

APPENDIX A  
BACKPROPAGATION FOR NON-DECOMPOSABLE OBJECTIVES

It worth to mention that the original backpropagation [73] is applicable to objective function, which value on the whole training dataset (or mini-batch)  $\mathbb{T}$  can be approximated by the averaged value on every single sample. This is possible due to the stochastic approximation theory [74]

$$E(\mathbb{T}) = \frac{1}{N} \sum_{(\mathbf{x}, \mathbf{y}) \in \mathbb{T}} E(\mathbf{x}, \mathbf{y}), \quad (23)$$

where  $E$  is the objective function, pair  $(\mathbf{x}, \mathbf{y})$  is a sample and its label. These objective functions (23) are known as the *decomposable* [75]. In machine learning we usually apply indirect optimization, regarding the objective function. Technically we try to improve a particular performance measure  $P$  (for some practical application), which is calculated on the test dataset. Whereas, during training stage, we optimize model parameters using different objective function  $E$  in the hope that in general it will improve performance  $P$ . This is indirect optimization, since such metrics are intractable or also known as *non-decomposable functions*. Another ambiguity is that we often optimize the same objective function (MSE, cross-entropy, KL-divergence, e.t.c.) for a different applications and for different performance metrics.

**Definition. 1** *Non-decomposable objective function is not decomposed into expectations over individual examples [76]. For example, these functions are F1-score, Precision/Recall Break Even Point (PRBEP), Precision at  $k$  (Prec@ $k$ ), ROCArea, EER, e.t.c.*

For a *non-decomposable* objective functions, parameters updating rule of backpropagation algorithm is assumed to be the same as for *decomposable* functions, except that *error signal*  $\delta^n$  is calculated on the whole mini-batch  $\mathbb{T}$

$$\begin{aligned} \mathbf{W}_n &\leftarrow \mathbf{W}_n - \eta \frac{\partial E(\mathbb{T})}{\partial \mathbf{W}_n} = \mathbf{W}_n - \eta \delta^n(\mathbb{T}) \frac{\partial \mathbf{z}^n}{\partial \mathbf{W}_n}, \quad \text{where} \\ \delta^n(\mathbb{T}) &= \frac{\partial E(\mathbb{T})}{\partial \mathbf{z}^n}, \end{aligned} \quad (24)$$

here  $\mathbf{W}_n$  are neural network parameters for a range of  $n = \overline{0, L}$  number of layers;  $E$  is a non-decomposable objective function;  $\mathbf{z}^n$  is a *pre-activation* output or output of  $n$ -th layer of network before applying activation function;  $\eta$  is a learning rate. In the *decomposable* objective function case (23), we can represent the gradient of the function as the averaged gradient calculated per-sample [74]. Whereas, for *non-decomposable* functions in (24), it is similar to the original backpropagation with SGD [73] calculated on a *single sample*, here as a single sample we treat the whole mini-batch  $\mathbb{T}$ . It is very rough assumption, though it works in practice with *adaptive learning rate methods*.

#### A. Backpropagation for MFoM-micro-F1

The micro-F1 objective function for multiclass case of  $M$  classes in terms of discrete counts *true positive* (TP), *false positive* (FP) and *false negative* (FN)

$$F_1 = \frac{2TP}{FP + 2TP + FN} = \frac{2 \sum_{k=1}^M TP_k}{\sum_{k=1}^M FP_k + 2 \sum_{k=1}^M TP_k + \sum_{k=1}^M FN_k}. \quad (25)$$

The problem with micro-F1 is that it is the discrete function, thus we can not directly differentiate and optimize it. After introducing the MFoM framework: the terms of *discriminative functions* (2), *misclassification measure* (4) and *smooth error function* (6), - we are ready to represent micro-F1 as the smooth continues function  $\hat{F}_1$  and then optimize it. We use *smooth error function* (6) in order to approximate the discrete counts on a mini-batch  $\mathbb{T}$

$$TP \approx \widehat{TP} \triangleq \sum_{k=1}^M \widehat{TP}_k = \sum_{k=1}^M \sum_{\mathbf{x} \in \mathbb{T}} (1 - l_k(\mathbf{z})) \cdot y_k, \quad (26)$$

