
Créer des schémas de données

Etalab

etalab^{gouv.fr}

07/03/2022

Table des matières

1	Introduction	4
1.1	À qui s'adresse ce guide?	4
1.2	À quoi sert-il?	5
1.3	Sources	5
2	Phase d'investigation	5
2.1	Étapes à suivre	6
2.2	Exemples	6
2.2.1	Situations favorables à la création d'un schéma de données	6
2.2.2	Situations où le référencement d'un schéma sur schema.data.gouv.fr ne semble pas nécessaire	6
2.3	Points de sortie	7
3	Phase de concertation	7
3.1	Procédure de collaboration	8
3.2	Grands principes	8
3.3	Points de sortie	9
4	Phase de construction	9
4.1	Choisir un standard technique pour la description de votre schéma de données	10
4.2	Créer votre schéma de données	10
4.3	Documenter votre schéma de données	11
4.4	Publier et diffuser votre schéma de données	11
4.5	Référencer votre schéma de données sur schema.data.gouv.fr	12
4.6	Faire évoluer votre schéma de données	12
4.7	Points de sortie	12
5	Intégration avec schema.data.gouv.fr	12
5.1	Qui peut référencer des schémas?	13
5.2	Quels schémas de données sont acceptés?	14
5.3	Quand référencer son schéma de données?	14
5.4	Quels schémas de données sont acceptés?	14
5.4.1	Standards techniques supportés	14
5.4.2	Prérequis de validation des schémas de données	15
5.5	Points de sortie	15

6	Jeux de données sur data.gouv.fr et schémas	16
6.1	Indiquer qu'une ressource respecte un schéma sur data.gouv.fr	16
7	Aide à la construction d'un schéma TableSchema	19
7.1	Recommandations pour le formatage des fichiers csv	19
7.1.1	Format de fichier csv	19
7.2	Recommandations de formatage des données	20
7.2.1	Données de type string	22
7.2.2	Données de type décimal	22
7.2.3	Données de type date	22
7.2.4	Données de type date avec heure	23
7.2.5	Données de type date avec heure de début et de fin	23
7.2.6	Données de type horaires d'ouverture	23
7.2.7	Données de type géolocalisation	24
7.2.8	Données de type adresse	25
7.3	Recommandations de champs obligatoires	25
7.3.1	Identification du producteur	26
7.4	Recommandations pour le nommage des fichiers	29
7.5	Recommandations pour le nommage des champs	29
7.6	Recommandations pour la mise en conformité	29

1 Introduction

Rédaction de ce guide

Ce guide a été publié initialement fin mars 2020 et est mis à jour de manière régulière. Il résulte d'une co-rédaction entre les équipes d'[Etalab](#) et d'[OpenDataFrance](#). Si vous souhaitez faire des propositions pour le faire évoluer, vous êtes invités à [entrer en contact avec nous](#).

Lexique : Schémas de données

Les schémas de données permettent de décrire la structure d'un jeu de données. Ils indiquent clairement quels sont les différents champs, comment sont représentées les données, quelles sont les valeurs possibles etc.

Synonymes : *modèle de données, modèle logique de données, schéma*.

La création d'un jeu de données en conformité avec un schéma de données existant apporte plusieurs bénéfices :

- Le jeu de données créé peut être facilement croisé avec d'autres jeux de données conformes au schéma de données utilisé ;
- L'interopérabilité des données et leur croisement est simplifié ;
- Si le jeu de données que vous créez est une agrégation de plusieurs fichiers produits par différents acteurs, la formalisation et le partage d'un schéma de données facilite le travail d'agrégation des données - ce schéma devient donc un standard pour votre communauté ;
- La formalisation d'un schéma de données assure une pérennité des fichiers dans le temps ;
- La documentation d'un schéma de données existant est déjà rédigée et accessible.

Il est également possible de vérifier la conformité d'un jeu de données vis-à-vis d'un schéma de données, ce qui permet de valider un premier niveau de qualité de votre jeu de données. Par ailleurs, il est aussi possible de générer des jeux de données d'exemple ou de proposer des formulaires de saisie standardisés.

[schema.data.gouv.fr](#)

Le site [schema.data.gouv.fr](#) est l'initiative de la plateforme [data.gouv.fr](#). L'objectif de ce site est de référencer les schémas de données publiques existants en France.

1.1 À qui s'adresse ce guide ?

Ce guide s'adresse à **des personnes susceptibles de créer des schémas de données**. Vous pouvez vous trouver dans cette situation si vous envisagez de partager des données avec des partenaires ou

à tout le monde en open data.

1.2 À quoi sert-il ?

Ce guide propose de vous accompagner lors des phases nécessaires à la création d'un schéma de données et à son référencement sur schema.data.gouv.fr le cas échéant.

1. **Phase d'investigation** : envisager de créer un schéma de données ;
2. **Phase de concertation** : rassembler plusieurs parties prenantes pour créer un schéma de données ;
3. **Phase de construction** : implémenter le schéma de données obtenu après la phase de concertation.

Il propose un processus à suivre, des bonnes pratiques et des outils.

Conseil de lecture

Nous vous recommandons de lire une première fois ce guide **en intégralité** afin de prendre connaissance des différentes phases. Vous pourrez ensuite vous référer aux pages pertinentes au fur et à mesure de votre avancée.

1.3 Sources

Ce guide s'inspire du contenu rédigé par de nombreux partenaires, listés par ordre alphabétique :

- [Charles Nepote](#)
- [Dataactivist](#)
- [La FING](#)
- [OpenDataFrance](#)

Merci à eux !

2 Phase d'investigation

Lexique : Phase d'investigation

La phase d'investigation est la première phase de la création d'un schéma de données. Cette phase a pour finalité de s'assurer que la création d'un schéma est pertinente et vise à aboutir à la décision de continuer ou de choisir une autre alternative.

2.1 Étapes à suivre

Afin de déterminer s'il est nécessaire de créer ou non un schéma de données, nous vous recommandons de suivre les étapes suivantes :

1. Lire attentivement les différentes sections de ce guide ;
2. Organiser une réunion réunissant des acteurs métiers, techniques et de potentiels réutilisateurs. Lors de cette réunion, vous débattrez de la pertinence de la création de votre schéma de données ;
3. Entrer en contact avec [les équipes d'Etalab](#) et leurs partenaires en [référénçant votre schéma](#), afin de bénéficier de conseils lors de la création de votre schéma de données, d'une visibilité accrue pour celui-ci et d'une assistance d'experts.

