

# PROGRAMMATIC REPRESENTATION LEARNING WITH LANGUAGE MODELS

**Anonymous authors**

Paper under double-blind review

## ABSTRACT

Classical models for supervised machine learning, such as decision trees, are efficient and interpretable predictors, but their quality is highly dependent on the particular choice of input features. Although neural networks can learn useful representations directly from raw data (e.g., images or text), this comes at the expense of interpretability and the need for specialized hardware to run them efficiently. In this paper, we explore a hypothesis class we call *Learned Programmatic Representations* (LeaPR) models, which stack arbitrary features represented as code (functions from data points to scalars) and decision tree predictors. We synthesize feature functions using Large Language Models (LLMs), which have rich prior knowledge in a wide range of domains and a remarkable ability to write code using existing domain-specific libraries. We propose two algorithms to learn LeaPR models from supervised data. First, we design an adaptation of FunSearch to learn *features* rather than directly generate predictors. Then, we develop a novel variant of the classical ID3 algorithm for decision tree learning, where new features are generated on demand when splitting leaf nodes. In experiments from chess position evaluation to image and text classification, our methods learn high-quality, neural network-free predictors often competitive with neural networks. Our work suggests a flexible paradigm for learning interpretable representations end-to-end where features and predictions can be readily inspected and understood.

## 1 INTRODUCTION

The central problem in supervised machine learning is to find a predictor  $h : X \rightarrow Y$  in a hypothesis class  $\mathcal{H}$  that minimizes a certain risk function  $\mathcal{R}(h)$ , such as 0 – 1 error in classification or mean-squared error in regression (Michalski et al., 2013). Classical choices for  $\mathcal{H}$  include linear models, decision trees, and ensembles thereof, which are compellingly simple to understand and debug, and are both compute- and data-efficient. However, their effectiveness is highly limited in domains with **unstructured**, high-dimensional inputs, such as images or text. For these domains, high-quality models are often best learned by first constructing a high level *representation* of an input  $x \in X$  using a set of features  $f_i : X \rightarrow \mathcal{R}$  that yield a higher-level encoding of the input that predictors can then rely on. While this offers great flexibility, in practice the effort and domain expertise required to *engineer* a good set of features for a particular learning task severely limits the quality of models that can be obtained with classical predictors in high-dimensional input domains without extensive human effort (Dong & Liu, 2018; Cheng & Camargo, 2023).

A remarkably successful paradigm that avoids the need for hand-designed feature engineering is *deep learning*, where  $\mathcal{H}$  is set to a parameterized family of neural networks of a domain-appropriate architecture. The core advantage of deep learning is the ability of gradient-based optimization to automatically learn useful representations from raw data (Bengio & LeCun, 2007; Damian et al., 2022). Indeed, deep neural networks can be seen as computing a set of complex, non-linear neural features, then applying a simple predictor on top (e.g., the last fully-connected layer, corresponding to a linear model). However, despite being highly effective for *prediction*, neural features have several drawbacks. First, deep neural networks **trained as end-to-end predictors are data-intensive**, and their ability to generalize drops drastically when in-domain data are scarce (Wang et al., 2023). Second, neural features are not easily interpretable: analyses of large-scale neural networks typically only find a fraction of neurons that seem to encode human-aligned concepts (Huben et al., 2023). This limits the potential of neural models to provide faithful *explanations* for their predictions (e.g., express *why* a given  $x$  is being classified as  $y$ ), which are important when experts rely on learned models to support high-stakes decision-making (Doshi-Velez & Kim, 2017).

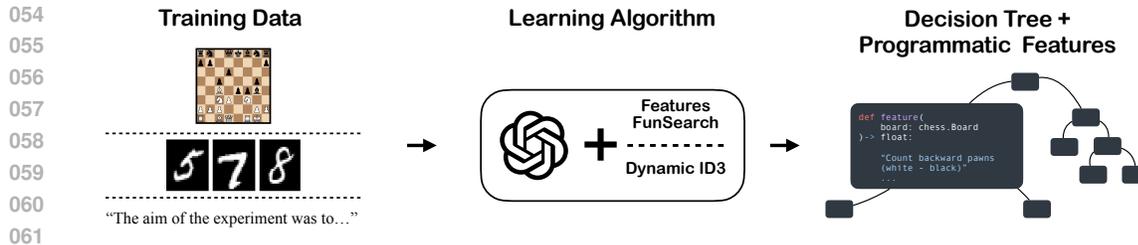


Figure 1: Learned Programmatic Representation models combine *programmatic features*, synthesized by LLMs as code, and decision tree predictors, yielding interpretable models. We give two algorithms for learning them end-to-end from supervised data in high-dimensional input domains.

In this paper, we seek to investigate alternative paradigms for representation learning that, like deep learning, do not require manual feature engineering from users and, like classical methods, yield interpretable, fast, and efficiently-learnable predictors. To that end, we propose learning *LeaPR* (Learned Programmatic Representation) models: a class of predictors with programmatic features paired with decision tree predictors, as illustrated in Figure 1. We leverage Large Language Models’ (LLMs) ability to generate code using domain-specific libraries to encode features as arbitrary Python functions from the input domain into the reals: these functions are generated during training with the goal of improving the empirical risk of the current predictor.

For learning from supervised data, we propose two alternative methods. First, we explore a variant of the FunSearch method (Romera-Paredes et al., 2024) called F2 (for **F**eatures **F**unSearch). In F2, an LLM iteratively generates functions that compute features from the input domain, which are then evaluated by training a Random Forest (Breiman, 2001) predictor and extracting importance weights. These scores are used in-context to steer the LLM to generate features with high importance that are thus as predictive as possible. F2 is a *black-box* procedure with respect to the underlying classic predictor, leveraging an off-the-shelf Random Forest training method. We further introduce a *white-box* training procedure for LeaPR models called Dynamic ID3, or D-ID3, inspired by the classical ID3 algorithm for training decision trees. D-ID3 follows an ID3-like training loop where leaf nodes of a growing decision tree are selected and “split” by introducing two new leaves and a decision based on the value of a specific feature. While ID3 operates with pre-existing features, D-ID3 queries an LLM to propose new programmatic features on the fly that can help split the node under consideration. D-ID3 gives the LLM the decision path leading to the node, with few-shot examples of samples in that branch **in domains where inputs can be encoded as text**, and asks for novel features that are useful for that specific context. The vast prior knowledge of modern LLMs about domain-specific Python libraries makes both approaches highly applicable across diverse, high-dimensional input domains.

We evaluate both methods on two LLMs across three domains: chess position evaluation, image classification (MNIST (LeCun, 1998) and Fashion-MNIST (Xiao et al., 2017)), and text classification (AI vs. human text classification on the Ghostbuster dataset (Verma et al., 2024)). Across domains, LeaPR methods learn high-quality neural network-free representations that extract empirically useful features, often spanning tens of thousands of lines of Python code spread across hundreds of functions, all stitched together by decision trees. The learned features tend to be highly intuitive yet specific enough to lead to high-quality predictors. Our main contributions are:

- We propose jointly learning *programmatic features*, represented as LLM-generated functions from the input domain to the reals, together with decision tree predictors, thus obtaining fast, interpretable predictors.
- We introduce F2 and D-ID3, two algorithms for learning LeaPR models, where features are generated on demand during decision tree training.
- We evaluate LeaPR models on three domains: chess positions, images and text, showing comparable accuracy and often favorable data efficiency compared to baseline methods. We analyze the learned features and show how programmatic features can be useful for data exploration and for understanding model failures.

## 2 RELATED WORK

**Code generation with LLMs** Our work is enabled by the capability of LLMs to generate code under flexible task specifications, including natural language instructions and examples (Chen et al., 2021). This ability has been explored in a variety of domains, including using code as an intermediate tool for reasoning Chen et al. (2022) and as an interface for agents to interact with an external environment Lv et al. (2024). We exploit code generation for synthesizing *features* from arbitrary input domains, building on LLMs’ prior knowledge of a rich set of existing libraries. **Feature generation, both with and without LLMs, has been widely explored in tabular datasets (Ko et al., 2025; Zhang & Liu, 2024), whereas we explore domains with more complex, unstructured input (such as images and text), where current methods do not apply. We provide an extensive discussion of existing feature generation methods in Appendix D.**

**Black-box code evolution with LLMs** LLMs can also be used to *optimize* a black-box objective function (Lange et al., 2024). This approach was pioneered in FunSearch (Romera-Paredes et al., 2024), the method that inspires our F2 method (Section 3.1). In FunSearch, the LLM is given the type signature of a function to synthesize, and it outputs candidate functions which are scored with a metric unknown to the model. LLMs can then see past candidate function examples and their scores in following rounds (Romera-Paredes et al., 2024), and can mutate and evolve previous attempts to propose new, more successful ones. AlphaEvolve (Novikov et al., 2025) explored this idea with larger models and novel methods to ensure diverse candidates. Our work differs from FunSearch and AlphaEvolve in that LeaPR models use LLM-generated code simply to generate *features*, and not the full predictor end-to-end. Using classical models to utilize LLM-generated modules allows our method to scale and simultaneously employ thousands of features at once: for instance, yielding predictors at much a larger scale magnitude than what AlphaEvolve has been used to generate (e.g., up to 50k lines in total in the chess features used in our largest predictor).

**Learning programmatic world models** A modular approach for using LLM-synthesized code has been recently employed in PoE-World (Piriyakulkij et al., 2025) for the problem of learning *world models*. In PoE-World, LLMs generate a set of small predictors that individually explain a specific behavior of the environment; a separate model combines these predictors with weights learned by gradient descent. This scales better than monolithic code representations, which World-Coder introduced earlier (Tang et al., 2024). Like LeaPR models, PoE-World can generate world models with thousands of lines since LLMs do not have to generate complete programs, only small modules that can be composed scalably. We propose a similarly modular architecture (scaling to tens of thousands of lines) aimed toward the general problem of supervised learning. **The observation that synthesis of large programs scales better by learning libraries has a long history in program synthesis, including in DreamCoder (Ellis et al. (2021), with libraries constructed via symbolic compression) and Lilo (Grand et al. (2023), with LLM-generated abstractions).**

**Interpretability of neural networks** The widespread deployment of deep neural networks strongly motivated the research community to understand the behavior of neural models (Mu & Andreas, 2021), often through the lens of *mechanistic interpretability* Conmy et al. (2023). One key hypothesis is that neural networks have interpretable “circuits” that drive specific behaviors (Marks et al., 2025) or encode learned knowledge (Yao et al., 2025). **Prior work has also explored neural architectures that are more amenable to interpretation by construction (Doumbouya et al., 2025; Hewitt et al., 2023), rather than post-hoc. Several works also explore reshaping the prediction pipeline to lead to more interpretable decisions, including ViperGPT (Surís et al., 2023), Tree Prompting (Singh et al., 2023), and FM+V-IP (Chan et al., 2023). In all of these, however, neural networks are still involved at inference time, and even in cases involving code generation, such as in ViperGPT, the code represents a full predictor, and not a feature as in LeaPR.**

## 3 LEARNING PROGRAMMATIC REPRESENTATIONS

Classical predictors considered in supervised machine learning, like decision trees, are *intrinsically explainable* in the sense that their structure is simple enough to warrant human inspection: this is highly desirable when humans might want to trust (e.g. in high-stakes decision making) or learn from (e.g., someone learning to play chess) machine learning models. However, this simplicity

comes at a steep cost: the performance of these methods is notoriously impacted by the particular choice of features given as their input.

Consider learning a board evaluation decision tree for the game of chess: a predictor for the likelihood that the player with the white pieces will win. During inference, a decision tree tests the value of one input variable at a time. If given only the raw information available on the board (e.g. the piece lying on each of the 8x8 squares), each input variable alone carries negligible information about the overall position. Leaf nodes are forced to be invariant to all features not tested on the path from the root to that node; thus, if any decision is made before testing *all* 64 squares, critical information can always be missed (e.g., a queen on one of the unobserved squares). Thus, **these individual features are not informative enough** for decision trees to encode effective predictors.

On the other hand, if we instead train the model on a *single useful feature* computed from the board, like the *material difference* between players (given by a weighted sum of the pieces of each color still on the board), even a very shallow decision tree can make predictions that are significantly better than random, since a large material advantage is highly correlated with one’s chance of winning. Adding more informative features will progressively allow a decision tree learner to improve.

In this paper, our starting point is the insight that LLMs have two capabilities that allow them to serve as highly effective *feature engineers* for classical ML models. First, many useful features, such as the material balance in chess described above, can be implemented in a few lines of code, and modern LLMs are capable of flexibly generating code to accomplish a variety of tasks. Second, broad pre-training of LLMs encodes useful prior knowledge about a very wide range of domains (Bommasani et al., 2021), equipping them with strong priors over what kinds of features could be useful for making predictions across these domains. Our main hypothesis is that we might be able to train high-quality classical models by leveraging LLM-generated features at scale. Our goal here is thus to learn a hypothesis class that consists of (a) *programmatically features* and (b) decision tree predictors. We call these *Learned Programmatic Representation* models, or *LeaPR* models. The main challenge we now tackle is how to elicit predictive features from language models.

### 3.1 BLACK-BOX FEATURE GENERATION

Given a supervised dataset  $\mathcal{D}$  and any scoring function  $S_{\mathcal{D}}(f)$  that measures the quality of a proposed *feature*  $f : X \rightarrow \mathbb{R}$ , a simple methodology to obtain increasingly good features is to apply a FunSearch-style procedure, where an LLM is used to propose candidates to try to maximize  $S$ . Modern LLMs are capable of generating complex functions even for high-dimensional input domains, such as images and text, partly due to their ability to write code that uses existing human-written libraries for each of these domains. Thus, the main component needed for this approach to work is to answer: what makes  $f$  a good feature?

In the standard FunSearch setup (Romera-Paredes et al., 2024), the scoring function  $S$  evaluates candidates independently: proposals are *self-contained solutions* of the task. Our setup here is different: for the purpose of learning a predictor from *all* generated features at once, a new candidate  $f_k$  is valuable only to the extent that it contributes predictive information when taking into account the existing feature set  $f_{1:k-1}$ . Assuming that we keep track of the set of features generated so far and propose new features in a loop, a naïve adaptation of FunSearch would thus score a new feature  $f_k$  on its *added* predictive power once it has been proposed. For that, we could train a new predictor using  $f_{1:k}$  and compare its risk against the previous predictor trained on  $f_{1:k-1}$ . However, this runs into the issue that early features receive disproportionately high scores simply because their baselines (initial predictors based on few features) are severely limited. In practice, with this approach, the highest-scoring features are essentially fixed after the first few iterations, which is undesirable: we would like to detect and reward powerful features even if they appear late during training.

To overcome this problem, we *simultaneously score all existing features*  $f_{1:k}$  independently of the order in which they were proposed. Given a learned decision tree, prior work has proposed several metrics of *importance* of each input feature (e.g., measuring decrease in “impurity” in decision nodes that use a given feature, (Breiman, 2001)). Importance metrics only depend on the final learned predictor, and decision tree learning methods are order-invariant with respect to input features.

Algorithm 1 (Figure 2, left) shows *Features FunSearch* (F2), our representation learning algorithm based on this FunSearch-style approach but with features scored as a set. Specifically, F2 takes

**Algorithm 1: Features FunSearch (F2)**


---

**Input** : Language model,  $LM$ ,  
Supervised dataset  $D \in 2^{X \times Y}$

**Output**: List of features  $F$ , where  $f_i : X \rightarrow \mathbb{R}$

$F \leftarrow []$ ;

**for**  $iteration \in [1, \dots, T]$  **do**

$r f \leftarrow \text{TrainRandomForest}(D, F)$  ;

$imp \leftarrow \text{FeatureImportances}(r f)$  ;

$top\_k \leftarrow \text{TopKFeatures}(F, imp)$  ;

$r \leftarrow \text{RandomKFeatures}(F \setminus top\_k, imp)$  ;

$p \leftarrow \text{ProposeFeatures}(LM, top\_k, r)$  ;

$F.extend(p)$  ;

**end**

**return**  $F$  ;

---

**Algorithm 2: Dynamic ID3 (D-ID3)**


---

**Input** : Language model  $LM$ ,  
Supervised dataset  $D \in 2^{X \times Y}$

**Output**: List of features  $F$ , where  $f_i : X \rightarrow \mathbb{R}$

$T \leftarrow \text{Leaf}(D_{train})$  ;

**for**  $iteration \in [1, \dots, T]$  **do**

$l \leftarrow \arg \max_{l \in \text{Leaf}(n)} \text{TotalError}(T, n.data)$  ;

$p \leftarrow \text{ProposeFeatures}(LM, l.path\_to\_root)$  ;

$\langle f, t \rangle \leftarrow$   
 $\arg \min_{f \in F, t} \text{SplitError}(f, t, n.data)$  ;

$l.split(f, t, \{x \in n.data | f(x) < t\},$   
 $\{x \in n.data | f(x) > t\})$  ;

**end**

**return**  $\{n.splitting\_feature | n \in T \wedge \text{Internal}(n)\}$  ;

---

Figure 2: Two learning algorithms for LeaPR models. F2 (left) uses a FunSearch-style loop that attempts to evolve features that are *globally useful* to train a Random Forest predictor, as estimated by feature importances. D-ID3 (right) runs an ID3-style decision tree training loop and attempts to propose new features that are *locally useful* for splitting specific leaf nodes, attempting to minimize their impurity (e.g., variance in regression, or entropy in classification). **Full subroutine definitions in Appendix A.**

a supervised dataset and learns a programmatic representation — i.e. a set of feature functions  $f_i : X \rightarrow \mathbb{R}$ , represented as executable code. Like FunSearch, F2 iteratively uses an LLM to make batches of proposals conditioned on a sample of existing features, which are shown to the model along with their assigned scores — the LLM’s task is to propose new features that will be assigned high importance score in a newly trained Random Forest predictor. These scores are a *global estimate* of the predictive power of each feature in a predictor trained with all of them.

