# Custom reporting with Mascot Distiller

**Patrick Emery**

**Matrix Science Ltd**

## Reports in Distiller

- **Mascot Distiller 2.8 ships with 14 reports**
  - ANOVA
  - Hierarchical and K-means clustering
  - PCA
  - Volcano plots
  - Table exports etc
- **Covers many use cases**
  - But not everything!

**MASCOT** : *Distiller custom reports* © 2024 Matrix Science — MATRIX SCIENCE

Mascot Distiller 2.8 ships with 14 different reports for analysing and exporting your quantitation results. These include reports such as Analysis of Variants (ANOVA), clustering reports such a hierarchical, K-means clustering and Principal Component Analysis. Standard reports such as the Volcano plot, box-plot are available and there are table exports which make it easy to get your data out of Distiller and into 3rd party packages such as Excel, R or Perseus from the Max Planck institute.

These cover many use cases, but we can't cover every possible analysis type ourselves and you may have some specific requirement that calls for a custom report.

## Reports in Distiller

- **Reports are written in Python**
- **Ships with embedded Python and many useful libraries**
  - Mascot Parser
  - Pandas
  - *And many more...*
- **Two files for every report**
  - Python source file
  - XML file describing inputs

In older versions of Distiller, reports were written in XSLT – an XML transformation language. This was powerful, but not commonly used and without access to many of the features you might need when creating a report.

In Mascot Distiller 2.8, reports are instead written in Python. Distiller ships with it's own embedded copy of Python with many useful libraries for data analysis.

Each report is comprised of two files – the Python report source file, and an XML file describing the report inputs and defining any Wizard to be displayed in the Distiller GUI.

## Getting started

- **Install you**
  - E.g. Visua
- **Install Dis**
- **Set the ID**
  **Python. I**

**MASCOT** : *Distiller custom reports © 2024 Matrix Science*

With a bit of knowledge of Python and the Mascot Parser library therefore you can write your own reports for use in Distiller. For this, you'll need an Integrated Development Environment (IDE) that supports Python, such as the free Visual Studio Code from Microsoft.

Install Distiller on your development PC. Without a licence, Distiller will run in viewer mode, but this will still install the embedded copy of Python and the various 3rd party libraries. Once you've done that, tell your IDE to use the version of Python installed with Distiller.

Here's an example of how to do that in Visual Studio Code

# Example: Top-3 Protein Intensity

- **Known as "Average" in Mascot**
- **Absolute quantitation method**
  - Similar to e.g. iBAQ
- **Average of the intensities of the 3 most intense peptides matching the protein**
  - Silva *et al*. MCP 5 114-156 (2006)
  - Average MS signal for the three most intense tryptic peptides/mole of protein constant within a CV < ±10%
- **Average method is label-free in Distiller, but could calculate for e.g. TMT channels**

**MASCOT** : *Distiller custom reports* © 2024 Matrix Science   MATRIX SCIENCE

As an example of the steps involved in, and the utility of, creating a custom report, we'll write a report to calculate absolute protein intensity using the top-3 method.

This is known as "Average" label free quantitation in Distiller and is an absolute quantitation method similar to iBAQ

The protein intensity is calculated from the average of the intensities of the 3 most intense peptides matching a protein. The observation from the original paper was that the average MS signal for the three most intense tryptic peptide/mole of protein is constant within Coefficient of Variation of +/- 10%

Average quantitation is label-free in Distiller, but there's no reason you couldn't calculate it for other quantitation techniques.

Example: Top-3 Protein Intensity

- XML file defines parameters and inputs for the report Wizard

Tells Distiller where to put the report in the menu tree

MASCOT : *Distiller custom reports* © 2024 Matrix Science

Once we've setup our IDE, we can start work on the report itself. As mentioned earlier, a report in Mascot Distiller is comprised of two files – the actual Python .py source file for the report and an XML file which defines the report inputs and Wizard for the Distiller GUI.

Starting with the XML file, here, we're telling Distiller where to put the report on the menu tree – in this case in a subfolder called Custom

# Example: Top-3 Protein Intensity

```
<Supports average="true" precursor="true" replicate="true"
reporter="true" multiplex="true"/>
```

Tells Distiller report supports all quantitation methods

And here we're telling Distiller that our report will support all quantitation methods.

## Example: Top-3 Protein Intensity

```
<ReportConfiguration>
    <Parameter name="exportFormat" label="Format" value="CSV"
    type="text" mapsTo="ExportFileType"/>
    <Parameter name="databaseCount" label="no databases"
    type="integer" value="@{DatabaseNames.Count}"/>
</ReportConfiguration>
```

Tells Distiller report will produce a CSV file
Set a variable called databaseCount

**MASCOT** : *Distiller custom reports* © 2024 Matrix Science

MATRIX SCIENCE

This is telling Distiller it will produce a CSV file and sets the number of sequence databases searched to a variable called databaseCount

# Example: Top-3 Protein Intensity

```
<Wizard>
    <WelcomeText>This report generates a CSV file containing
    protein component intensity values calculated using the
    top-3 method</WelcomeText>
```

**Run Top 3 protein intensity report**

**Welcome**

This report generates a CSV file containing protein component intensity values calculated using the top-3 method

Next >    Cancel

…o display a landing page in the report
…e specified text

**MASCOT** : *Distiller custom reports* © 2024 Matrix Science

MATRIX SCIENCE

A separate section in the XML can then be used to define the Wizard to display in the GUI. This can be used to capture input and settings information from the user before running the report.

