

BOOSTING 2: CLASSIFICATION

-STATISTICAL MACHINE LEARNING-

Lecturer: Darren Homrighausen, PhD

ADDITIVE MODELS (FOR CLASSIFICATION)

As squared error loss isn't quite right for classification, **additive logistic regression** is a popular approach

Suppose $Y \in \{-1, 1\}$

$$\log \left(\frac{\mathbb{P}(Y = 1|X)}{\mathbb{P}(Y = -1|X)} \right) = \sum_{j=1}^p f_j(X_j) = F(X)$$

This gets inverted in the usual way to acquire a probability estimate

$$\pi(X) = \mathbb{P}(Y = 1|X) = \frac{e^{F(X)}}{1 + e^{F(X)}}$$

($f(X) = X^\top \beta$ gives us (linear) logistic regression, with classifier $g(X) = \text{sgn}(F(X))$)

These models are usually fit by numerically maximizing the binomial likelihood, and hence enjoy all the asymptotic optimality features of MLEs

ADDITIVE MODELS (FOR CLASSIFICATION)

EXAMPLE: In **R**, this can be fit with the package **gam**

In the **gam** package there is a dataset **kyphosis**

This dataset examines a disorder of the spine

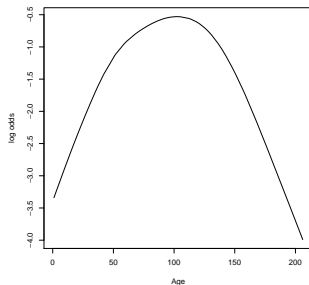
Let's look at two possible covariates **Age** and **Number**

(**Number** refers to the number of vertebrae that were involved in a surgery)

ADDITIVE MODELS (FOR CLASSIFICATION)

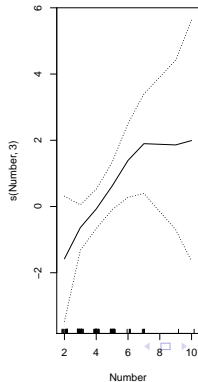
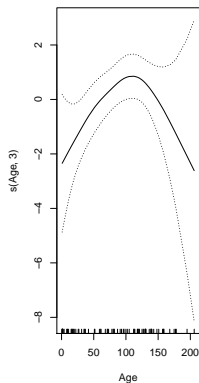
```
library(gam)
data(kyphosis)

out = gam(Kyphosis~s(Age,3),family=binomial,data=kyphosis)
out.pred = predict(out)
plot(sort(kyphosis$Age),out.pred[order(kyphosis$Age)],
     type='l',xlab='Age',ylab='log odds')
```



ADDITIVE MODELS (FOR CLASSIFICATION)

```
out = gam(Kyphosis ~ s(Age,3) + s(Number,3),  
          family = binomial, data=kyphosis)  
par(mfrow=c(1,2))  
plot(out,se=T)
```



ADDITIVE MODELS (FOR CLASSIFICATION)

FIT: backfitting approach within a Newton-Raphson iteration:

(This presumes $Y \in \{0, 1\}$)

Let \hat{f} be the current estimate of the GAM and

$$\hat{p}_i = (1 + \exp\{-\hat{f}(X_i)\})^{-1}$$

1. Form the “working response”

$$\tilde{Y}_i = \hat{f}(X_i) + \frac{Y_i - \hat{p}_i}{\hat{p}_i(1 - \hat{p}_i)}$$

(These working responses are)

2. Form weights $w_i = \hat{p}_i(1 - \hat{p}_i)$
3. **BACKFIT:** Smooth \tilde{Y}_i with a nonparametric procedure \mathcal{S} with weights w_i w.r.t. each feature

(Compare to usual update: $\hat{\beta} = (\mathbb{X}W\mathbb{X})^{-1}\mathbb{X}^T W(\mathbb{X}\hat{\beta} + W^{-1}(Y - \hat{\rho}))$)

Adaboost

ADABOOST OUTLINE

We give an overview of 'AdaBoost.M1.'

(Freund and Schapire (1997))

First, train the classifier as usual

(This is done by setting $w_i \equiv 1/n$)

At each step b , the misclassified observations have their weights increased

(Implicitly, this lowers the weight on correctly classified observations)

A new classifier is trained which emphasizes the previous mistakes

ADABOOST ALGORITHM

Assume $Y \in \{-1, 1\}$

1. Initialize $w_i \equiv 1/n$
2. For $b = 1, \dots, B$
 - 2.1 Fit $g_b(X)$ on \mathcal{D} , weighted by w_i
 - 2.2 Compute

$$R_b = \frac{\sum_{i=1}^n w_i \mathbf{1}(Y_i \neq g_b(X_i))}{\sum_{i=1}^n w_i}$$

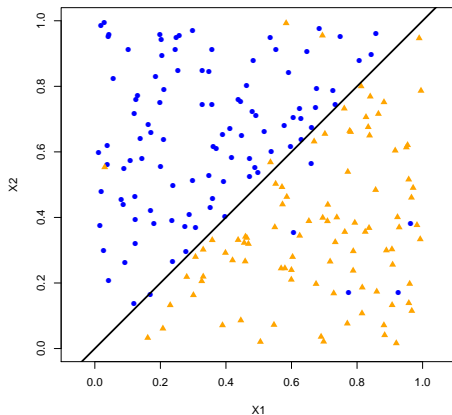
- 2.3 Find $\beta_b = \log((1 - R_b)/R_b)$
 - 2.4 Set $w_i \leftarrow w_i \exp\{\beta_b \mathbf{1}(Y_i \neq g_b(X_i))\}$
3. **OUTPUT:** $g(X) = \text{sgn}\left(\sum_{b=1}^B \beta_b g_b(X)\right)$

