

# Probabilistic Graphical Models

## Lecture 20: Gaussian Processes

Andrew Gordon Wilson

[www.cs.cmu.edu/~andrewgw](http://www.cs.cmu.edu/~andrewgw)  
Carnegie Mellon University



# What is Machine Learning?

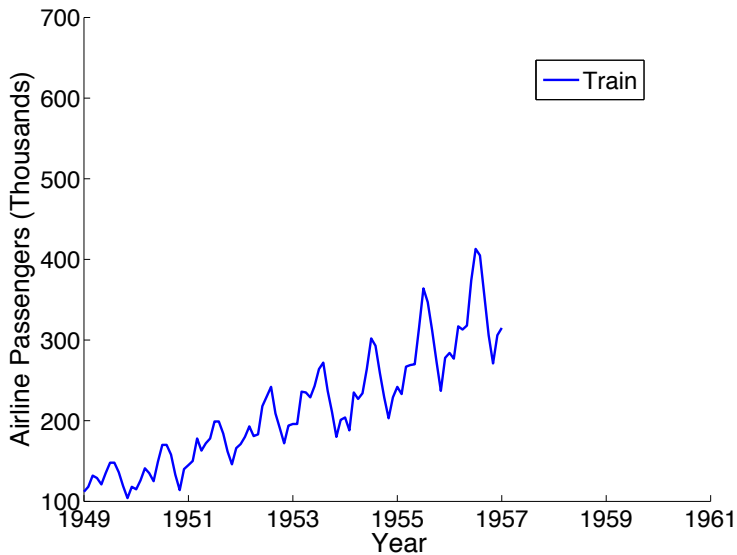
- ▶ Machine learning algorithms adapt with data versus having fixed decision rules.
- ▶ Machine learning aims not only to equip people with tools to analyse data, but to create algorithms which can learn and make decisions without human intervention.<sup>1,2</sup>
- ▶ In order for a model to automatically learn and make decisions, it must be able to discover patterns and extrapolate those patterns to new situations.

---

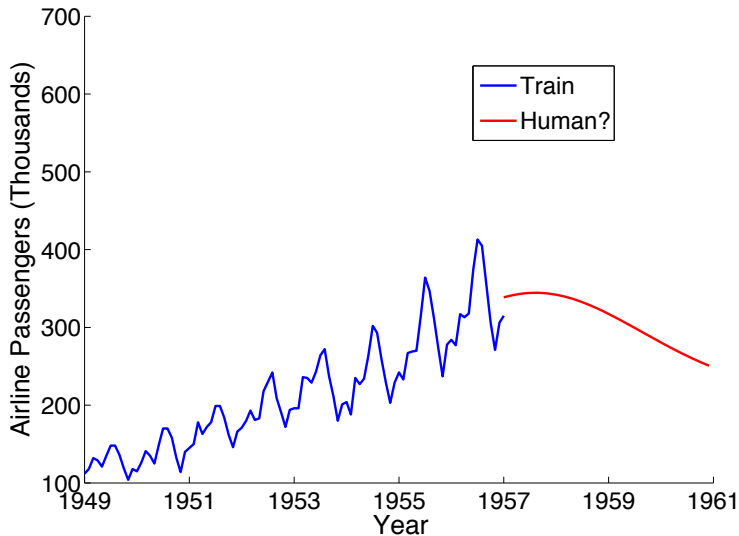
<sup>1</sup>E.g., N.D. Lawrence (2010), “What is Machine Learning?”

<sup>2</sup>T.M. Mitchell (2006), “What is Machine Learning and Where Is it Headed?”

