Graussian Processes for Regression References

- Graussian processes for ML (MIT Press 2006) [RW]
- PRML section 6.4
- Tensorflow probability





Graussian Processes Sor Regression



Quantify the uncertainty (or uncertainty quantify) Craussian Process to describe a distribution over functions.

A random process X(t) is a GP indexed by Xif for any $\{t_{1,...}, t_n\} \subset X$, a random vector formed by $X(t_1),... X(t_n)$ is jointly Gaussian (or $(X(t_1)... X(t_n))^T$ is a multivariate Gaussian)

Specification of GP
GP can be specified by
- mean function
$$\mu: X \rightarrow \mathbb{R}$$
 set $\mu(t) = \mathbb{E}[X(t)]$
- covariance function $k: X \times X \rightarrow \mathbb{R}$ set
 $k(t,s) = cov [X(t), X(s)]$ a.k.a Kernel function
The kernel function is required only to be symmetric and

positive definite.

$$X(t) \sim GP(\mu(t), K(t, s))$$
 t, s $\in X$

Example D-Jim multivariate Gaussian distribution is a GP intexet by {1,2,... D} See PRML section 2.3.2 Example X(t):= tx where x = N(x|0,1) and $t \in \mathbb{R}$ Then X(t) is a GP indexed by IR t.... to x(t.) x(t.) M(t) = E[X(t)] = t E[x] = 0 $k(t,s) = cov [X(t), X(s)] = E[tx \cdot sx] = ts E[x^2] = ts$



For $t_1 < t_2 < \cdots < t_n$, joint density $X(t_1), \dots, X(t_n)$ $X(t_1), X(t_2) - X(t_1), \dots, X(t_n) - X(t_{n-1})$ are independent and $X(t_k) - X(t_{k-1}) \sim N(o, t_k - t_{k-1})$

Eq. Brownian motion {B(t), t20} is a GP indexed by IR^t Each of B(tr)... B(tr) can be expressed as a linear combination of independent normal r.v B(tr), B(tr)-B(tr)..., B(tr)-B(tr) Gaussian Processes for Regression (GPR) Notations

X: intex set of GP

Z: collection of random variables indexed by X

 $Z(x) \sim GP(\mu(x), k(x, x')) \quad x, x' \in X$

For any finite set
$$\mathcal{K}_1, \mathcal{K}_2, \ldots \mathcal{K}_n \in \mathcal{X}_n$$

 $(\mathcal{Z}(\mathcal{K}_1), \ldots \mathcal{Z}(\mathcal{K}_n))^T = \mathcal{N}(\mathcal{M}_n, \mathcal{K})$

where
$$M_n = (M(X_1), M(X_2), \dots, M(X_n))^T$$
 $K_{\lambda i} = K(X_{\lambda}, X_{i})$
= $cov [Z(X_{\lambda}), Z(X_{i})]$

intex set of $GP \approx$ domain of function

GPR: non parametric Bayesian regression method using the properties of GP

Idea i GP as defining a distribution over functions and inference taking place directly in the space of functions. We dispense with the parametric model and instead define a prior distribution over functions directly.

Linear regression revisited
Training data set
$$D = (x_n, t_n)_{n=1...,N}$$

 $x : D - dim input vector, $t : tanget value$
 $y(x_i): prediction value at x$
predictive distribution
Goal : find $P(t^* | x^*, X, t_i)$ at some new input x^*
Fix basis functions $g_{1,...}, g_{M-1}$ or seature map $\overline{P}(x_i)$
Linear model$

 $Y(X, W) := W^T \overline{Q}(X)$ weight vector $W \in \mathbb{R}^M$

Model assumption

$$t \sim N(t) y(x, w), \beta')$$
 precision $\beta > 0$

prior over w

$$P(W|X) = N(W|0, x^{-1}L)$$

posterior over

P

$$(w_{I} \#) = N(w_{I} m_{N}, S_{N})$$
$$m_{N} = \beta S_{N} \Phi^{T} \#$$
$$S_{N}^{-1} = \alpha I + \beta \Phi^{T} \Phi$$

Example

Let
$$W \sim N(W|\Theta, \kappa^{T}I)$$
 be as a prior.
For $n=1,2,...,N$, let $Y_{n}:=Y(x_{n}, w_{r}) = W^{T} \oplus (x_{n})$. Then Y_{n}
can be seen a $I-\dim$ random variable $(I-\dim$ Graussian)
Let $Y:=(Y_{1}...,Y_{N})^{T}$. Then
 $Y = \bigoplus W_{r}$
MXM
 $Y = \bigoplus W_{r}$
MXI
is the $N-\dim$ random vector where $\bigoplus_{n \in I} = \mathscr{P}_{K}(X_{n})$.

