

Supplementary material

A. Implementation Details for SparseMAP Solvers

A.1. Conditional Gradient Variants

We adapt the presentation of vanilla, away-step and pairwise conditional gradient of [Lacoste-Julien & Jaggi \(2015\)](#).

Recall the SparseMAP optimization problem (Equation 5), which we rewrite below as a minimization, to align with the formulation in ([Lacoste-Julien & Jaggi, 2015](#))

$$\text{SparseMAP}_A(\eta) := \arg \min_{\mathbf{u}: [\mathbf{u}, \mathbf{v}] \in \mathcal{M}_A} f(\mathbf{u}, \mathbf{v}), \quad \text{where } f(\mathbf{u}, \mathbf{v}) := \frac{1}{2} \|\mathbf{u}\|_2^2 - \eta_U^\top \mathbf{u} - \eta_F^\top \mathbf{v}.$$

The gradients of the objective function f w.r.t. the two variables are

$$\nabla_{\mathbf{u}} f(\mathbf{u}', \mathbf{v}') = \mathbf{u}' - \eta_U, \quad \nabla_{\mathbf{v}} f(\mathbf{u}', \mathbf{v}') = -\eta_V.$$

The ingredients required to apply conditional gradient algorithms are solving linear minimization problem, selecting the away step, computing the Wolfe gap, and performing line search.

Linear minimization problem. For SparseMAP, this amounts to a MAP inference call, since

$$\begin{aligned} & \arg \min_{[\mathbf{u}, \mathbf{v}] \in \mathcal{M}_A} \langle \nabla_{\mathbf{u}} f(\mathbf{u}', \mathbf{v}'), \mathbf{u} \rangle + \langle \nabla_{\mathbf{v}} f(\mathbf{u}', \mathbf{v}'), \mathbf{v} \rangle \\ &= \arg \min_{[\mathbf{u}, \mathbf{v}] \in \mathcal{M}_A} (\mathbf{u}' - \eta_U)^\top \mathbf{u} - \eta_F^\top \mathbf{v} \\ &= \{[\mathbf{m}_s; \mathbf{n}_s] : s \in \text{MAP}_A(\eta_U - \mathbf{u}', \eta_F)\}. \end{aligned}$$

where we assume MAP_A yields the set of maximally-scoring structures.

Away step selection. This step involves searching the currently selected structures in the active set \mathcal{I} with the *opposite* goal: finding the structure *maximizing* the linearization

$$\begin{aligned} & \arg \max_{s \in \mathcal{I}} \langle \nabla_{\mathbf{u}} f(\mathbf{u}', \mathbf{v}'), \mathbf{m}_s \rangle + \langle \nabla_{\mathbf{v}} f(\mathbf{u}', \mathbf{v}'), \mathbf{n}_s \rangle \\ &= \arg \max_{s \in \mathcal{I}} (\mathbf{u}' - \eta_U)^\top \mathbf{m}_s - \eta_F^\top \mathbf{n}_s \end{aligned}$$

Wolfe gap. The gap at a point $\mathbf{d} = [\mathbf{d}_u; \mathbf{d}_v]$ is given by

$$\begin{aligned} \text{gap}(\mathbf{d}, \mathbf{u}') &:= \langle -\nabla_{\mathbf{u}} f(\mathbf{u}', \mathbf{v}'), \mathbf{d}_u \rangle + \langle -\nabla_{\mathbf{v}} f(\mathbf{u}', \mathbf{v}'), \mathbf{d}_v \rangle \\ &= \langle \eta_U - \mathbf{u}', \mathbf{d}_u \rangle + \langle \eta_F, \mathbf{d}_v \rangle. \end{aligned} \tag{7}$$

Line search. Once we have picked a direction $\mathbf{d} = [\mathbf{d}_u; \mathbf{d}_v]$, we can pick the optimal step size by solving a simple optimization problem. Let $\mathbf{u}_\gamma := \mathbf{u}' + \gamma \mathbf{d}_u$, and $\mathbf{v}_\gamma := \mathbf{v}' + \gamma \mathbf{d}_v$. We seek γ so as to optimize

$$\arg \min_{\gamma \in [0, \gamma_{\max}]} f(\mathbf{u}_\gamma, \mathbf{v}_\gamma)$$

