

# Bayesian statistics, Occam razor, and model-independent learning of continuous probability densities

Ilya Nemenman  
ITP, UCSB

Joint work with:  
William Bialek, Princeton University

# Bayesian statistics ...

This is how it compares to other ultimate answers:

# Bayesian statistics . . .

. . . claims that the the answer to the Great Question of Life, The Universe and Everything is not 42, but

This is how it compares to other ultimate answers:

# Bayesian statistics . . .

. . . claims that the the answer to the Great Question of Life, The Universe and Everything is not 42, but

$$P(\text{Model}|\text{Data}) = \frac{P(\text{Data}|\text{Model}) \mathcal{P}(\text{Model})}{P(\text{Data})}$$

This is how it compares to other ultimate answers:

# Bayesian statistics ...

... claims that the the answer to the Great Question of Life, The Universe and Everything is not 42, but

$$P(\text{Model}|\text{Data}) = \frac{P(\text{Data}|\text{Model}) \mathcal{P}(\text{Model})}{P(\text{Data})}$$

This is how it compares to other ultimate answers:

Bayes

Best others

	Bayes	Best others
consistency	✓	✓
convergence rates	optimal	optimal
model selection	?	✓ (disagreement remains)
use of prior knowledge	✓	?

Bayes

Best others

model selection  
use of prior knowledge

✓ (disagreement remains)

?

later



Bayes

Best others

model selection


?

today

# Bayesian model selection for finitely parameterizable distributions

# Bayesian model selection for finitely parameterizable distributions


$P(x)$   
unknown




# Bayesian model selection for finitely parameterizable distributions

$$P(x) \xrightarrow{\text{i.i.d.}} X = \{x_1 \cdots x_N\}$$

unknown




# Bayesian model selection for finitely parameterizable distributions

  $P(x) \xrightarrow{\text{i.i.d.}} X = \{x_1 \cdots x_N\}$   
unknown

Model family  $A$   
 $Q_A(x|\alpha)$   
 $\dim \alpha = K_A$   
 $\mathcal{P}_A(\alpha), Pr(A)$

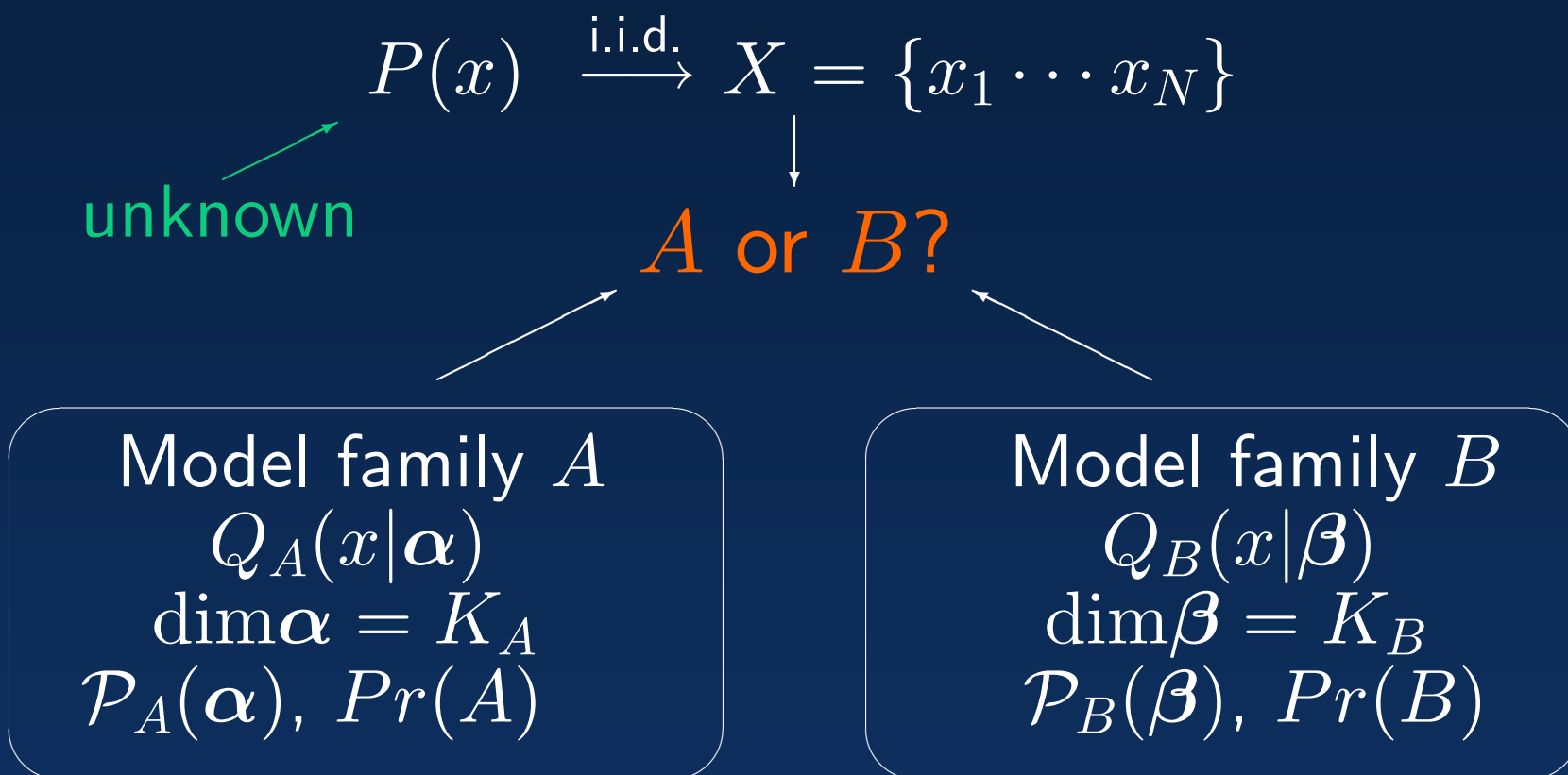
# Bayesian model selection for finitely parameterizable distributions

  $P(x) \xrightarrow{\text{i.i.d.}} X = \{x_1 \cdots x_N\}$   
unknown

Model family  $A$   
 $Q_A(x|\alpha)$   
 $\dim \alpha = K_A$   
 $\mathcal{P}_A(\alpha), Pr(A)$

Model family  $B$   
 $Q_B(x|\beta)$   
 $\dim \beta = K_B$   
 $\mathcal{P}_B(\beta), Pr(B)$

# Bayesian model selection for finitely parameterizable distributions



# Solution

Find the model with maximum posterior probability!



# Solution

Find the model with maximum posterior probability!

For example, for model  $A$ :

$$P(A|X) = \frac{P(X|A)Pr(A)}{P(X)} \quad \longleftarrow \quad P(X|A)Pr(A) + P(X|B)Pr(B) \equiv Z$$

$$P(X|A) = \int d\alpha \mathcal{P}_A(\alpha) P(X|\alpha) \sim P(X|\alpha_{\text{ML}}) \|\delta\alpha_{\text{ML}}\|$$

# Solution

Find the model with maximum posterior probability!

