

# Avoiding Pathologies in Very Deep Networks



David Duvenaud, Oren Rippel, Ryan Adams, Zoubin Ghahramani

February 18, 2014

# MOTIVATION 1: DEEP LEARNING IS COOL

---

- ▶ Don't you wish we were doing it?
- ▶ Zoubin (2011) “Do you guys ever wonder if this lab focuses too much on Gaussian processes? Like maybe we're going to miss the next big thing, like maybe, say, deep learning”
- ▶ But - GPs are just neural nets, we can make them deep!

## MOTIVATION 2: REGULARIZING NETS

---

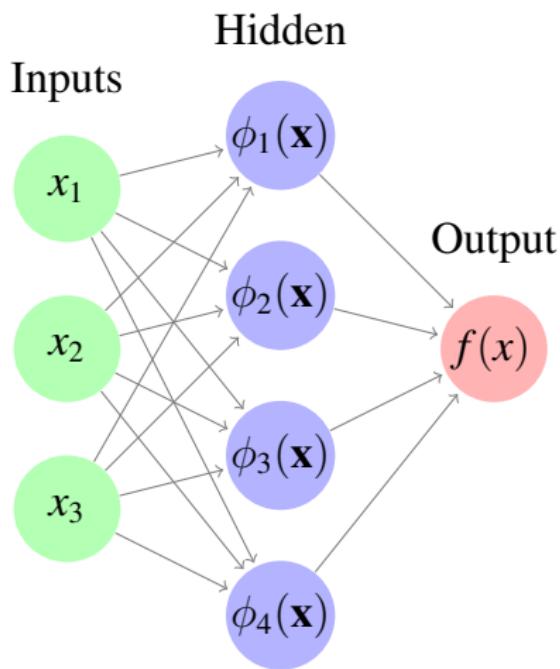
- ▶ Neural nets are getting larger
- ▶ How to regularize billions of parameters?
- ▶ Closely related to constructing priors
- ▶ Priors are easy to analyze - just sample from the prior and look and what sorts of things you get!
- ▶ Can we suggest new regularization schemes or network architectures?

# OUTLINE

---

- ▶ Relation between GPs and neural nets
- ▶ Two ways to deepness:
  - ▶ Deep kernels
  - ▶ Deep GPs
- ▶ What kind of prior on functions do we want?
  - ▶ problems with lots of independent layers
  - ▶ a simple fix
- ▶ Dropout for GPs
  - ▶ Dropping out features
  - ▶ Dropping out inputs

# GPS AS NEURAL NETS



A weighted sum of features,

$$f(\mathbf{x}) = \frac{1}{K} \sum_{i=1}^K \alpha_i \phi_i(\mathbf{x})$$

with any weight distribution,

$$\mathbb{E} [\alpha_i] = 0, \quad \mathbb{V} [\alpha_i] = \sigma^2, \quad i.i.d.$$

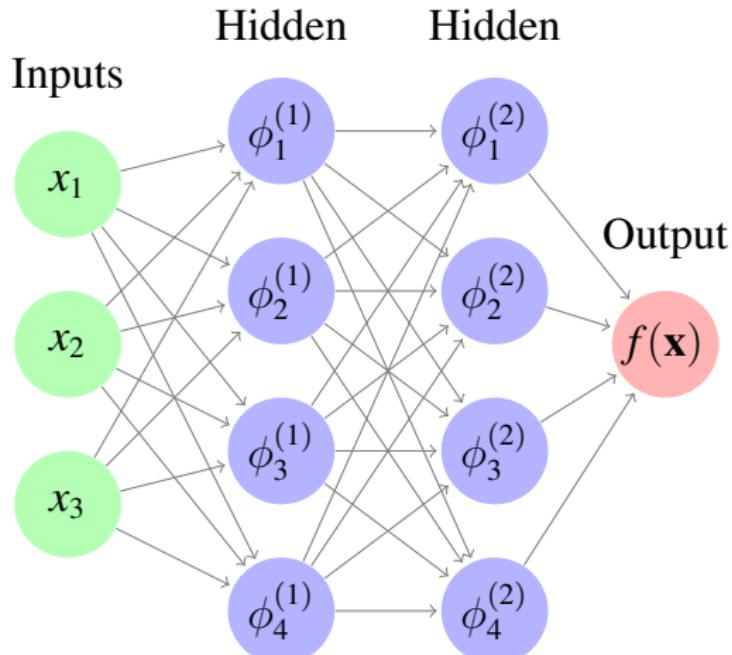
by CLT, gives a GP as  $K \rightarrow \infty$ !

$$\text{cov} \begin{bmatrix} f(\mathbf{x}) \\ f(\mathbf{x}') \end{bmatrix} \rightarrow \frac{\sigma^2}{K} \sum_{i=1}^K \phi_i(\mathbf{x}) \phi_i(\mathbf{x}')$$

# KERNEL LEARNING AS FEATURE LEARNING

- ▶ GPs have fixed features, integrate out feature weights.
- ▶ Mapping between kernels and features:  
 $k(\mathbf{x}, \mathbf{x}') = \Phi(\mathbf{x})^T \Phi(\mathbf{x}').$
- ▶ Any PSD kernel can be written as inner product of features. (Mercer's Theorem)
- ▶ Kernel learning = feature learning
  
- ▶ What if we make the GP neural network deep?

# DEEP NETS, DEEP KERNELS



Now our model is:

$$\begin{aligned}f(\mathbf{x}) &= \frac{1}{K} \sum_{i=1}^K \alpha_i \phi_i (\Phi^{(1)}(\mathbf{x})) \\&= \boldsymbol{\alpha}^\top \Phi^{(2)} (\Phi^{(1)}(\mathbf{x}))\end{aligned}$$

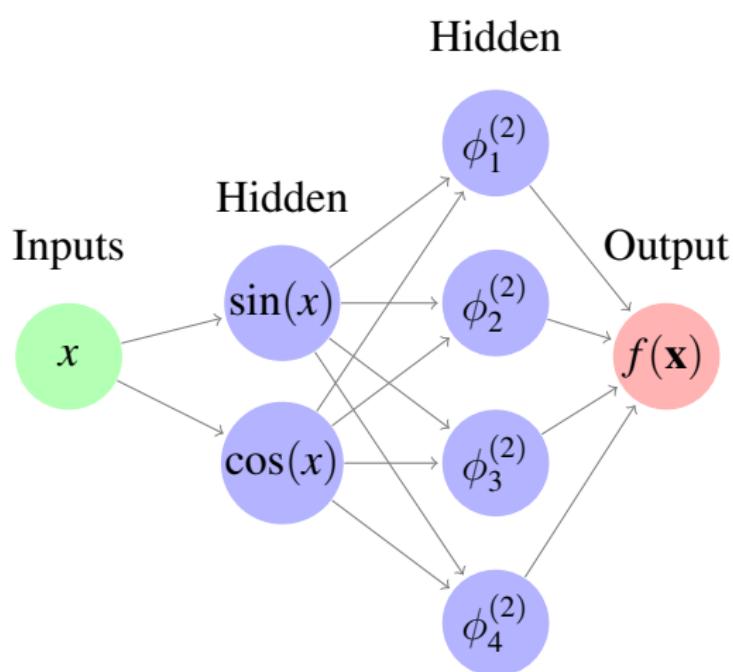
Instead of

$$k_1(\mathbf{x}, \mathbf{x}') = \Phi^{(1)}(\mathbf{x})^\top \Phi^{(1)}(\mathbf{x}'),$$

we have “deep kernel”:

$$\begin{aligned}k_2(\mathbf{x}, \mathbf{x}') &\\&= \Phi^{(2)}(\Phi^{(1)}(\mathbf{x}))^\top \Phi^{(2)}(\Phi^{(1)}(\mathbf{x}'))\end{aligned}$$

# EXAMPLE DEEP KERNEL: PERIODIC



Now our model is:

$$\Phi^1(\mathbf{x}) = [\sin(\mathbf{x}), \cos(\mathbf{x})]$$

we have “deep kernel”:

$$k_2(\mathbf{x}, \mathbf{x}')$$

$$= \exp\left(-\frac{1}{2} (\Phi^1(\mathbf{x}) - \Phi^1(\mathbf{x}'))\right)$$

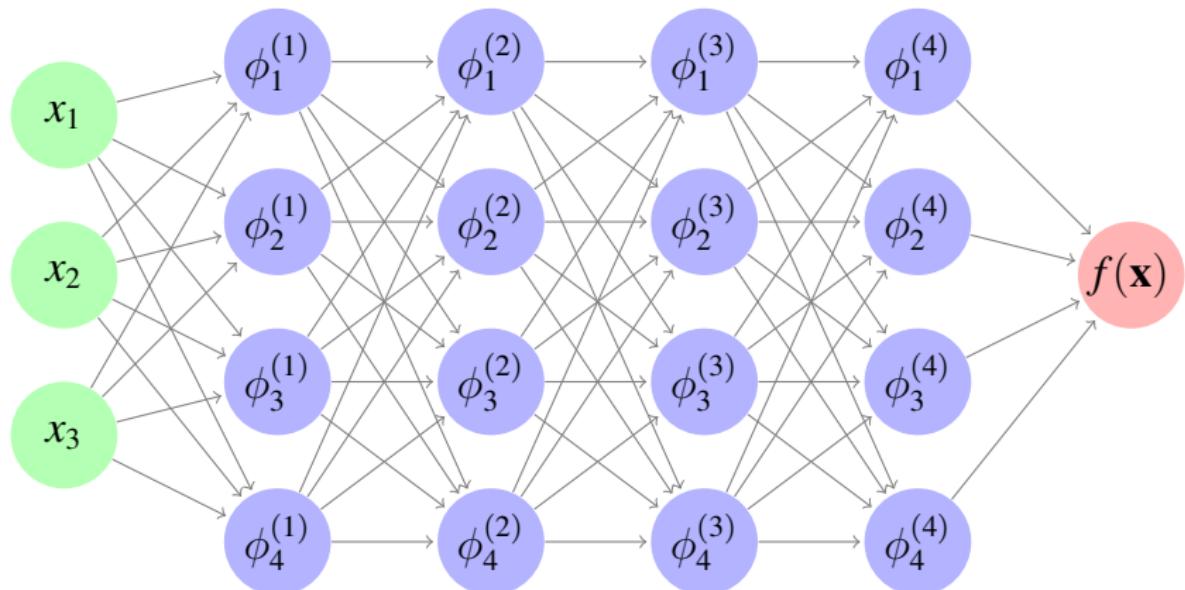
# DEEP KERNELS

- ▶ (Cho, 2012) built kernels by composing feature mappings.
- ▶ Composing any kernel  $k_1$  with a squared-exp kernel (SE):

$$\begin{aligned} k_2(\mathbf{x}, \mathbf{x}') &= \\ &= (\Phi^{SE}(\Phi^1(\mathbf{x})))^\top \Phi^{SE}(\Phi^1(\mathbf{x}')) \\ &= \exp\left(-\frac{1}{2}\|\Phi^1(\mathbf{x}) - \Phi^1(\mathbf{x}')\|_2^2\right) \\ &= \exp\left(-\frac{1}{2} [\Phi^1(\mathbf{x})^\top \Phi^1(\mathbf{x}) - 2\Phi^1(\mathbf{x})^\top \Phi^1(\mathbf{x}') + \Phi^1(\mathbf{x}')^\top \Phi^1(\mathbf{x}')]\right) \\ &= \exp\left(-\frac{1}{2} [k_1(\mathbf{x}, \mathbf{x}) - 2k_1(\mathbf{x}, \mathbf{x}') + k_1(\mathbf{x}', \mathbf{x}')]\right) \end{aligned}$$

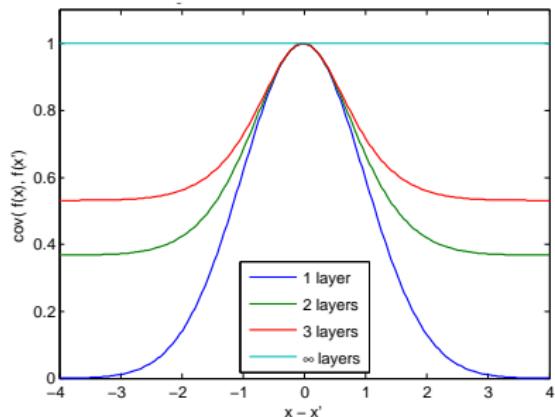
- ▶ A closed form... let's do it again!

# WE NEED TO GO DEEPER

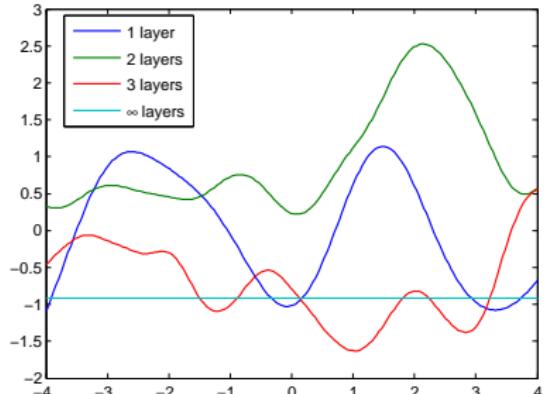


# INFINITELY DEEP KERNELS

- ▶ For SE kernel,  $k_{L+1}(\mathbf{x}, \mathbf{x}') = \exp(k_L(\mathbf{x}, \mathbf{x}') - 1)$ .
- ▶ What is the limit of composing SE features?



Kernel

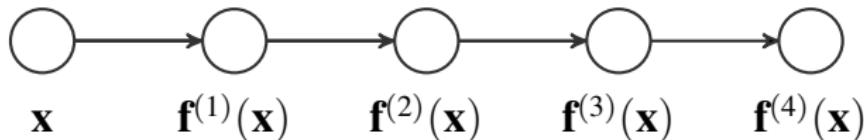


Draws from GP prior

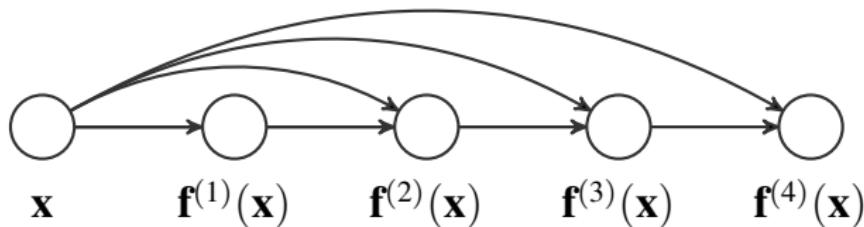
- ▶  $k_\infty(\mathbf{x}, \mathbf{x}') = 1$  everywhere. ☺

# A SIMPLE FIX...

- ▶ Following a suggestion from [Neal \(1995\)](#), we connect the inputs  $\mathbf{x}$  to each layer:



a) standard MLP architecture.



b) Input-connected architecture.

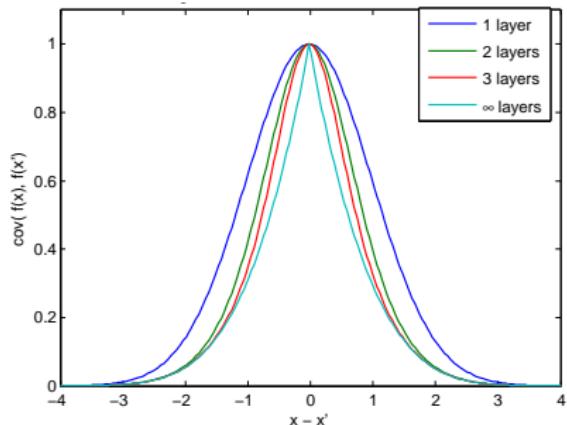
## A SIMPLE FIX...

- ▶ Following a suggestion from [Neal \(1995\)](#), we connect the inputs  $\mathbf{x}$  to each layer:

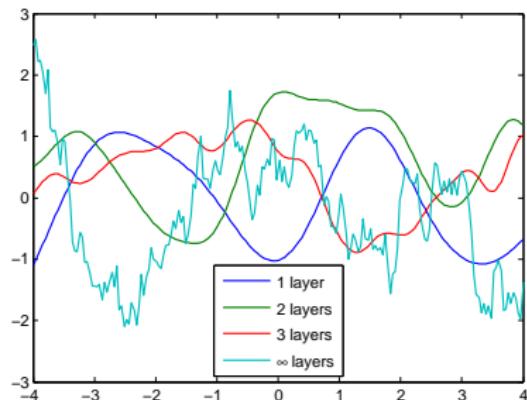
$$\begin{aligned} k_{L+1}(\mathbf{x}, \mathbf{x}') &= \\ &= \exp \left( -\frac{1}{2} \left\| \begin{bmatrix} \Phi^L(\mathbf{x}) \\ \mathbf{x} \end{bmatrix} - \begin{bmatrix} \Phi^L(\mathbf{x}') \\ \mathbf{x}' \end{bmatrix} \right\|_2^2 \right) \\ &= \exp \left( -\frac{1}{2} [k_L(\mathbf{x}, \mathbf{x}) - 2k_L(\mathbf{x}, \mathbf{x}') + k_L(\mathbf{x}', \mathbf{x}')] - \frac{1}{2} \|\mathbf{x} - \mathbf{x}'\|_2^2 \right) \end{aligned}$$

# INFINITELY DEEP KERNELS

- ▶ What is the limit of compositions of input-connected SE features?
- ▶  $k_{L+1}(\mathbf{x}, \mathbf{x}') = \exp \left( k_L(\mathbf{x}, \mathbf{x}') - 1 - \frac{1}{2} \|\mathbf{x} - \mathbf{x}'\|_2^2 \right)$ .



Kernels



Draws from GP priors

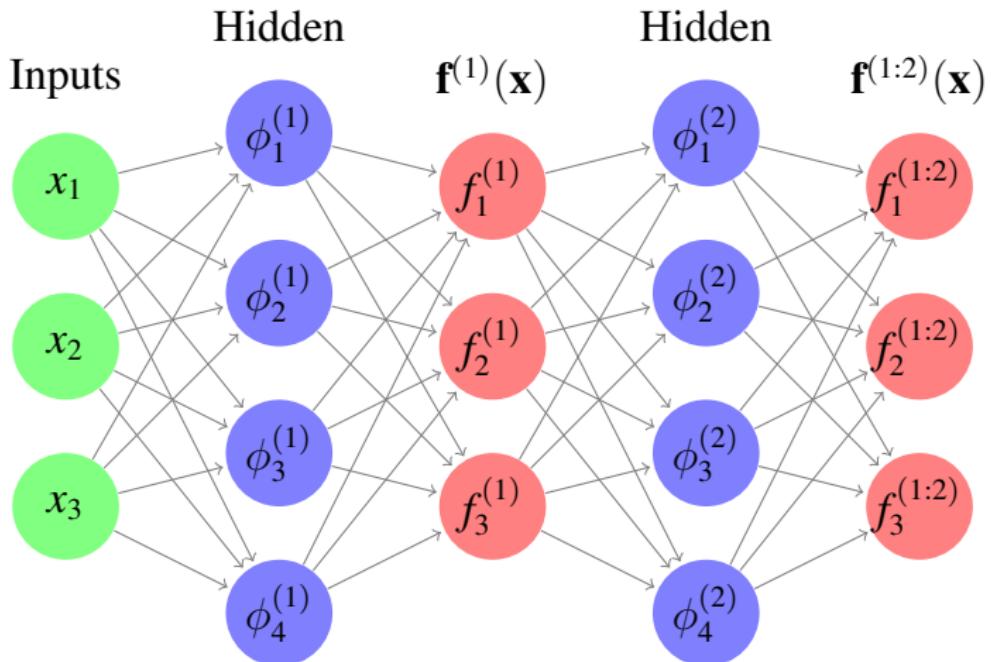
- ▶ Like an Ornstein-Uhlenbeck process with skinny tails
- ▶ Samples are non-differentiable (fractal).

# WHAT WENT WRONG?

---

- ▶ Fixed feature mapping, unlikely to be useful for anything
- ▶ power of neural nets comes from learning a custom representation
- ▶ Need to search over feature mappings!
- ▶ Can try to learn kernels, or even better, integrate over feature mappings

# DEEP GAUSSIAN PROCESSES



# DEEP GAUSSIAN PROCESSES

- ▶ a prior over compositions of functions:

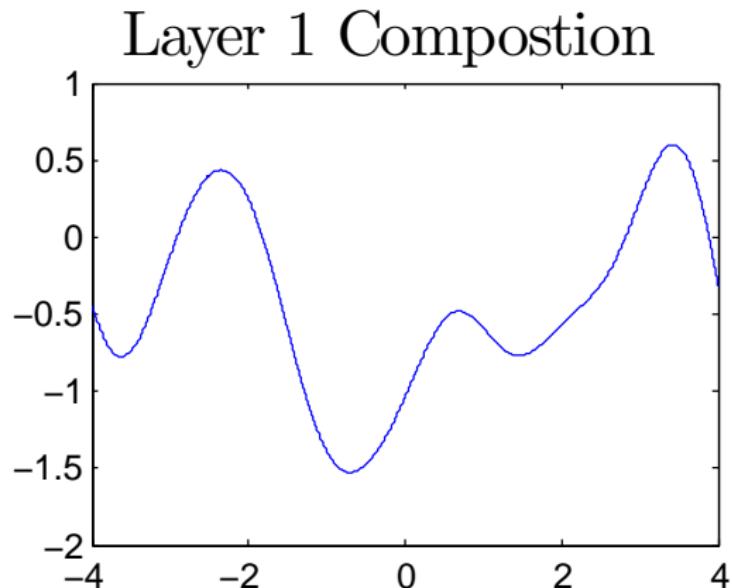
$$\mathbf{f}^{(1:L)}(\mathbf{x}) = \mathbf{f}^{(L)}\left(\mathbf{f}^{(L-1)}\left(\dots \mathbf{f}^{(2)}\left(\mathbf{f}^{(1)}(\mathbf{x})\right)\dots\right)\right) \quad (1)$$

where each  $\mathbf{f}_d^{(\ell)} \stackrel{\text{ind}}{\sim} \mathcal{GP}\left(0, k_d^\ell(\mathbf{x}, \mathbf{x}')\right)$ .

- ▶ Analogous to neural nets, where each neuron's activation function is an independent draw from a GP.
- ▶ inference is really hard.
- ▶ maybe we can learn something just from looking at draws?

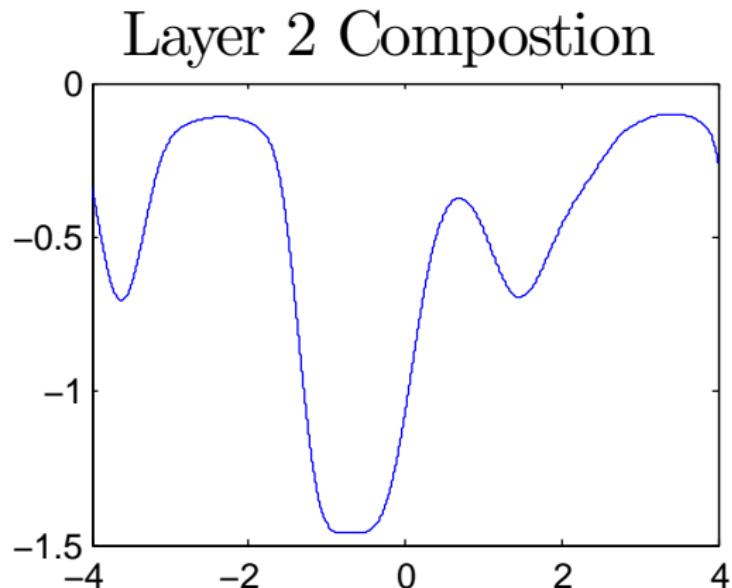
# DEEP GAUSSIAN PROCESSES

- ▶ Draws from one-dimensional deep GPs:



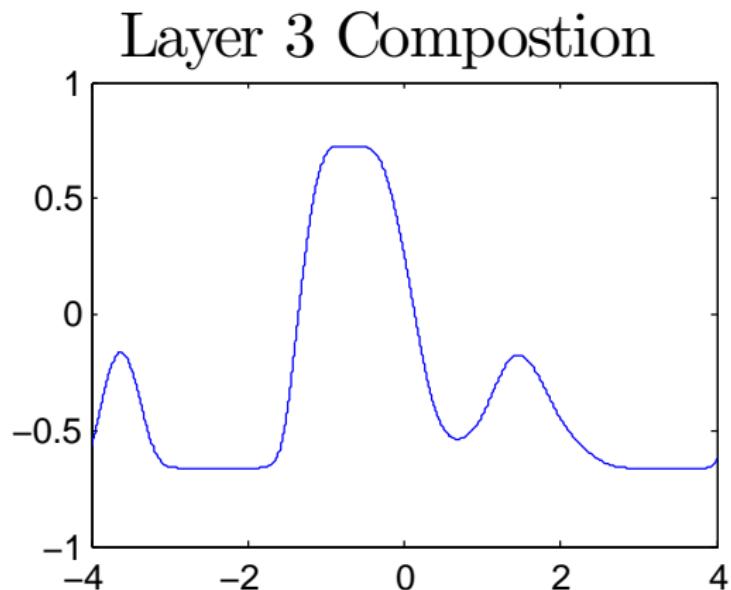
# DEEP GAUSSIAN PROCESSES

- ▶ Draws from one-dimensional deep GPs:



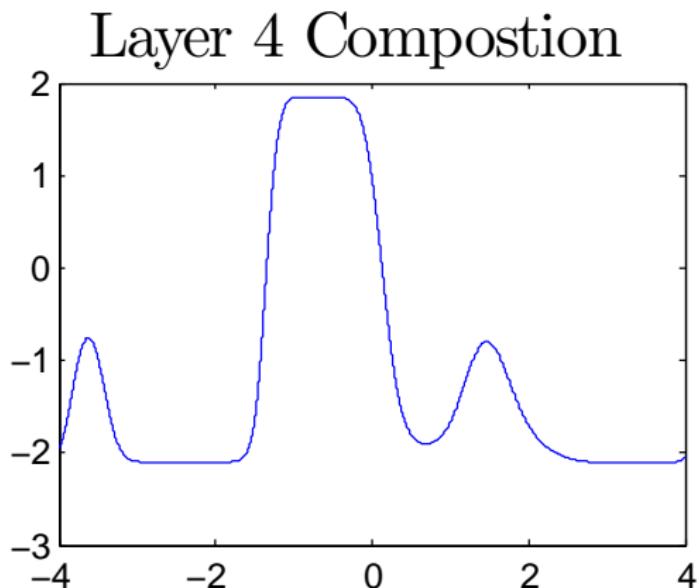
# DEEP GAUSSIAN PROCESSES

- ▶ Draws from one-dimensional deep GPs:



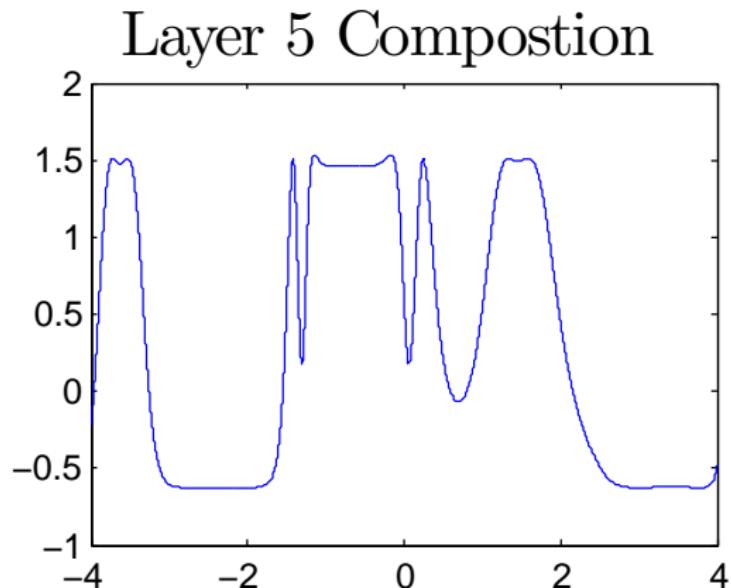
# DEEP GAUSSIAN PROCESSES

- ▶ Draws from one-dimensional deep GPs:



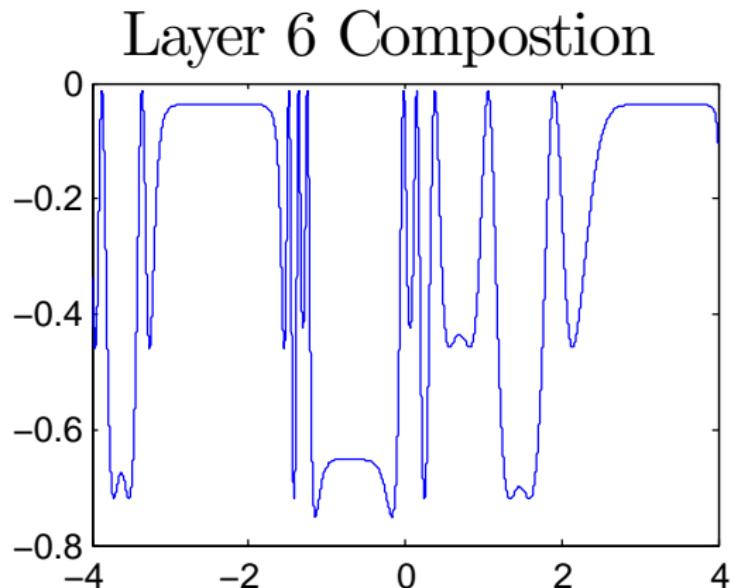
# DEEP GAUSSIAN PROCESSES

- ▶ Draws from one-dimensional deep GPs:



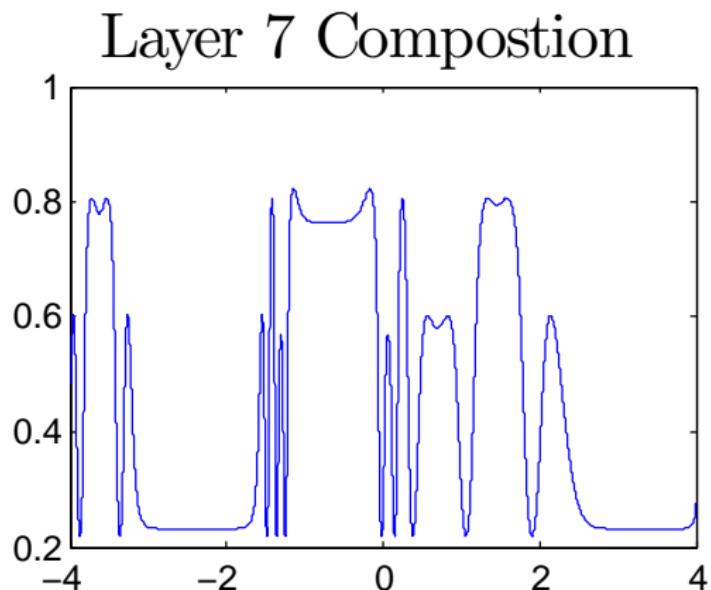
# DEEP GAUSSIAN PROCESSES

- ▶ Draws from one-dimensional deep GPs:



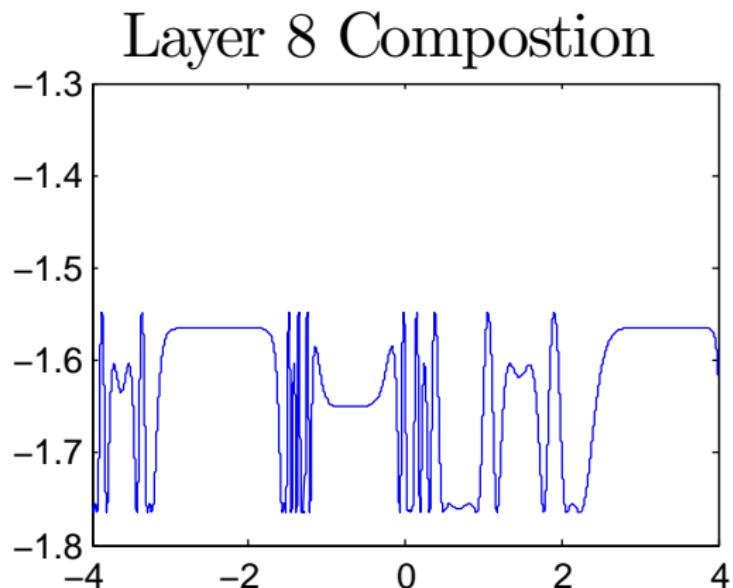
# DEEP GAUSSIAN PROCESSES

- ▶ Draws from one-dimensional deep GPs:



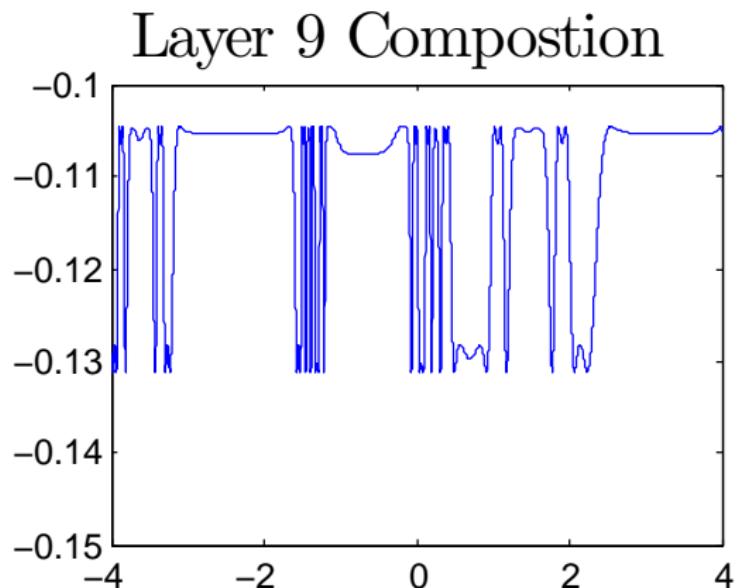
# DEEP GAUSSIAN PROCESSES

- ▶ Draws from one-dimensional deep GPs:



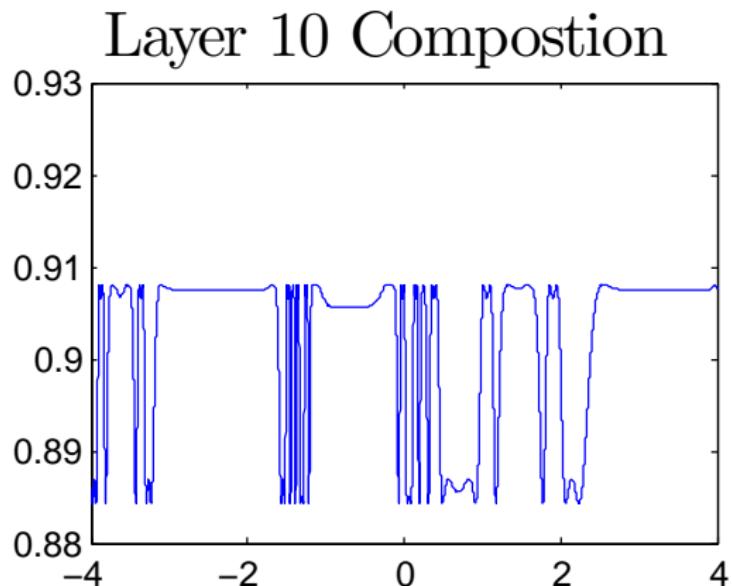
# DEEP GAUSSIAN PROCESSES

- ▶ Draws from one-dimensional deep GPs:



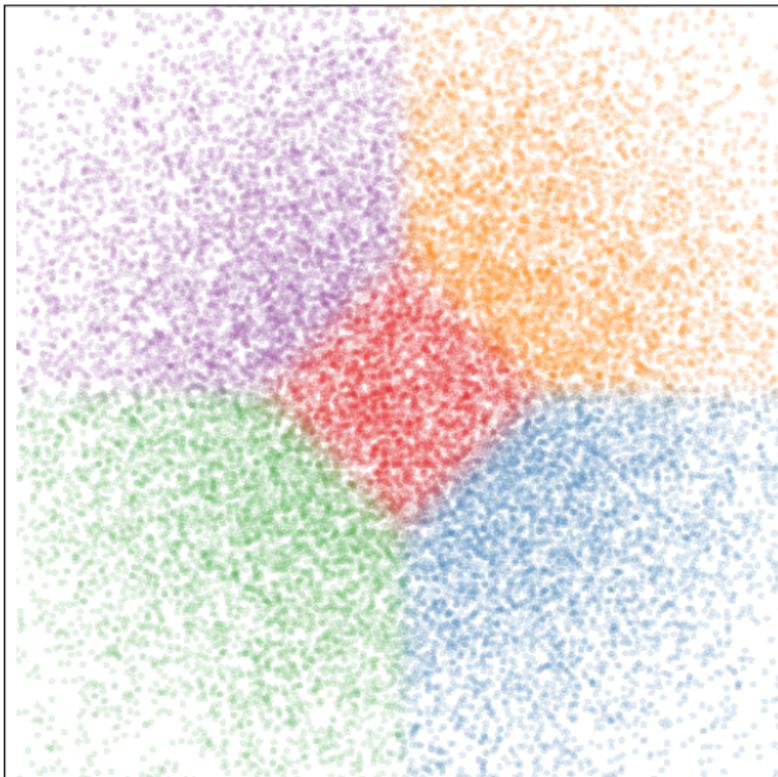
# DEEP GAUSSIAN PROCESSES

- ▶ Draws from one-dimensional deep GPs:



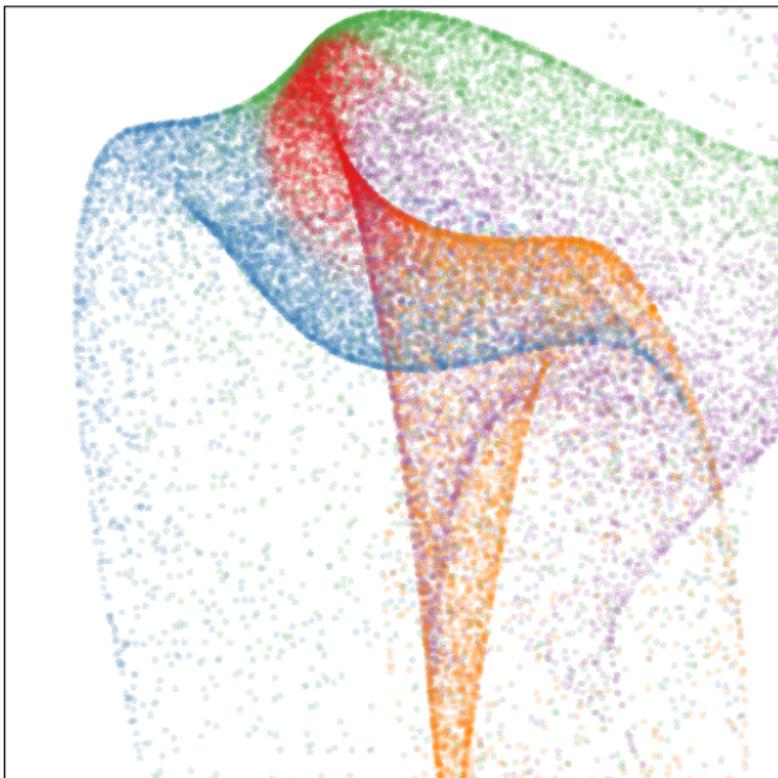
# DEEP GAUSSIAN PROCESSES

- ▶ 2D to 2D warpings of a Gaussian density:



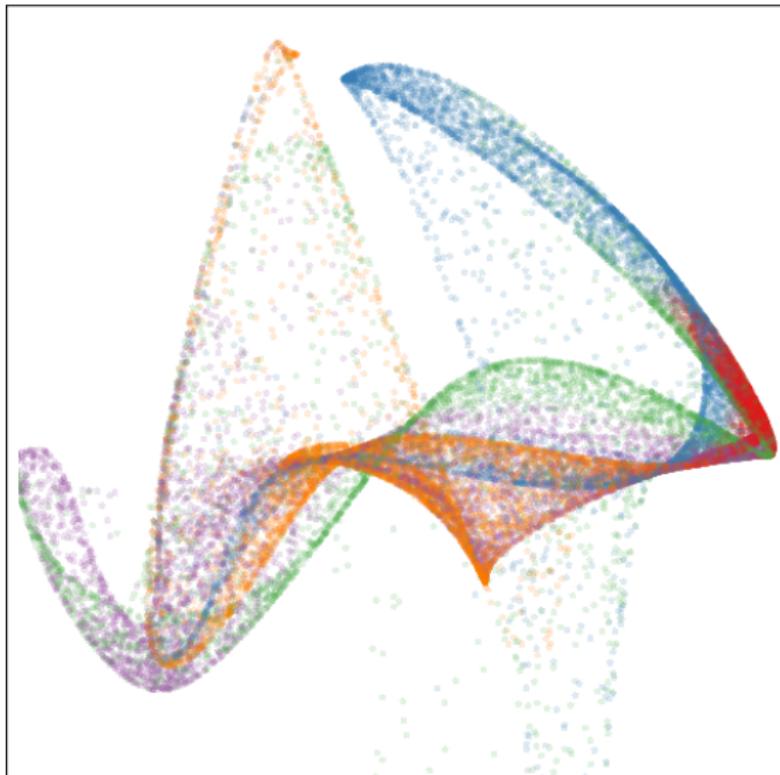
# DEEP GAUSSIAN PROCESSES

- ▶ 2D to 2D warpings of a Gaussian density:



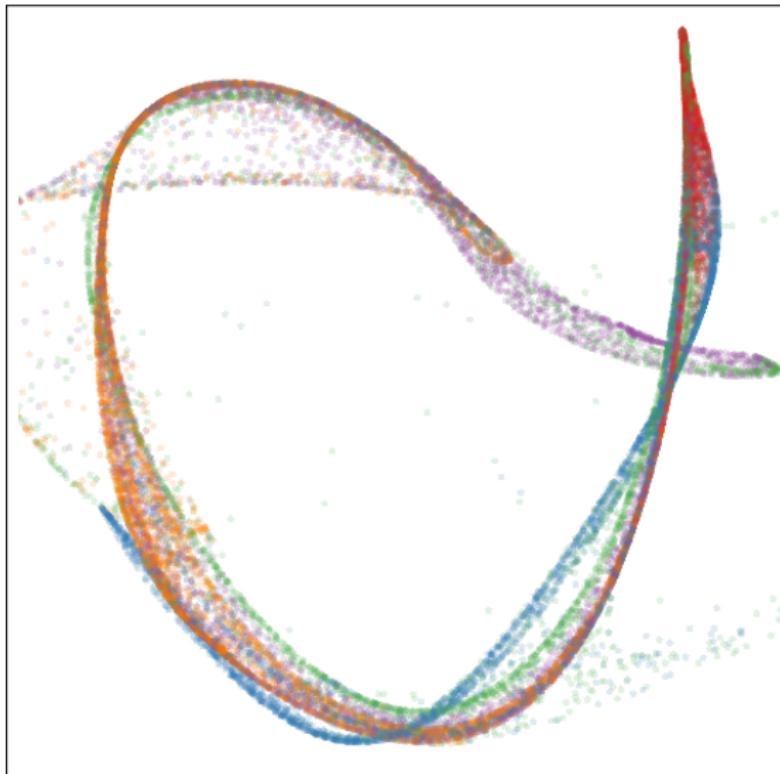
# DEEP GAUSSIAN PROCESSES

- ▶ 2D to 2D warpings of a Gaussian density:



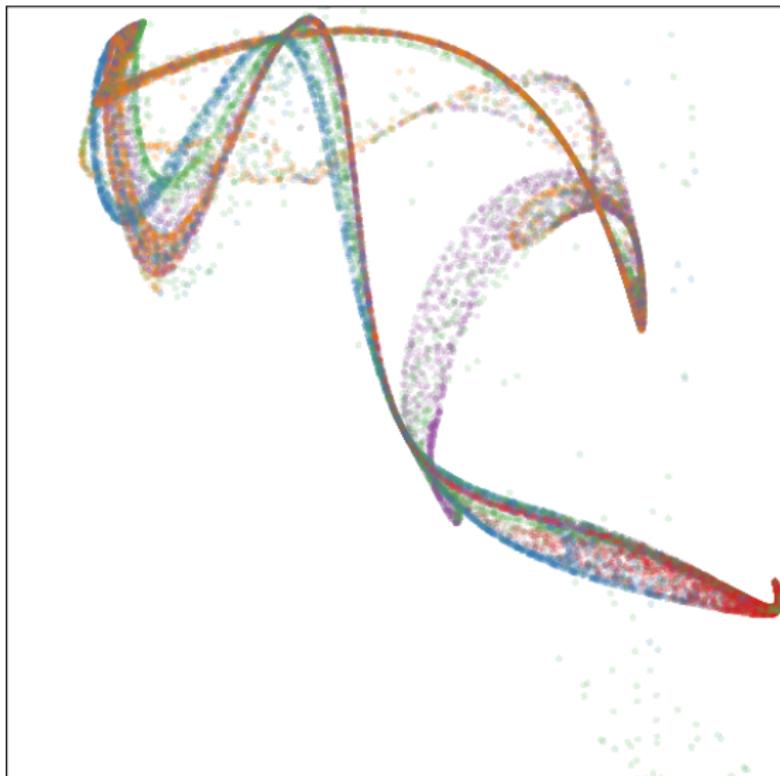
# DEEP GAUSSIAN PROCESSES

- ▶ 2D to 2D warpings of a Gaussian density:



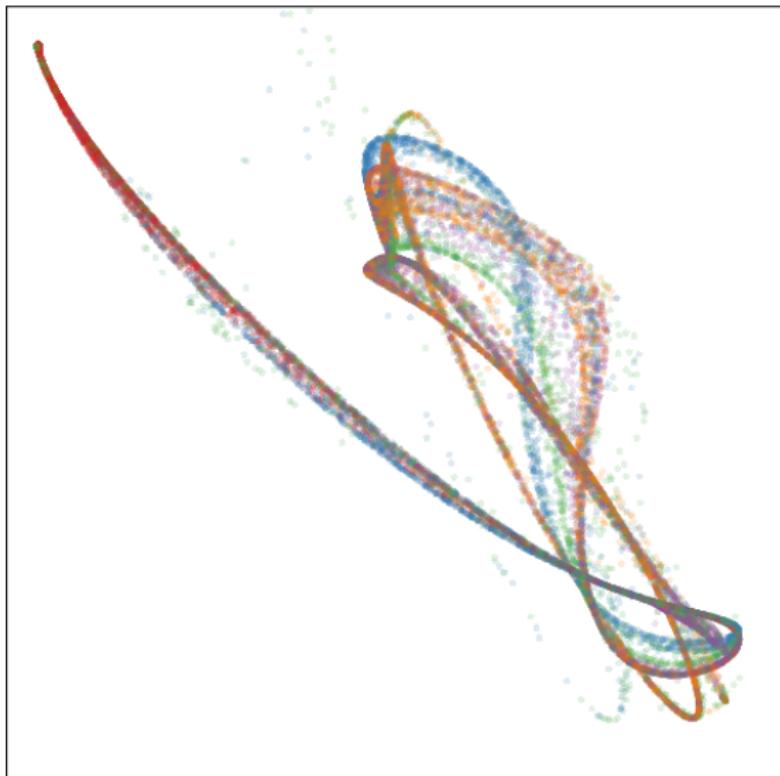
# DEEP GAUSSIAN PROCESSES

- ▶ 2D to 2D warpings of a Gaussian density:



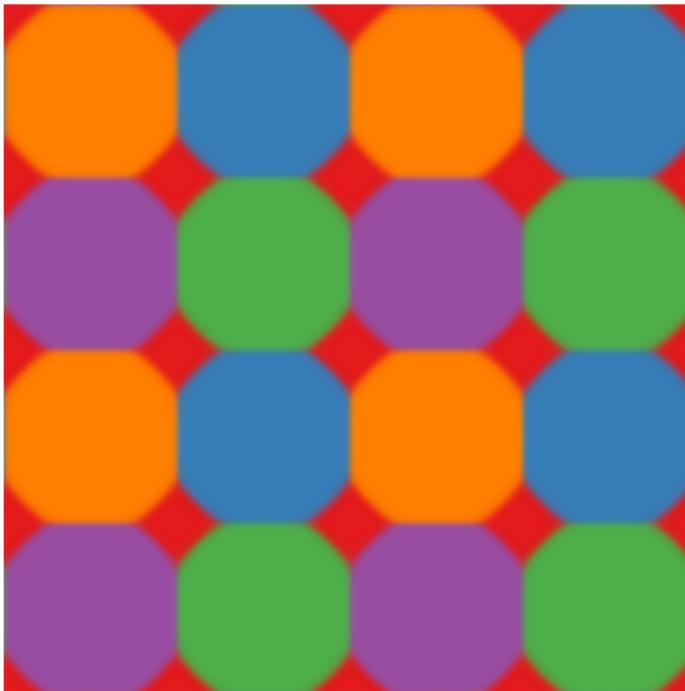
# DEEP GAUSSIAN PROCESSES

- ▶ 2D to 2D warpings of a Gaussian density:



# DEEP GAUSSIAN PROCESSES

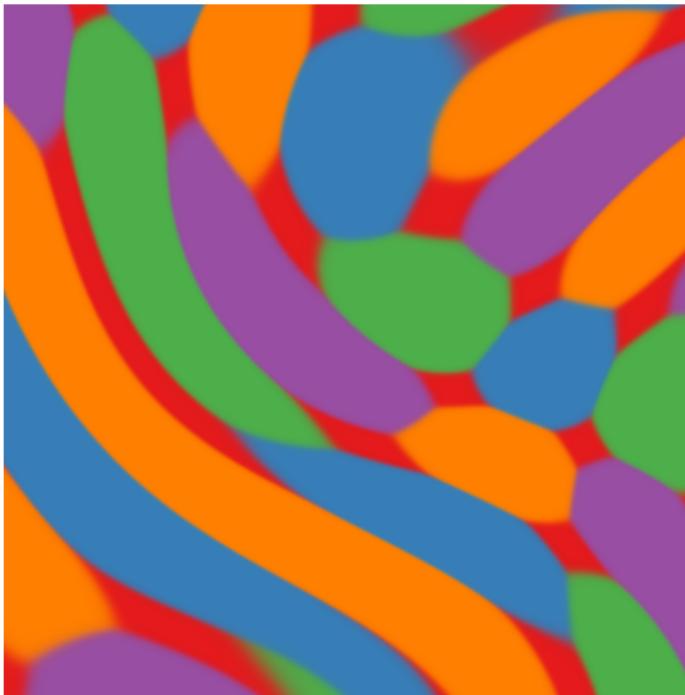
- ▶ Showing the  $x$  that gave rise to a particular  $y$
- ▶ (i.e. decision boundaries)



No warping

# DEEP GAUSSIAN PROCESSES

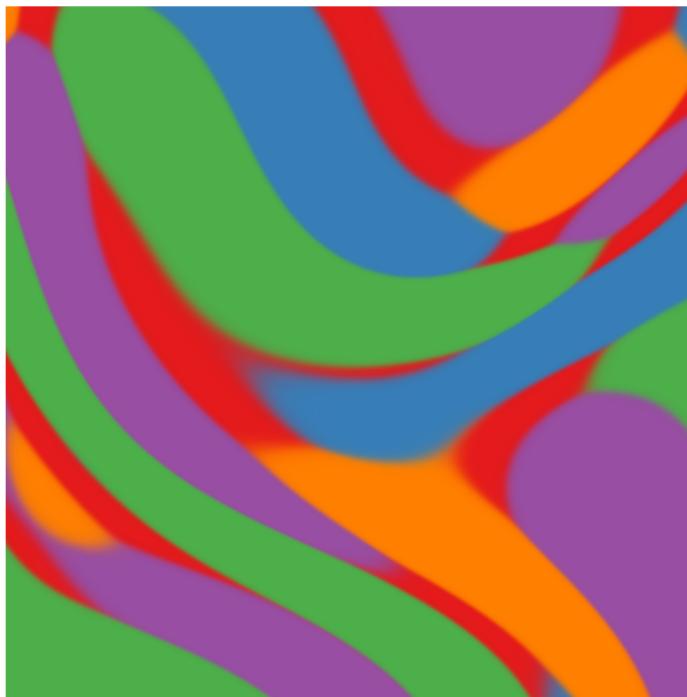
- ▶ Showing the  $x$  that gave rise to a particular  $y$
- ▶ (i.e. decision boundaries)



One Layers

# DEEP GAUSSIAN PROCESSES

- ▶ Showing the  $x$  that gave rise to a particular  $y$
- ▶ (i.e. decision boundaries)



Two Layers

# DEEP GAUSSIAN PROCESSES

- ▶ Showing the  $x$  that gave rise to a particular  $y$
- ▶ (i.e. decision boundaries)



Three Layers

# DEEP GAUSSIAN PROCESSES

- ▶ Showing the  $x$  that gave rise to a particular  $y$
- ▶ (i.e. decision boundaries)



Four Layers

# DEEP GAUSSIAN PROCESSES

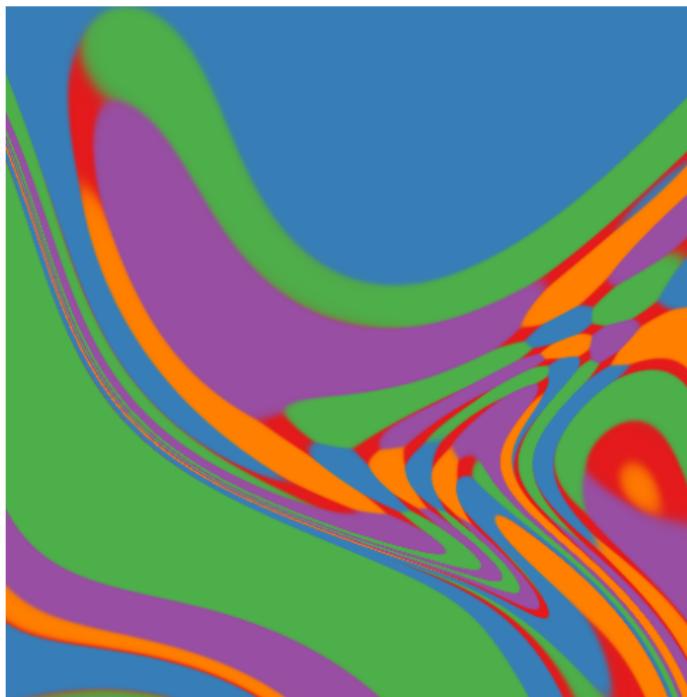
- ▶ Showing the  $x$  that gave rise to a particular  $y$
- ▶ (i.e. decision boundaries)



Five Layers

# DEEP GAUSSIAN PROCESSES

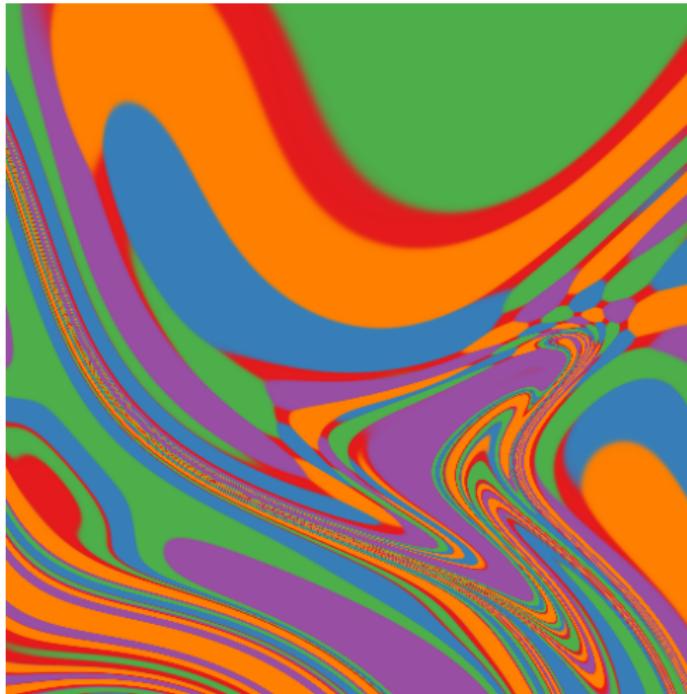
- ▶ Showing the  $x$  that gave rise to a particular  $y$
- ▶ (i.e. decision boundaries)



Ten Layers

# DEEP GAUSSIAN PROCESSES

- ▶ Showing the  $x$  that gave rise to a particular  $y$
- ▶ (i.e. decision boundaries)



Twenty Layers

# DEEP GAUSSIAN PROCESSES

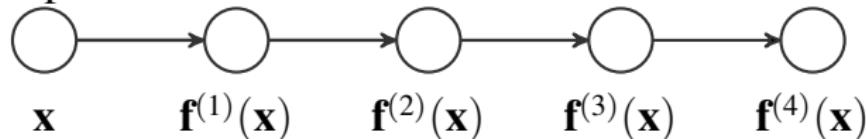
- ▶ Showing the  $x$  that gave rise to a particular  $y$
- ▶ (i.e. decision boundaries)



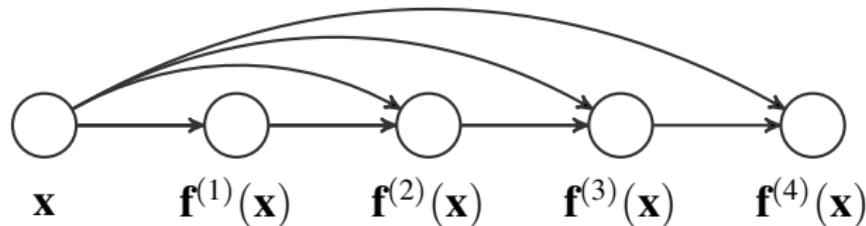
Forty Layers

# DEEP GAUSSIAN PROCESSES

- ▶ Only one degree of freedom in  $x$  is being captured!
- ▶ Again following Radford's thesis, connect every layer to input:



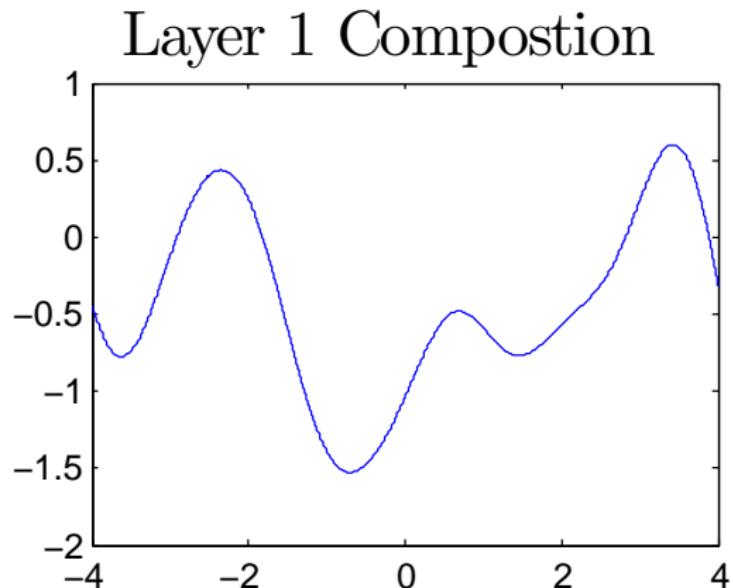
a) standard MLP architecture.



b) Input-connected architecture.

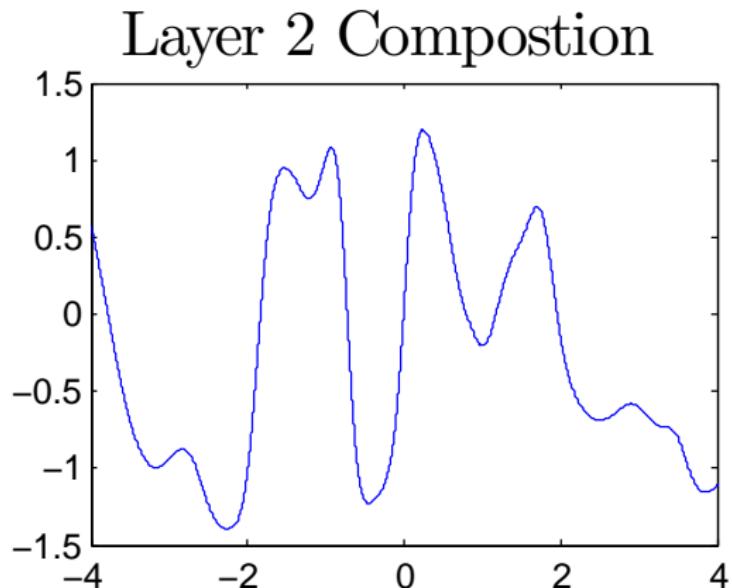
# DEEP GAUSSIAN PROCESSES

- ▶ Draws from input-connected one-dimensional deep GPs:



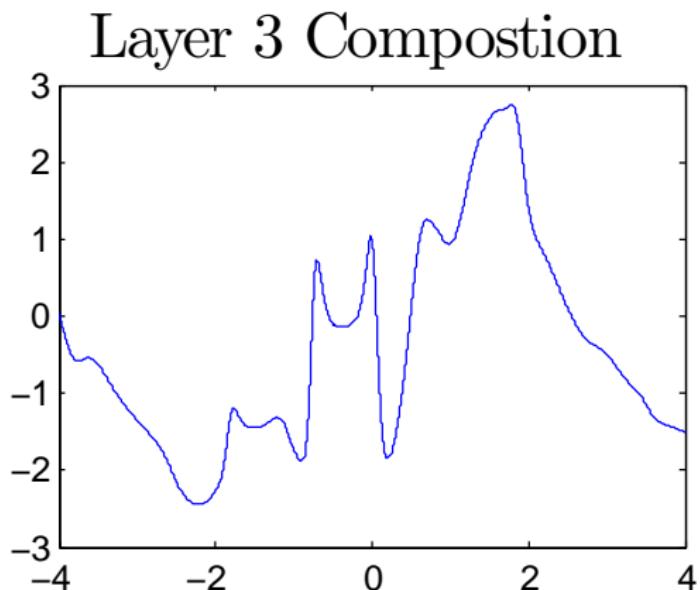
# DEEP GAUSSIAN PROCESSES

- ▶ Draws from input-connected one-dimensional deep GPs:



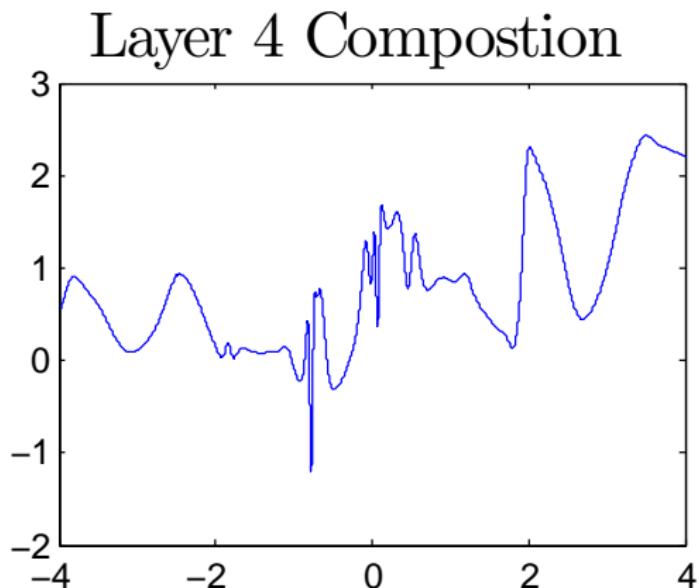
# DEEP GAUSSIAN PROCESSES

- ▶ Draws from input-connected one-dimensional deep GPs:



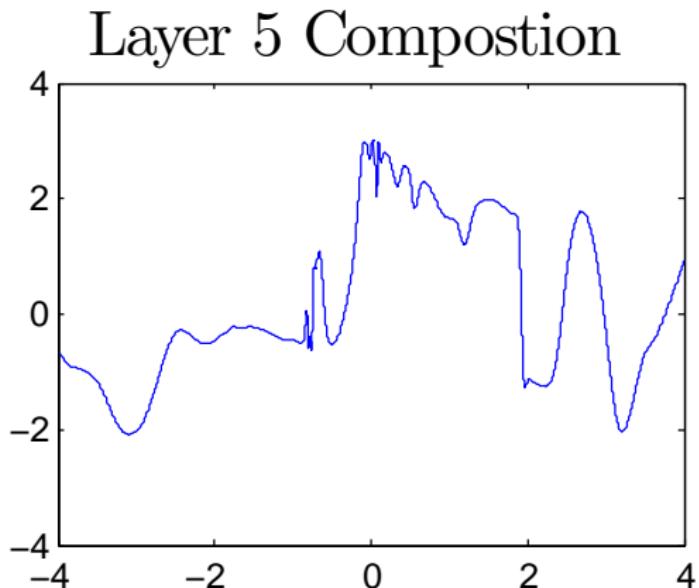
# DEEP GAUSSIAN PROCESSES

- ▶ Draws from input-connected one-dimensional deep GPs:



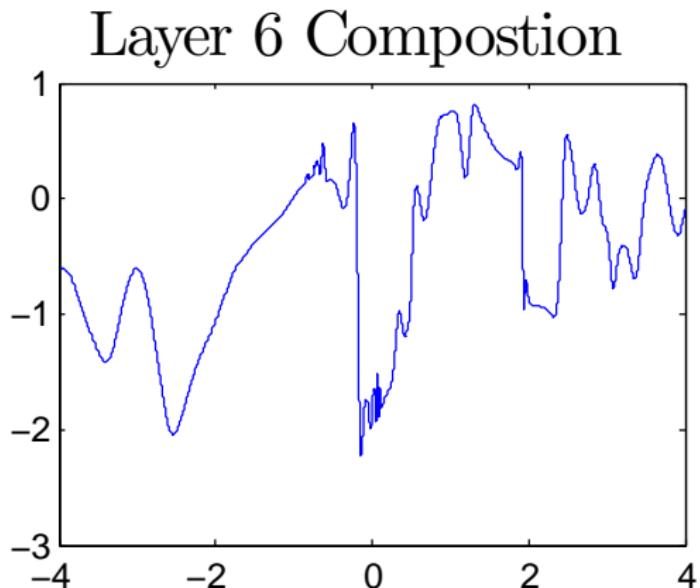
# DEEP GAUSSIAN PROCESSES

- ▶ Draws from input-connected one-dimensional deep GPs:



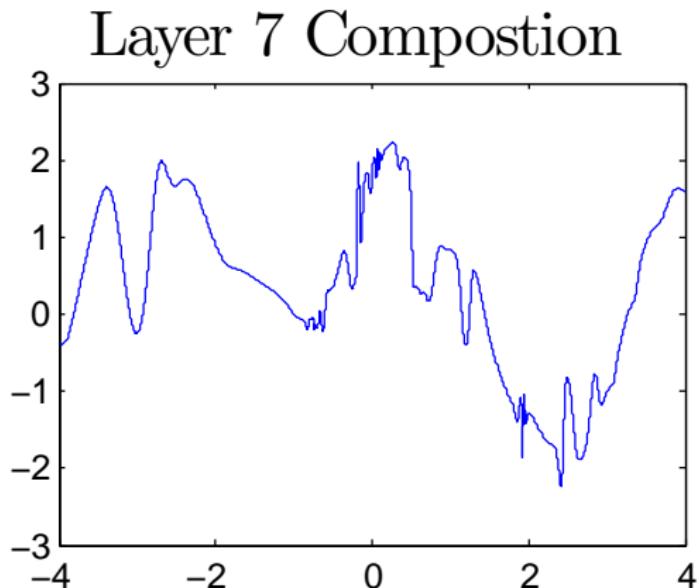
# DEEP GAUSSIAN PROCESSES

- ▶ Draws from input-connected one-dimensional deep GPs:



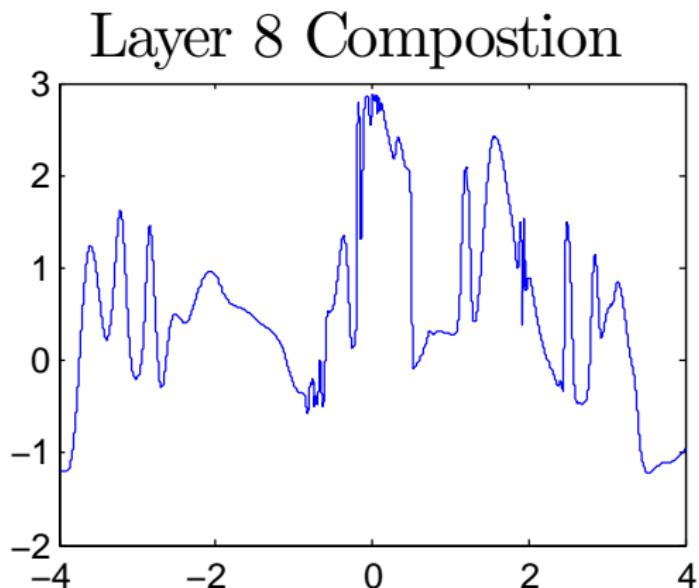
# DEEP GAUSSIAN PROCESSES

- ▶ Draws from input-connected one-dimensional deep GPs:



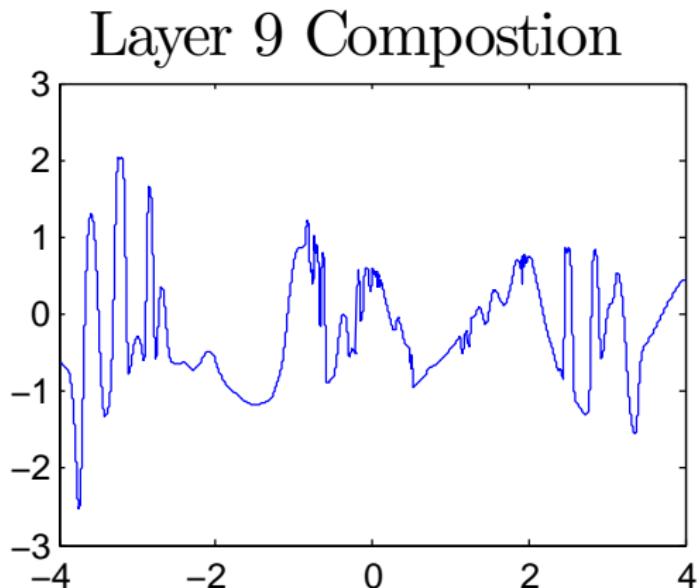
# DEEP GAUSSIAN PROCESSES

- ▶ Draws from input-connected one-dimensional deep GPs:



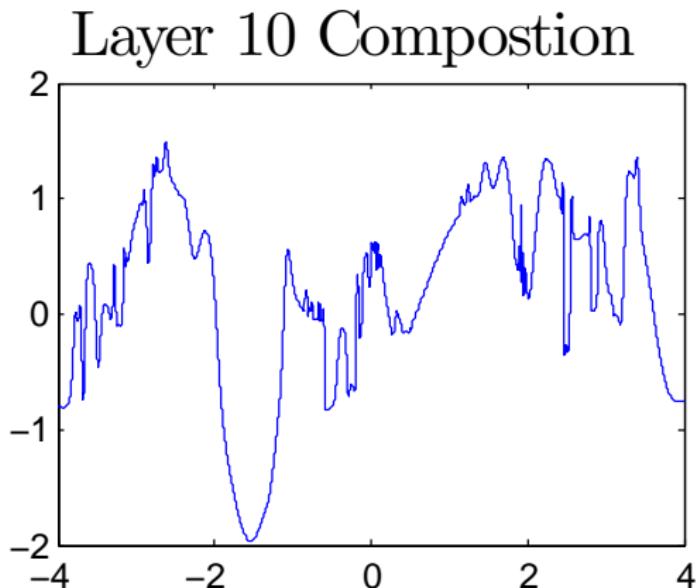
# DEEP GAUSSIAN PROCESSES

- ▶ Draws from input-connected one-dimensional deep GPs:



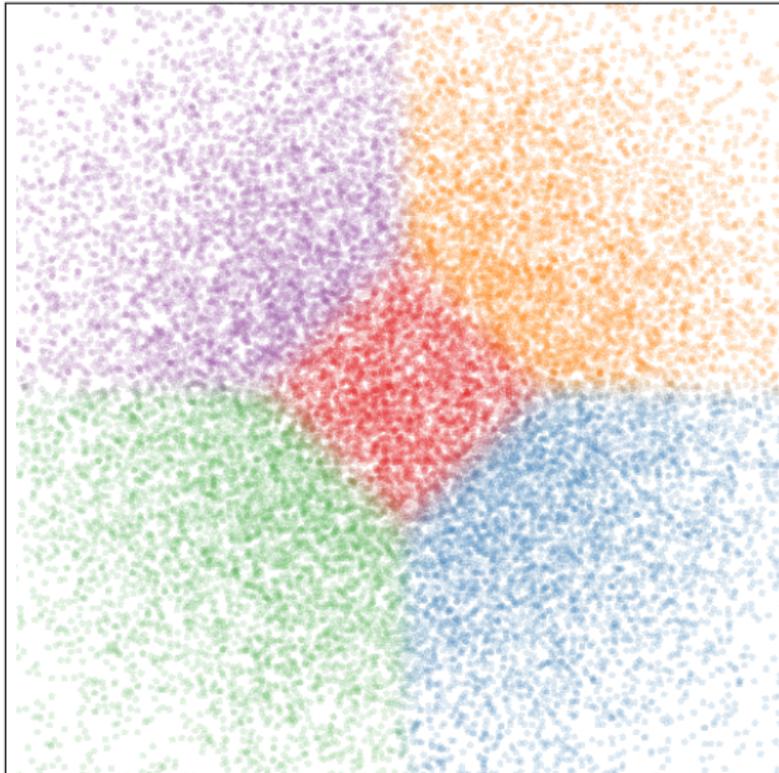
# DEEP GAUSSIAN PROCESSES

- ▶ Draws from input-connected one-dimensional deep GPs:



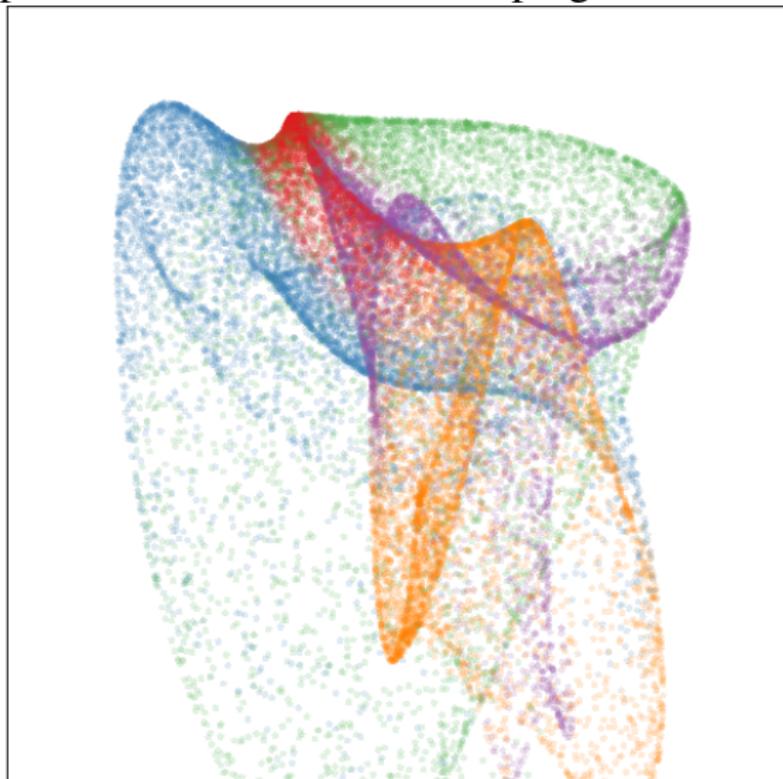
# DEEP GAUSSIAN PROCESSES

- ▶ input-connected 2D to 2D warpings of a Gaussian density:



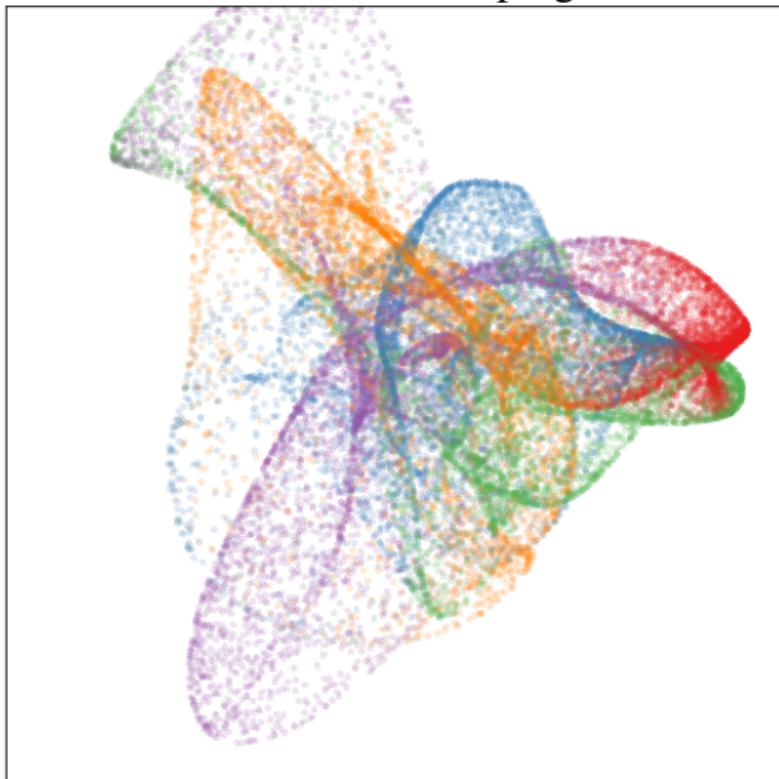
# DEEP GAUSSIAN PROCESSES

- ▶ input-connected 2D to 2D warpings of a Gaussian density:



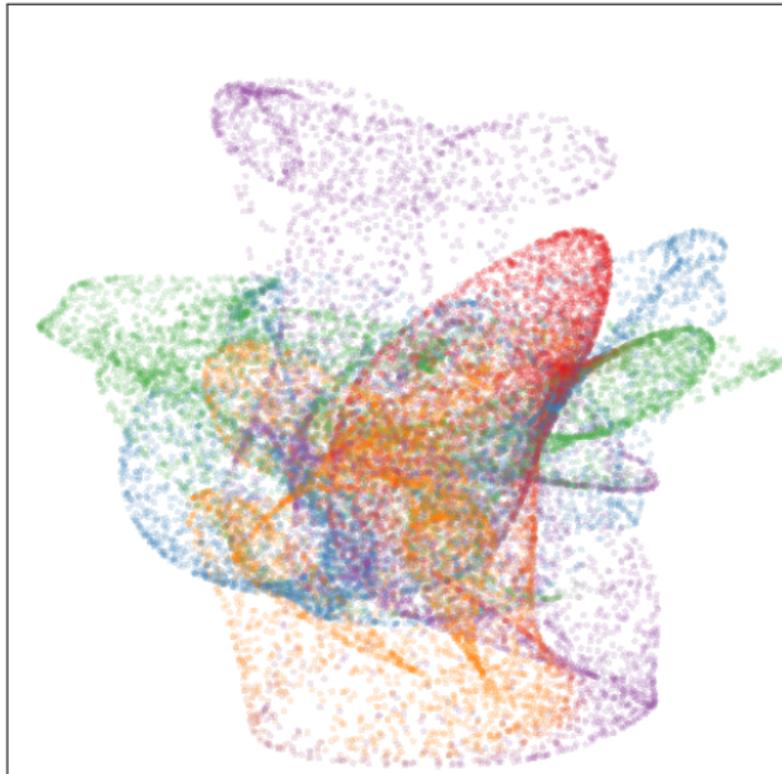
# DEEP GAUSSIAN PROCESSES

- ▶ input-connected 2D to 2D warpings of a Gaussian density:



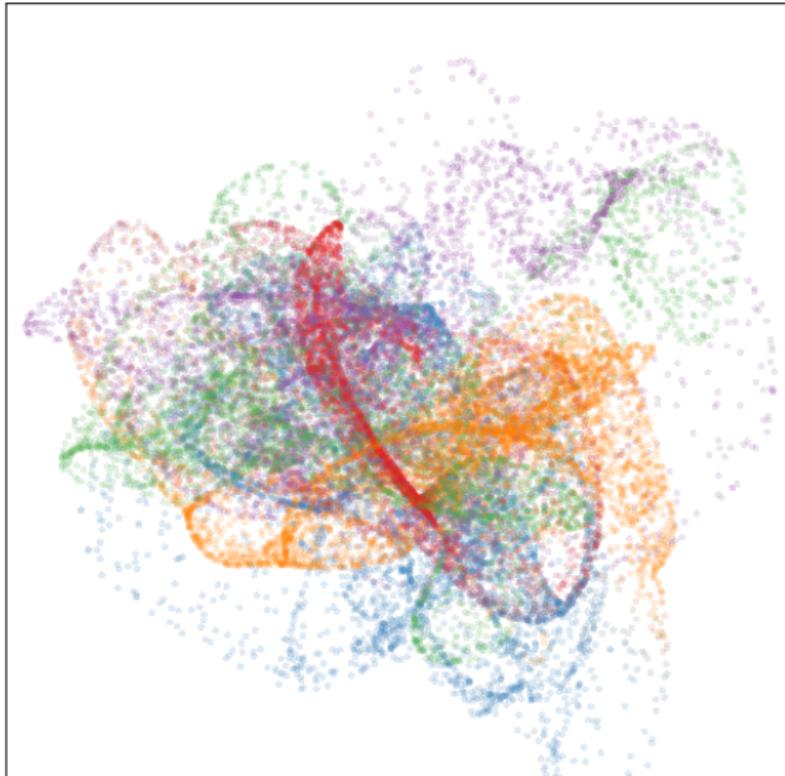
# DEEP GAUSSIAN PROCESSES

- ▶ input-connected 2D to 2D warpings of a Gaussian density:



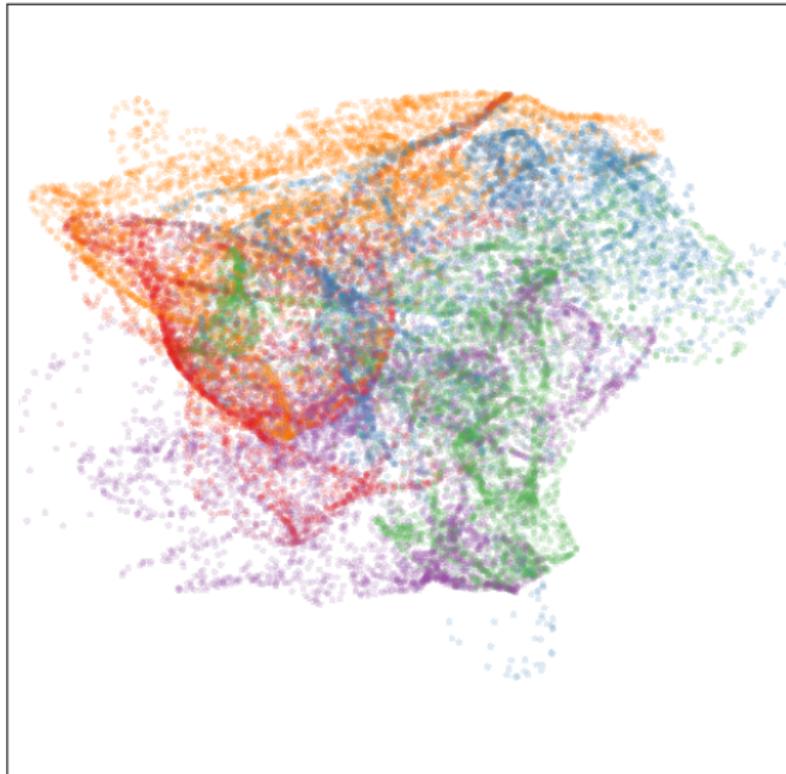
# DEEP GAUSSIAN PROCESSES

- ▶ input-connected 2D to 2D warpings of a Gaussian density:



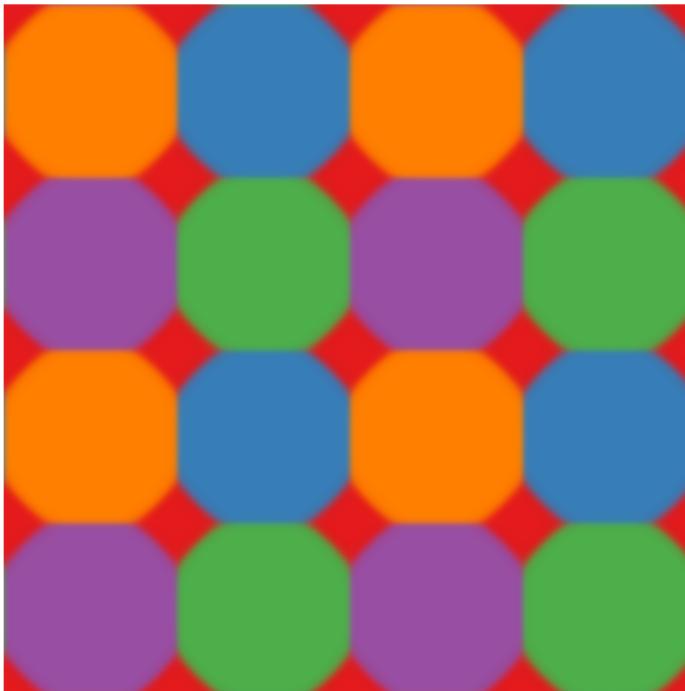
# DEEP GAUSSIAN PROCESSES

- ▶ input-connected 2D to 2D warpings of a Gaussian density:



# DEEP GAUSSIAN PROCESSES

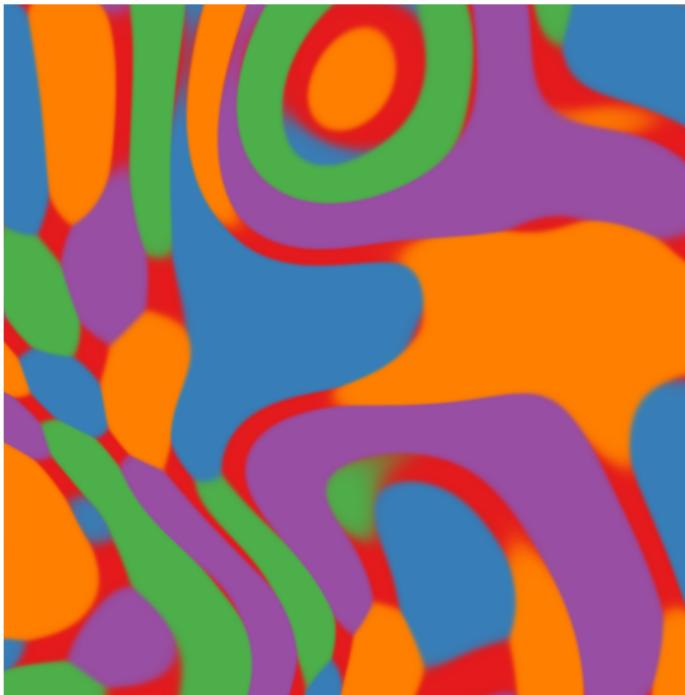
- ▶ Showing the  $x$  that gave rise to a particular  $y$
- ▶ (i.e. decision boundaries)



No warping

# DEEP GAUSSIAN PROCESSES

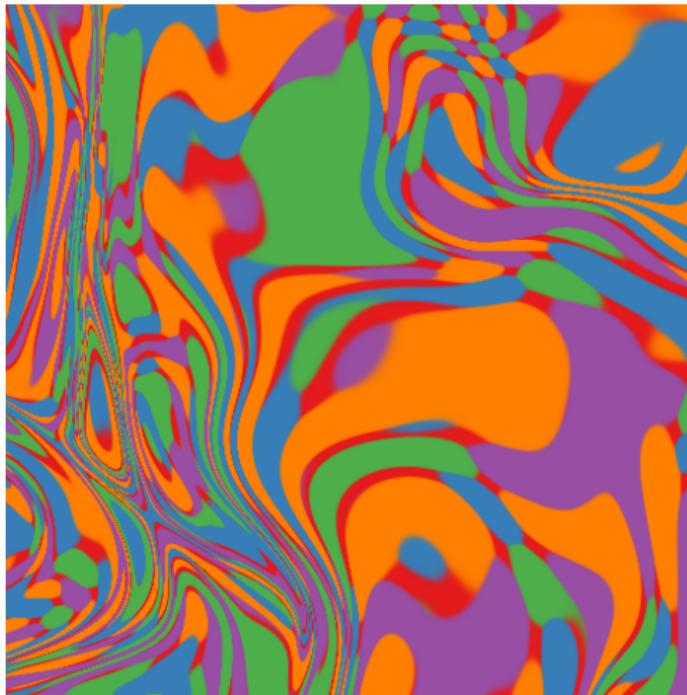
- ▶ Showing the  $x$  that gave rise to a particular  $y$
- ▶ (i.e. decision boundaries)



Two Layers

# DEEP GAUSSIAN PROCESSES

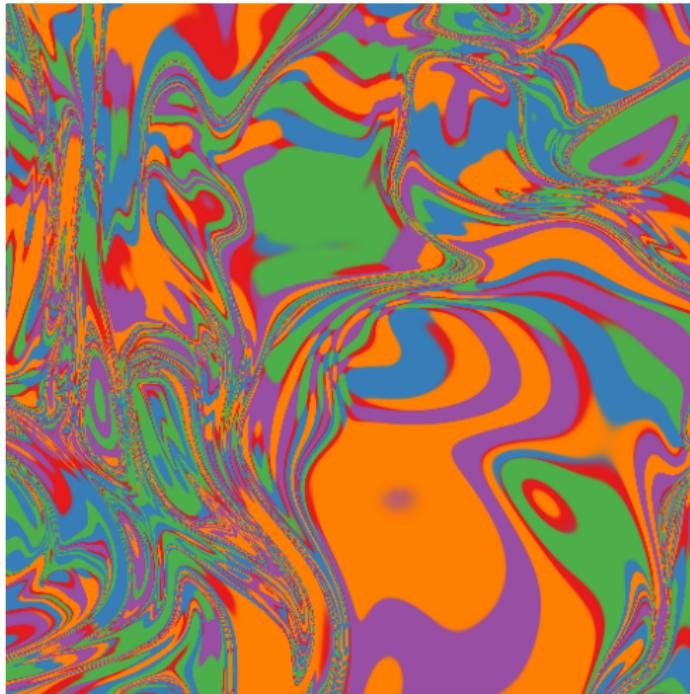
- ▶ Showing the  $x$  that gave rise to a particular  $y$
- ▶ (i.e. decision boundaries)



