

HW12 CS 189

1.

Who else did you work with on this homework? In case of course events, just describe the group. How did you work on this homework? Any comments about the homework?

I worked on this homework with Ehimare Okoyomon, Prashanth Ganeth, and Daniel Mockaitis. We worked by getting together throughout the week and communicating on facebook.

I certify that all solutions are entirely in my words and that I have not looked at another student's solutions. I have credited all external sources in this write up}

Nicholas Lorio, 26089160

Question 5 Own Question.

Q: When is random forest a poor choice for an algorithm?

Accurate ensembles can require more trees, which means using the model becomes slower. In most situations random forest is fast enough, but there are a lot of situations where run-time performance is a necessary constraint and therefore a different algorithm should be used.

Another reason random forest is a poor choice is when we are trying to find the relationships in the data. If it's important to find the relationships in the data then a different algorithm should be used. There is also the problem of overfitting when there are too many high cardinality categorical variables. It is common to run into overfitting with random forests.

<https://www.quora.com/When-is-a-random-forest-a-poor-choicerelative-to-other-algorithms>

189 - 12

$$2.a. \underset{\mathbf{w}}{\text{Min}} \frac{1}{2} \|\mathbf{y} - \mathbf{x}\mathbf{w}\|^2 + \lambda \|\mathbf{w}\|,$$

$$= \underset{\mathbf{w}}{\text{Min}} \frac{1}{2} \left\| \mathbf{y} - \sum_{i=1}^d w_i x_i \right\|^2 + \lambda \sum_{i=1}^d |w_i| \quad \text{as } \mathbf{y} = \sum_{i=1}^d (x_i) w_i$$

$$= \underset{\mathbf{w}}{\text{Min}} \frac{1}{2} \left(\mathbf{y}^T \mathbf{y} + \sum_{i=1}^d (x_i^T \mathbf{w})^2 x_i - 2 \mathbf{y}^T \mathbf{w} \cdot \mathbf{x}_i + \lambda \|\mathbf{w}\| \right)$$

$$= \underset{\mathbf{w}}{\text{Min}} \frac{1}{2} \|\mathbf{y}\|^2 + \sum_{i=1}^d (x_i^T \mathbf{w})^2 x_i - 2 \mathbf{y}^T \mathbf{w} \cdot \mathbf{x}_i + \lambda \|\mathbf{w}\|$$

thus $\mathbf{g}(\mathbf{y}) = \frac{1}{2} \|\mathbf{y}\|^2 - \mathbf{f}(\mathbf{x}, \mathbf{y}, \mathbf{w}, \lambda) = \mathbf{x}^T \mathbf{w} \cdot \mathbf{x}_i - 2 \mathbf{y}^T \mathbf{w} \cdot \mathbf{x}_i + \lambda \|\mathbf{w}\|$

$$2.b. \hat{w}_i > 0 \quad \hat{w}_i = \underset{\mathbf{w}_i}{\text{Min}} \frac{1}{2} \|\mathbf{y} - \mathbf{x}\mathbf{w}\|^2 + \lambda \|\mathbf{w}\|,$$

$$\underset{\mathbf{w}_i}{\text{Min}} \frac{1}{2} \|\mathbf{y}\|^2 + \sum_{j \neq i}^d (x_j^T \mathbf{w}_i)^2 x_i - 2 \mathbf{y}^T \mathbf{w}_i \cdot \mathbf{x}_i + \lambda \|\mathbf{w}_i\|$$

• Convex Objeciton function \rightarrow use line search method
• separate terms that depend on w_i

$$L(\mathbf{w}) = \frac{1}{2} \|\mathbf{y} - \mathbf{x}\mathbf{w}\|^2 + \lambda \|\mathbf{w}\| = \frac{1}{2} \left\| \sum_{j=1}^d w_j x_j - \mathbf{y} \right\|^2 + \lambda \|\mathbf{w}\| + \lambda \sum_{j \neq i}^d |w_j|$$

let us denote $C^{(1)} = \sum_{j \neq i}^d w_j x_j - \mathbf{y}$ & $C^{(2)} = \lambda \sum_{j \neq i}^d |w_j|$

$$= (\lambda \|\mathbf{w}\| + C^{(2)}) + \frac{1}{2} \|w_i x_i + C^{(1)}\|_2^2$$

$$= \lambda \|\mathbf{w}\| + C^{(2)} + \frac{1}{2} \sum_{j=1}^d (w_j x_j + C_j^{(1)})^2$$

$\hat{\omega}_i > 0$

$$\frac{\delta L}{\delta \omega_i} = \lambda \cdot 1 + 0 + \frac{1}{2} \sum_{j=1}^J (1 \cdot x_{j,i} + 0) \cdot 2(w_i x_{j,i} + c_j^{(1)}) = 0$$

$$\lambda + 2 \sum_{j=1}^J x_{j,i} w_i x_{j,i} + x_{j,i} c_j^{(1)} = 0$$

$$\left[\begin{array}{l} \hat{\omega}_i = -\lambda + 2 \sum_{j=1}^J x_{j,i} c_j^{(1)} \\ \hline 2 \sum_{j=1}^J x_{j,i}^2 \end{array} \right] \text{ but } a = 2 \sum_{j=1}^J x_{j,i} c_j^{(1)}$$

$$b = 2 \sum_{j=1}^J x_{j,i}^2$$

$$\hat{\omega}_i^* = \frac{-\lambda + a}{b} \quad \text{for } \hat{\omega}_i > 0$$

c. $\hat{\omega}_i < 0$

$$\frac{\delta L}{\delta \omega_i} = -\lambda + 0 + \frac{1}{2} \sum_{j=1}^J (-w_i x_{j,i} + c_j^{(1)}) = 0$$

$$\hat{\omega}_i = \frac{\lambda + a}{b} \quad \text{for } \hat{\omega}_i < 0$$

d. $\hat{\omega}_i = 0$ if i) $\frac{\lambda + a}{b} \geq 0$ or if ii) $\frac{\lambda + a}{b} \leq 0$

$$\lambda + 2 \sum_{j=1}^J x_{j,i} c_j^{(1)} \geq 0 \quad \lambda \geq -2 \sum_{j=1}^J x_{j,i} c_j^{(1)} \quad \text{i) } \lambda \geq -a$$

$$\lambda \geq -2 \sum_{j=1}^J x_{j,i} c_j^{(1)} \quad \text{ii) } \lambda \leq a$$

so $\hat{\omega}_i = 0$ if $\lambda = \pm a$

• if neither $\frac{\lambda + a}{b} > 0$ or $\frac{\lambda + a}{b} < 0$ then \exists no optimum