### 3.2 DYNAMIC SPLITTING

While F2 uses an underlying decision tree learner as a black-box, the insight that LLMs can be used to generate features on demand can serve as the basis for designing the decision tree learner itself. Recall that during inference in a decision tree, we start at the root and repeatedly follow the “decisions” associated with each node until we reach a leaf. Each such decision consists of testing the value of a particular feature: if this node splits on feature  $f_k$ , we compare  $f_k(x)$  with a threshold  $t$  learned during training. If  $f_k(x) < t$ , the node recursively returns the prediction made by its left child (or right child if  $f_k(x) \geq t$ ). Leaf nodes return a fixed prediction defined during training, e.g., the most common class label (classification), or average value (regression) for training points that fall on that leaf. For training, classical decision tree learning algorithms (e.g., ID3 or CART (Quinlan, 1986)) start with a single node and repeatedly improve the current decision tree predictor by (a) choosing a leaf node and (b) partitioning it into a new decision node with two new leaves as its children. Partitioning searches for a feature and comparison threshold that minimize the “impurity” (e.g., variance in the continuous case, or entropy of class labels in classification) of data points falling on both sides. For instance, in classification, the best-case scenario would be to find a partition where all training data points falling on each new leaf belong to the same class.

However, the ability of classical algorithms to find good splits is limited by the predictive power of preexisting dataset features. Here, we revisit this recursive splitting strategy considering that we can attempt to generate new features *on demand* for the purpose of successfully partitioning a particular leaf node. When we decide to split a leaf, we have significant local context aside from global dataset information: in particular, we know the specific path of decisions that leads to that leaf, and we have a corresponding set of training examples, with their labels, falling onto that node. To be *locally useful*, a feature only needs to help distinguish between examples in that set. Indeed, for informing the proposals of potentially useful features, we can even leverage the ability of LLMs to perform inductive reasoning, by presenting actual examples (if possible in text), along with their labels, in the model’s context: its objective, then, is to propose a feature that would explain the variation in the labels between those examples and others that reach the same leaf.

Algorithm 2, Dynamic ID3 (D-ID3), realizes this idea. In each iteration, D-ID3 selects the current leaf in the tree that accounts for the largest portion of training error (e.g. number of misclassified training examples). D-ID3 then generates new candidate features with an LLM on the fly to split that particular leaf. In modalities where we can easily represent examples in text, the LLM receives a sample of examples and their labels that fall in this branch (in our experiments, this only excludes image classification, where we only show a sample of image *class labels* in the prompt). D-ID3 considers these features, as well as all candidate features generated for ancestor nodes, and finds the best split for this leaf according to a user-defined impurity metric (all metrics available for classical methods are also possible here). This process repeats for a number of iterations. At the end, like F2, D-ID3 returns a learned *representation*: the set of programmatic features for the input domain that were used in the resulting decision tree. We note several practical considerations for both F2 and D-ID3, as well as other implementation details, in Appendix A.2.

## 4 EXPERIMENTS

We now evaluate programmatic representation models on three tasks with complex input domains where the standard practice is to train neural networks: chess position evaluation (given a chess board, predict the probability that White wins), image classification on MNIST (LeCun, 1998) and Fashion-MNIST (Xiao et al., 2017), and text classification (detecting whether a piece of text is human- or AI-generated) on Ghostbuster (Verma et al., 2024). In all domains, we compare LeaPR against standard neural network baselines; additionally, in chess and image classification, we also include the Random Forest baseline where we feed a simple “raw” encoding of the input (the piece in each square for chess boards, or pixel values for images). **We also compare against a baseline FunSearch method that learns the full predictor with code evolution, as well as a Program of Thoughts baselines – both of these using GPT-5.1 as their underlying language model (details in Appendix B.3). These baselines can also use the rich prior knowledge of language models and their code generation capabilities, but without the modularity that the LeaPR paradigm provides.** We discuss the features our methods learn in each domain, and finally conduct a case study debugging a classifier that has learned to rely on a spurious feature in the Waterbird dataset (Sagawa et al., 2020).

We run F2 and D-ID3 using two OpenAI models, for a total of 4 LeaPR models per task: GPT 4o-mini gpt-4o-mini-2024-07-21 (Hurst et al., 2024) and GPT 5-mini gpt-5-mini-2025-08-07 (OpenAI, 2025). We run both methods so that they output a maximum of 1000 features — this means using 1000 iterations of D-ID3, and 100 iterations of F2 with a proposal batch size of 10 features in each call. We sometimes end with fewer than 1000 features because we discard features that fail validation (see Section A.2) Using the features learned by either algorithm, we then train a Random Forest model using the standard Scikit-Learn (Pedregosa et al., 2011) implementation, with 500 trees and a maximum depth of 50.

### 4.1 CHESS POSITION EVALUATION

First, we train models on the regression task of state-value prediction in the game of chess: given the board position, predict the win probability for each player. We use a publicly available dataset of games from the Lichess online platform (Lichess.org), and hold out 1000 random board positions for evaluation. The dataset comes with state values estimated by Stockfish (Romstad et al., 2008), the strongest publicly available chess engine. We use Stockfish’s prediction value as the ground truth (Stockfish outputs values in “centipawns”, which we convert to win percentages using the standard formula used by Lichess and other prior work). To represent and manipulate chess boards, we use the popular `python-chess` library, with a standard API that facilitates iterating through the board, locating pieces, generating available moves, and testing for various pieces of game state (e.g., a player’s turn, whether the current player is in check, etc). Our prompts contain a short listing of the main API classes, methods, and functions available in the library. The models are instructed to generate features that receive an argument of type `chess.Board` and return a `float` value. We provide full prompts in Appendix F.

**Transformer baseline.** As a neural baseline, we train the 270M parameter Transformer architecture proposed in Ruoss et al. (2024) (their largest model) to predict the discretized win-probability for White (128 buckets) given a position encoded in the standard FEN format. Ruoss et al. (2024) compared models that predict both state values and state-action values; when controlled for the num-

ber of data points, models trained on state-value slightly outperformed when playing games against each other. Their strongest model (trained on 15.3B state-action values) achieved grandmaster-level play. Due to computational constraints, we reproduce their training of a state-value prediction Transformer run only up to 50M data points (10x less than their total state-value dataset of 500M data points; though we approximately match their number of epochs over the training data at 2.5).

Table 1 summarizes the results for this regression task. Here, we show both root mean square error (RMSE) and Pearson correlation ( $\rho$ ) between model predictions and Stockfish’s estimate. Our LeaPR models, trained on 200k board positions, compare favorably with the Transformer predictor trained on 250x more data. In contrast, as expected, Random Forests trained on the raw board struggle. LeaPR models benefit from the significant prior knowledge that LLMs have about useful chess concepts. We see basic features such as one that “Calculates the total piece value of both sides” (the model’s own function documentation string) proposed and implemented by GPT 4o-mini in 6 lines of code early during training, as well as significantly more complex, specific features such as “Pawn promotion pressure: sum over pawns of  $1/(1+\text{steps\_to\_promotion})$  weighted by being passed (white minus black). Encourages advanced, passed pawns.”, implemented by GPT 5-mini late in the D-ID3 run (with 51 lines of code). Generally, D-ID3 features appear to become more specific as training progresses, likely because the LLM is asked to distinguish only a subset of board positions that already share many similarities (due to falling on a specific leaf node), yielding better models than F2 even with the same number of total features.

We also compare models in terms of Top-1 move accuracy compared to Stockfish at its maximum strength. Since we only estimate state-values, to use our predictors we select the move that leads to the best successor value from the point of view of the current player (i.e., the move leading to the highest or lowest win-probability for White depending on who plays), and measure how often this matches Stockfish’s top move. Interestingly, we find that regression performance is not necessarily predictive of move accuracy: LeaPR models trained with GPT 4o-mini are significantly worse when used to select moves. Our best action predictors achieve non-trivial move selection accuracy: whereas random performance for this task is 11.4%, the D-ID3 model trained with GPT 5-mini predicts the top Stockfish move in 33.5% of the cases, with the Transformer baseline underperforming at 30.3%. Ruoss et al. (2024) trained this same state-value model with up to 500M data points and managed to achieve a move accuracy of 58.5%, showing that the Transformer keeps improving for much longer. Their most accurate model for action prediction is trained to directly predict action-values from 15B training data points, achieving an accuracy of 63.5% and a grandmaster-level ELO rating when playing with humans online. Although there remains a significant gap between their best results achieved with a Transformer and what we demonstrate here, LeaPR models still get surprisingly far in this challenging regression task. If the scalability challenges associated with LeaPR models can be overcome (e.g., we see negligible benefits from training Random Forests beyond 200k training data points), we might be able to obtain chess policies that

Predictor	Training Size	RMSE	$\rho$	Acc.
Random policy	0			11.4%
Transformer	$5 \times 10^7$	.161	.795	30.3%
Transformer (Ruoss et al., 2024)	$5 \times 10^8$			58.5%
Random Forest (raw board)	200k	.248	.306	14.5%
PoT + GPT-5.1	200k	.316	.611	23.10%
FunSearch + GPT-5.1	200k	.202	.652	27.4%
LeaPR F2 + GPT 5-mini	200k	.169	.762	31.4%
F2 + GPT 4o-mini	200k	.163	.783	16.7%
D-ID3 + GPT 5-mini	200k	.160	.789	33.5%
D-ID3 + GPT 4o-mini	200k	.156	.806	17.2%

Table 1: Performance in state-value prediction models in chess positions from Lichess. We train the 270M-parameter Transformer from Ruoss et al. (2024) with up to 50M data points, and report their results for their full run on 10x more data.

Predictor	MNIST	Fashion
ResNet-50	98.71%	89.54%
EfficientNetV2	98.8%	90.94%
Random Forest (raw pixels)	95.6%	88.29%
PoT + GPT-5.1	23.75%	22.39%
FunSearch + GPT-5.1	25.80%	25.96%
LeaPR F2 + GPT 5 mini	92.54%	85.77%
F2 + GPT 4o mini	89.26%	80.26%
D-ID3 + GPT 5 mini	96.91%	88.51%
D-ID3 + GPT 4o mini	93.71%	83.80%

Table 2: Top-1 accuracy on image classification on MNIST and Fashion-MNIST.

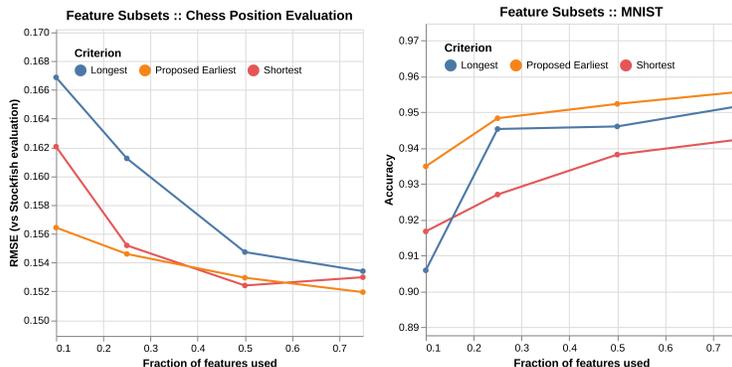


Figure 3: Random Forest predictor performance when trained on subsets of varying sizes of D-ID3-generated features (see Appendix D for remaining domains). Consistently, more D-ID3 features tends to improve the predictor’s performance, and the most impactful features are found early.

not only play at a high level but can also “explain” their moves — a feature that no existing chess engine possesses.

#### 4.2 IMAGE CLASSIFICATION

We now evaluate LeaPR models on two image classification datasets: MNIST (handwritten digit classification) from LeCun (1998) and Fashion-MNIST (grayscale fashion products classification) from Xiao et al. (2017). We train standard ResNet-50 and EfficientNetV2 baselines to convergence on the same datasets (training details in B.2). Unlike in Section 4.1, D-ID3 does not add *images* in the prompt when calling the LLM, but only a textual description of class labels of a set of training examples belonging to the leaf being split (e.g., “digit 0” in MNIST, or “T-shirt” in Fashion-MNIST). This is a significant limitation for this domain, since the features need to rely solely on the LLM’s prior knowledge about what the described objects might look like and hypotheses about how they will be detectable in a small grayscale image. Still, the best LeaPR models in this task achieve comparable accuracy to the neural baselines, even without the ability to construct features by directly observing the training data. Again, especially with D-ID3 we observe specific features that attempt to distinguish between particular classes, such as “*Count of ink endpoints*”: *ink pixels with only one ink neighbor (8-connected). Loops like 0/8 have few endpoints; open strokes like 5 have endpoints*” being the feature with the best split (thus selected) in a leaf where the majority classes were 8 and 5. GPT 5-mini correctly implemented the above feature in 30 lines of Python using `numpy`. Though MNIST and Fashion-MNIST generally only serve as “sanity checks” for computer vision models (with even Random Forests trained on the raw pixels performing near the neural baselines, given the small image and training set sizes), we believe that this result presents an encouraging signal towards the ability of LeaPR models to achieve comparable accuracy while proposing simple and interpretable programmatic features.

#### 4.3 TEXT CLASSIFICATION

We now evaluate LeaPR models on a binary text classification task: detecting whether the input was written by an LLM or a human. We use the Ghostbuster dataset (Verma et al., 2024), which contains a collection of student essays, creative writing, and news articles written by humans and by ChatGPT and Claude (Anthropic, 2024) given the same or similar prompts. In Table 3 we compare LeaPR models the Ghostbuster model and other neural baselines reported in Verma et al. (2024) for the “in distribution” setting with all domains combined (we omit DetectGPT, which achieves a low F1 score of 51.6% due to having been trained to detect another LLM). Here, LeaPR models are the only neural network-free predictors. Still, our models perform competitively when evaluated in F1 score: they outperform or match all baselines, with the best LeaPR model (with features obtained by D-ID3 with GPT 5-Mini) closely matching Ghostbuster (98.8 vs 99.0 in F1 score).

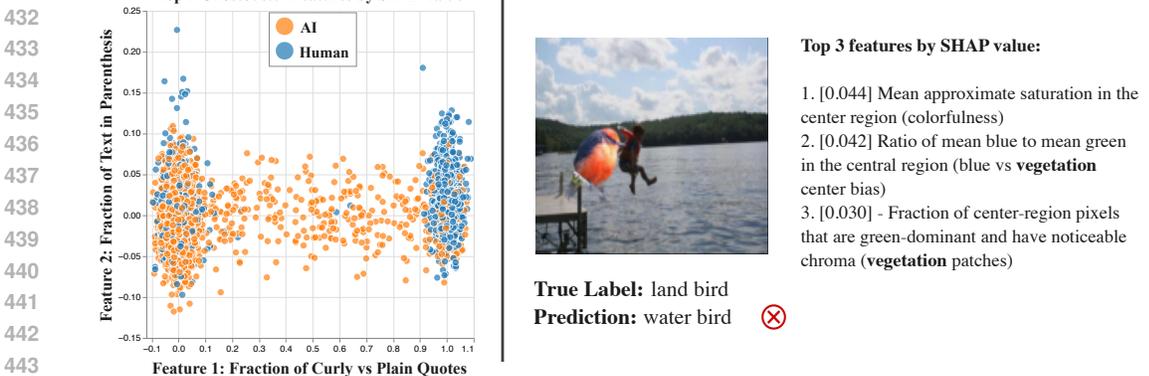


Figure 4: (Left) Distribution of the top 2 LeaPR-learned features with highest SHAP values on the Ghostbuster dataset (gaussian jitter added to aid visualization). On Feature 1 (which measures fraction of “curly”, or typographic characters, versus plain ASCII quotes), human text tends to cluster at the extremes, while AI-generated text features mid-range values. Feature 2 shows high values primarily for human-written text. (Right) A land bird misclassified by a LeaPR model on the Waterbird dataset; the top SHAP-valued features for this example show a clear reliance on the background.

