

AI6102 Machine Learning Methodologies & Applications

Chen Yongquan (G2002341D)

Nanyang Technological University

Assignment 1

## Question 1

There are two approaches to multi-class classification using logistic regression. The first approach is the one-vs-rest (OvR) scheme while the other is multinomial logistic regression.

### One-vs-rest

For the OvR scheme with regularization, we still use the cross-entropy loss defined in logistic regression for binary classification, defined as:

$$\ell(h(\mathbf{x}; \mathbf{w}), y) = - \sum_{i=1}^N [y_i \ln(h(\mathbf{x}_i; \mathbf{w})) - (1 - y_i) \ln(1 - h(\mathbf{x}_i; \mathbf{w}))] + \frac{\lambda}{2} \|\mathbf{w}\|_2^2$$

Where  $N$  is the number of training samples per batch while  $h$  is the logistic function.

The gradient descent update rule for binary classification is then defined as:

$$\mathbf{w}_{t+1} = \mathbf{w}_t + \rho \left( \sum_{i=1}^N (y_i - h(\mathbf{x}_i; \mathbf{w})) \mathbf{x}_i - \lambda \mathbf{w} \right)$$

Where  $t$  is the training iteration and  $\rho$  is the learning rate.

For a multi-class classification problem of  $C$  classes,  $\{0, 1, \dots, C - 1\}$ , we will then train  $C$  binary classifiers separately for all classes, each predicting the probability  $P(y = c | \mathbf{x})$  that  $\mathbf{x}$  belongs to the class  $c$  or not. Thus, we have  $C$  classifiers that give the predicted value for:

$$P(y = c | \mathbf{x}) = \frac{\exp(-\mathbf{w}_c^T \mathbf{x})}{1 + \exp(-\mathbf{w}_c^T \mathbf{x})}$$

Thus, for the OvR scheme, given an input  $\mathbf{x}$ , we use  $C$  classifiers to separately predict the probability that  $\mathbf{x}$  belongs to a class  $c$ . The predicted class is then the class with the highest probability out of all the classes:

$$\hat{y} = \underset{c}{\operatorname{argmax}} (P(y = c | \mathbf{x}))$$

Such an approach can result in probabilities for all classifiers summing up to more than 1.0. For single label cases, scikit-learn normalizes the probabilities for all classes such that they sum up to 1.0.

### Multinomial Logistic Regression

The second approach for multi-class classification uses the softmax function,  $g(\mathbf{w}^T \mathbf{x})$ , instead of the logistic function,  $h(\mathbf{w}^T \mathbf{x})$ , for generating the multi-class probability distribution. The softmax function is a generalization of the logistic function for more than 2 classes, and is defined as:

$$g(\mathbf{z}_i) = \frac{\exp(\mathbf{z}_i)}{\sum_{c=0}^{C-1} (\exp(\mathbf{z}_c))}$$

Where  $\mathbf{z}_i = -\mathbf{w}_c^T \mathbf{x}$  and the probabilities for all classes sum to 1.

If we fix the weight for class 0 to 0 and learn  $C - 1$  weights  $\mathbf{w}_c$  for classes  $c \in \{1, \dots, C - 1\}$ , then we can represent the predicted conditional probabilities in parametric form as:

$$\begin{aligned} \text{For } c > 0: \quad P(y = c|\mathbf{x}) &= \frac{\exp(-\mathbf{w}_c^T \mathbf{x})}{1 + \sum_{c=1}^{C-1} \exp(-\mathbf{w}_c^T \mathbf{x})} \\ \text{For } c = 0: \quad P(y = 0|\mathbf{x}) &= \frac{1}{1 + \sum_{c=1}^{C-1} \exp(-\mathbf{w}_c^T \mathbf{x})} \end{aligned}$$

Because we do not implement  $\mathbf{w}_c$  as a learnable parameter, we can assume its weight is 0 and  $\exp(-\mathbf{w}_0^T \mathbf{x}) = 1$ , giving us the 1 in the numerator and denominator in the second equation. Also, since the probabilities for the softmax sum to 1, we can get the conditional probability for class zero from the difference between 1 and the sums of the probabilities for all other classes:

$$\begin{aligned} \sum_{c=0}^{C-1} P(y = c|\mathbf{x}) &= 1 \\ P(y = 0|\mathbf{x}) &= 1 - \sum_{c=1}^{C-1} P(y = c|\mathbf{x}) \end{aligned}$$

