

Probabilistic models for neural data: From single neurons to population dynamics

NEUROBIO 316QC

Jan Drugowitsch

jan_drugowitsch@hms.harvard.edu

Session 7: Kernel methods

Today

Q&A about previous session

Paper discussion (~1h)

Kernel methods and Gaussian Processes (~30min)

Overview

Kernel methods

Gaussian processes

What makes them possible: marginalization of multivariate Gaussians

Gaussian process regression

Learning hyperparameters, extensions, etc.

Overview

Kernel methods

Gaussian processes

What makes them possible: marginalization of multivariate Gaussians

Gaussian process regression

Learning hyperparameters, extensions, etc.

Linear regression by kernel methods

Assume $p(y_n|\mathbf{x}_n) = N(y_n|\mathbf{w}^T \boldsymbol{\phi}(\mathbf{x}_n), \sigma^2)$ and $p(\mathbf{w}) = (\mathbf{w}|\mathbf{0}, \sigma^2 \lambda^{-1})$ (i.e., L2 regularization)

Negative log-posterior & MAP solution

$$J(\mathbf{w}) = \frac{1}{2\sigma^2} \|\boldsymbol{\Phi} \mathbf{w} - \mathbf{y}\|^2 + \frac{\lambda}{2\sigma^2} \mathbf{w}^T \mathbf{w} + \text{const.} \quad \longrightarrow \quad \hat{\mathbf{w}}_{MAP} = -\frac{1}{\lambda} \boldsymbol{\Phi}^T (\boldsymbol{\Phi} \hat{\mathbf{w}}_{MAP} - \mathbf{y}) = \boldsymbol{\Phi}^T \mathbf{a}$$

Substitute $\hat{\mathbf{w}}_{MAP}$ into negative log-posterior, becomes function of \mathbf{a}

$$J(\mathbf{a}) = \frac{1}{2\sigma^2} \mathbf{a}^T \mathbf{K} \mathbf{K} \mathbf{a} - \frac{1}{\sigma^2} \mathbf{a}^T \mathbf{K} \mathbf{y} + \frac{1}{2\sigma^2} \mathbf{y}^T \mathbf{y} + \frac{\lambda}{2\sigma^2} \mathbf{a}^T \mathbf{K} \mathbf{a} \quad \text{with Gram matrix } \mathbf{K} = \boldsymbol{\Phi} \boldsymbol{\Phi}^T$$

$$\hat{\mathbf{a}}_{MAP} = (\mathbf{K} + \lambda \mathbf{I})^{-1} \mathbf{y}$$

$$K_{nm} = \underbrace{\boldsymbol{\phi}(\mathbf{x}_n)^T \boldsymbol{\phi}(\mathbf{x}_m)}_{\sim \text{similarity between } \mathbf{x}_n \text{ and } \mathbf{x}_m} = k(\mathbf{x}_n, \mathbf{x}_m)$$

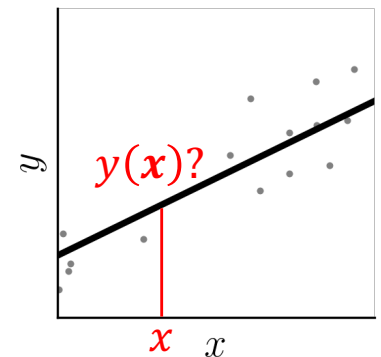
\sim similarity between \mathbf{x}_n and $\mathbf{x}_m \rightarrow$ kernel $k(\cdot, \cdot)$

Predict y for new x

similarity of x to training “inputs”

$$y(\mathbf{x}) = \hat{\mathbf{w}}_{MAP}^T \boldsymbol{\phi}(\mathbf{x}) = \boldsymbol{\phi}^T(\mathbf{x}) \boldsymbol{\Phi}^T \mathbf{a} = \underbrace{k(\mathbf{x})^T (\mathbf{K} + \lambda \mathbf{I})^{-1} \mathbf{y}}_{\text{weighted combination of training “outputs”}}$$

weighted combination of training “outputs”



Valid kernel functions & computational complexity

$K_{nm} = k(\mathbf{x}_n, \mathbf{x}_m) = \boldsymbol{\phi}(\mathbf{x}_n)^T \boldsymbol{\phi}(\mathbf{x}_m)$, but can specify $k(\mathbf{x}_n, \mathbf{x}_m)$ without $\boldsymbol{\phi}(\cdot)$

Gram matrix needs to be positive definite, i.e., $\mathbf{a}^T \mathbf{K} \mathbf{a} \geq 0$ for all \mathbf{a}

Examples

Polynomial kernel

$$k(\mathbf{x}, \mathbf{x}') = (\mathbf{x}^T \mathbf{x}' + c)^M$$

“Gaussian” kernel

$$k(\mathbf{x}, \mathbf{x}') = \exp(-\|\mathbf{x} - \mathbf{x}'\|^2 / 2\sigma^2)$$

