

Code Frisk

Generated by Doxygen 1.8.20

1 Hierarchical Index	1
1.1 Class Hierarchy	1
2 Class Index	3
2.1 Class List	3
3 Namespace Documentation	5
3.1 models Namespace Reference	5
3.1.1 Function Documentation	5
3.1.1.1 csv_file()	5
3.1.1.2 preprocessing()	6
3.1.1.3 remove_comments()	7
3.1.1.4 remove_comments_pythonfile()	7
3.1.1.5 remove_macros()	8
3.1.1.6 remove_redundant_functions()	9
3.1.1.7 similarity()	10
3.1.1.8 tf_idf()	10
3.1.1.9 txt_file()	11
3.1.1.10 visualizer()	12
3.2 views Namespace Reference	12
3.2.1 Function Documentation	12
3.2.1.1 create_user()	13
4 Class Documentation	15
4.1 models.add_data Class Reference	15
4.1.1 Detailed Description	15
4.1.2 Member Data Documentation	15
4.1.2.1 data	15
4.1.2.2 label	16
4.1.2.3 username	16
4.2 views.DataUpload Class Reference	16
4.2.1 Member Function Documentation	16
4.2.1.1 post()	17
4.2.1.2 view_files()	17
4.2.2 Member Data Documentation	17
4.2.2.1 authentication_classes	17
4.2.2.2 parser_classes	18
4.2.2.3 permission_classes	18
4.2.2.4 serializer_class	18
4.3 views.UserViewSet Class Reference	18
4.3.1 Member Data Documentation	18
4.3.1.1 authentication_classes	18
4.3.1.2 permission_classes	19

4.3.1.3 queryset	19
4.3.1.4 serializer_class	19

Index	21
------------------------------	-----------

Chapter 1

Hierarchical Index

1.1 Class Hierarchy

This inheritance list is sorted roughly, but not completely, alphabetically:

Model	
models.add_data	15
ModelViewSet	
views.DataUpload	16
views.UserViewSet	18

Chapter 2

Class Index

2.1 Class List

Here are the classes, structs, unions and interfaces with brief descriptions:

models.add_data	15
views.DataUpload	16
views.UserViewSet	18

Chapter 3

Namespace Documentation

3.1 models Namespace Reference

Classes

- class [add_data](#)

Functions

- def [remove_redundant_functions](#) (content)
- def [remove_comments](#) (string)
- def [visualizer](#) (list_of_files, similarity_matrix)
- def [remove_macros](#) (content)
- def [remove_comments_pythonfile](#) (file_content)
- def [preprocessing](#) (list_of_paths, list_of_files)
- def [tf_idf](#) (word_count_in_each_file, word_count_across_documents, list_of_paths, list_of_files)
- def [txt_file](#) (similarity_matrix, list_of_paths, list_of_files)
- def [csv_file](#) (list_of_files, similarity_matrix)
- def [similarity](#) (s, t)

3.1.1 Function Documentation

3.1.1.1 csv_file()

```
def models.csv_file (
    list_of_files,
    similarity_matrix )
```

Arguments :

- list_of_files :list of source code file names
- similarity_matrix:2-dimensional matrix representing mutual similarity between each pair of files

Functionality:

Plotting the output similarity_matrix and saving it as an csv file

```

296 def csv_file(list_of_files,similarity_matrix):
297     """ Interpreting the Output data as a CSV file ,
298     #where each element represent the percentage matching between the file
299     #corresponding to a row and column"""
300     """
301     Arguments      :
302         list_of_files      :list of source code file names
303         similarity_matrix:2-dimensional matrix representing mutual similarity between each pair of files
304     Functionality:
305         Plotting the output similarity_matrix and saving it as an csv file
306     """
307     f=similarity_matrix.tolist()
308     files=['filenames']+list_of_files
309     for x in range(len(list_of_files)):
310         f[x]=[list_of_files[x]]+f[x]
311     f=[files]+f
312     with open("media/result.csv", "w+") as myCsv:
313         csvWriter = csv.writer(myCsv, delimiter=',')
314         csvWriter.writerows(f)
315         visualizer(list_of_files,similarity_matrix)

```

3.1.1.2 preprocessing()

```

def models.preprocessing (
    list_of_paths,
    list_of_files )

