



Bordeaux, le 29 Novembre 2023

## CRITÈRES DE QUALITÉ DONNÉES

### Critères retenus pour les deux niveaux de la qualité des Métadonnées des jeux de données du Catalogue PIGMA

#### Membres de l'équipe PIGMA

Guillaume Blanchard – Administrateur de la plateforme PIGMA

Héloïse Deschamps – Animatrice technique de la plateforme PIGMA

## 1. Analyse des critères de qualité de la donnée analysables

Les jeux de données (ou JDD) du catalogue PIGMA se caractérisent par une fiche métadonnées, avec des champs obligatoires à renseigner et d'autres optionnels et une ou plusieurs ressources associées. La plupart des jeux de données ont une ressource principale ainsi que des ressources. Différents défauts peuvent affecter la qualité de ces données : topologies invalides, nommage des champs incorrects, absence d'identifiant, incohérence des attributs temporels, omission d'éléments, etc... Afin de déterminer la qualité des données du catalogue PIGMA, on analyse, pour chacun des 5 critères de la qualité des données géographiques, les sous-critères analysables et les analyses qui peuvent être menées.

### 1.1 La norme ISO 19157

L'Organisation Internationale de Normalisation, l'ISO, a publié en 2013 la norme 19157 sur la qualité des données géographiques. Cette norme établit des principes et définit des critères de description des données et de leur qualité, en spécifiant des méthodes de contrôle et en décrivant les procédures d'évaluation de cette qualité. Plus spécifiquement, elle pose 5 grands critères de qualité, que le Cerema a synthétisé sous forme de fiches en 2017 :

- **Cohérence logique** : Degré de cohérence interne des données, attestant que le jeu de données est exploitable
- **Exhaustivité** : Présence en excès ou absence d'éléments dans le jeu de données
- **Précision thématique** : Exactitude des informations portées par les entités
- **Précision de position** : Justesse de la position géographique, absolue ou relative
- **Qualité temporelle** : Qualité des attributs et des relations temporels entre les entités

Chacun de ces critères comporte plusieurs sous-critères, qu'on analyse un par un, en précisant les défauts qu'on peut analyser, la métrique utilisée pour mesurer ce défaut (nombre d'erreurs, taux d'erreurs ou conformité stricte) et les remarques qui peuvent être faites quant à la faisabilité de l'analyse.



## 1.2 Cohérence logique

Nom du critère	Défauts analysables	Représentation	Remarques
<b>Cohérence Conceptuelle</b>	<b>Non-conformité au schéma conceptuel :</b> <ul style="list-style-type: none"> <li>Identifiant non-unique</li> <li>Absence de champ identifiant</li> <li>Identifiant non-renseigné pour une entité</li> <li>Nom de champs incorrect (espace, accent, caractère spécial, etc.)</li> </ul>	Nombre d'erreurs, taux d'erreurs ou conformité stricte	<b>Critère important</b> pour l'exploitation. Identification automatique des champs identifiants difficile mais <b>faisable</b> . Respect du nommage des champs faisable via des expressions régulières.
<b>Cohérence au Domaine de Valeurs</b>	<b>Non-conformité au domaine de valeurs :</b> <ul style="list-style-type: none"> <li>Valeurs en dehors de la plage définie (pH d'un sol &gt; 15)</li> <li>Type erroné (population de 15 250,3 hab)</li> <li>Valeur d'attribut non prévue par les spécifications</li> </ul>	Nombre d'erreurs, taux d'erreurs	Hors-domaine de valeurs <b>difficile à établir</b> pour des valeurs numériques → détermination par les écarts-types uniquement si les valeurs sont très concentrées autour d'une valeur moyenne. Hors-domaine impossible pour les valeurs textuelles, sans disposer des spécifications. D'une manière général, nécessité de disposer des plages de valeurs admises pour contrôler ce critère.
<b>Cohérence de Format</b>	<b>Conflit de structure physique :</b> <ul style="list-style-type: none"> <li>Format de fichier incohérent avec le format renseigné</li> <li>Absence de fichiers obligatoires (ex : shx dans un shapefile)</li> <li>Format non-validé par un validateur (JSON invalide)</li> <li>Format de date erroné (JJ-AAAA-MM)</li> <li>Format de code incomplet (ex : INSEE_COM = 1300 au lieu de 01300)</li> </ul>	Nombre d'erreurs, taux d'erreurs ou conformité stricte	Cohérence entre format réel/format renseigné facilement analysable. Idem pour les fichiers du Shapefile. Usage d'un validateur possible pour certains formats (CSV, XML, JSON) Analyse du format des codes possibles uniquement pour des champs dont le type ou le nom permet d'être quasi-certain qu'il s'agisse de code (ex : type varchar(5), nom : code_commune, etc.)
<b>Cohérence Topologique</b>	<b>Défauts de topologie :</b> <ol style="list-style-type: none"> <li>Connexion arc-nœud erroné</li> <li>Erreur de connexion (undershoot, overshoot)</li> <li>Micro-polygones</li> <li>Pointe (spike)</li> </ol>	Nombre d'erreurs, taux d'erreurs ou conformité stricte	Tous les défauts de topologie cités sont analysables. Pour les couches ponctuelles, il reste à déterminer si l'on souhaite que deux points proches soient accrochés (coordonnées



