

HOW TO ANNOY A STATISTICIAN

Announcements

Course schedule updated.

In person lectures - PSYC G8

Week 8 Wed - Gaussian Process 2

Week 10 Wed - Guest lecture: ML for Cyber security

- No tutorial in week 7 – assignment 1 walk-through during drop-in, will be recorded.
 - No tutorial in week 8 – enjoy quiz 2 :)
-
- Video assignment released, due in week 12.
 - Assignment 2 out soon.
 - Quiz 2 released, 9 multiple choice questions.
 - Assignment 1 raw grades out, regrade-request open on Gradescope for 2 weeks.

Thank You for anonymous survey input

IML → SML

Student A: Math too hard!

Student B: “Content is too easy. Does not feel challenging enough to be a 4000 level course.
Make assessment more difficult, and material more in depth.”

? skipped derivations / proofs

Thanks. Let's clarify some expectations for derivations / proofs.

Things that are important and appearing for the first time are generally covered; when a proof is skipped it usually means that we expect students to be able to follow the overall logic, and the model design arguments without the proof.

? motivating the material

Motivations tend to be mathematical in this course. We're happy to help in a number of ways.

? which parts are important

Good question. We'll try to clarify + reemphasize this in the intro/outro slides. If still unclear, do ask questions during class Q&A, during tuts, or on piazza.

Gaussian Processes - Regression

Motivation

Defining Gaussian Processes (GP)

Kernel functions, sampling

GP regression

GP regression - predictive distribution

Sampling algorithm and computational costs

in assignment 2

Bishop, Chap 6.4 (6.4.1-6.4.3)

GP book

<http://gaussianprocess.org/gpml/chapters/> Chap 1, 2.1, 2.2, 2.5

Bayesian Joint Inversions for the Exploration of Earth Resources

Alistair Reid¹, Simon O'Callaghan¹, Edwin V. Bonilla¹, Lachlan McCalman¹, Tim Rawling² and Fabio Ramos³

1. NICTA, 2. University of Melbourne, 3. University of Sydney

{ alistair.reid, simon.ocallaghan, edwin.bonilla, lachlan.mccalman }@nicta.com.au

tim.rawling@unimelb.edu.au, f.ramos@acfr.usyd.edu.au

In a geological inversion problem, properties such as temperature, conductivity, density, magnetic susceptibility and permeability are inferred from related observations such as gravity, magnetics and seismic reflexion. ... In this paper we formulate geophysical inversion as a machine learning problem, and propose an approach based on *Gaussian processes regression* that naturally provides both a predictive distribution over the inverted quantities and a principled method to fuse different types of observations. We apply our method to a real dataset from South Australia containing gravity and drill-hole data with the goal of characterizing rock densities for geothermal target exploration, and also to simulated validation data involving gravity, drill-hole and magnetic observations.

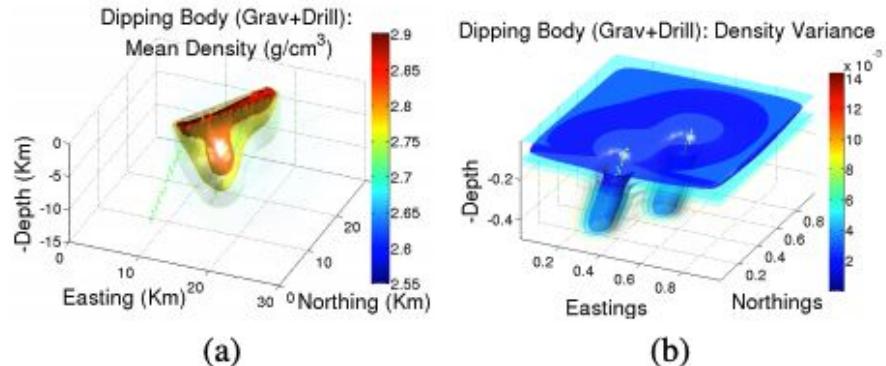


Figure 5: Outputs of the Gaussian process inversion algorithm after fusing gravity and drill observations of the dipping body.

Two ways to learn a function (regression)

[GPbook, intro]

Motivating scenarios

- SARCOS robotic arm
 - x : (7 joint positions, 7 joint velocities, 7 joint accelerations)
 - y : 7 joint torques
 - Learning: $y \sim f(x)$
 - Prediction: compute the torque needed to move the arm along a given trajectory
- Predict the yield of a chemical plant (y), given temperature, pressure, amount of catalyst, etc (x)
 - Hmm, we don't know a functional form of y (in either input or kernel space)
 - With full uncertainty estimate in y

Representation choices

- Restrict the class of functions we consider, e.g. linear
- Give prior probability to every possible function

we know how to do this :)

huh?

Un-countably infinite number of functions

Un-countably infinite number of input points (e.g. \mathbb{R} , \mathbb{R}^2 , ...)

Rely on parameter of the function, which can be a function itself

"Pretend" that it's a very long vector

Prediction function for regression

$$\mathbf{w}^* = (\lambda \mathbf{I} + \Phi^\top \Phi)^{-1} \Phi^\top \mathbf{t}$$

← normal eq.

Prediction function for new \mathbf{x}

$$y(\mathbf{x}) = \phi(\mathbf{x})^\top \mathbf{w}^* = \phi(\mathbf{x})^\top (\lambda \mathbf{I} + \Phi^\top \Phi)^{-1} \Phi^\top \mathbf{t}$$

Kernelized version

$$\begin{aligned} \mathbf{w} &= \Phi^\top \mathbf{a} \\ \mathbf{a} &= (\mathbf{K} + \lambda \mathbf{I}_N)^{-1} \mathbf{t}. \end{aligned} \tag{6.8}$$

$$y(\mathbf{x}) = \mathbf{w}^\top \phi(\mathbf{x}) = \mathbf{a}^\top \Phi \phi(\mathbf{x}) = \mathbf{k}(\mathbf{x})^\top (\mathbf{K} + \lambda \mathbf{I}_N)^{-1} \mathbf{t} \tag{6.9}$$

What is the Bayesian version of this?

point estimate of y
(Bayesian) view desires
 $P(y_{n+1} | X_{n+1}, X_{1:N}, Y_{1:N})$

→ expressed in \mathbf{K}

Bayesian linear regression (isotropic prior)

if

$$p(\mathbf{t}|\mathbf{X}, \mathbf{w}, \beta) = \prod_{n=1}^N \mathcal{N}(t_n | \mathbf{w}^T \phi(\mathbf{x}_n), \underline{\beta^{-1}}) \quad (3.10)$$

$$p(\mathbf{w}|\alpha) = \mathcal{N}(\mathbf{w}|\mathbf{0}, \alpha^{-1}\mathbf{I}) \quad (3.52)$$

then

$$\underline{p(\mathbf{w}|\mathbf{t})} = \mathcal{N}(\mathbf{w}|\mathbf{m}_N, \mathbf{S}_N) \quad (3.49)$$

$$\begin{cases} \mathbf{m}_N &= \beta \mathbf{S}_N \Phi^T \mathbf{t} \\ \mathbf{S}_N^{-1} &= \alpha \mathbf{I} + \beta \Phi^T \Phi. \end{cases} \quad (3.53) \quad (3.54)$$

$$\begin{cases} \mathbf{m}_N &= \beta \mathbf{S}_N \Phi^T \mathbf{t} \\ \mathbf{S}_N^{-1} &= \alpha \mathbf{I} + \beta \Phi^T \Phi. \end{cases} \quad (3.53) \quad (3.54)$$

Prediction function
for new \mathbf{x}

$$p(t|\mathbf{x}, \mathbf{t}, \alpha, \beta) = \mathcal{N}(t|\mathbf{m}_N^T \phi(\mathbf{x}), \sigma_N^2(\mathbf{x})) \quad (3.58)$$

$$\sigma_N^2(\mathbf{x}) = \frac{1}{\beta} + \phi(\mathbf{x})^T \mathbf{S}_N \phi(\mathbf{x}). \quad (3.59)$$

Is there a kernelized version of this? – can I have both kernels and prior

Gp book
2.1

"weight space"

regression prob;

learning : weights
etc,

prediction?
 (x_{N+1}, t_{N+1})
?

