Learning continuous distributions: Simulations with field theoretic priors

Ilya Nemenman^{1,2}, William Bialek¹

¹NEC Research Institute ²Princeton University

August 27, 2003



$$\log P(A|X) \rightarrow \sum_{i} \underbrace{\log Q_A(x_i | \alpha_{\mathsf{ML}})}_{\text{goodness of fit}} - \underbrace{\frac{K_A}{2} \log N}_{\text{generalization error}} + \dots$$

Bayesian inference penalizes for complexity (large K).

Fight between the goodness of fit and the complexity selects an optimal model family.

This is a Bayesian analogue of the MDL principle.

Does this generalize to infinite-dimensional models?

Bayesian learning for $K \to \infty$

(Bialek, Callan, Strong 1996)

Finite	Infinite			
α	$\phi(x) = -\log \ell_0 Q(x)$			
$\mathcal{P}(oldsymbollpha)$	$\mathcal{P}[Q] \propto \exp\left[-rac{\ell^{2\eta-1}}{2} \int dx (\partial_x^\eta \phi)^2 ight]$			
	smoothness penalty			
$\{A, K_A\}$	$\{\ell, \eta(?)\}$ – index continuum of families			
Pr(A)	$Pr(\ell, \eta(?))$			

Fix ℓ and η :

$$P[Q|X] = \frac{P(X|Q)\mathcal{P}[Q]}{P(X)}$$

$$\langle Q \rangle = \frac{\int [dQ] \mathcal{P}[Q] Q(x) \prod_{i=1}^{N} Q(x_i)}{\int [dQ] P[Q] \prod_{i=1}^{N} Q(x_i)}$$

$$= \frac{\langle Q(x)Q(x_1)\cdots Q(x_N) \rangle^0}{\langle Q(x_1)\cdots Q(x_N) \rangle^0}$$

Correlation function in a QFT defined by $\mathcal{P}[Q]$

Correlation functions:

C. F.
$$\equiv \int [dQ] \mathcal{P}[Q] \prod_{i=1}^{N} Q(x_i)$$
$$= \int [d\phi] \frac{1}{\ell_0^N} e^{-S[\phi]} \delta \left[\int dx \frac{1}{\ell_0} e^{-\phi} - 1 \right]$$
$$\underbrace{S[\phi]}_{\text{action}} = \underbrace{\frac{\ell}{2} \int dx (\partial_x^{\eta} \phi)^2}_{\text{kinetic term}} + \underbrace{\sum_{i=1}^{i} \phi(x_i)}_{\text{random potential}}$$

Large N approximation for $\eta = 1$

ML (classical, saddle point) solution dominates converges to $-\log \ell_0 P(x)$ $\ell \partial_x^2 \phi_{\text{Cl}}(x) + \frac{N}{\ell_0} e^{-\phi_{\text{Cl}}(x)} = \sum_j \delta(x - x_j)$ C. F. $\approx (1/\ell_0)^N e^{-S_{\text{eff}}[\phi_{\text{cl}}(x)]}$ $S_{\text{eff}}[\phi_{\text{cl}}] = \frac{\ell}{2} \int dx (\partial \phi_{\text{cl}})^2 + \sum_{\text{goodness of fit}} \phi_{\text{cl}}(x_i)$ $prior, \text{smoothness}}$

fluctuations, complexity, error

How do we measure performance? For $x \in [0, L)$ the *universal* learning curve is

$$\Lambda(N) \to \langle D_{\mathsf{KL}}(P||Q_{\mathsf{CI}}) \rangle_{\{x_i\}}^{\mathsf{0}} \sim \sqrt{\frac{L}{\ell N}}$$

For a different η :

$$\Lambda(N) \sim \left(\frac{L}{\ell}\right)^{1/2\eta} N^{1/2\eta-1}$$

5

Learning curves for fixed $\ell\text{, }\eta=1$



No overfits!

Smoothness scale selection

Allow a prior over $\ell,$ but keep $\eta=1$

C. F.
$$\rightarrow \langle \mathsf{C}. \mathsf{F}. \rangle_{\ell} = \int d\ell \ Pr(\ell) \ \mathrm{e}^{-S_{\mathsf{eff}}[\phi_{\mathsf{CI}}(\phi, \ell)]}$$

 $S_{\rm eff}[\phi_{\rm CI}] = \underbrace{{\rm smoothing} + {\rm data}_{\ell} + \underbrace{{\rm fluctuations}_{\ell}}_{{\rm grows \ with \ \ell}} + \underbrace{{\rm fluctuations}_{\ell}}_{{\rm grows \ with \ 1/\ell}}$

Some ℓ^* always dominates the C. F. and $\langle Q \rangle$!

What is ℓ^* for η_a and ℓ_a ?

If	$\eta =$	η_a ,	then	$\ell^* =$	= ℓ_a .	Otherwise:
----	----------	------------	------	------------	--------------	------------

$0.5 < \eta_a \le 1.5$	$1.5 < \eta_a$
data > smoothing	smoothing > data
$\ell^* \sim N^{(\eta_a - 1)/\eta_a}$	$\ell^* \sim N^{1/3}$
$\Lambda \sim N^{1/2\eta_a-1}$	$\Lambda \sim N^{-2/3}$
best possible	better, but not
performance	best performance

Averaging over ℓ and allowing $\ell^* = \ell^*(N)$ deals with *qualitatively* wrong smoothness $\eta_a \neq 1$!

MDL analogies.

Bayesian smoothness (model) selection works for nonparametric models!

Open questions

- constant factor or constant summand?
- what to do with $\eta_a > 1.5$?
- reparameterization invariance
- information theoretic meaningful priors
- higher dimensions

There is hope that all of this problems are resolvable in a single formulation.