Setting the gradient w.r.t. γ to 0 yields

$$\begin{aligned} 0 &= \frac{\partial}{\partial \gamma} f(\mathbf{u}_\gamma, \mathbf{v}_\gamma) \\ &= \langle \mathbf{d}_u, \nabla_{\mathbf{u}} f(\mathbf{u}_\gamma, \mathbf{v}_\gamma) \rangle + \langle \mathbf{d}_v, \nabla_{\mathbf{v}} f(\mathbf{u}_\gamma, \mathbf{v}_\gamma) \rangle \\ &= \langle \mathbf{d}_u, \mathbf{u}' + \gamma \mathbf{d}_u - \eta_U \rangle + \langle \mathbf{d}_v, -\eta_F \rangle \\ &= \gamma \|\mathbf{d}_u\|_2^2 + \mathbf{u}'^\top \mathbf{d}_u - \eta^\top \mathbf{d} \end{aligned}$$

We may therefore compute the optimal step size γ as

$$\gamma = \max \left(0, \min \left(\gamma_{\max}, \frac{\boldsymbol{\eta}^\top \mathbf{d} - \mathbf{u}'^\top \mathbf{d}_u}{\|\mathbf{d}_u\|_2^2} \right) \right) \quad (8)$$

Algorithm 1 Conditional gradient for SparseMAP

```

1: Initialization:  $s^{(0)} \leftarrow \text{MAP}_{\mathcal{A}}(\boldsymbol{\eta}_U, \boldsymbol{\eta}_F)$ ;  $\mathcal{I}^{(0)} = \{s^{(0)}\}$ ;  $\mathbf{y}^{(0)} = \mathbf{e}_{s^{(0)}}$ ;  $[\mathbf{u}^{(0)}; \mathbf{v}^{(0)}] = \mathbf{a}_{s^{(0)}}$ 
2: for  $t = 0 \dots t_{\max}$  do
3:    $s \leftarrow \text{MAP}_{\mathcal{A}}(\boldsymbol{\eta}_U - \mathbf{u}^{(t)}, \boldsymbol{\eta}_F)$ ;
4:    $w \leftarrow \arg \max_{w \in \mathcal{I}^{(t)}} (\boldsymbol{\eta}_U - \mathbf{u}^{(t)})^\top \mathbf{m}_w + \boldsymbol{\eta}_F^\top \mathbf{n}_w$ ;
5:    $\mathbf{d}^F \leftarrow \mathbf{a}_s - [\mathbf{u}^{(t)}; \mathbf{v}^{(t)}]$  (forward direction)
6:    $\mathbf{d}^W \leftarrow [\mathbf{u}^{(t)}; \mathbf{v}^{(t)}] - \mathbf{a}_w$  (away direction)
7:   if  $\text{gap}(\mathbf{d}^F, \mathbf{u}^{(t)}) < \epsilon$  then
8:     return  $\mathbf{u}^{(t)}$  (Equation 7)
9:   end if
10:  if variant = vanilla then
11:     $\mathbf{d} \leftarrow \mathbf{d}^F$ ;  $\gamma_{\max} \leftarrow 1$ 
12:  else if variant = pairwise then
13:     $\mathbf{d} \leftarrow \mathbf{d}^F + \mathbf{d}^W$ ;  $\gamma_{\max} \leftarrow y_w$ 
14:  else if variant = away-step then
15:    if  $\text{gap}(\mathbf{d}^F, \mathbf{u}^{(t)}) \geq \text{gap}(\mathbf{d}^W, \mathbf{u}^{(t)})$  then
16:       $\mathbf{d} \leftarrow \mathbf{d}^F$ ;  $\gamma_{\max} \leftarrow 1$ 
17:    else
18:       $\mathbf{d} \leftarrow \mathbf{d}^A$ ;  $\gamma_{\max} \leftarrow y_w / (1 - y_w)$ 
19:    end if
20:  end if
21:  Compute step size  $\gamma$  (Equation 8)
22:   $[\mathbf{u}^{(t+1)}; \mathbf{v}^{(t+1)}] \leftarrow [\mathbf{u}^{(t)}; \mathbf{v}^{(t)}] + \gamma \mathbf{d}$ 
23:  Update  $\mathcal{I}^{(t+1)}$  and  $\mathbf{y}^{(t+1)}$  accordingly.
24: end for
    
```

