



Dátová kvalita vo verejnej správe

*Koncept riadenia, merania a
zlepšovania kvality dát*

ITAPA 2019, Martin Florián

28.5.2019



Obsah

1

**Prečo dátová
kvalita?**

2

Fázovanie projektu

3

**Náš prístup k
riadeniu dátovej
kvality**

4

**Náš prístup k
meraniu a
monitorovaniu
dátovej kvality**

5

**Náš prístup k
zlepšeniu dátovej
kvality**

6

**Parametre dátovej
kvality**

7

**Príklad merania
dátovej kvality**



1. Prečo Dátová kvalita?



Prečo riešime dátovú kvalitu?

Prečo riešime dátovú kvalitu?

Cieľom zlepšovania dátovej kvality nie je dátová kvalita sama o sebe, ale **vyššia kvalita života a lepšie kvalitnejšie služby verejnosti** - ktoré sú silno prepojené a závislé práve na primeranej dátovej kvalite. Dátová kvalita nie je abstraktný izolovaný cieľ, ale cesta a základňa na zvyšovanie kvality služieb, produktov a komunikácie s ňou spojených.

Aký je zmysel dátovej kvality?

Dáta sú realita a holé fakty. Čím kvalitnejšie dáta máme, tým **hodnovernejšie sú informácie a znalosti z nich odvodené**. Cieľom je nielen zvýšenie prvotných ukazovateľov dátovej kvality, ako sú presnosť, aktuálnosť a úplnosť dát, ale aj ich celková konzistentnosť, vierohodnosť, použiteľnosť, spoľahlivosť a integrita.

Prečo je potrebné meranie dátovej kvality?

Prečo meriame dátovú kvalitu?

To, čo chceme systematicky zlepšovať, musíme aj nejakým spôsobom **vyhodnocovať** **merať** a **kvantifikovať**. Aby sme mohli **objektívizovať** súčasný aj plánovaný stav a vedeli čo najobjektívnejšie sledovať trendy a dynamiku vývoja.

Prečo práve teraz?

Dátová kvalita posledné roky získava stále väčšiu pozornosť a dôležitosť pre **exponenciálny nárast objemu štruktúrovaných aj neštruktúrovaných dát** a pre práve prebiehajúcu **digitálnu transformáciu**. V tomto prostredí už nestačí doteraz obvyklý intuitívny subjektívny ad hoc prístup ku kvalite dát, ale je nevyhnutné zaviesť **koncepčný prístup založený na overenej metodike**.

