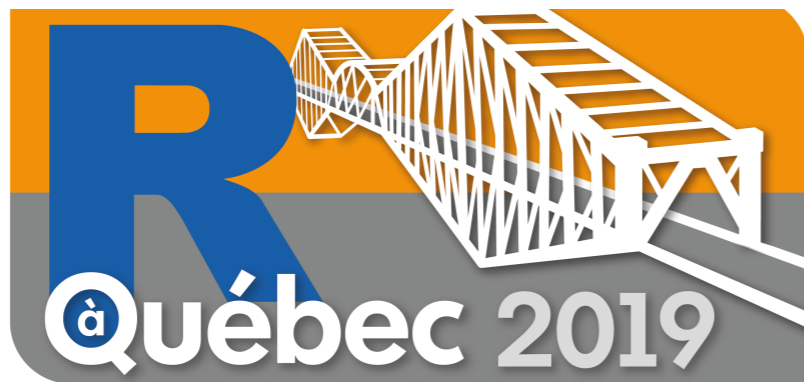


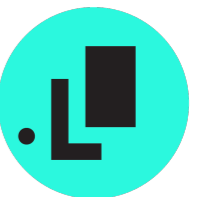
Modélisation prédictive avec R dans un contexte de production - De l'extraction au déploiement



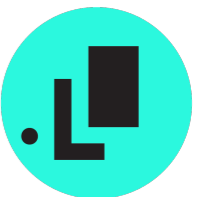
R à Québec 2019

Bienvenue

- Mettre nos noms, experiences, etc (pourrait être funny un peu les descriptions)

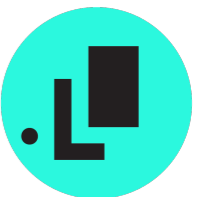


Objectifs de l'atelier



Informations

- Expliquer le concept du livre (ressource pour eux)
- Lien vers le repo GitHub de set up



Une première fois !

- Grande première planétaire

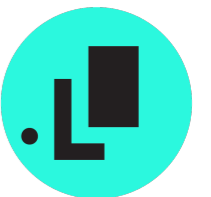


- On est pas orgueilleux! On veut du feedback.
- Disclaimer : Jeu de données simple (mais gros!)
 - Emphase sur le processus

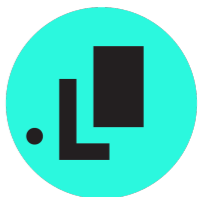


Processus de modélisation

- Parler du concept de back-and-forth
- Pour atelier on va s'en tenir à une ligne toute tracée



BIXI



Problématique

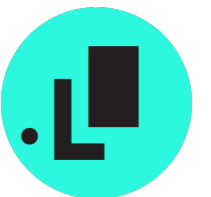
- Pour des raisons inconnues, on nous demande :
 1. Quels utilisateurs sont susceptibles de revenir à la même station?
 2. Quelle sera la durée du trajet d'un utilisateur?



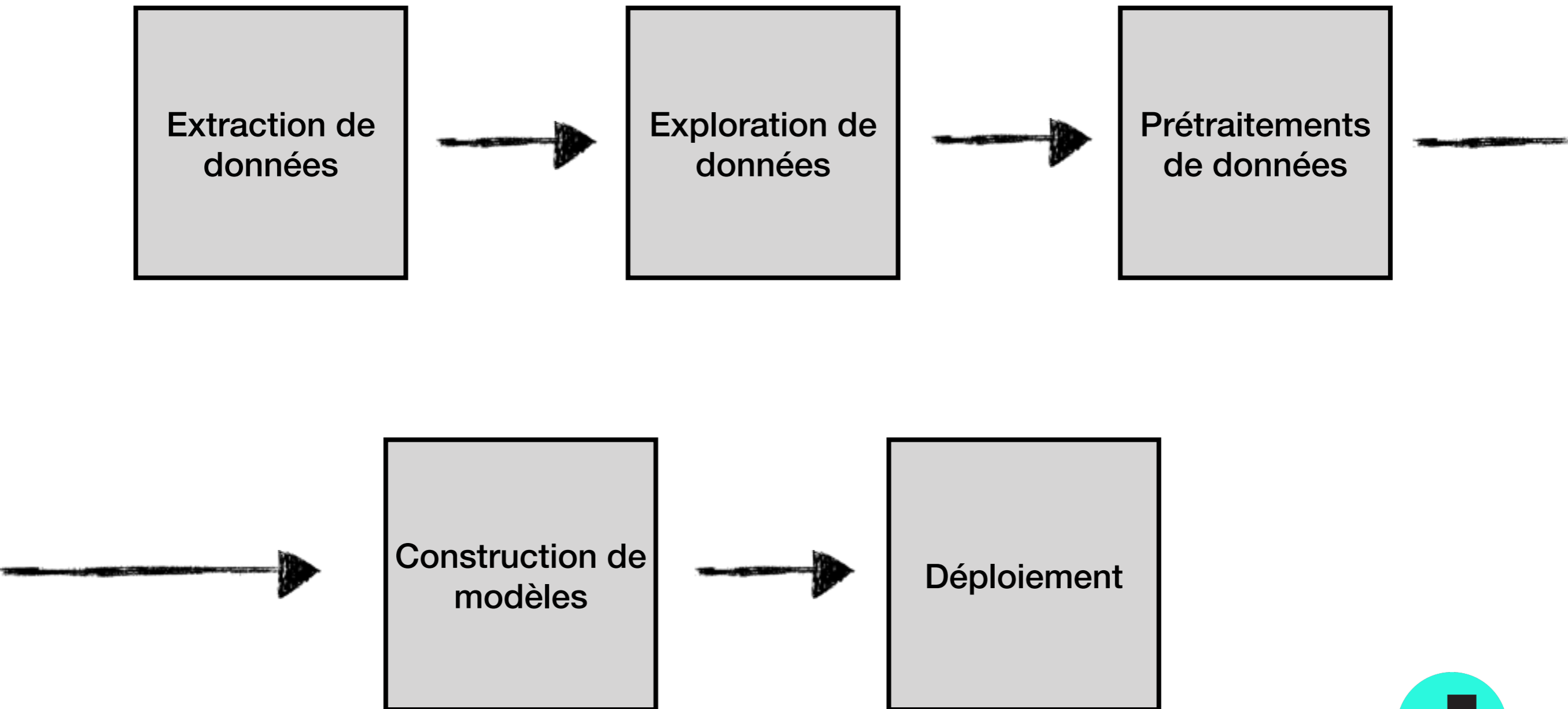
Définition de la tâche

Régression vs Classification

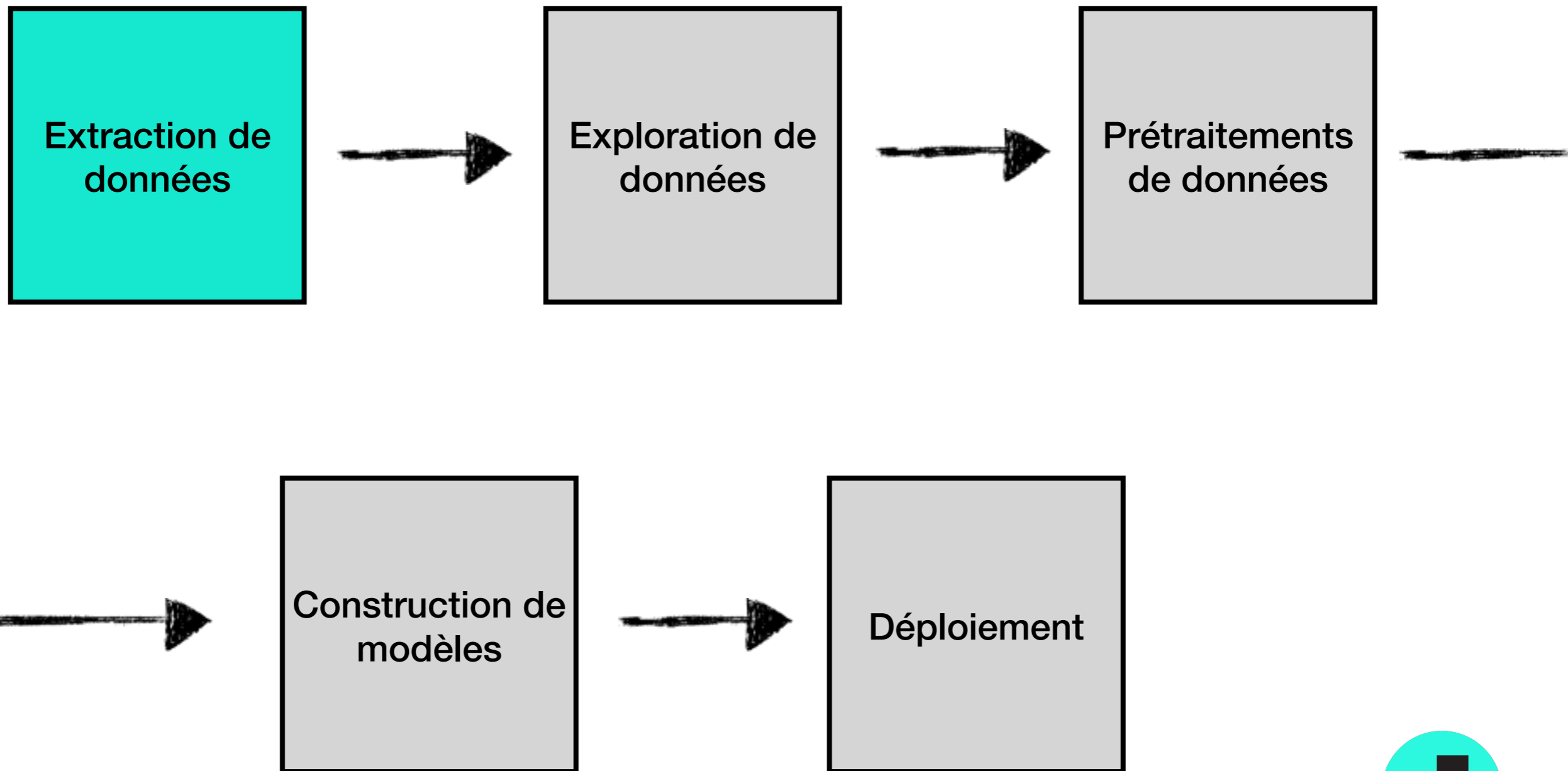
$$y \approx f(x)$$



Étapes à accomplir

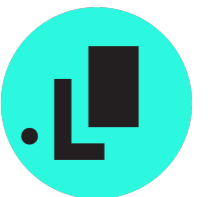


Collectons !