→ lasso is convex so \exists is an opt somewhere, for this case $w_i^* = 0$ thus $-\frac{\lambda + a}{b} \geq 0 \Leftrightarrow \frac{\lambda + a}{b} \leq 0$ which is equivalent to $|a| \leq \lambda$

$$(w_1^* = \hat{w}_1)$$

$$\text{e.g. } \hat{w} = \arg \min_w \frac{1}{2} \|y - Xw\|_2^2 + \lambda \|w\|_2^2$$

$$\arg \min_{w_i} \frac{1}{2} \left\| \sum_{j=1}^d w_j x_j - y \right\|_2^2 + \lambda \sum_{j=1}^d w_j^2 \quad \text{CONVEX.}$$

$$\arg \min_{w_i} \frac{1}{2} \underbrace{\left\| w_i x_i + \sum_{j \neq i} w_j x_j - y \right\|_2^2}_{C^{(1)}} + \underbrace{\lambda w_i^2 + \sum_{j \neq i} \lambda w_j^2}_{C^{(2)}}$$

$$\arg \min_{w_i} \frac{1}{2} \sum_{j \neq i} (w_i x_{ji} + C_j^{(1)})^2 + \lambda w_i^2 + C^{(2)}$$

$$\delta L / \delta w_i = \sum_{j=1}^d x_{ji} (w_i^* x_{ji} + C_j^{(1)}) + \lambda \cdot 2 w_i^* + 0 = 0$$

$$= 2 \sum_{j=1}^d x_{ji} w_i^* x_{ji} + \lambda \cdot 2 w_i^* = 0$$

$$= 2 w_i^* \sum_{j=1}^d x_{ji}^2 + 2 \sum_{j=1}^d x_{ji} C_j^{(1)} + \lambda \cdot 2 w_i^* = 0$$

$$= w_i^* \left(2 \sum_{j=1}^d x_{ji}^2 + 2 \lambda \right) + \sum_{j=1}^d x_{ji} C_j^{(1)} = 0$$

$$\left[\begin{array}{l} w_i^* = - \sum_{j=1}^d x_{ji} C_j^{(1)} \\ \sum_{j=1}^d x_{ji}^2 + 2 \lambda \end{array} \right] = \hat{w}_i = 0 \quad \text{when} \quad - \sum_{j=1}^d x_{ji} C_j^{(1)} = 0$$

$\text{so } \hat{w}_i = 0 \text{ when } a = 0$
or when $\lambda \rightarrow \infty$

$$= -a / (b + 2\lambda)$$

• l1 norm promotes sparsity as $\hat{w}_i = 0$ when $|a| \geq |b|$, by increasing λ we increase the threshold for $|a|$ to be set to zero and our solution becomes more sparse.

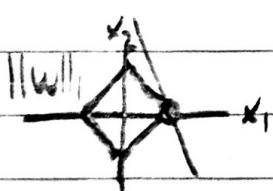
• a is some function of the data & other weights.

- the regularization term in ℓ_1 norm tends to pull the least squares update towards zero.
- ℓ_1 norm sparsity for when $|a| \leq \lambda$ we set $w_i = 0$. Because sparse as described previously.
- ℓ_2 norm does not promote sparse solutions. In contrast $\hat{w}_i = 0$ is not dependent on A in the same manner as ℓ_1 for sparsity.

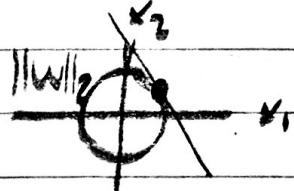
$\hat{w}_i = 0$ when $\lambda \rightarrow \infty$

or
 $\hat{w}_i = 0$ when $\lambda = 0$

thus ℓ_1 leads to frequently sparse solutions, many coeff w/ zero values, & ℓ_2 has weaker requirements for $\hat{w}_i = 0$.



ℓ_1



ℓ_2

2.5.

(from piazza: "Reporting the hyper parameter
is sufficient". we don't have to include
the ipython notebook for this part)

$\lambda = 0.0001$ is a sufficient hyperparameter
for L_1 reg to extract the image

$$3.b. P(\max_i |Z_i| \geq 2\sigma(\log(d))^{1/2})$$

$$= 1 - P(\max_i |Z_i| < t) = 1 - P(|Z_1| < t, \dots, |Z_d| < t)$$

$$\geq 1 - \sum_{i=1}^d P(|Z_i| < t) \quad \text{By Hint 2}$$

$$\geq 1 - \sum_{i=1}^d (1 - P(|Z_i| \geq t))$$

$$\cdot P(Z_i \geq t) \leq \exp\left(-\frac{(2\sigma(\log(d))^{1/2})^2}{2\sigma^2}\right)$$

$$\leq \exp\left(-\frac{2\sigma^2 \log(d)}{2\sigma^2}\right) = \exp(-2\log(d)) = \frac{1}{d^2}$$

$$\geq 1 - \sum_{i=1}^d \underbrace{(1 - P(Z_i \geq t))}_{\leq 1/d^2}$$

$$\geq 1 - 1/d^2$$

$$\geq d(1 - 1/d^2)$$

$$\leq 1 - d(1 - 1/d^2)$$

$$\leq \underbrace{1 - d + 1/d}_{\leq 1/d} \leq 1/d \quad \text{as required}$$

as $d \geq 1$ thus $1 - d \leq 0$

3.e. $X \in \mathbb{R}^{n \times d}$ (restricted to orthonormal columns)

Show $\hat{w}_{top}(s)$ returns top s entries of vector

$w^* + z'$ in absolute value. • $z' \sim N(0, \sigma^2)$

$\min ||y - Xw||^2$ gives \hat{w}_{top} for $y = Xw + z$

$$z = y - Xw$$

$$\sim N(0, \sigma^2)$$

$\hat{w}_{top}(s) = T_s(\hat{w}_{top})$

$$X = \begin{bmatrix} x_1 & \dots & x_d \end{bmatrix}, \quad Xw = \begin{bmatrix} x_1 & \dots & x_d \end{bmatrix} \begin{bmatrix} w_1 \\ \vdots \\ w_d \end{bmatrix} + \begin{bmatrix} z_1 \\ \vdots \\ z_d \end{bmatrix}$$

$\hat{w}_{top} = X^T y$ • zero out all but top s entries of w_{top} in absolute value.

$$(X^T X)^{-1} X^T y \quad \hat{w}_{top} = X^T (Xw^* + z') \quad w/ z \sim N(0, \sigma^2) \quad w/ \text{var } \sigma^2$$

$$= I^{-1} X^T y \quad \therefore X^T X = I$$

$$= X^T y \quad \therefore \text{let } z' = X^T z$$

$$\hat{w}_{top} = w^* + z'$$

$$\begin{bmatrix} z_1 \\ \vdots \\ z_d \end{bmatrix} = X^T \begin{bmatrix} z_1 \\ \vdots \\ z_d \end{bmatrix}$$

• Let C be the covariance matrix for z :

→ C is a diagonal matrix as z_i are independent.

→ z' has a covariance matrix $\hat{C} = X^T C X^* = X^T z z^T X$

\hat{C} is also a diagonal matrix as X has orthonormal columns & C is a diagonal matrix (as established early)

$$\rightarrow \text{as } \text{cov}(Az) = z^T \text{cov}(A) z$$

3.e. cont.

• \mathbf{z}_i are jointly Gaussian Random Variables, each has covar matrix given by $\hat{\mathbf{C}} = \mathbf{X}^T \mathbf{C} \mathbf{X}$,
→ $\hat{\mathbf{C}}$ is diagonal
thus \mathbf{z}_i are Marginal R.V.

We can ∴ conclude that $\hat{w}_{top}(s)$ returns the s top alns of $w^* + \mathbf{z}^i = \hat{w}_{top}$ in abs val w/ \mathbf{z}^i iid gauss w/ var σ^2 .

