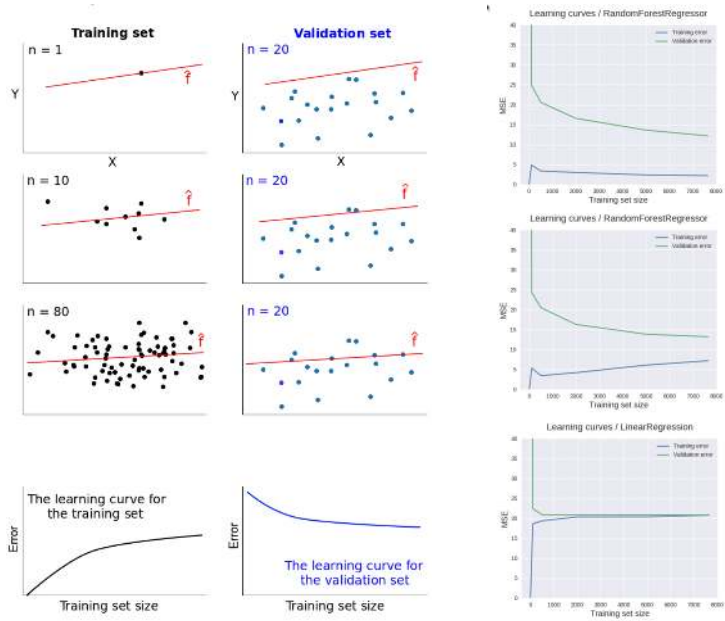
et

$$\text{Var} [\hat{f}(x)] = E \left[(\hat{f}(x) - E[\hat{f}(x)])^2 \right]$$

- $\text{Var}[X] = E[X^2] - (E[X])^2 \Leftrightarrow E[X^2] = \text{Var}[X] + (E[X])^2$.
- puisque f est déterministe : $E[f] = f$.
- comme $y = f + \epsilon$ et $E[\epsilon] = 0$, on a $E[y] = E[f + \epsilon] = E[f] = f$.
- puisque $\text{Var}[\epsilon] = \sigma^2$,
 $\text{Var}[y] = E[(y - E[y])^2] = E[(y - f)^2] = E[(f + \epsilon - f)^2]$
 $= E[\epsilon^2] = \text{Var}[\epsilon] + (E[\epsilon])^2 = \sigma^2$
- ainsi, puisque ϵ et \hat{f} sont indépendants, on peut écrire :

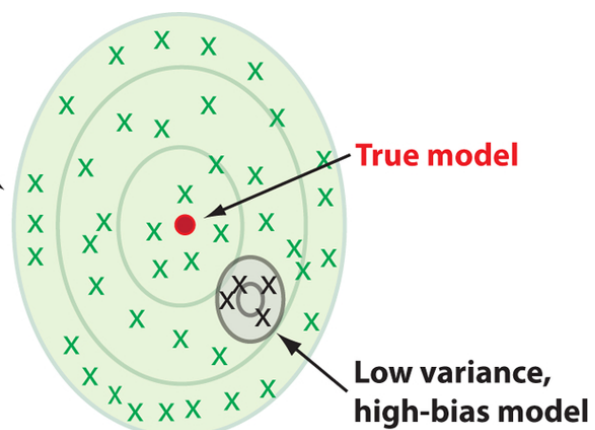
$$\begin{aligned} E[(y - \hat{f})^2] &= E[(f + \epsilon - \hat{f})^2] = E[(f + \epsilon - \hat{f} + E[\hat{f}] - E[\hat{f}])^2] \\ &= E[(f - E[\hat{f}])^2] + E[\epsilon^2] + E[(E[\hat{f}] - \hat{f})^2] \\ &\quad + 2E[(f - E[\hat{f}])\epsilon] + 2E[\epsilon(E[\hat{f}] - \hat{f})] + 2E[(E[\hat{f}] - \hat{f})(f - E[\hat{f}])] \\ &= \underbrace{(f - E[\hat{f}])^2}_{=0} + E[\epsilon^2] + E[(E[\hat{f}] - \hat{f})^2] \\ &\quad + \underbrace{2E[\epsilon(E[\hat{f}] - \hat{f})]}_{=0} + \underbrace{2E[(E[\hat{f}] - \hat{f})(f - E[\hat{f}])]}_{=0} \\ &= (f - E[\hat{f}])^2 + E[\epsilon^2] + E[(E[\hat{f}] - \hat{f})^2] \\ &= (f - E[\hat{f}])^2 + \text{Var}[y] + \text{Var}[\hat{f}] \\ &= \text{Bias}[\hat{f}]^2 + \text{Var}[y] + \text{Var}[\hat{f}] \\ &= \text{Bias}[\hat{f}]^2 + \sigma^2 + \text{Var}[\hat{f}] \end{aligned}$$





- Hypothesis : we cannot exactly learn the function $f^*(x)$
- As the number of data points gets large, the sampling noise decreases and the training data set becomes more representative of the true distribution
- in the infinite data limit, the in-sample and out-of-sample errors must approach the same value, which is called the bias of our model.
- The bias : the best our model could do if we had an infinite amount of training data to beat down sampling noise.

