

General concept of regularization

Regularization is a technique that imposes constraints on model parameters to control the solution space. By limiting where parameters can be found, it effectively shrinks the space of possible solutions. This reduction in parameter flexibility not only enhances model generalizability but also helps prevent overfitting by focusing on simpler solutions.

The term derives from Latin “regula” (rule) and “regularis” (in accordance with rules), reflecting its role in establishing systematic constraints on model behavior.

Through these controlled parameter constraints and reduced solution space, regularization helps create simpler, more robust models by reducing their sensitivity to noise in the training data.

The regularization method traces back to A.N. Tikhonov’s work in 1963, who proposed it for solving ill-posed problems where formal mathematical solutions are meaningless.

A linear system $\mathbf{y} = X\boldsymbol{\beta}$ has no solution when X is singular ($\det X = 0$), rank-deficient ($\text{rank } X < k$), or when data contains errors preventing $\mathbf{y} = X\boldsymbol{\beta}$ from being satisfied:

$$\|\mathbf{y} - X\boldsymbol{\beta}\|_2^2 \rightarrow \min_{\boldsymbol{\beta}}.$$

Adding constraints on model parameters $\boldsymbol{\beta}$ through a regularizer $R(\boldsymbol{\beta})$ term shrinks the solution space, making it possible to find (sometimes) a practically useful approximate solution:

$$\|\mathbf{y} - X\boldsymbol{\beta}\|_2^2 + R(\boldsymbol{\beta}) \xrightarrow{\text{OK}} \min_{\boldsymbol{\beta}}.$$

Probabilistic interpretation of regularization

Probabilistic framework. Consider a joint distribution of data $\mathbf{x} \in \mathbb{R}^k, y \in \mathbb{R}$ and model's parameters $\boldsymbol{\theta} \in \mathbb{R}$:

$$\mathbf{x}, y, \boldsymbol{\theta} \sim X, Y, \Theta. \quad (1)$$

1. The prior distribution of \mathbf{x} is independent of parameters $\boldsymbol{\theta}$ and can be assumed to be uniform and ignored in the model:

$$X \sim f_X(\mathbf{x}|\boldsymbol{\theta}) \Rightarrow f_X(\mathbf{x}) \sim \mathcal{U}. \quad (2)$$

2. The posterior distribution of responses y depends on parameters $\boldsymbol{\theta}$ and specific data point \mathbf{x}' , following a semi-probabilistic model formalism. The model is specified by defining the conditional distribution of responses $\mathbb{P}[y|\mathbf{x} = \mathbf{x}^*, \boldsymbol{\theta}]$ given a specific $\mathbf{x} = \mathbf{x}^*$ and model parameters $\boldsymbol{\theta}$. When the parameters are fitted, we make predictions for new data points \mathbf{x}' by maximizing the probability of a response y given \mathbf{x}' :

$$\mathbf{a}_\theta(\mathbf{x}') = \arg \max_{y \in \text{supp } Y} \underbrace{f_Y(y|\mathbf{x} = \mathbf{x}', \boldsymbol{\theta} = \hat{\boldsymbol{\theta}})}_{\text{model}}. \quad (3)$$

Support of a random variable Y is the set of all possible values y_1^*, y_2^*, \dots that Y can take with non-zero probability:

$$\text{supp } Y = \{y_1^*, y_2^*, \dots\}$$

3. The prior distribution of parameters $\boldsymbol{\theta}$ is assumed to be known and defined by the hyperparameter vector $\boldsymbol{\gamma}$:

$$f_\Theta(\boldsymbol{\theta}) \Rightarrow f_\Theta(\boldsymbol{\theta}|\boldsymbol{\gamma}) \sim \Theta(\boldsymbol{\gamma}) \quad (4)$$

Applying MAP. The joint distribution of data and parameters can be rewritten as a product of conditional pdfs:

$$\begin{aligned} f(\mathbf{x}, y, \boldsymbol{\theta}) &= f(y|\mathbf{x}, \boldsymbol{\theta}) \cdot f(\mathbf{x}, \boldsymbol{\theta}) \\ &= f(y|\mathbf{x}, \boldsymbol{\theta}) \cdot \cancel{f(\mathbf{x}|\boldsymbol{\theta})} \cdot f(\boldsymbol{\theta}|\boldsymbol{\gamma}) \\ &= f(y|\mathbf{x}, \boldsymbol{\theta}) \cdot f(\boldsymbol{\theta}|\boldsymbol{\gamma}) \end{aligned} \quad (5)$$

As it was mentioned, the canceled prior distribution of data $f(\mathbf{x}|\boldsymbol{\theta})$ is independent of the model parameters $\boldsymbol{\theta}$. We ignore it (or assume uniform).

Still, we didn't ignore the prior distribution of parameters $\boldsymbol{\theta}$, which is $f(\boldsymbol{\theta}|\boldsymbol{\gamma})$. Because of that, it's MAP (Maximum a Posteriori) estimation, not MLE (Maximum Likelihood Estimation).

We omit random variables X, Y, Θ in the pdf's underscripts for brevity. Just look at the arguments before the bar to understand to which random variable the pdf refers: e.g., $f(\mathbf{x}, y|\boldsymbol{\theta})$ means $f_{X,Y}(\mathbf{x}, y|\boldsymbol{\theta})$.

$$\begin{aligned} \mathbb{P}[y|\mathbf{x}, \boldsymbol{\theta}] &= \mathbb{P}\{Y = y\} \{X = \mathbf{x}, \Theta = \boldsymbol{\theta}\} \\ &= \mathbb{P}\{Y = y\} \{X = \mathbf{x}\} \{\Theta = \boldsymbol{\theta}\} \\ &= \frac{\mathbb{P}\{Y = y\} \{X = \mathbf{x}\} \{\Theta = \boldsymbol{\theta}\}}{\mathbb{P}\{X = \mathbf{x}\} \{\Theta = \boldsymbol{\theta}\}} \\ &= \frac{\mathbb{P}[x, y, \boldsymbol{\theta}]}{\mathbb{P}[x, \boldsymbol{\theta}]} \end{aligned}$$

Finding parameters. For specific training samples y^*, \mathbf{x}^* and predefined hyperparameters $\boldsymbol{\gamma}^*$, we write the joint likelihood of data and model parameters and maximize it:

$$\begin{aligned} \ell(\boldsymbol{\theta}) &= \log \prod_{(\mathbf{x}^*, y^*) \in (X, Y)^\ell} \underbrace{f(y = y^*|\mathbf{x} = \mathbf{x}^*, \boldsymbol{\theta}) \cdot f(\boldsymbol{\theta}|\boldsymbol{\gamma} = \boldsymbol{\gamma}^*)}_{\text{MAP}} \\ &= \sum_{(\mathbf{x}^*, y^*) \in (X, Y)^\ell} \{\log f(y = y^*|\mathbf{x} = \mathbf{x}^*, \boldsymbol{\theta}) + \log f(\boldsymbol{\theta}|\boldsymbol{\gamma} = \boldsymbol{\gamma}^*)\} \\ &= \sum_{(\mathbf{x}^*, y^*) \in (X, Y)^\ell} \underbrace{\log f(y = y^*|\mathbf{x} = \mathbf{x}^*, \boldsymbol{\theta})}_{\text{log-likelihood}} + \lambda \cdot \underbrace{\log f(\boldsymbol{\theta}|\boldsymbol{\gamma} = \boldsymbol{\gamma}^*)}_{\text{prior regularizer}} \rightarrow \max_{\boldsymbol{\theta}} \end{aligned} \quad (6)$$

The second term is the regularizer, its strength is defined by constant λ and hyperparameters $\boldsymbol{\gamma}$. Regularizer narrows the space in which the parameters can be found. The more narrow the space, the more constrained the model is.

After finding the parameter vector estimate $\hat{\boldsymbol{\theta}}$, predictions for a new data point \mathbf{x}' can be made by substituting the estimate $\hat{\boldsymbol{\theta}}$ into the model $f_Y(y|\mathbf{x} = \mathbf{x}', \boldsymbol{\theta} = \hat{\boldsymbol{\theta}})$:

$$\mathbf{a}_{\hat{\boldsymbol{\theta}}}(\mathbf{x}') = \arg \max_{y \in \text{supp } Y} f_Y(y|\mathbf{x} = \mathbf{x}', \boldsymbol{\theta} = \hat{\boldsymbol{\theta}}) \quad (7)$$

Loss-function. Probabilistic regularizer (6) can be rewritten as the empirical risk where it becomes an additional loss function:

A pdf $f(y|\mathbf{x}, \boldsymbol{\theta})$ becomes a likelihood function when we consider it as a function of arguments behind the bar, e.g.
 * $h(y) := f(y|\mathbf{x} = \mathbf{x}^*, \boldsymbol{\theta} = \boldsymbol{\theta}^*)$ is still a pdf of y given $\mathbf{x} = \mathbf{x}^*$ and $\boldsymbol{\theta} = \boldsymbol{\theta}^*$.
 * $g(\boldsymbol{\theta}) := f(y = y^*|\mathbf{x} = \mathbf{x}^*, \boldsymbol{\theta})$ is already a likelihood function of $\boldsymbol{\theta}$ given $\mathbf{x} = \mathbf{x}^*$ and $y = y^*$.

$$R(\theta) = \underbrace{\sum_{\mathbf{x} \in \mathcal{X}^\ell} \mathcal{L}(\mathbf{x})}_{\text{data fitting}} + \underbrace{\lambda \cdot \mathcal{L}_{\text{reg}}(\theta)}_{\substack{\text{parameters} \\ \text{regularization}}} \rightarrow \min_{\theta} \quad (8)$$

L_2 -norm regularization (Tikhonov regularization)

On a previous step, we found the general form of the regularizer assuming prior distribution of parameters $f_{\Theta}(\theta|\gamma)$:

$$R(\theta) = \sum_{\mathbf{x} \in X^\ell} \mathcal{L}(\mathbf{x}) + \lambda \cdot \mathcal{L}_{\text{reg}}(\theta) \rightarrow \min_{\theta} . \quad (9)$$

Model. Here we make specific choices of the prior distributions parameters $f_{\Theta}(\theta|\gamma)$ and the data:

1. All parameters $\theta \rightarrow \beta$ are independent and linear, so the joint distribution is a product of the individual distributions:

$$f_{\Theta}(\theta = \beta|\gamma) = \prod_{j=1}^k f(\beta_j|\gamma) \quad (10)$$

These distributions impose prior constraints on the model coefficients, effectively reducing the solution space

2. Each parameter β_j follows a Gaussian distribution with two hyperparameters common for all individual distributions: mean $\gamma_1 \rightarrow \mu = 0$ and standard deviation $\gamma_2 \rightarrow \tau$:

$$f(\beta_j|\gamma) \rightarrow f(\beta_j|\mu, \tau) = \frac{1}{\sqrt{2\pi\tau}} e^{-\beta_j^2/2\tau^2} \sim N(\mu, \tau). \quad (11)$$

3. The data is generated by a linear model with Gaussian noise $N(0, \sigma)$:

$$y(\mathbf{x}) = \beta^T \mathbf{x} + \varepsilon(\mathbf{x}), \quad \varepsilon(\mathbf{x}) \sim N(0, \sigma). \quad (12)$$

Both errors $\varepsilon \sim N(0, \sigma I)$ and model's parameters $\beta \sim N(0, \tau I)$ follow multivariate Gaussian distributions with zero mean and different covariance matrices σI and τI respectively.

Applying MAP. For any arbitrary model we can estimate the error term $\hat{\varepsilon}$ as difference between the predicted $\hat{\mathbf{y}}$ and the actual \mathbf{y} responses. The posterior distribution of parameters β :

$$\begin{aligned} f_{\beta}(\beta|\gamma, \varepsilon = \hat{\varepsilon}) &= \frac{f_{\beta, \gamma, \varepsilon}(\beta, \gamma, \varepsilon = \hat{\varepsilon})}{f_{\gamma, \varepsilon}(\gamma, \varepsilon = \hat{\varepsilon})} \\ &= \frac{f_{\varepsilon}(\varepsilon = \hat{\varepsilon}|\beta, \gamma) \cdot f_{\beta, \gamma}(\beta, \gamma)}{f_{\gamma, \varepsilon}(\gamma, \varepsilon = \hat{\varepsilon})} \\ &= \frac{f_{\varepsilon}(\varepsilon = \hat{\varepsilon}|\beta, \gamma) \cdot f_{\beta}(\beta|\gamma) \cdot f_{\gamma}(\gamma)}{f_{\gamma, \varepsilon}(\gamma, \varepsilon = \hat{\varepsilon})} \\ &= \frac{f_{\varepsilon}(\varepsilon = \hat{\varepsilon}|\beta) \cdot f_{\beta}(\beta|\gamma) \cdot \cancel{f_{\gamma}(\gamma)}}{\cancel{f_{\gamma}(\gamma)} \cdot f_{\varepsilon}(\varepsilon = \hat{\varepsilon})} \rightarrow \max_{\beta} \end{aligned} \quad (13)$$

Dimensions of vectors ε and β are different:

$$\dim \varepsilon = \ell, \quad \dim \beta = k.$$

All corresponding distribution parameters have appropriate dimensions:

$$\begin{aligned} \dim \mathbf{0}_{\varepsilon} &= \ell, & \dim \mathbf{0}_{\beta} &= k, \\ \dim \sigma I &= \ell \times \ell, & \dim \tau I &= k \times k. \end{aligned}$$

Here we denote pdf underscripts with letters corresponding to random variables

As $f_{\varepsilon}(\varepsilon = \hat{\varepsilon})$ is independent of β , we cancel it out:

$$\boxed{f_{\varepsilon}(\varepsilon = \hat{\varepsilon}|\beta) \cdot f_{\beta}(\beta) \rightarrow \max_{\beta}.} \quad (14)$$