	<ul style="list-style-type: none"> <li>5. Polygones troués</li> <li>6. Auto-intersection (polygones papillons)</li> <li>7. Auto-chevauchement (intersection entre entités mitoyennes)</li> </ul>		<p>identiques) ou parfaitement distincts (coordonnées différentes).          Les auto-chevauchements se détectent bien ; les lacunes (creux entre entités mitoyennes) sont plus difficiles à détecter.          Les défauts 1 et 2 concernent davantage les données de réseaux, les défauts 3 et 4 faussent les calculs de surface. Les défauts 5, 6 et 7 gênent voire invalident les analyses spatiales qu'on peut réaliser.</p>
--	--	--	---

La Cohérence au Domaine de Valeurs ne pourra vraisemblablement pas être contrôlée. Une notation et une pondération de ces différents sous-critères restent à établir.

Y:



### 1.3 Exhaustivité

Nom du critère	Défauts analysables	Représentation	Remarques
<b>Excédent</b>	<b>Éléments en excès dans le jeu de données :</b> <ul style="list-style-type: none"> <li>Éléments non-concernés par le jeu de données</li> <li>Éléments obsolètes</li> </ul>	Nombre d'erreurs, taux d'erreurs	Comparaison par rapport à une donnée référentielle, faisant foi, par rapport à une étendue géographique définie (si un point d'eau du SDIS 16 est situé dans le sud-Dordogne → excédent) <b>Impossibilité</b> de déterminer et <b>d'accéder automatiquement à une donnée de référence</b> pour un jeu de données donné. Éventuellement détermination des excédents par comparaison spatiale
<b>Omission</b>	<b>Éléments manquants dans le jeu de données :</b> <ul style="list-style-type: none"> <li>Éléments concernés par le jeu de données mais non-inclus</li> <li>Éléments supprimés par erreur</li> <li>Comparaison d'éléments entre deux versions d'un même jeu de données</li> </ul>	Nombre d'erreurs, taux d'erreurs	<b>Nécessité absolue de disposer d'une donnée référentielle.</b> La donnée OSM n'est pas toujours assez exhaustive et/ou précise pour servir de référence générique Éventuellement comparaison par rapport à une version antérieure de la donnée et rapport de différence remonté au producteur

L'Exhaustivité ne sera vraisemblablement pas analysable sauf sur deux éléments :

- La **comparaison par rapport à une zone géographique définie** (via les métadonnées ou l'étendue spatiale des données), avec un tampon de largeur à définir
- La comparaison par rapport à une version antérieure d'un jeu de données, permettant de détecter des entités potentiellement supprimées par erreur

Une notation et une pondération de ces différents sous-critères devraient également être établies.

Y:



## 1.4 Précision thématique

Nom du critère	Défauts analysables	Représentation	Remarques
<b>Justesse du Classement</b>	<b>Classement erroné :</b> <ul style="list-style-type: none"> <li>Valeurs d'attribut qui ne correspond pas à la valeur réelle (ex : Caserne référencée comme Hôpital)</li> </ul>	Nombre d'erreurs, taux d'erreurs ou conformité stricte	Comparaison par rapport à la réalité donc impossible à déterminer automatiquement. <b>Nécessité d'une donnée référentielle, d'un dire d'expert, etc.</b> Seul le producteur ou un utilisateur attentif peut attester de la Justesse du Classement.
<b>Justesse des Attributs Non-Quantitatifs</b>	<b>Attributs qualitatifs erronés :</b> <ul style="list-style-type: none"> <li>Valeurs d'attribut qui ne correspond pas à la valeur réelle (ex : Caserne référencée comme Hôpital)</li> <li>Orthographe de l'attribut incorrecte (ex : surface de type « <i>phorêt</i> »)</li> </ul>	Nombre d'erreurs, taux d'erreurs	Déterminer la justesse nécessite une donnée référentielle ou un dire d'expert. Seule la <b>vérification de l'orthographe</b> peut être faite automatiquement, via des bibliothèques Python (ex : pypellchecker)
<b>Justesse des Attributs Quantitatifs</b>	<b>Conflit de structure physique :</b> <ul style="list-style-type: none"> <li>Valeur d'attribut qui ne correspond pas à la valeur réelle (ex : pommier de hauteur 150m)</li> </ul>	Nombre d'erreurs, taux d'erreurs ou conformité stricte	Comparaison par rapport à la réalité donc impossible à déterminer automatiquement. <b>Nécessité d'une donnée référentielle, d'un dire d'expert, etc.</b> Seul le producteur ou un utilisateur attentif peut attester de la Justesse des Attributs Quantitatifs