$p(t_{N+1} | \dots)$

Motivation

In week 6, we talked about kernels in non-probabilistic settings:

test data ,
↓

$$y(\mathbf{x}) = \mathbf{k}(\mathbf{x})^\top (\mathbf{K} + \lambda \mathbf{I}_N)^{-1} \mathbf{t}$$

↓
 $\phi\phi^\top$ *training ,*

$$y(\mathbf{x}) = \mathbf{w}^\top \phi(\mathbf{x}) + b = \sum_{n=1}^N a_n t_n k(\mathbf{x}, \mathbf{x}_n) + b$$

Can one use kernels in a probabilistic setting?

“noiseless”

Starting again from linear regression ...

fixed

$$y(\mathbf{x}) = \mathbf{w}^T \phi(\mathbf{x}) \quad (6.49)$$

$$p(\mathbf{w}) = \mathcal{N}(\mathbf{w} | \mathbf{0}, \alpha^{-1} \mathbf{I}) \quad (6.50)$$

matrix form

fixed

Gaussian

$$\mathbf{y} = \Phi \mathbf{w} \quad (6.51)$$

$$\Phi_{nk} = \phi_k(\mathbf{x}_n).$$

Claim: \mathbf{Y} is Gaussian ... only its mean and covariance matters.

$$\mathbb{E}[\mathbf{y}] = \Phi \mathbb{E}[\mathbf{w}] = \mathbf{0} \quad \nearrow \alpha^{-1} \mathbf{I} \quad (6.52)$$

$$\text{cov}[\mathbf{y}] = \mathbb{E}[\mathbf{y}\mathbf{y}^T] = \Phi \mathbb{E}[\mathbf{w}\mathbf{w}^T] \Phi^T = \frac{1}{\alpha} \Phi \Phi^T = \mathbf{K} \quad (6.53) \quad p(\mathbf{y}) = \mathcal{N}(\mathbf{y} | \mathbf{0}, \mathbf{K})$$

$$K_{nm} = k(\mathbf{x}_n, \mathbf{x}_m) = \frac{1}{\alpha} \phi(\mathbf{x}_n)^T \phi(\mathbf{x}_m) \quad (6.54)$$

$$\mathbf{y} \sim \mathcal{N}(\mathbf{0}, \frac{1}{\alpha} \Phi \Phi^T)$$

Stochastic proc. $p(y(t) | \theta)$

Gaussian process (GP)

A Gaussian Process is defined as a probability distribution over functions $y(x)$ such that the set of values of $y(x)$ evaluated at an arbitrary set of points x_1, \dots, x_N jointly have a Gaussian distribution.

- More generally, a stochastic process $y(x)$ is specified by giving the joint probability distribution for any finite set of values $y(x_1), \dots, y(x_N)$ in a consistent manner.
- Common choice: set mean to zero, due to no prior knowledge of $y(x)$ $E[y] = 0$
- The joint distribution of any y_1, \dots, y_N Fully specified by their second-order statistics
- Input x in 2-D – Gaussian random field (also called Kriging in statistics)

$x \in \mathbb{D}$ – easy to visualise

$$y(\mathbf{x}) = \mathbf{w}^T \phi(\mathbf{x}) \quad (6.49)$$

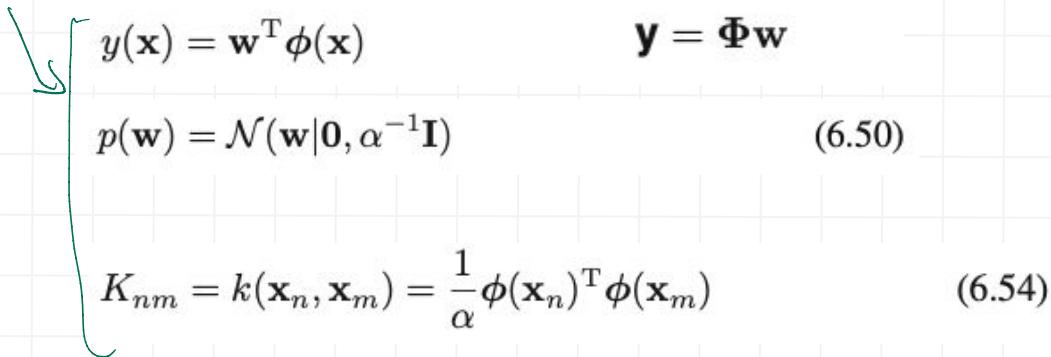
$$p(\mathbf{w}) = \mathcal{N}(\mathbf{w} | \mathbf{0}, \alpha^{-1} \mathbf{I}) \quad (6.50)$$

$$\mathbf{y} = \Phi \mathbf{w} \quad (6.51)$$

$$\vec{y} \sim N(\vec{\mu}, \vec{K})$$

Kernel functions in GP

$$p(\mathbf{y}) = \mathcal{N}(\mathbf{y} | \mathbf{0}, \mathbf{K})$$


$$\begin{aligned} y(\mathbf{x}) &= \mathbf{w}^T \phi(\mathbf{x}) \\ \mathbf{y} &= \Phi \mathbf{w} \\ p(\mathbf{w}) &= \mathcal{N}(\mathbf{w} | \mathbf{0}, \alpha^{-1} \mathbf{I}) \end{aligned} \tag{6.50}$$

Implicitly defined by $\phi(\mathbf{x})$

$$K_{nm} = k(\mathbf{x}_n, \mathbf{x}_m) = \frac{1}{\alpha} \phi(\mathbf{x}_n)^T \phi(\mathbf{x}_m) \tag{6.54}$$

