

CS 360: Machine Learning

Sara Mathieson, Sorelle Friedler

Spring 2024



HVERFORD
COLLEGE

Admin

- **Lab 3** due TONIGHT
- **Sorelle office hours** today 4-5pm in H110
- **Lab 4** posted tonight
- **Reading:** Geron Chap 7 (Ensembles) + optional reading up

Outline for Feb 15

- Finish AdaBoost
- Entropy and weighted entropy
- Gradient Boosting

Ensemble methods quiz (discuss with partner)

1. Briefly describe one advantage of an ensemble method over a single classifier.

Usually decreases our testing error

2. When using ensemble methods, we want the base classifiers that have:

- low variance, high bias
- high variance, low bias
- high variance, high bias

3. Ensembles often require us to run a base classifier on different training datasets. Name and briefly describe one method for generating multiple training datasets.

Bootstrap (sampling with replacement)

4. (Bonus) for decision trees, is it possible to use the same feature twice on different paths from the root?

Yes! (but not on the *same* path)

Outline for Feb 15

- Finish AdaBoost
- Entropy and weighted entropy
- Gradient Boosting

AdaBoost (adaptive boosting)

- **Train**

- Start with all examples weighted equally ($1/n$)
- for T iterations
 - Learn base classifier using weights
 - Change examples weights (up-weight incorrectly classified examples, down-weight correctly classified examples)

- **Test**

- Get predictions from all base classifiers
- Vote based on how well each classifier did during training

Ada Boost

set $w_i^{(1)} = \frac{1}{n}$ for $i=1, \dots, n$

for $t = 1 \dots T$:

(a) fit classifier to weighted examples

true $Y \in \{-1, 1\}$

binary classification

$h^{(t)}(\vec{x})$
 $\{-1, 1\}$

(b) compute weighted classification error

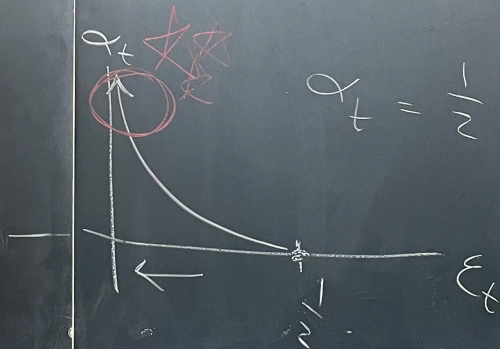
$$\epsilon_t = \sum_{i=1}^n w_i^{(t)} \mathbb{1}(y_i \neq h^{(t)}(\vec{x}_i))$$

$$0 \leq \epsilon_t \leq 1$$

(c) compute classifier score

$$\alpha_t = \frac{1}{2} \log \left(\frac{1 - \epsilon_t}{\epsilon_t} \right)$$

as $\epsilon_t \rightarrow 0, \alpha_t \rightarrow \infty$

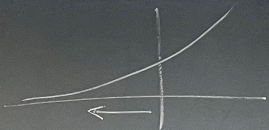


d) update weights

$$w_i^{(t+1)} = C_t w_i^{(t)} \exp(-y_i \alpha_t h^{(t)}(\vec{x}_i))$$

normalize
so weights
sum to 1

$\{-1, 1\}$ $\{-1, 1\}$



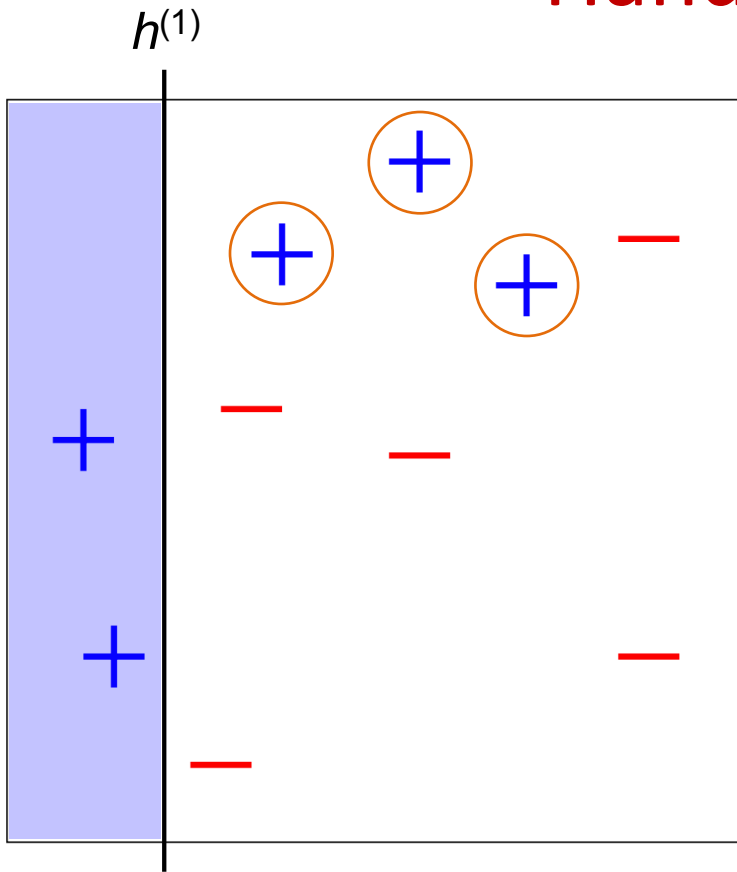
$$\exp(-y_i \alpha_t h^{(t)}(\vec{x}_i)) = \begin{cases} e^{-\alpha_t} & \text{if } y_i = h^{(t)}(\vec{x}_i) \\ & \text{(downweight)} \\ e^{\alpha_t} & \text{if } y_i \neq h^{(t)}(\vec{x}_i) \\ & \text{(upweight)} \end{cases}$$

