

Homework 1 - Linear Regression, Cross Validation, and Nonlinear Regression

EE/CS5841, Spring 2024

```
1 # For this Homework I have used the following refernces :
2 # CS229 : https://www.youtube.com/watch?v=jGwO_UgTS7I&list=PLoROMvody4rMiGQp3WXShTMGgzqpfVfbU
3 # LLM : I have used chatgpt, perplexity ai, to understand the code and how to do
  programming in julia from what I have known by doing mathematics on paper.
4 # Other resources :
5 # 1. https://web.mit.edu/zoya/www/linearRegression.pdf
6 # 2. https://discourse.julialang.org/t/simple-tool-for-train-test-split/473
```

Part 1: Linear Regression by least-squares (10 points)

```
1 begin
2     using LinearAlgebra
3     using Random
4     using DelimitedFiles
5 end
```

Function definitions:

linear_regression_train

Creates weights for a least-squares solution for linear regression.

Inputs:

x - n samples \times m features

y - n samples \times 1 (for now we can assume that we have only a single output)

Outputs:

w - m features \times 1 (again, assuming single output variable)

linear_regression_infer

Predicts an output based on weights and inputs

Inputs:

x - n samples \times m features

w - m features \times 1

Outputs $y_{\text{predicted}}$ - n samples \times 1

Hints:

- The LinearAlgebra library includes useful things like pseudoinverse
- We can transpose with `transpose()` or using `'`
- We can check dimensions of matrices with `size()`
- Some useful matrices can be made with `ones()`, `zeros()`, and `eye()`
- Julia expects us to specify type for the above matrices, just use `Float64` for most things unless you want an `Int`
- We can force arrays to have singleton dimensions with a `;;` - see the simple test below for an example

`linear_regression_train` (generic function with 1 method)

```
1 function linear_regression_train(x,y)
2     # your code here!
3     x= hcat(ones(size(x, 1)), x)
4     weights = inv(transpose(x) * x) * transpose(x) * y
5     return weights
6 end
```

linear_regression_infer (generic function with 1 method)

```
1 function linear_regression_infer(x,w)
2   # your code here
3   x= hcat(ones(size(x, 1)), x)
4   predicted_y = x * w
5   return predicted_y
6 end
7
```

Here's a simple test

```
1 begin #we have to use these for Pluto notebooks when we execute multiple things
2 in the same cell
3   x_demo = [1.0; 2.0; 3.0 ;;]
4   y_demo = [3.0; 5.0; 7.0 ;;] # the ";;" makes it a column vector
5   try
6     w = linear_regression_train(x_demo,y_demo)
7     y_demo_pred = linear_regression_infer(x_demo,w)
8
9     if norm(y_demo - y_demo_pred) < 1e-10
10      print("It works!")
11    else
12      print("Not working yet")
13    end
14
15  catch
16    print("Something's wrong, perhaps dimensions aren't matching in your
17      function?")
18  end
19 end
20
```

It works!



Part 2, test train splitting and cross validation (10 points)

Function definitions:

mse

Mean squared error

Inputs:

y - a vector of true values y_{est} - a vector of estimated values

Outputs:

error - the mean squared error between the two vectors

test_train_split

Generates some indices for randomly splitting data. Assume an 80/20 split can just be hardcoded in. It is also acceptable to have this function take in the data instead of the number of points and the first version of the assignment did this.

Inputs:

n - number of data points

Outputs:

test_train_indices - a vector of 1's and 0's where 1's represent the entries to be used for test and the 0's represent entries to be used for training

cross_validate_split

Generate some indices for a cross validation

Inputs:

n - number of data points

k - number of splits for cross validation

Outputs:

split_indices - a vector of values representing which validation split that entry belongs to

eval_model

Function that allows you to pass in data, test/train split indices, model and loss functions, and runs everything for you

Inputs:

x - n samples \times m features

y - n samples \times 1 (for now we can assume that we have only a single output)

test_train_indices - a vector of 1's and 0's where 1's represent the entries to be used for test and the 0's represent entries to be used for training

model_train - name of the function that you are using to create the model, assumes it uses the same model representation as model_infer

model_infer - name of the function you are using to run inference on the model

loss - function used to compute loss

Outputs:

error - the error for the

run_cross_validation

Function that allows you to pass in data, cross validation split indices, model and loss functions, and runs everything for you

