

Differential Dynamic Programming for Structured Prediction and Attention

A short report on [Mensch and Blondel, 2018]

Nazarov Ivan

Skoltech

August 3, 2018

Introduction

Mina idea

- ▶ cast a dynamic programming problem as a linear problem
- ▶ somehow make the linear problem smooth by relaxation
- ▶ recast the smoothed problem back into dynamic programming setting

This enables

- ▶ take gradients w.r.t. parameters
- ▶ compute hessian-vector products w.r.t parameters

Dynamic Programming Problem

Typical setting for a dynamic programming problem

- ▶ $G = (V, E)$ DAG with unique source v_* and sink v^* nodes
- ▶ weights $\theta \in \mathbb{R}^{V \times V}$ with $\theta_{v_* v_*} = 1$ and $\theta_{uv} = -\infty$ if $uv \notin E$

Goal

Get a path with the highest score among all paths $v_* \rightarrow v^*$

Solution

Identify each $v \in V$ with its number $i = 1 \dots m$ in topological order

$$F_j(\theta) \leftarrow \max_{i: j \in G_i} \theta_{ij} + F_j(\theta), \quad F_1(\theta) \leftarrow 0,$$

$$DP(\theta) \leftarrow F_m(\theta).$$

Find the optimal path by *backtracking* through $(F_j(\theta))_{j=1}^m$

Dynamic Programming Problems as Linear Problems

[Bellman, 1952] showed that $DP(\theta) = LP(\theta)$

$$LP(\theta) = \max_{\pi} \sum_{u \in V} \sum_{v \in V} \theta_{uv} 1_{uv \in \pi} = \max_{\mathbf{y} \in \mathcal{Y}} \langle \theta, \mathbf{y} \rangle$$

where \mathcal{Y} is the set of binary matrices representing paths $v_* \rightarrow v^*$

However ...

- ▶ $LP(\theta)$ is not differentiable unless its solution is unique
- ▶ the optimal solution $\mathbf{y}^*(\theta) = \arg LP(\theta)$ is a discontinuous map

What is a “maximum”?

The largest element of $\theta \in \mathbb{R}^d$ is $\max(\theta)$

$$\begin{aligned} \max: \quad \mathbb{R}^d &\longrightarrow \mathbb{R}, \\ \theta &\longmapsto \max_{i=1}^d \theta_i = \sup_{x \in \Delta^d} \langle x, \theta \rangle, \end{aligned} \quad (\text{Max})$$

where Δ^d is the *unit simplex* in \mathbb{R}^d , $\Delta^d = \{x: \|x\|_1 = 1, x \geq 0\}$.

- ▶ used in every optimization problem
- ▶ differentiable almost everywhere (except on negligible sets)
- ▶ non-differentiable solution

Making maxima smooth

Let $\Omega: \mathbb{R}^d \rightarrow \mathbb{R}$ be a strongly convex regularizer on Δ^d

$$\begin{aligned} \max_{\Omega}: \quad \mathbb{R}^d &\longrightarrow \mathbb{R}, \\ \theta &\longmapsto \sup_{x \in \Delta^d} \langle x, \theta \rangle - \Omega(x), \end{aligned} \quad (\text{Smooth-Max})$$

Properties from strong convexity of Ω

- ▶ $x^*(\theta) = \arg \sup_{x \in \Delta^d} \langle x, \theta \rangle - \Omega(x)$ exists and unique
- ▶ $\nabla \max_{\Omega}(\theta) = x^*(\theta)$ and is Lipschitz-continuous
- ▶ $\nabla^2 \max_{\Omega}(\theta)$ exists almost everywhere

Natural generalization of max

Table 1: Various types of \max_{Ω}

	<i>regular-max</i>	<i>soft-max</i>	<i>sparse soft-max</i>
Ω	0	$-\sum_i x_i \log x_i$	$\frac{1}{2} \ \cdot\ ^2$

- ▶ $\max_{\Omega}(\theta_1) \leq \max_{\Omega}(\theta_2)$ whenever $\theta_1 \leq \theta_2$
- ▶ $\max_{\Omega}(\mathbf{1}c + \theta) = c + \max_{\Omega}(\theta)$ for any $c \in \mathbb{R}$
- ▶ $\max_{\Omega}(\pi\theta) = \max_{\Omega}(\theta)$ for any permutation P with $\Omega \circ P = \Omega$
- ▶ if $\theta_j = -\infty$ then $(\nabla \max_{\Omega}(\theta))_j = 0$
- ▶ $\max_{\Omega}(\theta)$ is not far from $\max(\theta)$

Smoothed LP and DP

For any $f: \mathcal{Y} \rightarrow \mathbb{R}$ denote

$$\max_{\mathbf{y} \in \mathcal{Y}} f(\mathbf{y}) \triangleq \max_{\Omega} (f(\mathcal{Y})), \quad f(\mathcal{Y}) = (f(\mathbf{y}))_{\mathbf{y} \in \mathcal{Y}} \in \mathbb{R}^{|\mathcal{Y}|}.$$