Some supporting simulations

SIMULATION: EQUAL PROBABILITY

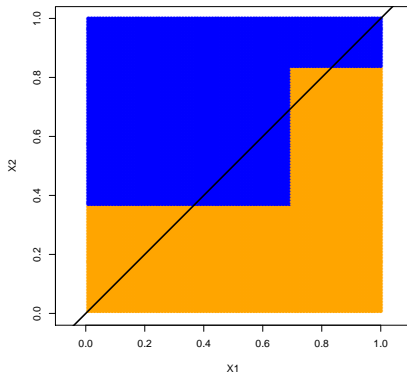
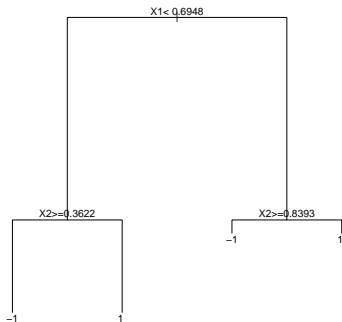
BASE LEARNER: 'depth 2-stumps'

(These are trees, but constrained to have no more than 4 terminal nodes)



SIMULATION: EQUAL PROBABILITY

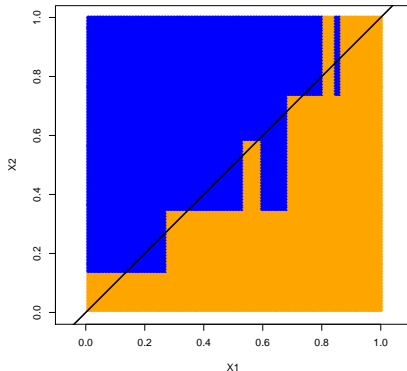
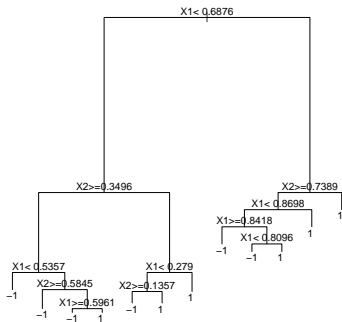
Using a depth-2 stump:



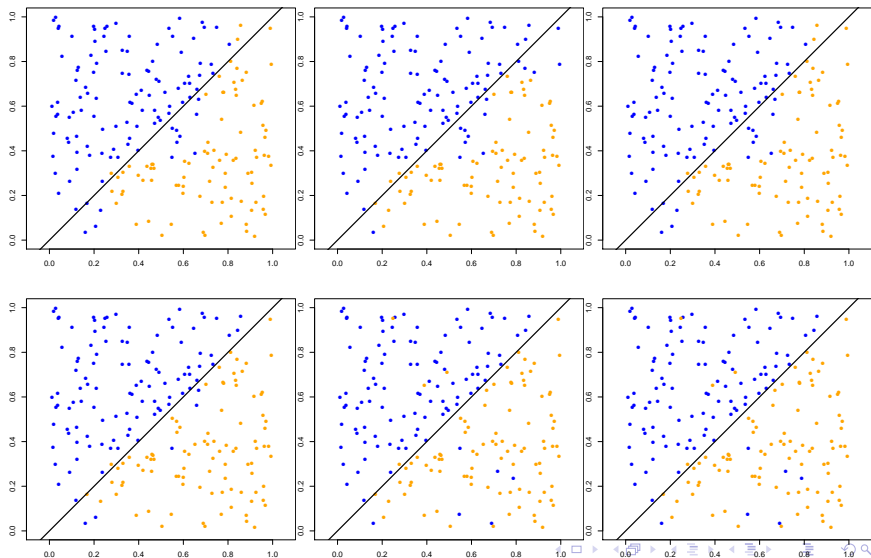
SIMULATION: EQUAL PROBABILITY

Using an unpruned tree:

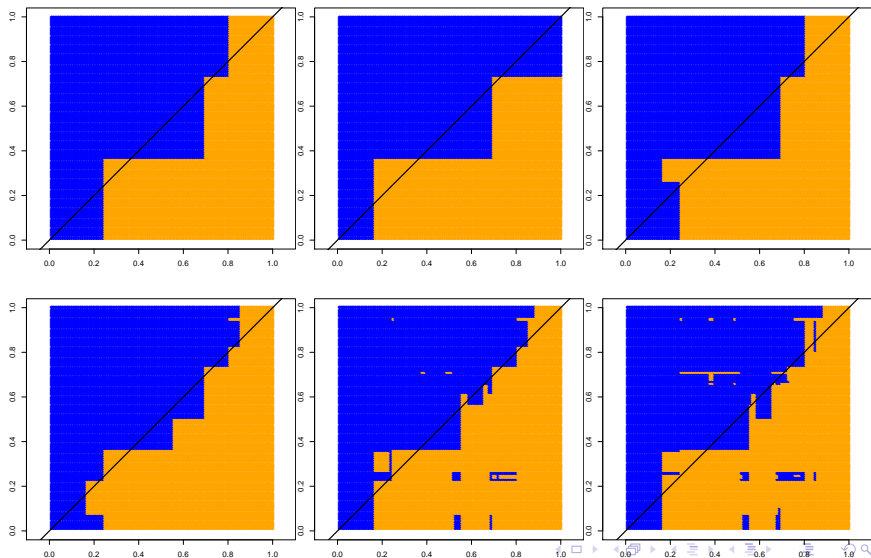
(Note that I used **rpart**, which parameterizes splits differently than **tree**)