One can also explicitly define a kernel

$$k(\mathbf{x}, \mathbf{x}') = \exp \left(-\|\mathbf{x} - \mathbf{x}'\|^2 / 2\sigma^2 \right) \tag{6.23}$$

$$k(x, x') = \exp(-\theta |x - x'|) \tag{6.56}$$

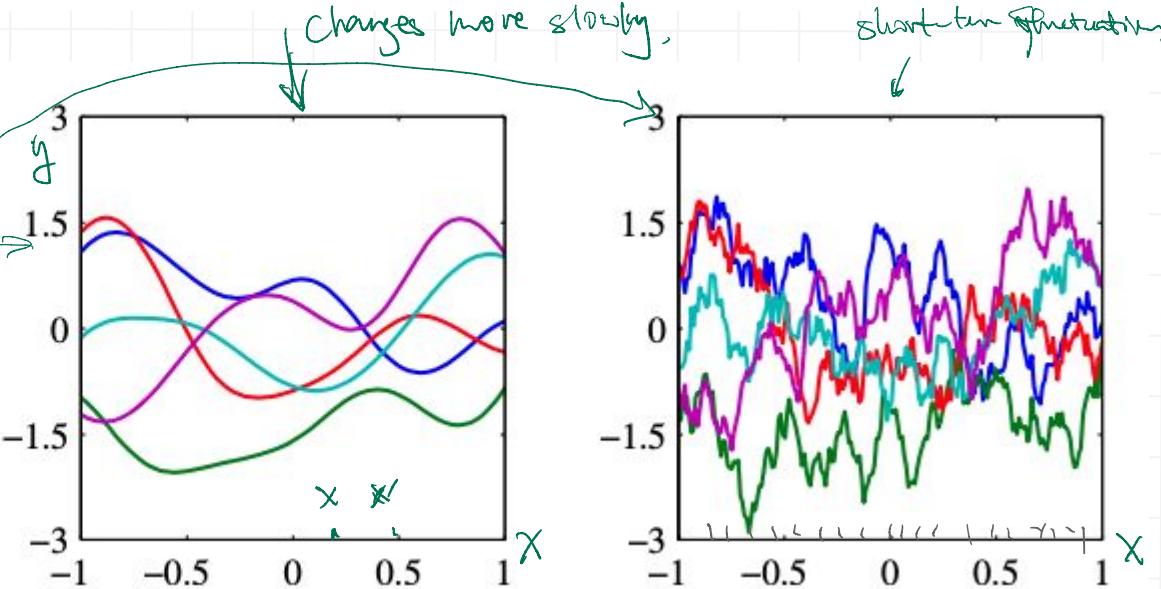
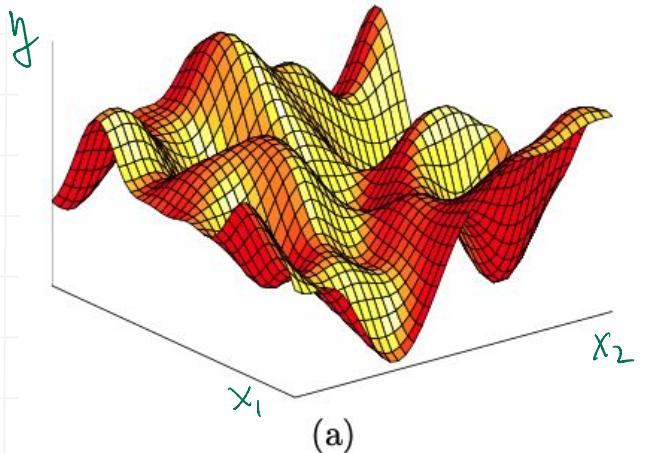
$k(\mathbf{x}, \mathbf{x}') = \mathbf{x}^\top \mathbf{x}'$	linear kernel
$k(\mathbf{x}, \mathbf{x}') = k(\mathbf{x} - \mathbf{x}')$	stationary kernel
$k(\mathbf{x}, \mathbf{x}') = k(\ \mathbf{x} - \mathbf{x}'\)$	homogeneous kernel

The whole heap of kernel construction tricks apply ...

Figure 6.4 Samples from Gaussian processes for a ‘Gaussian’ kernel (left) and an exponential kernel (right).

$$k(\mathbf{x}, \mathbf{x}') \quad \text{cov}(y(\mathbf{x}), y(\mathbf{x}'))$$

[GP book, Fig 1.2(a)]



$$k(\mathbf{x}, \mathbf{x}') = \exp(-\|\mathbf{x} - \mathbf{x}'\|^2 / 2\sigma^2) \quad (6.23)$$

$$k(x, x') = \exp(-\theta |x - x'|) \quad (6.56)$$

$$\tilde{y} \sim N(\bar{\mu}, K)$$

Q: How to sample from a GP?
How are these figures produced?

$$\{x_1, x_2, \dots, x_n\} \subset \mathbb{R}^{n \times n}$$

To generate samples $\mathbf{x} \sim \mathcal{N}(\mathbf{m}, K)$ with arbitrary mean \mathbf{m} and covariance matrix K using a scalar Gaussian generator (which is readily available in many programming environments) we proceed as follows: first, compute the Cholesky decomposition (also known as the “matrix square root”) L of the positive definite symmetric covariance matrix $K = LL^\top$, where L is a lower triangular matrix, see section A.4. Then generate $\mathbf{u} \sim \mathcal{N}(\mathbf{0}, I)$ by multiple separate calls to the scalar Gaussian generator. Compute $\mathbf{x} = \mathbf{m} + L\mathbf{u}$, which has the desired distribution with mean \mathbf{m} and covariance $L\mathbb{E}[\mathbf{u}\mathbf{u}^\top]L^\top = LL^\top = K$ (by the independence of the elements of \mathbf{u}).

In practice it may be necessary to add a small multiple of the identity matrix ϵI to the covariance matrix for numerical reasons. This is because the eigenvalues of the matrix K can decay very rapidly (see section 4.3.1 for a closely related analytical result) and without this stabilization the Cholesky decomposition fails. The effect on the generated samples is to add additional independent noise of variance ϵ . From the context ϵ can usually be chosen to have inconsequential effects on the samples, while ensuring numerical stability.

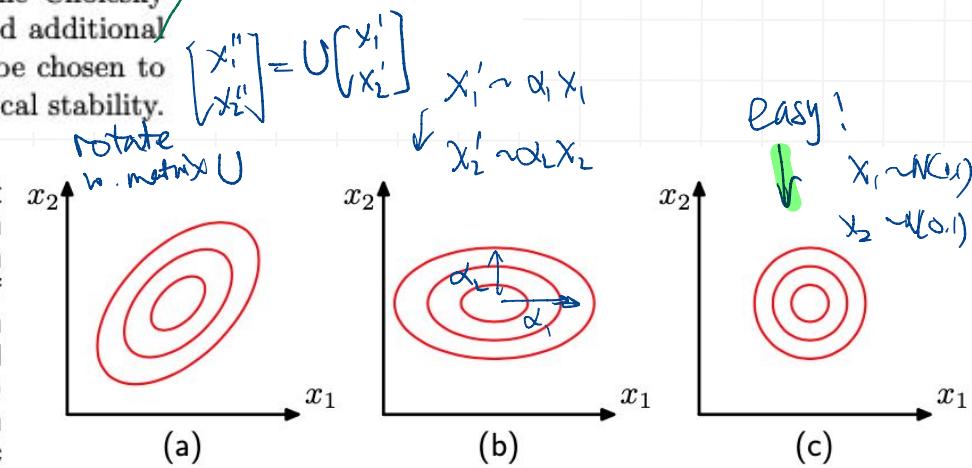
Intuition:

Figure 2.8 Contours of constant probability density for a Gaussian distribution in two dimensions in which the covariance matrix is (a) of general form, (b) diagonal, in which the elliptical contours are aligned with the coordinate axes, and (c) proportional to the identity matrix, in which the contours are concentric circles.

generating multivariate Gaussian samples

`numpy.random.standard_normal`

`numpy.random.multivariate_normal`



Why Cholesky rather than eigen decomposition?

faster to compute!