$$FP \approx \widehat{FP} \triangleq \sum_{k=1}^M \widehat{FP}_k = \sum_{k=1}^M \sum_{\mathbf{x} \in \mathbb{T}} (1 - l_k(\mathbf{z})) \cdot \bar{y}_k, \quad (27)$$

$$FN \approx \widehat{FN} \triangleq \sum_{k=1}^M \widehat{FN}_k = \sum_{k=1}^M \sum_{\mathbf{x} \in \mathbb{T}} l_k(\mathbf{z}) \cdot y_k, \quad (28)$$

where  $y_k$  and  $\bar{y}_k$  are labels of the corresponding sample  $\mathbf{x}$  or

$$y_k = \mathbf{1}(x \in C_k) = \begin{cases} 1, & \text{if } x \in C_k, \\ 0, & \text{if } x \notin C_k. \end{cases} \quad (29)$$

and

$$\bar{y}_k = \mathbf{1}(x \notin C_k) = \begin{cases} 0, & \text{if } x \in C_k, \\ 1, & \text{if } x \notin C_k. \end{cases} \quad (30)$$

where  $C_k$  is the set of samples where  $k^{\text{th}}$  label is ‘‘on’’ (label  $k$  equals 1), indicator functions are  $\mathbf{1}(\cdot)$ . We denote the sum of  $\widehat{\text{TP}}$  and  $\widehat{\text{FN}}$  as

$$\begin{aligned} |C| &\triangleq \sum_{k=1}^M |C_k| = \sum_{k=1}^M \widehat{\text{TP}}_k + \widehat{\text{FN}}_k = \\ &= \sum_{k=1}^M \sum_{\mathbf{x} \in \mathbb{T}} (1 - l_k(\mathbf{z})) \cdot y_k + l_k(\mathbf{z}) \cdot y_k = \sum_{k=1}^M \sum_{\mathbf{x} \in \mathbb{T}} y_k, \end{aligned} \quad (31)$$

Therefore, the term  $|C|$  is the count of unit ‘‘1’’ labels of samples  $\mathbf{x}$  across all mini-batch  $\mathbb{T}$ , i.e. the number of positive labels. The  $|C|$  is a constant value for the current mini-batch  $\mathbb{T}$ , because we calculate it one time for the whole  $\mathbb{T}$ .

Finally, we get the approximation of discrete micro-F1 from (25), we call it *MFoM-micro-F1* objective function

$$F_1 \approx \widehat{F}_1 = \frac{2\widehat{\text{TP}}}{\widehat{\text{FP}} + \widehat{\text{TP}} + |C|} \quad (32)$$

Our task is to minimize the objective function

$$E = 1 - \widehat{F}_1 \rightarrow \arg \min_{\mathbb{W}} \quad (33)$$

where a parameter set  $\mathbb{W} = \{\mathbf{W}_n | n = \overline{0, L}\}$  of the network consisting of  $L + 1$  layers. We are able to calculate the partial derivatives of all network parameters  $\mathbb{W}$ , if we find the network *error signal* term  $\delta^n(\mathbb{T})$ , starting from the output layer  $n = L$

$$\begin{aligned} \delta^L(\mathbb{T}) &= \frac{\partial E}{\partial \mathbf{z}} = -\frac{\partial \widehat{F}_1}{\partial \mathbf{z}} = -\left[ \frac{2\widehat{\text{TP}}}{\widehat{\text{FP}} + \widehat{\text{TP}} + |C|} \right]'_{\mathbf{z}} = \\ &= -\frac{2\widehat{\text{TP}}' \cdot (\widehat{\text{FP}} + \widehat{\text{TP}} + |C|) - 2\widehat{\text{TP}} \cdot (\widehat{\text{FP}}' + \widehat{\text{TP}}')}{(\widehat{\text{FP}} + \widehat{\text{TP}} + |C|)^2} = \\ &= -\frac{2\widehat{\text{TP}}' \cdot \widehat{\text{FP}} + 2\widehat{\text{TP}}' \cdot |C| - 2\widehat{\text{TP}} \cdot \widehat{\text{FP}}'}{(\widehat{\text{FP}} + \widehat{\text{TP}} + |C|)^2} = \\ &= -\frac{-2\widehat{\text{FN}}' \cdot \widehat{\text{FP}} - 2\widehat{\text{FN}}' \cdot |C| - 2\widehat{\text{TP}} \cdot \widehat{\text{FP}}'}{(\widehat{\text{FP}} + \widehat{\text{TP}} + |C|)^2} = \\ &= \frac{2 \cdot [(|C| + \widehat{\text{FP}}) \cdot \widehat{\text{FN}}' + \widehat{\text{TP}} \cdot \widehat{\text{FP}}']}{(\widehat{\text{FP}} + \widehat{\text{TP}} + |C|)^2}. \end{aligned}$$

