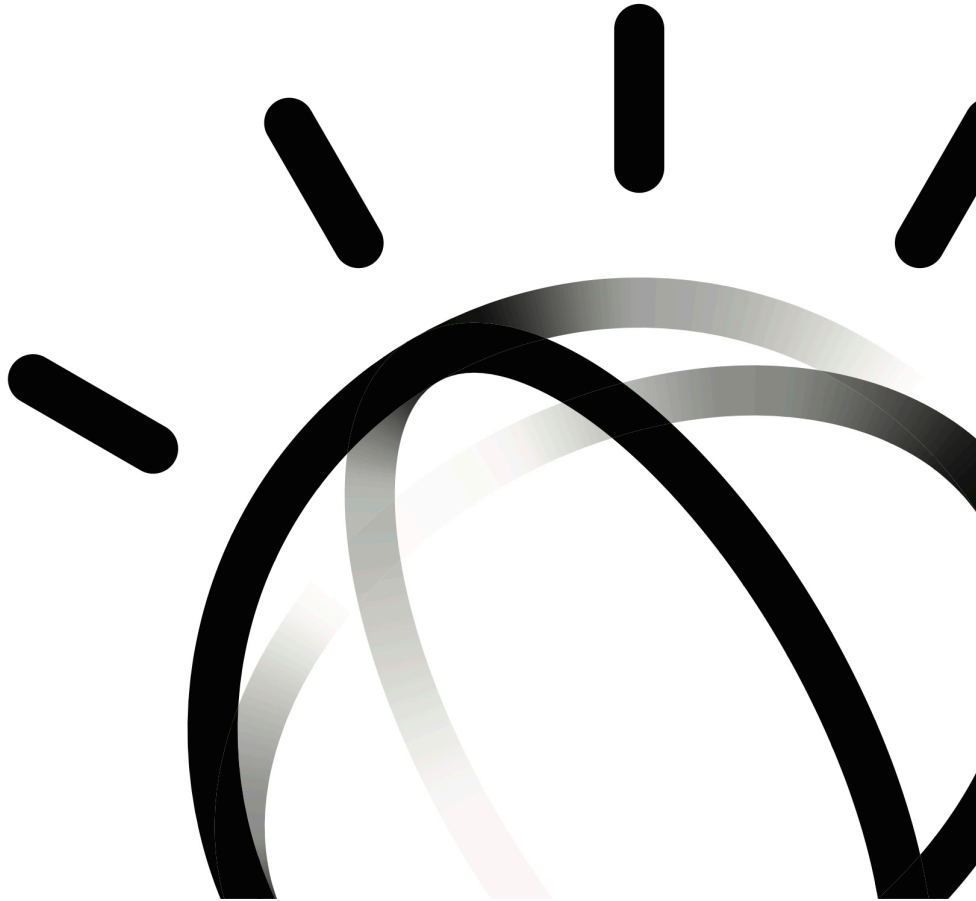


# **IBM Watson Solutions**

## **Business and Academic Partners**



## **Ingest, Convert, Enrich and Query with Watson Discovery Service**

**Prepared by Armen Pischdotchian**

**Updated by Hélène Fourot Quillaud**

**Version 3.1 August 2018**

# Overview

With the IBM Watson™ Discovery service you build cognitive, cloud-based exploration applications that unlock actionable insights hidden in unstructured data - including your own proprietary data, as well as public and third party data. Creating your first discovery journey using the IBM Watson™ Discovery service entails the following steps:

- Convert, enrich and normalize data.
- Securely explore your proprietary content as well as free and licensed public content.
- Apply additional enrichments such as concepts, relations, and sentiment through natural language processing.
- Query and analyze your results.
- Simplify development while providing direct access to APIs.

IBM Watson Discovery service architecture is depicted below.



You can upload content and begin finding insights with the Discovery service by using either the Discovery Tooling or the Discovery API. This document shows you how to use both.

You might want to watch this and related videos from the playlist:

<https://www.youtube.com/watch?v=fmIPeopG-ys&t=1s>

This workshop will start by guiding you on how to upload and query documents. Then, it will show you how to configure your sample document with conversions and enrichments, so that you can leverage custom configurations for better insights

