Apprentissage statistique : modélisation décisionnelle et apprentissage profond (RCP209)

Introduction à l'apprentissage supervisé

Nicolas Audebert
nicolas.audebert@lecnam.net
http://cedric.cnam.fr/vertigo/Cours/m12/

Département Informatique Conservatoire National des Arts & Métiers, Paris, France

21 septembre 2023

#### Plan du cours

- 1 Objectifs et contenu de l'enseignement
  - Organisation de l'enseignemen
  - Modélisation décisionnelle
  - Types de problèmes de décision
  - Modélisation à partir de données
  - Modélisation à partir de données : un cadre plus préd
  - Etapes generales
  - Choix d'une fonction de perte
  - Choix des familles paramétriques
  - Estimation du modele
  - Comment mesurer la
  - Évaluation de mod
  - Courbes ROC
    - Sélection de modèle
    - Grid search pour le choix des hyperparamètre
    - Randomized parameter optimization

## "Prediction can be very difficult, especially about the future."

(Niels Bohr, dans Teaching and Learning Elementary Social Studies, Arthur K. Ellis, 1970, p. 431.)

## → Modélisation décisionnelle (ou prédictive) à partir de données

- Données : observations caractérisées par les valeurs prises par un ensemble de variables
- Modélisation à partir de données : construction (en grande partie automatique) de modèle(s) qui « explique(nt) » les données
- Modélisation décisionnelle : capacité à prédire, pour chaque nouvelle observation, la valeur (inconnue) d'une variable expliquée à partir des valeurs (connues) de variables explicatives

#### Objectifs applicatifs de la modélisation décisionnelle

- Reconnaissance des formes (pattern recognition): (sens strict) identifier à quelle catégorie appartient une « forme » décrite par des données brutes
- Fouille de données (data mining) : (sens strict) recherche de régularités ou de relations inconnues a priori dans de (grands) volumes de données

## "Prediction can be very difficult, especially about the future."

- → Modélisation décisionnelle (ou prédictive) à partir de données
  - Données : observations caractérisées par les valeurs prises par un ensemble de variables
  - Modélisation à partir de données : construction (en grande partie automatique) de modèle(s) qui « explique(nt) » les données
  - Modélisation décisionnelle : capacité à prédire, pour chaque nouvelle observation, la valeur (inconnue) d'une variable expliquée à partir des valeurs (connues) de variables explicatives
  - Objectifs applicatifs de la modélisation décisionnelle
    - Reconnaissance des formes (pattern recognition): (sens strict) identifier à quelle catégorie appartient une «forme» décrite par des données brutes
    - Fouille de données (data mining) : (sens strict) recherche de régularités ou de relations inconnues a priori dans de (grands) volumes de données

## "Prediction can be very difficult, especially about the future."

- → Modélisation décisionnelle (ou prédictive) à partir de données
  - Données : observations caractérisées par les valeurs prises par un ensemble de variables
  - Modélisation à partir de données : construction (en grande partie automatique) de modèle(s) qui « explique(nt) » les données
  - Modelisation decisionnelle : capacité à predire, pour chaque nouvelle observation, la valeur (inconnue) d'une variable expliquée à partir des valeurs (connues) de variables explicatives
  - Objectifs applicatifs de la modélisation décisionnelle
    - Reconnaissance des formes (pattern recognition) : (sens strict) identifier à quelle catégorie appartient une « forme » décrite par des données brutes
    - Fouille de données (data mining) : (sens strict) recherche de régularités ou de relations inconnues a priori dans de (grands) volumes de données

## "Prediction can be very difficult, especially about the future."

- → Modélisation décisionnelle (ou prédictive) à partir de données
  - Données : observations caractérisées par les valeurs prises par un ensemble de variables
  - Modélisation à partir de données : construction (en grande partie automatique) de modèle(s) qui « explique(nt) » les données
  - Modélisation décisionnelle : capacité à prédire, pour chaque nouvelle observation, la valeur (inconnue) d'une variable expliquée à partir des valeurs (connues) de variables explicatives
  - Objectifs applicatifs de la modélisation décisionnelle
    - Reconnaissance des formes (pattern recognition): (sens strict) identifier à quelle catégorie appartient une «forme» décrite par des données brutes
    - Fouille de données (data mining) : (sens strict) recherche de régularités ou de relations inconques a priori dans de (grande) volumes de données

## "Prediction can be very difficult, especially about the future."

- → Modélisation décisionnelle (ou prédictive) à partir de données
  - Données : observations caractérisées par les valeurs prises par un ensemble de variables
  - Modélisation à partir de données : construction (en grande partie automatique) de modèle(s) qui «explique(nt) » les données
  - Modélisation décisionnelle : capacité à prédire, pour chaque nouvelle observation, la valeur (inconnue) d'une variable expliquée à partir des valeurs (connues) de variables explicatives
  - Objectifs applicatifs de la modélisation décisionnelle
    - Reconnaissance des formes (pattern recognition) : (sens strict) identifier à quelle catégorie appartient une «forme» décrite par des données brutes
    - Fouille de données (data mining) : (sens strict) recherche de régularités ou de relations inconnues a priori dans de (grands) volumes de données

"Prediction can be very difficult, especially about the future."

- → Modélisation décisionnelle (ou prédictive) à partir de données
  - Données : observations caractérisées par les valeurs prises par un ensemble de variables
  - Modélisation à partir de données : construction (en grande partie automatique) de modèle(s) qui «explique(nt) » les données
  - Modélisation décisionnelle : capacité à prédire, pour chaque nouvelle observation, la valeur (inconnue) d'une variable expliquée à partir des valeurs (connues) de variables explicatives
  - Objectifs applicatifs de la modélisation décisionnelle
    - Reconnaissance des formes (pattern recognition): (sens strict) identifier à quelle catégorie appartient une « forme » décrite par des données brutes
    - Fouille de données (data mining) : (sens strict) recherche de régularités ou de relations inconnues a priori dans de (grands) volumes de données

### "Prediction can be very difficult, especially about the future."

- → Modélisation décisionnelle (ou prédictive) à partir de données
  - Données : observations caractérisées par les valeurs prises par un ensemble de variables
  - Modélisation à partir de données : construction (en grande partie automatique) de modèle(s) qui «explique(nt) » les données
  - Modélisation décisionnelle : capacité à prédire, pour chaque nouvelle observation, la valeur (inconnue) d'une variable expliquée à partir des valeurs (connues) de variables explicatives
  - Objectifs applicatifs de la modélisation décisionnelle
    - Reconnaissance des formes (pattern recognition): (sens strict) identifier à quelle catégorie appartient une « forme » décrite par des données brutes
    - Fouille de données (data mining): (sens strict) recherche de régularités ou de relations inconnues a priori dans de (grands) volumes de données

## Problématiques abordées

- Comprendre la nature des problèmes de modélisation à partir de données
- Maîtriser la méthodologie générale de construction, évaluation et sélection de modèles décisionnels
- Maîtriser plusieurs outils de modélisation décisionnelle actuels : forêts d'arbres de décision, machines à vecteurs support (SVM), réseaux de neurones profonds
- Capacité à mettre œuvre des méthodes de modélisation décisionnelle à partir de données

- En mathématiques : connaissances de base en algèbre linéaire, probabilités, analyse
- En informatique : connaissances de base en programmation avec Python
  - bases de NumPy
  - RCP208 : Apprentissage statistique I fortement conseillé mais pas indispensable

## Problématiques abordées

- Comprendre la nature des problèmes de modélisation à partir de données
- Maîtriser la méthodologie générale de construction, évaluation et sélection de modèles décisionnels
- Maîtriser plusieurs outils de modélisation décisionnelle actuels : forêts d'arbres de décision, machines à vecteurs support (SVM), réseaux de neurones profonds
- Capacité à mettre œuvre des méthodes de modélisation décisionnelle à partir de données

- En mathématiques : connaissances de base en algèbre linéaire, probabilités, analyse
- En informatique : connaissances de base en programmation avec Python
  - bases de NumPy
  - RCP208 : Apprentissage statistique I fortement conseillé mais pas indispensable

## Problématiques abordées

- Comprendre la nature des problèmes de modélisation à partir de données
- Maîtriser la méthodologie générale de construction, évaluation et sélection de modèles décisionnels
- Maîtriser plusieurs outils de modélisation décisionnelle actuels : forêts d'arbres de décision, machines à vecteurs support (SVM), réseaux de neurones profonds
- Capacité à mettre œuvre des méthodes de modélisation décisionnelle à partir de données

- En mathématiques : connaissances de base en algèbre linéaire, probabilités, analyse
- En informatique : connaissances de base en programmation avec Python
  - bases de NumPy
  - RCP208 : Apprentissage statistique I fortement conseillé mais pas indispensable

### Problématiques abordées

- Comprendre la nature des problèmes de modélisation à partir de données
- Maîtriser la méthodologie générale de construction, évaluation et sélection de modèles décisionnels
- Maîtriser plusieurs outils de modélisation décisionnelle actuels : forêts d'arbres de décision, machines à vecteurs support (SVM), réseaux de neurones profonds
- Capacité à mettre œuvre des méthodes de modélisation décisionnelle à partir de données

- En mathématiques : connaissances de base en algèbre linéaire, probabilités, analyse
- En informatique : connaissances de base en programmation avec Python
  - bases de NumPy
  - RCP208 : Apprentissage statistique I fortement conseillé mais pas indispensable

#### Contenu détaillé

- Apprentissage supervisé : classification, régression; généralisation. Évaluation et sélection de modèles : validation croisée, grid search. (2 séances, Nicolas Audebert)
- Arbres de décision et forêts d'arbres de décision. (2 séances, Marin Ferecatu)
- SVM : maximisation de la marge, astuce des noyaux, classement (discrimination), régression, estimation du support d'une distribution, ingénierie des noyaux. (3 séances, Marin Ferecatu)
- Apprentissage profond (deep learning): réseaux convolutifs profonds, réseaux de neurones récurrents. (8 séances, Nicolas Audebert)
- = 15 séances cours + travaux pratiques (TP, Wafa Aissa)
  - ⇒ environ 60 heures de travail (6 ECTS)

#### Contenu détaillé

- Apprentissage supervisé : classification, régression; généralisation. Évaluation et sélection de modèles : validation croisée, grid search. (2 séances, Nicolas Audebert)
- Arbres de décision et forêts d'arbres de décision. (2 séances, Marin Ferecatu)
- SVM : maximisation de la marge, astuce des noyaux, classement (discrimination), régression, estimation du support d'une distribution, ingénierie des noyaux. (3 séances, Marin Ferecatu)
- Apprentissage profond (deep learning): réseaux convolutifs profonds, réseaux de neurones récurrents. (8 séances, Nicolas Audebert)
- = 15 séances cours + travaux pratiques (TP, Wafa Aissa)
  - ⇒ environ 60 heures de travail (6 ECTS)

## Travaux pratiques

- Mise en œuvre de la méthodologie de construction, évaluation et sélection de modèles décisionnels
- Mise en œuvre d'outils de modélisation actuels :
  - Forêts d'arbres de décision
  - Machines à vecteurs support (SVM)
  - Réseaux de neurones profonds
- Emploi de Scikit-learn (http://scikit-learn.org), outil libre et ouvert, en Python, déjà employé dans RCP208 (http://cedric.cnam.fr/vertigo/Cours/ml/); emploi de Keras (en Python aussi) pour l'apprentissage profond





# Quelques références bibliographiques

- C.-A. Azencott, Introduction au Machine Learning, Éditions Eyrolles, 2018.
- A. Géron, Machine Learning avec scikit-learn; Deep Learning avec Keras, Éditions Dunod, 2021.
- B. Schölkopf, A. Smola. Learning with Kernels. MIT Press, 2002. [4]
- I. Goodfellow, Y. Bengio, A. Courville. *Deep Learning*. MIT Press, 2016. http://www.deeplearningbook.org. [3]
- A. Amor, L. Estève, O. Grisel, G. Lemaître, G. Varoquaux, T. Schmitt, MOOC Machine learning with scikit-learn, France Université Numérique, https: //www.fun-mooc.fr/fr/cours/machine-learning-python-scikit-learn/.

D'autres références vous seront suggérées dans les différents chapitres du cours pour approfondir spécifiquement certaines parties.

## Plan du cours

- 1 Objectifs et contenu de l'enseigneme
  - Organisation de l'enseignement
  - Wodelisation decisionnell
    - Types de problèmes de décision
    - Modélisation à partir de données
    - Modélisation à partir de données : un
    - Etapes générales
    - Choix d'une fonction de perte
    - Choix des familles paramétriques
    - Estimation du modele
    - Comment mesurer la
    - Évaluation de mod
  - Courbes ROC
    - Sélection de modèle
    - Grid search pour le choix des hyperparamètre
    - Randomized parameter optimization

## Organisation

- Semestre 1 : formation à distance (FOD)
- Supports en accès ouvert (et pouvant évoluer à tout moment) : http://cedric.cnam.fr/vertigo/Cours/m12/
  - lacktriangle Cours : transparents (PDF) + explications en HTML ou vidéo
  - TP : contenu détaillé en HTML + notebooks Jupyter
  - Forum : Moodle (https://lecnam.net)
- Responsable du cours : Nicolas Audebert
- Enseignants : Nicolas Audebert, Marin Ferecatu, Clément Rambour, Arnaud Breloy



Nicolas Audebert



Marin Ferecatu



Clément Rambour



Arnaud Breloy

### Évaluation

- $\blacksquare$  S1 (FOD) : examen écrit : session 1 fin janvier/début février et session 2 en avril
  - Planification des examens : http://www.cnam-paris.fr/suivre-ma-scolarite/ rubrique Examens
- Mini-projet : analyse des données et construction de modèle(s) décisionnel(s) pour un problème proposé par vous ou par les enseignants
  - Choix à faire sur Moodle au mois de novembre
  - À rendre en même temps que l'examen de session 1 ou de session 2, au choix
- Note finale = moyenne non pondérée entre la note d'examen et la note de projet

## Plan du cours

- 1 Objectifs et contenu de l'enseignemer
- 2 Organisation de l'enseignemen
  - Modélisation décisionnelle
  - Types de problèmes de décision
  - Modélisation à partir de données
  - Modélisation à partir de données
  - Étapes générales
  - Quelques définitions
  - Choix d'une fonction de perteChoix des familles paramétriques
  - Estimation du modèle

  - Comment mesurer la
  - Validation croisée
  - Courbes ROC
    - Sélection de modèles
    - Grid search pour le choix des hyperparai
      - otimization

### Modèle décisionnel

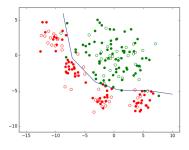
- Observations décrites par les valeurs prises par un ensemble de variables
- → Objectif: prédire, pour chaque donnée, la valeur d'une variable (expliquée ou « dépendante » ou « de sortie ») à partir des valeurs des autres variables (explicatives ou « d'entrée »)
  - Exemples
    - Une région d'une image représente un visage ou non?
    - Les symptômes correspondent à la maladie A ou B ou C ou aucune?
    - Quel est le volume d'algues vertes attendu en mai sur les plages de la commune?
    - 4 Quel sera le débit de la Loire à Tours dans 48h?
    - 5 Quelle est l'entité nommée dans « La Maison Blanche a démenti ces informations. » ?
    - Quelle est la région d'une image correspondant aux pantalons?

