

Assignment 1

Date _____
Page No. 1

Bagging and boosting

→ Bagging and Boosting are ensemble learning techniques used in machine learning. Ensemble learning means combining multiple models to improve overall performance. Both methods help in increasing accuracy and reducing errors, but they work in different ways.

1) Bagging (Bootstrap Aggregating)

Bagging is an ensemble technique in which multiple models are trained independently using different random samples of the dataset, and their results are combined.

Working of Bagging

- i) → From the original dataset, multiple random samples are created using bootstrap sampling (sampling with replacement).
- ii) → A separate model is trained on each sample.
- iii) → All models work in parallel.
- iv) → The final output is obtained by:
 - voting (for classification)
 - averaging (for regression)

Advantages of Bagging

- Reduces overfitting
- Improves model stability
- Works well with high-variance models like decision trees

Example: Random forest is a popular algorithm based on bagging.

Disadvantages of Bagging

- Requires more computation due to multiple models.

2) Boosting

Boosting is an ensemble technique where models are trained sequentially, and each new model focuses on correcting the errors of the previous model.

Working of Boosting

- Initially all data points are given equal importance.
- A model is trained and errors are identified.
- More weights are given to misclassified data.
- The next model focuses more on these difficult data points.
- Final prediction is made by weighted voting of all models.

II Advantages of Boosting

- Reduces bias
- Improves prediction accuracy
- Works well with weak learners

Examples: AdaBoost, Gradient Boosting,
• XGBoost etc

III Disadvantages of Boosting:

- Can overfit if data is noisy
- Slower than bagging due to sequential training.

Conclusion:

Both Bagging and Boosting are important ensemble techniques. Bagging is mainly used to reduce variance and avoid overfitting, while Boosting focuses on improving accuracy by learning from previous mistakes.

Choosing the right model method depends on the dataset and problem type.

2) K-means problem!

Given points are: $A_1(2, 10)$, $A_2(2, 5)$, $A_3(8, 4)$, $B_1(5, 8)$,
 $B_2(7, 5)$, $B_3(6, 4)$, $C_1(1, 2)$ and $C_2(4, 9)$

Let's suppose $k=3$ and initial centroids are
 $A(2, 10)$, $B(9, 5)$, $B(\frac{18}{3}, \frac{18}{3})$ and $C(1, 2)$

Points	X	Y	Dist. to A	Dist. to B	Dist. to C	Cluster	New cluster
A_1	2	10	0	$\sqrt{3} = 3.61$	$\sqrt{65} = 8.06$		A
A_2	2	5	5	$\sqrt{18} = 4.24$	$\sqrt{105} = 3.16$		C
A_3	8	4	$\sqrt{72} = 8.48$	5	$\sqrt{53} = 7.28$		B
B_1	5	8	$\sqrt{3} = 3.61$	0	$\sqrt{52} = 7.21$		B
B_2	7	5	$\sqrt{50} = 7.07$	$\sqrt{13} = 3.61$	$\sqrt{65} = 8.06$		B
B_3	6	4	$\sqrt{52} = 7.21$	$\sqrt{20} = 4.47$	$\sqrt{29} = 5.38$		B
C_1	1	2	$\sqrt{65} = 8.06$	$\sqrt{52} = 7.21$	0		C
C_2	4	9	$\sqrt{5} = 2.24$	$\sqrt{2} = 1.41$	$\sqrt{8} = 2.82$		B

Cluster assignment

A $\rightarrow A_1$

B $\rightarrow A_3, B_1, B_2, B_3, C_2$

C $\rightarrow A_2, C_1$

New centroids

A $\rightarrow (2, 10)$

$$B \rightarrow \left(\frac{8+7+7+6+4}{5}, \frac{4+8+5+4+9}{5} \right) \\ = (6, 6)$$

$$C \rightarrow \left(\frac{2+1}{2}, \frac{5+2}{3} \right) \\ = (1.5, 3.5)$$

Point	X	Y	Dist. to A	Dist. to B	Dist. to C	Cluster	New cluster
A ₁	2	10	0	5.61	6.52	A	A
A ₂	2	5	5	4.12	1.58	C	C
A ₃	8	4	8.68	2.93	6.67	B	B
B ₁	5	8	3.61	2.24	5.52	B	B
B ₂	7	5	7.07	1.41	15.87	B	B
B ₃	6	4	7.21	2	11.74	B	B
C ₁	1	2	8.06	6.00	1.58	C	C
C ₂	4	9	2.29	3.61	6.20	B	A