#### 4.4 IMPACT OF FEATURE COMPLEXITY AND QUANTITY

We now analyze the impact of both the quantity of features generated, as well as their complexity, in the quality of predictors we obtain. Here we look at D-ID3 paired with GPT-5 Mini, our best performing combination across domains. Ideally, we would like D-ID3 to steadily improve the trained predictor with a larger budget, rather than saturating early. Moreover, we would expect to obtain the most impactful features earlier, so as to make the most out of a smaller run budget.

Figure 3 shows the Random Forest predictor’s performance on MNIST and Chess when trained on varying subsets of the learned features selected according to three criteria: features generated *earliest*, the longest (in lines of code), and the shortest. (Results for Fashion-MNIST and Ghostbuster, showing similar trends, are in Appendix D. Here, lines of code serves as a crude proxy for feature complexity. Across domains, we observe that (1) using more of the D-ID3-proposed features consistently improves performance, and (2) D-ID3 tends to find the most impactful features earlier: using the first features is generally better than using either the simplest (less lines of code) or most complex. Whether feature complexity is generally positive depends on the domain: for instance, in chess, simpler features seem to have a higher impact, which does not hold in the image domains. Overall, however, we see a positive scaling trend with respect to features: up to the scale of our runs, we still see improving performance the larger the budget we allow D-ID3 to use.

Predictor	F1
Perplexity only	81.5
GPTZero	93.1
RoBERTa	98.1
Ghostbuster	99.0
PoT + GPT-5.1	56.9
FunSearch + GPT-5.1	74.2
LeaPR F2 + GPT 5-mini	97.7
F2 + GPT 4o-mini	98.6
D-ID3 + GPT 5-mini	98.8
D-ID3 + GPT 4o-mini	98.6

Table 3: F1 score on Ghostbuster: classifying text as human or AI written (Verma et al., 2024).

#### 4.5 CASE STUDIES: UNDERSTANDING FEATURES AND MODEL PREDICTIONS

The interpretable representations that LeaPR models learn have potential uses beyond their predictive power. We now describe two case studies using SHAP values (Lundberg & Lee, 2017) as a lens into the patterns that LeaPR models find in their training data, and into why a particular prediction — especially if erroneous — was made. SHAP values are a metric of feature importance for a model that can be applied both at a dataset-level or to a particular prediction; we refer to Lundberg & Lee (2017) for details. This is an especially compelling tool for understanding LeaPR models given that their features already come with natural language descriptions, in the form of documentation strings.

486 First, we compute SHAP values in the Ghostbuster training set to understand what features the  
 487 LeaPR model trained with D-ID3 and GPT 5-mini has learned to use to identify AI-generated text.  
 488 Sorting features by their SHAP values on a sample of 150 training examples, we find that two of  
 489 the top-3 most important features for the model are (1) the “Fraction of quotation marks that are  
 490 curly/typographic quotes (e.g., ‘ ’ “ ”) vs plain ASCII quotes, indicating published/edited text” and  
 491 the (2) “*Proportion of characters that lie inside parentheses (measures parenthetical/planned con-*  
 492 *tent like ”(50 words)”*” (the other top-3 feature also looks for kinds of quotation characters and is  
 493 strongly correlated with the first). Together, these two features already capture distinctive patterns  
 494 in human- and AI-written samples in Ghostbuster. Figure 4 (left) shows training samples projected  
 495 on these two features, with small Gaussian jitter added to both coordinates to aid visualization. For  
 496 feature 1, human-generated text either has value 1.0, meaning *all* quote characters are typographi-  
 497 cal, or “curly” quotes, or 0.0, meaning *all* quotes in the text are plain ASCII quotes (this could, for  
 498 instance, reflect default settings in the user device, with using curly quotes being much more fre-  
 499 quent). This is in stark contrast with AI-generated text, which often mixes both kinds of characters  
 500 in the same text: AI-generated text displays the full range from 0 to 1 in this feature. For Feature 2,  
 501 humans seem much more likely to wrap a significant fraction of the text in parentheses: almost all  
 502 samples with value over 0.1 in this feature (over 10% in parenthesis) were human-written. LeaPR  
 503 models allowed us to quickly discover these patterns without having to formulate specific a priori  
 504 assumptions: combined, our models contain thousands of automatically generated domain-relevant  
 features, thus serving as a highly useful tool for data understanding.

505 Finally, we conduct a case study on the Waterbird dataset (Sagawa et al., 2020) showing how pro-  
 506 grammatic features can serve to debug model failures. This dataset contains images with two classes  
 507 of birds: land birds and water birds — which are the two target classification labels. However, a  
 508 trained classifier might learn instead to rely on the background, rather than use features of the bird  
 509 itself. The dataset contains a subset of land birds placed on water backgrounds, and vice-versa:  
 510 typically, classification accuracy drops significantly across groups when models learn to predict bird  
 511 classes based on the (spuriously correlated) background. When we train a LeaPR model with F2  
 512 and GPT 5-mini on Waterbird, it achieves 100% validation accuracy when evaluated on *land birds*  
 513 *on land background*, but it drops to 84% when evaluated on *land birds on water background*. Again,  
 514 SHAP values can help us understand why this happens in a particular case. Figure 4 (right) shows  
 515 the first validation example of a land bird on water background that is misclassified. When we show  
 516 the top 3 features by their SHAP values, the second and third features explicitly indicate that their  
 517 goal is to detect *vegetation* — a spurious feature. We can indeed find several examples of features  
 518 that *attempt* to characterize the bird, such as “Fraction of warm (red-dominant) pixels in the center  
 519 region (bird color cue)”, that are however ignored by the trained predictor. This example shows how  
 520 LeaPR models can make their failures transparent. Since model failures often reflect properties of  
 521 the training data, we believe that LeaPR can serve as debugging tool for both models and datasets,  
 applicable to a wide range of domains.

## 522 5 LIMITATIONS AND CONCLUSION

523 We introduced *Learned Programmatic Representation* models, a class of neural network-free models  
 524 that combine programmatic LLM-generated features and decision tree predictors. We experiment  
 525 with chess boards, images and text, seeing encouraging initial results exploring this class of models.  
 526 Our learned features tend to be easy to interpret: we explore various examples across each of these  
 527 domains, including using SHAP values to explain individual predictions and overall learned models.  
 528

529 However, several limitations remain to be tackled by future work. Our methods do not learn deep  
 530 hierarchical features: each feature is directly computed from the input. Learning deep feature hierar-  
 531 chies is the main advantage of training *deep* neural networks, and we believe that to be an important  
 532 capability for LeaPR models to perform well in more complex domains. Moreover, our experiments  
 533 were done at a small scale, and there are scalability challenges — both in representation learning  
 534 as well as in training predictors — to be overcome to achieve competitive performance in data-rich  
 535 domains, like chess, where neural networks improve predictably with more data and compute.

536 Despite these limitations, we find our results encouraging for further exploring novel learning  
 537 paradigms that yield interpretable models *by construction*. With a rapidly advancing AI toolbox,  
 538 future tools might allow us to learn interpretable models just as easily as we can train neural net-  
 539 works today, with little to no sacrifice in quality. Overcoming the limitations in the LeaPR paradigm  
 can thus be a path to incorporating interpretable models into the practical AI toolbox.

## REFERENCES

- 540  
541 Anthropic. Claude, 2024. URL <https://www.anthropic.com>.
- 542  
543 Yoshua Bengio and Yann LeCun. Scaling learning algorithms towards AI. In *Large Scale Kernel*  
544 *Machines*. MIT Press, 2007.
- 545  
546 Rishi Bommasani, Drew A Hudson, Ehsan Adeli, Russ Altman, Simran Arora, Sydney von Arx,  
547 Michael S Bernstein, Jeannette Bohg, Antoine Bosselut, Emma Brunskill, et al. On the opportu-  
548 nities and risks of foundation models. *arXiv e-prints*, pp. arXiv–2108, 2021.
- 549  
550 Samuel R. Bowman, Gabor Angeli, Christopher Potts, and Christopher D. Manning. A large an-  
551 notated corpus for learning natural language inference, 2015. URL <https://arxiv.org/abs/1508.05326>.
- 552  
553 Leo Breiman. Random forests. *Machine learning*, 45(1):5–32, 2001.
- 554  
555 Kwan Ho Ryan Chan, Aditya Chattopadhyay, Benjamin David Haeffele, and Rene Vidal. Variational  
556 information pursuit with large language and multimodal models for interpretable predic-  
557 tions. *arXiv preprint arXiv:2308.12562*, 2023.
- 558  
559 Mark Chen, Jerry Tworek, Heewoo Jun, Qiming Yuan, Henrique Ponde De Oliveira Pinto, Jared  
560 Kaplan, Harri Edwards, Yuri Burda, Nicholas Joseph, Greg Brockman, et al. Evaluating large  
561 language models trained on code. *arXiv preprint arXiv:2107.03374*, 2021.
- 562  
563 Wenhui Chen, Xueguang Ma, Xinyi Wang, and William W Cohen. Program of thoughts prompt-  
564 ing: Disentangling computation from reasoning for numerical reasoning tasks. *Transactions on*  
565 *Machine Learning Research*, 2022.
- 566  
567 Isaac Cheng and Chico Camargo. Machine learning to study patterns in chess games. Master’s  
568 thesis, University of Exeter, 2023.
- 569  
570 Arthur Conmy, Augustine Mavor-Parker, Aengus Lynch, Stefan Heimersheim, and Adrià Garriga-  
571 Alonso. Towards automated circuit discovery for mechanistic interpretability. *Advances in Neural*  
572 *Information Processing Systems*, 36:16318–16352, 2023.
- 573  
574 Alexandru Damian, Jason Lee, and Mahdi Soltanolkotabi. Neural networks can learn representations  
575 with gradient descent. In *Conference on Learning Theory*, pp. 5413–5452. PMLR, 2022.
- 576  
577 Guozhu Dong and Huan Liu. *Feature engineering for machine learning and data analytics*. CRC  
578 press, 2018.
- 579  
580 Finale Doshi-Velez and Been Kim. Towards a rigorous science of interpretable machine learning.  
581 *arXiv preprint arXiv:1702.08608*, 2017.
- 582  
583 Moussa Koulako Bala Doumbouya, Dan Jurafsky, and Christopher D. Manning. Tversky neural  
584 networks: Psychologically plausible deep learning with differentiable tversky similarity. 2025.  
585 URL <https://arxiv.org/abs/2506.11035>.
- 586  
587 Kevin Ellis, Catherine Wong, Maxwell Nye, Mathias Sablé-Meyer, Lucas Morales, Luke Hewitt,  
588 Luc Cary, Armando Solar-Lezama, and Joshua B Tenenbaum. Dreamcoder: Bootstrapping in-  
589 ductive program synthesis with wake-sleep library learning. In *Proceedings of the 42nd acm*  
590 *sigplan international conference on programming language design and implementation*, pp. 835–  
591 850, 2021.
- 592  
593 Gabriel Grand, Lionel Wong, Maddy Bowers, Theo X Olausson, Muxin Liu, Joshua B Tenenbaum,  
and Jacob Andreas. Lilo: Learning interpretable libraries by compressing and documenting code.  
*arXiv preprint arXiv:2310.19791*, 2023.
- John Hewitt, John Thickstun, Christopher D Manning, and Percy Liang. Backpack language models.  
*arXiv preprint arXiv:2305.16765*, 2023.
- Noah Hollmann, Samuel Müller, and Frank Hutter. Large language models for automated data  
science: Introducing caafe for context-aware automated feature engineering, 2023. URL <https://arxiv.org/abs/2305.03403>.

- 594 Franziska Horn, Robert Pack, and Michael Rieger. The autofeat python library for automated feature  
595 engineering and selection, 2020. URL <https://arxiv.org/abs/1901.07329>.  
596
- 597 Robert Huben, Hoagy Cunningham, Logan Riggs Smith, Aidan Ewart, and Lee Sharkey. Sparse  
598 autoencoders find highly interpretable features in language models. In *The Twelfth International  
599 Conference on Learning Representations*, 2023.
- 600 Aaron Hurst, Adam Lerer, Adam P Goucher, Adam Perelman, Aditya Ramesh, Aidan Clark, AJ Os-  
601 trow, Akila Welihinda, Alan Hayes, Alec Radford, et al. Gpt-4o system card. *arXiv preprint  
602 arXiv:2410.21276*, 2024.
- 603 Jeonghyun Ko, Gyeongyun Park, Donghoon Lee, and Kyunam Lee. Ferg-llm: Feature engineering  
604 by reason generation large language models. In *Findings of the Association for Computational  
605 Linguistics: NAACL 2025*, pp. 4211–4228, 2025.
- 606
- 607 Robert Lange, Yingtao Tian, and Yujin Tang. Large language models as evolution strategies. In  
608 *Proceedings of the Genetic and Evolutionary Computation Conference Companion*, pp. 579–582,  
609 2024.
- 610 Yann LeCun. The mnist database of handwritten digits. [http://yann.lecun.com/exdb/  
611 mnist/](http://yann.lecun.com/exdb/mnist/), 1998.
- 612
- 613 Lichess.org. Lichess. <https://lichess.org/about>.
- 614 Scott M Lundberg and Su-In Lee. A unified approach to interpreting model predictions. *Advances  
615 in neural information processing systems*, 30, 2017.
- 616
- 617 Weijie Lv, Xuan Xia, and Sheng-Jun Huang. Codeact: Code adaptive compute-efficient tuning  
618 framework for code llms. *arXiv preprint arXiv:2408.02193*, 2024.
- 619 Samuel Marks, Can Rager, Eric J. Michaud, Yonatan Belinkov, David Bau, and Aaron Mueller.  
620 Sparse feature circuits: Discovering and editing interpretable causal graphs in language models.  
621 2025. URL <https://arxiv.org/abs/2403.19647>.
- 622
- 623 Ryszard Stanislaw Michalski, Jaime Guillermo Carbonell, and Tom M Mitchell. *Machine learning:  
624 An artificial intelligence approach*. Springer Science & Business Media, 2013.
- 625
- 626 Jesse Mu and Jacob Andreas. Compositional explanations of neurons. 2021. URL [https://  
arxiv.org/abs/2006.14032](https://arxiv.org/abs/2006.14032).
- 627
- 628 Jaehyun Nam, Kyuyoung Kim, Seunghyuk Oh, Jihoon Tack, Jaehyung Kim, and Jinwoo Shin. Op-  
629 timized feature generation for tabular data via llms with decision tree reasoning, 2024. URL  
<https://arxiv.org/abs/2406.08527>.
- 630
- 631 Alexander Novikov, Ngân Vũ, Marvin Eisenberger, Emilien Dupont, Po-Sen Huang, Adam Zsolt  
632 Wagner, Sergey Shirobokov, Borislav Kozlovskii, Francisco JR Ruiz, Abbas Mehrabian,  
633 et al. Alphaevolve: A coding agent for scientific and algorithmic discovery. *arXiv preprint  
634 arXiv:2506.13131*, 2025.
- 635
- 636 OpenAI. GPT-5 System Card. <https://cdn.openai.com/gpt-5-system-card.pdf>,  
August 2025.
- 637
- 638 F. Pedregosa, G. Varoquaux, A. Gramfort, V. Michel, B. Thirion, O. Grisel, M. Blondel, P. Pretten-  
639 hofer, R. Weiss, V. Dubourg, J. Vanderplas, A. Passos, D. Cournapeau, M. Brucher, M. Perrot, and  
640 E. Duchesnay. Scikit-learn: Machine learning in Python. *Journal of Machine Learning Research*,  
12:2825–2830, 2011.
- 641
- 642 Karol J. Piczak. ESC: Dataset for Environmental Sound Classification. In *Proceedings of the 23rd  
643 Annual ACM Conference on Multimedia*, pp. 1015–1018. ACM Press. ISBN 978-1-4503-3459-  
644 4. doi: 10.1145/2733373.2806390. URL [http://dl.acm.org/citation.cfm?doid=  
2733373.2806390](http://dl.acm.org/citation.cfm?doid=2733373.2806390).
- 645
- 646 Wasu Top Piriyakulkij, Yichao Liang, Hao Tang, Adrian Weller, Marta Kryven, and Kevin Ellis.  
647 Poe-world: Compositional world modeling with products of programmatic experts. *Advances in  
Neural Information Processing Systems (to appear)*, 2025.

