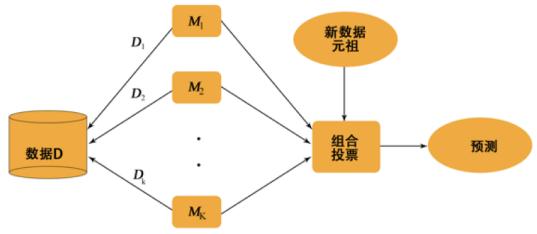


Ensemble Methods 集成化方法

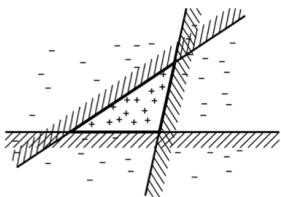
- Wisdom of Crowds ("三个臭皮匠，顶个诸葛亮")
- Multiple weak learners (base learners, may be heterogenous) can improve learning performance



基于不同的信息 (数据集划分、加权)

Why it can improve the performance

- More expressive, can approximate larger functional space
 - Single linear classifier (perceptron) does not work
 - Try multiple classifiers

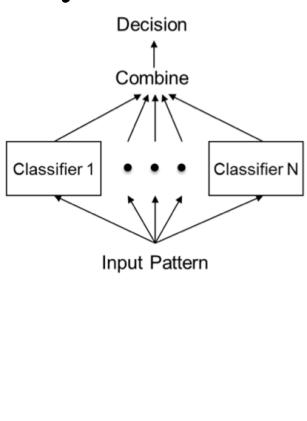


- Reduce misclassification rate
 - Misclassification rate of single classifier is p
 - Choose N classifiers, same type but independent (i.i.d.), voting
 - Error rate of majority vote = $\sum_{k>N/2} \binom{N}{k} p^k (1-p)^{N-k}$
 - When N = 5, p = 0.1, Error rate < 0.01

Two commonly used ensemble methods

• Bagging 袋装方法 (Bootstrap aggregation)

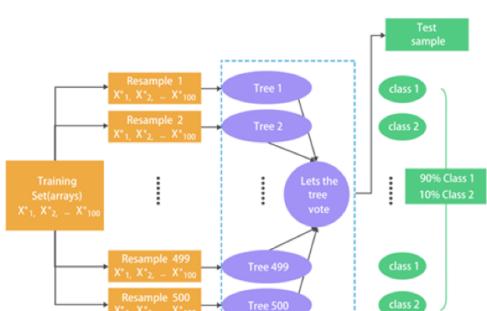
- Random sampling : generating independent models, and averaging for regressions (making majority vote for classifications)
- Reducing variances
- Example : Random forests



- Boosting
- Sequential training : training the subsequent models based on the errors of previous models
- Reducing bias
- Examples : AdaBoost and GBDT

Bagging

- Bagging is short for bootstrap aggregation
- Bagging generates a committee of predictors and combine them in a certain manner to the final model
- Single predictor suffers from instability, while bagging could improve the stability by majority vote (classification) or averaging (regression) over all single predictors



用多个分类器(或感知机)的结果作平均

基模型 → 集成模型：①多数投票
②平均值方法
③加权平均方法

每个分类器的错分率为p, 且互相独立.

要投票 (若10个分类器中有6个投1, 4个投0, 则投1)
对结果统计错误率 (N个分类器, 二分类)

Bagging: ①随机采样
②降低方差

(基模型常用: 决策树和神经网络)

如: 随机森林

Boosting: ①序列化训练: 使这一次的分类器比上一次更好.
②降低偏差

例: GBDT, AdaBoost

Boosting中的GBDT和XGBoost都可以并行.

(有效回)

对数据集中的数据重新采样, 对每个样本训练模型, 得到的结果通过多数投票或平均整合.
(分类) (回归)

便于并行化: 数据量大 / 基模型少时有优势.

Sampling

- Given a dataset D of n samples, at the iteration $m = 1, \dots, M$, the training set D_m is obtained by sampling from D with replacement. Then D_m is used to construct classifier $\hat{f}_m(x)$.
- Sampling with replacement : some samples in D may be missing in D_m , while some other samples may occur more than once
- On average, 63.2% of the samples in D could be selected into D_m . In fact, for each sample, the probability that it is not selected in one round is $1 - \frac{1}{n}$. Then it is not selected in all n rounds with probability $\lim_{n \rightarrow \infty} (1 - \frac{1}{n})^n = 0.368$.

