

# Open Data : Microcosme technologique

19 Octobre 2022 – Séance n°3

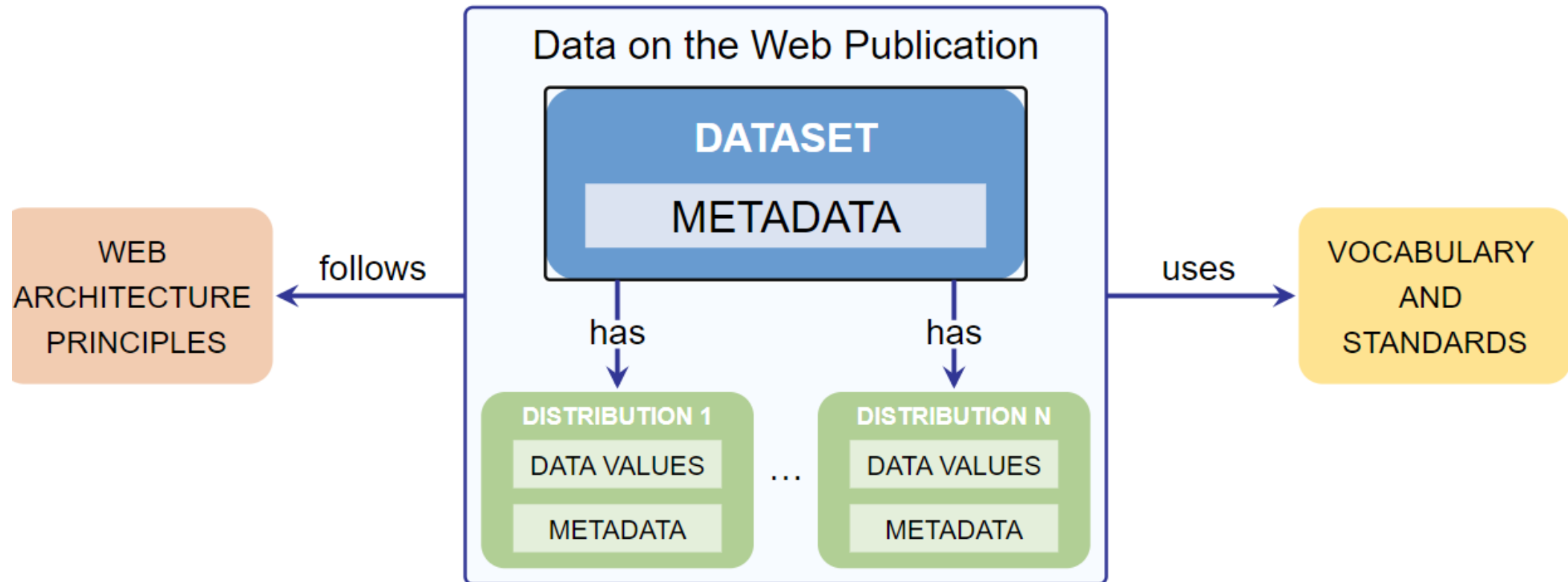


Shared Prosperity **Dignified Life**



# Définition d'un Jeu de données

W3C recommendation : Data on the Web Best Practices <https://www.w3.org/TR/dwbp/>



## Best Practices Summary

- [Best Practice 1](#): Provide metadata
- [Best Practice 2](#): Provide descriptive metadata
- [Best Practice 3](#): Provide structural metadata
- [Best Practice 4](#): Provide data license information
- [Best Practice 5](#): Provide data provenance information
- [Best Practice 6](#): Provide data quality information
- [Best Practice 7](#): Provide a version indicator
- [Best Practice 8](#): Provide version history
- [Best Practice 9](#): Use persistent URIs as identifiers of datasets
- [Best Practice 10](#): Use persistent URIs as identifiers within datasets
- [Best Practice 11](#): Assign URIs to dataset versions and series
- [Best Practice 12](#): Use machine-readable standardized data formats
- [Best Practice 13](#): Use locale-neutral data representations
- [Best Practice 14](#): Provide data in multiple formats
- [Best Practice 15](#): Reuse vocabularies, preferably standardized ones
- [Best Practice 16](#): Choose the right formalization level
- [Best Practice 17](#): Provide bulk download
- [Best Practice 18](#): Provide Subsets for Large Datasets
- [Best Practice 19](#): Use content negotiation for serving data available in multiple formats
- [Best Practice 20](#): Provide real-time access
- [Best Practice 21](#): Provide data up to date
- [Best Practice 22](#): Provide an explanation for data that is not available
- [Best Practice 23](#): Make data available through an API
- [Best Practice 24](#): Use Web Standards as the foundation of APIs
- [Best Practice 25](#): Provide complete documentation for your API
- [Best Practice 26](#): Avoid Breaking Changes to Your API
- [Best Practice 27](#): Preserve identifiers
- [Best Practice 28](#): Assess dataset coverage
- [Best Practice 29](#): Gather feedback from data consumers
- [Best Practice 30](#): Make feedback available
- [Best Practice 31](#): Enrich data by generating new data
- [Best Practice 32](#): Provide Complementary Presentations
- [Best Practice 33](#): Provide Feedback to the Original Publisher
- [Best Practice 34](#): Follow Licensing Terms
- [Best Practice 35](#): Cite the Original Publication

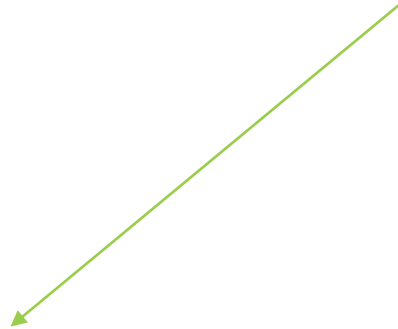
# Les bonnes pratiques W3C

« Les données doivent être **compréhensibles** et **faciles à découvrir** par les **humains** et les **machines** »

<https://www.w3.org/TR/dwbp>

"Data on the Web Best Practices" fournit les meilleures pratiques liées à la publication et à l'utilisation de données sur le Web conçues pour aider à soutenir un écosystème autonome.

