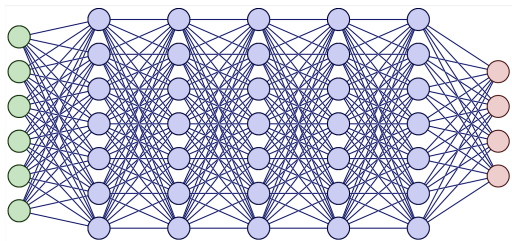


Field Theories and Machine Learning



L Del Debbio

Higgs Centre for Theoretical Physics
University of Edinburgh

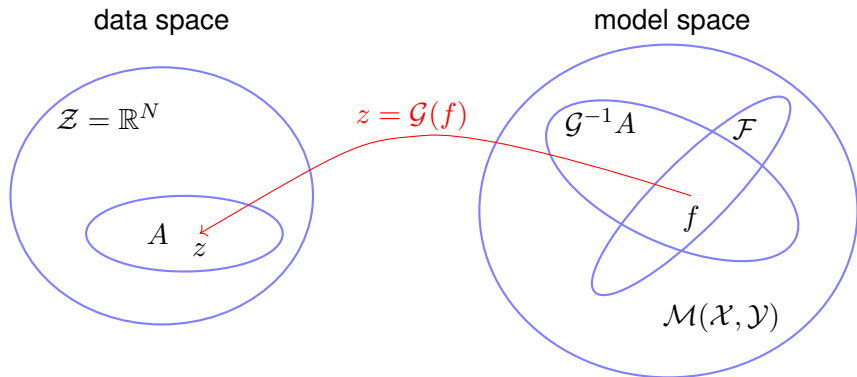
plan of the lectures

- Gaussian Processes and Inverse Problems
- FT for ML
- ML for FT

based on

- Idd et al 21
- Roberts et al 21
- Albantosa et al 23

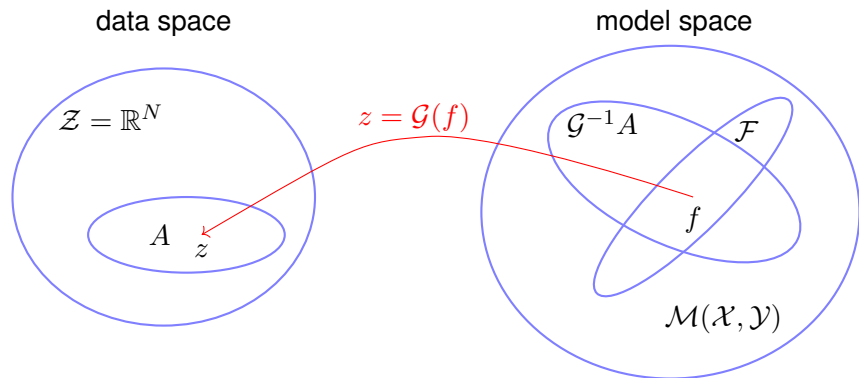
Inverse Problems



Bayes theorem

$$p(f|A) \propto p(A|f)p(f)$$

Inverse Problems



loss function $\mathcal{L} : \mathcal{M}(\mathcal{X}, \mathcal{Y}) \times \mathcal{Z} \rightarrow \mathbb{R}$

algorithm $\mathcal{A} : \mathcal{Z} \rightarrow \mathcal{F} \subset \mathcal{M}(\mathcal{X}, \mathcal{Y})$