每次选择中一个样本选中的概率 $\frac{1}{n}$.

有的样本可能会被选中多次，有的不会被选中。

在数据集很大时，有0.368的数据选不到。

Algorithm

- Input : training set $D = \{(x_1, y_1), \dots, (x_N, y_N)\}$

- Output : additive model $\hat{f}_{bag}(x)$

- For $m = 1$ to M :

- Sample from D with replacement to obtain D_m

- Train a model $\hat{f}_m(x)$ from the dataset D_m : for classification, $\hat{f}_m(x)$ returns a K-class 0-1 vector e_k ; for regression, it is just a value

- Compute bagging estimate $\hat{f}_{bag}(x)$: for classification, make majority vote $\hat{f}_{bag}(x) = \arg \max_k \sum_{k=1}^M \hat{f}_k(x)$; for regression, just return the average value $\hat{f}_{bag}(x) = \frac{1}{M} \sum_{m=1}^M \hat{f}_m(x)$

分类就加1投票(众数)
回归取平均值。

Variance Reduction

- In bagging, we use the same model to train different sample set in each iteration ; assume the models $\{\hat{f}_m(x)\}_{m=1}^M$ have the same variance $\sigma^2(x)$, while the correlation of each pair is $\rho(x)$
- Then the variance of the final model is :

$$\begin{aligned} \text{Var}(\hat{f}_{bag}(x)) &= \frac{1}{M^2} \left(\sum_{m=1}^M \text{Var}(\hat{f}_m(x)) + \sum_{t \neq m} \text{Cov}(\hat{f}_t(x), \hat{f}_m(x)) \right) \\ &= \rho(x)\sigma^2(x) + \frac{1 - \rho(x)}{M}\sigma^2(x) \end{aligned}$$

- As $M \rightarrow \infty$, $\text{Var}(\hat{f}_{bag}(x)) \rightarrow \rho(x)\sigma^2(x)$. This usually reduces the variance.
- If $\rho(x) = 0$, the variance could approach zero
- The random sampling in bagging is to reduce the correlation $\rho(x)$, i.e., make the sub-predictors as independent as possible

Bagging不能改变bias,但能降低variance.

一般来说, $\rho(x)$ 很难为0.
采样越随机, $\rho(x)$ 越小。

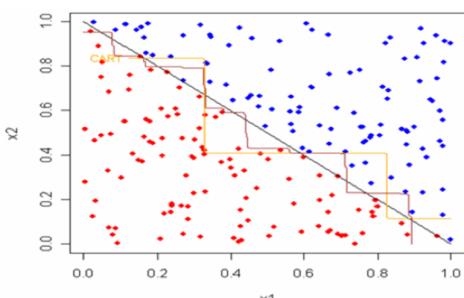
Bagging可以并行

Limitations of Decision Tree

- Stuck at local optimum : The greedy algorithm makes it stop at the local optimum, as it seeks the maximal information gain in each tree split
- Decision boundary : Use one feature in each split, the decision boundary is parallel to the coordinate axes
- Bad representability

决策树容易卡在局部最小值。

决策树的决策边界不好看出。
总是平行于坐标轴



若用Bagging可将多个决策树的决策边界合在一起。
找到一条直线边界

Random Forest

- Random Forest further reduces the variance by adding independency to the committee of decision trees
 - This is achieved by introducing more randomness.
 - More randomness :
 - Sampling on the training data with replacement
 - Select features at random
 - No pruning is needed.
 - Example : RF consisting of 3 independent trees, each with an error rate of 40%. Then the probability that more than one tree misclassify the samples is
$$0.4^3 + 3 * 0.4^2 * (1 - 0.4) = 0.352$$

Random Forest Algorithm

- Input : training set $D = \{(x_1, y_1), \dots, (x_N, y_N)\}$
 - Output : additive model $\hat{f}_{rf}(x)$