For example, for model  $A$ :

$$P(A|X) = \frac{P(X|A)Pr(A)}{P(X)} \quad \longleftarrow \quad P(X|A)Pr(A) + P(X|B)Pr(B) \equiv Z$$

$$P(X|A) = \int d\alpha \mathcal{P}_A(\alpha) P(X|\alpha) \sim P(X|\alpha_{\text{ML}}) \|\delta\alpha_{\text{ML}}\|$$

For large  $K_A$ ,  $\delta\alpha_{\text{ML}}$  (region of “good”  $\alpha$ ) decreases.

More complicated models are penalized!

# Solution

Find the model with maximum posterior probability!

For example, for model  $A$ :

$$P(A|X) = \frac{P(X|A)Pr(A)}{P(X)} \quad \leftarrow P(X|A)Pr(A) + P(X|B)Pr(B) \equiv Z$$

$$P(X|A) = \int d\alpha \mathcal{P}_A(\alpha) P(X|\alpha) \sim P(X|\alpha_{\text{ML}}) \|\delta\alpha_{\text{ML}}\|$$

For large  $K_A$ ,  $\delta\alpha_{\text{ML}}$  (region of “good”  $\alpha$ ) decreases.

More complicated models are penalized!

(See: Bayes factors, Occam factors; Jaynes 1968, 1979)

# Large $N$ expansion

Saddle point (large  $N$ ) expansion is almost always valid.

(See: Schwartz 1978, MacKay 1992, Balasubramanian 1996)

## Large $N$ expansion

Saddle point (large  $N$ ) expansion is almost always valid.

$$\log P(A|X) \rightarrow \sum_i \underbrace{\log Q_A(x_i|\boldsymbol{\alpha}_{\text{ML}})}$$

(See: Schwartz 1978, MacKay 1992, Balasubramanian 1996)

## Large $N$ expansion

Saddle point (large  $N$ ) expansion is almost always valid.

$$\log P(A|X) \rightarrow \sum_i \underbrace{\log Q_A(x_i|\alpha_{\text{ML}})}_{\substack{-\frac{K_A}{2} \log N - \log \det \partial^2_{\alpha_{\text{ML}}} \frac{\sum_i \log Q(x_i|\alpha_{\text{ML}})}{N}}}$$

(See: Schwartz 1978, MacKay 1992, Balasubramanian 1996)

## Large $N$ expansion

Saddle point (large  $N$ ) expansion is almost always valid.

$$\begin{aligned} \log P(A|X) \rightarrow & \sum_i \underbrace{\log Q_A(x_i|\alpha_{\text{ML}})} \\ & - \underbrace{\frac{K_A}{2} \log N - \log \det \partial^2_{\alpha_{\text{ML}}} \frac{\sum_i \log Q(x_i|\alpha_{\text{ML}})}{N}} \\ & + \log \mathcal{P}(\alpha_{\text{ML}}) + o(N^0) \end{aligned}$$

(See: Schwartz 1978, MacKay 1992, Balasubramanian 1996)



## Large $N$ expansion

Saddle point (large  $N$ ) expansion is almost always valid.

$$\begin{aligned} \log P(A|X) \rightarrow & \sum_i \underbrace{\log Q_A(x_i|\alpha_{\text{ML}})}_{\text{goodness of fit}} \\ & - \underbrace{\frac{K_A}{2} \log N - \log \det \partial^2_{\alpha_{\text{ML}}} \frac{\sum_i \log Q(x_i|\alpha_{\text{ML}})}{N}}_{\text{}} \\ & + \log \mathcal{P}(\alpha_{\text{ML}}) + o(N^0) \end{aligned}$$

(See: Schwartz 1978, MacKay 1992, Balasubramanian 1996)

# Large $N$ expansion

Saddle point (large  $N$ ) expansion is almost always valid.

$$\log P(A|X) \rightarrow \sum_i \underbrace{\log Q_A(x_i|\alpha_{\text{ML}})}_{\text{goodness of fit}}$$

$$- \underbrace{\frac{K_A}{2} \log N - \log \det \partial^2_{\alpha_{\text{ML}}} \frac{\sum_i \log Q(x_i|\alpha_{\text{ML}})}{N}}_{\text{generalization error, fluctuations, complexity; weak dependence on priors}}$$

$$+ \log \mathcal{P}(\alpha_{\text{ML}}) + o(N^0)$$

(See: Schwartz 1978, MacKay 1992, Balasubramanian 1996)

# Large $N$ expansion

Saddle point (large  $N$ ) expansion is almost always valid.

$$\log P(A|X) \rightarrow \sum_i \underbrace{\log Q_A(x_i|\alpha_{\text{ML}})}_{\text{goodness of fit}}$$

$$- \underbrace{\frac{K_A}{2} \log N - \log \det \partial^2_{\alpha_{\text{ML}}} \frac{\sum_i \log Q(x_i|\alpha_{\text{ML}})}{N}}_{\text{generalization error, fluctuations, complexity; weak dependence on priors}}$$

$$+ \log \mathcal{P}(\alpha_{\text{ML}}) + o(N^0)$$

(See: Schwartz 1978, MacKay 1992, Balasubramanian 1996)



# Conclusions

# Conclusions

- Bayesian inference penalizes for complexity (large  $K$ )

# Conclusions

- Bayesian inference penalizes for complexity (large  $K$ )
- Fight between the goodness of fit and the complexity selects an optimal model family.

# Conclusions

- 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.



# Conclusions

- 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?

# Estimating density

# Estimating density

Standard setting  
(solving IE)

Fisher–Wald setting  
(minimizing risk)

# Estimating density

Standard setting  
(solving IE)

$$F(t) = \int_{-\infty}^t Q(x) dx$$

Fisher–Wald setting  
(minimizing risk)

$$R[Q] = - \int_{-\infty}^{+\infty} \log Q(x) dF(x)$$

# Estimating density

Standard setting  
(solving IE)

$$F(t) = \int_{-\infty}^t Q(x) dx$$

$$\frac{1}{N} \sum_{x_i} \Theta(x_i - t) = \int_{-\infty}^t Q(x) dx$$

Fisher–Wald setting  
(minimizing risk)

$$R[Q] = - \int_{-\infty}^{+\infty} \log Q(x) dF(x)$$

$$R_{\text{emp}}[Q] = - \sum_{x_i} \log Q(x_i)$$

# Estimating density

Standard setting  
(solving IE)

$$F(t) = \int_{-\infty}^t Q(x) dx$$

$$\frac{1}{N} \sum_{x_i} \Theta(x_i - t) = \int_{-\infty}^t Q(x) dx$$

Fisher–Wald setting  
(minimizing risk)

$$R[Q] = - \int_{-\infty}^{+\infty} \log Q(x) dF(x)$$

$$R_{\text{emp}}[Q] = - \sum_{x_i} \log Q(x_i)$$

Both settings hypersensitive to fluctuations in  $F(t)$ .

