Systematic Trading Strategies with Machine Learning Algorithms

Mock Exam

June 5, 2025

The purpose of this exam is to evaluate your understanding of systematic trading strategies using machine learning techniques and your ability to apply the concepts covered in this course. The exam consists of two independent sections: primary model development using tree-based methods with trend scanning, and meta-model implementation using Variable Selection Networks for trade filtering. **Instructions:**

- Read all questions carefully before answering
- Show all mathematical derivations and algorithmic steps to receive full credit
- Answer questions in order as some concepts build upon previous ones
- Calculators are permitted for numerical computations

Exam Structure: (Total marks: 100)

- Section 1: Creating a Primary Model using Tree-Based Models and Trend Scanning (60 marks)
 - Exercise 1: Labeling Methodology using Trend Scanning (3 questions)
 - Exercise 2: Tree-Based Models for the Primary Model (5 questions)
 - Exercise 3: Cluster-Based Feature Importance Analysis (7 questions)
- Section 2: Meta-model using Variable Selection Networks (40 marks)
 - Exercise 1: Meta-labeling and Threshold-Based Trade Filtering (3 questions)
 - Exercise 2: Variable Selection Networks (7 questions)

Note: Each question is worth 4 marks. Time allowed: 2 hours

Section 1: Creating a Primary Model using tree based models and Trend Scanning

Introduction

Financial markets exhibit complex temporal patterns that can be systematically identified and exploited through machine learning techniques. Effective trading strategies require both accurate signal generation and intelligent trade filtering to maximize profitability while minimizing risk.

In this section, we explore a comprehensive approach to systematic trading that consists of three major components:

- 1. Labeling Methodology: We begin by implementing a trend detection framework that identifies directional movements in financial time series using trend scanning techniques.
- 2. Tree-Based Models for the Primary Model: We develop supervised learning models using Random Forest algorithms to map features to trend labels.
- 3. Cluster-Based Feature Importance Analysis: We implement advanced feature importance analysis that groups correlated features and provides robust importance estimates.

The combination of these three components creates a comprehensive framework for systematic trading strategy development.

Exercise 1 (Labeling Methodology).

Trend scanning is a data-driven labeling method that identifies significant directional movements in financial time series by fitting linear regression models over multiple forward-looking horizons. Unlike fixed-horizon approaches, trend scanning adapts to market dynamics by selecting the optimal look-ahead period that exhibits the strongest statistical trend.

In this exercise, we implement the trend scanning methodology as covered in the course. The approach has two main hyperparameters that define the search space for optimal trend detection:

- H_{\min} : Minimum look-ahead horizon
- H_{max} : Maximum look-ahead horizon

We want to label time steps from 1 to T. For a specific time step t, we evaluate all horizons $H \in \{H_{\min}, H_{\min} + 1, \ldots, H_{\max}\}$ to find the one that maximizes the statistical significance of the observed trend, as illustrated in Figure 1.

The trend scanning approach works by fitting a linear regression model for each possible horizon H and selecting the horizon that produces the most statistically significant trend. Let $\{P_s\}_{s=1}^T$ be a sequence of asset prices.

For a specific time step t and horizon H, we have prices from P_t to P_{t+H} . We fit the following linear model to the forward-looking price series:

$$P_{t+h} = \beta_0 + \beta_1 h + \varepsilon_{t+h}, \quad h = 0, 1, \dots, H$$

where β_1 represents the trend coefficient (slope) and ε_{t+h} is the error term. The statistical significance of the trend is measured by the t-statistic:

$$t_{\hat{\beta}_1} = \frac{\hat{\beta}_1}{\hat{\sigma}_{\hat{\beta}_1}}$$

where $\hat{\beta}_1$ is the estimated slope coefficient and $\hat{\sigma}_{\hat{\beta}_1}$ is its standard error.



Figure 1: Trend scanning example: evaluating multiple horizons to find optimal trend

As illustrated in the figure above, for each time step, we compute t-statistics across different horizons and select the one that indicates the most significant trend.

Q1. Based on the trend scanning methodology:

• For each horizon H, we compute a t-statistic $t_{\hat{\beta}_1}(H)$. How do we select the optimal horizon H^* among all the candidates $H \in \{H_{\min}, \ldots, H_{\max}\}$?

