

CS480/680, Spring 2025

# Review Notes

Student: Hongxu Xu (h445xu@uwaterloo.ca)

August 5, 2025

# Part I

## For Mid-Term

### 1 Perceptron

**Question 1** (Linear Function).

$$\begin{aligned} & \forall \alpha, \beta \in \mathbb{R}, \forall \mathbf{x}, \mathbf{z} \in \mathbb{R}^d, f(\alpha \mathbf{x} + \beta \mathbf{z}) = \alpha \cdot f(\mathbf{x}) + \beta \cdot f(\mathbf{z}) \\ & \iff \exists \mathbf{w} \in \mathbb{R}^d, f(\mathbf{x}) = \langle \mathbf{x}, \mathbf{w} \rangle \end{aligned} \quad (1)$$

**Proof  $\Rightarrow$ :** Let  $\mathbf{w} := [f(\mathbf{e}_1), \dots, f(\mathbf{e}_d)]^T$ , where  $\mathbf{e}_i$  is the  $i$ -th coordinate vector.

$$\begin{aligned} f(\mathbf{x}) &= f(x_1 \mathbf{e}_1 + \dots + x_d \mathbf{e}_d) \\ &= x_1 f(\mathbf{e}_1) + \dots + x_d f(\mathbf{e}_d) \\ &= \langle \mathbf{x}, \mathbf{w} \rangle \end{aligned} \quad (2)$$

**Proof  $\Leftarrow$ :**

$$\begin{aligned} f(\alpha \mathbf{x} + \beta \mathbf{z}) &= \langle \alpha \mathbf{x} + \beta \mathbf{z}, \mathbf{w} \rangle \\ &= \langle \alpha \mathbf{x}, \mathbf{w} \rangle + \langle \beta \mathbf{z}, \mathbf{w} \rangle \\ &= \alpha \langle \mathbf{x}, \mathbf{w} \rangle + \beta \langle \mathbf{z}, \mathbf{w} \rangle \\ &= \alpha f(\mathbf{x}) + \beta f(\mathbf{z}) \end{aligned} \quad (3)$$

**Question 2** ( $\mathbf{w}$  is Orthogonal to Decision Boundary  $H$ ). Any vector on  $H$  can be written as  $\overrightarrow{\mathbf{x}\mathbf{x}'} = \mathbf{x}' - \mathbf{x}$ ,

$$\langle \mathbf{x}' - \mathbf{x}, \mathbf{w} \rangle = \langle \mathbf{x}', \mathbf{w} \rangle - \langle \mathbf{x}, \mathbf{w} \rangle = -b - (-b) = 0 \quad (4)$$

**Question 3** (Update Rule for Perceptron).

$$\begin{aligned} \mathbf{w} &\leftarrow \mathbf{w} + y_i \mathbf{x}_i \\ b &\leftarrow b + y_i \end{aligned} \quad (5)$$

**Question 4** (Feasibility of Perceptron). The goal is to find  $\mathbf{w} \in \mathbb{R}^d, b \in \mathbb{R}$ , such that  $\forall i, y_i (\langle \mathbf{x}_i, \mathbf{w} \rangle + b) > 0$ . According to the update rule 5,

$$\begin{aligned} y [\langle \mathbf{x}, \mathbf{w}_{new} \rangle + b_{new}] &= y [\langle \mathbf{x}, \mathbf{w}_{old} + y \mathbf{x} \rangle + b_{old} + y] \\ &= y [\langle \mathbf{x}, \mathbf{w}_{old} \rangle + b_{old}] + y [\langle \mathbf{x}, y \mathbf{x} \rangle + y] \\ &= y [\langle \mathbf{x}, \mathbf{w}_{old} \rangle + b_{old}] + y [y \|\mathbf{x}\|_2^2 + y] \\ &= y [\langle \mathbf{x}, \mathbf{w}_{old} \rangle + b_{old}] + y^2 \|\mathbf{x}\|_2^2 + y^2 \\ &= y [\langle \mathbf{x}, \mathbf{w}_{old} \rangle + b_{old}] + \underbrace{\|\mathbf{x}\|_2^2 + 1}_{\text{always positive}} \end{aligned} \quad (6)$$

Notice that  $y \in \{\pm 1\} \Rightarrow y^2 = 1$ .

$\|\mathbf{x}\|_2^2 + 1$  is always positive, which means we always increase the confidence  $y\hat{y}$  after the update.

**Question 5** (Trick for Hiding the Bias Term – Padding).

$$\begin{aligned} \langle \mathbf{x}, \mathbf{w} \rangle + b &= \left\langle \begin{pmatrix} \mathbf{x} \\ 1 \end{pmatrix}, \begin{pmatrix} \mathbf{w} \\ b \end{pmatrix} \right\rangle \\ &= \langle \mathbf{x}_{pad}, \mathbf{w}_{pad} \rangle \end{aligned} \quad (7)$$

Correspondingly, the update rule can be written as:

$$\mathbf{w}_{pad} \leftarrow \mathbf{w}_{pad} + y \mathbf{x}_{pad} \quad (8)$$

**Question 6** (Margin). Suppose  $\exists \mathbf{w}^*$  such that  $\forall i, y_i \langle \mathbf{x}_i, \mathbf{w}^* \rangle > 0$ .

We normalize  $\mathbf{w}^*$  such that  $\|\mathbf{w}^*\|_2 = 1$ .

In other words,  $\mathbf{w}^*$  is the normalized weight for the decision boundary.

$$\text{Margin } \gamma := \min_i |\langle \mathbf{x}_i, \mathbf{w}^* \rangle| \quad (9)$$

**Question 7** (Convergence Theorem – Linearly Separable Case). Assume that  $\forall i, \|\mathbf{x}_i\|_2 \leq C$  (i.e. within a circle with radius  $C$ ). Then the Perceptron algorithm converges after  $\frac{C^2}{\gamma^2}$  mistakes.

**Proof:**

Suppose  $\mathbf{w}$  is the updating weight, and  $\theta$  is the angle between  $\mathbf{w}$  and  $\mathbf{w}^*$ .

We have  $\langle \mathbf{w}, \mathbf{w}^* \rangle = \|\mathbf{w}\|_2 \|\mathbf{w}^*\|_2 \cos \theta = \|\mathbf{w}\|_2 \cos \theta$ .

After an update,  $\|\mathbf{w}_{new}\|_2 \cos \theta_{new}$  will be

$$\begin{aligned} \langle \mathbf{w} + y\mathbf{x}, \mathbf{w}^* \rangle &= \langle \mathbf{w}, \mathbf{w}^* \rangle + y \langle \mathbf{x}, \mathbf{w}^* \rangle \\ &= \langle \mathbf{w}, \mathbf{w}^* \rangle + |\langle \mathbf{x}, \mathbf{w}^* \rangle| \\ &\geq \langle \mathbf{w}, \mathbf{w}^* \rangle + \gamma \end{aligned} \quad (10)$$

Let's see the change of  $\langle \mathbf{w}_{new}, \mathbf{w}_{new} \rangle = \|\mathbf{w}_{new}\|_2^2$ ,

$$\begin{aligned} \langle \mathbf{w} + y\mathbf{x}, \mathbf{w} + y\mathbf{x} \rangle &= \langle \mathbf{w}, \mathbf{w} \rangle + 2y \langle \mathbf{w}, \mathbf{x} \rangle + y^2 \langle \mathbf{x}, \mathbf{x} \rangle \\ &= \|\mathbf{w}\|_2^2 + 2y \langle \mathbf{w}, \mathbf{x} \rangle + \|\mathbf{x}\|_2^2 \end{aligned} \quad (11)$$

Because  $y \langle \mathbf{w}, \mathbf{x} \rangle < 0$  and  $\|\mathbf{x}\|_2 \leq C$ ,

$$\begin{aligned} \langle \mathbf{w} + y\mathbf{x}, \mathbf{w} + y\mathbf{x} \rangle &= \|\mathbf{w}\|_2^2 + 2y \langle \mathbf{w}, \mathbf{x} \rangle + \|\mathbf{x}\|_2^2 \\ &\leq \|\mathbf{w}\|_2^2 + C^2 \end{aligned} \quad (12)$$

Finally, suppose it converges after  $M$  updates, we have  $\langle \mathbf{w}, \mathbf{w}^* \rangle \geq M\gamma$  and  $\|\mathbf{w}\|_2^2 \leq MC^2$

$$\begin{aligned} 1 = \cos \theta &= \frac{\langle \mathbf{w}, \mathbf{w}^* \rangle}{\|\mathbf{w}\|_2 \|\mathbf{w}^*\|_2} \\ &\geq \frac{M\gamma}{\sqrt{MC^2} \times 1} \\ &= \sqrt{M} \frac{\gamma}{C} \end{aligned} \quad (13)$$

which means  $M \leq \frac{C^2}{\gamma^2}$ .

**Question 8** (Perceptron Loss).

$$\begin{aligned} l(\mathbf{w}, \mathbf{x}_t, y_t) &= -y_t \langle \mathbf{w}, \mathbf{x}_t \rangle \mathbb{I}[\text{mistake on } \mathbf{x}_t] \\ &= -\min \{y_t \langle \mathbf{w}, \mathbf{x}_t \rangle, 0\} \end{aligned} \quad (14)$$

$$L(\mathbf{w}) = -\frac{1}{n} \sum_{t=1}^n y_t \langle \mathbf{w}, \mathbf{x}_t \rangle \mathbb{I}[\text{mistake on } \mathbf{x}_t] \quad (15)$$

## 2 Linear Regression

**Question 9** (Least Square Regression).

$$\min_{f: \mathcal{X} \rightarrow \mathcal{Y}} \mathbb{E} \|f(\mathbf{X}) - Y\|_2^2 \quad (16)$$

The optimal regression function is

$$f^*(\mathbf{x}) = m(x) = \mathbb{E}[Y | \mathbf{X} = \mathbf{x}] \quad (17)$$

Calculating it needs to know the distribution, i.e., all pairs  $(\mathbf{X}, Y)$ .

