

TAM 598 Lecture 22:

Gaussian Process Regression

Announcements:

- HW b covers lectures 21-23; due on Fri May 2

where we are:

- * we have a Gaussian Process prior (an "infinite variate" gaussian) for unknown function $f(\cdot)$, given by $p(f(\cdot))$
- * we are going to collect some data and use Bayes' rule to update our prior and obtain the posterior

$$p(f(\cdot) | D) \propto p(D | f(\cdot)) p(f(\cdot))$$

→ (we'll do this for an arbitrary n -variate Gaussian;
the approach generalizes to the continuous infinite case)

$$p(f_1, f_2, \dots, f_n | D)$$

Gaussian Process ; a distribution over functions

$$D = \{ (x_i, f_i), i=1, \dots, N \}$$

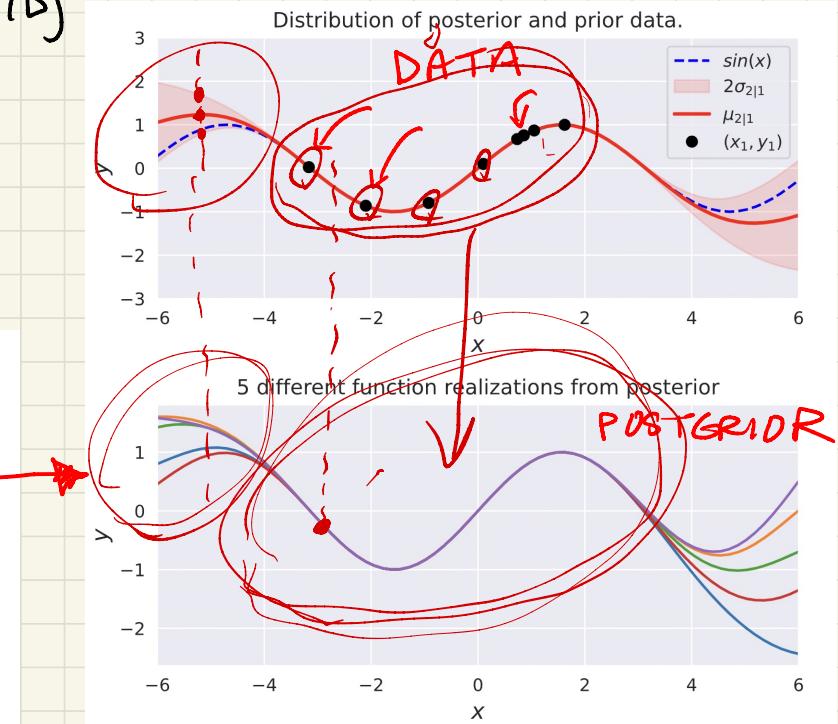
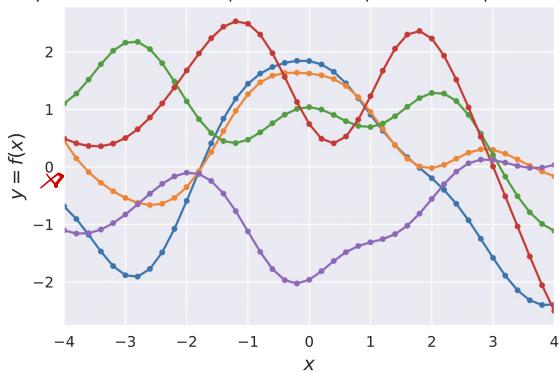
$$p(f|D) = \frac{p(D|f) p(f)}{p(D)}$$

posterior

prior

PRIOR

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



$$f(\cdot) \sim GP(m(\cdot), k(\cdot, \cdot)) \quad \text{so}$$

for any points $\underline{x}_{1:n} = (\underline{x}_1, \underline{x}_2, \dots, \underline{x}_n)$ the joint pdf of $f_{1:n} = (f(\underline{x}_1), \dots, f(\underline{x}_n))$ is the multivariate gaussian

$$\underline{f}_{1:n} \mid \underline{x}_{1:n} \sim N\left(\underline{m}(\underline{x}_{1:n}), \underline{\Sigma}_{1:n} \right)$$

