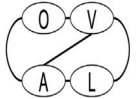




Partial Linearization based Optimization for Multi-Class SVM



Prithiv Mohapatra¹, Puneet K. Dokania², C.V. Jawahar¹, M. Pawan Kumar³

¹International Institute of Information Technology, Hyderabad

²Centrale Supélec & INRIA Saclay, ³University of Oxford



1. Short Summary

Contribution

A novel partial linearization [1] based approach for optimizing the multi-class SVM learning problem. Our method is an intuitive **generalization of the Frank-Wolfe [2] and the Exponentiated-Gradient [3] algorithms.**

Advantages

• Informative descent direction

computed using the marginals over each output class, similar to Exponentiated-Gradient.

• Optimal step-size in the descent direction

that guarantees an increase in the dual objective, similar to Frank-Wolfe.

• Block coordinate formulation

similar to the one proposed for Frank-Wolfe, which allows us to solve large-scale problems.

2. Multi-class SVM

- Input: $x \in X$; Output: $y \in Y$, where $Y = \{1, \dots, c\}$ & c is the no. of classes
- Feature representation for sample: $\phi(x)$
- Joint Feature map for sample and candidate class:

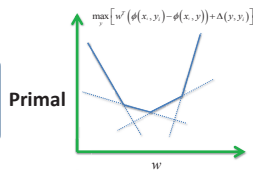
$$\phi(x, y) = [v_1^T \dots v_j^T \dots v_c^T] \quad \text{where, } v_j = \begin{cases} \phi(x) & \text{if } j = y. \\ 0 & \text{otherwise.} \end{cases}$$

• **Prediction:** $y_{opt} = \arg \max_y w^T \phi(x, y)$

• Learning:

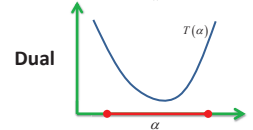
$$\min_{w, \xi} \frac{\lambda}{2} \|w\|^2 + \frac{1}{n} \sum_{i=1}^n \xi_i$$

s.t. $\forall y, i: w^T \phi(x_i, y_i) - w^T \phi(x_i, y) \geq \Delta(y, y_i) - \xi_i$



$$\min_{\alpha} T(\alpha) = -b^T \alpha + \frac{\lambda}{2} \alpha^T A^T A \alpha$$

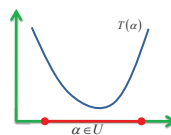
s.t. $\sum_{y \in Y} \alpha_y = 1; \forall i \in [n]$



3. Partial Linearization framework

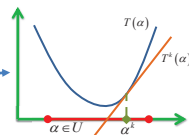
- Consider Optimization problem: $\min_{\alpha \in U} T(\alpha)$

Where, $T(\alpha)$ is a convex and continuously differentiable function and U is a compact & convex set.



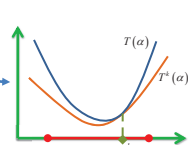
Optimization using Partial Linearization

- $k=0$, Initialize α^k .
- Repeat until convergence:
 - Construct local approximation $T^k(\alpha)$
 - Compute: $s^k = \arg \min_{\alpha \in U} T^k(\alpha)$
 - Compute optimal step size γ
 - Update $\alpha, \alpha^{k+1} \leftarrow (1-\gamma)\alpha^k + (\gamma)s^k$
 - $k \leftarrow k+1$
- Return optimal α .



Local Linear approximation (as in Frank-Wolfe)

$$T^k(\alpha) = T(\alpha^k) + [\nabla T(\alpha^k)]^T (\alpha - \alpha^k)$$



Local Partial-linear approximation

$$T(\alpha) = f(\alpha, \alpha^k) + (T(\alpha) - f(\alpha, \alpha^k))$$

$$T^k(\alpha) = f(\alpha, \alpha^k) + \left[(T(\alpha^k) - f(\alpha^k, \alpha^k)) + [\nabla T(\alpha^k) - \nabla_{\alpha} f(\alpha^k, \alpha^k)]^T (\alpha - \alpha^k) \right]$$

Linearization of only the Error term

6. References

- [1] M. Patriksson, "Partial linearization methods in nonlinear programming. Journal of Optimization Theory and Applications", Springer, 1993.
- [2] S. Lacoste-Julien, M. Jaggi, M. Schmidt, P. Pletscher, "Block-coordinate frank-wolfe for structural SVMs", ICML, 2012.
- [3] M. Collins, A. Globerson, T. Koo, X. Carreras, P. L. Bartlett, "Exponentiated gradient algorithms for conditional random fields and max-margin markov networks", The Journal of Machine Learning Research, 2008.