- For $m = 1$ to M :
 - Sample from D with replacement to obtain D_m
 - Grow a random-forest tree T_m to the dataset D_m : by recursively repeating the following steps for each terminal node of the tree, until the minimum node size n_{min} is reached
 - Select q features at random from the p features
 - Pick the best feature/split-point among the q
 - Split the node into two daughter nodes
 - Output the ensemble of trees $\{T_m\}_{m=1}^M$: for regression,

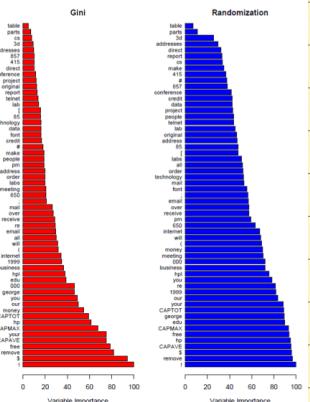
$$\hat{f}_{rf}(x) = \frac{1}{M} \sum_{m=1}^M T_m(x)$$
 : for classification, make majority vote
 - Small value of q increases the independency of trees;

Model Evaluation

- Margins : The difference between the percentage of decision trees that correctly classify each sample and the percentage of trees misclassifying it; margin is defined as the average difference for all samples
 - Out-of-bag (OOB) errors : The observation is called out-of-bag sample to some trees if it is not sampled for those trees. Denote the training set in the m -th sampling by D_m . OOB error is computed as :
 - For each observation (x_i, y_i) , find the trees which treat it as OOB sample : $\{\hat{T}_m(\mathbf{x}) : (x_i, y_i) \notin D_m\}$
 - Use those trees to classify this observation and make majority vote as the label of this observation :
$$\hat{f}_{oob}(\mathbf{x}_i) = \arg \max_{y \in \mathcal{Y}} \sum_{m=1}^M I(\hat{T}_m(\mathbf{x}_i) = y) I((x_i, y_i) \notin D_m)$$
 - Compute the number of misclassified samples, and take the ratio of this number to the total number of samples as OOB error : $Err_{oob} = \frac{1}{N} \sum_{i=1}^N I(\hat{f}_{oob}(\mathbf{x}_i) \neq y_i)$

Feature Importance

- Using split criteria
 - The improvement in the split-criterion as feature importance
 - It is accumulated over all the trees for each variable
 - Using OOB randomization
 - Randomly permute the values of each feature in the OOB samples, and compute the prediction accuracy
 - The decrease in accuracy as a result of this permutation is averaged over all trees as feature importance



Bagging 只要求有放回的随机采样，而 Random Forest 还会随机抽取特征。

通过引入随机性减少决策树的贪心和方差。

不需要減枝.

每次从 P 个特征中选 q 个，再用贪心分节点。

$q = \log_2 P + 1$ 是经验公式.

岱外错误

对每个样本，找到没用到这个样本的树
用这些树对这个样本作预测。

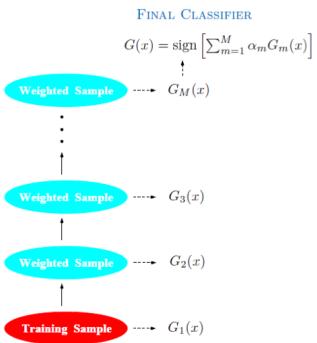
OOD Error 即是预测值不等于真实值的频率。
(当N趋于无穷大时, OOD Error = cross validation error)
需用大数定律证明

Pros and Cons

- Where it is good
 - Bagging or random forest (RF) work for models with high variance but low bias
 - Better for nonlinear estimators
 - RF works for very high-dimensional data, and no need to do feature selection as RF gives the feature importance
 - Easy to do parallel computing
- Disadvantage
 - Overfitting when the samples are large-sized with great noise, or when the dimension of data is low
 - Slow computing performance comparing to single tree
 - Hard to interpret

Boosting

- Boosting : combines the outputs of many "weak" classifiers to produce a powerful "committee"
- Weak classifier : error rate < 0.5 (random guessing)
- Sequentially apply the weak classifiers to the repeatedly modified data, emphasizing the misclassified samples
- Combine weak classifiers through a weighted majority vote or averaging to produce the final prediction



