

Q1

a) Myself

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.

Q5

What is the main downfall of decision trees and how do random forests solve it?

OVERFITTING.

That can be seen in the ^{training} accuracy results of a single decision tree vs its validation accuracy.

Q2

$$\hat{w} = \arg \min_{\underline{w}} \left\{ \frac{1}{2} \|y - Xw\|_2^2 + \lambda \|w\|_1 \right\}$$

$$\begin{aligned} J_{\lambda}(w) &= \frac{1}{2} \|y\|^2 + \frac{1}{2} \|Xw\|^2 - y^T Xw + \lambda \|w\|_1 \\ &= \frac{1}{2} \|y\|^2 + \frac{1}{2} w^T X^T X w - y^T Xw + \lambda \|w\|_1 \\ &= \frac{1}{2} \|y\|^2 + \frac{1}{2} w^T \left(\sum_{i,j} X_i^T X_j \right) w - y^T Xw + \lambda \|w\|_1 \\ &= \frac{1}{2} \|y\|^2 + \frac{n}{2} w^T w - y^T Xw + \lambda \|w\|_1 \end{aligned}$$

Uncorrelated X_i 's

$$\rightarrow E[X_i X_j^T] = E[X_i] E[X_j^T] = 0 \text{ for } i \neq j$$

~~$$J_{\lambda}(w) = \frac{1}{2} \|y\|^2 + \frac{1}{2} \sum_i w_i^2 \|X_i\|^2 + \lambda \|w\|_1 - y^T Xw$$~~

$$J_{\lambda}(w) = \frac{1}{2} \|y\|^2 + \frac{n}{2} \|w\|_2^2 + \lambda \|w\|_1 - y^T Xw$$

Sum All together

$$\begin{aligned} \frac{\partial}{\partial w} (J_{\lambda}(w)) &= n \underline{w} + \frac{\partial}{\partial w} \|w\|_1 - (X^T y) = 0 \\ &= n \underline{w} + \text{sgn}(w) - X^T y = 0 \\ \underline{w} &= \frac{1}{n} [X^T y - \text{sgn}(w)(\lambda)] \end{aligned}$$

Solve for w_i Separately

$$\begin{aligned} \frac{\partial}{\partial w_i} J_{\lambda}(w_i) &= \frac{\partial}{\partial w_i} \left[\frac{n}{2} w_i^2 - (y^T X_i) w_i + \lambda |w_i| \right] \\ &= n w_i - y^T X_i + \lambda \text{sgn}(w_i) = 0 \end{aligned}$$

$$w_i = \frac{1}{n} [y^T X_i - \lambda \text{sgn}(w_i)]$$

positive

$$b) \hat{w}_i = \frac{1}{n} [y_i^T x_i + \lambda] = + \frac{y_i^T x_i}{n} - \frac{\lambda}{n}$$

$$c) \hat{w}_i = \frac{1}{n} [\underbrace{y_i^T x_i}_{- \text{negative}} + \lambda] = + \frac{y_i^T x_i}{n} + \frac{\lambda}{n}$$

$$d) \text{ if } |y_i^T x_i| < |\lambda| \text{ then } \hat{w}_i = 0$$

$$e) \hat{Q} = (X^T X + \lambda I)^{-1} X^T y$$

$$\hat{w}_i = \frac{1}{n + \lambda} (x_i^T y_i) \quad x_i^T y_i = 0 \text{ for } \hat{w}_i = 0$$

$$f) k - 5 = \lambda$$

Q3

d) $z_i \sim N(0, \sigma^2) \quad \Pr(|z_i| \geq t) \leq e^{-t^2/2\sigma^2}$

$$\Pr(|z_i| \geq 2\sigma\sqrt{\log d}) \leq e^{-\frac{(2\sigma\sqrt{\log d})^2}{2\sigma^2}} = e^{-\frac{4\sigma^2 \log d}{2\sigma^2}} = e^{-2\log d} = \frac{1}{d^2}$$