4. Partial Linearization based Optimization for Multi-class SVM

Appropriate design of Surrogate function:

$$f(\alpha, \alpha^k) = \frac{\tau}{n} \sum_{i \in [n]} \sum_{y \in Y} \alpha_y \log(\alpha_y)$$

- Resulting problem is efficiently optimizable
- Captures information about violation of each primal constraint

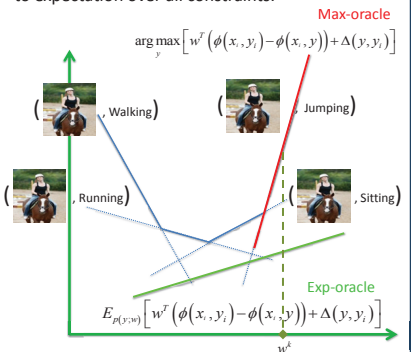
Temperature: Non-negative hyper-parameter

Corresponding update direction: $s_y^k = \frac{\exp\left(\log(\alpha_y^{k-1}) + \frac{1}{\tau} (\Delta_i(y) - W^{k-1 T} (\phi(x_i, y_i) - \phi(x_i, y)))\right)}{z_i}$

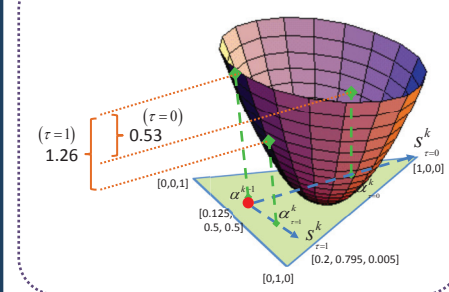
Partial Linearization for Optimizing Multi-Class SVM

- $k=0$, Initialize α^k .
- Repeat until convergence:
 - for all $i \in [n]$ do:
 - $\forall y \in Y$
 - $s_y^k = \frac{\exp\left(\log(\alpha_y^{k-1}) + \frac{1}{\tau} (\Delta_i(y) - W^{k-1 T} (\phi(x_i, y_i) - \phi(x_i, y)))\right)}{z_i}$
 - Compute optimal step size, $\gamma \leftarrow \frac{(\alpha^{k+1} - s^k - b + A^T A \alpha^{k+1})}{\|A(\alpha^{k+1} - s^k)\|}$
 - Update $\alpha, \alpha^{k+1} \leftarrow (1-\gamma)\alpha^k + (\gamma)s^k$
 - Update $w, w^k \leftarrow A\alpha^k$
 - $k \leftarrow k+1$
- Return optimal w .

In primal space, new update direction corresponds to expectation over all constraints.



Illustrative example: Update steps for $\tau=0$ and $\tau>0$



- For Temperature: $\tau=0$, Partial-Linearization \leftrightarrow Frank-Wolfe
- For $\tau>0$, can have potentially better update direction
- For Step-size: $\gamma=1$, Partial-Linearization \leftrightarrow Exponentiated-Gradient
- For $\gamma<1$, can have potentially better update step

Block Coordinate Partial Linearization

- Useful/Necessary for very large datasets.
- A coordinate descent like approach.
- Updates model after every single sample and is more efficient than the batch version.

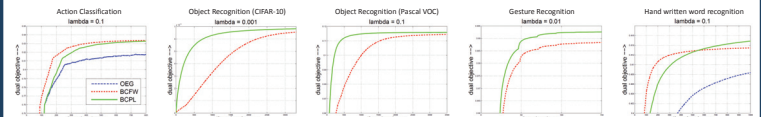
Partial Linearization based optimization for Structured-SVM

- Can be used for tree-structured output spaces.
- Message passing for computing exact marginals for the update step.

5. Experiments and Results

- **Problems and Datasets:** Training Multi-Class SVM or Structured SVM models for Vision tasks like Classification (Pascal VOC), Object Recognition (CIFAR-10, Pascal VOC), Gesture Recognition (MSRC-12) and Hand written word recognition (OCR).
- **Methods:** We compare our Block-Coordinate Partial Linearization algorithm (BCPL) with Block-coordinate Frank-Wolfe (BCFW) algorithm [2] and Online exponentiated gradient (OEG) algorithm. We use a fixed temperature of $\tau=0.01$ for our BCPL algorithm in all the experiments.
- **Results:**

Example training curves (Dual objective Vs. time (sec)) :



Mean training time (sec) of Multi-Class SVM or Structured-SVM models for different tasks:

