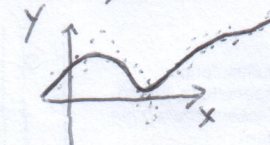


Support Vector Regression

Data: $\left\{ \left(\underset{\substack{\uparrow \\ \mathbb{R}^N}}{\underline{x}^{(\alpha)}}, \underset{\substack{\uparrow \\ \mathbb{R}}}{y_T^{(\alpha)}} \right) \right\}$

nonlinear regression:



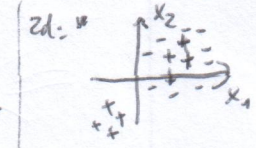
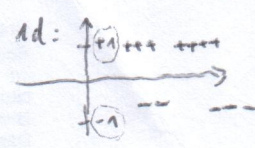
linear regression:



Model: linear model $y(\underline{x}) = \underline{w}^T \underline{x} + b$

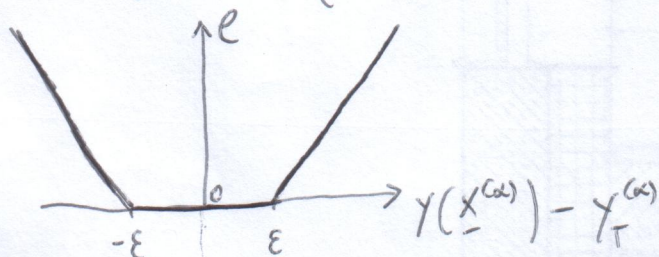
\uparrow normal vector of hyper plane
 \uparrow offset of hyper plane

reminder: classification problems in contrast have labels with only finitely many possible values, e.g., $y_T^{(\alpha)} \in \{-1, +1\}$ for binary classification



a) ϵ -sensitive cost function

$$e(\underline{x}^{(\alpha)}, y_T^{(\alpha)}) = \max(0, |y(\underline{x}^{(\alpha)}) - y_T^{(\alpha)}| - \epsilon)$$



\Rightarrow cost function parametrized by $\epsilon \geq 0$

\Rightarrow no error for model output $y(\underline{x}^{(\alpha)}) \in [y_T^{(\alpha)} - \epsilon, y_T^{(\alpha)} + \epsilon]$

\Rightarrow linear error outside this interval

$y_T^{(\alpha)} - y(\underline{x}^{(\alpha)}) > \epsilon$ $y(\underline{x}^{(\alpha)}) - y_T^{(\alpha)} > \epsilon$

$0 < y_T^{(\alpha)} - y(\underline{x}^{(\alpha)}) \leq \epsilon$ $0 \leq y(\underline{x}^{(\alpha)}) - y_T^{(\alpha)} \leq \epsilon$

b) Derivation of the primal optimization problem for the " ϵ -SVR":

i) simplifying assumptions (later relaxed): ⁽¹⁾ perfect regression solution (0 error) and ⁽²⁾ linear regression problem

min $\frac{1}{2} \|\underline{w}\|^2$

s.t. $y(\underline{x}^{(\alpha)}) - y_T^{(\alpha)} \leq \epsilon$

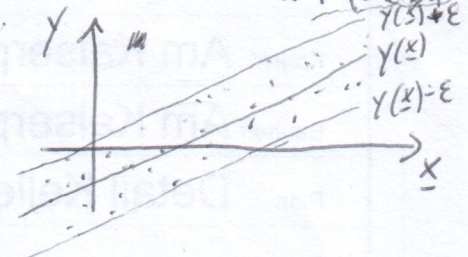
$y_T^{(\alpha)} - y(\underline{x}^{(\alpha)}) \leq \epsilon$ for $\alpha = 1, \dots, P$

\nwarrow parameters defining the zero error interval

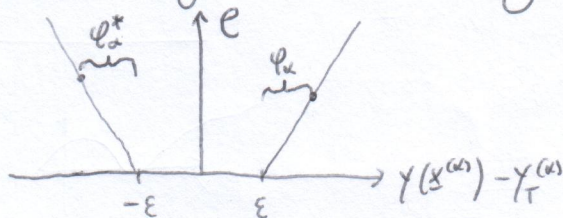
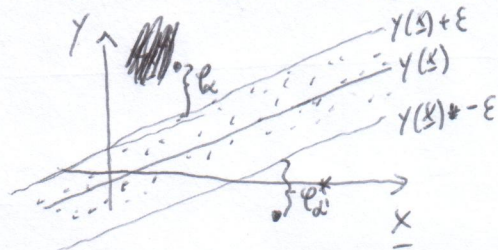
\Leftrightarrow min $\frac{1}{2} \|\underline{w}\|^2$

s.t. $\underline{w}^T \underline{x}^{(\alpha)} + b - y_T^{(\alpha)} \leq \epsilon$

$y_T^{(\alpha)} - \underline{w}^T \underline{x}^{(\alpha)} - b \leq \epsilon$



b) ii) extension: allow some error to reduce overfitting to noise in training data (cf. C-SVM)



ϕ_α : "positive" slack variable
 ϕ_α^* : "negative" slack variable

primal optimization problem ("ε-SVR"):

$$\begin{aligned} \min_{\underline{w}, b, \{\phi_\alpha\}, \{\phi_\alpha^*\}} \quad & \frac{1}{2} \|\underline{w}\|^2 + C \frac{1}{P} \sum_{\alpha=1}^P (\phi_\alpha + \phi_\alpha^*) \\ \text{s.t.} \quad & \underline{w}^T \underline{x}^{(\alpha)} + b - y_T^{(\alpha)} \leq \epsilon + \phi_\alpha \\ & y_T^{(\alpha)} - \underline{w}^T \underline{x}^{(\alpha)} - b \leq \epsilon + \phi_\alpha^* \\ & \phi_\alpha \geq 0, \phi_\alpha^* \geq 0 \end{aligned}$$

for $\alpha=1, \dots, P$

← new constraints

ε-SVR

Derivation of the primal problem for the "ν-SVR"

c) Choosing ε properly can be difficult as it depends on the unknown data noise
 • extension: make ε a ~~scalar~~ variable with respect to which we optimize:

$$\begin{aligned} \min_{\underline{w}, b, \{\phi_\alpha\}, \{\phi_\alpha^*\}, \epsilon} \quad & \frac{1}{2} \|\underline{w}\|^2 + C \left[\underset{\substack{\uparrow \\ \text{new parameter}}}{\nu \epsilon} + \frac{1}{P} \sum_{\alpha=1}^P (\phi_\alpha + \phi_\alpha^*) \right] \\ \text{s.t.} \quad & \underline{w}^T \underline{x}^{(\alpha)} + b - y_T^{(\alpha)} \leq \epsilon + \phi_\alpha \\ & y_T^{(\alpha)} - \underline{w}^T \underline{x}^{(\alpha)} - b \leq \epsilon + \phi_\alpha^* \\ & \phi_\alpha \geq 0, \phi_\alpha^* \geq 0, \epsilon \geq 0 \end{aligned}$$

for $\alpha=1, \dots, P$

← new constraint

primal optimization
 problem of the
 ν-SVR

- it can be shown that $\nu \leq \frac{\pm 5\nu}{P}$, i.e., ν lower bounds the fraction of support vectors (on $\pm \epsilon$ interval border or in $\pm \epsilon$ error region)
- ν-SVR is standard SV-based regression method when combined with a suitable kernel (e.g. RBF)
- ν has still to be chosen (e.g. cross-validation based hyperparam. selection)

d) Lagrangian:

$$\mathcal{L}(\underbrace{\underline{w}, b, \{\varphi_\alpha\}, \{\varphi_\alpha^*\}}_{\text{primal variables}}, \underbrace{\varepsilon, \{\lambda_\alpha\}, \{\lambda_\alpha^*\}, \{\eta_\alpha\}, \{\eta_\alpha^*\}, \delta}_{\text{dual variables (Lagrange multipliers)}})$$

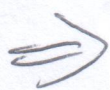
$$\begin{aligned} &= \frac{1}{2} \|\underline{w}\|^2 + C \left[\nu \varepsilon + \frac{1}{P} \sum_{\alpha=1}^P (\varphi_\alpha + \varphi_\alpha^*) \right] \\ &\quad - \sum_{\alpha=1}^P \lambda_\alpha \left(\varphi_\alpha + \varepsilon + y_T^{(\alpha)} - \underline{w}^T \underline{x}^{(\alpha)} - b \right) \\ &\quad - \sum_{\alpha=1}^P \lambda_\alpha^* \left(\varphi_\alpha^* + \varepsilon + \underline{w}^T \underline{x}^{(\alpha)} + b - y_T^{(\alpha)} \right) \\ &\quad - \sum_{\alpha=1}^P \eta_\alpha \varphi_\alpha - \sum_{\alpha=1}^P \eta_\alpha^* \varphi_\alpha^* - \delta \varepsilon \end{aligned}$$

$$\text{s.t. } \lambda_\alpha \geq 0, \lambda_\alpha^* \geq 0, \eta_\alpha \geq 0, \eta_\alpha^* \geq 0, \delta \geq 0 \quad \forall \alpha$$

in the homework you will see/show that

both ^{ai)} the dual ~~pr~~ optimization problem as well as ii) the model ~~to~~ evaluation ^{for} new data \underline{x} , i.e. $y(\underline{x})$, depends on the data exclusively via scalar products $\underline{x}^{(\alpha)T} \underline{x}^{(\beta)}$ [and $\underline{x}^{(\alpha)T} \underline{x}$ respectively]

e.g.
final hyper-
param opt:
via
C, ν , λ



~~kernel~~ kernel usage: replace these scalar products by kernel evaluation $k(\underline{x}^{(\alpha)}, \underline{x}^{(\beta)})$ which then implicitly does a linear regression on the ~~nonlinearly separated~~ data set in feature space, $\{k(\phi(\underline{x}^{(\alpha)}), y_T^{(\alpha)})\}$, i.e. \underline{x} ~~is~~ operated by nonlinear transform $\phi(\underline{x})$ with $\phi(\underline{x}^{(\alpha)T})^T \phi(\underline{x}^{(\beta)}) = k(\underline{x}^{(\alpha)}, \underline{x}^{(\beta)})$

yields
nonlinear regions
in data space
⇕