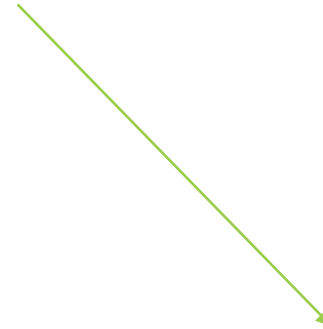
# Google Dataset Search Bêta



Schema.org



Data Catalog Vocabulary (DCAT)



CSV on the Web

# Microcosme technologique

Identifier  
& Prioriser

Extraire,  
Convertir &  
Anonymiser

Documenter &  
Fournir les  
métadonnées

Publier le  
JDD

Actualiser  
le JDD

Valoriser &  
Socialiser



# Identifier les données & Prioriser l'ouverture

# Identifier les données & Prioriser l'ouverture

## Identifier ?

La réussite d'une initiative d'ouverture des données se mesure par le degrés d'ouverture des données **à fort impact**, dont le potentiel de **réutilisation** est **très élevé** !

**La théorie** : Plus il y a de données ouvertes, mieux c'est.

**Le pragmatisme** : On n'a pas assez de ressources pour maintenir ouvert un nombre important de jeux de données!

**La solution** ➔ Il faut **prioriser** !

Pour **prioriser**, il est important de raisonner

- ☐ de point de vue **producteurs** de données,
- ☐ et de point de vue des **ré-utilisateurs** des données.

# Identifier les données & Prioriser l'ouverture

## Prioriser ?

### Acteur public

- ✓ Accroître la transparence pour plus de redevabilité,
- ✓ Atteinte des objectifs stratégiques/opérationnels,
- ✓ Satisfaire les obligations et exigences légales,
- ✓ Optimisation et Modernisation des services publics,
- ✓ Rationalisation des dépenses publiques,

### Ré-utilisateur

- ✓ C'est la génération d'une **valeur économique et/ou sociale**.
- ✓ La taille de la communauté des utilisateurs finaux.
- ✓ La nature des systèmes et/ou services qui vont consommer ce jeu de données.





# Extraire, Convertir & Anonymiser

# Extraire, Convertir & Anonymiser

**L'extraction de données** : est un processus automatisé de collecte ou de récupération de données à partir d'une ou bien de différentes sources, dont certaines peuvent être non structurées ou mal structurées. L'extraction des données permet de consolider les données pour faciliter leur transformation.

L'extraction des données est la première étape des processus ETL (Extraction, Transformation, Chargement)

