

# Évaluer la précision des données d'entrée et du modèle résultant

Saeid Molladavoudi (Scientifique principal des données)  
Division de la science des données, Statistiques Canada

Une conférence du gouvernement du Canada

Le 18 février 2021



Éclairer grâce aux données pour bâtir un Canada meilleur.

# Considérations générales

- Chaque projet d'AA commence par une discussion avec les personnes concernées et les propriétaires des données.
- Les points suivants doivent être clairs dès le début :
  1. **Problème** : identifier le besoin des entreprises
  2. **Champ d'application** : Approches d'Apprentissage Automatique supervisées et non supervisées
  3. **Objectifs** : Classification, Régression, Regroupement ...
  4. **Contraintes** : Préoccupations concernant les données (collecte, montant, ...)
  5. **Produits livrables** : Solutions consommables pour les machines et solutions consommables pour l'homme
  6. **Critères de réussite** : Économie d'heures de travail manuel ou en temps de traitement, augmentation de la qualité, etc.

Cycle de vie du développement de l'apprentissage automatique



<https://www.jeremyjordan.me/ml-projects-guide/>

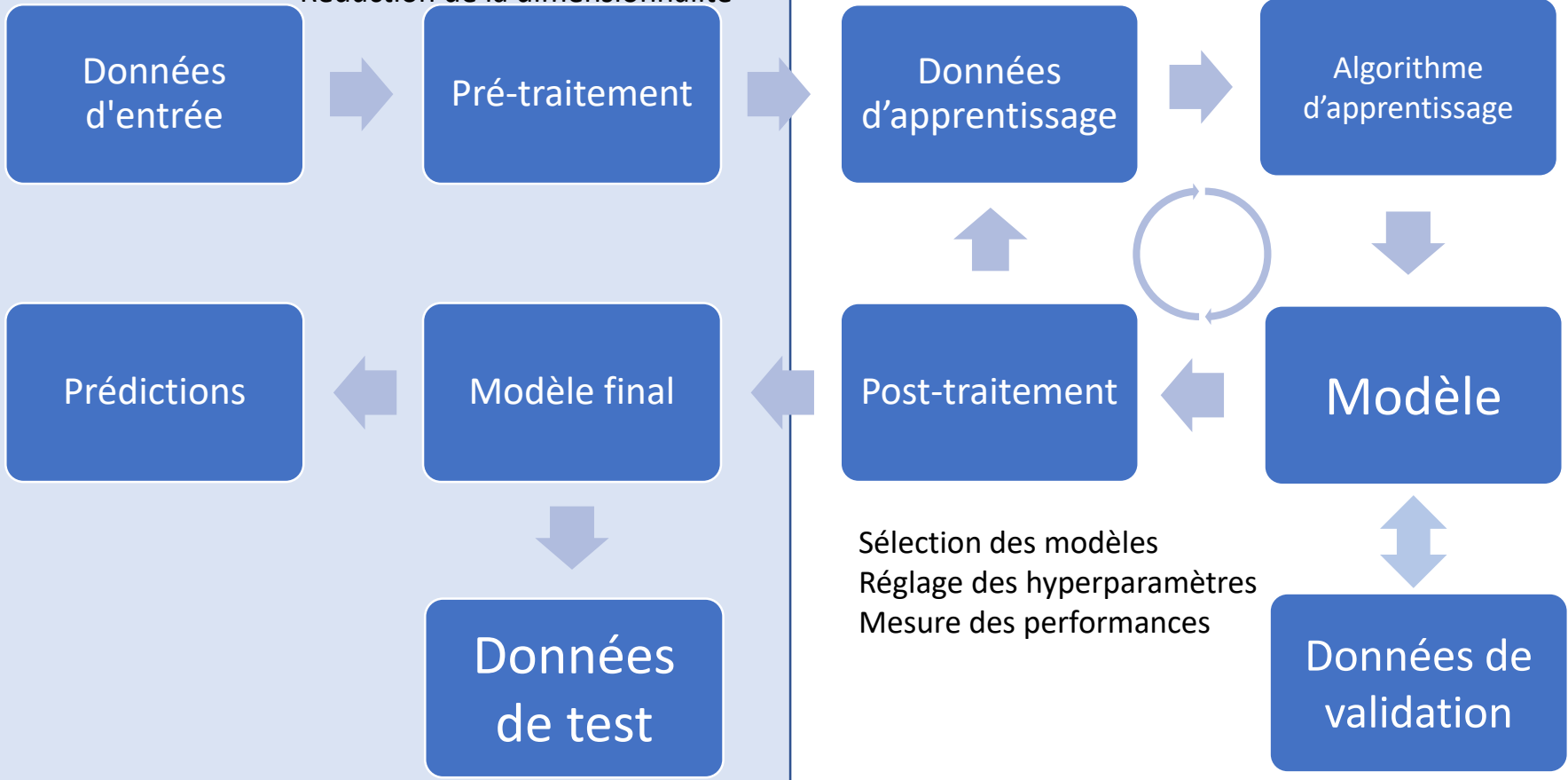
1

# Flux d'exploration de modèles supervisés

100

Nettoyage  
Données manquantes  
Extraction des caractéristiques  
Mise à l'échelle des caractéristiques  
Réduction de la dimensionnalité

Répartition des données(train/val/test)



Sélection des modèles  
Réglage des hyperparamètres  
Mesure des performances

# Les données d'entrée : la partie la plus importante d'un projet d'AA



<b>Principaux défis</b>	<ul style="list-style-type: none"><li>• Manque de données</li><li>• Mauvaise qualité</li><li>• Données d'apprentissage non représentatives de la situation</li><li>• Pas assez d'éléments pertinents</li><li>• Trop de caractéristiques non pertinentes</li></ul>
