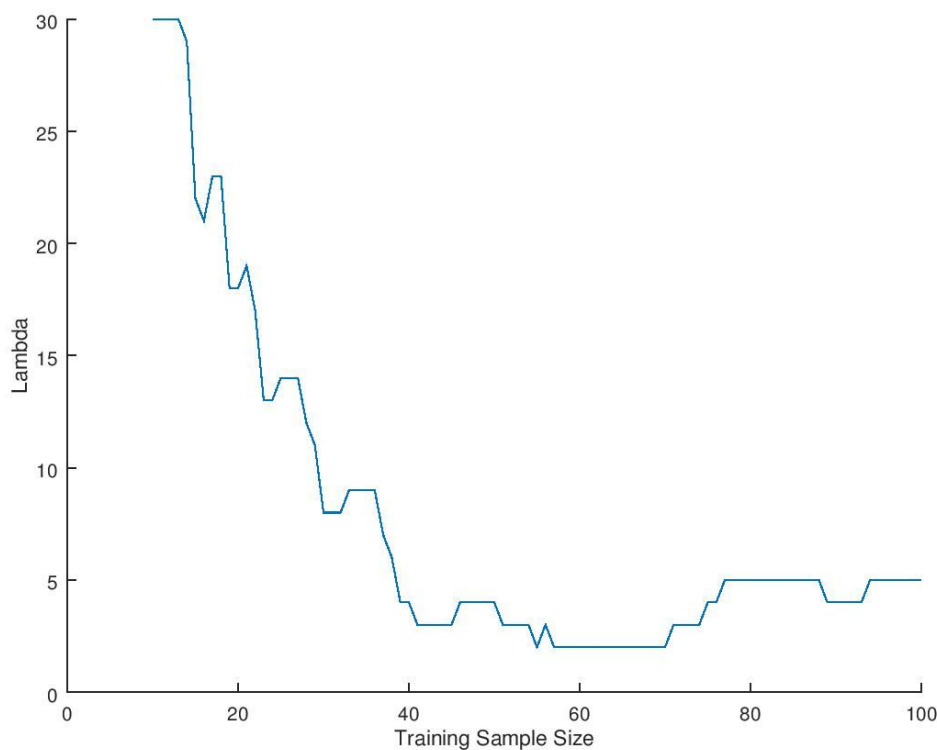


### Assignment 3

**Question 1.** Implement the ridge-regression algorithm. **no need to submit your code.** Run the algorithm on the dataset `regdata.mat` provided on the course web page. Run the regression using  $\lambda \in \{0, 1, 2, \dots, 30\}$  on the training set  $X, Y$  provided in the data file. Try training-set sizes between 10 and 100. For each training set size that you try, find the value of  $\lambda$  that obtains the smallest mean-squared-error (the average squared loss) on the test set provided in the data file.

- (a) (10pts) Submit a plot of the value of  $\lambda$  that minimizes the mean squared error on the test set as a function of the training set size  $m$ .



- (b) (10pts) What trend to you expect in the plot based on what we learned in class? Explain.

Lambda is the parameter representing the trade off between the norm of  $W$  and the loss function. The higher lambda is the more we are encouraged to pick  $W$  with smaller norm.

We would expect that the more samples we have the less the norm should be taken into consideration and a smaller lambda should be used. A wrong predictor, meaning a predictor that gives wrong weight to features, will eventually “pay” on larger training set. This is because in high probability there will be one or more training samples that will “punish” him. Therefore we don’t really need to consider the norm.

In the other hand, if we have a smaller training set, we might pick a predictor that supplied good results to the training set, but from the wrong intentions (wrong features) - the classic overfitting. Punish a predictor for its norm will encourage him to take into consideration only the right features.

(c) (5pts) Did you get this trend in the plot you submitted? If there are any differences, explain why they could occur.

Yes, this graph is decreasing, not monotonically but obviously decreasing.

After about 40 samples the required lambda is low(2-5). The results supports the theory we learned in class - the more samples we used a smaller lambda was needed to get smaller error.

There isn't any difference beside the noise in the end, it might have happened because all the lambdas in the region 2-5 gave small error rate.

