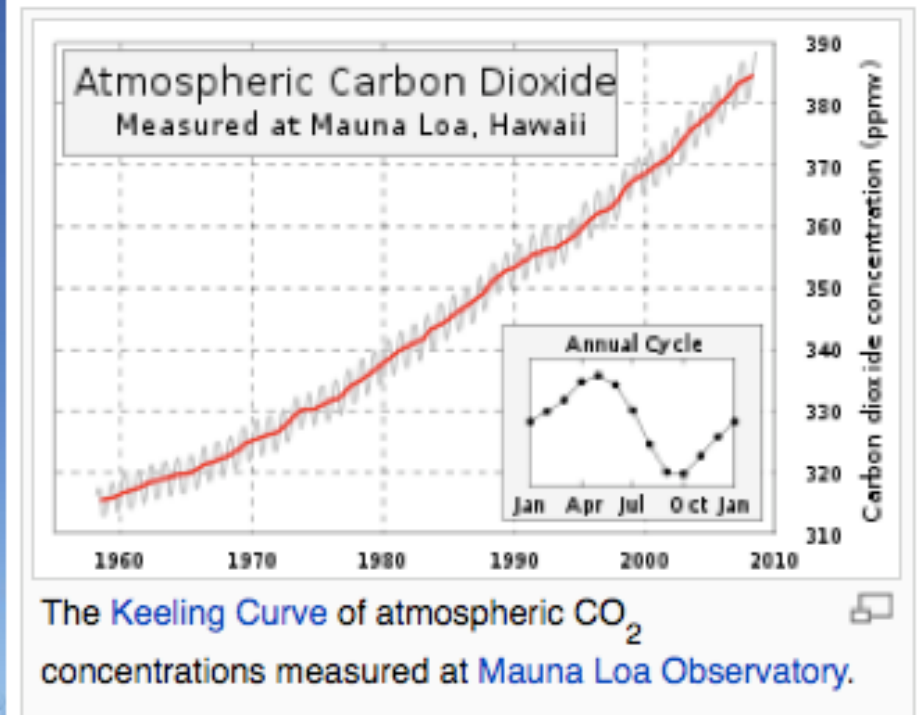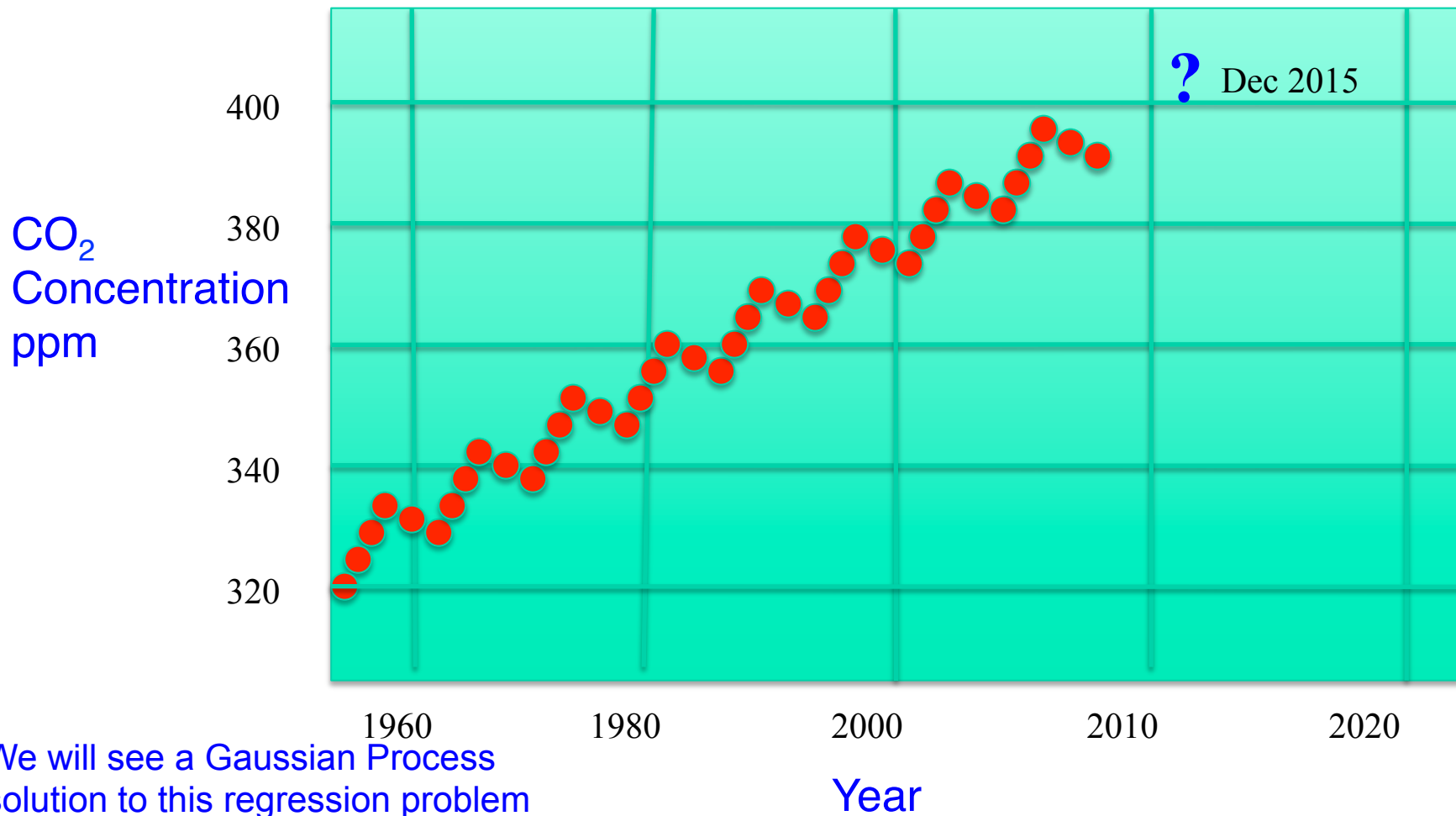# Gaussian Processes

## Sargur Srihari

# Topics in Gaussian Processes

1. Examples of use of GP
2. Duality: From Basis Functions to Kernel Functions
3. GP Definition and Intuition
4. Linear regression revisited
5. Gaussian processes for regression
6. Learning the hyperparameters
   Automatic Relevance Determination
7. Gaussian processes for classification
   Laplace approximation
8. Connection to neural networks

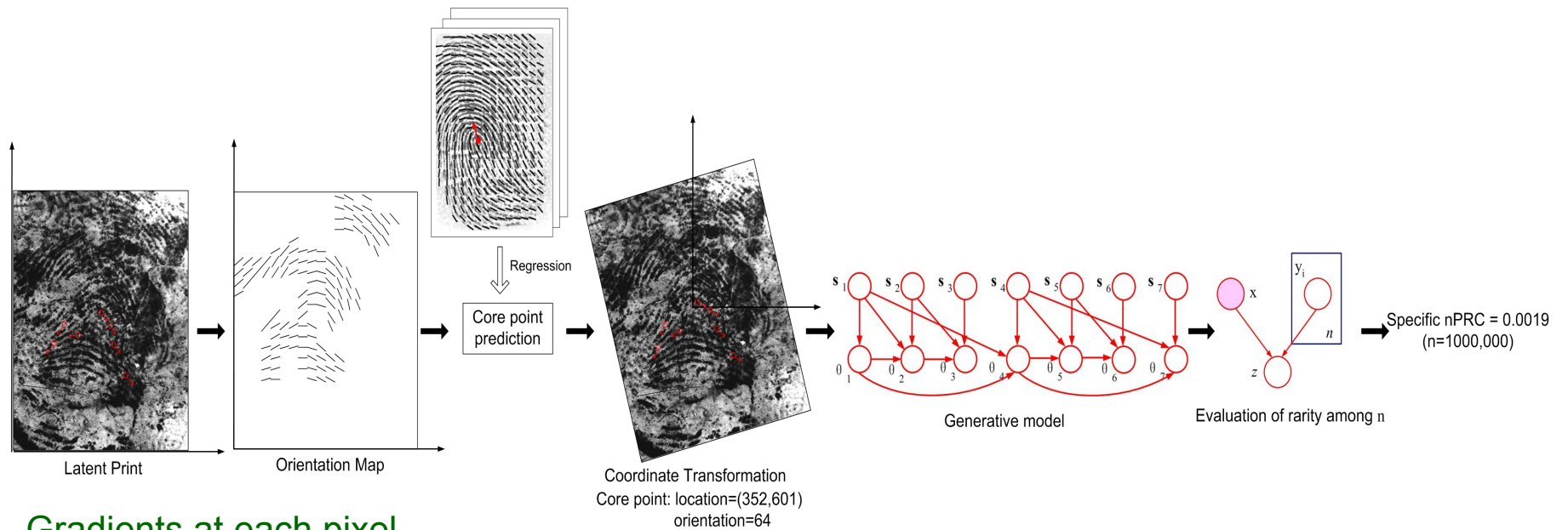# National Academy of Sciences: Keck Center



The Keeling Curve of atmospheric $CO_2$ concentrations measured at Mauna Loa Observatory.

