# A Comprehensive Review of Model-Based Reinforcement Learning: From Planning with Perfect Models to Navigating Uncertainty

A Detailed Explanation By Taha Majlesi

July 17, 2025

### Abstract

This review provides a structured and comprehensive overview of Model-Based Reinforcement Learning (MBRL). We begin by establishing the fundamental dichotomy between model-free and model-based approaches, highlighting their core principles and inherent trade-offs. The document then explores the idealized scenario of planning with perfect environmental knowledge, contrasting deterministic trajectory optimization with the challenges of stochastic environments. This lays the groundwork for understanding the necessity of closed-loop policies. Subsequently, we delve into practical planning algorithms, including derivative-free optimization methods like the Cross-Entropy Method (CEM) and search-based techniques like Monte Carlo Tree Search (MCTS).

A significant portion of this review is dedicated to the critical challenge of imperfect models. We explain the problem of distributional shift and detail mitigation strategies such as iterative data aggregation and Model Predictive Control (MPC). The review then transitions to the modern paradigm of uncertainty-aware MBRL. We provide a formal taxonomy of uncertainty, distinguishing between aleatoric (inherent) and epistemic (model-based) uncertainty, and explain their distinct implications for agent behavior. We then describe methods for building uncertainty-aware models, focusing on Bayesian Neural Networks and Deep Ensembles. Finally, we synthesize these concepts to illustrate how planning algorithms can be adapted to leverage uncertainty for robust decision-making, exploration, and safety, culminating in a case study of the PETS algorithm. The aim is to provide a clear, accessible, and in-depth understanding of the evolution and current state of MBRL.

## 1 The Model-Based Reinforcement Learning Framework

Reinforcement Learning (RL) offers a mathematical foundation for decision-making through learning, where an agent interacts with an environment to maximize a cumulative reward. A central distinction within RL is between **model-free** and **model-based** methods, which hinges on whether the agent explicitly learns a model of the environment's dynamics.

## 1.1 The Core Dichotomy: Learning a Policy vs. Learning a World

### 1.1.1 Model-Free RL (MFRL): The Direct Approach

In MFRL, the agent learns a behavior directly, without building a model of the world. The environment is treated as a "black box." The agent learns which actions are effective in which states through extensive trial and error, but not necessarily the underlying reasons.

Formally, the agent learns a policy, $\pi_\theta(a_t|s_t)$, parameterized by $\theta$. The goal is to find the optimal parameters $\theta^*$ that maximize the expected cumulative reward:

$$\theta^* = \arg\max_\theta E_{\tau \sim p_\theta(\tau)} \left[ \sum_t r(s_t, a_t) \right]$$

The trajectory distribution is given by $p_\theta(\tau) = p(s_1) \prod_{t=1}^{T} \pi_\theta(a_t|s_t) p(s_{t+1}|s_t, a_t)$. A key aspect is that the true transition dynamics, $p(s_{t+1}|s_t, a_t)$, are unknown and the algorithm does not attempt to learn them.

### 1.1.2 Model-Based RL (MBRL): The Indirect Approach

MBRL takes a more deliberative route, separating the problem into two phases:

1. **Model Learning:** The agent uses its experience to learn an explicit model of the environment's dynamics, $\hat{p}(s'|s, a)$, and often the reward function, $\hat{r}(s, a)$. This is akin to forming a "mental map."

2. **Planning:** The agent uses the learned model to simulate "what-if" scenarios, planning actions by predicting their outcomes without real-world interaction.

For a deterministic environment, model learning can be framed as a supervised learning problem to find a function $s_{t+1} = f(s_t, a_t)$ by minimizing a loss, such as mean squared error, on a dataset of transitions $\mathcal{D} = \{(s, a, s')_i\}$:

$$\min_f \sum_i ||f(s_i, a_i) - s'_i||^2$$

In the stochastic case, the model learns a probability distribution over the next states, $p(s'|s, a)$.

## 1.2 A Comparative Analysis: The Fundamental Trade-Offs

The choice between MFRL and MBRL involves several critical trade-offs, summarized in Table 1.

The distinction is not a rigid binary but a spectrum. Hybrid algorithms like **Dyna-Q** use a learned model to generate simulated experiences to augment the data for a model-free algorithm, blurring the lines between the two paradigms.

