

starter

November 29, 2018

```
In [20]: from urllib.request import urlopen
import re
from collections import defaultdict
import numpy as np
import chardet
import pandas as pd
import csv
import json
from bs4 import BeautifulSoup
import matplotlib.pyplot as plt
%matplotlib inline

In [2]: def get_variable_names():
    #contains all variables corresponding to a string
    varNames_dict = {}

    # Grab the page from the web
    f = urlopen("http://scdb.wustl.edu/documentation.php?var=chief")
    html = str(f.read())

    reviews = html.split('<li class="toc">')[1:]
    for variable in reviews:
        str_1 = variable.split("</a>")[0].split('class="toc')
        if (re.split(".var=",str_1[0])[1].replace('" ','') == 'intro'):
            continue
        if (re.split(".var=",str_1[0])[1].replace('" ','') == 'cite'):
            continue
        varName_dict[str_1[1].replace('>','')] = re.split(".var=",str_1[0])[1].replace(

    return varName_dict

In [3]: def get_all_var_values(varName_dict):
    #contains all values possible for a var
    varValues_dict = {}

    for var in varName_dict.keys():
        #print("-----")
```

```

# print(str(var)+":"+str(varName_dict[var]))
f = urlopen("http://scdb.wustl.edu/documentation.php?var="+str(varName_dict[var]))
html = str(f.read())

# for variables which are not encoded and stright values are used
if ('<strong>Values:</strong><br />' not in html):
    continue

# print(html)
# Extract the (30) review elements
reviews = html.split('<strong>Values:</strong><br />')[1]
reviews = reviews.split('<div class="toc" onClick="toggleBlock')[0]
# print(reviews)

soup = BeautifulSoup(reviews)
# print(soup.text)

final_str = soup.text.replace("\t", "#").replace("\n", "#")
# print(final_str)
list_val = re.split("#+", final_str)
list_val = [x for x in list_val if x not in ["", " "]]

assert(len(list_val)%2 == 0 or len(list_val)%3 == 0)

# choose which format it is
pattern = re.compile("\d+")

# print(len(list_val))
if (len(list_val)>4):
    if (pattern.match(list_val[0]) and pattern.match(list_val[2]) and pattern.match(list_val[4])):
        div = 2
    elif (pattern.match(list_val[0]) and pattern.match(list_val[3]) and pattern.match(list_val[6])):
        div = 3
    else:
        print("diff pattern found")
elif (len(list_val)<4):
    if (len(list_val)%2==0):
        div = 2
    elif (len(list_val)%3==0):
        div = 3
# print(div)

varValues_dict[str(varName_dict[var])] = {}

if (div == 2):
    for x in np.arange(0, len(list_val)//2):
        # print("varValues_dict["+str(varName_dict[var])+"]["+str(list_val[x*2])+"]")
        varValues_dict[str(varName_dict[var])][list_val[x*2]] = str(list_val[x*2+1])

```

```

        elif (div == 3):
            for x in np.arange(0,len(list_val)//3):
                #print("varValues_dict["+str(varName_dict[var])+"]["+str(list_val[x*3])
                varValues_dict[str(varName_dict[var))][list_val[x*3]] = str(list_val[(x
    return varValues_dict

In [4]: #fill dicts
varName_dict = {}
varName_dict = get_variable_names()

varValues_dict = {}
varValues_dict = get_all_var_values(varName_dict)

In [5]: def fill_csv_with_var_values(varValues_dict,fileName,readLinesCount=None):
    #entireListDict has values instead of keys
    entireListDict = []

    #csv file to read
    csvfile = open(fileName, 'r')

    #fieldnames = ("caseId","sctCite")
    reader = csv.DictReader( csvfile)

    count = 0
    for row in reader:
        entireListDict.append(row)
        for column in row:
            if column in varValues_dict.keys():
                if (row[column] in varValues_dict[column].keys()):
                    entireListDict[count][column] = str(varValues_dict[column][row[column]
                else:
                    entireListDict[count][column] = str(row[column])
            else:
                entireListDict[count][column] = str(row[column])
        count += 1

    #how many lines to read
    if readLinesCount is not None:
        if count >= 1:
            break
    return entireListDict

In [6]: def write_csv_with_values(entireListDict,outputFileName):
    #write the output replacing keys with values
    with open(outputFileName,'w') as f:
        wr = csv.writer(f)
        for rowNum in np.arange(len(entireListDict)):
            if rowNum == 0:

```

```

        wr.writerow([key for key,val in entireListDict[rowNum].items()])
        wr.writerow([val for key,val in entireListDict[rowNum].items()])
    print("File "+str(outputFileName)+" is generated!")

```

```

In [7]: ##for reading csv as dict format
        # with open('sample.csv') as csv_file:
        #     csv_reader = csv.reader(csv_file, delimiter=',')
        #     line_count = 0
        #     for row in csv_reader:
        #         if line_count == 0:
        #             print(f'Column names are {", ".join(row)}')
        #             line_count += 1
        #             #break
        #         else:
        #             print(row)
        #             line_count += 1
        #     print(f'Processed {line_count} lines.')

```

```

In [8]: #test dicts
        for var in varName_dict.keys():
            print(str(var)+":"+str(varName_dict[var]))

        x = 'chief'
        for y in varValues_dict[x]:
            print(str(y)+"\t"+str(varValues_dict[x][y]))

        #get list of dicts aka csv file
        entireListDict = fill_csv_with_var_values(varValues_dict,"JusticeData.csv")
        print(len(entireListDict))
        #write_csv_with_values(entireListDict,"sample_output_2.csv")

```

```

SCDB Case ID:caseId
SCDB Docket ID:docketId
SCDB Issues ID:caseIssuesId
SCDB Vote ID:voteId
U.S. Reporter Citation:usCite
Supreme Court Citation:sctCite
Lawyers Edition Citation:ledCite
LEXIS Citation:lexisCite
Docket Number:docket
Case Name:caseName
Petitioner:petitioner
Petitioner State:petitionerState
Respondent:respondent
Respondent State:respondentState
Manner in which the Court takes Jurisdiction:jurisdiction
Administrative Action Preceeding Litigation:adminAction
Administrative Action Preceeding Litigation State:adminActionState

```

Three-Judge District Court:threeJudgeFdc
 Origin of Case:caseOrigin
 Origin of Case State:caseOriginState
 Source of Case:caseSource
 Source of Case State:caseSourceState
 Lower Court Disagreement:lcDisagreement
 Reason for Granting Cert:certReason
 Lower Court Disposition:lcDisposition
 Lower Court Disposition Direction:lcDispositionDirection
 Date of Decision:dateDecision
 Term of Court:term
 Natural Court:naturalCourt
 OnChief Justice:chief
 Date of Oral Argument:dateArgument
 Date of Reargument:dateRearg
 Issue:issue
 Issue Area:issueArea
 Decision Direction:decisionDirection
 Decision Direction Dissent:decisionDirectionDissent
 Authority for Decision 1:authorityDecision1
 Authority for Decision 2:authorityDecision2
 Legal Provisions Considered by the Court:lawType
 Legal Provision Supplement:lawSupp
 Legal Provision Minor Supplement:lawMinor
 Decision Type:decisionType
 Declaration of Unconstitutionality:declarationUncon
 Disposition of Case:caseDisposition
 Unusual Disposition:caseDispositionUnusual
 Winning Party:partyWinning
 Formal Alteration of Precedent:precedentAlteration
 Vote Not Clearly Specified:voteUnclear
 Majority Opinion Writer:majOpinWriter
 Majority Opinion Assigner:majOpinAssigner
 Split Vote:splitVote
 Majority Votes:majVotes
 Minority Votes:minVotes
 Justice ID:justice
 Justice Name:justiceName
 The Vote in the Case:vote
 Opinion:opinion
 Direction of the Individual Justice\'s Votes:direction
 Majority and Minority Voting by Justice:majority
 First Agreement:firstAgreement
 Second Agreement:secondAgreement
 1 Jay
 2 Rutledge
 3 Ellsworth
 4 Marshall

```

5      Taney
6      Chase
7      Waite
8      Fuller
9      White
10     Taft
11     Hughes
12     Stone
13     Vinson
14     Warren
15     Burger
16     Rehnquist
17     Roberts
119838

```

```
In [9]: df = pd.DataFrame(entireListDict)
```

```
In [10]: df['firstAgreement'].value_counts(normalize=True)
```

```

Out[10]:
WJBrennan      0.872495
JPStevens      0.011040
HLBlack        0.010706
0              0.010297
0              0.008203
BRWhite        0.007093
WHRehnquist    0.007093
AScalia        0.005975
HABlackmun     0.005741
TMarshall      0.005449
JHarlan2       0.005441
AMKennedy      0.004673
LFPowell       0.004673
PStewart       0.004548
SGBreyer       0.004389
SDOConnor      0.004281
FFrankfurter   0.004131
RBGinsburg     0.004039
CThomas        0.003655
DHSouter       0.003355
TCClark        0.003329
WODouglas      0.003121
JGRoberts      0.002153
SAAlito        0.001769
WEBurger       0.001679
SFReed         0.001369
RHJackson      0.001285
SSotomayor     0.001210

```

HHBurton	0.001018
EWarren	0.000851
WBRutledge	0.000826
AJGoldberg	0.000726
CEWhittaker	0.000676
AFortas	0.000651
SMinton	0.000651
FMurphy	0.000618
EKagan	0.000551
FMVinson	0.000317
NMGorsuch	0.000117

Name: firstAgreement, dtype: float64

```
In [16]: # for var in entireListDict[0].keys():
#         print("make_graphs_var.append("+str(var)+")")
make_graphs_var = []

make_graphs_var.append('decisionType')
make_graphs_var.append('term')
make_graphs_var.append('naturalCourt')
make_graphs_var.append('chief')
make_graphs_var.append('docket')

make_graphs_var.append('petitioner')
make_graphs_var.append('petitionerState')
make_graphs_var.append('respondent')
make_graphs_var.append('respondentState')
make_graphs_var.append('jurisdiction')
make_graphs_var.append('adminAction')
make_graphs_var.append('adminActionState')
make_graphs_var.append('threeJudgeFdc')
make_graphs_var.append('caseOrigin')
make_graphs_var.append('caseOriginState')
make_graphs_var.append('caseSource')
make_graphs_var.append('caseSourceState')
make_graphs_var.append('lcDisagreement')
make_graphs_var.append('certReason')
make_graphs_var.append('lcDisposition')
make_graphs_var.append('lcDispositionDirection')
make_graphs_var.append('declarationUncon')
make_graphs_var.append('caseDisposition')
make_graphs_var.append('caseDispositionUnusual')
make_graphs_var.append('partyWinning')
make_graphs_var.append('precedentAlteration')
make_graphs_var.append('voteUnclear')
make_graphs_var.append('issue')
make_graphs_var.append('issueArea')
make_graphs_var.append('decisionDirection')
```

```

make_graphs_var.append('decisionDirectionDissent')
make_graphs_var.append('authorityDecision1')
make_graphs_var.append('authorityDecision2')
make_graphs_var.append('lawType')
make_graphs_var.append('lawSupp')

```

```

make_graphs_var.append('majOpinWriter')
make_graphs_var.append('majOpinAssigner')
make_graphs_var.append('splitVote')
make_graphs_var.append('majVotes')
make_graphs_var.append('minVotes')
make_graphs_var.append('justice')
make_graphs_var.append('justiceName')
make_graphs_var.append('vote')
make_graphs_var.append('opinion')
make_graphs_var.append('direction')
make_graphs_var.append('majority')
make_graphs_var.append('firstAgreement')
make_graphs_var.append('secondAgreement')

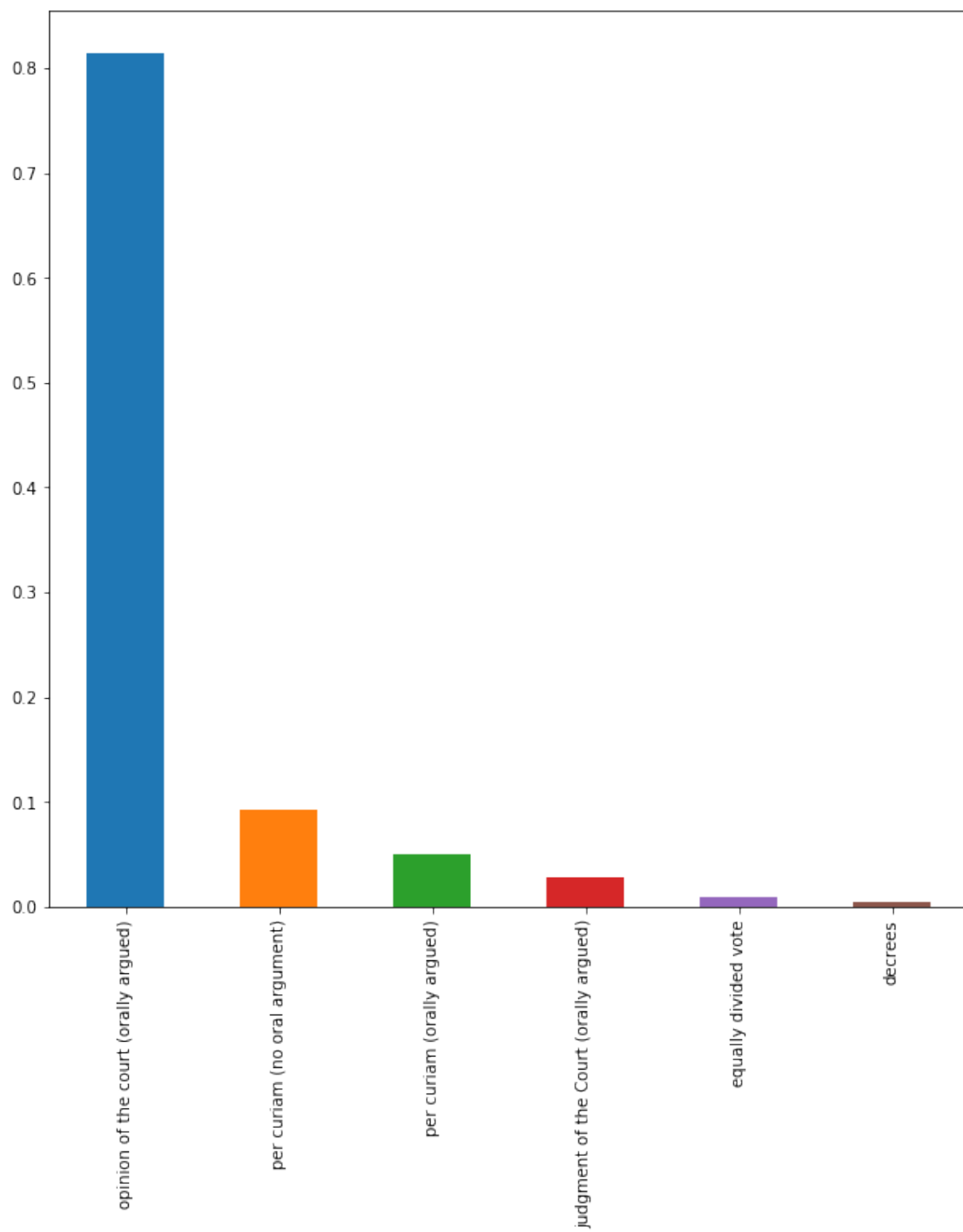
```

```

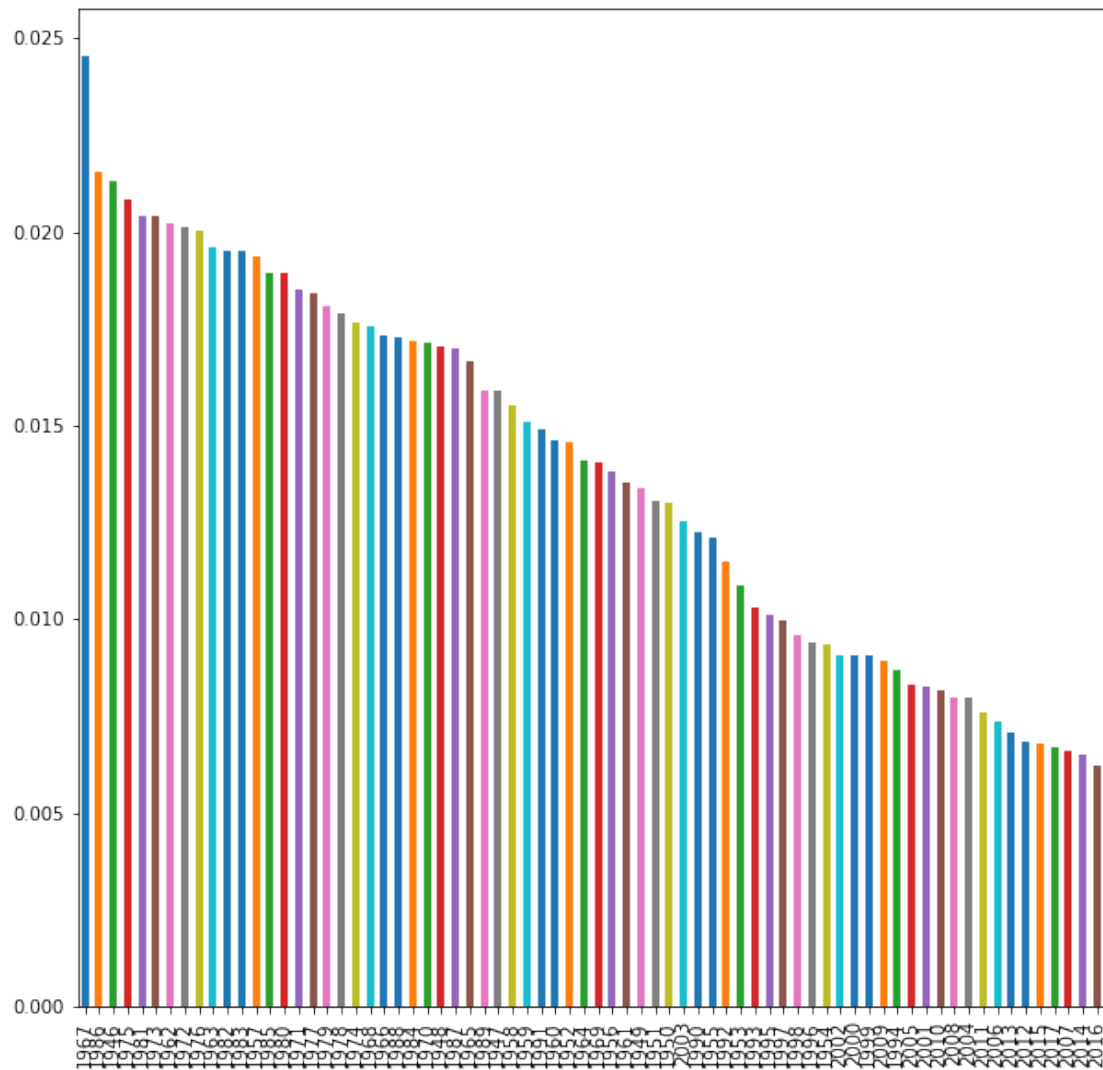
In [21]: for var in make_graphs_var:
          print(str(var))
          plot = df[var].value_counts(normalize=True).plot(kind="bar", figsize=(10, 10))
          plt.show()