$\sum_{\text{all } i \text{ s.t. } |z_i| \geq 2\sigma\sqrt{\log d}} \Pr(|z_i| \geq 2\sigma\sqrt{\log d}) \leq \left(\frac{1}{d^2}\right) \cdot d \leq \frac{1}{d} \Rightarrow \Pr(\max_i |z_i| \geq 2\sigma\sqrt{\log d}) \leq \frac{1}{d}$

e) $w = (X^T X)^{-1} X^T (y) = (X^T X)^{-1} X^T (y^* + z) = w^* + (X^T X)^{-1} X^T z$
 $= w^* + z'$

$z' \sim (0, \sigma^2 I)$ transformed cov matrix
 because \leftarrow

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

$\hat{w}_{\text{top}}(s) = T_s(\hat{w}) = T_s(w^* + z')$

$\sigma^2 (X^T X)^{-1} X^T X (X^T X)^{-1} I$
 I because has orthogonal columns

f) $e = \hat{w}_{\text{top}}(s) - w^*$ (assuming w^* is s-sparse)

Max error occurs when $\hat{w}_{\text{top}}(s)$ predicts s-originally 0 entries in w^* to have some value.

$\max(\hat{w}_{\text{top}}(s) - w^*) \Rightarrow$ it will give s-wrongly predicted non-zero entries

2- negative s entries of original weights

$\max(\text{Sparsity}(e)) = 2s$

g) Given $\max_i |z_i| \leq 2\sigma\sqrt{\log d}$

$w_{\text{top}}(s) = T_s(w^* + z_i)$

Consider max entry

$e = T_s(w^* + z_i) - w^*$

$e_i = w_i^* + z_i' - w_i^* = z_i \quad \text{or} \quad -w_i^*$

$T_s(w^* + z_i)$

$\left\{ \begin{array}{l} |w_i^* + z_i| \text{ if } i \text{ is top } s \\ 0 \text{ otherwise} \end{array} \right.$

$$e_i = \begin{cases} -z_i & \max |e_i| \leq 26\sqrt{\log d} \\ -W_i^* & \leftarrow \text{Assume } \rightarrow 0 \text{ because we pr} \end{cases}$$

$$(iv) \Pr(\|\hat{w}_{top(s)} - w^*\|^2 \leq 326^2 s \log d)$$

$$\begin{aligned} \text{hence } \sum_i |e_i|^2 &\leq (46\sqrt{\log d})^2 \times 25 \text{ because } |e_i| \neq 0 \text{ only } 25 \text{ times max} \\ &\leq (166^2 \log d) \times 25 \\ &\leq 326^2 \log d \end{aligned}$$

Given that ϵ
 from (d) $\Rightarrow \Pr(\text{not } \epsilon) \leq 1/d$
 $\Rightarrow \Pr(\epsilon) \geq 1 - 1/d$

$$(i) \|XV\|^2 = \|V\|^2 \text{ since } X \text{ has orthonormal columns}$$

$$|x_i| = 1 \quad \forall x_i \in \text{Spec}(X)$$

$$\frac{1}{n} \|X(\hat{w}_{top(s)} - w^*)\| \leq \frac{1}{n} 326^2 s \log d \quad \text{w.h.p. } (1 - 1/d)$$

(j) When Sparsity comes into play

$$\left(\frac{\|K_{top(s)}\|^2}{n}\right) \leq \frac{326^2 s \log d}{n} \quad \text{compared } \frac{\|K_{ms}\|}{n} \leq \sigma^2 \frac{d}{n}$$

w.h.p. $(1 - 1/d)$
 \Rightarrow high probability

$$\Rightarrow \frac{1}{n} \|K_{top(s)}\|^2 \leq \frac{1}{n} \|K_{ms}\|^2$$

$$\frac{326^2 s \log d}{n} \leq \sigma^2 \frac{d}{n}$$

when

$$\boxed{s \leq \frac{d}{326^2 \log d}}$$

(ix) The variance would only be 6^2 because that's the variance of the noise.

there is a factor of $32 \log d$ that we pay for not knowing which features to use.

Q4

b) Missing missing features:

the values are filled with the mode of the data feature

Not Numerical Values:

If the string is common enough, we make a one-hot-feature column for it.