- 648 J. Ross Quinlan. Induction of decision trees. *Machine learning*, 1(1):81–106, 1986.
- 649
- 650 Bernardino Romera-Paredes, Mohammadamin Barekatin, Alexander Novikov, Matej Balog,  
651 M Pawan Kumar, Emilien Dupont, Francisco JR Ruiz, Jordan S Ellenberg, Pengming Wang,  
652 Omar Fawzi, et al. Mathematical discoveries from program search with large language models.  
653 *Nature*, 625(7995):468–475, 2024.
- 654 Tord Romstad, Marco Costalba, Joonas Kiiski, Gary Linscott, Yu Nasu, Motohiro Isozaki, Hisayori  
655 Noda, et al. Stockfish, 2008. URL <https://stockfishchess.org>.
- 656
- 657 Anian Ruoss, Grégoire Delétang, Sourabh Medapati, Jordi Grau-Moya, Li K Wenliang, Elliot Catt,  
658 John Reid, Cannada A Lewis, Joel Veness, and Tim Genewein. Amortized planning with large-  
659 scale transformers: A case study on chess. *Advances in Neural Information Processing Systems*,  
660 37:65765–65790, 2024.
- 661 Shiori Sagawa, Pang Wei Koh, Tatsunori B Hashimoto, and Percy Liang. Distributionally robust  
662 neural networks for group shifts: On the importance of regularization for worst-case generaliza-  
663 tion. *International Conference on Learning Representations*, 2020.
- 664 Chandan Singh, John Morris, Alexander M Rush, Jianfeng Gao, and Yuntian Deng. Tree prompting:  
665 Efficient task adaptation without fine-tuning. In *Proceedings of the 2023 Conference on Empirical*  
666 *Methods in Natural Language Processing*, pp. 6253–6267, 2023.
- 667
- 668 Dídac Surís, Sachit Menon, and Carl Vondrick. Vipergpt: Visual inference via python execution  
669 for reasoning. In *Proceedings of the IEEE/CVF international conference on computer vision*, pp.  
670 11888–11898, 2023.
- 671 Hao Tang, Darren Key, and Kevin Ellis. Worldcoder, a model-based llm agent: Building world  
672 models by writing code and interacting with the environment. *Advances in Neural Information*  
673 *Processing Systems*, 37:70148–70212, 2024.
- 674
- 675 Vivek Verma, Eve Fleisig, Nicholas Tomlin, and Dan Klein. Ghostbuster: Detecting text ghostwrit-  
676 ten by large language models. In *Proceedings of the 2024 Conference of the North American*  
677 *Chapter of the Association for Computational Linguistics: Human Language Technologies (Vol-*  
678 *ume 1: Long Papers)*, pp. 1702–1717, 2024.
- 679 Jindong Wang, Cuiling Lan, Chang Liu, Yidong Ouyang, Tao Qin, Wang Lu, Yiqiang Chen, Wenjun  
680 Zeng, and Philip S. Yu. Generalizing to unseen domains: A survey on domain generalization.  
681 *IEEE Transactions on Knowledge and Data Engineering*, 35(8):8052–8072, 2023. doi: 10.1109/  
682 TKDE.2022.3178128.
- 683 Han Xiao, Kashif Rasul, and Roland Vollgraf. Fashion-mnist: a novel image dataset for benchmark-  
684 ing machine learning algorithms. *arXiv preprint arXiv:1708.07747*, 2017.
- 685
- 686 Yunzhi Yao, Ningyu Zhang, Zekun Xi, Mengru Wang, Ziwen Xu, Shumin Deng, and Huajun Chen.  
687 Knowledge circuits in pretrained transformers. 2025. URL [https://arxiv.org/abs/  
688 2405.17969](https://arxiv.org/abs/2405.17969).
- 689 Tianping Zhang, Zheyu Zhang, Zhiyuan Fan, Haoyan Luo, Fengyuan Liu, Qian Liu, Wei Cao,  
690 and Jian Li. Openfe: Automated feature generation with expert-level performance, 2023. URL  
691 <https://arxiv.org/abs/2211.12507>.
- 692
- 693 Xinhao Zhang and Kunpeng Liu. Tifg: Text-informed feature generation with large language mod-  
694 els. In *2024 IEEE International Conference on Big Data (BigData)*, pp. 8256–8258. IEEE, 2024.
- 695
- 696
- 697
- 698
- 699
- 700
- 701

## A IMPLEMENTATION DETAILS FOR F2 AND D-ID3

In this section, we first expand on the description of Algorithms 1 and 2 by detailing the subroutines they rely on, and then discuss practical implementation details of both.

### A.1 ALGORITHM SUBROUTINES

**TrainRandomForest**( $D, F$ ): Computes feature vectors for all examples in  $D$  by evaluating each  $f \in F$  (parallelized via `ProcessPoolExecutor`), then trains a `scikit-learn RandomForestRegressor` or `RandomForestClassifier`.

**FeatureImportances**( $rf$ ): Returns  $rf.feature\_importances$ , the mean decrease in impurity for each feature.

**TopKFeatures**( $F, imp$ ): Sorts pairs  $(f, imp[f])$  for  $f \in F$  by importance scores (descending) and returns the top  $k$  pairs as tuples.

**RandomKFeatures**( $F, imp$ ): Uniformly samples  $k$  features from  $F$  using `random.sample` (no importance weighting).

**ProposeFeatures**( $LM, context$ ): Constructs a domain-specific prompt with context (feature examples + importances for F2; decision path + example datapoints for D-ID3), calls  $LM$  API, parses response to extract Python function definitions, validates each in separate process on test subset (10k examples), and returns only those that execute without errors/timeouts/non-finite values.

**TotalError**( $T, D$ ): For leaf node, returns domain-specific leaf error (MAE for regression or error rate for classification); for internal node, returns weighted average of children’s total errors.

**SplitError**( $f, t, D$ ): Partitions  $D$  by threshold  $t$  on feature  $f$ , computes impurity of each partition (variance for regression via `RunningMedianAbs`, entropy for classification via histogram), and returns weighted sum.

### A.2 PRACTICAL CONSIDERATIONS

Both F2 and D-ID3 run for a user-specified number of iterations; this number is exactly equal to the number of LLM calls that the algorithm will perform, allowing users to budget for LLM usage. In a sense, both algorithms are “anytime algorithms” — they can always return their latest set of learned features. The algorithms return a *representation*, rather than a predictor (e.g. the decision tree constructed by D-ID3), to allow for separation of concerns: having a representation, users can later iterate on learning predictors (which need not be decision trees) without additional LLM calls. Most of the time in F2 and D-ID3 is generally spent computing features; luckily, feature computation for all relevant examples is embarrassingly parallel, and we exploit this in our implementation. During training, we always validate proposed features on a subset of the training set (we use 10k examples in our experiments), and discard features that throw exceptions, timeout, or return non-finite values for some example (e.g., NaN or  $\pm\infty$ ).

Our training runs for D-ID3 were the most expensive, with cost ranging from 0.5 to 5 US dollars per run with GPT 5-mini (1000 iterations). Runs with F2 were around 10x cheaper, due to performing 10x less LLM calls. Runs took from 5 to 24h on a CPU-only commodity machine.

## B EXPERIMENTAL DETAILS

### B.1 TRANSFORMER TRAINING

We followed the architectural hyperparameters in Ruoss et al. (2024), and trained the 270M model with cross-entropy loss on their same tokenization scheme, and the same learning rate of  $10^{-5}$ . We train for 300k steps with a batch size of 400 on a machine with 4x H100 NVIDIA GPUs. This gives around 2.4 epochs over 50M data points. We find training to be generally stable, with top-1 move accuracy slowly but monotonically improving across checkpoints (whereas regression metrics, like Pearson correlation, seem less stable, often temporarily decreasing before improving again).

## B.2 MNIST AND FASHION-MNIST TRAINING

We train standard ResNet-50 and EfficientNet-V2 models for 4000 steps and a batch size of 1024 with Adam on a single H100 GPU. We use the default random initialization from PyTorch. For ResNet-50, we use a learning rate of 0.03; for EfficientNet-V2, we use 0.001, which we tuned using the validation set.

## B.3 FUNSEARCH AND PROGRAM OF THOUGHTS BASELINES

For our FunSearch (Novikov et al., 2025) baseline, we based our implementation on Google DeepMind’s (<https://github.com/google-deepmind/funsearch/>) while adapting it to our domain. We used GPT-5.1 as the underlying base model – a much stronger model than what we used with LeaPR (GPT-5 mini). In FunSearch, we evolved 4 “islands” for 50 iterations: each island is a collection of programs. Each program is a complete task-specific predictor: for instance, in Ghostbuster, the program takes a string and outputs a class label (human or AI-generated). Unlike in the original FunSearch paper, for fairness we included the same task description and API reference that we give in our LeaPR prompts (FunSearch used a domain-agnostic prompt instead). The programs are scored during evolution in the same training set we used for LeaPR, and the best program at the end is evaluated in the test set, and its performance reported in our main results tables.

For our Program of Thoughts (Chen et al., 2022) baseline, we also gave GPT-5.1 the same task and API description as for FunSearch and LeaPR. However, here there is no evolution: we simply sample 64 full predictors from the model (divided as 8 predictors).

## C ADDITIONAL RESULTS

We conducted two preliminary experiments on additional datasets to demonstrate: (1) how LeaPR-discovered features compare with domain expert features, and (2) a limitation regarding appropriate library selection. The first experiment was on audio classification, and the second on natural language inference, a more semantic task compared to Ghostbuster. For both experiments, we evaluated D-ID3 with GPT-5-mini, using the same configurations as our main experiments.

### C.1 AUDIO CLASSIFICATION ON ESC-50

The ESC-50 dataset (Piczak) contains 2,000 5-second environmental audio recordings across 50 classes, including animal sounds, natural soundscapes, and urban noises.

This dataset showcases LeaPR working on yet another input modality (raw audio), but mainly it allows us to compare LeaPR’s features with existing expert-generated features, since the original ESC-50 dataset paper included a baseline Random Forest trained on standard audio features. Note that, for audio, the input is very high-dimensional (at 44100 kHz for 5 seconds, each input in ESC-50 is 220,500-dimensional). The standard random forest baseline, trained on hand-designed features for audio, achieved 44.30% top-1 accuracy on ESC-50. When training LeaPR with the same configurations as in our main experiments, and with D-ID3 and gpt-5-mini, we achieve an accuracy of 64.1% with a trained random forest classifier on top. Both underperform the state-of-the-art neural networks, which generally achieve over 90% accuracy on ESC-50 by leveraging large-scale audio pre-training beyond ESC-50 (which only has 2000 audio samples in total). However, this allows us to directly compare LeaPR-learned features with hand-designed ones, showing a case where LeaPR improves on top of human-written features crafted for the same domain.

### C.2 SNLI

The Stanford Natural Language Inference (SNLI) dataset (Bowman et al., 2015) consists of 570k sentence pairs labeled as entailment, contradiction, or neutral. Although this is a text classification task like Ghostbuster, it requires deeper semantic understanding of the text: thus we expect the same LeaPR setup that worked well on Ghostbuster — using only the Python standard library

for manipulating strings — to work less well for tasks requiring semantics. Indeed, we find that D-ID3 with gpt-5-mini achieves an accuracy of 66.3% on SNLI, generally underperforming baselines based on fine-tuning pre-trained Transformers (which achieve 90%+ accuracy). In this setting, LeaPR features cannot reliably perform useful operations like checking if two words are synonyms or antonyms, which are highly useful for this task. We hypothesize that having access to better libraries, such as NLTK or spaCy, could be significantly beneficial here for LeaPR. Overall, this study shows that LeaPR is limited by the libraries it has access to, and it should not be expected to perform well in domains where good libraries are not available to the model.

## D ADDITIONAL RESULTS ON FEATURE SUBSETS

Figure 5 shows the same analyses we discussed in Section 4.4 for the remaining domains. Again, features proposed earliest tend to be most impactful, with performance consistently improving with more features.

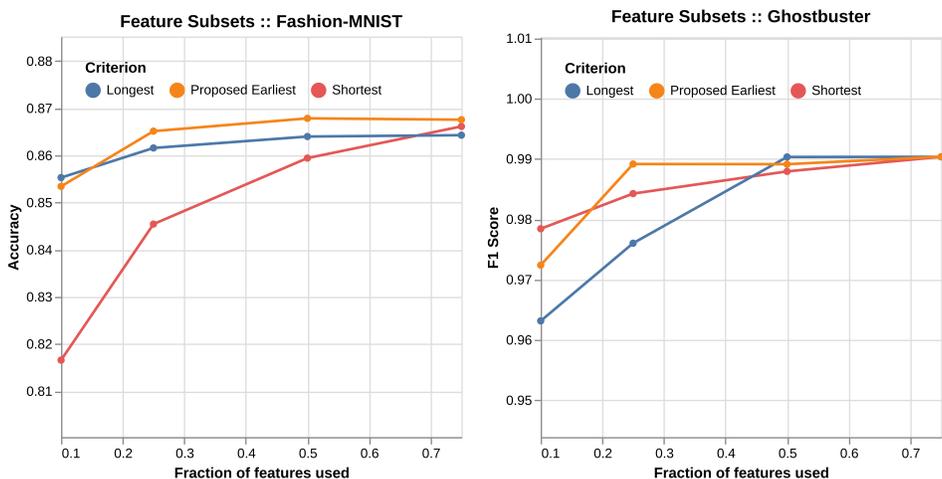


Figure 5: Same experiment from Figure 3 on Fashion-MNIST and Ghostbuster.

## E AUTOMATED FEATURE GENERATION METHODS FOR TABULAR DATA

**LLM-based Feature Generation for Tabular Data** Recently, OCTree (Nam et al., 2024) proposed using LLMs to generate *transformations of existing features* for *tabular data* by leveraging decision tree reasoning, where the LLM synthesizes features specifically designed to improve splits at nodes based on existing tree paths. While their approach shares our motivation of creating interpretable programmatic features through LLM synthesis, it differs in two key aspects: (1) OCTree operates on *structured tabular data with pre-existing features*, while LeaPR is more general and works with raw inputs like images or text, and (2) OCTree uses the *structure* of already-trained trees to guide feature generation, whereas our D-ID3 algorithm generates features during tree construction based on a subset of training examples that reached that node and our FunSearch variant uses evolutionary selection independent of tree structure.

Similarly to OCTree, CAAFE (Hollmann et al., 2023) also uses LLMs to generate *feature transformations* for *tabular datasets* by leveraging dataset metadata and column descriptions, then iteratively selects the most promising features through rounds of model training and feature importance analysis. While our FunSearch variant also follows a generate-then-select pipeline, it differs from both

864 OCTree and CAAFE in generating features from raw inputs rather than tabular data with pre-existing  
865 features.  
866

867  
868 **Rule-based Automated Feature Engineering for Tabular Data** Alternative approaches to LLM  
869 synthesis use deterministic, rule-based methods for automated feature engineering. OpenFE (Zhang  
870 et al., 2023) proposes a two-stage approach for *tabular data*: it generates candidate features through  
871 combinatorial expansion of predefined operators (arithmetic operations, group-by aggregations, etc.)  
872 organized as operator trees, then prunes candidates using gradient boosting feature importance.  
873 While OpenFE shares our FunSearch variant’s generate-and-select approach, it uses determinis-  
874 tic combinatorial search over a fixed operator vocabulary rather than LLM-based synthesis, limiting  
875 it to predefined transformations.  
876

877 Similarly, `autofeat` (Horn et al., 2020) systematically generates polynomial and non-linear trans-  
878 formations of existing *tabular features* (products, ratios, logarithms) up to a specified complexity,  
879 then selects useful features via Lasso regression. While both `autofeat` and OpenFE produce in-  
880 terpretable mathematical expressions and achieve strong performance on *tabular benchmarks*, they  
881 differ fundamentally from our LLM-based approach in relying on *predefined transformation rules*.  
882 This limits them to standard mathematical operations and tabular data with pre-existing features,  
883 while LeapR’s LLM-based generation enables domain-specific reasoning and works directly from  
884 raw inputs across diverse modalities.