# Types de problèmes de décision

- Classement (ou discrimination) : la variable expliquée est une variable nominale, chaque observation possède une modalité (appelée en général classe)
  - → quel est le chiffre représenté par cette image?
- lacksquare Régression : la variable expliquée est une variable quantitative (domaine  $\subset \mathbb{R}$ )
  - → combien vaudra le CAC 40 dans une semaine?
- Prédiction structurée : la variable expliquée prend des valeurs dans un domaine de données structurées (les relations entre parties comptent)
  - ightarrow quelle est la forme de la protéine compte-tenu des molécules qui la composent ?

# Qu'est-ce qu'un modèle?

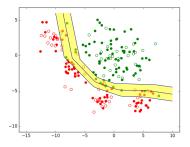
- Modèle = règle de décision
- Exemple : frontière de discrimination pour problème de classement à 2 classes



- Éventuellement complété par des critères de rejet (refus d'affectation
  - Refus de classer les données trop proches de la frontière (rejet d'ambiguïté
  - Refus de classer les données trop éloignées des données connues (rejet de nom représentativité)

# Qu'est-ce qu'un modèle?

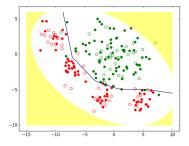
- Modèle = règle de décision
- Exemple : frontière de discrimination pour problème de classement à 2 classes



- Éventuellement complété par des critères de rejet (refus d'affectation)
  - Refus de classer les données trop proches de la frontière (rejet d'ambiguïté)
  - Refus de classer les données trop éloignées des données connues (rejet de non représentativité)

# Qu'est-ce qu'un modèle?

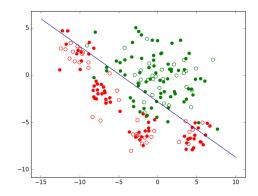
- Modèle = règle de décision
- Exemple : frontière de discrimination pour problème de classement à 2 classes



- Éventuellement complété par des critères de rejet (refus d'affectation)
  - Refus de classer les données trop proches de la frontière (rejet d'ambiguïté)
  - Refus de classer les données trop éloignées des données connues (rejet de non représentativité)

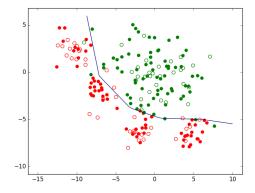
#### Classement

- Modèle : règle de classement, par ex. frontière de discrimination (trait bleu foncé)
- $\blacksquare$  Exemple : (2 var. explicatives pour chaque observation : abscisse X et ordonnée Y)



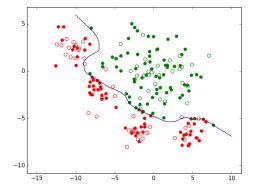
#### Classement

- Modèle : règle de classement, par ex. frontière de discrimination (trait bleu foncé)
- $\blacksquare$  Exemple : (2 var. explicatives pour chaque observation : abscisse X et ordonnée Y)



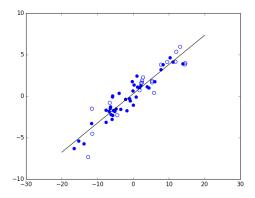
#### Classement

- Modèle : règle de classement, par ex. frontière de discrimination (trait bleu foncé)
- $\blacksquare$  Exemple : (2 var. explicatives pour chaque observation : abscisse X et ordonnée Y)



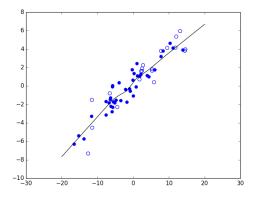
## Régression

- Modèle : règle de prédiction (trait noir dans la figure)
  - $\blacksquare$  Par ex. y = ax + b pour modèle linéaire
- lacktriangle Exemple : (variable explicative X en abscisse, variable expliquée Y en ordonnée)



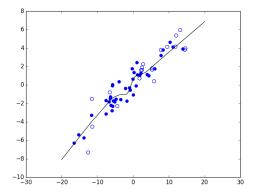
## Régression

- Modèle : règle de prédiction (trait noir dans la figure)
  - Par ex. y = ax + b pour modèle linéaire
- lacktriangle Exemple : (variable explicative X en abscisse, variable expliquée Y en ordonnée)



## Régression

- Modèle : règle de prédiction (trait noir dans la figure)
  - Par ex. y = ax + b pour modèle linéaire
- lacktriangle Exemple : (variable explicative X en abscisse, variable expliquée Y en ordonnée)



#### Prédiction structurée

- Modèle : règle de prédiction
- Exemples :
  - Déterminer que l'entité nommée de la phrase « La Maison Blanche a démenti ces informations. » est La Maison Blanche
    - Les classements des mots composant l'entité nommée ne sont pas indépendants
  - Délimiter la région correspondant aux pantalons dans l'image [5]





- Les affectations des pixels composant la région ne sont pas indépendantes
- Trier des photographies d'extérieur de la plus illuminée à la plus sombre
  - Problème d'ordonnancement relatif : le rang d'une image dépend des autres échantillons

#### Prédiction structurée

- Modèle : règle de prédiction
- Exemples :
  - Déterminer que l'entité nommée de la phrase « La Maison Blanche a démenti ces informations. » est La Maison Blanche
    - Les classements des mots composant l'entité nommée ne sont pas indépendants
  - Délimiter la région correspondant aux pantalons dans l'image [5]





- Les affectations des pixels composant la région ne sont pas indépendantes
- Trier des photographies d'extérieur de la plus illuminée à la plus sombre
  - Problème d'ordonnancement relatif : le rang d'une image dépend des autres échantillons

#### Prédiction structurée

- Modèle : règle de prédiction
- Exemples :
  - Déterminer que l'entité nommée de la phrase « La Maison Blanche a démenti ces informations. » est La Maison Blanche
    - Les classements des mots composant l'entité nommée ne sont pas indépendants
  - Délimiter la région correspondant aux pantalons dans l'image [5]

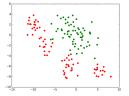




- Les affectations des pixels composant la région ne sont pas indépendantes
- Trier des photographies d'extérieur de la plus illuminée à la plus sombre.
  - Problème d'ordonnancement relatif : le rang d'une image dépend des autres échantillons.

#### Comment obtenir un modèle décisionnel

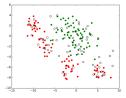
- Construction analytique, à partir d'une parfaite connaissance du phénomène
  - Exemples :
    - $\blacksquare$  Temps de vol  $\leftarrow$  distance et vitesse
    - lacktriangle Concentration de produit de réaction  $\leftarrow$  concentration de réactif et température
  - Néglige souvent l'impact de variables non contrôlables!
- A partir de données : ensemble d'observations pour lesquelles les valeurs des variables explicatives et des variables expliquées sont en général connues
  - → Apprentissage supervisé : à partir d'observations pour lesquelles les valeurs des variables explicatives et de la variable expliquée sont connues



 Apprentissage semi-supervisé (voir [2]): tient compte aussi des observations pour lesquelles les valeurs de la variable expliquée sont inconnues

#### Comment obtenir un modèle décisionnel

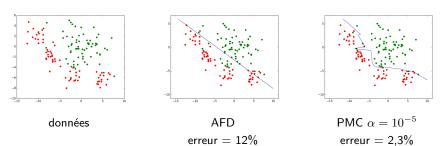
- Construction analytique, à partir d'une parfaite connaissance du phénomène
  - Exemples :
    - Temps de vol ← distance et vitesse
    - lacktriangle Concentration de produit de réaction  $\leftarrow$  concentration de réactif et température
  - Néglige souvent l'impact de variables non contrôlables!
- À partir de données : ensemble d'observations pour lesquelles les valeurs des variables explicatives et des variables expliquées sont en général connues
  - → Apprentissage supervisé : à partir d'observations pour lesquelles les valeurs des variables explicatives et de la variable expliquée sont connues



 Apprentissage semi-supervisé (voir [2]): tient compte aussi des observations pour lesquelles les valeurs de la variable expliquée sont inconnues

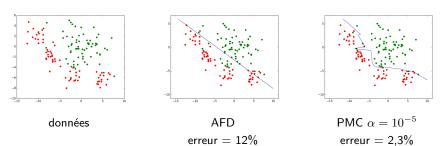
### Apprentissage et généralisation

- (Information de) Supervision = valeur de la variable expliquée
- Modélisation à partir de données (observations) d'apprentissage, qui disposent de l'information de supervision
- lacktriangle Choix famille paramétrique, puis optimisation des paramètres ightarrow modèle
- Erreur du modèle sur ces données = erreur d'apprentissage ou risque empirique



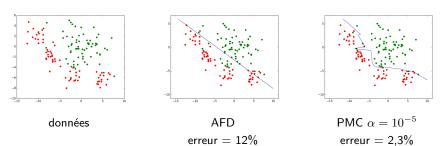
### Apprentissage et généralisation

- (Information de) Supervision = valeur de la variable expliquée
- Modélisation à partir de données (observations) d'apprentissage, qui disposent de l'information de supervision
- lacktriangle Choix famille paramétrique, puis optimisation des paramètres ightarrow modèle
- Erreur du modèle sur ces données = erreur d'apprentissage ou risque empirique



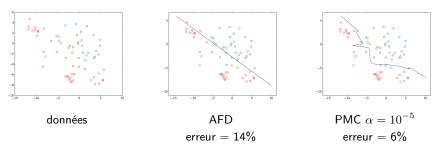
### Apprentissage et généralisation

- (Information de) Supervision = valeur de la variable expliquée
- Modélisation à partir de données (observations) d'apprentissage, qui disposent de l'information de supervision
- lacktriangle Choix famille paramétrique, puis optimisation des paramètres ightarrow modèle
- Erreur du modèle sur ces données = erreur d'apprentissage ou risque empirique



# Apprentissage et généralisation (2)

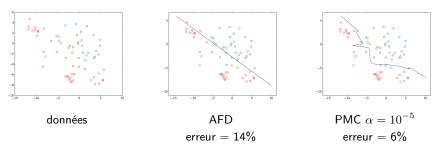
- Le modèle permet de prendre des décisions pour de futures (nouvelles) données
- Erreur du modèle sur ces futures données = erreur de généralisation ou risque espéré



ightarrow Objectif : avoir la meilleure généralisation (le risque espéré le plus faible)

# Apprentissage et généralisation (2)

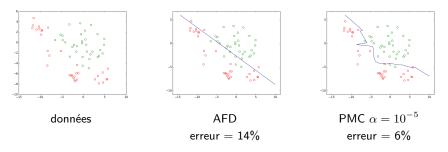
- Le modèle permet de prendre des décisions pour de futures (nouvelles) données
- Erreur du modèle sur ces futures données = erreur de généralisation ou risque espéré



ightarrow Objectif : avoir la meilleure généralisation (le risque espéré le plus faible)

# Apprentissage et généralisation (2)

- Le modèle permet de prendre des décisions pour de futures (nouvelles) données
- Erreur du modèle sur ces futures données = erreur de généralisation ou risque espéré



ightarrow Objectif : avoir la meilleure généralisation (le risque espéré le plus faible)

- Erreur d'apprentissage (facilement) mesurable car ces données sont disponibles
- Données futures inconnues ⇒ erreur de généralisation ne peut pas être mesurée
- Hypothèse importante : la distribution des données d'apprentissage est représentative de celle des données futures!
  - Or, on constate souvent que la distribution évolue dans le temps (n'est pas stationnaire) ⇒ il est nécessaire d'adapter régulièrement le modèle
- ightarrow Minimiser l'erreur d'apprentissage permet de minimiser l'erreur de généralisation ?
  - Considérons des données de test, non utilisées pour l'apprentissage mais disposant de l'information de supervision
  - Comparons trois modèles différents
  - If the control of th

- Erreur d'apprentissage (facilement) mesurable car ces données sont disponibles
- Données futures inconnues ⇒ erreur de généralisation ne peut pas être mesurée
- Hypothèse importante : la distribution des données d'apprentissage est représentative de celle des données futures!
  - Or, on constate souvent que la distribution évolue dans le temps (n'est pas stationnaire) ⇒ il est nécessaire d'adapter régulièrement le modèle
- ightarrow Minimiser l'erreur d'apprentissage permet de minimiser l'erreur de généralisation ?
  - Considérons des données de test, non utilisées pour l'apprentissage mais disposant de l'information de supervision
    - Comparons trois modèles différents
  - in Modèle Inézite obtenu par analyse factoriale discriminante (AFO) in Parceptina multicourties (PMC) pare un coefficient a d'outrit a (vielgia decay)  $\alpha = 10^{-12}$

- Erreur d'apprentissage (facilement) mesurable car ces données sont disponibles
- $lue{}$  Données futures inconnues  $\Rightarrow$  erreur de généralisation ne peut pas être mesurée
- Hypothèse importante : la distribution des données d'apprentissage est représentative de celle des données futures!
  - Or, on constate souvent que la distribution évolue dans le temps (n'est pas stationnaire) ⇒ il est nécessaire d'adapter régulièrement le modèle
- ightarrow Minimiser l'erreur d'apprentissage permet de minimiser l'erreur de généralisation ?
  - Considérons des données de test, non utilisées pour l'apprentissage mais disposant de l'information de supervision
    - Comparons trois modèles différents

