

Gaussian processes for PDF determination

Tommaso Giani

Based on Eur.Phys.J.C 84 (2024)

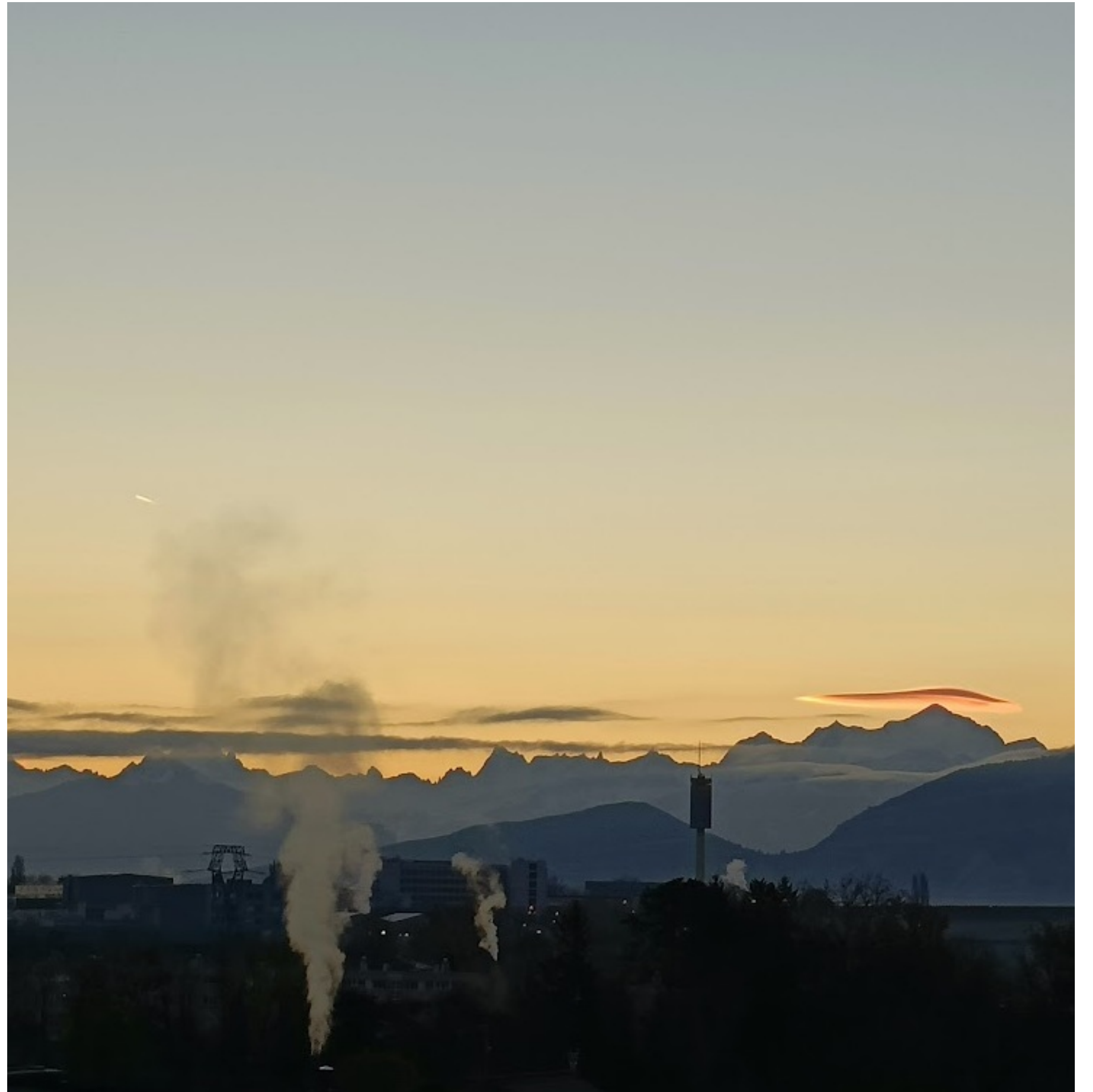
+

WIP with Maria Ubiali, Mark Constantini, Luigi Del
debbio, Luca Mantani

PDF4LHC24, Cern,
03/12/2024



MAX-PLANCK-INSTITUT
FÜR PHYSIK

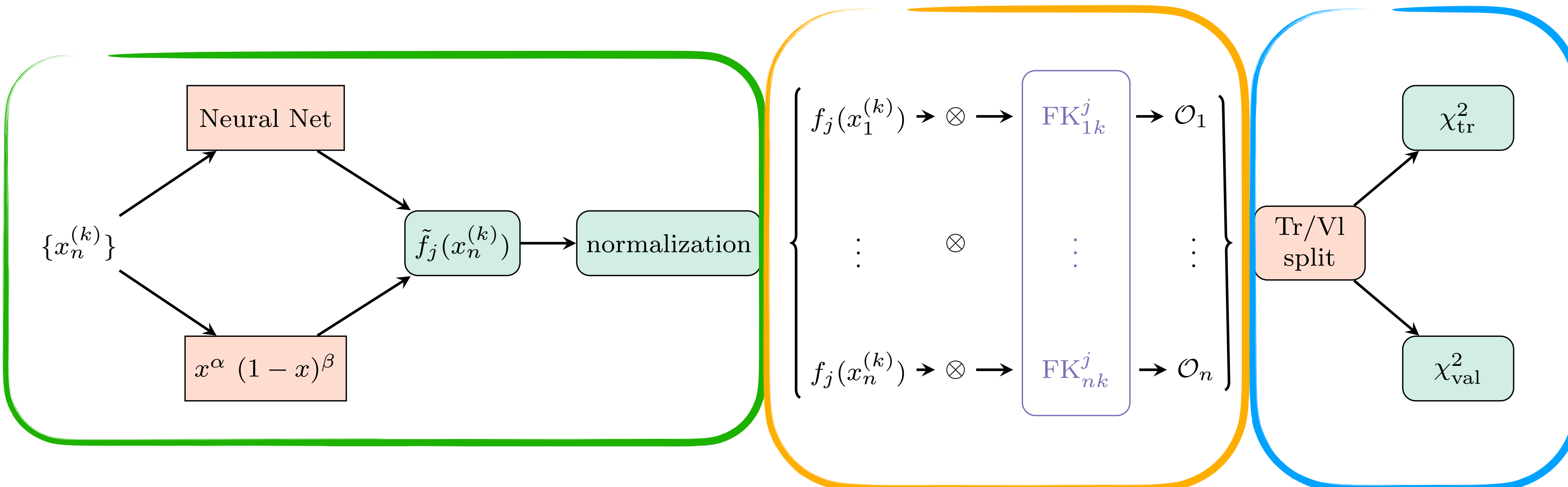


Parametric regression

- PDFs are parametrised at some initial scale Q_0
- Theory prediction are computed as function of the free parameters
- Use data to build χ^2 and determine best fit

Error propagation:

- Montecarlo replicas
- Hessian approach



Non-parametric regression and Bayesian approach

- Start from a prior on the model $p(f)$
- Look at the data
- Get the posterior $p(f|D)$

$$p(f|D) = \frac{p(D|f)p(f)}{p(D)}$$

Prior on the model

Posterior of model given the data

- Introduce probability distribution on a space of functions
- Build a suitable prior
- Use Bayes' theorem

Gaussian Processes

$$\mathbf{f} = \begin{pmatrix} f(x_1) \\ \vdots \\ f(x_N) \end{pmatrix} \in \mathbb{R}^N$$

Parameters \mathbf{f} : stochastic variables representing values of the PDF on a grid of points

Kernel \mathbf{K} and mean function \mathbf{m} : functions modelling the correlation between parameters

$$m(x_i; \theta) = \mathbb{E} \left(f(x_i) \right)$$

$$k(x_i, x_j; \theta) = \text{COV} \left(f(x_i), f(x_j) \right)$$

Hyperparameters θ : set of parameters entering the definition of the kernel (they control some specific feature of the prior)

Joint probability distribution of \mathbf{f} and θ : target of the analysis

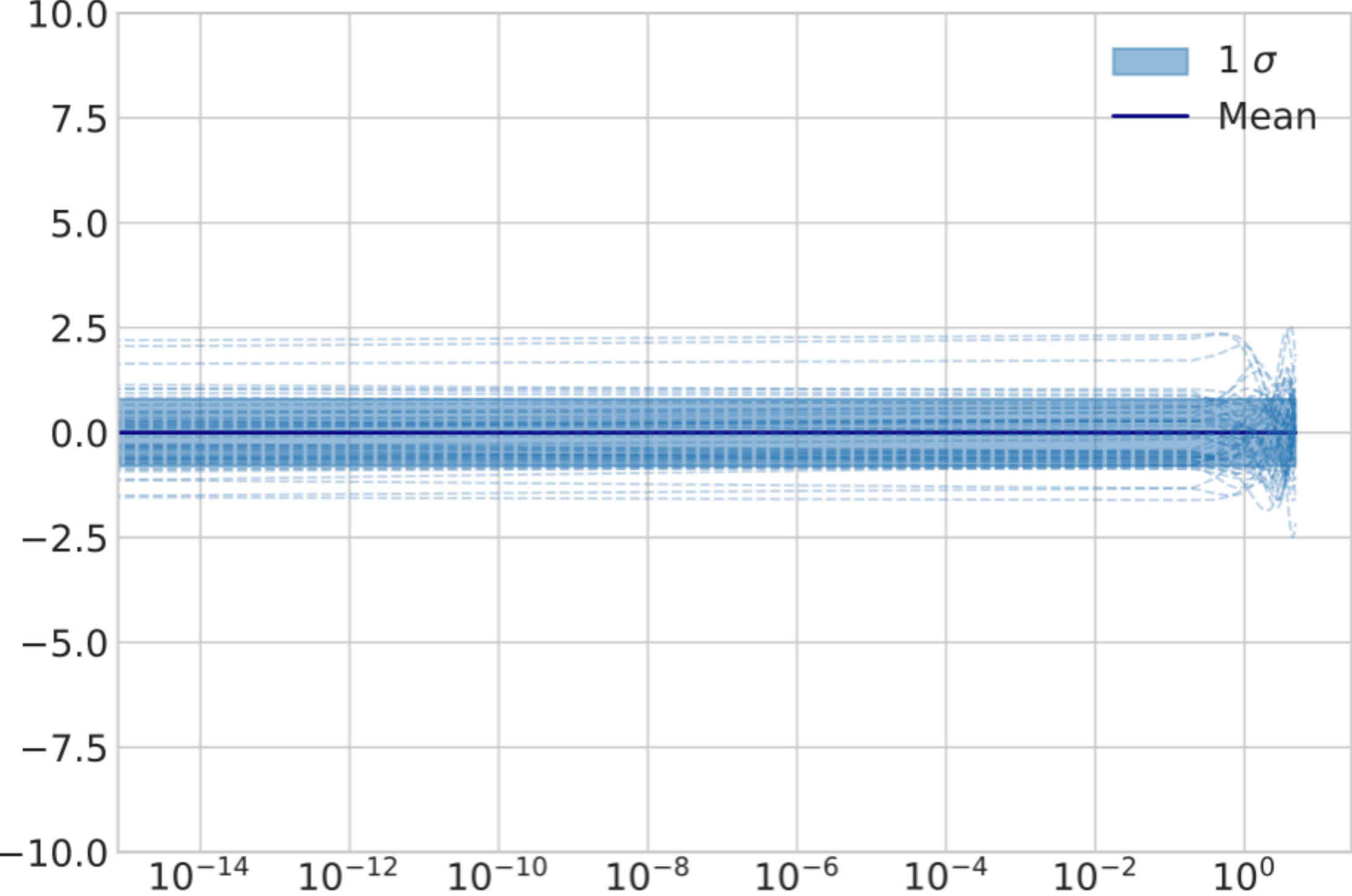
$$p(\mathbf{f}, \theta | \text{data})$$

Prior for PDF

Exponential quadratic

$$m(x) = 0$$

$$k(x, y) = \sigma^2 \exp \left[-\frac{(x - y)^2}{l^2} \right]$$



Prior for PDF

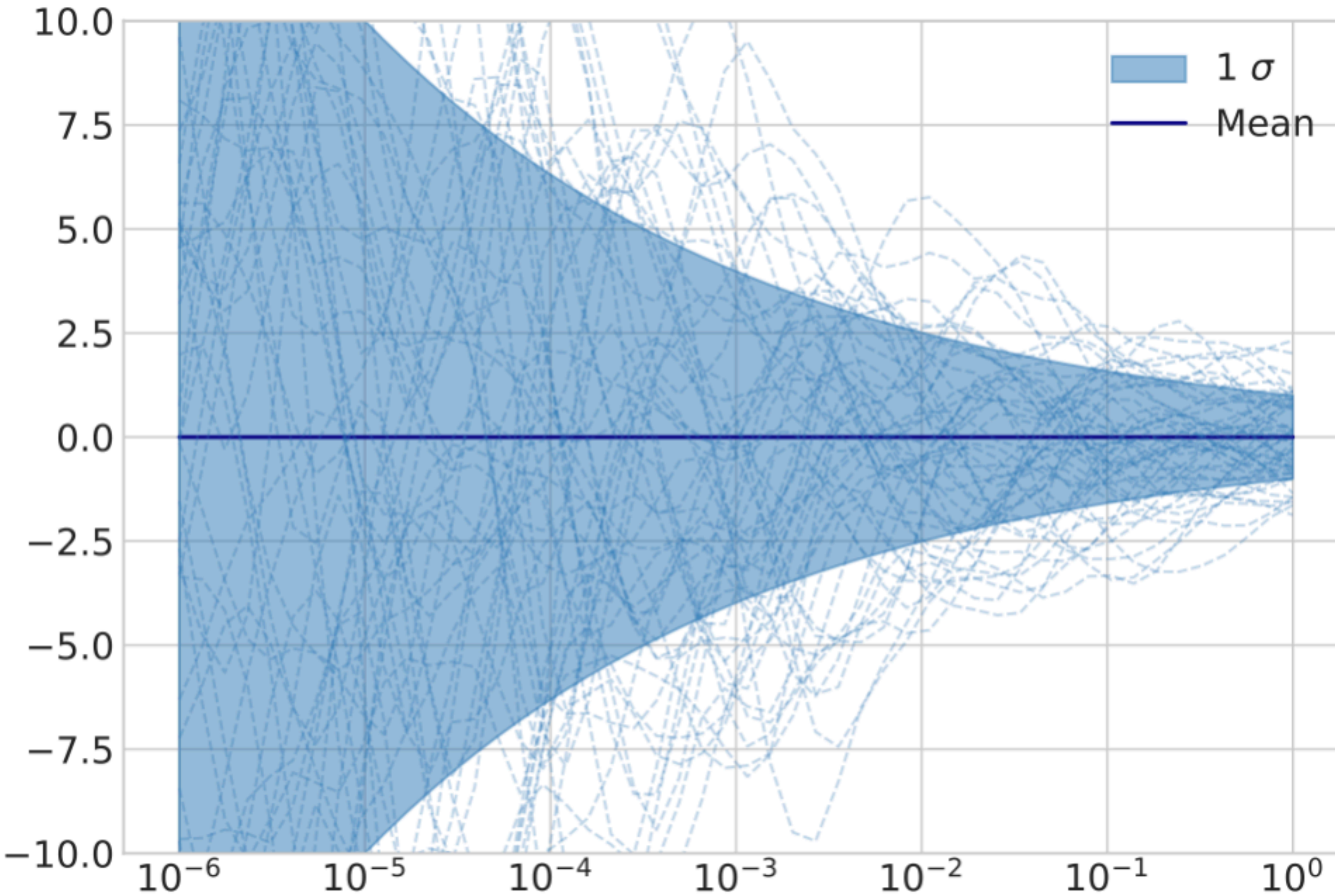
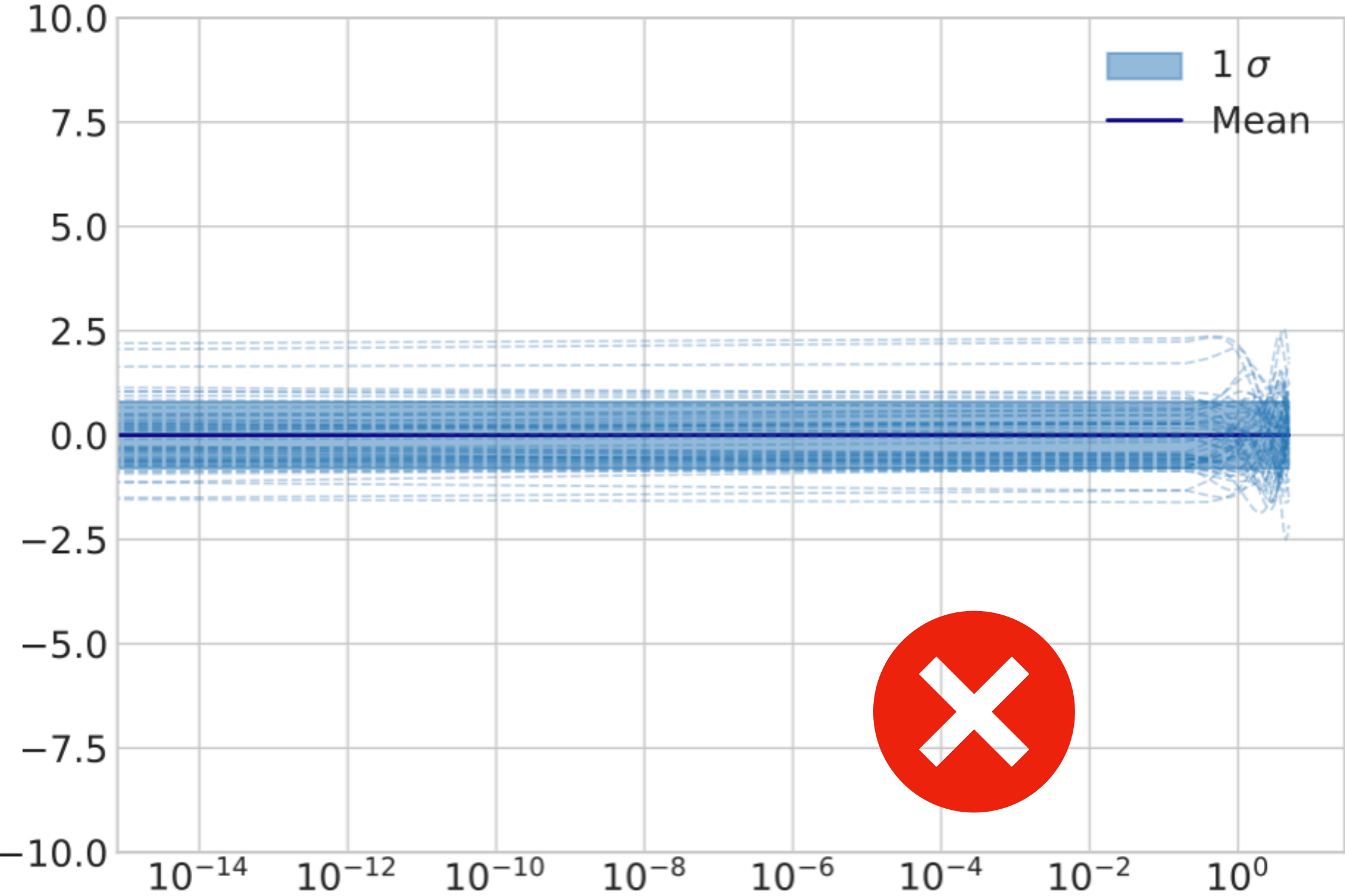
Exponential quadratic

$$m(x) = 0$$

$$k(x, y) = \sigma^2 \exp \left[-\frac{(x - y)^2}{l^2} \right]$$

Gibbs Kernel

$$\tilde{k}(x, y) \propto x^\alpha y^\alpha \sigma^2 \exp \left[-\frac{(x - y)^2}{l^2(x) + l^2(y)} \right] \quad \text{with} \quad l(x) = l_0 x$$