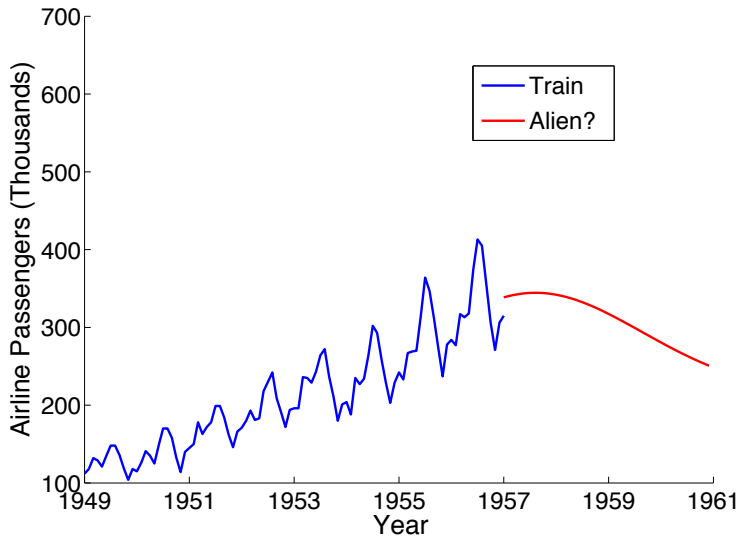
# Function Learning Example



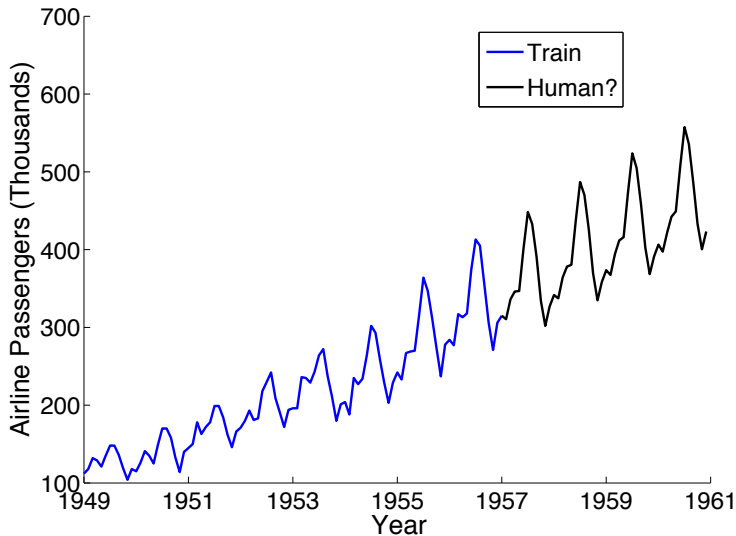
# Function Learning Example



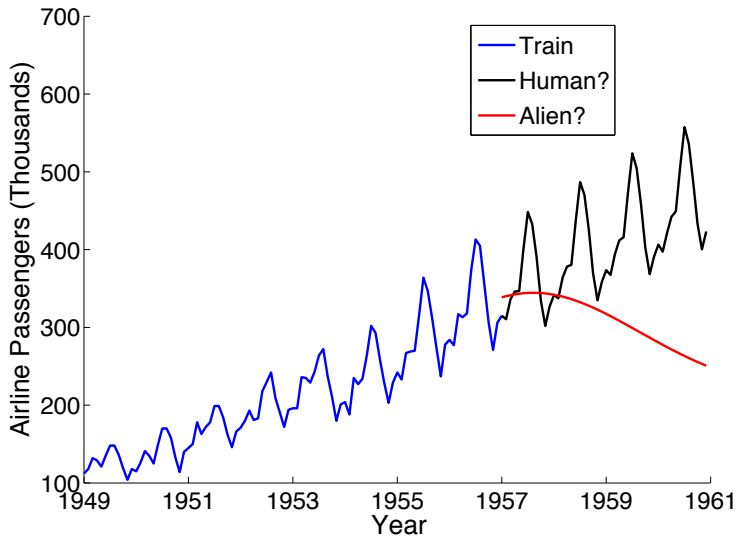
# Function Learning Example



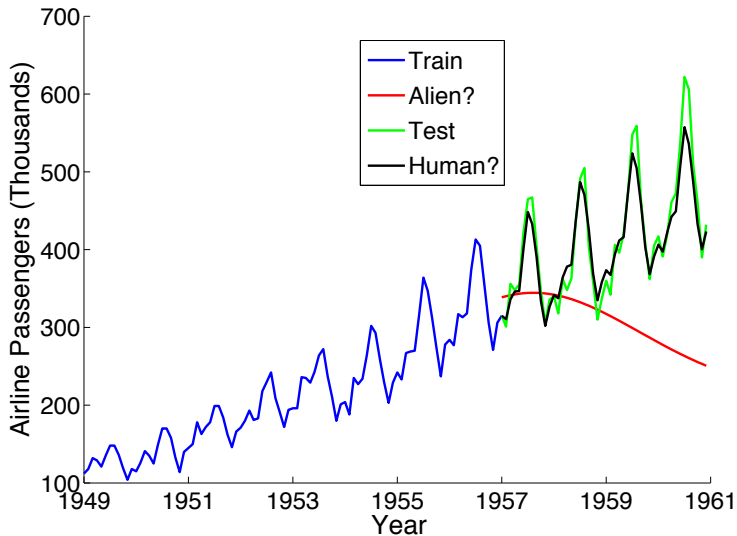
# Function Learning Example



# Function Learning Example



# Function Learning Example





# Building an Intelligent Model

The ability for a model to learn from data depends on its:

1. Support: what solutions we think are a priori possible.
2. Inductive biases: what solutions we think are a priori likely.
  - ▶ Examples: Function Learning, Character Recognition
  - ▶ Human ability to make remarkable generalisations from data could derive from an expressive prior combined with Bayesian inference.

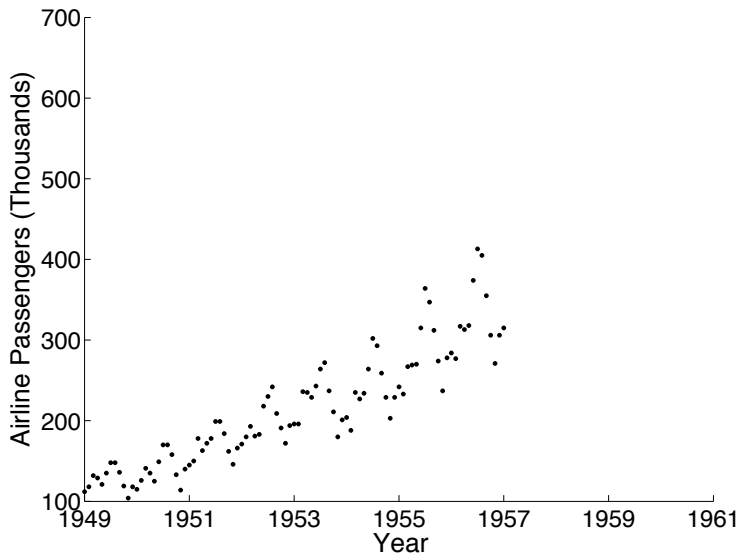
## Basic Regression Problem

- ▶ Training set of  $N$  targets (observations)  $\mathbf{y} = (y(x_1), \dots, y(x_N))^T$ .
- ▶ Observations evaluated at inputs  $X = (x_1, \dots, x_N)^T$ .
- ▶ Want to predict the value of  $y(x_*)$  at a test input  $x_*$ .

For example: Given CO<sub>2</sub> concentrations  $\mathbf{y}$  measured at times  $X$ , what will the CO<sub>2</sub> concentration be for  $x_* = 2024$ , 10 years from now?

Just knowing high school math, what might you try?

# Statistics from Scratch



Guess the parametric form of a function that could fit the data

- ▶  $f(x, \mathbf{w}) = \mathbf{w}^T x$  [Linear function of  $\mathbf{w}$  and  $x$ ]
- ▶  $f(x, \mathbf{w}) = \mathbf{w}^T \phi(x)$  [Linear function of  $\mathbf{w}$ ] (Linear Basis Function Model)
- ▶  $f(x, \mathbf{w}) = g(\mathbf{w}^T \phi(x))$  [Non-linear in  $x$  and  $\mathbf{w}$ ] (E.g., Neural Network)

$\phi(x)$  is a vector of basis functions. For example, if  $\phi(x) = (1, x, x^2)$  and  $x \in \mathbb{R}^1$  then  $f(x, \mathbf{w}) = w_0 + w_1x + w_2x^2$  is a quadratic function.

Choose an error measure  $E(\mathbf{w})$ , minimize with respect to  $\mathbf{w}$

- ▶  $E(\mathbf{w}) = \sum_{i=1}^N [f(x_i, \mathbf{w}) - y(x_i)]^2$

## A probabilistic approach

We could explicitly account for noise in our model.

- ▶  $y(x) = f(x, \mathbf{w}) + \epsilon(x)$ , where  $\epsilon(x)$  is a noise function.