Ten Layers

# DEEP GAUSSIAN PROCESSES

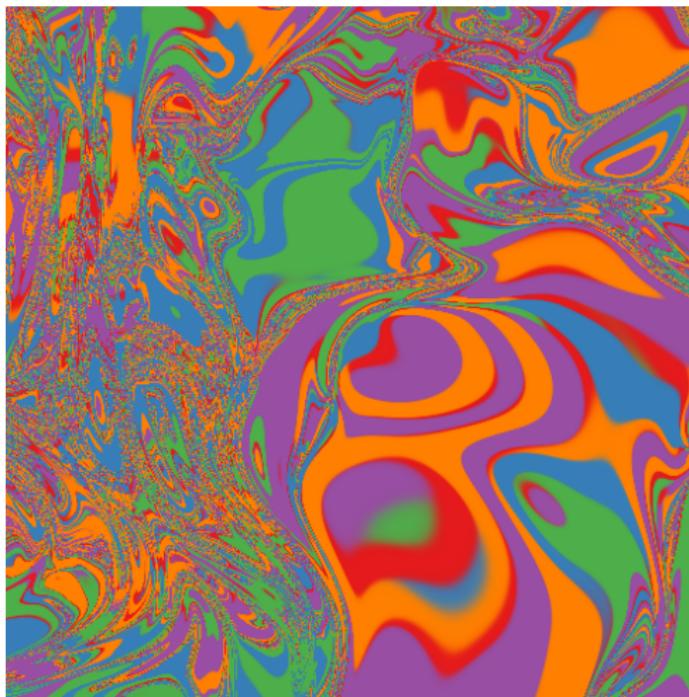
- ▶ Showing the  $x$  that gave rise to a particular  $y$
- ▶ (i.e. decision boundaries)



Twenty Layers

# DEEP GAUSSIAN PROCESSES

- ▶ Showing the  $x$  that gave rise to a particular  $y$
- ▶ (i.e. decision boundaries)



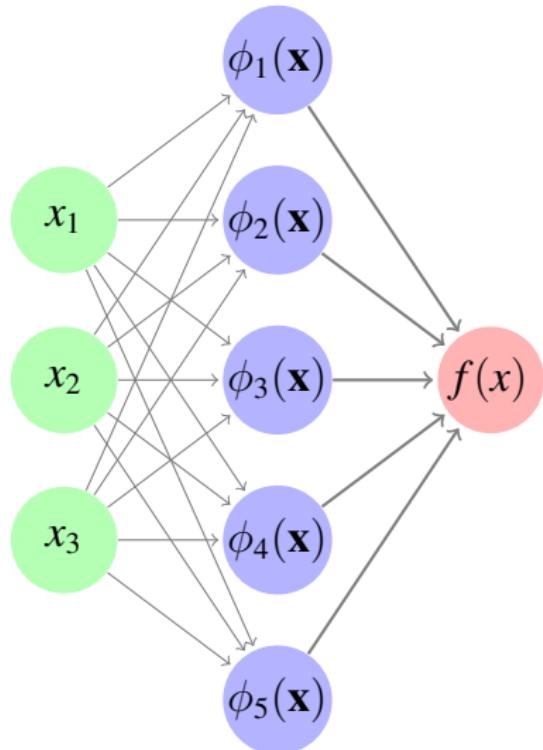
Forty Layers

# DROPOUT

---

- ▶ Dropout is a method for regularizing neural networks.
- ▶ Recipe:
  1. randomly set half of feature activations to zero
  2. Double the remaining features activations
  3. Average over all possible ways of dropping out
- ▶ Gives robustness, since neurons can't depend on one another.
- ▶ What do we get when we do dropout in GPs?

# DROPOUT ON FEATURE ACTIVATIONS



Original formulation:

$$f(\mathbf{x}) = \frac{1}{K} \sum_{i=1}^K \alpha_i \phi_i(\mathbf{x})$$

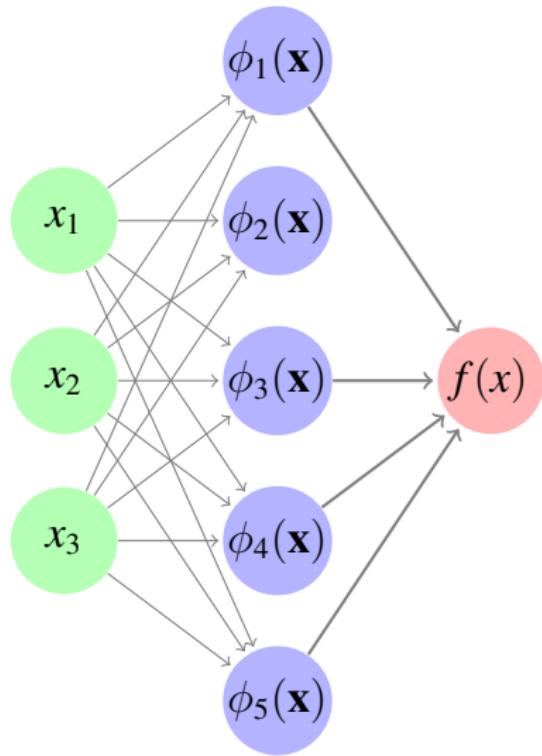
with any weight distribution,

$$\mathbb{E} [\alpha_i] = 0, \quad \mathbb{V} [\alpha_i] = \sigma^2$$

by CLT, gives a GP as  $K \rightarrow \infty$

$$\text{cov} \left[ \begin{array}{c} f(\mathbf{x}) \\ f(\mathbf{x}') \end{array} \right] \rightarrow \frac{\sigma^2}{K} \sum_{i=1}^K \phi_i(\mathbf{x}) \phi_i(\mathbf{x}')$$

# DROPOUT ON FEATURE ACTIVATIONS



Remove units with probability  $1/2$ :

$$f(\mathbf{x}) = \frac{1}{K} \sum_{i=1}^K \mathbf{r}_i \alpha_i \phi_i(\mathbf{x}) \quad \mathbf{r}_i \stackrel{\text{iid}}{\sim} \text{Ber}(1/2)$$

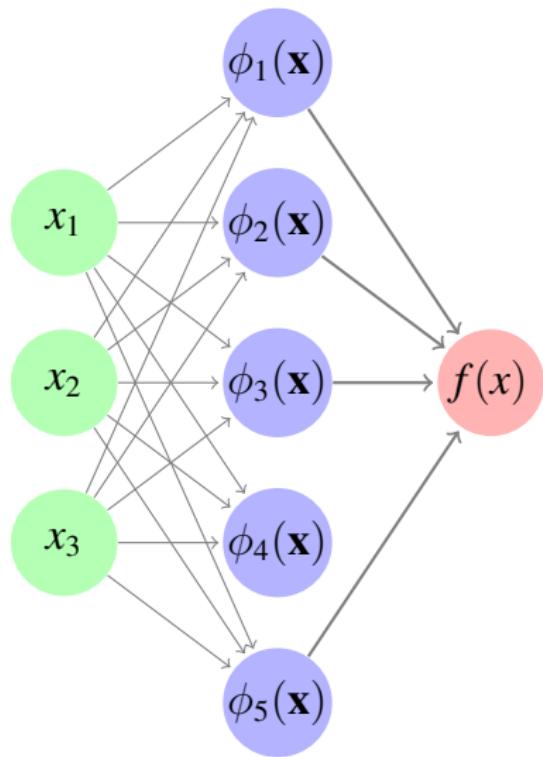
with any weight distribution,

$$\mathbb{E} [\mathbf{r}_i \alpha_i] = 0, \quad \mathbb{V} [\mathbf{r}_i \alpha_i] = \frac{1}{4} \sigma^2$$

by CLT, gives a GP as  $K \rightarrow \infty$

$$\text{cov} \begin{bmatrix} f(\mathbf{x}) \\ f(\mathbf{x}') \end{bmatrix} \rightarrow \frac{1}{4} \frac{\sigma^2}{K} \sum_{i=1}^K \phi_i(\mathbf{x}) \phi_i(\mathbf{x}')$$

# DROPOUT ON FEATURE ACTIVATIONS



Remove units with probability  $1/2$ :

$$f(\mathbf{x}) = \frac{1}{K} \sum_{i=1}^K \mathbf{r}_i \alpha_i \phi_i(\mathbf{x}) \quad \mathbf{r}_i \stackrel{\text{iid}}{\sim} \text{Ber}(1/2)$$

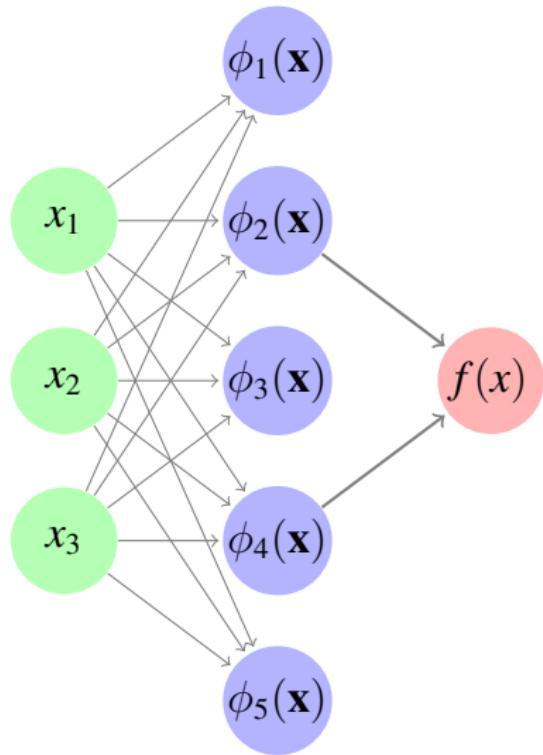
with any weight distribution,

$$\mathbb{E} [\mathbf{r}_i \alpha_i] = 0, \quad \mathbb{V} [\mathbf{r}_i \alpha_i] = \frac{1}{4} \sigma^2$$

by CLT, gives a GP as  $K \rightarrow \infty$

$$\text{cov} \begin{bmatrix} f(\mathbf{x}) \\ f(\mathbf{x}') \end{bmatrix} \rightarrow \frac{1}{4} \frac{\sigma^2}{K} \sum_{i=1}^K \phi_i(\mathbf{x}) \phi_i(\mathbf{x}')$$

# DROPOUT ON FEATURE ACTIVATIONS



Remove units with probability  $1/2$ :

$$f(\mathbf{x}) = \frac{1}{K} \sum_{i=1}^K \mathbf{r}_i \alpha_i \phi_i(\mathbf{x}) \quad \mathbf{r}_i \stackrel{\text{iid}}{\sim} \text{Ber}(1/2)$$

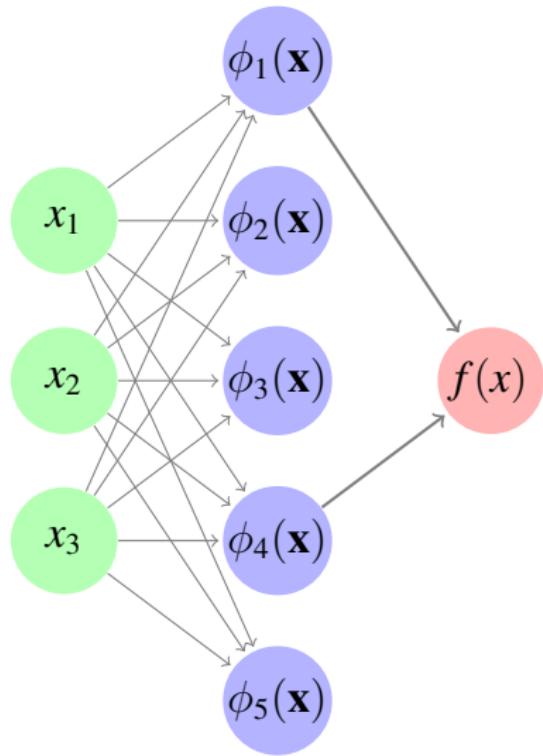
with any weight distribution,

$$\mathbb{E} [\mathbf{r}_i \alpha_i] = 0, \quad \mathbb{V} [\mathbf{r}_i \alpha_i] = \frac{1}{4} \sigma^2$$

by CLT, gives a GP as  $K \rightarrow \infty$

$$\text{cov} \begin{bmatrix} f(\mathbf{x}) \\ f(\mathbf{x}') \end{bmatrix} \rightarrow \frac{1}{4} \frac{\sigma^2}{K} \sum_{i=1}^K \phi_i(\mathbf{x}) \phi_i(\mathbf{x}')$$

# DROPOUT ON FEATURE ACTIVATIONS



Remove units with probability  $1/2$ :

$$f(\mathbf{x}) = \frac{1}{K} \sum_{i=1}^K \mathbf{r}_i \alpha_i \phi_i(\mathbf{x}) \quad \mathbf{r}_i \stackrel{\text{iid}}{\sim} \text{Ber}(1/2)$$

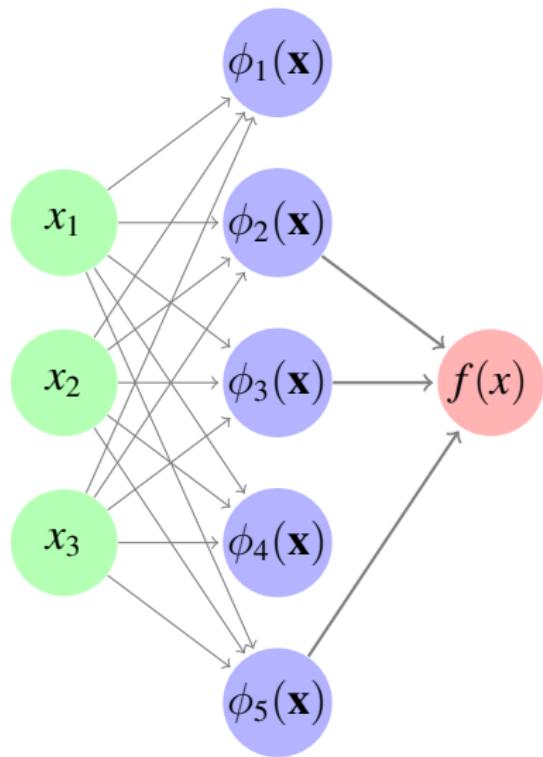
with any weight distribution,

$$\mathbb{E} [\mathbf{r}_i \alpha_i] = 0, \quad \mathbb{V} [\mathbf{r}_i \alpha_i] = \frac{1}{4} \sigma^2$$

by CLT, gives a GP as  $K \rightarrow \infty$

$$\text{cov} \begin{bmatrix} f(\mathbf{x}) \\ f(\mathbf{x}') \end{bmatrix} \rightarrow \frac{1}{4} \frac{\sigma^2}{K} \sum_{i=1}^K \phi_i(\mathbf{x}) \phi_i(\mathbf{x}')$$

# DROPOUT ON FEATURE ACTIVATIONS



Remove units with probability  $1/2$ :

$$f(\mathbf{x}) = \frac{1}{K} \sum_{i=1}^K \mathbf{r}_i \alpha_i \phi_i(\mathbf{x}) \quad \mathbf{r}_i \stackrel{\text{iid}}{\sim} \text{Ber}(1/2)$$

