,	<ul> <li>Find the most actively modified module?</li> <li>How many commits occurred during the studied period?</li> <li>How much churn occurred during the studied period? Churn is defined as the sum of added and removed lines by all commits.</li> <li>NB: This workflow is responsible for the pre-processing, analysis, and generation of insight from the collected data. It is assumed the automated collection of the data via the script accessible in thesame folder with this notebook has been completed. The collected data will be leaded here before the other process in the workflow executes.</li> </ul>
	will be loaded here before the other process in the workflow executes.  Required imports:  # Built-in libraries import json import os  # The normal data science ecosystem libraries
	<pre># pandas for data wrangling import pandas as pd  # Plotting modules and libraries required import matplotlib as mpl import matplotlib.pyplot as plt  Required settings: # Settings: # 1. Command needed to make plots appear in the Jupyter Notebook</pre>
	<pre>%matplotlib inline  # 2. Command needed to make plots bigger in the Jupyter Notebook plt.rcParams['figure.figsize']= (12, 10)  # 3. Command needed to make 'ggplot' styled plots- professional and yet good looking theme. plt.style.use('ggplot')  # 4. This will make the plot zoomable # mpld3.enable_notebook()</pre>
	Other utility functions for data manipulation  # Utility data manipulation functions  # 1. Extract path parameters from filename  def get_path_parameters(dframe):     filename = os.path.basename(dframe["filename"])     filetype = os.path.splitext(dframe["filename"])[1]     directory = os.path.dirname(dframe["filename"])  return directory, filename, filetype
	<pre>1. Loading the data # Open and load json file with open('data.json', encoding="utf8") as file:     data = json.load(file)     print("data loaded successfully")  data loaded successfully</pre> Data normalization
	The collected commit data is a semi-structured json which has nested data similar to the image below. Files is a list of file objects. The loaded data will be normalized into a flat table using pandas.json_normalize.  {  "sha":"232f8275ec00767d1f100cacae4823e6f77e04ef",     "node_id":"C_kwDDAAw0D9oAKDIzMmY4Mjc1ZWMwMDc2N2QxZjEwMGNhY2F1NDgyM2U2Zjc3ZTA0ZWY",     "commit":{  }  "url":"https://api.github.com/repos/openstack/nova/commits/232f8275ec00767d1f100cacae4823e6f77e04ef",     "html_url":"https://github.com/openstack/nova/commits/232f8275ec00767d1f100cacae4823e6f77e04ef",     "comments_url":"https://api.github.com/repos/openstack/nova/commits/232f8275ec00767d1f100cacae4823e6f77e04ef",     "comments_url":"https://api.github.com/repos/openstack/nova/commits/232f8275ec00767d1f100cacae4823e6f77e04ef',     "comments_url":"https://api.github.com/repos/openstack/nova/commits/232f8275ec00767d1f100cacae4823e6f77e04ef',
	<pre>"comments_url":"https://api.github.com/repos/openstack/nova/commits/232f8275ec00767d1f100cacae4823e6f77e04ef/comments", "author":null, "committer":{\(\Overline{\</pre>
	<pre>df = pd.json_normalize(data, "files", ["commit_node_id", "commit_sha", "commit_html_url", "commit_date"</pre> 2. Displaying current state of the data
	# The first 3 rows df.head(3)  sha filename status additions deletions changes  0 0438913de34e2751410da7a035e21b7e73325760 nova/exception.py modified 4 0 4 https://gi
	2       d71d13ab3721422c013cff586e7c31326ae50e1c       nova/tests/unit/virt/libvirt/test_host.py       modified       8       0       8       https://github.com/github.co
	1971 48adcb07d6af188ced7acb6d7ddcbfbe8fec0489 nova/tests/unit/virt/vmwareapi/fake.py modified 4 2  1972 0d8639f9ae83bb81df262d98d58a7ebd2109de41 nova/tests/unit/virt/vmwareapi/test_driver_api.py modified 1 1  # Summary of the dataframe
	<pre>df.info()  <class 'pandas.core.frame.dataframe'=""> RangeIndex: 1973 entries, 0 to 1972 Data columns (total 15 columns):     # Column</class></pre>
	5 changes 1973 non-null int64 6 blob_url 1973 non-null object 7 raw_url 1973 non-null object 8 contents_url 1973 non-null object 9 patch 1965 non-null object 10 previous_filename 8 non-null object 11 commit_node_id 1973 non-null object 12 commit_sha 1973 non-null object 13 commit_html_url 1973 non-null object 14 commit_date 1973 non-null object dtypes: int64(3), object(12) memory usage: 231.3+ KB
	<pre># Let us manually examine atleast one commit and see if the present rows are correct. # We use the most recent commit as at the development time. # Pls note that this commit will not be part of commits after 6 month from today February 17th, 2022 commit = '3a14c1a4277a9f44b67e080138b28b680e5e6824' df[df["commit_sha"] == commit]</pre> <pre>sha</pre> filename status additions deletions changes