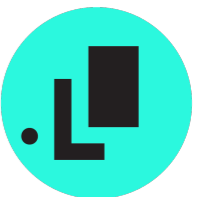
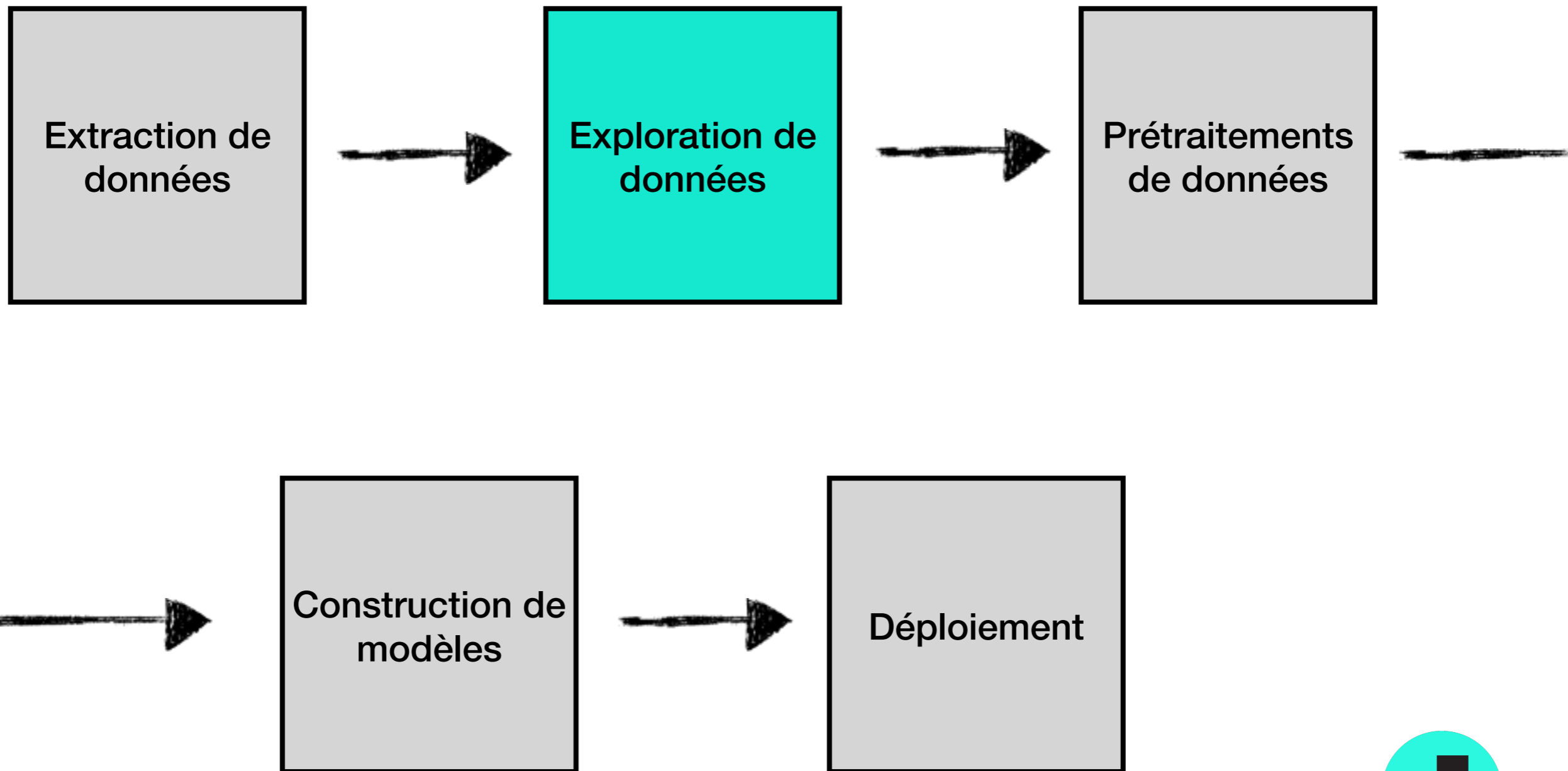


Extraction de données

- Lau



Explorons !

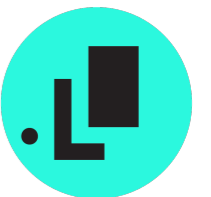


Exploration des données



Objectifs

1. S'appropriier le jeu de données;
2. Suggérer des transformations pertinentes pour le prétraitement des données.



Appropriation

Pour chaque variable d'entrée...

1. Observer la distribution;
2. Observer la corrélation avec les autres variables d'entrée;
3. Observer l'effet unidimensionnel sur la variable réponse;
4. Observer les effets multidimensionnels avec les autres variables d'entrée sur la variable réponse.



Transformations

- Exclusion;
- Identité;
- Regroupements;
- Tout autre fonction.

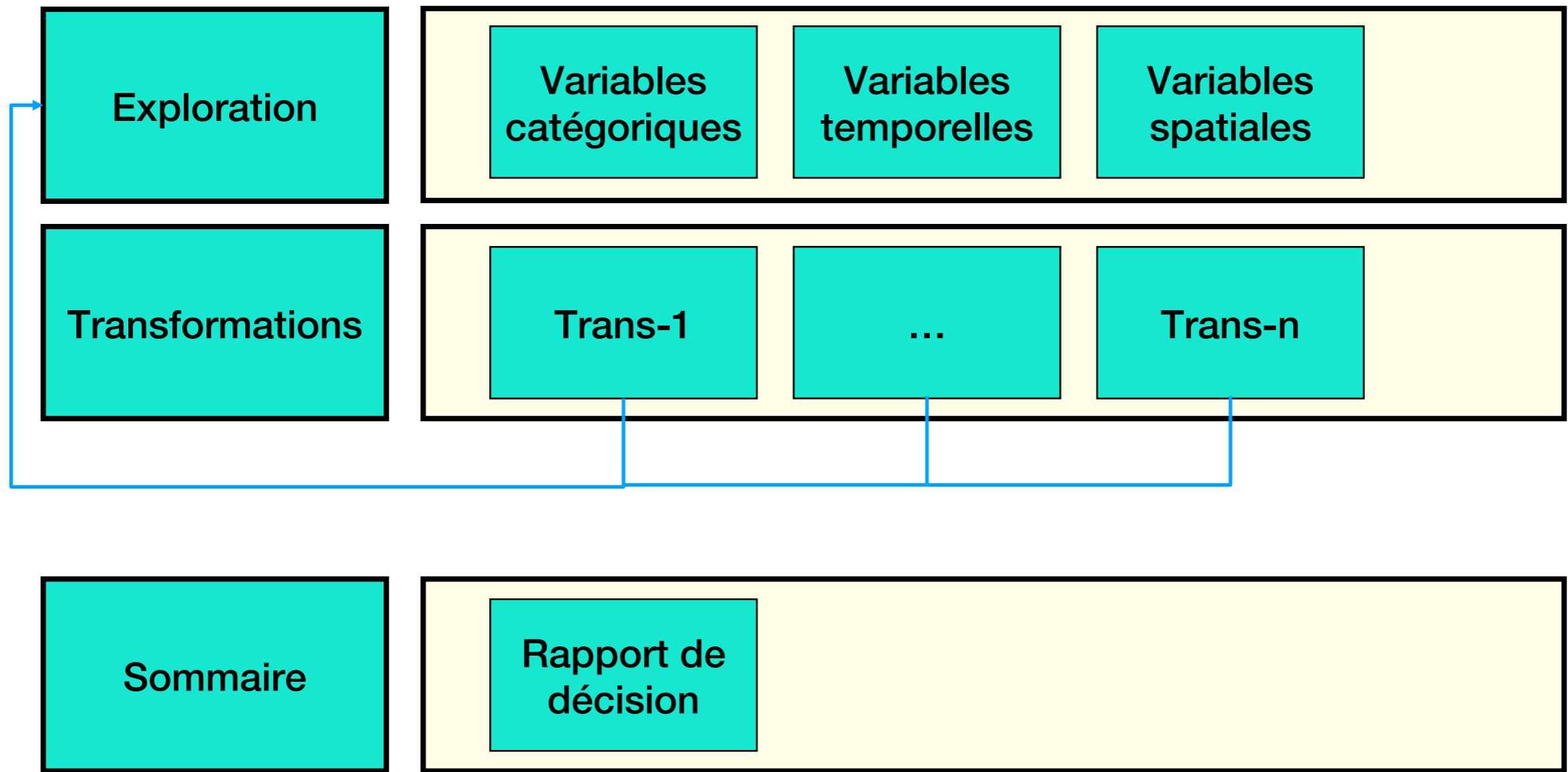


Variables d'entrée

- Numérique;
- Catégorique;
- Temporelle;
- Spatiale.



Récits à implémenter



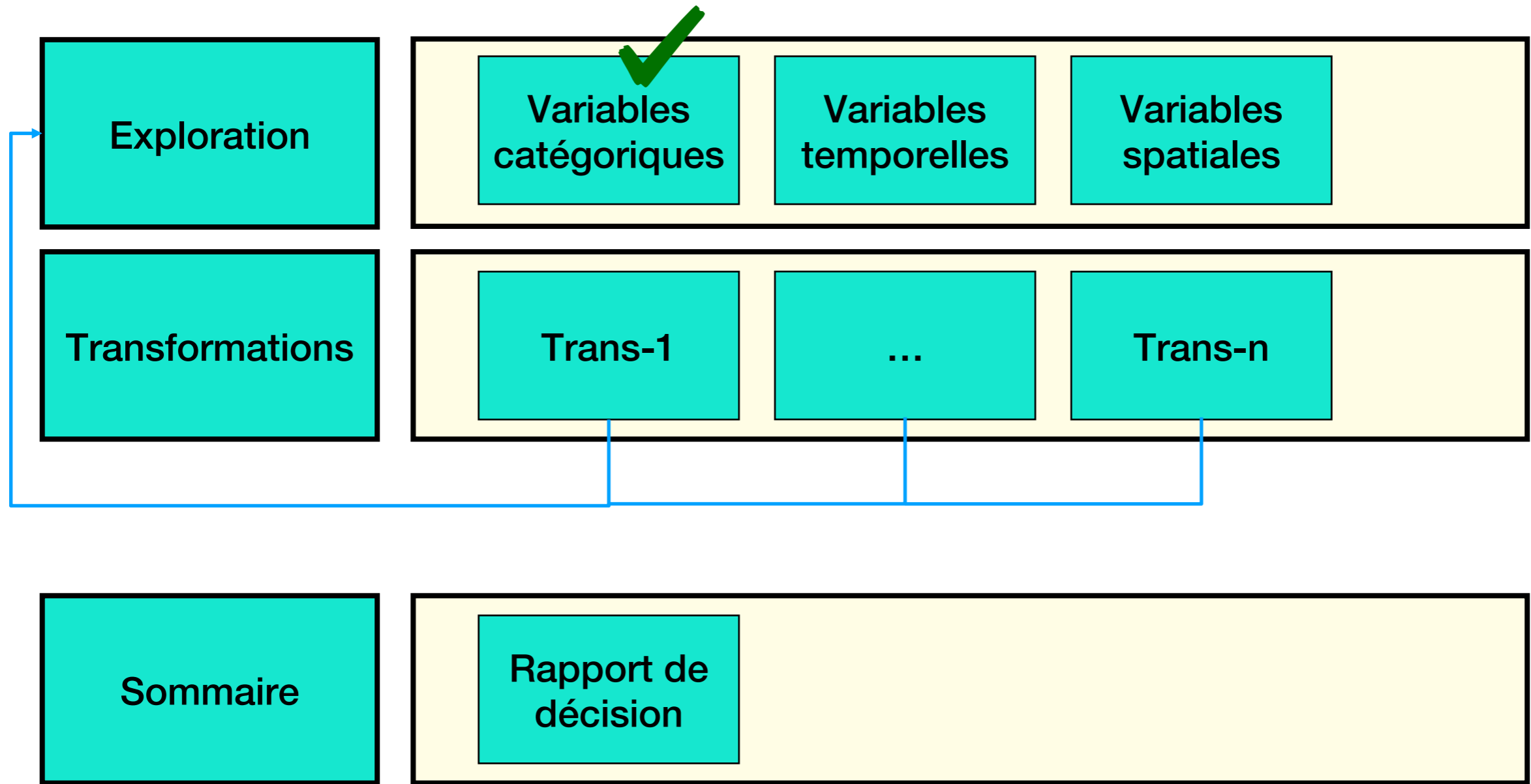
Nom du récit	Variables catégoriques
ID	1
Section	Exploration des données

Description
<p>Créer des visuels pour comprendre les distributions des variables d'entrée</p> <p><u>Packages proposés</u></p> <p>Agrégation des données :</p> <ul style="list-style-type: none"> • base • dplyr • data.table <p>Packages proposés pour les graphiques :</p> <ul style="list-style-type: none"> • base • ggplot2 • plotly • lattice

Entrée
<ol style="list-style-type: none"> 1. Données brutes 2. Données utilisées pour le prétraitement 3. Données externes
Sortie
<ol style="list-style-type: none"> 1. Certains des graphiques suivants : <ul style="list-style-type: none"> - Distribution - Corrélations - One-way - Interactions