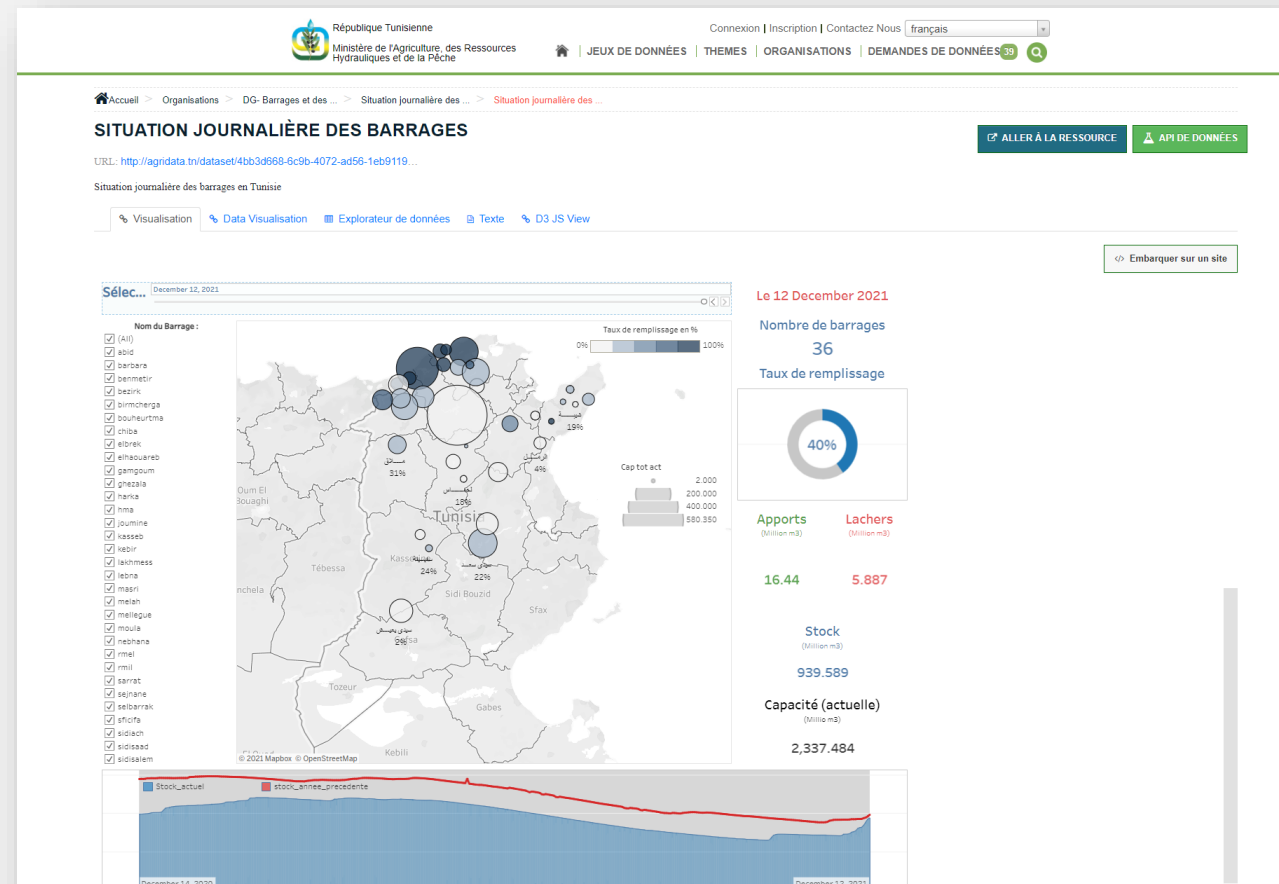
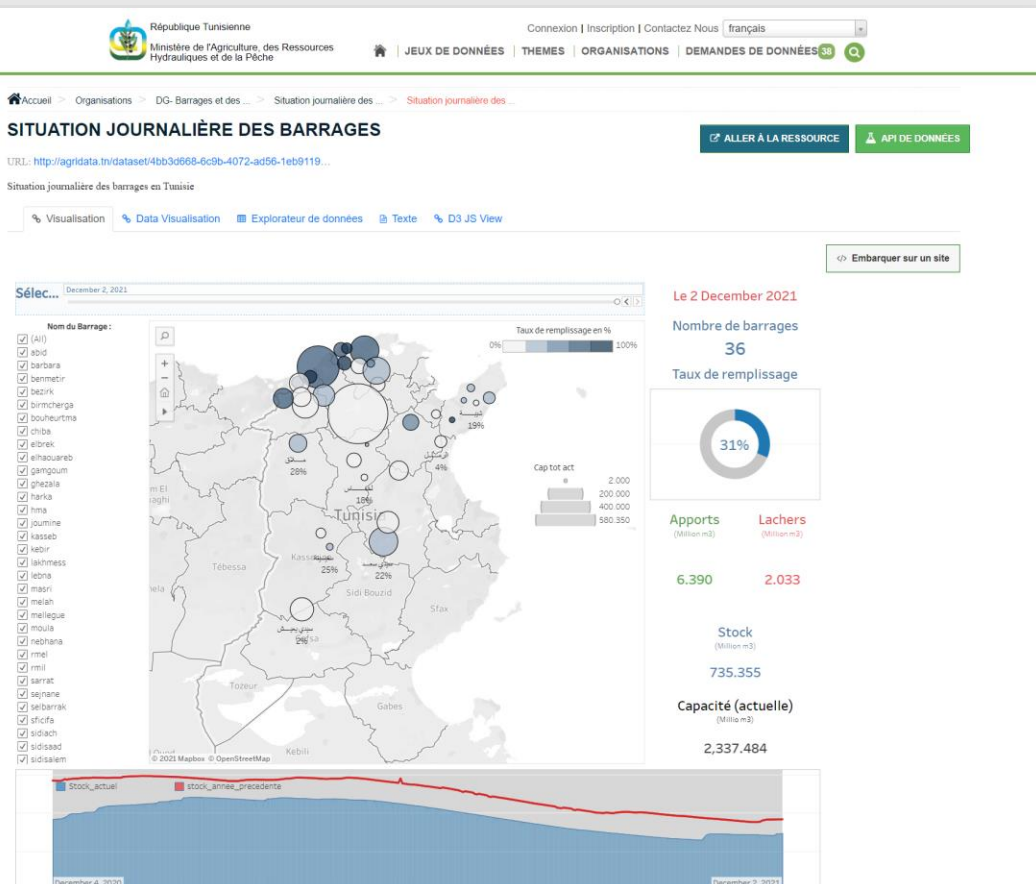
En open data on parle souvent d'ETP (Extraction, Transformation, Publication)



La solution Apache Airflow permet de créer, de planifier et de surveiller des workflows de données. Il s'agit d'une solution totalement open source, très utile pour l'architecture et l'orchestration de pipelines de données complexes et pour lancer de tâches.

# Extraire, Convertir & Anonymiser

## Extraire les données ?



# Extraire, Convertir & Anonymiser

**Conversion** : Transformation des données d'une structure initiale vers une structure recommandée!

Il est important de :

- ☐ S'assurer que le processus de transformation ne conduit pas à une perte de qualité des données.
- ☐ Automatiser l'opération d'extraction → Cela permet de faciliter par la suite la mise à jour (pérenniser l'ouverture).
- ☐ Nettoyer: Utiliser un formatage standardisé de données (sans entrées manquantes, des données dans chaque champ et avec le moins d'erreurs possibles).
- ☐ Standardiser : Permet d'assurer une **interopérabilité** optimale des données, en vue de leur réutilisation. Elle permet aussi d'automatiser le traitement des données,

# Extraire, Convertir & Anonymiser

## Nettoyer les données :

The screenshot shows an Excel spreadsheet with the following annotations:

- 1: Yellow background area at the top of the sheet.
- 2: Yellow background area on the left side of the sheet.
- 3: Arabic title 'السنة' in the header row.
- 4: Space in the header row between 'année 2019' and 'annee\_2019'.
- 5: Space in the header row between 'annee\_2019' and 'minimum dinar la balle'.
- 6: Formula '3+4+6' in a cell.
- 7: Orange background area for data fusion.
- 8: Multiple columns with the same title 'annee'.
- 9: Comma as a decimal separator '3,5'.
- 10: Empty cell.
- 11: Multiple columns with the same title 'rendement'.

