

ID	N° d'ordre	Domaine de problèmes concernés	Type de données	Type de problème	Exemple	Détecter le problème	Ligne de commande pour le contrôle	Réduit considérablement l'usage des données	Niveau d'importance du problème	Temps d'analyse en min	Temps d'analyse automatisable	Solution	Temps de traitement	Remarques		
1	1	Jeu de données & métadonnées	Toutes	Le jeu de données est dans un format "image" ne permettant pas de manipuler les données	Le jeu de données est un fichier image au format JPEG ou PDF.	Ouvrir le fichier et tenter de copier/coller les données		oui	1	1	1	* Demander au producteur une version qui permette de manipuler les données (CSV, Excel, etc.) * Essayer une phase d'OCR du document				
2	2	Jeu de données & métadonnées	Toutes	Le jeu de données est dans un format non spécifiquement adapté aux données - PDF, Word, ODF, epub, HTML, SVG, etc.	Le jeu de données est un fichier HTML.	Déterminer le format du fichier		oui	1	1	1	* Dans certains cas la méthode du scrapping est une solution.	Les formats PDF ou de traitement de texte rendent l'exploitation des données difficiles.			
3	3	Jeu de données & métadonnées	Toutes	Le format du jeu de données n'est pas précisé (fichier CSV, TXT, etc.)	L'extension du jeu de données ne permet pas de savoir quel logiciel permet de l'ouvrir et l'obtenir ni pas fournir d'indication complémentaire	Essayer d'ouvrir le fichier ?			1	1	1	0	* Demander au producteur * Rétro-documenter le format			
31	4	Jeu de données & métadonnées	Toutes	La documentation et les métadonnées sont quasi inexistantes voire absentes	La documentation tient sur 5 lignes alors que le fichier est très complexe	Lecture des métadonnées		oui	1	1	1	0	* Produire les métadonnées * Faire renseigner ou valider les métadonnées par le producteur	C'est important de voir ça très en amont. Si l'on veut corriger cela, certains points de contrôle participent à la documentation		
115	5	Jeu de données & métadonnées	Toutes	La documentation et les métadonnées sont d'un usage officieux (doc papier, doc au format PDF image, doc uniquement en anglais, etc.)	La documentation est fournie sous forme de PDF image, les usagers ne peuvent pas rechercher des termes pour y naviguer rapidement	Consultation de la doc et des métadonnées			1	1	1	0	* Demander au producteur			
4	6	Jeu de données & métadonnées	Toutes	La licence du jeu de données ne nous permet pas de l'utiliser	Le jeu de données est un fichier commercial que l'on n'a pas acheté	En cas de doute, demander au producteur d'où viennent les données			1	15	0	0				
5	7	Jeu de données & métadonnées	Toutes	Le format du jeu de données n'est pas ouvert	Le fichier est au format .xls ou .xlsx	Ne pas seulement se baser sur l'extension mais ouvrir également le fichier	file nom du fichier		1	1	1	1	1	* Vérifier que le document existe dans un format ouvert * Convertir le document dans un format ouvert		
6	8	Jeu de données & métadonnées	Toutes	Le format du jeu de données ne permet pas d'ouvrir le fichier dans des outils très répandus (Excel, Notepad, ...)	Le fichier au format csv d'ouvre mal dans Excel, outil le plus répandu pour ouvrir des tableaux	Essayer d'ouvrir le fichier ?			1	1	1	1	1			
7	9	Jeu de données & métadonnées	Toutes	L'encodage du fichier n'est pas spécifique (ISO-8859-1, UTF-8, etc.)	Le fichier contient des caractères exotiques mais on ne sait pas s'il s'agit d'un problème d'encodage	Lecture des métadonnées			1	1	1	1	1	0	* Produire les métadonnées * Faire renseigner ou valider les métadonnées par le producteur * Rétro-documenter le format iconv -f ISO-8859-1 -t UTF-8 -c caractéristiques_2015.csv -c caractéristiques_2015_rev1.csv	
8	10	Jeu de données & métadonnées	Toutes	L'encodage n'est pas en UTF-8 : ce dernier devient le format du fichier et autres encodages peuvent engendrer des problèmes	L'encodage est en ISO-8859-1	Ouvrir le fichier avec un éditeur qui spécifie l'encodage	perl -f //; alias 'print' 'f. 1 print if (m{[a-z]}); cat 'file.csv'   csvlook		1	1	1	1	1	1		
9	11	Jeu de données & métadonnées	Toutes	L'encodage n'est pas homogène	Certains données sont correctement encodées et d'autres contiennent des caractères exotiques	Ouvrir le fichier et parcourir visuellement les données ; rechercher quelques chaînes comme "A0r" ou "r" ou "T0r" ou "n" ou "na" etc.			1	1	1	1	1	1		
10	12	Jeu de données & métadonnées	Toutes	Le fichier est mal formé	Pour certaines lignes, parfois une colonne manque ou le fichier csv comporte des "virgule" non formatées et empêche l'ouverture correcte du fichier	Ouvrir le fichier, lire la dernière colonne du fichier et regarder le résultat	catview -dry -n file.csv	oui	1	1	1	1	1	1		
11	13	Jeu de données & métadonnées	Toutes	Le jeu de données concernant des horaires de mode de transport ne possède pas de version au format GTFS	Le fichier n'est pas au format GTFS				2	1	1	1	1	1		
12	14	Jeu de données & métadonnées	Toutes	Le jeu de données concernant des sauveurs n'est pas au format Dublin Core	Le fichier n'est pas au format Dublin Core				2	1	1	1	1	1		
13	15	Jeu de données & métadonnées	Toutes	Le jeu de données utilise une norme peu accessible au plus grand nombre (outil, complexe)	Le jeu de données est au format TRIDENT				2	1	0	0	0	0		
14	16	Jeu de données & métadonnées	Toutes	Le processus d'acquisition n'est pas connu	Wiki leaks	Lecture des métadonnées			2	1	0	0	0	0	Évaluer le processus et vérifier un échantillon de données	Redondant avec "Métadonnées imprécises : process et contexte de production non explicites" ?
