# Prepare, Secure, and Publish

*by Mark Rittman*

**Automate ingesting, profiling, and transforming data so you can use it in big data cloud and other environments.**

If you are working on an initiative involving big data, analytics, or data warehousing in your organization, you know that adding new datasources and making data suitable for use by analysts can be complicated, largely manual tasks. In the past, these tasks required technical skills and would have been carried out by the IT department. Now, with Oracle Big Data Preparation Cloud Service, available in Oracle Cloud, business users and data domain experts can perform these tasks by using just a web browser.

Oracle Big Data Preparation Cloud Service automates many of the tasks involved in ingesting, profiling, and transforming data, so you can use it in big data and other environments. Oracle Big Data Preparation Cloud Service uses machine learning and natural language processing to detect common data patterns and to alert you to sensitive data items (such as credit card numbers) so that you can easily obfuscate them. You can apply transformations and enrichments by using a point-and-click graphical interface. All you need to do then is apply your knowledge of the data domain you're working in to confirm the recommendations the tool makes and to add your own enhancements before publishing the data for use downstream in other applications.
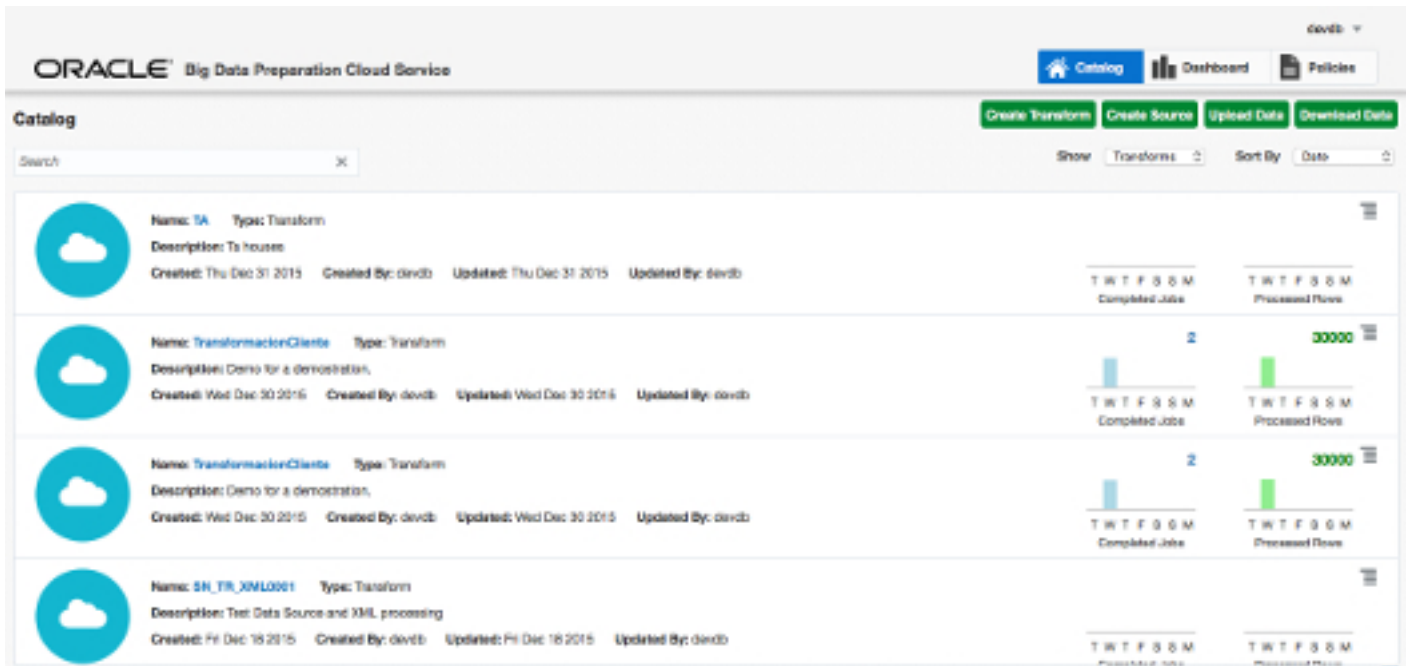
For the example in this article, I'll take a spreadsheet file of customer data that needs to be prepared for use before being made available for other users in an organization. As is often the case, the file comes with no information about the data it contains, and my job is to properly name the columns in the file, make sure all sensitive data is obfuscated, and add whatever enrichments I can to make the contents more useful downstream.

