# Statistical Methods in AI (CSE 471)
# Ensemble Methods

Vineet Gandhi



INTERNATIONAL INSTITUTE OF
INFORMATION TECHNOLOGY
H Y D E R A B A D

# Bias Variance trade-off

Courtesy — Chris Albon

- Fitting a polynomial to 1d input (linear regression) for various regularization values.

- We can sample different datasets and see the variance in predictions and bias (loss of averaged prediction)



- Left: Predictions trained on various sampled datasets. Low variance
- Right: The mean prediction. High bias

$\ln \lambda = -0.31$

- Left: Predictions trained on various sampled datasets. high variance
- Right: The mean prediction. Low bias

The plot on the left shows the label $\ln \lambda = -2.4$ with $t$ on the vertical axis and $x$ on the horizontal axis. The plot on the right shows the mean prediction with $t$ on the vertical axis and $x$ on the horizontal axis.

- Left: Predictions trained on various sampled datasets. Higher variance

- Right: The mean prediction. Lower bias

# Boosting

# The idea of probabilistic sampling

| Chest Pain | Blocked Arteries | Patient Weight | Heart Disease | Sample Weight |
|------------|------------------|----------------|---------------|---------------|
| Yes | Yes | 205 | Yes | 1/8 |
| No | Yes | 180 | Yes | 1/8 |
| Yes | No | 210 | Yes | 1/8 |
| Yes | Yes | 167 | Yes | 1/8 |
| No | Yes | 156 | No | 1/8 |
| No | Yes | 125 | No | 1/8 |
| Yes | No | 168 | No | 1/8 |
| Yes | Yes | 172 | No | 1/8 |

In a **Random Forest**, each tree has an equal vote on the final classification.

In contrast, in a **Forest of Stumps** made with **AdaBoost**, some stumps get more say in the final classification than others.

Lastly, in a **Random Forest**, each decision tree is made independently of the others.



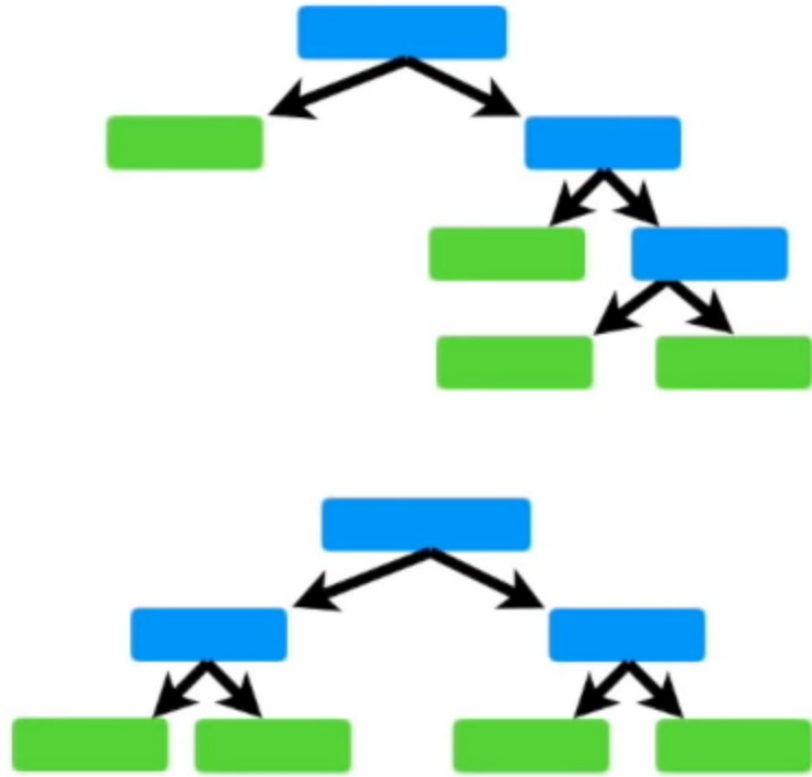In contrast, in a **Forest of Stumps** made with **AdaBoost**, order is important.

In a **Random Forest**, each time you make a tree, you make a full sized tree.

In contrast, in a **Forest of Trees** made with **AdaBoost**, the trees are usually just a **node** and two **leaves**.

