

	Number of lines in the table : 273226	CHARACTERISTICS						
# Column	Name of the Column	Description	Variable's type	Percentage of missing values	Categorical / Quantitative	Distribution	Comments	Keep?
		What does this variable represent (from a business perspective ?)	<code>int64, float etc...</code> If "object", develop.	in %		<code>For categorical variables with less than 10 categories, list all categories.</code> <code>For quantitative variables, detail the distribution (basic descriptive statistics)</code>	Free text	Should we keep or drop the column?
1	Num_Acc	Unique accident identifier number	<code>integer</code>	0,00%	Unique Value			Keep
2	mois (rename to month)	Month of the accident	<code>integer</code>	0,00%	Categorical	<code>Integer from 1 to 12</code>		Keep
3	jour (rename to day)	Day of the accident	<code>integer</code>	0,00%	Categorical	<code>Integer from 1 to 31</code>		Keep
4	hrmn	Time (hour and minute) of the accident	<code>object</code>	0,00%	Datetime	<code>Datetime given in military time</code>		Keep
5			<code>integer</code>	0,00%	Categorical - 3 to 5 categories	<code>1 – Daylight</code> <code>2 – Dawn or dusk</code> <code>3 – Night without street lighting</code> <code>4 – Night with unlit street lighting</code> <code>5 – Night with lit street lighting</code>		Keep
	<i>lum (rename to light_cond)</i>	<i>Light conditions</i>						
6	agg (rename to loc_type)	Location type	<code>integer</code>	0,00%	Categorical - Binary	<code>1 – Outside built-up area</code> <code>2 – In built-up area</code>		Keep
7	<i>int (rename to intersc_type)</i>	<i>Intersection type</i>	<code>integer</code>	0,00%	Categorical - 6 to 10 categories	<code>1 – Not at intersection</code> <code>2 – Crossroads (X)</code> <code>3 – T-junction</code> <code>4 – Y-junction</code> <code>5 – Intersection with more than 4 branches</code> <code>6 – Roundabout</code> <code>7 – Square</code> <code>8 – Level crossing</code> <code>9 – Other intersection</code>		Keep
8	<i>atm (rename to weather_cond)</i>	<i>Weather conditions</i>	<code>integer</code>	0,01%	Categorical - 6 to 10 categories	<code>-1 – Not recorded</code> <code>1 – Normal</code> <code>2 – Light rain</code> <code>3 – Heavy rain</code> <code>4 – Snow or hail</code> <code>5 – Fog or smoke</code> <code>6 – Strong wind / storm</code> <code>7 – Dazzling conditions (sun glare, etc.)</code> <code>8 – Overcast</code> <code>9 – Other</code>		Keep
9	<i>col (rename to collision_type)</i>	<i>Type of collision</i>	<code>integer</code>	0,00%	Categorical - 6 to 10 categories	<code>-1 – Not recorded</code> <code>1 – Head-on (two vehicles)</code> <code>2 – Rear-end (two vehicles)</code> <code>3 – Side impact (two vehicles)</code> <code>4 – Chain collision (three or more vehicles)</code> <code>5 – Multiple collisions (three or more vehicles)</code> <code>6 – Other collision</code> <code>7 – No collision</code>		Keep
10	com	Commune – INSEE commune code	<code>object</code>	0,00%	Categorical			Keep
11	adr	Postal address (for accidents in built-up areas)	<code>object</code>	1,46%	Categorical			Drop
12	lat	Latitude	<code>object</code>	0,00%	Quantitative			Keep
13	long	Longitude	<code>object</code>	0,00%	Quantitative			Keep
14	dep	Department – INSEE department code	<code>object</code>	0,00%	Categorical		Could be a way to differentiate between mainland and overseas departments	Keep
15	year	Year of the accident	<code>integer</code>	0,00%	Categorical	<code>Integers from 2019 to 2023</code>		Keep

	Number of lines in the table : 467169	VEHICLES						
# Column	Name of the Column	Description	Variable's type	Percentage of missing values	Categorical / Quantitative	Distribution	Comments	Keep?
