# CPSC340A2

Qinglan Huang n6v9a        Liu Yang v4d9

November 25, 2017

## 1 Random Forests

### 1.1 Implementation

1. If we change the depth parameters to Inf the training process will terminate not only when depth <1 but also when we reach baseSplit where the remaining objection of each leaf<1.

2.Because each random tree randomly pick only sqrt(d) features to train, it does not cover all data set and will never overfit no matter how deep the tree is.

3. check code at decisionTree.jl

```
294    function randomForest(X,y,depth,nTrees)
295        subModels = Array{GenericModel}(nTrees)
296
297        for i in 1: nTrees
298            subModels[i] = randomTree(X,y,depth)
299        end
300        return subModels
301
302    end
303
304
305
306    function predictA(subModels, Xhat)
307        (t,d) = size(Xhat)
308        yhat = zeros(t)
309        y = zeros(t,nTrees)
310
311        for i in 1: nTrees
312            model = subModels[i]
313
314            y[:,i] = model.predict(Xhat)
315        end
316        for i in 1:t
317            yhat[i] = mode(y[i,:])
318        end
319
320        return yhat
321    end
322
```

    4. Test result is:
Train Error with depth-Inf decision tree:0.000
Test Error with depth-Inf decision tree:0.367
Train Error with depth-5 decision tree:0.311
Test Error with depth-5 decision tree:0.504
Train Error with depth-Inf decision forest:0.000
Test Error with depth-Inf decision forest:0.170

A decision tree cannot avoid overfitting, therefore it has a training error of
0.However, there is a relatively large test error. A single random tree, since
only pick part of features, technically does not finish training process and has a
bad performance on both train error and test error. But when we increase the
number of random tree ( random forest), we can avoid overfitting and cover all
d features. Therefore the random forest has a low train error and test error as

2

expected

## 1.2 Very-Short Answer Questions

1. It is going to be expensive because of the runtime.

2. b c f
b. Decreasing the depth is going to decrease the probability of the overfitting.
c. larger amount of data can improve accuracy
f.When we have many strong features, random tree can be correlated, the potential of overfitting increases .

3.Adding more dataset for random forests, and building more random tree, control depth of the each random tree in order to control number of feature for clear overall accuracy, avioding strong feature.

# 2 K-Means Clustering

## 2.1 Selecting among k-means Initialization

1. check code at kMeans.jl

```
81
82   function kMeansError(X,y,W)
83       (n,d) = size(X)
84       (k,d2) = size(W)
85       assert(d == d2)
86       temp =0
87       for i in 1:n
88           for j in 1:d
89               temp += (X[i,j] - W[Int(y[i]),j]).^2
90           end
91       end
92       return temp
93   end
```
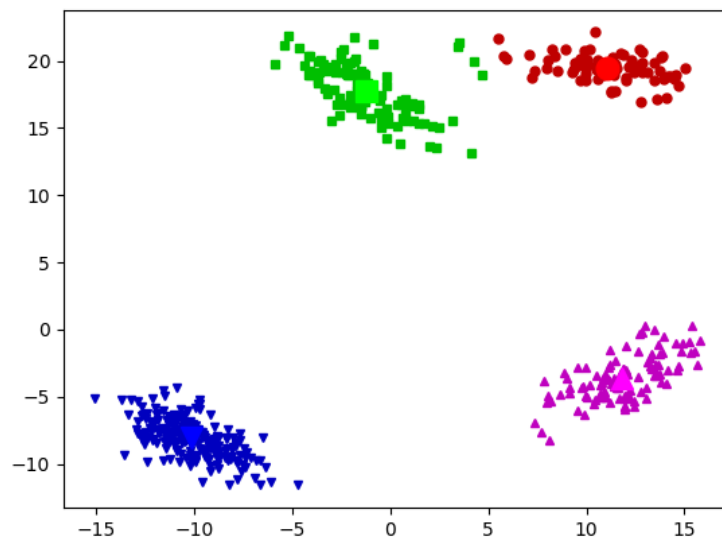
2.
kMeansError monotonically decreases until stabilization.

```
julia> include("example_Kmeans.jl")
Running k-means, changes = 500
Running k-means, kMeansError= 9609.749047512
Running k-means, changes = 32
Running k-means, kMeansError= 9466.879704072
Running k-means, changes = 20
Running k-means, kMeansError= 9404.445409588
Running k-means, changes = 11
Running k-means, kMeansError= 9386.083251910
Running k-means, changes = 10
Running k-means, kMeansError= 9376.066372056
Running k-means, changes = 5
Running k-means, kMeansError= 9371.944601192
Running k-means, changes = 3
Running k-means, kMeansError= 9370.487700096
Running k-means, changes = 0
Running k-means, kMeansError= 9370.487700096
```
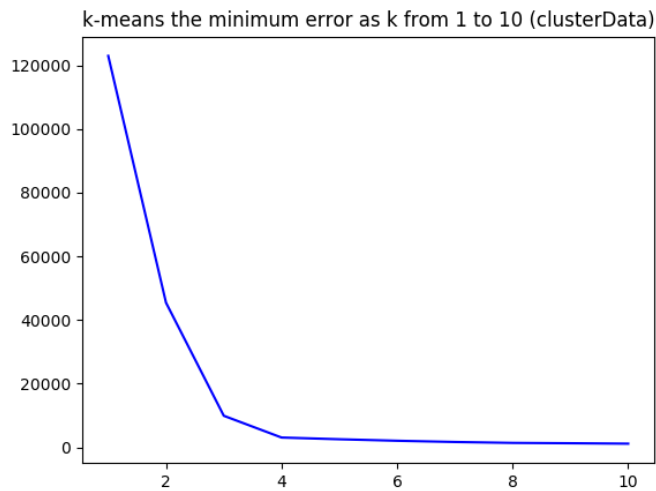
error.png