Thus,

$$\delta^L(\mathbb{T}) = \frac{\partial E}{\partial \mathbf{z}} = A \left( w_1 \frac{\partial \widehat{\text{FN}}}{\partial \mathbf{z}} + w_2 \frac{\partial \widehat{\text{FP}}}{\partial \mathbf{z}} \right), \quad (34)$$

where constants, which are calculated on the whole current mini-batch  $\mathbb{T}$

$$A = \frac{2}{(\widehat{\text{FP}} + \widehat{\text{TP}} + |C|)^2}, \quad w_1 = |C| + \widehat{\text{FP}}, \quad w_2 = \widehat{\text{TP}}. \quad (35)$$

If we plug in the  $\widehat{\text{FP}}$  and  $\widehat{\text{FN}}$  from (27) and (28) to the  $\delta^L$  in eq. (34)

$$\begin{aligned} \delta^L(\mathbb{T}) &= A \left( w_1 \frac{\partial \widehat{\text{FN}}}{\partial \mathbf{z}} + w_2 \frac{\partial \widehat{\text{FP}}}{\partial \mathbf{z}} \right) = \\ &= A \left( w_1 \cdot \sum_{k=1}^M \sum_{\mathbf{x} \in \mathbb{T}} \frac{\partial l_k(\mathbf{z})}{\partial \mathbf{z}} \cdot y_k - w_2 \cdot \sum_{k=1}^M \sum_{\mathbf{x} \in \mathbb{T}} \frac{\partial l_k(\mathbf{z})}{\partial \mathbf{z}} \cdot \bar{y}_k \right) = \end{aligned}$$

$$= A \left( \sum_{\mathbf{x} \in \mathbb{T}} \sum_{k=1}^M \frac{\partial l_k(\mathbf{z})}{\partial \mathbf{z}} [w_1 \cdot y_k - w_2 \cdot \bar{y}_k] \right),$$

It means that if  $k^{th}$  label of a sample  $\mathbf{x}$  is unit (i.e.,  $y_k = 1$  and  $\bar{y}_k = 0$ ), then we multiply by the weight constant  $w_1$ , otherwise by the constant  $-w_2$ . Then we rewrite in the vector form the last result

$$\delta^L(\mathbb{T}) = A \sum_{\mathbf{x} \in \mathbb{T}} \left( \frac{\partial l_1(\mathbf{z})}{\partial \mathbf{z}}, \dots, \frac{\partial l_k(\mathbf{z})}{\partial \mathbf{z}}, \dots, \frac{\partial l_M(\mathbf{z})}{\partial \mathbf{z}} \right) \cdot \begin{pmatrix} w_1 \cdot y_1 - w_2 \cdot \bar{y}_1 \\ \dots \\ w_1 \cdot y_k - w_2 \cdot \bar{y}_k \\ \dots \\ w_1 \cdot y_M - w_2 \cdot \bar{y}_M \end{pmatrix}. \quad (36)$$

In (36), the left vector under the sum is the transposed Jacobian, because we find the partial derivatives with respect to the vector of multiple variables  $\mathbf{z} = (z_1, \dots, z_m, \dots, z_M)^\top$

$$\begin{aligned} J_{\mathbf{z}}^\top \mathcal{L} &= \left( \frac{\partial l_1(\mathbf{z})}{\partial \mathbf{z}}, \dots, \frac{\partial l_k(\mathbf{z})}{\partial \mathbf{z}}, \dots, \frac{\partial l_M(\mathbf{z})}{\partial \mathbf{z}} \right) = \\ &= \begin{pmatrix} \frac{\partial l_1(\mathbf{z})}{\partial z_1} & \dots & \frac{\partial l_M(\mathbf{z})}{\partial z_1} \\ \dots & \frac{\partial l_k(\mathbf{z})}{\partial z_m} & \dots \\ \frac{\partial l_1(\mathbf{z})}{\partial z_M} & \dots & \frac{\partial l_M(\mathbf{z})}{\partial z_M} \end{pmatrix}, \end{aligned} \quad (37)$$

where  $m = \overline{1, M}$  is the row index of the variables  $z_m$  and  $k = \overline{1, M}$  is the column index for the functions  $l_k(\mathbf{z})$ .

$$\delta^L(\mathbb{T}) = A \cdot \sum_{(\mathbf{x}, \mathbf{y}) \in \mathbb{T}} J_{\mathbf{z}}^\top \mathcal{L}(\mathbf{x}) \cdot (w_1 \cdot \mathbf{y} - w_2 \cdot \bar{\mathbf{y}}). \quad (38)$$

**NOTE:** we can interpret (38) as for every pair  $(\mathbf{x}, \mathbf{y})$  from mini-batch  $\mathbb{T}$  we calculate the value of the transposed Jacobian and find the weighted linear combination of its columns with  $w_1$  or  $-w_2$ . The weight constants  $w_1$  or  $w_2$  are defined by the labels  $y_k$  or  $\bar{y}_k$  of sample. Then we sum up all weighted Jacobians across all samples in a mini-batch  $\mathbb{T}$ . Thus, we get “decomposed” error signal  $\delta^L$  per each sample  $\mathbf{x}$ , i.e. MFoM framework allowed us to archive it.

### B. Jacobian for the Units-vs-Zeros Misclassification Measure

In this section, we infer a Jacobian for the *units-vs-zeros* misclassification measure (4). Units-vs-zeros misclassification measure was adopted for multi-label classification and proposed in [32]. First, we consider the example A.1 of the *units-vs-zeros* misclassification measure with sample  $\mathbf{x}$  having multiple label vector  $\mathbf{y}$ .

**Example A.1** Let we have training pair  $(\mathbf{x}, \mathbf{y})$  and  $\mathbf{y} = (1, 1, 0, 0)^\top$ , see Fig. 8. Then, we have *units-vs-zeros* misclassification as

$$\begin{aligned} \psi_1(\mathbf{z}) &= -g_1 + \ln \left[ \frac{1}{2} (e^{g_3} + e^{g_4}) \right], \\ \psi_2(\mathbf{z}) &= -g_2 + \ln \left[ \frac{1}{2} (e^{g_3} + e^{g_4}) \right], \\ \psi_3(\mathbf{z}) &= -g_3 + \ln \left[ \frac{1}{2} (e^{g_1} + e^{g_2}) \right], \\ \psi_4(\mathbf{z}) &= -g_4 + \ln \left[ \frac{1}{2} (e^{g_1} + e^{g_2}) \right], \end{aligned}$$

where  $g_k$  is the output score of a neural network.

The *units-vs-zeros* misclassification measure (4) can be rewritten in the convenient form for inference of derivatives

$$\psi_k(\mathbf{z}) = -g_k + \bar{y}_k \cdot u + y_k \cdot v, \quad (39)$$

where Kolmogorov *f-mean* of *unit models*

$$u = \ln \left[ \frac{1}{\sum_{i=1}^M y_i} \cdot \langle \mathbf{y}, \exp(\mathbf{g}) \rangle \right] = \ln \left[ \frac{1}{\sum_{i=1}^M y_i} \cdot (y_1 e^{g_1} + \dots + y_k e^{g_k} + \dots + y_M e^{g_M}) \right], \quad (40)$$

and Kolmogorov *f-mean* of *zero models*

$$v = \ln \left[ \frac{1}{\sum_{i=1}^M \bar{y}_i} \cdot \langle \bar{\mathbf{y}}, \exp(\mathbf{g}) \rangle \right] = \ln \left[ \frac{1}{\sum_{i=1}^M \bar{y}_i} \cdot (\bar{y}_1 e^{g_1} + \dots + \bar{y}_k e^{g_k} + \dots + \bar{y}_M e^{g_M}) \right], \quad (41)$$