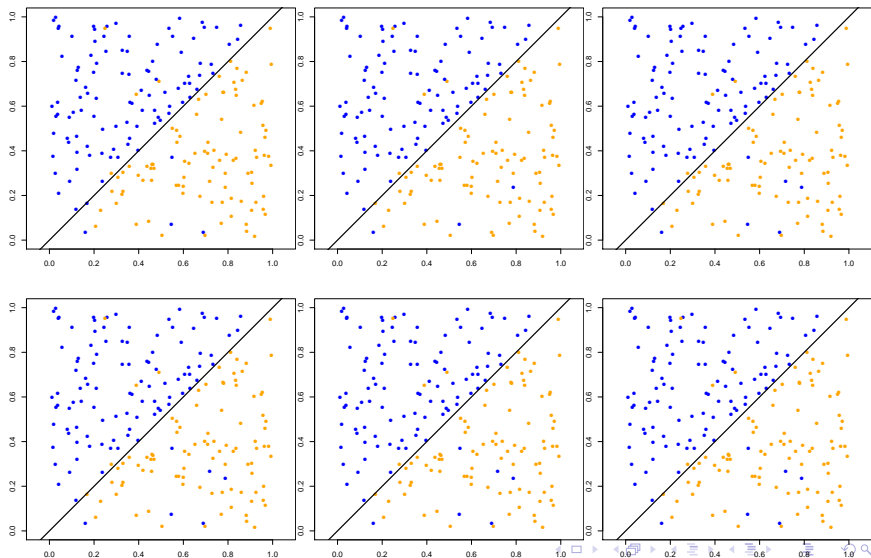
ADABOOST: INCREASING B (TRAIN)



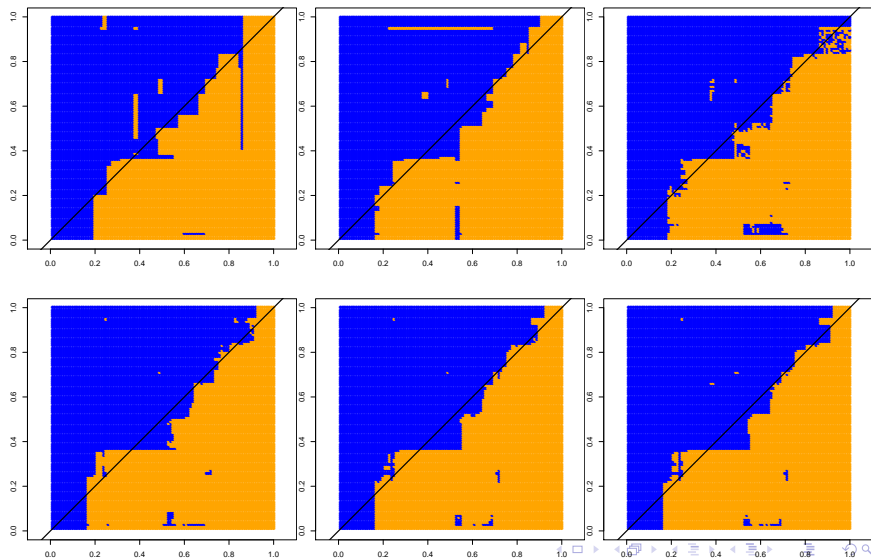
ADABOOST: INCREASING B (TEST)



RANDOM FOREST: INCREASING B (TRAIN)



RANDOM FOREST: INCREASING B (TEST)



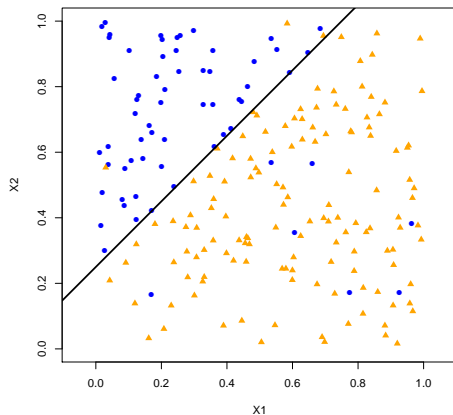
RESULTS: CONFUSION MATRICES

(These are at best B solution)

		Truth		Mis-Class	
		-1	1		
Our Preds	UNPRUNED	-1	84	18	14.5%
		1	11	87	
	STUMP	-1	77	16	17%
		1	18	89	
	BOOST	-1	92	8	8%
		1	5	92	
	RF	-1	87	9	8.5%
		1	8	96	

