

32.1 Ensemble Models

Ensembles:

↳ group of things

(87)

Collection of things

In ML: multiple models used together

$m_1, m_2, m_3, \dots, m_k \rightarrow$ base models

(SVM, logistic, DT, kNN etc.)

Combine these models to create a "more powerful model"

4 types

↳ Bagging (Bootstrapped Aggregation)

↳ Boosting \rightarrow high performing

↳ Stacking \rightarrow v. powerful

↳ Cascading \rightarrow kaggle competition

Key aspect $\vdots \underbrace{m_1, m_2, m_3, \dots, m_k}$

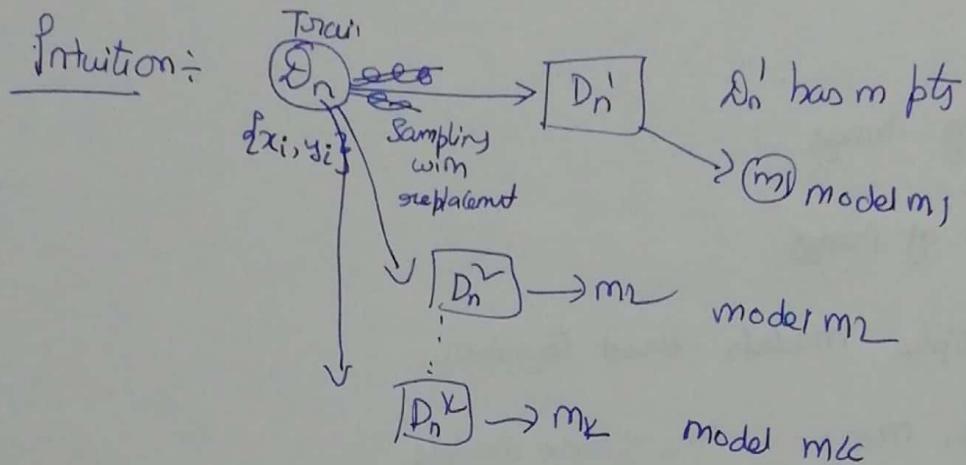
more different the models are the better you can combine them.

of Problem $\vdash (m_i)$ - Expert

32.2

Bootstrapped Aggregation (Bagging)

Bagging

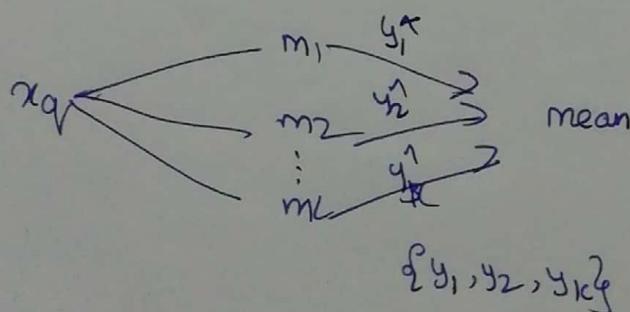
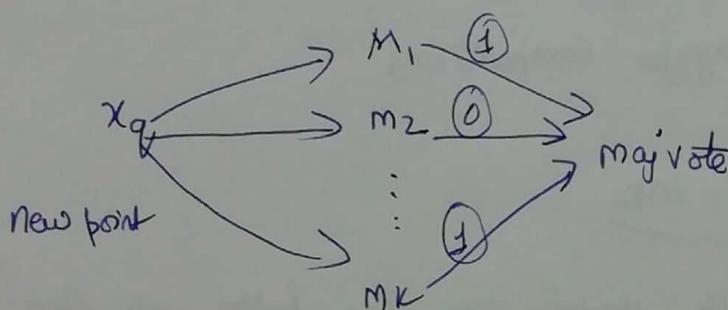


Models are built using D_n' of size m ($m \leq n$)

- \Rightarrow Each model M_i has seen a different subset of data
- \rightarrow Each sample is called bootstrap sample

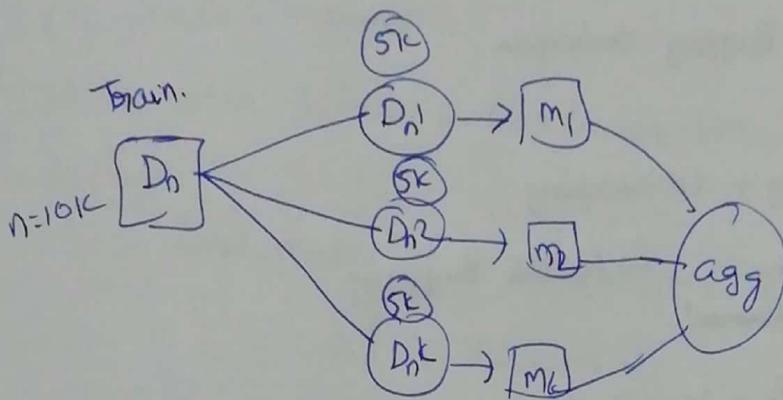
Aggregation: classification \div majority vote

Regression \div mean / median



mean of $\{y_1^*, y_2^*, y_K^*\}$ is our output

Bootstrap Sampling



Variance: How much does a model change with change in training data.

Let's assume we first we have $n=10K$, now we have changed 100 pts in Train. These 100 pts might effect ~~few~~ ^{few} subset of the ~~models~~ sample. It will effect few models not all.

