

1. (a)

I did this homework with Luning Zhao.

We work on our homework separately, but we discuss when meeting problems.

Homework is fine.

(b) 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.

Siya Jia

$$\begin{aligned}
 J_\lambda(\vec{w}) &= \frac{1}{2} (\vec{y} - \vec{x}\vec{w})^T (\vec{y} - \vec{x}\vec{w}) + \lambda \|\vec{w}\|^2 \\
 &= \frac{1}{2} \|\vec{y}\|^2 - \vec{y}^T \vec{x} \vec{w} + \vec{w}^T \vec{x}^T \vec{x} \vec{w} + \lambda \|\vec{w}\|^2 \\
 &= \frac{1}{2} \|\vec{y}\|^2 - \sum_{i=1}^d \vec{y}^T \vec{x}_i w_i + \|\vec{w}\|^2 + \lambda \|\vec{w}\|^2 \\
 &= \frac{1}{2} \|\vec{y}\|^2 + \sum_{i=1}^d (w_i^2 + \lambda |w_i| - \vec{y}^T \vec{x}_i w_i) \\
 &= g(\vec{y}) + \sum_{i=1}^d f(\vec{x}_i, \vec{y}, w_i, \lambda) \quad \text{where } g(\vec{y}) = \frac{1}{2} \|\vec{y}\|^2 \\
 &\qquad\qquad\qquad f(\vec{x}_i, \vec{y}, w_i, \lambda) = w_i^2 + \lambda |w_i| - \vec{y}^T \vec{x}_i w_i
 \end{aligned}$$

$$\text{So } \hat{\vec{w}} = \arg \min_{\vec{w}} \{J_\lambda(\vec{w})\} = \arg \min_{\vec{w}} \sum_{i=1}^d f(\vec{x}_i, \vec{y}, w_i, \lambda)$$

so the joint optimization can be decomposed into individual optimizations.

$$\text{so we can learn } w_i \text{ by } \arg \min_{w_i} f(\vec{x}_i, \vec{y}, w_i, \lambda)$$

$$\text{(b) If } w_i > 0, \quad f(\vec{x}_i, \vec{y}, w_i, \lambda) = w_i^2 + \lambda w_i - \vec{y}^T \vec{x}_i w_i$$

$$\text{f is minimized when } \hat{w}_i = -\frac{\lambda - \vec{y}^T \vec{x}_i}{2} = \frac{-\lambda + \vec{y}^T \vec{x}_i}{2} \quad (\text{when } \lambda < \vec{y}^T \vec{x}_i)$$

$$\text{(c) If } w_i < 0, \quad f(\vec{x}_i, \vec{y}, w_i, \lambda) = w_i^2 - \lambda w_i - \vec{y}^T \vec{x}_i w_i$$

$$\text{f is minimized when } \hat{w}_i = -\frac{-\lambda - \vec{y}^T \vec{x}_i}{2} = \frac{\lambda + \vec{y}^T \vec{x}_i}{2} \quad (\text{when } \lambda < -\vec{y}^T \vec{x}_i)$$

$$\text{(d) when } \lambda \in [-\vec{y}^T \vec{x}_i, \vec{y}^T \vec{x}_i], \text{ neither (b) or (c) works, so } \hat{w}_i = 0$$