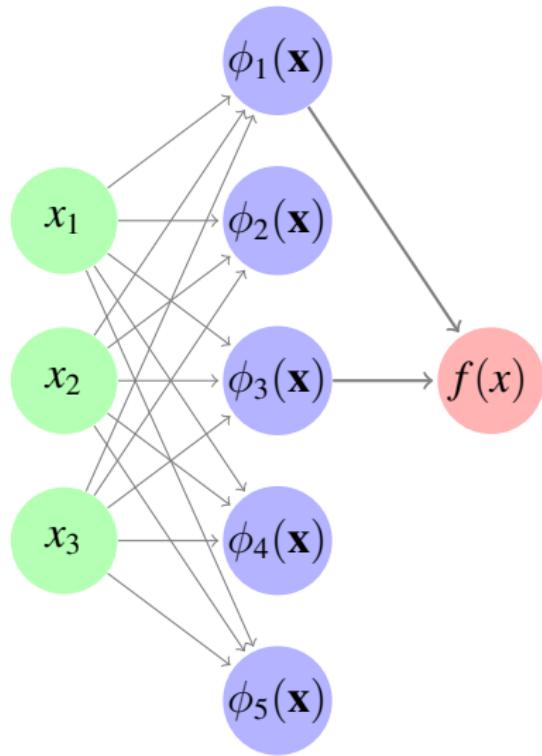
with any weight distribution,

$$\mathbb{E} [\mathbf{r}_i \alpha_i] = 0, \quad \mathbb{V} [\mathbf{r}_i \alpha_i] = \frac{1}{4} \sigma^2$$

by CLT, gives a GP as  $K \rightarrow \infty$

$$\text{cov} \begin{bmatrix} f(\mathbf{x}) \\ f(\mathbf{x}') \end{bmatrix} \rightarrow \frac{1}{4} \frac{\sigma^2}{K} \sum_{i=1}^K \phi_i(\mathbf{x}) \phi_i(\mathbf{x}')$$

# DROPOUT ON FEATURE ACTIVATIONS



Remove units with probability  $1/2$ :

$$f(\mathbf{x}) = \frac{1}{K} \sum_{i=1}^K \mathbf{r}_i \alpha_i \phi_i(\mathbf{x}) \quad \mathbf{r}_i \stackrel{\text{iid}}{\sim} \text{Ber}(1/2)$$

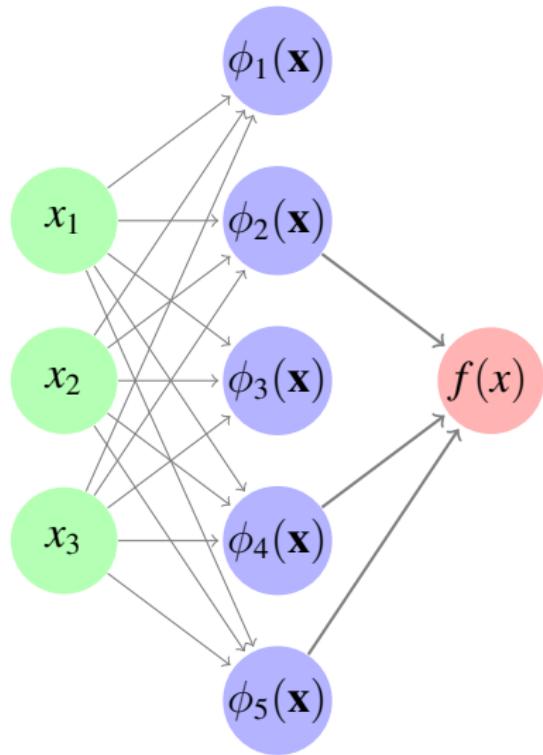
with any weight distribution,

$$\mathbb{E} [\mathbf{r}_i \alpha_i] = 0, \quad \mathbb{V} [\mathbf{r}_i \alpha_i] = \frac{1}{4} \sigma^2$$

by CLT, gives a GP as  $K \rightarrow \infty$

$$\text{cov} \begin{bmatrix} f(\mathbf{x}) \\ f(\mathbf{x}') \end{bmatrix} \rightarrow \frac{1}{4} \frac{\sigma^2}{K} \sum_{i=1}^K \phi_i(\mathbf{x}) \phi_i(\mathbf{x}')$$

# DROPOUT ON FEATURE ACTIVATIONS



Remove units with probability  $1/2$ :

$$f(\mathbf{x}) = \frac{1}{K} \sum_{i=1}^K \mathbf{r}_i \alpha_i \phi_i(\mathbf{x}) \quad \mathbf{r}_i \stackrel{\text{iid}}{\sim} \text{Ber}(1/2)$$

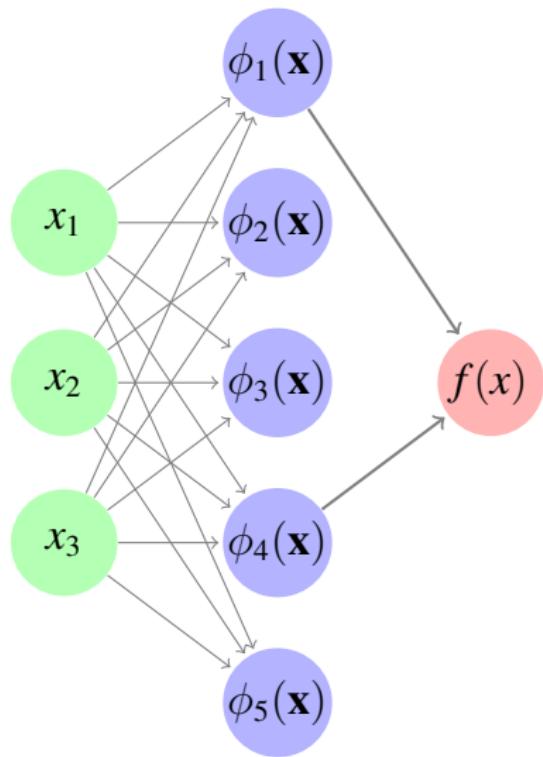
with any weight distribution,

$$\mathbb{E} [\mathbf{r}_i \alpha_i] = 0, \quad \mathbb{V} [\mathbf{r}_i \alpha_i] = \frac{1}{4} \sigma^2$$

by CLT, gives a GP as  $K \rightarrow \infty$

$$\text{cov} \begin{bmatrix} f(\mathbf{x}) \\ f(\mathbf{x}') \end{bmatrix} \rightarrow \frac{1}{4} \frac{\sigma^2}{K} \sum_{i=1}^K \phi_i(\mathbf{x}) \phi_i(\mathbf{x}')$$

# DROPOUT ON FEATURE ACTIVATIONS



Remove units with probability  $1/2$ :

$$f(\mathbf{x}) = \frac{1}{K} \sum_{i=1}^K \mathbf{r}_i \alpha_i \phi_i(\mathbf{x}) \quad \mathbf{r}_i \stackrel{\text{iid}}{\sim} \text{Ber}(1/2)$$

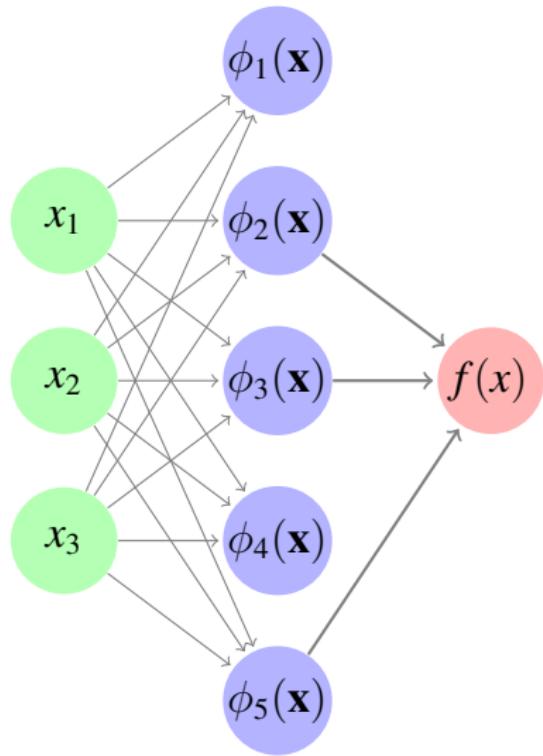
with any weight distribution,

$$\mathbb{E} [\mathbf{r}_i \alpha_i] = 0, \quad \mathbb{V} [\mathbf{r}_i \alpha_i] = \frac{1}{4} \sigma^2$$

by CLT, gives a GP as  $K \rightarrow \infty$

$$\text{cov} \begin{bmatrix} f(\mathbf{x}) \\ f(\mathbf{x}') \end{bmatrix} \rightarrow \frac{1}{4} \frac{\sigma^2}{K} \sum_{i=1}^K \phi_i(\mathbf{x}) \phi_i(\mathbf{x}')$$

# DROPOUT ON FEATURE ACTIVATIONS



Remove units with probability  $1/2$ :

$$f(\mathbf{x}) = \frac{1}{K} \sum_{i=1}^K \mathbf{r}_i \alpha_i \phi_i(\mathbf{x}) \quad \mathbf{r}_i \stackrel{\text{iid}}{\sim} \text{Ber}(1/2)$$

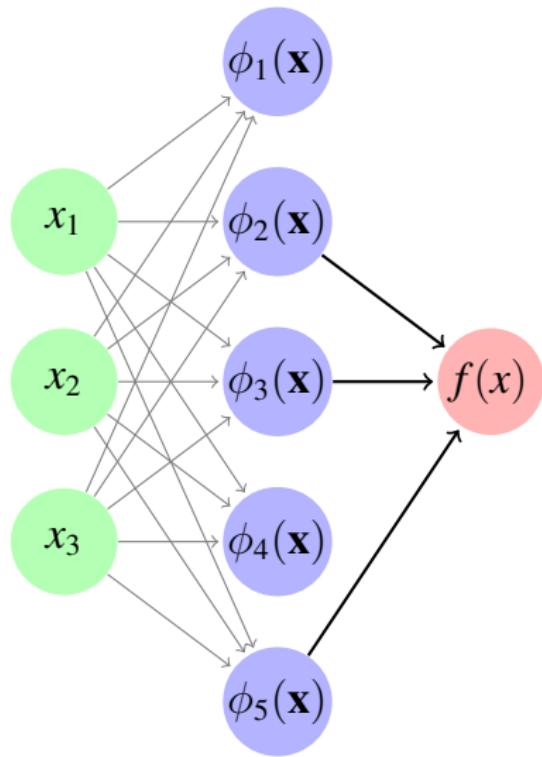
with any weight distribution,

$$\mathbb{E} [\mathbf{r}_i \alpha_i] = 0, \quad \mathbb{V} [\mathbf{r}_i \alpha_i] = \frac{1}{4} \sigma^2$$

by CLT, gives a GP as  $K \rightarrow \infty$

$$\text{cov} \begin{bmatrix} f(\mathbf{x}) \\ f(\mathbf{x}') \end{bmatrix} \rightarrow \frac{1}{4} \frac{\sigma^2}{K} \sum_{i=1}^K \phi_i(\mathbf{x}) \phi_i(\mathbf{x}')$$

# DROPOUT ON FEATURE ACTIVATIONS



Double output variance:

$$f(\mathbf{x}) = \frac{2}{K} \sum_{i=1}^K r_i \alpha_i \phi_i(\mathbf{x}) \quad r_i \stackrel{\text{iid}}{\sim} \text{Ber}(1/2)$$

with any weight distribution,

$$\mathbb{E}[2r_i \alpha_i] = 0, \quad \mathbb{V}[2r_i \alpha_i] = \frac{4}{4} \sigma^2$$

by CLT, gives a GP as  $K \rightarrow \infty$

$$\text{cov} \begin{bmatrix} f(\mathbf{x}) \\ f(\mathbf{x}') \end{bmatrix} \rightarrow \frac{4}{4} \frac{\sigma^2}{K} \sum_{i=1}^K \phi_i(\mathbf{x}) \phi_i(\mathbf{x}')$$

# DROPOUT ON FEATURE ACTIVATIONS

- ▶ Dropout on feature activations gives same GP
  - ▶ Averaging the same model doesn't do anything
- ▶ GPs were doing dropout all along? ☺
- ▶ GPs strange because any one feature doesn't matter.
- ▶ Is there a better way to drop out features that would lead to robustness?

# DROPOUT ON INPUTS

- ▶ Let  $k(\mathbf{x}, \mathbf{x}') = \prod_{d=1}^D k_d(\mathbf{x}_d, \mathbf{x}'_d)$
- ▶ Exact averaging over all dropouts is a mixture of GPs:

$$p(f(\mathbf{x})) = \frac{1}{2^D} \sum_{\mathbf{r} \in \{0,1\}^D} \text{GP} \left( 0, \prod_{d=1}^D k_d(\mathbf{x}_d, \mathbf{x}'_d)^{r_d} \right)$$

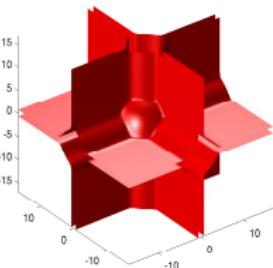
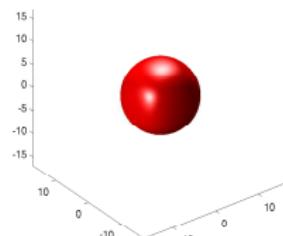
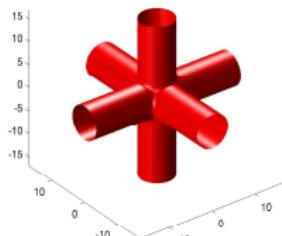
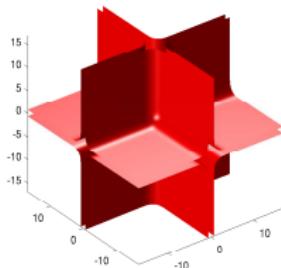
- ▶ Probably a nice model, but presumably intractable.
- ▶ In neural net literature, exponentially-many dropout models are averaged over in tractable ways.
- ▶ For instance, dropout mixture has same covariance as

$$f(\mathbf{x}) \sim \text{GP} \left( 0, \frac{1}{2^D} \sum_{\mathbf{r} \in \{0,1\}^D} \prod_{d=1}^D k_d(\mathbf{x}_d, \mathbf{x}'_d)^{r_d} \right)$$

# DROPOUT ON INPUTS

$$f(\mathbf{x}) \sim \text{GP} \left( 0, \sum_{\mathbf{r} \in \{0,1\}^D} \prod_{d=1}^D k_d(\mathbf{x}_d, \mathbf{x}'_d)^{r_d} \right)$$

- That's exactly additive GPs! Kernel isocontours look like:



$k_1 + k_2 + k_3$

$k_1 k_2 + k_2 k_3 + k_1 k_3$

$k_1 k_2 k_3$

all terms

- Can compute all  $2^D$  terms in  $\mathcal{O}(D^2)$
- lots of functions, each only depends on some of the inputs

# SUMMARY

- ▶ At least two different ways to make GPs deep:
  - ▶ composing kernels (inference easy, have to specify kernel)
  - ▶ composing GPs (inference is hard, but can learn structure)
- ▶ Kernel learning is a form of representation learning
- ▶ Building priors over functions can shed light on architectures choices or initialization strategies in a data-independent way.
- ▶ What sorts of structures do we want to be able to learn?
- ▶ We can bring tricks from the neural net literature back to GPs. (still need to try translation-invariant image kernels!)

# SUMMARY

- ▶ At least two different ways to make GPs deep:
  - ▶ composing kernels (inference easy, have to specify kernel)
  - ▶ composing GPs (inference is hard, but can learn structure)
- ▶ Kernel learning is a form of representation learning
- ▶ Building priors over functions can shed light on architectures choices or initialization strategies in a data-independent way.
- ▶ What sorts of structures do we want to be able to learn?
- ▶ We can bring tricks from the neural net literature back to GPs. (still need to try translation-invariant image kernels!)

Thanks!