To follow along with the example, request a trial subscription for Oracle Big Data Preparation Cloud Service on Oracle Cloud and download the example file.

## Uploading the File and Profiling Its Contents

Let's start by logging in to Oracle Big Data Preparation Cloud Service, uploading my data file, and then profiling the data so it is ready for transformation.

1. Using your web browser, navigate to the URL given to you for your instance of Oracle Big Data Preparation Cloud Service and log in with the username, password, and identity domain provided to you. After successful authentication, you will see the Catalog page, shown in **Figure 1**, which lists the transforms and sources created by you and others in this cloud environment.

**Figure 1:** Catalog page

**2.** Unzip the o62ba-2867973.zip file. Click the **Upload Data** button to begin the process of uploading the unzipped spreadsheet file to Oracle Storage Cloud Service. When prompted, select **External Client** as the source type.

**3.** Use the **Browse** button to select the file to upload, and then in the Oracle Storage Cloud Service directory browser, select a directory, such as **/Data**, in which to save the file. When you are done, click the **Upload File** button to upload your spreadsheet file to the cloud environment.
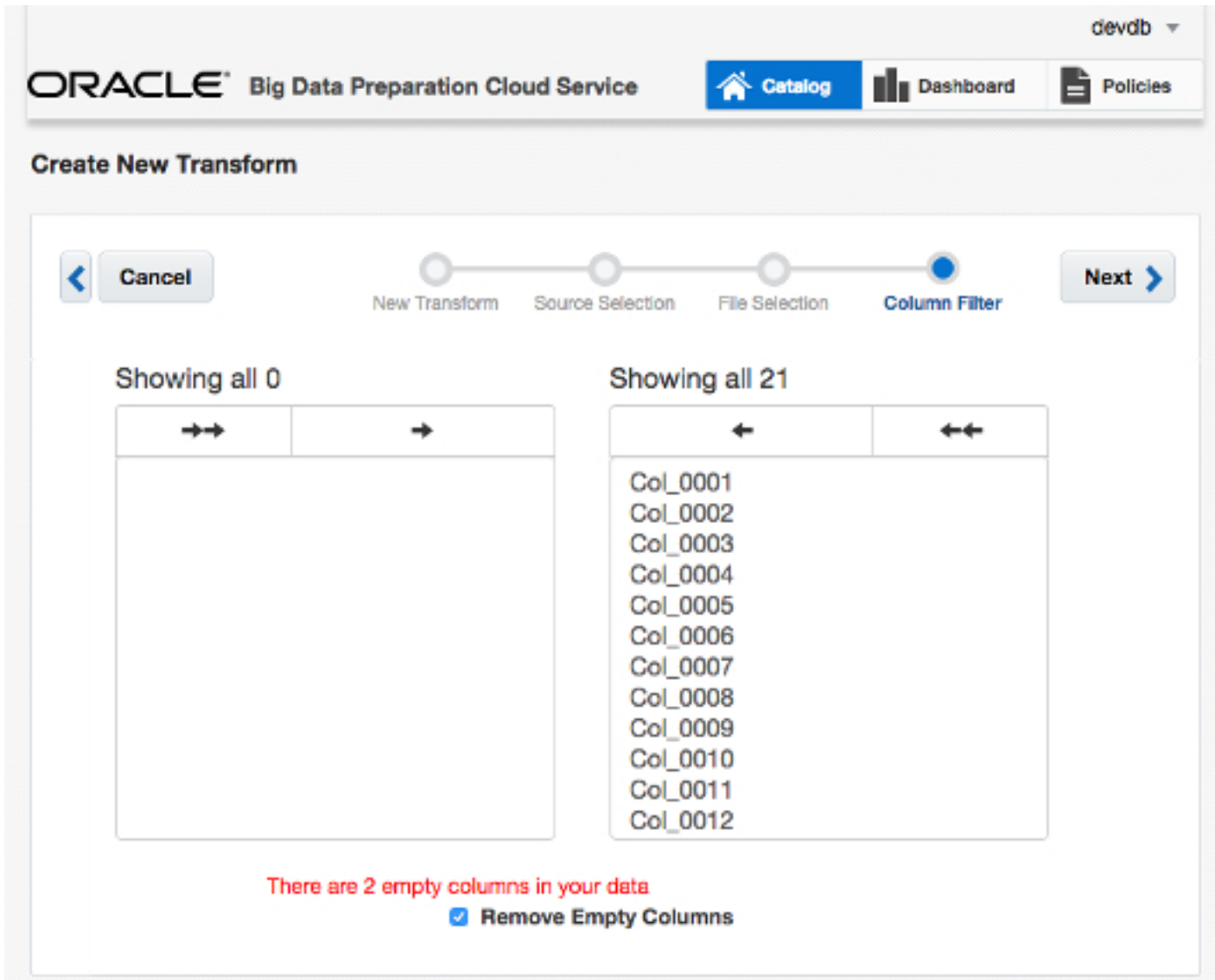
**4.** Navigate back to the Catalog page in Oracle Big Data Preparation Cloud Service, and click the **Create Transform** button to start the process that first profiles the data and then enables you to transform the data. When prompted, enter the following details for the name and description of the transform, and then click **Next** to proceed.