One commonly takes  $\epsilon(x) = \mathcal{N}(0, \sigma^2)$  for i.i.d. additive Gaussian noise, in which case

$$p(y(x)|x, \mathbf{w}, \sigma^2) = \mathcal{N}(y(x); f(x, \mathbf{w}), \sigma^2) \quad \text{Observation Model} \quad (1)$$

$$p(\mathbf{y}|x, \mathbf{w}, \sigma^2) = \prod_{i=1}^N \mathcal{N}(y(x_i); f(x_i, \mathbf{w}), \sigma^2) \quad \text{Likelihood} \quad (2)$$

- ▶ Maximize the likelihood of the data  $p(\mathbf{y}|x, \mathbf{w}, \sigma^2)$  with respect to  $\sigma^2, \mathbf{w}$ .

For a Gaussian noise model, this approach will make the same predictions as using a squared loss error function:

$$\log p(\mathbf{y}|X, \mathbf{w}, \sigma^2) \propto -\frac{1}{2\sigma^2} \sum_{i=1}^N [f(x_i, \mathbf{w}) - y(x_i)]^2 \quad (3)$$

# Statistics from Scratch

- ▶ The probabilistic approach helps us interpret the error measure in a deterministic approach, and gives us a sense of the noise level  $\sigma^2$ .
- ▶ Probabilistic methods thus provide an intuitive framework for representing uncertainty, and model development.
- ▶ Both approaches are prone to *over-fitting* for flexible  $f(x, \mathbf{w})$ : low error on the training data, high error on the test set.

## Regularization

- ▶ Use a penalized log likelihood (or error function), such as

$$\log p(\mathbf{y}|X, \mathbf{w}) \propto \underbrace{-\frac{1}{2\sigma^2} \sum_{i=1}^n (f(x_i, \mathbf{w}) - y(x_i))^2}_{\text{model fit}} \underbrace{-\lambda \mathbf{w}^T \mathbf{w}}_{\text{complexity penalty}}. \quad (4)$$

- ▶ **But how should we define complexity, and how much should we penalize complexity?**
- ▶ Can set  $\lambda$  using *cross-validation*.

## Bayes' Rule

$$p(a|b) = p(b|a)p(a)/p(b), \quad p(a|b) \propto p(b|a)p(a). \quad (5)$$

$$\text{posterior} = \frac{\text{likelihood} \times \text{prior}}{\text{marginal likelihood}}, \quad p(\mathbf{w}|\mathbf{y}, X, \sigma^2) = \frac{p(\mathbf{y}|X, \mathbf{w}, \sigma^2)p(\mathbf{w})}{p(\mathbf{y}|X, \sigma^2)}. \quad (6)$$

## Predictive Distribution

$$p(y|x_*, \mathbf{y}, X) = \int p(y|x_*, \mathbf{w})p(\mathbf{w}|\mathbf{y}, X)d\mathbf{w}. \quad (7)$$

- ▶ Average of infinitely many models weighted by their posterior probabilities.
- ▶ No over-fitting, automatically calibrated complexity.
- ▶ Typically more interested in distribution over functions than in parameters  $\mathbf{w}$ .

# Representing Uncertainty

Different types of uncertainty:

- ▶ Uncertainty through lack of knowledge
- ▶ Intrinsic uncertainty; e.g., radioactive decay.

Uncertainty through lack of knowledge can seem like intrinsic uncertainty (e.g., rolling dice).

Regardless of whether or not the universe is deterministic – whether there is some underlying true answer – we will always have uncertainty. We can represent this belief using probability distributions (Bayesian methods, probabilistic modelling).



# Parametric Regression Review

## Deterministic

$$E(\mathbf{w}) = \sum_{i=1}^N (f(x_i, \mathbf{w}) - y_i)^2. \quad (8)$$

## Maximum Likelihood

$$p(y(x)|x, \mathbf{w}) = \mathcal{N}(y(x); f(x, \mathbf{w}), \sigma_n^2), \quad (9)$$

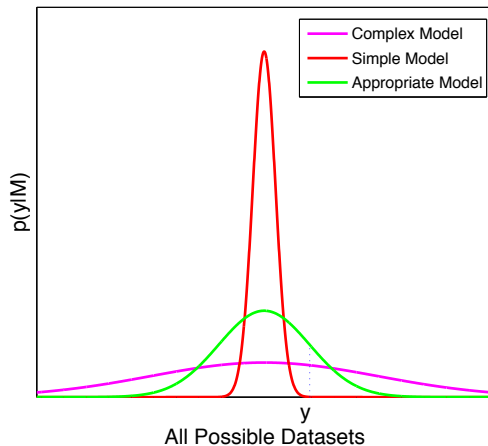
$$p(\mathbf{y}|X, \mathbf{w}) = \prod_{i=1}^N \mathcal{N}(y(x_i); f(x_i, \mathbf{w}), \sigma_n^2). \quad (10)$$

## Bayesian

$$\text{posterior} = \frac{\text{likelihood} \times \text{prior}}{\text{marginal likelihood}}, \quad p(\mathbf{w}|\mathbf{y}, X) = \frac{p(\mathbf{y}|X, \mathbf{w})p(\mathbf{w})}{p(\mathbf{y}|X)}. \quad (11)$$

# Model Selection and Marginal Likelihood

$$p(\mathbf{y}|\mathcal{M}_1, X) = \int p(\mathbf{y}|f_1(x, \mathbf{w}))p(\mathbf{w})d\mathbf{w} \quad (13)$$



## Blackboard: Examples of Occam's Razor in Everyday Inferences

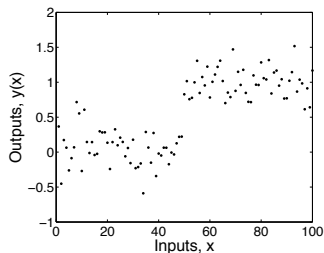
For further reading, see MacKay (2003) textbook, *Information Theory, Inference, and Learning Algorithms*.

## Occam's Razor Example

-1, 3, 7, 11, ??, ??