we observe $D = (\underline{x}_{1:n}, \underline{y}_{1:n})$

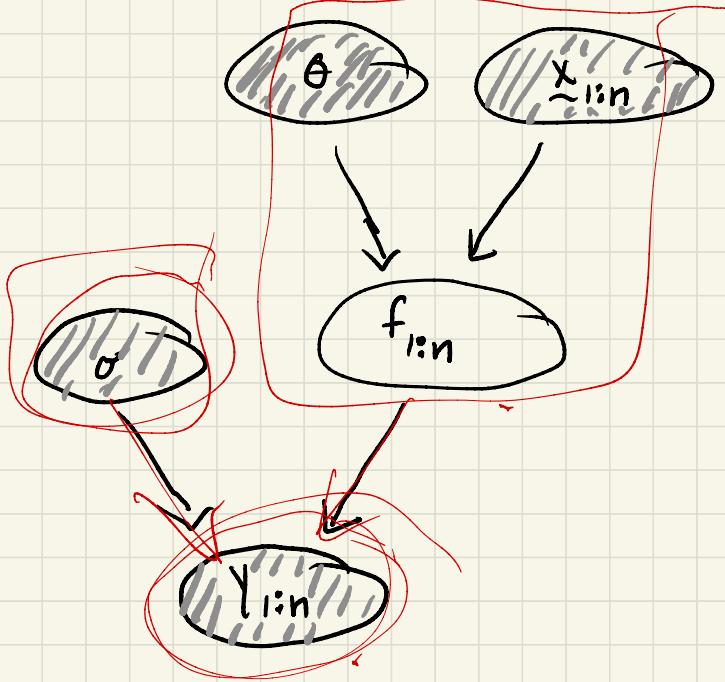
\hookrightarrow noisy measure of $f(\underline{x}_{1:n})$, gaussian noise, variance σ^2

single observation

$$(y_i) \mid f(\underline{x}_i) \sim N(f(\underline{x}_i), \sigma^2)$$

all observations:

$$\underline{y}_{1:n} \mid f(\underline{x}_i) \sim N\left(f_{1:n}, \sigma^2 \underline{I}_n\right)$$



- fixed hyperparameters θ and σ^2
- arbitrary collection of test points, densely cover input space
 $\underline{x}_{1:n^*}^* = (\underline{x}_1^*, \underline{x}_2^*, \dots, \underline{x}_{n^*}^*)$
and function values
 $\underline{f}_{1:n^*}^* = (f(\underline{x}_1^*), \dots, f(\underline{x}_{n^*}^*))$

Consider the joint pdf of \underline{f} and \underline{f}^* , a MVG,

$$p(\underline{f}_{1:n}, \underline{f}_{1:n^*}^*) = N\left(\begin{pmatrix} \underline{f}_{1:n} \\ \underline{f}_{1:n^*}^* \end{pmatrix} \middle| \begin{pmatrix} \underline{m}(\underline{x}_{1:n}) \\ \underline{m}(\underline{x}_{1:n}^*) \end{pmatrix}, \begin{pmatrix} \underline{\Sigma}(\underline{x}_{1:n}, \underline{x}_{1:n}) \\ \underline{\Sigma}(\underline{x}_{1:n}, \underline{x}_{1:n}^*) \end{pmatrix}\right)$$

$$\begin{aligned}
 P(f_{-1:n}^* | \bar{x}_{1:n}^*, D) &= P(f_{-1:n}^* | \bar{x}_{1:n}^*, \bar{x}_{1:n}, y_{1:n}) \\
 &= \int P(f_{1:n}, f_{1:n}^* | \bar{x}_{1:n}^*, \bar{x}_{1:n}, y_{1:n}) df_{1:n} \\
 &\propto \underbrace{\int P(y_{1:n} | f_{1:n})}_{\text{gaussian}} \underbrace{\int P(f_{-1:n}, f_{1:n}^* | \bar{x}_{1:n}^*, \bar{x}_{1:n}) df_{1:n}}_{\text{gaussian}}
 \end{aligned}$$

