

Title: Matrix regularization

URL: [https://en.wikipedia.org/wiki/Matrix\\_regularization](https://en.wikipedia.org/wiki/Matrix_regularization)

PageID: 44628821

Categories: Category:Estimation theory, Category:Machine learning, Category:Matrices (mathematics)

Source: Wikipedia (CC BY-SA 4.0).

-----

In the field of statistical learning theory, matrix regularization generalizes notions of vector regularization to cases where the object to be learned is a matrix. The purpose of regularization is to enforce conditions, for example sparsity or smoothness, that can produce stable predictive functions. For example, in the more common vector framework, Tikhonov regularization optimizes over  $\min_x \|Ax - y\|^2 + \lambda \|x\|^2$  to find a vector  $x$  that is a stable solution to the regression problem. When the system is described by a matrix rather than a vector, this problem can be written as  $\min_X \|AX - Y\|^2 + \lambda \|X\|^2$ , where the vector norm enforcing a regularization penalty on  $x$  has been extended to a matrix norm on  $X$ .

Matrix regularization has applications in matrix completion, multivariate regression, and multi-task learning. Ideas of feature and group selection can also be extended to matrices, and these can be generalized to the nonparametric case of multiple kernel learning.

#### Basic definition

Consider a matrix  $W$  to be learned from a set of examples,  $S = (X_i, y_i)$ , where  $i$  goes from 1 to  $n$ , and  $t$  goes from 1 to  $T$ . Let each input matrix  $X_i$  be  $\in \mathbb{R}^{D \times T}$ , and let  $W$  be of size  $D \times T$ . A general model for the output  $y$  can be posed as  $y_i = \langle W, X_i \rangle F$ , where the inner product is the Frobenius inner product. For different applications the matrices  $X_i$  will have different forms, but for each of these the optimization problem to infer  $W$  can be written as  $\min_{W \in \mathcal{H}} E(W) + R(W)$ , where  $E$  defines the empirical error for a given  $W$ , and  $R(W)$  is a matrix regularization penalty. The function  $R(W)$  is typically chosen to be convex and is often selected to enforce sparsity (using  $\ell_1$ -norms) and/or smoothness (using  $\ell_2$ -norms). Finally,  $W$  is in the space of matrices  $\mathcal{H}$  with Frobenius inner product  $\langle \dots \rangle_F$ .

#### General applications

##### Matrix completion

In the problem of matrix completion, the matrix  $X_i$  takes the form  $X_i = e_t \otimes e_i'$ , where  $(e_t)_t$  and  $(e_i')_i$  are the canonical basis in  $\mathbb{R}^T$  and  $\mathbb{R}^D$ . In this case the role of the Frobenius inner product is to select individual elements  $w_{it}$  from the matrix  $W$ . Thus, the output  $y$  is a sampling of entries from the matrix  $W$ .

The problem of reconstructing  $W$  from a small set of sampled entries is possible only under certain restrictions on the matrix, and these restrictions can be enforced by a regularization function. For example, it might be assumed that  $W$  is low-rank, in which case the regularization penalty can take the form of a nuclear norm.  $R(W) = \lambda \|W\|_*$

$= \lambda \sum_i |\sigma_i|$ ,  $\{\displaystyle R(W)=\lambda \left\|W\right\|_*=\lambda \sum_i \left\|\sigma_{ij}\right\|_1\}$ , where  $\sigma_i$   $\{\displaystyle \sigma_{ij}\}$ , with  $i$   $\{\displaystyle i\}$  from 1  $\{\displaystyle 1\}$  to  $\min D, T$   $\{\displaystyle \min D, T\}$ , are the singular values of  $W$   $\{\displaystyle W\}$ .

### Multivariate regression

Models used in multivariate regression are parameterized by a matrix of coefficients. In the Frobenius inner product above, each matrix  $X$   $\{\displaystyle X\}$  is  $X_{it} = e_t \otimes x_i$   $\{\displaystyle X_{it}=e_t \otimes x_i\}$  such that the output of the inner product is the dot product of one row of the input with one column of the coefficient matrix. The familiar form of such models is  $Y = XW + b$   $\{\displaystyle Y=XW+b\}$

Many of the vector norms used in single variable regression can be extended to the multivariate case. One example is the squared Frobenius norm, which can be viewed as an  $\ell^2$   $\{\displaystyle \ell^2\}$ -norm acting either entrywise, or on the singular values of the matrix:  $R(W) = \lambda \|W\|_F^2 = \lambda \sum_i \sum_j |w_{ij}|^2 = \lambda \text{Tr}(W^* W) = \lambda \sum_i \sigma_i^2$ .  $\{\displaystyle R(W)=\lambda \left\|W\right\|_F^2=\lambda \sum_i \sum_j \left|w_{ij}\right|^2=\lambda \operatorname{Tr}\left(W^* W\right)=\lambda \sum_i \sigma_i^2.\}$

In the multivariate case the effect of regularizing with the Frobenius norm is the same as the vector case; very complex models will have larger norms, and, thus, will be penalized more.

### Multi-task learning

The setup for multi-task learning is almost the same as the setup for multivariate regression. The primary difference is that the input variables are also indexed by task (columns of  $Y$   $\{\displaystyle Y\}$ ). The representation with the Frobenius inner product is then  $X_{it} = e_t \otimes x_i$ .  $\{\displaystyle X_{it}=e_t \otimes x_i\}$

The role of matrix regularization in this setting can be the same as in multivariate regression, but matrix norms can also be used to couple learning problems across tasks. In particular, note that for the optimization problem  $\min_W \|XW - Y\|_F^2 + \lambda \|W\|_F^2$   $\{\displaystyle \min_W \left\|XW-Y\right\|_F^2+\lambda \left\|W\right\|_F^2\}$  the solutions corresponding to each column of  $Y$   $\{\displaystyle Y\}$  are decoupled. That is, the same solution can be found by solving the joint problem, or by solving an isolated regression problem for each column. The problems can be coupled by adding an additional regularization penalty on the covariance of solutions  $\min_{W, \Omega} \|XW - Y\|_F^2 + \lambda_1 \|W\|_F^2 + \lambda_2 \text{Tr}(\Omega^{-1} W^T W)$   $\{\displaystyle \min_{W, \Omega} \left\|XW-Y\right\|_F^2+\lambda_1 \left\|W\right\|_F^2+\lambda_2 \operatorname{Tr}\left(W^T \Omega^{-1} W\right)\}$  where  $\Omega$   $\{\displaystyle \Omega\}$  models the relationship between tasks. This scheme can be used to both enforce similarity of solutions across tasks, and to learn the specific structure of task similarity by alternating between optimizations of  $W$   $\{\displaystyle W\}$  and  $\Omega$   $\{\displaystyle \Omega\}$ . [ 3 ] When the relationship between tasks is known to lie on a graph, the Laplacian matrix of the graph can be used to couple the learning problems.

### Spectral regularization

Regularization by spectral filtering has been used to find stable solutions to problems such as those discussed above by addressing ill-posed matrix inversions (see for example Filter function for Tikhonov regularization). In many cases the regularization function acts on the input (or kernel) to ensure a bounded inverse by eliminating small singular values, but it can also be useful to have spectral norms that act on the matrix that is to be learned.

There are a number of matrix norms that act on the singular values of the matrix. Frequently used examples include the Schatten p-norms, with  $p = 1$  or  $2$ . For example, matrix regularization with a Schatten 1-norm, also called the nuclear norm, can be used to enforce sparsity in the spectrum of a matrix. This has been used in the context of matrix completion when the matrix in question is believed to have a restricted rank. [ 2 ] In this case the optimization problem becomes:  $\min_W \|W\|_*$  subject to  $W_{i,j} = Y_{i,j}$ .  $\{\displaystyle \min \left\|W\right\|_* \sim \text{subject to } W_{i,j}=Y_{i,j}.\}$

Spectral Regularization is also used to enforce a reduced rank coefficient matrix in multivariate regression. [ 4 ] In this setting, a reduced rank coefficient matrix can be found by keeping just the top  $n$   $\{\displaystyle n\}$  singular values, but this can be extended to keep any reduced set of singular

values and vectors.

### Structured sparsity

Sparse optimization has become the focus of much research interest as a way to find solutions that depend on a small number of variables (see e.g. the Lasso method). In principle, entry-wise sparsity can be enforced by penalizing the entry-wise  $\ell^0$ -norm of the matrix, but the  $\ell^0$ -norm is not convex. In practice this can be implemented by convex relaxation to the  $\ell^1$ -norm. While entry-wise regularization with an  $\ell^1$ -norm will find solutions with a small number of nonzero elements, applying an  $\ell^1$ -norm to different groups of variables can enforce structure in the sparsity of solutions. [ 5 ]

The most straightforward example of structured sparsity uses the  $\ell_{p,q}$  norm with  $p = 2$  and  $q = 1$ :  $\|W\|_{2,1} = \sum_i \left( \sum_j |w_{ij}|^2 \right)^{1/2}$ .

For example, the  $\ell_{2,1}$  norm is used in multi-task learning to group features across tasks, such that all the elements in a given row of the coefficient matrix can be forced to zero as a group. [ 6 ] The grouping effect is achieved by taking the  $\ell^2$ -norm of each row, and then taking the total penalty to be the sum of these row-wise norms. This regularization results in rows that will tend to be all zeros, or dense. The same type of regularization can be used to enforce sparsity column-wise by taking the  $\ell^2$ -norms of each column.

More generally, the  $\ell_{2,1}$  norm can be applied to arbitrary groups of variables:  $R(W) = \lambda \sum_g \left( \sum_{j \in G_g} |w_{gj}|^2 \right)^{1/2} = \lambda \sum_g \left( \sum_{j \in G_g} |w_{gj}|^2 \right)^{1/2}$  where the index  $g$  is across groups of variables, and  $|G_g|$  indicates the cardinality of group  $g$ .

Algorithms for solving these group sparsity problems extend the more well-known Lasso and group Lasso methods by allowing overlapping groups, for example, and have been implemented via matching pursuit [ 7 ] and proximal gradient methods [ 8 ]. By writing the proximal gradient with respect to a given coefficient,  $w_{gi}$ , it can be seen that this norm enforces a group-wise soft threshold [ 1 ]  $\text{prox}_{\lambda} (w_{gi}) = \begin{cases} w_{gi} - \lambda & \text{if } |w_{gi}| \geq \lambda \\ 0 & \text{if } |w_{gi}| < \lambda \end{cases}$  where  $\mathbf{1}_{|w_{gi}| \geq \lambda}$  is the indicator function for group norms  $\geq \lambda$ .

Thus, using  $\ell_{2,1}$  norms it is straightforward to enforce structure in the sparsity of a matrix either row-wise, column-wise, or in arbitrary blocks. By enforcing group norms on blocks in multivariate or multi-task regression, for example, it is possible to find groups of input and output variables, such that defined subsets of output variables (columns in the matrix  $Y$ ) will depend on the same sparse set of input variables.

### Multiple kernel selection

The ideas of structured sparsity and feature selection can be extended to the nonparametric case of multiple kernel learning [ 9 ]. This can be useful when there are multiple types of input data (color and texture, for example) with different appropriate kernels for each, or when the appropriate kernel is unknown. If there are two kernels, for example, with feature maps  $A$  and  $B$  that lie in corresponding reproducing kernel Hilbert spaces  $H_A, H_B$ , then a larger space,  $H_D$ , can be created as the sum of two spaces:  $H_D : f = h + h'; h \in H_A, h' \in H_B$  assuming linear independence in  $A$  and  $B$ . In this case the  $\ell_{2,1}$ -norm is again the sum of norms:  $\|f\|_{2,1} = \|h\|_A + \|h'\|_B$ .

Thus, by choosing a matrix regularization function as this type of norm, it is possible to find a solution that is sparse in terms of which kernels are used, but dense in the coefficient of each used kernel. Multiple kernel learning can also be used as a form of nonlinear variable selection, or as a model aggregation technique (e.g. by taking the sum of squared norms and relaxing sparsity constraints). For example, each kernel can be taken to be the Gaussian kernel with a different width.

See also

Regularization (mathematics)

References