after loop

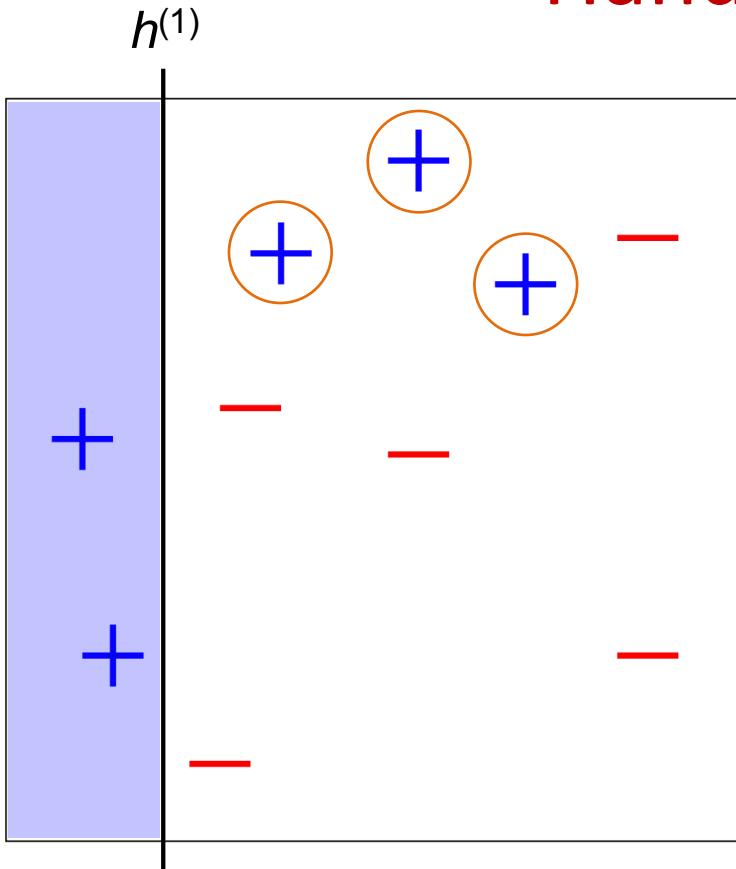
$$h(\vec{x}) = \text{sign} \left(\sum_{t=1}^T \alpha_t h^{(t)}(\vec{x}) \right)$$

Handout 8

Handout 8: Round 1



Handout 8: Round 1



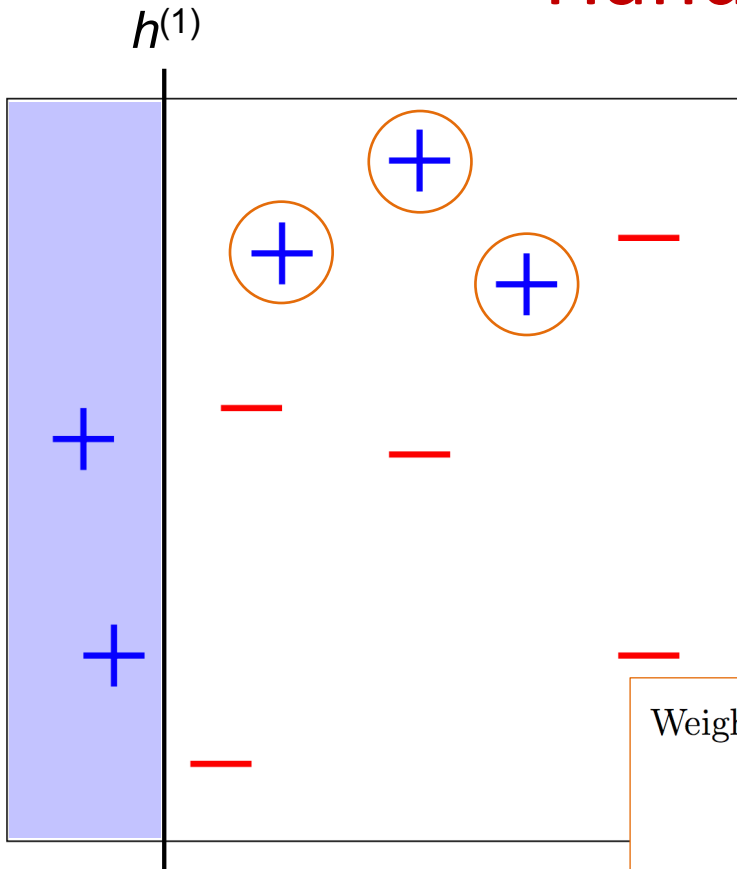
$$w_i^{(1)} = \frac{1}{10} \text{ for all } i = 1, 2, \dots, 10.$$

