

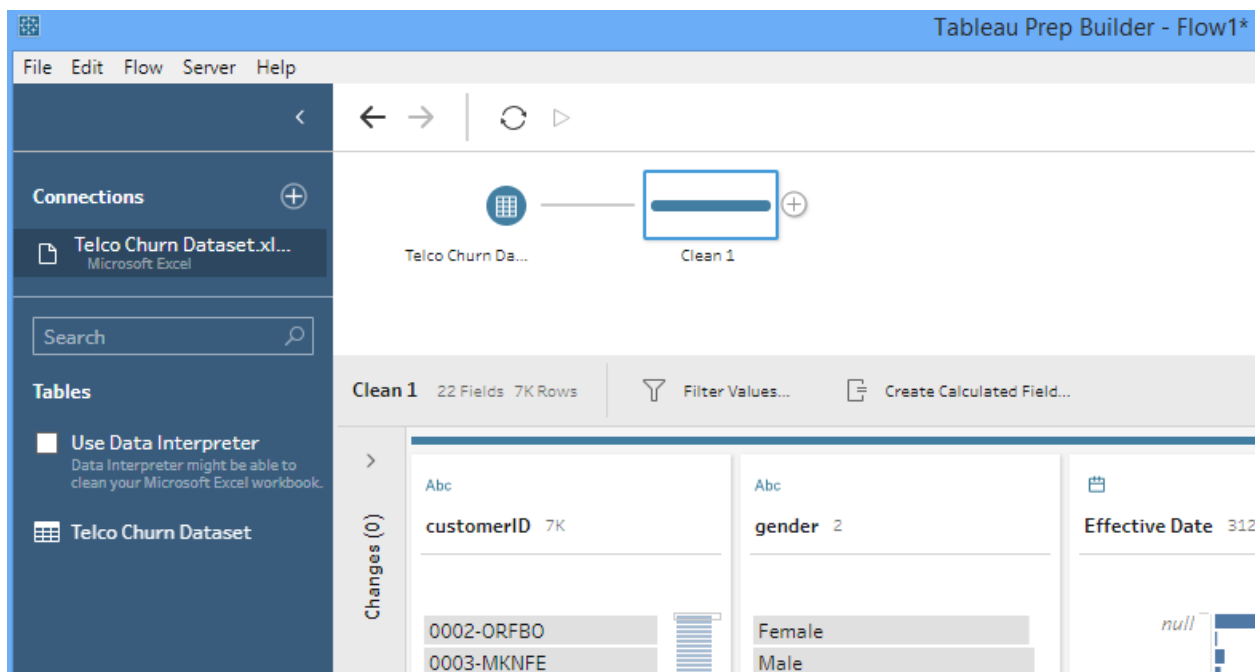
## Tableau Prep Builder – Exerciții de curățare a datelor

Tableau Prep Builder oferă funcționalități în pregătirea datelor pentru analiză, ajutând la creșterea performanței și a optimizării extragerii datelor.

Vom utiliza pentru exemplificare setul de date *Telco Churn Dataset.xlsx*.