		What does this variable represent (from a business perspective ?) If "object", develop.	int64, float etc... If "object", develop.	in %		For categorical variables with less than 10 categories, list all categories. For quantitative variables, detail the distribution (basic descriptive statistics)	Free text	Should we keep or drop the column?
1	Num_Acc	Foreign key for Characteristics table	integer	0.00%	Categorical	-1 – Not recorded 0 – Unknown 1 – Increasing PK/PR/address number 2 – Decreasing PK/PR/address number 3 – No reference point		Keep
2	senc (rename to direc_travel)	Direction of travel	integer	0.00%	Categorical - 3 to 5 categories	00 – Undetermined 01 – Bicycle 02 – Moped <50 cm3 03 – Micrcar (body quadricycle, formerly motor tricycle) 07 – Passenger car (VL) 10 – Light utility vehicle (1.5T <= GVWR <= 3.5T) 13 – HGV (3.5T < GVWR <= 7.5T) 14 – HGV (> 7.5T) 15 – HGV + trailer 16 – Road tractor 17 – Road tractor + semi-trailer 20 – Special vehicle 21 – Agricultural tractor 30 – Scooter <50 cm3 31 – Motorcycle >50 cm3 and <=125 cm3 32 – Scooter >50 cm3 and <=125 cm3 33 – Motorcycle >125 cm3 34 – Scooter >125 cm3 35 – Light quad <=50 cm3 (non-body) 36 – Heavy quad >50 cm3 (non-body) 37 – Bus 38 – Coach 39 – Train 40 – Tramway 41 – 3RM <=50 cm3 42 – 3RM >50 cm3 and <=125 cm3 43 – 3RM >125 cm3 50 – Motorized personal mobility device 60 – Non-motorized personal mobility device 80 – Electrically assisted bicycle 99 – Other vehicle	Group categories.	Drop
3	catv (rename to vehicle_cat)	Vehicle category	integer	0.00%	Categorical - more than 10 categories	mean = 0.08; sd = 2.2; min = 0; max = 9; Q1, Q2 and Q3 = 0		Keep
4	occutc (rename to num_passenger_pub)	Number of passengers in public transport	float	99.16%	Quantitative	-1 – Not recorded 0 – None 1 – Parked vehicle 2 – Tree 3 – Metal guardrail 4 – Concrete barrier 5 – Other guardrail 6 – Building/wall/bridge pier 7 – Sign support/emergency phone post 8 – Pole 9 – Street furniture 10 – Parapet 11 – Island/refuge/high bollard 12 – Kerb 13 – Ditch/embankment/rock face 14 – Other fixed obstacle on carriageway 15 – Other fixed obstacle on pavement/shoulder 16 – Off-road departure without obstacle 17 – Culvert/aqueduct		Drop
5	obs (rename to fixed_obs_hit)	Fixed obstacle hit	integer	0.00%	Categorical - more than 10 categories	-1 – Not recorded 0 – None 1 – Pedestrian 2 – Vehicle 4 – Rail vehicle 5 – Domestic animal 6 – Wild animal 9 – Other		Keep
6	obsm (rename to mobile_obs_hit)	Mobile obstacle hit	integer	0.00%	Categorical - 6 to 10 categories			Keep

	Number of lines in the table : 467169	VEHICLES						
# Column	Name of the Column	Description	Variable's type	Percentage of missing values	Categorical / Quantitative	Distribution	Comments	Keep?