3.f. $e = \hat{w}_{top}(s) - w^*$ is (@ Most) 2S-Sparse

$$\hat{w}_{top}(s) = \begin{bmatrix} a + z_1 \\ b + z_2 \\ \vdots \\ 0 \\ \vdots \\ 0 \end{bmatrix} \quad \begin{array}{l} s \text{ non-zero} \\ \text{assume the top } s \text{ return} \\ \text{from } \hat{w}_{top}(s) \text{ are the } s \\ \text{largest in magnitude of } \hat{w}_{top}. \end{array}$$

$$w^* = \begin{bmatrix} 0 \\ \vdots \\ 0 \\ \vdots \\ 0 \end{bmatrix} \quad \begin{array}{l} s \text{ NonZero} \\ s = s \text{-sparse} \end{array}$$

$$e = \begin{bmatrix} a + z_1 \\ b + z_2 \\ \vdots \\ 0 \\ \vdots \\ 0 \\ \vdots \\ -y \end{bmatrix} \quad \begin{array}{l} s \text{ NonZero} \\ s = s \text{-sparse} \end{array} \quad \begin{array}{l} \text{thus } e \text{ has @} \\ \text{most } 2s \text{ non zero &} \\ \text{is 2S-Sparse (@ most)} \end{array}$$

3.g. Show $|e_i| \leq 4\sigma(\log(d))^{1/2}$ for each i

$$\mathcal{E} = \{e_i\}_{i=1,\dots}$$

$$\hat{\omega}_{top}(s) = \begin{bmatrix} z_1 \\ z_2 \\ \vdots \\ 0 \\ \vdots \\ 0 \end{bmatrix} s \text{ nonzero } \quad \omega^* = \begin{bmatrix} 0 \\ 0 \\ \vdots \\ 0 \end{bmatrix} s \text{ nonzero}$$

$$e_i = \begin{bmatrix} z_1 \\ z_2 \\ \vdots \\ 0 \end{bmatrix} s \text{ nonzero } z_i \quad \left| \begin{array}{l} \text{Max error for one is } 2\sigma(\log(d))^{1/2} \\ \text{thus} \end{array} \right.$$

$$e = (\hat{\omega}_{top}(s) - \omega^*) = (z + \bar{z}) - \omega^* = z' \in$$

$$\text{Max } |z'_i| \leq 2\sigma\sqrt{\log(d)} \quad \leftarrow \text{Lemma}$$

if top s entries are different than ω^* :

$$\text{Max}(e_i) - \text{Min}(e_i) = 2\sigma\sqrt{\log(d)} - (-2\sigma\sqrt{\log(d)})$$

$$\text{so } |e_i| \leq 4\sigma\sqrt{\log(d)}$$

$$3.H. P\left(\|\hat{w}_{top}(s) - w^*\|_2^2 \leq 32\sigma^2 s \log(d)\right) \geq 1 - 1/d$$

$$= P\left(\|e\|_2^2 \leq 32\sigma^2 \cdot s \cdot \log(d)\right) = P\left(\sum_i e_i^2 \leq 32\sigma^2 s \log(d)\right)$$

• e is $\in \mathbb{R}^{2s}$ most sparse so

$$\leq P\left(2s \cdot \max_{i \in \{1, \dots, d\}} |e_i|^2 \leq 32\sigma^2 s \log(d)\right)$$

$$\leq P\left(\max_i |e_i|^2 \leq 16\sigma^2 \log(d)\right)$$

$$= P\left(\max_i |e_i| \leq 4\sigma(\log(d))^{1/2}\right) \quad \text{this implies}$$

$$\Rightarrow P\left(\max_i |z_i| \leq 2\sigma \sqrt{\log(d)}\right) = 1 - P\left(\max_i |z_i| > 2\sigma (\log(d))^{1/2}\right) \leq 1/d$$

$$\geq 1 - 1/d \quad \text{as required.}$$

$$3.I. P\left(\frac{1}{n} \|\mathcal{X}(\hat{w}_{top}(s) - w^*)\|_2^2 \leq \underbrace{32\sigma^2 s \log(d)}_n\right) \geq 1 - 1/d$$

$$\bullet \|\mathcal{X}(\hat{w}_{top}(s) - w^*)\|_2^2 = \|Xe\|_2^2 = (xe)^T (xe) = e^T X^T X e = e^T e = \|e\|_2^2$$

$$= P\left(\frac{1}{n} \|e\|^2 \leq 32\sigma^2 s \log(d)\right) = P\left(\|e\|^2 \leq 32\sigma^2 s \log(d)\right) \geq 1 - 1/d$$

(From part 3.H)

$$3.5. \text{ given } E\left[\frac{1}{n} \|\chi(\hat{\omega}_{top} - \omega^*)\|_2^2\right] = \sigma^2 \frac{d}{n}$$

$$\& E\left[\|\hat{\omega}_{top} - \omega^*\|_2^2\right] = \sigma^2 d$$

• From previous parts ① $P(\|\epsilon\|_2^2 \leq 32\sigma^2 s \log(d)) \geq 1 - 1/d$

w.p. $\geq 1 - \frac{1}{d}$ ① ② $P\left(\frac{1}{n} \|\chi\epsilon\|_2^2 \leq 32\sigma^2 \frac{s \log(d)}{n}\right) \geq 1 - 1/d$

$$E\left[\|\hat{\omega}_{top}(s) - \omega^*\|_2^2\right] \leq E[32\sigma^2 \log(d)] = 32\sigma^2 s \log(d)$$

Similarly w.p. $\geq 1 - 1/d$ ②

$$E\left[\frac{1}{n} \|\chi\epsilon\|_2^2\right] \leq E[32\sigma^2 \frac{s \log(d)}{n}] = 32\sigma^2 \frac{s \log(d)}{n}$$

$$s \leq \frac{d}{32 \log(d)} \quad \text{then} \begin{cases} E\left[\|\hat{\omega}_{top}(s) - \omega^*\|_2^2\right] \leq \sigma^2 d \\ E\left[\frac{1}{n} \|\hat{\omega}_{top}(s) - \omega^*\|_2^2\right] \leq \sigma^2 \frac{d}{n} \end{cases}$$

3.6. Given that we know the important s factors

then $\epsilon = \hat{\omega}_{top}(s) - \omega^*$ is s -sparse

$$\Rightarrow P(\|\epsilon\|_2^2 \leq 16\sigma^2 s \log(d)) = P(\sum_i |\epsilon_i|^2 \leq 16\sigma^2 s \log(d))$$

$$\leq P(s \max_i |\epsilon_i|^2 \leq 16\sigma^2 s \log(d))$$

$$= P(\max_i |\epsilon_i| \leq 4\sigma(\log(d))^{1/2})$$

$$\geq (1 - 1/d) \quad (\text{similar to 3.4.})$$

3.k. cont.

Now,

- Similar to j: if $[s \leq \frac{d}{16\log(d)}]$

$$\rightarrow E\left[\frac{1}{n} \| \chi(\hat{\omega}_{top}(s) - \omega^*) \|_2^2 \right] \leq 16 \underbrace{\sigma^2 s \log d}_{\uparrow} \leq \frac{\sigma^2 d}{n}$$

$$\rightarrow E\left[\| \hat{\omega}_{top}(s) - \omega^* \|_2^2 \right] \leq 16 \sigma^2 s \log(d) \leq \sigma^2 d$$

- this is lower than the behavior for the sparse $\hat{\omega}_{top}(s)$, in which we did not know the important s Features, derived earlier.

\rightarrow So, if we don't know the important Features then we @ least double the variance.