Récits à implémenter



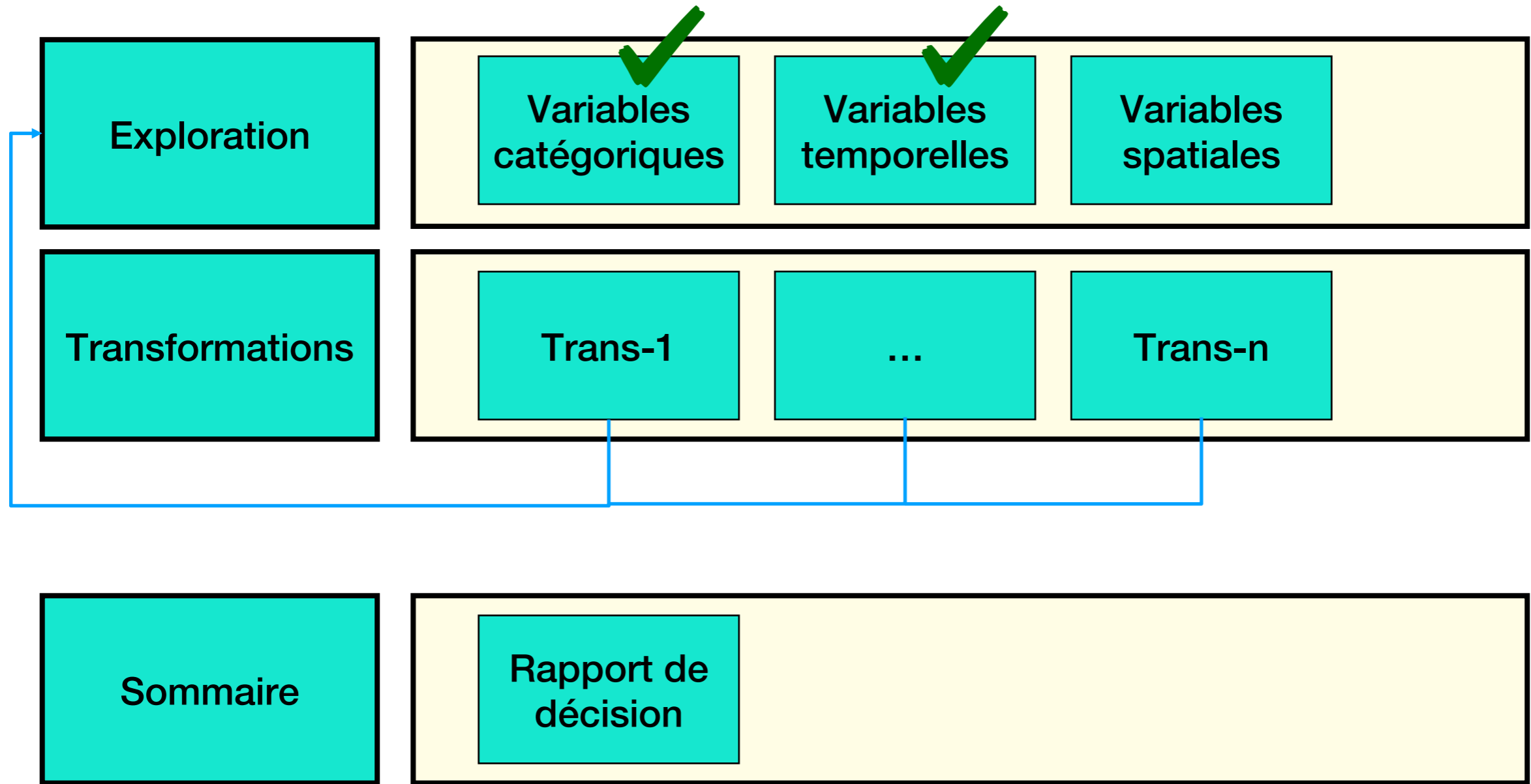
Nom du récit	Variables temporelles
ID	2
Section	Exploration des données

Description
Créer des visuels pour comprendre les distributions des variables d'entrée

Entrée
<ol style="list-style-type: none"> 1. Données brutes 2. Données utilisées pour le prétraitement 3. Données externes
Sortie
<ol style="list-style-type: none"> 1. Certains des graphiques suivants : <ul style="list-style-type: none"> - Distribution - Corrélations - One-way - Interactions



Récits à implémenter



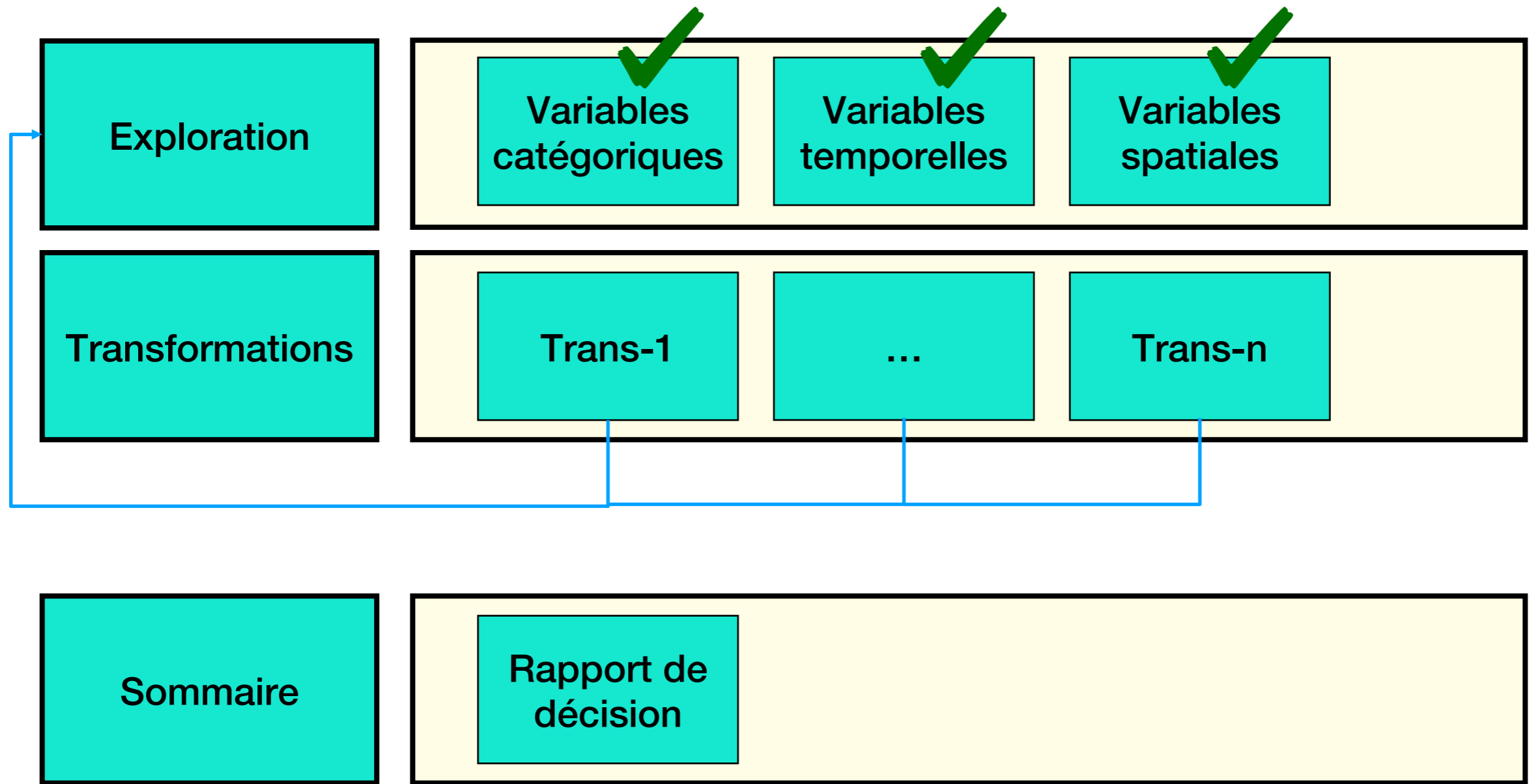
Nom du récit	Variables spatiales
ID	3
Section	Exploration des données

Description
<p>Créer des visuels pour comprendre les distributions des variables d'entrée</p> <p><u>Packages proposés</u></p> <p>Création de cartes :</p> <ul style="list-style-type: none"> • ggmap • plotly • leaflet • tmap

Entrée
<ol style="list-style-type: none"> 1. Données brutes 2. Données utilisées pour le prétraitement 3. Données externes
Sortie
<ol style="list-style-type: none"> 1. Certains des graphiques suivants : <ul style="list-style-type: none"> - Distribution - Corrélations - One-way - Interactions



Récits à implémenter



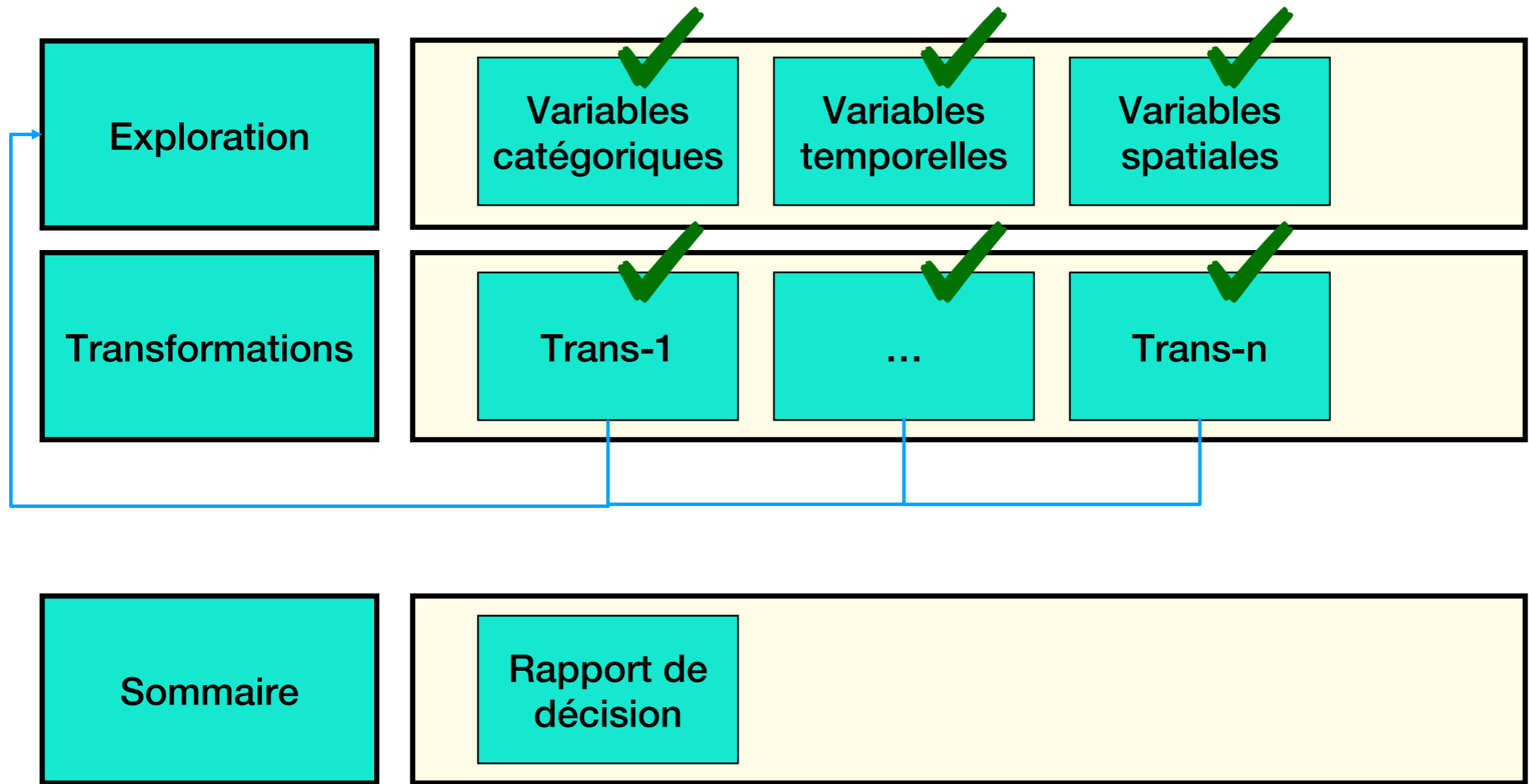
Nom du récit	Transformations
ID	4
Section	Exploration des données

Description
Analyser des transformations potentielles suite à l'analyse des variables disponibles dans la collecte des données

Entrée
<ol style="list-style-type: none"> 1. Données brutes 2. Données utilisées pour le prétraitement 3. Données externes
Sortie
<ol style="list-style-type: none"> 1. Certains des graphiques suivants : <ul style="list-style-type: none"> - Distribution - Corrélations - One-way - Interactions



Récits à implémenter



Nom du récit	Rapport de décision
ID	5
Section	Exploration des données

Description
Créer un rapport de décision concernant les transformations suggérées afin que l'étape du prétraitement des données puisse incorporer l'information

Entrée
1. Les graphiques créés depuis le début de la section
Sortie
1. Un rapport de recommandations



Recommandations (1/2)

- Moments de la journée
 - Matin : 6h à 11h
 - Journée : 11h à 16h
 - Soir : 16h à 23h
 - Nuit : 23h à 6h
- Semaine/Fin de semaine
 - Semaine : Lundi, Mardi, Mercredi, Jeudi, Vendredi
 - Fin de semaine : Samedi, Dimanche



