

Entrenamiento de una red neuronal para la identificación de los parámetros de la señal que caracterizan la señal l_{cm} y σ a partir de la representación del autoencoder

- Ahora que cambió la representación en baja dimensionalidad porque cambiamos la forma en que se hacía tenemos que reentrenar una red neuronal para poder identificar los parámetros de las distribuciones de tamaños de donde provienen las señales. Para esto se hizo una exploración de hiperparámetros teniendo en cuenta distintas arquitecturas y funciones de activación. Al final se llegó a una red bastante similar a la anterior con respecto a la arquitectura pero esta vez con la función de activación tahn.

Multiple expressions in one cell.

How would you like to fix it?

- Split this cell into 14 cells, or
- Wrap all code in a `begin ... end` block.

```
1 # Importamos las librerías necesarias
2 using Flux
3 using Statistics
4 using Flux: train!
5 using Plots
6 using DataFrames
7 using CSV
8 using StatsPlots
9 using LaTeXStrings
10 using LinearAlgebra
11 using PlotlyJS
12 using CUDA
13 using Random
14 using Measures
15 # Importamos los parámetros del experimento
16 include("../1-GeneracionDatos/Parameters.jl");
```

Multiple expressions in one cell.

How would you like to fix it?

- Split this cell into 13 cells, or
- Wrap all code in a `begin ... end` block.

```
1 # Distribucion de probabilidad log-normal que se puede agregar a la función costo
  de la red neuronal, lleva mucho tiempo de entrenamiento
2
3 function Pln(lcm, σ)
4     return [(exp(-(log(lc) - log(lcm))^2 / (2σ^2))) / (lc * σ * sqrt(2π)) for lc in
  lcs]
5 end
6
7 #####
  #####
8 # Funciones de pre procesamiento para escalar los datos y estandarizarlos
9
10 # Normalización Max-Min
11 function MaxMin(data)
12     min_vals = minimum(data, dims=1)
13     max_vals = maximum(data, dims=1)
14     scaled_data = (data .- min_vals) ./ (max_vals .- min_vals)
15     return scaled_data
16
17 end
18
19 # Estandarización Z
20 function Standarize(data)
21     mean_vals = mean(data, dims=1)
22     std_devs = std(data, dims=1)
23     standardized_data = (data .- mean_vals) ./ std_devs
24     return standardized_data
25 end
26
27 #####
  #####
28 # Metricas de validacion de la red neuronal, solo utilice RMAE
29
30 # Root Mean Squared Error
31 function RMSE(predicted, real)
32     return sqrt(Flux.mse(predicted, real))
33 end
34
35 # Mean Absolute Error
36 function MAE(predicted, real)
37     return sum(abs.(predicted .- real)) / length(predicted)
38 end
39
40 # R2 score
41 function R2_score(predicted, real)
42     return 1 - sum((predicted .- real).^2) / sum((real .- mean(real)).^2)
43 end
44
```

```

45 # Realitive Root Mean Squared Error
46 function RRMSE(predicted, real)
47     return sqrt(mean((predicted .- real).^2)) / mean(real)
48 end
49
50 # Relative Mean Absolute Error
51 function RMAE(predicted, real)
52     return mean(abs.(predicted .- real)) / mean(real)
53 end
54
55 # Mean Absolute Percentaje Error
56 function MAPE(predicted, real)
57     return mean(abs.((predicted .- real) ./ real))
58 end
59
60 #####
61 #####
62
63 # Regularizaciones L1 y L2
64 pen_l2(x::AbstractArray) = Float32.(sum(abs2, x) / 2)
65 pen_l1(x::AbstractArray) = Float32.(sum(abs, x) / 2)
66 #####
67 #####
68
69 # Funciones de pre procesamiento para escalar los datos y estandarizarlos
70
71 # Normalización Max-Min
72 function MaxMin(data)
73     min_vals = minimum(data, dims=1)
74     max_vals = maximum(data, dims=1)
75     scaled_data = (data .- min_vals) ./ (max_vals .- min_vals)
76     return scaled_data
77 end
78
79 # Estandarización Z
80 function Standarize(data)
81     mean_vals = mean(data, dims=1)
82     print(mean_vals)
83     std_devs = std(data, dims=1)
84     print(std_devs)
85     standardized_data = (data .- mean_vals) ./ std_devs
86     return standardized_data
87 end:

```

Multiple expressions in one cell.