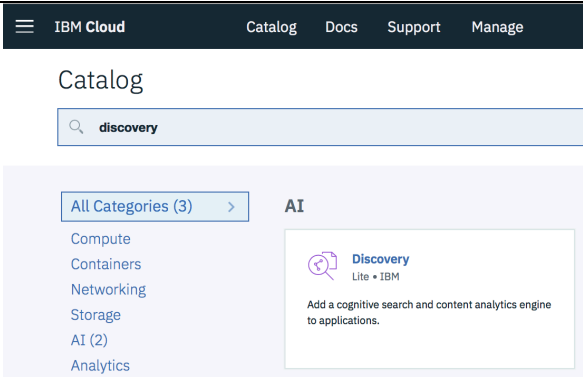
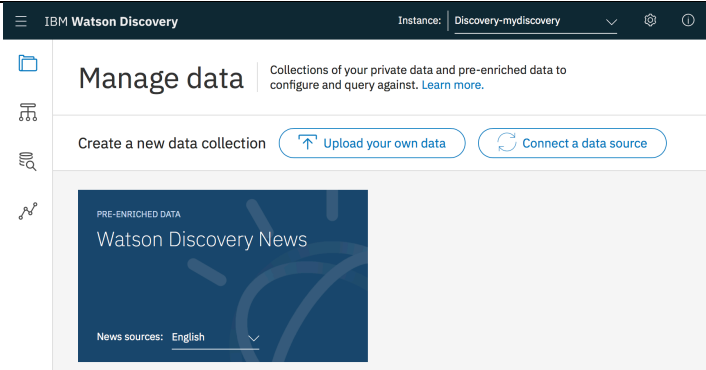
At the end of this workshop, you may want to spend some time and consider coupling the Discovery service with the Watson Assistant service and think about how you could use the Watson Knowledge Studio (a SaaS offering, not on IBM Cloud) to further edit the annotators used with some of the cognitive language microservices used within the Discovery service.

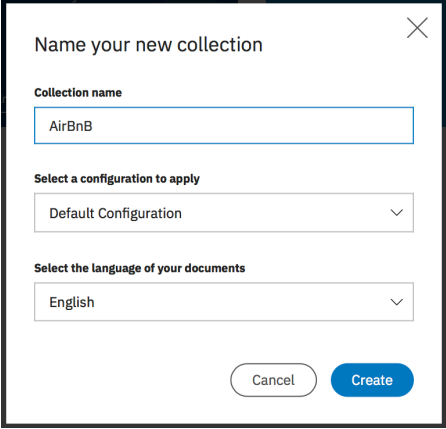
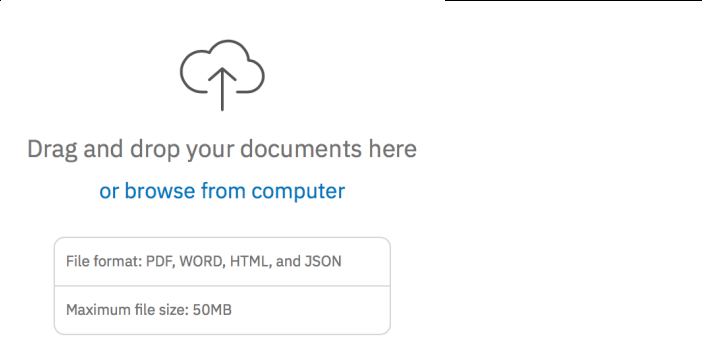
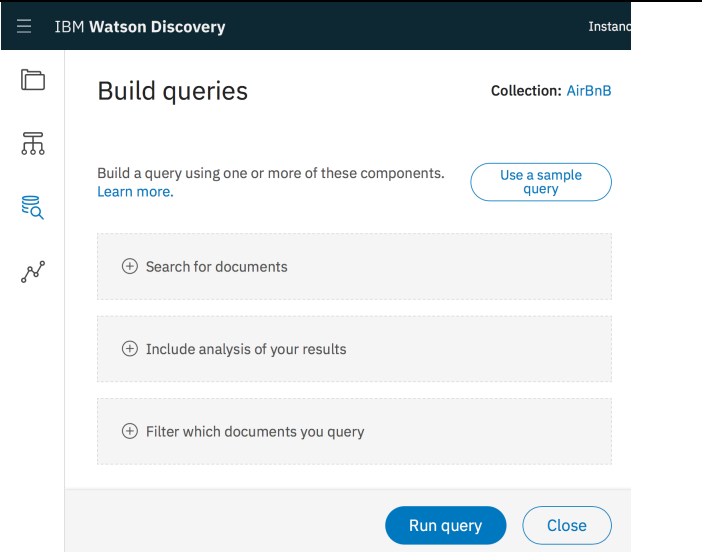
## Prerequisite

Prior to starting this lab, you must obtain IBM Cloud credentials.