3. Bias and Variance of Sparse Linear Regression

```
In [6]: import numpy as np
import matplotlib.pyplot as plt
%matplotlib inline
```

```
In [7]: def ground_truth(n, d, s):
    """
        Input: Two positive integers n, d. Requires n >= d >= s.
        Output: A tuple containing i) random matrix of dimension n X d with orthonormal
columns. and
                ii) a d-dimensional, s-sparse wstar with (large) Gaussian entries
    """
    if d > n:
        print("Too many dimensions")
        return None

    if d < s:
        s = d
    A = np.random.randn(n, d) #random Gaussian matrix
    U, S, V = np.linalg.svd(A, full_matrices=False) #reduced SVD of Gaussian matrix
    wstar = np.zeros(d)
    wstar[:s] = 10 * np.random.randn(s)

    np.random.shuffle(wstar)
    return U, wstar

def get_obs(U, wstar):
    """
        Input: U is an n X d matrix and wstar is a d X 1 vector.
        Output: Returns the n-dimensional noisy observation y = U * wstar + z.
    """
    n, d = np.shape(U)
    z = np.random.randn(n) #i.i.d. noise of variance 1
    y = np.dot(U, wstar) + z
    return y
```

```
In [8]: def LS(U, y):
    """
    Input: U is an n X d matrix with orthonormal columns and y is an n X 1 vector.
    Output: The OLS estimate what_{LS}, a d X 1 vector.
    """
    wls = np.dot(U.T, y) #pseudoinverse of orthonormal matrix is its transpose
    return wls

def thresh(U, y, lmbda):
    """
    Input: U is an n X d matrix and y is an n X 1 vector; lambda is a scalar threshold
    of the entries.
    Output: The estimate what_{T}(lambda), a d X 1 vector that is hard-thresholded
    (in absolute value) at level lambda.
    When code is unfilled, returns the all-zero d-vector.
    """
    n, d = np.shape(U)
    wls = LS(U, y)
    what = np.zeros(d)

    #print np.shape(wls)
    #####
    #TODO: Fill in thresholding function; store result in what
    #####
    #YOUR CODE HERE:

    for i in range(len(what)):
        if np.abs(wls[i]) > lmbda:
            what[i] = wls[i]

    #####
    return what

def topk(U, y, s):
    """
    Input: U is an n X d matrix and y is an n X 1 vector; s is a positive integer.
    Output: The estimate what_{top}(s), a d X 1 vector that has at most s non-zero
    entries.
    When code is unfilled, returns the all-zero d-vector.
    """
    n, d = np.shape(U)
    what = np.zeros(d)
    wls = LS(U, y)

    #####
    #TODO: Fill in thresholding function; store result in what
    #####
    #YOUR CODE HERE: Remember the absolute value!
    inds = np.argpartition(np.abs(wls), -s)[-s:]
    what[inds] = wls[inds]
    #####
    return what
```

```
In [9]: def error_calc(num_iters=10, param='n', n=1000, d=100, s=5, s_model=True, true_s=5):
    """
        Plots the prediction error  $1/n + U(\text{what} - \text{wstar})^2 = 1/n + \text{what} - \text{wstar}$  ||^2
        for the three estimators averaged over num_iter experiments.

        Input:
        Output: 4 arrays -- range of parameters, errors of LS, topk, and thresh estimator, respectively. If thresh and topk
                functions have not been implemented yet, then these errors are simply the norm of wstar.
    """
    wls_error = []
    wtopk_error = []
    wthresh_error = []

    if param == 'n':
        arg_range = np.arange(100, 2000, 50)
        lmbda = 2 * np.sqrt(np.log(d))
        for n in arg_range:
            U, wstar = ground_truth(n, d, s) if s_model else ground_truth(n, d, true_s)
            error_wls = 0
            error_wtopk = 0
            error_wthresh = 0
            for count in range(num_iters):
                y = get_obs(U, wstar)
                wls = LS(U, y)
                wtopk = topk(U, y, s)
                wthresh = thresh(U, y, lmbda)
                error_wls += np.linalg.norm(wstar - wls)**2
                error_wtopk += np.linalg.norm(wstar - wtopk)**2
                error_wthresh += np.linalg.norm(wstar - wthresh)**2
            wls_error.append(float(error_wls)/ n / num_iters)
            wtopk_error.append(float(error_wtopk)/ n / num_iters)
            wthresh_error.append(float(error_wthresh)/ n / num_iters)

    elif param == 'd':
        arg_range = np.arange(10, 1000, 50)
        for d in arg_range:
            lmbda = 2 * np.sqrt(np.log(d))
            U, wstar = ground_truth(n, d, s) if s_model else ground_truth(n, d, true_s)
            error_wls = 0
            error_wtopk = 0
            error_wthresh = 0
            for count in range(num_iters):
                y = get_obs(U, wstar)
                wls = LS(U, y)
                wtopk = topk(U, y, s)
                wthresh = thresh(U, y, lmbda)
                error_wls += np.linalg.norm(wstar - wls)**2
                error_wtopk += np.linalg.norm(wstar - wtopk)**2
                error_wthresh += np.linalg.norm(wstar - wthresh)**2
            wls_error.append(float(error_wls)/ n / num_iters)
            wtopk_error.append(float(error_wtopk)/ n / num_iters)
            wthresh_error.append(float(error_wthresh)/ n / num_iters)

    elif param == 's':
        arg_range = np.arange(5, 55, 5)
        lmbda = 2 * np.sqrt(np.log(d))
        for s in arg_range:
            U, wstar = ground_truth(n, d, s) if s_model else ground_truth(n, d, true_s)
```

a

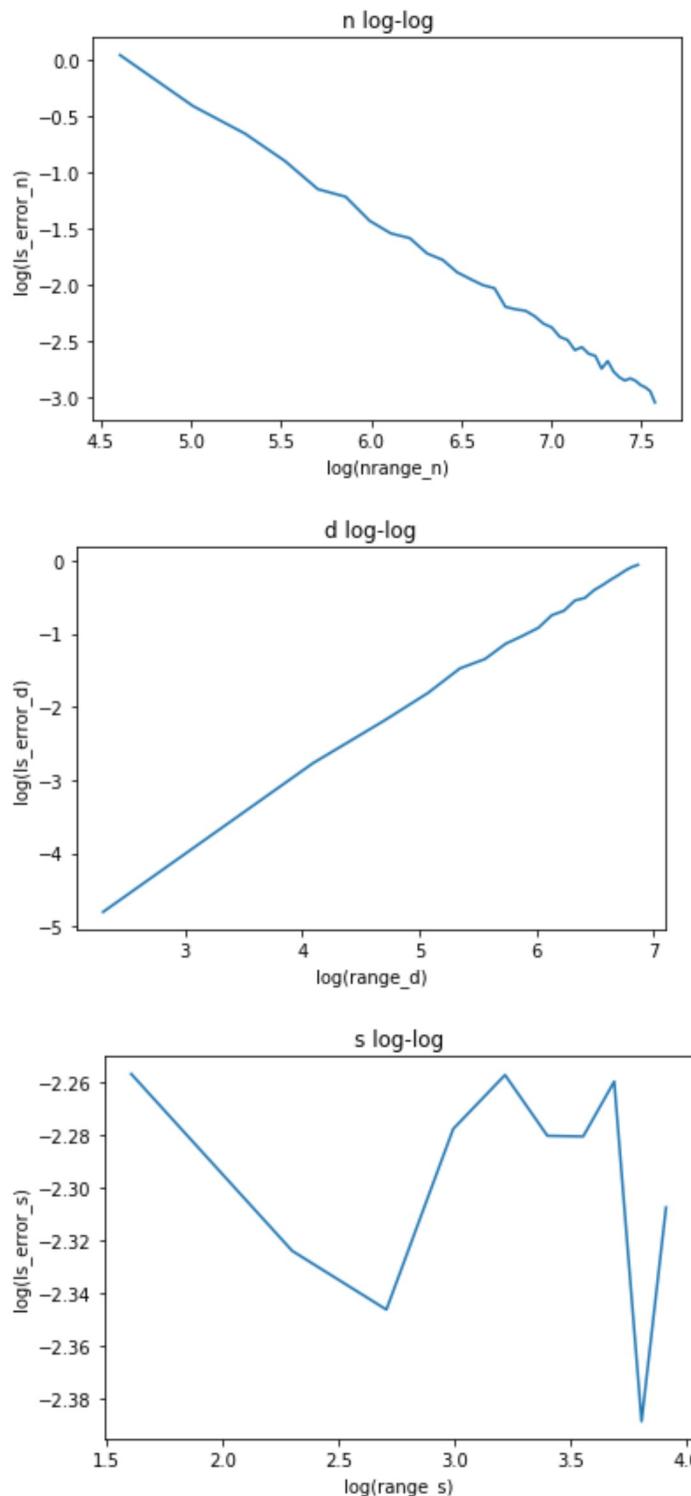
```
In [10]: #nrange contains the range of n used, ls_error the corresponding errors for the OLS
estimate
nrange, ls_error, _, _ = error_calc(num_iters=10, param='n', n=1000, d=100, s=5, s_
model=True, true_s=5)
#####
#TODO: Your code here: call the helper function for d and s, and plot everything
nrange_d, ls_error_d, _, _ = error_calc(num_iters=10, param='d', n=1000, d=100, s=5
, s_model=True, true_s=5)
nrange_s, ls_error_s, _, _ = error_calc(num_iters=10, param='s', n=1000, d=100, s=5
, s_model=True, true_s=5)

#####
#YOUR CODE HERE:

plt.title("n log-log")
plt.plot(np.log(nrange), np.log(ls_error))
plt.xlabel('log(nrange_n)')
plt.ylabel('log(ls_error_n)')
plt.show()

plt.title("d log-log")
plt.plot(np.log(nrange_d), np.log(ls_error_d))
plt.xlabel('log(range_d)')
plt.ylabel('log(ls_error_d)')
plt.show()

plt.title("s log-log")
plt.plot(np.log(nrange_s), np.log(ls_error_s))
plt.xlabel('log(range_s)')
plt.ylabel('log(ls_error_s)')
plt.show()
```



