

Ensemble Models

Multiple models used together.

$\{M_1, M_2, M_3, M_4, \dots, M_k\}$ $\xrightarrow{\text{combine}}$ More powerful Models.

4 types of ensembles

① Bagging (Bootstrapped Aggregation).

② Boosting

③ Stacking

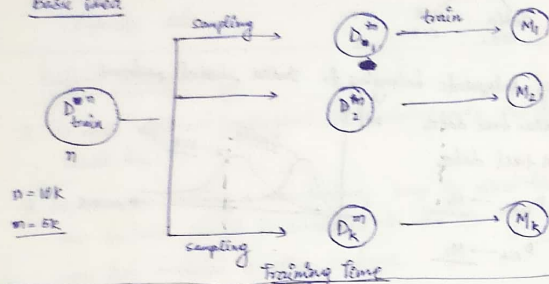
④ Cascading

Mostly used in Kaggle as well as Real world scenarios.

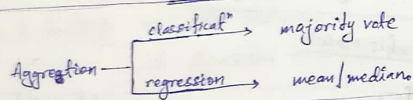
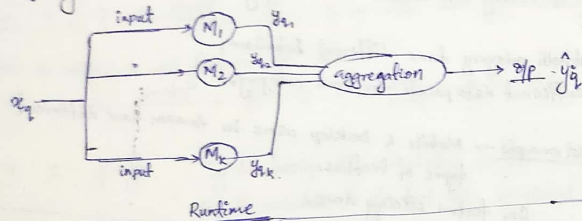
Key Aspect \rightarrow More different the models are, the better & powerful it becomes.

Bagging (Bootstrap Aggregation)

Basic Idea



* sampling done with replacement.



variance \rightarrow degree of change in model with change in data.

- * Say we changed 100 pts. of 10K.
- * Only a subset of the dataset will get affected.
- * That also to very little extent.
- * So, overall model doesn't change much.

Bagging \rightarrow can reduce variance in a model without impacting the bias.

$$\text{Model error} = \text{Bias}^2 + \text{Variance}$$

If the base models have low bias + high variance, then bagging is just the perfect ensemble.

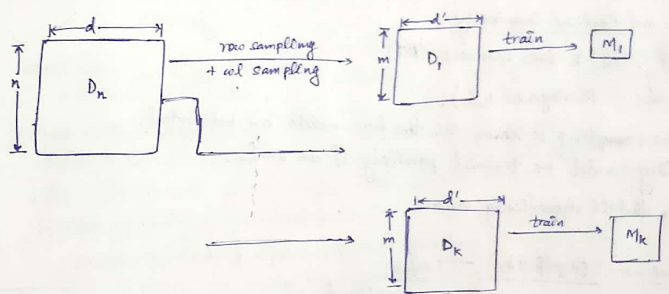
Bagging (MIs) \rightarrow low-bias; reduced variance.

E.g. \rightarrow DT of good depth \rightarrow low bias & high variance.
 \downarrow bagging \rightarrow Random forests

Random forests

Random bootstrap sampling
 \rightarrow collection of decision trees

Random Forest \rightarrow DT (base model) + Bagging + column-sampling \rightarrow feature bagging.



* Sampled datasets now have different sets of features each.
 MIs \rightarrow Decision trees of reasonable depth \rightarrow (low bias + high variance).

Model M_i uses D_i for training.

So, $(D_n - D_i)$ used as cross-validation set.
 \rightarrow Out of Bag (OOB) samples.

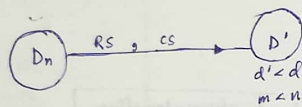
RF \rightarrow very powerful.

Bias-variance trade-off

$$M = \text{aggregat}^n(m_1, m_2, \dots, m_k)$$

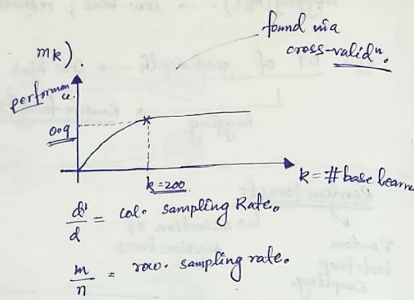
* $K \uparrow$, variance \downarrow .

$$\text{bias}(M) \approx \text{bias}(m_i)$$



low cal. SR \downarrow \rightarrow variance \downarrow .
row SR \downarrow \rightarrow variance \downarrow .