*Smoothing is required.*

# Bayesian learning for $K \rightarrow \infty$

Finite	Infinite
--------	----------

# Bayesian learning for $K \rightarrow \infty$

Finite	Infinite
$\alpha$	$\phi(x) = -\log \ell_0 Q(x)$



# Bayesian learning for $K \rightarrow \infty$

Finite	Infinite
$\alpha$	$\phi(x) = -\log \ell_0 Q(x)$
$\mathcal{P}(\alpha)$	$\mathcal{P}[Q] \propto \exp \left[ -\frac{\ell^{2\eta-1}}{2} \underbrace{\int dx (\partial_x^\eta \phi)^2}_{\text{}} \right]$

# Bayesian learning for $K \rightarrow \infty$

Finite	Infinite
$\alpha$	$\phi(x) = -\log \ell_0 Q(x)$
$\mathcal{P}(\alpha)$	$\mathcal{P}[Q] \propto \exp \left[ -\frac{\ell^{2\eta-1}}{2} \underbrace{\int dx (\partial_x^\eta \phi)^2}_{\text{smoothness penalty}} \right]$

# Bayesian learning for $K \rightarrow \infty$

Finite	Infinite
$\alpha$	$\phi(x) = -\log \ell_0 Q(x)$
$\mathcal{P}(\alpha)$	$\mathcal{P}[Q] \propto \exp \left[ -\frac{\ell^{2\eta-1}}{2} \underbrace{\int dx (\partial_x^\eta \phi)^2}_{\text{smoothness penalty}} \right]$ <p>spline prior of order <math>2\eta - 1</math></p>

# Bayesian learning for $K \rightarrow \infty$

Finite	Infinite
$\alpha$	$\phi(x) = -\log \ell_0 Q(x)$
$\mathcal{P}(\alpha)$	$\mathcal{P}[Q] \propto \exp \left[ -\frac{\ell^{2\eta-1}}{2} \underbrace{\int dx (\partial_x^\eta \phi)^2}_{\text{smoothness penalty}} \right]$ <p>spline prior of order <math>2\eta - 1</math></p>
$\{A, K_A\}$	$\{\ell, \eta(?)\}$ — index continuum of families

# Bayesian learning for $K \rightarrow \infty$

Finite	Infinite
$\alpha$	$\phi(x) = -\log \ell_0 Q(x)$
$\mathcal{P}(\alpha)$	$\mathcal{P}[Q] \propto \exp \left[ -\frac{\ell^{2\eta-1}}{2} \underbrace{\int dx (\partial_x^\eta \phi)^2}_{\text{smoothness penalty}} \right]$ <p>spline prior of order <math>2\eta - 1</math></p>
$\{A, K_A\}$	$\{\ell, \eta(?)\}$ – index continuum of families
$Pr(A)$	$Pr(\ell, \eta(?))$

# Bayesian learning for $K \rightarrow \infty$

Finite	Infinite
$\alpha$	$\phi(x) = -\log \ell_0 Q(x)$
$\mathcal{P}(\alpha)$	$\mathcal{P}[Q] \propto \exp \left[ -\frac{\ell^{2\eta-1}}{2} \underbrace{\int dx (\partial_x^\eta \phi)^2}_{\text{smoothness penalty}} \right]$ <p style="text-align: center;">spline prior of order <math>2\eta - 1</math></p>
$\{A, K_A\}$	$\{\ell, \eta(?)\}$ – index continuum of families
$Pr(A)$	$Pr(\ell, \eta(?))$

(See: Bialek, Callan, Strong, 1996)

# Quantum Field Theory analogy

Fix  $\ell$  and  $\eta$ :

$$= \frac{\langle Q(x)Q(x_1) \cdots Q(x_N) \rangle^0}{\underbrace{\langle Q(x_1) \cdots Q(x_N) \rangle^0}_{\text{Correlation function in a QFT defined by } \mathcal{P}[Q]}}$$

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

# Quantum Field Theory analogy

Fix  $\ell$  and  $\eta$ :

$$P[Q|X] = \frac{P(X|Q)\mathcal{P}[Q]}{P(X)}$$
$$= \frac{\langle Q(x)Q(x_1) \cdots Q(x_N) \rangle^0}{\underbrace{\langle Q(x_1) \cdots Q(x_N) \rangle^0}_{\text{Correlation function in a QFT defined by } \mathcal{P}[Q]}}$$

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



# Quantum Field Theory analogy

Fix  $\ell$  and  $\eta$ :

$$\begin{aligned}
 P[Q|X] &= \frac{P(X|Q)\mathcal{P}[Q]}{P(X)} \\
 \langle Q \rangle &= \frac{\int [dQ] \mathcal{P}[Q] Q(x) \prod_{i=1}^N Q(x_i)}{\int [dQ] P[Q] \prod_{i=1}^N Q(x_i)} \\
 &= \frac{\langle Q(x) Q(x_1) \cdots Q(x_N) \rangle^0}{\underbrace{\langle Q(x_1) \cdots Q(x_N) \rangle^0}_{\text{Correlation function in a QFT defined by } \mathcal{P}[Q]}}
 \end{aligned}$$

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



# Explicit form of correlation functions

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



# Large $N$ approximation for $\eta = 1$

ML (classical, saddle point) solution dominates

## Large $N$ approximation for $\eta = 1$

ML (classical, saddle point) solution dominates

$$\ell \partial_x^2 \phi_{\text{cl}}(x) + \frac{N}{\ell_0} e^{-\phi_{\text{cl}}(x)} = \sum_j \delta(x - x_j)$$

# Large $N$ approximation for $\eta = 1$

ML (classical, saddle point) solution dominates

converges to  $-\log \ell_0 P(x)$       changes on scale  $\delta x \sim \sqrt{\ell/NP(x)}$

$$\ell \partial_x^2 \phi_{\text{cl}}(x) + \frac{N}{\ell_0} e^{-\phi_{\text{cl}}(x)} = \sum_j \delta(x - x_j)$$

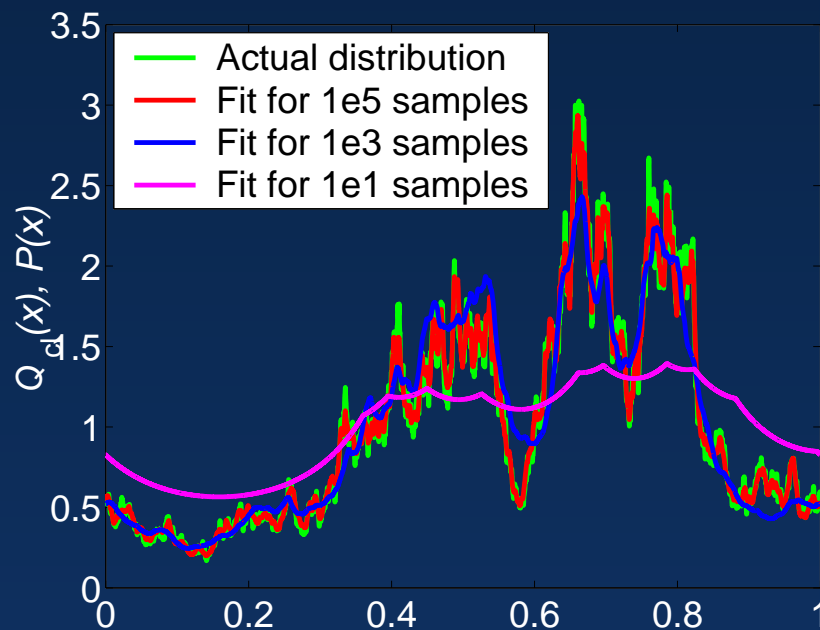
# Large $N$ approximation for $\eta = 1$

