

Advanced Machine Learning Subsidiary Notes

Lecture 17: Gaussian Processes

Adam Prügel-Bennett

April 19, 2021

1 Keywords

- Gaussian Processes, regression

2 Main Points

2.1 The Big Picture

- Gaussian processes are a machine learning method for regression
 - You can use them for classification but it is regression where they excel
- They put a prior on functions (preferring smooth functions)
- Conceptually they are easy
- They are easy to use
- But they are mathematically fiddly to understand
- Gaussian processes were first used for spatial modelling where the technique was known as *kriging*
- They are one of the most powerful methods we have of doing regression
- They are quite general purpose as the prior we use is reasonable for a large number of applications

2.2 Gaussian Processes

- Gaussian Processes are also mathematical objects
- A Gaussian Process assigns a probability to functions depending on their spatial smoothness
 - We assume that the probability of $f(\mathbf{x})$ taking a particular value is normally/Gaussian distributed in a way that depends on the value of the function at other points
- This will be governed by a *covariance* or *kernel* function $k(\mathbf{x}, \mathbf{y})$ which tells us the covariance of the prior

$$\mathbb{E}_f [(f(\mathbf{x}) - m(\mathbf{x})) (f(\mathbf{y}) - m(\mathbf{y}))] = k(\mathbf{x}, \mathbf{y})$$

- This is an expectation over all functions from our probability distribution
- $m(\mathbf{x})$ encodes a prior belief in the value of $f(\mathbf{x})$ —normally we choose this to be 0 as most often we don't have any prior knowledge of the function
- Because the covariance is quadratic it will be positive semi-definite

- The kernel function is critical to determining the type of function we will learn
- Choosing this kernel function correctly is essential to make GP work
- The values at each point \mathbf{x} is normally or Gaussianly distributed

$$p(f|m, k) \propto e^{-\frac{1}{2} \int (f(\mathbf{x}) - m(\mathbf{x})) k^{-1}(\mathbf{x}, \mathbf{y}) (f(\mathbf{y}) - m(\mathbf{y})) d\mathbf{x} d\mathbf{y}}$$

- It is a little unclear exactly what this means for functions: working with function spaces one has to be very careful (e.g. to normalise this we have to sum over all functions)
- Nevertheless it is clear that we can use this for comparing the likelihood of different functions
- If you discretise the domain of \mathbf{x} then

$$p(f|m, k) \propto e^{-\frac{1}{2} \sum_{i,j} (f(\mathbf{x}_i) - m(\mathbf{x}_i)) k^{-1}(\mathbf{x}_i, \mathbf{x}_j) (f(\mathbf{x}_j) - m(\mathbf{x}_j))}$$

- * the probability function is now a well-defined multi-dimensional normal distribution
- * here it is clear that the value at each point is normally distributed
- * to get back to the integral form we have to take a limit where our discrete points become every closer together (this stresses out mathematicians who point out that nasty things can happen, but for the rest of us we just go ahead)
- We can sample Gaussian processes from this distribution
 - * in low dimensions this allows us to visualise the type of functions we are likely to learn
 - * you can try this out in the experiment section

2.3 Bayesian Inference

- We can use Gaussian Processes as a prior for performing regression
- We assume we have a set of data points $\mathcal{D} = ((\mathbf{x}_i, y_i) | i = 1, \dots, m)$
- The likelihood for data is taken to be a normal distribution

$$p(\mathcal{D}|f) = \prod_{i=1}^m \mathcal{N}(y_i | f(\mathbf{x}_i), \sigma^2)$$

- We assume that $f(\mathbf{x}_i)$ is the true value of the function at \mathbf{x}_i , but, because of measurement error, y_i will be distributed according to a normal distribution
- σ is the accuracy of our measured data (it is another hyper-parameter we must choose)
- Note we are assuming the data points are independent of each other (which is usually a good approximation)
- We can now compute a posterior

$$p(f|\mathcal{D}, m, k, \sigma) = \frac{p(\mathcal{D}|f, \sigma) p(f|m, k)}{p(\mathcal{D}|m, k, \sigma)}$$

- I have explicitly written in the dependencies on the hyper-parameters m , k and σ
- The Gaussian Process prior is a conjugate distribution for the normal likelihood
- Computing probabilities over functions is slightly wild
- It is more useful to compute the probability at a particular point \mathbf{x}^*

- To compute this we marginalise over all other points
 - * This is non-trivial but doable
 - * If you have a multivariate normal distribution

$$\mathcal{N}(\mathbf{x}|\boldsymbol{\mu}, \mathbf{C}) = \frac{1}{\sqrt{|2\pi\mathbf{C}|}} e^{-\frac{1}{2}(\mathbf{x}-\boldsymbol{\mu})^\top \mathbf{C}(\mathbf{x}-\boldsymbol{\mu})}$$

and you integrate over x_i (marginalise it out) then

$$\int_{-\infty}^{\infty} \mathcal{N}(\mathbf{x}|\boldsymbol{\mu}, \mathbf{C}) dx_i = \mathcal{N}(\hat{\mathbf{x}}|\hat{\boldsymbol{\mu}}, \hat{\mathbf{C}})$$

where $\hat{\mathbf{x}}$ and $\hat{\boldsymbol{\mu}}$ are identical to \mathbf{x} and $\boldsymbol{\mu}$ except with their i^{th} component removed and $\hat{\mathbf{C}}$ is identical to \mathbf{C} except with the i^{th} row and column removed