Recommandations (2/2)

- Regroupement des quartiers
 - Groupe 1 : Plateau-Mont-Royal
 - Groupe 2 : Ville-Marie
 - Groupe 3 : Ahuntsic-Cartierville, Villeray-Saint-Michel-Parc-Extension, Rosemont-La Petite-Patrie, Mercier-Hochelaga-Maisonneuve
 - Groupe 4 : Outremont, Côte-des-Neiges-Notre-Dame-de-Grâce, Westmount, Le Sud-Ouest, Verdun, LaSalle
 - Groupe 5 : Autre



Données externes à explorer

- Météo
- FSAs au lieu de quartiers
- Réseau de Métro
- Jours fériés
- Dates de festivals

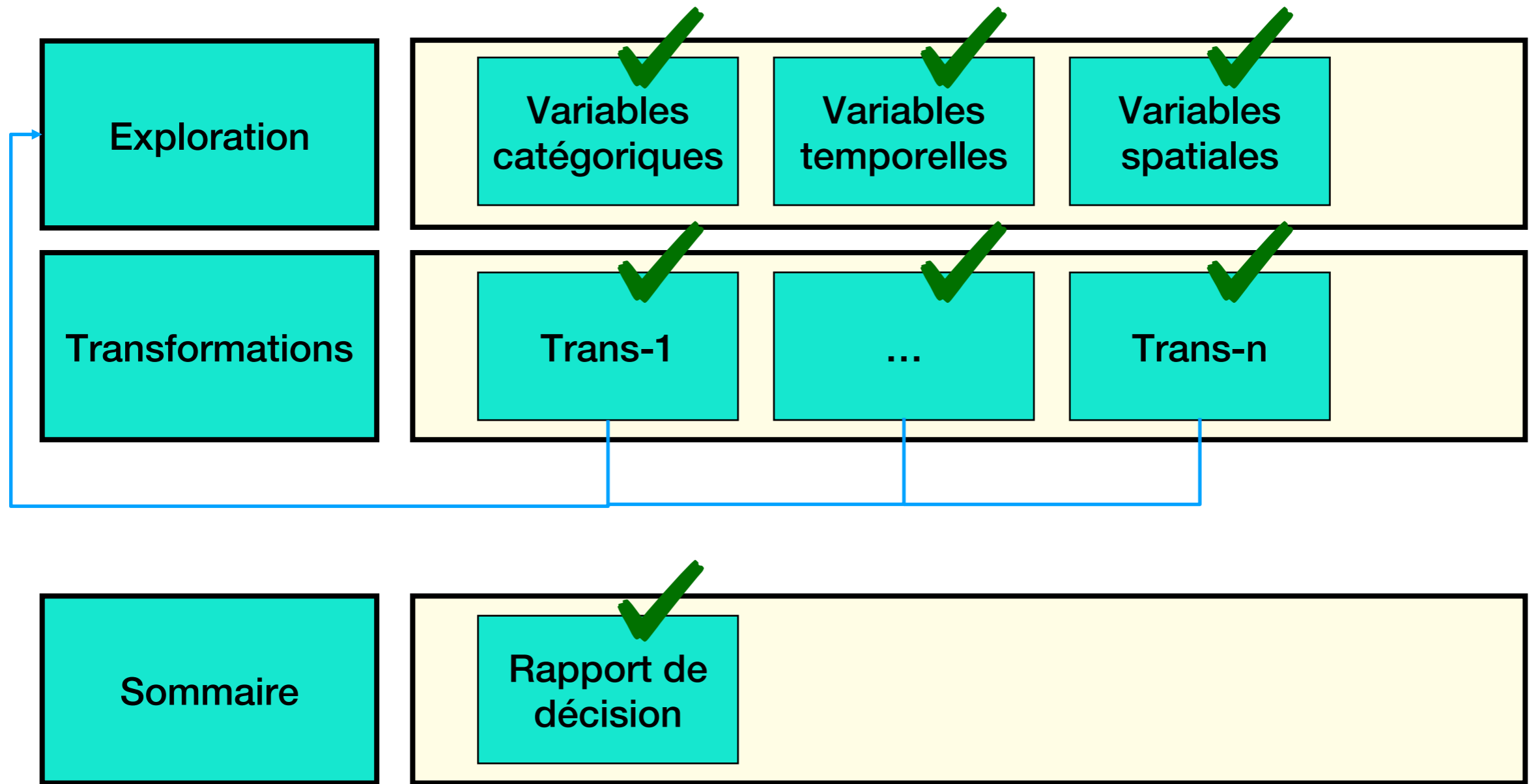


Transformations à explorer

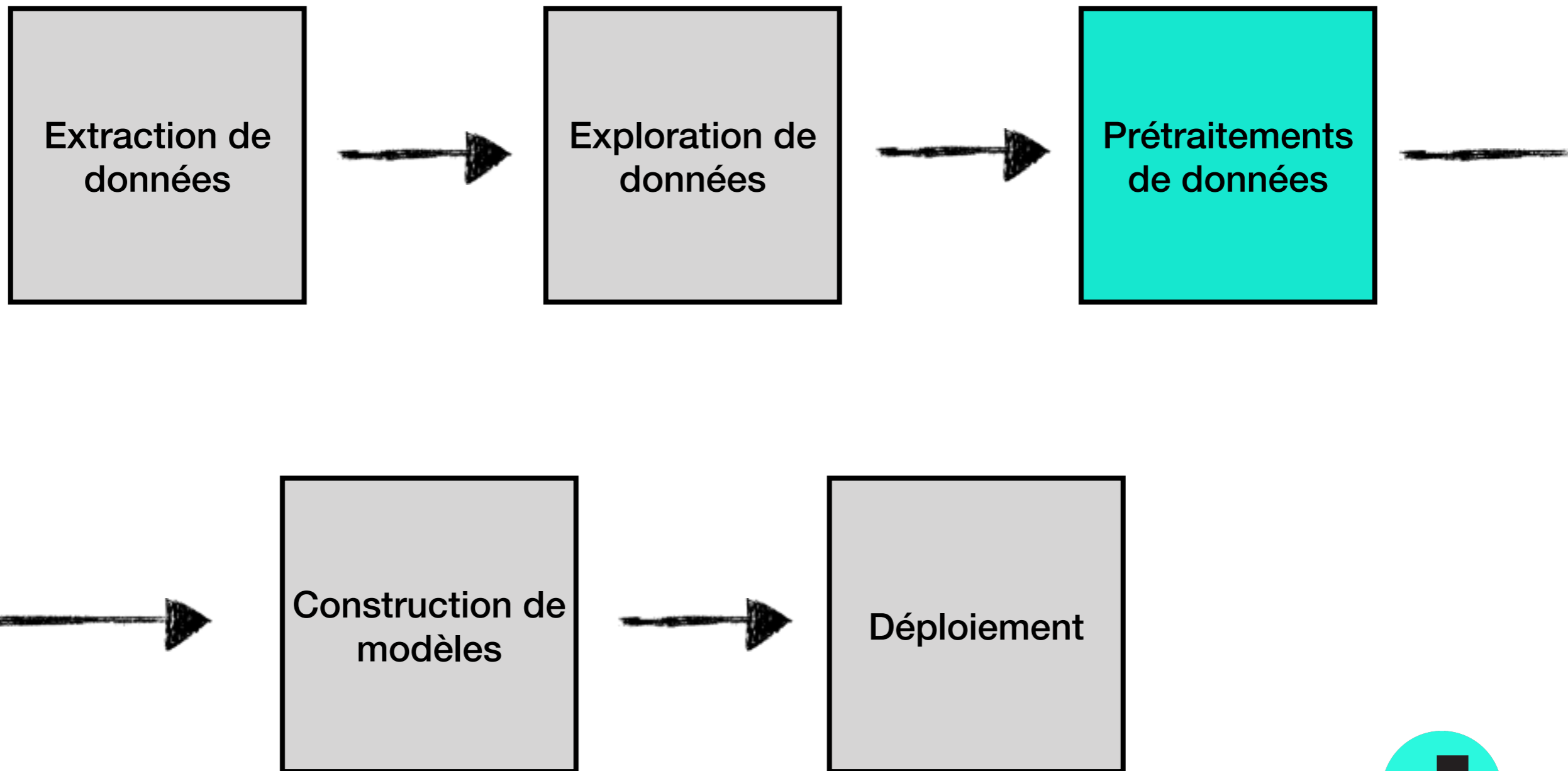
- Meilleur regroupement des moments de la journée
- Distance avec les autres quartiers
- Distance avec le fleuve (\approx distance avec la bordure de la carte des quartiers)



Récits à implémenter



Nettoyons !



Prétraitements de données



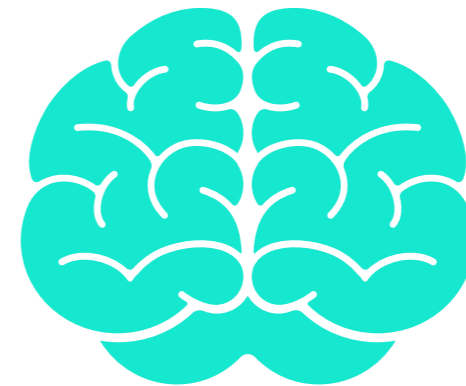
+



=



+



Objectifs

2 objectifs derrière le prétraitement de données:

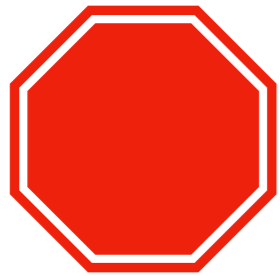
1. Transformer les données dans un format compatible pour l'algorithme
2. Transformer les données de manière à faciliter l'apprentissage



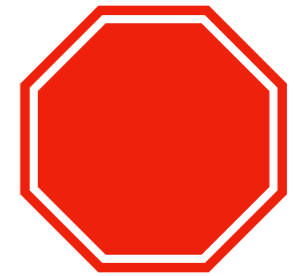
Prétraitements en 3 étapes

1. Nettoyage de données
2. Réduction de données
3. Transformations de données





Séparation



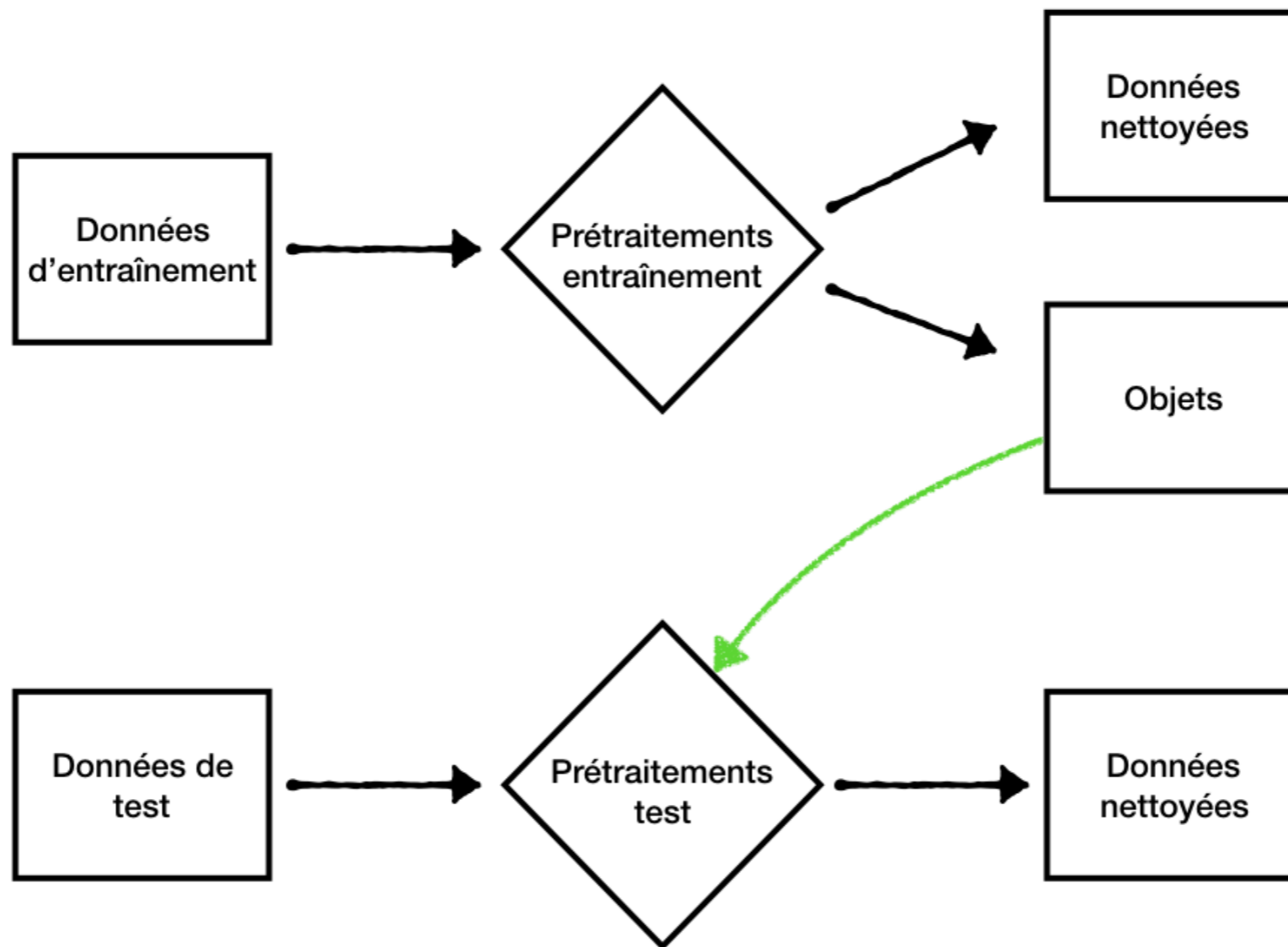
Tout d'abord, il faut séparer notre jeu de données ...