- Erreur d'apprentissage (facilement) mesurable car ces données sont disponibles
- $lue{}$  Données futures inconnues  $\Rightarrow$  erreur de généralisation ne peut pas être mesurée
- Hypothèse importante : la distribution des données d'apprentissage est représentative de celle des données futures!
  - Or, on constate souvent que la distribution évolue dans le temps (n'est pas stationnaire) ⇒ il est nécessaire d'adapter régulièrement le modèle
- ightarrow Minimiser l'erreur d'apprentissage permet de minimiser l'erreur de généralisation ?
- Considérons des données de test, non utilisées pour l'apprentissage mais disposant de supervision
  - Comparons trois modèles différents

- Erreur d'apprentissage (facilement) mesurable car ces données sont disponibles
- Données futures inconnues ⇒ erreur de généralisation ne peut pas être mesurée
- Hypothèse importante : la distribution des données d'apprentissage est représentative de celle des données futures!
  - Or, on constate souvent que la distribution évolue dans le temps (n'est pas stationnaire) ⇒ il est nécessaire d'adapter régulièrement le modèle
- → Minimiser l'erreur d'apprentissage permet de minimiser l'erreur de généralisation ?
  - Considérons des données de test, non utilisées pour l'apprentissage mais disposant de l'information de supervision
  - Comparons trois modèles différents
  - Modèle linéaire obtenu par analyse factorielle discriminante (AFD)
  - Perceptron multicouches (PMC) avec un coefficient « d'oubli » (weight decay)  $\alpha=10^-$
  - $\blacksquare$  Perceptron multicouches (PMC) avec un coefficient « d oubli »  $lpha=\pm$

- Erreur d'apprentissage (facilement) mesurable car ces données sont disponibles
- $lue{}$  Données futures inconnues  $\Rightarrow$  erreur de généralisation ne peut pas être mesurée
- Hypothèse importante : la distribution des données d'apprentissage est représentative de celle des données futures!
  - Or, on constate souvent que la distribution évolue dans le temps (n'est pas stationnaire) ⇒ il est nécessaire d'adapter régulièrement le modèle
- → Minimiser l'erreur d'apprentissage permet de minimiser l'erreur de généralisation?
  - Considérons des données de test, non utilisées pour l'apprentissage mais disposant de l'information de supervision
  - Comparons trois modèles différents
    - Modèle linéaire obtenu par analyse factorielle discriminante (AFD)

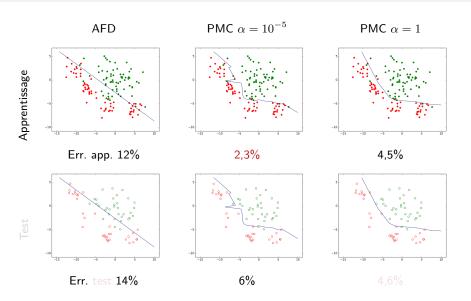
      Perceptron multicouches (PMC) avec un coefficient « d'oubli » (weight decay)  $\alpha=10^{-1}$ Perceptron multicouches (PMC) avec un coefficient « d'oubli »  $\alpha=1$

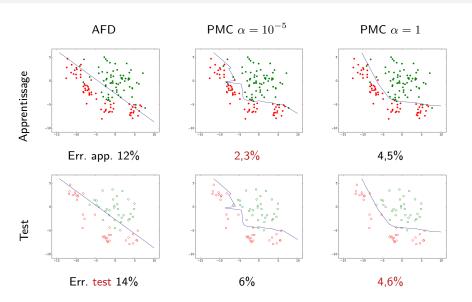
- Erreur d'apprentissage (facilement) mesurable car ces données sont disponibles
- Données futures inconnues ⇒ erreur de généralisation ne peut pas être mesurée
- Hypothèse importante : la distribution des données d'apprentissage est représentative de celle des données futures!
  - Or, on constate souvent que la distribution évolue dans le temps (n'est pas stationnaire) ⇒ il est nécessaire d'adapter régulièrement le modèle
- → Minimiser l'erreur d'apprentissage permet de minimiser l'erreur de généralisation ?
  - Considérons des données de test, non utilisées pour l'apprentissage mais disposant de l'information de supervision
  - Comparons trois modèles différents :
    - Modèle linéaire obtenu par analyse factorielle discriminante (AFD)
    - Perceptron multicouches (PMC) avec un coefficient «d'oubli» (weight decay)  $\alpha = 10^{-5}$
    - Perceptron multicouches (PMC) avec un coefficient « d'oubli »  $\alpha = 1$

- Erreur d'apprentissage (facilement) mesurable car ces données sont disponibles
- Données futures inconnues ⇒ erreur de généralisation ne peut pas être mesurée
- Hypothèse importante : la distribution des données d'apprentissage est représentative de celle des données futures!
  - Or, on constate souvent que la distribution évolue dans le temps (n'est pas stationnaire) ⇒ il est nécessaire d'adapter régulièrement le modèle
- → Minimiser l'erreur d'apprentissage permet de minimiser l'erreur de généralisation?
  - Considérons des données de test, non utilisées pour l'apprentissage mais disposant de l'information de supervision
  - Comparons trois modèles différents :
    - Modèle linéaire obtenu par analyse factorielle discriminante (AFD)
    - Perceptron multicouches (PMC) avec un coefficient «d'oubli» (weight decay)  $\alpha = 10^{-5}$
    - Perceptron multicouches (PMC) avec un coefficient « d'oubli »  $\alpha=1$

- Erreur d'apprentissage (facilement) mesurable car ces données sont disponibles
- Données futures inconnues ⇒ erreur de généralisation ne peut pas être mesurée
- Hypothèse importante : la distribution des données d'apprentissage est représentative de celle des données futures!
  - Or, on constate souvent que la distribution évolue dans le temps (n'est pas stationnaire) ⇒ il est nécessaire d'adapter régulièrement le modèle
- → Minimiser l'erreur d'apprentissage permet de minimiser l'erreur de généralisation ?
  - Considérons des données de test, non utilisées pour l'apprentissage mais disposant de l'information de supervision
  - Comparons trois modèles différents :
    - Modèle linéaire obtenu par analyse factorielle discriminante (AFD)
    - Perceptron multicouches (PMC) avec un coefficient «d'oubli» (weight decay)  $\alpha = 10^{-5}$
    - Perceptron multicouches (PMC) avec un coefficient « d'oubli »  $\alpha = 1$

- Erreur d'apprentissage (facilement) mesurable car ces données sont disponibles
- Données futures inconnues ⇒ erreur de généralisation ne peut pas être mesurée
- Hypothèse importante : la distribution des données d'apprentissage est représentative de celle des données futures!
  - Or, on constate souvent que la distribution évolue dans le temps (n'est pas stationnaire) ⇒ il est nécessaire d'adapter régulièrement le modèle
- → Minimiser l'erreur d'apprentissage permet de minimiser l'erreur de généralisation ?
  - Considérons des données de test, non utilisées pour l'apprentissage mais disposant de l'information de supervision
  - Comparons trois modèles différents :
    - Modèle linéaire obtenu par analyse factorielle discriminante (AFD)
    - Perceptron multicouches (PMC) avec un coefficient «d'oubli» (weight decay)  $\alpha = 10^{-5}$
    - Perceptron multicouches (PMC) avec un coefficient « d'oubli »  $\alpha=1$





- Le modèle qui a la plus faible erreur d'apprentissage n'a pas la plus faible erreur de test

  Cela reste valable si on compare des modèles issus de la même famille, par ex. par arrêt
- L'erreur d'apprentissage est en général une estimation optimiste de l'erreur de test
- L'écart entre erreur d'apprentissage et erreur de test dépend de la famille de modèles
- Si on ne peut pas mesurer l'erreur de généralisation, comment l'estimer?
  - Par l'erreur sur des données de test, non utilisées pour l'apprentissag
  - Grâce à une éventuelle borne supérieure sur l'écart entre erreur d'apprentissage et
  - erreur de généralisation : erreur généralisation ≤ erreur apprentissage + borne

- Le modèle qui a la plus faible erreur d'apprentissage n'a pas la plus faible erreur de test
  - Cela reste valable si on compare des modèles issus de la même famille, par ex. par arrêt précoce de la procédure d'optimisation
- L'erreur d'apprentissage est en général une estimation optimiste de l'erreur de test
- L'écart entre erreur d'apprentissage et erreur de test dépend de la famille de modèle
- Si on ne peut pas mesurer l'erreur de généralisation, comment l'estimer?
  - 🔟 Par l'erreur sur des données de test, non utilisées pour l'apprentissage
  - Grâce à une éventuelle borne supérieure sur l'écart entre erreur d'apprentissage et erreur de généralisation : erreur généralisation < erreur apprentissage + home</p>

- Le modèle qui a la plus faible erreur d'apprentissage n'a pas la plus faible erreur de test
  - Cela reste valable si on compare des modèles issus de la même famille, par ex. par arrêt précoce de la procédure d'optimisation
- L'erreur d'apprentissage est en général une estimation optimiste de l'erreur de test
- L'écart entre erreur d'apprentissage et erreur de test dépend de la famille de modèles
- Si on ne peut pas mesurer l'erreur de généralisation, comment l'estimer?
  - Par l'erreur sur des données de test, non utilisées pour l'apprentissage
  - Grâce à une éventuelle borne supérieure sur l'écart entre erreur d'apprentissage e

- Le modèle qui a la plus faible erreur d'apprentissage n'a pas la plus faible erreur de test
  - Cela reste valable si on compare des modèles issus de la même famille, par ex. par arrêt précoce de la procédure d'optimisation
- L'erreur d'apprentissage est en général une estimation optimiste de l'erreur de test
- L'écart entre erreur d'apprentissage et erreur de test dépend de la famille de modèles
- Si on ne peut pas mesurer l'erreur de généralisation, comment l'estimer?
  - Crâce à une éventuelle borne que évent entre erroux d'apprenties que l'écert entre erroux d'apprenties que l'
  - ☐ Grace a une eventuelle borne superieure sur l'ecart entre erreur d'apprentissage et erreur de généralisation : erreur généralisation ≤ erreur apprentissage + borne

#### Constats

- Le modèle qui a la plus faible erreur d'apprentissage n'a pas la plus faible erreur de test
  - Cela reste valable si on compare des modèles issus de la même famille, par ex. par arrêt précoce de la procédure d'optimisation
- L'erreur d'apprentissage est en général une estimation optimiste de l'erreur de test
- L'écart entre erreur d'apprentissage et erreur de test dépend de la famille de modèles
- Si on ne peut pas mesurer l'erreur de généralisation, comment l'estimer?

■ Grâce à une éventuelle borne supérieure sur l'écart entre erreur d'apprentissage et erreur de généralisation : erreur généralisation < erreur apprentissage + borne

- Le modèle qui a la plus faible erreur d'apprentissage n'a pas la plus faible erreur de test
  - Cela reste valable si on compare des modèles issus de la même famille, par ex. par arrêt précoce de la procédure d'optimisation
- L'erreur d'apprentissage est en général une estimation optimiste de l'erreur de test
- L'écart entre erreur d'apprentissage et erreur de test dépend de la famille de modèles
- Si on ne peut pas mesurer l'erreur de généralisation, comment l'estimer?

  - ☑ Grâce à une éventuelle borne supérieure sur l'écart entre erreur d'apprentissage et erreur de généralisation : erreur généralisation ≤ erreur apprentissage + borne

#### Constats

- Le modèle qui a la plus faible erreur d'apprentissage n'a pas la plus faible erreur de test
  - Cela reste valable si on compare des modèles issus de la même famille, par ex. par arrêt précoce de la procédure d'optimisation
- L'erreur d'apprentissage est en général une estimation optimiste de l'erreur de test
- L'écart entre erreur d'apprentissage et erreur de test dépend de la famille de modèles
- Si on ne peut pas mesurer l'erreur de généralisation, comment l'estimer?
  - Par l'erreur sur des données de test, non utilisées pour l'apprentissage
    - → Les observations disponibles avec information de supervision sont séparées en données d'apprentissage (→ obtenir le modèle) et données de test (→ estimer la généralisation
  - ☑ Grâce à une éventuelle borne supérieure sur l'écart entre erreur d'apprentissage et erreur de généralisation : erreur généralisation ≤ erreur apprentissage + borne

→ Lorsqu'elle existe, la borne peut être trop élevée pour être exploitable

- Le modèle qui a la plus faible erreur d'apprentissage n'a pas la plus faible erreur de test
  - Cela reste valable si on compare des modèles issus de la même famille, par ex. par arrêt précoce de la procédure d'optimisation
- L'erreur d'apprentissage est en général une estimation optimiste de l'erreur de test
- L'écart entre erreur d'apprentissage et erreur de test dépend de la famille de modèles
- Si on ne peut pas mesurer l'erreur de généralisation, comment l'estimer?
  - Par l'erreur sur des données de test, non utilisées pour l'apprentissage
    - → Les observations disponibles avec information de supervision sont séparées en données d'apprentissage (→ obtenir le modèle) et données de test (→ estimer la généralisation)
  - ☑ Grâce à une éventuelle borne supérieure sur l'écart entre erreur d'apprentissage et erreur de généralisation : erreur généralisation ≤ erreur apprentissage + borne

- Le modèle qui a la plus faible erreur d'apprentissage n'a pas la plus faible erreur de test
  - Cela reste valable si on compare des modèles issus de la même famille, par ex. par arrêt précoce de la procédure d'optimisation
- L'erreur d'apprentissage est en général une estimation optimiste de l'erreur de test
- L'écart entre erreur d'apprentissage et erreur de test dépend de la famille de modèles
- Si on ne peut pas mesurer l'erreur de généralisation, comment l'estimer?
  - Par l'erreur sur des données de test, non utilisées pour l'apprentissage
    - → Les observations disponibles avec information de supervision sont séparées en données d'apprentissage (→ obtenir le modèle) et données de test (→ estimer la généralisation)
  - ☑ Grâce à une éventuelle borne supérieure sur l'écart entre erreur d'apprentissage et erreur de généralisation : erreur généralisation ≤ erreur apprentissage + borne
    - → Lorsqu'elle existe. la borne peut être trop élevée pour être exploitable

- ILe modèle qui a la plus faible erreur d'apprentissage n'a pas la plus faible erreur de test
  - Cela reste valable si on compare des modèles issus de la même famille, par ex. par arrêt précoce de la procédure d'optimisation
- L'erreur d'apprentissage est en général une estimation optimiste de l'erreur de test
- L'écart entre erreur d'apprentissage et erreur de test dépend de la famille de modèles
- Si on ne peut pas mesurer l'erreur de généralisation, comment l'estimer?
  - Par l'erreur sur des données de test, non utilisées pour l'apprentissage
    - → Les observations disponibles avec information de supervision sont séparées en données d'apprentissage (→ obtenir le modèle) et données de test (→ estimer la généralisation)
  - ☑ Grâce à une éventuelle borne supérieure sur l'écart entre erreur d'apprentissage et erreur de généralisation : erreur généralisation ≤ erreur apprentissage + borne
    - → Lorsqu'elle existe, la borne peut être trop élevée pour être exploitable

