Gaussian processes for data-driven modelling and uncertainty quantification: a hands-on tutorial

Andreas Damianou

Department of Computer Science, University of Sheffield, UK

Brown University, 16/02/2016

Sheffield



Outline

- 1. Gaussian processes as infinite dimensional Gaussian distributions
 - 1.1 Gaussian distribution
 - 1.2 Intuition by sampling and plotting
 - 1.3 Mean and covariance functions
 - 1.4 Marginalization and conditioning properties
- 2. Noise model
- 3. Covariance functions, aka kernels
- 4. "Full" GP implementation!
- 5. Running our GP
 - 5.1 Fitting and overfitting
- 6. GPy
- 7. Classification













- A Gaussian distribution depends on a mean and a covariance matrix.
- A Gaussian process depends on a mean and a covariance function.




























































GO TO NOTEBOOK (start)

Infinite model... but we always work with finite sets!

Let's start with a multivariate Gaussian:

$$p(\underbrace{f_1, f_2, \cdots, f_s}_{\mathbf{f}_A}, \underbrace{f_{s+1}, f_{s+2}, \cdots, f_N}_{\mathbf{f}_B}) \sim \mathcal{N}(\boldsymbol{\mu}, \mathbf{K}).$$

with:

$$oldsymbol{\mu} = egin{bmatrix} oldsymbol{\mu}_A \ oldsymbol{\mu}_B \end{bmatrix}$$
 and $oldsymbol{K} = egin{bmatrix} oldsymbol{K}_{AA} & oldsymbol{K}_{AB} \ oldsymbol{K}_{BA} & oldsymbol{K}_{BB} \end{bmatrix}$

Marginalisation property:

$$p(\mathbf{f}_A, \mathbf{f}_B) \sim \mathcal{N}(\boldsymbol{\mu}, \mathbf{K}).$$
 Then:
 $p(\mathbf{f}_A) = \int_{\mathbf{f}_B} p(\mathbf{f}_A, \mathbf{f}_B) d\mathbf{f}_B = \mathcal{N}(\boldsymbol{\mu}_A, \mathbf{K}_{AA})$

Infinite model... but we always work with finite sets!

Let's start with a multivariate Gaussian:

$$p(\underbrace{f_1, f_2, \cdots, f_s}_{\mathbf{f}_A}, \underbrace{f_{s+1}, f_{s+2}, \cdots, f_N}_{\mathbf{f}_B}) \sim \mathcal{N}(\boldsymbol{\mu}, \mathbf{K}).$$

with:

$$\boldsymbol{\mu} = \begin{bmatrix} \boldsymbol{\mu}_A \\ \boldsymbol{\mu}_B \end{bmatrix}$$
 and $\mathbf{K} = \begin{bmatrix} \mathbf{K}_{AA} & \mathbf{K}_{AB} \\ \mathbf{K}_{BA} & \mathbf{K}_{BB} \end{bmatrix}$

Marginalisation property:

$$p(\mathbf{f}_A, \mathbf{f}_B) \sim \mathcal{N}(\boldsymbol{\mu}, \mathbf{K}).$$
 Then:
 $p(\mathbf{f}_A) = \int_{\mathbf{f}_B} p(\mathbf{f}_A, \mathbf{f}_B) d\mathbf{f}_B = \mathcal{N}(\boldsymbol{\mu}_A, \mathbf{K}_{AA})$

Infinite model... but we always work with finite sets!

In the GP context $\mathbf{f} = f(x)$:

$$\boldsymbol{\mu}_{\infty} = \begin{bmatrix} \boldsymbol{\mu}_{\rm f} \\ \cdots \\ \cdots \end{bmatrix} \text{ and } \mathbf{K}_{\infty} = \begin{bmatrix} \mathbf{K}_{\rm ff} & \cdots \\ \cdots & \cdots \end{bmatrix}$$

Covariance function: Maps locations x_i, x_j of the input domain \mathcal{X} to an entry in the covariance matrix:

$$\mathbf{K}_{i,j} = k(\mathbf{x}_i, \mathbf{x}_j)$$

For all available inputs:

$$\mathbf{K} = \mathbf{K}_{\!\!\mathbf{f} \!\!\mathbf{f}} = k(\mathbf{X}, \mathbf{X})$$