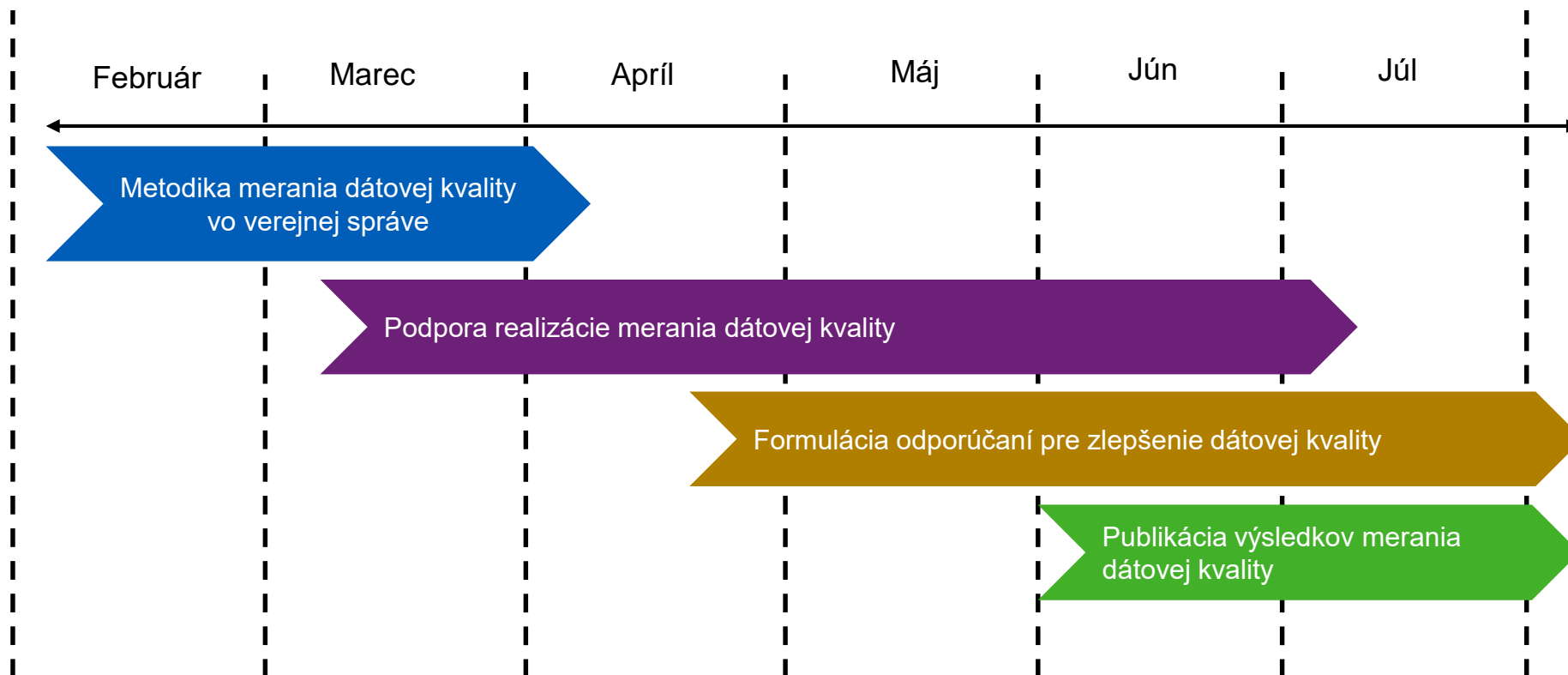


2. Fázovanie projektu



Fázovanie projektu merania dátovej kvality vo verejnej správe

Časové obdobie projektu: 8.2. - 7.8.2019





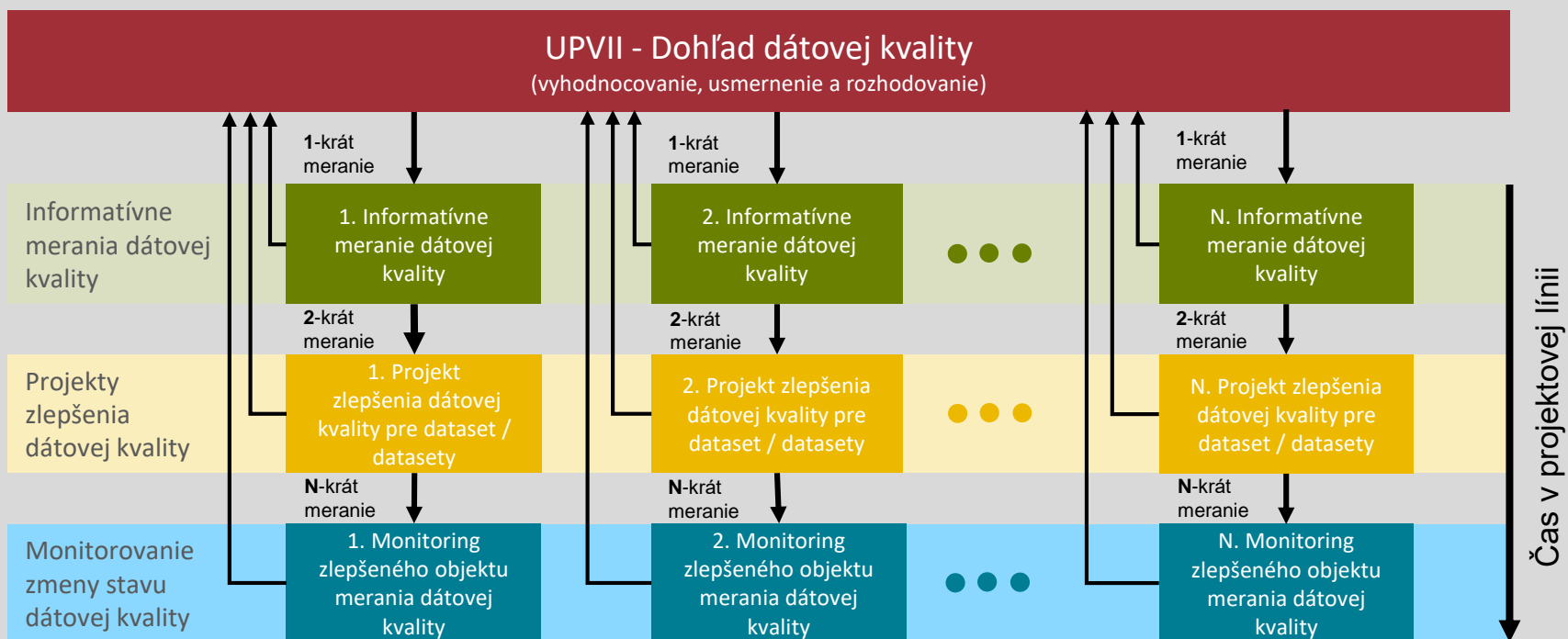
3. Náš prístup k riadeniu dátovej kvality



Zavádzanie dátovej kvality vo verejnej správe

Stratégia dátovej kvality (definovanie)

Program dátovej kvality (definovanie a riadenie)





