

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

$$2.(a) -\nabla_{\vec{v}} L(\vec{v}) = \nabla_{\vec{v}} \frac{1}{2} (\vec{y}^T \vec{y} - 2\vec{y}^T \vec{v} + \vec{v}^T \vec{v}) \\ = \vec{v} - \vec{y}$$

$$-\nabla^{t+1} = \vec{v}^t - \text{dt} \nabla_{\vec{v}} L(\vec{v}) |_{\vec{v}=\vec{v}^t}$$

$$= \vec{v}^t - \text{dt} (\vec{v}^t - \vec{y})$$

$$= (1 - \text{dt}) \vec{v}^t + \text{dt} \vec{y}$$

$$-\nabla_{\vec{v}} L(\vec{v}) |_{\vec{v}=x\vec{w}} = x\vec{w} - \vec{y}$$

$$-\vec{v}^{t+1} = (1 - \text{dt}) x\vec{w}^t + \text{dt} \vec{y}$$

$$-x\vec{w}^{t+1} = x\vec{w}^t - \text{dt} X X^T (x\vec{w}^t - \vec{y})$$

The difference between \vec{v}^{t+1} and $x\vec{w}^{t+1}$ is because \vec{v} and \vec{w} lie in different space, so $x\vec{w}^{t+1}$ has an extra XX^T term.

$$(b) \tilde{v}^t = \underset{\tilde{v} \in \tilde{V}}{\operatorname{argmax}} \langle -\nabla L(\vec{v}^t), \tilde{v} \rangle$$

$$= \underset{\tilde{v} = x\vec{w}, \| \vec{w} \|_2^2 \leq 1}{\operatorname{argmax}} \langle -x\vec{w}^t + \vec{y}, \tilde{v} \rangle$$

$$= \underset{\| \vec{w} \|_2^2 \leq 1}{X \operatorname{argmax}} \langle -x\vec{w}^t + \vec{y}, x\vec{w} \rangle$$

$$= \underset{\| \vec{w} \|_2^2 \leq 1}{X \operatorname{argmax}} (-\vec{w}^{t+1} X^T x\vec{w} + \vec{y}^T x\vec{w})$$

$$\text{Let } \mathcal{L} = -\vec{w}^{t+1} X^T x\vec{w} + \vec{y}^T x\vec{w} + \lambda (1 - \vec{w}^T \vec{w})$$

$$\frac{\partial \mathcal{L}}{\partial \vec{w}} = 0 \Rightarrow -\vec{w}^{t+1} X^T X + \vec{y}^T X - 2\lambda \vec{w}^T = 0$$

$$\Rightarrow \vec{w} \propto -X^T X \vec{w}^t + X^T \vec{y}$$

$$\Rightarrow \tilde{v}^t = \underset{\| \vec{w} \|_2^2 \leq 1}{X \operatorname{argmax}} (-\vec{w}^{t+1} X^T x\vec{w} + \vec{y}^T x\vec{w})$$

$$\propto X X^T (x\vec{w}^t - \vec{y})$$

$$\begin{aligned}\vec{v}^{t+1} &= \vec{v}^t + \alpha_t \vec{v}^t \\ &= X\vec{w}^t - \alpha_t X X^\top (X\vec{w}^t - \vec{y}) \\ &= X\vec{w}^{t+1}\end{aligned}$$

(c) $\vec{r}^t = \underset{i=1, \dots, d}{\operatorname{argmax}} | \langle \vec{r}^t, \phi_i(x) \rangle |$

$$\vec{v}^t = \underset{\vec{v} \in \tilde{\mathcal{V}}}{\operatorname{argmax}} \langle -\nabla L(\vec{v}^t), \vec{v} \rangle$$

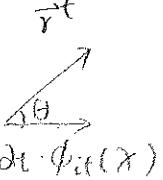
$$= \underset{\vec{v} \in \tilde{\mathcal{V}}}{\operatorname{argmax}} \langle -X\vec{w}^t + \vec{y}, \vec{v} \rangle$$

$$= \underset{\vec{v} \in \tilde{\mathcal{V}}}{\operatorname{argmax}} \langle \vec{r}^t, \vec{v} \rangle$$

$$\Rightarrow \tilde{\mathcal{V}} = \{ \vec{v} = \pm \phi_i(x), i=1, \dots, d \}$$

In this case, $\vec{v}^t = \underset{i=1, \dots, d}{\operatorname{argmax}} | \langle \vec{r}^t, \phi_i(x) \rangle |$

(d)



