# The Smart Supplier: Optimizing Orders in a Fluctuating Market - 6 Marks

## Group ID:

Group ID: Group 278

## Group Members Name with Student ID:

1. Rahul Prashar – 2024AB05058 – 100%
2. Rahul Agarwal – 2024AA05676 – 100%
3. Rahul Sinha – 2024AA05036 – 100%
4. Raghavendra Sathakarni A S – 2024AB05118 – 100%

## Assginment 1 – Part 2

Develop a reinforcement learning agent using dynamic programming to help a Smart Supplier decide which products to manufacture and sell each day to maximize profit. The agent must learn the optimal policy for choosing daily production quantities, considering its limited raw materials and the unpredictable daily demand and selling prices for different products.

### Scenario

A small Smart Supplier manufactures two simple products: Product A and Product B. Each day, the supplier has a limited amount of raw material. The challenge is that the market demand and selling price for Product A and Product B change randomly each day, making some products more profitable than others at different times. The supplier needs to decide how much of each product to produce to maximize profit while managing their limited raw material.

### Objective

The Smart Supplier's agent must learn the optimal policy $\pi*$ using dynamic programming (Value Iteration or Policy Iteration) to decide how many units of Product A and Product B to produce each day to maximize the total profit over the fixed number of days, given the daily changing market conditions and limited raw material.

## 1. Custom Environment Creation (SmartSupplierEnv)

```
In [8]:  import random
```

```
In [9]:  # Define market states and their product prices
         # Structure: {Market_State_ID: {'A_price': X, 'B_price': Y}}
         # Define product raw material costs
         # Define actions: (num_A, num_B, raw_material_cost_precalculated)
                 # Action ID mapping:
                 # 0: Produce_2A_0B
```

```python
        # 1: Produce_1A_2B
        # 2: Produce_0A_5B
        # 3: Produce_3A_0B
        # 4: Do_Nothing

 # Define state space dimensions
        # Current Day: 1 to num_days
        # Current Raw Material: 0 to initial_raw_material
        # Current Market State: 1 or 2

# get reward function
class SmartSupplierEnv:
    """
    Environment modeling a Smart Supplier with fluctuating daily market p
    and limited raw material to produce two products A and B.
    """
    def __init__(self, num_days=5, init_rm=10):
        """
        Initialize the environment parameters.
        Args:
            num_days (int): Total number of days in an episode.
            init_rm (int): Initial units of raw material available each d
        """
        self.num_days = num_days
        self.init_rm = init_rm
        # Define market states and their selling prices
        # Market state 1: High demand for A; state 2: High demand for B
        self.market_states = [1, 2]
        self.prices = {
            1: {'A': 8, 'B': 2},   # Market State 1 prices
            2: {'A': 3, 'B': 5},   # Market State 2 prices
        }
        # Define actions mapping: action_id -> (units_A, units_B, rm_cost
        self.actions = {
            0: (2, 0, 2*2),        # Produce 2A, 0B costing 4 RM
            1: (1, 2, 1*2 + 2*1),  # Produce 1A, 2B costing 4 RM
            2: (0, 5, 5*1),        # Produce 0A, 5B costing 5 RM
            3: (3, 0, 3*2),        # Produce 3A, 0B costing 6 RM
            4: (0, 0, 0),          # Do nothing
        }

    def get_reward(self, market, action):
        """
        Compute the profit reward for taking a given action in a market s
        Args:
            market (int): Current market state (1 or 2).
            action (int): Action ID.
        Returns:
            int: Profit earned (zero if invalid or do-nothing).
        """
        units_A, units_B, cost = self.actions[action]
        # Price per unit from current market state
        price_A = self.prices[market]['A']
        price_B = self.prices[market]['B']
        return units_A * price_A + units_B * price_B
```

## 2. Dynamic Programming Implementation (Value Iteration or Policy Iteration) (2 Mark)

```python
# Value Iteration function
def value_iteration(env, gamma=1.0, theta=1e-6):
    """
    Perform Value Iteration over the state space to compute optimal V* an
    Args:
        env (SmartSupplierEnv): The Smart Supplier environment.
        gamma (float): Discount factor (1.0 for episodic sum maximization
        theta (float): Convergence threshold.
    Returns:
        V (dict): Mapping from (day, rm, market) -> value.
        policy (dict): Mapping from (day, rm, market) -> best action ID.
    """
    V = {}        # State-value function
    policy = {}  # Optimal policy mapping
    # Initialize V(s) = 0 for all states
    for day in range(1, env.num_days + 1):
        for rm in range(env.init_rm + 1):
            for m in env.market_states:
                V[(day, rm, m)] = 0.0

    while True:
        delta = 0.0
        # Sweep over all states
        for day in range(1, env.num_days + 1):
            for rm in range(env.init_rm + 1):
                for m in env.market_states:
                    state = (day, rm, m)
                    v_old = V[state]
                    action_values = {}
                    # Evaluate each action
                    for a, (_, _, cost) in env.actions.items():
                        # If not enough raw material, reward is zero
                        if cost > rm:
                            reward = 0
                        else:
                            reward = env.get_reward(m, a)
                        # Next state transition
                        next_day = day + 1
                        if next_day > env.num_days:
                            future = 0.0
                        else:
                            # Expectation over next day's market shift (5
                            future = 0.0
                            for m2 in env.market_states:
                                future += 0.5 * V[(next_day, env.init_rm,
                        action_values[a] = reward + gamma * future
                    # Pick best action
                    best_action = max(action_values, key=action_values.ge
                    V[state] = action_values[best_action]
                    policy[state] = best_action
                    delta = max(delta, abs(v_old - V[state]))
        if delta < theta:
            break
    return V, policy
```

