

Uporaba Excel Power Query-a za združevanje podatkov iz Si-Stat

Author	Version	Status	Date
mz	0.1.0	draft	2023-04-13
mz	1.0.0	draft	2023-04-17

Da ne pozabiš vsega naučenega na tečaju napredne uporabe Excela, imaš tukaj osvežitveno domačo nalogo, ki vključuje tudi precej navodil in poglavje Troubleshooting.

Kazalo

1	Nagradna igra	2
1.1	Osnovna verzija	2
1.2	Napredna verzija	2
1.3	Še najnaprednejšja verzija	3
2	Priprava poizvedbe na Si-Stat	3
3	Priprava povezave v Excel Power Query-u.	4
4	Nalaganje podatkov	5
4.1	Osveževanje podatkov	6
5	Transformacije podatkov	7
5.1	Pogoste transformacije	8
5.1.1	Column split	8
5.1.2	Dodajanje novega stolpca	8
5.1.3	Preimenovanje stolpca	8
5.1.4	Agregiranje	8
5.1.5	Filtriranje vrstic	8
5.1.6	Združevanje tabel	9
5.2	Vrstni red transformacij	9
6	Troubleshooting	10
6.1	Troubleshooting manjkajočih stolpcev	10
6.2	Troubleshooting manjkajočih decimal	11
6.3	Troubleshooting neobičajnih manjkajočih vrednosti	12

1 Nagradna igra

Naloga ima več verzij, izberi si tisto, ki ti ustreza, ali pa začni z osnovno in jo potem nadgradi. Vsa navodila potrebna za osnovno in napredno verzijo so opisana v nadaljevanju. Za najnaprednejšjo pa bo potrebno malce eksperimentiranja.

Prvi trije, ki končate nalogo, dobite nagrado :)

1.1 Osnovna verzija

1. Izberi si dve časovni seriji iz dveh različnih tabel na Si-Stat-u z enakim časovnim intervalom (recimo obe četrtletni). (Zaradi točke 4. premisli kateri spremenljivke izbereš.)
2. Odpri nov Excel fajl in s pomočjo spodnjih navodil vzpostavi podatkovni povezavi z obema viroma podatkov.
3. Ko imaš vzpostavljeni obe povezavi (ni treba naložiti podatkov, dovolj je vzpostaviti “connection only” povezavo), združi obe tabeli glede na časovno dimenzijo. (Glej razdelek *Združevanje tabel* spodaj.) To bo seveda delovalo samo, če imata obe seriji enake časovne intervale. Če ne, si v napredni verziji naloge in kar pogumno nadaljuj tam!
4. Nato dodaj nov stolpec, ki je izračunan iz obeh časovnih serij. (Glej razdelek *Dodajanje novega stolpca* spodaj.)
5. Ko si zadovoljen s tem, kako si pripravil poizvedbo, naloži podatke na nov zavihek in jih prikaži z enostavnim linijskim grafom.
6. Zdaj se vrni v Power Query Editor in poizvedbi dodaj filter npr. izberi samo podatke po letu 2010. Preveri kaj se je zgodilo z grafom.

1.2 Napredna verzija

1. Izberi si dve časovni seriji iz dveh različnih tabel na Si-Stat-u z *različnim* časovnim intervalom (recimo ena mesečna in ena letna). (Zaradi točke 5. premisli kateri spremenljivke izbereš.)
2. Odpri nov Excel fajl in s pomočjo spodnjih navodil vzpostavi podatkovni povezavi z obema viroma podatkov.
3. Tisto serijo, ki ima bolj podrobne podatke, agregiraj na raven druge serije (npr. mesečne podatke agregiraj na letne). Seveda pazi, da izbereš ustrezno funkcijo glede na tip podatka. (Glej razdelek *Agregiranje* spodaj.)
4. Ko imaš vzpostavljeni obe povezavi (ni treba naložiti podatkov, dovolj je vzpostaviti “connection only” povezavo), združi obe tabeli glede na časovno dimenzijo. (Glej razdelek *Združevanje tabel* spodaj.)
5. Nato dodaj nov stolpec, ki je izračunan iz obeh časovnih serij. (Glej razdelek *Dodajanje novega stolpca* spodaj.)
6. Ko si zadovoljen s tem, kako si pripravil poizvedbo, naloži podatke na nov zavihek in jih prikaži z enostavnim linijskim grafom.
7. Zdaj se vrni v Power Query Editor in poizvedbi dodaj filter npr. izberi samo podatke po letu 2010. Preveri kaj se je zgodilo z grafom.