2.2 Exemples

2.2.1 Situations favorables à la création d'un schéma de données

Ces situations sont des exemples où il est pertinent de créer un schéma de données :

- Le ministère chargé des transports souhaite consolider une base nationale des lieux pouvant servir de points de covoiturage. Les collectivités territoriales sont en charge de la création, recensement et aménagement de ces lieux.

Il est pertinent de créer un schéma de données car un grand nombre de producteurs de données doivent produire le même jeu de données. Un schéma facilitera la diffusion des prérequis, permettra la validation des données et facilitera l'agrégation nationale.

- L'INSEE souhaite diffuser le Code Officiel Géographique. Il rassemble des données sur des communes, des cantons, des arrondissements, des départements, des régions et des pays. Ce fichier est actualisé tous les ans.

Il est pertinent de créer un schéma car ces données sont des données de référence. Un grand nombre de réutilisateurs est susceptible d'utiliser ces données. Il est primordial que ces réutilisateurs aient accès à une documentation de qualité, que la structure des fichiers des données reste stable dans le temps et que les données mises à disposition soient de bonne qualité.

2.2.2 Situations où le référencement d'un schéma sur [schema.data.gouv.fr](#) ne semble pas nécessaire

Ces situations sont des exemples où il ne semble pas pertinent de créer ou diffuser un schéma :

- Une administration centrale diffuse des statistiques d'activité d'un bureau, en open data, de manière annuelle.

Avec ces seules informations, la création d'un schéma ne semble pas nécessaire. En effet, il n'y a qu'un seul producteur et le potentiel de réutilisation semble limité.

Bénéfices des schémas de données en interne

Bien qu'il ne paraisse pas nécessaire dans certaines situations de créer et diffuser un schéma, vous pouvez choisir de le faire. En effet, les schémas de données comportent de nombreux avantages (documentation, montée en qualité, réutilisations, etc.) qui sont bénéfiques, même lorsque les données sont utilisées uniquement en interne.

2.3 Points de sortie

À l'issue de cette phase, vous devriez :

- Connaître les schémas de données;
- Être en mesure de décider si votre projet requiert la création d'un schéma de données;
- Savoir si votre schéma de données devra être référencé à terme sur schema.data.gouv.fr.

3 Phase de concertation

Lexique : Phase de concertation

La phase de concertation est la phase centrale de la création d'un schéma de données. C'est l'étape où plusieurs parties prenantes (producteurs, réutilisateurs, experts métiers et techniques) se rassemblent pour définir et spécifier les éléments essentiels à la constitution de ce schéma.

Pour spécifier un schéma de données, il est nécessaire de définir :

- les champs;
- les types associés de ces champs (une date, un nombre, une chaîne de caractère etc.);
- les contraintes de chaque champ (entier positif, texte dans une liste fermée etc.);
- la description de chaque champ;
- une documentation associée au schéma de données décrivant le contexte, les acteurs, les cas d'usages.

3.1 Procédure de collaboration

Nous conseillons de mener cette phase de concertation en travaillant sur un document partagé, accessible en ligne, tel qu'un [Framapad](#) ou Google Doc. L'important est que plusieurs contributeurs puissent contribuer (modifier ou mettre des commentaires) sans avoir besoin d'être présents physiquement ou de recevoir des versions intermédiaires par e-mails.

En complément de ce document partagé, nous vous conseillons d'organiser plusieurs réunions afin de débattre du schéma de données à produire. L'implication d'une multitude d'acteurs est clé : vous devez rassembler des producteurs, experts métiers, experts techniques et réutilisateurs. La richesse des profils et des enjeux permettra d'aboutir à une solution la plus adaptée.

Référencer votre schéma

Référencer votre schéma sur schema.data.gouv.fr vous permettra de bénéficier de conseils de la part d'Etalab et de ses partenaires institutionnels et associatifs. Découvrez comment [référencer votre schéma en cours de concertation](#).

3.2 Grands principes

Nous avons listé ci-dessous plusieurs conseils qui vous permettront de construire un schéma de données de qualité.

- **Profiter de l'existant.** De nombreux standards existent déjà, qu'ils concernent des formats de données ou des formats de champs. Certains standards sont devenus incontournables aujourd'hui, comme [ISO-8601](#) pour les dates ou [WGS 84](#) pour les coordonnées géographiques.
- **Identifier et associer l'écosystème.** Les personnes/organisations que vous associez sont la meilleure garantie d'un schéma de données efficace et largement adopté, permettant d'aboutir à un véritable standard.
 - Les producteurs d'une part qui connaissent la réalité de leurs données, de la collecte, etc. et qui ont leurs propres usages.
 - Les usagers d'autre part, leurs besoins et leurs difficultés d'autres part, qu'ils soient déjà connus, « sous le radar » ou en devenir.
- **Prendre le temps.** Un schéma de données est susceptible de concerner beaucoup de producteurs et d'usagers. Sa modification peut avoir un impact important. Il est donc crucial de prendre le temps d'obtenir tous les retours avant de publier un schéma utilisable par le plus grand nombre. Un schéma de données devrait être publié quand il est prêt, non pas en fonction d'un impératif de délai.

- **Lever les implicites et les ambiguïtés.** Le diable est dans les détails... Toutes les spécifications d'un schéma de données doivent être les plus claires possibles, y compris pour des cas/données qui n'existent pas encore mais pourraient apparaître à l'avenir.
- **Éviter la redondance mais sans l'exclure absolument.** Trois champs pour définir une latitude et une longitude (latitude, longitude, lat-lon) est inutilement redondant. Toutefois, préciser le nom d'une commune en plus de son code INSEE rend les données plus faciles à lire et à exploiter.
- **Utiliser des données pivot relevant d'un référentiel ouvert** pour relier les données à d'autres données, par exemple l'utilisation du numéro SIREN pour identifier des organisations. Ce principe permet aussi d'éviter l'abondance de détails et d'aller à l'essentiel : l'obtention d'informations complémentaires se fera par le biais d'un autre référentiel.