$$\overline{\Phi} = \begin{pmatrix} \varphi_{c}(x_{1}) & \varphi_{c}(x_{2}) & \cdots & \varphi_{M-1}(x_{N}) \\ \varphi_{c}(x_{2}) & \varphi_{c}(x_{2}) & \varphi_{M-1}(x_{N}) \\ \vdots & \vdots & \vdots \\ \varphi_{c}(x_{N}) & \vdots & \varphi_{M-1}(x_{N}) \end{pmatrix} = \begin{pmatrix} \overline{\varphi}(x_{1})^{T} \\ \overline{\varphi}(x_{N})^{T} \\ \vdots \\ \overline{\varphi}(x_{N})^{T} \end{pmatrix}$$
Since Y is a linear combination of Wr (Graussian dist.)
 Y is $N - \dim$ Graussian. (Y is determined by the first
and the second moments)
 $E[Y] = E[\overline{\Phi} W] = \overline{\Phi} E[W] = \Theta$
 $Cov [Y] = E[\overline{\Psi} Y^{T}] = \overline{\Phi} E[W W]^{T} = \overline{\Phi}$

I is called design matrix so that

where IK is the gram matrix

$$K_{nm} = K(\mathcal{X}_n, \mathcal{X}_m) = \frac{1}{\alpha} \overline{\Phi}(\mathcal{X}_n)^{\mathsf{T}} \overline{\Phi}(\mathcal{X}_m)$$

In general, GP is refined as a probability distribution over functions y(x) sit the set of values of y(x)evaluated at an arbitrary set $x_{1,...}$ x_{N} jointly have a Gaussian Distribution.

A key point about Graussian processes is that the distribution over N variables $y_1 \dots y_N$ is specified joint completly by the second - order statistics. O Since we have not any knowledge about the mean of ycx), we take it to be zero. This is equivalent to choosing the mean vector of prior pcw/(x) to be 0 $(\mathbf{2})$ The specification of GP is then completed by giving covarience of y(x) evaluated at any two values the of X with $[\text{kemel} \mathbb{E}[Y(X_n), Y(X_m)] = K(X_n, X_m).$

Let f(x)(Y(x)) be a zero mean GP indexed by its domain space. IK I.e. $\begin{pmatrix} Y(x_{i}) \\ \vdots \\ Y(x_{N}) \end{pmatrix} \sim N\left(\bigotimes_{i} \begin{bmatrix} k(x_{i}, x_{i}) \cdots k(x_{i}, x_{N}) \\ \vdots \\ k(x_{N}, x_{i}) \cdots k(x_{N}, x_{N}) \end{bmatrix} \right)$

and

$$\begin{pmatrix} \gamma(\mathcal{X}_{i}) \\ \vdots \\ \gamma(\mathcal{X}_{N}) \\ \gamma(\mathcal{X}_{N}) \end{pmatrix} \sim \mathcal{N} \left(\emptyset, \begin{bmatrix} k(\mathcal{X}_{i}, \mathcal{X}_{i}) & \cdots & k(\mathcal{X}_{i}, \mathcal{X}_{N}) & k(\mathcal{X}_{i}, \mathcal{X}_{i}) \\ \vdots & \ddots & \vdots \\ k(\mathcal{X}_{N}, \mathcal{X}_{i}) & \cdots & k(\mathcal{X}_{N}, \mathcal{X}) & \vdots \\ k(\mathcal{X}, \mathcal{X}_{i}) & \cdots & k(\mathcal{X}, \mathcal{X}) & k(\mathcal{X}, \mathcal{X}_{i}) & - \end{pmatrix} \right)$$

Assume
$$t = y + \varepsilon_n$$
 with $\varepsilon_{noise} \sim N(o, \beta^{-1})$

 $\Rightarrow \begin{pmatrix} t_i \\ \vdots \\ t_n \end{pmatrix} \sim N(o, |k + \beta^{-1} I_n)$
and $\begin{pmatrix} t_i \\ \vdots \\ t_n \\ t_n^* \end{pmatrix} \sim N(o, |k^* + \beta^{-1} I_{n+1})$
Since $P(t, t^*)$ is Gaussian so is $P(t^* + t^*)$
 $P(t^* + t^*) = N(t^* + m(x^*), \sigma^{-2}(x^*))$

GPR in tetail We shall consider the noise on the observed target value Assumption: Graussian noisy model.