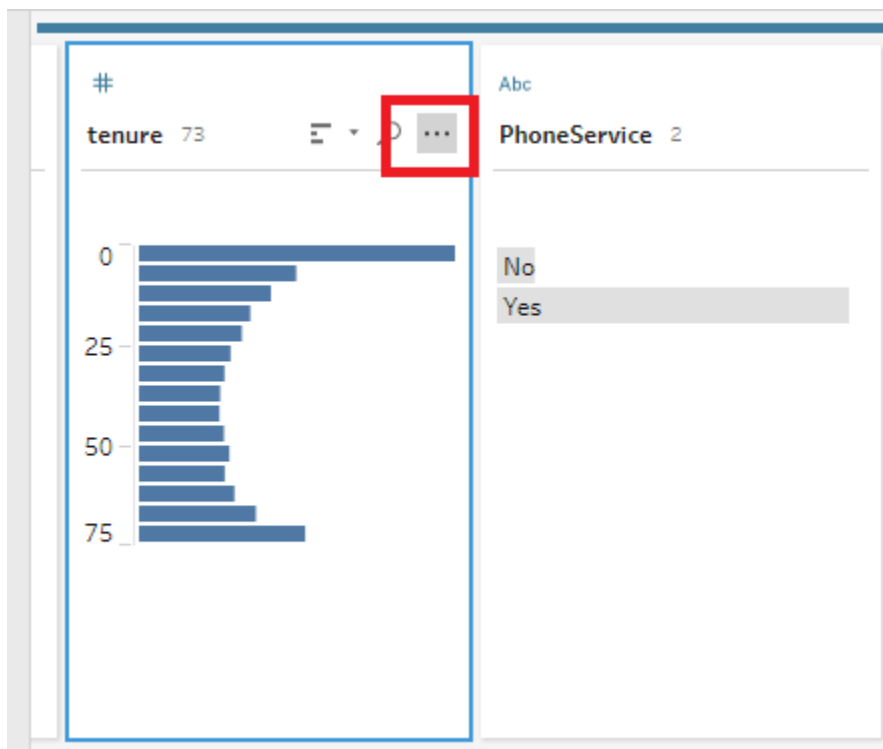
Setul de date folosit se referă la una din problemele comune în domeniul telecomunicațiilor, respectiv migrarea (portarea) clienților. Conform mai multor studii, a fost demonstrat că un client nou costă mult mai mult decât păstrarea unui vechi. Așadar companiile sunt interesate în dezvoltarea unor modele de identificare a clienților cu probabilitate ridicată să renunțe la servicii.

Înainte de a realiza operații în Tableau Prep, este recomandat să vizualizăm datele, respectiv structura lor și să înțelegem semnificația variabilelor. După ce conectăm setul de date vom adăuga un pas “Clean” pentru a avea posibilitatea de vizualizare a structurii datelor.



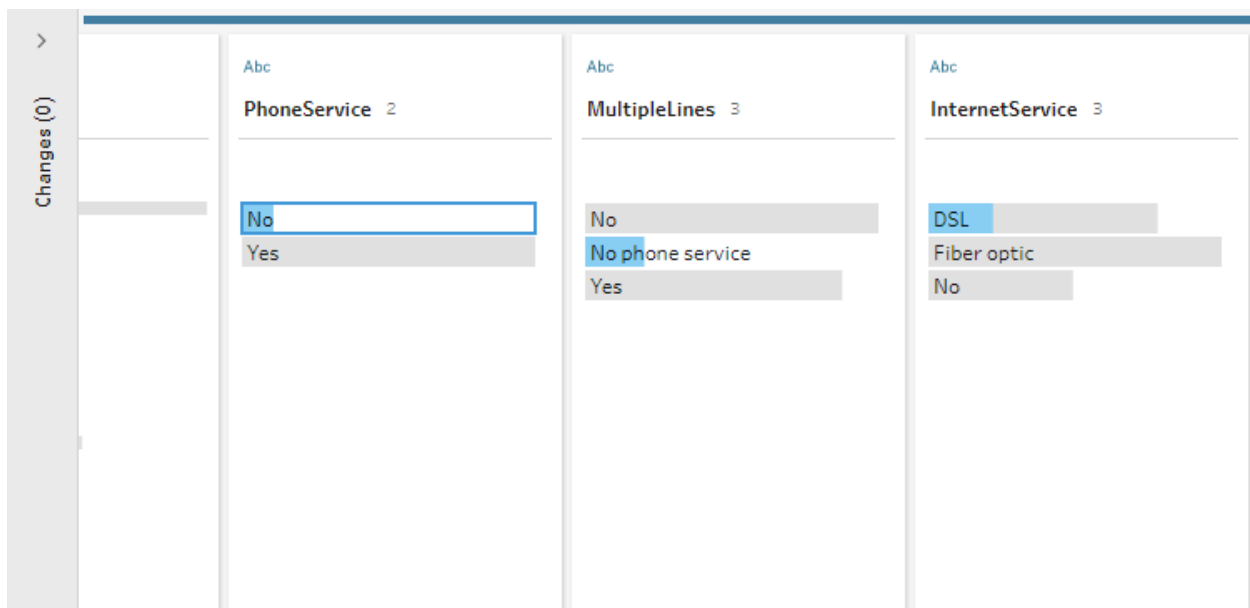
Observați modul de reprezentare a variabilelor discrete din setul de date (Ex. PhoneService), în toate valorile sunt afișate (No/Yes), iar frecvența lor este reprezentată de mărimea barelor orizontale gri.

