Plotting alignment data In [1]: %matplotlib inline import pandas as pd import utils.db utils as db import utils.plot utils as plot targetLang = 'en' bibleType = 'en ult' dbPath = f'./data/{bibleType} alignments.sqlite' connection = db.initAlignmentDB(dbPath) Connection to SQLite DB successful def getDataFrameForOriginalWords(connection, word, searchLemma = True): # print(f"searchLemma = {searchLemma}") words = [word] alignments = db.getAlignmentsForOriginalWords(connection, words, searchLemma) alignments = pd.DataFrame(alignments [word]) return alignments In [3]: # find all alignments for this lemma word = 'θεός' lemmaAlignments = getDataFrameForOriginalWords(connection, word, searchLemma = **True**) lemmaAlignments updating 'θεός' splitLemmasAndAddData - found original word ' $\Theta \epsilon \delta \varsigma$ ' for lemma splitLemmasAndAddData - found original word ' $\Theta\epsilon\delta\varsigma$ ' for lemma  ${\tt splitLemmasAndAddData - found\ original\ word\ '\Thetaεο\~υ'\ for\ lemma}$ Replacing "God s" with "God's" Replacing "the things that are God's" with "the things that are God's" Replacing "God s" with "God's" Replacing "God s" with "God's" Replacing "of God s" with "of God's" Replacing "God s" with "God's" Replacing "of God s" with "of God's" Replacing "God s" with "God's" Replacing "with God's words" with "with God's words" Replacing "God s" with "God's" Replacing "of God s" with "of God's" Replacing "God s" with "God's" splitLemmasAndAddData - found original word ' $\Theta\epsilon\delta\nu$ ' for lemma  $\verb|splitLemmasAndAddData - found original word '0εον' for lemma|\\$ splitLemmasAndAddData - found original word ' $\Theta\epsilon\tilde{\phi}$ ' for lemma  $\texttt{splitLemmasAndAddData - found original word '} \Theta \epsilon \epsilon \text{' for lemma}$  $\verb|splitLemmasAndAddData - found original word '\theta \verb|soi'| for lemma|\\$ splitLemmasAndAddData - found original word ' $\theta\epsilon$ oùç' for lemma  $\texttt{splitLemmasAndAddData-found original word '}\theta \texttt{so\~u'} \texttt{ for lemma}$  $\verb|splitLemmasAndAddData - found original word '\theta \verb|sol' for lemma||$ splitLemmasAndAddData - found original word ' $\theta\epsilon\delta\nu$ ' for lemma <code>splitLemmasAndAddData</code> - found original word ' $\theta\epsilon\delta\nu$ ' for lemma splitLemmasAndAddData - found original word ' $\theta\epsilon \delta\varsigma$ ' for lemma splitLemmasAndAddData - found original word ' $\theta\epsilon$ oĩς' for lemma id book\_id chapter verse alignment\_num orig\_lang\_keys target\_lang\_keys origSpan origWords origWordsTxt alig Out[3]: [{'id': 402, 'book\_id': 0 326 23 mat 16 ,402,403, ,518, 1 'mat', ό Θεός 'chapter': '1',... [{'id': 19162, 'book\_id': 7 17427 2 10 ,19162,19163, ,24951, mrk ό Θεός 'mrk', 'chapter': '2... [{'id': 24891, 'book\_id': 22758 mrk 10 18 10 ,24891,24892, ,32410, ό Θεός 'mrk', 'chapter': '1... [{'id': 29156, 'book\_id': 26691 15 34 17 ,29156,29157, mrk ,38064, ό Θεός 'mrk', 'chapter': '1... [{'id': 29159, 'book\_id': 26693 mrk 15 34 19 ,29159,29160, ,38066, ό Θεός 'mrk', 'chapter': '1... [{'id': 77510, 'book\_id': 1349 72801 19 37 ,101074, 11 ,77510, θεον act 'act', 'chapter': '1... [{'id': 111284, 'book\_id': **1350** 105300 2th 2 4 6 ,111284, ,147759, θεον '2th', 'chapter': [{'id': 82744, 'book\_id': 1351 77911 28 27 6 ,82744, ,108220,108221, act θεόν 'act', 'chapter': '2... [{'id': 98291, 'book\_id': 92604 2co 1352 3 ,98291, ,130171, θεος '2co', 'chapter': [{'id': 102952, 'book\_id': **1353** 