Méthodes:

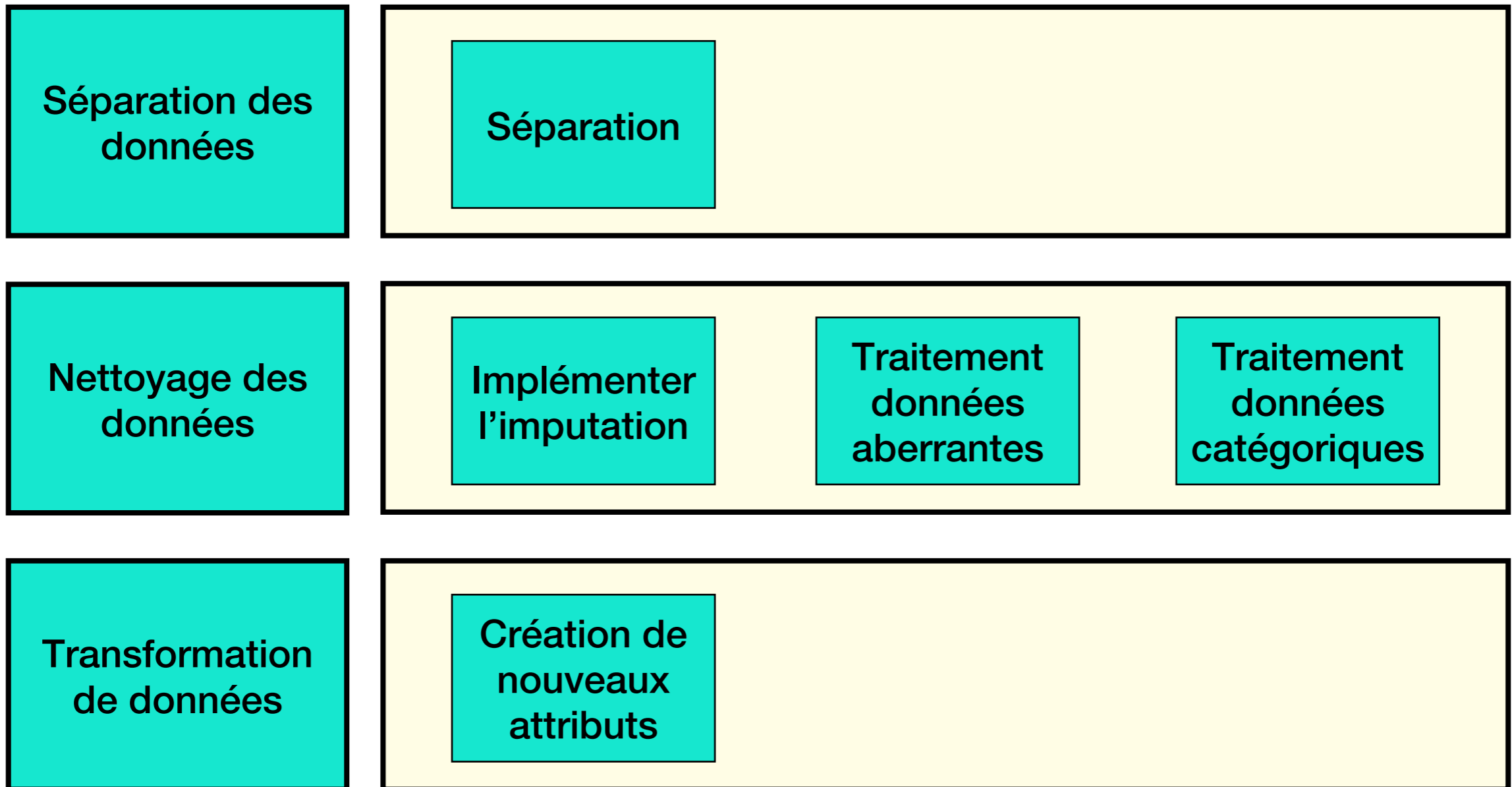
- Aléatoire
- Échantillonnage stratifié



Vu d'ensemble



Récits à implémenter



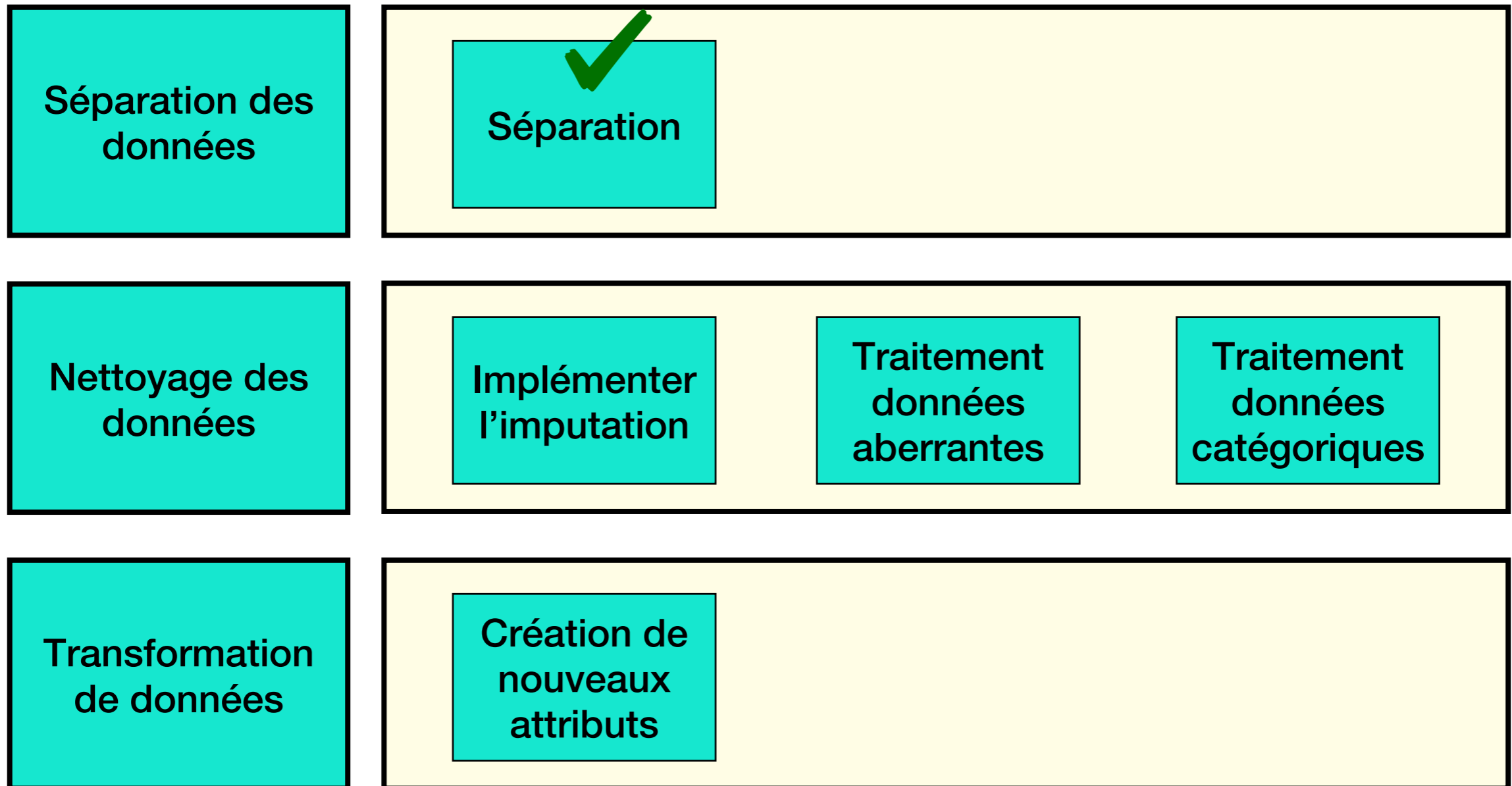
Nom du récit	Séparation du jeu de données
ID	1
Section	Prétraitements de données

Description
<p>Implémenter une procédure pour séparer les données en 2 ensembles :</p> <ul style="list-style-type: none"> ● Entraînement ● Test <p>Utiliser une séparation aléatoire (voir fonction `sample` de base R)</p>

Entrée
<ol style="list-style-type: none"> 1. Données brutes pour l'ensemble du jeu de données
Sortie
<ol style="list-style-type: none"> 1. Un vecteur d'index pour les observations d'entraînement 2. Sauvegarder les données test



Récits à implémenter



Nettoyage des données



Imputations de données manquantes

1. Faire le **constat** sur la quantité de données manquante.
2. Identifier le mécanisme de non-réponse
3. Traiter les données manquantes



Constat sur les données manquantes

- Quelles variables ont des données manquantes?
- Quelle est la proportion de données manquantes par variable?



Mécanisme de non-réponse

- Données manquantes complètement au hasard (MCAR)
- Données manquantes au hasard (MAR)
- Données manquantes pas au hasard (NMAR)

Voir la section 5.1.1 du livre pour des exemples.



Faire l'imputation

- Analyse des cas complets: Conserver uniquement les observations pour lesquelles toutes les variables sont présentes.
- Imputation par une mesure de centralité: Utiliser la moyenne, la médiane ou le mode pour remplacer les données manquantes.
- Imputation par régression: Remplacer les données manquantes par la prévision de modèle de régression entraîné sur les observations pour lesquelles cette variable est présente.
- Imputation par régression stochastique: Même chose que la méthode par régression, mais on ajoute un résidu aléatoire à la prévision.