Variabilele continue (Ex. tenure) sunt reprezentate sub forma unor grafice de tip histogramă, unde valorile sunt grupate pe intervale, iar distribuția lor este reprezentată prin bare orizontale albastre.



Acest format de vizualizare a variabilelor continue poate fi modificat prin selectarea butonului More Options si apoi a optiunii Detail, caz în care vor fi afișate toate valorile, similar cu variabilele discrete.

Afișarea din acest panel ne poate ajuta să observăm imediat anumite informații din setul de date, prin selectarea unor categorii. De exemplu selectând valoarea No pentru variabila PhoneService, putem observa că toate aceste persoane au InternetService de tip DSL.

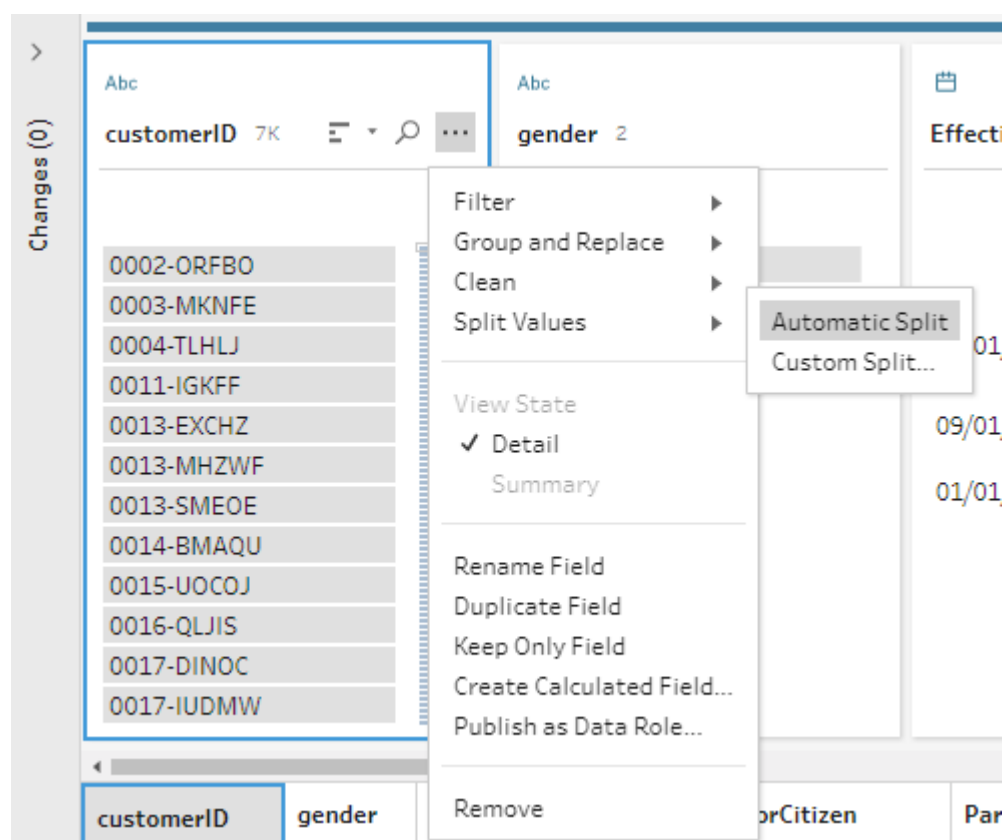


## Funcții

Avem posibilitatea de a realiza operații de curățare a datelor în urma vizualizării lor, incluzând Fitrare, Splitare, Redenumirea câmpurilor, Eliminarea câmpurilor, Schimbarea tipului datelor etc.

