

GPFADs Problems: generalize beyond latent planes

# Bayesian Inference & Predictions in finite regression Models

Data:  $\underline{x}, \underline{y}$

Model M:  $y = f_w(x) + \epsilon$

Gaussian likelihood:

$$p(\underline{y} | \underline{x}, \underline{w}, M) \propto \prod_{n=1}^N \exp\left(-\frac{1}{2} (y_n - f_w(x_n))^2 / \sigma_{\text{noise}}^2\right)$$

Maximum Likelihood:

$$\underline{w}_{\text{ML}} = \arg \max_w p(\underline{y} | \underline{x}, \underline{w}, M)$$

data, model & weights  $\& x$

new data / new  $x$ , trained weights  $\underline{w}_{\text{ML}}$  & model)

Predictions:  
(Max. likelihood)

$$p(y_* | \underline{x}_*, \underline{w}_{\text{ML}}, M)$$

Posterior

(Bayes rule for 3+ variables)

$$p(\underline{w} | \underline{x}, \underline{y}, M) = \frac{p(\underline{w} | M) p(\underline{y} | \underline{x}, \underline{w}, M)}{p(\underline{y} | \underline{x}, M)}$$

likelihood

(Marginalizing out the parameters)

Predictions:

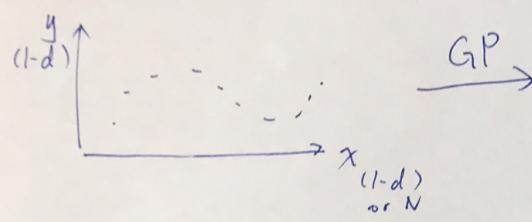
$$\begin{aligned} p(y_* | \underline{x}_*, \underline{x}, \underline{y}, M) &= \int p(y_*, \underline{w} | \underline{x}, \underline{y}, \underline{x}_*, M) d\underline{w} \\ &= \int \underbrace{p(y_* | \underline{w}, \underline{x}_*, M)}_{\text{likelihood}} \underbrace{\frac{p(\underline{w} | \underline{x}, \underline{y}, M)}{p(\underline{y} | \underline{x}, M)}}_{\text{posterior}} d\underline{w} \end{aligned}$$

Marginal Likelihood:

$$p(\underline{y} | \underline{x}, M) = \int \underbrace{p(\underline{w} | \underline{x}, M)}_{\text{prior}} \underbrace{\frac{p(\underline{y} | \underline{x}, \underline{w}, M)}{p(\underline{y} | \underline{x}, M)}}_{\text{likelihood}} d\underline{w}$$

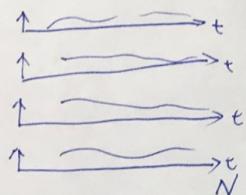
Howard

Coursework



Tensorflow  $\rightarrow$  GPflow (GP Intro)  
Multi-O/P GP  $\rightarrow$  Covariance fn.  
design

Project ① ② Non-rev GP Regression



② GPFADS (GP Factor Analysis with DS)

① Plot kernel  $\checkmark \checkmark \checkmark$   
 $\uparrow$   
build kernel  $\Rightarrow$

$$\left[ \begin{array}{c} t_1 \\ [f(z)] \\ -H[f](z) \end{array} \right] \left[ \begin{array}{c} H[f](z) \\ f(z) \end{array} \right]$$

② Plot trajectories from kernel

$$\frac{\tau}{\tau_2} = \left( \frac{1 + \tau^2}{\tau} \right)$$

$$K\left(\left(\frac{t_1}{t_2}\right)^l, \left(\frac{t_1}{t_2}\right)^r\right)$$

$$= \begin{bmatrix} f^{(+)}(t_1^l - t_2^r) & \alpha f^{(-)}(t_1^l - t_2^r) \\ -\alpha f^{(-)}(t_2^l - t_1^r) & f^{(+)}(t_2^l - t_2^r) \end{bmatrix}$$

Hilbert transform

Scalar cov function (RQ)

should I add "-" here? K not symmetric?

$$K_{xx} \in \mathbb{R}^{MT \times MT}$$

$$\tilde{F}_{\text{circulant}} = \begin{bmatrix} f_0 & & & & & & \\ f_1 & & & & & & \\ \vdots & & & & & & \\ f_{T-2} & & & & & & \\ f_{T-1} & & & & & & \\ f_{T-2} & & & & & & \\ \vdots & & & & & & \\ f_1 & & & & & & \\ \uparrow & & & & & & \end{bmatrix} \triangleq \begin{bmatrix} F & S \\ S^T & F' \end{bmatrix}$$

contains all the info

so for each dimension we have a  $\begin{bmatrix} \bar{v}_m \\ 0 \end{bmatrix}$  so that  $\begin{bmatrix} F & S \\ S^T & F' \end{bmatrix} \begin{bmatrix} \bar{v}_m \\ 0 \end{bmatrix} = \begin{bmatrix} F\bar{v}_m \\ S^T\bar{v}_m \end{bmatrix}$