3. If we just running k-means once, we potentially get a terrible clustering, with kMeansError stabilizes at 9370 (worser case is possible).However, when we run k-means 50 times, we can ensure a relatively small kMeansError about 3000, a relatively meaningful clustering.

Here is the "best" clustering we obtained after running k-means 50 times.

## 2.2 Selecting k in k-means

1. Because the kMeansError decreases strictly monodically along with k increasing, kMeansError will lead us to choose k as big as the size of dataset but sacrificing running complexity ( When k =n, the kMeansError =0)

2.Violate golden rule.
3. The plot graph is:

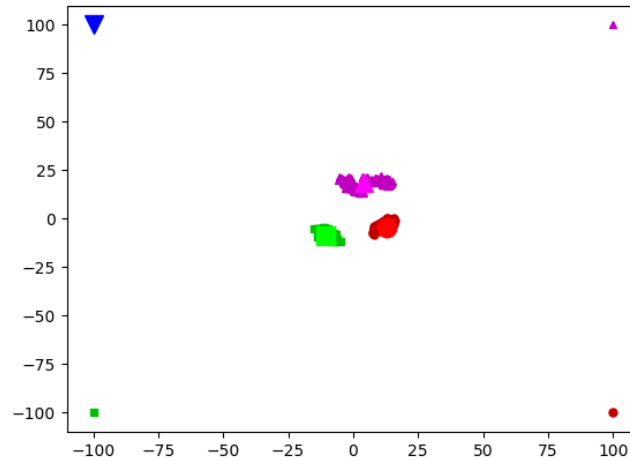

k-means the minimum error as k from 1 to 10 (clusterData)

check code at example_kMeans2.jl
4. I would choose 4 in this case. There is tradeoff between run time and accuracy. Although, the sharpest "elbow" occurs at k =2, the minimum error a t k =2 is too high. But when we sacrifice runtime a little bit, I say a little bit because increasing k from 2 to 4 does not make a large difference on runtime ,but almost decreasing minimum error by 8 times, therefore there is huge difference
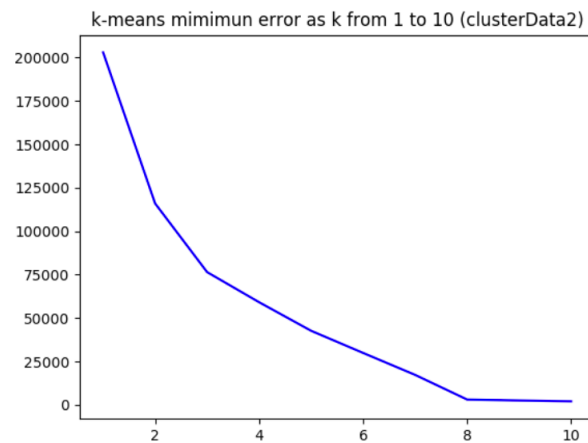
## 2.3 K-Medians

1.



check code at example_kMedians.jl
The clustering2Dplot function gives a bad clustering because it cannot detect outliers, so has a high error.

