

Interpréter et explorer

Approfondissement en analyse de données avec R

Hubert Cadieux



Plan de l'atelier

9h à 12h

- Qui suis-je?
- R et RStudio

- Analyse textuelle
- PCA et clustering

12h à 15h

- Atelier libre avec des mentors de la CLESSN
- Exercices avec R

Qui suis-je?

- Maîtrise en science politique
- Membre de la Chaire de leadership en enseignement des sciences sociales numériques (CLESSN)
- Membre du Centre d'analyse des politiques publiques (CAPP)
- Membre du Groupe de recherche en communication politique (GRCP)
- Membre du Centre d'étude sur la citoyenneté démocratique (CECD)

La CLESSN



enseignement des sciences sociales numériques (CLESSN)

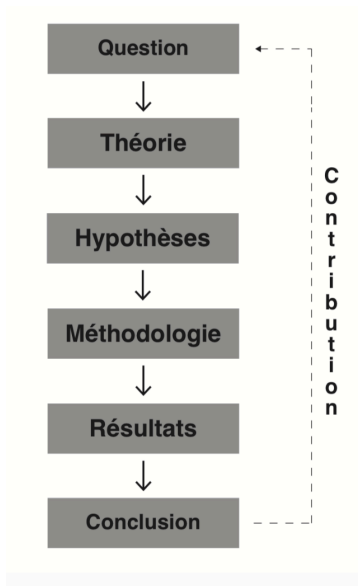
Chaire de leadership en

- Titulaire: Yannick Dufresne
- Objectif: Développer des compétences en science des données pour les étudiants en sciences sociales à travers des projets de recherche et des formations

Pourquoi suivre cet atelier?

- Base solide pour la recherche académique

Permet également de comprendre la pertinence de l'utilisation d'un logiciel de traitement de données et d'obtenir des outils pour dédramatiser la science des données



Pourquoi suivre cet atelier?

Comprendre le monde autour de nous à travers des données





Les données sont partout, mais pas nécessairement compréhensible dans leur forme « raw ». L'analyse de données permet de simplifier ces données pour comprendre et expliquer des phénomènes ou des comportements dans le monde qui nous entoure.

Pouvoir de prédiction: prédire comment certaines variables ou facteurs influencent d'autres variables ou événements.

Qu'est-ce que l'analyse de données?

L'analyse de données vise à explorer, décrire et interpréter des informations recueillies sur les comportements humains, les relations sociales et les phénomènes sociétaux. **Ce processus peut inclure des données quantitatives (comme des données issus de sondages) et des données qualitatives (comme des entretiens ou des observations).**

Principales étapes de l'analyse de données :