* Usually people fix CSR & RSR and keep increasing no of models to reduce variance.



Train and Runtime Complexity

Ex \rightarrow RF with k -base learners (DT)

train - $O(n \log n \times d \times k)$.

* Once sampling is done, as the base models are independent, they could be trained parallelly if we want.
So, that'll dramatically reduce.

Runtime $\rightarrow O(\text{depth} \times k)$ - fast

This can also be parallelized.

Space complexity $\rightarrow O(DT \times k)$ \leftarrow convenient as $DT \Rightarrow$ if else condⁿ.

Extremely Randomized Tree

Only change for RF \rightarrow while handling real value features, instead of trying each & every value after sorting, ~~try~~ for threshold, try from a random sample of numbers. to make a bit faster & introduce some randomization.

Boosting

Bagging - high variance, low bias + randomisation + aggregation (CS, RS)

Boosting - low-variance, high bias + additively combine. reduce bias while keeping var low

core idea - reduces bias.

Step 0: $D_{\text{train}} \rightarrow M_0 \rightarrow$ (high bias) i.e. high train error.
 $\{x_i, y_i\}_{i=1}^n$ $y_i = h_0(x_i)$ $\text{loss error}_i = y_i - F_0(x_i)$
 $F_0(x) = h_0(x)$
 \hookrightarrow Model at end of step 0.

for any step k :

$D \rightarrow M_k$
 $\{x_i, y_i\}_{i=1}^n$ $y_i = h_k(x_i)$
where $e_i = y_i - F_{k-1}(x_i)$ $F_k(x) = \sum_{j=0}^k \alpha_j h_j(x)$ \rightarrow Model at end of stage k .

* since at each stage we are training the base model on the residual error of prev stage, it ends up having low residual error.

examples of boosting algo \rightarrow Gradient Boosting, AdaBoost, etc.

Residuals, loss functions and gradients

Residual at end of stage $k \Rightarrow e_i = y_i - F_k(x_i)$. Train model M_{k+1} by fitting its function h_{k+1} to $\{x_i, e_i\}_{i=1}^n$.

* Assuming regression, we define loss at end of stage k as

$$L(y_i, F_k(x_i)) = (y_i - F_k(x_i))^2$$

$$\frac{\partial L}{\partial F_k(x_i)} = 2(y_i - F_k(x_i)) \times (-1)$$

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

* negative-gradient is pseudo-residual
 \hookrightarrow Gradient Boosting utilises this.

So, what gradient boosting does is

$$D \xrightarrow{\quad} M_K$$
$$\{x_i, e_i\}_{i=1}^n \quad e_i = h_K(x_i).$$

where e_i = pseudo-residual i.e. $-\frac{\partial L}{\partial F_{K-1}(x)}$ * This allows us to use any differentiable loss function.

So, we can minimize hinge loss using ~~random forests~~ decision trees.

* So, Gradient Boosted Decision Trees (GBDT) tend to perform better than random forests.

Refer \rightarrow Gradient Boosting (wikipedia article) for exact algo.

Input: Training set $\{(x_i, y_i)\}_{i=1}^n$, a differentiable loss function $L(y, F(x))$, number of iterations M (no. of base models).

Algorithm:-

1. Initialize model with constant value:-

$$F_0(x) = \underset{s}{\operatorname{argmin}} \sum_{i=1}^n L(y_i, s). \quad \text{If } L \text{ is sq-loss } s = \operatorname{mean}(y_i).$$

2. For $m=1$ to M :

1. compute so called pseudo-residuals

$$r_{im} = - \left[\frac{\partial L(y_i, F(x_i))}{\partial F(x_i)} \right] \quad \text{for } i=1, 2, \dots, n.$$
$$F(x) = F_{m-1}(x).$$

2. Fit a base learner (e.g. tree) $h_m(x)$ to pseudo-residuals, i.e. train it using the training set $\{(x_i, r_{im})\}_{i=1}^n$

3. Compute multiplier δ_m by solving the following 1-dimension optimization problem

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

4. update the model:-

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

3. Output $F_M(x)$

Regularization and Shrinkage

At the end of the algorithm, we'll get

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

M : # base models (hyperparameter).

