

random forest models

multiple CART models fitted on bootstrapped data sets drawn from the sample

example of "bagging" — bootstrap aggregating

Bishop: "Committees" of model

committee's prediction  $y_{com}(x)$

$$= \frac{1}{M} \sum_{m=1}^M y_m(x)$$

where  $m$  indexes each model

$M = \#$  of models

$h(x)$ : DGP       $y_m(x) = h(x) + \epsilon_m(x)$

$$E_x[SSE_m] = E_x[(y_m(x) - h(x))^2] = E_x[\epsilon_m(x)^2]$$

the expectation here is over  $x$ .

For the entire committee,

$$E[SSE_{com}] = E\left[\left(\frac{1}{M} \sum_{m=1}^M y_m(x) - h(x)\right)^2\right] = E\left[\left(\frac{1}{M} \sum_{m=1}^M \epsilon_m(x)\right)^2\right]$$

start omitting dependency of  $\epsilon$  on  $x$ .

$$E \left[ \frac{1}{M^2} \cdot (\varepsilon_1 + \varepsilon_2 + \dots + \varepsilon_M) (\varepsilon_1 + \varepsilon_2 + \dots + \varepsilon_M) \right]$$

If we assume  $\mu_{\varepsilon_m} = 0$  and  $\varepsilon_m, \varepsilon_k$  are statistically independent for all  $m, k \in 1 \dots M$  then

$$\frac{1}{M^2} E \left[ \left( \sum_{m=1}^M \varepsilon_m - \mu_{\sum \varepsilon} \right) \cdot \left( \sum_{m=1}^M \varepsilon_m - \mu_{\sum \varepsilon} \right) \right]$$

$$\mu_{\sum \varepsilon} = \frac{1}{M} \cdot E[\varepsilon_1 + \varepsilon_2 + \dots + \varepsilon_M]$$

by independence,

$$\frac{1}{M} \cdot E[\varepsilon_1] + E[\varepsilon_2] + \dots + E[\varepsilon_M]$$

by assumption,  $E[\varepsilon_m] = 0 \quad \forall m \in 1 \dots M$

ergo,

$$\mu_{\sum \varepsilon} = 0$$

$$\frac{1}{M^2} \cdot \text{var} \left[ \sum_{m=1}^M \varepsilon_m \right] = \frac{1}{M^2} \cdot \sum_{m=1}^M \text{var}(\varepsilon_m)$$

b/c  $\varepsilon$ s are statistically independent.

so, because  $\mu_m = 0 \quad \forall \varepsilon_m, m \in 1 \dots M,$

$$\text{var}(\varepsilon_m) = E[\varepsilon_m^2]$$

← putting dependency on  $x$  back in.

$$E[\text{SSE}_{\text{com}}] = \frac{1}{M^2} \sum_{m=1}^M E[\varepsilon_m(x)^2]$$

↑ ?

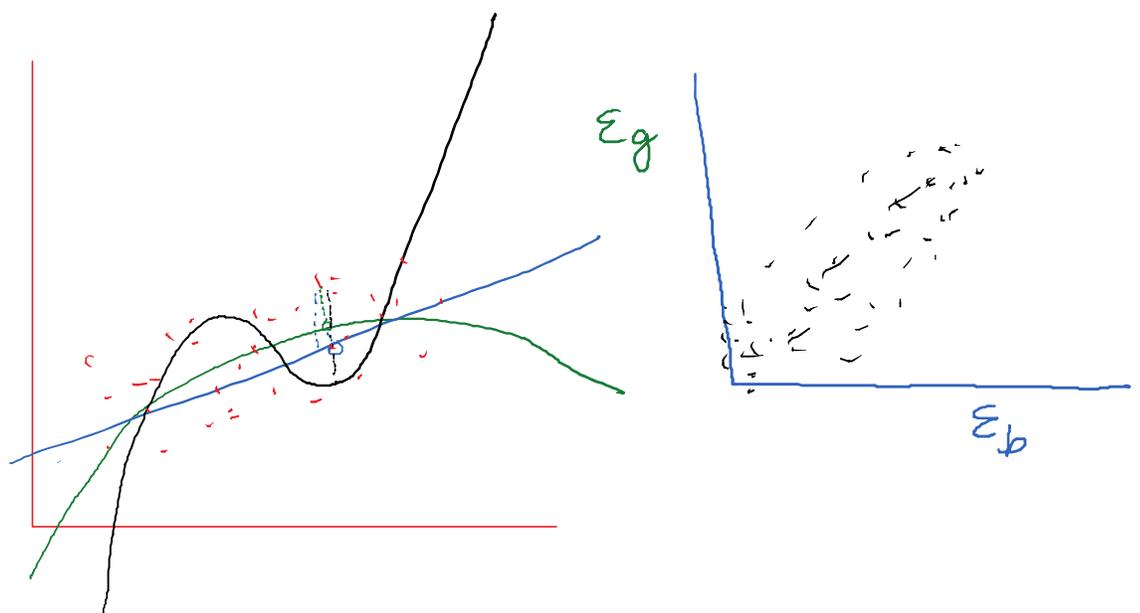
$$E[SSE_m] = E_x[\varepsilon_m(x)^2] \quad \swarrow ?$$

well,  $E[SSE_{com}] = \frac{1}{M^2} \cdot M E[SSE_m]$

$$E[SSE_{com}] = \frac{1}{M} E[SSE_m]$$

the expected SSE of an ensemble estimator is equal to  $\frac{1}{M}$  times the expected SSE of its  $M$  many constituents...