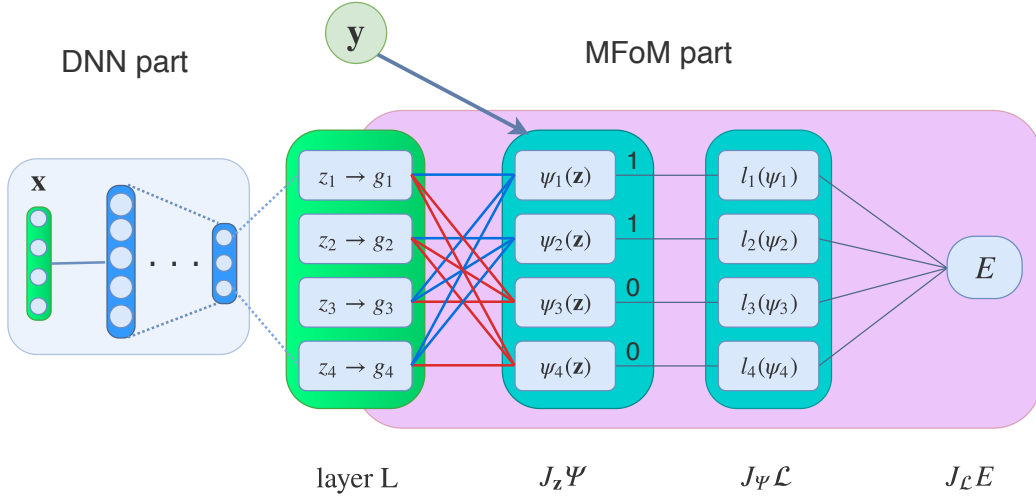


Fig. 8. Example of the extended DNN with the MFoM objective function and *units-vs-zeros* misclassification measure:  $\Psi$  is the vector of misclassification measure,  $\mathcal{L}$  is the vector of smooth error count,  $E$  is the smoothed MFoM-based objective;  $J_{\mathcal{L}}E$ ,  $J_{\Psi}\mathcal{L}$ ,  $J_{\mathbf{z}}\Psi$  are Jacobians.

where binary vector of labels is  $\mathbf{y}$  and its inverse is  $\bar{\mathbf{y}}$ . Thus, if the current sample  $\mathbf{x}$  is labeled as 1 in the ground-truth for the class  $C_k$  (i.e.  $y_k = 1$  or  $\bar{y}_k = 0$ ), the competing models will be considered only those with labels 0, and we find an average of these competing zero models, i.e. it is the equation (41). The misclassification measure for that sample is calculated as

$$\psi_k(\mathbf{z}) = -g_k + y_k \cdot v. \quad (42)$$

Otherwise, if the sample  $\mathbf{x}$  is labeled as 0 in the ground-truth (i.e.  $y_k = 0$  or  $\bar{y}_k = 1$ ) for the class  $C_k$ , the competing models are with labels 1, and we find average of these, i.e. the equation (40) and

$$\psi_k(\mathbf{z}) = -g_k + \bar{y}_k \cdot u. \quad (43)$$

We denote Jacobian matrix for *units-vs-zeros* misclassification measure as  $J_{\mathbf{z}}\Psi$ . Further, we find the partial derivatives  $\frac{\partial \psi_k(\mathbf{z})}{\partial z_m}$  of the Jacobian  $J_{\mathbf{z}}\Psi$ , we have two cases:

a) if  $m = k$ , diagonal elements of the Jacobian

$$\begin{aligned} \frac{\partial \psi_k}{\partial z_k} &= -g'_k + \bar{y}_k \frac{\partial u}{\partial z_k} + y_k \frac{\partial v}{\partial z_k}, \\ \frac{\partial u}{\partial z_k} &= \frac{1}{\langle \mathbf{y}, \exp(\mathbf{g}) \rangle} \cdot y_k g'_k \exp(g_k), \\ \frac{\partial v}{\partial z_k} &= \frac{1}{\langle \bar{\mathbf{y}}, \exp(\mathbf{g}) \rangle} \cdot \bar{y}_k g'_k \exp(g_k). \end{aligned}$$

then we get

$$\frac{\partial \psi_k}{\partial z_k} = -g'_k, \quad (44)$$

because  $y_k \bar{y}_k = 0$ .

b) if  $m \neq k$ , off-diagonal elements

$$\begin{aligned} \frac{\partial \psi_k}{\partial z_m} &= \bar{y}_k \frac{\partial u}{\partial z_m} + y_k \frac{\partial v}{\partial z_m} = \frac{1}{\langle \mathbf{y}, \exp(\mathbf{g}) \rangle} \cdot \bar{y}_k y_m g'_m \exp(g_m) + \frac{1}{\langle \bar{\mathbf{y}}, \exp(\mathbf{g}) \rangle} y_k \bar{y}_m g'_m \exp(g_m) = \\ &= g'_m \exp(g_m) \left[ \frac{\bar{y}_k y_m}{\langle \mathbf{y}, \exp(\mathbf{g}) \rangle} + \frac{y_k \bar{y}_m}{\langle \bar{\mathbf{y}}, \exp(\mathbf{g}) \rangle} \right] \end{aligned}$$

