



CS-F415: Data Mining

10

Boosting Bagging



Dr. Kamlesh Tiwari
Associate Professor, Department of CSIS,
BITS Pilani, Pilani Campus, Rajasthan-333031 INDIA
Feb 23, 2024 **MW/F 4:00pm** 6101 @ BITS-Pilani [Jan-May 2024]

<http://ktiware.in/dm>

Boosting

- Can an **ensemble** of weak learners get a strong one

$$H(x) = \text{sign}(h^1(x) + h^2(x) + h^3(x) + \dots)$$

- What is error rate?

$$\epsilon = \frac{1}{m} \sum_{i=1}^m \mathbb{1}_{H(x^{(i)}) \neq y^{(i)}} = \sum_{\text{wrong}} \frac{1}{m}$$

Three weak learners, having error on **different data points** is prefect

- But, most likely it is not going to be the case
- We can do the following
 - Get weak learner h^1 (that is best)
 - Exaggerate the data where h^1 errs and get weak learner h^2
 - Exaggerate again the data for $h^1 \neq h^2$ and get weak learner h^3
- Will this work?

Data Mining (CS-F415)

MW/F (4:00pm) 6101@BITS-Pilani

Lecture-10 (Feb 23, 2024)

3/24

Boosting

- Let **each data point** have a **weight** w_i associated with them. Initially $w_i = \frac{1}{m}$ so that error rate becomes

$$\epsilon = \sum_{\text{wrong}} \sum_{i=1}^m \frac{1}{m} = \sum_{i \in \text{wrong}} w_i$$

- Sum of weights is 1

$$\sum w_i = 1$$
- Weights are modified in every round

Data Mining (CS-F415)

MW/F (4:00pm) 6101@BITS-Pilani

Lecture-10 (Feb 23, 2024)

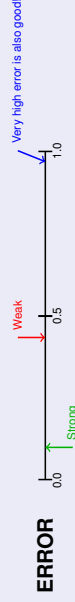
5/24

Classification

- Binary classification approaches
 - Linear classifiers
 - Quadratic classifiers
 - Bayesian classification
 - Support vector machines
 - Decision trees
 - Neural networks
 - k-nearest neighbor

Which one to use?

Weak and strong classifier?



Data Mining (CS-F415)

MW/F (4:00pm) 6101@BITS-Pilani

Lecture-10 (Feb 23, 2024)

2/24

Boosting

- Boosting uses **ensemble** of learners (weak ones)
- AdaBoost¹ is one of the widely used boosting algorithm
- Given with m labeled training examples $(x^{(1)}, y^{(1)}), \dots, (x^{(m)}, y^{(m)})$ where the $x^{(i)} \in \mathcal{X}$, and the labels $y^{(i)} \in \{-1, +1\}$
- Round t** computes a distribution D_t over training examples. Weak learning algorithm is applied to find a suitable (low weighted error ϵ_t relative to D_t) weak hypothesis $h_t : \mathcal{X} \rightarrow \{-1, +1\}$
- Final or combined hypothesis H computes the sign of a weighted combination of weak hypotheses

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

¹ A decision-theoretic generalization of on-line learning and an application to boosting. Freund, Yoav, and Schapire, Robert E. European conference on computational learning theory, pp.23-37. Springer(1995).

Data Mining (CS-F415)

MW/F (4:00pm) 6101@BITS-Pilani

Lecture-10 (Feb 23, 2024)

4/24

AdaBoost Algorithm

Algorithm 1: AdaBoost

- Initialize** w_1, \dots, w_m to $\frac{1}{m}$
- for** round $t = 1$ **to** T **do**
- Choose hypothesis** h^t minimizing error ϵ^t on **current distribution**
- Compute $\alpha_t = \frac{1}{2} \ln\left(\frac{1-\epsilon^t}{\epsilon^t}\right)$
- Update** all w_1, \dots, w_m values
- return** h and α

- Update of weights w_i at time t is done using

$$w_i^{t+1} = \frac{w_i^t}{z_t} e^{-\alpha_t h^t(x^{(i)}) y^{(i)}}$$

for all data point i . e. value of i varies from 1 to m

- Here z_t is a normalization factor

Data Mining (CS-F415)