a few definitions

empirical risk

$$\hat{\mathcal{R}}_s(f) = \mathcal{L}(f, s)$$

erm algorithm

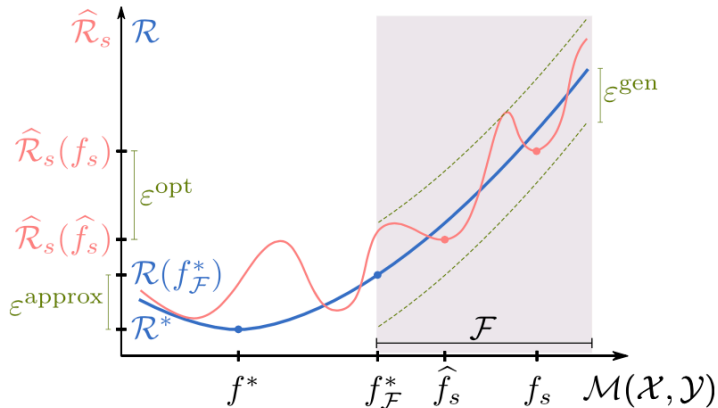
$$\mathcal{A}^{\text{erm}}(s) = \arg \min_{f \in \mathcal{F}} \hat{\mathcal{R}}_s(f) = \hat{f}_s$$

risk

$$\begin{aligned} \mathcal{R}(f) &= \mathbb{E}_Z[\mathcal{L}(f, Z)] = \int dP_Z(z) \mathcal{L}(f, z) \\ &= \frac{1}{N_{\text{rep}}} \sum_{k=1}^{N_{\text{rep}}} \mathcal{L}(f, z^{(k)}) \end{aligned}$$

Bayes optimal function (generalization)

$$\mathcal{R}^* = \mathcal{R}(f^*) = \inf_{f \in \mathcal{M}(\mathcal{X}, \mathcal{Y})} \mathcal{R}(f)$$



robust definitions are important for precision physics!

Linear Problems & Bayesian inference

$$y_I = \int dx C_I(x) f(x), \quad I = 1, \dots, N_{\text{dat}}$$

↔ ill-defined problem (e.g. PDFs from DIS, lectures by Forte)

simple underdetermined problem

model: $u \in \mathbb{R}^{N_{\text{mod}}}$, $N_{\text{mod}} > 1$

theory prediction: $N_{\text{dat}} = 1$, $g \in \mathbb{R}^{N_{\text{mod}}}$, $y = (g^T u)$

prior: $p(u) = \frac{1}{\sqrt{\det(2\pi K)}} \exp\left(-\frac{1}{2}|u - m|_K^2\right)$

$$|u - m|_K^2 = (u - m)^T K^{-1} (u - m)$$

fluctuations of the experimental data: $\eta \sim p_\eta = \mathcal{N}(0, \gamma^2)$

then

$$\begin{aligned} p(y|u) &= \mathbf{Prob}(g^T u + \eta = y) \\ &= \int d\eta p_\eta(\eta) \delta(g^T u + \eta - y) \\ &= p_\eta(y - g^T u) \end{aligned}$$

Bayes theorem:

$$\begin{aligned} p(u|y) &\propto \exp\left(-\frac{1}{2\gamma^2} |y - g^T u|^2 - \frac{1}{2} |u - m|_K^2\right) \\ &\propto \exp\left(-\frac{1}{2} |u - \tilde{m}|_{\tilde{K}}^2\right) \end{aligned}$$

$$\tilde{m} = m + (Kg) \frac{y - g^T m}{\gamma^2 + (g^T Kg)}$$
$$\tilde{K} = K - \frac{(Kg)(Kg)^T}{\gamma^2 + (g^T Kg)}$$

hence:

$$\tilde{m}^* = \lim_{\gamma \rightarrow 0} \tilde{m} = \frac{(Kg)y}{(g^T Kg)}$$
$$\implies g^T \tilde{m}^* = y$$

and

$$\tilde{K}^* = \lim_{\gamma \rightarrow 0} \tilde{K} = K - \frac{(Kg)(Kg)^T}{(g^T Kg)}$$
$$\implies (\tilde{K}^* g) = 0$$

removing the prior

set $K = \sigma \mathbb{1}$ and take the limit $\sigma \rightarrow \infty$

yields

$$\begin{aligned}\tilde{K}^* &= \sigma \left(\mathbb{1} - \frac{gg^T}{(g^T g)} \right) \\ &= \begin{pmatrix} 0 & \\ & \sigma \end{pmatrix}\end{aligned}$$

the error in the g direction vanishes

the error in the orthogonal subspace diverges

Bayesian Inference

$$y_I = \int dx C_I(x) f(x)$$

Bayesian solution of the form:

$$f_i = f(x_i) \quad \longrightarrow \quad p(f_1 \dots f_N | y, \mathcal{H}) = \frac{p(y|f, \mathcal{H})p(f|\mathcal{H})}{p(y|\mathcal{H})}$$