Exemples à votre disposition

Vous pouvez parcourir des fichiers de schémas sur schema.data.gouv.fr pour faciliter votre travail. Consultez par exemple [le schéma des lieux de stationnement](#).

En complément, nous avons rédigé [un guide dédié à la préparation de jeux de données](#) qui pourrait vous être utile pour définir votre schéma.

3.3 Points de sortie

À l'issue de cette phase, vous devriez :

- Avoir réuni divers partenaires afin de collaborer sur votre schéma de données ;
- Avoir décidé des différents champs de votre schéma de données, leurs types et définitions et produit une documentation associée.

4 Phase de construction

Lexique : Phase de construction

La phase de construction consiste à implémenter techniquement le schéma de données obtenu après la phase de concertation. Pour cela, il est nécessaire de choisir un standard technique, créer les fichiers requis, les tester et les diffuser.

Durant cette phase, vous devez mobiliser des personnes possédant des compétences techniques. Cette phase consiste à transcrire les décisions prises lors de la phase de concertation pour un schéma de données.

4.1 Choisir un standard technique pour la description de votre schéma de données

Lexique : Standard

On utilise les termes « normes » et « standards » pour décrire un référentiel commun et documenté destiné à harmoniser l'activité d'un secteur.

Il existe plusieurs standards techniques pour les schémas de données. Le standard est à choisir en fonction de la nature des données concernées et des habitudes de l'écosystème produisant ou réutilisant les données liées au schéma.

Les principaux standards techniques sont les suivants :

- **Table Schema** : adapté pour la description de données tabulaires (sous forme de tableurs ou de CSV). Ce standard technique utilise le format JSON
- **JSON Schema** : adapté pour la description de données avec une notion de hiérarchie. Ce standard utilise le format JSON
- **XML Schema Definition (XSD)** : adapté pour la description de données avec une notion de hiérarchie. Ce standard utilise le format XML

Notez que tous ces standards techniques sont supportés par schema.data.gouv.fr.

Aller au-delà de la documentation texte

Un schéma de données décrit uniquement par du texte ou par un tableau se prive de nombreux avantages, notamment celui de l'interopérabilité entre différents systèmes informatiques.

Les schémas de données décrits par des standards techniques permettent, en plus d'une documentation textuelle ou sous forme d'un tableau, de valider que des données correspondent à un modèle de données, d'agréger des données similaires, de générer automatiquement des données respectant un schéma.

4.2 Créer votre schéma de données

Une fois un standard technique choisi, il faudra créer les fichiers requis pour modéliser vos données. La documentation de chaque standard technique décrit le contenu des fichiers à renseigner. Reportez-vous aux documentations respectives pour tirer parti des fonctionnalités avancées offertes : types de données et contraintes sur les valeurs en particulier.

Il est souvent possible de vérifier qu'un fichier correspond à un standard à l'aide d'outils en ligne ou en ligne de commande. Utilisez ces outils pour vérifier que vos productions correspondent au standard.

Exemples à votre disposition

Pour un schéma au format Table Schema, nous mettons à votre disposition [un modèle de départ](#) pour créer un dépôt Git contenant un schéma au format Table Schema.

Pour les autres formats de schémas, nous vous recommandons de consulter les schémas et dépôts Git listés sur schema.data.gouv.fr.

4.3 Documenter votre schéma de données

En complément du fichier du schéma de données, nous vous conseillons de rédiger a minima deux documents complémentaires :

- **une documentation générale** : vous indiquerez le contexte, les modalités de production des données, le cadre juridique, la finalité, les cas d'usage etc. Ce fichier est traditionnellement rédigé en Markdown et nommé `README.md` ;
- **un fichier répertoriant les changements** : permettant de suivre les modifications, d'une version à une autre. Ce fichier est traditionnellement rédigé en Markdown et nommé `CHANGE-LOG.md`.

La présence de ces fichiers représente un package complet (documentation, liste des changements et schéma de données décrit dans un standard technique), apprécié des réutilisateurs. schema.data.gouv.fr se repose sur ces éléments pour intégrer votre documentation et votre liste de changements sur une page web.

Exemples à votre disposition

Vous pouvez consulter [la documentation](#) et [la liste des changements](#) du schéma des lieux de stationnement.

4.4 Publier et diffuser votre schéma de données

Une fois votre schéma de données créé, il est nécessaire de le publier et de le diffuser pour que d'autres personnes puissent en bénéficier. Nous vous recommandons de publier vos schémas de données en tant que logiciels libres, sur votre forge de développement ou par le biais de [GitLab](#) ou [GitHub](#).

Vous bénéficierez alors des avantages habituels des dépôts de code Git en ligne : historique des modifications, fonctionnalités de tickets ou de demandes de modifications. Utilisez un compte d'organisation (dédié à votre entreprise, direction, service, ministère) et non votre compte personnel afin d'assurer une URL stable dans le temps.

Exemples à votre disposition

Vous trouverez plusieurs dépôts Git de schémas sur schema.data.gouv.fr. Consultez par exemple [le dépôt Git décrivant les lieux de stationnement](#) à l'aide d'un schéma TableSchema sur GitHub.

4.5 Référencer votre schéma de données sur schema.data.gouv.fr

Pour faciliter la découverte de votre schéma de données et des données sous-jacentes, nous vous recommandons de le faire référencer sur schema.data.gouv.fr. Nous avons rédigé [une page dédiée](#) à ce sujet décrivant les plus-values, prérequis et démarches à suivre.

4.6 Faire évoluer votre schéma de données

Une fois votre schéma de données défini et implémenté, le travail ne s'arrête pas là. Au-delà du besoin de diffusion et de promotion, il est probable que vous deviez faire des modifications : clarifications de la documentation, corrections d'erreurs, évolutions du cadre réglementaire, etc. Autant de raisons où il est nécessaire de mettre en œuvre une nouvelle version.

Posséder un dépôt Git pour votre schéma de données vous permettra d'avoir plusieurs versions et tags. Notez que schema.data.gouv.fr supporte plusieurs versions pour un même schéma de données et affiche les modifications effectuées au fur et à mesure, dès lors que ces modifications sont renseignées dans un fichier dédié.