$O(n^3) \leftarrow$ Cholesky
eigen val inverse

Cost of some matrix factorizations and decompositions. A is $n \times n$, except for QR as $\sim VD$, where it is $m \times n$ ($m \geq n$).

Factorization/decomposition	Number of flops
LU factorization with partial pivoting ($PA = LU$)	$2n^3/3$
LU factorization with partial pivoting of upper Hessenberg matrix ($PA = LU$)	n^2
Cholesky factorization ($A = R^* R$)	$n^3/3$
Householder QR factorization ($A = QR$)	$2n^2(m - n/3)$ for R ; $4(m^2n - mn^2 + n^3/3)$ for $m \times m Q$; $2n^2(m - n/3)$ for $m \times n Q$;
SVD ^a ($A = P\Sigma Q^*$)	$2np(2m - n)$ for QB with $m \times p B$ and Q held in factored form.
Hessenberg decomposition ($A = QHQ^*$)	$14mn^2 + 8n^3$ ($P(:, 1:n)$, Σ , and Q) ^b
Schur decomposition ^a ($A = QTQ^*$)	$6mn^2 + 20n^3$ ($P(:, 1:n)$, Σ , and Q) ^c
For Hermitian A :	
Tridiagonal reduction ($A = QTQ^*$)	$14n^3/3$ (Q and H), $10n^3/3$ (H only)
Spectral decomposition ($A = QDQ^*$)	$25n^3$ (Q and T), $10n^3$ (T only)
	$8n^3/3$ (Q and T), $4n^3/3$ (T only)
	$9n^3$ (Q and D), $4n^3/3$ (D only)

^aThese costs are estimates taken from Golub and Van Loan [224, 1996].

^bGolub–Reinsch SVD.

^cGolub–Reinsch SVD with preliminary QR factorization.

Gaussian Processes - Regression

Motivation

Defining Gaussian Processes (GP)

Kernel functions, sampling

GP regression

GP regression - predictive distribution

Sampling algorithm and computational costs

Bishop, Chap 6.4 (6.4.1-6.4.3)

GP book

<http://gaussianprocess.org/gpml/chapters/> Chap 1, 2.1, 2.2, 2.5

GP for regression (w noise)

“noiseless”

$$\mathbf{y} \sim \mathcal{N}(\mathbf{0}, \frac{1}{\alpha} \Phi \Phi^T)$$

$$\sim \mathcal{N}(0, k)$$

$$y_n = y(\mathbf{x}_n)$$

$$t_n = y_n + \epsilon_n$$

$$\epsilon_n \sim \mathcal{N}(0, \beta^{-1}) \quad (6.57)$$

$$p(t_n|y_n) = \mathcal{N}(t_n|y_n, \beta^{-1}) \quad (6.58)$$

$$p(\mathbf{t}|\mathbf{y}) = \mathcal{N}(\mathbf{t}|\mathbf{y}, \beta^{-1} \mathbf{I}_N) \quad (6.59)$$

$$p(\mathbf{y}) = \mathcal{N}(\mathbf{y}|\mathbf{0}, \mathbf{K}). \quad (6.60)$$

$$\mathbf{C} = \begin{bmatrix} k_{11} + \beta^{-1} & k_{12} & \cdots & \cdots \\ k_{21} & \ddots & \ddots & \vdots \\ \vdots & \ddots & \ddots & k_{nn} + \beta^{-1} \end{bmatrix}$$

$$p(\mathbf{t}) = \int p(\mathbf{t}|\mathbf{y})p(\mathbf{y}) d\mathbf{y} = \mathcal{N}(\mathbf{t}|\mathbf{0}, \mathbf{C}) \quad (6.61)$$

$$C(\mathbf{x}_n, \mathbf{x}_m) = \underbrace{k(\mathbf{x}_n, \mathbf{x}_m)}_{\mathbf{x}_n} + \underbrace{\beta^{-1} \delta_{nm}}_{\mathbf{x}_m} \quad (6.62)$$

$$\mathbf{t} \sim \mathcal{N}(\mathbf{0}, \frac{1}{\beta} I_N + K)$$

GP for regression

$$p(\mathbf{y}) = \mathcal{N}(\mathbf{y}|\mathbf{0}, \mathbf{K}). \quad (6.60)$$

$$p(\mathbf{t}) = \int p(\mathbf{t}|\mathbf{y})p(\mathbf{y}) d\mathbf{y} = \mathcal{N}(\mathbf{t}|\mathbf{0}, \mathbf{C}) \quad (6.61)$$

$$C(\mathbf{x}_n, \mathbf{x}_m) = k(\mathbf{x}_n, \mathbf{x}_m) + \beta^{-1}\delta_{nm}. \quad (6.62)$$

training

Let's explicitly define a kernel

$$k(\mathbf{x}_n, \mathbf{x}_m) = \underbrace{\theta_0 \exp \left\{ -\frac{\theta_1}{2} \|\mathbf{x}_n - \mathbf{x}_m\|^2 \right\}}_{\text{sq-exponential}} + \underbrace{\theta_2}_{\text{const}} + \underbrace{\theta_3 \mathbf{x}_n^T \mathbf{x}_m}_{\text{linear}}. \quad (6.63)$$

We covered representation and "training" (implicit) in GP.
via \mathbf{C}

$$k(\mathbf{x}_n, \mathbf{x}_m) = \theta_0 \exp \left\{ -\frac{\theta_1}{2} \|\mathbf{x}_n - \mathbf{x}_m\|^2 \right\} + \theta_2 + \theta_3 \mathbf{x}_n^T \mathbf{x}_m. \quad (6.63)$$

$\theta_1 \sim \mathcal{F}_{\theta_1} \sim \theta^{**} 0.5 \text{ scale} \downarrow$