### 1) Splitare

În general în seturile de date există situații în care un câmp conține mai multe informații (ex. Data calendaristică, CNP). Aceste informații pot fi extrase sub forma unor câmpuri noi atunci când considerăm că ne sunt utile. În acest set de date, câmpul `customerID` este format din două părți, separate prin caracterul "-". Prima parte reprezintă un cod intern, iar a doua codul de client. Putem realiza splitarea în mod automat prin intermediu funcției `Split Values`.



Vizualizați rezultatul splitării și diferența dintre cele două câmpuri nou create.

### 2) Redenumirea câmpurilor

Uneori câmpurile trebuie redenumite pentru a avea un înțeles mai bun pentru utilizatori. Modificați denumirile celor două câmpuri create mai sus prin splitare, folosind opțiunea `Rename Field`.

### 3) Eliminarea câmpurilor

Pentru a elimina anumite câmpuri ne putem folosi de opțiunea Remove. Eliminarea câmpurilor este utilă pentru a crea un set de date mai simplu, cu mai puțină memorie ocupată inutil.

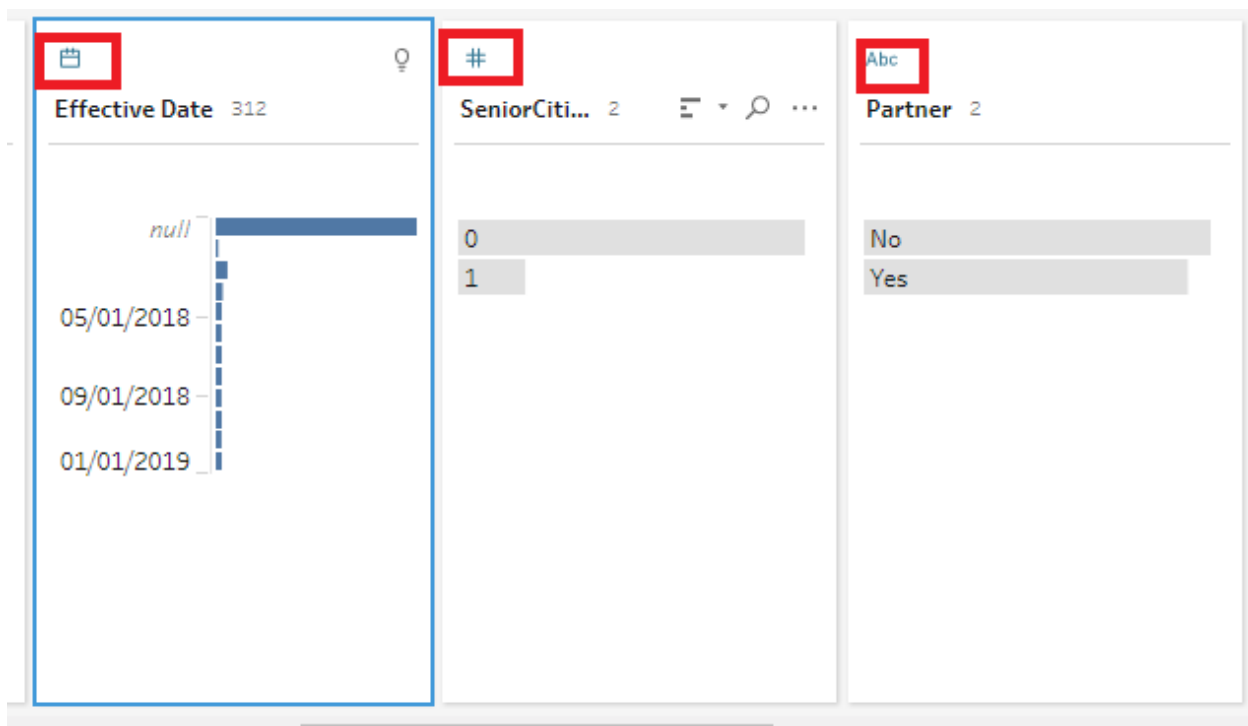
Eliminați câmpul customerID din exemplul de mai sus.

