# Table of Contents

# Sentiment analysis using a transformer

## Model training

> **Tip**
>
> Hidden below is a useful snippet of HTML to setup a `restart` button in case training gets out of hand.

Restart

## Download module from GitHub

Because we're working from a GitHub repo and not the standard Julia repository, we have to manage the installation and use of all packages rather than rely on Pluto.

```julia
1  begin
2      import Pkg
3      Pkg.activate(mktempdir())
4      Pkg.develop(url="https://github.com/rgreilly/Transformers")
5      Pkg.add(["Revise", "PlutoUI", "Flux", "DataFrames", "Printf",
6          "BSON", "JSON", "Arrow", "StatsBase", "Unicode", "Random",
7          "DataStructures", "ProgressMeter", "RemoteFiles"])
8
9      using Revise
10     using TransformersLite
11     using PlutoUI
12     using Flux
13     using Flux.CUDA
14     using Flux: DataLoader
15     using DataFrames
16     using BSON, JSON
17     using Arrow
18     using Printf
19     using StatsBase
20     using StatsBase: mean
21     using Dates
22     using Unicode
23     using Random
24     using DataStructures
25     using TransformersLite
26     using RemoteFiles
27 end;
```

```
    Activating new project at `C:\Users\alexm\AppData\Local\Temp\jl_BQHHFa
`
      Cloning git-repo `https://github.com/rgreilly/Transformers`
Path `C:\Users\alexm\.julia\dev\TransformersLite` exists and looks like the
correct repo. Using existing path.
    Resolving package versions...
    Updating `C:\Users\alexm\AppData\Local\Temp\jl_BQHHFa\Project.toml`
  [6579f8b0] + TransformersLite v0.1.0 `C:\Users\alexm\.julia\dev\Transforme
rsLite`
    Updating `C:\Users\alexm\AppData\Local\Temp\jl_BQHHFa\Manifest.toml`
  [621f4979] + AbstractFFTs v1.5.0
 ⊼ [79e6a3ab] + Adapt v3.7.2
  [dce04be8] + ArgCheck v2.3.0
  [69666777] + Arrow v2.7.0
  [31f734f8] + ArrowTypes v2.3.0
  [a9b6321e] + Atomix v0.1.0
  [ab4f0b2a] + BFloat16s v0.4.2
  [fbb218c0] + BSON v0.3.7
  [198e06fe] + BangBang v0.3.39
  [9718e550] + Baselet v0.1.1
  [c3b6d118] + BitIntegers v0.3.1
 ⊼ [fa961155] + CEnum v0.4.2
 ⊼ [052768ef] + CUDA v4.4.1
  [1af6417a] + CUDA_Runtime_Discovery v0.2.2
  [082447d4] + ChainRules v1.58.1
  [d360d2e6] + ChainRulesCore v1.19.0
  [5ba52731] + CodecLz4 v0.4.1
  [6b39b394] + CodecZstd v0.8.1
  [bbf7d656] + CommonSubexpressions v0.3.0
  [34da2185] + Compat v4.10.1
  [a33af91c] + CompositionsBase v0.1.2
  [f0e56b4a] + ConcurrentUtilities v2.3.0
  [187b0558] + ConstructionBase v1.5.4
  [6add18c4] + ContextVariablesX v0.1.3
```

In addition to the list of modules, we also need to include individual Julia files from the repo. This is done using the `RemoteFiles` module. However, this downloads them as `JSON` objects, which we need to convert back to regular `.jl` files.

```julia
1  begin
2      @RemoteFileSet FILES "Transformer utilities" begin
3
4          utilities = @RemoteFile
             "https://github.com/rgreilly/Transformers/blob/main/examples/utilities.jl"
             dir="utilities" file="utilities.jl.json"
5
6          training = @RemoteFile
             "https://github.com/rgreilly/Transformers/blob/main/examples/training.jl"
             dir="utilities" file="training.jl.json"
7      end
8
9      download(FILES) # Files downloaded in JSON format
10 end
```