		What does this variable represent (from a business perspective ?)	<i>int64, float etc... If "object", develop.</i>	<i>in %</i>		<i>For categorical variables with less than 10 categories, list all categories. For quantitative variables, detail the distribution (basic descriptive statistics)</i>	<i>Free text</i>	<i>Should we keep or drop the column?</i>
7	<i>choc (rename to int_impact_point)</i>	<i>Initial impact point</i>	<i>integer</i>	<i>0,00%</i>	Categorical - 6 to 10 categories	<i>-1 – Not recorded 0 – None 1 – Front 2 – Front right 3 – Front left 4 – Rear 5 – Rear right 6 – Rear left 7 – Right side 8 – Left side 9 – Multiple impacts (rollover)</i>		<i>Keep</i>
8	<i>manv (rename to maneuvre)</i>	<i>Main maneuver before accident</i>	<i>integer</i>	<i>0,00%</i>	Categorical - more than 10 categories	<i>-1 – Not recorded 0 – Unknown 1 – No change of direction 2 – Same lane, same direction 3 – Between two lanes 4 – Reversing 5 – Wrong way 6 – Crossing central reservation 7 – Bus lane, same direction 8 – Bus lane, opposite direction 9 – Merging 10 – U-turn 11 – Lane change left 12 – Lane change right 13 – Drift left 14 – Drift right 15 – Left turn 16 – Right turn 17 – Overtaking left 18 – Overtaking right 19 – Crossing carriageway 20 – Parking maneuver 21 – Avoidance maneuver 22 – Opening door 23 – Stopped (not parked) 24 – Parked (with occupants) 25 – Driving on pavement 26 – Other maneuvers</i>		<i>Keep</i>
9	<i>num_veh</i>	<i>Vehicle identifier (alphanumeric)</i>	<i>object</i>	<i>0,00%</i>	Categorical		<i>Identification of the vehicle taken back for each user occupying this vehicle (including pedestrians who are attached to vehicles that hit them) - alphanumeric code - change to string</i>	<i>Keep</i>
10	<i>year</i>	<i>Year of accident</i>	<i>integer</i>	<i>0,00%</i>	Categorical	<i>Integers from 2019 to 2023</i>		<i>Keep</i>
11	<i>id_vehicle</i>	<i>Unique vehicle identifier (numeric)</i>	<i>object (integer)</i>	<i>0,00%</i>	Unique Value		<i>Together with num_veh used to merge to the characteristics file - change to integer</i>	<i>Keep</i>
12	<i>motor</i>	<i>Vehicle engine type</i>	<i>integer</i>	<i>0,00%</i>	Categorical - 6 to 10 categories	<i>-1 – Not recorded 0 – Unknown 1 – Fossil fuels 2 – Hybrid electric 3 – Electric 4 – Hydrogen 5 – Human powered 6 – Other</i>		<i>Keep</i>

Number of lines in the table : 619807																																													
# Column	Name of the Column	Description	Variable's type	Percentage of missing values	Categorical / Quantitative	Distribution	Comments	Keep?																																					
		What does this variable represent (from a business perspective ?) If "object", develop.	int64, float etc... If "object", develop.	in %		For categorical variables with less than 10 categories, list all categories. For quantitative variables, detail the distribution (basic descriptive statistics)	Free text	Should we keep or drop the column?																																					
