

Ensemble methods

Bagging

Herman Kamper

<http://www.kamperh.com/>

Bagging:

$$\underline{x} = \begin{bmatrix} - (x^{(1)})^T - \\ - (x^{(2)})^T - \\ \vdots \\ - (x^{(n)})^T - \end{bmatrix} \quad \underline{y} = \begin{bmatrix} y^{(1)} \\ y^{(2)} \\ \vdots \\ y^{(n)} \end{bmatrix}$$

$$\begin{matrix} \uparrow \\ \downarrow \end{matrix} \begin{bmatrix} \underline{x}_1 \\ \vdots \\ \underline{x}_B \end{bmatrix} ; \begin{bmatrix} \underline{y}_1 \\ \vdots \\ \underline{y}_B \end{bmatrix}$$

$$\Rightarrow f_1(x; \Theta_1)$$

$$\Rightarrow f_2(x; \Theta_2)$$

$$\Rightarrow f_B(x; \Theta_B)$$

$$f(x; \Theta) = \frac{1}{B} \sum_{b=1}^B f_b(x; \Theta_b)$$

$$\underline{x} = \begin{bmatrix} \text{orange} \\ \text{blue} \\ \text{green} \end{bmatrix} ; \underline{y} = \begin{bmatrix} \text{orange} \\ \text{blue} \\ \text{green} \end{bmatrix}$$

$$\begin{matrix} \uparrow \\ \downarrow \end{matrix} \begin{bmatrix} \underline{x}_1 \\ \vdots \\ \underline{x}_2 \end{bmatrix} ; \begin{bmatrix} \underline{y}_1 \\ \vdots \\ \underline{y}_2 \end{bmatrix}$$

$$\Rightarrow f_1(x; \Theta_1)$$

$$\Rightarrow f_2(x; \Theta_2)$$

$$\vdots \\ f_B(x; \Theta_B)$$

$$f(x; \Theta) = \frac{1}{B} \sum_{b=1}^B f_b(x; \Theta_b)$$

↑ Regression

Classification: Majority or weighted voting

Ensemble methods

Random forests

Herman Kamper

<http://www.kamperh.com/>

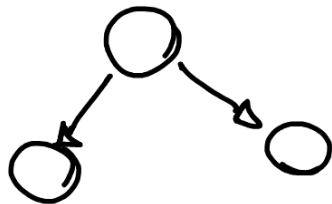
Random forests:

- Specifically used with decision and regression trees.
- Use bagging: Train each tree on different bootstrap sample. But then also...
- Every time we split, only consider $M < D$ random features
- $M = \sqrt{D}$ is often used

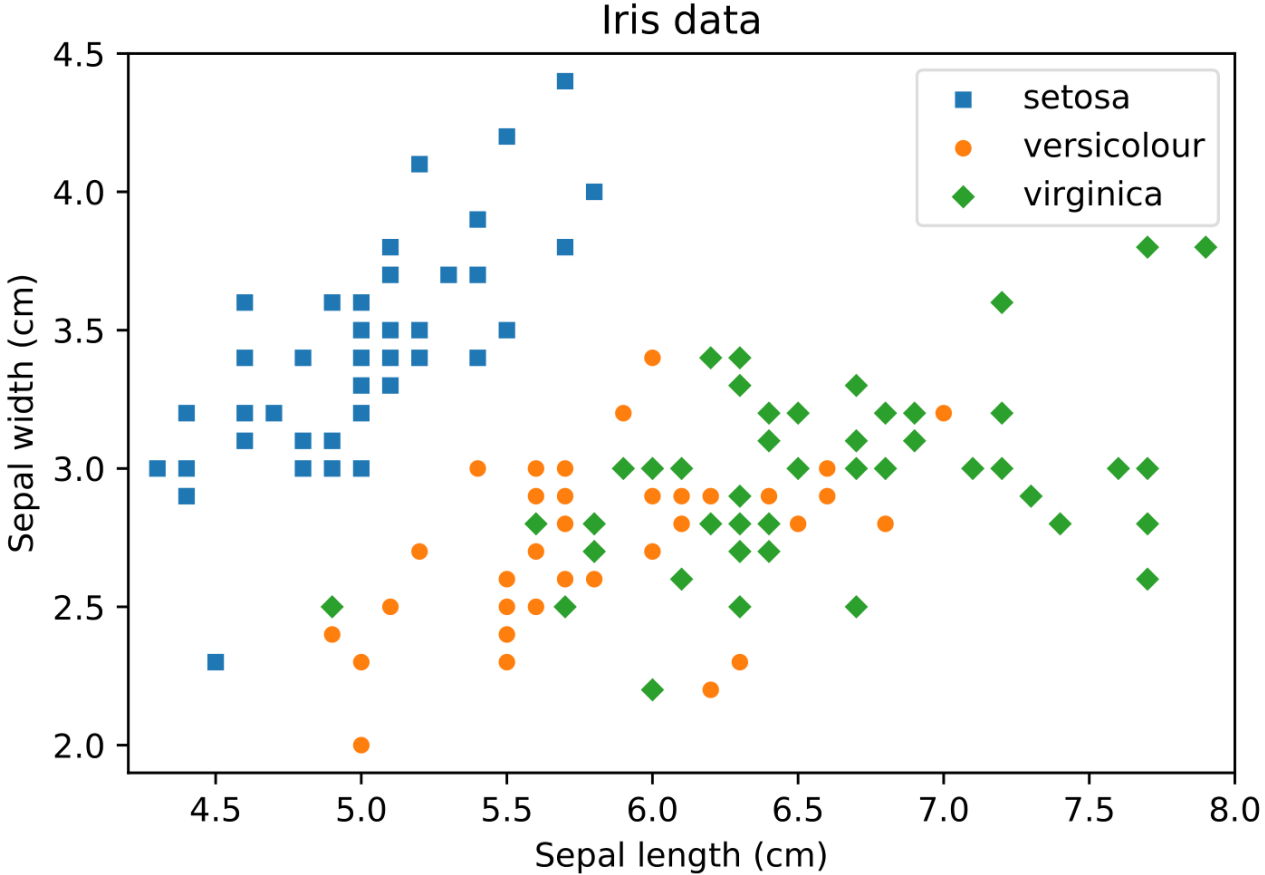
$$f(\underline{x}; \underline{\Theta}) = \frac{1}{B} \sum_{b=1}^B f_b(\underline{x}; \underline{\Theta}_b)$$