15	17	Jeu de données & métadonnées	Toutes	L'échantillon n'est pas documenté	L'échantillon semble représentatif mais on ne peut pas vérifier qu'il le soit bien, puisque ce dernier n'est pas documenté	Lecture des métadonnées			2	1	0	0	0	0		
16	18	Jeu de données & métadonnées	Toutes	Le format d'un des champs n'est pas documenté, si on ne peut pas comprendre ce qu'il est difficile ou bien contrôler ses valeurs	"La date est parfois exprimée par le nombre de secondes depuis 1970 ; cette donnée est difficile à comprendre." "Un jeu de données contient un champ "image" en binaire, dont le format n'est pas spécifique.	Lecture des métadonnées et ouverture du fichier : le format champ binaire est-il documenté ?			2	3	0	0	0	0	* Demander au producteur * Rétro-documenter le format	
17	19	Jeu de données & métadonnées	Toutes	La taille maximale d'un champ n'est pas documentée	On ne sait pas à un code peu dépasser 10 caractères et si certaines valeurs sont donc fausses	Lecture des métadonnées			2	3	0	0	0	0	* Demander au producteur * Rétro-documenter le format	
18	20	Jeu de données & métadonnées	Toutes	Pour tel champ, l'incertitude de la mesure n'est pas connue (pas de marge d'erreur (précision à 0 m, à 100 m ?) ; la précision d'une mesure de température n'est pas explicitée (-0.1, 0.1, -1, +1 ?))	Des coordonnées GPS sont indiquées mais on ne connaît pas leur marge d'erreur (précision à 0 m, à 100 m ?) ; la précision d'une mesure de température n'est pas explicitée (-0.1, 0.1, -1, +1 ?)	Lecture des métadonnées			3	1	0	0	0	0		
117	21	Jeu de données & métadonnées	Toutes	L'incertitude de la mesure n'est pas connue par le producteur	Le producteur des données ne connaît pas la précision de ses mesures	Si l'incertitude de la mesure n'est pas documentée (D18), demander au producteur			3	15	0	0	0	0		
22	22	Jeu de données & métadonnées	Toutes	La précision n'est pas cohérente avec la granularité : l'incertitude de la mesure est 100 fois supérieure à la granularité	Des coordonnées géographiques annoncent une granularité au cm alors que l'incertitude des appareils de mesure est de +/- 5 mètres	Lecture des métadonnées			3	0	0	0	0	0		
19	23	Jeu de données & métadonnées	Toutes	L'origine de certaines données est une entrée manuelle non contrôlée	Le risque est d'obtenir 25 orthographe de "Saint-André-des-Arts"	Ouvrir le jeu de données et parcourir : des données sont-elles manifestement entrées à la main ?			3	3	0	0	0	0		
20	24	Jeu de données & métadonnées	Toutes	Les données proviennent d'un processus de reconnaissance automatique dont la marge d'erreur est globalement bonne mais localement problématique (OCR, reconnaissance de forme, géocodage, etc.)	OCR : reconnaissance automatique des voyages (va dépendre de la qualité de la lumière de la prise de vue, de la couleur des personnes concernées (c'est encore un problème en 2016)) etc.	Ouvrir le jeu de données et parcourir : des données sont-elles manifestement issues d'un processus de reconnaissance automatique ?			3	3	0	0	0	0		
21	25	Jeu de données & métadonnées	Toutes	L'échantillon est biaisé	Certaines populations sont absentes, sur-représentées ou sous-représentées ; les données subissent une forte variation saisonnière	WTF.csv			3	0	0	0	0	0		
23	26	Jeu de données & métadonnées	Toutes	Le processus de signalement d'erreur et d'échange avec le producteur n'est pas explicite	Aucune forme de contact n'est donnée	Lecture des métadonnées			2	1	0	0	0	0	0	* Produire les métadonnées * Faire renseigner ou valider les métadonnées par le producteur
27	27	Jeu de données & métadonnées	Toutes	Le processus de signalement d'erreur et d'échange avec le producteur n'est pas ou bien il est défilant	Le producteur ne répond pas aux questions	Demander au producteur			15	0	0	0	0	0		
25	28	Jeu de données & métadonnées	Toutes	La disponibilité de la donnée n'est pas documentée (temps pendant lequel la donnée est accessible par rapport au temps total souhaité, généralement exprimé en pourcentage)	L'utilisateur ne sait pas si la qualité de service est de 95% ou 99.9% ; si le système qui héberge la donnée est régulièrement inaccessible (maintenance, etc.), les usagers devaient en être informés pour savoir si leur usage est impacté	Lecture des métadonnées			15	0	0	0	0	0	<a href="https://fr.wikipedia.org/wiki/Disponibilit%C3%A9">https://fr.wikipedia.org/wiki/Disponibilit%C3%A9</a>	