4.7 Points de sortie

À l'issue de cette phase, vous devriez :

- Avoir implémenté votre schéma de données dans un des standards reconnus ;
- Avoir publié votre travail en ligne, dans un répertoire Git dédié ;
- Avoir pris contact avec les équipes de schema.data.gouv.fr dans le but de référencer votre schéma de données si nécessaire.

5 Intégration avec schema.data.gouv.fr

schema.data.gouv.fr est l'initiative de data.gouv.fr de référencement des schémas de données publiques pour la France. Cette plateforme de référencement nationale permet un accès aux schémas produits par différents acteurs et facilite l'intégration avec des systèmes informatiques par le biais de standards, d'URLs stables, de processus de validation et d'API.

Vous trouverez ci-dessous une capture d'écran de l'interface de schema.data.gouv.fr pour [le schéma dédié aux lieux de covoiturage](#).



FIG. 1 : Capture d'écran de l'interface de schema.data.gouv.fr

5.1 Qui peut référencer des schémas ?

Tout acteur est libre de proposer le référencement de schémas.

Concrètement, vous pouvez être une administration, une entreprise privée, une association, un citoyen etc.

5.2 Quels schémas de données sont acceptés ?

schema.data.gouv.fr accepte des schémas de données décrivant des données publiques.

Les schémas de données sont acceptés dès lors que leur l'existence est justifiée par voie :

- **réglementaire** : c'est une disposition réglementaire qui est à l'origine de la définition du schéma de données ;
- **d'usage** : la réutilisation des données décrites par le schéma bénéficie à un grand nombre ou de nombreux producteurs sont amenés à utiliser ce schéma de données.

Etalab se réserve le droit de refuser l'ajout de schémas en motivant son refus. Nous vous encourageons à [initier une discussion](#) préalablement à l'ouverture d'une *pull request*.

5.3 Quand référencer son schéma de données ?

Nous vous invitons à référencer votre schéma de données le plus tôt possible, **dès la phase d'investigation**. En référençant celui-ci en amont, vous bénéficierez de l'accompagnement d'Etalab et de partenaires tout au long de la création de votre schéma de données : de l'investigation à la publication sur schema.data.gouv.fr.

Vous pouvez référencer votre schéma de données en ouvrant un ticket sur GitHub ou en entrant en contact [avec notre équipe par e-mail](#) Nous avons créé [une page dédiée pour détailler la procédure](#). Nous tenons à jour une liste de schémas actuellement en phase d'investigation ou de construction sur cette même page.

5.4 Quels schémas de données sont acceptés ?

schema.data.gouv.fr accepte des schémas de données décrit par un standard technique (voir la page ["Phase de construction"](#) de ce présent guide). Les schémas de données décrits uniquement par de la documentation textuelle ou des tableaux ne sont pas acceptés.

5.4.1 Standards techniques supportés

Les standards techniques de schémas de données actuellement supportés sont les suivants :

- **Table Schema** : adapté pour la description de données tabulaires (sous forme de tableurs ou de CSV). Ce standard technique utilise le format JSON
- **JSON Schema** : adapté pour la description de données avec une notion de hiérarchie. Ce standard utilise le format JSON

- [XML Schema Definition \(XSD\)](#) : adapté pour la description de données avec une notion de hiérarchie. Ce standard utilise le format XML

5.4.2 Prérequis de validation des schémas de données

Lexique : Validation d'un schéma de données

La validation d'un schéma de données est l'étape qui permet de vérifier si celui-ci est conforme au standard technique sélectionné et aux prérequis de schema.data.gouv.fr. Cette étape s'intéresse uniquement au schéma de données et à la façon dont il est publié.

Il ne faut pas confondre la validation d'un schéma avec le fait de vérifier que des données correspondent à un schéma.

Pour tous les types de schéma de données, il faut que :

- votre schéma de données soit sur un dépôt Git, à raison d'un dépôt par schéma. Ce dépôt doit pouvoir être cloné depuis Internet sans authentification préalable ;
- votre dépôt Git doit comporter des tags indiquant les versions de votre schéma de données. Ces versions doivent respecter la [gestion sémantique de version semver](#), sous la forme 1 . 3 . 2 par exemple ;
- votre dépôt doit comporter un fichier README . md à la racine contenant une documentation du schéma de données indiquant par exemple le contexte de production, la gouvernance ;
- passer avec succès les tests spécifiques au type de schéma de données que votre dépôt contient.

Critères complets de validation

Cette page présente les grands principes de validation des schémas de données. Pour connaître en détail les prérequis propres à chaque type de schéma de données et accéder à des exemples, consultez [la page dédiée à la validation des schémas de données](#).

5.5 Points de sortie

À l'issue de cette phase, vous devriez :

- Avoir pris connaissance des procédures de validation en place sur schema.data.gouv.fr ;
- Avoir un dépôt Git conforme aux prérequis de schema.data.gouv.fr ;
- Avoir effectué votre demande de référencement.

6 Jeux de données sur data.gouv.fr et schémas

Une fois que votre schéma est finalisé et publié schema.data.gouv.fr, il est temps de produire des données conformes à ce schéma.

data.gouv.fr propose de multiples intégrations avec schema.data.gouv.fr permettant de spécifier qu'une ressource présente sur data.gouv.fr est censée être conforme à un schéma. Il est ensuite possible de procéder à la validation de la ressource par rapport au schéma ou de consulter la documentation du schéma à partir de la page d'un jeu de données.

6.1 Indiquer qu'une ressource respecte un schéma sur data.gouv.fr

Vous pouvez indiquer qu'une ressource d'un jeu de données correspond à un schéma depuis l'interface d'administration de data.gouv.fr. Lorsque vous déposez ou éditez une ressource, vous pouvez sélectionner le schéma correspondant à vos données depuis une liste déroulante.

