

# Optimization and Regularization Under Arbitrary Objectives

Jared Lakhani  
Supervisor: Dr Etienne Pienaar

Introduction

Metropolis-Hastings

The Objective

The Likelihood

The Navigation Problem

XO

References

# Introduction

Introduction

Metropolis-  
Hastings

The Objective

The  
Likelihood

The  
Navigation  
Problem

XO

References

## Update-Step:

- ▶ Given current state  $\boldsymbol{\theta}^{(j)} \in \mathbb{R}^S$ .
- ▶ Propose candidate  $\boldsymbol{\theta}^*$  obtained from  $\boldsymbol{\theta}^* = \boldsymbol{\theta}^{(j)} + \mathbf{Q}$ .

## Acceptance step:

$$\boldsymbol{\theta}^{(j+1)} = \begin{cases} \boldsymbol{\theta}^*, & \text{if } U < \alpha \\ \boldsymbol{\theta}^{(j)}, & \text{otherwise.} \end{cases}$$

- ▶  $U \sim \text{Unif}(0, 1)$  induces randomness, so not all “better” proposals are accepted = better exploration of the parameter space.

## Acceptance probability $\alpha$ :

$$\alpha = \min \left( \frac{p(\boldsymbol{\theta}^* | \mathcal{D})}{p(\boldsymbol{\theta}^{(j)} | \mathcal{D})} \cdot \frac{Q(\boldsymbol{\theta}^{(j)} | \cancel{\boldsymbol{\theta}^*})}{\cancel{Q(\boldsymbol{\theta}^* | \boldsymbol{\theta}^{(j)})}}, 1 \right).$$

For symmetric proposals, e.g.  $\mathbf{Q} \sim \mathcal{N}(\mathbf{0}, \boldsymbol{\Sigma}_s)$ :

$$\begin{aligned} Q(\boldsymbol{\theta}^* | \boldsymbol{\theta}^{(j)}) &= \frac{1}{(2\pi)^{S/2} |\boldsymbol{\Sigma}|^{1/2}} \exp\left(-\frac{1}{2} (\boldsymbol{\theta}^* - \boldsymbol{\theta}^{(j)})^\top \boldsymbol{\Sigma}^{-1} (\boldsymbol{\theta}^* - \boldsymbol{\theta}^{(j)})\right) \\ &= Q(\boldsymbol{\theta}^{(j)} | \boldsymbol{\theta}^*). \end{aligned}$$

$$\alpha = \min \left( \underbrace{\frac{p(\mathcal{D} | \boldsymbol{\theta}^*)}{p(\mathcal{D} | \boldsymbol{\theta}^{(j)})}}_{\text{Likelihood ratio}} \cdot \underbrace{\frac{p(\boldsymbol{\theta}^*)}{p(\boldsymbol{\theta}^{(j)})}}_{\text{Prior ratio}}, 1 \right).$$

- ▶ Higher values  $\Rightarrow \boldsymbol{\theta}^*$  fits the data better (via the likelihood).
- ▶ Higher values  $\Rightarrow \boldsymbol{\theta}^*$  is more consistent with prior beliefs.
- ▶  $\boldsymbol{\theta}^*$  may not score highly on both the likelihood and prior.
- ▶ Even if both ratios are high,  $\boldsymbol{\theta}^*$  may still be rejected if  $U > \alpha$ .
- ▶ If  $p(\boldsymbol{\theta}^* | \mathcal{D}) \geq p(\boldsymbol{\theta}^{(j)} | \mathcal{D})$ , then  $\boldsymbol{\theta}^*$  is accepted automatically.

- ▶ Define: Maximum a Posteriori estimate:  $\hat{\theta}^{MAP} = \operatorname{argmax}_{\theta} p(\theta \mid \mathcal{D})$  = the mode of the posterior.

$\theta^{(j)}$  vs. MCMC iterations, and formation of  $p(\theta \mid \mathcal{D})$ .

