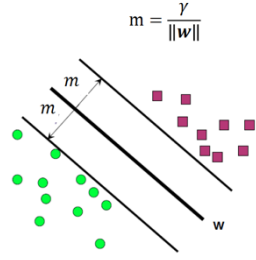


$$\mathbf{w} = \sum_{i=1}^N \alpha_i y_i \mathbf{x}_i$$



Define **hinge loss** or **soft margin**:

$$L_i(\mathbf{w}, b) = \max(0, 1 - y_i z_i)$$

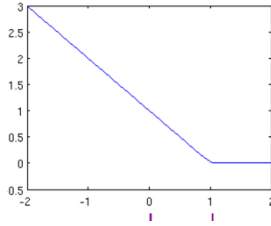
Define **Lagrangian**: $L(\mathbf{w}, \lambda) = f(\mathbf{w}) + \lambda g(\mathbf{w})$

\mathbf{w} is called the **primal** variable

λ is called the **dual** variable

KKT Conditions: $\hat{\mathbf{w}}, \hat{\lambda}$ satisfy:

- $\hat{\mathbf{w}}$ minimizes the Lagrangian: $\hat{\mathbf{w}} = \arg \min_{\mathbf{w}} L(\mathbf{w}, \hat{\lambda})$
- Either
 - $g(\hat{\mathbf{w}}) = 0$ and $\hat{\lambda} \geq 0$ [active constraint]
 - $g(\hat{\mathbf{w}}) < 0$ and $\hat{\lambda} = 0$ [inactive constraint]



- $y_i z_i \geq 1 \Rightarrow$ Sample meets margin target,
- $y_i z_i \in [0, 1] \Rightarrow$ Sample margin too small
- $y_i z_i \leq 0 \Rightarrow$ Sample misclassified

$$J(\mathbf{w}, b) = C \sum_{i=1}^N \max(0, 1 - y_i(\mathbf{w}^T \phi(\mathbf{x}_i) + b)) + \frac{1}{2} \|\mathbf{w}\|^2$$

$$\mathbf{w} = \sum_{i=1}^N \alpha_i y_i \phi(\mathbf{x}_i)$$

$$z = b + \mathbf{w}^T \mathbf{x} = b + \sum_{i=1}^N \alpha_i y_i \phi(\mathbf{x}_i)^T \phi(\mathbf{x})$$

$K(\mathbf{x}_i, \mathbf{x}) = \text{"kernel"}$

Perfectly linearly separable if there exists a $\theta = (b, w_1, \dots, w_d)$

- $b + w_1 x_{i1} + \dots + w_d x_{id} > \gamma$ when $y_i = 1$
- $b + w_1 x_{i1} + \dots + w_d x_{id} < -\gamma$ when $y_i = -1$
- $y_i(b + w_1 x_{i1} + \dots + w_d x_{id}) > \gamma$

Support vectors: Samples that either:

- Are exactly on margin: $y_i(\mathbf{w}^T \mathbf{x}_i + b) = 1$
- Or, on wrong side of margin: $y_i(\mathbf{w}^T \mathbf{x}_i + b) \leq 1$

$$J(\mathbf{w}, b) = C \sum_{i=1}^N \max(0, 1 - y_i(\mathbf{w}^T \mathbf{x}_i + b)) + \frac{1}{2} \|\mathbf{w}\|^2$$

Parameter ϵ_i called the **slack variable**

- $\epsilon_i = 0 \Rightarrow$ Sample on correct side of margin
- $\epsilon_i \geq 0 \Rightarrow$ Sample violates the margin
- $\epsilon_i \geq 1 \Rightarrow$ Sample misclassified (wrong side of hyperplane)

Parameter C :

- C large: Forces minimum number of violations. Highly fit to data. Low bias, higher variance
- C small: Enables more samples violations. Higher bias, lower variance

$$\text{Hidden layer: } z_j^H = \sum_{k=1}^{N_i} W_{jk}^H x_k + b_j^H, \quad u_j^H = g_{\text{act}}(z_j^H), \quad j = 1, \dots, N_h$$

Gradient tensors: Suppose that $\mathbf{Y} = f(\mathbf{X})$ where:

$$\text{Output layer: } z_j^O = \sum_{k=1}^{N_h} W_{jk}^O u_k^H + b_j^O, \quad u^O = g_{\text{out}}(\mathbf{z}^O), \quad j = 1, \dots, N_o.$$

– The input \mathbf{X} is a tensor of size (N_1, \dots, N_r) ,

The output \mathbf{Y} is a tensor of size (M_1, \dots, M_s) .

In the hidden layer, the function $g_{\text{act}}(z)$ is called the **activation function**.

Hard threshold:

$$g_{\text{act}}(z) = \begin{cases} 1, & \text{if } z \geq 0 \\ 0, & \text{if } z < 0. \end{cases}$$

Sigmoid: $g_{\text{act}}(z) = 1/(1 + e^{-z})$.

Rectified linear unit (ReLU): $g(z) = \max\{0, z\}$.

output map $g_{\text{out}}(z)$

$$P(y = 1|\mathbf{x}) = u^O = g_{\text{out}}(z^O) = \frac{1}{1 + e^{-z^O}}.$$

$$P(y = k|\mathbf{x}) = u_k^O = g_{\text{out},k}(z^O) = \frac{e^{z_k^O}}{\sum_{\ell=1}^K e^{z_\ell^O}}.$$

$$\hat{\mathbf{y}} = \mathbf{u}^O = g_{\text{out}}(\mathbf{z}^O) = \mathbf{z}^O.$$

The input \mathbf{X} is a tensor of size (K_1, \dots, K_t) ;

The intermediate variable $\mathbf{Y} = g(\mathbf{X})$ is a tensor of size (N_1, \dots, N_r) ;

The output $\mathbf{Z} = h(\mathbf{Y})$ is a tensor of size (M_1, \dots, M_s) .

$$\frac{\partial Z_{i_1, \dots, i_r}}{\partial X_{k_1, \dots, k_t}} = \sum_{j_1=0}^{N_1-1} \dots \sum_{j_r=0}^{N_r-1} \frac{\partial Z_{i_1, \dots, i_r}}{\partial Y_{j_1, \dots, j_s}} \frac{\partial Y_{j_1, \dots, j_s}}{\partial X_{k_1, \dots, k_t}}.$$

$$\mathbf{u} = f(\mathbf{z}) = (f(z_1), \dots, f(z_N)), \quad f(z_i) = \frac{1}{1 + e^{-z_i}},$$

$$\frac{\partial \mathbf{u}}{\partial \mathbf{z}} = \begin{bmatrix} f'(z_1) & 0 & \dots & 0 \\ 0 & f'(z_2) & \dots & 0 \\ \vdots & \vdots & \ddots & \vdots \\ 0 & 0 & \dots & f'(z_N) \end{bmatrix}$$

$$\frac{\partial u_i}{\partial W_{ij}} = \frac{\partial u_i}{\partial z_i} \frac{\partial z_i}{\partial W_{ij}} = f'(z_i) x_j,$$

$$\frac{\partial u_i}{\partial b_i} = \frac{\partial u_i}{\partial z_i} \frac{\partial z_i}{\partial b_i} = f'(z_i).$$

$$\frac{\partial u_i}{\partial x_j} = \frac{\partial u_i}{\partial z_i} \frac{\partial z_i}{\partial x_j} = f'(z_i) W_{ij}.$$

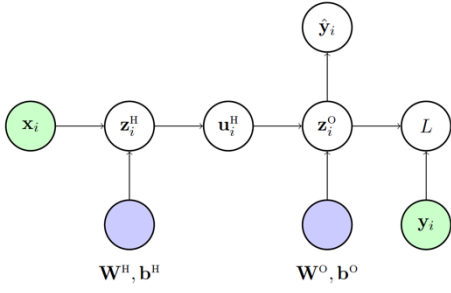
The loss function factor $g_{\text{loss}}(z_i^O, y_i)$

$$g_{\text{loss}}(z_i^O, y_i) = \ln [1 + e^{z_i^O}] - y_i z_i.$$

$$g_{\text{loss}}(\mathbf{z}_i^O, y_i) := \ln \left[\sum_{\ell=1}^K e^{z_{i\ell}^O} \right] - \sum_{k=1}^K r_{ik} z_{ik}^O,$$

$$r_{ik} = \begin{cases} 1 & \text{if } y_i = k, \\ 0 & \text{if } y_i \neq k. \end{cases}$$

$$g_{\text{loss}}(\mathbf{z}_i^O, \mathbf{y}_i) = \|\mathbf{z}_i^O - \mathbf{y}_i\|^2 = \sum_{k=1}^d (z_{ik}^O - y_{ik})^2.$$



$$\mathbf{z}_i^H = \mathbf{W}^H \mathbf{x}_i + \mathbf{b}^H, \quad \mathbf{u}_i^H = g_{\text{act}}(\mathbf{z}_i^H),$$

$$\mathbf{z}_i^O = \mathbf{W}^O \mathbf{u}_i^H + \mathbf{b}^O, \quad \hat{y}_i = g_{\text{out}}(\mathbf{z}_i^O).$$

$$L = \sum_{i=1}^N g_{\text{loss}}(\mathbf{z}_i^O, y_i).$$

$$u_{ij}^H = g_{\text{act}}(z_{ij}^H) = \max(0, z_{ij}^H).$$

$$\frac{\partial u_{ij}^H}{\partial z_{ij}^H} = \begin{cases} 1, & \text{if } z_{ij}^H > 0 \\ 0, & \text{if } z_{ij}^H < 0. \end{cases}$$

$$\frac{\partial L}{\partial z_{ij}^H} = \frac{\partial L}{\partial u_{ij}^H} \frac{\partial u_{ij}^H}{\partial z_{ij}^H}.$$