Inputs:

x - n samples \times m features

y - n samples \times 1 (for now we can assume that we have only a single output)

split_indices - a vector of 1's and 0's where 1's represent the entries to be used for test and the 0's represent entries to be used for training

model_train - name of the function that you are using to create the model, assumes it uses the same model representation as model_infer

model_infer - name of the function you are using to run inference on the model

loss - function used to compute loss

Outputs:

loss_per_run - a vector that is the same length as the number of cross validation folds that has the validation loss for that fold. if you're really motivated, you could make it output a training and validation loss separately

Hints:

- use vectorized dot operations, for example if we wanted to find everything equal to 1 in an array, we could use '`== 1`' and it would return a 'true' for every element that had a 1

mse (generic function with 1 method)

```
1 function mse(y,y_est)
2     # Your code here VV
3     n = length(y)
4     error = sum((y .- y_est).^2) / n
5     return error
6 end
7
```

Let's write a simple test for this

```
1 begin
2     y1 = [1 0 0]
3     y2 = [1 2 0]
4     try
5         error = mse(y1,y2)
6         if error - 1.33333333 < 1e-3
7             print("Yay, it works!")
8         else
9             print("Something isn't right")
10        end
11    catch
12        print("Something isn't working")
13    end
14 end
```

Yay, it works!



test_train_split (generic function with 1 method)

```
1 function test_train_split(n)
2     train_size = Int(0.8 * n)
3     test_size = n - train_size
4     test_train_indices = vcat(zeros(train_size), ones(test_size))
5     test_train_indices = shuffle(test_train_indices)
6     return test_train_indices
7 end
```


2x1 Matrix{Int64}:

5
8

```

1  #lazy visual test, needs asserts/condition checking for a proper test
2  #we can select from data using this one-hot by making it into a boolean
3  begin
4      try
5
6          example_split = test_train_split(10)
7          selection_test_x = [1:10 ;;] #force it to be an nx1 if we only have 1
8          feature
9          example_train_indices = iszero.(example_split)
10         test_indices = .!example_train_indices
11         selection_test_x[test_indices,:]
12     catch
13     end
14 end

```

cross_validate_split (generic function with 1 method)

```

1  function cross_validate_split(n,k)
2      # Your code here
3      split_size = div(n, k)
4      split_indices = zeros(Int, n)
5      for i in 1:k
6          start_idx = (i - 1) * split_size + 1
7          end_idx = min(i * split_size, n)
8          split_indices[start_idx:end_idx] .= i
9      end
10     split_indices = shuffle(split_indices)
11     return split_indices
12 end

```

[4, 3, 5, 3, 1, 4, 2, 5, 2, 1]

```

1  #little visual test - it will update once the above code works
2  try
3      cross_validate_split(10,5)
4  catch
5  end
6

```

```

1  #some demo code to show you how you could use the split indices to separate some
2  data
3  #this may be useful a little bit further down...
4  begin
5      try
6
7          example_split = cross_validate_split(10,5)
8          selection_test_x = [1:10 ;;] #force it to be an nx1 if we only have 1
9          feature
10         splits = zeros(5,2)
11         for split = 1:5
12             val_set_indices = findall(==(split),example_split)
13             splits[split,:] = selection_test_x[val_set_indices,:]
14         end
15         print(splits)
16     catch
17     end
18 end

```

```
[2.0 5.0; 6.0 7.0; 4.0 10.0; 3.0 8.0; 1.0 9.0]
```



This next part is just combining stuff from Part 1 with the loss function (MSE) from Part 2. It's just to see how we can pass functions into functions to do things for us.

eval_model (generic function with 1 method)

```

1  function eval_model(x,y,test_train_indices,model_train,model_infer,loss)
2      x = [x ;;] # force it to be an nx1 instead of just an n-element array if x is
3      1D
4      x_train = x[test_train_indices .== 0,:]
5      y_train = y[test_train_indices .== 0,:]
6      x_test = x[test_train_indices .== 1,:]
7      y_test = y[test_train_indices .==1,:]
8      trained_model = model_train(x_train, y_train)
9      y_pred = model_infer(x_test, trained_model)
10     error = loss(y_test, y_pred)
11     return error
12 end

```

Let's reuse our test code from Part 1, modified slightly. This isn't a thorough test, as it's a zero-error case with almost no data, but we could write a second test with a higher known error if we wanted to include that as well.

```

1 begin
2   try
3     test_train_indices = [0,1,0]
4     loss =
       eval_model(x_demo,y_demo,test_train_indices,linear_regression_train,linear_regression_infer,mse)
5
6     if loss < 1e-3
7       print("It works!")
8     else
9       print("Not working yet")
10    end
11
12  catch
13    print("Something's wrong, perhaps dimensions aren't matching in your function?")
14  end
15
16 end