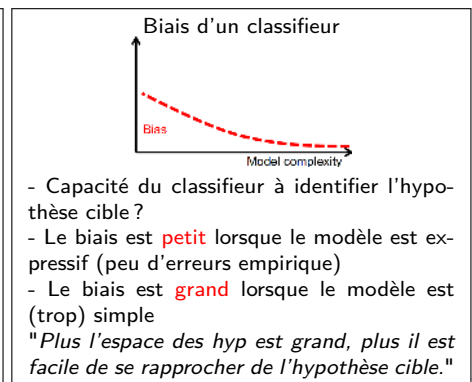
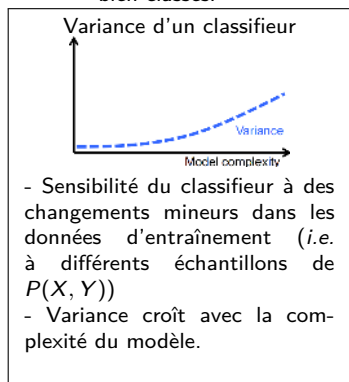
High variance, low-bias model

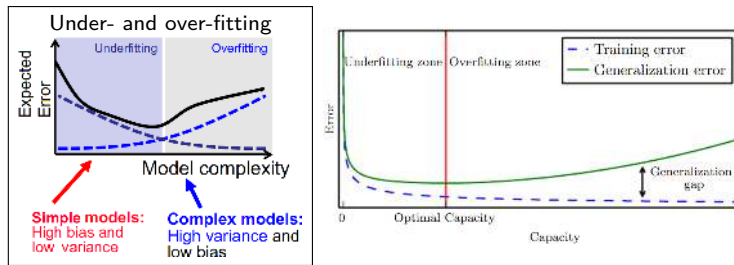
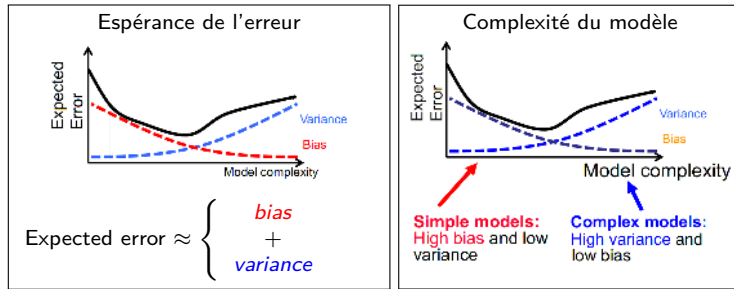


- Hypothesis : we cannot exactly learn the function $f^*(x)$
- As the number of data points gets large, the sampling noise decreases and the training data set becomes more representative of the true distribution
- in the infinite data limit, the in-sample and out-of-sample errors must approach the same value, which is called the bias of our model.
- The bias : the best our model could do if we had an infinite amount of training data to beat down sampling noise.

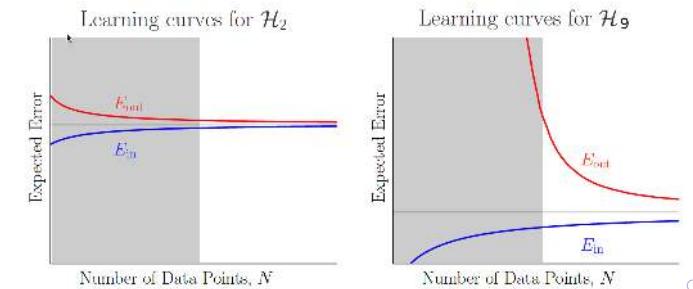
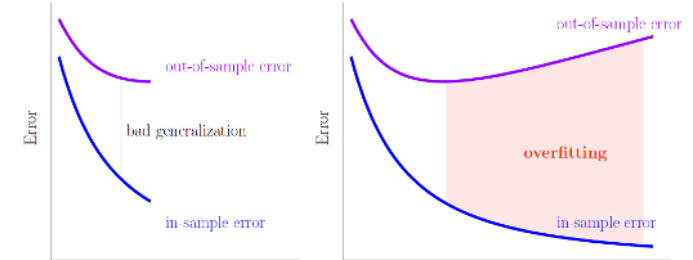


- Un modèle **sur-apprend** l'ensemble d'apprentissage quand la **justesse (accuracy)** sur l'ensemble d'apprentissage **augmente** alors que la **justesse sur des nouvelles données diminue**.
- **Erreur empirique** (sur un ensemble de données) : % d'items qui ne sont pas bien classés.