97161 gal 4 8 11 ,102952, ,136824, θεοῖς 'gal', 'chapter': 1354 rows × 20 columns In [4]: # find all alignments for this original word # word = ' $\Theta \varepsilon \delta \varsigma$ ' # found 69 # word = ' $\Theta \varepsilon \delta c$ ' # found 239 word =  $'\Theta \epsilon o \tilde{v}'$  # found 712 origAlignments = getDataFrameForOriginalWords(connection, word, searchLemma = False) origAlignments updating 'Θεοῦ' Replacing "God s" with "God's" Replacing "the things that are God's" with "the things that are God's" Replacing "God s" with "God's" Replacing "God s" with "God's" Replacing "of God s" with "of God's" Replacing "God s" with "God's" Replacing "of God s" with "of God's" Replacing "God s" with "God's" Replacing "with God's words" with "with God's words" Replacing "God s" with "God's" Replacing "of God s" with "of God's" Replacing "God s" with "God's" id book\_id chapter verse alignment\_num orig\_lang\_keys target\_lang\_keys origSpan origWords origWordsTxt align Out[4]: [{'id': 1202, 'book\_id': 1047 3 16 16 ,1202,1203, ,1537,1538, mat τοῦ Θεοῦ 'mat', 'chapter': [{'id': 1259, 'book\_id': 1099 ,1611,1612, mat ,1259,1260, τοῦ Θεοῦ 'mat', 'chapter': '4'... [{'id': 1287, 'book\_id': 2 1125 mat 4 19 ,1287, ,1645,1646, Θεοû 'mat', 'chapter': '4'... [{'id': 1311, 'book\_id': ,1311,1312, 3 1149 4 6 6 mat ,1676,1677, 'mat', τοῦ Θεοῦ 'chapter': '4'... [{'id': 1738, 'book\_id': 1534 5 9 ,1738, ,2227,2228, Θεοû 'mat', 'chapter': '5'... [{'id': 137161, 'book\_id': **707** 129469 23 ,137161,137162, ,182348, 10 τοῦ Θεοῦ 'rev', 'chapter': [{'id': 137166, 'book\_id': **708** 129473 rev 21 11 3 ,137166,137167, ,182353,182354, τοῦ Θεοῦ 'rev', 'chapter': [{'id': 137431, 'book\_id': ,182674,182675, **709** 129716 21 23 16 ,137431,137432, τοῦ Θεοῦ rev 'rev', 'chapter': [{'id': 137522, 'book\_id': ,182789,182790, **710** 129796 22 1 13 ,137522,137523, τοῦ Θεοῦ rev 'rev', 'chapter': [{'id': 137568, 'book\_id': **711** 129835 22 3 9 ,137568,137569, ,182847,182848, τοῦ Θεοῦ rev 'rev', 'chapter': 712 rows × 20 columns db.describeAlignments(origAlignments) Alignments description: origSpan alignmentOrigWords targetSpan alignmentTargetWords \ count 712.000000 712.000000 712.000000 712.000000 0.623596 1.620787 0.599719 1.599719 mean 0.501928 0.494146 0.515471 0.515471 0.000000 1.000000 0.000000 1.000000 min 25% 0.000000 1.000000 0.000000 1.000000 50% 1.000000 2.000000 1.000000 2.000000 75% 2.000000 1.000000 1.000000 2.000000 3.000000 3.000000 4.000000 5.000000 max frequency matchCount origWordsBetween targetWordsBetween count 712.000000 712.000000 712.000000 0.241384 171.865169 mean 0.002809 0.0 0.147977 105.359575 0.074953 0.0 std min 0.001404 1.000000 0.000000 0.0 98.000000 25% 0.137640 0.000000 0.0 0.176966 126.000000 50% 0.000000 0.0 0.410112 292.000000 75% 0.000000 0.0 0.410112 292.000000 2.000000 0.0 fields = ['origSpan', 'alignmentOrigWords', 'targetSpan', 'alignmentTargetWords', 'frequency', 'matchCount', 'origWordsBetween', 'targetWordsBetween'] Frequency of origSpan: 435 273 0 2 3 Name: origSpan, dtype: int64 Frequency of alignmentOrigWords: 1 273 3 3 Name: alignmentOrigWords, dtype: int64 Frequency of targetSpan: 1 417 0 291 2 1 Name: targetSpan, dtype: int64 Frequency of alignmentTargetWords: 417 1 291 3 3 5 1 Name: alignmentTargetWords, dtype: int64 Frequency of matchCount: 292 126 126 102 102 98 98 23 23 21 21 1 2.0 11 11 5 5 4 4 2 2 Name: matchCount, dtype: int64 Frequency of origWordsBetween: 0 711 Name: origWordsBetween, dtype: int64 Frequency of