1. Collecte des données
2. Préparation des données
3. Analyse statistique

Les exercices de cet atelier

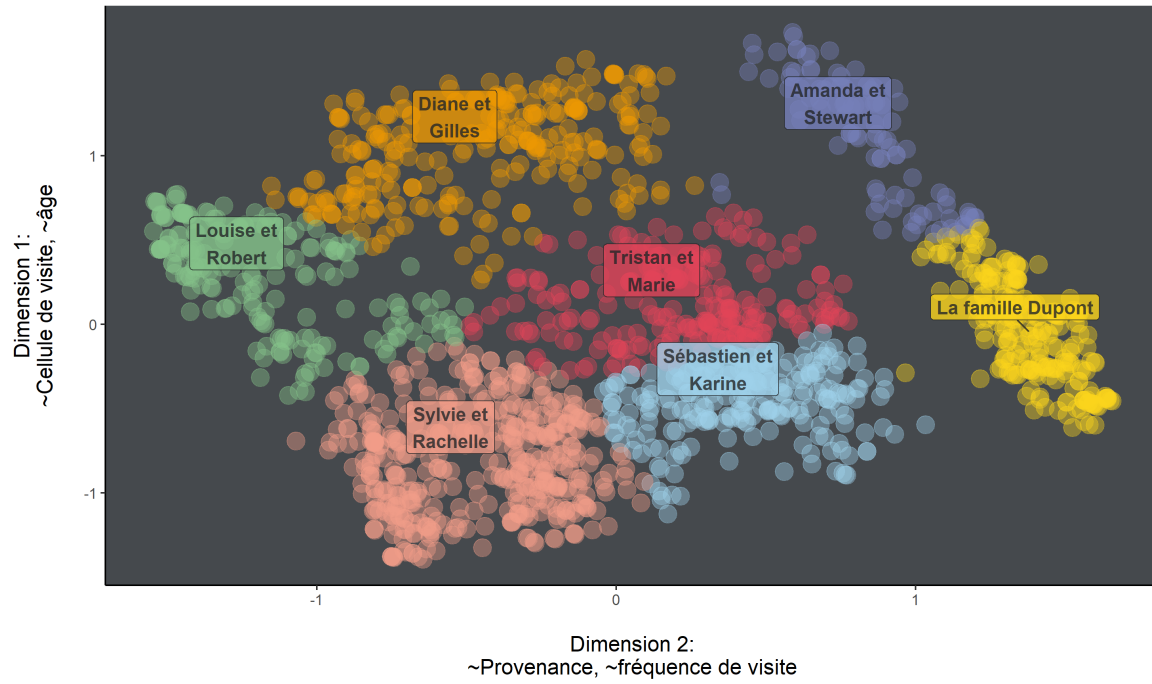
- Analyse textuelle
- PCA et clustering

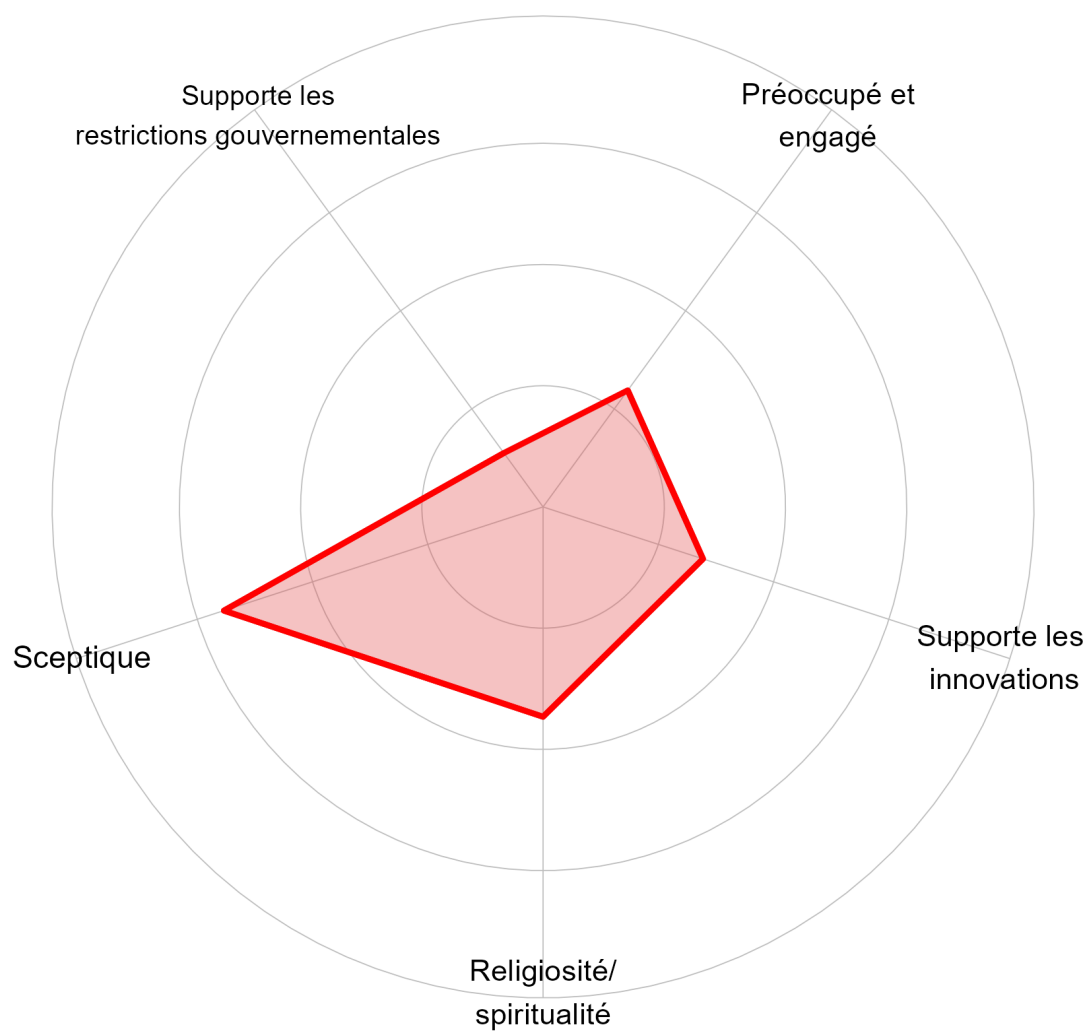
Voici quelques exemples de ce que vous pourrez faire avec R