We can also start by define the model as a set of logits:

$$\begin{aligned} \ln \frac{P(y = 1|\mathbf{x})}{P(y = 0|\mathbf{x})} &= -\mathbf{w}_1^T \mathbf{x} \\ &\vdots \\ \ln \frac{P(y = C - 1|\mathbf{x})}{P(y = 0|\mathbf{x})} &= -\mathbf{w}_{C-1}^T \mathbf{x} \end{aligned}$$

This set of logits then gives us the equations:

$$\begin{aligned} P(y = 1|\mathbf{x}) &= P(y = 0|\mathbf{x}) \exp(-\mathbf{w}_1^T \mathbf{x}) \\ &\vdots \\ P(y = C - 1|\mathbf{x}) &= P(y = 0|\mathbf{x}) \exp(-\mathbf{w}_{C-1}^T \mathbf{x}) \end{aligned}$$

Since all the conditional probabilities sum to 1, we have the equation:

$$\begin{aligned} P(y = 0|\mathbf{x}) &= 1 - \sum_{c=1}^{C-1} P(y = c|\mathbf{x}) = 1 - \sum_{c=1}^{C-1} (P(y = 0|\mathbf{x}) \exp(-\mathbf{w}_c^T \mathbf{x})) \\ &= 1 - P(y = 0|\mathbf{x}) \sum_{c=1}^{C-1} \exp(-\mathbf{w}_c^T \mathbf{x}) \end{aligned}$$

Then:

$$\begin{aligned} 1 &= P(y = 0|\mathbf{x}) + P(y = 0|\mathbf{x}) \sum_{c=1}^{C-1} \exp(-\mathbf{w}_c^T \mathbf{x}) \\ 1 &= P(y = 0|\mathbf{x}) + \left( 1 + \sum_{c=1}^{C-1} \exp(-\mathbf{w}_c^T \mathbf{x}) \right) P(y = 0|\mathbf{x}) \\ P(y = 0|\mathbf{x}) &= \frac{1}{1 + \sum_{c=1}^{C-1} \exp(-\mathbf{w}_c^T \mathbf{x})} \end{aligned}$$

Thus, substituting into the set of logits we defined at the beginning, we arrive at the same parametric form for multinomial logistic regression:

$$P(y = c|\mathbf{x}) = \frac{\exp(-\mathbf{w}_c^T \mathbf{x})}{1 + \sum_{c=1}^{C-1} \exp(-\mathbf{w}_c^T \mathbf{x})}$$

We can also represent the parametric form more compactly as:

$$P(y|\mathbf{x}; \mathbf{w}) = \prod_{c=0}^{C-1} [g(\mathbf{x}; \mathbf{w}_c)]^{y_c} = \sum_{c=0}^{C-1} y_c \ln(g(\mathbf{x}; \mathbf{w}_c))$$

Thus, given a set of  $N$  training samples, in order to maximize the likelihood of  $P(y_i|\mathbf{x}_i; \mathbf{w}_c)$  with respect to  $\mathbf{w}_c$ , where  $i \in \{1, \dots, N\}$ , we have to find a solution for  $\mathbf{w}_c$ :

$$\widehat{\mathbf{w}}_c = \underset{\mathbf{w}_c}{\operatorname{argmax}} \prod_{i=1}^N P(y_i|\mathbf{x}_i; \mathbf{w}_c) = \underset{\mathbf{w}_c}{\operatorname{argmax}} \sum_{i=1}^N \sum_{c=0}^{C-1} (y_i \ln(g(\mathbf{x}_i; \mathbf{w}_c)))$$

Then, the cross-entropy loss for a batch of  $N$  training samples using softmax function for prediction is given by:

$$\ell(g(\mathbf{x}; \mathbf{w}), y_i) = - \sum_{i=1}^N \sum_{c=0}^{C-1} y_i \ln(g(\mathbf{x}_i; \mathbf{w}_c))$$