Dernière version consolidée

Titre : Dernière version consolidée

Type : Fichier principal

Description : Base adresse locale
Budget des collectivités et établissements publics locaux
Catalogue simplifié
Equipements
Lieux de covoiturage
Lieux de stationnement
✓ Schéma IRVE
Schéma SCDL Délibérations
Schéma SCDL Subventions
Spécification du fichier de déclaration de profil d'acheteur

Date de publication

Schema

URL

Taille : Taille

Format : csv

Type Mime : text/csv

Somme de contrôle : sha1 Somme de contrôle

Annuler **Enregistrer**

FIG. 2 : Capture d'écran de la sélection d'un schéma depuis l'interface d'administration de data.gouv.fr

Le fait d'indiquer que votre ressource est censée respecter un schéma permet de bénéficier de vérifications de la qualité des données et d'indiquer aux réutilisateurs que vos données respectent un référentiel.

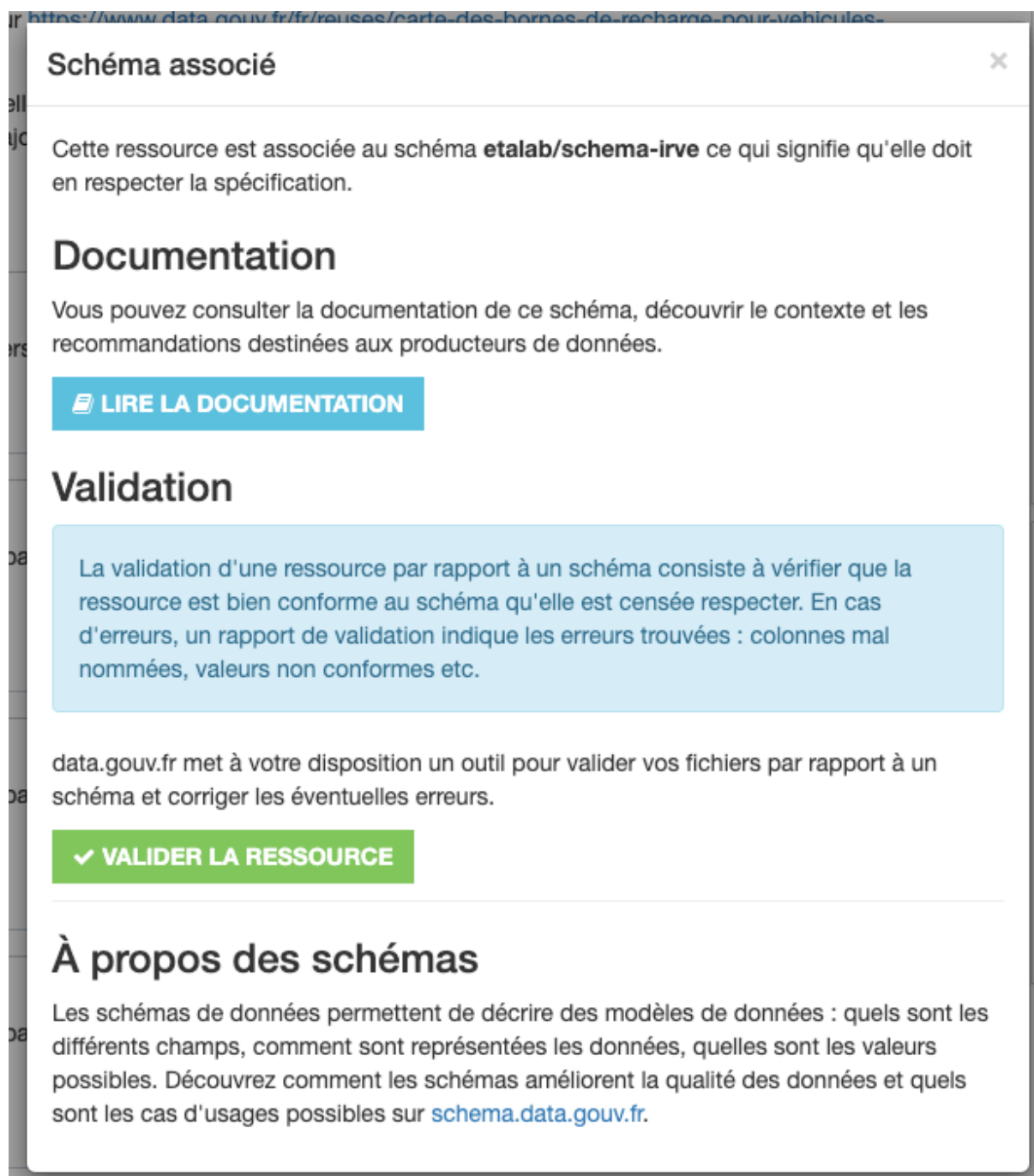


Schéma associé

Cette ressource est associée au schéma **etalab/schema-irve** ce qui signifie qu'elle doit en respecter la spécification.

Documentation

Vous pouvez consulter la documentation de ce schéma, découvrir le contexte et les recommandations destinées aux producteurs de données.

[LIRE LA DOCUMENTATION](#)

Validation

La validation d'une ressource par rapport à un schéma consiste à vérifier que la ressource est bien conforme au schéma qu'elle est censée respecter. En cas d'erreurs, un rapport de validation indique les erreurs trouvées : colonnes mal nommées, valeurs non conformes etc.

data.gouv.fr met à votre disposition un outil pour valider vos fichiers par rapport à un schéma et corriger les éventuelles erreurs.

[✓ VALIDER LA RESSOURCE](#)

À propos des schémas

Les schémas de données permettent de décrire des modèles de données : quels sont les différents champs, comment sont représentées les données, quelles sont les valeurs possibles. Découvrez comment les schémas améliorent la qualité des données et quels sont les cas d'usages possibles sur schema.data.gouv.fr.

FIG. 3 : Capture d'écran de data.gouv.fr des informations disponibles sur la page d'un jeu de données lorsqu'un schéma est spécifié sur une ressource

7 Aide à la construction d'un schéma TableSchema

La pertinence de la mise en place d'un standard de données réside dans son adéquation entre les capacités de sa mise en oeuvre par les producteurs de données et les outils permettant l'automatisation des jeux de données valides par rapport à cette spécification. Cette standardisation doit permettre de **faciliter la mise en relation des jeux de données** issus de différents producteurs.

Il ne s'agit donc pas de règles mais de recommandations, visant à faciliter la création de nouveaux schémas et **leur intégration dans une chaîne de validation et de publication généralisable**.

