

# Gaussian Process (Bayesian Learning meets Kernels)

CS772A: Probabilistic Machine Learning

Piyush Rai

# Discriminative Models

The function  $f$  gives a real-valued **score**  $f(\mathbf{x})$  to each input  $\mathbf{x}$ . Parameters of the likelihood model  $p(\mathbf{y}|f, \mathbf{x})$  depend on this score

- Discriminative models learn a function  $f$  that maps inputs  $\mathbf{x}$  to outputs  $\mathbf{y}$

$$p(\mathbf{y}|f, \mathbf{x}) = \mathcal{N}(\mathbf{y}|f(\mathbf{x}), \beta^{-1})$$

$$p(\mathbf{y}|f, \mathbf{x}) = [\sigma(f(\mathbf{x}))]^y [1 - \sigma(f(\mathbf{x}))]^{1-y}$$

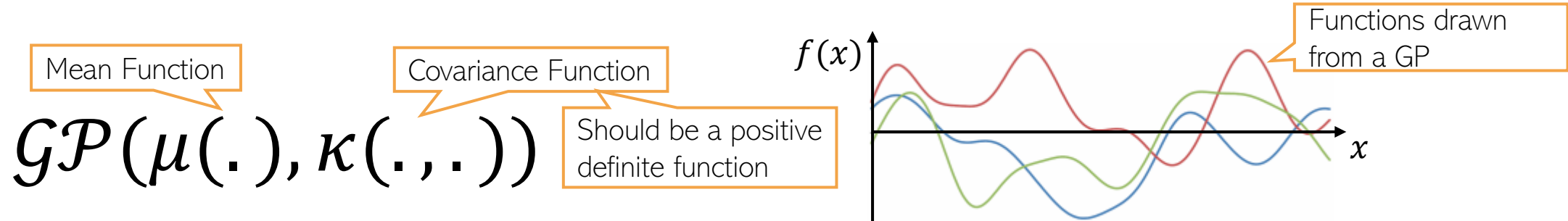
$$p(\mathbf{y}|f, \mathbf{x}) = \text{ExpFam}(f(\mathbf{x}))$$

- We usually define the function  $f$  using weights (i.e., the model parameters), e.g.,
  - Linear:  $f(\mathbf{x}) = \mathbf{w}^T \mathbf{x}$ , used in linear/logistic regression and GLM
  - Nonlinear:  $f(\mathbf{x}) = \text{NN}(\mathbf{x}, \mathbf{w})$ , used in deep neural nets based models
- Since these models use parameters, we call them “**parametric**” models
- We can also define the function  $f$  “directly” without using parameters
  - Essentially, this can be done using a “**nonparametric**” approach
  - It would be similar to nearest neighbors or kernel SVMs which are also “nonparametric”
- Gaussian Process (GP) is such a **Bayesian nonparametric** approach



# Gaussian Process (GP)

- A Gaussian Process (GP) defines a **distribution over functions** and is denoted as



- Mean function defines what functions drawn from this GP look like on average

$$\mu(x) = \mathbb{E}_{f \sim \mathcal{GP}(\mu, \kappa)}[f(x)]$$

- Covariance/kernel function defines the similarity between a pair of inputs

$$K_{ij} = \kappa(x_i, x_j)$$

- Covariance/kernel function controls the shape of the functions drawn from GP
  - If  $\kappa(x_i, x_j)$  is high then  $f(x_i)$  and  $f(x_j)$  will have similar values



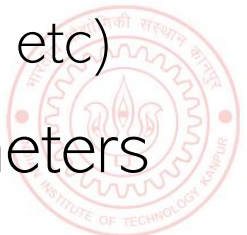
# Covariance/kernel functions

- Kernel functions are popular in learning nonlinear functions (e.g., kernel SVMs)
- Using a kernel corresponds to applying a nonlinear mapping function  $\phi$  on inputs, s.t.,

Kernel function  $\kappa(\cdot, \cdot)$  gives the pairwise similarity between **two inputs  $x$  and  $x'$**  in the new feature space defined by mapping function  $\phi(\cdot)$