### 4) Schimbarea tipurilor datelor

Tableau Prep are definite următoarele tipuri de date:

- Number(decimal)
- Number(whole)
- Date and time
- Date
- String

La conectarea unei surse de date, tipurile sunt alocate în mod automat în funcție de conținutul fiecărui câmp. În situația în care alocarea automată a tipului nu este satisfăcătoare, avem posibilitatea de a modifica în mod manual tipul. Aceasta operațiune se face apăsând pe pictograma din stanga sus a câmpului.



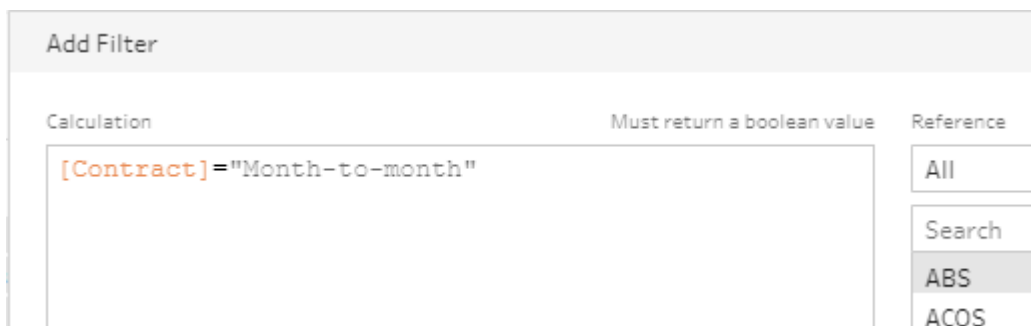
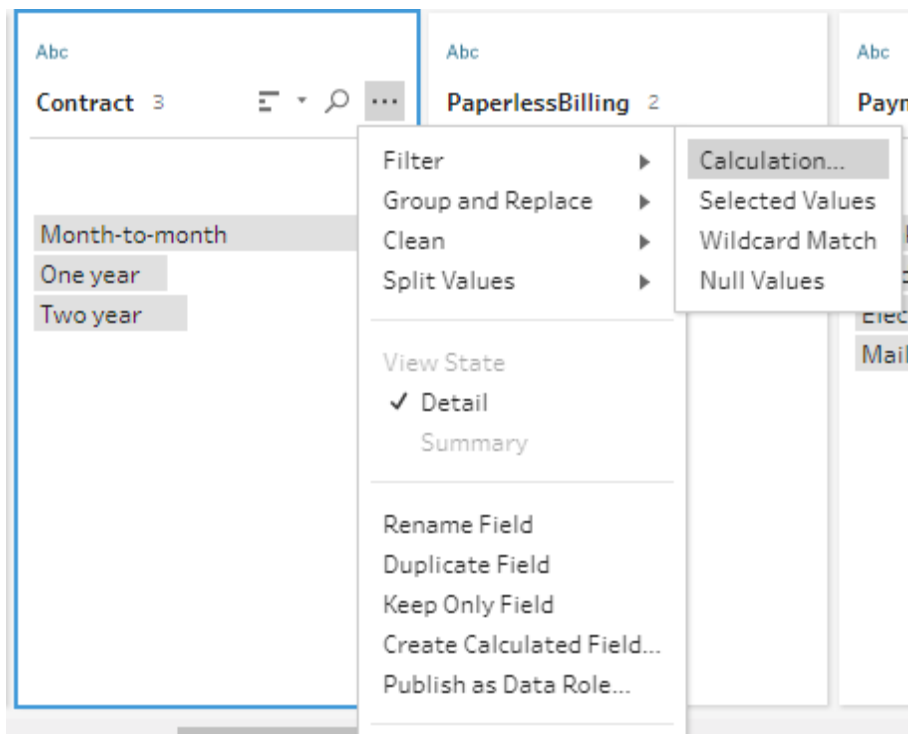
### 5) Filtrarea datelor

Filtrarea datelor ajută utilizatorul în a elimina informațiile care nu îi sunt necesare și a le păstra numai pe cele relevante.