Here we're defining the landing page for the Wizard – this should define some welcome text to explain to the user running the report what it is going to calculate.

# Example: Top-3 Protein Intensity

```
<Page title="Peptide selection criteria">
    <HelpText>The selection type determines whether the n peptides
    must have different sequences (unique_sequence) or whether to
    accept different modification states of same sequence (unique_mr),
    or                                                    e and modifications
    in                                                    Text>
    <Pa                                                    tion type" type="select">
                                                          tring="Unique sequence"
                                                          "Unique Mr"/>
    </P
</Page>
```

**Run Top 3 protein intensity report** ×

**Peptide selection criteria**

The selection type determines whether the n peptides must have different sequences (unique_sequence) or whether to accept different modification states of same sequence (unique_mr), or even to accept peptides with same sequence and modifications in different charge states (unique_mz)

Selection type | Unique sequence ▾

Unique sequence
Unique Mr
Unique M/Z

zard page with a drop-down
e user chooses the rule for
ntense peptides

< Back | Next > | Cancel

**MASCOT** : *Distiller custom reports* © *2024 Matrix Science*

Here we're defining an actual input page with a parameter.  In this case, we're defining a drop-down selection control where the end user will set the criteria for selecting the top-3 most intense peptides.

The default is Unique Sequence – where different charge and modification states of the same peptide are treated as a single match.  Selecting Unique Mr would mean different modification states are counted as separate peptide matches, while Unique m/z means the different charge and modification states of the same peptide sequence are treated as different matches.

Here's the Wizard page defined by that XML displayed in Distiller.  We can define as many input pages as required to capture the inputs for our report.  Entered values are then output to the Python script as comma separate values.

# Example: Top-3 Protein Intensity

```python
185    # \param isAverage boolean - True if the quantResults are the Average [MD] method
186    # \param isMS1 boolean - True is the quantResults are ms_ms1quantitation results
187    # \param peptideSelection The peptide selection/grouping criteria selected by the user on the report Wizard
188    # \return The calculated intensity, or 0 if the intensity cannot be calculated
189    def calculateProteinIntensity(protein, peptideSummary, quantResults, componentName, isAverage, isMS1, peptideSelectionType):
190        if isAverage:
191            return calculateProteinIntensityAverage(protein, quantResults)
192        elif isMS1:
193            return calculateProteinIntensityMS1(protein, quantResults, componentName, peptideSelectionType)
194        else:
195            return calculateProteinIntensityMS2(protein, peptideSummary, quantResults, componentName, peptideSelectionType)
196
197    ## Uses msparser to calculate the top-3 average protein intensity for the passed ms_protein
198    # \param protein The ms_protein to calculate the average intensity for
199    # \param quantResults The msparser quantitation results
200    # \return The calculated intensity, or 0 if the intensity cannot be calculated
201    def calculateProteinIntensityAverage(protein, quantResults):
202        # we can use msparser to calculate the values
203        include = msparser.ms_peptide_quant_key_vector()
204        exclude = msparser.ms_peptide_quant_key_vector()
205        sampleSize = 0
206        # getAveragePeptideIntensity only works if the quantitation method is Average
207        # In Python this results an array, the first value is True or False depending on whether or not
208        # the calculation was successful (sufficient peptides matching the criteria)
209        # the second value in the intensity
210        intensity = quantResults.getAveragePeptideIntensity(protein.getAccession(), protein.getDB(), include, exclude)
211        if intensity[0] == True:
212            return intensity[1]
213        else:
214            return 0
```

**MASCOT** : *Distiller custom reports* © 2024 Matrix Science

MATRIX SCIENCE

Once we've defined our inputs in the XML, we need to create our Python report file. This is a standard Python script – we're using Mascot Parser to access the search results, and also some useful script libraries we've provided to help with loading and formatting the data. You'll need to be familiar with the Python 3 programming language to create your own reports, and also spend a bit of time getting to know Mascot Parser and any other libraries you want to use.

Here we're looking at a snippet of the code involved in calculating the top-3 protein intensity - a slightly different method being required depending on the quantitation method used in Mascot and Distiller.

