

# Overview of Recent Developments in Penalized Likelihood and MDL

Andrew R. Barron

YALE UNIVERSITY

Collaborators

Penal. Like: **Sabyasachi Chatterjee, Cong Huang, Jonathan Li, Xi Luo**

NML vs Bayes: **Teemu Roos, Kazuho Watanabe**

Workshop on Information Theory Methods in Science & Engineering  
WITMSE, 5 July 2014, Wakiki, Hawaii

# Recent Developments in Penalized Likelihood and MDL

Andrew Barron

# Outline

- 1 Some Principles of Information Theory and Statistics
  - Data Compression
  - Minimum Description Length Principle
  - Two-stage Codes, Mixture Codes, Normalized Max Like.
- 2 Penalized Likelihood
  - Redundancy, Resolvability, and Risk
  - Risk-Valid Penalties
  - Penalized Likelihood Analysis for Continuous Parameters
  - $\ell_0$  and  $\ell_1$  Penalties are Codelength-Valid and Risk-Valid
- 3 Summary

# Outline

- 1 Some Principles of Information Theory and Statistics
  - Data Compression
  - Minimum Description Length Principle
  - Two-stage Codes, Mixture Codes, Normalized Max Like.
- 2 Penalized Likelihood
  - Redundancy, Resolvability, and Risk
  - Risk-Valid Penalties
  - Penalized Likelihood Analysis for Continuous Parameters
  - $\ell_0$  and  $\ell_1$  Penalties are Codelength-Valid and Risk-Valid
- 3 Summary

# Shannon, Kraft, McMillan

- Characterization of uniquely decodeable codelengths

$$L(\underline{x}), \quad \underline{x} \in \underline{\mathcal{X}}, \quad \sum_{\underline{x}} 2^{-L(\underline{x})} \leq 1$$

$$L(\underline{x}) = \log 1/q(\underline{x}) \quad q(\underline{x}) = 2^{-L(\underline{x})}$$

- Operational meaning of probability:

A distribution is given by a choice of code

# Codelength Comparison

- Targets  $p$  are possible distributions
- Compare codelength  $\log 1/q(\underline{x})$  to targets  $\log 1/p(\underline{x})$
- Redundancy or regret

$$\left[ \log 1/q(\underline{x}) - \log 1/p(\underline{x}) \right]$$

- Expected redundancy

$$D(P_{\underline{X}} \| Q_{\underline{X}}) = E_P \left[ \log \frac{p(\underline{X})}{q(\underline{X})} \right]$$

# Universal Codes

- MODELS

Family of coding strategies  $\Leftrightarrow$  Family of prob. distributions

Indexed by parameters or functions:

$$\{L_{\theta}(\underline{x}) : \theta \in \Theta\} \Leftrightarrow \{p_{\theta}(\underline{x}) : \theta \in \Theta\}$$

$$\{L_f(\underline{x}) : f \in \mathcal{F}\} \Leftrightarrow \{p_f(\underline{x}) : f \in \mathcal{F}\}$$

- Universal codes  $\Leftrightarrow$  Universal probabilities  $q(\underline{x})$

$$L(\underline{x}) = \log 1/q(\underline{x})$$

- Redundancy:  $[\log 1/q(\underline{x}) - \log 1/p_{\theta}(\underline{x})]$

Want it small either uniformly in  $\underline{x}, \theta$  or in expectation

# Statistical Aim

- Training data  $\underline{x}$   $\Rightarrow$  estimator  $\hat{p} = p_{\hat{\theta}}$
- Subsequent data  $\underline{x}'$
- Want  $\log 1/\hat{p}(\underline{x}')$  to compare favorably to  $\log 1/p(\underline{x}')$
- For targets  $p$  close to or in the families



# Loss

- Kullback Information-divergence:

$$D(P_{\underline{X}'} \| Q_{\underline{X}'}) = E \left[ \log p(\underline{X}') / q(\underline{X}') \right]$$

- Bhattacharyya, Hellinger, Chernoff, Rényi divergence:

$$d(P_{\underline{X}'}, Q_{\underline{X}'}) = 2 \log 1 / E[q(\underline{X}') / p(\underline{X}')]^{1/2}$$