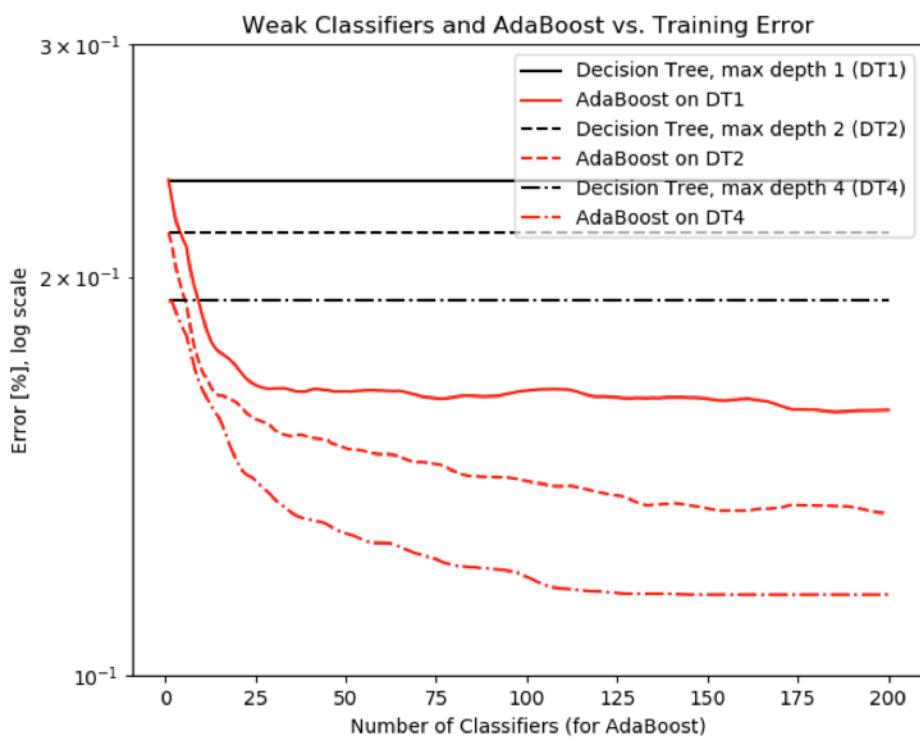
To pick a step-size α_t that reduces the loss most, we can project \vec{r}^t into the direction of $\phi_{it}(x)$.

Then $\vec{r}^t \cdot \cos \theta = \alpha_t \phi_{it}(x)$

$$\alpha_t = \frac{\langle \vec{r}^t, \phi_{it}(x) \rangle}{\|\phi_{it}(x)\|^2}$$

so we can see that obtaining α_t via linesearch is the same as doing MP.

2(h)