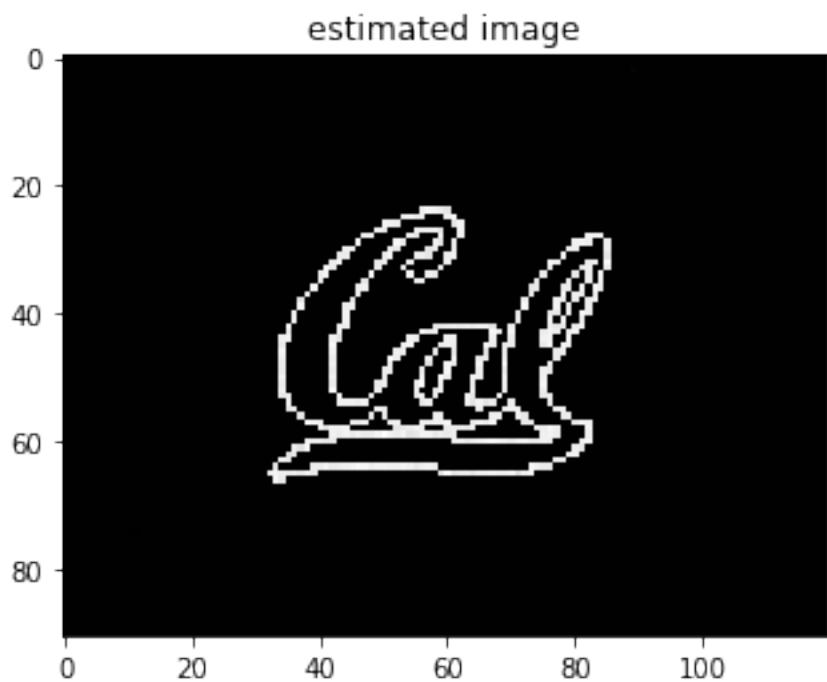
$$\begin{aligned}
 \text{(e) Similarly, } f(\vec{x}_i, \vec{y}, w_i, \lambda) &= w_i^2 + \lambda w_i^2 - \vec{y}^T \vec{x}_i w_i \\
 &= (1+\lambda) w_i^2 - \vec{y}^T \vec{x}_i w_i
 \end{aligned}$$

$$f(\vec{x}_i, \vec{y}, w_i, \lambda) \text{ is minimized when } w_i = -\frac{\vec{y}^T \vec{x}_i}{2(1+\lambda)} = \frac{\vec{y}^T \vec{x}_i}{2(1+\lambda)}$$

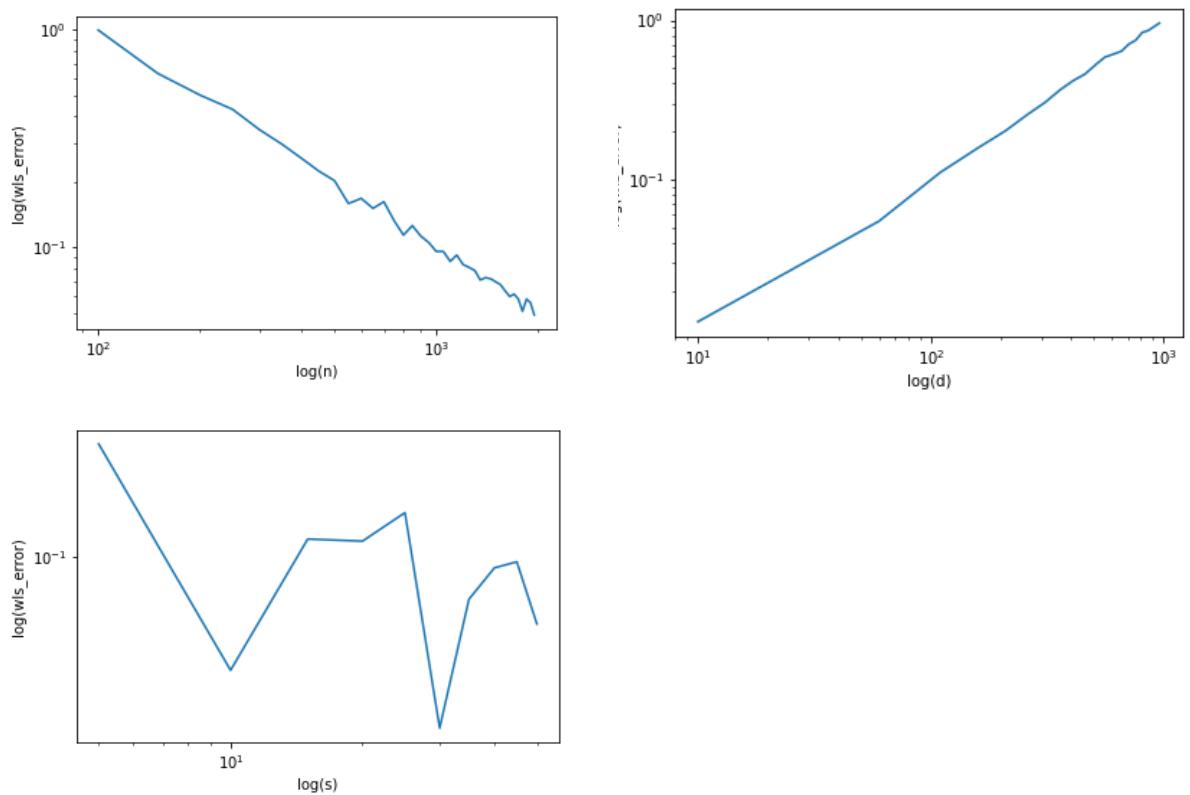
$\hat{w}_i = 0$ only when $\vec{y}^T \vec{x}_i = 0$, which is very hard to obtain compare with (d).

