

Optimizing Solution-Samplers for Combinatorial Problems: The Landscape of Policy-Gradient Methods

NeurIPS'23

Constantine Caramanis, Dimitris Fotakis, Alkis Kalavasis, Vasilis
Kontonis, Christos Tzamos

[Stanford CS/MS&E 331](#)

Gradient descent for combinatorial opt

- **Common pipeline** (e.g., several papers from this quarter):
 1. Train NN whose parameters define a **distribution** over solutions
 2. Optimize expected cost via **gradients**
- Works well, but **lacks theory** explaining *why*
- **Challenges:**
 - Naïve **exponential-size** simplexes guarantee convexity
 - But compact parameterizations often have **bad landscapes**
- Q: Design polynomial-size, well-behaved distribution families where gradient descent **provably** finds near-optimal solutions?
- This paper: **yes**, with applications to several CO problems

Setup

- Instances $I \in \mathcal{I}$ with common solution space S
 - Prior distribution \mathcal{R} over train/test inputs
- Cost function $L(\cdot ; I): S \rightarrow \mathbb{R}$ (assume efficient to evaluate)
 - E.g., for max cut: $L(s ; I) = -(\text{cut weight})$
 - Results hold even for **blackbox/oracle access** to L
- Distribution $p(\cdot ; I, \mathbf{w})$ over solutions with trainable $\mathbf{w} \in \mathcal{W}$
- Loss function $\mathcal{L}(\mathbf{w}) = \mathbb{E}_{I \sim \mathcal{R}, s \sim p(\cdot ; I, \mathbf{w})} [L(s ; I)]$
- **Goal:** find \mathbf{w} so that $\mathcal{L}(\mathbf{w})$ is close to $\text{opt} = \mathbb{E}_{I \sim \mathcal{R}} \left[\min_{s \in S} L(s ; I) \right]$

Gradient descent dynamics

- **Algorithm is fixed:** gradient descent
 - Many other algorithms you could imagine (see, e.g., Remark 7)
 - Choose gradient descent because it's very common

- **At each iteration (t):**

1. Sample an instance $I \sim R$ (or a minibatch)
2. Sample solutions $s_1, \dots, s_m \sim p(\cdot; I, \mathbf{w}_t)$
3. Query costs $c_k = L(I, s_k)$
4. Form a policy-gradient estimate \mathbf{g}_t and update $\mathbf{w}_{t+1} \leftarrow \mathbf{w}_t + \eta_t \mathbf{g}_t$

This paper: Form of the distribution

$$\mathbf{g}_t \approx \frac{1}{m} \sum_{k=1}^m c_k \nabla_{\mathbf{w}} \log p(s_k; I, \mathbf{w}_t)$$

This paper: How to regularize loss

Key desiderata

- Distribution over all (s, I) needs exponentially many params
- Neural parameterizations are **compressed representations**
 - $[I]$ = # bits to represent instance
- **Three desiderata:**
 - **Complete:** some \bar{w} achieves $\mathcal{L}(\bar{w}) \leq \text{opt} + \epsilon$
 - **Compressed:** Description size of \mathcal{W} is $[\mathcal{W}] = \text{poly}\left([I], \frac{1}{\epsilon}\right)$
 - **Efficiency:** first-order methods reach 2ϵ -opt loss in $\text{poly}\left([\mathcal{W}], \frac{1}{\epsilon}\right)$ steps
- Does **not** imply P=NP: sampling may still be hard
 - In practice, **approximate samplers** often work well

Key assumptions

Feature mappings $\psi_S(s) \in \mathbb{R}^{n_X}, \psi_J(I) \in \mathbb{R}^{n_Z}$

1. *Boundedness*: $\|\psi_S(s)\|_2 \leq D_S, \|\psi_J(I)\|_2 \leq D_J$ for all s, I
2. *Bilinear cost structure* $L(s, I) = \psi_S(s)^\top M \psi_J(I)$ with $\|M\|_F \leq C$
 - Example in future slides
 - M is unknown (results hold even for blackbox/oracle access to L)
3. *Variance preserves features* when $s \sim \text{Unif}(S)$
 - For all $\mathbf{v} \in \mathbb{R}^{n_X}, \text{Var}_{s \sim \text{Unif}(S)}[\mathbf{v} \cdot \psi_S(s)] \geq \alpha \|\mathbf{v}\|_2$
 - Ensures gradients of optimization objective are not vanishing
4. *Polynomial scaling* $n_X, n_Z, D_S, D_J, C \leq \text{poly}([J]), \alpha \geq 1/\text{poly}([J])$

Example: Max-cut

Diagonal degree matrix:

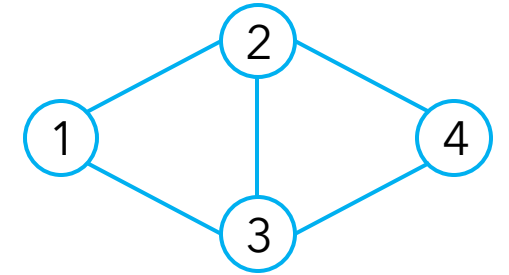
$$D = \begin{pmatrix} 2 & 0 & 0 & 0 \\ 0 & 3 & 0 & 0 \\ 0 & 0 & 3 & 0 \\ 0 & 0 & 0 & 2 \end{pmatrix}$$

Adjacency matrix:

$$A = \begin{pmatrix} 0 & 1 & 1 & 0 \\ 1 & 0 & 1 & 1 \\ 1 & 1 & 0 & 1 \\ 0 & 1 & 1 & 0 \end{pmatrix}$$

Laplacian:

$$L_G = D - A = \begin{pmatrix} 2 & -1 & -1 & 0 \\ -1 & 3 & -1 & -1 \\ -1 & -1 & 3 & -1 \\ 0 & -1 & -1 & 2 \end{pmatrix}$$



Example: Max-cut

- For n -node graph, max-cut equivalent to: $\max_{\mathbf{s} \in \{\pm 1\}^n} \frac{1}{4} \mathbf{s}^\top L_G \mathbf{s}$
- For any $\mathbf{s} \in \{\pm 1\}^n$, $\psi_S(\mathbf{s}) = (\mathbf{s}\mathbf{s}^\top)^\flat \in \mathbb{R}^{n^2}$

Flatten to \mathbb{R}^{n^2}
- $\psi_J(G) = L_G^\flat \in \mathbb{R}^{n^2}$
 1. *Boundedness*: $\|\psi_S(\mathbf{s})\|_2, \|\psi_J(G)\|_2 \leq \text{poly}([J]) = \text{poly}(n)$
 2. *Bilinear cost*: Define $M = \frac{-1}{4} \cdot I_{n^2}$
$$L(s, I) = \psi_J(G)^\top M \psi_S(\mathbf{s}) = -\frac{1}{4} (L_G^\flat)^\top (\mathbf{s}\mathbf{s}^\top)^\flat$$
 3. *Variance preservation*: $\text{Var}_{\mathbf{s} \sim \text{Unif}(\{\pm 1\}^n)} [\mathbf{v} \cdot \mathbf{s}] = \|\mathbf{v}\|_2^2$

Solution sampler: Obstacle 1

- **Goal in mind:** sample $\propto \exp(-\tau \cdot \psi_S(s)^\top M \psi_J(I))$
 - τ : temperature; as $\tau \rightarrow \infty$, samples solutions with small loss
 - But M is **unknown** (results hold even for blackbox/oracle access to L)
- **Natural candidate:** $p(s ; I, W) \propto \exp(\psi_S(s)^\top W \psi_J(I))$
- In GD, is W_t moving **in the direction** of $\bar{W} = -\tau M$ with $\tau \rightarrow \infty$?
 - Yes: $\nabla_W \mathcal{L}(W) \cdot M \geq 0$, so moving in the direction of $-M$ decreases loss
- But **no finite optimizer**; infimum only reached as $\|W\| \rightarrow \infty$
- Conflicts with “efficiency” desideratum:
“First-order methods reach 2ϵ -opt loss in $\text{poly}\left([\mathcal{W}], \frac{1}{\epsilon}\right)$ steps”

Solution 1: Add entropy regularization

- Natural candidate: $p(s ; I, W) \propto \exp(\psi_S(s)^\top W \psi_J(I))$
- Make loss landscape **more benign** by adding regularizer
 - **Goal:** make the landscape “quasar-convex”
- f is γ -quasar-convex wrt minimizer x^* on domain D if, $\forall x \in D$
$$\nabla f(x) \cdot (x - x^*) \geq \gamma(f(x) - f(x^*)), \quad \gamma \in [0,1]$$
- Gradient always points somewhat **toward** x^*
- Role of γ :
 - Measures how strongly the gradient “leans” toward the optimum
 - **Larger $\gamma \Rightarrow$ faster progress** for gradient descent