- * Intuitively this seems very natural (if x_i can take any value it doesn't change the mean or covariance between other random variables)
- * You can prove this algebraically but it is not a trivial calculation
- * Because of this for a Gaussian Process we only care about those set of points where we have data values
- For all point \mathbf{x} (except the point \mathbf{x}^*) we integrate over the possible values that $f(\mathbf{x})$ can take
- Where we have no prior this integral is just a constant—and will cancel with the denominator $p(\mathcal{D}|m, k, \sigma)$
- We are left with

$$\begin{aligned} p(f(\mathbf{x}^*)|\mathcal{D}) \propto & \left(\prod_i \int df(\mathbf{x}_i) \right) e^{-\frac{1}{2\sigma^2} \sum_{i=1}^m (y_i - f(\mathbf{x}_i))^2} e^{-\frac{1}{2} \sum_{i,j=1}^m (f(\mathbf{x}_i) - m(\mathbf{x}_i)) k^{-1}(\mathbf{x}_i, \mathbf{x}_j) (f(\mathbf{x}_j) - m(\mathbf{x}_j))} \\ & e^{-\sum_{j=1}^m (f(\mathbf{x}^*) - m(\mathbf{x}^*)) k^{-1}(\mathbf{x}^*, \mathbf{x}_j) (f(\mathbf{x}_j) - m(\mathbf{x}_j))} \\ & e^{-\frac{1}{2} (f(\mathbf{x}^*) - m(\mathbf{x}^*)) k^{-1}(\mathbf{x}^*, \mathbf{x}^*) (f(\mathbf{x}^*) - m(\mathbf{x}^*))} \end{aligned}$$

- * Writing $f^* = f(\mathbf{x}^*)$, $f_i = f(\mathbf{x}_i)$, $m_i = m(\mathbf{x}_i)$, $k_{ij}^{-1} = k^{-1}(\mathbf{x}_i, \mathbf{x}_j)$, $k_j^{*-1} = k^{-1}(\mathbf{x}^*, \mathbf{x}_j)$ and $k^{*-1} = k^{-1}(\mathbf{x}^*, \mathbf{x}^*)$ then

$$\begin{aligned} p(f^*|\mathcal{D}) \propto & e^{-\frac{1}{2} k^{*-1} (f^* - m^*)^2} \left(\prod_i \int df_i \right) e^{-\frac{1}{2\sigma^2} \sum_{i=1}^m (y_i - f_i)^2 - \frac{1}{2} \sum_{i,j=1}^m (f_i - m_i) k_{ij}^{-1} (f_j - m_j)} \\ & e^{-\sum_{j=1}^m (f^* - m^*) k_j^{*-1} (f_j - m_j)} \end{aligned}$$

- * Doing Gaussian integrals is a pain—I don't expect you to do this
- * But if you are brave enough to try let $g_i = f_i - m_i$, $g^* = f^* - m^*$ and $\hat{y}_i = y_i - m_i$ then we can rewrite this as

$$p(f^*|\mathcal{D}) \propto e^{-\frac{1}{2} k^{*-1} g^{*2}} \left(\prod_i \int dg_i \right) e^{-\frac{1}{2\sigma^2} \sum_{i=1}^m (\hat{y}_i - g_i)^2 - \frac{1}{2} \sum_{i,j=1}^m g_i k_{ij}^{-1} g_j - g^* \sum_{j=1}^m k_j^{*-1} g_j}$$

· where we made a change of variable in the integral

- * In matrix form we can write this as

$$p(f^*|\mathcal{D}) \propto e^{-\frac{1}{2} k^{*-1} g^{*2}} \left(\prod_i \int dg_i \right) e^{(\sigma^{-2} \hat{\mathbf{y}} - g^* \mathbf{k}^{*-1})^\top \mathbf{g} - \frac{1}{2} \mathbf{g}^\top (\mathbf{k}^{-1} - \sigma^{-2} \mathbf{I}) \mathbf{g}}$$

- I've dropped a term $e^{-\frac{\|\hat{\mathbf{y}}\|^2}{2\sigma^2}}$ since this is a constant that can be absorbed into the constant of proportionality
- * To perform the integrals over g_i we *complete the square*. To make the algebra easier to follow we define $\mathbf{M} = \mathbf{k}^{-1} - \sigma^{-2}\mathbf{I}$ and $\mathbf{b} = \sigma^{-2}\hat{\mathbf{y}} - g^*\mathbf{k}^{*-1}$. We now rewrite the exponent, E , of the integral as

$$\begin{aligned} E &= (\sigma^{-2}\hat{\mathbf{y}} - g^*\mathbf{k}^{*-1})^\top \mathbf{g} - \frac{1}{2}\mathbf{g}^\top (\mathbf{k}^{-1} - \sigma^{-2}\mathbf{I}) \mathbf{g} = \mathbf{b}^\top \mathbf{g} - \frac{1}{2}\mathbf{g}^\top \mathbf{M} \mathbf{g} \\ &= -\frac{1}{2}(\mathbf{g} - \mathbf{M}^{-1}\mathbf{b})^\top \mathbf{M} (\mathbf{g} - \mathbf{M}^{-1}\mathbf{b}) + \frac{1}{2}\mathbf{b}^\top \mathbf{M}^{-1}\mathbf{b} \end{aligned}$$

- You simply have to expand out all the terms to see these are the same
- * Defining $\mathbf{h} = \mathbf{g} - \mathbf{M}^{-1}\mathbf{b}$ and substituting this back into $p(f^*|\mathcal{D})$ we get

$$p(f^*|\mathcal{D}) \propto e^{-\frac{1}{2}k^{*-1}g^{*2} + \frac{1}{2}\mathbf{b}^\top \mathbf{M}^{-1}\mathbf{b}} \left(\prod_i \int d\mathbf{h}_i \right) e^{-\frac{1}{2}\mathbf{h}^\top \mathbf{M} \mathbf{h}}$$