7.1 Recommandations pour le formatage des fichiers csv

Un des formats privilégiés pour les standards de données est le [CSV](#) (Comma Separated Values, valeurs séparées par des virgules). Il s'agit d'un format de données "à plat", **adéquat pour les structures de données simples**. Cependant ce format simple ne dispose pas de spécifications contraignant la saisie des données. Pour cela un schéma en Json est ajouté dont la structure est définie par le standard [TableSchema](#). TableSchema permet d'indiquer les formats des données attendus, de spécifier des contraintes (types de valeurs, cardinalité) et de documenter les différents champs composant le schéma.

L'outil de validation utilisé pour vérifier la conformité d'un fichier csv au standard auquel il fait référence s'appuie sur la structure tabulaire des données. Elles peuvent donc être contenues dans un tableur numérique au format .xls, .xlsx ou .ods ou dans un fichier texte au format .csv, .txt ou autre.

La question du séparateur utilisé pour séparer deux champs de données dans un fichier .csv n'est donc pas essentielle. Cependant, certains outils se basent sur la valeur de ce séparateur pour traiter et publier des jeux de données. Nous vous proposons donc un certain nombre de recommandations afin de favoriser la généralisation d'un usage contribuant à l'interopérabilité des données produites.

7.1.1 Format de fichier csv

Bien que de nombreux jeux de données en CSV utilisent le point-virgule comme séparateur de champs, il a été décidé de privilégier le **séparateur virgule** car plus conforme à l'esprit du format csv.

Les tableurs numériques courants (Excel et Calc) peuvent produire et lire des fichiers csv. Lors de l'enregistrement d'un fichier créé avec l'outil Calc, l'utilisatrice ou utilisateur doit spécifier le format d'encodage des données ainsi que le séparateur de champs. Lorsque le séparateur de champs retenu est

la virgule, il est recommandé d'utiliser les guillemets double " comme séparateur de chaîne de caractères. De cette manière, si une virgule est présente à l'intérieur d'une cellule elle ne sera pas considérée comme un séparateur de champs.

Lors de l'ouverture d'un fichier csv dans Calc, une fenêtre modale propose plusieurs options permettant de spécifier un caractère de séparation et un encodage des données.

Dans Excel, il faut aller dans l'onglet données et sélectionner l'option Fichier texte pour accéder à l'assistant d'import des données.

L'encodage des caractères à privilégier est l'[UTF-8](#) de manière à garantir une **meilleure interopérabilité des données**.

Pour faciliter la lecture des fichiers publiés en CSV il est recommandé d'y associer dans les outils de publication le **type MIME ou Content-Type "text/csv"**.

Chaque ligne du fichier doit avoir le même nombre de champs, ce qui signifie que lorsqu'une cellule est vide elle doit quand même être présente soit avec la valeur Null, soit avec des crochets vides [] dans le cas des champs de type tableau (array), soit laissée vide mais apparaître à l'export avec 2 virgules qui se suivent „ .

7.2 Recommandations de formatage des données

Les recommandations de formatage pour les données sont généralement issues du standard [TableSchema](#), lui-même inspiré des spécifications du format [Json](#), dans lequel sont exprimés les schémas de données permettant l'automatisation de leur validation.

Ce standard dispose des types de données suivants :

- **string** : s'applique pour toutes les chaînes de caractères
- **number** : s'applique pour les chiffres et nombres contenant éventuellement des décimales
- **integer** : s'applique pour les chiffres et nombres entiers
- **boolean** : s'applique pour indiquer que la valeur d'un champs ne peut être égale qu'à "vrai" ou "faux" (ou "1" et "0" ou "oui" ou "non")
- **object** : s'applique pour les données de type objet
- **array** : s'applique pour les tableaux de données

Les types de données peuvent être assortis de formats de données facilitant l'automatisation de leur validation.

Pour déclarer un format de données dans un schéma JSON il est possible d'utiliser différentes propriétés permettant de le caractériser :

- **name** : le nom du champ
- **title** : le titre du champ
- **description** : la description des valeurs attendues dans ce champ
- **format** : le format du champ
- **type** : le type du champ

Il est également possible de contraindre les valeurs autorisées dans ce champ à l'aide de plusieurs propriétés :

- **required** : indique l'obligation de la présence d'une valeur pour ce champ dans toutes les lignes du fichier
- **unique** : indique que chaque valeur de ce champ à l'intérieur du fichier doit être unique
- **minLength** : indique la taille minimale des valeurs de ce champ
- **maxLength** : indique la taille maximale des valeurs de ce champ
- **minimum** : indique la valeur minimum autorisée pour ce champ (par exemple pour une date on peut indiquer une année en deça de laquelle les valeurs ne sont pas autorisées)
- **maximum** : indique la valeur maximale autorisée pour ce champ
- **pattern** : indique une expression régulière à laquelle doivent être conforme les valeurs de ce champ (par exemple pour un numéro SIRET on peut indiquer `^\d{14}$` ce qui signifie que les valeurs de ce champ doivent contenir exactement 14 chiffres)
- **enum** : indique une liste de valeurs autorisées pour ce champ

Ci-dessous quelques exemples tirés du schéma des menus de la restauration collective.

Le champ permettant d'indiquer le numéro SIRET d'une collectivité est spécifiée de la manière suivante

```
{
  "name": "menuCollSiret",
  "title": "Code SIRET de la collectivité qui produit les données.",
  "description": "Identifiant du Système d'Identification du Répertoire des Etablissements",
  "type": "string",
  "examples": "21330063500017",
  "constraints": {
    "required": true,
    "pattern": "^\d{14}$"
  }
}
```

Le champ permettant d'indiquer la date de publication d'un enregistrement du jeu de données est spécifié de la manière suivante :

```
{
  "name": "menuPublicationDate",
  "title": "Date de publication de l'enregistrement d'un menu",
  "description": "Lors de la publication ce champ d'horodatage permet d'indiquer la",
  "type": "datetime",
  "examples": "2020-05-11T14:08:32Z",
  "constraints": {
    "required": true
  }
}
```

Les informations ci-dessous décrivent les différents types de champs disponibles dans la spécification TableSchema.

