

Plotting alignment data

```
In [1]: %matplotlib inline
import matplotlib.pyplot as plt
import numpy as np
import json
import csv
import math
import pandas as pd
import utils.db_utils as db
import utils.plot_utils as plot
import utils.file_utils as file
import config

#####

cfig = config.getConfig() # configure values in config.js
#####

targetLang = cfig['targetLang']
bibleType = cfig['targetBibleType']
tWordsTypeList = cfig['tWordsTypeList']
dbPath = cfig['dbPath']
trainingDataPath = cfig['trainingDataPath']
testamentStr = cfig['testamentStr']
baseDataPath = cfig['baseDataPath']
```

```
In [2]: # get alignments for tWords

minAlignments = 20
remove = ['ó', 'ró', 'rà', 'òùòç', 'Àÿu', 'ùç', 'pé', 'éç']

alignmentsForWord, filteredAlignmentsForWord = db.fetchAlignmentDataForAllTWordsCached(trainingDataPath, bib
print(f"Original Language Alignments: {len(filteredAlignmentsForWord)}")

Using cached Alignments
Unfiltered Alignments: 4368
filtered alignments by original list count is 243
Size of filtered alignments by original ./data/en/ult/TrainingData/kt_en_ult_NT_alignments_by_orig_20.json i
7.628 MB
Size of filtered alignments by original ./data/en/ult/TrainingData/kt_en_ult_NT_alignments_by_orig_20.csv is
2.076 MB
Filtered Alignments: 243
Using cached Alignments
Unfiltered Alignments: 538
filtered alignments by original list count is 33
Size of filtered alignments by original ./data/en/ult/TrainingData/names_en_ult_NT_alignments_by_orig_20.jso
n is 1.298 MB
Size of filtered alignments by original ./data/en/ult/TrainingData/names_en_ult_NT_alignments_by_orig_20.csv
is 0.363 MB
Filtered Alignments: 33
Using cached Alignments
Unfiltered Alignments: 7380
filtered alignments by original list count is 250
Size of filtered alignments by original ./data/en/ult/TrainingData/other_en_ult_NT_alignments_by_orig_20.jso
n is 6.275 MB
Size of filtered alignments by original ./data/en/ult/TrainingData/other_en_ult_NT_alignments_by_orig_20.csv
is 1.670 MB
Filtered Alignments: 250
Original Language Alignments: 429
```