#### Plan du cours

- 1 Objectifs et contenu de l'enseignemen
- 2 Organisation de l'enseignement
- Modélisation décisionnelle
  - Modélisation à partir de données
  - Modélisation à partir de données : un cadre plus précis
  - f. A de de domices : un caure plus precis
  - Choix d'une fonction de perte
  - Choix des familles paramétriques
  - Estimation du modèle
  - Commont magginer la
  - Comment mesurer la cap
  - Validation croisée
  - Courbes ROC
    - Sélection de modèle
    - Grid search pour le choix des hyperparamètre
    - Randomized parameter optimization

#### Modélisation à partir de données : étapes générales

- Préparation des données et choix d'une fonction de perte (loss ou erreur)
- Choix des familles paramétriques dans lesquelles chercher des modèles
- Dans chaque famille, estimation du « meilleur » modèle intra-famille
- Choix du meilleur modèle entre familles
- Évaluation des performances de généralisation du modèle retenu

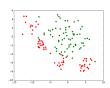
#### Modélisation à partir de données : un cadre plus précis

- Domaine des variables explicatives (ou espace d'entrée) :  $\mathcal{X}$  (par ex.  $\mathbb{R}^p$ )
- Domaine de la variable expliquée (ou espace de sortie) :  $\mathcal{Y}$  (par ex.  $\{-1,1\}$ ,  $\mathbb{R}$ )
- Données à modéliser décrites par variables aléatoires  $(X, Y) \in \mathcal{X} \times \mathcal{Y}$  suivant la distribution inconnue P
- Exemples

Classement:

$$\mathcal{X} \subset \mathbb{R}^2$$

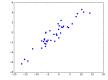




Régression:  $\mathcal{X} \subset \mathbb{R}$ 

$$\mathcal{Y} \subset \mathbb{R}$$





## Modélisation à partir de données : un cadre plus précis (2)

- Observations (données) avec information de supervision :  $\mathcal{D}_N = \{(\mathbf{x}_i, y_i)\}_{1 \leq i \leq N}$  correspondant à des tirages identiquement distribués suivant P
  - Supervision :  $\{y_i\}_{1 \le i \le N}$
  - Sauf cas particuliers (par ex. séries temporelles) on considère les données de  $\mathcal{D}_N$  issues de tirages indépendants

- o Objectif : trouver, dans une famille  $\mathcal{F}$ , une fonction  $f:\mathcal{X}\to\mathcal{Y}$  qui predit y a partir de x et présente le risque espéré  $R(f)=E_P[L(X,Y,f)]$  le plus faible
  - *L*() est la fonction de perte (ou d'erreur)
  - $\blacksquare$   $E_P$  est l'espérance par rapport à la distribution inconnue P
- Le choix d'une fonction de perte dépend de
  - La nature du problème de modélisation : classement, régression, prédiction structurée
  - lacksquare Le choix de la famille  ${\mathcal F}$  et de la procédure d'optimisation associé

## Modélisation à partir de données : un cadre plus précis (2)

- Observations (données) avec information de supervision :  $\mathcal{D}_N = \{(\mathbf{x}_i, y_i)\}_{1 \leq i \leq N}$  correspondant à des tirages identiquement distribués suivant P
  - Supervision :  $\{y_i\}_{1 \le i \le N}$
  - Sauf cas particuliers (par ex. séries temporelles) on considère les données de  $\mathcal{D}_N$  issues de tirages indépendants
- ightarrow Objectif : trouver, dans une famille  $\mathcal{F}$ , une fonction  $f: \mathcal{X} \to \mathcal{Y}$  qui prédit y à partir de  $\mathbf{x}$  et présente le risque espéré  $R(f) = E_P[L(X,Y,f)]$  le plus faible
  - L() est la fonction de perte (ou d'erreur)
  - $\blacksquare$   $E_P$  est l'espérance par rapport à la distribution inconnue P
  - Le choix d'une fonction de perte dépend de
    - La nature du problème de modélisation : classement, régression, prédiction structurée
    - lacksquare Le choix de la famille  ${\mathcal F}$  et de la procédure d'optimisation associée

#### Fonctions de perte pour problèmes de classement

- Perte 0-1 :  $L_{01}(\mathbf{x}, y, f) = \mathbf{1}_{f(\mathbf{x}) \neq y}$ 
  - $f(\mathbf{x}), y \in \mathcal{Y}$  ensemble fini
  - Perte nulle si prédiction correcte, perte unitaire si prédiction incorrecte
  - lacksquare Si  $\mathit{f}(x) \in \mathbb{R}$  alors  $L_{01}(x,y,\mathit{f}) = 1_{\mathit{H}(\mathit{f}(x)) 
    eq y}$ , avec  $\mathit{H}()$  fonction échelon adéquate

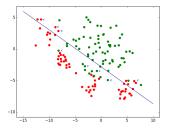


FIG. – Les flèches bleues indiquent quelques données mal classées par le modèle (frontière de discrimination linéaire, dans ce cas)

## Fonctions de perte pour problèmes de classement (2)

- *Hinge loss* pour la discrimination entre 2 classes en maximisant la marge (voir chapitre SVM) :  $L_h(\mathbf{x}, y, f) = \max\{0, 1 yf(\mathbf{x})\}$  (pour  $f(\mathbf{x}) \in \mathbb{R}$ )
  - lacksquare  $L_h$  n'est pas différentiable par rapport à f mais admet un sous-gradient
  - Des extensions existent pour le cas multi-classe et la prédiction structurée

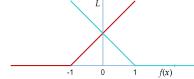


Fig. – Hinge loss pour y = -1 (en rouge) et y = 1 (en bleu)

#### Fonctions de perte pour problèmes de régression

- Perte quadratique :  $L_q(\mathbf{x}, y, f) = [f(\mathbf{x}) y]^2$ 
  - lacksquare f(x) est la prédiction du modèle f pour l'entrée x
  - lacksquare y est l'information de supervision (prédiction désirée) pour l'entrée x
  - $\blacksquare$  Différentiable par rapport à  $f(\mathbf{x}) \Rightarrow$  une optimisation basée sur le gradient peut être appliquée directement

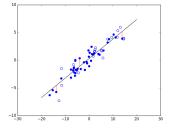
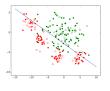


FIG. – Les traits rouges représentent des écarts entre trois prédictions d'un modèle (linéaire, da ce cas) et les prédictions désirées correspondantes

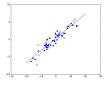
#### Familles paramétriques

- Modèles linéaires : prédiction = combinaison linéaire des variables explicatives
  - Exemples :

Classement : 
$$f(\mathbf{x}) = \mathbf{w}^T \mathbf{x} + \mathbf{w}_0$$
  $H(f(\mathbf{x})) \in \{-1, 1\}$ 



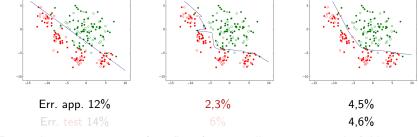
Régression :  $f(x) = w_1x + w_0$ 



- Peuvent s'avérer insuffisants (voir ci-dessus l'ex. de classes non linéairement séparables)
- Utile de commencer par un modèle linéaire, ne serait-ce que pour pouvoir comparer
- Modèles polynomiaux de degré borné : la capacité d'approximation (d'une frontière pour le classement, d'une dépendance pour la régression) augmente avec le degré
- Diverses familles de modèles non linéaires, par ex. perceptrons multicouches (PMC)
   d'architecture donnée, etc.

# Comment choisir la famille paramétrique?

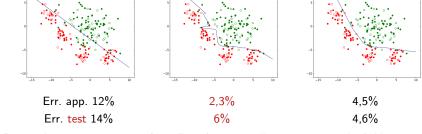
 $lue{}$  Modèles linéaires souvent insuffisants ightarrow pourquoi ne pas choisir systématiquement une famille de capacité d'approximation aussi grande que possible?



- → Risque de sur-apprentissage (overfitting) : erreur d'apprentissage très faible mais erreur de test comparativement élevée
- $\Rightarrow$  Ce n'est pas avec la capacité la plus grande qu'on obtient la meilleure généralisation
  - → Quel lien entre capacité et généralisation ?

# Comment choisir la famille paramétrique?

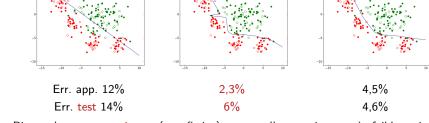
 $lue{}$  Modèles linéaires souvent insuffisants ightarrow pourquoi ne pas choisir systématiquement une famille de capacité d'approximation aussi grande que possible?



- → Risque de sur-apprentissage (overfitting) : erreur d'apprentissage très faible mais erreur de test comparativement élevée
- $\Rightarrow$  Ce n'est pas avec la capacité la plus grande qu'on obtient la meilleure généralisation
  - → Quel lien entre capacité et généralisation ?

# Comment choisir la famille paramétrique?

■ Modèles linéaires souvent insuffisants → pourquoi ne pas choisir systématiquement une famille de capacité d'approximation aussi grande que possible?



- → Risque de sur-apprentissage (overfitting) : erreur d'apprentissage très faible mais erreur de test comparativement élevée
- $\Rightarrow$  Ce n'est pas avec la capacité la plus grande qu'on obtient la meilleure généralisation
  - → Quel lien entre capacité et généralisation?

- Rappel de l'objectif : trouver, dans une famille  $\mathcal{F}$  choisie, une fonction (un modèle)  $f: \mathcal{X} \to \mathcal{Y}$  qui prédit y à partir de x et présente le risque espéré (ou théorique)  $R(f) = E_P[L(X, Y, f)]$  le plus faible
- R(f) ne peut pas être évalué car P est inconnue, mais on peut mesurer le risque empirique  $R_{\mathcal{D}_N}(f) = \frac{1}{N} \sum_{i=1}^N L(\mathbf{x}_i, y_i, f)$
- Si R(f) est inaccessible, comment répondre à l'objectif f
  - Minimisation du risque empirique (MRE) : considérer le modèle qui minimise l'erreux d'apprentissage,  $f_{\mathcal{D}_N} = \arg\min_{f \in \mathcal{F}} R_{\mathcal{D}_N}(f)$
  - Minimisation du risque empirique régularisé (MRER)  $= \arg\min_{\theta \in \mathcal{T}} |R_{D_{\theta}}(\theta) + \alpha G(\theta)|$
  - Minimisation du risque structurel (MRS) : séquence de familles de capacité qui augmente, estimation MRE dans chaque famille, choix tenant compte à la fois de  $\mathcal{D}_N$  et de la capacité

- Rappel de l'objectif : trouver, dans une famille  $\mathcal{F}$  choisie, une fonction (un modèle)  $f \colon \mathcal{X} \to \mathcal{Y}$  qui prédit y à partir de x et présente le risque espéré (ou théorique)  $R(f) = E_P[L(X,Y,f)]$  le plus faible
- R(f) ne peut pas être évalué car P est inconnue, mais on peut mesurer le risque empirique  $R_{\mathcal{D}_N}(f) = \frac{1}{N} \sum_{i=1}^N L(\mathbf{x}_i, y_i, f)$
- Si R(f) est inaccessible, comment répondre à l'objectif
  - Minimisation du risque empirique (MRE) : considérer le modèle qui minimise l'erreur d'apprentissage  $f_{rr} = \arg\min_{r \in \mathcal{F}} R_{rr}(f)$ 
    - Minimisation du risque empirique régularisé (MRER)
    - $f_{\mathcal{D}}^* = \arg\min_{f \in \mathcal{F}} [R_{\mathcal{D},i}(f) + \alpha G(f)]$
  - Minimisation du risque structurel (MRS) : séquence de familles de capacité qui augmente, estimation MRE dans chaque famille, choix tenant compte à la fois de  $\mathcal{D}$ 
    - et de la capacite

- Rappel de l'objectif : trouver, dans une famille  $\mathcal{F}$  choisie, une fonction (un modèle)  $f \colon \mathcal{X} \to \mathcal{Y}$  qui prédit y à partir de x et présente le risque espéré (ou théorique)  $R(f) = E_P[L(X,Y,f)]$  le plus faible
- R(f) ne peut pas être évalué car P est inconnue, mais on peut mesurer le risque empirique  $R_{\mathcal{D}_N}(f) = \frac{1}{N} \sum_{i=1}^N L(\mathbf{x}_i, y_i, f)$
- Si R(f) est inaccessible, comment répondre à l'objectif?
  - Minimisation du risque empirique (MRE) : considérer le modèle qui minimise l'erreur d'apprentissage, f<sup>\*</sup><sub>D...</sub> = arg min<sub>fe,F</sub> R<sub>D.u</sub>(f)
  - Minimisation du risque empirique régularisé (MRER) :  $f_{\mathcal{D}_N}^* = \arg\min_{f \in \mathcal{F}} [R_{\mathcal{D}_N}(f) + \alpha \mathcal{G}(f)]$
  - Minimisation du risque structurel (MRS) : séquence de familles de capacité qui augmente, estimation MRE dans chaque famille, choix tenant compte à la fois de  $\mathcal{D}_N$  et de la capacité

- Rappel de l'objectif : trouver, dans une famille  $\mathcal{F}$  choisie, une fonction (un modèle)  $f: \mathcal{X} \to \mathcal{Y}$  qui prédit y à partir de x et présente le risque espéré (ou théorique)  $R(f) = E_P[L(X, Y, f)]$  le plus faible
- R(f) ne peut pas être évalué car P est inconnue, mais on peut mesurer le risque empirique  $R_{\mathcal{D}_N}(f) = \frac{1}{N} \sum_{i=1}^N L(\mathbf{x}_i, y_i, f)$
- Si R(f) est inaccessible, comment répondre à l'objectif?
  - III Minimisation du risque empirique (MRE) : considérer le modèle qui minimise l'erreur d'apprentissage,  $f_{\mathcal{D}_N}^e = \arg\min_{f \in \mathcal{F}} R_{\mathcal{D}_N}(f)$
  - Minimisation du risque empirique régularisé (MRER) :  $f_{r_0}^* = \arg \min_{t \in \mathcal{T}} [Rp_{r_0}(t) + \alpha G(t)]$
  - Minimisation du risque structurel (MRS) : séquence de familles de capacité qui augmente, estimation MRE dans chaque famille, choix tenant compte à la fois de  $\mathcal{D}_N$  et de la capacité

