

```
In [47]: # 1

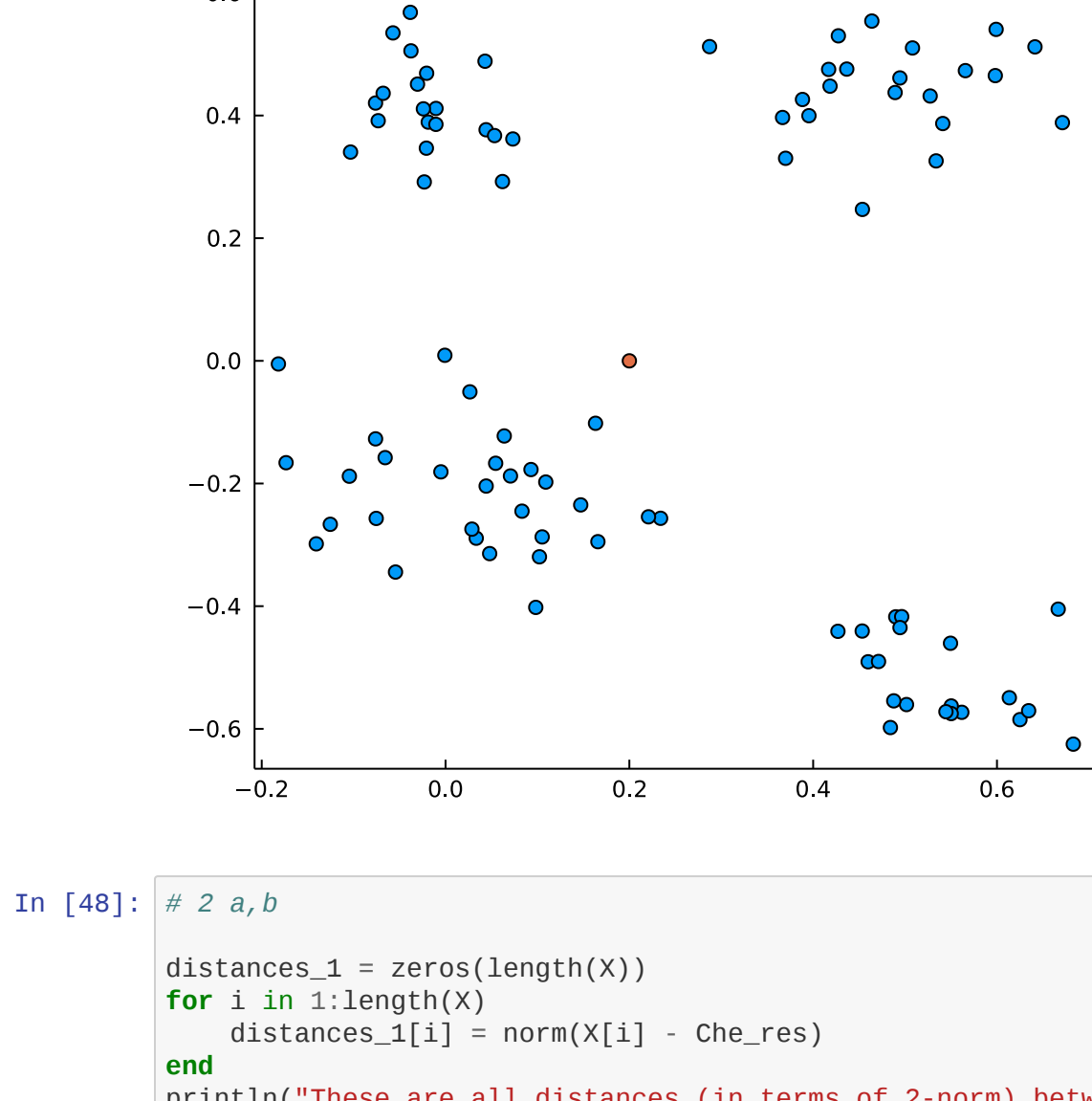
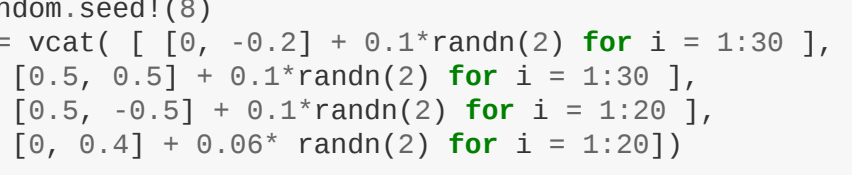
using LinearAlgebra
using Statistics
import Random
using Plots
using WMLS

Random.seed!(8)
X = vcat([ [ 0., -0.2] + 0.1*randn(2) for i = 1:30 ],
          [ [ 0.5, 0.5] + 0.1*randn(2) for i = 1:30 ],
          [ [ 0.5, -0.5] + 0.1*randn(2) for i = 1:20 ],
          [ [ 0., 0.4] + 0.06* randn(2) for i = 1:120])

Che_res = [0.2, 0]

scatter([x[1] for x in X], [x[2] for x in X])
scatter!([Che_res[1]], [Che_res[2]])
plot!(legend = false, grid = false, size = (500, 500))

Out[47]:
```