Generalized “Gaussian” kernel

$$k(\mathbf{x}, \mathbf{x}') = \exp\left(-\frac{1}{2\sigma^2} \left(\tilde{k}(\mathbf{x}, \mathbf{x}) + \tilde{k}(\mathbf{x}', \mathbf{x}') - 2\tilde{k}(\mathbf{x}, \mathbf{x}')\right)\right)$$

Computational complexity

Assume N “inputs”, $\mathbf{x}_1, \dots, \mathbf{x}_N$, D -dimensional basis function $\boldsymbol{\phi}(\cdot)$ and \mathbf{w} (usually $N \gg D$)

“standard” linear regression

$$y(\mathbf{x}) = \boldsymbol{\phi}^T(\mathbf{x}) \hat{\mathbf{w}}_{MAP} = \boldsymbol{\phi}^T(\mathbf{x}) \underbrace{(\boldsymbol{\Phi}^T \boldsymbol{\Phi} + \lambda \mathbf{I})^{-1}}_{D \times D} \boldsymbol{\Phi}^T \mathbf{y}$$

Kernel regression

$$y(\mathbf{x}) = \boldsymbol{\phi}^T(\mathbf{x}) \boldsymbol{\Phi}^T \hat{\mathbf{a}}_{MAP} = k(\mathbf{x})^T \underbrace{(\boldsymbol{\Phi} \boldsymbol{\Phi}^T + \lambda \mathbf{I})^{-1}}_{N \times N} \mathbf{y}$$

Overview

Kernel methods

Gaussian processes

What makes them possible: marginalization of multivariate Gaussians

Gaussian process regression

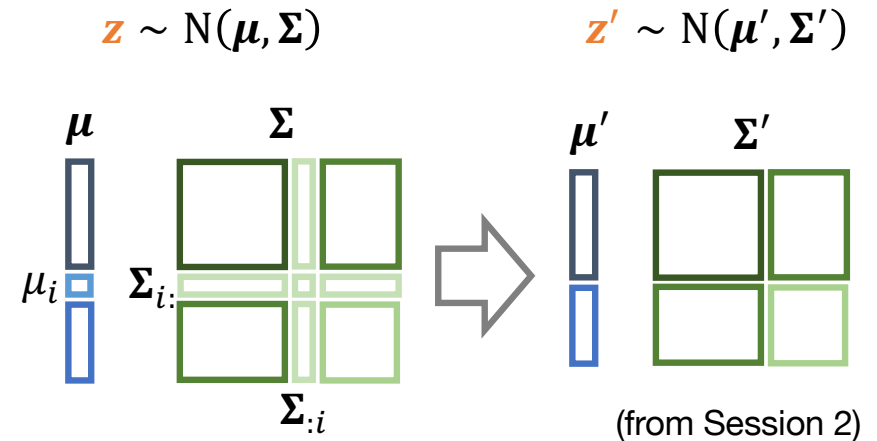
Learning hyperparameters, extensions, etc.

Marginalizing multivariate Gaussians → Gaussian process

Marginalization: ‘removing’ z_i from \mathbf{z} by

$$p(\mathbf{z}') = p(z_1, \dots, z_{i-1}, z_{i+1}, \dots, z_N) = \int p(\mathbf{z}) dz_i$$

does not change moments of $z_{j \neq i}$
(as moments of z_j only depend on μ_j and \mathbf{a}_j .)



Gaussian process $z(\mathbf{x})$ (continuous function)

Distribution over any $z(\mathbf{x}_1), \dots, z(\mathbf{x}_N)$ is multivariate Gaussian

$$\mathbf{z} \equiv z(\mathbf{x}_1), \dots, z(\mathbf{x}_N) \sim N(0, \mathbf{K}) \quad \text{with covariance } K_{nm} = \underbrace{k(\mathbf{x}_n, \mathbf{x}_m)}_{\text{expected similarity of } z(\mathbf{x}_n) \text{ and } z(\mathbf{x}_m)}$$