## 2-Block MCMC (Bayesian Hierarchical)

- ▶ Gaussian prior for parameters  $\boldsymbol{\theta} \in \mathbb{R}^S \implies \boldsymbol{\theta} \sim \mathcal{N}(\mathbf{o}, \sigma_{\theta}^2 \mathbf{I}_S)$ .
- ▶ Express uncertainty about dispersion  $\sigma_{\theta}^2 \rightarrow$  give it its own distribution to sample from.
- ▶  $\sigma_{\theta}^2$  now also needs a prior  $\rightarrow$ , assume  $\sigma_{\theta}^2 \sim \text{Inv-Gamma}(a, b)$  with  $a, b \approx 0 \implies$  nearly uninformative.
- ▶  $(\sigma_{\theta}^2)^{(j+1)} \mid \boldsymbol{\theta}^{(j+1)}, \mathcal{D} \sim \text{Inv-Gamma}\left(a + \frac{S}{2}, b + \frac{\|\boldsymbol{\theta}^{(j+1)}\|^2}{2}\right) \rightarrow$  Gibbs sampling i.e. directly sample from known distribution.

Introduction

Metropolis-  
Hastings

The Objective

The  
Likelihood

The  
Navigation  
Problem

XO

References

## 2-Block MCMC (Bayesian Hierarchical)

Introduction

Metropolis-Hastings

The Objective

The Likelihood

The Navigation Problem

XO

References

- ▶ 2-Block MCMC samples from the joint,  $p(\boldsymbol{\theta}, \sigma_\theta^2 \mid \mathcal{D})$ .
- ▶ By **alternately** drawing samples from the conditionals  $p(\boldsymbol{\theta} \mid \sigma_\theta^2, \mathcal{D})$  and  $p(\sigma_\theta^2 \mid \boldsymbol{\theta}, \mathcal{D})$ .
- ▶ I.e.  $p(\boldsymbol{\theta} \mid \sigma_\theta^2, \mathcal{D}) = \frac{p(\boldsymbol{\theta}, \sigma_\theta^2 \mid \mathcal{D})}{p(\sigma_\theta^2 \mid \mathcal{D})}$ , when sampling  $\boldsymbol{\theta}$  from  $p(\boldsymbol{\theta} \mid \sigma_\theta^2, \mathcal{D})$ ,  $\sigma_\theta^2$  is fixed  $\implies p(\boldsymbol{\theta} \mid \sigma_\theta^2, \mathcal{D}) \propto p(\boldsymbol{\theta}, \sigma_\theta^2 \mid \mathcal{D})$ .
- ▶ I.e.  $p(\sigma_\theta^2 \mid \boldsymbol{\theta}, \mathcal{D}) = \frac{p(\boldsymbol{\theta}, \sigma_\theta^2 \mid \mathcal{D})}{p(\boldsymbol{\theta} \mid \mathcal{D})}$ , when sampling  $\sigma_\theta^2$  from  $p(\sigma_\theta^2 \mid \boldsymbol{\theta}, \mathcal{D})$ ,  $\boldsymbol{\theta}$  is fixed  $\implies p(\sigma_\theta^2 \mid \boldsymbol{\theta}, \mathcal{D}) \propto p(\boldsymbol{\theta}, \sigma_\theta^2 \mid \mathcal{D})$ .
- ▶ Later we'll show  $\nu \propto \frac{1}{\sigma_\theta^2} \implies \hat{\boldsymbol{\theta}}^{MAP}$  inherently has degree of regularization 'baked in' since  
 $p(\boldsymbol{\theta} \mid \mathcal{D}) = \int p(\boldsymbol{\theta}, \sigma_\theta^2 \mid \mathcal{D}) d\sigma_\theta^2 \rightarrow$  training data infers regularization!

- ▶ To improve mixing (reduced autocorrelation between samples) → better match the geometry of the stationary distribution with proposal distribution.
- ▶ Proposal for  $\boldsymbol{\theta} \in \mathbb{R}^S$  at current iteration  $j \rightarrow \mathbf{Q}_{\boldsymbol{\theta}} \sim \mathcal{N}(\boldsymbol{\theta}^{(j)}, \Sigma_j)$ .
- ▶  $\Sigma_j = \text{Cov}(f(\boldsymbol{\theta}^{(1)}, \boldsymbol{\theta}^{(2)}, \dots, \boldsymbol{\theta}^{(j)})) + \epsilon \mathbf{I}$ .