## 885 F PROMPTS AND EXAMPLE FEATURES

### 886 F.1 CHESS POSITION EVALUATION

#### 887 F.1.1 F2

```
890 1 You are an expert chess programmer creating feature functions to help a machine learning model
891   predict chess position evaluations.
892 2
893 3 Your task is to write a feature function that helps discriminate between board positions with
894   different evaluations (e.g., probability that white or black wins).
895 4 A feature function is a Python function that takes a chess Board and computes a feature out of
896   the board. It should return a float, but note that a feature could also be effectively
897   boolean-valued (0.0 or 1.0), or integer-valued, even if its type is float.
898 5
899 6 You have access to the following API from the 'chess' library:
900 7
901 8 # Chess API Documentation
902 9 ## Class chess.Board Methods
903 10 - board.turn: True if White to move, False if Black
904 11 - board.fullmove_number: Current move number
905 12 - board.halfmove_clock: Halfmove clock for 50-move rule
906 13 - board.is_check(): True if current player is in check
907 14 - board.is_checkmate(): True if current player is checkmated
908 15 - board.is_stalemate(): True if stalemate
909 16 - board.is_insufficient_material(): Returns True if insufficient material
910 17 - board.piece_at(square): Returns piece at given square (or None)
911 18 - board.piece_map(): Returns dict mapping squares to pieces
912 19 - board.legal_moves: Iterator over legal moves
913 20 - board.attackers(color, square): Returns set of squares attacking given square
914 21 - board.is_attacked_by(color, square): Returns True if square is attacked by color
915 22
916 23 ## Chess Squares and Pieces
917 24 - chess.A1, chess.A2, ..., chess.H8: Square constants
918 25 - chess.PAWN, chess.KNIGHT, chess.BISHOP, chess.ROOK, chess.QUEEN, chess.KING: Piece types
919 26 - chess.WHITE, chess.BLACK: Colors
920 27 - piece.piece_type: Type of piece (PAWN, KNIGHT, etc.)
921 28 - piece.color: Color of piece (WHITE or BLACK)
922 29
923 30 ## Useful Functions
924 31 - chess.square_name(square): Convert square index to name (e.g., "e4")
925 32 - chess.parse_square(name): Convert square name to index
926 33 - chess.square_file(square): Get file (0-7) of square
927 34 - chess.square_rank(square): Get rank (0-7) of square
928 35 - chess.square_distance(sq1, sq2): Manhattan distance between squares
929 36
930 37
931 38 ## Current Feature Database
```

```

918 39 Here are some of our existing features and their importance to the current model (higher
919     importance means this is a more useful feature for the current model):
920 40
921 41 Feature:
922 42 def feature(board: chess.Board) -> float:
923 43     'King centralization in endgames: positive if White king is closer to center than Black
924     king (only active when low material)'
925 44     total_non_king_pieces = sum(1 for _, p in board.piece_map().items() if p.piece_type !=
926     chess.KING)
927 45     # Activation threshold: small material (pawns + pieces <= 6)
928 46     if total_non_king_pieces > 6:
929 47         return 0.0
930 48     center_sq = [chess.parse_square(s) for s in ('d4','e4','d5','e5')]
931 49     center_sq_avg = sum(center_sq) # not used directly; use center coords (3.5,3.5)
932 50     def king_dist(color):
933 51         for sq, p in board.piece_map().items():
934 52             if p.piece_type == chess.KING and p.color == color:
935 53                 # distance to center by min over central squares
936 54                 return min(chess.square_distance(sq, c) for c in center_sq)
937 55         return 8.0
938 56     wd = king_dist(chess.WHITE)
939 57     bd = king_dist(chess.BLACK)
940 58     # Normalize in range roughly -1..1
941 59     return float((bd - wd) / 8.0)
942 60
943 61
944 62 Importance: 0.014
945 63 ---
946 64
947 65
948 66 Feature:
949 67 def feature(board: chess.Board) -> float:
950 68     'Bishop-pair and minor-piece balance: (White advantage) bishops and minor piece
951     composition bonus'
952 69     w_bishops = 0
953 70     b_bishops = 0
954 71     w_minors = 0
955 72     b_minors = 0
956 73     for _, p in board.piece_map().items():
957 74         if p.piece_type == chess.BISHOP:
958 75             if p.color == chess.WHITE:
959 76                 w_bishops += 1
960 77             else:
961 78                 b_bishops += 1
962 79         if p.piece_type in (chess.BISHOP, chess.KNIGHT):
963 80             if p.color == chess.WHITE:
964 81                 w_minors += 1
965 82             else:
966 83                 b_minors += 1
967 84     score = 0.0
968 85     # bishop pair bonus
969 86     if w_bishops >= 2:
970 87         score += 0.6
971 88     if b_bishops >= 2:
972 89         score -= 0.6
973 90     # minor piece imbalance small weight
974 91     score += 0.12 * (w_minors - b_minors)
975 92     return float(score)
976 93
977 94
978 95 Importance: 0.039
979 96 ---
980 97
981 98
982 99 Feature:
983 100 def feature(board: chess.Board) -> float:
984 101     'Center control: difference in control of d4,e4,d5,e5 (occupied=1, attacked=0.5)'
985 102     center_squares = [chess.parse_square(s) for s in ('d4','e4','d5','e5')]
986 103     def control_for(color):
987 104         c = 0.0
988 105         for sq in center_squares:
989 106             occ = board.piece_at(sq)
990 107             if occ is not None and occ.color == color:
991 108                 c += 1.0
992 109             # attacked by color
993 110             if board.is_attacked_by(color, sq):
994 111                 c += 0.5
995 112         return c
996 113     wc = control_for(chess.WHITE)
997 114     bc = control_for(chess.BLACK)
998 115     return float(wc - bc)

```

```

972 116
973 117
974 118 Importance: 0.044
975 119 ---
976 120
977 121
978 122 Feature:
979 123 def feature(board: chess.Board) -> float:
980 124     'Piece activity squares: difference in number of unique squares attacked by non-pawn, non-
981 125     king pieces (White - Black) normalized by 64'
982 126     def active_squares(color):
983 127         count = 0
984 128         for sq in range(64):
985 129             attackers = board.attackers(color, sq)
986 130             found = False
987 131             for a in attackers:
988 132                 p = board.piece_at(a)
989 133                 if p is None:
990 134                     continue
991 135                 if p.piece_type not in (chess.PAWN, chess.KING):
992 136                     found = True
993 137                     break
994 138             if found:
995 139                 count += 1
996 140         return count
997 141     w = active_squares(chess.WHITE)
998 142     b = active_squares(chess.BLACK)
999 143     return float((w - b) / 64.0)
1000 144
1001 145 Importance: 0.329
1002 146 ---
1003 147
1004 148
1005 149 Feature:
1006 150 def feature(board: chess.Board) -> float:
1007 151     'Pawn structure weakness: (Black penalties - White penalties), positive if Black worse (
1008 152     good for White). Penalty = doubled*0.5 + isolated*0.7'
1009 153     def pawn_penalty(color):
1010 154         files = {f:0 for f in range(8)}
1011 155         pawn_sqs = []
1012 156         for sq, p in board.piece_map().items():
1013 157             if p.piece_type == chess.PAWN and p.color == color:
1014 158                 f = chess.square_file(sq)
1015 159                 files[f] += 1
1016 160                 pawn_sqs.append(sq)
1017 161         doubled = sum(max(0, cnt-1) for cnt in files.values())
1018 162         isolated = 0
1019 163         for sq in pawn_sqs:
1020 164             f = chess.square_file(sq)
1021 165             if files.get(f-1,0) == 0 and files.get(f+1,0) == 0:
1022 166                 isolated += 1
1023 167         return doubled * 0.5 + isolated * 0.7
1024 168     bp = pawn_penalty(chess.BLACK)
1025 169     wp = pawn_penalty(chess.WHITE)
1026 170     return float(bp - wp)
1027 171
1028 172 Importance: 0.065
1029 173 ---
1030 174
1031 175
1032 176 Feature:
1033 177 def feature(board: chess.Board) -> float:
1034 178     'King safety pressure: weighted sum of attackers near each king (positive = pressure on
1035 179     Black king > White king)'
1036 180     def king_square(color):
1037 181         for sq, p in board.piece_map().items():
1038 182             if p.piece_type == chess.KING and p.color == color:
1039 183                 return sq
1040 184         return None
1041 185     wk = king_square(chess.WHITE)
1042 186     bk = king_square(chess.BLACK)
1043 187     values = {chess.PAWN:1.0, chess.KNIGHT:3.0, chess.BISHOP:3.25, chess.ROOK:5.0, chess.QUEEN
1044 188     :9.0, chess.KING:0.5}
1045 189     def pressure_on(king_sq, attacker_color):
1046 190         if king_sq is None:
1047 191             return 0.0
1048 192         attackers = board.attackers(attacker_color, king_sq)
1049 193         s = 0.0
1050 194         for a in attackers:

```

```

1026
1027 193         p = board.piece_at(a)
1028 194         if p is None:
1029 195             continue
1030 196         dist = chess.square_distance(a, king_sq)
1031 197         s += values.get(p.piece_type, 0.0) / (1.0 + dist)
1032 198         return s
1033 199         p_on_black = pressure_on(bk, chess.WHITE)
1034 200         p_on_white = pressure_on(wk, chess.BLACK)
1035 201         return float(p_on_black - p_on_white)
1036 202
1037 203
1038 204 Importance: 0.005
1039 205 ---
1040 206
1041 207
1042 208 Feature:
1043 209 def feature(board: chess.Board) -> float:
1044 210     'Undefended high-value threats: (value of Black pieces under more attackers than defenders
1045 211     ) - (value of White pieces similarly threatened)'
1046 212     values = {chess.PAWN:1.0, chess.KNIGHT:3.0, chess.BISHOP:3.25, chess.ROOK:5.0, chess.QUEEN
1047 213     :9.0, chess.KING:0.0}
1048 214     threat_black = 0.0
1049 215     threat_white = 0.0
1050 216     for sq, p in board.piece_map().items():
1051 217         if p.piece_type == chess.PAWN:
1052 218             continue
1053 219         attackers = board.attackers(not p.color, sq)
1054 220         defenders = board.attackers(p.color, sq)
1055 221         atk = sum(1 for _ in attackers)
1056 222         defn = sum(1 for _ in defenders)
1057 223         if atk > defn and atk > 0:
1058 224             score = values.get(p.piece_type, 0.0) * (atk - defn)
1059 225             if p.color == chess.BLACK:
1060 226                 threat_black += score
1061 227             else:
1062 228                 threat_white += score
1063 229         return float(threat_black - threat_white)
1064 230
1065 231 Importance: 0.045
1066 232 ---
1067 233
1068 234 Feature:
1069 235 def feature(board: chess.Board) -> float:
1070 236     'Material balance (White minus Black) using common piece values: P=1,N=3,B=3.25,R=5,Q=9'
1071 237     values = {chess.PAWN:1.0, chess.KNIGHT:3.0, chess.BISHOP:3.25, chess.ROOK:5.0, chess.QUEEN
1072 238     :9.0, chess.KING:0.0}
1073 239     total = 0.0
1074 240     for sq, piece in board.piece_map().items():
1075 241         val = values.get(piece.piece_type, 0.0)
1076 242         total += val if piece.color == chess.WHITE else -val
1077 243     return float(total)
1078 244
1079 245 Importance: 0.199
1080 246 ---
1081 247
1082 248 Feature:
1083 249 def feature(board: chess.Board) -> float:
1084 250     'Normalized mobility difference: (White legal moves - Black legal moves) / 100'
1085 251     try:
1086 252         white_moves = 0
1087 253         black_moves = 0
1088 254         # count current side moves
1089 255         white_board = board.copy()
1090 256         white_board.turn = chess.WHITE
1091 257         white_moves = sum(1 for _ in white_board.legal_moves)
1092 258         black_board = board.copy()
1093 259         black_board.turn = chess.BLACK
1094 260         black_moves = sum(1 for _ in black_board.legal_moves)
1095 261         return float((white_moves - black_moves) / 100.0)
1096 262     except Exception:
1097 263         return 0.0
1098 264
1099 265
1100 266 Importance: 0.165
1101 267 ---
1102 268
1103 269
1104 270 Feature:

```

```

1080
1081 271 def feature(board: chess.Board) -> float:
1082 272     'Passed pawns score: sum of passed-pawn strengths (White minus Black), advanced pawns
1083 273     weighted more'
1084 274     def is_passed(sq, color):
1085 275         f = chess.square_file(sq)
1086 276         r = chess.square_rank(sq)
1087 277         if color == chess.WHITE:
1088 278             ahead_ranks = range(r+1, 8)
1089 279             opp_color = chess.BLACK
1090 280             for ar in ahead_ranks:
1091 281                 for df in (-1,0,1):
1092 282                     ff = f + df
1093 283                     if 0 <= ff < 8:
1094 284                         sq2 = chess.square(ff, ar)
1095 285                         p = board.piece_at(sq2)
1096 286                         if p is not None and p.color == opp_color and p.piece_type == chess.
1097 287                             PAWN:
1098 288                                 return False
1099 289                             return True
1100 290                         else:
1101 291                             ahead_ranks = range(r-1, -1, -1)
1102 292                             opp_color = chess.WHITE
1103 293                             for ar in ahead_ranks:
1104 294                                 for df in (-1,0,1):
1105 295                                     ff = f + df
1106 296                                     if 0 <= ff < 8:
1107 297                                         sq2 = chess.square(ff, ar)
1108 298                                         p = board.piece_at(sq2)
1109 299                                         if p is not None and p.color == opp_color and p.piece_type == chess.
1110 300                                             PAWN:
1111 301                                                 return False
1112 302                                                 return True
1113 303         score_w = 0.0
1114 304         score_b = 0.0
1115 305         for sq, p in board.piece_map().items():
1116 306             if p.piece_type != chess.PAWN:
1117 307                 continue
1118 308             rank = chess.square_rank(sq)
1119 309             if p.color == chess.WHITE:
1120 310                 if is_passed(sq, chess.WHITE):
1121 311                     advancement = (rank - 1) / 6.0 if rank >= 1 else 0.0
1122 312                     score_w += 1.0 + advancement
1123 313             else:
1124 314                 if is_passed(sq, chess.BLACK):
1125 315                     advancement = (6 - rank) / 6.0 if rank <= 6 else 0.0
1126 316                     score_b += 1.0 + advancement
1127 317         return float(score_w - score_b)
1128 318
1129 319 Importance: 0.094
1130 320 ---
1131 321 # Task
1132 322 Generate 10 new chess board feature functions in Python that:
1133 323
1134 324 1. Help us discriminate between strong and weak board positions, hopefully with positions
1135 325     before and after the optimal split point having the lowest possible variance between
1136 326     their evaluations.
1137 327
1138 328 2. Return a float value given a board position.
1139 329
1140 330 3. Handle edge cases gracefully - won't crash on unusual positions
1141 331
1142 332 Your task is to generate diverse, creative features that are relevant to explain the
1143 333     evaluations for the board positions shown above. Focus on features that would help
1144 334     distinguish between positions of different strengths. These features will be used in this
1145 335     decision tree that will predict the evaluation of a given board position in estimated %
1146 336     win probability for white (e.g., 20 means Black winning with around 80% probability).
1147 337     Think about new features that would help such a predictor in the particular cases above,
1148 338     trying to add information that the already existing features shown above are missing.
1149 339
1150 340 # Code Requirements
1151 341
1152 342 - Use single quotes for docstrings: "description here"
1153 343 - No markdown code blocks
1154 344 - No explanatory text after the function
1155 345 - Each function should be complete and standalone, and return a float
1156 346
1157 347 # Output Format
1158 348 Generate exactly 10 features in this format:
1159 349
1160 350 def feature(board: chess.Board) -> float:

```

```

1134
1135 341 "Simple, clear description of what this feature measures"
1136 342 # ... Calculate and return the feature value
1137 343 return result
1138 344
1139 345 def feature(board: chess.Board) -> float:
1140 346     "Another feature description"
1141 347     # ... Calculate and return the feature value
1142 348     return result
1143 349
1144 350 The body of the function can be anything, but the first line (function declaration) should be
1145 351 identical to those examples above, and the second line should be a one-line docstring.
1146 352 Don't output explanatory text - just the function definitions as shown above.
1147
1148 351
1149 352 Optimize for producing discriminant features that are novel compared to the existing features
1150 353 and that are likely to achieve a high importance for scoring positions, once we retrain
1151 354 the model using your new features combined with the existing ones.
1152
1153
1154
1155
1156
1157
1158
1159
1160
1161
1162
1163
1164
1165
1166
1167
1168
1169
1170
1171
1172
1173
1174
1175
1176
1177
1178
1179
1180
1181
1182
1183
1184
1185
1186
1187
1188
1189