- ▶  $H_1$ : the sequence is an arithmetic progression, add  $n$ , where  $n$  is an integer.
- ▶  $H_2$ : the sequence is generated by a cubic function of the form  $cx^3 + dx^2 + e$ , where  $c$ ,  $d$ , and  $e$  are fractions.  $(-\frac{1}{11}x^3 + \frac{9}{11}x^2 + \frac{23}{11})$

# Model Selection



Observations  $y(x)$ . Assume  $p(y(x)|f(x)) \sim \mathcal{N}(y(x); f(x), \sigma^2)$ . Consider polynomials of different orders. As always, observations are out of the chosen model class! Which model should we choose?

$$f_0(x) = a_0, \quad (14)$$

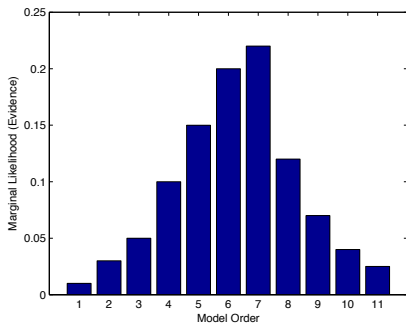
$$f_1(x) = a_0 + a_1x, \quad (15)$$

$$f_2(x) = a_0 + a_1x + a_2x^2, \quad (16)$$

$$\vdots \quad (17)$$

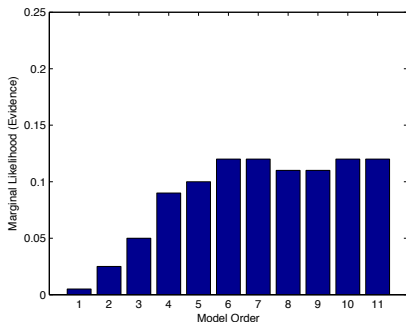
$$f_J(x) = a_0 + a_1x + a_2x^2 + \cdots + a_Jx^J. \quad (18)$$

# Model Selection: Occam's Hill



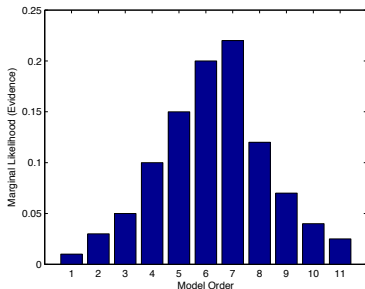
Marginal likelihood (evidence) as a function of model order, using an isotropic prior  $p(a) = \mathcal{N}(0, \sigma^2 I)$ .

# Model Selection: Occam's Asymptote

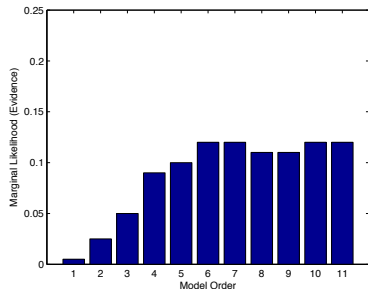


Marginal likelihood (evidence) as a function of model order, using an anisotropic prior  $p(a_i) = \mathcal{N}(0, \gamma^{-i})$ , with  $\gamma$  learned from the data.

# Occam's Razor



(a) Isotropic Gaussian Prior



(b) Anisotropic Gaussian Prior

For further reading, see Rasmussen and Ghahramani (2001) (*Occam's Razor*) and Kass and Raftery (1995) (*Bayes Factors*)

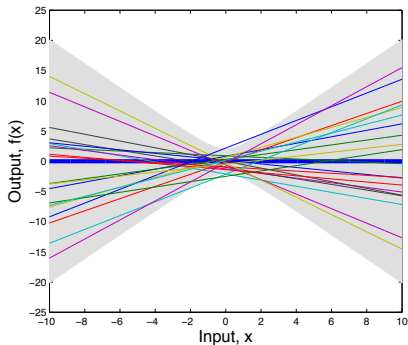


# Linear Basis Models

Consider the simple linear model,

$$f(x) = a_0 + a_1x, \quad (19)$$

$$a_0, a_1 \sim \mathcal{N}(0, 1). \quad (20)$$



# Linear Models

We are interested in the induced distribution over functions, not the parameters...

Let's characterise the properties of these functions directly:

$$f(x|a_0, a_1) = a_0 + a_1x, \quad a_0, a_1 \sim \mathcal{N}(0, 1). \quad (21)$$

$$\mathbb{E}[f(x)] = \mathbb{E}[a_0] + \mathbb{E}[a_1]x = 0. \quad (22)$$

$$\text{cov}[f(x_b), f(x_c)] = \mathbb{E}[f(x_b)f(x_c)] - \mathbb{E}[f(x_b)]\mathbb{E}[f(x_c)] \quad (23)$$

$$= \mathbb{E}[a_0^2 + a_0a_1(x_b + x_c) + a_1^2x_bx_c] - 0 \quad (24)$$

$$= \mathbb{E}[a_0^2] + \mathbb{E}[a_1^2x_bx_c] + \mathbb{E}[a_0a_1(x_b + x_c)] \quad (25)$$

$$= 1 + x_bx_c + 0 \quad (26)$$

$$= 1 + x_bx_c. \quad (27)$$

# Linear Models

Therefore any collection of values has a joint Gaussian distribution

$$[f(x_1), \dots, f(x_N)] \sim \mathcal{N}(\mathbf{0}, \mathbf{K}), \quad (28)$$

$$K_{ij} = \text{cov}(f(x_i), f(x_j)) = k(x_i, x_j) = 1 + x_b x_c. \quad (29)$$

By definition,  $f(x)$  is a Gaussian process.

## Definition

A Gaussian process (GP) is a collection of random variables, any finite number of which have a joint Gaussian distribution. We write

$f(x) \sim \mathcal{GP}(m, k)$  to mean

$$[f(x_1), \dots, f(x_N)] \sim \mathcal{N}(\boldsymbol{\mu}, \mathbf{K}) \quad (30)$$

$$\boldsymbol{\mu}_i = m(x_i) \quad (31)$$

$$K_{ij} = k(x_i, x_j), \quad (32)$$

for any collection of input values  $x_1, \dots, x_N$ . In other words,  $f$  is a GP with mean function  $m(x)$  and *covariance kernel*  $k(x_i, x_j)$ .

# Linear Basis Function Models

## Model Specification

$$f(x, \mathbf{w}) = \mathbf{w}^T \phi(x) \quad (33)$$

$$p(\mathbf{w}) = \mathcal{N}(0, \Sigma_w) \quad (34)$$

## Moments of Induced Distribution over Functions

$$\mathbb{E}[f(x, \mathbf{w})] = m(x) = \mathbb{E}[\mathbf{w}^T] \phi(x) = 0 \quad (35)$$

$$\text{cov}(f(x_i), f(x_j)) = k(x_i, x_j) = \mathbb{E}[f(x_i)f(x_j)] - \mathbb{E}[f(x_i)]\mathbb{E}[f(x_j)] \quad (36)$$

$$= \phi(x_i)^T \mathbb{E}[\mathbf{w}\mathbf{w}^T] \phi(x_j) - 0 \quad (37)$$

$$= \phi(x_i)^T \Sigma_w \phi(x_j) \quad (38)$$

- ▶  $f(x, \mathbf{w})$  is a Gaussian process,  $f(x) \sim \mathcal{N}(m, k)$  with mean function  $m(x) = 0$  and covariance kernel  $k(x_i, x_j) = \phi(x_i)^T \Sigma_w \phi(x_j)$ .
- ▶ The entire basis function model of Eqs. (33) and (34) is encapsulated as a distribution over functions with kernel  $k(x, x')$ .

