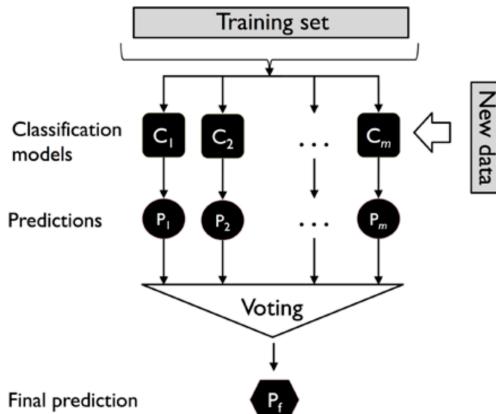
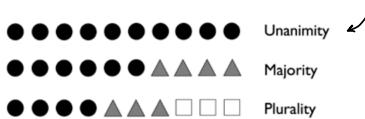


CHAPTER 7 - COMBINING DIFFERENT MODELS FOR ENSEMBLE LEARNING

LEARNING WITH ENSEMBLE

- ENSEMBLE METHOD → COMBINE CLASSIFIERS TO META-CLASSIFIER → BETTER GENERALIZATION PERFORMANCE.
- MAJORITY VOTING PRINCIPLE → CLASS LABEL WITH MORE THAN 50% OF VOTES → REFERS TO BINARY CLASS SETTINGS ONLY.
↳ GENERALISE TO MULTICLASS SETTINGS → PLURALITY VOTING (UK → 'ABSOLUTE' + 'RELATIVE' VOTING) → CLASS LABEL WITH MOST VOTES



- TRAINING DATASET → TRAINING m DIFFERENT CLASSIFIERS
- DEPENDING ON TECHNIQUE → USE DIFFERENT CLASSIFICATION ALGORITHM
 - ↳ DECISION TREES, SVM, LOGISTIC REGRESSION
- CAN ALSO USE SAME BASE CLASSIFICATION ALGORITHM → FIT DIFFERENT SUBSETS OF TRAINING DATASET
 - ↳ BAGGING + RANDOM FOREST → COMBINING DIFFERENT DECISION TREES.
- ⇐ GENERAL ENSEMBLE APPROACH.

SEBASTIAN RASCHKA

- PREDICT CLASS LABEL WITH MAJORITY OR PLURALITY VOTING → COMBINE PREDICTED CLASS LABELS CLASSIFIER, C_i , + SELECTED CLASS LABEL, \hat{y} .

$$\hat{y} = \text{mode}\{C_1(x), C_2(x), \dots, C_m(x)\} \rightarrow \text{MOST FREQUENT EVENT, mode = STATS.}$$

- BINARY CLASSIFICATION + MAJORITY VOTE:

$$C(x) = \text{sign}\left[\sum_i^n C_i(x)\right] = \begin{cases} 1 & \text{if } \sum_i^n C_i(x) \geq 0 \\ -1 & \text{otherwise} \end{cases}$$

- ENSEMBLE METHOD WORKS BETTER THAN INDIVIDUAL CLASSIFIERS → PROVEN WITH COMBINATORICS.

- ASSUMPTIONS FOR BINARY CLASSIFICATION TASK:

↳ n -BASE CLASSIFIERS → EQUAL ERROR RATE, ϵ .

↳ CLASSIFIERS INDEPENDENT + ERROR RATE NOT CORRELATED.

$$P(y \geq k) = \sum_{k=0}^n \binom{n}{k} \epsilon^k (1-\epsilon)^{n-k} = E_{\text{ensemble}} \rightarrow \text{ERROR PROBABILITY OF ENSEMBLE OF BASE CLASSIFIERS AS PROBABILITY MASS FUNCTION OF BINOMIAL DISTRIBUTION}$$

↳ BINOMIAL COEFFICIENT

- E.G. 11 CLASSIFIERS ($n=11$) + ERROR RATE 0.25 ($\epsilon=0.25$)