ML (classical, saddle point) solution dominates

converges to  
 $-\log \ell_0 P(x)$

changes on scale  
 $\delta x \sim \sqrt{\ell/NP(x)}$

$$\ell \partial_x^2 \phi_{cl}(x) + \frac{N}{\ell_0} e^{-\phi_{cl}(x)} = \sum_j \delta(x - x_j)$$





# Large $N$ approximation for $\eta = 1$ , continued

Van Vleck calculation of functional determinant:

# Large $N$ approximation for $\eta = 1$ , continued

Van Vleck calculation of functional determinant:

$$\text{C. F.} \approx (1/\ell_0)^N e^{-S_{\text{eff}}[\phi_{\text{cl}}(x)]}$$

# Large $N$ approximation for $\eta = 1$ , continued

Van Vleck calculation of functional determinant:

$$\begin{aligned}
 \text{C. F.} &\approx (1/\ell_0)^N e^{-S_{\text{eff}}[\phi_{\text{cl}}(x)]} \\
 S_{\text{eff}}[\phi_{\text{cl}}] &= \underbrace{\frac{\ell}{2} \int dx (\partial \phi_{\text{cl}})^2}_{\text{kinetic}} + \underbrace{\sum \phi_{\text{cl}}(x_i)}_{\text{potential}} \\
 &\quad + \underbrace{\frac{1}{2} \sqrt{\frac{N}{\ell \ell_0}} \int dx e^{-\phi_{\text{cl}}(x)/2}}_{\text{quantum correction}}
 \end{aligned}$$

# Large $N$ approximation for $\eta = 1$ , continued

Van Vleck calculation of functional determinant:

$$\begin{aligned}
 \text{C. F.} &\approx (1/\ell_0)^N e^{-S_{\text{eff}}[\phi_{\text{cl}}(x)]} \\
 S_{\text{eff}}[\phi_{\text{cl}}] &= \underbrace{\frac{\ell}{2} \int dx (\partial \phi_{\text{cl}})^2}_{\text{prior, smoothness}} + \underbrace{\sum \phi_{\text{cl}}(x_i)}_{\text{goodness of fit}} \\
 &+ \underbrace{\frac{1}{2} \sqrt{\frac{N}{\ell \ell_0}} \int dx e^{-\phi_{\text{cl}}(x)/2}}_{\text{fluctuations, complexity, error}}
 \end{aligned}$$

# How do we measure performance?

# How do we measure performance?

For  $x \in [0, L)$  the *universal* learning curve is

$$\Lambda(N) \rightarrow \langle D_{\text{KL}}(P || Q_{\text{cl}}) \rangle_{\{x_i\}}^0 \sim \sqrt{\frac{L}{\ell N}}$$

# How do we measure performance?

For  $x \in [0, L)$  the *universal* learning curve is

$$\Lambda(N) \rightarrow \langle D_{\text{KL}}(P || Q_{\text{cl}}) \rangle_{\{x_i\}}^0 \sim \sqrt{\frac{L}{\ell N}}$$

For a different  $\eta$ :

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

# Learning curves for fixed $\ell$ , $\eta = 1$



# Learning curves for fixed $\ell$ , $\eta = 1$

Learner's assumptions  $\mathcal{P}_{\ell, \eta=1}[Q]$

# Learning curves for fixed $\ell$ , $\eta = 1$

Learner's assumptions	$\mathcal{P}_{\ell, \eta=1}[Q]$
Actual target distribution	$\mathcal{P}'_{\ell_a, \eta_a}[Q]$

## Learning curves for fixed $\ell$ , $\eta = 1$

Learner's assumptions  $\mathcal{P}_{\ell, \eta=1}[Q]$

Actual target distribution  $\mathcal{P}'_{\ell_a, \eta_a}[Q]$

$\eta = \eta_a$ ,  $\ell = \ell_a$  learning typical cases,  $\mathcal{P} = \mathcal{P}'$

# Learning curves for fixed $\ell$ , $\eta = 1$

Learner's assumptions  $\mathcal{P}_{\ell, \eta=1}[Q]$

Actual target distribution  $\mathcal{P}'_{\ell_a, \eta_a}[Q]$

$\eta = \eta_a$ ,  $\ell = \ell_a$  learning typical cases,  $\mathcal{P} = \mathcal{P}'$

$\eta = \eta_a$ ,  $\ell \neq \ell_a$  marginal outliers of  $\mathcal{P}$

# Learning curves for fixed $\ell$ , $\eta = 1$

Learner's assumptions  $\mathcal{P}_{\ell, \eta=1}[Q]$

Actual target distribution  $\mathcal{P}'_{\ell_a, \eta_a}[Q]$

$\eta = \eta_a$ ,  $\ell = \ell_a$  learning typical cases,  $\mathcal{P} = \mathcal{P}'$

$\eta = \eta_a$ ,  $\ell \neq \ell_a$  marginal outliers of  $\mathcal{P}$

$\eta > \eta_a$  extremely rough outliers

# Learning curves for fixed $\ell$ , $\eta = 1$

Learner's assumptions  $\mathcal{P}_{\ell, \eta=1}[Q]$

Actual target distribution  $\mathcal{P}'_{\ell_a, \eta_a}[Q]$

$\eta = \eta_a$ ,  $\ell = \ell_a$  learning typical cases,  $\mathcal{P} = \mathcal{P}'$

$\eta = \eta_a$ ,  $\ell \neq \ell_a$  marginal outliers of  $\mathcal{P}$

$\eta > \eta_a$  extremely rough outliers

$\eta < \eta_a$  extremely smooth outliers

# Learning curves for fixed $\ell$ , $\eta = 1$

Learner's assumptions  $\mathcal{P}_{\ell, \eta=1}[Q]$

Actual target distribution  $\mathcal{P}'_{\ell_a, \eta_a}[Q]$

$\eta = \eta_a$ ,  $\ell = \ell_a$  learning typical cases,  $\mathcal{P} = \mathcal{P}'$