- Product model case:  $p(\underline{x}') = \prod_{i=1}^n p(x'_i)$

$$D(P_{\underline{X}'} \| Q_{\underline{X}'}) = n D(P \| Q)$$

Likewise

$$d(P_{\underline{X}'}, Q_{\underline{X}'}) = n d(P, Q)$$

# Minimum Description Length (MDL) Principle

- Universal coding brought into statistical play
- Minimum Description Length Principle:

The shortest code for data gives the best statistical model

## MDL: Two-stage Version

- Two-stage codelength:

$$L(\underline{x}) = \min_{\theta \in \Theta} \left[ \log 1/p_{\theta}(\underline{x}) + L(\theta) \right]$$

bits for  $\underline{x}$  given  $\theta$  + bits for  $\theta$

- The corresponding statistical estimator is  $\hat{p} = p_{\hat{\theta}}$
- Typically in  $d$ -dimensional families  $L(\theta)$  is of the form

$$\frac{d}{2} \log n + C_d(\theta)$$

## MDL: Mixture Versions

- Codelength based on a Bayes mixture

$$L(\underline{x}) = \log \frac{1}{q(\underline{x})}$$

where

$$q(\underline{x}) = \int p(\underline{x}|\theta)w(\theta)d\theta \text{ or } \sum_{\theta} p(\underline{x}|\theta)w(\theta)$$

minimax optimal with least favorable  $w$  (capacity achieving)

- Codelength approximation (Barron 1985, Clarke and Barron 1990,1994)

$$\log \frac{1}{p(\underline{x}|\hat{\theta})} + \frac{d}{2} \log \frac{n}{2\pi} + \log \frac{|\hat{I}(\hat{\theta})|^{1/2}}{w(\hat{\theta})}$$

where  $\hat{I}(\hat{\theta})$  is the empirical Fisher Information at the MLE

## MDL: Mixture Versions

- Codelength based on a Bayes mixture

$$q(\underline{x}) = \int p(\underline{x}|\theta)w(\theta)d\theta$$

- Corresponding statistical estimator is the predictive distrib

$$\hat{p}(\underline{x}') = q(\underline{x}'|\underline{x}) = \frac{\int p(\underline{x}'|\theta)p(\underline{x}|\theta)w(\theta)d\theta}{\int p(\underline{x}|\theta)w(\theta)d\theta}$$

- It has a clean relative entropy risk bound

$$ED(P\|\hat{P}) \leq \text{Resolvability}$$

$$= \text{Kullback Approx Error} + \log 1/(\text{Posterior prob of neigh.})$$

# NML Version

- Codelength via normalized maximum likelihood (NML)

$$NML(\underline{x}) = \frac{\max_{\theta} p(\underline{x}|\theta)}{C_{Shtarkov}}$$

- **minimax optimal for pointwise regret**
- It has the same codelength approximation as Bayes with asymp. least favorable prior (Takeuchi & B. 1998, B., Rissanen, Yu 1998)
- Is NML exactly Bayes in finite samples? (B., Roos, Watanabe 2014)

# NML Version

- Codelength via normalized maximum likelihood (NML)

$$NML(\underline{x}) = \frac{\max_{\theta} p(\underline{x}|\theta)}{C_{Shtarkov}}$$

- minimax optimal for pointwise regret
- It has the same codelength approximation as Bayes with asymp. least favorable prior (Takeuchi & B. 1998, B., Rissanen, Yu 1998)
- Is NML exactly Bayes in finite samples? (B., Roos, Watanabe 2014)

# NML Version

- Codelength via normalized maximum likelihood (NML)

$$NML(\underline{x}) = \frac{\max_{\theta} p(\underline{x}|\theta)}{C_{Shtarkov}}$$

- minimax optimal for pointwise regret
- It has the same codelength approximation as Bayes with asymp. least favorable prior (Takeuchi & B. 1998, B., Rissanen, Yu 1998)
- **Is NML exactly Bayes in finite samples?** (B., Roos, Watanabe 2014)