Input: $\underline{x} \in \mathbb{R}^D$

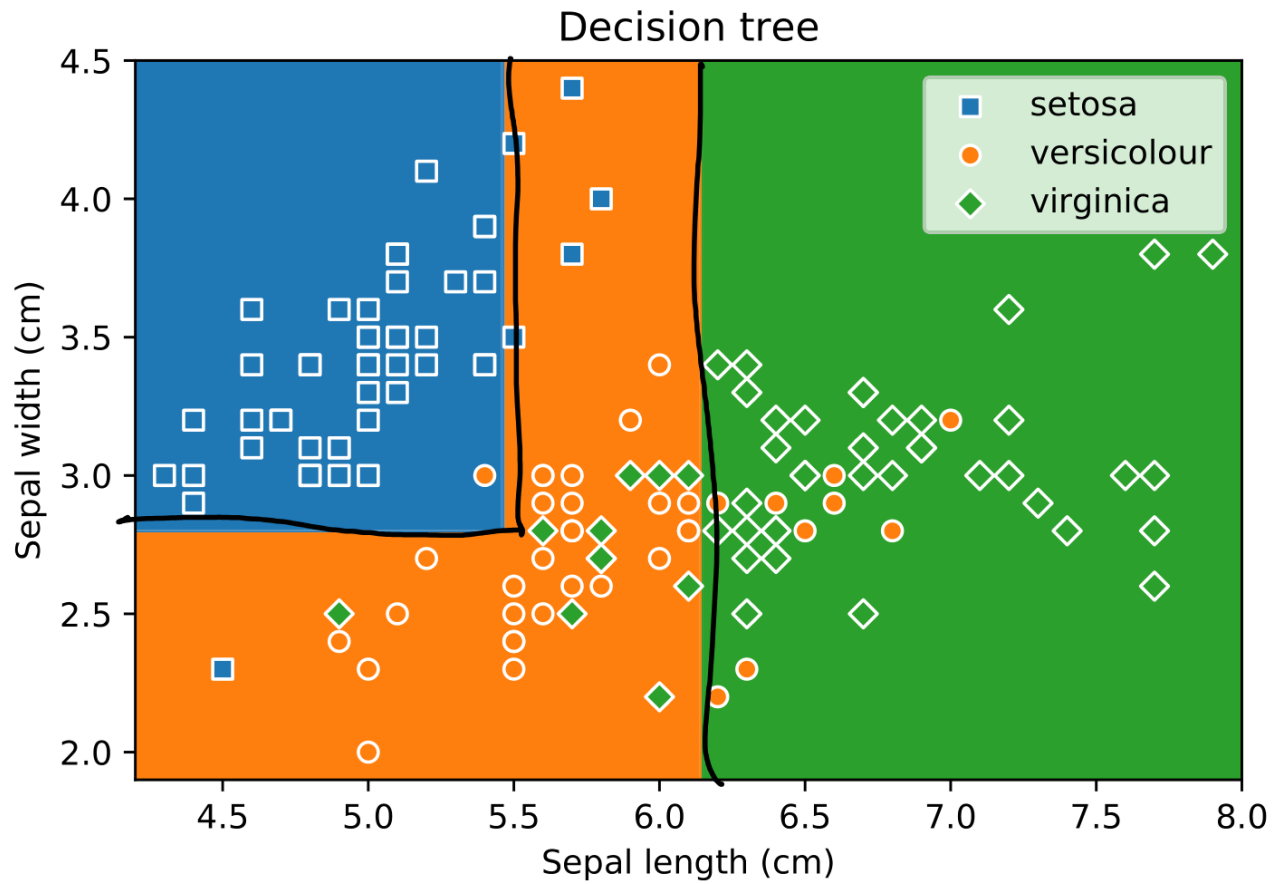
$x_1, x_2, x_3, \dots, x_D$



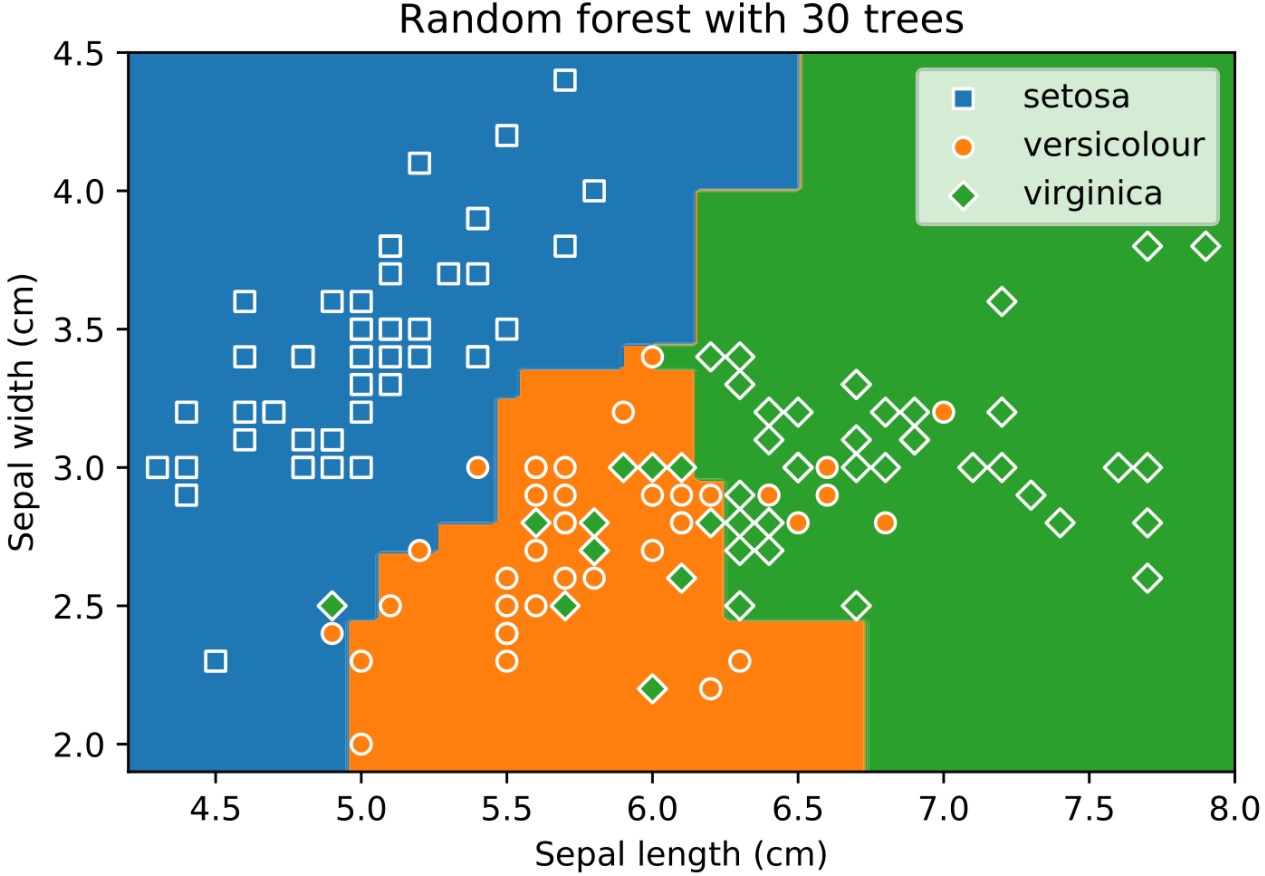
Random forest on Iris data



Random forest on Iris data



Random forest on Iris data



Ensemble methods

Boosting for regression

Herman Kamper

<http://www.kamperh.com/>

Can we combine multiple weak models (just a little bit better than random) into a "strong" model.

- Different from bagging in that we train one model, look at mistakes, only then train next
- Can be used with any weak models.

Boosting for regression

1. Initialise $r^{(n)} = y^{(n)}$ for all N training items and set $f(\mathbf{x}; \boldsymbol{\theta}) = 0$.

2. for iteration $b = 1$ to B :

(a) Fit model $f_b(\mathbf{x}; \boldsymbol{\theta}_b)$ to inputs \mathbf{X} , outputs \mathbf{r} .

(b) Update model by adding shrunken version:

$$f(\mathbf{x}; \boldsymbol{\theta}) \leftarrow f(\mathbf{x}; \boldsymbol{\theta}) + \lambda f_b(\mathbf{x}; \boldsymbol{\theta}_b)$$

(c) Update the residuals:

$$r^{(n)} \leftarrow r^{(n)} - \lambda f_b(\mathbf{x}; \boldsymbol{\theta}_b)$$

3. Final model: $f(\mathbf{x}; \boldsymbol{\theta}) = \sum_{b=1}^B \lambda f_b(\mathbf{x}; \boldsymbol{\theta}_b)$

$$\underline{\mathbf{X}} = \begin{bmatrix} - (\underline{\mathbf{x}}^{(1)})^T - \\ - (\underline{\mathbf{x}}^{(2)})^T - \\ \vdots \\ - (\underline{\mathbf{x}}^{(N)})^T - \end{bmatrix}; \quad \underline{\mathbf{r}} = \begin{bmatrix} r^{(1)} \\ r^{(2)} \\ \vdots \\ r^{(N)} \end{bmatrix}$$

At $b=1$:

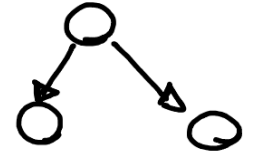
$\underline{\mathbf{r}} = \underline{\mathbf{y}}$, so we are just fitting a model to inputs $\underline{\mathbf{X}}$, outputs $\underline{\mathbf{y}}$

At $b > 1$:

Fitting a model to the residuals.

Boosting for regression

Weak learner:



Decision tree stub (only one split)

1. Initialise $r^{(n)} = y^{(n)}$ for all N training items and set $f(\mathbf{x}; \boldsymbol{\theta}) = 0$.

At $b=1$: $\boldsymbol{\Gamma} = \mathbf{y}$

2. for iteration $b = 1$ to B :

(a) Fit model $f_b(\mathbf{x}; \boldsymbol{\theta}_b)$ to inputs \mathbf{X} , outputs \mathbf{r} .

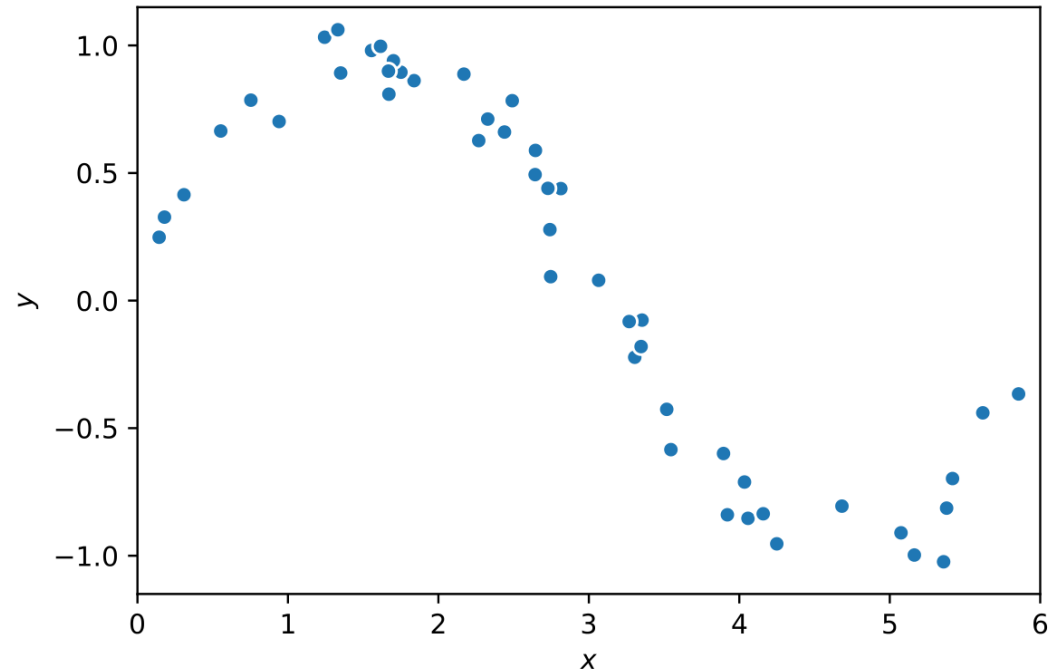
(b) Update model by adding shrunken version:

$$f(\mathbf{x}; \boldsymbol{\theta}) \leftarrow f(\mathbf{x}; \boldsymbol{\theta}) + \lambda f_b(\mathbf{x}; \boldsymbol{\theta}_b)$$

(c) Update the residuals:

$$r^{(n)} \leftarrow r^{(n)} - \lambda f_b(\mathbf{x}; \boldsymbol{\theta}_b)$$

3. Final model: $f(\mathbf{x}; \boldsymbol{\theta}) = \sum_{b=1}^B \lambda f_b(\mathbf{x}; \boldsymbol{\theta}_b)$



Boosting for regression

1. Initialise $r^{(n)} = y^{(n)}$ for all N training items and set $f(\mathbf{x}; \boldsymbol{\theta}) = 0$.

2. for iteration $b = 1$ to B :

(a) Fit model $f_b(\mathbf{x}; \boldsymbol{\theta}_b)$ to inputs \mathbf{X} , outputs \mathbf{r} .

