
Supplementary Material: Visual Semantic Search: Retrieving Videos via Complex Textual Queries

Dahua Lin¹ **Sanja Fidler**^{1,2} **Chen Kong**³ **Raquel Urtasun**^{1,2}
 TTI Chicago¹ University of Toronto² Tsinghua University³
 dhlin@ttic.edu, kcl10@mails.tsinghua.edu.cn, {fidler, urtasun}@cs.toronto.edu

Abstract

This document provides some technical details related to the learning problem presented in the paper [1]. In particular, we review the concept of conciseness, and provide the proof to the Proposition 1 in [1], which establishes the fact that our learning problem is concise, and finally give the detailed derivation of the simplified optimization problem given in Eq.(7).

1 The Learning Problem

The problem of learning the optimal combination weights of scores was formulated in section 4.5.1 of the paper. For self-containedness, we briefly revisit the problem below.

$$\begin{aligned}
 & \text{minimize } \frac{1}{2} \|\mathbf{w}\|^2 + C \sum_i \xi_i & (1) \\
 & \text{s.t. } \mathbf{w}^T \phi_i(\mathbf{y}^{(i)}) \geq \mathbf{w}^T \phi_i(\mathbf{y}) + \Delta(\mathbf{y}, \mathbf{y}^{(i)}) - \xi_i, \quad \forall \mathbf{y}^{(i)} \in \mathcal{Y}^{(i)}. \\
 & \quad \xi_i \geq 0, \quad \forall i = 1, \dots, N.
 \end{aligned}$$

Here, $\mathbf{y}^{(i)}$ is the ground-truth matching for the i -th instance, $\phi_i(\mathbf{y})$ is a vector of matching scores for \mathbf{y} , and $\Delta(\mathbf{y}, \mathbf{y}^{(i)})$ the loss function. In particular, $\phi_i(\mathbf{y})$ can be expressed as

$$\phi_i(\mathbf{y}) = [\phi_i^{(1)}(\mathbf{y}), \dots, \phi_i^{(K)}(\mathbf{y})], \quad \text{with } \phi_i^{(k)}(\mathbf{y}) = \sum_{uv} f_{uv}^{(ik)} y_{uv}. \quad (2)$$

We use the Hamming loss, as

$$l(\mathbf{y}; \mathbf{y}^{(i)}) = \sum_{uv} \mathbf{1}(y_{uv} \neq y_{uv}^{(i)}) = a^{(i)} - \sum_{uv} y_{uv} y_{uv}^{(i)}, \quad (3)$$

where $a^{(i)} = \sum_u s_u^{(i)}$ is the total number of matching edges, which is a constant.

The domain $\mathcal{Y}^{(i)}$ depends on particular instance, and can be written as

$$\mathcal{Y}^{(i)} = \left\{ \mathbf{y} : \sum_v y_{uv} = s_u^{(i)}, \sum_u y_{uv} \leq t_v^{(i)}, 0 \leq y_{uv} \leq c_{uv}^{(i)} \right\}. \quad (4)$$

2 The Notion of Conciseness

The learning problem given by Eq.(1) can be re-written as

$$\begin{aligned}
 & \text{minimize } \frac{1}{2} \|\mathbf{w}\|^2 + C \sum_i \xi_i & (5) \\
 & \text{s.t. } \mathbf{w}^T \phi_i(\mathbf{y}^{(i)}) \geq \max_{\mathbf{y} \in \mathcal{Y}^{(i)}} \left(\mathbf{w}^T \phi_i(\mathbf{y}) + \Delta(\mathbf{y}, \mathbf{y}^{(i)}) \right) - \xi_i, \quad \xi_i \geq 0, \quad \forall i = 1, \dots, N.
 \end{aligned}$$

This model is called *concise* if there exists a function \tilde{f}_i that is concave in $\boldsymbol{\mu}$ and a convex set $\mathcal{U}^{(i)}$ for each i such that

$$\max_{\mathbf{y} \in \mathcal{Y}^{(i)}} \left(\mathbf{w}^T \boldsymbol{\phi}_i(\mathbf{y}) + \Delta(\mathbf{y}; \mathbf{y}^{(i)}) \right) = \max_{\boldsymbol{\mu} \in \mathcal{U}^{(i)}} \tilde{f}_i(\mathbf{w}, \boldsymbol{\mu}). \quad (6)$$