It's a bit beyond the scope of a presentation like this to run through all of this, but you can find a detailed tutorial pdf on our public website here:

Example: Top-3 Protein Intensity

C:\ProgramData\Matrix Science\Mascot Distiller\reports

MASCOT : *Distiller custom reports* *© 2024 Matrix Science*

The name of your XML file should be the same as the Python .py report script with .XML appended to the end. When both files are written, copy them into the c:\ProgramData\Matrix Science\Mascot Distiller\reports directory and restart Mascot Distiller – your new report should now appear in the Analysis->Reports menu in Distiller in the subfolder that you've specified in the XML.

Note, the ProgramData folder is hidden in Windows by default, but if you type the path into the navigation bar in Explorer, it will open.

Once you've copied the files into place and restarted Distiller, you can update the Python file without restarting Distiller, but changes to the XML file may require Distiller to be restarted in order to take effect.

**Example: Top-3 Protein Intensity**

- **Example dataset – two files from PXD001385**
- **HeLa background 1:1**
- **Spiked in *E.coli* proteins at:**
  - 3ng, 7.5ng, 10ng,15ng
- **Two files, 15ng & 3ng**
- **Processed with Distiller, search with Mascot 2.8**
- **Quantitation (Replicate LFQ) in Distiller**
  - Enables 'match between runs'

**MASCOT** : *Distiller custom reports* © 2024 Matrix Science

MATRIX SCIENCE

Let's use our new report on an example dataset. We've taken two files from the following dataset in the PRIDE repository. This is a benchmarking dataset, where we have a background of human proteins from a HeLa extract at 1:1, and then a series of samples where E.coli protein extracts were spiked in at 3ng, 7.5ng, 10ng and 15ng. We took two of the files, one of the 15ng and the matched 3ng file, processed them with Mascot Distiler 2.8 and search against the Uniprot Ecoli proteome with Mascot 2.8, before running quantitation in Distiller using the "Replicate" LFQ method (which enabled 'match between runs' type behaviour)

# Example: Top-3 Protein Intensity

Here's our complete Wizard, as defined by the report XML file. When we're ready, click 'Finish' on the final page and Distiller will run the report. When it completes, the generated CSV file will be automatically opened in our spreadsheet application (Excel in this case)

Example: Total Intensity

MASCOT : *Distiller custom reports* © 2024 Matrix Science

And here's our generated report – at the top we have some basic header information about the raw files and search results. Below that is our table of results, with calculated intensity values for the 3ng and 15ng samples. Within Excel you could calculate relative abouts of each protein against any control protein of known amount, or to the total intensity of all proteins etc. In this case, we have values for one of the 3ng Ecoli and 15ng Ecoli samples – so the protein intensities for the 15ng sample should be 5x greater than the intensities from the 3ng sample – obviously, I could have had the report calculate the ratio, but it's very easy to do in Excel.

So, there's a bit of spread (as you'd expect), but overall, those look like pretty good numbers – as confirmed by the Median and Mean ratios for the dataset.

Example: Top-3 Protein Intensity

How does that compare to another technique such as iBAQ? Well, looks like the ratio of the iBAQ values for these proteins aren't bad, but they're actually a bit too high, as shown by the Mean and Median ratios. So, in this case Top-3 has out performed iBAQ and is a simpler calculation to boot. Unlike iBAQ you don't need the protein sequence – you could calculate it yourself very easily in a spreadsheet with just the protein's peptide intensity values.

While top-3 has done better in this instance, that isn't always the case of course. In general both techniques perform similarly, and generally surprisingly well. The main disadvantage of top-3 compared with iBAQ is the requirement for 3 distinct peptide matches to calculate the average intensity from – iBAQ allows you to calculate intensities for proteins with fewer distinct matches – which may be why top-3 out performs iBAQ here; we're only using the most intense peptides in the calculation.

## Conclusions

- **Mascot Distiller ships with 14 reports**
  - Covers many use cases
- **Custom reporting with Python**
  - Commonly used
  - Many powerful data analysis libraries available
- **Fully functional programming language**
  - Download additional data (e.g. IntAct)
  - Call other software
  - Etc etc etc

Mascot Distiller 2.8 ships with 14 standard reports. These cover many use cases, but they can't cover every type of analysis you could ever want to carry out.

Reports are written in Python – this is a commonly used programming language for data analysis and it has many powerful libraries available to help with that.

Unlike in earlier versions of Distiller, we're using a fully function programming language for reporting. That means you could get very creative with your reports – for example, you could download protein interaction data from the EBI's IntAct database and incorporate that into a report, or you could call another tool, carry out some downstream analysis and then incorporate that into a report.

## Finally…

**Tutorial: Creating custom reports in Mascot Distiller**
**https://www.matrixscience.com/blog/tutorial-creating-custom-reports-in-mascot-distiller.html**

**Includes links to a tutorial PDF and zip file with the report source code and XML**

Creating_a_custom_Distiller_report.pdf        top-3.py.zip

Finally, you can find a blog article about this here

A detailed tutorial pdf about creating the top-3 report here

And the Python source code and XML file for the report can be downloaded from here