2. I would choose k = 8 in this case.As we can see the picture down blew that. Also by choosing k =8 we can assign 4 outliers to a separate cluster ensuring the main data is not influenced by them.



k-means mimimun error as k from 1 to 10 (clusterData2)

3. check code at kMedians.j or here:

```julia
include("misc.jl")
include("clustering2Dplot.jl")

type PartitionModel
    predict # Function for clustering new points
    y # Cluster assignments
    W # Prototype points
end

function kMedian(X,k;doPlot=false)
    # K-means clustering
    (n,d) = size(X)

    # Choos random points to initialize means
    W = zeros(k,d)
    perm = randperm(n)

    for c = 1:k
        # initially pick points from X
        W[c,:] = X[perm[c],:]
    end

    # Initialize cluster assignment vector
    y = zeros(n)
    changes = n
    D=zeros(n,k)
    while changes != 0
        # Compute L1 norm distance between each point and each median
        for i in 1:k
            for j in 1:n
                D[j,i] = sum(abs.(X[j,:]−W[i,:]))
            end
        end

        # Assign each data point to closest mean (track number of changes labels)
```

```julia
35        # Assign each data point to closest mean (track number of changes labels)
36        changes = 0
37        for i in 1:n
38            #y_new is the nearest cluster
39            (~,y_new) = findmin(D[i,:])
40            changes += (y_new != y[i])
41            #y records corresponding cluster of every datapoint
42            y[i] = y_new
43        end
44
45        # Optionally visualize the algorithm steps
46        if doPlot && d == 2
47            clustering2Dplot(X,y,W)
48            sleep(.1)
49        end
50
51        # Find median of each cluster
52
53        for c in 1:k
54            W[c,:] = median(X[y.==c,:],1)
55        end
56
57        # Optionally visualize the algorithm steps
58        if doPlot && d == 2
59            clustering2Dplot(X,y,W)
60            sleep(.1)
61        end
62
63        #@printf("Running k-means, changes = %d\n",changes)
64        #@printf("Running k-means, kMeansError= %.3f\n",kMeansError(X,y,W))
65    end
66
67    function predict(Xhat)
68        (t,d) = size(Xhat)
69
```
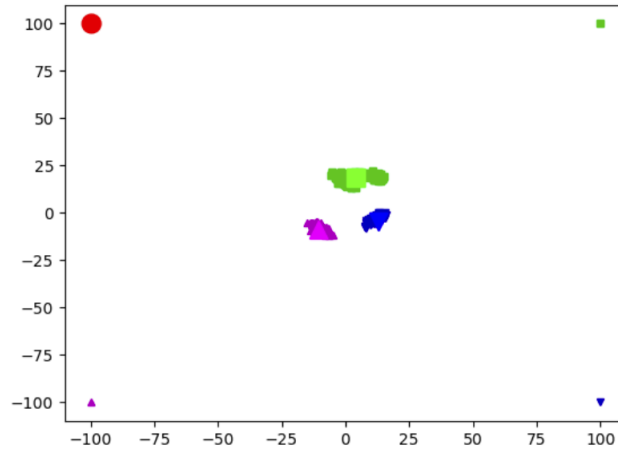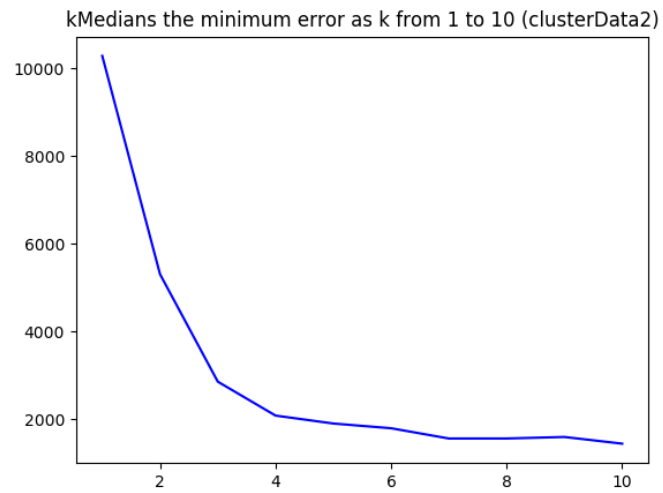
```julia
    function predict(Xhat)
        (t,d) = size(Xhat)

        for i in 1:k
            for j in 1:n
                D[j,i] = sum(abs.(X[j,:]-W[i,:]))
            end
        end

        yhat = zeros(Int64,t)
        for i in 1:t
            (~,yhat[i]) = findmin(D[i,:])
        end
        return yhat
    end

    return PartitionModel(predict,y,W)
end


function kMeansError(X,y,W)
    (n,d) = size(X)
    (k,d2) = size(W)
    assert(d == d2)
    temp =0
    for i in 1:n
        for j in 1:d
            #temp += abs.(X[i,j] - W[Int(y[i]),j])
            temp += (X[i,j] - W[Int(y[i]),j]).^2
        end
    end
    return temp
end
```