$$\epsilon_1 = \frac{3}{10} \text{ (three points incorrectly classified, all with weight } \frac{1}{10}\text{)}$$

$$\alpha_1 = \frac{1}{2} \ln \left(\frac{1 - \frac{3}{10}}{\frac{3}{10}} \right) = \ln \sqrt{\frac{7}{3}} \approx 0.42$$

- correctly classified: $w_i^{(2)} = c_1 \cdot \frac{1}{10} \exp \left(-\ln \sqrt{\frac{7}{3}} \right)$
- incorrectly classified: $w_i^{(2)} = c_1 \cdot \frac{1}{10} \exp \left(\ln \sqrt{\frac{7}{3}} \right)$

Handout 8: Round 1



$$w_i^{(1)} = \frac{1}{10} \text{ for all } i = 1, 2, \dots, 10.$$

$$\epsilon_1 = \frac{3}{10} \text{ (three points incorrectly classified, all with weight } \frac{1}{10})$$

$$\alpha_1 = \frac{1}{2} \ln \left(\frac{1 - \frac{3}{10}}{\frac{3}{10}} \right) = \ln \sqrt{\frac{7}{3}} \approx 0.42$$

- correctly classified: $w_i^{(2)} = c_1 \cdot \frac{1}{10} \exp \left(-\ln \sqrt{\frac{7}{3}} \right)$
- incorrectly classified: $w_i^{(2)} = c_1 \cdot \frac{1}{10} \exp \left(\ln \sqrt{\frac{7}{3}} \right)$