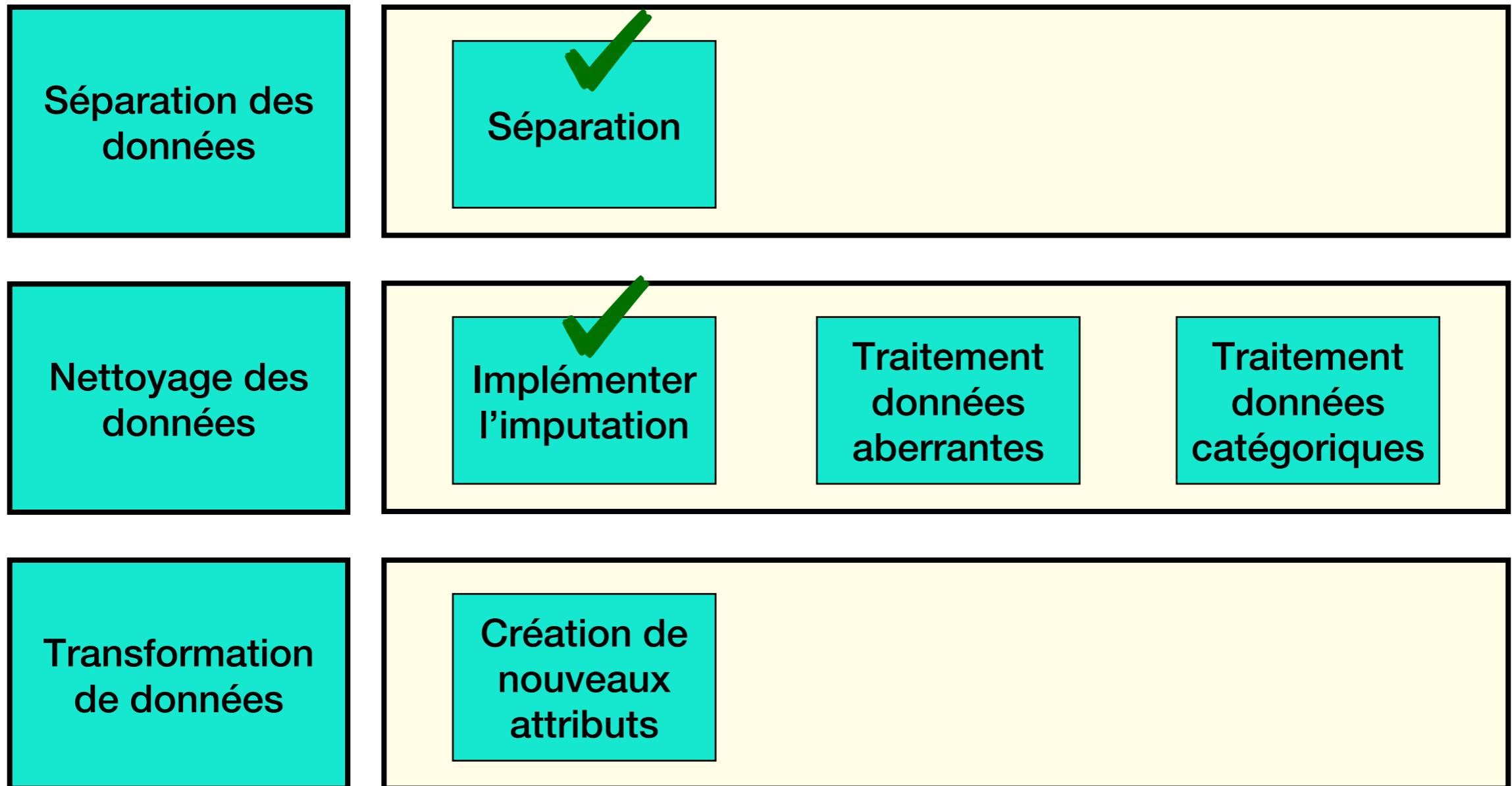
Nom du récit	Implémenter l'imputation
ID	2
Section	Prétraitements de données

Description
<p>Implémenter l'imputation de données manquantes. Il faut diviser l'imputation en 2 sections:</p> <ol style="list-style-type: none"> 1. Une qui aura comme but de définir les valeurs d'imputations. 2. Une qui aura comme but d'appliquer ces valeurs aux données manquants. <p>Le 1. sera fait en entraînement seulement. Le 2. sera fait autant en entraînement qu'en inférence.</p>

Entrée
<ol style="list-style-type: none"> 1. Données brutes 2. Données brutes + liste de valeurs d'imputations
Sortie
<ol style="list-style-type: none"> 1. Liste de valeurs d'imputaiton 2. Données imputées



Récits à implémenter



Traitement des données aberrantes

Pourquoi?: Effets importants dans le calcul des estimateurs

Comment?

- Plus ou moins 3 écarts-types de la moyenne
- Plus ou moins 1.5 EI (écarts interquartile)
- Partitionnement (*clustering*)



Données catégoriques

Identifier le type de données catégoriques:

- Attribut ordinal
- Attribut nominal

Le traitement n'est évidemment pas le même :

- Assigner une valeur numérique
- Encodage *un-chaud* (*one-hot encoding*)



Attribut ordinal

```
##      Id condition_station
## 1:   1          moyen
## 2:   2          mauvais
## 3:   3          excellent
## 4:   4             bon
```

```
##      Id condition_station
## 1:   1             1
## 2:   2             0
## 3:   3             3
## 4:   4             2
```

Super !



Attribut nominal

```
##      Id start_quartier
## 1:    1  Ville-Marie
## 2:    2      Verdun
## 3:    3  Westmount
## 4:    4    LaSalle
```

```
##      Id start_quartier
## 1:    1          1
## 2:    2          2
## 3:    3          3
## 4:    4          4
```

lssh ! ...



Attribut nominal

```
##      Id quartier_centre.ville quartier_plateau.mont.royal quartier_verdun
## 1:   1                        1                        0                0
## 2:   2                        0                        1                0
## 3:   3                        0                        0                1
## 4:   4                        0                        0                0
##      quartier_rosemont
## 1:           0
## 2:           0
## 3:           0
## 4:           1
```

Ahh c'est mieux ...



Nom du récit	Traitement des données aberrantes
ID	3
Section	Prétraitements de données

Description
<p>Implémenter une opération pour gérer les données abberantes.</p> <p>Choisir une méthode de détection de données aberrantes et appliquer un traitement pour ces données.</p> <p>Appliquer ce traitement dans le prétraitement des données d'entraînement seulement.</p>

Entrée
<ol style="list-style-type: none"> 1. Données brutes
Sortie
<ol style="list-style-type: none"> 1. Données brutes sans les données abberantes 2. "Print" du nombre de données traitées (fins de documentation)



Nom du récit	Encodage des données catégoriques
ID	4
Section	Prétraitements de données

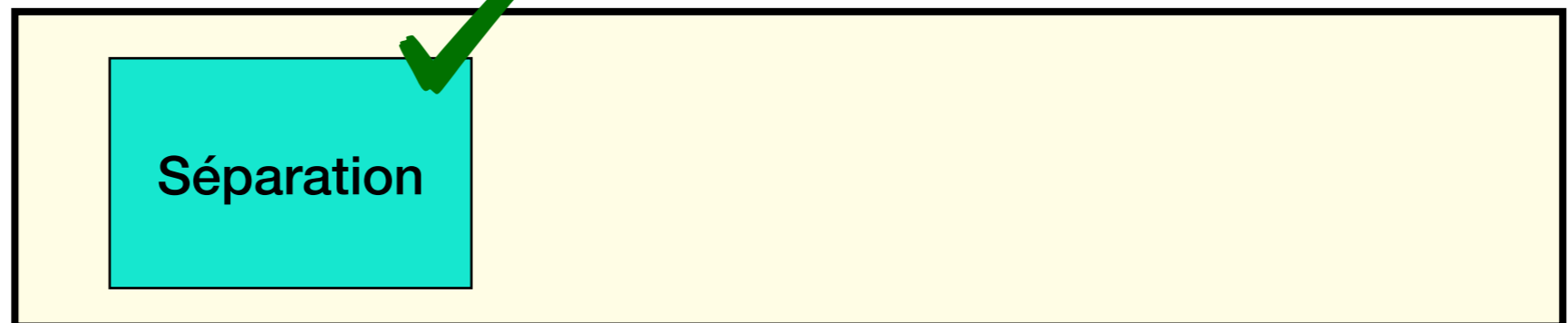
Description
<p>Implémenter une opération pour gérer les données catégoriques.</p> <ol style="list-style-type: none"> 1. Identifier les attributs catégoriques et déterminer leur type. 2. Faire l'assignation numérique pour les attributs ordinaux. 3. Lire la documentation de la fonction caret::dummyVars 4. Implémenter l'encodage <i>un-chaud</i> <p>Conseils: Il faut séparer l'encodage un-chaud pour l'entraînement et l'inférence. Prévoir un traitement pour les classes inconnus (défaut).</p>

Entrée
<ol style="list-style-type: none"> 1. Données brutes avec des catégories
Sortie
<ol style="list-style-type: none"> 1. Données sans catégories 2. Objets nécessaire pour faire l'encodage <i>un-chaud</i>

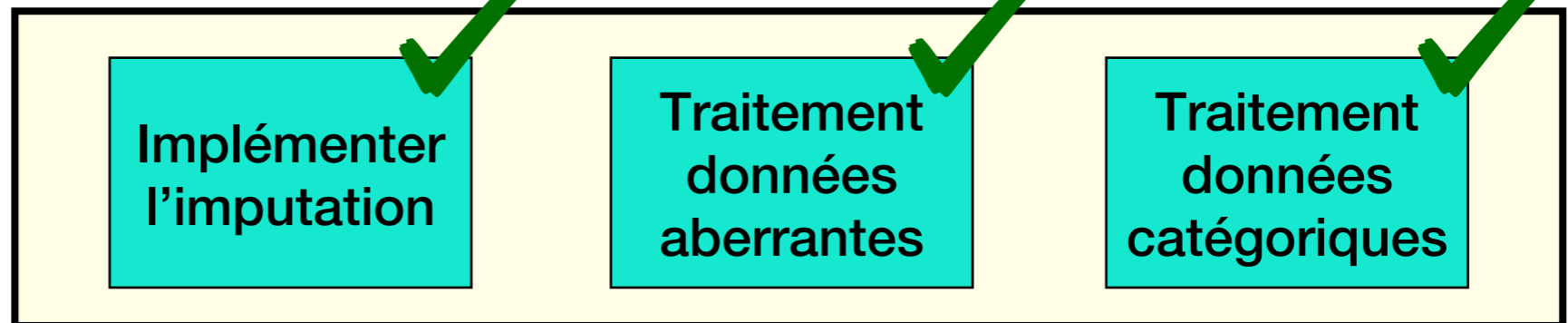


Récits à implémenter

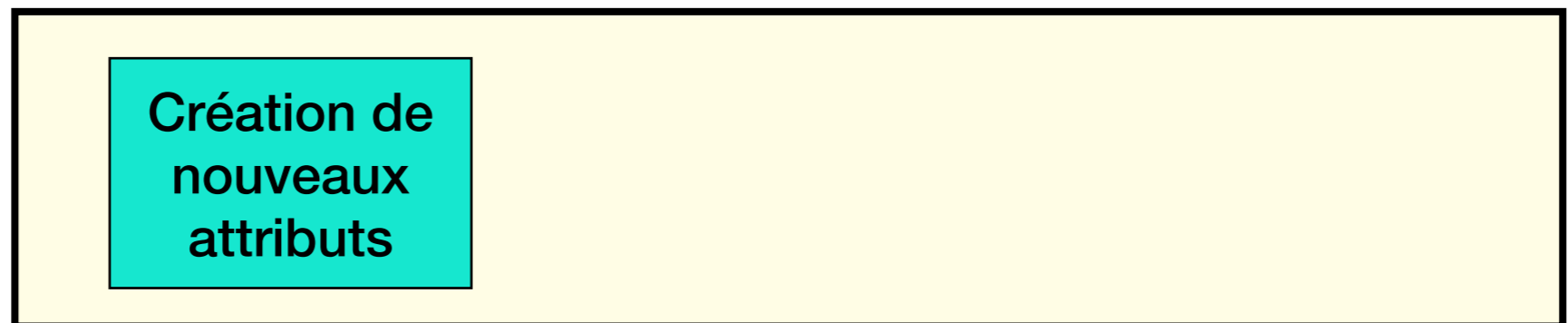
Séparation des données



Nettoyage des données



Transformation de données





Réduction de données



Fléau de la dimensionnalité

- Survient lorsque le nombre de dimensions est très élevée ($p \gg n$)
- Éloigne les données les unes des autres
- Rend l'apprentissage plus difficile



Fléau de la dimensionnalité

On suppose qu'on a un jeu de données avec $p = 1$ attribut de n observations où $x_1, \dots, x_n \stackrel{iid}{\sim} U(0, 1)$.
Combien d'observations en moyenne se trouveront dans l'intervalle $[0; 0.1]$?

La réponse : $\frac{n}{10}$ observations.

Maintenant, supposons que notre jeu de données est plus complexe et possède $p = 10$ attributs au lieu d'un seul attribut. Les observations suivent toujours une loi uniforme où $x_1, \dots, x_n \stackrel{iid}{\sim} U([0, 1]^{10})$.
Combien d'observations en moyenne se trouveront dans l'intervalle $[0; 0.1]^{10}$?

La réponse : $n\left(\frac{1}{10}\right)^{10}$ observations.