# Working with the Discovery Tooling

## Uploading and querying documents

Steps	Example screen capture
<ol style="list-style-type: none"> <li>1. Login into IBM Cloud at <a href="https://console.ng.bluemix.net">https://console.ng.bluemix.net</a></li> <li>2. Click the <b>Catalog</b> tab.</li> <li>3. Search for the <b>Discovery</b> service and click that tile.</li> <li>4. Edit the Service name to something meaningful to you (for example: Discovery-<b>mydiscovery</b>) and click <b>Create</b> (If you have just created your account and accessed it from the confirmation email, you may need to log into IBM Cloud once again, then you can see the Create button in the bottom right corner).</li> <li>5. Click the <b>Launch Tool</b> button.</li> </ol>	 <p>The screenshot shows the IBM Cloud Catalog interface. At the top, there's a navigation bar with 'IBM Cloud', 'Catalog', 'Docs', 'Support', and 'Manage'. Below this, the 'Catalog' section has a search bar containing 'discovery'. On the left, there's a list of categories: 'All Categories (3)', 'Compute', 'Containers', 'Networking', 'Storage', 'AI (2)', and 'Analytics'. On the right, the 'Discovery' service tile is highlighted, showing 'Discovery Lite • IBM' and a description: 'Add a cognitive search and content analytics engine to applications.'</p>
<p>Watson Discovery News, a public data set that has been pre-enriched with cognitive insights, is included with Discovery. You can also use this public, unstructured data set to query for insights that you can integrate into your applications. Watson News is a dataset of primarily English language news sources that is updated continuously, with approximately 300,000 new articles and blogs added daily.</p> <p>In order to experiment with the robust querying capabilities of the Discovery service, you will create a new Collection and upload the AirBnB reviews that are already in JSON format as final documents to your collection.</p> <ol style="list-style-type: none"> <li>6. Click the folders icon in the top right corner of the page.</li> <li>7. Create a new collection by clicking <b>Upload your own data</b>.</li> </ol>	 <p>The screenshot shows the 'Manage data' page in the IBM Watson Discovery console. The header includes 'IBM Watson Discovery' and 'Instance: Discovery-mydiscovery'. The main heading is 'Manage data' with a subtitle: 'Collections of your private data and pre-enriched data to configure and query against. <a href="#">Learn more.</a>'. Below this, there are two buttons: 'Create a new data collection' and 'Upload your own data' (which is highlighted). To the right of these buttons is a 'Connect a data source' button. Below the buttons, there's a card for 'PRE-ENRICHED DATA' titled 'Watson Discovery News' with a dropdown menu for 'News sources: English'.</p>

<p>8. Name your collection; for example: <b>AirBnB</b>.</p> <p>9. Click <b>Create</b>.</p> <p>Each collection you create is a logical division of your data in the environment. Each collection will be queried independently when you get to the point of delivering results.</p> <p><i>Why would I want more than one collection?</i> There are a few reasons, including:</p> <ul style="list-style-type: none"> <li>You may want multiple collections in order to separate results for different audiences</li> <li>The data may be so different that it doesn't make sense for it all to be queried at once</li> </ul>	
<p>10. Extract the contents of the rental_reviews_easy_ingest_handpicked.zip file that you can download from Github at <a href="https://github.com/apischdo/Bluemix-workshop-assets/blob/master/Discovery.zip">https://github.com/apischdo/Bluemix-workshop-assets/blob/master/Discovery.zip</a>.</p> <p>11. Select all JSON files and upload them to the Discovery Service. This may take a few minutes.</p>	
<p>12. Now that you have uploaded documents in the AirBnB collection, you are ready to query the data. Click the magnifying glass icon on the left.</p> <p>13. Without running any specific queries, just click <b>Run query</b>.</p> <p>14. Take a moment and view the results. There are two ways to look at results, a <b>Summary</b> view and a <b>JSON</b> view. First, take a look at the Summary view.</p>	

15. Select **JSON** and look at the **enriched\_text** section below the first passage.
16. Notice the hierarchy: sentiment, entities, concepts. You will note that these are the enrichments that are applied by Watson Discovery to the ingested documents. In essence, they are information extracted from the unstructured text and now available as structured meta-data that you can query.
17. You are now ready to run some queries.

## Summary JSON

[Train Watson to improve results](#)

Query URL <https://gateway-syd.watsonplatform.net/discovery>

```
and it was so interesting exchanging
conversations and learning about each other's
cultures and experiences. The neighbourhood is
friendly, a walk away from Soho, Little Italy,
China Town and convenience to subway and bus
routes. ",
  "listing_longitude": "-73.99241186363034",
  "reviewer_id": "36205569",
  "host_id": "36656552",
  "date": "2015-09-25",
  "listing_latitude": "40.71892051069551",
  "reviewer_name": "Cindy",
  "enriched_text": {
    "sentiment": {...},
    "entities": [...],
    "concepts": [],
    "categories": [...]
```

18. Let's build a basic query using natural language. In the *Search for Documents*, type: **room with a great view**. On the right-hand panel, notice that results are presented in two sections.
  - a. **Passages** are relevant paragraphs extracted from documents
  - b. **Results** are enriched documents
19. In the **More options** section, you may set **Include relevant passages** to **No** and run your query again. You now only see the **Results** section.
20. You may click **Learn more** and take a moment and read the guidelines for querying. Basic queries can also be run using query language.
21. Clear all fields, and click **Run query** again.

