

Feature Selection and Sparsity

Makoto Yamada
makoto.yamada@riken.jp

Kyoto University/RIKEN AIP

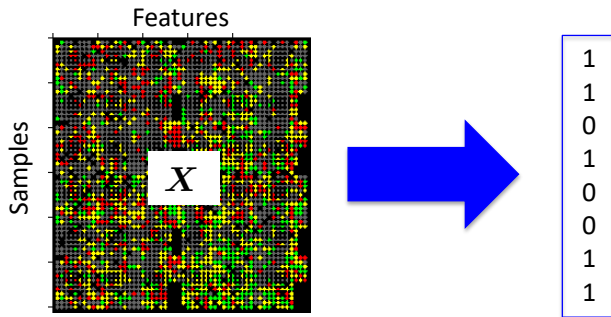
June/3/2019



Introduction

Feature selection is important for handling **high-dimensional** data:

- User data ($d > 100$)
e.g., e-mail spam detection.
- Gene expression data ($d > 20000$)
e.g., cancer classification.
- Text based feature such as TF-IDF ($d > 100,000$)
e.g., Sentiment analysis



Motivation1

The purpose of feature selection is

- to **improve the prediction** accuracy by getting rid of non-important features.
- to make the prediction **faster**.
- to **interpret** data.
- to handle **high-dimensional** data.

Let us think about a least-squared regression problems:

$$\min_{\mathbf{w} \in \mathbb{R}^d} \|\mathbf{y} - \mathbf{X}^\top \mathbf{w}\|_2^2$$

where $\mathbf{x} = (x_1, x_2, \dots, x_d)^\top \in \mathbb{R}^d$ and $\mathbf{y} \in \mathbb{R}^n$, and $\|\cdot\|_2^2$ is the ℓ_2 norm.

Question:

- $d < n$ and the rank of \mathbf{X} is d . Please derive the analytical solution of \mathbf{w} .

Motivation2

Take the objective function with respect to \mathbf{w} and set it to zero:

$$\frac{\partial}{\partial \mathbf{w}} \|\mathbf{y} - \mathbf{X}^\top \mathbf{w}\|_2^2 = -2\mathbf{X}(\mathbf{y} - \mathbf{X}^\top \mathbf{w}) = \mathbf{0}$$

Use Eq. (84) of [1]. The solution is given as

$$\hat{\mathbf{w}} = (\mathbf{X}\mathbf{X}^\top)^{-1}\mathbf{X}\mathbf{y}.$$

If the rank of \mathbf{X} is d , the rank of $\mathbf{X}\mathbf{X}^\top$ is also d and it is [invertible](#).

What happens if the rank of \mathbf{X} is less than d ?

- $\mathbf{X}\mathbf{X}^\top$ is [not invertible](#).
- Maybe, we can add a regularizer (or use pseudo-inverse). we get a [dense](#) solution and numerically unstable :(

A possible solution is to use [feature selection](#)! If we select $r < d$ features, we can compute \mathbf{w} .

Problem formulation

Problem formulation of feature selection

- Input vector: $\mathbf{x} = (x_1, x_2, \dots, x_d)^\top \in \mathbb{R}^d$
- Output: $y \in \mathbb{R}$
- Paired data: $\{(\mathbf{x}_i, y_i)\}_{i=1}^n$

Goal: Select r ($r < d$) features of input \mathbf{x} that are responsible for output y .

Problems: There is 2^d combinations :(It is hard even if d is 100.

Feature Selection Algorithms

The feature selection algorithms are categorized into three types:

- **Wrapper Method**

Use a predictive model to select features.

- **Filter Method**

Use a proxy measure (such as mutual information) instead of the error rate to select features.

- **Embedded Method**

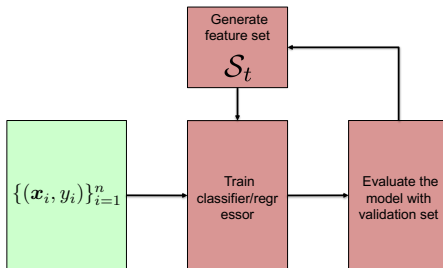
Features are selected as part of the model construction process.

