

$$Z Z^T = V \text{diag} \left(\frac{\sigma_i^2}{(\sigma_i^2 + \lambda I)^2} \right) V^T$$

→ Bagging: Randomforest

Bootstrapping: Sampling technique where sample are derived from whole population using replacement procedure.

* Aggregation: in bagging is done to incorporate all possible outcomes of prediction & randomize the outcomes.

- Reduce the variance.

* Boosting: Combine weak learning to form strong learner.

Adaboost used boost (XGBoost)

- misclassified point are more important.

Aduboost
Ziken

Weighted Mode ($M_1, M_2 \dots M_n$)

$$L = \sum_{i=1}^n w_i M_i$$

→ use decision stump

• XkBoost

- use verify & learn.

* Class imbalance

- Means # no. of points in some class is higher than other classes
e.g. credit card fraud detection.
- Bias towards majority class
- Acc. may not be good performance measure for classifiers for imbalanced data

Soln

* Random Sampling

Adv: simple to implement

- Efficient

Dis: throws away lot of data

- not represent original data

* Random oversampling

- Dis: overfit on minority class.

Adv: no loss data

- Empirically seen perform better than understanding

* SMOTE: Synthetic Minority Over Sampling Technique

- Avoid overfitting due to explicit replicas of minority class.
- Subject minority class is taken.

Dis: No good for data is in higher dimension.

- may introduce noise

Adv: Reduce a chance of overfitting.

* Other techniques: Penalization based model.

- Class weight based model.
- try to solve diff. evaluation measure
- Ensemble methods.

*

* Ada Boost :

Setting : data : $\{(x_1, y_1), (x_2, y_2) \dots (x_n, y_n)\}$

$y_i \in \{-1, +1\} \rightarrow$ binary class problem

weight $\rightarrow w_i = \frac{1}{n} \quad \forall i$

① learner which can work with weighted data points to give

predicted op $h(x) \in \{-1, +1\}$

AdaBoost : ① initialize $\{w_i\}$ as $w_i^{(1)} = \frac{1}{n}$

for $i = 1 \dots n$

② for $t = 1 \dots T \rightarrow$ this how many learner we are actually apply.

(a) fit $h_t(x)$ to training data by minimizing

loss function \downarrow the misclassified points.

$$J_t = \sum_{i=1}^n w_i^{(t)} I(h_t^{(x)} \neq y_i) \quad \text{--- (A)}$$

I is indicator which give op 0/1 for class labels.

fraction \downarrow miss classify points

$$\epsilon_t = \frac{\sum_{i=1}^n w_i^{(t)} I(h_t^{(x)} \neq y_i)}{\sum_{i=1}^n w_i^{(t)}}$$

$$\text{set } \alpha_t = \ln \left(\frac{1 - \epsilon_t}{\epsilon_t} \right) \quad \text{--- (B)}$$

(c) Update weight

$$w_i^{(t+1)} = w_i^{(t)} \exp(\alpha_t I(h_t^{(x_i)} \neq y_i))$$

this means if correctly classify points then $I(h_t^{(x_i)} \neq y_i)$ is zero then we did not change the weight. when it is incorrectly classify then weight will increase. (c)

3). Final Model

$$H_T(x) = \text{sign} \left(\sum_{t=1}^T d_t h_t(x) \right)$$

→ why does it work
- Boosting as sequential minⁿ of exponential error funⁿ

$$E = \sum_{i=1}^n \exp(-y_i f_t(x_i))$$

where $f_t(x) = \frac{1}{2} \sum_{t=1}^t d_t h_t(x)$

Goal: Minimize E wr. to both d_t and $h_t(x)$ (parameters)

Assume $h_1(x) \dots h_{t-1}(x)$ & $d_1 \dots d_{t-1}$ all fixed minimize only wr. to $d_t, h_t(x)$

$$\left\{ \because f_t(x) = f_{t-1}(x) + f_t(x) \right\}$$

$$\therefore E = \sum_{i=1}^n \exp \left\{ \underbrace{-y_i f_{t-1}(x_i)}_{\text{constant}} - \frac{1}{2} y_i d_t h_t(x) \right\}$$

$$(e^{-a-b} = e^{-a} \cdot e^{-b})$$

↓
this const.

$$= \sum_{i=1}^n w_i^{(t)} \exp \left(-\frac{1}{2} y_i d_t h_t(x) \right)$$

$$\left(\because w_i^{(t)} = \exp(-y_i f_{t-1}(x_i)) \right)$$

- Consider $w_i^{(t)}$ as constants.

can be considered constants.