$$t_n = Y_n + \varepsilon_\eta$$
 $n = 1, 2, ... N$

where E_n followed zero mean Graussian with precision β . So for an N-dim vector $\gamma := (\gamma_1 \dots \gamma_N)^T$, the probability distribution over $# := (t_1, \dots t_N)^T$ can be expressed as

$$P(\#|\forall) = N(\#|\forall, \beta'I_N)$$

where I_N is the $N \times N$ identity matrix. From the definition of G_P , (prior dist. over Y) P(Y) = N(Y | O, |K) N-dim Graussian

The kernel function determining the gram matrix IK is typically chosen to express the property that $x_n \sim x_m \implies y(x_n)$ and $y(x_m)$ are strongly correlated

One widely used kernel function for GPR is given

$$K(\mathcal{X}_{n}, \mathcal{X}_{m}) := O_{0} \exp \left(-\frac{O_{1}}{2} \| \mathcal{X}_{n} - \mathcal{X}_{m} \|^{2} \left(+ O_{2} + O_{3} \mathcal{X}_{n}^{T} \mathcal{X}_{m} \right) \right)$$

Now we consider the distribution over
$$#$$
. By PRML
section 2.3.3 (linear Granssian model), we obtain
 $P(#) = \int P(#19) P(9) d9 = N(4100, C)$ (6.61)
where $C = 1k + \beta^{-1}I_N$ with $C(x_n, x_m) = k(x_n, x_m) + \beta^{-1}S_{nm}$

Using.	Graussian	process	view point,	we s	hall build	a model
os the	e joint	probability	distributi	on over	sets of	data
points.						
Groal:	Griven	training d	lata points	make	prediction	fo
	target	value 50	r new	input		
* :	some	new inpu	it, t'	+ : its	target	value
(※,世) set	of training	g Jata			
Find	the pre	edictive dist	cribution P	(++)	= P Ct* (**	*, %, +)

Let
$$\#^* := (t_1, \dots, t_N, t^*)^T$$
 (N+1 - dim random vector)
To find $p(t^*| t_1)$, we first consider $p(t^*)$. From C6.61),
 $p(t^*) = N(t^*| o, C^*)$ N+1 - dim
Graussian
The covariance matrix C^* can be decomposed as
NXN

$$\mathbb{C}^{*} = \begin{pmatrix} \mathbb{C} & \mathbb{K} \\ \mathbb{K}^{\mathsf{T}} & \mathbb{C} \end{pmatrix} \qquad (\mathbb{N} + 1) \times (\mathbb{N} + 1)$$

where IK is the N-dim vector whose nth element is $K(x_n, x^*)$ and $c = K(x_n^*, x_n^*) + \beta^{-1}$.

Since the conditional Gaussian distribution is again
a Gaussian,
$$p(t^*|t)$$
 is Gaussian. $t^* = \begin{pmatrix} t \\ t^* \end{pmatrix}$
From (2.81) and (2.82), $\wedge t1$ dim

$$P(t^*|t) = N(t^*|m(x^*), \sigma^2(x^*))$$

where

$$m(x^*) = k^T C^{-1} \#$$
) depending on x^*
 $G^2(x^*) = c - k^T C^{-1} k$

The expectation of predictive distribution can be written as

$$m(x^*) = \sum_{n=1}^{N} a_n \in (x_n, x^*)$$

- where an is the nth component of C¹t Remark:
- Computational costs of GPR: the invertion of $N \times N$ matrix requires $O(N^3)$ once. For each new point, GPR has the cost $O(N^2)$ from the matrix multiplication

GPR in the real world (learning for hyperparameters) Hyperparameter 0:= 100, 01, 02, 03 (and B Log likelihoot function $\ln P(\# | Q) = -\frac{1}{2} \ln | Q | - \frac{1}{2} \#^{T} C^{-1} \# - \frac{N}{2} \ln (2\pi)$ $\frac{\partial}{\partial Q_{i}} \ln P(\#|Q) = -\frac{1}{2} \left(C^{-1} \frac{\partial C}{\partial Q_{i}} \right) + \frac{1}{2} \#^{T} C^{-1} \frac{\partial C}{\partial Q_{i}} C^{-1} \#$ Generally, ln P(#10) is not convex.

Tensorflow Probability

Exponentiated Quadratic

$$(K_{n}, K_{m}) := O_{0}^{2} \exp\left(\frac{\|K_{n} - K_{m}\|^{2}}{-2 O_{1}^{2}}\right)$$

Oo: amplitude

O₁: length scale

1/2: observation noise variance.

$$K(\mathcal{K}_{n}, \mathcal{K}_{m}) := O_{0} \exp \left(-\frac{O_{1}}{2} \| \mathcal{K}_{n} - \mathcal{K}_{m} \|^{2} \right) + O_{2} + O_{3} \mathcal{K}_{n}^{T} \mathcal{K}_{m}$$