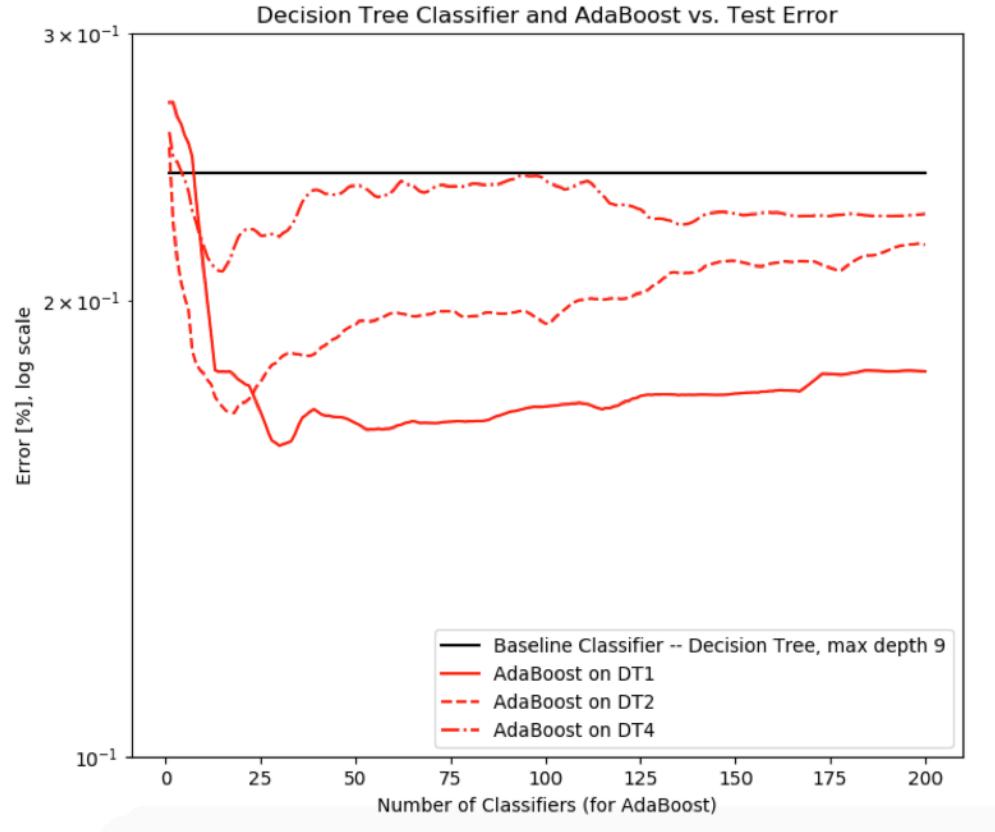


Performance is better with deeper trees.

The training error decreases with deeper trees and more classifiers for AdaBoost.

This is because more classifiers give more accurate models, as the trees are improved along more classifiers.

2(i)



There is a difference between training error and test error.

Training error keeps decreasing when we increase DT, but test error increases from DT1 to DT4.

DT1 works best for AdaBoost.

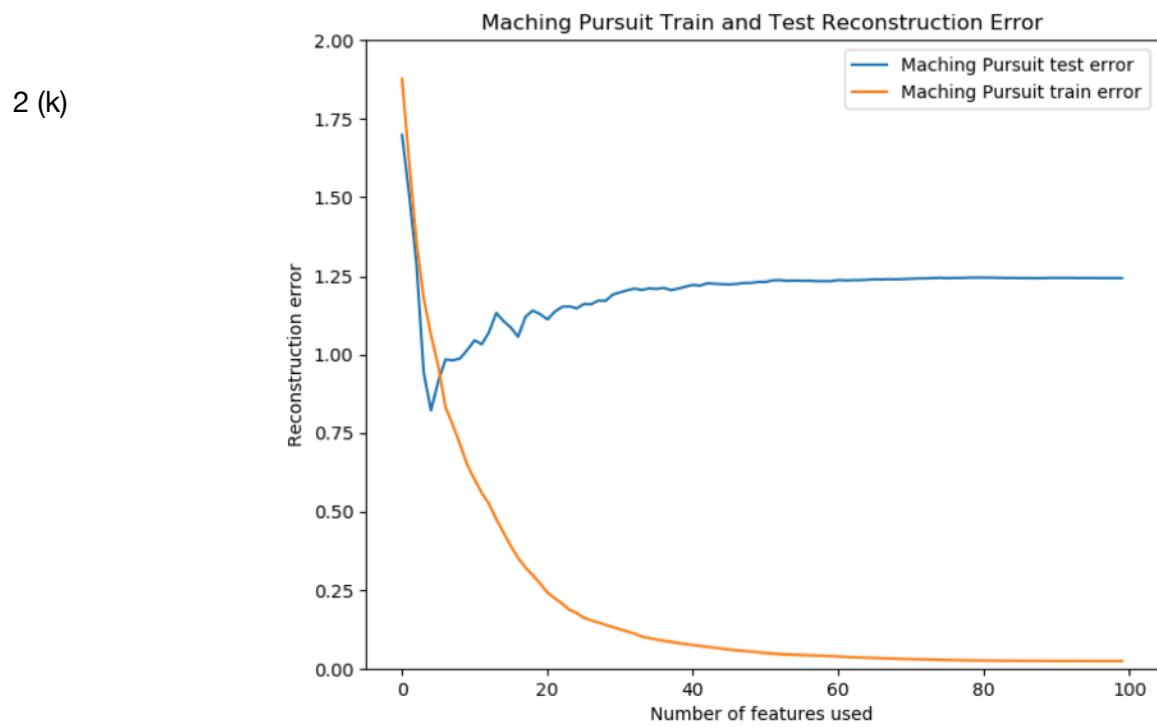
This is because DT2 and DT4 overfit data.

2 (j)

Limiting the number of base classifiers will help for Adaboost.

From plot in (i), we can run more boosting iterations on base classifier of 20.

Because we have less classifiers, when we use deeper trees, there will be less free parameters, so less likely to have overfitting problem.



Training error decreases with number of features, because we are using more features, so we are fitting the observed y better and better.

Test error in this plot is similar to plot in part(i), also because of overfitting.