Missing Class label:

Remove it.

c) $\text{exclusion} < 0.5$ for all

g) Spam:

money < 0.5
private < 0.5
bunk < 0.5

More common

Most common
= 'money'

Titanic:

sex < 0.5

class

Sibsp < 0.5

Most common
= Sex

j) When x has very little counts of anything, meaning we don't have much information about x when $w > 1e-15$

When w is not spam it is easy to classify because they mostly only have characters and no keywords of the ones above.

Then

$w < 1e-35$

Bias and Variance of Sparse Linear Regression

In this notebook, you will explore numerically how sparse vectors change the rate at which we can estimate the underlying model. This corresponds to parts (a), (b), (c) of Homework 12.

First, some setup. We will only be using basic libraries.

```
In [5]: import numpy as np
import matplotlib.pyplot as plt
%matplotlib inline
```

The following functions produce the ground truth matrix $A \in \mathbb{R}^{n \times d}$ (denoted by U since it is unitary), as well as the vector $w^* \in \mathbb{R}^d$ and observations $y \in \mathbb{R}^n$. They have been implemented for you, but it is worth going through the code to observe its limitations.


```

In [6]: def ground_truth(n, d, s):
        """
        Input: Two positive integers n, d. Requires n >= d >= s. If d < s, we let s = d
        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

```

We now implement the estimators that we will simulate. The least squares estimator has already been implemented for you. You will be implementing the top k and threshold estimators in part (b), but it is fine to skip this for now and compile.

```

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

    return arg_range, wls_error, wtopk_error, wthresh_error

```

We are now ready to perform the parts of the question.

Part (a)

As an example, in the following cell, we run the helper function above to return error values of the OLS estimate for various values of n . You are required to:

- 1) Plot the error as a function of n . You may find a log-log plot useful to see the expected behavior.
- 2) Run the helper function to return the error as a function of d and s , and plot your results.

