# CPSC340A3

Qinglan Huang n6v9a      Liu Yang v4d9

November 6, 2017

## 1  Convex Functions

1.$\frac{d}{dx}f'(w) = 2\alpha - \beta \geq 0$
Since the we can use second derivative and make the second derivative always greater or equal to zero for this entire domain therefore the function is always convex.

2.$\frac{d}{dx}f'(w) = \frac{1}{w} \geq 0$
Computing the second derivative for one-dimension function, the second derivative always greater or equal to zero for entire domain therefore the function is always convex.

3.$f(w) = \|Xw - y\|^2 + \lambda\|w\|$
$Since \|Xw - y\|^2$ is convex function since it is norm, and $\|w\|$ is convex since L1-regularized norm is tend to set variable to 0, sum of two convex function is convex.

4.Since $-y_i w^T x_i$ is linear function and scalar. Log(1+ $e^{-y_i w^T x_i}$) can be written as $\log(1+e^{n_i})$, and second derivative of linear function is always greater or equal to zero thus $n_i$ is convex.Second derivative of $\log(1+ e^{n_i})$ is convex, since $\frac{d}{dx}f'(w) = \frac{e^{n_i}}{(1+e^{n_i})^2}$ is always greater or equal to zero thus it is convex. Therefore the sum of the convex function is convex function as well.

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

Let $g(x) = \log(1 + \exp(x))$

$$g'(x) = \frac{(1 + \exp(x))'}{1 + \exp(x)} = \frac{\exp(x)}{1 + \exp(x)}$$

$$g''(x) = \frac{\exp(x)(1 + \exp(x)) - \exp(x)(1 + \exp(x))'}{(1 + \exp(x))^2}$$

$$= \frac{\exp(x)(1 + \exp(x)) - \exp(x)\exp(x)}{(1 + \exp(x))^2}$$

$$= \frac{\exp(x)}{(\exp(x) + 1)^2} \quad \text{for all } x.$$

math is here:

5.Since $w_0 - w^T x_i$ is linear function and convex because of the second derivative is greater or equal to 0, then the max of two convex function is convex and sum of the convex is convex. And $\|w\|^2$ is convex, thus the f(w) is convex which is sum of two convex function.
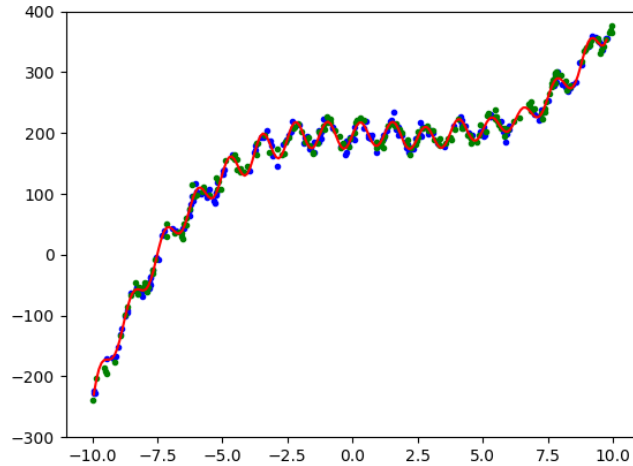
## 2 Gaussian RBFs and Regularization

### 2.1 Regularization

```julia
65
66    function leastSquaresRBF(X,y,sigma,lambda)
67        (n,d) = size(X)
68
69        Z = rbf(X,X,sigma)
70
71
72        v=ones(n)
73        v=Diagonal(v)
74
75        w = (Z'*Z + lambda * v)\(Z'*y)
76
77        predict(Xhat) = rbf(Xhat,X,sigma)*w
78
79        return LinearModel(predict,w)
80    end
81
82    function rbf(Xhat,X,sigma)
83        (t,d) = size(Xhat)
84        n = size(X,1)
85        D = distancesSquared(Xhat,X)
86        return (1/sqrt(2pi*sigma^2))exp.(-D/(2sigma^2))
87    end
```

## 2.2 Cross-Validation

```julia
using JLD
data = load("basisData.jld")
(X,y,Xtest,ytest) = (data["X"],data["y"],data["Xtest"],data["ytest"])

# Data is sorted, so *randomly* split into train and validation:
n = size(X,1)
perm = randperm(n)
#validStart = Int64(n/2+1) # Start of validation indices
#validEnd = Int64(n) # End of validation incides
#validNdx = perm[validStart:validEnd] # Indices of validation examples
#trainNdx = perm[setdiff(1:n,validNdx)] # Indices of training examples

stop = floor(Int,(n / 10))
println(stop)

validIndex = Matrix{Int64}(10, stop)
trainIndex = Matrix{Int64}(10,n-stop)

# Find best value of RBF variance parameter,
#   training on the train set and validating on the test set
include("leastSquares.jl")
minErr = Inf
bestSigma = []

for i in 1: 10
    validEnd = stop * i
    validStart = validEnd - stop + 1
    validIndex[i,:] = perm[validStart:validEnd]
    trainIndex[i,:] = perm[setdiff(1:n,validStart:validEnd)]
end


    for sigma in 2.0.^(-15:15)
        sumErr = 0
        for i in 1:10
            Xtrain = X[trainIndex[i,:],:]
            ytrain = y[trainIndex[i,:],:]
            Xvalid = X[validIndex[i,:],:]
            yvalid = y[validIndex[i,:],:]
            # Compute the error on the validation set
            model = leastSquaresRBF(Xtrain,ytrain,sigma,10^-12.0)
            yhat = model.predict(Xvalid)
            validError = sum((yhat - yvalid).^2)/(9*n/10)
            #@printf("With sigma = %.3f, validError = %.2f\n",sigma,validError)
            # Keep track of the lowest validation error
            sumErr += validError
        end

        avgErr = sumErr/10
        if avgErr < minErr
            minErr = avgErr
            bestSigma = sigma
        end
    end
```