Then, we get the Jacobian matrix

$$J_{\mathbf{z}}^{\top} \Psi = \begin{pmatrix} -g'_1 & \cdots & g'_1 \exp(g_1) \left[ \frac{\bar{y}_k y_1}{P} + \frac{y_k \bar{y}_1}{B} \right] & \cdots & g'_1 \exp(g_1) \left[ \frac{\bar{y}_M y_1}{P} + \frac{y_M \bar{y}_1}{B} \right] \\ \vdots & \ddots & \vdots & \vdots & \vdots \\ g'_m \exp(g_m) \left[ \frac{\bar{y}_1 y_m}{P} + \frac{y_1 \bar{y}_m}{B} \right] & \cdots & -g'_k & \cdots & g'_m \exp(g_m) \left[ \frac{\bar{y}_M y_m}{P} + \frac{y_M \bar{y}_m}{B} \right] \\ \vdots & \vdots & \vdots & \ddots & \vdots \\ g'_M \exp(g_M) \left[ \frac{\bar{y}_1 y_M}{P} + \frac{y_1 \bar{y}_M}{B} \right] & \cdots & g'_M \exp(g_M) \left[ \frac{\bar{y}_k y_M}{P} + \frac{y_k \bar{y}_M}{B} \right] & \cdots & -g'_M \end{pmatrix},$$

where  $P = \langle \mathbf{y}, \exp(\mathbf{g}) \rangle$  and  $B = \langle \bar{\mathbf{y}}, \exp(\mathbf{g}) \rangle$ . Extracting  $g'_k$ , we get

$$J_{\mathbf{z}}^{\top} \Psi = \text{diag}(g'_m) \cdot \begin{pmatrix} -1 & \cdots & \exp(g_1) \left[ \frac{\bar{y}_k y_1}{P} + \frac{y_k \bar{y}_1}{B} \right] & \cdots & \exp(g_1) \left[ \frac{\bar{y}_M y_1}{P} + \frac{y_M \bar{y}_1}{B} \right] \\ \vdots & \ddots & \vdots & \vdots & \vdots \\ \exp(g_m) \left[ \frac{\bar{y}_1 y_m}{P} + \frac{y_1 \bar{y}_m}{B} \right] & \cdots & -1 & \cdots & \exp(g_m) \left[ \frac{\bar{y}_M y_m}{P} + \frac{y_M \bar{y}_m}{B} \right] \\ \vdots & \vdots & \vdots & \ddots & \vdots \\ \exp(g_M) \left[ \frac{\bar{y}_1 y_M}{P} + \frac{y_1 \bar{y}_M}{B} \right] & \cdots & \exp(g_M) \left[ \frac{\bar{y}_k y_M}{P} + \frac{y_k \bar{y}_M}{B} \right] & \cdots & -1 \end{pmatrix} =$$

$$= \text{diag}(g'_m) \cdot \mathbf{G}$$

Thus, backpropagation network *error signal*, using *units-vs-zeros* misclassification measure, is

$$\delta^L(\mathbb{T}) = A \sum_{(\mathbf{x}, \mathbf{y}) \in \mathbb{T}} J_{\mathbf{z}}^{\top} E(\mathbf{x}) = A \sum_{(\mathbf{x}, \mathbf{y}) \in \mathbb{T}} J_{\mathbf{z}}^{\top} \Psi \cdot J_{\Psi}^{\top} \mathcal{L} \cdot J_{\mathcal{L}}^{\top} E, \quad (45)$$

where

$$\begin{aligned} J_{\mathbf{z}}^{\top} \Psi &= \text{diag}(g'_m) \cdot \mathbf{G} \\ J_{\Psi}^{\top} \mathcal{L} &= \text{diag}(l'_k) \\ J_{\mathcal{L}}^{\top} E &= (w_1 \mathbf{y} - w_2 \bar{\mathbf{y}}) \end{aligned} \quad (46)$$