

# INFORMATION THEORY AND FLEXIBLE HIGH-DIMENSIONAL NON-LINEAR FUNCTION ESTIMATION

Andrew R. Barron

YALE UNIVERSITY  
DEPARTMENT OF STATISTICS

Presentation, November 12, 2011

at the Info-Metrics Institute, American University, Washington, DC

- Data, Model
- Combining Non-linearly Parameterized Terms
- Penalized Likelihood Criteria, Minimum Description Length
- Statistical Risk Determination
- Computation
- Adaptive Annealing
- Summary

- Data:  $(X_i, Y_i), i = 1, 2, \dots, n$
- Inputs: explanatory variables  $X_i$  in a unit cube in  $R^d$
- Random design: independent  $X_i \sim P$
- Output: response variable  $Y_i$  in  $R$
- Relationship:  $Y_i = f(X_i) + \epsilon_i$
- Noise:  $\epsilon_i$  independent  $N(0, \sigma^2)$
- Function:  $f$  unknown

# Non-linear Dictionaries

- Build functions  $f_m(x) = \sum_{j=1}^m c_j \phi_d(\theta_j, x)$  in the span of a dictionary  $\Phi = \{\phi_d(\theta, \cdot) : \theta \in \Theta\}$

- Product Bases

$$\phi_d(\theta, x) = \phi_1(\theta_1, x_1) \phi_1(\theta_2, x_2) \cdots \phi_1(\theta_d, x_d)$$

- Ridge Bases (as in projection pursuit regression)

$$\phi_d(\theta, x) = \phi_1(\theta_0 + \theta_1 x_1 + \theta_2 x_2 + \dots + \theta_d x_d)$$

- Examples of activation functions  $\phi(z) = \phi_1(z)$  for

- Perceptron networks:  $\phi(z) = 1_{\{z > 0\}}$
- Sigmoidal networks:  $e^z / (1 + e^z)$
- Sinusoidal models:  $\cos(z)$
- Hinging hyperplanes:  $(z)_+$
- Quadratic splines:  $1, z, (z)_+^2$
- Cubic splines:  $1, z, z^2, (z)_+^3$
- Polynomials:  $(z)^q$

- Response vector:  $Y = (Y_i)_{i=1}^n$  in  $R^n$
- Dictionary vectors:  $\Phi_{(n)} = \{(\phi_d(\theta, X_i))_{i=1}^n : \theta \in \Theta\}$
- Sample squared norm:  $\|f\|_{(n)}^2 = \frac{1}{n} \sum_{i=1}^n f^2(X_i)$
- Population squared norm:  $\|f\|^2 = \int f^2(x)P(dx)$
- Normalized dictionary condition:  $\|\phi\| \leq 1$  for  $\phi \in \Phi$

# Functions represented in span of non-linear dictionary

- Variation of  $f$  w.r.t.  $\Phi$

$$V_{\Phi}(f) = \|f\|_{\Phi} = \inf\{V : f/V \in \text{conv}(\pm\Phi)\}$$

- E.g.  $f(x) = \sum_j c_j \phi(\theta_j, x)$  has  $\|f\|_{\Phi} = \sum_j |c_j|$ , the  $\ell_1$  norm of the coefficients in representation of  $f$  in the span of  $\Phi$
- E.g.  $f(x) = \int e^{i\theta^T x} \tilde{f}(\theta) d\theta$  (Fourier representation)  
Then  $\|f\|_{\Phi}$  is given by an  $L_1$  spectral norm:

$$V_{\text{cos}}(f) = \int_{R^d} |\tilde{f}(\theta)| d\theta$$

$$V_{\text{step}}(f) = \int |\tilde{f}(\theta)| \|\theta\|_1 d\theta$$

$$V_{q\text{-spline}}(f) = \int |\tilde{f}(\theta)| \|\theta\|_1^{q+1} d\theta$$

# Penalized Likelihood, Minimum Description Length

- **Description Length** divided by sample size:  
 $\frac{1}{n} [ \log 1/\text{likelihood} + \text{Model Complexity Penalty} ]$
- **Control of number of terms:**