## Summary JSON

[Train Watson to improve results](#)

Query URL <https://gateway-syd.watsonplatform.net/discovery/api/v1/en>

### Passages

"...Great experience renting with Paul! Clean room, comfortable apartment, great atmosphere, great location. I recommend it for anyone looking for a comfortable spot in close-in Brooklyn...."

"...Awesome apartment, cosy room, really clean, lot of space, and the room is on the other side of the house, do you'll be next to the bathroom and you'll get your intimacy. Achilles is a great host, really discreet but there if you need him...."

### Results

Showing 10 of 69 matching documents

▼ [airbnb\\_review99.json](#)

22. Now let's build an aggregation query using query language. Open the *Include analysis of your results* section, select **Edit in query language** and then click the question mark next to the input field.
23. Try out some aggregation formats replacing the example text to key words relevant to the AirBnB documents. Ensure that the hierarchy follows your classification scheme, not necessarily the example.
- ```
term(enriched_text.categories.label,count:5)
```
- Returns the five most frequent categories in the data set and the number of documents for each category
- ```
term(host_name,count:3)
```
- Returns the three most frequent host names for this collection.
24. Click the **JSON** tab for more details.
25. Clear all fields, and click **Run query** again.

### Build queries

Collection: AirBnB

Build a query using one or more of these components. [Learn more.](#)

Use a sample query

Search for documents

Include analysis of your results [Build in visual mode](#)

Write an aggregation query using the Discovery Query Language

```
term(enriched_text.categories.label,count:5)
```

Train Watson to improve results

**Summary JSON**

Query URL <https://gateway-syd.watson>

**Aggregations**

term(enriched\_text.categories.label)

- /travel/tourist facilities/hotel (44)
- /home and garden/bed and bath/bedroom (29)
- /real estate/apartments (24)
- /travel/tourist facilities/bed and breakfast (14)
- /pets/dogs (10)

26. And finally, let's build a filter query. Repeat the above action with the *Filter which documents you query* section by clicking the question mark and following the example format.
- ```
enriched_text.sentiment.document.label:"positive"
```
- ```
enriched_text.sentiment.document.score>0
```
- Results include documents with a positive sentiment.
- Filter queries run faster than basic queries because they are cached, but results are not sorted.

### Build queries

Collection: AirBnB

Build a query using one or more of these components. [Learn more.](#)

Use a sample query

Search for documents

Include analysis of your results

Filter which documents you query [Build in visual mode](#)

Write a filter to narrow down the document set using the Discovery Query Language

```
enriched_text.sentiment.document.label:"positive"
```

[More options](#)

Train Watson to improve results

**Summary JSON**

Query URL <https://gateway-syd.watson>

**Results**

Showing 10 of 90 matching documents

- airbnb\_review3.json
- airbnb\_review20.json
- airbnb\_review7.json
- airbnb\_review50.json
- airbnb\_review31.json
- airbnb\_review69.json
- airbnb\_review152.json
- airbnb\_review160.json
- airbnb\_review172.json
- airbnb\_review132.json

27. How about finding out apartments that are near tourist attractions? Clear all fields and run the following query in the *Filter which documents you query...* it's one line.

```
enriched_text.entities.disambiguation.subtype: "TouristAttraction"
```

The result set contains documents that have at least one entity with « TouristAttraction » in its list of disambiguation subtypes.

### Build queries

Collection: AirBnB

Build a query using one or more of these components. [Learn more.](#)

Use a sample query

Search for documents

Include analysis of your results

### Filter which documents you query

Build in visual mode

Write a filter to narrow down the document set using the Discovery Query Language

```
enriched_text.entities.disambiguation.subtype: "TouristAttraction"
```

Train Watson to improve results  
Summary **JSON**

Query URL: <https://gateway-syd.watson>

```
{
  "enriched_text": {
    "sentiment": { "...",
    "entities": [
      {
        "count": 1,
        "sentiment": {
          "score": 0,
          "label": "neutral"
        },
        "text": "Met Museum",
        "relevance": 0.77296,
        "type": "Facility",
        "disambiguation": {
          "subtype": [
            "Organization",
            "Location",
            "HistoricPlace",
            "TouristAttraction",
            "Building",
            "Museum"
          ]
        }
      }
    ]
  }
}
```

28. Notice the **Query URL** built based on your query. This URL will be used when you build an application that needs to query the document collection using Discovery. You may copy the URL into a temporary location, it will be used in the last section of this lab, **Working with the Discovery API**.

