

Exploratory data analysis

Для качественного исследования и хорошего результата требуется знание о данных.

Мы имеем следующие наборы данных:

1. overall_80K.csv - датасет из параллельных предложений рус↔манси, предоставленный кейсодателем.
2. mansi.csv - собранный нами монокорпус из [Луима Сэрипос](#).
3. dict.csv - распарсенный нами мансийско-русский словарь.
4. Синтетические данные.

Начнем по порядку.

overall_80K.csv

Данные параллельного корпуса, предоставляемые кейсодателем. Имеет 81,146 пар предложений рус↔манси.

- А. От кейсодателя была получена информация о наличии кейсов, где слова написаны через пробел (пример: Г Е Р М А Н И Я). Для начала мы решили проверить количество этих кейсов используя регулярные выражения. Итого у нас вышло 847 различных предложений с данным дефектом. Данные кейсы были удалены из-за невозможности восстановить автоматически предложение.

Из примеров видно, что далеко не всегда можно определить сколько слов написаны в таком стиле, чтобы правильно преобразовать предложение.

Примеры:

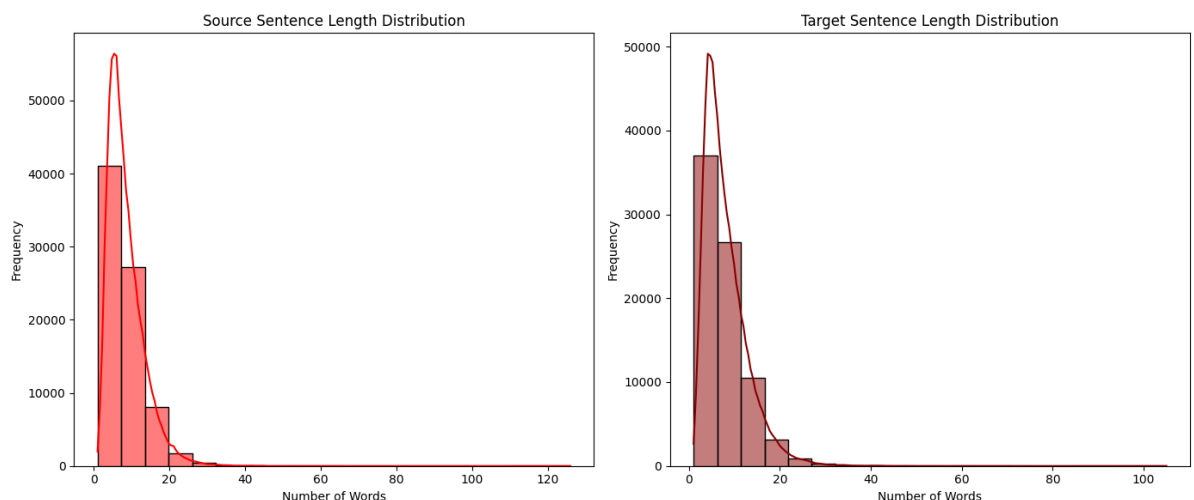
Тõн тйисмõйттыт ЮванНиколаевич такви щёпитасанэ.
Павлув коныпал ты ляпат õлнэ мǎнь пǎвлытныл Нумто, Юильск ос сǎлыу мǎхум нǎвраманыл тыг Касумн т оты глы яныл .
Т у в ы л ā г м ы џ м о с ы џ нǎврамыт ёт рўпитэв.
Ань Урай ўсныл ос К он д ин с к и й р а й о н ы л õ л н э м ā х у м в ā т и х а л м ā н ь щ и пǎвылт ўщлахтэгыт.
Маснут ёнтуңкве с а к а а т ā л ы м ё г ы т.
Мõлты тǎл «Сургутнефтегаз» 1 5 х а н т ы н ъ в р а м ўщлахтын мǎгсыл олн тǎстыглас.
А к в а н а т ы м ол н ы т ул вос тўлмантавет, м о р и ул в о с х ол т а в е т, о б щ е с т в е н н ы й организацият ты рўпата вǎрнэ мǎхум õс уральтаңкве патыяныл.
Т у в ы л И в д е л ь ў с мус машинал минм ы г т а с м ё н .
Ты т ā л с а к а с ā в хõтпа янытлавес, тǎн халанылт сǎв мǎньлат хõтпа блыс.

