

Mathematical Formulations for Synthetic Data Generation with Controlled Leakage

1. Feature Generation (\mathbf{X})

We generate the feature matrix $\mathbf{X} \in \mathbb{R}^{n \times d}$ by sampling n observations from a multivariate normal distribution:

$$\mathbf{x}_i \sim \mathcal{N}(\boldsymbol{\mu}_x, \boldsymbol{\Sigma}_x), \quad i = 1, \dots, n$$

- $\boldsymbol{\mu}_x \in \mathbb{R}^d$: Mean vector of the features.
- $\boldsymbol{\Sigma}_x \in \mathbb{R}^{d \times d}$: Covariance matrix of the features.

2. Ground Truth Concepts (\mathbf{c}_i)

We construct the ground truth concepts $\mathbf{c}_i \in \{0, 1\}^k$ as follows:

2.1. Constructing Matrix \mathbf{A}

Matrix $\mathbf{A} \in \mathbb{R}^{k \times d}$ is designed to project the first b features into the concept space:

$$\mathbf{A} = [\mathbf{R}_A \mid \mathbf{0}_{k \times (d-b)}]$$

- $\mathbf{R}_A \in \mathbb{R}^{k \times b}$: Random projection matrix with entries:

$$(\mathbf{R}_A)_{jp} \stackrel{\text{iid}}{\sim} \mathcal{N}(0, 1), \quad j = 1, \dots, k; \quad p = 1, \dots, b$$

- $\mathbf{0}_{k \times (d-b)}$: Zero matrix to exclude the remaining features.

2.2. Computing Success Probabilities π_i

We compute the logits for the ground truth concepts:

$$\boldsymbol{\eta}_i = \mathbf{A}\mathbf{x}_i + \boldsymbol{\epsilon}_c$$

- $\boldsymbol{\epsilon}_c \sim \mathcal{N}(\mathbf{0}, \boldsymbol{\Sigma}_c)$: Noise vector with covariance $\boldsymbol{\Sigma}_c \in \mathbb{R}^{k \times k}$.

The success probabilities are then obtained via the sigmoid function:

$$\pi_i = \sigma(\boldsymbol{\eta}_i)$$

- Sigmoid function $\sigma(z)$:

$$\sigma(z) = \frac{1}{1 + e^{-z}}$$

2.3. Sampling Ground Truth Concepts \mathbf{c}_i

Each concept c_{ij} is sampled from a Bernoulli distribution:

$$c_{ij} \sim \text{Bernoulli}(\pi_{ij}), \quad j = 1, \dots, k$$

3. Leakage Term (\mathbf{l}_i)

We construct the leakage term $\mathbf{l}_i \in \mathbb{R}^k$ as follows:

3.1. Constructing Matrix \mathbf{B}

Matrix $\mathbf{B} \in \mathbb{R}^{k \times d}$ projects selected features into the concept space:

$$\mathbf{B} = [\mathbf{0}_{k \times b} \mid \mathbf{R}_B \mid \mathbf{0}_{k \times l}]$$

- $\mathbf{0}_{k \times b}$: Zero matrix to exclude the first b features.
- $\mathbf{R}_B \in \mathbb{R}^{k \times (d-b-l)}$: Random projection matrix with entries:

$$(\mathbf{R}_B)_{jq} \stackrel{\text{iid}}{\sim} \mathcal{N}(0, 1), \quad j = 1, \dots, k; \quad q = 1, \dots, d - b - l$$

- $\mathbf{0}_{k \times l}$: Zero matrix to exclude the last l features.

3.2. Computing Leakage Term l_i

$$l_i = \mathbf{B} \mathbf{x}_i$$

4. Estimated Concepts ($\hat{\mathbf{c}}_i$)

We compute the estimated concepts $\hat{\mathbf{c}}_i \in [0, 1]^k$ as:

$$\hat{\mathbf{c}}_i = \sigma(\mathbf{A} \mathbf{x}_i + \mathbf{l}_i + \boldsymbol{\epsilon}_{\hat{\mathbf{c}}})$$

- $\boldsymbol{\epsilon}_{\hat{\mathbf{c}}} \sim \mathcal{N}(\mathbf{0}, \boldsymbol{\Sigma}_{\hat{\mathbf{c}}})$: Noise vector with covariance $\boldsymbol{\Sigma}_{\hat{\mathbf{c}}} \in \mathbb{R}^{k \times k}$.
- σ : Sigmoid function as defined earlier.