You need to have 3 plots in your answer. Make sure to label axes properly, and to make the plotting visible in general. Feel free to play with the parameters, but ensure that your answer describes your parameter choices. At this point, `s_model` is `True`, since we are only interested in the variance of the model.

```
In [14]: #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)

plt.figure()
plt.plot(nrange, ls_error)
plt.xlabel('Data Point Number')
plt.title('Param = n')
plt.ylabel('Error')

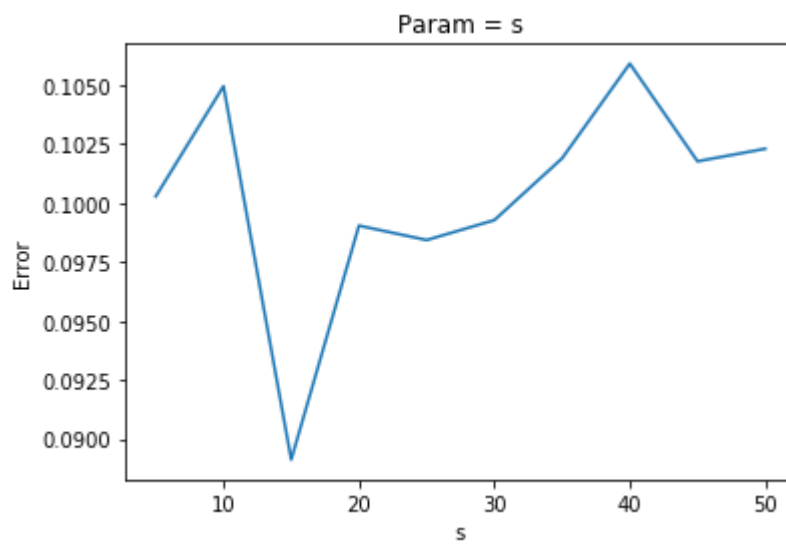
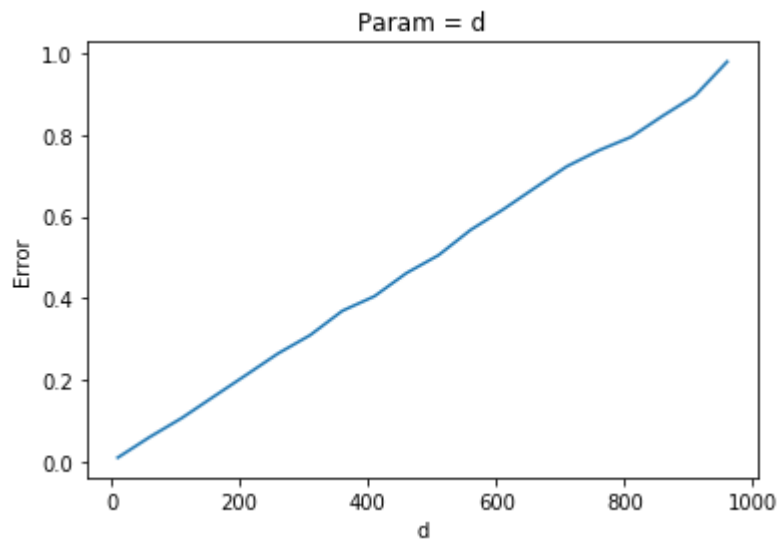
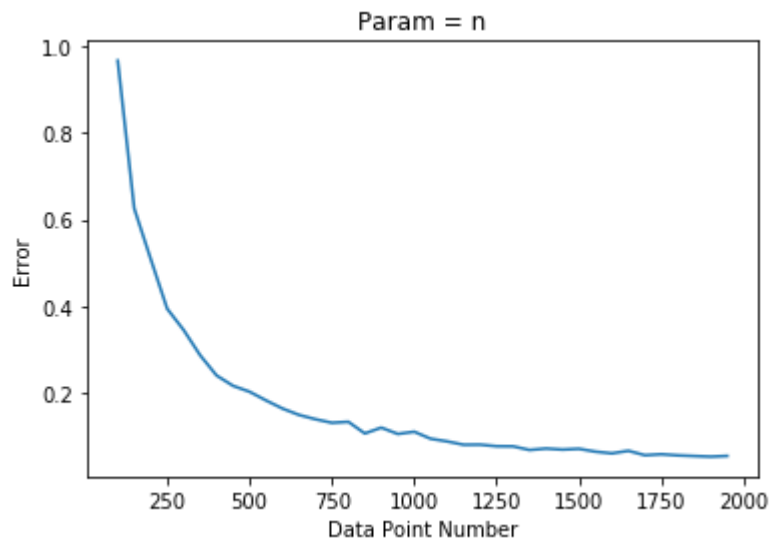
nrange, ls_error, _, _ = error_calc(num_iters=10, param='d', n=1000,
d=100, s=5, s_model=True, true_s=5)

plt.figure()
plt.plot(nrange, ls_error)
plt.xlabel('d')
plt.title('Param = d')
plt.ylabel('Error')

nrange, ls_error, _, _ = error_calc(num_iters=10, param='s', n=1000,
d=100, s=5, s_model=True, true_s=5)

plt.figure()
plt.plot(nrange, ls_error)
plt.xlabel('s')
plt.title('Param = s')
plt.ylabel('Error')
```

Out[14]: Text(0,0.5,'Error')



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.

They are as expected. For LS, the higher the dimensions, the more the bias error. Moreover, the higher the complexity of our Model. S , the sparsity is unrelated to our OLS solver, so as expected, the error is uncorrelated to the sparsity s . Also, as n increases, as expected, the error decreases, as the effect of noise has less effect on our most likely weights. Optimal: $d = 1$, $n = \text{infinity}$, $s = \text{anything}$

Part (b)

Now fill out the functions implementing the sparsity-seeking estimators: thresh, and topk in the above cells. You should be able to test these functions using some straightforward examples.

We will now simulate the error of all the estimators, as a function of n , d , and s . An example of this for n is given below. You must:

- 1) Plot the error of all estimators as a function of n .
- 2) Run the helper function to sweep over d and s , and plot the behavior of all three estimators.

