Hangman AI: Analysis ReportBidirectional HMM + Q-Learning Reinforcement Learning

**Project:** UE23CS352A Machine Learning Hackathon

**Team:** Team Number 4: AM236,AM227,AM197,AM214

**Date:** November 2024-----1.
Key Observations
Most Challenging Parts
1.1 Overfitting Problem

The most significant challenge was the **severe overfitting** observed during development:

- **Training Performance:** ~95-96% win rate on training corpus
- **Test Performance:** ~32-33% win rate on test set
- **Gap:** 60%+ performance drop between training and test data

**Root Cause Analysis:**

- The RL agent learned to memorize patterns from the training corpus (50,000 words)
- Test set contains different words, many rare/technical terms not well-represented in training
- Q-table became too specific to training patterns, failing to generalize

**Insight Gained:**

- In sequence prediction tasks, RL can easily overfit to training vocabulary
- HMM as the foundation (with heavy weighting) helps prevent overfitting
- Need for better regularization and more diverse training strategies

1.2 State Space Explosion

Another major challenge was the **exponentially large state space**:

- Each state: `masked_word:guessed_letters:lives`
- Example: `"APP_E:{a,p,e}:4"`
- With 26 letters, various word lengths, and 6 lives, state space ≈ 10^15+ states

**Insight:**

- Q-table approach becomes infeasible for large state spaces
- Heavy reliance on HMM (expert oracle) was necessary
- Considered Deep Q-Network (DQN) but time constraints limited implementation

1.3 HMM-RL Integration Balance

Finding the right balance between HMM guidance and RL learning was critical:

- Too much RL → overfitting, poor generalization
- Too much HMM → no learning, static strategy
- Found optimal: HMM weight × 50, RL weight × 1

**Key Insight:**

- In domain-specific tasks, trust the expert (HMM) more than the learner (RL)
- RL should refine HMM predictions, not replace them

-----2. Strategies2.1 Hidden Markov Model (HMM) Design ChoicesArchitecture Decision: Bidirectional First-Order HMM

**Choice:** Implemented bidirectional HMM with first-order transitions

**Rationale:**

1. **Bidirectional Context:**
   - Traditional HMM only uses forward context (previous letters)
   - Bidirectional uses both forward AND backward context
   - Example: Pattern `"_PPLE"` - forward context is limited (start of word), but backward context ($? \rightarrow P$) is highly informative
   - Result: Better predictions when pattern is partially revealed
2. **First-Order vs Higher-Order:**
   - First-order: P(letter | previous_letter) - simpler, more generalizable
   - Higher-order (2nd, 3rd): More context but higher risk of overfitting
   - Compromise: Added trigram patterns as additional feature (not full 2nd-order HMM)
   - Result: Good balance between context and generalization

Key HMM Components:

**1. Vocabulary Pattern Matching (Most Accurate):**
# Direct lookup in corpus vocabulary
matching_words = find_words_matching_pattern(masked_word)
if matching_words:
   return letter_frequencies_from_matches()

- **Why:** Most accurate when corpus contains the word
- **When:** Pattern has few matches (< 100 words)
- **Accuracy:** ~90-95% when pattern is specific enough

**2. Position-Specific Probabilities:**
P(letter | word_length, position)

- **Why:** Letter frequencies vary by position (e.g., 'Q' never starts words)
- **Implementation:** Separate counters for each (length, position) combination
- **Result:** Better predictions for specific word positions

### 3. Bigram and Trigram Patterns:
Bigram: P(letter | previous_letter)
Trigram: P(letter | previous_2_letters)

- **Why:** Captures common letter sequences (e.g., "APP" → "L")
- **Usage:** Combined with bidirectional context for maximum accuracy
- **Weight:** Trigram has 2.5x boost when pattern is common (> 50 occurrences)

### 4. Reverse Transitions (Bidirectional):
Reverse: P(letter | next_letter)

- **Why:** Patterns like `"_PPLE"` benefit from knowing what comes after
- **Example:** "P" commonly follows "A" in English (APPLE, APPROACH)
- **Result:** 20-30% accuracy improvement on partially-revealed patterns

HMM Smoothing:

**Choice:** Laplace (Additive) Smoothing with $\alpha = 0.01$

**Rationale:**

- Prevents zero probabilities for unseen transitions
- Allows generalization to new words
- Small $\alpha$ value: conservative smoothing, maintains strong patterns
- Tested $\alpha = 0.1, 0.05, 0.01 \rightarrow 0.01$ gave best validation performance