0	Num_Acc (rename num_accident)	Accident identifier	int64	0,00%	Categorical - more than 10 categories	Foreign key - Characteristics table - used to merge	Keep																																						
1	id_vehicule (rename id_vehicle)	Unique numeric identifier for each vehicle involved in the accident	object	0,00%	Categorical - more than 10 categories	Foreign key - Vehicles table - used to merge	Change type to integer	Keep																																					
2	num_veh (rename num_vehicle)	Alphanumeric vehicle identifier used to link each user (driver, passenger, pedestrian) to the corresponding vehicle	object	0,00%	Categorical - more than 10 categories	Foreign key - Vehicle table - used to merge	Change type to string due to alphanumeric values	Keep																																					
3	place	Allows to locate the place occupied in the vehicle by the user at the time of the accident	int64	0,00%	Categorical - 6 to 10 categories	General codes (apply to all vehicles): -1 – Not specified NaN – Missing (no value recorded in dataset) Motorcycle / Sidecar 1 – Driver (rider) 2 – Passenger behind the rider (pillon) 3 – Sidecar passenger Car (Passenger car) 1 – Driver seat 2 – Front passenger seat 3 – Rear side seat (right) 4 – Rear side seat (left) 5 – Rear middle seat 6 – Front middle seat (bench, if present) 7 – Other side seat (3rd row) 8 – Other/unspecified passenger position 9 – Unknown seat Public transport (bus/coach, etc.) 1 – Driver 6 – Door/platform area 7 – Seated passenger 8 – Standing passenger (aisle/interior) 9 – Other/unspecified area Pedestrian 10 – Pedestrian (not applicable) - not found in dataset	1. Value 0 is in dataset, but no description what it means 2. Change type to int64 instead of float64	Drop		Transport en commun Moto / Side-car Voliture  <table border="1"><tr><td>4</td><td>7</td><td>7</td><td>7</td><td>7</td><td>7</td><td>1</td></tr><tr><td>5</td><td>8</td><td>8</td><td>8</td><td>8</td><td>8</td><td>6</td></tr><tr><td>5</td><td>8</td><td>8</td><td>8</td><td>8</td><td>8</td><td>6</td></tr><tr><td>5</td><td>8</td><td>8</td><td>8</td><td>8</td><td>8</td><td>6</td></tr><tr><td>3</td><td>9</td><td>9</td><td>9</td><td>9</td><td>9</td><td>2</td></tr></table>	4	7	7	7	7	7	1	5	8	8	8	8	8	6	5	8	8	8	8	8	6	5	8	8	8	8	8	6	3	9	9	9	9	9	2
4	7	7	7	7	7	1																																							
5	8	8	8	8	8	6																																							
5	8	8	8	8	8	6																																							
5	8	8	8	8	8	6																																							
3	9	9	9	9	9	2																																							
4	catu (rename category_user)	User category	int64	0,00%	Categorical - 3 to 5 categories	1 – Driver 2 – Passenger 3 – Pedestrian	Keep																																						
5	grav (rename injury_severity)	Severity of the accident. The injured users are classified into three categories of victims plus the uninjured	int64	0,00%	Categorical - 3 to 5 categories	-1 – Not specified 0 – Uninjured 2 – Killed 3 – Hospitalized wounded 4 – Light injury	Main indicator the measure for the project; -1 probably means not specified? There are 419 users with no grav value (0.01% who do not have effect when excluded)	Keep	count percentage grav 1 200753 42.06 4 247097 39.86 3 93257 15.36 2 16445 2.65 -1 419 0.07																																				
6	sexe (rename to gender)	Sex of the user	int64	0,00%	Categorical - Binary	-1 – Not specified 1 – Male 2 – Female	1. Decide if we want to leave the missing values out or replace it with mode or the distribution between male/female	Keep	count percentage sex 1 417865 67.40 2 193849 31.27 -1 607 1.33																																				
7	an_nais (rename year_of_birth)	Year of birth of the user	float64	1,38%	Categorical - more than 10 categories	Ranging from 1900.0 to 2023.0	1. Change the type to int64 instead of float64 2. Change or add Age (year of birth minus year) 3. Decide if we want to leave the missing values out or replace it with mean	Keep																																					
8	trajet (rename to trip_purpose)	Reason for traveling at the time of the accident	int64	0,00%	Categorical - 6 to 10 categories	-1 – Not specified 0 – ?? 1 – Home - work 2 – Home - school 3 – Shopping - Shopping 4 – Professional use 5 – Promenade - leisure 9 – Other	No check what 0 means due to dropping the column	Keep																																					

Number of lines in the table : 619807											
# Column	Name of the Column	Description	Variable's type	Percentage of missing values	Categorical / Quantitative	Distribution	Comments	Keep?			