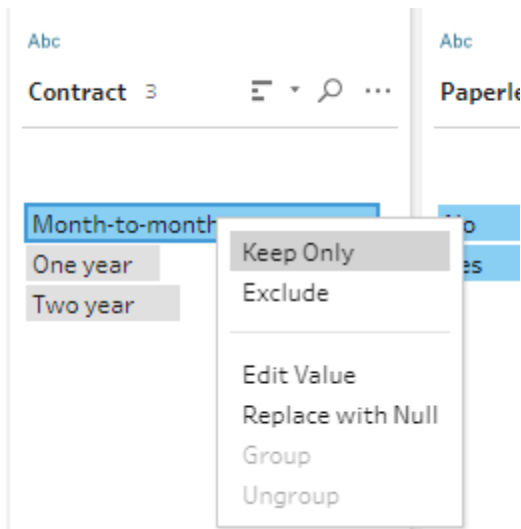
În exemplul folosit în acest laborator presupunem că ne dorim să identificăm în baza de date numai un grup de clienți care respectă anumite criterii. Aceste criterii sunt:

- Sa fie un client cu contract semnat in ultimele doua luni (din setul de date)
- Sa aiba un contract de tip month to month
- Sa genereze venituri mai mici de 80\$ pe luna

Filtrele pot fi adaugate pentru fiecare câmp sub formă de calculation. Spre exemplu, tipul contractului month to month.

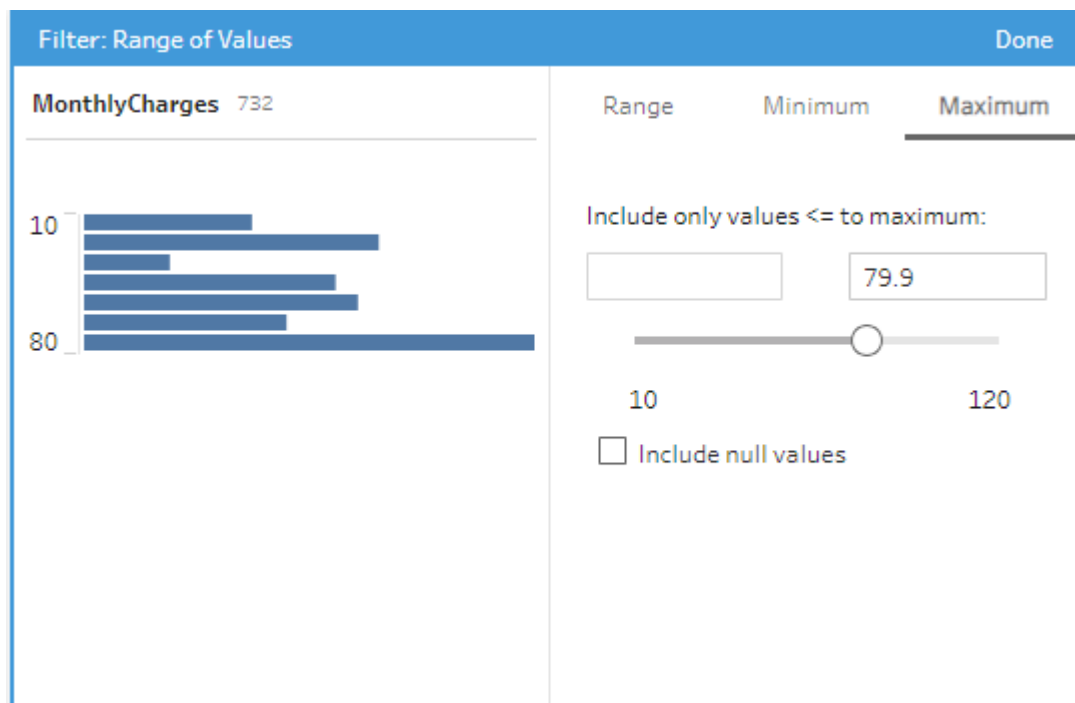


Un mod similar de a realiza același lucru este prin a selecta valoarea dorită pentru câmpul respectiv și apoi a selecta Keep only.



