

Homework-4

= Given

Design matrix \mathbf{X}

Covariance matrix $\mathbf{C} = \frac{1}{n} \mathbf{X}^T \mathbf{X}$

Eigenvalues $\mathbf{V}_1, \mathbf{V}_2, \dots, \mathbf{V}_k$

$\tilde{\mathbf{x}}_i^o$ is the orthogonal projection of \mathbf{x}_i^o , column vectors representing a data point, orthogonal to \mathbf{V}_1

$$\tilde{\mathbf{x}}_i^o = (\mathbf{I} - \mathbf{V}_1 \mathbf{V}_1^T) \mathbf{x}_i^o$$

$\tilde{\mathbf{X}}_T = [\tilde{\mathbf{x}}_1^o, \dots, \tilde{\mathbf{x}}_n^o]$ deflated matrix, lies in the direction of the first principal eigenvectors

$$\tilde{\mathbf{X}}_T = (\mathbf{I} - \mathbf{V}_1 \mathbf{V}_1^T) \mathbf{X}^T$$

$$\mathbf{C} = \frac{1}{n} \mathbf{X}^T \mathbf{X} = \mathbf{V} \Lambda_n \mathbf{V}^T$$

Let's Suppose that the kernel matrix is $K = \mathbf{X} \mathbf{X}^T$

$$\text{So, } K = \mathbf{X} \mathbf{X}^T = \mathbf{U} \Lambda_m \mathbf{V}^T$$

The columns \mathbf{v}_i^o of the orthonormal matrix \mathbf{V} are the eigenvectors of \mathbf{C} and the columns \mathbf{u}_i^o of the orthonormal

matrix U are the eigenvectors of K . Considering an eigenvector-eigenvalue pair v, λ of K .

$$C = \frac{1}{n} (X^T U) = X^T X X^T U = X^T K U = \lambda X^T U$$

$X^T U, \lambda$ is an eigenvector-eigenvalue pair for C .

$$\|X^T U\|^2 = U^T X X^T U = \lambda$$

So the normalized eigenvectors of C is $v_i = \lambda^{-\frac{1}{2}} X^T U$.

$$v_i = \lambda^{-\frac{1}{2}} X^T U$$

Deflating X ,

$$\tilde{X} = X - U U' X = X - \lambda^{\frac{1}{2}} U v_i v_i' = X - X v_i v_i'$$

$$\tilde{C} = \tilde{X}^T \tilde{X} = \frac{1}{n} (X - U U^T X)^T (X - U U^T X)$$

$$= \frac{1}{n} (X^T X - X^T U U^T X)$$

$$\boxed{= \frac{1}{n} (X^T X - \lambda V V^T)}$$

D) $C v_j^o = \lambda_j^o v_j^o$

$$\tilde{C} = C - \lambda_j^o v_j^o v_j^{oT}$$

$$(C - \lambda_j^o v_j^o v_j^{oT}) v_j^o$$

$$= C v_j^o - \lambda_j^o v_j^o v_j^{oT} v_j^o$$

$$= \lambda_j^o v_j^o - \lambda_j^o v_j^o (v_j^{oT} v_j^o)$$

$$\text{if } j \neq 1, \text{ then } (C - \lambda_j^o v_j^o v_j^{oT}) v_j^o = \lambda_j^o v_j^o - \lambda_j^o v_j^o (0) = \lambda_j^o v_j^o$$

Thus $(C - \lambda_j^o v_j^o v_j^{oT})$ has the same eigenvectors as C and the same eigenvalues as C except that the largest one has been replaced by 0.

i) Hotelling deflation assumes that the largest eigenvalue λ_1 and an associated eigenvector $v(1)$ of C have been obtained from its deflation matrix \tilde{C} , which has the same eigenvalues $\lambda_2 \dots \lambda_n$ as C except that \tilde{C} has eigenvalue 0 with eigenvector $v(1)$ instead of eigenvalue λ_1 .

So if v be the first principal eigenvectors of \tilde{A}
 $v = v_2$ because the first eigenvector is zero

(d)

After k iterations,

$$C^k v_0 = \lambda^k \left[c_1 v_1 + c_2 \left(\frac{\lambda_2}{\lambda_1} \right)^k v_2 + \dots + c_n \left(\frac{\lambda_n}{\lambda_1} \right)^k v_n \right] [v_0 = c_1 v_1 + \dots + c_n v_n]$$

Since here λ_1 dominates, the ratio $\left(\frac{\lambda_i}{\lambda_1} \right)^k \rightarrow 0$ as $k \rightarrow \infty$
for all i . Therefore, $C^k v_0 = \lambda_1^k c_1 v_1$ gets better as k
grows. Here c_1, \dots, c_n is constant.

So, for large powers of k we will obtain a good
approximation of the dominant eigenvector.

For k iteration, the computational complexity is

$O(k^2)$.