**New Transform Name :** `Customer_Info_Transform`
**Description :** `Transformation to prepare customer file and remove`
`sensitive details`

**5.** Now choose the Oracle Storage Cloud Service source you created earlier to store your data file, and then use the directory browser and file picker to locate and select the file you uploaded. Leave the **Contains Header** checkbox deselected and the **Smart Sample** checkbox selected, so that the profiling process will know that the first row in your file isn't column names and so it will use only a sample of the file's data to determine column contents rather than reading all the contents. Click **Next** to proceed.

**6.** You will then be presented with the list of columns shown in **Figure 2**. Note that at this point, the column names are generic, default values (because in Step 5, you indicated that the file didn't contain a header row that could be used to provide column names). However, you can set the column names later (in most cases, automatically, based on recommendations the tool will provide).
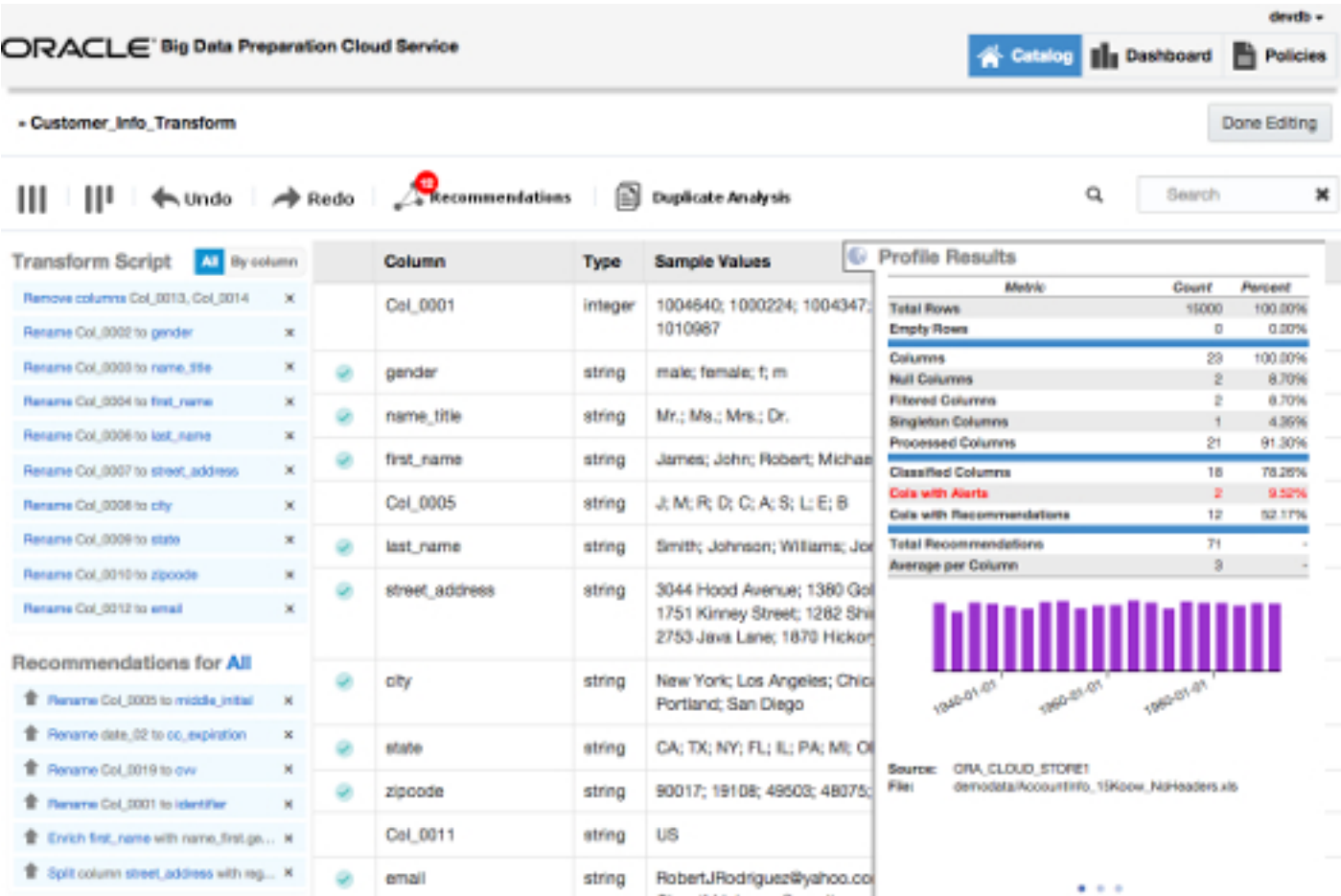
**Figure 2:** Create New Transform Page