Prior for PDF

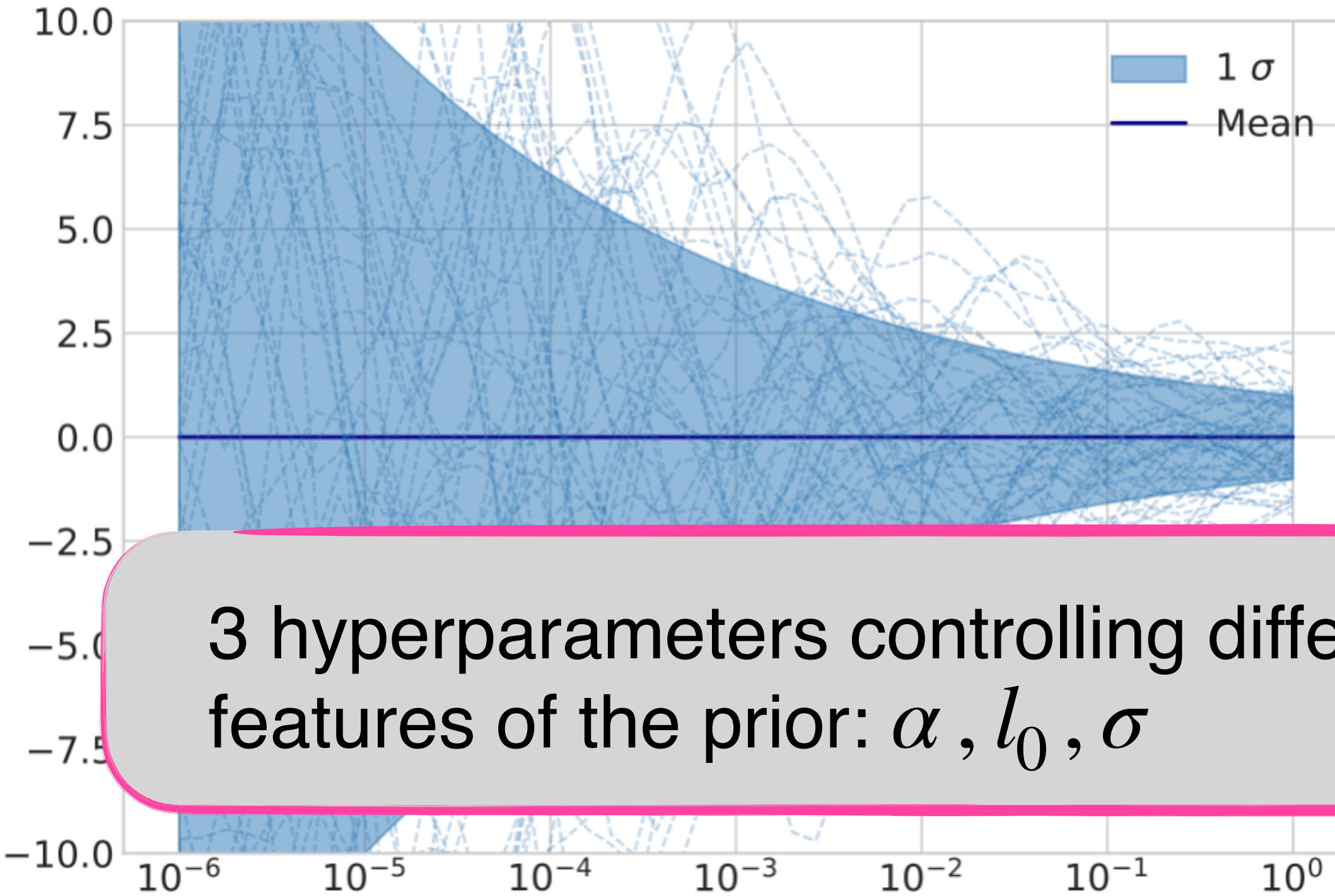
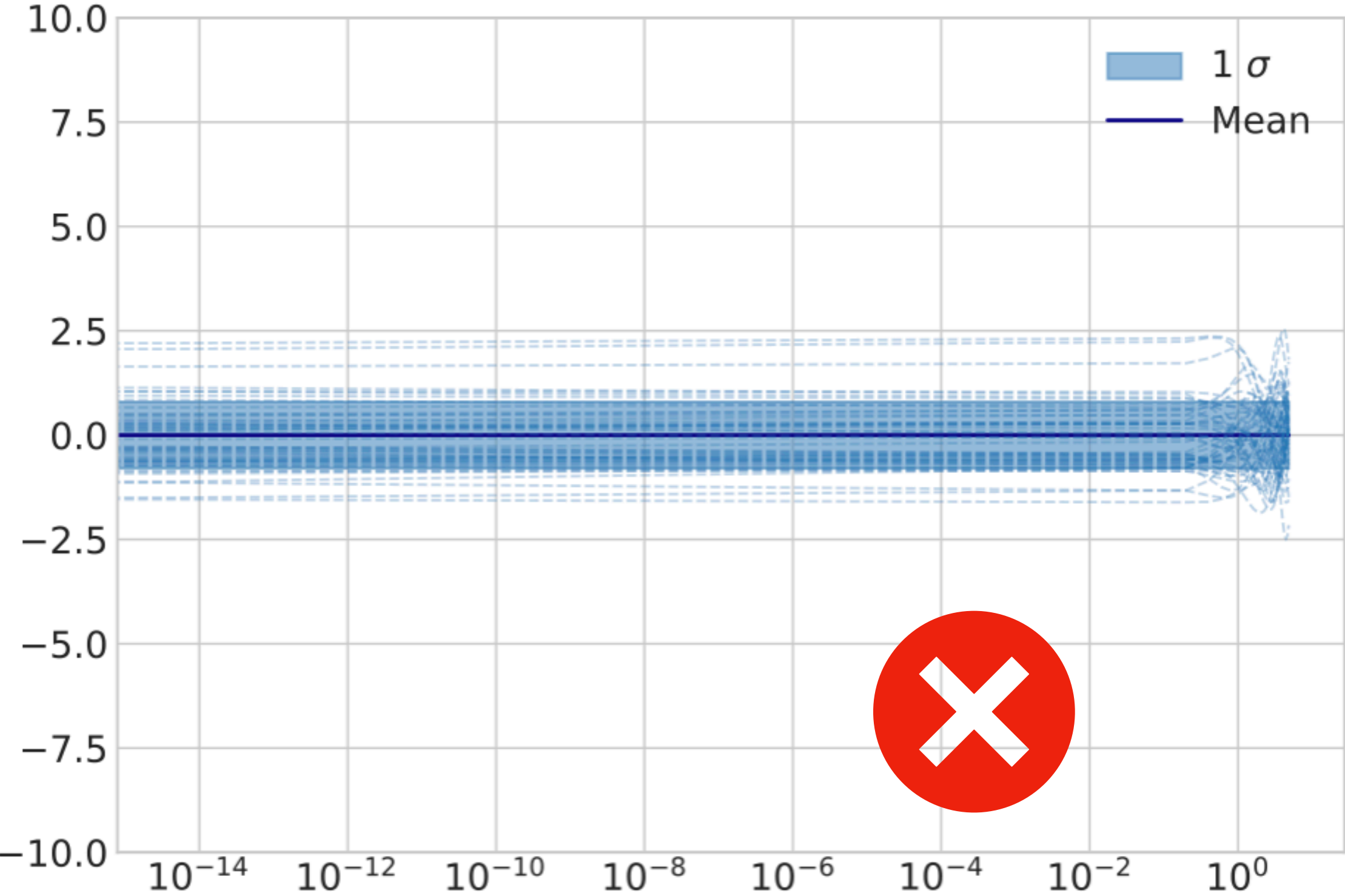
Exponential quadratic

$$m(x) = 0$$

$$k(x, y) = \sigma^2 \exp \left[-\frac{(x - y)^2}{l^2} \right]$$

Gibbs Kernel

$$\tilde{k}(x, y) \propto x^\alpha y^\alpha \sigma^2 \exp \left[-\frac{(x - y)^2}{l^2(x) + l^2(y)} \right] \quad \text{with} \quad l(x) = l_0 x$$



Example: PDFs from DIS

$$F(x, Q^2) = \sum_i \int_x^1 dy C_i \left(\frac{x}{y}, \frac{Q}{\mu}, \alpha_s \right) f_i(y, \mu) \quad \xrightarrow{\text{Introduce an interpolation basis for } f} \quad F_i = \sum_{\alpha} (FK)_{i\alpha} f(x_{\alpha}) = FK \mathbf{f}$$

Gaussian variable representing PDF on interpolation points \mathbf{x}

$$\mathcal{O} = FK \mathbf{f}$$

\mathbf{f}^*

Gaussian variable representing PDF on any set of points \mathbf{x}^*

$K(x, y; \theta)$

Function modelling correlation

$y, \epsilon \sim N(0, C_y)$

Data and corresponding experimental error

Gaussian inference (fixed hyperparameters)

Consider stochastic variable

$$\begin{pmatrix} \mathbf{f}^* \\ FK\mathbf{f} \end{pmatrix} \sim \mathcal{N} \left(0, \begin{pmatrix} K_{\mathbf{x}^*\mathbf{x}^*} & K_{\mathbf{x}^*\mathbf{x}} FK^T \\ FK K_{\mathbf{xx}^*} & FK K_{\mathbf{xx}} FK^T \end{pmatrix} \right)$$

Compute the distribution

$$p(\mathbf{f}^* | FK\mathbf{f} + \epsilon = y, \theta)$$

This is a gaussian distribution. Its mean and covariance can be computed analytically

$$\tilde{\mathbf{m}}^* = \mathbf{m} + K_{\mathbf{x}^*\mathbf{x}} FK^T \left(FK K_{\mathbf{xx}} FK^T + C_y \right)^+ (\mathbf{y} - \mathbf{m})$$

$$\tilde{K}^* = K_{\mathbf{x}^*\mathbf{x}^*} - K_{\mathbf{x}^*\mathbf{x}} FK^T \left(FK K_{\mathbf{xx}} FK^T + C_y \right)^+ FK K_{\mathbf{xx}^*}$$

- We generate pseudo-data for NNPDF4.0 DIS datasets
(NNPDF4.0 as underlying law)
- Gibbs Kernel with fixed hyperparameters

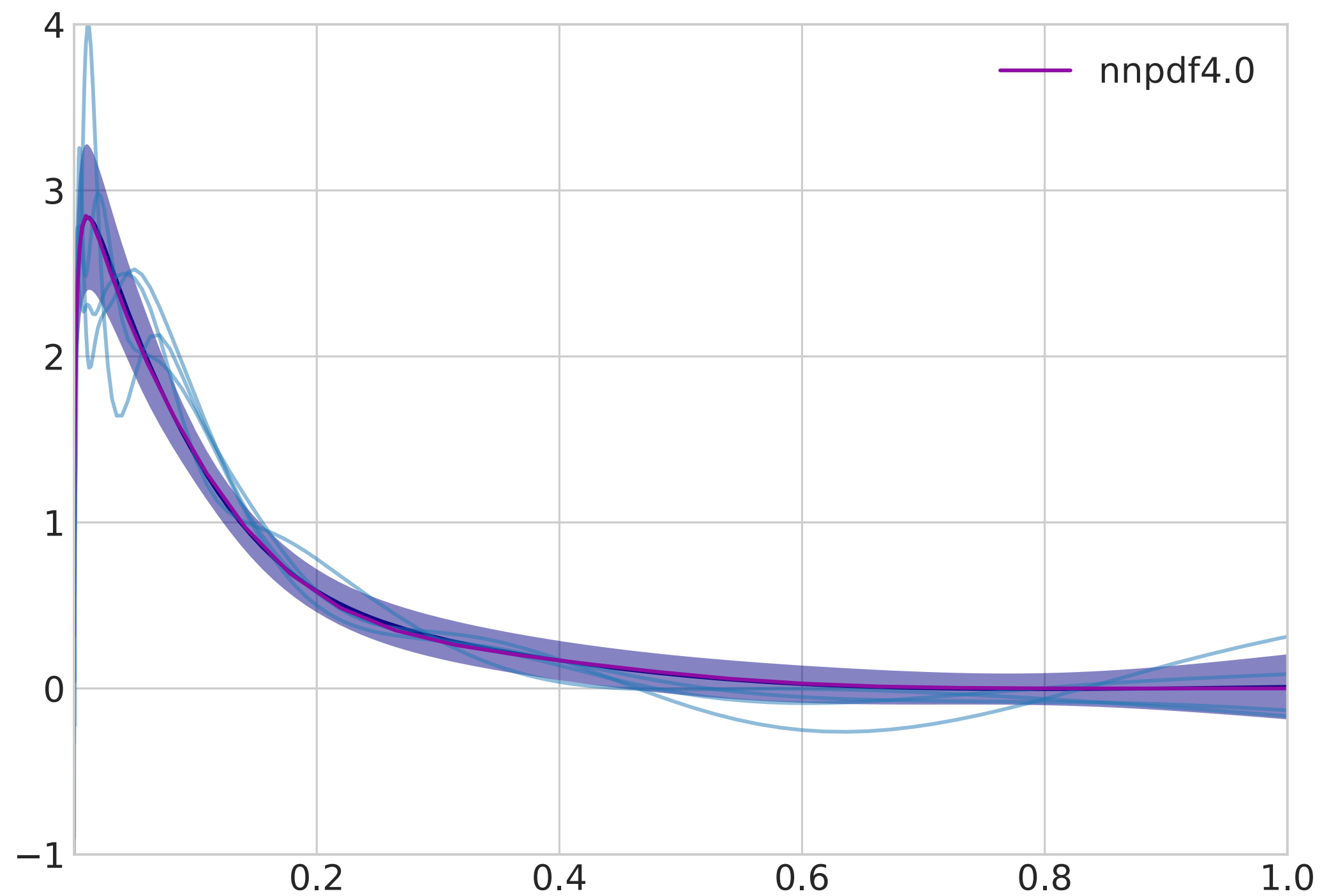
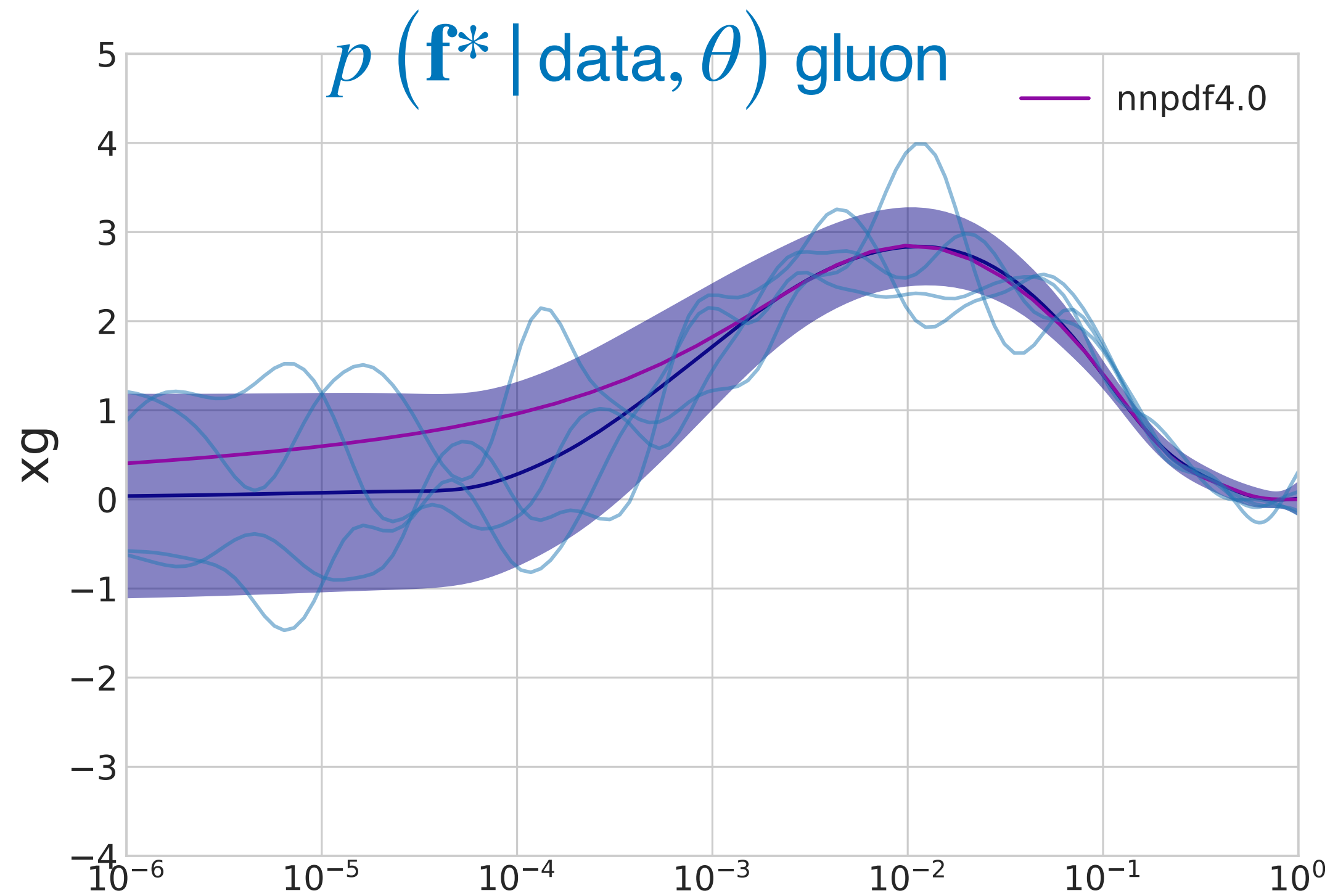
$$\tilde{\mathbf{m}}^* = \mathbf{m} + K_{\mathbf{x}^*\mathbf{x}} F K^T \left(F K K_{\mathbf{xx}} F K^T + C_y \right)^+ (\mathbf{y} - \mathbf{m})$$

$$\tilde{K}^* = K_{\mathbf{x}^*\mathbf{x}^*} - K_{\mathbf{x}^*\mathbf{x}} F K^T \left(F K K_{\mathbf{xx}} F K^T + C_y \right)^+ F K K_{\mathbf{xx}}^*$$