4.



kMedians the minimum error as k from 1 to 10 (clusterData2)

The graph has the sharpest change in slope between k =2 and k = 4, so the appropriate value for k could be 3.
It would give a relatively good clustering but cannot detect outliers

## 2.4 Very-Short Answer Questions

1.Nope, result depends on initial clusters
2.k=n, where n is the data size. when K = n, the distance is 0 no matter how.

3. clusters with different hierarchies

# 3 More Unsupervised Learning

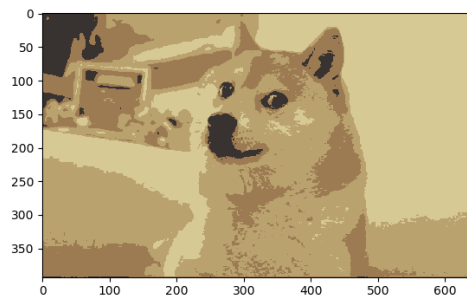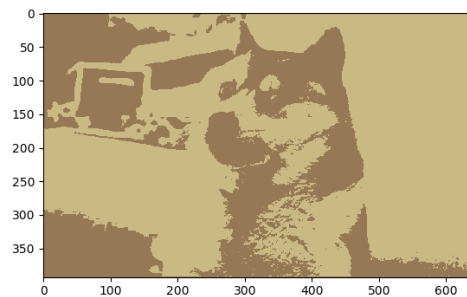## 3.1 Density-Based Clustering

1.
radius = 2, minPoints = 2
2.
radius = 4, minPoints = 2
3.
radius = 15, minPoints =2
4.
radius = 20, minPoints =2
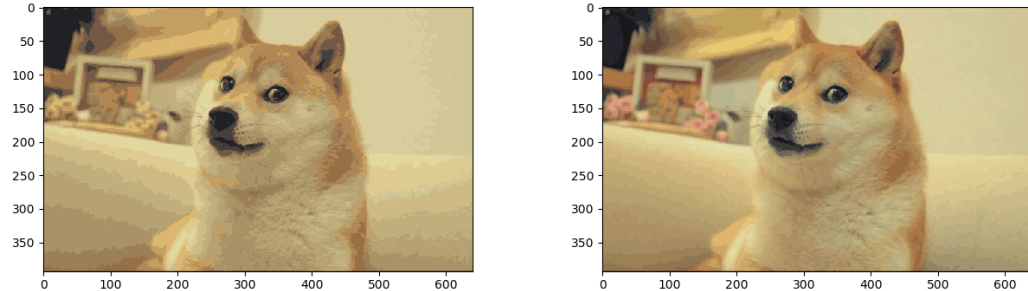
## 3.2 Vector Quantization