- В. Значений NaN в датасете нет, есть полные дубликаты - их 251. Дубликаты были удалены, оставив только 1 версию.
- С. Средняя длина предложений на мансийском составляет 8 слов, а максимальная целых 188, что слишком сильно разнится со средней длиной. Оказалось, что текст просто не был разбит на предложения:

“34. Ирина Константиновна Поята Хальӯс район Лӧпмус пӧвылт самын патыс. 35. Аще Константин Корнилович, омаӧ Ольга Максимовна Албиныг ӧлсӧг. 36. Экваг-ӧйкаг колтӧглӧнт китхуйплов няврам янмалтасӧг. 37. Ань сӧт хӧтпа хультыс. 38. Ирина школа ӧстламе юи-пӧлт Салехард ӯс медучилищан ханищтахтукве минас. 39. Ӑстламе юи-пӧлт Ямал мӧн Тарко-Сале ӯс пӧльнищан рӧпитанкве кӧтвес. 40. Тот мощ рӧпитас ос хум вӧрыс. 41. Ӑйкатӧн Молдавия мӧн ӧлуӧкве тотвес. 42. Тувыл 1990 тӧлт тӧн ювле Хальӯсн щӧмьяӧ тӧгыл вӧнтлысӧг. 43. Нӧ пӧльнищат терапевт-лӧккарн нӧтым нӧловхуйплов тӧл рӧпитас. 44. Та юи-пӧлт ос физиотерапия вӧрмалъ щирыл ӧлалъ ханищтахтас. 45. Ань ты пӧсмаӧтан вӧрмалъ щирыл Хальӯс пӧльнищат китхуйплов тӧл рӧпиты. 46. Ирина Константиновна ӧйкатӧнтыл кит няврам янмалтасӧг. 47. Пыгӧн юридический академия ӧстлас. 48. Ань Хальӯст ӧлы. 49. Ӑгитӧн ос Ханты-Мансийск ӯс педколледжит нилыт тӧл ханищтахты. 50. Мӧн мӧньщи щӧмьят йильпи тӧл кастыл янытлыянув! 51. Ӑлупсанын кӧпнитыг вос ӧлы. 52. Йильпи тӧлт нӧн щуниӧгыг вос ӧмтӧгын. 53. Рӧтанын, юртанын ӧт сӧв тӧл пустӧгыл ӧлӧн.”

И это не единичный случай. Всего таких кейсов около 463х. Для всех этих кейсов с помощью регулярных выражений были убраны порядковые номера, а для предложений, чья длина больше 128, были разбиты на части поменьше. Кол-во предложений, чья длина больше 128 - 1.