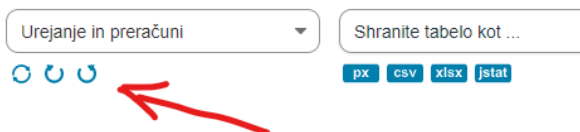
1.3 Še najnaprednejšaja verzija

1. Izberi si dve časovni seriji iz dveh poljubnih (različnih) virov. Npr Si-Stat, Eurostat, obstoječ Excel ali csv ali pdf fajl na tvojem disku - najbolje nekaj, kar sicer pogosto uporabljaš. Odvisno od vira, ki si si ga izbral, ti bodo spodnja navodila bolj ali manj v pomoč. V vsakem primeru jih preberi, potem pa uporabi google, sodelavce ali Majo Z., da najdeš rešitev kako vzpostaviti želene povezave.
2. Nadaljuj s točko 2. naprednih navodil.

2 Priprava poizvedbe na Si-Stat

1. Na strani Si-Stat poiščeš tabelo iz katere želiš dobiti podatke.
2. Izbereš zelene kategorije za vsako dimenzijo in klikneš “Izpis podatkov”.
3. Preveri, da imaš tabelo pravilno zavrteno in orientirano: to pomeni, (i) da imaš podatke za vsako obdobje v svoji vrstici, torej čas poteka od zgoraj navzdol in (ii) da imaš enovrstični header, torej je v stolpcih samo ena dimenzija, ostale pa vse v vrsticah. Da to dosežeš, uporabi gumbke za vrtenje tabele, dokler ne dobiš zelenega rezultata.

Prikaz tabele O tabeli Statistična znamenja



4. Potem klikni na tekst “Shrani poizvedbo” nad tabelo. Ker hočemo dinamično poizvedbo, ki se podaljšuje, preveri, da imaš izbrano prvo možnost “Stalni začetni in drseči zaključni časovni presek (podaljševanje serije)”.
5. Za izpis serije pa izberi “CSV (ločeno s podpičjem), brez glave (.csv)”. In klikni na “Shrani poizvedbo”

Shranjevanje poizvedb

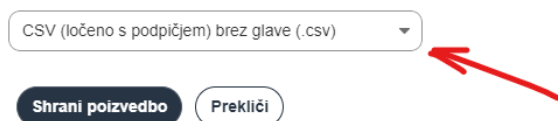
Izbor podatkov lahko shranite kot shranjeno poizvedbo z možnostjo izbora posodobitve časovne serije. Izberete vrsto izpisa oz. formata podatkov in potrdite vaš izbor s klikom na gumb Shrani poizvedbo.

Shranjevanje poizvedb - nastavitve

Kako želite nastaviti shranjeno poizvedbo, ko posodobimo tabelo z novo časovno točko?

- ☒ Stalni začetni in drseči zaključni časovni presek (podaljševanje serije)
- ☐ Drseča časovna vrsta (vedno enako število časovnih točk)
- ☐ Vedno uporabi izbrano obdobje (ne posodobil z novo časovno točko)

Izpis podatkov (izberi):



6. Prikazalo se ti bo polje z url-jem tvoje poizvedbe. To si skopiraj, ker ga boš rabil pri pripravi povezave v Power Query-u.