So L1 norm promotes sparsity.

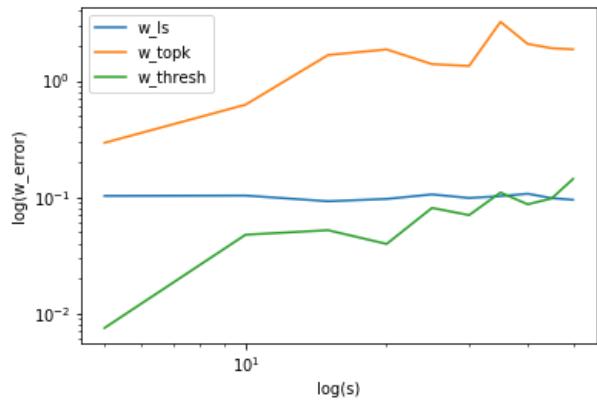
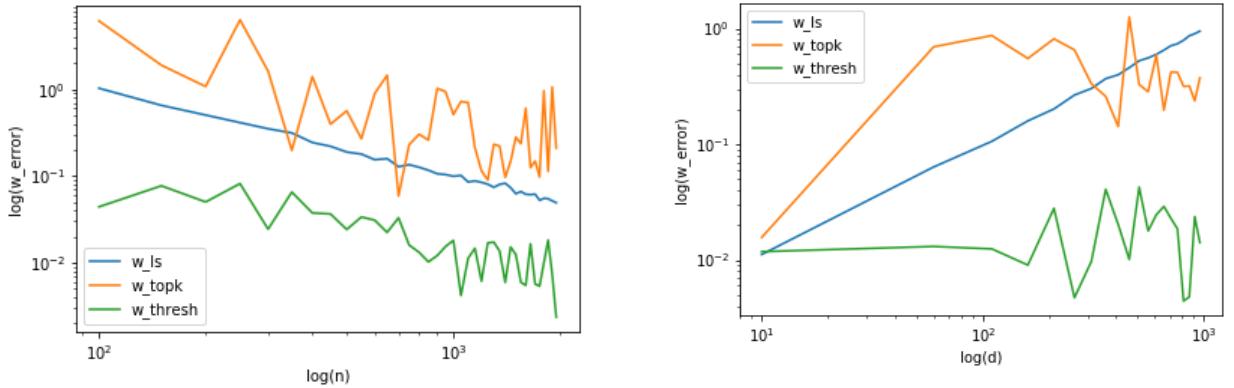
2 (f) let lambda = 0.00001, and run LASSO(imDims, measurements, X, 0.00001).