$= \begin{bmatrix} DFT^{-1}[DFT(F) \odot DFT(\bar{v}_m)] \\ \dots \end{bmatrix}$

0	0
1	1
2	2
3	3
4	4
5	5
6	6
7	7
8	8
9	9

$$\text{ell-half} = \frac{r^2}{d_{\text{ts}} \cdot (t_5 - t_1)^2}$$

$\downarrow$

nu: mean present parameter  $\uparrow$

$\boxed{\square} \quad \boxed{\square} \quad \boxed{\square}$  K-half: complete  $\boxed{\square}$  from  $\boxed{\square}$  of  $\boxed{\square} \quad \boxed{\square}$

K-half  $\xrightarrow[2 \times 180]{}$  Args:

$$\left( \begin{array}{c|c|c} \hline & 180 & 180 \\ \hline 1 & \boxed{1} & \boxed{1} \\ \hline 2 & \boxed{2} & \boxed{2} \\ \hline \end{array}, \quad \begin{array}{c|c|c} \hline & 180 & 180 \\ \hline v_1 & \boxed{v_1} & \boxed{v_1} \\ \hline v_2 & \boxed{v_2} & \boxed{v_2} \\ \hline \end{array} \right) \quad (k-1) \cdot 256 + 1, \quad (k-1) \cdot 256 + 256$$

~~$k=255, k$~~   
 $256k-255, 256k+1$

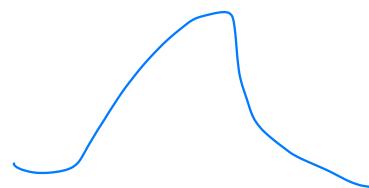
$\boxed{1} \quad \boxed{v_1} \quad 180 \times 1$

$\boxed{2} \quad \boxed{v_2} \quad 180 \times 1,$

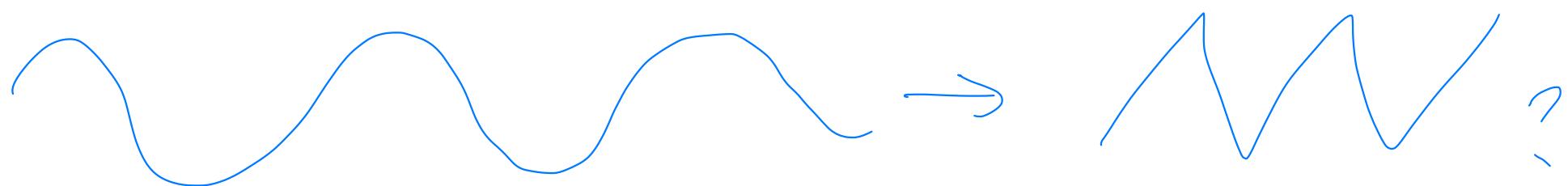
Continuous multidimensional data  $\Rightarrow$  discrete spike trains

what about reversing the process?

Spectral mixture



- time scales  
about neuron



objective?

Why do we need non-reversibility?