Wrapper Method

Use a predictive model (e.g., classifier) to select features.

The simplest approach would be...

- 1 Generate feature set \mathcal{S}_t
- 2 Train predictive model with \mathcal{S}_t and test the prediction accuracy with hold-out set.
- 3 Iterate 1 and 2 until all feature combination is examined.



Wrapper Method

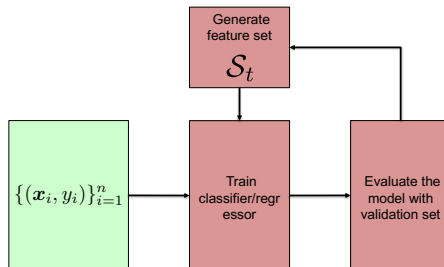
Pro:

- It can select features that have feature-feature interaction.

Cons:

- It can be overfitted if the number of samples is insufficient.
- Computationally expensive.

Wrapper method is not that popular compared to filter and embedded methods...



Filter Method

Use a proxy measure (such as mutual information) instead of the error rate to select features.

Pros:

- Easy to implement.
- It scales well (easy to implement with distributed computing).
- Can select features from high-dimensional data (both linear and nonlinear way).

Cons:

- The feature selection is **independent** of the model. The selected features may not be the best set to achieve highest accuracy.
- It is hard to detect select features with interaction. (Of course, we can somehow select them, but it increase computation cost.

Filter Method (Example)

Maximum Relevance Feature Selection (MR)

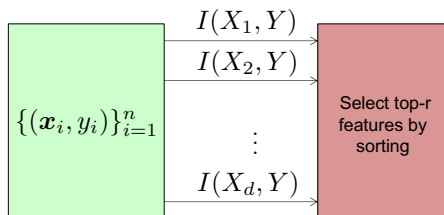
Compute association score between each feature and its output and rank them.

- Correlation, Mutual information, and the kernel based independence measures are used.
- Easy to implement and it scales well.

Optimization problem:

$$\max_{\beta \in \{0,1\}^d} \frac{1}{S} \sum_{k=1}^d \beta_k I(X_k, Y),$$

where $S = \beta_1 + \dots + \beta_d$.



Filter Method (Example)

Minimum Redundancy Maximum Relevance (mRMR) [2]

MR feature selection tends to select **redundant** features.

mRMR method is to

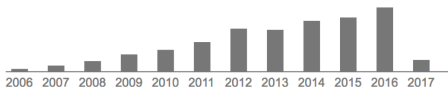
- select features that have high association to its output.
- select **independent** features.

Optimization problem:

$$\max_{\beta \in \{0,1\}^d} \frac{1}{S} \sum_{k=1}^d \beta_k I(X_k, Y) - \frac{1}{S^2} \sum_{k=1}^d \sum_{k'=1}^d \beta_k \beta_{k'} I(X_k, X_{k'}).$$

This optimization problem can be solved by using greedy algorithm.

[引用元 4361](#)



Filter Method (Mutual Information)

To optimize mRMR, we tend to use the **mutual information** as an association score.

Independence:

$$p(\mathbf{x}, \mathbf{y}) = p(\mathbf{x})p(\mathbf{y})$$

Mutual Information:

$$MI(X, Y) = \iint p(\mathbf{x}, \mathbf{y}) \log \frac{p(\mathbf{x}, \mathbf{y})}{p(\mathbf{x})p(\mathbf{y})} d\mathbf{x}d\mathbf{y}$$

Under independence:

$$MI(X, Y) = \iint p(\mathbf{x}, \mathbf{y}) \log \frac{p(\mathbf{x})p(\mathbf{y})}{p(\mathbf{x})p(\mathbf{y})} d\mathbf{x}d\mathbf{y} = 0$$

Filter Method (Example)

How to optimize mRMR? Suppose we use the mutual information as a dependency measure and we already have \mathcal{S}_{m-1} , which is the feature set with $m - 1$ features, then we can select the m -th feature by solving the following optimization problem:

$$\hat{j} = \operatorname{argmax}_{j \in \bar{\mathcal{S}}_{m-1}} \operatorname{MI}(X_j, Y) - \frac{1}{m-1} \sum_{i \in \mathcal{S}_{m-1}} \operatorname{MI}(X_j, X_i). \quad (1)$$

We select a feature that having high dependency with Y and independent of features in \mathcal{S}_{m-1} .