You should report 3 plots here once again. Make sure to make them fully readable.

```
In [34]: nrange, _, wtopk_error, wthresh_error= error_calc(num_iters=10, param=
'n', n=1000, d=100, s=5, s_model=True, true_s=5)

plt.figure()
plt.plot(nrange, wtopk_error, nrange, wthresh_error)
plt.xlabel('Data Point Number')
plt.title('Param = n')
plt.ylabel('Error')
plt.legend({'topk', 'thresh'})

nrange, _, wtopk_error, wthresh_error= error_calc(num_iters=10, param=
'd', n=1000, d=100, s=5, s_model=True, true_s=5)

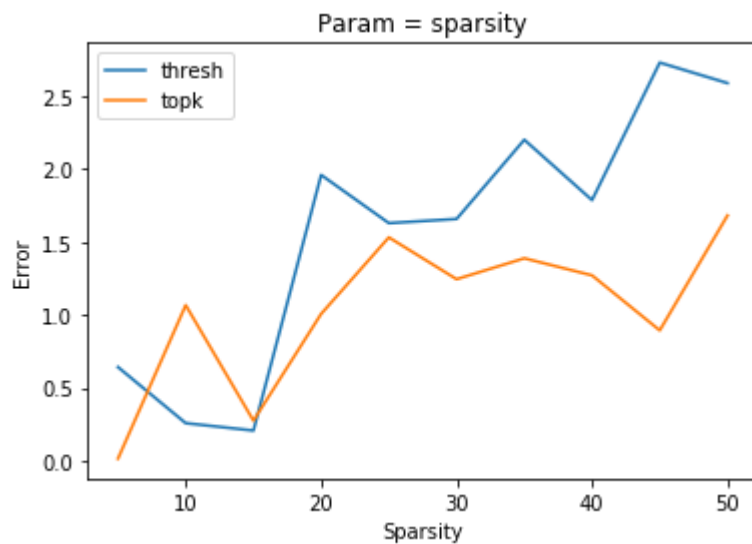
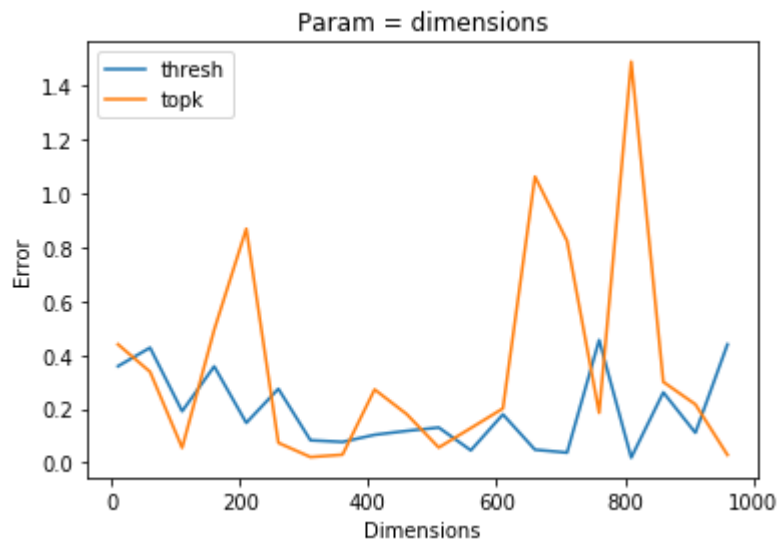
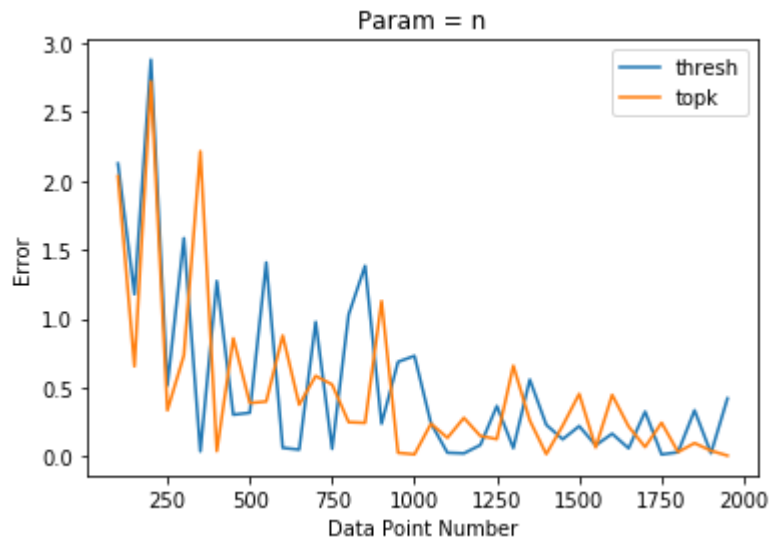
plt.figure()
plt.plot(nrange, wtopk_error, nrange, wthresh_error)
plt.xlabel('Dimensions')
plt.title('Param = dimensions')
plt.ylabel('Error')
plt.legend()
plt.legend({'topk', 'thresh'})

nrange, _, wtopk_error, wthresh_error= error_calc(num_iters=10, param=
's', n=1000, d=100, s=5, s_model=True, true_s=5)

plt.figure()
plt.plot(nrange, wtopk_error, nrange, wthresh_error)
plt.xlabel('Sparsity')
plt.title('Param = sparsity')
plt.ylabel('Error')
plt.legend()
plt.legend({'topk', 'thresh'})
```