MW/F (4:00pm) 6101@BITS-Pilani

Lecture-10 (Feb 23, 2024)

6/24

Normalization factor

- Here z_t is a normalization factor so

$$\begin{aligned}
 z_t &= \sum_{i=1}^m w_i^t e^{-(\alpha_t f(x^{(i)})) y^{(i)}} \\
 &= \sum_{i \in \text{Right}} w_i^t e^{-(\alpha_t f(x^{(i)})) y^{(i)}} + \sum_{i \in \text{Wrong}} w_i^t e^{-(\alpha_t f(x^{(i)})) y^{(i)}} \\
 &= \sum_{i \in \text{Right}} w_i^t e^{-\alpha_t} + \sum_{i \in \text{Wrong}} w_i^t e^{-\alpha_t(-1)} \\
 &= \sum_{i \in \text{Right}} w_i^t \sqrt{\frac{\epsilon^t}{1-\epsilon^t}} + \sum_{i \in \text{Wrong}} w_i^t \sqrt{\frac{1-\epsilon^t}{\epsilon^t}} \\
 &= \sqrt{\frac{\epsilon^t}{1-\epsilon^t}} \sum_{i \in \text{Right}} w_i^t + \sqrt{\frac{1-\epsilon^t}{\epsilon^t}} \sum_{i \in \text{Wrong}} w_i^t \\
 &= \sqrt{\frac{\epsilon^t}{1-\epsilon^t}} (1-\epsilon^t) + \sqrt{\frac{1-\epsilon^t}{\epsilon^t}} \epsilon^t \\
 &= 2\sqrt{\epsilon^t(1-\epsilon^t)}
 \end{aligned}$$

AdaBoost at work

Consider the following data set

(x y)	w_1	w_2	w_3	w_4	w_5	w_6	w_7	w_8	w_9	w_{10}
s_1 (1, -1)	0.1	0.1	0.1	0.1	0.1	0.1	0.1	0.1	0.1	0.1
s_2 (2, -1)	0.1	0.1	0.1	0.1	0.1	0.1	0.1	0.1	0.1	0.1
s_3 (3, +1)	0.1	0.1	0.1	0.1	0.1	0.1	0.1	0.1	0.1	0.1
s_4 (4, +1)	0.1	0.1	0.1	0.1	0.1	0.1	0.1	0.1	0.1	0.1
s_5 (5, +1)	0.1	0.1	0.1	0.1	0.1	0.1	0.1	0.1	0.1	0.1
s_6 (6, +1)	0.1	0.1	0.1	0.1	0.1	0.1	0.1	0.1	0.1	0.1
s_7 (7, +1)	0.1	0.1	0.1	0.1	0.1	0.1	0.1	0.1	0.1	0.1
s_8 (8, -1)	0.1	0.1	0.1	0.1	0.1	0.1	0.1	0.1	0.1	0.1
s_9 (9, -1)	0.1	0.1	0.1	0.1	0.1	0.1	0.1	0.1	0.1	0.1
s_{10} (10, -1)	0.1	0.1	0.1	0.1	0.1	0.1	0.1	0.1	0.1	0.1

- Round-00:** initialize weights

- Round-01:** Let sampling according to w produces (3, 10, 8, 2, 7, 5, 1, 3, 8, 3) so following sub-set of data is considered

s_1 (1, -1)
s_2 (2, -1)
s_3 (3, +1)
s_5 (5, +1)
s_7 (7, +1)
s_8 (8, -1)
s_{10} (10, -1)

Consider various hypothesis

AdaBoost at work (Round-01)

Threshold is 2.5

(x y)	h_1
s_1 (1, -1)	-1
s_2 (2, -1)	-1
s_3 (3, +1)	+1
s_4 (4, +1)	+1
s_5 (5, +1)	+1
s_6 (6, +1)	+1
s_7 (7, +1)	+1
s_8 (8, -1)	-1
s_9 (9, -1)	-1
s_{10} (10, -1)	-1

