

Supervised Learning :: CHEAT SHEET

Regularization

Regularization is an effective technique to reduce overfitting.

$$\mathcal{R}_{\text{reg}}(f) = \mathcal{R}_{\text{emp}}(f) + \lambda \cdot J(f) = \sum_{i=1}^n L\left(y^{(i)}, f\left(\mathbf{x}^{(i)}\right)\right) + \lambda \cdot J(f)$$

- $J(f)$: **complexity penalty, roughness penalty** or **regularizer**
- $\lambda \geq 0$: **complexity control** parameter
- The higher λ , the more we penalize complexity
- $\lambda = 0$: We just do simple ERM; $\lambda \rightarrow \infty$: we don't care about loss, models become as "simple" as possible
- λ is hard to set manually and is usually selected via CV
- As for \mathcal{R}_{emp} , \mathcal{R}_{reg} and J are often defined in terms of θ :

$$\mathcal{R}_{\text{reg}}(\theta) = \mathcal{R}_{\text{emp}}(\theta) + \lambda \cdot J(\theta)$$

Ridge Regression

Use L2 penalty in linear regression:

$$\begin{aligned} \hat{\theta}_{\text{ridge}} &= \arg \min_{\theta} \sum_{i=1}^n \left(y^{(i)} - \theta^T \mathbf{x}^{(i)}\right)^2 + \lambda \sum_{j=1}^p \theta_j^2 \\ &= \arg \min_{\theta} \|\mathbf{y} - \mathbf{X}\theta\|_2^2 + \lambda \|\theta\|_2^2 \end{aligned}$$

Can still analytically solve this:

$$\hat{\theta}_{\text{ridge}} = (\mathbf{X}^T \mathbf{X} + \lambda \mathbf{I})^{-1} \mathbf{X}^T \mathbf{y}$$

Equivalent to solving the following constrained optimization problem:

$$\begin{aligned} \min_{\theta} \quad & \sum_{i=1}^n \left(y^{(i)} - f\left(\mathbf{x}^{(i)} \mid \theta\right)\right)^2 \\ \text{s.t.} \quad & \|\theta\|_2^2 \leq t \end{aligned}$$

For special case of orthonormal design $\mathbf{X}^T \mathbf{X} = \mathbf{I}$, $\hat{\theta}_{\text{OLS}} = \mathbf{X}^T \mathbf{y}$:

$$\hat{\theta}_{\text{Ridge}} = (\mathbf{X}^T \mathbf{X} + \lambda \mathbf{I})^{-1} \mathbf{X}^T \mathbf{y} = ((1 + \lambda) \mathbf{I})^{-1} \hat{\theta}_{\text{OLS}} = \frac{\hat{\theta}_{\text{OLS}}}{1 + \lambda} \quad (\text{no sparsity})$$

Geometric Analysis

Quadratic Taylor approx of unregularized $\mathcal{R}_{\text{emp}}(\theta)$ around its minimizer $\hat{\theta}$, where \mathbf{H} is the Hessian of $\mathcal{R}_{\text{emp}}(\theta)$ at $\hat{\theta}$:

$$\tilde{\mathcal{R}}_{\text{emp}}(\theta) = \mathcal{R}_{\text{emp}}(\hat{\theta}) + \nabla_{\theta} \mathcal{R}_{\text{emp}}(\hat{\theta}) \cdot (\theta - \hat{\theta}) + \frac{1}{2} (\theta - \hat{\theta})^T \mathbf{H} (\theta - \hat{\theta})$$

Since we want a minimizer, first-order term is 0 and \mathbf{H} is positive semidefinite:

$$\tilde{\mathcal{R}}_{\text{emp}}(\theta) = \mathcal{R}_{\text{emp}}(\hat{\theta}) + \frac{1}{2} (\theta - \hat{\theta})^T \mathbf{H} (\theta - \hat{\theta})$$

$$\nabla_{\theta} \tilde{\mathcal{R}}_{\text{reg}}(\theta) = 0 \rightarrow \hat{\theta}_{\text{ridge}} = (\mathbf{H} + \lambda \mathbf{I})^{-1} \mathbf{H} \hat{\theta}$$

\mathbf{H} is a real symmetric matrix, it can be decomposed as $\mathbf{H} = \mathbf{Q} \Sigma \mathbf{Q}^T$:

$$\begin{aligned} \hat{\theta}_{\text{ridge}} &= (\mathbf{Q} \Sigma \mathbf{Q}^T + \lambda \mathbf{I})^{-1} \mathbf{Q} \Sigma \mathbf{Q}^T \hat{\theta} \\ &= [\mathbf{Q} (\Sigma + \lambda \mathbf{I}) \mathbf{Q}^T]^{-1} \mathbf{Q} \Sigma \mathbf{Q}^T \hat{\theta} \\ &= \mathbf{Q} (\Sigma + \lambda \mathbf{I})^{-1} \Sigma \mathbf{Q}^T \hat{\theta} \end{aligned}$$