```

### F.1.2 D-ID3 - PROMPT

```

1 You are an expert chess programmer creating feature functions to help a machine learning model
predict chess position evaluations.
2
3 Your task is to write a feature function that helps discriminate between the board positions
given below.
4 A feature function is a Python function that takes a chess Board and computes a feature out of
the board. It should return a float, but note that a feature could also be effectively
boolean-valued (0.0 or 1.0), or integer-valued, even if its type is float.
5
6 You have access to the following API from the 'chess' library:
7
8 # Chess API Documentation
9 ## Class chess.Board Methods
10 - board.turn: True if White to move, False if Black
11 - board.fullmove_number: Current move number
12 - board.halfmove_clock: Halfmove clock for 50-move rule
13 - board.is_check(): True if current player is in check
14 - board.is_checkmate(): True if current player is checkmated
15 - board.is_stalemate(): True if stalemate
16 - board.is_insufficient_material(): Returns True if insufficient material
17 - board.piece_at(square): Returns piece at given square (or None)
18 - board.piece_map(): Returns dict mapping squares to pieces
19 - board.legal_moves: Iterator over legal moves
20 - board.attackers(color, square): Returns set of squares attacking given square
21 - board.is_attacked_by(color, square): Returns True if square is attacked by color
22
23 ## Chess Squares and Pieces
24 - chess.A1, chess.A2, ..., chess.H8: Square constants
25 - chess.PAWN, chess.KNIGHT, chess.BISHOP, chess.ROOK, chess.QUEEN, chess.KING: Piece types
26 - chess.WHITE, chess.BLACK: Colors
27 - piece.piece_type: Type of piece (PAWN, KNIGHT, etc.)
28 - piece.color: Color of piece (WHITE or BLACK)
29
30 ## Useful Functions
31 - chess.square_name(square): Convert square index to name (e.g., "e4")
32 - chess.parse_square(name): Convert square name to index
33 - chess.square_file(square): Get file (0-7) of square
34 - chess.square_rank(square): Get rank (0-7) of square
35 - chess.square_distance(sq1, sq2): Manhattan distance between squares
36
37
38 # Task
39 Generate 10 new chess board feature functions in Python that:
40
41 1. Help us discriminate between strong and weak board positions, hopefully with positions
before and after the optimal split point having the lowest possible variance between
their evaluations.
42 2. Return a float value given a board position.
43 3. Handle edge cases gracefully - won't crash on unusual positions
44
45 Your task is to generate diverse, creative features that are relevant to explain the
evaluations for the board positions shown above. Focus on features that would help
distinguish between positions of different strengths. These features will be used in this
decision tree that will predict the evaluation of a given board position in estimated %
win probability for white (e.g., 20 means Black winning with around 80% probability).
Think about new features that would help such a predictor in the particular cases above,
trying to add information that the already existing features shown above are missing.
46
47 # Code Requirements
48
49 - Use single quotes for docstrings: "description here"

```

```

1188 50 - No markdown code blocks
1189 51 - No explanatory text after the function
1190 52 - Each function should be complete and standalone, and return a float
1191 53
1191 54 # Output Format
1192 55 Generate exactly 10 features in this format:
1193 56
1193 57 def feature(board: chess.Board) -> float:
1194 58     "Simple, clear description of what this feature measures"
1195 59     # ... Calculate and return the feature value
1196 60     return result
1197 61
1197 62 def feature(board: chess.Board) -> float:
1198 63     "Another feature description"
1199 64     # ... Calculate and return the feature value
1200 65     return result
1200 66
1200 67 The body of the function can be anything, but the first line (function declaration) should be
1201 68 identical to those examples above, and the second line should be a one-line docstring.
1202 69 Don't output explanatory text - just the function definitions as shown above.
1202 70
1202 71 # Current decision tree node
1203 72 You are currently focusing on features that explain the position evaluation of board positions
1204 73 in the following subtree of a decision tree:
1205 74
1205 75 [root]
1206 76 -> value < 3.000 for "Measure the material balance between both players" -> value < 1.182
1207 77 for "Counts the number of pieces for each player and computes the ratio of the piece
1208 78 counts." -> value > 0.650 for "Counts the number of squares attacked by pieces of each
1209 79 color to assess control of the board."
1210 80
1210 81 # Board positions
1211 82 Here are examples of board positions in this subtree, along with their position evaluations (
1212 83 computed by Stockfish):
1213 84
1213 85 Board:
1214 86 . . r r . . k .
1215 87 . . . . . p p p
1216 88 p b . . p . . .
1217 89 . p . . P . . .
1218 90 . P . . N P n q
1219 91 P . . . . . P .
1220 92 . B . . Q . . P
1221 93 R . . . . R . K
1222 94 Evaluation: 38.08281678856247
1223 95 ---
1224 96
1224 97 Board:
1225 98 r . . . . . k .
1226 99 . b . r . p p .
1227 100 p . . P . . . .
1228 101 . p . . . . . Q
1229 102 . P . b . . . .
1230 103 P . . . . . .
1231 104 . . . . . P P
1232 105 . . . . . R . K
1233 106 Evaluation: 15.383900101915987
1234 107 ---
1235 108
1235 109 Board:
1236 110 r n b . k . . r
1237 111 p p p . n p p p
1238 112 . . . b . . . .
1239 113 . . . B . . . .
1240 114 . . . P P p . q
1241 115 . . . . . . .
1242 116 P P P . . . P P
1243 117 R N B Q . K N R
1244 118 Evaluation: 47.700323643230064
1245 119 ---
1246 120
1246 121 Board:
1247 122 . . . r . . . k
1248 123 p p . . . p p .
1249 124 . . . . . p . . p
1250 125 . . b r P . . P
1251 126 . q . N R . . .
1252 127 . . . . . B . . .
1253 128 . . . . . K P .
1254 129 . . R Q . . . .
1255 130 Evaluation: 50.0

```

```

1242
1243 124 ---
1244 125
1245 126 Board:
1246 127 r n b q k b n r
1247 128 p p p p . p p p
1248 129 . . . . .
1249 130 . . . . p . . .
1250 131 . . . . P . . .
1251 132 . . . . . N . .
1252 133 P P P P . P P P
1253 134 R N B Q K B . R
1254 135 Evaluation: 51.19644756166626
1255 136 ---
1256 137
1257 138
1258 139 Optimize for producing discriminant features that are novel compared to the existing features
1259 used to arrive at this subtree. Focus on explaining the differences between the board
1260 positions shown above.

```

### 1255 F.1.3 D-ID3 - EXAMPLE FEATURES

```

1257 1 def feature(board: chess.Board) -> float:
1258 2     "Measure the material balance between both players"
1259 3     material_values = {chess.PAWN: 1, chess.KNIGHT: 3, chess.BISHOP: 3, chess.ROOK: 5, chess.QUEEN: 9, chess.KING: 0}
1260 4     white_material = sum(material_values[piece.piece_type] for piece in board.piece_map().values() if piece.color == chess.WHITE)
1261 5     black_material = sum(material_values[piece.piece_type] for piece in board.piece_map().values() if piece.color == chess.BLACK)
1262 6     return float(white_material - black_material)
1263 7
1264 8
1265 9 def feature(board: chess.Board) -> float:
1266 10     "Counts the number of pieces for each player and computes the ratio of the piece counts."
1267 11     piece_count_white = sum(1 for piece in board.piece_map().values() if piece.color == chess.WHITE)
1268 12     piece_count_black = sum(1 for piece in board.piece_map().values() if piece.color == chess.BLACK)
1269 13     if piece_count_black == 0:
1270 14         return float('inf') # Black has no pieces left
1271 15     return piece_count_white / piece_count_black
1272 16
1273 17 def feature(board: chess.Board) -> float:
1274 18     "Counts the number of squares attacked by pieces of each color to assess control of the board."
1275 19     white_attacks = sum(board.is_attacked_by(chess.WHITE, square) for square in chess.SQUARES)
1276 20     black_attacks = sum(board.is_attacked_by(chess.BLACK, square) for square in chess.SQUARES)
1277 21     control_ratio = white_attacks / (black_attacks + 1) # Avoid division by zero
1278 22     return float(control_ratio)
1279 23
1280 24 def feature(board: chess.Board) -> float:
1281 25     "Calculates the center control by counting pieces in the central squares."
1282 26     central_squares = [chess.D4, chess.D5, chess.E4, chess.E5]
1283 27     control = sum(1 for square in central_squares if board.piece_at(square) is not None)
1284 28     return float(control)
1285 29
1286 30 def feature(board: chess.Board) -> float:
1287 31     "Calculates the total piece value for each player based on standard chess piece values."
1288 32     piece_values = {
1289 33         chess.PAWN: 1,
1290 34         chess.KNIGHT: 3,
1291 35         chess.BISHOP: 3,
1292 36         chess.ROOK: 5,
1293 37         chess.QUEEN: 9,
1294 38         chess.KING: 0 # King is invaluable
1295 39     }
1296 40     white_value = sum(piece_values[piece.piece_type] for piece in board.piece_map().values() if piece.color == chess.WHITE)
1297 41     black_value = sum(piece_values[piece.piece_type] for piece in board.piece_map().values() if piece.color == chess.BLACK)
1298 42     return float(white_value - black_value)

```

## 1291 F.2 IMAGE CLASSIFICATION

### 1292 F.2.1 F2

```

1293 1 You are an expert computer vision programmer creating evaluation features for a machine
1294 learning model that predicts values from images.
1295 2

```

```

1296 3 This is a classification task with the following classes: 0: landbird, 1: waterbird.
1297 4
1298 5 Your task is to write a feature function that helps discriminate between the image classes
1299 6 above.
1300 7 A feature function is a Python function that takes an image array and computes a feature out
1301 8 of the image. It should return a float, but note that a feature could also be effectively
1302 9 boolean-valued (0.0 or 1.0), or integer-valued, even if its type is float.
1303 10
1304 11 You have access to the following API from image processing libraries:
1305 12
1306 13 # Image Processing API Documentation
1307 14
1308 15 The features receive an image as a numpy array, so you can use any numpy functions on it. For
1309 16 RGB images, shape is (height, width, 3). For grayscale, shape is (height, width).
1310 17
1311 18 ## Image Processing Methods
1312 19 - image.shape: Returns (height, width, channels) for RGB or (height, width) for grayscale
1313 20 - image.mean(): Average pixel intensity across all channels
1314 21 - image.std(): Standard deviation of pixel intensities
1315 22 - image.max(), image.min(): Maximum and minimum pixel values
1316 23 - np.sum(image): Sum of all pixel values
1317 24 - np.count_nonzero(image): Count of non-zero pixels
1318 25
1319 26 ## Handle Both Grayscale and RGB
1320 27 - Check format: len(image.shape) == 2 for grayscale, len(image.shape) == 3 for RGB
1321 28 - Unpack safely: h, w = image.shape[:2] # Works for both formats
1322 29 - For RGB only: image[:, :, 0] (red), image[:, :, 1] (green), image[:, :, 2] (blue)
1323 30
1324 31 ## Useful Functions
1325 32 - np.mean(image): Average intensity
1326 33 - np.std(image): Standard deviation
1327 34 - np.gradient(image): Image gradients - for RGB use on single channel: np.gradient(image
1328 35[:, :, 0])
1329 36 - np.where(condition, x, y): Conditional selection
1330 37 - np.argmax(image), np.argmin(image): Location of max/min values
1331 38 - np.percentile(image, q): Percentile values
1332 39 - np.histogram(image.flatten(), bins): Intensity histogram
1333 40
1334 41 ## Spatial Analysis
1335 42 - image[start_row:end_row, start_col:end_col]: Region selection
1336 43 - Center region: image[h//4:3*h//4, w//4:3*w//4]
1337 44 - Edge detection: np.gradient(np.mean(image, axis=2)) for RGB
1338 45 - Color channel differences: image[:, :, 0] - image[:, :, 1]
1339 46
1340 47 ## Example Feature Function
1341 48 def feature(image: np.ndarray) -> float:
1342 49     "Average pixel intensity in the center region"
1343 50     if len(image.shape) == 3:
1344 51         h, w, c = image.shape
1345 52         gray = np.mean(image, axis=2)
1346 53     else:
1347 54         h, w = image.shape
1348 55         gray = image
1349 56     center_h, center_w = h // 4, w // 4
1350 57     center_region = gray[center_h:3*center_h, center_w:3*center_w]
1351 58     return float(np.mean(center_region))
1352 59
1353 60 ## Current Feature Database
1354 61 Here are some existing features and their importances to the current classifier (importance =
1355 62 benefit from that feature, higher is better):
1356 63
1357 64 Feature:
1358 65 def feature(image: np.ndarray) -> float:
1359 66     'Relative difference in average blue intensity between bottom half and top half'
1360 67     import numpy as np
1361 68     h, w = image.shape[:2]
1362 69     if len(image.shape) != 3 or image.shape[2] < 3:
1363 70         return float(0.0)
1364 71     b = image[:, :, 2].astype(float)
1365 72     top_mean = np.mean(b[h // 2:, :]) if h // 2 > 0 else np.mean(b)
1366 73     bot_mean = np.mean(b[:h // 2, :]) if h - h // 2 > 0 else np.mean(b)
1367 74     denom = (np.mean(b) + 1e-8)
1368 75     return float((bot_mean - top_mean) / denom)
1369 76
1370 77 Importance: 0.081
1371 78 ---
1372 79
1373 80

```

```

1350 78 Feature:
1351 79 def feature(image: np.ndarray) -> float:
1352 80     'Aspect ratio (width/height) of bounding box of pixels significantly different from median
1353 81     intensity'
1354 82     import numpy as np
1355 83     h, w = image.shape[:2]
1356 84     if len(image.shape) == 3:
1357 85         gray = np.mean(image, axis=2).astype(float)
1358 86     else:
1359 87         gray = image.astype(float)
1360 88     med = np.median(gray)
1361 89     dynamic_range = np.max(gray) - np.min(gray)
1362 90     thresh = 0.15 * (dynamic_range + 1e-8)
1363 91     mask = np.abs(gray - med) > thresh
1364 92     if not np.any(mask):
1365 93         return float(0.5)
1366 94     rows = np.where(np.any(mask, axis=1))[0]
1367 95     cols = np.where(np.any(mask, axis=0))[0]
1368 96     if rows.size == 0 or cols.size == 0:
1369 97         return float(0.5)
1370 98     r0, r1 = rows[0], rows[-1]
1371 99     c0, c1 = cols[0], cols[-1]
1372 100     bbox_h = (r1 - r0 + 1)
1373 101     bbox_w = (c1 - c0 + 1)
1374 102     return float(bbox_w / (bbox_h + 1e-8))
1375 103 Importance: 0.082
1376 104 ---
1377 105
1378 106 Feature:
1379 107 def feature(image: np.ndarray) -> float:
1380 108     'Left-right symmetry score (1.0 = perfectly symmetric, lower = less symmetric)'
1381 109     import numpy as np
1382 110     if len(image.shape) == 3:
1383 111         gray = np.mean(image, axis=2).astype(float)
1384 112     else:
1385 113         gray = image.astype(float)
1386 114     flipped = np.fliplr(gray)
1387 115     diff = np.mean(np.abs(gray - flipped))
1388 116     # normalize by image contrast to avoid small-image issues
1389 117     denom = np.mean(np.abs(gray - np.mean(gray))) + 1e-8
1390 118     norm_diff = diff / denom
1391 119     score = 1.0 - np.tanh(norm_diff) # bounded between 0 and 1
1392 120     return float(np.clip(score, 0.0, 1.0))
1393 121
1394 122 Importance: 0.084
1395 123 ---
1396 124
1397 125 Feature:
1398 126 def feature(image: np.ndarray) -> float:
1399 127     'Ratio of average horizontal gradient magnitude to vertical gradient magnitude (texture
1400 128     orientation)'
1401 129     import numpy as np
1402 130     # compute gray
1403 131     if len(image.shape) == 3:
1404 132         gray = np.mean(image, axis=2).astype(float)
1405 133     else:
1406 134         gray = image.astype(float)
1407 135     gy, gx = np.gradient(gray)
1408 136     mean_dx = np.mean(np.abs(gx))
1409 137     mean_dy = np.mean(np.abs(gy))
1410 138     return float(mean_dx / (mean_dy + 1e-8))
1411 139
1412 140 Importance: 0.159
1413 141 ---
1414 142
1415 143 Feature:
1416 144 def feature(image: np.ndarray) -> float:
1417 145     'Proportion of pixels where green channel is significantly high (vegetation cue)'
1418 146     import numpy as np
1419 147     h, w = image.shape[:2]
1420 148     if len(image.shape) != 3 or image.shape[2] < 3:
1421 149         return float(0.0)
1422 150     img = image.astype(float)
1423 151     r, g, b = img[:, :, 0], img[:, :, 1], img[:, :, 2]
1424 152     # require green greater than both red and blue and at least above 60th percentile