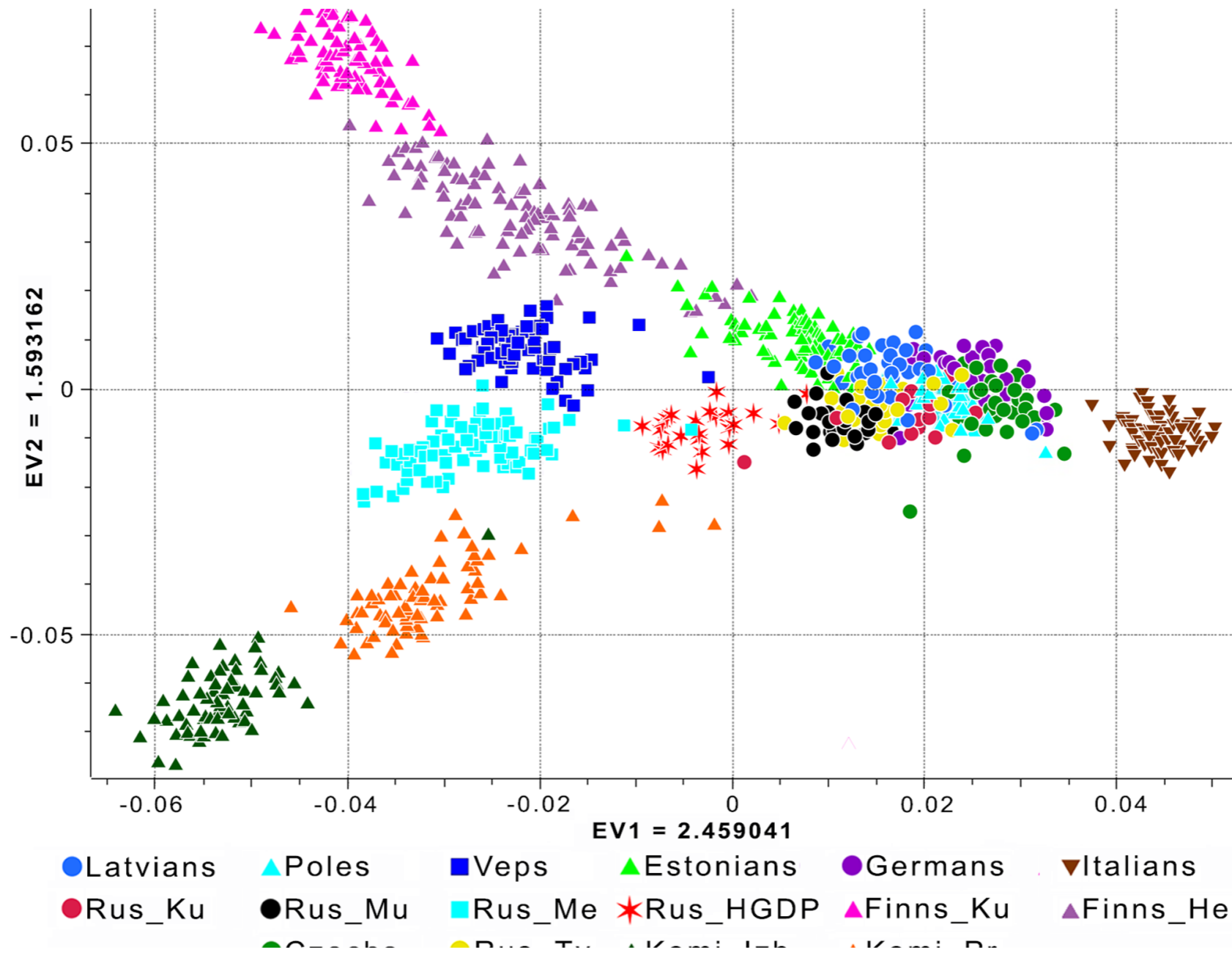


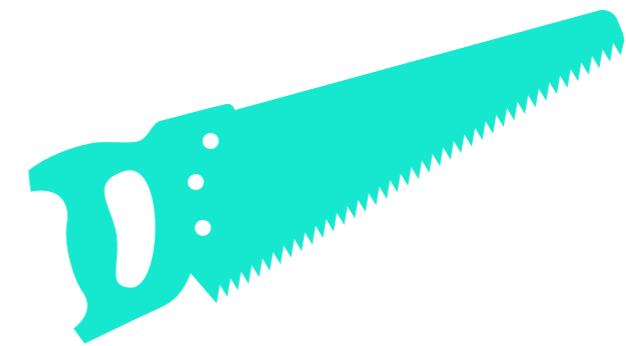
Méthodes de réduction de dimension

- Analyse en composantes principales
- Positionnement multidimensionnel
- Analyse factorielle
- Calcul de scores



ACP





Transformation de données



Normalisation

- Les échelles des attributs peuvent être différents
- Peut avoir un impact important dans certains algorithmes basés sur des distances (*k*-PPV, *clustering*)
- Permet de ramener les données autour d'une distribution plus "standard"



Méthodes de normalisation

- Normalisation centrée réduite
- Normalisation *min-max*
- Normalisation par décimation



Discrétisation et nouveaux attributs

- Communément appelé le *feature engineering*
- Grouper des observations à l'intérieur de *buckets* ou *bins*
- Transformer des attributs existants pour en créer des nouveaux



Nom du récit	Création de nouveaux attributs
ID	5
Section	Prétraitements de données

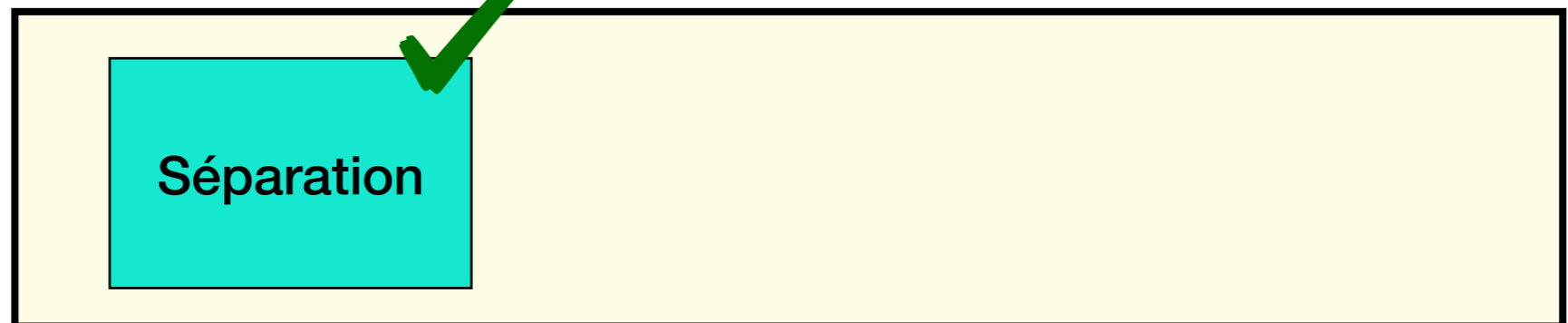
Description
<p>Implémenter le code qui permet de créer les nouveaux attributs pour le modèle. Utiliser les connaissances acquises lors de l'exploration de données.</p> <p>Créer une liste de variables à conserver dans le modèle.</p> <p>Conseils d'attributs:</p> <ul style="list-style-type: none"> - Grouper des quartiers - Moment de la journée - Fin de semaine versus semaine

Entrée
<ol style="list-style-type: none"> 1. Données brutes
Sortie
<ol style="list-style-type: none"> 1. Données avec nouveaux attributs 2. Liste de variables à conserver dans le modèle lors de l'inférence

