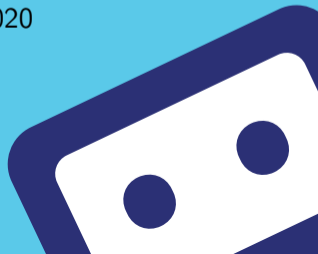


# Sparse Gaussian Processes with Spherical Harmonic Features

Vincent Dutordoir<sup>1</sup>, Nicolas Durrande<sup>1</sup> and James Hensman<sup>2</sup>

<sup>1</sup> PROWLER.io, <sup>2</sup>Amazon (Work completed while JH was at PROWLER.io)

International Conference of Machine Learning – 2020



# Contribution

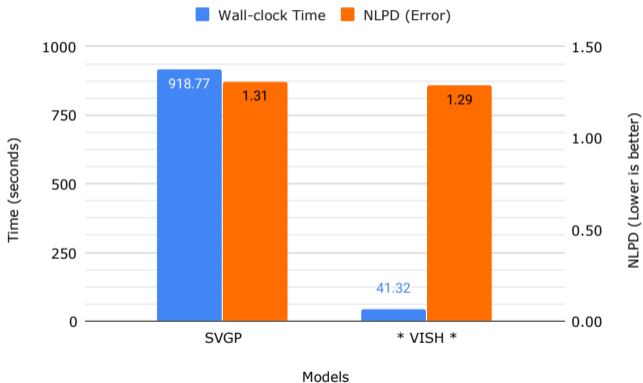
We improve the scaling of Sparse GPs with #datapoints and #inputs

## Airline dataset:

- Regression problem
- $6 \cdot 10^6$  datapoints
- 8 input dimensions

## Setup

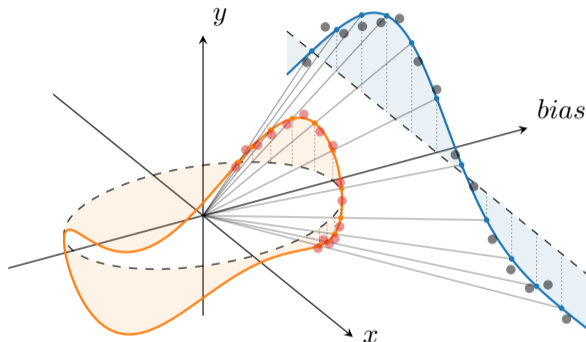
- GTX 1070 GPU



# Variational Inference with Spherical Harmonics (VISH)

## Gist of method:

- make inputs  $d + 1$  dimensional
- project data radially on  $\mathbb{S}^d$
- Fast SVGP on the sphere
- map predictions on  $\mathbb{S}^d$  back to the original space

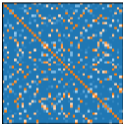


The efficiency of VISH comes from using *spherical harmonics as inducing functions* for the SVGP on the sphere.

## From inducing points to inducing features

Inducing Points

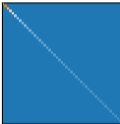
$$u_m = f(z_m)$$

$$K_{uu} =$$


$$K_{uu}^{-1} \text{ is } O(M^3)$$

VISH

$$u_m = hf; \quad m \in H$$

$$K_{uu} =$$


$$K_{uu}^{-1} \text{ is } O(M)$$

Orthogonality of the basisfunctions leads to diagonal  $K_{uu}$  and  $O(M)$  inversion

Deep-dive

# Sparse Variational Gaussian processes

Scalable and flexible

- Capture the GP by a set of inducing variables  $u = f(Z)$ , at locations  $z_1; \dots; z_M$ .

# Sparse Variational Gaussian processes

Scalable and flexible

- Capture the GP by a set of inducing variables  $u = f(Z)$ , at locations  $Z_1; \dots; Z_M$ .
- Minimise KL-divergence from  $p(f(\cdot) | y)$  to  $q(f(\cdot)) = GP(\cdot; (\cdot; \theta))$

$$\begin{aligned}
 (\cdot) &= k_u^>(\cdot) K_{uu}^{-1} m \\
 (\cdot; \theta) &= k(\cdot; \theta) - k_u^>(\cdot) K_{uu}^{-1} (K_{uu} - S) K_{uu}^{-1} k_u(\theta)
 \end{aligned}$$

where  $[K_{uu}]_{m;m^0} = \text{COV}(u_m; u_{m^0})$  and  $[k_u(\cdot)]_m = \text{COV}(u_m; f(\cdot))$ .

