

Optimization Machine learning

Nelly Pustelnik

CNRS, Laboratoire de Physique de l'ENS de Lyon, France

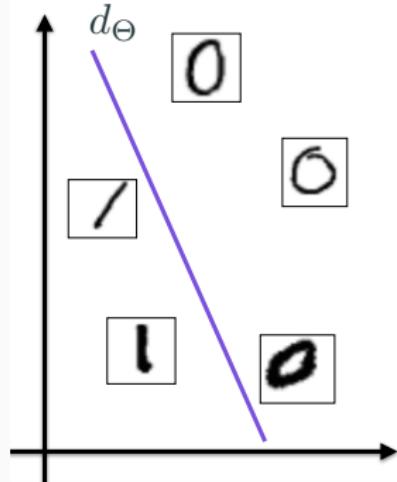


Example: Supervised learning

- **Database:** $\mathcal{S} = \{(u_\ell, c_\ell) \in \mathcal{H} \times \mathcal{G} \mid \ell \in \{1, \dots, L\}\}$

e.g. $u_\ell \in \underbrace{\mathbb{R}^N}_{\mathcal{H}}$ image and $c_\ell \in \underbrace{\{-1, +1\}}_{\mathcal{G}}$ classe

- **Goal:** Learn a prediction function $d_\Theta: \mathcal{H} \rightarrow \mathcal{G}$

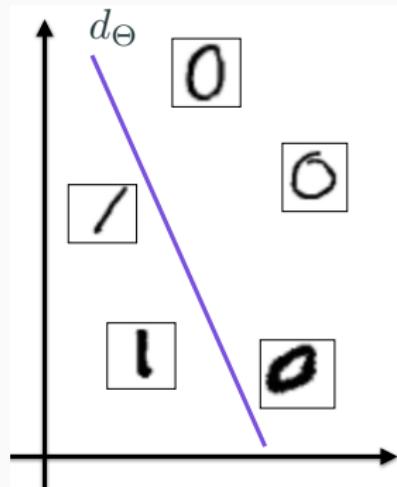


Example: Supervised learning

- **Database:** $\mathcal{S} = \{(u_\ell, c_\ell) \in \mathcal{H} \times \mathcal{G} \mid \ell \in \{1, \dots, L\}\}$

e.g. $u_\ell \in \underbrace{\mathbb{R}^N}_{\mathcal{H}}$ image and $c_\ell \in \underbrace{\{-1, +1\}}_{\mathcal{G}}$ classe

- **Goal:** Learn a prediction function $d_\Theta: \mathcal{H} \rightarrow \mathcal{G}$

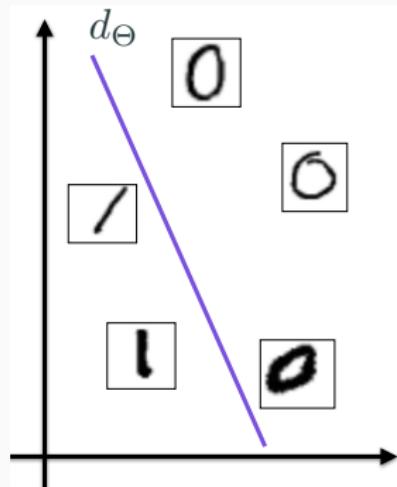


Example: Supervised learning

- **Database:** $\mathcal{S} = \{(u_\ell, c_\ell) \in \mathcal{H} \times \mathcal{G} \mid \ell \in \{1, \dots, L\}\}$

e.g. $u_\ell \in \underbrace{\mathbb{R}^N}_{\mathcal{H}}$ image and $c_\ell \in \underbrace{\{-1, +1\}}_{\mathcal{G}}$ classe

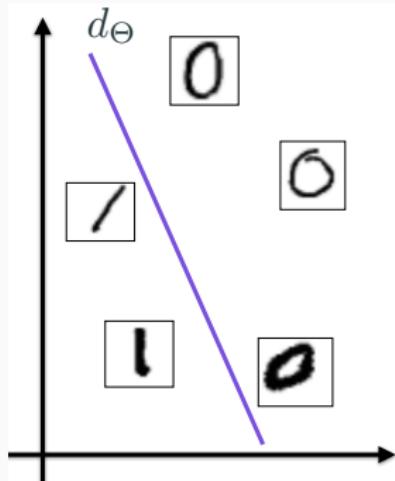
- **Goal:** Learn a prediction function $d_\Theta: \mathcal{H} \rightarrow \mathcal{G}$



Example: Supervised learning

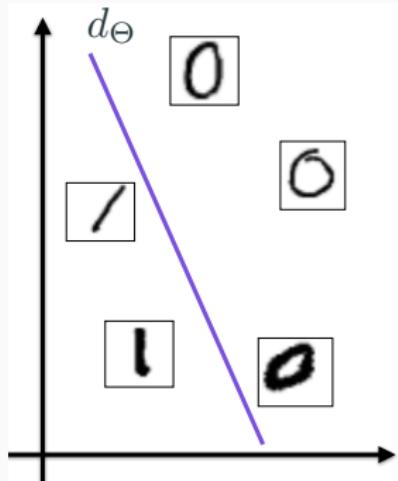
- **Database:** $\mathcal{S} = \{(u_\ell, c_\ell) \in \mathcal{H} \times \mathcal{G} \mid \ell \in \{1, \dots, L\}\}$
- **Goal:** Learn a prediction function $d_\Theta: \mathcal{H} \rightarrow \mathcal{G}$
- **Typical choices for \mathcal{H} :**

- $\mathcal{H} = \mathbb{R}^N$ for image of size $N = N_1 \times N_2$;
- $\mathcal{H} = \mathbb{R}^{N \times M}$ for multivariate images with N samples and M components;