```

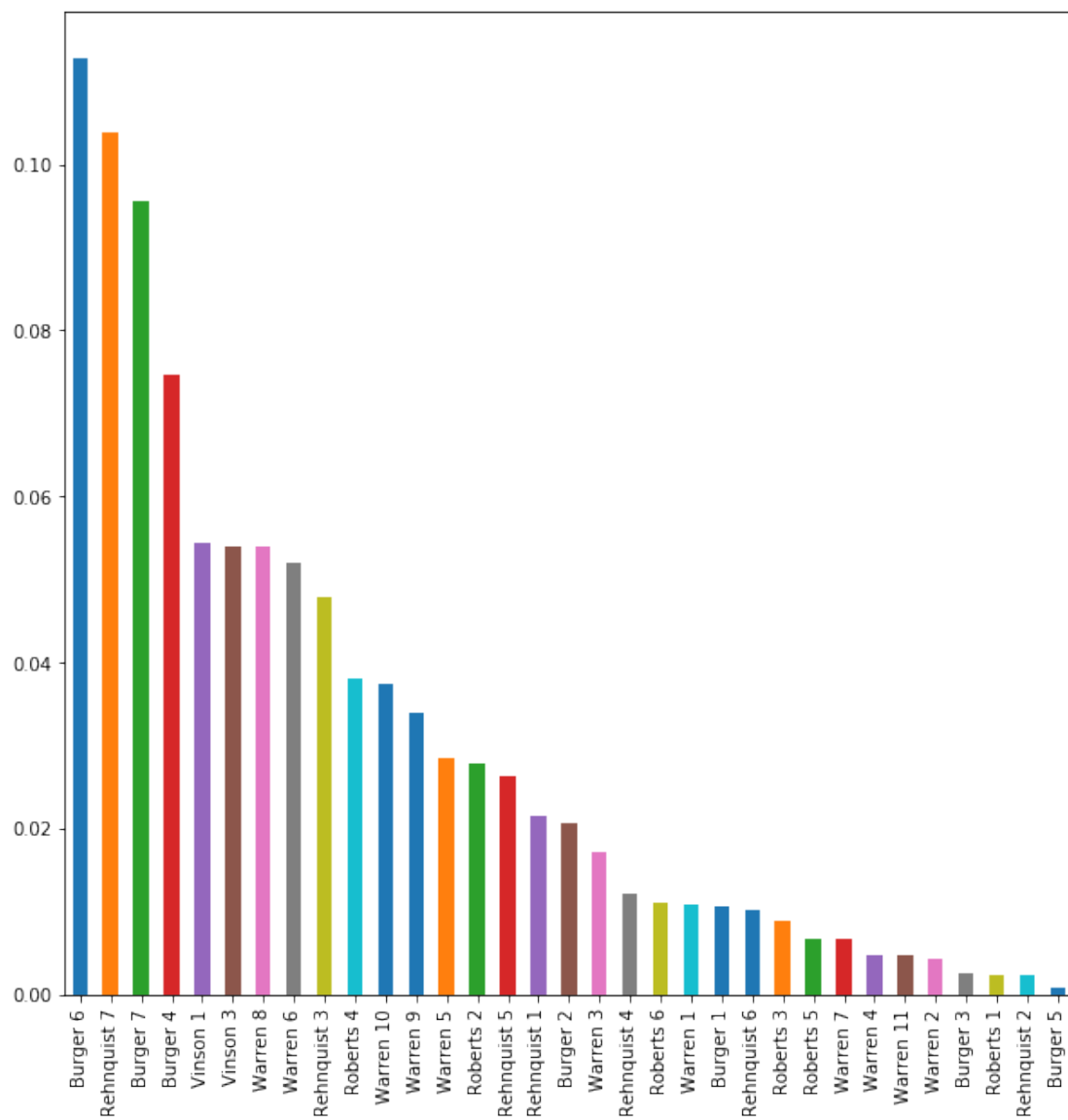
decisionType



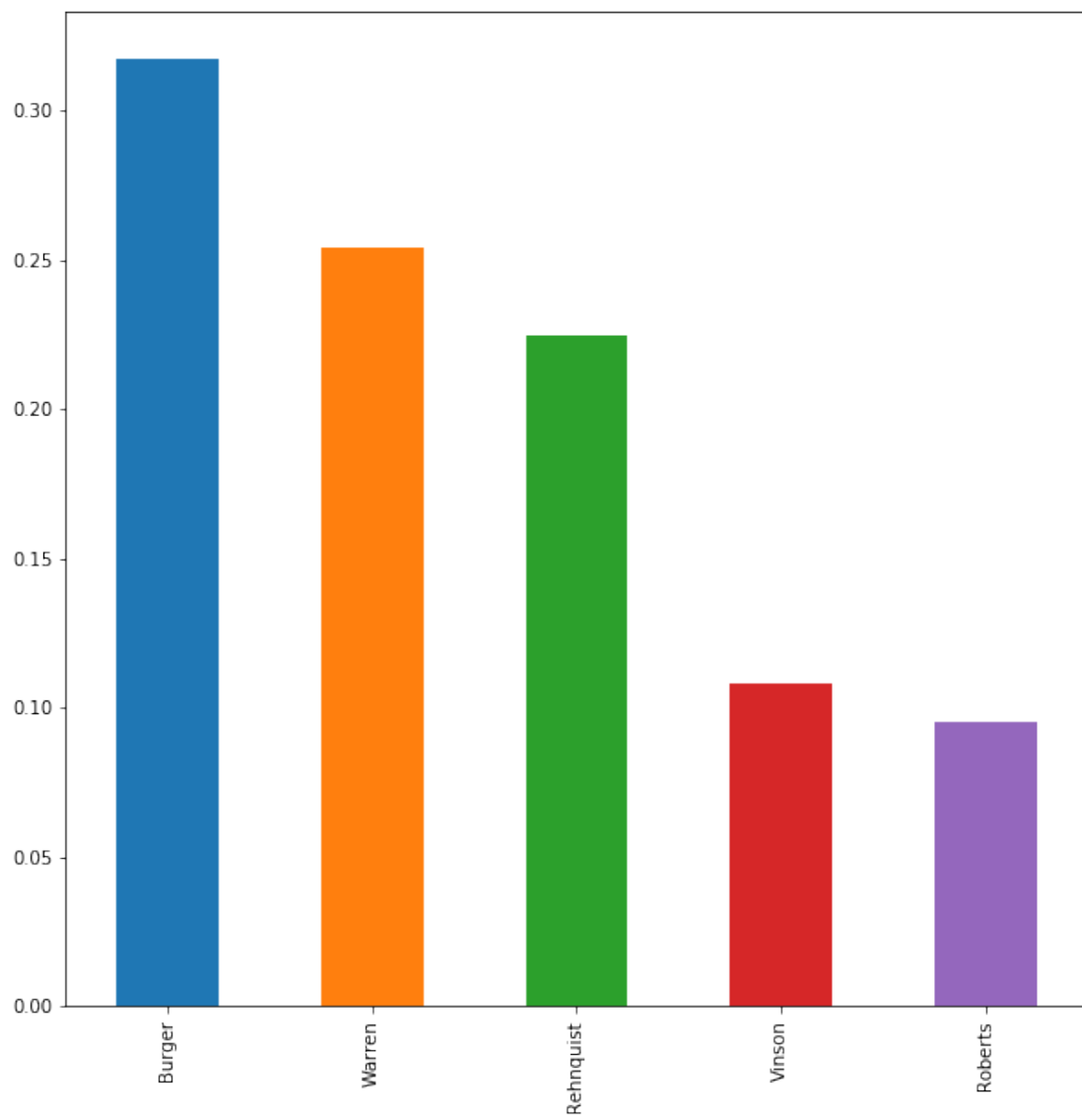
term



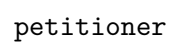
naturalCourt



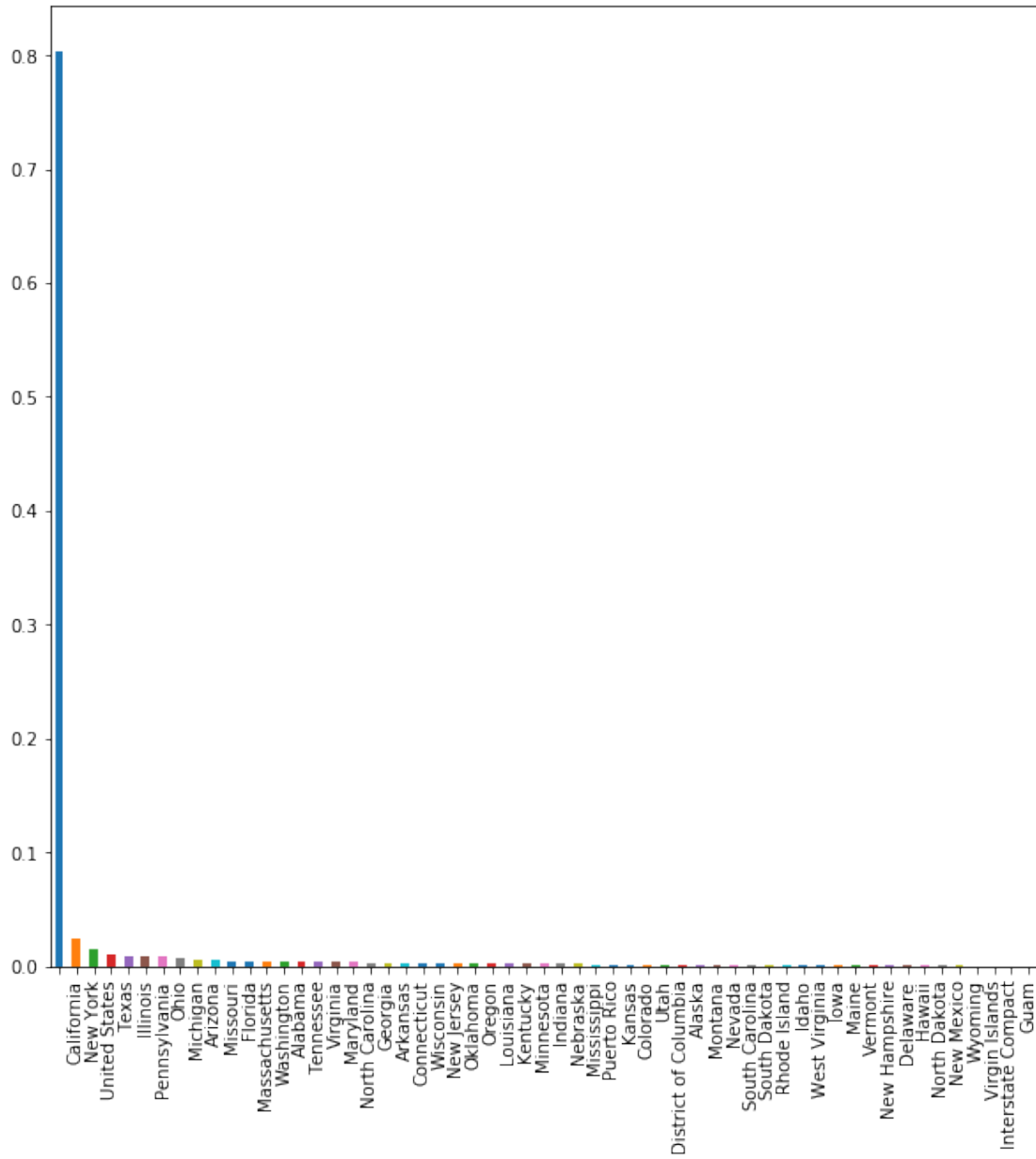
chief



docket

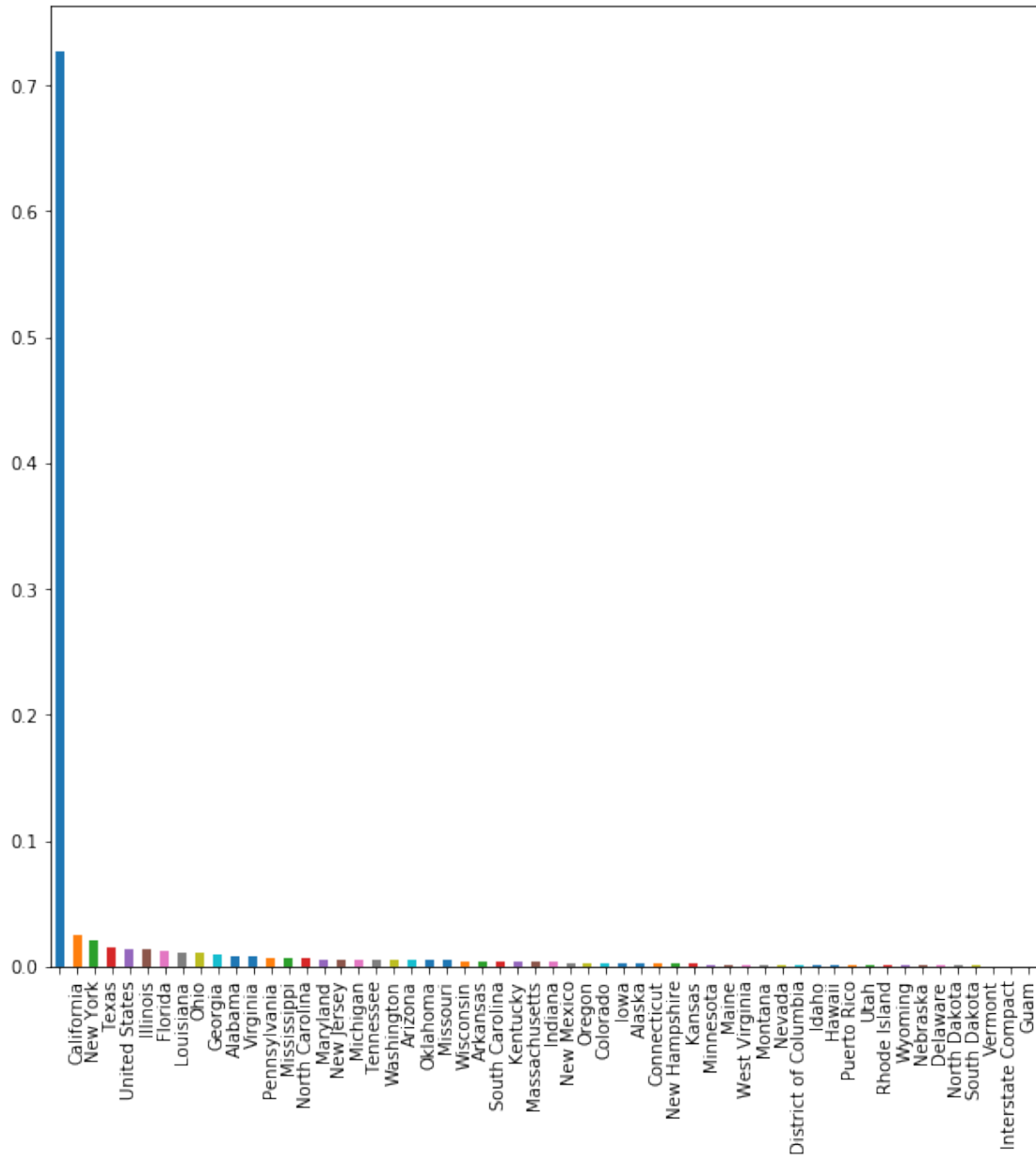


petitionerState

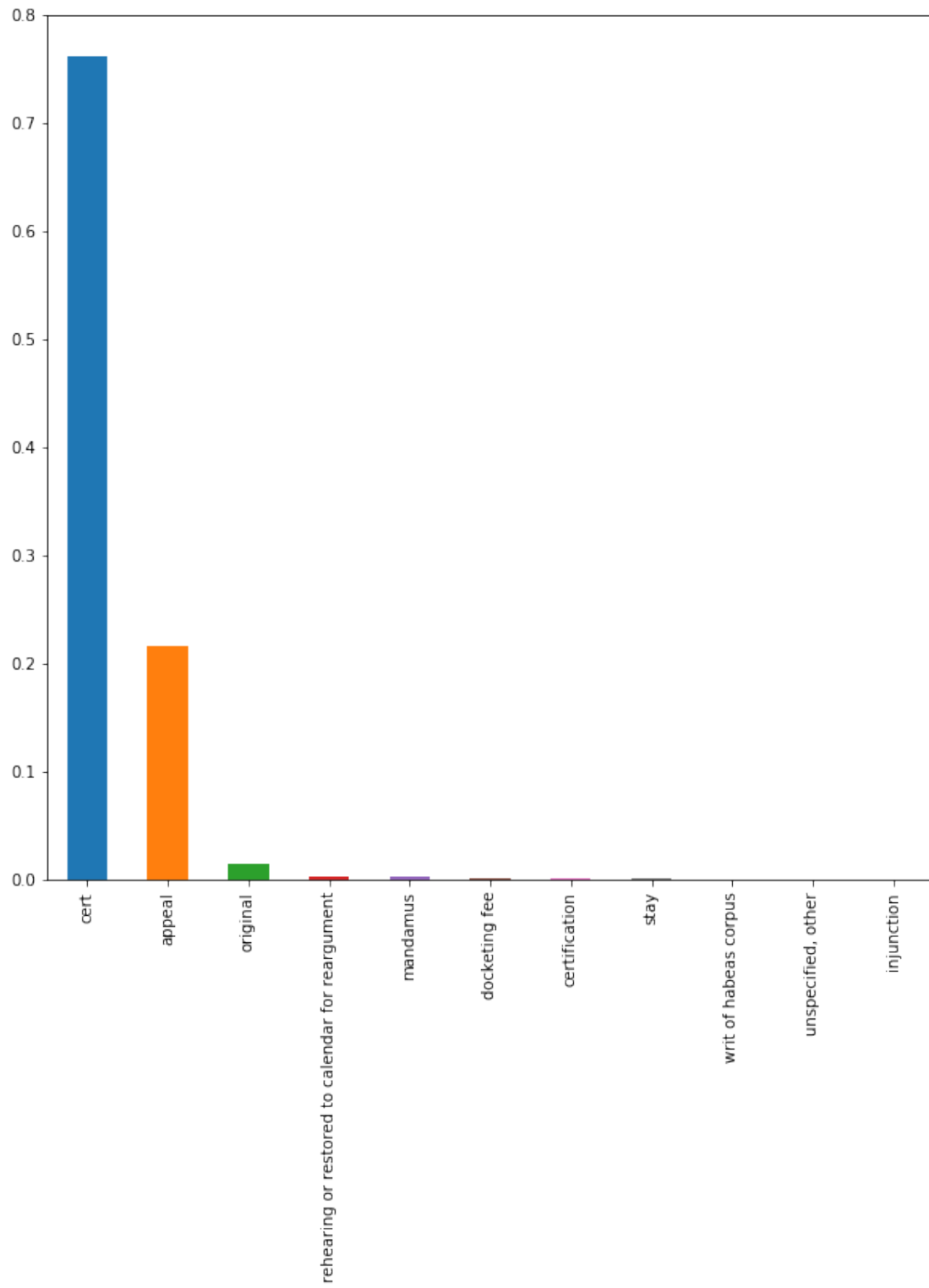


respondent

respondentState