$$\kappa(x, x') = \phi(x)^\top \phi(x')$$

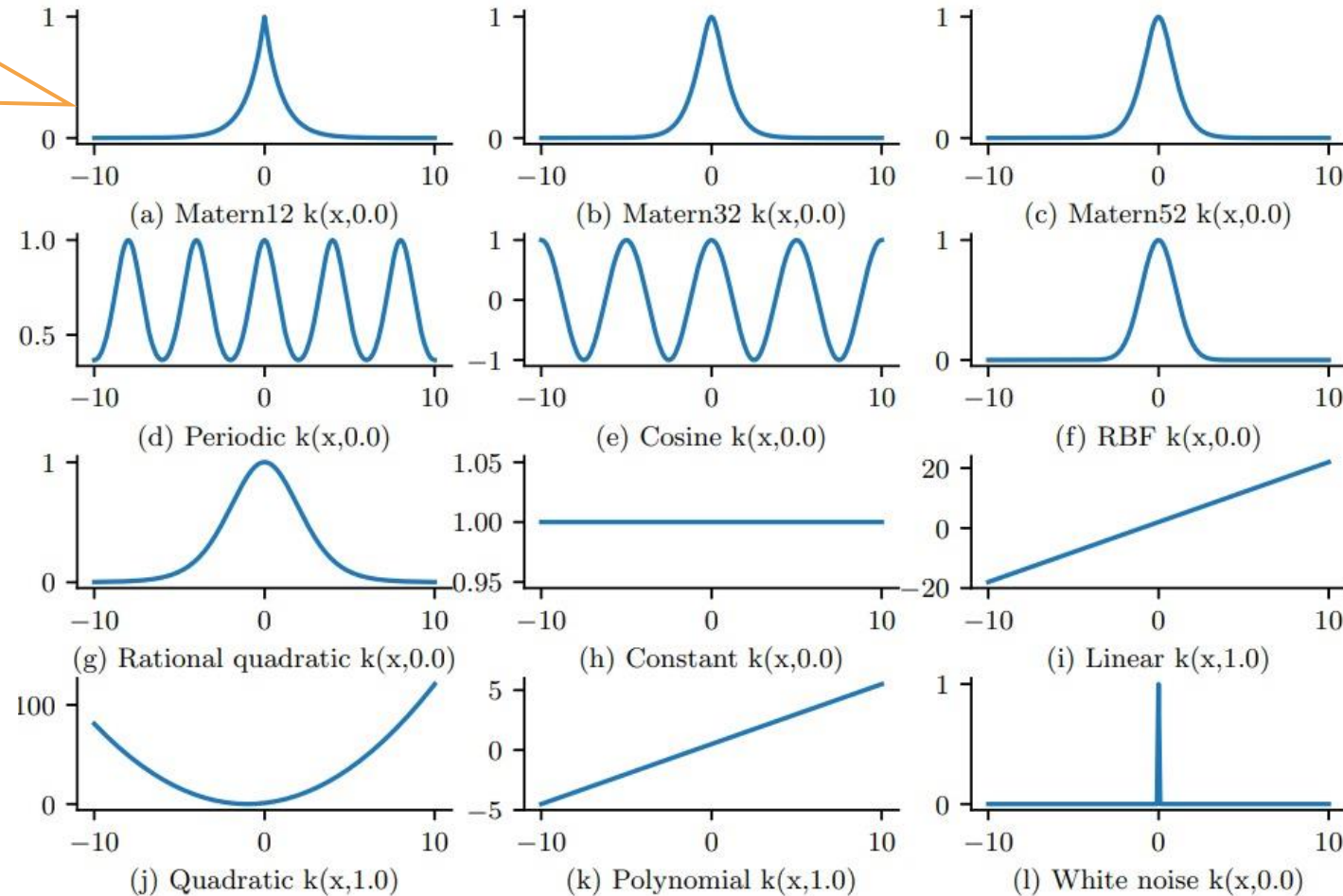
- A wide variety of kernel functions exists that suit different types of data
  - Linear kernel, polynomial kernel, Squared exponential (RBF) kernel
  - Automatic Relevance Determination (ARD) kernel
  - Matérn kernel, Periodic kernel, and many others
- We can combine multiple kernels and use them for GP, e.g., possible kernels
  - $\kappa_1 + \kappa_2$ ,  $\kappa_1 \times \kappa_2$ ,  $\alpha\kappa_1 + \kappa_2$ , etc (add, mult, positive scalar mult, or combinations, etc)
- We can learn how to combine multiple kernels and learn their hyperparameters
  - Possible naturally with the Bayesian approach



# Covariance/kernel functions

- Visualization of some kernel functions (how similarity changes with “distance”)

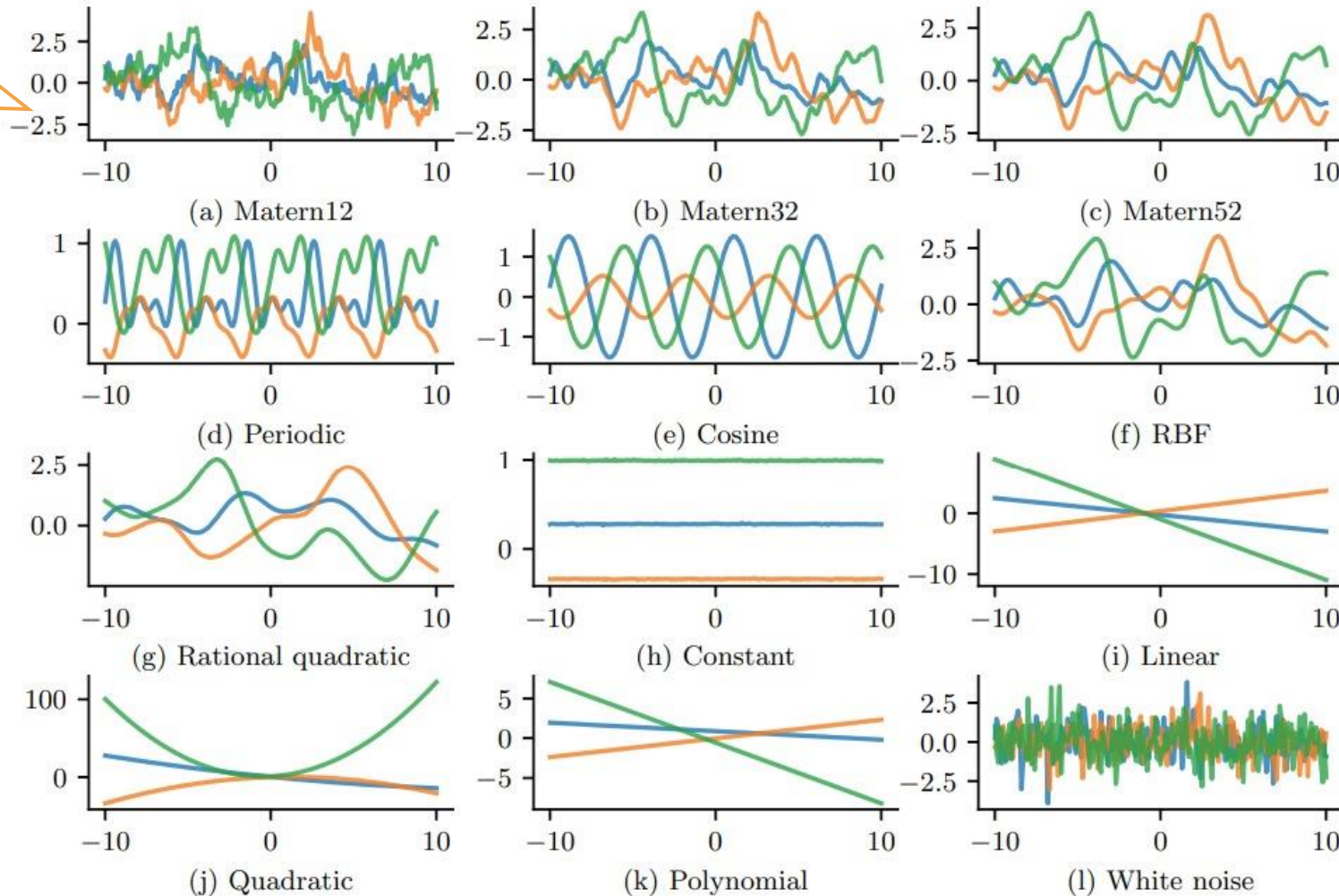
Kernel for all the plots is  $\kappa(x, x')$ , i.e., fixing one input as  $x'$  and varying the other input  $x$  (along the X axis)



# Covariance/kernel functions

- Examples of functions  $f \sim \text{GP}(\mu, \kappa)$  drawn from a GP using some standard kernels

Each plot shows 3 random functions drawn from the corresponding GP with the specified kernel function



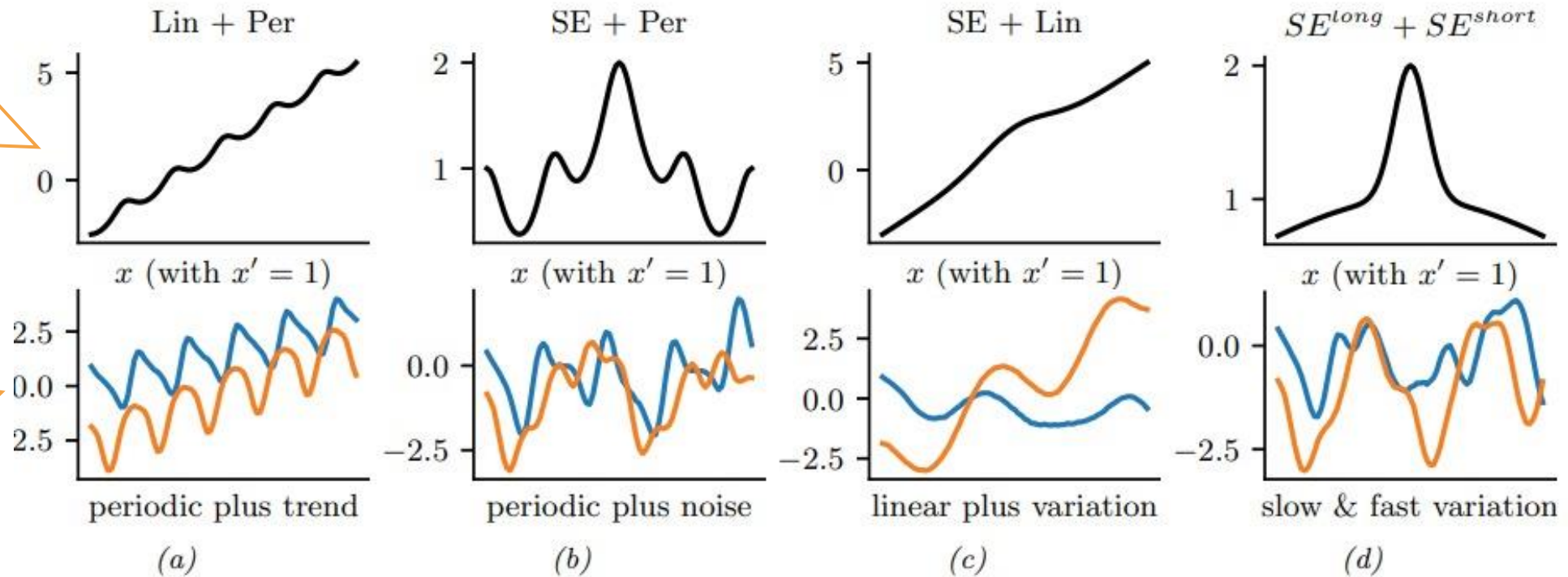


# Combining two (or more) covariance/kernel functions

- Adding two kernels and its effect on the function  $f$  defined by the resulting kernel

Kernel for all 4 plots is  $\kappa(x, x')$ , i.e., fixing  $x'$  as 1 and varying the other input  $x$  (along the X axis)

Two randomly drawn functions from the GP with this kernel combination

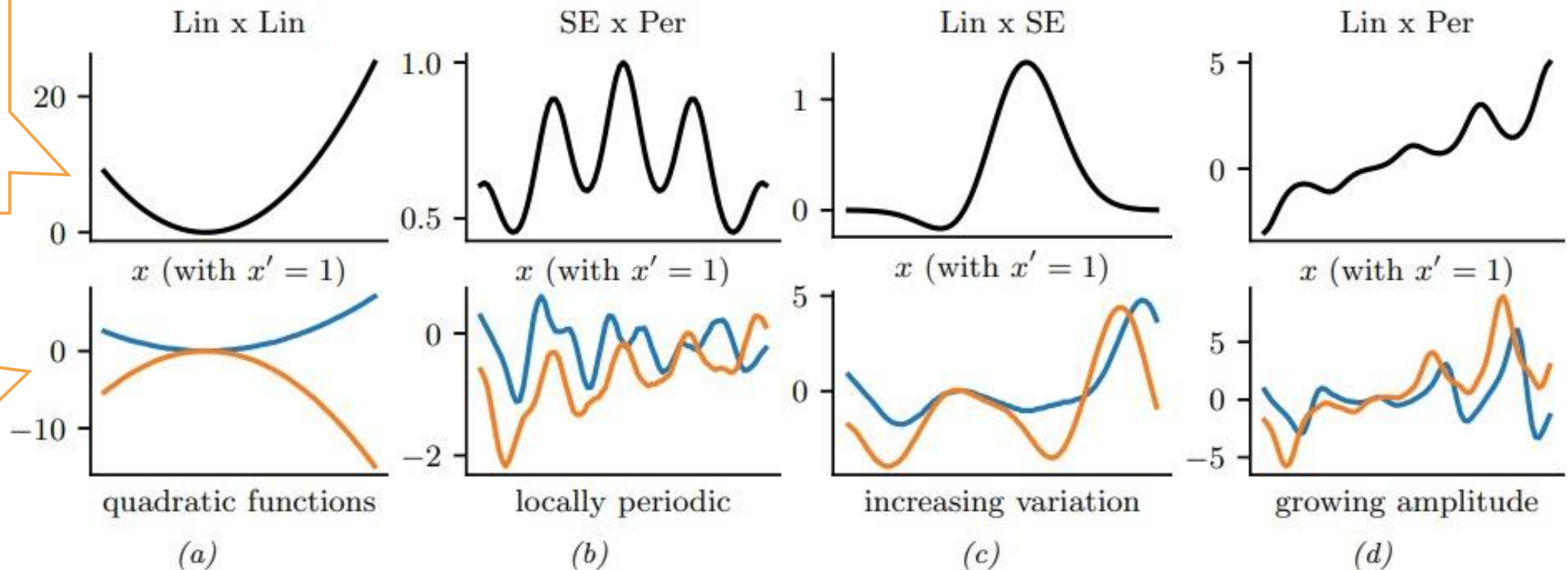


# Combining two (or more) covariance/kernel functions

- Multiplying two kernels and its effect on the function  $f$  defined by the resulting kernel

Kernel for all 4 plots is  $\kappa(x, x')$ , i.e., fixing  $x'$  as 1 and varying the other input  $x$  (along the X axis)

Two randomly drawn functions from the GP with this kernel combination





# Gaussian Process: The Predictive Model

- If  $f \sim \mathcal{GP}(\mu, \kappa)$  then  $f$ 's value at any finite set of inputs is jointly Gaussian

$$p\left(\begin{bmatrix} f(x_1) \\ f(x_2) \\ \vdots \\ f(x_N) \end{bmatrix}\right) = \mathcal{N}\left(\begin{bmatrix} \mu(x_1) \\ \mu(x_2) \\ \vdots \\ \mu(x_N) \end{bmatrix}, \begin{bmatrix} \kappa(x_1, x_1) & \dots & \kappa(x_1, x_N) \\ \kappa(x_2, x_1) & \dots & \kappa(x_2, x_N) \\ \vdots & \ddots & \vdots \\ \kappa(x_N, x_1) & \dots & \kappa(x_N, x_N) \end{bmatrix}\right) \longrightarrow p(\mathbf{f}) = \mathcal{N}(\mathbf{f} | \boldsymbol{\mu}, \mathbf{K})$$

$N \times 1$        $N \times 1$        $N \times N$

- Denoting  $f$ 's score for a new test input  $\mathbf{x}_*$  as  $f_* = f(\mathbf{x}_*)$ , we must also have

$$p\left(\begin{bmatrix} \mathbf{f} \\ f_* \end{bmatrix}\right) = \mathcal{N}\left(\begin{bmatrix} \boldsymbol{\mu} \\ \mu_* \end{bmatrix}, \begin{bmatrix} \mathbf{K} & \mathbf{k}_* \\ \mathbf{k}_*^\top & \kappa(x_*, x_*) \end{bmatrix}\right)$$

$\mathbf{k}_* = [\kappa(x_1, x_*), \kappa(x_2, x_*), \dots, \kappa(x_N, x_*)]^\top$   
 $N + 1$  dim Gaussian

- Very useful result: Easy to see that, given the above, the GP predictive distribution

$$p(f_* | \mathbf{f}) = \mathcal{N}(f_* | \mu_* + \mathbf{k}_*^\top \mathbf{K}^{-1}(\mathbf{f} - \boldsymbol{\mu}), \kappa(x_*, x_*) - \mathbf{k}_*^\top \mathbf{K}^{-1} \mathbf{k}_*)$$

- Thus score of  $f$  on  $\mathbf{x}_*$  given scores on training inputs has a Gaussian distribution

# Gaussian Process: The Predictive Model

- Assuming the mean function  $\mu(\mathbf{x}) = \mathbf{0}$ , the conditional distribution of score becomes

$$p(f_* | \mathbf{f}) = \mathcal{N}(f_* | \mathbf{k}_*^\top \mathbf{K}^{-1} \mathbf{f}, \kappa(x_*, x_*) - \mathbf{k}_*^\top \mathbf{K}^{-1} \mathbf{k}_*) = \mathcal{N}(f_* | \hat{\mu}, \hat{\sigma}^2)$$

- Note that the predictive mean  $\hat{\mu}$  can be written in the following two equivalent ways

Weighted sum of the scores of the  $N$  training inputs

$$\hat{\mu} = \sum_{i=1}^N \beta_i f_i$$

Like doing a “weighted” nearest neighbors using all the  $N$  training inputs as neighbors with weight  $\beta_i$  given to input  $\mathbf{x}_i$  with score  $f_i = f(\mathbf{x}_i)$

Weighted sum of kernel based similarities of  $\mathbf{x}_*$  with the  $N$  training inputs

$$\hat{\mu} = \sum_{i=1}^N \alpha_i \kappa(\mathbf{x}_i, \mathbf{x}_*)$$

Methods like kernel regression or kernel SVM have their predictions in this form

- Advantage: GP also gives the score's variance  $\hat{\sigma}^2 = \kappa(x_*, x_*) - \mathbf{k}_*^\top \mathbf{K}^{-1} \mathbf{k}_*$
- Thus GP can be viewed as a probabilistic/Bayesian version of kernel methods



# From GP Scores ( $\mathbf{f}$ ) to Actual Outputs ( $\mathbf{y}$ )

- Assume a supervised learning problem with  $N$  training examples  $(\mathbf{X}, \mathbf{y}) = \{(\mathbf{x}_i, y_i)\}_{i=1}^N$
- Denoting  $f_i = f(\mathbf{x}_i)$ , for regression with added noise  $\mathcal{N}(0, \beta^{-1})$

This GP score  $f_i = f(\mathbf{x}_i)$  is the mean of this Gaussian likelihood

$$y_i = f_i + \epsilon_i \quad \longrightarrow \quad p(y_i | f_i) = \mathcal{N}(y_i | f_i, \beta^{-1})$$

The **likelihood** function for all the training outputs (assuming i.i.d.)

$$p(\mathbf{y} | \mathbf{f}) = \mathcal{N}(\mathbf{y} | \mathbf{f}, \beta^{-1} \mathbf{I}_N)$$

For multi-class case with  $\mathcal{C}$  classes, we will use a multinoulli with probability vector  $\text{softmax}(f_i)$  where  $f_i$  will be a  $\mathcal{C}$ -dim vector of logits

- Likewise, for binary classification, **likelihood**  $p(\mathbf{y} | \mathbf{f}) = \prod_{i=1}^N \text{Bernoulli}(y_i | \sigma(f_i))$
- In general, when using GP, the PPD  $p(y_* | \mathbf{y})$  of output  $y_*$  for a new test input  $\mathbf{x}_*$

Averaging over  $f$

$$p(y_* | \mathbf{y}) = \int p(y_* | f_*) p(f_*, \mathbf{f} | \mathbf{y}) d\mathbf{f} df_*$$

The GP posterior

Likelihood function for training output

$$p(\mathbf{f} | \mathbf{y}) \propto p(\mathbf{f}) p(\mathbf{y} | \mathbf{f})$$

Skipping the training and test inputs from the PPD notation

Likelihood

GP predictive: Always a Gaussian

The GP posterior

$$= \int p(y_* | f_*) p(f_* | \mathbf{f}) p(\mathbf{f} | \mathbf{y}) d\mathbf{f} df_*$$

The GP prior: Gaussian  $p(\mathbf{f}) = \mathcal{N}(\mathbf{f} | \boldsymbol{\mu}, \mathbf{K})$  or  $p(\mathbf{f}) = \mathcal{N}(\mathbf{f} | \mathbf{0}, \mathbf{K})$  if mean function is zero

# GP Prediction with Gaussian Likelihood

- In general, the PPD when using GP is defined as

$$p(y_*|\mathbf{y}) = \int p(y_*|f_*)p(f_*|\mathbf{f})p(\mathbf{f}|\mathbf{y})d\mathbf{f}df_*$$

And don't even have to compute/use the posterior  $p(\mathbf{f}|\mathbf{y})$  (which in this case is a Gaussian by the way ☺) to get the PPD

- For Gaussian likelihood (and fixed hyperparams), we don't need to do above integral
- Reason: The marginal likelihood is Gaussian

Gaussian likelihood  
(assuming  $\beta$  is fixed)

$$p(\mathbf{y}|\mathbf{f}) = \mathcal{N}(\mathbf{y}|\mathbf{f}, \beta^{-1}\mathbf{I}_N)$$

GP prior

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

Assuming zero  
mean function

Marginal likelihood  
of training outputs

$$p(\mathbf{y}) = \int p(\mathbf{y}|\mathbf{f})p(\mathbf{f})d\mathbf{f} = \mathcal{N}(\mathbf{y}|\mathbf{0}, \mathbf{K} + \beta^{-1}\mathbf{I}_N) = \mathcal{N}(\mathbf{y}|\mathbf{0}, \mathbf{C}_N)$$

