

```
In [11]: 1 !pip install pyPDF2
```

Requirement already satisfied: pyPDF2 in c:\users\nihar\anaconda3\lib\site-packages (3.0.1)
Requirement already satisfied: typing_extensions>=3.10.0.0 in c:\users\nihar\anaconda3\lib\site-packages (from pyPDF2) (4.8.0)

```
In [12]: 1 #!pip uninstall pyPDF2
```

```
In [13]: 1 #pip install pyPDF2=3.0.1
```

```
In [17]: 1 import PyPDF2  
2 from PyPDF2 import PdfFileReader
```

```
In [18]: 1 PyPDF2.__version__
```

```
Out[18]: '3.0.1'
```

```
In [19]: 1 #Creating a pdf file object
2 pdf = open("file1pdf.pdf","rb")
3
4 #creating pdf reader object
5 pdf_reader = PyPDF2.PdfReader(pdf)
6
7 #checking number of pages in a pdf file
8 print("Number of pages:", len(pdf_reader.pages))
9
10 #creating a page object
11 page = pdf_reader.pages[1]
12
13 #finally extracting text from the page
14 print(page.extract_text())
15
16 #closing the pdf file
17 pdf.close()
```

Number of pages: 35

Development Plan for Greater Mumbai 2014-2034

Acknowledgements

The Consultant wishes to thank the following individuals from the Municipal Corporation of Greater Mumbai for their invaluable support, insights and contributions towards 'Working Paper 1

- Preparation of Base Map' for the preparation of the Development Plan for Greater Mumbai 2014-34.

☐ Mr. Subodh Kumar, IAS, Municipal Commissioner;

☐ Mr. Rajeev Kuknoor, Chief Engineer Development Plan;

☐ Mr. Sudhir Ghate, Deputy Chief Engineer Development Plan;

☐ Mr. A.G. Marathe, Deputy Chief Engineer Development Plan;

☐ Mr. R. Balachandran, Executive Engineer and Town Planning Officer, Development Plan.

Our gratitude to the following experts for their invaluable insights and support:

☐

Mr. V.K Phatak, Former Chief Town Planner (MMRDA);

☐ Mr. A.N Kale, Former Chief Engineer, (DP);

☐ Mr. A. S Jain Former Dy. Chief Engineer, (DP).

We wish to especially thank MCGM officers, Mr. Jagdish Talreja, Mr. Dinesh Naik, Mr. Hiren Daftardar, Ms. Anita Naik for their continual support since the beginning of the project and their

help towards familiarization and data collection. They have been instrumental in helping to

contact various MCGM departments as well as in helping to establish contact with personnel from

other government departments and organizations. Many thanks for the MCGM team, for deploying personnel, particularly Mr. Prasad Gharat, on extensive field visits that have helped in

understanding actual ground conditions.

We apologize if we have inadvertently omitted anyone to whom acknowledgement is due. We hope

and anticipate the work's usefulness for the intended purpose.

```
In [20]: 1 import PyPDF2, urllib , nltk
2 from io import BytesIO
3 from nltk.tokenize import word_tokenize
4 from nltk.corpus import stopwords
```

In [21]: 1 nltk.download('stopwords')

[nltk_data] Downloading package stopwords to
[nltk_data] C:\Users\nihar\AppData\Roaming\nltk_data...
[nltk_data] Package stopwords is already up-to-date!

Out[21]: True

In [22]: 1 wFile = urllib.request.urlopen('http://www.udri.org/pdf/02%20working%20paper%201.pdf')
2 pdfreader = PyPDF2.PdfReader(BytesIO(wFile.read()))

In [23]: 1 pageObj = pdfreader.pages[2]
2 page2 = pageObj.extract_text()
3 punctuations = ['(', ')', ';', ':', '[', ']', ', ', '...', '.']
4 tokens = word_tokenize(page2)
5 stop_words = stopwords.words('english')
6 keywords = [word for word in tokens if not word in stop_words and not word in punctuations]

In [24]: 1 keywords

Out[24]: ['Development',
'Plan',
'Greater',
'Mumbai',
'2014-2034',
'Table',
'Contents',
'The',
'Consultant',
'wishes',
'thank',
'following',
'individuals',
'Municipal',
'Corporation',
'Greater',
'Mumbai',
'invaluable',
'...']