$$f(\boldsymbol{\theta}^{(1)}, \boldsymbol{\theta}^{(2)}, \dots, \boldsymbol{\theta}^{(j)}) = \begin{cases} (\boldsymbol{\theta}^{(j-\delta\Delta)}, \boldsymbol{\theta}^{(j-\delta\Delta+\delta)}, \dots, \boldsymbol{\theta}^{(j-\delta\Delta+(\Delta-1)\delta)}, \boldsymbol{\theta}^{(j)}), & \text{if } j > \delta\Delta \\ (\boldsymbol{\theta}^{(j-\lfloor \frac{j}{\delta} \rfloor \delta)}, \boldsymbol{\theta}^{(j-\lfloor \frac{j}{\delta} \rfloor \delta+\delta)}, \dots, \boldsymbol{\theta}^{(j-\lfloor \frac{j}{\delta} \rfloor \delta+(\lfloor \frac{j}{\delta} \rfloor - 1)\delta)}, \boldsymbol{\theta}^{(j)}), & \text{if } \delta < j \leq \delta\Delta \\ (\boldsymbol{\theta}^{(1)}, \boldsymbol{\theta}^{(2)}, \dots, \boldsymbol{\theta}^{(j)}), & \text{if } j \leq \delta. \end{cases}$$

- ▶  $\delta$  serves as a stride parameter → non-consecutive  $\boldsymbol{\theta}$  used (thinning).
- ▶  $\Delta \rightarrow$  window of accepted proposals used (eases computation).
- ▶  $\Sigma_j$  now asymmetric → use Adaptive MH only in burn-in.

# Adaptive Scaling

Introduction

Metropolis-Hastings

The Objective

The Likelihood

The Navigation Problem

XO

References

- ▶ Want acceptance probability  $\alpha$  toward theoretical optimum of  $\approx 0.234$  per Roberts et al. 1997.
- ▶ Scale  $\Sigma_j \implies (\varsigma^2)^{(j)} \Sigma_j$ .

$$(\varsigma^2)^{(j+1)} = (\varsigma^2)^{(j)} \times \exp(\gamma_j \cdot (\alpha_j - 0.234)).$$

- ▶  $\gamma_j$  is a sequence of diminishing adaptation rates  $\rightarrow \gamma_j = \frac{1}{j^\kappa}$  per Roberts et al. 1997.
- ▶  $(\varsigma^2)^{(j)} \Sigma_j$  now asymmetric  $\rightarrow$  use adaptive scaling only in burn-in.

# The Objective & Regularization

Introduction

Metropolis-  
Hastings

The Objective

The  
Likelihood

The  
Navigation  
Problem

XO

References

- ▶ L2 penalized objective:  $\operatorname{argmax}_{\boldsymbol{\theta}} \left( \text{Obj}(\boldsymbol{\theta}) - \nu \sum_{i=1}^S \theta_i^2 \right)$ .
- ▶ For Gaussian prior:  $\boldsymbol{\theta} \sim \mathcal{N}(\mathbf{o}, \sigma_{\theta}^2 \mathbf{I}_S) \rightarrow \log[p(\boldsymbol{\theta} | \mathcal{D})] \propto \log [p(\mathcal{D} | \boldsymbol{\theta})] - \frac{1}{2\sigma_{\theta}^2} \sum_{i=1}^S \theta_i^2$ .
- ▶ If we assume  $p(\mathcal{D} | \boldsymbol{\theta}) \propto \text{Obj}(\boldsymbol{\theta}) \rightarrow \log[p(\boldsymbol{\theta} | \mathcal{D})] \propto \text{Obj}(\boldsymbol{\theta}) - \nu \sum_{i=1}^S \theta_i^2$ , with  $\nu \propto \frac{1}{\sigma_{\theta}^2}$ .
- ▶  $\implies \hat{\boldsymbol{\theta}}^{MAP} = \operatorname{argmax}_{\boldsymbol{\theta}} p(\boldsymbol{\theta} | \mathcal{D}) = \operatorname{argmax}_{\boldsymbol{\theta}} \left( \text{Obj}(\boldsymbol{\theta}) - \nu \sum_{i=1}^S \theta_i^2 \right)$ .

# The Objective & Regularization

Introduction

Metropolis-Hastings

The Objective

The Likelihood

The Navigation Problem

XO

References

- ▶ **2-Block MCMC:** since we sample  $\sigma_\theta^2$  at each iteration  $\implies$  L2 shrinkage  $\nu^{(j)} \propto \frac{1}{(\sigma_\theta^2)^{(j)}}$  on each proposal  $\boldsymbol{\theta}^*$ .
- ▶ Marginal  $p(\sigma_\theta^2 | \mathcal{D})$  reflects posterior uncertainty about shrinkage strength.
- ▶ Previously: fixed prior  $\boldsymbol{\theta} \sim \mathcal{N}(\mathbf{o}, \sigma_\theta^2 \mathbf{I}_S)$ .
- ▶ Now: hierarchical prior  $\boldsymbol{\theta} | \sigma_\theta^2 \sim \mathcal{N}(\mathbf{o}, \sigma_\theta^2 \mathbf{I}_S)$  with  $\sigma_\theta^2 \sim \text{Inv-Gamma}(a, b) \implies$  marginal prior  $\boldsymbol{\theta} \sim t_{2a}\left(\mathbf{o}, \frac{b}{a} \mathbf{I}_S\right)$ .
- ▶ MAP estimator:

$$\hat{\boldsymbol{\theta}}^{MAP} = \underset{\boldsymbol{\theta}}{\operatorname{argmax}} \left( \text{Obj}(\boldsymbol{\theta}) - \left( a + \frac{S}{2} \right) \log \left( b + \frac{1}{2} \sum_{i=1}^S \theta_i^2 \right) \right).$$

- ▶ Student- $t$  prior  $\approx$  Gaussian prior near  $\mathbf{0}$  **but** much weaker shrinkage for large  $|\theta_i|$ .

- ▶ **Classical case:** Likelihood = joint density of data given parameters:

$$p(\mathcal{D} \mid \boldsymbol{\theta}) = f(x_1, \dots, x_n \mid \boldsymbol{\theta}) = \prod_{i=1}^n f(x_i \mid \boldsymbol{\theta}).$$

- ▶ **Arbitrary objective:** No assumptions on data-generating model → only an objective  $\text{Obj}(\boldsymbol{\theta})$  to maximise (e.g., accuracy, score, ROI).

- ▶ **Idea:** Define a **pseudo-likelihood**:

$$p(\mathcal{D} \mid \boldsymbol{\theta}) \propto \text{Obj}(\boldsymbol{\theta}).$$

- ▶ Target the conditional **pseudo-posterior**:

$$p(\boldsymbol{\theta} \mid \sigma_\theta^2, \mathcal{D}) \propto p(\mathcal{D} \mid \boldsymbol{\theta}) \cdot p(\boldsymbol{\theta} \mid \sigma_\theta^2).$$

repurposing MCMC as a mode-seeking algorithm → sample around dominant mode of conditional.

- ▶ Acceptance guided by  $\text{Obj}(\boldsymbol{\theta})$  → higher objective ⇒ higher acceptance.
- ▶ **Tempering:**  $[p(\boldsymbol{\theta} \mid \mathcal{D})]^\beta$  sharpens sensitivity to  $\text{Obj}(\boldsymbol{\theta})$ .

# The Navigation Problem

- ▶ **Arena setup:**  $T$  drones navigate a 2D annular region bounded by inner radius  $R_{inner}$  and outer radius  $R_{outer}$ .
- ▶ **Objective:** Each drone attempts to exit the arena within  $K$  steps, with fixed step length  $\delta$ .
- ▶ **Crash rule:** A drone crashes (and ceases navigation) if it comes within distance  $R_{crash}$  of any of the  $J$  orbiting obstacles.

Introduction

Metropolis-Hastings

The Objective

The Likelihood

The Navigation Problem

XO

References

- ▶ Denote: drone and obstacle coordinates  $(\mathbf{x}_t^{(k)}, \mathbf{o}_j^{(k)})$  and detection radius  $R_{detection} = \phi(\boldsymbol{\theta}_1)$ .
- ▶ Network inputs for the  $t^{\text{th}}$  drone:  $\Omega(\mathbf{x}_t^{(k)}, \mathbf{o}_j^{(k)}, \phi(\boldsymbol{\theta}_1))$ , defined as the sum of reciprocal Euclidean distances to obstacles within  $R_{detection}$ , capturing the caution level needed.
- ▶ Network outputs: position updates  $(\tilde{x}_{t_1}^{(k)}, \tilde{x}_{t_2}^{(k)})$  for the  $t^{\text{th}}$  drone at step  $k$ .
- ▶ Output activation:  $\sigma_L(\cdot) = \tanh$ , enabling movement in all planar directions.
- ▶ Full mapping:

$$(\mathbf{x}_t^{(k)}, \mathbf{o}_j^{(k)}, \boldsymbol{\theta}) \longrightarrow \text{model}(\Omega(\mathbf{x}_t^{(k)}, \mathbf{o}_j^{(k)}, \phi(\boldsymbol{\theta}_1)), \boldsymbol{\theta}_2) \xrightarrow{\sigma_L(\cdot)} (\tilde{x}_{t_1}^{(k)}, \tilde{x}_{t_2}^{(k)}) \in [-I, I]^2$$

where  $\boldsymbol{\theta}_1 \subset \boldsymbol{\theta}, \boldsymbol{\theta}_2 \subseteq \boldsymbol{\theta}$