# NML Version

- Codelength via normalized maximum likelihood (NML)

$$NML(\underline{x}) = \frac{\max_{\theta} p(\underline{x}|\theta)}{C_{Shtarkov}}$$

- minimax optimal for pointwise regret
- It has the same codelength approximation as Bayes with asymp. least favorable prior (Takeuchi & B. 1998, B., Rissanen, Yu 1998)
- **Is NML exactly Bayes in finite samples?** (B., Roos, Watanabe 2014)

Yes, with sufficiently many linearly independent  $p(\cdot|\theta)$ ,  
[though in esoteric cases some weights are negative]:

$$NML(\underline{x}) = \sum_{\theta} p(\underline{x}|\theta) w(\theta)$$

# NML Version

- Codelength via normalized maximum likelihood (NML)

$$NML(\underline{x}) = \frac{\max_{\theta} p(\underline{x}|\theta)}{C_{Shtarkov}}$$

- minimax optimal for pointwise regret
- It has the same codelength approximation as Bayes with asymp. least favorable prior (Takeuchi & B. 1998, B., Rissanen, Yu 1998)
- **NML is exactly Bayes in finite samples** (B., Roos, Watanabe 2014)

$$NML(\underline{x}) = \sum_{\theta} p(\underline{x}|\theta) w(\theta)$$

Allows fast computation of marginals and predictive distrib.

# Outline

- 1 Some Principles of Information Theory and Statistics
  - Data Compression
  - Minimum Description Length Principle
  - Two-stage Codes, Mixture Codes, Normalized Max Like.
- 2 Penalized Likelihood
  - Redundancy, Resolvability, and Risk
  - Risk-Valid Penalties
  - Penalized Likelihood Analysis for Continuous Parameters
  - $\ell_0$  and  $\ell_1$  Penalties are Codelength-Valid and Risk-Valid
- 3 Summary

# Penalized Likelihood

## Penalized likelihood estimators

$$\hat{\theta} = \operatorname{argmin}_{\theta} \{ \log 1/p(\underline{x}|\theta) + \operatorname{pen}(\theta) \}$$

where the penalty  $\operatorname{pen}(\theta)$  arises from

- prior density:  $\operatorname{pen}(\theta) = \log 1/w(\theta)$
- codelength:  $\operatorname{pen}(\theta) = L(\theta)$
- dimensionality:  $\operatorname{pen}(\theta) = \lambda_n \operatorname{dim}(\Theta)$
- sparsity control:  $\operatorname{pen}(\theta) = \lambda_n \|\theta\|_1$
- roughness penalty:  $\operatorname{pen}(\theta) = \text{norm of derivative}$
- maximum likelihood:  $\operatorname{pen}(\theta) = \text{constant}$

# Penalized Likelihood

## Penalized likelihood estimators

$$\hat{\theta} = \operatorname{argmin}_{\theta} \{ \log 1/p(\underline{x}|\theta) + \operatorname{pen}(\theta) \}$$

where the penalty  $\operatorname{pen}(\theta)$  arises from

- prior density:  $\operatorname{pen}(\theta) = \log 1/w(\theta)$
- codelength:  $\operatorname{pen}(\theta) = L(\theta)$
- dimensionality:  $\operatorname{pen}(\theta) = \lambda_n \operatorname{dim}(\Theta)$
- sparsity control:  $\operatorname{pen}(\theta) = \lambda_n \|\theta\|_1$
- roughness penalty:  $\operatorname{pen}(\theta) = \text{norm of derivative}$
- maximum likelihood:  $\operatorname{pen}(\theta) = \text{constant}$

MDL analysis reveals which size penalties are valid  
for good prediction and compression properties.