4. Náš prístup k meraniu a monitorovaniu dátovej kvality



7-krokový postup merania a monitorovania DK

Informatívne meranie



Monitorovacie meranie



7-krokový postup merania a monitorovania DK

Komplexné meranie súčasného stavu



Kontrolne meranie po implementácii zlepšení

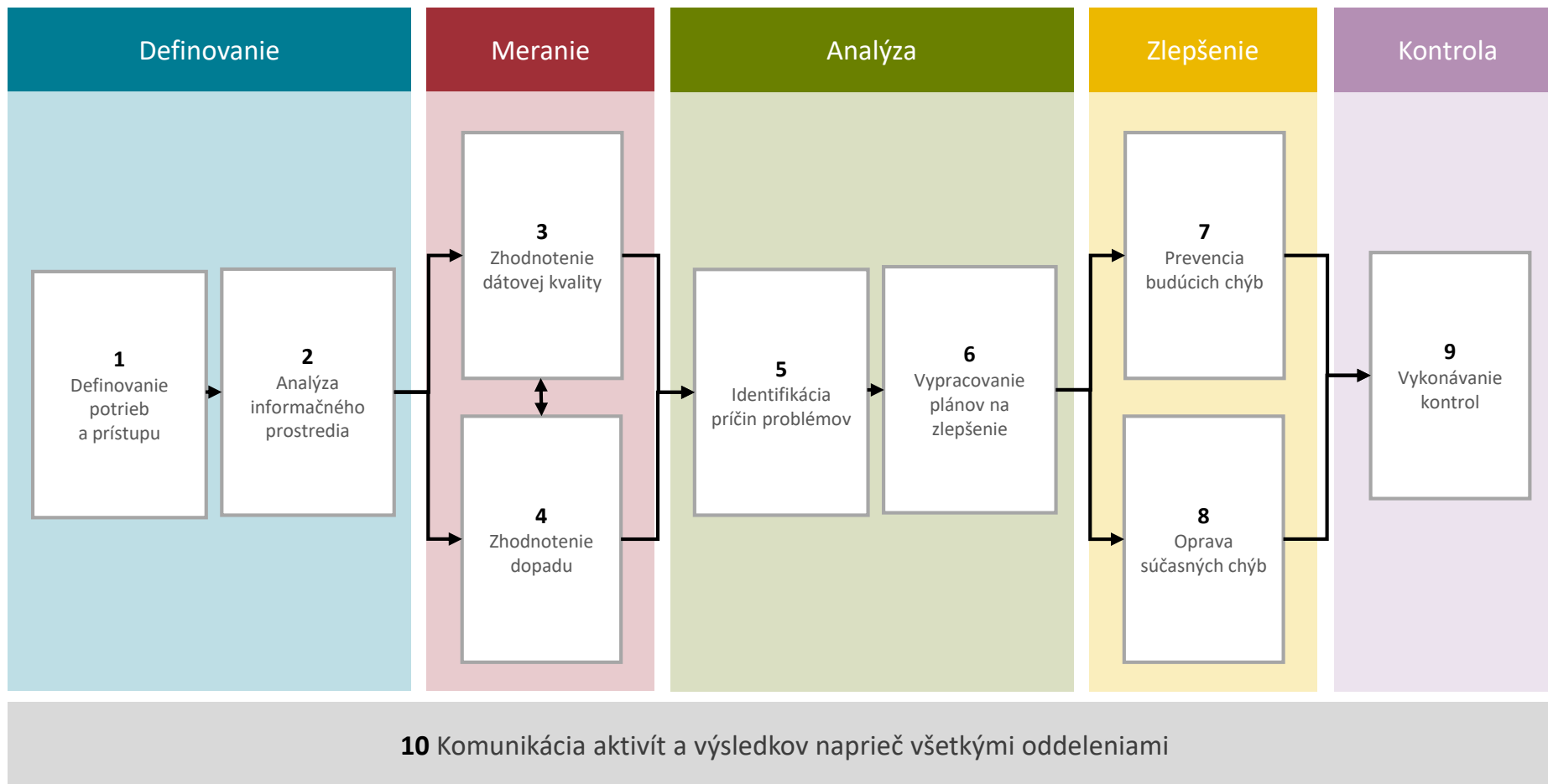




5. Náš prístup k zlepšeniu dátovej kvality



10-krokový postup zlepšenia dátovej kvality





6. Parametre dátovej kvality



Parametre dátovej kvality

Parameter	Popis
1. Presnosť	Ako dobre dáta zodpovedajú skutočnostiam v reálnom svete
2. Správnosť	Hovorí, či formát údajov zodpovedá definovaným pravidlám
3. Kompletnosť	Hovorí, či dáta obsahujú dostatok informácií
4. Unikátnosť	Hovorí o duplicitách v záznamoch
5. Aktuálnosť	Hovorí, či sú dáta aktuálne alebo staršieho dátumu
6. Strojová spracovateľnosť	Hovorí, ako jednoducho je možné údaje automatizovane spracovať
7. Referenčná integrita	Zameriava sa na použitie dát naprieč verejnou správou a previazanie údajov v rámci aj mimo nej pre ďalších konzumentov údajov