1. see at quantizeImage.jl and also can check here:

```julia
using PyPlot
include("kMeans.jl")
dog = imread("dog.png")

function quantizeImage(img,b)
    (nRows,nCols,a) = size(img)
    assert(a .== 3)
    ob = nRows * nCols
    X=reshape(img,ob, 3)

    model = kMeans(X,2.^b ,doPlot=false)
    y=model.predict(X)
    return deQuantizeImage(y,model.W,nRows,nCols)
end


function deQuantizeImage(y,W,nRows,nCols)
    y = reshape(y, nRows,nCols)
    biu = zeros(nRows,nCols,3)
    for i in 1:nRows
        for j in 1:nCols
            biu[i,j,:] = W[y[i,j],:]
        end
    end
    return biu
end

imshow(quantizeImage(dog,6))
```

2.
value =1 and value = 2



value = 4 and value =6

## 3.3 Very-Short Answer Questions

1.The cluster will be decided mainly by the distance of larger scale.

2.Advantage: We can label outliers in supervised model. Therefore it can predict similar outliers very well.

Disadvantage:If there is outliers in dataset which are different from training set, there it cannot detect this outlier and will behave horribly.

3. The cost of finding all rows i in X is going to be like: $O(n)$

The cost of runtime for finding all rows after given us a hash table that assigns rows of X to keys that divide the space into a 2D grid of squares with radius r, and we are using k to demote the maximum number of pointed to hashed to same key value: $O(k)$

# 4 Matrix Notation and Linear Regression

## 4.1 Converting to Matrix/Vector/Norm Notation

1. $\sum_{i=1}^{n} \|w^T x_i - y_i\| = \|Xw - y\|_1$

2. $\max_{1 \le i \le n} \|w^T x_i - y_i\| + \frac{\lambda}{2} \sum_{j=1}^{n} w_j^2 = \|Xw - y\|_\infty + \frac{\lambda}{2}\|w^2\|$

3. $\sum_{i=1}^{n} v_i(w^T x_i - y_i)^2 + \lambda \sum_{j=1}^{d} \|wj\| = v_i\|w^T x_i - y_i\|^2_2 + \lambda\|w_j\|_1$

## 4.2 Minimizing Quadratic Functions as Linear Systems

1.

f(w) $= \frac{1}{2}\|w - u\|^2$

f(w)$= \frac{1}{2}\|w\|^2 + \frac{1}{2}\|u\|^2 - w^T u$

13

$\nabla f(w) = w - u$

Thus, w = u

2.

f(w) $= \frac{1}{2}\|w\|^2 + w^T X^T y$

f(w) $= \frac{1}{2}w^T I w + w^T X^T y$

$\nabla f(w) = w + X^T y$

Thus, w $= -X^T y$

3.

f(w) $= \frac{1}{2}\|Xw - y\|^2 + \frac{1}{2}w^T \Lambda w$

f(w) $= \frac{1}{2}\|X\|^2\|w\|^2 + \frac{1}{2}\|y\|^2 - w^T X^T y + \frac{\Lambda}{2}\|w\|^2$

$\nabla f(w) = \|X\|^2 w - X^T y + \Lambda w = 0$

4.

f(w) $= \frac{1}{2}\sum\limits_{i=1}^{n} v_i(w^T x_i - y_i)^2$
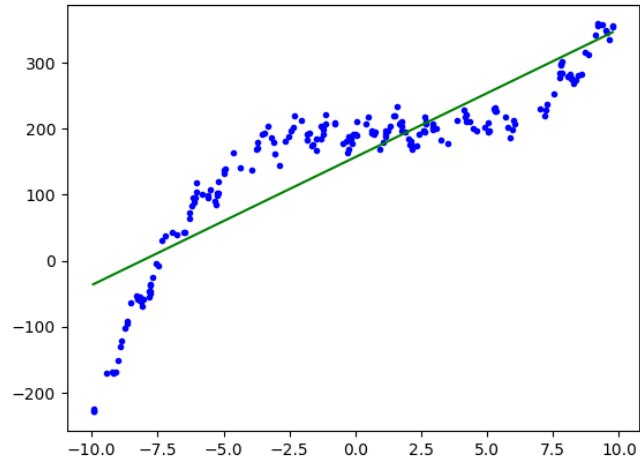
f(w) $= \frac{1}{2}(Xw - y)^T v_i(Xw - y)$