- Error rate $\epsilon = \sum_{i \in \text{Wrong}} w_i = w_8 + w_9 + w_{10} = 0.1 + 0.1 + 0.1 = 0.3$
- $\alpha = \frac{1}{2} \ln \left(\frac{1-\epsilon}{\epsilon} \right) = \frac{1}{2} \ln \left(\frac{1-0.3}{0.3} \right) = 0.4236$
- Weights are modified according to

$$w_i = w_i \times \begin{cases} \frac{1}{2(1-\epsilon)} & \text{correct} \\ \frac{1}{2\epsilon} & \text{wrong} \end{cases} = 0.7142 \quad \text{correct} \\ = 1.6666 \quad \text{wrong}$$

w_1	w_2	w_3	w_4	w_5	w_6	w_7	w_8	w_9	w_{10}
0.07	0.07	0.07	0.07	0.07	0.07	0.07	0.17	0.17	0.17

Boosting

- Weight update for correctly classified data point
- Weight update is therefore for wrongly classified data point

$$w_i^{t+1} = w_i^t \frac{1}{2(1-\epsilon^t)}$$

$$w_i^{t+1} = w_i^t \frac{1}{2\epsilon^t}$$

Final hypothesis

$$H(x) = \text{sign}(\alpha_1 h^1(x) + \alpha_2 h^2(x) + \alpha_3 h^3(x) + \dots)$$

$$\epsilon^t = \sum_{i \in \text{Wrong}} w_i^t$$

AdaBoost at work (Round-01)

	h_a	h_b	h_c	h_d	h_e	h_r	h_g	h_h
s_1 (1, -1)	+1	-1	-1	-1	-1	-1	-1	-1
s_2 (2, -1)	+1	+1	-1	-1	-1	-1	-1	-1
s_3 (3, +1)	+1	+1	+1	+1	-1	-1	-1	-1
s_5 (5, +1)	+1	+1	+1	+1	-1	-1	-1	-1
s_7 (7, +1)	+1	+1	+1	+1	+1	+1	-1	-1
s_8 (8, -1)	+1	+1	+1	+1	+1	+1	-1	-1
s_{10} (10, -1)	+1	+1	+1	+1	+1	+1	+1	+1

- Select h_c as h_1
- What is decision threshold? 2.5
- Compute error on whole dataset

AdaBoost at work (Round-02)

Consider the data set

(x y)	s_1	s_2	s_3	s_4	s_5	s_6	s_7	s_8	s_9	s_{10}
(1, -1)	(1, -1)									
(2, -1)	(2, -1)									
(3, +1)	(3, +1)									
(4, +1)	(4, +1)									
(5, +1)	(5, +1)									
(6, +1)	(6, +1)									
(7, +1)	(7, +1)									
(8, -1)	(8, -1)									
(9, -1)	(9, -1)									
(10, -1)	(10, -1)									

- Round-02:** Let sampling according to new w produces (8, 10, 7, 10, 3, 1, 10, 8, 4, 9) so following sub-set of data is considered

s_1	(1, -1)
s_3	(3, +1)
s_4	(4, +1)
s_7	(7, +1)
s_8	(8, -1)
s_9	(9, -1)
s_{10}	(10, -1)

Consider various hypothesis

AdaBoost at work (Round-02)

	h_a	h_b	h_c	h_d	h_e	h_f	h_g	h_h
S_1	(1, -1)	+1	-1	-1	-1	-1	-1	-1
S_3	(3, +1)	+1	+1	-1	-1	-1	-1	-1
S_4	(4, +1)	+1	+1	+1	-1	-1	-1	-1
S_7	(7, +1)	+1	+1	+1	+1	-1	-1	-1
S_8	(8, -1)	+1	+1	+1	+1	-1	-1	-1
S_9	(9, -1)	+1	+1	+1	+1	+1	+1	-1
S_{10}	(10, -1)	+1	+1	+1	+1	+1	+1	-1

- Select h_h as h_2
- What is decision threshold? 10.5
- Compute error on whole dataset

AdaBoost at work (Round-03)

Consider the data set

	(x y)
S_1	(1, -1)
S_2	(2, -1)
S_3	(3, +1)
S_4	(4, +1)
S_5	(5, +1)
S_6	(6, +1)
S_7	(7, +1)
S_8	(8, -1)
S_9	(9, -1)
S_{10}	(10, -1)