	36       08006e2b92a44b709b3ef6171cfb9a95519c8f5e       nova/compute/api.py       modified       18       4       22       https://githu         37       79a62da21a6d4e768b48f8be3c44e39b9bcd3a83       nova/tests/unit/compute/test_api.py       modified       30       0       30       https://githu         38       ef5582543a2ac953179e5dbc3e493e69c9bc84bf       releasenotes/notes/bug-1960401-504eb255253d966       added       8       0       8       https://githu
	<pre>"sha":"3a14c1a4277a9f44b67e080138b28b680e5e6824",     "node_id": "C_kwDDAAwOD9oAkDNhMTRjWWE0Mjc3YTlmNDRINjdIMDgwMTM4Yj14YjY4MGU1ZTY4MjQ",     "commit": {</pre>
	"additions":18,     "deletions":4,     "changes":22,     "blob_url": "https://github.com/openstack/nova/blob/3a14c1a4277a9f44b67e880138b28b680e5e6824/nova/compute/api.py",     "raw_url": "https://github.com/openstack/nova/raw/3a14c1a4277a9f44b67e880138b28b680e5e6824/nova/compute/api.py",     "contents_url": "https://api.github.com/repos/openstack/nova/comtents/nova/compute/api.py?ref=3a14c1a4277a9f44b67e880138b28b680e5e6824",     "patch": "@0 -4822,10 +4822,24 @0 def_attach_volume(self, context, instance, volume, device,\n
	"sha":"ef5582543a2ac953179e5dbc3e493e69c9bc84bf",
	<pre># Removing columns not needed for the analysis columns = ['previous_filename', 'patch', 'contents_url', 'raw_url', 'commit_node_id'] df.drop(columns, inplace=True, axis=1)  # Generating and adding extra columns df[["directory", "file_name", "file_type"]] = df.apply(lambda x: get_path_parameters(x), axis=1, result_ # Delete the previous filename column as it is no longer required df.drop("filename", inplace=True, axis=1)</pre>
	<pre>df.drop("filename", inplace=True, axis=1)  # Rename columns df.rename(columns={"sha": "file_sha", "status": "file_status", "additions":"no_of_additions", "deletions  # Optimising the data frame by correcting the data types. # This will also make more operations possible on the data frame  df = df.astype({'file_sha': 'str', 'file_status': 'category', 'no_of_additions':'int', 'no_of_deletions' df['commit_date'] = pd.to_datetime(df['commit_date'], infer_datetime_format=True)</pre>
1	df.info() <class 'pandas.core.frame.dataframe'=""> RangeIndex: 1973 entries, 0 to 1972  Data columns (total 12 columns):  # Column Non-Null Count Dtype</class>
]	4 changes 1973 non-null int32 5 blob_url 1973 non-null object 6 commit_sha 1973 non-null object 7 commit_html_url 1973 non-null object 8 commit_date 1973 non-null datetime64[ns, UTC] 9 directory 1973 non-null object 10 file_name 1973 non-null object 11 file_type 1973 non-null category dtypes: category(2), datetime64[ns, UTC](1), int32(3), object(6) memory usage: 135.8+ KB
	A. Basic Analysis and Visualization  1. Total number of commits that occured during the studied period.  # value_counts returns a series object counting all unique values.  # the 1st value being the most frequently occuring i.e. the commit with highest no of file changes.  df["commit_sha"].value_counts()  f3d48000b139ec38d92da276a43a8387f76cbc89 83 0e0196d979cf1b8e63b9656358116a36f1f09ede 83 ccef1940bf92da7441beb6b88fa9f998b1e9b2b2 36
	d2a5fe5621d6ff1ae8ba5087049e0c4347592cf6 21 ac21c6674c8444edc5afd25b7d63936182fe3580 21 e2b1581d8c4f03cecca770aabb8e6123a7bef93a 1 f318f822fcf6dec4c3cd9b7e5111f3e1371aa51a 1 61b169d40f7df3c9ea782b3cf2ac96d3fbdceef2 1 f024490e95c2bdb8072247e9907c6aa1475c80d8 1 e28afc564700a1a35e3bf0269687d5734251b88a 1 Name: commit_sha, Length: 467, dtype: int64  print("The total no. of processed commits is: {commits_total}".format(commits_total = len(df["commit_sha
	The total no. of processed commits is: 467  2. The 12 most modified files  df["file_name"].value_counts().head(12)  driver.py 77 api.py 53 test_driver.py 48 test_api.py 44 neutron.py 44
]	manager.py 38  test_servers_resource_request.py 37  test_migrations.py 36  test_neutron.py 36  servers.py 31  test_compute.py 30  Name: file_name, dtype: int64  df["file_name"].value_counts().head(12).sort_values().plot.barh(figsize=(10, 9)); plt.axhline(0, color=')
	The 12 most modified files  driver.py -  api.py -  test_driver.py -
	neutron.py -  test_api.py -  nova.py -  manager.py -  test_servers_resource_request.py -
	test_neutron.py -  test_migrations.py -  servers.py -  test_compute.py -  0 10 20 30 40 50 60 70 80