In [25]: 1 name_list = list()
2 check = ['Mr.', 'Mrs.', 'Ms.']
3 for idx, token in enumerate(tokens):
4 if token.startswith(tuple(check)) and idx < len(tokens)-1:
5 name = token + tokens[idx+1] + ' ' + tokens[idx+2]
6 name_list.append(name)
7 print(name_list)

['Mr.Jagdish Talreja', 'Mr.Dinesh Naik', 'Mr.Hiren Daftardar', 'Ms.Anita Naik', 'Mr.Prasad Gh
arat']

In [26]: 1 wFile.close()

In [27]: 1 !pip install python_docx

Requirement already satisfied: python_docx in c:\users\nihar\anaconda3\lib\site-packages (1.1.2)
Requirement already satisfied: lxml>=3.1.0 in c:\users\nihar\anaconda3\lib\site-packages (from python_docx) (4.9.1)
Requirement already satisfied: typing-extensions>=4.9.0 in c:\users\nihar\anaconda3\lib\site-packages (from python_docx) (4.12.2)

```
In [34]: 1 import docx
```

```
In [36]: 1 doc = open("Task-1-Answers.docx", "rb")  
2 document = docx.Document(doc)
```

```
In [37]: 1 docu=""  
2 for para in document.paragraphs:  
3     docu+=para.text  
4 print(docu)
```

```
In [38]: 1 !pip install bs4
```

Collecting bs4

Downloading bs4-0.0.2-py2.py3-none-any.whl (1.2 kB)

Requirement already satisfied: beautifulsoup4 in c:\users\nihar\anaconda3\lib\site-packages (from bs4) (4.11.1)

Requirement already satisfied: soupsieve>1.2 in c:\users\nihar\anaconda3\lib\site-packages (from beautifulsoup4->bs4) (2.3.1)