How would you like to fix it?

- Split this cell into 2 cells, or
- Wrap all code in a `begin ... end` block.

```
1 # Leemos los datos de las señales en PCA
2 path_read = "../1-GeneracionDatos/Data/"
3 df_datasignals =
  CSV.read("C:/Users/Propietario/OneDrive/Escritorio/ib/Tesis_V1/MLonNMR/2-2-
  Autoencoders/Models/ReducedDataNLData.csv", DataFrame)
```

K-folds

Usamos la técnica de validación cruzada K-folds para prevenir el sobreajuste de la red neuronal. Este método divide el conjunto de datos en K subconjuntos, y el modelo se entrena K veces, utilizando en cada iteración un subconjunto diferente como conjunto de validación y el resto como conjunto de entrenamiento. Al final, los resultados obtenidos en cada iteración se promedian, esto nos permite obtener una mejor estimación de cómo funcionará el modelo con datos no vistos previamente.

El proceso de separación de los datos es el siguiente:

Los datos se dividen en K subconjuntos de igual tamaño. Cada subconjunto requiere de una representación adecuada de los datos, eligiendo señales de validación basadas en diferentes valores de σ y l_{cm} de la siguiente manera: Se comienza con $l_{cm} = 0.5$, pero en cada fold se utiliza un valor distinto de σ . Posteriormente, se seleccionan valores equiespaciados de l_{cm} y σ de modo que el 20% de los datos se destinen a la validación. Así, al comenzar cada fold con un valor diferente de σ , se obtienen K conjuntos disjuntos para validar. Después de entrenar el modelo K veces, cambiando el conjunto de validación en cada iteración, se calcula el rendimiento promedio de la red. Para esto, se obtienen la media y la desviación estándar de los errores de validación RMSE de los K folds.

Multiple expressions in one cell.

How would you like to fix it?

- Split this cell into 11 cells, or
- Wrap all code in a `begin ... end` block.

```
1 # Utilizamos la técnica k-fold de validación cruzada para prevenir el overfitting  
2 # Definimos el número de folds  
3 folds = 5  
4 percent_valid = 0.2  
5 step_valid = Int(1 / percent_valid)  
6 num_datos = Int(size(df_datasignals, 1))  
7  
8 # Guardamos los datos de validacion de cada NN en cada fold  
9 out_of_sample_data = []  
10 out_of_sample_pred = []  
11  
12 # Guardamos la metrica de validación de cada NN en cada fold  
13 scores_RMSE = []  
14  
15 # Definimos desde donde empezamos los datos de testing  
16 idx_startTest = 6  
17  
18 # Primero sacamos los datos de testing de los datos de señales, estos seran un 3er  
19 conjunto de datos que no se usara para entrenar ni validar la red  
19 df_datasignals_out = df_datasignals[idx_startTest:step_valid:num_datos,:]  
20 # Datos de Testing  
21 df_datasignals_minus_out = df_datasignals[setdiff(1:num_datos,  
21 idx_startTest:step_valid:num_datos),:]  
22 # Nuevo numero de datos que tenemos para entrenamiento + validacion  
23 num_datos_new = Int(size(df_datasignals_minus_out, 1))
```

UnDefVarError: `folds` not defined

Stack trace

Here is what happened, the most recent locations are first:

1. (This cell: line 3

```
2
3 for k in 1:folds
4     # Usamos 5 conjuntos disjuntos de datos de validación del 5% total de
    los datos para cada fold
```

```
1 # Método de K-Folds
2
3 for k in 1:folds
4     # Usamos 5 conjuntos disjuntos de datos de validación del 5% total de los datos
    para cada fold
5
6     datasignals_valid = Float32.(Matrix(df_datasignals_minus_out[k^2 +
    10:step_valid:num_datos_new,1:3]))'
7     datasignals = Float32.(Matrix(df_datasignals_minus_out[setdiff(1:num_datos_new,
    k^2 + 10:step_valid:num_datos_new),1:3]))'
8
9     σ_valid = df_datasignals_minus_out[k^2 + 10:step_valid:num_datos_new,4]
10    lcm_valid = df_datasignals_minus_out[k^2 + 10:step_valid:num_datos_new,5]
11
12    σ_col = df_datasignals_minus_out[setdiff(1:num_datos_new, k^2 +
    10:step_valid:num_datos_new),4]
13    lcm_col = df_datasignals_minus_out[setdiff(1:num_datos_new, k^2 +
    10:step_valid:num_datos_new),5]
14
15    dataparams = hcat(lcm_col, σ_col)'
16    dataparams_valid = hcat(lcm_valid, σ_valid)'
17
18    # Uno de los modelos quedio buenos resultados en la exploración anterior
19    # 17,"[3, 32, 64, 32, 16,
    2]",tanh,ADAM,0.0,None,0,0.017810457567638147,0.01427131790632415
20
21    id = 17
22
23    # Definimos la red neuronal
24    model = Chain(
25        Dense(3, 32),
26        Dense(32, 64, tanh_fast),
27        Dense(64, 32, tanh_fast),
28        Dense(32, 16, tanh_fast),
29        Dense(16, 2, softplus),
30    )
31
32    # Cargamos los parámetros del modelo
33    path = "C:/Users/Propietario/OneDrive/Escritorio/ib/Tesis_V1/MLonNMR/2-
    NNIdentificaHahn/G_8.73e-7_TE_1_AE/Params/017_params.csv"
```

```

34
35     p, re = Flux.destructure(model)
36
37     theta = CSV.read(path, DataFrame)
38     p = @views Float32.(theta[:,1])
39
40     # Función de loss
41     function loss(x,y)
42         return Flux.mse(re(p)(x), y)
43     end
44
45     # Definimos el metodo de aprendizaje y la tasa de aprendizaje
46      $\eta$  = 1e-4
47     opt = ADAM( $\eta$ )
48
49     # Definimos el número de épocas
50     epochs = 5000
51
52     # Definimos el tamaño del batch
53     batch_size = 100
54
55     # Usamos dataloader para cargar los datos
56     data = Flux.DataLoader((datasignals, dataparams), batchsize = batch_size,
57 shuffle = true)
58     data_valid = Flux.DataLoader((datasignals_valid, dataparams_valid), batchsize =
59 batch_size, shuffle = true)
60
61     # Definimos una funcion de callback para ver el progreso del entrenamiento
62     global iter = 0
63     cb = function()
64         global iter += 1
65         if iter % length(data) == 0
66             epoch = iter ÷ length(data)
67             if epoch % 500 == 0
68                 actual_loss = loss(data.data[1], data.data[2])
69                 actual_valid_loss = loss(data_valid.data[1], data_valid.data[2])
70                 println("Epoch $epoch || Loss = $actual_loss || Valid Loss =
71 $actual_valid_loss")
72             end
73         end
74     end;
75
76     # Entrenamos la red neuronal con el loss mse variando la tasa de aprendizaje
77 cada 500 épocas
78     for epoch in 1:epochs
79         Flux.train!(loss, Flux.params(p), data, opt, cb=cb)
80         if epoch % 500 == 0
81              $\eta$  =  $\eta$  * 0.2
82             opt = ADAM( $\eta$ )
83         end
84     end
85
86     # Predicción de la red en la validacion
87     predictions_valid = re(p)(datasignals_valid)
88
89     # Métricas de validación de la red
90     RMSE_valid = RMSE(predictions_valid, dataparams_valid)
91

```

```

88     push!(scores_RMSE, RMSE_valid)
89
90     # Guardamos los datos de validación y las predicciones de la red
91     push!(out_of_sample_data, dataparams_valid)
92     push!(out_of_sample_pred, predictions_valid)
93
94     println("Fold $k terminado con score de validación RMSE = $RMSE_valid")
95

```

El promedio de la metrica RMSE en los conjuntos 5 folds de validación es 0.014 y el desvio estandar es 0.0011

- Al haber testado en todos estos datos y seguir obteniendo valores similares podemos asegurar que el modelo no está sobreajustando
- Podemos reentrenar el modelo con todos los datos y obtener un modelo más robusto

Multiple expressions in one cell.

How would you like to fix it?