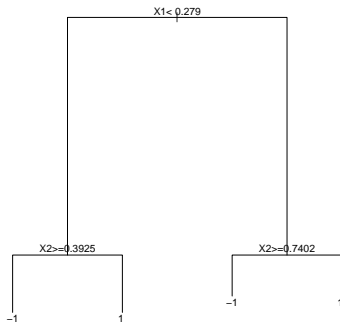
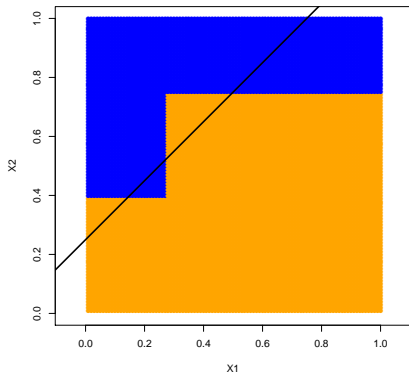
SIMULATION: UNEQUAL PROBABILITY

Let's change the simulation so that the class probabilities aren't the same



SIMULATION: UNEQUAL PROBABILITY

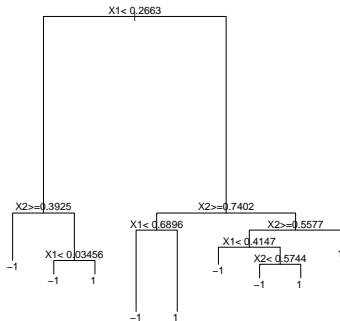
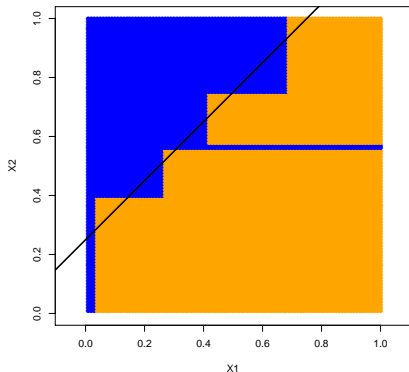
Using a depth-2 stump:



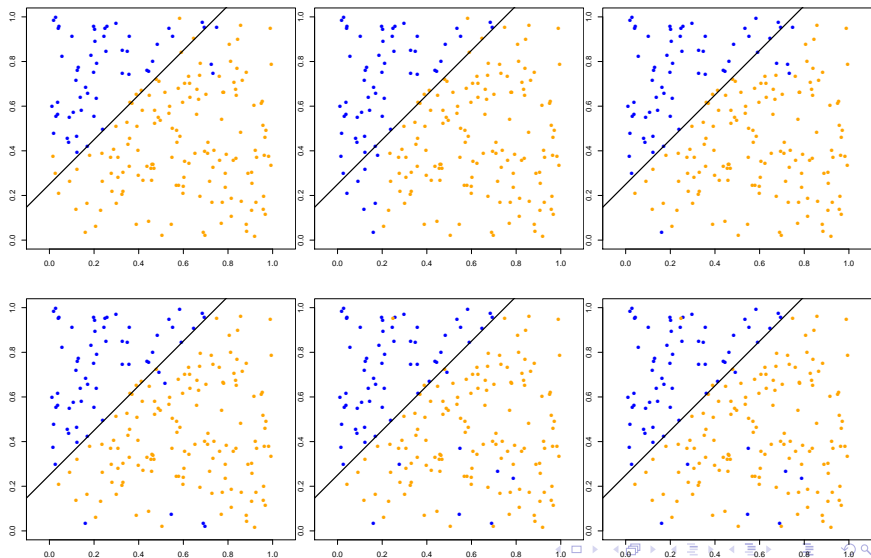
SIMULATION: UNEQUAL PROBABILITY

Using an unpruned tree:

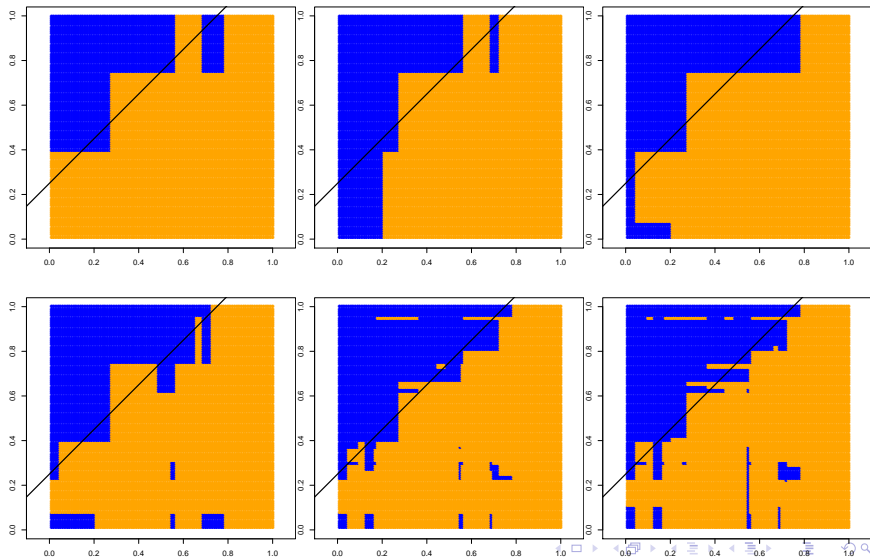
(Note that I used **rpart**, which parameterizes splits differently than **tree**)