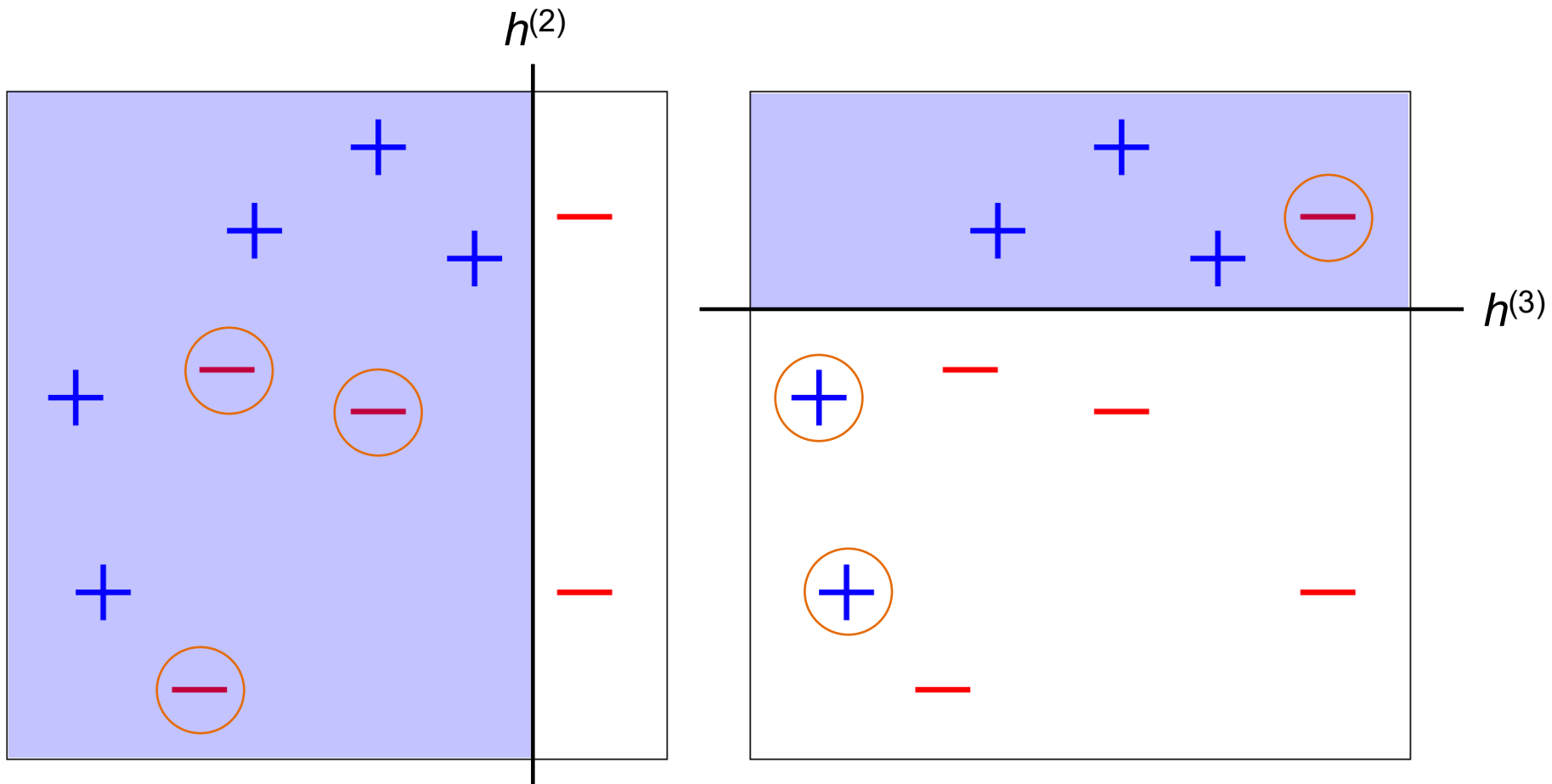
Weights must sum to 1, \Rightarrow

$$7 \cdot \frac{c_1}{10} \exp \left(-\ln \sqrt{\frac{7}{3}} \right) + 3 \cdot c_1 \cdot \frac{1}{10} \exp \left(\ln \sqrt{\frac{7}{3}} \right) = 1$$

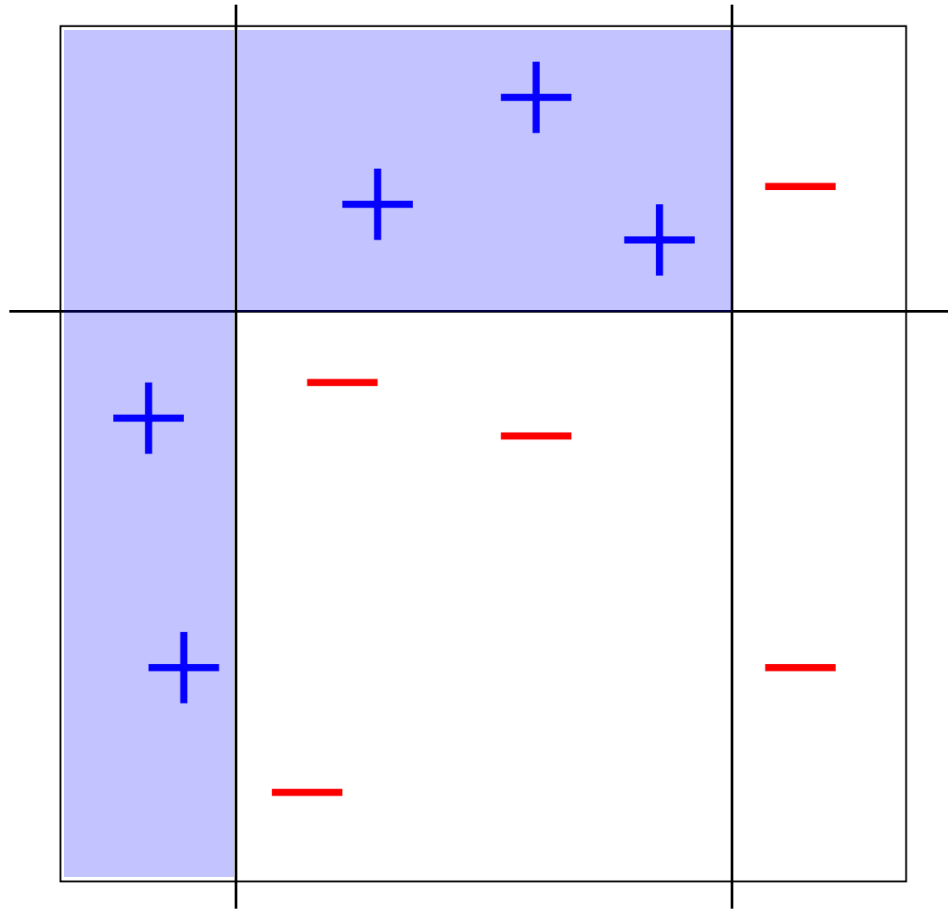
$$\Rightarrow c_1 = \frac{5}{\sqrt{21}}$$

- correctly classified: $w_i^{(2)} = \frac{5}{\sqrt{21}} \cdot \frac{1}{10} \sqrt{\frac{3}{7}} = \frac{1}{14}$ decrease!
- incorrectly classified: $w_i^{(2)} = \frac{5}{\sqrt{21}} \cdot \frac{1}{10} \sqrt{\frac{7}{3}} = \frac{1}{6}$ increase!

Handout 8: Round 2 & 3 (exercise!)