- Once we have identified the optimal horizon H^* , how do we generate the final trend label $y_t \in \{-1, +1\}$ for this time step?
- **Q2.** Complete the trend scanning algorithm for labeling a single time step t:

Algorithm 1 Trend Scanning for Single Time Step	
Require: Price series $\{P_s\}_{s=1}^T$, time step t, H_{\min} , H_{\max}	
Ensure: Trend label y_t and optimal horizon H_t^*	
1: Initialize $t_{\text{max}} = 0, H_t^* = H_{\text{min}}$	
2: for $H = H_{\min}$ to H_{\max} do	
3: [BLANK 1]	\triangleright Fill in the blank
4: [BLANK 2]	\triangleright Fill in the blank
5: Calculate t-statistic: $t_{\hat{\beta}_1} = \frac{\hat{\beta}_1}{\hat{\sigma}_{\hat{\beta}_1}}$	
6: if $ t_{\hat{\beta}_1} > t_{\max}$ then	
7: $t_{\text{max}} = [\text{BLANK 3}]$	\triangleright Fill in the blank
8: $H_t^* = H$	
9: end if	
10: end for	
11: Generate label: $y_t = [BLANK 4]$	\triangleright Fill in the blank
12: return y_t , H_t^*	

With this algorithm, we have generated y_t associated with time step t.

The trend scanning technique can also be used backward-looking to generate features at time step t.

Q3. What modifications should be made to the aforementioned algorithm to generate a new feature that captures the strength of historical trends ending at time t?

This process of labeling and feature generation (along with many other indicators) are used over the entire dataset to generate training data $\{(\mathbf{X}_t, y_t)\}$ for $t \in \{1, \ldots, T\}$. The next step is to train tree-based models to map the features to the targets.

Exercise 2 (Tree-Based Models for the Primary Model).

We will fit a Random Forest algorithm on the training data $\{(\mathbf{X}_t, y_t)\}_{t=1}^T$, where $y_t \in \{-1, +1\}$ represents the trend labels generated from the trend scanning methodology. The Random Forest consists of multiple decision trees, each trained on a bootstrap sample of the data.

Q1. Cite and briefly describe three key hyperparameters of a decision tree algorithm that should be tuned to optimize performance.

Recall that for binary classification with a set S containing classes $\{-1, +1\}$, the entropy is defined as:

$$H(S) = -p_{+1}\log_2(p_{+1}) - p_{-1}\log_2(p_{-1})$$

where $p_{+1} = \frac{|\{s \in S: class(s) = +1\}|}{|S|}$ and $p_{-1} = \frac{|\{s \in S: class(s) = -1\}|}{|S|}$ are the proportions of samples in S belonging to classes +1 and -1 respectively.

Consider a parent node containing 100 samples with equal distribution (50 samples with label +1 and 50 samples with label -1). A categorical feature splits this node based on feature values 0 and 1, resulting in:

Node	Feature Value	Label $+1$	Label -1
Parent	-	50	50
Left Child	0	50	0
Right Child	1	0	50

- **Q2.** Calculate the Information Gain (IG) for the split described in the example above using entropy as the impurity measure. Calculate the entropy before the split, the weighted entropy after the split, and the information gain.
- **Q3.** Complete the decision tree learning algorithm for our trend classification problem:

Algorithm 2 Decision Tree Learning Algorithm

Require: Training data $\{(\mathbf{X}_t, y_t)\}_{t=1}^T$ where $y_t \in \{-1, +1\}$, stopping criteria **Ensure:** Decision tree T

- 1: Initialize tree with single root node containing all data
- 2: while nodes can be split and stopping criteria not met do
- for each leaf node with region \mathcal{R} do 3:
- Find (j^*, τ^*) that maximizes: **[BLANK 1]** \triangleright Fill in the blank 4:

5:
$$IG(j,\tau) = I(\mathcal{R}) - \frac{|\mathcal{R}_L|}{|\mathcal{R}|}I(\mathcal{R}_L) - \frac{|\mathcal{R}_R|}{|\mathcal{R}|}I(\mathcal{R}_R)$$