targetWordsBetween: Name: targetWordsBetween, dtype: int64 Out[5]: {'desc': {'origSpan': {'count': 712.0, 'mean': 0.6235955056179775, 'std': 0.5019282025771814, 'min': 0.0, '25%': 0.0, '50%': 1.0, '75%': 1.0, 'max': 3.0}, 'alignmentOrigWords': {'count': 712.0, 'mean': 1.6207865168539326, 'std': 0.4941462676438909, 'min': 1.0, '25%': 1.0, '50%': 2.0, 175%': 2.0, 'max': 3.0}, 'targetSpan': {'count': 712.0, 'mean': 0.5997191011235955, 'std': 0.5154708412840628, min': 0.0, '25%': 0.0, '50%': 1.0, '75%': 1.0, 'max': 4.0}, 'alignmentTargetWords': {'count': 712.0, 'mean': 1.5997191011235956, 'std': 0.5154708412840628, 'min': 1.0, '25%': 1.0, '50%': 2.0, '75%': 2.0, 'max': 5.0}, 'frequency': {'count': 712.0, 'mean': 0.2413836636788284, 'std': 0.14797693103697737, 'min': 0.0014044943820224719, '25%': 0.13764044943820225, '50%': 0.17696629213483145, '75%': 0.4101123595505618, 'max': 0.4101123595505618}, 'matchCount': {'count': 712.0, 'mean': 171.86516853932585, 'std': 105.35957489832789, 'min': 1.0,
'25%': 98.0, '50%': 126.0, '75%': 292.0, 'max': 292.0}, 'origWordsBetween': {'count': 712.0, 'mean': 0.0028089887640449437, 'std': 0.07495316889958614, 'min': 0.0, '25%': 0.0, '50%': 0.0, 175%': 0.0, 'max': 2.0}, 'targetWordsBetween': {'count': 712.0, 'mean': 0.0, 'std': 0.0, 'min': 0.0, 125%!: 0.0, '50%': 0.0, 175%': 0.0, 'max': 0.0}}, 'fields': {'origSpan': [435, 273, 3, 1], 'alignmentOrigWords': [436, 273, 3], 'targetSpan': [417, 291, 3, 1], 'alignmentTargetWords': [417, 291, 3, 1], 'matchCount': [292, 126, 102, 98, 23, 21, 20, 11, 8, 5, 4, 2], 'origWordsBetween': [711, 1], 'targetWordsBetween': [712]}} Analysis of alignments for Θεοῦ in the en\_ult: Frequency of alignments: In [6]: frequency = origAlignments['alignmentTxt'].value\_counts() print(frequency) τοῦ Θεοῦ = of God 292  $\Theta$ εοῦ = God 126 Θεοῦ = of God 102 τοῦ Θεοῦ = God 98 Θεοῦ = God's 23 τοῦ Θεοῦ = God's 21 τοῦ Θεοῦ = of 11 8 Θεοῦ = of  $\Theta$ εοῦ = of the 5 τοῦ Θεοῦ = from God τοῦ Θεοῦ = of God's λόγια Θεοῦ = with God's words 1 μὲν Θεοῦ = of God 1 Θεοῦ = of a god ότι Θεοῦ = of God 1  $\Theta$ εοῦ = Father 1 Θεοῦ = by God 1 τοῦ Θεοῦ = against God  $\Theta \epsilon \circ \tilde{\mathbf{u}} = \text{for God}$ 1 = and God ένώπιον τοῦ Θεοῦ 1 τὴν τοῦ Θεοῦ = of God's Θεοῦ = in God  $\Theta$ εοῦ = the Son  $\Theta \epsilon \circ \tilde{v} = the$ τοῦ Θεοῦ = about God τὰ τοῦ Θεοῦ = the things that are God's  $\Theta \epsilon \circ \tilde{v} = \text{because of God}$ τοῦ Θεοῦ = for God Θεοῦ = godly 1 τοῦ Θεοῦ = godly 1 τοῦ Θεοῦ = to God 1 Name: alignmentTxt, dtype: int64 Notes: the left column is the specific alignment, and the right column is the number of times that specific alignment has been made so far in the NT. alignments that contain more words are more suspect. in future will combine "God s" to "God's" before doing analysis plot.plotFieldFrequency(frequency, "", 'alignment', title="Frequency of Alignments", xNumbers=False, xShowTi Frequency of Alignments 300 250 200 150 150 100 50 alignment **Analysis: Analysis of numerical metrics:** descr = origAlignments.describe() print(f"Alignments description:\n{descr}") Alignments description: origSpan alignmentOrigWords targetSpan alignmentTargetWords count 712.000000 712.000000 712.000000 712.000000 mean 0.623596 1.620787 0.599719 1.599719 0.501928 0.494146 0.515471 0.515471 std 0.000000 1.000000 min 0.000000 1.000000 25% 0.000000 1.000000 0.000000 1.000000 50% 1.000000 2.000000 1.000000 2.000000 2.000000 75% 2.000000 1.000000 1.000000 3.000000 4.000000 5.000000 3.000000 