- Split this cell into 8 cells, or
- Wrap all code in a `begin ... end` block.

```

1  # Convertimos las señales de test y train a una matriz de Float32
2  datasignals_test = Float32.(Matrix(df_datasignals_out[:,1:3]))'
3  datasignals = Float32.(Matrix(df_datasignals_minus_out[:,1:3]))'
4
5  # Extraemos las columnas objetivo de test y train
6  σ_test = df_datasignals_out[:,4]
7  lcm_test = df_datasignals_out[:,5]
8
9  σ_col = df_datasignals_minus_out[:,4]
10 lcm_col = df_datasignals_minus_out[:,5]
11
12 # Las concatenamos en dataparams
13 dataparams = hcat(lcm_col, σ_col)'
14 dataparams_test = hcat(lcm_test, σ_test)';

```


Multiple expressions in one cell.

How would you like to fix it?

- Split this cell into 2 cells, or
- Wrap all code in a `begin ... end` block.

```
1 # Función de loss
2 function loss(x,y)
3     y_hat = re(p)(x)
4     return Flux.mse(y_hat, y)
5 end
6
7 # Loss compuesto para hacer un fine tune comparando la predicción de la red con la
  distribución log-normal
8 function composed_loss(x,y)
9     y_hat = model(x)
10    Pln_predicted = Pln.(y_hat[1,:], y_hat[2,:])
11    Pln_real = Pln.(y[1,:], y[2,:])
12    return mean(Flux.mse.(Pln_predicted,Pln_real)) + Flux.mse(y_hat, y)
13 end
```

Multiple expressions in one cell.

How would you like to fix it?

- Split this cell into 15 cells, or
- Wrap all code in a `begin ... end` block.

```
1 # Definimos la red neuronal
2 model = Chain(
3     Dense(3, 32),
4     Dense(32, 64, tanh_fast),
5     Dense(64, 32, tanh_fast),
6     Dense(32, 16, tanh_fast),
7     Dense(16, 2, softplus),
8 )
9
10 # Cargamos los parámetros del modelo
11 path = "C:/Users/Propietario/OneDrive/Escritorio/ib/Tesis_V1/MLonNMR/2-
    NNIdentificaHahn/G_8.73e-7_TE_1_AE/Params/017_params.csv"
12
13 p, re = Flux.destructure(model)
14
15 # Definimos el tamaño de los mini batchs para entrenar
16 batch_size = 100
17
18 # Usamos DataLoader para cargar los datos
19 data = Flux.DataLoader((datasignals, dataparams), batchsize = batch_size, shuffle =
    true)
20
21 # Parámetros de la red neuronal
22 params = Flux.params(model)
23
24 # Definimos el vector donde guardamos la pérdida
25 losses = []
26
27 # Definimos una funcion de callback para ver el progreso del entrenamiento
28 iter = 0
29 cb = function()
30     global iter += 1
31     if iter % length(data) == 0
32         epoch = iter ÷ length(data)
33         actual_loss = loss(data.data[1], data.data[2])
34         if epoch%100 == 0
35             println("Epoch $epoch || Loss = $actual_loss")
36         end
37         push!(losses, actual_loss)
38     end
39 end;
40
41 losses_composed = []
42 cb2 = function()
43     global iter += 1
44     epoch = iter ÷ length(data)
45     actual_loss = composed_loss(data.data[1], data.data[2])
46     println("Epoch $epoch || Loss = $actual_loss")
```

```

47     push!(losses_composed, actual_loss)
48 end;
49
50 # Definimos el modo de aprendizaje y la tasa de aprendizaje
51  $\eta$  = 1e-4
52 opt = ADAM( $\eta$ )
53
54 # Definimos el número de épocas
55 epochs = 5000;
56
57 println("Se va a entrenar una red con \n batch_size: $(batch_size) \n  $\eta$ : $( $\eta$ ) \n
Optimizer: $(opt) \n Epochs: $(epochs)")

```

UndefVarError: `epochs` not defined

Stack trace

Here is what happened, the most recent locations are first:

1. (This cell: line 2

```

1 # Entrenamos la red neuronal con el loss mse
2 for epoch in 1:epochs
3     Flux.train!(loss, Flux.params(p), data, opt, cb=cb)