In each iteration if we do $v_i^0 = \frac{C v_{i-1}^0}{\|C v_{i-1}^0\|_2}$,

this can be a better approach because the

Power method tends to produce approximations with large entries. In practice it is best to 'scale down' each approximation before proceeding to the next iteration. The way to accomplish this scaling is to determine the component of Cv_i that has the largest absolute value.

e) The computational complexity is $\propto kkm^2$

(See code)

problem 2

a) We have to set r_0 and c such that $\min_i \|q - x_i\| \leq r_0$ and $\|q - x'_i\| \leq cr_0$. More specifically, an r_0 -nearest query with constant factors c will return a j^* such that $\|q - x_j^*\| \leq cr_0$ if there is some i so that $\min_i \|q - x_i\| \leq r_0$.

we can do this by solving nearest neighbors checking r_0

with $r_0 = 1, 2, 4, 8, \dots, R$

we only need to give output if there is a point within r_0 , even if something within cr_0 .

b) Given,

$$h(x_i^0) = x_i^0[a]$$

For each point $x_i^0, x_j^0 \in \{0,1\}^d$

$$\Pr(h(x_i^0) = h(x_j^0)) \geq P_1$$

where $P_1 = 1 - \frac{P}{d} \times e^{-r/d}$ if $\text{dist}(x_i^0, x_j^0) \leq r$.

$$\Pr(h(x_i^0) = h(x_j^0)) \leq P_2$$

where $P_2 = 1 - \frac{cP}{d} \approx e^{-cr/d}$ if $\text{dist}(x_i^0, x_j^0) \geq cr$

Here $\boxed{P_1 > P_2}$

c) $\Pr(g(x_i^0) = g(x_j^0)) \geq P_1^K$ [lower bound] $[d(x_i^0, x_j^0) \leq r]$

$$\Pr(g(x_i^0) = g(x_j^0)) \leq P_2^K$$
 [upper bound] $[d(x_i^0, x_j^0) \geq cr]$

d) $\Pr(g_b(x_i^0) = g_b(x_j^0)) \geq 1 - \frac{1}{n}$ [LB] $[d(x_i^0, x_j^0) \leq r]$

$$\Pr(g_b(x_i^0) = g_b(x_j^0)) \leq \frac{1}{n}$$
 [UB] $[d(x_i^0, x_j^0) \geq cr]$

$$e) K = \frac{\ln(n)}{\ln(1/P_2)}$$

This implies that for fixed b , the expected number of x' that map to the same bucket as q is at most $n \times \frac{1}{n} = 1$.

Applying linearity of expectation over all d buckets, the expected number of false positive is at most d .

By Markov's inequality the probability there are more than $4d$ false positives is therefore at most $\frac{1}{4}$ and at least $\boxed{\frac{3}{4}}$.

According to first event, for any b

$$\Pr [g_b(x^*) \neq g_b(q)] \leq 1 - P_1^K = 1 - n^{-\rho} = 1 - \frac{1}{n^\rho} \quad [K = \frac{\ln(n)}{\ln(1/P_2)}]$$

$$\rho = \frac{\ln(P_1)}{\ln(P_2)}$$

$$\text{As } d = n^\rho,$$

$$\Pr [g_b(x^*) \neq g_b(q) \forall i] \leq \left(1 - \frac{1}{n^\rho}\right)^{n^\rho} \leq \frac{1}{e}$$

So, the first event holds with probability at least $\boxed{1 - \frac{1}{e}}$

With union bound the probability of both event holding at least:

$$1 - \frac{1}{4} - 1/e \geq \frac{1}{3}$$

2f) In expectation, there are $O(l)$ points in the data

Get which map to same hash functions as a for

Some K . we need to examine these points to

check if any of them within distance ϵr_0 from q

No, it is not guaranteed that there is a point with

distance at most ϵr_0 because the probability of

upper bound is very low (only $1/n$).

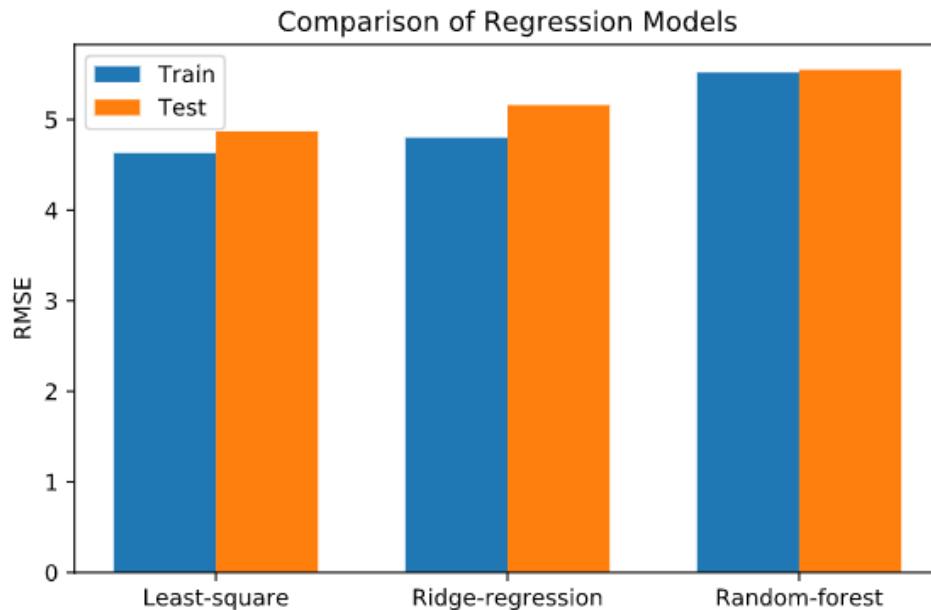
Problem 3(a)

