



Table of Contents

Sentiment analysis using a transformer

- Model training
 - Download module from GitHub
- Setup the training data
- Tokenisers
- Tokenise
- Model definition
- Training
- Save the model

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.

```

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  using TokenizersLite
28 end;

```

```

Activating new project at `tmp/jl_d6sXXg`
Cloning git-repo `https://github.com/rgreilly/Transformers`
Path `home/thomas/.julia/dev/TransformersLite` exists and looks like the cor
rect repo. Using existing path.
Resolving package versions...
Updating `tmp/jl_d6sXXg/Project.toml`
[6579f8b0] + TransformersLite v0.1.0 `~/.julia/dev/TransformersLite`
Updating `tmp/jl_d6sXXg/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
[a8cc5b0e] + Crayons v4.1.1
[9a962f9c] + DataAPI v1.15.0
[a93c6f00] + DataFrames v1.6.1
[864edb3b] + DataStructures v0.18.15
[e2d170a0] + DataValueInterfaces v1.0.0
[244e2a9f] + DefineSingletons v0.1.2
[8bb1440f] + DelimitedFiles v1.9.1
[163ba53b] + DiffResults v1.1.0
[b552c78f] + DiffRules v1.15.1
[ffbed154] + DocStringExtensions v0.9.3
[4e289a0a] + EnumX v1.0.4
[e2ba6199] + ExprTools v0.1.10
[cc61a311] + FLoops v0.2.1
[b9860ae5] + FLoopsBase v0.1.1
[1a297f60] + FillArrays v1.9.3
[587475ba] + Flux v0.13.17
[f6260f11] + ForwardDiff v0.10.26

```

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.

```

1 begin
2   @RemoteFileSet FILES "Transformer utilities" begin
3
4     utilities = @RemoteFile
5       "https://github.com/rgreilly/Transformers/blob/main/examples/utilities.jl"
6       dir="utilities" file="utilities.jl.json"
7
8     training = @RemoteFile
9       "https://github.com/rgreilly/Transformers/blob/main/examples/training.jl"
10      dir="utilities" file="training.jl.json"
11   end
12 end

```

convertJSON (generic function with 1 method)

```

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

```

```

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-parameters and store them in a dictionary.
- Set number of training epochs

Changes

The hyperparameter changes I made are

- Setting the tokenizer to byte pair encoding
- Setting the dropout rate to 0.2
- Setting the size of the embeddings to 64
- I added the TokenizersLite package to be able to use byte pair encoding

```

1 md"""
2 !!! changes
3   The hyperparameter changes I made are
4   - Setting the tokenizer to byte pair encoding
5   - Setting the dropout rate to 0.2
6   - Setting the size of the embeddings to 64
7   - I added the TokenizersLite package to be able to use byte pair encoding
8 """

```

```

1 begin
2   path = normpath(joinpath(@__DIR__, "..", "examples/datasets",
3     "amazon_reviews_multi", "en", "1.0.0"))
4   filename = "train.arrow"
5   to_device = cpu # gpu or cpu
6
7   filepath = joinpath(path, filename)
8
9   df = DataFrame(Arrow.Table(filepath))
10  display(first(df, 20))
11  println("")
12
13  hyperparameters = Dict(
14    "seed" => 314159,
15    "tokenizer" => "bpe", # options: none bpe affixes
16    "nlabels" => 5,
17    "pdrop" => 0.2,
18    "dim_embedding" => 64
19  )
20  nlabels = hyperparameters["nlabels"]
21  n_epochs = 10
22 end;

```

20x9 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

```

```

BytePairEncoder{String}(length(rules)=8000, length(cache)=0, unksym=[UNK], s
tartsym=., symbols=["#a", "#b", "#c", "#d", "#e", "#f", "#g", "#h", "#i",
"#j", "#k", "#l", "#m", "#n", "#o", "#p", "#q", "#r", "#s", "#t", "#u", "#v",
"#w", "#x", "#y", "#z", "a", "b", "c", "d", "e", "f", "g", "h", "i", "j", "k", "l",
m", "n", "o", "p", "q", "r", "s", "t", "u", "v", "w", "x", "y", "z", " ", "?", "!",
",", ":", ";", "<", ">", "&"])

IndexTokenizer{String}(length(vocabulary)=7455, unksym=[UNK])

```

Tokenise

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

```

1 begin
2   documents = df[:, :review_body]
3   labels = df[:, :stars]
4   max_length = 50
5   indices_path = joinpath(@__DIR__, "outputs", "indices_" *
6     hyperparameters["tokenizer"] * ".bson")
7   @time tokens = map(d->preprocess(d, tokenizer, max_length=max_length), 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("")
28 end

```

```

48.452142 seconds (690.92 M allocations: 18.200 GiB, 13.53% gc time, 1.46% ②
compilation time)
21.753756 seconds (31.81 k allocations: 81.911 MiB, 0.15% gc time, 0.42% compi
lation time)
train samples:      (50, 184500) (5, 184500)
validation samples: (50, 20500) (5, 20500)

```

Model definition

Assemble the model's components.