\hookrightarrow dependence on the choice of the prior

Gaussian processes: $f \sim \mathcal{GP}(m, k)$

$$\mathbb{E}[f_i] = m(x_i), \quad \text{Cov}[f_i, f_j] = k(x_i, x_j)$$

prior : $\mathcal{GP}(m, k) \quad \Longrightarrow \quad$ posterior : $\mathcal{GP}(\tilde{m}, \tilde{k})$

Gaussian Processes in Action

$$\mathbf{f} \in \mathbb{R}^N, \quad \mathbf{f}^* \in \mathbb{R}^M,$$

prior distribution

$$\frac{1}{\sqrt{\det(2\pi K)}} \exp \left\{ -\frac{1}{2} \left((\mathbf{f} - \mathbf{m})^T, (\mathbf{f}^* - \mathbf{m}^*)^T \right) K^{-1} \begin{pmatrix} \mathbf{f} - \mathbf{m} \\ \mathbf{f}^* - \mathbf{m}^* \end{pmatrix} \right\}$$

where K is an $(N + M) \times (N + M)$ matrix

$$K = \begin{pmatrix} k(\mathbf{x}, \mathbf{x}^T) & k(\mathbf{x}, \mathbf{x}^{*T}) \\ k(\mathbf{x}^*, \mathbf{x}^T) & k(\mathbf{x}^*, \mathbf{x}^{*T}) \end{pmatrix} = \begin{pmatrix} K_{\mathbf{xx}} & K_{\mathbf{xx}^*} \\ K_{\mathbf{x}^*\mathbf{x}} & K_{\mathbf{x}^*\mathbf{x}^*} \end{pmatrix}.$$

$$T_I = \sum_{i=1}^N (\text{FK})_{Ii} f_i$$

where $(\text{FK})_{Ii} = b_I(x_i)$.

$$E[T_I] = (\text{FK})_{Ij} m_j ,$$

$$\text{Cov}[T_I, T_J] = (\text{FK})_{Ii} (K_{\mathbf{xx}})_{ij} (\text{FK})_{jJ}^T .$$

experimental error on data

$$\epsilon \sim \mathcal{N}(0, C_Y) ,$$

we have three stochastic variables $\mathbf{f}, \mathbf{f}^*, \epsilon$ with

$$\text{Cov} = \begin{pmatrix} K & 0 \\ 0 & C_Y \end{pmatrix} ,$$

and we will impose that

$$y = (\text{FK})\mathbf{f} + \epsilon ,$$

Posterior Distribution

$$(\mathbf{f} | (\text{FK})\mathbf{f} + \epsilon = \mathbf{y}) \sim \mathcal{N}(\tilde{\mathbf{m}}, \tilde{K}_{\mathbf{xx}}),$$

where

$$\begin{aligned}\tilde{\mathbf{m}} &= \mathbf{m} + K_{\mathbf{xx}}(\text{FK})^T C_{YT}^+ (\mathbf{y} - (\text{FK})\mathbf{m}), \\ \tilde{K}_{\mathbf{xx}} &= K_{\mathbf{xx}} - K_{\mathbf{xx}}(\text{FK})^T C_{YT}^+ (\text{FK})K_{\mathbf{xx}},\end{aligned}$$

$$C_{YT} = (\text{FK})K_{\mathbf{xx}}(\text{FK})^T + C_Y.$$

exercise:

$$\tilde{K}_{\mathbf{xx}}^{-1} = K_{\mathbf{xx}}^{-1} + (\text{FK})^T C_Y^{-1} (\text{FK}) \quad \rightarrow \quad (\text{FK})^T C_Y^{-1} (\text{FK})$$