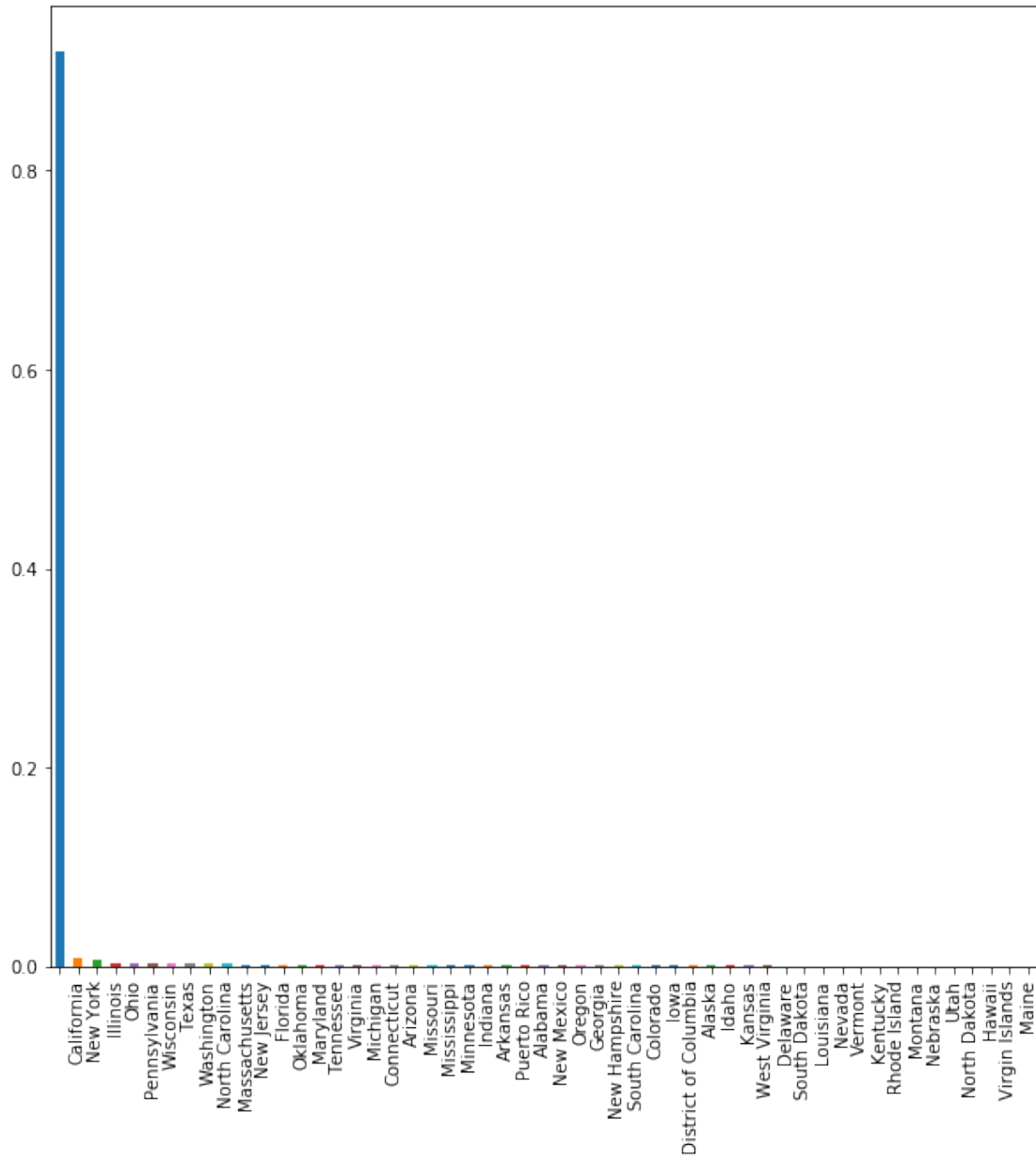


jurisdiction

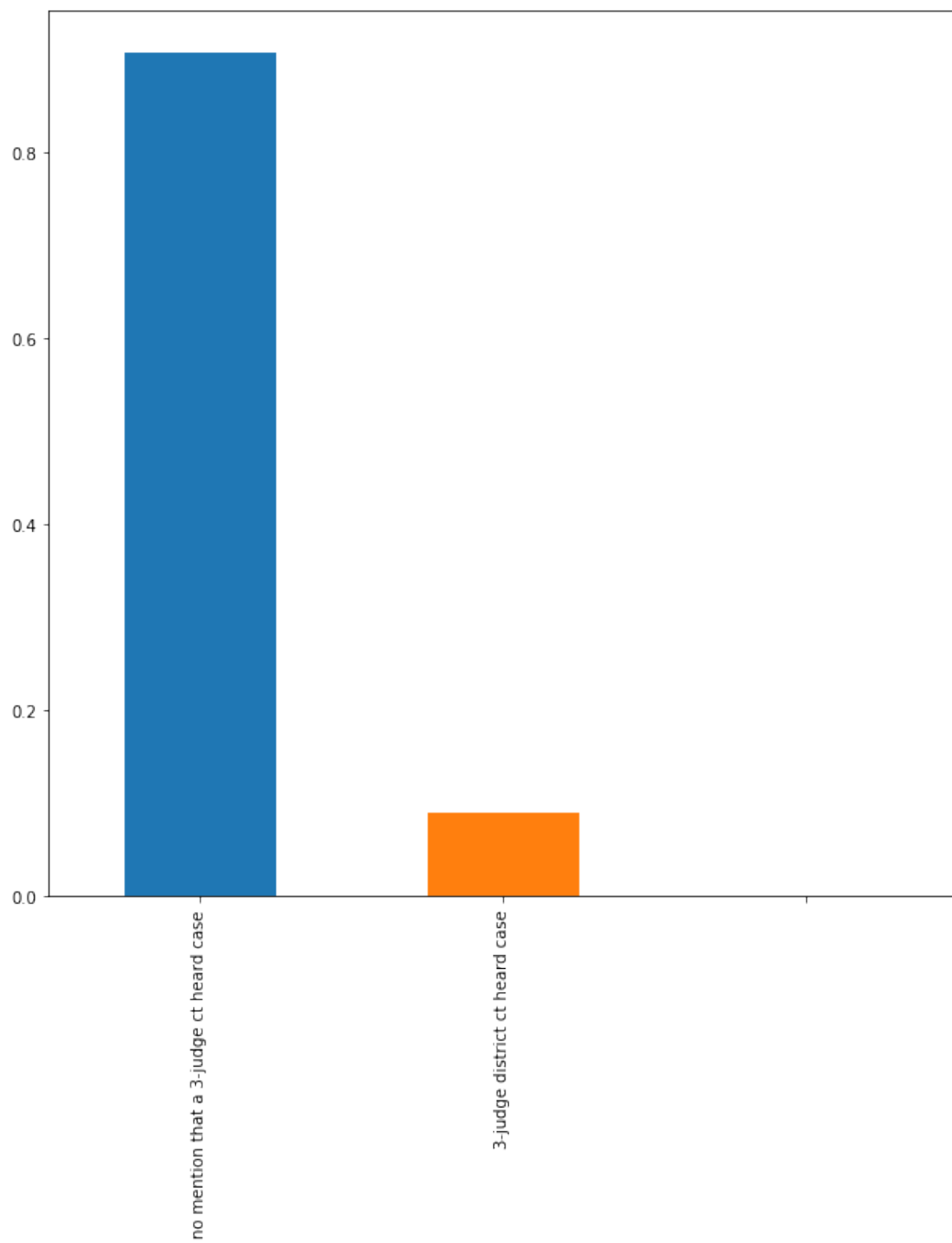


adminAction

adminActionState

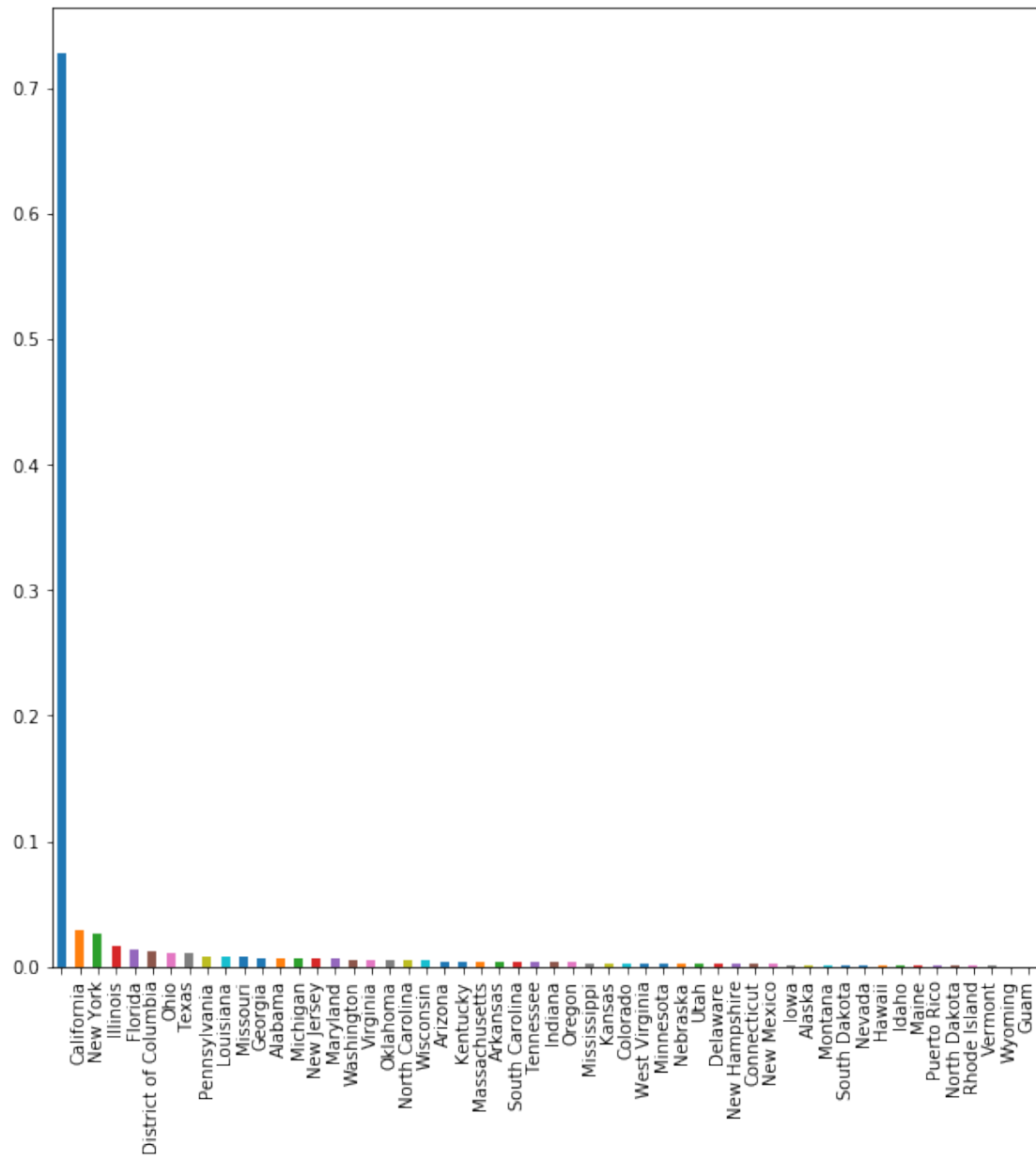


threeJudgeFdc

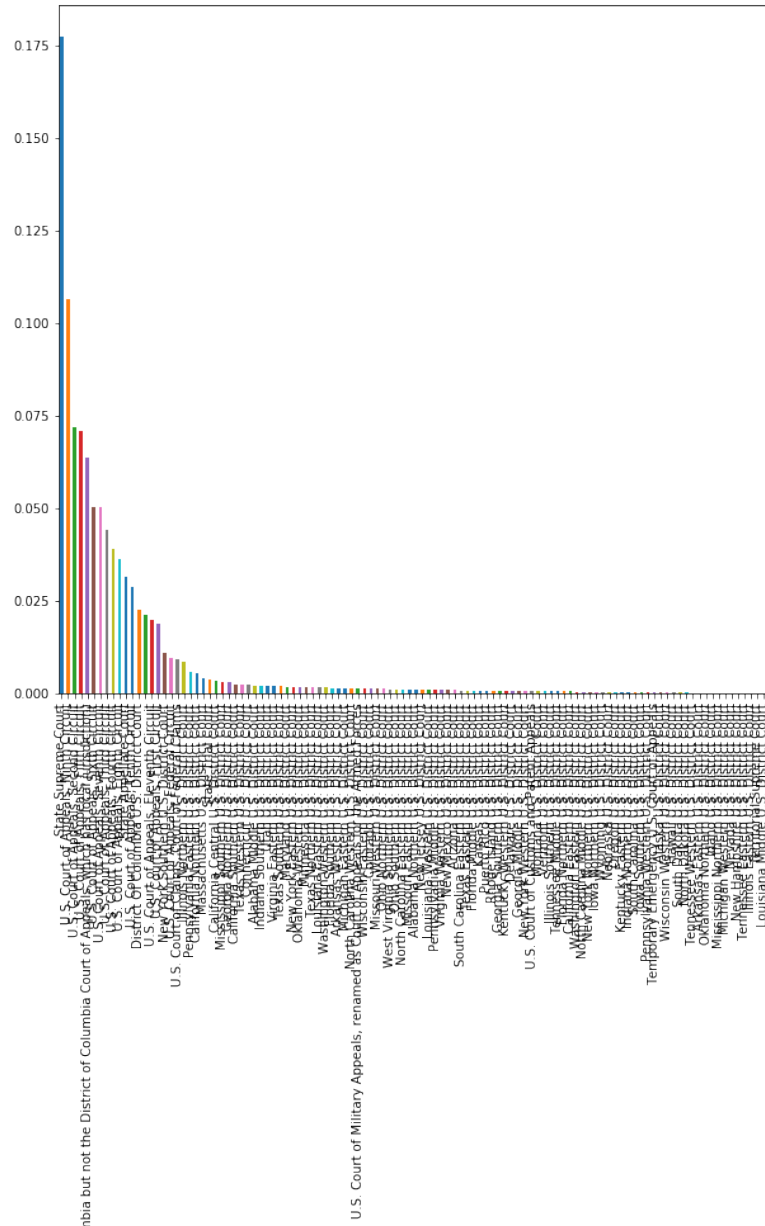


caseOrigin

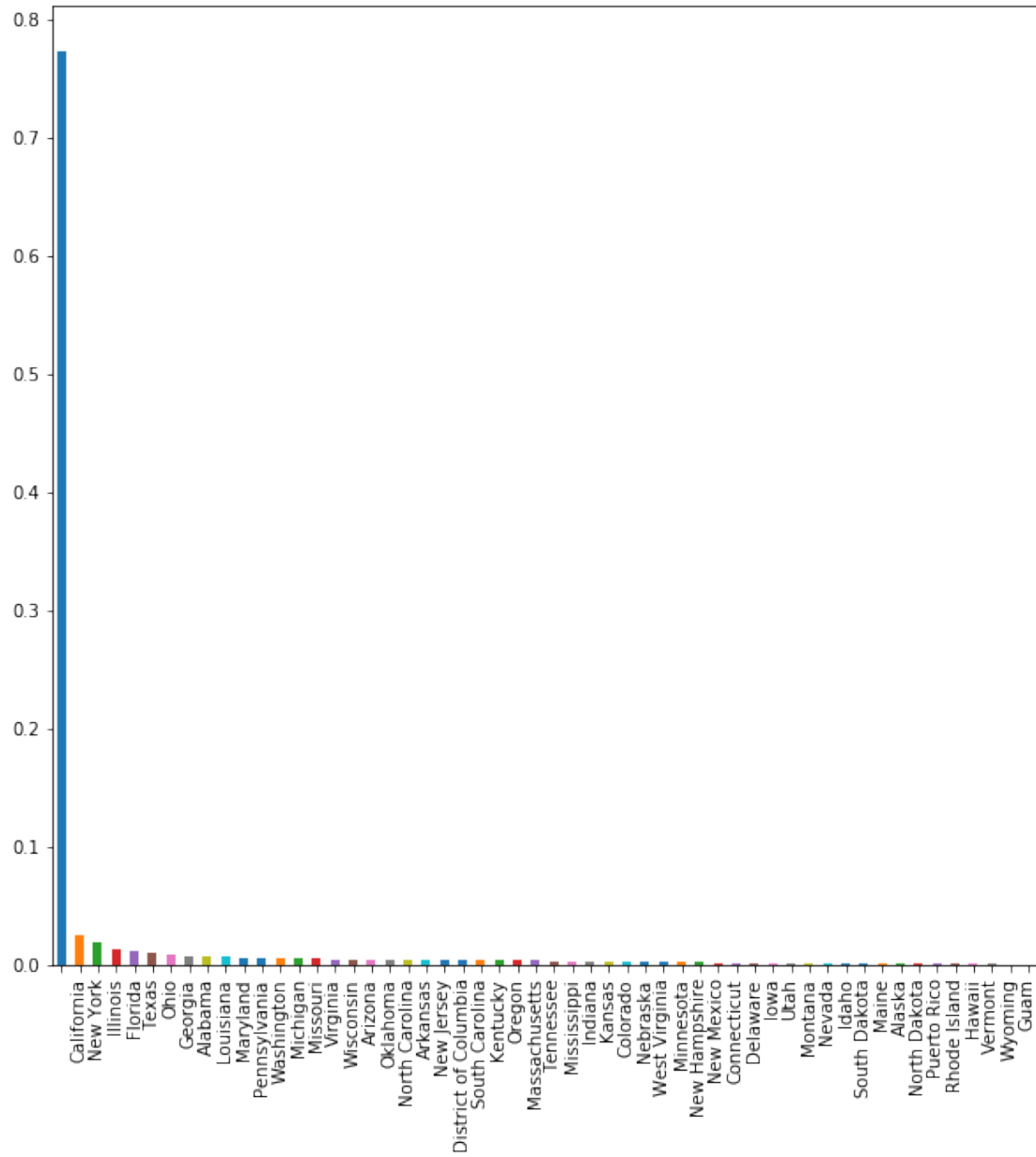
caseOriginState



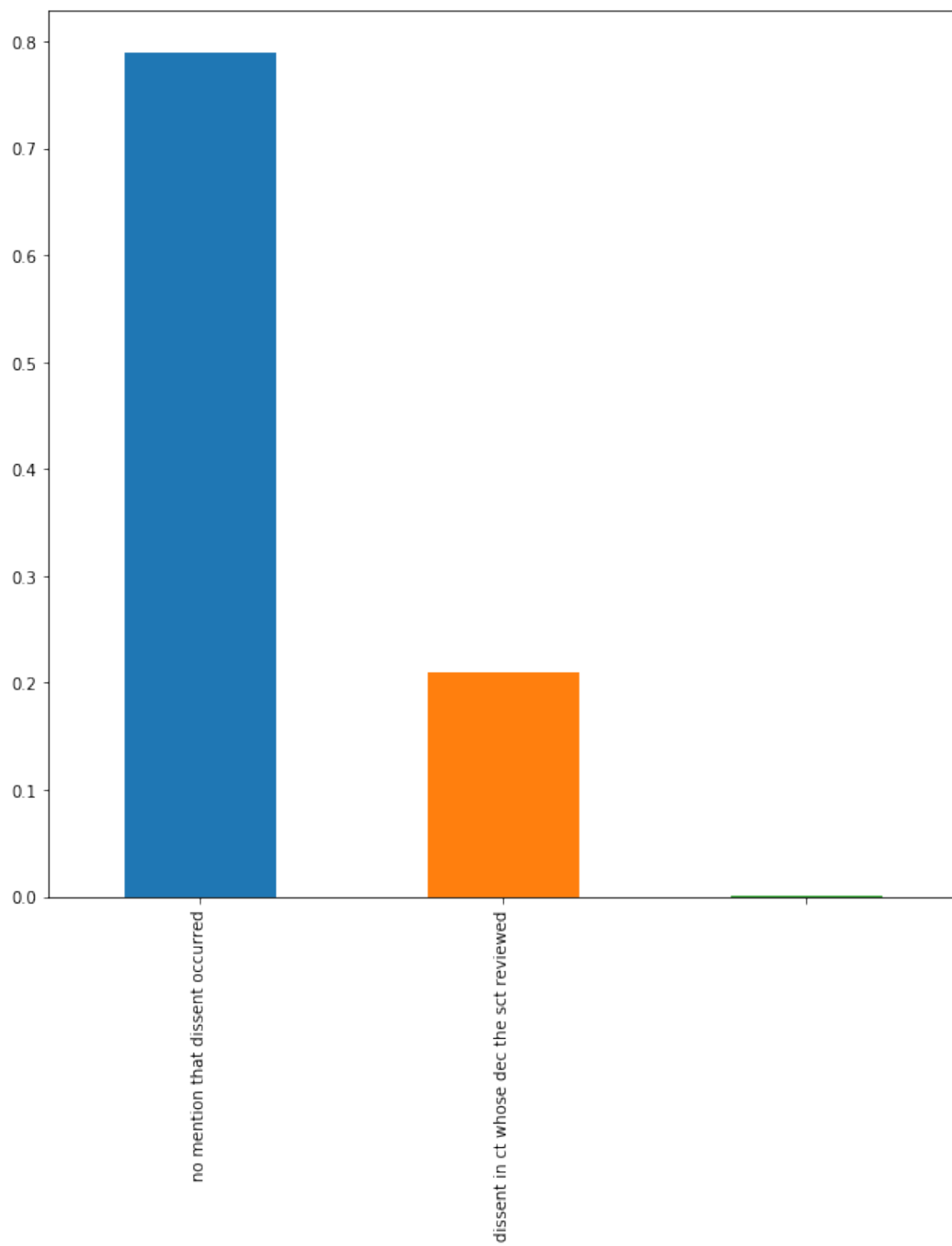
caseSource



