

General Linear Model

GLM in general refers to conventional linear regression models for a continuous response variable given continuous and categorical predictors. It includes multiple linear regression as well as ANOVA and ANCOVA.

$$y_i \sim N(x_i^T \beta, \sigma^2)$$

$x_i \rightarrow$ contains known covariates.

$\beta \rightarrow$ coefficient to be estimated

These models are fit by least squares and weighted least squares using.

for Ex: SAS's GLM procedure or R's `lm()` function

Linear Regression

Prediction is done by simply computing a weighted sum of input features plus a constant called the bias term. / Intercept

→ Gradient Descent is used to tune parameters iteratively in order to minimize a cost function. \Rightarrow size of step is important
↓
Learning rate.

Logistic Regression (can be used for classification as well)
↓
binary classifier. Is used to estimate whether the probability that an instance belongs to a particular class
if the estimated probability $\geq 50\%$. then the instance belongs to class

→ way to reduce overfitting.

Regularized Linear models. Regularization

Regularization → typically achieved by constraining the weights of the model.

a) Ridge regression → regularized version of linear Regression.
term $\lambda(\text{slope})^2$ → it not only fits the data but also keep

a) Lasso regression regularized term. the model weighs as small as possd.

↳ least absolute shrinkage & selection operator.

it tends to completely eliminate the weights

it performs feature selection.

of the least important features.

c) Elastic net → middle btm Ridge & Lasso

Ensemble Techniques

Two most popular ensemble methods

↳ Bagging & Boosting (here a single base learning algorithm is used)

Bagging is a technique for reducing prediction variance by producing additional data for training from dataset by combining repetitions with combinations to create multi-sets of the original data.

Boosting is an iterative strategy for adjusting an observation's weight based on the previous classification.

Ensemble method is used in ML to train multiple models or weak learners to solve the same problem & integrated to give desired result.

↳ weak models combined rightly give accurate models.

The Ensemble model made by Bagging & Boosting is known as homogeneous model.

There are some methods in which different types of base learning algorithms are also implied with heterogeneous weak learners to make a 'heterogeneous ensemble models'.

Ensemble Techniques

Bagging → A homogeneous weak learners' model that learns from each other independently in parallel & combines them for determining the model average.

Boosting → A homogeneous weak learners' model but works differently from Bagging. Here learners learn sequentially & adaptively to improve model predictions of a learning algorithm.

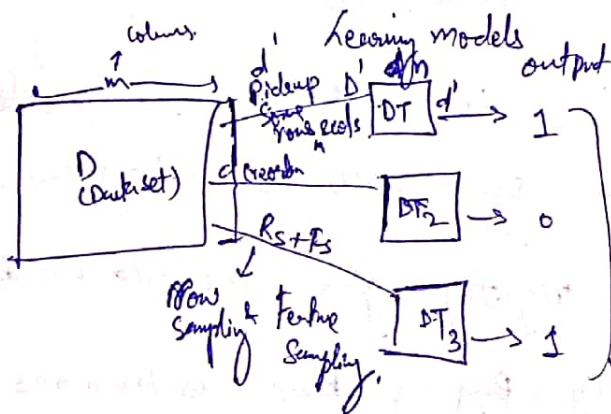
Ensemble Techniques

→ Combining multiple models.

- ① Bagging (Bootstrap aggregation)
 - ② Random Forest
- Row sampling with replacement.

- ③ Boosting
- ④ Ada Boost
- ⑤ Gradient Boosting
- ⑥ XG Boost

Random Forest



$$\begin{aligned} D' &< D \\ n &< m \\ d' &< d \end{aligned}$$

DT₁ & DT₂ etc are different sets
Shuffling was done.

majority 1

Decision Tree Properties when we let DT to go all the depth.

- ① Low Bias
- ② High Variance

Care is Random Forest as it get trained for specific records & gives better test accuracy.

In Regression problem in RF we take either average or median of all the results, Classifier takes majority of votes.

Ensemble methods :- Bagging, boosting, stacking etc.
hard voting classifier - The class which gets the most votes

Bagging :- Uses the same training algorithm for every predictor but to train them on different random subsets of the training set.
(Bootstrap aggregating) \rightarrow Here sampling is performed with replacement.

Pasting \rightarrow when sample is performed without replacement.