Handout 8: final classifier



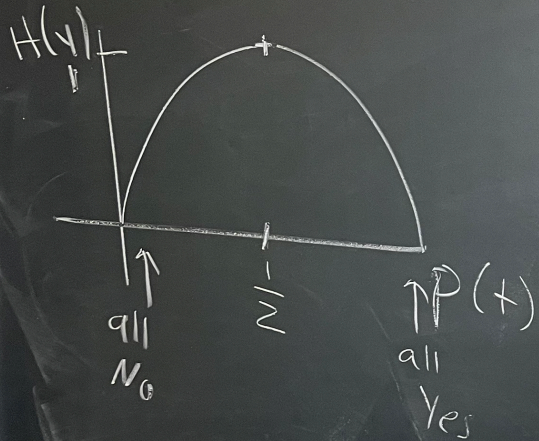
$$h(\mathbf{x}) = \text{sign}\left(0.42 \cdot h^{(1)}(\mathbf{x}) + 0.65 \cdot h^{(2)}(\mathbf{x}) + 0.92 \cdot h^{(3)}(\mathbf{x})\right)$$

Outline for Feb 15

- Finish AdaBoost
- Entropy and weighted entropy
- Gradient Boosting

Entropy & Info Gain

$$H(Y) = - \sum_{c \in \text{vals}(Y)} P(c) \log_2 P(c) = - \left(\frac{5}{14} \log_2 \frac{5}{14} + \frac{9}{14} \log_2 \frac{9}{14} \right)$$



Conditional Entropy

$$H(Y|X)$$

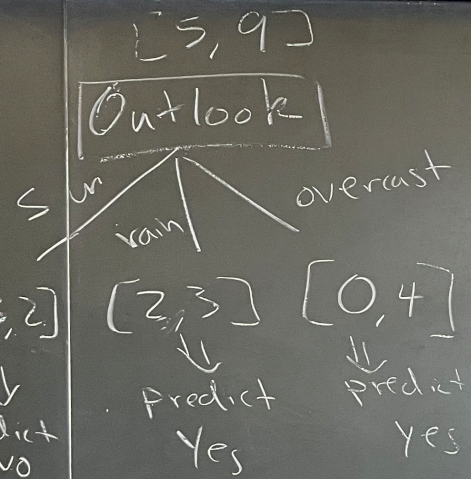
↑
one feature

$$= \sum_{v \in \text{vals}(X)} P(X=v) H(Y|X=v) = \frac{5}{14} H(Y|X=)$$

No Yes
[5, 9]

[3, 2]

↓
predict
No



$$H(Y|X=v) = - \sum_{c \in \text{vals}(Y)} P(Y=c|X=v) \log_2 P(Y=c|X=v)$$

$$P(Y=c|X=v) = \frac{\text{count}(Y=c, X=v)}{\text{count}(X=v)} = \frac{3}{5}$$

$$P(B|A) = \frac{P(A, B)}{P(A)}$$

$$\left. \begin{array}{l} c = \text{no} \\ v = \text{sun} \end{array} \right\}$$

$$H(Y|X=\text{sun}) + \frac{5}{14} H(Y|X=\text{rain}) + \frac{4}{14} H(Y|X=\text{overcast})$$

Weighted Examples

$$P(Y=c | X_j=v) = \frac{\sum_{i=1}^n w_i^{(t)} \mathbb{1}(y_i=c, x_{ij}=v)}{\sum_{i=1}^n w_i^{(t)} \mathbb{1}(x_{ij}=v)}$$

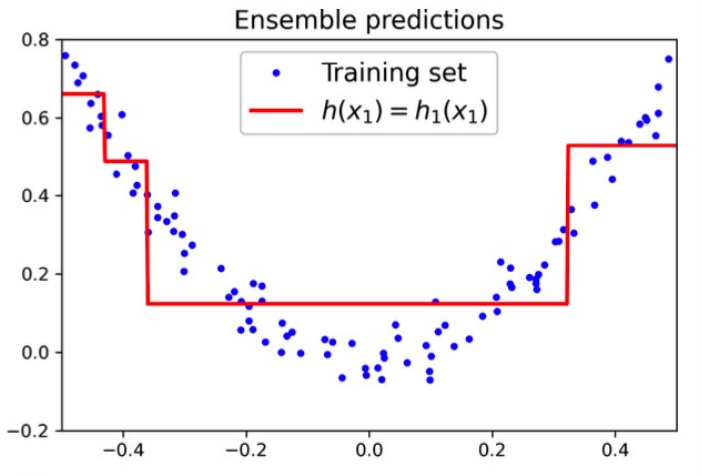
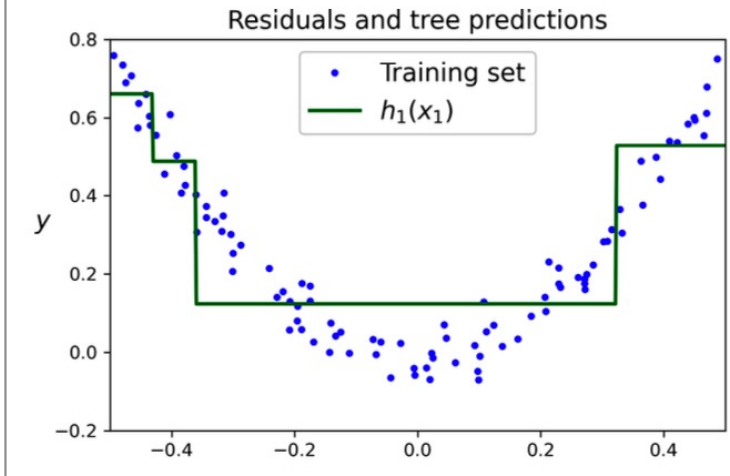
leaves

$$P(\text{leaf label is } 1) = \frac{\sum_{i=1}^n w_i^{(t)} \mathbb{1}(y_i=1)}{\sum_{i \text{ in leaf}} w_i^{(t)}}$$