26	29	Jeu de données & métadonnées	Toutes	La disponibilité de la donnée n'est pas mesurée	Le producteur ne sait pas si la qualité de service est de 95% ou 99.9% alors que le futur usage est critique	Demander au producteur			15	0	0	0	0	0	0	<a href="https://fr.wikipedia.org/wiki/Disponibilit%C3%A9">https://fr.wikipedia.org/wiki/Disponibilit%C3%A9</a>
116	30	Jeu de données & métadonnées	Toutes	Le mode d'accès à la donnée est un frein à l'usage (accès directs, outil d'accès long et complexe, droit d'accès limité)	La requête d'une donnée "temps réel" met plus de 40 secondes ; l'accès à la donnée nécessite un certificat de sécurité long à obtenir ; l'architecture du site ne permet pas à un robot de télécharger les actualisations des données	* Tester l'accès aux données * Tester la récupération des données via un outil automatisé (commande wget par exemple)			2	5	1	1	1	1		
27	31	Jeu de données & métadonnées	Toutes	La mesure de la qualité n'est pas documentée	Des contrôles qualité existent (omni ou avari) mais ils ne sont pas explicités si bien qu'on ne peut savoir si tel champ est fiable ou non	Lecture des métadonnées			3	15	0	0	0	0		
28	32	Jeu de données & métadonnées	Toutes	La qualité de la donnée n'est pas mesurable à travers des données formées	Il n'existe pas de méthode de contrôle permettant de dire si la syntaxe de ce champ est bonne	Demander au producteur			3	15	0	0	0	0		
29	33	Jeu de données & métadonnées	Toutes	La qualité de la données n'est pas mesurée	Aucune méthode de contrôle n'est mise en oeuvre pour mesurer la qualité des données	Demander au producteur ou à l'éditeur			3	15	0	0	0	0		
30	34	Jeu de données & métadonnées	Toutes	Une entité possède plusieurs identifiants					3	15	0	0	0	0	0	Exemple que me prend Simon sur les Assu qui ont à la fois un numéro d'asso et un code SIREN.
32	35	Jeu de données & métadonnées	Toutes	Le nom ou titre du jeu de données est vague, ambigu ou trop complexe : titre de la notice éditoriale, nom donné dans les métadonnées ou dans la documentation (pas le nom du fichier)	"Résultat des élections" : lesquelles ? où ? quand ? "Résultats des élections à Montréal" : il existe 6 communes appelées Montréal dans le monde	Lecture des métadonnées, de la documentation et/ou de la fiche de présentation			1	1	0	0	0	0		
33	36	Jeu de données & métadonnées	Date	Manque de métadonnées : fourchette temporelle non explicitée	Des dates figurent dans le jeu mais aucune métadonnée ne peut confirmer la fourchette attendue de ces dates. Exemple : Trésorerie du 01/02/2010 au 24/11/2016.	Lecture des métadonnées			2	1	0	0	0	0	0	* Produire les métadonnées * Faire renseigner ou valider les métadonnées par le producteur
34	37	Jeu de données & métadonnées	Date	Manque de métadonnées : zone spatiale non explicitée	Des coordonnées figurent dans le jeu mais aucune métadonnée ne peut confirmer la zone d'appartenance attendue pour ces points.	Lecture des métadonnées			2	1	0	0	0	0	0	* Produire les métadonnées * Faire renseigner ou valider les métadonnées par le producteur
35	38	Jeu de données & métadonnées	Nombre	Manque de métadonnées : fourchette non spécifiée	On peut attendre d'un nombre qu'il soit compris entre une valeur minimum et une valeur maximum ; par exemple l'âge d'une personne devrait toujours être entre 0 et 130 voir 18 et 70 selon les cas.	Lecture des métadonnées			2	1	0	0	0	0	0	* Produire les métadonnées * Faire renseigner ou valider les métadonnées par le producteur
36	39	Jeu de données & métadonnées	Boolean	Manque de métadonnées : fait que le champ soit un boolean n'est pas spécifique		Lecture des métadonnées			2	1	0	0	0	0	0	* Produire les métadonnées * Faire renseigner ou valider les métadonnées par le producteur
37	40	Jeu de données & métadonnées	Boolean	Manque de métadonnées : le format du boolean n'est pas spécifique	On ne sait pas à quelles valeurs s'attendre : "vrai", "faux" ou "oui", "non"	Lecture des métadonnées			2	1	0	0	0	0	0	* Produire les métadonnées * Faire renseigner ou valider les métadonnées par le producteur
38	41	Jeu de données & métadonnées	Toutes	Manque de métadonnées : processus et contexte de production non explicités	On ne sait pas si une mesure vient d'un capteur ou d'une mesure manuelle	Lecture des métadonnées		oui	1	1	0	0	0	0	0	* Produire les métadonnées * Faire renseigner ou valider les métadonnées par le producteur
