

CSE463: Neural Networks

Neural Architectures for Multiclass Models

by:

Hossam Abd El Munim

Computer & Systems Engineering Dept.,

Ain Shams University,

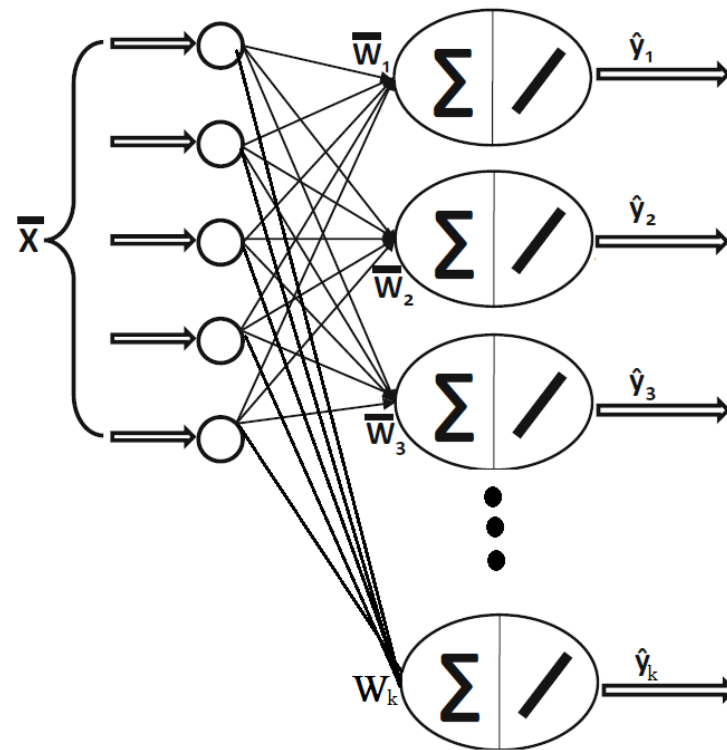
1 El-Sarayut Street, Abbassia, Cairo 11517

Outline

- Shallow neural network for multi-classification
- Problem definition.
- Class score
- Perceptron criterion loss function
- Multi-class SVM (MCSVM)
- Action recognition and MCSVM
- Hierarchical MCSVM
- Multinomial regression (Soft-max classification)
- Initialization of parameters
- Data normalization

Multiclass Problem Definition

Given a data set, $D = \{(\mathbf{X}_1, y_1), (\mathbf{X}_2, y_2) \dots (\mathbf{X}_N, y_N)\}$ of labelled feature vectors where $y \in \{1, 2, \dots, k\}$. Here, k represents the number of categories. We to estimate the vectors $\mathbf{W}_1, \mathbf{W}_2, \dots, \mathbf{W}_k$ that minimize a certain objective/loss function. What is it?



Score and Classification Decision

For a feature vector input \mathbf{X}_i , it must have the maximum score of class y_i for a correct classification and zero loss.

Class#	1	2	...	y_i	...	k
Score	$\mathbf{W}_1 \cdot \mathbf{X}_i$	$\mathbf{W}_2 \cdot \mathbf{X}_i$...	$\mathbf{W}_{y_i} \cdot \mathbf{X}_i$...	$\mathbf{W}_k \cdot \mathbf{X}_i$

Multiclass Perceptron: Loss Function

Given a data set, $D = \{(\mathbf{X}_1, y_1), (\mathbf{X}_2, y_2) \dots (\mathbf{X}_N, y_N)\}$ of labelled feature vectors where $y \in \{1, 2, \dots, k\}$. Here, k represents the number of categories. We to estimate the vectors $\mathbf{W}_1, \mathbf{W}_2, \dots, \mathbf{W}_k$ that minimize a certain objective/loss function. What is it?

$$L = \sum_{i=1}^N L_i$$

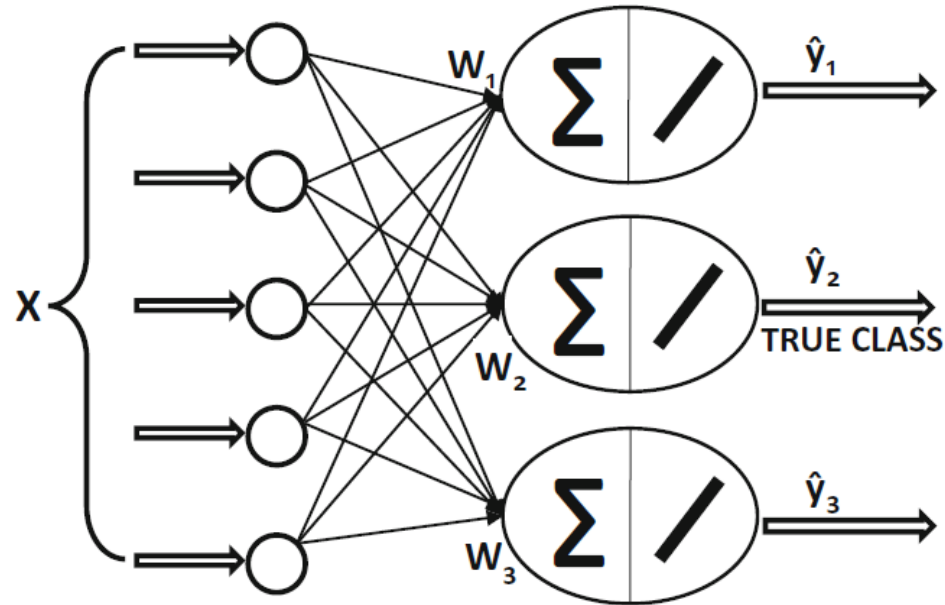
$$d_r = \max(\mathbf{W}_r \cdot \mathbf{X}_i - \mathbf{W}_{y_i} \cdot \mathbf{X}_i, 0)$$

