# MALA and Barker's Proposal: Gradient-Based MCMC Methods

# From RWM to more advanced methods

**Random Walk Metropolis (RWM):**

$$q^* = q + \sigma W, \quad W \sim N(0, I_d)$$

**Fundamental Trade-off:**
- ▶ Large step-size $\sigma$: Low acceptance
- ▶ Small $\sigma$: Slow exploration
- ▶ Optimal: $\sigma = \mathcal{O}(d^{-1})$

**Problem:** In high dimensions, RWM becomes inefficient
- ▶ Optimal acceptance rate: 0.234
- ▶ Curse of dimensionality: step size $\propto 1/d$

# From Langevin Diffusion to MALA

Use gradient to **move toward modes of** $\pi$
**Continuous Langevin Diffusion:**

$$dX_t = \frac{1}{2}\nabla \log \pi(X_t)dt + dB_t$$

▶ Has $\pi$ as stationary distribution
▶ Gradient provides moves toward high-probability regions

**Unadjusted Langevin Algorithm (ULA):**

$$X^{(t)} = X^{(t-1)} + \frac{\sigma^2}{2}\nabla \log \pi(X^{(t-1)}) + \sigma W$$

**Problem** $\pi$ is **not** the invariant distribution of ULA!
**Solution:** Add Metropolis-Hastings correction $\Rightarrow$ MALA

# Metropolis-Adjusted Langevin Algorithm

just for reference. do not write down algorithm during exam

---

**Algorithm 1** MALA

---

    **Input:** Initial $X^{(0)}$
    **for** $t = 1, 2, \ldots$ **do**
        Propose: $X^* = X^{(t-1)} + \frac{\epsilon}{2} \nabla \log \pi(X^{(t-1)}) + \sqrt{\epsilon} W$
        Compute acceptance ratio:

$$\alpha = \min \left\{ 1, \frac{\pi(X^*) q(X^{(t-1)}|X^*)}{\pi(X^{(t-1)}) q(X^*|X^{(t-1)})} \right\}$$

        Accept $X^{(t)} = X^*$ with probability $\alpha$, else $X^{(t)} = X^{(t-1)}$
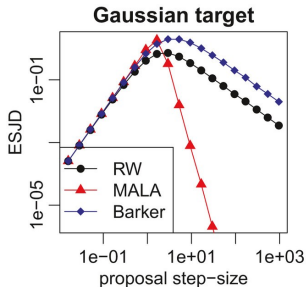    **end for**

---

# Optimal Scaling Theory

**Maximizing Expected Squared Jump Distance (ESJD)** $\mathbb{E}\left[\|X^{(t+1)} - X^{(t)}\|^2\right]$

**Dimension Scaling:**
- ▶ RWM: $\sigma = \mathcal{O}(d^{-1})$
- ▶ MALA: $\sigma = \mathcal{O}(d^{-1/3})$

**Optimal Acceptance:**
- ▶ RWM: 0.234
- ▶ MALA: 0.574



Gaussian target

**Implication:** MALA maintains larger step sizes in high dimensions
- ▶ Better exploration efficiency
- ▶ Faster convergence to target distribution
- ▶ MALA hard to tune
- ▶ **Catch** - requires gradient computation

# Building Informed Proposals

RWM does not use target information to guide proposals.
**Core idea:** Use target density $\pi$ to guide proposals

**General Framework:** $Q_g(x, y) = \frac{g\left(\frac{\pi(y)}{\pi(x)}\right)K(x,y)}{Z_g(x)}$

where $g$ is a "balancing function".

**Remarkable result:**, $Q_g$ is **locally balanced** if $g(t) = tg(1/t)$.

Two Special Cases:

1. MALA: $g(t) = \sqrt{t}$
2. Barkers Proposal: $g(t) = \frac{t}{1+t}$

Both cases satisfy this property. This framework unifies MALA and Barker as different choices of the balancing function $g$.

# Barker's Proposal

Consider the family of informed proposals:

$$Q_g(x, dy) = \frac{g\left(e^{(\nabla \log \pi(x))^T(y-x)}\right) K(x, dy)}{Z_g(x)}$$

and use $g(t) = t/(1+t)$. This gives $Z_g(x) = 1/2$ and **Proposal Density:**

$$Q_B(x, dy) = \frac{2}{1 + e^{-(\nabla \log \pi(x))^T(y-x)}} K(x, dy)$$

**Key Idea:** Use gradient to stochastically bias proposal direction. Barker uses gradient to flip coin for direction, MALA uses it as deterministic drift.

---

**Algorithm 2** 1D case with Gaussian kernel

    Sample $Z \sim N(0, \sigma^2)$

    Propose $Y = x + Z$ with probability $1/(1 + \exp(-Z^T \nabla \log \pi(x)))$

    Propose $Y = x - Z$ with residual probability

---

# Just for my own reference

This help to understand why Barker proposal use gradient to stochastically bias proposal direction.

**Weight of proposal:**

$$w = \frac{1}{1 + e^{-(\nabla \log \pi(x))^T (y-x)}} = \frac{1}{1 + e^{-g(y-x)}}$$

where $g = \nabla \log \pi(x)$.

**Behavior analysis:** Consider four scenarios based on the signs of $g$ and $(y - x)$:

| Scenario | $(y - x)$ | $g(y-x)$ | $e^{-g(y-x)}$ | Weight | Meaning |
|---|---|---|---|---|---|
| $g > 0$, move right | $> 0$ | $> 0$ | $\approx 0$ | $\approx 1$ | Favored |
| $g > 0$, move left | $< 0$ | $< 0$ | $\to \infty$ | $\approx 0$ | Penalized |
| $g < 0$, move left | $< 0$ | $> 0$ | $\approx 0$ | $\approx 1$ | Favored |
| $g < 0$, move right | $> 0$ | $< 0$ | $\to \infty$ | $\approx 0$ | Penalized |

Table: Barker proposal weight behavior

# MALA vs Barker's Proposal

| Aspect | MALA | Barker |
|---|---|---|
| Proposal | $Y = x + \frac{\sigma^2}{2}\nabla \log \pi(x) + \sigma Z$ | $Y = x \pm Z$ with directional prob |
| Gradient use | Drift term (deterministic shift) | Direction selection (probabilistic) |
| Robustness | Sensitive to step size | More robust to large gradients |
| Scaling | $\mathcal{O}(d^{-1/3})$ | $\mathcal{O}(d^{-1/3})$ |