

Bayesian Linear Reg. w/ features

$$\phi(x) = \begin{bmatrix} x^2 \\ x \\ 1 \end{bmatrix}$$

$$w \sim \mathcal{N}(0, \Sigma_p)$$

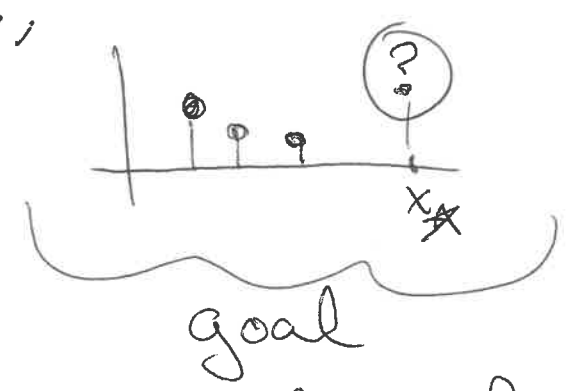
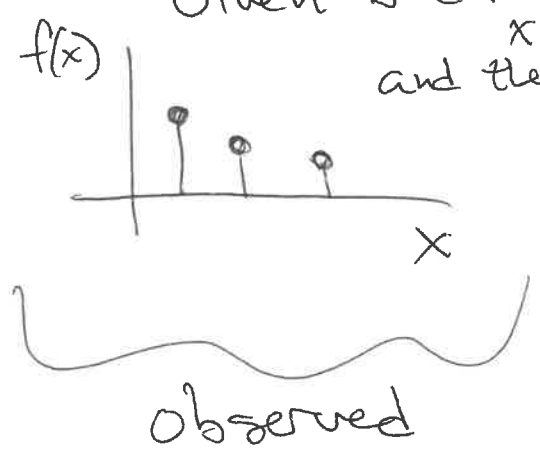
noise w/ variance  $\sigma_n^2$

$$y_n \leftarrow \underbrace{w^T \phi(x_n)}_{f(x_n)} + \epsilon_n$$

Day 2 01

Two ways to write predictive distrib.

Given 3 data points, what is function value at a new point  $x^*$  and their function values,



Look at posterior predictive (formulas from R&W)

$$f^* | x^*, X, y \sim \mathcal{N} \left( \frac{1}{\sigma_n^2} \phi^{*T} A^{-1} \Phi y, \phi^{*T} A^{-1} \phi^* \right)$$

$$A = \frac{1}{\sigma_n^2} \Phi \Phi^T + \Sigma_p^{-1}$$

above is feature view  
needs to invert  $F \times F$  matrix  $A$

$$F = \# \text{ feature dims in } \phi(x), \Phi = [\phi(x_1) \dots \phi(x_n)]$$

$$f^* | x^*, X, y \sim \mathcal{N} \left( k(x^*, X) [k(X, X) + \sigma_n^2 I_N]^{-1} y, k(x^*, x^*) - k(x^*, X) [k(X, X) + \sigma_n^2 I_N]^{-1} k(X, x^*) \right)$$

above is the kernel view  
needs to invert  $N \times N$  matrix  $\sigma_n^2 I + k(X, X)$

# Kernel Functions must be

Suppose  $x \in \mathbb{R}^2 = \begin{bmatrix} x_1 \\ x_2 \end{bmatrix}$

- (1) symmetric
- (2) pos. semi-definite

$$k(x, z) = (1 + x^T z)^2$$

$$= (1 + x_1 z_1 + x_2 z_2)^2$$

$$= 1 + x_1^2 z_1^2 + x_2^2 z_2^2 + 2x_1 z_1 + 2x_2 z_2 + 2x_1 x_2 z_1 z_2$$

Can we write as a dot-product  $\phi(x) \phi(z)$ ?

Yes, if  $\phi(x) = \begin{bmatrix} 1 \\ x_1 z \\ x_2 z \\ \sqrt{2} x_1 \\ \sqrt{2} x_2 \\ x_1 x_2 \end{bmatrix}$