39	42	Jeu de données & métadonnées	Toutes	Manque de métadonnées : la fraîcheur des données n'est pas explicitée "le délai entre le réel et la mise en base de la donnée" "le délai entre le réel et la publication de la donnée"	Il n'est pas dit si telle information sur une grosseuse va mettre plus de neuf avant d'arriver au réalisateur	Lecture des métadonnées			1	1	0	0	0	0	0	* Produire les métadonnées * Faire renseigner ou valider les métadonnées par le producteur
40	43	Jeu de données & métadonnées	Toutes	Manque de métadonnées : la langue des textes n'est pas spécifiée		Lecture des métadonnées			1	1	0	0	0	0	0	* Produire les métadonnées * Faire renseigner ou valider les métadonnées par le producteur
41	44	Jeu de données & métadonnées	Date	Métadonnées imprécises : le format de date n'est pas spécifique	Format américain ? anglais ? européen ? etc.	Lecture des métadonnées			1	1	0	0	0	0	0	* Produire les métadonnées * Faire renseigner ou valider les métadonnées par le producteur
42	45	Jeu de données & métadonnées	Nombre	Métadonnées imprécises : unités non spécifiées	On ne dit pas si colonne "hauteur" est en cm ou dm	Lecture des métadonnées		oui	1	1	0	0	0	0	0	* Produire les métadonnées * Faire renseigner ou valider les métadonnées par le producteur
43	46	Jeu de données & métadonnées	Coordonnées	Métadonnées imprécises : système de coordonnées non spécifié	La documentation n'indique pas si les coordonnées sont en WGS 84, Lambert ou un autre système	Lecture des métadonnées			1	1	0	0	0	0	0	* Produire les métadonnées * Faire renseigner ou valider les métadonnées par le producteur
44	47	Jeu de données & métadonnées	Toutes	Métadonnées imprécises : nom de colonnes ambigu	"Emploicement" ne dit rien sur la donnée attendue : une adresse ? "en haut" ? "devant" ? etc.	Lecture des métadonnées		oui	1	1	0	0	0	0	0	* Produire les métadonnées * Faire renseigner ou valider les métadonnées par le producteur
45	48	Jeu de données & métadonnées	Toutes	Métadonnées fautes		Lecture des métadonnées			1	1	0	0	0	0	0	
46	49	Jeu de données & métadonnées	Toutes	La taille maximale d'un champ dépasse celle qui est spécifiée dans la documentation	La colonne "age" spécifie une longueur de 3 caractères maximum et certaines valeurs sont des 4 caractères ou plus	Lecture des métadonnées et des données	csvgrep -c colonne_x + '[25,]' file.csv   csvlook --query 'SELECT MAX (LENGTH(mois)) FROM file' file.csv		3	3	1	1	1	1		
47	50	Jeu de données & métadonnées	Toutes	L'ordre des colonnes ne correspond pas à l'ordre donné dans la documentation	La documentation donne Prénom,Nom,Age,Profession alors que le jeu se présente sous la forme Nom,Prénom,Age,Profession	Lecture des métadonnées et des données			3	3	0	0	0	0	5	
48	51	Jeu de données & métadonnées	Coordonnées	Les coordonnées ne sont pas au format WGS 84	Les coordonnées sont au format Lambert II nécessitant une conversion des points pour des usages mobile liés à des GPS grand public	Lecture des métadonnées et des données			3	3	0	0	0	0		
49	52	Jeu de données & métadonnées	Chaîne alpha	Les codes pays ne sont pas au format ISO 3166	L'Allemagne est notée "ALL" alors qu'il existe un code ISO employé internationalement	Lecture des métadonnées et des données			2	3	1	1	1	1		
50	53	Jeu de données & métadonnées	Chaîne alpha	Les codes de langues ne sont pas au format ISO 639	Le français est noté "fr" ou "français" en lieu et place de "fr"	Lecture des métadonnées et des données			2	3	1	1	1	1		
51	54	Jeu de données & métadonnées	Date	La date n'est pas au format ISO 8601	La date est notée "01/01/2016"	Lecture des métadonnées et des données			2	3	1	1	1	1		
52	55	Jeu de données & métadonnées	Chaîne alpha	Les normales ne sont pas au format ISO 4217	Le franc suisse est noté FS	Lecture des métadonnées et des données			3	3	1	1	1	1		
53	56	Syntaxe	Chaîne alpha	Erreurs syntaxiques : espaces au début ou à la fin du champ	"Pieme" au lieu de "Pieme"	Rechercher les cas à partir d'une regexp	perl -f //; alias 'print' 'f. 1 print if (m{^\s ^\s\$}); cat 'file.csv'   csvlook		3	1	0	0	0	0	0	<a href="https://fr.wikipedia.org/wiki/ISO_4217">https://fr.wikipedia.org/wiki/ISO_4217</a>
54	57	Syntaxe	Chaîne alpha	Erreurs syntaxiques : bug syntaxique dans les strates du SI : le cas de l'apostrophe	"N'Diaye" à la place de N'Diaye	Rechercher "y"			1	1	1	1	1	1	1	Remplacer "y" par "
55	58	Syntaxe	Nombre	Erreurs syntaxiques : syntaxe des numéros ou nombres en trois genres	"456" au lieu de "45"; "1,000,000" au lieu de "1000000"				30	0	0	0	0	0	0	