Change

- The dense layer is fully connected mapping the output of the transformer encoder to 64-dimensional space
- The output goes through the Flux.relu activation function.
- It appears the TransformerClassifier() constructor has rigid parameters, trying to add another dense layer (between the encoder block or flatten layer) or another dropout layer causes an error with the parameters passed, hence I could not add another dense or dropout layer
- The final dense layers input size is 64 times the length, to match the dimensional space

```

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 * 2; pdrop=pdrop)
10    ],
11     Dense(dim_embedding, 64, Flux.relu),
12     FlattenLayer(),
13     Dense(64 * max_length, nlabels)
14   )
15   display(model)
16   println("")
17   model = to_device(model)
18
19   hyperparameters["model"] = "$(typeof(model).name.wrapper)"
20   hyperparameters["trainable parameters"] = sum(length, Flux.params(model));
21
22   if nlabels == 1
23     loss(x, y) = Flux.logitbinarycrossentropy(x, y)
24     accuracy(ŷ, y) = mean((Flux.sigmoid(ŷ) .> 0.5) .== y)
25   else
26     loss(x, y) = Flux.logitcrossentropy(x, y)
27     accuracy(ŷ, y) = mean(Flux.onecold(ŷ) .== Flux.onecold(y))
28   end
29 end;

```

```

TransformerClassifier(
  Embed((64, 7455)),           # 477_120 parameters
  PositionEncoding(64),
  Dropout(0.2),
  TransformerEncoderBlock(
    MultiheadAttention(num_heads=4, head_size=16, 64=>64)(
      denseQ = Dense(64 => 64),   # 4_160 parameters
      denseK = Dense(64 => 64),   # 4_160 parameters
      denseV = Dense(64 => 64),   # 4_160 parameters
      denseO = Dense(64 => 64),   # 4_160 parameters
    ),
    Dropout(0.2),
    LayerNorm(64),              # 128 parameters
    Dense(64 => 128, relu),       # 8_320 parameters
    Dense(128 => 64),             # 8_256 parameters
    Dropout(0.2),
    LayerNorm(64),              # 128 parameters
  ),
  Dense(64 => 64, relu),         # 4_160 parameters
  FlattenLayer(),
  Dense(3200 => 5),              # 16_005 parameters
)
# Total: 21 trainable arrays, 530_757 parameters,
# plus 1 non-trainable, 64_000 parameters, summarysize 2.270 MiB.

```

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.

```

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  output_path = joinpath(directory2, "model.bson")
21  history_path = joinpath(directory2, "history.json")
22
23  hyperparameter_path = joinpath(directory2, "hyperparameters.json")
24  open(hyperparameter_path, "w") do f
25    JSON.print(f, hyperparameters)
26  end
27  println("saved hyperparameters to $(hyperparameter_path).")
28  println("")
29
30  start_time = time_ns()
31  history = train!(
32    loss, model, train_data_loader, opt_state, val_data_loader;
33    num_epochs=n_epochs)
34  end_time = time_ns() - start_time
35
36  println("done training")
37  @printf "time taken: %.2fs\n" end_time/1e9
38 end

```

```

epoch 10/10 100%|██████████████████████████████████████████| ETA: 0:00:01
mean_loss: 1.0768053779938915
batch_loss: 1.123481

epoch 10/10 100%|██████████████████████████████████████████| ETA: 0:00:00
mean_loss: 1.0769139246160395
batch_loss: 1.1102967

epoch 10/10 100%|██████████████████████████████████████████| ETA: 0:00:00
mean_loss: 1.0769191571052545
batch_loss: 1.050536

epoch 10/10 100%|██████████████████████████████████████████| ETA: 0:00:00
mean_loss: 1.0769764501242447
batch_loss: 1.3724704

epoch 10/10 100%|██████████████████████████████████████████| ETA: 0:00:00
mean_loss: 1.0770129046257266
batch_loss: 0.95596945

epoch 10/10 100%|██████████████████████████████████████████| Time: 0:03:02
mean_loss: 1.0771362375683806
batch_loss: 1.505089
train_acc=55.1946%; train_loss=1.0277; val_acc=51.8341%; val_loss=1.1053 ;
done training
time taken: 3180.33s

```

accuracy (generic function with 1 method)

1 accuracy

Save the model

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

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

```
saved model to /home/thomas/.julia/outputs/20231231_1706/model.bson.
saved history to /home/thomas/.julia/outputs/20231231_1706/history.json.
```

Tip

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