# Adaboost

- In Adaboost we assign (non-negative) weights to points in the data set, which are then normalised so that they sum to one
- Iteratively learn new classifier
- In each iteration, we generate a training set by sampling from the data using the weights
- After learning the current classifier, we increase the (relative) weights of the data points which are misclassified by the current classifier
- The final classifier is the weighted majority voting by all classifiers

# Adaboost

- Let $\{(X_1, y_1),\ldots,(X_n, y_n)\}$ be the data. We take $y_i$ in $\{-1, +1\}$
- Let $w_i(k)$ denote the weight for the ith data point at kth iteration
- Let $h_k$ be the classifier learnt at kth iteration, we take $h_k(X)$ in $\{-1, +1\}$
- We assume error rate of each classifier on its training data is less than 0.5

Given: $(x_1, y_1), \ldots, (x_m, y_m)$ where $x_i \in \mathcal{X}$, $y_i \in \{-1, +1\}$.

Initialize: $D_1(i) = 1/m$ for $i = 1, \ldots, m$.

For $t = 1, \ldots, T$:

- Train weak learner using distribution $D_t$.
- Get weak hypothesis $h_t : \mathcal{X} \to \{-1, +1\}$.
- Aim: select $h_t$ with low weighted error:

$$\varepsilon_t = \Pr_{i \sim D_t}\left[h_t(x_i) \neq y_i\right].$$

- Choose $\alpha_t = \frac{1}{2}\ln\left(\frac{1 - \varepsilon_t}{\varepsilon_t}\right)$.
- Update, for $i = 1, \ldots, m$:

$$D_{t+1}(i) = \frac{D_t(i)\exp(-\alpha_t y_i h_t(x_i))}{Z_t}$$

where $Z_t$ is a normalization factor (chosen so that $D_{t+1}$ will be a distribution).
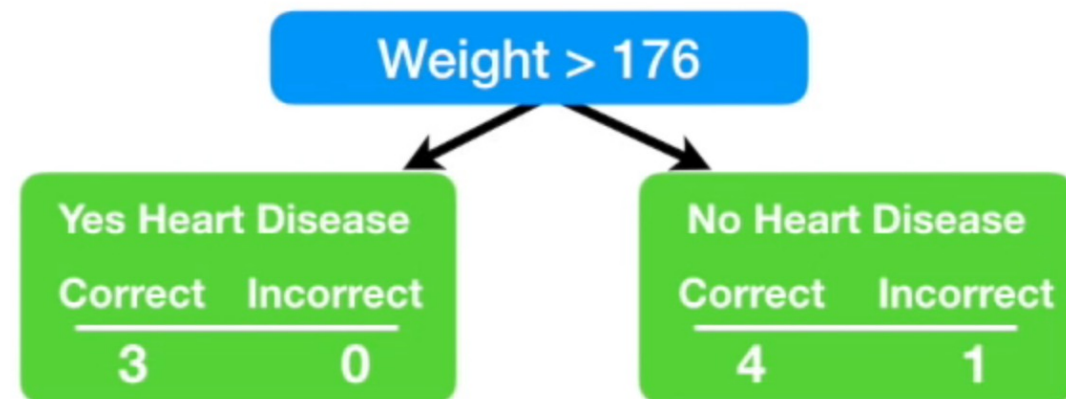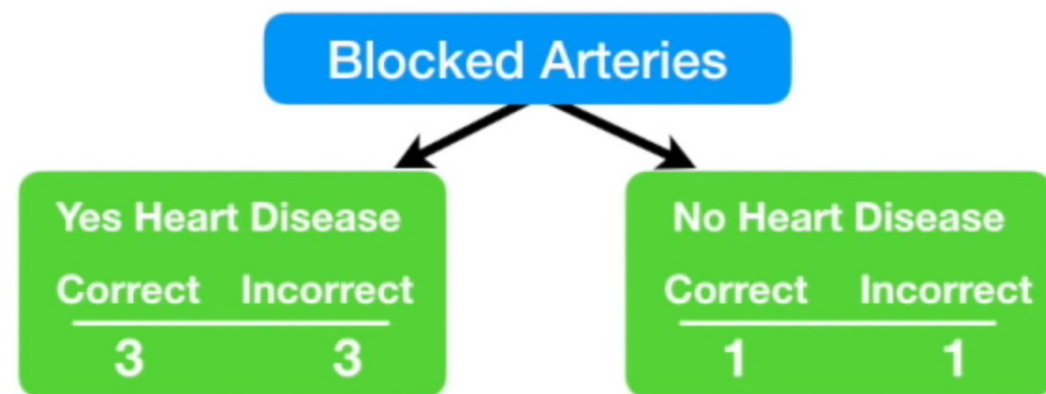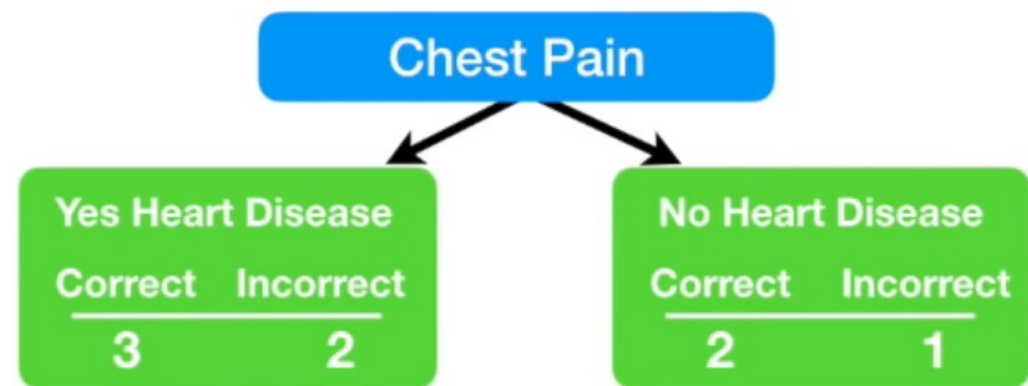
Output the final hypothesis:

