

Overview

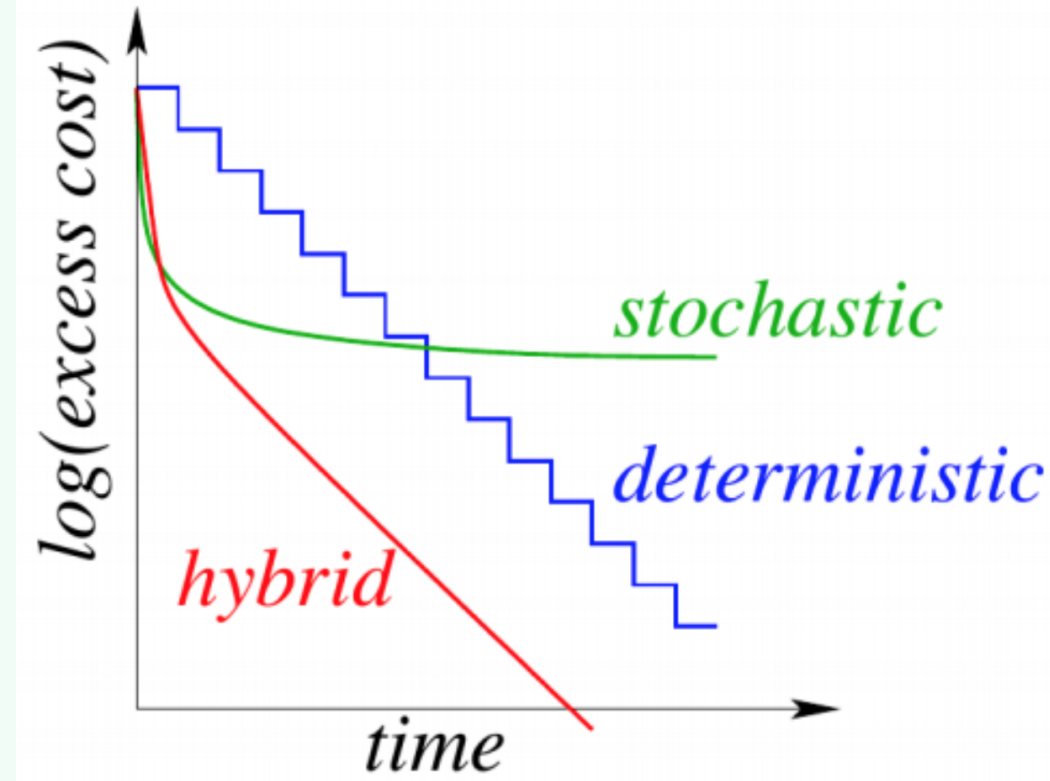
Goal

Fast and exact optimization of Conditional Random Fields

State of the Art

Variance reduced methods (hybrid):

- Stochastic Average Gradient (SAG) [3]
- Online Exponentiated Gradient (OEG) [1]
- Stochastic Dual Coordinate Ascent (SDCA) [this work]