# Sparse Variational Gaussian processes

Scalable and flexible

- Capture the GP by a set of inducing variables  $u = f(Z)$ , at locations  $Z_1; \dots; Z_M$ .
- Minimise KL-divergence from  $p(f(\cdot) | y)$  to  $q(f(\cdot)) = GP(\cdot; (\cdot; \theta))$

$$\begin{aligned} q(\cdot) &= k_u^>(\cdot) K_{uu}^{-1} m \\ q(\cdot; \theta) &= k(\cdot; \theta) \left[ k_u^>(\cdot) K_{uu}^{-1} (K_{uu} \quad S) K_{uu}^{-1} k_u(\cdot; \theta) \right] \end{aligned}$$

where  $[K_{uu}]_{m;m^0} = \text{COV}(u_m; u_{m^0})$  and  $[k_u(\cdot)]_m = \text{COV}(u_m; f(\cdot))$ .

- A more flexible (e.g. non-Gaussian likelihoods) and scalable (e.g. mini-batching) model at a cost of  $O(M^3 + M^2 N)$ .



# Sparse Variational Gaussian processes

Scalable and flexible

- Capture the GP by a set of inducing variables  $u = f(Z)$ , at locations  $z_1; \dots; z_M$ .
- Minimise KL-divergence from  $p(f(\cdot) | y)$  to  $q(f(\cdot)) = GP(\cdot; \cdot; \theta)$

$$\begin{aligned} q(\cdot) &= k_u^>(\cdot) K_{uu}^{-1} m \\ q(\cdot; \theta) &= k(\cdot; \theta) \left[ k_u^>(\cdot) K_{uu}^{-1} (K_{uu} - S) K_{uu}^{-1} k_u(\cdot; \theta) \right] \end{aligned}$$

where  $[K_{uu}]_{m;m^0} = \text{COV}(u_m; u_{m^0})$  and  $[k_u(\cdot)]_m = \text{COV}(u_m; f(\cdot))$ .

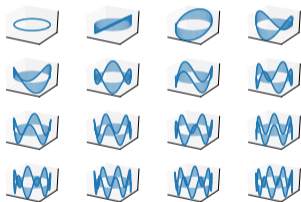
- A more flexible (e.g. non-Gaussian likelihoods) and scalable (e.g. mini-batching) model at a cost of  $O(M^3 + M^2 N)$ .
- Speedup through structure in the  $K_{uu}$  matrix (e.g. Hensman et al 2017, VFF).

# Outline

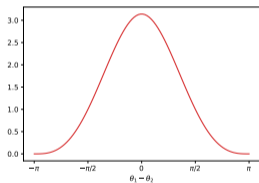
- Gaussian processes on the circle and hypersphere
- Spherical harmonics as inducing features
- Linear projection data on the hyper-sphere

# Gaussian processes on the circle

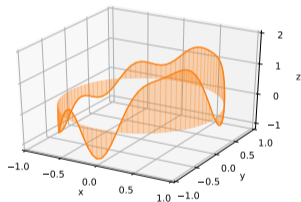
$$k(\theta_1; \theta_2) = [\cos(i\theta_1); \sin(i\theta_1)]^T [\cos(i\theta_2); \sin(i\theta_2)]_{i=0}^7$$



$$k(\theta_1; \theta_2) = \sum_{i=0}^7 \cos(i(\theta_1 - \theta_2))$$

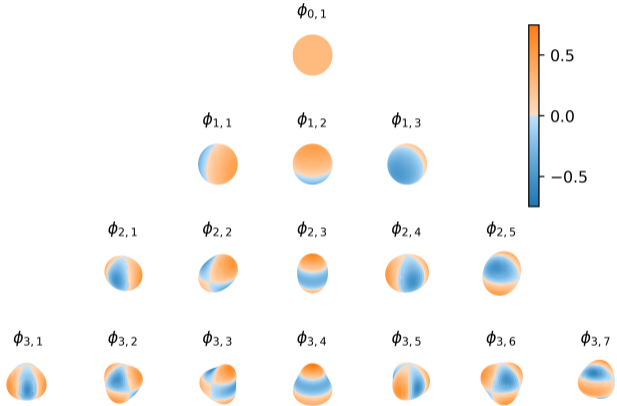


$$f = \sum_{i=0}^7 \beta_i \cos(i\theta); \text{ with } \beta_i \sim N(0; \sigma_i^2)$$