(d) (15pts) In this data set, the label  $y$  of each example  $x$  was generated by setting  $y = \langle w, x \rangle + \eta$ , where  $w$  is a fixed vector which is the same for all examples in the data set, and  $\eta$  is a standard Gaussian random variable,  $\eta \sim N(0, \sigma)$  for some  $\sigma > 0$ .  $\eta$  is drawn independently for each example in the data set. What is the Bayes-optimal predictor for this problem with respect to the squared loss? And how about the absolute loss? Prove your claims.

squared loss:

$h^*(x) = E[Y|X=x]$  (We learn in class that Bayes-optimal predictor for squared loss is the **expected value**)

$$= E(\langle w, x \rangle + \eta | X = x) = E(\langle w, x \rangle | X = x) + E(\eta)$$

Since  $\langle w, x \rangle$  is constant and  $\eta \sim N(\mu = 0, \sigma) \Rightarrow E[\eta] = 0$

So:  $h^*(x) = \langle w, x \rangle$  for this problem with respect to the squared loss.

absolute loss:

$h^*(x) = \text{MEDIAN}[Y|X=x]$  (We learn in class that Bayes-optimal predictor for absolute loss is the **MEDIAN value**).

$$= M(\langle w, x \rangle + \eta) \text{ (The operator } M \text{ is homogeneous and keeps shifting)} = M(\langle w, x \rangle) + M(\eta)$$

Since  $\langle w, x \rangle$  is constant and  $\eta \sim N(\mu = 0, \sigma) \Rightarrow M[\eta] = \mu = 0$ , by the properties of Normal distribution.

So:  $h^*(x) = \langle w, x \rangle$  for this problem with respect to the absolute loss.

**Question 2.** (10pts) We saw in class that the LASSO regression algorithm minimizes the following objective function:

$$\text{Minimize}_{w \in \mathbb{R}^d} \lambda \|w\|_1 + \sum_{i=1}^m (\langle w, x_i \rangle - y)^2.$$

Write a quadratic program that solves this minimization program.

**Answers 2.** We will show how to solve LASSO regression using quadratic program.

First we will recall the settings of quadratic program:

$$\text{Minimize}_{z \in \mathbb{R}^n} \frac{1}{2} z^T * H * z + \langle u, z \rangle \quad \text{subject to} \quad Az \geq v$$

We will define:

$$A_{2d+2m \times 2d+m} = \begin{bmatrix} Id_{d \times d} & 0_{d \times m} & -Id_{d \times d} \\ -Id_{d \times d} & 0_{d \times m} & Id_{d \times d} \\ X_{m \times d} & -Id_{m \times m} & 0_{m \times d} \\ -X_{m \times d} & Id_{m \times m} & 0_{m \times d} \end{bmatrix}$$

$$v = \begin{bmatrix} 0_{2d \times 1} \\ Y_{m \times 1} \\ -Y_{m \times 1} \end{bmatrix}$$

The way A and v are defined forces z to look like:

$$z = [w_{d \times 1}; (X * w - Y)_{m \times 1}; |w|_{d \times 1}]$$

Explanation:

$w_{d \times 1}$  - d “free” coordinates, the w we are looking for in LASSO.

$(X * W - Y)_{m \times 1}$  - notice that for each coordinate  $i$   $d < i \leq d + m$  we force

$$z_i \geq x_{i-d} * w - y_{i-d} \wedge z_i \leq x_{i-d} * w - y_{i-d} \rightarrow z_i = x_{i-d} * w - y_{i-d}.$$

$|w|_{d \times 1}$  - for the last d coordinates we forced  $z_{d+m+i} \geq z_i \wedge z_{d+m+i} \leq -z_{d+m+i} \Rightarrow z_{d+m+i} = |z_i|$ .

Since this is a minimization problem the vector z that achieve the minimum has to apply equality.

We will now show how to define H and u:

$$H_{2d+2m \times 2d+m} = \text{if } (d < i, j \leq d + m) : H_{ij} = 2, \text{ else } : 0$$

$$u_{2d+2m \times 1} = [0_{d+m \times 1}; \lambda_{d \times 1}]$$

We will notice now:

$$\frac{1}{2} z^T * H * z = \frac{1}{2} [w_{d \times 1}; (X * w - Y)_{m \times 1}; |w|_{d \times 1}]^T * H * [w_{d \times 1}; (X * w - Y)_{m \times 1}; |w|_{d \times 1}] =$$

$$[0_{d \times 1}; (X * w - Y)_{m \times 1}; 0_{d \times 1}]^T * [w_{d \times 1}; (X * W - Y)_{m \times 1}; |w|_{d \times 1}] = z = (X * W - Y)^2$$

$$\langle u, z \rangle = \langle [0_{d+m \times 1}; \lambda_{d \times 1}], [w_{d \times 1}; (X * w - Y)_{m \times 1}; |w|_{d \times 1}] \rangle = |w|$$

Therefore we managed to transform the LASSO problem to quadratic program problem. Given z, the output of the quadratic program solver, the first d coordinates of z are the vector w we are looking for in LASSO.