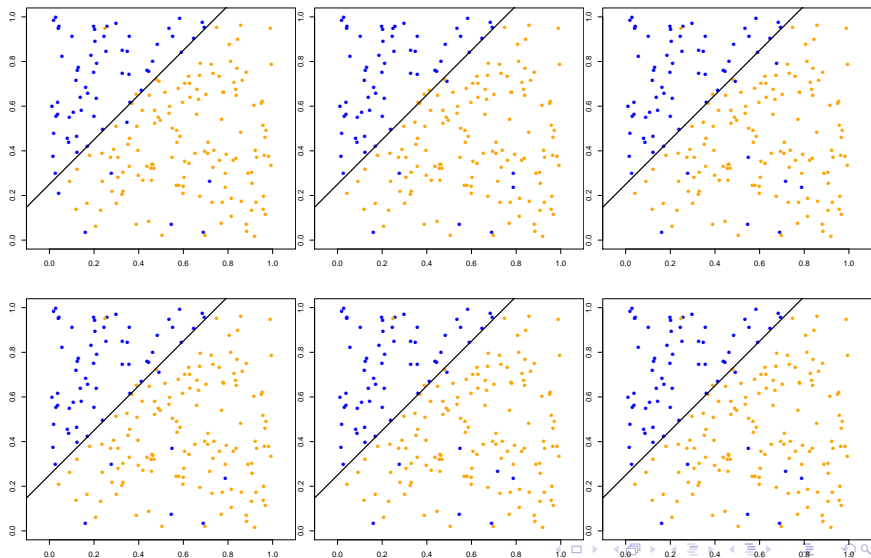
ADABOOST: INCREASING B (TRAIN)



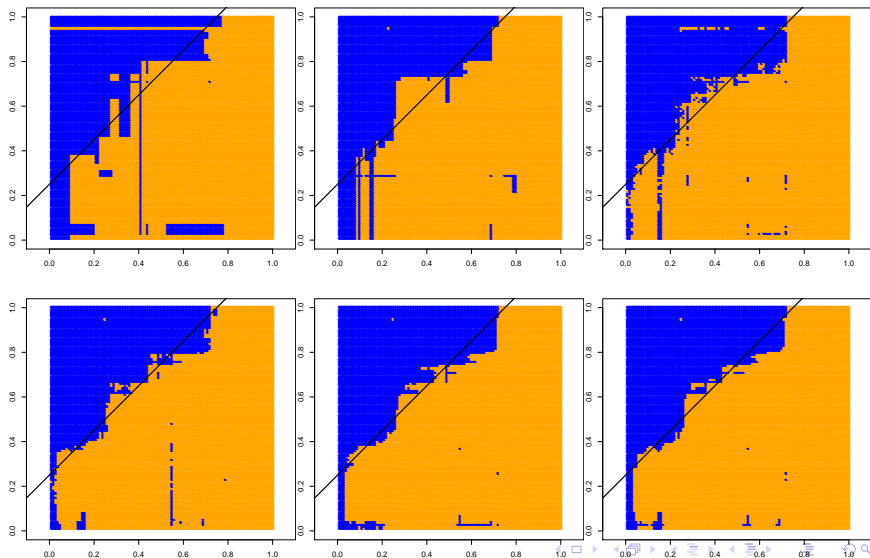
ADABOOST: INCREASING B (TEST)



RANDOM FOREST: INCREASING B (TRAIN)



RANDOM FOREST: INCREASING B (TEST)



RESULTS: CONFUSION MATRICES

(These are at best B solution)

		Truth		
		-1	1	Mis-Class
Our Preds	UNPRUNED	-1	51 10	10.5%
		1	11 128	
	STUMP	-1	50 22	17%
		1	12 116	
	BOOST	-1	51 12	7.5%
		1	3 134	
	RF	-1	53 5	7%
		1	9 133	

Back to Algorithms

ADABOOST

This algorithm became known as ‘discrete AdaBoost’

(This is due to the base classifier returning a discrete label)

⋮

1. Compute

$$R_b = \frac{\sum_{i=1}^n w_i \mathbf{1}(Y_i \neq g_b(X_i))}{\sum_{i=1}^n w_i}$$

⋮

2. Set

$$w_i \leftarrow w_i \exp\{\beta_b \mathbf{1}(Y_i \neq g_b(X_i))\}$$

This was adapted to real-valued predictions in Real AdaBoost

(In particular, probability estimates)

(This terminology was introduced Friedman “Functional Gradient Boosting” (2001))

REAL ADABOOST

1. Initialize $w_i \equiv 1/n$
2. For $b = 1, \dots, B$
 - 2.1 Fit the classifier on \mathcal{D} , weighted by w_i and produce $p_b(X) = \hat{\mathbb{P}}_w(Y = 1|X)$
 - 2.2 Set $\hat{f}_b(X) \leftarrow \frac{1}{2} \log(p_b(X)/(1 - p_b(X)))$
 - 2.3 Set $w_i \leftarrow w_i \exp\{-Y_i \hat{f}_b(X_i)\}$
3. **OUTPUT:** $g(X) = \text{sgn} \left(\sum_{b=1}^B \hat{f}_b(X) \right)$

This is referred to as **Real AdaBoost** and it used the class probability estimates to construct the contribution of the b^{th} classifier, instead of the estimated label