Are these plots as expected? Discuss. Also put down your parameter choices (either here or in plot captions.) It's fine to use the default values, but put them down nonetheless.

Answer Here

Paramaters for plot of error as function of n:

n=1000, d=100, s=5, s_model=True, true_s=5

Paramaters for plot of error as function of d:

n=1000, d=100, s=5, s_model=True, true_s=5

Paramaters for plot of error as function of s:

n=1000, d=100, s=5, s_model=True, true_s=5

- **First Plot:** Error decreases as we increase the number of data points n. First plot makes intutive sense. The bias decreases as a function of n increasing.
- **Second Plot:** Error increases as the degree increases. The number of data points n is fixed. The complexity increases as the degree of the polynomial increases. Increasing complexity increases the variance which contributes to the increased error.
- **Third Plot:** Error is not correlated with the sparcity.

b

In [17]:

```
#TODO: Part (b)
#####
#YOUR CODE HERE:

nrange, ls_error, thresh_error_n, topk_error_n = error_calc(num_iters=10, param='n'
, n=1000, d=100, s=5, s_model=True, true_s=5)
nrange_d, ls_error_d, thresh_error_d, topk_error_d = error_calc(num_iters=10, param
='d', n=1000, d=100, s=5, s_model=True, true_s=5)
nrange_s, ls_error_s, thresh_error_s, topk_error_s = error_calc(num_iters=10, param
='s', n=1000, d=100, s=5, s_model=True, true_s=5)

plt.title("Error of all estimators as a function of n")
plt.plot(np.log(nrange), np.log(ls_error), label = 'OLS')
plt.plot(np.log(nrange), np.log(thresh_error_n), label = 'thresh')
plt.plot(np.log(nrange), np.log(topk_error_n), label = 'topk')
plt.xlabel('log(n)')
plt.ylabel('log(Error)')
plt.legend(loc='upper right', shadow=True)

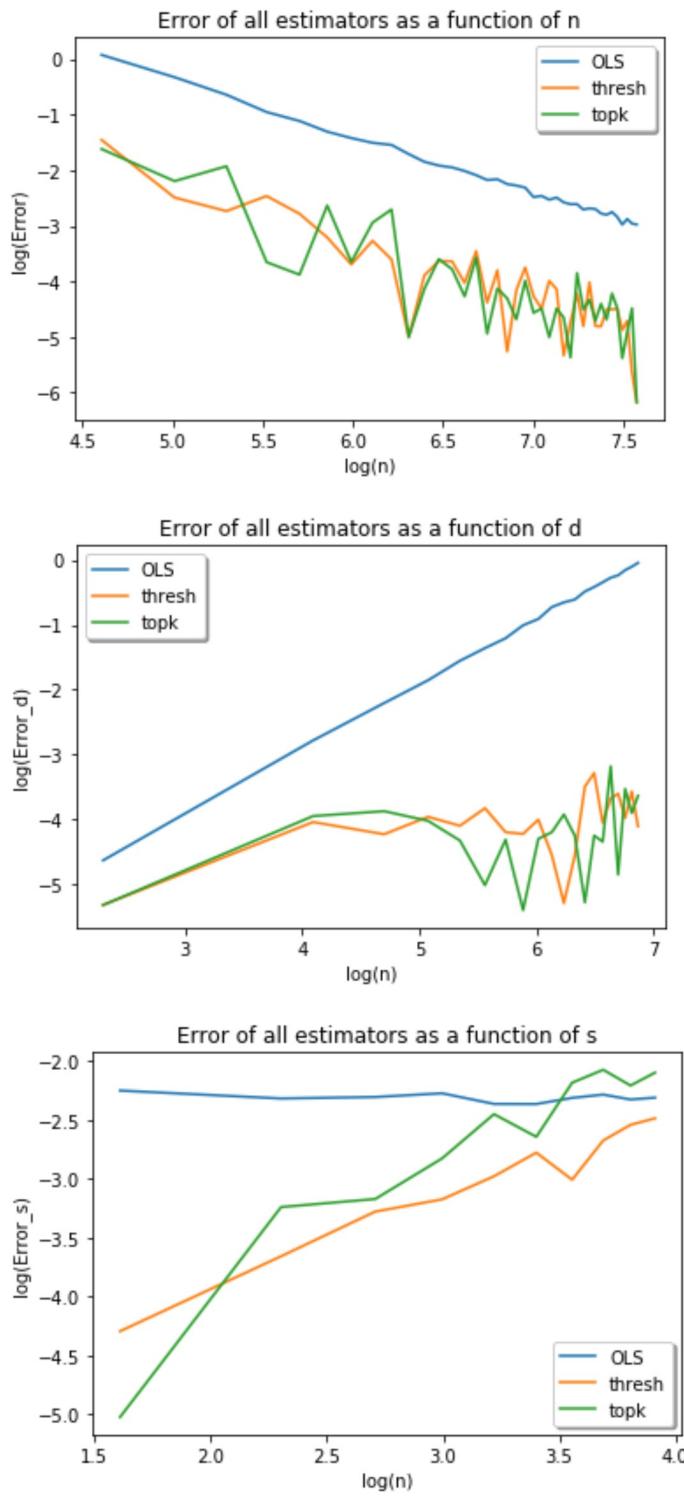
plt.show()

plt.title("Error of all estimators as a function of d")
plt.plot(np.log(nrange_d), np.log(ls_error_d), label = 'OLS')
plt.plot(np.log(nrange_d), np.log(thresh_error_d), label = 'thresh')
plt.plot(np.log(nrange_d), np.log(topk_error_d), label = 'topk')
plt.xlabel('log(n)')
plt.ylabel('log(Error_d)')
plt.legend(loc='upper left', shadow=True)

plt.show()

plt.title("Error of all estimators as a function of s")
plt.plot(np.log(nrange_s), np.log(ls_error_s), label = 'OLS')
plt.plot(np.log(nrange_s), np.log(thresh_error_s), label = 'thresh')
plt.plot(np.log(nrange_s), np.log(topk_error_s), label = 'topk')
plt.xlabel('log(n)')
plt.ylabel('log(Error_s)')
plt.legend(loc='lower right', shadow=True)

plt.show()
```

