

Sentiment analysis using a transformer

Model evaluation

Table of Contents

Sentiment analysis using a transformer

Model evaluation

Load data

Tokenise

Evaluate

Test data

Examples

Probabilities

Single sample

Conclusion

Quick disclaimer, I was unable to run the training data after changing it as the training was taking o...

Restart

```

1 begin
2     import Pkg
3     Pkg.activate(mktempdir())
4     Pkg.develop(url="https://github.com/rgreilly/Transformers")
5     Pkg.add(["Revise", "PlutoUI", "Flux", "Plots", "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: DataLoader
14    using Plots
15    using DataFrames
16    using Printf
17    using BSON, JSON
18    using Arrow
19    using StatsBase
20    using Unicode
21    using Random
22    using DataStructures
23    using ProgressMeter
24    using RemoteFiles
25 end;

```

```

Activating new project at `C:\Users\alexm\AppData\Local\Temp\jl_xnjEWG`
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_xnjEWG\Project.toml`
[6579f8b0] + TransformersLite v0.1.0 `C:\Users\alexm\.julia\dev\Transforme
rsLite`
Updating `C:\Users\alexm\AppData\Local\Temp\jl_xnjEWG\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] + ContextVariables v0.1.3

```

```

1 begin
2   @RemoteFileSet FILES "Transformer utilities" begin
3     reporting = @RemoteFile
4       "https://github.com/rgreilly/Transformers/blob/main/examples/reporting.jl"
5     dir="utilities" file="reporting.jl.json"
6     utilities = @RemoteFile
7       "https://github.com/rgreilly/Transformers/blob/main/examples/utilities.jl"
8     dir="utilities" file="utilities.jl.json"
9     training = @RemoteFile
10      "https://github.com/rgreilly/Transformers/blob/main/examples/training.jl"
11    dir="utilities" file="training.jl.json"
12  end
13
14  download(FILES) # Files downloaded in JSON format
15 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/reporting.jl.json", "utilities/reporting.jl")
3   convertJSON("utilities/utilities.jl.json", "utilities/utilities.jl")
4   convertJSON("utilities/training.jl.json", "utilities/training.jl")
5
6   include("utilities/reporting.jl")
7   include("utilities/utilities.jl")
8   include("utilities/training.jl")
9 end;

```

Multi-class classification: stars from 1 to 5

Load data

The original CSV data format has been converted to arrow format for faster loading.

```

1 begin
2   path = "datasets/amazon_reviews_multi/en/1.0.0/"
3   file_train = "train.arrow"
4   file_test = "test.arrow"
5   nlabels = 5
6 end;

```

Load training and test data into dataframes and about the size of both.

```
(205000, 5000)
```

```
1 begin
2   filepath = joinpath(path, file_train)
3   df = DataFrame(Arrow.Table(filepath))
4
5   filepath = joinpath(path, file_test)
6   df_test = DataFrame(Arrow.Table(filepath))
7
8   (nrow(df), nrow(df_test))
9 end
```

Extract just the review text and the star rating

```
1 begin
2   documents = df[:, "review_body"]
3   labels = df[:, "stars"]
4
5   println("training samples: ", size(documents), " ", size(labels))
6 end
```

```
training samples: (205000,) (205000,)
```



Do the same for the test data

```
1 begin
2   documents_test = df_test[:, "review_body"]
3   labels_test = df_test[:, "stars"];
4
5   println("test samples: ", size(documents_test), " ", size(labels_test))
6 end
```

```
test samples: (5000,) (5000,)
```



Load the already trained and saved model. Note that models and associated details are stored in the outputs directory under a sub-directory name generated from the time it was saved in the format: `yyyymmdd_hhmm`.

```

1 begin
2     directory = "../outputs/20231219_1608/" #CHANGE HERE
3     saved_objects = BSON.load(joinpath(directory, "model.bson"))
4     tokenizer = saved_objects[:tokenizer]
5     @show tokenizer
6     indexer = saved_objects[:indexer]
7     @show indexer
8     model = saved_objects[:model]
9     display(model)
10 end;