Filter Method (Continuous optimization)

The MR and mRMR feature selection algorithms are **discrete** optimization problem. In feature selection, **continuous** optimization based approach is also popular.

The key idea is to **relax** the condition (i.e., allow to take continuous number).

Quadratic Programming Feature Selection (QPFS) [3]:

$$\begin{aligned} \max_{\alpha \in \mathbb{R}^d} \quad & \sum_{k=1}^d \alpha_k I(X_k, Y) - \frac{1}{2} \sum_{k=1}^d \sum_{k'=1}^d \alpha_k \alpha_{k'} I(X_k, X_{k'}), \\ \text{s.t.} \quad & \alpha_1 + \alpha_2 + \dots + \alpha_d = 1, \alpha_1, \dots, \alpha_d \geq 0 \end{aligned}$$

Filter Method (Continuous optimization)

Let us denote:

$$\begin{aligned}\mathbf{h}_k &= I(X_k, Y), \\ \mathbf{H}_{kk'} &= I(X_k, X_{k'})\end{aligned}$$

where $\mathbf{h} \in \mathbb{R}^d$ and $\mathbf{H} \in \mathbb{R}^{d \times d}$.

We have

$$\begin{aligned}\min_{\alpha \in \mathbb{R}^d} \quad & \frac{1}{2} \alpha^\top \mathbf{H} \alpha - \mathbf{h}^\top \alpha \\ \text{s.t.} \quad & \alpha^\top \mathbf{1} = 1, \alpha_1, \dots, \alpha_d \geq 0.\end{aligned}$$

This is a **quadratic programming** with simplex constraint (can be solved by using an off-the-shelf package).

Note: For mutual information, \mathbf{H} may not be positive definite. **It can be non-convex optimization.**

Embedded Method

Features are selected as part of the model construction process.
Embedded method can be regarded as an intermediate method between wrapper and filter methods.

Pros:

- Can select features with high prediction accuracy.
- Computationally efficient than wrapper method.

Cons:

- Computationally expensive than filter method.
- If the input output relationship are nonlinear, it is computationally expensive. It is more suited for **linear** method.

Embedded Method (Lasso)

Least Absolute Shrinkage and Selection Operator (Lasso)

The optimization problem of Lasso can be written as

$$\min_{\mathbf{w}} \frac{1}{2} \|\mathbf{y} - \mathbf{X}^\top \mathbf{w}\|_2^2 + \lambda \|\mathbf{w}\|_1,$$

where $\mathbf{X} = [\mathbf{x}_1, \dots, \mathbf{x}_n] \in \mathbb{R}^{d \times n}$ is the input matrix and $\mathbf{y} = [y_1, \dots, y_n]^\top \in \mathbb{R}^n$ is the output vector.

$$\|\mathbf{w}\|_1 = \sum_{k=1}^d |w_k|$$

is an ℓ_1 norm.

Lasso is a convex method: The first term is a convex function w.r.t. \mathbf{w} . ℓ_1 norm (all norm) is convex:

$$\begin{aligned} \|\alpha \mathbf{w} + (1 - \alpha) \mathbf{v}\|_1 &\leq \|\alpha \mathbf{w}\|_1 + \|(1 - \alpha) \mathbf{v}\|_1 && (\text{triangle inequality}) \\ &= \alpha \|\mathbf{w}\|_1 + (1 - \alpha) \|\mathbf{v}\|_1 && (\text{absolutely scalable}), \end{aligned}$$

where $0 \leq \alpha \leq 1$. The sum of two convex functions is convex.