Marginal likelihood of  
training and test outputs

$$p\left(\begin{bmatrix} \mathbf{y} \\ y_* \end{bmatrix}\right) = \mathcal{N}\left(\begin{bmatrix} \mathbf{y} \\ y_* \end{bmatrix} \middle| \begin{bmatrix} \mathbf{0} \\ 0 \end{bmatrix}, \begin{bmatrix} \mathbf{C}_N & \mathbf{k}_* \\ \mathbf{k}_*^\top & \kappa(x_*, x_*) + \beta^{-1} \end{bmatrix}\right)$$

PPD obtained using  
joint to conditional  
results of Gaussians

$$p(y_*|\mathbf{y}) = \mathcal{N}(y_*|\mathbf{k}_*^\top \mathbf{C}_N^{-1} \mathbf{y}, \kappa(x_*, x_*) - \mathbf{k}_*^\top \mathbf{C}_N^{-1} \mathbf{k}_* + \beta^{-1})$$

- $p(y_*|\mathbf{y})$  is almost identical to  $p(f_*|\mathbf{f})$  with  $\mathbf{K}$  replaced by  $\mathbf{C}_N + \text{extra } \beta^{-1} \text{ noise variance}$



# Learning Hyperparameters in GP based Models

- Can learn the hyperparameters of the GP prior as well as of the likelihood model
- Assuming  $\mu = \mathbf{0}$ , the hyperparams of GP are cov/kernel function hyperparams

$$\kappa(\mathbf{x}_n, \mathbf{x}_m) = \exp\left(-\frac{\|\mathbf{x}_n - \mathbf{x}_m\|^2}{\gamma}\right) \quad \text{(RBF kernel)}$$

$$\kappa(\mathbf{x}_n, \mathbf{x}_m) = \exp\left(-\sum_{d=1}^D \frac{(\mathbf{x}_{nd} - \mathbf{x}_{md})^2}{\gamma_d}\right) \quad \text{(ARD kernel)}$$

$$\kappa(\mathbf{x}_n, \mathbf{x}_m) = \kappa_{\theta_1}(\mathbf{x}_n, \mathbf{x}_m) + \kappa_{\theta_2}(\mathbf{x}_n, \mathbf{x}_m) + \dots + \kappa_{\theta_M}(\mathbf{x}_n, \mathbf{x}_m) \quad \text{(flexible composition of multiple kernels)}$$

Can help in feature selection (irrelevant features will tend to have very large  $\gamma_d$ )

Different RBF kernel bandwidth  $\gamma_d$  for each feature

Ability to learn kernel hyperparams (without cross-valid) is another very appealing property of GP



- MLE-II is a popular choice for learning these hyperparams (otherwise MCMC, VI, etc)
- Denoting the covariance/kernel matrix as  $\mathbf{K}_\theta$ , for Gaussian likelihood case, the marg-lik

$$p(\mathbf{y}|\theta, \beta^{-1}) = \mathcal{N}(\mathbf{y}|\mathbf{0}, \mathbf{K}_\theta + \beta^{-1}\mathbf{I}_N)$$

- This can be maximized to learn  $\theta$  and  $\beta$
- For **non-Gaussian likelihoods**, the marg-lik itself will need to be approximated





# Weight Space View vs Function Space View

- GPs are defined w.r.t. a **function space** that models input-output relationship
- In contrast, we have seen models that are defined w.r.t. a **weight space**, e.g.,

$$p(\mathbf{y}|\mathbf{X}, \mathbf{w}) = \mathcal{N}(\mathbf{y}|\mathbf{X}\mathbf{w}, \beta^{-1}\mathbf{I}_N)$$

Likelihood

$$p(\mathbf{w}) = \mathcal{N}(\mathbf{w}|\boldsymbol{\mu}_0, \boldsymbol{\Sigma}_0)$$

Prior over weight vector

$$p(\mathbf{y}|\mathbf{X}) = \int p(\mathbf{y}|\mathbf{X}, \mathbf{w})p(\mathbf{w})d\mathbf{w} = \mathcal{N}(\mathbf{y}|\mathbf{X}\boldsymbol{\mu}_0, \beta^{-1}\mathbf{I}_N + \mathbf{X}\boldsymbol{\Sigma}_0\mathbf{X}^\top)$$

Marginal likelihood  
after integrating  
out the weights

$$p(\mathbf{y}|\mathbf{X}) = \mathcal{N}(\mathbf{y}|\mathbf{0}, \beta^{-1}\mathbf{I}_N + \mathbf{X}\mathbf{X}^\top)$$

Marginal likelihood assuming  $\boldsymbol{\mu}_0 = \mathbf{0}$  and  $\boldsymbol{\Sigma}_0 = \mathbf{I}$ 

$$p(\mathbf{y}|\mathbf{X}) = \mathcal{N}(\mathbf{y}|\mathbf{0}, \mathbf{X}\mathbf{X}^\top)$$

Assuming noise-free likelihood

- Thus the **joint marginal** of the  $N$  responses  $\mathbf{y}_1, \mathbf{y}_2, \dots, \mathbf{y}_N$  is a **multivariate Gaussian**

This equivalence also shows  
that Bayesian linear regression  
is a special case of GP with  
linear kernel

$$p\left(\begin{bmatrix} y_1 \\ y_2 \\ \vdots \\ y_N \end{bmatrix}\right) = \mathcal{N}\left(\begin{bmatrix} 0 \\ 0 \\ \vdots \\ 0 \end{bmatrix}, \begin{bmatrix} x_1^\top x_1 & \dots & x_1^\top x_N \\ x_2^\top x_1 & \dots & x_2^\top x_N \\ \vdots & \ddots & \vdots \\ x_N^\top x_1 & \dots & x_N^\top x_N \end{bmatrix}\right)$$