Train Watson to improve results  
Summary **JSON**

Query URL: <https://gateway-syd.watson>

```
https://gateway-syd.watsonplatform.net/discovery/api/v1/environments/7d560384-6c90-49a6-8711-f873e10d326e/collections/03b487cf-73d3-47ad-8cd9-5a9ae58bb5a0/query?version=2018-08-01&filter=enriched_text.entities.disambiguation.subtype%3A%22TouristAttraction%22&deduplicate=false&highlight=true&passages=true&passages.count=5&query={
  "enriched_text": {
    "sentiment": { "...",
    "entities": [
      {
        "count": 1,
        "sentiment": {
          "score": 0,
          "label": "neutral"
        },
        "text": "Met Museum",
        "relevance": 0.77296,
        "type": "Facility",
        "disambiguation": {
          "subtype": [
            "Organization",
            "Location",
            "HistoricPlace",
            "TouristAttraction",
            "Building",
            "Museum"
          ]
        }
      }
    ]
  }
}
```

## Building a custom configuration

You have been adding content to your collection just after it was created. When a collection is created, a set of default configurations are automatically provided. If you are happy with these defaults, you can proceed to uploading your content. This is what we have done so far. If you need to configure the service to process the content the way that you want, then the best practice is to define a **custom configuration** before uploading your content.

29. From the “Manage Data” page, in the Configuration section, Click **Switch**.
30. Click **Create a new configuration**.
31. Name your configuration, for example **Config01**. In order to set up this configuration, you will:

- Identify some sample content (documents that are representative of your files)
- Upload the sample content
- Adjust the conversion
- Define enrichments
- Normalize the results

IBM Watson Discovery

Cookie Preferences Instance: Discovery-mydiscovery

Manage data > AirBnB

View data schema

Overview Errors and warnings (0)

Document count	Errors and warnings	Configuration	Collection info
101	0 documents failed <a href="#">View details</a>	Default Configuration <a href="#">Switch</a>	Created on 8/2/2018 9:49:14 am EDT Last updated 8/2/2018 10:51:50 am EDT <a href="#">Use this collection in API</a>

To make the configuration process more efficient, you can upload of Microsoft Word, HTML, JSON, or PDF files that are representative of your document set. These are called sample documents. Sample documents are not added to your collection, they are only used to identify fields that are common to your documents and customize those fields to your requirements.

When creating a new configuration file in the Discovery tooling, you can upload sample documents via drag and drop or browse. Click on the file name in the Upload Sample Documents pane to preview each file.

Remember the following items when uploading sample documents:

- All of your documents are converted to JSON before they are enriched and indexed.
- Microsoft Word and PDF documents are converted to HTML first, then JSON.
- HTML documents are converted directly to JSON.

Note: Sample documents are automatically deleted after 1 month, but you can upload the same documents again if you would like to make additional changes to your configuration.

32. For this exercise, upload a PDF sample document from your local drive. Here we are using the following file:  
<https://public.dhe.ibm.com/common/ssi/ecm/41/en/41015641usen/ibmwatsonstudioflyer.pdf>

33. Once the document is uploaded it will appear in the right pane. Click the document name.

#### Upload sample documents

You can preview your configuration by using sample documents that are representative of your data collection. [Learn more.](#)



Drag and drop up to 10 documents  
or [browse from computer](#)

Sample documents are stored locally

File format: PDF, WORD, HTML, and JSON

Maximum file size: 1MB

[ibmwatsonstudioflyer.pdf](#)

34. Click the **Convert** link.

35. Select **PDF**.

36. For example, this particular document uses font size **24** and **bold** for **h1**, font size **14** and **bold** for **h2** and font size **12** and **bold** for **h3**. On the left panel, you can specify html text conversion features that are relevant to the document.

