

TAM 598 Lecture 21:

Gaussian Process Regression

Announcements:

- HW 5 covers lectures 17-20; due on Fri Apr 18
- HW 6 covers lectures 21-23; due on Fri May 2

Gaussian Process Regression - fully non-parametric, Bayesian regression

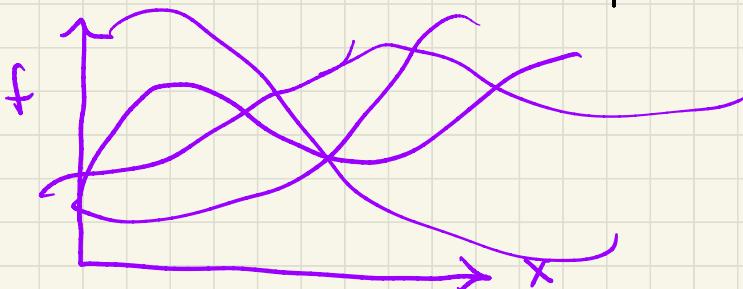
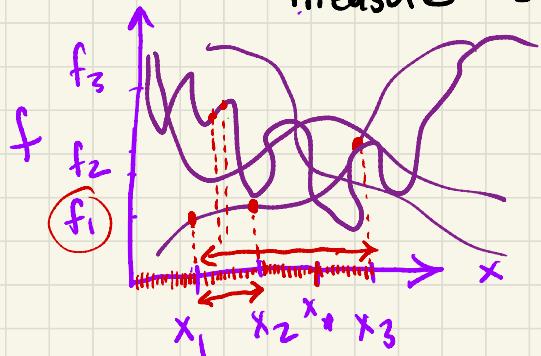
goal: given some data, learn a function $f(\cdot) : \mathcal{X} \rightarrow \mathbb{R}$

approach:

(1) before gathering any data, use your beliefs to generate a probability measure $p(f(\cdot))$ so that we can sample possible f 's. This is your prior

{ (2) gather data D and model the likelihood of the data $p(D | f(\cdot))$

(3) Use Bayes' rule to come up with your posterior probability measure over the space of functions $p(f(\cdot) | D)$



Stochastic process - collection of random variables $\{\tilde{x}_i\}$

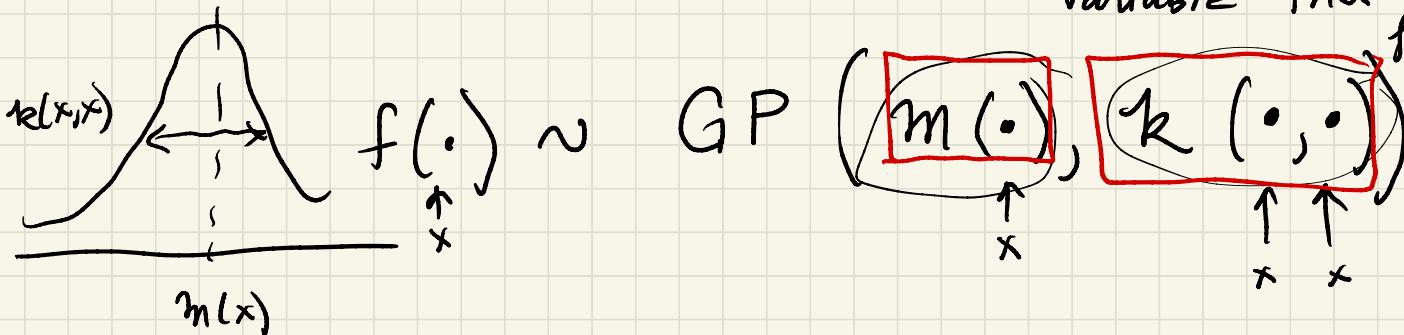
for some i in set I .

e.g. $\tilde{X}_t = \tilde{X}(t)$ is a stochastic process, parametrized by time

usually discrete

Gaussian process - a generalization of a multivariate Gaussian to infinite dimensions, so a continuous stochastic process

" $f(\cdot)$ is a Gaussian process" \Rightarrow f is a random variable that is a function



difference
w/ multivariate
Gaussian

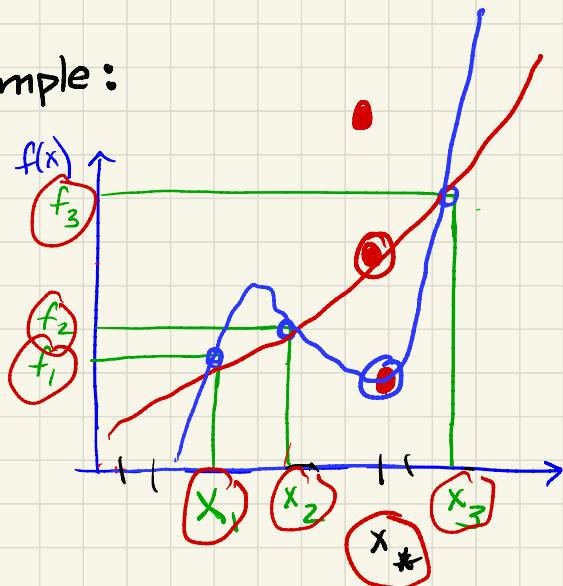
$$f(\cdot) \sim GP(m(\cdot), k(\cdot, \cdot))$$