7.2.1 Données de type string

Pour le type string, les formats de données suivants sont disponibles :

- **default** : n'importe quelle chaîne de caractère
- **email** : une adresse email valide.
 - motif de validation :
- **uri** : une URI valide
- **binary** : une chaîne de caractère encodées en base 64 représentant des données binaires.
- **uuid** : une chaîne de caractère représentant un identifiant unique.

7.2.2 Données de type décimal

- **Description** : Les valeurs décimales doivent utiliser le point afin d'être plus facilement exploitables par les tableurs numériques.
- **Type** : number
- **Exemple** : 3900.50

7.2.3 Données de type date

- **Description** : date au format AAAA-MM-JJ suivant la norme internationale [ISO 8601](#).
- **Type** : date
- **Exemple** : 2017-10-15

- **Format** : “%Y-%m-%d”
- **Nommage** : abbreviation-du-schemaDate

7.2.4 Données de type date avec heure

- **Description** : date au format aaaa-mm-jjThh :mi :ssZZZZZZ suivant la norme internationale [ISO 8601](#). On considérera que ZZZZZZ (+ou- décalage horaire GMT), est par défaut +01 :00 en France et qu’il est inutile de le préciser dans les formats.
- **Type** : datetime
- **Exemple** : 1997-07-16T19 :20 :00

7.2.5 Données de type date avec heure de début et de fin

- **Description** : date au format aaaa-mm-jjThh :mi/hh :mi suivant la norme internationale ISO 8601. Ce type de données s’applique pour un créneau horaire dans la même journée, sans les secondes. Pour une extension de ces conditions, voir la norme [ISO 8601](#).
- **Type** : datetime
- **Exemple** : 1997-07-16T08 :30/17 :30

7.2.6 Données de type horaires d’ouverture

- **Description** : horaires indiquant les heures d’ouverture d’un service ou d’un commerce. Ce type de données permet de préciser les différents horaires d’ouverture pour les différents jours de la semaine. Il s’agit donc d’un type de données multi-valeur au sein duquel le nom du jour de la semaine est abrégé et suivi par les heures d’ouvertures. Les abréviations pour les jours sont en anglais (Mo, Tu, We, Th, Fr, Sa, Su) et les horaires sont sous la forme HH :MM

Un assistant graphique en ligne [yohours](#) permet de générer simplement cette structure de données

- **Type** : string (chaîne de caractères)
- **Exemple** : Mo 08 :15-13 :15; Tu 03 :15-06 :15; We 03 :15-09 :30; Th 02 :30-07 :15; Fr 01 :30-05 :45; Sa 00 :30-05 :00; Su 02 :45-08 :30
- **Nommage** : abbreviation-du-schemaHoraires

7.2.7 Données de type géolocalisation

La possibilité est laissée de décrire les points de géolocalisation d'une donnée à l'intérieur d'un champ unique (geopoint) ou à l'aide de 2 champs (latitude et longitude).

7.2.7.1 Latitude

- **Description** : ce type de données permet de saisir la coordonnée de latitude exprimée en [WGS 84](#) permettant de localiser un équipement. Le signe de séparation entre les parties entière et décimale du nombre est le point. Précision : 6 décimales maximum.
- **Type** : number
- **Exemple** : 48.563433
- **Nommage** : abbreviation-du-schemaLat

7.2.7.2 Longitude

- **Description** : ce type de données permet de saisir la coordonnée de longitude exprimée en [WGS 84](#) permettant de localiser un équipement. Le signe de séparation entre les parties entière et décimale du nombre est le point. Précision : 6 décimales max.
- **Type** : number
- **Exemple** : 2.572875
- **Nommage** : abbreviation-du-schemaLon

7.2.7.3 Geopoint

- **Description** : ce type de données permet de saisir les coordonnées de latitude et de longitude exprimée en [WGS 84](#) permettant de localiser un équipement. Le signe de séparation entre les parties entière et décimale du nombre est le point. Précision : 6 décimales max. Le séparateur de valeur est la virgule. Il est donc nécessaire d'entourer ces valeurs de guillemets. La première valeur est la latitude
- **Type** : number
- **Exemple** : "48.563433, 2.572875"
- **Nommage** : abbreviation-du-schemaGeo

7.2.7.4 Geoshape

- **Description** : ce type de données permet de décrire la forme géographique d'un équipement. La forme est décrite à l'aide de paires de coordonnées, séparées par un espace vide et chaque paire séparée par une virgule. La description d'une ligne est exprimée à l'aide de 2 ou plus paires de

points séparés par des virgules. La description d'un polygone est exprimée par 4 ou plus paires de points séparés par des virgules dont la dernière est identique à la première.

- **Type** : string
- **Exemple** : "48.563433 2.572875, 49.234933 2.134432, 49.885311 2.134003, 48.974635 2.1134567, 48.563433 2.572875"

7.2.8 Données de type adresse

Ce type de champ permet de décrire l'adresse postale d'un équipement. Il est décomposé entre 3 champs permettant de distinguer et de faciliter le tri à l'intérieur des informations de voirie, de code postal et de commune. Le numéro et le nom de la voie sont séparés par une virgule.

7.2.8.1 Voie

- **Description** : ce type de champs permet de saisir le numéro et le nom de la voie
- **Type** : string
- **Exemple** : 34, rue de Latresne
- **Nommage** : abreviation-du-schemaVoie ##### Code postal ou Code INSEE
- **Description** : ce type de champs permet de saisir le code postal (ou le code INSEE) de la commune
- **Type** : number
- **Exemple** : 45800
- **Nommage** : abreviation-du-schemaCodePostal

7.2.8.2 Commune

- **Description** : ce type de champs permet de saisir le nom de la commune
- **Type** : string
- **Exemple** : Saint-Jean-de-Braye
- **Nommage** : abreviation-du-schemaCommune

7.3 Recommandations de champs obligatoires

Afin d'unifier la description des données au travers des différentes thématiques abordées par le propositions de standard de données, **il est fortement recommandé de rendre obligatoire la présence d'un certains nombre de champs**. Ceux-ci contribuent à la **portabilité des données** (qui produit la donnée) ou à **leur fiabilité** (quand a été produite la donnée)

7.3.1 Identification du producteur

Pour l'identification des autorités publiques à l'origine de la production et de la publication des jeux de données, il est recommandé d'indiquer le nom et le numéro de siret sur chaque ligne de chaque jeu de données.