f(w) $= \frac{1}{2}w^T X^T v X w - w^T x^T v y + \frac{1}{2}y^T v y$

$\nabla f(w) = X^T v X w - X^T v y$

Thus, $X^T v X w = X^T v y$

## 4.3    Linear Regresion with Bias Variable

The graph plot is going to be like



The function code is:

```julia
2
3    function leastSquaresBias(X,y)
4
5        # Find regression weights minimizing squared error
6        (n,d) = size(X)
7        #w = (X'*X)\(X'*y)
8
9        X_0 = ones(n)
10       biaX = hcat(X_0 , X)
11       biaw = (biaX'*biaX)\(biaX'*y)
12
13       # Make linear prediction function
14
15       function predict(Xpredict)
16           (n2,) = size(Xpredict)
17           X_1 = ones(n2)
18           Xpredict = hcat(X_1,Xpredict)
19       return Xpredict * biaw
20   end
21
22
23       # Return model
24       return GenericModel(predict)
25   end
```

The test result is:
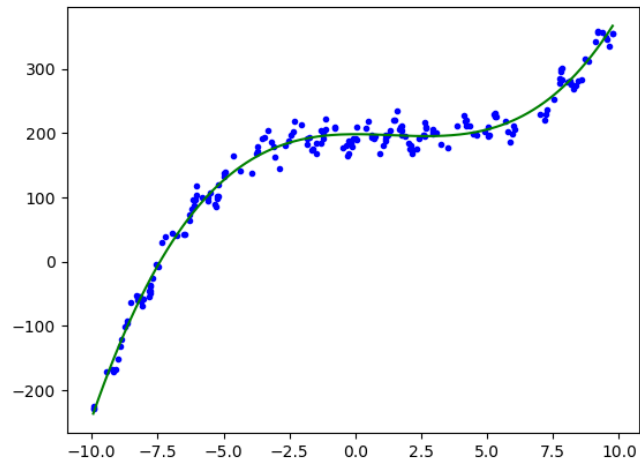
```
julia> include("example_nonLinear.jl")
Squared train Error with least squares: 3551.346
Squared test Error with least squares: 3393.869
1-element Array{Any,1}:
 PyObject <matplotlib.lines.Line2D object at 0x12d9e7ed0>
```

15

## 4.4 Linear Regression with Polynomial Basis

The plot graph is:



The code is:

```
27  function polyBasis(x,p)
28      (n,) = size(x)
29      Z = ones(n,p+1)
30      for i in 2:(p+1)
31          Z[:,i] = x.^(i-1)
32      end
33      return Z
34  end
35
36
37  function leastSquaresBasis(x,y,p)
38      Z = polyBasis(x,p)
39      w = (Z' * Z) \ (Z' * y)
40
41      function predict(Xhat)
42          bang = polyBasis(Xhat,p)
43          return bang * w
44      end
45
46      return GenericModel(predict)
47  end
48
```
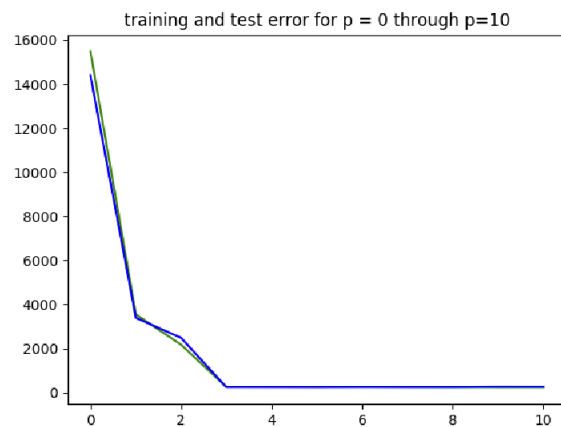
The test is:

```
6   # Fit a least squares model
7   include("leastSquares.jl")
8
9   trainError =[]
10  testError = []
11
12  for p in 0:10
13      model = leastSquaresBasis(X,y,p)
14
15      # Evaluate training error
16      yhat = model.predict(X)
17      push!(trainError, mean((yhat - y).^2))
18      #@printf("Squared train Error with least squares Basis: %.3f\n",trainError)
19
20      # Evaluate test error
21      yhat = model.predict(Xtest)
22      push!(testError, mean((yhat - ytest).^2))
23      #@printf("Squared test Error with least squares Basis: %.3f\n",mean((yhat - ytest).^2))
24  end
25
26  using PyPlot
27  figure()
28  title("training and test error for p = 0 through p=10")
29  plot(0:10,trainError,"g")
30  plot(0:10,testError,"b")
31  |
32
33  model = leastSquaresBasis(X,y,10)
34  # Plot model
35  figure()
36  plot(X,y,"b.")
37  Xhat = minimum(X):.1:maximum(X)
38  yhat = model.predict(Xhat)
39  plot(Xhat,yhat,"g")
```