<b>Questions auxquelles nous devons répondre</b>	<ul style="list-style-type: none"><li>• Y a-t-il "suffisamment" de données ?</li><li>• Combien de sources de données existent ?</li><li>• Combien de données provenant de chaque source ?</li><li>• Est-il possible d'acquérir davantage de données ?</li><li>• Les données sont-elles structurées (par exemple via l'API) ou non structurées (par exemple des images) ?</li></ul>
<b>Par exemple, si le problème est un cas d'utilisation supervisée d'AA</b>	<ul style="list-style-type: none"><li>• Est-ce qu'il y a des labels ?</li><li>• Quelle est la quantité de données avec des labels ?</li><li>• Est-il possible d'acquérir d'autres labels ?</li></ul>

# Analyse exploratoire des données (pré-traitement)

A considérer	Détails
Types de données	Hétérogènes (continues, catégoriques, textes, images, ...)
Transformations	Ingénierie des caractéristiques, catégorisation à un seul coup, normalisation, intégration de mots, matrice de termes de documents, ...
Extraction	Pipeline s'il s'agit de données non structurées, par exemple vidéo, audio, pdfs, images ...
Labels	Classification (binaire, multi-classe, multi-label) vs labels de régression
Quantité de données	Existe-t-il "suffisamment" de données pour créer et entraîner, faire la validation des groupes d'essai ?
Déséquilibre des classes	Distribution des labels (déséquilibre de classe, évolution dans le temps)
Visuels	Outils visuels, par exemple diagrammes de dispersion et histogrammes
Statistiques de synthèse	En particulier pour les données structurées, par exemple pour trouver les corrélations entre les caractéristiques
Valeurs manquantes	Imputation ou retrait ?
Dimensionnalité	Réduction de la dimensionnalité des données d'entrée, par exemple l'analyse en composantes principales
Autres types de tâches	Réseaux, séries chronologiques, regroupement, détection des anomalies ....