```

```

tokenizer = identity
indexer = IndexTokenizer{String}(length(vocabulary)=6654, unksym=[UNK])
TransformerClassifier(
  Embed((32, 6654)),           # 212_928 parameters
  PositionEncoding(32),
  Dropout(0.1),
  TransformerEncoderBlock(
    MultiheadAttention(num_heads=4, head_size=8, 32=>32)(
      denseQ = Dense(32 => 32),      # 1_056 parameters
      denseK = Dense(32 => 32),      # 1_056 parameters
      denseV = Dense(32 => 32),      # 1_056 parameters
      denseO = Dense(32 => 32),      # 1_056 parameters
    ),
    Dropout(0.1),
    LayerNorm(32),                # 64 parameters
    Dense(32 => 128, relu),         # 4_224 parameters
    Dense(128 => 32),               # 4_128 parameters
    Dropout(0.1),
    LayerNorm(32),                # 64 parameters
  ),
  Dense(32 => 1),                  # 33 parameters
  FlattenLayer(),
  Dense(50 => 5),                  # 255 parameters
)
# Total: 21 trainable arrays, 225_920 parameters,
# plus 1 non-trainable, 32_000 parameters, summarysize 1009.188 KiB.

```

Tokenise

Tokenise the training and test data

```

1 begin
2   max_length = size(model.classifier.weight, 2)
3   @time tokens = map(d->preprocess(d, tokenizer, max_length=max_length),
4   documents) #takes about 30 seconds for all documents
5   @time indices = indexer(tokens) #takes about 12 seconds for all documents
6
7   y_train = copy(labels)
8   idxs = Base.OneTo(length(labels))
9   X_train, y_train = indices[:, idxs], y_train[idxs];
10  y_train = Flux.onehotbatch(y_train, 1:5) # multi-class
11  train_data, val_data = split_validation(X_train, y_train;
12    rng=MersenneTwister(2718))
13
14  println("train samples:      ", size(train_data[1]), " ", size(train_data[2]))
15  println("validation samples: ", size(val_data[1]), " ", size(val_data[2]))
16 end

```

```

4.125783 seconds (28.48 M allocations: 1.786 GiB, 16.44% gc time)
12.979427 seconds (4 allocations: 79.765 MiB)
train samples:      (50, 184500) (5, 184500)
validation samples: (50, 20500) (5, 20500)

```



```

1 begin
2   y_test = copy(labels_test)
3   y_test = Flux.onehotbatch(y_test, 1:5);
4
5   @time tokens_test = map(d->preprocess(d, tokenizer, max_length=max_length),
6   documents_test)
7   @time indices_test = indexer(tokens_test)
8
9   X_test = indices_test
10
11  println("test indices: ", size(indices_test))
12  println("test samples: ", size(X_test), " ", size(y_test))
13 end

```

```

0.100063 seconds (718.85 k allocations: 46.182 MiB, 26.17% compilation time)
0.307723 seconds (19 allocations: 1.946 MiB, 1.12% compilation time)
test indices: (50, 5000)
test samples: (50, 5000) (5, 5000)

```



Create the training and validation data loaders

321-element DataLoader(::Tuple{Matrix{Int64}, OneHotArrays.OneHotMatrix{UInt32, Vector{UInt32}}, Vector{UInt32}} with first element:
(50×64 Matrix{Int64}, 5×64 OneHotMatrix{::Vector{UInt32}} with eltype Bool,)

```

1 begin
2   train_data_loader = DataLoader(train_data; batchsize=64, shuffle=false);
3   val_data_loader   = DataLoader(val_data; batchsize=64, shuffle=false);
4 end

```

Evaluate

accuracy (generic function with 1 method)

```
1 begin
2     loss(x, y) = Flux.logitcrossentropy(model(x), y)
3     loss(x::Tuple) = loss(x[1], x[2])
4     accuracy(ŷ, y) = mean(Flux.onecold(ŷ) .== Flux.onecold(y))
5 end
```

0.5520271002710027

```
1 @time batched_metric(model, accuracy, train_data_loader)
```

39.119026 seconds (6.73 M allocations: 56.692 GiB, 8.86% gc time, 10.61% compilation time) ?