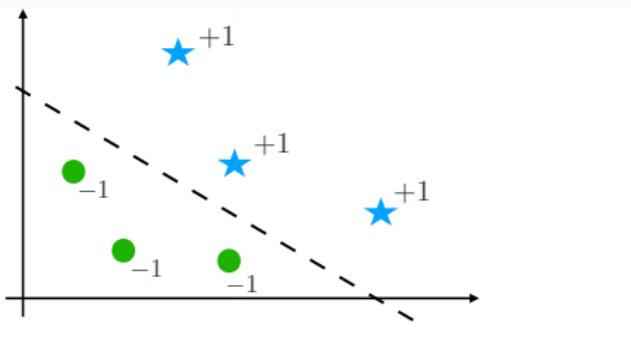
Example: Supervised learning

- **Database:** $\mathcal{S} = \{(u_\ell, c_\ell) \in \mathcal{H} \times \mathcal{G} \mid \ell \in \{1, \dots, L\}\}$
- **Goal:** Learn a prediction function $d_\Theta: \mathcal{H} \rightarrow \mathcal{G}$
- **Typical choices for \mathcal{H} :**
 - $\mathcal{H} = \mathbb{R}^N$ for image of size $N = N_1 \times N_2$;
 - $\mathcal{H} = \mathbb{R}^{N \times M}$ for multivariate images with N samples and M components;
- **Typical choices for \mathcal{G} :**
 - $\mathcal{G} = \{-1, +1\}$ for binary classification;
 - $\mathcal{G} = \{1, \dots, K\}$ for multiclass classification;
 - $\mathcal{G} = \mathbb{R}$ for regression;
 - $\mathcal{G} = \mathbb{R}^K$ for multivariate regression;



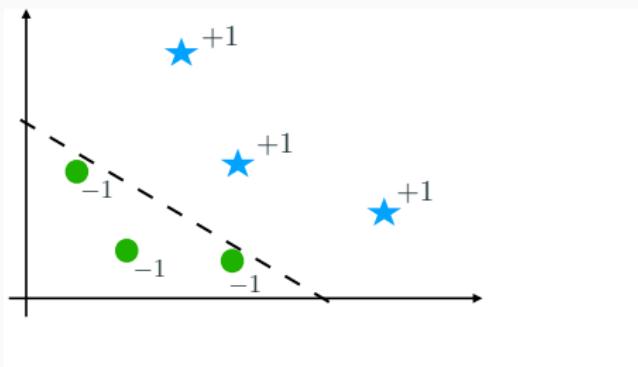
Support-vector machine (SVM)

- SVM – **Vladimir Vapnik** (early 90's) – Supervised classification.
- Example : dataset $\mathcal{S} = \{(u_\ell, c_\ell) \in \mathbb{R}^2 \times \{-1; +1\} \mid \ell \in \{1, \dots, 6\}\}$.



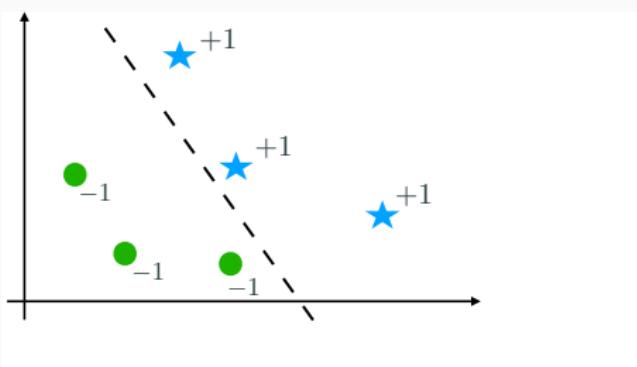
Support-vector machine (SVM)

- SVM – **Vladimir Vapnik** (early 90's) – Supervised classification.
- Example : dataset $\mathcal{S} = \{(u_\ell, c_\ell) \in \mathbb{R}^2 \times \{-1; +1\} \mid \ell \in \{1, \dots, 6\}\}$.
- Many hyperplanes that might classify the data.



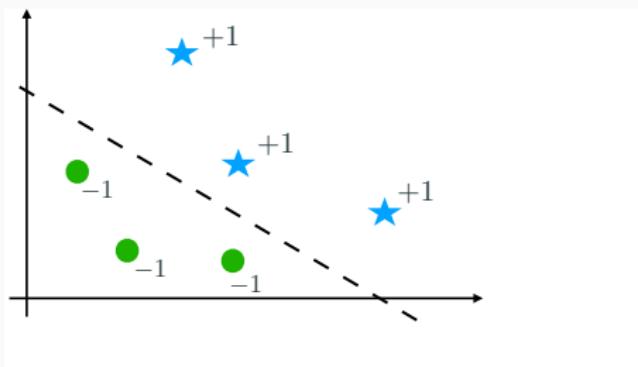
Support-vector machine (SVM)

- SVM – **Vladimir Vapnik** (early 90's) – Supervised classification.
- Example : dataset $\mathcal{S} = \{(u_\ell, c_\ell) \in \mathbb{R}^2 \times \{-1; +1\} \mid \ell \in \{1, \dots, 6\}\}$.
- Many hyperplanes that might classify the data.



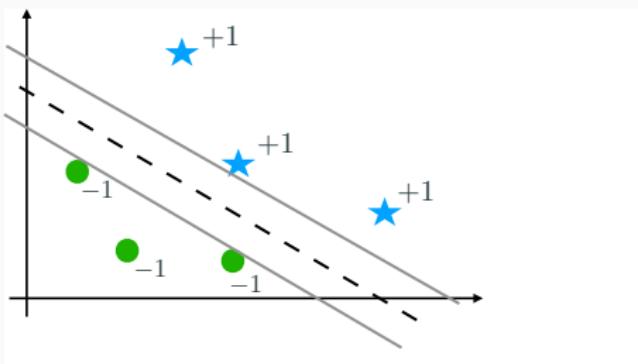
Support-vector machine (SVM)

- SVM – **Vladimir Vapnik** (early 90's) – Supervised classification.
- Example : dataset $\mathcal{S} = \{(u_\ell, c_\ell) \in \mathbb{R}^2 \times \{-1; +1\} \mid \ell \in \{1, \dots, 6\}\}$.
- Many hyperplanes that might classify the data.



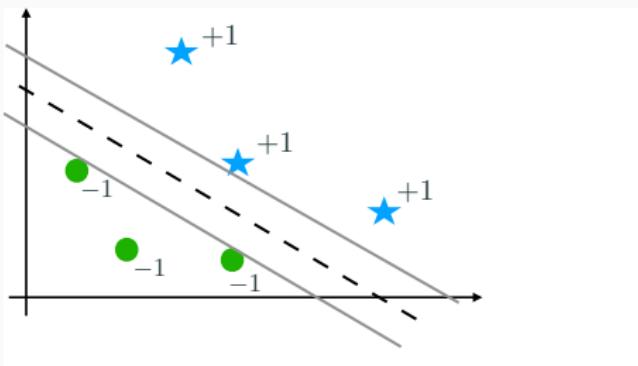
Support-vector machine (SVM)

- SVM – **Vladimir Vapnik** (early 90's) – Supervised classification.
- Example : dataset $\mathcal{S} = \{(u_\ell, c_\ell) \in \mathbb{R}^2 \times \{-1; +1\} \mid \ell \in \{1, \dots, 6\}\}$.
- Many hyperplanes that might classify the data.
- **Best hyperplane = largest margin between the two classes.**



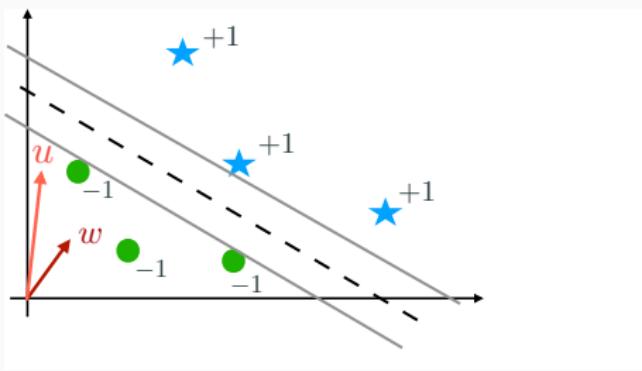
Support-vector machine (SVM)

- Best hyperplane = largest margin between the two classes.



Support-vector machine (SVM)

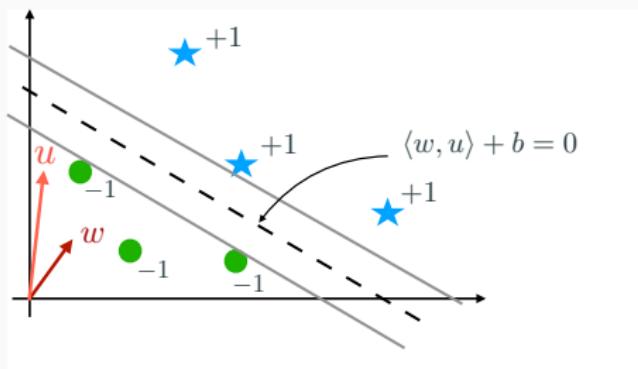
- Best hyperplane = largest margin between the two classes.



- Decision rule: $\langle w, u \rangle \geq cte$

Support-vector machine (SVM)

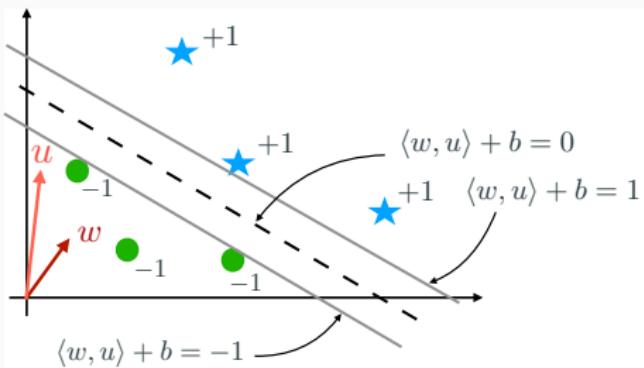
- Best hyperplane = largest margin between the two classes.



- **Decision rule:** $\langle w, u \rangle \geq \text{cte} \Leftrightarrow \boxed{\langle w, u \rangle + b \geq 0}$

Support-vector machine (SVM)

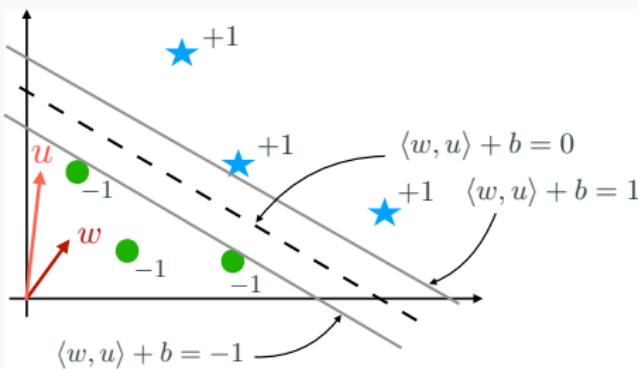
- Best hyperplane = largest margin between the two classes.



- **Decision rule:** $\langle w, u \rangle \geq \text{cte} \Leftrightarrow \boxed{\langle w, u \rangle + b \geq 0}$
- Deal with margin: $c(\langle w, u \rangle + b) \geq 1 \Rightarrow \begin{cases} \langle w, u_+ \rangle + b \geq 1 \\ \langle w, u_- \rangle + b \leq -1 \end{cases}$

Support-vector machine (SVM)

- Best hyperplane = largest margin between the two classes.

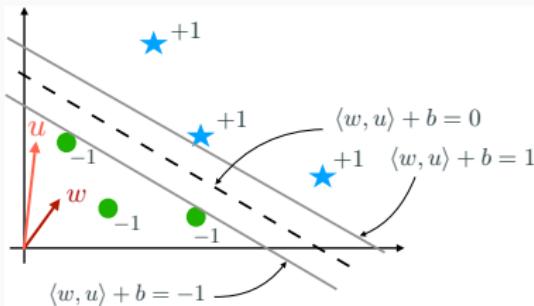


- Decision rule: $\langle w, u \rangle \geq \text{cte} \Leftrightarrow \boxed{\langle w, u \rangle + b \geq 0}$
- Deal with margin: $c(\langle w, u \rangle + b) \geq 1 \Rightarrow \begin{cases} \langle w, u_+ \rangle + b \geq 1 \\ \langle w, u_- \rangle + b \leq -1 \end{cases}$

→ Margin delimitation: $\boxed{c(\langle w, u \rangle + b) - 1 = 0}$

Support-vector machine (SVM)

- Margin delimitation: $c(\langle w, u \rangle + b) - 1 = 0$

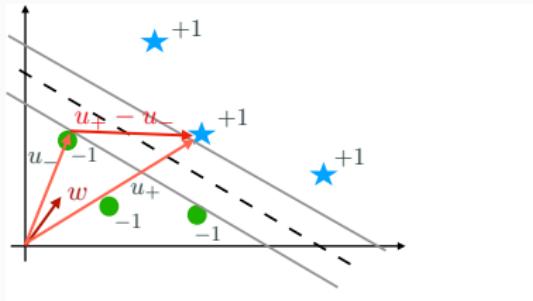


- Objective: Maximize the width of the margin

$$\begin{aligned}\text{Width} &= \frac{1}{\|w\|} \langle u_+ - u_-, w \rangle \\ &= \frac{1}{\|w\|} (1 - b + 1 + b) \\ &= \frac{2}{\|w\|}\end{aligned}$$

Support-vector machine (SVM)

- Margin delimitation: $c(\langle w, u \rangle + b) - 1 = 0$

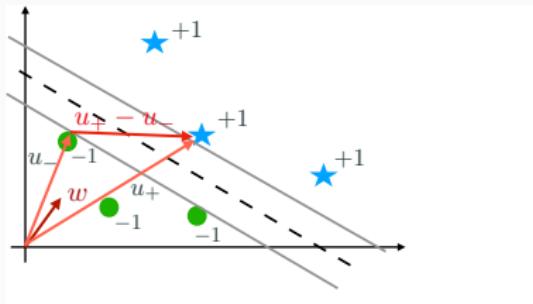


- Objective: Maximize the width of the margin

$$\begin{aligned}\text{Width} &= \frac{1}{\|w\|} \langle u_+ - u_-, w \rangle \\ &= \frac{1}{\|w\|} (1 - b + 1 + b) \\ &= \frac{2}{\|w\|}\end{aligned}$$

Support-vector machine (SVM)

- Margin delimitation: $c \left(\langle w, u \rangle + b \right) - 1 = 0$

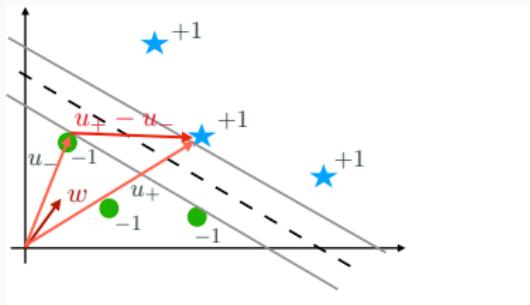


- Objective: Maximize the width of the margin

$$\begin{aligned} \text{Width} &= \frac{1}{\|w\|} \langle u_+ - u_-, w \rangle \\ &= \frac{1}{\|w\|} (1 - b + 1 + b) \quad \Rightarrow \boxed{\min_w \|w\| \quad \text{s.t.} \quad (\forall \ell) c_\ell \left(\langle w, u_\ell \rangle + b \right) \geq 1} \\ &= \frac{2}{\|w\|} \end{aligned}$$

Support-vector machine (SVM)

- Margin delimitation: $c \left(\langle w, u \rangle + b \right) - 1 = 0$



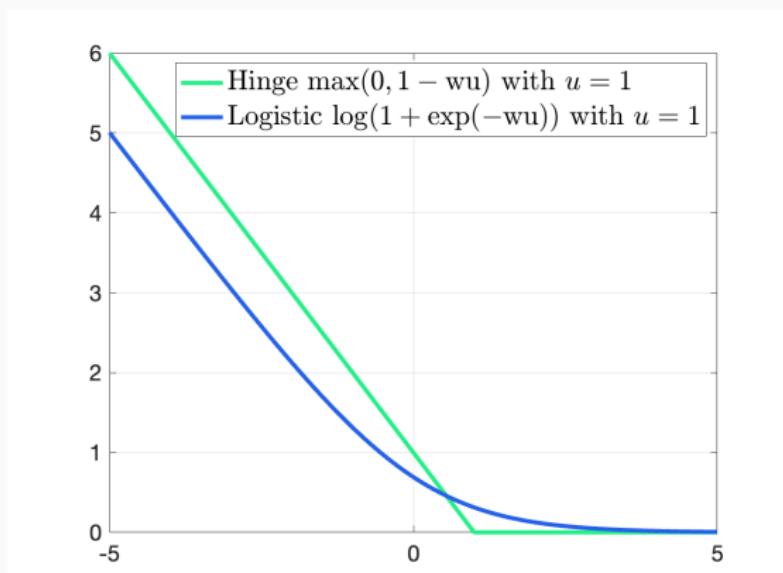
- Objective: Maximize the width of the margin

$$\begin{aligned} \text{Width} &= \frac{1}{\|w\|} \langle u_+ - u_-, w \rangle \\ &= \frac{1}{\|w\|} (1 - b + 1 + b) \quad \Rightarrow \boxed{\min_w \frac{1}{2} \|w\|^2 \text{ s.t. } (\forall \ell) c_\ell \left(\langle w, u_\ell \rangle + b \right) \geq 1} \\ &= \frac{2}{\|w\|} \end{aligned}$$

Losses

- **SVM** : $\min_w \frac{1}{2} \|w\|^2$ s.t. $(\forall \ell) c_\ell \left(\langle w, u_\ell \rangle + b \right) \geq 1$
- **Hinge** + L2 : $\min_w \frac{1}{2} \|w\|^2 + \frac{\lambda}{L} \sum_\ell \max \left(0, 1 - c_\ell (\langle w, u_\ell \rangle + b) \right)$
- **Logistic** + L2 : $\min_w \frac{1}{2} \|w\|^2 + \frac{\lambda}{L} \sum_\ell \log \left(1 + \exp(-c_\ell (\langle w, u_\ell \rangle + b)) \right)$

where $\lambda > 0$.



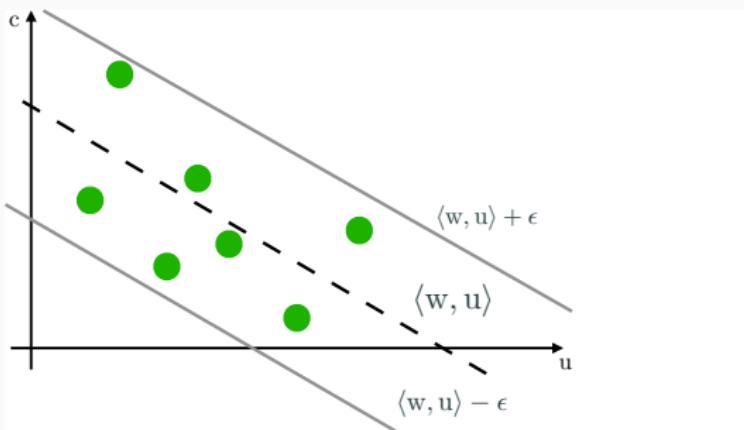
From linear regression to SVR

- Database : dataset $\mathcal{S} = \{(u_\ell, c_\ell) \in \mathbb{R}^N \times \mathbb{R} \mid \ell \in \{1, \dots, L\}\}$.
- **Linear regression:**

$$\min_w \sum_{\ell} (c_\ell - \langle w, u_\ell \rangle)^2$$

- **Support Vector Regression:**

$$\min_w \frac{1}{2} \|w\|^2 \quad \text{s.t.} \quad (\forall \ell) \quad |c_\ell - \langle w, u_\ell \rangle| \leq \epsilon$$



Linearization

- **Database:** $\mathcal{S} = \{(u_\ell, c_\ell) \in \mathcal{H} \times \mathcal{G} \mid \ell \in \{1, \dots, L\}\}$

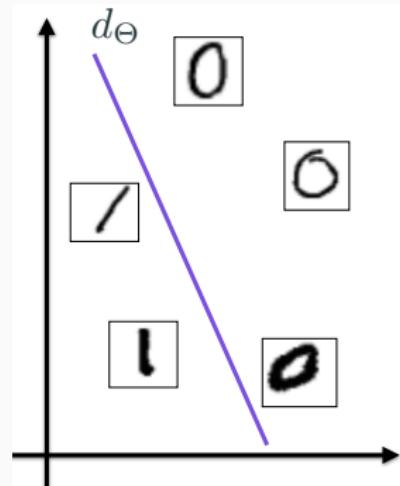
examples: $u_\ell = \boxed{1}$ and $z_\ell = 2$

$u_\ell = \boxed{8}$ and $z_\ell = 9$

- $\varphi(u): \mathbb{R}^N \rightarrow \mathbb{R}^M$: mapping from the input space onto an arbitrary feature space with $M > N$

⇒ **linearization**

examples: convolution networks.
scattering coefficients.



Objective

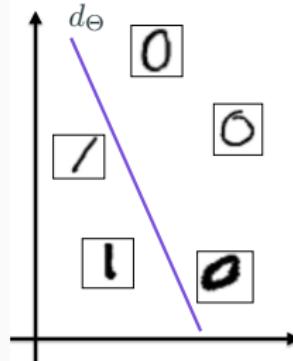
- The predictor relies on K different discriminating functions $D_{\theta^{(k)}} : \mathbb{R}^N \rightarrow \mathbb{R}$:

$$D_{\theta^{(k)}}(u) = (w^{(k)})^\top \varphi(u) + b^{(k)} \quad \text{where} \quad \theta_k = \{w^{(k)}, b^{(k)}\}$$

with $\phi(u) = [\varphi(u)^\top 1]^\top$ and $\Theta = [\underbrace{(w^{(1)})^\top, b^{(1)}}_{\theta^{(1)}}, \dots, \underbrace{(w^{(K)})^\top, b^{(K)}}_{\theta^{(K)}}]^\top$

- The predictor selects the class that best matches an observation

$$d_\Theta(u) = \arg \max_{1 \leq k \leq K} D_{\theta^{(k)}}(u)$$



Objective

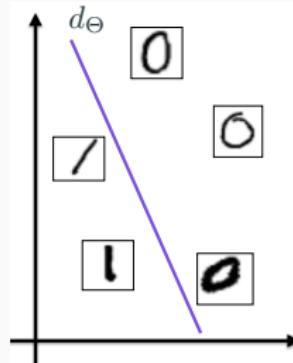
- The predictor relies on K different discriminating functions $D_{\theta^{(k)}} : \mathbb{R}^N \rightarrow \mathbb{R}$:

$$D_{\theta^{(k)}}(u) = (w^{(k)})^\top \varphi(u) + b^{(k)} \Leftrightarrow D_k(u) = (\theta^{(k)})^\top \phi(u)$$

with $\phi(u) = [\varphi(u)^\top 1]^\top$ and $\Theta = [\underbrace{(w^{(1)})^\top, b^{(1)}}_{\theta^{(1)}}, \dots, \underbrace{(w^{(K)})^\top, b^{(K)}}_{\theta^{(K)}}]^\top$

- The predictor selects the class that best matches an observation

$$d_\Theta(u) = \arg \max_{1 \leq k \leq K} D_{\theta^{(k)}}(u)$$



Objective

- Objective of the learning stage: estimate Θ to correctly predict the input-output pair $(u_\ell, c_\ell) \in \mathcal{S}$ for every $\ell \in \{1, \dots, L\}$,

$$d_\Theta(u) = \arg \max_{1 \leq k \leq K} D_{\theta^{(k)}}(u) : \quad c_\ell = \arg \max_{1 \leq k \leq K} (\theta^{(k)})^\top \phi(u_\ell)$$

Objective

- Objective of the learning stage: estimate Θ to correctly predict the input-output pair $(\mathbf{u}_\ell, c_\ell) \in \mathcal{S}$ for every $\ell \in \{1, \dots, L\}$,

$$\begin{aligned} d_\Theta(\mathbf{u}) &= \arg \max_{1 \leq k \leq K} D_{\theta^{(k)}}(\mathbf{u}) : & c_\ell &= \arg \max_{1 \leq k \leq K} (\theta^{(k)})^\top \phi(\mathbf{u}_\ell) \\ && \Leftrightarrow & \max_{k \neq c_\ell} (\theta^{(k)} - \theta^{(c_\ell)})^\top \phi(\mathbf{u}_\ell) < 0 \end{aligned}$$

Objective

- Objective of the learning stage: estimate Θ to correctly predict the input-output pair $(\mathbf{u}_\ell, c_\ell) \in \mathcal{S}$ for every $\ell \in \{1, \dots, L\}$,

$$\begin{aligned} d_\Theta(\mathbf{u}) = \arg \max_{1 \leq k \leq K} D_{\theta^{(k)}}(\mathbf{u}) : \quad c_\ell &= \arg \max_{1 \leq k \leq K} (\theta^{(k)})^\top \phi(\mathbf{u}_\ell) \\ \Leftrightarrow \quad \max_{k \neq c_\ell} (\theta^{(k)} - \theta^{(c_\ell)})^\top \phi(\mathbf{u}_\ell) &< 0 \\ [\text{relax the strict inequality}] \Leftrightarrow \quad \max_{k \neq c_\ell} (\theta^{(k)} - \theta^{(c_\ell)})^\top \phi(\mathbf{u}_\ell) &\leq -\mu_\ell \end{aligned}$$

with $\mu_\ell > 0$

Objective

- Objective of the learning stage: estimate Θ to correctly predict the input-output pair $(\mathbf{u}_\ell, c_\ell) \in \mathcal{S}$ for every $\ell \in \{1, \dots, L\}$,

$$\begin{aligned} d_\Theta(\mathbf{u}) = \arg \max_{1 \leq k \leq K} D_{\theta^{(k)}}(\mathbf{u}) : \quad & c_\ell = \arg \max_{1 \leq k \leq K} (\theta^{(k)})^\top \phi(\mathbf{u}_\ell) \\ \Leftrightarrow & \max_{k \neq c_\ell} (\theta^{(k)} - \theta^{(c_\ell)})^\top \phi(\mathbf{u}_\ell) < 0 \\ \text{[relax the strict inequality]} \Leftrightarrow & \max_{k \neq c_\ell} (\theta^{(k)} - \theta^{(c_\ell)})^\top \phi(\mathbf{u}_\ell) \leq -\mu_\ell \\ \text{[deal with unfeasible constraints]} \Leftrightarrow & \max_{k \neq c_\ell} (\theta^{(k)} - \theta^{(c_\ell)})^\top \phi(\mathbf{u}_\ell) \leq \zeta_\ell - \mu_\ell \end{aligned}$$

with $\mu_\ell > 0$ and $\zeta_\ell \geq 0$.

Multiclass SVM

$$\begin{aligned} & \underset{(\Theta, \xi) \in \mathbb{R}^{(M+1)K} \times \mathbb{R}^L}{\text{minimize}} && \sum_{k=1}^K \|\theta^{(k)}\|_2^2 + \lambda \sum_{\ell=1}^L \xi^{(\ell)} \quad \text{subj. to} \\ & && \left\{ \begin{array}{ll} (\forall \ell \in \{1, \dots, L\}) & \max_{k \neq c_\ell} (\theta^{(k)} - \theta^{(c_\ell)})^\top \phi(u_\ell) \leq \xi_\ell - \mu_\ell, \\ (\forall \ell \in \{1, \dots, L\}) & \xi_\ell \geq 0, \end{array} \right. \end{aligned}$$

or equivalently

$$\underset{\Theta \in \mathbb{R}^{(M+1)K}}{\text{minimize}} \sum_{k=1}^K \|\theta^{(k)}\|_2^2 + \lambda \sum_{\ell=1}^L \max \left\{ 0, \mu_\ell + \max_{k \neq c_\ell} (\theta^{(k)} - \theta^{(c_\ell)})^\top \phi(u_\ell) \right\}.$$

Alternative to standard SVM data-term

$$h(\Theta) = \sum_{\ell=1}^L \max \left\{ 0, \mu_\ell + \max_{k \neq c_\ell} (\theta^{(k)} - \theta^{(c_\ell)})^\top \phi(u_\ell) \right\}.$$

- Multiclass SVM [Blondel et al.]

$$h(\Theta) = \sum_{\ell=1}^L \sum_{k \neq c_\ell} \left(\max \left\{ 0, \mu_\ell + (\theta^{(k)} - \theta^{(c_\ell)})^\top \phi(u_\ell) \right\} \right)^2$$

- Multinomial logistic regression [Krishnapuram et al.]

$$h(\Theta) = \sum_{\ell=1}^L \log \left(1 + \sum_{k \neq c_\ell} \exp \left\{ \mu_\ell + (\theta^{(k)} - \theta^{(c_\ell)})^\top \phi(u_\ell) \right\} \right)$$

- “one-vs-all” strategy binary SVM [Laporte et al.] with \tilde{c}_ℓ being equal to 1 if $c_\ell = k$, and -1 otherwise $h(\Theta) = \sum_{\ell=1}^L \left(\max \left\{ 0, \mu_\ell + \tilde{c}_\ell (\theta^{(k)})^\top \phi(u_\ell) \right\} \right)^2$

Example: Supervised learning

- **Database:** $\mathcal{S} = \{(u_\ell, c_\ell) \in \mathcal{H} \times \mathcal{G} \mid \ell \in \{1, \dots, L\}\}$
- **Goal:** Learn a prediction function $d_\Theta: \mathcal{H} \rightarrow \mathcal{G}$
- Learning procedure relies on a minimization problem:

$$\underset{\Theta}{\text{minimize}} \quad \underbrace{\frac{1}{L} \sum_{\ell=1}^L F(c_\ell, d_\Theta(u_\ell))}_{\text{Data-term}} + \lambda \underbrace{R(\Theta)}_{\text{Prior}}$$

Example: Supervised learning

- **Database:** $\mathcal{S} = \{(u_\ell, c_\ell) \in \mathcal{H} \times \mathcal{G} \mid \ell \in \{1, \dots, L\}\}$
- **Goal:** Learn a prediction function $d_\Theta: \mathcal{H} \rightarrow \mathcal{G}$
- Learning procedure relies on a minimization problem:

$$\underset{\Theta}{\text{minimize}} \quad \underbrace{\frac{1}{L} \sum_{\ell=1}^L F(c_\ell, d_\Theta(u_\ell))}_{\text{Data-term}} + \lambda \underbrace{R(\Theta)}_{\text{Prior}}$$

- **Linear predictor:** $d_\Theta(u) = \Theta^\top u$

○ Ridge regression: $\underset{\Theta \in \mathcal{H}}{\text{minimize}} \quad \frac{1}{L} \sum_{\ell=1}^L (c_\ell - \Theta^\top u_\ell)^2 + \lambda \|\Theta\|_2^2$

○ Logistic classification: $\underset{\Theta \in \mathcal{H}}{\text{minimize}} \quad \frac{1}{L} \sum_{\ell=1}^L \log(1 + e^{-c_\ell \Theta^\top u_\ell}) + \lambda \|\Theta\|_2^2$

Example: Supervised learning

- **Database:** $\mathcal{S} = \{(u_\ell, c_\ell) \in \mathcal{H} \times \mathcal{G} \mid \ell \in \{1, \dots, L\}\}$
- **Goal:** Learn a prediction function $d_\Theta: \mathcal{H} \rightarrow \mathcal{G}$
- Learning procedure relies on a minimization problem:

$$\underset{\Theta}{\text{minimize}} \quad \underbrace{\frac{1}{L} \sum_{\ell=1}^L F(c_\ell, d_\Theta(u_\ell))}_{\text{Data-term}} + \lambda \underbrace{R(\Theta)}_{\text{Prior}}$$

- **Linear predictor:** $d_\Theta(u) = \Theta^\top u$

$$\odot \text{ Ridge regression: } \underset{\Theta \in \mathcal{H}}{\text{minimize}} \quad \frac{1}{L} \sum_{\ell=1}^L (c_\ell - \Theta^\top u_\ell)^2 + \lambda \|\Theta\|_2^2$$

$$\odot \text{ Logistic classification: } \underset{\Theta \in \mathcal{H}}{\text{minimize}} \quad \frac{1}{L} \sum_{\ell=1}^L \log(1 + e^{-c_\ell \Theta^\top u_\ell}) + \lambda \|\Theta\|_2^2$$

Example: Supervised learning

- **Database:** $\mathcal{S} = \{(u_\ell, c_\ell) \in \mathcal{H} \times \mathcal{G} \mid \ell \in \{1, \dots, L\}\}$
- **Goal:** Learn a prediction function $d_\Theta: \mathcal{H} \rightarrow \mathcal{G}$
- Learning procedure relies on a minimization problem:

$$\underset{\Theta}{\text{minimize}} \quad \underbrace{\frac{1}{L} \sum_{\ell=1}^L F(c_\ell, d_\Theta(u_\ell))}_{\text{Data-term}} + \lambda \underbrace{R(\Theta)}_{\text{Prior}}$$

- **Linear predictor:** $d_\Theta(u) = \Theta^\top u$
⇒ can be extended to $d_\Theta(u) = \Theta^\top \phi(u)$ (e.g. ϕ scattering transform)
 - Ridge regression: $\underset{\Theta \in \mathcal{H}}{\text{minimize}} \quad \frac{1}{L} \sum_{\ell=1}^L (c_\ell - \Theta^\top u_\ell)^2 + \lambda \|\Theta\|_2^2$
 - Logistic classification: $\underset{\Theta \in \mathcal{H}}{\text{minimize}} \quad \frac{1}{L} \sum_{\ell=1}^L \log(1 + e^{-c_\ell \Theta^\top u_\ell}) + \lambda \|\Theta\|_2^2$

Example: Supervised learning

- **Database:** $\mathcal{S} = \{(u_\ell, c_\ell) \in \mathcal{H} \times \mathcal{G} \mid \ell \in \{1, \dots, L\}\}$
- **Goal:** Learn a prediction function $d_\Theta: \mathcal{H} \rightarrow \mathcal{G}$
- Learning procedure relies on a minimization problem:

$$\underset{\Theta}{\text{minimize}} \quad \underbrace{\frac{1}{L} \sum_{\ell=1}^L F(c_\ell, d_\Theta(u_\ell))}_{\text{Data-term}} + \lambda \underbrace{R(\Theta)}_{\text{Prior}}$$

- **Linear predictor:** $d_\Theta(u) = \Theta^\top u$

○ Lasso:

$$\underset{\Theta \in \mathcal{H}}{\text{minimize}} \quad \frac{1}{L} \sum_{\ell=1}^L (c_\ell - \Theta^\top u_\ell)^2 + \lambda \|\Theta\|_1$$

○ Sparse logistic classification:

$$\underset{\Theta \in \mathcal{H}}{\text{minimize}} \quad \frac{1}{L} \sum_{\ell=1}^L \log(1 + e^{-c_\ell \Theta^\top u_\ell}) + \lambda \|\Theta\|_1$$

○ SVM classification:

$$\underset{\Theta \in \mathcal{H}}{\text{minimize}} \quad \frac{1}{L} \sum_{\ell=1}^L \max(0, 1 - c_\ell \Theta^\top u_\ell) + \lambda \|\Theta\|_2^2$$

Example: Supervised learning

- **Database:** $\mathcal{S} = \{(u_\ell, c_\ell) \in \mathcal{H} \times \mathcal{G} \mid \ell \in \{1, \dots, L\}\}$
- **Goal:** Learn a prediction function $d_\Theta: \mathcal{H} \rightarrow \mathcal{G}$
- Learning procedure relies on a minimization problem:

$$\underset{\Theta}{\text{minimize}} \quad \underbrace{\frac{1}{L} \sum_{\ell=1}^L F(c_\ell, d_\Theta(u_\ell))}_{\text{Data-term}} + \lambda \underbrace{R(\Theta)}_{\text{Prior}}$$

- **Linear predictor:** $d_\Theta(u) = \Theta^\top u \Rightarrow \text{Convex non-smooth problems}$

○ Lasso:

$$\underset{\Theta \in \mathcal{H}}{\text{minimize}} \quad \frac{1}{L} \sum_{\ell=1}^L (c_\ell - \Theta^\top u_\ell)^2 + \lambda \|\Theta\|_1$$

○ Sparse logistic classification:

$$\underset{\Theta \in \mathcal{H}}{\text{minimize}} \quad \frac{1}{L} \sum_{\ell=1}^L \log(1 + e^{-c_\ell \Theta^\top u_\ell}) + \lambda \|\Theta\|_1$$

○ SVM classification:

$$\underset{\Theta \in \mathcal{H}}{\text{minimize}} \quad \frac{1}{L} \sum_{\ell=1}^L \max(0, 1 - c_\ell \Theta^\top u_\ell) + \lambda \|\Theta\|_2^2$$

Experiment 2

MNIST database

- ◆ $N = 28 \times 28$
- ◆ $K = 10$
- ◆ 60000 training images
- ◆ 10000 test images

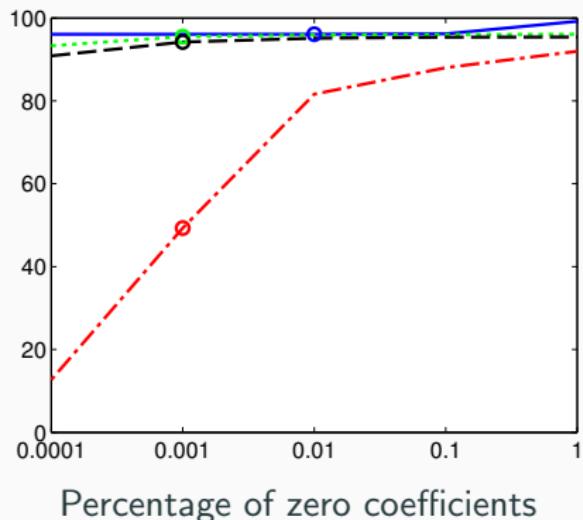
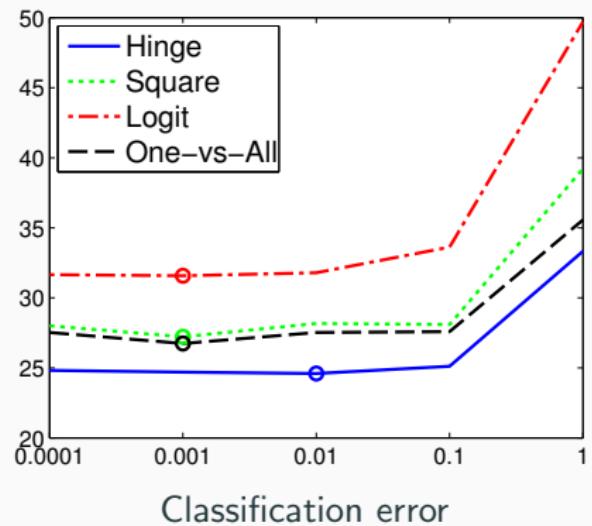


Scattering convolution network [Bruna & Mallat, 2013]

- ◆ 2 wavelet layers
- ◆ 4 scales DWT
- ◆ Feature mapping: $\phi: \mathbb{R}^{28 \times 28} \mapsto \mathbb{R}^{14 \times 14 \times 81}$

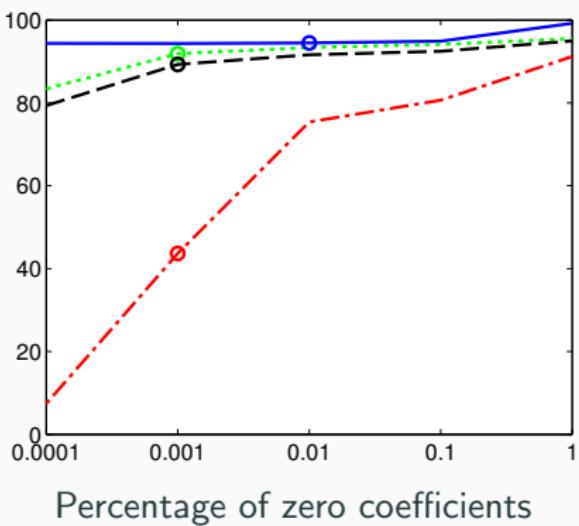
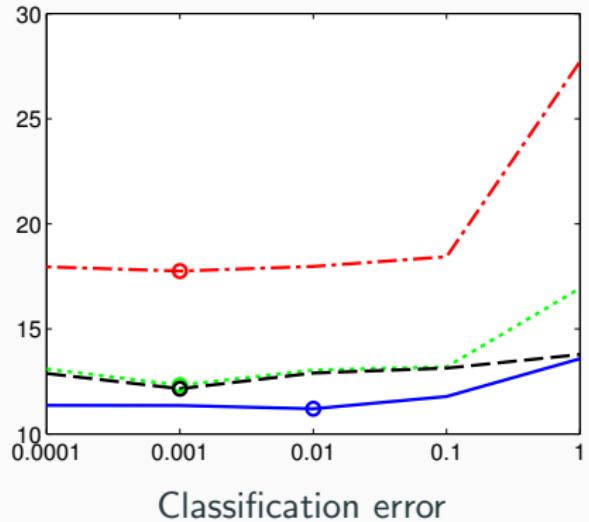
Experiment 2

Classification results for $L/K = 3$



Experiment 2

Classification results for $L/K = 10$



Example: Supervised learning

- **Database:** $\mathcal{S} = \{(u_\ell, c_\ell) \in \mathcal{H} \times \mathcal{G} \mid \ell \in \{1, \dots, L\}\}$
- **Goal:** Learn a prediction function $d_\Theta: \mathcal{H} \rightarrow \mathcal{G}$
- Learning procedure relies on a **minimization problem**:

$$\underset{\Theta}{\text{minimize}} \quad \underbrace{\frac{1}{L} \sum_{\ell=1}^L F(c_\ell, d_\Theta(u_\ell))}_{\text{Data-term}} + \lambda \underbrace{R(\Theta)}_{\text{Prior}}$$

- **Neural network predictor:**

$$d_\Theta(u) = \eta^{[K]}(W^{[K]}\eta^{[K-1]}(W^{[K-1]}\dots\eta^{[2]}(W^{[2]}\eta^{[1]}(W^{[1]}u))\dots))$$

- Linear operators: $W^{[1]}, W^{[2]}, \dots, W^{[K]}$
- Activation functions: $\eta^{[1]}, \eta^{[2]}, \dots, \eta^{[K]}$

Example: Supervised learning

- **Database:** $\mathcal{S} = \{(u_\ell, c_\ell) \in \mathcal{H} \times \mathcal{G} \mid \ell \in \{1, \dots, L\}\}$
- **Goal:** Learn a prediction function $d_\Theta: \mathcal{H} \rightarrow \mathcal{G}$
- Learning procedure relies on a **minimization problem**:

$$\underset{\Theta}{\text{minimize}} \quad \underbrace{\frac{1}{L} \sum_{\ell=1}^L F(c_\ell, d_\Theta(u_\ell))}_{\text{Data-term}} + \lambda \underbrace{R(\Theta)}_{\text{Prior}}$$

- **Neural network predictor:** \Rightarrow Non-convex problems

$$d_\Theta(u) = \eta^{[K]}(W^{[K]}\eta^{[K-1]}(W^{[K-1]}\dots\eta^{[2]}(W^{[2]}\eta^{[1]}(W^{[1]}u))\dots))$$

- Linear operators: $W^{[1]}, W^{[2]}, \dots, W^{[K]}$
- Activation functions: $\eta^{[1]}, \eta^{[2]}, \dots, \eta^{[K]}$