



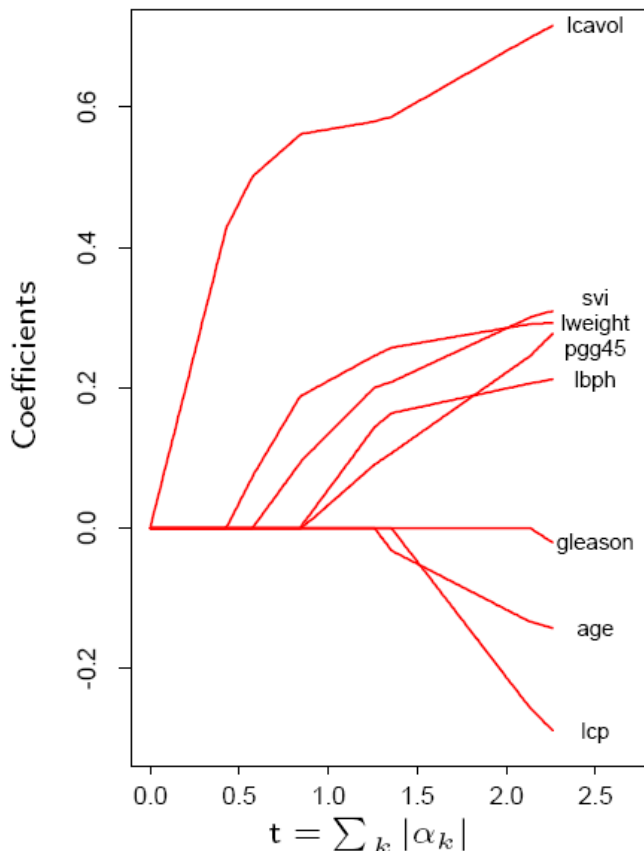
本不用调参数。第三，几乎不 Overfitting。我估计当时 Breiman 和 Friedman 肯定高兴坏了，因为眼看着他们提出的 CART 正在被 SVM 比下去的时候，AdaBoost 让决策树起死回生！Breiman 情不自禁地在他的论文里赞扬 AdaBoost 是最好的现货方法（off-the-shelf，即“拿下了就可以用”的意思）。其实在 90 年代末的时候，大家对 AdaBoost 为什么有如此神奇的性能迷惑不解。1999 年，Friedman 的一篇技术报告“Additive logistic regression: a statistical view of boosting”解释了大部分的疑惑（没有解释 AdaBoost 为什么不容易 Overfitting，这个问题好像至今还没有定论），即搞清楚了 AdaBoost 在优化什么指标以及如何优化的。基于此，Friedman 提出了他的 GBM（Gradient Boosting Machine，也叫 MART 或者 TreeNet）。几乎在同时，Breiman 另辟蹊径，结合他的 Bagging（Bootstrap aggregating）提出了 Random Forest（今天微软的 Kinect 里面就采用了 Random Forest，相关论文 Real-time Human Pose Recognition in Parts from Single Depth Images 是 CVPR2011 的 best paper）。