Embedded Method (Lasso) Some intuitive explanation

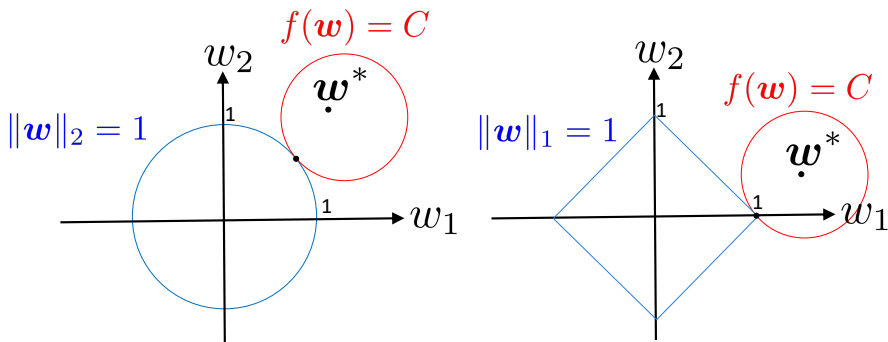
Using the ℓ_1 regularizer, we can make \mathbf{w} sparse.

The ℓ_1 regularization is equivalent to ℓ_1 norm constraint:

$$\min_{\mathbf{w}} f(\mathbf{w}) + \lambda \|\mathbf{w}\|_1 \longrightarrow \min_{\mathbf{w}} f(\mathbf{w}), \text{ s.t. } \|\mathbf{w}\|_1 \leq \eta.$$

If we consider the Lagrange function of the ℓ_1 norm constraint, there exists the same solution of the ℓ_1 norm constraint with an arbitrary λ .

Level curves of norms and loss:



Embedded Method (Lasso) When Lasso helpful?

Let us think about a least-squared regression problems:

$$\min_{\mathbf{w} \in \mathbb{R}^d} \|\mathbf{y} - \mathbf{X}^\top \mathbf{w}\|_2^2$$

where $\mathbf{x} = (x_1, x_2, \dots, x_d)^\top \in \mathbb{R}^d$ and $\mathbf{y} \in \mathbb{R}^n$, and $\|\cdot\|_2^2$ is the ℓ_2 norm. Take the objective function with respect to \mathbf{w} and set it to zero:

$$\frac{\partial}{\partial \mathbf{w}} \|\mathbf{y} - \mathbf{X}^\top \mathbf{w}\|_2^2 = -2\mathbf{X}(\mathbf{y} - \mathbf{X}^\top \mathbf{w}) = \mathbf{0}$$

Use Eq. (84) of [1]. The solution is given as

$$\hat{\mathbf{w}} = (\mathbf{X}\mathbf{X}^\top)^{-1}\mathbf{X}\mathbf{y}.$$

If the rank of \mathbf{X} is d , the rank of $\mathbf{X}\mathbf{X}^\top$ is also d and it is [invertible](#).

What happens if the rank of \mathbf{X} is less than d ?

- $\mathbf{X}\mathbf{X}^\top$ is **not invertible** \rightarrow Adding ℓ_1 regularizer to make \mathbf{w} sparse
- The number of nonzero elements of \mathbf{w} should be smaller than n .

Embedded Method (Lasso)

Lasso has no closed form solution. Thus, we need to iteratively optimize the problem.

Here, we introduce the [Alternating Direction Method of Multipliers \(ADMM\)](#) [4].

We can rewrite the Lasso optimization problem as

$$\begin{aligned} \min_{\mathbf{w}, \mathbf{z}} \quad & \frac{1}{2} \|\mathbf{y} - \mathbf{X}^T \mathbf{w}\|_2^2 + \lambda \|\mathbf{z}\|_1 + \frac{\rho}{2} \|\mathbf{w} - \mathbf{z}\|_2^2 \\ \text{s.t.} \quad & \mathbf{w} = \mathbf{z} \end{aligned}$$

The key idea here is to split the main objective and the non-differentiable regularization term. Since the last term $\frac{\rho}{2} \|\mathbf{w} - \mathbf{z}\|_2^2$ is zero if the constraint is satisfied, this problem is equivalent to the original Lasso problem.