$$H(x) = \operatorname{sign}\left(\sum_{t=1}^{T} \alpha_t h_t(x)\right).$$

**Fig. 1** The boosting algorithm AdaBoost.

# Lets take an example

| Chest Pain | Blocked Arteries | Patient Weight | Heart Disease | Sample Weight |
|---|---|---|---|---|
| Yes | Yes | 205 | Yes | 1/8 |
| No | Yes | 180 | Yes | 1/8 |
| Yes | No | 210 | Yes | 1/8 |
| Yes | Yes | 167 | Yes | 1/8 |
| No | Yes | 156 | No | 1/8 |
| No | Yes | 125 | No | 1/8 |
| Yes | No | 168 | No | 1/8 |
| Yes | Yes | 172 | No | 1/8 |

**Chest Pain**

Yes Heart Disease
Correct 3 | Incorrect 2

No Heart Disease
Correct 2 | Incorrect 1

**Blocked Arteries**

Yes Heart Disease
Correct 3 | Incorrect 3

No Heart Disease
Correct 1 | Incorrect 1

**Weight > 176**

Yes Heart Disease
Correct 3 | Incorrect 0

No Heart Disease
Correct 4 | Incorrect 1

| Chest Pain | Blocked Arteries | Patient Weight | Heart Disease | Sample Weight |
|---|---|---|---|---|
| Yes | Yes | 205 | Yes | 1/8 |
| No | Yes | 180 | Yes | 1/8 |
| Yes | No | 210 | Yes | 1/8 |
| Yes | Yes | 167 | Yes | 1/8 |
| No | Yes | 156 | No | 1/8 |
| No | Yes | 125 | No | 1/8 |
| Yes | No | 168 | No | 1/8 |
| Yes | Yes | 172 | No | 1/8 |