- * The last integral is a Gaussian integral that integrates to a constant (and can be absorbed into the constant of proportionality)
- * We are left with

$$\begin{aligned} p(f^*|\mathcal{D}) &\propto e^{-\frac{1}{2}k^{*-1}g^{*2} + \frac{1}{2}\mathbf{b}^\top \mathbf{M}^{-1}\mathbf{b}} \\ &= e^{-\frac{1}{2}k^{*-1}g^{*2} + \frac{1}{2}(\sigma^{-2}\hat{\mathbf{y}} - g^*\mathbf{k}^{*-1})^\top \mathbf{M}^{-1}(\sigma^{-2}\hat{\mathbf{y}} - g^*\mathbf{k}^{*-1})} \\ &\propto e^{-\frac{1}{2}(k^{*-1} - \mathbf{k}^{*-1\top} \mathbf{M}^{-1} \mathbf{k}^{*-1})g^{*2} - \frac{g^*}{\sigma^2} \hat{\mathbf{y}}^\top \mathbf{M}^{-1} \mathbf{k}^{*-1}} \end{aligned}$$

- where we have used $\mathbf{b} = \sigma^{-2}\hat{\mathbf{y}} - g^*\mathbf{k}^{*-1}$
- we have also thrown away terms not involving g^* (as they are constants that can be absorbed into the constant of proportionality)
- * We can again *complete the square* this time for g^*

$$p(f^*|\mathcal{D}) \propto e^{-\frac{1}{2}(k^{*-1} - \mathbf{k}^{*-1\top} \mathbf{M}^{-1} \mathbf{k}^{*-1})(g^* - \frac{\hat{\mathbf{y}}^\top \mathbf{M}^{-1} \mathbf{k}^{*-1}}{\sigma^2(k^{*-1} - \mathbf{k}^{*-1\top} \mathbf{M}^{-1} \mathbf{k}^{*-1})})^2}$$

- We've again dropped terms that don't contain g^*
- * But $g^* = f^* - m^*$ so we can write this as

$$\begin{aligned} p(f^*|\mathcal{D}) &\propto e^{-\frac{1}{2}(k^{*-1} - \mathbf{k}^{*-1\top} \mathbf{M}^{-1} \mathbf{k}^{*-1})(f^* - m^* - \frac{\hat{\mathbf{y}}^\top \mathbf{M}^{-1} \mathbf{k}^{*-1}}{\sigma^2(k^{*-1} - \mathbf{k}^{*-1\top} \mathbf{M}^{-1} \mathbf{k}^{*-1})})^2} \\ &\propto \mathcal{N}\left(f^* \middle| m^* + \frac{\hat{\mathbf{y}}^\top \mathbf{M}^{-1} \mathbf{k}^{*-1}}{\sigma^2(k^{*-1} - \mathbf{k}^{*-1\top} \mathbf{M}^{-1} \mathbf{k}^{*-1})}, k^{*-1} - \mathbf{k}^{*-1\top} \mathbf{M}^{-1} \mathbf{k}^{*-1}\right) \end{aligned}$$

- where we use the fact that we end up with a term that has the form of a normal distribution
- because the posterior term is normalise, in fact it has to be exactly equal to this normal distribution
- * Using $\mathbf{M} = \mathbf{k}^{-1} - \sigma^{-2}\mathbf{I}$

$$p(f^*|\mathcal{D}) = \mathcal{N}\left(f^* \middle| m^* + \frac{\hat{\mathbf{y}}^\top (\mathbf{k}^{-1} - \sigma^{-2}\mathbf{I})^{-1} \mathbf{k}^{*-1}}{\sigma^2(k^{*-1} - \mathbf{k}^{*-1\top} (\mathbf{k}^{-1} - \sigma^{-2}\mathbf{I})^{-1} \mathbf{k}^{*-1})}, \frac{1}{k^{*-1} - \mathbf{k}^{*-1\top} (\mathbf{k}^{-1} - \sigma^{-2}\mathbf{I})^{-1} \mathbf{k}^{*-1}}\right)$$

- * In other words the expected mean value of $f^* = f(\mathbf{x}^*)$ is

$$\mathbb{E}[f(\mathbf{x}^*)] = m(\mathbf{x}^*) + \frac{\hat{\mathbf{y}}^\top (\mathbf{k}^{-1} - \sigma^{-2}\mathbf{I})^{-1} \mathbf{k}^{*-1}}{\sigma^2(k^{*-1} - \mathbf{k}^{*-1\top} (\mathbf{k}^{-1} - \sigma^{-2}\mathbf{I})^{-1} \mathbf{k}^{*-1})}$$

(note that usually $m(\mathbf{x}^*) = 0$) and the expected variance in the value is

$$\mathbb{V}\text{ar}[f(\mathbf{x}^*)] = \frac{1}{k^{*-1} - \mathbf{k}^{*-1\top} (\mathbf{k}^{-1} - \sigma^{-2}\mathbf{I})^{-1} \mathbf{k}^{*-1}}$$