Boostrapping a fit

fluctuations of the data encoded in an ensemble of "replicas"

$$y^{(k)} = y + \epsilon^{(k)}, \quad k = 1, \dots, N_{\text{rep}}$$

for each replica minimize

$$\begin{aligned}\chi^{2,(k)} &= \frac{1}{2} \left(y^{(k)} - (\text{FK})\mathbf{f} \right)^T C_Y^{-1} \left(y^{(k)} - (\text{FK})\mathbf{f} \right) \\ f^{(k)} &= \left((\text{FK})^T C_Y^{-1} (\text{FK}) \right)^{-1} (\text{FK})^T C_Y^{-1} y^{(k)}\end{aligned}$$

$f^{(k)}$ are normally distributed with mean $\tilde{\mathbf{m}}$ and covariance $\tilde{K}_{\mathbf{xx}}$

for $\mathbf{m} = 0$, $K_{\mathbf{xx}}^{-1} = 0$

More on the Posterior

for the *unseen* points

$$(\mathbf{f}^* | (\mathbf{F}\mathbf{K})\mathbf{f} + \epsilon = y) \sim \mathcal{N}(\tilde{\mathbf{m}}^*, \tilde{K}_{\mathbf{x}^*\mathbf{x}^*}) ,$$

where in this case

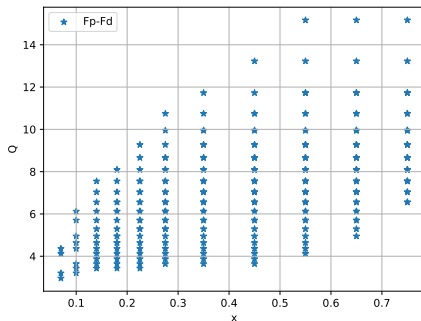
$$\begin{aligned}\tilde{\mathbf{m}}^* &= \mathbf{m}^* + K_{\mathbf{x}^*\mathbf{x}}(\mathbf{F}\mathbf{K})^T C_{YT}^+ (\mathbf{y} - (\mathbf{F}\mathbf{K})\mathbf{m}) , \\ \tilde{K}_{\mathbf{x}^*\mathbf{x}^*} &= K_{\mathbf{x}^*\mathbf{x}^*} - K_{\mathbf{x}^*\mathbf{x}}(\mathbf{F}\mathbf{K})^T C_{YT}^+ (\mathbf{F}\mathbf{K})K_{\mathbf{x}\mathbf{x}^*} .\end{aligned}$$

Example

$$F_2^p - F_2^d = C_{T_3} \otimes T_3$$

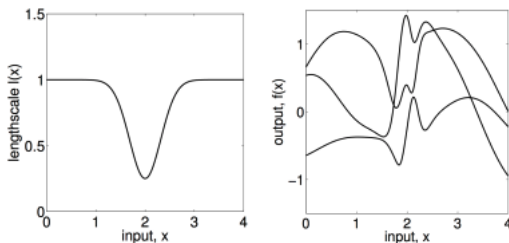
data covariance

$$C_Y = \text{Cov} [F_2^p, F_2^p] + \text{Cov} [F_2^d, F_2^d] - 2\text{Cov} [F_2^p, F_2^d]$$