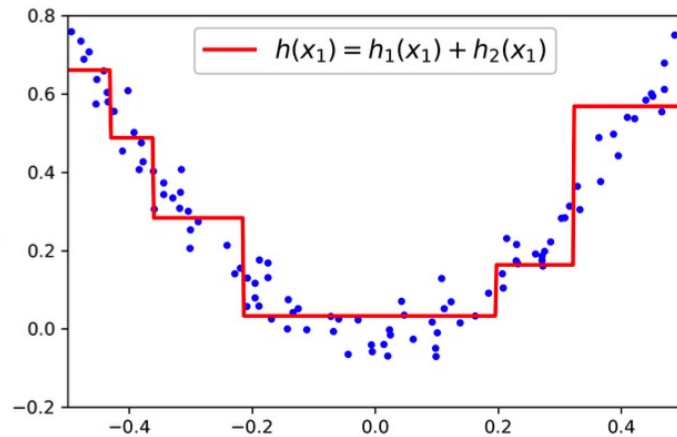
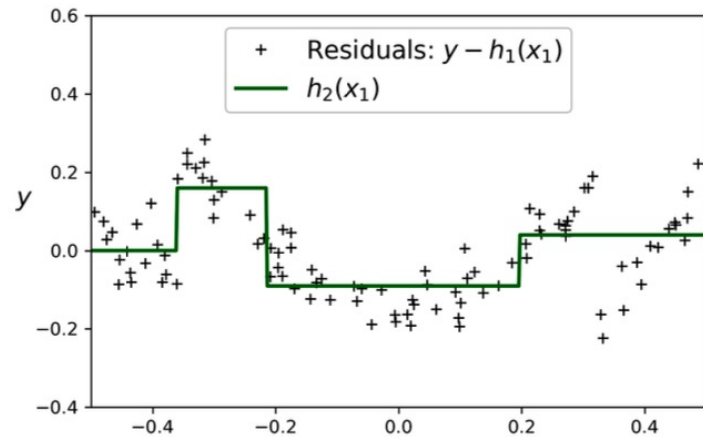
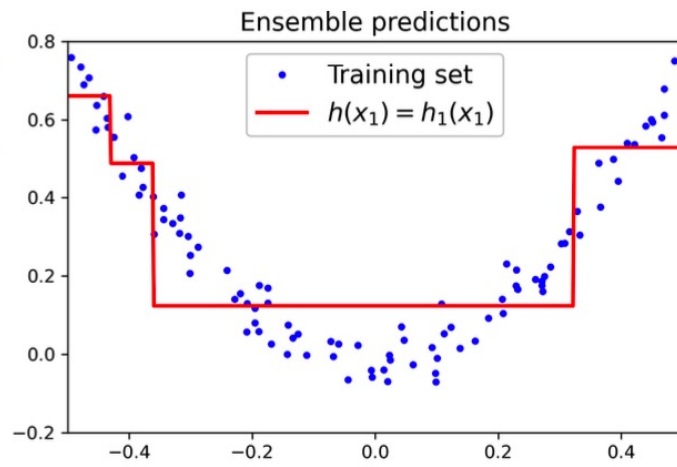
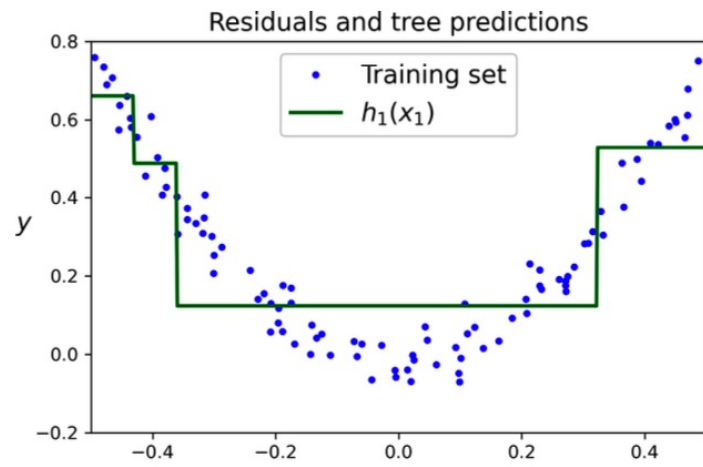