# Spherical Harmonics

- Orthonormal basis on the hyper sphere
- Eigenfunctions the Laplace-Beltrami operator  $\mathbb{S}^{d-1}$   $i = i_1 i_2$
- Eigenfunction of zonal kernels



## Mercer's theorem for zonal kernels on the sphere

- Zonal kernels are the spherical counterpart of stationary kernels  $k(x; x^{\theta}) = k^{\theta}(\text{distance}(x; x^{\theta}))$ .

## Mercer's theorem for zonal kernels on the sphere

- Zonal kernels are the spherical counterpart of stationary kernels  $k(x; x^\theta) = k^\theta(\text{distance}(x; x^\theta))$ .
- Mercer's decomposition: Any zonal kernel  $k$  on the hypersphere can be decomposed as

$$k(x; x^\theta) = \sum_{i=0}^{\infty} \lambda_i \phi_i(x) \phi_i(x^\theta):$$

## Mercer's theorem for zonal kernels on the sphere

- Zonal kernels are the spherical counterpart of stationary kernels  $k(x; x^\theta) = k^\theta(\text{distance}(x; x^\theta))$ .
- Mercer's decomposition: Any zonal kernel  $k$  on the hypersphere can be decomposed as

$$k(x; x^\theta) = \sum_{i=0}^{\infty} \lambda_i \phi_i(x) \phi_i(x^\theta):$$

- Karhunen–Loève expansion: A GP  $f$  on the hypersphere with zonal covariance  $k$  can be written  $f = \sum_i \lambda_i \phi_i$  with  $\phi_i \sim N(0; \lambda_i)$ :

## Spherical harmonics as inducing features in SVGPs

- Define the kernel's RKHS  $H$  with reproducing inner-product:

$$\langle k(x; \cdot); h(\cdot) \rangle_H = h(x)$$



## Spherical harmonics as inducing features in SVGPs

- Define the kernel's RKHS  $H$  with reproducing inner-product:

$$\langle k(x; \cdot); h(\cdot) \rangle_H = h(x)$$

- Approximate posterior constructed out of inducing features

$$u_m = \langle f; \phi_m \rangle_H$$

## Spherical harmonics as inducing features in SVGPs

- Define the kernel's RKHS  $H$  with reproducing inner-product:

$$\langle h(x; \cdot); h(\cdot) \rangle_H = h(x)$$

- Approximate posterior constructed out of inducing features

$$u_m = hf; \quad m \text{ i}_H$$

=) Diagonal covariance matrix:  $[K_{uu}]_{m;m^0} = \text{COV}(u_m; u_{m^0}) = \langle h_{m^0}; h_m \rangle_H = \delta_{m^0 m}^1$

## Spherical harmonics as inducing features in SVGPs

- Define the kernel's RKHS  $H$  with reproducing inner-product:

$$\langle h(\cdot; x); h(\cdot; x) \rangle_H = h(x)$$

- Approximate posterior constructed out of inducing features

$$u_m = \langle h(\cdot; x); f(\cdot) \rangle_H$$

- => Diagonal covariance matrix:  $[K_{uu}]_{m;m^0} = \text{COV}(u_m; u_{m^0}) = \langle h(\cdot; x); h(\cdot; x) \rangle_H = \delta_{mm^0}$
- => Spherical Harmonics as features  $[k_u(\cdot)]_m = \text{COV}(u_m; f(\cdot)) = \langle h(\cdot; x); f(\cdot) \rangle_H$