```

```

1 # Entrenamos la red neuronal con el loss mse
2 for epoch in 1:epochs
3     Flux.train!(loss, Flux.params(p), data, opt, cb=cb)
4 end
5
6 # Entrenamiento con loss compuesto
7 # for epoch in 1:1
8 #     Flux.train!(composed_loss, Flux.params(model, opt), data, opt, cb=cb2)
9 # end

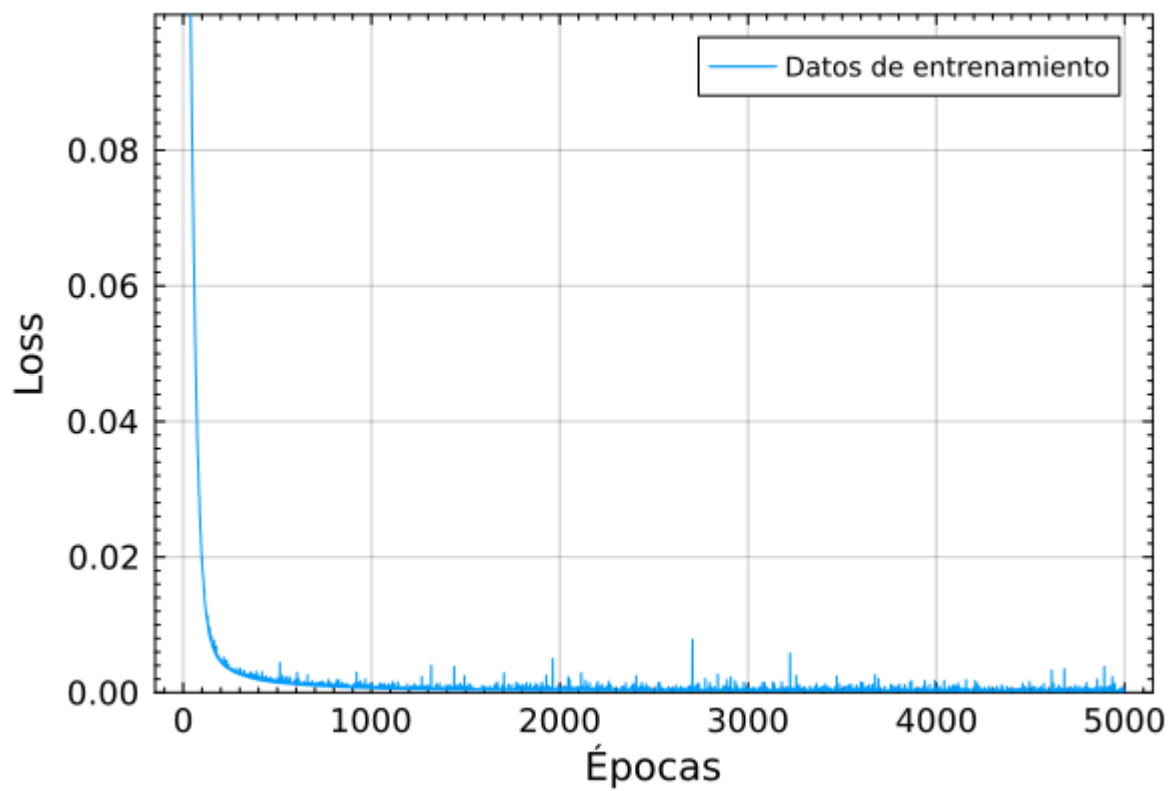
```