- ▶ We are ultimately more interested in – and have stronger intuitions about – the *functions* that model our data than weights  $\mathbf{w}$  in a parametric model, and we can express those intuitions using a covariance kernel.
- ▶ The kernel controls the support and inductive biases of our model, and thus its ability to generalise.

## Example: RBF Kernel

$$k_{\text{RBF}}(x, x') = \text{cov}(f(x), f(x')) = a^2 \exp\left(-\frac{\|x - x'\|^2}{2\ell^2}\right) \quad (39)$$

- ▶ Far and above the most popular kernel.
- ▶ Expresses the intuition that function values at nearby inputs are more correlated than function values at far away inputs.
- ▶ The kernel *hyperparameters*  $a$  and  $\ell$  control amplitudes and wiggleness of these functions.
- ▶ GPs with an RBF kernel have large support and are *universal approximators*.

# Sampling from a GP with an RBF Kernel

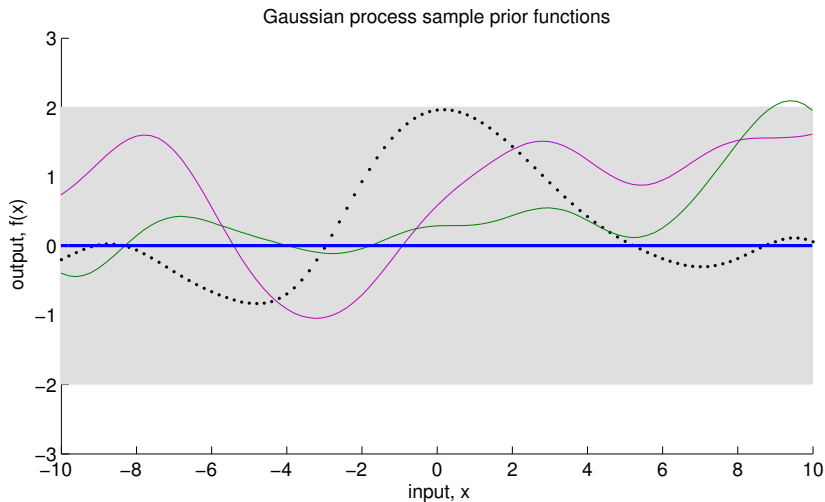
```
x = [-10:0.2:10]'; % inputs (where we query the GP)
N = numel(x); % number of inputs
K = zeros(N,N); % covariance matrix

% very inefficient way of creating K in Matlab
for i=1:N
    for j=1:N
        K(i,j) = k_rbf(x(i),x(j));
    end
end

K = K + 1e-6*eye(N); % add jitter for conditioning
CK = chol(K);
f = CK'*randn(N,1); % draws from N(0,K)

plot(x,f);
```

# Samples from a GP with an RBF Kernel





# 1D RBF Kernel with Different Length-scales

$$k_{\text{RBF}}(x, x') = \text{cov}(f(x), f(x')) = a^2 \exp\left(-\frac{\|x - x'\|^2}{2\ell^2}\right) \quad (40)$$

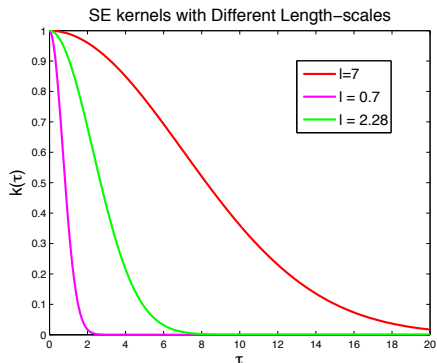
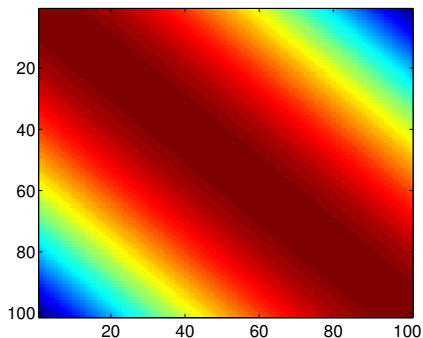


Figure: SE kernels with different length-scales, as a function of  $\tau = x - x'$ .

# RBF Kernel Covariance Matrix

$$k_{\text{RBF}}(x, x') = \text{cov}(f(x), f(x')) = a^2 \exp\left(-\frac{\|x - x'\|^2}{2\ell^2}\right) \quad (41)$$

The covariance matrix  $K$  for ordered inputs on a 1D grid.  $K_{ij} = k_{\text{RBF}}(x_i, x_j)$ .



# Gaussian Process Inference

- ▶ Observed noisy data  $\mathbf{y} = (y(x_1), \dots, y(x_N))^T$  at input locations  $X$ .
- ▶ Start with the standard regression assumption:  $\mathcal{N}(y(x); f(x), \sigma^2)$ .
- ▶ Place a Gaussian process distribution over noise free functions  $f(x) \sim \mathcal{GP}(0, k_\theta)$ . The kernel  $k$  is parametrized by  $\theta$ .
- ▶ Infer  $p(\mathbf{f}_* | \mathbf{y}, X, X_*)$  for the noise free function  $f$  evaluated at test points  $X_*$ .