No handles with labels found to put in legend.
No handles with labels found to put in legend.

Out[34]: <matplotlib.legend.Legend at 0x7f9d9652d748>



Part (c)

Now, call the helper function with the true sparsity being greater than the sparsity assumed by the top-k estimator. Remember to set `s_model` to `False`! Plot the behavior of all three estimators once again, as a function of n , d , s , where s is the assumed sparsity of the top-k model.

You should return 3 plots, and explain what you see in terms of the bias variance tradeoff.


```
In [38]: nrange, _, wtopk_error, wthresh_error= error_calc(num_iters=10, param=
'n', n=1000, d=100, s=5, s_model=False, true_s=5)

plt.figure()
plt.plot(nrange, wtopk_error, nrange, wthresh_error)
plt.xlabel('Data Point Number')
plt.title('Param = n')
plt.ylabel('Error')
plt.legend({'topk', 'thresh'})

nrange, _, wtopk_error, wthresh_error= error_calc(num_iters=10, param=
'd', n=1000, d=100, s=5, s_model=False, true_s=5)

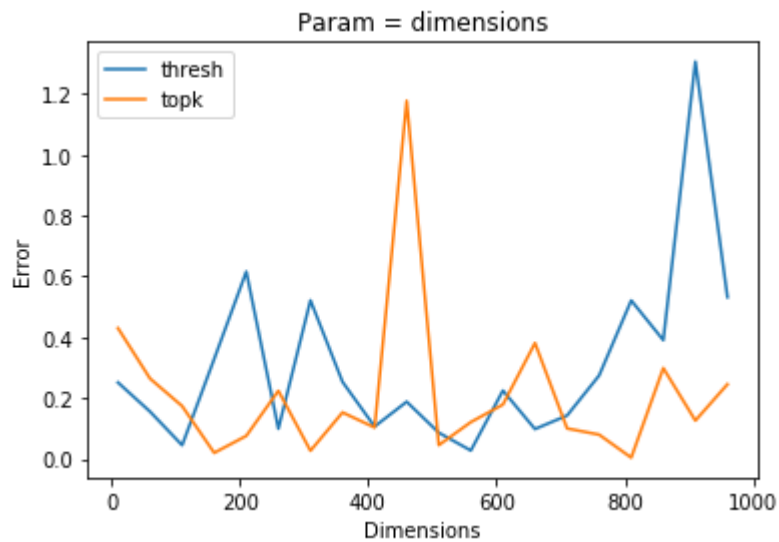
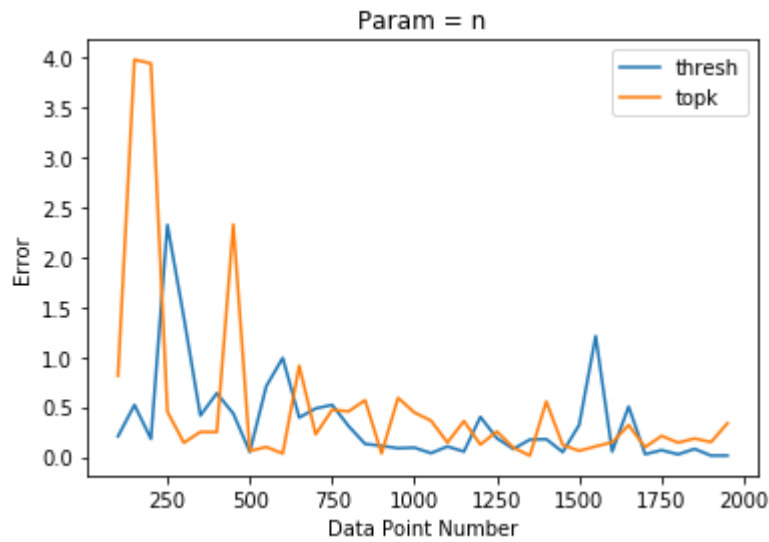
plt.figure()
plt.plot(nrange, wtopk_error, nrange, wthresh_error)
plt.xlabel('Dimensions')
plt.title('Param = dimensions')
plt.ylabel('Error')
plt.legend()
plt.legend({'topk', 'thresh'})

nrange, _, wtopk_error, wthresh_error= error_calc(num_iters=10, param=
's', n=1000, d=100, s=5, s_model=False, true_s=5)

plt.figure()
plt.plot(nrange, wtopk_error, nrange, wthresh_error)
plt.xlabel('Sparsity')
plt.title('Param = sparsity')
plt.ylabel('Error')
plt.legend()
plt.legend({'topk', 'thresh'})
```