↑
positive (semi-)definite

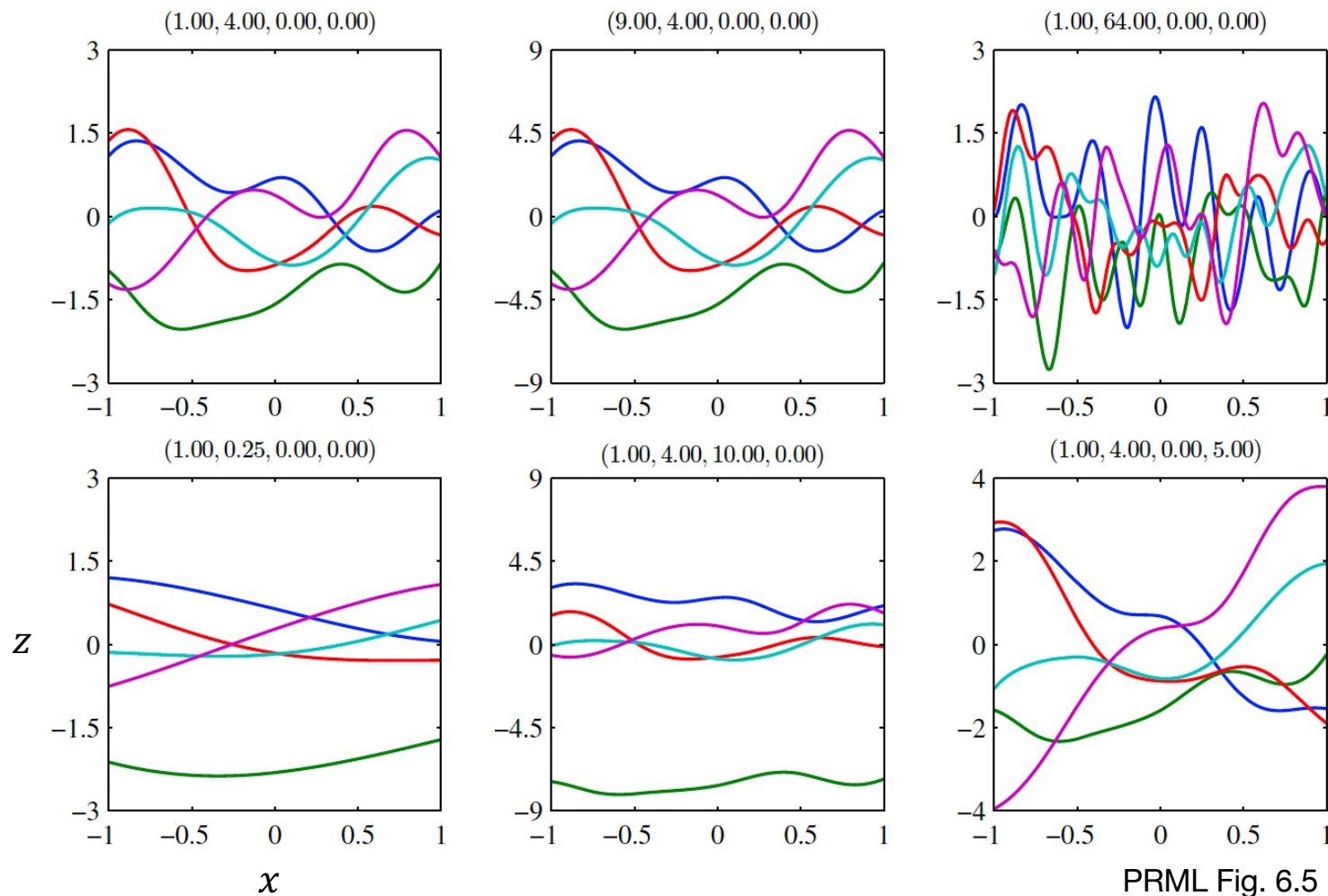
Construction possible as “marginalizing out” unobserved $z(\mathbf{x})$ ’s yields Gaussian

Also: adding another $z(\mathbf{x}_{N+1})$ again yields Gaussian

Draws from Gaussian processes

magnitude length-scale long-range correlations linear in \mathbf{x} and \mathbf{x}'

$$k(\mathbf{x}, \mathbf{x}') = \theta_0 \exp\left(-\frac{\theta_1}{2} \|\mathbf{x} - \mathbf{x}'\|^2\right) + \theta_2 + \theta_3 \mathbf{x}^T \mathbf{x}'$$