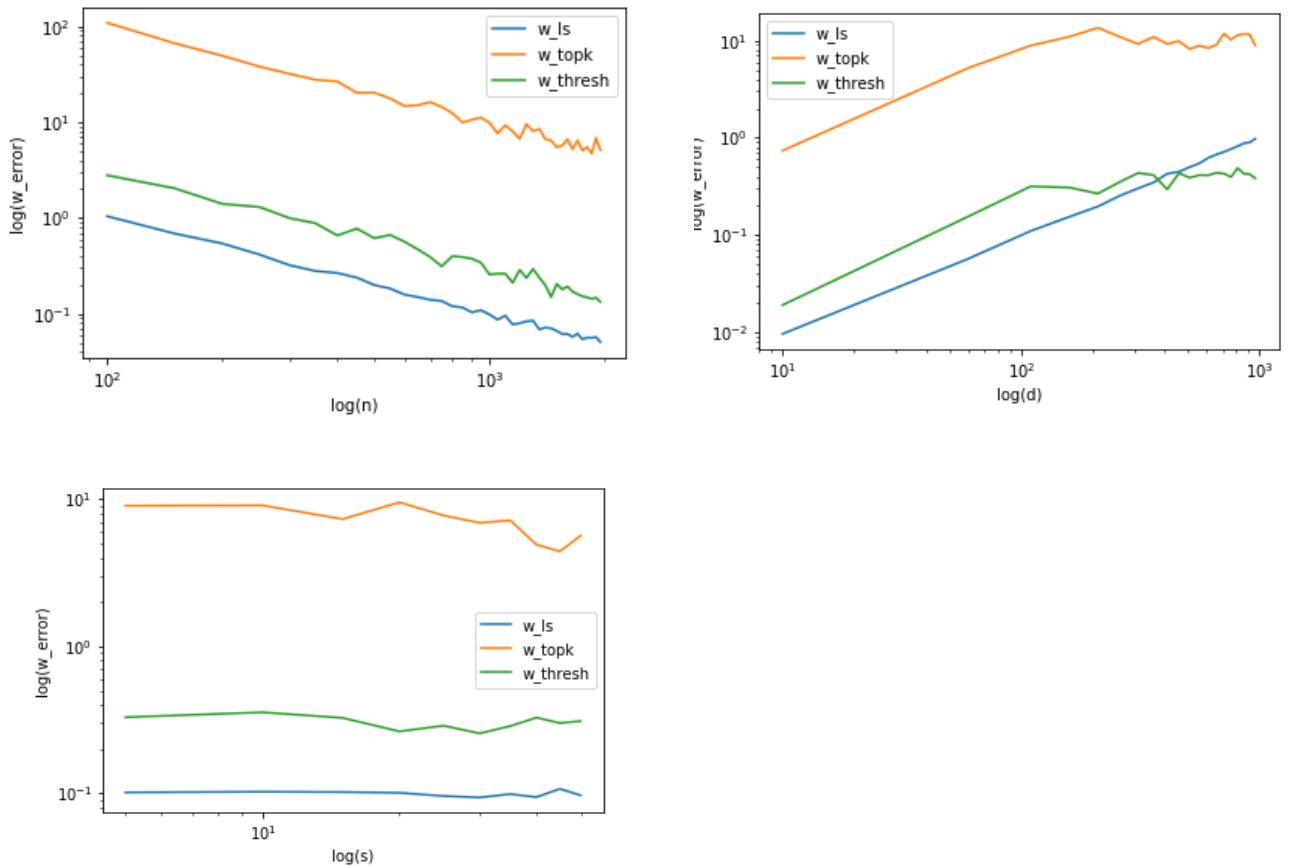
3 (a)



(b)



(c)



For sparse linear model, threshold/topk sparsity estimator has smaller variance, and its bias is not very large, so it has small total errors as shown in (b).

But for non-sparse linear model, threshold/topk sparsity estimator has very large bias, which results in large total errors as shown in (c).

Topk model has so large bias that it always has larger error than OLS, but threshold sparsity model works even better than OLS for sparse linear model.

$$\begin{aligned}
 3.(d) \quad & \Pr \left\{ \max_i |z_i| \geq 20\sqrt{\log d} \right\} \\
 & \geq \sum_{i=1}^d \Pr \left\{ |z_i| \geq 20\sqrt{\log d} \right\} \\
 & \geq d e^{-\frac{40^2 \log d}{20^2}} \\
 & = d \cdot (e^{\log d})^{-2} \\
 & = \frac{1}{d}
 \end{aligned}$$

So we prove that $\Pr \left\{ \max_i |z_i| \geq 20\sqrt{\log d} \right\} \leq \frac{1}{d}$.

$$(e) \quad \hat{w}_{LS} = X^T y = X^T(Xw^* + z) = X^T X w^* + X^T z$$

Because X is orthonormal, $X^T X = I$

$$\text{cov}(X^T z) = X \text{cov}(z) X^T = \sigma^2 X X^T = \sigma^2 I$$

$$\text{so } \hat{w}_{LS} = w^* + z' \text{ where } z' \sim N(0, \sigma^2 I_n)$$

$\hat{w}_{top}(s) = T_s(\hat{w}_{LS})$, so $\hat{w}_{top}(s)$ returns the top s entries of the vector $w^* + z'$