# 2 Planning with Perfect World Knowledge

To understand the complexities of planning with a learned (and thus imperfect) model, it is useful to first consider the idealized case where the agent has a perfect model of the environment.

Table 1: Comparative Analysis of Model-Free vs. Model-Based RL

| Characteristic | Model-Free RL (MFRL) | Model-Based RL (MBRL) |
|---|---|---|
| **Core Principle** | Learns a policy or value function directly from experience. | Learns a model of the environment, then uses the model to plan. |
| **Sample Efficiency** | Low (sample-inefficient). Requires many real-world interactions. | High (sample-efficient). Can leverage the model for simulated experience. |
| **Asymptotic Performance** | Generally higher. Not limited by model accuracy. | Capped by the accuracy of the learned model (model bias). |
| **Computational Profile** | Simpler architecture. Fast inference once trained. | More complex. Planning can be computationally intensive at decision time. |
| **Adaptability to Change** | Low. Policy is often task-specific and requires retraining. | High. Can often adapt to new tasks by re-planning with the existing model. |
| **Exploration Strategy** | Often relies on simple heuristics (e.g., $\epsilon$-greedy). | Can use the model to guide exploration towards uncertain states. |
| **Primary Challenge** | High sample complexity. | Overcoming model bias and compounding errors. |

## 2.1 Trajectory Optimization in Deterministic Systems

In a deterministic world, where $s_{t+1} = f(s_t, a_t)$ is known, planning becomes a trajectory optimization problem. The goal is to find a sequence of actions $A = (a_1, ..., a_T)$ that maximizes the cumulative reward:

$$\max_{a_1, ..., a_T} \sum_{t=1}^{T} r(s_t, a_t) \quad \text{s.t.} \quad s_{t+1} = f(s_t, a_t)$$

This is a well-studied problem in optimal control, with solution methods like shooting or collocation.

## 2.2 The Folly of Open-Loop Planning in a Stochastic World

When the dynamics are stochastic, $p(s_{t+1}|s_t, a_t)$, a naive approach is **open-loop planning**, where a fixed sequence of actions is determined beforehand. This approach is brittle because it cannot react to unexpected outcomes. An illustrative Persian anecdote tells of a deaf man who rehearses a conversation with a sick friend. His plan fails comically and tragically when the friend's answers deviate from his expectations, highlighting the failure of non-adaptive planning.

## 2.3 The Necessity of Closed-Loop (State-Contingent) Policies

The solution is a **closed-loop policy**, $\pi(a_t|s_t)$, which maps the current state to an action. This allows the agent to react dynamically to the actual evolution of the environment. The objective of RL is to find an optimal policy $\pi^*$ that maximizes the expected reward over all possible trajectories:

$$\pi^* = \arg\max_\pi E_{\tau \sim p(\tau)} \left[ \sum_t r(s_t, a_t) \right]$$

where the trajectory distribution is $p(\tau) = p(s_1) \prod_{t=1}^{T} \pi(a_t|s_t) p(s_{t+1}|s_t, a_t)$.

# 3 Planning as Optimization: Derivative-Free Methods

When a dynamics model is available but its gradients are not, planning can be framed as a black-box optimization problem.

## 3.1 The "Guess and Check" Paradigm: Random Shooting

This simple method involves:

1. Sampling $N$ candidate action sequences.

2. Simulating each sequence with the model to get its total reward.

3. Executing the sequence with the highest reward.

Its efficiency degrades exponentially with the planning horizon and action space size.

## 3.2 The Cross-Entropy Method (CEM): An Iterative Refinement

CEM is a more sophisticated derivative-free method that iteratively refines a sampling distribution over action sequences. The algorithm proceeds as follows:

1. **Initialization:** Start with a parameterized distribution over action sequences, typically a Gaussian $p_\phi(A) = \mathcal{N}(\mu, \Sigma)$.

2. **Sample:** Draw $N$ candidate action sequences $\{A_i\}$ from $p_\phi(A)$.

3. **Evaluate:** Calculate the total reward $J(A_i)$ for each sequence using the model.

4. **Select Elites:** Identify the top-performing fraction of sequences (the "elites").

5. **Refit Distribution:** Update the distribution parameters $\phi$ to maximize the likelihood of the elite samples. For a Gaussian, this means calculating the sample mean and covariance of the elites.

6. **Iteration:** Repeat until convergence.

CEM is simple and parallelizable but, as an open-loop method, is vulnerable to stochasticity and can converge to local optima.

# 4 Planning as Search: Monte Carlo Tree Search (MCTS)

For problems with discrete action spaces, MCTS is a powerful heuristic search algorithm that builds a search tree to find the best action.

## 4.1 The Four Phases of an MCTS Iteration

MCTS iteratively refines its value estimates through a four-step cycle:

1. **Selection:** Traverse the existing tree from the root by selecting children that balance exploration and exploitation until a leaf node is reached.

2. **Expansion:** Create new child nodes for untried actions from the leaf node.

3. **Simulation (Rollout):** From a new node, run a simulation with a simple default policy (e.g., random actions) to the end of the episode to get a reward estimate.

4. **Backpropagation:** Update the visit counts ($N$) and value estimates ($Q$) of the nodes on the selection path with the result of the simulation.

## 4.2 The UCT Formula: Principled Exploration-Exploitation

The selection phase is often guided by the **Upper Confidence Bound 1 for Trees (UCT)** formula, which selects the action that maximizes:

$$\text{Score}(s_{t+1}) = \underbrace{\frac{Q(s_{t+1})}{N(s_{t+1})}}_{\text{Exploitation: Mean Reward}} + \underbrace{C\sqrt{\frac{\ln N(s_t)}{N(s_{t+1})}}}_{\text{Exploration: Uncertainty Bonus}}$$

This formula elegantly balances choosing actions that have been successful in the past (exploitation) with trying less-visited actions (exploration).

# 5 The Challenge of Imperfect Models: Distributional Shift

In practice, the model is learned from finite data and is therefore imperfect. This introduces the critical problem of **distributional shift**.

## 5.1 The Inevitable Failure of Naïve MBRL

A naive approach of collecting a static dataset, learning a model, and then planning with that model is doomed to fail. The model is trained on states visited by an initial policy, but the planner derives a new policy that visits a different distribution of states. The model is thus evaluated on out-of-distribution data, where its predictions are unreliable. This creates a vicious cycle where the planner exploits model inaccuracies, leading the agent further into unseen states and causing errors to compound.

## 5.2 Mitigation Strategy 1: Iterative Data Aggregation (DAgger for RL)

To combat distributional shift, the learning process can be made iterative. This approach, analogous to the DAgger algorithm, involves:

1. Learning a model from the current dataset.

2. Using the model to derive an improved policy.

3. Using the new policy to collect more data in the real environment.

4. Aggregating the new data into the dataset and repeating.

## 5.3 Mitigation Strategy 2: Model Predictive Control (MPC) for Error Correction

To prevent the accumulation of small prediction errors over long horizons, **Model Predictive Control (MPC)** uses a replanning strategy:

1. From the current real state $s_t$, plan an optimal action sequence over a finite horizon $H$.

2. Execute only the **first** action, $a_t$.

3. Observe the true next state, $s_{t+1}$, from the environment.

4. Discard the rest of the plan and repeat the process from $s_{t+1}$.

By constantly re-grounding the plan in reality, MPC makes the system much more robust to model inaccuracies.

# 6 Quantifying the Unknown: Uncertainty in Model-Based RL

The key to high-performance MBRL is to build models that can quantify and reason about their own uncertainty.

## 6.1 A Formal Taxonomy of Uncertainty

Uncertainty can be decomposed into two types:

- **Aleatoric Uncertainty (Data Uncertainty):** Inherent, irreducible randomness in the environment (e.g., a coin flip). It cannot be reduced by collecting more data. It signals **risk** and should drive cautious behavior.

- **Epistemic Uncertainty (Model Uncertainty):** Uncertainty due to the model's lack of knowledge from limited data. It is reducible by collecting more data in uncertain regions. It signals **ignorance** and should drive exploration.