- ▶ For a parameter configuration  $\boldsymbol{\theta} \in \mathbb{R}^S$ , where  $\boldsymbol{\theta}_2$  denotes network weights and biases and  $R_{detection} = \phi(\boldsymbol{\theta}_1)$ , record relative success frequency of the  $T$  drones,  $k(\boldsymbol{\theta})$ :

$$\hat{\boldsymbol{\theta}} = \underset{\boldsymbol{\theta}}{\operatorname{argmax}} \frac{I}{T} k(\boldsymbol{\theta}).$$

## Navigation: The effects of regularization

- ▶ **Environment initialization:** At iteration  $k$ , obstacles are placed at coordinates  $\mathbf{o}_j^{(k)}, j = 1, \dots, J$ , determined by angular frequencies  $\omega_j$ , phase shifts  $\phi_j$ , and orbital radii  $r_j$ . Also specifies  $T$  initial drone positions  $\mathbf{x}_t^{(o)}, t = 1, \dots, T$ .
- ▶ **Training & testing setup:** Initialization is controlled by a training seed  $\omega_0^{\text{Train}}$ ; evaluation uses 1000 test initializations with seeds  $\{\omega_j^{\text{Test}}\}_{j=1}^{1000}$ .
- ▶ **Preliminary regularization assessment:** Estimate  $\hat{\theta}_\nu^{\text{GA}}$  using a GA with L2 regularization of strength  $\nu \rightarrow$  apply it to a newly initialized environment.
- ▶ **Result:** No consistent pattern in success rates observed → performance driven mainly by whether  $\hat{R}_{\text{detection}} \approx 0.07$ .

$\nu \times 10^6$	$k(\hat{\theta}_\nu^{\text{GA}})$	Out-of-Sample		
		$\tilde{k}(\hat{\theta}_\nu^{\text{GA}})$	$\bar{k}(\hat{\theta}_\nu^{\text{GA}})$	$\hat{R}_{\text{detection}}$
0	100	<b>87.00</b>	<b>78.67</b>	<b>0.0726</b>
1	98	54.00	50.61	0.1056
2	100	15.50	29.77	0.1206
3	98	60.00	54.6	0.0983
4	100	84.00	73.90	0.1075
5	99	66.00	58.56	0.0962
6	98	70.50	61.09	0.1016
7	100	<b>86.00</b>	<b>78.47</b>	<b>0.0721</b>
8	100	<b>87.00</b>	<b>79.34</b>	<b>0.0724</b>
9	100	80.00	69.20	0.1161

**Table:** Number of successes  $k(\hat{\theta}_\nu^{\text{GA}})$  for the in-sample initialization, median ( $\tilde{k}$ ) and mean ( $\bar{k}$ ) for the distributions of successes on the 1000 test initializations and estimated  $\hat{R}_{\text{detection}}$ , against varying  $\nu$  for  $T = 100$  and  $R_{\text{crash}} = 0.05$ .

Introduction  
Metropolis-Hastings  
The Objective  
The Likelihood  
The Navigation Problem

XO  
References



# Navigation: MCMC

- ▶ Recall: pseudo-likelihood  $\rightarrow p(\mathcal{D} \mid \boldsymbol{\theta}) \propto \text{Obj}(\boldsymbol{\theta})$ .
- ▶  $\rightarrow$  repurposed MCMC as a mode-seeking algorithm  $\implies$  samples concentrated around the dominant mode of conditional  $p(\boldsymbol{\theta} \mid \sigma_\theta^2, \mathcal{D})$ .
- ▶ Employ three variants of pseudo-likelihood, differing in their degree of sharpness.
- ▶ Recall: MH operates on ratios of likelihoods  $\rightarrow$  rates of change more important than absolute scale.