- Where $\mathcal{R}_L = \{ \mathbf{X}_t \in \mathcal{R} : \mathbf{X}_{t,j} \leq \tau \}$ and $\mathcal{R}_R = \{ \mathbf{X}_t \in \mathcal{R} : \mathbf{X}_{t,j} > \tau \}$ 6:
- Create left child node with samples: **[BLANK 2]** \triangleright Fill in the blank 7:
- Create right child node with samples: **[BLANK 3]** \triangleright Fill in the blank 8:
- end for 9:
- 10: end while

```
11: Assign prediction to each leaf node: [BLANK 4]
                                                                   \triangleright Fill in the blank
12: return T
```





Figure 2: Random Forest algorithm overview

- **Q4.** What are the two main sources of randomness introduced by Random Forest to avoid overfitting compared to a single decision tree?
- **Q5.** We would like to evaluate the Random Forest algorithm. Cite and briefly describe two evaluation metrics that don't depend on a specific threshold to turn the probability of positive label into hard predictions (-1/+1).

Exercise 3 (Cluster-Based Feature Importance Analysis).

Traditional feature importance methods like Mean Decrease Impurity (MDI) or Permutation Feature Importance (PFI) can be misleading when features are correlated or noisy. To address this limitation, we develop a cluster-based approach that groups similar features and provides more robust importance estimates.

Let $\mathbf{X} \in \mathbb{R}^{T \times D}$ be our feature matrix where T is the number of time steps We will use Random Forest to avoid overfitting, as illustrated in Figure 2. and D is the number of features. Each column $\mathbf{X}_{i,j}$ represents feature j across all time steps.

1. Step 1: Initial Feature Importance with MDI

We start by computing the Mean Decrease Impurity (MDI) for each feature using our trained Random Forest. Let MDI_j denote the importance score for feature j, where:

$$\sum_{j=1}^{D} \mathrm{MDI}_{j} = 1 \quad \mathrm{and} \quad \mathrm{MDI}_{j} \ge 0 \quad \forall j$$

However, MDI has known limitations that we must address.

Q1. Explain why a noisy feature can have $MDI_j > 0$ even if it doesn't bring any information to the prediction task.

2. Step 2: Out-of-Sample Feature Importance

To mitigate MDI's bias toward noisy features, we employ Permutation Feature Importance (PFI), which measures the decrease in model performance when a feature's values are randomly shuffled.

Let *m* be our trained model, $\mathcal{D} = \{(\mathbf{X}_i, y_i)\}_{i=1}^T$ be our validation dataset, and $s(\mathcal{D})$ be a performance metric (e.g., accuracy). For feature *j*, we create multiple corrupted datasets $\tilde{\mathcal{D}}_{k,j}$ for $k \in \{1, \ldots, K\}$ by randomly permuting the *j*-th feature column across different random seeds.

Q2. Complete the Permutation Feature Importance algorithm:

Algorithm 3 Permutation Feature Importance (PFI)

Require: Fitted model m, validation data \mathcal{D} , repetitions K**Ensure:** Feature importance scores $\{PFI_j\}_{j=1}^{D}$ with standard deviations $\{\sigma_j\}_{j=1}^{D}$ 1: Compute reference score $s_0 = s(m, \mathcal{D})$ 2: for each feature $j \in \{1, \ldots, D\}$ do Initialize array $scores_i$ of length K 3: for each repetition $k \in \{1, \ldots, K\}$ do 4: [BLANK 1] \triangleright Fill in the blank 5: Compute score $s_{k,j} = s(m, \tilde{\mathcal{D}}_{k,j})$ Store in array: $scores_j[k] = [\mathbf{BLANK 2}]$ 6: \triangleright Fill in the blank 7: 8: end for Compute $PFI_j = mean(scores_j)$ and $\sigma_j = std(scores_j)$ 9: 10: **end for** 11: return $\{PFI_i\}_{i=1}^D$ and $\{\sigma_i\}_{i=1}^D$

3. Step 3: Addressing Feature Correlation

In our trend scanning framework, many generated features are highly correlated (e.g., different moving averages, overlapping technical indicators). This correlation creates challenges for feature importance analysis.

Let $\mathbf{C} \in \mathbb{R}^{D \times D}$ be the Spearman correlation matrix of our features, where $C_{i,j}$ represents the correlation between features i and j.

Q3. Explain how feature correlation affects permutation feature importance analysis through the substitution effect. When one important feature is permuted, how can correlated features mask the performance drop? Why does this motivate cluster-based feature importance analysis?