(b) Update model by adding shrunken version:

$$f(\mathbf{x}; \boldsymbol{\theta}) \leftarrow f(\mathbf{x}; \boldsymbol{\theta}) + \lambda f_b(\mathbf{x}; \boldsymbol{\theta}_b)$$

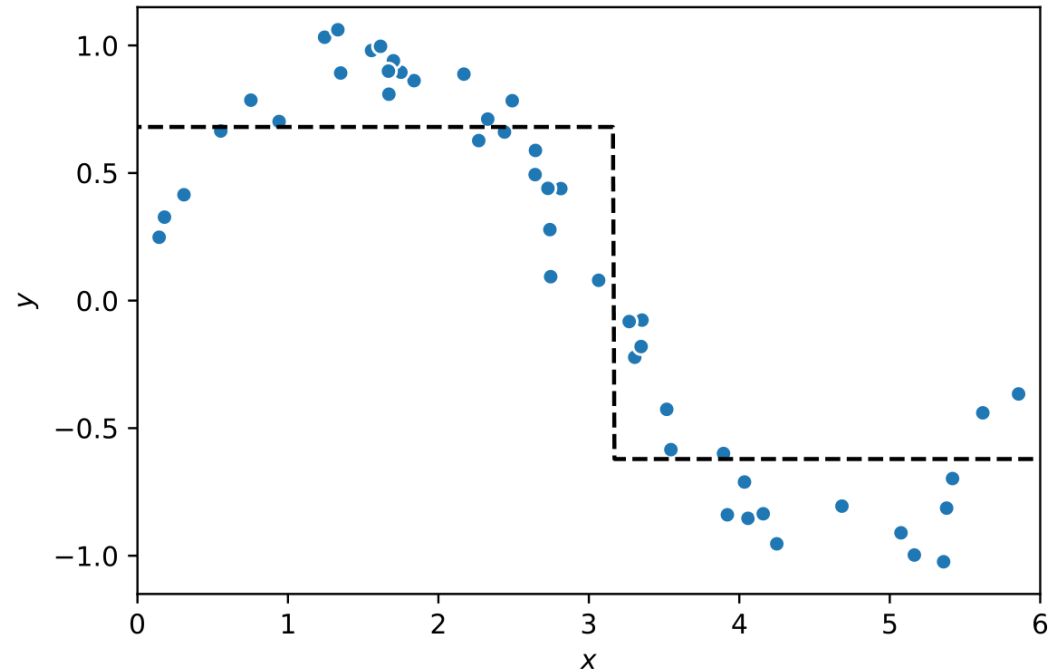
(c) Update the residuals:

$$r^{(n)} \leftarrow r^{(n)} - \lambda f_b(\mathbf{x}; \boldsymbol{\theta}_b)$$

3. Final model: $f(\mathbf{x}; \boldsymbol{\theta}) = \sum_{b=1}^B \lambda f_b(\mathbf{x}; \boldsymbol{\theta}_b)$

$$f_b(\mathbf{x}; \boldsymbol{\theta}_b)$$

Model trained at iteration $b = 1$



Boosting for regression

1. Initialise $r^{(n)} = y^{(n)}$ for all N training items and set $f(\mathbf{x}; \boldsymbol{\theta}) = 0$.

2. for iteration $b = 1$ to B :

(a) Fit model $f_b(\mathbf{x}; \boldsymbol{\theta}_b)$ to inputs \mathbf{X} , outputs \mathbf{r} .

(b) Update model by adding shrunken version:
 $f(\mathbf{x}; \boldsymbol{\theta}) \leftarrow f(\mathbf{x}; \boldsymbol{\theta}) + \lambda f_b(\mathbf{x}; \boldsymbol{\theta}_b)$

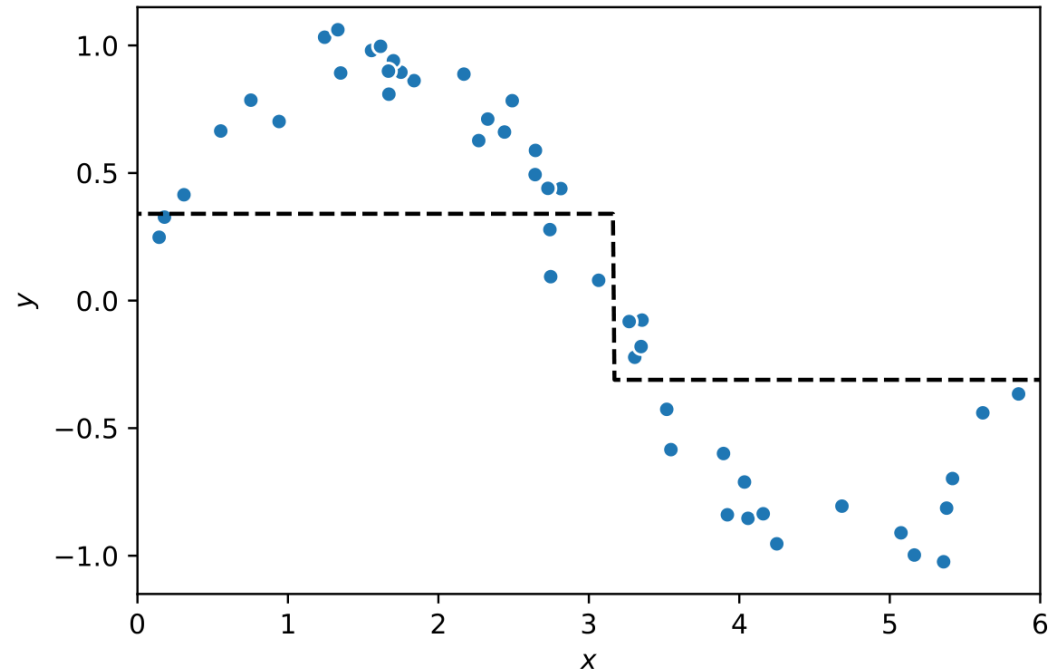
(c) Update the residuals:

$$r^{(n)} \leftarrow r^{(n)} - \lambda f_b(\mathbf{x}; \boldsymbol{\theta}_b)$$

3. Final model: $f(\mathbf{x}; \boldsymbol{\theta}) = \sum_{b=1}^B \lambda f_b(\mathbf{x}; \boldsymbol{\theta}_b)$

$$\lambda f_b(\mathbf{x}; \boldsymbol{\theta}_b)$$

Shrunken model output at iteration $b = 1$



Boosting for regression

1. Initialise $r^{(n)} = y^{(n)}$ for all N training items and set $f(\mathbf{x}; \boldsymbol{\theta}) = 0$.

2. for iteration $b = 1$ to B :

(a) Fit model $f_b(\mathbf{x}; \boldsymbol{\theta}_b)$ to inputs \mathbf{X} , outputs \mathbf{r} .

(b) Update model by adding shrunken version:

$$f(\mathbf{x}; \boldsymbol{\theta}) \leftarrow f(\mathbf{x}; \boldsymbol{\theta}) + \lambda f_b(\mathbf{x}; \boldsymbol{\theta}_b)$$

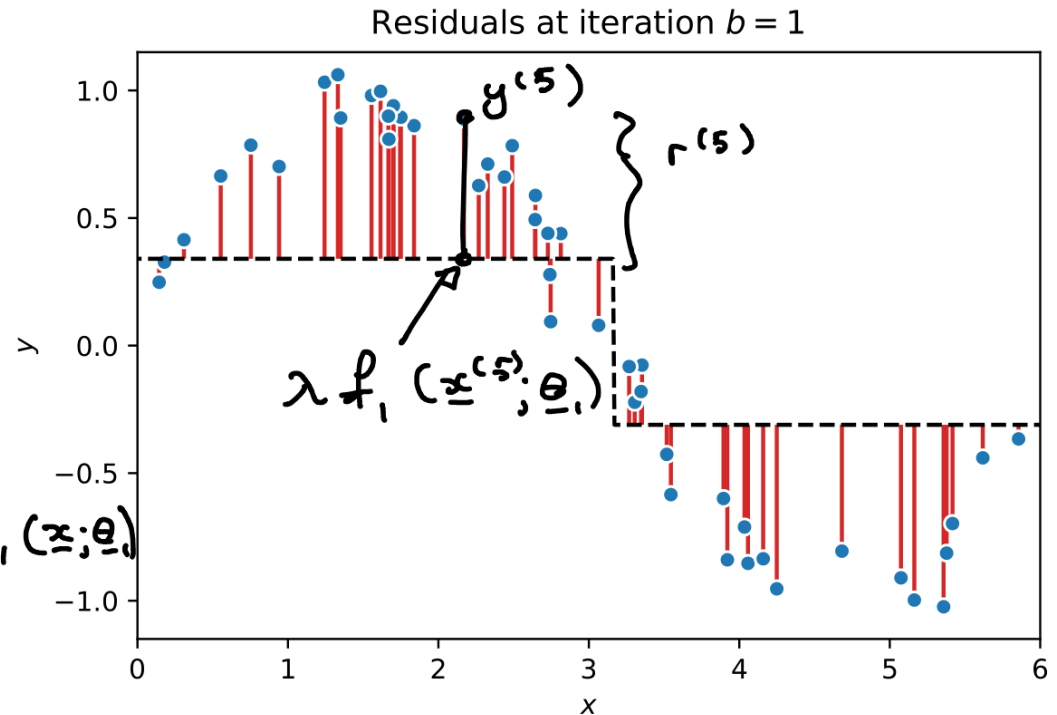
(c) Update the residuals:

$$r^{(n)} \leftarrow r^{(n)} - \lambda f_b(\mathbf{x}; \boldsymbol{\theta}_b)$$

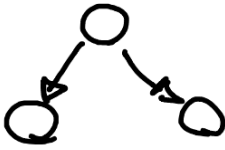
At $b=1$:

$$r^{(n)} \leftarrow y^{(n)} - \lambda f_1(\mathbf{x}; \boldsymbol{\theta}_1)$$

3. Final model: $f(\mathbf{x}; \boldsymbol{\theta}) = \sum_{b=1}^B \lambda f_b(\mathbf{x}; \boldsymbol{\theta}_b)$



Boosting for regression



1. Initialise $r^{(n)} = y^{(n)}$ for all N training items and set $f(\mathbf{x}; \boldsymbol{\theta}) = 0$.

2. for iteration $b = 1$ to B :

$b=2$

(a) Fit model $f_b(\mathbf{x}; \boldsymbol{\theta}_b)$ to inputs \mathbf{X} , outputs \mathbf{r} .

(b) Update model by adding shrunken version:

$$f(\mathbf{x}; \boldsymbol{\theta}) \leftarrow f(\mathbf{x}; \boldsymbol{\theta}) + \lambda f_b(\mathbf{x}; \boldsymbol{\theta}_b)$$

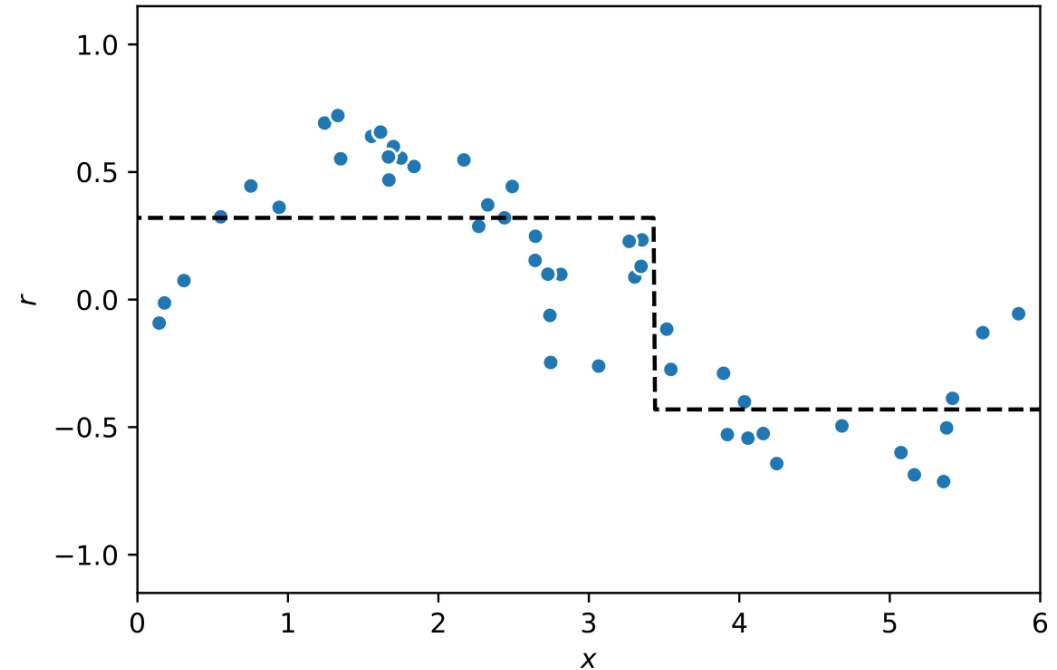
(c) Update the residuals:

$$r^{(n)} \leftarrow r^{(n)} - \lambda f_b(\mathbf{x}; \boldsymbol{\theta}_b)$$

3. Final model: $f(\mathbf{x}; \boldsymbol{\theta}) = \sum_{b=1}^B \lambda f_b(\mathbf{x}; \boldsymbol{\theta}_b)$

$$f_2(x; \boldsymbol{\theta}_2)$$

Model trained at iteration $b = 2$



Boosting for regression

1. Initialise $r^{(n)} = y^{(n)}$ for all N training items and set $f(\mathbf{x}; \boldsymbol{\theta}) = 0$.

2. for iteration $b = 1$ to B :

$b = 2$

(a) Fit model $f_b(\mathbf{x}; \boldsymbol{\theta}_b)$ to inputs \mathbf{X} , outputs \mathbf{r} .

(b) Update model by adding shrunken version:

$$f(\mathbf{x}; \boldsymbol{\theta}) \leftarrow f(\mathbf{x}; \boldsymbol{\theta}) + \lambda f_b(\mathbf{x}; \boldsymbol{\theta}_b)$$

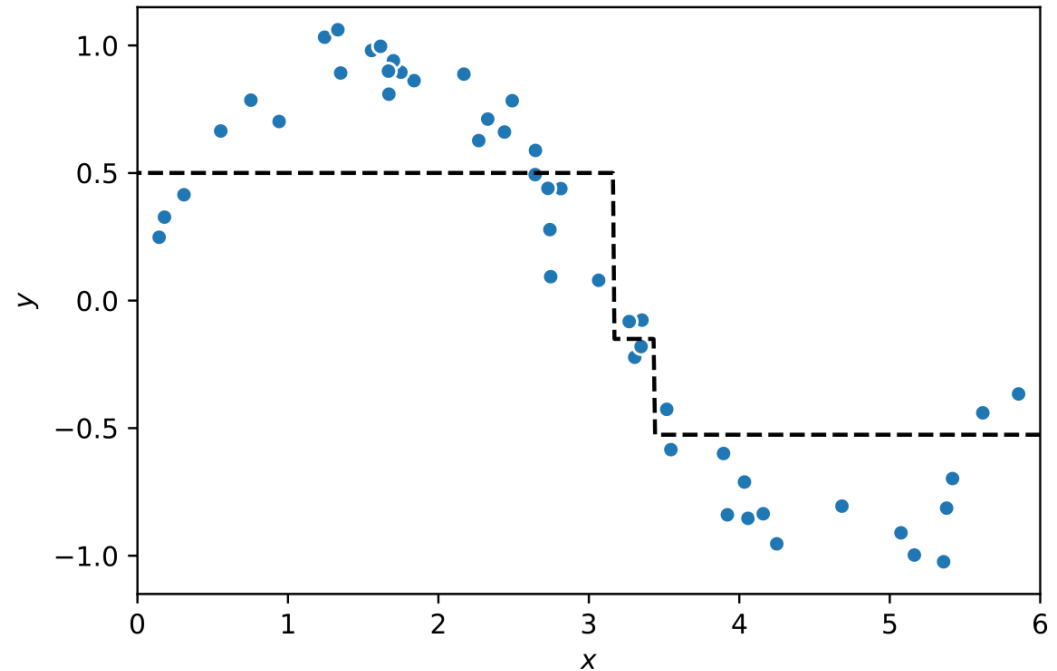
(c) Update the residuals:

$$r^{(n)} \leftarrow r^{(n)} - \lambda f_b(\mathbf{x}; \boldsymbol{\theta}_b)$$

3. Final model: $f(\mathbf{x}; \boldsymbol{\theta}) = \sum_{b=1}^B \lambda f_b(\mathbf{x}; \boldsymbol{\theta}_b)$

$$\lambda f_2(\mathbf{x}; \boldsymbol{\theta}_2) + \lambda f_1(\mathbf{x}; \boldsymbol{\theta}_1)$$

Combined model output at iteration $b = 2$



[Compare to 3 slides back]

Boosting for regression

1. Initialise $r^{(n)} = y^{(n)}$ for all N training items and set $f(\mathbf{x}; \boldsymbol{\theta}) = 0$.

2. for iteration $b = 1$ to B :

$b=3$

(a) Fit model $f_b(\mathbf{x}; \boldsymbol{\theta}_b)$ to inputs \mathbf{X} , outputs \mathbf{r} .