Gaussian progress regression

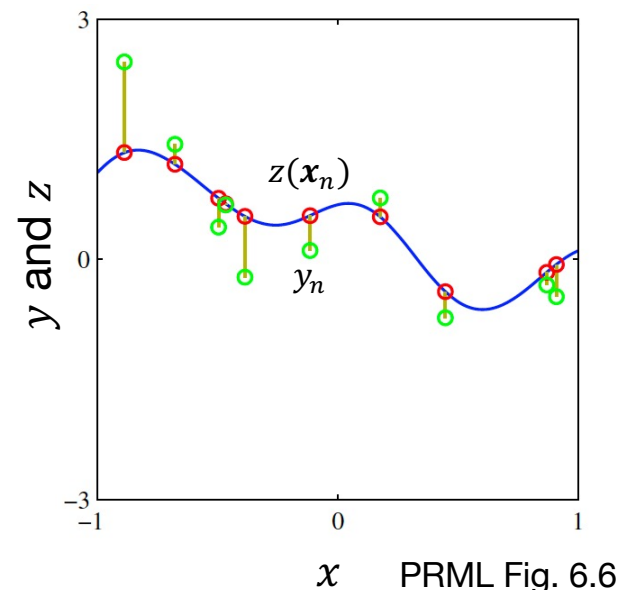
Model data as draws from GP + noise

$$\begin{array}{ll} \text{GP} & \mathbf{z}|\mathbf{x}_{1:N} \sim \mathcal{N}(\mathbf{0}, \mathbf{K}) \\ \text{noise} & y_n|z_n \sim \mathcal{N}(z_n, \lambda) \longrightarrow \mathbf{y}|\mathbf{z} \sim \mathcal{N}(\mathbf{y}|\mathbf{z}, \lambda\mathbf{I}) \end{array}$$

Yields marginal \mathbf{y}

$$p(\mathbf{y}|\mathbf{x}_{1:N}) = \int p(\mathbf{y}|\mathbf{z})p(\mathbf{z}|\mathbf{x}_{1:N})d\mathbf{z} = \mathcal{N}(\mathbf{y}|\mathbf{0}, \underbrace{\mathbf{K} + \lambda\mathbf{I}}_{\mathbf{C}_N})$$

\swarrow GP \swarrow indep. noise



Predict $y_{N+1:N+M}$ from $\mathbf{x}_{N+1:N+M}$

1. Joint distribution over $\mathbf{y}_{N+M} \equiv y_{1:N+M}$

$$p(\mathbf{y}_{N+M}|\mathbf{x}_{1:N+M}) = \mathcal{N}(\mathbf{y}_{N+M}|\mathbf{0}, \mathbf{C}_{N+M}) \quad \text{with } \mathbf{C}_{N+M} = \begin{pmatrix} \mathbf{C}_N & \mathbf{k} \\ \mathbf{k}^T & \mathbf{c} \end{pmatrix}$$

\swarrow GP cov. between $\mathbf{x}_{1:N}$ and $\mathbf{x}_{N+1:N+M}$
 \swarrow $\mathbf{x}_{N+1:N+M}$ GP cov. + noise

2. Condition on $\mathbf{y} \equiv y_{1:N}$

$$\begin{aligned} p(y_{N+1:N+M}|\mathbf{y}, \mathbf{x}_{1:N+M}) &= \mathcal{N}(y_{N+1:N+M} | \underbrace{\mathbf{k}^T \mathbf{C}_N^{-1} \mathbf{y}}_{\text{as in Kernel regression}}, \mathbf{c} - \mathbf{k}^T \mathbf{C}_N^{-1} \mathbf{k}) \\ &= \mathbf{k}^T (\mathbf{K} + \lambda\mathbf{I})^{-1} \mathbf{y} \end{aligned}$$

as in Kernel regression

Learning the GP hyperparameters, extensions & complexity

Hyperparameters = parameters θ of the GP kernel

Standard approach: maximum likelihood

$$\mathbf{C}_N = \mathbf{K} + \lambda \mathbf{I}, \text{ function of } \theta$$

$$\hat{\theta} = \operatorname{argmax}_{\theta} \log p(\mathbf{y} | \mathbf{x}_{1:N}, \theta) = \operatorname{argmax}_{\theta} \left(-\frac{1}{2} \log |\mathbf{C}_N| - \frac{1}{2} \mathbf{y}^T \mathbf{C}_N^{-1} \mathbf{y} \right)$$

Challenge: potentially many local maxima

Extensions

GP kernels matching AR(n) process, LDS, OU process, etc. (see Rasmussen & Williams, 2006)

Kernel classification: non-conjugate likelihood, requires (Laplace) approximation

General: use of non-conjugate likelihood (e.g., Poisson spiking) requires approximations

Computational complexity

Use of GP requires $N \times N$ matrix inversion, $O(N^3)$ complexity

Low-dimensional approximation of covariance matrix K (e.g., Fourier decomposition)

Sparse GP: replace data with representative examples

Overview

Kernel methods

Gaussian processes

What makes them possible: marginalization of multivariate Gaussians

Gaussian process regression

Learning hyperparameters, extensions, etc.

Summary

Kernel methods: perform regression in (dual) space of kernels rather than basis functions

Can handle (effectively) infinitely-dimensional basis functions, at cost of increased complexity

Akin to interpolation: predictions by similarity-based averaging of training data

Gaussian processes: priors over (smooth) functions

Made possible due to nice marginalization properties of multivariate Gaussians

Gaussian process regression like kernel regression + measure of uncertainty

Non-conjugate likelihoods require approximations

Until next week

Read paper and prepare presentation (see notes for Session 7)

Read statistical methods section (no separate PRML reading, only session notes)

Next session

Q & A for previous session

Paper discussions (~1h)

Brief introduction to variational autoencoders (~30min)

Note: one-week break, next session on March 23