4. Step 4: Distance Matrix Construction

To cluster correlated features, we transform the correlation matrix into a distance matrix. We define the distance between features i and j as:

$$d_{i,j} = 1 - |C_{i,j}|$$

This ensures that highly correlated features (positive or negative correlation) are close in distance space, while uncorrelated features are far apart.

5. Step 5: Dimensionality Reduction and K-means Clustering

We apply Principal Component Analysis (PCA) to the distance matrix to reduce dimensionality. Let $\mathbf{Z} \in \mathbb{R}^{D \times d}$ be the PCA-transformed data where $d \ll D$ is chosen to retain at least 95% of the variance.

Next, we determine the optimal number of clusters using silhouette analysis.

For each data point *i* in cluster C_k , the silhouette coefficient is:

$$s(i) = \frac{b(i) - a(i)}{\max(a(i), b(i))}$$

where a(i) is the average intra-cluster distance, and b(i) is the average distance to the nearest different cluster.

Q4. Complete the algorithm to find the optimal number of clusters:

Algorithm 4 Optimal Cluster Selection using Silhou	lette Score	Algorithm 5 EM Algorithm for GMM	
Require: PCA-reduced data $\mathbf{Z} \in \mathbb{R}^{D \times d}$, maximum c	lusters $K_{\rm max}$	Require: PCA data $\{\mathbf{z}_i\}_{i=1}^D$, number of components K	
Ensure: Optimal number of clusters K^*		Ensure: Parameters $\boldsymbol{\theta} = (\boldsymbol{\pi}, \boldsymbol{\mu}, \boldsymbol{\Sigma})$	
1: Initialize <i>silhouette</i> $scores = []$		1: Initialize parameters $\boldsymbol{\theta}^{(0)}$	
2: for $K = 2$ to $K_{\text{max}} \mathbf{d} \mathbf{o}$		2: $t \leftarrow 0$	
3: [BLANK 1]	\triangleright Fill in the blank	3: while not converged do	
4: [BLANK 2]	\triangleright Fill in the blank	K 4: E-step: Compute responsibilities $\tau_{i,k}^{(t)}$ for all $i \in \{1, \dots, D\}$, k	$k \in$
5: Append score to silhouette_scores		$\{1,\ldots,K\}$	
6: end for		5: M-step: Update parameters:	
7: $K^* = [\mathbf{BLANK 3}]$	\triangleright Fill in the blank	^K 6: $\pi_{1}^{(t+1)} = [\text{BLANK 1}] \qquad \triangleright \text{ Fill in the bl}$	lank
8: return K*		- 7: μ_1^{κ} = [BLANK 2] \triangleright Fill in the bl	lank

8:

9:

6. Step 6: Gaussian Mixture Model for Soft Clustering

While K-means provides hard cluster assignments, we use Gaussian Mixture Models (GMM) to obtain probabilistic cluster memberships. This allows features to belong to multiple clusters with different probabilities.

Let $\boldsymbol{\theta} = (\boldsymbol{\pi}, \boldsymbol{\mu}, \boldsymbol{\Sigma})$ be the GMM parameters where:

- $\boldsymbol{\pi} = (\pi_1, \dots, \pi_K)$ are mixture weights with $\sum_{k=1}^K \pi_k = 1$
- $\boldsymbol{\mu} = (\boldsymbol{\mu}_1, \dots, \boldsymbol{\mu}_K)$ are component means where $\boldsymbol{\mu}_k \in \mathbb{R}^d$
- $\Sigma = (\Sigma_1, \dots, \Sigma_K)$ are covariance matrices where $\Sigma_k \in \mathbb{R}^{d \times d}$
- **Q5.** For a GMM with K components modeling d-dimensional data, calculate the total number of parameters. Consider:
 - Mixture weights: π_1, \ldots, π_K .
 - Component means: μ_1, \ldots, μ_K .
 - Full covariance matrices: $\Sigma_1, \ldots, \Sigma_K$. (symmetric positive definite)

The GMM is trained using the Expectation-Maximization (EM) algorithm, which alternates between computing posterior probabilities (E-step) and updating parameters (M-step).

Q6. Complete the EM algorithm for GMM parameter estimation. The posterior probabilities (responsibilities) are given by:

$$\tau_{i,k} = \frac{\pi_k \mathcal{N}(\mathbf{z}_i | \boldsymbol{\mu}_k, \boldsymbol{\Sigma}_k)}{\sum_{j=1}^K \pi_j \mathcal{N}(\mathbf{z}_i | \boldsymbol{\mu}_j, \boldsymbol{\Sigma}_j)}$$

where \mathbf{z}_i is the PCA representation of feature *i*.

7. Step 7: Hard Cluster Assignment

 $\Sigma_{k}^{(t+1)} = [\text{BLANK 3}]$

 $t \leftarrow t + 1$

10: end while 11: return $\theta^{(t)}$

> After training the GMM, we obtain a probability matrix $\mathbf{P} \in \mathbb{R}^{D \times K}$ where $P_{i,k}$ represents the probability that feature *i* belongs to cluster *k*.

 \triangleright Fill in the blank

Q7. Given the soft probability matrix **P**, derive the hard cluster assignment for each feature.

Section 2: Meta-model using Variable Selection Networks

Introduction

While the primary Random Forest model from Section 1 successfully generates trading signals based on trend scanning labels, not every signal translates into profitable trades. Market conditions, volatility regimes, and structural breaks can significantly impact the reliability of these predictions. To enhance the overall trading strategy performance, we implement a meta-modeling framework that learns to identify when the primary model's signals are most likely to succeed.

This section introduces a comprehensive approach to trade filtering through meta-learning, consisting of three critical components:

1. Meta-labeling Framework: We develop a systematic approach to generate binary labels that indicate trade profitability using barrier-based exit strategies adapted to market volatility.

- 2. Feature Engineering for Meta-learning: We construct a distinct feature set that captures market microstructure, regime characteristics, and environmental conditions rather than price patterns used by the primary model.
- 3. Variable Selection Networks: We implement advanced neural network architectures specifically designed for feature selection and binary classification in high-dimensional financial data.

The meta-model operates as an intelligent filter that evaluates market conditions and predicts the probability that the primary model's signal will generate alpha. This probabilistic framework allows for flexible threshold-based trade filtering, enabling systematic control over the precision-recall trade-off.

Exercise 4 (Meta-labeling and Threshold-Based Trade Filtering).

The primary Random Forest model generates trend signals $s_t \in \{-1, +1\}$ based on the trend scanning methodology. However, these signals exhibit varying success rates across different market regimes. Our meta-model M_{meta} learns to predict the probability $p_t = P(\text{Profitable Trade}|\text{Market Conditions}_t)$ that a trade initiated at time t will be profitable.

Architecture Overview:

- The **Primary Model** processes price-based features to generate directional signals $s_t \in \{-1, +1\}$.
- The Meta-model analyzes market environment features to predict signal reliability p_t .
- The **Final Decision** trades only if $p_t > \tau$ for threshold τ .

The meta-model framework, illustrated in Figure 3, processes a distinct feature space to avoid information leakage and redundancy with the primary model.



Figure 3: Metamodel Framework

1. Meta-label Generation using Adaptive Barriers

To train our meta-model, we require binary labels $y_t^{\text{meta}} \in \{0, 1\}$ indicating trade profitability. We employ an adaptive triple-barrier method that dynamically adjusts barrier widths based on realized volatility.

For a trade initiated at time t with signal s_t , we define:

- Entry price: P_t
- Volatility estimate: $\hat{\sigma}_t$ (computed from historical price returns)
- Upper barrier: $P_t \times (1 + \alpha \times \hat{\sigma}_t)$
- Lower barrier: $P_t \times (1 \alpha \times \hat{\sigma}_t)$
- Time barrier: $t + H_{\text{max}}$ (maximum holding period)

The volatility scaling factor α ensures that profit targets and stop-losses adapt to market conditions.

Q1. Consider the adaptive barrier method for meta-label generation at time t.

Algorithm 6 Adaptive Meta-label Generation for Single Time Step

Require: Price series $\{P_s\}_{s=t}^{t+H_{\text{max}}}$, signal s_t , volatility estimate $\hat{\sigma}_t$, scaling factor α , max holding period H_{max} **Ensure:** Meta-label y_t^{meta} 1: Set entry price: $P_{\text{entry}} = P_t$ 2: Set barriers: $P_{upper} = [BLANK 1], P_{lower} = [BLANK 2]$ 3: Initialize exit $time = t + H_{max}$, exit $price = P_{t+H_{max}}$ 4: for h = 1 to H_{max} do if $P_{t+h} \geq P_{upper}$ then 5:exit time = t + h, exit price = P_{t+h} 6: Break 7: else if $P_{t+h} \leq P_{\text{lower}}$ then 8: exit time = t + h, exit price = P_{t+h} 9: Break 10: end if 11:12: end for 13: Compute trade return: $R_t = s_t \times \frac{exit_price-P_{entry}}{P_{entry}}$ 14: Generate meta-label: $y_t^{\text{meta}} = [\text{BLANK 3}]$ \triangleright Fill in the blank 15: return y_t^{meta}

