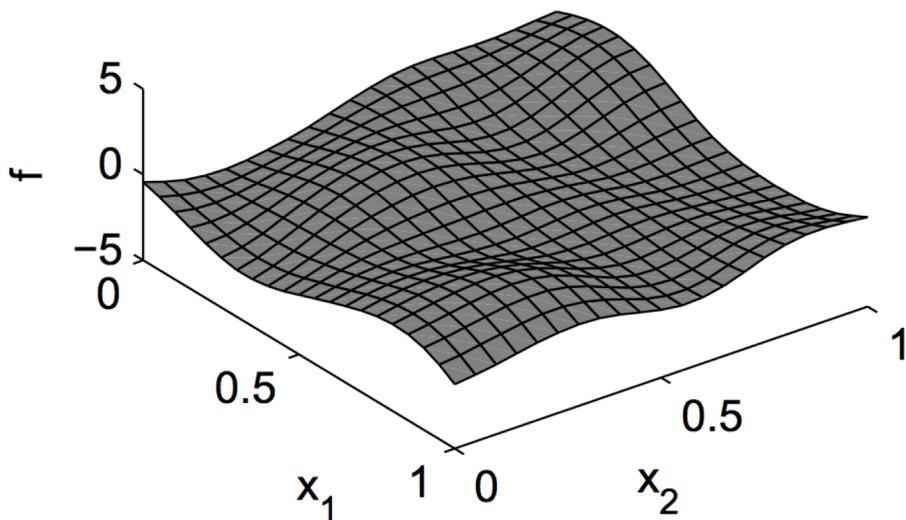
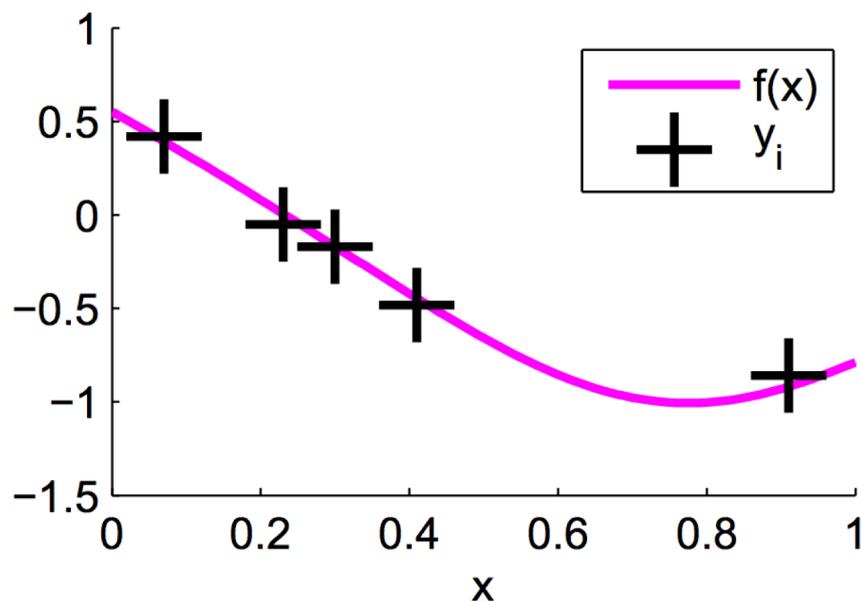


Gaussian Processes

Kayhan Batmanghelich

Goal: Learning a Function

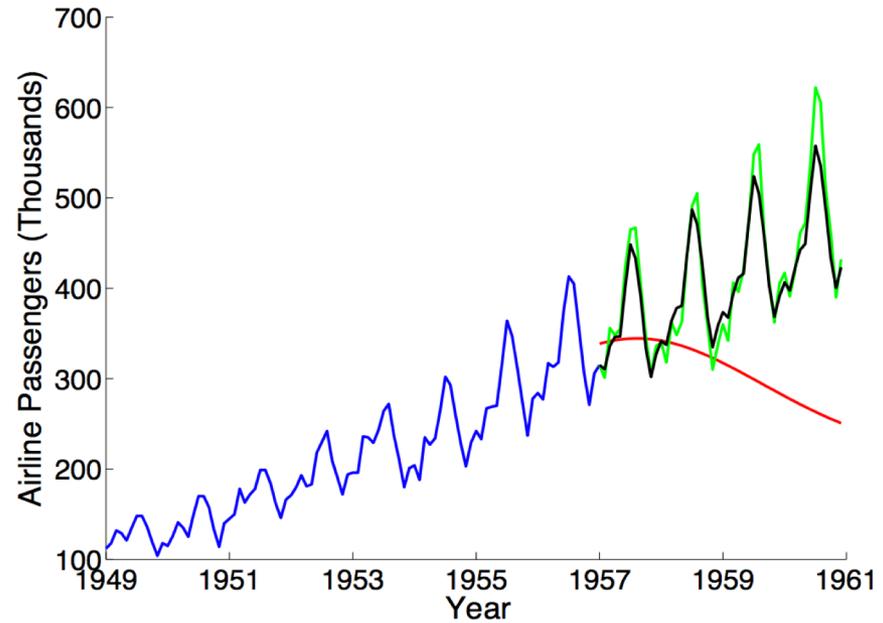
Learn scalar function of vector values $f(\mathbf{x})$



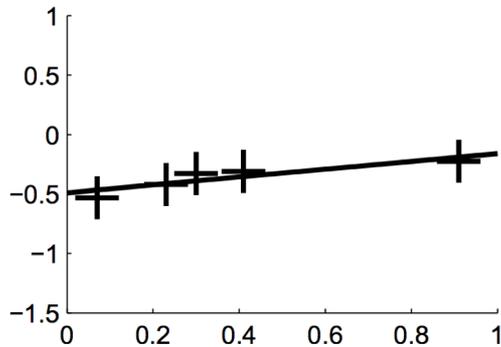
We have (possibly noisy) observations $\{\mathbf{x}_i, y_i\}_{i=1}^n$

A Regression Example

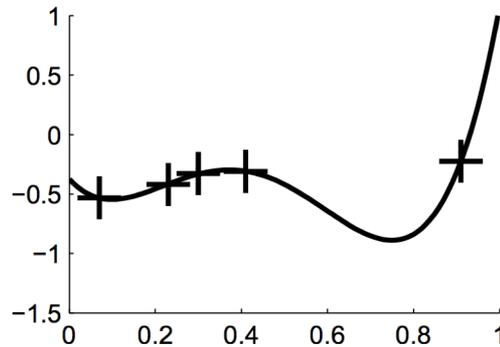
Predict the future:



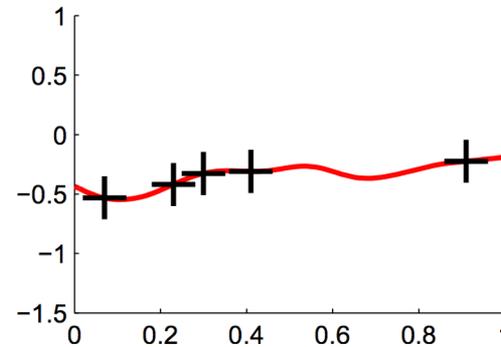
The world is often complicated:



simple fit



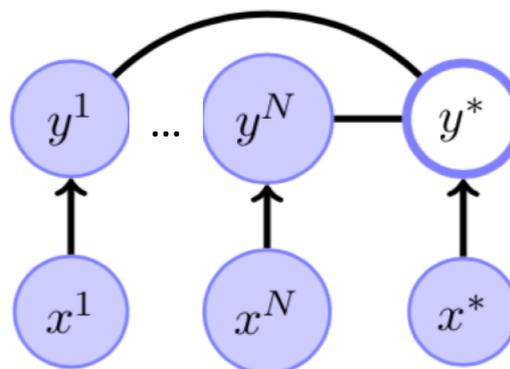
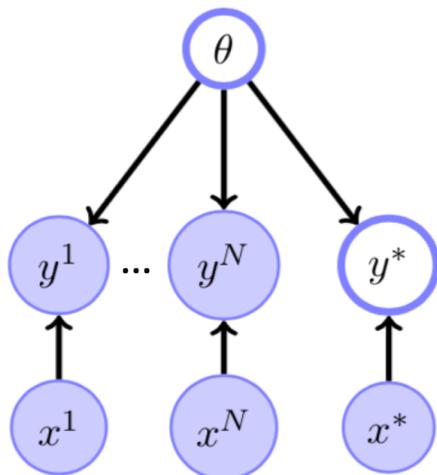
complex fit



truth

Recap: Non-Parametric Approach

Let's revisit prediction (at test time):



Non-parametric methods directly model the joint conditional distribution!

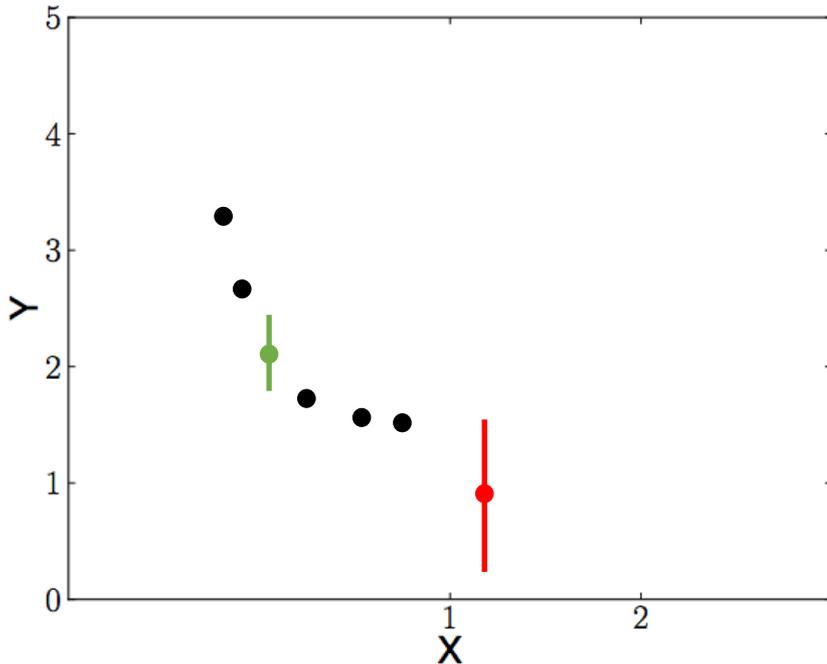
$$p(y^*, \mathcal{Y} | x^*, \mathcal{X}) = \int_{\theta} p(y^* | x^*, \theta) p(\theta) \prod_n p(y^n | \theta, x^n)$$

Recap: Revisiting Linear Regression

Linear Regression with no noise:

$$y = \mathbf{w}^T \underbrace{\phi(\mathbf{x})}_{\text{Feature vector}}$$

$$p(\mathbf{w}) = \mathcal{N}(\mathbf{w} | \mathbf{0}, \Sigma_{\mathbf{w}})$$



- In many applications, one relies on generic **smoothness** assumptions:
 - If two inputs x and x' then outputs (y and y') should be similar.

Recap: Revisiting Linear Regression

Linear Regression with no noise:

$$y = \mathbf{w}^T \underbrace{\phi(\mathbf{x})}_{\text{Feature vector}} \quad p(\mathbf{w}) = \mathcal{N}(\mathbf{w} | \mathbf{0}, \Sigma_{\mathbf{w}})$$

Let's integrate w out to get $p(y^*, \mathcal{Y} | x^*, \mathcal{X})$:

$$\mathbf{y} = [y^1, \dots, y^N] \quad \Phi = [\phi(x^1), \dots, \phi(x^N)]^T$$

$$p(\mathbf{y} | \mathbf{x}) = \mathcal{N}(\mathbf{y} | \mathbf{0}, \mathbf{K})$$


$$[\mathbf{K}]_{n,n'} = \phi(x^n)^T \phi(x^{n'}) = k(x^n, x^{n'}) \quad \text{covariance function}$$

What if instead of $\phi(x)$, we specify $k(x', x)$ directly?

Covariance Function

- The **covariance function** controls the behavior of the prediction function ($x \rightarrow y$) **implicitly**.
- Examples of $k(x, x')$:

$$v_0 \exp \left\{ -\frac{1}{2} \sum_{l=1}^D \lambda_l (x_l - x'_l)^2 \right\} \quad \|x - x'\|^\nu K_\nu(\|x - x'\|)$$

Squared Exponential

Matern

$$\mathbf{x}^\top \mathbf{x}'$$

Linear

Stationary

Non-Stationary

Jupyter Demo!