0.5595121951219513

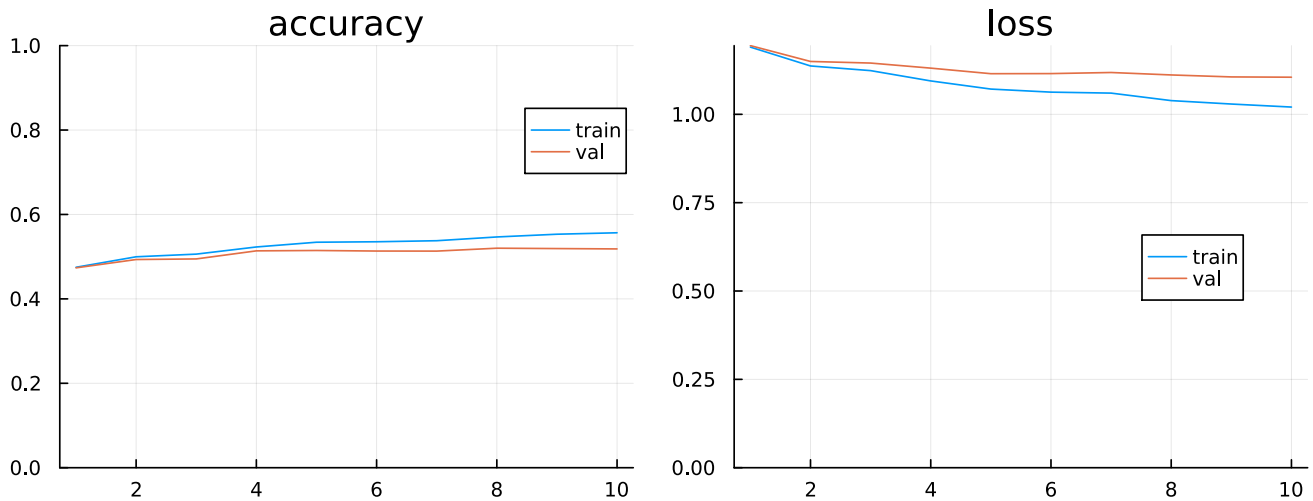
```
1 @time batched_metric(model, accuracy, val_data_loader)
```

3.695079 seconds (100.34 k allocations: 6.257 GiB, 9.76% gc time) ?

history =

Dict{"train_loss" => [1.19026, 1.13687, 1.12365, 1.0945, 1.07144, 1.06278, 1.06003, 1.03

```
1 history = open(joinpath(directory, "history.json"), "r") do f
2     JSON.parse(read(f, String))
3 end
```

```

1 begin
2   epochs = 1:length(history["train_acc"])
3   p1 = plot(epochs, history["train_acc"], label="train")
4   plot!(p1, epochs, history["val_acc"], label="val")
5   plot!(p1, ylims=[0, 1], title="accuracy", legend=(0.9, 0.8))
6
7   p2 = plot(epochs, history["train_loss"], label="train")
8   plot!(p2, epochs, history["val_loss"], label="val")
9   plot!(p2, title="loss", ylims=[0, Inf], legend=(0.8, 0.5))
10
11  p3 = plot(p1, p2, layout=grid(1, 2), size=(800, 300))
12  savefig(p3, joinpath(directory, "history.png"))
13  p3
14 end

```

Test data

0.5242

```

1 begin
2   logits = model(X_test)
3   accuracy(logits, y_test)
4 end

```

[1, 1, 1, 2, 2, 1, 4, 1, 2, 1, 3, 1, 1, 2, 1, 1, 2, 1, 2, 2, more ,3, 5, 1, 5, 2, 3, 5, 5

```

1 begin
2   probs = softmax(logits, dims=1)
3   y_pred = Flux.onecold(probs);
4 end

```

```

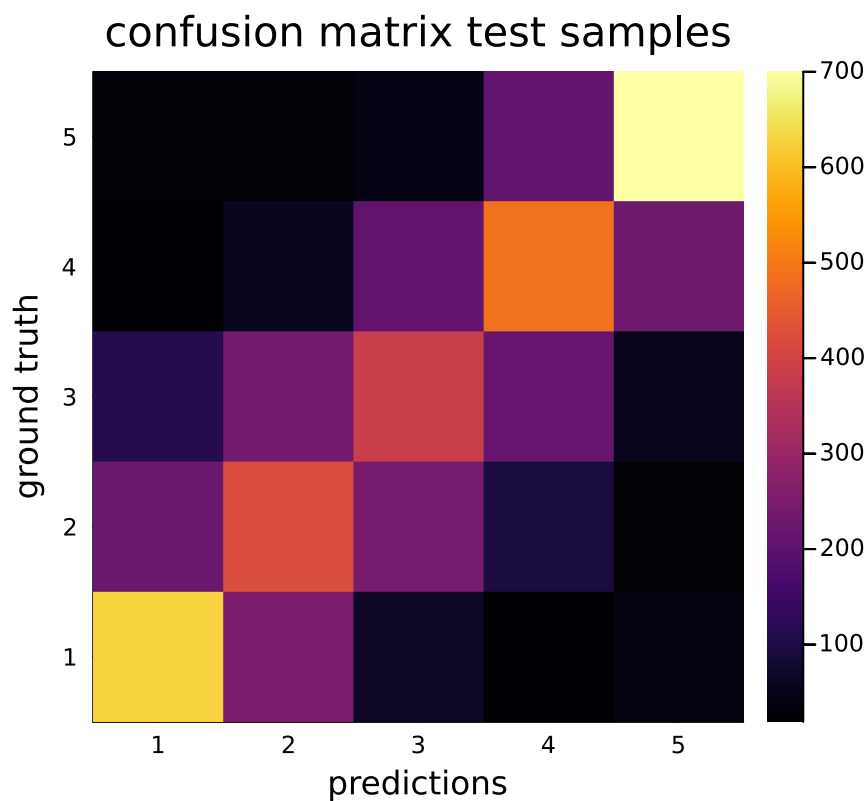
cm = 5×5 Matrix{Int64}:
 627  252   64   21   36
 221  421  242   92   24
 110  239  385  215   51
  19   57  206  488  230
  25   26   41  208  700

```

```

1 cm = confusion_matrix(vec(y_pred), Flux.onecold(y_test), 1:nlabels)

```



```

1 begin
2   p4 = heatmap(1:5, 1:5, cm, xlabel="predictions", ylabel="ground truth", xlims=
      (0.5, nlabels+0.5), aspectratio=1,
3     title="confusion matrix test samples", xticks=(1:5)) #, ["negative", "mix",
      "positive"]))
4   savefig(p4, joinpath(directory, "confusion_matrix.png"))
5   p4
6 end

```

```

1 classification_report(cm, 1:nlabels)

```

	precision	recall	f1-score	support
1	0.63	0.63	0.63	1000
2	0.42	0.42	0.42	1000
3	0.41	0.39	0.40	1000
4	0.48	0.49	0.48	1000
5	0.67	0.70	0.69	1000
weighted avg	0.52	0.52	0.52	5000



Examples

```
1 begin
2     println("star y  $\hat{y}$  prob")
3     for star in nlabels:-1:1
4         pos_max = argmax(probs[star, :])
5         @printf("    %1d %d %d %.4f\n %s\n\n",
6                 star, labels_test[pos_max], y_pred[pos_max], probs[star, pos_max],
7                 documents_test[pos_max])
8     end
9 end
```

```
star y  $\hat{y}$  prob
5 5 5 0.9818
Best USB cable I ever owned. Fit and finish are top notch. Highly recommended.

4 4 4 0.9047
I like the case, it works well. I wish it came in purple. The only complaint I have is that it doesn't charge well on the cordless charger.

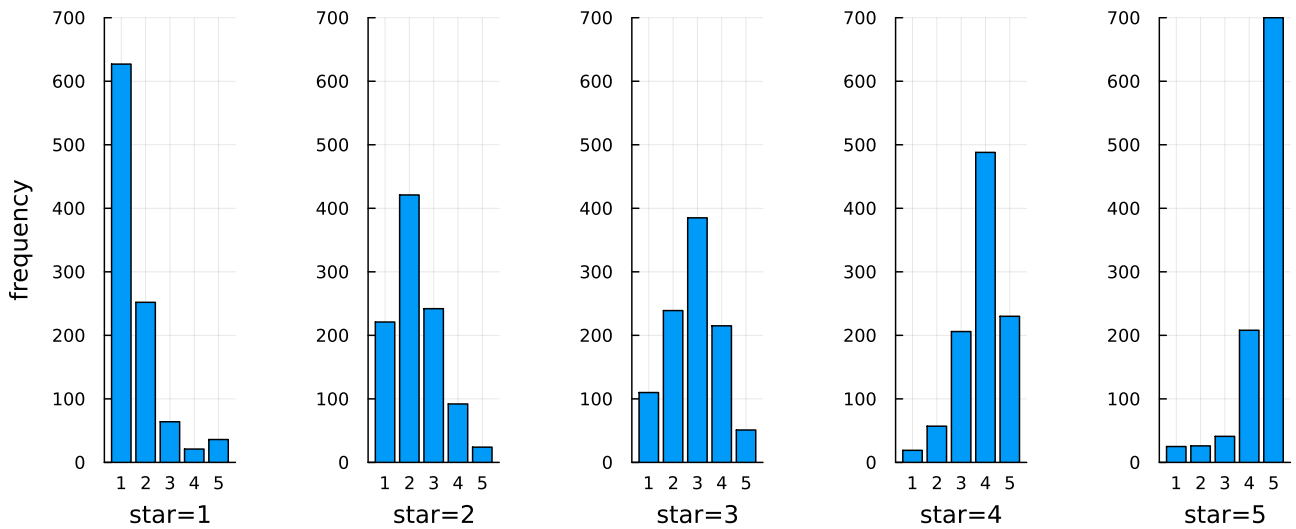
3 3 3 0.6972
Works ok kinda cheap

2 4 2 0.6562
The cover and shams look as pictured. Very soft and comfortable. Its very delicate, I tore a hole by the zipper very easily. I did not return cover because it was super cheap.

1 1 1 0.9905
I didn't even want to give this 1 star! This is a piece of JUNK! Terrible made! I am having. Such a HARD TIME RETURNING THIS PRODUCT! Don't WASTE YOUR MONEY. It looks like these toys were made for a doll house. Wrongly advertised!
```