Next, we review how *conciseness* can be exploited to simplify the learning problem. Without losing generality, we express $\boldsymbol{\mu} \in \mathcal{U}^{(i)}$ using a convex function \mathbf{g}_i as

$$\mathbf{g}_i(\boldsymbol{\mu}) \leq 0. \quad (7)$$

Then the Lagrangian for $\tilde{f}_i(\mathbf{w}, \boldsymbol{\mu})$ is

$$L_i(\boldsymbol{\mu}, \boldsymbol{\lambda}; \mathbf{w}) = \tilde{f}_i(\mathbf{w}, \boldsymbol{\mu}) - \boldsymbol{\lambda}^T \mathbf{g}_i(\boldsymbol{\mu}) \quad \text{with } \boldsymbol{\lambda} \geq 0. \quad (8)$$

This provides an upper bound for $\tilde{f}_i(\mathbf{w}, \boldsymbol{\mu})$ within $\mathcal{U}^{(i)}$. By strong duality (which can be easily verified), we have:

$$\begin{aligned} \max_{\boldsymbol{\mu} \in \mathcal{U}^{(i)}} \tilde{f}_i(\mathbf{w}, \boldsymbol{\mu}) &= \max_{\boldsymbol{\mu} \in \mathcal{U}^{(i)}} \min_{\boldsymbol{\lambda} \geq 0} L_i(\boldsymbol{\mu}, \boldsymbol{\lambda}; \mathbf{w}), \\ &= \min_{\boldsymbol{\lambda} \geq 0} \max_{\boldsymbol{\mu} \in \mathcal{U}^{(i)}} L_i(\boldsymbol{\mu}, \boldsymbol{\lambda}; \mathbf{w}). \end{aligned} \quad (9)$$

Suppose $\max_{\boldsymbol{\mu} \in \mathcal{U}^{(i)}} L_i(\boldsymbol{\mu}, \boldsymbol{\lambda}, \boldsymbol{\nu}; \mathbf{w})$ has a Lagrangian dual given by

$$\rho_i(\boldsymbol{\lambda}; \mathbf{w}) \quad \text{s.t.} \quad \boldsymbol{\eta}_i(\boldsymbol{\lambda}; \mathbf{w}) \leq 0. \quad (10)$$

Then, we have

$$\max_{\boldsymbol{\mu} \in \mathcal{U}^{(i)}} \tilde{f}_i(\mathbf{w}, \boldsymbol{\mu}) = \min_{\boldsymbol{\eta}^{(i)}(\boldsymbol{\lambda}; \mathbf{w}) \leq 0} \rho_i(\boldsymbol{\lambda}; \mathbf{w}) \quad (11)$$

For conciseness, the condition $\boldsymbol{\lambda} \geq 0$ is merged into $\boldsymbol{\eta}_i(\boldsymbol{\lambda}; \mathbf{w}) \leq 0$. Incorporating this into Eq.(5) results in

$$\begin{aligned} \text{minimize} \quad & \frac{1}{2} \|\mathbf{w}\|^2 + C \sum_i \xi_i \\ \text{s.t.} \quad & \mathbf{w}^T \boldsymbol{\phi}_i(\mathbf{y}^{(i)}) \geq \min_{\boldsymbol{\eta}^{(i)}(\boldsymbol{\lambda}; \mathbf{w}) \leq 0} \rho_i(\boldsymbol{\lambda}; \mathbf{w}) - \xi_i, \quad \xi_i \geq 0, \quad \forall i = 1, \dots, N. \end{aligned} \quad (12)$$

Combining the optimization over \mathbf{w} and that over $\boldsymbol{\lambda}$, we finally gets the following problem:

$$\begin{aligned} \text{minimize} \quad & \frac{1}{2} \|\mathbf{w}\|^2 + C \sum_i \xi_i \\ \text{s.t.} \quad & \mathbf{w}^T \boldsymbol{\phi}_i(\mathbf{y}^{(i)}) \geq \rho_i(\boldsymbol{\lambda}, \boldsymbol{\nu}; \mathbf{w}) - \xi_i, \quad \forall i = 1, \dots, N, \\ & \boldsymbol{\eta}_i(\boldsymbol{\lambda}, \boldsymbol{\nu}; \mathbf{w}) \leq 0, \quad \xi_i \geq 0, \quad \forall i = 1, \dots, N. \end{aligned} \quad (13)$$

3 Proof of Proposition 1

Proposition 1 in the paper establishes the fact that our learning problem is *concise*. Below, we prove this proposition.