... as long as  $\varepsilon_m$  are statistically independent.



it's hard to create models of the same data  
with independent  $\epsilon_m(x)$ .

under perfect error correlation,

$$\begin{aligned} E[SSE_{\text{com}}] &= E\left[\frac{1}{M^2} \cdot (\epsilon_1 + \epsilon_2 + \dots + \epsilon_M)(\epsilon_1 + \epsilon_2 + \dots + \epsilon_M)\right] \\ &= E\left[\frac{1}{M^2} \cdot M \cdot \epsilon_M \cdot M \cdot \epsilon_M\right] \\ &= E[\epsilon_M^2] = E[SSE_M]. \end{aligned}$$

... the same as any constituent model. Ergo there  
is no benefit to constructing an ensemble.

lower bound  $E[SSE_{\text{com}}] = \frac{1}{M} E[SSE_M]$

upper bound  $E[SSE_{\text{com}}] = E[SSE_M]$

... in other words, ensemble models can't do worse than  
their (average) constituent, but can do much better.

# Boosting

Sunday, April 14, 2013 10:24 PM

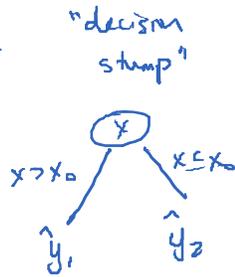
fitting a series of models sequentially to a data set, where each model improves upon the last by focusing upon the observations for which the previous model was most deficient.

Adaboost ("adaptive boosted") is one implementation of boosting.

Weights misclassified (poorly predicted) cases more heavily when fitting subsequent models.

## Adaboost

- ① set  $w_n^{(1)} = \frac{1}{N}$  for all observations  $n = 1, \dots, N$ .  
This is the observation weighting.



- ② start with your first classifier.

Fit this classifier to minimize:  
classifier vs prediction for  $y_n$  given  $x_n$  (IVs for obs  $n$ )

$$\sum_{n=1}^N w_n^{(m)} \cdot \mathbb{I} \left( y_m(x_n) \neq t_n \right)$$

target  $t_n$  for obs  $n$

indicator weight for classifier  $m$ .

indicator: 1 if true  
0 if false

in words: minimize weighted classification error.

③ Calculate

$$\epsilon_m = \frac{\sum_{n=1}^N w_n^{(m)} \cdot I(y_m(x_n) \neq t_n)}{\sum_{n=1}^N w_n^{(m)}}$$

if for  $m$ th we misclassified EVERY obs,

$$\epsilon_m = \frac{\sum_{n=1}^N \frac{1}{N} \cdot 1}{\sum_{n=1}^N \frac{1}{N}} = \frac{\frac{1}{N} \cdot N}{\frac{1}{N} \cdot N} = 1$$

and if we correctly classified EVERY case

$$\epsilon_m = \frac{\sum_{n=1}^N \frac{1}{N} \cdot 0}{\sum_{n=1}^N \frac{1}{N}} = 0.$$

④ update  $w_n^{(m+1)} = w_n^{(m)} \cdot \exp(\alpha_m \cdot I(y_m(x_n) \neq t_n))$

$$\alpha_m = \ln \left[ \frac{1 - \epsilon_m}{\epsilon_m} \right]$$

all for  
 $m=1$

under "expected awful" classifier (coin flip)

$$\epsilon_m \approx \frac{\frac{1}{N} \cdot 1 \cdot \frac{1}{2} N}{\frac{1}{N} \cdot N} \approx \frac{1}{2}$$

$$\alpha_m = \ln \left( \frac{1 - \frac{1}{2}}{\frac{1}{2}} \right) = \ln(1) = 0$$

→ no weight updating.

if we classify  $\frac{3}{4}$  correct:

$$\epsilon_m \approx \frac{\frac{1}{N} \cdot 1 \cdot \frac{1}{4} N}{\frac{1}{N} \cdot N} = \frac{3}{4}$$

$$\alpha_m = \ln \left( \frac{1 - \frac{1}{4}}{\frac{1}{4}} \right) = \ln \left( \frac{\frac{3}{4}}{\frac{1}{4}} \right)$$

$$= \ln \left( \frac{3}{1} \right) \approx 0.90$$

$$\rightarrow w_n^{m+1} = \frac{3}{2} \cdot w_n^m$$

IF  $n$  was  
misclassified.

$$w_n^{m+1} = w_n^m$$

IF  $n$  was  
correctly classified.

⑤ Repeat 2-4 for  $m=1 \dots M$ .