Solution 1: Add entropy regularization

- Natural candidate: $p(s ; I, W) \propto \exp(\psi_S(s)^\top W \psi_J(I))$
- Make loss landscape **more benign** by adding regularizer
- Negative entropy: $H(W) = \mathbb{E}_{I \sim \mathcal{R}, s \sim p(\cdot ; I, W)} [\log p(s ; I, W)]$
- **Regularized objective:** $\mathcal{L}_\lambda(W) = \mathcal{L}(W) + \lambda H(W)$
- Paper shows:
 - $\bar{W} = -\frac{M}{\lambda}$ is a minimizer of \mathcal{L}_λ
 - Sampling from $p(s ; I, \bar{W})$ yields **2 ϵ -opt loss** for $\lambda = \text{poly}\left(\epsilon, \frac{1}{D_S}, \frac{1}{D_J}, \dots\right)$
- So **GD will eventually converge** to a 2ϵ -opt loss
 - But quasr-convexity parameter γ may be very small \Rightarrow **slow rates**

Obstacle 2: Vanishing gradients

For quasar-convexity, we want

$$\nabla \mathcal{L}_\lambda(W) \cdot (W - \bar{W}) \geq \gamma (\mathcal{L}_\lambda(W) - \mathcal{L}_\lambda(\bar{W}))$$

Paper proves:

$$\nabla \mathcal{L}_\lambda(W) \cdot (W - \bar{W}) = \text{Var}[\psi_S(s)^\top (W - \bar{W}) \psi_J(I)] = \text{Var}(Y)$$

Problem: variance term can be **tiny** near \bar{W} , so γ may be small

- Sampler is a softmax over scores: $p(s ; I, W) \propto \exp(\psi_S(s)^\top W \psi_J(I))$
- As $\|W\|$ grows, softmax **concentrates on argmax** of $\psi_S(s)^\top W \psi_J(I)$
- So RV Y becomes **almost deterministic**, so $\text{Var}(Y) \rightarrow 0$

Solution 2: Fast/slow mixture generators

New generator family: mix two exponential-family samplers

$$\mathcal{P} = \{(1 - \beta^*)p(\cdot; I, W) + \beta^*p(\cdot; I, \rho^*W)\}$$

Fast component: $p(\cdot; I, W)$

- Drives convergence toward minimizer $(\bar{W} = -\frac{M}{\lambda})$

Slow component: $p(\cdot; I, \rho^*W)$ with small ρ^*

- Stays close to the **uniform** distribution over solutions
- Key use of assumption: variance preservation under uniform dist
 - Guarantees a lower bound on variance

Gradient descent dynamics

At each iteration (t):

1. Sample an instance $I \sim R$ (or a minibatch)
2. Sample solutions $s_1, \dots, s_m \sim p(\cdot; I, \mathbf{w}_t)$
3. Query costs $c_k = L(I, s_k)$
4. Form a policy-gradient estimate and update $\mathbf{w}_{t+1} \leftarrow \mathbf{w}_t + \eta_t \mathbf{g}_t$

This paper: Exponential mixture

$$\mathbf{g}_t \approx \frac{1}{m} \sum_{k=1}^m c_k \nabla_{\mathbf{w}} \log p(s_k; I, \mathbf{w}_t)$$

This paper: Replace with entropy-regularized loss

- ✓ **Complete:** some $\bar{\mathbf{w}}$ achieves $\mathcal{L}(\bar{\mathbf{w}}) \leq \text{opt} + \epsilon$
- ✓ **Compressed:** Description size of \mathcal{W} is $[\mathcal{W}] = \text{poly}([I], 1/\epsilon)$
- ✓ **Efficiency:** GD reaches 2ϵ -opt loss in $\text{poly}([\mathcal{W}], 1/\epsilon)$ steps

Overview

- (Under assumptions) exist solution samplers that are:
 - **Complete** (near-optimal solutions exist)
 - **Compressed** (poly-sized parameterization)
 - **Efficiently optimizable** (GD finds ϵ -opt in poly steps)
- Two key landscape **obstacles**:
 - Minimizers at infinity
 - Vanishing gradients near good solutions
- Two corresponding **fixes**:
 - Entropy-regularized loss \rightarrow finite, quasar-convex minimizer
 - Fast/slow exponential mixture \rightarrow non-vanishing gradients
- Results apply to several canonical CO problems