- We generate pseudo-data for NNPDF4.0 DIS datasets (NNPDF4.0 as underlying law)
- Gibbs Kernel with fixed hyperparameters

$$\tilde{\mathbf{m}}^* = \mathbf{m} + K_{\mathbf{x}^*\mathbf{x}} F K^T \left(F K K_{\mathbf{xx}} F K^T + C_y \right)^+ (\mathbf{y} - \mathbf{m})$$

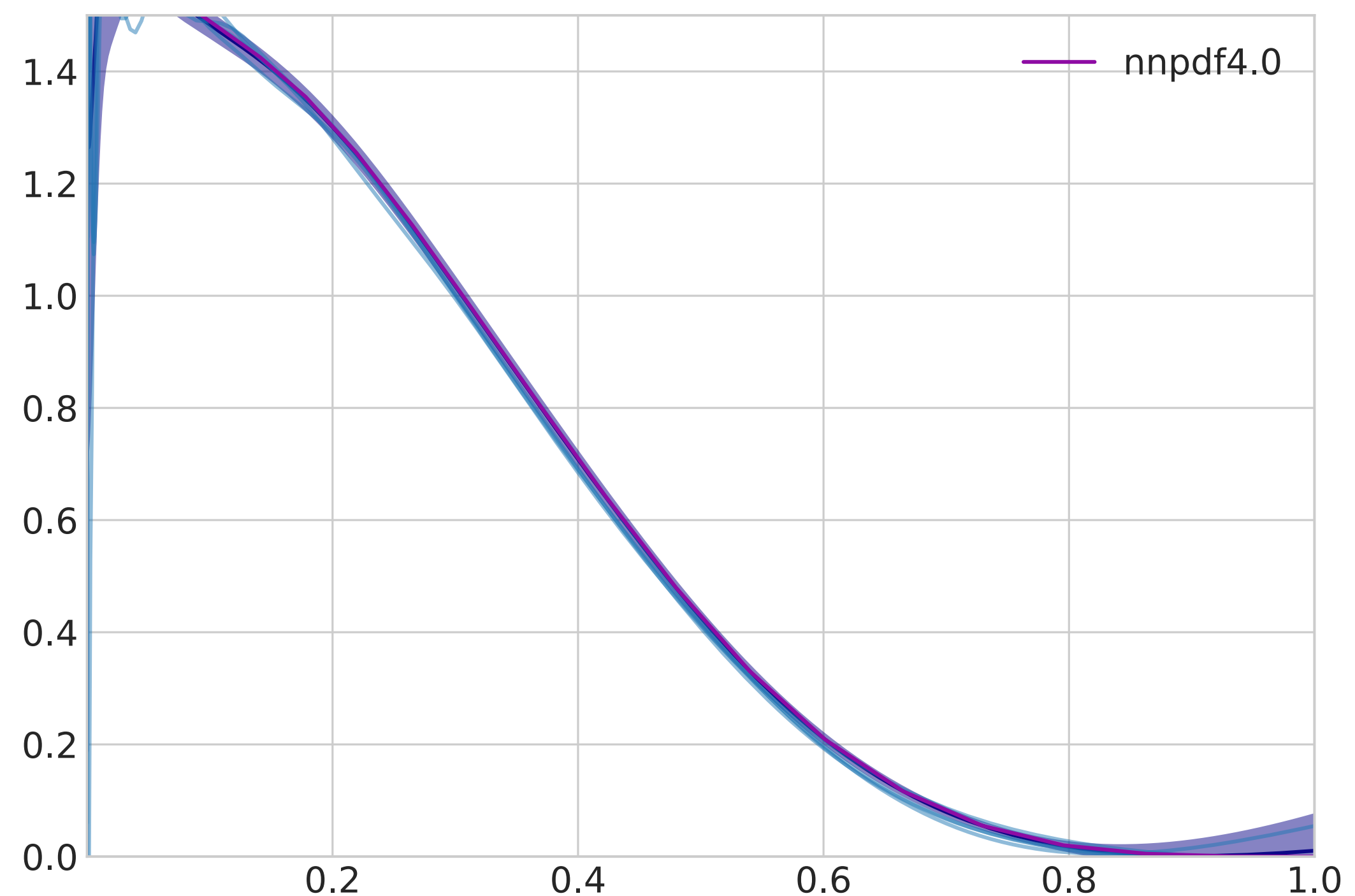
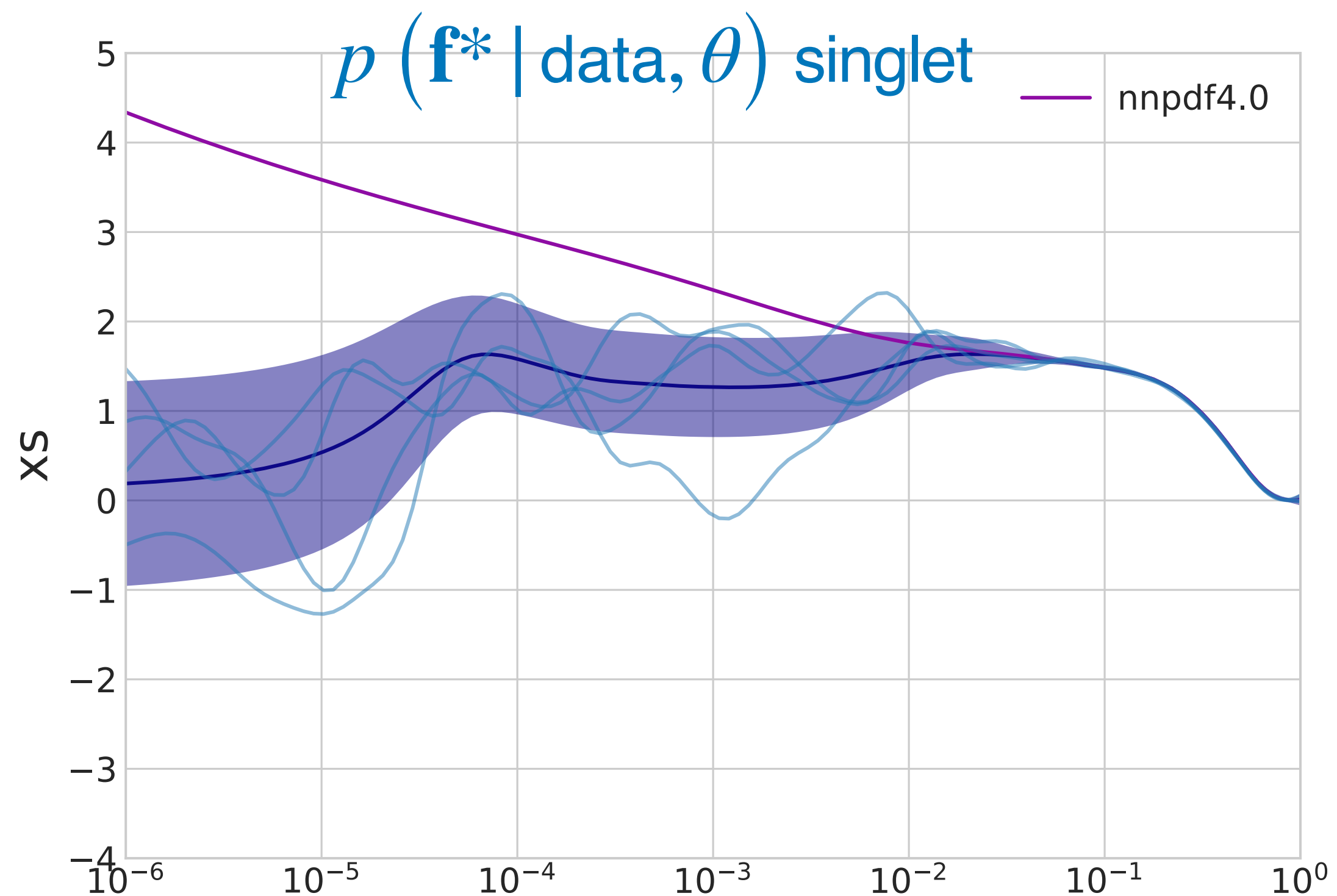
$$\tilde{K}^* = K_{\mathbf{x}^*\mathbf{x}^*} - K_{\mathbf{x}^*\mathbf{x}} F K^T \left(F K K_{\mathbf{xx}} F K^T + C_y \right)^+ F K K_{\mathbf{xx}}^*$$



- We generate pseudo-data for NNPDF4.0 DIS datasets (NNPDF4.0 as underlying law)
- Gibbs Kernel with fixed hyperparameters

$$\tilde{\mathbf{m}}^* = \mathbf{m} + K_{\mathbf{x}^*\mathbf{x}} F K^T \left(F K K_{\mathbf{xx}} F K^T + C_y \right)^+ (\mathbf{y} - \mathbf{m})$$

$$\tilde{K}^* = K_{\mathbf{x}^*\mathbf{x}^*} - K_{\mathbf{x}^*\mathbf{x}} F K^T \left(F K K_{\mathbf{xx}} F K^T + C_y \right)^+ F K K_{\mathbf{xx}}^*$$



