

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

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
10    using Revise
11    using TransformersLite
12    using PlutoUI
13    using Flux
14    using Flux: DataLoader
15    using Plots
16    using DataFrames
17    using Printf
18    using BSON, JSON
19    using Arrow
20    using StatsBase
21    using Unicode
22    using Random
23    using DataStructures
24    using ProgressMeter
25    using RemoteFiles
26 end;

```

```

Activating new project at 'C:\Users\u210148\AppData\Local\Temp\jl_9Sqjzt' ?
Cloning git-repo 'https://github.com/rgreilly/Transformers'
Path 'C:\Users\u210148\.julia\dev\TransformersLite' exists and looks like the
correct repo. Using existing path.
Resolving package versions...
Updating 'C:\Users\u210148\AppData\Local\Temp\jl_9Sqjzt\Project.toml'
[6579f8b0] + TransformersLite v0.1.0 'C:\Users\u210148\.julia\dev\Transforme
rsLite'
Updating 'C:\Users\u210148\AppData\Local\Temp\jl_9Sqjzt\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 v4.11.0
[a33af91c] + CompositionsBase v0.1.2

```

```

1 begin
2   @RemoteFileSet FILES "Transformer utilities" begin
3     reporting = @RemoteFile "https://github.com/rgreilly/Transformers/blob/main
4       /examples/reporting.jl" dir="utilities" file="reporting.jl.json"
5     utilities = @RemoteFile "https://github.com/rgreilly/Transformers/blob/main
6       /examples/utilities.jl" dir="utilities" file="utilities.jl.json"
7     training = @RemoteFile "https://github.com/rgreilly/Transformers/blob/main
8       /examples/training.jl" dir="utilities" file="training.jl.json"
9   end
10
11   download(FILES) # Files downloaded in JSON format
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/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/20240111_1629/"
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((20, 6654)),          # 133_080 parameters
    PositionEncoding(20),
    Dropout(0.3),
    TransformerEncoderBlock(
        MultiheadAttention(num_heads=4, head_size=5, 20=>20)(
            denseQ = Dense(20 => 20),      # 420 parameters
            denseK = Dense(20 => 20),      # 420 parameters
            denseV = Dense(20 => 20),      # 420 parameters
            denseO = Dense(20 => 20),      # 420 parameters
        ),
        Dropout(0.3),
        LayerNorm(20),            # 40 parameters
        Dense(20 => 80, relu),     # 1_680 parameters
        Dense(80 => 20),          # 1_620 parameters
        Dropout(0.3),
        LayerNorm(20),            # 40 parameters
    ),
    Dense(20 => 1),               # 21 parameters
    FlattenLayer(),
    Dense(50 => 5),               # 255 parameters
)
# Total: 21 trainable arrays, 138_416 parameters,
# plus 1 non-trainable, 20_000 parameters, summarysize 620.500 KiB.

```

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

```

```

4.467094 seconds (28.48 M allocations: 1.786 GiB, 15.92% gc time)
14.673765 seconds (4 allocations: 79.765 MiB)
train samples:      (50, 184500) (5, 184500)
validation samples: (50, 20500) (5, 20500)

```

Tokenise

```

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

```

```

0.114423 seconds (718.85 k allocations: 46.182 MiB, 26.68% compilation time)
0.357098 seconds (19 allocations: 1.946 MiB, 1.00% 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{UI
  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
```

0.5246124661246613

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

```
37.649349 seconds (6.50 M allocations: 43.466 GiB, 7.74% gc time, 10.66% c
ompilation time) ?
```

0.5284390243902439

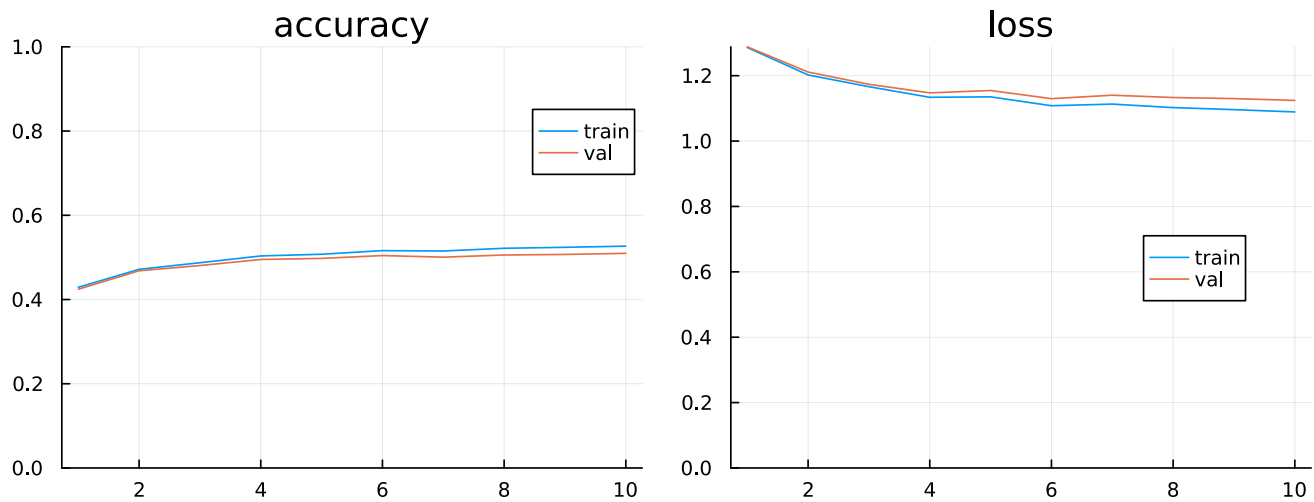
```
1 @time batched_metric(model, accuracy, val_data_loader)
```

```
3.564160 seconds (80.86 k allocations: 4.788 GiB, 8.48% gc time) ?
```

history =

```
Dict{"train_loss" => [1.2862, 1.20225, 1.16633, 1.13392, 1.13527, 1.10809, 1.11309, 1.10:
```

```
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

0.5104

```

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

```

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

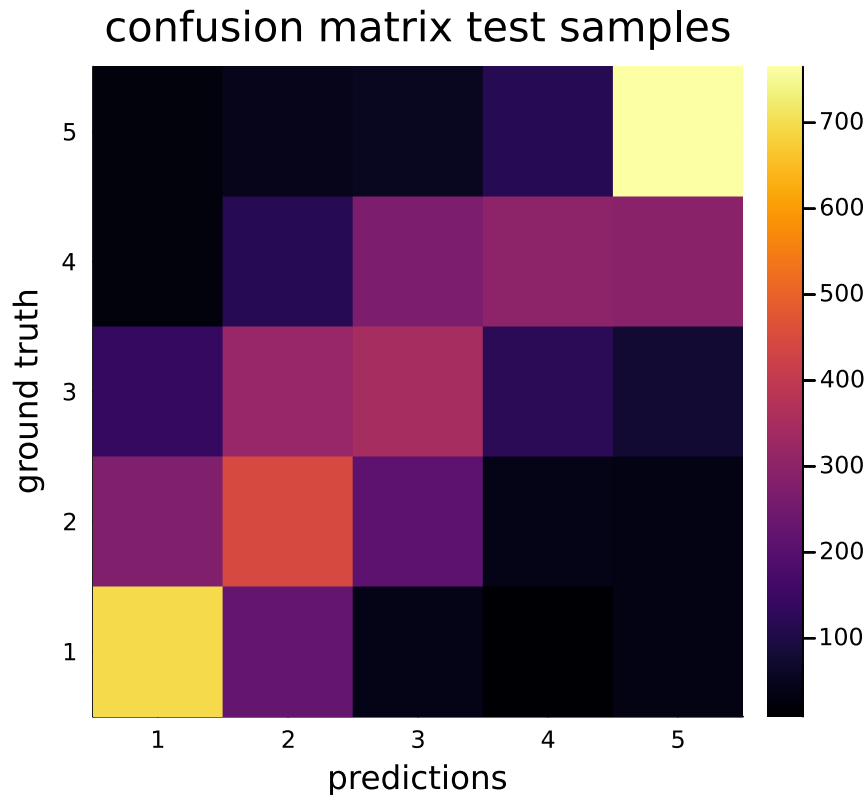
```