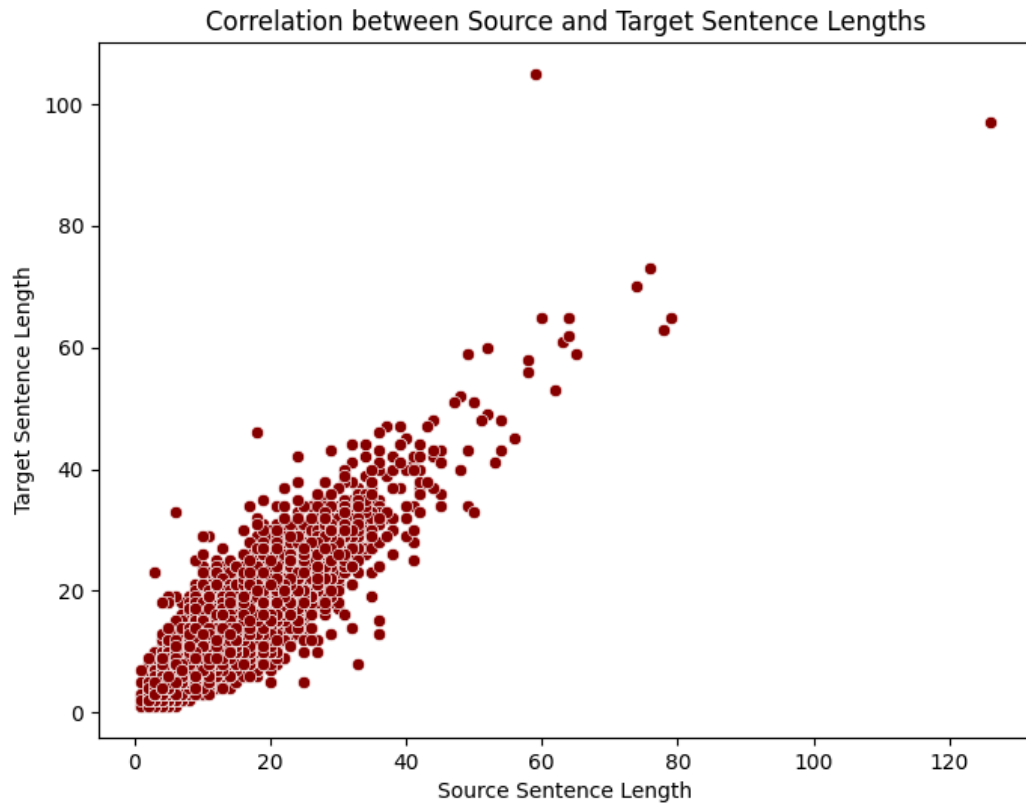
- Д. Теперь посмотрим на распределение длин предложений.



Из bar chart'a видно, что распределение похоже на экспоненциальное или

Хи-квадрат. Больше всего, конечно же, коротких предложений.

- Е. При визуализации корреляции длин источника и таргета было замечено сильное расхождение между семплами:



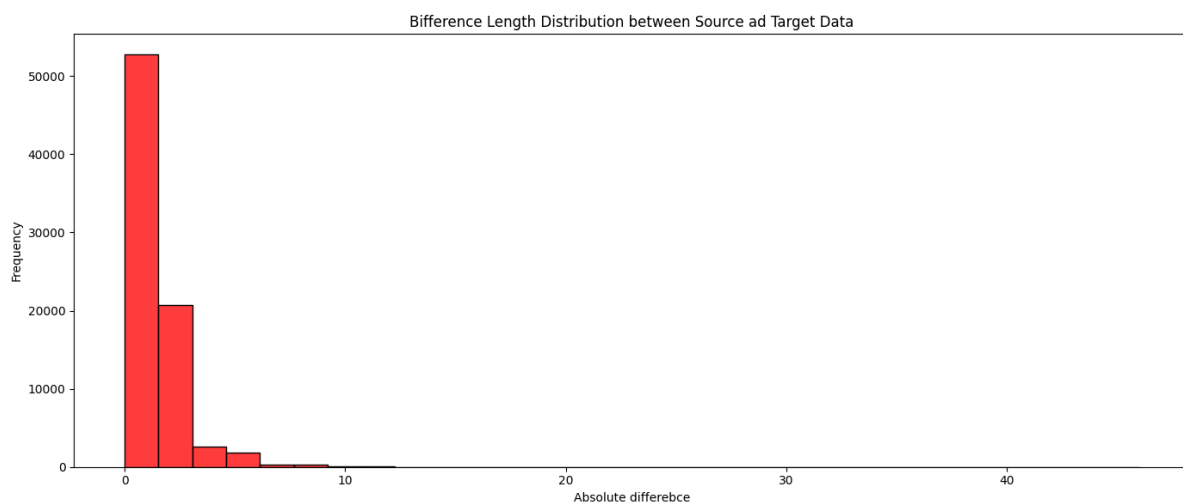
```
[]):
```

```
    differ.describe()
```