```

Core logic is based on the Bag_of_Words
and TF-IDF Strategy(Term frequency and Inverse Document Frequency)

Arguments :
list_of_paths :list of source code paths
list_of_files :it consists of list of file names
Functionality:
It finds the count of each word after removing comments and replacing macros and passes this vector to tf_idf

```

167 def preprocessing(list_of_paths,list_of_files):
168     """
169     Core logic is based on the Bag_of_Words
170     and TF-IDF Strategy( Term frequency and Inverse Document Frequency )
171     Arguments      :
172         list_of_paths      :list of source code paths
173         list_of_files      :it consists of list of file names
174     Functionality:
175         It finds the count of each word after removing comments and replacing macros and passes this
176         vector to tf_idf function.
177     """
178     word_count_across_documents={}
179     """ Maintains the word count of each element in a document """
180     word_count_in_each_file=[]
181     for files in list_of_paths:
182         filename='media/'+files
183         temp={}
184         myfile=open(filename,"r")
185         content=myfile.read()
186         myfile.close()
187         if(files[-4:]==' .cpp'):
188             content=remove_redundant_functions(content)
189             content=remove_macros(content)
190             content=remove_comments(content)
191         if(files[-5:]==' .java'):
192             content=remove_comments(content)
193         elif(files[-3:]==' .py'):
194             content=remove_comments_pythonfile(content)
195     sym=[" ", ",", ";", ":", "{", "}", "(", ")", "[", "]", "+", "-", "*", "/", "%", "&", "^", "!", "=", "<", ">", "?", "'", '"', '#', '.', '']
196     for i in sym:
197         content=content.replace(i, " "+i+" ")

```

```

203 if (files[-4:]==' .cpp' or files[-5:]==' .java') :
204     content=content.replace("while","for")
205     content=content.replace("switch","if")
206     content=content.replace("case","else if")
207     content=content.replace("default","else")
208     content=content.replace("unsigned long long int","double")
209     content=content.replace("unsigned long long","double")
210     content=content.replace("long long int","double")
211     content=content.replace("long long","double")
212     content=content.replace("float","double")
213     content=content.replace("int","double")
214     content=content.replace("for","double")
215     content=content.replace("+","+ = 1")
216     content=content.replace("- -","- = 1")
217     content=content.replace("< <","<<")
218     content=content.replace('> >','>>')
219     cont=""
220
221     for i in content.split('\n'):
222         if(i!='\n'):
223             continue
224         i=i.strip()
225         if(len(i)==0 or i[0]!='#'):
226             continue
227         else:
228             cont=cont+' '+i
229     content=cont
230
231 elif(files[-3:]==' .py') :
232     content=content.replace('while','for')
233     content=content.replace('switch','if')
234     content=content.replace('case','elif')
235     content=content.replace('default','else')
236     content=content.replace('do','')
237     content=re.sub(':|\'|\\"',",",content)
238 List_of_words=content.split()
239
240 for i in List_of_words:
241     temp[i]=temp.get(i,0)+1
242     word_count_across_documents[i]=word_count_across_documents.get(i,0)+1
243
244     word_count_in_each_file.append(temp)
245 tf_idf(word_count_in_each_file,word_count_across_documents,list_of_paths,list_of_files)
246

```

3.1.1.3 remove_comments()

```

def models.remove_comments (
    string )
66 def remove_comments(string):
67     pattern = r"(\".*?\"|'.'*?'|\/.*?\/|\/[^\r\n]*$)"
68     # first group captures quoted strings (double or single)
69     # second group captures comments (//single-line or /* multi-line */)
70     regex = re.compile(pattern, re.MULTILINE|re.DOTALL)
71     def _replacer(match):
72         # if the 2nd group (capturing comments) is not None,
73         # it means we have captured a non-quoted (real) comment string.
74         if match.group(2) is not None:
75             return "" # so we will return empty to remove the comment
76         else: # otherwise, we will return the 1st group
77             return match.group(1) # captured quoted-string
78     return regex.sub(_replacer, string)
79

```

3.1.1.4 remove_comments_pythonfile()

```
def models.remove_comments_pythonfile (
    file_content )
```



```

146         string=""
147         for i in range(1, len(j)-1):
148             string+=j[i]+' '
149         content=content.replace(last_word, string)
150     return content
151

```

3.1.1.6 remove_redundant_functions()

```

def models.remove_redundant_functions (
    content )

```

Arguments:

content: string storing the source code

Return type: updated string

Functionality:

redundant functions in a source code are removed

Logic Used:

Write a minimal parser that can identify functions.