**Question 10** (Bias-Variance Decomposition).

$$\begin{aligned} \mathbb{E} \|f(\mathbf{X}) - Y\|_2^2 &= \mathbb{E} \|f(\mathbf{X}) - m(x) + m(x) - Y\|_2^2 \\ &= \mathbb{E} \|f(\mathbf{X}) - m(x)\|_2^2 + \mathbb{E} \|m(x) - Y\|_2^2 + 2\mathbb{E} \langle f(\mathbf{X}) - m(x), m(x) - Y \rangle \\ &= \mathbb{E} \|f(\mathbf{X}) - m(x)\|_2^2 + \mathbb{E} \|m(x) - Y\|_2^2 + \mathbb{E} \mathbb{E}_{Y|\mathbf{X}} [\langle f(\mathbf{X}) - m(x), m(x) - Y \rangle] \\ &= \mathbb{E} \|f(\mathbf{X}) - m(x)\|_2^2 + \mathbb{E} \|m(x) - Y\|_2^2 + \mathbb{E} \langle f(\mathbf{X}) - m(x), m(x) - \mathbb{E}_{Y|\mathbf{X}}[Y] \rangle \\ &= \mathbb{E} \|f(\mathbf{X}) - m(x)\|_2^2 + \mathbb{E} \|m(x) - Y\|_2^2 + \mathbb{E} \langle f(\mathbf{X}) - m(x), m(x) - m(x) \rangle \\ &= \mathbb{E} \|f(\mathbf{X}) - m(x)\|_2^2 + \underbrace{\mathbb{E} \|m(x) - Y\|_2^2}_{\text{noise (variance)}} \end{aligned} \quad (18)$$

The last term is the noise (variance), irrelevant to  $f$ . So, to minimize the squared error, we need  $f \approx m$ . However,  $m(\mathbf{x})$  is incalculable, because  $\mathbb{E}[Y|\mathbf{X} = \mathbf{x}]$  is unknown. Let's learn  $f_D$  from the training data  $D$ . Define  $\bar{f}(\mathbf{X}) = \mathbb{E}_D[f_D(\mathbf{X})]$ .

$$\begin{aligned}
\underbrace{\mathbb{E}_{\mathbf{X}, Y, D} \|f_D(\mathbf{X}) - Y\|_2^2}_{\text{test error}} &= \mathbb{E}_{\mathbf{X}} \|f_D(\mathbf{X}) - m(x)\|_2^2 + \mathbb{E}_{\mathbf{X}, Y} \|m(x) - Y\|_2^2 \\
&= \mathbb{E}_{\mathbf{X}, D} \|f_D(\mathbf{X}) - \bar{f}(\mathbf{X}) + \bar{f}(\mathbf{X}) - m(x)\|_2^2 + \mathbb{E}_{\mathbf{X}, Y} \|m(x) - Y\|_2^2 \\
&= \mathbb{E}_{\mathbf{X}, D} \|f_D(\mathbf{X}) - \bar{f}(\mathbf{X})\|_2^2 + \mathbb{E}_{\mathbf{X}} \|\bar{f}(\mathbf{X}) - m(x)\|_2^2 \\
&\quad + 2\mathbb{E}_{\mathbf{X}, D} \langle f_D(\mathbf{X}) - \bar{f}(\mathbf{X}), \bar{f}(\mathbf{X}) - m(x) \rangle \\
&\quad + \mathbb{E}_{\mathbf{X}, Y} \|m(x) - Y\|_2^2 \\
&= \dots + 2\mathbb{E}_{\mathbf{X}} \mathbb{E}_D \langle f_D(\mathbf{X}) - \bar{f}(\mathbf{X}), \bar{f}(\mathbf{X}) - m(x) \rangle + \dots \\
&= \dots + 2\mathbb{E}_{\mathbf{X}} \langle \mathbb{E}_D[f_D(\mathbf{X})] - \bar{f}(\mathbf{X}), \bar{f}(\mathbf{X}) - m(x) \rangle + \dots \\
&= \dots + 2\mathbb{E}_{\mathbf{X}} \langle \bar{f}(\mathbf{X}) - \bar{f}(\mathbf{X}), \bar{f}(\mathbf{X}) - m(x) \rangle + \dots \\
&= \dots + 0 + \dots \\
&= \mathbb{E}_{\mathbf{X}, D} \|f_D(\mathbf{X}) - \bar{f}(\mathbf{X})\|_2^2 + \mathbb{E}_{\mathbf{X}} \|\bar{f}(\mathbf{X}) - m(x)\|_2^2 + \mathbb{E}_{\mathbf{X}, Y} \|m(x) - Y\|_2^2 \\
&= \underbrace{\mathbb{E}_{\mathbf{X}, D} \|f_D(\mathbf{X}) - \mathbb{E}_D[f_D(\mathbf{X})]\|_2^2}_{\text{variance}} + \underbrace{\mathbb{E}_{\mathbf{X}} \|\mathbb{E}_D[f_D(\mathbf{X})] - m(x)\|_2^2}_{\text{bias}^2} + \underbrace{\mathbb{E}_{\mathbf{X}, Y} \|m(x) - Y\|_2^2}_{\text{noise (variance)}}
\end{aligned} \tag{19}$$

**Question 11** (Sampling  $\rightarrow$  Training). Replace expectation with sample average:  $(\mathbf{X}_i, Y_i) \tilde{P}$ .

$$\min_{f: \mathcal{X} \rightarrow \mathcal{Y}} \hat{\mathbb{E}} \|f(\mathbf{X}) - Y\|_2^2 := \frac{1}{n} \sum_{i=1}^n \|f(\mathbf{X}_i) - Y_i\|_2^2 \tag{20}$$

Uniform law of large numbers: as training data size  $n \rightarrow \text{argmin } \mathbb{E}$ ,  $\hat{\mathbb{E}} \rightarrow \mathbb{E}$  and (hopefully)  $\text{argmin } \hat{\mathbb{E}} \rightarrow \mathbb{E}$ .

**Question 12** (Linear Regression). Padding:  $\mathbf{x} \leftarrow \begin{pmatrix} \mathbf{x} \\ 1 \end{pmatrix}$ ,  $W \leftarrow \begin{pmatrix} W \\ \mathbf{b} \end{pmatrix}$

$$\begin{aligned}
X &= [\mathbf{x}_1, \dots, \mathbf{x}_n] \in \mathbb{R}^{(d+1) \times n}, \\
Y &= [\mathbf{y}_1, \dots, \mathbf{y}_n] \in \mathbb{R}^{t \times n}, \\
W &\in \mathbb{R}^{t \times (d+1)}, \\
\|A\|_F &= \sqrt{\sum_{ij} a_{ij}^2}
\end{aligned}$$

Linear regression is:

$$\min_{W \in \mathbb{R}^{t \times (d+1)}} \frac{1}{n} \|WX - Y\|_F^2 \tag{21}$$

**Question 13** (Optimality Condition). If  $\mathbf{w}$  is a minimizer (or maximizer) of a differentiable function  $f$  **over an open set**, then  $f'(\mathbf{w}) = 0$ .

**Question 14** (Solving Linear Regression).

$$\begin{aligned}
L(W) &= \frac{1}{n} \|WX - Y\|_F^2 \\
\nabla_W L(W) &= \frac{2}{n} (WX - Y)X^T = 0 \\
&\Rightarrow WX X^T = Y X^T \\
&\Rightarrow W = Y X^T (X X^T)^{-1}
\end{aligned} \tag{22}$$

$$\tag{23}$$

**Question 15** (Ill-Conditioning). Slight perturbation leads to chaotic behavior, which happens whenever  $X$  is ill-conditioned, i.e., (close to) rank-deficient.

Rank-deficient  $X$  means:

1. two columns in  $X$  are linearly dependent (or simply the same)
2. but the corresponding  $y$  might be different

**Question 16** (Ridge Regression).

$$\min_W \frac{1}{n} \|WX - Y\|_F^2 + \lambda \|W\|_F^2 \quad (24)$$

$$\begin{aligned} \nabla_W L(W) &= \frac{2}{n} (WX - Y)X^T + 2\lambda W = 0 \\ \Rightarrow WX X^T - YX^T + \lambda W &= 0 \\ \Rightarrow W(XX^T + n\lambda I) &= YX^T \end{aligned} \quad (25)$$

$$\begin{aligned} X &= U\Sigma V^T \\ \Rightarrow XX^T &= U\Sigma(V^T V)\Sigma U^T = U\Sigma^2 U^T \\ \Rightarrow XX^T + n\lambda I &= U \underbrace{(\Sigma^2 + n\lambda I)}_{\text{strictly positive}} U^T \\ \Rightarrow XX^T + n\lambda I &\text{ is of full-rank} \end{aligned} \quad (26)$$

$\lambda$  is regularization parameter.  $\lambda = \infty \Rightarrow W \equiv \mathbf{0}$ .

**Question 17** (Regularization  $\equiv$  Data Augmentation).

$$\frac{1}{n} \|WX - Y\|_F^2 + \lambda \|W\|_F^2 = \frac{1}{n} \left\| W \begin{bmatrix} X & \sqrt{n\lambda} I \end{bmatrix} - \begin{bmatrix} Y & \mathbf{0} \end{bmatrix} \right\|_F^2 \quad (27)$$

### 3 Logistic Regression

**Question 18** (Max Likelihood Estimation). Let  $\mathcal{Y} = \{0, 1\}$ . Learn confidence  $p(\mathbf{x}; \mathbf{w}) := \Pr(Y = 1 | X = \mathbf{x})$ .

$$\begin{aligned} \max_{\mathbf{w}} \Pr(Y_1 = y_1, \dots, Y_n = y_n) &= \max_{\mathbf{w}} \prod_{i=1}^n \Pr(Y_i = y_i | X_i = x_i) \\ &\stackrel{\mathcal{Y}=\{0,1\}}{=} \max_{\mathbf{w}} \prod_{i=1}^n [p(\mathbf{x}_i; \mathbf{w})]^{y_i} [1 - p(\mathbf{x}_i; \mathbf{w})]^{1-y_i} \end{aligned} \quad (28)$$

Use negative log-likelihood:

$$\min_{\mathbf{w}} \sum_{i=1}^n [-y_i \log p(\mathbf{x}_i; \mathbf{w}) - (1 - y_i) \log(1 - p(\mathbf{x}_i; \mathbf{w}))] \quad (29)$$

**Question 19** (Odds Ratio and Sigmoid).

$$\text{Odds Ratio} = \frac{\Pr}{1 - \Pr} \quad (30)$$

Assume  $\log \frac{p(\mathbf{x}; \mathbf{w})}{1 - p(\mathbf{x}; \mathbf{w})} = \langle \mathbf{x}, \mathbf{w} \rangle$ .

The Sigmoid transformation is:

$$p(\mathbf{x}; \mathbf{w}) = \frac{1}{1 + \exp(-\langle \mathbf{x}, \mathbf{w} \rangle)} \quad (31)$$

**Question 20** (Logistic Regression). Plug the sigmoid in the negative log-likelihood:

$$\begin{aligned} &\min_{\mathbf{w}} \sum_{i=1}^n [-y_i \log p(\mathbf{x}_i; \mathbf{w}) - (1 - y_i) \log(1 - p(\mathbf{x}_i; \mathbf{w}))] \\ &= \min_{\mathbf{w}} \sum_{i=1}^n \left[ y_i \log[1 + \exp(-\langle \mathbf{x}_i, \mathbf{w} \rangle)] - (1 - y_i) \log \left( 1 - \frac{1}{1 + \exp(-\langle \mathbf{x}_i, \mathbf{w} \rangle)} \right) \right] \\ &= \min_{\mathbf{w}} \sum_{i=1}^n \left[ y_i \log[1 + \exp(-\langle \mathbf{x}_i, \mathbf{w} \rangle)] - (1 - y_i) \log \left( \frac{\exp(-\langle \mathbf{x}_i, \mathbf{w} \rangle)}{1 + \exp(-\langle \mathbf{x}_i, \mathbf{w} \rangle)} \right) \right] \\ &= \min_{\mathbf{w}} \sum_{i=1}^n [y_i \log[1 + \exp(-\langle \mathbf{x}_i, \mathbf{w} \rangle)] + (1 - y_i) [\langle \mathbf{x}_i, \mathbf{w} \rangle + \log(1 + \exp(-\langle \mathbf{x}_i, \mathbf{w} \rangle))] ] \\ &= \min_{\mathbf{w}} \sum_{i=1}^n [\log[1 + \exp(-\langle \mathbf{x}_i, \mathbf{w} \rangle)] + (1 - y_i) \langle \mathbf{x}_i, \mathbf{w} \rangle] \end{aligned} \quad (32)$$

Because  $y_i \in \{0, 1\}$ , let's map it to  $\{\pm 1\}$ .

$$\begin{aligned}
L(\mathbf{w}) &\stackrel{y_i \in \{0,1\}}{=} \log[1 + \exp(-\langle \mathbf{x}_i, \mathbf{w} \rangle)] + (1 - y_i) \langle \mathbf{x}_i, \mathbf{w} \rangle \\
&\stackrel{y_i \in \{0,1\}}{=} \log[1 + \exp(-\langle \mathbf{x}_i, \mathbf{w} \rangle)] + \log[\exp((1 - y_i) \langle \mathbf{x}_i, \mathbf{w} \rangle)] \\
&\stackrel{y_i \in \{0,1\}}{=} \log[1 + \exp(-\langle \mathbf{x}_i, \mathbf{w} \rangle) \cdot \exp((1 - y_i) \langle \mathbf{x}_i, \mathbf{w} \rangle)] \\
&\stackrel{y_i \in \{0,1\}}{=} \log[\exp((1 - y_i) \langle \mathbf{x}_i, \mathbf{w} \rangle) + \exp(-y_i \langle \mathbf{x}_i, \mathbf{w} \rangle)] \\
&\stackrel{y_i \in \{0,1\}}{=} \begin{cases} \log[\exp(\langle \mathbf{x}_i, \mathbf{w} \rangle) + 1] & y_i = 0 \\ \log[1 + \exp(-\langle \mathbf{x}_i, \mathbf{w} \rangle)] & y_i = 1 \end{cases} \\
&\stackrel{y_i \in \{\pm 1\}}{=} \log[1 + \exp(-y_i \langle \mathbf{x}_i, \mathbf{w} \rangle)]
\end{aligned} \tag{33}$$

**Question 21** (Multi-Class: Sigmoid  $\rightarrow$  Softmax).

$$\Pr(Y = k | X = \mathbf{x}; \mathbf{W} = [\mathbf{w}_1, \dots, \mathbf{w}_c]) = \frac{\exp(\langle \mathbf{x}, \mathbf{w}_k \rangle)}{\sum_{l=1}^c \exp(\langle \mathbf{x}, \mathbf{w}_l \rangle)} \tag{34}$$

Maximum likelihood estimation (log loss, cross-entropy loss):

$$\min_{\mathbf{W}} \sum_{i=1}^n \left[ -\log \frac{\exp(\langle \mathbf{x}, \mathbf{w}_k \rangle)}{\sum_{l=1}^c \exp(\langle \mathbf{x}, \mathbf{w}_l \rangle)} \right] \tag{35}$$

## 4 Hard-Margin Support Vector Machines

**Question 22** (Distance from a Point to a Hyperplane). Let  $H := \{\mathbf{x} : \langle \mathbf{x}, \mathbf{w} \rangle + b = 0\}$ ,  $\mathbf{x}$  be any vector in  $H$ .

$$\text{Distance}(\mathbf{x}_i, \mathbf{w}) = \frac{|\langle \mathbf{x}_i - \mathbf{x}, \mathbf{w} \rangle|}{\|\mathbf{w}\|_2} = \frac{|\langle \mathbf{x}_i, \mathbf{w} \rangle - \langle \mathbf{x}, \mathbf{w} \rangle|}{\|\mathbf{w}\|_2} \stackrel{\mathbf{x} \in H}{=} \frac{|\langle \mathbf{x}_i, \mathbf{w} \rangle + b|}{\|\mathbf{w}\|_2} \stackrel{y_i \hat{y}_i > 0}{=} \frac{y_i \hat{y}_i}{\|\mathbf{w}\|_2} \tag{36}$$

**Question 23** (Margin Maximization). Margin is the smallest distance to  $H$  among all separable data.

$$\max_{\mathbf{w}, b} \min_i \frac{y_i \hat{y}_i}{\|\mathbf{w}\|_2}, \text{ such that } \forall i, y_i \hat{y}_i > 0 \tag{37}$$

Let  $c > 0$ , then  $\mathbf{w} = c\mathbf{w}, b = cb$  keeps the loss same:

$$\begin{aligned}
\max_{\mathbf{w}, b} \min_i \frac{cy_i \hat{y}_i}{\|c\mathbf{w}\|_2} &= \max_{\mathbf{w}, b} \min_i \frac{y_i(\langle \mathbf{x}, c\mathbf{w} \rangle + cb)}{\|c\mathbf{w}\|_2} \\
&= \max_{\mathbf{w}, b} \min_i \frac{cy_i(\langle \mathbf{x}, \mathbf{w} \rangle + b)}{c\|\mathbf{w}\|_2} \\
&= \max_{\mathbf{w}, b} \min_i \frac{y_i(\langle \mathbf{x}, \mathbf{w} \rangle + b)}{\|\mathbf{w}\|_2} \\
&= \max_{\mathbf{w}, b} \min_i \frac{y_i \hat{y}_i}{\|\mathbf{w}\|_2}
\end{aligned} \tag{38}$$

Let  $c = \frac{1}{\min_i y_i \hat{y}_i}$ ,

$$\begin{aligned}
\max_{\mathbf{w}, b} \min_i \frac{cy_i \hat{y}_i}{c\|\mathbf{w}\|_2} &= \max_{\mathbf{w}, b} \frac{1}{c\|\mathbf{w}\|_2} \\
&= \max_{\mathbf{w}, b} \frac{1}{\|\mathbf{w}\|_2} \text{ s.t. } \min_i y_i \hat{y}_i = 1
\end{aligned} \tag{39}$$

Max  $\rightarrow$  Min:

$$\min_{\mathbf{w}, b} \frac{1}{2} \|\mathbf{w}\|_2^2 \text{ s.t. } \forall i, y_i \hat{y}_i \geq 1 \tag{40}$$

**Question 24** (Hard-Margin SVM v.s. Perceptron).

$$\text{Hard-Margin SVM: } \min_{\mathbf{w}, b} \frac{1}{2} \|\mathbf{w}\|_2^2 \quad \text{s.t. } \forall i, y_i \hat{y}_i \geq 1 \tag{41}$$

$$\text{Perceptron: } \min_{\mathbf{w}, b} 0 \quad \text{s.t. } \forall i, y_i \hat{y}_i \geq 1 \tag{42}$$

**Question 25** (Lagrangian Dual). Dual variables  $\alpha \in \mathbb{R}^n$ .

$$\begin{aligned}
\min_{\mathbf{w}, b} \max_{\alpha \geq 0} \frac{1}{2} \|\mathbf{w}\|_2^2 - \sum_i \alpha_i [y_i(\langle \mathbf{x}_i, \mathbf{w} \rangle + b) - 1] &= \min_{\mathbf{w}, b} \begin{cases} +\infty, & \text{if } \exists i, y_i(\langle \mathbf{x}_i, \mathbf{w} \rangle + b) < 1 \quad (\alpha_i = +\infty) \\ \frac{1}{2} \|\mathbf{w}\|_2^2, & \text{if } \forall i, y_i(\langle \mathbf{x}_i, \mathbf{w} \rangle + b) \geq 1 \quad (\forall i, \alpha_i = 0) \end{cases} \\
&= \min_{\mathbf{w}, b} \frac{1}{2} \|\mathbf{w}\|_2^2, \quad \text{s.t. } \forall i, y_i(\langle \mathbf{x}_i, \mathbf{w} \rangle + b) \geq 1 \\
&= \min_{\mathbf{w}, b} \frac{1}{2} \|\mathbf{w}\|_2^2 \quad \text{s.t. } \forall i, y_i \hat{y}_i \geq 1
\end{aligned} \tag{43}$$

Swap min and max:

$$\max_{\alpha \geq 0} \min_{\mathbf{w}, b} \frac{1}{2} \|\mathbf{w}\|_2^2 - \sum_i \alpha_i [y_i(\langle \mathbf{x}_i, \mathbf{w} \rangle + b) - 1] \tag{44}$$

Solve inner problem by setting derivative to 0:

$$\frac{\delta}{\delta \mathbf{w}} = \mathbf{w} - \sum_i \alpha_i y_i \mathbf{x}_i = 0, \quad \frac{\delta}{\delta b} = - \sum_i \alpha_i y_i = 0, \tag{45}$$

Plug them into the loss:

$$\begin{aligned}
L(\alpha) &= \min_{\mathbf{w}, b} \frac{1}{2} \|\mathbf{w}\|_2^2 - \sum_i \alpha_i [y_i(\langle \mathbf{x}_i, \mathbf{w} \rangle + b) - 1] \\
&= \frac{1}{2} \left\| \sum_i \alpha_i y_i \mathbf{x}_i \right\|_2^2 - \sum_i \alpha_i y_i \langle \mathbf{x}_i, \mathbf{w} \rangle - \sum_i \alpha_i y_i b + \sum_i \alpha_i \\
&= \frac{1}{2} \left\| \sum_i \alpha_i y_i \mathbf{x}_i \right\|_2^2 - \left\langle \sum_i \alpha_i y_i \mathbf{x}_i, \sum_i \alpha_i y_i \mathbf{x}_i \right\rangle - b \sum_i \alpha_i y_i + \sum_i \alpha_i \\
&= \frac{1}{2} \left\| \sum_i \alpha_i y_i \mathbf{x}_i \right\|_2^2 - \left\| \sum_i \alpha_i y_i \mathbf{x}_i \right\|_2^2 + \sum_i \alpha_i \\
&= \sum_i \alpha_i - \frac{1}{2} \left\| \sum_i \alpha_i y_i \mathbf{x}_i \right\|_2^2, \quad \text{s.t. } \sum_i \alpha_i y_i = 0
\end{aligned} \tag{46}$$

So, 44 is solved as:

$$\max_{\alpha \geq 0} \sum_i \alpha_i - \frac{1}{2} \left\| \sum_i \alpha_i y_i \mathbf{x}_i \right\|_2^2, \quad \text{s.t. } \sum_i \alpha_i y_i = 0 \tag{47}$$

Max  $\rightarrow$  min and expand the norm:

$$\min_{\alpha \geq 0} - \sum_i \alpha_i + \frac{1}{2} \sum_i \sum_j \alpha_i \alpha_j y_i y_j \underbrace{\langle \mathbf{x}_i, \mathbf{x}_j \rangle}_{\text{Kernel, closed form w.r.t. } \mathbf{x}_i, \mathbf{x}_j}, \quad \text{s.t. } \sum_i \alpha_i y_i = 0 \tag{48}$$

**Question 26** (Support Vectors). From 45, we know  $\mathbf{w} = \sum_i \alpha_i y_i \mathbf{x}_i$ . Vectors with  $\alpha_i \neq 0$  are support vectors, which lie on the margin.

## 5 Soft-Margin Support Vector Machines

**Question 27** (Goal). minimize over  $\mathbf{w}, b$ ,

$$\Pr(Y \neq \text{sign}(\hat{Y})) = \Pr(Y \hat{Y} \leq 0) = \mathbb{E} \underbrace{\mathbb{I}[Y \hat{Y} \leq 0]}_{\text{indicator function}} := \mathbb{E} l_{0-1}(Y \hat{Y}) \tag{49}$$

where  $\hat{Y} = \langle X, \mathbf{w} \rangle + b, Y = \pm 1$ .

$$\begin{aligned}
\min_{\hat{Y}: \mathcal{X} \rightarrow \mathbb{R}} \mathbb{E} l_{0-1}(Y \hat{Y}) &= \min_{\hat{Y}: \mathcal{X} \rightarrow \mathbb{R}} \mathbb{E}_X \mathbb{E}_{Y|X} l_{0-1}(Y \hat{Y}) \\
&= \mathbb{E}_X \min_{\hat{Y}: \mathcal{X} \rightarrow \mathbb{R}} \mathbb{E}_{Y|X} l_{0-1}(Y \hat{Y})
\end{aligned} \tag{50}$$

Minimizing the 0-1 error is **NP-hard**.

**Question 28** (Bayes Rule).

$$\eta(\mathbf{x}) := \operatorname{argmax}_{\hat{y} \in \mathbb{R}} \Pr(Y = \hat{y} | X = \mathbf{x}) \quad (51)$$

$$\eta(\mathbf{x}) := \operatorname{argmin}_{\hat{y} \in \mathbb{R}} \mathbb{E}_{Y|X=\mathbf{x}} l_{0-1}(Y\hat{y}) \quad (52)$$

**Question 29** (Classification Calibrated). A loss  $l(y\hat{y})$  is classification calibrated, iff  $\forall \mathbf{x}$ ,

$$\hat{y}(\mathbf{x}) := \operatorname{argmin}_{\hat{y} \in \mathbb{R}} \mathbb{E}_{Y|X=\mathbf{x}} l(Y\hat{y})$$

has the same sign as the Bayes rule  $\eta(\mathbf{x}) := \operatorname{argmin}_{\hat{y} \in \mathbb{R}} \mathbb{E}_{Y|X=\mathbf{x}} l_{0-1}(Y\hat{y})$

Notice:  $\eta(\mathbf{x}), \hat{y}(\mathbf{x})$  provide the **score**, their sign provides the prediction.

**Question 30** (Characterization under Convexity). Any **convex** loss  $l$  is classification calibrated iff  $l$  is differentiable at 0 and  $l'(0) < 0$ .

**Question 31** (Hinge Loss).

$$l_{\text{hinge}}(y\hat{y}) = (1 - y\hat{y})^+ := \max\{0, 1 - y\hat{y}\} = \begin{cases} 1 - y\hat{y}, & \text{if } y\hat{y} < 1 \\ 0, & \text{otherwise} \end{cases} \quad (53)$$

The classifier that minimizes the expected hinge loss minimizes the expected 0-1 loss.

**Question 32** (Soft-Margin SVM).

$$\min_{\mathbf{w}, b} \frac{1}{2} \|\mathbf{w}\|_2^2 + C \cdot \sum_i l_{\text{hinge}}(y_i \hat{y}_i), \quad \text{s.t. } \hat{y}_i = \langle \mathbf{x}_i, \mathbf{w} \rangle + b \quad (54)$$

**Question 33** (Lagrangian Dual). Apply  $C \cdot l_{\text{hinge}}(t) := \max\{0, C(1 - t)\} = \max_{0 \leq \alpha \leq C} \alpha(1 - t)$

$$\min_{\mathbf{w}, b} \max_{0 \leq \alpha \leq C} \frac{1}{2} \|\mathbf{w}\|_2^2 + \sum_i \alpha_i (1 - y_i \hat{y}_i) \quad (55)$$

Swap min with max:

$$\max_{0 \leq \alpha \leq C} \min_{\mathbf{w}, b} \frac{1}{2} \|\mathbf{w}\|_2^2 + \sum_i \alpha_i (1 - y_i \hat{y}_i) \quad (56)$$

Solve it by setting derivative to 0:

$$\frac{\delta}{\delta \mathbf{w}} = \mathbf{w} - \sum_i \alpha_i y_i \mathbf{x}_i = 0, \quad \frac{\delta}{\delta b} = - \sum_i \alpha_i y_i = 0, \quad (57)$$

Plug them into the loss:

$$\begin{aligned} \max_{0 \leq \alpha \leq C} \min_{\mathbf{w}, b} \frac{1}{2} \|\mathbf{w}\|_2^2 + \sum_i \alpha_i (1 - y_i \hat{y}_i) &= \max_{0 \leq \alpha \leq C} \min_{\mathbf{w}, b} \frac{1}{2} \|\mathbf{w}\|_2^2 + \sum_i \alpha_i [1 - y_i (\langle \mathbf{x}_i, \mathbf{w} \rangle + b)] \\ &= \max_{0 \leq \alpha \leq C} \sum_i \alpha_i - \frac{1}{2} \left\| \sum_i \alpha_i y_i \mathbf{x}_i \right\|_2^2, \quad \text{s.t. } \sum_i \alpha_i y_i = 0 \end{aligned} \quad (58)$$

Max  $\rightarrow$  min and expand the norm:

$$\min_{0 \leq \alpha \leq C} - \sum_i \alpha_i + \frac{1}{2} \sum_i \sum_j \alpha_i \alpha_j y_i y_j \underbrace{\langle \mathbf{x}_i, \mathbf{x}_j \rangle}_{\text{Kernel, closed form w.r.t. } \mathbf{x}_i, \mathbf{x}_j}, \quad \text{s.t. } \sum_i \alpha_i y_i = 0 \quad (59)$$

$C \rightarrow \infty \Rightarrow$  Hard-margin SVM,  $C \rightarrow 0 \Rightarrow$  a constant classifier

## 6 Reproducing Kernels

**Question 34** ((Reproducing) Kernels).  $k : (X) \times \mathcal{X} \rightarrow \mathbb{R}$  is a (reproducing) kernel iff there exists some  $\Phi : \mathcal{X} \rightarrow \mathcal{H}$  so that  $\langle \Phi(\mathbf{x}), \Phi(\mathbf{z}) \rangle = k(\mathbf{x}, \mathbf{z})$ .

- A feature transform  $\Phi$  determines the corresponding kernel  $k$ .
- A kernel  $k$  determines some feature transforms  $\Phi$ , but may not be unique.  
E.g.  $\langle \phi(\mathbf{x}), \phi(\mathbf{z}) \rangle = \langle \phi'(\mathbf{x}), \phi'(\mathbf{z}) \rangle$

1.  $\phi(x) := [x_1^2, \sqrt{2}x_1x_2] \in \mathbb{R}^2$
2.  $\phi'(x) := [x_1^2, x_1x_2, x_1x_2] \in \mathbb{R}^3$



**Question 35** (Mercer's Theorem).  $k : \mathcal{X} \times \mathcal{X} \rightarrow \mathbb{R}$  is a kernel,  
iff  $\forall n \in \mathbb{N}, \forall \mathbf{x}_1, \dots, \mathbf{x}_n \in \mathcal{X}$ , the kernel matrix  $K$  such that  $K_{ij} := k(\mathbf{x}_i, \mathbf{x}_j)$  is symmetric and positive semi-definite (PSD).

$$k \text{ is a kernel} \Leftrightarrow \begin{cases} K_{ij} = K_{ji} & \text{(symmetric)} \\ \langle \boldsymbol{\alpha}, K \boldsymbol{\alpha} \rangle = \sum_{i=1}^n \sum_{j=1}^n \alpha_i \alpha_j K_{ij} \geq 0 & \forall \boldsymbol{\alpha} \in \mathbb{R}^n \quad \text{(PSD)} \end{cases} \quad (60)$$

**Question 36** (Symmetric PSD). For a symmetric matrix  $A$ , the following conditions are equivalent.

1.  $A \succeq 0$
2.  $A = U^T U$  for some matrix  $U$
3.  $x^T A x \geq 0$  for every  $x \in \mathbb{R}^n$
4. All principal minors of  $A$  are nonnegative

## 7 Gradient Descent

**Question 37** (Gradient Descent Template). Choose initial point  $x^{(0)} \in \mathbb{R}^d$  and repeat:

$$x^{(k)} = x^{(k-1)} - \underbrace{\eta}_{\text{step size}} \nabla f(x^{(k-1)}), \quad k = 1, 2, \dots \quad (61)$$

**Question 38** (Interpretation from Taylor Expansion). Expand  $f$  locally at  $x$ :

$$\begin{aligned} f(y) &\approx f(x) + \nabla f(x)^T (y - x) + \frac{1}{2t} \|y - x\|_2^2 \\ \Rightarrow \min_y f(y) &\approx \min_y \left[ f(x) + \nabla f(x)^T (y - x) + \frac{1}{2t} \|y - x\|_2^2 \right] \end{aligned} \quad (62)$$

When  $y - x = \frac{\nabla f(x)}{-\frac{1}{2t}} = -t \nabla f(x) \implies y = x - t \nabla f(x)$ , it reaches the minimum.

**Question 39** ( $L$ -smooth or  $L$ -Lipschitz Continuous).  $f$  is convex and differentiable.  $\nabla f$  is  $L$ -Lipschitz continuous ( $L$ -smooth):

$$L \mathbf{I} \succeq \nabla^2 f(x), \forall x \quad (63)$$

**Question 40** (Convergence Rate for Convex Case). Assume  $f$  is  $L$ -smooth. Gradient descent with fixed step size  $t \leq \frac{1}{L}$  satisfies:

$$f(x^{(k)}) - f(x^*) \leq \frac{\|x^{(0)} - x^*\|_2^2}{2tk} \quad (64)$$

We say gradient descent has convergence rate  $O(\frac{1}{k})$ , i.e.  $f(x^{(k)}) - f(x^*) \leq \epsilon$  can be achieved using only  $O(\frac{1}{\epsilon})$  iterations.

**Proof**

$$\begin{aligned} f(y) &= f(x) + \nabla f(x)^T (y - x) + \frac{1}{2} (y - x)^T \nabla^2 f(\xi) (y - x) \\ &\leq f(x) + \nabla f(x)^T (y - x) + \frac{1}{2} L \|y - x\|_2^2 \quad (L\text{-smooth}, L \mathbf{I} \succeq \nabla^2 f(\xi)) \end{aligned} \quad (65)$$

Plug in gradient descent:

$$\begin{aligned} f(x^+) &= f(y) \\ &\leq f(x) + \nabla f(x)^T (x - t \nabla f(x) - x) + \frac{1}{2} L \|x - t \nabla f(x) - x\|_2^2 \\ &= f(x) - (1 - \frac{1}{2} L t) t \|\nabla f(x)\|_2^2 \\ &\leq f(x) - \frac{1}{2} t \|\nabla f(x)\|_2^2 \quad (t \leq \frac{1}{L}) \end{aligned} \quad (66)$$

$$f \text{ is convex} \Rightarrow f(x^*) \geq f(x) + \nabla f(x)^T (x^* - x) \Rightarrow f(x) \leq f(x^*) + \nabla f(x)^T (x - x^*)$$

Plug this into 66:

$$\begin{aligned}
f(x^+) &\leq f(x^*) + \nabla f(x)^T(x - x^*) - \frac{1}{2}t \|\nabla f(x)\|_2^2 \\
\Rightarrow f(x^+) - f(x^*) &\leq \frac{1}{2t} \left( 2t \nabla f(x)^T(x - x^*) - t^2 \|\nabla f(x)\|_2^2 \right) \\
\Rightarrow f(x^+) - f(x^*) &\leq \frac{1}{2t} \left( 2t \nabla f(x)^T(x - x^*) - t^2 \|\nabla f(x)\|_2^2 - \|x - x^*\|_2^2 + \|x - x^*\|_2^2 \right) \\
\Rightarrow f(x^+) - f(x^*) &\leq \frac{1}{2t} \left( \|x - x^*\|_2^2 - (\|x - x^*\|_2^2 + t^2 \|\nabla f(x)\|_2^2 - 2t \nabla f(x)^T(x - x^*)) \right) \\
\Rightarrow f(x^+) - f(x^*) &\leq \frac{1}{2t} \left( \|x - x^*\|_2^2 - \|x - x^* - t \nabla f(x)\|_2^2 \right) \\
\Rightarrow f(x^+) - f(x^*) &\leq \frac{1}{2t} \left( \|x - x^*\|_2^2 - \|x^+ - x^*\|_2^2 \right)
\end{aligned} \tag{67}$$

Viewing  $x^+$  as  $x^{(i)}$  and  $x$  as  $x^{(i-1)}$ :

$$\begin{aligned}
\sum_{i=1}^k \left( f(x^{(i)}) - f(x^*) \right) &\leq \sum_{i=1}^k \frac{1}{2t} \left( \|x^{(i-1)} - x^*\|_2^2 - \|x^{(i)} - x^*\|_2^2 \right) \\
&= \frac{1}{2t} \left( \|x^{(0)} - x^*\|_2^2 - \|x^{(k)} - x^*\|_2^2 \right) \\
&\leq \frac{1}{2t} \|x^{(0)} - x^*\|_2^2
\end{aligned} \tag{68}$$

which implies

$$f(x^{(k)}) \leq \frac{1}{k} \sum_{i=1}^k f(x^{(i)}) \leq f(x^*) + \frac{\|x^{(0)} - x^*\|_2^2}{2tk} \tag{69}$$

**Question 41** (Convergence Rate for Strong Convexity).  $f$  is differentiable,  $L$ -smooth, and  $m$ -strongly convex.  $m$ -strong convexity of  $f$  means  $f(x) - \frac{m}{2} \|x\|_2^2$  is convex, i.e.  $\nabla^2 f(x) \succeq m\mathbf{I}$ . Then, there is a constant  $0 < \gamma < 1$  such that gradient descent with fixed step size  $t \leq \frac{2}{m+L}$  satisfies:

$$f(x^{(k)}) - f(x^*) \leq \gamma^k \frac{L}{2} \|x^{(0)} - x^*\|_2^2 \tag{70}$$

Rate is  $O(\gamma^k)$ . Only  $O(\log_{\frac{1}{\gamma}}(\frac{1}{\epsilon}))$  iterations needed.

**Question 42** (Convergence Rate for Non-Convex Case).  $f$  is differentiable and  $L$ -smooth, but non-convex. Gradient descent with fixed step size  $t \leq \frac{1}{L}$  satisfies:

$$\min_{i=0, \dots, k} \|\nabla f(x^{(i)})\|_2 \leq \sqrt{\frac{2(f(x^{(0)}) - f^*)}{t(k+1)}} \tag{71}$$

Rate is  $O(\frac{1}{\sqrt{k}})$  for finding stationary point.  $O(\frac{1}{\epsilon^2})$  iterations are needed.

**Question 43** (Convergence Rate for Stochastic Gradient Descent). For convex and  $L$ -smooth  $f$ ,

- Gradient Descent

$$\mathbf{w}^+ = \mathbf{w} - t \cdot \frac{1}{n} \sum_{i=1}^n \nabla f_i(\mathbf{w}) \tag{72}$$

- Step size  $t \leq \frac{1}{L}$
- Time complexity  $O(\frac{n}{\epsilon})$

- Stochastic Gradient Descent

$$\mathbf{w}^+ = \mathbf{w} - t \cdot \nabla f_{I_{\text{random}}}(\mathbf{w}) \tag{73}$$

- Step size  $t = \frac{1}{k}, k = 1, 2, 3, \dots$  (adaptive step size)
- Time complexity  $O(\frac{1}{\epsilon^2})$

## 8 Fully-Connected Neural Networks

**Question 44** (Forward and Backward Pass of a 2-Layer MLP). A 2-layer MLP ( $k$  is the NN width,  $c$  is the output dim):

$$\mathbf{x} = \text{input} \quad (\mathbf{x} \in \mathbb{R}^d) \quad (74)$$

$$\mathbf{z} = \mathbf{W}\mathbf{x} + \mathbf{b}_1 \quad (\mathbf{W} \in \mathbb{R}^{k \times d}, \mathbf{z}, \mathbf{b} \in \mathbb{R}^k) \quad (75)$$

$$\mathbf{h} = \text{ReLU}(\mathbf{z}) \quad (\mathbf{h} \in \mathbb{R}^k) \quad (76)$$

$$\boldsymbol{\theta} = \mathbf{U}\mathbf{h} + \mathbf{b}_2 \quad (\mathbf{U} \in \mathbb{R}^{c \times k}, \boldsymbol{\theta}, \mathbf{b}_2 \in \mathbb{R}^c) \quad (77)$$

$$J = \frac{1}{2} \|\boldsymbol{\theta} - \mathbf{y}\|_2^2 \quad (\mathbf{y} \in \mathbb{R}^c, J \in \mathbb{R}) \quad (78)$$

$$\text{ReLU} = \begin{cases} x & x > 0 \\ 0 & x \leq 0 \end{cases} \quad (79)$$

$$\text{ReLU}' = \begin{cases} 1 & x > 0 \\ 0 & x \leq 0 \end{cases} \quad (80)$$

Backward pass ( $\odot$  is the Hadamard product, i.e. element-wise product):

$$\frac{\delta J}{\delta \boldsymbol{\theta}} = \boldsymbol{\theta} - \mathbf{y} \quad (81)$$

$$\frac{\delta J}{\delta \mathbf{U}} = \frac{\delta J}{\delta \boldsymbol{\theta}} \odot \frac{\delta \boldsymbol{\theta}}{\delta \mathbf{U}} = (\boldsymbol{\theta} - \mathbf{y})\mathbf{h}^T \quad (82)$$

$$\frac{\delta J}{\delta \mathbf{b}_2} = \frac{\delta J}{\delta \boldsymbol{\theta}} \odot \frac{\delta \boldsymbol{\theta}}{\delta \mathbf{b}_2} = \boldsymbol{\theta} - \mathbf{y} \quad (83)$$

$$\frac{\delta J}{\delta \mathbf{h}} = \frac{\delta J}{\delta \boldsymbol{\theta}} \odot \frac{\delta \boldsymbol{\theta}}{\delta \mathbf{h}} = \mathbf{U}^T(\boldsymbol{\theta} - \mathbf{y}) \quad (84)$$

$$\frac{\delta J}{\delta \mathbf{z}} = \frac{\delta J}{\delta \mathbf{h}} \odot \frac{\delta \mathbf{h}}{\delta \mathbf{z}} = \mathbf{U}^T(\boldsymbol{\theta} - \mathbf{y}) \odot \text{ReLU}'(\mathbf{z}) \quad (85)$$

$$\frac{\delta J}{\delta \mathbf{W}} = \frac{\delta J}{\delta \mathbf{z}} \odot \frac{\delta \mathbf{z}}{\delta \mathbf{W}} = \mathbf{U}^T(\boldsymbol{\theta} - \mathbf{y}) \odot \text{ReLU}'(\mathbf{z})\mathbf{x}^T \quad (86)$$

$$\frac{\delta J}{\delta \mathbf{b}_1} = \frac{\delta J}{\delta \mathbf{z}} \odot \frac{\delta \mathbf{z}}{\delta \mathbf{b}_1} = \mathbf{U}^T(\boldsymbol{\theta} - \mathbf{y}) \odot \text{ReLU}'(\mathbf{z}) \quad (87)$$

$$(88)$$

**Question 45** (Universal Approximation Theorem by 2-Layer NNs). For any continuous function  $f : \mathbb{R}^d \rightarrow \mathbb{R}^c$  and any  $\epsilon > 0$ , there exists  $k \in \mathbb{N}$ ,  $\mathbf{W} \in \mathbb{R}^{k \times d}$ ,  $\mathbf{b} \in \mathbb{R}^k$ ,  $\mathbf{U} \in \mathbb{R}^{c \times k}$  such that

$$\sup_{\mathbf{x}} \|f(\mathbf{x}) - g(\mathbf{x})\|_2 < \epsilon \quad (89)$$

where  $g(\mathbf{x}) = \mathbf{U}(\sigma(\mathbf{W}\mathbf{x} + \mathbf{b}))$  and  $\sigma$  is the element-wise ReLU operation.

As long as the 2-layer MLP is wide enough, it can approximate any continuous function arbitrarily closely.

## 9 Convolutional Neural Networks

**Question 46** (Controlling the Convolution). Hyperparameters.

- **Filter (kernel) size:** width *times* height.
- **Number of filters (kernels).**  
Weights are not shared between different filters (kernels)
- **Stride:** how many pixels the filter moves each time.
- **Padding:** add zeros around the boundary of the input.

**Question 47** (Size Calculation).

Input size:  $m \times n \times c_{in}$

Filter size:  $a \times b \times c_{in}$

Stride:  $s \times t$

Padding:  $p \times q$

Output size:

$$\left\lfloor 1 + \frac{m + 2p - a}{s} \right\rfloor \times \left\lfloor 1 + \frac{n + 2q - b}{t} \right\rfloor \quad (90)$$

## Part II

# For Final

## 10 Transformer

**Question 48** (Attention Layer Inputs and Outputs). Inputs:  $V \in \mathcal{R}^{n \times d}$ ,  $K \in \mathcal{R}^{n \times d}$ ,  $Q \in \mathcal{R}^{m \times d}$ , Outputs: an  $m \times d$  matrix.

- **Self Attention:**  $m = n$ ,
- **Cross Attention:**  $m \neq n$  where  $m$  is the sequence length of decoder,  $n$  is the sequence length of encoder.

**Question 49** (Attention Layer Calculation).

$$\text{Attention}(Q, K, V) = \text{softmax}\left(\frac{QK^T}{\sqrt{d}}\right)V \quad (91)$$

Softmax is row-wise, i.e. for each row of  $QK^T$ , it is normalized to sum to 1.

**Question 50** (Learnable Attention Layer).

$$\text{Attention}(XW^v, XW^k, XW^q) = \text{softmax}\left(\frac{XW^q(XW^k)^T}{\sqrt{d}}\right)XW^v \quad (92)$$

**Question 51** (RMSNorm (LLaMA's Choice)).

$$\bar{a}_i = \frac{a_i}{\text{RMS}(a)}\gamma = \frac{a_i}{\sqrt{\frac{1}{d} \sum_{j=1}^d a_j^2}}\gamma \quad (93)$$

**Question 52** (Transformer Loss).

$$\min_W \hat{\mathbb{E}} \left[ - \left\langle Y, \log \hat{Y} \right\rangle \right] \quad (94)$$

$Y$  is output sequence, one-hot;

$\hat{Y}$  is the predicted probabilities

**Question 53** (Transformer Implementation). As following.

```
import torch
import torch.nn as nn
import torch.nn.functional as F
import math

class RMSNorm(nn.Module):
    def __init__(self, hidden_dim, eps = 1e-6):
        super().__init__()
        self.eps = eps
        self.weight = nn.Parameter(torch.ones(hidden_dim))

    def forward(self, hidden_state):
        norm = hidden_state.pow(2).mean(-1, keepdim = True)
        output = hidden_state * self.weight * torch.rsqrt(norm + self.eps)
        return output

class MultiHeadAttention(nn.Module):
    def __init__(self, hidden_dim, num_heads):
        super().__init__()
        self.hidden_dim = hidden_dim
        self.num_heads = num_heads
        self.head_dim = hidden_dim // num_heads
        self.q_linear = nn.Linear(hidden_dim, hidden_dim)
        self.k_linear = nn.Linear(hidden_dim, hidden_dim)
        self.v_linear = nn.Linear(hidden_dim, hidden_dim)
        self.o_linear = nn.Linear(hidden_dim, hidden_dim)
        self.norm = RMSNorm(hidden_dim)

    def forward(self, hidden_state, mask, past_kv = None, use_cache = True):
        bs = hidden_state.shape[0]

        residual = hidden_state

        hidden_state = self.norm(hidden_state) # LLaMA style normalization

        q = self.q_linear(hidden_state) # (bs, seq_len, hidden_dim)
```

```

k = self.k_linear(hidden_state) # (bs, seqlen, hidden_dim)
v = self.v_linear(hidden_state) # (bs, seqlen, hidden_dim)

q = q.view(bs, -1, self.num_heads, self.head_dim).transpose(1, 2)
k = k.view(bs, -1, self.num_heads, self.head_dim).transpose(1, 2)
v = v.view(bs, -1, self.num_heads, self.head_dim).transpose(1, 2)
# (bs, nums_head, seqlen, head_dim)

q, k = apply_rope(q, k)

# kv cache
if past_kv is not None:
    past_k, past_v = past_kv
    k = torch.cat([past_k, k], dim = 2)
    v = torch.cat([past_v, v], dim = 2)
new_past_kv = (k, v) if use_cache else None

# compute attention
attention_scores = torch.matmul(q, k.transpose(-1, -2)) / math.sqrt(self.head_dim)
attention_scores += mask * -1e9
attention_scores = F.softmax(attention_scores, dim = -1)
output = torch.matmul(attention_scores, v)

# concat
output = output.transpose(1, 2).contiguous().view(bs, -1, self.hidden_dim)

# o.linear
output = self.o_linear(output)

output += residual

return output, new_past_kv if use_cache else output

```

## 11 Large Language Models

**Question 54** (BERT v.s. GPT). BERT is encoder; GPT is decoder.

- BERT predicts middle words; GPT predicts the next word.
- BERT is **NOT** auto-regressive; GPT is auto-regressive.

**Question 55** (GPT – Generative Pre-Training).

$$\min_{\Theta} \hat{\mathbb{E}} - \log \prod_{j=1}^m \Pr(x_j | x_1, \dots, x_{j-1}; \Theta) \quad (95)$$

**Question 56** (Fine-Tuning Tasks). Supervised fine-tuning tasks:

$$\min_{\Theta} \underbrace{-\hat{\mathbb{E}} \log \Pr(y | X_{1:m}; \Theta)}_{\text{task-aware supervised loss}} - \lambda \underbrace{\hat{\mathbb{E}} \log \prod_{j=1}^m \Pr(x_j | X_{1:j-1}; \Theta)}_{\text{pretraining loss}} \quad (96)$$

**Question 57** (BERT → RoBERTa). Training longer, with bigger batches, over more data and longer sequence. Removing the next sentence prediction objective.

**Question 58** (Sentence-BERT). a twin network architecture that uses BERT to derive sentence embeddings.

**Question 59** (GPT-2). 1.5B parameters.

- 10x larger than GPT-1
- Training method is same as GPT-1.
- Performs on par with BERT on fine-tuning tasks.
- Good at zero-shot learning.
- Open-source.

**Question 60** (GPT-3). 175B parameters.

- 100x larger than GPT-2.

- Training method is same as GPT-2.
- New phenomenon: **in-context learning** (ICL, or few-shot learning) and **chain-of-thoughts** (CoT).

**Question 61** (GPT-3.5 – RLHF). Reinforcement Learning from Human Feedback (RLHF).

- state = prompt
- action = model output
- policy function = LLM
- reward = levels of matching human feedback

Pari-wise comparison loss to train reward model  $r_\theta$ :

$$\mathcal{L}_{\text{pair}}(\theta) = -\frac{1}{\binom{K}{2}} \mathbb{E}_{(x, y_w, y_l)} [\log(\sigma(r_\theta(x, y_w) - r_\theta(x, y_l)))] \quad (97)$$

Proximal Policy Optimization (PPO) to maximize objective:

$$\max_{\Phi} \mathbb{E}_{(x, y)} \left[ \underbrace{r_\theta(x, y)}_{\text{maximize reward}} - \beta \underbrace{\log \left( \frac{\pi_{\Phi}^{\text{RL}}(y|x)}{\pi^{\text{SFT}}(y|x)} \right)}_{\text{model is close to SFT model}} + \gamma \underbrace{\mathbb{E}[\log(\pi_{\Phi}^{\text{RL}}(x))]}_{\text{pretraining loss}} \right] \quad (98)$$

## 12 Speculative Sampling

**Question 62** (Reject Sampling for Check). Check in parallel.

- $r \sim U(0, 1)$ , if  $r < \underbrace{\min \left( 1, \frac{p(t)}{q(t)} \right)}_{\text{accept rate}}$ , next token =  $t$ .
- else: next token =  $t' \sim \underbrace{\text{norm}(\max(0, p - q))}_{\text{residual distribution}}.$

**Question 63** (Proof: Reject Sampling  $\equiv t \sim p$ ).

$$\begin{aligned} \min(p(t), q(t)) + \max(0, p(t) - q(t)) &= \begin{cases} p(t) + 0 & \text{if } p(t) < q(t) \\ q(t) + p(t) - q(t) & \text{if } p(t) \geq q(t) \end{cases} \\ &= p(t) \end{aligned} \quad (99)$$

$$\Rightarrow \sum_t (\min(p(t), q(t)) + \max(0, p(t) - q(t))) = \sum_t p(t) = 1 \quad (100)$$

$$\Rightarrow 1 - \sum_t \min(p(t), q(t)) = \sum_t \max(0, p(t) - q(t)) \quad (101)$$

$$\begin{aligned} \Pr(X = t) &= \Pr(\tilde{X} = t) \Pr(\tilde{X} \text{ accept} | \tilde{X} = t) + \Pr(\tilde{X} \text{ reject}) \Pr(\tilde{X} = t | \tilde{X} \text{ reject}) \\ &= q(t) \cdot \min \left( 1, \frac{p(t)}{q(t)} \right) + (1 - \Pr(\tilde{X} \text{ accept})) \cdot \text{norm}(\max(0, p(t) - q(t))) \\ &= \min(q(t), p(t)) + (1 - \sum_t \min(p(t), q(t))) \cdot \frac{\max(0, p(t) - q(t))}{\sum_t \max(0, p(t) - q(t))} \\ &= \min(q(t), p(t)) + \max(0, p(t) - q(t)) \\ &= p(t) \end{aligned} \quad (102)$$

## 13 Generative Adversarial Networks

**Question 64** (Representation through Push-Forward). Let  $r$  be any continuous distribution on  $\mathbb{R}^h$ . For any distribution  $p$  on  $\mathbb{R}^d$ , there exists push-forward maps  $G : \mathbb{R}^h \rightarrow \mathbb{R}^d$  such that

$$z \sim r \implies G(z) \sim p \quad (103)$$

**Question 65** (Discriminator's Goal). For a fixed generator  $G$ , minimize a log loss over  $D$ :

- If  $x$  is real, minimize  $-\log D(x)$ ;
- If  $x$  is fake, minimize  $-\log(1 - D(x))$ .

$$\min_D -\frac{1}{2}\mathbb{E}_{x \sim p_{\text{data}}} [\log D(x)] - \frac{1}{2}\mathbb{E}_{z \sim r} [\log(1 - D(G(z)))] \quad (104)$$

**Question 66** (Generator's Goal). For a fixed discriminator  $D$ , maximize a log loss over  $G$ :

$$\max_G -\frac{1}{2}\mathbb{E}_{x \sim p_{\text{data}}} [\log D(x)] - \frac{1}{2}\mathbb{E}_{z \sim r} [\log(1 - D(G(z)))] \quad (105)$$

**Question 67** (Solver).

$$\min_G \max_D V(G, D) = \mathbb{E}_{x \sim p_{\text{data}}} [\log D(x)] + \mathbb{E}_{z \sim r} [\log(1 - D(G(z)))] \quad (106)$$

Solved by alternative minimization-maximization:

- G step: Fix  $D$  and update  $G$  by one-step gradient descent
- D step: Fix  $G$  and update  $D$  by one-step gradient descent
- Repeat until the algorithm reaches an approximate equilibrium

**Question 68** (Solution of  $D^*$ ). Let  $p_g(x)$  be the density of  $x$  estimated by the generator  $G$ . For  $G$  fixed, the optimal discriminator  $D$  is  $D_G^*(x) = \frac{p_{\text{data}}(x)}{p_{\text{data}}(x) + p_g(x)}$

**Proof:**

$$\begin{aligned} V(G, D) &:= \mathbb{E}_{x \sim p_{\text{data}}} [\log D(x)] + \mathbb{E}_{z \sim r} [\log(1 - D(G(z)))] \\ &= \int \log D(x) p_{\text{data}}(x) dx + \int_z \log(1 - D(G(z))) p_z(z) dz \\ &= \int \underbrace{\log D(x) p_{\text{data}}(x) + p_g(x) \log(1 - D(x))}_{f(D(x))} dx \end{aligned} \quad (107)$$

For any fixed  $x$ , taking derivative = 0:

$$\begin{aligned} f'(D(x)) &= \frac{p_{\text{data}}(x)}{D(x)} - \frac{p_g(x)}{1 - D(x)} = 0 \\ D_G^*(x) &= \frac{p_{\text{data}}(x)}{p_{\text{data}}(x) + p_g(x)} \end{aligned} \quad (108)$$

**Question 69** (Solution of  $G^*$ ). The global minimum of  $\min_G \max_D V(G, D)$  is achieved if and only if  $p_g = p_{\text{data}}$ . The optimal objective value is  $-\log 4$ .

**Proof:**

$$\begin{aligned} V(G, D_G^*) &= \mathbb{E}_{x \sim p_{\text{data}}} [\log D_G^*(x)] + \mathbb{E}_{z \sim r} [\log(1 - D_G^*(G(z)))] \\ &= \mathbb{E}_{x \sim p_{\text{data}}} [\log D_G^*(x)] + \mathbb{E}_{x \sim p_g} [\log(1 - D_G^*(x))] \\ &= \mathbb{E}_{x \sim p_{\text{data}}} \left[ \log \frac{p_{\text{data}}(x)}{p_{\text{data}}(x) + p_g(x)} \right] + \mathbb{E}_{x \sim p_g} \left[ \log \frac{p_g(x)}{p_{\text{data}}(x) + p_g(x)} \right] \end{aligned} \quad (109)$$

By definition of KL divergence  $\text{KL}(P\|Q) = \mathbb{E}_{x \sim P} \left[ \log \frac{p(x)}{q(x)} \right]$ , we have:

$$\begin{aligned} V(G, D_G^*) &= \mathbb{E}_{x \sim p_{\text{data}}} \left[ \log \frac{p_{\text{data}}(x)}{p_{\text{data}}(x) + p_g(x)} \right] + \mathbb{E}_{x \sim p_g} \left[ \log \frac{p_g(x)}{p_{\text{data}}(x) + p_g(x)} \right] \\ &= -\log 4 + \text{KL} \left( p_{\text{data}} \parallel \frac{p_{\text{data}} + p_g}{2} \right) + \text{KL} \left( p_g \parallel \frac{p_{\text{data}} + p_g}{2} \right) \\ &= -\log 4 + 2 \cdot \text{JSD}(p_{\text{data}} \parallel p_g) \\ &\geq -\log 4 \end{aligned} \quad (110)$$

The equality holds if and only if  $p_{\text{data}} = p_g$ .

## 14 Adversarial Attacks

**Question 70** (Principle of Generating Adversarial Attacks).

$$\max_{\|x_{\text{adv}} - x\|_{\infty} \leq \epsilon} \mathcal{L}(C(x_{\text{adv}}), y) \quad (111)$$

where  $C$  is the composition of  $h$  and  $f$ .

**Question 71** (Different Solvers). to optimize the adversarial attack.

- **Zero-Order Solvers** (only access to the output of NN)
  - Black-box attack
- **First-Order Solvers** (access to the gradient of NN)
  - White-box attack
  - Fast Gradient Sign Method (FGSM), BIM, PGD, CW attack, ...
- **Second-Order Solvers** (access to the Hessian matrix)
  - White-box attack
  - L-BFGS attack

**Question 72** (Holder Inequality). For any  $p, q \geq 1$  such that  $\frac{1}{p} + \frac{1}{q} = 1$ ,

$$\|x\|_p \cdot \|y\|_q \geq \langle x, y \rangle \quad (112)$$

where  $\langle x, y \rangle$  is the inner product.

$\|\cdot\|_p$  and  $\|\cdot\|_q$  are also known as dual norms.

- $\|\cdot\|_2$  is self-dual.
- $\|\cdot\|_{\infty}$  and  $\|\cdot\|_1$  are dual norms.

**Question 73** (FGSM – Fast Gradient Sign Method). White-box and non-targeted (maximize the loss w.r.t. the true label). Do linear expansion at  $x$ :

$$\mathcal{L}(C(x + \delta), y) \approx \underbrace{\mathcal{L}(C(x), y)}_{\text{constant}} + \nabla_x \mathcal{L}(C(x), y) \cdot \delta \quad (113)$$

The problem reduces to:

$$\max_{\|\delta\|_{\infty} \leq \epsilon} \nabla_x \mathcal{L}(C(x), y) \cdot \delta \quad (114)$$

Because of holder inequality (112), we have:

$$\nabla_x \mathcal{L}(C(x), y) \cdot \delta \leq \|\delta\|_{\infty} \cdot \|\nabla_x \mathcal{L}(C(x), y)\|_1 \leq \epsilon \cdot \|\nabla_x \mathcal{L}(C(x), y)\|_1 \quad (115)$$

Thus, the adversarial example is generated by:

$$x_{\text{adv}} = x + \epsilon \cdot \text{sign}(\nabla_x \mathcal{L}(C(x), y)) \quad (116)$$

where  $\epsilon$  is the perturbation size.

**Question 74** (BIM – Basic Iterative Method). BIM is an iterative version of FGSM.

- Initialize  $x^{(0)} = x$ .
- For  $k = 1, 2, \dots, K$ :

$$x^{(k)} = x^{(k-1)} + \gamma \cdot \text{sign}(\nabla_x \mathcal{L}(C(x^{(k-1)}), y)) \quad (117)$$

Issues:

- By repeating, the perturbation size  $\epsilon$  will become larger.
- For a pre-defined  $\epsilon$ ,  $x^{(k)}$  may not satisfy  $\|x^{(k)} - x\|_{\infty} \leq \epsilon$ .

**Question 75** (PGD – Projected Gradient Descent). To resolve the issue of BIM, PGD involves a truncation operation:

- Initialize  $x^{(0)} = x + \delta$ , where  $\delta \in (-\epsilon, \epsilon)$ .
- For  $k = 1, 2, \dots, K$ :

$$x^{(k)} = \text{clip}_{(-\epsilon, \epsilon)}(x^{(k-1)} + \gamma \cdot \text{sign}(\nabla_x \mathcal{L}(C(x^{(k-1)}), y))) \quad (118)$$



where  $\text{clip}_{(-\epsilon, \epsilon)}(x)$  projects  $x$  back to the  $\ell_\infty$  ball of radius  $\epsilon$  around  $x$ .

**Question 76** (Targeted PGD Attack). Objective:

- **Untargeted:**

$$\max_{\|\delta\|_\infty \leq \epsilon} \mathcal{L}(C(x + \delta), y_{\text{true}}) \quad (119)$$

- **Targeted:**

$$\min_{\|\delta\|_\infty \leq \epsilon} \mathcal{L}(C(x + \delta), y_{\text{target}}) \quad (120)$$

## 15 Adversarial Robustness

**Question 77** (Defense Mechanisms). Categorized into two types:

- **Gradient Masking:** hide the gradients and make first-order attacks fail
  - **Shattered Gradients** By applying a non-smooth or non-differentiable preprocessor  $g$  to the inputs, and then training a DNN model  $f$  on the preprocessed inputs  $g(x)$ .
  - **Stochastic/Randomized Gradients** Apply some form of randomization of the DNN model. E.g. train a set of classifiers and during the testing phase randomly select one classifier to predict the labels.
- **Adversarial Training:**

$$\min_C \mathbb{E}_{(x,y) \sim \mathcal{D}^n} \left[ \max_{\|\delta\|_\infty \leq \epsilon} \mathcal{L}(C(x + \delta), y) \right] \quad (121)$$

**Question 78** (Trade-Off between Natural and Robust Error).

$$\min_f R_{\text{nat}}(f) + R_{\text{rob}}(f)/\lambda \quad (122)$$

$$\begin{aligned} R_{\text{nat}}(f) &:= \Pr_{x,y \sim \mathcal{D}} \{f(x)y \leq 0\} \\ &= \mathbb{E}_{(x,y) \sim \mathcal{D}} [\mathbb{I}(f(x)y \leq 0)] \end{aligned} \quad (123)$$

$$\begin{aligned} R_{\text{rob}}(f) &:= \Pr_{x,y \sim \mathcal{D}} \{\exists \delta \in B_\epsilon(x) \text{ s.t. } f(x + \delta)y \leq 0\} \\ &= \mathbb{E}_{(x,y) \sim \mathcal{D}} \left[ \max_{\|\delta\|_\infty \leq \epsilon} \mathbb{I}(f(x + \delta)y \leq 0) \right] \end{aligned} \quad (124)$$

Approximate by a differentiable surrogate loss  $\Phi$ :

$$R_{\text{nat}}(f) \approx \mathbb{E}_{(x,y) \sim \mathcal{D}} [\Phi(f(x)y)] \quad (125)$$

**Question 79** (TRADES).

$$\min_f \left[ \underbrace{\underbrace{\mathbb{E}_{(x,y) \sim \mathcal{D}} [\Phi(f(x)y)]}_{\text{minimize diff btw } f(x) \text{ and } y \text{ for accuracy}} + \underbrace{\mathbb{E}_{(x,y) \sim \mathcal{D}} \left[ \max_{\|\delta\|_\infty \leq \epsilon} \Phi(f(x + \delta)f(x)) \right]}_{\text{minimize diff btw } f(x) \text{ and } f(x+\delta) \text{ for robustness}}}_{\text{TRADES Loss}} / \lambda \right] \quad (126)$$

## 16 Self-Supervised Learning

**Question 80** (Contrastive Learning). Loss:

$$\max_{\Theta} \Pr_1 = \frac{\exp(z_1)}{\sum_j \exp(z_j)} \quad (127)$$