random function
not a random vector

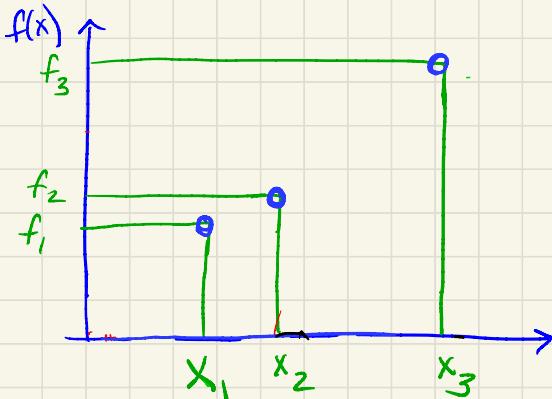
mean function,
not mean
vector

covariance
function, not
cov. matrix

example:



say we know $f(x_1) = f_1$,
 $f(x_2) = f_2$, $f(x_3) = f_3$. We
use a GP to predict f at
any other point $x = x^*$



Approach:

① imagine that f_1, f_2, f_3 are components of a 3D random vector, generated by a 3-variate gaussian.

$$\text{R.V. } \tilde{f} = \begin{pmatrix} f(x_1) \\ f(x_2) \\ f(x_3) \end{pmatrix} = \begin{pmatrix} f_1 \\ f_2 \\ f_3 \end{pmatrix}$$

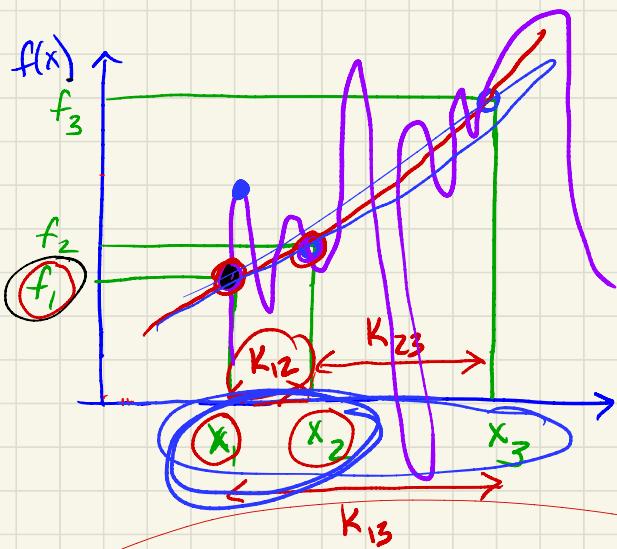
② Then

$$f \sim \mathcal{N} \left(\tilde{m}(\tilde{x}_{1:3}), \tilde{k}(\tilde{x}_{1:3}, \tilde{x}_{1:3}) \right)$$

where $\tilde{m}(\tilde{x}_{1:3}) = \begin{pmatrix} m(\tilde{x}_1) \\ m(\tilde{x}_2) \\ m(\tilde{x}_3) \end{pmatrix}$ and $\tilde{k}(\tilde{x}_{1:3}, \tilde{x}_{1:3}) = \begin{pmatrix} k(x_1, x_1) & k(x_1, x_2) & k(x_1, x_3) \\ k(x_2, x_1) & k(x_2, x_2) & k(x_2, x_3) \\ k(x_3, x_1) & k(x_3, x_2) & k(x_3, x_3) \end{pmatrix}$

Regression

M



→ x_i 's are given

→ model f_i 's as if pulled from a multivariate gaussian

$$\left[\begin{array}{c} f_1 \\ f_2 \\ f_3 \end{array} \right] \sim \left(\begin{array}{c} 0 \\ 0 \\ 0 \end{array} \right),$$

$$\left(\begin{array}{ccc} K_{11} & K_{12} & K_{13} \\ K_{21} & K_{22} & K_{23} \\ K_{31} & K_{32} & K_{33} \end{array} \right)$$

$$\sim \left(\begin{array}{c} 0 \\ 0 \\ 0 \end{array} \right), \quad \left(\begin{array}{ccc} 1 & 0.7 & 0.2 \\ 0.7 & 1 & 0.5 \\ 0.2 & 0.5 & 1 \end{array} \right)$$

$$K_{ij} = \exp\left(-\frac{1}{\lambda_i} \|x_i - x_j\|^2\right)$$

$$= \begin{cases} 0 & \|x_i - x_j\| \rightarrow \infty \\ 1 & x_i = x_j \end{cases}$$

$$f \sim N(0, K)$$

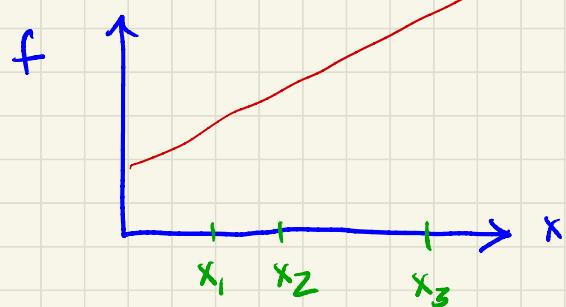
How to interpret the mean $m(\cdot)$?