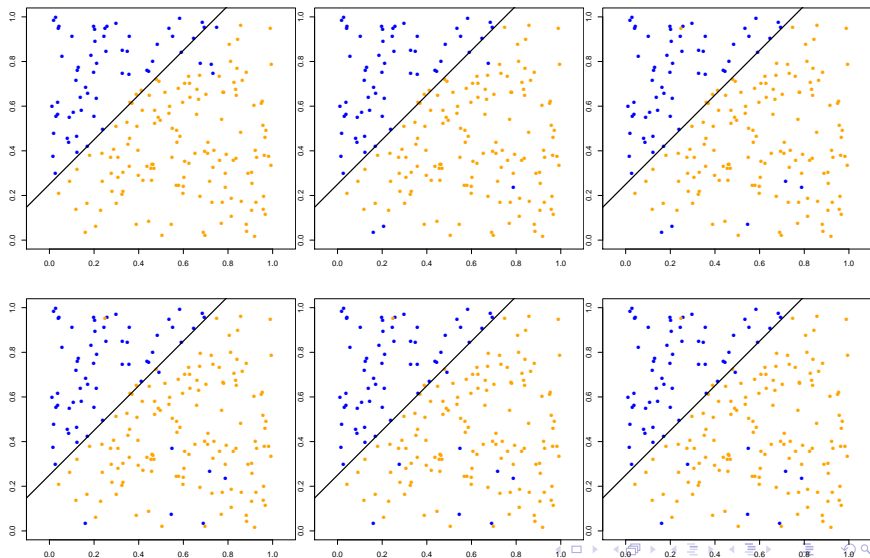
REAL ADABOOST

1. Initialize $w_i \equiv 1/n$
2. For $b = 1, \dots, B$
 - 2.1 Fit the classifier on \mathcal{D} , weighted by w_i and produce $p_b(X) = \hat{\mathbb{P}}_w(Y = 1|X)$
 - 2.2 Set $\hat{f}_b(X) \leftarrow \frac{1}{2} \log(p_b(X)/(1 - p_b(X)))$
 - 2.3 Set $w_i \leftarrow w_i \exp\{-Y_i \hat{f}_b(X_i)\}$
3. **OUTPUT:** $g(X) = \text{sgn}\left(\sum_{b=1}^B \hat{f}_b(X)\right)$

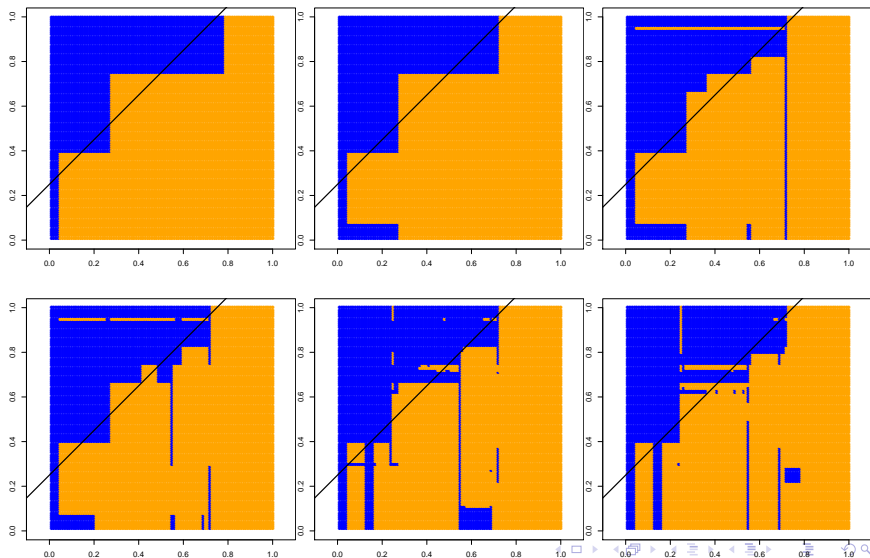
This is referred to as **Real AdaBoost** and it used the class probability estimates to construct the contribution of the b^{th} classifier, instead of the estimated label

(Care is needed when computing \hat{f}_b in practice due to numerical issues with probabilities near 0 or 1. We need to be sure that induced weights are nonnegative.)

REAL ADABOOST: INCREASING B (TRAIN)



REAL ADABOOST: INCREASING B (TEST)



ADABOOST INTUITION

QUESTION: Why does this work?

ONE ANSWER: Boosting fits an additive model

$$F_B(X) = \sum_{b=1}^B \beta_b \phi(X, \theta_b)$$

where

- β are weights
- ϕ is some **base learner** that depends on parameters θ

(EXAMPLE: Trees with all of its splits and terminal node values)

OVERALL: Both discrete and real AdaBoost can be interpreted as stage wise estimation procedures for fitting additive logistic regression models

(DISCRETE) ADABOOST INTERPRETATION

Forward stagewise additive modeling:

(Using a general loss ℓ)

1. $\beta_b, \theta_b = \operatorname{argmin}_{\beta, \theta} \sum_{i=1}^n \ell(Y_i, F_{b-1}(X_i) + \beta \phi(X_i, \theta))$
2. Set $F_b(X) = F_{b-1}(X) + \beta_b \phi(X; \theta_b)$

AdaBoost implicitly does this by use of the **exponential loss function**

$$\ell(Y, F) = \exp\{-YF(X)\}$$

and basis functions $\phi(x, \theta) = g_b(X)$

ADABOOST INTUITION

Suppose we minimize exponential loss in a forward stagewise manner

Doing the forward selection for this loss, we get

$$(\beta_b, g_b) = \operatorname{argmin}_{\beta, g} \sum_{i=1}^n \exp\{-Y_i(F_{b-1}(X_i) + \beta g(X_i))\}$$