$y \sim N(0, k)$

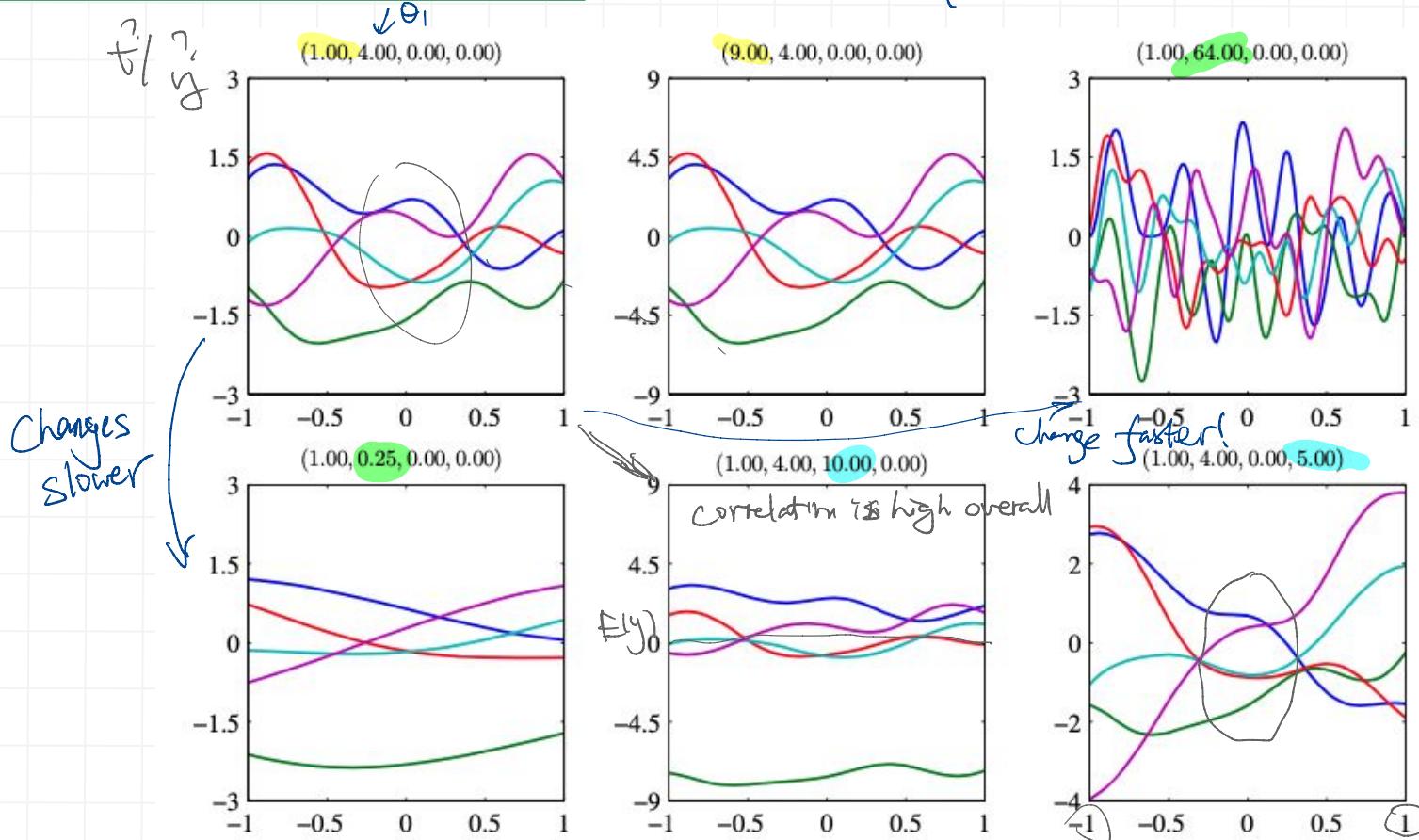


Figure 6.5 Samples from a Gaussian process prior defined by the covariance function (6.63). The title above each plot denotes $(\theta_0, \theta_1, \theta_2, \theta_3)$.

$$P(\vec{E}_N) \sim \mathcal{N}(0, C_N)$$

Predictive distribution $p(t_{N+1} | \mathbf{t}_N)$

$$p(\mathbf{t}_{N+1}) = \mathcal{N}(\mathbf{t}_{N+1} | \mathbf{0}, \mathbf{C}_{N+1})$$

$$\mathbf{C}_{N+1} = \begin{pmatrix} \mathbf{C}_N & \mathbf{k} \\ \mathbf{k}^T & c \end{pmatrix} \quad (6.64)$$

$\xrightarrow{\text{K+ } \beta^{-1} I_N}$ Kernel val of test pt x_{N+1}

vector \mathbf{k} has elements $k(\mathbf{x}_n, \mathbf{x}_{N+1})$

$$c = k(\mathbf{x}_{N+1}, \mathbf{x}_{N+1}) + \beta^{-1}$$

$$\mu_{a|b} = \mu_a + \Sigma_{ab}\Sigma_{bb}^{-1}(\mathbf{x}_b - \mu_b) \quad (2.81)$$

$$\Sigma_{a|b} = \Sigma_{aa} - \Sigma_{ab}\Sigma_{bb}^{-1}\Sigma_{ba}. \quad (2.82)$$

The conditional mean and variance of t_{N+1} given $\mathbf{t}_{1:N}$, expressed as a function of \mathbf{x}_{N+1}

$$m(\mathbf{x}_{N+1}) = \mathbf{k}^T \mathbf{C}_N^{-1} \mathbf{t} \quad \text{"training"} \quad (6.66)$$

$$\sigma^2(\mathbf{x}_{N+1}) = c - \mathbf{k}^T \mathbf{C}_N^{-1} \mathbf{k} \quad (6.67)$$

Variance reduction due to "training"

$$m(\mathbf{x}_{N+1}) = \sum_{n=1}^N a_n k(\mathbf{x}_n, \mathbf{x}_{N+1}) \quad (6.68)$$

where a_n is the n^{th} component of $\mathbf{C}_N^{-1} \mathbf{t}$.

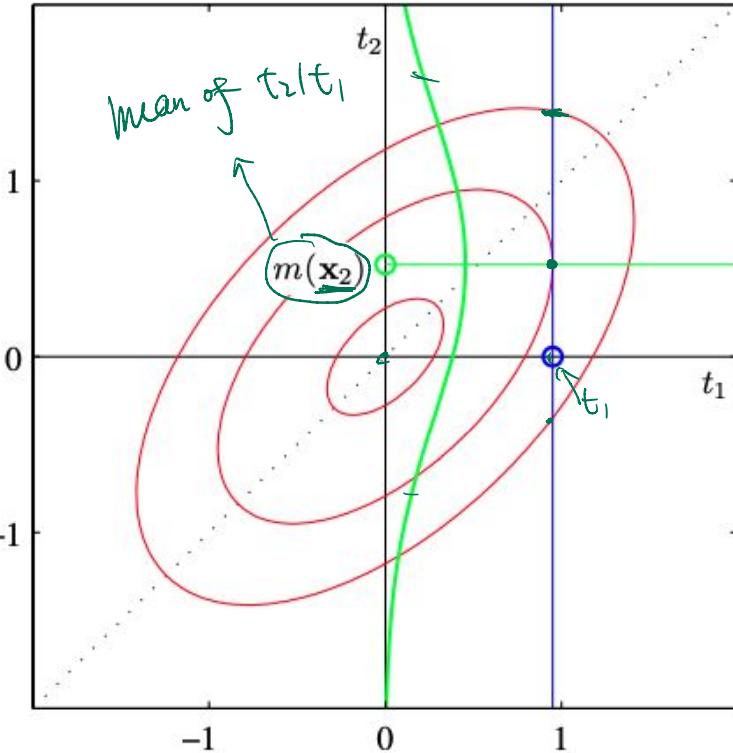
Figure 6.7

Illustration of the mechanism of Gaussian process regression for the case of one training point and one test point, in which the red ellipses show contours of the joint distribution $p(t_1, t_2)$. Here t_1 is the training data point, and conditioning on the value of t_1 , corresponding to the vertical blue line, we obtain $p(t_2|t_1)$ shown as a function of t_2 by the green curve.

$$p(\mathbf{t}_{N+1}) = \mathcal{N}(\mathbf{t}_{N+1} | \mathbf{0}, \mathbf{C}_{N+1})$$