	Total number of changes  3. The 12 most modified directories  # the term directory is used in place of module df["directory"].value_counts().head(12)  nova/tests/unit/compute 101 nova/tests/functional 91 nova/compute 79 nova/api/openstack/compute 79
1	nova/tests/fixtures 75 releasenotes/notes 74 nova/virt/libvirt 73 doc/source/admin 61 nova/tests/unit/virt/libvirt 61 nova/tests/functional/regressions 59 nova/network 52 Name: directory, dtype: int64  df["directory"].value_counts().head(12).sort_values().plot.pie(autopct='%1.1f%%', figsize=(20,8),
]	Text(0.5, 1.0, 'The 12 most modified modules')  The 12 most modified modules  Pova/network
]	nova/tests/functional/regressions  nova/tests/functional/regressions
1	nova/tests/functional/regressions  doc/source/admin 6.7% 5.9% 11.5% nova/tests/functional  nova/tests/functional  nova/tests/unit/virt/libvirt 6.9% 9.0% nova/api/openstack/compute
	nova/tests/functional/regressions  doc/source/admin 6.7% 5.9% 11.5% nova/tests/functional  nova/tests/functional  nova/tests/unit/compute  nova/tests/functional  nova/tests/functional  nova/tests/unit/compute  nova/tests/functional  nova/tests/functional  nova/tests/functional  nova/tests/functional  nova/tests/functional  nova/tests/functional  nova/tests/functional  nova/tests/functional  nova/tests/functional
	nova/tests/functional/regressions  doc/source/admin 6.7% 5.9% 11.5% nova/tests/functional  nova/tests/unit/virt/libvirt 6.9% 9.0% nova/api/openstack/compute  8.3% 9.0% nova/api/openstack/compute  8.3% 8.5% 8.7% nova/compute  4. The most modified file types  df ["file_type"] . value_counts () .py 1557 .rst 167 .yam 114 .]son 31 .txt 31 .tp1 21 .in1 18 .tp1 21 .in1 18
	nova/tests/functional/regressions  doc/source/admin 6.7% 5.9% 10.3% nova/tests/functional  nova/tests/functional  6.9% 10.3% 10.3% nova/tests/functional  nova/tests/functional  6.9% 10.3% nova/tests/functional  nova/tests/functional  10.3% nova/tests/functional  nova/tests/functional  10.3% nova/tests/functional  nova/tests/functional  nova/tests/functional  10.3% nova/tests/functional
	nova/tests/unit/ornpute  doc/source/admin  6.7%  5.9%  10.3%  nova/tests/unit/ornpute  nova/tests/unit/ornpute  6.9%  9.0%  nova/dests/functional  nova/qei/openstack/compute  8.3%  8.3%  8.3%  8.5%  8.7%  nova/compute  1.04  1.05  1.0
	nova/tests/unctional/regressions  doc/source/admin  67% 59% 103% nova/tests/functional  68% 90% nova/qpilopenstack/compute  4. The most modified file types  def["file_type"].value_counts()
	novalvests/unctional repression  decharactional repression  decharactional repression  63%  63%  63%  63%  63%  63%  63%  63
	### Continue of the second of
	novabetachundrottende  de de construit de la c
	## 130 ##
	## A The most modified file types  ## A The most mo
	A. The most modified file types  4. The most modified file types  5. Churn  5. Churn  6. Churn  7. Inter of 21 colored countries for the colored countries for the colored countries file types  4. The total number of the modifications by directory i.e. no. of rows per directory. A row in directories a file damage & the countries for the colored colored countries for the colored countries for the colored colore
	### According to the property of the property
	### A. The most modified file types  4. The most modified file typ
	4. The most modified file types  or sentential to the part of the modification by directory is not of records a file change of the committee and the committ
	4. The most modified file types  4. The most modified file types  5. Class (2) of the control of
	A. The most modified file types  After a control modified file modified file types  After a control modified file types  After a control modified file types  After a control modified file modified file types  After a control modified file modified file types  After a control modified file types  After a control modified file modified file types  After a control modified file types  After a control modified file modified file modified file types  After a control modified file mo
	4. The most modified file types  5. The most modified file types  4. The most modified file types  5. The most modified file types  6. The most model
	4. The most modified file types  3. The cost modified file types  4. The most modified file types  5. Chun  6.
	4. The most modified fit oppose  2. Churn  2.
	A. The most modified file types  So Chum  So Chum  So Chum  The state of the stage course of the stage of the state of the stage of the state of the stage of the
	A. The most modeline file oppose  ### The most mode
	A Communication of the property of the control of the property
	## 1
	A. The macron of fine by discourse of the control o