7. Príklad merania

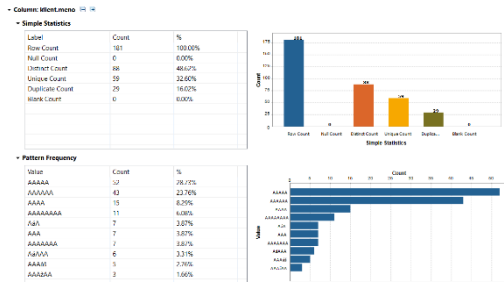


Unikátnosť - unikátnosť atribútu

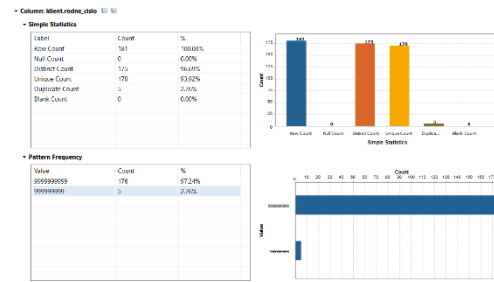
Názov ukazovateľa	Unikátnosť atribútu
Definícia	Nič sa nezaznamená viac ako 1 krát na základe toho, ako sa vec identifikuje
Metrika	Analýza počtu vecí hodnotených v reálnom svete v porovnávaní s počtom záznamov v datasete
Rozsah	Meraný voči všetkým záznamom v rámci jedného datasetu
Jednotka metriky	Percento
Príbuzné parametre	Konzistencia
Voliteľnosť	Záleží na okolnostiach
Príklad	Škola ma 120 súčasných študentov a 380 bývalých študentov (t.j. celkovo 500). Avšak databáza študenta ukazuje 520 odlišných študentských záznamov. Môže obsahovať meno Jožo Mrkvička a Jožko Mrkvička ako samostatný záznam aj napriek skutočnosti, že v škole existuje len 1 študent menom Jožko Mrkvička. To znamená, že unikátnosť je $500/520 = 96,2\%$.
Pseudokód (výpočtový vzorec)	$(\text{Number of things in real world}) / (\text{Number of records describing different things})$

Register Úpadcov - Klient

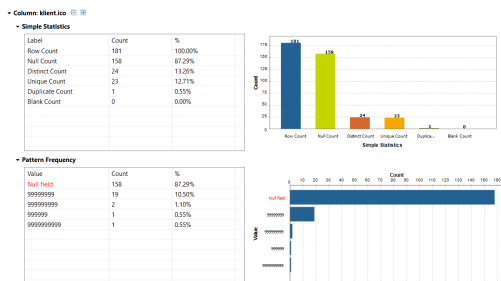
1. Meno



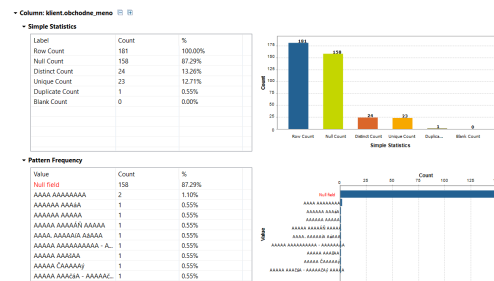
2. Rodné číslo



3. IČO



4. Obchodné meno



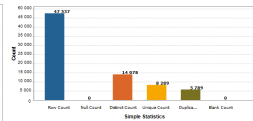
Register Adries - Street Name

1. Street Name

Column: metadata.StreetName

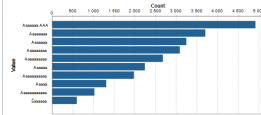
Simple Statistics

Label	Count	%
Row Count	47537	100.00%
Null Count	0	0.00%
Distinct Count	14038	29.74%
Unique Count	8289	17.44%
Duplicate Count	5789	12.18%
Blank Count	0	0.00%



Pattern Frequency

Value	Count	%
Adress AAA	4922	10.35%
AdressAA	3715	7.82%
Adressaa	3350	7.05%
Adressaaa	3093	6.51%
Adressaaaa	2962	6.23%
Adressaa	2347	4.94%
Adressaaaaaa	1983	4.17%
Adress	1932	4.07%
Adressaaaaaa	1007	2.12%
Straaaa	601	1.27%

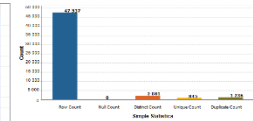


2. Valid From

Column: metadata.validFrom

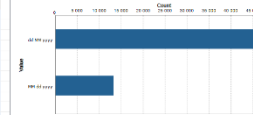
Simple Statistics

Label	Count	%
Row Count	47537	100.00%
Null Count	0	0.00%
Distinct Count	3983	8.38%
Unique Count	845	1.78%
Duplicate Count	1236	2.62%



Date Pattern Frequency

Value	Count	%
dd.MM.yyyy	47117	99.13%
MM.dd.yyyy	1320	2.78%

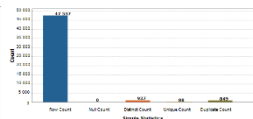


3. Valid To

Column: metadata.validTo

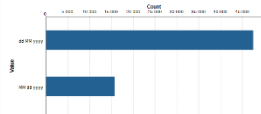
Simple Statistics

Label	Count	%
Row Count	47537	100.00%
Null Count	0	0.00%
Distinct Count	317	0.67%
Unique Count	88	0.19%
Duplicate Count	849	1.79%



Date Pattern Frequency

Value	Count	%
dd.MM.yyyy	47317	99.53%
MM.dd.yyyy	1576	3.32%

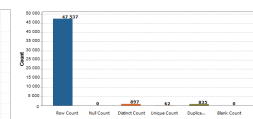


4. District Identifier

Column: metadata.districtIdentifier

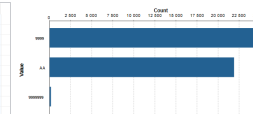
Simple Statistics

Label	Count	%
Row Count	47537	100.00%
Null Count	0	0.00%
Distinct Count	897	1.89%
Unique Count	62	0.13%
Duplicate Count	835	1.76%
Blank Count	0	0.00%



