

# Feature Selection and Sparsity

Makoto Yamada  
myamada@i.kyoto-u.ac.jp  
Kyoto University

June/15/2020

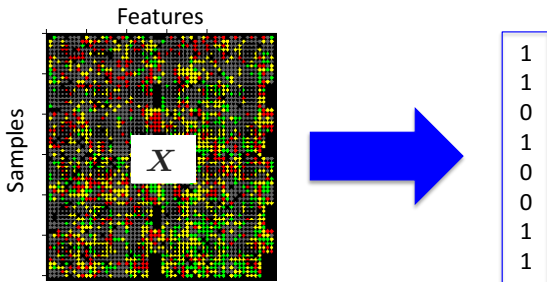
1 Introduction

2 Feature Selection Algorithms

# Introduction

Feature selection is important for **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 ( $d > 100,000$ )



# 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.

# Motivation2

Let us think about the least-squared regression problem:

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

where  $\mathbf{x} = (x_1, x_2, \dots, x_d)^\top \in \mathbb{R}^d$ ,

$\mathbf{X} = (\mathbf{x}_1, \mathbf{x}_2, \dots, \mathbf{x}_n) \in \mathbb{R}^{d \times n}$ ,  $\mathbf{w} = (w_1, w_2, \dots, w_d)^\top \in \mathbb{R}^d$ ,  
 $\mathbf{y} \in \mathbb{R}^n$ , and  $\|\cdot\|_2^2$  is the  $\ell_2$  norm.

Question:

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

# Motivation2

Take the derivative 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  $d$ ,  $\mathbf{X}\mathbf{X}^\top$  is **invertible**.

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

- $\mathbf{X}\mathbf{X}^\top$  is **not invertible**.

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.

① Introduction

② Feature Selection Algorithms



# Feature Selection Algorithms

The feature selection algorithms can be 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**

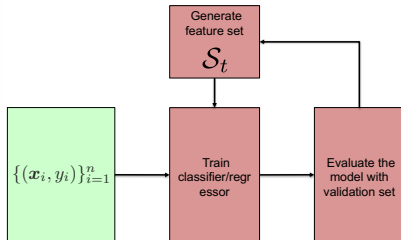
Features are selected as part of the model construction process.

# 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.



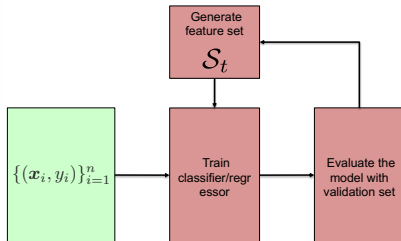
# Wrapper Method

Pro:

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

Cons:

- Computationally expensive ( $2^d$  combination).



# Filter Method

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

## Pros:

- It scales well.
- 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.

# Filter Method (Example)

## Maximum Relevance Feature Selection (MR)

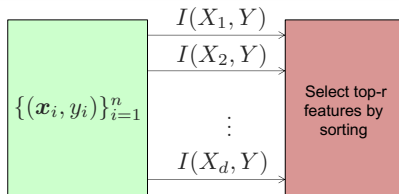
Compute association score between each feature and its output and rank them.

- 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.

# 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 (Linear Correlation)

To optimize mRMR, we may be able to use the Pearson's correlation coefficient

Pearson's correlation coefficient:

$$\text{PCC}(X, Y) = \frac{\text{Cov}(X, Y)}{\sigma_X \sigma_Y},$$

$$\text{Cov}(X, Y) = \mathbb{E}[(X - \mu_X)(Y - \mu_Y)]$$

where  $\mu_X = \mathbb{E}[X]$ ,  $\mu_Y = \mathbb{E}[Y]$ ,  $\sigma_X^2 = \mathbb{E}[(X - \mu_X)^2]$ , and  $\sigma_Y^2 = \mathbb{E}[(Y - \mu_Y)^2]$ .

The cross-covariance can be written as

$$\text{Cov}(X, Y) = \mathbb{E}[(X - \mu_X)(Y - \mu_Y)] = \mathbb{E}[XY] - \mathbb{E}[X]\mathbb{E}[Y].$$

That is, if  $\text{PCC}(X, Y) = 0$ ,  $\mathbb{E}[XY] = \mathbb{E}[X]\mathbb{E}[Y]$



# The relationship between independence and correlation

If  $X$  and  $Y$  are independent, we can write

$$\begin{aligned}\mathbb{E}[XY] &= \iint xy p(x, y) dx dy, \\ &= \iint xy p(x)p(y) dx dy, \text{ (independence)} \\ &= \left( \int x p(x) dx \right) \left( \int y p(y) dy \right) \\ &= \mathbb{E}[X]\mathbb{E}[Y]\end{aligned}$$

That is, if  $X$  and  $Y$  are independent,  $\mathbb{E}[XY] = \mathbb{E}[X]\mathbb{E}[Y]$ .

Note that, even if  $\mathbb{E}[XY] = \mathbb{E}[X]\mathbb{E}[Y]$ ,  $X$  and  $Y$  can be dependent.

# Empirical estimation of Cross-covariance

To optimize mRMR, we may be able to use the Pearson's correlation coefficient

Cross-Covariance (population):

$$\text{Cov}(X, Y) = \mathbb{E}[(X - \mu_X)(Y - \mu_Y)]$$

Cross-Covariance estimation:

$$\widehat{\text{Cov}}(X, Y) = \frac{1}{n} \sum_{i=1}^n (x_i - \hat{\mu}_X)(y_i - \hat{\mu}_Y)$$
$$\hat{\mu}_X = \frac{1}{n} \sum_{i=1}^n x_i = \frac{1}{n} \mathbf{x}^\top \mathbf{1}_n, \quad \hat{\mu}_Y = \frac{1}{n} \sum_{i=1}^n y_i = \frac{1}{n} \mathbf{y}^\top \mathbf{1}_n,$$

where  $\mathbf{1}_n = (1, 1, \dots, 1)^\top \in \mathbb{R}^n$  is the vector with all ones.

# Empirical estimation of cross-covariance

Cross-Covariance estimation:

$$\begin{aligned}\widehat{\text{Cov}}(X, Y) &= \frac{1}{n} \sum_{i=1}^n (x_i - \frac{1}{n} \mathbf{x}^\top \mathbf{1}_n) (y_i - \frac{1}{n} \mathbf{y}^\top \mathbf{1}_n) \\ &= \frac{1}{n} \left( \sum_{i=1}^n x_i y_i - \frac{1}{n} \mathbf{x}^\top \mathbf{1}_n \mathbf{1}_n^\top \mathbf{y} \right) \\ &= \frac{1}{n} \left( \mathbf{x}^\top \mathbf{y} - \frac{1}{n} \mathbf{x}^\top \mathbf{1}_n \mathbf{1}_n^\top \mathbf{y} \right) \\ &= \frac{1}{n} \mathbf{x}^\top \left( \mathbf{I}_n - \frac{1}{n} \mathbf{1}_n \mathbf{1}_n^\top \right) \mathbf{y} \\ &= \frac{1}{n} \mathbf{x}^\top \mathbf{H} \mathbf{y},\end{aligned}$$

where  $\mathbf{H} = \mathbf{I}_n - \frac{1}{n} \mathbf{1}_n \mathbf{1}_n^\top$  is the centering matrix and  $\mathbf{I}_n$  is the identity matrix. (Note  $\mathbf{H}\mathbf{H} = \mathbf{H}$ ).

# Empirical estimation of covariance

Covariance estimation:

$$\begin{aligned}\widehat{\text{Cov}}(X, Y)^2 &= \frac{1}{n^2} \mathbf{x}^\top \mathbf{H} \mathbf{y} \mathbf{x}^\top \mathbf{H} \mathbf{y}, \\ &= \frac{1}{n^2} \text{tr} \left( \mathbf{x}^\top \mathbf{H} \mathbf{y} \mathbf{y}^\top \mathbf{H} \mathbf{x} \right) \\ &= \frac{1}{n^2} \text{tr} \left( \mathbf{x} \mathbf{x}^\top \mathbf{H} \mathbf{y} \mathbf{y}^\top \mathbf{H} \right) \\ &= \frac{1}{n^2} \text{tr} \left( \mathbf{K} \mathbf{H} \mathbf{L} \mathbf{H} \right),\end{aligned}$$

where  $\mathbf{K} = \mathbf{x} \mathbf{x}^\top \in \mathbb{R}^{n \times n}$  and  $\mathbf{L} = \mathbf{y} \mathbf{y}^\top \in \mathbb{R}^{n \times n}$ .

# Advanced Topic (Hilbert-Schmidt Independence Criterion)

Hilbert Schmidt Independence Criterion (HSIC) [3]

Empirical V-statistics of HSIC is given as

$$\text{HSIC}(X, Y) = \frac{1}{n^2} \text{tr}(\mathbf{K} \mathbf{H} \mathbf{L} \mathbf{H}),$$

where we use the Gaussian kernel:

$$\mathbf{K}_{ij} = \exp\left(-\frac{\|\mathbf{x}_i - \mathbf{x}_j\|_2^2}{2\sigma^2}\right), \quad \mathbf{L}_{ij} = \exp\left(-\frac{\|\mathbf{y}_i - \mathbf{y}_j\|_2^2}{2\sigma^2}\right).$$

HSIC takes 0 if and only if  $X$  and  $Y$  are independent.