**7.** For now, leave the **Remove Empty Columns** checkbox selected so the next stage of the process will consider only columns that contain data. Then click Next to proceed to the main transform authoring page.

# Reviewing, Transforming, and Obfuscating Data in Your File

**Figure 3** shows the main transform authoring page you'll now use to review the data and metadata that was determined during the initial profiling process. Then you'll specify transformations and other "data wrangling" actions to apply to the data set before you make it available to users in your big data environment.



**Figure 3:** Transform Authoring Page

Transformations that will be applied automatically to columns in the data set are listed at the top of the **Transform Script** panel at the left side of the page. Typically the transformations include renaming columns for which there is a clear match with known data patterns.

Recommendations—which are transformations you might want to apply but that are not run automatically, because the pattern match is weak or the transformation alters data values—are shown in the **Recommendations for All** panel at the bottom left of the page.

Now let's start working through the transformations you might want to apply. First you'll look at the highest-priority items on the page, which are the columns associated with the red Cols with Alerts metric in the **Profile Results** panel at the right. Then you'll apply one of the recommended transformations listed in the **Recommendations for All** panel.

> **1.** In the **Profile Results** panel, click the **2** next to the Cols with Alerts metric. Clicking the count next to a metric filters the list of columns displayed in the center of the page. So, in this case, only the **creditcard** and **us_ssn** columns (automatically renamed for you) will be shown.