Inference on the hyperparameters

Hyperparameters are also stochastic variables

Posterior on the hyperparameters given the data

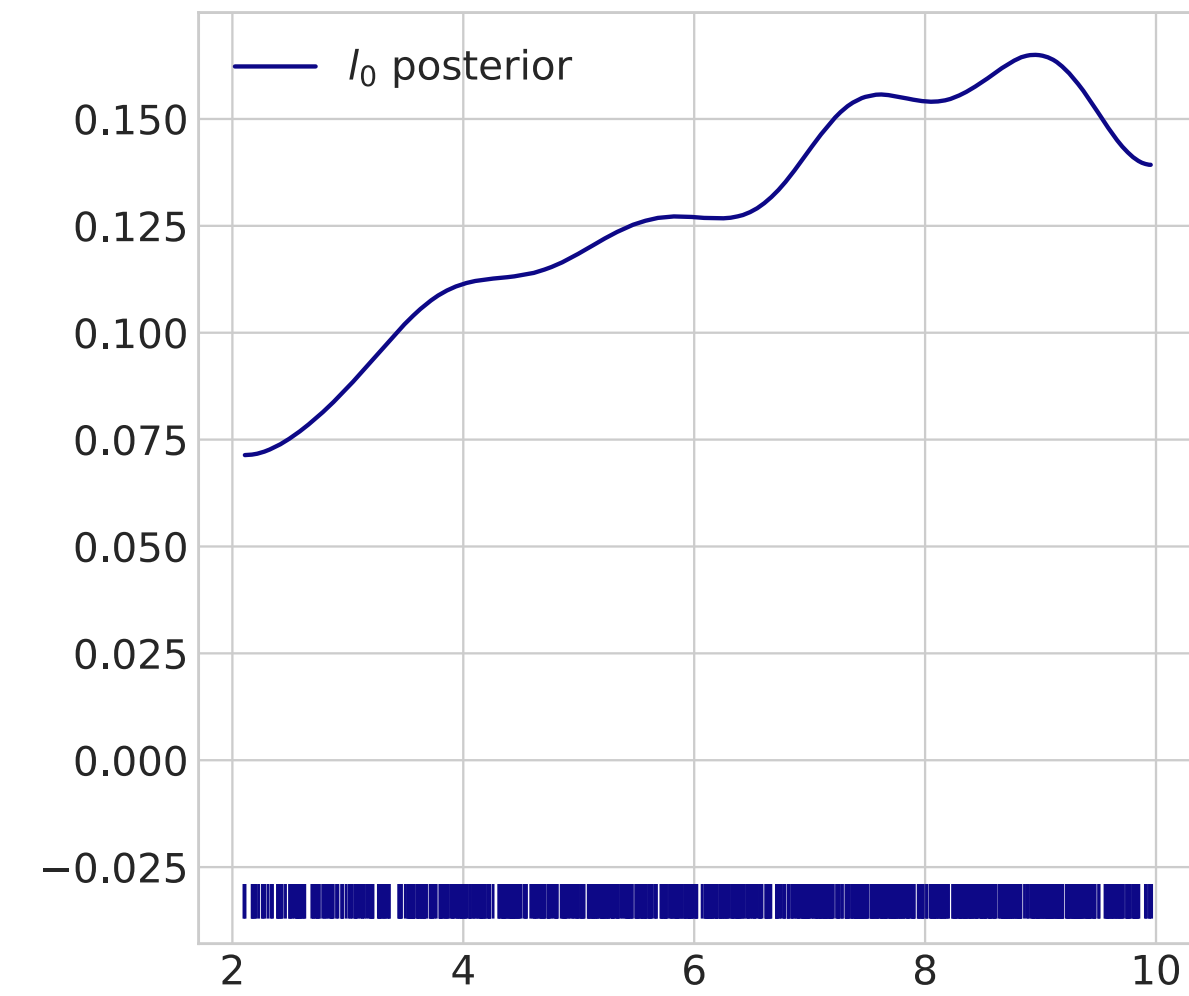
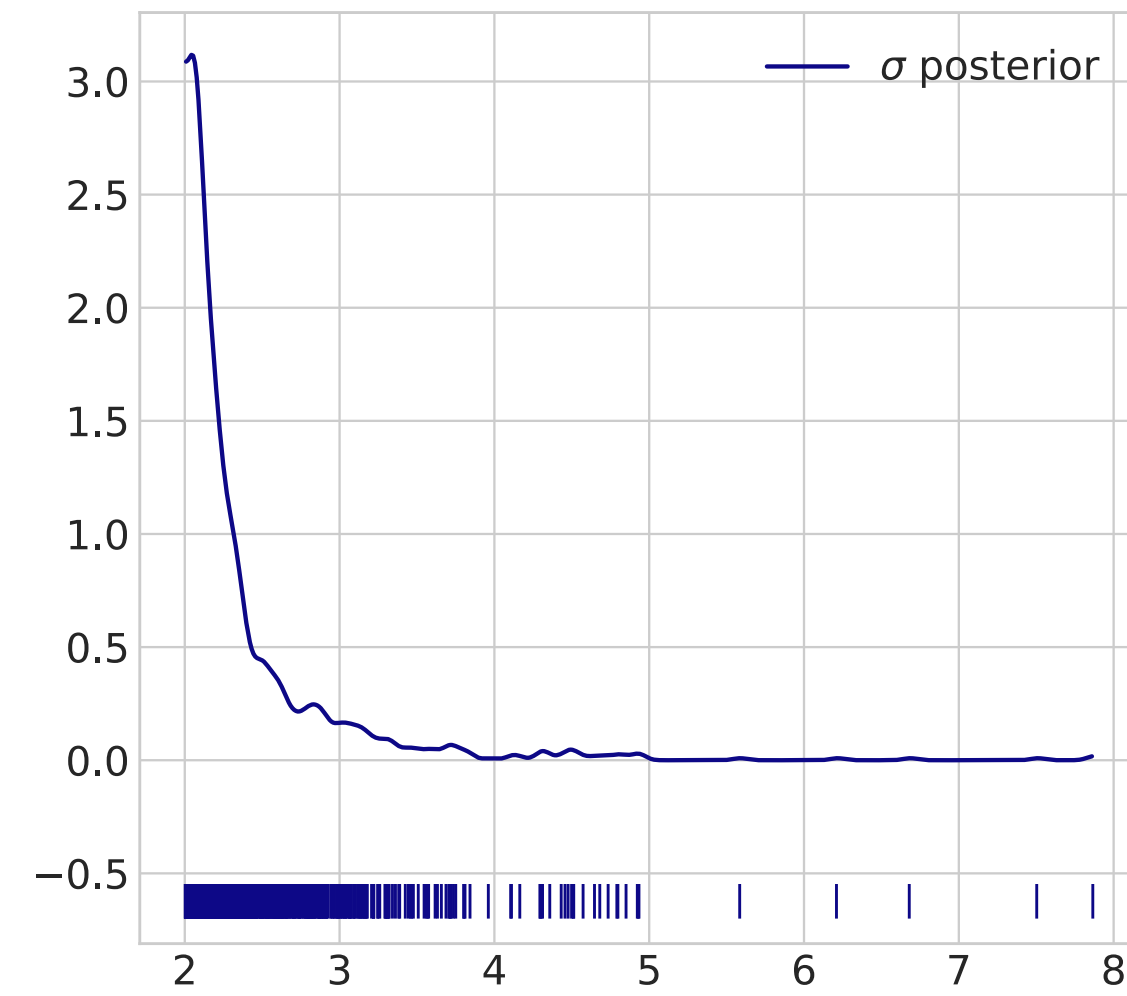
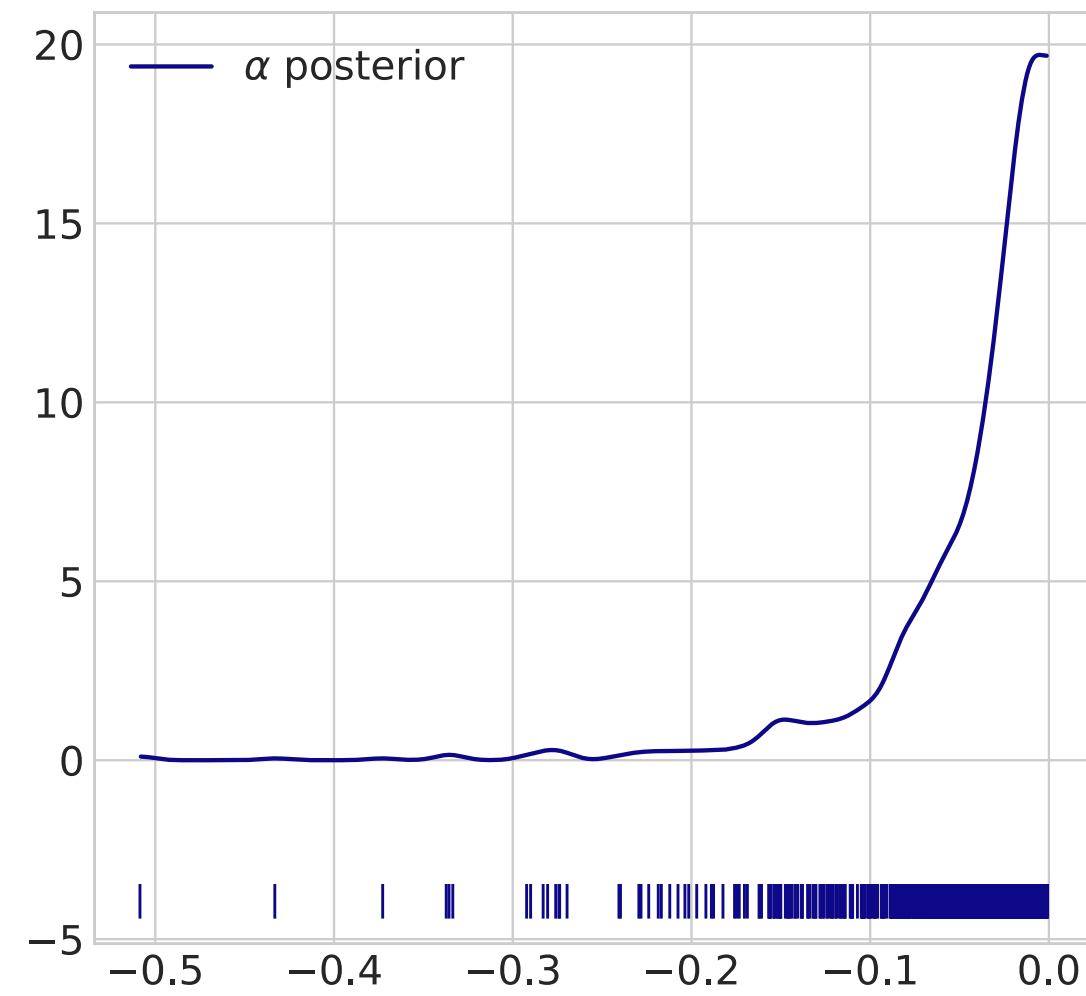
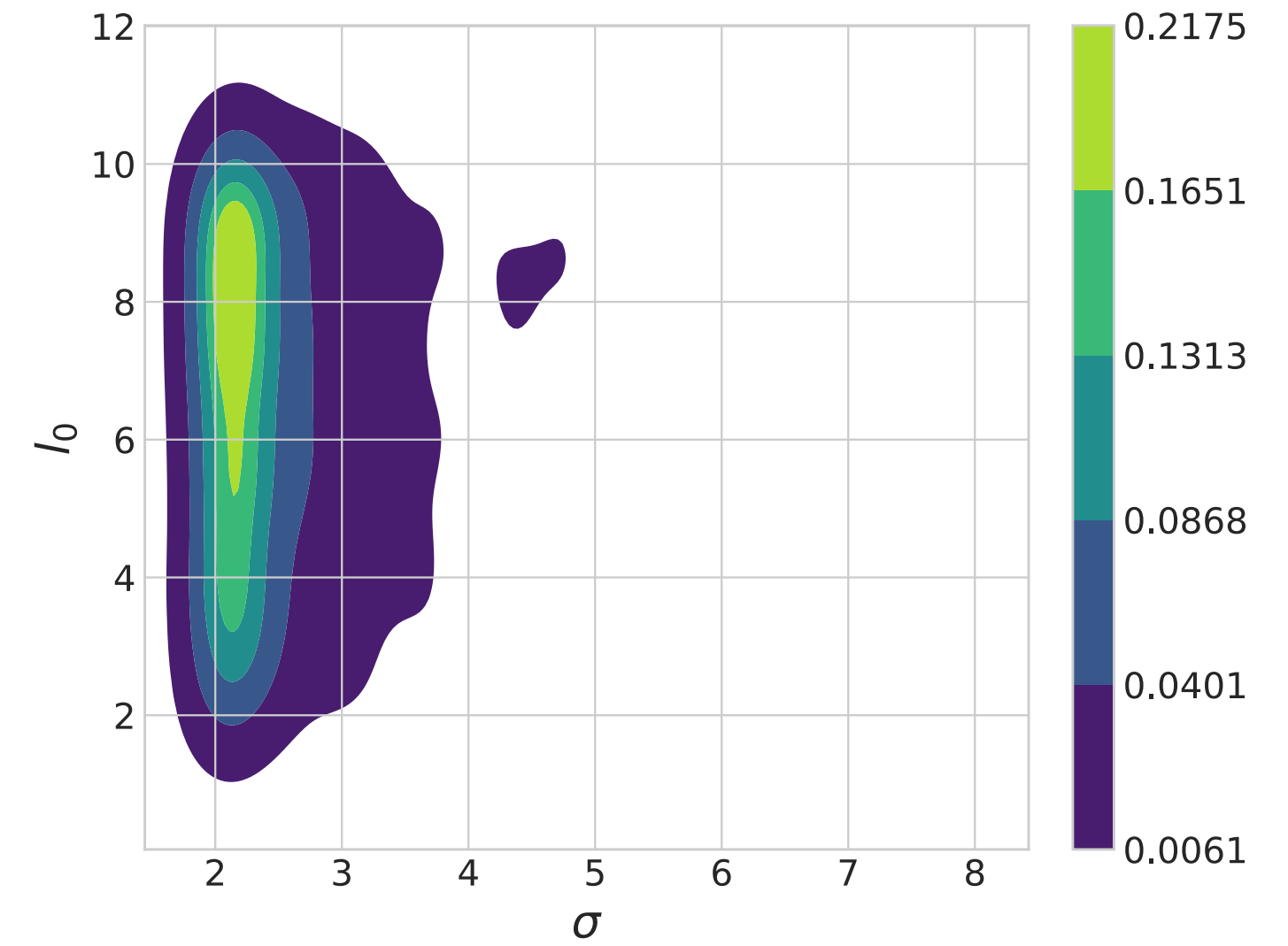
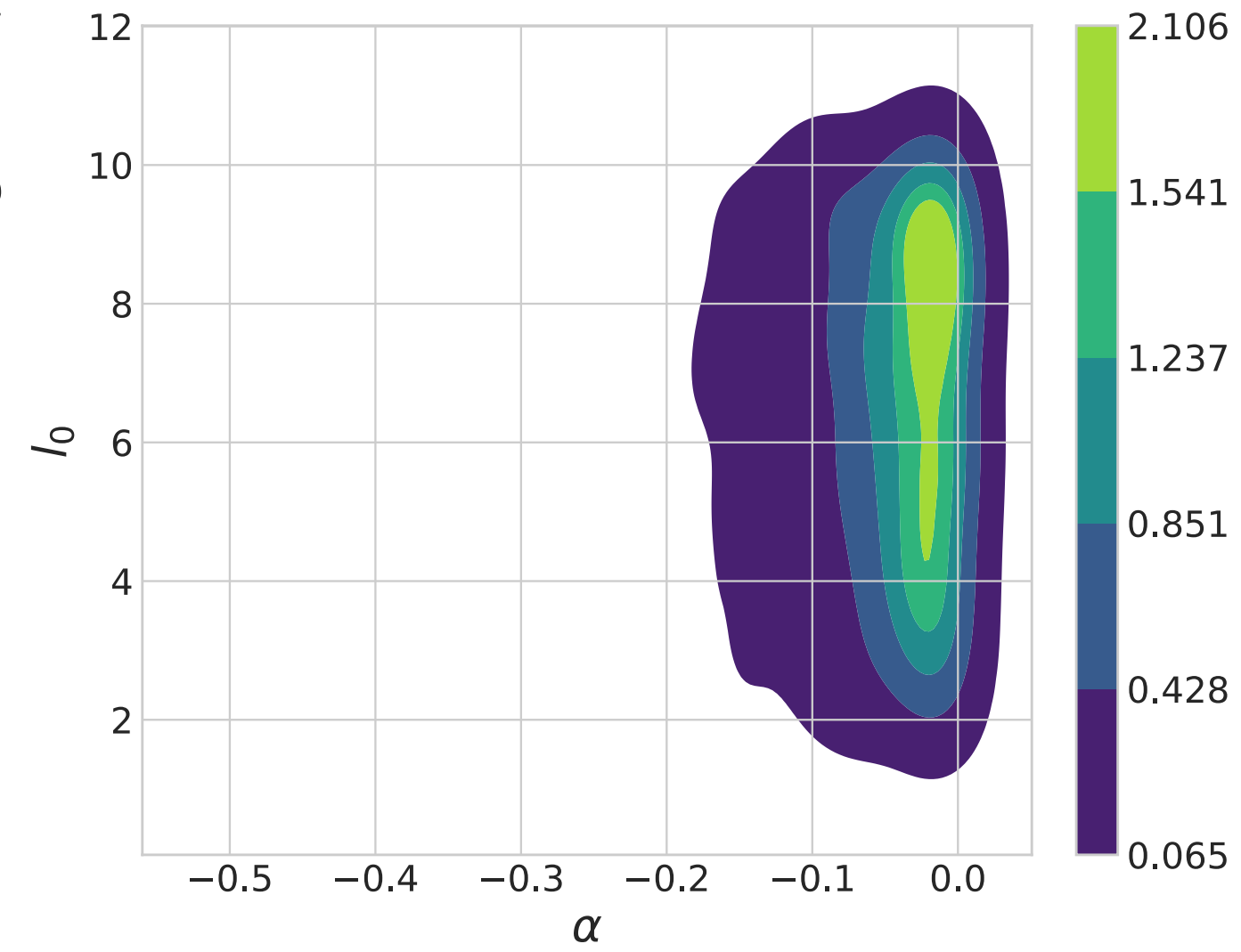
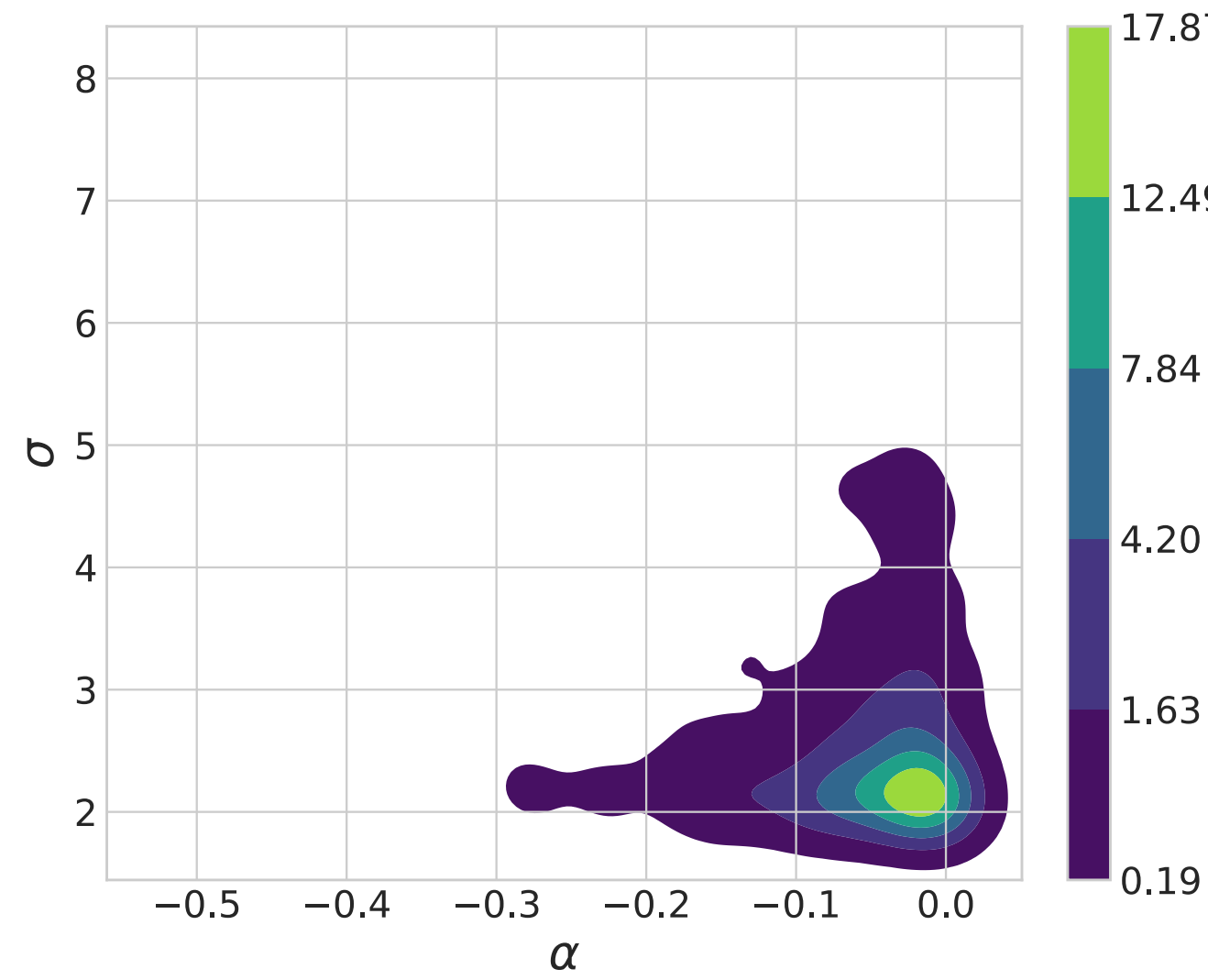
$$p(\mathbf{f}^*, \theta | \text{data}) = p(\mathbf{f}^* | \theta, \text{data}) p(\theta | \text{data})$$