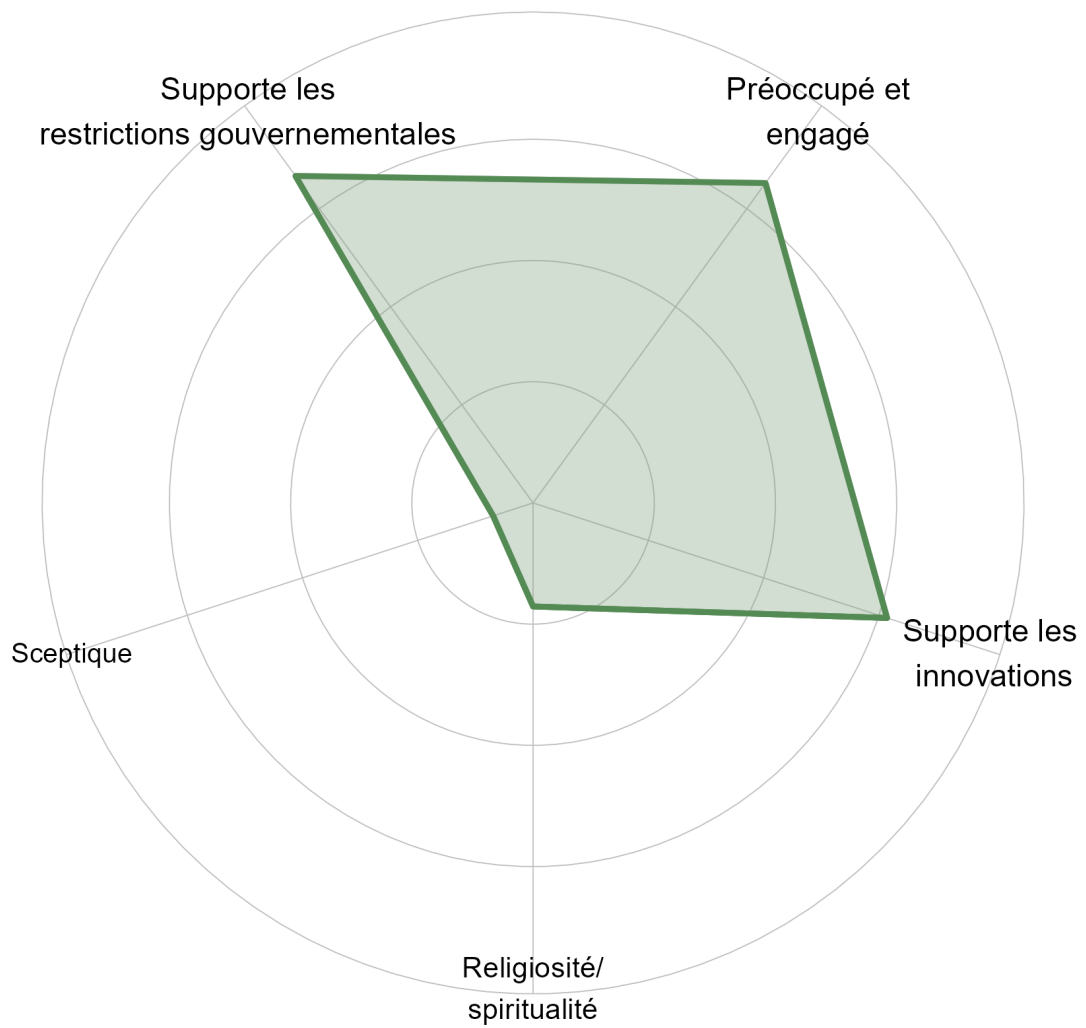


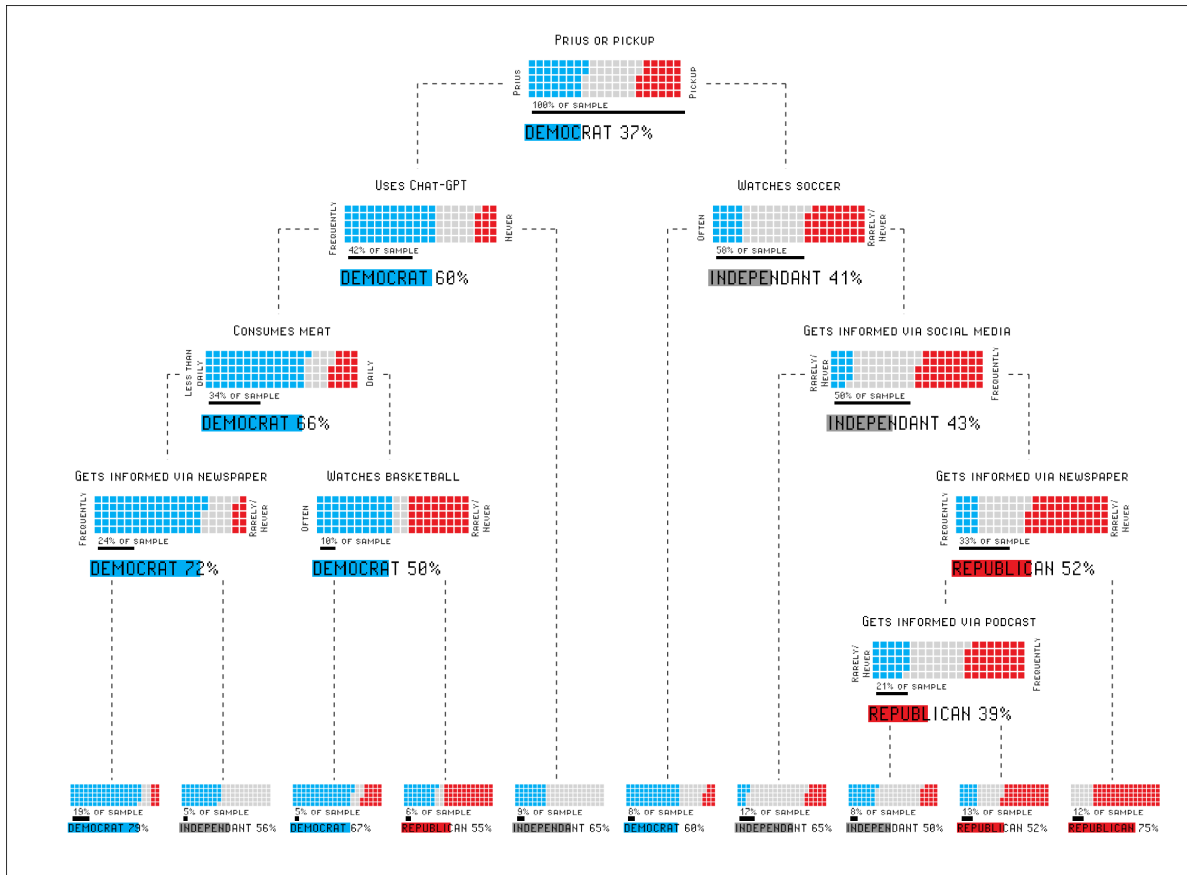
Diamonds indicate the parties' positions on the issue.
Vertical lines indicate the 95% confidence interval.