**2.** To obfuscate the values in the **us_ssn** column so that no sensitive customer data makes its way into downstream applications or analysis environments, start by hovering your cursor next to the us_ssn column's name, clicking the menu icon next to it, and selecting **Obfuscate** from the menu, as shown in **Figure 4**.
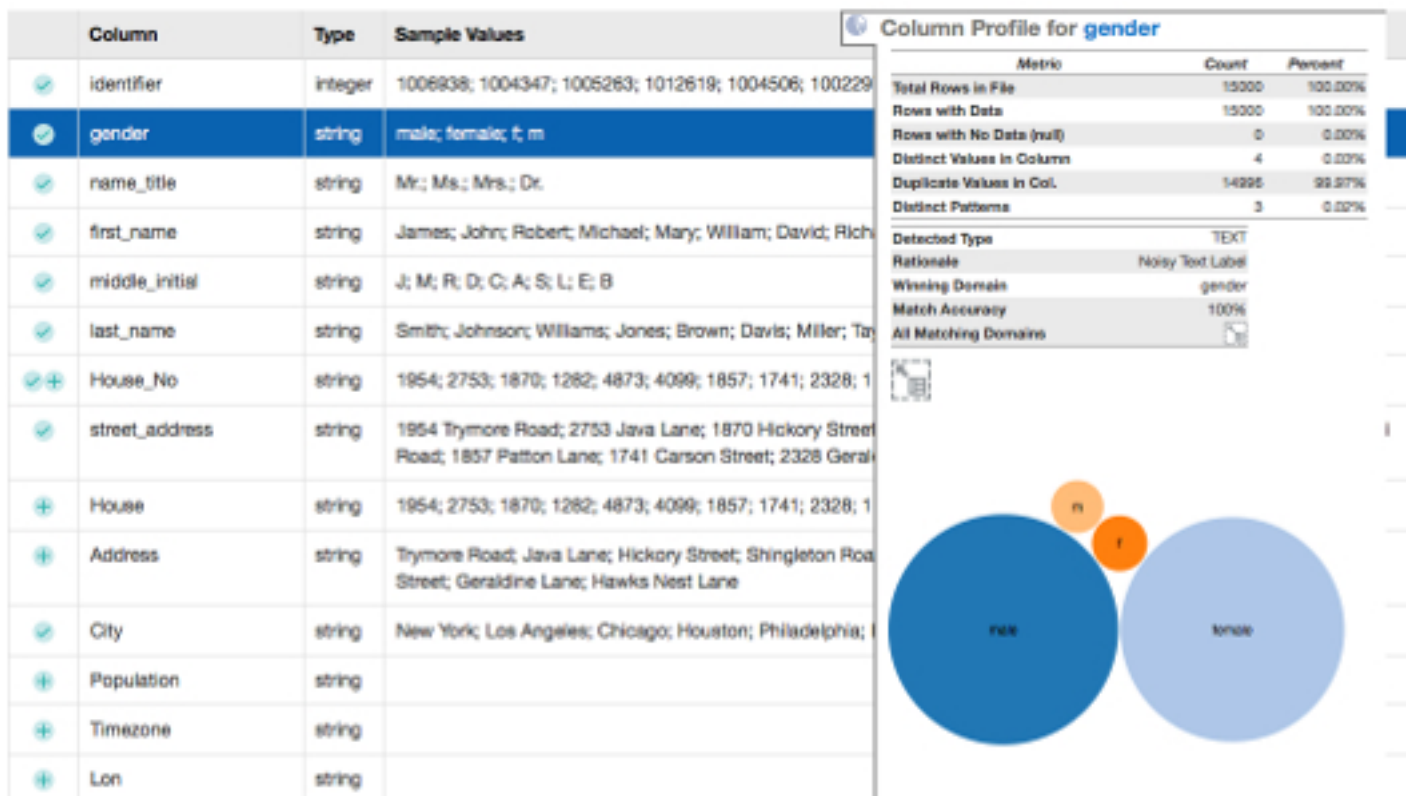


| | Column | Type | Sample Values |
|---|---|---|---|
| | creditcard | string | ****************** |
| | us_ssn | string | ###-##-5702; ###-##-2467; ###-##-5366; ###-##-8700 ##-1600; ###-##-6328 |

Rename

Extract Entities

Split on Delimiter

Extract with Expression

Obfuscate

Date

Change Case

Table Replace

**Figure 4:** Obfuscate Menu Item

**3.** Now, repeat Step 2 to obfuscate the values in the **creditcard** column.

**4.** If you now check the list of transformations listed in the **Transform Script** panel, you should see the two obfuscation transformations added to the list of transformations that were automatically created by the profiling process. When you're done with the obfuscations, click the **Metadata View** button in the top left corner of the page to redisplay the full list of columns and their accompanying metadata items. Now you are ready to set up another transformation—this time, to standardize values in the gender column.

**5.** Locate the gender column in the list, and click to select it, as shown in **Figure 5**. Note that "male" and "female" are domain values but that there are some incidences of "m" and "f" values that were probably entered into the data file in error. These values are likely to cause problems for anyone looking to analyze the values in that column.

| | Column | Type | Sample Values |
|---|---|---|---|
| ✔ | identifier | integer | 1006938; 1004347; 1005263; 1012619; 1004506; 100229 |
| ✔ | gender | string | male; female; f; m |
| ✔ | name_title | string | Mr.; Ms.; Mrs.; Dr. |
| ✔ | first_name | string | James; John; Robert; Michael; Mary; William; David; Rich |
| ✔ | middle_initial | string | J; M; R; D; C; A; S; L; E; B |
| ✔ | last_name | string | Smith; Johnson; Williams; Jones; Brown; Davis; Miller; Tay |
| ✔➕ | House_No | string | 1954; 2753; 1870; 1282; 4873; 4099; 1857; 1741; 2328; 1 |
| ✔ | street_address | string | 1954 Trymore Road; 2753 Java Lane; 1870 Hickory Street Road; 1857 Patton Lane; 1741 Carson Street; 2328 Geral |
| ➕ | House | string | 1954; 2753; 1870; 1282; 4873; 4099; 1857; 1741; 2328; 1 |
| ➕ | Address | string | Trymore Road; Java Lane; Hickory Street; Shingleton Roa Street; Geraldine Lane; Hawks Nest Lane |
| ✔ | City | string | New York; Los Angeles; Chicago; Houston; Philadelphia; |
| ➕ | Population | string | |
| ➕ | Timezone | string | |
| ➕ | Lon | string | |