GP: joint Gaussian distribution of the training and the (potentially infinite!) test data:

$$\begin{bmatrix} \mathbf{f} \\ \mathbf{f}^* \end{bmatrix} \sim \mathcal{N} \left(\mathbf{0}, \begin{bmatrix} \mathbf{K} & \mathbf{K}_* \\ \mathbf{K}_*^\top & \mathbf{K}_{*,*} \end{bmatrix} \right)$$

 \mathbf{K}_* is the (cross)-covariance matrix obtained by evaluating the covariance function in pairs of training inputs \mathbf{X} and test inputs $\mathbf{X}_*,$ ie.

$$\mathbf{f}_* = k(\mathbf{X}, \mathbf{X}_*).$$

Similarly:

$$\mathbf{K}_{**} = k(\mathbf{X}_*, \mathbf{X}_*).$$

$$\begin{split} p(\mathbf{f}_A, \mathbf{f}_B) &\sim \mathcal{N}(\boldsymbol{\mu}, \mathbf{K}). \quad \text{Then:} \\ p(\mathbf{f}_A | \mathbf{f}_B) &= \mathcal{N}(\boldsymbol{\mu}_A + \mathbf{K}_{AB} \mathbf{K}_{BB}^{-1}(\mathbf{f}_B - \boldsymbol{\mu}_B), \mathbf{K}_{AA} - \mathbf{K}_{AB} \mathbf{K}_{BB}^{-1} \mathbf{K}_{BA}) \end{split}$$

In the GP context this can be used for inter/extrapolation:

$$p(f_*|f_1, \cdots, f_N) = p(f(x_*)|f(x_1), \cdots, f(x_N)) \sim \mathcal{N}$$
$$p(\mathbf{f}_*|\mathbf{f}_1, \cdots, \mathbf{f}_N) = p(f(x_*)|f(x_1), \cdots, f(x_N))$$
$$\sim \mathcal{N}(\mathbf{K}_*^{\top}\mathbf{K}^{-1}\mathbf{f}, \mathbf{K}_{*,*} - \mathbf{K}_*^{\top}\mathbf{K}^{-1}\mathbf{K}_*)$$

 $p(f(x_*)|f(x_1),\cdots,f(x_N))$ is a posterior **process**!

$$\begin{split} p(\mathbf{f}_A, \mathbf{f}_B) &\sim \mathcal{N}(\boldsymbol{\mu}, \mathbf{K}). \quad \text{Then:} \\ p(\mathbf{f}_A | \mathbf{f}_B) &= \mathcal{N}(\boldsymbol{\mu}_A + \mathbf{K}_{AB} \mathbf{K}_{BB}^{-1}(\mathbf{f}_B - \boldsymbol{\mu}_B), \mathbf{K}_{AA} - \mathbf{K}_{AB} \mathbf{K}_{BB}^{-1} \mathbf{K}_{BA}) \end{split}$$

In the GP context this can be used for inter/extrapolation:

$$p(f_*|f_1, \cdots, f_N) = p(f(x_*)|f(x_1), \cdots, f(x_N)) \sim \mathcal{N}$$
$$p(\mathbf{f}_*|\mathbf{f}_1, \cdots, \mathbf{f}_N) = p(f(x_*)|f(x_1), \cdots, f(x_N))$$
$$\sim \mathcal{N}(\mathbf{K}_*^\top \mathbf{K}^{-1} \mathbf{f} , \mathbf{K}_{*,*} - \mathbf{K}_*^\top \mathbf{K}^{-1} \mathbf{K}_*)$$

 $p(f(x_*)|f(x_1), \cdots, f(x_N))$ is a posterior **process**!