$$r \in \{\{1, 2, \dots, k\} - \{y_i\}\}$$

$$L_i = \max(\dots, d_r, \dots)$$

Multiclass Perceptron: Loss Function Example

Assume $k = 3$, $y_i = 2$, then L_i will be:-



$$d_1 = \max(\mathbf{w}_1 \cdot \mathbf{x}_i - \mathbf{w}_2 \cdot \mathbf{x}_i, 0)$$

$$d_3 = \max(\mathbf{w}_3 \cdot \mathbf{x}_i - \mathbf{w}_2 \cdot \mathbf{x}_i, 0)$$

$$r \in \{1, 3\}$$

$$L_i = \max(d_1, d_3)$$

Multiclass Perceptron: Learning Rule

- For correctly classified features, there will be no update as the loss is zero. Its derivative will be zero as well.
- For miss-classified features, there will be a derivative computed as follows:-

$$d_1 = \max(\mathbf{W}_1 \cdot \mathbf{X}_i - \mathbf{W}_2 \cdot \mathbf{X}_i, 0)$$

$$d_3 = \max(\mathbf{W}_3 \cdot \mathbf{X}_i - \mathbf{W}_2 \cdot \mathbf{X}_i, 0)$$

$$r \in \{1, 3\}$$

$$L_i = \max(d_1, d_3)$$

$$\frac{\partial}{\partial \mathbf{W}_r} L_i = \begin{cases} -\mathbf{X}_i & \text{if } r = y_i \\ \mathbf{X}_i & \text{if } r \neq y_i \text{ is the most miss-classified prediction.} \\ 0 & \text{otherwise} \end{cases}$$

Multiclass Perceptron: Learning Rule

Disadvantages

- The multiclass perceptron only updates the linear separator of a class that is predicted most incorrectly along with the linear separator of the true class.

Weston-Watkins SVM

- On the other hand, the Weston-Watkins SVM updates the separator of any class that is predicted more favorably than the true class. In both cases, the separator of the observed class is updated by the same aggregate amount as the incorrect classes (but in the opposite direction).
- Not only does the Weston-Watkins SVM update the separator in the case of misclassification, it updates the separators in cases where an incorrect class gets a prediction that is “uncomfortably close” to the true class. This is based on the notion of margin.

Weston-Watkins SVM Loss Function

$$L = \sum_{i=1}^N L_i$$

$$d_r = \max(1 + \mathbf{W}_r \cdot \mathbf{X}_i - \mathbf{W}_{y_i} \cdot \mathbf{X}_i, 0)$$

$$r \in \{\{1, 2, \dots, k\} - \{y_i\}\}$$

$$L_i = \sum_{r=1, r \neq y_i}^k d_r$$

Weston-Watkins SVM Learning Rule

$$\frac{\partial}{\partial \mathbf{W}_r} L_i = \begin{cases} -\mathbf{X}_i \left(\sum_{j=1, j \neq r}^k \delta_{ji} \right) & \text{if } r = y_i \\ \mathbf{X}_i (\delta_{ri}) & \text{if } r \neq y_i \end{cases}$$