## Two-stage Code Redundancy

- Two-stage codes: the start of penalized likelihood analysis
- Expected codelength minus target at  $p_{\theta^*}$

$$\text{Redundancy} = E \left[ \min_{\theta \in \Theta} \left\{ \log \frac{1}{p_{\theta}(\underline{X})} + L(\theta) \right\} - \log \frac{1}{p_{\theta^*}(\underline{X})} \right]$$

- Redundancy approx in smooth families of dimension  $d$

$$\frac{d}{2} \log n + C_d(\theta)$$

## Redundancy and Resolvability

- Redundancy =  $E \min_{\theta \in \Theta} \left[ \log \frac{p_{\theta^*}(\underline{x})}{p_{\theta}(\underline{x})} + L(\theta) \right]$
- Resolvability =  $\min_{\theta \in \Theta} E \left[ \log \frac{p_{\theta^*}(\underline{x})}{p_{\theta}(\underline{x})} + L(\theta) \right]$   
=  $\min_{\theta \in \Theta} \left[ D(P_{\underline{X}|\theta^*} \| P_{\underline{X}|\theta}) + L(\theta) \right]$
- Ideal tradeoff of Kullback approximation error & complexity
- Population analogue of the two-stage code MDL criterion
- Divide by  $n$  to express as a rate. In the i.i.d. case

$$R_n(\theta^*) = \min_{\theta \in \Theta} \left[ D(\theta^* \| \theta) + \frac{L(\theta)}{n} \right]$$

## Risk of Estimator based on Two-stage Code

- Estimator  $\hat{\theta}$  is the choice achieving the minimization

$$\min_{\theta \in \Theta} \left\{ \log \frac{1}{p_{\theta}(\underline{x})} + \mathcal{L}(\theta) \right\}$$

- Codelengths for  $\theta$  are  $\mathcal{L}(\theta) = 2L(\theta)$  with  $\sum_{\theta \in \Theta} 2^{-L(\theta)} \leq 1$ .
- Total loss  $d_n(\theta^*, \hat{\theta})$  with  $d_n(\theta^*, \theta) = d(P_{\underline{X}'|\theta^*}, P_{\underline{X}'|\theta})$

$$\text{Risk} = E[d_n(\theta^*, \hat{\theta})]$$

- Info-Thy bound on risk: (Barron 1985, Barron and Cover 1991, Jonathan Li 1999)

$$\text{Risk} \leq \text{Redundancy} \leq \text{Resolvability}$$



## Risk of Estimator based on Two-stage Code

- Estimator  $\hat{\theta}$  is the choice achieving the minimization

$$\min_{\theta \in \Theta} \left\{ \log \frac{1}{p_{\theta}(\underline{x})} + \mathcal{L}(\theta) \right\}$$

- Codelengths for  $\theta$  are  $\mathcal{L}(\theta) = 2L(\theta)$  with  $\sum_{\theta \in \Theta} 2^{-L(\theta)} \leq 1$ .
- Total loss  $d_n(\theta^*, \hat{\theta})$  with  $d_n(\theta^*, \theta) = d(P_{\underline{X}'|\theta^*}, P_{\underline{X}'|\theta})$

$$\text{Risk} = E[d_n(\theta^*, \hat{\theta})]$$

- Info-Thy bound on risk: (Barron 1985, Barron and Cover 1991, Jonathan Li 1999)

$$\text{Risk} \leq \text{Redundancy} \leq \text{Resolvability}$$

- Drawback: Two-part code interpretation needs countable  $\Theta$

## Key to Risk Analysis (in the countable $\Theta$ case)

- log likelihood-ratio discrepancy at training  $\underline{x}$  and future  $\underline{x}'$

$$\left[ \log \frac{p_{\theta^*}(\underline{x})}{p_{\theta}(\underline{x})} - d_n(\theta^*, \theta) \right]$$

- Proof shows, for  $L(\theta) = Lcal(\theta)/2$  satisfying Kraft, that

$$\min_{\theta \in \Theta} \left\{ \left[ \log \frac{p_{\theta^*}(\underline{x})}{p_{\theta}(\underline{x})} - d_n(\theta^*, \theta) \right] + \mathcal{L}(\theta) \right\}$$

has expectation  $\geq 0$ . From which the risk bound follows.