- * The result is pretty horrible because it involves inverting a matrix \mathbf{K} with components $k(\mathbf{x}, \mathbf{y})$ evaluated at the set of points $\{\mathbf{x}_1, \mathbf{x}_2, \dots, \mathbf{x}_m, \mathbf{x}^*\}$
- * This is usually feasible but there is a simpler form that occurs because of identities involving inverse matrices
- * This is best obtained using a second way of deriving the results
- * I went to the pain of deriving the result this way because it is just a consequence of Bayes' rule
- In the second derivation we consider joint probability of the observations $\{y_i | i = 1, 2, \dots, m\}$ and $f^* = f(\mathbf{x}^*)$
 - Even this derivation is a pain and you are **not** expected to know it
 - Now we assume $y_i = f(\mathbf{x}_i) + \epsilon_i$ where $\epsilon_i \sim \mathcal{N}(0, \sigma^2)$ (that is our observations have a normally distributed error) then

$$\mathbb{E}[y_i y_j] = \mathbb{E}[(f(\mathbf{x}_i) + \epsilon_i)(f(\mathbf{x}_j) + \epsilon_j)] = k(\mathbf{x}_i, \mathbf{x}_j) + \sigma^2 \delta_{ij}$$

- * since the observation error, ϵ_i , is assumed independent of the function value, $f(\mathbf{x}_i)$
- * here we assume $m(\mathbf{x}) = 0$ just to make life simple

- Similarly

$$\mathbb{E}[y_i f^*] = \mathbb{E}[(f(\mathbf{x}_i) + \epsilon_i)f(\mathbf{x}^*)] = k(\mathbf{x}_i, \mathbf{x}^*)$$

and

$$\mathbb{E}[(f^*)^2] = k(\mathbf{x}^*, \mathbf{x}^*)$$

- Thus

$$p(f^*, \mathbf{y}) = \mathcal{N}\left(\begin{pmatrix} \mathbf{y} \\ f^* \end{pmatrix} \middle| \mathbf{0}, \begin{pmatrix} \mathbf{K} + \sigma^2 \mathbf{I} & \mathbf{k}^* \\ \mathbf{k}^{*\top} & k^* \end{pmatrix}\right)$$

- * where \mathbf{K} is a matrix with components $k(\mathbf{x}_i, \mathbf{x}_j)$, \mathbf{k}^* is a vector with components $k(\mathbf{x}_i, \mathbf{x}^*)$ and $k^* = k(\mathbf{x}^*, \mathbf{x}^*)$

- We now use an identity involving matrix inverses

$$\begin{pmatrix} \mathbf{K} + \sigma^2 \mathbf{I} & \mathbf{k}^* \\ \mathbf{k}^{*\top} & k^* \end{pmatrix}^{-1} = \begin{pmatrix} \left(\mathbf{K} + \sigma^2 \mathbf{I} - \frac{\mathbf{k}^{*\top} \mathbf{k}^*}{k^*}\right)^{-1} & -\frac{(\mathbf{K} + \sigma^2 \mathbf{I})^{-1} \mathbf{k}^*}{k^* - \mathbf{k}^{*\top} (\mathbf{K} + \sigma^2 \mathbf{I})^{-1} \mathbf{k}^*} \\ -\frac{\mathbf{k}^{*\top} (\mathbf{K} + \sigma^2 \mathbf{I})^{-1}}{k^* - \mathbf{k}^{*\top} (\mathbf{K} + \sigma^2 \mathbf{I})^{-1} \mathbf{k}^*} & \frac{1}{k^* - \mathbf{k}^{*\top} (\mathbf{K} + \sigma^2 \mathbf{I})^{-1} \mathbf{k}^*} \end{pmatrix}$$

- * these identities involving inverses of matrices seem to come from nowhere—they make working with normal distributions a real pain
- * you can prove the identity by multiply the right-hand side by $\begin{pmatrix} \mathbf{K} + \sigma^2 \mathbf{I} & \mathbf{k}^* \\ \mathbf{k}^{*\top} & k^* \end{pmatrix}$ and showing you get the identity
- * if you are tempted to do this then set $\sigma^2 = 0$ and rename \mathbf{k}^* and k^* so its easier to do the algebra—it is a tedious calculation (see exercises)

- Now we use $p(f^* | \mathcal{D}) = p(f^* | \mathbf{y}) = p(f^*, \mathbf{y}) / p(\mathbf{y})$
- All you need to do is collect the terms in $p(f^*, \mathbf{y})$ involving f^*

$$p(f^* | \mathcal{D}) \propto e^{-\frac{f^{*2}}{2(k^* - \mathbf{k}^{*\top} (\mathbf{K} + \sigma^2 \mathbf{I})^{-1} \mathbf{k}^*)}} + \frac{f^* \mathbf{k}^{*\top} (\mathbf{K} + \sigma^2 \mathbf{I})^{-1} \mathbf{y}}{k^* - \mathbf{k}^{*\top} (\mathbf{K} + \sigma^2 \mathbf{I})^{-1} \mathbf{k}^*}$$

- *Completing the square*

$$p(f^* | \mathcal{D}) \propto e^{-\frac{1}{2(k^* - \mathbf{k}^{*\top} (\mathbf{K} + \sigma^2 \mathbf{I})^{-1} \mathbf{k}^*)}} (f^* - \mathbf{k}^{*\top} (\mathbf{K} + \sigma^2 \mathbf{I})^{-1} \mathbf{y})^2$$

- But as $p(f^* | \mathcal{D})$ must be normalised

$$p(f(\mathbf{x}^*) | \mathcal{D}) = \mathcal{N}\left(f(\mathbf{x}^*) \middle| \mathbf{k}^{*\top} (\mathbf{K} + \sigma^2 \mathbf{I})^{-1} \mathbf{y}, k^* - \mathbf{k}^{*\top} (\mathbf{K} + \sigma^2 \mathbf{I})^{-1} \mathbf{k}^*\right)$$