Cluster assignment

$$A \rightarrow A_1, C_2$$

$$B \rightarrow A_3, B_1, B_2, B_3$$

$$C \rightarrow A_2, C_1$$

New Centroids

$$A \rightarrow \left(\frac{2+4}{2}, \frac{10+9}{2} \right) = (3, 9.5)$$

$$B \rightarrow \left(\frac{8+5+7+6}{4}, \frac{4+8+5+4}{4} \right) = (6.5, 5.25)$$

$$C \rightarrow \left(\frac{2+1}{2}, \frac{5+2}{2} \right) = (1.5, 3.5)$$

Points	X	Y	Dist. to A	Dist. to B	Dist. to C	Cluster	New cluster
A ₁	2	10	1.50	6.53	6.52	A	A
A ₂	2	5	4.50	4.50	1.58	C	C
A ₃	8	4	7.43	1.85	6.52	B	B
B ₁	5	8	2.50	3.13	5.70	B	B
B ₂	7	5	6.53	0.56	5.70	B	B
B ₃	6	4	6.26	1.35	4.54	B	B
C ₁	1	2	7.76	6.39	1.58	C	C
C ₂	4	9	1.11	13.51	6.204	A	A

Cluster assignment

A → A₁, C₂, B₁

B → A₃, B₁, B₂, B₃

C → A₂, C₁

(Cluster assignment did not change from the previous step so, the final clusters are)

clusters	Centroid	Points
A	(3, 9.5)	A ₁ , (A ₂), (A ₃), (C ₁)
B	(6.5, 5.25)	B ₁ , B ₂ , B ₃ , (B ₄)
C	(1.5, 3.5)	(C ₂)

New centroids

$$A = \left(\frac{2+5+4}{3}, \frac{10+8+9}{3} \right) = (3.6, 9)$$

$$B = \left(\frac{8+7+6}{3}, \frac{4+5+4}{3} \right) = (7, 4.3)$$

$$C = \left(\frac{2+1}{2}, \frac{5+2}{2} \right) = (1.5, 3.5)$$

Point	X	Y	Dist. to A	Dist. to B	Dist. to C	Cluster	New cluster
A ₁	2	10	1.95	7.55	6.52	A	A
A ₂	2	5	4.34	3.65	1.58	C	C
A ₃	8	4	6.61	1.05	6.52	B	B
B ₁	5	8	1.66	4.10	5.70	A	A
B ₂	7	5	5.21	0.67	6.70	B	B
B ₃	6	4	5.52	1.05	4.63	B	B
C ₁	1	2	7.49	6.44	1.58	C	C
C ₂	4	9	0.33	5.56	6.04	A	A

^{new}
Since a cluster did change from previous cluster.
So, final clusters are

Clusters	Centroid	Points
A	(3.6, 9)	A ₁ (2, 10), B ₁ (5, 8), C ₂ (4, 9)
B	(7, 4.3)	A ₃ (8, 4), B ₂ (7, 5), B ₃ (6, 4)
C	(1.5, 3.5)	A ₂ (2, 5), C ₁ (1, 2)

3) Supervised learning models

⇒ Classification Model

II Advantages:

- High prediction accuracy when sufficient labeled data is available.
- Clear performance evaluation using metrics (accuracy, precision, recall, F1-score)
- Effective for tasks like spam detection, image classification, disease prediction etc.
- Many supervised algorithms like decision trees provide insights into feature importance, making it easier to understand model decisions and iterate on improvements.

II Disadvantages:

- Requires large amount of labeled data, which is costly and time consuming.
- Overfitting can occur (especially Decision Trees)
- Performance drops if data is noisy or imbalanced.

→ Regression models

Advantages

- Predicts continuous values effectively
- Easy to interpret (especially linear regression)
- Useful for forecasting tasks (house prices, stock prices)

Disadvantages

- Sensitive to outliers (linear regression)
- Assumes relationship between variables may not hold true
- Complex models may overfit

Unsupervised learning models

② Clustering models

Advantages

- No labeled data is required
- Useful for discovering hidden patterns
- Simple and computationally efficient
- Widely used for customer segmentation and anomaly detection

Disadvantages

- Need to predifine number of clusters (k)
- Sensitive to initial centroid selection
- Struggles with non-spherical or uneven cluster sizes

③ Dimensionality Reduction

Advantages

- Reduces complexity of high dimensional data
- Improves visualization and performance
- Removes redundant features

Disadvantages

- Possible loss of important information
- Reduced interpretability