السنة	année 2019	annee_2019	minimum dinar la balle	rendement	rendement	poids
						14
						16
						0
						2
						13
						14
						16
						0
						2
						16
						18
						0
						12
						12
						12
						12

- 1- éviter l'espace vide avant le tableau
- 2- une seule feuille par fichier
- 3- éviter le titre en arabe dans les colonnes
- 4- éviter les accents et l'espace dans les titres
- 5- remplacer l'espace par le tiret (8)
- 6- éviter les données traitées et les formules
- 7- éviter la fusion des cases
- 8- un seul tableau par fichier
- 9- remplacer la virgule par un point
- 10- éviter les vides dans les cases (mais ne pas mettre zéro)
- 11- éviter les colonnes avec le même titre
- 12- Pas de titre ou bien du texte avant le tableau

# Extraire, Convertir & Anonymiser

## Standardiser les données :

Les codes de pays sont au format ISO 3166

 Morocco	The Kingdom of Morocco	UN member state	MA	MAR	504	ISO 3166-2:MA	.ma
---	------------------------	-----------------	----	-----	-----	---------------	-----

Codes de langue sont au format ISO 639

Les devises sont au format ISO 4217

MAD	504	2	Moroccan dirham	 Morocco
-----	-----	---	-----------------	---

Les dates sont au format ISO 8601 : « 2018-01-01T15:00:00.15+02:00 »

Les coordonnées géographiques sont au format WGS 84 : « +37.5665,+126.9780 »

## Un standard ouvert ?

GTFS « General Transit Feed Specification », OCDS « Open Contracting Data Standard », etc...

Plus de 100 standards ouverts :

<https://docs.google.com/spreadsheets/d/1r7OByH4leFH Zot43nayjlpIgEHHV9114uBIUn59SKgU/edit#gid=0>

# Extraire, Convertir & Anonymiser

## Standardiser les données :

Name	Abbrev.	URL	Publisher	Topic	Subtopic
General Bikeshare Feed Specification	GBFS	<a href="https://github.com/NABSA/gbfs">https://github.com/NABSA/gbfs</a>	North American Bike Share Association	Transportation	Bike Sharing
Service Interface for Real Time Information	SIRI	<a href="http://www.transmodel-cen.eu/standards/siri/">http://www.transmodel-cen.eu/standards/siri/</a>	CEN/TS	Transportation	Real-Time Transit
GTFS-realtime		<a href="https://developers.google.com/transit/gtfs-realtime/">https://developers.google.com/transit/gtfs-realtime/</a>	Google	Transportation	Real-Time Transit
Transit Communications Interface Profiles	TCIP	<a href="https://www.apta.com/research-technical-resources/sta">https://www.apta.com/research-technical-resources/sta</a>	American Public Transportation Association	Transportation	Real-Time Transit
Open511		<a href="http://www.open511.org/">http://www.open511.org/</a>	Open North	Transportation	Road Construction
Traffic Management Data Dictionary	TMDD	<a href="http://www.ite.org/standards/TMDD/">http://www.ite.org/standards/TMDD/</a>	Institute of Transportation Engineers	Transportation	Road Construction
Vocabulario para la representación de datos sobre tráfico		<a href="http://vocab.linkeddata.es/datosabiertos/def/transporte/">http://vocab.linkeddata.es/datosabiertos/def/transporte/</a>		Transportation	Traffic
General Transit Feed Specification	GTFS	<a href="https://developers.google.com/transit/gtfs/">https://developers.google.com/transit/gtfs/</a>	Google	Transportation	Transit Schedules
Vocabulario para la representación de medios de transporte		<a href="http://vocab.linkeddata.es/datosabiertos/def/urbanismo-">http://vocab.linkeddata.es/datosabiertos/def/urbanismo-</a>		Transportation	
Vocabulario para la representación de medios de transporte públicos		<a href="http://vocab.linkeddata.es/datosabiertos/def/transporte/">http://vocab.linkeddata.es/datosabiertos/def/transporte/</a>		Transportation	
Vocabulario para la representación de direcciones postales en España		<a href="http://vocab.linkeddata.es/datosabiertos/def/urbanismo-">http://vocab.linkeddata.es/datosabiertos/def/urbanismo-</a>		Urban Planning	Addresses
OpenAddresses		<a href="https://github.com/openaddresses/openaddresses/blob">https://github.com/openaddresses/openaddresses/blob</a>	OpenAddresses	Urban Planning	Addresses
Addresses		<a href="http://localgovdigital.github.io/localo/specifications/addr">http://localgovdigital.github.io/localo/specifications/addr</a>	LocalGov Digital	Urban Planning	Addresses
Building & Land Development Specification	BLDS	<a href="http://permitdata.org/">http://permitdata.org/</a>	OpenPermit Foundation	Urban Planning	Permits
Application Tracking Data Interchange Specification	ATDIS	<a href="https://www.planningalerts.org.au/atdis/specification">https://www.planningalerts.org.au/atdis/specification</a>	OpenAustralia Foundation	Urban Planning	Permits

# Extraire, Convertir & Anonymiser

## Standardiser les données :

### Comment choisir un bon standard ?

- Ce standard peut-il être utilisé par quiconque (existe-t-il une licence ouverte)?
- Ce standard est-il conçu pour répondre à mes besoins?
- Ce standard est-il activement maintenu ?
- Y a-t-il suffisamment de guides pour comprendre son utilisation ?

➔ Il faut savoir quand utiliser des standards ouverts et quand **ne pas créer de nouveaux standards**.

➔ Le guide intitulé « **Comment une proposition de normes ouvertes est évaluée** » [\*], publié par le gouvernement britannique est très utile pour répondre à certaines questions.