- That is the mean value of f at a point \mathbf{x}^* is

$$\mathbb{E}[f(\mathbf{x}^*)] = m(\mathbf{x}^*) + \mathbf{k}^{*\top}(\mathbf{K} + \sigma^2 \mathbf{I})^{-1} \mathbf{y}$$

* I've re-instated $m(\mathbf{x}^*)$ (even though it is usually taken to be zero)

- The variance is given by

$$\text{Var}[f(\mathbf{x}^*)] = k^* - \mathbf{k}^{*\top}(\mathbf{K} + \sigma^2 \mathbf{I})^{-1} \mathbf{k}^*$$

- This may look completely different from the previous expression but surprisingly it is identical (you can demonstrate this numerically)
- The advantage over the previous methods is that we only need to do one matrix inversion (and we get a slightly easier expression)
- Although the mathematics to get here was horrible the expression is extremely easy to compute
- One of the advantages of Bayesian inference is it provides an estimate of its own uncertainty (in this case $\text{Var}[f(\mathbf{x}^*)]$)

2.4 Learning Hyperparameters

• Hyperparameters in the Bayesian Framework

- As I have said many times to get machine learning to work it is vital to choose the hyperparameter properly (we will see this is true for GPs)
- In the Bayesian framework if ϕ are hyperparameters then Bayes' rules says

$$p(\mathbf{x}|\mathcal{D}, \phi) = \frac{p(\mathcal{D}|\mathbf{x}, \phi) p(\mathbf{x}|\phi)}{p(\mathcal{D}|\phi)}$$

* this is just Bayes' rule with everything conditioned on the hyperparameters ϕ

- We can select hyperparameters by considering the **evidence**, $p(\mathcal{D}|\phi)$

* when the hyperparameters are viewed as different models this is known as **model selection**

* this is also called the evidence framework

* We can view the relative likelihood of one model, ϕ_1 compare to a second, ϕ_2 by examining

$$\frac{p(\mathcal{D}|\phi_1)}{p(\mathcal{D}|\phi_2)}$$

- If we want to be hyper-Bayesian then we can put a prior, $p(\phi)$, over our hyperparameters and then calculate the joint posterior for the parameters of the likelihood θ and the hyper-parameters ϕ

$$p(\theta, \phi|\mathcal{D}) = \frac{p(\mathcal{D}|\theta, \phi) p(\theta|\phi) p(\phi)}{p(\mathcal{D})}$$

* to compute the posterior we are interested in we marginalise out the hyperparameters

$$p(\theta|\mathcal{D}) = \int p(\theta, \phi|\mathcal{D}) d\phi$$

* Often this integral is not computable in closed form and we are forced to estimate it using Monte-Carlo techniques

• Hyperparameter for Gaussian Processes

- For Gaussian Processes this means choosing the right kernel function

- Here we are in a better position than in SVMs in that the kernel represents the covariance function
- Given data we could estimate the kernel function using

$$\mathbb{E}[y_i, y_j] = k(\mathbf{x}_i, \mathbf{x}_j) + \sigma^2 \delta_{ij}$$

- It would be an inverse problem to estimate $k(\mathbf{x}_i, \mathbf{x}_j)$
- Often we start with a guess of the form of the kernel
- A very common kernel function is the Gaussian kernel

$$k(\mathbf{x}, \mathbf{y}) = e^{-\frac{\|\mathbf{x}-\mathbf{y}\|^2}{2\ell^2}}$$

- Now ℓ is a *scale parameter* that has to be determined
- However, a great advantage of Gaussian Processes is that we can compute the evidence in closed form

$$\log(p(\mathcal{D}|\boldsymbol{\theta})) = -\frac{1}{2}\mathbf{y}^\top(\mathbf{K} + \sigma^2\mathbf{I})^{-1}\mathbf{y} - \frac{1}{2}\log(|\mathbf{K} + \sigma^2\mathbf{I}|) - \frac{m}{2}\log(2\pi)$$

- * This allows us to rapidly compare different hyperparameters (kernels and noise level)
- * We are not doing this in a fully Bayesian way (we are maximising our evidence rather than updating a distribution for the hyperparameters) therefore we could overfit
- * However, usually we have relatively few hyperparameters so overfitting is less severe

3 Exercises

3.1 Working with Gaussians

This is optional. Learning to work with Gaussians allows you to understand a lot of machine learning techniques, but it is not going to be examined.

1. Integrals involving Gaussians can be done in closed form. This means that they are heavily used in machine learning. It does, however, take a lot of practice to learn how to do these integrals. The starting point is the integral

$$I_1 = \int_{-\infty}^{\infty} e^{-x^2/2} dx = \sqrt{2\pi}.$$

This integral isn't easy. The trick is to turn this into a two dimensional integral and perform the integral in polar coordinates. Have a go.

2. Show by a change of variables

$$I_2 = \int_{-\infty}^{\infty} e^{-x^2/(2\sigma^2)} dx = \sqrt{2\pi}\sigma.$$

3. Show that

$$I_3 = \int_{-\infty}^{\infty} e^{-x^2/2+ax} \frac{dx}{\sqrt{2\pi}} = e^{a^2/2}$$

This involves completing the square and using a change of variables. It is probably the most frequently used trick working with Gaussian integrals.