Boosting Fits an Additive Model

- Additive model : $f(x) = \sum_{m=1}^M \beta_m b(x; \gamma_m)$
- Possible choices for basis function $b(x; \gamma)$:
 - Neural networks : $\sigma(\gamma_0 + \gamma_1^T x)$, where $\sigma(t) = 1/(1 + e^{-t})$
 - Wavelets
 - Cubic spline basis
 - Trees
 - Eigenfunctions in reproducing kernel Hilbert space (RKHS)
- Parameter fitting : $\min_{\{\beta_m, \gamma_m\}} \sum_{i=1}^N L(y_i, \sum_{m=1}^M \beta_m b(x_i; \gamma_m))$
- Loss function : squared error $L(y, f(x)) = (y - f(x))^2$ or likelihood-based loss

Forward Stagewise Additive Modeling

- Input : training set $D = \{(x_1, y_1), \dots, (x_N, y_N)\}$
- Output : additive model $f_M(x)$

1. Initialize $f_0(x) = 0$
2. For $m = 1$ to M :
 - 2.1 Compute $(\beta_m, \gamma_m) = \arg \min_{\beta, \gamma} \sum_{i=1}^N L(y_i, f_{m-1}(x_i) + \beta b(x_i; \gamma))$
 - 2.2 Update $f_m(x) = f_{m-1}(x) + \beta_m b(x_i; \gamma_m)$
- Squared error loss : in step 2.1,

$$L(y_i, f_{m-1}(x_i) + \beta b(x_i; \gamma)) = (\underbrace{y_i - f_{m-1}(x_i)}_{\text{residual}} - \beta b(x_i; \gamma))^2$$

Bagging / Random Forest 对高方差低偏差模型效果好。 (因为 Bagging 缩降方差)

数据噪音大时会过拟合。

解释性差 (集成化方法的通病)

弱分类器：错误率<0.5 (仅比随机猜测稍好)

序列化训练：每次的样本都根据权重更新。

通过改变权重提升分类器性能

将弱分类器通过加权多数投票或平均来集成

基函数一般都用同一种分类器。

混合着选效率较低，因为不同分类器参数不同。

Boosting 要预先知道弱分类模型的误差上界，但 AdaBoost 不用。

Exponential Loss and AdaBoost

- Exponential loss : $L(y, f(x)) = \exp(-yf(x))$
- Classifier as basis function : $b(x; \gamma) = G(x) \in \{-1, 1\}$
- Let $w_i^{(m)} = \exp(-y_i f_{m-1}(x_i))$, then step 2.1 turns to be :

$$(\beta_m, G_m) = \arg \min_{\beta, G} \sum_{i=1}^n w_i^{(m)} \exp(-\beta y_i G(x_i))$$

$$= \arg \min_{\beta, G} \left[\sum_{y_i \neq G(x_i)} w_i^{(m)} (e^\beta - e^{-\beta}) + e^{-\beta} \sum_{i=1}^n w_i^{(m)} \right]$$

交替方向优化

- $G_m = \arg \min_G \sum_{i=1}^n w_i^{(m)} I(y_i \neq G(x_i))$.
- $\beta_m = \arg \min_\beta \left[\epsilon_m (e^\beta - e^{-\beta}) + e^{-\beta} \right] = \frac{1}{2} \log \frac{1-\epsilon_m}{\epsilon_m}$ where
 $\epsilon_m = (\sum_{i=1}^n w_i^{(m)} I(y_i \neq G(x_i))) / \sum_{i=1}^n w_i^{(m)}$ is weighted error rate

AdaBoost Algorithm

- Input : training set $D = \{(x_1, y_1), \dots, (x_N, y_N)\}$, loss function $L(y, f(x))$
- Output : Weighted classifier $G(x)$

1. Initialize $w_i = 1/N$, $i = 1, \dots, N$

2. For $m = 1$ to M :

- 2.1 Fit a classifier $G_m(x)$ to the training data D with weight $\{w_i\}$
- 2.2 Compute the error $\epsilon_m = (\sum_{i=1}^n w_i^{(m)} I(y_i \neq G_m(x_i))) / \sum_{i=1}^n w_i^{(m)}$
- 2.3 Compute $\alpha_m = \log \frac{1-\epsilon_m}{\epsilon_m}$ ($\alpha_m = 2\beta_m > 0$)
- 2.4 Update the weight $w_i^{(m+1)} = w_i^{(m)} \exp(\alpha_m I(y_i \neq G_m(x_i)))$, for $i = 1, \dots, N$

3. Output $G(x) = \text{sign} \left[\sum_{m=1}^M \alpha_m G_m(x) \right]$

初始化数据有相同权重.