⊖ Shrani poizvedbo

Časovna dimenzija:

Stalni začetni in drseči zaključni časovni presek (podaljševanje serije)

Izpis/Format:

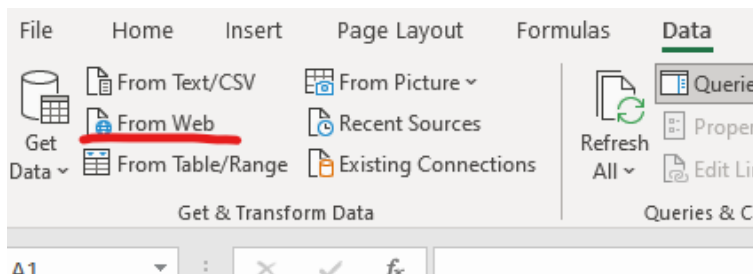
CSV (ločeno s podpičjem) brez glave (.csv)

Do shranjene poizvedbe se vrnete s klikom na to povezavo

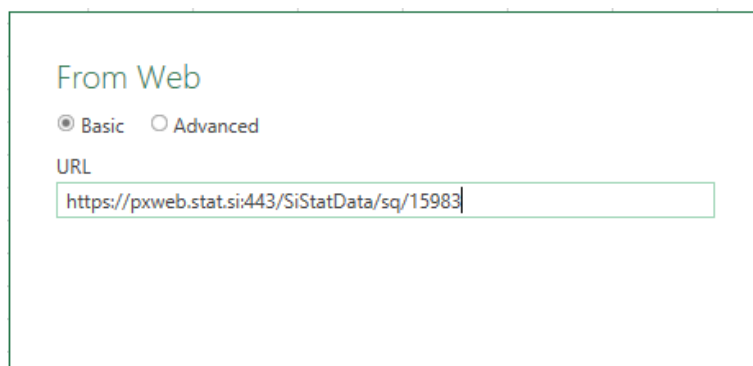
<https://pxweb.stat.si:443/SiStatData/sq/15983>

3 Priprava povezave v Excel Power Query-u.

1. Odpri nov svež prazen Excel fajl.
2. Na zavihku Data izberi From Web



3. V polje URL skopiraj url, ki si ga dobil v zadnji točki priprave poizvedbe in klikni OK.



4. če je šlo vse po sreči, bi moral videti nekaj podobnega kot je na spodnji sliki. Najbolj pomembno je, da se v tem predogledu vidi več stolpcev. Glej troubleshooting poglavje, če ne gre povsem gladko.

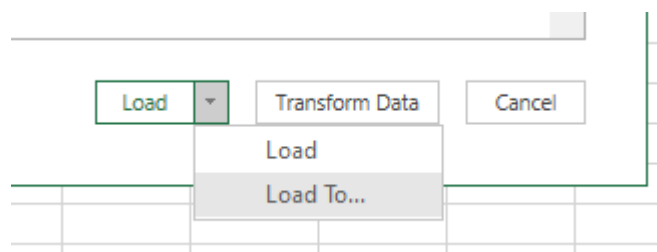
<https://pxweb.stat.si/SiStatData/sq/15991>

File Origin		Delimiter	Data Type Detection	
1250: Central European (Windows)		Semicolon	Based on first 200 rows	
MERITVE	ČETRLETJE	Notranji / mednarodni prevoz - SKUPAJ		Mednarodni prevoz - blago naloženo v Sloveniji
Tone (1000)	2001Q1	11403		845
Tone (1000)	2001Q2	15598		866
Tone (1000)	2001Q3	18270		879
Tone (1000)	2001Q4	12539		629
Tone (1000)	2002Q1	13423		729
Tone (1000)	2002Q2	17224		811

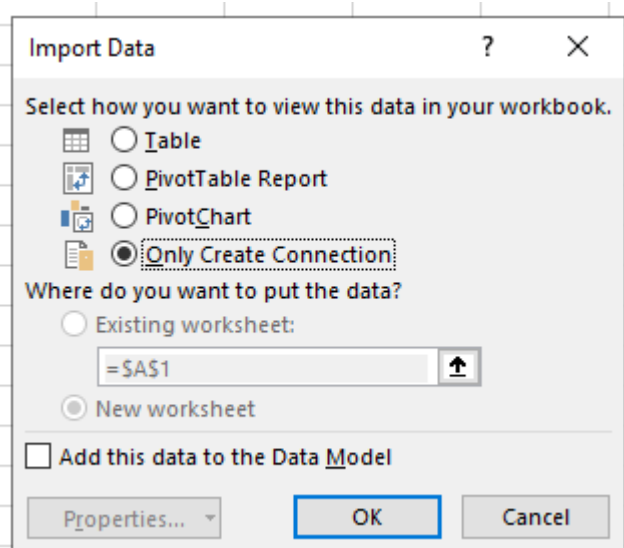
5. Zdaj sta na voljo dve možnosti, odvisno od tega kaj nameravaš delati s podatki: podatke lahko naložiš (*Load*) ali pa se lotiš nadaljnjih transformacij (*Transform*). Nobena odločitev pa ni dokončna - če podatke naložiš, jih še vedno lahko kasneje spreminjaš in seveda obratno.

