

Regularization

a) Approach to overfitting reduction: L_1/L_2 regularization.

if optimization of E^T yields too large $\hat{E}^G \Rightarrow$ regularize by optimizing risk $R(\underline{w}) = E^T(\underline{w}) + \lambda E^R(\underline{w})$ instead!

\uparrow training error \uparrow regularization strength (hyperparameter) \uparrow regularization term, $\lambda > 0$

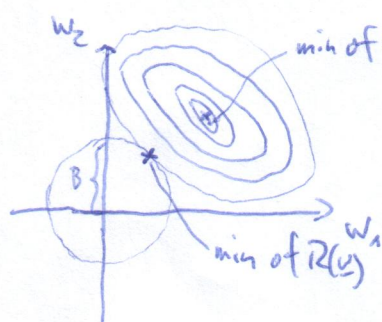
L_2 norm regularization: "weight decay": $E^R(\underline{w}) = \frac{1}{2P} \|\underline{w}\|_2^2 = \frac{1}{2P} \sum_{i=1}^d w_i^2$

\rightarrow standard regularizer

\rightarrow differentiable (gradient descent straight): $\frac{\partial R}{\partial \underline{w}} = \frac{\partial E^T}{\partial \underline{w}} + \lambda \frac{\partial E^R}{\partial \underline{w}}$, $\frac{\partial E^R}{\partial \underline{w}} = \frac{1}{P} \underline{w}$

\rightarrow robust to noise (distributed model: redundancy in weights allowed)

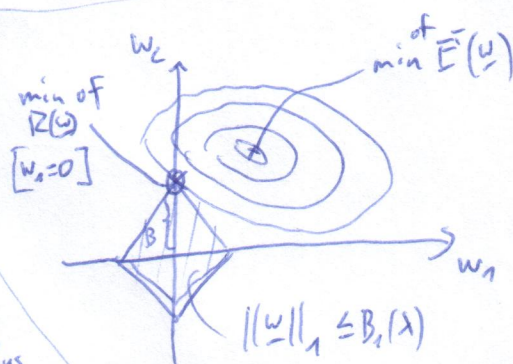
\rightarrow non-sparse



since $\min_{\underline{w}} R(\underline{w})$

$$E) \begin{cases} \min_{\underline{w}} E^T \\ \text{s.t. } \|\underline{w}\|_2 \leq B_2(\lambda) \end{cases}$$

\uparrow ball of radius B_2 (dep. on λ)



L_1 norm regularization: "Lasso" / "sparsify": $E^R(\underline{w}) = \frac{1}{P} \|\underline{w}\|_1 = \frac{1}{P} \sum_{i=1}^d |w_i|$

\rightarrow sparse (good for interpretable models, e.g. in medical applications)

\rightarrow non-differentiable at $w_i = 0 \forall i \Rightarrow$ gradient descent not straight applicable (but solved via subgradient / least angle regression)

\rightarrow more sensitive to ~~noise~~ noise since much less redundancies

b) Analytical solution for linear neuron of quadratic cost function and L_2 regularization.

cost function $R(\underline{w}) = \underbrace{\frac{1}{P} \sum_{\alpha=1}^P \frac{1}{2} (\underline{w}^T \underline{x}^{(\alpha)} - y_T^{(\alpha)})^2}_{\text{quadratic cost of data pt. } \alpha} + \lambda \frac{1}{2P} \underline{w}^T \underline{w}$

Gradient $\underline{g} = \frac{1}{P} (\underline{X} \underline{X}^T \underline{w} - \underline{X} \underline{y}^T) + \frac{\lambda}{P} \underline{w} \stackrel{\underline{g}=0}{\Rightarrow} \underline{w}^* = (\underline{X} \underline{X}^T + \lambda \underline{I})^{-1} \underline{X} \underline{y}^T$

Hessian $\underline{H} = \frac{1}{P} (\underline{X} \underline{X}^T + \lambda \underline{I})$ still positive definite $\Rightarrow \underline{w}^*$ unique global minimum

$\underline{X} = (\underline{x}^{(1)}, \dots, \underline{x}^{(P)})$ data columnwise