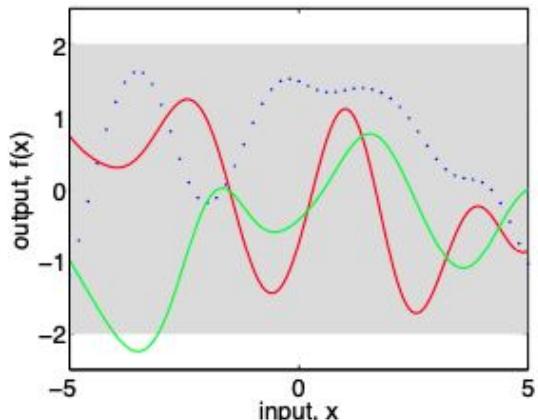
$$p\left(\begin{bmatrix} t_1 \\ t_2 \end{bmatrix}\right) \sim \mathcal{N}\left(\begin{pmatrix} 0 \\ 0 \end{pmatrix}, \begin{bmatrix} K_{11} & K_{12} \\ K_{12} & K_{22} \end{bmatrix}\right)$$

$$p(t_2|t_1)$$

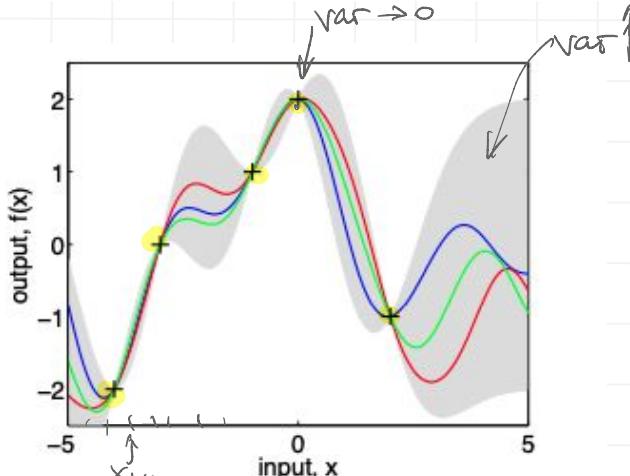


Notations

	Bishop book	GP book
"train" input	$x, \phi(x)$ ϕ	x, X
Output	$y = \phi w$ $t = y + \epsilon$	$f(x) = \phi^T(x) w$ $y \sim N(f, \sigma_n^2)$
Kernel / cov-fn	$K, L_N, \vec{K} = [k(x_{N+1}, x_{i:N})]$	$K(x, x_*), K(x, X), k(x, x')$
"test" input	x_{N+1}	x_x, x_*
Output	t_{N+1} (scalar) $w, m(x_{N+1}), f(x_{N+1})$	$\frac{f_*}{w, f_x}$ $V[f_*]$
Noise cov	β^{-1}	σ_n^2



(a), prior

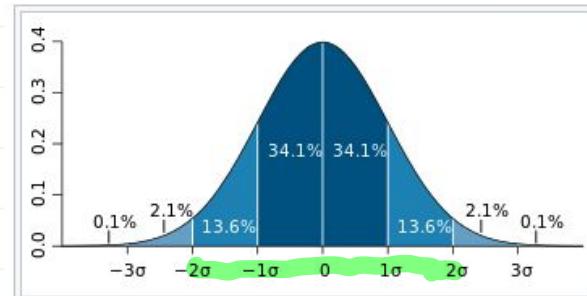


(b), posterior

Figure 2.2: Panel (a) shows three functions drawn at random from a GP prior; the dots indicate values of y actually generated; the two other functions have (less correctly) been drawn as lines by joining a large number of evaluated points. Panel (b) shows three random functions drawn from the posterior, i.e. the prior conditioned on the five noise free observations indicated. In both plots the shaded area represents the pointwise mean plus and minus two times the standard deviation for each input value (corresponding to the 95% confidence region), for the prior and posterior respectively.

$$\begin{aligned} \mathbf{f}_* | \mathbf{X}_*, \mathbf{X}, \mathbf{f} &\sim \mathcal{N}(K(\mathbf{X}_*, \mathbf{X}) K(\mathbf{X}, \mathbf{X})^{-1} \mathbf{f}, \\ &K(\mathbf{X}_*, \mathbf{X}_*) - K(\mathbf{X}_*, \mathbf{X}) K(\mathbf{X}, \mathbf{X})^{-1} K(\mathbf{X}, \mathbf{X}_*)) \end{aligned} \quad (2.19)$$

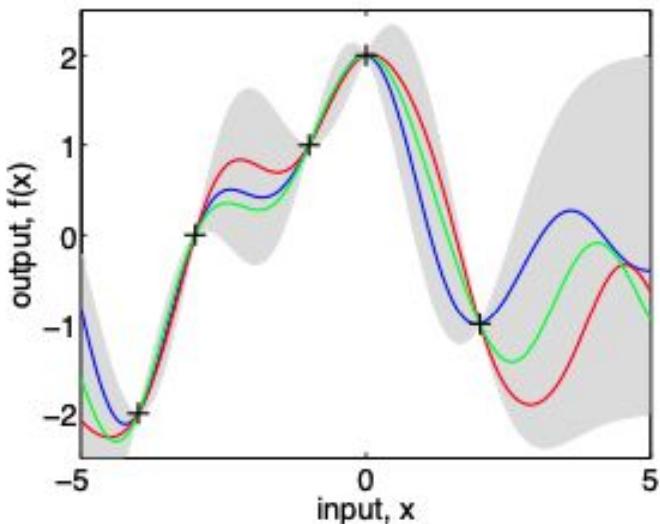
$$\begin{aligned} m(\mathbf{x}_{N+1}) &= \mathbf{k}^T \mathbf{C}_N^{-1} \mathbf{t} \\ \sigma^2(\mathbf{x}_{N+1}) &= c - \mathbf{k}^T \mathbf{C}_N^{-1} \mathbf{k} \end{aligned} \quad (6.66) \quad (6.67)$$

Noise-less: $\mathbf{C}_N = \mathbf{K}_N$ 

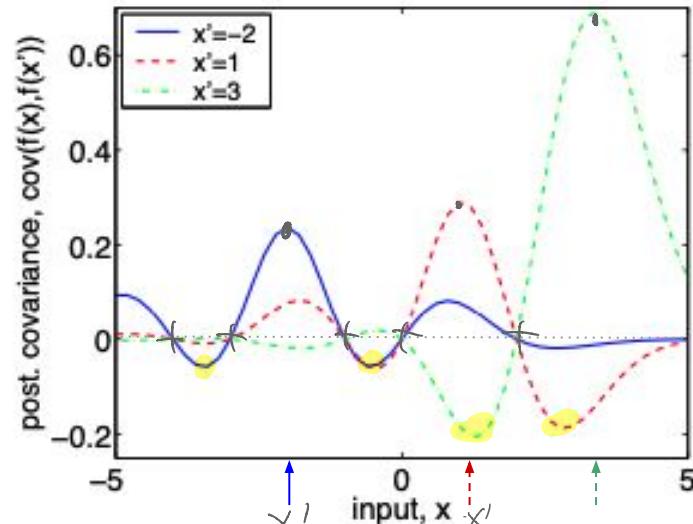
For the normal distribution, the values less than one standard deviation away from the mean account for 68.27% of the set; while two standard deviations from the mean account for 95.45%; and three standard deviations account for 99.73%.