		What does this variable represent (from a business perspective?) If "object", develop.	int64, float etc... If "object", develop.	in %		For categorical variables with less than 10 categories, list all categories. For quantitative variables, detail the distribution (basic descriptive statistics)	Free text	Should we keep or drop the column?			
9	secu1 (rename safety_equipment1)	If safety equipment was used and which one	int64	0.00%	Categorical - more than 10 categories	-1 – Not specified 0 – Do not use 1 – Seatbelt 2 – Helmet 3 – Child restraint system 4 – Reflective vest 5 – Gloves (2-wheel/3-wheel vehicles) 6 – Gloves (2-wheel/3-wheel vehicles) 7 – Gloves + Airbag (2-wheel/3-wheel vehicles) 8 – Not determinable 9 – Other	Since 2010, the variables secu1, secu2, and secu3 allow recording up to three types of safety equipment per user (e.g., helmet, gloves, airbag), replacing the previous single-variable (kind of equipment + if used) system. Here we have not only not specified (kinda like missing value) but also not determinable (it was not possible to determine the safety equipment with 1 % not specified and 12 % not determinable, we have to decide how to work with this variable	Keep	count percentage secu1 1 365320 58.93 2 113314 18.28 8 74985 12.03 0 50120 8.57 -1 7883 1.27 3 3884 0.63 9 658 0.11 6 534 0.09 4 401 0.06 5 246 0.04 7 16 0.00		
10	secu2 (rename safety_equipment2)	If second safety equipment was used and which one	int64	0.00%	Categorical - more than 10 categories	see secu1	see secu1	Keep			
11	secu3 (rename safety_equipment3)	If third safety equipment was used and which one	int64	0.00%	Categorical - more than 10 categories	see secu1	see secu1	Keep			
12	locp (rename location_pedestrian)	Location of the pedestrian	int64	0.00%	Categorical - 6 to 10 categories	-1 – Not specified 0 – Not applicable (no pedestrian) On the roadway: 1 – More than 50 m from a pedestrian crossing 2 – Within 50 m of a pedestrian crossing On pedestrian crossing: 3 – Without traffic light signal 4 – With traffic light signal Other: 5 – On the sidewalk 6 – On the roadside shoulder 7 – On a refuge island or emergency lane (BAU) 8 – On a service road 9 – Unknown	with 45 % not specified, probably drop the column	Drop	count percentage locp 0 291911 47.08 -1 280083 45.16 3 15048 2.43 2 10413 1.68 4 7207 1.16 1 6668 1.12 5 3770 0.61 9 2375 0.38 6 1445 0.23 8 698 0.11 7 53 0.01		
13	actp (rename action_pedestrian)	Action of the pedestrian	object	0.00%	Categorical - 6 to 10 categories	-1 – Not specified 0 – Not applicable (no pedestrian) While moving: 1 – In the direction of the oncoming vehicle 2 – In the opposite direction of the vehicle Other: 3 – Crossing 4 – Hidden (obstructed view) 5 – Playing / running 6 – With an animal 9 – Other A – Getting on / off a vehicle B – Unknown	with 40 % not specified, would probably drop the column	Drop	count percentage actp 0 322934 52.09 -1 248919 40.15 3 3533 5.70 9 3228 0.52 1 3133 0.51 5 1878 0.30 B 1608 0.26 2 1479 0.24 4 621 0.10 A 540 0.09 6 135 0.02 7 84 0.01 8 61 0.01		
14	etatp (rename company_pedestrian)	Specify whether the injured pedestrian was alone or not	int64	0.00%	Categorical - 3 to 5 categories	-1 – Not specified 1 – Alone 2 – Accompanied 3 – In a group	with 92 % not specified, probably drop the column	Drop	count percentage etatp -1 571936 92.25		
15	id_usager (rename id_user)	Unique user identifier (numeric)	string	38.43%	Categorical - more than 10 categories	First time added in year 2021, therefore missing values for the years, 2019-2020	Drop due to high Percentage of missing values	Drop	count percentage id_usager 1 36457 5.88 2 9856 1.56 3 1922 0.31		
16	year	year of the accident	int64	0.00%	Categorical - 3 to 5 categories	Ranging from 2019 to 2023		Keep			

# Column	Name of the Column	PLACES	Description	Variable's type	Percentage of missing values	Categorical / Quantitative	Distribution	Comments	Keep?