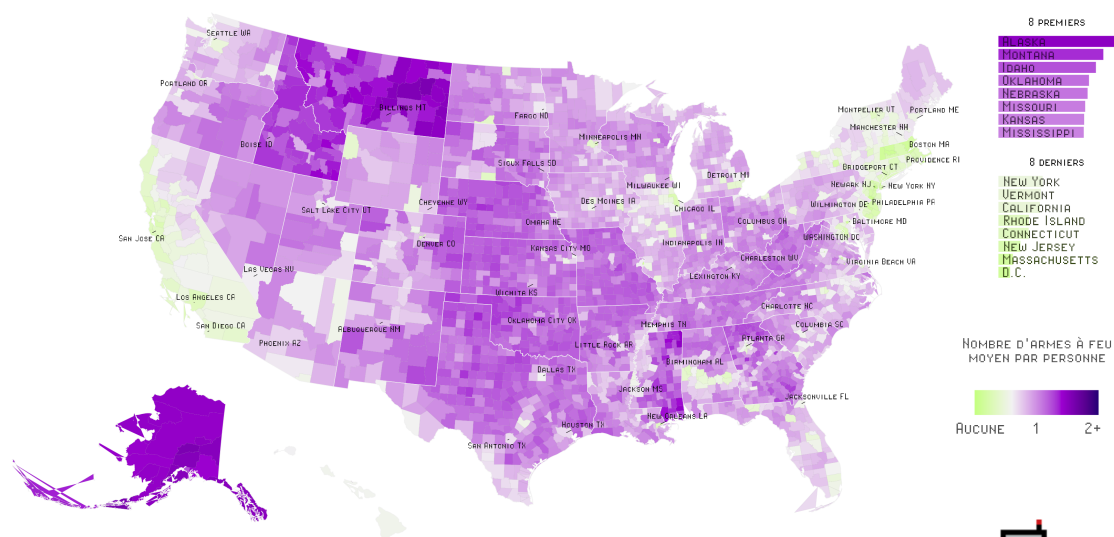
##







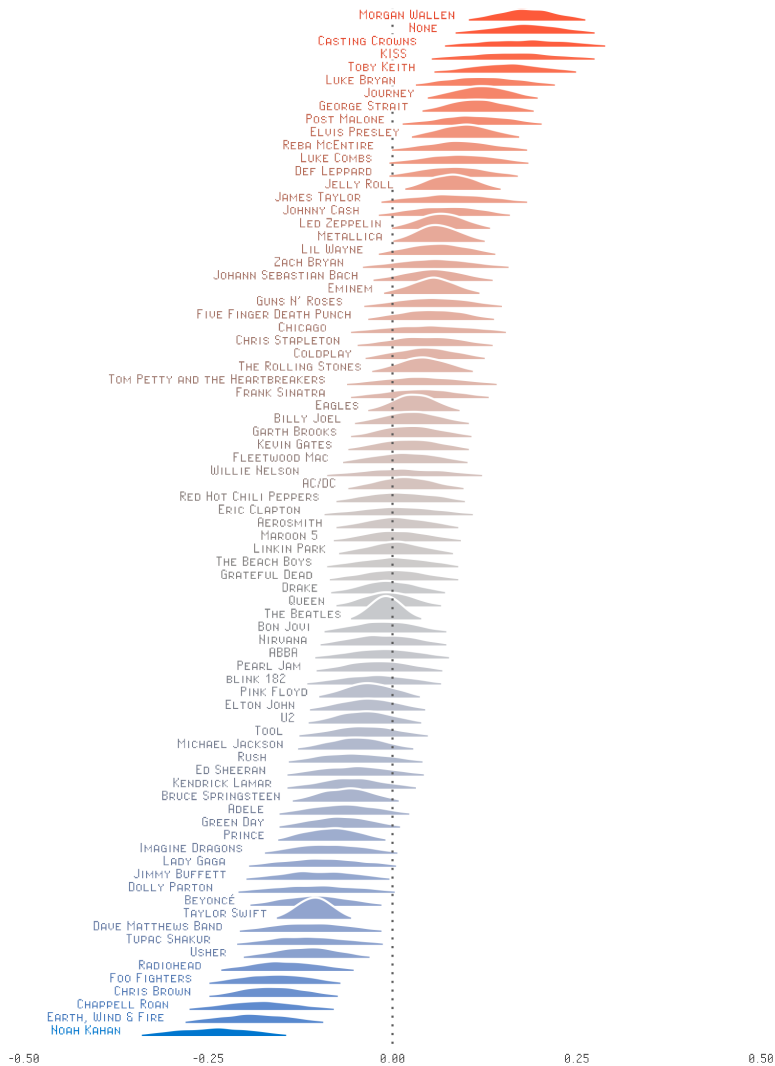




À PARTIR DES DONNÉES DE L'APPLICATION DATAGOTCHI (N = 14 329), LE NOMBRE D'ARMES À FEU PAR PERSONNE A ÉTÉ PRÉDIT POUR CHAQUE COMTÉ EN UTILISANT UNE RÉGRESSION MULTINIVEAU ET UNE POST-STRATIFICATION (MRP). LE DÉGRADÉ DE COULEUR SUR LA CARTE REFLÈTE LA PRÉDICTION DU NOMBRE D'ARMES D'UN RÉPONDANT AGRÉGÉE POUR CHAQUE COMTÉ. CETTE VALEUR EST BASÉE SUR LES PRÉDICTIONS DU MODÈLE DE RÉGRESSION MULTINIVEAU.