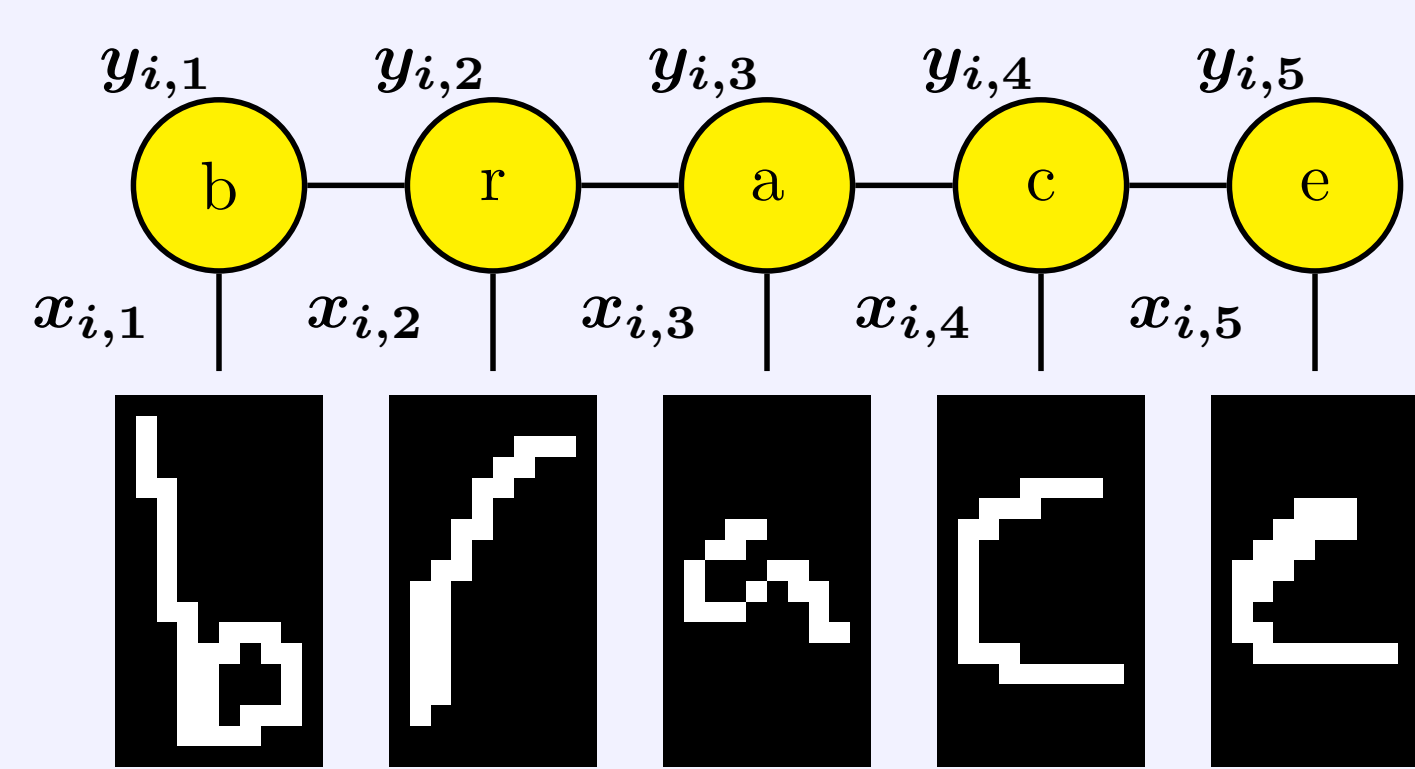


Contributions

- Adapt the algorithm SDCA to the CRF setting.
- Accelerate SDCA with an adaptive sampling strategy (with proof).
- Get state of the art optimization speed on sparse datasets.

Background

Structured Prediction



sample x	The	Chicago	Bulls	won.
label y	0	B-Organization	I-Organization	0
independent guess	0	B-Location	0	0

a) OCR

b) NER

Data: inputs $x_i \in \mathcal{X}$ and structured labels $y_i \in \mathcal{Y}$ (e.g. sequence) for $i \in \{1, \dots, n\}$.