Same as a GP  $f(\mathbf{x}_i) = y_i$ ,  $\mu(\mathbf{x}) = \mathbf{0}$  and linear  
covariance/kernel function  $\kappa(\mathbf{x}_i, \mathbf{x}_j) = \mathbf{x}_i^\top \mathbf{x}_j$

- Thus GPs can be seen as bypassing the weight space and directly defining the model using a marginal likelihood via a function space defined by the GP



# Scalability of GPs

- Computational costs in some steps of GP models scale in the size of training data
- For example, prediction cost is  $O(N)$

$O(N)$  cost assuming  $\mathbf{C}_N$  is already inverted

$$p(y_* | \mathbf{y}) = \mathcal{N}(y_* | \hat{\mu}, \hat{\sigma}^2) \quad \hat{\mu} = \mathbf{k}_*^\top \mathbf{C}_N^{-1} \mathbf{y} \quad \hat{\sigma}^2 = \kappa(x_*, x_*) - \mathbf{k}_*^\top \mathbf{C}_N^{-1} \mathbf{k}_* + \beta^{-1}$$

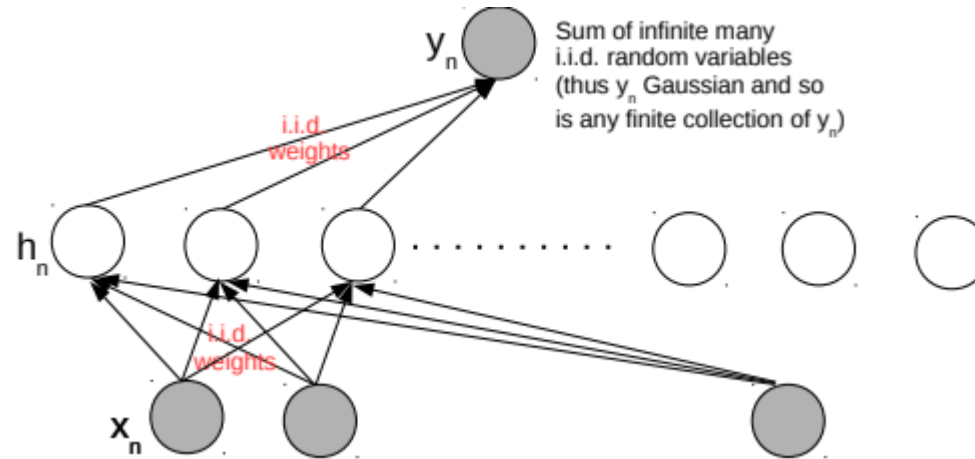
- GP models often require matrix inversions (e.g., in marg-lik computation when estimating hyperparameters) – takes  $O(N^3)$
- Storage also requires  $O(N^2)$  since need to store the covariance matrix
- A lot of work on speeding up GPs<sup>1</sup>. Some prominent approaches include
  - **Inducing Point Methods** (condition predictions only on a small set of “learnable” points)
    - Divide-and-Conquer (learn GP on small subsets of data and aggregate predictions)
    - Kernel approximations
- Note that nearest neighbor methods and kernel methods also face similar issues
  - Many tricks to speed up kernel methods can be used for speeding up GPs too

$M \ll N$  pseudo-inputs and pseudo-outputs



# Neural Networks and Gaussian Process

- An infinitely-wide single hidden layer NN with i.i.d. priors on weights = GP
- Shown formally by (Neal<sup>2</sup>, 1994). Based on applying the central limit theorem



- This equivalence is useful for several reasons
  - Can use a GP instead of an **infinitely wide** Bayesian NN (which is impractical anyway)
  - With GPs, inference is easy (at least for regression and with known hyperparams)
  - A proof that GPs can also learn any function (just like infinitely wide neural nets - Hornik's theorem)
- Connection generalized to infinitely wide multiple hidden layer NN (Lee et al<sup>3</sup>, 2018)



<sup>2</sup>Priors for infinite networks, Tech Report, 1994

<sup>3</sup>Deep Neural Networks as Gaussian Processes (ICLR 2018)

# GP: Some other comments

- GPs can be thought of as Bayesian analogues of kernel methods
- Can get estimate of the uncertainty in the function and its predictions
- Can learn the kernel (by learning the hyperparameters of the kernels)
- In some ways, GPs and (Bayesian/ensembles of) deep neural nets have same goals
  - These methods are also very related (though appear different based on their formulation)
  - Several recent papers have investigated these connections
- GP can be a nice alternative to (Bayesian/ensembles of) deep neural networks
  - GP may be preferable if we don't have that much training data (deep networks requires lots of data to train well)
  - When we have lots of training data, training and test speed may be an issue for GP (but faster versions exist)
- Not limited to supervised learning problems
  - $f$  could even define a mapping of low-dim latent variable  $\mathbf{z}_n$  to an observation  $\mathbf{x}_n$

$$\mathbf{x}_n = f(\mathbf{z}_n) + \text{"noise"}$$

GP latent variable model for dimensionality reduction  
(like a kernel version of probabilistic PCA)