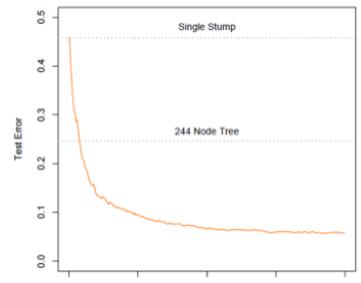
因为弱分类器有 $\epsilon_m < 0.5$, 故 $\alpha_m = \log \frac{1-\epsilon_m}{\epsilon_m} > 0$.

核心即提高分错样本的权重.

当每次迭代生成的集成模型 $\sum_m \alpha_m G_m(x)$ 没有出现分类错误时停止迭代.

Illustration

- Weights of weak classifiers : the better the classifier is, the larger its weight is
- Weights of samples : Re-weighting after each step, increase the weights for misclassified samples
- Simulation : 2-class classification, 1000 training samples from each class, 10,000 test samples ; two-leaf classification tree (stump) as base learner

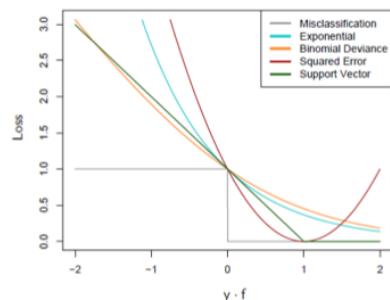


Stump: 只有两个叶子节点的决策树.
(也就是只找到一个feature分裂)

越好的弱分类器有越大的权重
在每一步后更新权重, 增加错分样本的权重.
减小正确预测样本的权重.

Loss Functions

- For classification, exponential loss and binomial negative log-likelihood (deviance) loss $\log(1 + \exp(-2yf))$ share the same population minimizer; thus it is equivalent to MLE rule
- For classification, squared error loss is not good (not monotonically decreasing); the exponential loss is good and binomial deviance is better (less penalty for large $-yf$)



square loss 在分类错误时惩罚太大，对异常值敏感。robust 较差，一般少选。
0-1-loss 不能求导

Pros and Cons

- Where it is good
 - AdaBoost improve the classification performance comparing to weak classifiers
 - Many choices for weak classifiers: trees, SVMs, kNNs, etc.
 - Only one tuning parameter M : # of weak classifiers
 - prevent overfitting suffered by single weak classifiers (e.g. complex decision tree)
- Disadvantage
 - Weak interpretability
 - Overfitting when using very bad weak classifiers
 - Sensitive to outliers
 - Not easy for parallel computing

Boosting Tree

- Using classification trees or regression trees as base learners
- $f_M(x) = \sum_{m=1}^M T(x; \Theta_m)$ where $T(x; \Theta) = \sum_{j=1}^J \gamma_j I(x \in R_j)$
- Parameter set $\Theta = \{R_j, \gamma_j\}_{j=1}^J$ → 一棵树
- Parameter finding: minimizing the empirical risk

$$\hat{\Theta} = \arg \min_{\Theta} \sum_{j=1}^J \sum_{x_i \in R_j} L(y_i, \gamma_j) \quad (\text{Combinatorial optimization})$$

最小二乘回归问题

- Approximate suboptimal solutions:
 - Finding γ_j given R_j : $\gamma_j = \bar{y}_j = \frac{1}{|R_j|} \sum_{y_i \in R_j} y_i$ for L^2 loss; and
 γ_j = modal class in R_j for misclassification loss
 - Finding R_j given γ_j : Difficult, need to estimate γ_j as well; → 贪心，降低不纯度
greedy, top-down recursive partitioning algorithm