The test graph from p=0 to p=10 is:



training and test error for p = 0 through p=10

As in the graph that we can see when p value increase, test error and train error decreasing and getting more stable and we also notice that test error and train error is almost same.
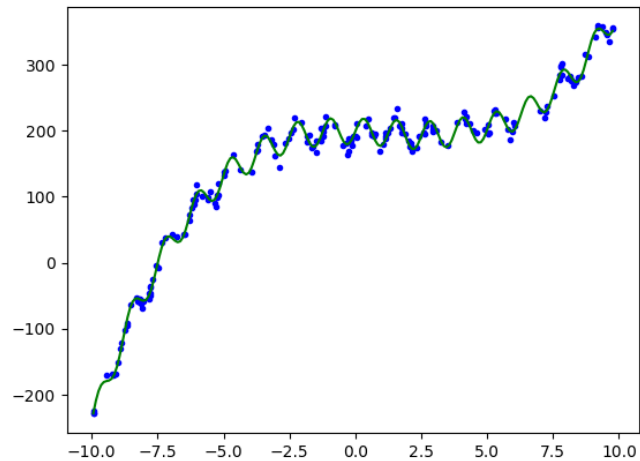
## 4.5 Manual Search for Optimal Basis

The code is like:

```julia
59  function manualSearch(X,p)
60      (n,) = size(X)
61      Z = ones(n,1+p*2)
62      for i in 1:p
63          Z[:,2*i] = X.^(i)
64          Z[:,2*i+1] =sin.(i*X)
65      end
66      return Z
67  end
68
69  function  leastSquaresManual(X,y,p)
70      Z = manualSearch(X,p)
71      w = (Z' * Z) \ (Z' * y)
72      function predict(Xhat)
73          bang = manualSearch(Xhat,p)
74          return bang * w
75      end
76
77      return GenericModel(predict)
78  end
79
```

The result is:

```
julia> include("example_manualSearch.jl")
Squared train Error with least squares Basis: 46.077
Squared test Error with least squares Basis: 50.865
1-element Array{Any,1}:
 PyObject <matplotlib.lines.Line2D object at 0x14e6af050>
```
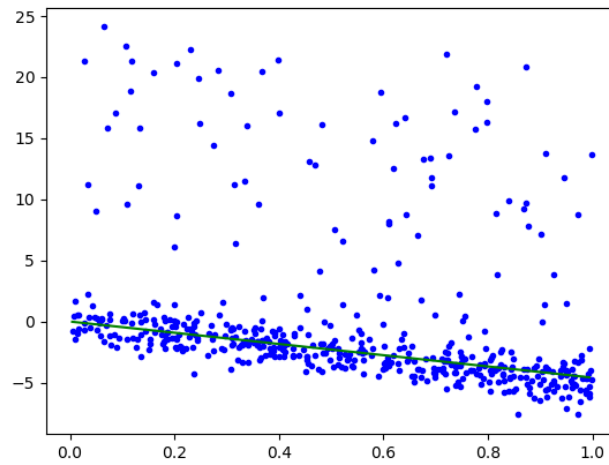
The plot graph is:

## 4.6   Very-Short Answer Questions

1. Because $y_i$ is continuous variable in this case

2.When there is a perfect linear relationship between two features, then least squares estimate is not going to be unique.