## 3. Simulation and Policy Analysis ( 1 Mark)

In [11]:
```python
# simulate policy function – Simulates the learned policy over multiple r

# analyze policy function – Analyzes and prints snippets of the learned o
def simulate_policy(env, policy, episodes=1000, seed=42):
    """
    Simulate the learned policy to estimate average total profit.
    Args:
        env (SmartSupplierEnv): Environment instance.
        policy (dict): Optimal policy mapping states to actions.
        episodes (int): Number of simulation runs.
        seed (int): Random seed for reproducibility.
    Returns:
        float: Average total profit over all episodes.
    """
    random.seed(seed)  # Ensure reproducibility
    total_profits = []
    for _ in range(episodes):
        profit = 0
        # Random initial market state each episode
        m = random.choice(env.market_states)
        # Run through days
        for day in range(1, env.num_days + 1):
            state = (day, env.init_rm, m)
            a = policy[state]
            # Apply action if RM constraint allows
            if env.actions[a][2] <= env.init_rm:
                profit += env.get_reward(m, a)
            # Market shifts for next day
            m = random.choice(env.market_states)
        total_profits.append(profit)
    # Return mean profit
    return sum(total_profits) / episodes


def analyze_policy(policy, env):
    """
    Print key insights of the optimal policy for selected states.
    """
    print("\nPolicy snippets at Day 1:")
    for market in env.market_states:
        for rm in [0, env.init_rm//2, env.init_rm]:
            action = policy[(1, rm, market)]
            units = env.actions[action][:2]
            print(f"Market={market}, RM={rm} -> Action {action}, produce
```

## Performance Evaluation: (1 Mark)

In [12]:
```python
# --- Main Execution ---
if __name__ == "__main__":
    # Create environment
    env2 = SmartSupplierEnv(num_days=5, init_rm=10)

    # Compute optimal values and policy via Value Iteration
    V_opt, policy_opt = value_iteration(env2)

    # Display key value
    print(f"V*(Day1, RM=10, Market=1) = {V_opt[(1,10,1)]:.2f}")
```

```python
        # Show policy examples
        analyze_policy(policy_opt, env2)

        # Simulate policy performance
        avg_profit = simulate_policy(env2, policy_opt, episodes=1000)
        print(f"Average profit over 1000 runs: {avg_profit:.2f}")
```

```
V*(Day1, RM=10, Market=1) = 122.00

Policy snippets at Day 1:
Market=1, RM=0 -> Action 0, produce A=2, B=0
Market=1, RM=5 -> Action 0, produce A=2, B=0
Market=1, RM=10 -> Action 3, produce A=3, B=0
Market=2, RM=0 -> Action 0, produce A=2, B=0
Market=2, RM=5 -> Action 2, produce A=0, B=5
Market=2, RM=10 -> Action 2, produce A=0, B=5
Average profit over 1000 runs: 122.50
```

## 5. Impact of Dynamics Analysis (1 Mark)

Discusses the impact of dynamic market prices on the optimal policy.

The learned optimal policy adapts strategically to daily market fluctuations:

- 1. **Market Sensitivity** – When the market state favors Product A (State 1), the policy allocates more raw material to produce A, maximizing high-margin returns, and vice versa for Product B in State 2. This dynamic switching captures the core benefit of reacting to price signals.
- 2. **Resource Conservation** – On days with limited remaining raw material or earlier in the horizon, the policy may choose smaller bundles or even Do_Nothing to preserve capacity, anticipating a potential shift to a more lucrative market.
- 3. **Horizon Effects** – As the final day approaches, the policy becomes more aggressive, exhausting all available raw material regardless of state, since there is no future value to conserve. This results in maximal immediate profit extraction.
- 4. **Robustness to Uncertainty** – By embedding a 50/50 expectation over next-day market states in the value iteration, the policy inherently balances expected gains across both possibilities, avoiding overcommitment to one product when future states are ambiguous.

Overall, dynamic programming produces a policy that intelligently balances exploitation of current high-price conditions with cautious preservation of resources for uncertain future opportunities, thereby maximizing cumulative profit over time.