The total predictive variance can be decomposed as:

$$\text{Var}(y|x, \mathcal{D}) = \underbrace{E_{p(\theta|\mathcal{D})}[\text{Var}(y|x, \theta)]}_{\text{Aleatoric}} + \underbrace{\text{Var}_{p(\theta|\mathcal{D})}[E(y|x, \theta)]}_{\text{Epistemic}}$$

Table 2: A Taxonomy of Uncertainty in Reinforcement Learning

| Feature | Aleatoric Uncertainty | Epistemic Uncertainty |
|---|---|---|
| **Also Known As** | Statistical, Data, Irreducible Uncertainty | Model, Structural, Reducible Uncertainty |
| **Definition** | Inherent randomness in the data-generating process. | Uncertainty in model parameters due to lack of knowledge. |
| **Source** | Stochasticity of the environment, sensor noise. | Limited or sparse training data. |
| **Reducibility** | No. Cannot be reduced by collecting more data. | Yes. Can be reduced by collecting more data. |
| **Implication for Agent** | Signals inherent risk. Should drive caution. | Signals ignorance. Should drive exploration. |

## 6.2 Building Uncertainty-Aware Dynamics Models

Two primary methods are used to build deep learning models that quantify uncertainty.

### 6.2.1 Method 1: Bayesian Neural Networks (BNNs)

BNNs learn a probability distribution over each network weight instead of a single point estimate. Predictions are made by averaging over all possible models, weighted by their posterior probability. This is computationally expensive and complex to train, often relying on approximate inference techniques like Variational Inference (VI).

### 6.2.2 Method 2: Deep Ensembles

A simpler, often more effective approach is to train an ensemble of $N$ individual neural networks. Diversity is encouraged through different random initializations and bootstrapped data.

- The **mean** of the ensemble's predictions is used as the final prediction.

- The **variance** of the predictions across the ensemble members quantifies the epistemic uncertainty.

Deep ensembles are simple to implement and highly parallelizable, and they produce high-quality uncertainty estimates.

# 7 Planning Under Uncertainty

With uncertainty-aware models, planning algorithms can be made more robust and intelligent.

## 7.1 Uncertainty-Aware Trajectory Optimization

Trajectory optimization methods like CEM can be adapted by modifying the evaluation step. Instead of a single rollout, an action sequence is evaluated by simulating it across all

$N$ models in a deep ensemble. The objective becomes maximizing the **average reward** across this ensemble of predicted futures:

$$J(A) = \frac{1}{N} \sum_{i=1}^{N} \left( \sum_{t=1}^{H} r(s_{t,i}, a_t) \right), \quad \text{where } s_{t+1,i} = f_i(s_{t,i}, a_t)$$

This implicitly penalizes action sequences that lead to high disagreement among the models (high epistemic uncertainty), thus avoiding the exploitation of model flaws.

## 7.2 Case Study: The PETS Algorithm

The **Probabilistic Ensembles with Trajectory Sampling (PETS)** algorithm successfully integrates these concepts:

- It learns an ensemble of probabilistic neural networks as its dynamics model.

- It uses an MPC framework, replanning at each step.

- It uses CEM to optimize the ensemble-averaged reward.

PETS has been shown to match the high performance of top model-free algorithms while requiring orders of magnitude less data.

## 7.3 The Dual-Objective Planner: Balancing Safety and Exploration

A sophisticated planner can use the disentangled uncertainty signals to optimize a dual objective:

- **Uncertainty-Driven Exploration:** Add an exploration bonus proportional to epistemic uncertainty to incentivize gathering informative data.

- **Risk-Averse Planning:** Add a penalty proportional to aleatoric uncertainty to avoid inherently risky states, crucial for safety-critical applications.

# 8 Conclusion

Model-Based Reinforcement Learning has evolved from brittle methods plagued by distributional shift and compounding errors to sophisticated frameworks that achieve high performance. Key developments like iterative data aggregation and Model Predictive Control provided necessary robustness. The true breakthrough, however, has been the explicit incorporation of uncertainty. By building models that quantify their own epistemic and the environment's aleatoric uncertainty, modern MBRL algorithms can plan more intelligently. This allows them to avoid exploiting model flaws, drive efficient exploration, and behave cautiously in the face of inherent risk. Algorithms like PETS demonstrate that it is possible to combine the sample efficiency of model-based learning with the high asymptotic performance of model-free methods, paving the way for more capable and reliable autonomous systems.