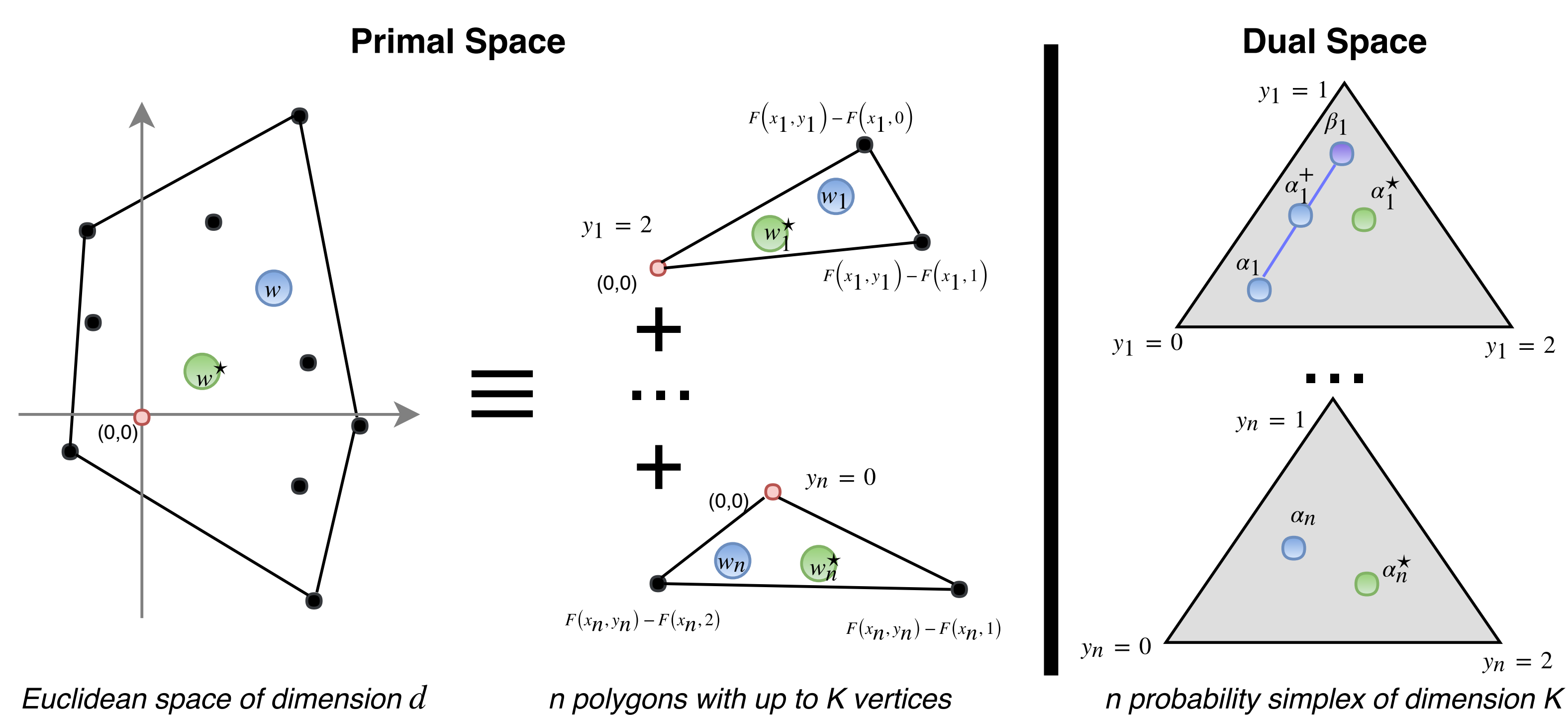
Conditional Random Fields

- **Model:** exponential family with sufficient statistic F : $p(y|x; \mathbf{w}) \propto \exp(\mathbf{w}^T F(x, y))$.
- **Goal:** learn the optimal parameter \mathbf{w}^* .
- **Problem:** partition function $= \sum_y e^{\mathbf{w}^T F(x, y)}$ = intractable sum over \mathcal{Y} (huge set)
- **Solution:** graphical model of $y|x$ so that F decomposes as a sum over the cliques of the graph. For a sequential graph with T nodes as in OCR:

$$F(x, y) = \sum_{t=1}^T F_t(x_t, y_t) + \sum_{t=1}^{T-1} F_t(y_t, y_{t+1}) \quad (1)$$

Use message passing to evaluate the partition function and marginals.

