# BPE\_training1

#### **→ input**

hyperparams.csv

### output →

data/concatenated UniProt.txt

## **BPE\_training2**

### **→** input

data/concatenated\_UniProt.txt
hyperparams.csv

### output →

results/encoded\_sequence/BPE\_model.bpe

## generate\_tokens

### **→ input**

hyperparams.csv

results/encoded\_sequence/BPE\_model.bpe

#### output →

results/encoded\_sequence/model\_vocab.csv
results/encoded\_sequence/words.csv

# seq2vec\_skipgrams

### **→** input

hyperparams.csv

results/encoded\_sequence/words.csv

### output →

results/embedded\_sequence/seq2vec\_ids.csv
results/embedded\_sequence/skipgrams.txt

## seq2vec\_training

## **→ input**

hyperparams.csv

results/embedded\_sequence/seq2vec\_ids.csv
results/embedded\_sequence/skipgrams.txt

#### output →

results/embedded\_sequence/model.h5
results/embedded\_sequence/model\_metrics.txt
results/embedded\_sequence/seq2vec\_weights.csv

## TF IDF

#### **→** input

results/encoded\_sequence/words.csv

#### output →

results/encoded\_sequence/TF\_IDF.csv

# model\_metrics

#### **→** input

results/embedded sequence/model metrics.txt

#### output →

results/metrics/model\_acc.png
results/metrics/model loss.png

## sequence repres

## **→ input**

hyperparams.csv

results/embedded\_sequence/seq2vec\_ids.csv
results/embedded\_sequence/seq2vec\_weights.csv
results/encoded\_sequence/TF\_IDF.csv
results/encoded\_sequence/words.csv

#### output →

results/embedded\_sequence/sequence\_repres.csv

## all

### **→ input**

results/embedded\_sequence/model.h5
results/embedded\_sequence/sequence\_repres.csv
results/metrics/model\_acc.png
results/metrics/model\_loss.png