In order to perform gradient descent, we need the first derivative of the loss function with respect to parameter  $\mathbf{w}_c$ :

$$\begin{aligned}\mathbf{w}_{c,t+1} &= \mathbf{w}_{c,t} - \rho \frac{\partial E(\mathbf{w}_c)}{\partial \mathbf{w}_c} \\ \frac{\partial E(\mathbf{w}_c)}{\partial \mathbf{w}_c} &= \frac{\partial (-\sum_{i=1}^N \sum_{c=0}^{C-1} (y_i \ln(g(\mathbf{x}_i; \mathbf{w}_c))))}{\partial \mathbf{w}_c} \\ &= -\sum_{i=1}^N \sum_{c=0}^{C-1} \left( \frac{\partial y_i \ln(g(\mathbf{x}_i; \mathbf{w}_c))}{\partial \mathbf{w}_c} \right) \quad (1)\end{aligned}$$

$$\frac{\partial y_i \ln(g(\mathbf{x}_i; \mathbf{w}_c))}{\partial \mathbf{w}_c} = y_i \left( \frac{1}{g(\mathbf{x}_i; \mathbf{w}_c)} \right) \left( \frac{\partial g(\mathbf{x}_i; \mathbf{w}_c)}{\partial \mathbf{w}_c} \right) \quad (2)$$

Let  $z_{id} = -\mathbf{w}_d^T \mathbf{x}_i$ , then when  $c = d$ :

$$\begin{aligned}\frac{\partial g(z_{id})}{\partial z_{ic}} &= \frac{\partial \frac{\exp(z_{ic})}{\sum_{j=0}^{C-1} \exp(z_{ij})}}{\partial z_{ic}} \\ &= \frac{(\sum_{j=0}^{C-1} \exp(z_{ij}))(\exp(z_{ic})) - (\exp(z_{ic}))(\exp(z_{ic}))}{(\sum_{j=0}^{C-1} \exp(z_{ij}))^2} \\ &= \frac{(\exp(z_{ic}))(\sum_{j=0}^{C-1} \exp(z_{ij}) - \exp(z_{ic}))}{(\sum_{j=0}^{C-1} \exp(z_{ij}))^2} \\ &= \left( \frac{\exp(z_{ic})}{\sum_{j=0}^{C-1} \exp(z_{ij})} \right) \left( \frac{\sum_{j=0}^{C-1} \exp(z_{ij}) - \exp(z_{ic})}{\sum_{j=0}^{C-1} \exp(z_{ij})} \right) \\ &= \left( \frac{\exp(z_{ic})}{\sum_{j=0}^{C-1} \exp(z_{ij})} \right) \left( 1 - \frac{\exp(z_{ic})}{\sum_{j=0}^{C-1} \exp(z_{ij})} \right) \\ &= g(\mathbf{x}_i; \mathbf{w}_c)(1 - g(\mathbf{x}_i; \mathbf{w}_c)) \quad (3)\end{aligned}$$

When  $c \neq d$ :

$$\begin{aligned}\frac{\partial g(z_{id})}{\partial z_{ic}} &= \frac{0 - (\exp(z_{ic}))(\exp(z_{id}))}{(\sum_{j=0}^{C-1} \exp(z_{ij}))^2} = -\frac{\exp(z_{ic})}{\sum_{j=0}^{C-1} \exp(z_{ij})} \frac{\exp(z_{id})}{\sum_{j=0}^{C-1} \exp(z_{ij})} \\ &= -g(\mathbf{x}_i; \mathbf{w}_c)g(\mathbf{x}_i; \mathbf{w}_d) \quad (4)\end{aligned}$$

Then, by chain rule:

$$\frac{\partial g(z_{id})}{\partial \mathbf{w}_c} = \frac{\partial g(z_{id})}{\partial z_{ic}} \cdot \frac{\partial z_{ic}}{\partial \mathbf{w}_c} = \frac{\partial g(z_{id})}{\partial z_{ic}} \frac{\partial (-\mathbf{w}_c^T \mathbf{x}_i)}{\partial \mathbf{w}_c} = \frac{\partial g(z_{id})}{\partial z_{ic}} (-\mathbf{x}_i)$$