[illegible]

```
distances_2 = zeros(length(X))
for i in 1:length(X)
    distances_2[i] = abs(X[i][1] - Che_res[1]) + abs(X[i][2] - Che_res[2])
end
# println(distances_2)
println("Customer ", findmax(distances_2)[2], " is the farthest customer to Chelsea's restaurant in terms of 1-norm.")
```

Customer 65 is the farthest customer to Chelsea's restaurant in terms of 1-norm.

```
In [50]: # 2 d 2-norm

distances_1 = zeros(length(X))
for i in 1:length(X)
    distances_1[i] = norm(X[i] - Che_res)
end

copy_distances_1 = copy(distances_1)

small_index_1 = []

for i in 1:5
    fake_index = findmin(copy_distances_1)[2]
    real_index = findall(x -> x==copy_distances_1[fake_index], distances_1)
    append!(small_index_1, real_index)
    deleteat!(copy_distances_1, findmin(copy_distances_1)[2])
end

println(small_index_1)

Any{30, 15, 5, 1, 29}
```

```

In [51]: # 2 0 1-norm
distances_2 = zeros(length(X))
for i in 1:length(X)
    distances_2[i] = abs(X[i][1] - Che_res[1]) + abs(X[i][2] - Che_res[2])
end
copy_distances_2 = copy(distances_2)
small_index_2 = []

for i in 1:5
    fake_index = findmin(copy_distances_2)[2]
    real_index = findall(x -> x==copy_distances_2[fake_index], distances_2)
    append!(small_index_2, real_index)
    deleteat!(copy_distances_2, findmin(copy_distances_2)[2])
end

println(small_index_2)

Any{30, 1, 15, 5, 16}