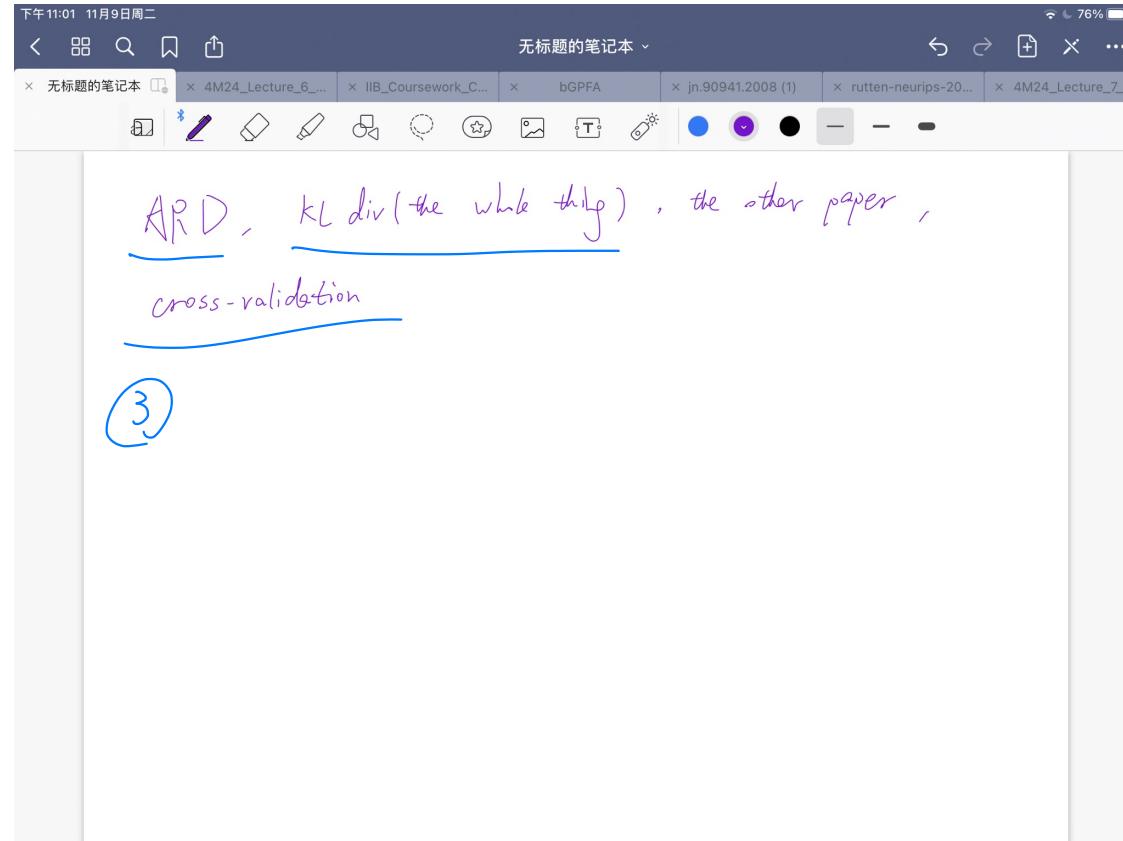
apply to neural data!

20/10/2021

- ① get to familiarize GP quickly !!!
- ② temporal reversibility !!
- ③ pendulum dynamics !
- ④ 2nd order non-reversibility ?

Nov 9, 2021

① ELBO



③

b GPFA:

$$p(Y|t) = \int p(Y|F) p(F|x) p(x|t) dF dx$$

neuron activity      neuron-specific variables  
(firing rates, etc)

latent variables  
time

$$p(x|Y, t) = \frac{p(x|t)p(Y|t, x)}{p(Y|t)}$$

+ F                  + F  
marginal likelihood

Bayesian GPFA is defined by the following generative model:

$$\mathbf{Y} \in \mathbb{R}^{N \times T} \text{ (data)}$$

$$\mathbf{X} \in \mathbb{R}^{D \times T} \text{ (latent variables)}$$

$$\mathbf{C} \in \mathbb{R}^{N \times D} \text{ (readout matrix)}$$

$$\mathbf{x}_d \sim \mathcal{GP}(0, K_{\text{RBF}})$$

$$c_{nd} \sim \mathcal{N}(c_{nd}; \mu = 0, \sigma^2 = s_d^2)$$

$$y_{nt|\mathbf{C}, \mathbf{X}} = p(y_{nt} | (\mathbf{C}\mathbf{X})_{nt})$$

We learn the scales  $\{s_d\}$  by maximizing a lower bound (ELBO) on the model log marginal likelihood:

$$\mathcal{L} \leq \log p(\mathbf{Y}) = \log \int p(\mathbf{Y}|\mathbf{C}, \mathbf{X})p(\mathbf{C})p(\mathbf{X}) d\mathbf{C} d\mathbf{X}$$

This is achieved using variational inference which also provides an estimate of the posterior distribution over our latent variables that we can use for further analyses:  $q(\mathbf{X}) \approx p(\mathbf{X}|\mathbf{Y})$ .

See Jensen & Kao et al. (2021) for further details about the generative model and inference procedure.

We start by downloading an example dataset which was originally recorded by O'Doherty et al. (2018). Here we consider a single recording session where we have binned the data in 25 ms bins in advance. We have put this data on google drive for ease of access in this tutorial; note that the original dataset is available from <https://zenodo.org/record/3854034#.YNCEy5P0nUI>.

$$p(\mathbf{Y}|\mathbf{F}) = \prod_{n,t} \mathcal{N}(y_{nt}; f_{nt}, \sigma_n^2) \quad (2)$$

$$p(\mathbf{F}|\mathbf{X}) = \delta(\mathbf{F} - \mathbf{C}\mathbf{X}) \quad (3)$$

$$p(\mathbf{X}|\mathbf{t}) = \prod_d \mathcal{N}(\mathbf{x}_d; \mathbf{0}, \mathbf{K}_d) \quad \text{with } \mathbf{K}_d = k_d(\mathbf{t}, \mathbf{t}) \quad (4)$$

That is, the prior over the  $d^{\text{th}}$  latent function  $x_d(t)$  is a Gaussian process (Rasmussen and Williams, 1996) with covariance function  $k_d(\cdot, \cdot)$  (usually a radial basis function), and the observation model  $p(\mathbf{Y}|\mathbf{X})$  is given by a parametric linear transformation with independent Gaussian noise.