(f) Because \hat{w}^* is s sparsity, and $\hat{w}_{top}(s)$ also only keep s entries and also s sparsity, therefore $\hat{e} = \hat{w}_{top}(s) - \hat{w}^*$ is at most $2s$ -sparsity.

$$(g) \quad |e_i| = |\hat{w}_{top}(s)_i - w^*_i|$$

$$\text{If } \hat{w}_{top}(s)_i = 0, \quad w^*_i = 0, \quad |e_i| = 0 \leq 40\sqrt{\log d}$$

$$\text{If } \hat{w}_{top}(s)_i \neq 0, \quad w^*_i = 0, \quad |e_i| = |\hat{w}_{top}(s)_i| = |z'_i + w^*_i| = |z'_i| \leq 20\sqrt{\log d} \leq 40\sqrt{\log d}$$

$$\text{If } \hat{w}_{top}(s)_i \neq 0, \quad w^*_i \neq 0, \quad |e_i| = |\hat{w}_{top}(s)_i - w^*_i| = |z'_i| \leq 20\sqrt{\log d} \leq 40\sqrt{\log d}.$$

$$\text{If } \hat{w}_{top}(s)_i = 0, \quad w^*_i \neq 0, \quad \text{meaning there is } j \text{ so that } |w^*_i + z'_i| \leq |z'_j + w^*_j| = |z'_j|$$

$$\text{so } |e_i| = |w^*_i| \leq 2|z'_i| = 40\sqrt{\log d}.$$

So we prove that $|e_i| \leq 40\sqrt{\log d}$

$$(h) \|\hat{w}_{top}(s) - w^*\|^2 = \sum_{i=1}^d \|e_i\|^2 \leq 2s(4\sigma \log d)^2 = 32\sigma^2 s \log d$$

$$\text{So } \Pr \{ \|\hat{w}_{top}(s) - w^*\|^2 \leq 32\sigma^2 s \log d \}$$

$$= \Pr \{ \max_i |e_i| \leq 4\sigma \log d \} = \Pr \{ \max_i |z_i| \leq 2\sigma \log d \}$$

$$\geq 1 - \Pr \{ \max_i |z_i| \geq 2\sigma \log d \}$$

$$\geq 1 - \frac{1}{d}$$

$$(i) \frac{1}{n} \|X(\hat{w}_{top}(s) - w^*)\|^2$$

$$= \frac{1}{n} (X(\hat{w}_{top} - w^*))^T (X(\hat{w}_{top} - w^*)) \rightarrow \bar{a}$$

$$= \frac{1}{n} \bar{a}^T \bar{a}$$

$$= \frac{1}{n} \bar{a}^T \bar{a}$$

$$= \frac{1}{n} \|\hat{w}_{top}(s) - w^*\|^2 \leq \frac{32\sigma^2 s \log d}{n}$$

(j) sparsity predictor gives variance of $32\sigma^2 s \log d / n$

and OLS gives variance of $\sigma^2 d / n$

To make the first term smaller than the 2nd term $\Rightarrow 32\sigma^2 s \log d / n \leq \sigma^2 d / n$

$$\Rightarrow s \leq \frac{d}{32 \log d}$$

so by leveraging the sparsity assumption, we have smaller variance.

(k) OLS gives variance of $\sigma^2 s / n$

sparse $\hat{w}_{top}(s)$ gives variance of $32\sigma^2 s \log d / n = \frac{\sigma^2 s}{n} \cdot 32 \log d > \frac{\sigma^2 s}{n}$

So when knowing important s features, OLS gives smaller variance

compared with sparse $\hat{w}_{top}(s)$ above.