ADABOOST INTUITION

Rewriting:

$$\begin{aligned}(\beta_b, g_b) &= \operatorname{argmin}_{\beta, g} \sum_{i=1}^n \exp\{-Y_i(F_{b-1}(X_i) + \beta g(X_i))\} \\&= \operatorname{argmin}_{\beta, g} \sum_{i=1}^n \exp\{-Y_i F_{b-1}(X_i)\} \exp\{-Y_i \beta g(X_i)\} \\&= \operatorname{argmin}_{\beta, g} \sum_{i=1}^n w_i \exp\{-Y_i \beta g(X_i)\}\end{aligned}$$

Where

- Define $w_i = \exp\{-Y_i F_{b-1}(X_i)\}$
(This is independent of β, g)
- $\sum_{i=1}^n w_i \exp\{-Y_i \beta g_b(X_i)\}$ needs to be optimized

ADABOOST INTUITION

Note that

$$\begin{aligned}\sum_{i=1}^n w_i \exp\{-\beta Y_i g(X_i)\} &= e^{-\beta} \sum_{i: Y_i = g(X_i)} w_i + e^{\beta} \sum_{i: Y_i \neq g(X_i)} w_i \\ &= (e^{\beta} - e^{-\beta}) \sum_{i=1}^n w_i \mathbf{1}(Y_i \neq g(X_i)) + \\ &\quad + e^{-\beta} \sum_{i=1}^n w_i\end{aligned}$$

As long as $(e^{\beta} - e^{-\beta}) \geq 0$, we can find

$$g_b = \operatorname{argmin}_g \sum_{i=1}^n w_i \mathbf{1}(Y_i \neq g(X_i))$$

(Note: If $(e^{\beta} - e^{-\beta}) < 0$, then $\beta < 0$. However, as $\beta_b = \log((1 - R_b)/R_b)$, this implies $R > 1/2$. Hence, we would flip the labels and get $R \leq 1/2$.)

REMINDER: ADABOOST

1. Initialize $w_i \equiv 1/n$

2. For $b = 1, \dots, B$

2.1 Fit $g_b(x)$ on \mathcal{D} , weighted by w_i

(In this case, we would be growing the (heavily pruned) tree via minimizing misclassifications)

2.2 Compute

$$R_b = \frac{\sum_{i=1}^n w_i \mathbf{1}(Y_i \neq g_b(X_i))}{\sum_{i=1}^n w_i}$$

2.3 Find $\beta_b = \log((1 - R_b)/R_b)$

2.4 Set $w_i \leftarrow w_i \exp\{\beta_b \mathbf{1}(Y_i \neq g_b(X_i))\}$

3. **OUTPUT:** $g(X) = \text{sgn} \left(\sum_{b=1}^B \beta_b g_b(X) \right)$

REMINDER: ADABOOST

1. Initialize $w_i \equiv 1/n$
2. For $b = 1, \dots, B$
 - 2.1 Fit $g_b(x)$ on \mathcal{D} , weighted by w_i
(In this case, we would be growing the (heavily pruned) tree via minimizing misclassifications)
 - 2.2 Compute

$$R_b = \frac{\sum_{i=1}^n w_i \mathbf{1}(Y_i \neq g_b(X_i))}{\sum_{i=1}^n w_i}$$

- 2.3 Find $\beta_b = \log((1 - R_b)/R_b)$
 - 2.4 Set $w_i \leftarrow w_i \exp\{\beta_b \mathbf{1}(Y_i \neq g_b(X_i))\}$
3. **OUTPUT:** $g(X) = \text{sgn} \left(\sum_{b=1}^B \beta_b g_b(X) \right)$

ADABOOST INTUITION

GOAL: Minimize

$$\sum_{i=1}^n w_i \exp\{-\beta Y_i g_b(X_i)\}$$

(Here, we have fixed $g = g_b$)

We showed this can be written

$$\sum_{i=1}^n w_i \exp\{-\beta Y_i g_b(X_i)\} = (e^{\beta} - e^{-\beta}) R_b W + e^{-\beta} W \quad (W = \sum w_i)$$

Take derivative with respect to β

$$(e^{\beta} + e^{-\beta}) R_b W - e^{-\beta} W \stackrel{\text{set}}{=} 0 \stackrel{\text{set}}{=} e^{\beta} R_b + e^{-\beta} (R_b - 1)$$

Solve for β to find $\beta_b = 1/2 \log[(1 - R_b)/R_b]$

REMINDER: ADABOOST

1. Initialize $w_i \equiv 1/n$

2. For $b = 1, \dots, B$

2.1 Fit $g_b(x)$ on \mathcal{D} , weighted by w_i

(This step is finding the next best version of the classifier, trained on weighted data and added to the previous classifiers)

2.2 Compute

$$R_b = \frac{\sum_{i=1}^n w_i \mathbf{1}(Y_i \neq g_b(X_i))}{\sum_{i=1}^n w_i}$$

2.3 Find $\beta_b = \log((1 - R_b)/R_b)$

2.4 Set $w_i \leftarrow w_i \exp\{\beta_b \mathbf{1}(Y_i \neq g_b(X_i))\}$

3. **OUTPUT:** $g(x) = \text{sgn}\left(\sum_{b=1}^B \beta_b g_b(x)\right)$

REMINDER: ADABOOST

1. Initialize $w_i \equiv 1/n$

2. For $b = 1, \dots, B$

2.1 Fit $g_b(x)$ on \mathcal{D} , weighted by w_i

(This step is finding the next best version of the classifier, trained on weighted data and added to the previous classifiers)

