



Deep networks for system identification: a survey

Gianluigi Pillonetto♣, Aleksandr Aravkin◇, **Daniel Gedon**♡, Lennart Ljung♠,
Antônio H. Ribeiro♡, Thomas B. Schön♡

♣University of Padova, Italy

◇University of Washington, USA

♡Uppsala University, Sweden

♠Linköping University, Sweden

ERNSI Workshop 2023

Stockholm, September 26, 2023

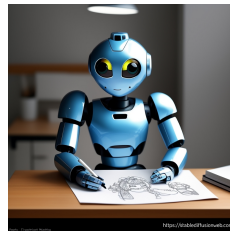
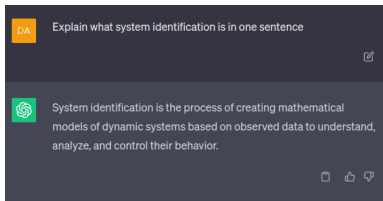
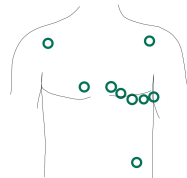
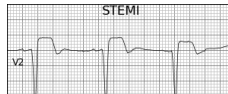
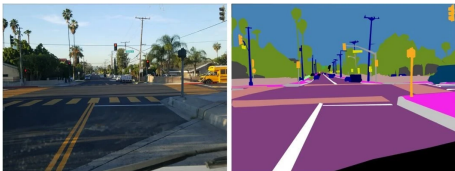
arXiv link:



System identification with long history



Deep neural networks with recent success



→ Innovate system identification with power of deep neural networks

1. Modeling of dynamical systems
2. Deep neural network architectures
3. Optimization
4. Deep kernel-based learning
5. Theoretical development
6. Applications
7. Conclusion

Three main players:

1. Family of parameterized models

$$Z = \{x(t), y(t)\}_{t=1}^{\#train}$$
$$g_{\theta} : Z(t) \mapsto \hat{y}(t+1), \quad \theta \in D_{\theta}$$

2. Parameter estimation method

$$\hat{\theta} = \arg \min_{\theta \in D_{\theta}} \mathcal{L}_N(\theta, Z_e)$$

3. Validation process

- residual analysis
- cross-validation

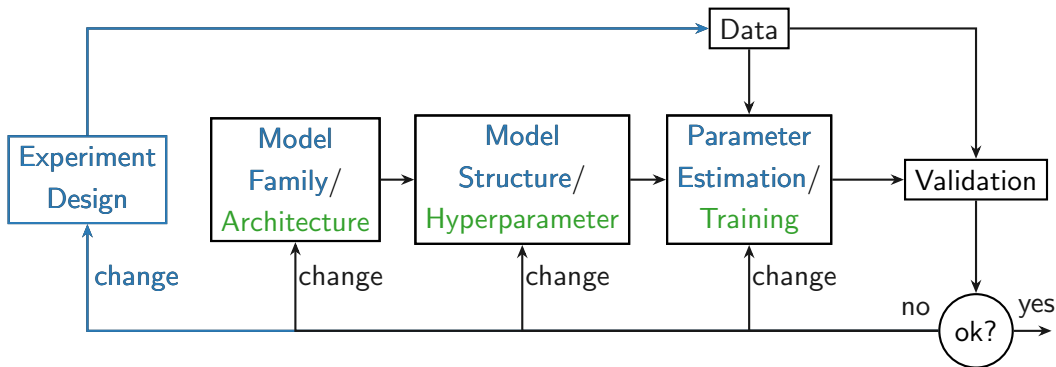
$$\#features = \dim \theta$$

$$\mathcal{L}_{emp} = \mathcal{L}(\hat{\theta}, Z_e)$$

overfitting $\mathcal{L}_{emp} = 0$ typically for $\#features = \#train$.

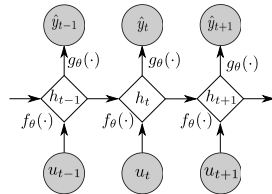
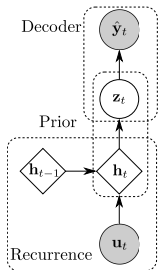
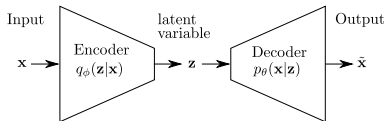
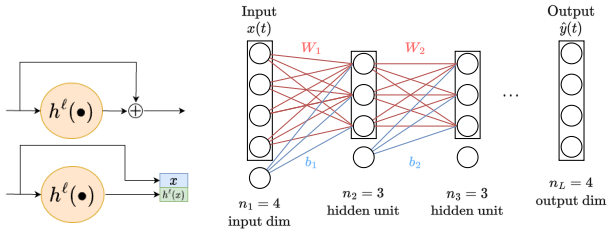
Modeling procedure:

System identification vs deep learning



1. Modeling of dynamical systems
2. Deep neural network architectures
3. Optimization
4. Deep kernel-based learning
5. Theoretical development
6. Applications
7. Conclusion

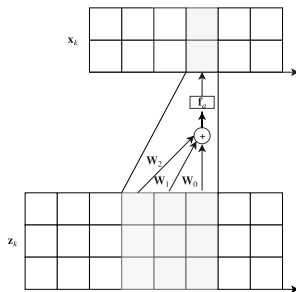
- Fully-connected networks
- Skip and direct connections
- **Convolutional networks**
- Recurrent neural networks
- Latent variable models
 - Autoencoder
 - Variational autoencoder
 - Deep state-space models
- **Energy-based models**



Convolutional networks

Basic building block: convolutional layer

$$w(t) * z(t) = \sum_{j=0}^{k-1} w(j)^\top z(t-j)$$



Not just one filter but many: $W = \{w^1, \dots, w^b\}$.

Then, i th output: $x^i(t) = w^i(t) * z(t)$ for $i = 1, \dots, b$

Formulating regression problems

Find predictive distribution $p(y(t)|x(t))$.

Example: NARX model

$$y(t) = f_{\theta}(x(t)) + e(t), \quad \text{with} \quad e(t) \sim \mathcal{N}(0, \sigma^2)$$

→ Implicit assumption: $p(y(t)|x(t))$ is Gaussian → neural network models the mean.

Energy-based models

$$p_{\theta}(y(t) | x(t)) = \frac{e^{g_{\theta}(y(t), x(t))}}{Z_{\theta}(x(t))} \quad \text{with} \quad Z_{\theta}(x(t)) = \int e^{g_{\theta}(z, x(t))} dz$$

- Neural network mapping $g_{\theta} : (y(t), x(t)) \mapsto \mathbb{R}$
- Generalize implicit Gaussian assumption

→ asymmetric, heavy-tailed, multimodal, ... distributions possible

System identification:

$$\min_{\theta} \sum_{t=1}^{\#train} \mathcal{L}(y(t), f_{\theta}(z(t)))$$

Deep learning:

$$\min_{\theta_1, \dots, \theta_L} \sum_{t=1}^{\#train} \mathcal{L}(y(t), f_{\theta_L}^L \circ f_{\theta_{L-1}}^{L-1} \circ \dots \circ f_{\theta_1}^1(z(t)))$$

Optimization: Newton's method $\mathcal{O}(\#train \#param^2 + \#param^3)$ \downarrow

→ first-order methods

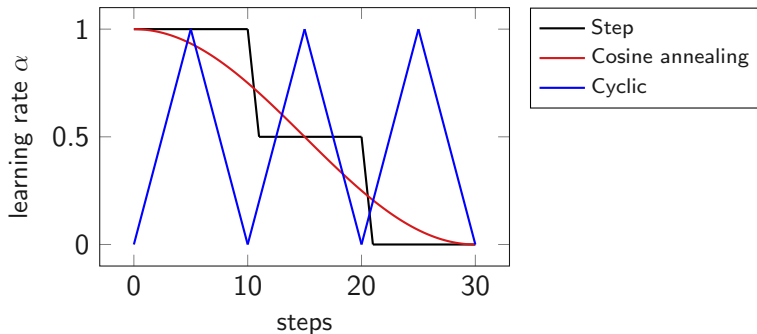
- Large $\dim(\theta)$, nested structure → gradient w.r.t. each layer + chain rule
→ Backpropagation
- Large datasets → stochastic methods

Gradient decent optimization:

$$\theta^{i+1} = \theta^i - \alpha \nabla V(\theta^i) \quad \text{with } \alpha \text{ as learning rate}$$

Stochastic gradient descent with fixed α does not converge ζ

Solution: Learning rate scheduler \rightarrow reduce α to zero



1. Modeling of dynamical systems
2. Deep neural network architectures
3. Optimization
4. Deep kernel-based learning
5. Theoretical development
6. Applications
7. Conclusion

Kernels for modeling dynamical systems

- Linear kernel

$$K(x_i, x_j) = x_i^\top P x_j \quad \text{with positive semidefinite } P$$

induces linear functions $f(x) = \theta^\top x$ \rightarrow FIR models

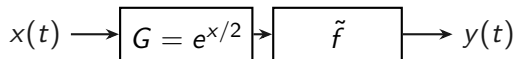
- Linear kernel with $P_{ij} = \varphi^{\max(i,j)}$ with $0 \leq \varphi < 1$ \rightarrow stable spline/TC kernel

- Gaussian kernel $K(x_i, x_j) = \exp\left(-\frac{\|x_i - x_j\|^2}{\rho}\right)$ with $\rho > 0$ \rightarrow NFIR models

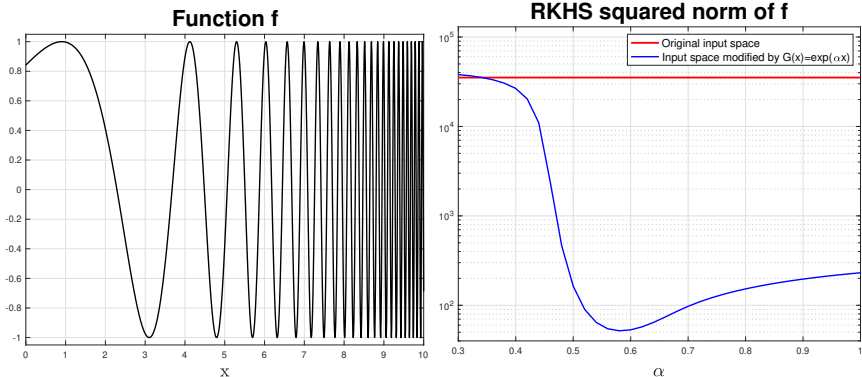
Choice of kernel \rightarrow encode high level assumptions

Example: $f = \sin(e^{x/2}) \rightarrow$ complicated frequency content

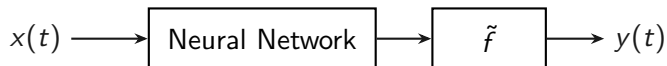
- Gaussian kernel: high RKHS norm \rightarrow biased estimator
- Idea: transform data $f = \tilde{f} \circ G$



Choose $G = e^{x/2} \rightarrow \tilde{f} = \sin(x)$ with single frequency



Consider idea: $f = \tilde{f} \circ G$



→ manifold Gaussian process with

$$K(x_i, x_j) := \tilde{K}(\tilde{x}_i, \tilde{x}_j) = \tilde{K}(G(x_i), G(x_j))$$

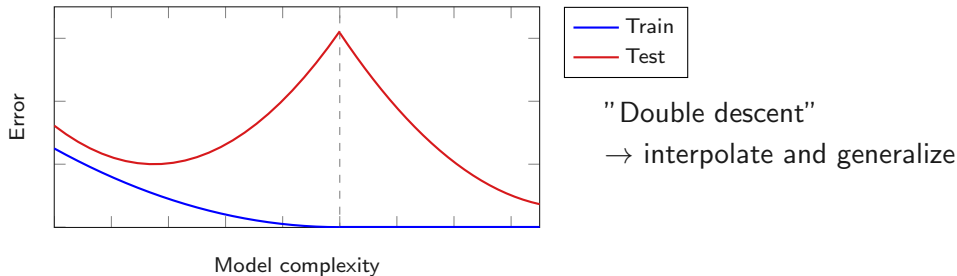
Previously: Gaussian kernel K with one scale parameter $\rho > 0$

Now: Manifold Gaussian kernel K with many parameters $\eta = [\rho, \theta]$

→ Optimize by marginal likelihood of joint density $p(Y, f|\eta)$

Why are deep models so successful?

- 2-layer ConvNet on MNIST: 1.2m parameters vs 60k data points
- AlexNet on ImageNet: 62.3m parameters vs 1.2m data points
- ...



"Double descent"

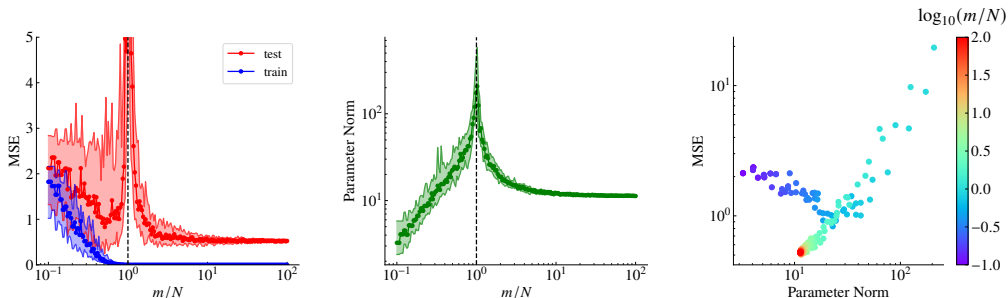
→ interpolate and generalize

Theoretical development:

1. **interplay of overparameterization and generalization**
2. simplification of non-convex optimization problem

System identification example:

- NARX model: $\hat{y}(t) = \sum_{i=1}^{\#features} \theta_i \phi_i(x(t))$
- Data from: $y(t) = f(x(t)) + v(t)$
- $\#train = 100$ samples
- 1-step ahead prediction



- Nonlinear transformation $\phi(x)$, input to feature space

$$\phi : \mathbb{R}^{\#inputs} \mapsto \mathbb{R}^{\#features}$$

- Linear model:

$$\hat{y} = \hat{\theta}^\top \phi(x)$$

- Estimation procedure:

$$\min_{\theta} \sum_{i=1}^{\#train} (y_i - \hat{\theta}^\top \phi(x_i))^2$$

- Optimization procedure: Gradient descent starting from zero

$$\theta^{i+1} = \theta^i - \alpha \nabla V(\theta^i)$$

Solutions of a linear system

$$X\theta = y$$

Three scenarios:

1. no solution if *#features* < *#train*
2. one unique solution if *#features* = *#train*
3. multiple solution if *#features* > *#train*

Gradient descent:

$$\min_{\theta} \|\theta\|_2 \quad \text{subject to} \quad X\theta = y$$

converges to the minimum-norm solution

→ Implicit regularization of gradient descent

Implicit Regularization

Gradient descent step: $\theta^{i+1} = \theta^i - \alpha \nabla V(\theta^i)$

→ does not follow continuous gradient flow

Gradient descent follows more closely

$$\dot{\theta} = -\nabla \tilde{V}(\theta)$$

with modified cost

$$\tilde{V}(\theta) = V(\theta) + \lambda R(\theta)$$

$$\lambda = \frac{\alpha \#features}{4}, \quad R(\theta) = \frac{1}{\#features} \sum_{j=1}^{\#features} (\nabla_j V(\theta))^2$$

→ gradient descent penalizes directions j with large cost $V(\theta)$

2. Simplification of non-convex optimization problem

Setup:

- wide neural network with large $\theta \in \mathbb{R}^{\#\text{features}}$
- each update changes θ just by small amount

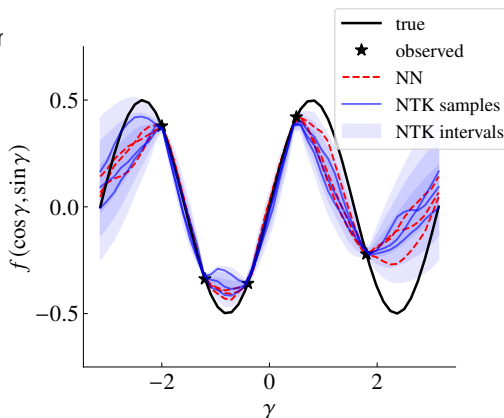
→ linearize model around θ_0

$$f_{\theta}(x) \approx f_{\theta_0}(x) + \nabla f_{\theta_0}(x)^{\top} (\theta - \theta_0)$$

Neural tangent kernel

$$K(x, z; \theta_0) = \nabla f_{\theta_0}(x)^{\top} \nabla f_{\theta_0}(z)$$

→ convex optimization problem

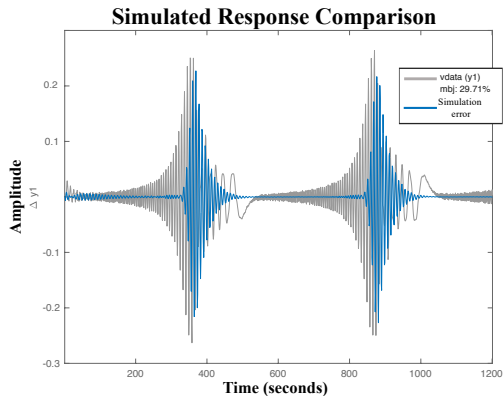


1. Modeling of dynamical systems
2. Deep neural network architectures
3. Optimization
4. Deep kernel-based learning
5. Theoretical development
- 6. Applications**
7. Conclusion

Matlab example: forced duffing oscillator (silverbox benchmark)

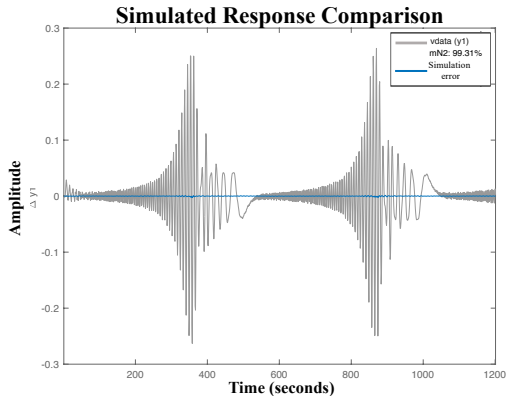
Linear Box-Jenkins type model

→ Fit is 29.7%



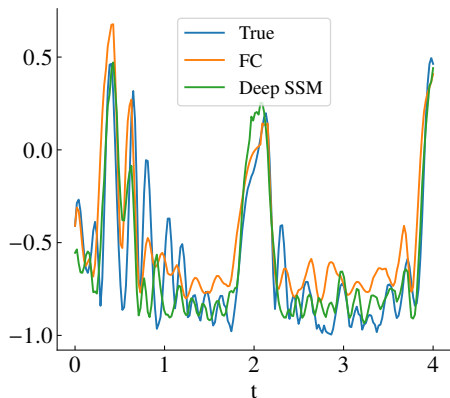
Cascaded feedforward network

→ Fit is 99.2%



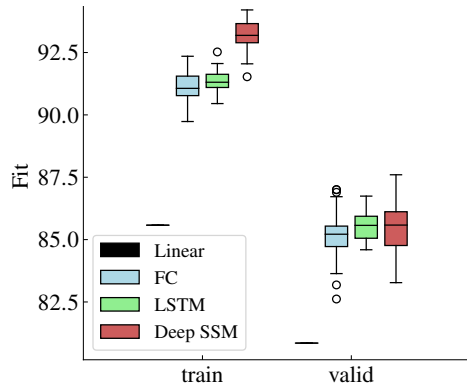
Pytorch example: Coupled electronic drives benchmark

- Baseline: linear ARX model
- Feedforward model
- LSTM
- Deep state-space model



Good fit of deep models despite $\#train = 300$

- $\dim(\theta_{FF}) = 184,200$
- $\dim(\theta_{LSTM}) = 169,801$
- $\dim(\theta_{DSSM}) = 111,902$



Essential for using neural networks:

- many parameters → overparameterization
- many layers → deep architectures

Open problems:

- Successful architectures:
 - Attention models and transformers
 - Flow-based models
 - Generative adversarial models (GANs) and diffusion models
 - Graph neural networks
- Robustness issues
- Theoretical development
- ...

Thank you!

Daniel Gedon, Uppsala University

E-mail: `daniel.gedon@it.uu.se`

Web: `dgedon.github.io`

Twitter: @danigedon

arXiv link:

