

4/15/15

1

# ENSEMBLE LEARNING

## ① Main Idea

→ Weak / Simple learners (NB, LR, Linear SVMs, decision stumps)

→ Underfit

→ High bias

→ Low Variance

→ Strong learners (Kernel SVMs, <sup>Deep</sup> Neural Nets, <sup>Deep</sup> Decision Trees)

→ Overfit

→ Low Bias

→ High Variance

Ensemble Methods: Can we get the best of both?

Bagging

(multiple)

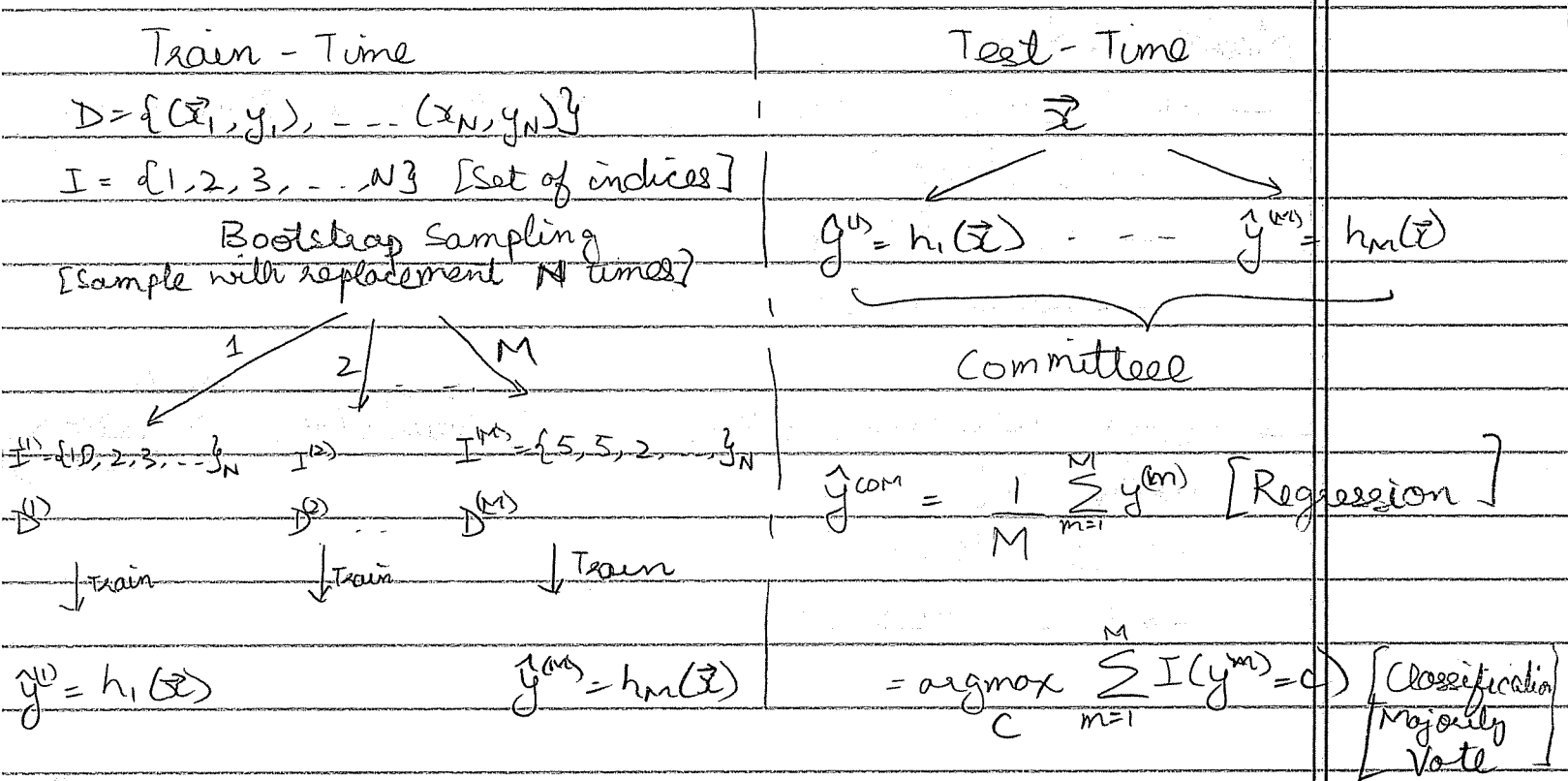
→ Reducing the variance of strong learners