$\phi(x)$  has size 6

so  $\phi^T \phi$  has cost  $O(6)$

but  $k(x, z)$  costs only  $O(2)$

$$k(x, z) = \exp\{- (x-z)^2\}$$

$$= \exp\{-x^2 - z^2 + 2xz\}$$

$$e^x = 1 + x + \frac{x^2}{2!} + \dots$$

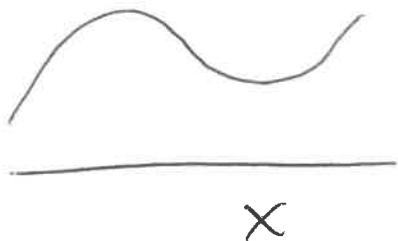
Taylor on  $e^{2xz}$

$$= \exp\{-x^2\} \exp\{-z^2\} \sum_{k=0}^{\infty} \frac{2^k}{k!} x^k z^k$$

$$= \sum_{k=0}^{\infty} \left( \sqrt{\frac{2^k}{k!}} e^{-x^2} x^k \right) \left( \sqrt{\frac{2^k}{k!}} e^{-z^2} z^k \right)$$

inner product

How do we represent functions?



Algebra

$$f(x) = x^2$$

most useful for GPs

Table

| x   | f(x) |
|-----|------|
| 1.1 | 4.38 |
| 2.1 | 2.16 |
| 3.3 | 8.11 |

each input has one output

GP says:

Given any finite set of inputs  $x_1, \dots, x_N$ , we can write  $f(x_1), \dots, f(x_N)$  as a vector

$$\vec{f} | \vec{x} \sim \mathcal{N}(m(\vec{x}), k(\vec{x}, \vec{x}))$$

and ask what is the joint distribution of this vector

$m(x)$  can be any function  $x \rightarrow \mathbb{R}$

$k(\cdot, \cdot)$  needs to specify valid covariance

must be positive semi-definite & symmetric

$k(x, x)$  is  $N \times N$  valid kernel matrix if

all eigenvalues  $\geq 0$

invertible (if  $> 0$ )

implies

Exercise:

With neighbor, come up with Python pseudocode to sample function values from a GP w/

mean  $m(x) = 1$

kernel  $k(x, x') = \sigma(x, x')$

# Useful Gaussian Properties

Day 2 04

Consider a joint distribution that is Gaussian

Marginalization yields Gaussians

$$\begin{bmatrix} X_A \\ X_B \end{bmatrix} \sim N \left( \begin{bmatrix} \mu_A \\ \mu_B \end{bmatrix}, \begin{bmatrix} \Sigma_{AA} & \Sigma_{AB} \\ \Sigma_{BA} & \Sigma_{BB} \end{bmatrix} \right)$$

$$X_A \sim N(\mu_A, \Sigma_{AA}) \quad p(X_A) = \int p(X_A, X_B) dx_B$$

$$X_B \sim N(\mu_B, \Sigma_{BB})$$

integral "marginalizes away"  $X_B$

Conditioning yields Gaussians

$$X_A | X_B \sim N \left( \mu_A + \Sigma_{AB} \Sigma_{BB}^{-1} (X_B - \mu_B), \right. \\ \left. \Sigma_{AA} - \Sigma_{AB} \Sigma_{BB}^{-1} \Sigma_{BA} \right)$$

Summation yields Gaussians if  $X_A$  same size as  $X_B$

$$X_A + X_B \sim N(\mu_A + \mu_B, \Sigma_{AA} + \Sigma_{BB})$$

Big Idea:

if you know joint distr. of  $\begin{bmatrix} X_A \\ X_B \end{bmatrix}$  is Gaussian,

there are super easy formulas to look up for

- conditional  $p(X_A | X_B)$
- marginal  $p(X_A)$
- or sum  $p(X_A + X_B)$

} all Gaussian

Weights space

prior is on weight vector  
of size  $D$  (#input dims)  
features

Func space

prior is directly on  
values of function

GP prediction

easier than Linear regr.  
to derive

$$\begin{bmatrix} \vec{y} \\ \vec{y}_* \end{bmatrix} \mid x, x_* \sim \mathcal{N} \left( \begin{bmatrix} 0 \\ 0 \end{bmatrix}, \begin{bmatrix} \sigma^2 + k(x, x) & k(x, x_*) \\ k(x, x_*) & \sigma^2 + k(x_*, x_*) \end{bmatrix} \right)$$

assumes  $m(x)=0$ . For general case, see R&W.

Use conditioning rule from useful Gaussian properties

$$\vec{y}_* \mid y, x, x_* \sim \mathcal{N}(\mu^*, \Sigma^*)$$

$$\mu^* = k(x_*, x) [k(x, x) + \sigma^2 I_N]^{-1} \vec{y}$$

$$\Sigma^* = k(x_*, x_*) + \sigma^2 I - k(x_*, x)^T [k(x, x) + \sigma^2 I]^{-1} k(x, x_*)$$