ID	N° d'ordre	Domaine du problème	Type de données concernées	Type de problème	Exemple	Détecter le problème	Ligne de commande pour le contrôle	Réduit considérablement l'usage des données	Niveau d'importance du problème	Temps d'analyse en min	Analyse semi-automatisable	Solution	Temps de traitement	Remarques
56	59	Syntaxe	Chaîne alpha	Erreurs syntaxiques : codes (code INSEE, code postal, SIRET, SIREN, n° de Sécu, ISBN, ISSN, IBAN, BIC, code ROM, indicatif du pays, code APE, code NAF, etc.)	7100 au lieu de 07100 pour un code postal	Bâtir la regex relative au code attendu et tester.	cvgrep -c colonne_x + "^(dd/dd/ddd ) file.csv   cvsort   cvslook # code postal				5	1		
57	60	Syntaxe	Chaîne alpha	Erreurs syntaxiques : sigles et abréviations	"SMCF", "S. N. C. F.", "S. C. N. F." ? "Bou" ou "Bil" ?	Classer la colonne date par ordre alphabétique permet de rapidement voir les problèmes de syntaxe					5			
58	61	Syntaxe	Booléen	Erreurs syntaxiques : booléen	"Y" au lieu de "1" selon la spécification du booléen	Regarder WTF.csv sur les colonnes de booléens	cvstat -c col_booléen -freq file.csv				1	1		
59	62	Syntaxe	Chaîne alpha	Erreurs syntaxiques : email, url	laurent.dupont@wanadoo.fr						1	1		
60	63	Syntaxe	Date	Erreur syntaxique sur la date	2016/09/30 au lieu de 2016-09-30 attendu	Classer la colonne date par ordre alphabétique permet de rapidement voir les problèmes de syntaxe	cvlook -c colonne_x file.csv   cvsort   cvslook				1	1		
61	64	Syntaxe	Chaîne alpha	Incohérences syntaxiques : syntaxe des noms propres	"de La Tour" ou "La Tour (de) ?"						5	1		
62	65	Syntaxe	Nombre	Incohérences syntaxiques : homogénéité de la syntaxe des numéros ou nombres en tous genres	Dans le même fichier nous avons pour des chiffres parfois "1000,00" et parfois "100.000,00"						5	1		
63	66	Syntaxe	Chaîne alpha	Incohérences syntaxiques : usage du pluriel ou du singulier		Rechercher les pluriels à l'aide de motifs d'expressions régulières (?)					5	1		
64	67	Sémantique	Chaîne alpha	Plusieurs termes sont utilisés pour un même sens	Parfois on lit "Dashed", parfois "tas" et parfois "E1" ou bien "agent" ou "commercial", etc.	Trier le champ concerné par ordre alphabétique et regarder les valeurs (?)					5	1		<a href="https://fr.wikipedia.org/wiki/R%C3%A9partition_de_fonction">https://fr.wikipedia.org/wiki/R%C3%A9partition_de_fonction</a>
65	68	Sémantique	Chaîne alpha	Certains termes sont mal régionalisés ou traduits dans la langue attendue	Dans un fichier ou tout est en français, si l'on a "Grande-Bretagne" on devrait avoir "Etats-Unis" et pas "USA" qui est un terme anglais	Rechercher les chaînes et lancer le correcteur d'orthographe dans la langue désirée (?)					5	1		<a href="https://fr.wikipedia.org/wiki/R%C3%A9partition_de_fonction">https://fr.wikipedia.org/wiki/R%C3%A9partition_de_fonction</a>
66	69	Sémantique	Chaîne alpha	Certains termes, valeurs utilisées sont vieilles, inusitées, cryptiques ou incompréhensibles		Trier le champ concerné par ordre alphabétique et regarder les valeurs (?)					3	0		
67	70	Sémantique	Chaîne alpha	Les abréviations ou sigles ne sont pas explicités	Wikipédia fournit des listes de très nombreux sigles : https://fr.wikipedia.org/wiki/Sigle	Trier le champ concerné par ordre alphabétique et regarder les valeurs (?)					3	0		
68	71	Sémantique	Chaîne alpha/Nombre	La valeur nulle est remplacée par une autre chaîne : zero ou "0" ou "null" ou "1970-00-00" ou "0'00'00.0'N+0'00'00.0'E"	0'00'00.0'N+0'00'00.0'E est un problème car ce point existe mais il est placé en plein Atlantique	Trier le champ concerné par ordre alphabétique et regarder les valeurs (?)					3	1		
69	72	Sémantique	Toutes	Inversion dans un couple de données	"Dupont Jean" au lieu de "Jean Dupont"	Repérer les couples de données et classer les colonnes par ordre alphabétique pour repérer une éventuelle inversion (?)	TODD : si une chaîne de la colonne_x est présent 3 fois dans la colonne_y et inversement alors il y a suspicion d'inversion (?)				3	0		Le producteur peut avoir saisi Prénom Nom en étant persuadé de cet ordre. Ce problème survient également pour des Prénom-Noms d'origine culturelle différente (les chinois utilisent Nom-Prénom)
70	73	Sémantique	Chaîne alpha	L'absence de lettres accentuées peut poser des problèmes de sens	"JUPE TUE LA FRANCE GAGNE"	Rechercher des colonnes alpha qui ne possèdent pas d'accent	egrep [àâéèëé] file.csv   wc -l   uniq				1	1		les accents sont significatifs en Français