Independence. We applied MAP and wrote the optimization problem, now continue with substituting the specific distributions:

$$f_{\varepsilon}(\varepsilon = \hat{\varepsilon}|\beta) = \prod_{\mathbf{x}^* \in X^\ell} e^{-\varepsilon(\mathbf{x}=\mathbf{x}^*|\beta)^2/2\sigma^2} \quad (15)$$

The error estimates are directly related to the data:

$$\hat{\varepsilon}(\mathbf{x} = \mathbf{x}^*|\beta) = \hat{\mathbf{y}}(\mathbf{x}^*|\beta) - y(\mathbf{x}^*) \quad (16)$$

Let's write the prior distribution of parameters:

$$f_{\beta}(\beta) = \prod_{j=1}^k e^{-\frac{\beta_j^2}{2}\tau^2} \quad (17)$$

The data distribution can be written through the error distribution:

$$f_Y(\mathbf{y} = \mathbf{y}^*|\beta) = f_{\varepsilon}(\varepsilon = \mathbf{y} - \mathbf{y}^*|\beta) \quad (18)$$

Let's write the posterior distribution of parameters:

$$f_{\beta}(\beta|\varepsilon = \hat{\varepsilon}) := \prod_{\mathbf{x} \in X^{\ell}} e^{-\frac{\varepsilon(\mathbf{x}|\beta)^2}{2\sigma^2}} \cdot \prod_{j=1}^k e^{-\frac{\beta_j^2}{2\tau^2}} \quad (19)$$

Let's write the likelihood function (log-loss):

$$\ell(\varepsilon, \beta|X) := - \sum_{\mathbf{x} \in X^{\ell}} \frac{\varepsilon(\mathbf{x}|\beta)^2}{2\sigma^2} - \sum_{j=1}^k \frac{\beta_j^2}{2\tau^2} \rightarrow \max_{\beta} \quad (20)$$

Let's rewrite it as empirical risk minimization:

$$Q(\beta) = \sum_{\mathbf{x} \in X^{\ell}} (\hat{y}(\mathbf{x}) - y(\mathbf{x}))^2 + \frac{2\sigma^2}{2\tau^2} \sum_{j=1}^k \beta_j^2 \rightarrow \min_{\beta} \quad (21)$$

$$Q(\beta) = \|\mathbf{y} - X\beta\|_2^2 + \lambda \cdot \|\beta\|_2^2 \rightarrow \min_{\beta} \quad (22)$$

In L_1 regularization, everything is similar, but the errors are described by the Laplace distribution.

L_1 -norm regularization

$$Q(\beta) = \|\mathbf{y} - X\beta\|_2^2 + \lambda \cdot \|\beta\|_1 \rightarrow \min_{\beta} \quad (23)$$

Unlike L_2 regularization, LASSO assumes model errors follow the Laplace distribution, characterized by heavy tails and a sharp peak:

$$\varepsilon \sim \text{pdf}(\varepsilon|\mu, b) = \frac{1}{2b} \exp \frac{-|\varepsilon - \mu|}{b} \quad (24)$$

Using MAP for parameter estimation:

$$\begin{aligned} f_{\beta}(\beta|\varepsilon = \hat{\varepsilon}) &= \frac{f_{\varepsilon}(\varepsilon = \hat{\varepsilon}|\beta) \cdot f_{\beta}(\beta)}{f_{\varepsilon}(\varepsilon = \hat{\varepsilon})} \\ &= \prod_{\mathbf{x} \in X^{\ell}} e^{-\frac{\varepsilon(\mathbf{x}|\beta)^2}{2\sigma^2}} \cdot \prod_{j=1}^k \frac{1}{2b} e^{-\frac{|\beta_j|}{b}} \\ &\rightarrow \max_{\beta} \end{aligned} \quad (25)$$

Let's write the likelihood function (log-loss):

$$\ell(\varepsilon, \beta|X) := - \sum_{\mathbf{x} \in X^{\ell}} \frac{\varepsilon(\mathbf{x}|\beta)^2}{2\sigma^2} - \sum_{j=1}^k \frac{|\beta_j|}{b} - k \ln 2b \rightarrow \max_{\beta} \quad (26)$$

Let's rewrite it as empirical risk minimization:

$$Q(\beta) = \sum_{\mathbf{x} \in X^{\ell}} (\hat{y}(\mathbf{x}) - y(\mathbf{x}))^2 + \frac{2\sigma^2}{b} \sum_{j=1}^k |\beta_j| \rightarrow \min_{\beta} \quad (27)$$