$M \uparrow \Rightarrow \text{overfit} \uparrow \Rightarrow \text{variance} \uparrow$.

Optimal 'M' could be found by cross validation and bruteforce or RandomSearch.

Alternatively Shrinkage

$$F_M(x) = F_{M-1}(x) + \lambda \sum_{m=1}^M h_m(x) \quad 0 < \lambda \leq 1$$

$\lambda \rightarrow \text{learning rate}$

$\lambda \uparrow \Rightarrow \text{overfit} \uparrow \Rightarrow \text{variance} \uparrow$.

$\lambda \downarrow \Rightarrow \text{variance} \downarrow$.

So, we have 2 hyperparameters: M & λ .

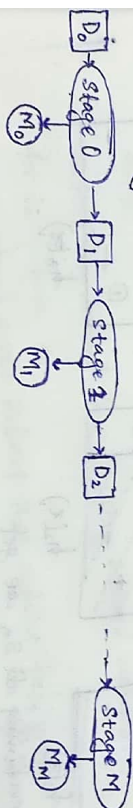
CV with bruteforce / RandomSearch (M, λ).

Train & Runtime Complexity (GBDT)

Train: $O(n \log n * d * M)$.

M : # base learners.

unlike RF, GBDT not parallelizable.
boosting is a serial algorithm



Runtime: $O(\text{depth} * M)$

Space: $O(\text{each tree} * M + \sum \lambda_m)$

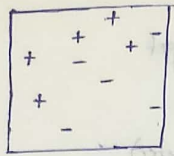
Adaboost

Special case of Gradient Boosting.

Application: face detection.

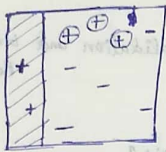
Core idea: Increase # of misclassified points.

See - universality of perceptrons example.



original training data
lets say we train DT of
depth=1 (high bias, low variance)

1st stage



$h_1(x)$



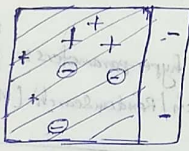
$h_2(x)$

exponentially
increased no.
of misclassified
pts.

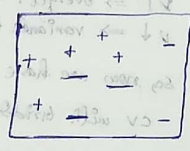
2nd stage



$h_1(x)$

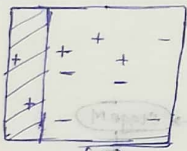


$h_2(x)$



$h_3(x)$

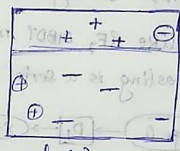
3rd stage



$h_1(x)$



$h_2(x)$



$h_3(x)$

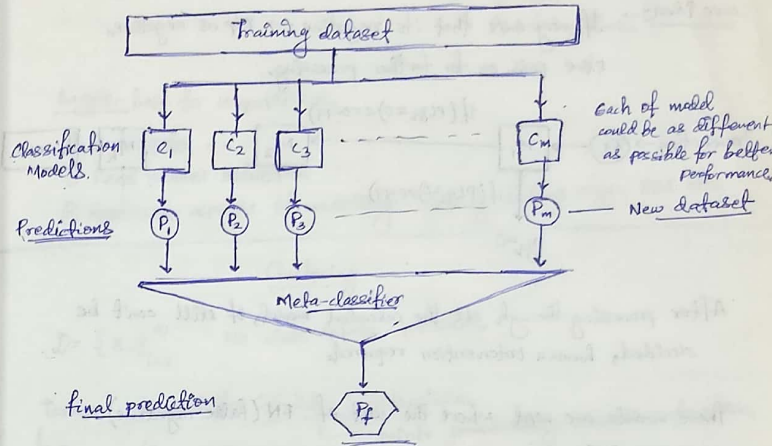
linearly combining all 3, we get



GBDT more often used in internet companies than AdaBoost.

Stacking Models

library - `mlxtend.classifier` - see official docs for details.



Algorithm: Stacking

Input:- Training data $D = \{x_i, y_i\}_{i=1}^m$ ($x_i \in \mathbb{R}^n, y_i \in Y$)

Output:- An ensemble classifier H .