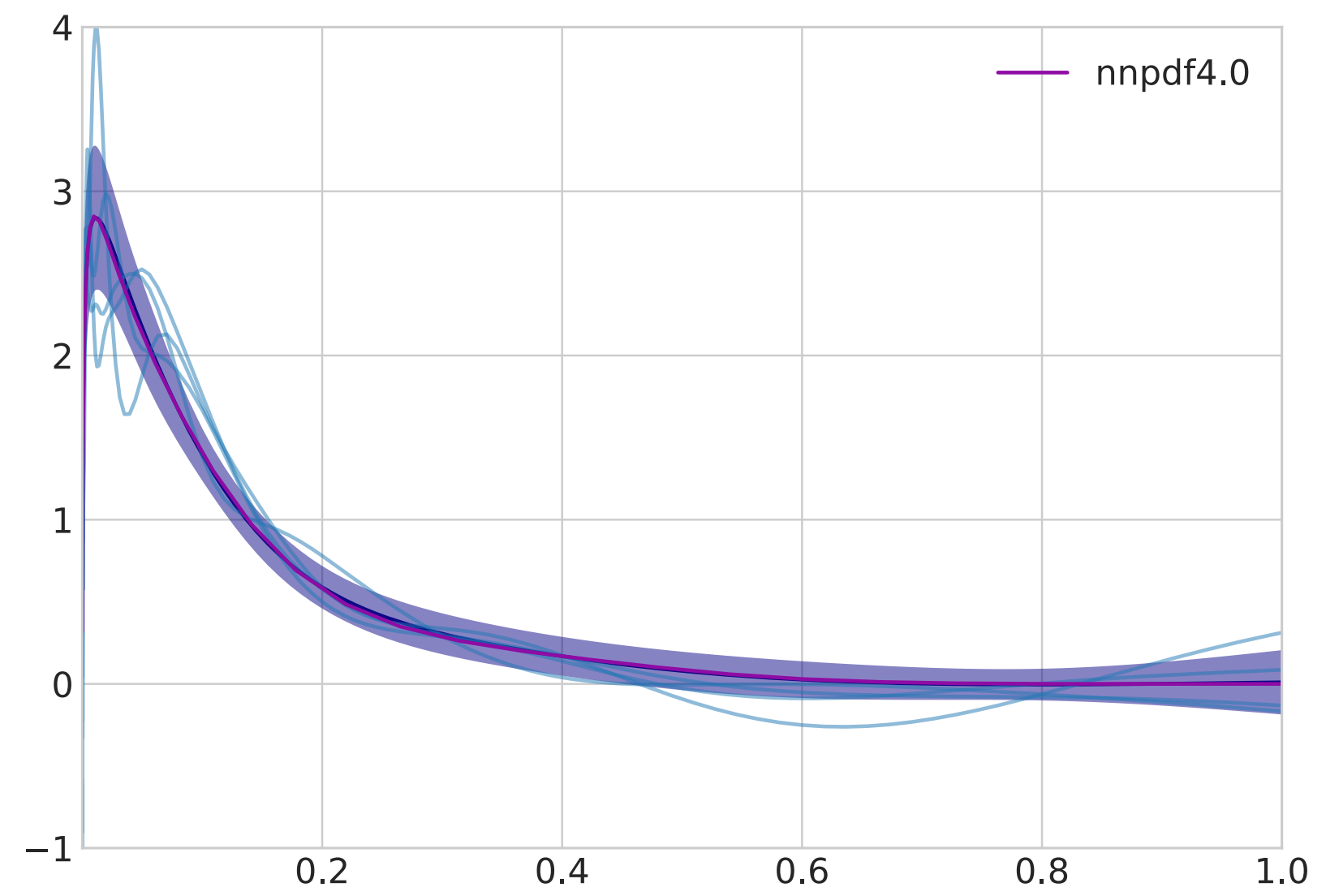
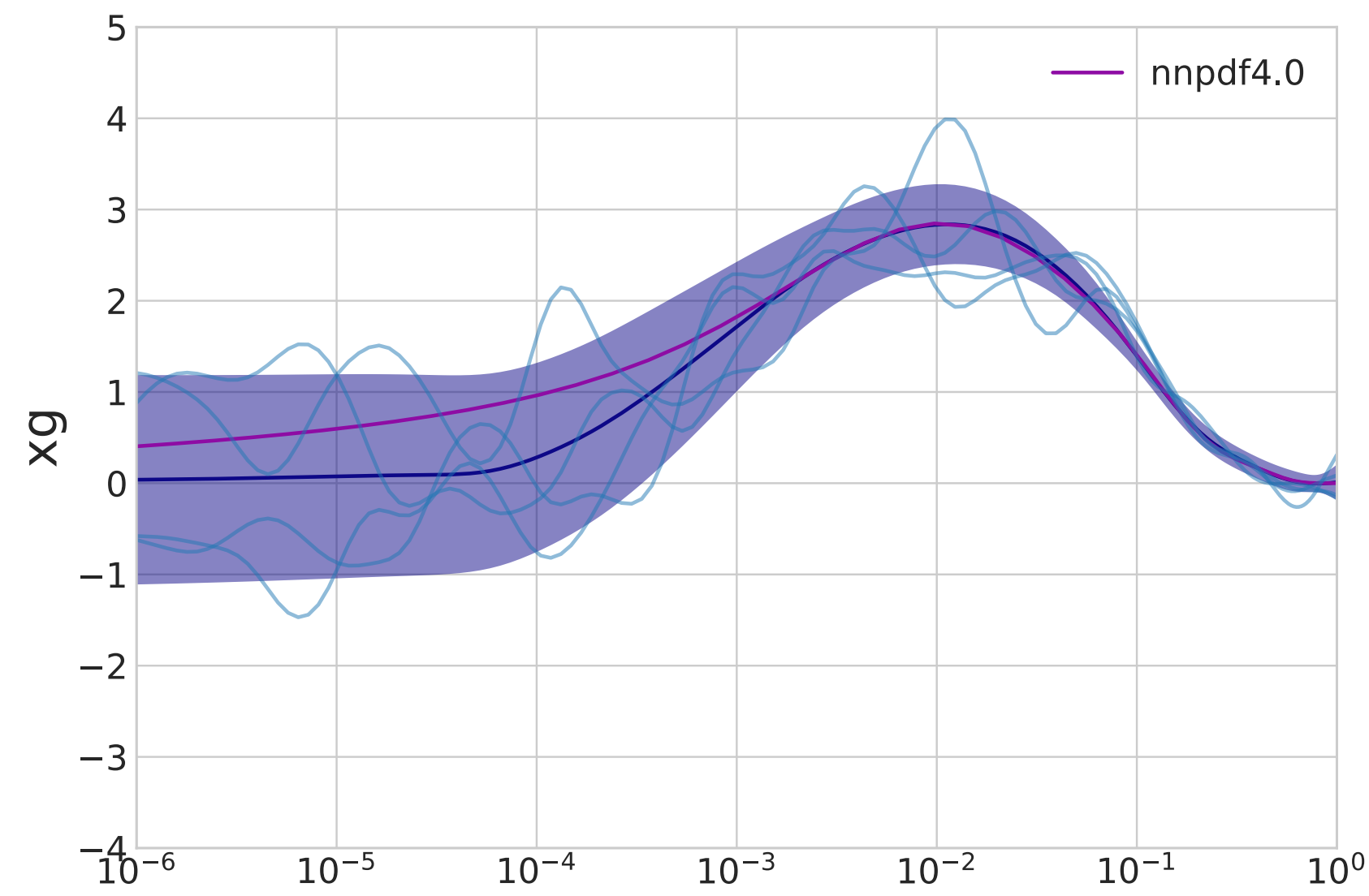
Joint probability distribution of \mathbf{f}^* and θ

$$\propto p(\text{data} | \theta) p_{\theta}(\theta)$$

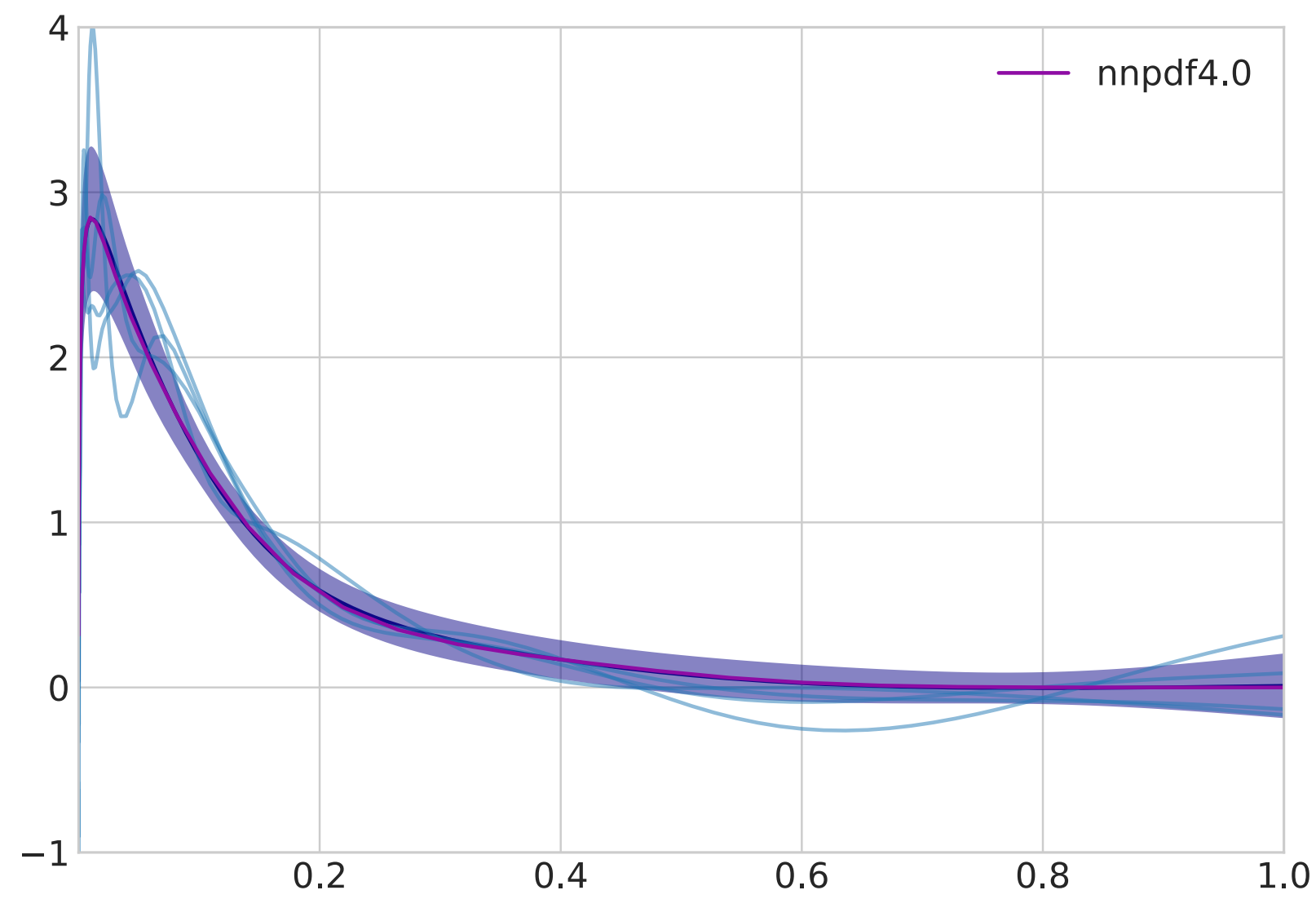
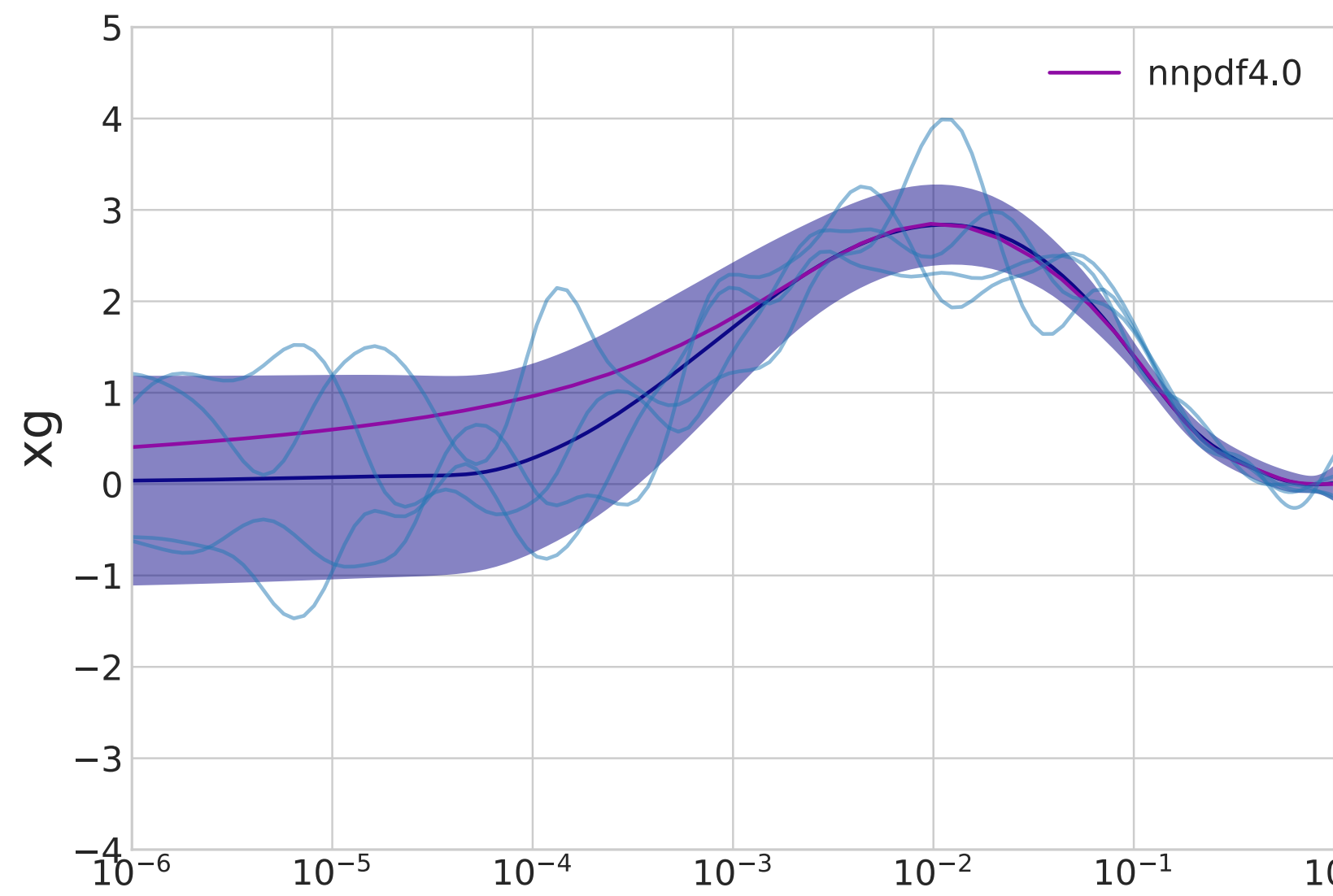
We can sample from $p(\theta | \text{data})$ running a MCMC algorithm

Posterior for hyperparameters (gluon)



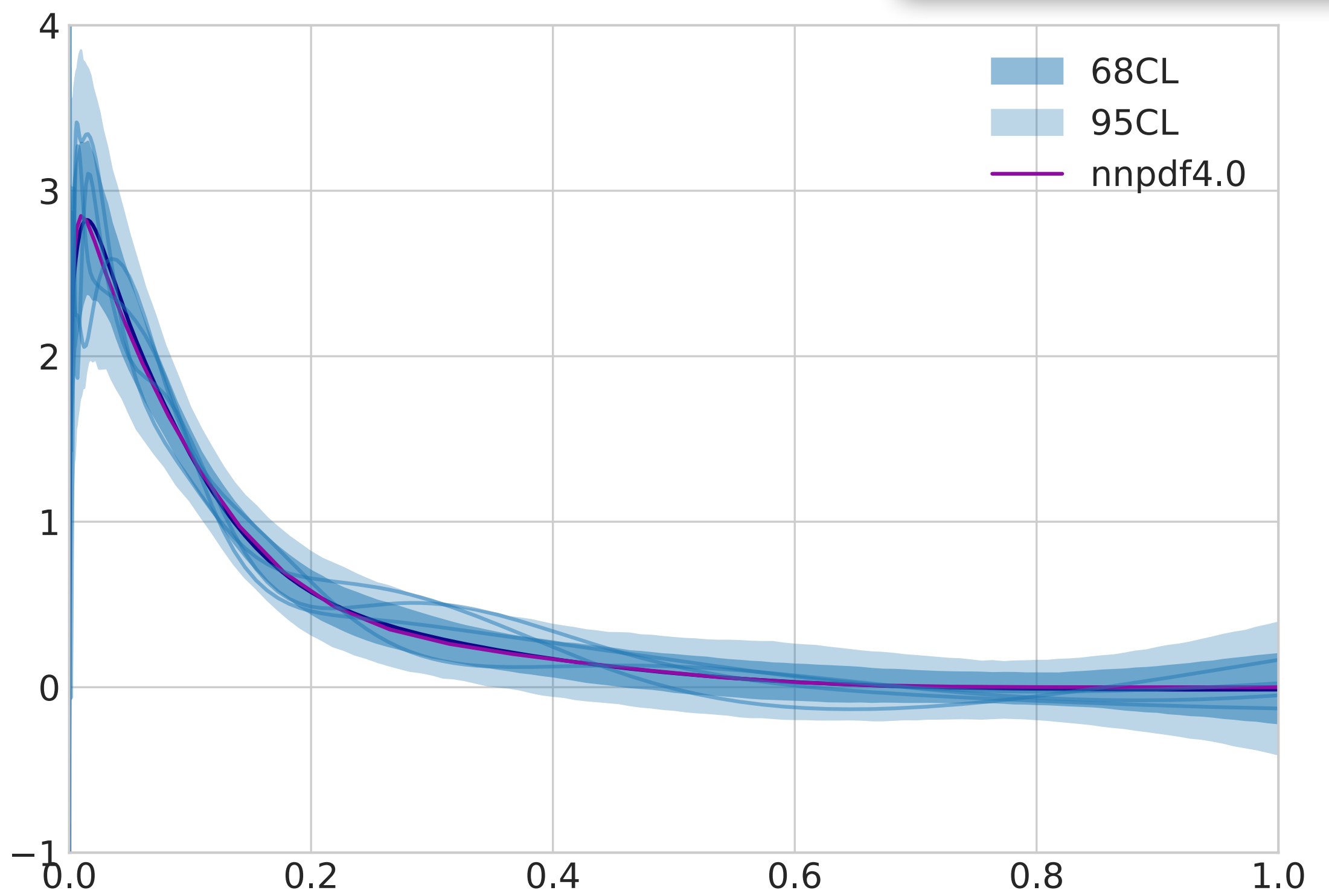
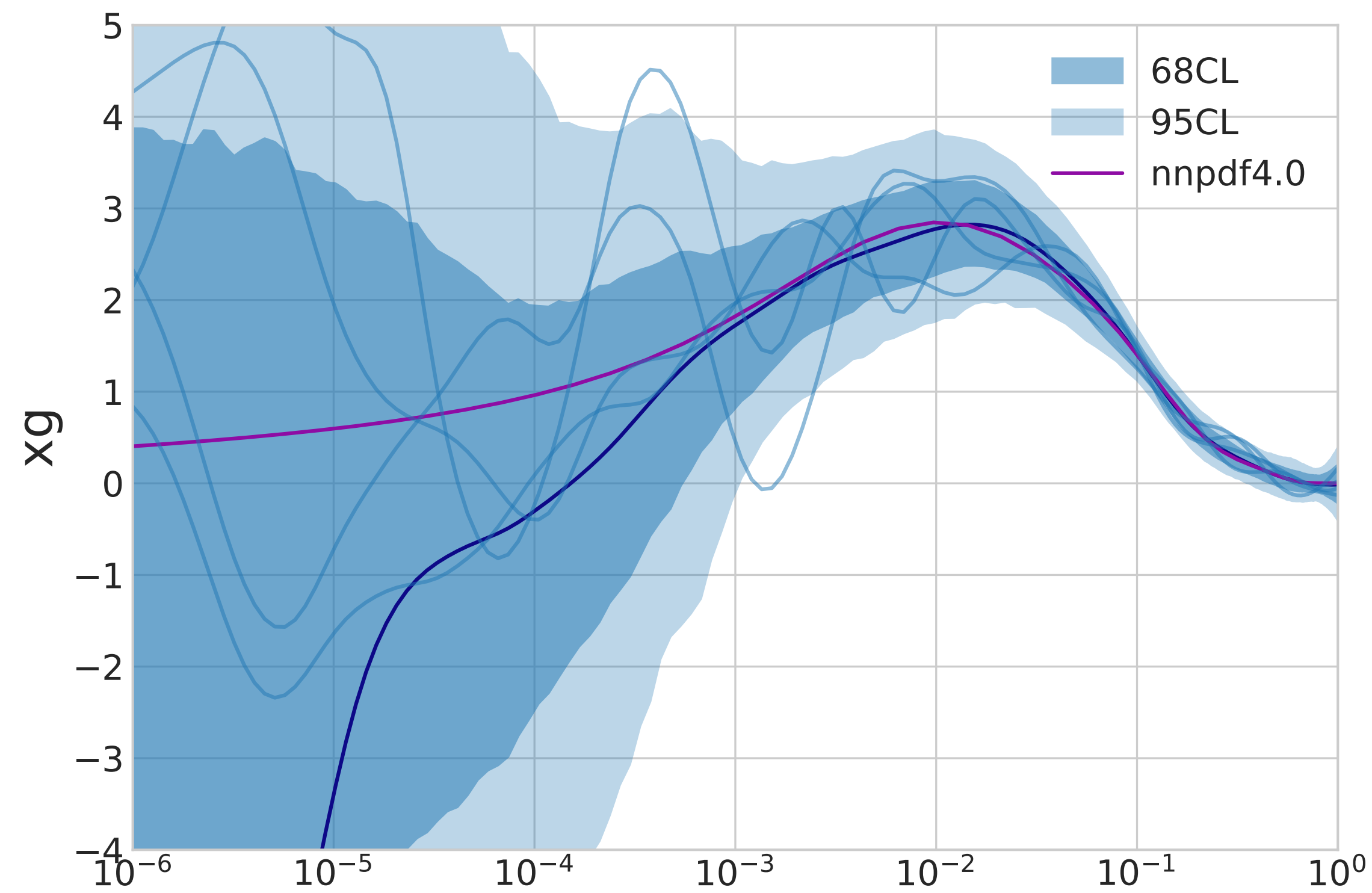


