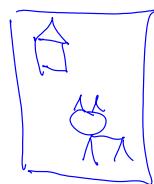


- today:
- latent SVM struct
 - kernels
 - deep learning?

I) Latent variables

motivation: semantic segmentation \rightarrow finding boundary of different objects



segmentation is expensive \rightarrow z "latent variable"
only have class labels $\rightarrow y$

also: [Felzenszwalb et al TPAMI 2010
"Deformable Part models" for object recognition
 $\rightarrow z$ there was an object configuration]

before, we had $s(x, y; w) = \langle w, \varphi(x, y) \rangle$

now, consider $s(x, y, z; w) = \langle w, \varphi(x, y, z) \rangle$

as before, add product with $\operatorname{argmax}_{y \in \mathcal{Y}, z \in \mathcal{Z}} s(x, y, z; w)$

* CRF ($p(y|x)$) \rightarrow hidden CRF similar to latent variable modeling with graphical model
ML \rightarrow Expectation-Maximization

\hookrightarrow analog for SVMstruct
is CCCP

Latent SVM struct:

$$\ell(y, (\hat{y}, \hat{z}))$$

generalization of structured hinge loss:

$$g(x, y, w) \triangleq \max_{(\hat{y}, \hat{z})} \langle w, \ell(x, \hat{y}, \hat{z}) \rangle + \ell(y, (\hat{y}, \hat{z})) - \underbrace{\max_{\hat{z}} \langle w, \ell(x, y, \hat{z}) \rangle}_{\text{concave function of } w} \geq \ell(y, h_w(x))$$

recall: sup of convex function \Rightarrow always convex
inf of jointly convex function is convex i.e. $\inf_z g(w, z)$ where g is jointly convex in $w \setminus z$
 $\hat{g}(w)$ is convex in w

$$\text{here } f(x, y, w) = u(w) - v(w)$$

where $u \setminus v$ are convex functions of w "difference of convex functions"

\hookrightarrow CCCP procedure
is used to minimize this

CCCP procedures

- linearize $v(w)$ at w_t to get upper bound
- minimize the upper bound
- repeat

\rightarrow majorization-minimization procedure
(EM is another example of that)

$$\begin{cases} \mathcal{L}_t(w) = u(w) - [v(w_t) + \langle \nabla v(w_t), w - w_t \rangle] \\ \quad \text{(or subgradient)} \\ w_{t+1} = \underset{w}{\operatorname{arg\min}} \mathcal{L}_t(w) \end{cases}$$

properties of this procedure:

- like EM, descent procedure i.e. $\mathcal{L}(w_{t+1}) \leq \mathcal{L}(w_t)$

$$g(w_t) = \mathcal{L}_t(w_t) \geq \mathcal{L}_t(w_{t+1}) \geq \mathcal{L}(w_{t+1})$$

- Normalization converges to solution much

- local linear convergence to stationary point
for latent SVMstruct [see NIPS OPT 2012 paper]

for SVMstruct

$$V(w) = \max_{\tilde{z}} \langle w, \ell(x, y, \tilde{z}) \rangle$$

$\xrightarrow{\text{argmax}_{\tilde{z}} \langle w, \ell(x, y, \tilde{z}) \rangle}$

$$\partial V(w_t) = \ell(x, y, \hat{z}(x, w_t))$$

$$\Rightarrow f_t(x, y, w) = \max_{(\tilde{y}, \tilde{z})} \langle w, \ell(x, \tilde{y}, \tilde{z}) \rangle + \ell(y, \tilde{y}, \tilde{z}) \rightarrow \langle w, \ell(x, y, \hat{z}_t) \rangle + \text{cst.}$$

→ like SVMstruct objective

CCCP for Latent SVM struct repeat:

- fill in $\hat{z}_t^{(i)}$ for all ground truth $y^{(i)}$ using w_t
- solve standard SVMstruct b get w_t
- repeat

Kernels?

$$\text{so far } s(x, y; w) = \langle w, \ell(x, y) \rangle \quad \nabla \ell(x^{(i)}, y^{(i)}) - \ell(x^{(i)}, \tilde{y})$$

$$\text{recall for both SVMstruct } w(\alpha) = \frac{1}{\lambda n} \sum_{i, \tilde{y}} \alpha_i \psi_i(\tilde{y}) \psi_i(\tilde{y})^\top$$

$$\langle w, \ell(x, y) \rangle = \frac{1}{\lambda n} \sum_{i, \tilde{y}} \alpha_i \psi_i(\tilde{y}) \underbrace{\langle \psi_i(\tilde{y}), \ell(x, y) \rangle}_{K((x^{(i)}, y^{(i)}); (x, y)) - K((x^{(i)}, \tilde{y}); (x, y))}$$

$$K((x, y), (x', y')) \triangleq \langle \ell(x, y), \ell(x', y') \rangle$$

$$\text{examples: } K(\dots) = k_x(x, x') \cdot k_y(y, y')$$

$$\text{OCR} \quad \ell(x, y) = \sum_p \ell_n(x_p, y_p) + \sum_p \ell_e(y_p, y_{p+1})$$

$$\langle \ell(x, y), \ell(x', y') \rangle = \sum_{p \neq p'} (\langle \ell_n(x_p, y_p), \ell_n(x'_p, y'_p) \rangle + \langle \ell_e(y_p, y_{p+1}), \ell_e(y'_p, y'_{p+1}) \rangle)$$