4

For plot picture is like:

## 2.3 Cost of Non-Parametric Bases

1.Computing distance between the $x_i$ and $x_j$ cost $O(n^2d)$.Then computing the function $(Z^T Z + \lambda)w = Z^T y$ which cost $O(n^3)$, sum of two cost is $O(n^2d + n^3)$

2.First,Compute the distance between t objects and $x_i$ with d features cost O(tnd).Then, Computing prediction with the cost with Z*w with txn matrix from Z and nx1 matrix from w, we can get tx1 matrix, for one example cost only O(tn). Therefore the cost of classifying t new example with this model cost O(tnd).

3.
Runtime for building RBFs: $O(n^3 + n^2d)$
Runtime for building linear basis:$O(nd^2 + d^3)$.
math:
$n^3 + n^2d > nd^2 + d^3$
$n^2 > d^2$
Thus, when $n > d$ we choose RBFs than linear basis

4.
Runtime for testing RBFs: $O(tnd)$
Runtime for testing linear basis: $O(td)$
Thus, linear basis is always cheaper than RBFs.

# 3 Logistic Regression with Sparse Regularization

## 3.1 Logistic Regression

According to data report that we know with the condition that the number of features is the same. TrainError is decreasing and the validError is decreasing.

## 3.2 L2-Regularization

```
31
32   function logRegL2(X,y,lambda)
33       (n, d) =size(X)
34
35       #Initial guess
36       w = zeros(d,1)
37
38       #Function we'ra going to minimize(and that computes gradient)
39       funObj(w) = logisticObjL2(w,X,y)
40
41       # Solve least squares problem
42       w = findMin(funObj,w,derivativeCheck=true)
43       # Make linear prediction function
44       predict(Xhat) = sign.(Xhat*w)
45
46       # Return model
47       return LinearModel(predict,w)
48   end
49
50
51   function logisticObjL2(w,X,y)
52       yXw = y.*(X*w)
53       f = sum(log.(1 + exp.(-yXw))) + lambda/2 * norm(w)^2
54       g = -X'*(y./(1+exp.(yXw))) + lambda* w
55       return (f,g)
56   end
```

Training error increases to 0.002 and validation error decreasing to 0.074 with the same number of the feature, there are 101 features in the L2-regularization model. And the largest gradient descent iteration is going to be t=O(max(eig(X'X+$\lambda$))/min(eig(X'X+$\lambda$))) = 6

## 3.3  L1-Regularization

```
57
58   function logRegL1(X,y,lambda)
59       (n, d) =size(X)
60
61       #Initial guess
62       w = zeros(d,1)
63
64       #Function we'ra going to minimize(and that computes gradient)
65       funObj(w) = logisticObjL1(w,X,y)
66
67       # Solve least squares problem
68       w = findMinL1(funObj,w,lambda)
69       # Make linear prediction function
70       predict(Xhat) = sign.(Xhat*w)
71
72       # Return model
73       return LinearModel(predict,w)
74   end
75
76
77   function logisticObjL1(w,X,y)
78       yXw = y.*(X*w)
79       f = sum(log.(1 + exp.(-yXw)))
80       g = -X'*(y./(1+exp.(yXw)))
81       return (f,g)
82   end
83
```

Training error decrease to 0 and validation error decreasing to 0.052, the number of the feature in L1-regularization model is 71.

### 3.4 L0-Regularization

```
116
117            for j in setdiff(1:d,S)
118                # Fit the model with 'j' added to the feature set 'S'
119                # then compute the score and update 'minScore' and 'minS'
120                Sj = [S;j]
121                Xs = X[:,Sj]
122                w = zeros(length(Sj), 1)
123                w = findMin(funObj,w,verbose= false)
124                (f,~) = funObj(w)
125                score = f + lambda*length(Sj)
126
127                if score < minScore
128                    minScore = score
129                    minS = Sj
130                end
131                # PUT YOUR CODE HERE
132            end
133            S = minS
134        end
135
```

Training error maintain to 0 and validation error decreasing to 0.018,the number of the feature in L0-regularization model is 24

# 4 Logistic Regression with Sparse Regularization

1. Because validation error will try to eliminate noise in order to achieve a better score but noise should be part of true model. In one word, validation error tends to choose a larger parameter or a more complicated model.Also BIC can find a true model as n goes to infinity large.

2. Exhaustive search has $2^d$ models and is hard to find an efficient score to choose a best model (optimization bias is huge if we choose validation error). Therefor we use forward selection which substantially decreases number of possible models. Also, forward selection reduces false positives.

3. When $\lambda$ is small, model we get will be complicate and train error is small but approximation error is large.
When $\lambda$ is large or, model is simple and therefore large train error and small approximation error.

4. When there is bunch of irrelevant features, L1-regularization is preferred because it is robust to irrelevant features and can do feature selection. L2-regularization is preferred over L1 because it is differentiable and therefore can be computed easily, also L2 has a unique solution.

5. Generate several bootstrap samples, run feature selection using L1 as score, record selected features, find the union of these selected features.

6. Least squares penalize for being too right

7. SVM will found a classifier with largest margin, but perceptron will be satisfied as long as a classifier is found.

8. When d is large, polynomial basis ($Z$) would be too large to be stored. Kernel trick avoids storing Z but stores $K = ZZ^T$ because K is what really needed to make prediction and K's size is independent of d.

9. The perceptron algorithm, SVMs,logistic regression.