Fixed hyperparameters
 $p(\mathbf{f}^* | \text{data}, \theta)$

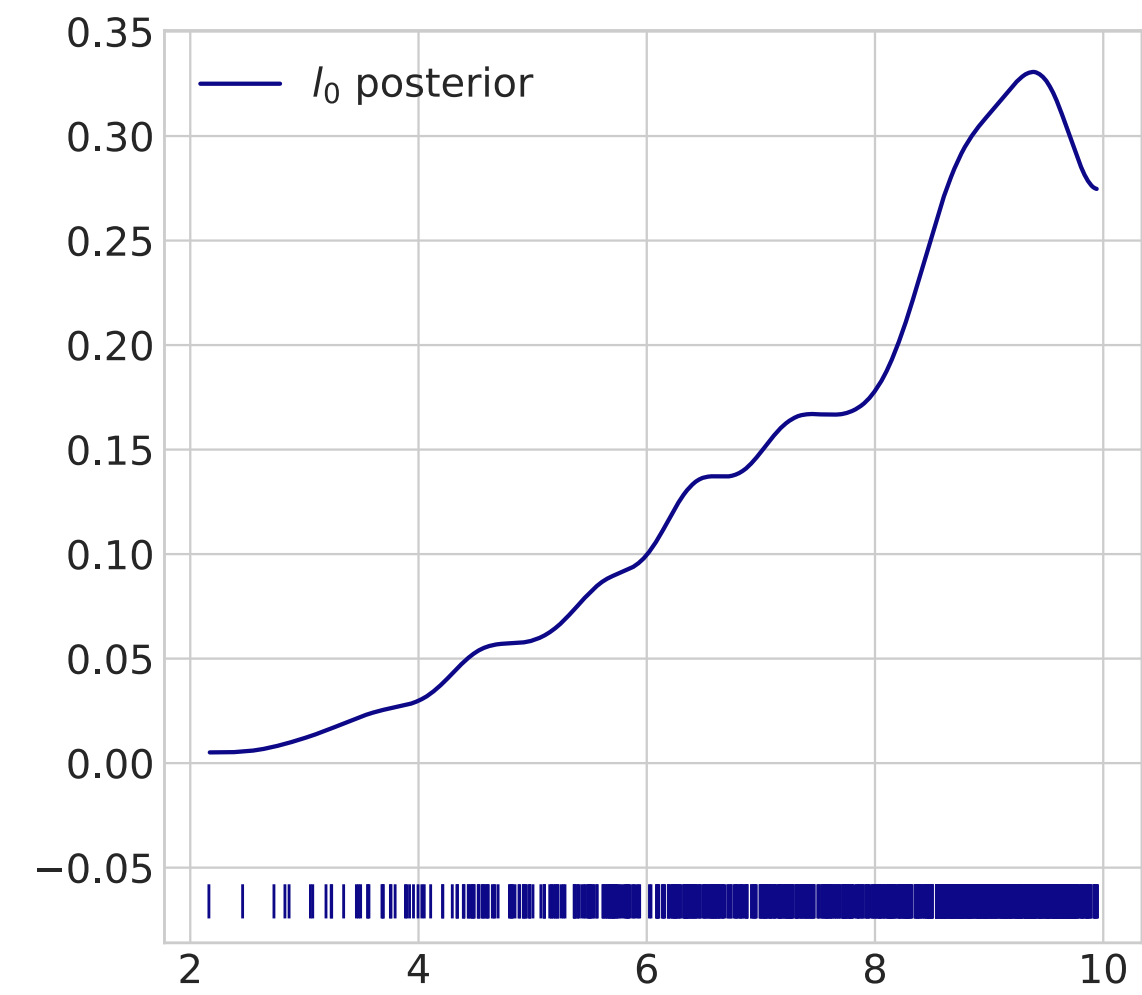
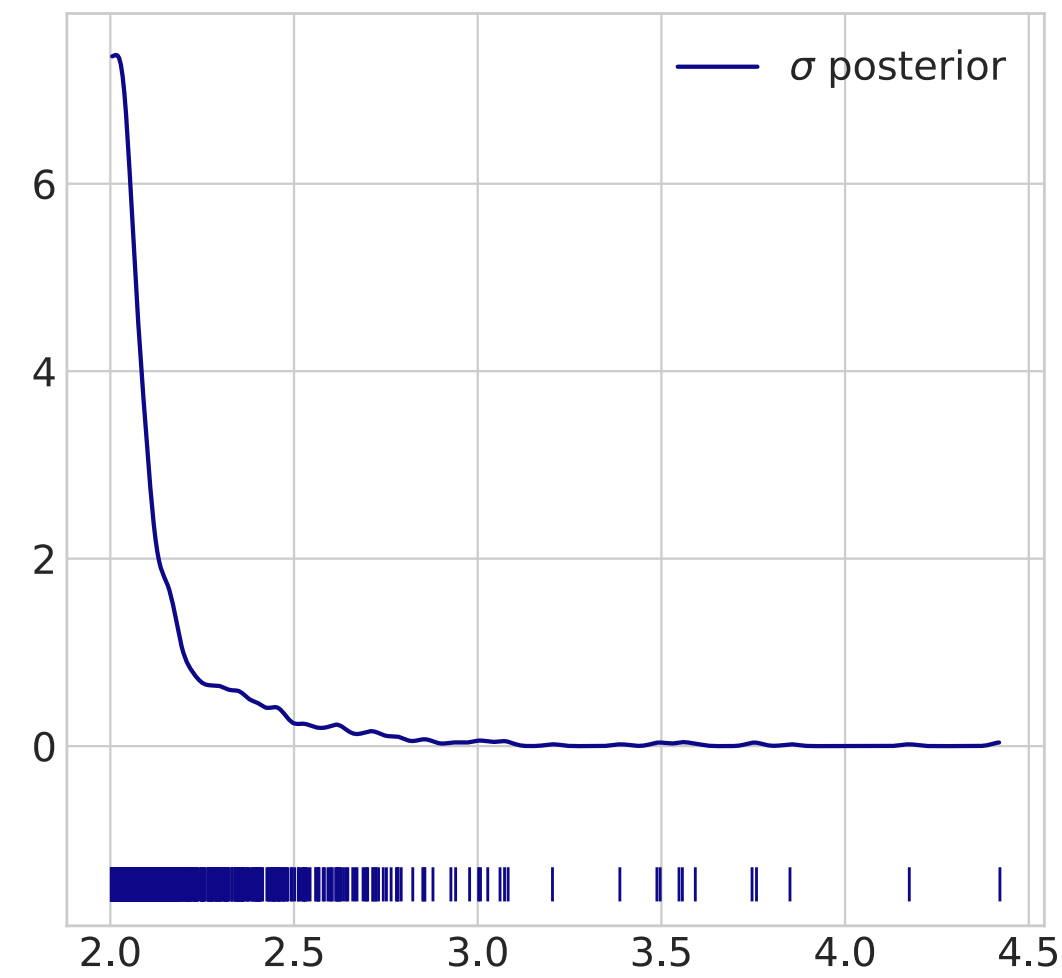
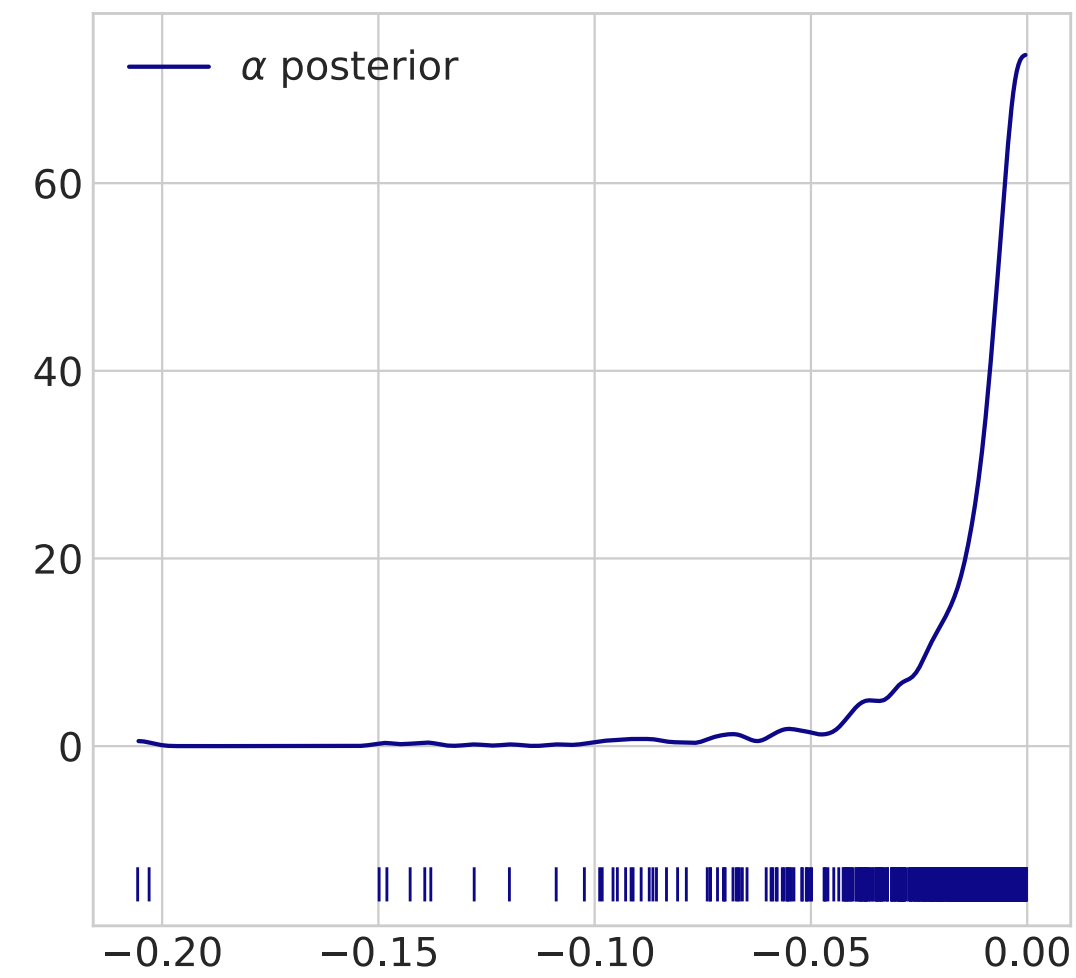
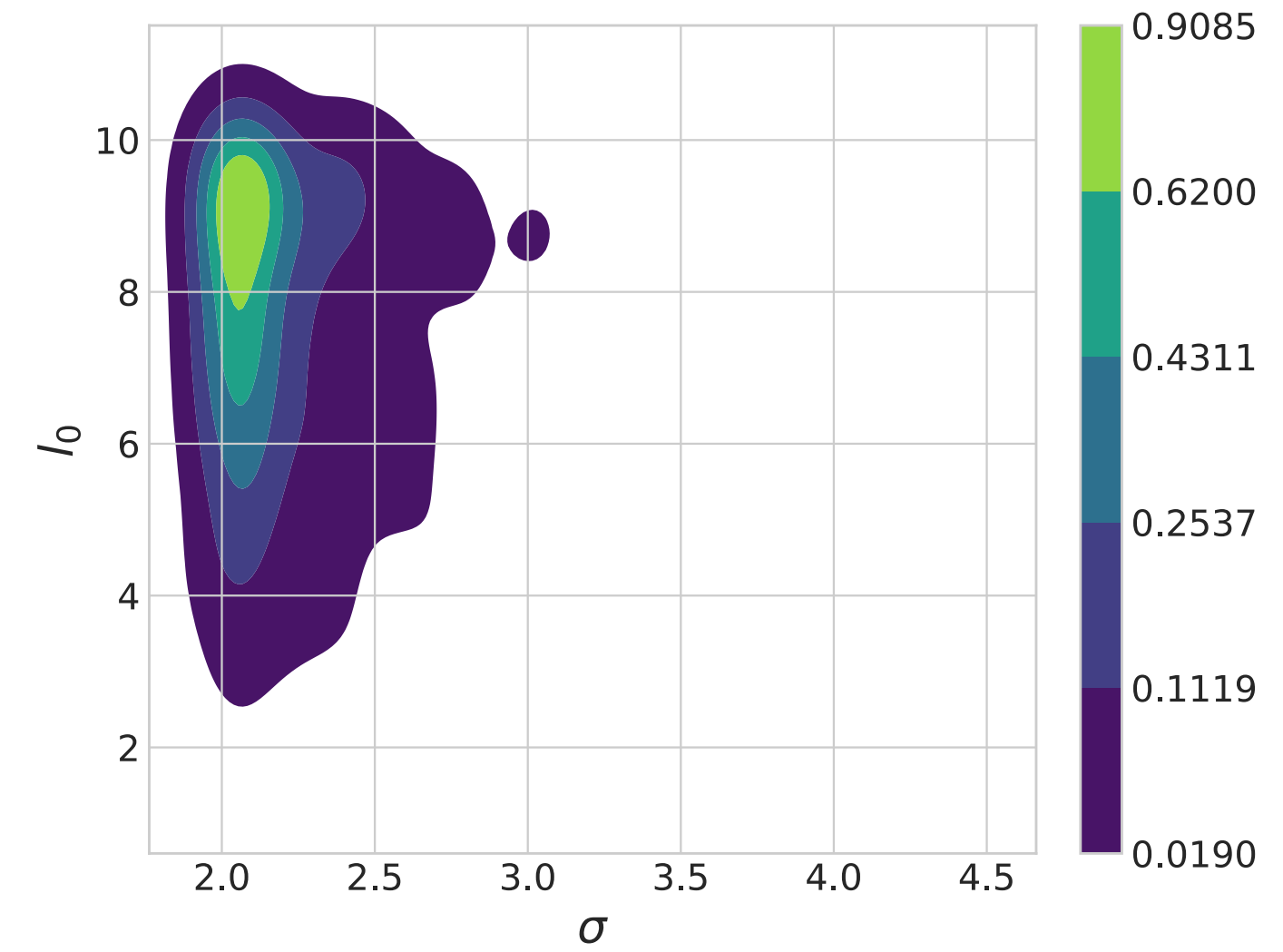
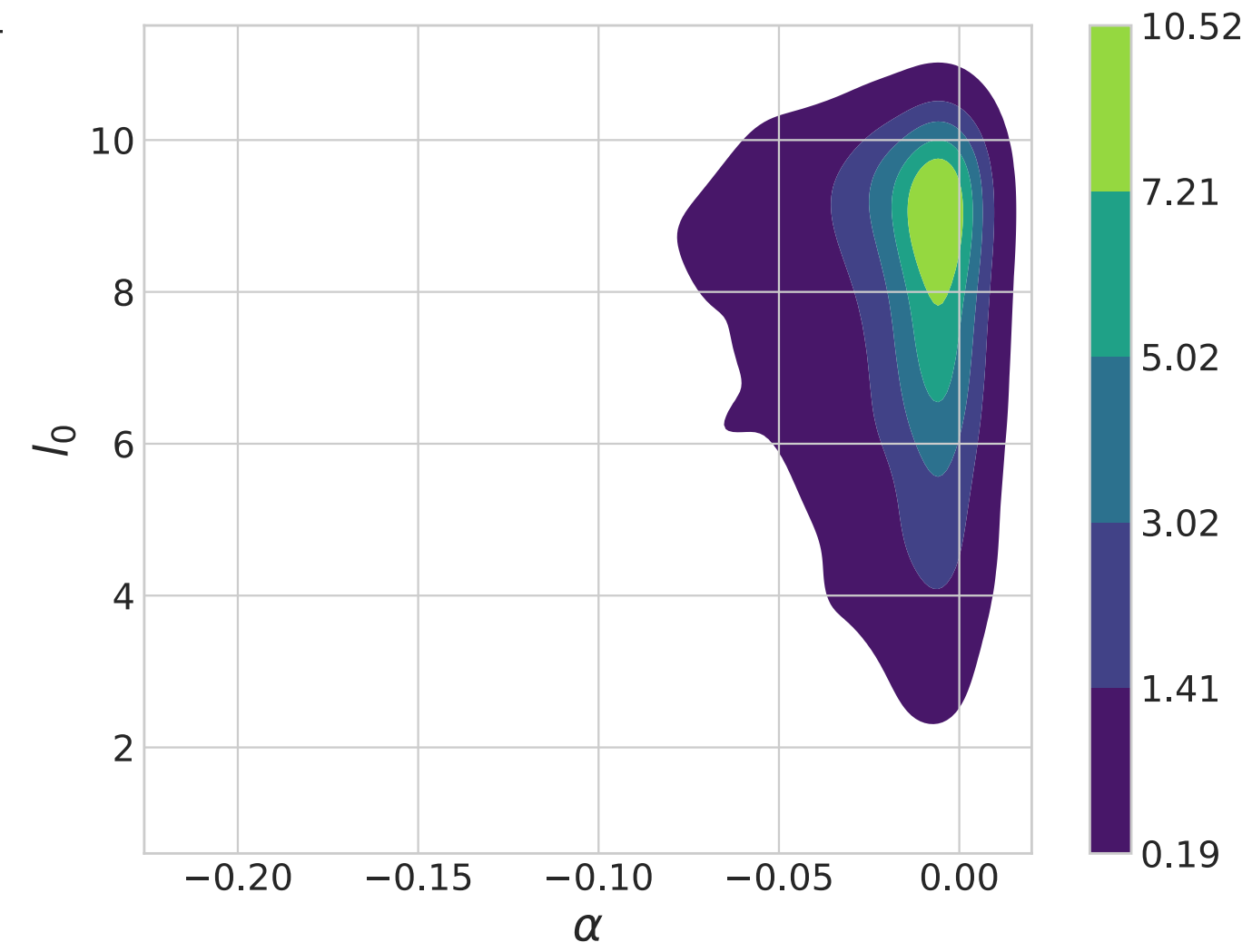
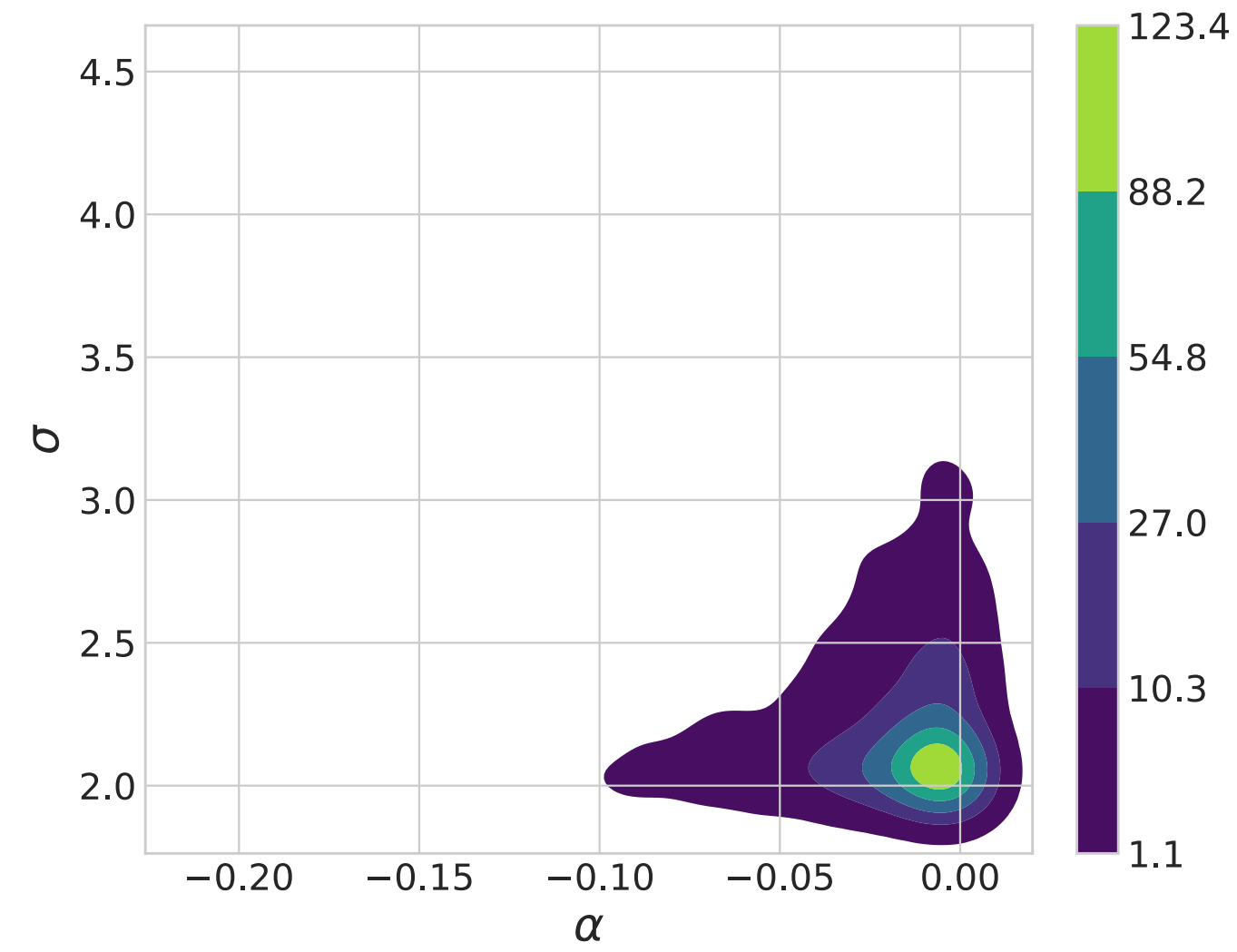