(b) Update model by adding shrunken version:

$$f(\mathbf{x}; \boldsymbol{\theta}) \leftarrow f(\mathbf{x}; \boldsymbol{\theta}) + \lambda f_b(\mathbf{x}; \boldsymbol{\theta}_b)$$

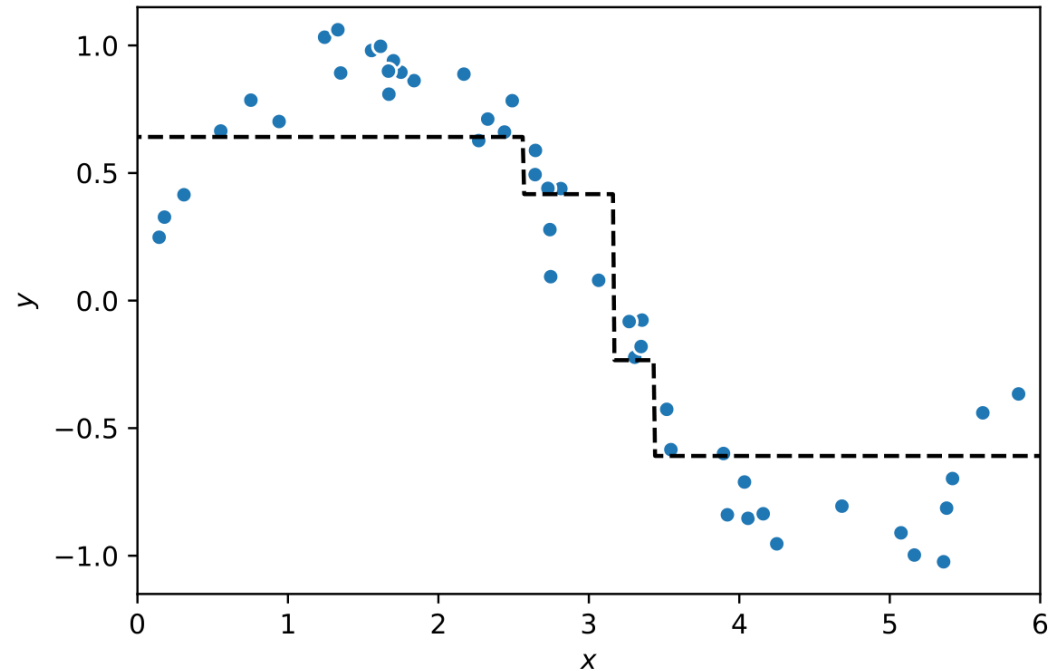
(c) Update the residuals:

$$r^{(n)} \leftarrow r^{(n)} - \lambda f_b(\mathbf{x}; \boldsymbol{\theta}_b)$$

3. Final model: $f(\mathbf{x}; \boldsymbol{\theta}) = \sum_{b=1}^B \lambda f_b(\mathbf{x}; \boldsymbol{\theta}_b)$

$$\lambda f_3(\mathbf{x}; \boldsymbol{\theta}_3) + \lambda f_2(\mathbf{x}; \boldsymbol{\theta}_2) + \lambda f_1(\mathbf{x}; \boldsymbol{\theta}_1)$$

Combined model output at iteration $b = 3$



Boosting for regression

1. Initialise $r^{(n)} = y^{(n)}$ for all N training items and set $f(\mathbf{x}; \boldsymbol{\theta}) = 0$.

2. for iteration $b = 1$ to B :

(a) Fit model $f_b(\mathbf{x}; \boldsymbol{\theta}_b)$ to inputs \mathbf{X} , outputs \mathbf{r} .

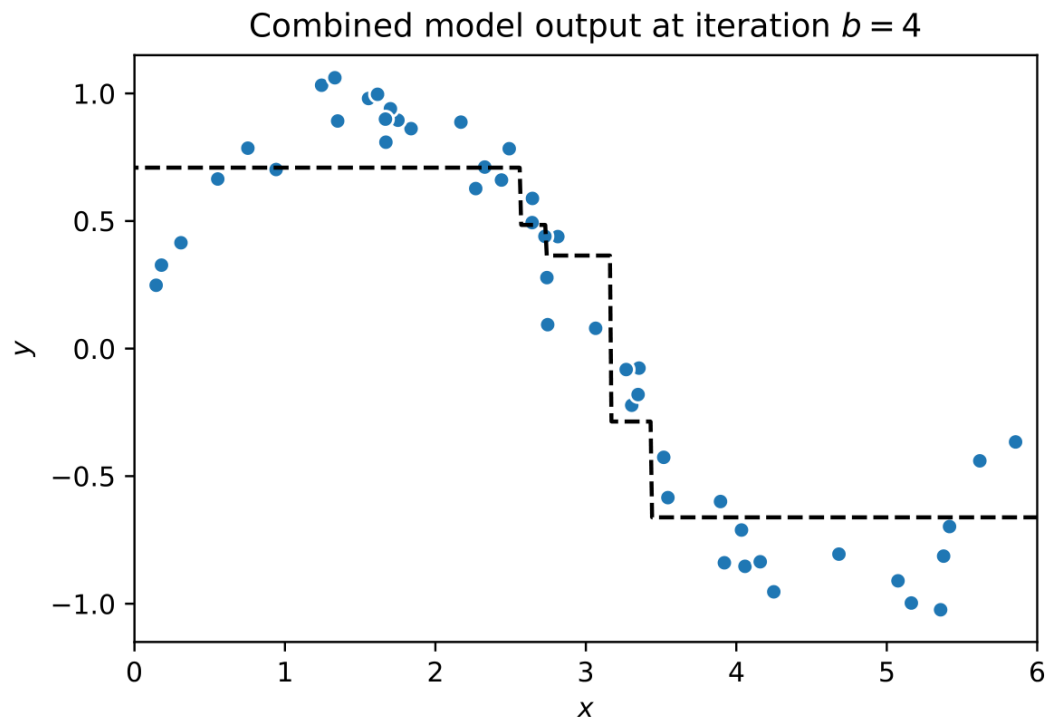
(b) Update model by adding shrunken version:

$$f(\mathbf{x}; \boldsymbol{\theta}) \leftarrow f(\mathbf{x}; \boldsymbol{\theta}) + \lambda f_b(\mathbf{x}; \boldsymbol{\theta}_b)$$

(c) Update the residuals:

$$r^{(n)} \leftarrow r^{(n)} - \lambda f_b(\mathbf{x}; \boldsymbol{\theta}_b)$$

3. Final model: $f(\mathbf{x}; \boldsymbol{\theta}) = \sum_{b=1}^B \lambda f_b(\mathbf{x}; \boldsymbol{\theta}_b)$



Boosting for regression

1. Initialise $r^{(n)} = y^{(n)}$ for all N training items and set $f(\mathbf{x}; \boldsymbol{\theta}) = 0$.

2. for iteration $b = 1$ to B :

(a) Fit model $f_b(\mathbf{x}; \boldsymbol{\theta}_b)$ to inputs \mathbf{X} , outputs \mathbf{r} .

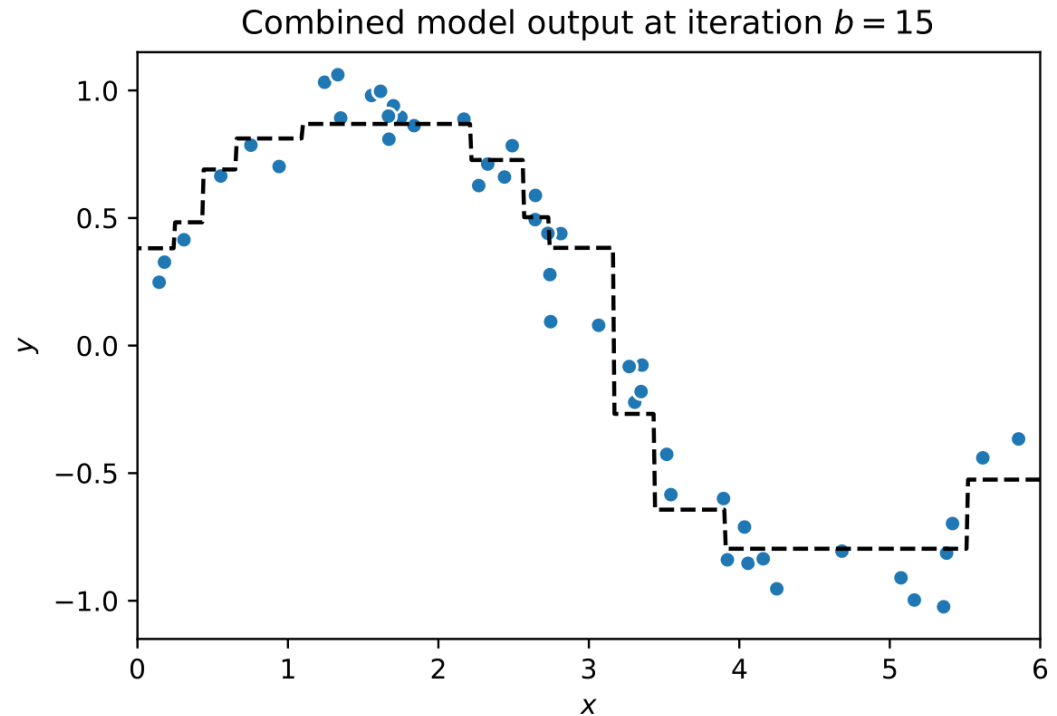
(b) Update model by adding shrunken version:

$$f(\mathbf{x}; \boldsymbol{\theta}) \leftarrow f(\mathbf{x}; \boldsymbol{\theta}) + \lambda f_b(\mathbf{x}; \boldsymbol{\theta}_b)$$

(c) Update the residuals:

$$r^{(n)} \leftarrow r^{(n)} - \lambda f_b(\mathbf{x}; \boldsymbol{\theta}_b)$$

3. Final model: $f(\mathbf{x}; \boldsymbol{\theta}) = \sum_{b=1}^B \lambda f_b(\mathbf{x}; \boldsymbol{\theta}_b)$



Ensemble methods

AdaBoost: Boosting for classification

Herman Kamper

<http://www.kamperh.com/>

Combine weak
(slightly better
than random)

- Binary classification
- Building block (principles) for classification models often used in practice

AdaBoost: Boosting for classification

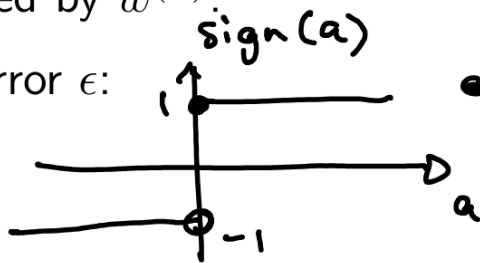
1. Initialise training item weights $w^{(n)} = 1$ for all N training items.