**Question 3.** Let  $\mathcal{X}$  be the set of all the undirected graphs over  $n$  vertices and let  $\mathcal{Y} = \{0, 1\}$ .

For a graph  $x \in \mathcal{X}$ , define the mapping  $g : \mathcal{X} \rightarrow \mathbb{N}^n$ , where coordinate  $i$  in the vector  $g(x)$  is the number of  $i$ -cliques in the graph  $x$ .

Let  $\mathcal{H}$  be a hypothesis class which includes all the functions  $h : \mathcal{X} \rightarrow \mathcal{Y}$  that have the form:  $h(x) \equiv \mathbb{I}[g(x) = v]$ , for some vector  $v \in \mathbb{N}^n$ .

Suppose that  $\mathcal{D}$  is a distribution over  $\mathcal{X} \times \mathcal{Y}$ , and suppose that in this distribution, the label of a graph is completely determined by the numbers of  $i$ -cliques that the graph has for the various  $i$ .

- (a) (10pts) Use the PAC-learning upper bounds that we showed in class to show that the sample complexity of learning  $\mathcal{H}$  is  $O(n^2 \log(n))$ .

The PAC learning upper bound claims that any ERM algorithm with training sample size  $m \geq \frac{2 \log(|H|) + 2 \log(2/\delta)}{\epsilon^2}$  gets an error of at most  $\epsilon$ , with a probability of at least  $1 - \delta$  over the random training samples ( $\epsilon, \delta \in (0, 1)$ ). Given  $\epsilon, \delta$ , we find out that the sample size needed is at least  $\frac{2}{\epsilon^2} \log(|H|) + \frac{2 \log(2/\delta)}{\epsilon^2}$ . Since  $\frac{2}{\epsilon^2}, \frac{2 \log(2/\delta)}{\epsilon^2}$  are constant values, it terms of big O the sample complexity depends on  $\log(|H|)$ . We will now show an upper bound for the size of  $H$ . If we could find the number of different outputs of  $g(x), x \in \mathcal{X}$ , we could find the number of different predictors. The  $i$  coordinate in  $g(x)$  is the number of  $i$ -clique in the graph  $x$ . The number of  $i$ -clique possible is at most  $\binom{n}{i}$ . Since  $i$  is between 1 to  $n$ , the coordinate with the highest value is for  $i = \frac{n}{2}$ , for this  $i$  there are at most  $\frac{n!}{(\frac{n}{2})! * (\frac{n}{2})!} \leq n^{\frac{n}{2}}$ . Because there are  $n$  coordinates and each coordinate is a natural number which isn't greater than  $n^{\frac{n}{2}}$ , the number of different  $g(x)$  is at most  $(n^{\frac{n}{2}})^n$ , hence  $|H| = (n^{\frac{n}{2}})^n = (n^{\frac{n^2}{2}})$

This means, the sample complexity is  $O(\log(|H|)) = O(\log(n^{\frac{n^2}{2}})) = O(n^2 * \log(n^{\frac{n}{2}})) = O(n^2 * \log(n))$ , as required.

- (b) (10pts) Adam and Ronnie want to learn a classifier that is guaranteed a low error with a probability of at least 95%. Ronnie used a training set of size  $m$  drawn from  $\mathcal{D}$  to learn a classifier from  $\mathcal{H}$ . How many examples should Adam feed his classifier, so that his classifier has an error guarantee which is half the error guarantee of Ronnie's? Explain.

Lets define  $\epsilon$  the error Ronnie gets with here training sample  $m$ . From PAC upper bound in probability of 0.95,

$m \geq \frac{2(\log(|H|) * \log(\frac{2}{0.95}))}{\epsilon^2} \Rightarrow \epsilon \geq \sqrt{\frac{2(\log(|H|) * \log(40))}{m}}$ . Adam wants his error to be half the error Ronnie has. Therefore Adam wants to get with probability of 0.95 an error of  $\frac{\epsilon}{2}$ . From PAC this means  $m'$ , the sample complexity of Adam, is:

$m' \geq \frac{2(\log(|H|) * \log(\frac{2}{0.95}))}{(\frac{\epsilon}{2})^2} = 4 * \frac{2(\log(|H|) * \log(\frac{2}{0.95}))}{\epsilon^2}$ . Since epsilon is at least  $\sqrt{\frac{2(\log(|H|) * \log(40))}{m}}$ , we get:

$4 * \frac{2(\log(|H|) * \log(\frac{2}{0.95}))}{(\frac{\epsilon}{2})^2} \leq 4 * \frac{2(\log(|H|) * \log(\frac{2}{0.95}))}{\frac{2(\log(|H|) * \log(40))}{m}} = 4 * m$ . Therefore if  $m' \geq 4m$  in probability of 0.95 Adam's error guaranteed to be half the error guaranteed for Ronnie.