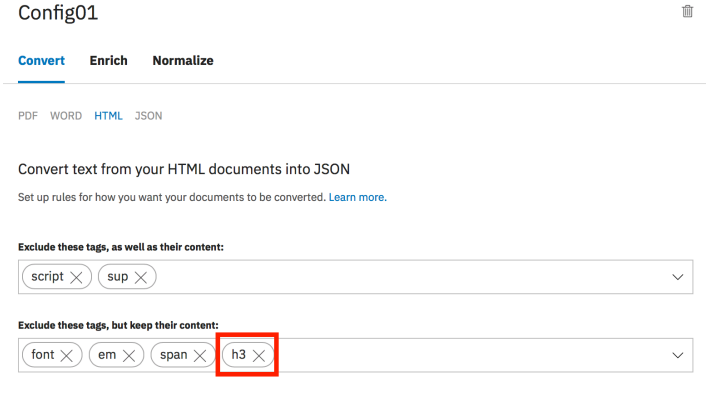
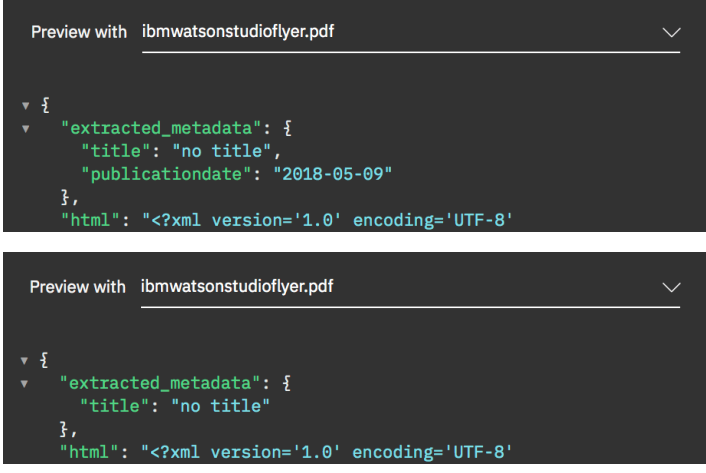
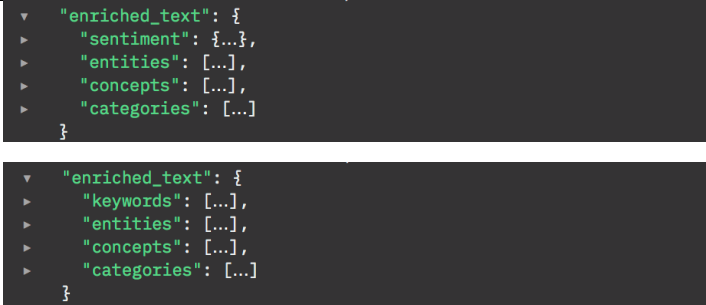
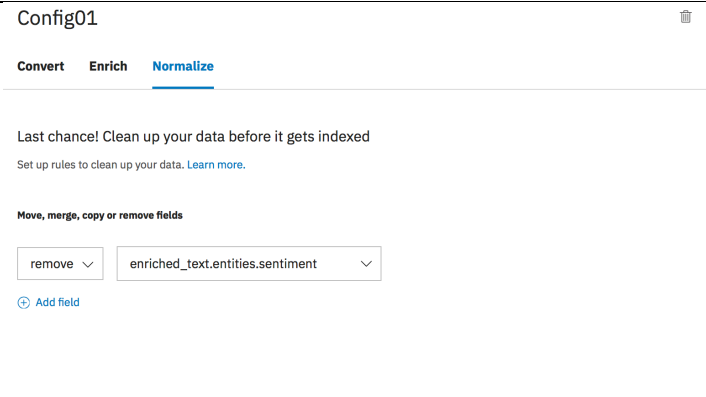
37. Click **Apply and Save**. Notice the updates in the right panel.

38. Click the **Learn more** link and take a moment to read the details behind text conversion.

Preview with [ibmwatsonstudioflyer.pdf](#)

```
<?xml version='1.0' encoding='UTF-8' standalone='yes'?>
<html>\n<head>\n  <meta content=\"text/html; charset=UTF-8\"
http-equiv=\"Content-Type\"/><meta content=\"2018-05-09\"
name=\"publicationdate\"/><meta content=\"2\"
name=\"numPages\"/><title>no title</title></head>\n<body><h1>
<p>IBM Watson Studio\n</p></h1><p>Build, train, deploy and
manage AI models, and prepare and analyze data, in a single,
integrated environment.\n</p><h2><p>What is Watson Studio?\n</p>
</h2><p>IBM Watson((R)) Studio accelerates the machine and deep
learning workflows required to infuse AI into your business to
drive innovation. It provides a suite of tools for data
scientists, application developers and subject matter experts to
collaboratively and easily work with data to build, train and
deploy models at scale.\n</p><h3><p>Embedded Watson tools\n</p>
</h3><p>Train with embedded AI services, like Watson Visual
Recognition, and deploy models as APIs or CoreML.\n</p><h3>
<p>Built on open source\n</p></h3><p>Use familiar open source
data science and machi..."
```



<p>39. Click the <b>HTML</b> link.</p> <p>40. For the purposes of this exercise, type <b>h3</b> in the <i>Exclude these tags, but keep their content</i> box. Click <b>Apply &amp; Save</b> and look at the right panel to see that your <code>&lt;h3&gt;</code> tags are now gone.</p> <p>Important to note, that the conversions you specify for the html conversion will apply to both your PDF and MS Word uploads. If you want your PDF documents to have a different configuration from your MS Word documents, then you have to place them in separate collections with separate custom configurations.</p>	
<p>41. Click the <b>JSON</b> link.</p> <p>You will now have an opportunity to do some housekeeping and clean the output of your documents before they are indexed.</p> <p>42. For the purpose of this exercise, select <b>Add field</b>, then <b>remove</b> and <code>extracted_metadata.publicationdate</code> (or any other field applicable to your document). Turn on the <b>Remove empty fields</b> option</p> <p>43. Click <b>Apply and save</b>.</p> <p>44. Look at the right panel and notice that the removed field is now gone.</p>	
<p>45. Now that we have configured the Conversion, let's move on to configuring the enrichment step. Click the <b>Enrich</b> link.</p> <p>46. Take a moment and observe the enrichments that are applied to the document. In the right panel, "enriched_text".</p> <p>Let's assume that you do not need the <b>sentiment</b> enrichment based on this sample document.</p> <p>47. Remove the <b>sentiment</b> enrichment, add the <b>Keywords</b> enrichment and click <b>Apply and Save</b>. Check changes in the enriched sample documents in the right panel.</p>	
<p>48. Click <b>Normalize</b>. This is where you get an opportunity to clean and normalize the JSON output.</p> <p>49. In this example <b>remove</b> the field <code>enriched_text.entities.sentiment</code>.</p> <p>50. Click <b>Apply and Save</b>.</p> <p>51. Congratulations! Now that you have created your custom configuration based on sample data, you are ready to upload and query your actual (not sample) data. The query does not run on sample documents. The sample documents reside in a temporary repository and are not indexed for query.</p>	