- Rappel de l'objectif : trouver, dans une famille  $\mathcal{F}$  choisie, une fonction (un modèle)  $f \colon \mathcal{X} \to \mathcal{Y}$  qui prédit y à partir de x et présente le risque espéré (ou théorique)  $R(f) = E_P[L(X,Y,f)]$  le plus faible
- R(f) ne peut pas être évalué car P est inconnue, mais on peut mesurer le risque empirique  $R_{\mathcal{D}_N}(f) = \frac{1}{N} \sum_{i=1}^N L(\mathbf{x}_i, y_i, f)$
- Si R(f) est inaccessible, comment répondre à l'objectif?
  - III Minimisation du risque empirique (MRE) : considérer le modèle qui minimise l'erreur d'apprentissage,  $f_{\mathcal{D}_N}^e = \arg\min_{f \in \mathcal{F}} R_{\mathcal{D}_N}(f)$
  - Minimisation du risque empirique régularisé (MRER) :  $f_{\mathcal{D}_N}^* = \arg\min_{f \in \mathcal{F}} [R_{\mathcal{D}_N}(f) + \alpha G(f)]$
  - Minimisation du risque structurel (MRS) : séquence de familles de capacité qui augmente, estimation MRE dans chaque famille, choix tenant compte à la fois de  $\mathcal{D}_N$  et de la capacité

- Rappel de l'objectif : trouver, dans une famille  $\mathcal{F}$  choisie, une fonction (un modèle)  $f: \mathcal{X} \to \mathcal{Y}$  qui prédit y à partir de x et présente le risque espéré (ou théorique)  $R(f) = E_P[L(X, Y, f)]$  le plus faible
- R(f) ne peut pas être évalué car P est inconnue, mais on peut mesurer le risque empirique  $R_{\mathcal{D}_N}(f) = \frac{1}{N} \sum_{i=1}^N L(\mathbf{x}_i, y_i, f)$
- Si R(f) est inaccessible, comment répondre à l'objectif?
  - III Minimisation du risque empirique (MRE) : considérer le modèle qui minimise l'erreur d'apprentissage,  $f_{D_N}^e = \arg\min_{f \in \mathcal{F}} R_{D_N}(f)$
  - Minimisation du risque empirique régularisé (MRER) :  $f_{\mathcal{D}_N}^* = \arg\min_{f \in \mathcal{F}} [R_{\mathcal{D}_N}(f) + \alpha G(f)]$
  - Minimisation du risque structurel (MRS) : séquence de familles de capacité qui augmente, estimation MRE dans chaque famille, choix tenant compte à la fois de  $\mathcal{D}_N$  et de la capacité

#### Analyse des composantes du risque espéré

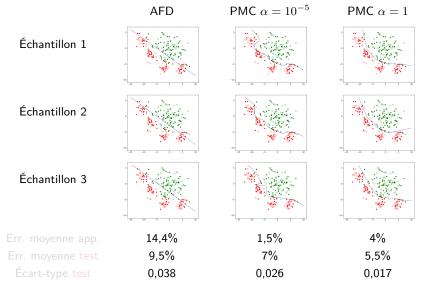
- Considérons
  - lacksquare  $f_{\mathcal{D}_N}^*$  la fonction de  $\mathcal{F}$  qui minimise le risque empirique  $R_{\mathcal{D}_N}$
  - lacksquare la fonction de  $\mathcal F$  qui minimise le risque espéré R, alors

$$R(f_{\mathcal{D}_{N}}^{*}) = R^{*} + [R(f^{*}) - R^{*}] + [R(f_{\mathcal{D}_{N}}^{*}) - R(f^{*})]$$

- R\* est le risque résiduel (ou risque de Bayes), borne inférieure
  - Strictement positif en présence de bruit : suivant le bruit, à un même x peuvent correspondre plusieurs valeurs de y
- $\mathbb{Z}[R(f^*)-R^*]$  est l'erreur d'approximation  $(\geq 0)$  car  $\mathcal{F}$  ne contient pas nécessairement la « vraie » fonction de décision
  - Nulle seulement si  $R^*$  peut être atteint par une fonction de  $\mathcal{F}$
- $[R(f_{D_N}^*) R(f^*)]$  est l'erreur d'estimation  $(\geq 0)$ 
  - La fonction de F qui minimise le risque empirique n'est pas nécessairement celle qui minimise le risque espéré

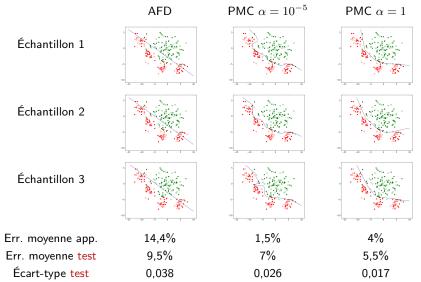
### Capacité, erreur d'approximation et erreur d'estimation

lacktriangle Résultats obtenus à partir de 3 familles sur 3 échantillons différents de  $\mathcal{D}_N$  :



### Capacité, erreur d'approximation et erreur d'estimation

 $\blacksquare$  Résultats obtenus à partir de 3 familles sur 3 échantillons différents de  $\mathcal{D}_{\textit{N}}$  :



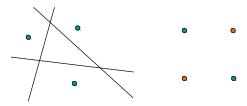
# Capacité, erreur d'approximation et erreur d'estimation (2)

Capacité famille linéaire (AFD) < capacité PMC  $\alpha=1$  < capacité PMC  $\alpha=10^{-5}$ 

- Famille linéaire (modèles obtenus ici par AFD)
  - Erreur d'apprentissage élevée donc capacité insuffisante pour ce problème
  - ⇒ Erreur d'approximation élevée (fort biais)
- $\blacksquare$  Famille définie par PMC 1 couche cachée de 100 neurones, avec coefficient « d'oubli »  $\alpha=10^{-5}$ 
  - Erreur d'approximation probablement faible car erreur d'apprentissage faible ⇒ capacité suffisante
  - lacktriangle Erreur de test bien plus élevée, variance supérieure à PMC lpha=1
  - ⇒ Erreur d'estimation élevée
- $\blacksquare$  Famille définie par PMC 1 couche cachée de 100 neurones, avec coefficient « d'oubli »  $\alpha=1$ 
  - Somme assez faible entre erreur d'approximation et erreur d'estimation, meilleure généralisation que les deux autres familles
  - Erreur de test assez faible et proche de l'erreur d'apprentissage

#### Comment mesurer la capacité?

- Considérons un ensemble de N vecteurs  $\{x_i\}_{1 \leq i \leq N} \in \mathbb{R}^p \to \mathsf{il}$  y a  $2^N$  façons différentes de le séparer en 2 parties
- **Définition**: la famille  $\mathcal F$  de fonctions  $f: \mathbb R^p \to \{-1,1\}$  pulvérise  $\{\mathbf x_i\}_{1 \le i \le N}$  si toutes les  $2^N$  séparations peuvent être construites avec des fonctions de  $\mathcal F$
- **Définition** (Vapnik-Chervonenkis) : l'ensemble  $\mathcal{F}$  est de VC-dimension h s'il pulvérise au moins un ensemble de h vecteurs et aucun ensemble de h+1 vecteurs
- Exemple : la VC-dimension de l'ensemble des hyperplans de  $\mathbb{R}^p$  est h=p+1
  - Dans R<sup>2</sup>, l'ensemble des droites pulvérise le triplet de points à gauche mais aucun quadruplet (par ex., aucune droite ne peut séparer les points bleus des rouges)



#### Lien entre capacité et généralisation

- La VC-dimension est une mesure intéressante de la capacité car elle permet d'obtenir une borne pour l'écart entre risque théorique et risque empirique
- Théorème [1] : soit  $R_{\mathcal{D}_N}(f)$  le risque empirique défini par la fonction de perte  $L_{01}(\mathbf{x},y,f)=\mathbf{1}_{f(\mathbf{x})\neq y}$ ; si la VC-dimension de  $\mathcal{F}$  est  $h<\infty$  alors pour toute  $f\in\mathcal{F}$ , avec une probabilité au mois égale à  $1-\delta$   $(0<\delta<1)$ , on a

$$R(f) \le R_{\mathcal{D}_N}(f) + \underbrace{\sqrt{\frac{h\left(\log \frac{2N}{h} + 1\right) - \log \frac{\delta}{4}}{N}}}_{B(N,\mathcal{F})} \quad \text{pour} \quad N > h$$

- $B(N, \mathcal{F})$  diminue quand  $N \uparrow$ , quand  $h \downarrow$  et quand  $\delta \uparrow$
- $B(N, \mathcal{F})$  ne fait pas intervenir le nombre de variables
- $B(N, \mathcal{F})$  ne fait pas intervenir la loi conjointe P
- → résultat dans le pire des cas, intéressant d'un point de vue théorique bien que peu utile en pratique

# Lien entre capacité et généralisation (2)

Conséquences de l'existence d'une borne

$$R(f) \leq R_{\mathcal{D}_N}(f) + B(N, \mathcal{F})$$

et de la forme de  $B(N, \mathcal{F})$ :

- lacksquare Famille  ${\mathcal F}$  de capacité trop faible (par ex. ici modèles linéaires)
  - $\Rightarrow$   $B(N, \mathcal{F})$  faible mais  $R_{\mathcal{D}_N}(f)$  (erreur d'apprentissage) élevé(e)
  - $\Rightarrow$  absence de garantie intéressante pour R(f)
- Famille  $\mathcal{F}$  de capacité trop élevée (par ex. ici PMC  $\alpha=10^-5$ )
  - $\Rightarrow R_{\mathcal{D}_N}(f)$  probablement faible mais  $B(N, \mathcal{F})$  élevée
  - $\Rightarrow$  absence de garantie intéressante pour R(f)
- lacksquare Famille  ${\mathcal F}$  de capacité « adéquate » (par ex. ici PMC lpha=1)
  - $\Rightarrow$   $R_{\mathcal{D}_{N}}(f)$  probablement faible et  $B(N,\mathcal{F})$  plutôt faible
  - $\Rightarrow$  garantie intéressante pour R(f)!

# Minimisation du risque empirique régularisé (MRER)

- La minimisation du risque empirique ne suffit pas à assurer une bonne généralisation, il faut maîtriser la capacité de  $\mathcal{F}$  (ou la complexité du modèle)
- La régularisation est une des solutions : le modèle est obtenu en minimisant la somme entre le risque empirique  $R_{\mathcal{D}_N}(f)$  et un terme G(f) qui pénalise (indirectement) la capacité

$$f_{\mathcal{D}_N}^* = \arg\min_{f \in \mathcal{F}} [R_{\mathcal{D}_N}(f) + \alpha G(f)]$$

lpha : hyperparamètre qui pondère le terme de régularisation

- lacksquare Différentes formes pour G(f), en rapport aussi avec le choix de la famille  $\mathcal F$ , par ex. :
  - $G(f) = \|\mathbf{w}\|_2^2$ , w étant le vecteur de paramètres du modèle; par ex. pour PMC terme « d'oubli » (weight decay)
  - Implicite : par ex., toujours pour PMC, terme G(f) absent mais arrêt précoce (early stopping) de l'algorithme d'optimisation non linéaire

### Minimisation du risque structurel (MRS)

minimisation du risque structurel [1]

Une solution de maîtrise explicite de la capacité de la famille de modèles est la

- Définition d'une séquence  $\mathcal{F}_1 \subset \mathcal{F}_2 \subset \mathcal{F}_3 \ldots$  de familles de capacités croissantes, c'est à dire pour lesquelles  $h_1 < h_2 < h_3 \ldots$
- Pour  $i \in \{1, 2, 3 \dots\}$ , minimisation dans chaque famille du risque empirique  $f_{\mathcal{D}_N}^{(i)*} = \arg\min_{f \in \mathcal{F}_i} R_{\mathcal{D}_N}(f)$
- Tenant compte de la borne trouvée pour le risque espéré R, sélection de  $f_{\mathcal{D}_N}^{(i)*}, i \in \{1,2,3\ldots\}$ , qui minimise  $R_{\mathcal{D}_N}(f_{\mathcal{D}_N}^{(i)*}) + B(N,\mathcal{F}_i)$

# Comment minimiser le risque empirique (régularisé)?

■ Dans une famille paramétrique  $\mathcal{F}$ , un modèle est défini par les valeurs d'un ensemble de paramètres, par ex.