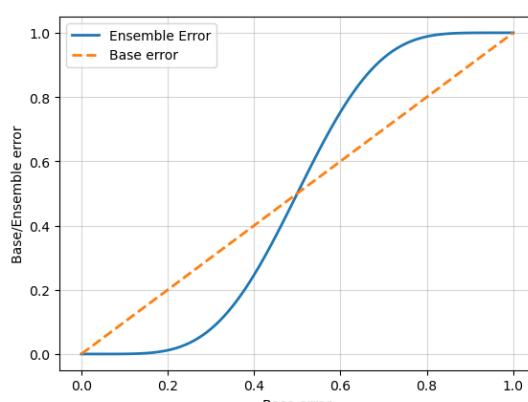
$$P(y \geq k) = \sum_{k=0}^n \binom{n}{k} 0.25^k (1-0.25)^{n-k} = 0.034 \rightarrow \text{ERROR RATE LOWER THAN OF EACH INDIVIDUAL CLASSIFIER.}$$

BINOMIAL COEFFICIENT

- NUMBER OF WAYS WE CHOOSE k FROM SET OF SIZE n ⇒ " n choose k "

$$\frac{n!}{(n-k)! k!}$$

ERROR PROBABILITY OF ENSEMBLE
ALWAYS BETTER



COMBINING CLASSIFIERS VIA MAJORITY VOTE

SINGLE MAJORITY VOTE CLASSIFIER

$$\hat{y} = \arg_{\hat{z}} \max \sum_{j=1}^m w_j \chi_A(C_j(x) = i)$$

w_j = WEIGHT ASSOCIATED WITH BASE CLASSIFIER C_j
 \hat{y} = PREDICTED CLASS LABEL OF THE ENSEMBLE
 A = SET OF UNIQUE CLASS LABELS
 χ_A = CHARACTERISTIC OR INDICATOR FUNCTION

- FOR EQUAL WEIGHTS \rightarrow SIMPLIFIED EQUATION:

$$\hat{y} = \text{mode}\{C_1(x), C_2(x), \dots, C_m(x)\}$$

EXAMPLE TO UNDERSTAND CONCEPT OF WEIGHTS:

- $C_1(x) \rightarrow 0, C_2(x) \rightarrow 0, C_3(x) \rightarrow 1 \quad \hat{y} = \text{mode}\{0, 0, 1\} = 0$

- WEIGHTS = C_1 AND $C_2 = 0.2, C_3 = 0.6$

$$\hat{y} = \arg_{\hat{z}} \max [0.2 \times i_0 + 0.2 \times i_0, 0.6 \times i_1] = 1$$

- $3 \times 0.2 = 0.6 \Rightarrow C_3$ HAS THREE TIMES THE WEIGHT THAN C_1 AND $C_2 \rightarrow \hat{y} = \text{mode}\{0, 0, 1, 1, 1\} = 1$

- MODIFIED VERSION OF MAJORITY VOTE FOR PREDICTING CLASS LABELS: ← USING 'predict_proba' METHOD FROM LOGISTIC REGRESSION.

$$\hat{y} = \arg_{\hat{z}} \max \sum_{j=1}^m w_j p_{i,j} \quad p_{i,j} = \text{PREDICTED PROBABILITY OF } j\text{TH CLASS LABEL } i$$

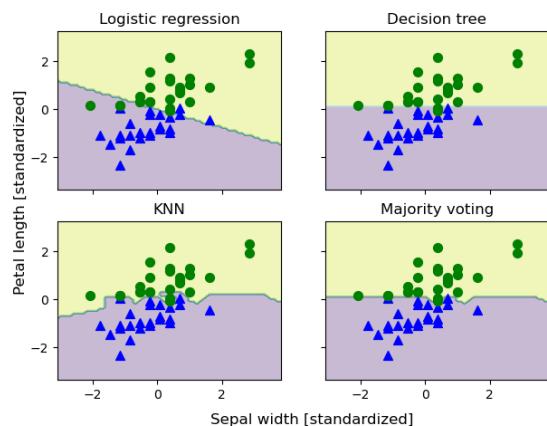
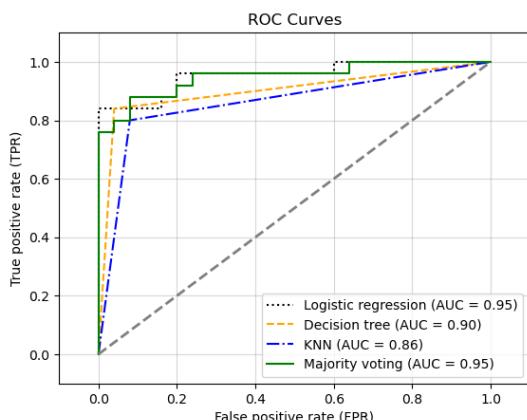
BINARY CLASSIFICATION PROBLEM + SAME WEIGHTS

- $C_1(x) \rightarrow [0.9, 0.1], C_2(x) \rightarrow [0.8, 0.2], C_3(x) \rightarrow [0.4, 0.6]$

$$p(i_0|x) = 0.2 \times 0.9 + 0.2 \times 0.8 + 0.6 \times 0.4 = 0.58$$

$$p(i_1|x) = 0.2 \times 0.1 + 0.2 \times 0.2 + 0.6 \times 0.6 = 0.42$$

$$\hat{y} = \arg_{\hat{z}} \max [p(i_0|x), p(i_1|x)] = 0$$



- ROC CURVE COMPUTED TO TEST IF 'MAJORITY VOTE CLASSIFIER' DOES GENERALISE WELL WITH UNSEEN DATA.

- DECISION REGIONS OF ENSEMBLE CLASSIFIER (STANDARDISED)

STACKING (ENSEMBLE METHODS)

- DAVID H. WOLPERT - "STACKED GENERALIZATION." 1992

INPUT: TRAINING DATA $D = \{(x_i, y_i)\}_{i=1}^n$ ($x_i \in \mathbb{R}^m$, $y_i \in \mathcal{Y}$)

OUTPUT: ENSEMBLE CLASSIFIER H

STEP 1: LEARN FIRST-LEVEL CLASSIFIERS.

```
for t ← 1 to T do
    LEARN A BASE CLASSIFIER  $h_t$  BASED ON D
end for
```

$$\{h_1, h_2, \dots, h_T\}$$

GREEN BOX = SAME AS
MAJORITY VOTING

STEP 2: CONSTRUCT NEW DATA SETS FROM D

for i ← 1 to n do

CONSTRUCT A NEW DATA SET THAT CONTAINS $\{x'_i, y'_i\}$, WHERE $x'_i = \{h_1(x_i), h_2(x_i), \dots, h_T(x_i)\}$

end for

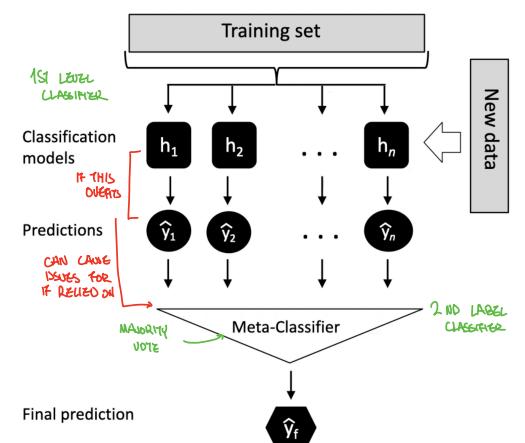
FOR EACH TRAINING EXAMPLE

$\xrightarrow{\text{MODIFIED FEATURE VECTOR}}$ $\xleftarrow{\text{PREDICTED CLASS LABELS VIA SKIT-LEARN}}$
mode

STEP 3: LEARN A SECOND-LEVEL CLASSIFIER → LEARN NEW CLASSIFIER BASED ON PREDICTIONS

LEARN A NEW CLASSIFIER H' BASED ON RECENTLY CONSTRUCTED DATA SET.

return $H(x) = h'(h_1(x), h_2(x), \dots, h_T(x))$

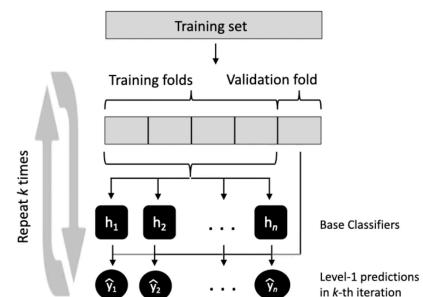


PROBLEM WITH STACKING ??
PROVE TO OVERFITTING

IMPROVEMENT OF STACKING → CROSS-VALIDATION

STACKING → K-FOLD CLASSIFICATION FOR 2ND LEVEL CLASSIFIER

↳ USE K PREDICTIONS FOR



MODIFIED STACKING ALGORITHM

- INPUT + OUTPUT STAYS SAME

STEP 1: CROSS VALIDATION APPROACH PREPARATION FOR 2ND-LEVEL CLASSIFIER

RANDOMLY SPLIT D INTO K EQUAL-SIZED SUBSETS → $D = \{D_1, D_2, \dots, D_K\}$

for $k \leftarrow 1$ to K do

STEP 1.1 LEARN 1ST-LEVEL CLASSIFIERS

for $t \leftarrow 1$ to T do

LEARN CLASSIFIER h_{kt} FROM D/D_k

end for

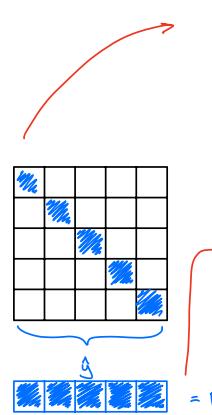
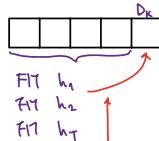
STEP 1.2 CONSTRUCT TRAINING SET FOR 2ND-LEVEL CLASSIFIER

for $x_i \in D_k$ do

GET $\{x'_i, y'_i\}$ WHERE $x'_i = \{h_1(x_i), h_2(x_i), \dots, h_T(x_i)\}$

end for

end for



STEP 2: FIT 2ND-LEVEL CLASSIFIER

STEP 3: RE-LEARN 1ST-LEVEL CLASSIFIER
FIT ON WHOLE DATASET

All level-1 predictions

↓ Train

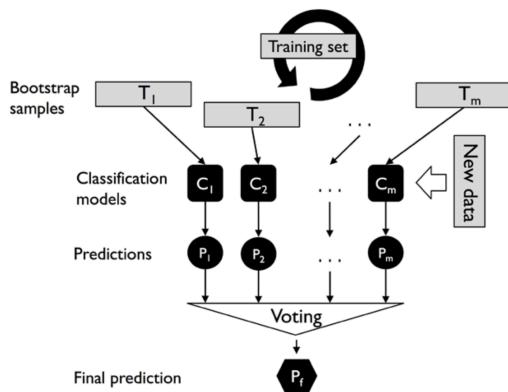
Meta-Classifier

Final prediction

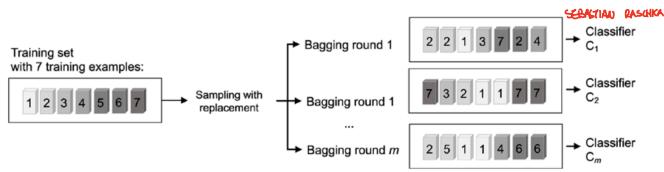
SEBASTIAN RASCHKA

BAGGING - ENSEMBLE CLASSIFIERS FROM BOOTSTRAP SAMPLES

- CLOSELY RELATED TO 'MAJORITY/VOTE CLASSIFIER'
- INSTEAD OF SAME TRAINING DATASET TO FIT INDIVIDUAL CLASSIFIERS → DRAW BOOTSTRAP SAMPLES (RANDOM SAMPLES WITH REPLACEMENT) FROM INITIAL TRAINING DATASET
↳ BOOTSTRAP AGGREGATING



SEBASTIAN RASCHKA



- RANDOM SAMPLES OBTAINED VIA BAGGING → BAGGING ROUND ...
- EACH SUBSET CONTAINS SAME DUPLICATES + SOME ORIGINAL EXAMPLES DON'T APPEAR
↳ SAMPLING WITH REPLACEMENT
- INDIVIDUAL CLASSIFIERS FIT TO BOOTSTRAP SAMPLES → PREDICTIONS COMBINED WITH MAJORITY VOTE.
- BAGGING RELATED TO RANDOM FOREST CLASSIFIER
↳ RANDOM FOREST → SPECIAL EXAMPLE → RANDOM FEATURE SUBSETS ALSO USED WHEN FITTING INDIVIDUAL DECISION TREES.

- COMPLEX CLASSIFICATION TASKS + DATASET'S HIGH DIMENSIONALITY CAN LEAD TO OVERFITTING IN SINGLE DECISION TREE → BAGGING COMES IN HANDY
- BAGGING EFFECTIVE APPROACH TO REDUCING VARIANCE OF MODEL
↳ INEFFECTIVE IN REDUCING MODEL BIAS → MODELS TOO SIMPLE TO CAPTURE TRENDS IN DATA.
↳ WHY WE WANT TO PERFORM BAGGING ON ENSEMBLE OF CLASSIFIERS WITH LOW BIAS → E.G. UNPRUNED DECISION TREES

ADAPTIVE BOOSTING - LEVERAGING WEAK LEARNERS

- ADAPTIVE BOOSTING (AdaBoost)
- ↳ ROBERT E. SCHAPIRA (1990) → 'THE STRENGTH OF WEAK LEARNABILITY.'
- BOOSTING - ENSEMBLE CONSISTS OF VERY SIMPLE BASE CLASSIFIERS → REFERRED TO AS 'WEAK LEARNERS' → SLIGHT EDGE OVER RANDOM GUESSING
- ↳ KEY CONCEPT → FOCUS ON TRAINING EXAMPLES THAT ARE HARD TO CLASSIFY
WEAK LEARNERS SUBSEQUENTLY LEARN FROM MISCLASSIFIED EXAMPLES TO IMPROVE PERFORMANCE OF ENSEMBLE

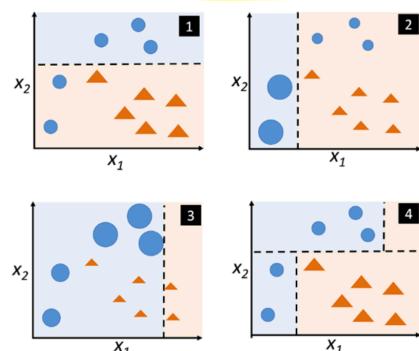
How Adaptive Boosting works:

- INITIAL FORMULATION - ALSO USES RANDOM SUBSETS OF TRAINING EXAMPLES FROM TRAINING DATASET WITHOUT REPLACEMENT.

FOUR KEY STEPS:

1. DRAW RANDOM SUBSET (SAMPLE) FROM TRAINING EXAMPLES, D WITHOUT REPLACEMENT FROM TRAINING DATASET, D FOR WEAKER LEARNER C_1
2. SECOND RANDOM TRAINING SUBSET, D_2 WITHOUT REPLACEMENT + 50% OF EXAMPLES FROM PREVIOUSLY MISCLASSIFIED TO TRAIN WEAKER LEARNER C_2
3. FIND TRAINING EXAMPLES, D_3 IN TRAINING DATASET, D WHICH $C_1 + C_2$ DISAGREE WITH, TO TRAIN C_3
4. COMBINE C_1, C_2, C_3 VIA MAJORITY VOTING

- LEO BREIMAN → BOOSTING CAN LEAD TO DECREASE IN BIAS + VARIANCE - COMPARED TO BAGGING
↳ IN PRACTICE → AdaBoost known for HIGH VARIANCE → TEND TO OVERFIT TRAINING DATA.
- ADABoost → USES THE COMPLETE TRAINING DATASET TO TRAIN WEAK LEARNERS → TRAINING EXAMPLES REWEIGHTED AFTER EACH ITERATION TO LEARN FROM MISTAKES



1. SET WEIGHT VECTOR, w , TO UNIFORM WEIGHTS, $\sum_i w_i = 1$.
2. FOR j IN m BOOSTING ROUNDS:
 - a. TRAIN WEIGHTED WEAK LEARNER, $C_j = \text{train}(X, y, w)$
 - b. PREDICT CLASS LABELS, $\hat{y} = \text{pred}(C_j, X)$
 - c. COMPUTE WEIGHTED ERROR RATE, $E = w \cdot (\hat{y} \neq y)$
 - d. COMPUTE COEFFICIENT, $\alpha_j = 0.5 \log \frac{1-E}{E}$
 - e. UPDATE WEIGHTS, $w := w \times \exp(-\alpha_j \times \hat{y} \times y)$
 - f. NORMALISE WEIGHTS TO SUM TO 1, $w := w / \sum_i w_i$
3. COMPUTE FINAL PREDICTION, $\hat{y} = (\sum_{j=1}^m (\alpha_j \times \text{pred}(C_j, X)) > 0)$

SEBASTIAN RASCHKA - ALL + PyTorch Book.

- CONSIDERED BAD PRACTICE TO SELECT A MODEL BASED ON THE REPEATED USAGE OF THE TEST DATASET. \rightarrow GENERALISATION PERFORMANCE MAY BE OVEROPTIMISTIC
- ENSEMBLE LEARNING INCREASES COMPUTATIONAL COMPLEXITY COMPARED TO INDIVIDUAL CLASSIFIERS
- \hookrightarrow NO PRACTICE ASK = IS IT WORTH THE LARGEST IMPROVEMENT IN PREDICTIVE PERFORMANCE FOR INCREASED COMPUTATIONAL COST.

BOOSTING VS. STACKING

■ = BOOSTING ■ = STACKING

TRAINING ORDER - SEQUENTIAL TRAINING, PARALLEL TRAINING

CORE GOAL - REDUCE BIAS BY REPEATEDLY CORRECTING MISTAKES, EXPLOIT DIVERSITY - LEARN HOW TO MIX DIFFERENT STRONG MODELS

NEW MODEL FITTING - EMPHASIZES MISCLASSIFIED/LARGE ERROR SAMPLES, EACH BASE MODEL SEES SAME DATA, META-MODEL TRAINED ON PREDICTIONS.

COMBINATION METHOD - FINAL PREDICTIONS = WEIGHTED SUM/LARGE ERROR SAMPLES, META-MODEL OUTPUT

PARALLELISATION - LIMITED (ITERATION DEPENDANT), BASE LEARNERS EASILY PARALLELISABLE

SENSITIVITY - CAN OVERFIT/HOUSE-SENSITIVE IF NOT REGULARISED, CAN OVERFIT IF META-MODEL IS TOO FLEXIBLE \rightarrow DEPENDS ON GOOD CV SETUP.

Typical USE - GRADIENT BOOSTING TREES, XGBOOST, ADABoost, COMPETITION/PRODUCTION 'BLENDER' \rightarrow COMBINE TREES, LINEAR MODELS, SVMS, NEURAL NETS, ETC ...