It just needs to detect the start and ending line of a function.

Programmatically comment out the first function, save to a temp file.

Try to compile the file by invoking the compiler.

Check if there are compile errors, if yes, the function is called, if not, it is unused.

Continue with the next function.

Reference:<https://stackoverflow.com/questions/33209302/removal-of-unused-or-redundant-code>

```

24 def remove_redundant_functions(content):
25     """
26     Arguments:
27         content: string storing the source code
28     Return type: updated string
29     Functionality:
30         redundant functions in a source code are removed
31     Logic Used:
32         Write a minimal parser that can identify functions.
33         It just needs to detect the start and ending line of a function.
34         Programmatically comment out the first function, save to a temp file.
35         Try to compile the file by invoking the compiler.
36         Check if there are compile errors, if yes, the function is called, if not, it is unused.
37         Continue with the next function.
38         Reference:https://stackoverflow.com/questions/33209302/removal-of-unused-or-redundant-code
39     """
40     type_list=['int', 'void', 'char', 'string']
41     t=""
42     for type in type_list:
43         L=re.findall(type+"s*[a-z0-9_]\s\([a-z0-9_ \n\t, \r\f\v]\)\s\{", content)
44         for y in L:
45             y=y.replace("(", "\(")
46             y=y.replace(")", "\)")
47             x=re.search(y, content)
48             first=x.span()[0]
49             last=x.span()[1]
50             count=1
51             while(count!=0):
52                 if(content[last]=='{'):
53                     count+=1
54                 if(content[last]=='}')':
55                     count-=1
56                 last+=1
57             t=content[0:first]+content[last:]
58             temp = open("temp.cpp", "w")
59             temp.write(t)
60             temp.close()
61             g = subprocess.getstatusoutput("g++ temp.cpp")
62             if(g[1]==""):
63                 content=t
64     return content
65

```

3.1.1.7 similarity()

```
def models.similarity (
    s,
    t )
```

Arguments :

- s : (sorted) list of numbers
- t : (sorted) list of numbers

Return type :

- returns a number in the range(0,1)

Functionality:

- Evaluates the cosine product of the two vectors

```
317 def similarity(s,t):
318     """
319     Arguments :
320         s : (sorted) list of numbers
321         t : (sorted) list of numbers
322     Return type :
323         returns a number in the range(0,1)
324     Functionality:
325         Evaluates the cosine product of the two vectors
326     """
327     x=np.zeros(abs(s.size-t.size))
328     ""
329     if(s.size>t.size):
330         t=np.concatenate((x,t))
331     else:
332         s=np.concatenate((x,s))
333     ""
334     x=min(s.size,t.size)
335     s=s[-x:]
336     t=t[-x:]
337
338
339     #s=(s-np.mean(s))/np.std(s)
340     #t=(t-np.mean(t))/np.std(t)
341     return 1-np.linalg.norm(s-t)/(np.linalg.norm(s)+np.linalg.norm(t))
342     return np.dot(s,t)/(np.linalg.norm(s)*np.linalg.norm(t))
```

3.1.1.8 tf_idf()

```
def models.tf_idf (
    word_count_in_each_file,
    word_count_across_documents,
    list_of_paths,
    list_of_files )
```

Arguments :

- list_of_paths : list of source code files
- list_of_files : It consists data of all the Users who have been SignedUp
- word_count_in_each_file : Frequency of word corresponding to each file as an array of dictionary
- word_count_across_documents: Frequency of each word across as files corresponding to a particular assignment

Functionality:

- It computes tf_idf vector corresponding to each file.
- The tf_idf function is somewhat different from the original one
- If we use the bag of words strategy then similarity is determined mostly by the variables which have maximum frequency.
- But similarity should depend more on core logic like number of functions, operators loops etc.
- The weight added for each word say 'x' in file 'f' is $\log(\text{freq of } x \text{ across all files corresponding to assignment } f)$
- Words which have low frequency in a file than average frequency across all files are given +ve weightage
- Words which have high frequency in a file than average frequency across all files are given -ve weightage
- Uniqueness is determined by high weightage words.

```

247 def tf_idf(word_count_in_each_file, word_count_across_documents, list_of_paths, list_of_files):
248     """ Arguments :
249         list_of_paths      :list of source code files
250         list_of_files      :It consists data of all the Users who have been SignedUp
251         word_count_in_each_file :Frequency of word corresponding to each file as an array of
dictionary
252         word_count_across_documents:Frequency of each word across as files corresponding to a particular
assignment as dictionary
253     Functionality:
254         It computes tf_idf vector corresponding to each file.
255         The tf_idf function is somewhat different from the original one
256         If we use the bag of words strategy then similarity is determined mostly by the variables which
have maximum count in a file.
257         But similarity should depend more on core logic like number of functions, operators loops etc.
258         The weight added for each word say 'x' in file 'f' is log(freq of x across all files
corresponding to assignment/(freq of x in f*number of files))
259         Words which have low frequency in a file than average frequency across all files are given +ve
weightage
260         Words which have high frequency in a file than average frequency across all files are given -ve
weightage
261         Uniqueness is determined by high weightage words.
262     """
263     similarity_matrix=np.zeros((len(list_of_paths),len(list_of_paths)))
264     tf_idf_vec=[]
265     for i in range(len(list_of_paths)):
266         temp=[]
267         for j in word_count_in_each_file[i]:
268
269             temp.append(word_count_in_each_file[i].get(j)*((math.log(word_count_across_documents.get(j)/word_count_in_each_file[i]
temp.sort()
270             tf_idf_vec.append(temp)
271
272     for i in range(len(list_of_paths)):
273         similarity_matrix[i,i]=0
274         for j in range(i+1,len(list_of_paths)):
275             similarity_matrix[i,j]=similarity(np.array(tf_idf_vec[i]),np.array(tf_idf_vec[j]))
276             similarity_matrix[j,i]=similarity_matrix[i,j]
277
278     txt_file(similarity_matrix,list_of_paths,list_of_files)
279

```

3.1.1.9 txt_file()

```

def models.txt_file (
    similarity_matrix,
    list_of_paths,
    list_of_files )

```

Arguments :

- list_of_paths :list of source code file names
- similarity_matrix:2-dimensional matrix representing mutual similarity between each pair of files
- list_of_files :It consists data of all the Users who have been SignedUp

Functionality:

Displaying the Percentage matching among files in text format and saving it as a csv file

```

280 def txt_file(similarity_matrix,list_of_paths,list_of_files):
281     """
282     Arguments :
283         list_of_paths      :list of source code file names
284         similarity_matrix:2-dimensional matrix representing mutual similarity between each pair of files
285         list_of_files      :It consists data of all the Users who have been SignedUp
286     Functionality:
287         Displaying the Percentage matching among files in text format and saving it as a csv file """
288     result=open("media/result.txt","w")
289     res=""
290     for i in range(len(list_of_paths)):
291         for j in range(i+1,len(list_of_paths)):
292             res+="similarity between "+ list_of_files[i]+" and "+list_of_files[j]+" =
"+str(similarity_matrix[i][j])+"\n"
293     result.write(res)
294     csv_file(list_of_files,similarity_matrix)
295

```

3.1.1.10 visualizer()

```
def models.visualizer (
    list_of_files,
    similarity_matrix )
```

Arguments :

- list_of_files :list of source code file names
- similarity_matrix:2-dimensional matrix representing mutual similarity between each pair of files

Return type :

- Path of the saved image

Functionality:

- Plotting the output similarity_matrix and saving it as an image

```
80 def visualizer(list_of_files,similarity_matrix):
81     """
82     Arguments      :
83         list_of_files      :list of source code file names
84         similarity_matrix:2-dimensional matrix representing mutual similarity between each pair of files
85     Return type    :
86         Path of the saved image
87     Functionality:
88         Plotting the output similarity_matrix and saving it as an image """
89     x=range(len(list_of_files))
90     y=range(len(list_of_files))
91     xx,yy=np.meshgrid(x,y)
92     z=similarity_matrix[xx,yy]
93     z=np.round(z*100)
94     cmap = matplotlib.colors.LinearSegmentedColormap.from_list("", ["green","yellow","red"])
95     fig=plt.figure()
96     ax=fig.add_subplot(111)
97     im=ax.matshow(z,cmap=cmap,vmin=0,vmax=100,origin='lower')
98     for i in range(len(list_of_files)):
99         for j in range(i+1,len(list_of_files)):
100             ax.text(j, i, int(similarity_matrix[i,j]*100),ha="center", va="center", color="k")
101             ax.text(i, j, int(similarity_matrix[i,j]*100),ha="center", va="center", color="k")
102     fig.colorbar(im,shrink=0.5)
103     ax.set_xticks(range(len(list_of_files)))
104     ax.set_yticks(range(len(list_of_files)))
105     ax.set_xticklabels(list_of_files,rotation=90)
106     ax.set_yticklabels(list_of_files)
107     random=np.random.randint(1,100)
108     path='result.png'
109     plt.savefig('media/'+path)
110
```