			What does the variable represent (from business perspective?)	int64 float64	0.00%		For categorical variables with less than 10 categories, order of categories matters.	For quantitative variables, order of distribution matters.	For categorical variables with less than 10 categories, order of categories matters.
			"object" develop	in %			For categorical variables with less than 10 categories, order of categories matters.	For quantitative variables, order of distribution matters.	For categorical variables with less than 10 categories, order of categories matters.
0	Num_Acc	Accident ID		int64	0.00%				Keep
1	car (rename: road_category)	Category of road		float64	0.00%	Categorical categories	6 to 10		Drop
2	vie (rename: road_id)	Road number	object	12.3%	Quantitative			Roads in France can contain several numbers for it is an 'Open Street'.	Drop
3	v1 (rename: num_index)	Numeric index of route	float64	3.71%	Quantitative			e.g. Road 2 in Road 3	Drop
4	v2 (rename: left_index)	Letter index of route	object	91.90%	Quantitative			Only letters as routes are not allowed ('Open Street'). e.g. Road 05A would be allowed but not 05B and left_index = 'A'	Drop
5	crc (rename: traff_msp)	Traffic regime	float64	0.00%	Categorical categories	3 to 5			Keep
6	rdv (rename: traff_lane_nb)	Number of traffic lane	float64	2.22%	Quantitative	X			Keep
7	veip (rename: reserved_lane)	Reserved lane	float64	0.00%	Categorical categories	3 to 5		 	Keep
8	prof (rename: road_dif)	Road slope	float64	0.00%	Categorical categories	3 to 5		-1 - Not specified 1 - One way 2 - Bidirectional 3 - With vertical separation channels 4 - With variable separation channels	Keep
9	pr (rename: ref_point)	object	0.00%	Quantitative				48 Point kilométrique	Drop
10	prt (rename: ref_point_dms)	Distance from reference point (metres)	object	0.00%	Quantitative			48 Point kilométrique	Drop
11	plan (rename: road_align)	Road alignment	float64	0.00%	Categorical categories	3 to 5		-1 - Not specified 1 - Straight path 2 - Curved on the left 3 - Curved right 4 - In "Z"	Keep
12	lrgc	Width of central median (metres)	float64	99.80%	Quantitative			This describes the width of the central median (the space between the two directions of a divided road).	Drop
13	lrmcl (rename: road_width)	Roadway width (metres)	float64	20.20%	Quantitative			Width of the roadway DOES NOT TAKE INTO ACCOUNT THE EMERGENCY LANE PREVIOUSLY SHOWN	Drop
14	surf	Road surface condition	float64	0.00%	Categorical categories	6 to 10		-1 - Not specified 1 - normal 2 - wet 3 - puddles 4 - flooded 5 - snow 6 - mud 7 - dry 8 - ice - cold 9 - other -1 - Not specified 0 - None 1 - Disconnection - tunnel 2 - Bridge - autospot 3 - Exchange or connection bridge 4 - On the road 5 - On the road 6 - On the emergency stop 7 - On the emergency stop 8 - Other	Keep
15	infra	Infrastructure	float64	0.00%	Categorical categories	more than 10 categories			Keep
16	slu	Accident situation	float64	0.00%	Categorical categories	6 to 10		3 - On the verge 4 - On the shoulder 5 - On bike path 6 - On another special track 8 - Other	Keep
17	vma (rename: speed_lim)	Speed limit	float64	0.00%	Quantitative	X		Started to be taken into account from 2019	Keep
18	arcs (rename: year)	Year	int64	0.00%	Quantitative				Keep