A.2. The Active Set Algorithm

We use a variant of the active set algorithm (Nocedal & Wright, 1999, Ch. 16.4 & 16.5) as proposed for the quadratic subproblems of the AD³ algorithm; our presentation follows (Martins et al., 2015, Algorithm 3). At each step, the active set algorithm solves a relaxed variant of the SparseMAP QP, relaxing the non-negativity constraint on \mathbf{y} , and restricting the solution to the current active set \mathcal{I}

$$\text{minimize}_{\mathbf{y}_{\mathcal{I}} \in \mathbb{R}^{|\mathcal{I}|}} \frac{1}{2} \|\mathbf{M}_{\mathcal{I}} \mathbf{y}_{\mathcal{I}}\|_2^2 - \boldsymbol{\eta}^\top \mathbf{A}_{\mathcal{I}} \mathbf{y}_{\mathcal{I}} \quad \text{subject to} \quad \mathbf{1}^\top \mathbf{y}_{\mathcal{I}} = 1$$

whose solution can be found by solving the KKT system

$$\begin{bmatrix} \mathbf{M}_{\mathcal{I}}^\top \mathbf{M}_{\mathcal{I}} & \mathbf{1} \\ \mathbf{1}^\top & 0 \end{bmatrix} \begin{bmatrix} \mathbf{y}_{\mathcal{I}} \\ \tau \end{bmatrix} = \begin{bmatrix} \mathbf{A}_{\mathcal{I}}^\top \boldsymbol{\eta} \\ 1 \end{bmatrix}. \quad (9)$$

At each iteration, the (symmetric) design matrix in Equation 9 is updated by adding or removing a row and a column; therefore its inverse (or a decomposition) may be efficiently maintained and updated.

Line search. The optimal step size for moving a feasible current estimate \mathbf{y}' toward a solution $\hat{\mathbf{y}}$ of Equation 9, while keeping feasibility, is given by (Martins et al., 2015, Equation 31)

$$\gamma = \min \left(1, \min_{s \in \mathcal{I}, y'_s > \hat{y}_s} \frac{y'_s}{y'_s - \hat{y}_s} \right) \quad (10)$$

When $\gamma \leq 1$ this update zeros out a coordinate of \mathbf{y}' ; otherwise, \mathcal{I} remains the same.

Algorithm 2 Active Set algorithm for SparseMAP

```

1: Initialization:  $s^{(0)} \leftarrow \text{MAP}_A(\eta_U, \eta_F)$ ;  $\mathcal{I}^{(0)} = \{s^{(0)}\}$ ;  $\mathbf{y}^{(0)} = \mathbf{e}_{s^{(0)}}$ ;  $[\mathbf{u}^{(0)}; \mathbf{v}^{(0)}] = \mathbf{a}_{s^{(0)}}$ 
2: for  $t = 0 \dots t_{\max}$  do
3:   Solve the relaxed QP restricted to  $\mathcal{I}^{(t)}$ ; get  $\hat{\mathbf{y}}, \hat{\tau}, \hat{\mathbf{u}} = M\hat{\mathbf{y}}$  (Equation 9)
4:   if  $\hat{\mathbf{y}} = \mathbf{y}^{(t)}$  then
5:      $s \leftarrow \text{MAP}_A(\eta_U - \hat{\mathbf{u}}, \eta_F)$ 
6:     if  $\text{gap}(\mathbf{a}_s, \hat{\mathbf{u}}) \leq \hat{\tau}$  then
7:       return  $\mathbf{u}^{(t)}$  (Equation 7)
8:     else
9:        $\mathcal{I}^{(t+1)} \leftarrow \mathcal{I}^{(t)} \cup \{s\}$ 
10:    end if
11:  else
12:    Compute step size  $\gamma$  (Equation 10)
13:     $\mathbf{y}^{(t+1)} \leftarrow (1 - \gamma)\mathbf{y}^{(t)} + \gamma\hat{\mathbf{y}}$  (sparse update)
14:    Update  $\mathcal{S}^{(t+1)}$  if necessary
15:  end if
16: end for

```

B. Computing the SparseMAP Jacobian: Proof of Proposition 1

Recall that SparseMAP is defined as the \mathbf{u}^* that maximizes the value of the quadratic program (Equation 5),

$$g(\eta_U, \eta_F) := \max_{[\mathbf{u}; \mathbf{v}] \in \mathcal{M}_A} \eta_U^\top \mathbf{u} + \eta_F^\top \mathbf{v} - \frac{1}{2} \|\mathbf{u}\|_2^2. \quad (11)$$

As the ℓ_2^2 norm is strongly convex, there is always a unique minimizer \mathbf{u}^* (implying that SparseMAP is well-defined), and the convex conjugate of the QP in (11), $g^*(\mathbf{u}, \mathbf{v}) = \{\frac{1}{2} \|\mathbf{u}\|_2^2, [\mathbf{u}; \mathbf{v}] \in \mathcal{M}_A; -\infty \text{ otherwise}\}$ is smooth in \mathbf{u} , implying that SparseMAP (which only returns \mathbf{u}) is Lipschitz-continuous and thus differentiable almost everywhere.

We now rewrite the QP in Equation 11 in terms of the convex combination of vertices of the marginal polytope

$$\min_{\mathbf{y} \in \Delta^D} \frac{1}{2} \|M\mathbf{y}\|_2^2 - \boldsymbol{\theta}^\top \mathbf{y} \quad \text{where } \boldsymbol{\theta} := A^\top \boldsymbol{\eta} \quad (12)$$

We use the optimality conditions of problem 12 to derive an explicit relationship between \mathbf{u}^* and \mathbf{x} . At an optimum, the following KKT conditions hold

$$M^\top M\mathbf{y}^* - \boldsymbol{\lambda}^* + \tau^* \mathbf{1} = \mathbf{0} \quad (13)$$

$$\mathbf{1}^\top \mathbf{y}^* = 1 \quad (14)$$

$$\mathbf{y}^* \geq \mathbf{0} \quad (15)$$

$$\boldsymbol{\lambda}^* \geq \mathbf{0} \quad (16)$$

$$\boldsymbol{\lambda}^{*\top} \mathbf{y}^* = 0 \quad (17)$$

Let \mathcal{I} denote the support of \mathbf{y}^* , i.e., $\mathcal{I} = \{s : y_s^* > 0\}$. From Equation 17 we have $\boldsymbol{\lambda}_{\mathcal{I}} = \mathbf{0}$ and therefore

$$M_{\mathcal{I}}^\top M_{\mathcal{I}} \mathbf{y}_{\mathcal{I}}^* + \tau^* \mathbf{1} = \boldsymbol{\theta}_{\mathcal{I}} \quad (18)$$

$$\mathbf{1}^\top \mathbf{y}_{\mathcal{I}}^* = 1 \quad (19)$$

Solving for $\mathbf{y}_{\mathcal{I}}^*$ in Equation 18 we get a direct expression

$$\mathbf{y}_{\mathcal{I}}^* = (M_{\mathcal{I}}^\top M_{\mathcal{I}})^{-1} (\boldsymbol{\theta}_{\mathcal{I}} - \tau^* \mathbf{1}) = \mathbf{Z} (\boldsymbol{\theta}_{\mathcal{I}} - \tau^* \mathbf{1}).$$

where we introduced $\mathbf{Z} = (\mathbf{M}^\top \mathbf{M})^{-1}$. Solving for τ^* yields

$$\tau^* = \frac{1}{\mathbf{1}^\top \mathbf{Z} \mathbf{1}} (\mathbf{1}^\top \mathbf{Z} \boldsymbol{\theta}_{\mathcal{I}} - 1)$$

Plugging this back and left-multiplying by $\mathbf{M}_{\mathcal{I}}$ we get

$$\mathbf{u}^* = \mathbf{M}_{\mathcal{I}} \mathbf{y}_{\mathcal{I}}^* = \mathbf{M}_{\mathcal{I}} \mathbf{Z} \left(\boldsymbol{\theta}_{\mathcal{I}} - \frac{1}{\mathbf{1}^\top \mathbf{Z} \mathbf{1}} \mathbf{1}^\top \mathbf{Z} \boldsymbol{\theta}_{\mathcal{I}} \mathbf{1} + \frac{1}{\mathbf{1}^\top \mathbf{Z} \mathbf{1}} \mathbf{1} \right)$$

Note that, in a neighborhood of $\boldsymbol{\eta}$, the support of the solution \mathcal{I} is constant. (On the measure-zero set of points where the support changes, SparseMAP is subdifferentiable and our assumption yields a generalized Jacobian (Clarke, 1990).) Differentiating w.r.t. the score of a configuration θ_s , we get the expression

$$\frac{\partial \mathbf{u}^*}{\partial \theta_s} = \begin{cases} \mathbf{M} \left(\mathbf{I} - \frac{1}{\mathbf{1}^\top \mathbf{Z} \mathbf{1}} \mathbf{Z} \mathbf{1} \mathbf{1}^\top \right) \mathbf{z}_s & s \in \mathcal{I} \\ \mathbf{0} & s \notin \mathcal{I} \end{cases} \quad (20)$$

Since $\theta_s = \mathbf{a}_s^\top \boldsymbol{\eta}$, by the chain rule, we get the desired result

$$\frac{\partial \mathbf{u}^*}{\partial \boldsymbol{\eta}} = \frac{\partial \mathbf{u}^*}{\partial \boldsymbol{\theta}} \mathbf{A}^\top. \quad (21)$$

C. Fenchel-Young Losses: Proof of Proposition 2

We recall that the structured Fenchel-Young loss defined by a convex $\Omega : \mathbb{R}^D \rightarrow \mathbb{R}$ and a matrix \mathbf{A} is defined as

$$\ell_{\Omega, \mathbf{A}} : \mathbb{R}^k \times \Delta^D \rightarrow \mathbb{R}, \quad \ell_{\Omega, \mathbf{A}}(\boldsymbol{\eta}, \mathbf{y}) := \Omega_\Delta^*(\mathbf{A}^\top \boldsymbol{\eta}) + \Omega_\Delta(\mathbf{y}) - \boldsymbol{\eta}^\top \mathbf{A} \mathbf{y}.$$

Since Ω_Δ is the restriction of a convex function to a convex set, it is convex (Boyd & Vandenberghe, 2004, Section 3.1.2).