GP($m(\cdot)$, $k(\cdot, \cdot)$)

- For any point \underline{x} , $m(\underline{x})$ should give us the expected value

$$\mathbb{E}[f(\underline{x})]$$

- prior* • before we measure any data, options for $m(\underline{x})$ are



① $m(\underline{x}) = 0$ or $m(\underline{x}) = \text{constant}$

② $m(\underline{x}) = c_0 + \sum_{i=1}^d c_i \underline{x}_i$ linear

③ use basis functions

$$m(\underline{x}) = \sum_{i=1}^m c_i \phi_i(\underline{x})$$

- ④ generalized polynomial chaos (gPC): use d polynomial basis functions up to degree R : $m(\underline{x}) = \sum_{i=1}^d c_i \phi_i(\underline{x})$ with $\int \phi_i(\underline{x}) \phi_j(\underline{x}) d\mu(\underline{x}) = \delta_{ij}$

How to interpret the covariance?

- ④ diagonal terms $k(\underline{x}, \underline{x})$ are the variance of the R.V.

$$f(\underline{x})$$

$$\mathbb{V}[f(\underline{x})] = \mathbb{E}[(f(\underline{x}) - m(\underline{x}))^2]$$

so that w/ 95% confidence, $f(\underline{x})$ falls in

$$m(\underline{x}) \pm 2\sqrt{k(\underline{x}, \underline{x})}$$

- ④ off-diagonal terms $\underline{k(\underline{x}, \underline{x}')}$ measure the correlation of $f(\underline{x})$ and $f(\underline{x}')$

$$k(\underline{x}, \underline{x}') = \mathbb{C}[f(\underline{x}), f(\underline{x}')]$$

$$= \mathbb{E}[(f(\underline{x}) - m(\underline{x}))(f(\underline{x}') - m(\underline{x}'))]$$

is a measure of similarity

properties of covariance function

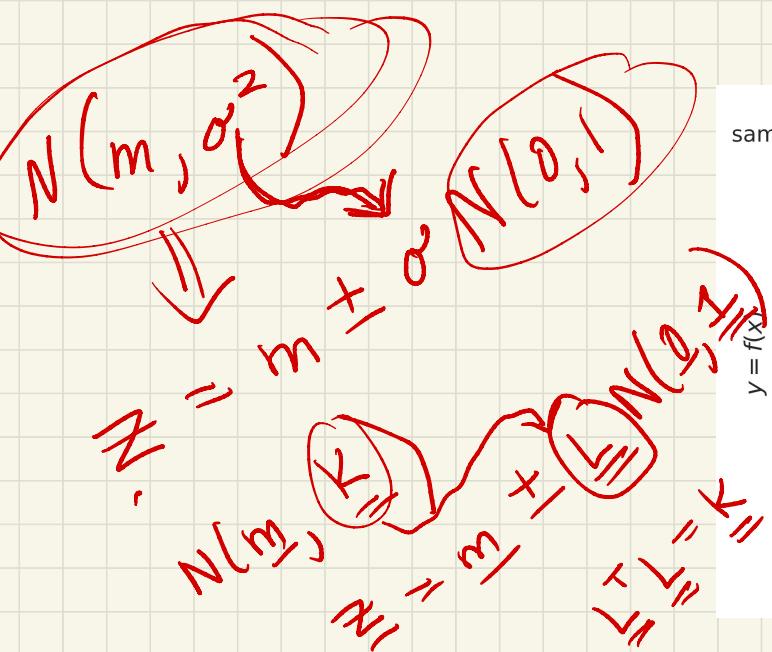
- 1) $k(\underline{x}, \underline{x}) > 0$ since it is a variance
 - 2) $\underline{k}(\underline{x}, \underline{x}')$ becomes smaller as the distance between $\underline{x}, \underline{x}'$ grows
 - 3) for any choice of specific points $\underline{x}_{1:n}$, the covariance matrix K needs to be positive definite
- covariance functions govern smoothness
- they can be designed to model invariances, eg if $k(T\underline{x}, \underline{x}') = k(\underline{x}, T\underline{x}') = k(\underline{x}, \underline{x}')$ then the GP is also invariant wrt T

Gaussian Process ; a distribution over functions

$$f(x) \sim GP(m(x), K(x, x'))$$

$$m(x) = \mathbb{E}[f(x)]$$

$$K(x, x') = \mathbb{E}[(f(x) - m(x))(f(x') - m(x'))^\top]$$



$$K(x, x') = \exp\left(-\frac{1}{2}(x-x')^2\right)$$

5 different function realizations at 41 points
sampled from a Gaussian process with exponentiated quadratic kernel