Substituting back into (2), (3) and (4) back into (1):

$$\begin{aligned} \frac{\partial E(\mathbf{w}_c)}{\partial \mathbf{w}_c} &= - \sum_{i=1}^N \sum_{c=0}^{C-1} \left( \frac{\partial y_i \ln(g(z_{ic}))}{\partial \mathbf{w}_c} \right) \\ &= - \sum_{i=1}^N \sum_{c=0}^{C-1} \left( y_i \left( \frac{1}{g(z_{ic})} \right) \left( \frac{\partial g(z_{ic})}{\partial \mathbf{w}_c} \right) \right) \\ &= - \sum_{i=1}^N \left( \left( \frac{y_i}{g(z_{ic})} \right) \left( \frac{\partial g(z_{ic})}{\partial \mathbf{w}_c} \right) + \sum_{d \neq c}^{C-1} \left( \frac{y_i}{g(z_{ic})} \frac{\partial g(z_{id})}{\partial \mathbf{w}_c} \right) \right) \\ &= - \sum_{i=1}^N \left( \left( \frac{y_i}{g(z_{ic})} \right) \left( \frac{\partial g(z_{ic})}{\partial z_{ic}} \right) + \sum_{d \neq c}^{C-1} \left( \frac{y_i}{g(z_{ic})} \frac{\partial g(z_{id})}{\partial z_{ic}} \right) \right) \cdot \frac{\partial z_{ic}}{\partial \mathbf{w}_c} \\ &= - \sum_{i=1}^N \left( \left( \frac{y_i g(z_{ic}) (1 - g(z_{ic}))}{g(z_{ic})} \right) + \sum_{d \neq c}^{C-1} \left( \frac{-y_i g(z_{ic}) g(z_{id})}{g(z_{ic})} \right) \right) (-\mathbf{x}_i) \\ &= - \sum_{i=1}^N \left( \left( y_i (1 - g(z_{ic})) \right) + \sum_{d \neq c}^{C-1} (-y_i g(z_{id})) \right) (-\mathbf{x}_i) \\ &= - \sum_{i=1}^N \left( y_i - y_i g(z_{ic}) - \sum_{d \neq c}^{C-1} (y_i g(z_{id})) \right) (-\mathbf{x}_i) \\ &= - \sum_{i=1}^N \left( y_i - \sum_{c=0}^{C-1} (y_i g(z_{ic})) \right) (-\mathbf{x}_i) \\ &= - \sum_{i=1}^N \left( y_i - g(z_{ic}) \sum_{c=0}^{C-1} y_i \right) (-\mathbf{x}_i) \\ &= - \sum_{i=1}^N (g(\mathbf{x}_i; \mathbf{w}_c) - y_i) (\mathbf{x}_i) \end{aligned}$$

Finally, substituting everything back into the gradient descent rule:

$$\mathbf{w}_{t+1,c} = \mathbf{w}_{t,c} + \rho \sum_{i=1}^N (g(\mathbf{x}_i; \mathbf{w}_c) - y_i) \mathbf{x}_i$$

With L2 regularization, the loss function is augmented by a regularization parameter as shown:

$$\ell(g(\mathbf{x}; \mathbf{w}), y_i) = - \sum_{i=1}^N \left[ \sum_{c=0}^{C-1} y_i \ln(g(\mathbf{x}_i; \mathbf{w}_c)) \right] + \frac{\lambda}{2} \|\mathbf{w}_c\|_2^2$$

Consequently, the partial derivative for the loss function with respect to  $\mathbf{w}$  becomes:

$$\frac{\partial E(\mathbf{w}_c)}{\partial \mathbf{w}_c} = - \sum_{i=1}^N (g(\mathbf{x}_i; \mathbf{w}_c) - y_i) (\mathbf{x}_i) + \lambda \mathbf{w}_c$$

And the gradient descent rule:

$$\mathbf{w}_{t+1,c} = \mathbf{w}_{t,c} + \rho \left[ \sum_{i=1}^N (g(\mathbf{x}_i; \mathbf{w}_c) - y_i) \mathbf{x}_i - \lambda \mathbf{w}_c \right]$$

Depending on the convention we use for the softmax function,

$$\text{Maximum convention: } g(z_i) = \frac{\exp(\mathbf{w}_i^T \mathbf{x})}{\sum_{c=0}^{C-1} \exp(\mathbf{w}_c^T \mathbf{x})}$$

$$\text{Minimum convention: } g(z_i) = \frac{\exp(-\mathbf{w}_i^T \mathbf{x})}{\sum_{c=0}^{C-1} \exp(-\mathbf{w}_c^T \mathbf{x})}$$

Where  $z_{ic} = \pm \mathbf{w}_c^T \mathbf{x}_i$ , the resulting partial derivative would also follow the change in sign,  $\frac{\partial z_{ic}}{\partial \mathbf{w}_c} = \pm \mathbf{x}_i$ . Thus, if we follow the positive convention, the update rule would then be as such:

$$\mathbf{w}_{t+1,c} = \mathbf{w}_{t,c} + \rho \left[ \sum_{i=1}^N (y_i - g(\mathbf{x}_i; \mathbf{w}_c)) \mathbf{x}_i - \lambda \mathbf{w}_c \right]$$

## Question 2

1.

This is the code for importing the training and test data and setting up the basic environment.

```
from sklearn.datasets import load_svmlight_files
from sklearn.model_selection import KFold
from sklearn.svm import LinearSVC
from sklearn import datasets
import numpy as np
import pandas as pd
import sklearn as sk
np.set_printoptions(linewidth = 140)

C = [0.01, 0.1 , 1.0, 10, 100]
train_X, train_y, test_X, test_y =
sk.datasets.load_svmlight_files(['a5a', 'a5a.t'])
```

This is the code used for performing the 3-fold cross-validation and for creating the LinearSVC instance:

```
def kfLSVC(X, y, C = 1, n_splits = 3, max_iter = 1000):
    svc = sk.svm.LinearSVC(C = C, max_iter = max_iter)
    kf = sk.model_selection.KFold(n_splits)
    scores = []
    for cv_train, cv_test in kf.split(X):
        svc.fit(X[cv_train], y[cv_train])
        scores.append(svc.score(X[cv_test], y[cv_test]))
    return svc, scores

svc_set = []
scores = []
for C_i in C:
    output = kfLSVC(train_X, train_y, C_i)
    svc_set.append(output[0])
    scores.append(output[1])
    print("C: %-5s Scores: [%-20s%-20s%-20s], Average: %-20s n_iter_: %s" %
          (str(C_i)+', ',
            str(output[1][0])+', ', str(output[1][1])+', ', str(output[1][2]),
            str(np.average(output[1]))+', ', output[0].n_iter_))

print('')

for i, svc in enumerate(svc_set):
    print("C: " + str(C[i]) + " TScore: " + str(svc.score(test_X, test_y)))
```



**2.**

The 3-fold cross-validation results of the set of 5 candidate values of  $C$  on the a5a training set is as shown below:

	Accuracy of linear SVMs
$C = 0.01$	0.839569691
$C = 0.1$	0.842531961
$C = 1$	0.839101964
$C = 10$	0.839101964
$C = 100$	0.797318366

The results above were obtained with default parameters for LinearSVC and the KFold model except for the values of  $C$ . We note that the maximum number of iterations was reached by the classifier when  $C$  was set to 10 and 100. The default maximum number of iterations is 100.

**3.**

Using  $C = 0.1$ , which gave the best mean accuracy in the 3-fold cross validation, we train a LinearSVC model using the entire training set. The code used is as shown below:

```
idx = 1
final_svc, final_svc_score = kfLSVC(train_X, train_y, C[idx])
print("C: %-5s Scores: [%-20s%-20s%-20s], Average: %-20s n_iter_: %s" %
      (str(C[idx])+' ',
       str(final_svc_score[0])+' ',
       str(final_svc_score[1])+' ',
       str(final_svc_score[2]),
       str(np.average(final_svc_score))+' ',
       final_svc.n_iter_))
print("C: %-5s TScore: %-20s" %
      (str(C[idx])+' ',
       final_svc.score(test_X, test_y)))
```

The test results obtained from this classifier is as shown:

	Accuracy of linear SVMs
$C = 0.1$	0.8451065131755077

We note that by setting max iterations to 100000, we can observe convergence in the classifiers for when  $C = 10$  and  $C = 100$ .