**Column Profile for gender**

| Metric | Count | Percent |
|---|---|---|
| Total Rows in File | 15000 | 100.00% |
| Rows with Data | 15000 | 100.00% |
| Rows with No Data (null) | 0 | 0.00% |
| Distinct Values in Column | 4 | 0.03% |
| Duplicate Values in Col. | 14996 | 99.97% |
| Distinct Patterns | 3 | 0.02% |

| | |
|---|---|
| Detected Type | TEXT |
| Rationale | Noisy Text Label |
| Winning Domain | gender |
| Match Accuracy | 100% |
| All Matching Domains | |

**Figure 5:** Values in the Gender Column

**6.** To begin standardizing the values in the gender column so that only "male" and "female" are used, click the menu icon next to the column's name and choose Table Replace from the menu that is displayed.

**7.** Then, in the Table Replace dialog box, first click the **Samples** button to show the set of column values. Then remove the "male" and "female" values by clicking the **x** next to them. To replace the "f" and "m" values that remain, type female and male, respectively, into the **Replace With This** column, as shown in **Figure 6**. Then click **Apply** to close the dialog box and return to the transformation editor screen.

| Table Replace - gender | | ✖ |
| --- | --- | --- |
| **Find This** | **Replace With This** | |
| f | female | ✖ |
| m | m | ✖ |
| | | ✖ |

Import CSV  Samples  Apply

Import CSV  Samples  Apply

**Figure 6:** Table Replace Dialog Box

**8.** As the final transformation, you'll use some of the enrichment recommendations to add population, latitude, longitude, and other data discovered earlier. To add the enrichment recommendations shown in **Figure 7** to the file you publish, locate the recommendations in the **Recommendations for All** panel and then click the up arrow next to each recommendation you want to add to the main set of transformations for this file.

## Recommendations for All

| | |
|---|---|
| ⬆ Enrich City with City.Lat | ✕ |
| ⬆ Enrich City with City.Lon | ✕ |
| ⬆ Enrich City with City.Country | ✕ |
| ⬆ Enrich City with City.Province | ✕ |
| ⬆ Enrich City with City.Jurisdiction | ✕ |
| ⬆ Enrich City with City.Population | ✕ |
| ⬆ Enrich City with City.Elevation_Meters | ✕ |
| ⬆ Enrich City with City.Timezone | ✕ |
| ⬆ Enrich City with City.Feature_Type | ✕ |

**Figure 7:** Enrichment Recommendations

Once you're finished, click the **Done Editing** button in the top-right corner of the screen to close this page and return to the Catalog page.

## Publishing and Accessing the Prepared Data File

Now that you have defined the transformations, enrichments, and column obfuscations for this data set, you can publish the prepared data file back to Oracle Storage Cloud Service for later download to your workstation or you can transfer the prepared file to another environment. To publish your prepared file, take the following steps:

**1.** On the Catalog page in Oracle Big Data Preparation Cloud Service, navigate to the transformation you created a moment ago and click the context menu icon to the right of it. When the menu is displayed, select **Publish** from the listed options.

**2.** On the Publish page, select **Default_Cloud_Storage** for **Target** and then click the **Publish** button to proceed. You should then see the output from the publish process displayed in your web browser and, when the process is complete, you should see a message indicating that the publish process succeeded.

**3.** To download the published file and view its contents, return to the Catalog page and click the **Download Data** button. On the Select Source Type page, select the Oracle Storage Cloud Service source type and then **Default_Cloud_Storage** as the source. Navigate to the /Publish directory, and locate the published file. (Your username is appended to the name of the transformation file.) Once you have located the file, select it and click the **Download** button to download it from Oracle Cloud to your workstation.

**4.** Finally, open the prepared file with a text editor and check that the prepared file contains not only the data you selected from the original data file but also the columns added by the enrichment process and the standardized gender values. Also check that all the sensitive customer data is obfuscated, as shown in **Figure 8**.



**Figure 7:** Enrichment Recommendations

**Conclusion**

Much of the work involved in big data and other analytic projects involves preparing, obfuscating, and enriching datasources so they are useful for analysis and do not contain customer data that is sensitive or private. Oracle Big Data Preparation Cloud Service puts this capability into the hands of data domain experts and makes the process quick and easy to perform, by automating much of the work and running as a service in Oracle Cloud.