Pattern Frequency

Value	Count	%
9999	25415	53.46%
AA	21924	46.12%
9999999	198	0.42%





Ďakujeme za
pozornosť

Klient - Meno

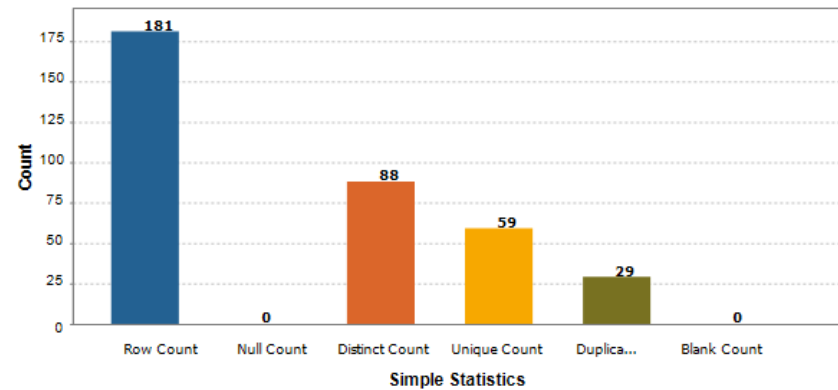
Spät



▼ Column: klient.meno

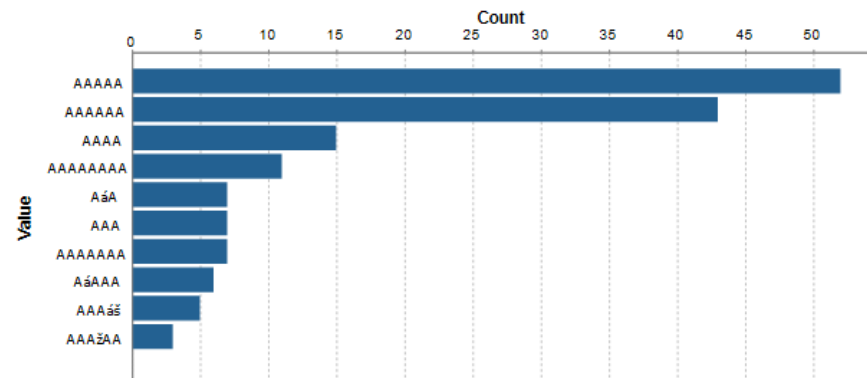
▼ Simple Statistics

Label	Count	%
Row Count	181	100.00%
Null Count	0	0.00%
Distinct Count	88	48.62%
Unique Count	59	32.60%
Duplicate Count	29	16.02%
Blank Count	0	0.00%



▼ Pattern Frequency

Value	Count	%
AAAAA	52	28.73%
AAAAAA	43	23.76%
AAAA	15	8.29%
AAAAAAA	11	6.08%
AáA	7	3.87%
AAA	7	3.87%
AAAAAAA	7	3.87%
AáAAA	6	3.31%
AAAás	5	2.76%
AAAŽAA	3	1.66%



Klient - Rodné číslo

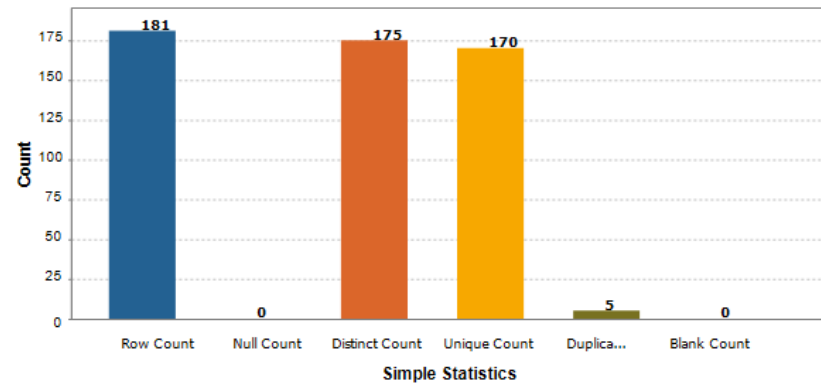
Spät



▼ Column: klient.rodne_cislo

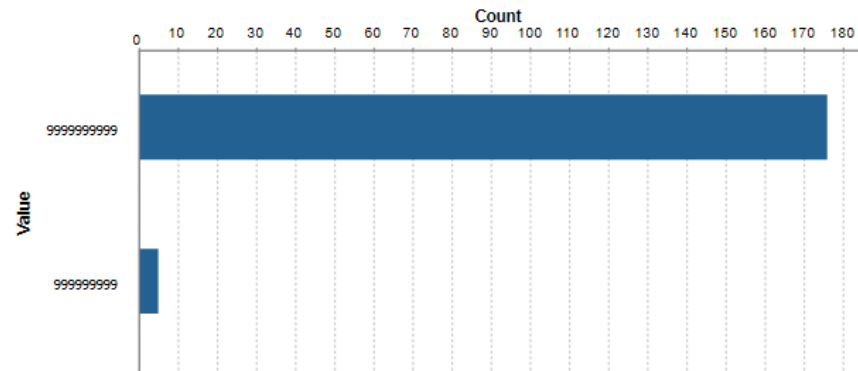
▼ Simple Statistics

Label	Count	%
Row Count	181	100.00%
Null Count	0	0.00%
Distinct Count	175	96.69%
Unique Count	170	93.92%
Duplicate Count	5	2.76%
Blank Count	0	0.00%