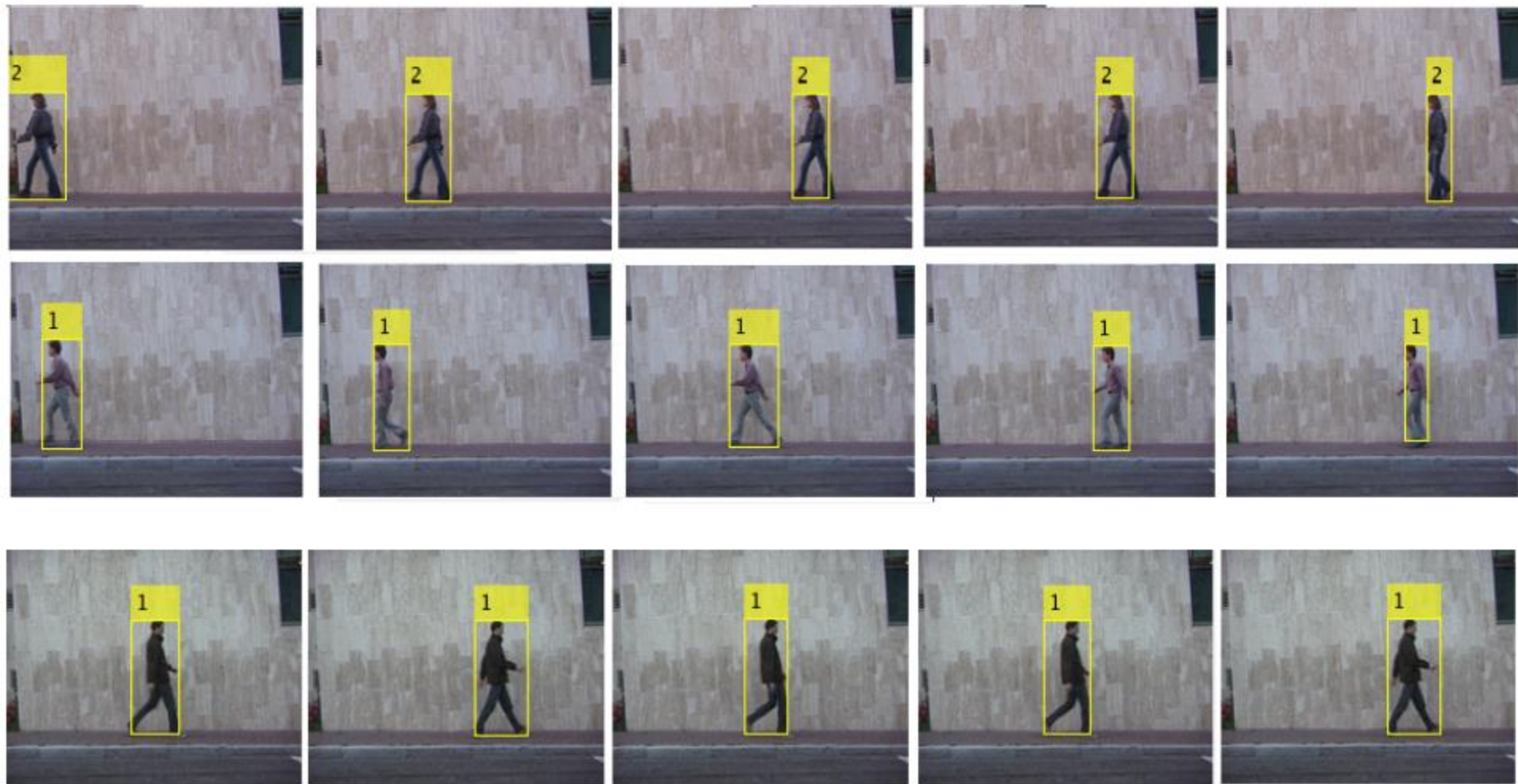
$$\delta_{ji} = \begin{cases} 1 & \text{if } 1 + \mathbf{W}_j \cdot \mathbf{X}_i - \mathbf{W}_{y_i} \cdot \mathbf{X}_i > 0 \\ 0 & \text{Otherwise} \end{cases}$$

Weston-Watkins SVM Learning Rule with Regularization

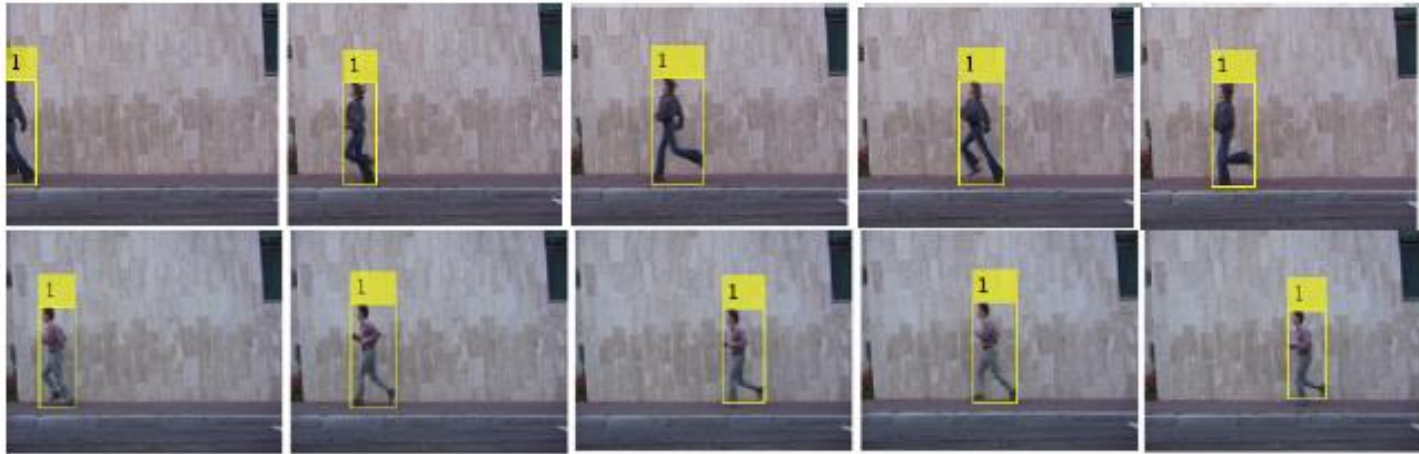
$$L = \sum_{i=1}^N \sum_{r=1, r \neq y_i}^k d_r + (\lambda/2) \sum_{r=1}^k \mathbf{w}_r \cdot \mathbf{w}_r$$

Regularization is considered
essential to the
proper functioning of a support
vector machine.

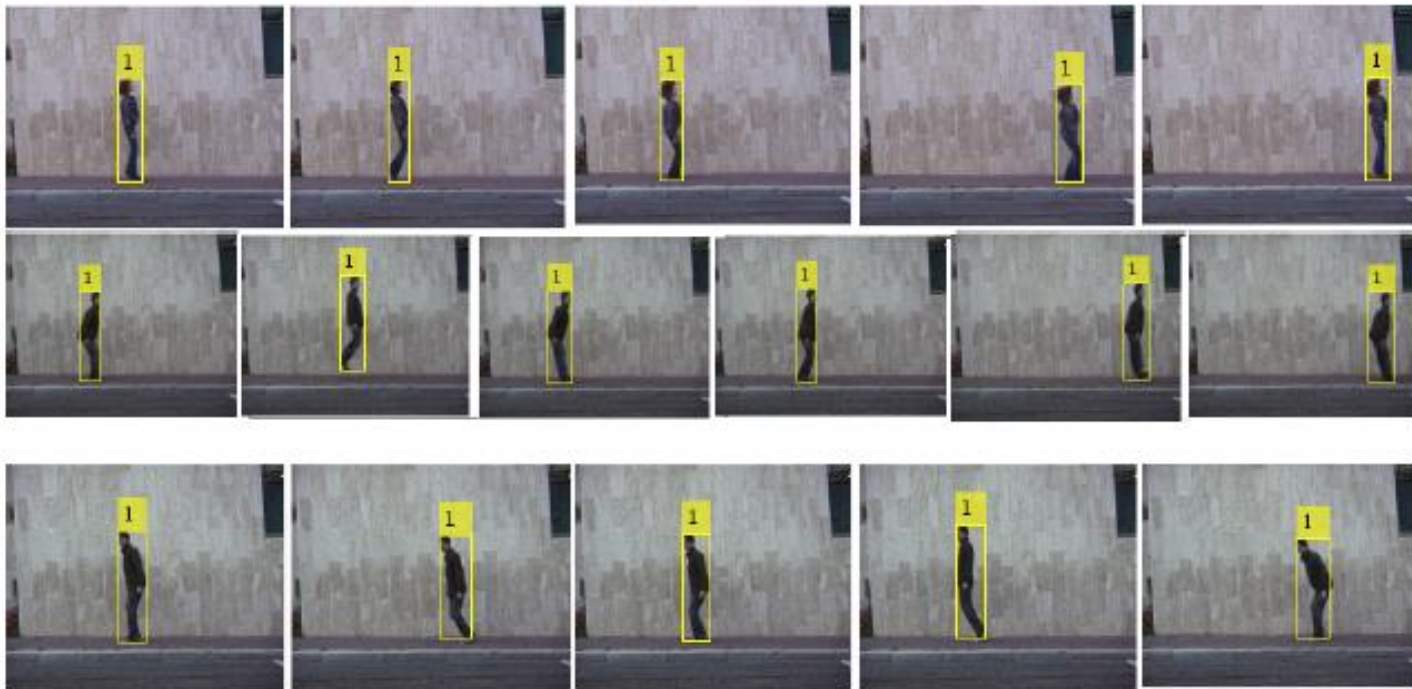
Action Recognition (I)



Action Recognition (II)