3

# Regression Problem:
# Carbon Dioxide in Atmosphere



We will see a Gaussian Process
solution to this regression problem

# Fingerprint Core Point using GP Regression



Regression

Core point prediction

Latent Print

Orientation Map

Coordinate Transformation
Core point: location=(352,601)
orientation=64

Generative model

Evaluation of rarity among $n$
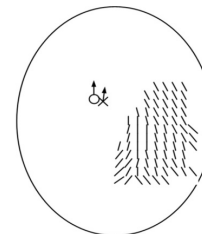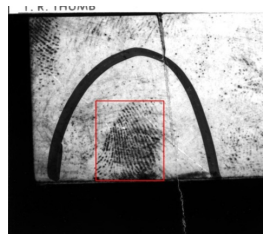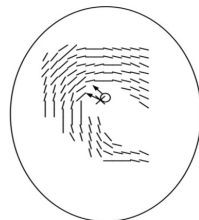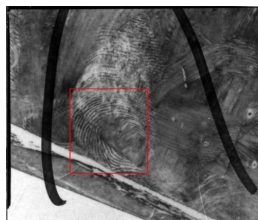
Specific nPRC = 0.0019
(n=1000,000)

### Gradients at each pixel

0 0 0 34 32 33 32 38 42 35 29 32 24 24 22 21 17 10 5 169 143 131 135  131 131 127 139
146 0 0 0 0 0 0 32 33 33 30 34 39 45 27 33 31 22 22 29 15 11 7 176 150 141 134  130
136 119 129 131 118 0 0 0 0 0 29 30 31 35 36 35 39 40 33 33 36 24 21 34 25 17 10 1 173
151 140  132 136 122 129 113 110 0 0 0 0 0 34 32 33 31 30 39 41 43 37 32 31 19 23 26
16 9 8 178 161 146 136  125 129 116 134 127 121 0 0 0 0 0 41 38 36 30 33 44 43 43 39
35 31 26 23 26 19 5 1 176 164 156 136 121 129 114 121 110 102 0 0 0 · · · · · ·

Core point $(x, y, \theta)$:
(253, 221, 85)

# Paper at NIPS 2010

**NIPS 2010**

## Evaluation of Rarity of Fingerprints in Forensics

Chang Su and Sargur N. Srihari
*University at Buffalo, The State University of New York*

**UB**

### ABSTRACT

Since the earliest days of forensics, importance of considering the rarity of features used in the comparison of prints has been known. Yet methods to compute rarity of features has been elusive due to the large number of variables and the complexity of the distributions. When a latent print, typically found in a crime scene, is compared to a known (inked or live scan) in a database, the degree of uncertainty involves both rarity and similarity. It has become necessary to quantify this uncertainty in courtroom situations.

Proposed method for rarity, based on level 2 detail (minutiae), determines probability of random match among n knowns:
1. The coordinate system is transformed into standard position with core point as origin; Gaussian Process regression is used to determine the core point. Method better than standard topology-based method.
2. Generative model, that takes into account inter-minutia dependencies and minutia confidences determines evidence probability; model is validated using a goodness-of-fit test.
3. Specific probability of match among n is evaluated.
The rarity of several examples, including latent fingerprints and minutia configurations are given.
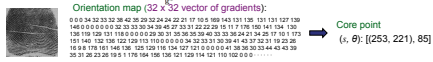
## 1. Core Point Prediction (Gaussian Process)

Given orientation map of latent print, core point is predicted by

$$\hat{y}^* = k(\mathbf{g}^*_{MAX}, G)[K + \sigma^2 I]^{-1}\mathbf{y} \quad \text{where} \quad k(\mathbf{g}, \mathbf{g}') = \theta_1 \exp(-\tfrac{\theta_2}{2}|\mathbf{g} - \mathbf{g}'|^2)$$

where $\mathbf{g}^*_{MAX}$ is orientation map with the maximum predictive probability

$$\mathbf{g}^*_{MAX} = \operatorname{argmax} p(m(y^*)|\mathbf{g}^*, G, \mathbf{y})$$

Orientation map (32 x 32 vector of gradients):

Core point
$(s, \theta)$: [(253, 221), 85]

Orientation map estimated in blocks of pixels using gradients

$$\theta_s = \tfrac{1}{2}\tan^{-1}\left(\frac{\sum_i^w \sum_j^w 2G_x(i,j)G_y(i,j)}{\sum_i^w \sum_j^w (G_x^2(i,j) - G_y^2(i,j))}\right)$$

Latent print with predicted core point *within* print    Latent print with predicted core point *outside* print

Performance Comparison with Standard Method (Poincare topology)

|  | Poincare Index | Gaussian Processes |
|---|---|---|
| Good | 90.6% | 93.1% |
| Bad | 68.2% | 87.1% |
| Ugly | 46.6% | 72.7% |
| Overall | 68.6% | 84.5% |

Training data : NIST4(4000 fingerprints)
Testing data: NIST27
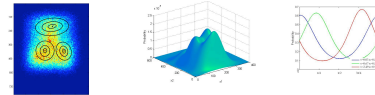(258 latent prints: 88 good, 85 bad, 85 ugly)

## 2. Coordinate Transformation

### Process Flow for a Latent Print (Four Steps)

Latent Print — Orientation Map — Core point prediction — Coordinate Transformation Core point: locations(352,601) orientation=64 — Generative model — Evaluation of rarity among n — Specific nPRC = 0.0019 (n=1000,000)

## 3. A Generative Model

Distribution of minutiae: mixture model with bivariate Gaussian for location and von Mises for Orientation

Distribution of fingerprint captures dependency between minutiae: define a unique sequence for given minutiae set. Graphical model (Bayesian network) used to represent the minutiae set.

(a) minutiae sequencing and (b) dependency

Graphical model

For minutia set **X** joint distribution is

$$p(\mathbf{X}) = p(\mathbf{s}_1)p(\theta_1|\mathbf{s}_1)\prod_{n=2}^{N} p(\mathbf{s}_n)p(\theta_n|\mathbf{s}_n, \mathbf{s}_{\psi(n)}, \theta_{\psi(n)})$$

s =location, θ is orientation, ψ(n) are parent nodes

which involves three mixture models

$$f(\mathbf{s}) = \sum_{k_1=1}^{K_1} \pi_{k_1} \mathcal{N}(\mathbf{s}|\mu_{k_1}, \Sigma_{k_1})$$

$$f(\mathbf{s}, \theta) = \sum_{k_2}^{K_2} \pi_{k_2} \mathcal{N}(\mathbf{s}|\mu_{k_2}, \Sigma_{k_2})\mathcal{V}(\theta|\nu_{k_2}, \kappa_{k_2})$$

Within a mixture component vonMises models orientation given the location

$$f(\theta_n|\mathbf{s}_n, \mathbf{s}_{\psi(n)}, \theta_{\psi(n)}) = \sum_{k_3=1}^{K_3} \pi_{k_3} \mathcal{V}(\theta_n|\nu_{k_3}, \kappa_{k_3})$$

Goodness of Fit: Chi-square test for three generative models.

| Generative models | Dataset sizes | Model accepted | Model rejected |
|---|---|---|---|
| $f(\mathbf{s})$ | 4000 | 3387 | 613 |
| $f(\mathbf{s}, \theta)$ | 4000 | 3216 | 784 |
| $f(\theta_n|\mathbf{s}_n, \mathbf{s}_{\psi(n)}, \theta_{\psi(n)})$ | 4096 | 3558 | 538 |

## 4. Specific nPRC Calculation with Minutiae Confidence

Specifying minutia correspondence
Minutiae are within tolerance $\varepsilon = [\varepsilon_s, \varepsilon_\theta]$.
$$\|\mathbf{s}_a - \mathbf{s}_b\| \le \varepsilon_s \wedge |\theta_a - \theta_b| \le \varepsilon_\theta$$

Modeling minutia quality
Confidence of minutia $\mathbf{x}_n$ is $(d_{sn}, d_{\theta n})$, where $d_{sn}$ is location confidence and $d_{\theta n}$ is direction confidence. Given minutiae $x_n = (s_n, \theta_n)$ distributions of location s' and direction θ' are modeled by Gaussian and von-Mises:

$$c(\mathbf{s}'|\mathbf{s}_n, d_{s_n}) = \mathcal{N}(\mathbf{s}'|\mathbf{s}_n, d_{s_n}^{-1}) \qquad c(\theta'|\theta_n, d_{\theta_n}) = \mathcal{V}(\theta'|\theta_n, d_{\theta_n})$$

Specific nPRC given by

$$p_c(\mathbf{X}, \hat{m}, n) = 1 - (1 - p_c(\mathbf{X}, \hat{m}))^{n-1}$$

where $p_c(X, \hat{m})$ is the probability that $\hat{m}$ minutiae match between latent print and a random one among n knowns

$$p_c(\mathbf{X}, \hat{m}) = \sum_{m' \in M} p(m')\binom{m'}{\hat{m}} \cdot \sum_{i=1}^{\binom{N}{\hat{m}}} p_c(\widetilde{\mathbf{X}}_i)$$

and probability that there is a one-to-one correspondence between ^X and ^X' is given by

$$p_c(\widetilde{\mathbf{X}}) = p_c(\mathbf{s}_1, \theta_1)\prod_{n=2}^{\hat{m}} p_c(\mathbf{s}_n)p_c(\theta_n|\mathbf{s}_n, \mathbf{s}_{\psi(n)}, \theta_{\psi(n)})$$

$$p_c(\mathbf{s}_n, \theta_n)$$
$$= \int\int\int\int_{|\mathbf{s}-\mathbf{s}'|\le\varepsilon_s} c(\mathbf{s}'|\mathbf{s}_n, d_{s_n})c(\theta'|\theta_n, d_{\theta_n})f(\mathbf{s}, \theta)d\mathbf{s}'d\theta'd\mathbf{s}d\theta$$
$$p_c(\mathbf{s}_n) = \int\int_{|\mathbf{s}-\mathbf{s}'|\le\varepsilon_s} c(\mathbf{s}'|\mathbf{s}_n, d_{s_n})f(\mathbf{s})d\mathbf{s}'d\mathbf{s}$$

$$p_c(\theta_n|\mathbf{s}_n, \mathbf{s}_{\psi(n)}, \theta_{\psi(n)})$$
$$= \int\int_{|\theta-\theta'|\le\varepsilon_\theta} c(\theta'|\theta_n, d_{\theta_n})f(\theta|\mathbf{s}_n, \mathbf{s}_{\psi(n)}, \theta_{\psi(n)})d\theta'd\theta$$

## Examples of Rarity Evaluation

- Three Minutiae Configurations, specific nPRC for n=1000

| 1.2 E-2, | 7.9 E-4 | 2.3 E-6 |

- Latent Print: Madrid bombing case

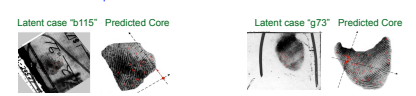Original latent print level 2 details.    Mayfield inked print    Charted features on latent print

Probability of at least one print in FBI IAFIS database (470 million prints) that shares the same minutiae is 0.93.

Probability of falsely identifying source of latent print is $7.8 \times 10^{-7}$

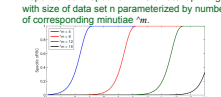- Latent Prints: NIST special dataset 27

Latent case "b115"  Predicted Core     Latent case "g73"  Predicted Core

Minutiae confidence manually assigned by visual inspection. Tolerance set at 10 pixels and π/8.