▼ Pattern Frequency

Value	Count	%
999999999	176	97.24%
99999999	5	2.76%



Klient - IČO

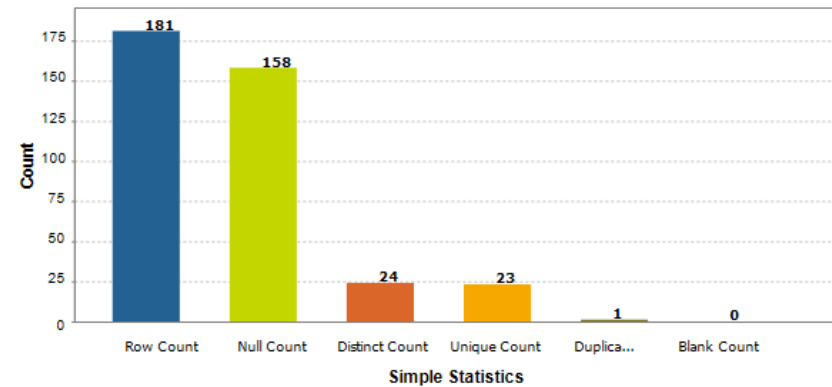
Spät



▼ Column: klient.ico

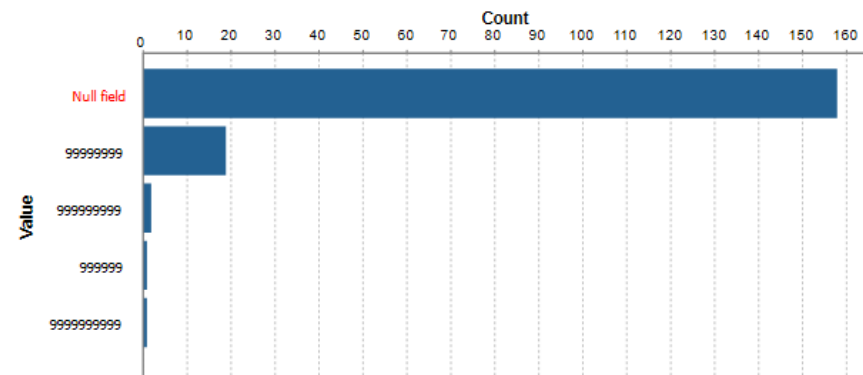
▼ Simple Statistics

Label	Count	%
Row Count	181	100.00%
Null Count	158	87.29%
Distinct Count	24	13.26%
Unique Count	23	12.71%
Duplicate Count	1	0.55%
Blank Count	0	0.00%



▼ Pattern Frequency

Value	Count	%
Null field	158	87.29%
99999999	19	10.50%
999999999	2	1.10%
999999	1	0.55%
9999999999	1	0.55%



Klient - Obchodné meno

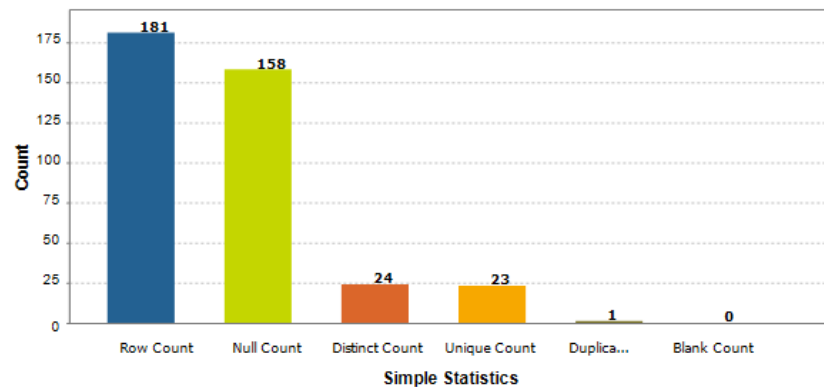
Spät'



▼ Column: klient.obchodne_meno

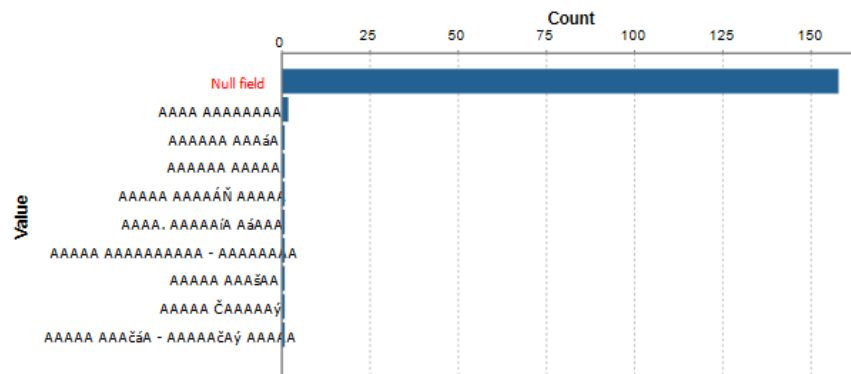
▼ Simple Statistics

Label	Count	%
Row Count	181	100.00%
Null Count	158	87.29%
Distinct Count	24	13.26%
Unique Count	23	12.71%
Duplicate Count	1	0.55%
Blank Count	0	0.00%