(\*) <https://www.gov.uk/guidance/choosing-open-standards-for-government>



# Extraire, Convertir & Anonymiser

## Anonymiser les données :

L'anonymisation est un traitement qui consiste à utiliser un ensemble de techniques de manière à rendre impossible, en pratique, toute identification de la personne par quelque moyen que ce soit et de manière irréversible. [\*]

### **Anonymiser != Pseudonymiser**

Remplacer les données directement identifiantes (cin, matricule unique, nom, prénom, etc.) d'un jeu de données par des données indirectement identifiantes (alias, numéro séquentiel, etc.).

Cette technique permet le traitement des données d'individus sans pouvoir identifier ceux-ci de façon directe. Toutefois il est bien souvent possible de retrouver l'identité de ceux-ci grâce à des données tierces (en faisant un recouplement avec d'autres sources de données).

C'est une mesure préventive mais non suffisante !

(\*) <https://www.cnil.fr/fr/lanonymisation-de-donnees-personnelles>

# Extraire, Convertir & Anonymiser

## Anonymiser les données :

L'objectif étant de garantir à la fois :

- La **non sensibilité** des données à ouvrir,
- La **préservation de l'utilité** des données à ouvrir.

## Comment concevoir un procédé d'anonymisation :

**L'individualisation** : il ne doit pas être possible d'isoler un individu dans un jeu de données,

**La corrélation** : il ne doit pas être possible de relier entre eux des jeux de données distincts concernant un même individu,

**L'inférence** : il ne doit pas être possible de déduire, de façon quasi certaine, de nouvelles informations sur un individu.

(\*) <https://www.cnil.fr/fr/le-g29-publie-un-avis-sur-les-techniques-danonymisation>

# Extraire, Convertir & Anonymiser

## Anonymiser les données :

Les données sensibles ?

Matricule Unique	Date de naissance	Sexe	Adresse	Salaire en MAD
3694521	01/01/1971	Femme	1 Rue Sidi Saleh El Kram 2015, <b>Préfecture d'Oujda-Angad</b>	1751
3683561	01/03/1971	Femme	Immeuble El nozha, Appart. N°4, Carthage 2089 <b>Préfecture d'Oujda-Angad</b>	1809
3694511	13/05/1971	Femme	4 Cité administrative Cité El Khadhra, 1003 <b>Préfecture d'Oujda-Angad</b>	2450
3683572	05/06/1971	Femme	1 Rue Ammar El Hajji El Menzah 1013 <b>Préfecture d'Oujda-Angad</b>	2015
3694522	13/05/1975	Femme	15, Avenue des Finances Cité administrative 2080 <b>Province de Berkane</b>	1089
3683563	11/05/1975	Femme	4, Rue Chakib Arselène Cité Etadhamen 241 <b>Province de Berkane</b>	2400
3694512	16/06/1975	Femme	5, Rue du 18 janvier 1952 Ariana Center 2080 <b>Province de Berkane</b>	2010
3683573	13/04/1977	Homme	15, Avenue des Finances Cité administrative 2080 <b>Préfecture de Fès</b>	1876
3694523	11/07/1977	Homme	4, Rue Chakib Arselène Cité Etadhamen 241 <b>Préfecture de Fès</b>	1845
3683563	18/08/1977	Homme	5, Rue du 18 janvier 1952 Ariana Center 2080 <b>Préfecture de Fès</b>	1956
3683574	13/05/1978	Homme	15, Avenue Hedi Chaker Route Tunis Km 7 Sakiet Ezzit, 3021 <b>Province d'El Hajeb</b>	1655
3694524	11/06/1978	Homme	10, Rue Arbi Zarrouk , 3029 <b>Province d'El Hajeb</b>	1655
3683564	16/09/1978	Homme	3, Avenue de l'environnement Bir Ali Ben Khelifa, 3040 <b>Province d'El Hajeb</b>	1655

# Extraire, Convertir & Anonymiser

## Anonymiser les données :

### Anonymiser ?

Ré-identification par LATANYA SWEENEY en 1997, des données médicales du gouverneur du Massachusetts « William Weld » .

*K-anonymity*

***K = 3***

*L-diversity*

***L = 1***

Pseudonyme	Année de naissance	Sexe	Préfecture / Province	Salaire en MAD
1000001	1971	Femme	Préfecture d'Oujda-Angad	1751
1000002	1971	Femme	Préfecture d'Oujda-Angad	1809
1000003	1971	Femme	Préfecture d'Oujda-Angad	2450
1000004	1971	Femme	Préfecture d'Oujda-Angad	2015
1000005	1974	Femme	Province de Berkane	1089
1000006	1974	Femme	Province de Berkane	2400
1000007	1974	Femme	Province de Berkane	2010
1000008	1977	Homme	Préfecture de Fès	1876
1000009	1977	Homme	Préfecture de Fès	1845
1000010	1977	Homme	Préfecture de Fès	1956
1000011	1978	Homme	Province d'El Hajeb	1655
1000012	1978	Homme	Province d'El Hajeb	1655
1000013	1978	Homme	Province d'El Hajeb	1655

[https://epic.org/wp-content/uploads/privacy/reidentification/Sweeney\\_Article.pdf](https://epic.org/wp-content/uploads/privacy/reidentification/Sweeney_Article.pdf)

# Extraire, Convertir & Anonymiser

Anonymiser les données :

Anonymiser ?