Covariance Function

- The **covariance function** controls the behavior of the prediction function ($x \rightarrow y$) **implicitly**.
- The $k(x, x')$ specifies the prior over functions.
- Potentially, the $\phi(x)$ function can be infinite dimensional which can be viewed as **basis** function for $f(x)$.
- k is a covariance $\Rightarrow k$ is a positive semi-definite function

$$\int_{\mathcal{X}^2} k(x, x') f(x) f(x') d\mu(x) d\mu(x') \geq 0$$

Theorem (Loeve)

k corresponds to the covariance of a GP



k is a (symmetric) positive definite function

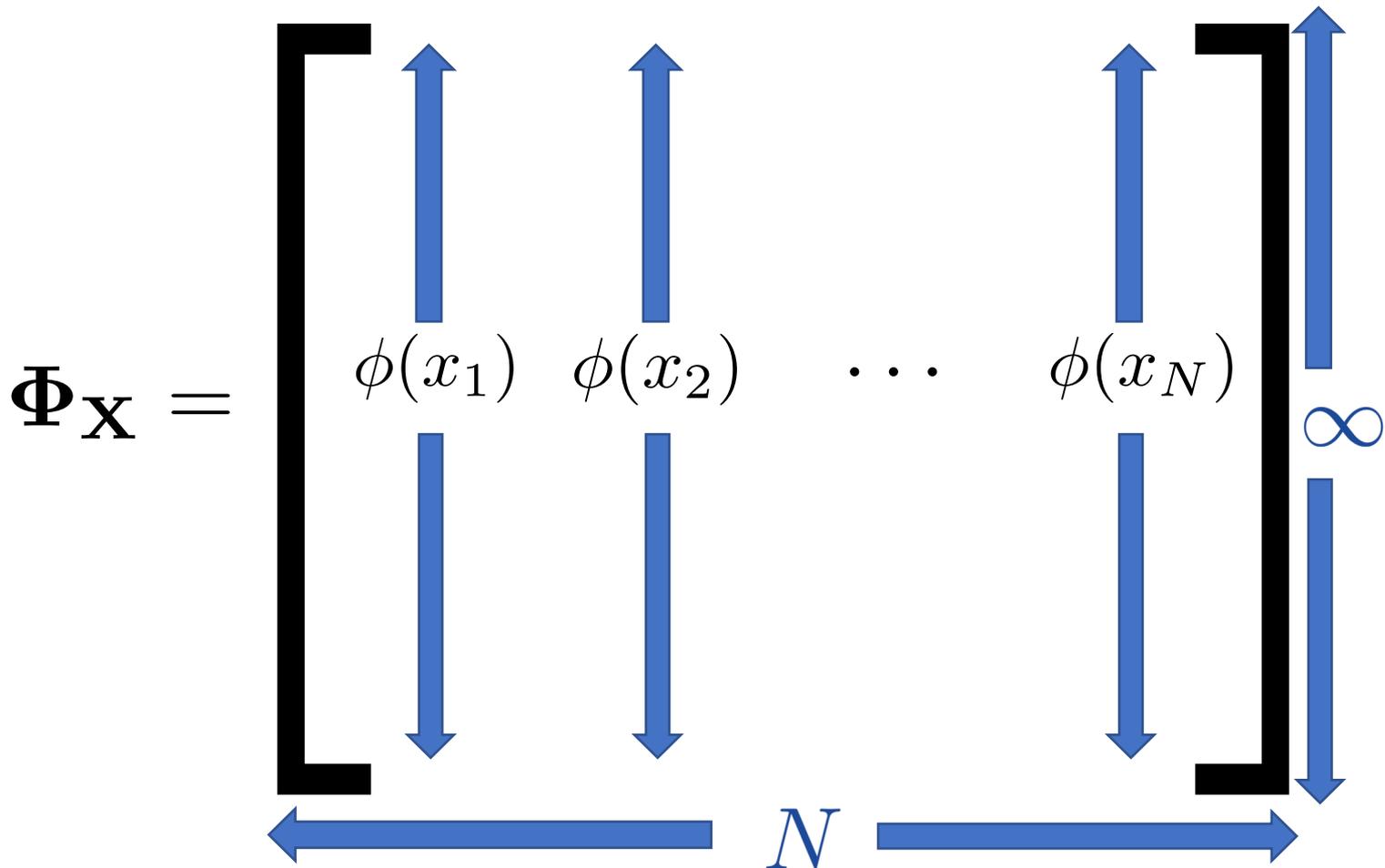
Kernel Matrix

$$K_{XX} = \begin{bmatrix} \langle \phi(x_1), \phi(x_1) \rangle & \dots & \langle \phi(x_1), \phi(x_n) \rangle \\ \vdots & \ddots & \vdots \\ \langle \phi(x_n), \phi(x_1) \rangle & \dots & \langle \phi(x_n), \phi(x_n) \rangle \end{bmatrix}$$

The diagram illustrates the structure of the Kernel Matrix K_{XX} . The matrix is square and symmetric, with elements representing inner products of feature vectors $\phi(x_i)$. The diagonal elements are $\langle \phi(x_i), \phi(x_i) \rangle$. The off-diagonal elements are $\langle \phi(x_i), \phi(x_j) \rangle$. The matrix is enclosed in large black L-shaped brackets. A red dashed box highlights the element $\langle \phi(x_1), \phi(x_n) \rangle$, with a red arrow pointing to it from the label $k(x_1, x_n)$ above. Blue arrows indicate the dimensions: a vertical arrow on the right labeled N and a horizontal arrow at the bottom labeled N .

Implicit Feature Vector

- ϕ_X is an **operator** that maps R^N to functions \mathcal{F} .
- $K_{XX} = \Phi_X^T \Phi_X$



Gaussian Process

- The Gaussian process (GP) as a **prior** on **functions**.
- **Covariance function** and hyperparameters reflect the prior belief on function smoothness, length scales etc.
- A GP is a collection of random variables, any **finite number** of which have a joint **Gaussian** distribution.

Infinite dimension:

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

Finite dimension:

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

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

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

Gaussian Process

- The Gaussian process (GP) as a **prior** on **functions**.
- *Covariance function* and hyperparameters reflect the prior belief on function smoothness, length scales etc.
- A GP is a collection of random variables, any **finite number** of which have a joint **Gaussian** distribution.
- Easy inference for the regression.

Inference in GP: Regression

Given observed **noisy** data $\mathcal{D} = \{\mathbf{X}, \mathbf{y}\}$, the joint probability over latent function values \mathbf{f} and \mathbf{f}^* given \mathbf{y} is

Regression Case:

$$p([\mathbf{f}, \mathbf{f}^*] | \mathbf{X}, \mathbf{X}^*, \mathbf{y}, \boldsymbol{\theta}_K, \sigma^2) \propto \underbrace{\mathcal{N}([\mathbf{f}, \mathbf{f}^*] | \mathbf{0}, \begin{bmatrix} \mathbf{K}_{\mathbf{X}, \mathbf{X}} & \mathbf{K}_{\mathbf{X}, \mathbf{X}^*} \\ \mathbf{K}_{\mathbf{X}^*, \mathbf{X}} & \mathbf{K}_{\mathbf{X}^*, \mathbf{X}^*} \end{bmatrix})}_{\text{Prior}} \times \underbrace{\prod_{n=1}^N \mathcal{N}(y_n | f_n, \sigma^2)}_{\text{Likelihood}},$$

$$\mathbf{f}(\mathbf{X}) = \begin{bmatrix} f(\mathbf{x}_1) \\ f(\mathbf{x}_2) \\ \vdots \end{bmatrix}$$

$$\mathbf{f}^*(\mathbf{X}^*) = \begin{bmatrix} f^*(\mathbf{x}_1^*) \\ f^*(\mathbf{x}_2^*) \\ \vdots \end{bmatrix}$$

Inference in GP: Regression

Given observed **noisy** data $\mathcal{D} = \{\mathbf{X}, \mathbf{y}\}$, the joint probability over latent function values \mathbf{f} and \mathbf{f}^* given \mathbf{y} is

Regression Case:

$$p([\mathbf{f}, \mathbf{f}^*] | \mathbf{X}, \mathbf{X}^*, \mathbf{y}, \boldsymbol{\theta}_K, \sigma^2) \propto \mathcal{N} \left([\mathbf{y}, \mathbf{f}^*] \mid \mathbf{0}, \underbrace{\begin{bmatrix} \mathbf{K}_{\mathbf{X}, \mathbf{X}} + \sigma^2 \mathbf{I} & \mathbf{K}_{\mathbf{X}, \mathbf{X}^*} \\ \mathbf{K}_{\mathbf{X}^*, \mathbf{X}} & \mathbf{K}_{\mathbf{X}^*, \mathbf{X}^*} \end{bmatrix}}_{\text{Prior}} \right) \\ \times \underbrace{\prod_{n=1}^N \mathcal{N}(y_n | f_n, \sigma^2)}_{\text{Likelihood}},$$

$$p(\mathbf{f}^* | \mathbf{X}, \mathbf{y}, \mathbf{X}^*, \boldsymbol{\theta}_K, \sigma^2) = \mathcal{N}(\mathbf{f}^* | \boldsymbol{\mu}^*, \boldsymbol{\Sigma}^*) \text{ with}$$

$$\boldsymbol{\mu}^* = \mathbf{K}_{\mathbf{X}^*, \mathbf{X}} [\mathbf{K}_{\mathbf{X}, \mathbf{X}} + \sigma^2 \mathbf{I}]^{-1} \mathbf{y}$$

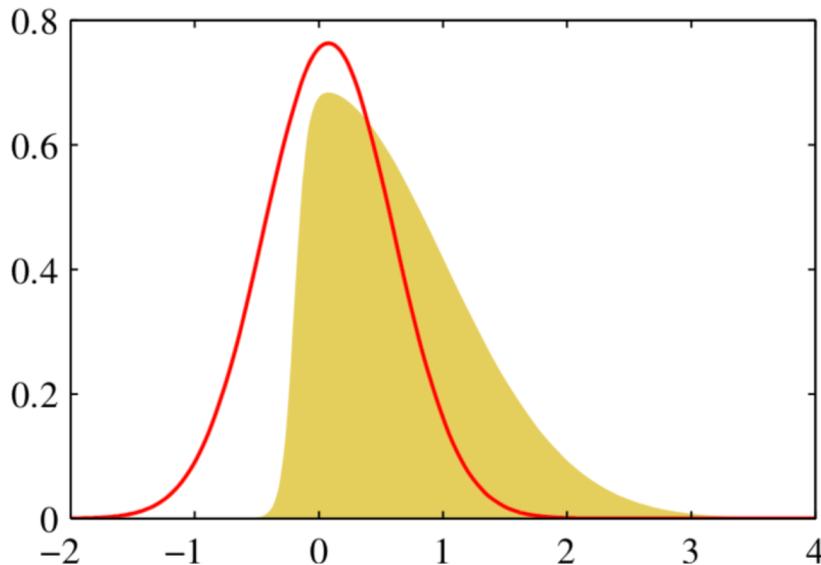
$$\boldsymbol{\Sigma}^* = \mathbf{K}_{\mathbf{X}^*, \mathbf{X}^*} - \mathbf{K}_{\mathbf{X}^*, \mathbf{X}} [\mathbf{K}_{\mathbf{X}, \mathbf{X}} + \sigma^2 \mathbf{I}]^{-1} \mathbf{K}_{\mathbf{X}, \mathbf{X}^*}$$

Inference in GP: Classification

Given observed **noisy** data $\mathcal{D} = \{\mathbf{X}, \mathbf{y}\}$, the joint probability over latent function values \mathbf{f} and \mathbf{f}^* given \mathbf{y} is

Classification Case:

$$p([\mathbf{f}, \mathbf{f}^*] | \mathbf{X}, \mathbf{X}^*, \mathbf{y}, \boldsymbol{\theta}_K, \sigma^2) \propto \mathcal{N} \left([\mathbf{f}, \mathbf{f}^*] \mid \mathbf{0}, \overbrace{\begin{bmatrix} \mathbf{K}_{\mathbf{X}, \mathbf{X}} & \mathbf{K}_{\mathbf{X}, \mathbf{X}^*} \\ \mathbf{K}_{\mathbf{X}^*, \mathbf{X}} & \mathbf{K}_{\mathbf{X}^*, \mathbf{X}^*} \end{bmatrix}}^{\text{Prior}} \right)$$



$$\times \prod_{n=1}^N \underbrace{\mathcal{N}(y_n | f_n, \sigma^2)}_{\text{Likelihood}},$$

$$p(y_n = 1 | f_n) = \frac{1}{1 + \exp(-f_n)}$$

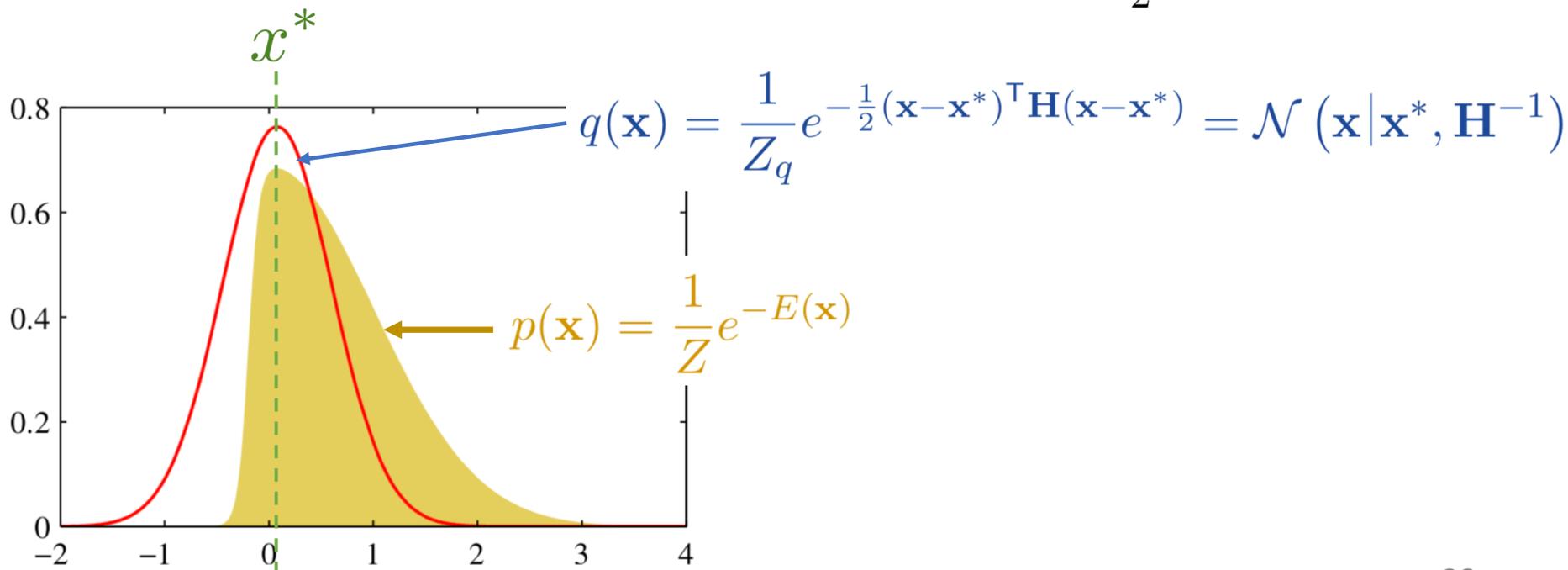
What if we approximate the likelihood term to look like a Gaussian? **Not closed form!**

Laplace Method

Given observed **noisy** data $\mathcal{D} = \{\mathbf{X}, \mathbf{y}\}$, the joint probability over latent function values \mathbf{f} and \mathbf{f}^* given \mathbf{y} is

Use **Taylor** expansion around \mathbf{x}^* :

$$E(\mathbf{x}) \approx E(\mathbf{x}^*) + (\mathbf{x} - \mathbf{x}^*)^\top \nabla E|_{\mathbf{x}^*} + \frac{1}{2} (\mathbf{x} - \mathbf{x}^*)^\top \mathbf{H} (\mathbf{x} - \mathbf{x}^*)$$



Learning

Using maximum likelihood to learn the parameters of the kernel (θ)?

$$p(\mathbf{y}|\mathbf{X}) = \int p(\mathbf{y}|\mathbf{f}, \mathbf{X}) p(\mathbf{f}|\mathbf{X}) d\mathbf{f}$$

$$p(\mathbf{y}|\mathbf{f}) = \prod_i \mathcal{N}(y_i|f_i, \sigma_y^2) \quad p(\mathbf{f}|\mathbf{X}) = \mathcal{N}(\mathbf{f}|\mathbf{0}, \mathbf{K})$$

$$\log p(\mathbf{y}|\mathbf{X}) = \log \mathcal{N}(\mathbf{y}|\mathbf{0}, \mathbf{K}_y) = -\frac{1}{2} \mathbf{y} \mathbf{K}_y^{-1} \mathbf{y} - \frac{1}{2} \log |\mathbf{K}_y| - \frac{N}{2} \log(2\pi)$$

\mathbf{K}_y absorbs the effect of noise variance

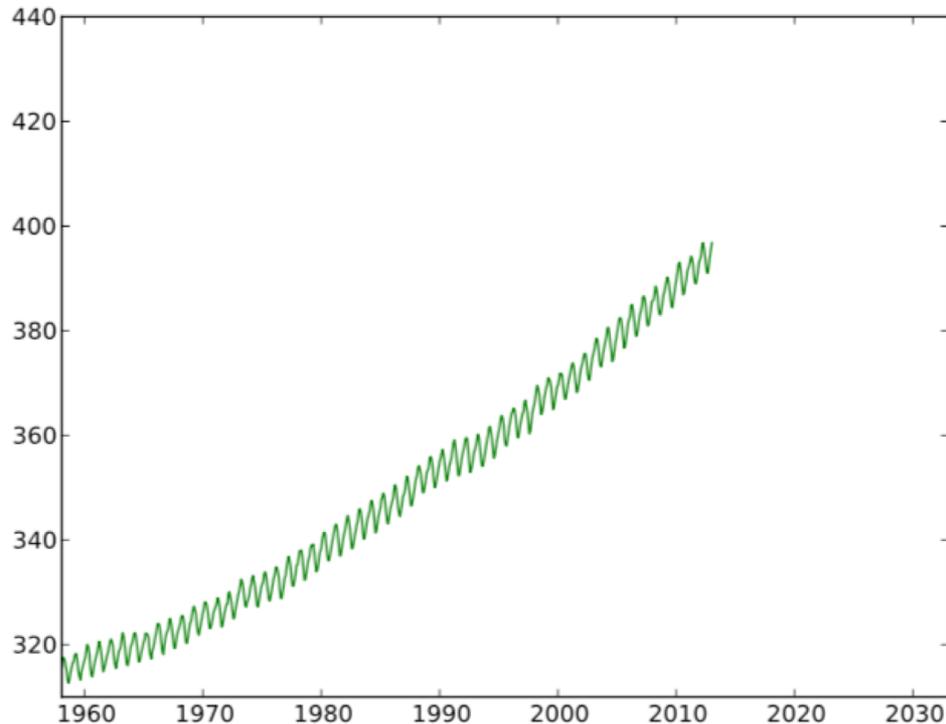
$$\frac{\partial}{\partial \theta_j} \log p(\mathbf{y}|\mathbf{X}) = \frac{1}{2} \mathbf{y}^T \mathbf{K}_y^{-1} \frac{\partial \mathbf{K}_y}{\partial \theta_j} \mathbf{K}_y^{-1} \mathbf{y} - \frac{1}{2} \text{tr}(\mathbf{K}_y^{-1} \frac{\partial \mathbf{K}_y}{\partial \theta_j})$$

How to Choose $k(x, x')$?

Example:

Data: This famous dataset compiles the monthly CO_2 concentration in Hawaii since 1958.

Task: Predict the concentration for the next 20 years.

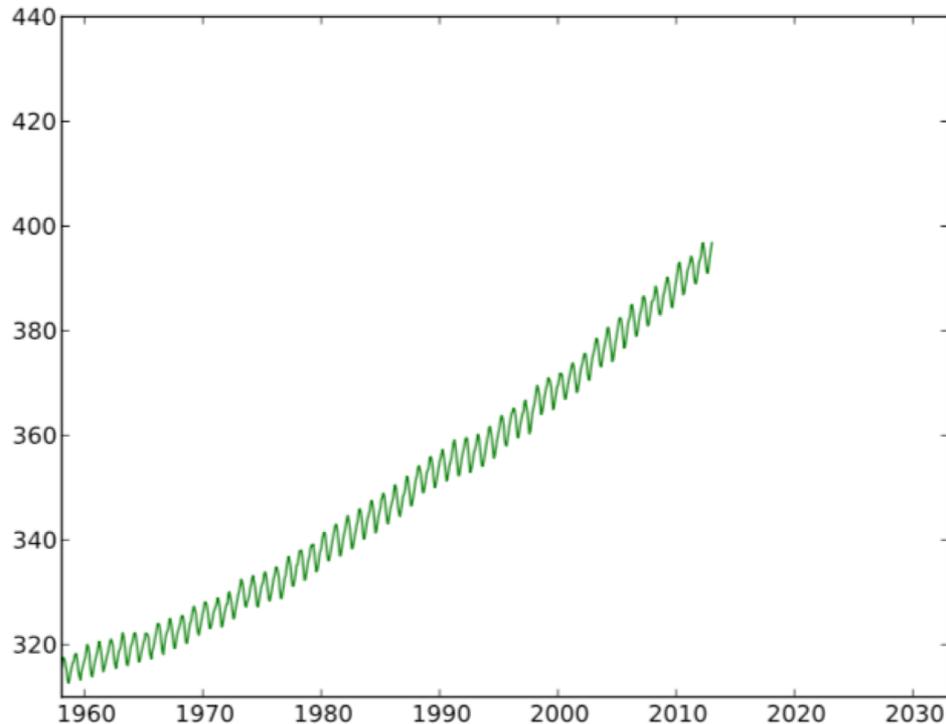


How to Choose $k(x, x')$?

Example:

Data: This famous dataset compiles the monthly CO_2 concentration in Hawaii since 1958.

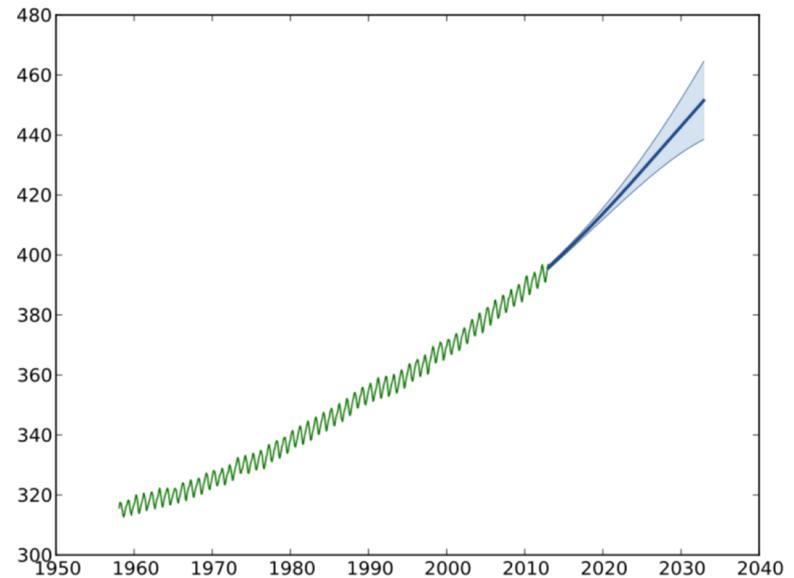
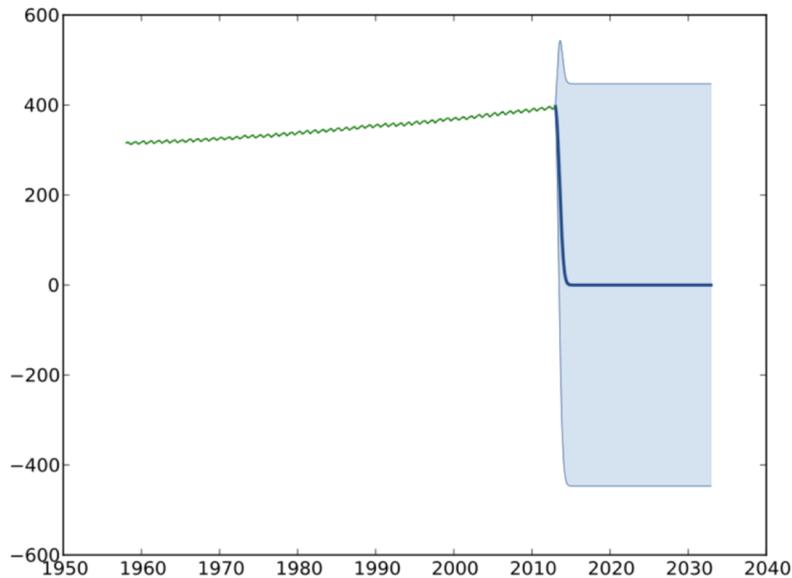
Task: Predict the concentration for the next 20 years.



How to Choose $k(x, x')$?

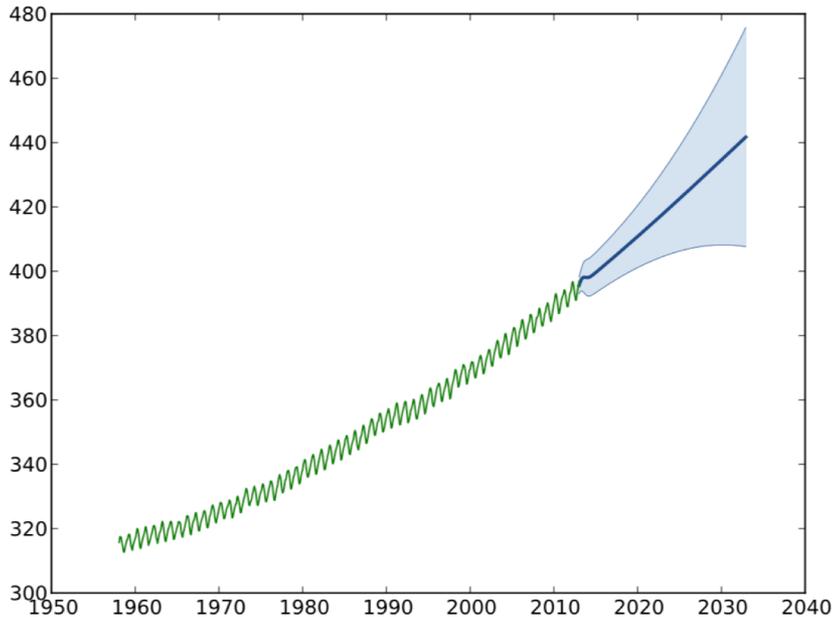
$$k(x, y) = \sigma^2 \exp\left(-\frac{(x - y)^2}{\theta^2}\right)$$

Terrible!



How to Choose $k(x, x')$?

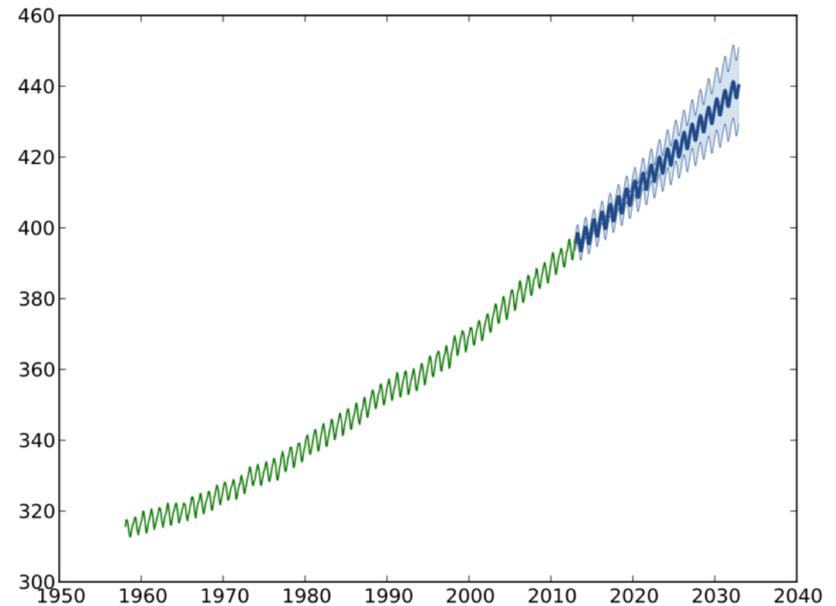
$$k(x, y) = k_{rbf1}(x, y) + k_{rbf2}(x, y)$$



Meh...!

What about the oscillations?

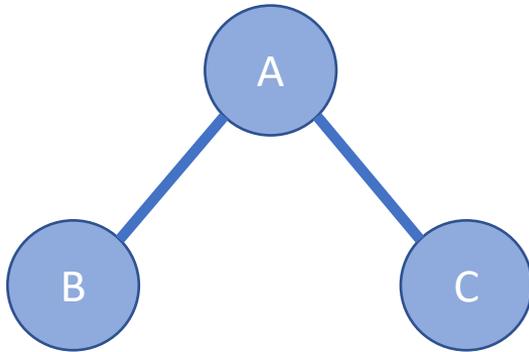
$$k(x, y) = \sigma_0^2 x^2 y^2 + k_{rbf1}(x, y) + k_{rbf2}(x, y) + k_{per}(x, y)$$



Much better!

How can we extend this idea to
Graphical Model?

General Ideas



- Parametric forms of conditional

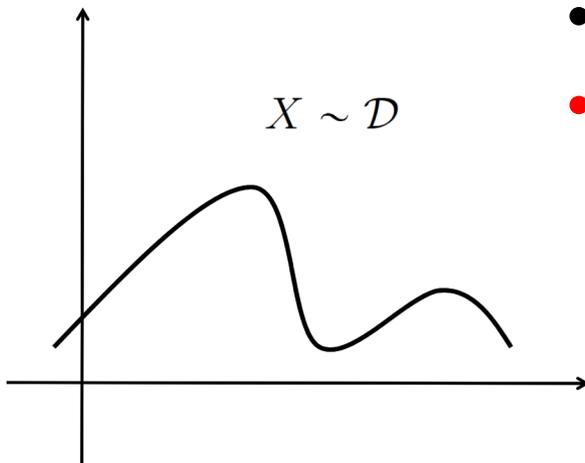
Discrete

$\mathbb{P}[C = 0 A = 0]$	$\mathbb{P}[C = 0 A = 1]$
$\mathbb{P}[C = 1 A = 0]$	$\mathbb{P}[C = 1 A = 1]$

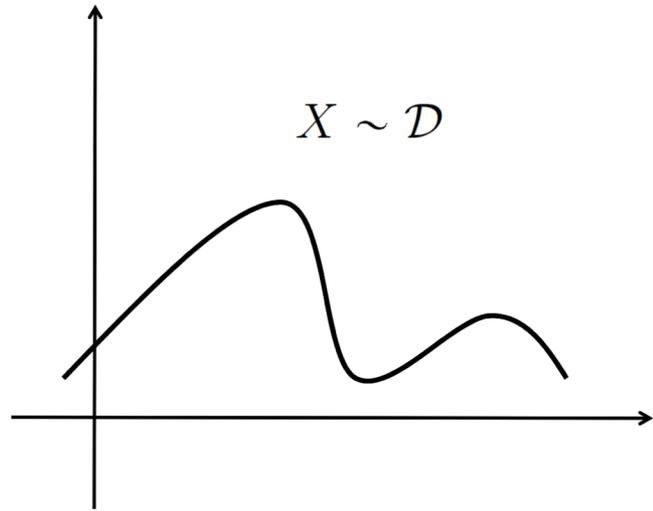
Continuous

$$\mathcal{N}(\mu, \sigma)$$

- Can we make it **non-parametric**?
- **Key idea**: represent this distribution with a small vector μ_X



Key Idea



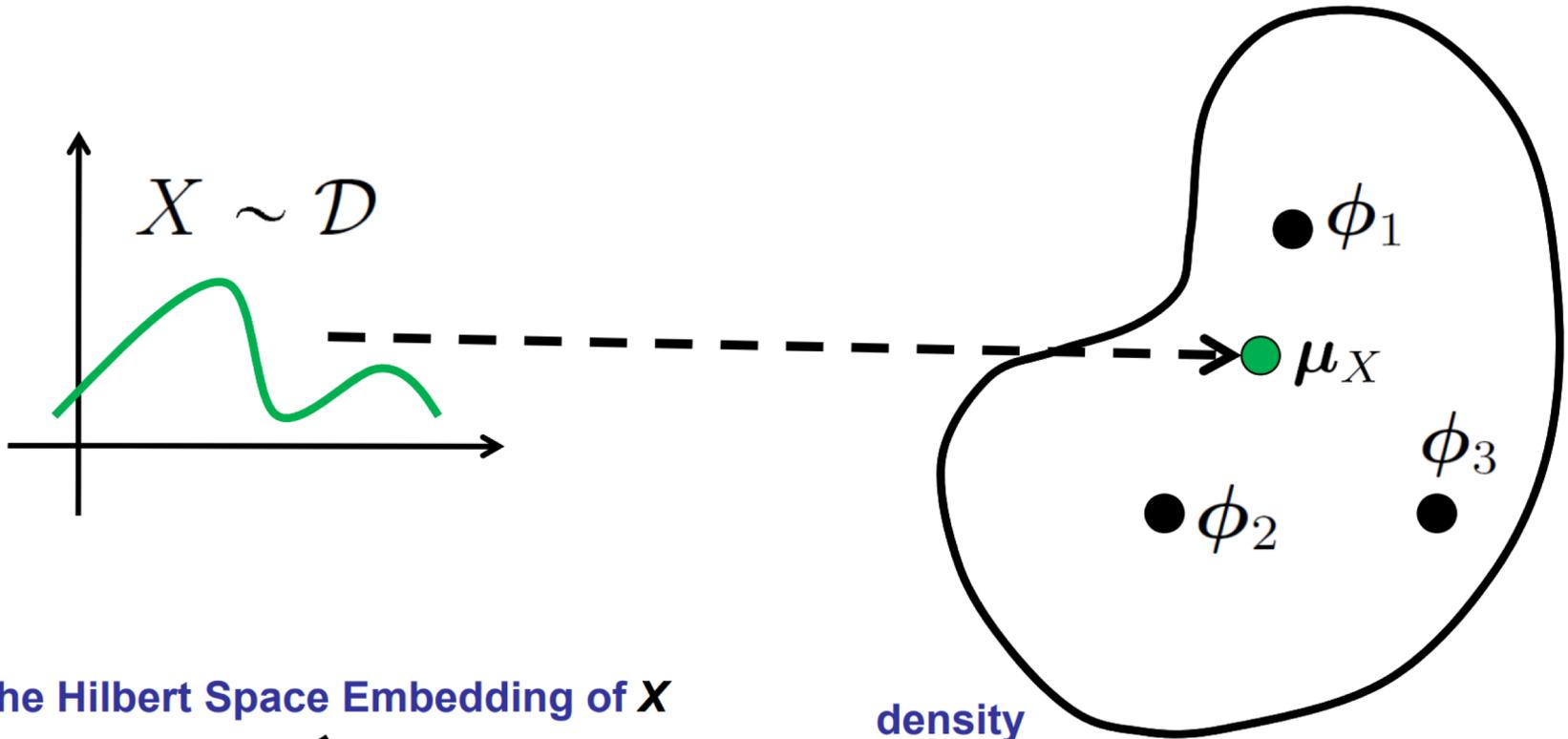
$$\mu_X = \begin{pmatrix} \mathbb{E}[X] \\ \mathbb{E}[X^2] \\ \mathbb{E}[X^3] \\ \vdots \end{pmatrix}$$