4 Nalaganje podatkov

Pri nalaganju podatkov so relevantne naslednje možnosti:



1. *Load*: default nalaganje naloži tabelo v nov zavihek.
2. *Load to*: ti da več možnosti, kjer lahko recimo določiš lokacijo tabele.
3. *Load to*: in izbira *Only create connection* je tudi zelo uporabna možnost, kadar ne rabimo tabele neposredno, ampak jo bomo uporabili kasneje, kot input za novo poizvedbo.



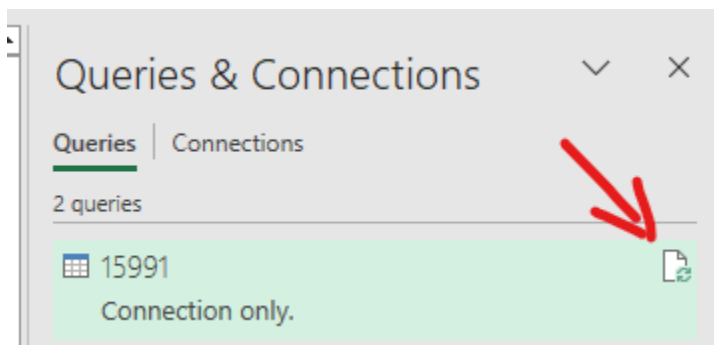
Ne glede na to katero možnost si izbral, imaš zdaj v Excelu aktivno povezavo na izbrani vir podatkov. Ta vir lahko urejaš tako, da zavihku *Data* izbereš *Queries & Connections*, kar ti bo odprlo seznam odprtih povezav oziroma poizvedb na desni strani ekrana.

4.1 Osveževanje podatkov

Povezava na podatke oz. poizvedba, ki si jo pripravil, se ne osvežuje avtomatično kadar pride do spremembe pri viru.

Podatke lahko osvežiš na več načinov:

1. Z ikonico desno od vsake posamezne povezave.



2. Z ikonico *Refresh all* na vrhu zavihka *Data*, ki osveži vse poizvedbe.
3. Lahko pa si nastaviš lastnosti vsake povezave, kjer določiš kdaj se poizvedba samodejno posodobi. Do tega okenca prideš tako da z desno klikneš na poizvedbo v seznamu na desni strani ekrana in izbereš možnost *Properties* na dnu.

Query Properties

Query name: Promet

Description: Tonaža mednarodnega prevoza blaga naloženega v Sloveniji - iz Si-stat tabele 22077015

Usage Definition Used In

Refresh control

Last Refreshed:

☒ Enable background refresh

☐ Refresh every 60 minutes

☒ Refresh data when opening the file

☐ Remove data from the external data range before saving the workbook

☒ Refresh this connection on Refresh All

☐ Enable Fast Data Load

OLAP Server Formatting

Retrieve the following formats from the server when using this connection:

☐ Number Format ☐ Fill Color

☐ Font Style ☐ Text Color

OLAP Drill Through

Maximum number of records to retrieve:

Language

☐ Retrieve data and errors in the Office display language when available

OK Cancel

4. Najbolj uporabna kljukica tukaj je *Refresh data when opening file the file*, ki naredi to, kar piše ;)
5. Velja biti pozoren na opcijo *Enable background refresh*: to se zdi uporabno, ker pomeni, da med osveževanjem lahko uporabljaš Excel in delaš kaj drugega. Ampak pomeni pa tudi, da Excel ne čaka z drugimi stvarmi, medtem ko se v ozadju posodabljaajo podatki, in v praksi to lahko pomeni, da recimo vrtilna tabela iz teh podatkov ne bo počakala, da se podatki najprej posodobijo - v tem primeru jo bo treba še enkrat osvežit. Samo omenim, če se komu to zgodi, da ve zakaj.
6. Poleg tega ni slaba ideja, da povezavo poimenuješ bolj opisno, kar tudi narediš v tem pogovornem okencu, na vrhu v polju *Query name* in po želji še *Description*.