Boosting Tree as Forward Stagewise Algorithm

- $\hat{\Theta}_m = \arg \min_{\Theta_m} \sum_{i=1}^N L(y_i, f_{m-1}(x_i) + T(x_i; \Theta_m))$
 - $\hat{\gamma}_{jm} = \arg \min_{\gamma_{jm}} \sum_{x_i \in R_{jm}} L(y_i, f_{m-1}(x_i) + \gamma_{jm})$
 - Finding R_{jm} is more difficult than for a single tree in general.
- Squared-error loss: fit a tree to the residual

$$L(y_i, f_{m-1}(x_i) + T(x_i; \Theta_m)) = (\underbrace{y_i - f_{m-1}(x_i)}_{\text{residual}} - T(x_i; \Theta_m))^2$$
- Two-class classification and exponential loss: AdaBoost for trees, $\hat{\Theta}_m = \arg \min_{\Theta_m} \sum_{i=1}^N w_i^{(m)} \exp[-y_i T(x_i; \Theta_m)]$
 - $\hat{\gamma}_{jm} = \log \frac{\sum_{x_i \in R_{jm}} w_i^{(m)} I(y_i=1)}{\sum_{x_i \in R_{jm}} w_i^{(m)} I(y_i=-1)}$
- Absolute error or the Huber loss: robust but slow

迭代次数 M 也是一个超参数，可通过 CV 来寻找。
用非常差的弱分类器容易过拟合。
解释性差
不能并行
对离群值敏感

→ 只有 L_2 -loss 才能这么转

用绝对误差和 Huber loss 虽然鲁棒但速度慢
更一般的 loss 只能用一般的优化方法

Gradient Descent for General Loss

- Supervised learning is equivalent to the optimization problem

$$\min_f L(f) = \min_f \sum_{i=1}^N L(y_i, f(x_i))$$

- Numerical optimization : $\hat{\mathbf{f}} = \arg \min_{\mathbf{f}} L(\mathbf{f})$ where $\mathbf{f} = \{f(x_1), f(x_2), \dots, f(x_N)\}$,
- Approximate $\hat{\mathbf{f}}$ by $\mathbf{f}_M = \sum_{m=0}^M \mathbf{h}_m$, where $\mathbf{f}_0 = \mathbf{h}_0$ is initial guess
- Gradient descent method : $\mathbf{f}_m = \mathbf{f}_{m-1} - \rho_m \mathbf{g}_m$, where $\mathbf{g}_m = \left[\frac{\partial L(y_i, f(x_i))}{\partial f(x_i)} \right]_{f(x_i)=f_{m-1}(x_i)}$, and $\mathbf{h}_m = -\rho_m \mathbf{g}_m$

$$\vec{f}_m = \vec{f}_{m-1} - \rho_m \vec{g}_m \quad \vec{g}_m = \left(\frac{\partial L}{\partial f(x_i)} \right)$$

通过 regression $T_m(x)$ over \vec{g}_m 来拟合.

$$f_m(x) = f_{m-1}(x) + T_m(x)$$

Gradient Boosting Decision Tree (GBDT)

- Find a tree $T(x; \Theta_m)$ by minimization problem

$$\tilde{\Theta}_m = \arg \min_{\Theta_m} \sum_{i=1}^N (-g_{im} - T(x_i; \Theta_m))^2$$

Setting	Loss Function	$-\partial L(y_i, f(x_i))/\partial f(x_i)$
Regression	$\frac{1}{2}[y_i - f(x_i)]^2$	$y_i - f(x_i)$
Regression	$ y_i - f(x_i) $	$\text{sign}[y_i - f(x_i)]$
Regression	Huber	$y_i - f(x_i)$ for $ y_i - f(x_i) \leq \delta_m$ $\delta_m \text{sign}[y_i - f(x_i)]$ for $ y_i - f(x_i) > \delta_m$ where $\delta_m = \alpha\text{th-quantile}\{ y_i - f(x_i) \}$
Classification	Deviance	kth component: $I(y_i = g_k) - p_k(x_i)$

GBDT Algorithm

- Input : training set $D = \{(x_1, y_1), \dots, (x_N, y_N)\}$, loss function $L(y, f(x))$
- Output : boosting tree $\hat{f}(x)$