4. Show that

$$I_4 = \int_{-\infty}^{\infty} \cdots \int_{-\infty}^{\infty} e^{-\frac{1}{2}\|\mathbf{x}\|^2} d\mathbf{x} = (2\pi)^{\frac{n}{2}}$$

Often we don't bother writing all the integral signs (one will do)

5. By writing $\mathbf{M} = \mathbf{V}\mathbf{\Lambda}\mathbf{V}^T$ show that

$$I_5 = \int_{-\infty}^{\infty} e^{-\frac{1}{2}\mathbf{x}^T\mathbf{M}\mathbf{x}} d\mathbf{x} = (2\pi)^{\frac{n}{2}} \frac{1}{\sqrt{|\mathbf{M}|}}$$

where $|\mathbf{M}|$ is the determinant of \mathbf{M} . Note that $|\mathbf{M}| = |\mathbf{\Lambda}| = \prod_i \lambda_i$.

6. Show that

$$I_6 = \int_{-\infty}^{\infty} e^{-\frac{1}{2}\mathbf{x}^T\mathbf{M}^{-1}\mathbf{x}} d\mathbf{x} = (2\pi)^{\frac{n}{2}} \sqrt{|\mathbf{M}|}$$

3.2 Working with Matrix Inverses

- Modelling using normal (aka Gaussian) distributions is very powerful because we can usually computer everything in closed form—integrals involving Gaussians are always doable
- However, there is a twist which makes it extremely painful
 - The normal distribution uses the inverse of the covariance matrix
 - The inverse of matrices are a pain to work with
- To get nice results we often consider partition matrices into four blocks—the inverse of the full matrix can be written in terms of the inverse of each block
- Let's look at a very simple example
 - Show that

$$\begin{pmatrix} \mathbf{K} & \ell \\ \ell^T & m \end{pmatrix}^{-1} = \begin{pmatrix} \left(\mathbf{K} - \frac{\ell\ell^T}{m}\right)^{-1} & -\frac{\mathbf{K}^{-1}\ell}{m - \ell^T\mathbf{K}^{-1}\ell} \\ -\frac{\ell^T\mathbf{K}^{-1}}{m - \ell^T\mathbf{K}^{-1}\ell} & \frac{1}{m - \ell^T\mathbf{K}^{-1}\ell} \end{pmatrix}$$

* \mathbf{K} is a symmetric matrix, ℓ is a vector and m a scalar

* Similar identities exist when ℓ and m are matrices, but let's keep thing simple

- Now we can show this by multiply both sides by $\begin{pmatrix} \mathbf{K} & \ell \\ \ell^T & m \end{pmatrix}$ (the algebra is easier if we multiply on the right)
- This is still very fiddly to show, so break it down
 - * Show that

$$\begin{pmatrix} \mathbf{A} & \mathbf{b} \\ \mathbf{b}^T & c \end{pmatrix} \begin{pmatrix} \mathbf{K} & \ell \\ \ell^T & m \end{pmatrix} = \begin{pmatrix} \mathbf{A}\mathbf{K} + \mathbf{b}\ell^T & \mathbf{A}\ell + m\mathbf{b} \\ \mathbf{b}^T\mathbf{K} + c\ell^T & \mathbf{b}^T\ell + cm \end{pmatrix}$$

* Now you have to show that when $\mathbf{A} = \left(\mathbf{K} - \frac{\ell\ell^T}{m}\right)^{-1}$, $\mathbf{b} = -\frac{\mathbf{K}^{-1}\ell}{m - \ell^T\mathbf{K}^{-1}\ell}$ and $c = \frac{1}{m - \ell^T\mathbf{K}^{-1}\ell}$ that

1. $\mathbf{A}\mathbf{K} + \mathbf{b}\ell^T = \mathbf{I}$
2. $\mathbf{A}\ell + m\mathbf{b} = \mathbf{0}$
3. $\mathbf{b}^T\mathbf{K} + c\ell^T = \mathbf{0}^T$
4. $\mathbf{b}^T\ell + cm = 1$

* To show 1. and 2. take out \mathbf{A} as a common factor using $\mathbf{I} = \mathbf{A}\mathbf{A}^{-1}$ (answer below)

4 Experiments

4.1 Generate Samples from a Gaussian Processes

- To generate a Gaussian Process, $\mathcal{GP}(0, k)$, with $m(\mathbf{x}) = 0$ and covariance

$$k(\mathbf{x}, \mathbf{y}) = e^{-\frac{\|\mathbf{x}-\mathbf{y}\|^2}{2\ell^2}}$$

a points $\{\mathbf{x}_i | i = 1, 2, \dots, n\}$

1. We use first compute the covariance matrix \mathbf{K} with elements $K_{ij} = k(\mathbf{x}_i, \mathbf{x}_j)$
2. Compute the *Cholesky decomposition* \mathbf{L} such that $\mathbf{K} = \mathbf{L}\mathbf{L}^\top$
 - For any positive definite matrices we can always compute a Cholesky decomposition (although beware, due to numerical rounding some matrices that should be positive definite, appear not to have Cholesky decompositions—you can $\epsilon \mathbf{I}$ to your matrix to help)
 - The Cholesky decomposition is a lower diagonal matrix
 - It is used to efficiently solve linear problems involving positive definite matrices (it is more efficient and more stable than Gaussian elimination)
3. Now let $\mathbf{f} = \mathbf{L}\boldsymbol{\eta}$ where $\boldsymbol{\eta} \sim \mathcal{N}(\mathbf{0}, \mathbf{I})$
 - We note that

$$\mathbb{E}[\mathbf{f}\mathbf{f}^\top] = \mathbb{E}[\mathbf{L}\boldsymbol{\eta}\boldsymbol{\eta}^\top\mathbf{L}^\top] = \mathbf{L}\mathbb{E}[\boldsymbol{\eta}\boldsymbol{\eta}^\top]\mathbf{L}^\top = \mathbf{L}\mathbf{I}\mathbf{L}^\top = \mathbf{L}\mathbf{L}^\top = \mathbf{K}$$