Specific nPRCs for the latent prints, with n= 1,000,000.

| Latent Print "b115" | | Latent Print "g73" | |
|---|---|---|---|
| $N$ | $\hat{m}$ | $p_c(\hat{m}, \mathbf{X})$ | $N$ | $\hat{m}$ | $p_c(\hat{m}, \mathbf{X})$ |
|  |  | 0.73 |  | 4 | 1 |
|  | 2 | $9.04 \times 10^{-10}$ |  | 8 | $3.11 \times 10^{-18}$ |
| 16 | 8 | $2.46 \times 10^{-19}$ | 39 | 12 | $2.56 \times 10^{-25}$ |
|  | 12 | $6.13 \times 10^{-33}$ |  | 24 | $3.10 \times 10^{-52}$ |
|  | 16 | $1.82 \times 10^{-46}$ |  | 39 | $7.51 \times 10^{-79}$ |

Dependence of specific nPRC of latent print "g73" with size of data set n parameterized by number of corresponding minutiae ^m.

Values of specific nPRC are largely dependent on given latent print

When ^m decreases and n increases, PRC increases

6

# Dual Representation

- ## Linear regression model
  $$y(\mathbf{x},\mathbf{w}) = \sum_{j=0}^{M-1} w_j \phi_j(\mathbf{x}) = \mathbf{w}^T\phi(\mathbf{x})$$

  $\phi_j$ are basis functions or features

  - Parameters obtained by minimizing regularized *sum-of-squares:*

  $$J(\mathbf{w}) = \frac{1}{2}\sum_{n=1}^{N}\left\{\mathbf{w}^T\phi(\mathbf{x}_n) - t_n\right\}^2 + \frac{\lambda}{2}\mathbf{w}^T\mathbf{w}$$

  where $\mathbf{w}=(w_0,..,w_{M-1})^T$, $\phi=(\phi_0,..\phi_{M-1})^T$, $N$ samples $\{x_1,..x_N\}$, $\lambda$ = regularization coefficient

  - Solution obtained by setting gradient of $J(\mathbf{w})$ wrt $\mathbf{w}$ equal to zero

  $$\mathbf{w} = -\frac{1}{\lambda}\sum_{n=1}^{N}\left\{\mathbf{w}^T\phi(\mathbf{x}_n) - t_n\right\}\phi(\mathbf{x}_n) = \sum_{n=1}^{N}a_n\phi(\mathbf{x}_n) = \Phi^T\mathbf{a}$$

  where $\Phi$ is the design matrix whose $n^{th}$ row is given by $\phi(\mathbf{x}_n)^T$

  $$\Phi = \begin{bmatrix} \phi_0(x_1) & . & . & \phi_{M-1}(x_1) \\ . & & & . \\ \phi_0(x_n) & . & . & \phi_{M-1}(x_n) \\ . & & & . \\ \phi_0(x_N) & . & . & \phi_{M-1}(x_N) \end{bmatrix}$$ is a $N \times M$ matrix

  vector $\mathbf{a}=(a_1,..,a_N)^T$ with

  $$a_n = -\frac{1}{\lambda}\left\{\mathbf{w}^T\phi(\mathbf{x}_n) - t_n\right\}$$

- Solution is a linear combination of vectors $\phi(\mathbf{x}_n)$ with weights $a_n$
- Instead of working with parameter vector $\mathbf{w}$ we work with vector $\mathbf{a}$

# Transformation from Basis to Kernel

- Substitute $w = \Phi^T a$ into $J(w)$

  | Gram Matrix Definition: |
  |---|
  | Given $N$ vectors, it is an $N$ x $N$ matrix of inner products |

- with Gram matrix $K = \Phi\Phi^T$ elements

$$K_{nm} = \phi(x_n)^T \phi(x_m) = k(x_n, x_m)$$

  - where we introduce the kernel function $k(x,x') = \phi(x)^T \phi(x')$

    Example is Gaussian kernel $k(x,x') = \exp(-\|x-x'\|^2/2\sigma^2)$



- Sum of squares error written in terms of Gram matrix as

$$J(a) = \frac{1}{2}a^T KKa - a^T Kt + \frac{1}{2}t^T t + \frac{\lambda}{2}a^T Ka$$

  whose solution is    $a = (K + \lambda I_N)^{-1}t$    Solution for least squares problem expressed entirely in terms of kernel $k(x,x')$

- Prediction for new input $x$

$$y(x) = w^T \phi(x) = k(x)^T (K + \lambda I_N)^{-1}t$$
$$\text{where } k(x) \text{ has elements } k_n(x) = k(x_n, x)$$

  - Prediction is linear combination of kernel evaluated at training points
    - Need to invert $N$ x $N$ matrix
    - but avoid problems of very high (even infinite) dimensionality of $x$

# Role of Gaussian Processes

1. ## As a kernel method
   - Duality leads to a non-probabilistic model for linear regression
   - Extending role of kernels to probabilistic discriminative models leads to Gaussian Processes (kernels arise naturally in a Bayesian setting)

2. ## As a Bayesian method
   - In Bayesian linear regression
     - Model has the form $y(\mathrm{x},\mathrm{w})=\mathrm{w}^T\phi(\mathrm{x})$ where a prior distribution over $\mathrm{w}$ induces a prior distribution over functions $y(\mathrm{x},\mathrm{w})$
     - Given training data set, we evaluate the posterior distribution over $\mathrm{w}$ which gives posterior distribution over regression functions
     - With noise it implies a predictive distribution $p(t|\mathrm{x})$ for a new input vector $\mathrm{x}$
   - In Gaussian Processes
     - Dispense with parametric model and instead define a prior probability distribution over functions directly
     - For a finite training set we only need to consider values of the function at discrete set of input values $\mathrm{x}_n$

9

# Probabilistic Linear Regression Revisited

- Re-derive the predictive distribution by working in terms of distribution over functions $y(\mathrm{x},\mathrm{w})$
  - It will provide a specific example of a Gaussian Process
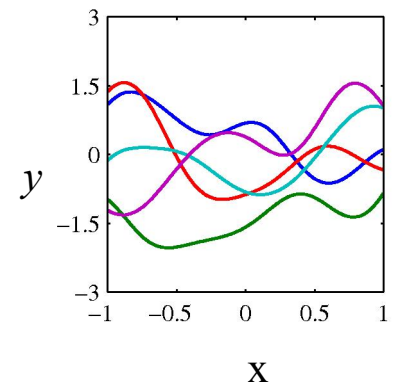
- Consider model with $M$ fixed basis functions

$$y(\mathrm{x}) = \mathrm{w}^{\mathrm{T}}\phi(\mathrm{x})$$

  where $\mathrm{x}$ is the input vector and $\mathrm{w}$ is the $M$-dimensional weight vector

- Assume a Gaussian distribution of weight $\mathrm{w}$

$$p(\mathrm{w}) = \mathcal{N}(\mathrm{w}|0,\alpha^{-1}\mathrm{I})$$

- Probability distribution over $\mathrm{w}$ induces a probability distribution over functions $y(\mathrm{x})$

# Probabilistic Linear Regression is a Gaussian Process

- We wish to evaluate $y$ (x) at training points $x_1,.., x_N$

- Our interest is the joint distribution of values $y= [y(x_1),..y(x_N)]$
  Since   $y(x)=w^T\phi(x)$   we can write

  $$y = \Phi w$$              N x M times M x 1 yields a N x 1

  - where $\Phi$ is the design matrix with elements $\Phi_{nk}=\phi_k(x_n)$

- Since y is a linear combination of elements of $w$ which are
  Gaussian distributed as $p(w) = \mathcal{N}(w|0,\alpha^{-1}I)$

  y is itself Gaussian with
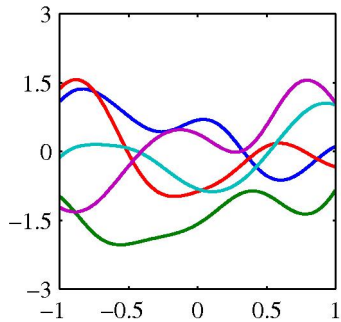
  mean:                    $E[y] =\Phi E[w]=0$  and

  variance:                    $Cov[y]=E[yy^T]=\Phi E[ww^T]\Phi^T=(1/\alpha)\Phi\Phi^T=K$

  where K is the *N x N* Gram Matrix with elements

  $$K_{nm}= k (x_n,x_m) = (1/\alpha)\phi(x_n)^T\phi(x_m)$$

  and $k$ (x,x$'$ ) is the kernel function

Covariance expressed by Kernel function

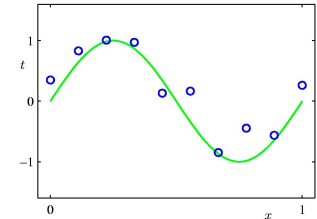Thus Gaussian distributed weight vector induces a Gaussian joint distribution over training samples

# General definition of Gaussian Process

- ## We saw a particular example of a Gaussian process
  - Linear regression using the parametric model   $y(\mathbf{x})=\mathbf{w}^T\phi(\mathbf{x})$
  - With samples y= $[y(\mathbf{x}_1),..y(\mathbf{x}_N)]$
  - Assume  $p(\mathbf{w}) = \mathcal{N}(\mathbf{w}|0,\alpha^{-1}\mathbf{I})$
  - Then E[y] =0  and  Cov[y]=K, the Gram matrix which is equivalent to  pairwise kernel values

- ## More generally, a Gaussian process is a probability distribution over functions $y$ (x)
  - Such that the set of values of $y(\mathbf{x})$ evaluated at arbitrary points $\mathbf{x}_1,..,\mathbf{x}_N$ jointly have a multivariate Gaussian distribution
    - For a single input $\mathbf{x}_1$, output $y_1$ is univariate Gaussian. For two inputs $\mathbf{x}_{1,}\,\mathbf{x}_2$ the output $y_1,y_2$ is bivariate Gaussian, etc.12

# Stochastic Process

- A *stochastic process* $y(\mathbf{x})$ is specified by giving the joint probability distribution for any finite set of values $y(\mathbf{x}_1),\ldots,y(\mathbf{x}_N)$ in a consistent manner
  - The random variables typically develop over time