- **Round-03:** Let sampling according to new w produces (8, 7, 4, 8, 6, 9, 5, 8, 3, 4) so following sub-set of data is considered

	(x y)
S_3	(3, +1)
S_4	(4, +1)
S_5	(5, +1)
S_6	(6, +1)
S_7	(7, +1)
S_8	(8, -1)
S_9	(9, -1)

AdaBoost at work (Round-03)

Threshold is 0.5

	(x y)	h_3
S_1	(1, -1)	+1
S_2	(2, -1)	+1
S_3	(3, +1)	+1
S_4	(4, +1)	+1
S_5	(5, +1)	+1
S_6	(6, +1)	+1
S_7	(7, +1)	+1
S_8	(8, -1)	+1
S_9	(9, -1)	+1
S_{10}	(10, -1)	+1

- Error rate $\epsilon = \sum_{j \in \text{wrong}} w_j = 2 \times 0.037 + 3 \times 0.091 = 0.34$
- $\alpha = \frac{1}{2} \ln \left(\frac{1-\epsilon}{0.34} \right) = \frac{1}{2} \ln \left(\frac{1-0.34}{0.34} \right) = 0.3316$
- We may continue for next round like that
- Let use see our accuracy now

AdaBoost at work (Round-02)

Threshold is 10.5

(x y)	h_2
S_1	(1, -1) -1
S_2	(2, -1) -1
S_3	(3, +1) -1
S_4	(4, +1) -1
S_5	(5, +1) -1
S_6	(6, +1) -1
S_7	(7, +1) -1
S_8	(8, -1) -1
S_9	(9, -1) -1
S_{10}	(10, -1) -1

- Error rate $\epsilon = \sum_{j \in \text{wrong}} w_j = w_3 + w_4 + w_5 + w_6 + w_7 = 5 \times 0.07 = 0.35$
- $\alpha = \frac{1}{2} \ln \left(\frac{1-\epsilon}{0.35} \right) = \frac{1}{2} \ln \left(\frac{1-0.35}{0.35} \right) = 0.3095$
- Weights are modified according to

$$w_j = w_j \times \begin{cases} \frac{1}{2(1-\epsilon)} = 0.7692 & \text{correct} \\ \frac{1}{2\epsilon} = 1.4285 & \text{wrong} \end{cases}$$

	w_1	w_2	w_3	w_4	w_5	w_6	w_7	w_8	w_9	w_{10}
	0.1	0.1	0.1	0.1	0.1	0.1	0.1	0.1	0.1	0.1
	0.07	0.07	0.07	0.07	0.07	0.07	0.07	0.17	0.17	0.17
	0.037	0.037	0.177	0.177	0.177	0.177	0.177	0.091	0.091	0.091

AdaBoost at work (Round-03)

	h_a	h_b	h_c	h_d	h_e	h_f	h_g	h_h
S_3	(3, +1)	+1	-1	-1	-1	-1	-1	-1
S_4	(4, +1)	+1	+1	-1	-1	-1	-1	-1
S_5	(5, +1)	+1	+1	+1	-1	-1	-1	-1
S_6	(6, +1)	+1	+1	+1	+1	-1	-1	-1
S_7	(7, +1)	+1	+1	+1	+1	+1	-1	-1
S_8	(8, -1)	+1	+1	+1	+1	+1	+1	-1
S_9	(9, -1)	+1	+1	+1	+1	+1	+1	-1

- Select h_a as h_3
- What is decision threshold? 0.5
- Compute error on whole dataset

AdaBoost at work (Round-03)

$(\alpha_1, \alpha_2, \alpha_3) = (0.4236, 0.3095, 0.3316)$

(x y)	h_1	h_2	h_3	$\text{sign}(\sum_j \alpha_j (x_j \times h_j))$
S_1	(1, -1)	-1	-1	+1
S_2	(2, -1)	-1	-1	+1
S_3	(3, +1)	+1	-1	+1
S_4	(4, +1)	+1	-1	+1
S_5	(5, +1)	+1	-1	+1
S_6	(6, +1)	+1	-1	+1
S_7	(7, +1)	+1	-1	+1
S_8	(8, -1)	+1	+1	+1
S_9	(9, -1)	+1	+1	+1
S_{10}	(10, -1)	+1	-1	+1