Bagging → Can reduce the variance in the model without impacting the bias

$$\text{Model Error} = \text{Bias}^2 + \text{Var.}$$

bare-model (m_i) → low bias, & high-var model

⊕⊕ Bagging (m_i) → low-bias & reduced variance.

Bagging says take a bunch of low-bias, high-var models

with
DT of depth of 10
↳ high variance
↳ low bias

Combine them using bagging
↓
(low bias, reduced-variance)

32.3 Random Forest & Their Construction

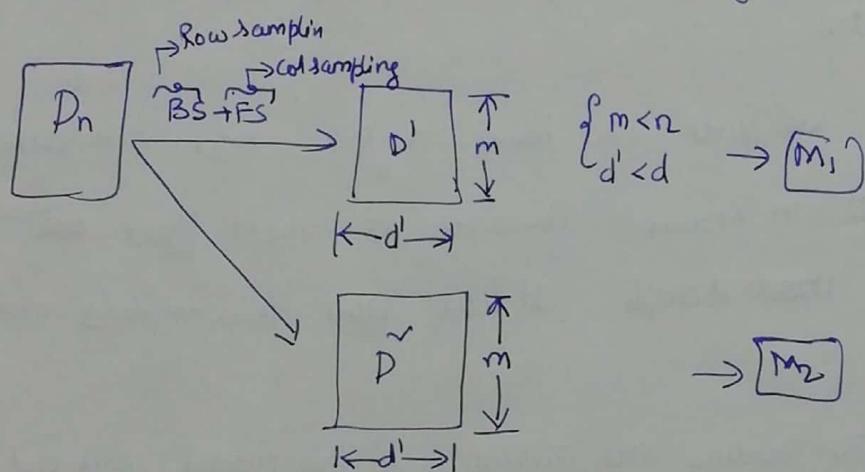
→ RF → most popular Bagging technique

RF : $D^T + \text{Bagging} + \text{Col. sampling}$

base Row sampling \hookrightarrow feature Bagging
with replacement

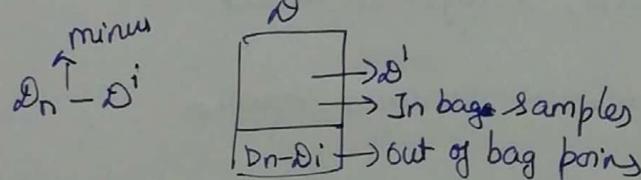
Bootstrap Sampling : Row sampling / Bagging

column sampling : Feature sampling / Bagging



M_i 's : D^T of reasonable depth

$M_i \rightarrow D^i \rightarrow d^i$ columns
m rows



\hookrightarrow we can be used as a CV dataset for model M_i

Summary

RF : D^T base model \rightarrow of reasonable depth (or) fully grown tree

+ Row sampling with replacement

+ col. sampling

+ agg (majority vote) or (mean & median)

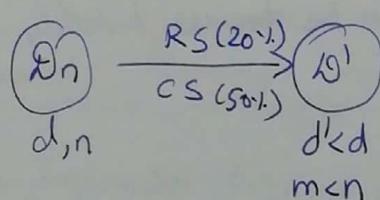
32.4

Bias and Variance Trade off.

RF: reduce variance

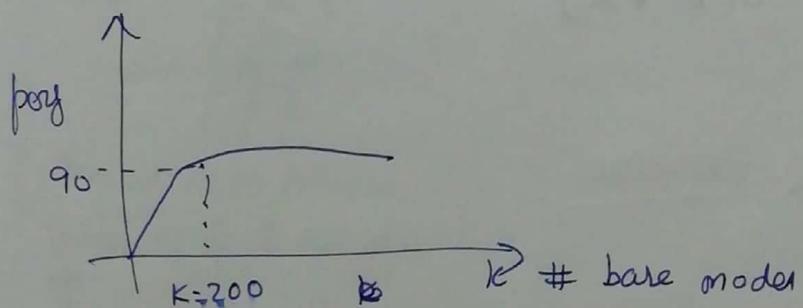
↳ low bias because, base learners (M_i 's) are low biased

↳ (Var)

Final model $M = \text{agg}(M_1, M_2, M_3, \dots, M_k)$ $K \uparrow$; Variance \downarrow $K \downarrow$; Variance \uparrow → we have seen 2000 learners $\frac{d'}{d} = \text{col sampling rate}$ $\frac{m}{n} = \text{Row sampling rate}$

col. S.R \downarrow → low variance model or Variance \downarrow
 row S.R \downarrow

$D_n \rightarrow$ Training data $\rightarrow D_1 - M_1$
 $\quad \quad \quad D_2 \rightarrow M_2$

Hyper parameters $K \rightarrow \# \text{ base learners}$

col. sampling rate

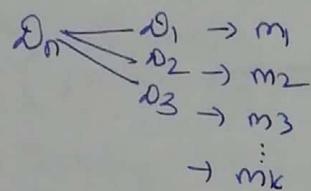
Row sampling rate

32.5

Train & run time Complexity

Imagine we are training RF with K base learners (DT) (fully grown)

Train: $O(\underbrace{n \lg n}_2 * k)$



→ In modern Computing, we will have multicore, Each model can be trained on a single core

If 4 core system, four models can be trained parallelly

→ Trivially parallelizable

→ When we have large amount of data with reasonable # features

Runtime : $O(\text{depth} * k)$
 \downarrow
 large (10 to 20)

Space Complexity : $O(\overset{\text{not}}{\text{DT}} * k)$ *if else*

32.6 Bagging code sample

→ sklearn.ensemble.RandomForestClassifier.

→ n_estimators = K (# of models)

→ max_depth = None

↳ we build DT to deep

→ bootstrap = True

→ oob_score = False

↳ If we make it true

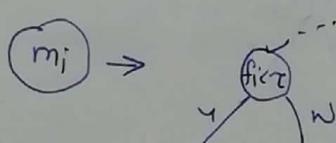
for every model that is build, it will return the oob-score.

32.7 Extremely Randomized Trees

→ sklearn implemented ExtraTreesClassifier and ExtraTreesRegressor

→ In RF we do col.sampling, row sampling and agg

→ Extreme Trees →



RP/DT $\left\{ \begin{array}{l} f_i: \text{real valued values} \\ \text{sort } f_i \\ \begin{array}{c|c} 1 & \leftarrow x \\ 2 & \leftarrow x \\ 3 & \leftarrow x \\ 4 & \leftarrow x \\ \vdots & \leftarrow x \end{array} \end{array} \right.$

Try out all the possible values of f_i to determine γ .

?

→ Instead of trying all the values, what if we try out random sample of possible values to determine γ

→ Extreme Trees: col.sampling, row sampling + agg +

Randomization when selecting γ

→ Randomization is a way to reduce variance

→ Extreme trees better than RF in term of variance

32.8 Random Trees: Cares

Random Forest :-

$DT + RS + CS + agg$

to reduce variance

do not handle large-dim

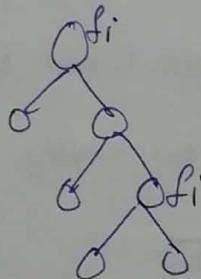
cat-features with many categories

Bias-variance of RF \rightarrow $DT \rightarrow (\text{depth})$
 $\rightarrow \text{R} = \# \text{ of base learners}$
 \hookrightarrow deep trees

② Feature Imp :-

$DT \doteq f_i \doteq$ overall reduction in Entropy $\& I_G$ because
of this feature @ various levels in the DT

RF



RF :- overall reduction in $H(Y)$ & I_G
because of feature f_i @ various levels of each of m trees

32.9 Boosting Iteration

32.9 : Boosting Intuition.

reduce-var.

Bagging $\hat{=}$ high variance, low bias + standardization + aggregation
(CS, RS)

Boosting $\hat{=}$ low variance, high bias + additively combine

$$E_{\text{end}} = \text{bias} + \text{var} + \epsilon$$

reduce ~~var~~ bias

while keeping over variance low

Core idea $\hat{=}$ reduce bias

$$D_{\text{train}} = \{x_i, y_i\}_{i=1}^n$$

Step 0.

$D_{\text{train}} \rightarrow$ Build a model (M_0)

\hookrightarrow high bias \leftrightarrow low variance.

Eg: DT which is shallow/less depth
(high variance, low bias.)

high-bias \rightarrow large training error

② $D_{\text{train}} \rightarrow M_0$

$$(x_i, y_i)_{i=1}^n \quad y = h_0(x)$$

③

$$y_i - h_0(x_i) = \epsilon_{0,i} \quad \leftarrow \text{Eq. Error}$$

\leftarrow log Err

\leftarrow hinge Err.

$$\{x_i, y_i, \epsilon_{0,i}\}_{i=1}^n$$

Step + 1 : Instead of taking D_{train} we will take $\{x_i, \epsilon_{0,i}\}_{i=1}^n$

$$M_1 \leftarrow \{x_i, \epsilon_{0,i}\}_{i=1}^n$$

here $\epsilon_{0,i}$ is ~~what~~ the error which we calculated in step 0

$$h_1(x) = h_1(x_i) = \epsilon_{0,i} - y_i - h_0(x_i) \quad \} \text{ it tries to predict } \epsilon_{\text{end}}.$$

$F_1(x)$: model at $\overset{\text{no end of}}{\underset{\curvearrowleft}{\text{stage 1}}}$

$$F_1(x) = \alpha_0 h_0(x) + \alpha_1 h_1(x) \quad \} \div \text{we can think as weighted sum of 2-base models}$$

$$\textcircled{2} \quad \{x_i, y_{\text{end}_i}\} \rightarrow \textcircled{m_2} \quad h_2(x)$$

$$\uparrow y_i - F_1(x_i) \rightarrow \text{here } F_1(x) = \alpha_0 h_0(x) + \alpha_1 h_1(x)$$

$$F_2(x) = \alpha_0 h_0(x) + \alpha_2 h_2(x) + \alpha_1 h_1(x)$$

At the end of Stage $K \vdash$

$$F_K(x) = \sum_{i=0}^K \alpha_i h_i(x)$$

Trained to fit the residual error @
End of the previous stage.

additive weighted model

$$h_i(x) \leftarrow \{x_i, y_{\text{end}_i}\} \quad \overbrace{y_i - F_{i-1}(x)}^{\text{Residual error} @ \text{End of stage } (i-1)}$$

$$F_K(x) = \sum_{i=0}^K \alpha_i h_i(x)$$

→ Since we are calculating consecutive models based on the error from the previous stage.

Let $K=10$

→ it ends up having a low residual error.

→ At every stage we are reducing error and bias in model is also reduces

There are many boosting Techniques, we will see

→ Gradient Boosting DT

→ Ada Boosting

↳ Face Detection

32.10 Residuals, loss functions, gradients

At the end we got

$$F_K(x) = \sum_{i=1}^K \alpha_i h_i(x)$$

Residual $\rightarrow \underbrace{y_i - F_K(x)}_{\text{Residual}}$
 @ End of Stage K

If we want to train a model M_{K+1} , we will train using

$$M_{K+1} \leftarrow \{x_i, \text{Resid}_i\}$$

loss-minimization

In Logistic Regression we minimized logistic loss

In Linear

"

Square loss

SVM

"

Hinge loss.

$$\underbrace{L(y_i, F_K(x_i))}_{\text{sq. loss}} = (y_i - F_K(x_i))^2$$

Let's call
 $F_K(x_i) = z_i$

$$\frac{\partial L}{\partial F_K(x_i)} = \frac{\partial L}{\partial z_i} = \frac{\partial}{\partial z_i} (y_i - z_i)^2$$

$$= (-1) * 2 * (y_i - z_i)$$

$$\frac{\partial L}{\partial z_i} = -2(y_i - z_i)$$

$\frac{\partial L}{\partial F_K(x_i)}$ = $2(y_i - F_K(x_i))$

negative gradient
 negative derivative
 of loss function
 w.r.t $F_K(x_i)$

$\underbrace{2(y_i - F_K(x_i))}_{\text{Residual}}$

* Negative gradient/derivative \approx Residual.

Reg. gradient is a pseudo residual (i) Proj. gradient
↳ false

$$h_i(x) \leftarrow x_i, \underbrace{e_{\text{proj},i}}_{\downarrow} = \underbrace{y_i - F_{i-1}(x)}_{\downarrow}$$

Instead of having a $e_{\text{proj},i}$ at any point i as a residual off point,
what if we replace $e_{\text{proj},i}$ with pseudo residual

$$\downarrow, -\frac{\partial L}{\partial F_{i-1}(x)}$$

→ Why we have to train on pseudo residuals why not only
graduals

Let us have any loss function which is differentiable

RF → min hinge loss

↳ we cannot perform using hinge loss, RF uses Entropy loss

GB → min (any loss) function
↳ which is differentiable
Super powerfull

At any stage :

$$F_{i-1}(x)$$

$$(x_i, e_{\text{proj},i}) \rightarrow (m_i) \uparrow \quad h_i(x)$$

pseudo residual

$$-\frac{\partial L}{\partial F_{i-1}(x)}$$

GB : If we have base models with high bias → we will combine
them using pseudo residuals $\left(\frac{-\partial L}{\partial F_{i-1}(x)} \right)$

32:11 Gradient Boosting

Input: Training set $\{(x_i, y_i)\}_{i=1}^n$, a differentiable loss function

$L(y, F(x))$, number of iterations M

↑ Logistic Loss

↳ Sq. Loss

↳ Hinge Loss

Algorithm

1. Initialize model with a constant value:

$$h_0(x) = F_0(x) = \underset{\gamma}{\operatorname{argmin}} \sum_{i=1}^n L(y_i, \gamma)$$

Find a constant γ which minimize my loss y_i over Y

If we are doing regression problem, we will use squared loss

$$\underset{\gamma}{\operatorname{argmin}} \sum_{i=1}^n (y_i - \gamma)^2 \quad \gamma = \text{mean}(y_i)$$

2. For $m = 1$ to $M \rightarrow h_1(x), h_2(x), h_3(x) \dots h_m(x)$

1. Compute so-called pseudo-residuals:

$$\gamma_{im} = - \left[\frac{\partial L(y_i, F(x_i))}{\partial F(x_i)} \right] \quad \text{for } i = 1, \dots, n$$

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

2. Fit a base learner (e.g. tree), $h_m(x)$ to pseudo-residuals.

i.e., train it using the training set $\{(x_i, \gamma_{im})\}_{i=1}^n$

$$F_m(x) = h_0(x) + \gamma_1 h_1(x) + \gamma_2 h_2(x) + \dots + \gamma_m h_m(x)$$

3. Compute multiple γ_m by solving the following one-dimensional optimization problem.

$$\gamma_m = \underset{\gamma}{\operatorname{argmin}} \sum_{i=1}^n L(y_i, F_{m-1}(x_i) + \gamma h_m(x_i))$$

4. update the model

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

32.12 Regularization by shrinkage

$$F_m(x) = h_0(x) + \sum_{m=1}^M \gamma_m h_m(x)$$

$m = \# \text{ of base model}$ \rightarrow Cross Validation

as $m \uparrow = \text{overfit} \uparrow = \text{Var} \uparrow$
 $\hookrightarrow \text{bias} \downarrow$

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

\rightarrow In addition to what ever weightage we found using gradient boosting, we are going to produce that weight by multiplying with a constant " ϑ "

\rightarrow A ϑ ^{increase} ~~reduce~~ chance of overfitting \uparrow , variance \uparrow

$\vartheta \uparrow \rightarrow \text{variance} \uparrow$

$\vartheta \downarrow \rightarrow \text{variance} \downarrow$

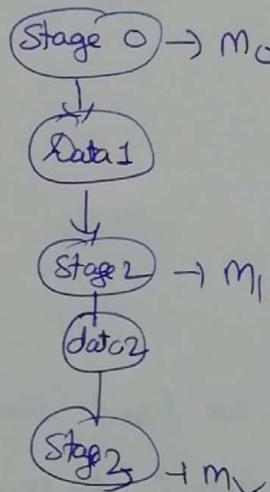
$$\vartheta: [0 \text{ to } 1]$$

\rightarrow PBL will generally overfit using GBT

32.13 Train & Run time Complexity

Train : $O(n \lg n + m)$ $m = \# \text{ bare-learners}$

→ GB DT is not easy to parallelize



Runtime

$$\text{GBDT} = O(\text{depth} * m)$$

↑
Small in GBDT

Space Complexity

↳ Space it would take to store each tree & γ_m 's for each bare learner

- Good for low latency as the depth of the trees are very less
- We can choose any loss function.

Analytics Vidhya

The accuracy of a predictive model can be boosted in two ways.
Either by Embracing Feature Engineering & by applying boosting algorithms straight away.

There are multiple boosting algorithms

Gradient Boosting

XG Boosting

Ada Boost

Gentle Boost etc.,

Square Error (<http://blog.kaggle.com/2017/01/23/a-kaggle-master-explains-gradient-boosting/>)

Age	F0	PseudoResidual0	h0	gamma0	F1	PseudoResidual1	h1	gamma1	F2
13	40.33	-27.33	-21.08	1	19.25	-6.25	-3.567	1	15.68
14	"	-26.33	-21.08	1	19.25	-5.25	-3.567	1	15.68
15	"	-25.33	-21.08	1	19.25	-4.25	"	1	15.68
25	"	-15.33	16.87	1	57.2	-32.2	"	1	53.63
35	"	-5.333	-21.08	1	19.25	15.75	"	1	15.68
49	"	8.667	16.87	1	57.2	-8.2	7.133	1	64.33
68	"	27.67	16.87	1	57.2	10.8	-3.567	1	53.63
71	"	30.67	16.87	1	57.2	13.8	7.133	1	64.33
73	"	32.67	16.87	1	57.2	15.8	7.133	1	"

Root
[-27.33, -26.33, -25.33 ... 30.67, 32.67]

LikesGarden == F
{-27.3, -26.3 ...}

LikesGarden == T
{-15.3, 8.7}

PlayVideoGame == F
{-8.2, 13.8, 15.87}
PlayVideoGame == T
{-6.2, 5.2, -4.2, -32.2, 15.8, 10.87}

32.14 XGBoost: Boosting + Randomization

Sklearn: GBDT + $\gamma \cdot S$

XGBoost: GBDT + $\gamma \cdot S + c \cdot S$

→ pip3 install XGBoost (81) sudo pip3 install XGBoost

<https://github.com/dmlc/xgboost/blob/master/demo/guide-python/sklearn.Example.py>

32.15 Ada Boost: Geometric Intuition

Adaptive Boosting

→ Typically used in Image processing, face detection.

	+	+	+	-
+	-	-	-	-
+	-	-	-	-

The first round:

h_1	-	-	-	-
+	⊕	⊕	⊕	-
+	-	-	-	-
-	-	-	-	-

$$\epsilon_1 = 0.30$$

$$\alpha_1 = 0.42$$

D_2	-	-	-	-
+	+	+	+	-
+	-	-	-	-

The second round:

$h_1(x)$	-	-	-	-
-	-	-	-	-
-	-	-	-	-

A_2	-	-	-	-
+	+	+	+	-
+	-	-	-	-
+	-	-	-	-

$$\epsilon_2 = 0.71$$

$$\alpha_2 = 0.68$$

D_3	-	-	-	-
+	+	+	+	-
+	-	-	-	-

The third round:

+	+	+	-
+	-	-	
-	-	-	

+	+	+	-
+	-	-	-
-	-	-	-

+	+	+	-
+	-	-	-
-	-	-	-

$$\sum \varepsilon_3 = 0.10$$

$$A_3 = 0.92$$

The final classifier

$$\alpha_1 h_1 + \alpha_2 h_2 + \alpha_3 h_3 = F_3(x)$$

$$h_{\text{final}} = \text{Sign} \left(0.42 \begin{array}{|c|} \hline \square \\ \hline \end{array} + 0.65 \begin{array}{|c|} \hline \square \\ \hline \end{array} + 0.92 \begin{array}{|c|} \hline \square \\ \hline \end{array} \right)$$
$$= \begin{array}{|c|} \hline \begin{array}{|c|c|c|} \hline & + & - \\ \hline + & - & - \\ \hline + & - & - \\ \hline \end{array} \\ \hline \end{array}$$

$$\alpha_i \rightarrow \gamma_i$$

In the case of GBDT it uses residuals (computed from the negative of the gradient from loss function)

→ misclassified points are given more weightage

↑ At every stage we are adapting ~~to~~ the errors made in the previous stage.

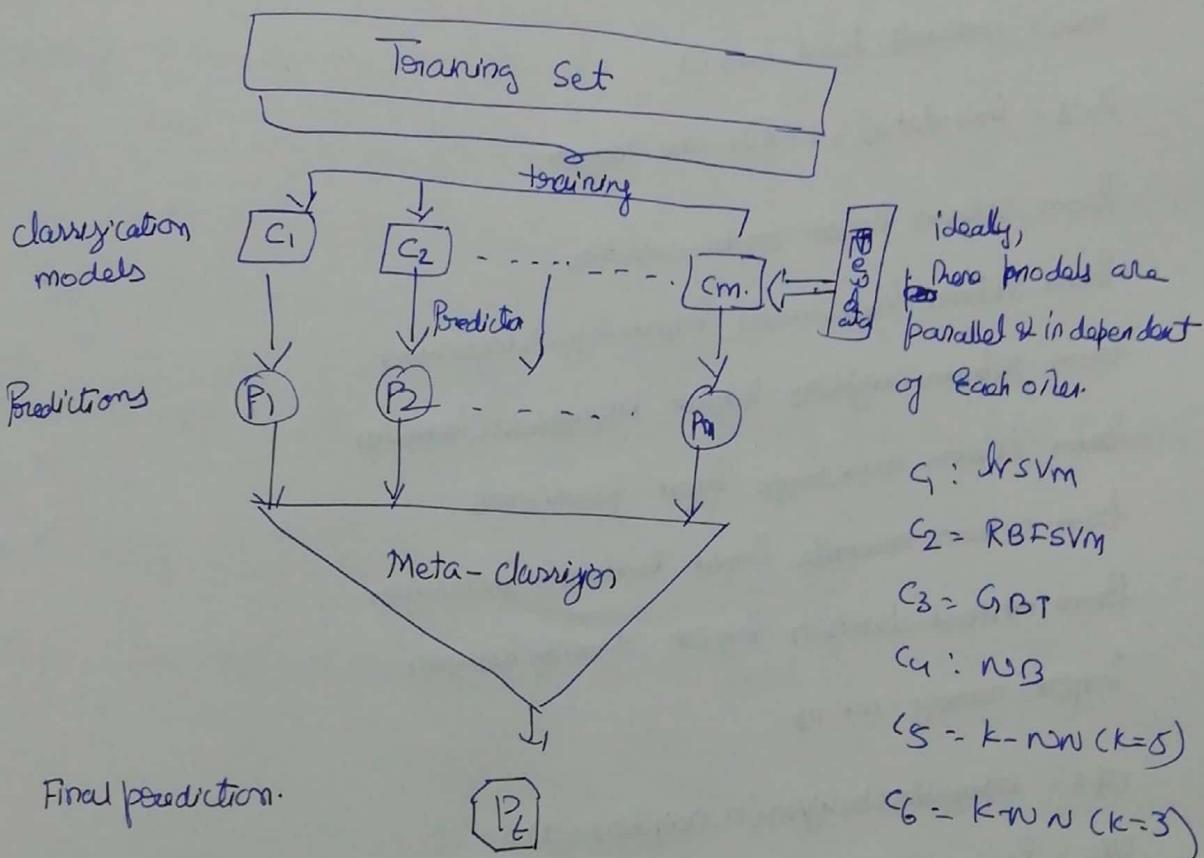
32.16

Stacked models

→ One of the form of Ensembling

→ https://rasbt.github.io/mlxtend/user-guide/classifier/stacking_classifier/

From mlxtend.classifier import StackingClassifier

Algorithm

Input: Training Data $D = \{x_i, y_i\}_{i=1}^m$ ($x_i \in \mathbb{R}^d$, $y_i \in \mathcal{Y}$)

Output: An Ensemble classifier H

- 1: Step 1 : learn first-level classifiers.
- 2: for $t \leftarrow 1$ to T do
- 3: learn a base classifier h_t based on D $\{h_1, h_2, h_3, \dots, h_T\}$
(with good bias-var trade off)
- 4: End for
- 5: Step 2 : Construct new data set from D
- 6: for $i \leftarrow 1$ to m do
- 7: construct a new dataset that contains $\{x_i', y_i\}$ $x_i' = [h_1(x), h_2(x), \dots, h_T(x)]$
- 8: End for

9: Steps: Learn a second-level function

16: Learn a new classifier b' based on the Newly constructed data set

11: $\text{return } H(x) = h'(h_1(x), h_2(x), \dots, h_7(x))$

Code : Pip3 install mixfend

```
from sklearn import datasets
```

```
iris = datasets.load_iris()
```

```
X, y = iris.data[:, 1:3], iris.target
```

```
from sklearn import model_selection
```

from sklearn.linear_model import LogisticRegression

from Sklearn. neighbors import KNeighborsClassifier

```
from sklearn.naive_bayes import GaussianNB
```

```
from sklearn.ensemble import RandomForestClassifier
```

from mixed-classifiers import StackingClassifiers

Import numpy as np

Clf1 = knoeighborsClassifier (n_neighbors = 1)

clf2 = Random Forest classifier (random_state = 1)

clf3 = GaussianNB()

lr = Logistic Regression()

$Sclf = \text{Stacking classifiers}$ ($\text{classifiers} = [clf_1, clf_2, \dots, clf_n]$)

meta-danijen = dr)

fout('3-fold Cross Validation: \n')

for clf, label in ([clf1, clf2, clf3, scf], ['know',
'RandomForest',
'Naive Baye',
'Stacking classifier']) :

`scores = model_selection.cross_val_score(clf, X, y, cv=3,`

`Point('Accuracy': ~0.2f (+/- ~0.2f) [~s],
scoring='accuracy')
* (scores.mean(), scores.std(), label))`

O/P

Accuracy : 0.91 (+/- 0.01) [KNN]

" : 0.91 (+/- 0.06) [Random Forest]

" : 0.92 (+/- 0.03) [Naive Bayes]

: 0.95 (+/- 0.03) [Stacking classifier]

→ Here stacking classifier has the highest accuracy over the bare learners

→ We can use probabilities?

`clf = StackingClassifier(classifiers=[clf1, clf2, clf3],`

`use_probas=True,`

`average_probas=False,`

`meta_classifier=clf4)`

① Key points

① Each of the bare models $h_1(x), h_2(x), h_3(x)$, more different the models are, the better it is

② Kaggle: stacking is generally used

③ Real world: cost of stacking is actually lot

Run/Evaluation time ↑

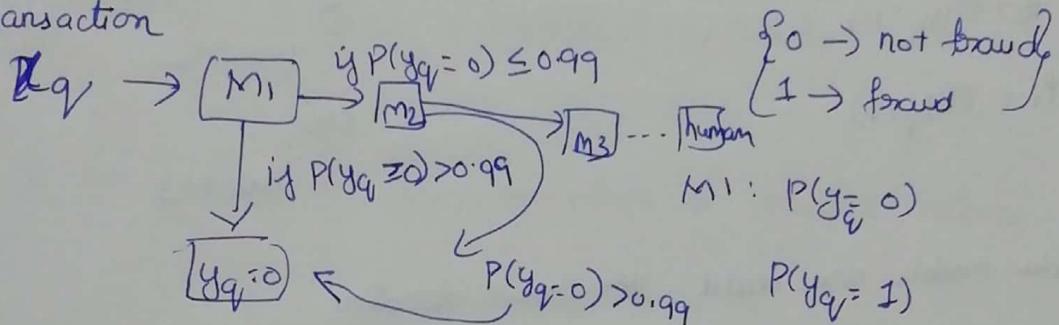
Training time ↑

In real world stacking being used less often than boosting & bagging (RF) (GB)

32.17 Cascading classifiers

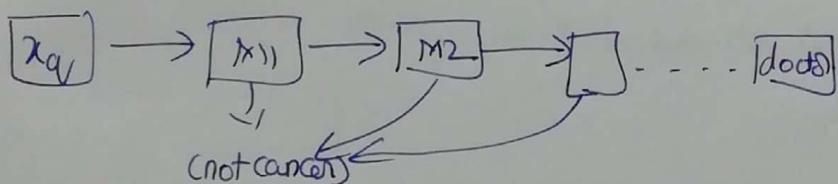
- we have seen Bagging, boosting, stacking
- Eg: Predicts Credit Card transaction is fraudulent or not.

Transaction

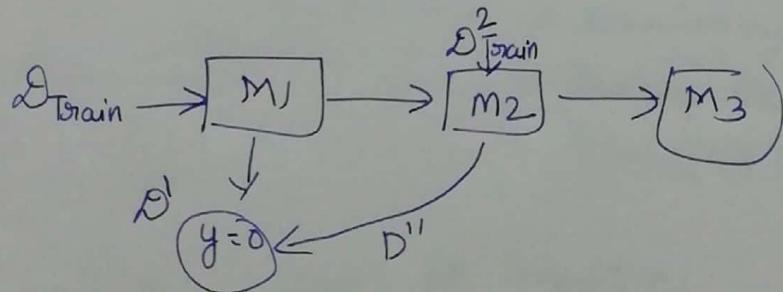


Cascade-model:

Typically used when the cost of making a mistake is high



- we should be slightly careful.



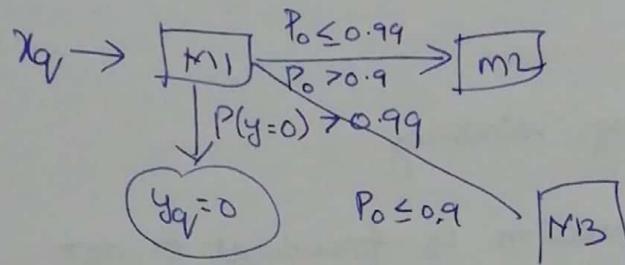
$D_{\text{train}} \rightarrow M_1$

$D_{\text{train}}^2 \rightarrow M_2$

$D_{\text{train}}^3 = D_{\text{train}} - D^1 \text{ points}$

$D_{\text{train}}^3 \rightarrow M_2$

$D_{\text{train}}^3 \rightarrow D_{\text{train}}^2 - D^1 \text{ points}$



$$P_0 = P(Y_{q_1} = 0)$$

$$P_1 = P(Y_{q_1} = 1)$$

If $P_0 > 0.99$

$$Y_{q_1} = 0$$

Else

If $P_0 > 0.9 \text{ & } P_0 \leq 0.99$

use M_2

Else

use M_3

→ Cascade models are build for fraud detection

→ In some case studies 20-levels are used

team will work for your
medical / fraud.

32.18

Kaggle Vs Real World

Ensembles: Kaggle, Competition.

① Kaggle cares about only one metric

Competition → 1st (Log-loss, AUC, F1-score)
→ 2nd

In the real world we have look at multiple metrics

→ primary metric
↳ secondary metric.

Problem: Improve Sales of Amazon.

↳ which metric we have reduce.

② { Business metric → Translate it to one or more
ML-metrics.

In Kaggle

② Very Complex Ensembles → to Improve one metric

↑
impractical → low-latency
↳ Training time
↳ Interpretability

③ Kaggle: great for Competitions

↳ data cleaning
data preprocessing
feature engineering

Random Forest

i) a) Max_features :

↳ Auto / None
↳ Sq, rt
↳ log₂
↳ o.2

b) n_estimators : higher the better

c) min_sample_leaf :

Smaller leaf makes the model more prone to capturing noise in train data.

Generally I prefer a minimum leaf size of more than 50.

2. Features which will make the model training earlier :

n_jobs : This parameter tells the engine how many processes is it allowed to use

'-' : no restriction

1 : use 1 one process

Random state

oob-score : This is a random forest CV method. It is similar to leave one out validation technique, however this is so much faster.

Gradient Boosting

Tree Specific Parameters

→ `min_samples_split`

↳ used to control overfitting

↳ higher values prevent a model from learning relations which might be highly specific to the particular sample selected for a tree.

↳ To high will lead to underfitting.

→ `min_samples_leaf` / `minimum_weight_fraction_leaf`:

↳ Control overfitting

↳ Generally lower values should be chosen for imbalanced class problems because the regions in which the minority class will be in majority will be very small.

→ ~~→~~ `max_depth` / `max_leaf_nodes`

↳ Control overfitting

↳ higher depth will allow model to learn relations very specific to a particular sample

↳ should be tuned using CV

→ ↳ If `max_leaf_nodes` are defined, it will ignore `max_depth`.

→ `max_features`

↳ `sqrt`

↳ `log2`

↳ `auto`

↳ `0.3 0.4`

↳ higher values lead to overfitting

Boosting Parameters

→ Learning rate



- * This determines the impact of each tree on the final outcome
- * GBM works by starting with an initial estimate which is updated using the output of each tree.
- * The learning parameter controls the magnitude of this change in the estimate.
- * Lower values are preferred as they make the model robust to the specific characteristics of tree and thus allowing it to generalize ~~the~~.
- * Lower values would need higher number of trees. Computationally expensive.

n_estimators

- * The # of sequential trees to be modeled
- * Though the GBM is fairly robust at higher # of trees but it can still overfit a point
- * CX should be used for tuning

sub sample

- * The fraction of observations to be selected for each tree. Selection is done by random sampling
- * Typical values ~0.8 generally work fine but can be fine tuned further

Fix learning rate & number of estimators for tuning tree based parameters

min_samples_split : 500 : Should be no. 5 to 1% of total values. Since this is a imbalanced class problem, we'll take a small value from the range.

min_samples_leaf : 50 : Can be selected based on intuition

max_depth : 8 : Should be chosen (5-8) based on # of observations & predictions. This has 87K rows & 49 columns.

max_features : $\sqrt{49}$

subsample : 0.8 Commonly used start value.

XG Boosting

Boosting Parameters

- eta
- * Analogous to learning rate in GBM
 - * Makes the model more robust by shrinking the weights on each step
 - * Typically final values to be used: $0.01 \rightarrow 0.2$

min_child_weight

- * Defines the minimum sum of weights of all observations required in a child
- * This is similar to minchild_leaf in GBM but not exactly. This refers to min "sum of weights" of observations while GBM has min "number of observation"
 - ↑ overfitting
 - ↑ underfitting (use CV)

max_depth \vdots / max_leaf_nodes

→ Typically 3 to 10

max_leaf_nodes \vdots / max_depth

→ maximum number of terminal nodes (or leaves) in a tree.
→

gamma

- A node is split only when resulting split gives a positive reduction in the loss function.
- Gamma specifies minimum loss reduction required to make a split.
- Make algo conservative
- Can vary depending on the loss function.

Subsample

Typically : $0.5 - 1$

Col sample - by tree

- Similar to max_features in GBM

typically : $0.5 - 1$

lambda

→ L2 regularisation term on weights (Ridge Regression)

→ This used to handle the regularisation part of XGBoost.

Though many data scientists don't use it often, it should be explored to reduce overfitting

alpha

→ L1 regularisation term on weight (Lasso regression)

→ Can be use in case of very high dimensionality so that the algorithm runs faster when implemented

Scale pos_weight

A value greater than 0 should be used in case of high class

imbalance as it helps in faster convergence.

Learning Task Parameters:

These parameters are used to define the optimization objective & the metric to be calculated at each step.

1. Objective [default = Reg: linear]

This defines the loss function to be minimized. Mostly used values are

- * binary: logistic : logistic regression, returns predicted probability (not class)
- * multi: softmax : multiclass classification using the softmax objective, returns predicted class (not probability)
- * multi: softprob : same as softmax, but returns predicted probability of each data point belonging to each class.

2. eval-metric : [default according to objective]

- The metric to be used for validation data
- The default values are rmse for regression & err for classification
- Typical values are
 - rmse
 - mae
 - logloss : negative log-likelihood
 - err : Binary classification error rate (0.5 threshold)
 - merror : Multiclass classification error rate
 - mlogloss : || logloss
 - auc : area under curve.

3. seed

- * The random number seed.
- * Can be used for generating reproducible results and also for parameter tuning

Pipe line

```
from sklearn.feature_selection import ColumnSelector
```

```
from sklearn.pipeline import make_pipeline
```

```
pipe1 = make_pipeline (ColumnSelector (cols = (0, 2)), LogisticRegression ())
```

```
pipe2 = " " ( " " (cols = (1, 2, 3)), " " )
```

```
scf = StackingClassifier (classifiers = [pipe1, pipe2],
```

```
meta_classifier = LogisticRegression ())
```

```
scf.fit (X, y)
```