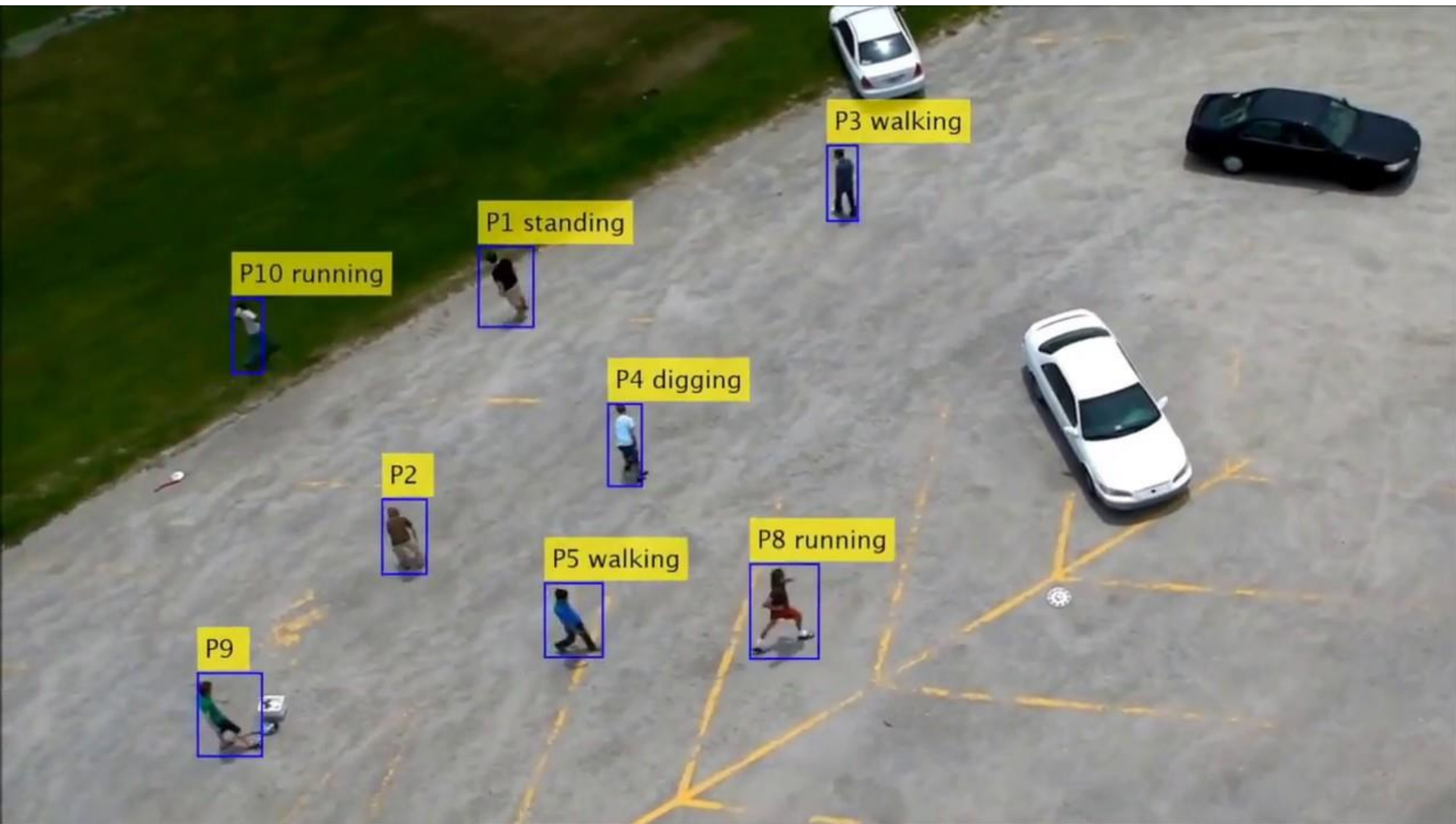


Samples of Running activity for tracked persons



Samples of Jumping activity for tracked persons

Action Recognition (III)



Action Recognition: Hierarchical MCSVM



(a) Bend

(b) Jack

(c) Jump

(d) Pjump

(e) Run



(f) Side

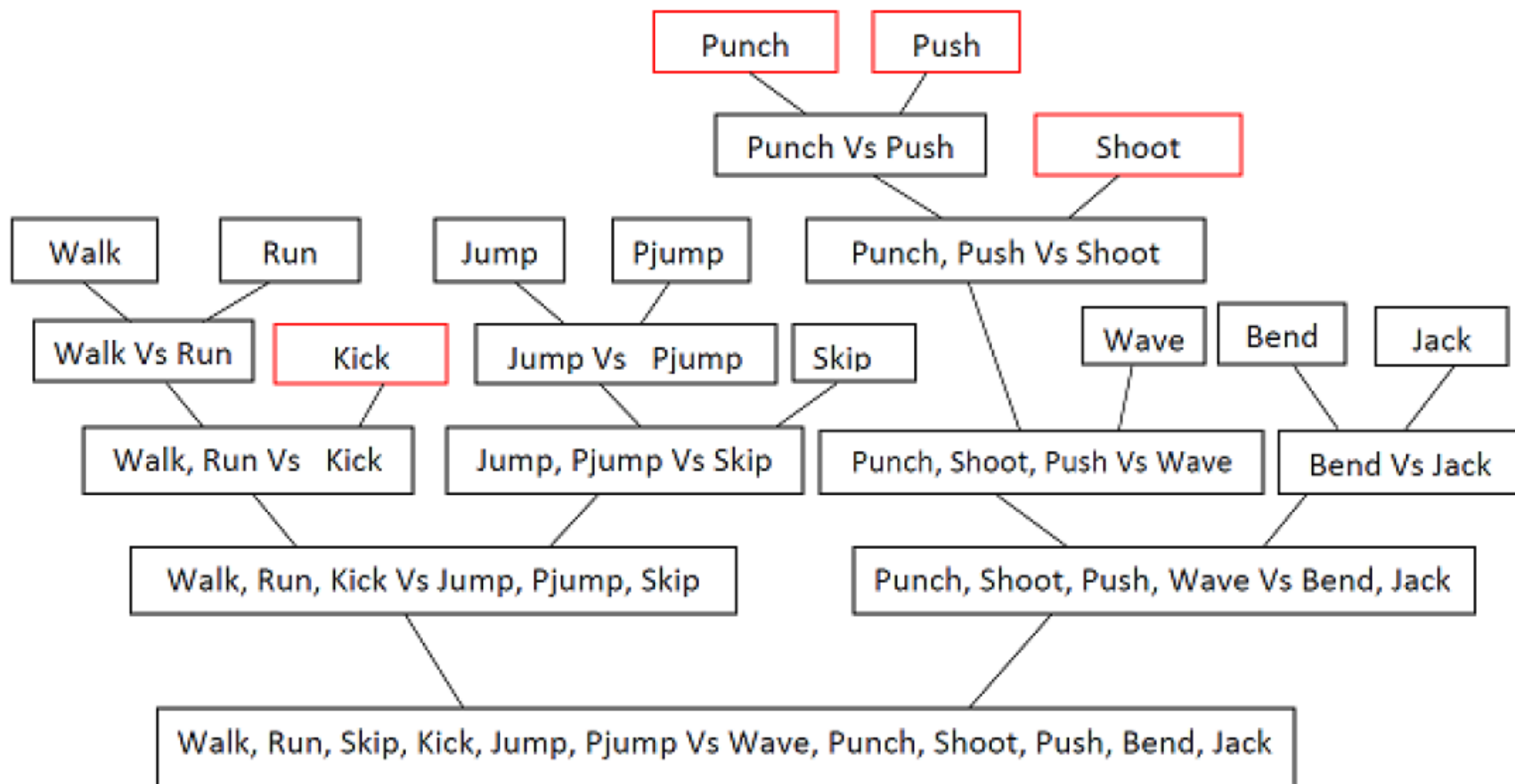
(g) Skip

(h) Walk

(i) Wave1

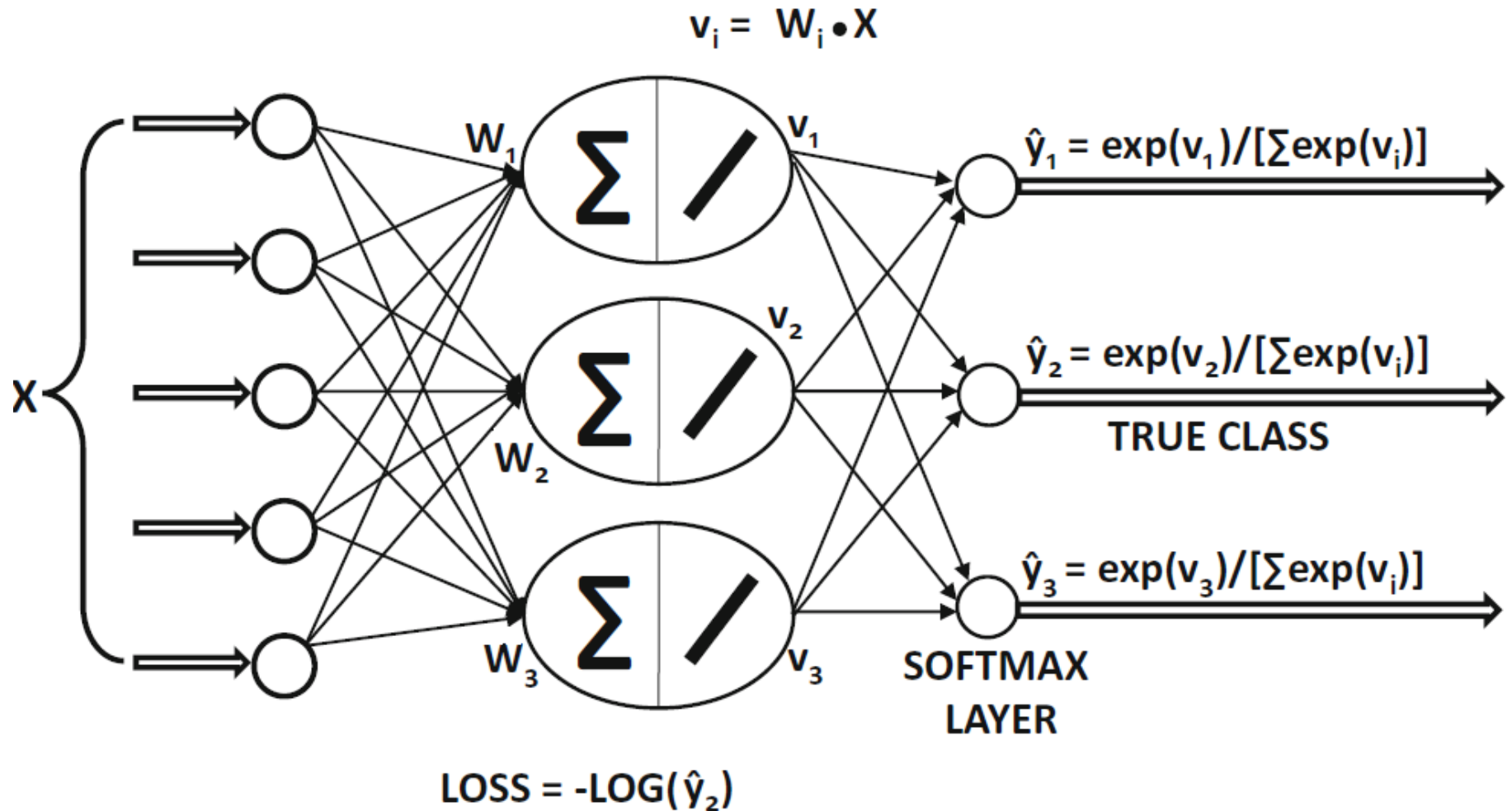
(j) Wave2

Action Recognition: Hierarchical MCSVM



Inverted binary tree hierarchical multiclass SVM.

