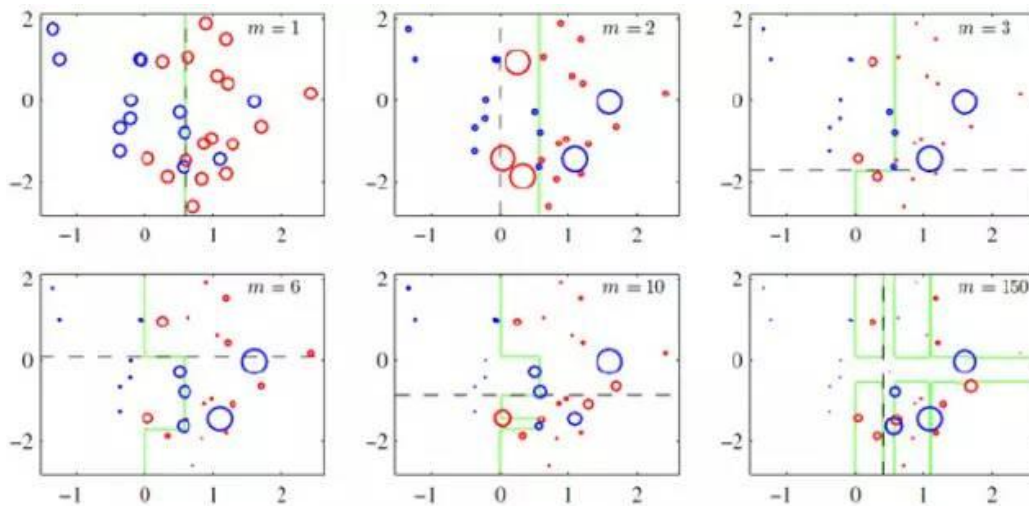


XGBoost

➤ 什麼是 XGBoost ?

1. 是由 Tianqi Chen <http://homes.cs.washington.edu/~tqchen/> 最初開發的實現可擴展，便攜，分佈式 gradient boosting (GBDT, GBRT or GBM) 算法的一個庫，可以下載安裝並應用於 C++，Python，R，Julia，Java，Scala，Hadoop，現在有很多協作者共同開發維護。
2. XGBoost 所應用的算法就是 gradient boosting decision tree，既可以用於分類也可以用於迴歸問題中。
3. Boosting
 - a. 基本思想：不同的訓練集是通過調整每個樣本對應的權重實現的，不同的權重對應不同的樣本分布，而這個權重為分類器不斷增加對錯樣本的重視程度。
 - b. 步驟：
 - i. 1. 首先賦予每個訓練樣本相同的初始化權重，在此訓練樣本分布下訓練出一個弱分類器；
 - ii. 2. 利用該弱分類器更新每個樣本的權重，分類錯誤的樣本認為是分類困難樣本，權重增加，反之權重降低，得到一個新的樣本分布；
 - iii. 3. 在新的樣本分布下，在訓練一個新的弱分類器，並且更新樣本權重，重複以上過程 T 次，得到 T 個弱分類器。
 - c. 通過改變樣本分布，使得分類器聚集在那些很難分的樣本上，對那些容易錯分的數據加強學習，增加錯分數據的權重。這樣錯分的數據再下一輪的疊代就有更大的作用（對錯分數據進行懲罰）。對於這些權重，一方面可以使用它們作為抽樣分布，進行對數據的抽樣；另一方面，可以使用權值學習有利於高權重樣本的分類器，把一個弱分類器提升為一個強分類器。
決策樹。
 - d. 對於 Boosting 來說，有兩個問題需要回答：一是在每一輪如何如何改變訓練數據的機率分布；二是如何將多個弱分類器組合成一個強分類器。



e. 上圖（圖片來自 prml p660）就是一個 Boosting 的過程，绿色的线表示目前取得的模型（模型是由前 m 次得到的模型合并得到的），虚线表示当前这次模型。每次分类的时候，会更关注分错的资料，上图中，红色和蓝色的点就是资料，点越大表示权重越高，看看右下角的图片，当 $m=150$ 的时候，获取的模型已经几乎能够将红色和蓝色的点区分开了。

f. 参考网址：<https://www.readhouse.net/articles/132582732/>

4. Gradient Boosting

- a. Gradient boosting 是 boosting 的其中一种方法
- b. Gradient boosting 就是通过加入新的弱学习器，来努力纠正前面所有弱学习器的残差，最终这样多个学习器相加在一起用来进行最终预测，准确率就会比单独的一个要高。之所以称为 Gradient，是因为在添加新模型时使用了梯度下降算法来最小化的损失。

5. 机器学习算法中 GBDT 和 XGBOOST 的区别有哪些

- a. 基分类器的选择：传统 GBDT 以 CART 作为基分类器，XGBoost 还支持线性分类器，这个时候 XGBoost 相当于带 L1 和 L2 正则化项的逻辑斯蒂回归（分类问题）或者线性回归（回归问题）。
- b. 二阶泰勒展开：传统 GBDT 在优化时只用到一阶导数信息，XGBoost 则对代价函数进行了二阶泰勒展开，同时用到了一阶和二阶导数。顺便提一下，XGBoost 工具支持自定义损失函数，只要函数可一阶和二阶求导。
- c. 方差-方差权衡：XGBoost 在目标函数里加入了正则项，用于控制模型的复杂度。正则项里包含了树的叶子节点个数 T 、每个叶子节点

上輸出分數的 L2 模的平方和。從 Bias-variance tradeoff 角度來講，正則項降低了模型的 variance，使學習出來的模型更加簡單，防止過擬合，這也是 XGBoost 優於傳統 GBDT 的一個特性。

d. 其他特點參考

<http://blog.csdn.net/chengfulukou/article/details/76906710>