convertJSON (generic function with 1 method)

```julia
1  function convertJSON(inFile, outFile)
2      body = JSON.parsefile(inFile)["payload"]["blob"]["rawLines"]
3      open(outFile, "w") do f
4        for i in body
5          println(f, i)
6        end
7      end
8  end
```

```julia
1  begin
2      convertJSON("utilities/utilities.jl.json", "utilities/utilities.jl")
3      convertJSON("utilities/training.jl.json", "utilities/training.jl")
4
5      include("utilities/utilities.jl")
6      include("utilities/training.jl")
7  end;
```

# Setup the training data

- Setup the file path to the Kaggle Amazon reviews dataset
- Assign values to various hyper-paraemeters and store them in a `dictionary`.
- Set number of training epochs

> **Tip**
>
> Here is where you can manipulate various training parameters - `pdrop`: proportion of weights to dropout (i.e., set to zero); `dim-embedding`: size of embedding; `n_epoch`: number of epochs.

```
 1  begin
 2      path = normpath(joinpath(@__DIR__, "..", "assignment6/datasets",
        "amazon_reviews_multi", "en", "1.0.0"))
 3      filename = "train.arrow"
 4      to_device = cpu # gpu or cpu
 5
 6      filepath = joinpath(path, filename)
 7
 8      df = DataFrame(Arrow.Table(filepath))
 9      display(first(df, 20))
10      println("")
11
12      hyperparameters = Dict(
13          "seed" => 314159,
14          "tokenizer" => "none", # options: none bpe affixes
15          "nlabels" => 5,
16          "pdrop" => 0.75, #CHANGE HERE FROM 0.1->0.5->0.75
17          "dim_embedding" => 128 #CHANGE HERE FROM 32->64->128
18      )
19      nlabels = hyperparameters["nlabels"]
20      n_epochs = 10
21  end;
```

```
20×9 DataFrame
 Row   Column1   review_id    product_id           reviewer_id           star  ⋯
       Int64     String       String               String                Int6  ⋯
─────────────────────────────────────────────────────────────────────────────────
   1   200000    en_0964290   product_en_0740675   reviewer_en_0342986         ⋯
   2   200001    en_0690095   product_en_0440378   reviewer_en_0133349
   3   200002    en_0311558   product_en_0399702   reviewer_en_0152034
   4   200003    en_0044972   product_en_0444063   reviewer_en_0656967
   5   200004    en_0784379   product_en_0139353   reviewer_en_0757638         ⋯
   ⋮      ⋮          ⋮             ⋮                     ⋮               ⋮   ⋱
  17   200016    en_0619473   product_en_0250211   reviewer_en_0056679
  18   200017    en_0533035   product_en_0566399   reviewer_en_0488191
  19   200018    en_0832890   product_en_0304984   reviewer_en_0667005
  20   200019    en_0550306   product_en_0387159   reviewer_en_0627216         ⋯
                                              5 columns and 11 rows omitted
```

# Tokenisers

Select a tokeniser. In this case, `none`, which just uses the various inflected word forms.

```
 1  begin
 2      if hyperparameters["tokenizer"] == "bpe"
 3          directory = joinpath("vocab", "bpe")
 4          path_rules = joinpath(directory, "amazon_reviews_train_en_rules.txt")
 5          path_vocab = joinpath(directory, "amazon_reviews_train_en_vocab.txt")
 6          tokenizer = load_bpe(path_rules, startsym="•")
 7      elseif hyperparameters["tokenizer"] == "affixes"
 8          directory = joinpath("vocab","affixes")
 9          path_vocab = joinpath(directory, "amazon_reviews_train_en_vocab.txt")
10          tokenizer = load_affix_tokenizer(path_vocab)
11      elseif hyperparameters["tokenizer"] == "none"
12          path_vocab = joinpath("vocab", "amazon_reviews_train_en.txt")
13          tokenizer = identity
14      end
15
16      vocab = load_vocab(joinpath(@__DIR__, path_vocab))
17      indexer = IndexTokenizer(vocab, "[UNK]")
18
19      display(tokenizer)
20      println("")
21      display(indexer)
22      println("")
23
24  end
```

```
identity (generic function with 1 method)                              ⑦

IndexTokenizer{String}(length(vocabulary)=6654, unksym=[UNK])
```