Weight > 176

Yes Heart Disease

Correct    Incorrect
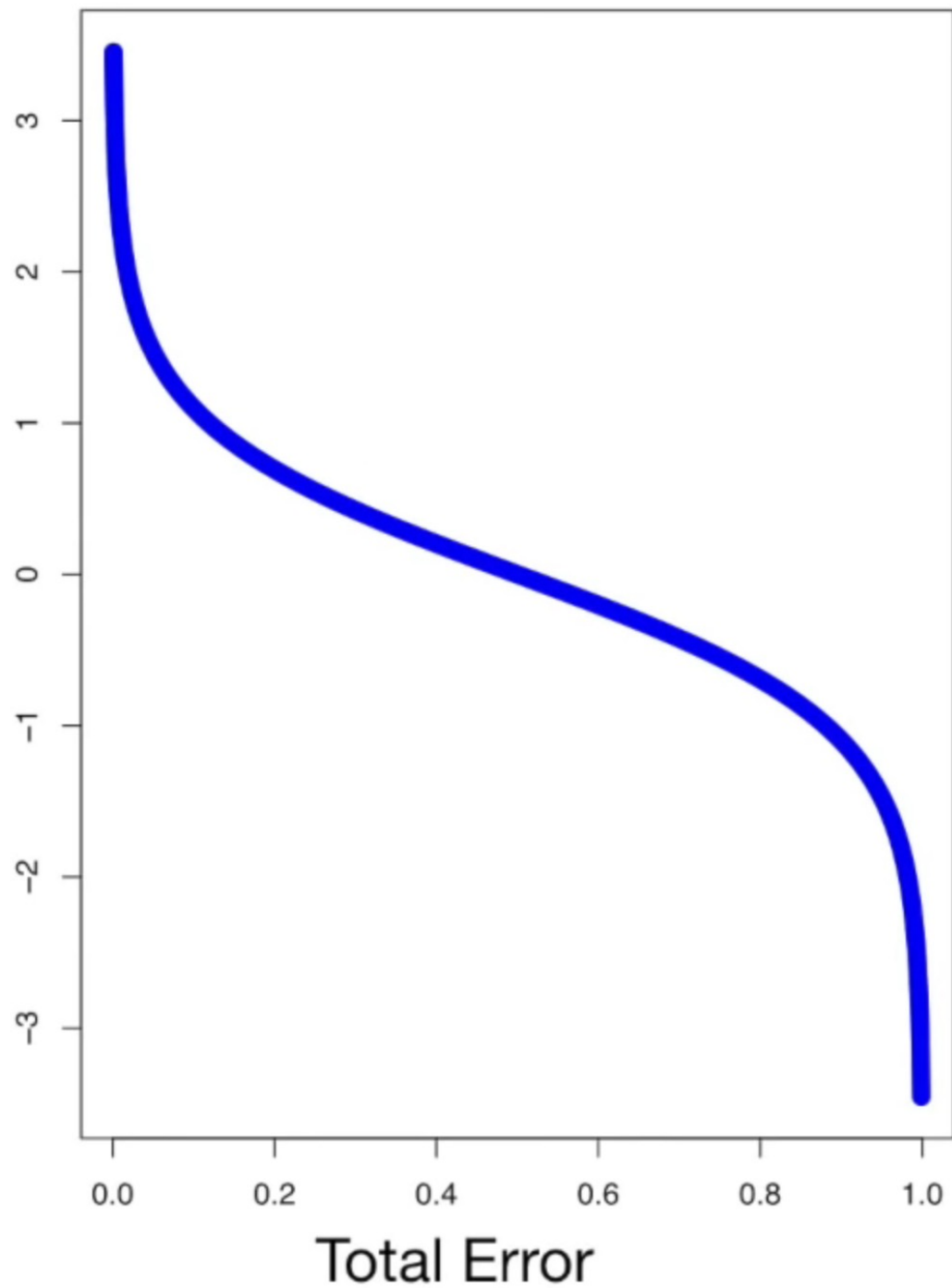3          0

No Heart Disease

Correct    Incorrect
4          1

Thus, in this case, the **Total Error** is **1/8**.

| Chest Pain | Blocked Arteries | Patient Weight | Heart Disease | Sample Weight |
|---|---|---|---|---|
| Yes | Yes | 205 | Yes | 1/8 |
| No | Yes | 180 | Yes | 1/8 |
| Yes | No | 210 | Yes | 1/8 |
| Yes | Yes | 167 | Yes | 1/8 |
| No | Yes | 156 | No | 1/8 |
| No | Yes | 125 | No | 1/8 |
| Yes | No | 168 | No | 1/8 |
| Yes | Yes | 172 | No | 1/8 |

**Weight > 176**

**Yes Heart Disease**

| Correct | Incorrect |
|---|---|
| 3 | 0 |

**No Heart Disease**

| Correct | Incorrect |
|---|---|
| 4 | 1 |

Thus, in this case, the **Total Error** is **1/8**.