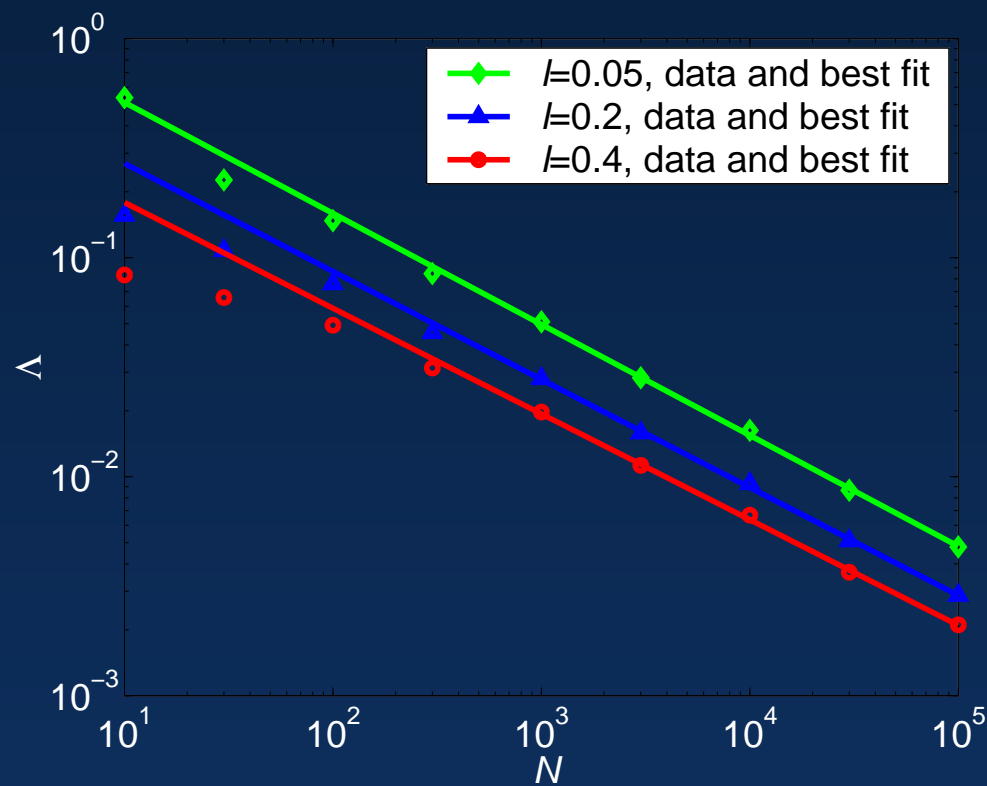
$\eta = \eta_a$ ,  $\ell \neq \ell_a$  marginal outliers of  $\mathcal{P}$

$\eta > \eta_a$  extremely rough outliers

$\eta < \eta_a$  extremely smooth outliers

Note: we must have  $\eta > 1/2$  for convergence of the integrals.

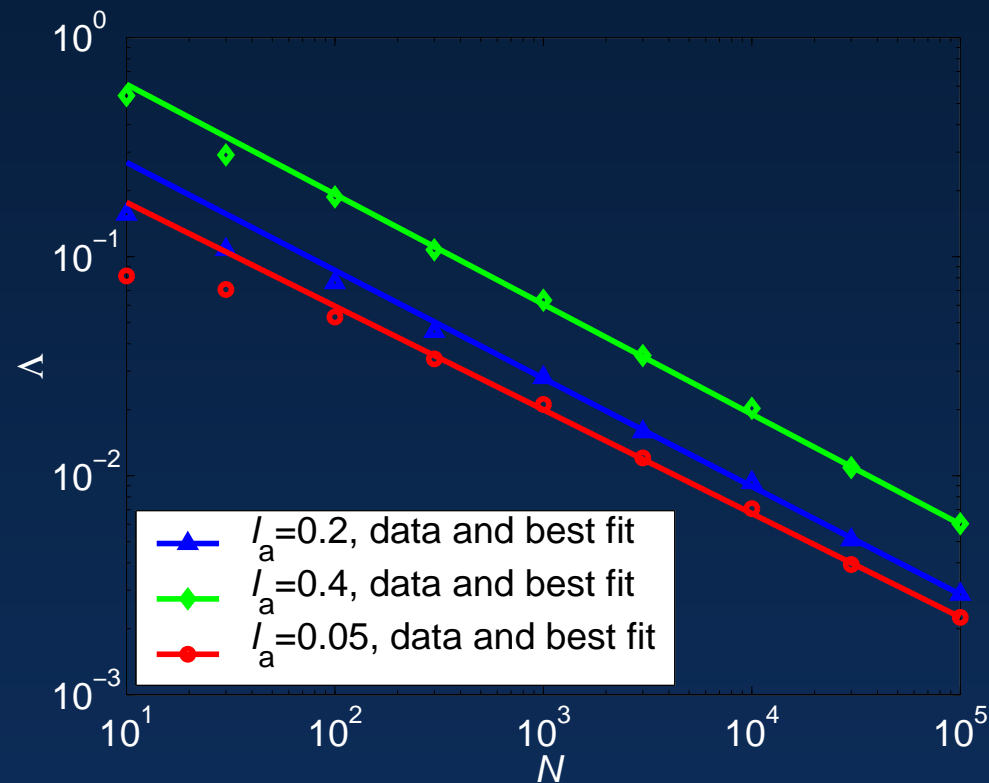
# Learning typical cases





$$\begin{aligned}\ell = 0.4, \quad \Lambda &= (0.54 \pm 0.07) N^{-0.483 \pm 0.014} \\ \ell = 0.2, \quad \Lambda &= (0.83 \pm 0.08) N^{-0.493 \pm 0.09} \\ \ell = 0.05, \quad \Lambda &= (1.64 \pm 0.16) N^{-0.507 \pm 0.09}\end{aligned}$$

# Learning marginal outliers

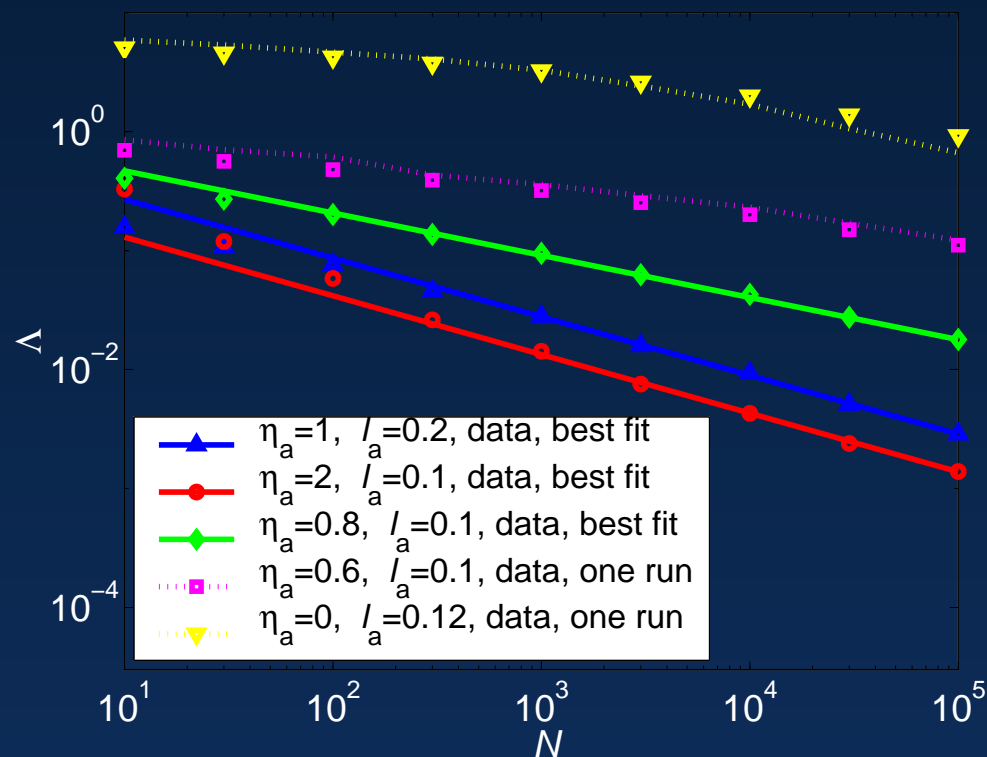


$$l_a = 0.4, \quad \Lambda = (0.56 \pm 0.08) N^{-0.477 \pm 0.015}$$

$$l_a = 0.05, \quad \Lambda = (1.90 \pm 0.16) N^{-0.502 \pm 0.008}$$

Learning at  $\ell = 0.2$ .