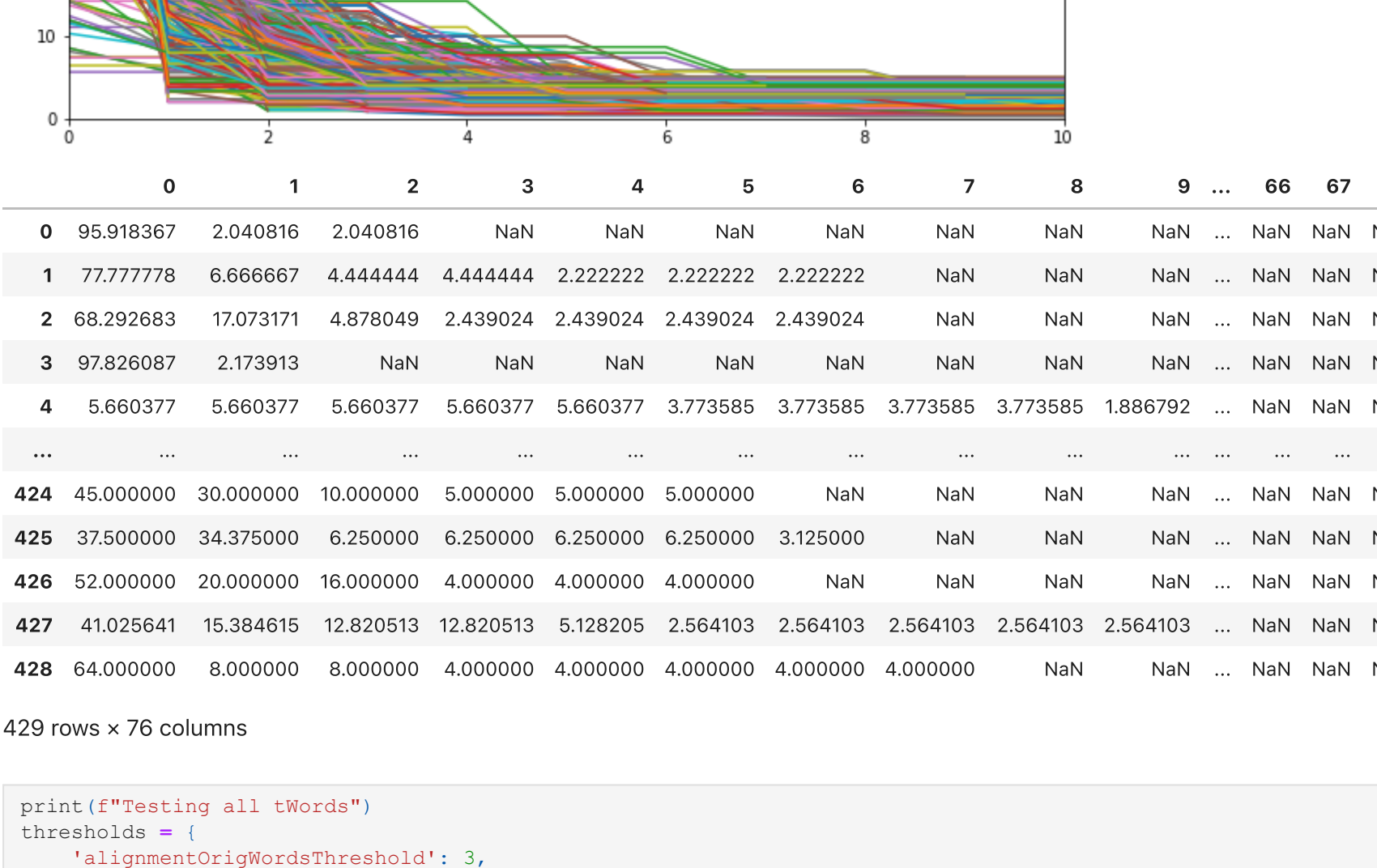
Analysis of alignments for tWords in the en_ult:

Frequency of alignments:

***Note that each line on the graphs below represents an alignment for a specific word. For example we have separate lines for 'Θεός', 'Θεός', or 'Θεοί' even though they have the same lemma. It made sense to group the alignments this way since aligners are likely to choose different target language words based on morphology of the word.

```
In [3]: frequenciesOfAlignments, stats = db.getFrequenciesOfFieldInAlignments(filteredAlignmentsForWord, 'alignmentT
print(f"Plotting of {len(filteredAlignmentsForWord)} tWord Alignments")
title = f"Plot of Variability of Specific Alignments in tWords"
ylabel = "Percent of Specific Alignments"
xlim = [0, 10]
ylim = [0, 75]
outputTable = plot.plotFrequencies(frequenciesOfAlignments, title, ylabel, showXValues=False, xlim=xlim,
csvPath = "PlotData_freqOfTWords.csv"
db.saveListToCSV(csvPath, outputTable)

Plotting of 429 tWord Alignments
```



```
In [4]: print(f"Testing all tWords")
thresholds = {
    'alignmentOrigWordsThreshold': 3,
    'alignmentTargetWordsThreshold': 5,
    'origWordsBetweenThreshold': 1,
    'targetWordsBetweenThreshold': 1,
    'alignmentFrequencyMinThreshold': 5
}

type = 'all twords'
warningData.summary = db.generateWarningsAndSummary(baseDataPath, type, bibleType, testamentStr, filteredAl
print(f"Found {len(warningData)} alignments to check - min threshold {minAlignments}")

frequencyWarnings = warningData[warningData['frequencyWarning'].str.len() > 0]
print(f"Found {len(frequencyWarnings)} frequencyWarnings")
frequencyWarningsByOrigWords = frequencyWarnings['originalWord'].value_counts()
print(f"FrequencyWarnings by original word:")
frequencyWarningsByOrigWords
```