Realizati filtrarea pentru conditia ca MonthlyCharges sa fie mai mici de 80\$ folosind Calculation, ca in exemplul de mai sus.

O alta modalitate pentru aceasta filtrare este aceea de a utiliza Range of Values si a selecta in mod vizual valorile dorite.



Valorile nule din setul de date pot fi sau nu incluse în rezultate, în funcție de decizia utilizatorului. Spre exemplu, ar trebui să ne întrebăm dacă în această situație valorile nule reprezintă clienți care nu aduc venituri sau clienți pentru care nu se cunoaște venitul.

Pentru a filtra in funcție de Data contractului, respectiv ultimele două luni din setul de date, vom folosi filtrarea pentru câmpul Effective Date si vom selecta Range of Dates.

Filter: Range of Dates Done

**Effective Date** 51

10/29/2018

12/03/2018

01/07/2019

**Range** Minimum Maximum

Include only values in this range:

11/01/2018 01/01/2019

02/01/2018 01/01/2019

☐ Include null values

### Vizualizarea modificărilor realizate

Vizualizați modificările realizate în panelul Changes. Dacă dorim să anulăm unele dintre modificările făcute putem să le facem din această interfață.

**Atentionare: Ordinea modificărilor făcute este aceea indicată de sus în jos, adică rezultatele pot fi diferite în funcție de aceasta.**