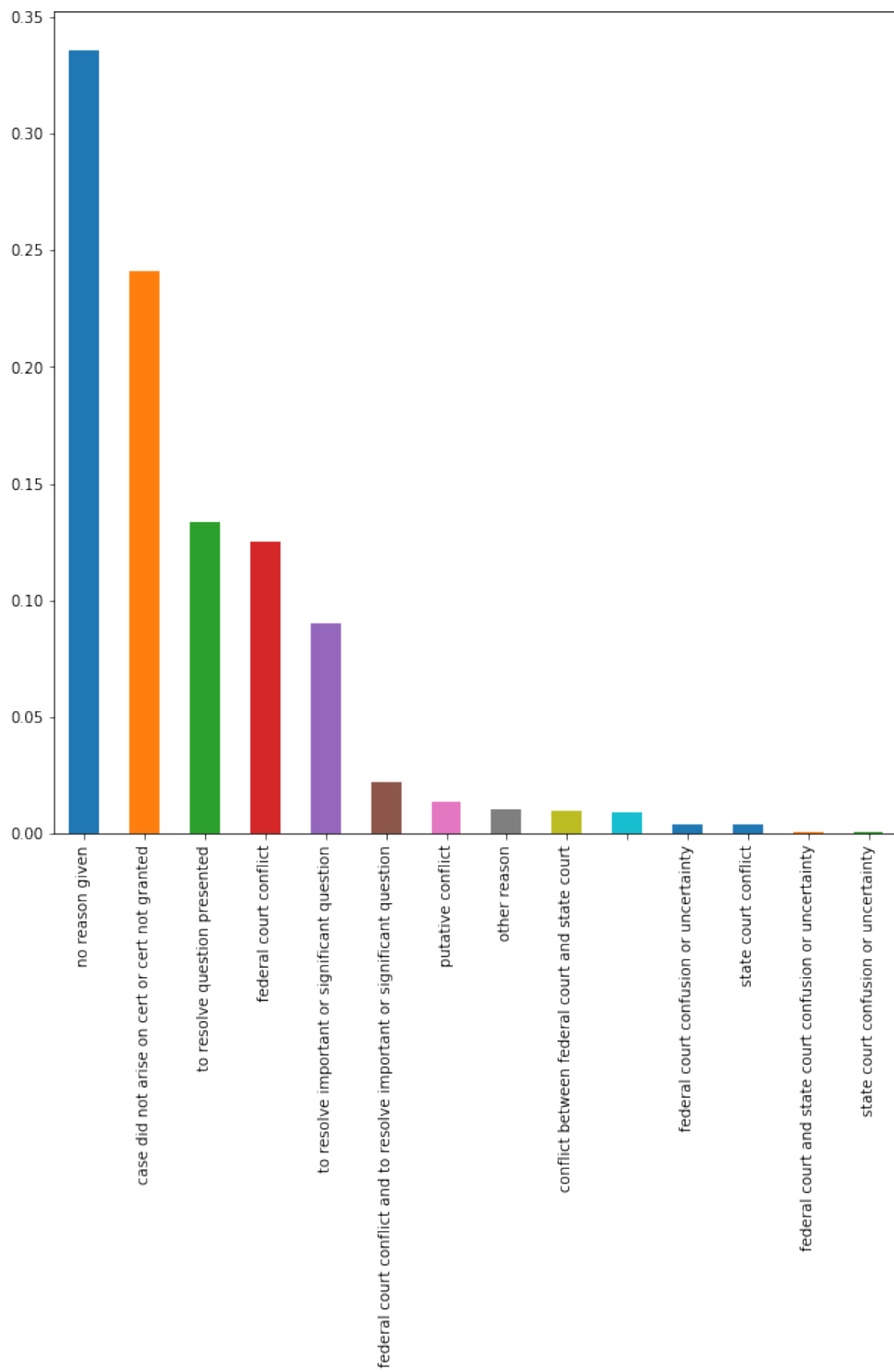
caseSourceState



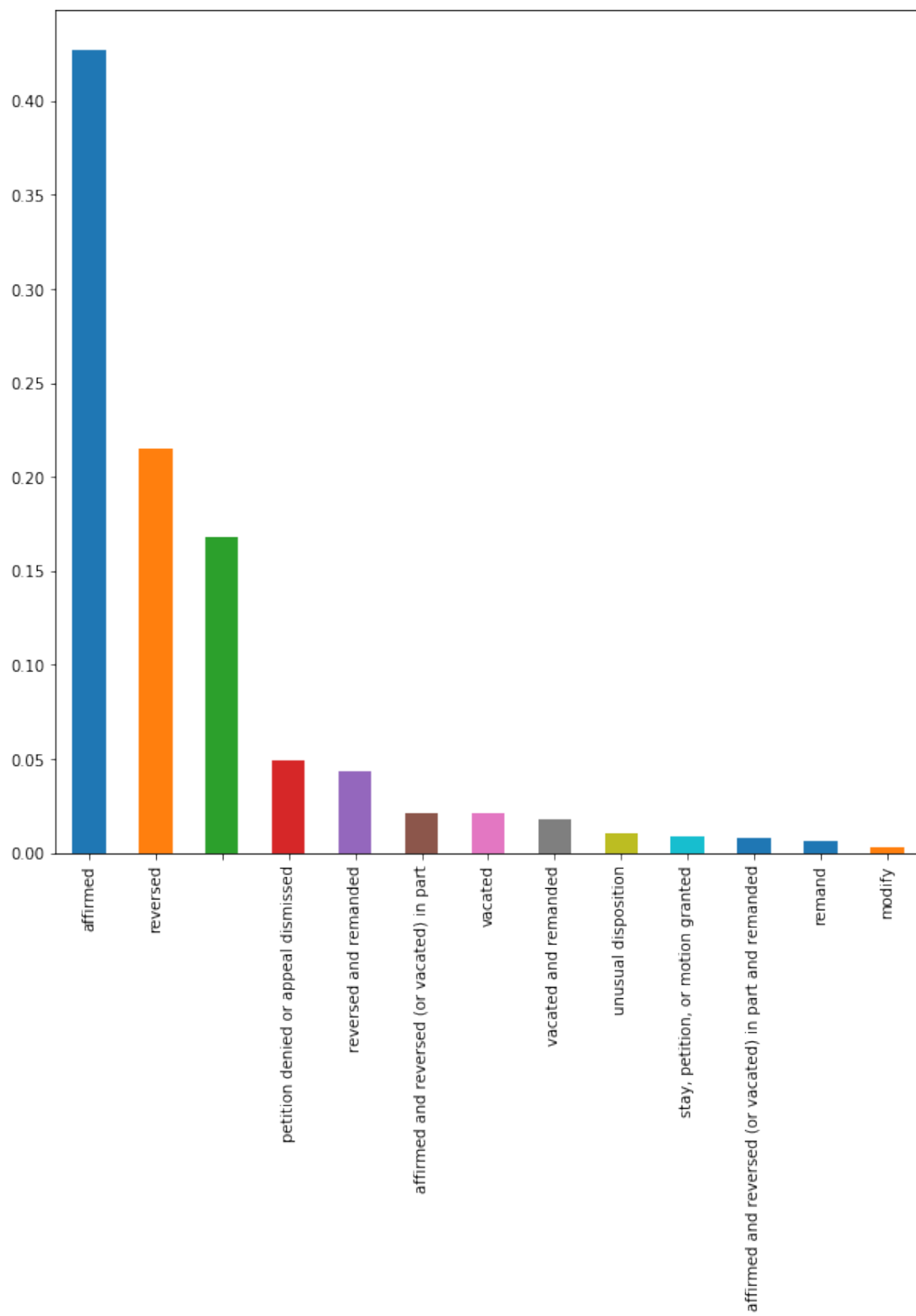
lcDisagreement



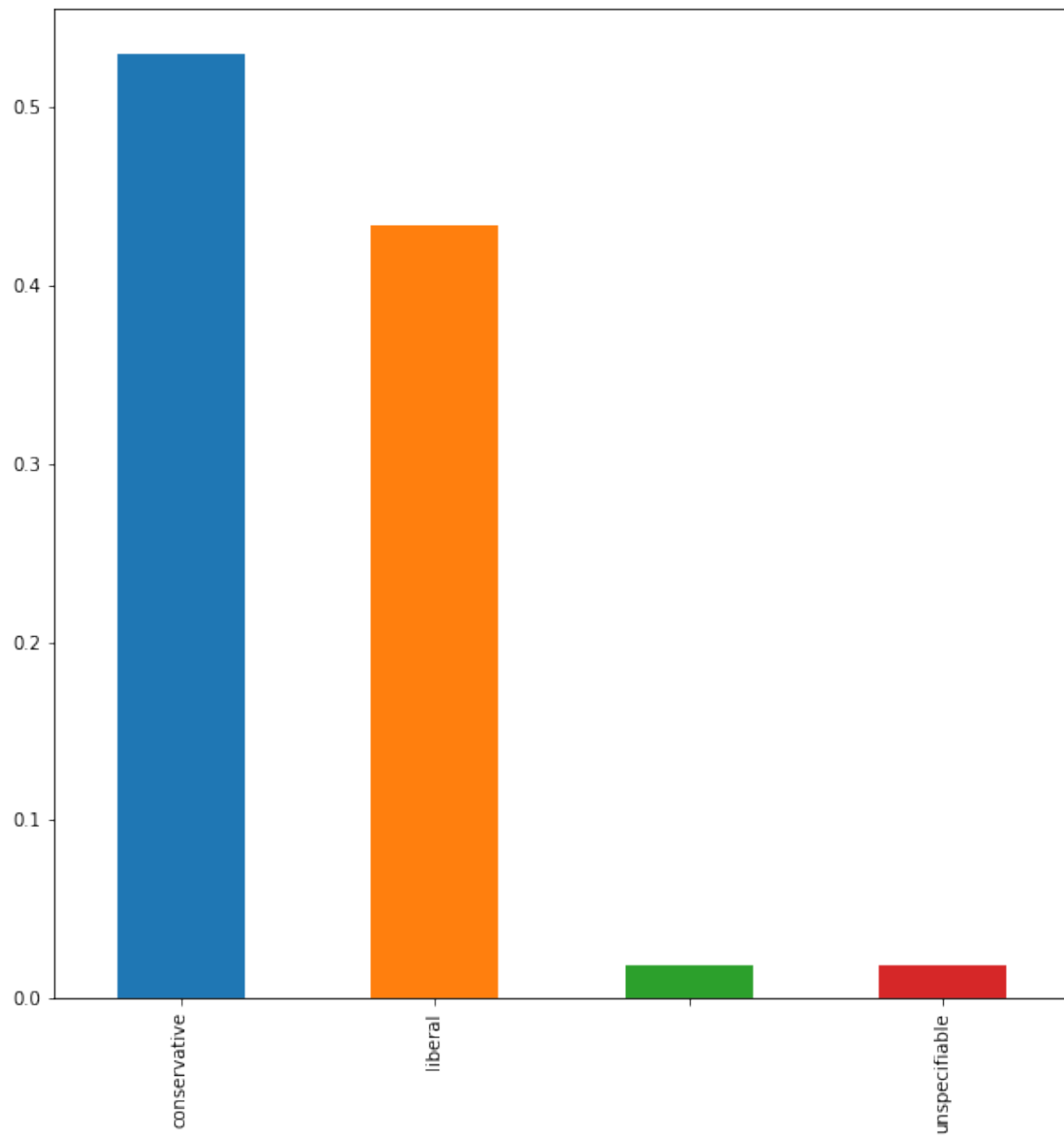
certReason



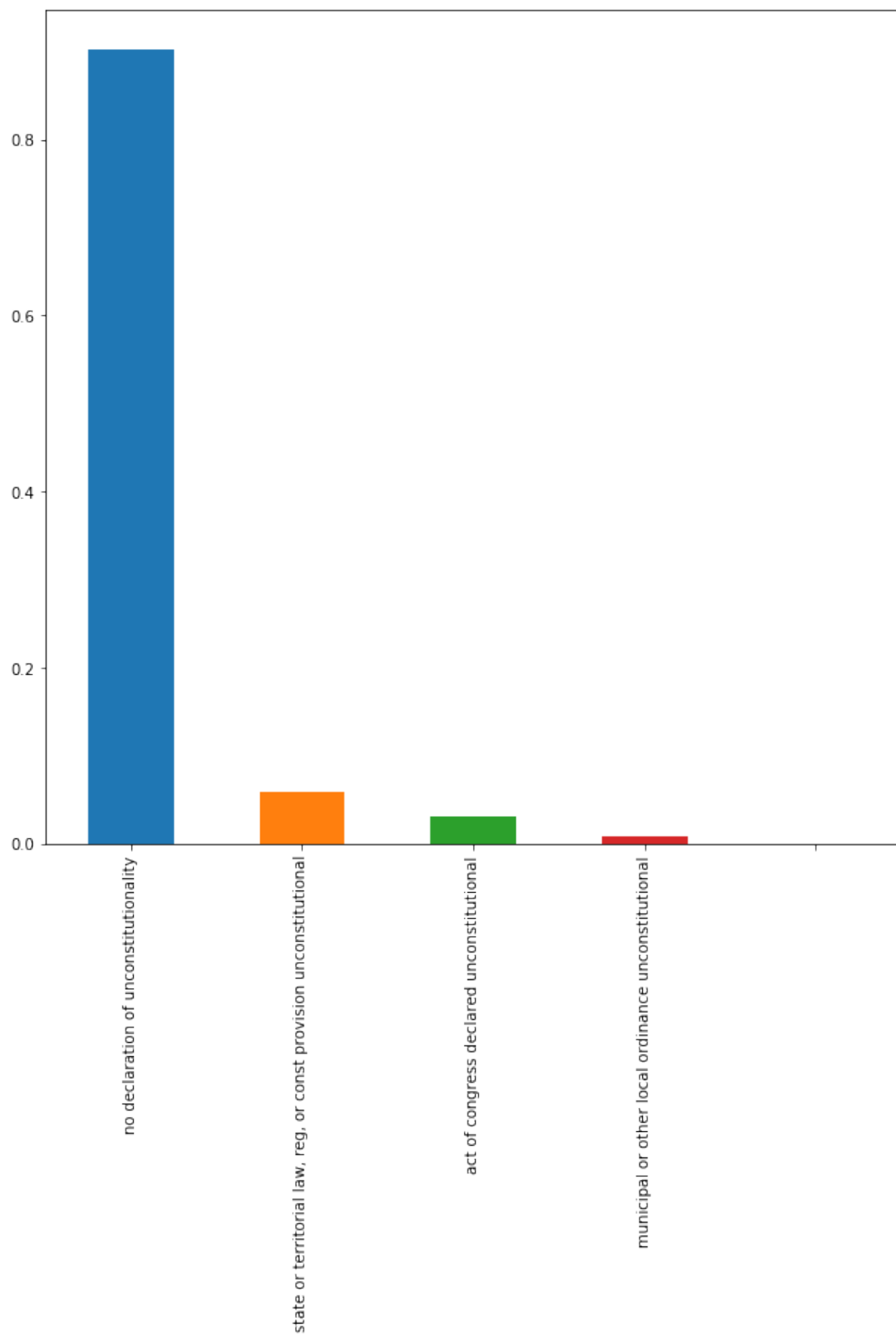
lcDisposition



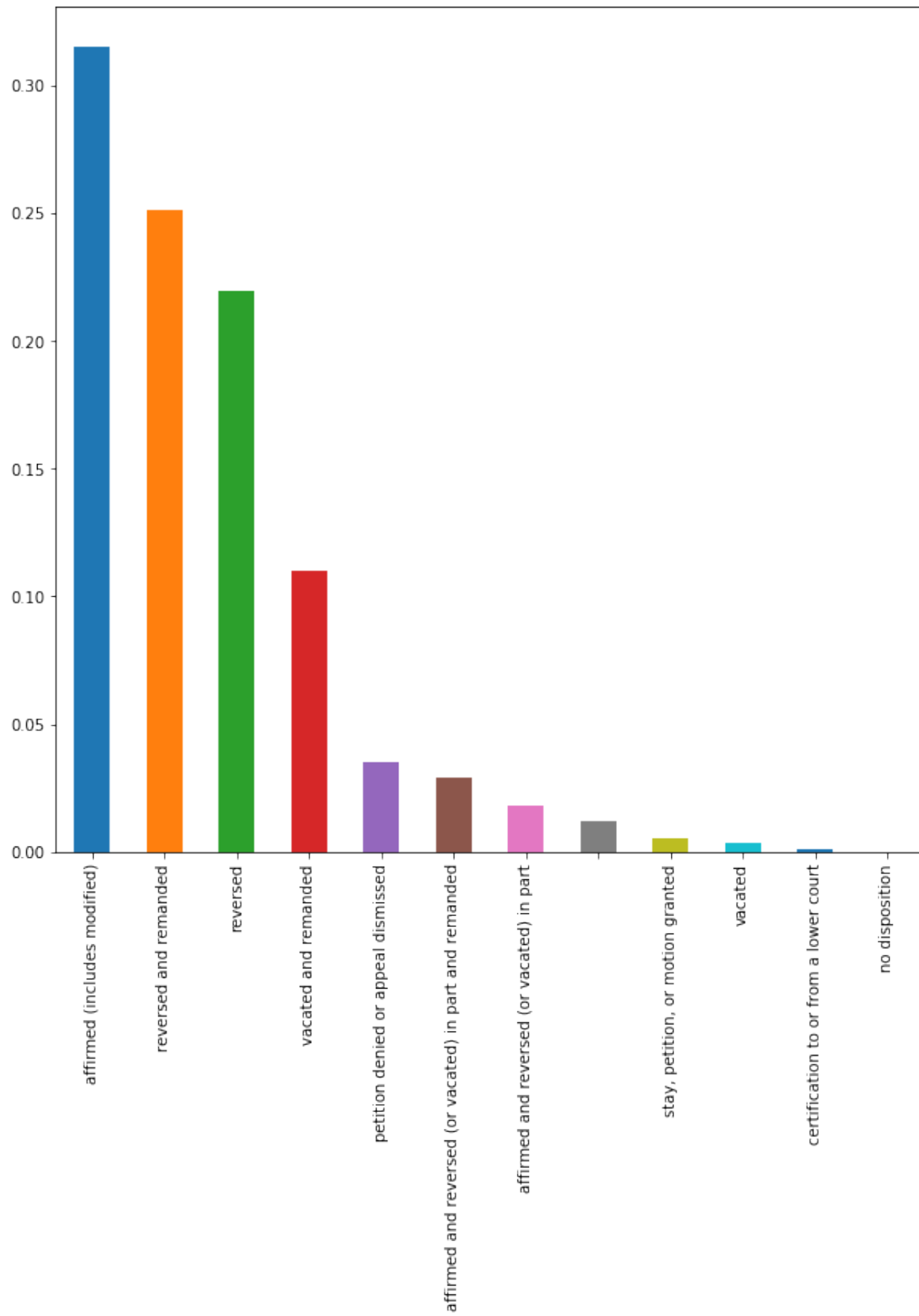
lcDispositionDirection



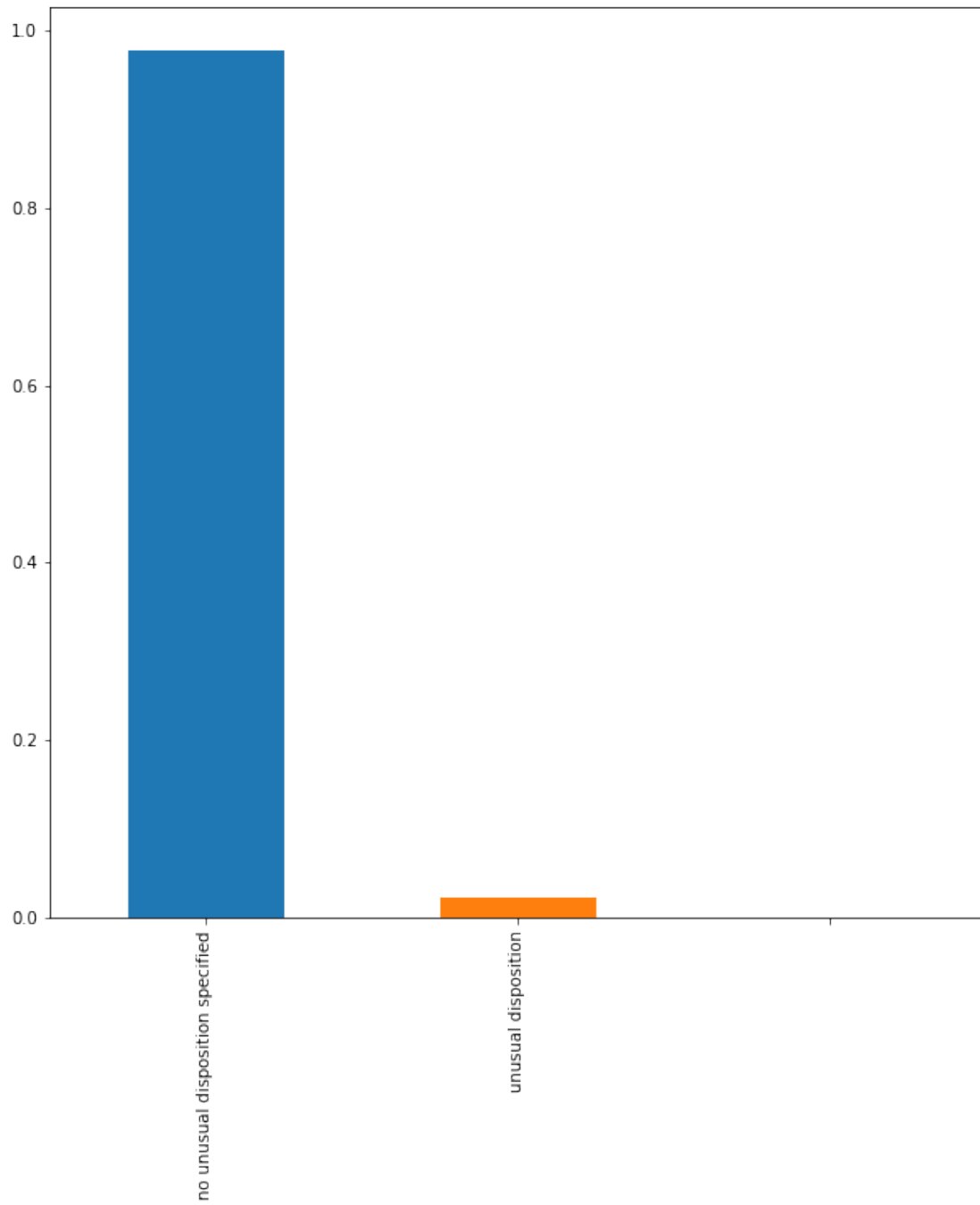
declarationUncon