```

It works!



run_cross_validation (generic function with 1 method)

```

1 function run_cross_validation(x,y,split_indices,model_train,model_infer,loss)
2   # Your code here
3   x = [x ;;] #force it to be an nx1 if we only have 1 feature
4   n_folds = maximum(split_indices)
5   loss_per_run = zeros(n_folds)
6   for fold in 1:n_folds
7     test_train_indices = split_indices .== fold
8     loss_per_run[fold] = eval_model(x, y, test_train_indices, model_train,
model_infer, loss)
9   end
10
11   return loss_per_run
12 end

```

Test it out on some synthetic data. Create a line with a slope of 2 and an intercept of 1 and add in some Gaussian noise

Run 5-fold cross validation on it and don't worry about a test/train split for now.

1.072133458475335

```
1 begin
2   try
3     x = collect(1:100)
4     y = x*2+ones(Float64,100)+randn(Float64,100)
5     split_indices = cross_validate_split(100,5)
6     loss_per_run =
7       run_cross_validation(x,y,split_indices,linear_regression_train,linear_regression_infer,mse)
8     mean_error = sum(loss_per_run)/5
9   catch
10  end
```


Part 3 Regularized Linear Regression (10 points)

For now we will just use L2 regularization, so modify your code from the first set of linear regression functions to include an L2 regularization term

regularized_linear_regression_train

Performs L2-regularized linear regression

Inputs:

x - n samples \times m features

y - n samples \times 1 (for now we can assume that we have only a single output)

λ - regularization constant

Outputs:

w - m features \times 1 (again, assuming single output variable)

regularized_linear_regression_train (generic function with 1 method)

```
1 function regularized_linear_regression_train(x,y,lambda)
2     # your code here!
3     x = hcat(ones(size(x, 1)), x)
4     m = size(x, 2)
5     regularization_matrix = lambda * I(m)
6     weights = inv(transpose(x) * x + regularization_matrix) * transpose(x) * y
7     return weights
8 end
```

1 Enter cell code...

```
3×1 Matrix{Float64}:
1.357803640704271
0.674425544613559
0.6744255446135217
```

```
1 #let's do a simple test
2 try
3     begin
4         x_demo_reg = 1.0*[1 1; 3 3; 4 4]
5         y_demo_reg = [3.0; 5.0; 7.0 ;;] # the ";" makes it a column vector
6         w_reg01 = regularized_linear_regression_train(x_demo_reg,y_demo_reg,0.1)
7     end
8 catch
9 end
```

```
3x1 Matrix{Float64}:  
0.7023809523809517  
0.7619047619047639  
0.7619047619047632
```

```
1 #for fun we can vary the regularization constant  
2 try  
3     begin  
4         x_demo_reg = 1.0*[1 1; 3 3; 4 4]  
5         y_demo_reg = [3.0; 5.0; 7.0 ;;]  
6         w_reg1 = regularized_linear_regression_train(x_demo_reg,y_demo_reg,1)  
7     end  
8 catch  
9 end  
10
```

```
1 #we can compare with no regularization as well  
2 try  
3     begin  
4         x_demo_reg = 1.0*[1 1; 3 3; 4 4]  
5         y_demo_reg = [3.0; 5.0; 7.0 ;;]  
6         w_noreg = linear_regression_train(x_demo_reg,y_demo_reg)  
7     end  
8 catch  
9 end  
10
```


Part 4 Basis functions (10 points)

Modify the function below to create n polynomials for the input features x in increasing power

polynomial_basis

Input:

x - n samples \times m features n - highest polynomial to create

Output:

xns - n samples \times $m \times n$ features. features should be ordered as $x_{1n1} \ x_{2n1} \ \dots \ x_{mn1} \ x_{1n2} \ x_{2n2} \ \dots \ x_{mn2} \ \dots \ x_{mnn}$ for each row of samples

Hints:

- We can concatenate arrays using square brackets
- We can do elementwise operations with the dot operator. ie. $x.^2$ squares each element

polynomial_basis (generic function with 1 method)

```
1 function polynomial_basis(x,n)
2     # Your code here
3     xns = []
4     for power in 1:n
5         xns = hcat([x .^ power for power in 1:n]...)
6     end
7
8     return xns #xns should be {x,x^2,...x^n}
9 end
10
```

2×6 Matrix{Int64}:

```
1 2 1 4 1 8
3 4 9 16 27 64
```

```
1 #simple visual test
2 begin
3     try
4         x_mat = [1 2; 3 4]
5         xmat_n2= polynomial_basis(x_mat,3)
6     catch
7     end
8 end
```


Part 5 Put it all together! (10 points)

Download and load the following dataset: Yacht Hydrodynamics: <https://archive.ics.uci.edu/ml/datasets/Yacht+Hydrodynamics> Load the data and run some experiments with it using regularized linear or polynomial regression.