**Question 4.** Consider the following optimization objective for a regression problem, where  $\mathcal{X} = \mathbb{R}^d$  and  $\mathcal{Y} = \mathbb{R}$ :

$$\text{Minimize}_{w \in \mathbb{R}^d} \lambda \|w\|_2^4 + \sum_{i=1}^m \exp(|\langle w, x_i \rangle - y_i|)$$

(a) (7pts) Does the Representer Theorem apply for this objective function? Prove your claim.

**Answers 4.a.** We will prove that this objective function, the Representer Theorem applies. Given a feature map  $\psi : \mathcal{X} \rightarrow \mathcal{H}$ , the representer theorem claims that if the objective can be represented as:

*Minimize*  $f(\langle w, \psi(x_1) \rangle, \dots, \langle w, \psi(x_m) \rangle) + R(\|w\|)$  s.t.  $f : \mathcal{H}^m \rightarrow \mathbb{R}$  and  $R : \mathbb{R}_+ \rightarrow \mathbb{R}$  monotonic non decreasing,

Then the Representer Theorem applies and the optimal solution to the objective can be written as:

$$w = \sum_{i=1}^m \alpha_i \psi(x_i) \text{ where } \alpha \in \mathbb{R}^m.$$

Therefore, we need to show that there exists  $f$  and  $R$  that satisfy the conditions and

$$\text{Minimize}_{w \in \mathcal{H}} \lambda \|w\|_2^4 + \sum_{i=1}^m \exp(|\langle w, \psi(x_i) \rangle - y_i|) = \text{Minimize}_{w \in \mathcal{H}} f(\langle w, \psi(x_1) \rangle, \dots, \langle w, \psi(x_m) \rangle) + R(\|w\|).$$

Lets define  $f(t_1, \dots, t_m) = \sum_{i=1}^m \exp(|t_i - y_i|)$ ,  $R(t) = \lambda t^4$ . we will now tonic that:

$$\text{Minimize}_{w \in \mathcal{H}} f(\langle w, \psi(x_1) \rangle, \dots, \langle w, \psi(x_m) \rangle) + R(\|w\|) = \text{Minimize}_{w \in \mathcal{H}} \lambda \|w\|_2^4 + \sum_{i=1}^m \exp(|\langle w, \psi(x_i) \rangle - y_i|).$$

We just need to stasify the conditions  $f : \mathcal{H}^m \rightarrow \mathbb{R}$  and  $R : \mathbb{R}_+ \rightarrow \mathbb{R}$  monotonic non decreasing.  $f$  is a function that gets  $m$  real numbers and returns a single real number, hence  $f : \mathcal{H}^m \rightarrow \mathbb{R}$ .  $R$  is a function that for every input  $t$  returns a real number. We just left to show that  $R$  is monotonic non decreasing On  $\mathbb{R}_+$ .

Let  $x, y > 0, x \leq y$  we will show that  $f(x) \leq f(y)$  :

$$f(x) = \lambda x^4, \lambda \geq 0 \Rightarrow \lambda x^4 \leq \lambda y^4 = f(y) \Rightarrow f(x) \leq f(y).$$

**Answers 4.b.**

From the section above, we know that the following problem applies the Representer Theorem. Therefore, the optimal solution to the objective can be written as  $w = \sum_{i=1}^m \alpha_i \psi(x_i)$ ,  $\forall i, \alpha_i \in \mathbb{R}$ . Let  $w$  be the optimal solution for the problem,

therefore  $\lambda \|w\|_2^4 + \sum_{i=1}^m \exp(|\langle w, \psi(x_i) \rangle - y_i|)$  is minimal for  $w$  over all  $w' \in \mathcal{H}$ .

$$\lambda \left\| \sum_{i=1}^m \alpha_i \psi(x_i) \right\|_2^4 + \sum_{i=1}^m \exp(|\langle \sum_{j=1}^m \alpha_j \psi(x_j), \psi(x_i) \rangle - y_i|) \leq \lambda \|w\|_2^4 + \sum_{i=1}^m \exp(|\langle w, \psi(x_i) \rangle - y_i|)$$

