# **BPE\_training1**

#### **→** input

hyperparams.csv

#### output →

data/concatenated\_UniProt\_hp.txt

## **BPE\_training2**

#### **→ input**

data/concatenated\_UniProt\_hp.txt
hyperparams.csv

#### output →

results/encoded\_sequence/BPE\_model\_hp.bpe

## generate\_tokens

#### **→** input

hyperparams.csv

results/encoded\_sequence/BPE\_model\_hp.bpe

#### output →

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

## seq2vec\_skipgrams

#### **→** input

hyperparams.csv

results/encoded\_sequence/words\_hp.csv

#### output →

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

### TF IDF

#### **→ input**

results/encoded\_sequence/words\_hp.csv

### output →

results/encoded\_sequence/TF\_IDF\_hp.csv

## 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

#### sequence\_repres

## **→ input**

hyperparams.csv

results/embedded\_sequence/model.h5
results/embedded\_sequence/seq2vec\_ids.csv
results/encoded\_sequence/TF\_IDF\_hp.csv

results/encoded\_sequence/words\_hp.csv

#### output →

results/embedded\_sequence/sequence\_repres\_seq2vec-SIF.csv results/embedded\_sequence/sequence\_repres\_seq2vec-TFIDF.csv results/embedded\_sequence/sequence\_repres\_seq2vec.csv

### **CCR**

### **→ input**

results/embedded\_sequence/sequence\_repres\_seq2vec-SIF.csv results/embedded\_sequence/sequence\_repres\_seq2vec-TFIDF.csv results/embedded\_sequence/sequence\_repres\_seq2vec.csv

#### output →

results/embedded\_sequence/sequence\_repres\_seq2vec-SIF\_CCR.csv results/embedded\_sequence/sequence\_repres\_seq2vec-TFIDF\_CCR.csv results/embedded\_sequence/sequence\_repres\_seq2vec\_CCR.csv

#### all

### **→ input**

results/embedded\_sequence/model.h5
results/embedded\_sequence/sequence\_repres\_seq2vec.csv
results/embedded\_sequence/sequence\_repres\_seq2vec\_CCR.csv