■ Modèle linéaire pour la régression y = ax + b: a et b



- Perceptron multi-couches d'architecture donnée : poids des connexions de la (des) couche(s) cachée(s) et de la couche de sortie
- → Optimisation pour trouver les valeurs qui minimisent le critère (MRE, MRER)
  - Solution analytique directe : cas assez rare, par ex. certains modèles linéaires
  - Algorithmes itératifs, par ex.
    - Optimisation quadratique sous contraintes d'inégalité : SVM
    - Optimisation non linéaire plus générale : PMC, réseaux profonds

# Exemple : régression linéaire

- lacksquare Problème de régression avec  $\mathcal{X}=\mathbb{R}^p$ ,  $\mathcal{Y}=\mathbb{R}$ ,  $\mathcal{D}_N=\{(\mathbf{x}_i,y_i)\}_{1\leq i\leq N}$
- Famille de modèles linéaires  $\hat{y} = w_0 + \sum_{j=1}^p w_j x_{ji}$ , où  $\hat{y}$  est la prédiction du modèle
- Sous forme matricielle :  $\hat{\mathbf{y}} = \mathbf{X}\mathbf{w}$ , où  $\mathbf{X}$  est la matrice  $\mathbf{N} \times (p+1)$  dont les lignes sont les observations de  $\mathcal{D}_N$  et les colonnes correspondent aux variables (sauf pour la dernière qui est une colonne de 1 et permet d'inclure  $w_0$  dans  $\mathbf{w}$ )
- lacksquare On cherche le modèle (défini par le vecteur de paramètres lacksquare ) qui minimise
  - lacksquare MRE : l'erreur quadratique totale  $\sum_{i=1}^{N}(\hat{y}_i-y_i)^2$  sur  $\mathcal{D}_N$ 

    - $\blacksquare$  Si  $\textbf{X}^T\textbf{X}$  est inversible, alors  $\textbf{X}^+ = (\textbf{X}^T\textbf{X})^{-1}\textbf{X}^T$



- MRER : la somme entre l'erreur quadratique sur  $\mathcal{D}_N$  et un terme de régularisation, par ex. (cas particulier de régularisation Tikhonov),  $\sum_{i=1}^N (\hat{y}_i y_i)^2 + \|\mathbf{w}\|_2^2$ 
  - $\rightarrow$  Solution  $\mathbf{w}^* = (\mathbf{X}^T\mathbf{X} + \mathbf{I}_{p+1})^{-1}\mathbf{X}^T\mathbf{y}$ , où  $\mathbf{I}_{p+1}$  est la matrice unité de rang p+1

#### Modélisation décisionnelle : que faut-il retenir?

- Construire un modèle décisionnel à partir de données : supervision nécessaire
- Objectif : obtenir le modèle qui présente la meilleure généralisation
- Estimer la généralisation : non à partir de l'erreur d'apprentissage
- Chercher le bon compromis entre minimisation de la capacité de la famille de modèles et minimisation de l'erreur d'apprentissage

#### Plan du cours

- Objectifs et contenu de l'enseignemer
- 2 Organisation de l'enseignemen
  - Modélisation décisionnelle
    - Types de problèmes de décision
    - Modélisation à partir de données
    - Modélisation à partir de données : u
    - Etapes générales
    - Choix d'une fonction de perte
    - Choix des familles paramétriques
    - Estimation du modele
    - Comment mesurer la
- 5 Évaluation de modèles
  - Validation crois
  - Courbes NOC
  - = Crid search nour le chaix des hypernarar
  - Randomized parameter optimization

- A partir du risque empirique et en tenant compte de bornes de généralisation :  $R(f_{D_N}^*) \leq R_{D_N}(f_{D_N}^*) + B(N, \mathcal{F})$ 
  - ightarrow Lorsqu'elle existe, la borne est en général trop élevée pour être utile en pratique
- Par l'erreur sur des données de test, non utilisées pour l'apprentissage
  - Les observations disponibles avec information de supervision sont partitionnées (par échantillonnage uniforme, en général) en données d'apprentissage (70-80%) et données de test (20-30%)
    - Apprentissage (estimation) du modèle sur les données d'apprentissage
    - Estimation du risque espéré par l'erreur de ce modèle sur les données de tes
  - → Difficultés de cette approche
    - La mise de coté des données de test reduit le nombre de données utilisées pour l'apprentisses
       Cet estimateur du risque espéré a une variance élevée (un autre partitionnement produira d'autres ensembles d'apprentissage et de test)
  - → Validation croisée (cross-validation): plusieurs partitionnements apprentissage | test, obtenir à chaque fois un modèle sur les données d'apprentissage et l'évaluer sur les données de test associées, employer la moyenne comme estimation du risque espéré
    - tout en utilisant mieux les données disponibles!

- A partir du risque empirique et en tenant compte de bornes de généralisation :  $R(f_{D_N}^*) \leq R_{D_N}(f_{D_N}^*) + B(N, \mathcal{F})$ 
  - ightarrow Lorsqu'elle existe, la borne est en général trop élevée pour être utile en pratique
- Par l'erreur sur des données de test, non utilisées pour l'apprentissage
  - Les observations disponibles avec information de supervision sont partitionnées (par échantillonnage uniforme, en général) en données d'apprentissage (70-80%) et données de test (20-30%)
    - Apprentissage (estimation) du modèle sur les données d'apprentissage
    - Estimation du risque espéré par l'erreur de ce modèle sur les données de tes
  - → Difficultés de cette approche
    - La mise de coté des données de test reduit le nombre de données utilisées pour l'apprentis
       Cet estimateur du risque espéré a une variance élevée (un autre partitionnement produira d'autres ensembles d'apprentissage et de test)
  - → Validation croisée (cross-validation): plusieurs partitionnements apprentissage | test, obtenir à chaque fois un modèle sur les données d'apprentissage et l'évaluer sur les données de test associées, employer la moyenne comme estimation du risque espéré
    - tout an utilizant minus les dennées dispenible

- A partir du risque empirique et en tenant compte de bornes de généralisation :  $R(f_{D_N}^*) \leq R_{D_N}(f_{D_N}^*) + B(N, \mathcal{F})$ 
  - ightarrow Lorsqu'elle existe, la borne est en général trop élevée pour être utile en pratique
- Par l'erreur sur des données de test, non utilisées pour l'apprentissage
  - Les observations disponibles avec information de supervision sont partitionnées (par échantillonnage uniforme, en général) en données d'apprentissage (70-80%) et données de test (20-30%)
    - Apprentissage (estimation) du modèle sur les données d'apprentissage
    - $\blacksquare$  Estimation du risque espéré par l'erreur de ce modèle sur les données de test
  - → Difficultés de cette approche
    - La mise de coté des données de test reduit le nombre de données utilisées pour l'apprentis
       Cet estimateur du risque espéré a une variance élevée (un autre partitionnement produira d'autres ensembles d'apprentissage et de test)
  - → Validation croisée (cross-validation): plusieurs partitionnements apprentissage | test, obtenir à chaque fois un modèle sur les données d'apprentissage et l'évaluer sur les données de test associées, employer la moyenne comme estimation du risque espéré
    - ... tout en utilisant mieux les données disponibles !

- A partir du risque empirique et en tenant compte de bornes de généralisation :  $R(f_{D_N}^*) \leq R_{D_N}(f_{D_N}^*) + B(N, \mathcal{F})$ 
  - ightarrow Lorsqu'elle existe, la borne est en général trop élevée pour être utile en pratique
- Par l'erreur sur des données de test, non utilisées pour l'apprentissage
  - Les observations disponibles avec information de supervision sont partitionnées (par échantillonnage uniforme, en général) en données d'apprentissage (70-80%) et données de test (20-30%)
    - Apprentissage (estimation) du modèle sur les données d'apprentissage
    - Estimation du risque espéré par l'erreur de ce modèle sur les données de test
  - → Difficultés de cette approche :
    - La mise de côté des données de test réduit le nombre de données utilisées pour l'apprentissage
    - Cet estimateur du risque espéré a une variance élevée (un autre partitionnement produira d'autres ensembles d'apprentissage et de test)
  - → Validation croisée (cross-validation): plusieurs partitionnements apprentissage | test, obtenir à chaque fois un modèle sur les données d'apprentissage et l'évaluer sur les données de test associées, employer la moyenne comme estimation du risque espéré
    - ⇒ estimateur de variance plus faible,
      - ... tout en utilisant mieux les données disponibles!

- A partir du risque empirique et en tenant compte de bornes de généralisation :  $R(f_{D_N}^*) \leq R_{D_N}(f_{D_N}^*) + B(N, \mathcal{F})$ 
  - ightarrow Lorsqu'elle existe, la borne est en général trop élevée pour être utile en pratique
- Par l'erreur sur des données de test, non utilisées pour l'apprentissage
  - Les observations disponibles avec information de supervision sont partitionnées (par échantillonnage uniforme, en général) en données d'apprentissage (70-80%) et données de test (20-30%)
    - Apprentissage (estimation) du modèle sur les données d'apprentissage
    - Estimation du risque espéré par l'erreur de ce modèle sur les données de test
  - $\rightarrow \ \, \text{Difficult\'es de cette approche} :$ 
    - La mise de côté des données de test réduit le nombre de données utilisées pour l'apprentissage
    - Cet estimateur du risque espéré a une variance élevée (un autre partitionnement produira d'autres ensembles d'apprentissage et de test)
  - → Validation croisée (cross-validation): plusieurs partitionnements apprentissage | test, obtenir à chaque fois un modèle sur les données d'apprentissage et l'évaluer sur les données de test associées, employer la moyenne comme estimation du risque espéré
    - $\Rightarrow$  estimateur de variance plus faible,
      - ... tout en utilisant mieux les données disponibles!

### Comment estimer le risque espéré

- A partir du risque empirique et en tenant compte de bornes de généralisation :  $R(f_{D_N}^*) \leq R_{D_N}(f_{D_N}^*) + B(N, \mathcal{F})$ 
  - ightarrow Lorsqu'elle existe, la borne est en général trop élevée pour être utile en pratique
- Par l'erreur sur des données de test, non utilisées pour l'apprentissage
  - Les observations disponibles avec information de supervision sont partitionnées (par échantillonnage uniforme, en général) en données d'apprentissage (70-80%) et données de test (20-30%)
    - Apprentissage (estimation) du modèle sur les données d'apprentissage
    - Estimation du risque espéré par l'erreur de ce modèle sur les données de test
  - $\rightarrow \ \, \text{Difficult\'es de cette approche} :$ 
    - La mise de côté des données de test réduit le nombre de données utilisées pour l'apprentissage
    - Cet estimateur du risque espéré a une variance élevée (un autre partitionnement produira d'autres ensembles d'apprentissage et de test)
  - → Validation croisée (cross-validation): plusieurs partitionnements apprentissage | test, obtenir à chaque fois un modèle sur les données d'apprentissage et l'évaluer sur les données de test associées, employer la moyenne comme estimation du risque espéré

⇒ estimateur de variance plus faible,

tout en utilisant mieux les données disponibles!

- A partir du risque empirique et en tenant compte de bornes de généralisation :  $R(f_{D_N}^*) \leq R_{D_N}(f_{D_N}^*) + B(N, \mathcal{F})$ 
  - ightarrow Lorsqu'elle existe, la borne est en général trop élevée pour être utile en pratique
- Par l'erreur sur des données de test, non utilisées pour l'apprentissage
  - Les observations disponibles avec information de supervision sont partitionnées (par échantillonnage uniforme, en général) en données d'apprentissage (70-80%) et données de test (20-30%)
    - Apprentissage (estimation) du modèle sur les données d'apprentissage
    - Estimation du risque espéré par l'erreur de ce modèle sur les données de test
  - → Difficultés de cette approche :
    - La mise de côté des données de test réduit le nombre de données utilisées pour l'apprentissage
    - Cet estimateur du risque espéré a une variance élevée (un autre partitionnement produira d'autres ensembles d'apprentissage et de test)
  - → Validation croisée (cross-validation): plusieurs partitionnements apprentissage | test, obtenir à chaque fois un modèle sur les données d'apprentissage et l'évaluer sur les données de test associées, employer la moyenne comme estimation du risque espéré
    - ⇒ estimateur de variance plus faible
    - tout en utilisant mieux les données disponibles l

- A partir du risque empirique et en tenant compte de bornes de généralisation :  $R(f_{D_N}^*) \leq R_{D_N}(f_{D_N}^*) + B(N, \mathcal{F})$ 
  - ightarrow Lorsqu'elle existe, la borne est en général trop élevée pour être utile en pratique
- Par l'erreur sur des données de test, non utilisées pour l'apprentissage
  - Les observations disponibles avec information de supervision sont partitionnées (par échantillonnage uniforme, en général) en données d'apprentissage (70-80%) et données de test (20-30%)
    - Apprentissage (estimation) du modèle sur les données d'apprentissage
    - Estimation du risque espéré par l'erreur de ce modèle sur les données de test
  - $\rightarrow \ \, \text{Difficult\'es de cette approche} :$ 
    - La mise de côté des données de test réduit le nombre de données utilisées pour l'apprentissage
    - Cet estimateur du risque espéré a une variance élevée (un autre partitionnement produira d'autres ensembles d'apprentissage et de test)
  - → Validation croisée (cross-validation): plusieurs partitionnements apprentissage | test, obtenir à chaque fois un modèle sur les données d'apprentissage et l'évaluer sur les données de test associées, employer la moyenne comme estimation du risque espéré
    - ⇒ estimateur de variance plus faible,

- A partir du risque empirique et en tenant compte de bornes de généralisation :  $R(f_{D_N}^*) \leq R_{D_N}(f_{D_N}^*) + B(N, \mathcal{F})$ 
  - ightarrow Lorsqu'elle existe, la borne est en général trop élevée pour être utile en pratique
- Par l'erreur sur des données de test, non utilisées pour l'apprentissage
  - Les observations disponibles avec information de supervision sont partitionnées (par échantillonnage uniforme, en général) en données d'apprentissage (70-80%) et données de test (20-30%)
    - Apprentissage (estimation) du modèle sur les données d'apprentissage
    - Estimation du risque espéré par l'erreur de ce modèle sur les données de test
  - → Difficultés de cette approche :
    - La mise de côté des données de test réduit le nombre de données utilisées pour l'apprentissage
    - Cet estimateur du risque espéré a une variance élevée (un autre partitionnement produira d'autres ensembles d'apprentissage et de test)
  - → Validation croisée (cross-validation): plusieurs partitionnements apprentissage | test, obtenir à chaque fois un modèle sur les données d'apprentissage et l'évaluer sur les données de test associées, employer la moyenne comme estimation du risque espéré
    - ⇒ estimateur de variance plus faible,
    - ... tout en utilisant mieux les données disponibles!