2. for iteration $b = 1$ to B :

(a) Fit model $f_b(\mathbf{x}; \theta_b)$ so that it minimises classification error weighted by $w^{(n)}$

(b) Set model weight using error ϵ :

$$\lambda_b = \frac{1}{2} \log \left(\frac{1 - \epsilon}{\epsilon} \right)$$



(c) Update training item weights:

$w^{(n)} \leftarrow w^{(n)} e^{-\lambda_b}$ if $f_b(\mathbf{x}^{(n)}; \theta_b)$ correct

$w^{(n)} \leftarrow w^{(n)} e^{\lambda_b}$ if $f_b(\mathbf{x}^{(n)}; \theta_b)$ incorrect

3. Final model: $f(\mathbf{x}; \theta) = \text{sign} \left[\sum_{b=1}^B \lambda_b f_b(\mathbf{x}; \theta_b) \right]$

Setting:

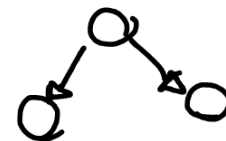
- Binary classification: $y \in \{-1, 1\}$
- Going to train B models, each $f_b(\mathbf{x}; \theta_b) \in \{-1, 1\}$
- Going to combine weighted votes:

$$f(\mathbf{x}; \theta) = \text{sign} \left[\sum_{b=1}^B \lambda_b f_b(\mathbf{x}; \theta_b) \right]$$

$$\begin{matrix} w^{(1)} \\ w^{(2)} \\ \vdots \\ w^{(N)} \end{matrix} \begin{bmatrix} -(\mathbf{x}^{(1)})^T & - & y^{(1)} \\ -(\mathbf{x}^{(2)}) & - & y^{(2)} \\ \vdots & & \vdots \end{bmatrix}$$

AdaBoost: Boosting for classification

Weak classifier:



1. Initialise training item weights $w^{(n)} = 1$ for all N training items.

2. for iteration $b = 1$ to B :

(a) Fit model $f_b(\mathbf{x}; \theta_b)$ so that it minimises classification error weighted by $w^{(n)}$.

(b) Set model weight using error ϵ :

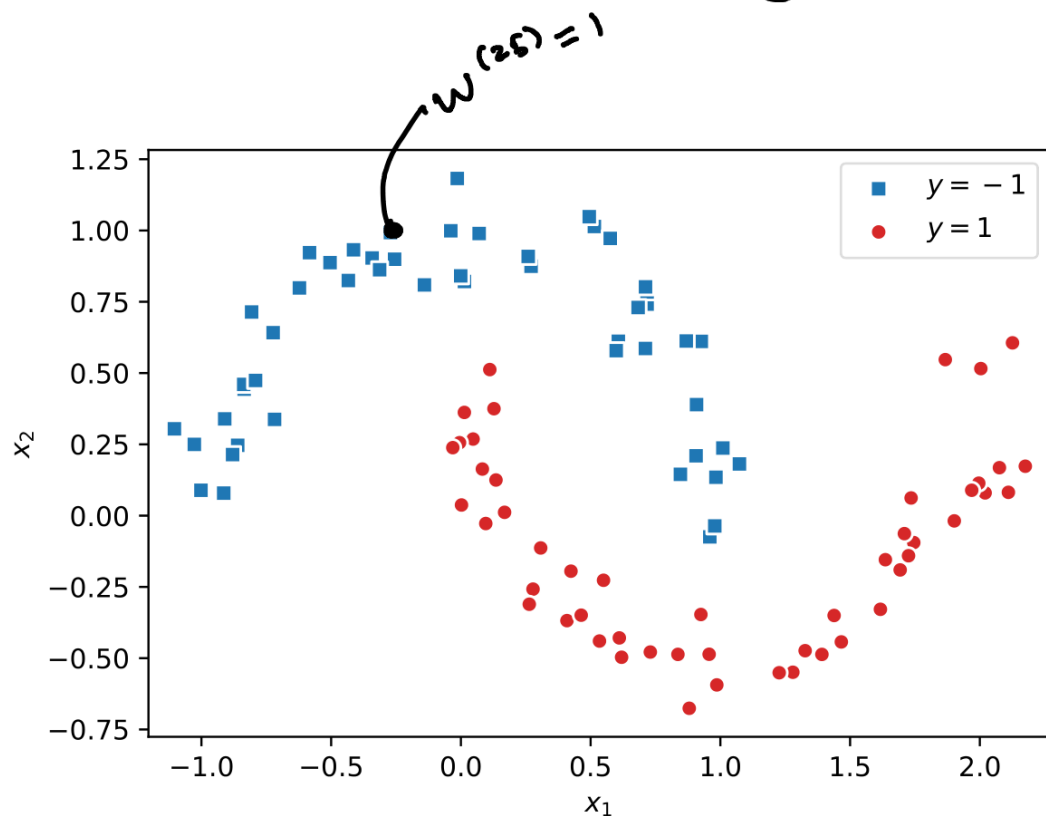
$$\lambda_b = \frac{1}{2} \log \left(\frac{1 - \epsilon}{\epsilon} \right)$$

(c) Update training item weights:

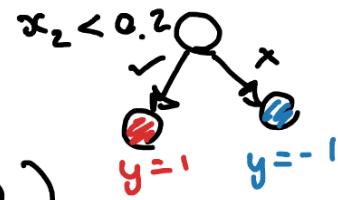
$$w^{(n)} \leftarrow w^{(n)} e^{-\lambda_b} \text{ if } f_b(\mathbf{x}^{(n)}; \theta_b) \text{ correct}$$

$$w^{(n)} \leftarrow w^{(n)} e^{\lambda_b} \text{ if } f_b(\mathbf{x}^{(n)}; \theta_b) \text{ incorrect}$$

3. Final model: $f(\mathbf{x}; \theta) = \text{sign} \left[\sum_{b=1}^B \lambda_b f_b(\mathbf{x}; \theta_b) \right]$



AdaBoost: Boosting for classification



$$f_1(\mathbf{x}; \theta_1)$$

1. Initialise training item weights $w^{(n)} = 1$ for all N training items.

2. for iteration $b = 1$ to B :

✓ (a) Fit model $f_b(\mathbf{x}; \theta_b)$ so that it minimises classification error weighted by $w^{(n)}$.

✓ (b) Set model weight using error ϵ :

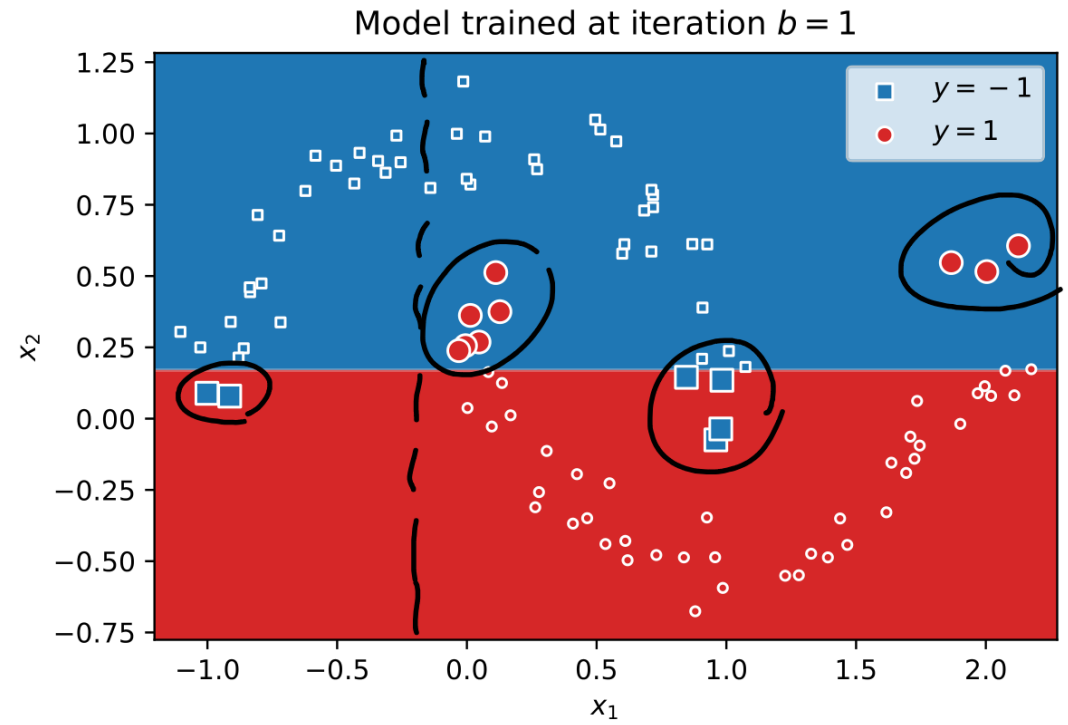
$$\lambda_b = \frac{1}{2} \log \left(\frac{1 - \epsilon}{\epsilon} \right) \quad \lambda_1 = 0.87$$

✓ (c) Update training item weights:

$$w^{(n)} \leftarrow w^{(n)} e^{-\lambda_b} \text{ if } f_b(\mathbf{x}^{(n)}; \theta_b) \text{ correct}$$

$$w^{(n)} \leftarrow w^{(n)} e^{\lambda_b} \text{ if } f_b(\mathbf{x}^{(n)}; \theta_b) \text{ incorrect}$$

3. Final model: $f(\mathbf{x}; \theta) = \text{sign} \left[\sum_{b=1}^B \lambda_b f_b(\mathbf{x}; \theta_b) \right]$



AdaBoost: Boosting for classification

1. Initialise training item weights $w^{(n)} = 1$ for all N training items.

2. for iteration $b = 1$ to B :

$b=2$

(a) Fit model $f_b(\mathbf{x}; \theta_b)$ so that it minimises classification error weighted by $w^{(n)}$.

