

On Deep Learning for Low-Dimensional Representations

Daniel Gedon

Uppsala University

Disputation

Uppsala, June 14, 2024

Deep Learning for Low-Dimensional Representation

Fundamental

Paper I

Paper II

Paper III

Applications

Paper VI

Paper V

Paper IV

- I. Uncertainty Estimation with Recursive Feature Machines, *UAI 2024*
- II. Invertible Kernel PCA with Random Fourier Features, *IEEE Signal Processing Letters 2023*
- III. No Double Descent in Principal Component Regression: A High-Dimensional Analysis, *ICML 2024*
- IV. Deep State Space Model for Nonlinear System Identification, *SYSID 2021*
- V. First Steps Towards Self-Supervised Pretraining of the 12-Lead ECG, *CinC 2021*
- VI. Development and Validation of Deep Learning ECG-Based Prediction of Myocardial Infarction in Emergency Department Patients, *Scientific Reports 2022*

Part I. Problem formulation

Part II. Fundamental ML

Part III. Applications

1. Machine Learning

- What is it?
- When do we need it?

2. Low-dimensional representations

- What are representations?
- Why are they low-dimensional?



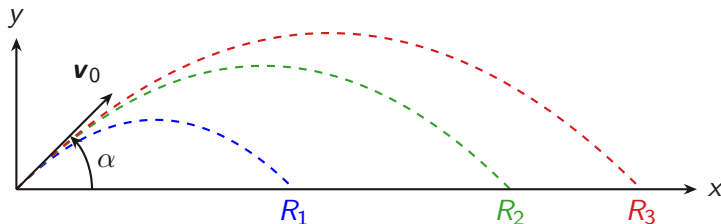
Part II: Machine learning for low-dimensional data structure

Part III: Low-dimensional representations for deep learning application

Scenario:

Real-world observations \rightarrow model \rightarrow model \rightarrow explanation / prediction

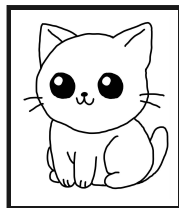
Example 1: Trajectory motion



Model $f : \mathbf{v}_0 \mapsto R$ given α, g :
$$R = f(\mathbf{v}_0) = \frac{v_0^2 \sin 2\alpha}{g}$$

\rightarrow Occam's razor / Principle of parsimony

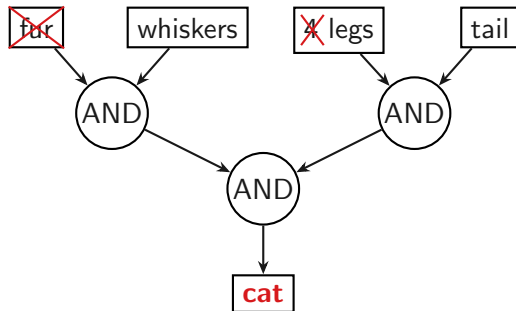
Example 2: Image classification



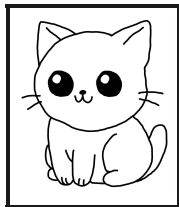
Structure/
Features?

Model

cat



⇒ Fails in many real-world scenarios



which
features?

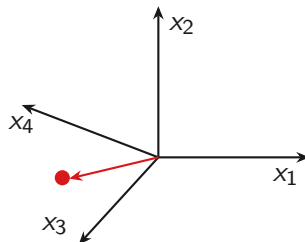
Electrocardiogram (ECG)

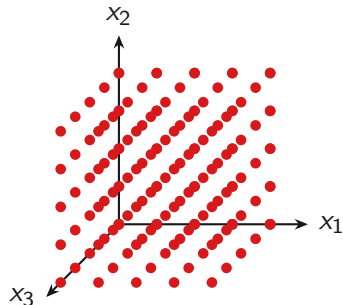


Input → **Model** → Prediction

-0.34
0.78
0.21
-0.09
0.89
-0.17

high-dimensional
features

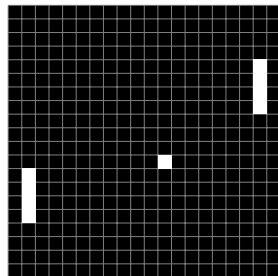




1-dim space takes volume: $V = s^1$ 2-dim
space takes volume: $V = s^2$ 3-dim space
takes volume: $V = s^3$ p -dim space takes
volume: $V = s^p$

→ exponentially more samples needed!