Installing collected packages: bs4

Successfully installed bs4-0.0.2

```
In [39]: 1 import urllib.request as urllib2  
2 from bs4 import BeautifulSoup
```

```
In [40]: 1 response = urllib2.urlopen('https://en.wikipedia.org/wiki/Natural_language_processing')  
2 html_doc = response.read()
```

```
In [41]: 1 soup = BeautifulSoup(html_doc, 'html.parser')
2 strhtm = soup.prettify()
3 print (strhtm[:5000])
```

```
<!DOCTYPE html>
<html class="client-nojs vector-feature-language-in-header-enabled vector-feature-language-in-main-page-header-disabled vector-feature-page-tools-pinned-disabled vector-feature-toc-pinned-d-clientpref-1 vector-feature-main-menu-pinned-disabled vector-feature-limited-width-clientpref-1 vector-feature-limited-width-content-enabled vector-feature-custom-font-size-clientpref-1 vector-feature-appearance-pinned-clientpref-1 vector-feature-night-mode-enabled skin-theme-clientpref-day vector-sticky-header-enabled vector-toc-available" dir="ltr" lang="en">
<head>
  <meta charset="utf-8"/>
  <title>
    Natural language processing - Wikipedia
  </title>
  <script>
    (function(){var className="client-js vector-feature-language-in-header-enabled vector-feature-language-in-main-page-header-disabled vector-feature-page-tools-pinned-disabled vector-feature-toc-pinned-clientpref-1 vector-feature-main-menu-pinned-disabled vector-feature-limited-width-clientpref-1 vector-feature-limited-width-content-enabled vector-feature-custom-font-size-clientpref-1 vector-feature-appearance-pinned-clientpref-1 vector-feature-night-mode-enabled skin-theme-clientpref-day vector-sticky-header-enabled vector-toc-available";var cookie=document.cookie.match(/(?:^|; )enwikimwclientpreferences=([^\;]+)/);if(cookie){cookie[1].split('%2C').forEach(function(pref){className=className.replace(new RegExp('(\\'+pref.replace(/-clientpref-\\w+$|\\w-|/g, '\\\\')+\\'+pref+'$2')),'$1'+pref+'$2');});}document.documentElement.className=className;})();RLCONF={"wgBreakFrames":false,"wgSeparatorTransformTable":["",""],"wgDigitTransformTable":["",""],"wgDefaultDateFormat":"dmy","wgMonthNames":["","January","February","March","April","May","June","July","August","September","October","November","December"],"wgRequestId":"bb8ccbb-931f-4281-aaf8-3aff8c9effb8","wgCanonicalNamespace":"","wgCanonicalSpecialPageName":false,"wgNamespaceNumber":0,"wgPageName":"Natural_language_processing","wgTitle":"Natural language processing","wgCurRevisionId":1274942014,"wgRevisionId":1274942014,"wgArticleId":21652,"wgIsArticle":true,"wgIsRedirect":false,"wgAction":"view","wgUserName":null,"wgUserGroups":["*"],"wgCategories":["All accuracy disputes","Accuracy disputes from December 2013","Harv and Sfn no-target errors","CS1 errors: periodical ignored","CS1 maint: location","Articles with short description","Short description is different from Wikidata","Articles needing additional references from May 2024","All articles needing additional references","All articles with unsourced statements","Articles with unsourced statements from May 2024","Commons category link from Wikidata","Natural language processing","Computational fields of study","Computational linguistics","Speech recognition"],"wgPageViewLanguage":"en","wgPageContentLanguage":"en","wgPageContentModel":"wikitext","wgRelevantPageName":"Natural_language_processing","wgRelevantArticleId":21652,"wgIsProbablyEditable":true,"wgRelevantPageIsProbablyEditable":true,"wgRestrictionEdit":[],"wgRestrictionMove":[],"wgNoticeProject":"wikipedia","wgCiteReferencePreviewsActive":false,"wgFlaggedRevsParams":{"tags":{"status":{"levels":1}}},"wgMediaViewerOnClick":true,"wgMediaViewerEnabledByDefault":true,"wgPopupsFlags":0,"wgVisualEditor":{"pageLanguageCode":"en","pageLanguageDir":"ltr","pageVariantFallbacks":"en"},"wgMFDisplayWikibaseDescriptions":{"search":true,"watchlist":true,"tagline":false,"nearby":true},"wgWMESchemaEditAttemptStepOversample":false,"wgWMEPageLength":60000,"wgEditSubmitButtonLabelPublish":true,"wgULSPosition":"interlanguage","wgULSisCompactLinksEnabled":false,"wgVector2022LanguageInHeader":true,"wgULSisLanguageSelectorEmpty":false,"wgWikibaseItemId":"Q30642","wgCheckUserClientHintsHeadersJsApi":["brands","architecture","bitness","fullVersionList","mobile","model","platform","platformVersion"],"GEOHomepageSuggestedEditsEnableTopics":true,"wgGETopicsMatchModeEnabled":false,"wgGESTructuredTaskRejectionReasonTextInputEnabled":false,"wgGELevelingUpEnabledForUser":false};RLSTATE={"ext.globalCssJs.user.styles":"ready","site.styles":"ready","user.styles":"ready","ext.globalCssJs.user":"ready","user":"ready","user.options":"loading","ext.cite.styles":"ready","ext.math.styles":"ready","skins.vector.search.codex.styles":"ready","skins.vector.styles":"ready","skins.vector.icons":"ready","jquery.makeCollapsible.styles":"ready","ext.wikimediamesages.styles":"ready","ext.visualEditor.desktopArticleTarget.noscript":"ready","ext.uls.interlanguage":"ready","wikibase.client.init":"ready","ext.wikimediaBadges":"ready"};RLPAGEMODULES=["ext.cite.ux-enhancements","ext.scribunto.logs","site","mediawiki.page.ready","jquery.makeCollapsible","mediawiki.toc","skins.vector.js","ext.centralNotice.geoIP","ext.centralNotice.startUp","ext.gadget.ReferenceTooltips","ext.gadget.switcher","ext.urlShortener.toolbar","ext.centralauth.centralautologin","mmv.bootstrap","ext.popups","ext.visualEditor.desktopArticleTarget.init","ext.visu
```

```
In [ ]: 1
```