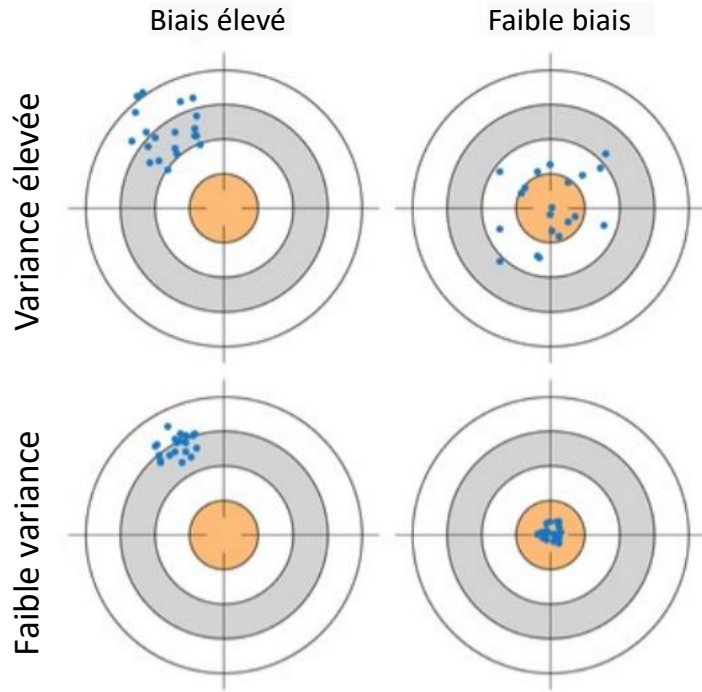
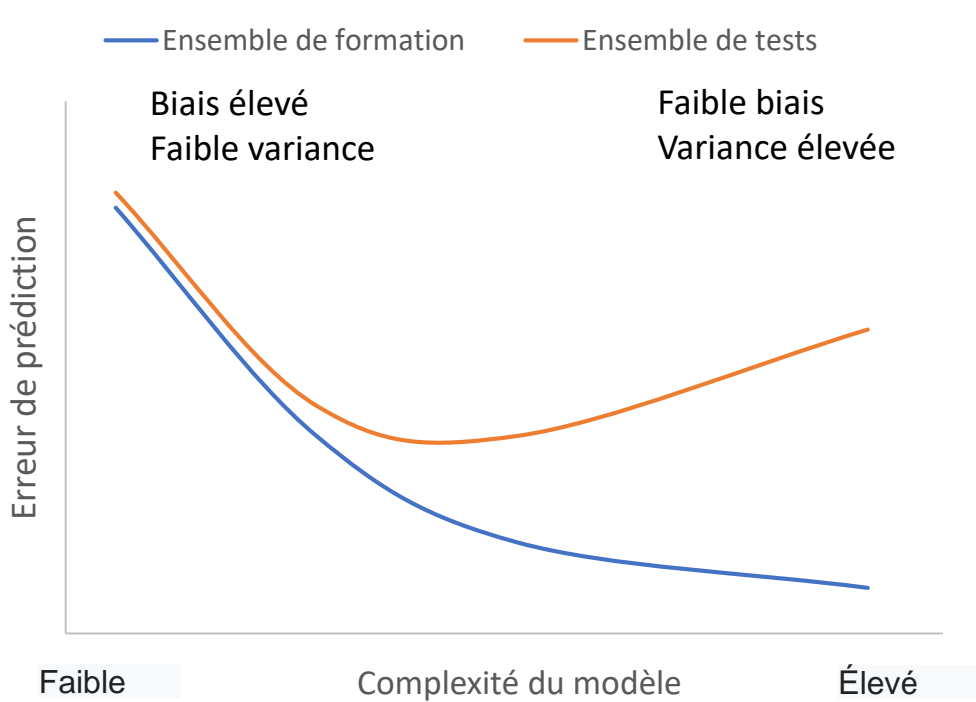
4

# Principe de formation, de validation et d'essai



Ensemble de données	Objectifs
<b>Ensemble de formation</b>	<ul style="list-style-type: none"><li>- L'objectif de l'AA est de minimiser la perte attendue (minimisation empirique du risque).</li><li>- Les modèles s'inspirent d'exemples de formation.</li></ul>
<b>Ensemble de validation</b>	<ul style="list-style-type: none"><li>- La validation est utilisée pour trouver les meilleurs paramètres du modèle (réglage des hyperparamètres).</li></ul>
<b>Ensemble de tests</b>	<ul style="list-style-type: none"><li>- Le test est utilisé pour simuler ce qui se passera lorsque nous appliquerons le modèle final à de nouvelles données invisibles provenant de la même distribution que les données de formation (c'est-à-dire le monde réel).</li><li>- Il donne une estimation non biaisée de l'erreur de généralisation du modèle.</li></ul>

# Compromis biais-variance



[https://www.variancejournal.org/articlespress/articles/Bias-Variance\\_Brady-Brockmeier.pdf](https://www.variancejournal.org/articlespress/articles/Bias-Variance_Brady-Brockmeier.pdf)