## Joint distribution

$$\begin{bmatrix} \mathbf{y} \\ \mathbf{f}_* \end{bmatrix} \sim \mathcal{N} \left( \mathbf{0}, \begin{bmatrix} K_\theta(X, X) + \sigma^2 I & K_\theta(X, X_*) \\ K_\theta(X_*, X) & K_\theta(X_*, X_*) \end{bmatrix} \right). \quad (42)$$

## Conditional predictive distribution

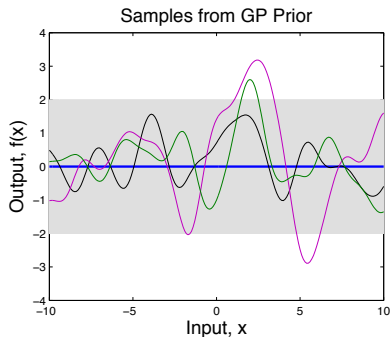
$$\mathbf{f}_* | X_*, X, \mathbf{y}, \theta \sim \mathcal{N}(\bar{\mathbf{f}}_*, \text{cov}(\mathbf{f}_*)), \quad (43)$$

$$\bar{\mathbf{f}}_* = K_\theta(X_*, X) [K_\theta(X, X) + \sigma^2 I]^{-1} \mathbf{y}, \quad (44)$$

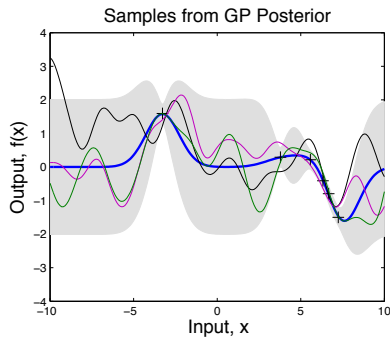
$$\text{cov}(\mathbf{f}_*) = K_\theta(X_*, X_*) - K_\theta(X_*, X) [K_\theta(X, X) + \sigma^2 I]^{-1} K_\theta(X, X_*). \quad (45)$$

# Inference using an RBF kernel

- ▶ Specify  $f(x) \sim \mathcal{GP}(0, k)$ .
- ▶ Choose  $k_{\text{RBF}}(x, x') = a_0^2 \exp(-\frac{\|x-x'\|^2}{2\ell_0^2})$ . Choose values for  $a_0$  and  $\ell_0$ .
- ▶ Observe data, look at the prior and posterior over functions.



(a)

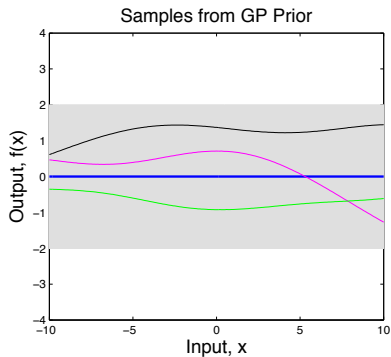


(b)

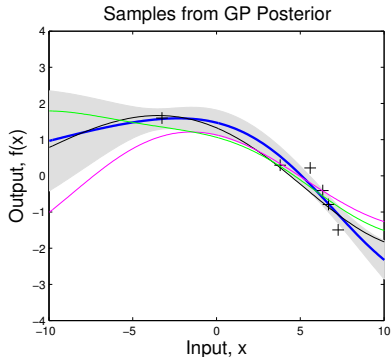
- ▶ Does something look strange about these functions?

# Inference using an RBF kernel

Increase the length-scale  $\ell$ .



(a)



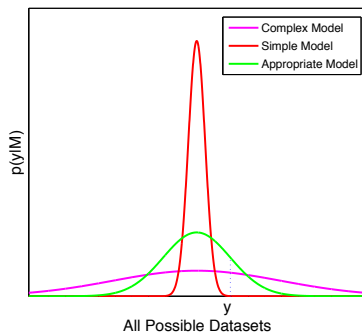
(b)

# Learning and Model Selection

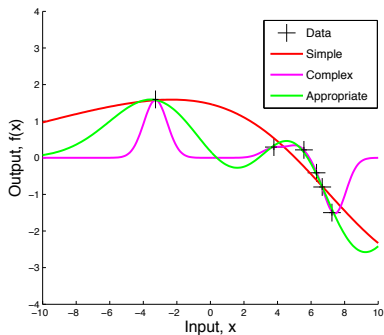
$$p(\mathcal{M}_i|\mathbf{y}) = \frac{p(\mathbf{y}|\mathcal{M}_i)p(\mathcal{M}_i)}{p(\mathbf{y})} \quad (46)$$

We can write the *evidence* of the model as

$$p(\mathbf{y}|\mathcal{M}_i) = \int p(\mathbf{y}|\mathbf{f}, \mathcal{M}_i)p(\mathbf{f})d\mathbf{f}, \quad (47)$$



(a)



(b)

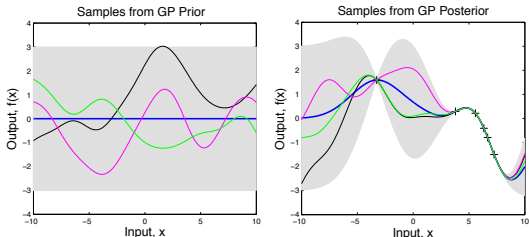
# Learning and Model Selection

- ▶ We can integrate away the entire Gaussian process  $f(x)$  to obtain the marginal likelihood, as a function of kernel hyperparameters  $\theta$  alone.

$$p(\mathbf{y}|\boldsymbol{\theta}, X) = \int p(\mathbf{y}|\mathbf{f}, X)p(\mathbf{f}|\boldsymbol{\theta}, X)d\mathbf{f}. \quad (48)$$

$$\log p(\mathbf{y}|\boldsymbol{\theta}, X) = \underbrace{-\frac{1}{2}\mathbf{y}^T(K_{\boldsymbol{\theta}} + \sigma^2\mathbf{I})^{-1}\mathbf{y}}_{\text{model fit}} - \underbrace{\frac{1}{2}\log |K_{\boldsymbol{\theta}} + \sigma^2\mathbf{I}|}_{\text{complexity penalty}} - \frac{N}{2}\log(2\pi). \quad (49)$$

- ▶ An extremely powerful mechanism for kernel learning.



# Learning and Model Selection

- ▶ A fully Bayesian treatment would integrate away kernel hyperparameters  $\theta$ .

$$p(\mathbf{f}_* | X_*, X, \mathbf{y}) = \int p(\mathbf{f}_* | X_*, X, \mathbf{y}, \boldsymbol{\theta}) p(\boldsymbol{\theta} | \mathbf{y}) d\boldsymbol{\theta} \quad (50)$$

- ▶ For example, we could specify a prior  $p(\boldsymbol{\theta})$ , use MCMC to take  $J$  samples from  $p(\boldsymbol{\theta} | \mathbf{y}) \propto p(\mathbf{y} | \boldsymbol{\theta}) p(\boldsymbol{\theta})$ , and then find

$$p(\mathbf{f}_* | X_*, X, \mathbf{y}) \approx \frac{1}{J} \sum_{i=1}^J p(\mathbf{f}_* | X_*, X, \mathbf{y}, \boldsymbol{\theta}^{(i)}), \quad \boldsymbol{\theta}^{(i)} \sim p(\boldsymbol{\theta} | \mathbf{y}). \quad (51)$$

- ▶ If we have a non-Gaussian noise model, and thus cannot integrate away  $\mathbf{f}$ , the strong dependencies between Gaussian process  $\mathbf{f}$  and hyperparameters  $\boldsymbol{\theta}$  make sampling extremely difficult. In my experience, the most effective solution is to use a deterministic approximation for the posterior  $p(\mathbf{f} | \mathbf{y})$  which enables one to work with an approximate marginal likelihood.



# Gaussian Process Covariance Kernels

Let  $\tau = x - x'$ :

$$k_{\text{SE}}(\tau) = \exp(-0.5\tau^2/\ell^2) \quad (52)$$

$$k_{\text{MA}}(\tau) = a\left(1 + \frac{\sqrt{3}\tau}{\ell}\right) \exp\left(-\frac{\sqrt{3}\tau}{\ell}\right) \quad (53)$$

$$k_{\text{RQ}}(\tau) = \left(1 + \frac{\tau^2}{2\alpha\ell^2}\right)^{-\alpha} \quad (54)$$

$$k_{\text{PE}}(\tau) = \exp(-2 \sin^2(\pi \tau \omega)/\ell^2) \quad (55)$$

1. **Learning:** Optimize marginal likelihood,

$$\log p(\mathbf{y}|\boldsymbol{\theta}, X) = \underbrace{-\frac{1}{2}\mathbf{y}^T(K_{\boldsymbol{\theta}} + \sigma^2 I)^{-1}\mathbf{y}}_{\text{model fit}} - \underbrace{\frac{1}{2}\log |K_{\boldsymbol{\theta}} + \sigma^2 I|}_{\text{complexity penalty}} - \frac{N}{2}\log(2\pi),$$

with respect to kernel hyperparameters  $\boldsymbol{\theta}$ .

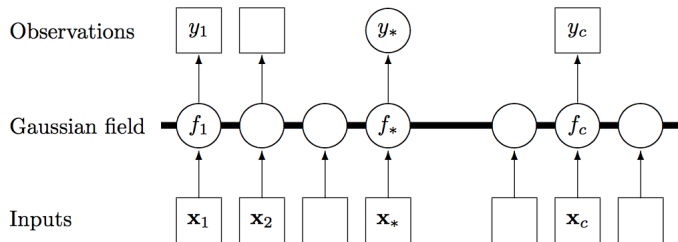
2. **Inference:** Conditioned on kernel hyperparameters  $\boldsymbol{\theta}$ , form the predictive distribution for test inputs  $X_*$ :

$$\mathbf{f}_*|X_*, X, \mathbf{y}, \boldsymbol{\theta} \sim \mathcal{N}(\bar{\mathbf{f}}_*, \text{cov}(\mathbf{f}_*)),$$

$$\bar{\mathbf{f}}_* = K_{\boldsymbol{\theta}}(X_*, X)[K_{\boldsymbol{\theta}}(X, X) + \sigma^2 I]^{-1}\mathbf{y},$$

$$\text{cov}(\mathbf{f}_*) = K_{\boldsymbol{\theta}}(X_*, X_*) - K_{\boldsymbol{\theta}}(X_*, X)[K_{\boldsymbol{\theta}}(X, X) + \sigma^2 I]^{-1}K_{\boldsymbol{\theta}}(X, X_*).$$

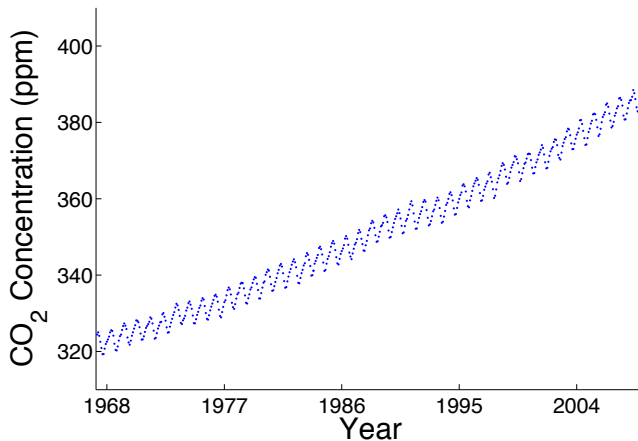
# Gaussian process graphical model



- ▶ Squared are observed, circles are latent, the thick bar is a set of fully connected nodes.
- ▶ Each  $y_i$  is conditionally independent given  $f_i$ .
- ▶ Because of the marginalization property of a GP, addition of further inputs  $x_*$  and unobserved targets  $y_*$  does not change the distribution of any other variables.

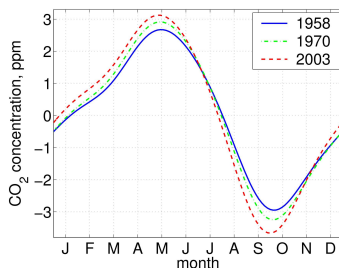
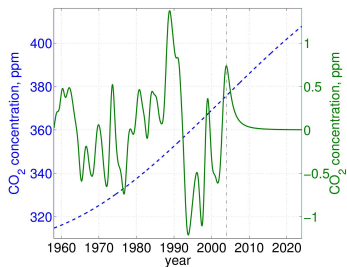
Figure from GPML, Rasmussen and Williams (2006)

## Worked Example: Combining Kernels, CO<sub>2</sub> Data

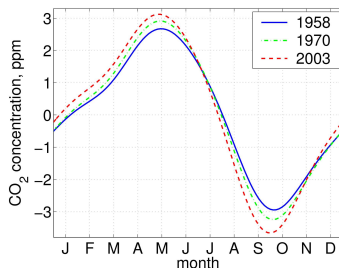
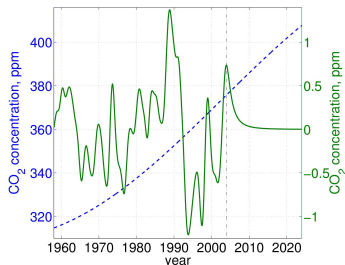


Example from Rasmussen and Williams (2006), *Gaussian Processes for Machine Learning*.