3.2 views Namespace Reference

Classes

- class [DataUpload](#)
- class [UserViewSet](#)

Functions

- def [create_user](#) (request)

3.2.1 Function Documentation

3.2.1.1 create_user()

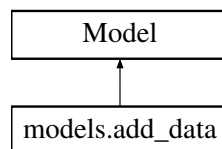
```
def views.create_user (
    request )
20 def create_user(request):
21     validated_data=JSONParser().parse(request)
22     if(validated_data['password']!=validated_data['password2']):
23         return JsonResponse("Passwords donot match",safe=False)
24     elif(User.objects.filter(email=validated_data['email']).exists()):
25         return JsonResponse("There exists an account with this email",safe=False)
26     elif(User.objects.filter(username=validated_data['username']).exists()):
27         return JsonResponse("There exists an account with this username",safe=False)
28     else:
29         try:
30             new_data={}
31             new_data['username']=validated_data['username']
32             new_data['password']=validated_data['password']
33             new_data['email']=validated_data['email']
34         except:
35             return JsonResponse("failed",safe=False)
36     user_serializer=UserSerializer(data=new_data)
37     if user_serializer.is_valid():
38         user_serializer.save()
39         return JsonResponse("success",safe=False)
40     return JsonResponse("Invalid credentials",safe=False)
41
```


Chapter 4

Class Documentation

4.1 models.add_data Class Reference

Inheritance diagram for models.add_data:



Static Public Attributes

- `username` = `models.CharField(max_length=50)`
- `label` = `models.CharField(max_length=50)`
- `data` = `models.FileField()`

4.1.1 Detailed Description

4.1.2 Member Data Documentation

4.1.2.1 data

```
models.add_data.data = models.FileField() [static]
```

4.1.2.2 label

```
models.add_data.label = models.CharField(max_length=50) [static]
```

4.1.2.3 username

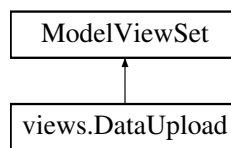
```
models.add_data.username = models.CharField(max_length=50) [static]
```

The documentation for this class was generated from the following file:

- models.py

4.2 views.DataUpload Class Reference

Inheritance diagram for views.DataUpload:



Public Member Functions

- def [post](#) (request, *args, **kwargs)
- def [view_files](#) (request)

Static Public Attributes

- tuple [authentication_classes](#) = (TokenAuthentication,)
- tuple [permission_classes](#) = (IsAuthenticated,)
- tuple [parser_classes](#) = (MultiPartParser, FormParser)
- [serializer_class](#) = DataSerializer

4.2.1 Member Function Documentation

4.2.1.1 post()

```
def views.DataUpload.post (
    request,
    * args,
    ** kwargs )
48     def post(request, *args, **kwargs):
49         label=request.POST['label']
50         username=request.POST['username']
51         your_file = request.FILES['data']
52         add_data.objects.create( username=username,label=label,data=your_file)
53         if(your_file.name.find('.tar')!=-1):
54             t=tarfile.open('media/'+str(your_file),'r')
55             L=t.getnames()
56             t.extractall(path='media/')
57         else:
58             pass
59         return JsonResponse("success",safe=False)
60
```