No handles with labels found to put in legend.
No handles with labels found to put in legend.

Out[38]: <matplotlib.legend.Legend at 0x7f9d97c40b70>



In the Model we see that with higher dimensional data we don't get more error, that is because we are setting our complexity within our sparsity to be set. Therefore, our variance is set, and our bias is already high. However when looking at the assumed sparsity of our model (which should be 100, since we don't have any sparsity), we see that the error stays the same. The randomness comes from our noise, however there seems to be no visible trend. We expect that as we increase sparsity, the more our error will decrease because we are approaching optimal variance and bias tradeoff as we approach the true model.

```
In [1]: import matplotlib.pyplot as plt
import matplotlib.image as mpimg

%matplotlib notebook
```

PART A

```
In [2]: # @staticmethod
# def information_gain(X, y, thresh):
#     X0,y0,X1,y1 = split(X,y, thresh)
#     def g(y):
#         pyk_2 = 0
#         for i in range(2):
#             try:
#                 pyeqk = len(y[y == i])/float(len(y))
#                 pyk_2 = pyeqk**2
#             except ZeroDivisionError:
#                 pyeqk = 0
#         return 1 - pyk_2
#
#     return len(X0)*g(y0)/float(len(X)) + len(X1)*g(y1)/float(len(X))

# @staticmethod
# def gini_impurity(X, y, thresh):
#     X0,y0,X1,y1 = split(X,y, thresh)
#
#     def h(y) :
#         hy = 0
#         for i in range(2) :
#             try :
#                 pyeqk = len(y[y == i])/float(len(y))
#                 hy = hy - pyeqk*np.log(pyeqk)
#             except ZeroDivisionError :
#                 pyeqk = 0
#
#         return hy
#
#     return len(X0)*h(y0)/float(len(X)) + len(X1)*h(y1)/float(len(X))
```

PART C



PART D (Bagged Trees)

```
In [ ]: # 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, **s
self.params)
#                 for i in range(self.n)
#             ]

#         def fit(self, X, y):
#             self.mask = []
#             for tree in self.decision_trees:
#                 mask = np.random.randint(0, high = len(X), size= len
(X))
#                 Xsampling = X[mask,:]
#                 ysampling = y[mask]
#                 tree.fit(Xsampling,ysampling)
#                 self.mask.append(mask)

#         def predict(self, X):
#             preds = []
#             for tree in self.decision_trees:
#                 preds.append(tree.predict(X))
#             return stats.mode(np.array(preds), axis = 0 )[0].reshape(len
(X))
```