# Codelength Valid Penalties (for uncountable $\Theta$ )

- Penalized Likelihood

$$\min_{\theta \in \Theta} \left\{ \log \frac{1}{p_{\theta}(\underline{x})} + \text{Pen}(\theta) \right\}$$

- Possibly uncountable  $\Theta$
- It is still a codelength if there exists a countable  $\tilde{\Theta}$  and  $L$  satisfying Kraft such that the above is not less than

$$\min_{\tilde{\theta} \in \tilde{\Theta}} \left\{ \log \frac{1}{p_{\tilde{\theta}}(\underline{x})} + L(\tilde{\theta}) \right\}$$

## Codelength Valid Penalties (for uncountable $\Theta$ )

- Equivalently,  $Pen(\theta)$  is valid for penalized likelihood to be a **codelength** if there is such a countable  $\tilde{\Theta}$  and Kraft summable  $L(\tilde{\theta})$ , such that, for every  $\theta$  in  $\Theta$ , there is a representor  $\tilde{\theta}$  in  $\tilde{\Theta}$  such that

$$Pen(\theta) \geq L(\tilde{\theta}) + \log \frac{p_{\theta}(\underline{x})}{p_{\tilde{\theta}}(\underline{x})}$$

- This is the link between uncountable and countable cases

# Statistical-Risk Valid Penalties

- Penalized Likelihood

$$\hat{\theta} = \operatorname{argmin}_{\theta \in \Theta} \left\{ \log \frac{1}{p_{\theta}(\underline{X})} + \operatorname{Pen}(\theta) \right\}$$

- Again: possibly uncountable  $\Theta$
- Task: determine a condition on  $\operatorname{Pen}(\theta)$  such that the risk is captured by the population analogue

$$Ed_n(\theta^*, \hat{\theta}) \leq \inf_{\theta \in \Theta} \left\{ E \log \frac{p_{\theta^*}(\underline{X})}{p_{\theta}(\underline{X})} + \operatorname{Pen}(\theta) \right\}$$

## Statistical-Risk Valid Penalty

- For an uncountable  $\Theta$  and a penalty  $Pen(\theta)$ ,  $\theta \in \Theta$ , suppose there is a countable  $\tilde{\Theta}$  and  $\mathcal{L}(\tilde{\theta}) = 2L(\tilde{\theta})$  where  $L(\tilde{\theta})$  satisfies Kraft, such that, for all  $\underline{x}, \theta^*$ ,

$$\begin{aligned} & \min_{\theta \in \Theta} \left\{ \left[ \log \frac{p_{\theta^*}(\underline{x})}{p_{\theta}(\underline{x})} - d_n(\theta^*, \theta) \right] + Pen(\theta) \right\} \\ & \geq \min_{\tilde{\theta} \in \tilde{\Theta}} \left\{ \left[ \log \frac{p_{\theta^*}(\underline{x})}{p_{\tilde{\theta}}(\underline{x})} - d_n(\theta^*, \tilde{\theta}) \right] + \mathcal{L}(\tilde{\theta}) \right\} \end{aligned}$$

- Proof of the risk conclusion:  
The second expression has expectation  $\geq 0$ ,  
so the first expression does too.
- This condition and result is obtained with J. Li and X. Luo (in Rissanen Festschrift 2008)

## Variable Complexity, Variable Distortion Cover

- **Equivalent statement of the condition:**  $Pen(\theta)$  is a valid penalty if for each  $\theta$  in  $\Theta$  there is a representer  $\tilde{\theta}$  in  $\tilde{\Theta}$  with complexity  $L(\tilde{\theta})$ , distortion  $\Delta_n(\tilde{\theta}, \theta)$  and

$$Pen(\theta) \geq \mathcal{L}(\tilde{\theta}) + \Delta_n(\tilde{\theta}, \theta)$$

where the distortion  $\Delta_n(\tilde{\theta}, \theta)$  is the difference in the discrepancies at  $\tilde{\theta}$  and  $\theta$

$$\Delta_n(\tilde{\theta}, \theta) = \log \frac{p_\theta(\underline{x})}{p_{\tilde{\theta}}(\underline{x})} + d_n(\theta, \theta^*) - d_n(\tilde{\theta}, \theta^*)$$