$$\frac{\|Y - f_m\|_{(n)}^2}{2\sigma^2} + \frac{m}{n} \log N_{\Phi, 1/n}$$

where the penalty is typically of order  $\frac{md}{n} \log n$

- **Control of the  $\ell_1$  norm of coefficients:**

$$\frac{\|Y - f\|_{(n)}^2}{2\sigma^2} + \lambda_n \|f\|_{\Phi}$$

$$\lambda_n = \sqrt{\frac{2 \log N_{\Phi, 1/n}}{n}}$$

- Optimize the above criteria to yield estimators  $\hat{f}$  and  $\hat{f}_{\hat{m}}$

# Statistical Risk Bounds

Bounds on the population accuracy of function estimates when the true function is  $f^*$

$$E\|\hat{f}_m - f^*\|^2 \leq \|f_m - f^*\|^2 + \frac{cmd}{n} \log n$$

$$E\|\hat{f}_{\hat{m}} - f^*\|^2 \leq \min_m \left\{ \|f_m - f^*\|^2 + \frac{cmd}{n} \log n \right\}$$

$$E\|\hat{f} - f^*\|^2 \leq \min_f \left\{ \|f - f^*\|^2 + \lambda_n \|f\|_\Phi \right\}$$

$$E\|\hat{f} - f^*\|^2 \leq \|f^*\|_\Phi \sqrt{\frac{2d \log n}{n}}$$



# Relaxed Greedy Algorithm and LASSO

- Initialize  $\hat{f}_0 = 0$ . For step  $k$ , have the previous fit  $\hat{f}_{k-1}$ .
- **Optimize the new term:** Maximize the inner product with the residuals  $res_i = Y_i - \hat{f}_{k-1}(X_i)$  to obtain the new  $\phi$  and its parameter vector  $\hat{\theta}_k$

$$\operatorname{argmax}_{\theta} \frac{1}{n} \sum_{i=1}^n res_i \phi(\theta, X_i)$$

- **Update the fit:**

$$\hat{f}_k = \alpha \hat{f}_{k-1} + \beta \phi$$

- **Obtain the coefficients  $\alpha, \beta$**  by either least squares or by  $\ell_1$  penalized least squares:
  - $\min \|Y - \alpha \hat{f}_{k-1} - \beta \phi_k\|^2$  for the relaxed greedy algorithm
  - $\min \|Y - \alpha \hat{f}_{k-1} - \beta \phi_k\|^2 + \lambda_n(|\alpha| V_{k-1} + |\beta|)$  for  $\ell_1$  pen pursuit (LASSO)

# Alternative: Forward Stepwise Regression or Orthogonal Matching Pursuit

- Initialize  $\hat{f}_0 = 0$ . For step  $k$ , have the previous fit  $\hat{f}_{k-1}$ .
- **Optimize the new term:** Maximize the inner product with the residuals  $res_i = Y_i - \hat{f}_{k-1}(X_i)$  to obtain the new  $\phi$  and its parameter vector  $\hat{\theta}_k$

$$\operatorname{argmax}_{\theta} \frac{1}{n} \sum_{i=1}^n res_i \phi(\theta, X_i)$$

- **Alternative fit update:**  $\hat{f}_k$  in  $\operatorname{span}\{\phi_1, \dots, \phi_{k-1}, \phi_k\}$
- with coefficients achieving

$$\min_{c_1, c_2, \dots, c_k} \left\| Y - \sum_{j=1}^k c_j \phi_j \right\|^2$$

## Choice of final $k$ :

- $k = m$  fixed, e.g.  $m$  equal a const multiple of  $\sqrt{n/(d \log n)}$
- $k = \hat{m}$  chosen by MDL with the relaxed greedy algorithm or forward stepwise regression
- May choose a larger final number of steps  $m$  between  $\hat{m}$  and  $n$ , for implementation of LASSO with control on closeness to the solution