有一个关于 Gradient Boosting 细节不得不提。

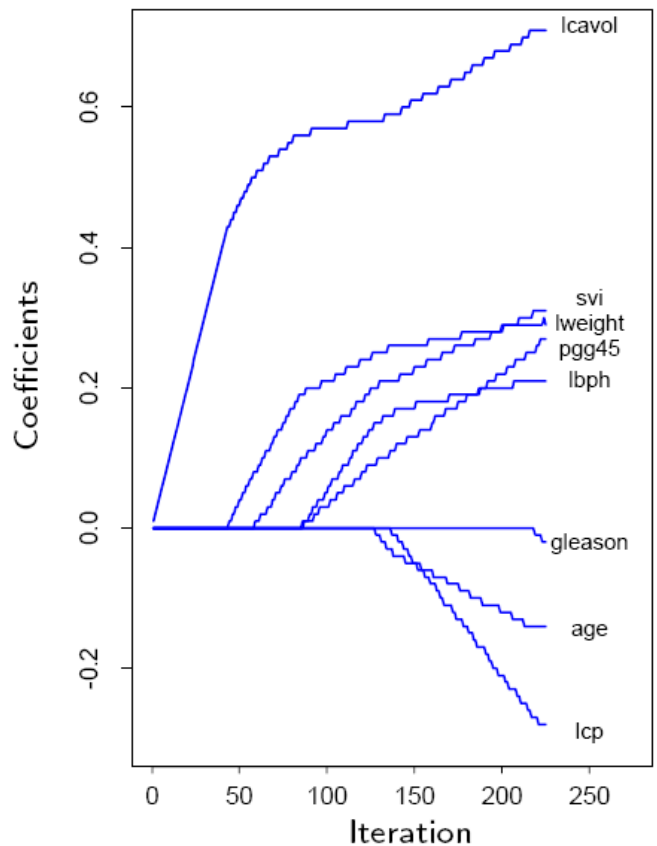
Friedman 在做实验的时候发现，把一棵新生成的决策树，记为  $f_m$ ，加到当前模型之前，在这棵决策树前乘以一个小的数，即  $v \times f_m$ （比如  $v=0.01$ ），再加入当前模型中，往往大大提高模型的准确度。他把这个叫做“Shrinkage”。接下来，Hastie, Tibshirani 和 Friedman 进一步发现（我发现大师们都是亲自动手写程序做实验的），如果把具有 Shrinkage 的 Gradient Boosting 应用到线性回归中时，得到的 Solution Path 与 Lasso 的 Solution Path 惊人地相似(如图所示)！他们把这一结果写在了 ESL 的第一版里，并推测这二者存在着某种紧密的联系，但精确的数学关系他们当时也不清楚。Tibshirani 说他们还请教了斯坦福的优化大师（我估计是 Stephen Boyd），但还是没有找到答案。

后来 Tibshirani 找到自己的恩师 Efron。Tibshirani 在“The Science of Bradley Efron”这本书的序言里写道，“**He sat down and pretty much single-handedly solved the problem. Along the way, he developed a new algorithm, ‘least angle regression,’ which is interesting in its own right, and sheds great statistical insight**

Lasso



Forward Stagewise



on the Lasso.” 我就不逐字逐句翻译了，大意是：Efron 独自摆平了这个问题，与此同时发明了“Least angle regression (LAR)” 。Efron 结论是 Lasso 和 Boosting 的确有很紧密的数学联系，它们都可以通过修改 LAR 得到。更令人惊叹的是 LAR 具有非常明确的几何意义。于是，Tibshirani 在序言中还有一句，“**In this work, Brad shows his great mathematical power - not the twentieth century, abstract kind of math, but the old-fashioned kind: geometric insight and analysis.**” 读 Prof Efron 的文章，可以感受到古典几何学与现代统计学的结合之美（推荐大家读读 Efron 教授 2010 年的一本新书 *Large-Scale Inference*）！总之，Efron 的这篇文章是现代统计学的里程碑，它结束了一个时代，开启了另一个时代。

这里，想补充说明一下 Lasso 的身世，它的全称是 The Least Absolute Shrinkage and Selection Operator，读音不是 [ˈlæso] 而是 [læˈsu:]，有中文翻译为“套索”，个人觉得这个翻译不好，太远离它本来的含义，不如就用 Lasso。Tibshirani 自己说他的 Lasso 是受到 Breiman 的 Non-Negative Garrote (NNG) 的启发。Lasso 把 NNG 的两步合并为一步，即 L1-norm regularization。Lasso 的巨大优势在于它所构造的模型是 Sparse 的，因为它会自动地选择很少一部分变量构造模型。现在，Lasso 已经家喻户晓了，但是 Lasso 出生后的头两年却很少有人问津。后来 Tibshirani 自己回忆时说，可能是由下面几个原因造成的：1. 速度问题：当时计算机求解 Lasso 的速度太慢；2. 理解问题：大家对 Lasso 模型的性质理解不够（直到 Efron 的 LAR 出来后大家才搞明白）；3. 需求问题：当时还没有遇到太多高维数据分析的问题，对 Sparsity 的需求似乎不足。Lasso 的遭遇似乎在阐释我们已经熟知的一些道理：1. 千里马常有，而伯乐不常有（没有 Efron 的 LAR，Lasso 可能很难有这么大的影响力）。2. 时势造英雄（高维数据分析的问题越来越多，比如 Bioinformatics 领域）。3. 金子总是会闪光的。