\rightarrow g. induced by $\mathcal{O}_{x \text{ (IN.)}} K(x, x')$



(a), posterior



(b), posterior covariance

Figure 2.4: Panel (a) is identical to Figure 2.2(b) showing three random functions drawn from the posterior. Panel (b) shows the posterior *co*-variance between $f(\mathbf{x})$ and $f(\mathbf{x}')$ for the same data for three different values of \mathbf{x}' . Note, that the covariance at close points is high, falling to zero at the training points (where there is no variance, since it is a noise-free process), then becomes negative, etc. This happens because if the smooth function happens to be less than the mean on one side of the data point, it tends to exceed the mean on the other side, causing a reversal of the sign of the covariance at the data points. Note for contrast that the *prior* covariance is simply of Gaussian shape and never negative.

GP book Chap 2.2

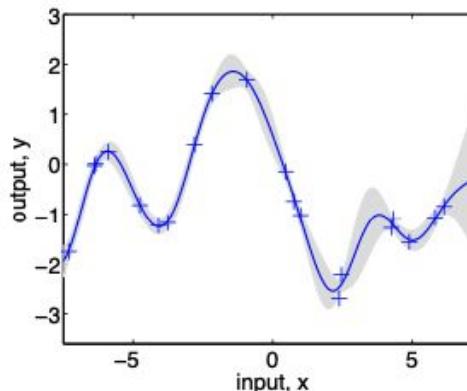
noise var.

$$k_y(x_p, x_q) = \sigma_f^2 \exp\left(-\frac{1}{2\ell^2}(x_p - x_q)^2\right) + \sigma_n^2 \delta_{pq}. \quad (2.31)$$

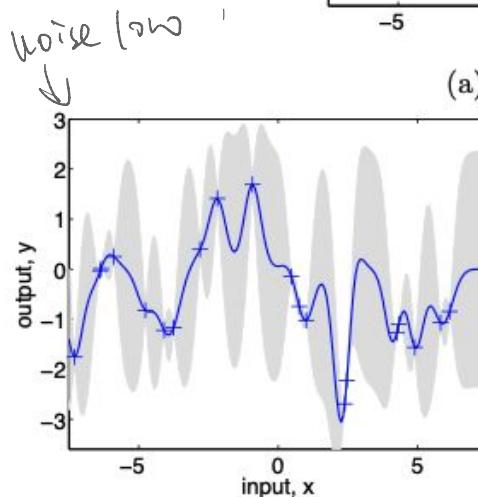
① Sampling w varying kernel param.

② posterior distribution

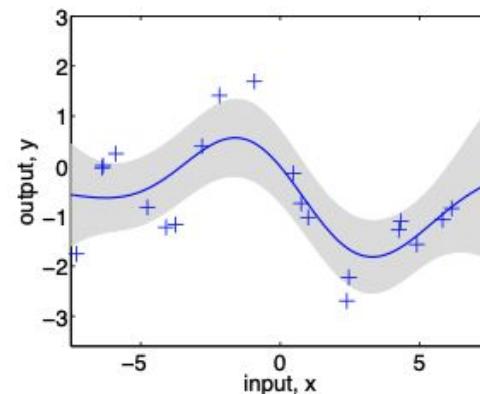
given input data.



(a), $\ell = 1$



(b), $\ell = 0.3$



(c), $\ell = 3$

Figure 2.5: (a) Data is generated from a GP with hyperparameters $(\ell, \sigma_f, \sigma_n) = (1, 1, 0.1)$, as shown by the + symbols. Using Gaussian process prediction with these hyperparameters we obtain a 95% confidence region for the underlying function f (shown in grey). Panels (b) and (c) again show the 95% confidence region, but this time for hyperparameter values $(0.3, 1.08, 0.00005)$ and $(3.0, 1.16, 0.89)$ respectively.

tion, shown in Figure A.6. The input values $\{x_n\}$ are generated uniformly in range $(0, 1)$, and the corresponding target values $\{t_n\}$ are obtained by first computing the corresponding values of the function $\sin(2\pi x)$, and then adding random noise with a Gaussian distribution having standard deviation 0.3. Various forms of this data set,

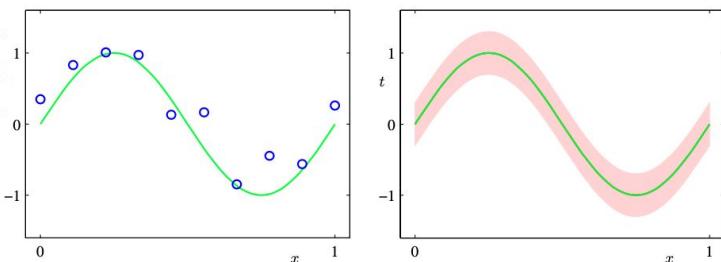
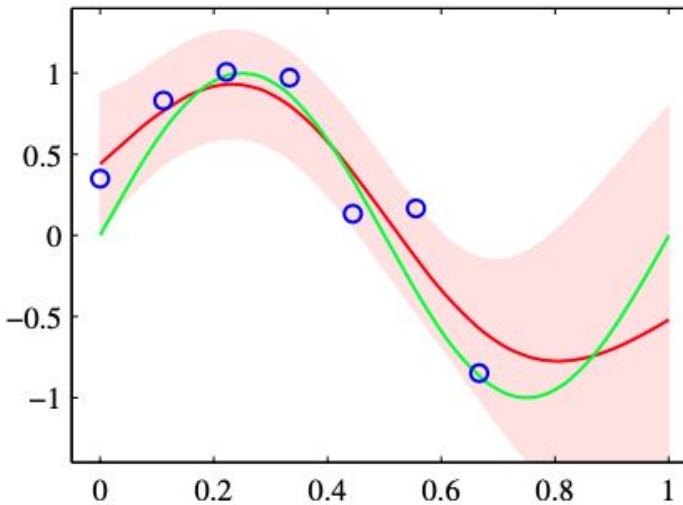


Figure A.6 The left-hand plot shows the synthetic regression data set along with the underlying sinusoidal function from which the data points were generated. The right-hand plot shows the true conditional distribution $p(t|x)$ from which the labels are generated, in which the green curve denotes the mean, and the shaded region spans one standard deviation on each side of the mean.