▼ Pattern Frequency

Value	Count	%
Null field	158	87.29%
AAAA AAAAAAAA	2	1.10%
AAAAAA AAAáA	1	0.55%
AAAAAA AAAAA	1	0.55%
AAAAA AAAAAĀ AAAAA	1	0.55%
AAAA. AAAAAíA AáAAA	1	0.55%
AAAAA AAAAAA - A...	1	0.55%
AAAAA AAAšAA	1	0.55%
AAAAA ČAAAAý	1	0.55%
AAAAA AAAčáA - AAAAAč...	1	0.55%



Street Name

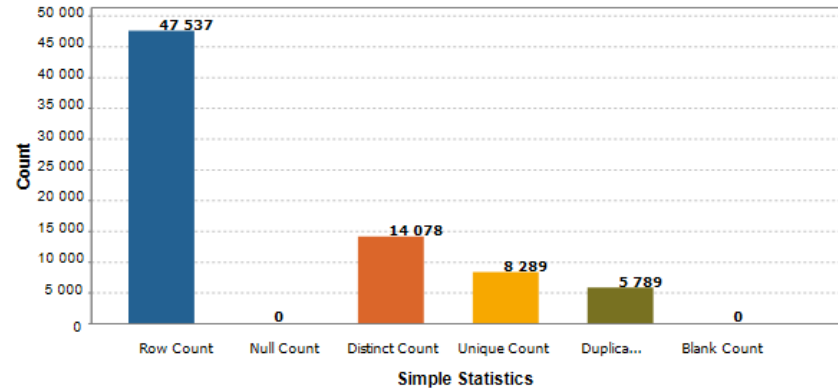
Spät



▼ Column: metadata.StreetName

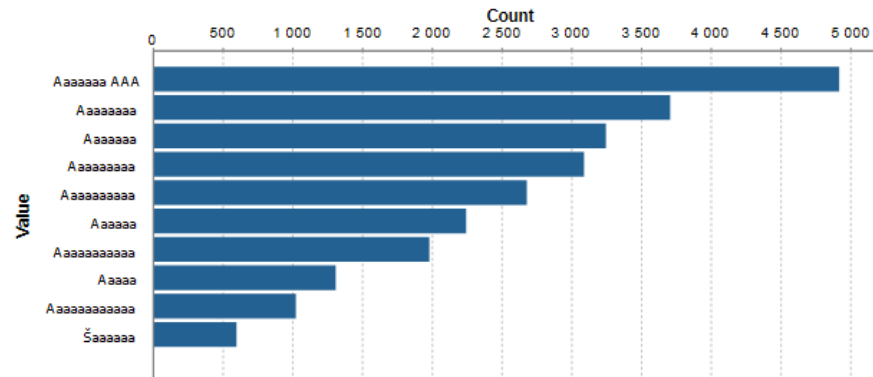
▼ Simple Statistics

Label	Count	%
Row Count	47537	100.00%
Null Count	0	0.00%
Distinct Count	14078	29.61%
Unique Count	8289	17.44%
Duplicate Count	5789	12.18%
Blank Count	0	0.00%



▼ Pattern Frequency

Value	Count	%
Aaaaaaa AAA	4922	10.35%
Aaaaaaaa	3710	7.80%
Aaaaaaa	3250	6.84%
Aaaaaaaaa	3093	6.51%
Aaaaaaaaaa	2682	5.64%
Aaaaaa	2247	4.73%
Aaaaaaaaaaaa	1983	4.17%
Aaaaa	1312	2.76%
Aaaaaaaaaaaaa	1027	2.16%
Šaaaaaa	601	1.26%



Valid From

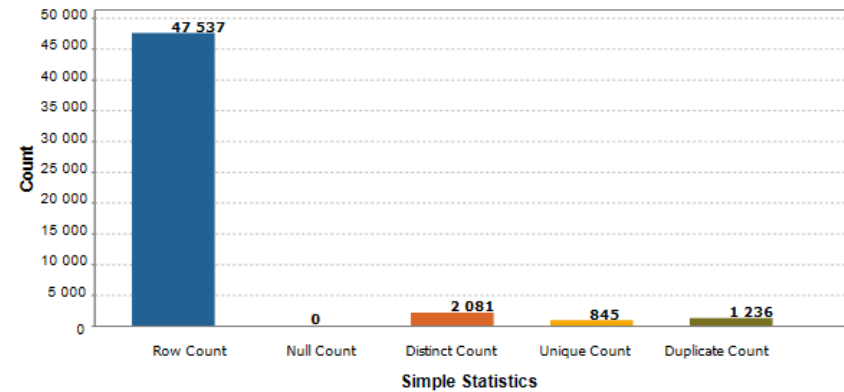
Spät



▼ Column: metadata.validFrom

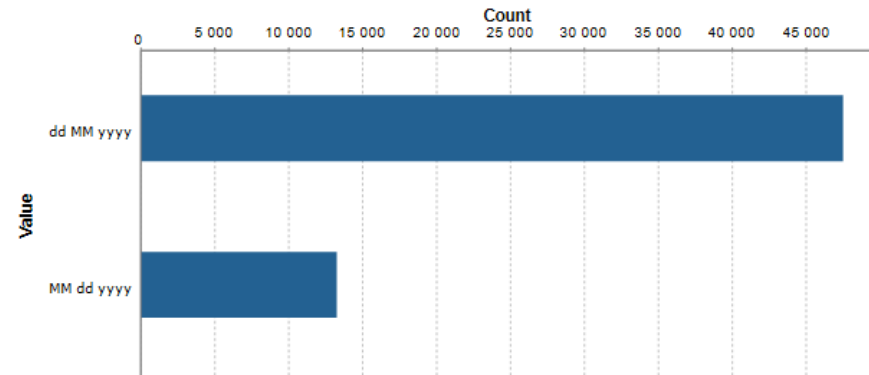
▼ Simple Statistics

Label	Count	%
Row Count	47537	100.00%
Null Count	0	0.00%
Distinct Count	2081	4.38%
Unique Count	845	1.78%
Duplicate Count	1236	2.60%



▼ Date Pattern Frequency

Value	Count	%
dd MM yyyy	47537	100.00%
MM dd yyyy	13289	27.96%



Valid To

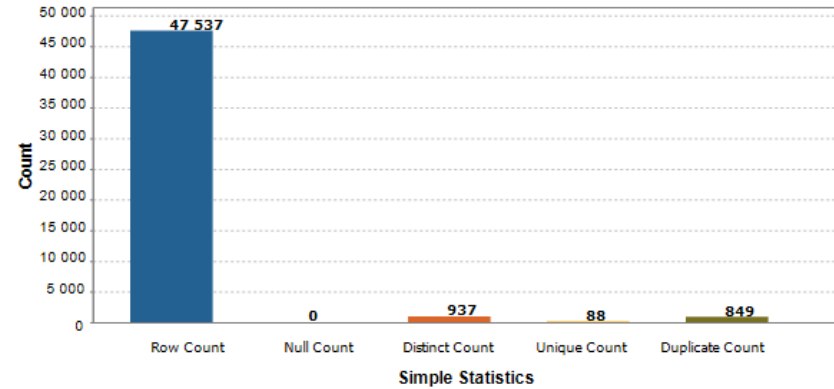
Spät



Column: metadata.validTo

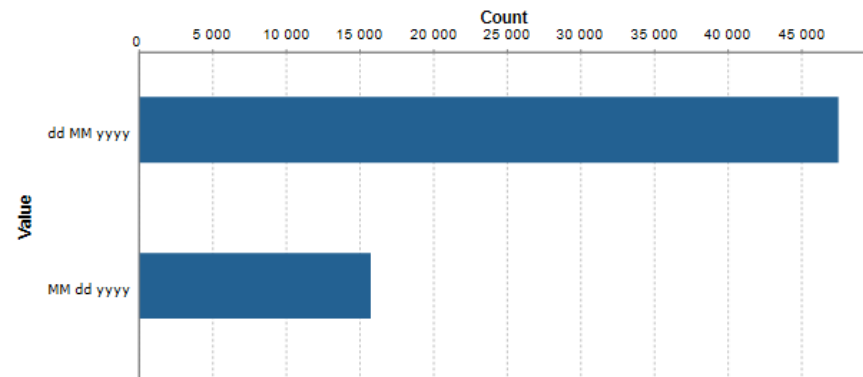
Simple Statistics

Label	Count	%
Row Count	47537	100.00%
Null Count	0	0.00%
Distinct Count	937	1.97%
Unique Count	88	0.19%
Duplicate Count	849	1.79%



Date Pattern Frequency

Value	Count	%
dd MM yyyy	47537	100.00%
MM dd yyyy	15776	33.19%



District Identifier

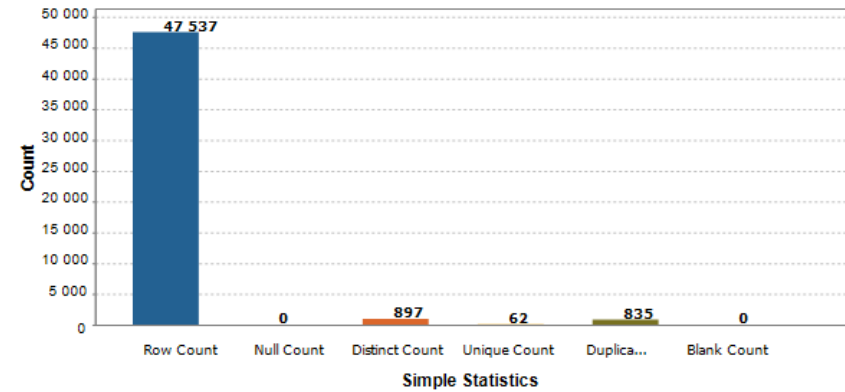
Spät



▼ Column: metadata.districtIdentifier

▼ Simple Statistics

Label	Count	%
Row Count	47537	100.00%
Null Count	0	0.00%
Distinct Count	897	1.89%
Unique Count	62	0.13%
Duplicate Count	835	1.76%
Blank Count	0	0.00%



▼ Pattern Frequency

Value	Count	%
9999	25415	53.46%
AA	21924	46.12%
9999999	198	0.42%