**Clean 1** 23 Fields 112 Rows Filter Values...

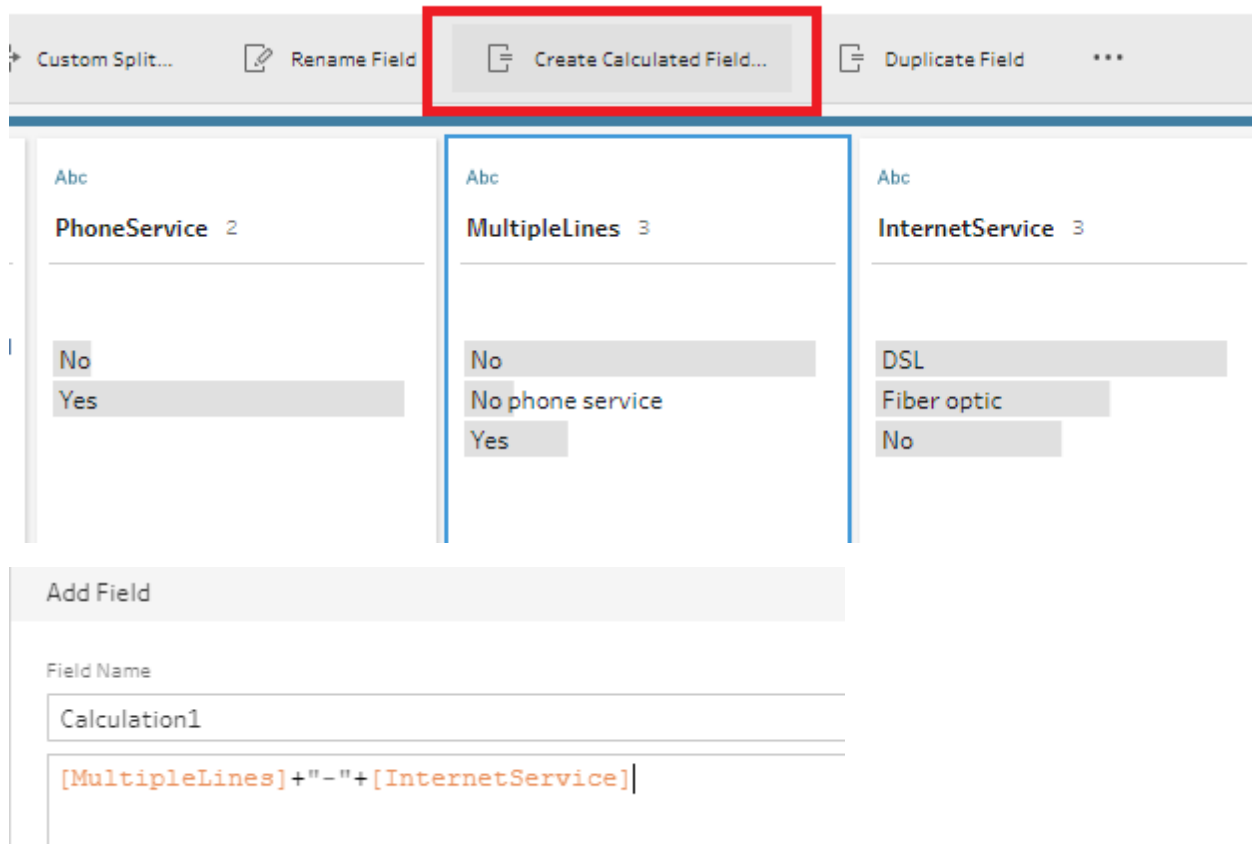
**Changes (8)**

- Calculated Field**  
[customerID - Split 1]  
INT(SPLIT([customerID], "-", 1))
- Calculated Field**  
[customerID - Split 2]  
TRIM(SPLIT([customerID], "-", 2))
- Remove Field**  
[customerID]
- Rename Field**  
[ID client]  
From [customerID - Split 1] to [ID client]
- Rename Field**  
[Cod client]  
From [customerID - Split 2] to [Cod client]
- Filter**  
[Contract]  
Keep-only: "Month-to-month"
- Filter: Range of Values**  
[MonthlyCharges]  
Keep only: values ≤ 79.9
- Filter: Range of Dates**  
[Effective Date]  
Keep only: 11/01/2018-01/01/2019

## Introducerea unor campuri calculate

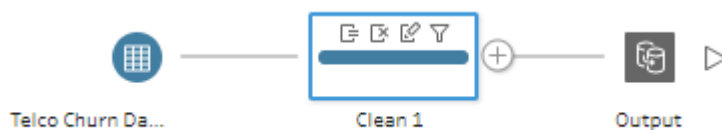
Este util uneori sa introducem noi câmpuri în setul de date, calculate pe baza altora existente. Să presupunem că dorim să introducem un nou câmp care să combine câmpurile MultipleLines și InternetService, astfel încât să vedem în același câmp dacă clienții utilizează cele două tipuri de servicii. Pentru aceasta vom utiliza opțiunea Create Calculated Field si vom realiza o operațiune de concatenare a valorilor celor două câmpuri.





Observați modul în care noul câmp creat a concatenate valorile celorlalte două.

După parcurgerea acestor operații adăugați un nou pas de Output pentru a putea salva datele în formatul dorit.



Observați deasupra pasului Clean anumite pictograme care reprezintă modificările realizate asupra setului de date în pașii descriși mai sus. Pentru a vizualiza oricare dintre modificări putem să selectăm pictograma dorită și vor fi afișate.

## Exerciții

Folosind setul de date Telco Churn Dataset si prelucrarile realizate in cadrul laboratorului, parcurgeți următoarele cerințe:

- a) Anulați filtrarea realizată pentru a selecta clientii cu contract Month to month
- b) Extrageți clienții care utilizează Fibra Optică si nu folosesc Paperless Billing
- c) Adaugati un pas de tip Aggregate pentru a observa media sumelor platite de clientii care s-au portat si de cei rămași (câmpurile Monthly Charges și Churn)
- d) Copiați și Adaugati acelasi pas Aggregate în paralel cu pasul Clean de la punctual b si comparati valorile pentru intreg setul de date cu valorile de la punctual c.
- e) Adaugati in cadrul pasului de la punctul d si campul Total Charges, pentru care sa afisati media. Ce puteti spune despre proporțiile celor doua categorii(Yes și No)?
- f) Folosind un alt pas Aggregate calculați media lunară plătită de bărbați si femei.
- g) Exportați fiecare din rezultatele de mai sus ale pasilor Aggregate in format Tableau Data Extract (.hyper) si importați-le în Tableau Desktop (4 fisiere).
- h) Realizați câte un grafic, la alegere, folosind datele importate.