In this work, we additionally introduce a prior distribution over the mixing matrix  $\mathbf{C} \in \mathbb{R}^{N \times D}$  with hyperparameters specific to each latent dimension. This allows us to *learn* an appropriate latent dimensionality for a given dataset using automatic relevance determination (ARD) similar to previous work in Bayesian PCA (Appendix H; Bishop, 1999) rather than relying on cross-validation or ad-hoc thresholds of variance explained. Unlike in standard GPFA, the log marginal likelihood (Equation 1) becomes intractable with this prior. We therefore develop a novel variational inference strategy (Wainwright and Jordan, 2008) which also (i) provides a scalable implementation appropriate for long continuous neural recordings, and (ii) extends the model to general non-Gaussian likelihoods better suited for discrete spike counts.

In this new framework, which we call Bayesian GPFA (bGPFA), we use a Gaussian prior over  $\mathbf{C}$  of the form  $c_{nd} \sim \mathcal{N}(0, s_d^2)$ , where  $s_d$  is a scale parameter associated with latent dimension  $d$ . Integrating  $\mathbf{C}$  out in Equation 3 then yields the following observation model:

$$p(\mathbf{F}|\mathbf{X}) = \prod_n \mathcal{N}(\mathbf{f}_n; \mathbf{0}, \mathbf{X}^T \mathbf{S}^2 \mathbf{X}), \quad \text{with } \mathbf{S} = \text{diag}(s_1, \dots, s_D). \quad (5)$$

Moreover, we use a general noise model  $p(\mathbf{Y}|\mathbf{F}) = \prod_{n,t} p(y_{nt}|f_{nt})$  where  $p(y_{nt}|f_{nt})$  is any distribution for which we can evaluate its density.

## 2.2 Variational inference and learning

To train the model and infer both  $\mathbf{X}$  and  $\mathbf{F}$  from the data  $\mathbf{Y}$ , we use a nested variational approach. It is

① GP Regression / GPEAD<sup>(1)</sup><sub>(2)</sub> ② (bGPFA) SVGP (GPflow) X

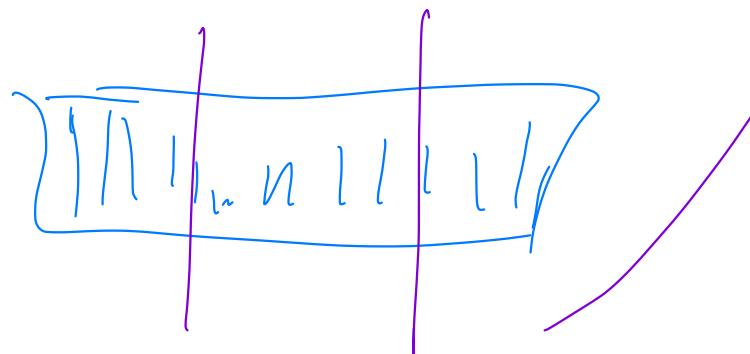
also / figure 2 GP bad or not?

has notebooks

multi-obj GP: GPflow? SVGP needed?

③ MIT - bGPFA → kinematic decoding meaning?

④ real-world application, GPFA: expensive for  $\geq$  a few 100 time bins



problems: ② maybe each chunk's represent different neural activities, experiment activities, cause difference in parameters?

why way to implement Kernel?

① How to train all these pieces of samples? Altogether? Separately?

Presentation :

1.5 Intro → Research's background, purpose of neural decoding,  
Contents in this presentation

GPFADs → GPTA, irrer-kernel: what is it, how to construct,  
example, comparison, higher-dimension (review, show  
understanding of the paper) | 2

Implementation → kernel plot, regression, GPFADs (all planar)

| Application to bGPTA → bGPTA intro (spatial, temporal?)

1.5 | Next plan, → bGPTA → bGPFADs ; other way to express irreversibility  
(GPFADs generate trajectories which )

difficulties... — coding!! Maths!!

Cross-over,  
What's the value of  
irrev index when don't over