[$\ell_n \perp \ell_e$]

for OCR, we had used $\ell_n(x_p, y_p) = \begin{cases} 0 & \text{vec}(x_p) \rightarrow y_p \\ \text{position} \end{cases}$

$$\Rightarrow K((x_p, y_p), (x'_p, y'_p)) = \mathbb{1}\{y_p = y'_p\} \langle x_p, x'_p \rangle$$

but instead, could use $\mathbb{1}\{y_p = y'_p\} \exp(-\frac{\|x_p - x'_p\|^2}{2\sigma^2})$

\downarrow
RBF kernel

* binary SVM with RBF kernel \rightarrow similar to nearest neighbor classifier

Computation:

BCFW: stored $w_i \rightarrow O(d)$ for each i

$$\langle w, \ell(x, y) \rangle \sim O(d)$$

with kernel: cannot store w , so maintain $\alpha_i(\tilde{y})$ instead ..

$$w = \sum_{i \in \mathcal{Y}} \alpha_i(\tilde{y}) \varphi_i(\tilde{y})$$

\mathcal{Y} (super)
Sparse when using FW

after t iterations of BCFW, only t non-zero variables

$$\langle w, \ell(x, y) \rangle \rightarrow O(t)$$

kernel methods $\sim O(n^2)$

$\sim n$ $\sim O(n^2)$

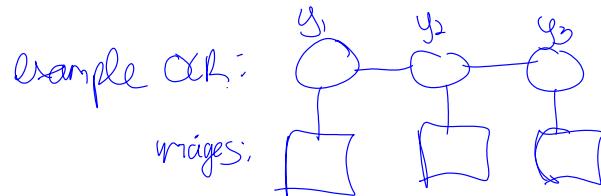
primal $\rightarrow O(d \cdot n)$

Deep learning

go from $\langle w, \phi(x, y) \rangle$ to $\langle w, \phi(x, y, \theta) \rangle$

can learn θ

I) plugging "deep learning" features in a structured prediction model



example:
[Vu et al., ICCV 2015] "context-aware CNNs for person head detection"

II) recurrent neural network (RNN)

$$p(y|x) = \prod_{t \leq T} p(y_t | y_{1:t-1}, x)$$

chain rule
RNN \rightarrow structured parameterization of $p(y_t | y_{1:t-1}, x)$

graphical model approach

conditional independence assumption

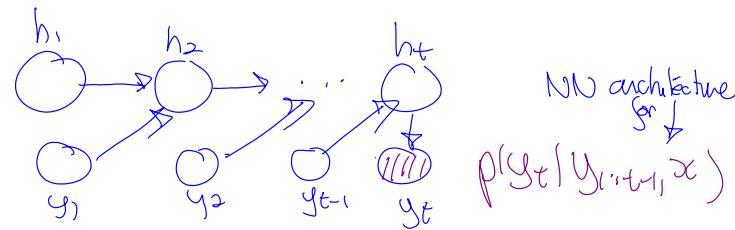
$$\prod_t p(y_t | y_{\pi_t}, x) \quad (\text{directed})$$
$$\frac{1}{Z(x)} \prod_t \pi^{\pi_t}(y_t, x) \quad (\text{undirected})$$

$$h_{t+1} = f(h_t, x, y_t, w)$$

$$p(y_t | y_{1:t-1}, x) \propto \exp \left(c(y_t) \tilde{W} h_t \right)$$

$$f(f(\dots, y_{t-2}), x, y_{t-1}, w)$$

$$\underline{h_1} \quad \underline{h_2} \quad \underline{h_t}$$



Learning: maximum likelihood $\log p(y|x) = \sum_t \log p(y_t | y_{1:t-1}, x)$ "teacher forcing"
 ("easy" apart the non-convexity)

decoding: $\arg \max_y \sum_t \log p(y_t | y_{1:t-1}, x)$ → need approximation
 e.g. "beam search" (see Friday)

Pointers

- latent variable SVMstruct:
 - Chun-Nam Yu, Thorsten Joachims. "Learning Structural SVMs with Latent Variables". [ICML 2009](#)
- others:
 - hidden CRF: A. Quattoni, S. Wang, L. Morency, M. Collins, and T. Darrell, "Hidden Conditional Random Fields," [TPAMI 2007](#).
 - deformable part models for object recognition (highly cited paper): Felzenszwalb, R. Girshick, D. McAllester, D. Ramanan, "[Object Detection with Discriminatively Trained Part Based Models](#)" TPAMI, Vol. 32, No. 9, September 2010
- CCCP procedure convergence rate:
 - Ian E.H. Yen, Nanyun Peng, Po-Wei Wang and Shou-de Lin, "On Convergence Rate of Concave-Convex Procedure", [NIPS 2012 OPT Workshop](#) (not considered a publication by the way)
- kernels:
 - example of early paper presenting kernels for structured SVM: Juho Rousu, Craig Saunders, Sandor Szemek, John Shawe-Taylor, "Kernel-Based Learning of Hierarchical Multilabel Classification Models", [JMLR 2006](#)
 - example of application: L. Bertelli, T. Yu, D. Vu, and B. Gokturk, "Kernelized structural SVM learning for supervised object segmentation," [CVPR 2011](#)
 - for computation in BCFW, see Appendix B.5 in the usual: S. Lacoste-Julien, M. Jaggi, M. Schmidt and P. Pletscher, "Block-Coordinate Frank-Wolfe Optimization for Structural SVMs", [ICML 2013](#)
 - the standard book for kernels: Bernhard Schölkopf and Alexander J. Smola, [Learning with Kernels](#), MIT Press 2001
- deep learning:

- see chapter 10 of the "[Deep learning book](#)" for RNNs
- head detection plug in example mentioned in class: Tuan-Hung Vu, Anton Osokin, and Ivan Laptev, "Context-aware CNNs for person head detection", [ICCV 2015](#)