Seule la vérification de l'orthographe dans certains champs peut être effectuée. Une notation et une pondération de ces différents sous-critères devraient également être établies.

Y:



## 1.5 Précision de position

Nom du critère	Défauts analysables	Représentation	Remarques
<b>Précision Absolue</b>	<b>Proximité des coordonnées :</b> <ul style="list-style-type: none"> <li>Écart des coordonnées absolues par rapport aux coordonnées réelles</li> </ul>	Nombre d'erreurs, taux d'erreurs, valeur moyenne des incertitudes de position	Comparaison par rapport à la réalité donc impossible à déterminer automatiquement. Nécessité d'une donnée référentielle, d'un dire d'expert, etc. Seul le producteur ou un utilisateur attentif peut attester de la Précision Absolue.
<b>Précision Relative</b>	<b>Proximité des coordonnées :</b> <ul style="list-style-type: none"> <li>Écart des coordonnées relatives par rapport aux coordonnées réelles</li> </ul>	Nombre d'erreurs, taux d'erreurs, valeur moyenne de l'erreur verticale/horizontale	Comparaison par rapport à la réalité donc impossible à déterminer automatiquement. Nécessité d'une donnée référentielle, d'un dire d'expert, etc. <b>Théoriquement faisable</b> pour des données routières ou de réseaux dont la distance entre les différentes bornes sont connues. Seul le producteur ou un utilisateur attentif peut attester de la Précision Relative.
<b>Présence dans une emprise</b>	<b>Comparaison par rapport à une emprise de référence :</b> <ul style="list-style-type: none"> <li>Éléments en dehors de l'emprise (avec ou sans tampon)</li> </ul>	Nombre d'erreurs, taux d'erreurs ou conformité stricte	<b>Comparaison spatiale faisable</b> → analyse similaire à la recherche d'excédents <b>mais difficulté à caractériser l'erreur</b> : si une entité est hors-zone → excédent ou écart de position ?

La présence dans une emprise peut être analysée, ce qui pose alors le problème de caractérisation de l'erreur (excédent ou écart de position). Une notation et une pondération de ce sous-critère devraient également être établies.

Y:



## 1.6 Qualité Temporelle

Nom du critère	Défauts analysables	Représentation	Remarques
<b>Exactitude de la précision temporelle</b>	<b>Proximité des attributs temporels :</b> <ul style="list-style-type: none"> <li>Écart des attributs temporels par rapport aux valeurs réelles</li> </ul>	Nombre d'erreurs, taux d'erreurs, valeur moyenne des écarts	Comparaison par rapport à la réalité donc impossible à déterminer automatiquement. Nécessité d'une donnée référentielle, d'un dire d'expert, etc. Seul le producteur ou un utilisateur attentif peut attester de l'Exactitude de la précision temporelle.
<b>Cohérence temporelle</b>	<b>Respect de la chronologie des dates :</b> <ul style="list-style-type: none"> <li>Chronologie des dates/horaires incohérentes (ex : jour de fermeture antérieur au jour d'ouverture)</li> </ul>	Nombre d'erreurs, taux d'erreurs, valeur moyenne de l'erreur verticale/horizontale	Comparaison entre deux attributs date, time ou varchar(4). Théoriquement possible de contrôler la chronologie. Avoir à disposition deux attributs date est peu fréquent dans les données du catalogue. En pratique, déterminer le champ date antérieur et le champ postérieur nécessite d'analyser sémantiquement le nom des champs (et y trouver des valeurs comme 'ouverture/fermeture', 'début/fin', 'entrée/sortie', etc.
<b>Validité temporelle</b>	<b>Dates invalides :</b> <ul style="list-style-type: none"> <li>Dates au mauvais format (JJ-AA-MM)</li> <li>Dates incongrues (13-13-0221, 30-02-2007, etc.)</li> </ul>	Nombre d'erreurs, taux d'erreurs ou conformité stricte	<b>Défauts de format temporels analysables</b> (DateTimeConverter, expressions régulières)