# Learning strong outliers

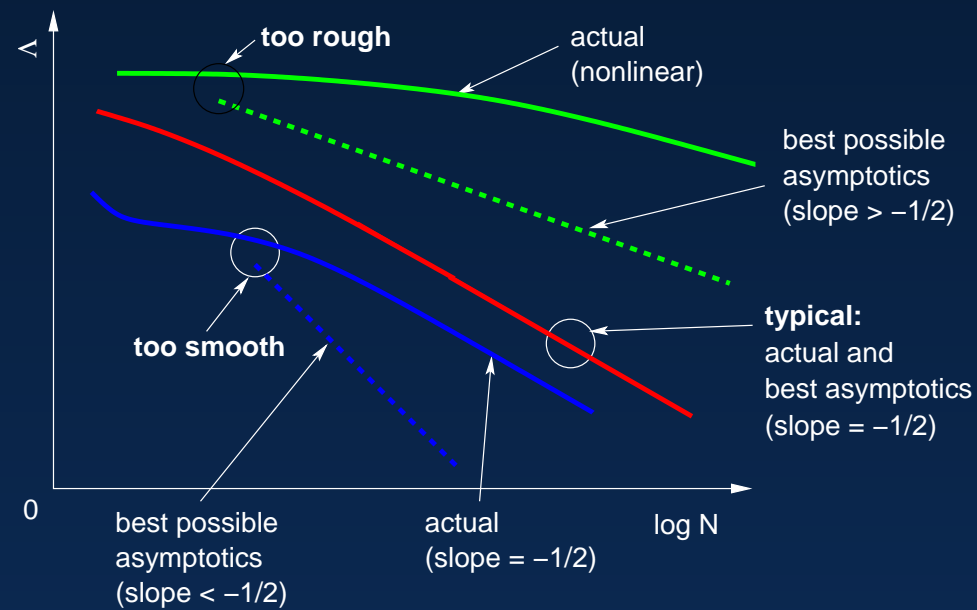


$$\eta_a = 2, \ell_a = 0.1, \quad \Lambda = (0.40 \pm 0.05) N^{-0.493 \pm 0.013}$$

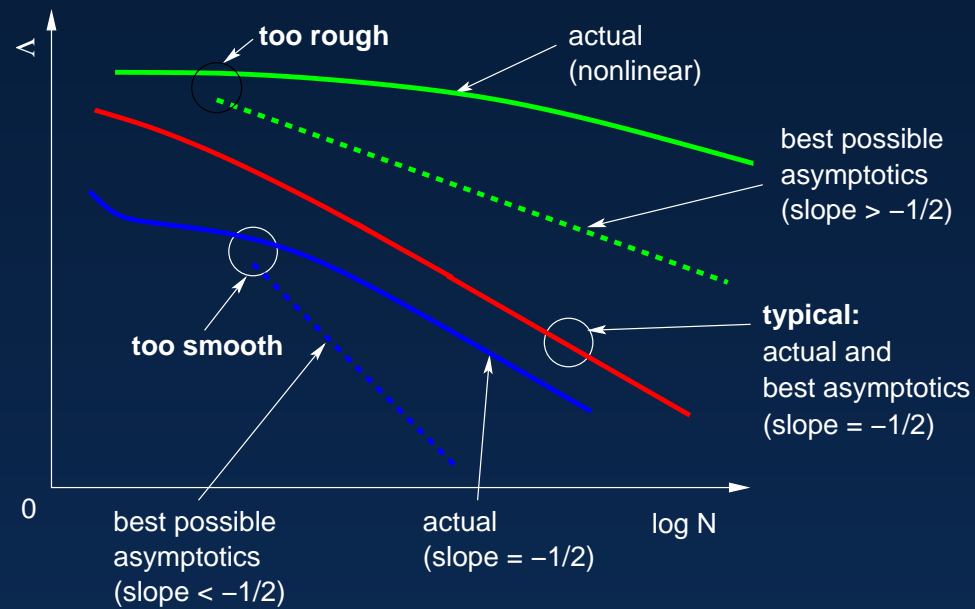
$$\eta_a = 0.8, \ell_a = 0.1, \quad \Lambda = (1.06 \pm 0.08) N^{-0.355 \pm 0.008}$$

$\ell = 0.1$  for  $\eta_a = 0$  and  $\ell = 0.2$  otherwise

# Conclusions for fixed $\eta$ and $\ell$

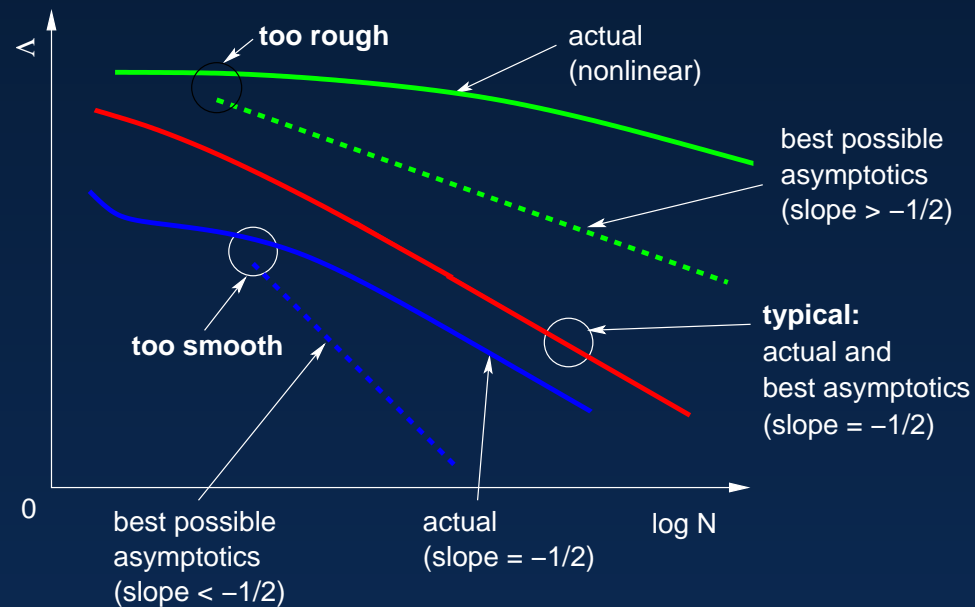


# Conclusions for fixed $\eta$ and $\ell$



- No overfits!