Random Forest Boston Housing

Train RMSE: 5.515574182181881

Test RMSE: 5.551093717718717

Performance of RF comparing to least square and regression



From the above plot it is evident that both in train and test cases for Boston housing data Least-Square is performing best than random forest and ridge regression.

Problem 3(b)

Credit-g-dataset:

Train Accuracy: 0.7114285714285714

Test Accuracy: 0.7066666666666667

Breast-cancer dataset:

Train Accuracy: 0.8844221105527639

Test Accuracy: 0.9064327485380117

4(a) The computational complexity of optimizing a tree of depth d in terms of m and n is $O(mn \log n)$

4(b) Applications with high dimensional sparse output requires most expensive computation in GBDT training.

GBDT builds a regression tree that fits the residual from the previous tree. So at least $O(MN)$ time and memory are required to build GBDT trees.

To solve this problem, we can solve a L₀ regularized optimization problem to enforce the prediction of each leaf node in each tree to have only a small number of non-zero elements/labels.

Problem 4(c)

We can train an individual decision tree in parallel in threads.

Training a decision tree is an iterative process that requires looping over all possible ways of splitting the current leaf nodes into new leaf nodes. So each individual decision tree takes a lot of time for computation. So these decision trees can be run in parallel threads.

Problem 4(d)

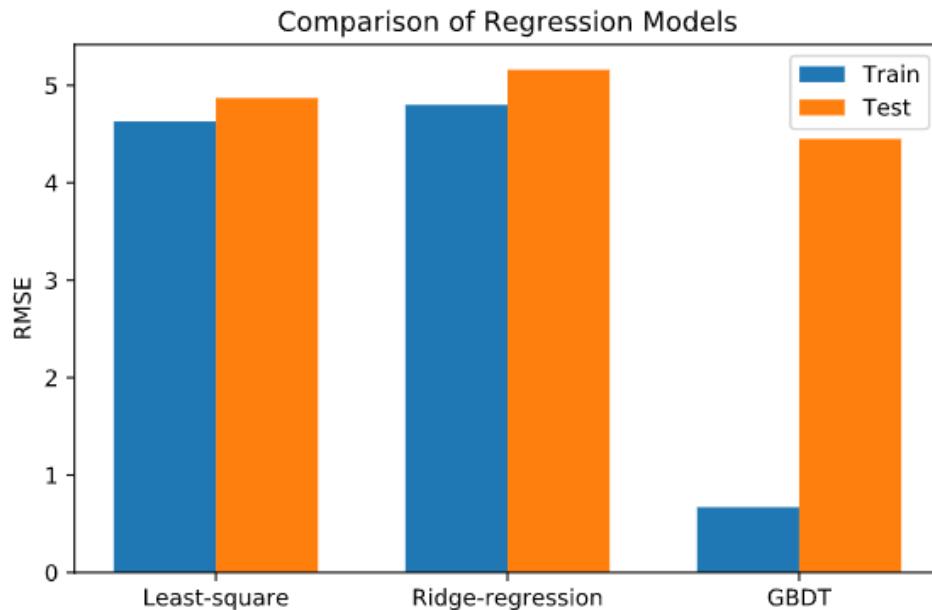
- 1) In gradient descent the goal of calculating gradients and updating θ accordingly is to optimize training loss. whereas in gradient boosting, one can fit a weak classifiers to the data and a loss function with respect to classifiers.
- 2) Gradient descent descends the gradient by introducing changes to parameters, whereas gradient boosting descends the gradient by introducing new models.

Problem 4(e)

Boston Housing Price GBDT

Train RMSE: 0.6733030286580176

Test RMSE: 4.448467175942173



From the above plot it is evident that both in train and test cases for Boston housing data GBDT is performing best than least-square and ridge regression.

Problem 4(f)

Credit-g-dataset:

Train Accuracy: 0.9742857142857143

Test Accuracy: 0.7466666666666667

Breast-cancer dataset:

Train Accuracy: 1.0

Test Accuracy: 0.9707602339181286

Problem 4(g)

For all cases, GBDT's accuracy and RMSE score is better than random forest. For RF each iteration the classifier is trained independently from the rest. Where GBDT is a boosting method, which builds on weak classifiers. It adds a classifier at a time, so that the next classifier is trained to improve the already trained ensemble, that makes the GBDT stronger than another classifier.