(b) Set model weight using error ϵ :

$$\lambda_b = \frac{1}{2} \log \left(\frac{1 - \epsilon}{\epsilon} \right) \quad \lambda_2 = 0.52$$

(c) Update training item weights:

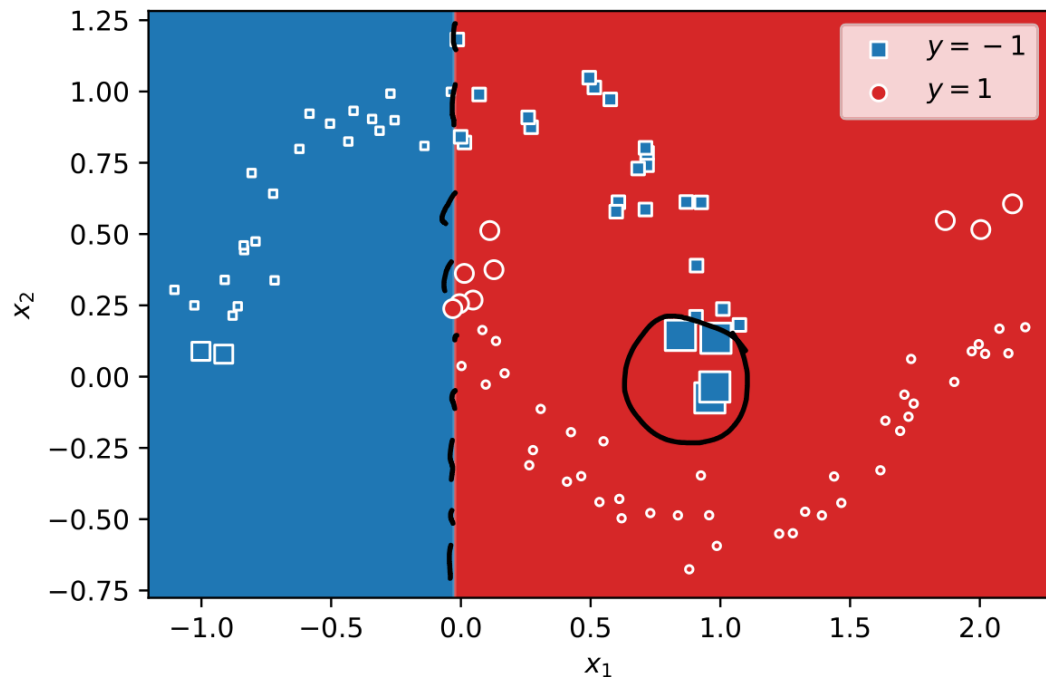
$$w^{(n)} \leftarrow w^{(n)} e^{-\lambda_b} \text{ if } f_b(\mathbf{x}^{(n)}; \theta_b) \text{ correct}$$

$$w^{(n)} \leftarrow w^{(n)} e^{\lambda_b} \text{ if } f_b(\mathbf{x}^{(n)}; \theta_b) \text{ incorrect}$$

3. Final model: $f(\mathbf{x}; \theta) = \text{sign} \left[\sum_{b=1}^B \lambda_b f_b(\mathbf{x}; \theta_b) \right]$

$$f_2(\mathbf{x}; \theta_2)$$

Model trained at iteration $b = 2$



AdaBoost: Boosting for classification

1. Initialise training item weights $w^{(n)} = 1$ for all N training items.

2. for iteration $b = 1$ to B :

$b=3$

(a) Fit model $f_b(\mathbf{x}; \boldsymbol{\theta}_b)$ so that it minimises classification error weighted by $w^{(n)}$.

(b) Set model weight using error ϵ :

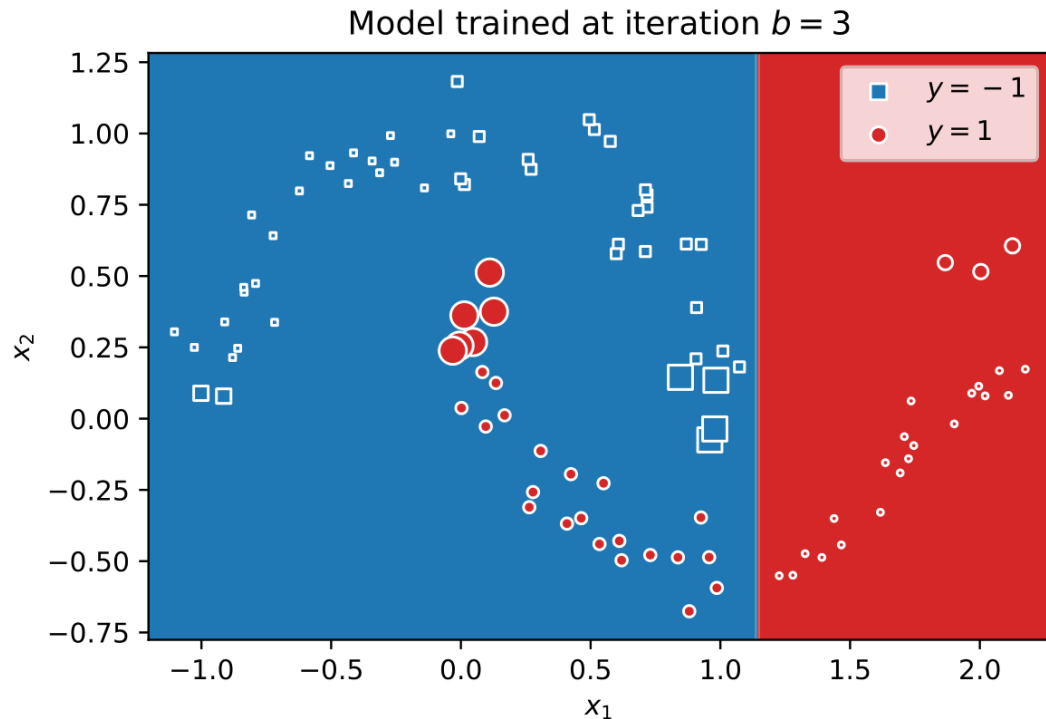
$$\lambda_b = \frac{1}{2} \log \left(\frac{1 - \epsilon}{\epsilon} \right)$$

(c) Update training item weights:

$$w^{(n)} \leftarrow w^{(n)} e^{-\lambda_b} \text{ if } f_b(\mathbf{x}^{(n)}; \boldsymbol{\theta}_b) \text{ correct}$$

$$w^{(n)} \leftarrow w^{(n)} e^{\lambda_b} \text{ if } f_b(\mathbf{x}^{(n)}; \boldsymbol{\theta}_b) \text{ incorrect}$$

3. Final model: $f(\mathbf{x}; \boldsymbol{\theta}) = \text{sign} \left[\sum_{b=1}^B \lambda_b f_b(\mathbf{x}; \boldsymbol{\theta}_b) \right]$



AdaBoost: Boosting for classification

1. Initialise training item weights $w^{(n)} = 1$ for all N training items.

2. for iteration $b = 1$ to B :

(a) Fit model $f_b(\mathbf{x}; \boldsymbol{\theta}_b)$ so that it minimises classification error weighted by $w^{(n)}$.