# Classification supervisée AA



Taille de l'ensemble de données	Labels	Stratégies d'atténuation
<b>Petit</b>	Peu de références pour les classes minoritaires.	<ul style="list-style-type: none"><li>• Utilisez la validation croisée en k.</li><li>• Ajoutez une régularisation au risque empirique pour éviter le sur-ajustement.</li><li>• Utiliser des modèles d'ensemble.</li><li>• Utiliser des méthodes d'augmentation des données.</li><li>• Transférer l'apprentissage.</li><li>• Réduction de la dimensionnalité.</li></ul>
<b>Moyen</b>	Peu de références pour les classes minoritaires.	<ul style="list-style-type: none"><li>• Utilisez la validation croisée en k.</li><li>• Ajoutez moins de régularisation.</li><li>• Utiliser des modèles d'ensemble.</li><li>• Utiliser des méthodes d'augmentation des données.</li><li>• Transférer l'apprentissage.</li><li>• Courbes d'apprentissage à rechercher en cas de sous-ajustement ou de surajustement.</li><li>• Outil d'annotation.</li></ul>
<b>Grand</b>	Petite quantité de données labellisées.	<ul style="list-style-type: none"><li>• Apprentissage semi-supervisé.</li><li>• Apprentissage actif.</li><li>• Outil d'annotation.</li></ul>



# Évaluation des performances

Tâche d'AA	Perte	Mesures
<b>Classification</b>	<ul style="list-style-type: none"><li>- Entropie croisée (binaire)</li><li>- Probabilité logarithmique négative (perte logarithmique)</li><li>- Perte de la charnière</li><li>- Divergence Kullback-Leibler</li><li>- Entropie relative</li><li>- Perte focale</li></ul>	<ul style="list-style-type: none"><li>- Exactitude</li><li>- Précision</li><li>- Rappeler</li><li>- F-score</li><li>- Courbe ROC</li><li>- AUC</li></ul>
<b>Régression</b>	<ul style="list-style-type: none"><li>- Erreur quadratique moyenne (EMS) ou perte quadratique</li><li>- (Lissage) Erreur Absolue Moyenne (EAM)</li><li>- Perte de bois</li></ul>	<ul style="list-style-type: none"><li>- EMS</li><li>- EAM</li></ul>
<b>Regroupement</b>	<ul style="list-style-type: none"><li>- somme des carrés à l'intérieur d'un groupe</li></ul>	<ul style="list-style-type: none"><li>- Indice Rand</li><li>- Score d'information mutuelle</li><li>- Coefficients de silhouette</li></ul>
<b>Classement</b>	<ul style="list-style-type: none"><li>- Perte de classement des marges</li><li>- Triple perte de marge</li></ul>	<ul style="list-style-type: none"><li>- Rang réciproque moyen</li><li>- Précision @ k</li><li>- Gain cumulé normalisé et actualisé</li></ul>

## Model Refinements:

- Utiliser la recherche par grille ou l'échantillonnage aléatoire pour le réglage des hyperparamètres.

# Assurance de la qualité

Éléments à prendre en considération	Détails
<b>Incertitude</b>	<ul style="list-style-type: none"><li>• Aléatoire : Bruit d'observation dans les données.</li><li>• Epistémique : Incertitude dans le modèle, les paramètres et la convergence.</li></ul>
<b>Biais dans les données</b>	<ul style="list-style-type: none"><li>• Biais humain dans la labellisation des données, par exemple, au sein d'un groupe, homogénéité et biais implicites.</li><li>• Le biais de collecte qui se réfère aux données qui ne reflètent pas les distributions réelles, par exemple le biais de convergence, le biais de déclaration, le biais de participation et le biais d'échantillonnage.</li></ul>
<b>Sécurité et vie privée</b>	<ul style="list-style-type: none"><li>• Les technologies de protection de la vie privée et leurs applications dans le domaine de l'AA, par exemple le cryptage homomorphe, les calculs multipartites sécurisés, ...</li><li>• Confidentialité des entrées et des sorties (confidentialité), par exemple, confidentialité différentielle</li></ul>
<b>Reproductibilité</b>	<ul style="list-style-type: none"><li>• Contrôle de la version du code (développement, test, mise en œuvre) et des données</li><li>• Pipeline automatisé</li></ul>
<b>Explicabilité</b>	<ul style="list-style-type: none"><li>• Transparence</li><li>• Documentation</li></ul>
<b>Maintenance</b>	<ul style="list-style-type: none"><li>• Surveiller la cohérence, l'exactitude, la couverture et la représentativité des données d'apprentissage dans le temps afin de détecter toute dérive ou tout dérapage de la qualité</li></ul>



# Merci/Thank you!

[saeid.molladavoudi@Canada.ca](mailto:saeid.molladavoudi@Canada.ca)