5. Target Labels (y_i)

We generate the target labels $y_i \in \{1, \dots, J\}$ as follows:

5.1. Defining Function f

If no custom function f is provided, we define $f : \mathbb{R}^k \times \mathbb{R}^k \rightarrow \mathbb{R}^J$ using a simple Multi-Layer Perceptron (MLP):

- **Inputs:** Concatenate \mathbf{c}_i and \mathbf{l}_i :

$$\mathbf{u}_i = \begin{bmatrix} \mathbf{c}_i \\ \mathbf{l}_i \end{bmatrix} \in \mathbb{R}^{2k}$$

- **First Layer:**

$$\mathbf{h}_i = \phi(\mathbf{W}_1 \mathbf{u}_i + \mathbf{b}_1) \in \mathbb{R}^h$$

- $\mathbf{W}_1 \in \mathbb{R}^{h \times 2k}$: Weight matrix with entries $\sim \mathcal{N}(0, 1)$.
- $\mathbf{b}_1 \in \mathbb{R}^h$: Bias vector initialized to zeros.
- $\phi(z)$: ReLU activation function:

$$\phi(z) = \max(0, z)$$

- **Output Layer:**

$$\mathbf{o}_i = \mathbf{W}_2 \mathbf{h}_i + \mathbf{b}_2 \in \mathbb{R}^J$$

- $\mathbf{W}_2 \in \mathbb{R}^{J \times h}$: Weight matrix with entries $\sim \mathcal{N}(0, 1)$.
- $\mathbf{b}_2 \in \mathbb{R}^J$: Bias vector initialized to zeros.

5.2. Computing Target Probabilities \mathbf{p}_i

We compute the logits for the target probabilities:

$$\mathbf{z}_i = f(\mathbf{c}_i, \mathbf{l}_i) + \boldsymbol{\epsilon}_y$$

- $\boldsymbol{\epsilon}_y \sim \mathcal{N}(\mathbf{0}, \boldsymbol{\Sigma}_y)$: Noise vector with covariance $\boldsymbol{\Sigma}_y \in \mathbb{R}^{J \times J}$.

We then apply the softmax function to obtain the probabilities:

$$\mathbf{p}_i = \text{softmax}(\mathbf{z}_i)$$

- Softmax function:

$$\text{softmax}(\mathbf{z}_i)_j = \frac{\exp(z_{ij})}{\sum_{k=1}^J \exp(z_{ik})}, \quad j = 1, \dots, J$$

5.3. Sampling Target Labels y_i

We sample y_i from a categorical distribution based on \mathbf{p}_i :

$$y_i \sim \text{Categorical}(\mathbf{p}_i)$$

This results in $y_i \in \{1, 2, \dots, J\}$.

6. Summary of Variables and Parameters

- n : Number of observations.
- d : Dimensionality of features.
- k : Number of concepts.
- J : Number of target classes.
- b : Number of features used in the ground truth concepts.
- l : Number of features excluded from leakage.
- $\boldsymbol{\mu}_x \in \mathbb{R}^d$: Mean vector of features.
- $\boldsymbol{\Sigma}_x \in \mathbb{R}^{d \times d}$: Covariance matrix of features.
- $\boldsymbol{\Sigma}_c \in \mathbb{R}^{k \times k}$: Covariance matrix of noise in ground truth concepts.
- $\boldsymbol{\Sigma}_{\hat{c}} \in \mathbb{R}^{k \times k}$: Covariance matrix of noise in estimated concepts.
- $\boldsymbol{\Sigma}_y \in \mathbb{R}^{J \times J}$: Covariance matrix of noise in target logits.
- f : Function mapping concepts and leakage to target logits.