

Introduction to Machine Learning

Lecture 12

Mehryar Mohri
Courant Institute and Google Research
mohri@cims.nyu.edu

Bagging

Ensemble Methods - Classification

- **Problem:** given T binary classification hypotheses (h_1, \dots, h_T) find combined classifier

$$f: x \mapsto \operatorname{sgn} \left(\sum_{t=1}^T \alpha_t h_t \right),$$

with better performance.

- When does it work? Need diversity (e.g., different features, different training sets).

→ use different subsets of the data for training.

Bagging - Classification

(Breiman, 1996)

Bagging (Bootstrap aggregating).

BAGGING($S = ((x_1, y_1), \dots, (x_m, y_m))$)

```
1 for  $t \leftarrow 1$  to  $T$  do
2    $S_t \leftarrow \text{BOOTSTRAP}(S)$   $\triangleright$  i.i.d. sampling with replacement from  $S$ .
3    $h_t \leftarrow \text{TRAINCLASSIFIER}(S_t)$ 
4 return  $h_S = x \mapsto \text{MAJORITYVOTE}((h_1(x), \dots, h_T(x)))$ 
```

Bagging - Regression

(Breiman, 1996)

Bagging (Bootstrap aggregating).

BAGGING($S = ((x_1, y_1), \dots, (x_m, y_m))$)

```
1 for  $t \leftarrow 1$  to  $T$  do
2    $S_t \leftarrow \text{BOOTSTRAP}(S)$   $\triangleright$  i.i.d. sampling with replacement from  $S$ .
3    $h_t \leftarrow \text{TRAINREGRESSIONALGORITHM}(S_t)$ 
4 return  $h_S = x \mapsto \text{MEAN}((h_1(x), \dots, h_T(x)))$ 
```

Bias-Variance Decomposition

- **Proposition:** for any hypothesis h_S , the following decomposition holds:

$$\begin{aligned} & \mathbb{E}_S \left[\mathbb{E}_{X,Y} [(h_S(X) - Y)^2] \right] \\ &= \underbrace{\mathbb{E}_{S,X} \left[\left(h_S(X) - \mathbb{E}_S[h_S(X)] \right)^2 \right]}_{\text{variance}} + \underbrace{\mathbb{E}_X \left[\left(\mathbb{E}_S[h_S(X)] - \mathbb{E}[Y|X] \right)^2 \right]}_{\text{bias}} + \underbrace{\mathbb{E}_X \left[(Y - \mathbb{E}[Y|X])^2 \right]}_{\text{noise}}. \end{aligned}$$

- Bias-variance minimization trade-off:
 - small S and large H : small bias, large variance.
 - large S and small H : large bias, small variance.

Bias-Variance Decomposition Proof

■ Observe that

$$\begin{aligned} & \mathbb{E}_{X,Y}[(h_S(X) - \mathbb{E}[Y|X])(\mathbb{E}[Y|X] - Y)] \\ &= \mathbb{E}_X \left[\mathbb{E}_{Y|X}[(h_S(X) - \mathbb{E}[Y|X])(\mathbb{E}[Y|X] - Y)] \right] \\ &= \mathbb{E}_X \left[(h_S(X) - \mathbb{E}[Y|X]) \mathbb{E}_{Y|X}[(\mathbb{E}[Y|X] - Y)] \right] = 0. \end{aligned}$$

■ Therefore,

$$\begin{aligned} \mathbb{E}_{X,Y}[(h_S(X) - Y)^2] &= \mathbb{E}_{X,Y} \left[[(h_S(X) - \mathbb{E}[Y|X]) + (\mathbb{E}[Y|X] - Y)]^2 \right] \\ &= \mathbb{E}_{X,Y}[(h_S(X) - \mathbb{E}[Y|X])^2] + \mathbb{E}_{X,Y}[(\mathbb{E}[Y|X] - Y)^2]. \end{aligned}$$

Bias-Variance Decomposition Proof

- Since $\mathbb{E}_S \left[h_S(X) - \mathbb{E}_S[h_S(X)] \right] = 0$, the following holds:

$$\begin{aligned} & \mathbb{E}_S \left[\mathbb{E}_{X,Y} \left[(h_S(X) - \mathbb{E}[Y|X])^2 \right] \right] \\ &= \mathbb{E}_{X,Y} \left[\mathbb{E}_S \left[\left(h_S(X) - \mathbb{E}_S[h_S(X)] + \mathbb{E}_S[h_S(X)] - \mathbb{E}[Y|X] \right)^2 \right] \right] \\ &= \mathbb{E}_{X,Y} \left[\mathbb{E}_S \left[\left(h_S(X) - \mathbb{E}_S[h_S(X)] \right)^2 \right] \right] + \mathbb{E}_X \left[\mathbb{E}_S \left[\left(\mathbb{E}_S[h_S(X)] - \mathbb{E}[Y|X] \right)^2 \right] \right]. \end{aligned}$$

Ensemble Bias

- Bias of averaged hypothesis in regression:

$$\begin{aligned}\text{Bias}(h_S, x) &= \mathbb{E}_{S \sim D^m} \left[\frac{1}{T} \sum_{t=1}^T h_t(x) \right] - \mathbb{E}[Y|x] \\ &= \frac{1}{T} \sum_{t=1}^T \left[\mathbb{E}_{S \sim D^m} [h_t(x)] - \mathbb{E}[Y|x] \right] \\ &= \frac{1}{T} \sum_{t=1}^T \text{Bias}(h_t, x).\end{aligned}$$

→ Thus, relatively unbiased base hypotheses lead to relatively unbiased ensemble.

Ensemble Variance

- **Proposition:** for any x , the variance of the ensemble hypothesis at x is given by

$$\text{Var}(h_S, x) = \frac{1}{T^2} \sum_{t=1}^T \text{Var}(h_t, x) + \frac{1}{T^2} \sum_{t \neq t'} \text{Cov}[h_t(x), h_{t'}(x)].$$

- thus, if approximately uncorrelated base hypotheses $\text{Cov}[h_t(x), h_{t'}(x)] \approx 0$.
- assume approximately equal variances, then

$$\text{Var}(h_S, x) \approx \frac{1}{T} \text{Var}(h_1, x).$$

→ reduction by $1/T$.

Bagging - Regression

- Regression properties:
 - small covariances (different subsets).
 - similar variances (on average).
 - similar biases.
- Classification: unclear explanation.

References

- Leo Breiman. Bagging Predictors. *Machine Learning*, 24:123-140, 1996.
- Peter Bühlmann and Bin Yu. Analyzing Bagging. *The Annals of Statistics*. 30(4): 927-961, 2002.
- Andreas Buja and Werner Stuetzle. “The effects of bagging on variance, bias, and mean squared error”, Preprint, AT&T Labs-Research, 2000.