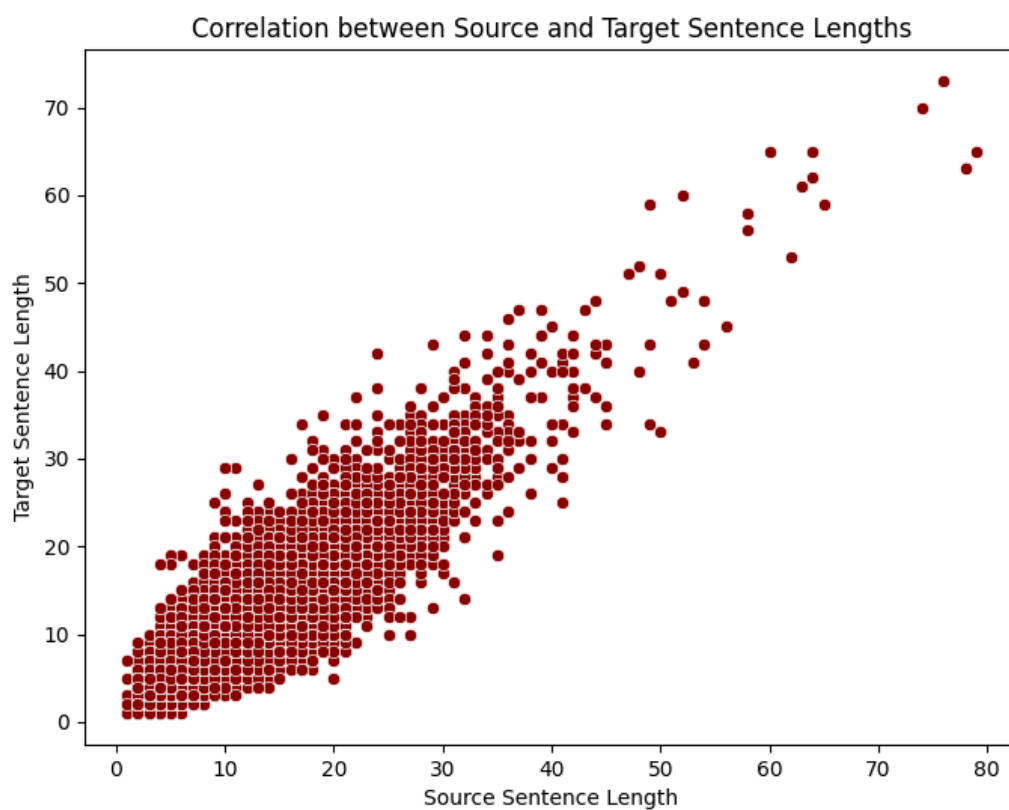
```
count    78702.000000
mean      1.322914
std       1.447477
min       0.000000
25%       0.000000
50%       1.000000
75%       2.000000
max       46.000000
dtype: float64
```

Был установлен трешхолд на разницу в длину равной 25. Видим следующие примеры:

Target	Source
<p>Ўлпыл , хӧлтан ӕлпылнува, ӱнтимен хӕпын - ам туп тарм, Маруся пӧсум сӧртнутын, Светлана ӕквӕт мӕннӕ хӧтпаг - и нӕтӕв ӕл тув, хӧтталь ӕри, лӕвегыт, ӕныг вӧр, хот ӕнымӕгыт ӕ вӕтат кӕтыг ӱми хӕлыг, хӧт хӧнтыс мӧл- хӧтал кӧл ӧнтсыл ӧлнӕ ӕгирись хӱрум ӕмас ӕнж лӕхс.</p>	<p>Отпустил коня Иванушка - дурачок и взял с него слово - пшеницы больше не есть и не топтать.</p>
<p>Тувле кос ӕмандас, луве лӕви: "Ул минӕн, тыг ӕиен! Наӱ тах хотум мори ты номылматӕн, ӕква ты сайкалы, ӧлум пасмен тӕй; ам нӕлссам асагумн сӕлтӕн, тот капак ӧлы, та капакта хурум сӕрка аен!"</p>	<p>Конь говорит: "Не уходи, иди сюда!"</p>
<p>купса Сибиряков (Щипрах) лӕнх оньшас - "Щипрах лӕнхыг" (Сибиряковский тракт) лавыглавес. Ялпын хоталт ӕхталас Халь кс куцай Фомин В.И., Хулюмсунт миркол куцай Ануфриев Я.В. Янитлавест махум хотит олсыт нӕлсат арыгкем тал. иильпи самын патум няврамыт. Янитлавест акван хасхатам (олмыгтам) махум. Сав потыртавест хоти махум олӕгыт хоса аквӕт (пурияныл- иив, келп, аргин, сорни). Няксимволь урыл ӕрыг хансум ӕква янитлавес. Няксимволь миркол куцай Волклва Т.К., потыртас, латын лавыс-хоти махум нетсыт ялпын хотал варункв. Няксимволь сав сунсылтаве Михаил Заплатин, Лев Вахитов - кинат, Игошев картинат мот хон мат тотыглавет, музей Ханты-Мансииск олӕгыт. карыс сип ватат- павлув,тагт я овтохти (хайти витӕ). Олнӕ вармалюв минанти, та олантев - павылприсювт!</p>	<p>Первым поселением была манси деревня. потом из-за урала в поисках лучших мест для жилья пришли коми-зыряне. Из воспоминаний Е.М.Носовой, которая была в числе первых переселенцев она рассказывала: шли из-за Урала от голода, надеялись на рыбные места, пушнину. Мы выжили благодаря Няксимволю. Сосьвинская пристань Сибирякова - это единственное в данной местности поселение русский пункт, отстоящий от Березово в 500 верстах.</p>



Было замечено, что предложения не были переведены до конца, а только первые части. Было сделано так: если $\text{target} > \text{source}$, то тогда берем первое предложение, иначе оставляем как есть. После данных манипуляций можем наблюдать следующую картину:



- Г. Также был проведен анализ частотности слов. Для русских предложений была проведена очистка от стоп слов и получилась такая картина:



Для мансийского языка мы не знаем стоп-слова, поэтому решили отобразить
напрямую частотность слов:



Видно, что здесь много “коротких” слов. Мы связались со специалистами, знающие мансийский язык и спросили их о наличии предлогов. Ответ говорил о том, что предлогов в мансийском нет, но есть превербы и послеслоги, скорее всего это они и есть в БОльшем случае.