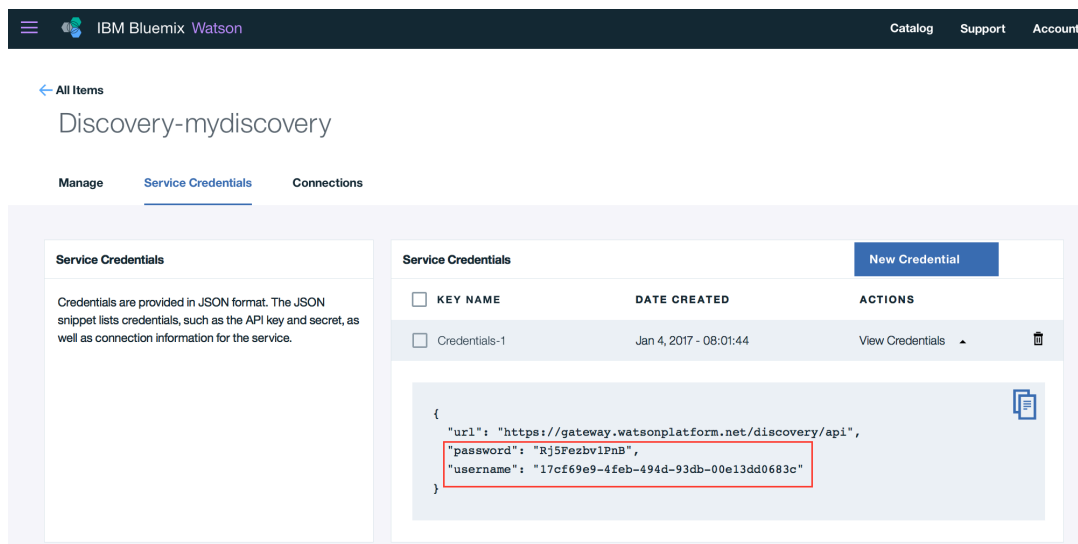
## Working with the Discovery API

So you had some fun with the tooling. But to do bulk uploads, crawling (static) and including custom models and annotations from Watson Knowledge Studio (not in the scope of this document), you would use the API approach. Let's begin by installing Postman, which is a recommended tool for developers who work with APIs.

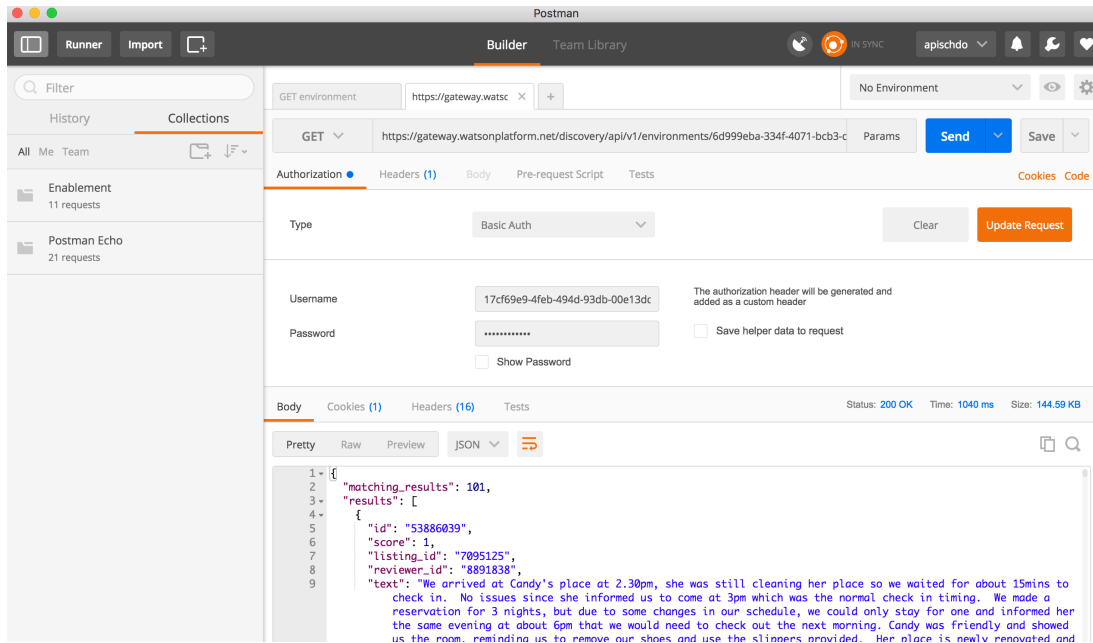
- 1) Direct your browser to this link: <https://www.getpostman.com/>
- 2) Download the free Postman app appropriate for your operating system, and install it.