# Worked Example: Combining Kernels, CO<sub>2</sub> Data

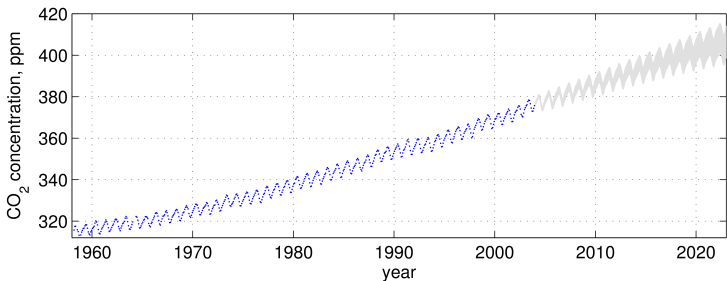


# Worked Example: Combining Kernels, CO<sub>2</sub> Data



- ▶ Long rising trend:  $k_1(x_p, x_q) = \theta_1^2 \exp\left(-\frac{(x_p - x_q)^2}{2\theta_2^2}\right)$
- ▶ Quasi-periodic seasonal changes:  $k_2(x_p, x_q) = k_{\text{RBF}}(x_p, x_q)k_{\text{PER}}(x_p, x_q) = \theta_3^2 \exp\left(-\frac{(x_p - x_q)^2}{2\theta_4^2} - \frac{2 \sin^2(\pi(x_p - x_q))}{\theta_5^2}\right)$
- ▶ Multi-scale medium term irregularities:  $k_3(x_p, x_q) = \theta_6^2 \left(1 + \frac{(x_p - x_q)^2}{2\theta_8\theta_7^2}\right)^{-\theta_8}$
- ▶ Correlated and i.i.d. noise:  $k_4(x_p, x_q) = \theta_9^2 \exp\left(-\frac{(x_p - x_q)^2}{2\theta_{10}^2}\right) + \theta_{11}^2 \delta_{pq}$
- ▶  $k_{\text{total}}(x_p, x_q) = k_1(x_p, x_q) + k_2(x_p, x_q) + k_3(x_p, x_q) + k_4(x_p, x_q)$

# Worked Example: Combining Kernels, CO<sub>2</sub> Data



- ▶ Hand crafted a kernel combination to perform extrapolation
- ▶ Confidence in the extrapolation is high (suggests that model is well specified).
- ▶ Can interpret the learned kernel hyperparameters  $\theta$  to learn information about our dataset.
- ▶ A lot of the interesting pattern recognition has been done by a human in this example. We would like to completely automate this modelling procedure.

# Non-Gaussian Likelihoods

We can no longer analytically integrate away the Gaussian process. But we can use a simple Monte carlo sum:

$$\begin{aligned} p(f_* | \mathbf{y}, X, x_*) &= \int p(f_* | \mathbf{f}, x_*) p(\mathbf{f} | \mathbf{y}) d\mathbf{f} \\ &\approx \frac{1}{J} \sum_{j=1}^J p(f_* | \mathbf{f}^{(j)}, x_*), \quad \mathbf{f}^{(j)} \sim p(\mathbf{f} | \mathbf{y}) \end{aligned}$$

But how do we sample from  $p(\mathbf{f} | \mathbf{y})$ ?



# Non-Gaussian Likelihoods

We can no longer analytically integrate away the Gaussian process. But we can use a simple Monte carlo sum:

$$\begin{aligned} p(f_* | \mathbf{y}, X, x_*) &= \int p(f_* | \mathbf{f}, x_*) p(\mathbf{f} | \mathbf{y}) d\mathbf{f} \\ &\approx \frac{1}{J} \sum_{j=1}^J p(f_* | \mathbf{f}^{(j)}, x_*), \quad \mathbf{f}^{(j)} \sim p(\mathbf{f} | \mathbf{y}) \end{aligned}$$

But how do we sample from  $p(\mathbf{f} | \mathbf{y})$ ?

*Elliptical slice sampling.* Murray et. al. AISTATS 2010.

But what about hyperparameters? It's easy to implement Gibbs sampling:

$$p(\mathbf{f}|\mathbf{y}, \theta) \propto p(\mathbf{y}|\mathbf{f})p(\mathbf{f}|\theta) \quad (56)$$

$$p(\theta|\mathbf{f}, \mathbf{y}) \propto p(\mathbf{f}|\theta)p(\theta). \quad (57)$$

But this won't work because of strong correlations between  $\mathbf{f}$  and  $\theta$ .

But what about hyperparameters? It's easy to implement Gibbs sampling:

$$p(\mathbf{f}|\mathbf{y}, \theta) \propto p(\mathbf{y}|\mathbf{f})p(\mathbf{f}|\theta) \quad (58)$$

$$p(\theta|\mathbf{f}, \mathbf{y}) \propto p(\mathbf{f}|\theta)p(\theta). \quad (59)$$

But this won't work because of strong correlations between  $\mathbf{f}$  and  $\theta$ .

- ▶ Transform into a *whitened* space,  $\mathbf{f} = L\boldsymbol{\nu}$ , and sample from  $\boldsymbol{\nu}$  and  $\theta$ , which decouples correlations.

# Non-Gaussian Likelihoods

But what about hyperparameters? It's easy to implement Gibbs sampling:

$$p(\mathbf{f}|\mathbf{y}, \theta) \propto p(\mathbf{y}|\mathbf{f})p(\mathbf{f}|\theta) \quad (60)$$

$$p(\theta|\mathbf{f}, \mathbf{y}) \propto p(\mathbf{f}|\theta)p(\theta). \quad (61)$$

But this won't work because of strong correlations between  $\mathbf{f}$  and  $\theta$ .

- ▶ Transform into a *whitened* space,  $\mathbf{f} = L\boldsymbol{\nu}$ , and sample from  $\boldsymbol{\nu}$  and  $\theta$ , which decouples correlations.
- ▶ Use a deterministic approach to approximately integrate away  $\mathbf{f}$  to access a marginal likelihood, conditioned only on kernel hyperparameters  $\theta$ :

$$p(\mathbf{y}|\theta) = \int p(\mathbf{y}|\mathbf{f})p(\mathbf{f}|\theta)d\mathbf{f} \quad (62)$$

- ▶ The *Laplace approximation*, for example, approximates  $p(\mathbf{f}|\mathbf{y})$  as a Gaussian.

# Readings for Next Time

- ▶ C. Rasmussen and C. Williams, GPML, Ch. 4, 5
- ▶ Y. Saatchi, PhD Thesis, 2011. Chapter 5
- ▶ J. Candela and C.E. Rasmussen, A unifying view of sparse approximation Gaussian process regression, JMLR 2005.
- ▶ A.G. Wilson and R.P. Adams. Gaussian process kernels for pattern discovery and extrapolation, ICML 2013.