71	74	Sémantique	Chaîne alpha/Nombre	Erreur sémantique manifeste	Utilisation de "M" en lieu et place de "H" pour signifier un homme ; 09 pour le département en lieu et place du nom "Rhodé"	Rechercher toutes les valeurs d'une colonne, les dédoubler et analyser les résultats	cvlook -c colonne1 file.csv   sort   uniq				5	0		
72	75	Sémantique	Coordonnées	Erreur de système de coordonnées	Coordonnées en Lambert II au lieu de WGS 84 spécifiée dans les métadonnées						5			
73	76	Sémantique	Coordonnées	Les coordonnées géographiques sont données en degrés, minutes, secondes et non en degrés décimaux, ce qui complique leur manipulation	23°56'33" ou bien 23°56'33"E en lieu et place de la forme décimale 23.9756	Regarder toutes les colonnes représentées des coordonnées					3	1		
74	77	Sémantique	Date	Le format de la date est celui d'un autre pays ou d'une autre culture	09/08/2016 au lieu de 08/09/2013 pour le 8 septembre 2016 (la syntaxe est correcte mais le sens est incorrect)	Rechercher toutes les valeurs d'une colonne, les dédoubler, les trier et analyser les résultats	cvgrep -c colonne1 + "^(d D)/3(d D) file.csv				3	1		
75	78	Sémantique	Chaîne alpha/booléen	Liste de réponses fermées mal conçue : réponse "rien" n'est pas inclusivement dans "sans réponse" ou autres pourraient convenir	"Vous êtes plutôt d'accord avec cette assertion : vrai/faux" / "Ne se prononce pas" devrait avoir une réponse pertinente	Détecter les colonnes ne possédant que deux valeurs et se poser la question	perl -F/, -alike 'print if 1..1,print if m/Au/Divers/;' file.csv   cvsort   cvslook				3	0		
76	79	Sémantique	Chaîne alpha	Liste de réponses fermées mal conçue : présence de la réponse "Autre" ou "Divers" très fréquente	"Quel est votre ville favorite : Marseille, Paris, Autre"	Rechercher les chaînes "Autre" et "Divers"					1	1		Dans certains cas, les réponses "Autre" ou "Divers" peut être parfaitement justifiées.
77	80	Morpho-syntaxique	Toutes	Exprimer une donnée à travers un code difficile à lire	Mise en forme pour exprimer une donnée : couleur, gras, etc.	Parcourir visuellement l'ensemble du fichier					5	0		
78	81	Morpho-syntaxique	Toutes	autres ?	Certains fichiers possèdent des cellules fusionnées des données sont ajoutées sous forme de commentaires ; etc.	Parcourir visuellement l'ensemble du fichier					5	0		
79	82	Peritence	Chaîne alpha/Nombre	Aberation	"197 ans (pour l'âge d'une personne)" "Général de Gaulle comme personne participant à un sondage"	* Un classement des champs par ordre alphabétique permet de localiser des grands aberrants. * Tester que les données vérifient la loi de Benford. * WTF.csv					30			
80	83	Peritence	Chaîne alpha/Nombre	Double très raisonnable, valeurs inexplicables	20 participants de plus de 110 ans	* Rechercher les valeurs extrêmes de chaque colonne et s'interroger : * WTF.csv ?		oui			5	1		
81	84	Peritence	Chaîne alpha	Certaines valeurs sont suspectes : 0000 ou xxxxxxxxxx (à compléter)	-	Rechercher des chaînes "999" et "123"	perl -F/, -alike 'print if 1..1,print if m/(000 xxx ) file.csv   cvsort   cvslook				3	1		
82	85	Peritence	Nombre	Certaines valeurs sont suspectes : suites de chiffres comme 9999 ou 12345	Des suites de 999 ; nombreuses valeurs "12345" (détaillez)	Recherche des chaînes "999" et "123"	perl -F/, -alike 'print if 1..1,print if m/999 12345 000/ file.csv   cvsort   cvslook				3	1		
83	86	Peritence	Date	Certaines valeurs sont suspectes : il existe des dates en 1900, 1904, 1909, 1970	-	Recherche des chaînes "1900", "1904", "1909", "1970"	perl -F/, -alike 'print if 1..1,print if m/1900 1904 1909 1970/ file.csv   cvsort   cvslook				3	1		
84	87	Peritence	Coordonnées	Certaines valeurs sont suspectes : il existe des coordonnées comme 0'00'00.0'N+0'00'00.0'E	0'00'00.0'N+0'00'00.0'E est une valeur suspecte car c'est un point en plein milieu de l'Atlantique	POI : placer tous les POI sur une carte pour voir si certains sont hors périmètre	cvgrep -c colonne_x + "0'00'00.0' file.csv   cvsort   cvslook				3	1		
85	88	Peritence	Toutes	La source n'est pas crédible (incompétent, juge et partie, etc.)	15000 manifestants selon les organisateurs	Questionner la crédibilité de la source - est-elle compétente pour collecter ces données ? A-t-elle un intérêt partisan à faire parler les données dans une certaine direction ?		oui			3	0		Autres exemples : chiffres du chômage (?), chiffres "hors du champ" par les politiciens, résultats d'audiences ou francs communiqués par l'acteur concerné par ces chiffres, ...