- A Gaussian process is a stochastic process which is Gaussian

- When input $\mathbf{x}$ is 2-D it is also known as a *Gaussian Random Field*
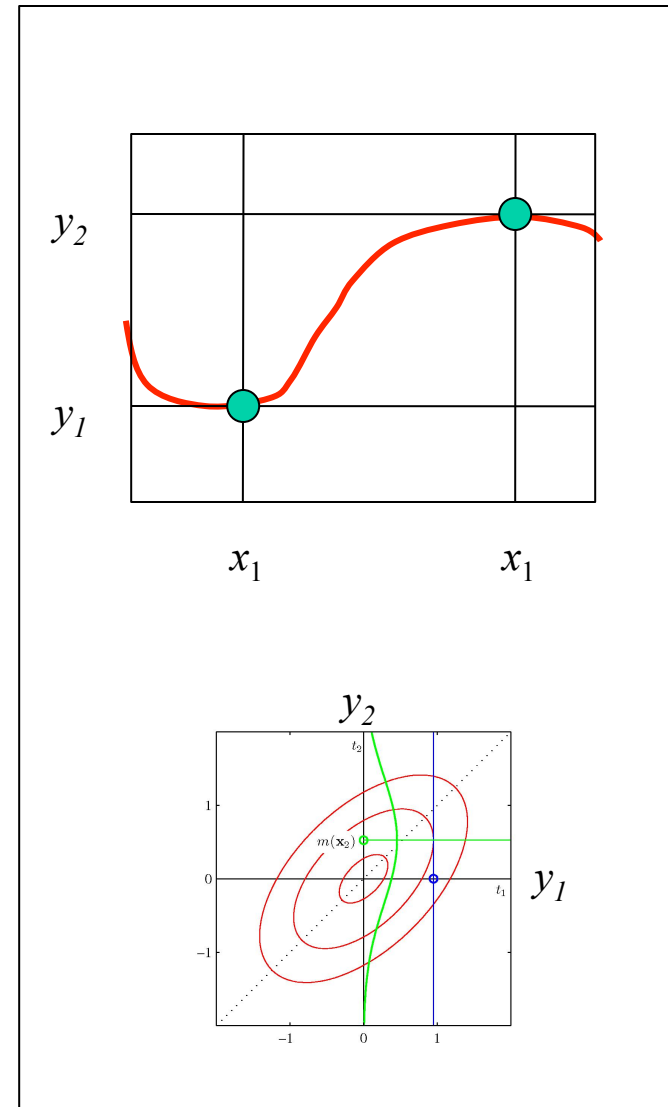
# Parametric Regression vs GP Regression

- In parametric regression we have several samples from which we learn the parameters

- In Gaussian processes we we view the samples as one huge input that has a Gaussian distribution
  - with a mean and a covariance matrix

- A **Gaussian process** is a stochastic process $X_t$, $t \in T$, for which any finite linear combination of samples has a joint Gaussian Distribution
  - any linear functional applied to the sample function $X_t$ will give a normally distributed result. Notation-wise, one can write $X \sim GP(m,K)$, meaning the random function $X$ is distributed as a GP with mean function $m$ and covariance function $K$.

14
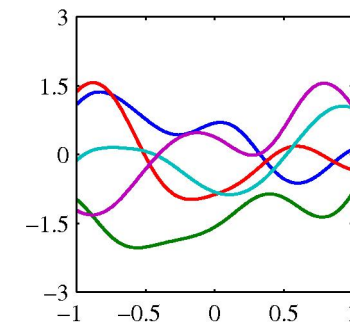
# Gaussian Process with Two Samples

- Let $y$ be a function (curve)
  - of a one-dimensional variable $x$
- We take two **samples** $y_1$ and $y_2$
  corresponding to $x_1$ and $x_2$
- Assume they have a bivariate
  Gaussian distribution $p(y_1, y_2)$
- Each point from this distribution
  - has an associated probability
  - It also defines a function $y(x)$
    - Assuming that two points are
      enough to define a curve
- More than two points will be
  needed to define a curve
  - Which leads to a higher
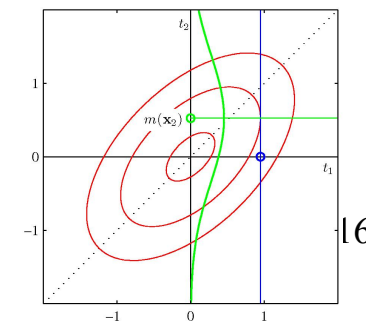    dimensional probability distribution

15

# Gaussian Process with $N$ samples

- A process that generates samples over time $\{y_1,..y_N\}$ is a GP iff every set of samples $Y=[y_1,..,y_N]$ is a vector-valued Gaussian random variable

- We define a distribution over all functions
  - Constrained by the samples

- The samples have a joint Gaussian distribution in $N$-dimensional space

Regression: Possible functions constrained by $N$ samples



Case of two samples $t_1$ and $t_2$ that are bivariate Gaussian



16

# Specifying a Gaussian Process

- Key point about Gaussian Stochastic Processes
  - Joint distribution over $N$ variables $y_1,.., y_N$ is completely specified by the second-order statistics,
    - i.e., mean and covariance

- With mean zero, it is completely specified by covariance of $y(\mathrm{x})$ evaluated at any two values of $\mathrm{x}$ which is given by a kernel function

  $$\mathrm{E}[y(\mathrm{x}_n)\, y(\mathrm{x}_m)] = k(\mathrm{x}_n,\mathrm{x}_m)$$

- For the Gaussian Process defined by the linear regression model $y(\mathrm{x},\mathrm{w})=\mathrm{w}^\mathrm{T}\phi(\mathrm{x})$ with prior $p(\mathrm{w}) = \mathcal{N}(\mathrm{w}|0,\alpha^{-1}\mathrm{I})$ the kernel function is

  $$\mathrm{K}_{nm}= k\,(\mathrm{x}_n,\mathrm{x}_m) = (1/\alpha)\,\phi(\mathrm{x}_n)^\mathrm{T}\phi(\mathrm{x}_m)$$
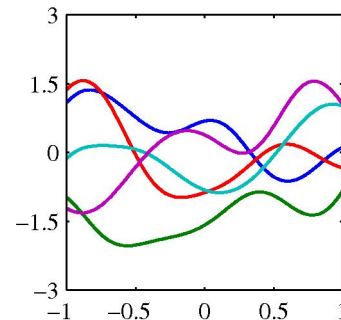
17

# Defining a Kernel Directly

- GP can also be specified directly by choice of kernel function (instead of indirectly by basis function)
- Samples of functions drawn from Gaussian processes for two different choices of kernel functions are shown next

# Samples from Gaussian Processes for two kernels

**Functions are drawn from**
**Gaussian processes**
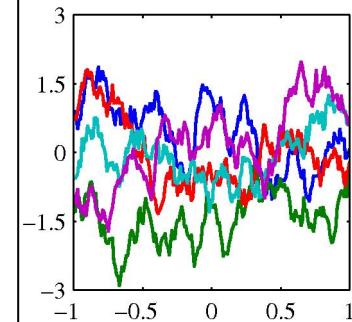Note that each sample
is a function

Gaussian Kernel

$$k(\mathrm{x},\mathrm{x}') = \exp\left(- \parallel \mathrm{x} - \mathrm{x}' \parallel^2 / 2\sigma^2\right)$$

Exponential Kernel

$$k(\mathrm{x},\mathrm{x}') = \exp\left(-\theta \mid \mathrm{x} - \mathrm{x}' \mid\right)$$





Ornstein-Uhlenbeck
process for
Brownian motion

19

# Probability Distributions and Stochastic Processes

- A probability distribution describes the distributions of scalars ($x$) and vectors ($\mathbf{x}$)

- A stochastic process describes distribution of functions *f(x)*

  - Can think of a function as a very long vector

    - Each entry in the vector is $f(x)$ which is the value of function at $x$

- A Gaussian process is a generalization of Gaussian distributions where it is describing the distribution of functions

- Since the vector is infinite-dimensional we constrain it to only those points $x$ corresponding to training and test points

20

# 3. Gaussian Process for Regression (Direct Approach)

- We specify Gaussian Process directly over functions
  - Abandon approach of defining a distribution over weights $\mathbf{w}$

- Take into account noise on observed target values as

$$t_n = y_n + \varepsilon_n \quad \text{where } y_n = y(\mathbf{x}_n)$$

  - Noise process has a Gaussian distribution $p(t_n|y_n) = N(t_n|y_n, \beta^{-1})$

- Note that target $t_n$ is output $y_n$ corrupted by noise
- Defining $\mathbf{t} = (t_1, .., t_N)^T$ our goal is to determine a distribution $p(\mathbf{t})$
  - Which is a distribution over functions

# 3. Gaussian Process for Regression (Direct Approach)

- Assuming noise is independent for each data point
  - joint distribution *of* $t = (t_1,..,t_N)^T$ on values $y = (y_1,..,y_N)^T$ is

  $$p(t|y) = N(t|y, \beta^{-1}I_N)$$

- From definition of GP, marginal distribution of y is given by
  - a Gaussian with zero mean, covariance matrix given by Gram matrix $K$

  $$p(y) = N(y|0, K)$$

- Where kernel function that determines $K$ is chosen to express:
  - property that for points $x_n$, $x_m$ that are similar corresponding values $y(x_n)$, $y(x_m)$ will be more strongly correlated than for dissimilar points
  - $K$ depends on application

# Regression: Marginal Distribution

- From distributions $p(t|y)$ and $p(y)$ we can get marginal $p(t)$ conditioned on input values $x_1,..,x_N$
- Applying Sum rule and product rule
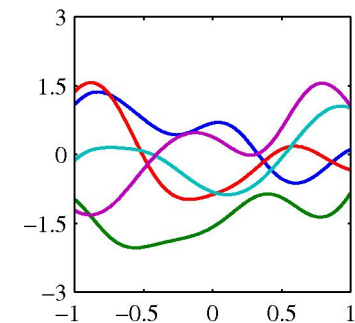
$$p(t) = \int p(t \mid y)p(y)dy$$

$$= N(t \mid 0, C)$$

From the result that when $p(y)$ and $p(t|y)$ are Gaussian $p(t)$ is also Gaussian

where covariance matrix $C$ has the elements

$$C(x_n, x_m) = k(x_n, x_m) + \beta^{-1}\delta_{nm}$$

Due to $y(x)$     Due to $\varepsilon$

- The two Gaussian sources of randomness, $y(x)$ and $\varepsilon$ are independent, so their covariances simply add