Loss luego de reentrenar con los datos de validación sabiendo que no hay overfitting

```

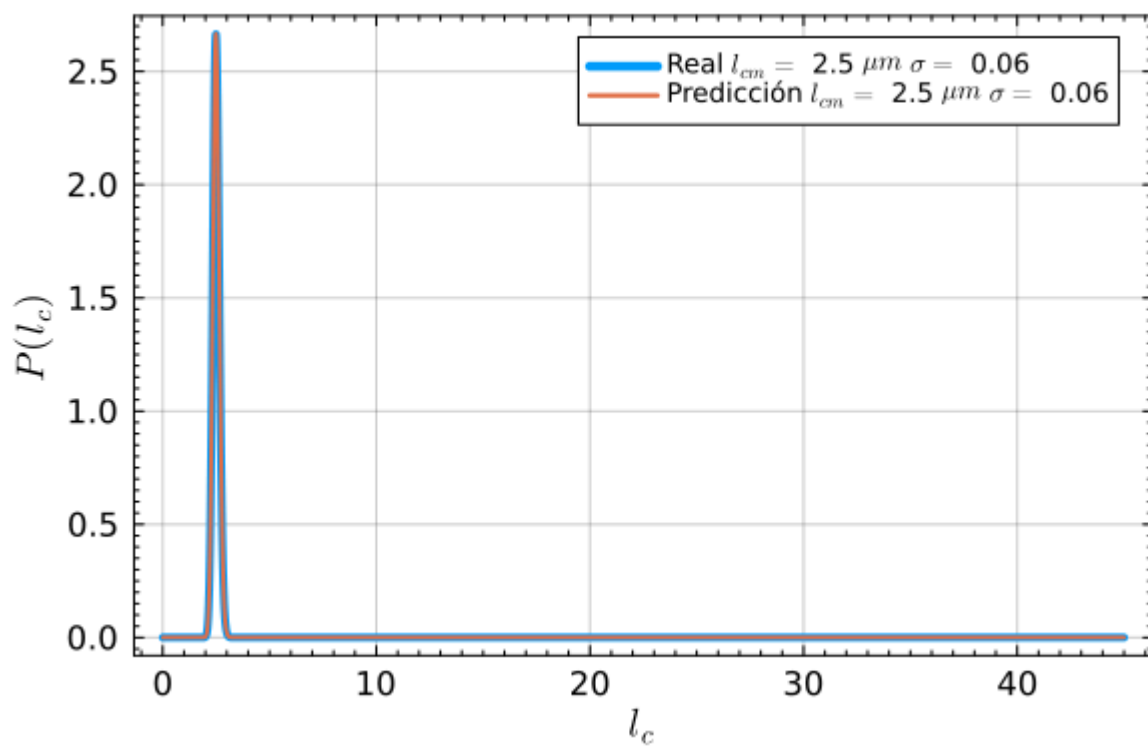
1 using PlutoUI

```

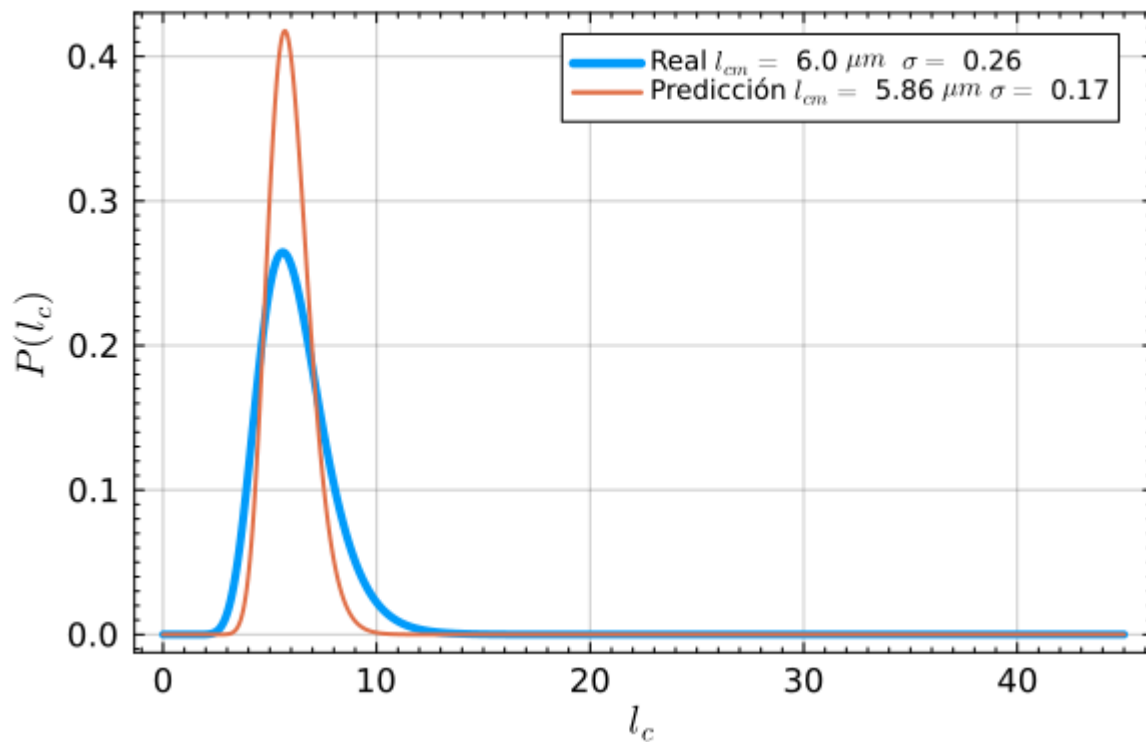


- El error raíz cuadrática media de la red neuronal en el conjunto de entrenamiento es 0.02