- Thus \mathbf{f} can be viewed as samples from a Gaussian Process
 - Draw many different samples and plot them
4. Experiment with different values of ℓ

1-d Gaussian Processes

```
x = [-10:0.2:10];    % define some points
n = length(x)
l = 1.0               % define length scale
K = zeros(n,n);       % define holder
for i = 1:n
    K(i,i) = 1;
    for j = 1:i
        K(j,i) = K(i,j) = exp(-(x(i)-x(j))^2/(2*l)); % define covariance matrix
    endfor
endfor
L = chol(K+0.0001*eye(n)); % cheat to persuade octave K is positive definite
f1 = L*randn(n,1);
f2 = L*randn(n,1);
f3 = L*randn(n,1);
plot(x,f1,x,f2,x,f3)
```

2-d Gaussian Process

```
range = [-4:0.2:4]; % define some points
n = length(range)
X = zeros(n*n,2);
for i = 1:n
    for j = 1:n
        X((i-1)*n+j,1) = range(i);
```

```

X((i-1)*n+j,2) = range(j);
endfor
endfor
l = 1.0 % define length scale
N = length(X)
K = zeros(N,N);
for i = 1:N
    K(i,i) = 1;
    for j = 1:i
        K(j,i) = K(i,j) = exp(-norm(X(i,:)-X(j,:))^2/(2*l)); % define covariance matrix
    endfor
endfor

L = chol(K+0.0001*eye(N)); % cheat to persuade octave K is positive definite
f = L*randn(N,1);
fm = reshape(f,n,n);
[xx, yy] = meshgrid(range, range);
mesh(xx, yy, fm);
xlabel ("x");
ylabel ("y");
zlabel ("f(x,y)");
title ("2-d Gaussian Process");

```

5 Answers

5.1 Working with Gaussians

1. The Gaussian Integral.

- We consider the two dimensional integral

$$\begin{aligned}
 I_1^2 &= \left(\int_{-\infty}^{\infty} e^{-x^2/2} dx \right)^2 \stackrel{(a)}{=} \int_{-\infty}^{\infty} \int_{-\infty}^{\infty} e^{-x^2/2-y^2/2} dx dy \\
 &\stackrel{(b)}{=} \int_0^{\infty} \int_0^{2\pi} r e^{-r^2} dr d\theta \stackrel{(c)}{=} \int_0^{\infty} r e^{-r^2} dr \int_0^{2\pi} d\theta = 2\pi
 \end{aligned}$$

(a) Changing the name of the dummy index from x to y

(b) Using the change of variable $x = r \cos(\theta)$ and $y = r \sin(\theta)$ so that $dx dy = r dr d\theta$

(c) Reordering the integrals and using $\int_0^{\infty} r e^{-r^2} dr = 1$ (which you can prove by a change of variables $u = r^2/2$ so that $r dr = du$)

Therefore $I_1 = \sqrt{2\pi}$.

2.

$$I_2 = \int_{-\infty}^{\infty} e^{-x^2/(2\sigma^2)} dx \stackrel{(a)}{=} \sigma \int_{-\infty}^{\infty} e^{-u^2/2} du \stackrel{(b)}{=} \sqrt{2\pi} \sigma.$$

(a) Using $u = x/\sigma$ so that $dx = \sigma du$

(b) Using I_1

3.

$$\begin{aligned}
 I_2 &= \int_{-\infty}^{\infty} e^{-x^2/2+ax} \frac{dx}{\sqrt{2\pi}} \stackrel{(a)}{=} \int_{-\infty}^{\infty} e^{-\frac{1}{2}(x-a)^2 + \frac{a^2}{2}} \frac{dx}{\sqrt{2\pi}} \\
 &\stackrel{(b)}{=} e^{a^2/2} \int_{-\infty}^{\infty} e^{-\frac{1}{2}(x-a)^2} \frac{dx}{\sqrt{2\pi}} \stackrel{(c)}{=} e^{a^2/2}
 \end{aligned}$$

- (a) Using $-x^2/2 + ax = -\frac{1}{2}(x-a)^2 + \frac{a^2}{2}$
 (b) Factoring out $e^{a^2/2}$
 (c) Making a change of variables $u = x - a$ and using integral I_1

4.

$$\begin{aligned} I_4 &= \int_{-\infty}^{\infty} \cdots \int_{-\infty}^{\infty} e^{-\frac{1}{2}\|\mathbf{x}\|^2} d\mathbf{x} \stackrel{(a)}{=} \int_{-\infty}^{\infty} \cdots \int_{-\infty}^{\infty} e^{-\frac{1}{2}\sum_{i=1}^n x_i^2} \prod_{i=1}^n dx_i \\ &\stackrel{(b)}{=} \prod_{i=1}^n \left(\int_{-\infty}^{\infty} e^{-\frac{x_i^2}{2}} dx_i \right) \stackrel{(c)}{=} (2\pi)^{\frac{n}{2}} \end{aligned}$$

- (a) Using $\|\mathbf{x}\|^2 = \sum_{i=1}^n x_i^2$
 (b) Using $\exp(\sum_i a_i) = \prod_i e^{a_i}$ and factoring out the different terms
 (c) Using integral I_1

5.