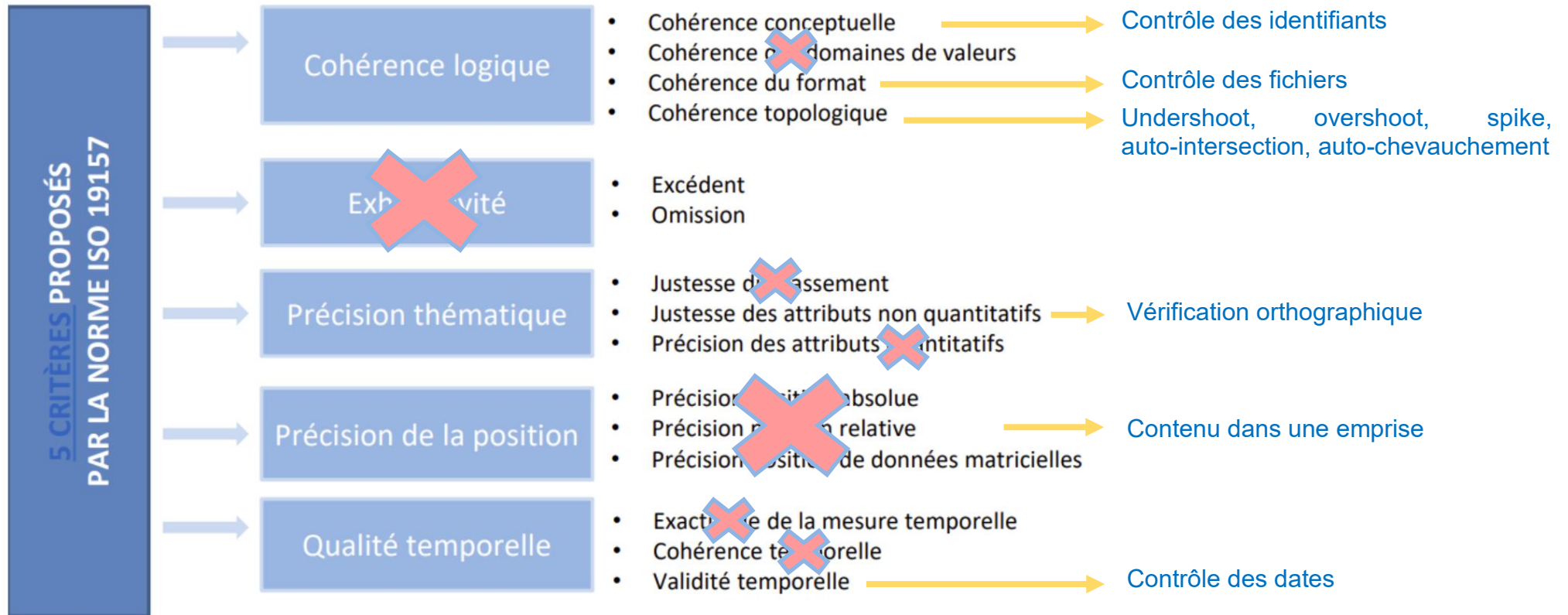
Le respect de la chronologie est théoriquement analysable mais reste difficile en pratique. Le critère de Validité temporelle peut être analysé, en contrôlant le format des attributs temporels. Une notation et une pondération de ce sous-critère devraient également être établies.

Y:



## 2. Synthèse des critères retenus

### Analyses possibles retenues



Y:





### 3. Notations et Pondération des critères retenus

#### 3.1 Cohérence logique

Nom du critère	Coeff	Type de critère	Utilité et Qualité attendue	Notation
<b>Identifiant unique</b>		Cohérence conceptuelle	La donnée doit posséder au moins un champ identifiant. Les identifiants doivent être uniques et renseignés pour toutes les entités, sans exception ni doublon.	<b>Critère exclusif :</b> La présence d'un doublon ou l'absence d'identifiant pour une entité de la donnée exclut le jeu de données du Label Optimal
<b>Noms des champs</b>		Cohérence conceptuelle	Le nommage des champs exclut les caractères spéciaux, les espaces et les signes de ponctuation.	<b>Critère exclusif :</b> La présence d'un ou plusieurs champs mal nommés exclut le jeu de données du Label Optimal
<b>Cohérence et Conformité du Format</b>	<b>5 et 5</b>	Cohérence du format	Le format déclaré de la donnée correspond au format réel (Cohérence). La donnée est conforme aux spécifications du format (JSON valide, fichiers indispensables du Shapefile).	<b>Note sur la cohérence :</b> <b>0 :</b> Les formats réel et déclaré ne concordent pas <b>10 :</b> Les formats réel et déclaré concordent + <b>Note sur la Conformité :</b> <b>ERREUR :</b> La donnée n'est pas conforme <b>5 :</b> La donnée est conforme mais ne contient que les fichiers indispensables <b>10 :</b> La donnée est conforme et contient plus de fichiers que ceux indispensables

Y:



<p><b>Codes valides</b></p>	<p><b>3</b></p>	<p>Cohérence du format</p>	<p>Toutes les valeurs des champs contenant des codes (type code culture, code insee, etc.) correspondent au format du code, avec un seuil d'erreur inférieur à 5%.</p>	<p><b>Note sur la cohérence :</b>  <b>0</b> : Au moins un champ code comporte des erreurs, avec un taux d'erreurs supérieur à 5%  <b>5</b> : Au moins un champ code comporte des erreurs, avec un taux d'erreurs inférieur à 5%  <b>OU</b> La donnée ne comporte aucun champ code  <b>10</b> : Tous les champs codes sont valides</p>
-----------------------------	-----------------	----------------------------	--	---

Y:



### 3.2 Cohérence logique - Topologie

Nom du critère	Coeff	Type de critère	Utilité et Qualité attendue	Notation
<b>Géométries nulles</b>		Cohérence Topologique	La donnée ne comporte aucune géométrie nulle. Chaque entité possède une géométrie.	<b>Critère exclusif :</b> La présence d'une géométrie nulle pour une entité de la donnée exclut le jeu de données du Label Optimal
<b>Micro-géométrie</b>	5	Cohérence Topologique	La donnée ne contient pas de micro-géométries (micro-polygones, scories de découpage).	<b>Note sur la cohérence :</b> <b>0 :</b> La donnée comporte plus de 5% de micro-géométries. <b>5 :</b> La donnée comporte moins de 5% de micro-géométries. <b>10 :</b> Aucune micro-géométrie
<b>Topologie</b>	5	Cohérence Topologique	La donnée ne possède pas de défaut de topologie. Les géométries sont correctement accrochées, sans écart, ni chevauchement. Lorsque deux entités linéaires se chevauchent, ces entités doivent partager un sommet commun. La donnée ne contient pas de pics, des erreurs de connexion des segments d'une géométrie.	<b>Note sur l'accrochage :</b> <b>0 :</b> La donnée contient au moins une erreur <b>10 :</b> La donnée ne contient aucune erreur

Y:



### 3.3 Précision de la position

Nom du critère	Coeff	Type de critère	Utilité et Qualité attendue	Notation
<b>Emprise de référence</b>	<b>3</b>	Précision absolue	Toutes les entités de la donnée sont contenues dans une emprise spatiale de référence (France ou Région Nouvelle-Aquitaine)	<b>Note sur la position :</b> <b>0</b> : La donnée comporte au moins une entité hors de l'emprise. <b>5</b> : La donnée est comprise en France. <b>10</b> : La donnée est comprise en Nouvelle-Aquitaine

### 3.4 Critères non déployés

Nom du critère	Coeff	Type de critère	Utilité et Qualité attendue	Notation
<b>Orthographe</b>	<b>1</b>	Précision thématique	Les valeurs dans les champs de typage (type, libellé, etc.) sont correctement orthographiées (ex : type « forêt », au lieu de « phorêt »)	<b>Note sur la position :</b> <b>0</b> : Au moins une erreur dans les champs de typage <b>5</b> : La donnée ne comporte pas de champs de typage <b>10</b> : Aucune erreur dans les champs de typage
<b>Dates</b>	<b>3</b>	Qualité temporelle	Les valeurs dans les champs « date » sont valides (pas de dates impossibles comme 30 février ou absurdes comme 01/01/0222) et correspondent aux formats suivants : JJMMAAAA, AAAAMMJJ (avec ou sans tiret ou barre oblique)	<b>Note sur la position :</b> <b>0</b> : Au moins une date est impossible ou absurde <b>5</b> : Les dates sont valides mais le format des dates n'est pas valide <b>10</b> : Le format des dates est valide et aucune erreur n'est présente

Y: