

Model driven machine learning and some applications

Tamás Dózsa ^{1,2}

¹Hungarian Research Network
Institute for Computer Science and Control,
Systems and Control Lab

²Department of Numerical Analysis,
Faculty of Informatics,
Eötvös Loránd University

Budapest, 2023.11.21

Project no. C1748701 has been implemented with the support provided by the Ministry of Culture and Innovation of Hungary from the National Research, Development and Innovation Fund, financed under the NVKDP-2021 funding scheme.

Outline

- 1** Variable Projection Operators
- 2** An example in system identification
- 3** Data driven modeling
- 4** VP-NET: Model driven Deep Learning
- 5** Road abnormality recognition with VP-NET
- 6** Variable Projection Support Vector Machines
- 7** Conclusion

Adaptive signal models

Variable subspaces in Hilbert spaces

- Let $(\mathcal{H}, \langle \cdot, \cdot \rangle)$ be a Hilbert space and $f \in \mathcal{H}$ arbitrary.
- Suppose $\varphi_k^\eta \in \mathcal{H}$ is a complete and orthogonal basis ($k \in \mathbb{N}$).
- $\eta \in \Gamma \subset \mathbb{R}^L$ ($L \in \mathbb{N}$) determines the system φ_k^η .
- For $m \in \mathbb{N}$, $\varphi_0^\eta, \dots, \varphi_{m-1}^\eta$ span the m -dimensional closed subspace $\mathcal{U} \subset \mathcal{H}$.
- For a fixed η , $\exists! \hat{f} = \sum_{k=0}^{m-1} c_k \varphi_k^\eta \in \mathcal{U}$ for which $\|\hat{f} - f\|$ is minimal and

$$\langle f - \hat{f}, g \rangle = 0 \quad (\forall g \in \mathcal{U}).$$

Adaptive signal models

Problem statement for applications

- Applications: $L_2(\mathbb{R})$, $H_2(\mathbb{D})$, etc.
- Some apriori information about f is usually known.
- Suppose $m \ll n$. We want to find a good approximation

$$f_i = f(t_i) \approx \sum_{k=0}^{m-1} c_k \varphi_k^\eta(t_i) = (\Phi(\eta)\mathbf{c})_i$$

$(i = 1, \dots, n, f \in L_2(\mathbb{R}))$

- Questions

- 1 How to choose the basis φ_k^η ?
- 2 How to determine the optimal η ?

Adaptive signal models

A nonlinear optimization problem

We look for the optimal parameter vector $\boldsymbol{\eta} \in \Gamma \subset \mathbb{R}^L$, for which

$$r_2(\boldsymbol{\eta}) = \|\mathbf{f} - \Phi(\boldsymbol{\eta})\Phi^+(\boldsymbol{\eta})\mathbf{f}\|_2^2 = \|\mathbf{f} - \mathbf{P}_{\Phi(\boldsymbol{\eta})}\mathbf{f}\|_2^2.$$

so-called variable projection functional is minimal.

Properties of variable projection operators

- $\mathbf{P}_{\Phi(\boldsymbol{\eta})}\mathbf{f}$ is the orthogonal projection of \mathbf{f} onto the column space of $\Phi(\boldsymbol{\eta})$.
- The gradient of $r_2(\boldsymbol{\eta})$ can be analytically calculated.¹
- Minimizing r_2 w.r.t. $\boldsymbol{\eta}$ is known as a **separable nonlinear least squares** (SNLLS) problem.

¹G., Golub, V., Pereyra. "The differentiation of pseudo-inverses and nonlinear least squares problems whose variables separate." SIAM Journal on numerical analysis (1973)

Adaptive signal models

Signal representation

- Let

$$\boldsymbol{\eta}^* = \arg \min_{\boldsymbol{\eta} \in \mathbb{R}^L} r_2(\boldsymbol{\eta}) = \arg \min_{\boldsymbol{\eta} \in \mathbb{R}^L} \|\mathbf{f} - \Phi(\boldsymbol{\eta})\Phi^+(\boldsymbol{\eta})\mathbf{f}\|_2^2.$$

- Depending on the application we can represent \mathbf{f} by

- 1 $\mathbf{f} \rightarrow \boldsymbol{\eta}^* \in \mathbb{R}^L$ (analysis)
- 2 $\mathbf{f} \rightarrow \mathbf{c} = \Phi^+(\boldsymbol{\eta}^*)\mathbf{f} \in \mathbb{R}^m$ (dimension reduction).
- 3 $\mathbf{f} \rightarrow \Phi(\boldsymbol{\eta}^*)\Phi^+(\boldsymbol{\eta}^*)\mathbf{f} \in \mathbb{R}^n$ (low pass filtering).
- 4 $\mathbf{f} \rightarrow \mathbf{f} - \Phi(\boldsymbol{\eta}^*)\Phi^+(\boldsymbol{\eta}^*)\mathbf{f} \in \mathbb{R}^n$ (high pass filtering).

System identification

Discrete time SISO-LTI systems

$$\mathbf{x} = \mathbf{h} * \mathbf{u} \xrightarrow{\mathcal{Z}} X(z) = H(z)U(z),$$

where

- \mathbf{u}, \mathbf{x} input and output sequences,
- \mathbf{h} impulse response,
- X, Y, H are the \mathcal{Z} -transforms of $\mathbf{x}, \mathbf{y}, \mathbf{h}$.
- Suppose system is causal and BIBO stable
 $\implies H \in H_\infty(\mathbb{D}) \subset H_2(\mathbb{D})$.
- **Identification task:** Find the (inverse) **poles**/zeros of the transfer function $H(z)$ in \mathbb{D} .

System identification

MT functions

- Approximate $H \in H_2(\mathbb{D})$ by a complete and orthogonal basis φ_k^η ($k = 0, \dots, m$).
- **Idea:** choose Malmquist-Takenaka (MT) functions as the basis functions:

$$\varphi_k^\eta(z) = R_{a_k} \prod_{j=0}^{k-1} B_{a_j}(z) = \frac{\sqrt{1 - |a_k|^2}}{1 - \bar{a}_k z} \prod_{j=0}^{k-1} \frac{z - a_j}{1 - \bar{a}_j z},$$

where $\eta := (a_0, \dots, a_{m-1}) \in \mathbb{D}^m$.

- The functions φ_k^η have poles at $1/\bar{a}_j$ ($j \leq k$, $a \in \mathbb{D}$).
- MT systems are complete and orthonormal in $H^2(\mathbb{D})$, provided η satisfies the *Szász condition*.
- Choosing $\eta = (0, 0, 0, \dots)$ we get the trigonometric system.

System identification

SNLLS formulation

- The frequency response $H|_{\mathbb{T}} \in H_2(\mathbb{T}) \subset L_2(\mathbb{T})$.
- Denote by $\mathbf{h} \in \mathbb{C}^n$ ($n \in \mathbb{N}$) a discrete sampling of $H|_{\mathbb{T}}$.
- SNLLS formulation:

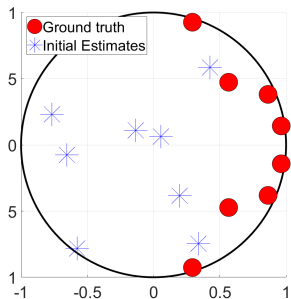
$$r_2(\boldsymbol{\eta}) = \|\Re \mathbf{h} - \Phi(\boldsymbol{\eta})\Phi^+(\boldsymbol{\eta})\Re \mathbf{h}\|_2^2,$$

where $\boldsymbol{\eta} = (r_0, \mu_0 \dots, r_{m-1}, \mu_{m-1}) \in \mathbb{R}^{2m}$, where $a_k = r_k e^{i\mu_k}$ and the columns of $\Phi(\boldsymbol{\eta})$ contain *real MT-functions*¹ sampled on \mathbb{T} .

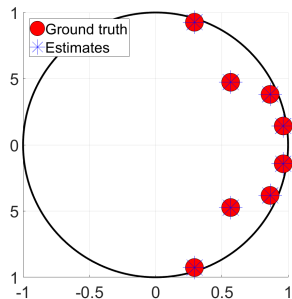
- It suffices to approximate the real part of \mathbf{h} (see e.g. *Titchmarsh theorem*).

¹F. Schipp, "Racionális ortogonális rendszerek a jel- és képfeldolgozásban", ELTE-IK Jegyzettár (2016)

System identification



(a) Random η inverse poles.



(b) Optimal η^* parameters.

Figure: A numerical example from¹.

¹T. Dózsa, M. Szabari, A. Soumelidis, P. Kovács, "Pole identification using discrete Laguerre expansion and variable projection", In: Proc. The 22nd World Congress of the International Federation of Automatic Control (IFAC2023), (2023)

Supervised learning

Problem formulation

- Approximate $G : \mathcal{X} \rightarrow \mathcal{Y}$, where \mathcal{Y} is a topological space.
- Common domains for G : $\mathcal{X} \subset \mathbb{R}^n$ or $\mathcal{X} \subset \mathbb{R}^{n \times s}$ ($n, s \in \mathbb{N}$).
- Common ranges for G :
 - $\mathcal{Y} \subset \{1, 2 \dots n\}$ ($n \in \mathbb{N}$) (**classification**, i.e. $|\mathcal{Y}| = 2$ binary classification).
 - $\mathcal{Y} \subset \mathbb{R}^n$ ($n \in \mathbb{N}$) (**Regression**).

Models

- Identify $G_\theta \approx G$, where θ is a vector of parameters (usually $\theta \in \mathbb{R}^P$, or $\theta \in \mathbb{C}^P$, $P \in \mathbb{N}$).
- **Goal:** find θ so that $E : \mathcal{Y} \times \mathcal{Y} \rightarrow [0, \infty)$

$$E(G(f), G_\theta(f)) \quad f \in \mathcal{X}$$

is minimized for *all* $f \in \mathcal{X}$.

Model training and evaluation

Training and test sets

- Suppose that the values $G(f_1), G(f_2), \dots, G(f_q)$ ($f_k \in \mathcal{T} \subset \mathcal{X}$, $k = 1, \dots, q$) are known.
- Training and test sets: $\mathcal{T}_{tr} \cup \mathcal{T}_{te} = \mathcal{T}$, $\mathcal{T}_{tr} \cap \mathcal{T}_{te} = \emptyset$.

Model training and evaluation

- **Training:** solve
$$\min_{\theta \in \mathbb{R}^P} E(G_\theta(f), G(f)), \quad (\forall f \in \mathcal{T}_{tr})$$

Denote by θ^* a solution.

- **Testing:** evaluate
$$\frac{1}{|\mathcal{T}_{te}|} \sum_{f \in \mathcal{T}_{te}} E(G_{\theta^*}(f), G(f)).$$

Classical supervised learning framework

Supervised learning steps

- 1 Identify $\mathcal{T} \subset \mathcal{X}$ and $G(f)$ ($f \in \mathcal{T}$). Measurements, expert input, data augmentation etc.
- 2 Feature extraction:
 - Classical approach: transform the dataset \mathcal{T} before training. (e.g. PCA, time frequency representations, etc.)
 - **Feature extraction transformations are incorporated into the model G_θ (e.g. convolutional neural networks).**
- 3 Choose model architecture G_θ . (e.g. SVM, Neural Networks, etc.).
- 4 Train model on (possibly transformed) data from \mathcal{T}_{tr} and evaluate on \mathcal{T}_{te} .

Open questions

Common issues

- Enough data? Labelled correctly?
- Loss function represents task to solve?
- Was the model architecture chosen correctly?
- Can extracted features be interpreted?
- Can we trust the trained model's predictions?

Model driven ML

Incorporate domain knowledge into supervised learning schemes through the use of interpretable mathematical models.

VP-NET

Variable Projection (VP) layers

- Suppose that for $\mathbf{f} \in \mathbb{R}^n$

$$\mathbf{f}_k = f(t_k) \quad (f \in L_2(\mathbb{R})).$$

Then, the mappings

- 1 $G_\eta(\mathbf{f}) = \mathbf{c} = \Phi^+(\eta)\mathbf{f} \in \mathbb{R}^m,$
- 2 $G_\eta(\mathbf{f}) = \Phi(\eta)\Phi^+(\eta)\mathbf{f} \in \mathbb{R}^n,$
- 3 $G_\eta(\mathbf{f}) = \mathbf{f} - \Phi(\eta)\Phi^+(\eta)\mathbf{f} \in \mathbb{R}^n,$

where $\eta \in \mathbb{R}^P$ and the columns $\Phi(\eta) \in \mathbb{R}^{n \times m}$ contain samplings of an orthogonal basis in $L_2(\mathbb{R})$ are called **variable projection (VP)** layers.

VP-NET

Properties of VP layers

- The gradients $\frac{\partial G_{\eta}}{\partial \eta}$ can be analytically calculated provided $\frac{\partial \Phi(\eta)}{\partial \eta}$ is known.¹
- G_{η} can be implemented as a layer in a neural network.
- If the basis functions in the columns of $\Phi(\eta)$ are chosen correctly, the parameters η can have physical meanings.²
- Usually η contains less parameters than equivalent convolution layer.

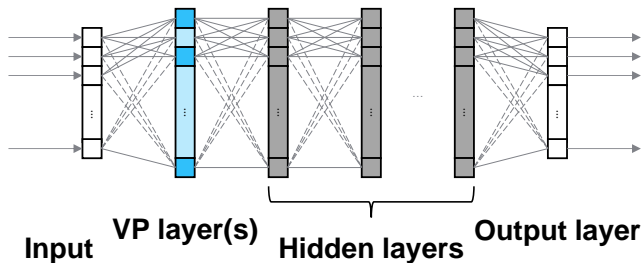
¹G., Golub, V., Pereyra. "The differentiation of pseudo-inverses and nonlinear least squares problems whose variables separate." SIAM Journal on numerical analysis (1973)

²P. Kovács, G. Bognár, C. Huber, and M. Huemer. "VPNET: Variable projection networks." International Journal of Neural Systems (2022)

VP-NET

Usual VP-NET architecture

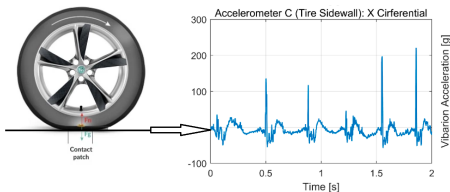
- First layers are VP-layers: these learn an appropriate representation of the data.
- Lower layers are fully connected: these solve the classification/regression task.



Tire sensor signal processing

Problem description

- Sensor: tire implanted 3D force sensors.
- Signals: changes in resistance due to mechanical forces.
- Task: surface abnormality detection.



Collaboration



NANOSENSORS
Centre for Energy Research
Institute of Technical Physics and Materials Science

Tire sensor signal processing

Signal properties

- Quasi periodic, quasi compactly supported.
- Width of support changes with vehicle speed.
- Tire revolutions occurring on abnormal surface lower SNR.
- **IDEA:** construct VP-layer using variable projection and adaptive Hermite functions. Adaptive high pass filtering.

Tire sensor signal processing

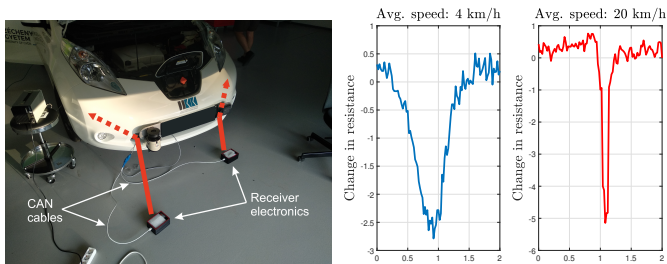


Figure: **LEFT:** test vehicle and readout electronics **RIGHT:** signal morphologies at different velocities.

Tire sensor signal processing

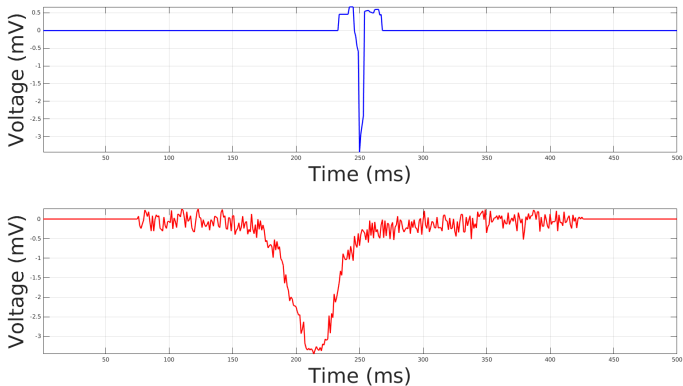


Figure: **TOP:** a tire revolution on normal surface. **BOTTOM:** a tire revolution on abnormal surface.

Tire sensor signal processing

Adaptive Hermite functions

- Classical Hermite polynomials: $\{h_k \mid k \in \mathbb{N}\}$.
- Hermite functions:

$$\varphi_k(t) = h_k(t) / \|h_k\|_2 \cdot \sqrt{w(t)} \quad (k \in \mathbb{N}),$$

where $w(t) = e^{-t^2}$.

- Adaptive Hermite functions:¹

$$\varphi_k^{(\tau, \lambda)}(t) := \sqrt{\lambda} \varphi_k(\lambda(t - \tau)) \quad (t, \tau \in \mathbb{R}, \lambda > 0).$$

- VP-layer: $G_{(\tau, \lambda)}(\mathbf{f}) := \mathbf{f} - \Phi(\tau, \lambda)\Phi^+(\tau, \lambda)\mathbf{f}$

¹T. Dózsa and P. Kovács, "ECG Signal Compression Using Adaptive Hermite Functions," ADVANCES IN INTELLIGENT SYSTEMS AND COMPUTING, vol. 399, pp. 245–254, 2015.

Tire sensor signal processing

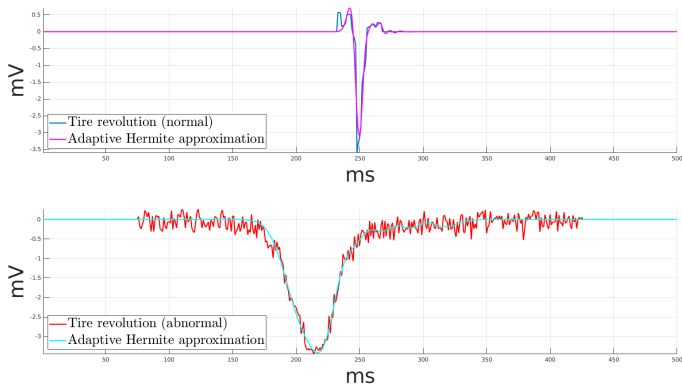


Figure: Adaptive Hermite approximations of tire signals

Tire sensor signal processing

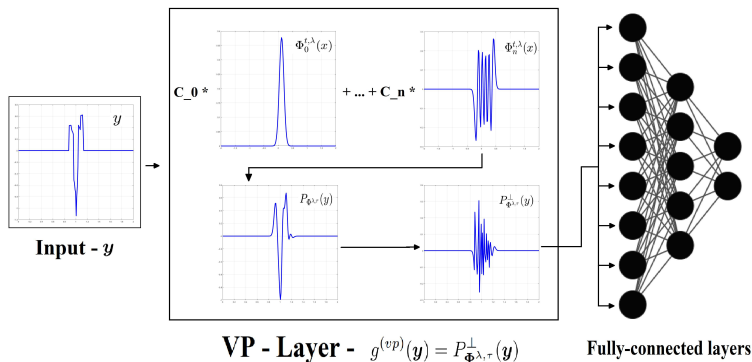


Figure: Proposed VP-NET architecture for road surface abnormality recognition

Tire sensor signal processing

Results

- Sensor: Nanosensors Laboratory, MFA
- Test vehicle: Nissan Leaf, SZTAKI-SCL
- Data: 282 normal and 235 abnormal tire revolutions.

Algorithm	Accuracy (on the test set)
SVM	94.23%
CNN	97.12%
FCNN	97.11%
VPNet (proposed)	98.08%

Table: Road surface abnormality recognition¹

¹T. Dózsa et. al., "Road Abnormality Detection Using Piezoresistive Force Sensors and Adaptive Signal Models," in IEEE Transactions on Instrumentation and Measurement, (2022)

Current research directions

Theoretical considerations

- Variable projection operators for other ML models (e.g. SVM^{1 2} and spiking networks³).
- Interpretable transformations using other new frameworks (e.g. hyperbolic convolution operators).

Applications

- Real-time road surface abnormality recognition.
- Wheel force estimation based on tire sensor signals.

¹T. Dózsa and P. Kovács, Variable projection support vector machines, Proc. 4th International Conference on Advances in Signal Processing and Artificial Intelligence (ASPAAI), (2022)

²T. Dózsa F. Deuschle, B. Cornelis, P. Kovács, "Variable projection support vector machines and some applications using adaptive Hermite expansions", International Journal of Neural Systems (2023) (Under review)

³P. Kovács, and K. Samiee. "Arrhythmia Detection Using Spiking Variable Projection Neural Networks." Computing in Cardiology (CinC) (2022)

VP-SVM

Support vector classification

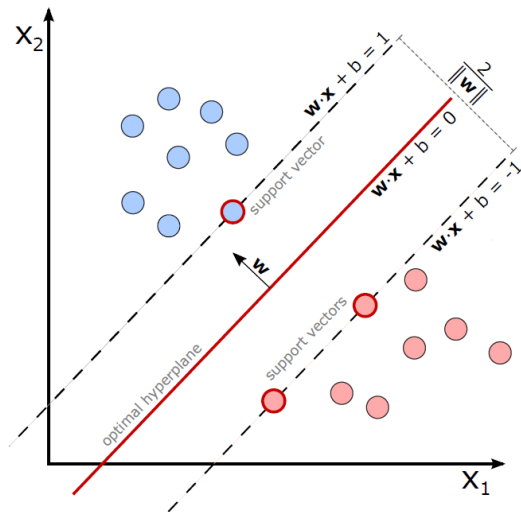
- Suppose $G : \mathbb{R}^n \supset \mathcal{X} \rightarrow \mathcal{Y} := \{-1, 1\}$.
- A Support Vector Machine (SVM) aims to identify an optimal hyperplane separating the examples in \mathcal{X} .
- $G_\theta(\mathbf{f}) := \text{sgn}(\mathbf{w}^T \mathbf{f} + b)$ ($\theta := [\mathbf{w}, b] \in \mathbb{R}^{n+1}$, $\mathbf{f} \in \mathcal{X} \subset \mathbb{R}^n$).
- Training (soft-margin SVC):

$$\min_{\mathbf{w} \in \mathbb{R}^n, b \in \mathbb{R}} \quad \frac{1}{2} \|\mathbf{w}\|^2 + \sum_{j=1}^q \xi_j,$$

$$\text{subject to} \quad y_k(\mathbf{w}^T x_k + b) \geq 1 - \xi_k,$$
$$\xi_k \geq 0 \quad (k = 1, \dots, q).$$

- Convex optimization can be used to solve for \mathbf{w} and b .

VP-SVM



VP-SVM

Unconstrained formulation

- Unconstrained SVM training formulation for a linearly separable \mathcal{X} :

$$\min_{\mathbf{w} \in \mathbb{R}^n, b \in \mathbb{R}} C \cdot \sum_{i=1}^q \max(0, 1 - y_i(\mathbf{w}^T \mathbf{f}_i + b)) + \|\mathbf{w}\|_2^2, \quad (C \in \mathbb{R}).$$

- Can solve for \mathbf{w} and b using (sub)gradient methods.
- Other efficient optimization algorithms exist.¹

¹J. Shawe-Taylor, and S. Shiliang. "A review of optimization methodologies in support vector machines." *Neurocomputing* 74.17 (2011)

VP-SVM

Variable Projection Support Vector Machines

- Suppose

$$G_{\theta}(\mathbf{f}) := \text{sgn}(\mathbf{w}^T(\Phi(\boldsymbol{\eta})^+ \mathbf{f}) + b)$$

$$(\theta := \{\mathbf{w}, \boldsymbol{\eta}, b\}, \mathbf{w} \in \mathbb{R}^m, \boldsymbol{\eta} \in \mathbb{R}^P, b \in \mathbb{R}).$$

- VP-SVM training objective:

$$C \sum_{i=1}^q \max(0, 1 - y_i(\mathbf{w}^T(\Phi^+(\boldsymbol{\eta})\mathbf{f}_i) + b)) +$$

$$\|\mathbf{w}\|_2^2 + R(\boldsymbol{\eta}),$$

where $\mathbf{w} \in \mathbb{R}^n$ and $R(\boldsymbol{\eta})$ is an added regulatory term:

$$R(\boldsymbol{\eta}) = \frac{\alpha}{q} \sum_{i=1}^q \frac{\|\mathbf{f}_i - \Phi(\boldsymbol{\eta})\Phi(\boldsymbol{\eta})^+ \mathbf{f}_i\|_2^2}{\|\mathbf{f}_i\|_2^2} \quad (\alpha \in \mathbb{R}).$$

VP-SVM

VP-SVM properties

- (Sub)gradient based methods can be used for training.
- $R(\eta)$ prevents the problem of vanishing gradients.
- More suitable for light-weight applications (less parameters than VP-NET).
- Can be used with Mercer kernels as well.^{1 2}

¹T. Dózsa and P. Kovács, Variable projection support vector machines, Proc. 4th International Conference on Advances in Signal Processing and Artificial Intelligence (ASPAI), (2022)

²T. Dózsa F. Deuschle, B. Cornelis, P. Kovács, "Variable projection support vector machines and some applications using adaptive Hermite expansions", International Journal of Neural Systems (2023) (Accepted)

Sensor fault detection

Problem description

- **GOAL:** identify peaks in accelerometer measurements appearing due to hardware failure.
- Sudden peaks can appear due to so-called shock events. These have similar morphology to sensor peaks.
- Peaks due to sensor failure and physical phenomena may overlap.

Collaboration



Sensor fault detection

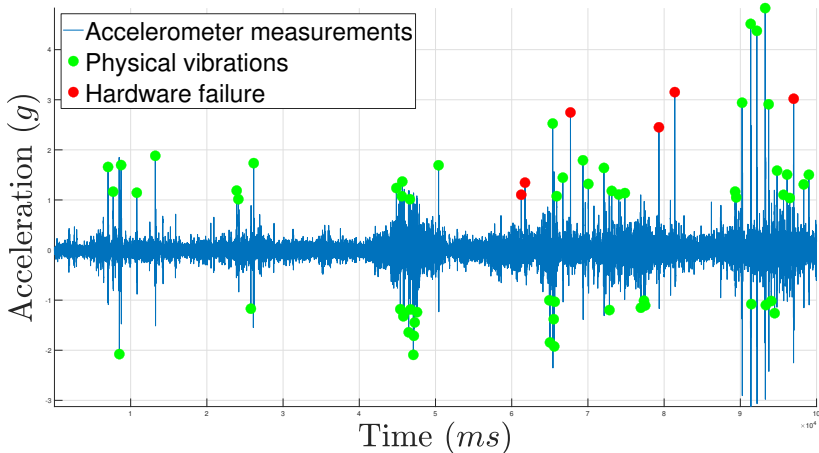


Figure: Accelerometer data to be classified

Sensor fault detection

Classification with VP-SVM

- **Difficulties and observations:**
 - Normal/abnormal examples have similar morphology.
 - Highly unbalanced dataset.
 - Examples have compact support, can be modelled efficiently with adaptive Hermite functions.
- **Methodology and preprocessing**
 - Transform measurement by truncated scalogram using complex Morlet wavelets.
 - Downsample normal examples: small training set, large test set.
 - Use VP-SVM with Gaussian kernel and adaptive Hermite functions.

Sensor fault detection

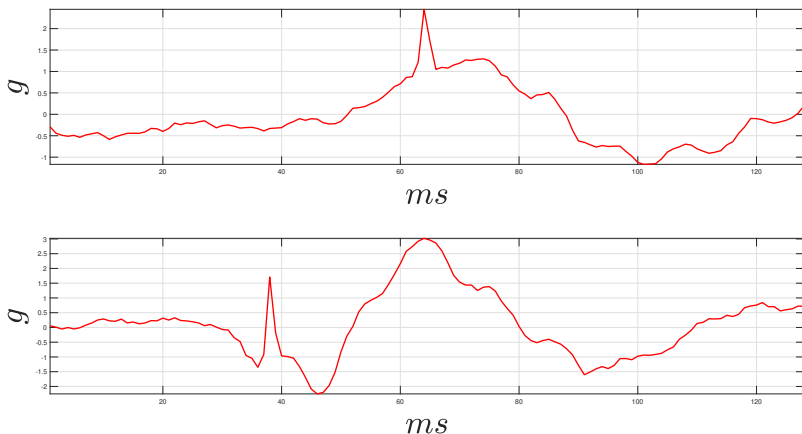


Figure: Peaks due to sensor failure can appear near vibration induced peaks

Sensor fault detection

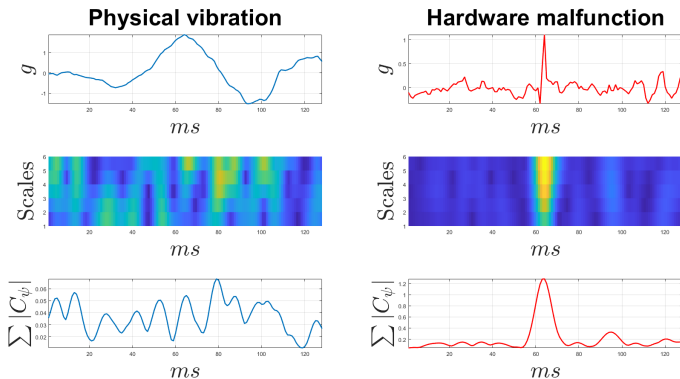


Figure: Preprocessing steps before classification with VP-SVM

Sensor fault detection

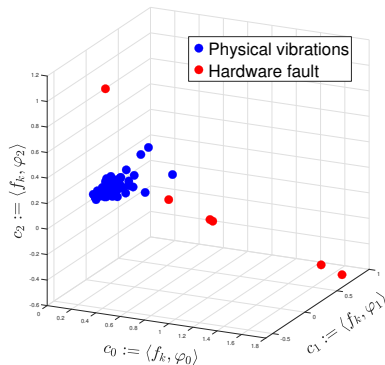


Figure: Output of the optimized adaptive Hermite transformation in the trained VP-SVM.

Conclusion

Summary

- Discussed variable projections in Hilbert spaces.
- Extension of neural networks with variable projection layers for interpretable feature extraction.
- Generalization of framework to other ML algorithms (e.g. SVM)
- Example applications:
 - Road abnormality detection using tire sensor signals.
 - Sensor fault detection in accelerometer measurements.

Conclusion

Thank you for your attention