1. learn first-level classifiers (base models).

for $i=1$ to T do
learn a base classifier h_i based on D .

end for.

2. construct new dataset from D using trained models.

for $i=1$ to m do

construct a new dataset that contains $\{x'_i, y_i\}$

where $x'_i = \langle h_1(x_i), h_2(x_i), h_3(x_i), \dots, h_T(x_i) \rangle$

end for.

3. learn a second-level classifier (meta-classifier).

learn a new classifier h' based on the newly constructed dataset.

return $H(x) = h'(\langle h_1(x), h_2(x), h_3(x), \dots, h_T(x) \rangle)$

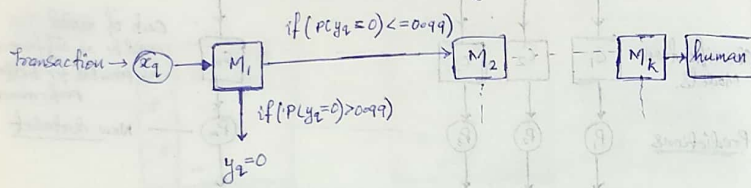
$h'_i(x)$ can be $\begin{cases} \text{off label} \\ \text{probability value} \end{cases}$

Kaggle uses
stacking extensively

Cascading Models

Ex → predicts credit card transaction is fraud or not.

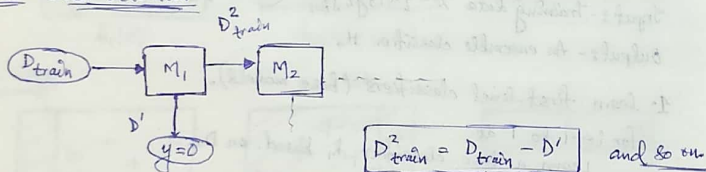
core ideas → If very sure that its negative → off as negative, else pass on for further processing.



After processing through all the cascaded models, if still can't be decided, human intervention required.

These models are used where the cost of FN (false negative) is very high.

Ex → Cancer detection



There could be more complex architecture also for cascade models.

Kaggle Vs Real World

① Kaggle cares about a single metric. - log-loss, F1-score, AUC, etc.

Real world - multiple metrics with different priorities.

② This is because in real-world we try to solve a business problem i.e. increase sales, etc.

So, we map this biz metric → multiple ML metrics.

③ Kaggle → very complex ensembles. → goal to improve 1 metric.

This might be impractical in real world due to: ① low latency ② interpretability ③ training time.

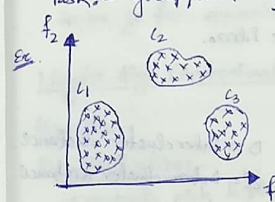
Kaggle - best for competitions.

- ① Participate
- ② Read winner solutions.
- ③ Learn aspects of modelling - cleaning, feature engg, EDA, etc.

Clustering

$D = \{x_i\}_{i=1}^n$ No class labels provided.

Task: - group/cluster "similar" datapoints.



- points geometrically close to each other.
- (a) points in a cluster are close together.
 - (b) points in different clusters are far away.

How to measure how well the algo performed? → All prev metrics are based on y_i .

Most used clustering →

- ① k-means
- ② Hierarchical.
- ③ DB-scan.

Unsupervised Learning

classification & Regression → Supervised learning.

$\{x_i, y_i\}$ the labels supervised/guided the learning process.

clustering → unsupervised learning, No y_i available.

Semi-supervised learning → $D = D_1 \cup D_2$

$|D_1| \ll |D_2|$

$\{x_i, y_i\}$ small

$\{x_i\}$ large

Application of clustering

data mining

- ① e-commerce → Group similar customers based on their purchasing behavior so as to offer deals to a group of similar customers who are more likely to take the deal.
- ② Image segmentation → grouping/clustering similar pixels.
forms base of object detection
- ③ Data labelling → cluster the data into 2 groups.
check some instance of 1 group to know whether its positive/negative class.
saves a lot of time.

Metrics for clustering

All the metrics have the same common base idea.

- ① Maximize inter-cluster distance.
- ② Minimize intra-cluster distance.

One example

Dunn Index →

$$D = \frac{\min D}{\max d}$$

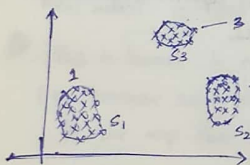
D = intercluster distance
 d = intra-cluster distance

D is high \Rightarrow good clustering.

K-means clustering

k : # clusters. - hyperparameter, found similar to cv.

* Each cluster has a centroid.



$k=3$

$$s_1 \cup s_2 \cup s_3 = D$$

$$s_1 \cap s_2 = \emptyset$$

$$s_2 \cap s_3 = \emptyset$$

$$s_1 \cap s_3 = \emptyset$$

centroid of S_j

$$c_j = \frac{1}{n} \sum_{x_j \in S_j} x_j$$

K-means - centroid based clustering.

Others are ① density based clustering.
② hierarchical clustering.

* Once we get the k -centroids, each point is assigned to the nearest centroid.

Mathematical formulation

Objective function →

$$\arg \min_{c_1, c_2, \dots, c_k} \sum_{i=1}^k \sum_{x \in S_i} \|x - c_i\|^2$$

It basically says that find centroids such that intracluster distance are minimized.

such that $S_i \cap S_j = \emptyset$.

* This problem is NP-Hard. Exponential time complexity.
so we go for approximation algorithms.

Lloyd's Algo (Approximation)

Algorithm →

① Initialization:-

→ randomly pick k pts from D and call them c_1, c_2, \dots, c_k .

② Assignment:-

for each point x_i in D
select the nearest c_j .
Add x_i to set S_j .

③ Recompute centroid/update

recalculate/update c_j 's as follows.

$$c_j = \frac{1}{|S_j|} \sum_{x_i \in S_j} x_i$$

④ Repeat ② & ③ until convergence.

→ recalculate/update

no significant difference in c_j 's

K-means — initialization sensitive.

↳ final clusters depend on initialization. — Refer slides.

One solution

→ Repeat k-means multiple times with different initializations.

Pick the best clustering based on

* smaller intracluster distance.

* larger intercluster distance.

Alternate solution (K-Means++)

① pick the first centroid randomly from D (c_1).

② for each non-centroid point in D , calculate distance of x_i to its closest centroid (d_i).

③ Pick a point out of all these points with a probability of d_i & nominate it as next centroid.

④ Repeat ② & ③ to get all the centroids.

* We don't directly select the pt with maximum distance so as to avoid impact of outliers.

Limitations of K-Means — refer slides.

* K-means tries to create clusters of same sizes (density may vary).

① clusters of different sizes.

② clusters of different density.

③ non-convex shapes.

Overcoming limitations

① clusters of different sizes → increase k and then again re-group clusters.

K-medoids

* centroids shouldn't compulsorily be one of the training set points.
* Also, they may end up having fractional components which may make them less interpretable.

Ex → BOM

* Idea of k-medoids is to select one of the training set points which is close to the actual centroid.

Partitioning Around Medoids (PAM)

① Initialization → similar to K-Means++ probabilistic methods.

② Assignment → same ✓ → assign to closest medoid.

③ Update → swap each medoid with each non-medoid point. If loss decreases, do nothing; else undo the swap.

loss is given by

$$\sum_{i=1}^k \sum_{x \in S_i} \|x - m_j\|^2$$

medoid

Advantage of medoids

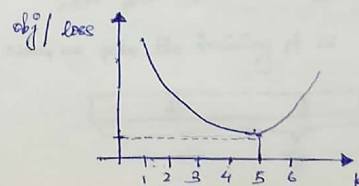
① Interpretable

② can use similarity matrix also, hence kernelization.

Determining the right 'k'

① domain knowledge — If we know only two classes possible — (t)me & (n)me. Then, no of clusters = 2.

② elbow method



Time complexity

$$\text{k-means} - O(\underbrace{n}_{\text{\#pts}} * \underbrace{k}_{\text{\#clusters}} * \underbrace{d}_{\text{dim}} * \underbrace{i}_{\text{\#iterations}})$$

$$\approx \underline{O(nd)}$$

$$(k \leq 10) \text{ typically}$$

$$(i \leq 300)$$

$$\text{Space complexity} - O(nd + kd)$$

$$\approx \underline{O(nd)}$$