## Spherical harmonics as inducing features in SVGPs

- Define the kernel's RKHS  $H$  with reproducing inner-product:

$$\langle k(x; \cdot); h(\cdot) \rangle_H = h(x)$$

- Approximate posterior constructed out of inducing features

$$u_m = hf; \quad m \text{ i}_H$$

- $\Rightarrow$  Diagonal covariance matrix:  $[K_{uu}]_{m;m^0} = \text{COV}(u_m; u_{m^0}) = \langle h_{m; \cdot}; h_{m^0; \cdot} \rangle_H = \delta_{mm^0}$
- $\Rightarrow$  Spherical Harmonics as features  $[k_u(\cdot)]_m = \text{COV}(u_m; f(\cdot)) = \langle h_m(\cdot); f(\cdot) \rangle$
- $\Rightarrow$  A  $O(M^2N)$  approximate GP  $q(f(\cdot))$

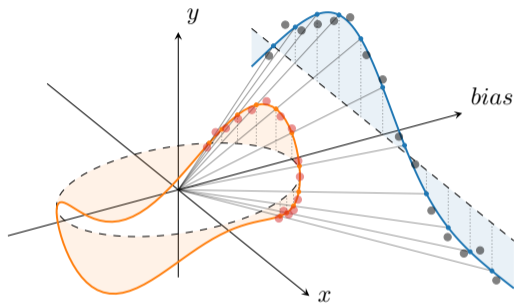
$$GP \quad \mathcal{Y}(\cdot); \quad k(\cdot; \cdot) \quad \mathcal{Y}(\cdot) \quad S \quad \mathcal{Y}(\cdot);$$

where  $\mathcal{Y} = \text{diag}(\mathcal{Y}_1; \dots; \mathcal{Y}_M)$  and  $\mathcal{Y}(\cdot) = [\mathcal{Y}_1(\cdot); \dots; \mathcal{Y}_M(\cdot)]$ .

## Linear mapping to the hypersphere

Most datasets do not correspond to data on a hypersphere...

The proposed solution is to augment the inputs with a constant variable (bias) before projecting it radially onto the hypersphere.

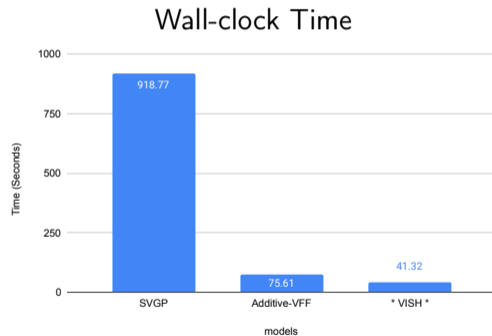
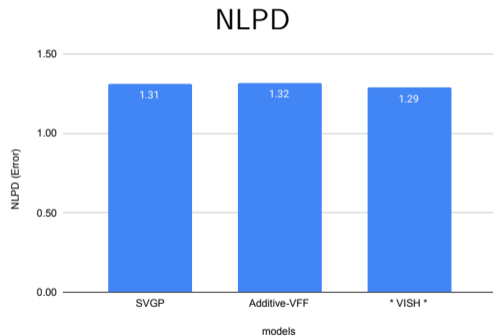


Although such construction may seem arbitrary, it is used implicitly in the Arc-Cosine kernel [Cho & Saul, 2009]:

$$k(x; x^0) = \underbrace{\frac{\|x\| \|x^0\|}{z}}_{\text{radial}} \left( \sin \left( \frac{\theta}{z} \right) + \underbrace{\left( \frac{\|x\| \|x^0\|}{z} \right)}_{\text{angular}} \cos \left( \frac{\theta}{z} \right) \right) \quad \text{with} \quad \theta = \arccos \frac{x \cdot x^0}{\|x\| \|x^0\|}$$

# Experiment

Airline dataset: 6,000,000 datapoints regression task fitted in 40 seconds on a single cheap GTX 1070 GPU



# Conclusion

## Summary of the advantages

- It is the fastest SVGP model to date
  - ) No need for expensive hardware
- The natural ordering of spherical harmonics makes our model scale nicely with the input dimension
  - ) Does not suffer from the curse of dimensionality as VFF
- Similarities with Arc-cosine kernel makes extrapolation properties similar to Neural Networks

Reach out to have a chat if you want to know more!