# A Setting for Regression and log-density Estimation: Linear Span of a Dictionary

- $\mathcal{G}$  is a dictionary of candidate basis functions  
E.g. wavelets, splines, polynomials, trigonometric terms, sigmoids, explanatory variables and their interactions
- Candidate functions in the linear span

$$f_{\theta}(x) = \sum_{g \in \mathcal{G}} \theta_g g(x)$$

- weighted  $\ell_1$  norm of coefficients

$$\|\theta\|_1 = \sum_g a_g |\theta_g|$$

- weights  $a_g = \|g\|_n$  where  $\|g\|_n^2 = \frac{1}{n} \sum_{i=1}^n g^2(x_i)$



## Example Models

- Regression (focus of current presentation)

$$p_{\theta}(y|x) = \text{Normal}(f_{\theta}(x), \sigma^2)$$

- Logistic regression with  $y \in \{0, 1\}$

$$p_{\theta}(y|x) = \text{Logistic}(f_{\theta}(x)) \quad \text{for } y = 1$$

- Log-density estimation (focus of Festschrift paper)

$$p_{\theta}(x) = \frac{p_0(x) \exp\{f_{\theta}(x)\}}{C_f}$$

- Gaussian graphical models

# $\ell_1$ Penalty

- $pen(\theta) = \lambda \|\theta\|_1$  where  $\theta$  are coeff of  $f_\theta(x) = \sum_{g \in \mathcal{G}} \theta_g g(x)$

## Regression with $\ell_1$ penalty

- $\ell_1$  penalized log-density estimation, i.i.d. case

$$\hat{\theta} = \operatorname{argmin}_{\theta} \left\{ \frac{1}{n} \log \frac{1}{p_{f_{\theta}}(\underline{x})} + \lambda_n \|\theta\|_1 \right\}$$

- Regression with Gaussian model, fixed  $\sigma^2$

$$\min_{\theta} \left\{ \frac{1}{2\sigma^2} \frac{1}{n} \sum_{i=1}^n (Y_i - f_{\theta}(x_i))^2 + \frac{1}{2} \log 2\pi\sigma^2 + \frac{\lambda_n}{\sigma} \|\theta\|_1 \right\}$$

- Valid for

$$\lambda_n \geq \sqrt{\frac{2 \log(2M_{\mathcal{G}})}{n}} \quad \text{with } M_{\mathcal{G}} = \operatorname{Card}(\mathcal{G})$$

## Adaptive risk bound

- For log density estimation with suitable  $\lambda_n$

$$Ed(f^*, f_{\hat{\theta}}) \leq \inf_{\theta} \left\{ D(f^* \| f_{\theta}) + \lambda_n \|\theta\|_1 \right\}$$

- For regression with fixed design points  $x_j$ , fixed  $\sigma$ , and

$$\lambda_n = \sqrt{\frac{2 \log(2M)}{n}},$$

$$\frac{E \|f^* - f_{\hat{\theta}}\|_n^2}{4\sigma^2} \leq \inf_{\theta} \left\{ \frac{\|f^* - f_{\theta}\|_n^2}{2\sigma^2} + \frac{\lambda_n}{\sigma} \|\theta\|_1 \right\}$$

## Adaptive risk bound specialized to regression

- Again for fixed design and  $\lambda_n = \sqrt{\frac{2 \log 2M}{n}}$ , multiplying through by  $4\sigma^2$ ,

$$E \|f^* - f_{\hat{\theta}}\|_n^2 \leq \inf_{\theta} \left\{ 2 \|f^* - f_{\theta}\|_n^2 + 4\sigma \lambda_n \|\theta\|_1 \right\}$$

- In particular for all targets  $f^* = f_{\theta^*}$  with finite  $\|\theta^*\|$  the risk bound  $4\sigma \lambda_n \|\theta^*\|$  is of order  $\sqrt{\frac{\log M}{n}}$
- Details in Barron, Luo (proceedings Workshop on Information Theory Methods in Science & Eng. 2008), Tampere, Finland

## Comments on proof

- Likelihood discrepancy plus complexity

$$\frac{1}{2\sigma^2} \sum_{i=1}^n (Y_i - f_{\tilde{\theta}}(x_i))^2 - \frac{1}{2\sigma^2} \sum_{i=1}^n (Y_i - f_{\theta}(x_i))^2 + K \log(2M)$$