Le sur-apprentissage n'est pas qu'une mauvaise erreur de généralisation



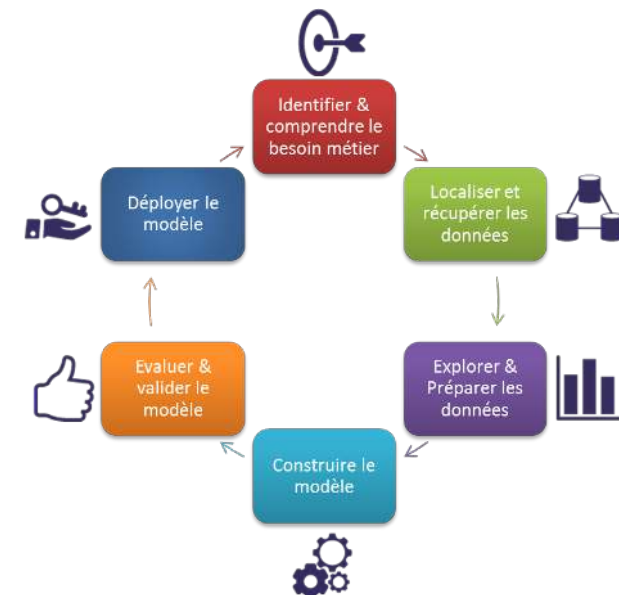
Ensembles d'apprentissage, de dev., de test

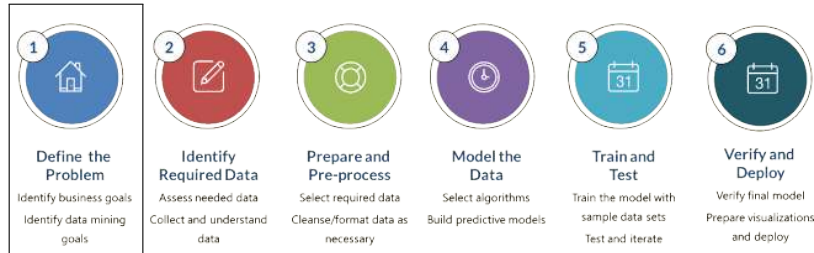
En pratique, on répartit les exemples en 3 ensembles distincts :

- Ensemble d'entraînement :
 - utilisé pour entraîner les modèles
 - e.g. les nœuds et branches d'un arbre de décision
- Ensemble de développement (tuning set/validation set/ held out data) :
 - Utilisé pour apprendre les **hyperparamètres** :
 - paramètres contrôlant les autres paramètres
 - e.g. max_depth d'un arbre de décision
- Ensemble Test :
 - utilisé pour évaluer la performance sur des données pas encore vues.

⇒ Never ever touch your test data !

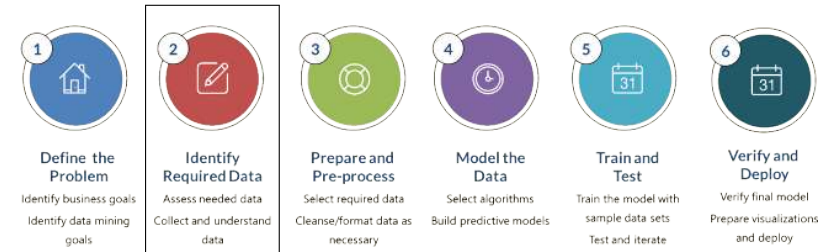
Déroulement d'un projet



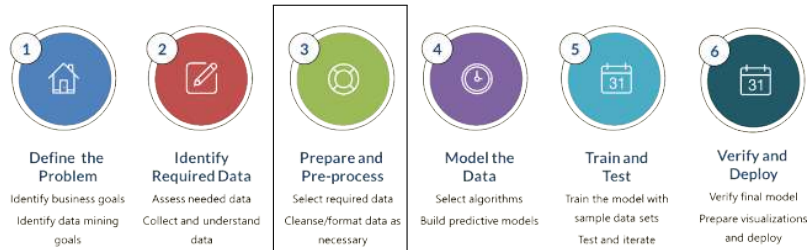


Quel est le problème que l'on essaie de résoudre ?