```

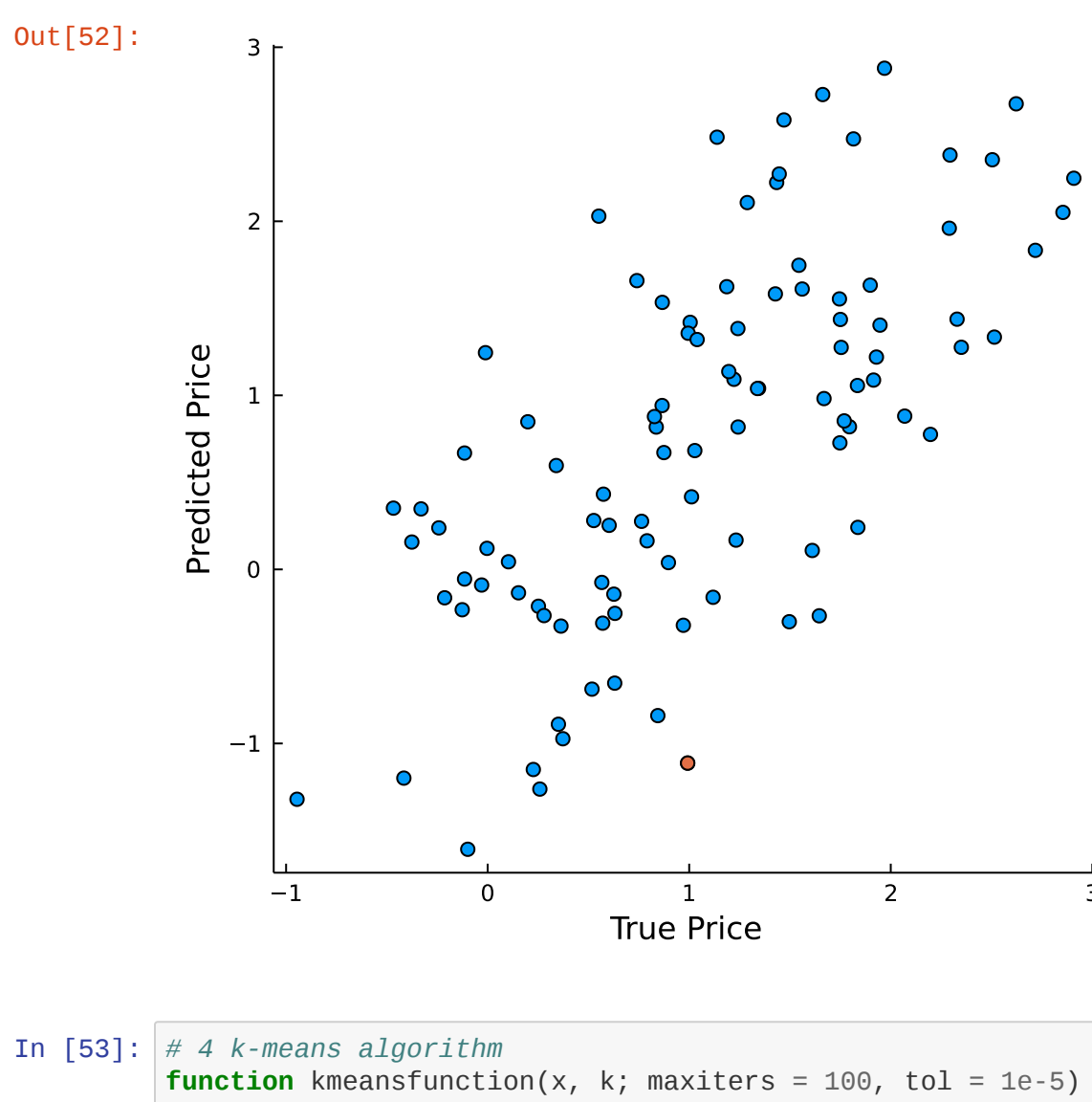
```
Random.seed!(8)
y = randn(100)
Random.seed!(8)
quality = y * 0.5 + rand(100) * 0.1
Random.seed!(10)
nutrition = y * 0.1 + rand(100) * 0.2
price = quality * 1 + nutrition * 0.1 + rand(100) * 2

predicted_price = 0.5*quality + 7.1*nutrition .* 0.1

println("Data point y, findmax([norm(price[i] - predicted_price[i]) for i = 1:length(price)])[2], ' is the farthest
(in terms of 2-norm.)' to its production.")

scatter(price, predicted_price)
scatter!([price[97]], [predicted_price[97]])
plot!(legend = false, grid = false, size = (500,500), xaxis = "True Price", yaxis = "Predicted Price")

Data point 97 is the farthest (in terms of 2-norm.) to its prediction.
```



```

distances = zeros(N) # used to store the distance of each point to the nearest representative.
reps = [zeros(k) for i=1:k] # used to store data representatives.

# 'assignment' is an array of N integers between 1 and k.
# The initial assignment is chosen randomly.
assignment = [rand(1,k) for i=1:N]
jprevius = Inf # used in stopping condition

for iter = 1:maxiters
    # cluster j representative is average of points in cluster j.
    for j = 1:k
        group = [i for i=1:N if assignment[i]==j] # get which point index belongs to the group j
        reps[j] = sum(x[group]) / length(group); # reps is the clustering center
    end;
    # For each x[i], find distance to the nearest representative
    # and its group index.
    for i = 1:N
        (distances[i], assignment[i]) = findmin(norm(x[i] - reps[j]) for j=1:k)
    end;

    # Compute clustering objective.
    J = norm(distances)^2 / N # distance to the clustering center

    # Show progress and terminate if J stopped decreasing.
    println("Iteration ", iter, ", J: ", J, ",")
    if (iter > 1) && (abs(J - jprevius) < tol)
        return assignment, reps
    end
    jprevius = J
end