LAR 把 Lasso (L1-norm regularization) 和 Boosting 真正的联系起来，如同打通了任督二脉（数学细节可以参考本人的一个小结<sup>1</sup>，当然最好还是亲自拜读 Efron 的原著）。LAR 结束了一个晦涩的时代：在 LAR 之前，有关 Sparsity 的模型几乎都是一个黑箱，它们的数学性质（更不要谈古典的几何性质了）几乎都是缺失。LAR 开启了一个光明的时代：有关 Sparsity 的好文章如雨后春笋般地涌现，比如 Candes

和 Tao 的 Dantzig Selector。伯克利大学的 Bin Yu 教授称“Lasso, Boosting and Dantzig are three cousins”。近年来兴起的 Compressed sensing (Candes & Tao, Donoho) 也与 LAR 一脉相承，只是更加强调 L1-norm regularization 其他方面的数学性质，比如 Exact Recovery。我觉得这是一个问题的多个方面，Lasso 关注的是构建模型的准确性，Compressed sensing 关注的是变量选择的准确性。由此引起的关于 Sparsity 的研究，犹如黄河泛滥，一发不可收拾。比如 Low-rank 逼近是把 L1-norm 从向量到矩阵的自然推广（现在流行的“**用户推荐系统**”用到的 Collaborative filtering 的数学原理源于此）。有兴趣的童鞋可以参考我个人的小结<sup>2</sup>。

还必须提到的是算法问题。我个人觉得，一个好的模型，如果没有一个快速准确的算法作为支撑的话，它最后可能什么也不是。看看 Lasso 头几年的冷遇就知道了。LAR 的成功除了漂亮的几何性质之外，还有它的快速算法。LAR 的算法复杂度相当于最小二乘法的复杂度，这几乎已经把 Lasso 问题的求解推向极致。这一记录在 2007 年被 Friedman 的 Coordinate Descent (CD) 刷新，至今没人打破。Hastie 教授趣称这个为“FFT (Friedman + Fortran + Tricks)”。因为 CD 对 Generalized Lasso 问题并不能一网打尽，许多凸优化解法应运而生，如 Gradient Projection, Proximal methods, ADMM (Alternating Direction Method of Multipliers), (Split) Bregman methods, Nesterov's method (一阶梯度法中最优的收敛速度，Candes 的很多软件包都根据这个方法设计等等。哪个方法更好呢？这个就像问“谁的武功天下第一”一样。我只能回答“王重阳以后再也没有天下第一了，东邪西毒南帝北丐，他们各有各的所长，有的功夫是这个人擅长一些，而另外几门功夫又是另一个人更擅长一些”。有关 L1 的算法可能还会大量涌现，正如优化大师 Stephen Boyd 所说（2010 年 9 月 28 日）：“God knows the last thing we need is another algorithm for the Lasso.”

最后我想以讨论“模糊系统”和“统计学习”来结尾。这个话题非常具有争议，我就冒天下之大不讳吧，谈一谈我这几年的学习体会。记得十年前，立新老师曾经写过一篇文章《模糊系统：挑战与机遇并存——十年研究之感悟》，发表在 2001 年《自动化学报》上。我 2005 年看到的时候，敬仰之情，犹如滔滔江水。立新老师曾经有这么一句话：“If a method works well in practice, there must be some theoretical reasons for its success.” 2005 年的

<sup>1</sup> [http://ihome.ust.hk/~eeyang/lars\\_Lasso\\_boost.pdf](http://ihome.ust.hk/~eeyang/lars_Lasso_boost.pdf)

<sup>2</sup> [http://ihome.ust.hk/~eeyang/Learning\\_from\\_sparsity.pdf](http://ihome.ust.hk/~eeyang/Learning_from_sparsity.pdf)

时候，我开始问自己什么使模糊系统的成功？立新老师认为有如下几个原因：1. 模糊系统的通用逼近性能（Universal Approximator）；2. 模糊系统快速的构造算法，比如他自己的 WM 方法，ANFIS 等等；3. 结果的可解释性；4 利用各种不同形式的信息。