$\delta$ specified by $\varepsilon$

# Widely used kernel function for Gaussian Process

- Exponential of a quadratic form

  with addition of constant and linear terms

$$k(\mathrm{x}_n, \mathrm{x}_m) = \theta_0 \exp\left\{-\frac{\theta_1}{2} \| \mathrm{x}_n - \mathrm{x}_m \|^2\right\} + \theta_2 + \theta_3 \mathrm{x}_n^T \mathrm{x}_m$$

  Corresponds to a parametric model that is
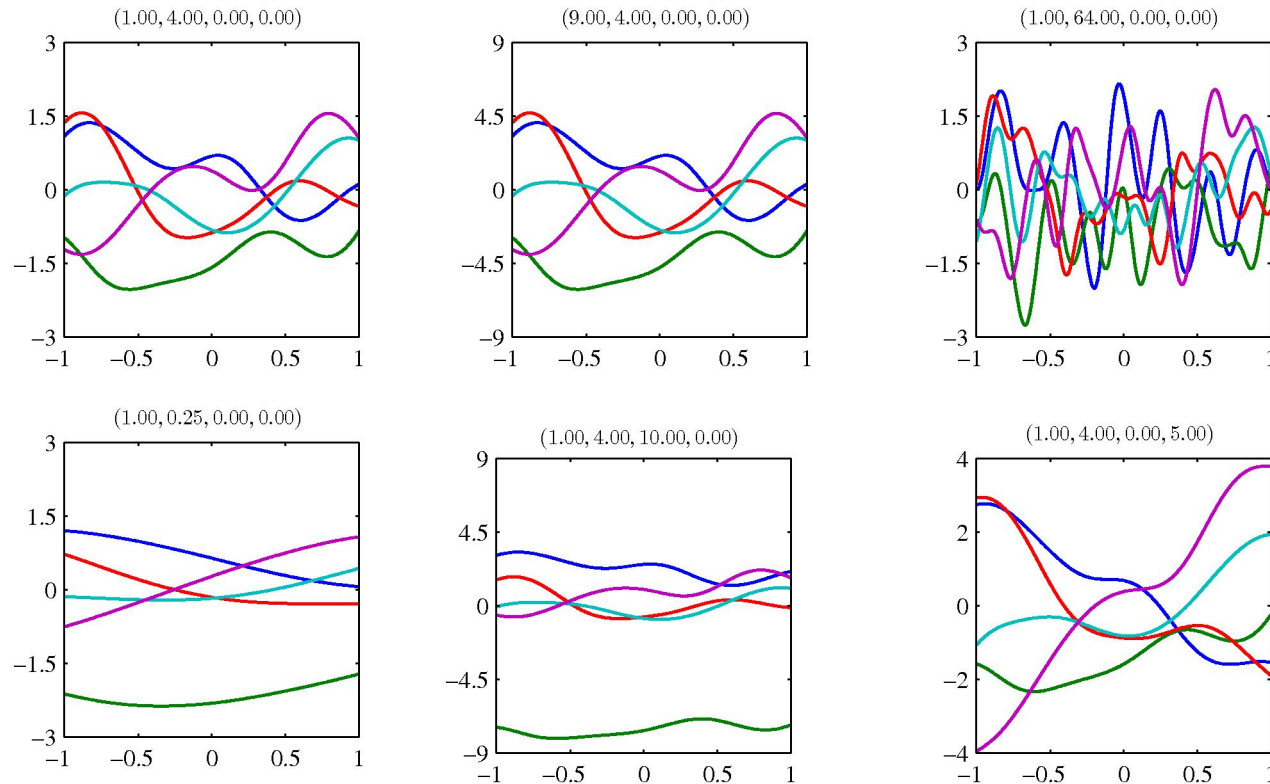  a linear function of input variables

- Samples from this prior are plotted for various values
  of the parameters $\theta_0$, $\theta_1$, $\theta_2$, and $\theta_3$

# Samples from a Gaussian Process Prior

$$k(\mathbf{x}_n, \mathbf{x}_m) = \theta_0 \exp\left\{ -\frac{\theta_1}{2} \| \mathbf{x}_n - \mathbf{x}_m \|^2 \right\} + \theta_2 + \theta_3 \mathbf{x}_n^T \mathbf{x}_m$$

for different settings of $\theta_0, \theta_1, \theta_2, \theta_3$



Each function is a sample

25

# Illustration of sampling data points from a Gaussian Process

- A sample from joint distribution $p(\mathrm{y}) = N(\mathrm{y}|0,K)$
  - where $\mathrm{y} = [y_1,.. y_N]^T$
  - Each $y_i$ is associated with a value of $x$ since the elements of $\mathrm{y}$ are ordered
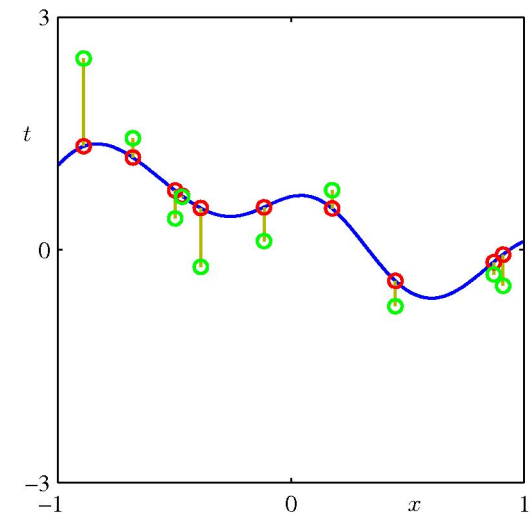- Noise values are added as defined by $p(\mathrm{t}) = N(\mathrm{t}|0,C)$

Blue curve:
a sample function from the Gaussian process prior over functions

Red points:
values of $y_n$ obtained by evaluating function at set of input values $\{x_n\}$

Green points:
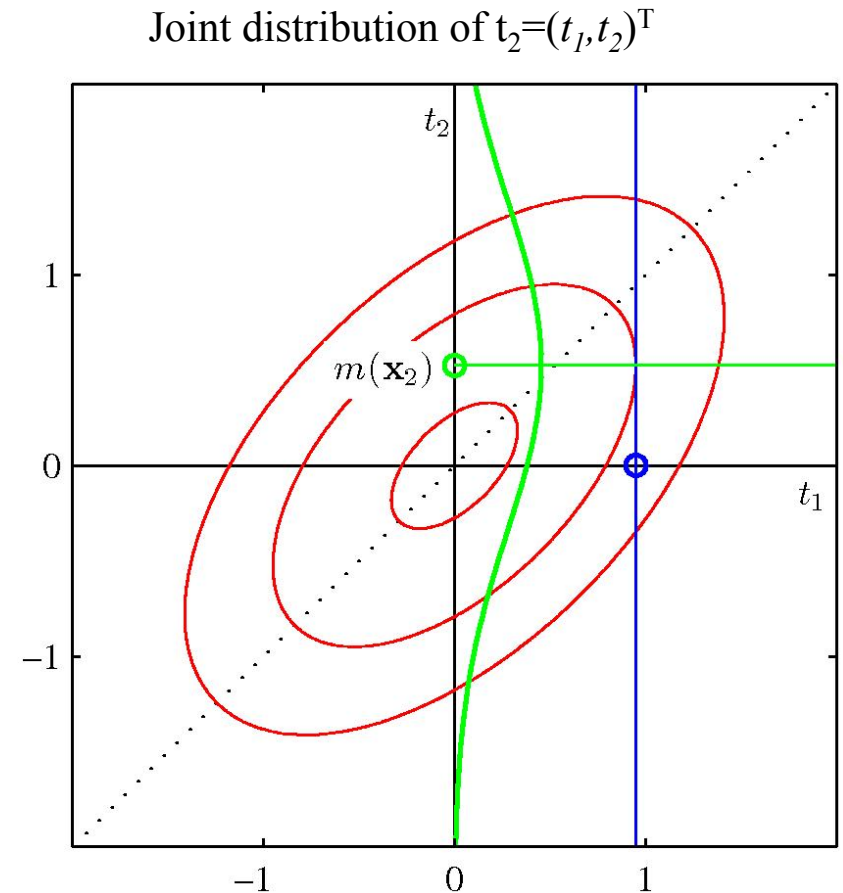values of $t_n$ obtained by adding independent Gaussian noise to each of $y_n$

# Making predictions using Gaussian Processes

- Goal of regression: predict target for new inputs
    - given training data $\mathbf{t}_N = (t_1,..,t_N)^T$
    - with corresponding input values $( \mathbf{x}_1,..\mathbf{x}_N)$
- Goal is to predict target variable $t_{N+1}$ given $\mathbf{x}_{N+1}$
- This requires we evaluate the predictive distribution

  $p(t_{N+1}|\mathbf{t}_N)$ which is also conditioned on $\mathbf{x}_1,..,\mathbf{x}_N$ and $\mathbf{x}_{N+1}$ which are not shown to keep notation simple

- To find conditional distribution $p(t_{N+1}|\mathbf{t}_N)$ we begin by writing down the joint distribution $p(\mathbf{t}_{N+1})$
- We then obtain the required conditional distribution

# Mechanism for Gaussian process regression

- One training point and one test point

- Red ellipses show contours of $p(t_1,t_2)$

- Green curve is $p(t_2|t_1)$

- $t_1$ is the training point

- Conditioning on $t_1$ corresponding to blue line

- Predictive distribution for $t_{N+1}$ is Gaussian whose mean and variance both depend on $\mathbf{x}_{N+1}$

Joint distribution of $t_2=(t_1,t_2)^T$



$t_2$

$m(\mathbf{x}_2)$

$t_1$

28

# Finding Conditional Distribution $p(t_{N+1}|\mathbf{t})$

- Joint distribution is given by

$$p(\mathbf{t}_{N+1})=\mathcal{N}(\mathbf{t}_{N+1}|\mathbf{0},\mathbf{C}_{N+1})$$

  - where $\mathbf{C}_{N+1}$ is the $(N+1)$ x $(N+1)$ covariance matrix with elements given by $C(\mathbf{x}_n,\mathbf{x}_m) = k(\mathbf{x}_n,\mathbf{x}_m) + \beta^{-1}d_{nm}$