let $C_t =$

* Consider a error

$$E = \sum_{i=1}^n \exp(-y_i f_t(x_i))$$

where $f_t(x) = \frac{1}{2} \sum_{l=1}^t d_l h_l(x)$
Hinge.

$$E = \sum_{i=1}^n w_i^{(t)} \exp\left(-\frac{1}{2} y_i d_t h_t(x_i)\right)$$

$$\left(\because w_i^{(t)} = \exp\left(\sum_{l=1}^{t-1} -y_i d_l h_l(x_i)\right) \right)$$

C_t = data point correctly classified point

M_t = " " in correctly " "

$$E = e^{-d_t/2} \sum_{i \in C_t} w_i^{(t)} + e^{d_t/2} \sum_{i \in M_t} w_i^{(t)}$$

$$E = \left(e^{d_t/2} - e^{-d_t/2} \right) \sum_{i=1}^n w_i^{(t)} \frac{I(h_t(x_i) \neq y_i)}{2}$$

which 1 for misclassified points
& 0 = for correctly class. points

\Rightarrow Minimize w.r.t. $h_t(x)$ Second term is Constant Equivalent to minimizing (A).

~~\rightarrow Similarly~~

\rightarrow Similarly minimize w.r.t. d_t we get (B)

$$\frac{\partial E}{\partial d_t} : e^{d_t/2} \left(\frac{1}{2}\right) - e^{-d_t/2} \left(-\frac{1}{2}\right) + e^{-d_t/2} \left(-\frac{1}{2}\right)$$

$$= e^{d_t/2} \left(\frac{1}{2}\right)$$

$$0 = \sum_{i=1}^n w_i^{(t)} I(h_t^{(x_i)} \neq y_i) e^{d_t/2} \cdot \left(\frac{1}{2}\right)$$

$$= \left(\sum_{i=1}^n w_i^{(t)} I(h_t^{(x_i)} \neq y) \cdot \left(\frac{1}{2}\right) \cdot e^{-d_t/2} \right.$$

$$\left. + \sum_{i=1}^n w_i^{(t)} e^{-d_t/2} \cdot \left(\frac{1}{2}\right) \right)$$

$$0 = \sum_{i=1}^n w_i^{(t)} I(h_t^{(x_i)} \neq y) \cdot e^{d_t/2} + \sum_{i=1}^n w_i^{(t)} \cdot e^{-d_t/2} + \frac{\left(\sum_{i=1}^n w_i^{(t)} I(h_t^{(x_i)} \neq y) \cdot e^{-d_t/2} + \sum_{i=1}^n w_i^{(t)} \cdot e^{-d_t/2} \right)}{e^{d_t/2}}$$

$$0 = \sum_{i=1}^n w_i^{(t)} I(h_t^{(x_i)} \neq y) (e^{d_t/2})^2 + \sum_{i=1}^n w_i^{(t)} I(h_t^{(x_i)} \neq y) + \sum_{i=1}^n w_i^{(t)}$$

$$(e^{d_t/2})^2 = - \left(\sum_{i=1}^n w_i^{(t)} I(h_t^{(x_i)} \neq y) + \sum_{i=1}^n w_i^{(t)} \right)$$

Derivative wrt

$$e^{d_t} + 1 = \frac{1}{\epsilon_t}$$

$$e^{d_t} = \frac{1 - \epsilon_t}{\epsilon_t}$$

$$d_t = \ln \left(\frac{1 - \epsilon_t}{\epsilon_t} \right)$$

Now $w_i(t+1) = w_i(t) \exp \left\{ -\frac{1}{2} y_i d_t h_t(x_i) \right\}$ put this here

Also $y_i h_t(x_i) = 1 - 2 I(h_t(x_i) \neq y_i)$

$$= w_i(t) \exp \left(-\frac{1}{2} d_t (1 - 2 I(h_t(x_i) \neq y_i)) \right)$$

⇒ in next iteration this term is independent from the data point then ignore it

$$w_i(t+1) = w_i(t) \exp(-d_t/2) \Rightarrow w_i(t) / \exp(d_t/2)$$

$$+ w_i(t) \exp(d_t I(h_t(x_i) \neq y_i))$$

$$w_i(t+1) = w_i(t) \exp(d_t I(h_t(x_i) \neq y_i))$$

this is (C)

* GMM: Gaussian distⁿ by itself may not enough in modeling real datasets.

→ sometime linear superposition of two or more Gaussians may combine does better.

Such superposition can be formulated as mixture

distⁿ
(this why it is in unsupervised learning.

→ ~~not~~ Mean combine two distⁿ