*K-anonymity*

$K = 6$

*L-diversity*

$L = 4$

Identifiant	Quasi-identifiant			Sensible
Pseudonyme	Année de naissance	Sexe	Région	Salaire en TND
1000001	[1971-1974]	Femme	Oriental	1751
1000002	[1971-1974]	Femme	Oriental	1809
1000003	[1971-1974]	Femme	Oriental	2450
1000004	[1971-1974]	Femme	Oriental	2015
1000005	[1971-1974]	Femme	Oriental	1089
1000006	[1971-1974]	Femme	Oriental	2400
1000007	[1971-1974]	Femme	Oriental	2010
1000008	[1975-1979]	Homme	Fès-Meknès	1959
1000009	[1975-1979]	Homme	Fès-Meknès	1617
1000010	[1975-1979]	Homme	Fès-Meknès	2011
1000011	[1975-1979]	Homme	Fès-Meknès	1655
1000012	[1975-1979]	Homme	Fès-Meknès	1655
1000013	[1975-1979]	Homme	Fès-Meknès	1655



# Documenter & Fournir les métadonnées

# Documenter & Fournir les métadonnées

## Fournir les métadonnées :

### Métadonnées descriptives pour un jeu de données :

Titre, description, mots clés, date de publication, producteur, point de contact, couverture spatiale, période temporelle, date de la dernière modification, thème/catégorie, etc...

Un contexte bilingue, multilingue (titre en arabe, en français, etc... )

<https://www.w3.org/TR/dwbp/#DescriptiveMetadata>

# Documenter & Fournir les métadonnées

## Machine-readable metadata

[DCTERMS] [VOCAB-DCAT] [DCAT-AP]

```
a dcat:Dataset;
dcterms:title "Bus stops of MyCity";
dcat:keyword "transport", "mobility", "bus";
dcterms:issued "2015-05-05"^^xsd:date;
dcat:contactPoint <http://data.mycity.example.com/transport/contact>;
dcterms:temporal <http://reference.data.gov.uk/id/year/2015>;
dcterms:spatial <http://www.geonames.org/3399415>;
dcterms:publisher <http://data.mycity.example.com/transport-agency-mycity>;
dcterms:accrualPeriodicity <http://purl.org/linked-data/sdmx/2009/code#freq-A>;
dcat:theme :mobility;
dcterms:language <http://id.loc.gov/vocabulary/iso639-1/en>, <http://id.loc.gov/vocabulary
dcterms:conformsTo "ISO 8601";
dcat:distribution <http://data.mycity.example.com/transport/dataset/bus/stops-2015-05-05.c
dcat:distribution <http://data.mycity.example.com/transport/dataset/bus/stops-2015-05-05.j
dcterms:creator <http://data.mycity.example.com/transport/people/john>;
owl:versionInfo "1.0";
pav:version "1.0".

# A mini SKOS concept scheme to cover the themes
:mobility a skos:Concept;
  skos:inScheme :themes;
  skos:prefLabel "Mobility".

:themes a skos:ConceptScheme ;
  skos:prefLabel "A set of domains to classify documents".

# Facts about the transport agency
<http://data.mycity.example.com/transport-agency-mycity> a foaf:Organization, prov:Agent;
  foaf:name "MyCity Transport Agency".
```

## Human-readable metadata

HTML / TXT file

### Bus stops of MyCity

This is the human-readable version of examples used in the [DWB document](#). Please note that, as this is a fictional example, hyperlinks to the data.mycity.example.com domain simply link to example.com which defers references to an explanation that the domain is designed for use in examples and serves no other function.

#### Dataset description

Title	Bus timetable of MyCity
URI	<a href="http://data.mycity.example.com/transport/dataset/bus/stops-2015-05-05">http://data.mycity.example.com/transport/dataset/bus/stops-2015-05-05</a>
Keywords	transport, mobility, bus
Publication date	2015-05-05
Publisher	Transport Agency MyCity
Creator	Adrian < <a href="mailto:adrian@mycitytransport.org">adrian@mycitytransport.org</a> >
Contact point	<a href="http://data.mycity.example.com/transport/contact">http://data.mycity.example.com/transport/contact</a>
Period that the dataset covers	The British calendar year of 2014
Spatial coverage	Fortaleza, Brazil
Update frequency	Annual
Theme	Mobility
Language	English, Portuguese
Date and time formats	ISO 8601
Current version	1.2

#### Dataset distributions

##### RDF Distribution

Title	RDF distribution of stops-2015-05-05 dataset
Description	RDF distribution of the stops dataset of MyCity.
Media type	text/turtle
License	CC BY-SA 3.0
Publication date	2015-05-05
Last modification	2015-05-05
Download URL	<a href="http://data.mycity.example.com/transport/dataset/bus/stops-2015-05-05.ttl">http://data.mycity.example.com/transport/dataset/bus/stops-2015-05-05.ttl</a>

##### CSV Distribution

Title	CSV distribution of stops-2015-05-05 dataset
-------	--



# Documenter & Fournir les métadonnées

## Fournir les métadonnées :

### Métadonnées structurelles pour les distributions :

Fournir des informations sur la structure interne d'une distribution → Dictionnaire des données

stop_id	stop_name	stop_description	stop_lat	stop_lon	zone_id	stop_url
F12	5 Av/53 St	Short text description here ...	40.760167	-73.975224	1	<a href="http://transport-com">http://transport-com</a>
E1	5 Av/53 St SW	It is a long established fact that	40.760474	-73.976099	2	<a href="http://transport-com">http://transport-com</a>
E2	5 Av/53 St NE	Contrary to popular belief, Lorem	40.76035	-73.97546	2	<a href="http://transport-com">http://transport-com</a>
E3	5 Av/53 St SE	There are many variations of pa	40.760212	-73.975512	2	<a href="http://transport-com">http://transport-com</a>
E4	Madison/53 St NE	Sed ut perspiciatis unde omnis is	40.759612	-73.973731	2	<a href="http://transport-com">http://transport-com</a>