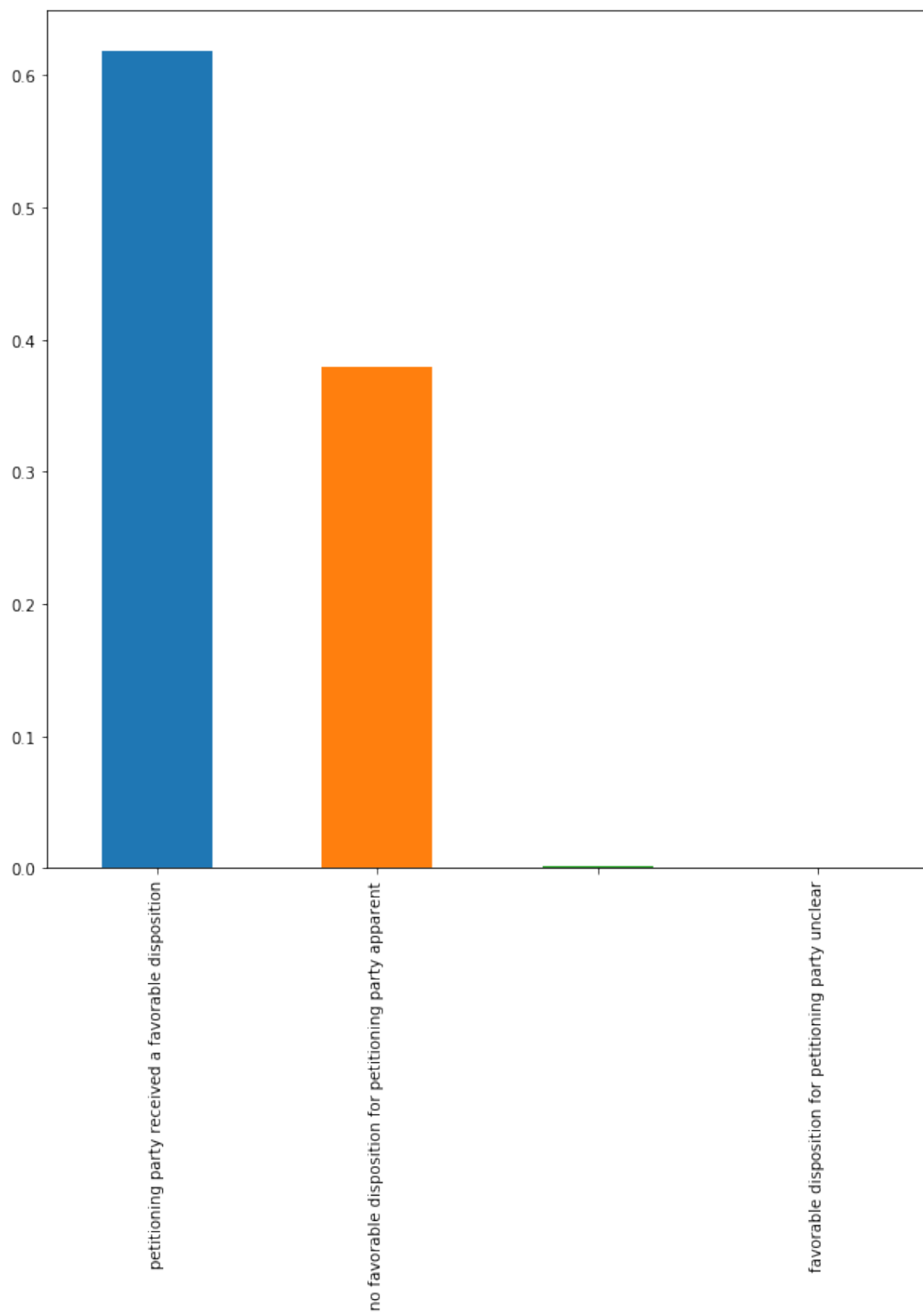
caseDisposition



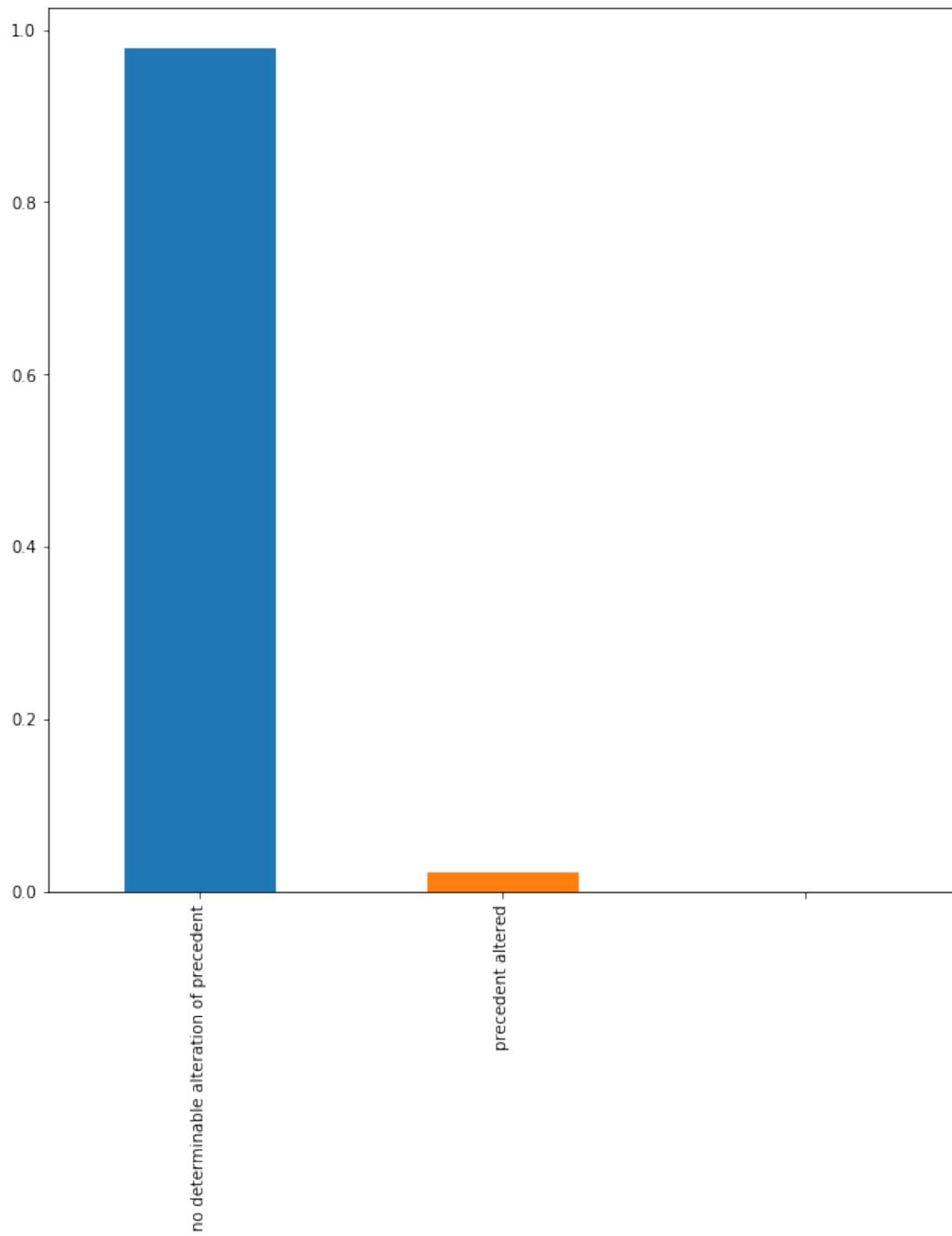
caseDispositionUnusual



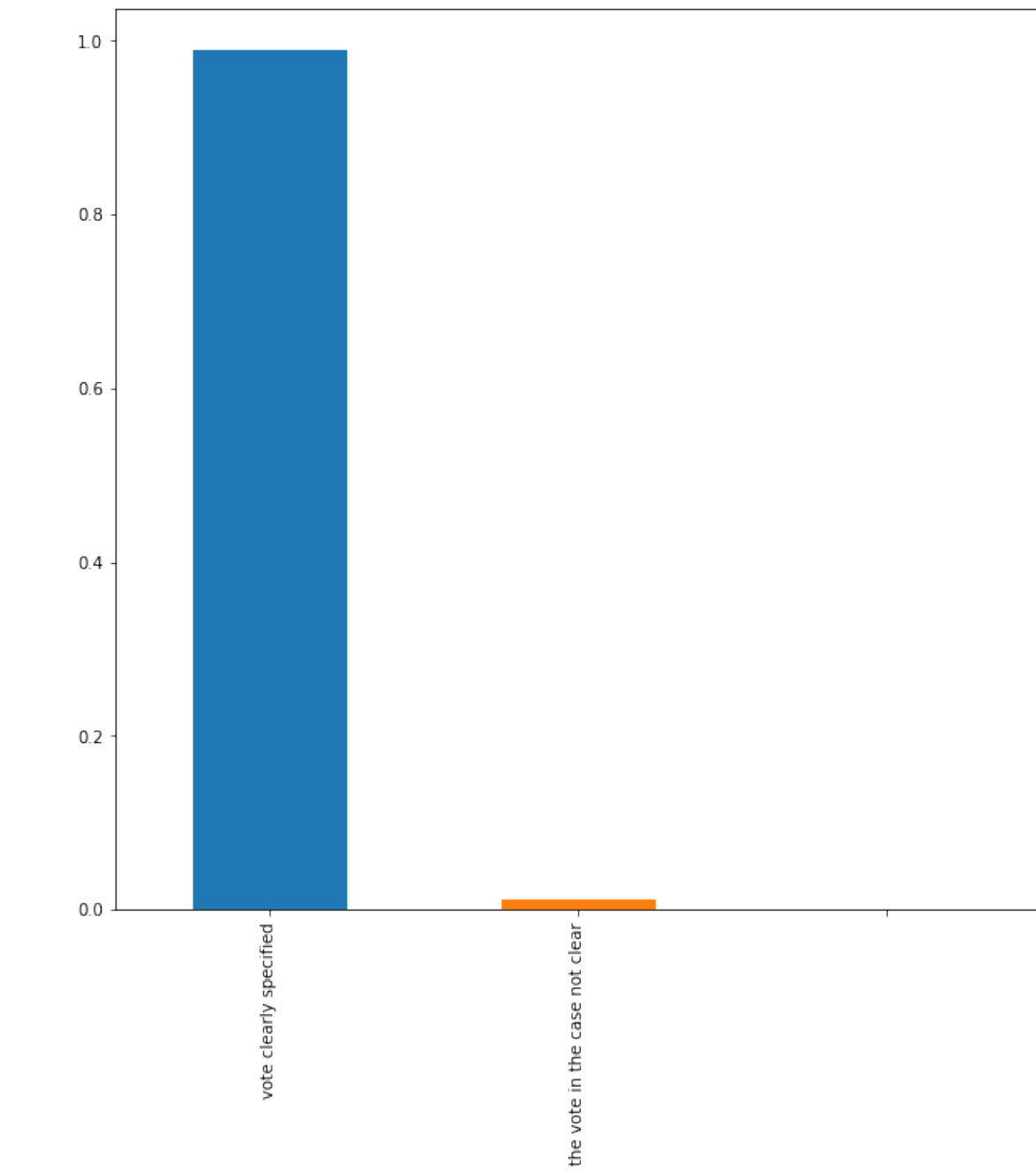
partyWinning



precedentAlteration

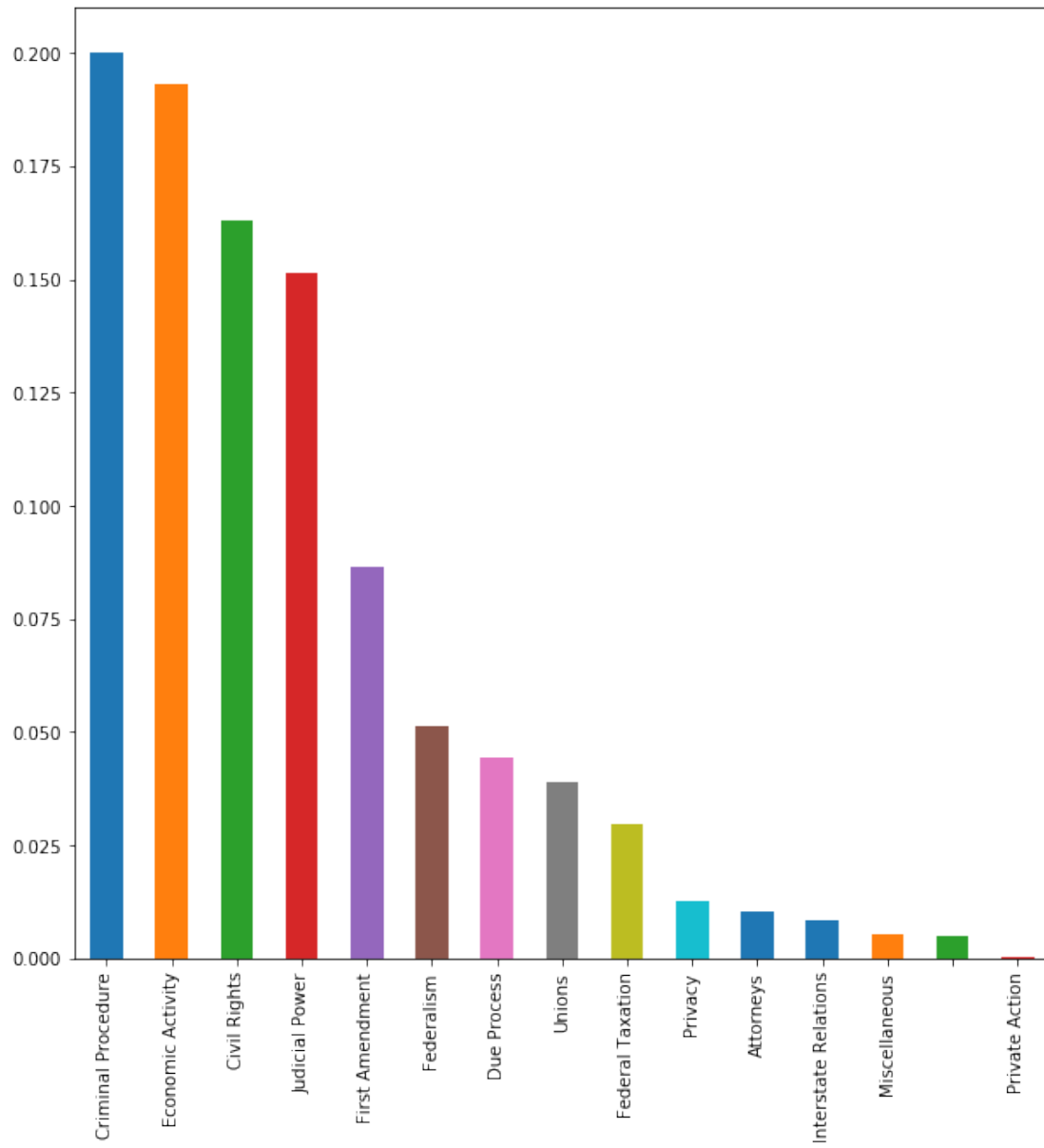


voteUnclear

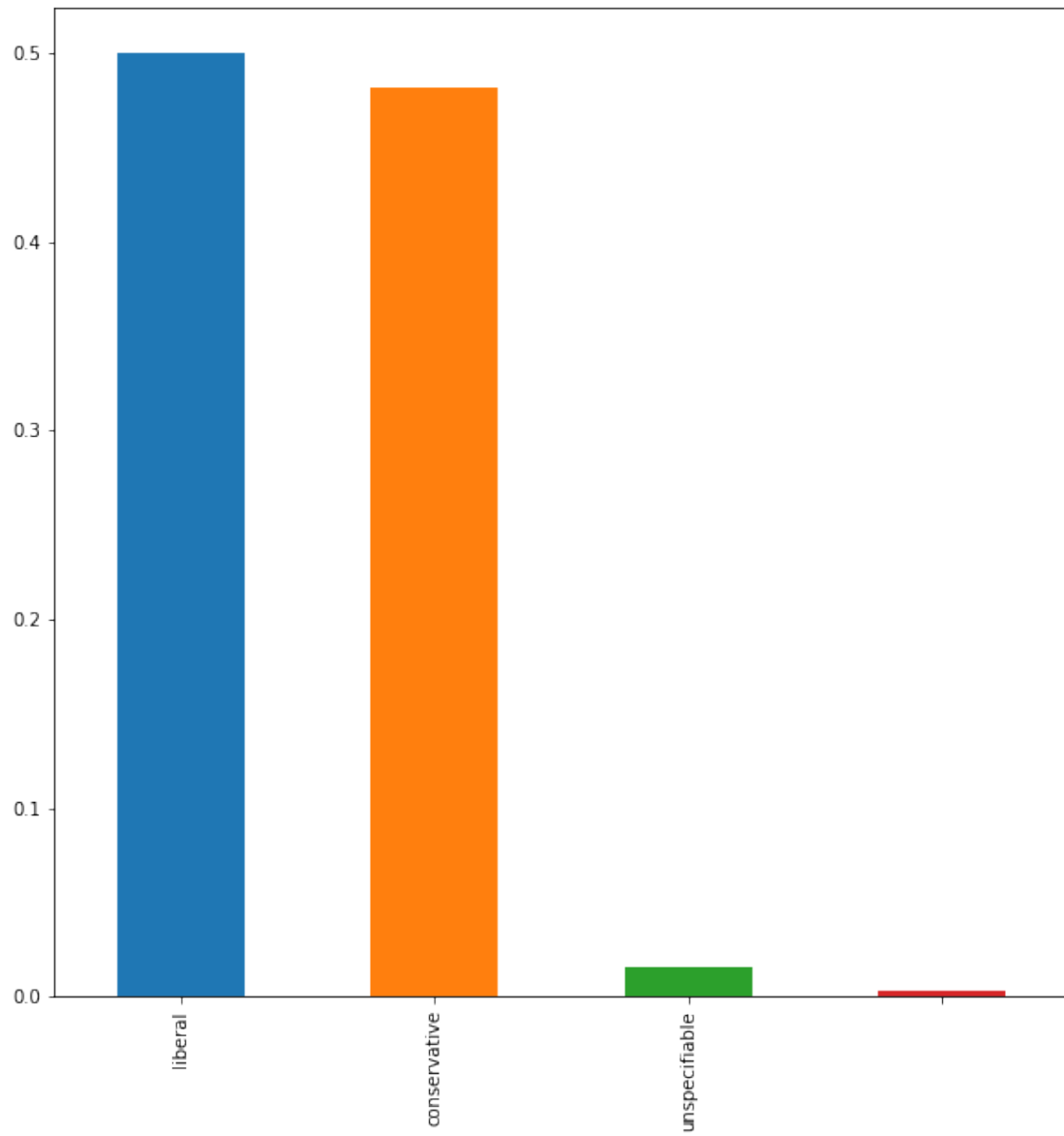


issue

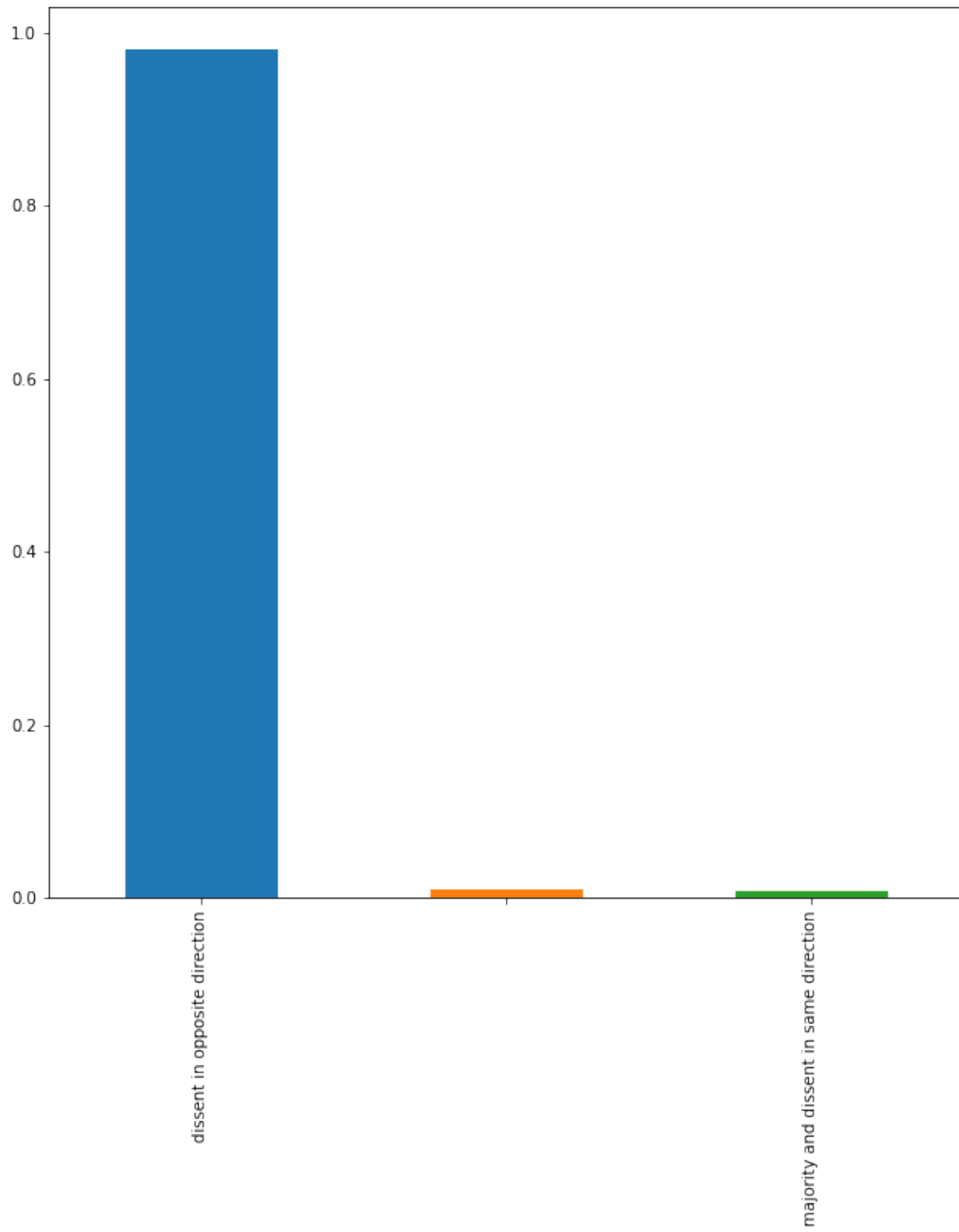