```

```
end
end

Out[53]: kmeansfunction (generic function with 1 method)

In [54]: # 4 a, b
assignment, reps = kmeansfunction(X, 3)
println("\n\n", reps, " are the best locations to open these three restaurants for Chelsea's franchisee.")
```

```

Iteration 1: Jclust = 0.21805080743372227.
Iteration 2: Jclust = 0.054753907794250544.
Iteration 3: Jclust = 0.04764446133736133.
Iteration 4: Jclust = 0.047870784266229665.
Iteration 5: Jclust = 0.047870784266229665.

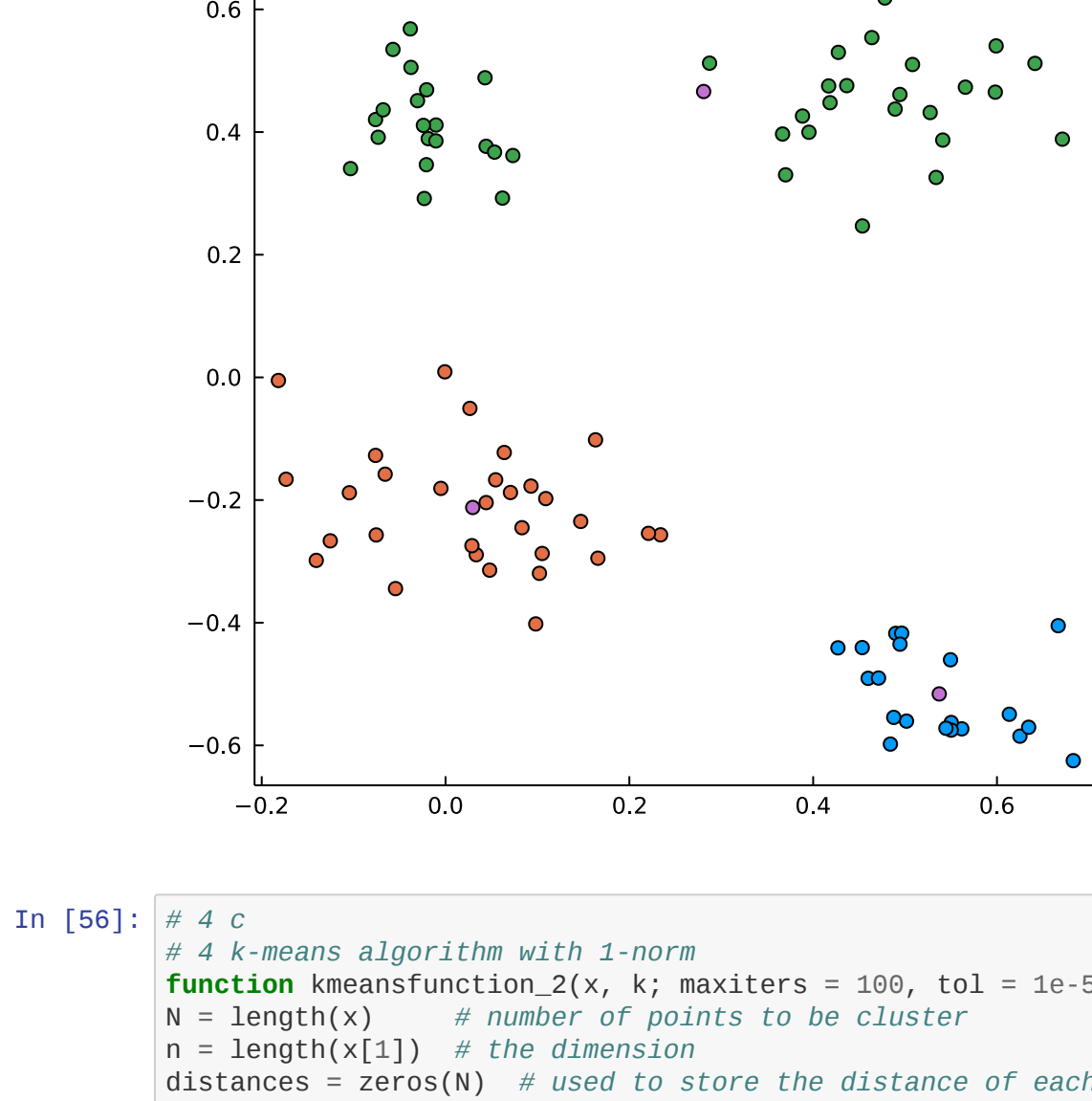
[0.53719099996616741, 0.5611717797970566], [0.02958765506726597, 0.21210461348896638], [0.2808296177102955, 0.46611678930286651], acc the best locations to open these three restaurants for Chelsea's franchisee