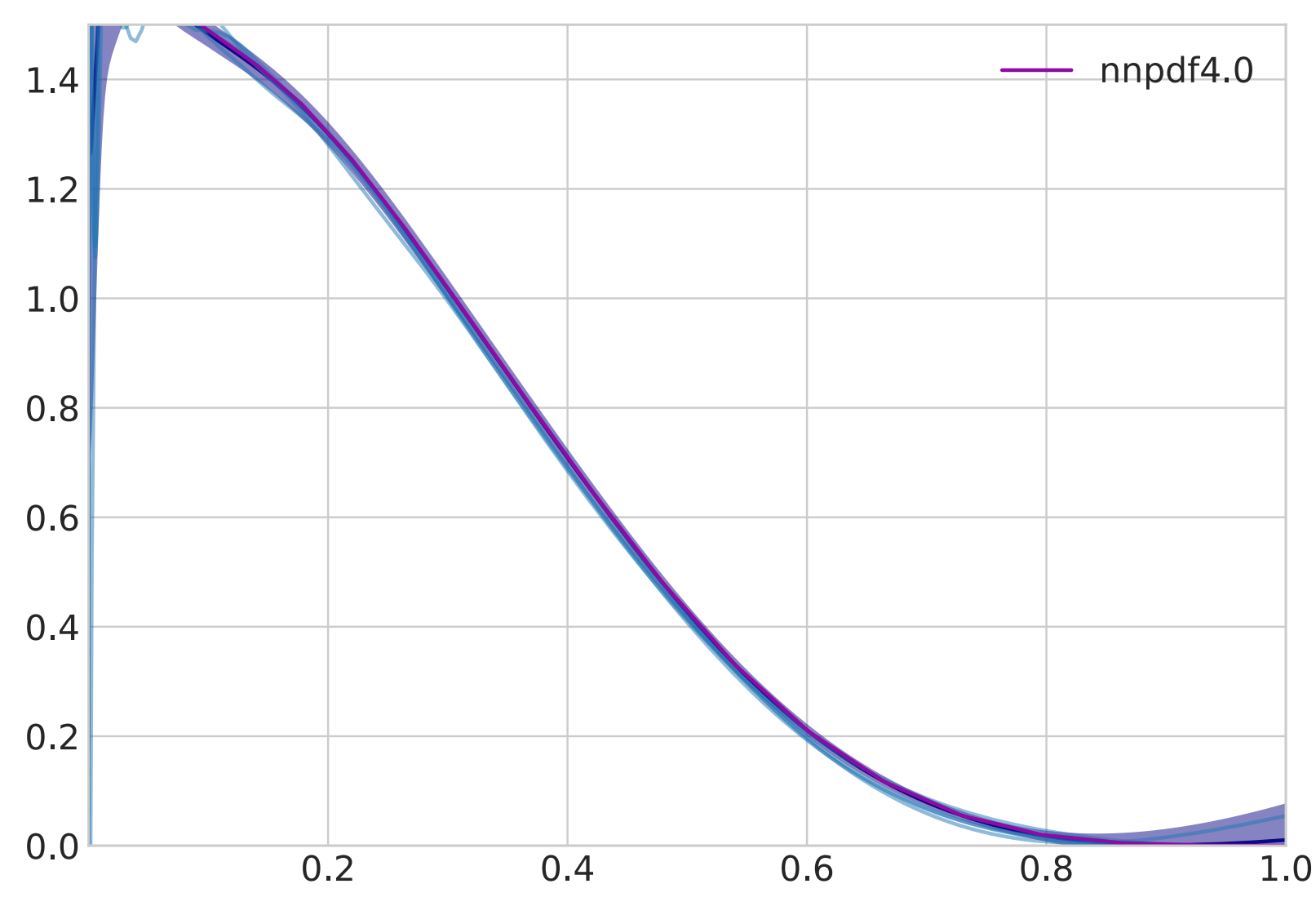
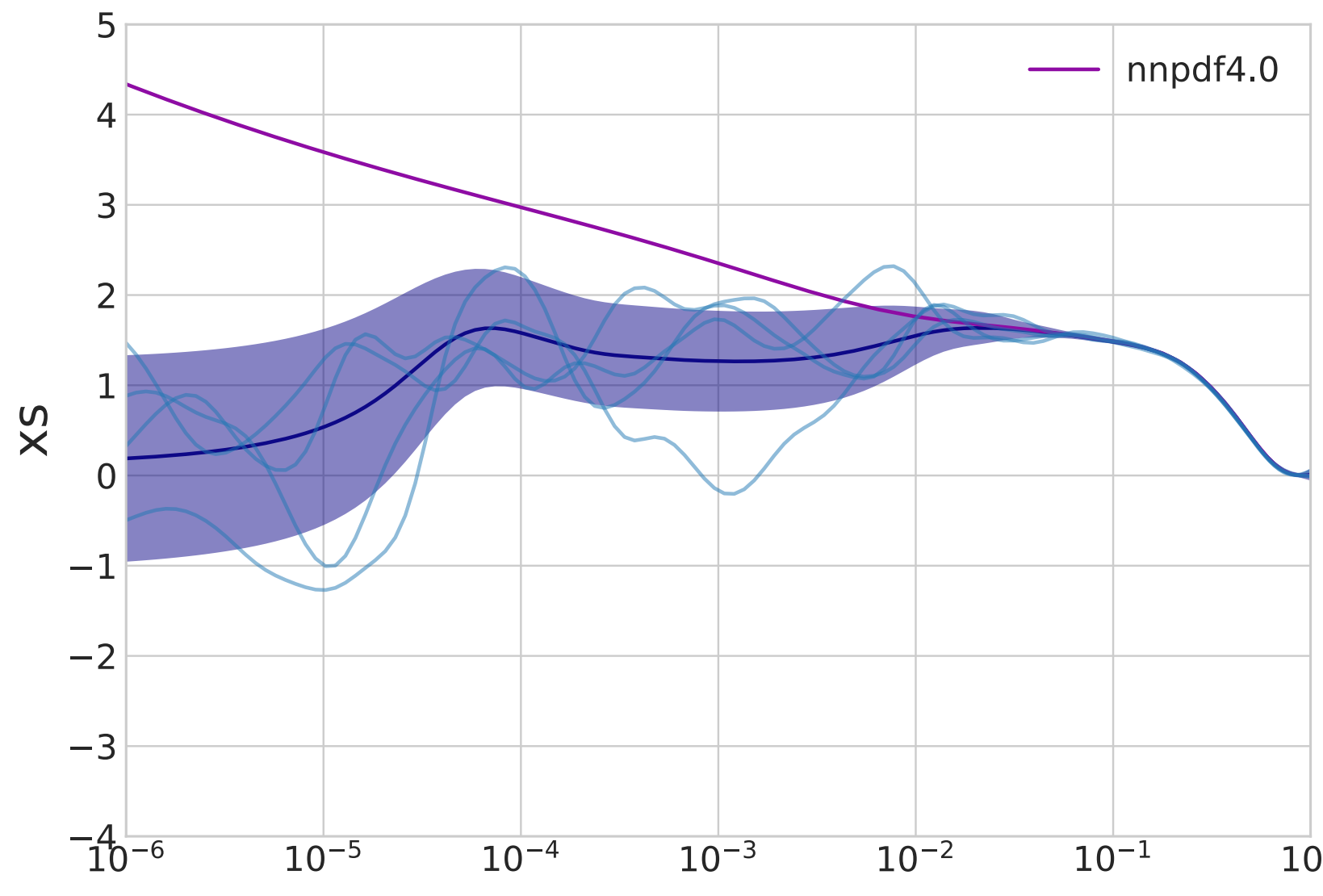
Fixed hyperparameters
 $p(\mathbf{f}^* | \text{data}, \theta)$

Full posterior $p(\mathbf{f}^*, \theta | \text{data})$



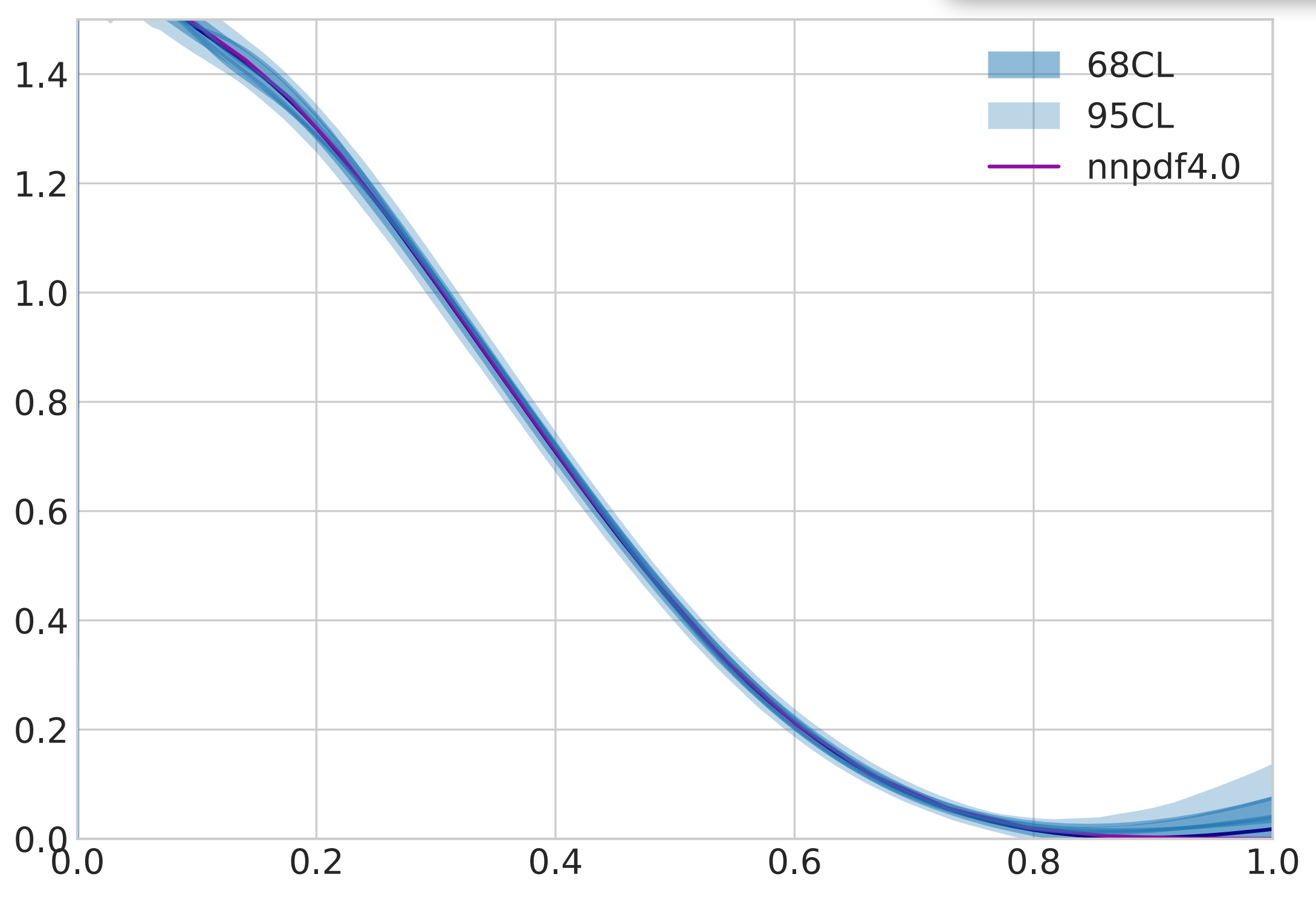
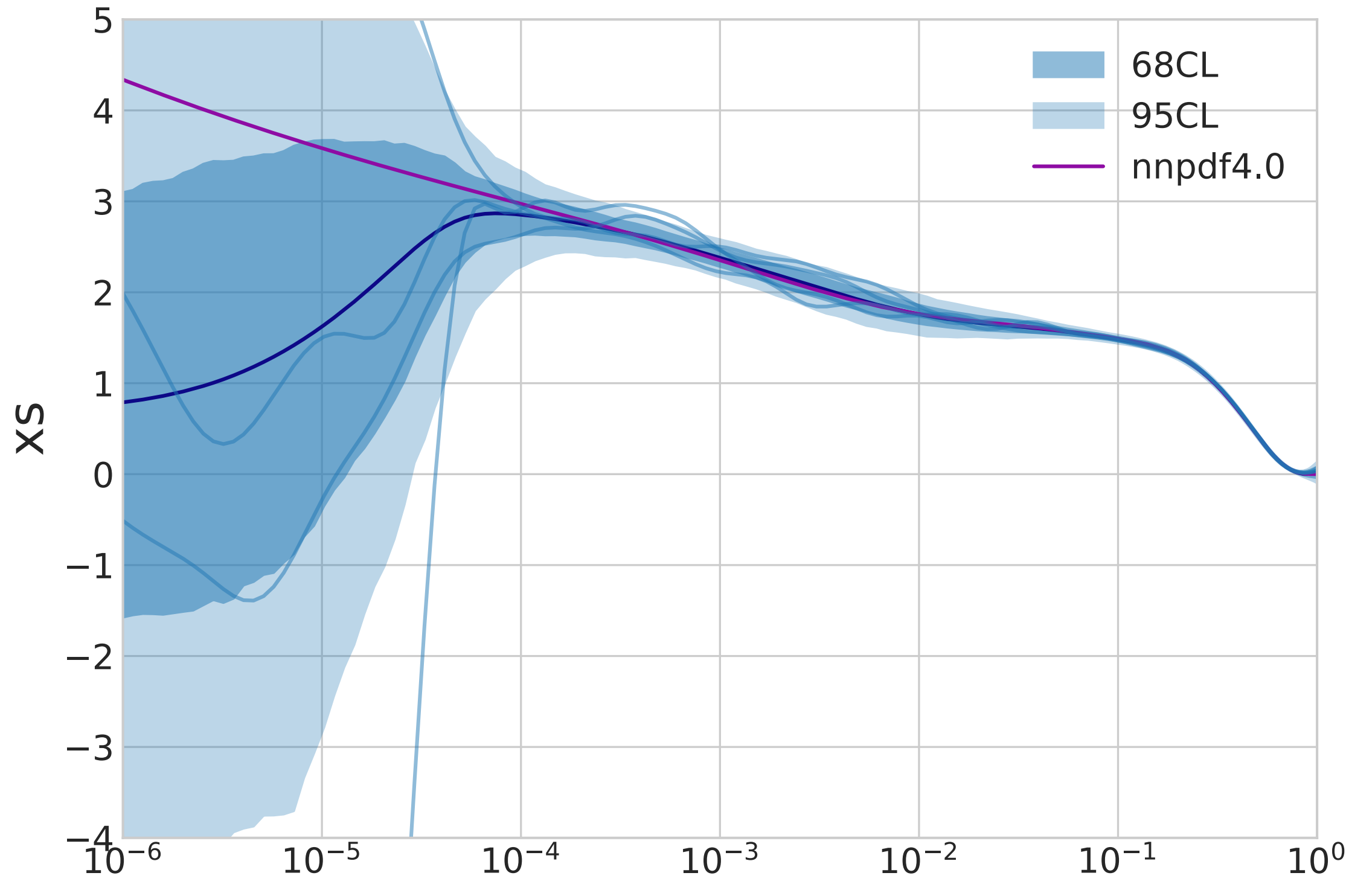
Posterior for hyperparameters (singlet)