Linear Program

$$\text{LP}_{\Omega}(\theta) = \max_{\mathbf{y} \in \mathcal{Y}} \langle \theta, \mathbf{y} \rangle.$$

Bellman Iterations

In topological order of G do

$$F_j(\theta) \leftarrow \max_{i: j \in G_i} (\theta_{ij} + F_i(\theta)), \quad F_1(\theta) \leftarrow 0,$$

$$\text{DP}_{\Omega}(\theta) \leftarrow F_m(\theta).$$

Behaviour of smoothed LP and DP

Both $LP_{\Omega}(\theta)$ and $DP_{\Omega}(\theta)$ are well-behaved

- ▶ convex and differentiable everywhere
- ▶ have Lipschitz continuous gradients
- ▶ gradients are differentiable almost everywhere

However ...

- ▶ $LP_{\Omega}(\theta)$ is intractable due to exponential size of \mathcal{Y}
- ▶ $DP_{\Omega}(\theta)$ is tractable and has complexity $\mathcal{O}(|E|)$

Properties of smoothed LP and DP

[Mensch and Blondel, 2018] propose and prove:

- ▶ $DP_{\Omega}(\theta)$ is convex w.r.t θ
- ▶ $DP_{\Omega}(\theta) = LP_{\Omega}(\theta)$ **if and only if** $\Omega(x) = -\gamma \sum_i x_i \log x_i$
- ▶ is “close” to LP: $|LP(\theta) - DP_{\Omega}(\theta)| \leq m M(\Omega, m)$ and

$$\lim_{\gamma \rightarrow 0} DP_{\gamma\Omega}(\theta) = LP(\theta),$$

- ▶ efficient recursion to compute $\nabla_{\theta} DP_{\Omega}(\theta) \in \mathbb{R}^{V \times V}$ and hessian-vector products $\nabla^2 DP_{\Omega}(\theta)Z$

Computing $\nabla_{\theta} \text{DP}_{\Omega}(\theta)$

Simply backpropagate along the reverse-topological order of G

- ▶ **forward-pass:** while computing $F_i(\theta)$ get

$$q_i(\theta) = \nabla \max_{\Omega} (\theta_i + F(\theta)) \in \Delta^m,$$

assuming $F_k(\theta) = -\infty$ for all k **after** i

- ▶ **backward-pass:** in reverse-topological order $j = m \dots 1$ do

$$\bar{w}_j \leftarrow \sum_{i \in G_j} w_{ij} \text{ if } j \neq m \text{ else } 1,$$

$$w_{ij} \leftarrow \bar{w}_i q_{ij} \text{ if } i \in G_j \text{ else } 0,$$

$$\nabla_{\theta} \text{DP}_{\Omega}(\theta) \leftarrow (w_{ij})_{i,j=1 \dots m} \in \mathbb{R}^{V \times V}.$$

Interpretation of $\nabla_{\theta} \text{DP}_{\Omega}(\theta)$

The matrix $Q(\theta) = (q_i(\theta))_{i=1}^m$

- ▶ a transition matrix for backward random walks from $v^* = m$ back to $v_* = 1$
- ▶ $\mathbb{P}(i \rightarrow j) = q_{ij}(\theta)$ if $i \in G_j$.

[Mensch and Blondel, 2018] demonstrate

- ▶ the gradient is the expected path of the random walk

$$\nabla_{\theta} \text{DP}_{\Omega} = \mathbb{E}_{\mathbf{y} \sim Q(\theta)} \mathbf{y}.$$

- ▶ convergence to the optimal solution

$$\nabla_{\theta} \text{DP}_{\gamma\Omega}(\theta) \xrightarrow{\gamma \rightarrow 0} \mathbf{y}^*(\theta) \in \partial \text{LP}(\theta).$$

Conclusion

[Mensch and Blondel, 2018] propose

- ▶ a theoretical framework for turning dynamic programs into convex, differentiable and tractable operators
- ▶ efficient way to embed the smoothed programs into models learnt by gradient descent

Applications: learning optimal cost parameters for end-to-end training in

- ▶ sequence prediction in part-of-speech tagging
- ▶ time series alignment in audio transcription
- ▶ attention mechanism in machine translation

References



Bellman, R. (1952).

On the theory of dynamic programming.

Proceedings of the National Academy of Sciences, 38(8):716–719.



Mensch, A. and Blondel, M. (2018).

Differentiable Dynamic Programming for Structured Prediction and Attention.

In *35th International Conference on Machine Learning*, volume 80 of *Proceedings of the 35th International Conference on Machine Learning*, Stockholm, Sweden.