With Eq.(2) and Eq.(3), we have

$$\begin{aligned} \mathbf{w}^T \boldsymbol{\phi}_i(\mathbf{y}) + \Delta(\mathbf{y}; \mathbf{y}^{(i)}) &= \sum_{k=1}^K w_k \sum_{uv} f_{uv}^{(ik)} y_{uv} + \left(a^{(i)} - \sum_{uv} y_{uv} y_{uv}^{(i)} \right) \\ &= a^{(i)} + \sum_{uv} \left(\sum_{k=1}^K w_k f_{uv}^{(ik)} - y_{uv}^{(i)} \right) y_{uv} \\ &= a^{(i)} + \left(\mathbf{F}^{(i)} \mathbf{w} - \mathbf{y}^{(i)} \right)^T \mathbf{y}. \end{aligned} \quad (14)$$

Here, each $\mathbf{F}^{(i)}$ is an mn -by- K matrix, where each row corresponding to a particular matching pair (u, v) and each column corresponds to a score channel. According to Eq.(6), we can conclude that this model is *concise*, with

$$\begin{aligned}\tilde{f}_i(\mathbf{w}, \boldsymbol{\mu}) &= a^{(i)} + \left(\mathbf{F}^{(i)} \mathbf{w} - \mathbf{y}^{(i)} \right)^T \boldsymbol{\mu} \\ &= a^{(i)} + \sum_{uv} \left(\mathbf{w}^T \mathbf{f}_{uv}^{(i)} - y_{uv}^{(i)} \right) \mu_{uv}.\end{aligned}\quad (15)$$

Here, $\mathbf{f}_{uv}^{(i)}$ is the uv -th row of $\mathbf{F}^{(i)}$, which is a K -dimensional vector. In addition, the constraint $\boldsymbol{\mu} \in \mathcal{U}^{(i)}$ can be written explicitly as

$$\sum_v \mu_{uv} = s_u^{(i)} \quad \forall u, \quad \sum_u \mu_{uv} \leq t_v^{(i)} \quad \forall v, \quad 0 \leq \mu_{uv} \leq c_{uv}^{(i)} \quad \forall u, v. \quad (16)$$

The proof is completed.

4 Simplified Optimization Problem

Then, we can derive the Lagrangian dual as follows

$$\rho^{(i)}(\boldsymbol{\lambda}, \boldsymbol{\eta}, \boldsymbol{\nu}, \mathbf{w}) = a^{(i)} + \sum_u \lambda_u s_u^{(i)} + \sum_v \eta_v t_v^{(i)} + \sum_{uv} \nu_{uv} c_{uv}^{(i)}, \quad (17)$$

with

$$\mathbf{w}^T \mathbf{f}_{uv}^{(i)} \leq y_{uv}^{(i)} + \lambda_u + \eta_v + \nu_{uv}, \quad \eta_v \geq 0, \quad \nu_{uv} \geq 0 \quad \forall u, v. \quad (18)$$

Finally, according to Eq.(13), the learning problem can be written as

$$\begin{aligned}\text{minimize} \quad & \frac{1}{2} \|\mathbf{w}\|^2 + C \sum_i \xi_i \\ \text{s.t.} \quad & \mathbf{w}^T \mathbf{z}^{(i)} \geq \rho^{(i)}(\boldsymbol{\lambda}, \boldsymbol{\eta}, \boldsymbol{\nu}, \mathbf{w}) - \xi_i, \quad \forall i = 1, \dots, N, \\ & \mathbf{w}^T \mathbf{f}_{uv}^{(i)} \leq y_{uv}^{(i)} + \lambda_u^{(i)} + \eta_v^{(i)} + \nu_{uv}^{(i)}, \quad \forall u, v, i \\ & \eta_v^{(i)} \geq 0, \quad \nu_{uv}^{(i)} \geq 0, \quad \xi^{(i)} \geq 0, \quad \forall u, v, i\end{aligned}\quad (19)$$

Here, $\mathbf{z}^{(i)} = [z_1^{(i)}, \dots, z_K^{(i)}]$ with $z_k^{(i)} = \sum_{uv} f_{uv}^{(ik)} y_{uv}^{(i)}$.

References

- [1] D. Lin, S. Fidler, C. Kong, and R. Urtasun. Visual semantic search: Retrieving videos via complex textual queries. In *CVPR*, 2014.