Embedded Method (Lasso)

Let us denote the Lagrange multipliers as $\gamma \in \mathbb{R}^d$, we can write a Lagrangian function (called Augmented Lagrangian function) as follows:

$$J(\mathbf{w}, \mathbf{z}, \gamma) = \frac{1}{2} \|\mathbf{y} - \mathbf{X}^\top \mathbf{w}\|_2^2 + \gamma^\top (\mathbf{w} - \mathbf{z}) + \lambda \|\mathbf{z}\|_1 + \frac{\rho}{2} \|\mathbf{w} - \mathbf{z}\|_2^2,$$

where $\rho > 0$ is a tuning parameter.

In ADMM, we consider the following optimization problem:

$$\max_{\gamma} \min_{\mathbf{w}, \mathbf{z}} J(\mathbf{w}, \mathbf{z}, \gamma) = \frac{1}{2} \|\mathbf{y} - \mathbf{X}^\top \mathbf{w}\|_2^2 + \gamma^\top (\mathbf{w} - \mathbf{z}) + \lambda \|\mathbf{z}\|_1 + \frac{\rho}{2} \|\mathbf{w} - \mathbf{z}\|_2^2,$$

Since we have the relationship,

$$\max_{\gamma} J(\mathbf{w}, \mathbf{z}, \gamma) = \begin{cases} \frac{1}{2} \|\mathbf{y} - \mathbf{X}^\top \mathbf{w}\|_2^2 + \lambda \|\mathbf{z}\|_1 & (\mathbf{w} = \mathbf{z}) \\ \infty & (\text{Otherwise}) \end{cases}$$

The optimization problem is equivalent to the original Lasso problem.

Embedded Method (Lasso)

Minimizing $J(\mathbf{w}, \mathbf{z}, \gamma)$ w.r.t. \mathbf{w} . If we fix \mathbf{z} and γ as $\mathbf{z}^{(t)}$ and $\gamma^{(t)}$, $J(\mathbf{w}, \mathbf{z}^{(t)}, \gamma^{(t)})$ is convex w.r.t. \mathbf{w} . That is,

$$\frac{\partial J(\mathbf{w}, \mathbf{z}, \gamma)}{\partial \mathbf{w}} = -\mathbf{X}(\mathbf{y} - \mathbf{X}^\top \mathbf{w}) + \gamma + \rho(\mathbf{w} - \mathbf{z}) = \mathbf{0}.$$

Here, we can use the following equation (see [1] Eq. (84)):

$$\frac{\partial \|\mathbf{y} - \mathbf{X}^\top \mathbf{w}\|_2^2}{\partial \mathbf{w}} = -2\mathbf{X}(\mathbf{y} - \mathbf{X}^\top \mathbf{w}).$$

Solving it for \mathbf{w} :

$$\begin{aligned} (\mathbf{X}\mathbf{X}^\top + \rho\mathbf{I})\mathbf{w} &= \mathbf{X}\mathbf{y} - \gamma^{(t)} + \rho\mathbf{z}^{(t)} \\ \mathbf{w}^{(t+1)} &= (\mathbf{X}\mathbf{X}^\top + \rho\mathbf{I})^{-1}(\mathbf{X}\mathbf{y} - \gamma^{(t)} + \rho\mathbf{z}^{(t)}). \end{aligned}$$

Embedded Method (Lasso)

Minimizing $J(\mathbf{w}, \mathbf{z}, \gamma)$ w.r.t. \mathbf{z} . If we fix \mathbf{w} and γ as $\mathbf{w}^{(t)}$ and $\gamma^{(t)}$, $J(\mathbf{w}^{(t)}, \mathbf{z}, \gamma^{(t)})$ is convex w.r.t. \mathbf{z} .