7.3.1.1 Nom de la collectivité

- **Description** : ce champs permet de saisir le nom de l'autorité publique responsable de la production des données
- **Type** : string
- **Exemple** : Conseil départemental de la Creuse
- **Nommage** : abreviation-du-schemaColl

Par exemple

```
{
  "name": "menuCollNom",
  "title": "Nom de la collectivité qui produit les données",
  "description": "Nom officiel de la collectivité ou de l'établissement public resp",
  "type": "string",
  "examples": "Grand Poitiers Communauté urbaine",
  "constraints": {
    "required": true
  }
}
```

7.3.1.1.1 Siret de la collectivité

- **Description** : ce champ permet d'indiquer le numéro d'identification de l'autorité publique au sein de la base nationale des établissements.
- **Type** : string
- **Exemple** : 21330063500017
- **Motif** : `^\d{14}$`
- **Nommage** : nom-ou-abreviation-du-schemaCollSiret

Par exemple :

```
{
  "name": "menuCollSiret",
```

```

    "title": "Code SIRET de la collectivité qui produit les données.",
    "description": "Identifiant du Système d'Identification du Répertoire des Etabli
    "type": "string",
    "examples": "21330063500017",
    "constraints": {
      "required": true,
      "pattern": "^\\d{14}$"
    }
  }
}

```

7.3.1.2 Horodatage des données Pour faciliter la réutilisation et la mise à jour des données, il est recommandé de fournir aux réutilisatrices et réutilisateurs potentiels des date de première publication et de dernière modification pour chaque entité du jeu de données.

Ces informations au format Date avec horaire peuvent correspondre à la date de première publication et faire apparaître les dates de dernière modification pour l'ensemble des lignes ou en cas de mise à jour partielle pour une ligne de données particulière.

Il est également recommandé d'y associer un champ permettant de décrire la raison ayant entraîné une mise à jour des données depuis leur publication

7.3.1.2.1 Date de création/publication

- **Description** : ce champs permet de décrire la date de première publication de la donnée
- **Type** : datetime
- **Exemple** : 2020-05-11T14:08:32Z
- **Nommage** : nom-ou-abreviation-du-schemaPublicationDate

Par exemple :

```

{
  "name": "menuPublicationDate",
  "title": "Date de publication de l'enregistrement d'un menu",
  "description": "Lors de la publication ce champ d'horodatage permet d'indiquer la
  "type": "datetime",
  "examples": "2020-05-11T14:08:32Z",
  "constraints": {
    "required": true
  }
}

```

7.3.1.2.2 Date de dernière modification

- **Description** : ce champs permet de décrire la date de dernière modification de la donnée
- **Type** : datetime
- **Exemple** : 2020-05-11T14:08:32Z
- **Nommage** : nom-ou-abreviation-du-schemaModificationDate

Par exemple :

```
{
  "name": "menuModificationDate",
  "title": "Date de dernière modification de l'enregistrement d'un menu",
  "description": "Lors de la modification ce champ d'horodatage permet d'indiquer l",
  "type": "datetime",
  "examples": "2020-05-11T14:08:32Z",
  "constraints": {
    "required": false
  }
}
```

7.3.1.2.3 Information sur les modifications

- **Description** : ce champs permet de décrire la raison d'une modification de la donnée depuis sa publication initiale
- **Type** : string
- **Exemple** : changement dû à un aléa de livraison
- **Nommage** : nom-ou-abreviation-du-schemaModificationInfo

Par exemple :

```
{
  "name": "menuModificationInfo",
  "title": "Information sur la modification ayant entraîné une mise à jour de la don",
  "description": "Afin de renseigner les usagers de la donnée, il est possible de pr",
  "type": "string",
  "examples": "changement dû à un aléa de livraison",
  "constraints": {
    "required": false
  }
}
```

7.4 Recommandations pour le nommage des fichiers

Les fichiers doivent, sauf exception et autant que possible, respecter les règles de nommage suivantes :

AAAAMMJJ_idProducteur_nom-du-fichier.extension

- **AAAAMMJJ** : Date de création du fichier
- **idProducteur** : Numéro [SIREN](#) sur 9 chiffres pour identifier le producteur
- **nom-du-fichier** Chaîne de caractères dont les termes, en minuscules non accentuées, sont séparés par un tiret du milieu
- **.extension** : Si les règles de formatage sont respectées, l'extension est .csv

Les 3 éléments constitutifs de la chaîne principale avant l'extension sont assemblés en un seul tenant et séparés par un tiret du bas.

- **Exemple** : '20180314_213502388_prenoms-nouveaux-nes-rennes-2017.csv'

7.5 Recommandations pour le nommage des champs

Afin d'uniformiser les fichiers produits dans le cadre de schémas de standardisation, il est recommandé de normaliser les intitulés des champs composant chaque standard.

La règle générale préconisée est l'utilisation de l'écriture camelCase où chaque mot composant l'intitulé du champ est écrit avec une majuscule à l'exception du premier. En complément il est recommandé d'utiliser un préfixe (mot complet ou abréviation) pour l'ensemble des champs d'un standard. En conséquence pour le standard des menus les intitulés des champs sont préfixés par le mot menu suivi des intitulés à proprement dit. Par exemple :

- menuCollNom
- menuRestaurantIdType
- menuRepasType

Aucun caractère accentué ou spécial ne doit être utilisé dans l'intitulé d'un champ. Il est également préconisé de ne pas dépasser 50 caractères pour l'intitulé d'un champ et d'utiliser le singulier pour les mots composant l'intitulé du champ.

7.6 Recommandations pour la mise en conformité

Pour garantir la conformité des jeux de données, il est demandé aux producteurs de s'assurer que la structure, les champs et les contenus attendus sont effectivement respectés.

De fait, les fichiers tabulaires doivent, autant que possible, contenir :

- **Toutes les colonnes**, y compris celles dont les cellules ne sont pas renseignées, dans le bon ordre, et avec des en-têtes correctement nommées sur la première ligne
- **Autant de lignes que nécessaire** comprenant des cellules dont les valeurs peuvent être **obligatoires** (elles doivent être impérativement renseignées) ou **optionnelles** (elles sont seulement recommandées ou soumises à condition de disponibilité / pertinence)