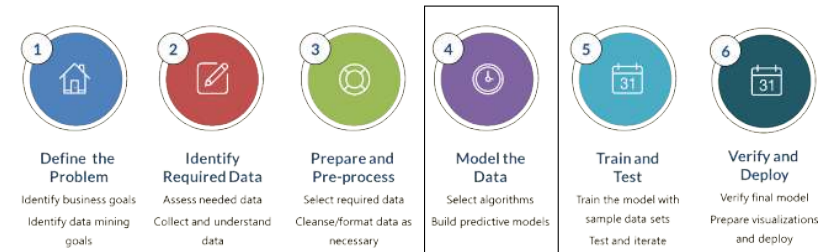
- Se poser les bonnes questions.
Quels types de relations faut-il trouver ?
- Définir ensuite **précisément** l'objectif **mesurable** et **quantifiable** du projet
- L'expression du besoin :
 - critère de réussite + condition d'arrêt
 - Compréhension commune de ce qui peut être espéré
 - Connaissances *a priori* sur les données ? Saisonnalité ? distribution *a priori* ?
 - éléments logistique du projet (budget, plan, deadline...)



- Les données peuvent être dispersées dans une entreprise et stockées dans différents formats, ou contenir des incohérences telles que des entrées incorrectes ou manquantes. Par exemple, les données peuvent indiquer qu'un client a acheté un produit avant que le produit ne soit offert sur le marché, ou que le client achète régulièrement dans un magasin situé à 3 000 km de son domicile.
 - Nettoyage : traiter les données incorrectes / manquantes ; rechercher des corrélations cachées ; identifier les sources de données les plus précises ; déterminer les colonnes les plus appropriées pour une analyse. Par exemple, utilisation de la date d'expédition ou de commande ?
- ⇒ identifier tous ces problèmes et déterminer comment les résoudre. Comme on ne peut passer en revue tous les enregistrements, on peut essayer d'abord de filtrer les données de manière automatique.

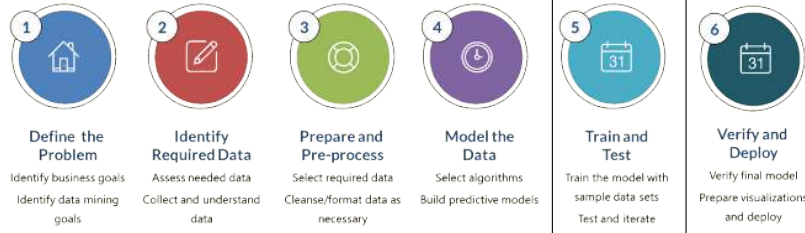


- comprendre les données afin de prendre les décisions appropriées lors de la création des modèles.
calcul des valeurs minimales et maximales, de moyennes et des écarts types, examen de la distribution des données.
- Est-ce que le jeu de données contient des données erronées ? Quelle stratégie pour résoudre les problèmes ou acquérir une compréhension plus approfondie des comportements typiques ?



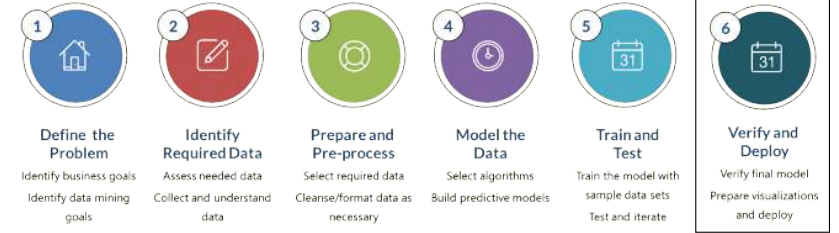
Trouver des schémas dans les données qui mènent vers les solutions

- des "va-et-vient" avec l'étape précédente peuvent être nécessaires
- Tâches de modélisation :
 - Classification
 - Prévion
 - Classement
 - Segmentation
 - trouver des associations/corrélations
- Exemple de méthodes de classification : arbre de décision, NN...



Est-ce que le modèle résout mon problème

- est-il précis selon les besoins ? Est-ce qu'il généralise bien ?
 - Est-il meilleur que le "résultat évident" ou modèle de base ?
 - Les résultats ont-ils du sens dans le contexte du problème ?
- ⇒ s'il y a plusieurs étapes, il faut faire une évaluation de chacune des tâches.



une fois que le modèle ait été construit :

- Utilisation du modèle pour faire des prédictions/décisions. (mise en place d'un serveur d'interrogation (DMX queries) ?)
- Création de requêtes sur les statistiques, les règles ou formules au sein d'un modèle (DMX).
- Intégration de la fonctionnalité d'exploration de données directement dans une application
- Mise à jour les modèles. Toute mise à jour nécessite le retraitement des modèles.
- Mise à jour automatique des modèles.