- Because the joint is Gaussian, conditional Gaussian distribution is given by partitioning the covariance matrix as

$$\mathbf{C}_{N+1} = \begin{pmatrix} \mathbf{C}_N & \mathbf{k} \\ \mathbf{k}^T & c \end{pmatrix}$$
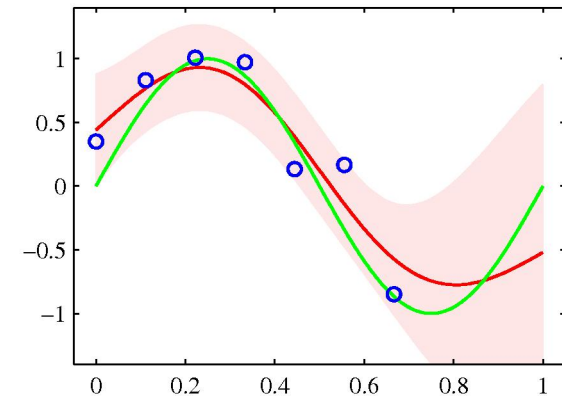
  where $\mathbf{C}_N$ is the $N$ x $N$ covariance matrix

  vector k has elements $k(\mathbf{x}_n,\mathbf{x}_{N+1})$ for $n=1,..,N$

  and scalar $c = k(\mathbf{x}_{N+1},\mathbf{x}_{N+1})+\beta^{-1}$

- Conditional distribution $p(t_{N+1}|\mathbf{t})$ is Gaussian with

  - Mean $m(\mathbf{x}_{N+1}) = \mathbf{k}^T\mathbf{C}_N^{-1}\mathbf{t}$
  - Variance $\sigma^2(\mathbf{x}_{N+1}) = c - \mathbf{k}^T\mathbf{C}_N^{-1}\mathbf{k}$

29

- Key results that define Gaussian Regression

# Gaussian process regression applied to sinusoidal data set

- Green: sinusoidal function from which data points are obtained

- Red: mean of the Gaussian process predictive distribution

- Shaded region: plus and minus two standard deviations

- Uncertainty increases to right of data points

# Summary of Gaussian Process Regression

- A Gaussian Process is a generalization of a multivariate Gaussian distribution to infinitely many variables

- A Gaussian *distribution* is fully specified by a mean vector $\mu$ and covariance matrix $\Sigma$

    $f = (f_1,..f_n)^T \sim \mathcal{N}(\mu,\Sigma)$        indexes $i = 1,...n$

- A Gaussian *process* is fully specified by a mean function $m(x)$ and covariance function $k(x,x')$

    $f(x) \sim \mathcal{GP}(m(x), k(x,x'))$ indexes   x  *(which are functions, or infinite vectors)*

- Note that the kernel function $k$ appears  in place of  the covariance matrix

    - Both express similarity of two points in multidimensional space
    - Note: correlation is same as covariance when all variables are zero-mean, unit variance

31

# The marginalization property

- Thinking of a GP as a Gaussian distribution with an infinitely long mean vector and an infinite by infinite covariance matrix may seem impractical

- Luckily we have the marginalization property

$$p(\mathbf{x}) = \int p(\mathbf{x}, \mathbf{y}) dy$$

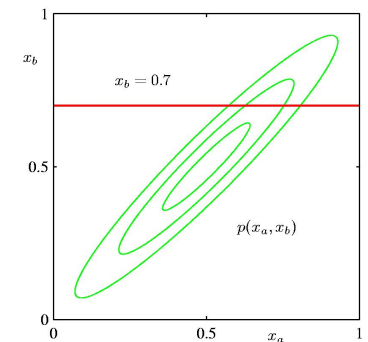$$\mathbf{x} = \begin{bmatrix} \mathbf{x}_a \\ \mathbf{x}_b \end{bmatrix}$$

$$p(\mathbf{x}_a, \mathbf{x}_b) = N\left( \begin{bmatrix} \mu_a \\ \mu_b \end{bmatrix}, \begin{bmatrix} A & B \\ B^T & C \end{bmatrix} \right)$$

$$\Rightarrow p(\mathbf{x}_a) = N(\mu_a, A)$$
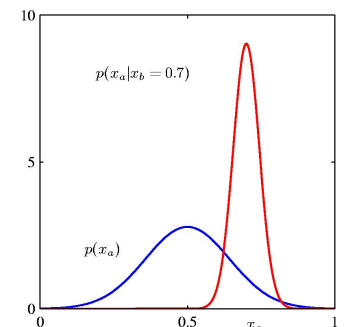
Joint is Gaussian



- Marginalizing over unknown function values $\mathbf{x}_b$
- To get distribution over known values $\mathbf{x}_a$
- Similarly conditional distribution

Marginal and
Conditional are Gaussian



$$\Rightarrow p(\mathbf{x}_a | \mathbf{x}_b) = N(\mu_a - A^{-1}B(x_b - \mu_b), A^{-1})$$
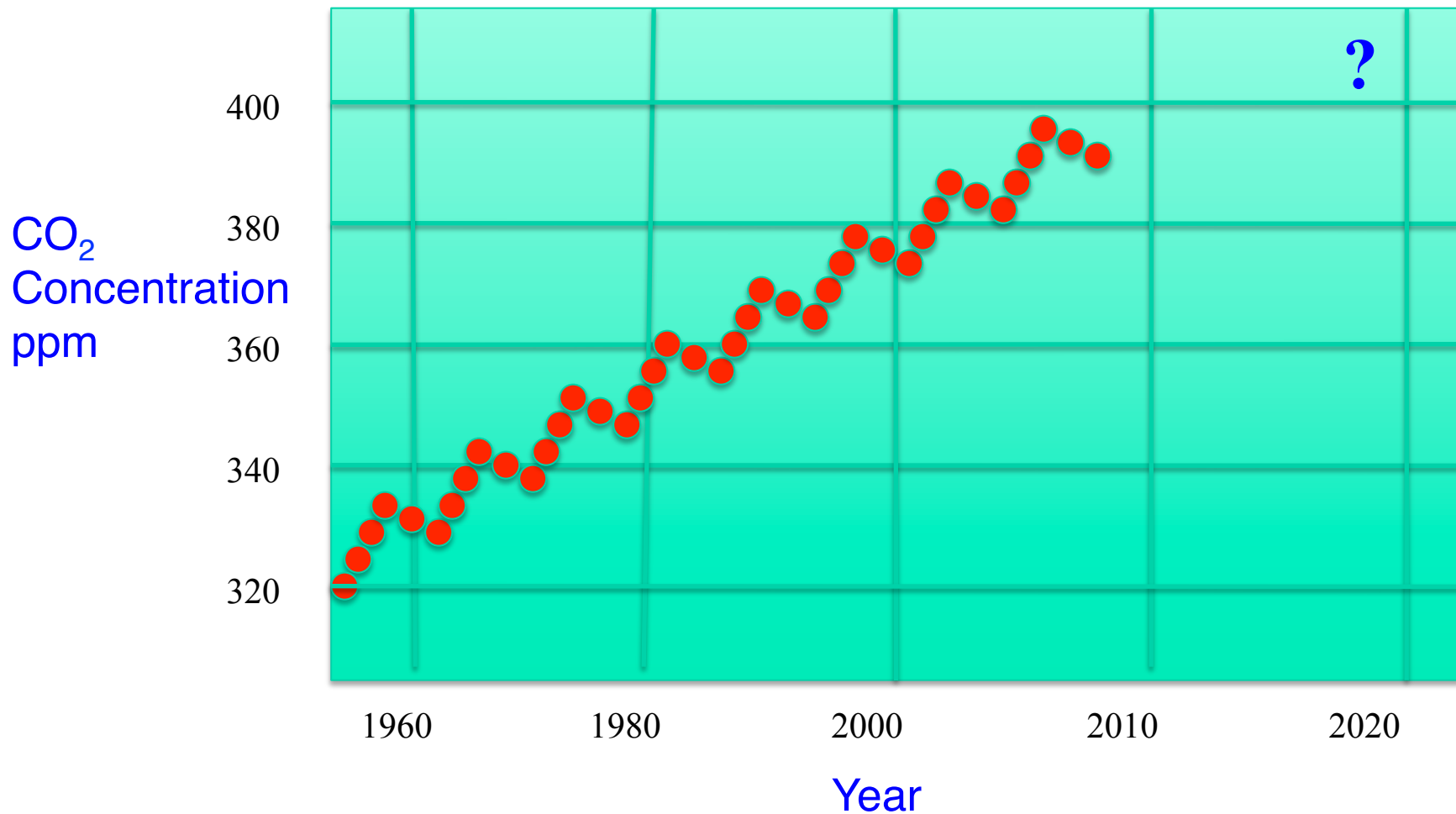
# Predictive distribution

- Mean $m(\mathrm{x}_{N+1}) = \mathrm{k}^\mathrm{T}\mathbf{C}_N^{-1}\mathbf{t}$

- Variance $\sigma^2(\mathrm{x}_{N+1}) = \mathrm{c} - \mathrm{k}^\mathrm{T}\mathbf{C}_N^{-1}\mathrm{k}$

  where $C_N$ is the $N \, x \, N$ covariance matrix with elements $k(\mathrm{x}_n,\mathrm{x}_m) + \beta^{-1}\delta_{nm}$

  vector k has elements $k(\mathrm{x}_n,\mathrm{x}_{N+1})$ for $n=1,..,\,N$

  and scalar $\mathrm{c} = k(\mathrm{x}_{N+1},\mathrm{x}_{N+1})+\delta^{-1}$

- Same restrictions on kernel function as in general kernel methods (positive definite)