下面我谈谈自己的看法，第一，通用逼近性能当然是一个好的性质，它表明模糊系统是很 flexible 的，但 flexible 的结构太多了，比如神经网络。问题往往不在 flexible，而在太 flexible 导致 overfitting。就如同 SVM 一样，没有 L2-norm regularization，实践中的性能就会变得很差。第二，快速算法，这是好的方法必备的，SVM，Boosting，Random Forest 的算法都很快，而且可以直接用到高维，这一点上，我没有看到模糊系统的优势。第三，可解释性：模糊系统对低维数据（比如 2-4 维）的确具有好的解释性（因为 IF-THEN 规则的前提和结论都很简洁），但这个时候其它工具也可以做得到，比如 Gradient Boosting 和 Random Forests（很多例子可以在 ESL 这本书里看到）。第四，充分的利用各种信息。立新老师指的是 IF-THEN 规则可以比较自由灵活的加入先验知识，并在他的书里面详细给出实例。遗憾的是，这些例子都在处理低维空间的问题。如何用 IF-THEN 规则解构高维空间呢？我个人看不到它们特殊的优势。然而，在统计学习里，利用不同的先验知识处理高维空间的例子比比皆是，比如

Sparsity, group-structure, smoothness 等等。现在举一个 Gradient Boosting machine (GBM, 也叫 MART) 的例子来说明我的观点。根据 Lasso 和 Boosting 的关系，可以知道 GBM 已经用到了 Sparsity 的性质 (L1-norm regularization)。GBM 有两个参数可以反映我们的先验知识。第一个参数是深度 (depth)，控制每棵决策树的深度。如果深度为 1，即树桩结构 (Stump)，表明 GBM 将采用加法模型 (Generalized Additive model)，即不考虑变量之间的交互式作用 (Interaction)；如果深度大于 1，则考虑交互式作用。因为交互式作用在非线性的建模中比较重要，如异或 (XOR) 问题，没有考虑交互式作用将失败得很惨，所以这个参数设置反映了对非线性建模的先验。第二个参数是 Shrinkage 的大小。假设深度选取是合理的，在噪声比较小的时候，没有 Shrinkage 会比较好；噪声比较大的时候，有 Shrinkage 会好一些。实践中，使用 GBM 对高维数据分析，试错法 (Trial and error) 很容易使用，因为就这两个参数（通常 depth=3~4；实际数据的噪声往往比较大，推荐设置 Shrinkage=0.01）。模型构建好之后，GBM 会告诉你哪些变量是重要的，变量之间的交互式作用如何等等，这样模型的结果也是比较容易理解。Random Forests 也有相似的功能。好了，最后借 Hastie 教授的一幅图(见下页)来总结一下，无疑，GBM (MART) 是他们的最爱，也是我的最爱。

## 结束语

问：世间是否此山最高，或者另有高处比天高？

答： 在世间自有山比此山更高，Open-mind 要比天高。

Some characteristics of different learning methods.

Key: ●= good, ●=fair, and ●=poor.

| Characteristic                                     | Neural Nets | SVM | CART | GAM | KNN, kernels | MART |
|--|-------------|-----|------|-----|--------------|------|
| Natural handling of data of “mixed” type           | ●           | ●   | ●    | ●   | ●            | ●    |
| Handling of missing values                         | ●           | ●   | ●    | ●   | ●            | ●    |
| Robustness to outliers in input space              | ●           | ●   | ●    | ●   | ●            | ●    |
| Insensitive to monotone transformations of inputs  | ●           | ●   | ●    | ●   | ●            | ●    |
| Computational scalability (large $N$ )             | ●           | ●   | ●    | ●   | ●            | ●    |
| Ability to deal with irrelevant inputs             | ●           | ●   | ●    | ●   | ●            | ●    |
| Ability to extract linear combinations of features | ●           | ●   | ●    | ●   | ●            | ●    |
| Interpretability                                   | ●           | ●   | ●    | ●   | ●            | ●    |
| Predictive power                                   | ●           | ●   | ●    | ●   | ●            | ●    |