Lasso Regression

Use L1 penalty in linear regression:

$$\begin{aligned} \hat{\theta}_{\text{lasso}} &= \arg \min_{\theta} \sum_{i=1}^n \left(y^{(i)} - \theta^T \mathbf{x}^{(i)}\right)^2 + \lambda \sum_{j=1}^p |\theta_j| \\ &= \arg \min_{\theta} (\mathbf{y} - \mathbf{X}\theta)^T (\mathbf{y} - \mathbf{X}\theta) + \lambda \|\theta\|_1 \end{aligned}$$

Lasso can shrink some coeffs to zero, which gives sparse solutions. However, it has difficulties handling correlated predictors.

Equivalent to solving the following constrained optimization problem:

$$\begin{aligned} \min_{\theta} \quad & \sum_{i=1}^n \left(y^{(i)} - f\left(\mathbf{x}^{(i)} \mid \theta\right)\right)^2 \\ \text{s.t.} \quad & \|\theta\|_1 \leq t \end{aligned}$$

For special case of orthonormal design $\mathbf{X}^T \mathbf{X} = \mathbf{I}$, $\hat{\theta}_{\text{OLS}} = \mathbf{X}^T \mathbf{y}$:

$$\hat{\theta}_{\text{lasso}} = \text{sign}(\hat{\theta}_{\text{OLS}})(|\hat{\theta}_{\text{OLS}}| - \lambda)_+ \quad (\text{sparsity}).$$

Function $S(\theta, \lambda) := \text{sign}(\theta)(|\theta| - \lambda)_+$ is called **soft thresholding** operator: for $|\theta| \leq \lambda$ it returns 0, whereas params $|\theta| > \lambda$ are shrunk toward 0 by λ .

Geometric Analysis

$$\tilde{\mathcal{R}}_{\text{emp}}(\theta) = \mathcal{R}_{\text{emp}}(\hat{\theta}) + \frac{1}{2} (\theta - \hat{\theta})^T \mathbf{H} (\theta - \hat{\theta})$$

We assume the \mathbf{H} is diagonal, with $H_{jj} \geq 0$

$$\tilde{\mathcal{R}}_{\text{reg}}(\theta) = \mathcal{R}_{\text{emp}}(\hat{\theta}) + \sum_j \left[\frac{1}{2} H_{jj} (\theta_j - \hat{\theta}_j)^2 \right] + \sum_j \lambda |\theta_j|$$

Minimize analytically:

$$\begin{aligned} \hat{\theta}_{\text{lasso},j} &= \text{sign}(\hat{\theta}_j) \max \left\{ |\hat{\theta}_j| - \frac{\lambda}{H_{jj}}, 0 \right\} \\ &= \begin{cases} \hat{\theta}_j + \frac{\lambda}{H_{jj}} & , \text{ if } \hat{\theta}_j < -\frac{\lambda}{H_{jj}} \\ 0 & , \text{ if } \hat{\theta}_j \in [-\frac{\lambda}{H_{jj}}, \frac{\lambda}{H_{jj}}] \\ \hat{\theta}_j - \frac{\lambda}{H_{jj}} & , \text{ if } \hat{\theta}_j > \frac{\lambda}{H_{jj}} \end{cases} \end{aligned}$$

If $H_{jj} = 0$ exactly, $\hat{\theta}_{\text{lasso},j} = 0$

More Regularization Methods

Elastic Net Regression

$$\begin{aligned} \mathcal{R}_{\text{elnet}}(\theta) &= \sum_{i=1}^n (y^{(i)} - \theta^T \mathbf{x}^{(i)})^2 + \lambda_1 \|\theta\|_1 + \lambda_2 \|\theta\|_2^2 \\ &= \sum_{i=1}^n (y^{(i)} - \theta^T \mathbf{x}^{(i)})^2 + \lambda \left((1 - \alpha) \|\theta\|_1 + \alpha \|\theta\|_2^2 \right), \end{aligned}$$

$$\alpha = \frac{\lambda_2}{\lambda_1 + \lambda_2}, \lambda = \lambda_1 + \lambda_2$$

Supervised Learning :: CHEAT SHEET

Other Examples

- **L0**: not continuous or convex, NP-hard

$$\lambda \|\boldsymbol{\theta}\|_0 = \lambda \sum_j |\theta_j|^0$$

- Smoothly Clipped Absolute Deviations (SCAD): non-convex, $\gamma > 2$ controls how fast penalty “tapers off”

$$\text{SCAD}(\theta \mid \lambda, \gamma) = \begin{cases} \lambda|\theta| & \text{if } |\theta| \leq \lambda \\ \frac{2\gamma\lambda|\theta| - \theta^2 - \lambda^2}{2(\gamma-1)} & \text{if } \lambda < |\theta| < \gamma\lambda \\ \frac{\lambda^2(\gamma+1)}{2} & \text{if } |\theta| \geq \gamma\lambda \end{cases}$$

- Minimax Concave Penalty (MCP): non-convex, $\gamma > 1$ controls how fast penalty “tapers off”

$$\text{MCP}(\theta \mid \lambda, \gamma) = \begin{cases} \lambda|\theta| - \frac{\theta^2}{2\gamma}, & \text{if } |\theta| \leq \gamma\lambda \\ \frac{1}{2}\gamma\lambda^2, & \text{if } |\theta| > \gamma\lambda \end{cases}$$

Equivalence of Regularization

RRM vs MAP

Regularized risk minimization (RRM) is the same as a maximum a posteriori (MAP) estimate in Bayes.

From Bayes theorem:

$$p(\boldsymbol{\theta} \mid \mathbf{x}, y) = \frac{p(y \mid \boldsymbol{\theta}, \mathbf{x})q(\boldsymbol{\theta})}{p(y \mid \mathbf{x})} \propto p(y \mid \boldsymbol{\theta}, \mathbf{x})q(\boldsymbol{\theta})$$

The maximum a posteriori (MAP) estimator of $\boldsymbol{\theta}$ is now the minimizer of

$$-\log p(y \mid \boldsymbol{\theta}, \mathbf{x}) - \log q(\boldsymbol{\theta}).$$

Identify the loss $L(y, f(\mathbf{x} \mid \boldsymbol{\theta}))$ with $-\log(p(y \mid \boldsymbol{\theta}, \mathbf{x}))$:

- If $q(\boldsymbol{\theta})$ is constant (i.e., we used a uniform, non-informative prior), the second term is irrelevant and we arrive at ERM.
- If not, we can identify $J(\boldsymbol{\theta}) \propto -\log(q(\boldsymbol{\theta}))$, i.e., the log-prior corresponds to the regularizer, and the additional λ , which controls the strength of our penalty, usually influences the peakedness / inverse variance / strength of our prior.

L2 vs Weight Decay

L2 regularization with GD is equivalent to weight decay.

Optimize L2-regularized risk of a model $f(\mathbf{x} \mid \boldsymbol{\theta})$ by GD:

$$\min_{\boldsymbol{\theta}} \mathcal{R}_{\text{reg}}(\boldsymbol{\theta}) = \min_{\boldsymbol{\theta}} \mathcal{R}_{\text{emp}}(\boldsymbol{\theta}) + \frac{\lambda}{2} \|\boldsymbol{\theta}\|_2^2$$

The gradient is

$$\nabla_{\boldsymbol{\theta}} \mathcal{R}_{\text{reg}}(\boldsymbol{\theta}) = \nabla_{\boldsymbol{\theta}} \mathcal{R}_{\text{emp}}(\boldsymbol{\theta}) + \lambda \boldsymbol{\theta}$$

We iteratively update $\boldsymbol{\theta}$ by step size α times the negative gradient:

$$\begin{aligned} \boldsymbol{\theta}^{[\text{new}]} &= \boldsymbol{\theta}^{[\text{old}]} - \alpha \left(\nabla_{\boldsymbol{\theta}} \mathcal{R}_{\text{emp}}(\boldsymbol{\theta}^{[\text{old}]}) + \lambda \boldsymbol{\theta}^{[\text{old}]} \right) \\ &= \boldsymbol{\theta}^{[\text{old}]}(1 - \alpha\lambda) - \alpha \nabla_{\boldsymbol{\theta}} \mathcal{R}_{\text{emp}}(\boldsymbol{\theta}^{[\text{old}]}) \end{aligned}$$

We see how $\boldsymbol{\theta}^{[\text{old}]}$ decays in magnitude – for small α and λ .

Early stopping

Early stopping is another technique to avoid overfitting, which makes training process stop when validation error stops decreasing.

1. Split training data $\mathcal{D}_{\text{train}}$ into $\mathcal{D}_{\text{subtrain}}$ and \mathcal{D}_{val} .
2. Train on $\mathcal{D}_{\text{subtrain}}$ and evaluate model using the validation set \mathcal{D}_{val} .
3. Stop training when validation error stops decreasing.
4. Use parameters of the previous step for the actual model.