**C**

In [18]:

```
#TODO: Part (c)
#####
#YOUR CODE HERE:
nrange, ls_error, thresh_error_n, topk_error_n = error_calc(num_iters=10, param='n'
, n=1000, d=100, s=5, s_model=False, true_s=100)
nrange_d, ls_error_d, thresh_error_d, topk_error_d = error_calc(num_iters=10, param
='d', n=1000, d=100, s=5, s_model=False, true_s=100)
nrange_s, ls_error_s, thresh_error_s, topk_error_s = error_calc(num_iters=10, param
='s', n=1000, d=100, s=5, s_model=False, true_s=100)

plt.title("Error of all estimators as a function of n")
plt.plot(np.log(nrange), np.log(ls_error), label = 'OLS')
plt.plot(np.log(nrange), np.log(thresh_error_n), label = 'thresh')
plt.plot(np.log(nrange), np.log(topk_error_n), label = 'topk')
plt.xlabel('log(n)')
plt.ylabel('log(Error)')
plt.legend(loc='upper right', shadow=True)

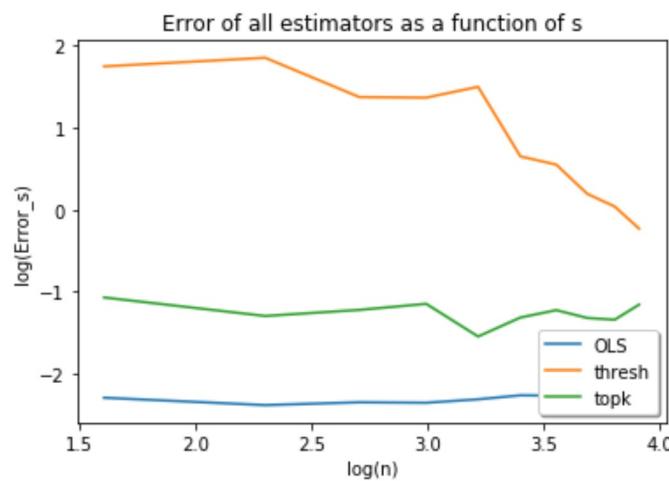
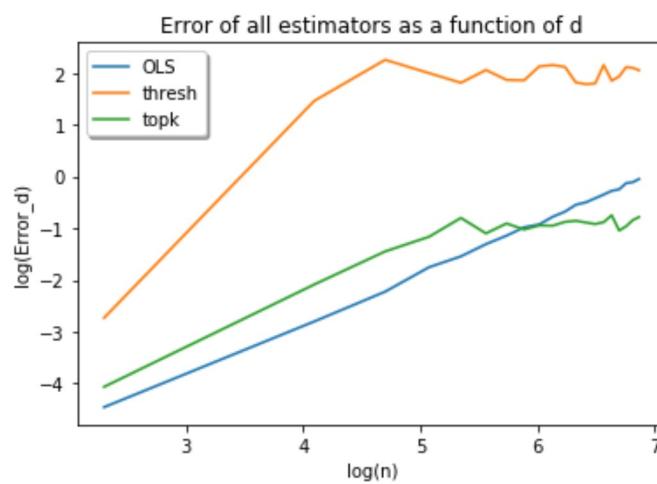
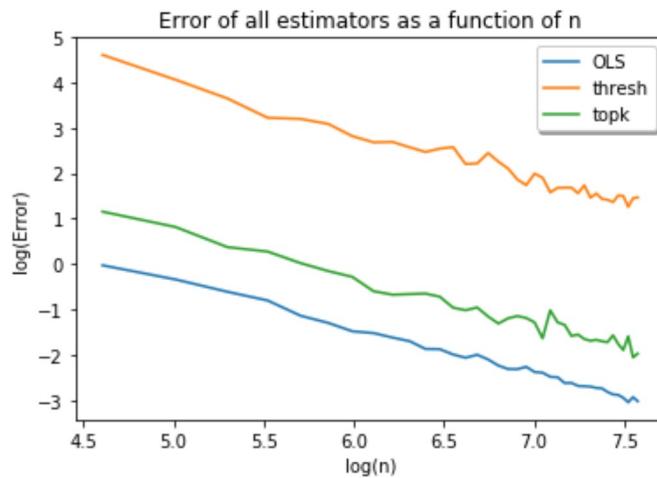
plt.show()

plt.title("Error of all estimators as a function of d")
plt.plot(np.log(nrange_d), np.log(ls_error_d), label = 'OLS')
plt.plot(np.log(nrange_d), np.log(thresh_error_d), label = 'thresh')
plt.plot(np.log(nrange_d), np.log(topk_error_d), label = 'topk')
plt.xlabel('log(n)')
plt.ylabel('log(Error_d)')
plt.legend(loc='upper left', shadow=True)

plt.show()

plt.title("Error of all estimators as a function of s")
plt.plot(np.log(nrange_s), np.log(ls_error_s), label = 'OLS')
plt.plot(np.log(nrange_s), np.log(thresh_error_s), label = 'thresh')
plt.plot(np.log(nrange_s), np.log(topk_error_s), label = 'topk')
plt.xlabel('log(n)')
plt.ylabel('log(Error_s)')
plt.legend(loc='lower right', shadow=True)

plt.show()
```



Param First Plot

```
error_calc(num_iters=10, param='n', n=1000, d=100, s=5, s_model=False, true_s=100)
```

Param Second Plot

```
error_calc(num_iters=10, param='d', n=1000, d=100, s=5, s_model=False, true_s=100)
```

Param Third Plot

```
error_calc(num_iters=10, param='s', n=1000, d=100, s=5, s_model=False, true_s=100)
```

Answer Here

- **First Plot:** Error as a function of n decreases as n increases. The bias decreases. For data from a non-sparse linear model: Same statements made earlier regarding the general behavior of the OLS plots are true here for the plotted error of the estimators and the OLS as a function of n.
- **Second Plot:** For data from a non-sparse linear model: Same statements made earlier regarding the general behavior of the OLS plots are true here for the plotted error of the estimators and the OLS as a function of d.

From Part a

"Error increases as the degree increases. The number of data points n is fixed. The complexity increases as the degree of the polynomial increases. Increasing complexity increases the variance which contributes to the increased error."

- **Third Plot:** As we increase the sparsity, we note that this reduces the error as the sparsity of the true model approaches the sparsity our estimators assume the data to have.

The true sparsity, s, of the model is greater than the s-sparsity assumed for our estimators. Our estimators want sparse data, the estimators will incorrectly estimate at most (true_s - s) weights in $w^{\hat{}}$ as zero when they are in fact non-zero/non-sparse coefficients of the weight vector. Thus our estimates have worse performance as they are assuming sparsity in the data when the true data is generated from a non-sparse linear model.

4. Decision Trees and Random Forests

```
In [7]: %matplotlib inline
```

a)

```
In [28]: # You may want to install "gprof2dot"
import io
from collections import Counter

import numpy as np
import scipy.io
import sklearn.model_selection
import sklearn.tree
from numpy import genfromtxt
from scipy import stats
from sklearn.base import BaseEstimator, ClassifierMixin

import pydot

eps = 1e-5 # a small number

class DecisionTree:
    def __init__(self, max_depth=3, feature_labels=None):
        self.max_depth = max_depth
        self.features = feature_labels
        self.left, self.right = None, None # for non-leaf nodes
        self.split_idx, self.thresh = None, None # for non-leaf nodes
        self.data, self.pred = None, None # for leaf nodes

    @staticmethod
    def information_gain(X, y, thresh):
        # TODO implement information gain function
        def entropy(ys):
            total = len(ys)
            if total == 0:
                return 0
            arr = [0, 0]
            for label in ys:
                arr[label] += 1
            return sum([- (pi/total)*np.log((pi/total)) for pi in arr if pi != 0])
        total = len(y)
        y_left = []
        y_right = []
        for i in range(total):
            if X[i] < thresh:
                y_left.append(y[i])
            else:
                y_right.append(y[i])
        return entropy(y) - (len(y_left)/total)*entropy(y_left) - (len(y_right)/total)*entropy(y_right)

    @staticmethod
    def gini_purification(X, y, thresh):
        # TODO implement gini_purification function
        def gini(ys):
            total = len(ys)
            if total == 0:
                return 0
            arr = [0, 0]
            for label in ys:
                arr[label] += 1
            return 1 - sum([(pi/total)**2 for pi in arr])
        total = len(y)
        y_left = []
        y_right = []
        for i in range(total):
            if X[i] < thresh:
                y_left.append(y[i])
            else:
```

b)

All answers can be found in preprocess function.