Vocabulary Storage Strategy:

**Choice:** Store entire vocabulary (50,000 words) for pattern matching

**Rationale:**

- Pattern matching is most accurate method (when matches exist)
- Memory cost acceptable (50K words ≈ 500KB)
- Dramatically improves accuracy when word is in corpus

**Trade-off:**

- Only works for words in corpus (test set overlaps but not identical)
- Alternative: Pattern matching only, no vocabulary → lower accuracy on known words

-----2.2 Reinforcement Learning Design2.2.1 State Representation

**State Design:**
state_key = f"{masked_word}:{sorted_guessed_letters}:{lives_left}"
**Example:** `"APP_E:{a,p,e}:4"`

**Components:**

1. **Masked Word Pattern:** Current revealed pattern (e.g., "APP_E")
    - Encodes position of revealed letters
    - Critical for HMM predictions
2. **Guessed Letters Set:** Sorted string of all guessed letters (e.g., "aep")
    - Prevents repeated guesses
    - Used to filter available actions
3. **Lives Remaining:** Integer (0-6)
    - Affects risk-taking behavior
    - Important for Q-value learning

**Rationale for This Design:**

- **String-based state:** Simple, interpretable, fast lookup
- **Included lives:** RL needs to know risk level (few lives = conservative)
- **Sorted guessed letters:** Deterministic state representation (order doesn't matter)
- **No word length:** Already encoded in masked_word

**Alternative Considered:**

- Vector representation (one-hot + probabilities)
- **Rejected:** Too complex, state space explosion, slower lookup

2.2.2 Action Space

**Design:** Discrete action space - 26 letters of alphabet
actions = {a, b, c, ..., z} - guessed_letters
**Why Simple:**

- Natural action space for Hangman
- No need for complex actions
- Easy to filter (remove already guessed letters)

2.2.3 Reward Function Design

**Final Reward Structure:**
Reward = {
    'correct_guess': +2.0 × occurrences,     # Per letter revealed
    'wrong_guess': -1.5,              # Lose a life

```
    'repeated_guess': -1.0,              # Inefficiency penalty
    'win_game': +25.0,                   # Large win bonus
    'lose_game': -8.0                    # Loss penalty
}
```
**Rationale:**

1. **Correct Guess: +2.0 × occurrences:**
   ○ Rewards discovering multiple letters (e.g., "E" appears twice = +4.0)
   ○ Encourages strategic guessing (common letters like "E", "A")
   ○ Higher weight than wrong guess penalty maintains exploration
2. **Wrong Guess: -1.5:**
   ○ Penalty for losing a life
   ○ Balanced: not too harsh (would discourage exploration) but not too lenient
   ○ Tested: -0.5 (too lenient), -2.0 (too harsh) → -1.5 optimal
3. **Repeated Guess: -1.0:**
   ○ Discourages inefficiency
   ○ Smaller than wrong guess (repeating is bad, but not as bad as wrong)
   ○ In practice: Environment prevents repeats, so this rarely triggers
4. **Win Bonus: +25.0:**
   ○ Large positive reward for success
   ○ Dominates other rewards (25 vs 2×5 = 10 for correct guesses)
   ○ Ensures agent prioritizes winning over minimizing guesses
5. **Lose Penalty: -8.0:**
   ○ Penalty for failure
   ○ Less than win bonus (encourages playing to win, not avoiding loss)
   ○ Combined with wrong guesses: -8.0 + (6×-1.5) = -17.0 total loss penalty

**Reward Balance:**

● Typical win: +25 + (5 correct × 2) = +35
● Typical loss: -8 + (6 wrong × -1.5) = -17
● Ratio: 2:1 win/loss ratio encourages risk-taking to win

**Iterations:**

● **Initial:** +10 win, -5 loss → too conservative
● **Version 2:** +50 win, -20 loss → too aggressive, overfitting
● **Final:** +25 win, -8 loss → balanced, good generalization signal

2.2.4 Q-Learning Algorithm

**Algorithm:** Tabular Q-Learning

**Update Rule:**
$Q(s, a) \leftarrow Q(s, a) + \alpha[r + \gamma \times \max_{a'} Q(s', a') - Q(s, a)]$

**Hyperparameters:**

- **Learning Rate (α):** 0.3
    - High enough for fast learning
    - Low enough for stability
- **Discount Factor (γ):** 0.98
    - High value: agent values long-term rewards (winning game)
    - Close to 1: considers future implications of each guess
- **Initial Q-values:** 0 (optimistic initialization via HMM)