$$\text{Amount of Say} = \frac{1}{2} \log\left(\frac{1 - \text{Total Error}}{\text{Total Error}}\right)$$

Amount of Say = $\frac{1}{2}$ log($\frac{1 - \text{Total Error}}{\text{Total Error}}$)

Amount of Say = $\frac{1}{2}$ log( 7 ) = 0.97

| Chest Pain | Blocked Arteries | Patient Weight | Heart Disease | Sample Weight |
|---|---|---|---|---|
| Yes | Yes | 205 | Yes | 1/8 |
| No | Yes | 180 | Yes | 1/8 |
| Yes | No | 210 | Yes | 1/8 |
| Yes | Yes | 167 | Yes | 1/8 |
| No | Yes | 156 | No | 1/8 |
| No | Yes | 125 | No | 1/8 |
| Yes | No | 168 | No | 1/8 |
| Yes | Yes | 172 | No | 1/8 |

New Sample Weight $=$ sample weight $\times e^{\text{amount of say}}$

$$= \frac{1}{8} e^{0.97} = \frac{1}{8} \times 2.64 = 0.33$$

New Sample Weight $=$ sample weight $\times e^{-\text{amount of say}}$

$$= \frac{1}{8} e^{-0.97} = \frac{1}{8} \times 0.38 = 0.05$$

| Chest Pain | Blocked Arteries | Patient Weight | Heart Disease | Sample Weight | New Weight | Norm. Weight |
|---|---|---|---|---|---|---|
| Yes | Yes | 205 | Yes | 1/8 | 0.05 | 0.07 |
| No | Yes | 180 | Yes | 1/8 | 0.05 | 0.07 |
| Yes | No | 210 | Yes | 1/8 | 0.05 | 0.07 |
| Yes | Yes | 167 | Yes | 1/8 | 0.33 | 0.49 |
| No | Yes | 156 | No | 1/8 | 0.05 | 0.07 |
| No | Yes | 125 | No | 1/8 | 0.05 | 0.07 |
| Yes | No | 168 | No | 1/8 | 0.05 | 0.07 |
| Yes | Yes | 172 | No | 1/8 | 0.05 | 0.07 |

| Chest Pain | Blocked Arteries | Patient Weight | Heart Disease | Sample Weight |
|---|---|---|---|---|
| Yes | Yes | 205 | Yes | 0.07 |
| No | Yes | 180 | Yes | 0.07 |
| Yes | No | 210 | Yes | 0.07 |
| Yes | Yes | 167 | Yes | 0.49 |
| No | Yes | 156 | No | 0.07 |
| No | Yes | 125 | No | 0.07 |
| Yes | No | 168 | No | 0.07 |
| Yes | Yes | 172 | No | 0.07 |

| Chest Pain | Blocked Arteries | Patient Weight | Heart Disease |
|---|---|---|---|
| No | Yes | 156 | No |
| Yes | Yes | 167 | Yes |
| No | Yes | 125 | No |
| Yes | Yes | 167 | Yes |
| Yes | Yes | 167 | Yes |
| Yes | Yes | 172 | No |
| Yes | Yes | 205 | Yes |
| Yes | Yes | 167 | Yes |

