#### Deep learning statistical models

#### D. Nychka, F. Gerber, M. Bailey and S. Sain\*

Colorado School of Mines and \*Jupiter

September 10, 2021



- The Gaussian CDF
- Climate model output
- Training Deep Nets on Gaussian Processes
- Comparison to Maximum likelihood
- Back to climate
- Hydrological models and the Generalized Pareto

#### Finding the normal CDF

$$erf(x) = F(x) = \int_{-\infty}^{x} \frac{1}{\sqrt{2\pi}} e^{-u^2/2} du$$

and the tail probability

$$erfc(x) = 1 - erf(x)$$

Compute F(.25)

# in R
> pnorm( .25)
[1] 0.59870632568292

### Finding the normal CDF

#### How is this really computed? (try to find pnorm.c within the R source code)

erf 
$$(x) \simeq x R_{lm}(x^2)$$
,  $|x| \le .5$ ,  
erfc  $(x) \simeq e^{-x^2} R_{lm}(x)$ , .46875  $\le x \le 4.0$ ,  
erfc  $(x) \simeq \frac{e^{-x^2}}{x} \left\{ \frac{1}{\sqrt{\pi}} + \frac{1}{x^2} R_{lm}(1/x^2) \right\}$ ,  $x \ge 4$ ,

Cody, W. D. (1993)

 $R_{lm}(x) = P_l(x)/Q_m(x)$  a ratio of a 5<sup>th</sup> and a 4<sup>th</sup> degree polynomial. Accurate to 14 digits!

Approximation works - but it is mysterious!

### Another computation



 $\rightarrow$  MLEs for covariance parameters

The MLE estimator in this case is the function

$$\mathcal{F}(\mathsf{field}) o \hat{ heta}, \hat{\lambda}$$
  
 $\Re^{256} o \Re^2$ 

In the fields package and for a Matern
fit <- spatialProcess( x, y, smoothness=.5)</pre>

Use a deep network (CNN or dense network) to approximate "maximum likelihood estimates".

The  $\mathcal{F}$  !

- Covariance parameters for a Gaussian spatial process.
- Parameters of the extreme value distribution

Approximation works – but it is mysterious!

# Climate model output

Local temperature sensitivity to global temperature First 8 out of 30 *centered* ensemble members



*Goal:* Simulate additional fields efficiently that match the spatial dependence in this 30 member ensemble.

### A Statistical Approach

- $\bullet\,\approx\,13 \text{K}$  grid boxes over N and S America
- Estimate a spatially varying covariance function by fitting stationary covariances to small windows.  $(16 \times 16)$
- Range and variance parameter for each window.
- Encode local estimates into a global model to simulate Gaussian random fields.

This is Ashton Wiens Ph D work.

Train a convolution neural net (CNN) on the "image" to estimate covariance parameters.

Or train a dense neural net on the variogram

We found a speedup by a factor of 100!

Covariance function:

$$k(\mathbf{x}_1, \mathbf{x}_2) = \sigma^2$$
Matern function  $_{\nu}(d)$ 

with  $d = ||\boldsymbol{x}_1 - \boldsymbol{x}_2||/\theta$ 

- Matérn function is a modified Bessel function.
- Smoothness  $\nu$  measures number of mean square derivatives and is equivalent to the polynomial tail behavior of the spectral density.
- θ is the range parameter.
- For  $\nu = .5$ : the classic exponential

$$k(\boldsymbol{x}_1, \boldsymbol{x}_2) = \sigma^2 e^{-||\boldsymbol{x}_1 - \boldsymbol{x}_2||/\theta}$$

#### Observational model

$$Y(\boldsymbol{x}) = g(\boldsymbol{x}) + \boldsymbol{e}(\boldsymbol{x})$$

with g following a Gaussian process with Matérn covariance, variance  $\sigma^2$  and **e** white noise, variance  $\tau^2$ .

We are interested in maximum likelihood estimates for  $\theta$  (range),  $\sigma^2$  and  $\tau^2$ .

- A useful short cut is to focus on  $\theta$  and  $\lambda = \sigma^2/\tau^2$ Can convert  $\lambda$  to equivalent degrees of freedom of the smoother for g.
- Analytical expressions for MLEs of  $\sigma^2$  and  $\tau^2$  based on MLE of  $\lambda$ .

## Finding the MLEs

log Likelihood for covariance parameters.

$$= -\frac{\mathbf{y}^{T}(\sigma^{2}C(\theta) + \tau^{2}I)^{-1}\mathbf{y}}{2} - (1/2)\ln|\sigma^{2}C(\theta) + \tau^{2}I| - (n/2)\ln(\pi)$$

or concentrating onto  $\lambda$  and  $\theta$ 

$$= -(n/2) - (1/2) \ln |\hat{\sigma}^2(\theta, \lambda)(C(\theta) + \lambda)I| - (n/2) \ln(\pi)$$

- $C(\theta)$  correlation matrix for observations.
- No closed form for maximum.
- Often hard to find good starting values for optimization.
- Evaluating inverse and determinant can be time consuming.

#### Examples of training fields



D. Nychka, F. Gerber, M. Bailey and S. Sain

# Variograms of training fields

increasing range ( $\theta$ )  $\downarrow$ 



D. Nychka, F. Gerber, M. Bailey and S. Sain

### Neural net setup

• A neural network is a composition of many simple functions to approximate arbitrary functions.

• It depends on estimating many parameters (weights) based on a large training sample. Training / testing sets

- $\bullet\,$  Input are 16  $\times\,$  16 Gaussian fields or their variograms
- $200 \times 201 = 40200$  values in covariance parameter space
- ullet pprox 1M fields generated for training.
- Tested on 10K fields from 2000 parameter combinations.



### Inputs and artificial neurons



Three neurons with four available inputs: Each neuron creates a linear combination of the inputs followed by a nonlinear transformation.

- Outputs from one layer become the inputs for another layer.
- Linear transformation is estimated (learned) for every neuron
- "Deep Learning" considers many neurons and mulitple layers.



Two hidden layers:

- $4 \times (3+1) = 16$  parameters in second layer
- $4 \times (4+1) = 20$  parameters in third layer
- 4 + 1 = 5 parameters in output layer

See Neural Networks and Deep Learning, Michael Nielsen, for a good introduction

# A convolution version CNN

Designed to work on images

A linear filter is applied to every  $3 \times 3$  block of the input field followed by a nonlinear transformation.



In this case the 5×5 image is reduced to a 2×2 output image.

- These filter results are then filtered again ... and again !
- Many filters (128) are considered. [ 🍝 🍝 ] [ 🍝 🍝 ] ... [ 🍝 🍝 ]
- Filter weights found by training (of course!).

## Basic functional step



$$y_j = \phi(b + w_1x_1 + \ldots + w_9x_9)$$

- $\{x_1, x_2, \dots, x_9\}$  pixel values from  $3 \times 3$  "image" (aka a tensor)
- $y_j$  the  $j^{th}$  pixel value for "image" (tensor) at next layer
- w the weights and b offset to be estimated/optimized
- $\phi(u) = u_+$  Rectified linear unit

### Net architecture

#### Form for finding MLEs based on an image

Layer	Operation	Size
1	Input	$16 \times 16$ image
2	2D convolution	128 7 $\times$ 7 filtered images
3	2D convolution	128 3 $ imes$ 3 filtered images
4	2D convolution	128 outputs
5	dense	500 outputs
6	dense	2 outputs ( $\theta$ and $\lambda$ )

#### 636K parameters!

E.g at layer 5 there are 500 neurons each with 128 inputs, (128  $\pm$ 1)X 500 parameters

• Simpler 2 layer dense network used for finding MLEs based on a variogram

#### Parameter estimates on 10K test samples

• Recall that the MLEs are "optimal" for large n based on Cramer-Rao lower bound.



Training:

- CNN , NN, and MLE estimates tend to track the red lines (truth)
- CNN and NN overall has comparable accuracy to the MLEs
- Potential tradeoff between bias in CNN estimates and variance in MLEs

# Climate model emulation

#### Estimated log Variance and range



Training the network 2-6 hours using cloud computing (but this can be shortened for a slightly less accuracy) For the climate model output

- 12,769 windowed estimates
- 10s of seconds using the CNN , 2 minutes using the NN variogram
- $\approx$  1.5 hours using standard MLE fitting in R (fields package)

In general we find a factor of 100 or more speedup.

Part of this may be due to the efficiency of the tensor flow libraries and low level coding.

#### More about irregular spatial data

Build a deep net based on the variogram statistic.

Can the variogram serve as an approximate "sufficient statistic" for a stationary covariance function?

log parameter estimates for  $\lambda$  and the range

log lambda



20

Range

- How to adjust variogram NN for different bin counts.?
- Train on a SAR model (LatticeKrig, SPDE) directly instead of Matérn .
- Train for the likelihood *surface* instead of just the estimates.

# Hydrologic models and extremes

#### Steve Sain, *Jupiter*, Maggie Bailey *Mines* Simulating flooding from extreme precipitation



Aerial view of study region and grid cell



Grid cell high water response



Simulated high water field for 60 cm



### Hydrologic models and extremes

Goal is to estimate the shape  $(\boldsymbol{\xi})$  and scale parameters  $(\boldsymbol{\sigma})$  of a generalized Pareto distribution at many (millions) of pixels.

$$f(x) = \frac{(1/\sigma)(1+\xi\frac{(x-\mu)}{\sigma})^{-\frac{\xi+1}{\xi}}}{(1/\sigma)e^{-\frac{(x-\mu)}{\sigma}}} \quad \xi \neq 0$$

and

 $\begin{array}{ll} x \ge \mu & \boldsymbol{\xi} \ge 0 \\ \mu \le x \le \mu - (\boldsymbol{\sigma}/\boldsymbol{\xi}) & \boldsymbol{\xi} \le 0 \end{array}$ 

• Train a neural net on finely binned histograms to obtain estimates comparable to the MLEs.

• Potential speedup will allow for data analysis on high resolution model output.

D. Nychka, F. Gerber, M. Bailey and S. Sain'

#### Prelimnary results

Sample size of 2500, 8000/2000 cases for training/testing,  $\mu=1$  i.e. threshold is fixed.



D. Nychka, F. Gerber, M. Bailey and S. Sain

Fall-19

- Exploit fast simulation of statistical samples to train a neural net.
- Neural nets can accurately reproduce statistical computations but evaluate much more quickly.
- Training and test samples provide a rigorous way to insure neural net approximations.

# Thank you



# An example of a Keras/R specification

 $\bullet~$  NxN image initial image, and 3X3 filters interspersed with max pooling reductions.

• Final step takes the last "image" and feeds to a dense neural network with 2 outputs

# Keras model summary

> modelMatern11 Model Model: "sequential"		
Layer (type)	Output Shape	Param #
conv2d (Conv2D)	(None, 14, 14, 32)	320
<pre>max_pooling2d (MaxPooling2D)</pre>	(None, 7, 7, 32)	0
conv2d_1 (Conv2D)	(None, 5, 5, 32)	9248
<pre>max_pooling2d_1 (MaxPooling2D)</pre>	(None, 2, 2, 32)	0
flatten (Flatten)	(None, 128)	0
dense (Dense)	(None, 2)	258
Total params: 9,826 Trainable params: 9,826 Non-trainable params: 0		