2. Meta-model Feature Engineering

The meta-model requires features that capture market environment and trading conditions, distinct from the price-based technical indicators used by the primary model. We construct features across several categories: volatility regime indicators, latent states from HMMs, and primary model confidence measures (Random Forest prediction probabilities, feature importance stability). Let $\mathbf{X}_t \in \mathbb{R}^{D_{\text{meta}}}$ denote the meta-feature vector at time t.

3. Variable Selection Networks for Binary Classification

We implement a Variable Selection Network (VSN) that simultaneously performs feature selection and binary classification. The VSN architecture consists of a feature selection layer that learns attention weights for each feature, and classification layers that process selected features to output probability $p_t = P(y_t^{\text{meta}} = 1 | \mathbf{X}_t)$.

The VSN produces probabilistic outputs that must be converted to binary decisions using threshold τ :

$$\hat{y}_t^{\text{meta}}(\tau) = \begin{cases} 1 & \text{if } p_t > \tau \\ 0 & \text{if } p_t \leq \tau \end{cases}$$

The final trading decision is:

Final Signal_t(
$$\tau$$
) =

$$\begin{cases} s_t & \text{if } p_t > \tau \\ 0 & \text{if } p_t \le \tau \end{cases}$$

Different threshold values τ control the trade-off between precision and recall, effectively filtering different proportions of the primary model's signals.

We evaluate three trading strategies on a test set of 2000 trading opportunities where the primary model generated signals. The ground truth shows 1200 profitable trades and 800 unprofitable trades. Consider the following confusion matrices:

• Strategy A (Primary Model Only): Always trade when primary model signals.

	Predicted Trade	Predicted No Trade
Actual Profitable	1200	0
Actual Unprofitable	800	0

• Strategy B (Primary + Meta-model, τ_1):

	Predicted Trade	Predicted No Trade
Actual Profitable	1080	120
Actual Unprofitable	600	200

• Strategy C (Primary + Meta-model, τ_2):

	Predicted Trade	Predicted No Trade
Actual Profitable	950	250
Actual Unprofitable	350	450

Q2. Calculate the precision and recall for each strategy.

Q3. Using the precision and recall values calculated in Q2:

- Compute the F1 score for each strategy
- Determine which meta-model threshold $(\tau_1 \text{ or } \tau_2)$ provides the optimal balance between precision and recall.

Exercise 5 (Variable Selection Networks).

We now implement a Variable Selection Network (VSN) to predict the metalabels $y_t^{\text{meta}} \in \{0, 1\}$ using the meta-features $\mathbf{X}_t \in \mathbb{R}^D$ constructed in Exercise 4. Our training dataset consists of N_T samples $\{(\mathbf{X}_i, y_i^{\text{meta}})\}_{i=1}^{N_T}$ where each feature vector \mathbf{X}_i contains both numerical and categorical variables.

The feature space is partitioned as $D = D_n + D_c$, where:

• D_n numerical features: x_i^d for $d \in \{1, \ldots, D_n\}$

- D_c categorical features: $x_i^{d'}$ for $d' \in \{D_n + 1, \dots, D_n + D_c\}$
- Each categorical variable $x_i^{d'}$ has $n_{d'}$ possible categories

The VSN architecture transforms each feature (numerical or categorical) into a D_e -dimensional embedding vector, then applies gating mechanisms for feature selection and final prediction.

1. Input Transformation Layer

The first layer processes numerical and categorical features differently:

- Dense Layers for numerical features: $x_i^d \to \boldsymbol{\xi}_i^d \in \mathbb{R}^{D_e}$
- Embedding Layers for categorical features: $x_i^{d'} \rightarrow \boldsymbol{\xi}_i^{d'} \in \mathbb{R}^{D_e}$

Q1. Consider the embedding layer for categorical features:

- For a categorical variable d' with $n_{d'}$ categories embedded into D_e dimensions, calculate the total number of parameters in the embedding layer.
- Explain two key advantages of embedding layers over one-hot encoding for categorical variables with high cardinality.

The Input Transformation layer, illustrated in Figure 4, converts the input matrix to embedded feature vectors.



Figure 4: Input Transformation Layer processing numerical and categorical features

After applying the Input Transformation layer to a batch of N samples with D_n numerical and D_c categorical features, each embedded into D_e dimensions.

Q2. What is the shape of the output tensor containing all embedded features?

2. Gated Residual Networks (GRN)

Each embedded feature $\boldsymbol{\xi}_i^d$ is processed through a Gated Residual Network, as shown in Figure 5, which applies gating mechanisms and residual connections to control information flow.



Figure 5: Gated Residual Network architecture with gating mechanism

The GRN transformation produces $\tilde{\xi}_i^d \in \mathbb{R}^{D_o}$ for each feature *d*. These transformed features are then processed through the feature selection mechanism illustrated in Figure 6.



Figure 6: GRN transformation and feature selection mechanism

- **Q3.** The concatenated feature vector $\boldsymbol{\xi}_i$ is processed through a GRN followed by a Softmax activation to produce attention weights $\boldsymbol{\alpha}_i$:
 - What is the purpose of applying the Softmax function to obtain α_i ?
 - What is the shape of α_i and what constraint does the Softmax impose?
 - What does each component α_i^d represent in the context of feature importance?

3. Feature Selection and Final Representation

The final feature representation is computed as a weighted sum:

$$ilde{oldsymbol{\xi}}_i = \sum_{d=1}^{D_n+D_c} lpha_i^d ilde{oldsymbol{\xi}}_i^d$$

The complete VSN architecture is illustrated in Figure 7.



Figure 7: Complete Variable Selection Network architecture

- **Q4.** Describe the complete forward propagation from an input batch matrix $\mathbf{X} \in \mathbb{R}^{N_b \times D}$ to the final representation $\tilde{\boldsymbol{\xi}} \in \mathbb{R}^{N_b \times D_o}$. For each transformation step, specify:
 - The operation performed
 - The shape of the data after the transformation
 - The key parameters involved

Assume a batch size of N_b samples.

4. Binary Classification Layer

Since we have a binary classification problem for predicting meta-labels, we add a final Dense layer to the VSN architecture.

Q5. Describe the final Dense layer for binary classification:

- What transformation does this layer perform on $\tilde{\xi}_i \in \mathbb{R}^{D_o}$?
- How many parameters does this layer contain?
- What is the shape of the final output and what does it represent?
- What activation function should be used and why?

5. Loss Function and Training

- **Q6.** What loss function should be used to train this VSN model for binary classification? Provide the mathematical formulation.
- **Q7.** Complete the Stochastic Gradient Descent training algorithm for the VSN model:

Algorithm 7 SGD Training for Variable Selection Network

Require: Training data, validation data, learning rate η , batch size B, epochs EEnsure: Optimal parameters θ^*

- 1: Initialize parameters $\theta^{(0)} = [BLANK \ 1]$ \triangleright Fill in the blank
- 2: Initialize train_losses = [], val_losses = []
- 3: $t \leftarrow 0$
- 4: for epoch e = 1 to E do
- 5: for each batch $\{(\mathbf{Z}_i, y_i^{\text{meta}})\}_{i \in \mathcal{B}}$ in training data do
- 6: Perform forward propagation: $\{p_i\}_{i \in \mathcal{B}} = \text{VSN}(\{\mathbf{Z}_i\}_{i \in \mathcal{B}}; \boldsymbol{\theta}^{(t)})$
- 7: Calculate loss function: $L_{\mathcal{B}} = \frac{1}{|\mathcal{B}|} \sum_{i \in \mathcal{B}} \ell(p_i, y_i^{\text{meta}})$
- 8: Update weights: $\theta^{(t+1)} = [BLANK 2]$ \triangleright Fill in the blank

9:
$$t \leftarrow t + 1$$

- end for
- 11: Compute training loss and append to *train_losses*
- 12: Compute validation loss and append to *val_losses*
- 13: end for

10:

- 14: Select $\boldsymbol{\theta}^*$ with minimum validation loss
- 15: return θ^*