issueArea



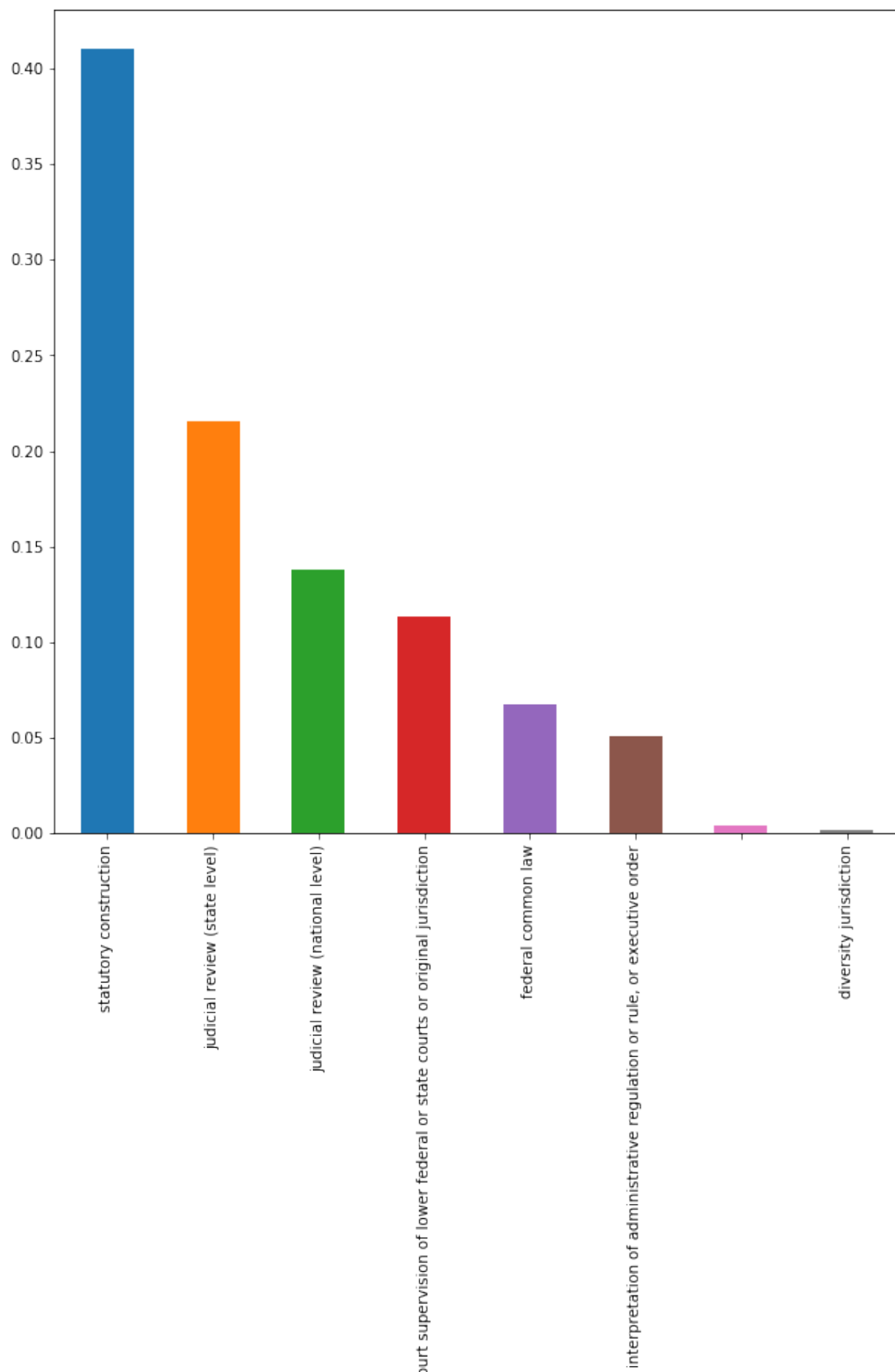
decisionDirection



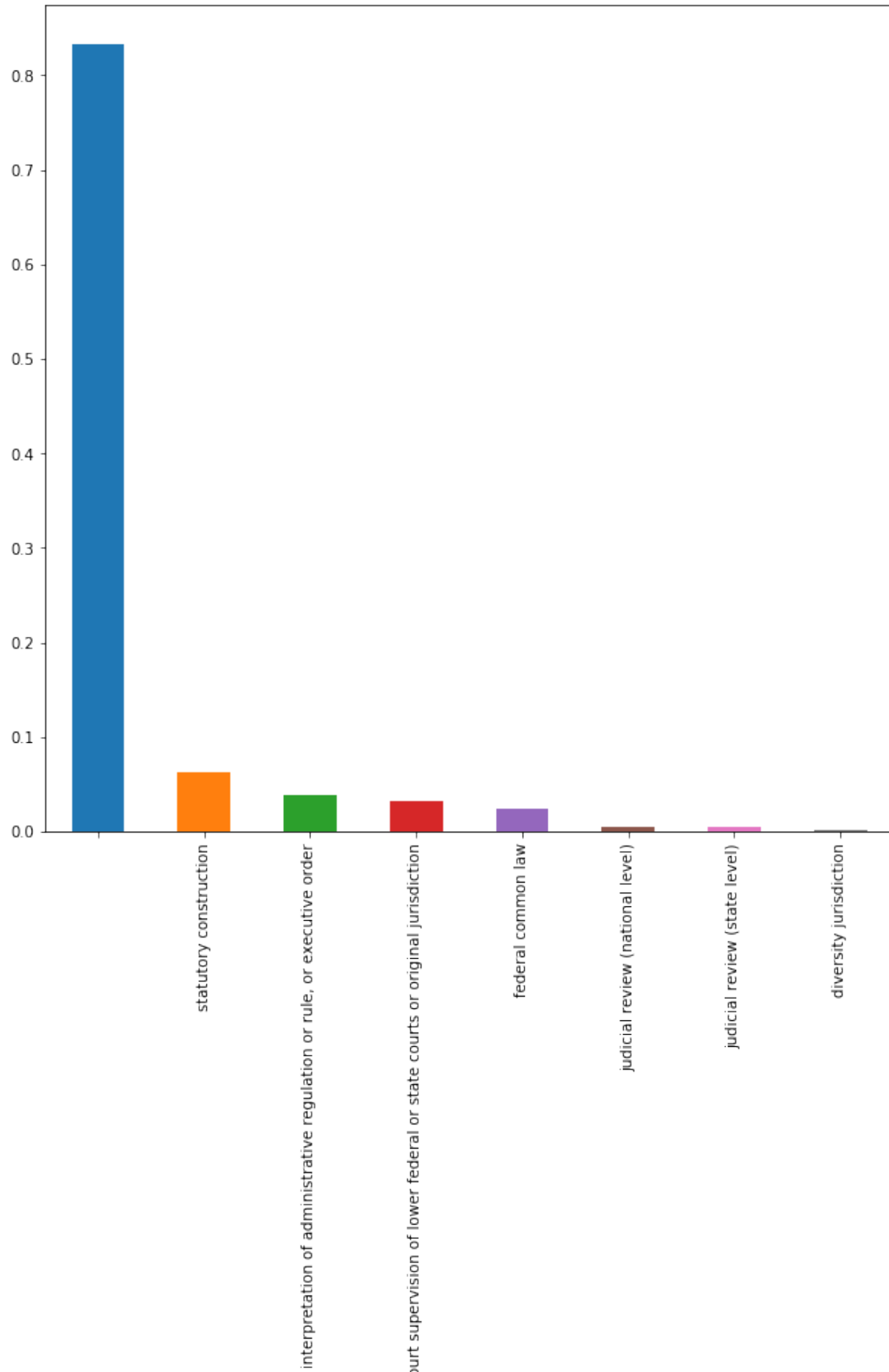
decisionDirectionDissent



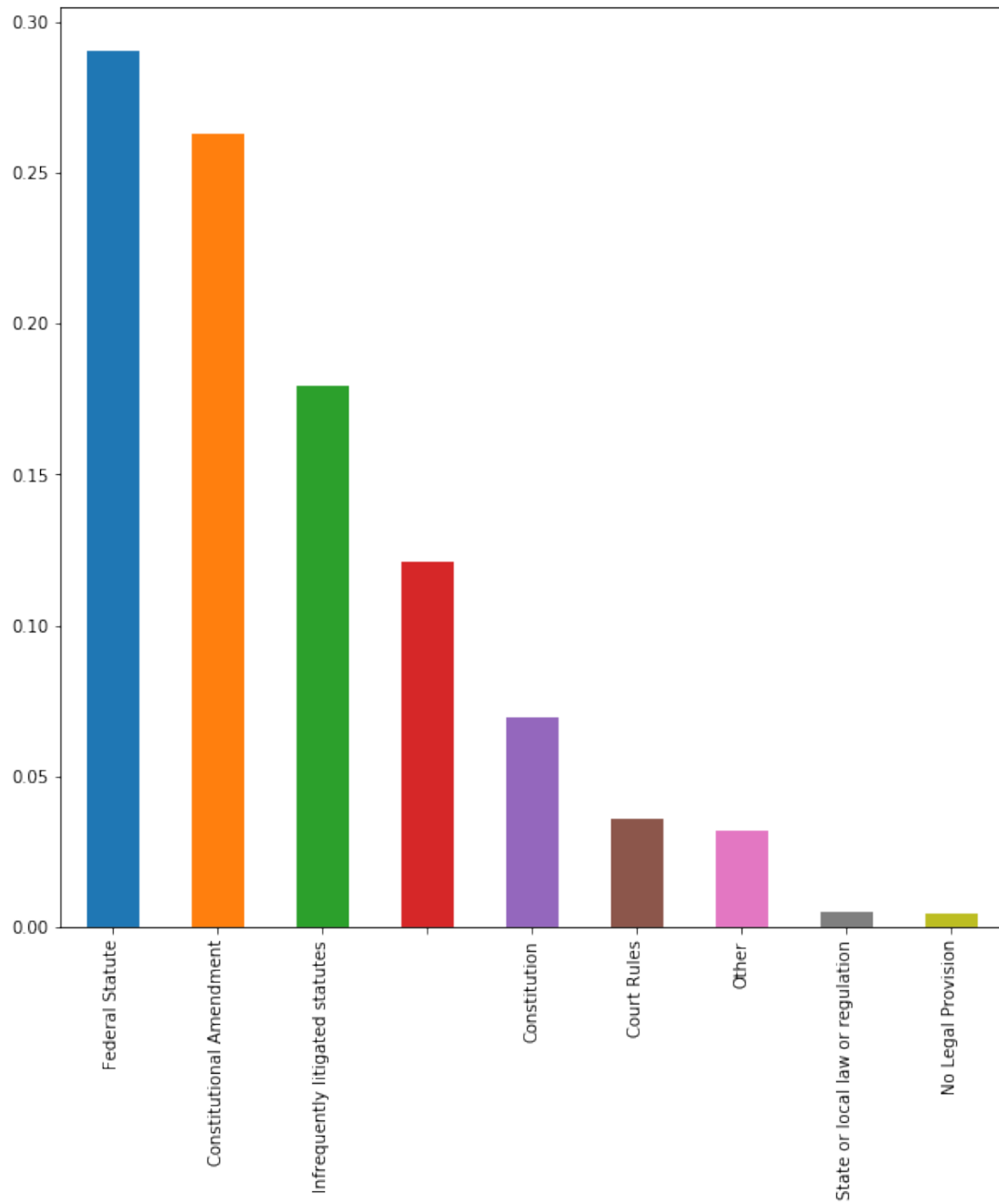
authorityDecision1



authorityDecision2

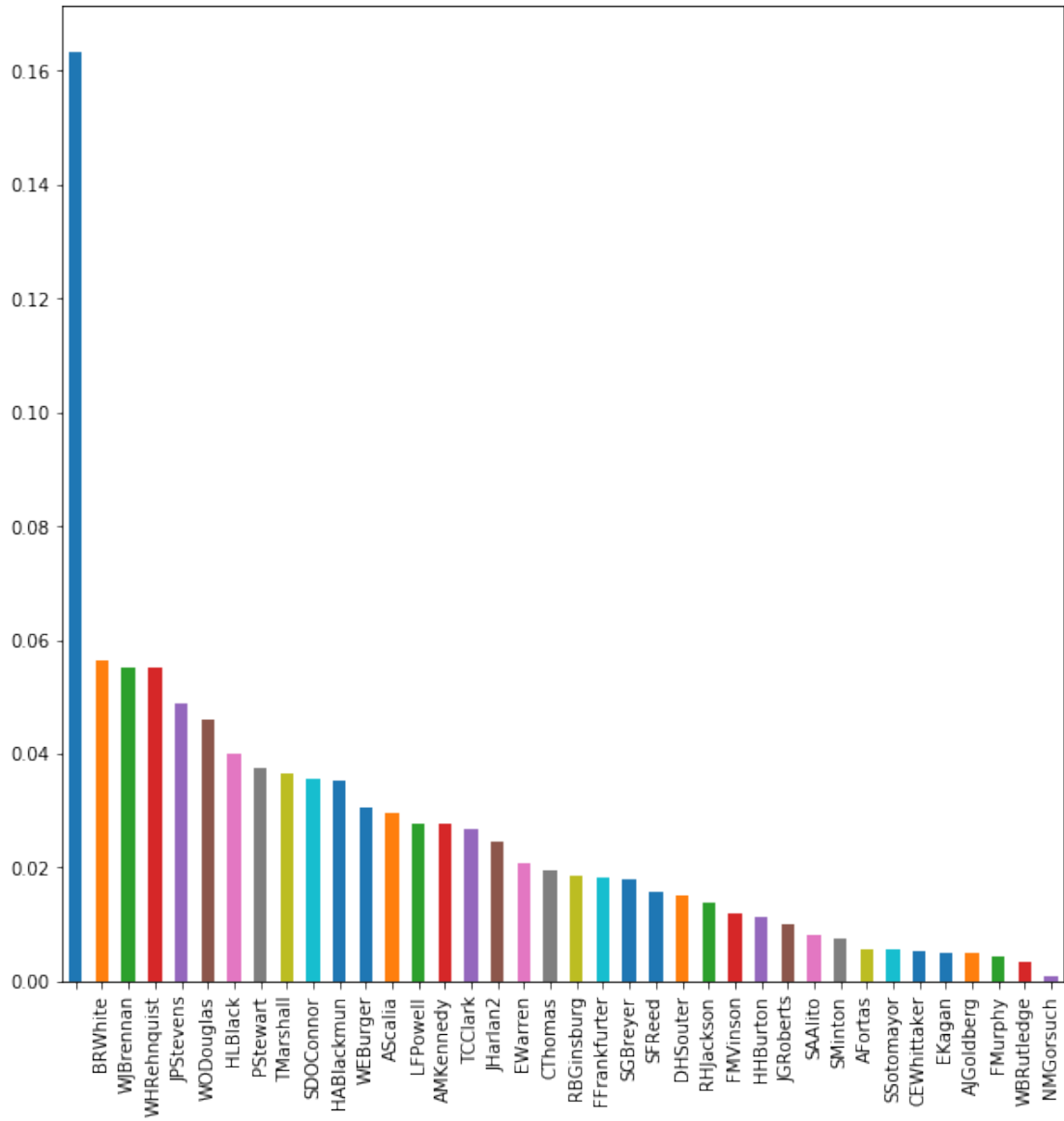


lawType

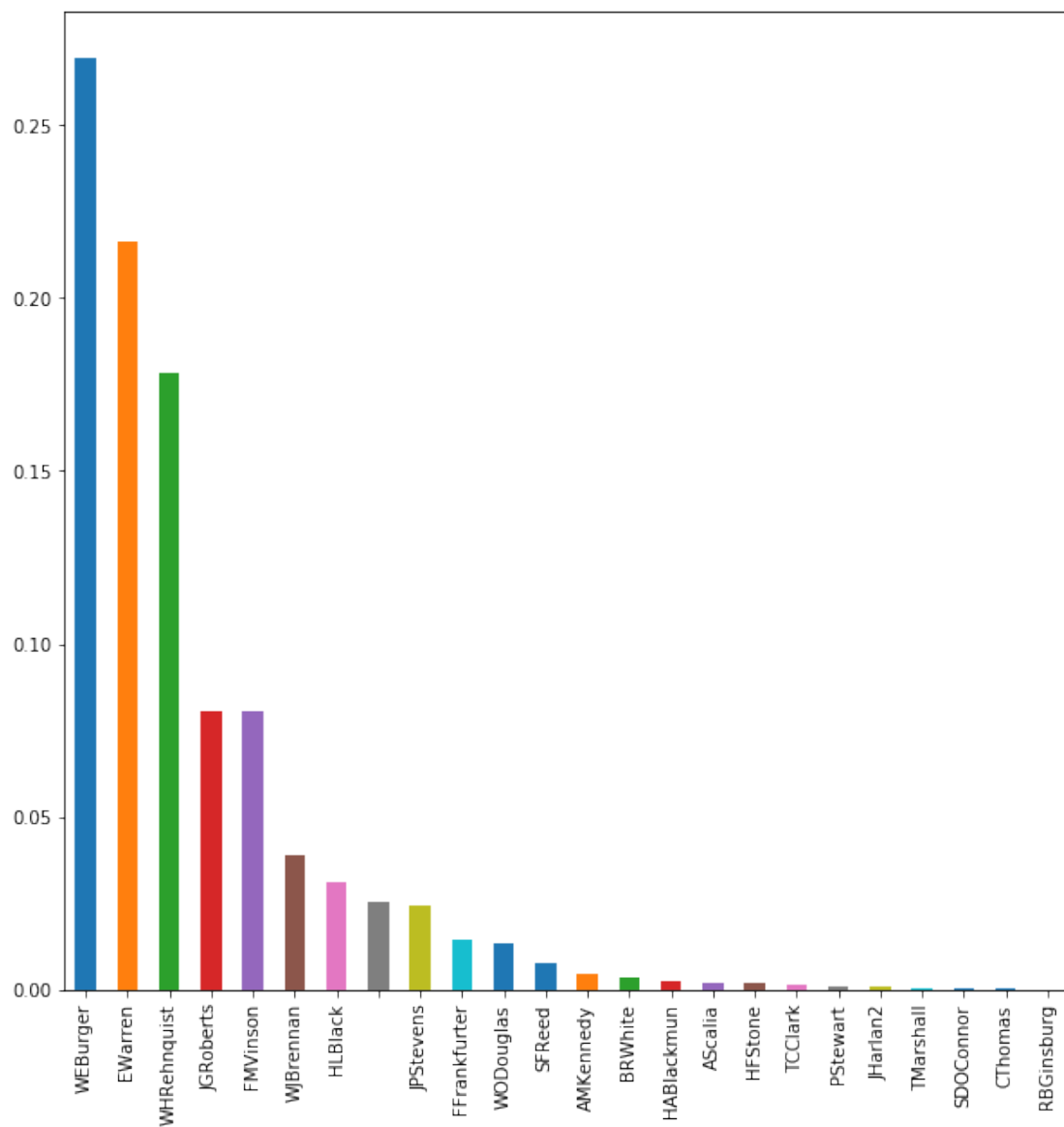


lawSupp

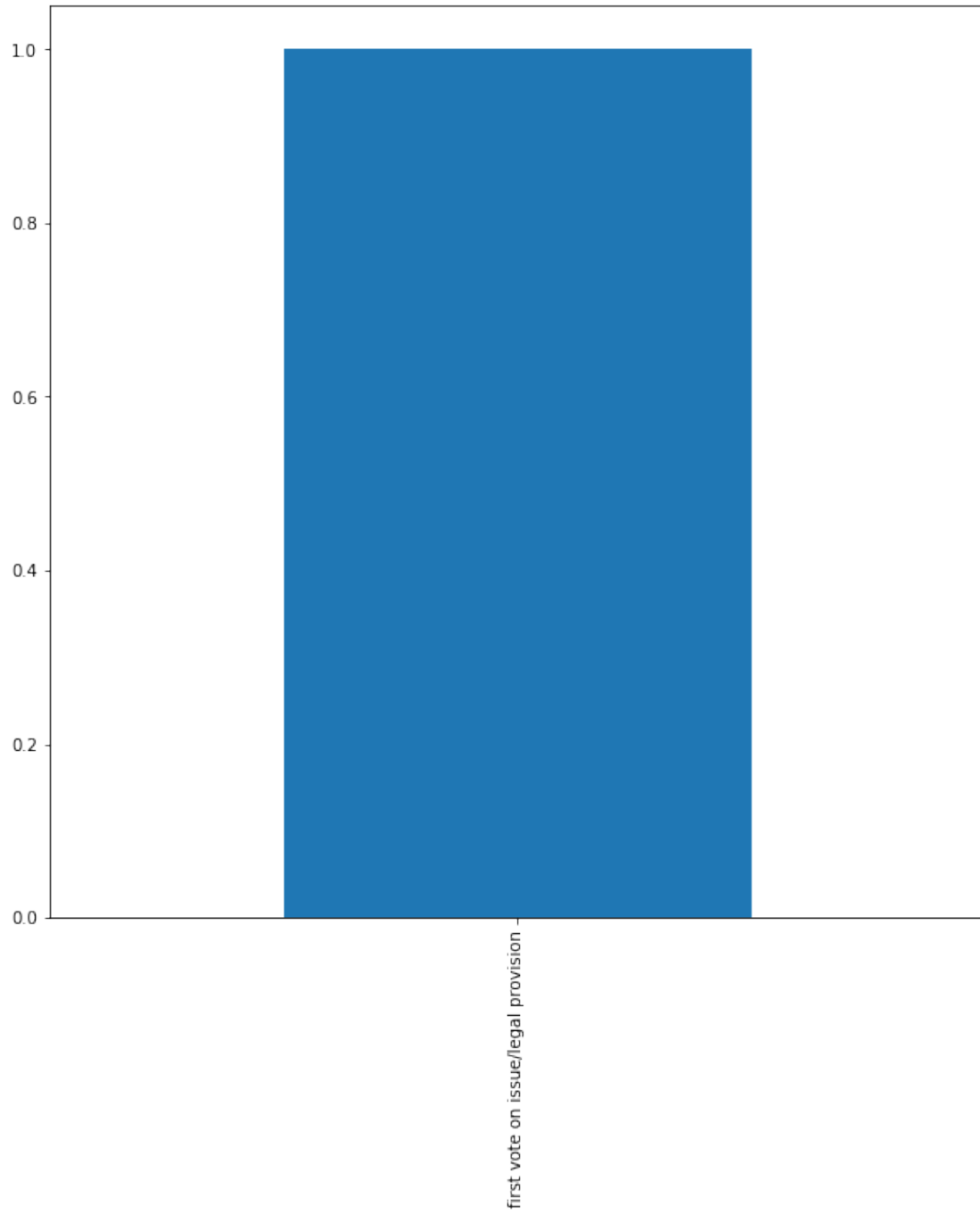