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

```

```
cm = 5×5 Matrix{Int64}:
 697  222   38    9   34
 276  443  208   39   34
 140  318  347  121   74
  24  115  267  300  294
  26   43   51  115  765
```

```
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",
3     xlims=(0.5, nlabels+0.5), aspectratio=1,
4     title="confusion matrix test samples", xticks=(1:5)) #, ["negative", "mix",
5     "positive"]))
6   savefig(p4, joinpath(directory, "confusion_matrix.png"))
7   p4
8 end
```

```
1 classification_report(cm, 1:nlabels)
```

	precision	recall	f1-score	support
1	0.60	0.70	0.64	1000
2	0.39	0.44	0.41	1000
3	0.38	0.35	0.36	1000
4	0.51	0.30	0.38	1000
5	0.64	0.77	0.70	1000
weighted avg	0.50	0.51	0.50	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 0.9788
The item ordered came exactly as advertised. I highly recommend this vendor and would order from them again.

4 4 4 0.8214
Love the shirt. I wish the fabric of the shirt was softer.... a little rough. Overall good shirt

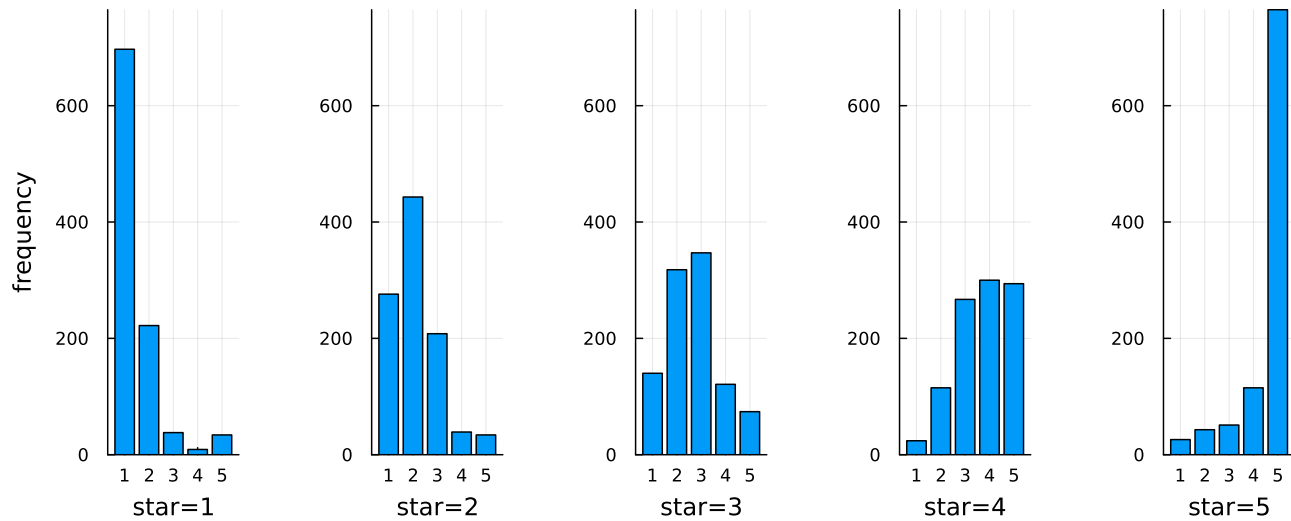
3 2 3 0.6738
I liked it okay while it lasted... I really loved this case when I bought it. However it cracked & leaked out. In addition, the lip needs to be a little bit bigger around the edge of the phone. I've gone through MANY, MANY Bodyguardz screen protectors with this case.

2 3 2 0.5584
It was great while it lasted, but it started falling apart after a couple of months. The brown inner lining ripped from the inner seam of the leather and it's not really fixable. Not great for long-term every day use.

1 1 1 0.9749
Scam, never came in!
```

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

[illegible]

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

11/01/2024, 17:16