$$J(\mathbf{w}^{(t)}, \mathbf{z}, \gamma^{(t)}) = \frac{\rho}{2} \|\mathbf{z} - \mathbf{w}^{(t)}\|_2^2 + \lambda \|\mathbf{z}\|_1 - \gamma^\top \mathbf{z} + \text{Const.}$$

$\|\mathbf{z}\|_1$ is not differentiable at 0. However, we can analytically solve the problem! Moreover, since there is no interaction in the elements of \mathbf{z} , we can solve it for each element.

$$J(\mathbf{w}^{(t)}, [z_1, \dots, z_\ell, \dots, z_d], \gamma^{(t)}) = \frac{\rho}{2} (z_\ell - w_\ell^{(t)})^2 + \lambda |z_\ell| - \gamma_\ell z_\ell + \text{Const.}$$

Embedded Method (Lasso)

$$J(\mathbf{w}^{(t)}, [z_1, \dots, z_\ell, \dots, z_d], \gamma^{(t)}) = \frac{\rho}{2}(z_\ell - w_\ell^{(t)})^2 + \lambda|z_\ell| - \gamma_\ell z_\ell + \text{Const.}$$

Case1: $z_\ell > 0, \rho(z_\ell - w_\ell^{(t)}) + \lambda - \gamma_\ell = 0 \longrightarrow z_\ell = w_\ell^{(t)} + \frac{1}{\rho}(\gamma_\ell - \lambda)$

That is, $z_\ell > 0$ if $w_\ell^{(t)} + \frac{1}{\rho}\gamma_\ell > \frac{\lambda}{\rho}$

Case2: $z_\ell < 0, \rho(z_\ell - w_\ell^{(t)}) - \lambda - \gamma_\ell = 0 \longrightarrow z_\ell = w_\ell^{(t)} + \frac{1}{\rho}(\gamma_\ell + \lambda)$

That is, $z_\ell < 0$ if $w_\ell^{(t)} + \frac{1}{\rho}\gamma_\ell < -\frac{\lambda}{\rho}$

Case3: $z_\ell = 0,$

$0 \in \rho(z_\ell - w_\ell^{(t)}) + \lambda[-1 \ 1] - \gamma_\ell \longrightarrow w_\ell + \frac{1}{\rho}\gamma_\ell \in [-\frac{\lambda}{\rho}, \frac{\lambda}{\rho}], (z_\ell = 0).$

Therefore, we have

$$z_\ell = \begin{cases} w_\ell^{(t)} + \frac{1}{\rho}\gamma_\ell - \frac{\lambda}{\rho} & (w_\ell^{(t)} + \frac{1}{\rho}\gamma_\ell > \frac{\lambda}{\rho}) \\ 0 & (w_\ell + \frac{1}{\rho}\gamma_\ell \in [-\frac{\lambda}{\rho}, \frac{\lambda}{\rho}]) \\ w_\ell^{(t)} + \frac{1}{\rho}\gamma_\ell + \frac{\lambda}{\rho} & (w_\ell^{(t)} + \frac{1}{\rho}\gamma_\ell < -\frac{\lambda}{\rho}) \end{cases}$$

Embedded Method (Lasso)

Let us introduce the **Soft-Thresholding function**:

$$S_{\lambda}(x) = \begin{cases} x - \lambda & (x > \lambda) \\ 0 & (x \in [-\lambda, \lambda]) \\ x + \lambda & (x < -\lambda) \end{cases},$$
$$= \text{sign}(x) \max(0, |x| - \lambda)$$

Therefore, the update of z_{ℓ} can be simply written by the soft-thresholding function as

$$\hat{z}_{\ell}^{(t+1)} = S_{\frac{\lambda}{\rho}}(w_{\ell}^{(t)} + \frac{1}{\rho}\gamma_{\ell}).$$

Embedded Method (Lasso)