4.2.1.2 view_files()

```
def views.DataUpload.view_files (
    request )
62     def view_files(request):
63         data = request.POST
64         your_files=add_data.objects.filter(label=data['label'])
65         if(not your_files.exists()):
66             return JsonResponse("You dont have any files",safe=False)
67         l=[]
68         h=[]
69         for i in your_files:
70             if(i.data.name.find('.tar')!=-1):
71                 t=tarfile.open('media/'+str(i.data.name),'r')
72                 L=t.getnames()
73                 for k in L:
74                     try:
75                         if(k.find('.')!=-1 or k.find('.DS_Store')!=-1):
76                             continue
77                         h.append(k.split('/')[1])
78                         l.append(k)
79                     except:
80                         pass
81             else:
82                 l.append(i.data.name)
83                 h.append(i.data.name)
84         preprocessing(l,h)
85         data={'png':'result.png','txt':'result.txt','csv':'result.csv'}
86         return JsonResponse(data, safe=False)
```

4.2.2 Member Data Documentation

4.2.2.1 authentication_classes

```
tuple views.DataUpload.authentication_classes = (TokenAuthentication,) [static]
```

4.2.2.2 parser_classes

```
tuple views.DataUpload.parser_classes = (MultiPartParser, FormParser) [static]
```

4.2.2.3 permission_classes

```
tuple views.DataUpload.permission_classes = (IsAuthenticated,) [static]
```

4.2.2.4 serializer_class

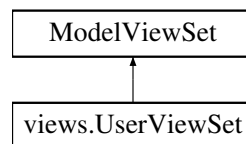
```
views.DataUpload.serializer_class = DataSerializer [static]
```

The documentation for this class was generated from the following file:

- views.py

4.3 views.UserViewSet Class Reference

Inheritance diagram for views.UserViewSet:



Static Public Attributes

- `queryset` = `User.objects.all().order_by("username")`
- `serializer_class` = `UserSerializer`
- tuple `authentication_classes` = `(TokenAuthentication,)`
- list `permission_classes` = `[IsAuthenticated]`

4.3.1 Member Data Documentation

4.3.1.1 authentication_classes

```
tuple views.UserViewSet.authentication_classes = (TokenAuthentication,) [static]
```

4.3.1.2 permission_classes

```
list views.UserViewSet.permission_classes = [IsAuthenticated] [static]
```

4.3.1.3 queryset

```
views.UserViewSet.queryset = User.objects.all().order_by('username') [static]
```

4.3.1.4 serializer_class

```
views.UserViewSet.serializer_class = UserSerializer [static]
```

The documentation for this class was generated from the following file:

- views.py

Index

- authentication_classes
 - views.DataUpload, [17](#)
 - views.UserViewSet, [18](#)
- create_user
 - views, [12](#)
- csv_file
 - models, [5](#)
- data
 - models.add_data, [15](#)
- label
 - models.add_data, [15](#)
- models, [5](#)
 - csv_file, [5](#)
 - preprocessing, [6](#)
 - remove_comments, [7](#)
 - remove_comments_pythonfile, [7](#)
 - remove_macros, [8](#)
 - remove_redundant_functions, [9](#)
 - similarity, [9](#)
 - tf_idf, [10](#)
 - txt_file, [11](#)
 - visualizer, [11](#)
- models.add_data, [15](#)
 - data, [15](#)
 - label, [15](#)
 - username, [16](#)
- parser_classes
 - views.DataUpload, [17](#)
- permission_classes
 - views.DataUpload, [18](#)
 - views.UserViewSet, [18](#)
- post
 - views.DataUpload, [16](#)
- preprocessing
 - models, [6](#)
- queryset
 - views.UserViewSet, [19](#)
- remove_comments
 - models, [7](#)
- remove_comments_pythonfile
 - models, [7](#)
- remove_macros
 - models, [8](#)
- remove_redundant_functions
 - models, [9](#)
- serializer_class
 - views.DataUpload, [18](#)
 - views.UserViewSet, [19](#)
- similarity
 - models, [9](#)
- tf_idf
 - models, [10](#)
- txt_file
 - models, [11](#)
- username
 - models.add_data, [16](#)
- view_files
 - views.DataUpload, [17](#)
- views, [12](#)
 - create_user, [12](#)
- views.DataUpload, [16](#)
 - authentication_classes, [17](#)
 - parser_classes, [17](#)
 - permission_classes, [18](#)
 - post, [16](#)
 - serializer_class, [18](#)
 - view_files, [17](#)
- views.UserViewSet, [18](#)
 - authentication_classes, [18](#)
 - permission_classes, [18](#)
 - queryset, [19](#)
 - serializer_class, [19](#)
- visualizer
 - models, [11](#)