86	89	Peritence	Toutes	Les données ont été hackées ou détournées	La source est crédible mais certains producteurs indirects ont pu agir pour que certaines données soient sur-représentées (sondage, etc.)	* La sur-représentation d'un profil ou des valeurs suspectes doivent conduire à s'interroger * Tester que les données vérifient la loi de Benford.					30	0		Exemple de l'affaire Clearstream. Affaire Hashley Madison (?). Voir le type de problème "Le process d'acquisition n'est pas connu"
87	90	Réglementation	Toutes	Identification explicite de personnes sans déclaration CNIL	Prénom Nom ou numéro de tél.	Détecter des pré-noms sur la base d'un dictionnaire est-il un bon indicateur ?		oui						
88	91	Réglementation	Toutes	Identification possible de personnes	Date et lieu de naissance	Parcourir le fichier dans son ensemble suffit-il ?		oui						
89	92	Réglementation	Chaînes alpha	Il existe des jugements de valeurs à propos d'individus	"Client chard", etc.	Rechercher des mots "interdits" comme "chard", "stupide", "idiot", "non-sensées", "ennui", etc.					5	1		
90	93	Réglementation	Chaînes alpha	Il existe des données de santé non anonymisées alors que les personnels qui les consultent n'y sont pas habilités	"Ne peut pas nous recevoir le mercredi matin car elle fait sa dialyse"	Rechercher des mots "interdits" comme "dialyse", "cancer", etc.					5	1		
91	94	Réglementation	Toutes	Données d'origine ethnique ou relative à la religion des personnes	"Ne répond pas au téléphone le samedi (shabbat) habilités"	Rechercher des mots qui peuvent être des indicateurs comme "caucasien", "ghébreu", "juif", "musulman", etc.		oui			3	0		
92	95	Réglementation	Toutes	Données relatives aux opinions politiques, philosophiques ou à l'appartenance syndicale	"Léti au parti pirate"	Rechercher des mots qui peuvent être des indicateurs comme le nom de partis politiques, de courants de pensée, etc.		oui			3	0		
93	96	Réglementation	Toutes	Données relatives à la vie sexuelle ou au meurs	"Ménage à 3"	Rechercher des mots qui peuvent être des indicateurs comme "sex", "homo", etc.		oui			3	0		
94	97	Réglementation	Toutes	Données tierces soumises à licence d'usage	Le fichier publié en Open Data utilise le géocodage de l'API de Grange	Lire les métadonnées ; en cas de doute, demander explicitement au producteur.		oui			5			
95	98	Réglementation	Chaîne alpha	Données relevant de la propriété littéraire et artistique sans autorisation d'usage : description textuelle	La description littéraire d'une chose est soumise à des droits	Rechercher les chaînes de plus de 100 (?) caractères et évaluer si la rédaction dépasse un simple caractère factuel (?)					3			
96	99	Réglementation	binaires	Données relevant de la propriété littéraire et artistique sans autorisation d'usage : images ou fichiers multimédia	Les images d'une base de données sont soumises à des droits	Le jeu de données comprend-il des images ? Le droit d'usage de ces images est-il explicite ? Ce droit pose-t-il problème pour des usages ultérieurs ?					3			
97	100	Réglementation	Toutes	Données sensibles du point de vue de la sécurité des biens et des personnes	Plan d'une base militaire	Parcourir le fichier dans son ensemble suffit-il ?		oui			3	0		
98	101	Réglementation	Toutes	Données sensibles du point de vue de l'éthique	Localisation de minéraux rares ou de zones d'habitat d'espèces protégées	?		oui			3	0		
99	102	Réglementation	Chaîne alpha/Nombre	divers : exemple le capital social d'une entreprise s'incrémentent annuellement et non la valeur inférieure								0		Les chiffres communiqués aux impôts comme la TVA sont arrondis : la raison sociale d'une entreprise : les pré-noms-noms dans un contexte d'identification officielle des tarifs (?), les cours des monnaies (?).