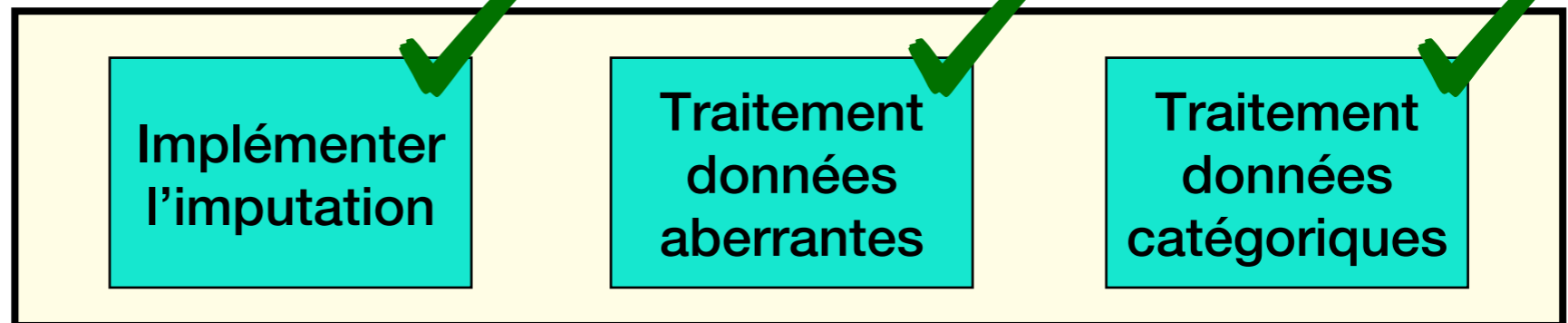


Récits à implémenter

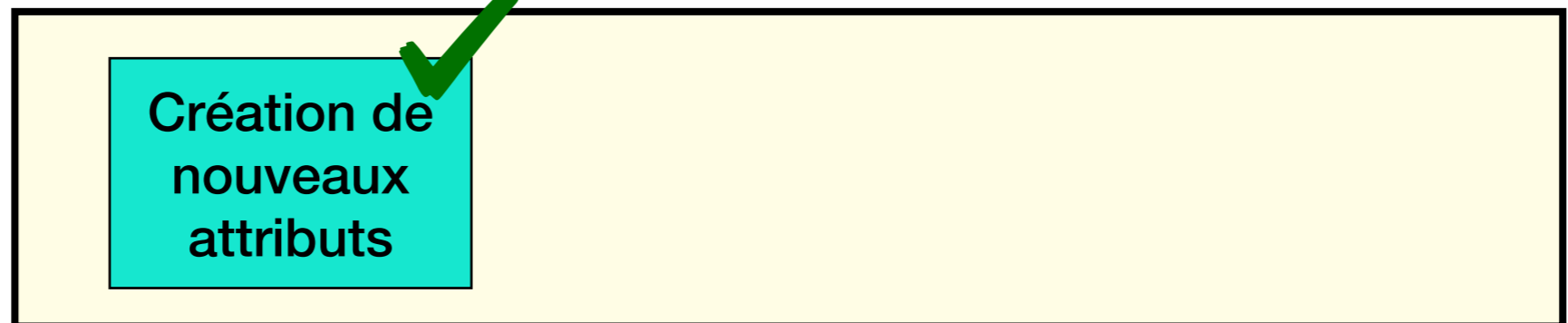
Séparation des données



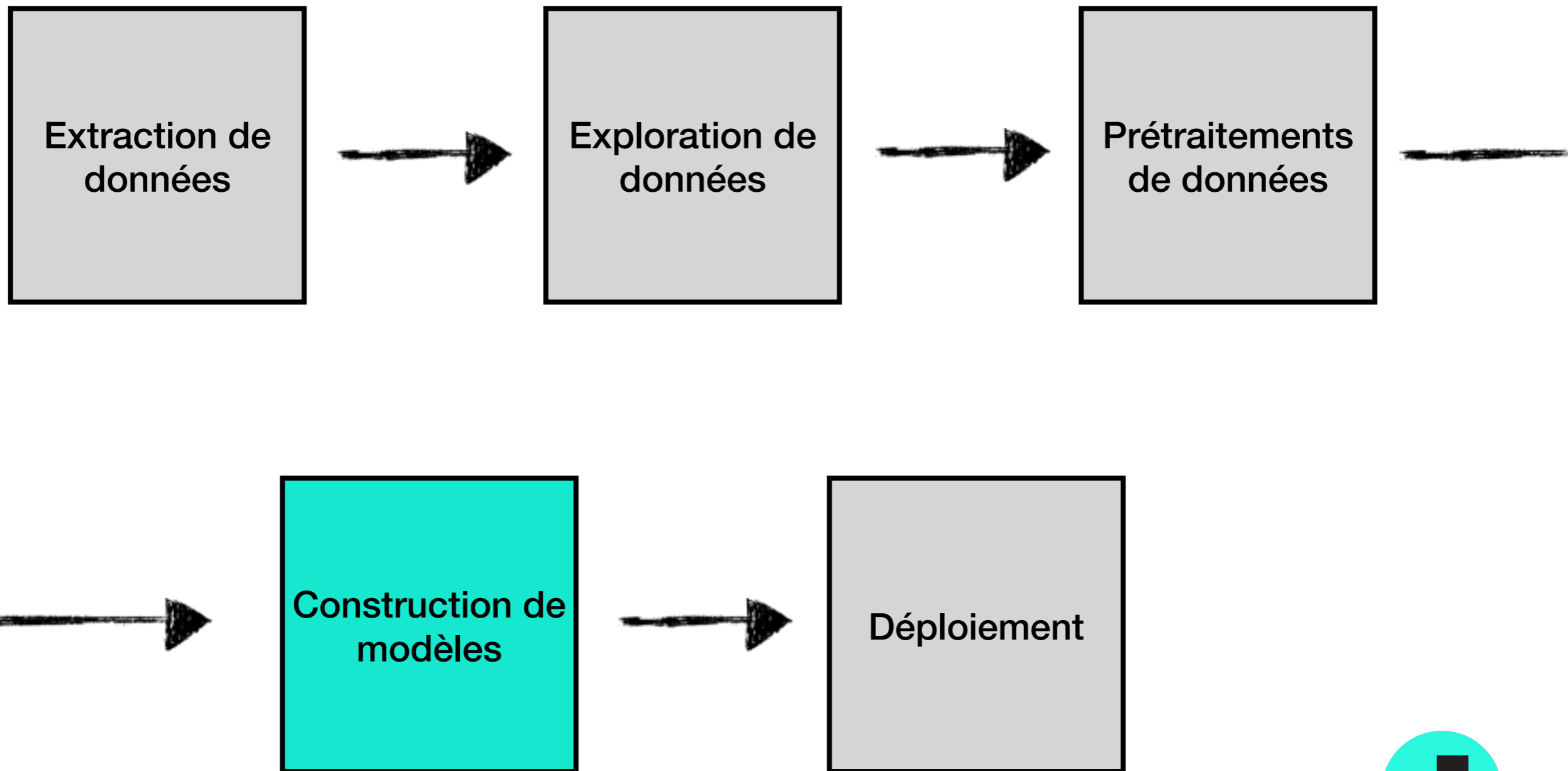
Nettoyage des données



Transformation de données



Modélisons !



Construction de modèles



Objectifs

1. Objectif de modélisation :

Prédire la réponse y pour de nouvelles données x

2. Objectif global :

Fournir un objet R à utiliser avec `predict`

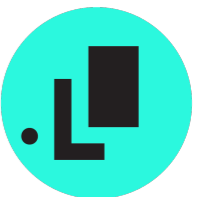


Modélisation en 3 étapes

1. Choix des modèles (on y a déjà pensé!)
2. Estimation des paramètres
3. Sélection du modèle final

Concepts clefs

entraînement/validation/test fonction de perte régularisation*
hyper-paramètres biais/variance validation croisée ...



Mais avant...

Un peu old school ça...

$$(\mathbf{X}|\mathbf{y}) = \left(\begin{array}{ccc|c} x_{11} & \dots & x_{1d} & y_1 \\ x_{21} & \dots & x_{2d} & y_2 \\ \vdots & & \vdots & \vdots \\ x_{n1} & \dots & x_{nd} & y_n \end{array} \right) \left. \begin{array}{l} \} \xrightarrow{\approx 1/2} (\mathbf{X}_{\text{train}}|\mathbf{y}_{\text{train}}) \\ \} \xrightarrow{\approx 1/4} (\mathbf{X}_{\text{val}}|\mathbf{y}_{\text{val}}) \\ \} \xrightarrow{\approx 1/4} (\mathbf{X}_{\text{test}}|\mathbf{y}_{\text{test}}) \end{array} \right.$$

Optionnel

Entrée	Sortie
Données test	Données validation + Données test

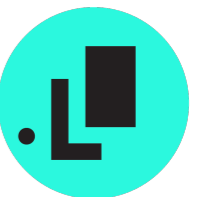


Description d'un modèle

Un peu de statistique...

ε

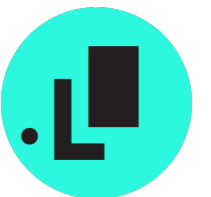
$$f(x) = \beta_0 + \beta_1 x_1 + \cdots + \beta_d x_d$$



Formulation d'un modèle

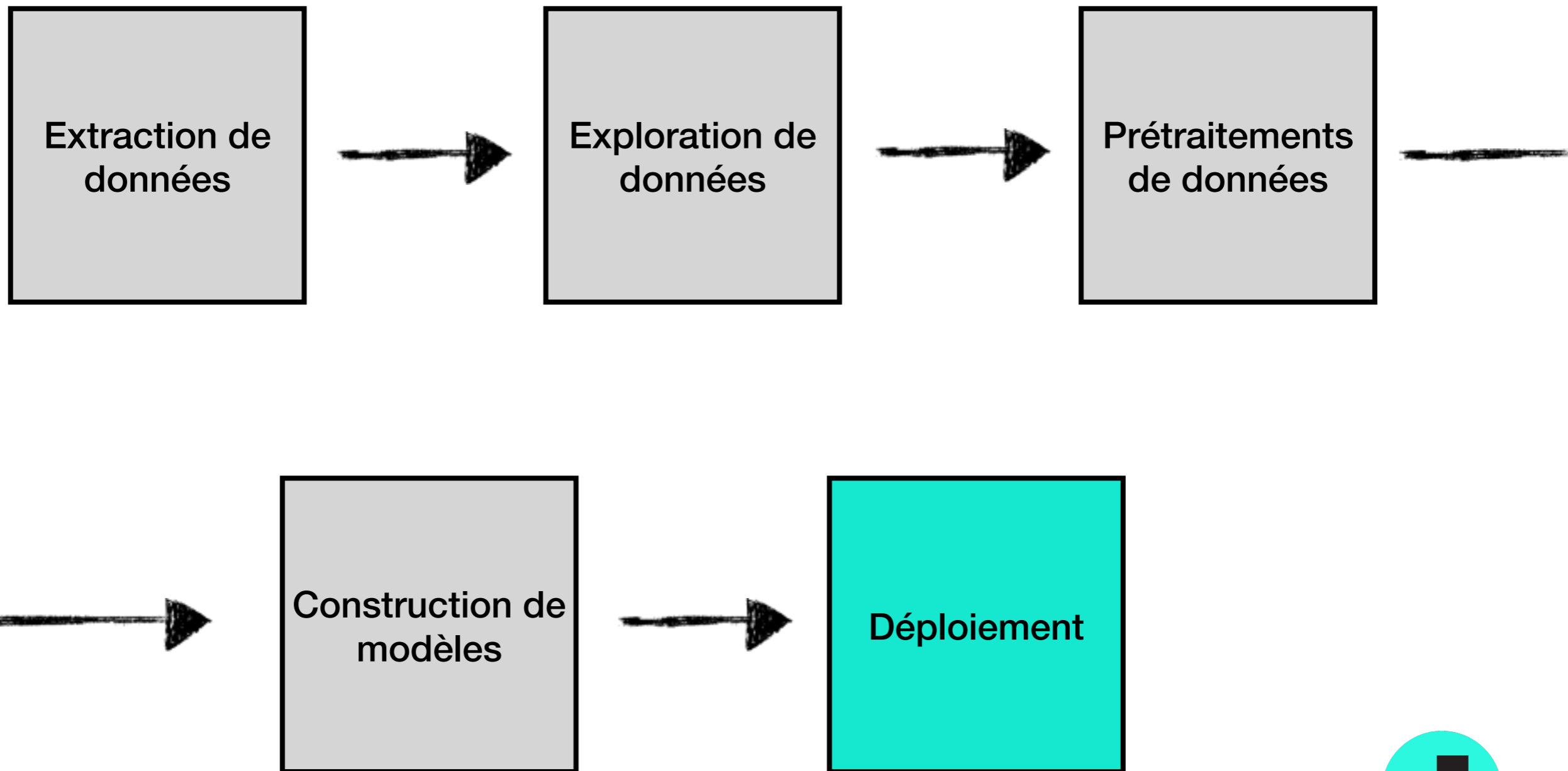
Classification...

ε



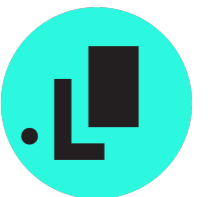
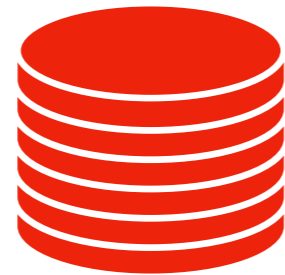
Choix des modèles

Déployons !

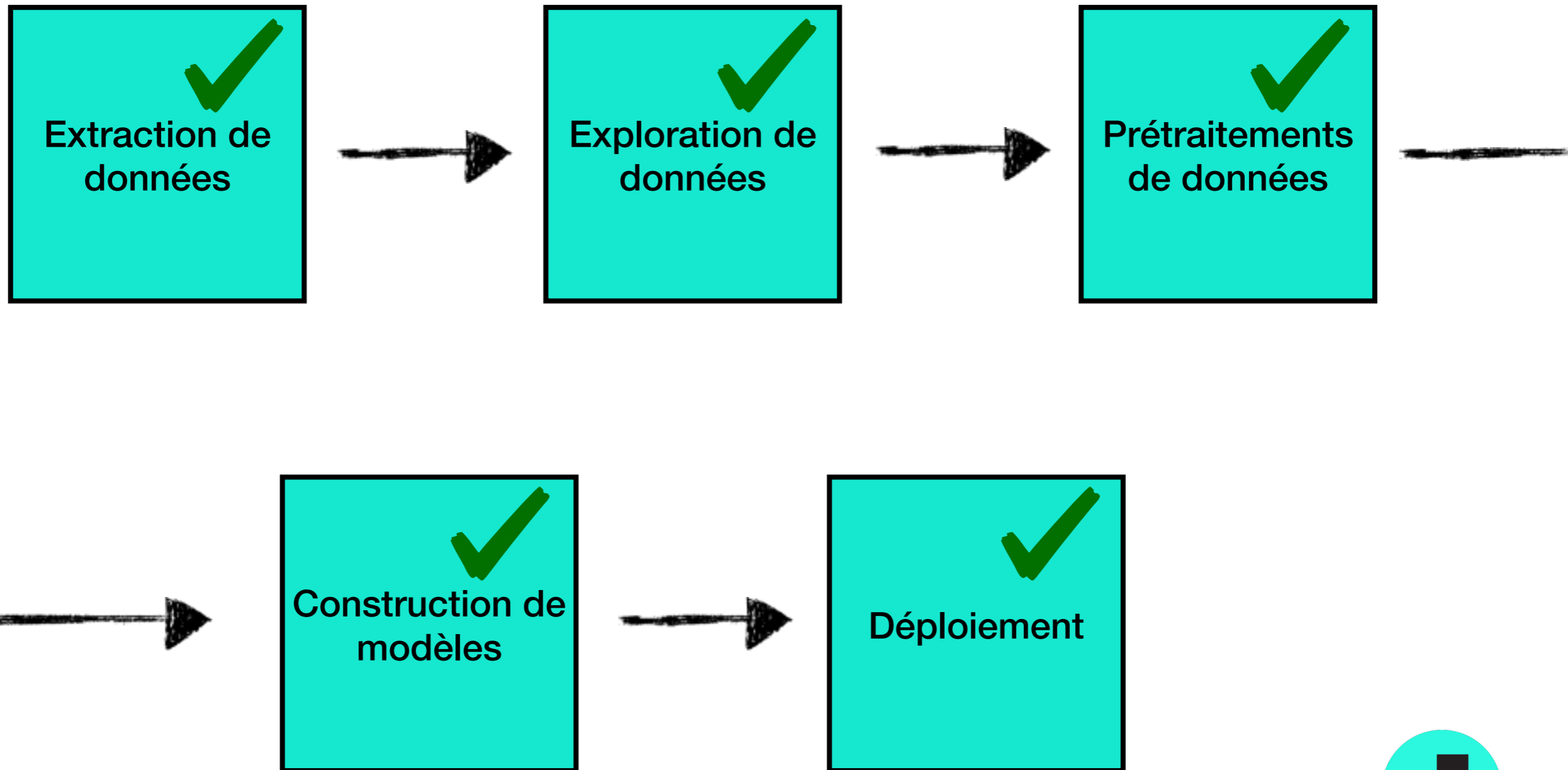


Déploiement

- Jé the man !



Mission accomplie !



Mot de la fin