(b) Set model weight using error ϵ :

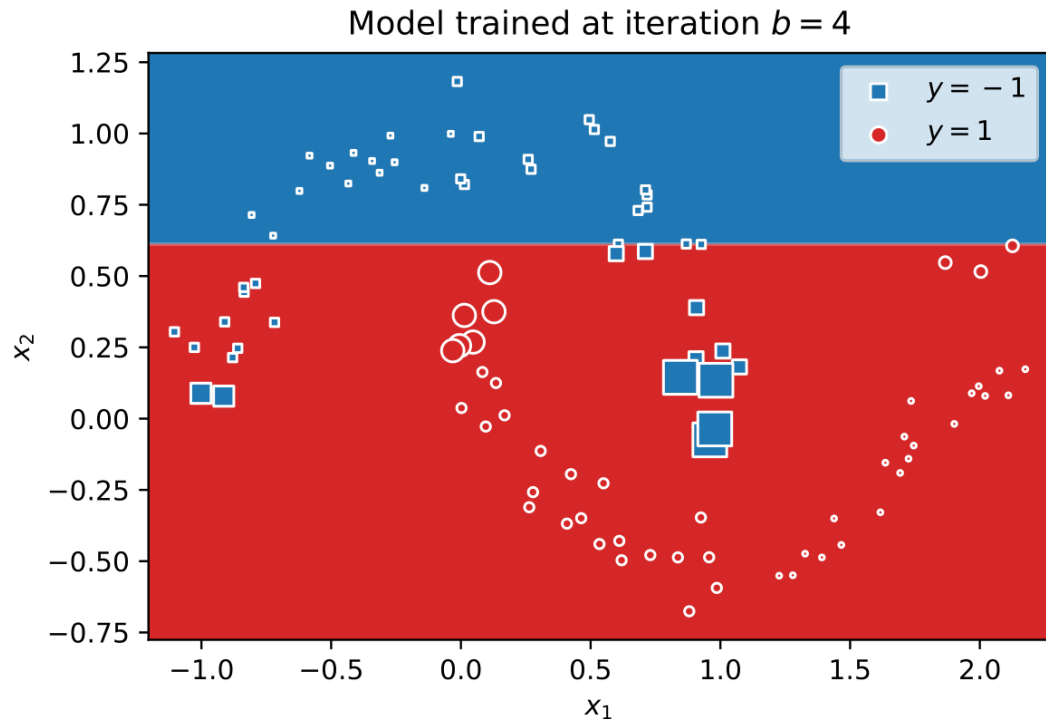
$$\lambda_b = \frac{1}{2} \log \left(\frac{1 - \epsilon}{\epsilon} \right)$$

(c) Update training item weights:

$$w^{(n)} \leftarrow w^{(n)} e^{-\lambda_b} \text{ if } f_b(\mathbf{x}^{(n)}; \boldsymbol{\theta}_b) \text{ correct}$$

$$w^{(n)} \leftarrow w^{(n)} e^{\lambda_b} \text{ if } f_b(\mathbf{x}^{(n)}; \boldsymbol{\theta}_b) \text{ incorrect}$$

3. Final model: $f(\mathbf{x}; \boldsymbol{\theta}) = \text{sign} \left[\sum_{b=1}^B \lambda_b f_b(\mathbf{x}; \boldsymbol{\theta}_b) \right]$



AdaBoost: Boosting for classification

$$f(\underline{x}; \underline{\theta}) = \lambda_1 f_1(\underline{x}; \underline{\theta}_1) + \lambda_2 f_2(\underline{x}; \underline{\theta}_2) + \dots + \lambda_{20} f_{20}(\underline{x}; \underline{\theta}_{20})$$

1. Initialise training item weights $w^{(n)} = 1$ for all N training items.

2. for iteration $b = 1$ to B :

(a) Fit model $f_b(\mathbf{x}; \boldsymbol{\theta}_b)$ so that it minimises classification error weighted by $w^{(n)}$.

(b) Set model weight using error ϵ :

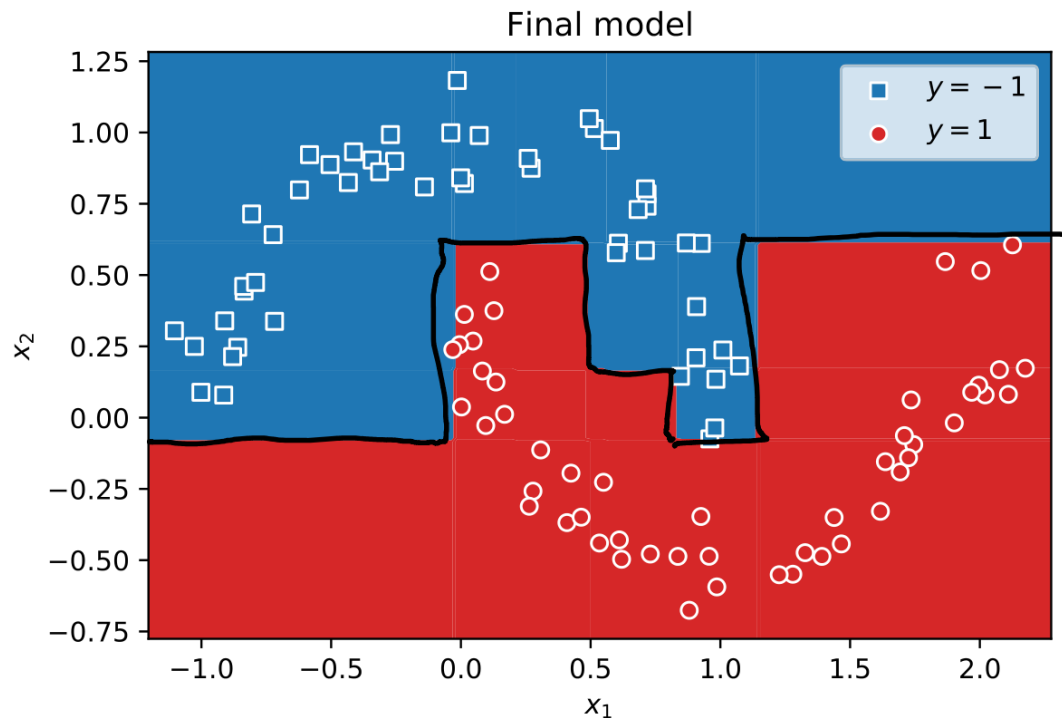
$$\lambda_b = \frac{1}{2} \log \left(\frac{1 - \epsilon}{\epsilon} \right)$$

(c) Update training item weights:

$$w^{(n)} \leftarrow w^{(n)} e^{-\lambda_b} \text{ if } f_b(\mathbf{x}^{(n)}; \boldsymbol{\theta}_b) \text{ correct}$$

$$w^{(n)} \leftarrow w^{(n)} e^{\lambda_b} \text{ if } f_b(\mathbf{x}^{(n)}; \boldsymbol{\theta}_b) \text{ incorrect}$$

3. Final model: $f(\mathbf{x}; \boldsymbol{\theta}) = \text{sign} \left[\sum_{b=1}^B \lambda_b f_b(\mathbf{x}; \boldsymbol{\theta}_b) \right]$



AdaBoost: Boosting for classification

1. Initialise training item weights $w^{(n)} = 1$ for all N training items.

(a) Minimise $\sum_{n=1}^N w^{(n)} I \{ y^{(n)} \neq f_b(x^{(n)}; \theta_b) \}$

2. for iteration $b = 1$ to B :

(a) Fit model $f_b(\mathbf{x}; \theta_b)$ so that it minimises classification error weighted by $w^{(n)}$.

(b)
$$\epsilon = \frac{\sum_{n=1}^N w^{(n)} I \{ y^{(n)} \neq f_b(x^{(n)}; \theta_b) \}}{\sum_{n=1}^N w^{(n)}}$$

(b) Set model weight using error ϵ :

$$\lambda_b = \frac{1}{2} \log \left(\frac{1 - \epsilon}{\epsilon} \right)$$

If $\epsilon \rightarrow 0$ (very good classifier)
 then λ_b big
 If $\epsilon \rightarrow \frac{1}{2}$ (bad/random classifier)
 then $\lambda_b \rightarrow 0$

(c) Update training item weights:

$$w^{(n)} \leftarrow w^{(n)} e^{-\lambda_b} \text{ if } f_b(\mathbf{x}^{(n)}; \theta_b) \text{ correct}$$

$$w^{(n)} \leftarrow w^{(n)} e^{\lambda_b} \text{ if } f_b(\mathbf{x}^{(n)}; \theta_b) \text{ incorrect}$$

(c) Correct: $w^{(n)} \leftarrow w^{(n)} \sqrt{\frac{\epsilon}{1 - \epsilon}}$

Incorrect: $w^{(n)} \leftarrow w^{(n)} \sqrt{\frac{1 - \epsilon}{\epsilon}}$

3. Final model: $f(\mathbf{x}; \theta) = \text{sign} \left[\sum_{b=1}^B \lambda_b f_b(\mathbf{x}; \theta_b) \right]$

AdaBoost: Boosting for classification

1. Initialise training item weights $w^{(n)} = 1$ for all N training items.

2. for iteration $b = 1$ to B :

(a) Fit model $f_b(\mathbf{x}; \boldsymbol{\theta}_b)$ so that it minimises classification error weighted by $w^{(n)}$.

(b) Set model weight using error ϵ :

$$\lambda_b = \frac{1}{2} \log \left(\frac{1 - \epsilon}{\epsilon} \right)$$

(c) Update training item weights:

$w^{(n)} \leftarrow w^{(n)} e^{-\lambda_b}$ if $f_b(\mathbf{x}^{(n)}; \boldsymbol{\theta}_b)$ correct

$w^{(n)} \leftarrow w^{(n)} e^{\lambda_b}$ if $f_b(\mathbf{x}^{(n)}; \boldsymbol{\theta}_b)$ incorrect

3. Final model: $f(\mathbf{x}; \boldsymbol{\theta}) = \text{sign} \left[\sum_{b=1}^B \lambda_b f_b(\mathbf{x}; \boldsymbol{\theta}_b) \right]$

Further reading:

- Raúl Rojas, "AdaBoost and the superbowl of classifiers: A tutorial introduction to adaptive boosting", 2009.
- ESL, Chapter 10