```

```
In [55]: k = 3 # Number of groups
```

```
N = length(X) # Total number of points to be clustered
grps = [[i|1] for i=1:N if assignment[i] == j] for j=1:k]
scatter([c1|1] for c in grps[1]), [c2|2] for c in grps[1]])
scatter([c1|1] for c in grps[2]), [c2|2] for c in grps[2]])
scatter([c1|1] for c in grps[3]), [c2|2] for c in grps[3]])
scatter([c1|1] for c in reps], [c2|2] for c in reps])
plot(legend = false, grid = false, size = (600,500))
```

Out[56]:



```
# 'assignment' is an array of N integers between 1 and k.
# The Initial assignment is chosen randomly.
assignment = [rand(1:k) for i=1:N]
jprevious = Inf # used in stopping condition

for iter = 1:maxiters
    # cluster j representative is average of points in cluster j.
    for j = 1:k
        group = [i for i=1:N if assignment[i]==j] # get which point index belongs to the group j
        reps[j] = sum(x[group]) / length(group); # reps is the clustering center
    end;
    # For each x[i,j], find distance to the nearest representative
    # and its group index.
    for i = 1:N
        (distances[i,j], assignment[i]) = ffindmin([abs(x[i][1] - reps[j][1]) + abs(x[i][2] - reps[j][2]) for j=1:k])
    end;
    # Compute clustering objective.
    J = norm(distances)^2 / N # distance to the clustering center

    # Show progress and terminate if J stopped decreasing.
    println("Iteration ", iter, ", J: ", J, ", ", N)
    if (iter > 1) && (abs(J - Jprevious) < tol)
        return assignment, reps
    end
    Jprevious = J
end

end
```

```
Out[56]: kmeansfunction_2 (generic function with 1 method)
```

```
In [57]: assignment, reps = kmeansfunction_2(X, 3)
```

```
Iteration 1: Jclust = 0.3403702703823343.
Iteration 2: Jclust = 0.09979403232802793.
Iteration 3: Jclust = 0.09823290083113136.
Iteration 4: Jclust = 0.097884909721599553.
Iteration 5: Jclust = 0.097395600044050234.
```

```
Iteration 5: JcIust = 0.09773526641653874.  
Iteration 6: JcIust = 0.09773526641653874.
```

Out[57]: ([1, 2, 2, 2, 2, 2, 2, 2 ... 1, 1, 1, 1, 1, 1, 1, 1], [[-0.02373170768873549, 0.33606929730232826], [0.25
790743657055338, -0.35895024685859395], [0.4792906462350338, 0.5022140724893905]])

```
In [58]: k = 3 # Number of groups
```

```

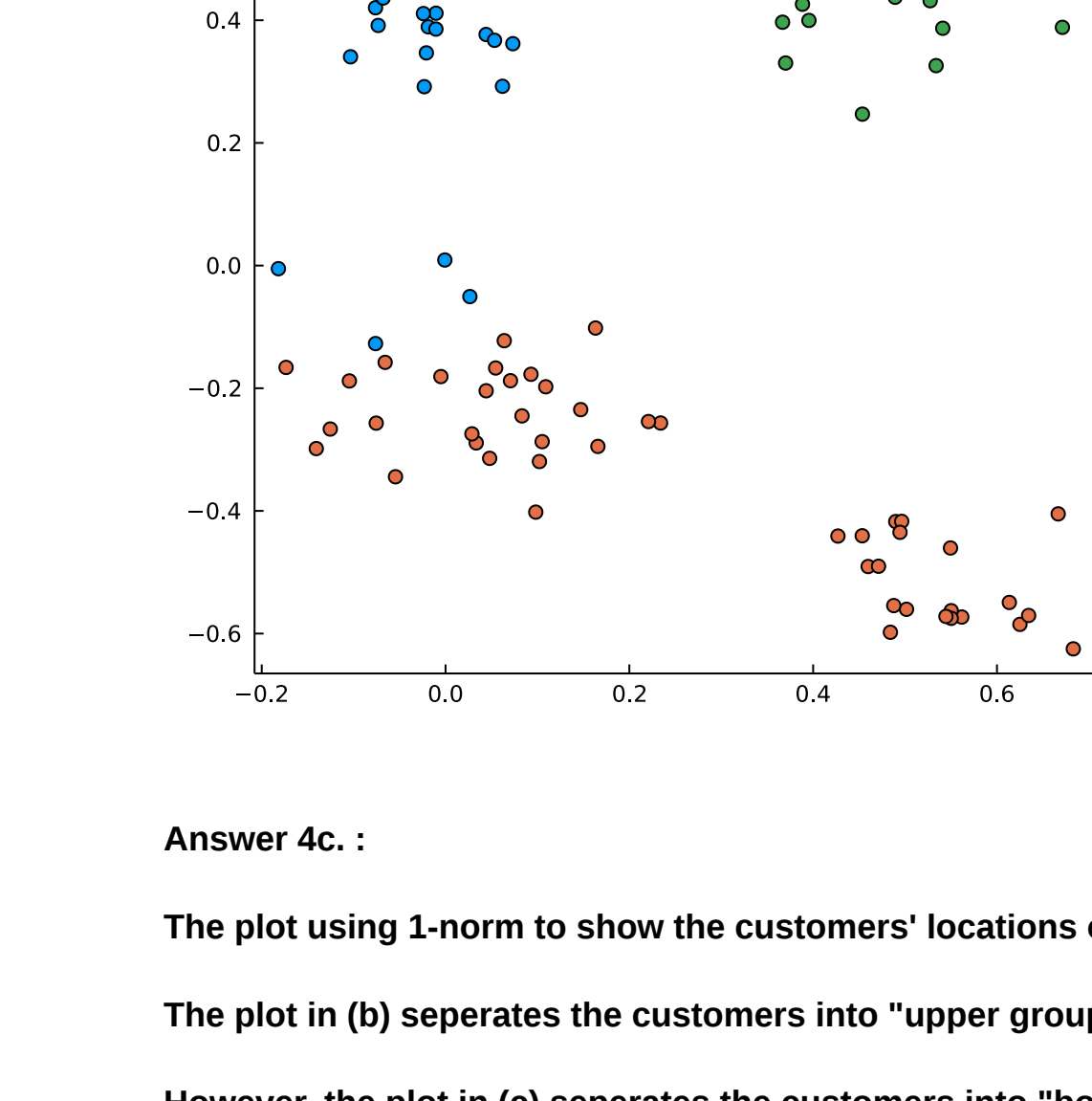
N = length(X) # Total number of points to be clustered
grps = [[X[i] for i=1:N if assignment[i] == j] for j=1:k]

```

```
scatter([c[1] for c in grps[1]], [c[2] for c in grps[1]])
scatter!([c[1] for c in grps[2]], [c[2] for c in grps[2]])
```

```
scatter!([c[1] for c in grps[3]], [c[2] for c in grps[3]])
plot!(legend = false, grid = false, size = (500,500))
```

Out[58]:



Answer 4d. :

3 is not a good clustering number because apparently there are "4" separated groups on the map.

4 would be a good number.

```
In [50]: # 5

function Gram_Schmidt(a; tol = 1e-10)
    copy_a = copy(a)
    copy_a = convert(Vector{Vector{Float64}}, copy_a)

    for i = 1:length(copy_a)
        for j = 1:i-1
            copy_a[i] -= (copy_a[j]'*copy_a[i]) * copy_a[j]
        end
    end
end
```

```

        printn("Vectors are linearly dependent.")
        return copy_a
    end

    copy_a[i] = copy_a[i]/norm(copy_a[i])
end

return copy_a
end

```