Figure 6.8 Illustration of Gaussian process regression applied to the sinusoidal data set in Figure A.6 in which the three right-most data points have been omitted. The green curve shows the sinusoidal function from which the data points, shown in blue, are obtained by sampling and addition of Gaussian noise. The red line shows the mean of the Gaussian process predictive distribution, and the shaded region corresponds to plus and minus two standard deviations. Notice how the uncertainty increases in the region to the right of the data points.



$$\sigma^2(\mathbf{x}_{N+1}) = c - \mathbf{k}^T \mathbf{C}_N^{-1} \mathbf{k}. \quad (6.67)$$

$$C(\mathbf{x}_n, \mathbf{x}_m) = k(\mathbf{x}_n, \mathbf{x}_m) + \beta^{-1} \delta_{nm}. \quad (6.62)$$

$$c = k(\mathbf{x}_{N+1}, \mathbf{x}_{N+1}) + \beta^{-1}$$

[Bishop 6.4]

Implementing

$$m(\mathbf{x}_{N+1}) = \mathbf{k}^T \mathbf{C}_N^{-1} \mathbf{t} \quad (6.66)$$

$$\sigma^2(\mathbf{x}_{N+1}) = c - \mathbf{k}^T \mathbf{C}_N^{-1} \mathbf{k}. \quad (6.67)$$

GP book Chap 2.2, A.4

input: X (inputs), \mathbf{y} (targets), k (covariance function), σ_n^2 (noise level),
 \mathbf{x}_* (test input)

- 2: $L := \text{cholesky}(K + \sigma_n^2 I) \quad \xrightarrow{\text{Cn}}$
- 3: $\alpha := L^T \backslash (L \backslash \mathbf{y}) \quad \leftarrow \xleftarrow{\text{Cn}^T}$
- 4: $f_* := \mathbf{k}_*^T \alpha \quad \} \text{ predictive mean eq. (2.25)}$
- 5: $\mathbf{v} := L \backslash \mathbf{k}_* \quad } \text{ predictive variance eq. (2.26)}$
- 6: $\mathbb{V}[f_*] := k(\mathbf{x}_*, \mathbf{x}_*) - \mathbf{v}^T \mathbf{v}$
- 7: $\log p(\mathbf{y}|X) := -\frac{1}{2} \mathbf{y}^T \alpha - \sum_i \log L_{ii} - \frac{n}{2} \log 2\pi \quad \text{eq. (2.30)}$
- 8: **return:** f_* (mean), $\mathbb{V}[f_*]$ (variance), $\log p(\mathbf{y}|X)$ (log marginal likelihood)

Algorithm 2.1: Predictions and log marginal likelihood for Gaussian process regression. The implementation addresses the matrix inversion required by eq. (2.25) and (2.26) using Cholesky factorization, see section A.4. For multiple test cases lines 4-6 are repeated. The log determinant required in eq. (2.30) is computed from the Cholesky factor (for large n it may not be possible to represent the determinant itself). The computational complexity is $n^3/6$ for the Cholesky decomposition in line 2, and $n^2/2$ for solving triangular systems in line 3 and (for each test case) in line 5.

Note also that the determinant of a positive definite symmetric matrix can be calculated efficiently by

$$|A| = \prod_{i=1}^n L_{ii}^2, \quad \text{or} \quad \log |A| = 2 \sum_{i=1}^n \log L_{ii}, \quad (\text{A.18})$$

where L is the Cholesky factor from A .

$A = LL^T$, L lower-triangular

To solve $Ax = b$ $\xrightarrow{\text{Cn}^T}$

i.e. $x = A^{-1}b$ or $\underline{A^{-1}b}$ for square A

First solve $Ly = b$, $y = L \backslash b$

Then $L^T x = y$

$x = L^T y = L^T \backslash (L \backslash b)$

$$L^T v = k^T *$$

$$v^T v = k^T C_n^{-1} k$$

GP vs Bayesian linear regression

$$m(\mathbf{x}_{N+1}) = \mathbf{k}^T \mathbf{C}_N^{-1} \mathbf{t} \quad (6.66)$$

$$\sigma^2(\mathbf{x}_{N+1}) = c - \mathbf{k}^T \mathbf{C}_N^{-1} \mathbf{k}. \quad (6.67)$$

Linear regression w basis functions

$$p(\mathbf{w}|\mathbf{t}) = \mathcal{N}(\mathbf{w}|\mathbf{m}_N, \mathbf{S}_N) \quad (3.49)$$

$$\mathbf{m}_N = \mathbf{S}_N (\mathbf{S}_0^{-1} \mathbf{m}_0 + \beta \Phi^T \mathbf{t}) \quad (3.50)$$

$$\mathbf{S}_N^{-1} = \mathbf{S}_0^{-1} + \beta \Phi^T \Phi. \quad (3.51)$$

Bottleneck:

compute \mathbf{C}_N^{-1} , $O(N^3)$ - at training time
 $O(N^2)$ for each test point

Bottleneck:

compute \mathbf{S}_N , $O(M^3)$ - at training time
 $O(M^2)$ for each test point

An example application

GP book Chap 2.5

- SARCOS robotic arm
 - x : (7 joint positions, 7 joint velocities, 7 joint accelerations)
 - y : 7 joint torques
 - Learning: $y \sim f(x)$
 - Prediction: compute the torque needed to move the arm along a given trajectory

Squared exponential covariance function

- Separate length scales for each input dim, σ_f^2, σ_n^2
- Optimise marginal likelihood on subset of data

~49K input-output pairs

- 44,484 for training
- 4,449 for testing

"inverting" (cholesky)
49K x 40K

Method	SMSE	MSLL
LR	0.075	-1.29
RBD	0.104	-
LWPR	0.040	-
GPR	0.011	-2.25

Table 2.1: Test results on the inverse dynamics problem for a number of different methods. The “-” denotes a missing entry, caused by two methods not producing full predictive distributions, so MSLL could not be evaluated.

Evaluation metrics

- SMSE - Standardized mean square error
- MSLL - Mean standardised log loss

*speed up with approximations, see Table 8.1

Learning hyperparameters

$$p(\mathbf{t}) = \int p(\mathbf{t}|\mathbf{y})p(\mathbf{y}) d\mathbf{y} = \mathcal{N}(\mathbf{t}|\mathbf{0}, \mathbf{C}) \quad (6.61)$$

$$\ln p(\mathbf{t}|\boldsymbol{\theta}) = -\frac{1}{2} \ln |\mathbf{C}_N| - \frac{1}{2} \mathbf{t}^T \mathbf{C}_N^{-1} \mathbf{t} - \frac{N}{2} \ln(2\pi). \quad (6.69)$$

Gaussian
(likelihood)

$$\frac{\partial}{\partial x} \ln |\mathbf{A}| = \text{Tr} \left(\mathbf{A}^{-1} \frac{\partial \mathbf{A}}{\partial x} \right) \quad (C.22)$$

$$\frac{\partial}{\partial x} (\mathbf{A}^{-1}) = -\mathbf{A}^{-1} \frac{\partial \mathbf{A}}{\partial x} \mathbf{A}^{-1} \quad (C.21)$$

$$\frac{\partial}{\partial \theta_i} \ln p(\mathbf{t}|\boldsymbol{\theta}) = -\frac{1}{2} \text{Tr} \left(\mathbf{C}_N^{-1} \frac{\partial \mathbf{C}_N}{\partial \theta_i} \right) + \frac{1}{2} \underbrace{\mathbf{t}^T \mathbf{C}_N^{-1}}_{\text{green bracket}} \underbrace{\frac{\partial \mathbf{C}_N}{\partial \theta_i}}_{\text{yellow circle}} \underbrace{\mathbf{C}_N^{-1} \mathbf{t}}_{\text{green bracket}}. \quad (6.70)$$

Gaussian Processes - Regression

Motivation

Defining Gaussian Processes (GP)

Kernel functions, sampling

GP regression

GP regression - predictive distribution

Sampling algorithm and computational costs

Bishop, Chap 6.4 (6.4.1-6.4.3)

GP book

<http://gaussianprocess.org/gpml/chapters/> Chap 1, 2.1, 2.2, 2.5