Since  $w$  is optimal, it can be represented as  $w = \sum_{i=1}^m \alpha_i \psi(x_i)$ ,  $\forall i, \alpha_i \in \mathbb{R}$ :

$$\begin{aligned} \lambda \left\| \sum_{i=1}^m \alpha_i \psi(x_i) \right\|_2^4 + \sum_{i=1}^m \exp(|\langle \sum_{j=1}^m \alpha_j \psi(x_j), \psi(x_i) \rangle - y_i|) &= \lambda \left( \sum_{i=1}^m \alpha_i \psi(x_i), \sum_{i=1}^m \alpha_i \psi(x_i) \right)^2 + \sum_{i=1}^m \exp(|\langle \sum_{j=1}^m \alpha_j \psi(x_j), \psi(x_i) \rangle - y_i|) \\ &= \lambda \left( \sum_{i=1}^m \sum_{j=1}^m \alpha_i \alpha_j \langle \psi(x_i), \psi(x_j) \rangle \right)^2 + \sum_{i=1}^m \exp(|\langle \sum_{j=1}^m \alpha_j \psi(x_j), \psi(x_i) \rangle - y_i|) \\ &= \lambda \left( \sum_{i=1}^m \sum_{j=1}^m \alpha_i \alpha_j K(x_i, x_j) \right)^2 + \sum_{i=1}^m \exp(|\langle \sum_{j=1}^m \alpha_j K(x_j, x_i) \rangle - y_i|) \end{aligned}$$

Therefore, we want to find a vector  $a \in \mathbb{R}^m$  that minimize the expression:

$$\lambda \left( \sum_{i=1}^m \sum_{j=1}^m \alpha_i \alpha_j K(x_i, x_j) \right)^2 + \sum_{i=1}^m \exp(|\langle \sum_{j=1}^m \alpha_j K(x_j, x_i) \rangle - y_i|).$$

We will notice that all vectors in this expression are 'm' length.

**Question 5.** (10pts) Consider the following loss function for regression:

$$\ell(y, y') = (\ln(y) - \ln(y'))^2.$$

Calculate the Bayes-optimal predictor for this loss,  $\operatorname{argmin}_{h: \mathcal{X} \rightarrow \mathcal{Y}} \ell(h, \mathcal{D})$ . Prove your claim using the definition of a Bayes-optimal predictor.

**Answers 5.**      **Claim:**  $h^*(x) = e^{\ln[Y|X=x]}$

**Proof** (for discrete distributions):

$$E_{(X,Y) \sim D}[(\ln(h(X)) - \ln(Y))^2] = E_X[E_Y[(\ln(h(X)) - \ln(Y))^2 | X = x]]$$

$\Rightarrow$  Bayes-optimal minimizes  $E_Y[(\ln(h(X)) - \ln(Y))^2 | X = x]$  for each  $x \in X$ .

Let  $v_1, v_2, \dots$  be the possible values of  $Y$ , and Let  $p_i = P[Y = v_i | X = x]$ .

**Then:**

$$h^*(x) = \operatorname{argmin}_b \sum_{i=1}^{\infty} p_i (\ln(v_i) - \ln(b))^2, \text{ of } b.$$

**Observe**  $f(b) = \sum_{i=1}^{\infty} p_i (\ln(v_i) - \ln(b))^2$ .

In order to find the **minimum** of the function we will have to derivate:

$$f'(b) = \sum_{i=1}^{\infty} 2p_i (\ln(v_i) - \ln(b)) \cdot \frac{1}{b} = \frac{1}{b} \cdot 2 \cdot \sum_{i=1}^{\infty} p_i (\ln(v_i) - \ln(b))$$

$$f'(b) = 0 \Rightarrow \frac{2}{b} \sum_{i=1}^{\infty} p_i (\ln(v_i) - \ln(b)) = 0 \Rightarrow \sum_{i=1}^{\infty} p_i (\ln(v_i) - \ln(b)) = 0$$

$$\Rightarrow \sum_{i=1}^{\infty} p_i \cdot \ln(v_i) = \ln(b) \cdot \sum_{i=1}^{\infty} p_i, \text{ (notice that } \sum_{i=1}^{\infty} p_i = 1) \Rightarrow \sum_{i=1}^{\infty} p_i \cdot \ln(v_i) = \ln(b)$$

$$\Rightarrow b = e^{\sum_{i=1}^{\infty} p_i \ln(v_i)}$$

So:  $h^*(x) = e^{E[\ln(Y) | X=x]}$