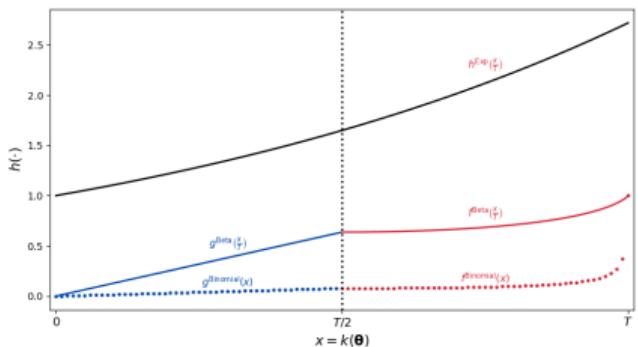


Figure:  $b^{\text{Binomial}}(x)$ ,  $b^{\text{Beta}}\left(\frac{x}{T}\right)$  and  $b^{\text{Exp}}\left(\frac{x}{T}\right)$  for  $x = k(\boldsymbol{\theta}) \in \{0, 1, \dots, T\}$ .

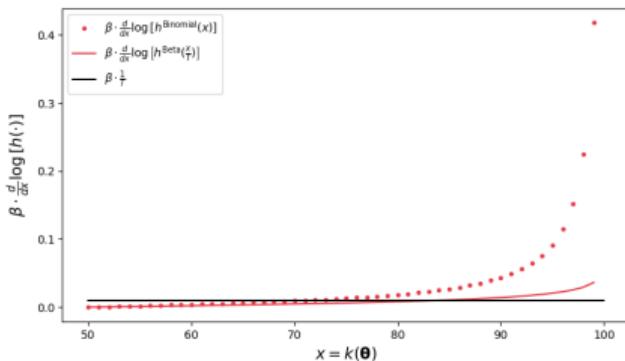


Figure:  $\beta \cdot \frac{d}{dx} \log [h(\cdot)]$  for  $x = k(\boldsymbol{\theta}) \in \{\frac{T}{2}, \frac{T}{2} + 1, \dots, T\}$  for  $\beta = 1$ .

Introduction

Metropolis-Hastings

The Objective

The Likelihood

The Navigation Problem

XO

References

- ▶ Sharpest likelihood = best in-sample performance.
- ▶ Out-of-sample performance remains highly sensitive to whether  $\hat{R}_{\text{detection}} \approx 0.07$ .
- ▶ ESS > 100  $\implies$  satisfactory mixing (expected, since adaptive MH used).

Likelihood Type	In-Sample		Out-of-Sample		
	$k(\hat{\theta}_\nu^{\text{GA}})$	$\tilde{k}(\hat{\theta}_\nu^{\text{GA}})$	$\bar{k}(\hat{\theta}_\nu^{\text{GA}})$	$\hat{R}_{\text{detection}}$	ESS
Binomial-based	<b>100</b>	13.00	24.45	0.3912	179.8550
Beta-based	50	8.00	12.72	0.5175	170.0205
Exponential-based	52	12.00	17.25	0.5601	317.2764

**Table:** Number of successes  $k(\hat{\theta}_\nu^{\text{GA}})$  for the in-sample initialization, summary statistics for the distributions of successes on the 1000 test initializations and estimated  $\hat{R}_{\text{detection}}$ , for the three likelihood types.

## Navigation: MCMC

- ▶ Indefinite contraction of  $\|\boldsymbol{\theta}^{(j)}\| \rightarrow$  for flat pseudo-likelihoods, it can be shown:

$$\log(\alpha_{\theta}) = \min \left( \frac{1}{z(\sigma_{\theta}^2)^{(j)}} (\|\boldsymbol{\theta}^{(j)}\|^2 - \|\boldsymbol{\theta}^*\|^2), 0 \right).$$

- ▶ Different pseudo-likelihoods  $\implies$  different  $p(\sigma_{\theta}^2 | \mathcal{D}) \implies$  different  $\nu$  inferred to MAP estimates.

Introduction

Metropolis-Hastings

The Objective

The Likelihood

The Navigation Problem

XO

References

Figure:  $\|\boldsymbol{\theta}^{(j)}\|^2$  for  $j = 1, \dots, 100,000$  (no burn-in) in the left panel, with distribution of marginal  $\sigma_{\theta}^2 | \mathcal{D}$  (using 20,000 burn-in for binomial-based and 60,000 burn-in for beta- and exponential-based) for the three likelihood types in the right panel (all plots use the same scale).