# GP: A Visualization

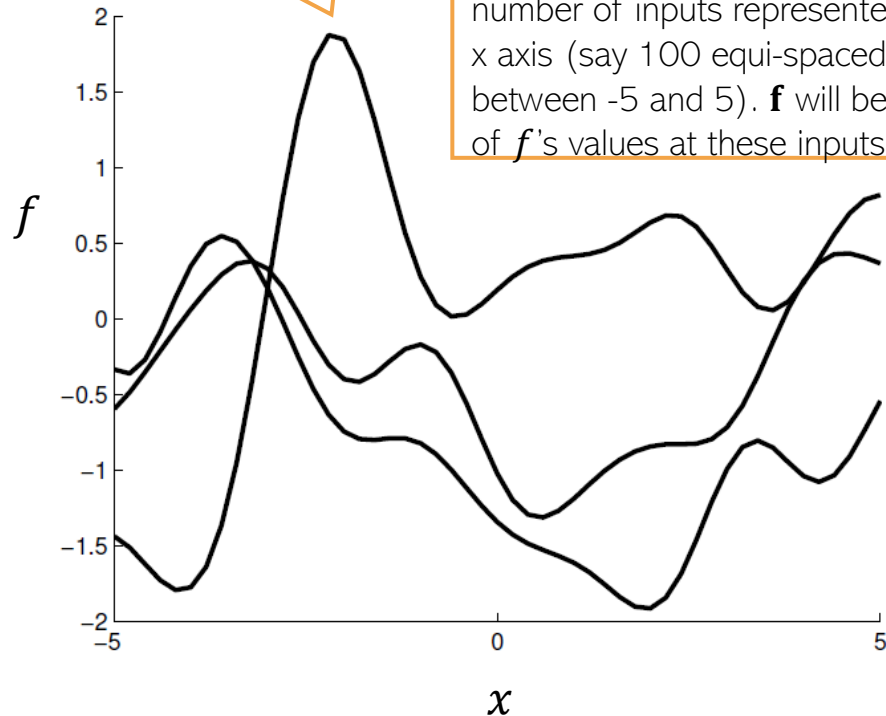
- Assumed zero mean function and a squared exponential kernel

$$k_{\text{SE}}(x, x') = \sigma^2 \exp\left(-\frac{(x-x')^2}{2\ell^2}\right)$$

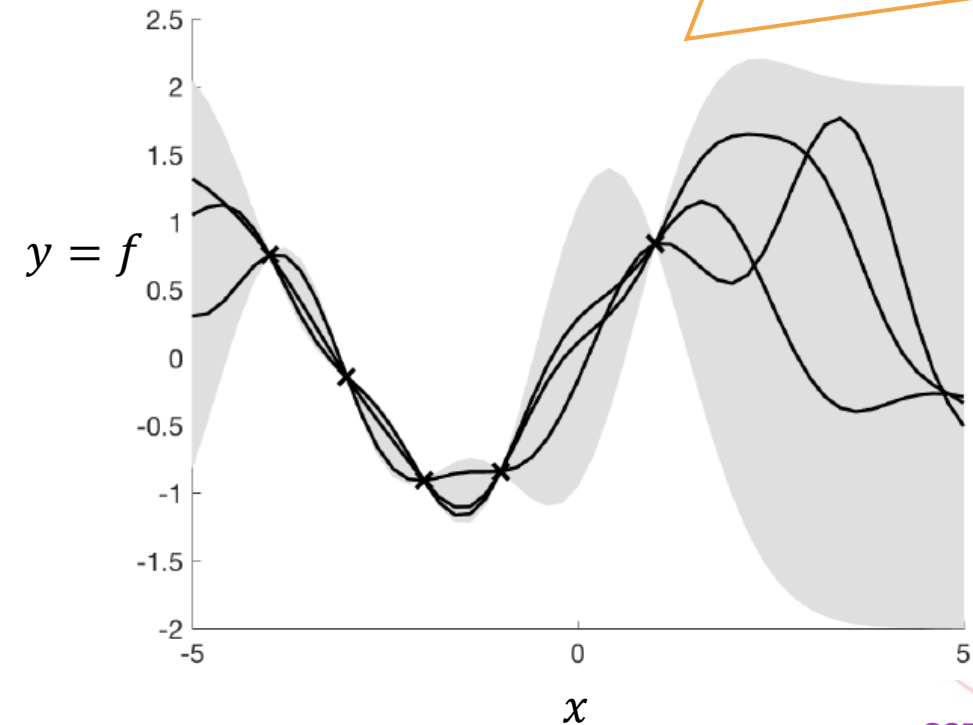
Each curve below is obtained by drawing a random  $\mathbf{f}$  from the GP prior  $p(\mathbf{f}) = \mathcal{N}(\mathbf{0}, \mathbf{K})$  and plotting it.

Shaded area shows the predictive uncertainty for each of the test inputs (+/- 2 std)

$\mathbf{K}$  is the kernel matrix of a finite number of inputs represented on the x axis (say 100 equi-spaced points between -5 and 5).  $\mathbf{f}$  will be a vector of  $f$ 's values at these inputs



Each curve below is obtained by drawing random  $\mathbf{f}$ 's from the GP posterior  $p(\mathbf{f}|\mathbf{y})$  which is also a Gaussian (The + symbols denote the training data and we assume noiseless outputs, i.e.,  $y_i = f_i$ ).





# GP packages

- Many mature implementations of GP exist. You may check out
  - GPyTorch (PyTorch), GPFlow (Tensorflow)
  - sklearn (Python with some basic GP implementations)
  - GPML (MATLAB), GPsuff (MATLAB/Octave)
  - Many others such as Stan, GPlax
- A comparison of the various packages:  
[https://en.wikipedia.org/wiki/Comparison\\_of\\_Gaussian\\_process\\_software](https://en.wikipedia.org/wiki/Comparison_of_Gaussian_process_software)

