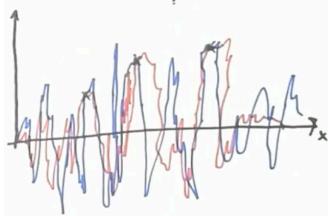


## Week 5a

Underfitting, model very certain,  
 check underfitting: observe residuals all positive -> all negative -> all positive;  
 observe residual highly correlated

problem similar to overfitting: model too complex (with lots of narrow RBF)  
 & not enough data  
 => posterior distribution very uncertain (similar to prior except at few data points)

e.g.



**Model choice:** represent belief about which model the features came from  
 e.g. narrow VS wide distribution (= lots of narrow RBFs); hyperparameters

marginal likelihood  $p(\mathbf{Y}_{\text{train}} | \mathbf{X}_{\text{train}}, \mathbf{M}) = p(\mathbf{D} | \mathbf{M})$

marginalized out weight  $\mathbf{w}$

$\mathbf{M}$  = model choice (implicitly always there): contains noise variance, prior mean, prior covariance (RBF position & width), activation functions

$x, y$ : training instance

Use sum & product rule, but can't solve the integral

$$p(\mathbf{y} | \mathbf{X}, \mathcal{M}) = \int p(\mathbf{y}, \mathbf{w} | \mathbf{X}, \mathcal{M}) d\mathbf{w} = \int p(\mathbf{y} | \mathbf{X}, \mathbf{w}, \mathcal{M}) p(\mathbf{w} | \mathcal{M}) d\mathbf{w},$$

Likelihood.      prior

Trick to get  $p(y | x, M)$ : use posterior weight  $p(\mathbf{w} | \mathbf{D})$  (= LHS) mentioned before:  
 denominator is marginal likelihood, nominator is standard likelihood & prior  
 solve marginal likelihood by rearranging terms  
 could use any  $\mathbf{w}$  (if  $\mathbf{w}=\mathbf{w}$  could cancel out terms), get the same answer

posterior weights:

$$p(\mathbf{w} | \mathbf{y}, \mathbf{X}, \mathcal{M}) = \frac{p(\mathbf{y} | \mathbf{w}, \mathbf{X}, \mathcal{M}) \cdot p(\mathbf{w} | \mathcal{M})}{p(\mathbf{y} | \mathbf{X}, \mathcal{M})}$$

Apply marginal likelihood (after having it):

apply Bayes rule, get  $p(D | M)$

or apply maximum marginal likelihood:

model with prior

$$p(\mathbf{w} | \sigma_w) = \mathcal{N}(\mathbf{w}; \mathbf{0}, \sigma_w^2 \mathbb{I}), \quad (2)$$

and likelihood

$$p(y | \mathbf{x}, \mathbf{w}, \sigma_y) = \mathcal{N}(y; \mathbf{w}^\top \mathbf{x}, \sigma_y^2), \quad (3)$$

we can fit the *hyperparameters*  $\sigma_w$  and  $\sigma_y$ , parameters which specify the model, to maximize their marginal likelihood:

$$p(\mathbf{y} | X, \sigma_w, \sigma_y) = \int p(\mathbf{y}, \mathbf{w} | X, \sigma_w, \sigma_y) d\mathbf{w} = \int p(\mathbf{y} | X, \mathbf{w}, \sigma_y) p(\mathbf{w} | \sigma_w) d\mathbf{w}. \quad (4)$$

To choose hyperparameters of M (**on training set only!!!**, since bayesian)

e.g. marginal likelihood of good model >> overfitting model (with narrow basis functions), also >> underfitting model

Since: Like the dice example, one model could only generate smooth curve, overfitting model could generate smooth & overfitted curve, marginal likelihood is much smaller

## Week5b

**Bayesian optimization:** use uncertainty to guide search (do the next experiment in highly uncertain area)

task: e.g. have uncertainty, find greatest y value (acquisition function)

limitation of bay linear regression: too certain if model too simple

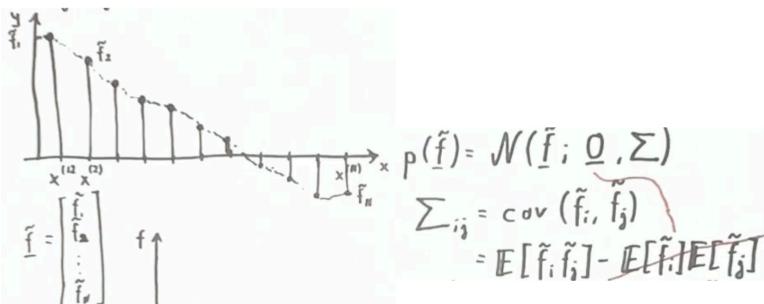
("underfitting")

(e.g. when move x to infinity, y of [increasing regression line & 3RBFs])

## Gaussian Process prior

Discretized function, represented by huge vector ( $\tilde{\mathbf{f}}$ , many points)

discrete points drawn from normal distribution, defined mainly by (very big) covariance matrix



set prior entries in covariance: defined by **kernel function** (= relation between 2 points  $f_i$   $f_j$ );