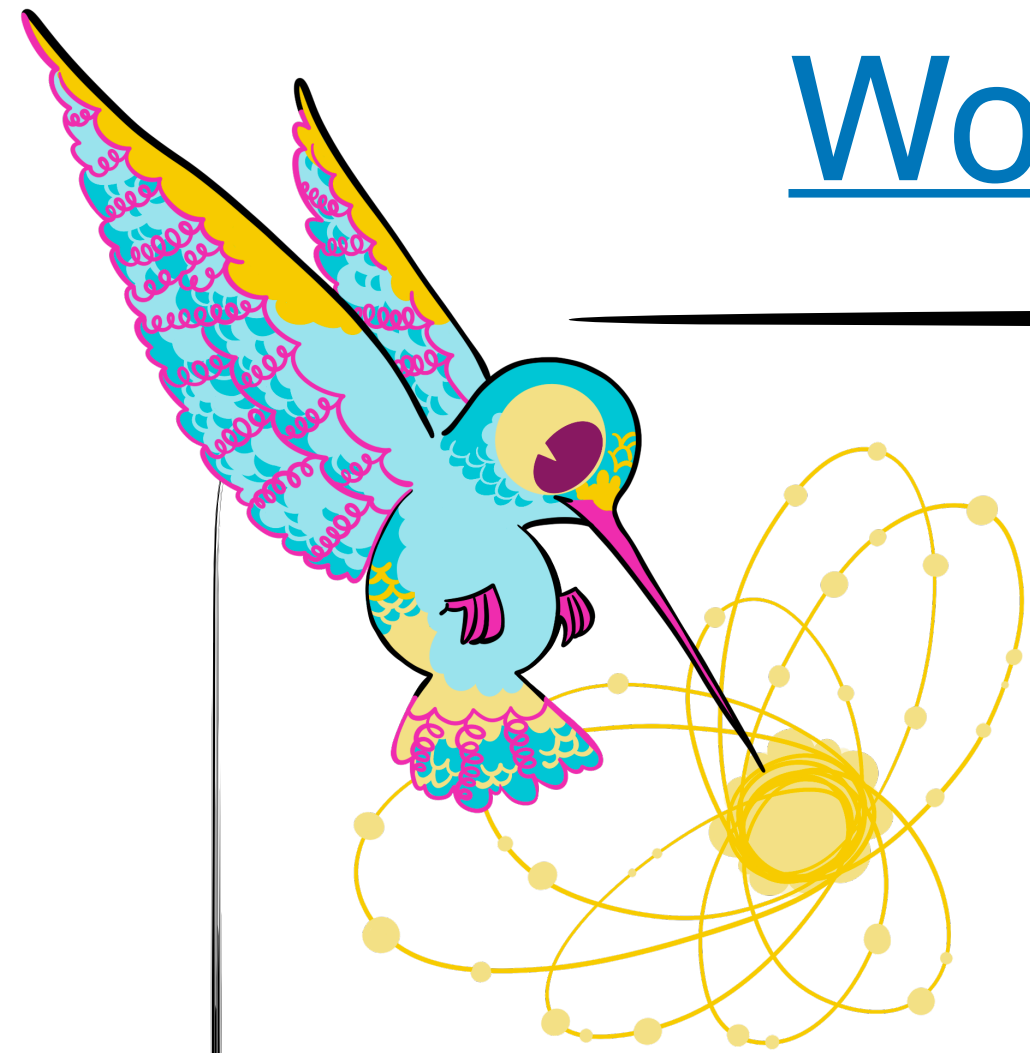
Fixed hyperparameters
 $p(\mathbf{f}^* | \text{data}, \theta)$

Full posterior $p(\mathbf{f}^*, \theta | \text{data})$



Workflow within colibri

WIP with Maria Ubiali, Mark Constantini, Luigi Del debbio, Luca Mantani



Collect data and FK tables



Build the prior as a function of hyperparameters:

- Choose kernel
- Encode theory constraints

Inference on hyperparameters



Inference on parameters

Theory constraints

Kinetic limit

$$f(1) = 0$$

Additional linear constrain on the PDF: can be implemented directly in the FK table

Sum rules

Sum rules can be implemented as additional linear constrains on the primitive of the PDF

$$g \sim GP\left(0, K(x, y)\right) \longrightarrow \frac{dg}{dx} \sim GP\left(0, \partial_x \partial_y K(x, y)\right)$$

Positivity, integrability

Penalty terms in the likelihood

Summary and future work

- Alternative methodology to fit PDFs, orthogonal to the ones currently used
- Well defined uncertainties
- Assumptions clearly defined in the prior
- Analytical understanding of what is going on during NN fits (NN dynamic, see Luigi's talk)

- Systematic study of different possible kernels
- Comparison with non Bayesian methodologies. Are there any differences?
- Implementation of a full global analysis



Artwork by @qftoons

Work in progress within colibri...

Backup slides

Global fits

$$\sigma = \sum_{i,j} \int dx_1 dx_2 f_i(x_1, \mu) f_j(x_2, \mu) \hat{\sigma} \left(x_1, x_2, \frac{Q}{\mu} \right) \times (1 + \mathcal{O}(\Lambda/M)^p)$$

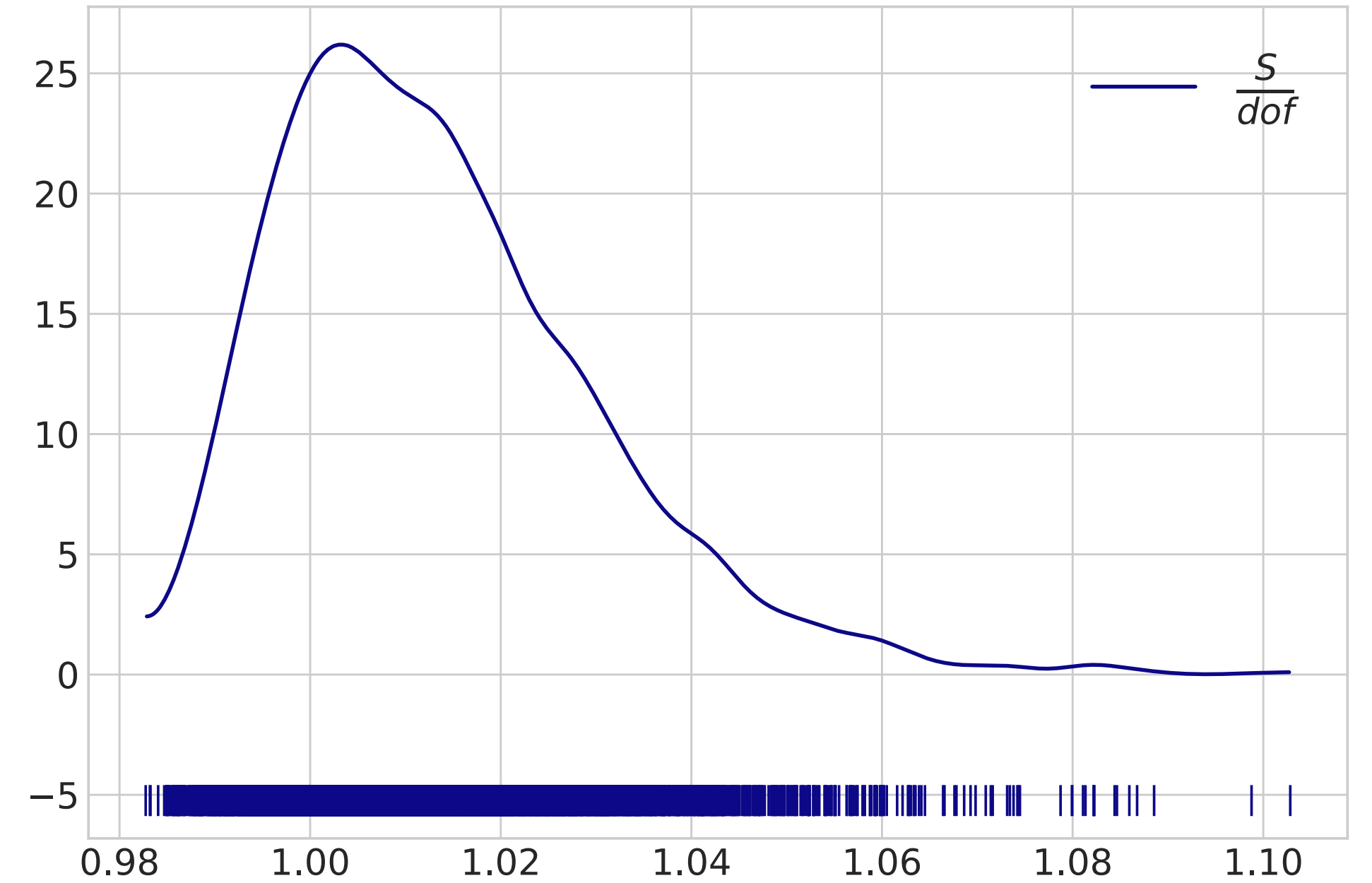
$$p(\mathbf{f}, \theta | \text{data}) = p(\mathbf{f} | \text{data}, \theta) p(\theta | \text{data})$$

This bit is not a gaussian
distribution any longer

To access the posterior we have
to run a MCMC having dimension
 $\dim \mathbf{f} + \dim \theta$

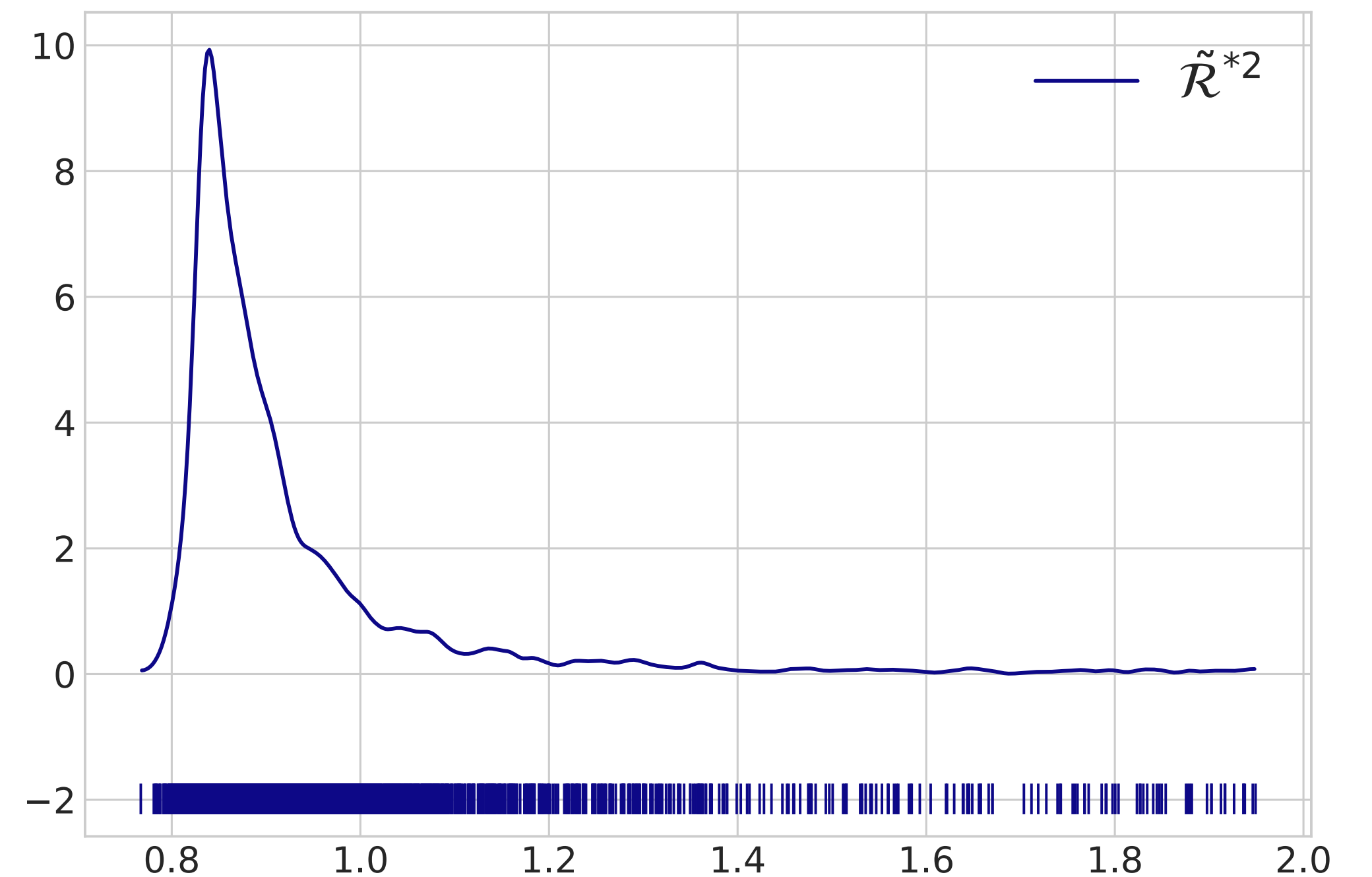
Fit quality

$$\frac{S}{dof} = \frac{1}{N_{\text{data}}} \left((\mathbf{m} - \tilde{\mathbf{m}})^T K_{xx}^{-1} (\mathbf{m} - \tilde{\mathbf{m}}) + (y - FK \tilde{\mathbf{m}})^T C_Y^{-1} (y - FK \tilde{\mathbf{m}}) \right)$$



Generalisation on unseen data

$$\tilde{\mathcal{R}}^{*2} = \frac{1}{\dim(y^* | y)} (FK^* \tilde{\mathbf{m}} - y^*)^T \left(FK^* \tilde{K}_{xx} FK^{*T} + C_Y^* \right)^+ (FK^* \tilde{\mathbf{m}} - y^*)$$



Decomposition of PDF uncertainty

$$\tilde{K} = \underbrace{(I - R_{xx}) K_{xx} (I - R_{xx})^T}_{\text{Methodology}} + \underbrace{a_{xx}^T C_y a_{xx}}_{\text{Experimental error}}$$

$$a_{xx}^T = K_{xx} F K^T \left(F K K_{xx} F K^T + C_y \right)^+$$

$$R_{xx} = a_{xx}^T F K$$