$$\begin{aligned} I_5 &= \int_{-\infty}^{\infty} e^{-\frac{1}{2}\mathbf{x}^T \mathbf{M} \mathbf{x}} d\mathbf{x} \stackrel{(a)}{=} \int_{-\infty}^{\infty} e^{-\frac{1}{2}\mathbf{x}^T \mathbf{V} \mathbf{\Lambda} \mathbf{V}^T \mathbf{x}} d\mathbf{x} \\ &\stackrel{(b)}{=} \int_{-\infty}^{\infty} e^{-\frac{1}{2}\mathbf{u}^T \mathbf{\Lambda} \mathbf{u}} d\mathbf{u} \stackrel{(c)}{=} \int_{-\infty}^{\infty} e^{-\frac{1}{2}\sum_i \lambda_i u_i^2} d\mathbf{u} \\ &\stackrel{(d)}{=} \prod_i \left(\int_{-\infty}^{\infty} e^{-\frac{1}{2}\lambda_i u_i^2} du_i \right) \stackrel{(e)}{=} \prod_i \left(\sqrt{\frac{2\pi}{\lambda_i}} \right) \stackrel{(f)}{=} (2\pi)^{n/2} \frac{1}{\sqrt{\prod_i \lambda_i}} \stackrel{(g)}{=} \frac{(2\pi)^{\frac{n}{2}}}{\sqrt{|\mathbf{M}|}} \end{aligned}$$

- (a) Using $\mathbf{M} = \mathbf{V} \mathbf{\Lambda} \mathbf{V}^T$ where
 (b) Writing $\mathbf{u} = \mathbf{V}^T \mathbf{x}$. Note that $d\mathbf{u} = |\mathbf{V}^T| d\mathbf{x}$, but \mathbf{V} is an orthogonal matrix so $|\mathbf{V}^T| = 1$. So $d\mathbf{u} = d\mathbf{x}$.
 (c) $\mathbf{\Lambda}$ is a diagonal matrix with elements $\Lambda_{ii} = \lambda_i$
 (d) Rearranging
 (e) Using integral I_2 with $\sigma = 1/\sqrt{\lambda_i}$
 (f) Rearranging
 (g) Using $\prod_i \lambda_i = |\mathbf{M}|$.

6. This follows using I_5 because $|\mathbf{M}^{-1}| = \frac{1}{|\mathbf{M}|}$.

5.2 Working with Matrix Inverses

- Using $\mathbf{A} = \left(\mathbf{K} - \frac{\ell \ell^T}{m} \right)^{-1}$, $\mathbf{b} = -\frac{\mathbf{K}^{-1} \ell}{m - \ell^T \mathbf{K}^{-1} \ell}$ and $c = \frac{1}{m - \ell^T \mathbf{K}^{-1} \ell}$, we show the four identities

1.

$$\begin{aligned} \mathbf{A} \mathbf{K} + \mathbf{b} \ell^T &= \mathbf{A} \left(\mathbf{K} + \mathbf{A}^{-1} \mathbf{b} \ell^T \right) \\ &= \left(\mathbf{K} - \frac{\ell \ell^T}{m} \right)^{-1} \left(\mathbf{K} + \left(\mathbf{K} - \frac{\ell \ell^T}{m} \right) \frac{(-\mathbf{K}^{-1} \ell \ell^T)}{m - \ell^T \mathbf{K}^{-1} \ell} \right) \\ &= \left(\mathbf{K} - \frac{\ell \ell^T}{m} \right)^{-1} \left(\mathbf{K} - \frac{\ell \ell^T}{m - \ell^T \mathbf{K}^{-1} \ell} + \frac{\ell \ell^T \mathbf{K}^{-1} \ell \ell^T / m}{m - \ell^T \mathbf{K}^{-1} \ell} \right) \\ &= \left(\mathbf{K} - \frac{\ell \ell^T}{m} \right)^{-1} \left(\mathbf{K} - \frac{\ell \ell^T}{m} \frac{m - \ell^T \mathbf{K}^{-1} \ell}{m - \ell^T \mathbf{K}^{-1} \ell} \right) \\ &= \left(\mathbf{K} - \frac{\ell \ell^T}{m} \right)^{-1} \left(\mathbf{K} - \frac{\ell \ell^T}{m} \right) = \mathbf{I} \end{aligned}$$

2.

$$\begin{aligned}
\mathbf{A}\ell + m\mathbf{b} &= \mathbf{A}(\ell + m\mathbf{A}^{-1}\mathbf{b}) \\
&= \mathbf{A}\left(\ell - m\left(\mathbf{K} - \frac{\ell\ell^\top}{m}\right)\frac{\mathbf{K}^{-1}\ell}{m - \ell^\top\mathbf{K}^{-1}\ell}\right) \\
&= \mathbf{A}\left(\ell - \frac{m\ell}{m - \ell^\top\mathbf{K}^{-1}\ell}\left(1 - \frac{\ell^\top\mathbf{K}^{-1}\ell}{m}\right)\right) \\
&= \mathbf{A}(\ell - \ell) = \mathbf{0}
\end{aligned}$$

3.

$$\mathbf{b}^\top\mathbf{K} + c\ell^\top = -\frac{\ell^\top}{m - \ell^\top\mathbf{K}^{-1}\ell} + \frac{\ell^\top}{m - \ell^\top\mathbf{K}^{-1}\ell} = \mathbf{0}^\top$$

4.

$$\mathbf{b}^\top\ell + cm = -\frac{\ell^\top\mathbf{K}^{-1}\ell}{m - \ell^\top\mathbf{K}^{-1}\ell} + \frac{m}{m - \ell^\top\mathbf{K}^{-1}\ell} = 1$$