4(a) The time complexity for computing distance is $\sim O(dn)$
 Finding the k shortest distance using quickselect is $\sim O(n)$.
 The time complexity for a single query is $\sim O(dn)$

(b) For 1-d, you need to search 3 cells.

For 2-d, you need to search $3^2 = 9$ cells

For d -d, you need to search 3^d cells

If each cell contains 1 point, the time complexity $O \sim (d 3^d)$

$$(c) \Pr(|\vec{v}^\ell, \vec{u}| \leq |\langle \vec{v}^\ell, \vec{u} \rangle|)$$

$$= \Pr(|\langle \vec{v}^\ell, \vec{u}'' \rangle| \leq |\langle \vec{v}^\ell, \vec{u}'' \rangle|)$$

$$= \Pr(|\|\vec{v}^\ell\| |\vec{u}''| | \cos \theta | \leq |\|\vec{v}^\ell\| |\vec{u}''| | \cos \theta |)$$

$$\leq \Pr(|\|\vec{v}^\ell\| | \cos \theta | \leq |\|\vec{v}^\ell\||)$$

$$= \Pr(| \cos \theta | \leq |\|\vec{v}^\ell\|| / |\|\vec{v}^\ell\||)$$

(d) Let $|\|\vec{v}^\ell\|| / |\|\vec{v}^\ell\|| = Vr$.

$$|\cos \theta| \leq Vr \Rightarrow -Vr \leq \cos \theta \leq Vr$$



$$\Rightarrow \theta \in [\cos^{-1} Vr, \pi - \cos^{-1} Vr] \cup [\pi + \cos^{-1} Vr, 2\pi - \cos^{-1} Vr]$$

$$\Pr(| \cos \theta | \leq Vr) = \frac{\pi - 2\cos^{-1} Vr + \pi - 2\cos^{-1} Vr}{2\pi}$$

$$= 1 - \frac{2}{\pi} \cos^{-1} Vr$$

$$\approx 1 - \frac{2}{\pi} \cos^{-1} \frac{|\|\vec{v}^\ell\||}{|\|\vec{v}^\ell\||}$$

(e) $\Pr(\text{failure})$

= $\Pr(\text{at least } k \text{ points are closer to } \vec{z} \text{ than } \vec{x}^{(i)} \text{ under projection } \vec{u})$

$\leq \frac{1}{k!} \sum_{i=2}^N \Pr(\vec{x}^{(i)} \text{ is closer to } \vec{z} \text{ than } \vec{x}^{(1)})$

$= \frac{1}{k!} \sum_{i=2}^k \Pr(|\langle \vec{x}^{(1)} - \vec{z}, \vec{u} \rangle| \leq |\langle \vec{x}^{(1)} - \vec{z}, \vec{u} \rangle|)$

$\leq \frac{1}{k!} \sum_{i=2}^k \left(1 - \frac{2}{\pi} \cos^{-1} \frac{\|\vec{x}^{(1)} - \vec{z}\|}{\|\vec{x}^{(1)} - \vec{z}\|}\right)$

$\leq \frac{1}{k!} \sum_{i=2}^N \frac{\|\vec{x}^{(1)} - \vec{z}\|}{\|\vec{x}^{(1)} - \vec{z}\|}$

(f) $\|\vec{x}^{(r)} - \vec{z}\| = r$

Let $r^{d'} = i \Rightarrow r = (\frac{i}{c})^{1/d'}$

$$\sum_{i=2}^N \|\vec{x}^{(1)} - \vec{z}\| / \|\vec{x}^{(1)} - \vec{z}\| = \sum_{i=2}^N \frac{(\frac{i}{c})^{1/d'}}{(\frac{i}{c})^{1/d'}} = \sum_{i=2}^N (1/c)^{1/d'}$$

(i) Blue: naive exhaustive search.

green: space partitioning

red: DCI

5. Q: prove $P(\bigvee A_i) \leq \sum P(A_i)$

A: When $i=1 \quad P(A_1) = . \quad P(A_1)$

Assume $P(\bigcup_{i=1-t}^t A_i) \leq \sum_{i=1}^t P(A_i)$

$$P(\bigcup_{i=t+1}^N A_i) = P(A_{t+1} \cup \bigcup_{i=t+1}^N A_i)$$

$$\leq P(A_{t+1}) + P(\bigcup_{i=t+1}^N A_i)$$

$$\leq P(A_{t+1}) + \sum_{i=1}^t P(A_i)$$

$$= \sum_{i=1}^{t+1} P(A_i)$$

So $P(\bigvee A_i) \leq \sum P(A_i)$