```

```

1404 157 thresh = np.percentile(g.flatten(), 60)
1405 158 mask = (g > r) & (g > b) & (g > thresh)
1406 159 return float(np.count_nonzero(mask) / (h * w + 1e-12))
1407 160
1408 161
1408 162 Importance: 0.094
1409 163 ---
1409 164
1410 165 Feature:
1411 166 def feature(image: np.ndarray) -> float:
1412 167     'Fraction of pixels brighter than the 90th percentile (bright spot proportion)'
1413 168     import numpy as np
1414 169     if len(image.shape) == 3:
1415 170         gray = np.mean(image, axis=2).astype(float)
1416 171     else:
1417 172         gray = image.astype(float)
1418 173     thresh = np.percentile(gray.flatten(), 90)
1419 174     frac = np.count_nonzero(gray > thresh) / (gray.size + 1e-12)
1420 175     return float(frac)
1421 176
1421 177
1421 178
1421 179 Importance: 0.080
1422 180 ---
1422 181
1422 182 Feature:
1423 183 def feature(image: np.ndarray) -> float:
1424 184     'Normalized difference between center region brightness and border brightness'
1425 185     import numpy as np
1426 186     h, w = image.shape[:2]
1427 187     if len(image.shape) == 3:
1428 188         gray = np.mean(image, axis=2).astype(float)
1429 189     else:
1430 190         gray = image.astype(float)
1431 191         ch0, ch1 = h // 4, w // 4
1432 192         center = gray[ch0:3 * ch0 or h, ch1:3 * ch1 or w]
1433 193         # border defined as whole minus center
1434 194         mask_border = np.ones_like(gray, dtype=bool)
1435 195         mask_border[ch0:3 * ch0 or h, ch1:3 * ch1 or w] = False
1436 196         center_mean = np.mean(center) if center.size > 0 else 0.0
1437 197         border_mean = np.mean(gray[mask_border]) if np.any(mask_border) else 0.0
1438 198         denom = np.mean(np.abs(gray)) + 1e-8
1439 199         return float((center_mean - border_mean) / denom)
1440 200
1440 201
1440 202
1440 203 Importance: 0.088
1441 204 ---
1441 205
1441 206 Feature:
1442 207 def feature(image: np.ndarray) -> float:
1443 208     'Proportion of pixels where blue channel is dominant (blue > red and blue > green)'
1444 209     import numpy as np
1445 210     h, w = image.shape[:2]
1446 211     if len(image.shape) != 3 or image.shape[2] < 3:
1447 212         return float(0.0)
1448 213     img = image.astype(float)
1449 214     r, g, b = img[:, :, 0], img[:, :, 1], img[:, :, 2]
1450 215     mask = (b > r) & (b > g)
1451 216     return float(np.count_nonzero(mask) / (h * w + 1e-12))
1452 217
1452 218
1452 219
1452 220 Importance: 0.117
1453 221 ---
1453 222
1453 223 Feature:
1454 224 def feature(image: np.ndarray) -> float:
1455 225     'Edge density: fraction of pixels with gradient magnitude above (mean+std)'
1456 226     import numpy as np
1457 227     if len(image.shape) == 3:
1458 228         gray = np.mean(image, axis=2).astype(float)
1459 229     else:
1460 230         gray = image.astype(float)
1461 231     gy, gx = np.gradient(gray)
1462 232     mag = np.hypot(gx, gy)
1463 233     thresh = np.mean(mag) + np.std(mag)
1464 234     count = np.count_nonzero(mag > thresh)
1465 235     total = gray.size
1466 236     return float(count / (total + 1e-12))
1467 237

```

```

1458
1459 238
1460 239 Importance: 0.114
1461 241 ---
1462 242
1463 243 Feature:
1464 244 def feature(image: np.ndarray) -> float:
1465 245     'Average per-pixel color "saturation" approximated by (max-min)/max across channels'
1466 246     import numpy as np
1467 247     if len(image.shape) != 3 or image.shape[2] < 3:
1468 248         return float(0.0)
1469 249     img = image.astype(float)
1470 250     mx = np.max(img, axis=2)
1471 251     mn = np.min(img, axis=2)
1472 252     # avoid division by zero
1473 253     sat = (mx - mn) / (mx + 1e-8)
1474 254     return float(np.mean(sat))
1475 255
1476 256 Importance: 0.100
1477 257 ---
1478 258
1479 259 ## Task
1480 260 Generate 10 NEW image features that:
1481 261
1482 262 1. Are different from existing features
1483 263 2. Capture useful visual patterns
1484 264 3. Return float values
1485 265 4. Handle edge cases gracefully - Won't crash on unusual images
1486 266 5. Use simple, short docstrings - Use single quotes, not triple quotes
1487 267 6. Are efficient to compute
1488 268
1489 269 Your task is to generate diverse, creative features that capture different aspects of image
1490 270 content for prediction. Focus on features that would help distinguish between different
1491 271 samples. These features will be used as input to a learned model that predicts target
1492 272 values from images.
1493 273
1494 274 ## IMPORTANT CODE REQUIREMENTS
1495 275 - Use SINGLE quotes for docstrings: "description here"
1496 276 - NO triple quotes (""" anywhere in the code
1497 277 - NO markdown code blocks
1498 278 - NO explanatory text after the function
1499 279 - Each function should be complete and standalone
1500 280
1501 281 ## Output Format
1502 282 Generate exactly 10 features in this format:
1503 283
1504 284 def feature(image: np.ndarray) -> float:
1505 285     "Clear description of what this feature measures"
1506 286     # ... Calculate and return the feature value
1507 287     return float(result)
1508 288
1509 289 def feature(image: np.ndarray) -> float:
1510 290     "Another feature description"
1511 291     # ... Calculate and return the feature value
1512 292     return float(result)
1513 293
1514 294 The body of the functions can be anything, but the first line (function declaration) should be
1515 295 identical to those examples above (always 'def feature(...)'), and the second line
1516 296 should be a one-line docstring. Don't output explanatory text - just the function
1517 297 definitions as shown above.

```

## 1501

### 1502 F.2.2 D-ID3 - PROMPT

```

1503 1 You are an expert image processing programmer creating feature functions to help a machine
1504 2 learning model perform image classification.
1505 3 This is a classification task with the following classes: 0: digit zero, 1: digit one, 2:
1506 4 digit two, 3: digit three, 4: digit four, 5: digit five, 6: digit six, 7: digit seven, 8:
1507 5 digit eight, 9: digit nine.
1508 6 Your task is to write a feature function that helps discriminate between the image samples
1509 7 given below.
1510 8 A feature function is a Python function that takes an image array and computes a feature out
1511 9 of the image. It should return a float, but note that a feature could also be effectively
1512 10 boolean-valued (0.0 or 1.0), or integer-valued, even if its type is float.
1513 11
1514 12 You have access to the following API from image processing libraries:

```

```

1512 9
1513 10
1514 11 # Image Processing API Documentation
1515 12
1516 13 The features receive an image as a numpy array, so you can use any numpy functions on it. For
1517 14     RGB images, shape is (height, width, 3). For grayscale, shape is (height, width).
1518 15 ## Image Processing Methods
1519 16 - image.shape: Returns (height, width, channels) for RGB or (height, width) for grayscale
1520 17 - image.mean(): Average pixel intensity across all channels
1521 18 - image.std(): Standard deviation of pixel intensities
1522 19 - image.max(), image.min(): Maximum and minimum pixel values
1523 20 - np.sum(image): Sum of all pixel values
1524 21 - np.count_nonzero(image): Count of non-zero pixels
1525 22
1526 23 ## Handle Both Grayscale and RGB
1527 24 - Check format: len(image.shape) == 2 for grayscale, len(image.shape) == 3 for RGB
1528 25 - Unpack safely: h, w = image.shape[:2] # Works for both formats
1529 26 - For RGB only: image[:, :, 0] (red), image[:, :, 1] (green), image[:, :, 2] (blue)
1530 27
1531 28 ## Useful Functions
1532 29 - np.mean(image): Average intensity
1533 30 - np.std(image): Standard deviation
1534 31 - np.gradient(image): Image gradients - for RGB use on single channel: np.gradient(image
1535 32     [ :, :, 0 ])
1536 33 - np.where(condition, x, y): Conditional selection
1537 34 - np.argmax(image), np.argmin(image): Location of max/min values
1538 35 - np.percentile(image, q): Percentile values
1539 36 - np.histogram(image.flatten(), bins): Intensity histogram
1540 37
1541 38 ## Spatial Analysis
1542 39 - image[start_row:end_row, start_col:end_col]: Region selection
1543 40 - Center region: image[h//4:3*h//4, w//4:3*w//4]
1544 41 - Edge detection: np.gradient(np.mean(image, axis=2)) for RGB
1545 42 - Color channel differences: image[:, :, 0] - image[:, :, 1]
1546 43
1547 44 ## Example Feature Function
1548 45 def feature(image: np.ndarray) -> float:
1549 46     "Average pixel intensity in the center region"
1550 47     if len(image.shape) == 3:
1551 48         h, w, c = image.shape
1552 49         gray = np.mean(image, axis=2)
1553 50     else:
1554 51         h, w = image.shape
1555 52         gray = image
1556 53         center_h, center_w = h // 4, w // 4
1557 54         center_region = gray[center_h:3*center_h, center_w:3*center_w]
1558 55         return float(np.mean(center_region))
1559 56
1560 57 # Task
1561 58 Generate 10 new image feature functions in Python that:
1562 59
1563 60 1. Help us discriminate between different image classes, hopefully with samples before and
1564 61     after the optimal split point having the lowest possible variance between their
1565 62     classifications.
1566 63 2. Return a float value given an image.
1567 64 3. Handle edge cases gracefully - won't crash on unusual images
1568 65
1569 66 Your task is to generate diverse, creative features that are relevant to explain the
1570 67     classifications for the image samples shown above. Focus on features that would help
1571 68     distinguish between samples of different classes. These features will be used in this
1572 69     decision tree that will predict the classification of a given image sample. Think about
1573 70     new features that would help such a predictor in the particular cases above, trying to
1574 71     add information that the already existing features shown above are missing.
1575 72
1576 73 # Code Requirements
1577 74 - Use single quotes for docstrings: "description here"
1578 75 - No markdown code blocks
1579 76 - No explanatory text after the function
1580 77 - Each function should be complete and standalone, and return a float
1581 78
1582 79 # Output Format
1583 80 Generate exactly 10 features in this format:
1584 81
1585 82 def feature(image: np.ndarray) -> float:
1586 83     "Simple, clear description of what this feature measures"
1587 84     # ... Calculate and return the feature value
1588 85     return result
1589 86
1590 87 def feature(image: np.ndarray) -> float:

```

```

1566 81 "Another feature description"
1567 82 # ... Calculate and return the feature value
1568 83 return result
1569 84
1570 85 The body of the function can be anything, but the first line (function declaration) should be
1571 86 identical to those examples above, and the second line should be a one-line docstring.
1572 87 Don't output explanatory text - just the function definitions as shown above.
1573 88 # Current decision tree node
1574 89 You are currently focusing on features that explain the image classifications in the following
1575 90 subtree of a decision tree:
1576 91 [root]
1577 92 -> value > 0.344 for "Ratio of edge pixels to the total number of pixels" -> value > 72.865
1578 93 for "Computes the brightness of the central region as a percentage of the whole image"
1579 94 -> value > 0.805 for "Determines the ratio of vertical to horizontal gradients in the
1580 95 image, indicating edge direction"
1581 96 # Image samples
1582 97 Here are examples of image samples in this subtree, along with their target classifications:
1583 98 Sample (Target: 9)
1584 99 Sample (Target: 6)
1585 100 Sample (Target: 7)
1586 101 Sample (Target: 4)
1587 102 Sample (Target: 4)
1588 103 Optimize for producing discriminant features that are novel compared to the existing features
1589 104 used to arrive at this subtree. Focus on explaining the differences between the image
1590 105 samples shown above.

```

### 1587 F.2.3 D-ID3 - EXAMPLE FEATURES

```

1589 1 def feature(image: np.ndarray) -> float:
1590 2     "Ratio of edge pixels to the total number of pixels"
1591 3     if len(image.shape) == 3:
1592 4         gray = np.mean(image, axis=2)
1593 5     else:
1594 6         gray = image
1595 7         edges = np.gradient(gray.astype(float))
1596 8         edge_pixels = np.count_nonzero(edges[0]) + np.count_nonzero(edges[1])
1597 9         total_pixels = gray.size
1598 10        return float(edge_pixels / total_pixels)
1599 11
1600 12
1601 13 def feature(image: np.ndarray) -> float:
1602 14     "Computes the brightness of the central region as a percentage of the whole image"
1603 15     if len(image.shape) == 3:
1604 16         gray = np.mean(image, axis=2)
1605 17     else:
1606 18         gray = image
1607 19         h, w = gray.shape
1608 20         center_region = gray[h//4:3*h//4, w//4:3*w//4]
1609 21         total_brightness = np.sum(gray)
1610 22         center_brightness = np.sum(center_region)
1611 23         return float(center_brightness) / (total_brightness + 1e-7) * 100
1612 24
1613 25
1614 26 def feature(image: np.ndarray) -> float:
1615 27     "Determines the ratio of vertical to horizontal gradients in the image, indicating edge
1616 28     direction"
1617 29     if len(image.shape) == 3:
1618 30         gray = np.mean(image, axis=2)
1619 31     else:
1620 32         gray = image
1621 33         grad_y, grad_x = np.gradient(gray)
1622 34         vertical_grad = np.sum(np.abs(grad_x))
1623 35         horizontal_grad = np.sum(np.abs(grad_y))
1624 36         return float(vertical_grad) / (horizontal_grad + 1e-7) # Avoid division by zero
1625 37
1626 38 def feature(image: np.ndarray) -> float:
1627 39     "Calculates the ratio of bright pixels (above a threshold) to total pixels in the image"
1628 40     if len(image.shape) == 3:
1629 41         gray = np.mean(image, axis=2)
1630 42     else:
1631 43         gray = image
1632 44         threshold = 200 # A threshold for bright pixels
1633 45         bright_pixels = np.count_nonzero(gray > threshold)
1634 46         total_pixels = gray.size
1635 47         return float(bright_pixels) / total_pixels

```

```

1620 47
1621 48 def feature(image: np.ndarray) -> float:
1622 49     "Measures the contrast of the image based on the standard deviation of pixel intensities"
1623 50     if len(image.shape) == 3:
1624 51         gray = np.mean(image, axis=2)
1625 52     else:
1626 53         gray = image
1627 54     return float(np.std(gray))

```

## 1627 F.3 TEXT CLASSIFICATION

### 1628 F.3.1 F2

```

1630 1 You are an expert text analysis programmer creating evaluation features for a machine learning
1631 2     model that classifies text.
1632 3
1633 4 # Text Processing API Documentation
1634 5
1635 6 The features receive text as a string, so you can use any string methods and text processing
1636 7     functions.
1637 8 ## String Methods
1638 9 - text.lower(), text.upper(): Case conversion
1639 10 - text.strip(): Remove whitespace
1640 11 - text.split(delimiter): Split into list
1641 12 - text.count(substring): Count occurrences
1642 13 - text.startswith(prefix), text.endswith(suffix): Check prefixes/suffixes
1643 14 - text.find(substring): Find position of substring
1644 15 - text.replace(old, new): Replace text
1645 16
1646 17 ## Text Analysis
1647 18 - len(text): Length of text
1648 19 - text.isdigit(), text.isalpha(), text.isalnum(): Character type checks
1649 20 - sum(1 for c in text if c.isupper()): Count uppercase letters
1650 21 - text.split(): Split on whitespace to get words
1651 22
1652 23 ## Regular Expressions (re module)
1653 24 - re.findall(pattern, text): Find all matches
1654 25 - re.search(pattern, text): Find first match
1655 26 - re.sub(pattern, replacement, text): Replace patterns
1656 27 - len(re.findall(r'\w+', text)): Count words
1657 28 - len(re.findall(r'[!?!]', text)): Count sentences
1658 29
1659 30 ## Useful Patterns
1660 31 - Word count: len(text.split())
1661 32 - Sentence count: text.count('.') + text.count('!?') + text.count('?!')
1662 33 - Average word length: sum(len(word) for word in text.split()) / len(text.split())
1663 34 - Punctuation density: sum(1 for c in text if not c.isalnum() and not c.isspace()) / len(text)
1664 35
1665 36 ## Example Feature Function
1666 37 def feature(text: str) -> float:
1667 38     "Average word length in the text"
1668 39     words = text.split()
1669 40     if not words:
1670 41         return 0.0
1671 42     return sum(len(word) for word in words) / len(words)
1672 43
1673 44 ## Current Feature Database
1674 45 Here are some existing features and their performance (Performance improvement = benefit from
1675 46     that feature, higher is better):
1676 47
1677 48 Feature:
1678 49 def feature(text: str) -> float:
1679 50     'Average number of words per sentence (sentences split on . ! ?)'
1680 51     import re
1681 52     if not text or not text.strip():
1682 53         return float(0.0)
1683 54     # Split on one or more sentence-ending punctuation and filter empties
1684 55     sentences = [s.strip() for s in re.split(r'[!?!]+', text) if s.strip()]
1685 56     if not sentences:
1686 57         return float(0.0)
1687 58     total_words = sum(len(s.split()) for s in sentences)
1688 59     return float(total_words / len(sentences))
1689 60
1690 61 Importance: 0.057
1691 62 ---
1692 63