Predicción con menos error

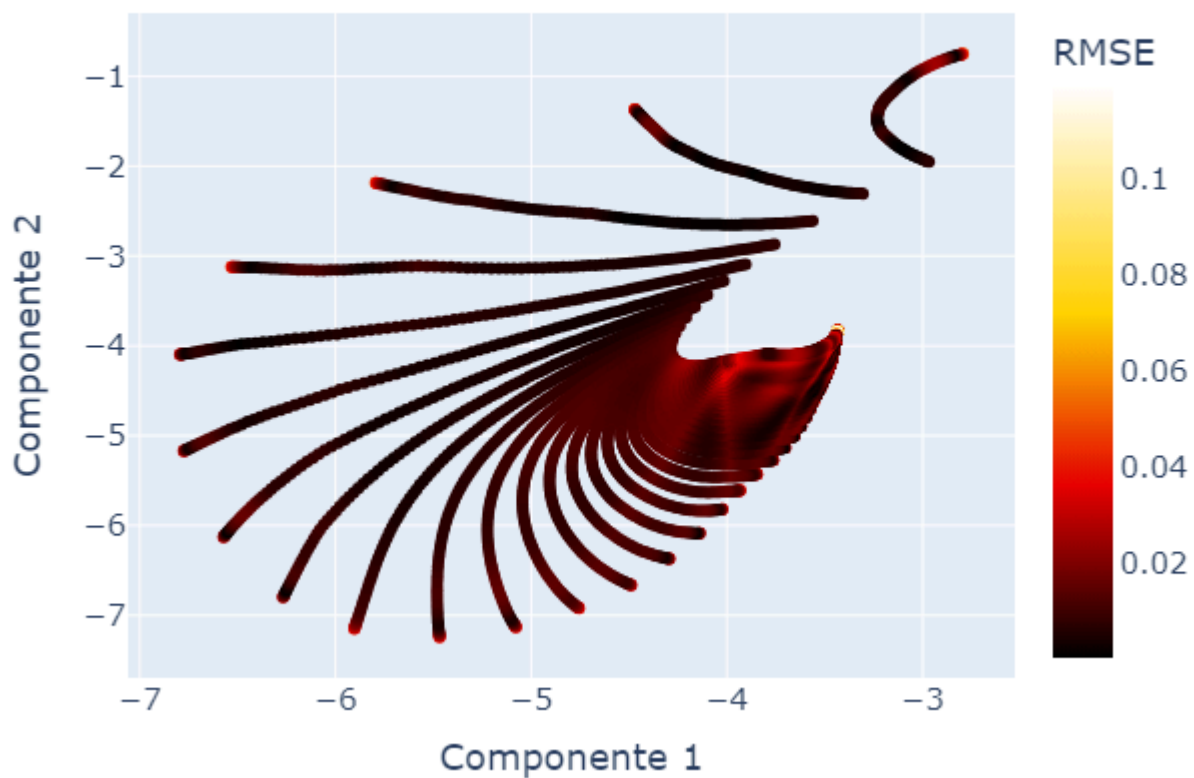


Predicción con mas error



- Este modelo funciona aún mejor que el anterior quizás porque esta vez la representación en baja dimensionalidad de donde entran los datos a esta red proviene de un autoencoder.
- El error RMSE en el conjunto de test es 0.02.

Todos los datos



Recordar que la representación aunque cambió sigue siendo similar en el hecho que σ y l_{cm} pequeños es donde mas separadas y van aumentando a medida que se acercan a la zona donde hay mas puntos acumulados. Dejo un HTML en el trello con el gráfico que indica cuando pasas el mouse el sigma y en color el lcm