Outline for Feb 15

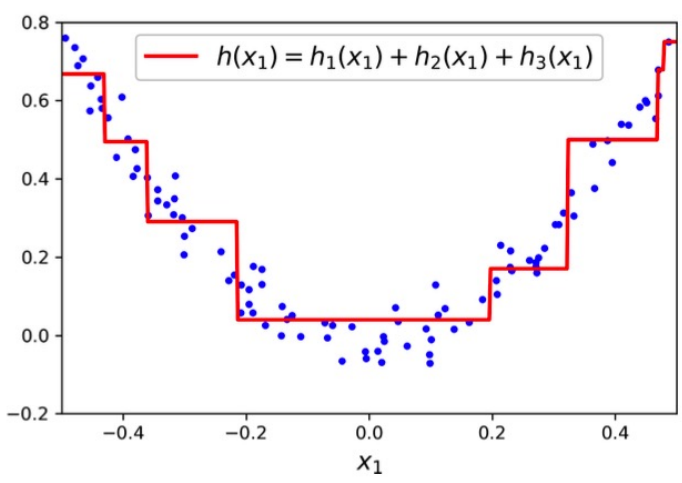
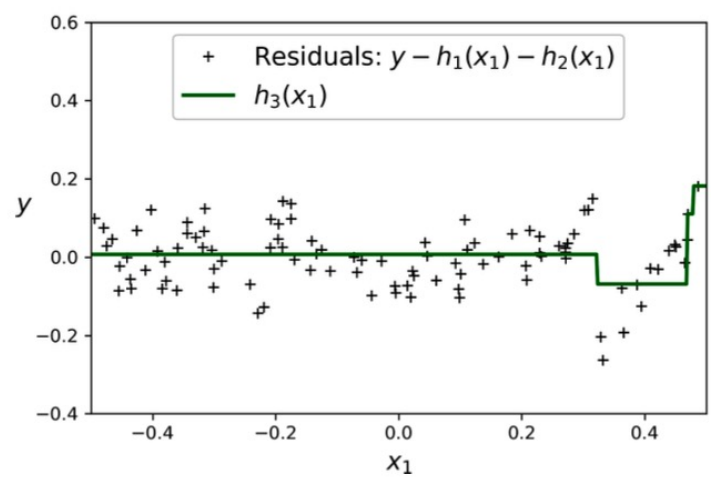
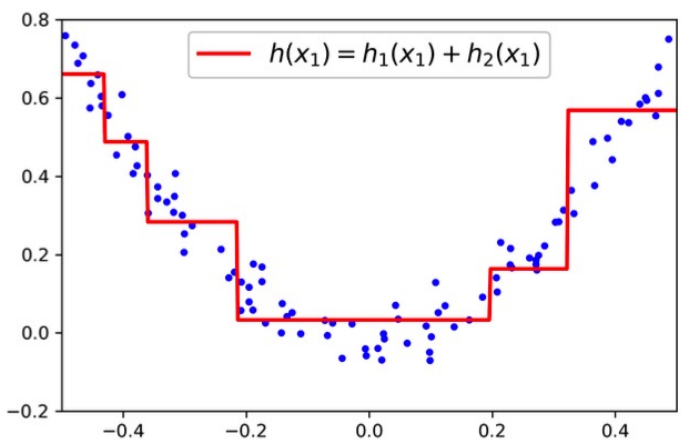
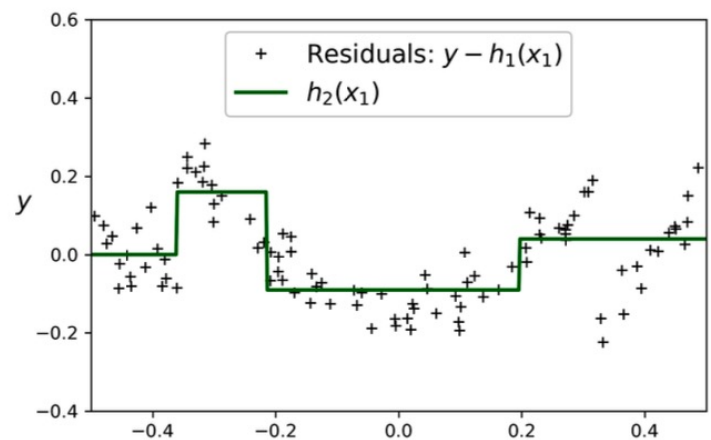
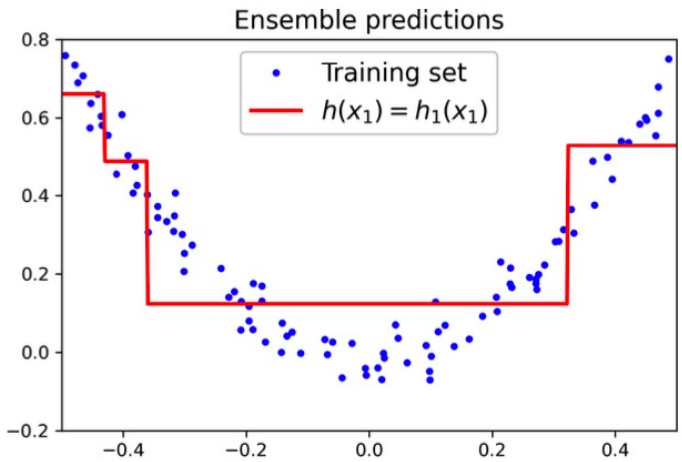
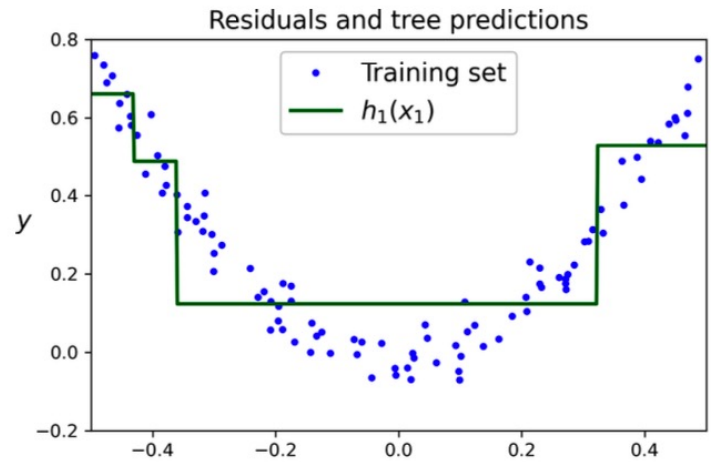
- Finish AdaBoost
- Entropy and weighted entropy
- Gradient Boosting