$$z_{ij}^H = \sum_k W_{jk}^H x_{ik} + b_j^H.$$

$$\frac{\partial z_{ij}^H}{\partial W_{jk}^H} = x_{ik}, \quad \frac{\partial z_{ij}^H}{\partial b_j^H} = 1.$$

$$\frac{\partial L}{\partial W_{jk}^H} = \sum_{i=1}^N \frac{\partial L}{\partial z_{ik}^H} \frac{\partial z_{ij}^H}{\partial W_{jk}^H} = \sum_{i=1}^N \frac{\partial L}{\partial z_{ik}^H} x_{ik},$$

$$\frac{\partial L}{\partial b_j^H} = \sum_{i=1}^N \frac{\partial L}{\partial z_{ik}^H} \frac{\partial z_{ij}^H}{\partial b_j^H} = \sum_{i=1}^N \frac{\partial L}{\partial z_{ik}^H}.$$

$$\frac{\partial L}{\partial z_{ij}^O} = \frac{\partial g_{\text{loss}}(\mathbf{z}_i^O, y_i)}{\partial z_{ij}^O} = \frac{e^{z_{ij}^O}}{\sum_{\ell=1}^K e^{z_{i\ell}^O}} - r_{ij}.$$

$$z_{ij}^O = \sum_k W_{jk}^O u_{ik}^H + b_j^O.$$

$$\frac{\partial z_{ij}^O}{\partial W_{jk}^O} = u_{ik}^H, \quad \frac{\partial z_{ij}^O}{\partial b_j^O} = 1, \quad \frac{\partial z_{ij}^O}{\partial u_{ik}^H} = W_{jk}.$$

$$\frac{\partial L}{\partial W_{jk}^O} = \sum_{i=1}^N \frac{\partial L}{\partial z_{ik}^O} \frac{\partial z_{ij}^O}{\partial W_{jk}^O} = \sum_{i=1}^N \frac{\partial L}{\partial z_{ik}^O} u_{ik}^H,$$

$$\frac{\partial L}{\partial b_j^O} = \sum_{i=1}^N \frac{\partial L}{\partial z_{ik}^O} \frac{\partial z_{ij}^O}{\partial b_j^O} = \sum_{i=1}^N \frac{\partial L}{\partial z_{ik}^O},$$

$$\frac{\partial L}{\partial u_{ik}^H} = \sum_{j=1}^{N_o} \frac{\partial L}{\partial z_{ij}^O} \frac{\partial z_{ij}^O}{\partial u_{ik}^H} \sum_{j=1}^{N_o} \frac{\partial L}{\partial z_{ij}^O} W_{jk}.$$

□ Suppose inputs are

- x , size $N_1 \times N_2$, w : size $K_1 \times K_2$, $K_1 \leq N_1$, $K_2 \leq N_2$
- $z = x * w$ (without reversal)

$$z[n_1, n_2] = \sum_{k_2=0}^{K_2-1} \sum_{k_1=0}^{K_1-1} w[k_1, k_2] x[n_1 + k_1, n_2 + k_2]$$

□ Valid mode: $0 \leq n_1 < N_1 - K_1 + 1$, $0 \leq n_2 < N_2 - K_2 + 1$

- Requires no zero padding

□ Same mode: Output size $N_1 \times N_2$

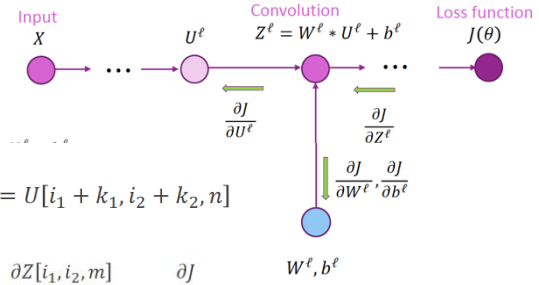
- Usually use zero padding for neural networks

□ Weight and bias:

- W : Weight tensor, size $(K_1, K_2, N_{in}, N_{out})$
- b : Bias vector, size N_{out}

$$Z[i_1, i_2, m] = \sum_{k_1=0}^{K_1-1} \sum_{k_2=0}^{K_2-1} \sum_{n=0}^{N_{in}-1} W[k_1, k_2, n, m] X[i_1 + k_1, i_2 + k_2, n] + b[m]$$

$$Z^\ell = W^\ell * U^\ell + b^\ell$$



$$\frac{\partial Z[i_1, i_2, m]}{\partial W[k_1, k_2, n, m]} = U[i_1 + k_1, i_2 + k_2, n]$$

$$\frac{\partial J}{\partial W[k_1, k_2, n, m]} = \sum_{i_1=1}^{N_1} \sum_{i_2=1}^{N_2} \frac{\partial Z[i_1, i_2, m]}{\partial W[k_1, k_2, n, m]} \frac{\partial J}{\partial Z[i_1, i_2, m]}$$

$$= \sum_{i_1=1}^{N_1} \sum_{i_2=1}^{N_2} U[i_1 + k_1, i_2 + k_2, n] \frac{\partial J}{\partial Z[i_1, i_2, m]}$$

$$W^\ell, b^\ell$$