# Tokenise

Extract the review body and star rating from the dataframe and create embeddings. Partition data into training and validation sets.

```julia
1  begin
2      documents = df[!, :review_body]
3      labels = df[!, :stars]
4      max_length = 50
5      indices_path = joinpath(@__DIR__, "outputs", "indices_" *
       hyperparameters["tokenizer"] * ".bson")
6      @time tokens = map(d->preprocess(d, tokenizer, max_length=max_length),
7      documents)
8      @time indices = indexer(tokens)
9
10     y_labels = Int.(labels)
11     if nlabels == 1
12         y_labels[labels .≤ 2] .= 0
13         y_labels[labels .≥ 4] .= 1
14         idxs = labels .!= 3
15         y_labels = reshape(y_labels, 1, :)
16     else
17         idxs = Base.OneTo(length(labels))
18         y_labels = Flux.onehotbatch(y_labels, 1:nlabels)
19     end
20
21     X_train, y_train = indices[:, idxs], y_labels[:, idxs];
22     rng = MersenneTwister(hyperparameters["seed"])
23     train_data, val_data = split_validation(X_train, y_train; rng=rng)
24
25     println("train samples:      ", size(train_data[1]), " ", size(train_data[2]))
26     println("validation samples: ", size(val_data[1]), " ", size(val_data[2]))
27     println("")
   end
```

```
   4.023848 seconds (28.85 M allocations: 1.836 GiB, 17.55% gc time, 3.22% c  (?)
ompilation time)
 12.477690 seconds (31.81 k allocations: 81.916 MiB, 0.33% gc time, 0.53% comp
ilation time)
train samples:      (50, 184500) (5, 184500)
validation samples: (50, 20500) (5, 20500)
```

# Model definition

Assemble the model's components.

> **Tip**
>
> Here's where you might want to adjust the number and nature of the `encoder` blocks (e.g., attention heads, dropout), number of `Dense` layers and their characteristics (e.g., activation function, dimensions), the number of `dropout` layers.

**MethodError: no method matching TransformersLite.TransformerEncoderBlock(::Int64, ::Int64, ::Int64; nheads::Int64, pdrop::Float64)**

Closest candidates are:

TransformersLite.TransformerEncoderBlock(::Int64, ::Int64, ::Int64; pdrop, act) got unsupported keyword argument "nheads"