So the price of not knowing important features will introduce larger variance ($\frac{d}{s}$)

4 (a)

```
@staticmethod
def information_gain(X, y, thresh):
    # TODO implement information gain function
    # calculate entropy before cut
    p0 = len(np.where(y==0)[0])/len(y)
    p1 = 1- p0
    entropy_bef = p0 * np.log(p0) + p1 * np.log(p1)

    # calculate entropy after cut
    # calculate for X> thresh
    idx_large = np.where(X>thresh)[0]
    y_large = y[idx_large]
    p0 = len(np.where(y_large==0)[0])/len(y_large)
    p1 = 1- p0
    entropy_aft_large = p0 * np.log(p0) + p1 * np.log(p1)

    idx_small = np.where(X<=thresh)[0]
    y_small = y[idx_small]
    p0 = len(np.where(y_small==0)[0])/len(y_small)
    p1 = 1 - p0
    entropy_aft_small = p0 * np.log(p0) + p1 * np.log(p1)

    # calculate the gain
    w_large = len(idx_large)/len(y)
    w_small = len(idx_small)/len(y)
    gain = entropy_bef - w_large*entropy_aft_large - w_small*entropy_aft_small
    return gain
```

```
@staticmethod
def gini_impurity(X, y, thresh):
    # TODO implement gini_impurity function
    # calculate G before cut
    p0 = len(np.where(y==0)[0])/len(y)
    p1 = 1- p0
    G_bef = 1 - p0**2 - p1**2

    # calculate G after cut
    # calculate for X> thresh
    idx_large = np.where(X>thresh)[0]
    y_large = y[idx_large]
    p0 = len(np.where(y_large==0)[0])/len(y_large)
    p1 = 1- p0
    G_aft_large = 1- p0**2 - p1**2

    idx_small = np.where(X<=thresh)[0]
    y_small = y[idx_small]
    p0 = len(np.where(y_small==0)[0])/len(y_small)
    p1 = 1 - p0
    G_aft_small = p0 * np.log(p0) + p1 * np.log(p1)

    # calculate the gain
    w_large = len(idx_large)/len(y)
    w_small = len(idx_small)/len(y)
    gain = G_bef - w_large*G_aft_large - w_small*G_aft_small

    return gain
```

4 (b)

- Miss class labels: ignore this data
- Features are not numerical values: divide into different groups where there are more than 10 data in this group. For each group, assign 1 to the data if it's the same as this group feature; assign 0 if it's different.
- Data miss features: For the hot columns, first assign -1 to it, so it's different than all other features, meaning it's 0 in every group features. For non-hot columns, use the mode of other data for missing features.

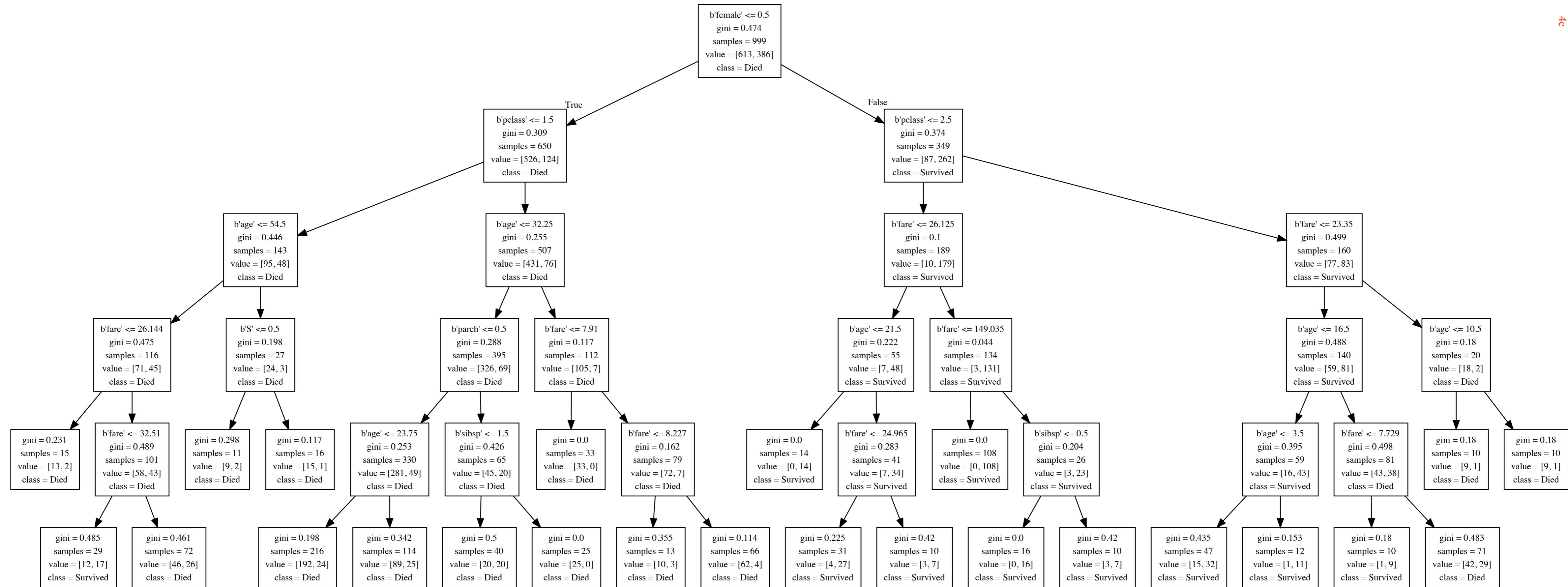
4 (d)

```
class BaggedTrees(BaseEstimator, ClassifierMixin):  
    def __init__(self, params=None, n=200):  
        if params is None:  
            params = {}  
        self.params = params  
        self.n = n  
        self.decision_trees = [  
            sklearn.tree.DecisionTreeClassifier(random_state=i, **self.params)  
            for i in range(self.n)  
        ]  
  
    def fit(self, X, y):  
        Ndata = len(y)  
        # TODO implement function  
        for dt in self.decision_trees:  
            idx = np.random.randint(0, Ndata-1, size=Ndata)  
            Xi = X[idx]  
            yi = y[idx]  
            dt.fit(Xi, yi)  
  
    def predict(self, X):  
        # TODO implement function  
        y = np.zeros(X.shape[0])  
        for dt in self.decision_trees:  
            y += dt.predict(X)  
        y /= self.n  
        return np.round(y)
```