Example: Atari game Pong



Problem: 20×20 pixels \rightarrow 400 dimensions

Solution: only need x/y-position of ball

$$\mathbf{z} = [11, 8]^\top \text{ with } \mathcal{Z} \in \mathbb{R}^2$$

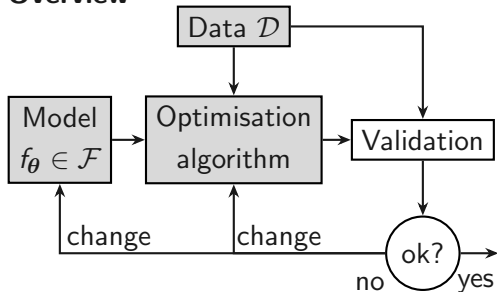
→ Obtain useful low-dim. representation

Part I. Problem formulation

Part II. Fundamental ML

Part III. Applications

Overview



$$\mathcal{D} = (\mathbf{x}, y) \mid \mathbf{x} \in \mathcal{X}, y \in \mathcal{Y}$$

Parameterised model $f_{\theta} : \mathcal{X} \mapsto \mathcal{Y}$

L layers $f_{\theta} = f_{\theta_L}^L \circ \dots \circ f_{\theta_2}^2 \circ f_{\theta_1}^1$

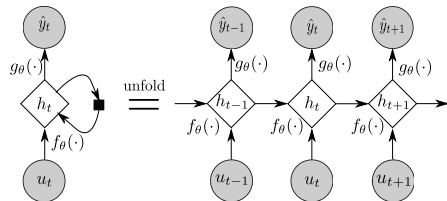
Empirical risk minimisation

$$\arg \min_{\theta} \sum_{i=1}^n \ell(f_{\theta}(\mathbf{x}_i), y_i)$$

Architectures for low-dim. representations

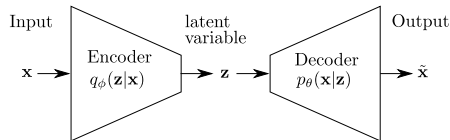
Recurrent neural network:

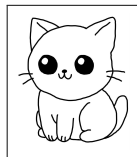
→ past information in state h



Variational autoencoder:

→ bottleneck layer





Similar data
↓
Same prediction

Define kernel function $k(\mathbf{x}, \mathbf{x}') \mapsto \mathbb{R}$

Prediction: $f(\mathbf{x}) = \sum_{i=1}^n \alpha_i k(\mathbf{x}, \mathbf{x}_i)$

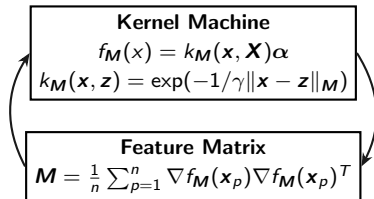
Which kernel to choose?

- Linear: $k_{lin}(\mathbf{x}, \mathbf{x}') := \langle \mathbf{x}, \mathbf{x}' \rangle$
- RBF: $k_{rbf}(\mathbf{x}, \mathbf{x}') := \exp\left(-\frac{\|\mathbf{x} - \mathbf{x}'\|^2}{2\ell^2}\right)$
- Mahalanobis:

$$k_M(\mathbf{x}, \mathbf{x}') := \exp\left(-\frac{\|\mathbf{x} - \mathbf{x}'\|_M}{\ell}\right)$$

$$\|\mathbf{x} - \mathbf{x}'\|_M^2 = (\mathbf{x} - \mathbf{x}')^\top \mathbf{M} (\mathbf{x} - \mathbf{x}')$$

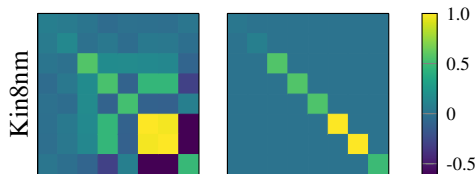
Recursive Feature Machines



Why Average Gradient Outer Product?

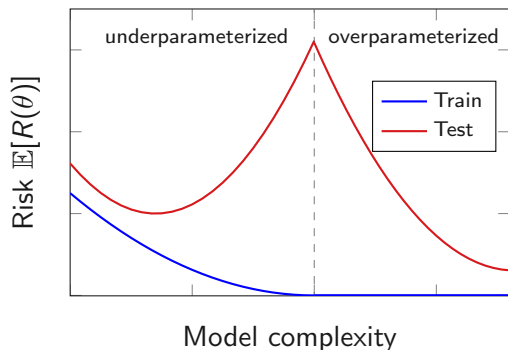
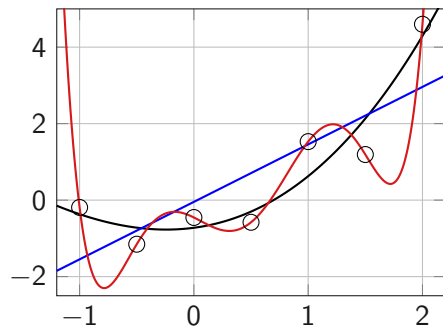
$$\mathbf{W}_l^\top \mathbf{W}_l \propto \frac{1}{n} \sum_{i=1}^n \nabla f^l(\tilde{\mathbf{x}}_i) \nabla f^l(\tilde{\mathbf{x}}_i)^\top$$

Why low-dimensional?



How well do machine learning models perform? → Generalisation risk

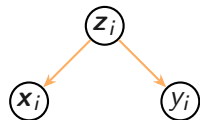
Example: regression



Model complexity:

- Parameters $p < \text{number of training data points } n$: underparameterized
- Parameters $p > \text{number of training data points } n$: overparameterized

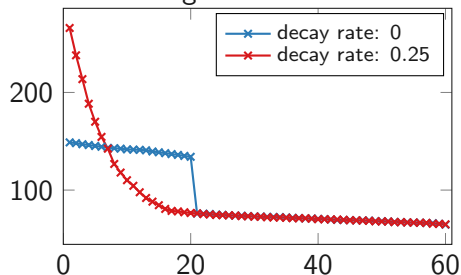
Data generator



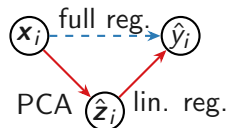
$$\mathbf{x}_i = \mathbf{W} \mathbf{r}_w \mathbf{z}_i + \mathbf{e}_i$$

$$y_i = \boldsymbol{\theta}^\top \mathbf{z}_i + \varepsilon_i$$

Eigenvalues of \mathbf{x}



Model

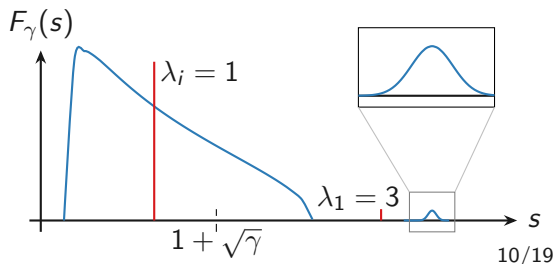


$$\text{SVD: } \mathbf{X} \approx \hat{\mathbf{U}} \hat{\boldsymbol{\Sigma}} \hat{\mathbf{V}}^\top,$$

$$\text{PCA: } \hat{\mathbf{z}}_i = \hat{\mathbf{V}}^\top \mathbf{x}_i,$$

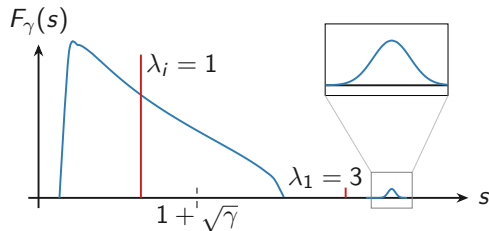
$$\text{lin. reg. } \hat{y}_i = \hat{\boldsymbol{\theta}}^\top \hat{\mathbf{z}}_i.$$

High-dim. random matrix theory



High dimensional analysis based on random matrix theory

Eigenvalue shift:



Risk:

$$R(\hat{\theta}) = \mathbb{E}_{(\mathbf{x}_0, y_0) \sim \mathcal{D}} [(y_0 - \hat{y}_0(\mathbf{x}_0))^2]$$

Asymptotic PCR risk:

$$\mathbb{E}_{\nu} [R(\hat{\theta})] \rightarrow \text{Bias}_{\gamma}(\hat{\theta})^2 + \text{Var}_{\gamma}(\hat{\theta}) + \sigma^2$$

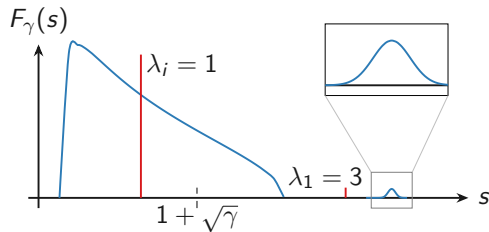
With the following:

$$\text{Bias}_{\gamma}(\hat{\theta})^2 = \bar{\beta}^{\top} (\Lambda_d - \Lambda_d \mathbf{P} - \mathbf{P} \Lambda_d + \mathbf{P} + \mathbf{P} r_w^2 \mathbf{C}_z \mathbf{P}) \bar{\beta}$$

$$\text{Var}_{\gamma}(\hat{\theta}) = \frac{\sigma^2}{n} \left(\text{Tr} \left[(\mathbf{P} r_w^2 \mathbf{C}_z + \mathbf{I}) \frac{1}{\mu(\Lambda, \gamma)} \right] + (p - d) \int_c^{(1+\sqrt{\gamma})^2} \frac{1}{s} dF_{\gamma}(s) \right)$$

High dimensional analysis based on random matrix theory

Eigenvalue shift:

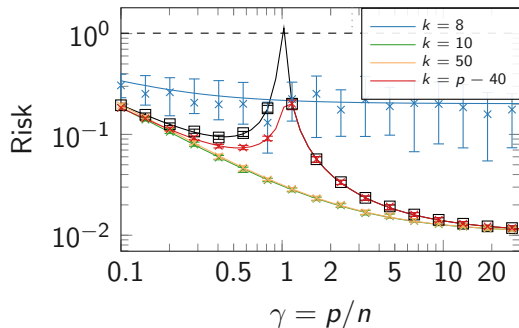


Risk:

$$R(\hat{\theta}) = \mathbb{E}_{(\mathbf{x}_0, y_0) \sim \mathcal{D}} [(y_0 - \hat{y}_0(\mathbf{x}_0))^2]$$

Asymptotic PCR risk:

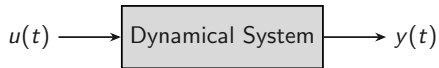
$$\mathbb{E}_\nu [R(\hat{\theta})] \rightarrow \text{Bias}_\gamma(\hat{\theta})^2 + \text{Var}_\gamma(\hat{\theta}) + \sigma^2$$



Part I. Problem formulation

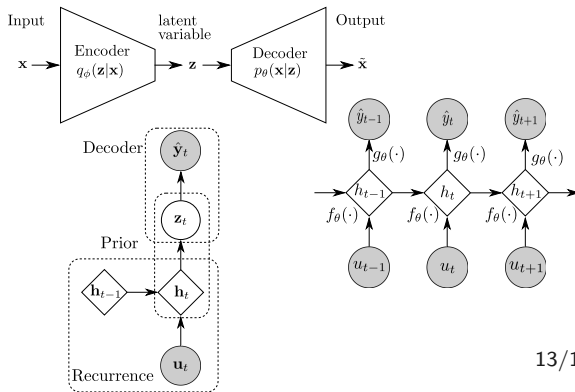
Part II. Fundamental ML

Part III. Applications



Modelling of systems:

- Fully-connected networks
- (Temporal) Convolutional networks
- Recurrent neural networks
- Latent variable models
 - Autoencoder
 - Variational autoencoder
 - **Deep state-space models**
- Energy-based models



Temporal extension of VAE \rightarrow time-varying prior for latent \mathbf{z}

Model definition:

- Joint model distribution

$$p_{\theta}(y_{1:T}, \mathbf{z}_{1:T} | \mathbf{x}_{1:T}, \mathbf{z}_0) = \prod_{t=1}^T p_{\theta}(y_t | \mathbf{z}_t) p_{\theta}(\mathbf{z}_t | \mathbf{z}_{t-1}, \mathbf{x}_t)$$

- Prior and decoder

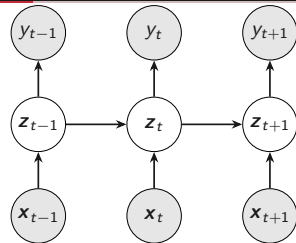
$$p_{\theta}(\mathbf{z}_t | \mathbf{z}_{t-1}, \mathbf{x}_t) = \mathcal{N}(\mathbf{z}_t | \mu_t^{\text{prior}}, \sigma_t^{\text{prior}} \mathbf{I})$$

$$p_{\theta}(y_t | \mathbf{z}_t) = \mathcal{N}(y_t | \mu_t^{\text{dec}}, \sigma_t^{\text{dec}} \mathbf{I})$$

- For training: variational encoder distribution

$$q_{\phi}(\mathbf{z}_{1:T} | y_{1:T}, \mathbf{x}_{1:T}, \mathbf{z}_0) = \prod_{t=1}^T q_{\phi}(\mathbf{z}_t | \mathbf{z}_{t-1}, y_{t:T}, \mathbf{x}_{t:T})$$

\rightarrow Requires nonlinear smoothing step for training



Temporal extension of VAE \rightarrow time-varying prior for latent \mathbf{z}

Model definition:

- Joint model distribution

$$p_{\theta}(y_{1:T}, \mathbf{z}_{1:T}, \mathbf{h}_{1:T} | \mathbf{x}_{1:T}, \mathbf{h}_0) = \prod_{t=1}^T p_{\theta}(y_t | \mathbf{z}_t) p_{\theta}(\mathbf{z}_t | \mathbf{h}_t) p_{\theta}(\mathbf{h}_t | \mathbf{h}_{t-1}, \mathbf{x}_t)$$

- Prior and decoder

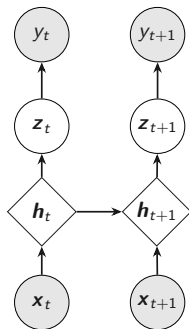
$$p_{\theta}(\mathbf{z}_t | \mathbf{z}_{t-1}, \mathbf{x}_t) = \mathcal{N}(\mathbf{z}_t | \mu_t^{\text{prior}}, \sigma_t^{\text{prior}} \mathbf{I})$$

$$p_{\theta}(y_t | \mathbf{z}_t) = \mathcal{N}(y_t | \mu_t^{\text{dec}}, \sigma_t^{\text{dec}} \mathbf{I})$$

- For training: variational encoder distribution

$$q_{\phi}(\mathbf{z}_{1:T}, \mathbf{h}_{1:T} | y_{1:T}, \mathbf{x}_{1:T}, \mathbf{h}_0) = \prod_{t=1}^T q_{\phi}(\mathbf{z}_t | y_t, \mathbf{h}_t) p_{\theta}(\mathbf{h}_t | \mathbf{h}_{t-1}, \mathbf{x}_t)$$

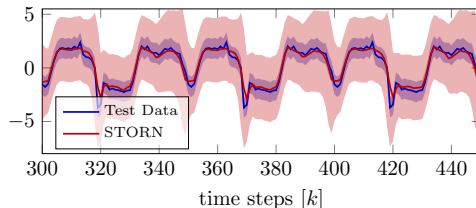
\rightarrow Solution: De-couple state \mathbf{z}_t with hidden state \mathbf{h}_t



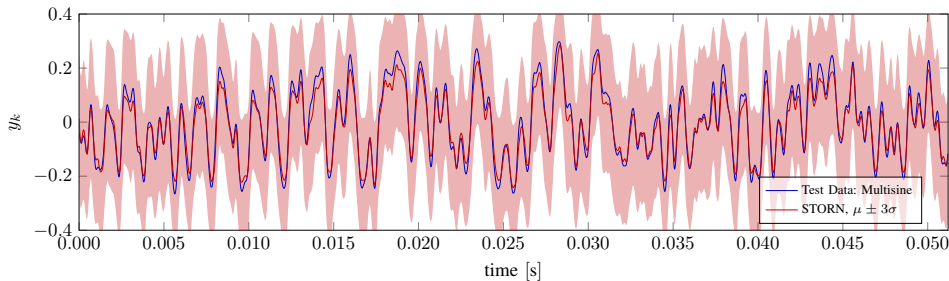
$$\begin{bmatrix} x_{k+1}^{(1)} \\ x_{k+1}^{(2)} \end{bmatrix} = \begin{bmatrix} \left(\frac{x_k^{(1)}}{1+(x_k^{(1)})^2} + 1 \right) \sin(x_k^{(2)}) \\ x_k^{(2)} \cos(x_k^{(2)}) + x_k^{(1)} \exp\left(-\frac{(x_k^{(1)})^2 + (x_k^{(2)})^2}{8}\right) + \dots \\ \dots + \frac{u_k^3}{1+u_k^2 + 0.5 \cos(x_k^{(1)} + x_k^{(2)})} \end{bmatrix}$$

$$y_k = \frac{x_k^{(1)}}{1 + 0.5 \sin(x_k^{(2)})} + \frac{x_k^{(2)}}{1 + 0.5 \sin(x_k^{(1)})} + e_k.$$

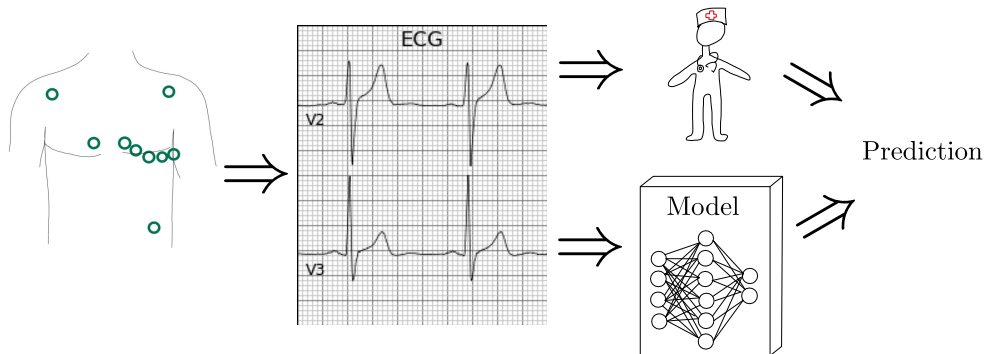
Narendra-Li Benchmark



Wiener-Hammerstein Benchmark: Multisine Test Data



Cardiovascular diseases: ≈ 18 million deaths in 2019 (32%)



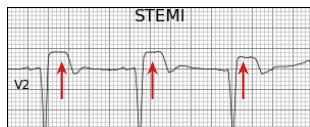
ECG is a major diagnostic tool:

- low-cost, safe, quick, non-invasive
- Can detect arrhythmias, myocardial infarctions, cardiomyopathy, ...

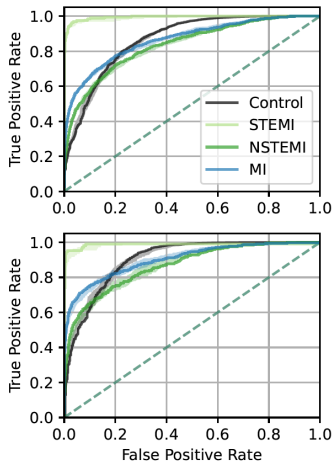
Self-supervised (**V**), NSTEMI (**VI**), atrial fibrillation, electrolyte regression, Chagas, multiclass challenge.

Myocardial Infarctions

- ST-elevation MI
→ ECG
- non-ST-elevation MI
→ blood test



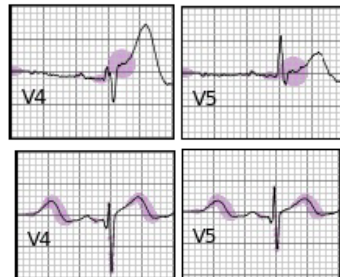
ROC curve



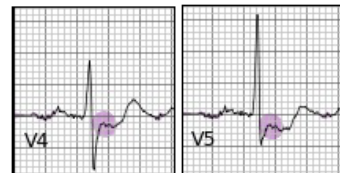
Top: temporal split.
Bottom: random split.

Grad-CAM

STEMI

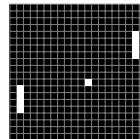


NSTEMI

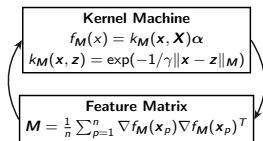




Data variations \rightarrow machine learning

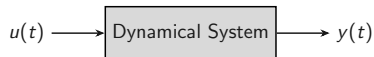
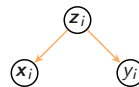
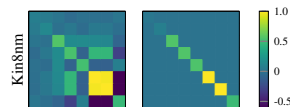
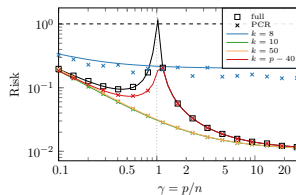


Real world data: low-dimensional



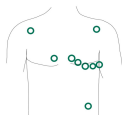
Recursive Feature Machine
 \rightarrow uncertainty quantification

High-dimensional risk analysis



Deep state-space model

Electrocardiogram modelling



Thank you!

Daniel Gedon, Uppsala University

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

Web: `dgedon.github.io`

Twitter / X: `@danigedon`