5 Transformacije podatkov

Z izbiro opcije *Transform* preden podatke naložiš, se odpre t.i. "Power Query Editor". Ravno tako se odpre, če dvakrat klikneš na poizvedbo na desni strani ekrana.

Power query editor je tako zelo powerful, da se tukaj niti približno ne moremo dotakniti velike večine funkcionalnosti - vedi le, da če si lahko zamisliš neko transformacijo, se jo verjetno da narediti. Zato vprašaj ali pogoogljaj :)

Pozor! Če iz Si-Stata uvažáš podatke z decimalnimi števkami in/ali podatke z manjkajočimi vrednostmi, potem si nujno poglej zadnja dva razdelka troubleshooting-a, kjer so razloženi koraki, kako podatke pravilno transformirati v tem primerih.

5.1 Pogoste transformacije

Tukaj bom naštela samo nekaj primerov transformacij, ki lahko pridejo prav:

5.1.1 Column split

Zavihek *home*, *Split column*.

- izberi stolpec in klikni na *Split column* in izberi metodo po kateri hočeš, da se stolpec razdeli.
- klasičen primer je recimo 2020Q04 razbiti na stolpca 2020 in Q04, tako da ga razdeliš na četrti poziciji.

5.1.2 Dodajanje novega stolpca

Zavihek *Add column*, gumb *Custom Column*.

- izberi ime novega stolpca in v spodnje okence vnesi formulo tako, da s pomočjo gumba *Insert* prestavljaš zelene stolpce iz desnega okenca.

5.1.3 Preimenovanje stolpca

Zavihek *Transform*, *Rename*

- Izberi stolpec in klikni *Rename* in preimenuj stolpec po želji.

5.1.4 Agregiranje

Zavihek *Home* ali *Transform*, *Group By*

- Izberi stolpec, po katerem želiš grupirati in klikni *Group By*
- v pogovornem okencu, ki se odpre, določi katere agregacije želiš
- če jih želiš več, klikni na *Advanced*
- npr, če imaš mesečne podatke, jih lahko združiš po letu in glede na tip podatka izbereš vsote ali povprečja za vsak stolpec, ki te zanima.

5.1.5 Filtriranje vrstic

S klikom na trikotniček desno od imena stolpca lahko stolpec filtriraš ali razvrstiš po želji.

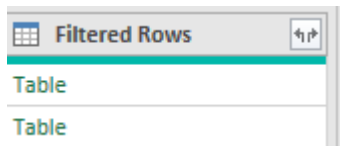
5.1.6 Združevanje tabel

Zavihek Home, *Merge queries*

Za združevanje tabel rabiš imeti vsaj dve tabeli oz. dve povezavi odprti. Poleg tega morata imeti obe tabeli vsaj en stolpec, po katerem se lahko združita.

Združevanja se lotiš tako, da odpreš Power Query Editor v eni od obeh tabel.

Možnih je več tipov združevanj - levo, desno, notranje ipd. najbolje, da malo eksperimentiraš. Ko potrdiš merge, se bo v tabeli pojavil not stolpec z ikonico z dvema puščicama. Na to ikonico klikneš in izbereš katere stolpce iz druge tabele želiš prikazane (če ne želiš vseh).