Noise model

- So far we assumed: $\mathbf{f} = f(\mathbf{X})$
- Assuming that we only observe noisy versions y of the true outputs f:

 $\begin{aligned} \mathbf{y} &= f(\mathbf{X}) + \epsilon, \text{ where:} \\ f &\sim \mathcal{GP}(0, k(x, x')) \\ \epsilon &\sim \mathcal{N}(0, \sigma^2 \mathbf{I}) \end{aligned}$

The above construction, gives us the following probabilities:

$$p(\mathbf{y}|\mathbf{f}) = \mathcal{N}(\mathbf{y}|\mathbf{f}, \sigma^2 \mathbf{I})$$

$$p(\mathbf{f}|\mathbf{x}) = \mathcal{N}(\mathbf{f}|\mathbf{0}, K_{ff}) = (2\pi)^{n/2} |K_{ff}|^{-1/2} \exp\left(-\frac{1}{2}\mathbf{f}^T K_{ff}\mathbf{f}\right)$$

$$p(\mathbf{y}|\mathbf{x}) = \int p(\mathbf{y}|\mathbf{f}) p(\mathbf{f}|\mathbf{x}) d\mathbf{f} = \mathcal{N}(\mathbf{y}|\mathbf{0}, K_{ff} + \sigma^2 \mathbf{I})$$

 $p(\mathbf{y}|\mathbf{x})$ is called the **marginal likelihood** and is tractable because of our choice for noise ϵ which is normally distributed.

Noise model

- So far we assumed: $\mathbf{f} = f(\mathbf{X})$
- Assuming that we only observe noisy versions y of the true outputs f:

 $\begin{aligned} \mathbf{y} &= f(\mathbf{X}) + \epsilon, \text{ where:} \\ f &\sim \mathcal{GP}(0, k(x, x')) \\ \epsilon &\sim \mathcal{N}(0, \sigma^2 \mathbf{I}) \end{aligned}$

The above construction, gives us the following probabilities:

$$p(\mathbf{y}|\mathbf{f}) = \mathcal{N}(\mathbf{y}|\mathbf{f}, \sigma^2 \mathbf{I})$$

$$p(\mathbf{f}|\mathbf{x}) = \mathcal{N}(\mathbf{f}|\mathbf{0}, K_{ff}) = (2\pi)^{n/2} |K_{ff}|^{-1/2} \exp\left(-\frac{1}{2}\mathbf{f}^T K_{ff}\mathbf{f}\right)$$

$$p(\mathbf{y}|\mathbf{x}) = \int p(\mathbf{y}|\mathbf{f}) p(\mathbf{f}|\mathbf{x}) d\mathbf{f} = \mathcal{N}(\mathbf{y}|\mathbf{0}, K_{ff} + \sigma^2 \mathbf{I})$$

 $p(\mathbf{y}|\mathbf{x})$ is called the **marginal likelihood** and is tractable because of our choice for noise ϵ which is normally distributed.

Noise model

- So far we assumed: $\mathbf{f} = f(\mathbf{X})$
- Assuming that we only observe noisy versions y of the true outputs f:

 $\begin{aligned} \mathbf{y} &= f(\mathbf{X}) + \epsilon, \text{ where:} \\ f &\sim \mathcal{GP}(0, k(x, x')) \\ \epsilon &\sim \mathcal{N}(0, \sigma^2 \mathbf{I}) \end{aligned}$

The above construction, gives us the following probabilities:

$$p(\mathbf{y}|\mathbf{f}) = \mathcal{N}(\mathbf{y}|\mathbf{f}, \sigma^2 \mathbf{I})$$

$$p(\mathbf{f}|\mathbf{x}) = \mathcal{N}(\mathbf{f}|\mathbf{0}, K_{ff}) = (2\pi)^{n/2} |K_{ff}|^{-1/2} \exp\left(-\frac{1}{2}\mathbf{f}^T K_{ff}\mathbf{f}\right)$$

$$p(\mathbf{y}|\mathbf{x}) = \int p(\mathbf{y}|\mathbf{f}) p(\mathbf{f}|\mathbf{x}) d\mathbf{f} = \mathcal{N}(\mathbf{y}|\mathbf{0}, K_{ff} + \sigma^2 \mathbf{I})$$

 $p(\mathbf{y}|\mathbf{x})$ is called the **marginal likelihood** and is tractable because of our choice for noise ϵ which is normally distributed.