$x_i$   $x_j$  could be multivariate (for a surface)

kernel: semi-positive definite and symmetric

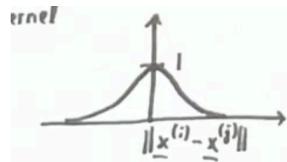
$$\text{cov}[\tilde{f}_i, \tilde{f}_j] = k(\mathbf{x}^{(i)}, \mathbf{x}^{(j)}).$$

functions for  $k$ : **Mercer kernel**, positive semi definite

$$k(\mathbf{x}^{(i)}, \mathbf{x}^{(j)}) = \exp(-\|\mathbf{x}^{(i)} - \mathbf{x}^{(j)}\|^2),$$

want smooth curve, so want neighboring points highly correlated, distant points less correlated

kernel function:



## Gaussian Process regression

Gaussian properties:

in covariance, A: cov just consider  $f$ ; B: cov just consider  $g$ ;

C and  $C^T$ : how  $f$  and  $g$  related

### Properties of Gaussians

$$\text{For a joint Gaussian } p(f, g) = \mathcal{N}\left(\begin{bmatrix} f \\ g \end{bmatrix}; \begin{bmatrix} a \\ b \end{bmatrix}, \begin{bmatrix} A & C \\ C^T & B \end{bmatrix}\right)$$

$D_1 \times D_1$     $D_2 \times D_2$     $D_1 \times D_2$   
 $D_1 + D_2 - d_{in}$     $D_1 + D_2 - d_{in}$     $D_2 \times D_2$

Marginals (sum rule):

$$\begin{aligned} p(f) &= \int p(f, g) dg \\ &= \mathcal{N}(f; a, A) \end{aligned}$$

Conditionals:

$$p(f|g) = \mathcal{N}(f; a + CB^{-1}(g-b), A - CB^{-1}C^T)$$

$$p(g|f) = \mathcal{N}(g; b + C^TA^{-1}(f-a), B - C^TA^{-1}C)$$

Terms:

- y: observed training labels
- $\tilde{f}$ : real training labels
- $\tilde{f}^*$ : real test labels
- $x, x^*$ : training location  $x$  & test location  $x^*$

have prior  $f \sim GP(k)$ ,

have noisy observations  $y$  (noise around function values  $\tilde{f}_n$ , as equation below),  
could update our belief (by likelihood)  
for n datapoints, product over n

$$y_n \sim \mathcal{N}(\tilde{f}_n, \sigma_y^2)$$

$$\Rightarrow \text{likelihood } p(y_n | \tilde{f}) = \mathcal{N}(y_n; \tilde{f}_n, \sigma_y^2)$$

have unknown true function, joint distribution of  $\tilde{f}$  and  $\tilde{f}^*$  :  $p(\tilde{f}, \tilde{f}^*)$

D: size of input vector  $x$

N: # of training locations

M: # of test locations

$$p(\tilde{f}, \tilde{f}^*) = \mathcal{N}\left(\begin{bmatrix} \tilde{f} \\ \tilde{f}^* \end{bmatrix}; \begin{bmatrix} 0 \\ 0 \end{bmatrix}, \begin{bmatrix} K(X, X) & K(X, X_*) \\ K(X_*, X) & K(X_*, X_*) \end{bmatrix}\right)$$

$$K(X, Z)_{ij} = k(x^{(i)}, z^{(j)})$$

**Prior:**

joint distribution of the observations  $y$  and  $\tilde{f}^*$  :  $p(y, \tilde{f}^*)$

(since  $\tilde{f}^*$  only related to  $\tilde{f}$ ;  $y$  is just observation noise on top of  $\tilde{f}$ )

$$p\left(\begin{bmatrix} y \\ f_* \end{bmatrix}\right) = \mathcal{N}\left(\begin{bmatrix} y \\ f_* \end{bmatrix}; 0, \begin{bmatrix} K(X, X) + \sigma_y^2 I & K(X, X_*) \\ K(X_*, X) & K(X_*, X_*) \end{bmatrix}\right),$$

= additional observation noise

Do **inference** by GP regression:

$p(\tilde{f}^* | y)$ : same as  $p(Y_{\text{test}} | X_{\text{test}}, D)$

just copy values from  $p(y, \tilde{f}^*)$  according to Gaussian properties above

**Advantage of GP:** just copy values from  $p(g | f)$ , no calculation needed

Inference

$$p(\tilde{f}_* | y) = \mathcal{N}(\tilde{f}_*; \underline{\quad}, \underline{\quad})$$

"empty values":

$$\text{mean}, \bar{\mathbf{f}}_* = K(X_*, X)(K(X, X) + \sigma_y^2 \mathbb{I})^{-1} \mathbf{y}$$

$$\text{cov}(\mathbf{f}_*) = K(X_*, X_*) - K(X_*, X)(K(X, X) + \sigma_y^2 \mathbb{I})^{-1} K(X, X_*)$$

cov matrix depends on input location  $x$ , but does not depend on training label  $y$   
so a surprising training label have little effect

Can we be more uncertain at prediction in response to a surprising training label?

=> change GP kernel function, learn the parameters in kernel (choose hyperparameter)

Visualizing the posterior