- Representer  $f_{\tilde{\theta}}$  of  $f_{\theta}$  of following form, with  $v$  near  $\|\theta\|_1$

$$f_{\tilde{\theta}}(x) = \frac{v}{K} \sum_{k=1}^K g_k(x) / \|g_k\|$$

- $g_1, \dots, g_K$  picked at random from  $\mathcal{G}$ , independently, where  $g$  arises with probability proportional to  $|\theta_g| a_g$
- Shows exists representer with like. discrep. + complexity

$$\frac{nv^2}{2K} + K \log(2M)$$

## Comments on proof

- Optimizing it yields penalty proportional to  $\ell_1$  norm
- Penalty  $\lambda \|\theta\|_1$  is valid for both data compression and statistical risk requirements for  $\lambda \geq \lambda_n$  where

$$\lambda_n = \sqrt{\frac{2 \log(2M)}{n}}$$

- Especially useful for very large dictionaries
- Improvement for small dictionaries gets rid of log factor:  
 $\log(2M)$  may be replaced by  $\log(2e \max\{\frac{M}{\sqrt{n}}, 1\})$

## Comments on proof

- Existence of representor shown by random draw is a Shannon-like demonstration of the variable cover (code)
- Similar approximation in analysis of greedy computation of  $\ell_1$  penalized least squares



## Fixed $\sigma$ versus unknown $\sigma$

- MDL with  $\ell_1$  penalty for each possible  $\sigma$ . Recall

$$\min_{\theta} \left\{ \frac{1}{2\sigma^2} \frac{1}{n} \sum_{i=1}^n (Y_i - f_{\theta}(x_i))^2 + \frac{1}{2} \log 2\pi\sigma^2 + \frac{\lambda_n}{\sigma} \|\theta\|_1 \right\}$$

- Provides a family of fits indexed by  $\sigma$ .
- For unknown  $\sigma$  suggest optimization over  $\sigma$  as well as  $\theta$

$$\min_{\theta, \sigma} \left\{ \frac{1}{2\sigma^2} \frac{1}{n} \sum_{i=1}^n (Y_i - f_{\theta}(x_i))^2 + \frac{1}{2} \log 2\pi\sigma^2 + \frac{\lambda_n}{\sigma} \|\theta\|_1 \right\}$$

- Slight modification of this does indeed satisfy our condition for an information-theoretically valid penalty and risk bound (details in the WITMSE 2008 proceedings)

## Best $\sigma$

- Best  $\sigma$  for each  $\theta$  solves the quadratic equation

$$\sigma^2 = \sigma \lambda_n \|\theta\|_1 + \frac{1}{n} \sum_{i=1}^n (Y_i - f_\theta(x_i))^2$$

- By the quadratic formula the solution is

$$\sigma = \frac{1}{2} \lambda_n \|\theta\|_1 + \sqrt{\left[ \frac{1}{2} \lambda_n \|\theta\|_1 \right]^2 + \frac{1}{n} \sum_{i=1}^n (Y_i - f_\theta(x_i))^2}$$

# Outline

- 1 Some Principles of Information Theory and Statistics
  - Data Compression
  - Minimum Description Length Principle
  - Two-stage Codes, Mixture Codes, Normalized Max Like.
- 2 Penalized Likelihood
  - Redundancy, Resolvability, and Risk
  - Risk-Valid Penalties
  - Penalized Likelihood Analysis for Continuous Parameters
  - $\ell_0$  and  $\ell_1$  Penalties are Codelength-Valid and Risk-Valid
- 3 Summary

# Summary

- We equated forms of MDL including NML and Bayes
- We related MDL and penalized likelihood
- Allow penalized likelihoods with continuous domains for  $\theta$
- Information-theoretically valid penalties exceed complexity plus distortion of optimized representor of  $\theta$
- Yields MDL interpretation and stat. risk  $\leq$  resolvability
- $\ell_1$  penalty  $\lambda_n \|\theta\|_1$  valid in regression and related problems for  $\lambda_n \geq \sqrt{2(\log 2M)/n}$
- $\ell_0$  penalty  $\frac{d}{2} \log n + C_d(\theta)$  is a classical codelength penalty.
- The next talk by Sabyasachi Chatterjee shows how to avoid the growth of  $C_d(\theta)$  with the size of  $\|\theta\|$ .
- He shows that  $\frac{d}{2} \log n + \text{constant}_d$  is conditionally codelength valid and risk valid.

