## 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 that?

OVERFITTING.
that can be seen in the training accuracy results of a single decision tree vs its validation accuracy.

## Q2

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

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

$$= \frac{1}{2}\|y\|^2 + \frac{1}{2}w^T X^T X\underline{w} - y^T Xw + \lambda\|w\|_1$$

$$= \frac{1}{2}\|y\|^2 + \frac{1}{2}w^T \left(\sum_{ij}\underline{x}_i^T \underline{x}_j\right)\underline{w} - y^T Xw + \lambda\|w\|_1$$

$$= \frac{1}{2}\|y\|^2 + \frac{n}{2}w^T w - y^T X\underline{w} + \lambda\|w\|$$

Uncorrelated $x_i$'s

$$\rightarrow \quad E[x_i x_j^T] - E[x_i]E[x_j] = \text{diag for } X$$

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

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

Solve All together

$$\frac{d}{dw}(J_\lambda(w)) = n\underline{w} + \frac{d}{dw}\|w\|_1 - (X^T y) = 0$$

$$= n\underline{w} + \text{sgn}(\underline{w}) - x^T y = 0$$

$$\underline{w} = \frac{1}{n}\left[X^T y - \text{sgn}(\underline{w})(\lambda)\right]$$

Solve for $w_i$ Seprated

$$\frac{d}{dw_r} J_{i\lambda}(w_i) = \frac{d}{dw_i}\left[\frac{n}{2}w_i^2 - (y^T \underline{x}_i)w_i + \lambda|w_i|\right]$$

$$= nw_i - y^T \underline{x}_i + \lambda \, \text{sgn}(w_i) = 0$$

$$w_i = \frac{1}{n}\left[+y^T \underline{x}_i - \lambda \, \text{sgn}(w_i)\right]$$

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

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

— negative

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

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

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

f) $1e-5 = \lambda$

## Q3

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

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

$$= + 1/d^2$$

$$d^* \left( Pr(|z_i| \geq 2\sigma\sqrt{\log d}) \right) \leq \left( \frac{1}{d^2} \right) \cdot d^1 \leq 1/d \Rightarrow P(max|z_i| \geq 2\sigma\sqrt{\log d})$$

all of them $< 2\sigma\sqrt{\log d}$   $\leq 1/d$

e) $w = (X^TX)^{-1}X^T(y) = (X^TX)^{-1}X^T(y^* + z) = w^* + (X^TX)^{-1}X^Tz$

$$= w^* + z'$$

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

because

$\sigma^2 (X^TX)^{-1} X^TX (X^TX)^{-1} I$

$I$ because has orthonormal columns

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

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

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

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

$max(\hat{w}_{top}(s) - w^*)$  ➔ 1- will give s-wrongly predicted non-zero entries

2- negative s entries of original weights

$max(Sparsity(e)) = 2s$

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

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

Consider max entry

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

$e_i = w_i^* + z_i' - w_i^*$

$\quad = z_i$   or   $-w_i^*$

$T_s(w_i^* + z_i)$

$\begin{cases} |w_i^* + z_i| & \text{if } i \text{ is top s} \\ \text{c} & \text{otherwise} \end{cases}$

$$e_i = \begin{cases} -z_i & \max|e_i| \le 2 \delta \sqrt{\log d} \\ - W_i^* & \leftarrow \text{Assume} \to 0 \\ & \quad \text{because} \quad \text{we pr} \end{cases}$$

(h) $\Pr\left( \|\hat{w}_{top(s)} - w^*\|^2 \le 32 \delta^2 s \log d \right)$

$\underset{\uparrow \text{here}}{\sum_{i}^{n} |e_i|^2} < \left(4\delta\sqrt{\log d}\right)^2 \times 2s \underset{\text{because } |e_i| \ne 0}{\uparrow} \quad \text{only } 2s \text{ times} \atop \max$

$\qquad\qquad < (16 \, \delta^2 \log d) \, 2s$

$\qquad\qquad \le 32 \delta^2 \log d$

G. Given that $\varepsilon$

from (d) $\Rightarrow \Pr(\text{not } \varepsilon) < \frac{1}{d}$

$\qquad\qquad \to \Pr(\varepsilon) \ge 1 - \frac{1}{d}$

(i) $\|Xv\|^2 = \|v\|^2$ since $X$ has orthonormal columns

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

$\frac{1}{n}\|X(\hat{w}_{top(s)} - w^*)\| \le \frac{1}{n} 32 \, \delta^2 s \log d \qquad \text{wm Prob } (1 - 1/d)$

(j) When Sparsity comes into play

$\left(\frac{\|x_{top}\|^2}{n}\right) \le 32 \delta^2 s \log d \qquad$ compared $\frac{\|x_{lns}\|}{n} \le \delta^2 \frac{d}{n}$

wm $\Pr(1 - 1/d)$

$\Rightarrow$ high probability

$\qquad\qquad \Rightarrow \frac{1}{n}\|x_{top}\|^2 \le \frac{1}{n}\|x_{lns}\|^2$

$\qquad\qquad \frac{32 \delta^2 s \log d}{\cancel{n}} \le \frac{\cancel{\delta^2} d}{\cancel{n}}$

$\qquad\qquad \qquad \underset{\text{when}}{\cancel{\delta}}$

$$\boxed{s \le \frac{d}{32 \log d}}$$

(4i) The variance would only be $\sigma^2$ because that's the variance of the noise.

there is a factor of $32\sigma^2\log d$ 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.

e) ~~exclamation~~ $< 0.5$ for all
    ~~< 0.5~~

g) Spam:
    money $< 0.5$
same common    private $< 0.5$
    bank $< 0.5$

    Most Common
    = 'money'

Titanic:
    sex $< 0.5$
    class
    Sibsp $< 0.5$

    Most Common
    = Sex

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

When label is not spam it is easy to classify because
they mostly only have characters and no keywords of
the one above. Then
    $w < 1e-35$