Property 1. From the Fenchel-Young inequality (Fenchel, 1949; Boyd & Vandenberghe, 2004, Section 3.3.2), we have

$$\boldsymbol{\theta}^\top \mathbf{y} \leq \Omega_\Delta^*(\boldsymbol{\theta}) + \Omega_\Delta(\mathbf{y}).$$

In particular, when $\boldsymbol{\theta} = \mathbf{A}^\top \boldsymbol{\eta}$,

$$\begin{aligned} 0 &\leq -\boldsymbol{\eta}^\top \mathbf{A} \mathbf{y} + \Omega_\Delta^*(\mathbf{A}^\top \boldsymbol{\eta}) + \Omega_\Delta(\mathbf{y}) \\ &= \ell_{\Omega, \mathbf{A}}(\boldsymbol{\eta}, \mathbf{y}). \end{aligned}$$

Equality is achieved when

$$\begin{aligned} \Omega_\Delta^*(\mathbf{A}^\top \boldsymbol{\eta}) &= \boldsymbol{\eta}^\top \mathbf{A} \mathbf{y} - \Omega_\Delta(\mathbf{y}) \iff \\ \max_{\mathbf{y}' \in \Delta^d} \boldsymbol{\eta}^\top \mathbf{A} \mathbf{y}' - \Omega(\mathbf{y}') &= \boldsymbol{\eta}^\top \mathbf{A} \mathbf{y} - \Omega(\mathbf{y}), \end{aligned}$$

where we used the fact that $\mathbf{y} \in \Delta^d$. The second part of the claim follows.

Property 2. To prove convexity in $\boldsymbol{\eta}$, we rewrite the loss, for fixed \mathbf{y} , as

$$\ell_{\Omega, \mathbf{A}}(\boldsymbol{\eta}) = h(\mathbf{A}^\top \boldsymbol{\eta}) + \text{const}, \quad \text{where } h(\boldsymbol{\theta}) = \Omega_\Delta^*(\boldsymbol{\theta}) - \boldsymbol{\theta}^\top \mathbf{y}.$$

Ω_Δ^* is a convex conjugate, and thus itself convex. Linear functions are convex, and the sum of two convex functions is convex, therefore h is convex. Finally, the composition of a convex function with a linear function is convex as well, thus the function $(h \mathbf{A}^\top)$ is convex. Convexity of $\ell_{\Omega, \mathbf{A}}$ in $\boldsymbol{\eta}$ directly follows. Convexity in \mathbf{y} is straightforward, as the sum of a convex and a linear function (Boyd & Vandenberghe, 2004, Sections 3.2.1, 3.2.2, 3.3.1).

Property 3. This follows from the scaling property of the convex conjugate (Boyd & Vandenberghe, 2004, Section 3.3.2)

$$(t\Omega)^*(\boldsymbol{\theta}) = t\Omega^*(t^{-1}\boldsymbol{\theta})$$

Denoting $\boldsymbol{\eta}' = t^{-1}\boldsymbol{\eta}$, we have that

$$\begin{aligned} \ell_{t\Omega, \mathbf{A}}(\boldsymbol{\eta}, \mathbf{y}) &= (t\Omega_\Delta)^*(\mathbf{A}^\top \boldsymbol{\eta}) + t\Omega_\Delta(\mathbf{y}) - \boldsymbol{\eta}^\top \mathbf{A} \mathbf{y} \\ &= t\Omega_\Delta^*(\mathbf{A}^\top \boldsymbol{\eta}') + t\Omega_\Delta(\mathbf{y}) - \boldsymbol{\eta}^\top \mathbf{A} \mathbf{y} \\ &= t(\Omega_\Delta^*(\mathbf{A}^\top \boldsymbol{\eta}') + \Omega_\Delta(\mathbf{y}) - \boldsymbol{\eta}'^\top \mathbf{A} \mathbf{y}) = t\ell_{\Omega, \mathbf{A}}(t^{-1}\boldsymbol{\eta}, \mathbf{y}). \end{aligned}$$