Parameters



Optimization Problem

Primal and Dual

- **Primal:** minimize the l_2 -regularized negative log-likelihood.

$$\min_{\mathbf{w} \in \mathbb{R}^d} \mathcal{P}(\mathbf{w}) := \frac{\lambda}{2} \|\mathbf{w}\|^2 + \frac{1}{n} \sum_{i=1}^n -\log p(y_i|x_i; \mathbf{w}) \quad (2)$$

- **Dual:** maximize the entropy-regularized negative squared loss.

$$\max_{\alpha} \mathcal{D}(\alpha) := -\frac{\lambda}{2} \|\hat{\mathbf{w}}(\alpha)\|^2 + \frac{1}{n} \sum_{i=1}^n H(\alpha_i) \quad (3)$$

- $\forall i, \alpha_i \in \Delta_{|\mathcal{Y}|}$ is a distribution over the labels \mathcal{Y} .

- Dual weights $\hat{\mathbf{w}}(\alpha)$ and entropy H are defined by:

$$\hat{\mathbf{w}}(\alpha) := \frac{1}{\lambda n} \sum_i (F(x_i, y_i) - \mathbb{E}_{y \sim \alpha_i}[F(x_i, y)]), \quad H(\alpha_i) := -\sum_{y \in \mathcal{Y}} \alpha_i(y) \log(\alpha_i(y)).$$

Properties

- Primal probabilities $\hat{\alpha}_i(\mathbf{w}) = p(\cdot|x_i; \mathbf{w})$. Dual weights $\hat{\mathbf{w}}(\alpha)$.
- **Fixed point property:** $\hat{\mathbf{w}}(\alpha^*) = \mathbf{w}^*$ and $\hat{\alpha}(\mathbf{w}^*) = \alpha^*$.
- **Duality gap:** $g(\mathbf{w}, \alpha) := \mathcal{P}(\mathbf{w}) - \mathcal{D}(\alpha) \geq 0$ and $g(\mathbf{w}^*, \alpha^*) = 0$.

$$g(\hat{\mathbf{w}}(\alpha), \alpha) = \frac{1}{n} \sum_i D_{KL}(\alpha_i || \hat{\alpha}_i(\hat{\mathbf{w}}(\alpha))) = \frac{1}{n} \sum_i g_i \quad (4)$$

- Idea: the *individual duality gap* g_i is the sub-optimality of data point i .

SDCA for CRF

Prox-SDCA algorithm

- Primal-dual method for Empirical Risk Minimization. Store $(\alpha, \hat{\mathbf{w}}(\alpha))$.
- At each step, update a random block α_i to maximize $\mathcal{D}(\alpha)$.
- Guaranteed ascent direction: $\delta_i = \hat{\alpha}_i(\hat{\mathbf{w}}(\alpha)) - \alpha_i$

$$\alpha_i^+ \leftarrow \alpha_i + \gamma \delta_i = (1 - \gamma) \alpha_i + \gamma \hat{\alpha}_i(\hat{\mathbf{w}}(\alpha)), \quad \gamma \in [0, 1] \quad (5)$$

- Exact line search on step-size $\gamma \in [0, 1]$ with Newton-Raphson.

Adaptation to CRFs

- **Problem:** α_i has dimension $|\mathcal{Y}| \gg 1$.
- **Hypothesis:** the CRF graph has a **junction tree** $(\mathcal{C}, \mathcal{S})$. (Similar to [1].)
- **Solution:** Replace joint probability α_i by marginal probabilities $\mu_{i,C}$ on maximal cliques of the graph $C \in \mathcal{C}$. Dual objective can be expressed only with these marginals.

$$\alpha_i(y) = \frac{\prod_{C \in \mathcal{C}} \mu_{i,C}(y_C)}{\prod_{S \in \mathcal{S}} \mu_{i,S}(y_S)} \quad \text{implies} \quad H(\alpha_i) = \sum_{C \in \mathcal{C}} H(\mu_{i,C}) - \sum_{S \in \mathcal{S}} H(\mu_{i,S}) =: \tilde{H}(\mu_i).$$

$$F(x, y) = \sum_{C \in \mathcal{C}} F_C(x, y_C) \quad \text{implies} \quad \hat{\mathbf{w}}(\alpha) = \sum_{C \in \mathcal{C}} \tilde{\mathbf{w}}_C(\mu_C) =: \tilde{\mathbf{w}}(\mu).$$

- Run message passing to infer $\hat{\mu}_{i,C}(\mathbf{w})(y_C) = p(y_C = \cdot | x_i; \mathbf{w})$.

Adaptive Sampling strategy

- Previous work [2] : sample proportionally to $\|\delta_i\|_1 = \|\alpha_i - \hat{\alpha}_i(\hat{\mathbf{w}}(\alpha))\|_1$.
- Problem: $\|\delta_i\|_1$ cannot be expressed with the marginals.
- Our strategy: sample proportionally to individual duality gaps g_i .
- Strong convexity: proof of acceleration by a factor $\chi(\mathbf{g})^2 = \frac{1}{n} \sum_i g_i^2 / (\frac{1}{n} \sum_i g_i)^2 \geq 1$.

Pseudocode

```

Initialize  $\mathbf{w} \leftarrow 0$  and  $\mu_{i,C}(y_C) \leftarrow \mathbb{1}_{\{y_C=y_{i,C}\}}, \forall i, \forall C \in \mathcal{C}, \forall y_C \in \mathcal{Y}_C$ .
(Optional) Let  $g_i \leftarrow 1000, \forall i$  and  $\bar{g} \leftarrow \frac{1}{n} \sum_i g_i$ . {duality gap estimate}
while  $\bar{g} > \text{required precision}$  do
    Sample  $i$  in  $\{1, \dots, n\}$  uniformly at random OR proportionally to  $g_i$ .
    Let  $\nu_{i,C}(y_C) \leftarrow p(y_C|x_i; \mathbf{w}), \forall C \in \mathcal{C}$  {message passing oracle}
    (Optional) Let  $g_i \leftarrow \tilde{D}(\mu_i || \nu_i)$  {individual duality gap}
    Let  $\delta_i \leftarrow \nu_i - \mu_i$  {dual ascent direction}
    Let  $\mathbf{v}_i \leftarrow \hat{\mathbf{w}}(\delta_i) = \frac{1}{\lambda n} \sum_{C \in \mathcal{C}} \mathbb{E}_{\mu_{i,C}}[F(x_i, y_C)] - \mathbb{E}_{\nu_{i,C}}[F(x_i, y_C)]$  {primal direction}
    Solve  $\gamma^* = \arg \max_{\gamma \in [0,1]} \tilde{H}(\mu_i + \gamma \delta_i) - \frac{\lambda n}{2} \|\mathbf{w} + \gamma \mathbf{v}_i\|^2$  {line search}
    Update  $\mu_i \leftarrow \mu_i + \gamma^* \delta_i$ 
    Update  $\mathbf{w} \leftarrow \hat{\mathbf{w}}(\mu) = \mathbf{w} + \gamma^* \mathbf{v}_i$ 
return  $\mathbf{w}$ 

```

Experiments

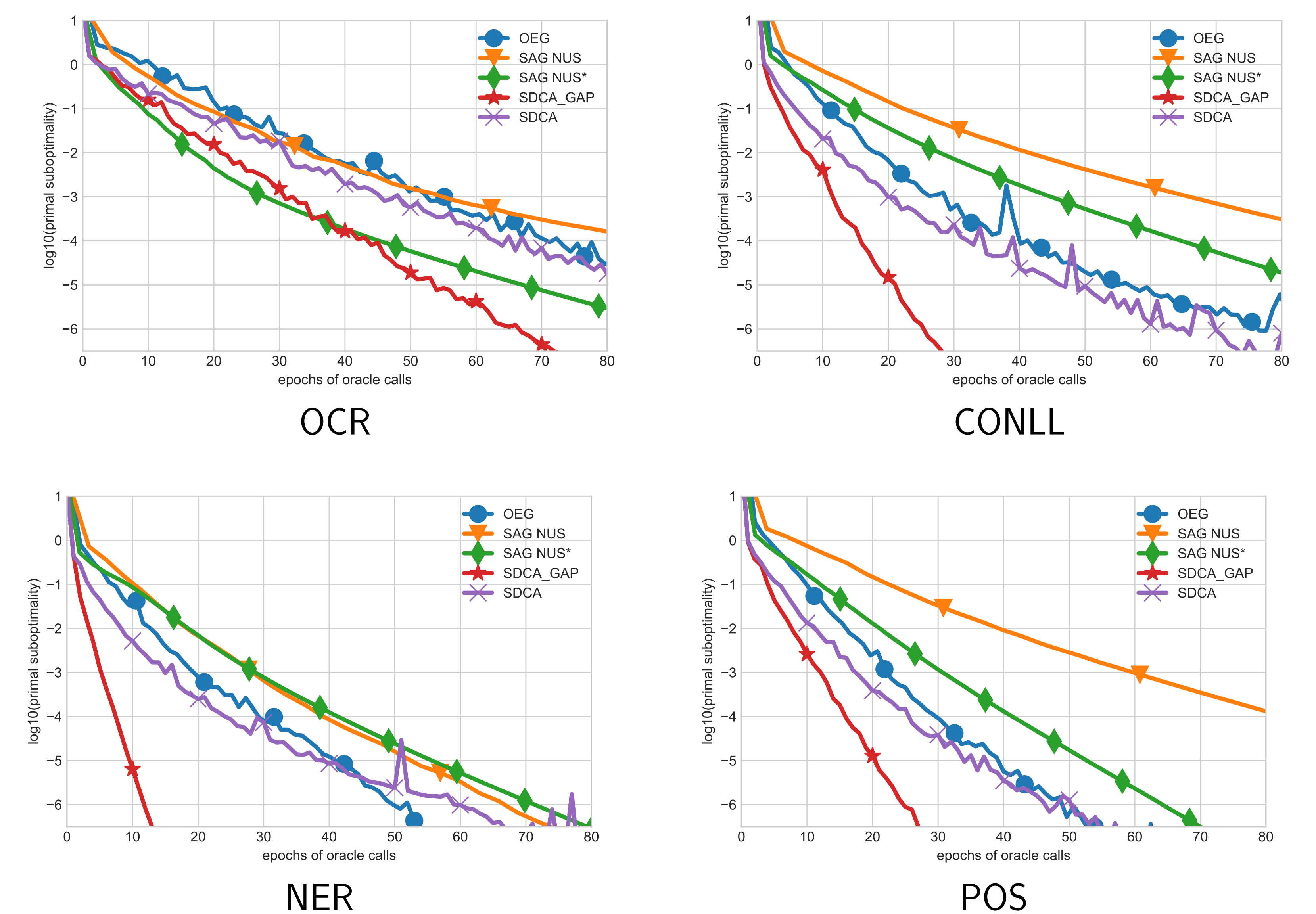


Figure 1: Primal sub-optimality as a function of the number of oracle calls (left). SDCA refers to uniform sampling. SDCA-GAP refers to sampling proportionally to the gaps 80% of the time. SAG-NUS performs a line search at every iteration.

- OCR has dense features (images). All methods perform comparably.
- CONLL, NER and POS are language understanding tasks with sparse features. Dual methods OEG and SDCA tend to perform better.
- The gap sampling strategy gives a significant advantage to SDCA.

Take-away

- Dual variance reduced methods can be applied to structured problems.
- Duality gap sampling and exact line search makes them fast on sparse datasets.

References

- [1] M. Collins, A. Globerson, T. Koo, X. Carreras, and P. L. Bartlett. Exponentiated gradient algorithms for conditional random fields and max-margin Markov networks. *Journal of Machine Learning Research*, 2008.
- [2] D. Csiba, Z. Qu, and P. Richtárik. Stochastic dual coordinate ascent with adaptive probabilities. In *ICML*, 2015.
- [3] M. Schmidt, R. Babanezhad, M. Ahmed, A. Defazio, A. Clifton, and A. Sarkar. Non-uniform stochastic average gradient method for training conditional random fields. In *AISTATS*, 2015.