

A Kind of Multi-Output Regression

Risk minimization for function-valued regression:

- \mathcal{X} input space (\mathbb{R}^d), Θ parameter space ($\subset \mathbb{R}$), \mathcal{Y} output space ($\subset \mathbb{R}$).
- Hypothesis space $\mathcal{H} \subset \mathcal{F}(\mathcal{X}; \mathcal{F}(\Theta; \mathcal{Y}))$, i.e. $h(x) \in \mathcal{F}(\Theta; \mathcal{Y})$.
- Parametrized cost $v: \Theta \times \mathcal{Y} \times \mathcal{Y} \rightarrow \mathbb{R}$.
- Local loss $V(y, h(x)) := \int_{\Theta} v(\theta, y, h(x)(\theta)) d\mu(\theta)$.

Minimizing population risk:

$$\arg \min_{h \in \mathcal{H}} R(h) := \mathbf{E}_{X, Y} [V(Y, h(X))]. \quad (1)$$

⇒ Extension of Multi-Task Learning to an infinite number of tasks [1].

Two examples

Quantile Regression (QR): Given $X, Y \in \mathcal{X} \times \mathcal{Y}$ random variables, estimate the quantile function of the conditional distribution $\mathbf{P}_{Y|X}$:

$$q(x)(\theta) = \inf \{t \in \mathcal{Y}, \mathbf{P}_{Y|X=x}[Y \leq t] \geq \theta\} \quad \forall (x, \theta) \in \mathcal{X} \times (0, 1). \quad (2)$$

Pinball loss:

$$v(\theta, y, h(x)) = |\theta - \mathbb{1}_{\mathbb{R}_-}(y - h(x))| |y - h(x)|. \quad (3)$$

Proposition. q defined in (2) minimizes (1) for the pinball loss (3).

Cost-Sensitive Classification (CSC): Support Vector Machine with asymmetric loss function

$$v(\theta, y, h(x)) = \left| \frac{\theta + 1}{2} - \mathbb{1}_{\{-1\}}(y) \right| |1 - y h(x)|_+.$$

The value of θ influences how the different classes are penalized.

Sampled Empirical Risk

Approximate expectation over $\mathbf{P}_{X, Y}$ and \int_{Θ}

- $(x_i, y_i)_{i=1}^n \stackrel{\text{i.i.d.}}{\sim} \mathbf{P}_{X, Y}$
- $(\theta_j)_{j=1}^m \sim \mu$ (Quasi-Monte Carlo)

Sampled empirical risk:

$$\tilde{R}_S(h) := \frac{1}{nm} \sum_{i=1}^n \sum_{j=1}^m v(\theta_j, y_i, h(x_i)(\theta_j)).$$

Regularized problem:

$$\arg \min_{h \in \mathcal{H}} \tilde{R}_S(h) + \lambda \Omega(h). \quad (4)$$

Vector-Valued RKHSs

Natural extension of RKHS for modelling outputs in any Hilbert space.

- $k_X: \mathcal{X} \times \mathcal{X} \rightarrow \mathbb{R}$ and $k_{\Theta}: \Theta \times \Theta \rightarrow \mathbb{R}$ two scalar-valued kernels.
- Operator-valued kernel $K(x, z) = k_X(x, z) I_{\mathcal{H}_{k_{\Theta}}}$ associated to \mathcal{H}_K a space of function-valued functions.
- $\mathcal{H}_K = \overline{\text{span}} \{K(\cdot, x)f \mid x \in \mathcal{X}, f \in \mathcal{H}_{k_{\Theta}}\} \cong \mathcal{H}_{k_X} \otimes \mathcal{H}_{k_{\Theta}}$.
- Hilbert norm $\frac{1}{2} \|h\|_{\mathcal{H}_K}^2$ as regularizer $\Omega(h)$.

Optimization

Proposition (Representer). If $\forall \theta \in \Theta, v(\theta, \cdot, \cdot)$ is proper lower semicontinuous with respect to its second argument, (4) has a unique solution $h^* \in \mathcal{H}_K$, and $\exists (\alpha_{ij})_{i,j=1}^{n,m} \in \mathbb{R}^{n \times m}$ such that $\forall (x, \theta) \in \mathcal{X} \times \Theta$

$$h^*(x)(\theta) = \sum_{i=1}^n \sum_{j=1}^m \alpha_{ij} k_X(x, x_i) k_{\Theta}(\theta, \theta_j).$$

- Solution shaped by k_X and k_{Θ} (Gaussian, Laplacian, ...)
- Infinite-dimensional problem \Rightarrow size $n \cdot m$
- In practice, solved via smoothing $v + L$ -BFGS.

Excess Risk Guarantees

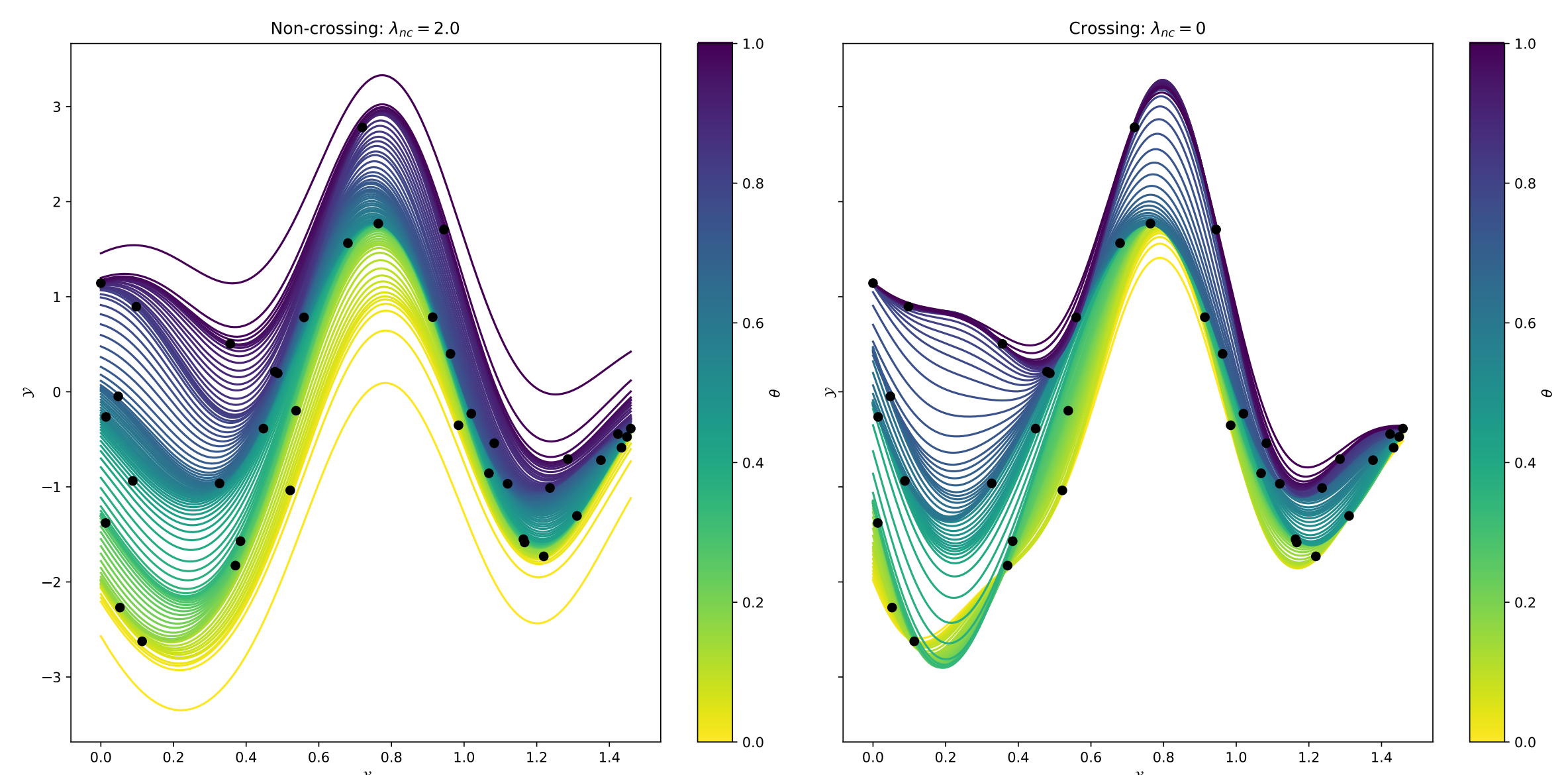
Framework of vv-RKHS allows for proper analysis [2], tradeoff n/m

$$R(h^*) \leq \tilde{R}_S(h^*) + \mathcal{O}_{\mathbf{P}_{X, Y}} \left(\frac{1}{\sqrt{\lambda n}} \right) + \mathcal{O} \left(\frac{\log(m)}{\sqrt{\lambda m}} \right).$$

Numerical Experiments

QR: Continuous model \Rightarrow new non-crossing constraint:

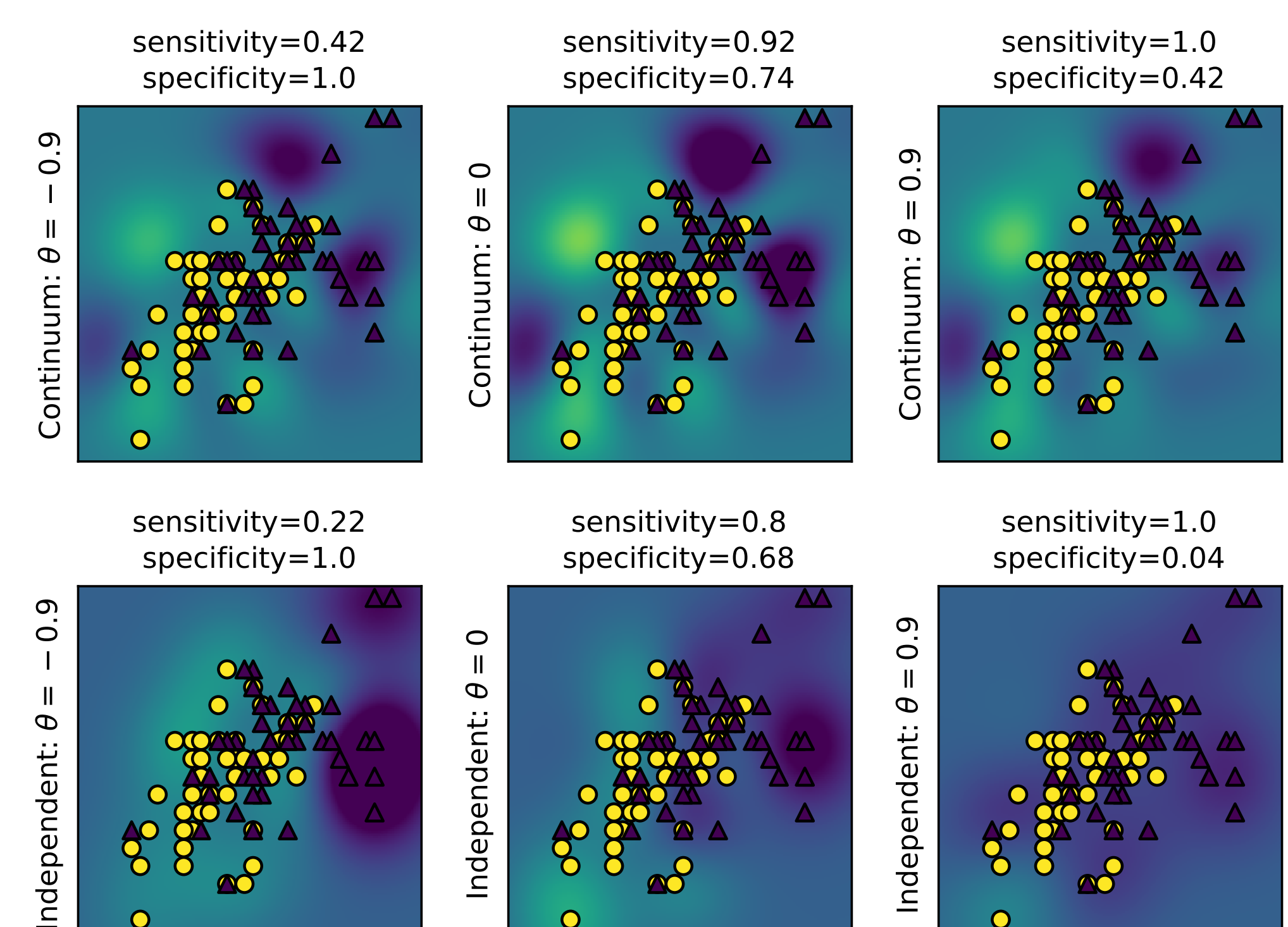
$$\tilde{\Omega}_{nc}(h) = \frac{\lambda_{nc}}{nm} \sum_{i=1}^n \sum_{j=1}^m \left| -\frac{\partial h}{\partial \theta}(x_i)(\theta_j) \right|_+.$$



Left: strong non-crossing penalty ($\lambda_{nc} = 2$). Right: no non-crossing penalty ($\lambda_{nc} = 0$). The plots show 100 quantiles of the continuum learned, linearly spaced between 0 (yellow) and 1 (purple).

⇒ Matches state of the art [3] on 20 UCI datasets

CSC: Improved performance:



Iris dataset. Top: infinite learning; bottom: independent learning for $\theta \in \{-0.9, 0, 0.9\}$.

Code available: <https://bitbucket.org/RomainBrault/itl/>