Testing all tWords
saved summary of 429 original words to ./data/en/ult/all_twords_en_ult_NT_summary_20.csv
Found 1410 alignments to check - min threshold 20

Found 1148 frequencyWarnings
frequencyWarnings by original word:

```
Out[4]: θεός      65
         ἐνέτω    50
         Ἰησοῦς   37
         Χριστοῦ   33
         Ἰησοῦ    31
         ἸησοῦςΑυτοῦ  1
         δοῦναι    1
         ἐρόναι    1
         ἐροῦναι    1
         ὁδῆσθαι    1
         Χριστός   1
Name: originalWord, Length: 166, dtype: int64
```

```
In [5]: frequencyWarningsByLemma = frequencyWarnings['lemma'].value_counts()
print(f"FrequencyWarnings by lemma:")
frequencyWarningsByLemma
```

```
Out[5]: θεός      104
         Ἰησοῦς   77
         γλυπτοῦ   65
         ὁδῶν     50
         ἡμεῖς     45
         ...
         δοῦναι    1
         ἐρόναι    1
         ἐροῦναι    1
         ἐροῦναι    1
         ὁδῆσθαι    1
         ἸησοῦςΑυτοῦ  1
Name: lemma, Length: 96, dtype: int64
```

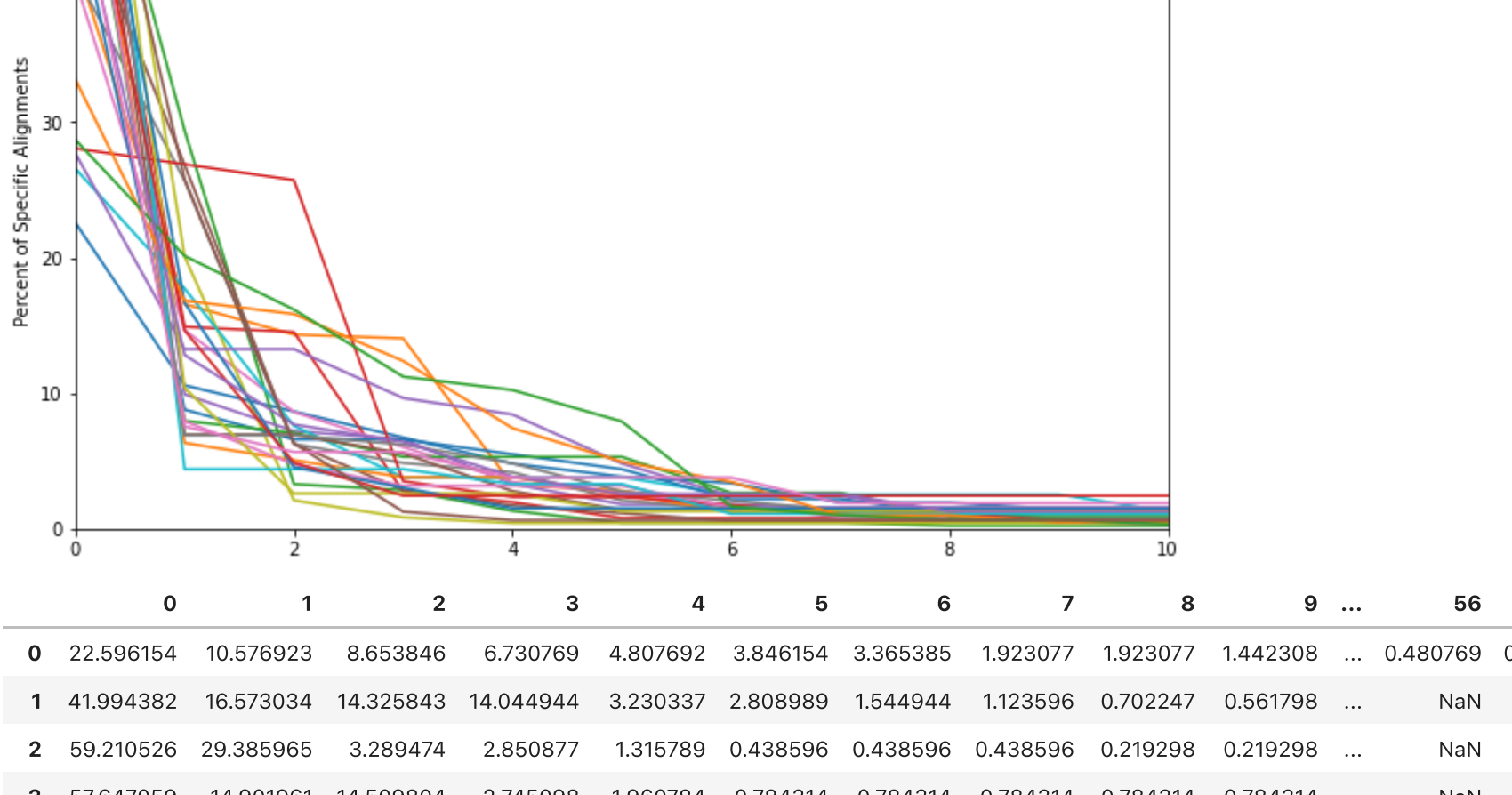
```
In [6]: minNumberOfWarnings = 10
origWordsWithWarnings = []
for key in frequencyWarningsByOrigWords.keys():
    count = frequencyWarningsByOrigWords[key]
    if count > minNumberOfWarnings:
        origWordsWithWarnings.append(key)

warningsAlignments = {}
for word in origWordsWithWarnings:
    warningsAlignments[word] = filteredAlignmentsForWord[word]
print(f"Found {len(origWordsWithWarnings)} original words with frequency warnings")

frequenciesOfAlignments, stats = db.getFrequenciesOfFieldInAlignments(warningsAlignments, 'alignmentText')

print(f"Plotting of {len(warningsAlignments)} tWord Alignments")
title = f"Plot of Variability of Alignments with Warnings"
ylabel = "Percent of Specific Alignments"
xlim = [0, 10]
ylim = [0, 75]
outputTable = plot.plotFrequencies(frequenciesOfAlignments, title, ylabel, showXValues=False, xlim=xlim,
csvPath = f"PlotData_freqOfTWords_minWarnings_{minNumberOfWarnings}.csv"
db.saveListToCSV(csvPath, outputTable)
```

Found 27 original words with frequency warnings
Plotting of 27 tWord Alignments



Analysis:

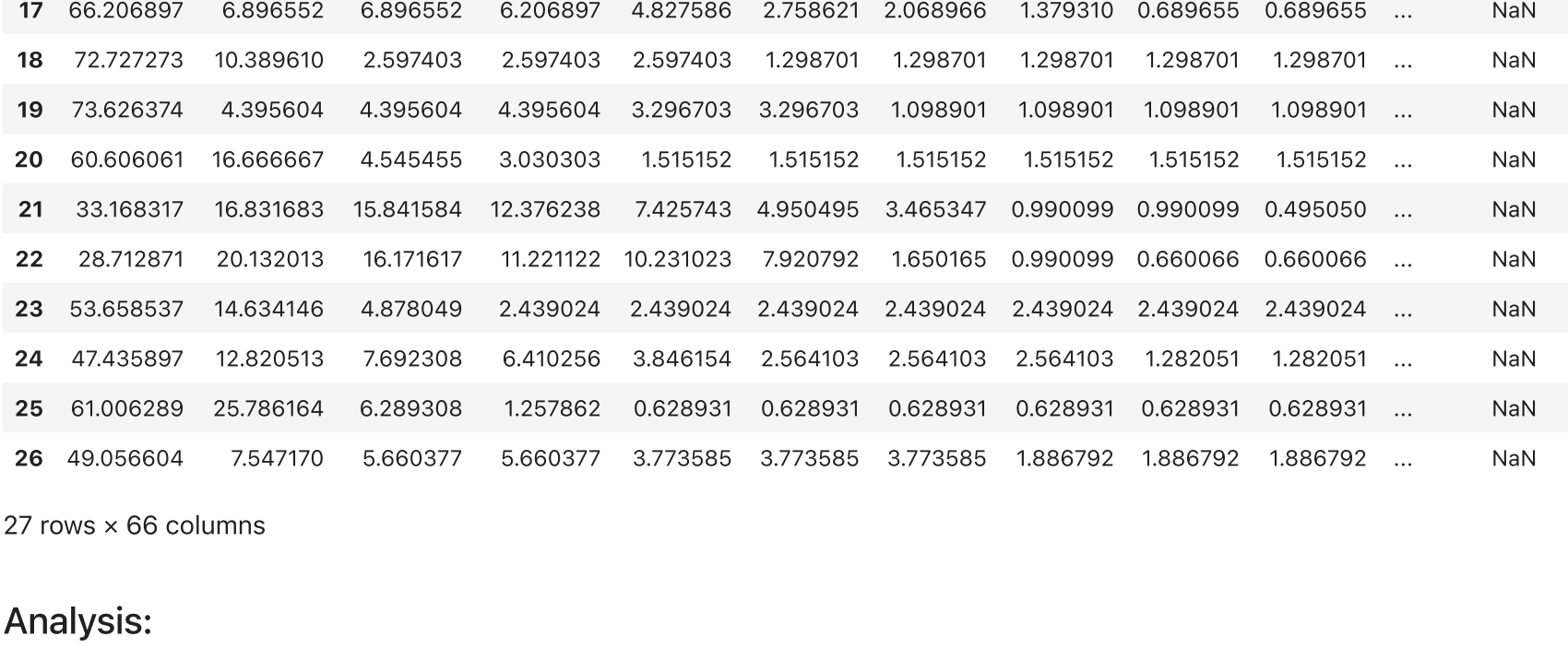
Analysis of numerical metrics:

Analysis of original language word count:

```
In [7]: type = 'all'
field = 'origWordsCount'
field_frequencies, stats = db.getFrequenciesOfFieldInAlignments(filteredAlignmentsForWord, field, sortIndex
filledFrequencies = db.zeroFillFrequencies(field_frequencies)

print(f"Found {len(field_frequencies)} original language words for tW type {type}")
title = f"Plot of number of Original Language Words in Specific Alignments in tW type {type}"
ylabel = "Percent of Specific Alignments"
xlabel = "Original Language Words"
plot.plotXYdataDict(filledFrequencies, title, ylabel, xlabel, showXValues=True, xlim=[1, 8])

Found 429 original language words for tW type all
```



Notes:

- this field analysis suggests that original word counts are tight - a threshold word count of 3 probably good for Greek to flag for review.

```
In [8]: threshold = 4
abnormalAlignments = {}
for origWord in field_frequencies:
    frequency = field_frequencies[origWord]
    count = len(frequency)
    if count >= threshold:
        abnormalAlignments[origWord] = frequency

print(f"Out of {len(field_frequencies)}, found {len(abnormalAlignments)} original language words that have i
filledFrequencies = db.zeroFillFrequencies(abnormalAlignments)

title = f"Plot of abnormal number of Original Language Words in Specific Alignments in tW KeyTerms"
ylabel = "Percent of Specific Alignments"
xlabel = "Original Language Words"
plot.plotXYdataDict(filledFrequencies, title, ylabel, xlabel, showXValues=True, xlim=[threshold, 10], ylim

Out of 429, found 5 original language words that have instances with over 4 words
```

Plot of abnormal number of Original Language Words in Specific Alignments in tW KeyTerms