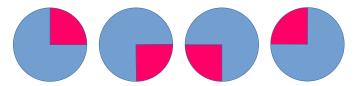
#### Validation croisée

#### Méthodes exhaustives :

- Leave p out (LPO) : N-p données pour l'apprentissage et p pour la validation  $\Rightarrow C_N^p$ découpages possibles donc  $C_N^p$  modèles à apprendre  $\Rightarrow$  coût excessif
- Leave one out (LOO): N-1 données pour l'apprentissage et 1 pour la validation  $\Rightarrow$  $C_N^1 = N$  découpages possibles (donc N modèles)  $\Rightarrow$  coût élevé

#### Méthodes non exhaustives :

 $\blacksquare$  k-fold: partitionnement fixé des N données en k parties, apprentissage sur k-1 parties et validation sur la k-ême  $\Rightarrow k$  modèles seulement (souvent k=5 ou k=10)



■ Échantillonnage répété (shuffle and split) : échantillon aléatoire de p données pour le test (les autres N-p pour l'apprentissage), on répète cela k fois  $\Rightarrow k$  modèles

#### Validation croisée : quelle méthode préférer?

- LPO très rarement employée car excessivement coûteuse
- LOO vs k-fold : k-fold préférée en général
  - LOO plus coûteuse car  $N \gg k$
  - Variance en général supérieure pour LOO
  - Estimation k-fold pessimiste car chaque modèle apprend sur  $\frac{k-1}{k}N < N-1$  données
- Shuffle and split vs k-fold
  - Pour k-fold le nombre de modèles (k) est lié à la proportion de données de test (1/k), shuffle and split moins contraignante
  - Pour shuffle and split certaines données ne sont dans aucun échantillon alors que d'autres sont dans plusieurs échantillons
- Quelle que soit la méthode, tous les partitionnements peuvent être explorés en parallèle (sur processeurs multi-cœur ou plateformes distribuées)

#### Validation croisée : précautions à prendre

- Problème de classement avec classes (très) déséquilibrées : pour s'assurer de conserver les rapports entre les classes dans tous les découpages, utiliser

- Problème de classement avec classes (très) déséquilibrées : pour s'assurer de conserver les rapports entre les classes dans tous les découpages, utiliser
  - Un partitionnement adapté pour k-fold (par ex. StratifiedKFold dans Scikit-learn)
  - Un échantillonnage stratifié pour shuffle and split (par ex. StratifiedShuffleSplit dans Scikit-learn)
  - LOO peut être employée telle quelle
- Observations qui ne sont pas indépendantes
  - Series temporelles: les observations successives sont correlees, le découpage doit être fait par séquences sur les observations ordonnées et non après shuffle sur les observations individuelles
  - Données groupées: dans un même groupe, les observations ne sont pas indépendantes les données de test doivent provenir de groupes différents de ceux dont sont issues les données d'apprentissage

- Problème de classement avec classes (très) déséquilibrées : pour s'assurer de conserver les rapports entre les classes dans tous les découpages, utiliser
  - Un partitionnement adapté pour k-fold (par ex. StratifiedKFold dans Scikit-learn)
  - Un échantillonnage stratifié pour shuffle and split (par ex. StratifiedShuffleSplit dans Scikit-learn)
  - LOO peut etre employee telle quelle
- Observations qui ne sont pas indépendante
  - Séries temporelles: les observations successives sont corrélées, le découpage doit être fait par séquences sur les observations ordonnées et non après shuffle sur les observations individuelles
  - Données groupées : dans un même groupe, les observations ne sont pas indépendantes les données de test doivent provenir de groupes différents de ceux dont sont issues les données d'apprentissage

- Problème de classement avec classes (très) déséquilibrées : pour s'assurer de conserver les rapports entre les classes dans tous les découpages, utiliser
  - Un partitionnement adapté pour k-fold (par ex. StratifiedKFold dans Scikit-learn)
  - Un échantillonnage stratifié pour shuffle and split (par ex. StratifiedShuffleSplit dans Scikit-learn)
  - LOO peut être employée telle quelle
- Observations qui ne sont pas indépendantes
  - Séries temporelles: les observations successives sont corrélées, le découpage doit être fait par séquences sur les observations ordonnées et non après shuffle sur les observations individuelles
  - Données groupées : dans un même groupe, les observations ne sont pas indépendantesseles données de test doivent provenir de groupes différents de ceux dont sont issues les données d'apprentissage

- Problème de classement avec classes (très) déséquilibrées : pour s'assurer de conserver les rapports entre les classes dans tous les découpages, utiliser
  - Un partitionnement adapté pour k-fold (par ex. StratifiedKFold dans Scikit-learn)
  - Un échantillonnage stratifié pour shuffle and split (par ex. StratifiedShuffleSplit dans Scikit-learn)
  - LOO peut être employée telle quelle
- Observations qui ne sont pas indépendantes
  - Séries temporelles : les observations successives sont corrélées, le découpage doit être fait par séquences sur les observations ordonnées et non après shuffle sur les observations individuelles
  - Données groupées: dans un même groupe, les observations ne sont pas indépendantes; les données de test doivent provenir de groupes différents de ceux dont sont issues les données d'apprentissage

Évaluation de modèles Validation croisée 44/56

- Problème de classement avec classes (très) déséquilibrées : pour s'assurer de conserver les rapports entre les classes dans tous les découpages, utiliser
  - Un partitionnement adapté pour k-fold (par ex. StratifiedKFold dans Scikit-learn)
  - Un échantillonnage stratifié pour shuffle and split (par ex. StratifiedShuffleSplit dans Scikit-learn)
  - LOO peut être employée telle quelle
- Observations qui ne sont pas indépendantes
  - Séries temporelles : les observations successives sont corrélées, le découpage doit être fait par séquences sur les observations ordonnées et non après shuffle sur les observations individuelles
  - Données groupées : dans un même groupe, les observations ne sont pas indépendantes les données de test doivent provenir de groupes différents de ceux dont sont issues les données d'apprentissage

Évaluation de modèles Validation croisée 44/56

- Problème de classement avec classes (très) déséquilibrées : pour s'assurer de conserver les rapports entre les classes dans tous les découpages, utiliser
  - Un partitionnement adapté pour k-fold (par ex. StratifiedKFold dans Scikit-learn)
  - Un échantillonnage stratifié pour shuffle and split (par ex. StratifiedShuffleSplit dans Scikit-learn)
  - LOO peut être employée telle quelle
- Observations qui ne sont pas indépendantes
  - Séries temporelles : les observations successives sont corrélées, le découpage doit être fait par séquences sur les observations ordonnées et non après shuffle sur les observations individuelles
  - Données groupées : dans un même groupe, les observations ne sont pas indépendantes; les données de test doivent provenir de groupes différents de ceux dont sont issues les données d'apprentissage

- Estimation du risque espéré d'un modèle de classement : taux de mauvais classement sur les données de test
  - Taux de mauvais classement  $\leftarrow$  fonction de perte  $L_0$
  - ightarrow coût symétrique : même coût si le modèle se trompe dans un sens ou dans l'autre
- De nombreux problèmes présentent des coûts asymétriques, par ex.
  - Pour un cargo, la non détection d'un autre navire par le radar peut mener à une collision, alors qu'une fausse alerte provoque seulement un ralentissement temporaire
  - La non détection de la maladie grave d'un patient est dramatique, alors que la détection erronée d'une telle maladie pour un patient sain est moins problématique
- ⇒ Comment examiner les caractéristiques de différents modèles lorsque les coûts sont asymétriques, sans fixer le «degré» d'asymétrie?

- Estimation du risque espéré d'un modèle de classement : taux de mauvais classement sur les données de test
  - Taux de mauvais classement  $\leftarrow$  fonction de perte  $L_{01}$
  - → coût symétrique : même coût si le modèle se trompe dans un sens ou dans l'autre
- De nombreux problèmes présentent des coûts asymétriques, par ex
  - Pour un cargo, la non détection d'un autre navire par le radar peut mener à une collision, alors qu'une fausse alerte provoque seulement un ralentissement temporaire
  - La non détection de la maladie grave d'un patient est dramatique, alors que la détection erronée d'une telle maladie pour un patient sain est moins problématique.
- ⇒ Comment examiner les caractéristiques de différents modèles lorsque les coûts sont asymétriques, sans fixer le «degré» d'asymétrie?

- Estimation du risque espéré d'un modèle de classement : taux de mauvais classement sur les données de test
  - lacktriangle Taux de mauvais classement  $\leftarrow$  fonction de perte  $L_{01}$
  - → coût symétrique : même coût si le modèle se trompe dans un sens ou dans l'autre
- De nombreux problèmes présentent des coûts asymétriques, par ex
  - Pour un cargo, la non detection d'un autre navire par le radar peut mener a une collision, alors qu'une fausse alerte provoque seulement un ralentissement temporaire la La non détection de la maladie grave d'un patient est dramatique, alors que la détection erronée d'une telle maladie pour un patient sain est moins problématique
- ⇒ Comment examiner les caractéristiques de différents modèles lorsque les coûts sont asymétriques, sans fixer le « degré » d'asymétrie?

- Estimation du risque espéré d'un modèle de classement : taux de mauvais classement sur les données de test
  - Taux de mauvais classement  $\leftarrow$  fonction de perte  $L_{01}$
  - ightarrow coût symétrique : même coût si le modèle se trompe dans un sens ou dans l'autre
- De nombreux problèmes présentent des coûts asymétriques, par ex.
  - Pour un cargo, la non détection d'un autre navire par le radar peut mener à une collision, alors qu'une fausse alerte provoque seulement un ralentissement temporaire
  - La non détection de la maladie grave d'un patient est dramatique, alors que la détection erronée d'une telle maladie pour un patient sain est moins problématique
- Comment examiner les caractéristiques de différents modèles lorsque les coûts sont asymétriques, sans fixer le « degré » d'asymétrie?

- Estimation du risque espéré d'un modèle de classement : taux de mauvais classement sur les données de test
  - Taux de mauvais classement  $\leftarrow$  fonction de perte  $L_{01}$
- ightarrow coût symétrique : même coût si le modèle se trompe dans un sens ou dans l'autre
- De nombreux problèmes présentent des coûts asymétriques, par ex.
  - Pour un cargo, la non détection d'un autre navire par le radar peut mener à une collision, alors qu'une fausse alerte provoque seulement un ralentissement temporaire
  - La non détection de la maladie grave d'un patient est dramatique, alors que la détection erronée d'une telle maladie pour un patient sain est moins problématique
- Comment examiner les caractéristiques de différents modèles lorsque les coûts sont asymétriques, sans fixer le « degré » d'asymétrie?

- Estimation du risque espéré d'un modèle de classement : taux de mauvais classement sur les données de test
  - Taux de mauvais classement  $\leftarrow$  fonction de perte  $L_{01}$
- ightarrow coût symétrique : même coût si le modèle se trompe dans un sens ou dans l'autre
- De nombreux problèmes présentent des coûts asymétriques, par ex.
  - Pour un cargo, la non détection d'un autre navire par le radar peut mener à une collision, alors qu'une fausse alerte provoque seulement un ralentissement temporaire
  - La non détection de la maladie grave d'un patient est dramatique, alors que la détection erronée d'une telle maladie pour un patient sain est moins problématique
- ⇒ Comment examiner les caractéristiques de différents modèles lorsque les coûts sont asymétriques, sans fixer le « degré » d'asymétrie?

- Estimation du risque espéré d'un modèle de classement : taux de mauvais classement sur les données de test
  - Taux de mauvais classement  $\leftarrow$  fonction de perte  $L_{01}$
- ightarrow coût symétrique : même coût si le modèle se trompe dans un sens ou dans l'autre
- De nombreux problèmes présentent des coûts asymétriques, par ex.
  - Pour un cargo, la non détection d'un autre navire par le radar peut mener à une collision, alors qu'une fausse alerte provoque seulement un ralentissement temporaire
  - La non détection de la maladie grave d'un patient est dramatique, alors que la détection erronée d'une telle maladie pour un patient sain est moins problématique
- ⇒ Comment examiner les caractéristiques de différents modèles lorsque les coûts sont asymétriques, sans fixer le « degré » d'asymétrie?

## Terminologie pour la discrimination entre 2 classes

- Une classe peut être considérée la classe « d'intérêt »
- Le modèle appris est vu comme le « détecteur » de la classe d'intérêt
- Pour un tel détecteur appris, les cas suivants peuvent être constatés :

	Classe présente	Classe absente
Classe détectée	Vrai Positif	Faux Positif
Classe non détectée	Faux Négatif	Vrai Négatif

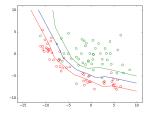
On définit les mesures suivantes :

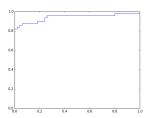
Taux de vrais positifs (ou sensibilité) = 
$$\frac{Vrais\ Positifs}{Total\ Positifs} = \frac{VP}{VP\ + FN}$$
  
Taux de faux positifs (ou  $1-$  spécificité) =  $\frac{Faux\ Positifs}{Total\ Négatifs} = \frac{FP}{VN\ + FP} = 1 - \frac{VN}{VN\ + FP}$ 

- Idéalement
  - Toutes les détections positives devraient correspondre à de vrais positifs : pas de faux négatifs (FN = 0), ou taux de vrais positifs = 1
  - Ce qui n'est pas détecté devrait correspondre aux seuls vrais négatifs : pas de faux positifs (FP = 0), ou taux de faux positifs = 0

## Courbes ROC pour discrimination entre 2 classes

- Modèle : en général décrit par un vecteur de paramètres w (par ex. poids connexions pour PMC) et un seuil b (par ex. sur la probabilité de la classe d'intérêt)
- Courbe ROC : taux de vrais positifs (en ordonnée) fonction du taux de faux positifs (en abscisse), la variable étant le seuil
- Pour un w fixé, peut-on réduire en même temps FN et FP en faisant varier le seuil?





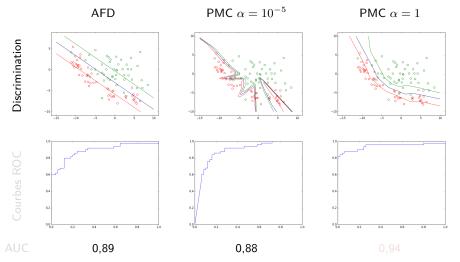
Frontières pour 3 valeurs du seuil de détection

Courbe ROC associée

si on augmente le taux de vrais positifs, le taux de faux positifs augmente également!