# Conclusions for fixed $\eta$ and $\ell$



- No overfits!
- but suboptimal performance for learning outliers



# Smoothness scale selection

# Smoothness scale selection

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

$$\text{C. F.} \rightarrow \langle \text{C. F.} \rangle_{\ell}$$

# Smoothness scale selection

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

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

## Smoothness scale selection

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

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

$$S_{\text{eff}}[\phi_{\text{cl}}] = \underbrace{\text{smoothing} + \text{data}} + \underbrace{\text{fluctuations}}$$

# Smoothness scale selection

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

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

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

## Smoothness scale selection

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

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

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

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

**Calculations: What is  $\ell^*$  for  $\eta_a$  and  $\ell_a$ ?**

Averaging over  $\ell$  and allowing  $\ell^* = \ell^*(N)$  deals with

## Calculations: What is $\ell^*$ for $\eta_a$ and $\ell_a$ ?

If  $\eta = \eta_a$ , then  $\ell^* = \ell_a$ .

Averaging over  $\ell$  and allowing  $\ell^* = \ell^*(N)$  deals with



# Calculations: What is $\ell^*$ for $\eta_a$ and $\ell_a$ ?

If  $\eta = \eta_a$ , then  $\ell^* = \ell_a$ . Otherwise:

$0.5 < \eta_a \leq 1.5$	$1.5 < \eta_a$
-------------------------	----------------

Averaging over  $\ell$  and allowing  $\ell^* = \ell^*(N)$  deals with

# Calculations: What is $\ell^*$ for $\eta_a$ and $\ell_a$ ?

If  $\eta = \eta_a$ , then  $\ell^* = \ell_a$ . Otherwise:

$0.5 < \eta_a \leq 1.5$	$1.5 < \eta_a$
data > smoothing	smoothing > data

Averaging over  $\ell$  and allowing  $\ell^* = \ell^*(N)$  deals with

# Calculations: What is $\ell^*$ for $\eta_a$ and $\ell_a$ ?

If  $\eta = \eta_a$ , then  $\ell^* = \ell_a$ . Otherwise:

$0.5 < \eta_a \leq 1.5$	$1.5 < \eta_a$
data > smoothing	smoothing > data
$\ell^* \sim N^{(\eta_a-1)/\eta_a}$	$\ell^* \sim N^{1/3}$

Averaging over  $\ell$  and allowing  $\ell^* = \ell^*(N)$  deals with

# Calculations: What is $\ell^*$ for $\eta_a$ and $\ell_a$ ?

If  $\eta = \eta_a$ , then  $\ell^* = \ell_a$ . Otherwise:

$0.5 < \eta_a \leq 1.5$	$1.5 < \eta_a$
data > smoothing	smoothing > data
$\ell^* \sim N^{(\eta_a-1)/\eta_a}$	$\ell^* \sim N^{1/3}$
$\Lambda \sim N^{1/2\eta_a-1}$	$\Lambda \sim N^{-2/3}$

Averaging over  $\ell$  and allowing  $\ell^* = \ell^*(N)$  deals with

# Calculations: What is $\ell^*$ for $\eta_a$ and $\ell_a$ ?

If  $\eta = \eta_a$ , then  $\ell^* = \ell_a$ . Otherwise:

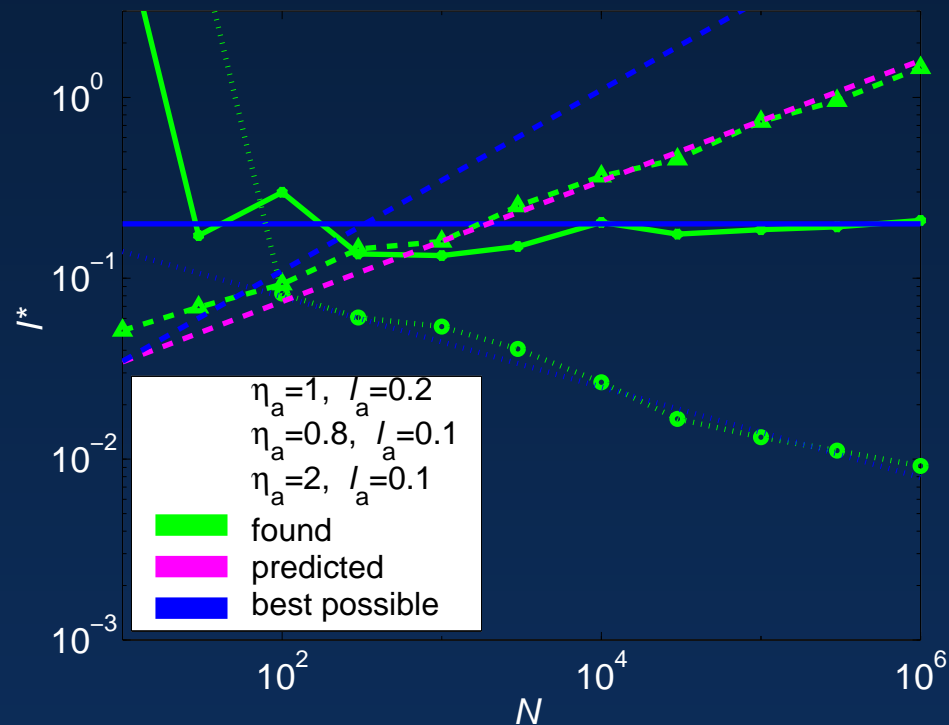
$0.5 < \eta_a \leq 1.5$	$1.5 < \eta_a$
data > smoothing	smoothing > data
$\ell^* \sim N^{(\eta_a-1)/\eta_a}$	$\ell^* \sim N^{1/3}$
$\Lambda \sim N^{1/2\eta_a-1}$	$\Lambda \sim N^{-2/3}$
best possible performance	better, but not best performance