<https://www.w3.org/TR/dwbp/#StructuralMetadata>

# Documenter & Fournir les métadonnées

## Fournir les métadonnées :

Métadonnées structurelles "*Human-readable*" :

### Structural metadata

Field	Titles	Description	Datatype	Primary key	Required
stop_id	Identifier	An identifier for the bus stop.	string	true	true
stop_name	Name	The name of the bus stop.	string		
stop_desc	Description	A description for the bus stop.	string		
stop_lat	Latitude	The latitude of the bus stop.	numeric		
stop_long	Longitude	The longitude of the bus stop.	numeric		
zone_id	Zone	An identifier for the zone where the bus stop is located.	string		
stop_url	URL	URL that identifies the bus stop.	string		

<https://www.w3.org/TR/dwbp/#StructuralMetadata>

# Documenter & Fournir les métadonnées

Fournir les métadonnées :

Métadonnées structurales

*Machine-readable*

```
"tableSchema": {  
  "columns": [{  
    "name": "stop_id",  
    "titles": "Identifiant",  
    "dct:description": "An identifier for the bus stop.",  
    "datatype": "string",  
    "required": true  
  }, {  
    "name": "stop_name",  
    "titles": "Name",  
    "dct:description": "The name of the bus stop.",  
    "datatype": "string"  
  }, {  
    "name": "stop_desc",  
    "titles": "Description",  
    "dct:description": "A description for the bus stop.",  
    "datatype": "string"  
  }, {  
    "name": "stop_lat",  
    "titles": "Latitude",  
    "dct:description": "The latitude of the bus stop.",  
    "datatype": "number"  
  }, {
```

# Documenter & Fournir les métadonnées

## Fournir les métadonnées :

### Fournir des informations sur la licence :

En open data si un jeu de données n'est pas associé à une licence ouverte, alors les données ne sont pas considérées véritablement **ouvertes**.

#### *Les licences Creative Commons « CC »*

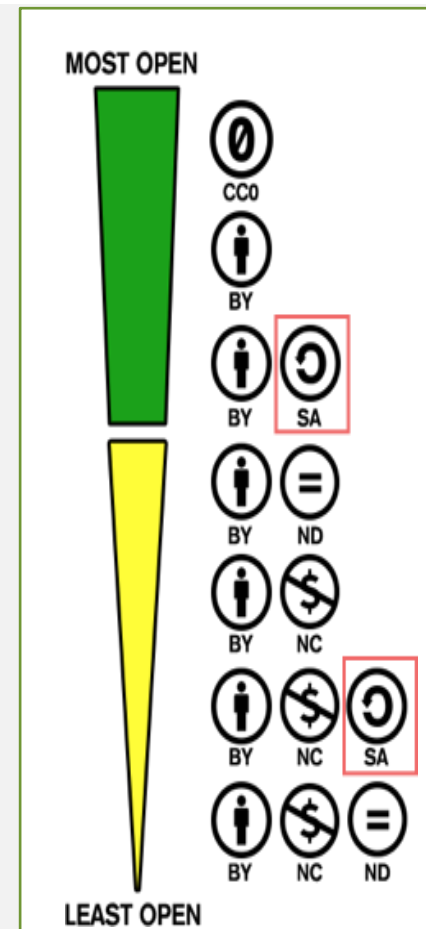
- Paternité « Attribution » (CC-BY v4.0) :
- Paternité et Partage à l'identique ('Attribution and share-alike'): (CC-BY-SA v4.0)

« Si vous souhaitez partager à l'identique non seulement les données dérivées, mais **aussi les œuvres produites à partir de ces données.** »

#### *L'Open Data base Licence « ODbL »*

« Elle **ne** couvre **que** les bases de données. Son partage à l'identique se propage **uniquement** à des bases de données, mais pas aux œuvres produites à partir de ces bases de données. »

Contrairement la licence **CC-BY-SA**





# Publier un jeu de données

# Publier un jeu de données JDD

## Variante de publication ?

Trois (03) variantes pour la publication d'un jeu de données :

- ⇒ De point de vue général, si vous disposez d'un faible nombre de jeux de données avec une faible fréquence de mise à jour, la saisie manuelle directe sur l'interface web de la plateforme nationale open data vous sera recommandée. Ce scénario est recommandé aussi si le producteur des données ne dispose pas de systèmes d'information ou bien si l'environnement ne favorise pas l'extraction et la mise à jour automatique.
- ⇒ Si vous disposez de plus que 10 jeux de données avec une fréquence de mise à jour régulière, l'approche automatisée vous sera conseillée. Cette approche se base sur l'exploitation des APIs (d'écriture) offerts par la plateforme nationale ou bien une autre plateforme sectorielle (à condition que cette dernière soit référencée sur la plateforme nationale).
- ⇒ Dans le cas où le nombre de jeux de données est important avec une actualisation fréquente, il serait recommandé d'étudier la mise en place d'une plateforme spécifique d'open data avec la mise en place de la technique de moissonnage des données vers la plateforme nationale.

# Publier un jeu de données JDD

## Format ouvert ?

Un jeu de données doit être enregistré dans un format ouvert !

⇒ Les données tabulaires

⇒ Le format le plus répandu est le CSV,


⇒ Il est recommandé de prévoir deux formats (wide & long)

⇒ Les données hiérarchiques


⇒ Les formats les plus répandus JSON, XML et GeoJSON, KML

⇒ Les données en réseau : RDF, etc...