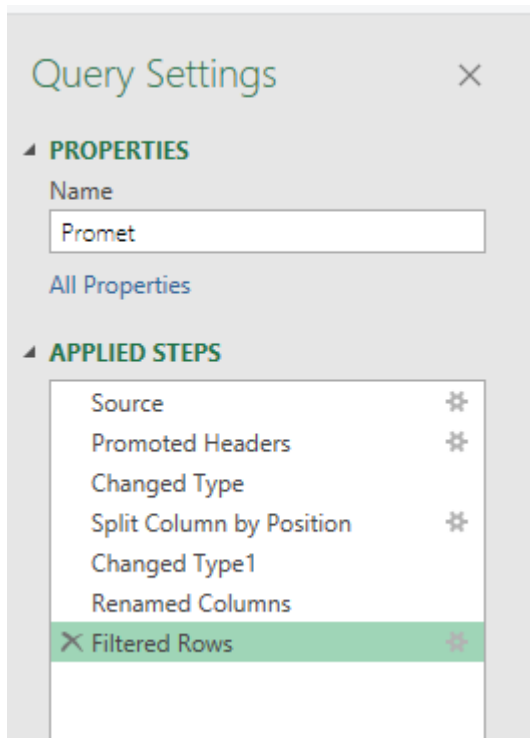


5.2 Vrstni red transformacij

Po vsaki transformaciji boš opazil, da se v okencu na desni strani dodajajo novi koraki tvoje poizvedbe.

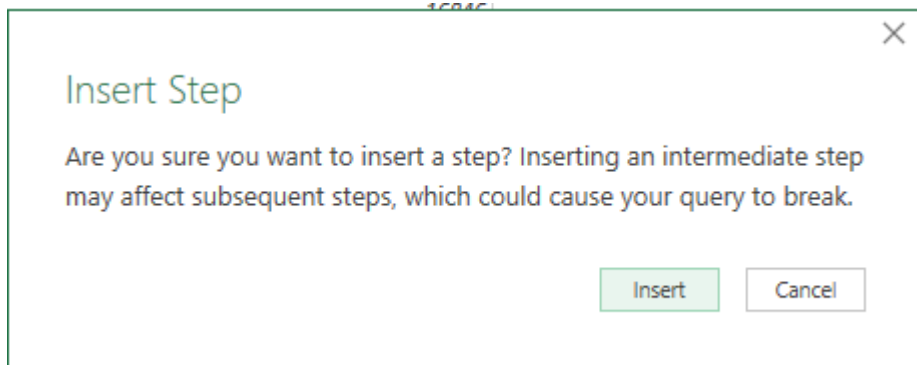
Tukaj lahko:

- zbrišeš korak s klikom na križec na levi strani pred zapisom
- spreminjaš vrstni red korakov z metodo “drag and drop”
- ponovno urejaš posamezen korak, če ga želiš popraviti, tako, da dvakrat klikneš na njega.



Prvi trije koraki: Source, Promoted Headers in Changed type se avtomatično pripravijo ob pripravi povezave. Te tri korake res predlagam, da pustiš pri miru, razen če veš, kaj delaš (cf. Troubleshooting). Ostale korake praviloma dodajaš po vrsti, dokler nisi zadovoljen s predogledom, ki ga vidiš v editorju.

Vrstni red korakov je pomemben! Recimo če najprej razdeliš stolpec s četrtnetji na leto in četrtnetje in potem grupiraš podatke glede na leto, ne moreš tega narediti v obratnem vrstnem redu. Zato se ne ustrašiti, če se ti pojavi spodnje opozorilo: Excel samo preverja, da res želiš nov korak vrniti med že obstoječe - ker ponavadi hočeš dodati korak na koncu.



Ko si zadovoljen s predogledom poizvedbe, pa izbereš **Close & Load in**, če še nisi, izbereš lokacijo, kamor želiš, da se ti tabela z rezultati poizvedbe naloži.

6 Troubleshooting

6.1 Troubleshooting manjkajočih stolpcev

Včasih se - iz meni misterijoznih razlogov - zgodi, da pri uvozu podatkov Excel noče pravilno razumeti stolpcev oz. delimiterjev. Težko reproduciram to situacijo, ampak po mojih izkušnjah se manj pogosto zgodi, če:

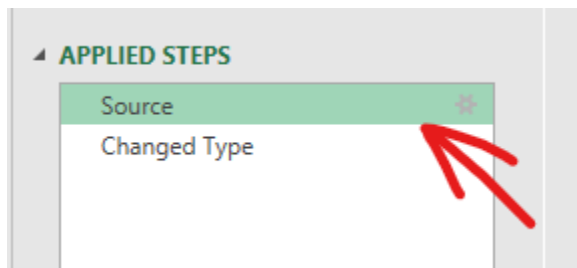
- startaš z novim, svežim Excel fajlom,
- Si-stat poizvedbo izvoziš brez glave (glava tukaj ne pomeni headerja, torej vrstice z imeni stolpcev, ampak tekst nad tem, ki ni del prave tabele),
- vse v prvo narediš prav :)

Včasih se navkljub vsemu zgodi naslednja situacija: čeprav imaš kot delimiter izbrano podpičje (*semicolon*), v predogledu vidiš samo prvi stolpec.

<https://pxweb.stat.si/SiStatData/sq/16067>

File Origin	Delimiter	Data Type Detection						
1250: Central European (Windows) ▾	Semicolon ▾	Based on first 200 rows ▾						
<table><thead><tr><th>Column1</th></tr></thead><tbody><tr><td>Bruto domači proizvod, Slovenija, letno</td></tr><tr><td>MERITVE</td></tr><tr><td>Stalne cene predhodnega leta (mio EUR)</td></tr><tr><td>Stalne cene, referenčno leto 2010 (mio EUR)</td></tr><tr><td>Letna sprememba obsega (%)</td></tr></tbody></table>			Column1	Bruto domači proizvod, Slovenija, letno	MERITVE	Stalne cene predhodnega leta (mio EUR)	Stalne cene, referenčno leto 2010 (mio EUR)	Letna sprememba obsega (%)
Column1								
Bruto domači proizvod, Slovenija, letno								
MERITVE								
Stalne cene predhodnega leta (mio EUR)								
Stalne cene, referenčno leto 2010 (mio EUR)								
Letna sprememba obsega (%)								

V tem primeru - če nočeš začeti spet od začetka - lahko zadevo rešiš takole: klikni na gumb *Transform*, da se odpre Power Query Editor. Tu zdaj vidiš samo prvi stolpec, zato *enkrat* klikni na prvi korak poizvedbe v desnem okencu: *Source*



Zdaj pa poglej v funkcijsko vrstico, kjer se ti je izpisala koda za ta korak poizvedbe, ki izgleda približno takole:

```
= Csv.Document(Web.Contents("https://pxweb.stat.si:443/SiStatData/sq/16067"),[Delimiter=";", Columns=1, Encoding=1250, QuoteStyle=QuoteStyle.None])
```

Težava je torej, da ima power query navodilo, da uvozi samo en stolpec. Zdaj pa ročno popravi število stolpcev za tekstom **Columns=** in pritisni enter. Če ne veš točno, koliko stolpcev je probaj in ko se posodobi predogled, boš videl ali imaš vse podatke ali je treba število stolpcev povečati.

6.2 Troubleshooting manjkajočih decimalk

Če boš na Si-Statu izbral sprejemljivko, ki ima decimalke, se bo zgodilo naslednje: Si-Stat pri izvozu uporablja angleški sistem oznak, torej decimalno piko in ne vejice, toda tvoj Excel je nastavljen na slovenske regionalne nastavitve in zato decimalne pike ne razume.

Če imaš v podatkih decimalke, potem ti ne bo všeč, kaj se je zgodilo v koraku *Change type*, ker je namreč ta korak verjetno narobe določil tip spremenljivke. Namesto decimalne številke je določil tip številke (integer – `Int64.Type`) in tako je 80.0 namesto 80,0, postala 800.

Query Settings

PROPERTIES

Name: 16068

APPLIED STEPS

- Source
- Promoted Headers
- Changed Type

Formula Bar:

```
= Table.TransformColumnTypes(#"Promoted Headers",{{"DEJAVNOST", type text}, {"ČETRTLETJE", type text}, {"Izkoriščenost proizvodnih zmogljivosti", Int64.Type}, {"Ustreznost proizvodnih zmogljivosti", Int64.Type}, {"Konkurenčni položaj na domačem trgu", type text}})
```

ČETRTLETJE	Izkoriščenost proizvo...	Ustreznost proizvodn...	Konkurenčni položaj...
Q1	800	26	...
Q2	799	37	...
Q3	779	41	...
Q4	796	31	...
Q1	791	25	...
Q2	786	25	...

