

Deep Kernel Learning (AISTATS'16)



2016. 08. 09.

Saehoon Kim (Ph. D. candidate)

Machine Learning Group

Outline

- **Deep Kernel Learning (AISTATS'16)**
 - A. G. Wilson (CMU), Z. Hu (CMU), R. Salakhutdinov (Univ. of Toronto) and E. P. Xing (CMU)
- Gaussian Process Kernels for Pattern Discovery and Extrapolation (ICML'13)
 - A. G. Wilson (Cambridge) and R. P. Adams (Harvard)
- Kernel Interpolation for Scalable Structured Gaussian Process (KISS-GP) (ICML'15)
 - A. G. Wilson (CMU) and H. Nickisch (Phillips Research)

Deep Kernel Learning

- They transform the inputs of a spectral mixture base kernel with a deep neural network
- For a scalable kernel representation, they use local kernel interpolation, inducing points, and structure exploiting (Kronecker and Toeplitz) algebra
- On diverse applications (including a 2 million dataset), they showed the improved performance

Gaussian Process (GP)

- Gaussian Process is a collection of random variables, any finite number of which have a joint Gaussian distribution
 - It is completely specified by its mean and covariance functions

$$\begin{aligned}m(\mathbf{x}) &= \mathbb{E}[f(\mathbf{x})] \\k(\mathbf{x}, \mathbf{x}') &= \mathbb{E}[(f(\mathbf{x}) - m(\mathbf{x}))(f(\mathbf{x}') - m(\mathbf{x}'))]\end{aligned}$$

- It can be considered as **a distribution over functions**

$$\begin{aligned}f(\mathbf{x}) &\sim GP(m(\mathbf{x}), k(\mathbf{x}, \mathbf{x}')) \\ \mathbf{f} &\sim \mathcal{N}(\mathbf{m}, \mathbf{K})\end{aligned}$$

$$\begin{aligned}\mathbf{f} &\triangleq [f(\mathbf{x}_1), \dots, f(\mathbf{x}_n)]^\top & \mathbf{m} &\triangleq [m(\mathbf{x}_1), \dots, m(\mathbf{x}_n)]^\top \\ & & [\mathbf{K}]_{ij} &\triangleq k(\mathbf{x}_i, \mathbf{x}_j)\end{aligned}$$

Gaussian Process (GP)

- Assuming a Gaussian noise, $y(\mathbf{x})|f(\mathbf{x}) \sim \mathcal{N}(y(\mathbf{x}); f(\mathbf{x}), \sigma^2)$, the predictive distribution is given by

$$\mathbf{f}_*|\mathbf{X}_*, \mathbf{X}, \mathbf{y}, \boldsymbol{\theta} \sim \mathcal{N}(\mathbb{E}[\mathbf{f}_*], \text{cov}(\mathbf{f}_*))$$

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

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

- The marginal log-likelihood is given by

$$\log p(\mathbf{y}|\mathbf{X}; \mathbf{w}, \boldsymbol{\theta}) \propto - [\mathbf{y}^\top (\mathbf{K} + \sigma^2 \mathbf{I})^{-1} \mathbf{y} + \log |\mathbf{K} + \sigma^2 \mathbf{I}|]$$

Gaussian Process vs. Neural Network

- Bayesian NN with infinitely many hidden units converged to GP with a particular kernel (covariance) function

| | Gaussian Process | Neural Network |
|---------------------------------|-------------------|----------------|
| Representation Learning | - | Possible |
| Uncertainty Measuring | Possible | - |
| Extrapolation | Possible | - |
| Time Complexity (Inference) | High ($O(N^2)$) | Low |
| Space Complexity (Inference) | High ($O(N^2)$) | Low |

Deep Kernel Learning

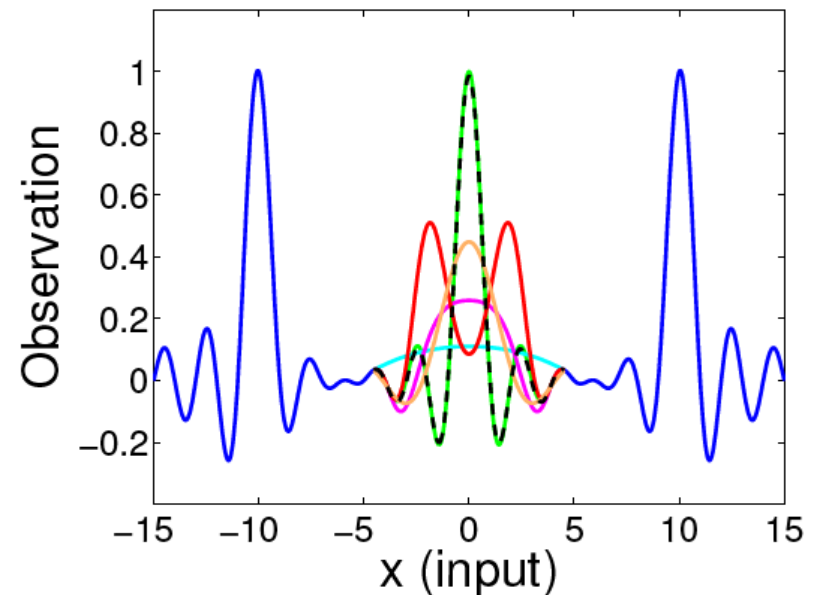
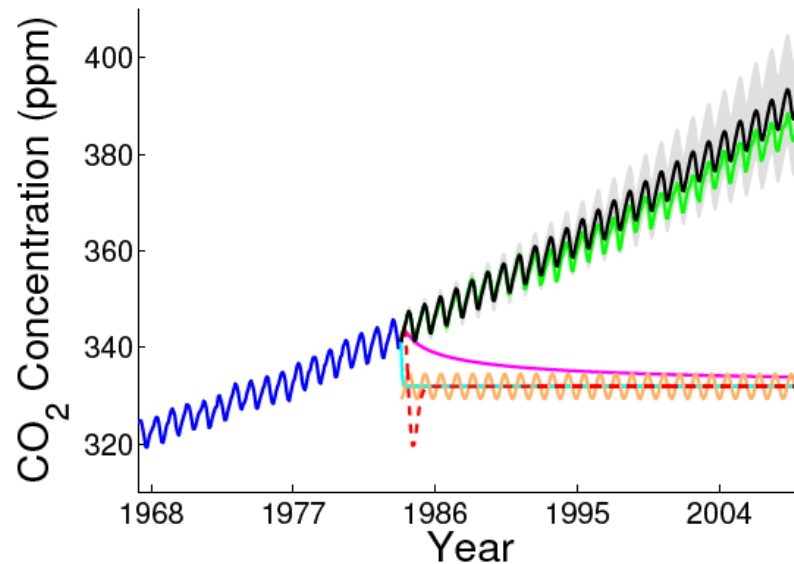
- They show ***how to construct kernels which encapsulate the expressive power of neural networks***, and **how to learn these kernels in a scalable GP framework**
- Starting from a base kernel with hyperparameters θ , the new kernel is derived as

$$k(\mathbf{x}_i, \mathbf{x}_j | \boldsymbol{\theta}) \rightarrow k(g(\mathbf{x}_i; \mathbf{w}), g(\mathbf{x}_j; \mathbf{w}) | \boldsymbol{\theta}, \mathbf{w})$$

- Which base kernel is used?
- How to learn the hyperparameters?

Why we need spectral mixture base kernels?

- To discover patterns and enable extrapolation, the authors explore flexible classes of kernels that support a broad class of stationary kernels



Bochner's theorem

Theorem 1. (Bochner) *A complex-valued function k on \mathbb{R}^d is the covariance function of a weakly stationary mean square continuous complex-valued random process on \mathbb{R}^d if and only if it can be represented as*

$$k(\tau) = \int_{\mathbb{R}^d} \exp(2\pi i s^\top \tau) p(s) ds,$$

where $p(s)$ is a probability density function.

It is often to call $p(s)$ as the *spectral density* or *power spectrum* of k , and k and p are Fourier duals:

$$\begin{aligned} k(\tau) &= \int p(s) \exp(2\pi i s^\top \tau) ds \\ p(s) &= \int k(\tau) \exp(-2\pi i s^\top \tau) d\tau \end{aligned}$$

Limitation of SE kernel

- The spectral density of squared exponential (SE) kernel is given as

$$k_{SE}(\mathbf{x}, \mathbf{x}') = \exp \left(-\frac{1}{2l^2} \|\mathbf{x} - \mathbf{x}'\|_2^2 \right)$$

$$S_{SE}(\mathbf{x}, \mathbf{x}') = (2\pi l^2)^{d/2} \exp(-2\pi^2 l^2 s^2)$$

- SE kernels and mixtures of SE kernels live in a very small region of the set of possible stationary kernels
 - They correspond only to Gaussian spectral densities centered on the origin!

Spectral Mixture base kernels

- They employ a mixture of Gaussian spectral densities that have non-zero means
- Consider a mixture model:

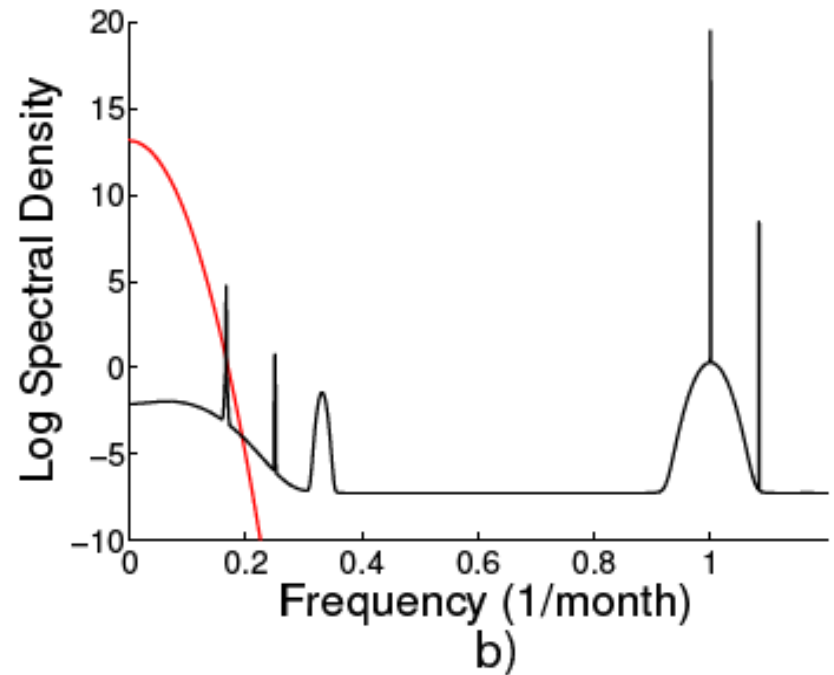
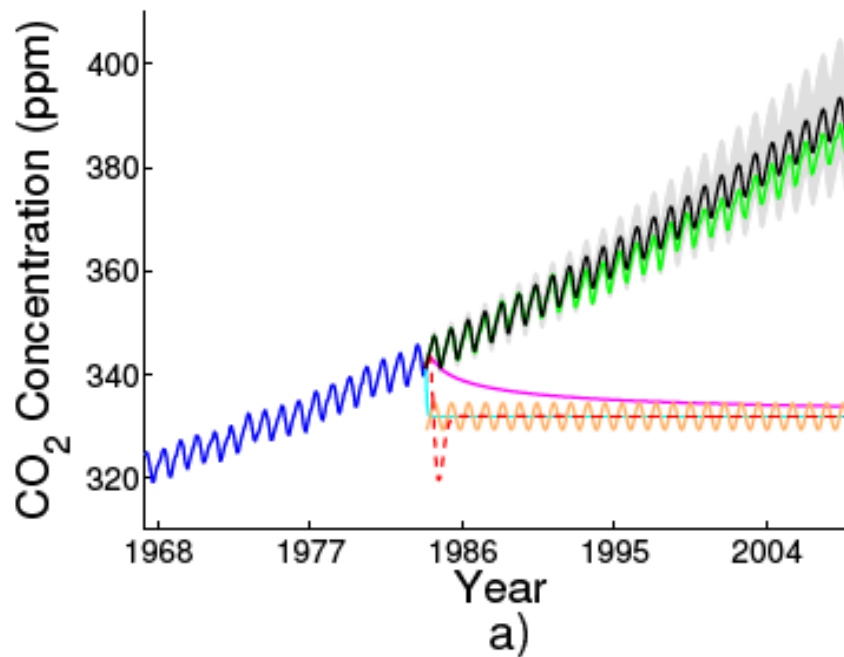
$$\phi(s; \mu, \sigma^2) = \frac{1}{\sqrt{2\pi\sigma^2}} \exp\left(-\frac{1}{2\sigma^2}(s - \mu)^2\right)$$

$$p(s) = \sum_{i=1}^Q w_i \phi(s; \boldsymbol{\mu}_i, \boldsymbol{\Sigma}_i)$$

$$k(\boldsymbol{\tau}) = \sum_{i=1}^Q w_i \frac{|\boldsymbol{\Sigma}_i|^{1/2}}{(2\pi)^{d/2}} \exp\left(-\frac{1}{2} \|\boldsymbol{\Sigma}_i^{1/2}(\boldsymbol{\tau})\|_2^2\right) \cos(\boldsymbol{\tau}^\top 2\pi \boldsymbol{\mu}_i)$$