**Why Q-Learning (not DQN):**

- State space is discrete and manageable (string keys)
- Tabular Q-learning is simpler, faster for this problem
- DQN would require neural network, more complexity, training time
- **Future work:** DQN could handle state space explosion better

**Integration with HMM:**
combined_score = Q(s, a) + (HMM_prob(a) × 50)
action = argmax_a combined_score

- HMM weight (50x) dominates Q-values
- RL learns refinements, not replacement of HMM strategy
- This design prevents overfitting to training vocabulary

-----3. ExplorationExploration vs Exploitation Trade-offStrategy: Epsilon-Greedy with HMM-Guided Exploration

**Implementation:**
ε-greedy with ε decay: 1.0 → 0.01 over 20,000 episodes
**Exploration Phase (High ε):**

- ε = 1.0 → 0.5: 100% exploration initially
- Explores HMM's top 5 predictions (not random!)
- **Smart exploration:** Even when exploring, uses HMM guidance
- Result: Learns from good guesses, not pure randomness

**Transition Phase (Medium ε):**

- ε = 0.5 → 0.1: Decreasing exploration
- Mixes exploitation (Q-values + HMM) with exploration
- Balances learning new states vs using learned knowledge

**Exploitation Phase (Low ε):**

- ε = 0.1 → 0.01: Mostly exploitation

- Uses learned Q-values combined with HMM
- Fine-tunes strategy on known patterns

**Key Innovation: HMM-Guided Exploration**

Instead of pure random exploration:
```
if exploring:
    # Bad: return random.choice(alphabet)  # Pure random
    # Good:
    top_5_hmm = sorted(hmm_probs, reverse=True)[:5]
    return random.choice(top_5_hmm)  # Explore HMM's suggestions
```
**Why This Works:**

- Exploration still discovers new states
- But explores **promising** states first (guided by HMM)
- Faster learning than pure random exploration
- Reduces wasted guesses on unlikely letters

**Epsilon Decay Schedule:**
$\varepsilon\_decay = 0.9995$ per episode
$episodes\_to\_50\% = \log(0.5) / \log(0.9995) \approx 1,386$ episodes
$episodes\_to\_10\% = \log(0.1) / \log(0.9995) \approx 4,605$ episodes
**Result:**

- ~1,400 episodes: 50% exploration/50% exploitation
- ~4,600 episodes: Mostly exploitation with occasional exploration
- Final: 1% random exploration maintains adaptability

**Validation of Exploration Strategy:**

- Training win rate increases from ~40% (early) to ~96% (final)
- Q-table grows from 0 to ~100,000+ states
- Test performance: 32% (indicates need for more exploration OR less overfitting)

-----4. Future Improvements

Given an additional week, here are the improvements I would prioritize:4.1 Immediate Fixes (Days 1-2)4.1.1 Deep Q-Network (DQN) Implementation

**Problem:** Tabular Q-learning can't scale to full state space

**Solution:** Replace Q-table with neural network

**Implementation:**

- Input: Encoded state vector (masked_word one-hot + HMM probabilities + lives)

- Architecture: 3-layer MLP (256 → 128 → 26)
- Output: Q-values for 26 actions
- Experience replay buffer: Store and sample diverse experiences

**Expected Impact:** +10-15% test accuracy (better generalization)4.1.2 Ensemble of HMMs by Word Length

**Problem:** Single HMM treats all word lengths equally

**Solution:** Train separate HMMs for different word lengths

**Implementation:**
```
HMMs = {
    3-5: ShortHMM(),      # Common words, simpler patterns
    6-8: MediumHMM(),     # Most common length range
    9-12: LongHMM(),      # Complex words
    13+: VeryLongHMM()    # Rare technical terms
}
```
**Expected Impact:** +5-8% accuracy (better pattern matching)4.2 Advanced Improvements (Days 3-4)4.2.1 Curriculum Learning

**Problem:** Training on random words doesn't match difficulty progression

**Solution:** Start with easy words, gradually increase difficulty

**Implementation:**

- Phase 1: Common 4-6 letter words
- Phase 2: Common 7-9 letter words
- Phase 3: Uncommon words
- Phase 4: All words (final fine-tuning)