For full credit, separate off 20% as a test set, explore at least 3 options (regularization and/or polynomial basis functions), use cross validation to select the best one, and report the error on the held out dataset. Most importantly, report your conclusions and justify it with a sentence or two!!


```
1 begin
2     #Download the data and load it using DelimitedFiles function
3     #update the code with your file location or use another loading method
4     local_data = readdlm("yacht_data.csv", ',', Any, '\n')
5     local_data_to_use = convert(Array{Float64,2}, local_data[:,1:7])
6     local_n = size(local_data_to_use, 1)
7     new_split_indices = cross_validate_split(local_n, 5)
8     x_new = local_data_to_use[:, 1:6]
9     x_new = [x_new ;;]
10    y_new = local_data_to_use[:, 7]
11    y_new = [y_new ;;]
12    lambda_values = 0.0001:0.5:10.0
13
14    function new_run_cross_validation(x_new, y_new,
15        new_split_indices, regularized_linear_regression_train,
16        linear_regression_train, linear_regression_infer, mse, lambda=nothing)
17        n_folds = maximum(new_split_indices)
18        loss_per_run = zeros(n_folds)
19
20        for fold in 1:n_folds
21            test_train_indices = new_split_indices .== fold
22            x_train = x_new[.!test_train_indices, :]
23            y_train = y_new[.!test_train_indices, :]
24            x_test = x_new[test_train_indices, :]
25            y_test = y_new[test_train_indices, :]
```

```
24
25     if lambda === nothing
26         w = linear_regression_train(x_train, y_train)
27     else
28
29         w = regularized_linear_regression_train(x_train, y_train, lambda)
30     end
31
32     y_pred = linear_regression_infer(x_test, w)
33     loss_per_run[fold] = mse(y_test, y_pred)
34 end
35
36 return loss_per_run
37 end
38
39
40 # Calculate errors for each lambda value
41 errors_reg = [new_run_cross_validation(x_new, y_new, new_split_indices,
regularized_linear_regression_train, linear_regression_train,
linear_regression_infer, mse, lambda_val) for lambda_val in lambda_values]
42
43
44 # Find the index of the minimum error
45 best_lambda_index_reg = argmin(errors_reg)
46 # Get the corresponding lambda value
47 best_lambda_reg = lambda_values[best_lambda_index_reg]
48
49 # Print the results
50 println("Mean Loss for Each Lambda: ", errors_reg)
51 println("Best Lambda Index: ", best_lambda_index_reg)
52 println("Best Lambda Value: ", best_lambda_reg)
53
54 println("Best Lambda for Regularized Regression: $best_lambda_reg, Mean Loss:
$(errors_reg[best_lambda_index_reg])")
55
56 # Run cross-validation for polynomial regression with degrees 1 to 5
57 errors_poly = []
58 for degree in 1:10
59     # Generate polynomial basis
60     x_poly = polynomial_basis(x_new, degree)
61
62     # Run cross-validation
63     errors = [new_run_cross_validation(x_poly, y_new, new_split_indices,
regularized_linear_regression_train, linear_regression_train,
linear_regression_infer, mse, lambda_val) for lambda_val in
lambda_values]
64
65
66     # Save the average error for this degree
67     push!(errors_poly, sum(errors) / 5)
68 end
69 # Find the degree with the minimum average error
70 best_degree_poly = argmin(errors_poly)
71
72 # Report results for polynomial regression
73 println("Best Degree for Polynomial Regression: ", best_degree_poly)
74 println("Average Cross-Validation Error for Polynomial Regression: ",
errors_poly[best_degree_poly])
75
```

```
75
76
77 function new_test_train_split(local_n)
78     train_size = round(Int, 0.8 * local_n)
79     test_size = local_n - train_size
80     train_test_indices = vcat(zeros(train_size), ones(test_size))
81     train_test_indices = shuffle(train_test_indices)
82     return train_test_indices
83 end
84
85
86 # Now, separate off 20% as a test set and report the error on the held-out
87 dataset for the best model