- Mean can also be written as

$$m(\mathrm{x}_{N+1}) = \sum_{n=1}^{N} a_n k(\mathrm{x}_n,\mathrm{x}_{N+1})$$

  - Where $a_n$ is the nth component of $C_N^{-1}t$

  - If $k(\mathrm{x}_n,\mathrm{x}_m)$ depends only on distance $\|\mathrm{x}_n\text{-}\mathrm{x}_m\|$ then we obtain an expansion in radial basis functions
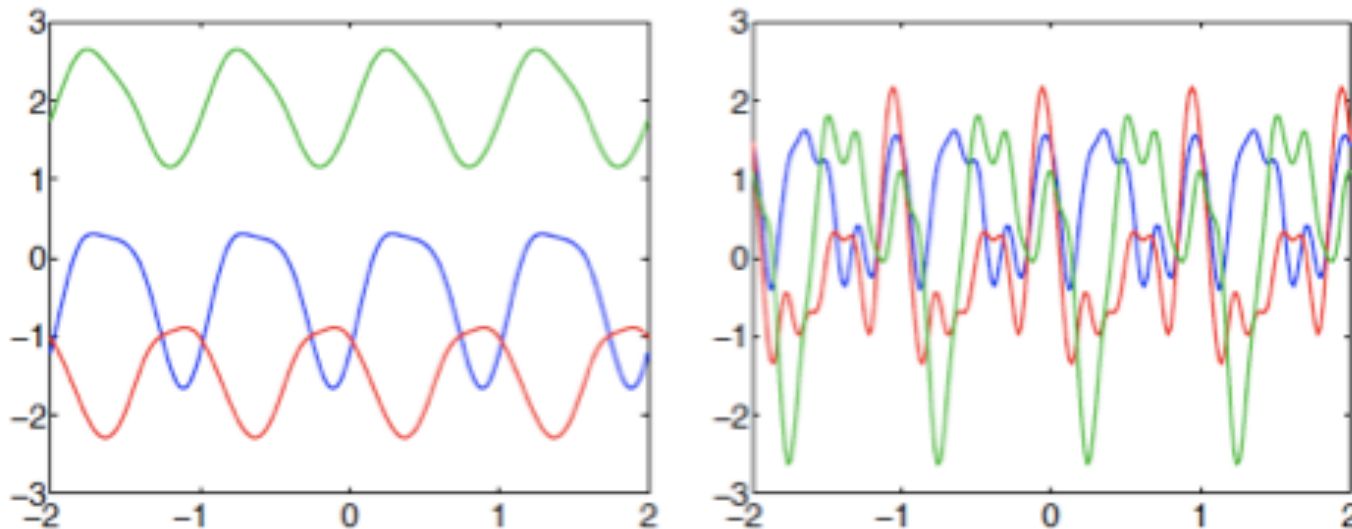
# Regression Problem:
# Carbon Dioxide in Atmosphere

# Periodic Smooth Functions

## To create distribution over periodic functions of $x$:

- First map inputs to $u=(\sin x, \cos x)^{\mathrm{T}}$
- Then measure distances in $u$-space
- Combined with squared exponential (SE) covariance function with characteristic length ale $l$ we get

$$k_{periodic}(x,x') = \exp\left(\frac{-2\sin^2(\pi(x-x'))}{l^2}\right)$$



35    Three functions drawn at random: left $l > 1$, right $l < 1$

# Covariance Function for Periodic Data

The covariance function consists of several terms, parameterized by 11 hyperparameters

$$k(x,x') = k_1(x,x') + k_2(x,x') + k_3(x,x') + k_4(x,x')$$

Where

Long-term smooth trend (squared exponential)

$$k_1(x,x') = q_1{}^2 \exp(-(x-x')^2/q_2{}^2)$$

Seasonal trend (quasi-periodic smooth)

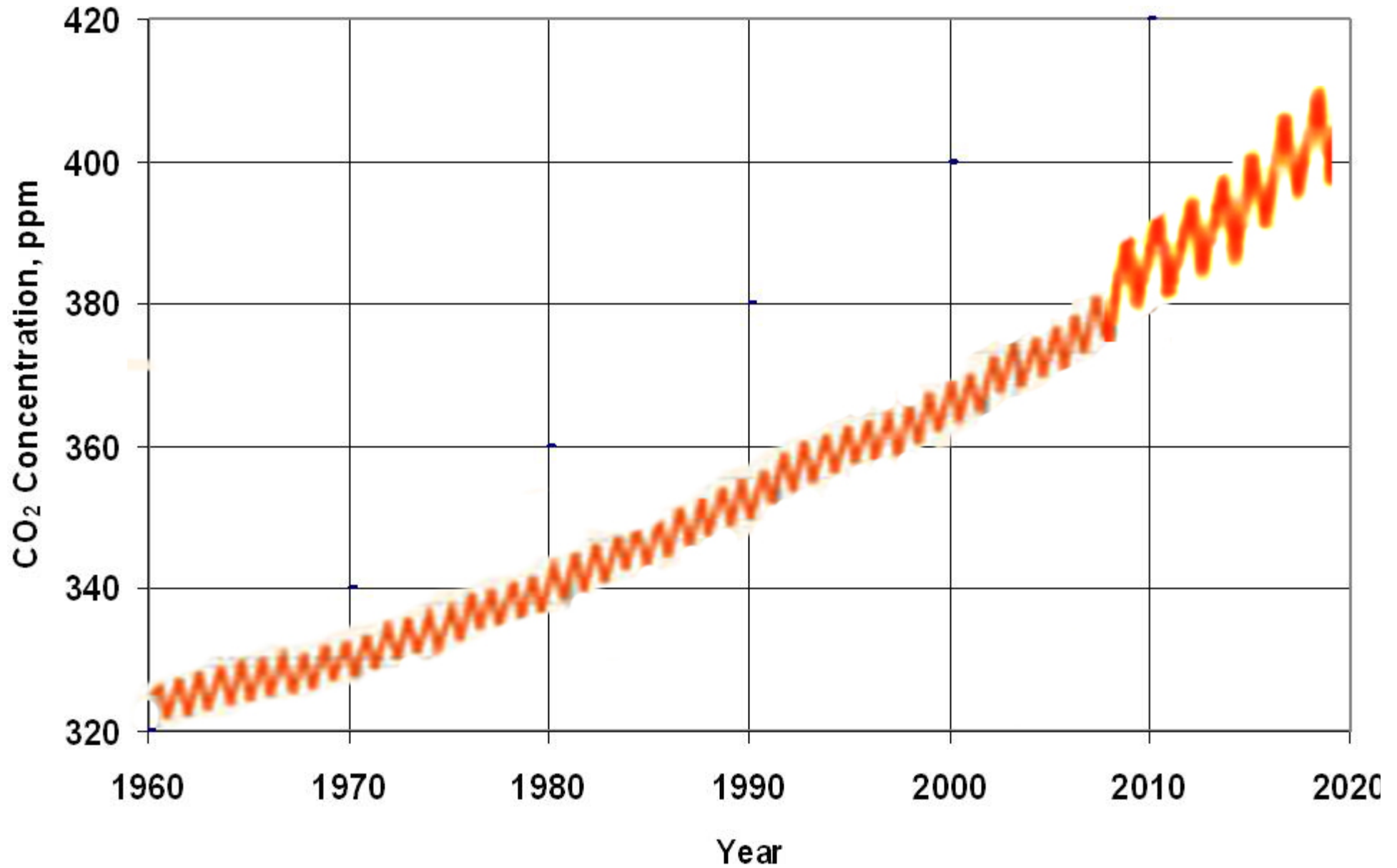$$k_2(x,x') = q_3{}^2 \exp(-2\sin^2[p(x-x')]/q_s{}^2 \text{ x } \exp(-(1/2)(x-x')^2/q_4{}^2),$$

Short-term and medium-term anomaly (rational quadratic)

$$k_3(x,x') = \theta_6^2\left(1 + \frac{(x-x')^2}{2\theta_8\theta_7}\right) - \theta_8$$

Noise (independent Gaussian and dependent)

$$k_4(x,x') = \theta_9^2 \exp\left(-\frac{(x-x')2}{2\theta_{10}}\right) + \theta_{11}^2\delta_{xx}{}'$$

36

# Gaussian Process Fit to $CO_2$ data

# Computational Comparison

- Method 1: Parameter Space viewpoint using linear regression
- Method 2: Function space viewpoint using Gaussian process
    - Inversion of matrix of size $N$ x $N$
        - Standard methods require $O(N^3)$ computations
    - In the basis function model we have to invert a matrix $S_N$ of size $M$ x $M$ which has $O(M^3)$ complexity
- When $M < N$ basis function approach is more efficient
    - For large training sets many approximate methods for GP
- Gaussian process allows considering covariance functions that can only be expressed in terms of an infinite number of basis functions

# 4. Learning the hyperparameters

- Prediction of a Gaussian process model will depend on the choice of covariance function

- Rather than fixing the covariance function we can use a parametric family of functions and then infer parameter values from the data

- Parameters govern
  - Length-scale of the correlations
    - Region over which variables are correlated
  - Precision of the noise

- They correspond to the hyper-parameters in a standard parametric model

# Techniques for Learning Hyperparameters

- Based on evaluation of likelihood function $p(\mathbf{t}|\theta)$

  - where $\theta$ denotes the hyperparameters of the Gaussian process model

- Point estimate of $\theta$ is obtained by maximizing log-likelihood

$$\ln p(\mathbf{t}\,|\,\theta) = -\frac{1}{2}\ln|C_N| - \frac{1}{2}\mathbf{t}^T C_N^{-1}\mathbf{t} - \frac{N}{2}\ln(2\pi)$$

- Gradient of log-likelihood

$$\frac{\partial}{\partial \theta_i}\ln p(\mathbf{t}\,|\,\theta) = -\frac{1}{2}Tr\left(C_N^{-1}\frac{\partial C_N}{\partial \theta_i}\right) + \frac{1}{2}\mathbf{t}^T C_N^{-1}\frac{\partial C_N}{\partial \theta_i}C_N^{-1}\mathbf{t}$$
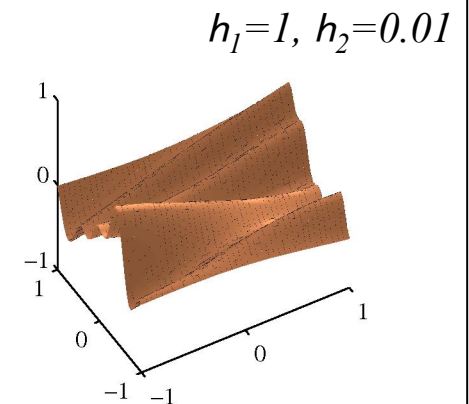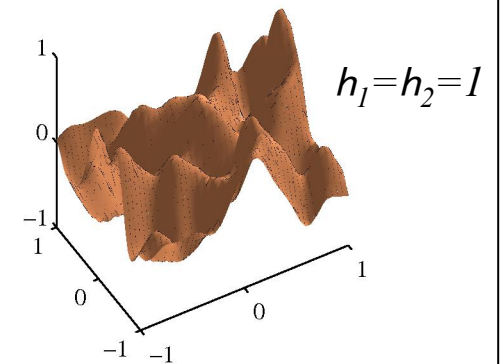
  - Can have multiple maxima