Probabilities

predicted class per ground truth class



```

1 begin
2   canvases1 = []
3   label_names = 1:5
4   for gt_star in 1:5
5     idxs = labels_test .== gt_star
6     value_counts = [sum((y_pred[idxs]) .== l) for l in 1:nlabels]
7     p = bar(value_counts, xlabel="star=$gt_star", legend=:none, xticks=
8       (1:nlabels, 1:5))#["neg", "mix", "pos"])
9     push!(canvases1, p)
10  end
11  plot!(canvases1[1], ylabel="frequency")
12  p5 = plot(canvases1..., layout=(1, 5), link=:y, size=(900, 400),
13    plot_title="predicted class per ground truth class",
14    margin=5Plots.mm)
15  savefig(p5, joinpath(directory, "prediction_star.png"))
16  p5
17 end

```

Single sample

```

1 begin
2     idx = 4600
3
4     d = documents_test[idx]
5     println(labels_test[idx])
6     println(d)
7     println("")
8
9     tokens2 = preprocess(d, tokenizer, max_length=50)
10    println(join(tokens2, "|"))
11    println("")
12
13    x = indexer(tokens2)
14    x = vcat(x, ones(Int, 50 - length(x)))
15    println(join(x, "|"))
16 end

```

[illegible]

```
5x1 Matrix{Float32}:
 0.0014079323
 0.00606461
 0.028641572
 0.2757954
 0.6880905
```

```
1 softmax(model(x))
```

Conclusion

Quick disclaimer, I was unable to run the training data after changing it as the training was taking over double the time it took to run it initially. Later on it wouldn't even get past the model creation stage without any changes that could possibly break it. The Train_multi.jl/pdf is the same as the file that was used to generate the output for this evaluation

The reason why I was able to slightly improve the model's performance was because I added some transformerEncoderBlocks, one of which made use of Attention heads to improve the model's capacity to capture complex patterns and to focus on different elements of the input. I also added some dense layers. I added dropout layers to prevent overfitting.

Apart from this I experimented with different values for the Dropout Rate to prevent overfitting as well as the size of the embedding to find a sweet spot.

```
1 md"""
2 # Conclusion
3
4 ## Quick disclaimer, I was unable to run the training data after changing it as the
   training was taking over double the time it took to run it initially. Later on it
   wouldn't even get past the model creation stage without any changes that could
   possibly break it. The Train_multi.jl/pdf is the same as the file that was used to
   generate the output for this evaluation
5
6 ###
7 The reason why I was able to slightly improve the model's performance was because I
   added some transformerEncoderBlocks, one of which made use of Attention heads to
   improve the model's capacity to capture complex patterns and to focus on different
   elements of the input.
8 I also added some dense layers.
9 I added dropout layers to prevent overfitting.
10
11 Apart from this I experimented with different values for the Dropout Rate to
   prevent overfitting as well as the size of the embedding to find a sweet spot.
12 """
```