Ultimately, the patient is classified
as **Has Heart Disease** because
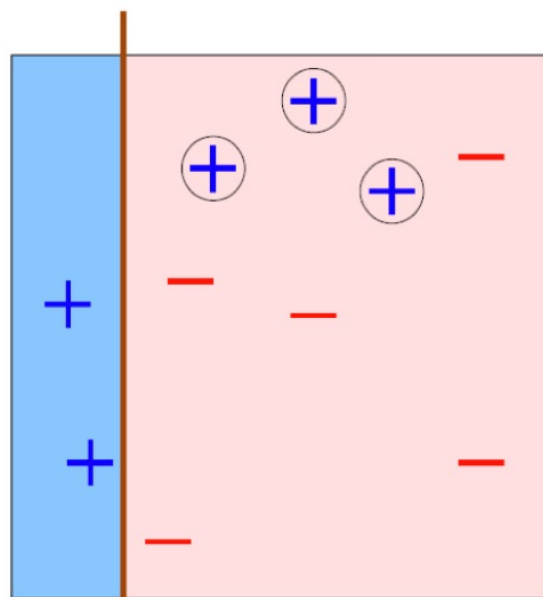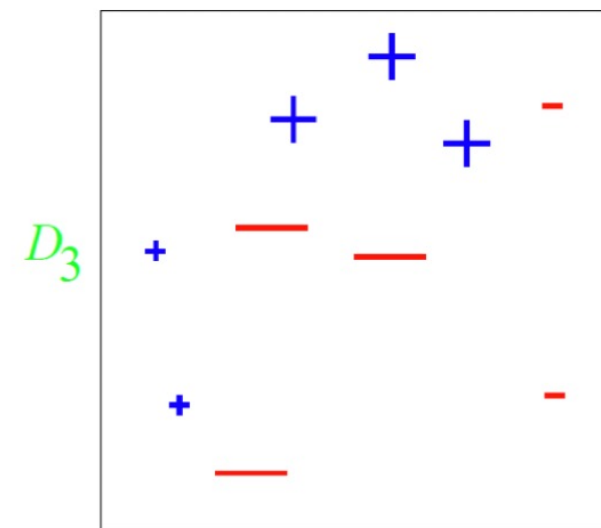this is the larger sum.

Has Heart Disease

Total = 2.7

Total = 1.23
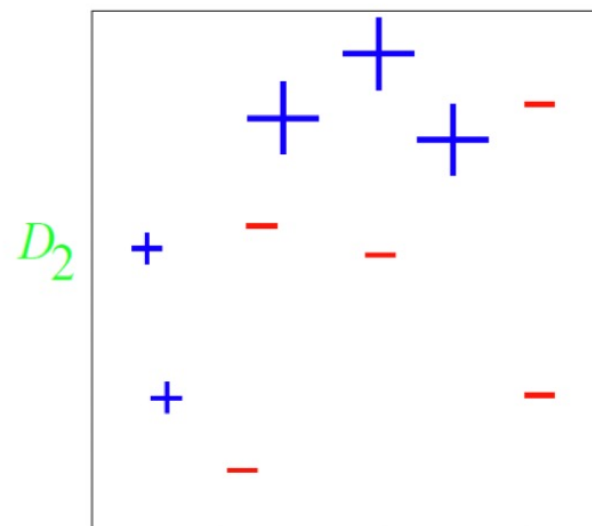
Does Not Have
Heart Disease

**Amount of Say**

→ 0.97

→ 0.32

→ 0.78

→ 0.63

**Amount of Say**

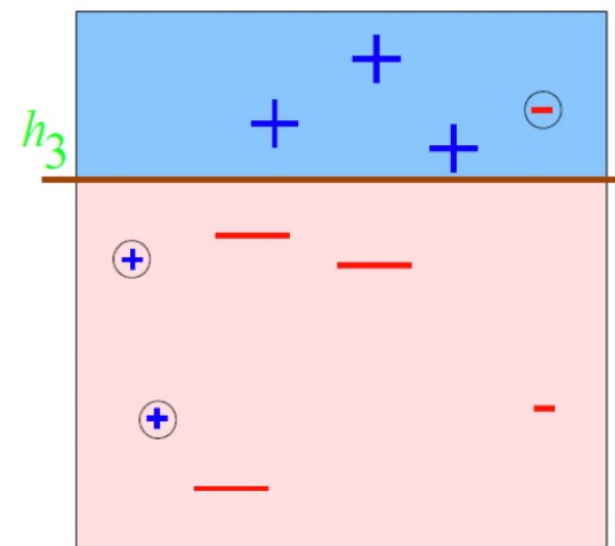0.41 ←

0.82 ←

$D_1$

$h_1$

$\varepsilon_1 = 0.30$

$\alpha_1 = 0.42$

$D_2$

$h_2$

$\varepsilon_2 = 0.21$

$\alpha_2 = 0.65$

$D_3$

$h_3$

$\varepsilon_3 = 0.14$

$\alpha_3 = 0.92$

13

$$H_{\text{final}} = \text{sign} \left( 0.42 \quad \blacksquare \quad + 0.65 \quad \blacksquare \quad + 0.92 \quad \blacksquare \right)$$