- Can be treated in a Bayesian manner by introducing a prior over $\theta$

- We have assumed a constant $\beta$. For heteroscedastic problems noise variance itself depends on $\mathrm{x}$

  - Can be handled by introducing a second Gaussian process
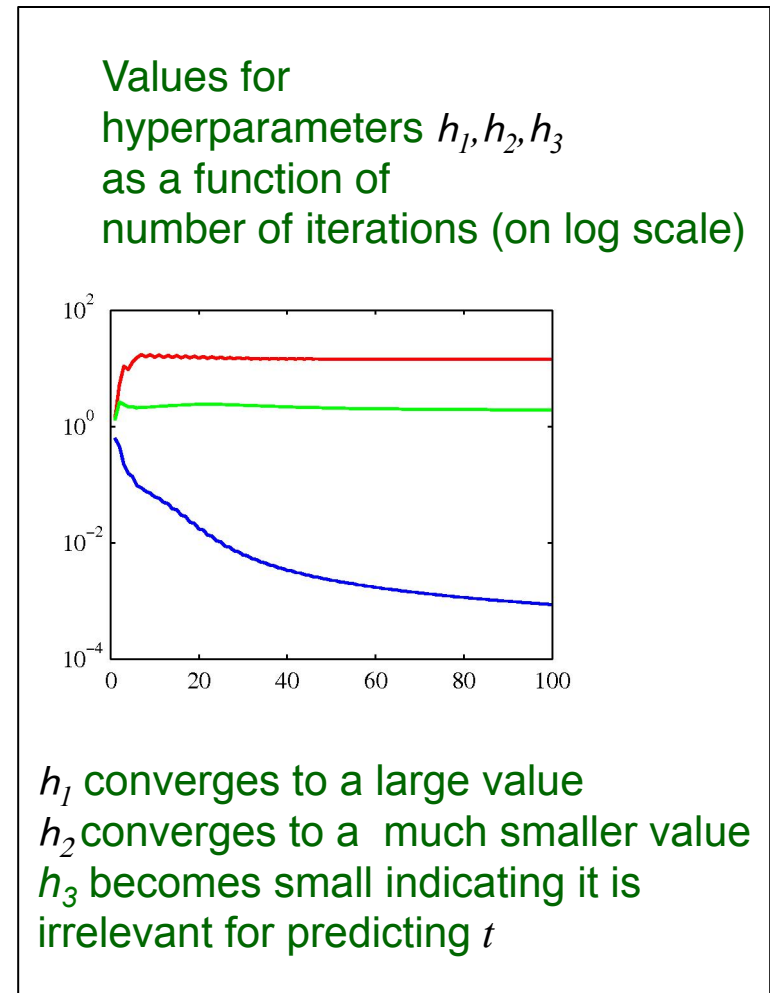
# 5. Automatic Relevance Determination

- Determining importance of variables
- Consider a Gaussian process with 2-D input space $\mathrm{x}=(x_1,x_2)$ having kernel function

$$k(x,x') = \theta_0 \exp\left\{-\frac{1}{2}\sum_{i=1}^{2}\eta_i(x_i - x_i')^2\right\}$$

- As particular parameter $h_i$ becomes small
  - function becomes insensitive to corresponding variable $x_i$
- By adapting these parameters to a data set by mle it becomes possible to detect variables that have little effect on predictive distribution
  - Such variables can be discarded

Samples from the ARD prior for Gaussian processes

$h_1=h_2=1$

$h_1=1, h_2=0.01$

41

# Automatic Relevance Determination in a Gaussian Process for Synthetic Problem

- Three inputs $x_1, x_2$ and $x_3$
- Target variable $t$ generated by sampling $100$ values of $x_1$ from a Gaussian, evaluating $sin\ (2\pi x_1)$ and adding Gaussian noise

- $x_2$ by adding noise to $x_1$

- $x_3$ are samples from an independent Gaussian

- Thus

  $x_1$ is a good predictor of t
  $x_2$ is a more noisy predictor
  $x_3$ has only chance correlations with t

Values for hyperparameters $h_1, h_2, h_3$ as a function of number of iterations (on log scale)



$h_1$ converges to a large value
$h_2$ converges to a much smaller value
$h_3$ becomes small indicating it is irrelevant for predicting $t$

# 6. Gaussian Processes for Classification

- Gaussian processes make predictions that lie on the entire real axis

- For two-class classification we need to model posterior probabilities of the target variable for a new input variable to lie in the interval $(0,1)$

- Can easily adapt Gaussian processes for classification

  - by transforming output of Gaussian process using appropriate nonlinear activation function

# Transforming Gaussian process output for classification

- Two class problem with target variable $t \in \{0,1\}$
- We define a Gaussian process over a function $a(\mathbf{x})$
- Then transform the function using a logistic sigmoid

  $$y = \sigma (a)$$

- Then we obtain a non-Gaussian stochastic process over functions $y(\mathbf{x})$ where $y \in \{0.1\}$
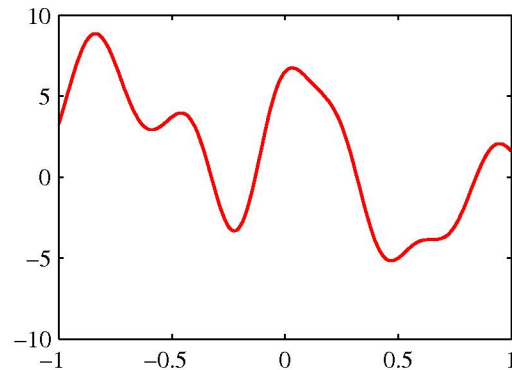
# Two-class Classification

## One-dimensional input space
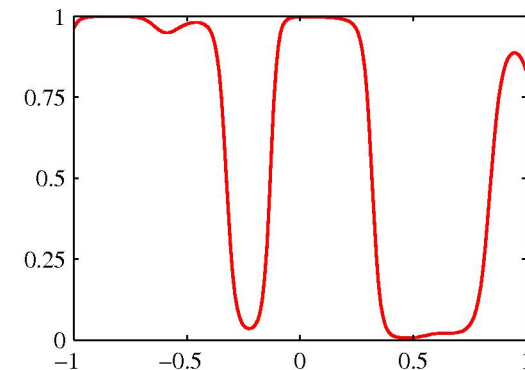
Bernoulli distribution

$$p(t \mid a) = \sigma(a)^t (1 - \sigma(a))^{1-t}$$

## Transforming output on real line to $(0,1)$ interval

A sample from a Gaussian process
prior over functions $a(\text{x})$

Result of transforming this sample
using a logistic sigmoid function

## Next: How to define the Gaussian Process

# Gaussian Process for Classification

- Training set samples $x_1 \ldots x_N$
- Corresponding target variables $t = (t_1, .., t_N)^T$
- Test point $x_{N+1}$ with target value $t_{N+1}$
- Goal is to determine the predictive distribution $p(t_{N+1}|t)$
    - where conditioning on input variable is left implicit
- Define a Gaussian process over function $a(x)$
- $a_{N+1}$ has the components $a(x_1), .., a(x_{N+1})$
- Gaussian process prior takes the form

  $p(a_{N+1}) = N(a_{N+1}|0, C_{N+1})$

  Unlike regression the covariance matrix no longer includes a noise term since all of the training data points are assumed to be correctly labeled

  Covariance matrix has elements $C(x_n, x_m) = k(x_n, x_m) + nd_{nm}$

- Predictive distribution is intractable

$$p(t_{N+1} = 1 \,|\, t_N) = \int p(t_{N+1} = 1 \,|\, a_{N+1}) p(a_{N+1} \,|\, t_N) \, da_{N+1}$$

# Three Approximation Methods for Gaussian Process Classification

1. Variational Inference
   - Uses local variational bound on logistic sigmoid
   - Allows product of sigmoids to be approximated by product of Gaussians thereby allowing marginalization over $a_N$ to be performed analytically

2. Expectation Propagation
   - Because true posterior is unimodal, can get good results
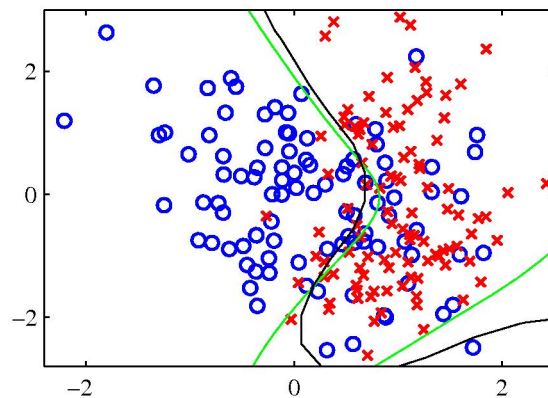
3. Laplacian Approximation

# 7. Laplace Approximation

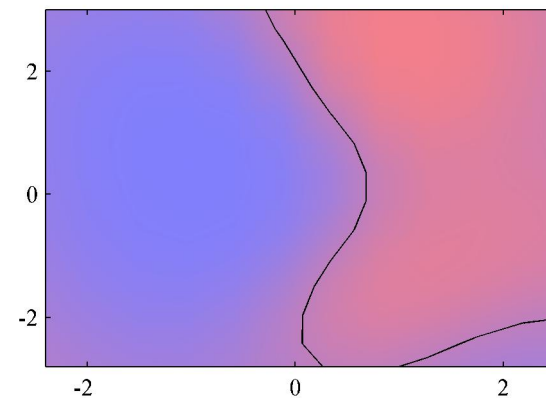- We seek a Gaussian approximation to the posterior distribution over $a_{N+1}$

# Illustration of Gaussian process for classification

Using Laplace Approximation

Green: Optimal boundary
Black: Gaussian Process Classifier

Predicted posterior probability together with Gaussian process decision boundary

# 8. Connection to Neural Networks

- Range of functions represented by a neural network is governed by $M$, number of hidden units

- In a Bayesian neural network
  - Prior distribution over paramter vector w produces a prior distributiion over functions from $y(x)$
  - As $M \to oo$ distribution of functions generated will tend to a Gaussian process

- Neural network outputs share the hidden units thus borrowing statistical strength from each other
  - i.e., weights associated with each hidden unit is influenced by all output variables, not just by one of them
  - This property is lost in the Gaussian process

- For non-stationary kernels
  - ie, $k(x,x')$ cannot be expressed as a function of $x-x'$
  - breaks translation invariance in neural network weight space