4 (e)

For Titanic: The most common splits made at root: First splits [(b'male', 101), (b'female', 99)]
For Spam: First splits [('exclamation', 200)]

```
# TODO implement and evaluate parts c-h  
bt = BaggedTrees()  
bt.fit(X,y)  
evaluate(bt)
```



4 (f)

```
class RandomForest(BaggedTrees):
    def __init__(self, params=None, n=200, m=1):
        if params is None:
            params = {}
        self.params = params
        self.n = n
        self.decision_trees = [
            sklearn.tree.DecisionTreeClassifier(random_state=i, max_features=m, **self.params)
            for i in range(self.n)
        ]
    # TODO implement function
    pass
```

4 (g)

```
bt = RandomForest(m=6)
bt.fit(X,y)
evaluate(bt)
```

For Titanic: First splits [(b'male', 71), (b'female', 66)
For spam: First splits [('exclamation', 38), ('money', 31)]

4(h)

```
class BoostedRandomForest(RandomForest):
    def fit(self, X, y):
        self.w = np.ones(X.shape[0]) / X.shape[0] # Weights on data
        self.a = np.zeros(self.n) # Weights on decision trees
        w = self.w

        # TODO implement function
        # record the number of data points
        ndata = len(y)

        # loop over decision trees
        for i in range(self.n):
            # sample randomly with weights
            idx = np.random.choice(ndata, ndata, p=w)
            Xi = X[idx]
            yi = y[idx]

            # train tree and predict
            self.decision_trees[i].fit(Xi, yi)
            yp = self.decision_trees[i].predict(X)

            # calculate ej and aj
            idx_wrong = np.where(yp!=y)[0]
            ej = np.sum(w[idx_wrong])/np.sum(w)
            aj = 0.5 * np.log((1-ej)/ej)
            self.a[i] = aj

            # update weight
            w = w * np.exp(-aj)
            w[idx_wrong] = w[idx_wrong] * np.exp(aj)

            # normalize w
            w /= np.sum(w)

        self.w = w
        return self

    def predict(self, X):
        # TODO implement function
        # initialize z for c=0/1
        z0 = np.zeros(len(X))
        z1 = np.zeros(len(X))

        # loop over decision trees
        for i in range(self.n):
            yp = self.decision_trees[i].predict(X)

            # calculate z
            idx0 = np.where(yp==0)[0]
            z0[idx0] += self.a[i]
            idx1 = np.where(yp==1)[0]
            z1[idx1] += self.a[i]

        labels = np.zeros(len(X))
        idx = np.where(z0<z1)[0]
        labels[idx] = 1
        return labels
```

Trees are being weighted by a . When trees makes 0 wrong prediction, $e=0$, corresponding to $a=\inf$; When trees makes all wrong prediction, $e=1$, corresponding to $a=-\inf$. So we give more weights to trees which makes more correct predictions.