# The Tic-Tac-Toe Problem

- ▶ Classical tic-tac-toe game is a two-player, deterministic, turn-based game in which the player and opponent alternately place their respective tokens  $X$  and  $O$  on a  $3 \times 3$  grid.
- ▶ The objective of the player (who plays first) is to be the first to align three of their tokens consecutively in a row, column, or diagonal, and likewise for the opponent.
- ▶ Agent trained to play against opponent whose decisions are random.

Introduction

Metropolis-Hastings

The Objective

The Likelihood

The Navigation Problem

XO

References

# XO: Feature engineering

Introduction

Metropolis-Hastings

The Objective

The Likelihood

The Navigation Problem

XO

References

- ▶ Represent  $3 \times 3$  board as matrix  $\mathbf{M}_{3 \times 3}$ , vectorized row-wise into  $\mathbf{m}$ .
- ▶ Terminal condition: compute  $\mathbf{m}' \mathbf{S}_{9 \times 8}$ , where each column of  $\mathbf{S}$  encodes a row, column, or diagonal.
  - ▶ If any entry =  $+3 \implies X$  wins.
  - ▶ If any entry =  $-3 \implies O$  wins.
  - ▶ Draws are identified from  $f(\mathbf{m}' \mathbf{S})$  when no moves remain.
- ▶ Network inputs at time  $\tau_k$ : (i) current board  $\mathbf{m}^{\tau_k-1}$ , (ii) evaluation state  $(\mathbf{m}^{\tau_k-1})' \mathbf{S}$ .
- ▶ Network output: action  $a^*$  from available board positions:  $\mathcal{A}_{\tau_k} = \{i \in \{1, \dots, 9\} : m_i^{\tau_k-1} = 0\}$ .
- ▶ Output activation:

$$\sigma_L(a \mid \mathbf{m}^{\tau_k-1}, \boldsymbol{\theta}) = \frac{\exp(\ell_a)}{\sum_{a' \in \mathcal{A}_{\tau_k}} \exp(\ell_{a'})},$$

with action  $a^*$  corresponding to action with the highest probability,  $a^* = \operatorname{argmax}_{a \in \mathcal{A}_{\tau_k}} \sigma_L(a \mid \cdot)$

- ▶ Objective: maximize win frequency of player  $O$  (encoded  $+1$ ) over  $K$  games:

$$\hat{\boldsymbol{\theta}} = \operatorname{argmax}_{\boldsymbol{\theta}} \frac{1}{K} \sum_{k=1}^K \mathbb{I}(\rho_{\tau_k}(\boldsymbol{\theta}) = +1).$$



# XO: The effects of regularization

- ▶ **Training setup:** Simulate 100 games against a random-policy opponent, using seed values  $\{\omega_i^{\text{Train}}\}_{i=1}^{100}$ .
- ▶ Learned solution  $\hat{\theta}$  influences opponent trajectories  $\implies$  each training run corresponds to a distinct environment.
- ▶  $\implies$  The optimization surface varies across training iterations.
- ▶ **Test setup:** Evaluate on 10,000 games with disjoint seeds  $\{\omega_i^{\text{Test}}\}$ , ensuring  $\omega_i^{\text{Test}} \notin \{\omega_j^{\text{Train}}\}$ .
- ▶ Question: Are all test opponent trajectories truly unseen compared to training?
- ▶ **Result:** Performance decreases (in- and out-of-sample) as regularization strength  $\nu$  increases.
  - ▶ Excessive regularization  $\implies$  underfitting.
  - ▶ Best out-of-sample % for  $\nu \in [10^{-6}, 10^{-4}]$  (beats random vs random hit-rate  $\approx 57\%$ ).

$\nu$	In-Sample (%)	Out-of-Sample (%)
$10^{-6}$	99	83.92
$10^{-5}$	99	74.84
$10^{-4}$	98	66.96
$10^{-3}$	93	59.96
$10^{-2}$	85	59.36
$10^{-1}$	74	69.10
1	52	46.04

**Table:** Normalized win percentage of  $O$  tokens across training ( $K = 100$ ) and test ( $K = 10,000$ ) games, for varying regularization strengths  $\nu$ .

Introduction  
Metropolis-Hastings  
The Objective  
The Likelihood  
The Navigation Problem  
XO

References

- ▶ **Recall:** pseudo-likelihood tempering  $[p(\boldsymbol{\theta} \mid \mathcal{D})]^\beta \implies$  sharpens sensitivity to  $\text{Obj}(\boldsymbol{\theta})$ .
- ▶  $\beta$  controls the trade-off between exploration and exploitation.
- ▶ As  $\beta$  increases, in-sample  $O$  wins improve:
  - ▶  $\beta$  amplifies the likelihood ratio:  $\left( \frac{p(\mathcal{D} \mid \boldsymbol{\theta}^*)}{p(\mathcal{D} \mid \boldsymbol{\theta}^{(j)})} \right)^\beta$ .
  - ▶ Larger  $\beta \implies$  chain favors higher-Obj proposals.
- ▶ Empirically:  $\beta \geq 100 \implies$  convergence to same dominant mode.
- ▶ Caution: overly large  $\beta \implies$  Markov chain stagnation.

Sharpness $\beta$	In-Sample (%)	Out-of-Sample (%)
0.1	56	50.08
1	70	59.86
10	63	56.68
50	86	72.98
100	98	72.62
1000	96	68.58

**Table:** Normalized  $O$  win rate across training ( $K = 100$ ) and test ( $K = 10,000$ ) games for varying sharpness  $\beta$ .

- ▶ For low  $\beta$  (flat likelihoods)  $\implies$  indefinite contraction of  $\|\boldsymbol{\theta}^{(j)}\|$ .
- ▶ Different  $\beta$ 's  $\implies$  different  $p(\sigma_\theta^2 \mid \mathcal{D})$   $\implies$  different  $\nu$  inferred to MAP estimates.

Introduction

Metropolis-Hastings

The Objective

The Likelihood

The Navigation Problem

XO

References

Figure:  $\|\boldsymbol{\theta}^{(j)}\|^2$  for  $j = 1, \dots, 80,000$  (post burn-in) with distribution of marginal  $\sigma_\theta^2 \mid \mathcal{D}$  for varying likelihood sharpness  $\beta$ .

- ▶ Markov chain initialized with variance  $\sigma_{\text{Init}}^2 \rightarrow \boldsymbol{\theta}^{(1)} \sim \mathcal{N}(\mathbf{o}_{S \times 1}, \sigma_{\text{Init}}^2 \mathbf{I}_S)$ .
- ▶ Marginals  $p(\sigma_\theta^2 | \mathcal{D})$  follow an inverse-gamma distribution  $\rightarrow$  shape parameters remain approximately constant, rate parameters grow with  $\sigma_{\text{Init}}^2 \implies$  both the mean and variance shift upward.

$\sigma_{\text{Init}}^2$	In-Sample (%)	Out-of-Sample (%)	Shape	Rate
0.1	92	71.18	41.75	3.14
1	88	65.40	46.23	41.31
10	98	72.62	44.93	413.30
100	92	68.92	47.28	5704.89

**Table:** The normalized number of  $O$  wins, as a percentage for in-sample ( $K = 100$ ) and out-of-sample ( $K = 10,000$ ) sets across various initial variances  $\sigma_{\text{Init}}^2$  accompanied by shape and rate parameters of marginal  $\sigma_\theta^2 | \mathcal{D} \sim \text{Inv-Gamma}$  for likelihood sharpness  $\beta = 100$ .

- ▶ **Motivation:** Introduced uncertainty about the dispersion parameter  $\sigma_\theta^2 \propto \frac{1}{\nu} \rightarrow$  training set informs  $\nu$ .
- ▶ **Hyperparameters:** pseudo-likelihood form, likelihood sharpness  $\beta$ , and initial variance  $\sigma_{\text{Init}}^2$  are user-specified  $\rightarrow$  strongly influences marginal  $p(\sigma_\theta^2 \mid \mathcal{D})$ .
- ▶ **Limitation:** Extent of regularization is not fully data-driven.
- ▶ **Conclusion:** Using a hierarchical Bayesian model solely to infer  $\sigma_\theta^2$  may be unnecessary  $\rightarrow$  instead fix  $\sigma_\theta^2$  (and thus  $\nu$ ) directly.
  - ▶ In this sense, the two-block MCMC procedure acts as the *user* inferring a specific regularization strength but with "extra steps".

Introduction

Metropolis-  
Hastings

The Objective

The  
Likelihood

The  
Navigation  
Problem

XO

References



# References I

Introduction

Metropolis-  
Hastings

The Objective

The  
Likelihood

The  
Navigation  
Problem

XO

References



Roberts, Gareth O et al. (1997). "Weak convergence and optimal scaling of random walk Metropolis algorithms". In: *The Annals of Applied Probability* 7.1, pp. 110–120.