Multinomial Logistic Regression (Softmax Classifier)



Multinomial Logistic Regression (Loss Function)

$$v_r = \mathbf{W}_r \cdot \mathbf{X}_i$$

$$\hat{y}_r = \frac{\exp(v_r)}{\sum_{j=1}^k \exp(v_j)}$$

$$L_i = -\log(\hat{y}_{y_i})$$

$$L_i = -v_{y_i} + \log\left(\sum_{j=1}^k \exp(v_j)\right)$$

Multinomial Logistic Regression (Learning Rule)

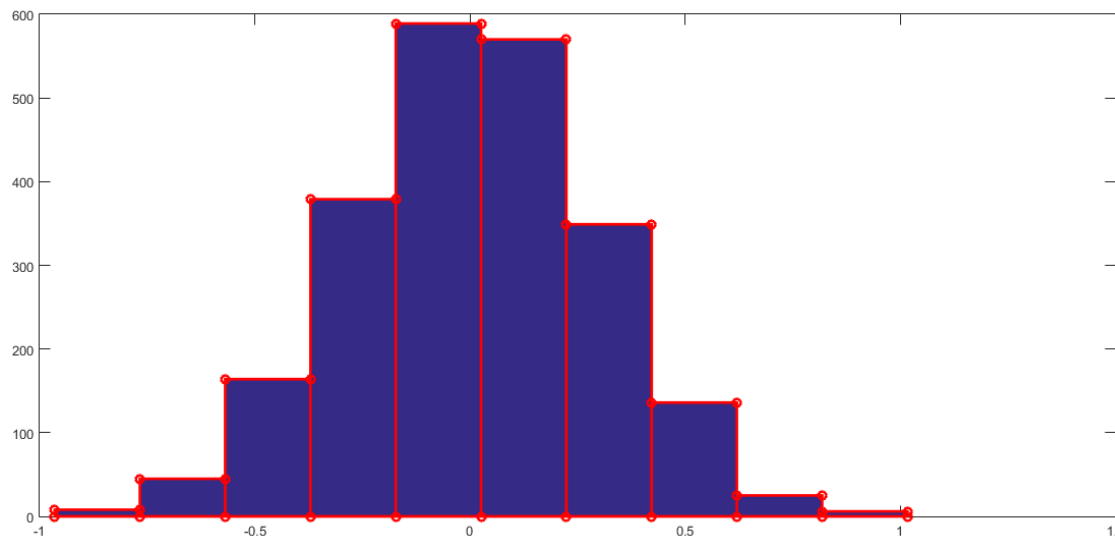
$$\frac{\partial L_i}{\partial \mathbf{W}_r} = (\partial L_i / \partial v_r)(\partial v_r / \partial \mathbf{W}_r)$$

$$\partial L_i / \partial v_r = \begin{cases} -(1 - \widehat{y}_{y_i}) & \text{if } y_i = r \\ \widehat{y}_{y_i} & \text{if } y_i \neq r \end{cases}$$

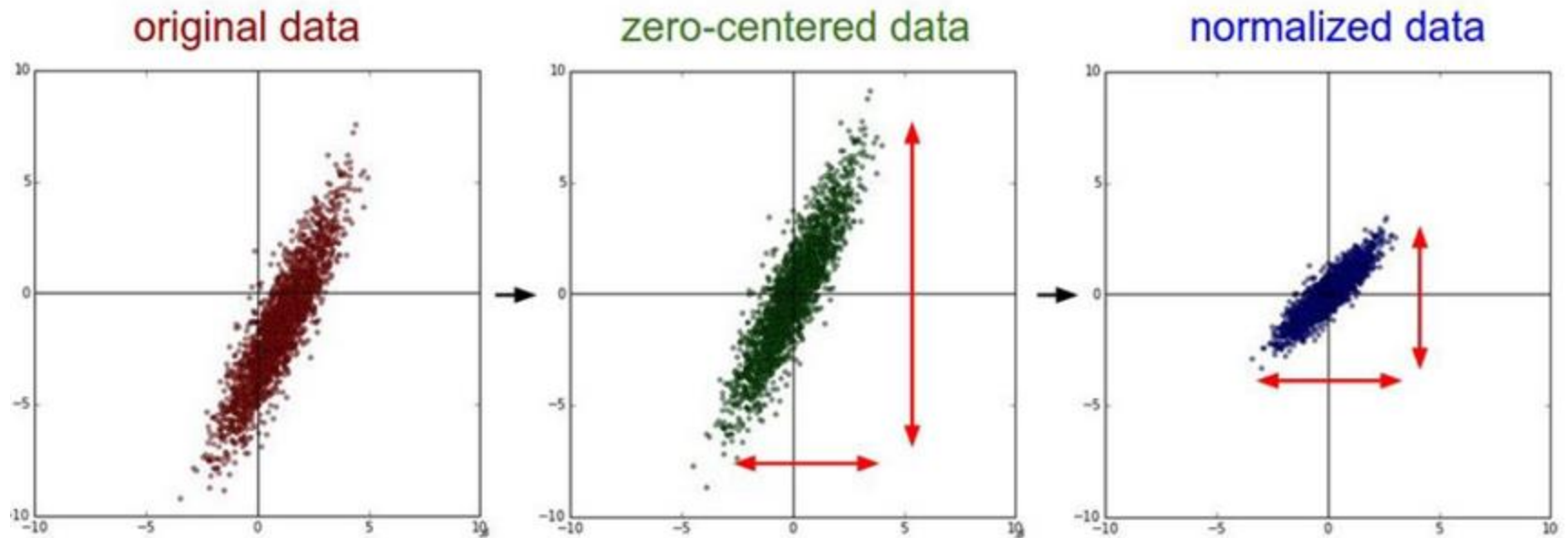
$$\frac{\partial v_r}{\partial \mathbf{W}_r} = \mathbf{X}_i$$