```

```

1674 64
1675 65
1676 66 Feature:
1677 67 def feature(text: str) -> float:
1678 68     'Average character-uniqueness per word (unique chars / word length), averaged over words'
1679 69     words = [w for w in text.split() if any(ch.isalnum() for ch in w)]
1680 70     if not words:
1681 71         return float(0.0)
1682 72     ratios = []
1683 73     for w in words:
1684 74         chars = [c for c in w if not c.isspace()]
1685 75         if not chars:
1686 76             continue
1687 77         unique = len(set(chars))
1688 78         ratios.append(unique / len(chars))
1689 79     if not ratios:
1690 80         return float(0.0)
1691 81     return float(sum(ratios) / len(ratios))
1692 82
1693 83
1694 84 Importance: 0.057
1695 85 ---
1696 86
1697 87
1698 88 Feature:
1699 89 def feature(text: str) -> float:
1700 90     'Proportion of alphabetic characters that are uppercase'
1701 91     if not text:
1702 92         return float(0.0)
1703 93     alpha_chars = [c for c in text if c.isalpha()]
1704 94     if not alpha_chars:
1705 95         return float(0.0)
1706 96     upper_count = sum(1 for c in alpha_chars if c.isupper())
1707 97     return float(upper_count / len(alpha_chars))
1708 98
1709 99
1710 100 Importance: 0.083
1711 101 ---
1712 102
1713 103
1714 104 Feature:
1715 105 def feature(text: str) -> float:
1716 106     'Proportion of long words (length > 7) among all words'
1717 107     words = [w for w in text.split() if w]
1718 108     if not words:
1719 109         return float(0.0)
1720 110     long_count = sum(1 for w in words if len(w) > 7)
1721 111     return float(long_count / len(words))
1722 112
1723 113
1724 114 Importance: 0.194
1725 115 ---
1726 116
1727 117
1728 118 Feature:
1729 119 def feature(text: str) -> float:
1730 120     'Punctuation characters per word (punctuation = not alnum and not whitespace)'
1731 121     if not text or not text.strip():
1732 122         return float(0.0)
1733 123     words = text.split()
1734 124     if not words:
1735 125         return float(0.0)
1736 126     punct_count = sum(1 for c in text if not c.isalnum() and not c.isspace())
1737 127     return float(punct_count / len(words))
1738 128
1739 129
1740 130 Importance: 0.096
1741 131 ---
1742 132
1743 133
1744 134 Feature:
1745 135 def feature(text: str) -> float:
1746 136     'Proportion of sentences that are questions (based on ? count over total sentence
1747 137     terminators)'
1748 138     if not text or not text.strip():
1749 139         return float(0.0)
1750 140     question_marks = text.count('?')
1751 141     sentence_terminators = text.count('.') + text.count('!') + text.count('?')
1752 142     if sentence_terminators == 0:
1753 143         return float(0.0)
1754 144     return float(question_marks / sentence_terminators)

```

```

1728
1729 144
1730 145 Importance: 0.023
1731 146 ---
1732 147
1733 148 Feature:
1734 149 def feature(text: str) -> float:
1735 150     'Type-token ratio: unique word tokens / total words (case-insensitive, alphanumeric tokens
1736 151     )'
1737 152     import re
1738 153     tokens = re.findall(r'\w+', text.lower())
1739 154     if not tokens:
1740 155         return float(0.0)
1741 156     unique = len(set(tokens))
1742 157     return float(unique / len(tokens))
1743 158
1744 159 Importance: 0.089
1745 160 ---
1746 161
1747 162 Feature:
1748 163 def feature(text: str) -> float:
1749 164     'Ratio of tokens that contain at least one digit'
1750 165     tokens = text.split()
1751 166     if not tokens:
1752 167         return float(0.0)
1753 168     num_with_digit = sum(1 for t in tokens if any(ch.isdigit() for ch in t))
1754 169     return float(num_with_digit / len(tokens))
1755 170
1756 171 Importance: 0.165
1757 172 ---
1758 173
1759 174 Feature:
1760 175 def feature(text: str) -> float:
1761 176     'Longest run of the same character normalized by text length'
1762 177     if not text:
1763 178         return float(0.0)
1764 179     max_run = 1
1765 180     current_run = 1
1766 181     prev = text[0]
1767 182     for c in text[1:]:
1768 183         if c == prev:
1769 184             current_run += 1
1770 185             if current_run > max_run:
1771 186                 max_run = current_run
1772 187         else:
1773 188             current_run = 1
1774 189             prev = c
1775 190     return float(max_run / max(1, len(text)))
1776 191
1777 192 Importance: 0.126
1778 193 ---
1779 194
1780 195 Feature:
1781 196 def feature(text: str) -> float:
1782 197     'Stopword density: fraction of tokens that are common English stopwords'
1783 198     stopwords = {
1784 199         'the', 'and', 'is', 'in', 'it', 'of', 'to', 'a', 'an', 'that', 'this', 'for', 'on', 'with',
1785 200         'as', 'by', 'at', 'from', 'or', 'be', 'are', 'was', 'were', 'has', 'have', 'not', 'but',
1786 201         'they', 'their', 'you', 'I'
1787 202     }
1788 203     tokens = [t.lower().strip(".,!?:\"'()[]") for t in text.split()]
1789 204     if not tokens:
1790 205         return float(0.0)
1791 206     stop_count = sum(1 for t in tokens if t and t in stopwords)
1792 207     return float(stop_count / len(tokens))
1793 208
1794 209 Importance: 0.111
1795 210 ---
1796 211
1797 212 ## Task
1798 213 Generate 10 NEW text features that:
1799 214
1800 215 1. Are different from existing features

```

```

1782 224 2. Capture useful textual patterns
1783 225 3. Return float values
1784 226 4. Handle edge cases gracefully - Won't crash on unusual texts
1785 227 5. Use simple, short docstrings - Use single quotes, not triple quotes
1786 228 6. Are efficient to compute
1787 229
1788 230 Your task is to generate diverse, creative features that capture different aspects of text
1789 231 content for classification. Focus on features that would help distinguish between
1790 232 different text classes. These features will be used as input to a learned model that
1791 233 predicts target values from text.
1792 234
1793 235 ## IMPORTANT CODE REQUIREMENTS
1794 236 - Use SINGLE quotes for docstrings: "description here"
1795 237 - NO triple quotes ("") anywhere in the code
1796 238 - NO markdown code blocks
1797 239 - NO explanatory text after the function
1798 240 - Each function should be complete and standalone
1799 241
1800 242 ## Output Format
1801 243 Generate exactly 10 features in this format:
1802 244
1803 245 def feature(text: str) -> float:
1804 246     "Clear description of what this feature measures"
1805 247     # ... Calculate and return the feature value
1806 248     return float(result)
1807 249
1808 250 def feature(text: str) -> float:
1809 251     "Another feature description"
1810 252     # ... Calculate and return the feature value
1811 253     return float(result)
1812 254
1813 255 The body of the functions can be anything, but the first line (function declaration) should be
1814 256 identical to those examples above (always 'def feature(...)'), and the second line
1815 257 should be a one-line docstring. Don't output explanatory text - just the function
1816 258 definitions as shown above.

```

### 1807 F.3.2 D-ID3 - PROMPT

```

1808 1 You are an expert text analysis programmer creating feature functions to help a machine
1809 2 learning model perform text classification.
1810 3
1811 4 This is a classification task with the following classes: 0: human-written text, 1: AI-
1812 5 generated text.
1813 6
1814 7 Your task is to write a feature function that helps discriminate between the text samples
1815 8 given below.
1816 9
1817 10 A feature function is a Python function that takes a text string and computes a feature out of
1818 11 the text. It should return a float, but note that a feature could also be effectively
1819 12 boolean-valued (0.0 or 1.0), or integer-valued, even if its type is float.
1820 13
1821 14 You have access to the following API from text processing libraries:
1822 15
1823 16 # Text Processing API Documentation
1824 17
1825 18 The features receive text as a string, so you can use any string methods and text processing
1826 19 functions.
1827 20
1828 21 ## String Methods
1829 22 - text.lower(), text.upper(): Case conversion
1830 23 - text.strip(): Remove whitespace
1831 24 - text.split(delimiter): Split into list
1832 25 - text.count(substring): Count occurrences
1833 26 - text.startswith(prefix), text.endswith(suffix): Check prefixes/suffixes
1834 27 - text.find(substring): Find position of substring
1835 28 - text.replace(old, new): Replace text
1836 29
1837 30 ## Text Analysis
1838 31 - len(text): Length of text
1839 32 - text.isdigit(), text.isalpha(), text.isalnum(): Character type checks
1840 33 - sum(1 for c in text if c.isupper()): Count uppercase letters
1841 34 - text.split(): Split on whitespace to get words
1842 35
1843 36 ## Regular Expressions (re module)
1844 37 - re.findall(pattern, text): Find all matches
1845 38 - re.search(pattern, text): Find first match
1846 39 - re.sub(pattern, replacement, text): Replace patterns
1847 40 - len(re.findall(r'\w+', text)): Count words
1848 41 - len(re.findall(r'[!?.]', text)): Count sentences

```

```

1836 37 ## Useful Patterns
1837 38 - Word count: len(text.split())
1838 39 - Sentence count: text.count('.') + text.count('!') + text.count('?')
1839 40 - Average word length: sum(len(word) for word in text.split()) / len(text.split())
1840 41 - Punctuation density: sum(1 for c in text if not c.isalnum() and not c.isspace()) / len(text)
1841 42
1842 43 ## Example Feature Function
1843 44 def feature(text: str) -> float:
1844 45     "Average word length in the text"
1845 46     words = text.split()
1846 47     if not words:
1847 48         return 0.0
1848 49     return sum(len(word) for word in words) / len(words)
1849 50
1850 51
1851 52 # Task
1852 53 Generate 10 new text feature functions in Python that:
1853 54
1854 55 1. Help us discriminate between different text classes, hopefully with samples before and
1855 56     after the optimal split point having the lowest possible variance between their
1856 57     classifications.
1857 58
1858 59 2. Return a float value given a text string.
1859 60
1860 61 3. Handle edge cases gracefully - won't crash on unusual texts
1861 62
1862 63 Your task is to generate diverse, creative features that are relevant to explain the
1863 64     classifications for the text samples shown above. Focus on features that would help
1864 65     distinguish between samples of different classes. These features will be used in this
1865 66     decision tree that will predict the classification of a given text sample. Think about
1866 67     new features that would help such a predictor in the particular cases above, trying to
1867 68     add information that the already existing features shown above are missing.
1868 69
1869 70
1870 71 # Code Requirements
1871 72 - Use single quotes for docstrings: "description here"
1872 73 - No markdown code blocks
1873 74 - No explanatory text after the function
1874 75 - Each function should be complete and standalone, and return a float
1875 76
1876 77 # Output Format
1877 78 Generate exactly 10 features in this format:
1878 79
1879 80 def feature(text: str) -> float:
1880 81     "Simple, clear description of what this feature measures"
1881 82     # ... Calculate and return the feature value
1882 83     return result
1883 84
1884 85 def feature(text: str) -> float:
1885 86     "Another feature description"
1886 87     # ... Calculate and return the feature value
1887 88     return result
1888 89
1889 90 The body of the function can be anything, but the first line (function declaration) should be
1890 91     identical to those examples above, and the second line should be a one-line docstring.
1891 92     Don't output explanatory text - just the function definitions as shown above.
1892 93
1893 94 # Current decision tree node
1894 95 You are currently focusing on features that explain the text classifications in the following
1895 96     subtree of a decision tree:
1896 97
1897 98 [root]
1898 99 -> value > 4.468 for "Average character length of words in the text" -> value > 0.024 for "
1899 100     Calculates the proportion of text that is in passive voice" -> value < 26.000 for "
1900 101     Assesses the use of passive voice constructions in the text"
1901 102
1902 103 # Text samples
1903 104 Here are examples of text samples in this subtree, along with their target classifications:
1904 105
1905 106 Sample: 'The Strange Tank
1906 107
1907 108 John awoke with a start. He was submerged in a tank of pink, viscous liquid. He thrashed
1908 109     around in a panic, trying to determine which way was up. He finally surfaced, gasping for
1909 110     ai...', (Target: 1)
1910 111
1911 112 Sample: 'The old map crumbled at the edges as Juan carefully unrolled it across the table. His
1912 113     grandfather had given him the map many years ago, telling him it showed the true history
1913 114     of this land. A history m...', (Target: 1)
1914 115
1915 116 Sample: 'The call came at midnight, just as I finished cleaning and assembling my sniper rifle
1916 117     . My burner phone buzzed, the one dedicated only to receiving instructions. A robotic
1917 118     female voice spoke: "Contract...', (Target: 1)
1918 119
1919 120 Sample: '"Bang Bang Baby" can be labeled a science fiction musical as it incorporates elements
1920 121     of both genres in an exciting plot. After watching the trailer first, I was surprised by
1921 122     the energetic nature of t...', (Target: 0)

```

1890 97 Sample: 'The red dust swirled around John as he took his first step onto the rust-colored soil  
 1891 of Mars. After decades of training and a perilous 9 month journey through the inky black  
 1892 void of space, he had fin...'. (Target: 1)  
 1893 98  
 1894 99 Optimize for producing discriminant features that are novel compared to the existing features  
 1895 used to arrive at this subtree. Focus on explaining the differences between the text  
 1896 samples shown above.

### 1896 F.3.3 D-ID3 - EXAMPLE FEATURES

```

1897
1898 1 def feature(text: str) -> float:
1899 2     "Average character length of words in the text"
1900 3     words = text.split()
1901 4     if not words:
1902 5         return 0.0
1903 6     return sum(len(word) for word in words) / len(words)
1904 7
1905 8 def feature(text: str) -> float:
1906 9     "Calculates the proportion of text that is in passive voice"
1907 10    passive_pattern = r'\b(?:is|was|were|be|being|been) \w+\b'
1908 11    passive_count = len(re.findall(passive_pattern, text))
1909 12    return float(passive_count) / len(text.split()) if text.split() else 0.0
1910 13
1911 14 def feature(text: str) -> float:
1912 15    "Assesses the use of passive voice constructions in the text"
1913 16    passive_voice = re.findall(r'\b(is|are|was|were|be|being|been)\s+\w+\b', text)
1914 17    return float(len(passive_voice))
1915 18
1916 19 def feature(text: str) -> float:
1917 20    "Counts the number of transition words to evaluate the flow of text"
1918 21    transition_words = set(['however', 'therefore', 'moreover', 'furthermore', 'nevertheless',
1919 22                           'consequently'])
1920 23    words = text.lower().split()
1921 24    count = sum(1 for word in words if word in transition_words)
1922 25    return float(count)
1923 26
1924 27 def feature(text: str) -> float:
1925 28    "Measures the proportion of first-person pronouns in the text"
1926 29    first_person_pronouns = set(['I', 'me', 'my', 'mine', 'we', 'us', 'our', 'ours'])
1927 30    words = text.lower().split()
1928 31    count = sum(1 for word in words if word in first_person_pronouns)
1929    return float(count) / len(words) if words else 0.0
1930
1931
1932
1933
1934
1935
1936
1937
1938
1939
1940
1941
1942
1943

```