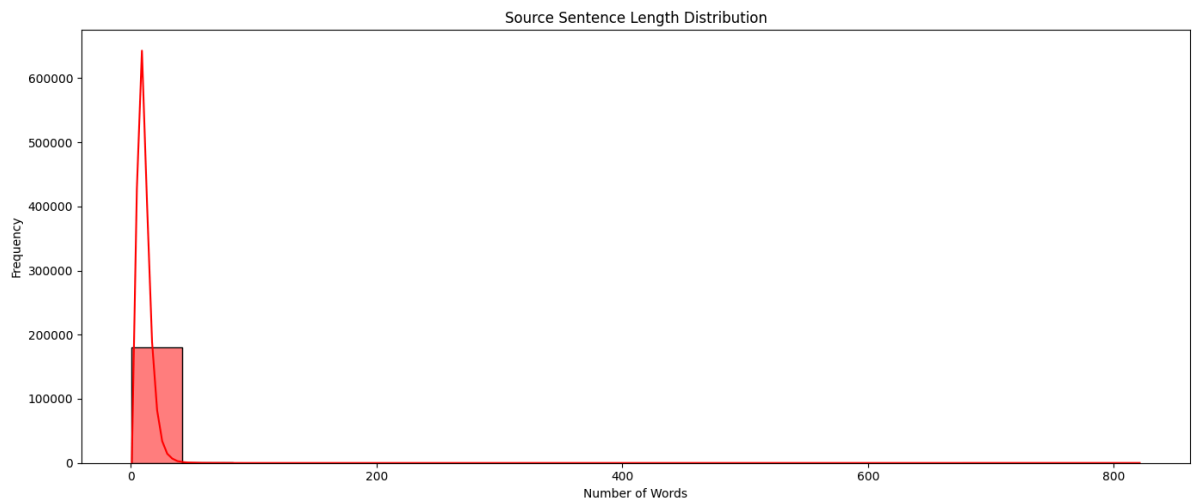
mansi.csv

Данный документ представляет из себя собранный монокорпус новостей на мансийском с сайта [Луима Сэрипос](#), а также данные ученого [Csilla Horvath](#).
Использовался для обучения модели m2m100-418 на задаче m1m и в последующем

для генерации синтетических данных для дообучения модели.

А. Датасет разбит на 183124 предложения с помощью модуля nltk, предобработан тем, что убран спец символ '\xad'. Таблица имеет 2670 дубликатов, которые были убраны.

В. Датасет имеет следующее распределение длин предложений:

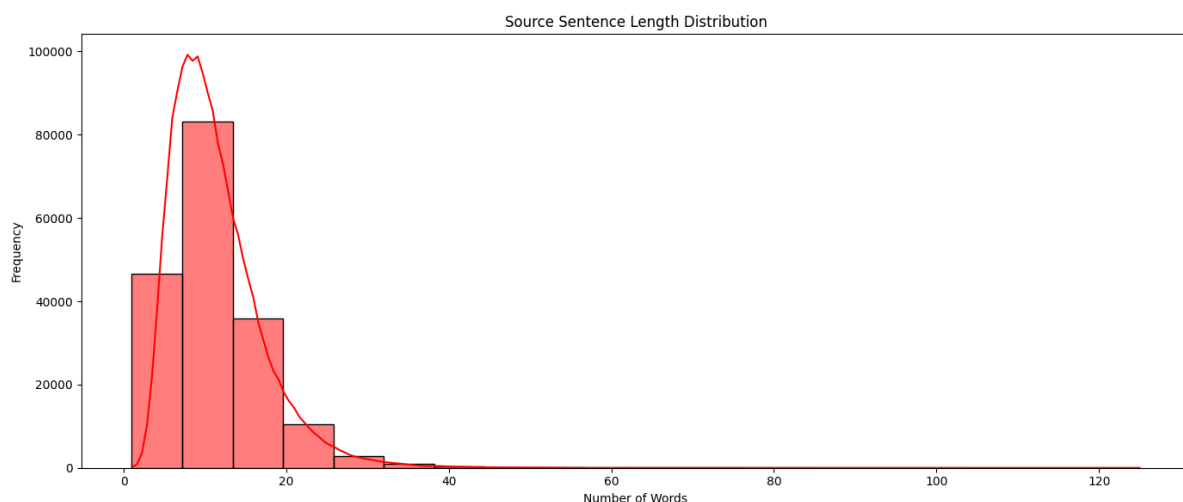


Summary of Mansi Text Lengths:

```
count      183124.000000
mean        11.417864
std          6.569508
min          1.000000
25%          7.000000
50%         10.000000
75%         14.000000
max         821.000000
Name: mns_length, dtype: float64
```

Предложение длиной в 821 слово действительно является одним предложением, так решили авторы статьи. Кол-во предложений с длиной больше 128 - 13.

Итоговое распределение для предложений, чья длина меньше 128:



C. 2063696 всего слов. Имеем следующую карту по популярности:



dict.csv

Данный файл представляет из себя объединение русско-мансийских словарей [mansi translator](#) и [fu-lab](#). Словарь содержит 11506 уникальных строк.