Gradient
Boosting
example
(regression)



Gradient
Boosting
example
(regression)



Gradient
Boosting
example
(regression)

Gradient Boosting

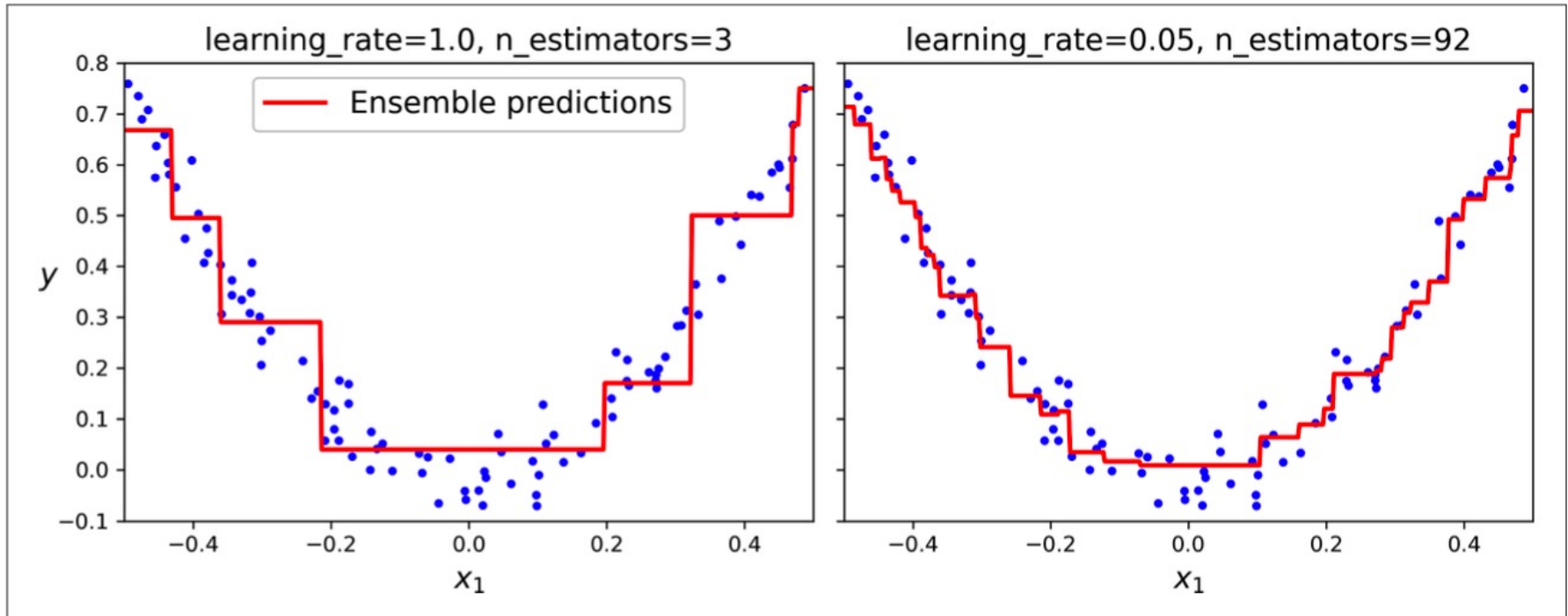


Figure 7-10. GBRT ensembles with not enough predictors (left) and just enough (right)