# Beyond Petrov-Galerkin projection by using « multi-space » priors

Cédric Herzet
Inria, France

Joint work with M. Diallo & P. Héas

# The target problem...

Find $\boldsymbol{h}^\star \in \mathcal{H}$ such that $a(\boldsymbol{h}^\star, \boldsymbol{h}) = b(\boldsymbol{h})$ $\quad \forall \boldsymbol{h} \in \mathcal{H}$

where $\quad \mathcal{H}$ is a Hilbert space ($\langle \cdot, \cdot \rangle$ and $\|\cdot\|$)

$a : \mathcal{H} \times \mathcal{H} \to \mathbb{R}$ is a bilinear operator

$b : \mathcal{H} \to \mathbb{R}$ is a linear operator

# ... and its Petrov-Galerkin approximation

Find $\hat{\boldsymbol{h}}_{\mathrm{PG}} \in V_n$ such that $\ a(\hat{\boldsymbol{h}}_{\mathrm{PG}}, \boldsymbol{h}) = b(\boldsymbol{h}) \quad \forall \boldsymbol{h} \in Z_n$

where $V_n \subset \mathcal{H}$, $Z_n \subset \mathcal{H}$ are $n$-dimensional subspaces

The precision of Petrov-Galerkin can be quantified by an « instance optimal property »

$$\left\| \boldsymbol{h}^\star - \hat{\boldsymbol{h}}_{\mathrm{PG}} \right\| \leq C(V_n, Z_n) \operatorname{dist}(\boldsymbol{h}^\star, V_n),$$

# Standard methods constructing $V_n$ often return a set of subspaces and their « widths »

Standard outputs of methods constructing $V_n$:

$$V_0 \subset V_1 \subset \ldots \subset V_n, \qquad \dim(V_k) = k$$

such that

$$\mathrm{dist}(\boldsymbol{h}^\star, V_k) \leq \hat{\epsilon}_k, \qquad k = 0 \ldots n.$$

*E.g.*, « reduced basis » methods

# The Petrov-Galerkin projection discards most of the available information

Standard outputs of methods constructing $V_n$:

$$\textcolor{red}{\mathbf{X}} \subset \textcolor{red}{\mathbf{X}} \subset \ldots \subset V_n, \qquad \dim(V_k) = k$$

such that

$$\mathrm{dist}(\boldsymbol{h}^\star, \textcolor{red}{\mathbf{X}}) \leq \hat{\epsilon}_k, \qquad k = 0 \ldots n.$$

*E.g.*, « reduced basis » methods

Can we use this information to improve the projection process?

# The Petrov-Galerkin projection can be reformulated as a variational problem

$$\hat{\boldsymbol{h}}_{\mathrm{PG}} = \arg\min_{\boldsymbol{h} \in V_n} \sum_{j=1}^{n} (b_j - \langle \boldsymbol{a}_j, \boldsymbol{h} \rangle)^2$$

where     $\mathrm{span}\left(\{\boldsymbol{z}_j\}_{j=1}^{n}\right) = Z_n$

$\boldsymbol{a}_j$ is the Riesz's representer of $a(\cdot, \boldsymbol{z}_j)$
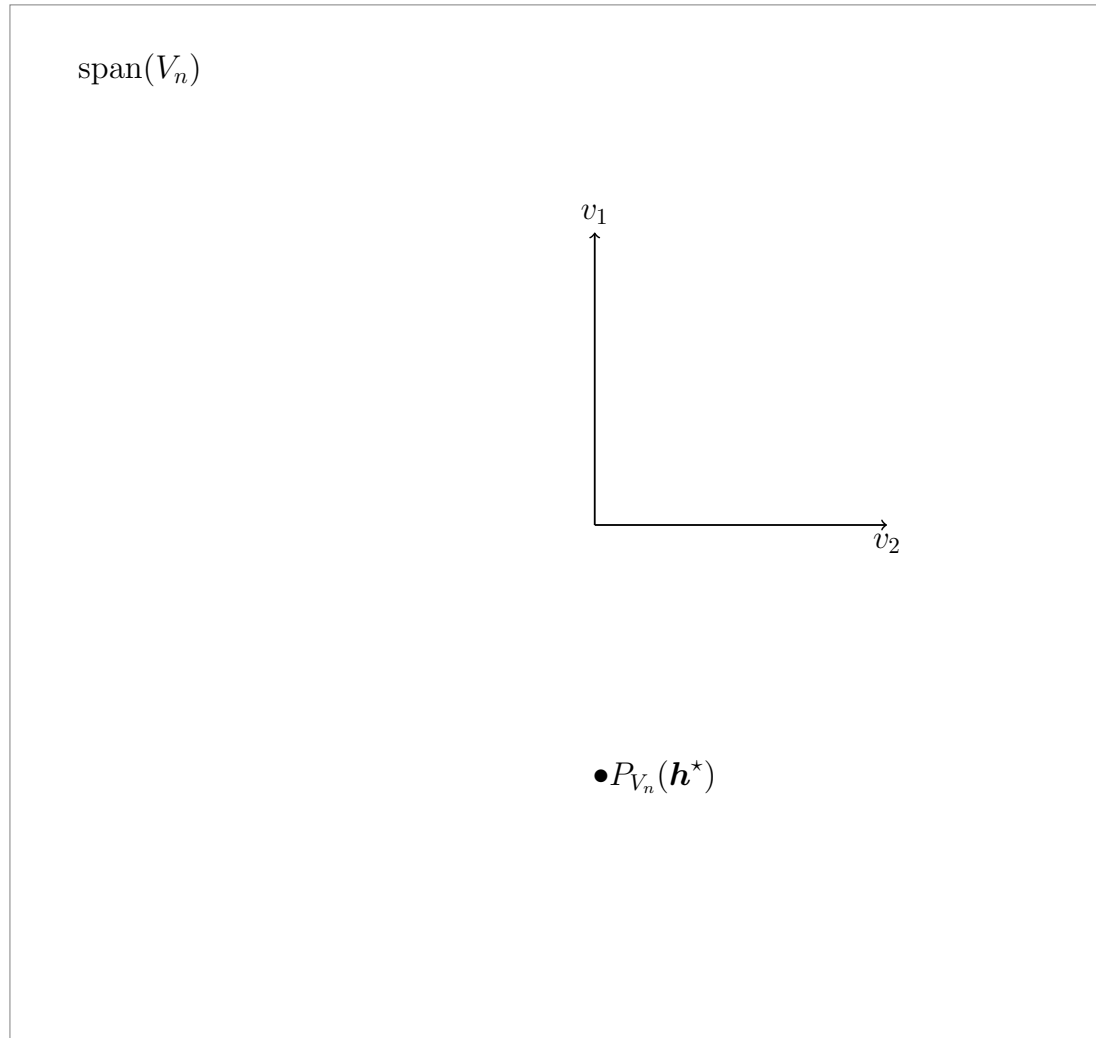
$b_j = b(\boldsymbol{z}_j)$

# The proposed « multi-space » decoder adds new constraints to the variational problem

$$\hat{\boldsymbol{h}}_{\mathrm{MS}} = \arg\min_{\boldsymbol{h}\in V_n} \sum_{j=1}^{n}(b_j - \langle \boldsymbol{a}_j, \boldsymbol{h}\rangle)^2,$$

$$\text{subject to } \mathrm{dist}(\boldsymbol{h}, V_k) \leq \hat{\epsilon}_k, \quad k = 0\dots n-1.$$
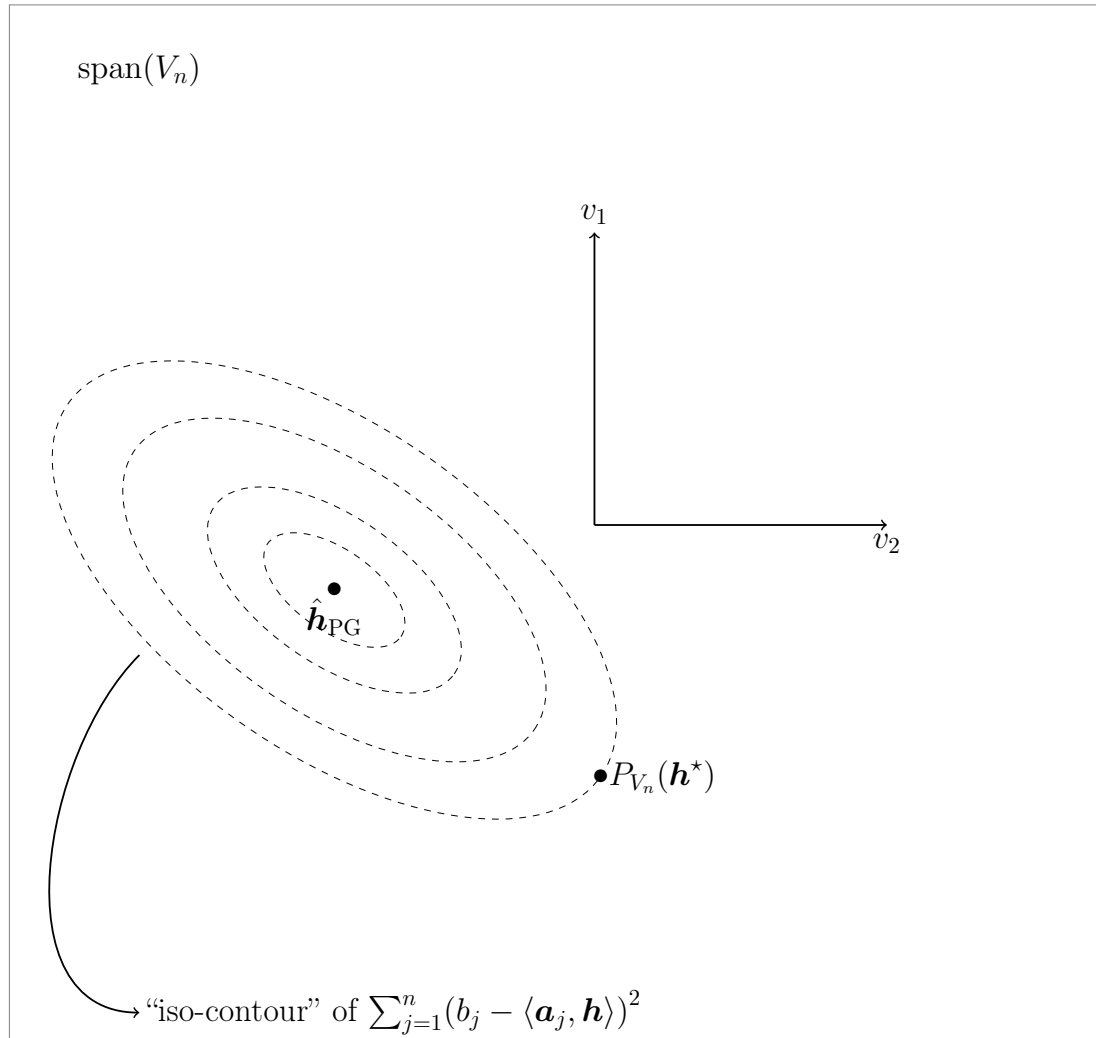
See *[Binev et al., SIAM JUQ 17 ]* for a related work.

# A graphical representation of the problem

$n = 2$

$V_k = \text{span}\left(\{v_i\}_{i=1}^k\right)$
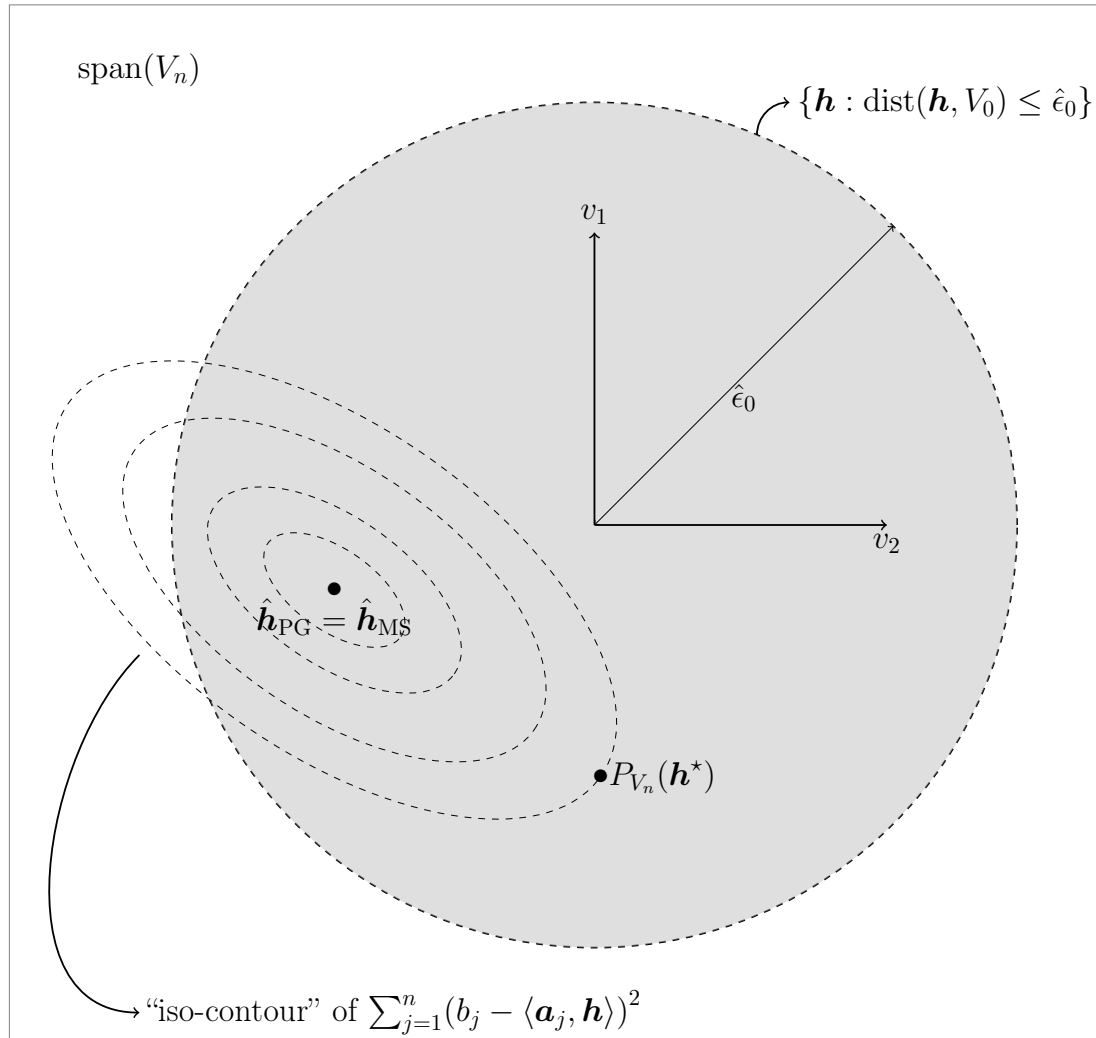
# A graphical representation of the problem



"iso-contour" of $\sum_{j=1}^{n}(b_j - \langle \boldsymbol{a}_j, \boldsymbol{h} \rangle)^2$

$n = 2$

$V_k = \operatorname{span}\left( \{\boldsymbol{v}_i\}_{i=1}^{k} \right)$

The shape of the iso-contours depends on $\mathbf{G} = [\langle \boldsymbol{a}_i, \boldsymbol{v}_j \rangle]_{ij}$
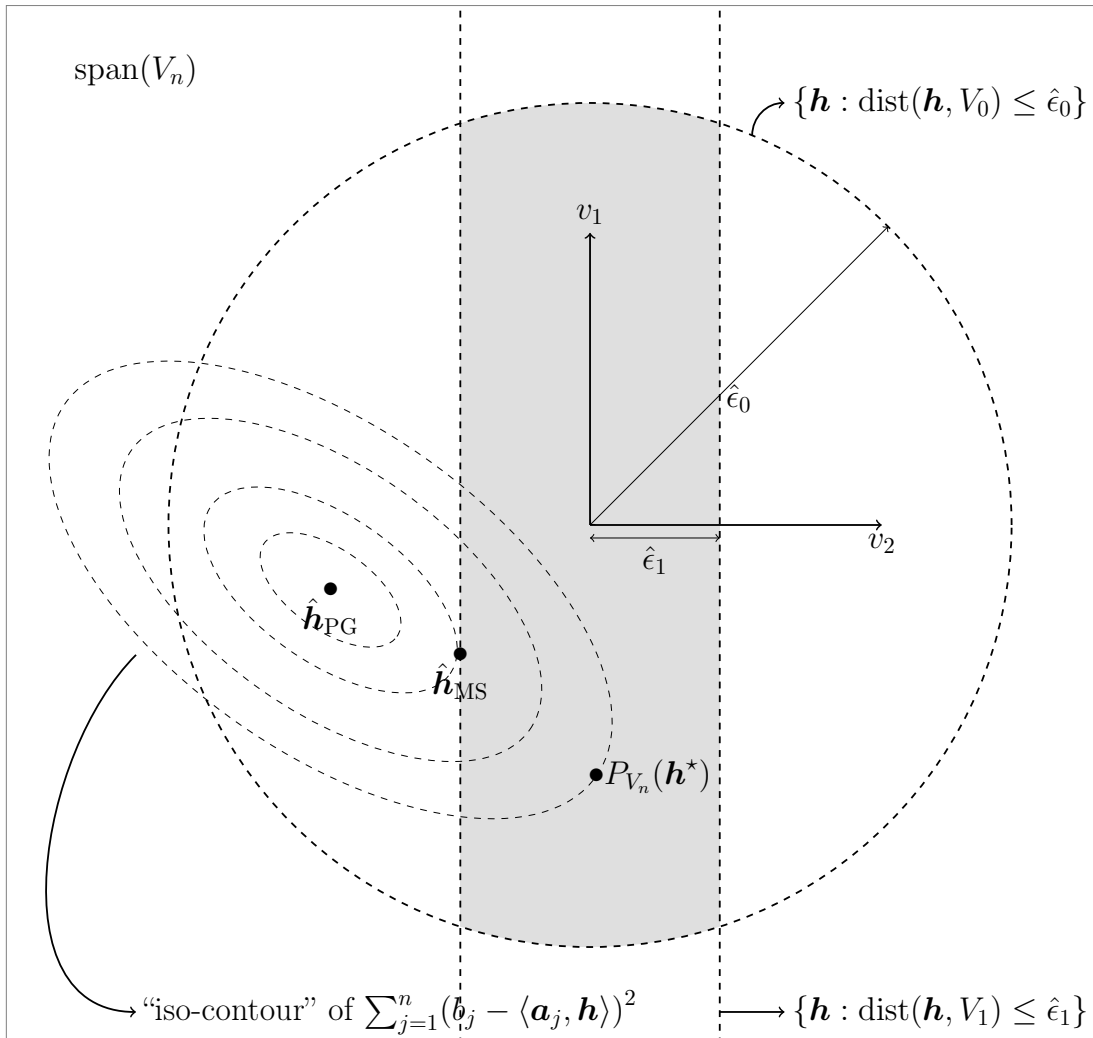
# A graphical representation of the problem



$n = 2$

$$V_k = \mathrm{span}\left(\{\boldsymbol{v}_i\}_{i=1}^k\right)$$

The shape of the iso-contours depends on $\mathbf{G} = [\langle \boldsymbol{a}_i, \boldsymbol{v}_j \rangle]_{ij}$

Labels within figure:

$\mathrm{span}(V_n)$

$\{\boldsymbol{h} : \mathrm{dist}(\boldsymbol{h}, V_0) \leq \hat{\epsilon}_0\}$

$v_1$

$\hat{\epsilon}_0$

$v_2$

$\hat{\boldsymbol{h}}_{\mathrm{PG}} = \hat{\boldsymbol{h}}_{\mathrm{MS}}$

$P_{V_n}(\boldsymbol{h}^\star)$

"iso-contour" of $\sum_{j=1}^n (b_j - \langle \boldsymbol{a}_j, \boldsymbol{h} \rangle)^2$

# A graphical representation of the problem



$n = 2$

$$V_k = \text{span}\left(\{\boldsymbol{v}_i\}_{i=1}^{k}\right)$$

The shape of the iso-contours depends on $\mathbf{G} = [\langle \boldsymbol{a}_i, \boldsymbol{v}_j \rangle]_{ij}$

The feasibility region depends on $\{V_k\}_{k=1}^{n}$ and $\{\hat{\epsilon}_k\}_{k=1}^{n}$

Can we give some guarantee on the performance of the « multi-space » decoder?
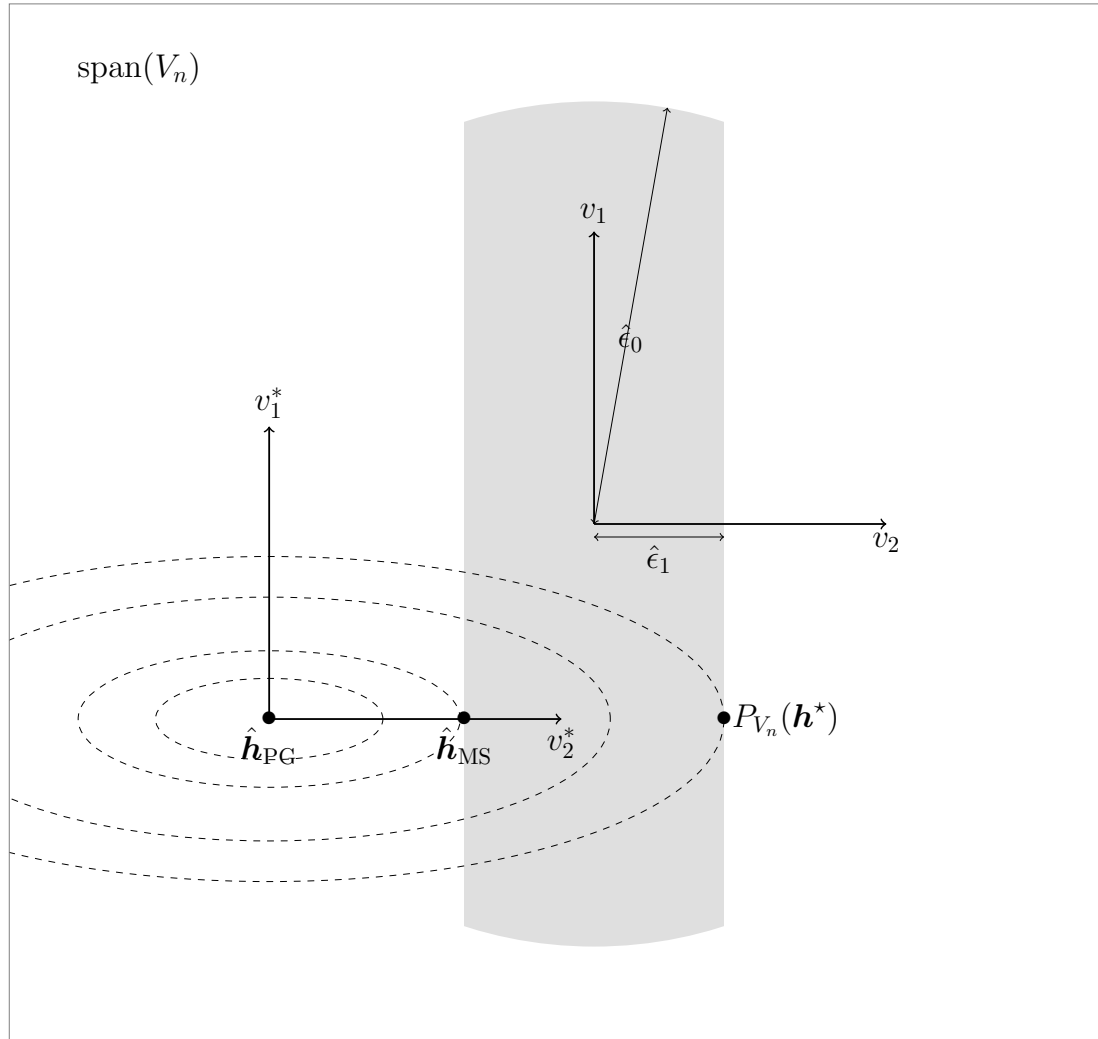
# Instance optimality properties

- *Petrov-Galerkin*: $\left\| \boldsymbol{h}^\star - \hat{\boldsymbol{h}}_{\mathrm{PG}} \right\| \leq C(\mathbf{G}) \, \mathrm{dist}(\boldsymbol{h}^\star, V_n),$

- *« Multi-space » decoder*:

$$\left\| \boldsymbol{h}^\star - \hat{\boldsymbol{h}}_{\mathrm{MS}} \right\| \leq \left( \sum_{j=\ell+1}^{n} \delta_j^2 + \rho \, \delta_\ell^2 + (\mathrm{dist}(\boldsymbol{h}^\star, V_n))^2 \right)^{\frac{1}{2}}$$
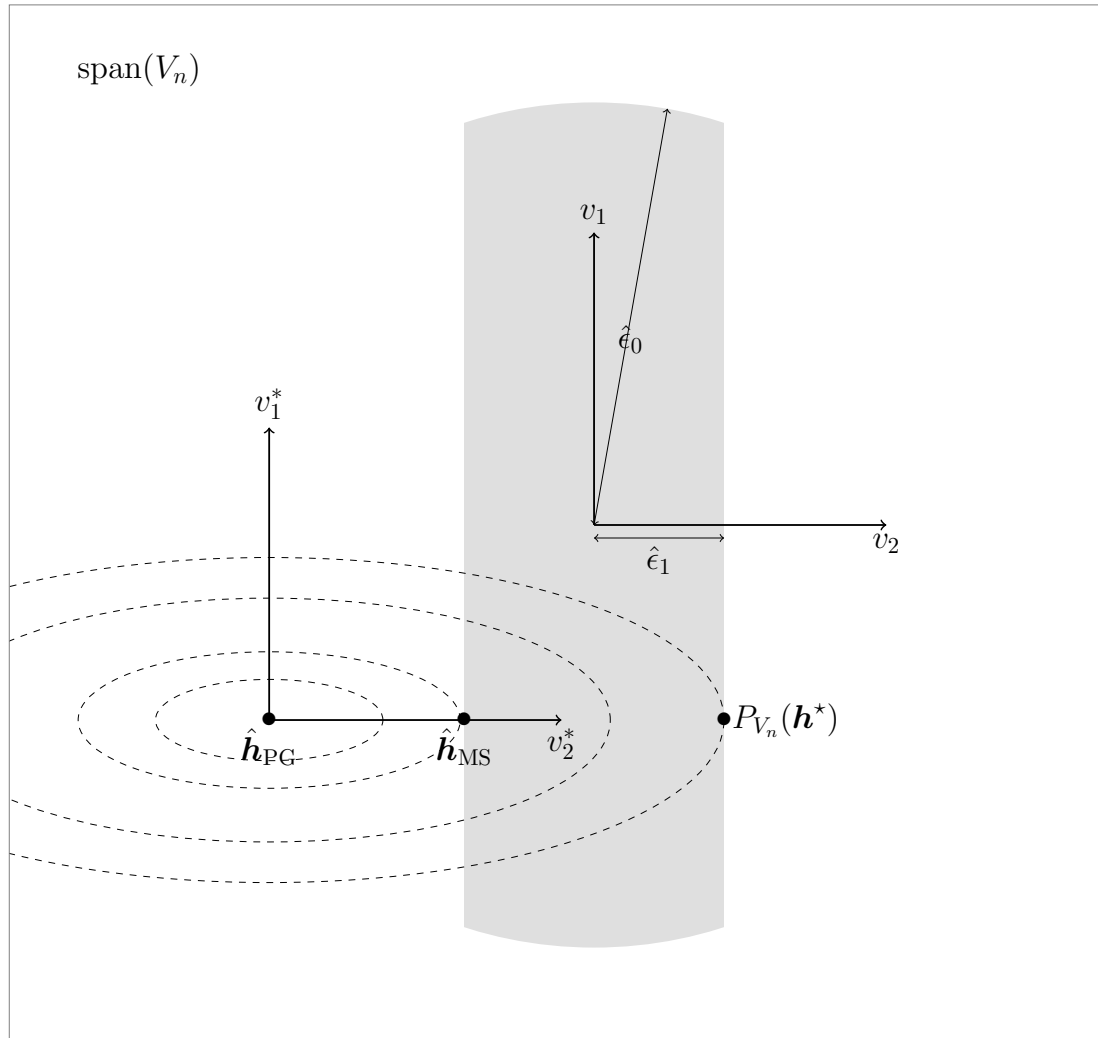
where $\ell$ and $\delta_j$'s are "easily-computable" quantities only depending on $\mathbf{G} = [\langle \boldsymbol{a}_i, \boldsymbol{v}_j \rangle]_{i,j}$, $\{\hat{\epsilon}_k\}_{k=1}^{n-1}$ and $\{\mathrm{dist}(\boldsymbol{h}^\star, V_k)\}_{k=1}^{n}$

# Particularization of the instance optimality bound to some examples

# Particularization of the instance optimality bound to some examples



$$\hat{\epsilon}_j = \begin{cases} 1 & j = 0 \ldots n-3, \\ \epsilon^{\frac{1}{2}} & j = n-2, n-1, \\ \epsilon & j = n, \end{cases}$$

$$\mathrm{dist}(\boldsymbol{h}^\star, V_k) = \hat{\epsilon}_j$$

$$\sigma_j = \mathrm{sing.\ values\ of\ } \mathbf{G}$$

$$= \begin{cases} 1 & j = 1 \ldots n-3, \\ \epsilon^{\frac{1}{2}} & j = n-2, n-1, \\ \epsilon & j = n. \end{cases}$$

$$\left\| \boldsymbol{h}^\star - \hat{\boldsymbol{h}}_{\mathrm{PG}} \right\| \leq 1$$

$$\left\| \boldsymbol{h}^\star - \hat{\boldsymbol{h}}_{\mathrm{MS}} \right\| \leq 3\epsilon^{\frac{1}{2}}$$

*Quid* of the computational complexity?

PG projection can been carried out efficiently with a complexity $\mathcal{O}(n^2)$ per iteration

$$\hat{\boldsymbol{h}}_{\mathrm{PG}} = \arg \min_{\boldsymbol{h} \in V_n} \sum_{j=1}^{n} (b_j - \langle \boldsymbol{a}_j, \boldsymbol{h} \rangle)^2$$

« Least square » problem: can be solved efficiently via gradient-based methods with a complexity $\mathcal{O}(n^2)$ per iteration.

# Our decoder can also be implemented with a complexity $\mathcal{O}(n^2)$ per iteration

Our problem can be rewritten as:

$$\hat{\boldsymbol{h}}_{\mathrm{MS}} = \arg\min_{\boldsymbol{h} \in V_n} \sum_{j=1}^{n} (b_j - \langle \boldsymbol{a}_j, \boldsymbol{h} \rangle)^2,$$

$$\text{subject to } \left\| P_{V_k}^{\perp}(\boldsymbol{h}) \right\| \leq \hat{\epsilon}_k, \quad k = 0 \ldots n-1.$$

We use the « Alternating Directions Method of Multipliers » to solve this convex problem with a complexity $\mathcal{O}(n^2)$ per iteration
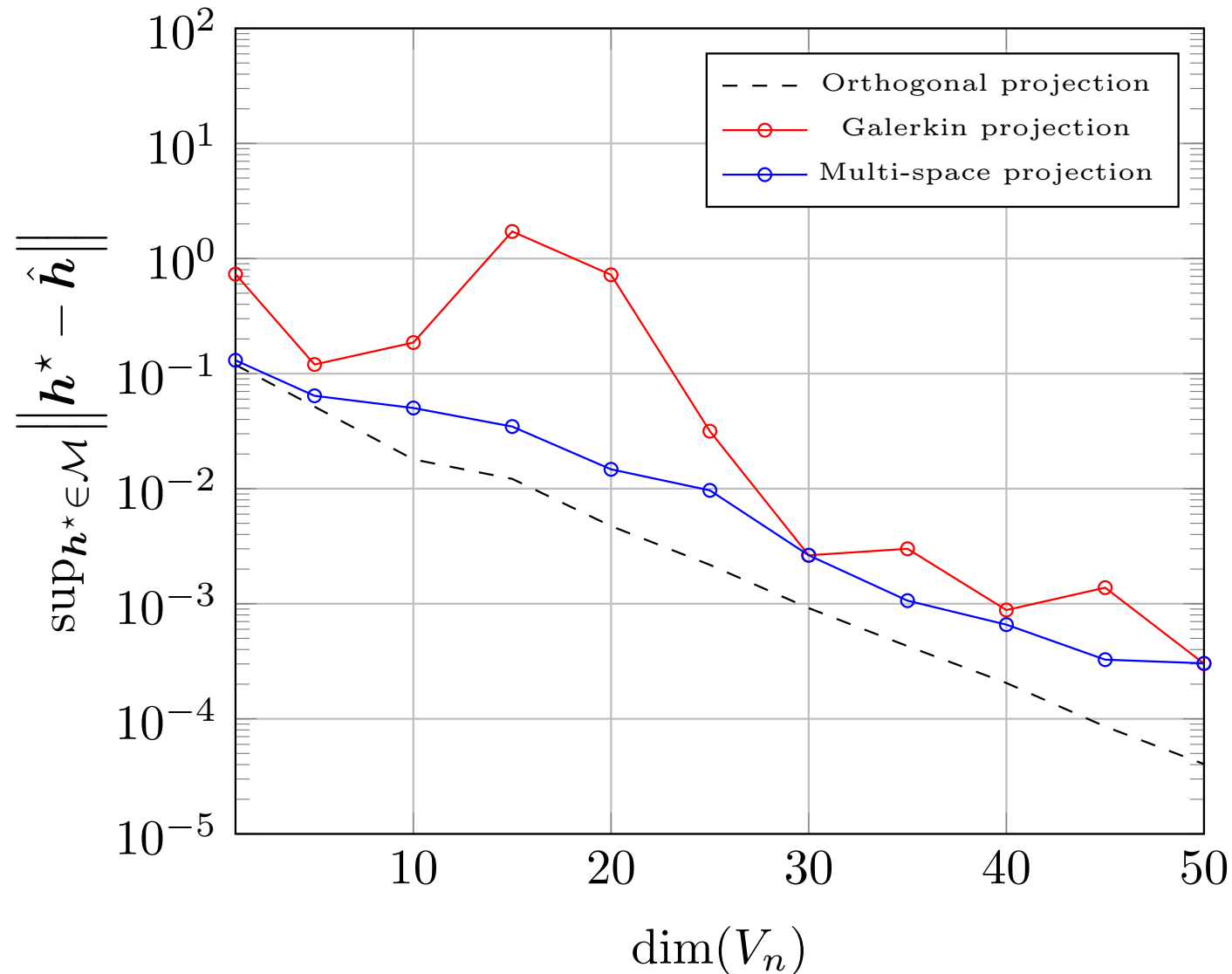
# Some results

# We consider a parametric mass transfert problem

$$\mu_1 \triangle \boldsymbol{h} + \mathbf{b}(\mu_2) \cdot \nabla \boldsymbol{h} = \boldsymbol{s} \quad \text{in } \Omega$$

$$\mu_1 \nabla \boldsymbol{h} \cdot \mathbf{n} = 0 \quad \text{in } \partial\Omega$$

where
$$\mathbf{b}(\mu_2) = [\cos(\mu_2) \sin(\mu_2)]^{\mathrm{T}}$$

$$\boldsymbol{s} = \exp\left(-\frac{\|\mathbf{x} - \mathbf{m}\|_2^2}{2\sigma^2}\right)$$

$$\mu_1 = [0.03, 0.05]$$

$$\mu_2 = [0, 2\pi]$$

# The multi-space decoder improves the worst-case approximation error over PG

# Our contributions

- We exploit additional information in the Petrov-Galerkin projection

- We derive an instance optimal guarantee for our « multi-space » decoder

- We propose an efficient implementation for the « multi-space » decoder

Thank you for your attention.

# Appendix

# Main Theorem (IOP for multi-space decoder)

**Theorem 1.** *The solution of the "multi-space" decoder satisfies:*

$$\left\| h^\star - \hat{h}_{\mathrm{MS}} \right\| \leq \left( \sum_{j=\ell+1}^{n} \delta_j^2 + \rho\, \delta_\ell^2 + (\mathrm{dist}(h^\star, V_n))^2 \right)^{\frac{1}{2}},$$

*where $\ell$ is the largest integer such that*

$$\sum_{j=\ell}^{n} \sigma_j^2 \delta_j^2 \geq 4\gamma^2 (\mathrm{dist}(h^\star, V_n))^2,$$

*and $\rho \in [0,1]$ is defined as*

$$\rho\sigma_\ell^2 \delta_\ell^2 + \sum_{j=\ell+1}^{n} \sigma_j^2 \delta_j^2 = 4\gamma^2 (\mathrm{dist}(h^\star, V_n))^2.$$

$$\gamma = \sup_{h \in V^\perp, \|h\|=1} \left( \sum_{j=1}^{n} \langle a_j, h \rangle^2 \right)^{\frac{1}{2}},$$

$$\mathbf{G} = [\langle a_i, v_j \rangle]_{i,j}$$

$$\mathbf{X} = \text{right singular vectors of } \mathbf{G},$$

$$\sigma_k = k\text{th singular value of } \mathbf{G},$$

$$\delta_j = \sum_{k=1}^{n} |x_{kj}|(\hat{\epsilon}_{k-1} + \mathrm{dist}(h^\star, V_k)).$$

# ADMM (first step): problem reformulation

$$\hat{\boldsymbol{h}}_{\mathrm{MS}} = \arg\min_{\boldsymbol{h}_n \in V_n} \sum_{j=1}^{n} (b_j - \langle \boldsymbol{a}_j, \boldsymbol{h}_n \rangle)^2,$$

$$\text{subject to } \begin{cases} \left\| P_{V_k}^{\perp}(\boldsymbol{h}_k) \right\| \leq \hat{\epsilon}_k, & k = 0 \ldots n-1, \\ \boldsymbol{h}_0 = \boldsymbol{h}_1 = \ldots = \boldsymbol{h}_n. \end{cases}$$

We add *n* new variables and *n* new constraints.

# ADMM (second step): main recursions

$$h_n^{(l+1)} = \underset{h_n \in V_n}{\arg\min} \sum_{j=1}^{n} (b_j - \langle a_j, h_n \rangle)^2 + \frac{\rho}{2} \sum_{k=1}^{n-1} \left\| h_n - h_k^{(l)} + u_k^{(l)} \right\|^2,$$

$$h_k^{(l+1)} = \underset{h_k}{\arg\min} \left\| h_k - h_n^{(l+1)} + u_k^{(l)} \right\|^2 \text{ subject to } \left\| P_{V_k}^{\perp}(h_k) \right\| \leq \hat{\epsilon}_k,$$

$$u_k^{(l+1)} = u_k^{(l)} + h_n^{(l+1)} - h_k^{(l+1)}$$

$h_k^{(l)} \in V_n$ and $u_k^{(l)} \in V_n \ \forall k, l$

They can thus be represented by $n$-dimensional vectors.