Not enough!

Right idea but not exactly....

- ϕ_X is an infinite dimensional vector space.

Embedding Distributions



The Hilbert Space Embedding of X

$$\mu_X(\cdot) = \mathbb{E}_{X \sim \mathcal{D}}[\phi_X] = \int \text{density} \, p_{\mathcal{D}}(X) \phi_X(\cdot) dX$$

Embedding Distributions

We leave detail discussion to a future. Here is a general idea:

Original Space

 X Y (X, Y) $\mathbb{P}(X)$ $\mathbb{P}(Y)$ $\mathbb{P}(X, Y)$

RKHS Space

 $\phi_X \in \mathcal{F}$ $\psi_Y \in \mathcal{G}$ $(\phi_X, \psi_Y) \in \mathcal{F} \times \mathcal{G}$ $\mu_X = \mathbb{E}[\phi_X]$ $\mu_Y = \mathbb{E}[\psi_Y]$ $\boxed{\mathcal{C}_{YX}} = \mathbb{E}[\psi_Y \otimes \phi_X]$

cross-covariance **operator**

Embedding Distributions

We leave detail discussion to a future. Here is a general idea:

Original Space

$$\mathbb{P}(X, Y)$$

$$\text{Diag}[P(X)]$$

$$\mathbb{P}[Y|X] = \mathbb{P}[Y, X] \times \text{Diag}(\mathbb{P}[X])^{-1}$$

$$\mathbb{P}[X] = \int_Y \mathbb{P}[X, Y] = \int_Y \mathbb{P}[X|Y]\mathbb{P}[Y]$$

$$\mathbb{P}[X, Y] = \mathbb{P}[X|Y]\mathbb{P}[Y]$$

RKHS Space

$$\mathcal{C}_{YX} = \mathbb{E}[\psi_Y \otimes \phi_X]$$

$$\mathcal{C}_{XX} = \mathbb{E}[\phi_X \otimes \phi_X]$$

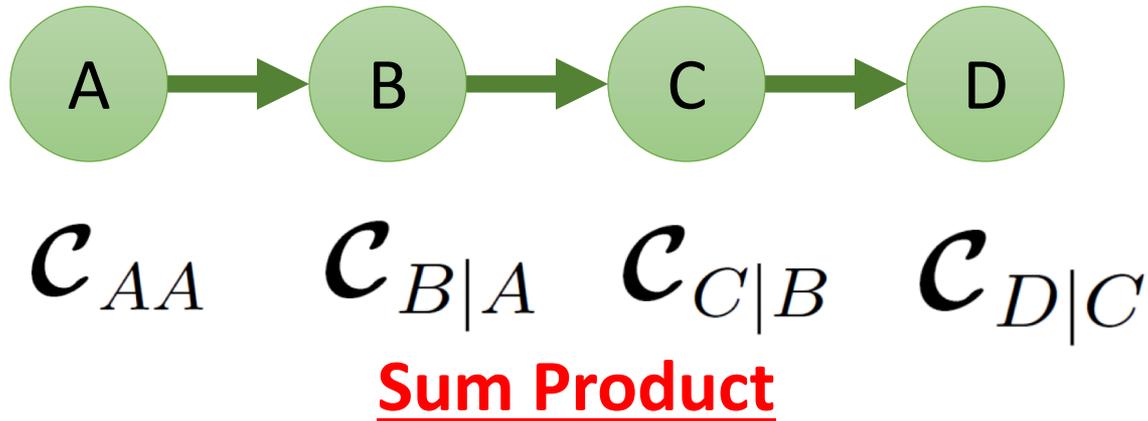
$$\mathcal{C}_{Y|X} = \mathcal{C}_{YX}\mathcal{C}_{XX}^{-1}$$

$$\mu_X = \mathcal{C}_{X|Y}\mu_Y$$

$$\mathcal{C}_{YX} = \mathcal{C}_{Y|X}\mathcal{C}_{XX}$$

Kernel Graphical Models

- Replacing CPT with RKHS operators/function



With discrete variables:

$$\mathcal{P}(A, D) = \mathcal{P}(\emptyset A) \mathcal{P}(B|A)^\top \mathcal{P}(C|B)^\top \mathcal{P}(D|C)^\top$$

With kernels:

$$\mathcal{C}_{AD} = \mathcal{C}_{AA} \mathcal{C}_{B|A}^\top \mathcal{C}_{B|C}^\top \mathcal{C}_{C|D}^\top$$

Summary

- Hilbert Space Embedding provides a way to create a “sufficient statistic” for an arbitrary distribution.
- Can embed marginal, joint, and conditional distributions into the RKHS.
- **More references:**
 - Smola, A. J., Gretton, A., Song, L., and Schölkopf, B., [A Hilbert Space Embedding for Distributions](#), Algorithmic Learning Theory, E. Takimoto (Eds.), Lecture Notes on Computer Science, Springer, 2007.
 - L. Song. [Learning via Hilbert space embedding of distributions](#). PhD Thesis 2008.
 - Song, L., Huang, J., Smola, A., and Fukumizu, K., [Hilbert space embeddings of conditional distributions](#), International Conference on Machine Learning, 2009.