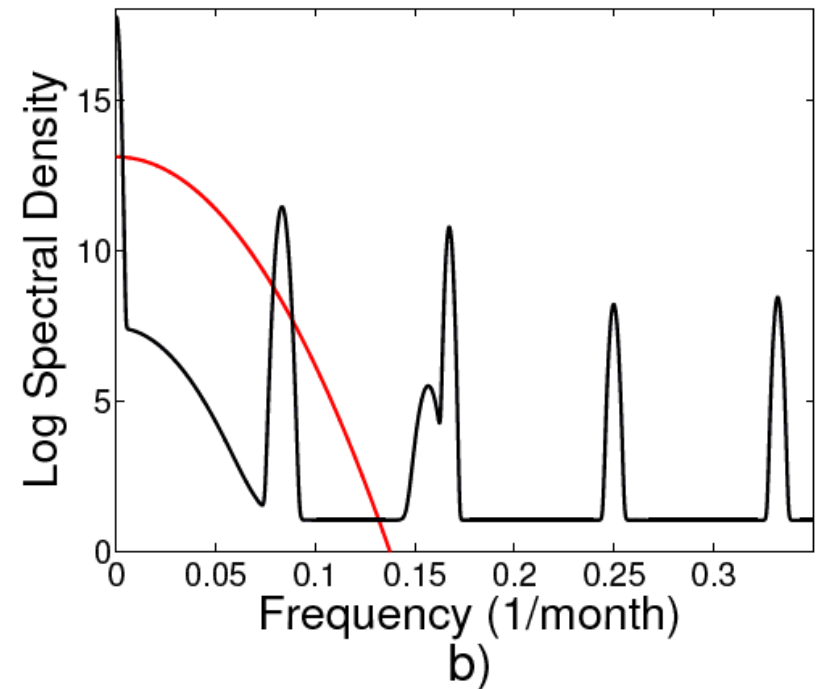
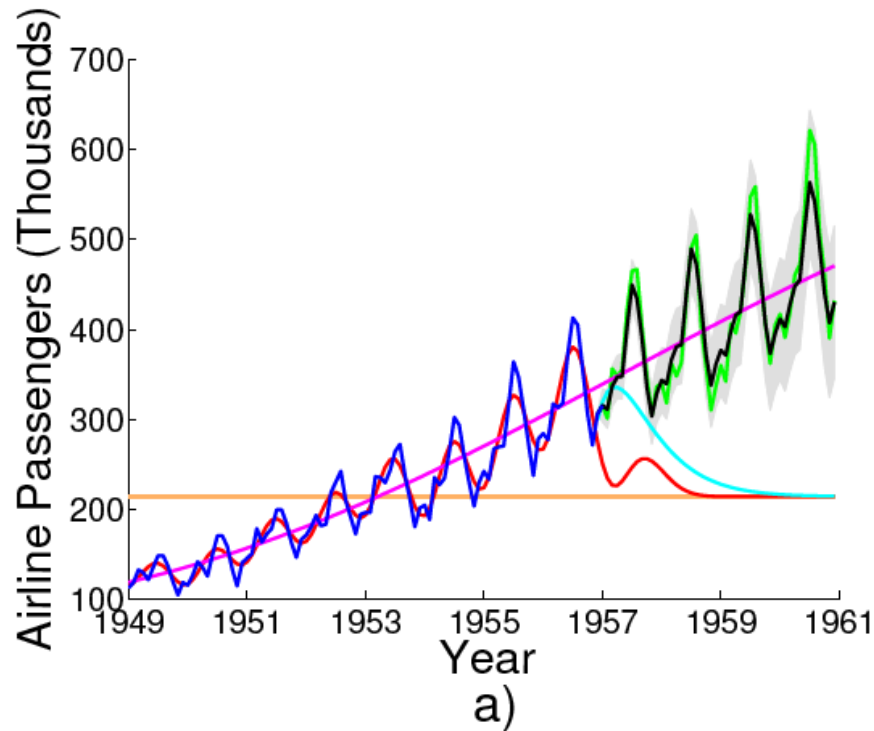
Some experimental results (1/2)

- SM kernels are tractable, allowing us to derive the closed form of marginal likelihood
- Extrapolation Atmospheric CO₂



Some experimental results (2/2)

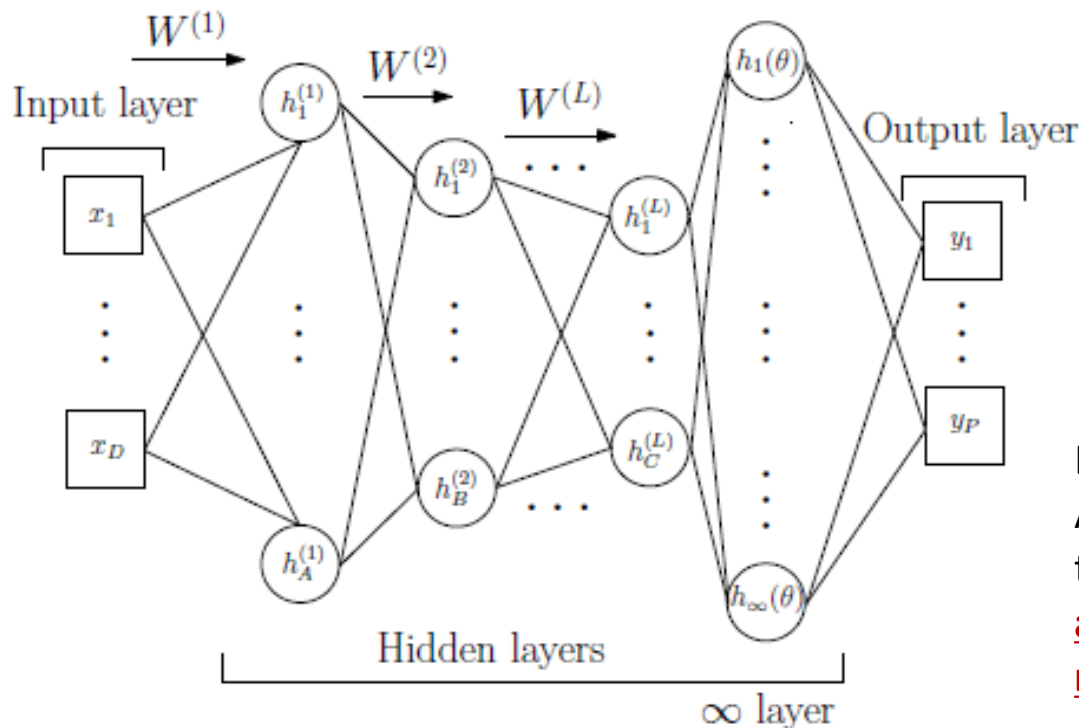
- Airline Passenger Data



Back to deep kernel learning

- Incorporating a neural network into SM kernels

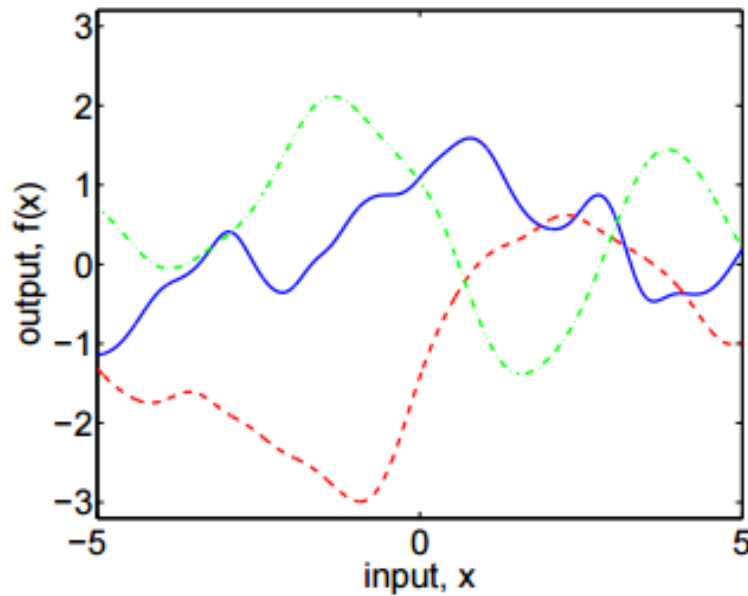
$$k(\mathbf{x}, \mathbf{x}') = \sum_{i=1}^Q w_i \frac{|\Sigma_i|^{1/2}}{(2\pi)^{d/2}} \exp\left(-\frac{1}{2} \|\Sigma_i^{1/2} (g(\mathbf{x}; \mathbf{w}) - g(\mathbf{x}'; \mathbf{w}))\|_2^2\right) \cos\left((g(\mathbf{x}; \mathbf{w}) - g(\mathbf{x}'; \mathbf{w}))^\top 2\pi \boldsymbol{\mu}_i\right)$$



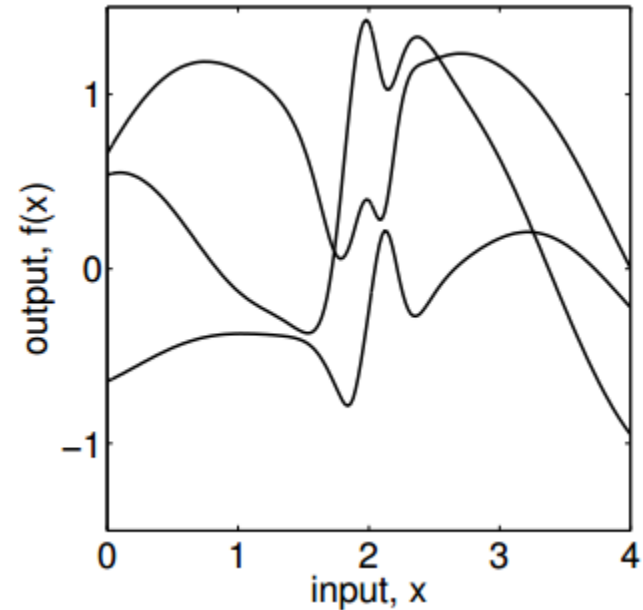
Deep kernel Learning:
A GP with a NN maps the inputs through **L hidden layers** followed by **a hidden layer with an infinite number of basis functions**

Nonstationary kernels

- Now, it is an example of nonstationary kernels..



Random functions drawn from GP
with a **stationary kernel**



Random functions drawn from GP
with a **non-stationary kernel**

How to learn hyperparameters?

- Maximize the marginal log-likelihood w.r.t. the kernel hyperparameter and the weights of the network
 - The weights of the network should be pre-trained!
 - Any type of network can be embedded
- A simple chain rule is used to update the parameters:

$$\log p(\mathbf{y}|\mathbf{X}; \mathbf{w}, \boldsymbol{\theta}) \propto - \left[\mathbf{y}^\top (\mathbf{K} + \sigma^2 \mathbf{I})^{-1} \mathbf{y} + \log |\mathbf{K} + \sigma^2 \mathbf{I}| \right]$$

$$\frac{\partial \mathcal{L}}{\partial \boldsymbol{\theta}} = \frac{\partial \mathcal{L}}{\partial \mathbf{K}} \frac{\partial \mathbf{K}}{\partial \boldsymbol{\theta}}, \quad \frac{\partial \mathcal{L}}{\partial \mathbf{w}} = \frac{\partial \mathcal{L}}{\partial \mathbf{K}} \frac{\partial \mathbf{K}}{\partial g(\mathbf{x}; \mathbf{w})} \frac{\partial g(\mathbf{x}; \mathbf{w})}{\partial \mathbf{w}}$$

$$\frac{\partial \mathcal{L}}{\partial \mathbf{K}} = \frac{1}{2} (\mathbf{K}^{-1} \mathbf{y} \mathbf{y}^\top \mathbf{K}^{-1} - \mathbf{K}^{-1})$$

A Scalable GP framework

- It is very time-consuming to compute the partial derivative required at every update step:

$$\frac{\partial \mathcal{L}}{\partial \mathbf{K}} = \frac{1}{2} (\mathbf{K}^{-1} \mathbf{y} \mathbf{y}^{\top} \mathbf{K}^{-1} - \mathbf{K}^{-1})$$

- Adopt the scalable GP framework to accelerate the update
 - Exploit the existing structure of the covariance matrix to make the inversion and matrix multiplication fast
 - The approximated covariance matrix should be similar to the original covariance matrix as much as possible

Inducing Point Methods

- A sort of low-rank approximation
- With m landmark points, the kernel is approximated as

$$\mathbf{K} \approx \mathbf{K}(\mathbf{X}, \mathbf{L})\mathbf{K}(\mathbf{L}, \mathbf{L})^{-1}\mathbf{K}(\mathbf{L}, \mathbf{X})$$

- It takes $\mathcal{O}(m^2n + m^3)$ computations and $\mathcal{O}(mn + m^2)$ storage for hyperparameter learning
- It takes $\mathcal{O}(m)$ and $\mathcal{O}(m^2)$ for predictive mean and variance

Fast Structure Exploiting Inference

- If we have inputs on a Cartesian grid, $\mathbf{x} \in \mathcal{X}_1 \times \cdots \times \mathcal{X}_P$,

$$k(\mathbf{x}_i, \mathbf{x}_j) = \prod_{p=1}^P k(\mathbf{x}_i^p, \mathbf{x}_j^p)$$

$$\mathbf{K} = \mathbf{K}_1 \otimes \cdots \otimes \mathbf{K}_p$$

- We can efficiently find the SVD of \mathbf{K} by separately computing SVD of each of \mathbf{K}_i .
 - Kernel matrix inversion (including matrix-vector multiplication) now becomes much faster!
 - A Cartesian grid can be approximated by k-means

Experimental Results

- Face orientation extraction
 - Predicting the orientation of a face
 - Constructing datasets of 28x28 face images by randomly rotating (-90 to + 90)



- Neural network architecture
 - 2 convolutional layers
 - 4 fully-connected layers
 - Two-dimensional outputs (that can be interpreted as the orientation of a face)

Experimental Results

- Better than CNN?

| Datasets | GP | DBN+GP | CNN | CNN+GP | DKL |
|----------|-------|--------|------|--------|-------------|
| Olivetti | 16.33 | 6.42 | 6.34 | 6.42 | 6.07 |
| MNIST | 1.25 | 1.03 | 0.59 | 0.56 | 0.53 |

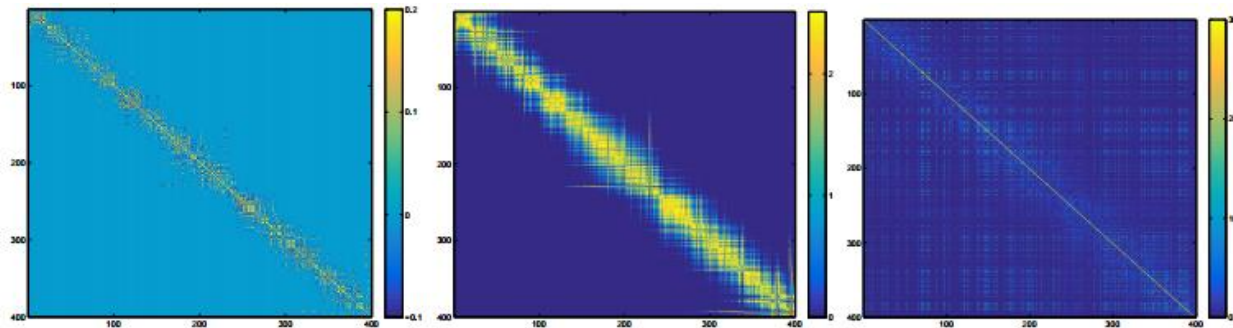


Figure 5: **Left:** The induced covariance matrix using DKL-SM kernel on a set of test cases, where the test samples are ordered by the *orientations* of the input faces. **Middle:** The respective covariance matrix using DKL-RBF kernel. **Right:** The respective covariance matrix using regular RBF kernel. The models are trained with $n = 12,000$, and $Q = 4$ for the SM base kernel.

Fooling images for CNN

- CNNs easily produce very high-confidence of unrecognizable images
- CNNs typically return many “fooling images” into one-of-classes with very high confidence (99.9%)
- It appears because it is a discriminative model?

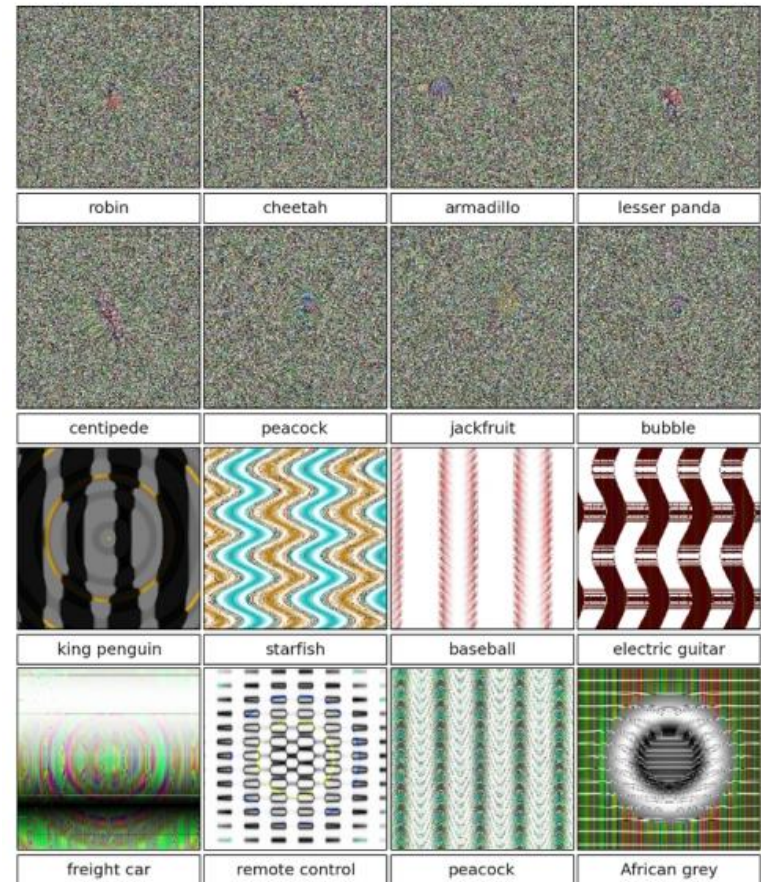


Figure 1. Evolved images that are unrecognizable to humans, but that state-of-the-art DNNs trained on ImageNet believe with $\geq 99.6\%$ certainty to be a familiar object. This result highlights differences between how DNNs and humans recognize objects. Images are either directly (*top*) or indirectly (*bottom*) encoded.

Fooling image examples

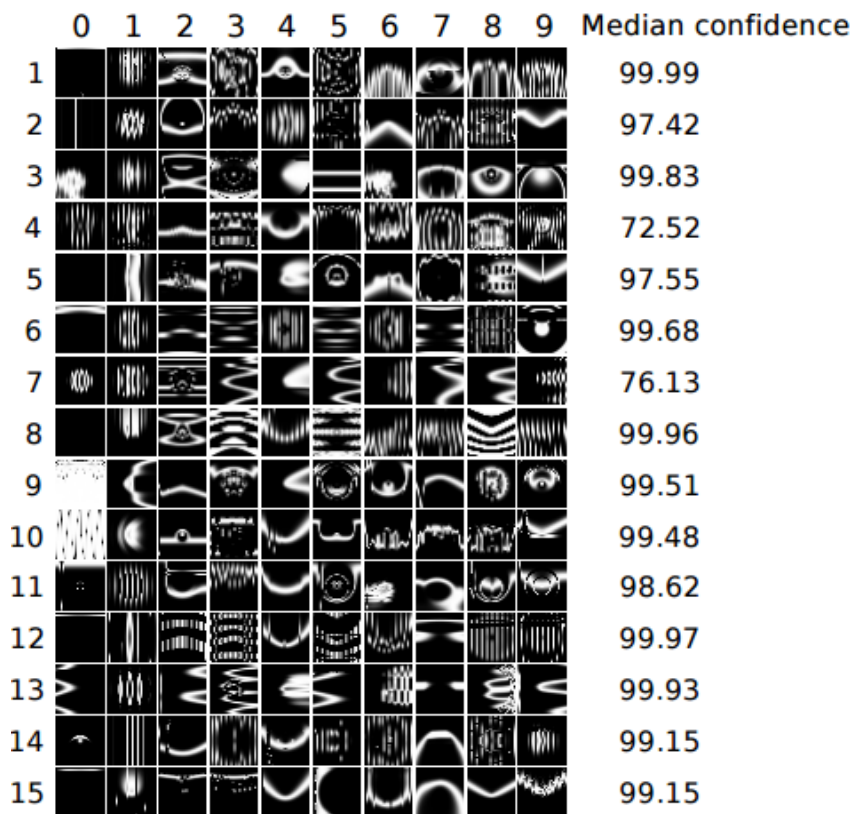


Figure 11. Training MNIST DNN_i with images that fooled MNIST DNN_1 through DNN_{i-1} does not prevent evolution from finding new fooling images for DNN_i . Columns are digits. Rows are DNN_i for $i = 1 \dots 15$. Each row shows the 10 final, evolved images from one randomly selected run (of 30) per iteration. Medians are taken from images from all 30 runs.

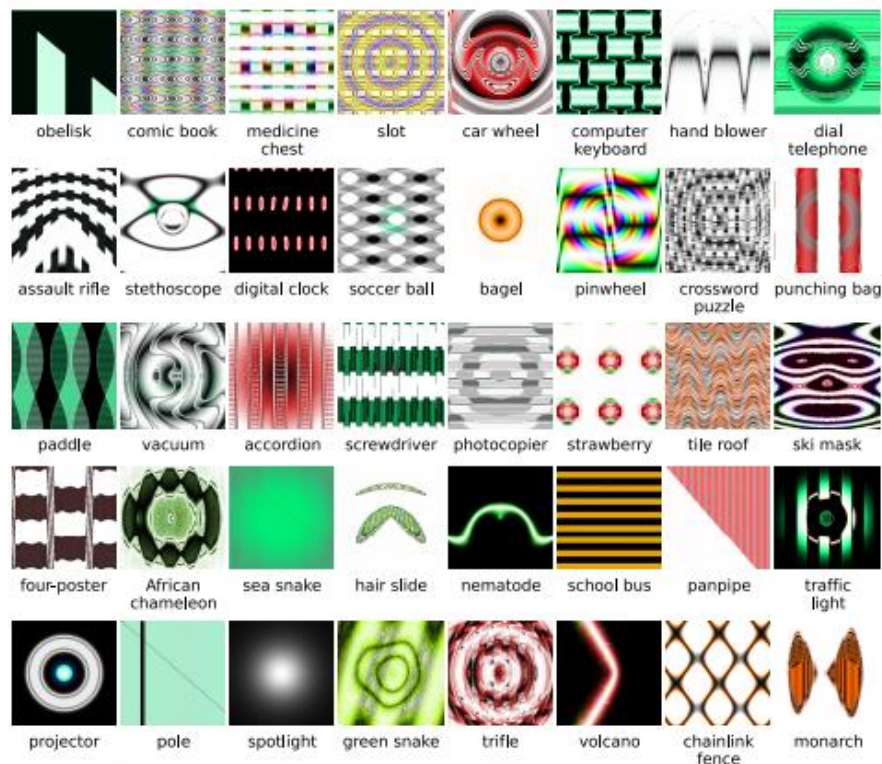


Figure 8. Evolving images to match DNN classes produces a tremendous diversity of images. Shown are images selected to showcase diversity from 5 evolutionary runs. The diversity suggests that the images are non-random, but that instead evolutions producing discriminative features of each target class. The mean DNN confidence scores for these images is 99.12%.