Predictions in the noise model

$$\begin{split} \mathbf{y}^* | \mathbf{y}, \mathbf{x}, \mathbf{x}_* \sim \mathcal{N}(\boldsymbol{\mu}_{\mathsf{pred}}, \mathbf{K}_{\mathsf{pred}}) \\ \text{with} \\ \boldsymbol{\mu}_{\mathsf{pred}} = \mathbf{K}_*^\top \left[\mathbf{K} + \sigma^2 \mathbf{I} \right]^{-1} \mathbf{y} \\ \text{and} \\ \mathbf{K}_{\mathsf{pred}} = \mathbf{K}_{*,*} - \mathbf{K}_*^\top \left[\mathbf{K} + \sigma^2 \mathbf{I} \right]^{-1} \mathbf{K}_*. \end{split}$$

Fitting the data (shaded area is uncertainty)



Fitting the data - Prior Samples







Fitting the data - more noise



Fitting the data - no noise



Fitting the data - Posterior samples















GO TO NOTEBOOK (Cov. functions)



Which curve fits the data better?



- Which curve fits the data better?
- ▶ Which curve is more "complex"?



- Which curve fits the data better?
- ▶ Which curve is more "complex"?
- Which curve is better overall?



- Which curve fits the data better?
- ▶ Which curve is more "complex"?
- Which curve is better overall?

Need a good balance between data fit vs overfitting!

(Bayesian) Occam's Razor

"A plurality is not to be posited without necessity". *W. of Ockham* "Everything should be made as simple as possible, but not simpler". *A. Einstein*



Evidence is higher for the model that is not "unnecessarily complex" but still "explains" the data D.

How do GPs solve the overfitting problem (i.e. regularize)?

- Answer: Integrate over the function itself!
- ► This is associated with the Bayesian methodology.
- So, we will average out all possible function forms, under a (GP) prior!

Recap:

$$\begin{array}{ll} \mathsf{ML:} & \operatorname*{argmax}_{\mathbf{w}} p(\mathbf{y} | \mathbf{w}, \phi(\mathbf{x})) & \text{e.g. } \mathbf{y} = \phi(\mathbf{x})^\top \mathbf{w} + \epsilon \\ \mathsf{Bayesian:} & \operatorname*{argmax}_{\theta} \int_{\mathbf{f}} p(\mathbf{y} | \mathbf{f}) \underbrace{p(\mathbf{f} | \mathbf{x}, \theta)}_{\mathsf{GP prior}} & \text{e.g. } \mathbf{y} = f(\mathbf{x}, \theta) + \epsilon \\ \end{array}$$

- θ are hyperparameters
- The Bayesian approach (GP) automatically balances the data-fitting with the complexity penalty.

How do GPs solve the overfitting problem (i.e. regularize)?

- Answer: Integrate over the function itself!
- ► This is associated with the Bayesian methodology.
- So, we will average out all possible function forms, under a (GP) prior!

Recap:

$$\begin{array}{ll} \mathsf{ML:} & \operatorname*{argmax}_{\mathbf{w}} p(\mathbf{y} | \mathbf{w}, \phi(\mathbf{x})) & \mathsf{e.g.} \ \mathbf{y} = \phi(\mathbf{x})^\top \mathbf{w} + \boldsymbol{\epsilon} \\ \mathsf{Bayesian:} & \operatorname*{argmax}_{\boldsymbol{\theta}} \ \int_{\mathbf{f}} p(\mathbf{y} | \mathbf{f}) \underbrace{p(\mathbf{f} | \mathbf{x}, \boldsymbol{\theta})}_{\mathsf{GP \ prior}} & \mathsf{e.g.} \ \mathbf{y} = f(\mathbf{x}, \boldsymbol{\theta}) + \boldsymbol{\epsilon} \end{array}$$

- θ are *hyper*parameters
- The Bayesian approach (GP) automatically balances the data-fitting with the complexity penalty.
GO TO NOTEBOOK

- Unsupervised learning with GPs (Bayesian non-linear dimensionality reduction)
- Deep GPs
- Dynamical systems
- ...and more...

Deep GP: Step function (credits for idea to J. Hensman)



Thanks to Prof. Neil Lawrence, his SheffieldML group and the GPy developers.

Thanks to George Karniadakis, Paris Perdikaris for hosting me

EOF