Najenostavnejši način, ki sem ga našla do zdaj, je, da na zavihku *Data* izbereš gumb *Get data* in na koncu seznama *Query options*. V okencu, ki se odpre, izbereš pod *Regional settings* eno – katerokoli – od angleških možnosti.

Query Options

GLOBAL

General

Data Load

Power Query Editor

Security

Privacy

Diagnostics

CURRENT WORKBOOK

Data Load

Regional Settings

Privacy

Locale ⓘ

English (Madagascar) ▼

Potem moraš najverjetneje ponoviti uvoz, da 'zagrabi'.

S to rešitvijo je seveda potencialna težava, če imaš v enem fajlu povezave na vire, ki uporabljajo različne sisteme: enega z decimalnimi vejicami in enega z decimalnimi pikami. Potem se je treba zadeve lotiti drugače, tako da popraviš samo poizvedbo v koraku *Change type*, ampak računam, da je to tako redko, da ne bom še tega malce bolj zakompliciranega postopka tukaj razlagala, seveda pa me lahko vedno pridete vprašati.

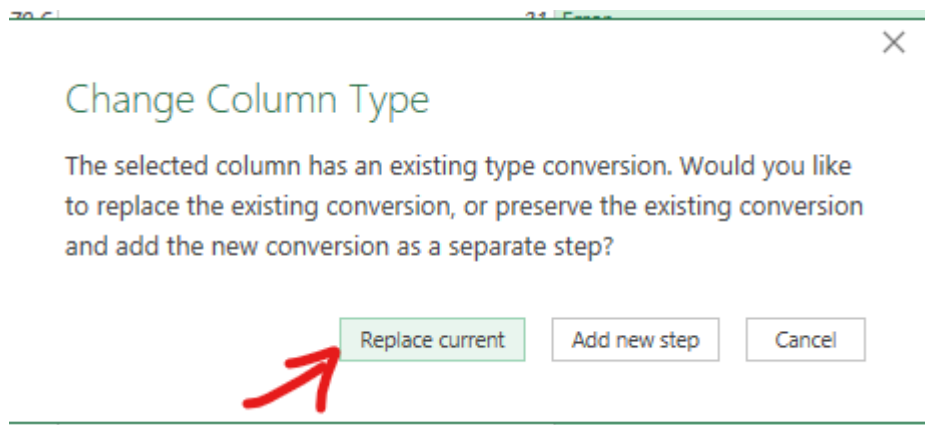
6.3 Troubleshooting neobičajnih manjkajočih vrednosti

Če boš na Si-Statu izbral spremenljivko, ki ima manjkajoče vrednosti, bodo te verjetno označene s tremi pikami (...) ali čem podobnim, česar Excel ne bo razpoznal kot manjkajočo vrednost.

V tem primeru bo stolpec uvozil kot tekst namesto kot številko, kar pa ti lahko povzroča težave kasneje, zato je to najbolje urediti že v Power Query-u.

Najprej se na desnem seznamu korakov postavi na obstoječ korak *Change Type*. Nato se postavi na dotični stolpec in iz zgornjega menija ali iz menija po desnem kliku izberi možnost *Data Type* ali *Change data type* in izberi ustrezen tip - verjetno *Decimal number* ali pa *Whole number*.

Power query te bo opozoril, da stolpec že ima spremembo tipa, in na to mu odgovoriš, da želiš zamenjati to spremembo: *Replace current*.



Po tej spremembi boš opazil naslednje: namesto ... je zdaj v teh poljih *Error*, ker seveda tri pike niso pravilnega tipa. Zato je potreben naslednji korak: errorje je treba zamenjati na *null*.

Klikni na mali trikotniček zraven gumba *Replace values* in izberi drugo možnost: *Replace errors*. V okencu, ki se ti odpre, v polje vpiši *null* in klikni OK.