### Question 3

The optimization problem statement for ridge regression is given by:

$$\hat{\mathbf{w}} = \underset{\mathbf{w}}{\operatorname{argmin}} \left[ \left( \frac{1}{2} \sum_{i=1}^N (\mathbf{w}^T \mathbf{x}_i - y_i)^2 \right) + \frac{\lambda}{2} \|\mathbf{w}\|_2^2 \right]$$

Applying the kernel trick to map  $\mathbf{x}$  to a higher dimension to solve non-linear regression problems by finding a linear function in the higher dimension:

$$k(\mathbf{x}_i, \mathbf{x}_j) = \varphi(\mathbf{x}_i) \cdot \varphi(\mathbf{x}_j)$$

$$\hat{\mathbf{w}} = \underset{\mathbf{w}}{\operatorname{argmin}} \left[ \left( \frac{1}{2} \sum_{i=1}^N (\mathbf{w}^T \varphi(\mathbf{x}_i) - y_i)^2 \right) + \frac{\lambda}{2} \|\mathbf{w}\|_2^2 \right]$$

Then, to find a closed-form solution, we have to solve for:

$$\frac{\partial \left[ \left( \frac{1}{2} \sum_{i=1}^N (\mathbf{w}^T \varphi(\mathbf{x}_i) - y_i)^2 \right) + \frac{\lambda}{2} \|\mathbf{w}\|_2^2 \right]}{\partial \mathbf{x}} = 0$$

$$\frac{\partial \left[ \frac{1}{2} \sum_{i=1}^N (\mathbf{w}^T \varphi(\mathbf{x}_i) - y_i)^2 \right]}{\partial \mathbf{x}} + \frac{\partial \left[ \frac{\lambda}{2} \|\mathbf{w}\|_2^2 \right]}{\partial \mathbf{x}} = 0$$

$$\frac{\partial \left[ \frac{1}{2} \sum_{i=1}^N (\mathbf{w}^T \varphi(\mathbf{x}_i) - y_i)^2 \right]}{\partial \mathbf{x}} + \lambda \mathbf{w} = 0$$

$$\sum_{i=1}^N (\mathbf{w}^T \varphi(\mathbf{x}_i) - y_i) \varphi(\mathbf{x}_i) + \lambda \mathbf{w} = 0$$

$$\sum_{i=1}^N (\mathbf{w}^T \varphi(\mathbf{x}_i)) \varphi(\mathbf{x}_i) - \sum_{i=1}^N y_i \varphi(\mathbf{x}_i) + \lambda \mathbf{w} = 0$$

$$\sum_{i=1}^N \varphi(\mathbf{x}_i) (\mathbf{w}^T \varphi(\mathbf{x}_i)) - \sum_{i=1}^N y_i \varphi(\mathbf{x}_i) + \lambda \mathbf{w} = 0$$

$$\left( \sum_{i=1}^N (\varphi(\mathbf{x}_i) \varphi(\mathbf{x}_i)^T) \right) \mathbf{w} - \sum_{i=1}^N y_i \varphi(\mathbf{x}_i) + \lambda \mathbf{w} = 0$$

Let  $\Phi$  represent the matrix for  $\varphi(x_i)$  for  $N$  data samples of  $m$  dimension as below:

$\Phi$			
$\varphi(x_{01})$	$\varphi(x_{02})$	$\dots$	$\varphi(x_{0N})$
$\varphi(x_{11})$	$\varphi(x_{12})$	$\dots$	$\varphi(x_{1N})$
$\dots$	$\dots$	$\dots$	$\dots$
$\varphi(x_{m1})$	$\varphi(x_{m2})$	$\dots$	$\varphi(x_{mN})$

Then:

$$\sum_{i=1}^N (\varphi(x_i) \varphi(x_i)^T) = \Phi \Phi^T$$

Substituting into our intermediate solution previously:

$$\Phi \Phi^T \mathbf{w} - \sum_{i=1}^N y_i \varphi(x_i) + \lambda \mathbf{w} = 0$$

$$\Phi \Phi^T \mathbf{w} - \Phi \mathbf{y} + \lambda \mathbf{w} = 0$$

$$\Phi \Phi^T \mathbf{w} + \lambda \mathbf{w} = \Phi \mathbf{y}$$

$$(\Phi \Phi^T + \lambda I) \mathbf{w} = \Phi \mathbf{y}$$

$$\mathbf{w} = (\Phi \Phi^T + \lambda I)^{-1} \Phi \mathbf{y}$$

$$\mathbf{w} = \Phi^T (\Phi^T \Phi + \lambda I)^{-1} \mathbf{y}$$

Let

$$K = k(x_i, x_j) = \varphi(x_i) \cdot \varphi(x_j)$$

Then,

$$\mathbf{w} = \Phi^T (K + \lambda I)^{-1} \mathbf{y}$$