- **Bound the accuracy of greedy computation** at step  $m$

For relaxed greedy and forward stepwise regression

$$\|Y - \hat{f}_m\|_{(n)}^2 \leq \inf_f \{ \|Y - f\|_{(n)}^2 + \frac{4\|f\|_{\Phi}^2}{m} \}$$

# Computation Bound ( $\ell_1$ penalized case)

- **Bound on accuracy of  $\ell_1$  penalized optimization** at step  $m$

For  $\ell_1$  penalized greedy pursuit (implementation of LASSO)

$$[\|Y - \hat{f}_m\|_{(n)}^2 + \lambda \|\hat{f}_m\|_{\Phi}] \leq \inf_f \left\{ [\|Y - f\|_{(n)}^2 + \lambda \|f\|_{\Phi}] + \frac{4\|f\|_{\Phi}^2}{m} \right\}$$

# Computation bound with rough choices of $\phi$

- Accuracy of computation with rough choice of  $\phi$  each step
- Choose  $\phi$  to achieve  $\frac{1}{n} \sum_{i=1}^n \text{res}_i \phi(X_i)$  at least  $(1/C)J_{max}$
- For relaxed greedy and forward stepwise regression

$$\|Y - \hat{f}_m\|_{(n)}^2 \leq \inf_f \{ \|Y - f\|_{(n)}^2 + \frac{4C^2 \|f\|_{\Phi}^2}{m} \}$$

- $C > 1$  is the approximate optimization factor (e.g.  $C = 2$ )

# Non-linear Optimization Step

- maximize

$$J(\theta) = \frac{1}{n} \sum_{i=1}^n r_i \phi(\theta, \mathbf{x}_i)$$

- There may be exponentially many peaks for  $\theta$  in  $R^{d+1}$
- Exact optimization is NP hard for certain dictionaries  $\Phi$  (e.g. the neural net case with the step activation function)
- Seek Choices of flexible non-linear dictionaries  $\Phi = \{\phi(\theta, \mathbf{x})\}$  for which optimization to within a constant factor is possible by stochastic search strategies

# Non-linear Optimization Step

- maximize

$$J(\theta) = \frac{1}{n} \sum_{i=1}^n r_i \phi(\theta, \mathbf{x}_i)$$

- Seek Choices of flexible non-linear dictionaries  $\Phi = \{\phi(\theta, \mathbf{x})\}$  for which optimization to within a constant factor is possible by stochastic search strategies
- Try running a Markov Chain, initialized with  $\theta \sim p_0(\theta)$  or  $p_\epsilon(\theta)$ , targeting having long-run distribution

$$p_\gamma(\theta) = \frac{1}{c_\gamma} \exp\{\gamma J(\theta)\}$$

- Gain  $\gamma$  of order  $d \log d$  would be sufficient for outcomes with  $J(\theta) \geq (1/2)J_{max}$  with high probability



- FAILURE in high-dimensions of methods that rely on transitions designed for invariance
  - Metropolis-Hastings
  - Simulated Annealing
  - Diffusion with gradient drift

$$d\theta(t) = \text{Drift}_t(\theta(t))dt + dB(t)$$

$$\theta(t + \delta) = \theta(t) + \text{Drift}_t(\theta(t))\delta + Z(t)\sqrt{\delta}$$

- Gradient drift

$$\text{Drift}_t(\theta) = \frac{1}{2}\nabla \log p_\gamma(\theta) = \frac{\gamma}{2}\nabla J(\theta)$$

- Time until distribution is near  $p_\gamma(\theta)$  is exponential in  $\gamma \times \text{depth}_J$