majOpinWriter



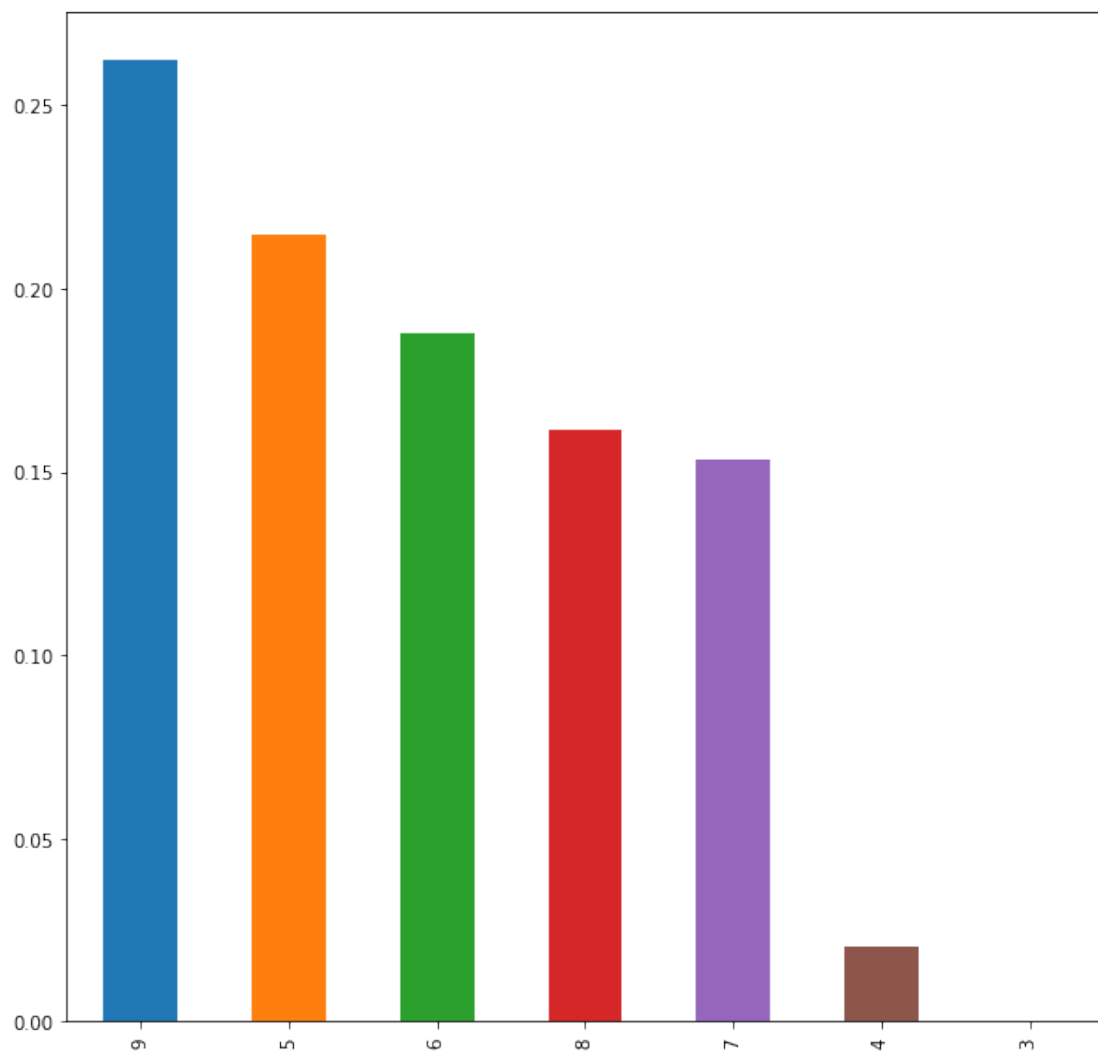
majOpinAssigner



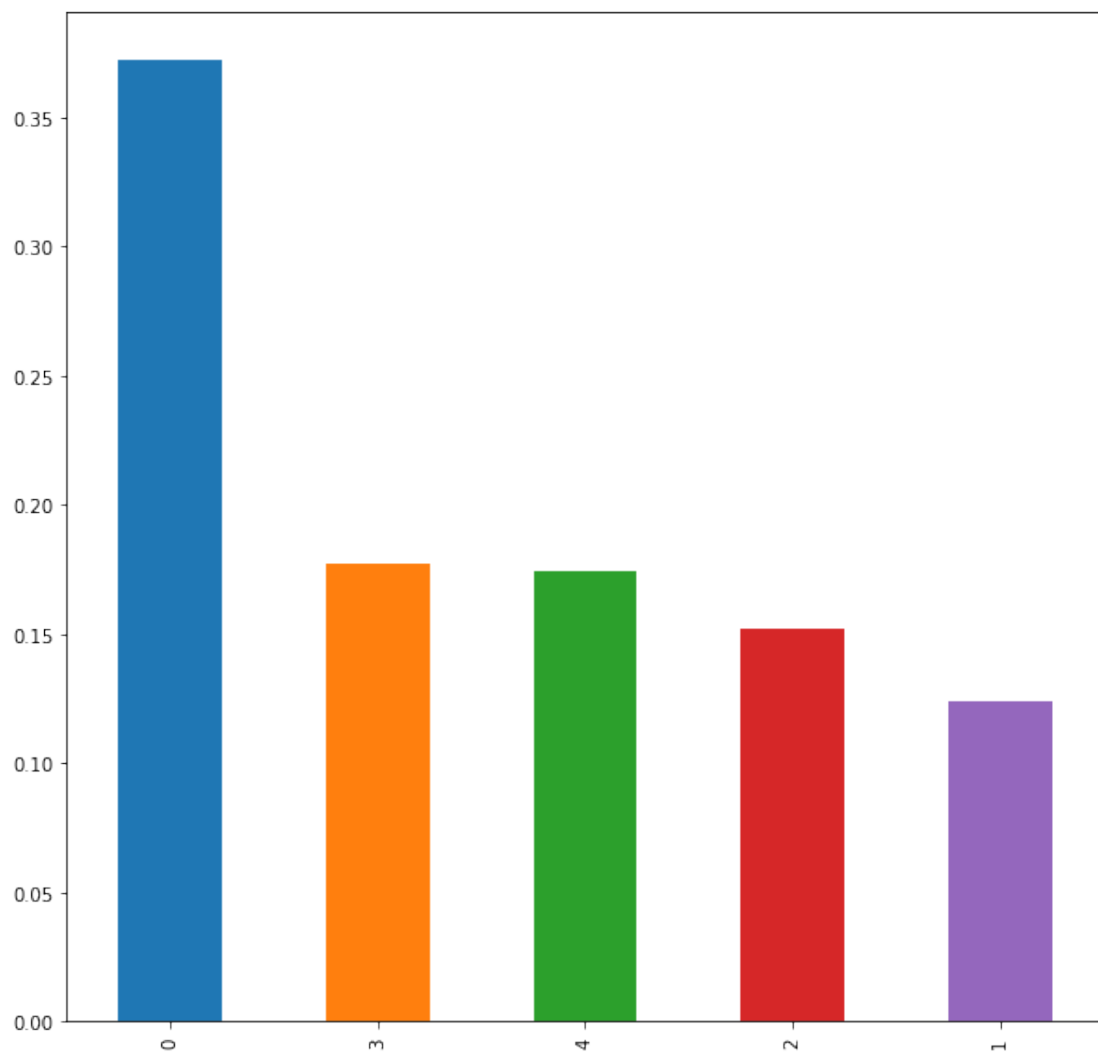
splitVote



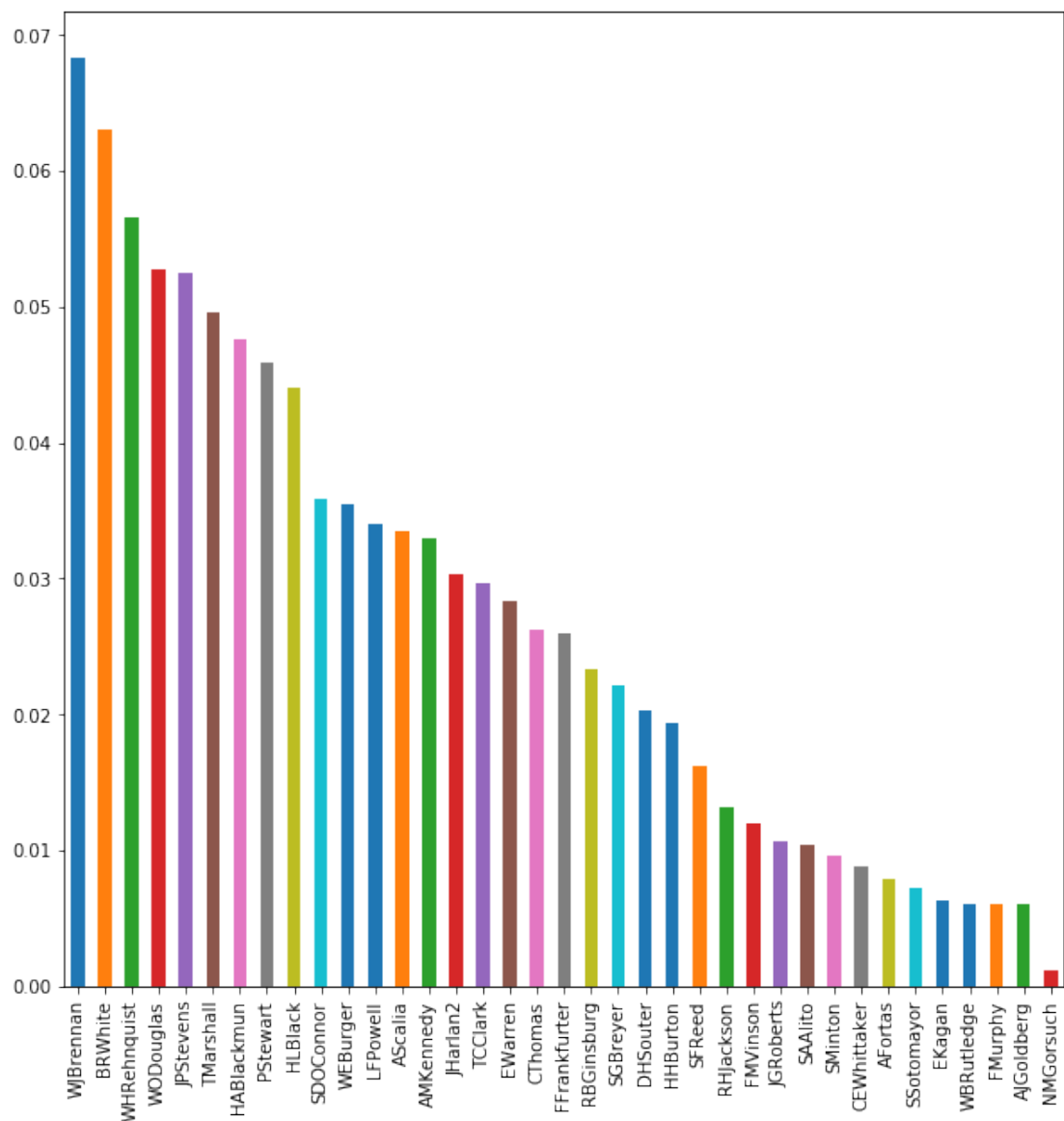
majVotes



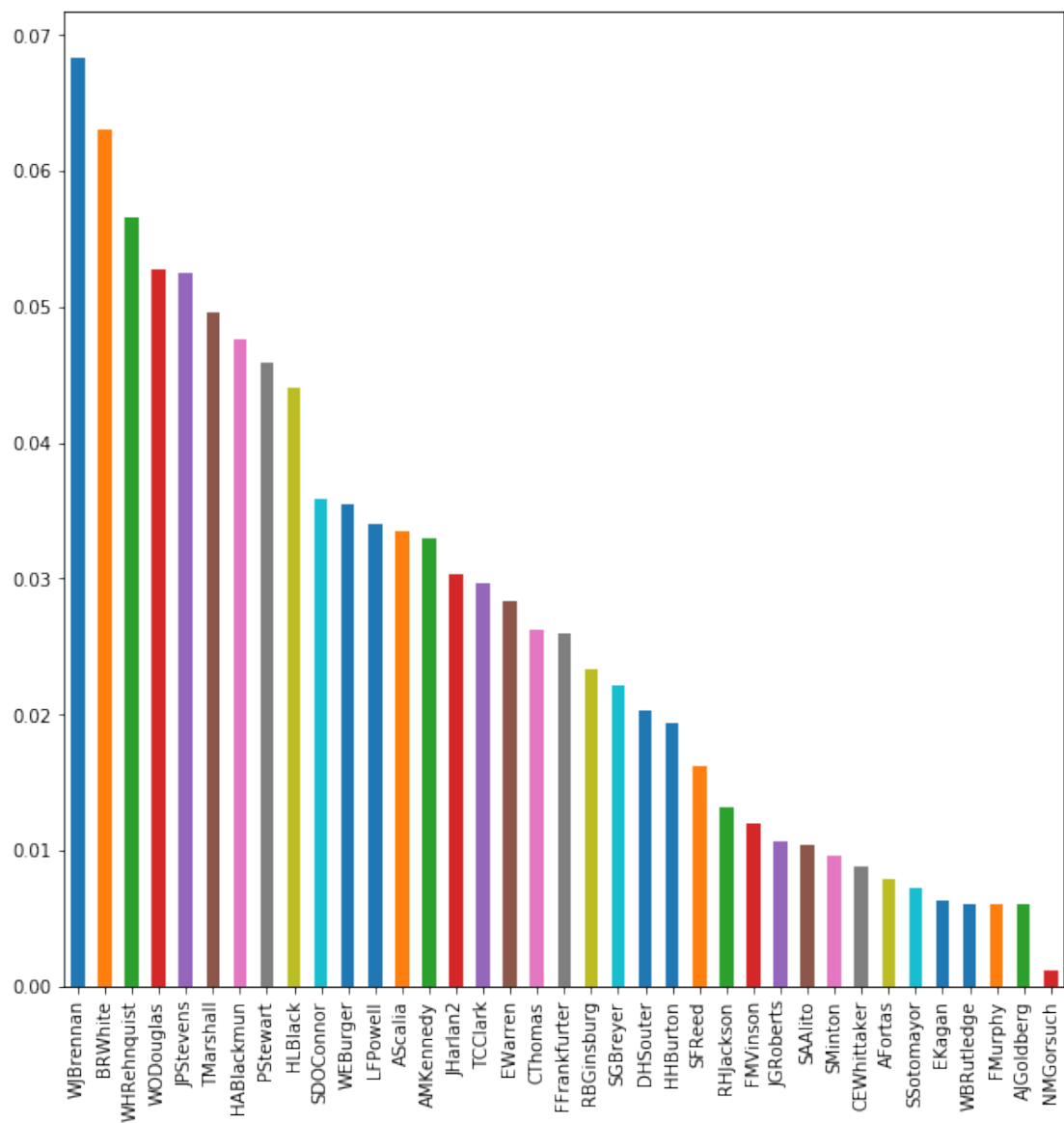
minVotes



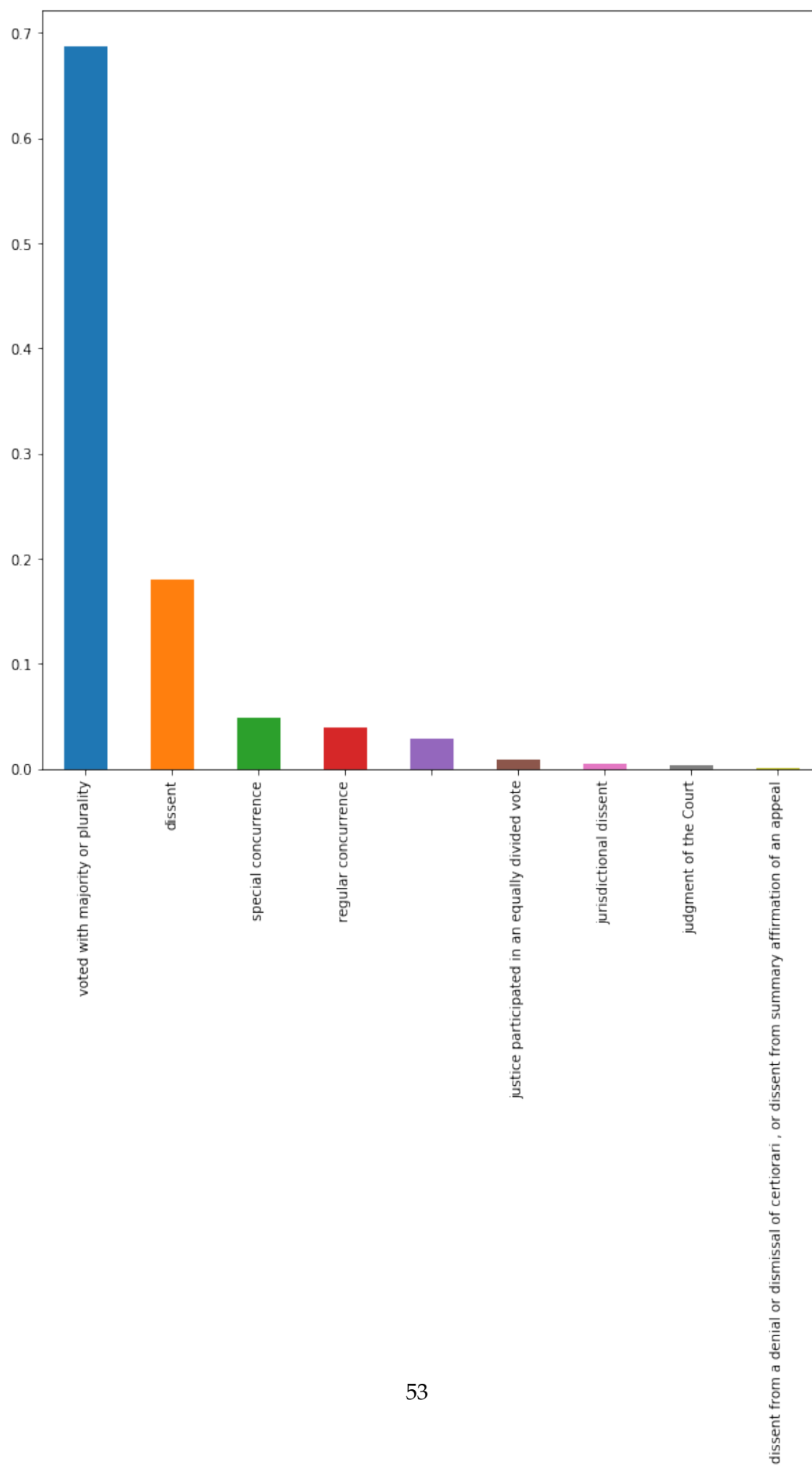
justice



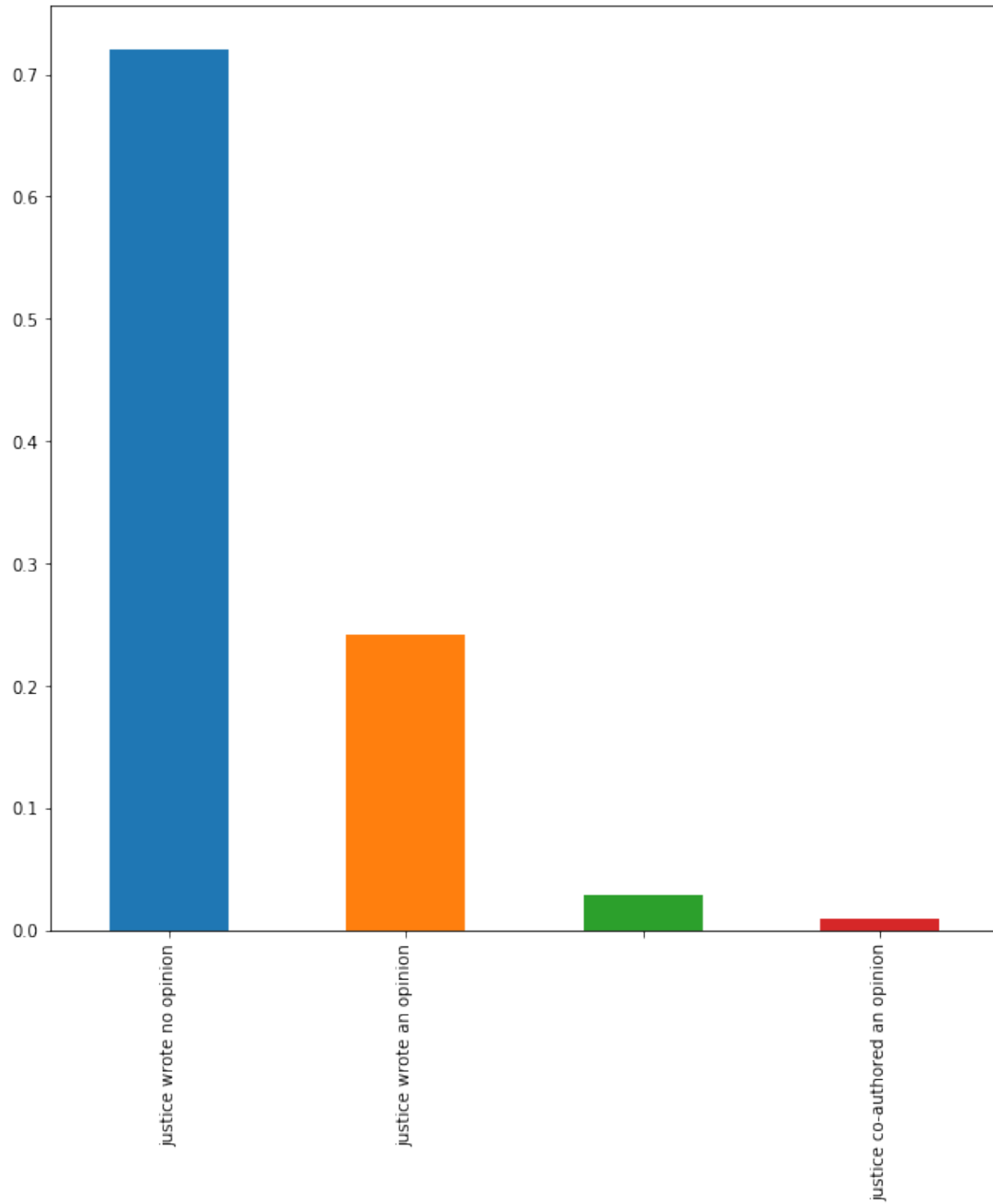
justiceName



