

Bagging + Boosting

Bagging: Bootstrap Aggregating

Bootstrap sampling

$X_i \sim N(\mu, \sigma^2)$. Goal: Est μ

$$\hat{\mu} = \frac{1}{n} \sum_{i=1}^n X_i, \quad E[\hat{\mu}] = \mu, \quad \text{Var}[\hat{\mu}] = \text{Var}\left[\frac{1}{n} \sum X_i\right] = \frac{1}{n^2} \text{Var}[\sum X_i]$$
$$= \frac{1}{n^2} \cdot n \cdot \text{Var}[X_i]$$
$$= \frac{1}{n} \cdot \sigma^2 = \frac{\sigma^2}{n}$$

Var too large...

Say we have B points $\sim N(\mu, \sigma^2)$

$$S_1 = \{X_1, \dots, X_n\}, \quad S_2 = \{X_{n+1}, \dots, X_{2n}\}, \dots, \quad S_B = \{X_{(B-1)n+1}, \dots, X_{Bn}\}$$

$$\hat{\mu}^{(j)} = \frac{1}{n} \sum_{z \in S_j} z$$

$$\hat{\mu}^{(\text{avg})} = \frac{1}{B} \sum \hat{\mu}^{(j)}, \quad E[\hat{\mu}^{(\text{avg})}] = \mu, \quad \text{Var}\left[\frac{1}{B} \sum \hat{\mu}^{(j)}\right] = \frac{1}{B^2} \cdot B \cdot \frac{\sigma^2}{n} = \frac{\sigma^2}{nB}$$

Needs $B \times$ more data

Bootstrap sampling \approx sampling w/ replacement

Dataset X_1, X_2, X_3, X_4, X_5

$$S_1 = \{X_3, X_4, X_1, X_1, X_4\}$$

$$S_2 = \{X_5, X_5, X_3, X_1, X_2\}$$

$$\vdots$$
$$S_B = \dots$$

$$\hat{\mu}^{(j)} = \frac{1}{n} \sum_{z \in S_j} z$$

$$\hat{\mu}^{(\text{avg})} = \frac{1}{B} \sum_{j=1}^B \hat{\mu}^{(j)}$$

$$E[\hat{\mu}^{(\text{avg})}] = \mu.$$

Var = ?

Not independent

Still works!

Application to ML

Some approaches: high variance

Decision trees

Running DT on two splits of same dataset \rightarrow different trees.

Use bagging.

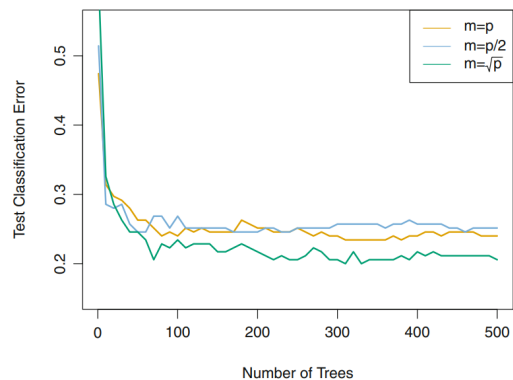
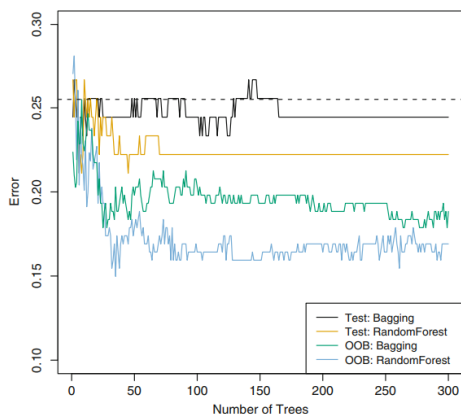
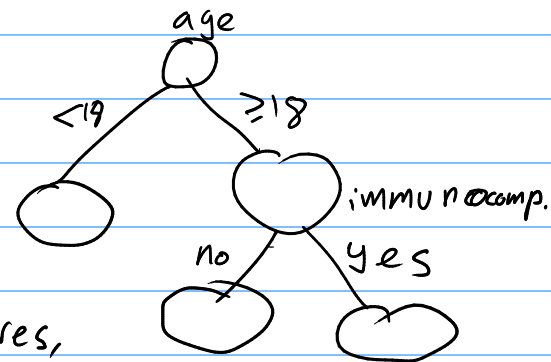
1. Bootstrap sample B datasets from training set S_1, \dots, S_B
2. Run learning algo on each $\hat{f}^{(1)}, \dots, \hat{f}^{(B)}$
3. Aggregate $\hat{f}^{(1)}, \dots, \hat{f}^{(B)}$. How? Regression: $\hat{f}(x) = \frac{1}{B} \sum \hat{f}^{(j)}(x)$
Classification: $\hat{f}(x) = \text{majority vote of } \hat{f}^{(j)}(x)$

Random Forests

Bagging - Correlations (some)
on DTs

Classic: Look at all d features, pick "best"

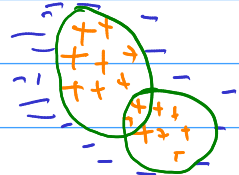
RFs: Choose random $m \approx \sqrt{d}$ features, split on "best"



Boosting

Weak Learner: ~55% accurate (better than rand guess)

Strong Learner: "~90%+ accurate"



Online Learning with Experts

Example: Horse racing. T races. n horses.

T rounds. At round t , algo picks weights $p_1^{(t)}, \dots, p_n^{(t)}$
 s.t. $\sum_{i=1}^n p_i^{(t)} = 1$. Algo experiences "loss" at time t : $\langle p^{(t)}, l^{(t)} \rangle$

Goal: $\min \sum_{t=1}^T \langle p^{(t)}, l^{(t)} \rangle$

Benchmark: $\min_{i=1, \dots, n} \sum_{t=1}^T l_i^{(t)}$ $p^{(t)} = \text{arg min}_i \sum_{t=1}^T l_i^{(t)}$ $\forall t$

↑
adversarially
picked loss
[0,1]^n

Hedge(β) $\beta \in [0, 1]$

Set $w^{(1)} = [1/n, \dots, 1/n] \in \mathbb{R}^n$

For $t=1, \dots, T$:

Set $p^{(t)} = \frac{w^{(t)}}{\sum_{i=1}^n w_i^{(t)}} \quad (\text{Normalize})$

Receive loss $\langle p^{(t)}, l^{(t)} \rangle$

Update: $w_i^{(t+1)} = w_i^{(t)} \beta^{l_i^{(t)}} \quad (\text{MWU})$ Best option in hindsight

$$\underbrace{\sum_{t=1}^T \langle p^{(t)}, l^{(t)} \rangle}_{\text{total loss}} \leq \frac{\log n + \log(1/\beta) \min_{i=1, \dots, n} \sum_{t=1}^T l_i^{(t)}}{1 - \beta}$$

Set β :

$$\underbrace{\frac{1}{T} \sum_{t=1}^T \langle p^{(t)}, l^{(t)} \rangle}_{\text{average loss}} \leq \underbrace{\frac{1}{T} \min_{i=1, \dots, n} \sum_{t=1}^T l_i^{(t)}}_{\text{avg best loss in hindsight}} + \underbrace{O\left(\sqrt{\frac{\log n}{T}}\right)}_{\text{"regret"}} \xrightarrow{\text{as } T \rightarrow \infty} 0$$

AdaBoost

Given algo WeakLearn that gets $> 5\%$ accuracy on a training set, Can we boost this to high accuracy?

Wrong way: algorithms are "horses" (experts)
"Put large weight on good algos"

Right way: Datapoints are experts
Large weight \rightarrow still must learn point.

Adaboost:

Set $w^{(0)} = [1/n, \dots, 1/n]$

For $t=1$ to T

$$\text{Set } p^{(t)} = \frac{w^{(t)}}{\sum_{i=1}^n w_i^{(t)}}$$

Run WeakLearn on training set (with weights $p^{(t)}$)

\hookrightarrow Get classifier $h^{(t)}: X \rightarrow [0, 1]$

$$\text{Calculate error: } \epsilon_t = \sum_{i=1}^n p_i^{(t)} |h^{(t)}(x_i) - y_i| \quad // \epsilon_t \leq 0.5$$

avg. weighted error on pt i

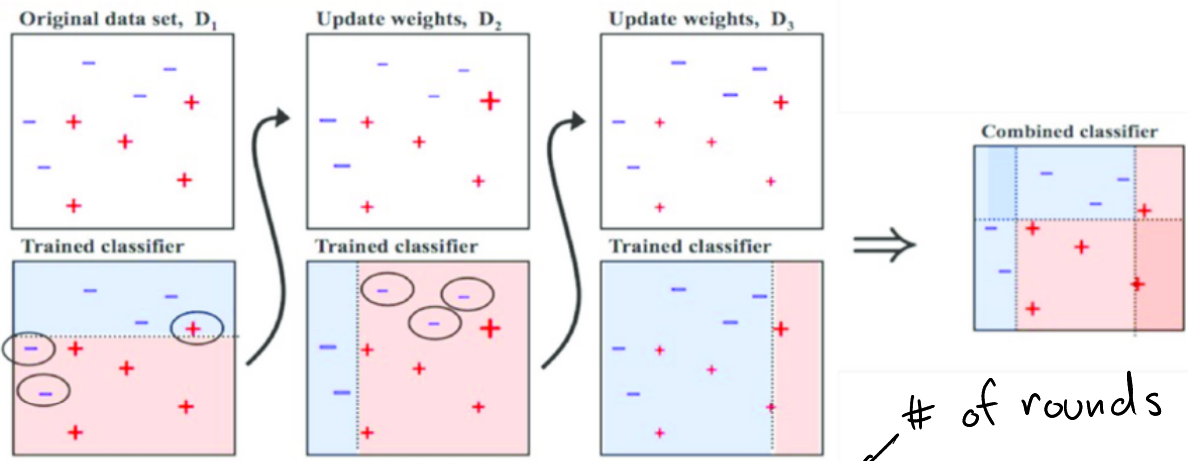
ϵ_t small $\rightarrow \beta_t$ small
 \rightarrow reduce weights a lot

$$\beta_t = \frac{\epsilon_t}{1 - \epsilon_t}$$

If $\epsilon_t \leq \frac{1}{2}$

$$\text{Set } w_i^{(t+1)} = w_i^{(t)} \beta_t^{1 - |h^{(t)}(x_i) - y_i|}$$

$$h(x) = \begin{cases} 1 & \text{if } \sum_{t=1}^T \log\left(\frac{1}{\beta_t}\right) h^{(t)}(x) \geq \frac{1}{2} \sum_{t=1}^T \log\left(\frac{1}{\beta_t}\right) \\ 0 & \text{else} \end{cases}$$

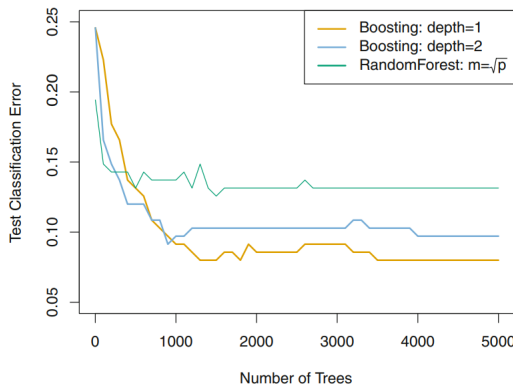


Training error: $\frac{1}{n} \sum_{i=1}^n \mathbb{1}\{h(x_i) \neq y_i\} \leq 2 \prod_{t=1}^T \sqrt{\epsilon_t(1-\epsilon_t)}$

Say $\epsilon_t \leq \frac{1}{2} - \gamma$.

$\leq \exp(-2T\gamma^2)$

Want: Good test error!
 Generalize well w simple base classifier.



Generic + Flexible

Less interpretable.

Bagging - Parallel

Boosting - Sequential