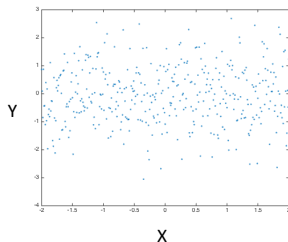
Since we can decompose  $\mathbf{K} = \Phi^\top \Phi$  and  $\mathbf{L} = \Psi^\top \Psi$ , we have

$$\text{HSIC}(X, Y) = \frac{1}{n^2} \text{tr}(\Phi^\top \Phi \mathbf{H} \Psi^\top \Psi \mathbf{H}) = \frac{1}{n^2} \|\text{vec}(\Psi \mathbf{H} \Phi^\top)\|_2^2 \geq 0$$

# Advanced Topic (HSIC)

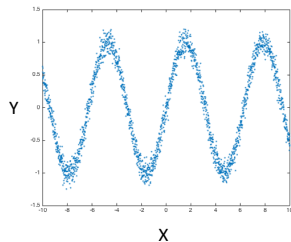
## Hilbert-Schmidt Independence Criterion (HSIC) experiments

X and Y are independent



NHSIC = 0.0031  
Pearson CC = 0.0343

X and Y are dependent



NHSIC = 0.2842  
Pearson CC = 0.1983

# 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

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

where  $\|\mathbf{w}\|_1 = \sum_{k=1}^d |w_k|$  is an  $\ell_1$  norm.

**Lasso is a convex method:** The first term is a convex function w.r.t.  $\mathbf{w}$ .  $\ell_1$  norm (all norm) is convex:

$$\begin{aligned} \|\alpha \mathbf{w} + (1 - \alpha) \mathbf{v}\|_1 &\leq \|\alpha \mathbf{w}\|_1 + \|(1 - \alpha) \mathbf{v}\|_1 \\ &= \alpha \|\mathbf{w}\|_1 + (1 - \alpha) \|\mathbf{v}\|_1 \end{aligned}$$

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



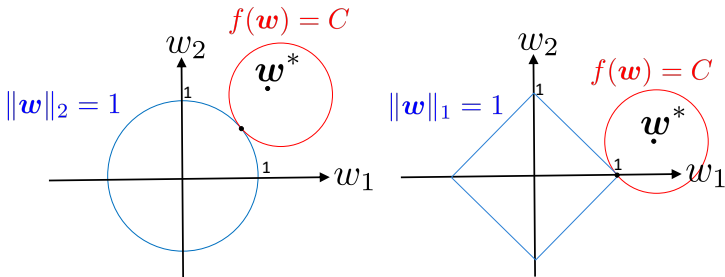
# Embedded Method (Lasso)

The  $\ell_1$  regularization is equivalent to  $\ell_1$  norm constraint:

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

There exists the same solution of the  $\ell_1$  norm constraint with an arbitrary  $\lambda$ .

Using the  $\ell_1$  regularizer, we can make  $\mathbf{w}$  sparse.



# When Lasso helpful?

Let us think about a least-squared regression problems:

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

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

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

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  $d$ ?

# Lasso with ADMM (1/8)

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

Here, we introduce the [Alternating Direction Method of Multipliers \(ADMM\)](#) [5].

We can rewrite the Lasso optimization problem as

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

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.

## Lasso with ADMM (2/8)

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.

# Lasso with ADMM (3/8)

In ADMM, we consider the following optimization problem:

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

Since we have the relationship,

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

The optimization problem is equivalent to the original Lasso problem.

# Lasso with ADMM (4/8)

Minimizing  $J(\mathbf{w}, \mathbf{z}, \gamma)$  w.r.t.  $\mathbf{w}$ . If we fix  $\mathbf{z}$  and  $\gamma$  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  $\mathbf{w}$ :

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

# Lasso with ADMM (5/8)

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

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

# Lasso with ADMM (6/8)

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



# Lasso with ADMM (7/8)

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  $z_\ell$  can be simply written by the soft-thresholding function as

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

# Lasso with ADMM (8/8)

Maximizing  $J(\mathbf{w}, \mathbf{z}, \gamma)$  w.r.t.  $\gamma$ . 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

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

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

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

# 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 [6]:

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

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.

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

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

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

Thanks to the  $\ell_2$  regularization,  $\mathbf{w}$  tends to be dense.

# Summary

- 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, etc.)
- Embedded method (Selecting features during training. e.g., Lasso)
- Alternating Direction Method of Multipliers (ADMM).

- [1] Kaare Brandt Petersen, Michael Syskind Pedersen, et al. The matrix cookbook. *Technical University of Denmark*, 7:15, 2008.
- [2] 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.
- [3] A. Gretton, O. Bousquet, Alex. Smola, and B. Schölkopf. Measuring statistical dependence with Hilbert-Schmidt norms. In *ALT*, 2005.
- [4] C. Cortes, M. Mohri, and A. Rostamizadeh. Algorithms for learning kernels based on centered alignment. *JMLR*, 13:795–828, 2012.
- [5] 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.

- [6] 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.