Gibbs Kernel

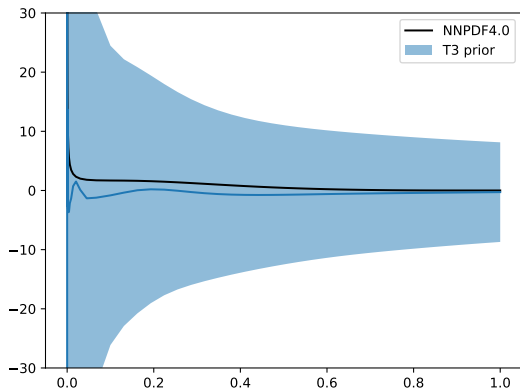
$$k(x, y) = \sigma^2 \sqrt{\frac{2l(x)l(y)}{l^2(x) + l^2(y)}} \exp \left[-\frac{(x - y)^2}{l^2(x) + l^2(y)} \right]$$



here we use

$$l(x) = l_0 \times (x + \epsilon) .$$

Prior Distribution



Closure Test

generate artificial data: $y = (\text{FK})\mathbf{f}_0$

introduce a *filter* function

$$R_{\mathbf{xx}} = K_{\mathbf{xx}}(\text{FK})^T C_{YT}^{-1}(\text{FK})$$
$$R_{\mathbf{x}^*\mathbf{x}} = K_{\mathbf{x}^*\mathbf{x}}(\text{FK})^T C_{YT}^{-1}(\text{FK})$$

then

$$\tilde{\mathbf{m}} - \mathbf{m} = R_{\mathbf{xx}}(\mathbf{f}_0 - \mathbf{m})$$
$$\tilde{\mathbf{m}}^* - \mathbf{m}^* = R_{\mathbf{x}^*\mathbf{x}}(\mathbf{f}_0 - \mathbf{m}) \implies \tilde{\mathbf{m}}^* = R_{\mathbf{x}^*\mathbf{x}}\mathbf{f}_0$$

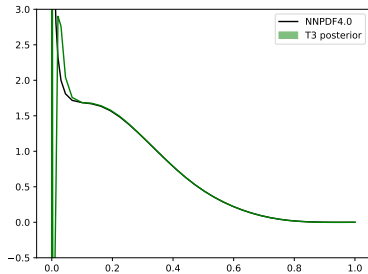
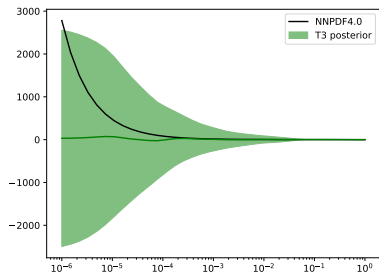
considering the case $\mathbf{x}^* = \mathbf{x}$ we define the bias

$$\tilde{\mathbf{m}} - \mathbf{f}_0 = [R_{\mathbf{xx}} - \mathbb{1}]\mathbf{f}_0,$$

and the covariance of the posterior can be written as

$$\tilde{K} = (\mathbb{1} - R_{\mathbf{xx}})K_{\mathbf{xx}}(\mathbb{1} - R_{\mathbf{xx}})^T + a_{\mathbf{xx}}^T C_Y a_{\mathbf{xx}}.$$

Posterior Distribution



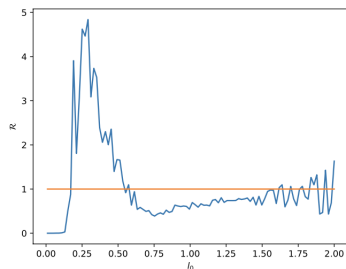
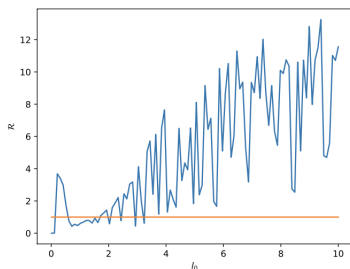
Bias & Variance

define

$$\mathcal{B} = (\text{FK})(\tilde{\mathbf{m}} - \mathbf{f}_0) = T - y$$

$$\mathcal{V} = (\text{FK})\tilde{K}(\text{FK})^T$$

$$\mathcal{R} = \sum_{I,J} \mathcal{B}_I (\mathcal{V}^{-1})_{IJ} \mathcal{B}_J$$



Outlook

- GP are an effective way to explore the space of functions
- prior is explicitly specified AND determines the posterior
- careful study of the bias/variance ratio in order to choose the prior