- Some data points are missing class labels;

```
# Temporarily assign -1 to missing data
```

- Some features are not numerical values;

```
# Hash the columns (used for handling strings)
```

- Some data points are missing some features.

```
# Replace missing data with the mode value. We use the mode instead of  
# the mean or median because this makes more sense for categorical  
# features such as gender or cabin type, which are not ordered.
```

Titanic Dataset

```
In [29]: if __name__ == "__main__":
    main("titanic")
```

```
/home/nlorio/.local/lib/python3.5/site-packages/ipykernel_launcher.py:287: VisibleDeprecationWarning: Reading unicode strings without specifying the encoding argument is deprecated. Set the encoding, use None for the system default.
/home/nlorio/.local/lib/python3.5/site-packages/ipykernel_launcher.py:289: VisibleDeprecationWarning: Reading unicode strings without specifying the encoding argument is deprecated. Set the encoding, use None for the system default.
```

```
Part (b): preprocessing the titanic dataset
Features: [b'pclass', b'sex', b'age', b'sibsp', b'parch', b'ticket', b'fare', b'cabin', b'embarked', b'male', b'female', b'S', b'C', b'Q']
Train/test size: (999, 14) (310, 14)
```

```
Part 0: constant classifier
Accuracy 0.6136136136136137
```

```
Part (a-b): simplified decision tree
Predictions [1. 1. 0. 0. 0. 1. 0. 0. 0. 0. 1. 1. 1. 1. 1. 1. 1. 0. 0. 0. 0. 1.
0. 0. 1. 1. 0. 0. 0. 0. 0. 0. 0. 0. 0. 0. 0. 0. 0. 1. 0. 0. 1. 0. 0. 0. 0.
1. 0. 0. 0. 0. 0. 1. 1. 0. 0. 0. 1. 0. 0. 0. 1. 1. 0. 0. 1. 0. 1. 1. 0. 0. 0.
1. 0. 0. 0. 0. 0. 0. 1. 0. 0. 0. 0. 0. 0. 0. 1. 0. 1. 0. 0. 1. 0. 0. 0. 1. 0. 0.
1. 0. 0. 0.]
```

```
Part (c): sklearn's decision tree
Cross validation [0.79640719 0.76576577 0.79819277]
```

```
Part (e): Our Bagged Tree for titanic
Cross validation [0.79041916 0.78378378 0.8253012 ]
First splits [(b'male', 101), (b'female', 99)]
BaggedTrees(n=200, params={'max_depth': 5, 'min_samples_leaf': 10})
```

```
Part (g): Our RandomForest for titanic
Cross validation [0.7994012 0.77477477 0.81927711]
First splits [(b'age', 24), (b'pclass', 23), (b'male', 23), (b'S', 22), (b'parch',
20), (b'fare', 20), (b'C', 20), (b'female', 18), (b'sibsp', 15), (b'Q', 15)]
RandomForest(m=None, n=200, params={'max_depth': 5, 'min_samples_leaf': 10})
```

```
Part (i, j): Our RandomForest for titanic
Cross validation [0.81137725 0.75375375 0.74698795]
First splits [(b'age', 24), (b'pclass', 23), (b'male', 23), (b'S', 22), (b'parch',
20), (b'fare', 20), (b'C', 20), (b'female', 18), (b'sibsp', 15), (b'Q', 15)]
BoostedRandomForest(m=None, n=200,
                    params={'max_depth': 5, 'min_samples_leaf': 10})
[1, 1, 0, 0, 0, 1, 0, 0, 0, 1, 1, 0, 0, 1, 1, 0, 0, 1, 0, 1, 1, 0, 0, 1,
1, 0, 0, 1, 0, 0, 0, 0, 0, 0, 1, 1, 0, 0, 1, 1, 0, 0, 1, 1, 1, 1, 0,
0, 0, 1, 0, 0, 0, 1, 1, 0, 0, 1, 1, 0, 1, 0, 1, 0, 0, 0, 0, 0, 0, 0, 0, 1,
0, 0, 1, 0, 0, 0, 1, 0, 1, 0, 0, 1, 0, 0, 1, 0, 0, 0, 0, 1, 0, 0, 0, 1, 1, 1,
1, 0, 0, 1, 0, 0, 1, 1, 1, 0, 1, 0, 1, 1, 0, 0, 0, 0, 0, 1, 0, 0, 0, 1, 0, 0, 1,
1, 0, 0, 1, 0, 0, 1, 1, 1, 0, 1, 0, 1, 1, 0, 1, 0, 0, 0, 0, 0, 1, 0, 0, 0, 1, 0, 0, 1,
0, 0, 1, 1, 0, 0, 0, 0, 0, 1, 0, 0, 0, 0, 1, 1, 0, 0, 0, 0, 1, 1, 0, 0, 0, 1, 0, 0, 0, 1,
1, 0, 0, 0, 0, 1, 0, 0, 1, 0, 0, 1, 0, 1, 0, 0, 1, 0, 0, 1, 0, 0, 1, 0, 0, 0, 1, 0, 0, 0, 1,
1, 0, 0, 0, 1, 1, 0, 0, 0, 1, 0, 0, 0, 0, 1, 0, 0, 0, 1, 0, 0, 0, 1, 0, 0, 0, 1, 0, 0, 0, 1,
1, 0, 0, 0, 1, 0, 0, 0, 0, 1, 0, 0, 0, 0, 1, 0, 0, 0, 1, 0, 0, 0, 0, 1, 0, 0, 0, 1, 0, 0, 0, 1,
1, 0, 0, 0, 1, 1, 0, 0, 0, 1, 0, 0, 0, 0, 1, 0, 0, 0, 1, 0, 0, 0, 0, 1, 0, 0, 0, 0, 1, 0, 0, 0, 1]
```

Spam Dataset

```
In [30]: if __name__ == "__main__":
    main(dataset="spam")
```

```
Features: ['pain', 'private', 'bank', 'money', 'drug', 'spam', 'prescription', 'creative', 'height', 'featured', 'differ', 'width', 'other', 'energy', 'business', 'message', 'volumes', 'revision', 'path', 'meter', 'memo', 'planning', 'pleased', 'record', 'out', 'semicolon', 'dollar', 'sharp', 'exclamation', 'parentheses', 'square_bracket', 'ampersand']  
Train/test size: (5172, 32) (5857, 32)
```

Part 0: constant classifier
Accuracy 0.7099767981438515

Part (c): sklearn's decision tree
Cross validation [0.80104408 0.80800464 0.77900232]

```
Part (e): Our Bagged Tree for spam  
Cross validation [0.80220418 0.80916473 0.78538283]  
First splits [('exclamation', 200)]  
BaggedTrees(n=200, params={'max_depth': 5, 'min_samples_leaf': 10})
```

```
Part (g): Our RandomForest for spam
Cross validation [0.79988399 0.80858469 0.78480278]
First splits [('money', 10), ('drug', 10), ('path', 10), ('semicolon', 10), ('parenthesis', 10), ('other', 9), ('memo', 9), ('featured', 8), ('pain', 7), ('private', 7), ('bank', 7), ('prescription', 7), ('differ', 7), ('business', 7), ('volumes', 7), ('meter', 7), ('ampersand', 7), ('creative', 6), ('record', 6), ('out', 6), ('dollar', 6), ('exclamation', 6), ('width', 5), ('square_bracket', 5), ('energy', 4), ('revision', 4), ('pleased', 4), ('sharp', 4), ('spam', 2), ('height', 1), ('message', 1), ('planning', 1)]
RandomForest(m=None, n=200, params={'max_depth': 5, 'min_samples_leaf': 10})
```

h)

Part H Observations:

Qualitatively describe what this algorithm is doing. What does it mean when $a_i < 0$, how does the algorithm handle such trees?