Boosting

(multiple)

→ Reducing the bias of weak learners

## ② Bagging = Bootstrap Averaging



Why? Because committees make better predictions [under certain assumptions]

$$E_{P(x,y)} [L(y, \hat{y}^{COM}(x))]$$

$$= E \left[ \left( y - \frac{1}{M} \sum_{m=1}^M \hat{y}^{(m)} \right)^2 \right] \quad \left[ \text{For Regression w/ } L_2 \text{-Loss} \right]$$

$$= \frac{1}{M^2} E \left[ \left( M y - \sum_{m=1}^M \hat{y}^{(m)} \right)^2 \right]$$

Let  $e_m = y - \hat{y}^{(m)}$  Error of  $m^{\text{th}}$  predictor

(2)

$$E_{com} = \frac{1}{M^2} E \left[ \left( \sum_{m=1}^M e_m \right)^2 \right]$$

$$= \frac{1}{M^2} E \left[ \sum_{m=1}^M e_m^2 + \sum_{m_1 \neq m_2} e_{m_1} e_{m_2} \right]$$

Now, Assume: ① zero mean errors  $E[e_m] = 0 \quad \forall m$   
 ② uncorrelated errors  $E[e_{m_1} e_{m_2}] = E[e_{m_1}] E[e_{m_2}]$

$$E_{com} = \frac{1}{M^2} \left[ \sum_{m=1}^M E[e_m^2] + \sum_{m_1 \neq m_2} \underbrace{E[e_{m_1}]}_0 \underbrace{E[e_{m_2}]}_0 \right]$$

$$= \frac{1}{M^2} \sum_{m=1}^M E[e_m^2] = \frac{1}{M^2} \sum_{m=1}^M E[(y - \hat{y}^{(m)})^2]$$

$$\Rightarrow E_{com} = \frac{1}{M} E_{av}$$

$\Rightarrow$  Large (uncorrelated-error) committees cut error!  $\frac{1}{M}$ th of average error!

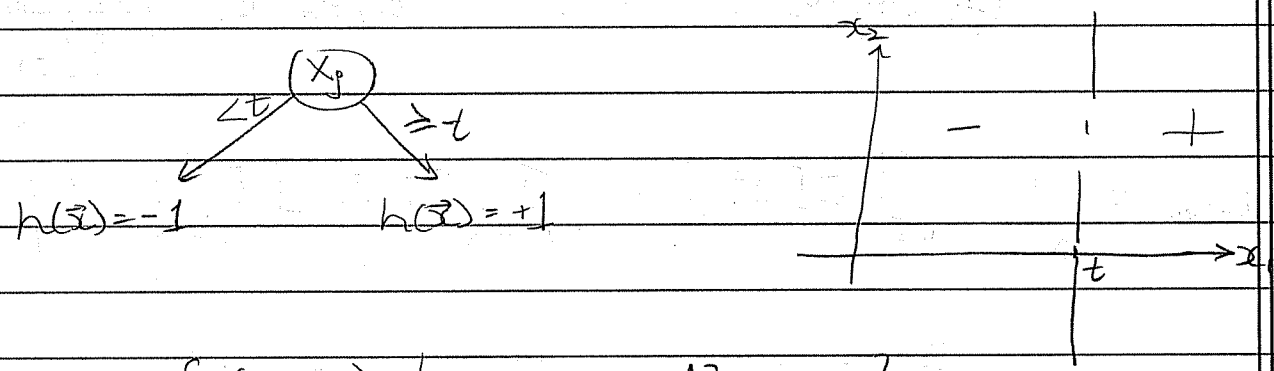
Assumption ① is okay. Why?

Assumption ② is difficult to satisfy in practice. True for bootstrap?

③ Boosting: Reducing Bias of Weak learners  
 i.e. making weak learners strong.