Region	2019	2020	2021
Nord	100 000	105 000	95 000
Centre	120 000	102 000	98 000
Sud	80 000	90 000	96 000



Region	Année	Production
Nord	2019	100 000
Centre	2019	120 000
Sud	2019	80 000
Nord	2020	105 000
Centre	2020	102 000
Sud	2020	90 000
Nord	2021	95 000
Centre	2021	98 000
Sud	2021	96 000



# Publier un jeu de données JDD

## Encodage du fichier ?

L'encodage d'un fichier est la norme utilisée pour coder chaque caractère par une suite compréhensible par une machine.

Lorsque l'encodage est mal choisi, le ré-utilisateur des données est souvent contraint de convertir le fichier, notamment afin de faire apparaître les accents et caractères spéciaux.

→ Il est conseillé d'utiliser l'encodage UTF-8.



# Publier un jeu de données JDD

## Automatisation ?

APIs d'écriture - Qu'est-ce que la plateforme permet et qu'est-ce qu'elle ne permet pas ?

- ☐ Les requêtes POST créent des ressources enfants sur une URI définie,
- ☐ Les requêtes PUT créent ou remplacent la ressource de l'URI définie par l'utilisateur,
- ☐ Les requêtes PATCH mettent à jour des parties de la ressource de l'URI définie par l'utilisateur.

Il ne faut pas changer les identifiants à chaque opération de mise à jour (persistance des URIs)

Si on utilise CKAN , on peut mettre à jour un fichier de données (ressource) :

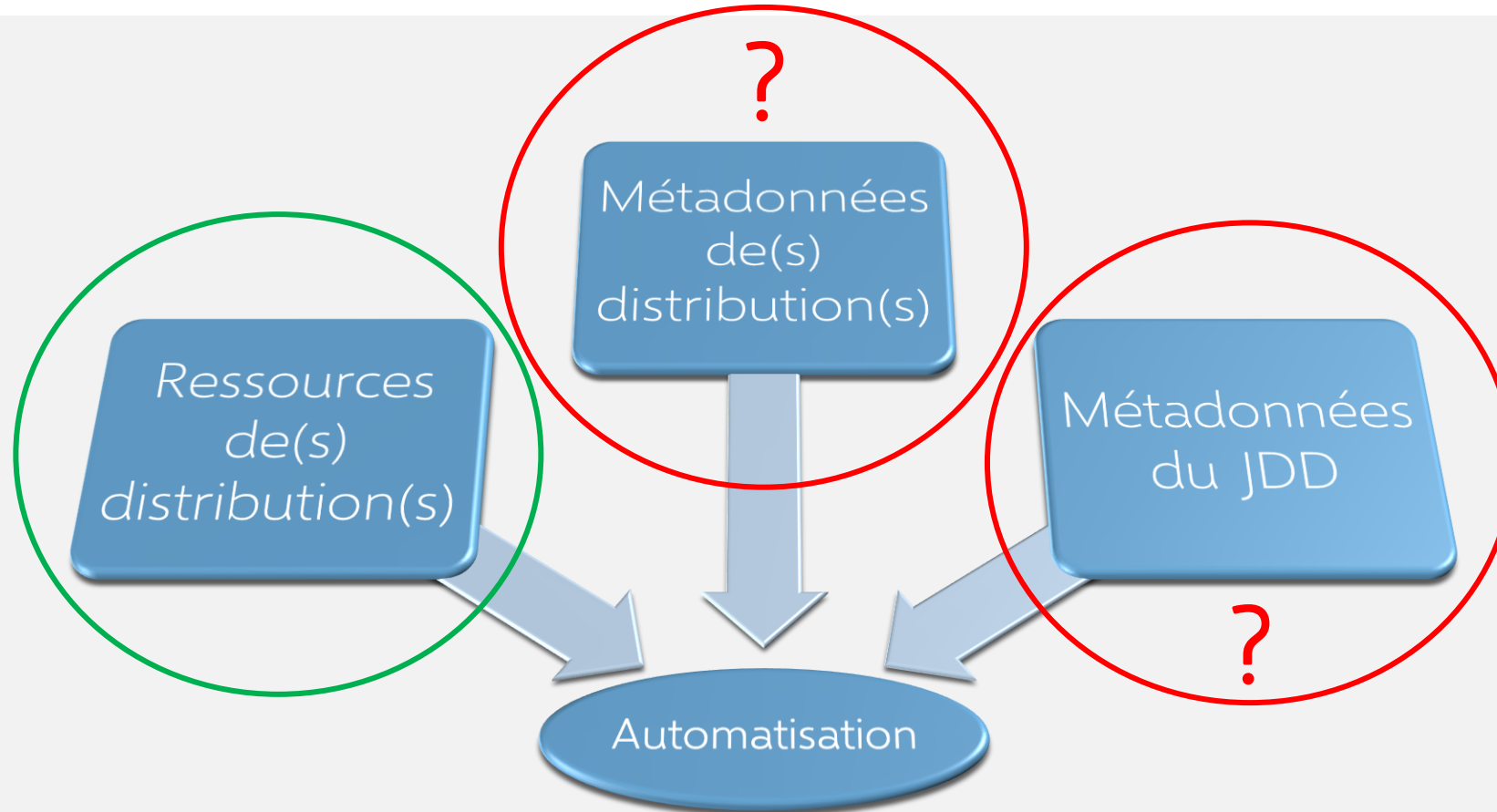
```
curl -X POST -H "Content-Type: multipart/form-data" -H "Authorization: XXXX" -F "id=<resource_id>" -F "upload=@updated_file.csv" https://<URL_platform>/api/3/action/resource_patch
```



# Actualiser un jeu de données

# Actualiser un jeu de données

Automatisation ? Ne pas se limiter seulement, à la mise à jour des fichiers des données !





Shared Prosperity **Dignified Life**



# Discussion