This result is a new Gaussian

$$P(f_{-1:n}^* | \bar{x}_{1:n}^*, D) = N \left(f_{-1:n}^* \mid \underline{m}_n(\bar{x}_{1:n}), \underline{\Sigma}_n(\bar{x}_{1:n}^*, \bar{x}_{1:n}) \right)$$

the result is a new Gaussian

prior: $\underline{m}(x), \underline{k}(x)$

$$\underbrace{p(f_{1:n^*}^* | \underline{x}_{1:n^*}^*, D)} = N\left(f_{1:n^*}^* | \underline{m}_n(\underline{x}_{1:n^*}), \underline{K}_n(\underline{x}_{1:n^*}, \underline{x}_{1:n^*})\right)$$

with posterior mean

$$\underline{m}_n(x) = \underline{m}(x) + \underline{k}(x, \underline{x}_{1:n}) (\underline{K}(\underline{x}_{1:n}, \underline{x}_{1:n}) + \sigma^2 \underline{I}_n)^{-1} (\underline{y}_{1:n} - \underline{m}(\underline{x}_{1:n}))$$

and posterior covariance

$$\underline{K}_n(x, x') = \underline{k}(x, x) - \underline{k}(x, \underline{x}_{1:n}) (\underline{K}(\underline{x}_{1:n}, \underline{x}_{1:n}) + \sigma^2 \underline{I}_n)^{-1} \underline{k}^T(x, \underline{x}_{1:n})$$

where $\underline{k}(x, \underline{x}_{1:n}) = (k(x, x_1), k(x, x_2), \dots, k(x, x_n))$ is
a cross-covariance.

since the test points are arbitrary \Rightarrow

$$f(\cdot) | D \sim GP(m_n(\cdot), k_n(\cdot, \cdot))$$

point predictive distribution - to predict $f(\cdot)$ at a single point, have your test points $\underline{x}_{\text{test}}^*$ just be a single point \underline{x}^* . Then

$$P(f(\underline{x}^*) | D) = N(f(\underline{x}^*) \mid \underbrace{\underline{m}_n(\underline{x}^*), \underline{k}_n(\underline{x}^*, \underline{x}^*)}_{\text{this is a predictive variance}})$$

this is a predictive variance

but the predicted measurement outcome y^* at \underline{x}^* is

$$P(y^* | \underline{x}^*, D) =$$

$$= N(f(\underline{x}^*) \mid \underline{m}_n(\underline{x}^*), \underline{k}_n(\underline{x}^*, \underline{x}^*) + \sigma^2)$$

need to include
noise variance

HYPERPARAMETER OPTIMIZATION

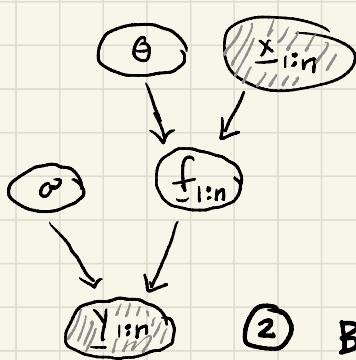
- covariance function: θ
- measurement variance: σ^2

If we don't know them:

- start with priors $p(\theta)$, $\underline{p(\sigma)}$
- Bayes Rule

$$\begin{aligned} p(\theta, \sigma | D) &\propto p(D | \theta, \sigma) p(\theta) p(\sigma) \\ &= \int p(\underline{y}_{1:n} | f_{1:n}, \sigma) p(f_{1:n} | \underline{x}_{1:n}, \theta) df_{1:n} p(\theta) p(\sigma) \end{aligned}$$

- Maximum a posteriori estimation to evaluate



① joint pdf:

② Bayes Rule:

③ marginalize out the unobserved variable $f_{1:n}$

Maximum a Posteriori (MAP) Estimate of Hyperparameters

$$p(\theta, \sigma | D) \approx \delta(\theta - \theta^*) \delta(\sigma - \sigma^*)$$

where θ^*, σ^* maximize $\log p(\theta, \sigma | D)$

so $\log p(\theta, \sigma | D) =$