Averaging over  $\ell$  and allowing  $\ell^* = \ell^*(N)$  deals with

*qualitatively* wrong smoothness  $\eta_a \neq 1!$

**Numerics: What is  $\ell^*$  for  $\eta_a$  and  $\ell_a$ ?**

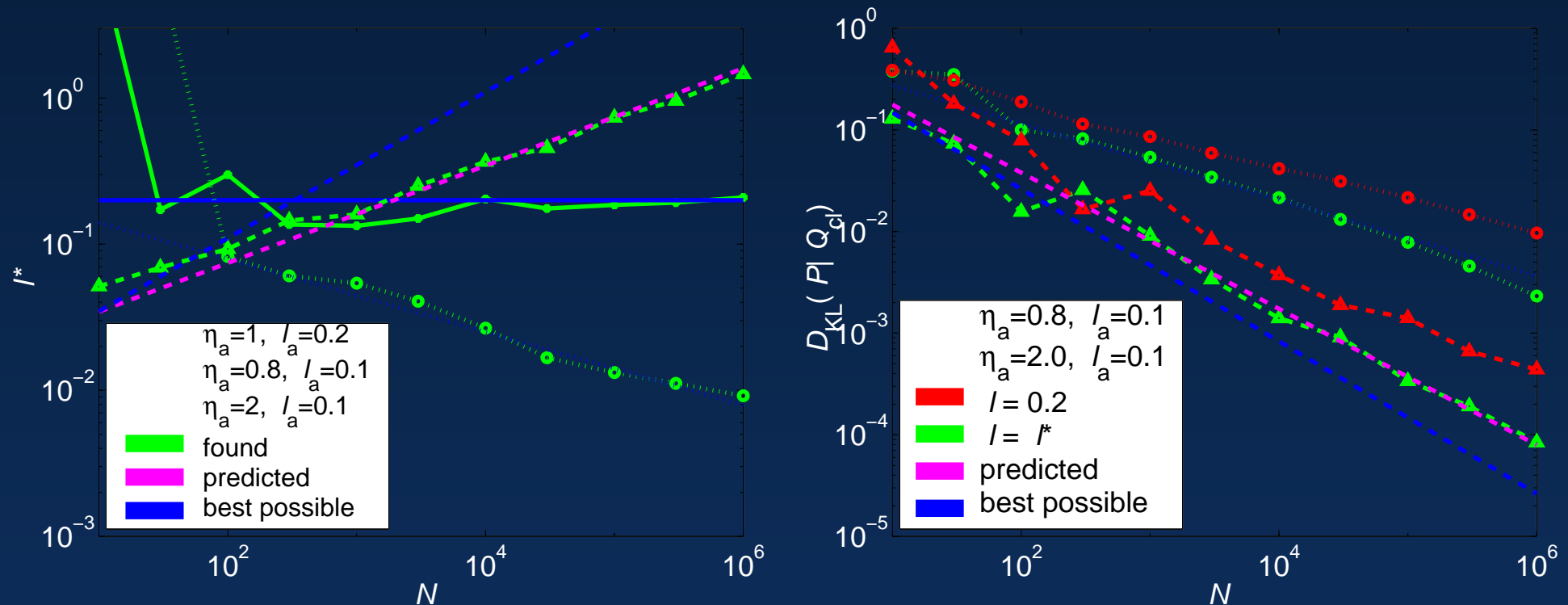
# Numerics: What is $\ell^*$ for $\eta_a$ and $\ell_a$ ?



Note: just single runs shown.

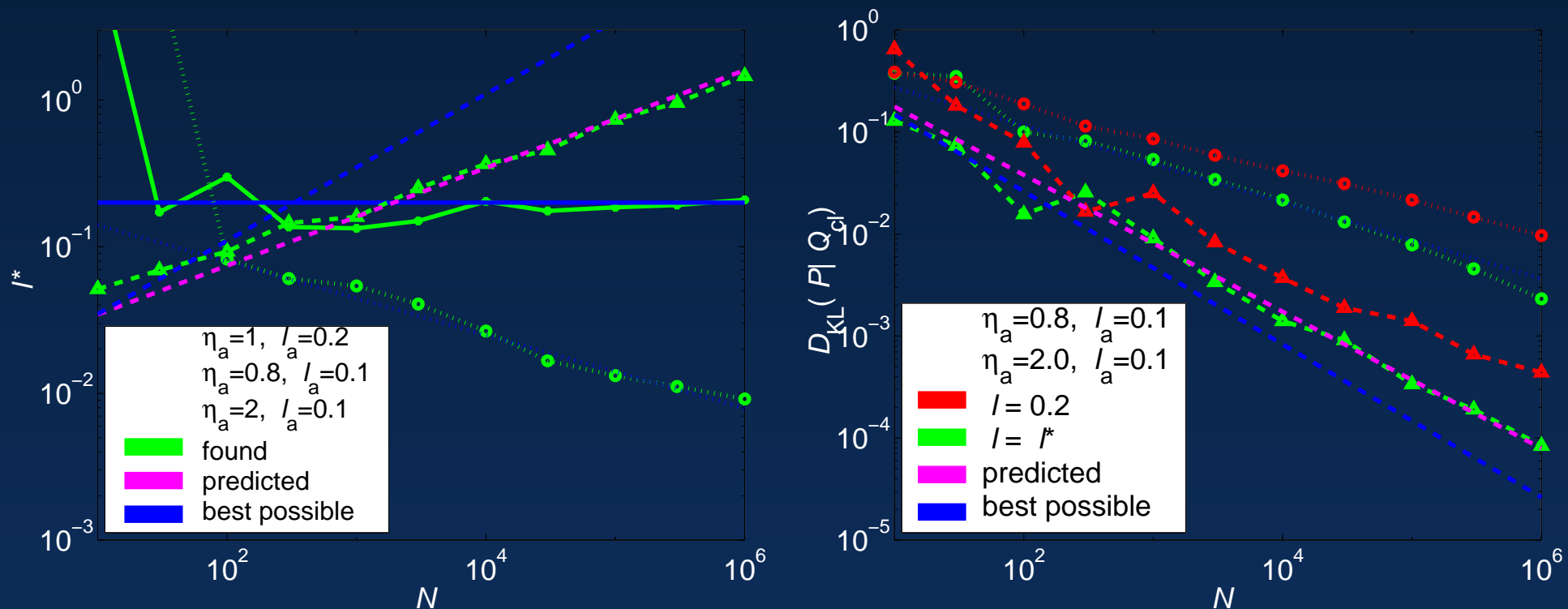


# Numerics: What is $\ell^*$ for $\eta_a$ and $\ell_a$ ?



Note: just single runs shown.

# Numerics: What is $\ell^*$ for $\eta_a$ and $\ell_a$ ?



Note: just single runs shown.

# Approaching model-independent optimal inference!

# Analogies

# Analogies

- choosing  $\ell^*$  corresponds to selection of a structure element with  $d_{VC} = \sqrt{NL/\ell^*}$  in Vapnik's SRM theory

# Analogies

- choosing  $\ell^*$  corresponds to selection of a structure element with  $d_{VC} = \sqrt{NL/\ell^*}$  in Vapnik's SRM theory
- maximizing  $P$  over model families ( $\ell$ 's) asymptotically corresponds to searching for MDL

# Analogies

- choosing  $\ell^*$  corresponds to selection of a structure element with  $d_{VC} = \sqrt{NL/\ell^*}$  in Vapnik's SRM theory
- maximizing  $P$  over model families ( $\ell$ 's) asymptotically corresponds to searching for MDL
- a lot in common with the Gaussian Processes theory; however normalization constraint is important

# Summary

**Bayesian smoothness (model) selection  
works for nonparametric spline priors!**



# Open questions

# Open questions

- constant factor or constant summand?

# Open questions

- constant factor or constant summand?
- what to do with  $\eta_a > 1.5$ ?

# Open questions

- constant factor or constant summand?
- what to do with  $\eta_a > 1.5$ ?
- reparameterization invariance

# Open questions

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

# Open questions

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

# Open questions

- constant factor or constant summand?
- what to do with  $\eta_a > 1.5$ ?
- reparameterization invariance
- information theoretic meaningful priors
- higher dimensions
- smooth transition from  $K = \text{const}$  to  $K \rightarrow \infty$

# Open questions

- constant factor or constant summand?
- what to do with  $\eta_a > 1.5$ ?
- reparameterization invariance
- information theoretic meaningful priors
- higher dimensions
- smooth transition from  $K = \text{const}$  to  $K \rightarrow \infty$
- which classes of priors are allowed?



# Open questions

- constant factor or constant summand?
- what to do with  $\eta_a > 1.5$ ?
- reparameterization invariance
- information theoretic meaningful priors
- higher dimensions
- smooth transition from  $K = \text{const}$  to  $K \rightarrow \infty$
- which classes of priors are allowed?

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