max frequency matchCount origWordsBetween targetWordsBetween count 712.000000 712.000000 712.000000 712.0 0.002809 mean 0.241384 171.865169 0.0 0.147977 105.359575 std 0.074953 0.0 min 0.001404 1.000000 0.000000 0.0 0.137640 98.000000 25% 0.000000 0.0 0.176966 126.000000 0.0 50% 0.000000 0.410112 292.000000 0.000000 75% 0.0 2.000000 0.410112 292.000000 Analysis of distance between first and last original language word: In [9]: field = 'origSpan' field frequency = origAlignments[field].value counts().sort index() print(f"\nFrequency of {field}:\n{field frequency}") Frequency of origSpan: 273 435 2 3 3 1 Name: origSpan, dtype: int64 Notes: • this field is less useful because it includes aligned words, so added originalWordsBetween as more useful normalized metric (see analysis below). Analysis of distance between first and last target language word: field = 'targetSpan' field frequency = origAlignments[field].value counts().sort index() print(f"\nFrequency of {field}:\n{field frequency}") Frequency of targetSpan: 0 291 417 2 3 1 Name: targetSpan, dtype: int64 Notes: this field is also less useful because it includes the aligned words, so added targetWordsBetween as more useful normalized metric (see analysis below). Analysis of original language word count: field = 'alignmentOrigWords' field\_frequency = origAlignments[field].value\_counts().sort\_index() print(f"\nFrequency of {field}:\n{field\_frequency}") Frequency of alignmentOrigWords: 273 436 3 Name: alignmentOrigWords, dtype: int64 Notes: this field analysis suggests for θεός nearly all the original language word counts are tight. The word counts of 3 may need review. So we could probaby use that as a threshold for to flag for review. Analysis of target language word count: field = 'alignmentTargetWords' field frequency = origAlignments[field].value counts().sort index() print(f"\nFrequency of {field}:\n{field\_frequency}") Frequency of alignmentTargetWords: 1 291 417 3 3 1 Name: alignmentTargetWords, dtype: int64 Notes: this field analysis suggests that for θεός likely all the target language word counts are tight. The word count of 3 probably good for English ( of a god ). But still we could probaby use that as a threshold for to flag for review. Analysis of count of extra unaligned words between aligned original language words: field = 'origWordsBetween' field\_frequency = origAlignments[field].value\_counts().sort\_index() print(f"\nFrequency of {field}:\n{field\_frequency}") Frequency of origWordsBetween: U  $/ \perp \perp$ 2 1 Name: origWordsBetween, dtype: int64 Notes: this field analysis suggests that most original language alignments probably good. Probably the cases of a word between (count > 0) aligned words should be reviewed. Analysis of count of extra unaligned words between aligned target language words: In [14]: field = 'targetWordsBetween' field\_frequency = origAlignments[field].value\_counts().sort\_index() print(f"\nFrequency of {field}:\n{field\_frequency}") plot.plotFieldFrequency(field\_frequency, field, f"Words Between", max=10) Frequency of targetWordsBetween: Name: targetWordsBetween, dtype: int64 Frequency of Words Between ('targetWordsBetween') 700 600 500 400 300 200 100 0 Words Between Notes: • this field analysis suggests that most target language alignments probably good. Large gaps between aligned words are likely due to wordmap suggesting wrong occurence of a word and the user selecting. Probably the cases of a word between (count > 0) aligned words should be reviewed.