# Optimization by Adaptive Annealing

- **SUCCESS** for certain  $\Phi$  of Adaptive Annealing
- Modify the Markov Chain so that the distribution tracks  $p_{\gamma_t}(\theta)$  with increasing  $\gamma_t$ 
  - A stochastic diffusion with modified drift accomplishes the desired evolution

$$d\theta(t) = \text{Drift}_t(\theta(t))dt + dB(t)$$

$$\theta(t + \delta) = \theta(t) + \text{Drift}_t(\theta(t))\delta + Z(t)\sqrt{\delta}$$

- The modified drift is a local gradient plus a simple global change function

$$\text{Drift}_t(\theta) = \frac{\gamma}{2}\nabla J(\theta) + \text{change}_t(\theta),$$

where we may set  $\text{change}_t(\theta) = a_t \theta$

- Starting from  $\gamma_0 = \epsilon$ , it tracks

$$p_{\gamma_t}(\theta) = (1/c_{\gamma_t}) \exp\{\gamma_t J(\theta)\}$$

# Optimization by Adaptive Annealing

- Adaptive Annealing tracks  $p_{\gamma_t}(\theta)$  with increasing  $\gamma_t$
- Stochastic diffusion using drift of the form

$$Drift_t(\theta) = \frac{\gamma}{2} \nabla J(\theta) + change_t(\theta)$$

solves the Kolmogorov, Fokker-Planck PDE governing the relationship between the drift and the desired evolution of the marginal density of the state  $\theta$ ,

$$\frac{\partial}{\partial t} p_{\gamma_t}(\theta) = -\nabla^T (Drift_t(\theta) p_{\gamma_t}(\theta)) + \frac{1}{2} \nabla^T \nabla p_{\gamma_t}(\theta)$$

when  $\phi(\theta, x)$  is a ridge polynomial or ridge spline

$$\phi(\theta, x) = \phi(\theta_0 + \theta_1 x_1 + \dots + \theta_d x_d)$$

with  $\phi(z)$  equal to  $(z)^q$  or  $(z)_+^q$  with  $q \geq 2$ .

- May set  $change_t(\theta) = a_t \theta$  with  $a_t = (\log \gamma_t)' / q$ .

# Optimization by Adaptive Annealing

Clarification:

- For normalizeability, e.g. with  $q = 2$ , may use

$$p_\gamma(\theta) = \frac{1}{c_\gamma} \exp \left\{ \gamma \left[ \frac{1}{n} \sum r_i (\theta^T X_i)_+^2 - \lambda \|\theta\|^2 \right] \right\}$$

- Both  $(\theta^T X_i)_+^2$  and  $\|\theta\|^2$  have the property that they are recovered by taking the inner product of  $\theta$  with their gradient. Likewise

$$\theta^T \nabla J(\theta) = q J(\theta)$$

for

$$J(\theta) = \frac{1}{n} \sum_{i=1}^n r_i (\theta^T X_i)_+^q - \lambda \|\theta\|^q$$

- These identities determine the suitability of a multiple of  $\theta$  as an ingredient in the drift to solve the PDE.

# Optimization by Adaptive Annealing

Refinement:

- More general change functions of the form

$$\text{change}_t(\theta) = a_t \theta + G_t(\theta) / p_{\gamma_t}(\theta)$$

where

$$\nabla G_t(\theta) = c_t p_{\gamma_t}(\theta)$$

also solve the PDE.

- For instance the following is an acceptable choice

$$G_t(\theta) = \int_0^{\theta_0} p_{\gamma_t}(\tilde{\theta}_0, \theta_1, \dots, \theta_d) d\theta_0$$

expressible as a sum of one-dimensional Gaussian integrals.

- These more general solutions provide more freedom in setting  $\gamma_t$  with favorable growth properties.

- Ridge splines with adaptive annealing
- Plus an information-theoretic criterion based on  $\ell_1$  penalized greedy pursuit
- Provides flexible high-dimensional function estimation that is fast and accurate