100	103	Manque	Nombre	La donnée est le résultat d'un calcul dont on n'a pas les données de départ	Le jeu de données contient un pourcentage, un rapport, une densité, etc.	Parcourir les métadonnées pour évaluer chaque champ					3	0		
101	104	Manque	Toutes	Les trous : manque des "enregistrements" : des données dont vous connaissez l'existence sont manquantes	Il manque 10 communes dans la liste des maires du département de la Savoie	Rechercher toutes les valeurs d'une colonne, les dédoubler, les trier et analyser les résultats					3	0		
102	105	Manque	Toutes	Les trous : manque des "enregistrements" : le tableau possède 65536 lignes	-	Ouvrir le jeu de données et regarder s'il contient 65536 lignes	cvstat -count file.csv cvsort file.csv				1	1		
103	106	Manque	Toutes	Les trous : les données d'un champs sont tronquées	"10, av. Général de Gaulle" est un exemple de champ tronqué à 25 caractères	* Pour une longueur de champ donnée, quel pourcentage d'enregistrements remplissent ce champ complètement ?	cvlook -c colonne_x file.csv   cgrep -c colonne_x + "[25]"   cvsort   c-adr   uniq   cvslook   head -20				3	1		
104	107	Manque	Toutes	Valeurs vides dans certains champs	-	Ouvrir le jeu de données et regarder s'il contient des valeurs vides	cvstat -null file.csv				1	1		
105	108	Manque	Toutes	La granularité n'est pas suffisante	On a des pays, là où il serait intéressant d'avoir des régions ; on a des mètres là où certains usages nécessiteraient des cm						5	0		
106	109	Manque	Date	Le fuseau horaire n'est pas précisé dans un contexte de données réparties sur des fuseaux horaires différents	Pour une heure locale donnée, l'absence du fuseau horaire oblige le développeur à tenter de calculer l'heure GMT pour comparer des durées						5	0		
107	110	Manque	Toutes	L'insuffisance en matière de fréquence	L'état de deux tricolores classiques est donné tous les jours à minute	Lire les métadonnées, examiner les données (un lit descendant des colonnes date peut aider) puis comparer avec des usages possibles					5	0		
108	111	Manque	Toutes	L'insuffisance en matière de maillage	La pollution dans Paris est mesurée avec un seul capteur						5	0		
109	112	Manque	Toutes	L'insuffisance en matière de fraîcheur	Les chiffres du recensement de cette espèce date de 1976	Lire les métadonnées, examiner les données (un lit descendant des colonnes date peut aider) puis comparer avec des usages possibles					5	0		
110	113	Surabondance	Toutes	Doublons	* Jean Martin 21/12/1956Lunville.moscontacts * Jean Martin 21/12/1956Lunville.lesestup.	Repérer les champs permettant d'identifier une chose de manière non ambiguë et analyser la présence de doublons à travers un tri alphabétique d'un des champs.	* Pour détecter les doublons stricts : - Pour détecter les doublons selon deux colonnes : cvsort -c colonne1.colonne2 file.xxx   sort   uniq -d				5	1		La notion de doublons inclut les valeurs en triple, quadruple, etc. Il existe d'une part les doublons stricts, qui peuvent être basés facilement ; d'autre part les doublons qui nécessitent une fusion des valeurs (plus difficiles à détecter et baser).

ID	N° d'ordre	Domaine du problème	Types de données concernées	Type de problème	Exemple	Détecter le problème	Ligne de commande pour le contrôle	Réduit considérablement l'usage des données	Niveau d'importance du problème	Temps d'analyse en min	Analyse semi-automatisable	Solution	Temps de traitement	Remarques
111	114	Surabondance	Toutes	Valeur obtenue par calcul sur la base de deux données déjà présentes	Le milieu de deux coordonnées géographiques est indiqué en plus des deux coordonnées.	Lire les métadonnées et vérifier le calcul sur un petit échantillon				3	0			
112	115	Surabondance	Toutes	Valeur renseignée alors qu'elle est calculable à partir de données prises ailleurs	Hors Paris et Marseille, pourquoi avoir un code postal quand la ville permet de le déduire ?					3	0			
113	116	Surabondance	Toutes	Surabondance de données en matière de : précision, fréquence, maillage ou fraîcheur	La localisation au mm d'une porte d'entrée ; ajout du fuseau horaire dans contexte où toutes les données sont en France métropolitaine ; etc.	* La précision est-elle de plus de 4 chiffres après la virgule ?	csvgrep -c colonne_x « "dd(,.)\d{4}\d{4}" file.csv   csvsort   csvlook			3	1			
114	117	Surabondance	Coordonnées	Pour des coordonnées géographiques, une précision supérieure à 8 unités après la virgule est inutile (précision de l'ordre du mm) ; 5 chiffres après la virgule donnent déjà une précision de l'ordre du mètre	23.73825619 positionne un objet à environ 1 mm	Regarder la longueur des champs représentant des coordonnées	csvgrep -c colonne_x « "dd(,.)\d{8}\d{4}\d{4}" file.csv   csvsort   csvlook			3	1	1 Arrondir à 7 chiffres après la virgule, voire à 6 chiffres.		
								25		425	39		26	

Sprint qualité	Sprint qualité basé sur la version 0.2 de la méthode "sprint qualité" : <a href="http://infolabs.io/sprint-qualite">http://infolabs.io/sprint-qualite</a>
<b>Titre initial du jeu de données</b>	Base de données accidents corporels de la circulation
URL de référence	<a href="https://www.data.gouv.fr/fr/datasets/base-de-donnees-accidents-corporels-de-la-circulation/">https://www.data.gouv.fr/fr/datasets/base-de-donnees-accidents-corporels-de-la-circulation/</a>
Date du sprint (aaaa-mm-jj)	10/11/2016
Auteurs	Animateurs : Armelle Gilliard <agilliard@fing.org>, Charles Nepote <charles.nepote@fing.org>
Binôme 1 Prénom-Nom <email>, Prénom2-Nom2 <email2>	Binôme 1 Romain Talès Christian Quest
Binôme 2 Prénom-Nom <email>, Prénom2-Nom2 <email2>	Binôme 2 Marie Heuze Vincent Bataille
etc.	Binôme 3 Sarah Labelle Fabien Antoine
Remarques	On étudie le fichier caracteristiques_2015.csv
URL de publication du résultat du sprint	





[1] \* très faible  
\* faible  
\* importante  
\* très importante

[2] \* négligeable  
\* faible  
\* important  
\* rend le jeu inutilisable

[3] Complète  
Partielle