L'EFFET DE L'ARTISTE PRÉFÉRÉ SUR LE VOTE AUX ÉLECTIONS AMÉRICAINES 2024



INTERVALLE DE CRÉDIBILITÉ DU COEFFICIENT

CET INTERVALLE DE CRÉDIBILITÉ BAYÉSIEN ILLUSTRÉ LA PLAGE DES VALEURS PROBABLES DU COEFFICIENT ASSOCIÉ À CHAQUE ARTISTE CALCULÉE À PARTIR DE LA DISTRIBUTION POSTÉRIEURE D'UN MODÈLE DE RÉGRESSION BAYÉSIEN. LE PRIOR A ÉTÉ DÉFINI SUR LA BASE D'UNE ENQUÊTE D'EXPERTS REGROUPANT 12 ÉTUDIANTS ET PROFESSEURS EN SCIENCE POLITIQUE. LES COEFFICIENTS SONT MESURÉS SUR L'ÉCHELLE SUIVANTE : STRONG DEMOCRAT, SOFT DEMOCRAT, INDEPENDENT, SOFT REPUBLICAN, STRONG REPUBLICAN, PAR INCRÉMENTS DE 0.25. UN COEFFICIENT DE 0.25 INDIQUE QU'UN ARTISTE PROGRESSE D'UN NIVEAU COMPLET SUR CETTE ÉCHELLE, TOUTES CHOSSES ÉTANT ÉGALES PAR AILLEURS, DANS UN MODÈLE AJUSTÉ POUR LE STATUT SOCIO-ÉCONOMIQUE (SES) ET CERTAINS INDICATEURS DE STYLE DE VIE.

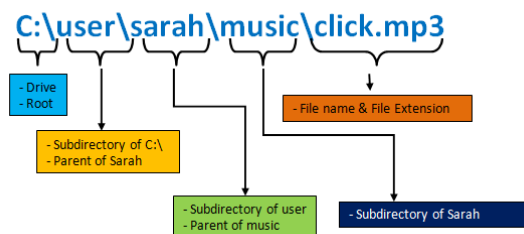


Les outils utilisés dans cet atelier

- R et RStudio
- C'est quoi la différence entre R et RStudio?
- R est le moteur, RStudio est l'interface
- Quarto
- GitHub

Concept important : Le chemin d'arborescence

- À tout moment vous devez savoir où vous êtes dans votre ordinateur pour pouvoir importer des données, exporter des graphiques ou mettre des fichiers en relation
- Votre R est toujours ouvert dans un dossier, et donc vous devez savoir où il est pour pouvoir importer des données
- La fonction `getwd()` dans R permet de savoir où vous êtes
- Un fichier quarto se réfère toujours à lui même



Concept important : Planifier avant de coder

La plus grosse erreur est de commencer à coder sans savoir ce que vous voulez faire

- Clarifier vos objectifs: Qu'est-ce que vous voulez faire?
 - Nettoyer des données?
 - Faire un graphique?

Les possibilités sont infinies, donc il est important de savoir où vous voulez aller

Concept important : Décomposer le problème

- Une fois que vous savez ce que vous voulez faire, il est important de décomposer le problème en petites étapes
- Un script R pour une seule tâche
 - Bien nommer vos scripts pour savoir ce qu'ils font
 - Exemples:
 - * nettoyage_donnees.R
 - * graphique.R
- Chaque script doit être clair et facile à comprendre
- Commenter votre code avec des #

Bonnes pratiques

- Nommer vos objets de façon explicite
- Commenter votre code
- Organiser votre code

Organiser votre répertoire et placer vos données

```
/MonProjet
data/
  raw_data.csv
  cleaned_data.csv
scripts/
  analyse.R
  visualization.R
results/
  summary_statistics.csv
plots/
  data_distribution.png
docs/
  methodology.md
  references.bib
  project_presentation.qmd
README.md
```

Apprendre plus de R

- swirl
- Datacamp
- R4DS (R for Data Science)
- Advanced R

Ressources

Quoi faire quand ça ne fonctionne pas?

- Google (Stack Overflow)
- ChatGPT / Claude
- Réessayer ChatGPT / Claude
- La documentation de R

Codons!