# Summary

- We equated forms of MDL including NML and Bayes
- We related MDL and penalized likelihood
- Allow penalized likelihoods with continuous domains for  $\theta$
- Information-theoretically valid penalties exceed complexity plus distortion of optimized representor of  $\theta$
- Yields MDL interpretation and stat. risk  $\leq$  resolvability
- $\ell_1$  penalty  $\lambda_n \|\theta\|_1$  valid in regression and related problems for  $\lambda_n \geq \sqrt{2(\log 2M)/n}$
- $\ell_0$  penalty  $\frac{d}{2} \log n + C_d(\theta)$  is a classical codelength penalty.
- The next talk by Sabyasachi Chatterjee shows how to avoid the growth of  $C_d(\theta)$  with the size of  $\|\theta\|$ .
- He shows that  $\frac{d}{2} \log n + \text{constant}_d$  is conditionally codelength valid and risk valid.

# Summary

- We equated forms of MDL including NML and Bayes
- We related MDL and penalized likelihood
- Allow penalized likelihoods with continuous domains for  $\theta$
- Information-theoretically valid penalties exceed complexity plus distortion of optimized representor of  $\theta$
- Yields MDL interpretation and stat. risk  $\leq$  resolvability
- $\ell_1$  penalty  $\lambda_n \|\theta\|_1$  valid in regression and related problems for  $\lambda_n \geq \sqrt{2(\log 2M)/n}$
- $\ell_0$  penalty  $\frac{d}{2} \log n + C_d(\theta)$  is a classical codelength penalty.
- The next talk by Sabyasachi Chatterjee shows how to avoid the growth of  $C_d(\theta)$  with the size of  $\|\theta\|$ .
- He shows that  $\frac{d}{2} \log n + \text{constant}_d$  is conditionally codelength valid and risk valid.

# Summary

- We equated forms of MDL including NML and Bayes
- We related MDL and penalized likelihood
- Allow penalized likelihoods with continuous domains for  $\theta$
- Information-theoretically valid penalties exceed complexity plus distortion of optimized representor of  $\theta$
- Yields MDL interpretation and stat. risk  $\leq$  resolvability
- $\ell_1$  penalty  $\lambda_n \|\theta\|_1$  valid in regression and related problems for  $\lambda_n \geq \sqrt{2(\log 2M)/n}$
- $\ell_0$  penalty  $\frac{d}{2} \log n + C_d(\theta)$  is a classical codelength penalty.
- The next talk by Sabyasachi Chatterjee shows how to avoid the growth of  $C_d(\theta)$  with the size of  $\|\theta\|$ .
- He shows that  $\frac{d}{2} \log n + \text{constant}_d$  is conditionally codelength valid and risk valid.

# Summary

- We equated forms of MDL including NML and Bayes
- We related forms of MDL and penalized likelihood
- Allow penalized likelihoods with continuous domains for  $\theta$
- Information-theoretically valid penalties exceed complexity plus distortion of optimized representor of  $\theta$
- Yields MDL interpretation and stat. risk  $\leq$  resolvability
- $\ell_1$  penalty  $\lambda_n \|\theta\|_1$  valid in regression and related problems for  $\lambda_n \geq \sqrt{2(\log 2M)/n}$
- $\ell_0$  penalty  $\frac{d}{2} \log n + C_d(\theta)$  is a classical codelength penalty.
- The next talk by Sabyasachi Chatterjee shows how to avoid the growth of  $C_d(\theta)$  with the size of  $\|\theta\|$ .
- He shows that  $\frac{d}{2} \log n + \text{constant}_d$  is conditionally codelength valid and risk valid.