GRADIENT BOOSTING - TRAINING ENSEMBLE BASED ON LOSS GRADIENTS

- ANOTHER VARIANT OF BOOSTING CONCEPT \rightarrow TRAINING WEAKER LEARNERS
- GRADIENT BOOSTING IS IMPORTANT \rightarrow FORMS BASIS OF POPULAR ML ALSO LIKE XGBOOST \rightarrow WELL KNOWN FOR KAGGLE COMPETITIONS
- GRADIENT BOOST FITS DECISION TREES IN AN ITERATIVE FASHION USING PREDICTION ERRORS
- \hookrightarrow DEEPER THAN DECISION TREE STUMPS + TYPICALLY A MAXIMUM DEPTH OF 3 TO 6. (MAX 8-16 LEAF NODES)
- \hookrightarrow DOES NOT USE PREDICTION ERRORS FOR ASSIGNING SOURCE WEIGHTS \rightarrow USED DIRECTLY TO FORM TARGET VARIABLE FOR FITTING NEXT TREE
- \hookrightarrow USES SAME GLOBAL LEARNING RATE FOR EACH TREE

OUTLINE OF GRADIENT BOOST ALGORITHM

- GRADIENT BOOST IS A GENERAL-PURPOSE SUPERVISED LEARNING METHOD \rightarrow HERE WE'LL LOOK AT CLASSIFICATION (BINARY CLASSIFICATION EXAMPLE)
- BUILDS SERIES OF TREES, WHERE EACH TREE IS FIT ON ERROR (DIFFERENCE BETWEEN LABEL + PREDICTED VALUE)
- EACH ROUND, TREE ENSEMBLE IMPROVES \rightarrow ADJUSTING TREE MORE IN THE RIGHT DIRECTION WITH SMALL UPDATES.
- UPDATES BASED ON LOSS GRADIENT \rightarrow GOT GRADIENT BOOSTING NAME.
- STEP-BY-STEP (GENERAL ALGORITHM):

1. INITIALISE MODEL TO RETURN CONSTANT PREDICTION VALUE, DECISION TREE ROOT NODE (DECISION TREE WITH SINGLE LEAF NODE)

$$F_0(x) = \underset{\hat{y}}{\operatorname{arg\,min}} \sum_{i=1}^n L(y_i, \hat{y}_i)$$

\hat{y} = VALUE RETURNED BY DECISION TREE
 L = LOSS FUNCTION
 n = n TRAINING EXAMPLES

2. FOR EACH TREE $m=1, \dots, M$, IS USER-SPECIFIED TOTAL NUMBER OF TREES

- a. COMPUTE DIFFERENCE BETWEEN PREDICTED VALUE $F(x_i) = \hat{y}_i$ AND CLASS LABEL y_i . SOMETIMES CALLED PSEUDO-RESPONSE OR PSEUDO-RESIDUAL

FORMALLY \rightarrow WRITE PSEUDO-RESIDUAL AS NEGATIVE OF LOSS FUNCTION WITH RESPECT TO PREDICTED VALUES:

$$r_{im} = - \left[\frac{\partial L(y_i, F(x_i))}{\partial F(x_i)} \right] \quad F(x) = \text{PREDICTION OF PREVIOUS TREE } F_{m-1}(x) \quad (\text{FIRST ROUND} \rightarrow \text{CONSTANT VALUE FROM SINGLE NODE TREE})$$

$$F(x) = F_{m-1}(x) + r_{im}$$

- b. FIT A TREE TO PSEUDO-RESIDUAL r_{im} . R_{jm} TO DENOTE $j=1 \dots J_m$ LEAF NODES OF RESULTING TREE IN ITERATION m.

- c. FOR EACH LEAF NODE R_{jm} \rightarrow COMPUTE OUTPUT VALUE:

$$\gamma_{jm} = \underset{\gamma}{\operatorname{arg\,min}} \sum_{x_i \in R_{jm}} L(y_i, F_{m-1}(x_i) + \gamma) \quad Y_{jm} = \text{COMPUTED BY MINIMIZING LOSS FUNCTION}$$

R_{jm} CAN CONTAIN MORE THAN ONE TRAINING EXAMPLE, HENCE SUMMATION

- d. UPDATE MODEL BY ADDING OUTPUT VALUE γ_{jm} TO PREVIOUS TREE:

$$F_m(x) = F_{m-1}(x) + \gamma_{jm}$$

INSTEAD OF ADDING FULL PREDICTED VALUES OF CURRENT TREE γ_{jm} TO PREVIOUS TREE F_{m-1} , SCALE γ_{jm} BY LEARNING RATE η (SMALL STEP \rightarrow TYPICALLY BETWEEN 0.01 - 1)

GRADIENT BOOST FOR CLASSIFICATION

- SINGLE TRAINING EXAMPLE, LOGISTIC LOSS:

$$L_i = -y_i \log p_i + (1-y_i) \log(1-p_i)$$

- LOG(ODDS):

$$\hat{y} = \log(\text{odds}) = \log\left(\frac{p}{1-p}\right)$$

- LOG(ODDS) TO REWRITE EQUATION:

$$L_i = \log(1+e^{\hat{y}}) - y_i \hat{y}$$

- PARTIAL DERIVATIVE OF LOSS FUNCTION WITH RESPECT TO LOG(ODDS), \hat{y} :

$$\frac{\partial L_i}{\partial \hat{y}_i} = \frac{e^{\hat{y}_i}}{1+e^{\hat{y}_i}} - y_i = p_i - y_i$$

- AFTER THESE STEPS - NOW YOU ADD GRADIENT BOOST STEPS:

1. CREATE ROOT NODE \rightarrow MINIMIZES LOGISTIC LOSS, LOSS MINIMIZED IF ROOT NODE RETURNS LOG(ODDS), \hat{y} .

2. FOR EACH TREE $M=1, \dots, M$, IS USER-SPECIFIED TOTAL NUMBER OF TREES

a. CONVERT LOG(ODDS) INTO PROBABILITY USING LOGISTIC FUNCTION (THAT WE USED IN LOGISTIC REGRESSION):

$$p = \frac{1}{1+e^{-\hat{y}}}$$

THEN COMPUTE PSEUDO-RESIDUAL - NEGATIVE PARTIAL DERIVATIVE OF LOSS WITH RESPECT TO LOG(ODDS)

\hookrightarrow DIFFERENCE BETWEEN CLASS LABEL AND PREDICTED PROBABILITY

b. FIT NEW TREE TO PSEUDO-RESIDUALS

c. FOR EACH R_{jm} , COMPUTE VALUE γ_{jm} - MINIMISES LOSS FUNCTION

$$\gamma_{jm} = \arg \min \sum_{x_i \in R_{jm}} L(y_i, F_{m-1}(x_i) + \gamma)$$

$$= \log(1+e^{\hat{y}_i+\gamma}) - y_i(\hat{y}_i+\gamma)$$

$$\text{RESULTS IN: } \gamma_{jm} = \frac{\sum_i y_i - p_i}{\sum_i p_i (1-p_i)} \rightarrow \text{THIS IS ONLY FOR } R_{jm}, \text{ AND NOT FULL TRAINING SET.}$$

d. UPDATE MODEL BY ADDING GAMMA VALUE FROM 2c WITH LEARNING RATE η :

$$F_m(x) = F_{m-1}(x) + \eta \gamma_m$$

- AFTER BOOSTING ALGO COMPLETE \rightarrow PREDICT CLASS LABELS BY THRESHOLDING PROBABILITY VALUES OF FINAL MODEL, $F_m(x)=0.5$, LIKE LOGISTIC REGRESSION.

\hookleftarrow WILL PRODUCE NON-LINEAR DECISION BOUNDARIES UNLIKE

XG BOOST

- GRADIENT BOOSTING IS A SEQUENTIAL PROCESS \rightarrow SLOW TO TRAIN (INEFFICIENT TIME WISE)

\hookrightarrow ALTERNATIVE \Rightarrow XG BOOST

- XGBOOST (EXTREME GRADIENT BOOST) - PROPOSED SEVERAL TRICKS + APPROXIMATIONS \rightarrow SPEED UP TRAINING PROCESS.