2.2 Compute

$$R_b = \frac{\sum_{i=1}^n w_i \mathbf{1}(Y_i \neq g_b(X_i))}{\sum_{i=1}^n w_i}$$

2.3 Find $\beta_b = \log((1 - R_b)/R_b)$

2.4 Set $w_i \leftarrow w_i \exp\{\beta_b \mathbf{1}(Y_i \neq g_b(X_i))\}$

3. **OUTPUT:** $g(x) = \text{sgn}\left(\sum_{b=1}^B \beta_b g_b(x)\right)$

ADABOOST INTUITION

The approximation is updated

$$F_b(X) = F_{b-1}(X) + \beta_b g_b(X)$$

This causes the weights

$$w_i^{(b+1)} = \exp\{-Y_i F_b(X_i)\} = w_i^{(b)} \exp\{-\beta_b Y_i g_b(X_i)\}$$

Using $-Y_i g_b(X_i) = 2\mathbf{1}(Y_i \neq g_b(X_i)) - 1$, this becomes

$$w_i^{(b+1)} \propto w_i^{(b)} \exp\{\beta_b \mathbf{1}(Y_i \neq g_b(X_i))\}$$

where $\beta_b \leftarrow 2\beta_b$, giving the last step of the algorithm

WHY EXPONENTIAL LOSS?

As previously alluded to, AdaBoost was originally the result of a constructive proof

Friedman et al. (2000) showed that it was equivalent to greedily fitting/learning a basis expansion with **exponential loss**

Note that this is compared to fitting:

$$\min_{(\beta_b), (g_b)} \sum_{i=1}^n \ell(Y_i, F(X_i))$$

(Which can be difficult or overfit, e.g. minimizing training error)

versus:

$$\min_{\beta, g} \sum_{i=1}^n \ell(Y_i, F_{b-1}(X_i) + \beta g(X_i))$$

WHY EXPONENTIAL LOSS?

It can be shown that

$$\operatorname{argmin}_F \mathbb{E}[\exp\{-YF(X)\}|X] = \frac{1}{2} \log \left(\frac{\mathbb{P}(Y = 1|X)}{\mathbb{P}(Y = -1|X)} \right)$$

Hence, we are estimating (half) the log odds

→ use the `sgn` rule for classification

Of course, using the **functional gradient descent** approach will work as well after specifying an alternative differentiable loss

REGULARIZATION IN GENERAL GBM

Regularization: prevent fitting to the training data too well

In boosting, this is usually done via

- Choosing the parameter B via a risk estimate
- Including a **learning rate**

$$F_B(X) = F_{B-1}(X) + \lambda f_b(X)$$

(Note that the functional gradient descent already includes this. AdaBoost needs this to be included)

Of course, there is a strong interaction between these two parameters

It has been observed repeatedly that

- setting λ to a small constant and choosing B via a risk estimate....
... provides better performance than:
- setting $\lambda = 1$ and choosing B via a risk estimate

BOOSTING: THE CONTROVERSY

CLAIM: Boosting is another version of bagging

The early versions of Boosting involved (weighted) resampling

Therefore, it was initially speculated that a connection with **bagging** explained its performance

However, boosting continues to work well when

- The algorithm is trained on weighted data rather than on sampling with weights

(This removes the randomization component that is essential to bagging)

- Weak learners are used that have high bias and low variance

(This is the **opposite** of what is prescribed for bagging)

BOOSTING: THE CONTROVERSY

CLAIM: Boosting fits an adaptive additive model which explains its effectiveness

The previous results appeared in Friedman et al. (2000) and claimed to have 'solved' the mystery of boosting

A crucial property of boosting is that it is difficult to overfit, especially with a small learning rate

However, the additive model view really should translate into intuition of 'over fitting is a major concern,' as it is with additive models

BOOSTING: THE CONTROVERSY

As adaBoost fits an additive model in the base classifier, it cannot have higher order interactions than the base classifier

For instance, a stump would provide a purely additive fit

(It only splits on one variable. In general, the complexity of a tree can be interpreted as the number of included interactions)

It stands to reason, then, if the Bayes' rule is additive in a similar fashion, stumps should perform well in Boosting

BOOSTING: THE CONTROVERSY

A recent paper investigating this property did substantial simulations using underlying purely additive models

(Mease, Wyner (2008))

Here is an example figure from their paper:

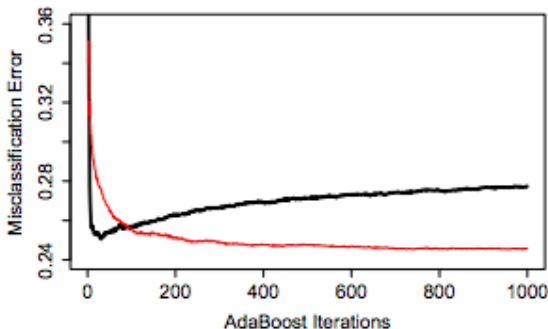


FIGURE: Black, bold line: Stumps. Red, thin line: 8-node trees

BOOSTING: THE CONTROVERSY CONTINUES

Ultimately, interpretations are just modes of human comprehension

The value of the insight is whether it provides fruitful thought about the idea

From this perspective, AdaBoost fits an additive model.

However, many of the other connections are still of debatable value

(For example, LogitBoost)

NEXT LECTURES

Discuss two current, popular algorithms and their **R** implementations

- **GBM**
- **XGBoost**