Given: Hypothesis Class  $\mathcal{H} = \{h \mid h: X \rightarrow Y\}$   
 collection of weak learners

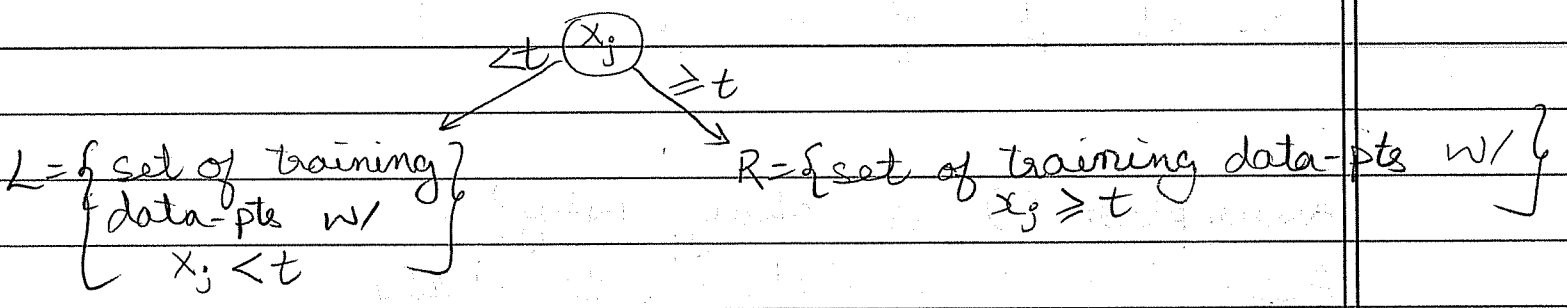
Example #1: Decision Stump [for binary classification]



$$\mathcal{H} = \{ (j, t) \mid j \in \{1, 2, \dots, d\}, t \in \mathbb{R} \}$$

↑ Pick a feature & Pick a threshold  
 And you have a decision-stump classifier

Example #2: Decision Stump [for Regression]



$$h(\vec{x}) = \begin{cases} \frac{1}{|L|} \sum_{i \in L} y_i & x_j < t \\ \frac{1}{|R|} \sum_{i \in R} y_i & x_j \geq t \end{cases}$$

[Bucket Average]  
 0<sup>th</sup>-order polynomial fit

$$h(\vec{x}) = \begin{cases} \vec{w}_L^T \vec{x} & x_j \leq t \\ \vec{w}_R^T \vec{x} & x_j > t \end{cases}$$

Assume: Access to a Black-Box Learning algorithm to pick a weak learner for a dataset, i.e.:

$$h^* = \operatorname{argmin}_{h \in H} \frac{1}{N} \sum_{i=1}^N L(y_i, h(\vec{x}_i))$$

can be solved

Goal of Boosting: Learn a committee of weak-learners with low-loss

$$f(\vec{x}) = \sum \alpha_t h_t(\vec{x})$$

$$\min_{\alpha_t, h_t} \frac{1}{N} \sum_{i=1}^N L(y_i, f(\vec{x}_i))$$

Problem: Joint optimization of all weak learners is difficult

Boosting involves greedy optimization.

At "time"  $t$  fixed

learned / Optimized

$$\alpha_1 h_1(x) + \alpha_2 h_2(x) + \dots + \alpha_{t-1} h_{t-1}(x) + \alpha_t h_t(x)$$

$$f_t(x) = f_{t-1}(x) + \alpha_t h_t(x)$$

So

$$\alpha_t^*, h_t^* = \min_{\substack{\alpha_t \in \mathbb{R} \\ h_t \in \mathcal{H}}} \frac{1}{N} \sum_{i=1}^N L(y_i, f_t(x_i))$$

$L_2$ -BOOST: Assume Regression &  $L_2$ -error

$$(\alpha_t^*, h_t^*) = \min_{\substack{\alpha_t \in \mathbb{R} \\ h_t \in \mathcal{H}}} \frac{1}{N} \sum_{i=1}^N [y_i - f_t(\tilde{x}_i)]^2$$

$$= \frac{1}{N} \sum_{i=1}^N \underbrace{[y_i - f_{t-1}(\tilde{x}_i) - \alpha_t h_t(\tilde{x}_i)]^2}_{\text{Current-Error}}$$

$$= \frac{1}{N} \sum_{i=1}^N [e_i - \alpha_t h_t(\tilde{x}_i)]^2$$

Usually  $\alpha_t$  scalar can be learned inside  $h_t$  so we set  $\alpha_t^* = 1$

$$= \min_{h_t \in \mathcal{H}} \frac{1}{N} \sum_{i=1}^N [e_i - h_t(x_i)]^2$$

Black-Box Call: Regress to Errors!

Train weak learner  $h(\tilde{x})$  with  $(\tilde{x}_i, y_i = e_i)$