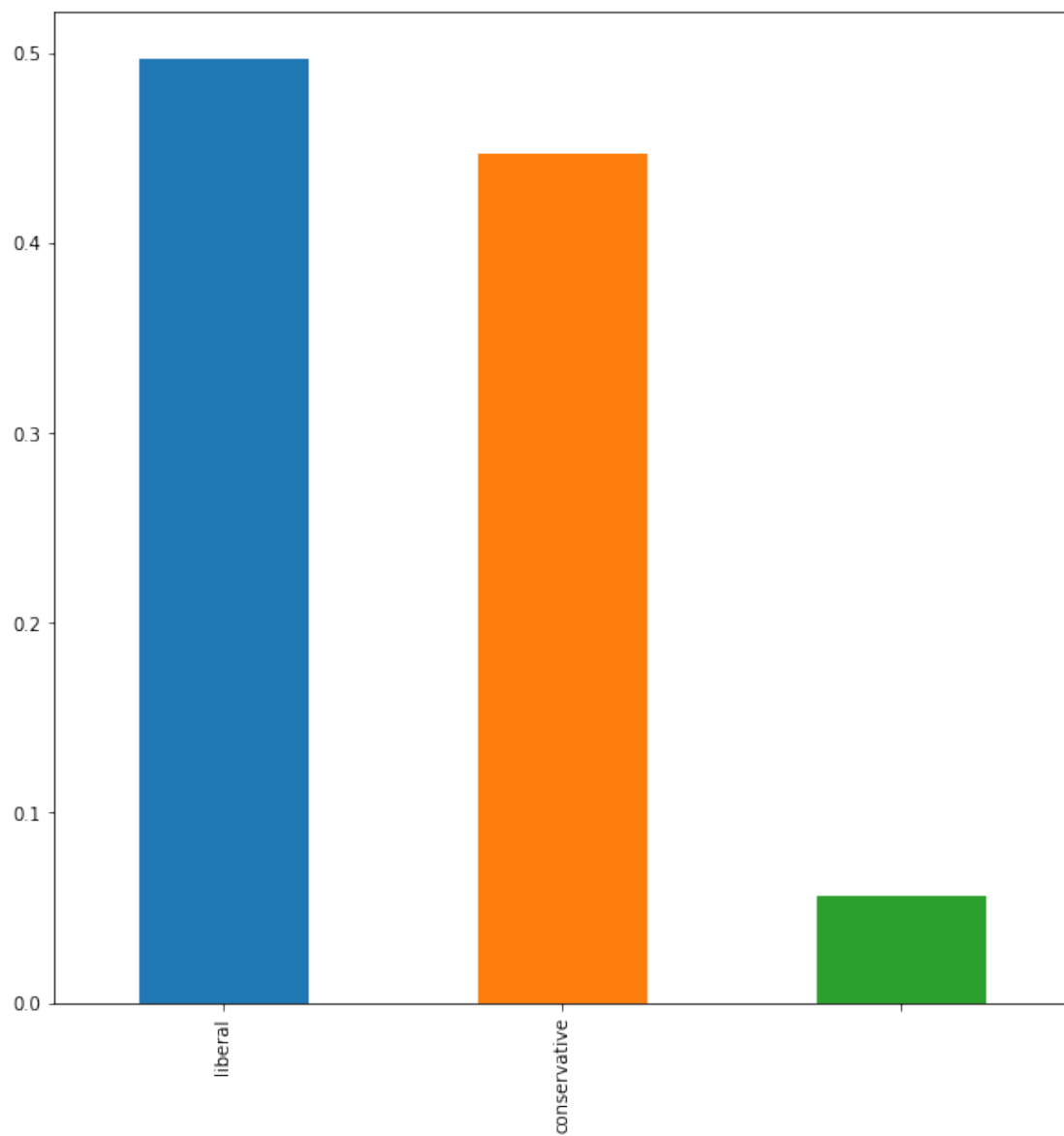
vote



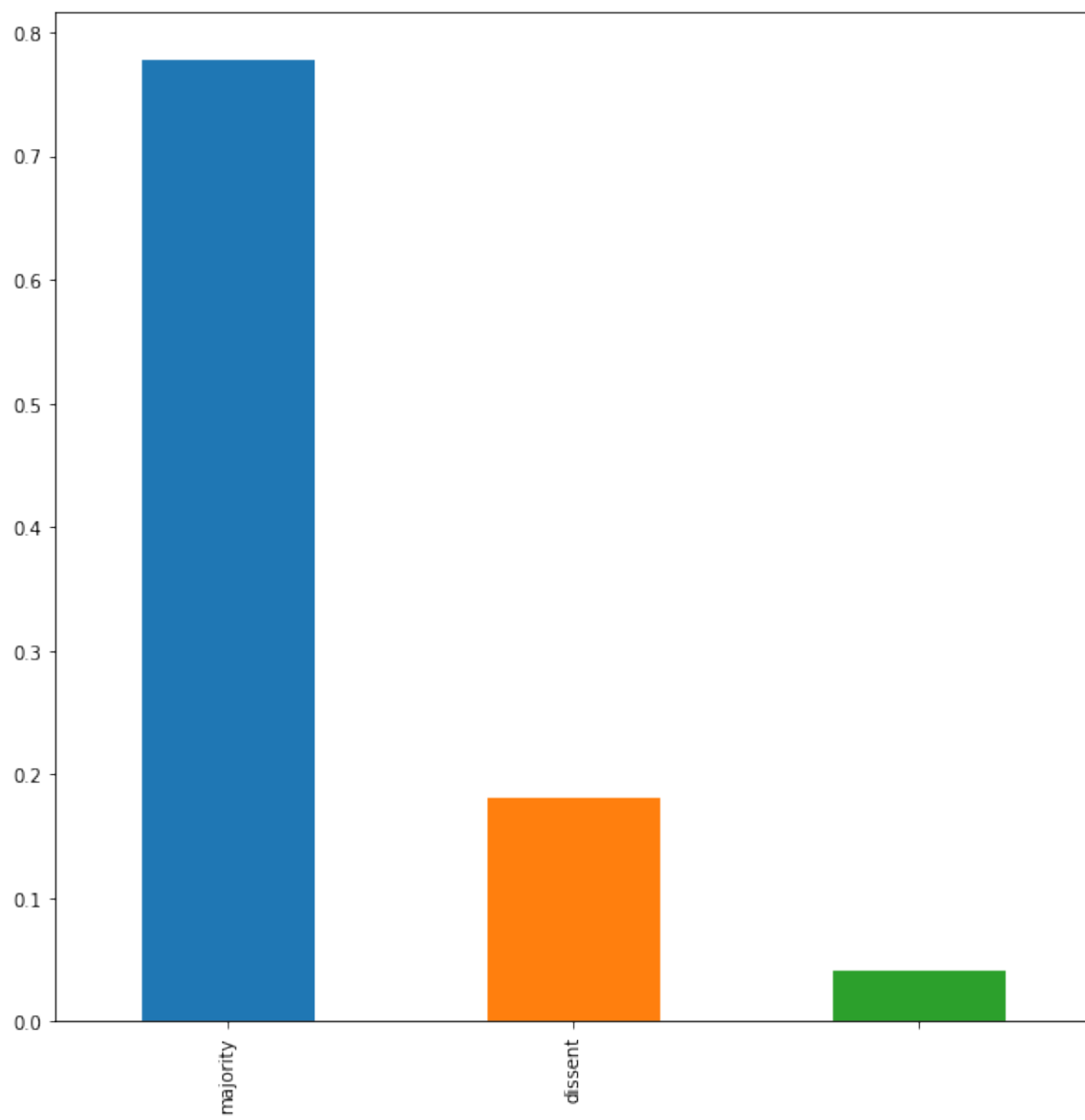
opinion



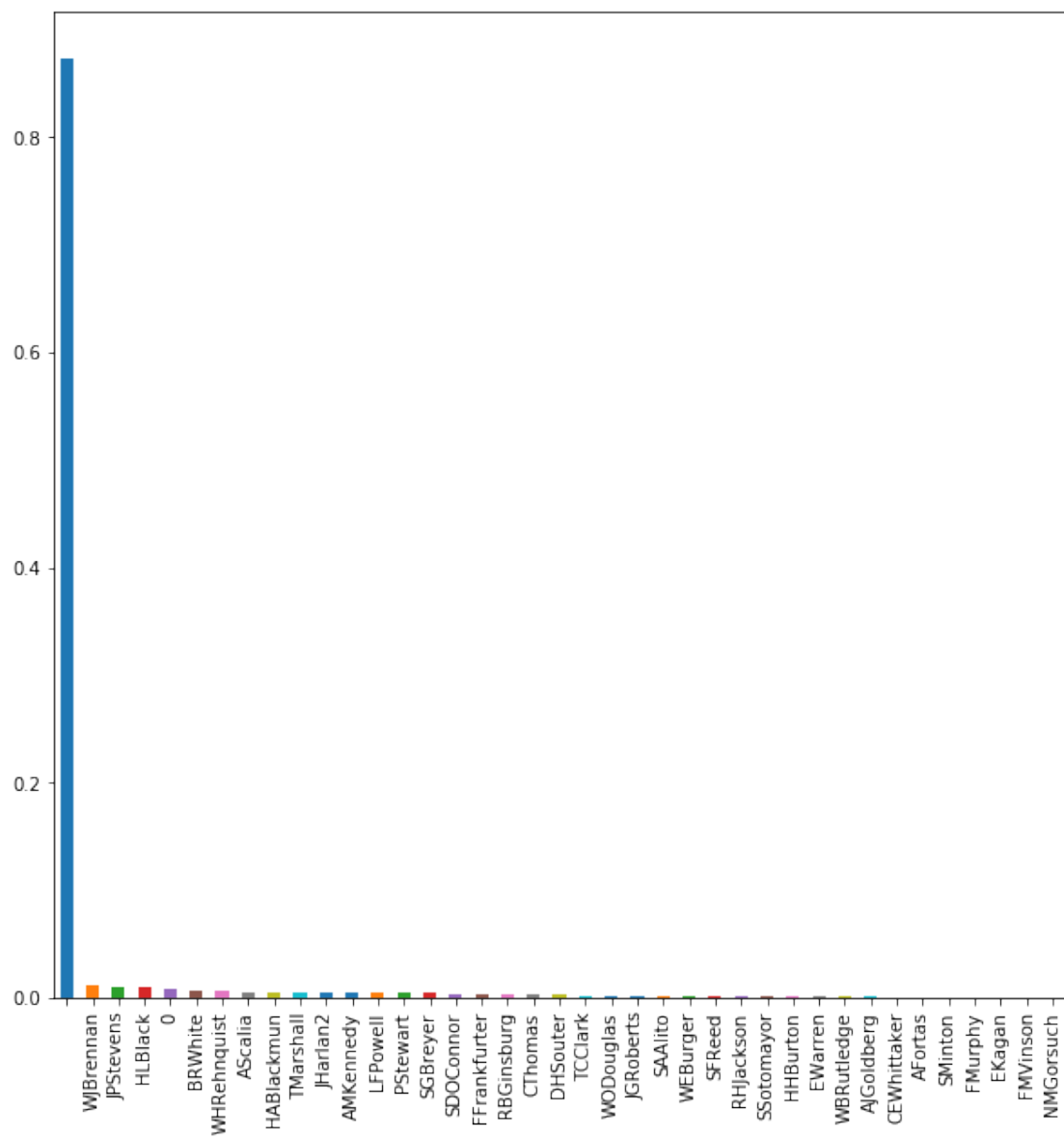
direction



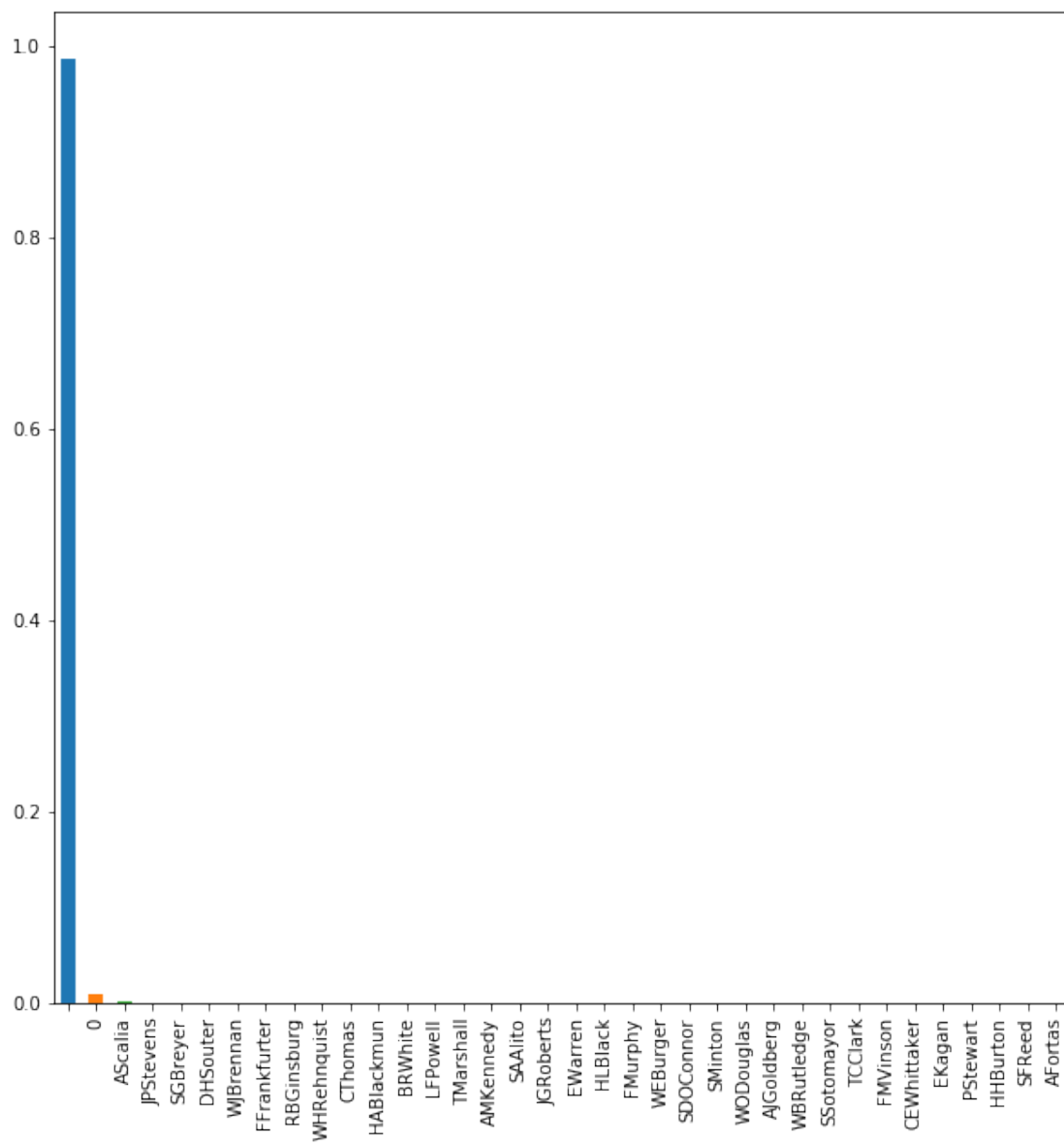
majority



firstAgreement



secondAgreement



In []:

```
df['firstAgreement'].value_counts(normalize=True).plot(kind="bar", figsize=(10, 10))
```

In []: