

## CHEAT SHEET

# Regularizers

When you look at regularizers, it helps to change the formulation of the optimization problem to obtain a better geometric intuition:

$$\min_{\mathbf{w}, b} \sum_{i=1}^n \ell(h_{\mathbf{w}}(\mathbf{x}_i)) + \lambda r(\mathbf{w}) \iff \min_{\mathbf{w}, b} \sum_{i=1}^n \ell(h_{\mathbf{w}}(\mathbf{x}_i))$$

subject to:  $r(\mathbf{w}) \leq B$

| Regularizers          |   | Details  |
|-----------------------|---|--|
| $l_2$ -Regularization | $r(\mathbf{w}) = \mathbf{w}^\top \mathbf{w} = \ \mathbf{w}\ _2^2 = \sum_{i=1}^d [\mathbf{w}]_i^2$ | <ul style="list-style-type: none"> <li>• ADVANTAGE: Strictly Convex</li> <li>• ADVANTAGE: Differentiable</li> <li>• DISADVANTAGE: Uses weights on all features, i.e. relies on all features to some degree (ideally we would like to avoid this) - these are known as Dense Solutions.</li> </ul>                |
| $l_1$ -Regularization | $r(\mathbf{w}) = \ \mathbf{w}\ _1 = \sum_{i=1}^d  [\mathbf{w}]_i $                                | <ul style="list-style-type: none"> <li>• Convex (but not strictly)</li> <li>• DISADVANTAGE: Not differentiable at 0</li> <li>• Effect: Sparse</li> </ul>   |
| $l_p$ -Norm           | $r(\mathbf{w}) = \ \mathbf{w}\ _p^p = \sum_{i=1}^d [\mathbf{w}]_i^p$                              | <ul style="list-style-type: none"> <li>• Often <math>0 &lt; p \leq 1</math></li> <li>• DISADVANTAGE: Non-convex</li> <li>• ADVANTAGE: Very sparse solutions</li> <li>• Initialization dependent</li> <li>• VERY sparse solutions (compared to <math>l_1</math> norm) if <math>0 &lt; p \leq 1</math>.</li> </ul> |