Boosting:

Fit:

- We determine the weighted error of the current iteration decision tree. a_i is computed based off of this weighted error in accordance with AdaBoost algorithm.
- a_i is always ≥ 0 . The sign assigned to a_i depends on whether or not its respective data point was correctly classified by the current iteration

For data points that are correctly classified we lower the precedence of said data point by reducing its weight, w_i , by a factor of $\exp(-a_i)$. Vice a versa for incorrectly classified data points, it increases the weight w_i on the data point by $\exp(+a_i)$. The new weights on the data points are used in the next iteration.

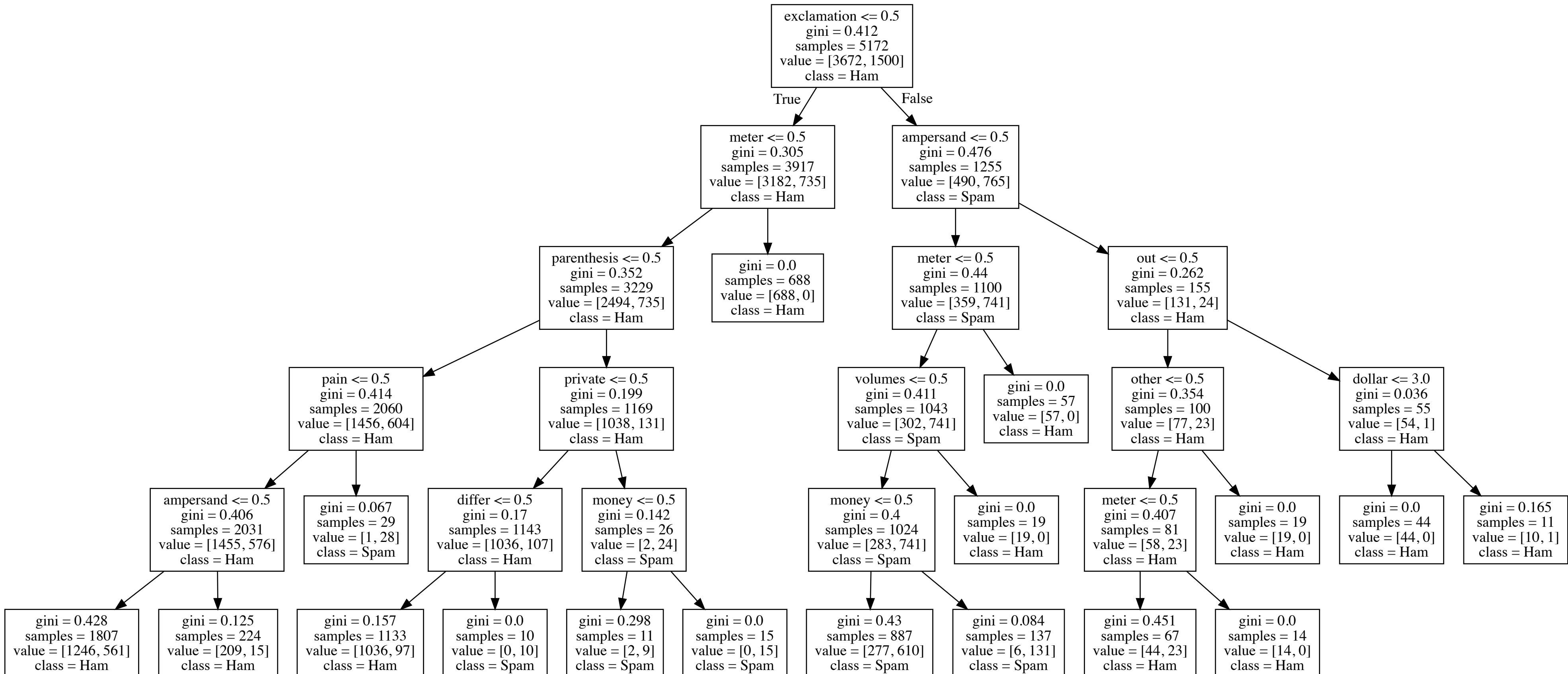
Prediction:

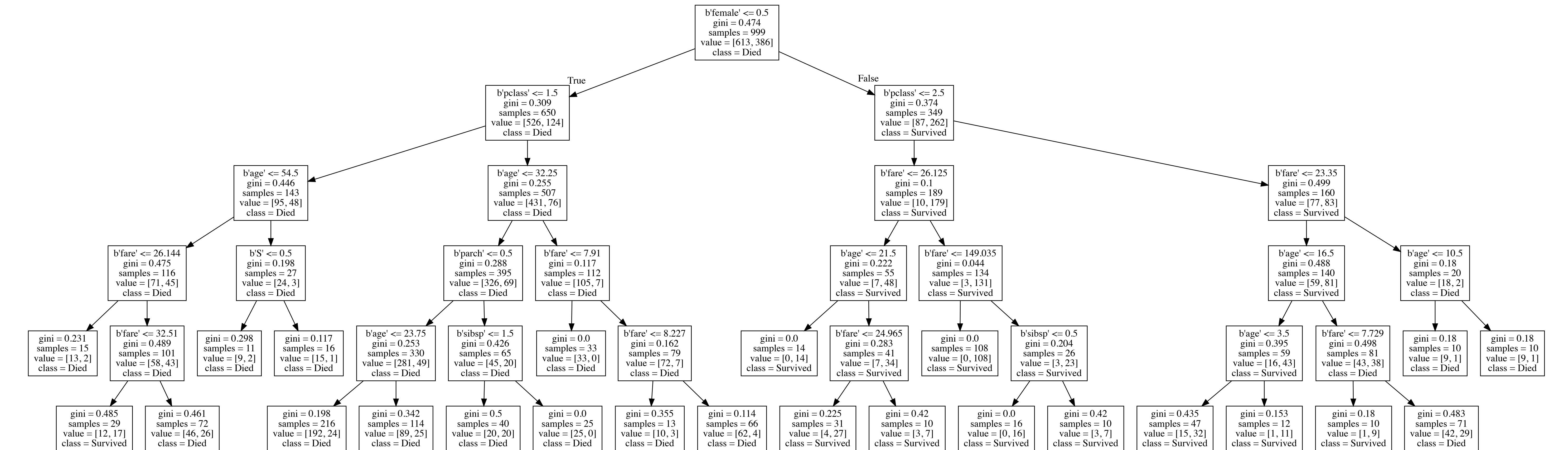
The performance of the classification in the current iteration is determined by the score.

Our algorithm is repeated until we have M trees and have fit and predicted our classification M times, with updated weights at each iteration.

i)

- Most challenging data to classify with decision trees are : Data with high variance, as for decision trees small variations in the data might result in a completely different tree being generated. Furthermore if some classes of data dominate than decision tree learners can create biased trees.





We see that boosted trees perform best and this is due to the fact they iteratively learn and can generate better future iterations. Boosting is based on weak learners with high bias and low variance. Random forests are based on low bias and high variance. Bias in our case performs better