## Comparaison de modèles à travers les courbes ROC

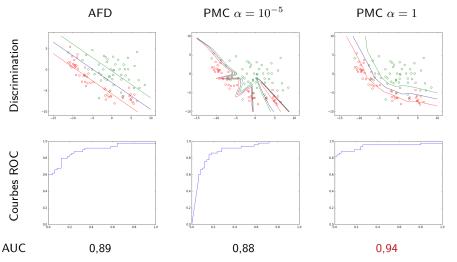
■ Comparaison globale par rapport au domaine de variation du seuil :



Courbes ROC

## Comparaison de modèles à travers les courbes ROC

■ Comparaison globale par rapport au domaine de variation du seuil :



# Comparaison de modèles à travers les courbes ROC (2)

- Un outil de comparaison globale est l'aire sous la courbe ROC (area under curve, AUC) : plus l'aire sous la courbe ROC est élevée, meilleur est le modèle
- Si valeurs AUC proches ou pour objectifs plus précis : comparaison des taux de vrais positifs (sensibilité) à taux de faux positifs (spécificité) donné(e)s

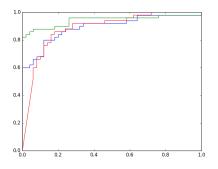


Fig. – Courbes ROC : AFD en bleu, PMC  $\alpha=10^{-5}$  en rouge, PMC  $\alpha=1$  en vert

### Plan du cours

- Objectifs et contenu de l'enseignemen
- 2 Organisation de l'enseignemen
  - Modélisation décisionnelle
  - Types de problèmes de décision
  - Modélisation à partir de données
  - Modélisation à partir de données : un cadre plus pr
  - Etapes générales
  - Choix d'une fonction de perte
  - Choix des familles paramétriques
  - Estimation du model
  - Comment mesurer la
  - Évaluation de mod
  - Validation crois
  - Sélection de modèles
    - Grid search pour le choix des hyperparamètres
    - Randomized parameter optimization

### Sélection de modèles

■ Dans l'estimation d'un modèle, par ex. par MRER

$$f_{\mathcal{D}_N}^* = \arg\min_{f \in \mathcal{F}} [R_{\mathcal{D}_N}(f) + \alpha G(f)]$$

- La pondération de la régularisation, α
- $\blacksquare$  Le critère de régularisation G(f)
- Autres paramètres qui ont un impact direct sur F, par ex. l'architecture pour un PMC le type de noyau (et la variance du noyau) pour une SVM, etc.
- Comment choisir de « bonnes » valeurs pour ces hyperparamètres
  - → Procédures de recherche qui explorent l'espace des valeurs des paramètres, de façor systématique ou aléatoire
  - Les modeles obtenus pour differentes valeurs des hyperparametres sont compares a travers leurs scores de validation croisée
  - Une fois trouvé le meilleur modèle, son risque espéré est estimé sur des données de test qui n'ont servi ni à la recherche des paramètres, ni à celle des hyperparamètres!

### Sélection de modèles

■ Dans l'estimation d'un modèle, par ex. par MRER

$$f_{\mathcal{D}_N}^* = \arg\min_{f \in \mathcal{F}} [R_{\mathcal{D}_N}(f) + \alpha G(f)]$$

- La pondération de la régularisation, d
- Le critère de régularisation G(f)
- le type de noyau (et la variance du noyau) pour une SVM, etc.
- Comment choisir de « bonnes » valeurs pour ces hyperparamètres
  - Procedures de recherche qui explorent i espace des valeurs des parametres, de façor systématique ou aléatoire
  - Les modeles obtenus pour differentes valeurs des hyperparametres sont compares a travers leurs scores de validation croisée
  - Une fois trouvé le meilleur modèle, son risque espéré est estimé sur des données de test qui n'ont servi ni à la recherche des paramètres, ni à celle des hyperparamètres!

### Sélection de modèles

■ Dans l'estimation d'un modèle, par ex. par MRER

$$f_{\mathcal{D}_N}^* = \arg\min_{f \in \mathcal{F}} [R_{\mathcal{D}_N}(f) + \alpha G(f)]$$

- $\blacksquare$  La pondération de la régularisation,  $\alpha$
- Le critère de régularisation G(f)
- le type de noyau (et la variance du noyau) pour une SVM, etc.
- Comment choisir de « bonnes » valeurs pour ces hyperparamètres
  - Procedures de recherche qui explorent i espace des valeurs des parametres, de façor systématique ou aléatoire
  - Les modeles obtenus pour differentes valeurs des hyperparametres sont compares a travers leurs scores de validation croisée
  - Une fois trouvé le meilleur modèle, son risque espéré est estimé sur des données de test qui n'ont servi ni à la recherche des paramètres, ni à celle des hyperparamètres!

### Sélection de modèles

■ Dans l'estimation d'un modèle, par ex. par MRER

$$f_{\mathcal{D}_N}^* = \arg\min_{f \in \mathcal{F}} [R_{\mathcal{D}_N}(f) + \alpha G(f)]$$

- $\blacksquare$  La pondération de la régularisation,  $\alpha$
- Le critère de régularisation G(f)
- le type de noyau (et la variance du noyau) pour une SVM, etc.
- Comment choisir de « bonnes » valeurs pour ces hyperparamètres
  - Procedures de recherche qui explorent i espace des valeurs des parametres, de façoi systématique ou aléatoire
  - Les modeles obtenus pour differentes valeurs des hyperparametres sont compares a travers leurs scores de validation croisée
  - Une fois trouvé le meilleur modèle, son risque espéré est estimé sur des données de test qui n'ont servi ni à la recherche des paramètres, ni à celle des hyperparamètres!

### Sélection de modèles

■ Dans l'estimation d'un modèle, par ex. par MRER

$$f_{\mathcal{D}_N}^* = \arg\min_{f \in \mathcal{F}} [R_{\mathcal{D}_N}(f) + \alpha G(f)]$$

- $\blacksquare$  La pondération de la régularisation,  $\alpha$
- Le critère de régularisation G(f)
- Autres paramètres qui ont un impact direct sur  $\mathcal{F}$ , par ex. l'architecture pour un PMC, le type de noyau (et la variance du noyau) pour une SVM, etc.
- Comment choisir de «bonnes» valeurs pour ces hyperparamètres
  - Procedures de recherche qui explorent i espace des valeurs des parametres, de façoi systématique ou aléatoire
  - Les modèles obtenus pour différentes valeurs des hyperparamètres sont comparés à travers leurs scores de validation croisée
  - Une fois trouvé le meilleur modèle, son risque espéré est estimé sur des données de test qui n'ont servi ni à la recherche des paramètres, ni à celle des hyperparamètres!

### Sélection de modèles

■ Dans l'estimation d'un modèle, par ex. par MRER

$$f_{\mathcal{D}_N}^* = \arg\min_{f \in \mathcal{F}} [R_{\mathcal{D}_N}(f) + \alpha G(f)]$$

- $\blacksquare$  La pondération de la régularisation,  $\alpha$
- Le critère de régularisation G(f)
- Autres paramètres qui ont un impact direct sur  $\mathcal{F}$ , par ex. l'architecture pour un PMC, le type de noyau (et la variance du noyau) pour une SVM, etc.
- Comment choisir de « bonnes » valeurs pour ces hyperparamètres ?
  - → Procédures de recherche qui explorent l'espace des valeurs des paramètres, de façon systématique ou aléatoire
  - Les modèles obtenus pour différentes valeurs des hyperparamètres sont comparés à travers leurs scores de validation croisée
  - Une fois trouvé le meilleur modèle, son risque espéré est estimé sur des données de test qui n'ont servi ni à la recherche des paramètres, ni à celle des hyperparamètres!

### Sélection de modèles

■ Dans l'estimation d'un modèle, par ex. par MRER

$$f_{\mathcal{D}_N}^* = \arg\min_{f \in \mathcal{F}} [R_{\mathcal{D}_N}(f) + \alpha G(f)]$$

- $\blacksquare$  La pondération de la régularisation,  $\alpha$
- Le critère de régularisation G(f)
- Autres paramètres qui ont un impact direct sur F, par ex. l'architecture pour un PMC, le type de noyau (et la variance du noyau) pour une SVM, etc.
- Comment choisir de « bonnes » valeurs pour ces hyperparamètres ?
  - → Procédures de recherche qui explorent l'espace des valeurs des paramètres, de façon systématique ou aléatoire
  - Les modèles obtenus pour différentes valeurs des hyperparamètres sont comparés à travers leurs scores de validation croisée
  - Une fois trouvé le meilleur modèle, son risque espéré est estimé sur des données de test qui n'ont servi ni à la recherche des paramètres, ni à celle des hyperparamètres!

### Sélection de modèles

■ Dans l'estimation d'un modèle, par ex. par MRER

$$f_{\mathcal{D}_N}^* = \arg\min_{f \in \mathcal{F}} [R_{\mathcal{D}_N}(f) + \alpha G(f)]$$

- $\blacksquare$  La pondération de la régularisation,  $\alpha$
- $\blacksquare$  Le critère de régularisation G(f)
- Autres paramètres qui ont un impact direct sur F, par ex. l'architecture pour un PMC, le type de noyau (et la variance du noyau) pour une SVM, etc.
- Comment choisir de « bonnes » valeurs pour ces hyperparamètres ?
  - → Procédures de recherche qui explorent l'espace des valeurs des paramètres, de façon systématique ou aléatoire
  - Les modèles obtenus pour différentes valeurs des hyperparamètres sont comparés à travers leurs scores de validation croisée
  - Une fois trouvé le meilleur modèle, son risque espéré est estimé sur des données de test qui n'ont servi ni à la recherche des paramètres, ni à celle des hyperparamètres!

#### Sélection de modèles

■ Dans l'estimation d'un modèle, par ex. par MRER

$$f_{\mathcal{D}_N}^* = \arg\min_{f \in \mathcal{F}} [R_{\mathcal{D}_N}(f) + \alpha G(f)]$$

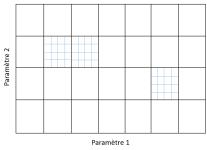
- $\blacksquare$  La pondération de la régularisation,  $\alpha$
- Le critère de régularisation G(f)
- Autres paramètres qui ont un impact direct sur F, par ex. l'architecture pour un PMC, le type de noyau (et la variance du noyau) pour une SVM, etc.
- Comment choisir de «bonnes» valeurs pour ces hyperparamètres?
  - → Procédures de recherche qui explorent l'espace des valeurs des paramètres, de façon systématique ou aléatoire
  - Les modèles obtenus pour différentes valeurs des hyperparamètres sont comparés à travers leurs scores de validation croisée
  - Une fois trouvé le meilleur modèle, son risque espéré est estimé sur des données de test qui n'ont servi ni à la recherche des paramètres, ni à celle des hyperparamètres!

## Recherche systématique : grid search

- Pour trouver les meilleures valeurs des hyperparamètres, une première possibilité est d'explorer l'espace des hyperparamètres de façon systématique
- Recherche en grille (grid search) :
  - 1. Définition d'intervalles et de pas de variation pour les hyperparamètres numériques (par ex. constante de régularisation  $\alpha$ , variance de noyau RBF)
  - 1. Définition d'ensembles de valeurs pour les hyperparamètres nominaux (par ex. architectures PMC, critères de régularisation, types noyaux SVM)
  - 2. Exploration systématique de l'espace des hyperparamètres
  - Choix des valeurs pour lesquelles le modèle obtenu présente les meilleures performances de validation croisée
- Estimation du risque espéré (erreur de généralisation) du modèle obtenu : sur des données non encore utilisées!

# Recherche systématique : grid search (2)

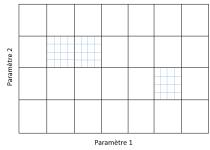
■ Lorsque seuls des hyperparamètres continus sont présents, on obtient une grille = combinaisons de valeurs à tester pour les m paramètres  $\Rightarrow$  grille de dimension m



- Tous les points de la grille peuvent être explorés en parallèle
- Plusieurs niveaux de « finesse » → recherche hiérarchique : exhaustive suivant la grille grossière, puis là où les résultats sont meilleurs on affine suivant le(s) niveau(x) plus fin(s) → augmentation du rapport qualité des résultats / coût

# Recherche systématique : grid search (2)

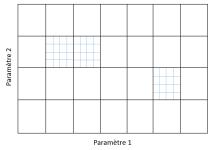
■ Lorsque seuls des hyperparamètres continus sont présents, on obtient une grille = combinaisons de valeurs à tester pour les m paramètres  $\Rightarrow$  grille de dimension m



- Tous les points de la grille peuvent être explorés en parallèle!
- Plusieurs niveaux de « finesse » → recherche hiérarchique : exhaustive suivant la grille grossière, puis là où les résultats sont meilleurs on affine suivant le(s) niveau(x) plus fin(s) → augmentation du rapport qualité des résultats / coût

## Recherche systématique : grid search (2)

■ Lorsque seuls des hyperparamètres continus sont présents, on obtient une grille = combinaisons de valeurs à tester pour les m paramètres  $\Rightarrow$  grille de dimension m



- Tous les points de la grille peuvent être explorés en parallèle!
- Plusieurs niveaux de «finesse»  $\rightarrow$  recherche hiérarchique : exhaustive suivant la grille grossière, puis là où les résultats sont meilleurs on affine suivant le(s) niveau(x) plus fin(s)  $\rightarrow$  augmentation du rapport qualité des résultats / coût

## Recherche aléatoire : randomized parameter optimization

- Des connaissances *a priori* permettent de privilégier certains intervalles de variation
  - ightarrow générer des valeurs conformes à ces connaissances ightarrow meilleure efficacité qu'avec grid search non hiérarchique
- Le coût peut être maîtrisé en fixant le nombre d'échantillons à générer
- Modalités d'échantillonnage
  - III Hyperparamètres numériques à valeurs continues (par ex.  $\alpha$ ) : loi d'échantillonnage (par ex. loi normale d'espérance et variance données)
  - Hyperparamètres numériques à valeurs discrètes (par ex. nombre de neurones cachés) : loi d'échantillonnage (par ex. loi uniforme sur intervalle donné)
  - $\blacksquare$  Hyperparamètres variables nominales : liste des valeurs (modalités) possibles  $\to$  loi uniforme sur ces valeurs
- Échantillons générés en considérant les hyperparamètres indépendants

## Évaluation et sélection de modèles : que faut-il retenir?

- Estimation du risque espéré (erreur de généralisation) sur des données non utilisées pour l'apprentissage
- Validation croisée : meilleure estimation qu'un seul découpage apprentissage | test
- Courbes ROC : comparaison plus globale de modèles de classement
- Meilleures valeurs pour les hyperparamètres : recherche systématique ou aléatoire, comparaison des modèles par validation croisée
- Si validation croisée employée pour sélectionner le meilleur modèle, estimation du risque espéré du modèle retenu sur des données non encore utilisées

#### Références I



O. Bousquet, S. Boucheron, and G. Lugosi.

*Introduction to Statistical Learning Theory*, volume Lecture Notes in Artificial Intelligence 3176, pages 169–207.

Springer, Heidelberg, Germany, 2004.



O. Chapelle, B. Schölkopf, and A. Zien, editors.

Semi-Supervised Learning.

MIT Press, Cambridge, MA, 2006.



I. Goodfellow, Y. Bengio, and A. Courville.

Deep Learning.

MIT Press, 2016.

http://www.deeplearningbook.org.



B. Schölkopf and A. Smola.

Learning with Kernels.

MIT Press, 2002.

#### Références II



L. Yang, H. Rodriguez, M. Crucianu, and M. Ferecatu.

Fully convolutional network with superpixel parsing for fashion web image segmentation.

In Proc. 23rd Intl. Conf. MultiMedia Modeling, Reykjavik, Iceland, pages 139–151, 2017.