@ TransformersLite C:\Users\alexm\.julia\dev\TransformersLite\src\encoder.jl:21

1. **kwerr(::NamedTuple{(:nheads, :pdrop), Tuple{Int64, Float64}}, ::Type, ::Int64, ::Int64, ::Int64)** @ *error.jl:165*
2. **top-level scope** @ | *Local: 4*

```
 1  begin
 2      dim_embedding = hyperparameters["dim_embedding"]
 3      pdrop = hyperparameters["pdrop"]
 4      model = TransformersLite.TransformerClassifier(
 5          Embed(dim_embedding, length(indexer)),
 6          PositionEncoding(dim_embedding),
 7          Dropout(pdrop),
 8          TransformerEncoderBlock[
 9              TransformerEncoderBlock(4, dim_embedding, dim_embedding * 4;
        pdrop=pdrop), #Added ,
10
11              TransformerEncoderBlock(4, dim_embedding, dim_embedding * 4;
                 pdrop=pdrop), #Added extra transformer block normal
12
13                  TransformerEncoderBlock(4, dim_embedding, dim_embedding * 4,
                 nheads=8; pdrop=pdrop), #Added extra encoder block with attention Head
14          ],
15          Dropout(pdrop), #Adding Dropout Layers
16          Dense(dim_embedding, 1),
17          FlattenLayer(),
18          Dense(max_length, nlabels),
19          Dropout(pdrop), #Adding Dropout Layers
20          Dense(dim_embedding, 256, relu), #Adding dense layers
21          Dropout(pdrop), #Adding Dropout Layers
22          Dense(256, 128, relu), #Adding dense layers
23          )
24      display(model)
25      println("")
26      model = to_device(model)
27
28      hyperparameters["model"] = "$(typeof(model).name.wrapper)"
29      hyperparameters["trainable parameters"] = sum(length, Flux.params(model));
30
31      if nlabels == 1
32          loss(x, y) = Flux.logitbinarycrossentropy(x, y)
33          accuracy(ŷ, y) = mean((Flux.sigmoid.(ŷ) .> 0.5) .== y)
34      else
35          loss(x, y) = Flux.logitcrossentropy(x, y)
36          accuracy(ŷ, y) = mean(Flux.onecold(ŷ) .== Flux.onecold(y))
37      end
38  end;
```

# Training

- Setup the dataloaders to batch and shuffle the training and validation data.
- Print out initial accuracy and loss values for the validation data.
- Setup a sub-directory in the outputs directory, based on date and time, to store the trained model and associated hyperparameters.
- call the `train!` method and log training progress.

**UndefVarError: `model` not defined**

```julia
1  begin
2      opt_state = Flux.setup(Adam(), model)
3      batch_size = 32
4
5      train_data_loader = DataLoader(train_data |> to_device; batchsize=batch_size,
6          shuffle=true)
7      val_data_loader = DataLoader(val_data |> to_device; batchsize=batch_size,
8          shuffle=false)
9
10     val_acc = batched_metric(model, accuracy, val_data_loader)
11     val_loss = batched_metric(model, loss, val_data_loader)
12
13     @printf "val_acc=%.4f%% ; " val_acc * 100
14     @printf "val_loss=%.4f \n" val_loss
15     println("")
16
17     directory2 = normpath( joinpath(@__DIR__, "..", "outputs",
18         Dates.format(now(), "yyyymmdd_HHMM")))
19     mkpath(directory2)
20
21     hyperparameter_path = joinpath(directory2, "hyperparameters.json")
22     open(hyperparameter_path, "w") do f
23         JSON.print(f, hyperparameters)
24     end
25     println("saved hyperparameters to $(hyperparameter_path).")
26     println("")
27
28     start_time = time_ns()
29     history = train!(
30         loss, model, train_data_loader, opt_state, val_data_loader;
31             num_epochs=n_epochs)
32     end_time = time_ns() - start_time
33
34     println("done training")
35     @printf "time taken: %.2fs\n" end_time/1e9
36 end
```

**UndefVarError: `accuracy` not defined**

```julia
1  accuracy
```

# Save the model

Save model, embeddings, and training history to the `outputs` sub-directory.

**UndefVarError: `model` not defined**

```
1  begin
2      model2 = model |> cpu
3      if hasproperty(tokenizer, :cache)
4          # empty cache
5          tokenizer2 = similar(tokenizer)
6      end
7      output_path = joinpath(directory2, "model.bson")
8      history_path = joinpath(directory2, "history.json")
9      BSON.bson(
10         output_path,
11         Dict(
12             :model=> model2,
13             :tokenizer=>tokenizer,
14             :indexer=>indexer
15         )
16     )
17     println("saved model to $(output_path).")
18
19     open(history_path,"w") do f
20       JSON.print(f, history)
21     end
22     println("saved history to $(history_path).")
23
24  end
```

> **Tip**
>
> Take note of the timestamped sub-directory so that you can load the saved model and parameters for use in the evaluation notebook.