Maximizing $J(\mathbf{w}, \mathbf{z}, \gamma)$ w.r.t. γ . That is the optimization problem can be written as

$$\max_{\gamma} J(\mathbf{w}, \mathbf{z}, \gamma) = \gamma^{\top}(\mathbf{w} - \mathbf{z}).$$

To optimize this problem, since we cannot get the analytical solution, we use the **gradient ascent** algorithm:

$$\gamma^{(t+1)} = \gamma^{(t)} + \rho(\mathbf{w}^{(t)} - \mathbf{z}^{(t)}).$$

Thus, the ADMM algorithm for Lasso can be summarized as

$$\mathbf{w}^{(t+1)} = (\mathbf{X}\mathbf{X}^{\top} + \rho\mathbf{I})^{-1}(\mathbf{X}\mathbf{y} - \gamma^{(t)} + \rho\mathbf{z}^{(t)})$$

$$\mathbf{z}_{\ell}^{(t+1)} = S_{\frac{\lambda}{\rho}}(\mathbf{w}^{(t+1)} + \frac{1}{\rho}\gamma)$$

$$\gamma^{(t+1)} = \gamma^{(t+1)} + \rho(\mathbf{w}^{(t+1)} - \mathbf{z}^{(t+1)}).$$

Embedded Method (Elastic-Net)

For Lasso, the number of non-zero features should be smaller than n .

How to select $r > n$ variables?

Ans: Use the elastic net regularization [5]:

$$\min_{\mathbf{w}} \|\mathbf{y} - \mathbf{X}^\top \mathbf{w}\|_2^2 + \lambda(\alpha \|\mathbf{w}\|_1 + (1 - \alpha) \|\mathbf{w}\|_2^2),$$

where $0 \leq \alpha \leq 1$ and $\lambda > 0$ is a regularization parameter.

$\|\mathbf{w}\|_2^2$ is differentiable; we can similarly solve it with ADMM.

$$\mathbf{w}^{(t+1)} = (\mathbf{X}\mathbf{X}^\top + 2\lambda(1 - \alpha)\mathbf{I} + \rho\mathbf{I})^{-1}(\mathbf{X}\mathbf{y} - \gamma^{(t)} + \rho\mathbf{z}^{(t)})$$

$$\mathbf{z}_\ell^{(t+1)} = S_{\frac{\lambda\alpha}{\rho}}(\mathbf{w}^{(t+1)} + \frac{1}{\rho}\gamma)$$

$$\gamma^{(t+1)} = \gamma^{(t)} + \rho(\mathbf{w}^{(t+1)} - \mathbf{z}^{(t+1)}).$$

Thanks to the ℓ_2 regularization, \mathbf{w} tends to be dense.

- Feature selection: Wrapper method, Filter method, and Embedded method
- Wrapper method (Selecting features that maximize prediction accuracy. **Computationally expensive.**)
- Filter method (Use mutual information to select features, e.g., MR, mRMR, QPFS, etc.)
- Embedded method (Selecting features during training. e.g., Lasso)
- Alternating Direction Method of Multipliers (ADMM).



Kaare Brandt Petersen, Michael Syskind Pedersen, et al.
The matrix cookbook.
Technical University of Denmark, 7:15, 2008.



H. Peng, F. Long, and C. Ding.
Feature selection based on mutual information: Criteria of max-dependency, max-relevance, and min-redundancy.
IEEE Transactions on Pattern Analysis and Machine Intelligence, 27:1226–1237, 2005.



I. Rodriguez-Lujan, R. Huerta, C. Elkan, and C. S. Cruz.
Quadratic programming feature selection.
JMLR, 11:1491–1516, 2010.



Stephen Boyd, Neal Parikh, Eric Chu, Borja Peleato, Jonathan Eckstein, et al.
Distributed optimization and statistical learning via the alternating direction method of multipliers.
Foundations and Trends® in Machine learning, 3(1):1–122, 2011.



Hui Zou and Trevor Hastie.

Regularization and variable selection via the elastic net.

Journal of the Royal Statistical Society: Series B (Statistical Methodology), 67(2):301–320, 2005.