This algorithm gives more weights to data that are not predicted correctly. At each step, if a data is not predicted correctly, it will be given weights of $\exp(a_j)$; On the contrary, if a data is predicted correctly, it will be given weights of $\exp(-a_j)$. So next time, those data who are not predicted correctly will have more probability when sampling.

$ai < 0$ means $ej > 0.5$, meaning this tree has large errors. but ada boots always assume $ej < 0.5$, so those trees will flip their prediction.

4(i)

For titanic: First splits [(b'fare', 42), (b'age', 37)]

For spam: First splits [('meter', 32), ('volumes', 28)]

The data that is hardest to predict is with large weight, like:

meaning they has no other keys except & or (

The data that is easiest to predict is with small weight, like:

meaning they have "meter" in it

4 (j)

For titanic:

For titanic

Cross validation [0.79640719 0.76576577 0.79819277]

accuracy: 0.822823

Bagged trees

Cross validation [0.79640719 0.75675676 0.78915663]

accuracy: 0.973974

Random forests

Cross validation [0.79041916 0.75675676 0.79216867]

accuracy: 0.973974

Ada Boost

Cross-validation [0.80538922 0.74774775 0.77108434]

accuracy: 0.973974

For spam
Single tree
Cross validation [0.80104408 0.80800464 0.77900232]
accuracy: 0.809938

Bagged trees
Cross validation [0.82192575 0.81786543 0.78016241]
accuracy: 0.885924

Random forests
Cross validation [0.82134571 0.81612529 0.78132251]
accuracy: 0.885924

Ada Boost
Cross validation [0.82656613 0.82424594 0.78828306]
accuracy: 0.885924

See submission for test data prediction

(k)
see submission

5. Q: In the m -th step of AdaBoost, we want to minimize

$$\underset{\alpha_m, G_m}{\operatorname{argmin}} \sum_{i=1}^n L(y_i, F_{m-1}(x_i) + \alpha_m G_m(x_i))$$

in which L is exponential loss: $L(y, h(x)) = e^{-y h(x)}$
 derive the expression for α_m .

A:

$$\alpha_m, G_m = \underset{\alpha, G}{\operatorname{argmin}} \sum_{i=1}^n L(y_i, F_{m-1}(x_i) + \alpha G(x_i))$$

$$\text{Let } W_m^i = e^{-y_i F_{m-1}(x_i)}$$

$$\alpha_m, G_m = \underset{\alpha, G}{\operatorname{argmin}} \sum W_m^i e^{-y_i \alpha G(x_i)}$$

$$= \underset{\alpha, G}{\operatorname{argmin}} \left(\sum_{y_i = G(x_i)} W_m^i e^{-\alpha} + \sum_{y_i \neq G(x_i)} W_m^i e^{\alpha} \right)$$

$$= \underset{\alpha, G}{\operatorname{argmin}} \sum_{i=1}^n W_m^i e^{-\alpha} + \sum_{y_i \neq G(x_i)} W_m^i (e^\alpha - e^{-\alpha})$$

$$= \sum_{i=1}^n W_m^i \underset{\alpha, G}{\operatorname{argmin}} e^{-\alpha} + \cancel{\sum_{y_i \neq G(x_i)} \frac{W_m^i}{\sum_{y_j \neq G(x_j)} W_m^j} (e^\alpha - e^{-\alpha})}$$

$$\text{let } \epsilon_m = \frac{\sum_{y_i \neq G(x_i)} W_m^i}{\sum W_m^i}$$

$$\alpha_m, G_m = \underset{\alpha, G}{\operatorname{argmin}} e^{-\alpha} + \epsilon_m (e^\alpha - e^{-\alpha}) = F(\alpha)$$

$$\frac{\partial F(\alpha)}{\partial \alpha} = 0 \Rightarrow \alpha = \frac{1}{2} \ln \left(\frac{1 - \epsilon_m}{\epsilon_m} \right)$$