With that installation out of the way, you are now ready to begin using the tool.

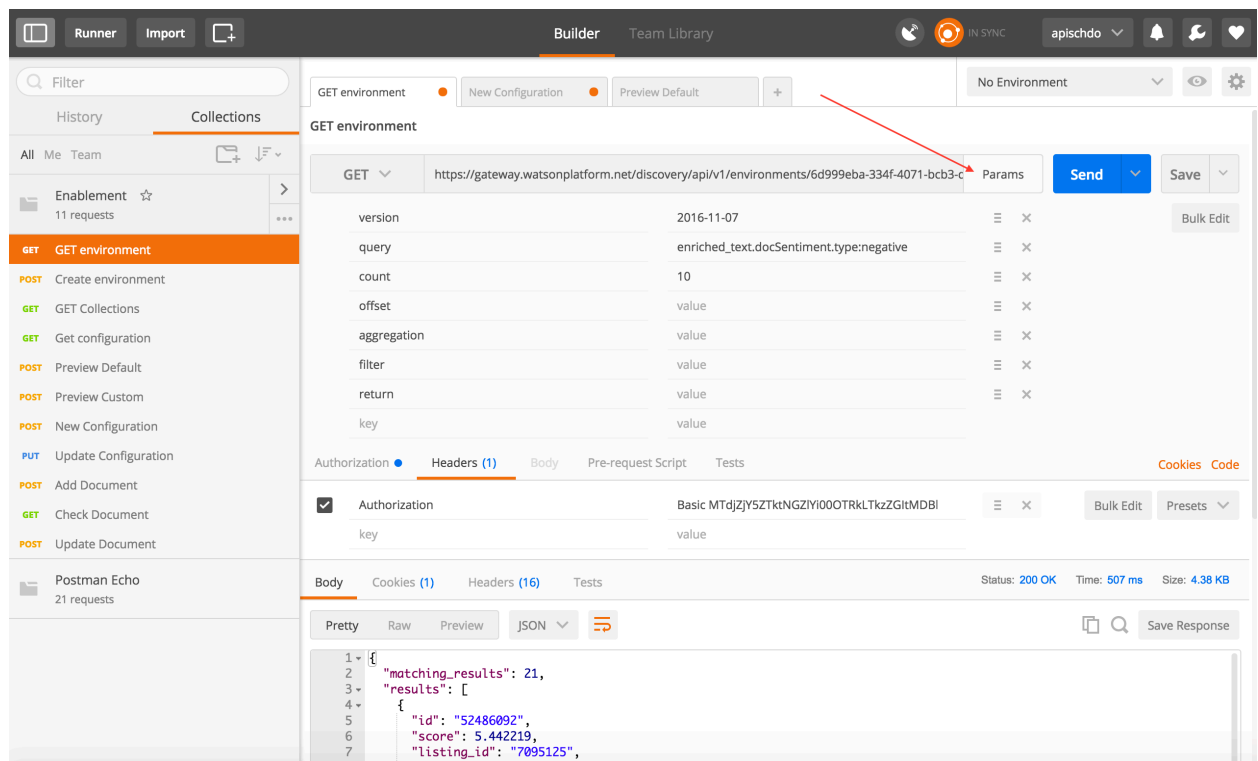
1. Go to IBM Cloud and access the Service Credentials panel of the Discovery service that you created in the beginning of this document.
2. Click the Service credentials link and note the username and password. You will be entering these in Postman along with the URL that copied earlier.



1. Open the Postman application: paste the URL as a Get method; select **Basic Auth** and enter the service credentials.
2. Click **Send**.



3. Click **Params**, enter a count of 10 as default and run some of the queries that you performed earlier. When you build your own application, you would not be using the Tooling we saw earlier, but APIs to integrate your query results within the frames of your application.



Use the documentation frequently by clicking Learn more from the tooling interface and explore the contents of the doc from the left panel. Enjoy your discovery journey.