88 new_train_test_indices = new_test_train_split(local_n)
89 x_train = x_new[new_train_test_indices .== 0, :]
90 y_train = y_new[new_train_test_indices .== 0, :]
91 x_test = x_new[new_train_test_indices .== 1, :]
92 y_test = y_new[new_train_test_indices .== 1, :]
93
94 # Train the best model on the training set
95 if best_degree_poly > 1
96     x_train = polynomial_basis(x_train, best_degree_poly)
97     x_test = polynomial_basis(x_test, best_degree_poly)
98 end
99
100 if best_lambda_reg > 0
101     w_best = regularized_linear_regression_train(x_train, y_train,
102     best_lambda_reg)
103 else
104     w_best = linear_regression_train(x_train, y_train)
105 end
106
107 # Make predictions on the test set
108 y_pred_test = linear_regression_infer(x_test, w_best)
109
110 # Calculate and report the test error
111 test_error = mse(y_test, y_pred_test)
112 println("Test Error for the Best Model: ", test_error)
113 #Ignore the 8th null column and reference the website for what each column
114 means
115
116 end
```



```
Mean Loss for Each Lambda: [[73.13597893840252, 84.79323201465455, 80.74 ②
59135429622, 85.73458126160574, 71.51706649145336], [75.13862315962967, 88.
01833148487924, 81.39312509227567, 87.72265969220771, 74.96478232708034],
[82.625104355779, 96.70165001363551, 86.97338935924103, 96.74432025020333,
82.52092483366155], [90.96468316808907, 106.62959979498498, 93.640236364392
49, 106.63385734637447, 90.30992088613199], [99.03799704313319, 116.3566595
8118247, 100.28379896761633, 116.14096666091069, 97.6676907148307], [106.50
006097699246, 125.3922519985465, 106.52568354686433, 124.88204827657839, 10
4.39868931599479], [113.27418921571079, 133.60857185132258, 112.25340635423
643, 132.78165441961707, 110.47898968933555], [119.38305979662273, 141.0182
7327419355, 117.45868664116126, 139.87735202081083, 115.94828036258471], [1
24.88403958480914, 147.6853638994207, 122.17360973019514, 146.2445118598747
6, 120.86679512768248], [129.84306759379595, 153.6883621888366, 126.4438379
6512841, 151.96628538590278, 125.29774203852472], [134.32414075232174, 159.
1051553834607, 130.31716209766043, 157.12184541014557, 129.30037996447197],
[138.3854191481204, 164.00718863897853, 133.83875265131195, 161.78228483975
84, 132.927668738518], [142.07818381599512, 168.45771759665584, 137.0494195
519167, 166.00974660504852, 136.22588006286463], [145.44699762170012, 172.5
1179970000618, 139.9852254141762, 169.85786838183998, 139.23501059999953],
[148.53033027838975, 176.2169771495621, 142.67768782795318, 173.37269844108
48, 141.98947755973023], [151.3613159957836, 179.6141786797402, 145.1542120
9326806, 176.5937138328199, 144.51886649949225], [153.9684993622348, 182.73
863171731264, 147.4385870286641, 179.5547868942413, 146.84863449427198], [1
56.3765119649631, 185.62069961699743, 149.5514682797893, 182.28504380708983
, 149.0007335347996], [158.60666237351506, 188.28661580137188, 151.51081798
307814, 184.80960273791345, 150.994147033487], [160.67744002445372, 190.759
11241328606, 153.33229097695605, 187.15019774195946, 152.84534436035028]]
Best Lambda Index: 1
Best Lambda Value: 0.0001
Best Lambda for Regularized Regression: 0.0001, Mean Loss: [73.135978938402
52, 84.79323201465455, 80.7459135429622, 85.73458126160574, 71.517066491453
36]
```

4.4171410930234085

1 test_error

- 1 *# Conclusions and Justify*
- 2 *#The test error for the best model was calculated to be 4.417, indicating that the model's performance on new data is relatively low, it will generalize well to new data.*
- 3 *# The mean loss for regularized regression was minimized at a lambda value of 0.0001, a very small regularization parameter was most effective in reducing overfitting. the best degree for polynomial regression is 4. The optimal balance between bias and variance in this model is good , it has the lowest average cross-validation error.*