3.construct Z:$\mathcal{O}(np)$

calculate model(w):$n * p^2 + p^3$

predict: construct +predict: $np$ +$np$

overall=$\mathcal{O}(np^2 + p^3)$

4. When data hase piecewise linear relationship

# 5 Robust Regression and Gradient Descent

## 5.1 Weighted Least Squares in One Dimension

The plot graph is:



The code is :

```julia
function weightedLeastSquares(X,y,v)

    w = (X'*v*X)\(X'*v*y)
    predict(Xhat) = Xhat*w

    return GenericModel(predict)
end
```

The test is like:

```
1   # Load X and y variable
2   using JLD
3   data = load("outliersData.jld")
4   (X,y,Xtest,ytest) = (data["X"],data["y"],data["Xtest"],data["ytest"])
5
6   # Fit a least squares model
7   include("leastSquares.jl")
8   #model = weightedLeastSquares(X,y,0)
9   v1 = ones(400)
10  v2 = ones(100)
11
12  V = vcat(v1,0.1*v2)
13  V = Diagonal(V)
14  model = weightedLeastSquares(X,y,V)
15  # Evaluate training error
16  yhat = model.predict(X)
17  trainError = mean((yhat - y).^2)
18  @printf("Squared train Error with least squares: %.3f\n",trainError)
19
20  # Evaluate test error
21  yhat = model.predict(Xtest)
22  testError = mean((yhat - ytest).^2)
23  @printf("Squared test Error with least squares: %.3f\n",testError)
24  # Plot model
25  using PyPlot
26  figure()
27  plot(X,y,"b.")
28  Xhat = minimum(X):.01:maximum(X)
29  yhat = model.predict(Xhat)
30  plot(Xhat,yhat,"g")
```

## 5.2   Smooth Approximation to the L1-Norm

f(w) $= \sum\limits_{i=1}^{n} log(exp(w^T x_i - y_i) + exp(y_i - w^T x_i))$

we make value r like:

$r_i = w^T x_i - y_i$

First we convert function to matrix notation.

$f(w) = log(exp(Xw - y) + ex(y - Xw))$

$\frac{d}{dw} f(w) = \frac{\frac{d}{dw}(exp(Xw-y)+exp(y-Xw))}{exp(Xw-y)+exp(y-Xw)}$

$\frac{d}{dw} f(w) = \frac{(\frac{d}{dw}(Xw-y))*exp(Xw-y)+\frac{d}{dw}(Xw-y))*exp(y-Xw)}{exp(Xw-y)+exp(y-Xw)}$

$\frac{d}{dw} f(w) = \frac{X^T*exp(Xw-y)-X^T*exp(y-Xw)}{exp(Xw-y)+exp(y-Xw)}$

21

$$\frac{d}{dw}f(w) = X^T \frac{exp(Xw-y)-exp(y-Xw)}{exp(Xw-y)+exp(y-Xw)}$$

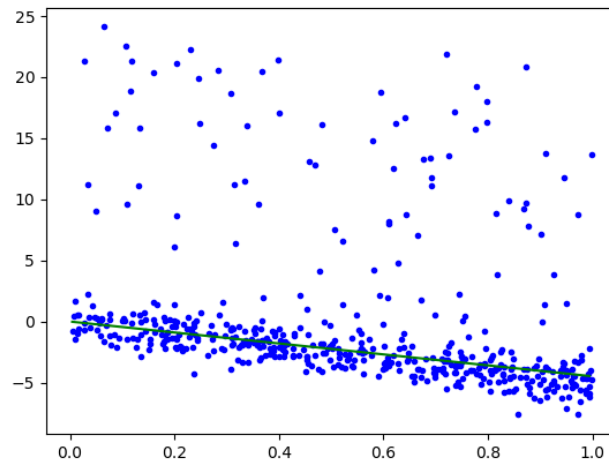## 5.3   Robust Regression

The code is like:

```
37   function robustRegressionObj(w,X,y)
38
39       a = exp.(X*w -y)
40       b = exp.(y-X*w)
41
42       f = sum(log.(a+b))
43
44       r = (a-b) ./ (a+b)
45
46       g = X' * r
47
48       return (f,g)
49   end
50
```

The dataset in graph is like:



## 5.4   Very Short Answer Questions

1. graph-based and distance-based Cluster-based cannot work because these outliers are disperse.

Model-based cannot work because we have too many outliers here and strongly influence our mean and standard deviatiol therefore the model we obtained is not accurate.

Graph-based and distance-based can detect outliers more effective in this dataset.

2.If there are lots of the features, then we need to use gradient descent to solve the problem.

3. Gradient is problem prone. Because model can be non-invertible or can not be differentible.