**Expected Impact:** +8-12% accuracy (better learning progression)4.2.2 Attention Mechanism for HMM

**Problem:** Current HMM treats all positions equally

**Solution:** Attention-weighted HMM predictions

**Implementation:**

- Learn attention weights for which positions/patterns matter most
- Weighted combination: position × bigram × trigram × global
- Attention learned via neural network or gradient descent

**Expected Impact:** +5-7% accuracy (focuses on important patterns)4.2.3 Multi-Step Lookahead

**Problem:** Greedy one-step decisions may not be optimal

**Solution:** Consider future implications of current guess

**Implementation:**

- Tree search (Monte Carlo Tree Search) for 2-3 steps ahead
- Evaluate: "If I guess E now, what will HMM predict next?"
- Choose action with best expected long-term outcome

**Expected Impact:** +7-10% accuracy (better strategic planning)4.3 Research-Level Improvements (Days 5-7)4.3.1 Transformer-Based HMM

**Problem:** Current HMM only uses local context (bigrams/trigrams)

**Solution:** Transformer encoder for full word context

**Implementation:**

- Pre-train transformer on corpus
- Fine-tune for Hangman predictions
- Captures long-range dependencies (e.g., "QU" pattern affects later letters)

**Expected Impact:** +10-15% accuracy (better context understanding)4.3.2 Meta-Learning for Few-Shot Adaptation

**Problem:** Model doesn't adapt to test set characteristics

**Solution:** Meta-learning to quickly adapt to new word distributions

**Implementation:**

- Learn to learn: optimize for fast adaptation
- Given small sample of test words, quickly adjust strategy
- MAML (Model-Agnostic Meta-Learning) algorithm

**Expected Impact:** +12-18% accuracy (adapts to test distribution)4.3.3 Uncertainty Quantification

**Problem:** Model doesn't know when it's uncertain

**Solution:** Bayesian RL or ensemble methods

**Implementation:**

- Track uncertainty in HMM predictions
- When uncertain, explore more; when confident, exploit
- Ensembles of multiple HMMs vote on predictions

**Expected Impact:** +8-12% accuracy (better decision-making)4.4 Practical Optimization (Any Time)4.4.1 Hyperparameter Optimization

**Current:** Manual tuning

**Future:** Automated search (Bayesian Optimization, Grid Search)

**Parameters to Optimize:**

- HMM smoothing ($\alpha$)
- RL learning rate, discount factor
- Reward function weights
- HMM weight multiplier
- Exploration schedule

**Expected Impact:** +3-5% accuracy (finds optimal configuration)4.4.2 Data Augmentation

**Problem:** Limited to 50,000 words in corpus

**Solution:** Generate synthetic words or use external dictionaries (if allowed)

**Implementation:**

- If rules allow: Add common English words
- Pattern-based word generation
- Cross-validation to prevent overfitting

**Expected Impact:** +5-8% accuracy (more training data)4.4.3 Transfer Learning

**Problem:** Starting from scratch each time

**Solution:** Pre-train on larger word lists, fine-tune on corpus

**Implementation:**

- Pre-train HMM on general English word lists (100K+ words)
- Fine-tune on specific corpus (50K words)
- Transfer learned patterns to new domain

**Expected Impact:** +6-10% accuracy (better base knowledge)-----5. Summary

This project successfully implemented a **Bidirectional HMM + Q-Learning** system for Hangman, achieving:

- **Training Accuracy:** ~96%
- **Test Accuracy:** ~32%
- **Key Innovation:** Heavy HMM weighting (50x) with RL refinement

**Main Challenges:**

1. Severe overfitting (60%+ gap)
2. State space explosion
3. HMM-RL balance

**Key Insights:**

1. Trust the expert (HMM) more than the learner (RL) in domain-specific tasks
2. Bidirectional context dramatically improves predictions
3. Pattern matching is most accurate when vocabulary matches
4. Exploration must be guided (by HMM), not random

**Future Work:**

- DQN for better state space handling
- Curriculum learning for progressive difficulty
- Transformer-based HMM for long-range context
- Meta-learning for test set adaptation

The system demonstrates solid understanding of HMM, RL, and their integration, with clear paths for improvement given more time.-----**Total Development Time:** ~40-50 hours
**Lines of Code:** ~1,500
**Key Files:** `final_hmm_rl.py`, `hangman_env.py`, `demo.py`