⑥ Calculate final predictions:

$$y_m(x) = \text{sign} \left[ \sum_{m=1}^M \alpha_m \cdot y_m(x) \right]$$

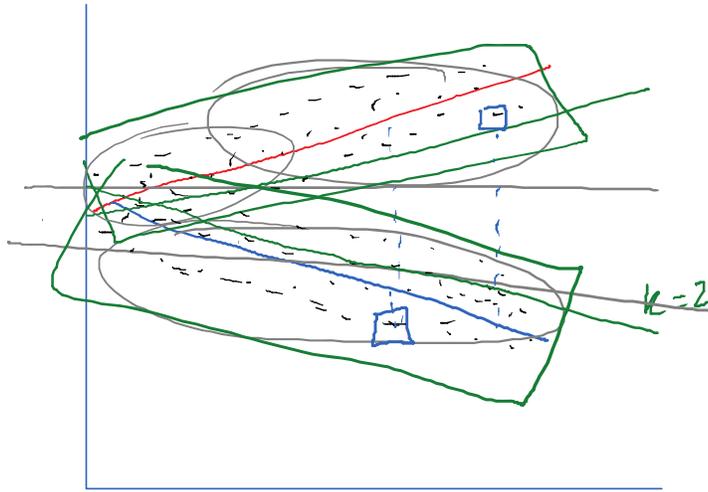
$$y_m \in \{-1, 1\}$$

Adaboost: designed for discrete DVs.

For continuous DVs, "gradient boosting" is a similar algorithm; each classifier is fit to residual errors  $t_n - f_{(m-1)}(x)$

# Conditional Mixture Models

Sunday, April 14, 2013 10:30 PM



$K = \#$  of distinct models

Conditional regression mixture model

$$\ln L = \sum_{n=1}^N \sum_{k=1}^K \underbrace{z_{nk}}_{\text{latent parameter}} \cdot \ln \left( \underbrace{\pi_k}_{\text{probability}} \underbrace{\Phi(t_n | \lambda \beta_k, \sigma_k^2)}_{\text{likelihood}} \right)$$

latent parameter  $\{0, 1\}$  indicating under obs.  $n$  belongs to model  $k$ .

probability that any given data point is in model  $k$

likelihood of  $t_n$  given model  $k$ 's slope, int., and  $\sigma_k^2$  parameters.

maximization is via the EM algorithm

$$E\text{-step: } \gamma_{nk} = E[z_{nk}] = \frac{\pi_k \Phi(t_n | x_n \beta_k, \sigma_k^2)}{\sum_{k=1}^K \pi_k \Phi(t_n | x_n \beta_k, \sigma_k^2)}$$

$$M\text{-step: } \pi_u = \frac{1}{N} \sum_{n=1}^N f_{nu} \quad \sum_j \pi_j \Phi(t_n | X_n \beta_j, \sigma_j^2)$$

regression step: WLS for each model  $u$ , weights given by  $\text{diag}(\gamma_{nu}) = R$

$$\beta^* = (X'RX)^{-1} X'RT$$

matrix of  $1u$ s  
including  
a constant.

target DV

BMA is directed at averaging predictions over a variety of models on the same data set, with the idea that we are not certain which of the models is "true" and we want to incorporate this uncertainty into our predictions.

$$\text{Ensemble: } \int y(x, z) f(z) dz = y(x)$$

$$\text{BMA: } \int y_m(x, z) f(m) dm = y(x)$$

$$y \sim x, z, \omega, \beta, \tau$$

Not sure about specification.

$$y \sim z$$

$$y \sim x + z$$

$$y \sim x + \omega$$

$$y \sim x + z + \omega$$

$$y \sim \beta_0 + \beta_1 z + \dots$$

$$\hat{y} \sim x + z + w$$

$$y \sim \beta_0 + \underline{\beta_1 z} + \dots$$

Run many models, and then combine results.

E.g. for  $\hat{\beta}_z = \sum_{\Delta} w^{(\Delta)} \hat{\beta}_z^{(\Delta)}$  for all candidate models  $\Delta$

but then  $w^{(\Delta)}$  is set according to the "likelihood" (posterior probability) the model is the best one.

BIC weights: 
$$w^{(\Delta)} = \frac{\exp\left(\frac{1}{2} \text{BIC}_{\Delta}^*\right)}{\sum_{\Delta} \exp\left(\frac{1}{2} \text{BIC}_{\Delta}\right)}$$

(assuming BIC is positively associated w/ model quality)

$$\text{BIC}_{\Delta} = 2 \ln L_{\Delta} - \log n \cdot k_{\Delta}$$

normalize BIC: 
$$\text{BIC}_{\Delta}^* = \max_{\Delta} (\text{BIC}_{\Delta}) = \text{BIC}_{\Delta}^{\#}$$

"leaps and bounds" algorithm \*

"Occam's window"

① only select models:

quality

$$\frac{\max_t (BIC_t)}{BIC_{\Delta}} \leq C = 20$$

② delete a model (t) if

parsimony

$$\frac{BIC_{\Delta}}{BIC_t} > 1 \quad \text{and} \quad \Delta C_t$$

in words: if  $t$  contains everything in  $\Delta$  but has a lower BIC .. exclude it.