PART F (Random Forest)

```

In [ ]: # class RandomForest(BaggedTrees):
#         def __init__(self, params=None, n=200, m=2):
#             super().__init__(params = params , n = n )
#             self.m = m

#         def fit(self, X,y):
#             self.mask = []
#             self.features = []
#             for tree in self.decision_trees:
#                 mask = np.random.randint(0, high = len(X), size= len
(X))
#                 features = np.random.choice( X.shape[1], size = self.
m )
#                 Xsampling = X[mask,:]
#                 Xsampling = Xsampling[:,features]
#                 ysampling = y[mask]
#                 tree.fit(Xsampling,ysampling)
#                 self.mask.append(mask)
#                 self.features.append(features)
#         def predict(self,X):
#             preds = []
#             k = 0
#             for tree in self.decision_trees:
#                 preds.append(tree.predict(X[:,self.features[k]]))
#             return stats.mode(np.array(preds), axis = 0 )[0].reshape(len
(X))

```

PART H (AdaBoost)

```

In [ ]: # 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
#             k = 0
#             self.features = []
#             for tree in self.decision_trees:
#                 mask = np.random.randint(0, high = len(X), size= len
(X))
#                 features = np.random.choice(X.shape[1], size = self.m
)
#                 Xsampling = X[mask,:]
#                 Xsampling = Xsampling[:,features]
#                 ysampling = y[mask]
#                 tree.fit(Xsampling,ysampling)
#                 self.features.append(features)
#
#                 ej = 0
#                 for j in range(len(Xsampling)):
#                     ej = checkXY(Xsampling[j,:], ysampling[j], tree)
#                 ej = ej/float(sum(self.w))
#
#                 self.a[k] = 0.5*np.log((1-ej)/float(ej))
#
#                 for i in range(len(Xsampling)):
#                     if checkXY(Xsampling[i,:],ysampling[i], tree)
> 0.5 :
#                         self.w[i] = self.w[i]*np.exp(self.a
[k])
#                     else :
#                         self.w[i] = self.w[i]*np.exp(-self.a
[k])
#                 k = k + 1
#
#         def predict(self, X):
#             classes = list(set(y))
#             preds_tot = []
#             for i in range(len(X)):
#                 preds = []
#                 for c in classes :
#                     zj = 0
#                     k = 0
#                     for tree in self.decision_trees:
#                         Xcheck = X[:,self.features[k]]
#                         Xcheck = Xcheck[i,:]
#                         zj = zj + self.a[k]*checkXY(Xcheck,
c, tree)
#                     k = k + 1
#                 preds.append(zj)
#                 preds_tot.append(classes[np.argmax(preds)])
#             return preds_tot

```

Part J (Results)

```
In [ ]: #TITANIC :
# accuracy = [kfold1, kfold2, kfold3, avg_training]

# DecisionTree
# [0.5993975903614458, 0.63855421686746983, 0.60240963855421692, 0.61369315342328834]
# BaggedTrees
# accuracy = [kfold1, kfold2, kfold3, avg_training]
# [0.79518072289156627, 0.76204819277108438, 0.77710843373493976, 0.97801099450274853]
# RandomForest
# accuracy = [kfold1, kfold2, kfold3, avg_training]
# [0.5993975903614458, 0.4006024096385542, 0.60240963855421692, 0.58020989505247378]
# Adaboost
# [0.21686746987951808, 0.63855421686746983, 0.21385542168674698, 0.24837581209395301]

# SPAM

# DecisionTree
# accuracy = [kfold1, kfold2, kfold3, avg_training]
# [0.5993975903614458, 0.63855421686746983, 0.60240963855421692, 0.61369315342328834]
# BaggedTrees
# accuracy = [kfold1, kfold2, kfold3, avg_training]
# [0.78915662650602414, 0.76204819277108438, 0.77409638554216864, 0.97801099450274853]
# RandomForest
# accuracy = [kfold1, kfold2, kfold3, avg_training]
# [0.5993975903614458, 0.6506024096385542, 0.61144578313253017, 0.62868565717141422]
# Adaboost
# accuracy = [kfold1, kfold2, kfold3, avg_training]
# [0.5993975903614458, 0.22289156626506024, 0.21987951807228914, 0.24887556221889059]
```