- Initialize $f_0(x) = \arg \min_{\gamma} \sum_{i=1}^N L(y_i, \gamma)$
- For $m = 1$ to M :
 - For $i = 1, 2, \dots, N$ compute $r_{im} = \left[\frac{\partial L(y_i, f(x_i))}{\partial f(x_i)} \right]_{f=f_{m-1}}$
 - Fit a regression tree to the target (residual) r_{im} , giving terminal regions R_{jm} , $j = 1, \dots, J_m$
 - For $j = 1, \dots, J_m$, compute $\gamma_{jm} = \arg \min_{\gamma} \sum_{x_i \in R_{jm}} L(y_i, f_{m-1}(x_i) + \gamma)$
 - Update $f_m(x) = f_{m-1}(x) + \sum_{j=1}^{J_m} \gamma_{jm} I(x_i \in R_{jm})$
- $\hat{f}(x) = f_M(x)$

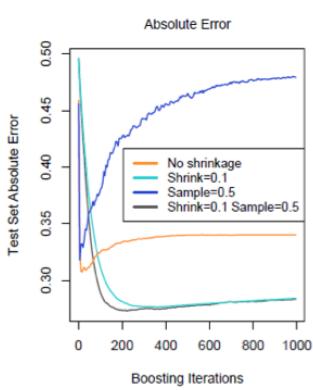
} 对一棵决策树的拟合

Regularization Techniques

- Shrinkage : the step 2.4 is modified as

$$f_m(x) = f_{m-1}(x) + \nu \sum_{j=1}^{J_m} \gamma_{jm} I(x_i \in R_{jm}) \quad \text{加一棵树}$$

- Subsampling : at each iteration, sample a fraction η of the training set and grow the next tree using the subsample
- Shrinkage + subsampling : best performance



收缩：每次新加的树都乘上一个小于1的权重，防止过拟合
(收缩系数)

子采样：在每次循环中，从训练集中采出一部分 η ，并从中生成树。

收缩+子采样一起用有更好的结果。

Feature importance and Partial Dependence Plots

• Feature importance

- When fitting a single tree T , at each node t , one feature $X_{v(t)}$ and one separate value $X_{v(t)} = c_{v(t)}$ are chosen to improve a certain quantity of criterion (e.g. GINI, entropy, squared error, etc.)
- Sum all these improvements i_t brought by each feature X_k over all internal nodes : $I_k(T) = \sum_{t=1}^{J-1} i_t I(v(t) = k)$
- Average the improvements of all trees \Rightarrow importance of that feature : $I_k = \frac{1}{M} \sum_{m=1}^M I_k(T_m)$

• Partial Dependence Plots

- Partial dependence of $f(X)$ on X_S : $f_S(X_S) = \mathbb{E}_{X_C} f(X_S, X_C)$
- Estimate by empirical mean : $\bar{f}_S(X_S) = \frac{1}{N} \sum_{i=1}^N f(X_S, x_{iC})$

Pros and Cons

• Where it is good

- For all regression problems
- Better for two-class classification, possible for multi-class problems (not suggested)
- Various nonlinearity, strong representability

• Disadvantage

- Sequential process, inconvenient for parallel computing
- High computational complexity, not suitable for high-dimensional problems with sparse features

• A powerful extension : XGBoost

Conclusions

- Ensemble methods have integrable abilities of single models, achieving better performance
- Easy to generalize to new data
- When there are strong noises, easy to overfit
- Computationally intensive

对回归、二分类很有效，多分类可能可以（但不建议）

不能并行、复杂度较高。

Python Example

• Random forest :

```
from sklearn.ensemble import RandomForestClassifier
rf = RandomForestClassifier(n_estimators=100)
# RandomForestClassifier(bootstrap=True, class_weight=None,
#                       criterion='gini', max_depth=None, max_features='auto',
#                       max_leaf_nodes=None, min_impurity_split=1e-07,
#                       min_samples_leaf=1, min_samples_split=2,
#                       min_weight_fraction_leaf=0.0, n_estimators=100,
#                       n_jobs=1, oob_score=False, random_state=None,
#                       verbose=0, warm_start=False)
# Feature importance in random forest
feature_imp = pd.Series(rf.feature_importances_)
rf.fit(X_train, Y_train)
Y_predict_rf = rf.predict(X_test)
oob_error = 1 - rf.oob_score_
```

• AdaBoost :

```
from sklearn.ensemble import AdaBoostClassifier
adaboost = AdaBoostClassifier(n_estimators = 50)
adaboost.fit(X_train, Y_train)
adaboost.staged_predict(X_train)
Y_predict_ada = adaboost.predict(X_test)
```