Feature Importance :- RF makes it easy to measure the relative importance of each feature.
Scikit learn measure/determines how much the tree nodes that use that feature reduce impurity on avg.

Boosting :- Ensemble method that combines several weak learners into a strong learner.
 \rightarrow general idea :- to train predictors sequentially \rightarrow each trying to correct its predecessor.

One way to do it is Ada Boost. \rightarrow Pays more attention to the training instances that the predecessor underfitted.
 \hookrightarrow results in new predictors focusing more and more on the hard cases.

Ada Boost:- ① Base Classifier (DT classifier) is trained & used to make predictions on the training set.

② Relative weight of misclassified training instances is then increased.

③ A second classifier is trained, using the updated weights & given it to make predictions on the training set.

④ weights are updated and so on.

⑤ Once all the predictors are trained, the ensemble makes predictions very much like bagging or pasting, except that predictors have different weights depending on their overall accuracy on the weighted training set.

Drawback:- It cannot be parallelized since each predictor can only be trained after the previous predictor has been trained.

→ result:- It does not scale as well as bagging or pasting.

Gradient Boosting:- Works by sequentially adding predictors to an ensemble each on corrected its predecessor ^{at every iteration}.

→ Instead of tweaking weights like AdaBoost it tries to fit the new predictor to the residual errors made by previous predictors.

XG Boost:- An optimized implementation of Gradient Boosting available in python library XGBoost.
↓
Extreme Gradient Boosting.

Aims:- being extremely fast, scalable & portable.

Stacking → Here instead of using trivial functions to aggregate the predictions of all predictors in an ensemble it trains a model to perform this aggregation.

→ the final predictor (Blender or meta learner) takes these predictions as inputs & makes the final predictions.

It is actually possible to train several blenders (one using linear regression, other using RF etc ...)

Trick is to split the training set into three subsets.

→ The first one is used to train first layer.

→ Second one is used to create the training set used to train the second layer. (using predictions made by predictors of the first layer)

→ third one is used to create the training set to train 3rd layer. (using predictions made by the predictors of the 2nd layer)

~~and so on~~
→ Once it is Done, we can make prediction for a new instance by going through each layer sequentially.

SVM

SVM → Capable of performing linear or ^{non}linear classification, regression and even outliers detection.

→ well suited for classification of complex but small or medium sized data.

Linear SVM classification

→ two classes can ~~be~~ ^{clearly be} separated easily with a straight line.

→ The st. line not only separates the two classes but also stays as far away from the closest training instances as possible.

Hard margin Classification

↳ Only work for data which is linearly separable & it is quite sensitive to outliers.

Soft margin Classification

↳ more flexible, to keep the street as large as possible & limiting the margin violation.

Non linear SVM classification

↳ to add more features.

slow for high complex features.
(poly, rbf)

Decision Tree (classification & regression task \rightarrow versatile) & multiclass tasks

\rightarrow capable of fitting complex data.

Start from root node (depth = 0, at the top)

root's child nodes (left & right) \rightarrow depth 1 and so on.
 \downarrow
if condition satisfies.

leaf node \rightarrow doesn't have any child node.

Decision trees require very little data preparation & do not require feature scaling at all in particular.

node's value attribute tells how many training instances of each class this node applies to.

When a node is pure (gini = 0) if all the training instances belong to same class.

$$Gini = 1 - \sum (P_i^2)$$

ADT can also estimate the probability that an instance belongs to a particular class.

reduction of Entropy is often called an information gain

$$H_i = - \sum P_{i,k} \log_2 (P_{i,k})$$

Random forest is an ensemble of Decision Trees.

here it is easy to measure the relative importance of each feature.

Boosting : To train predictors ~~separately~~ sequentially.

Ada Boost : firstly base classifier is ~~used for~~ trained & used to make predictions on training set.

Adaptive.

Second classifier is trained using the updated weights. & soon it dies tweaking the instance weights at every iteration.

Gradient Boosting

Boost & Random forest come under ensemble

↳ works sequentially like AdaBoost

Adding predictors each one correcting its predecessor.

This method tries to fit a new predictor to the residual errors made by previous predictor.

Learning rate is ~~steps~~ steps, low learning rate (0.1) needs more trees in the ensemble to fit the training set.

To find optimal no of trees early stoppers are used.