A Hint about Initialization

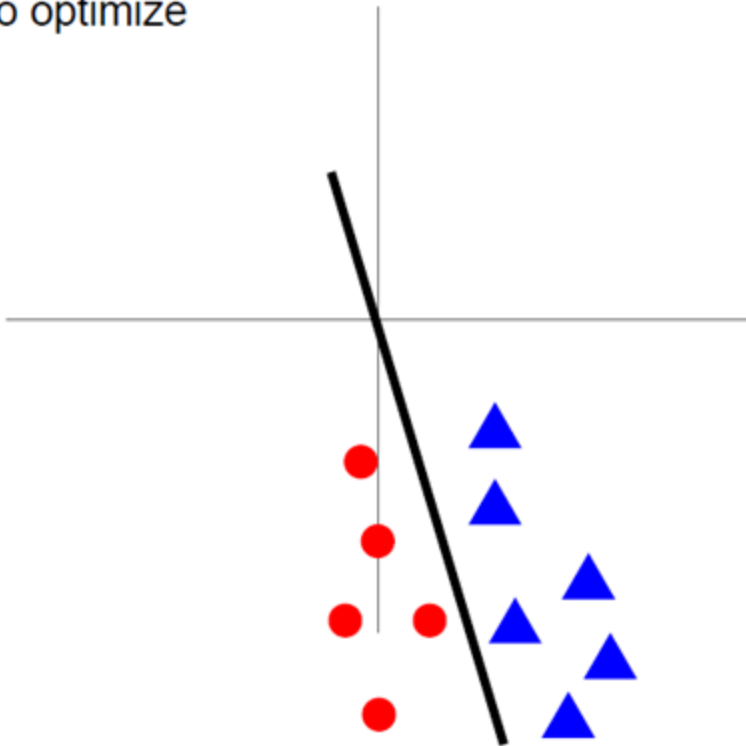
- If weights are very close to zero, you will get zeros functions and hence zero derivatives. There will be no learning.
- Very large values are not recommended to avoid numerical instability. If you are using tanh/sigmoid functions, you get saturated outputs with zero derivatives as well.
- Initialize with small values around zero, to get a nice distribution. Optimization/Learning will go smoothly.



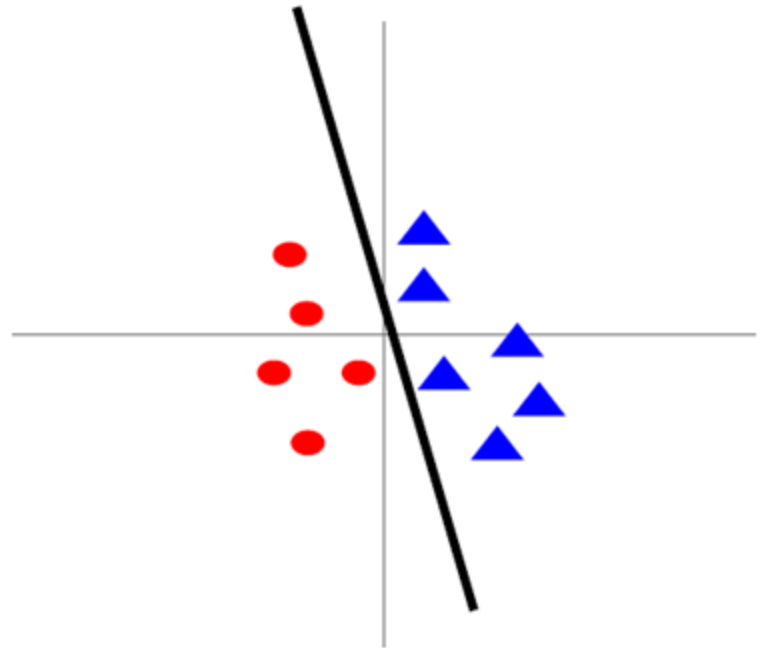
Training Data Preprocessing



Before normalization: classification loss very sensitive to changes in weight matrix; hard to optimize



After normalization: less sensitive to small changes in weights; easier to optimize



Batch Normalization

Input: $x : N \times D$

Intermediates: $\mu, \sigma : D$
 $\hat{x} : N \times D$

Output: $y : N \times D$

$$\mu_j = \frac{1}{N} \sum_{i=1}^N x_{i,j}$$

$$\sigma_j^2 = \frac{1}{N} \sum_{i=1}^N (x_{i,j} - \mu_j)^2$$

$$\hat{x}_{i,j} = \frac{x_{i,j} - \mu_j}{\sqrt{\sigma_j^2 + \varepsilon}}$$

$$y_{i,j} = \gamma_j \hat{x}_{i,j} + \beta_j$$