

Lecture 22 - latent SVMstruct

Thursday, April 4, 2019 13:33

today : • latent variable SVMstruct ~ CCCP
• deep learning

prox SAGA

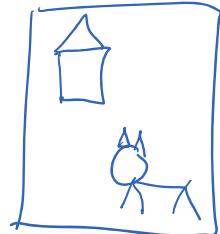
$$\min_w \frac{1}{n} \sum_i f_i(w) + \Omega(w)$$

$$w_{t+1} = \text{Prox}_\gamma(w_t - \gamma [\nabla f_{i_t}(w_t) - g_{i_t}^t + \frac{1}{n} \sum_j g_j^t])$$

→ what is used by default in Scikit-Learn for Lasso

latent variables

motivation: semantic segmentation → find boundary of different objects



Segmentation is expensive → Z "latent variable"
perhaps only have class labels → y

also: [Felzenszwalb et al., TPAMI 2010

"deformable part models" for object recognition

↳ Z there was an object

before, we had $s(x, y; \omega) = \langle \omega, \ell(x, y) \rangle$

configurations

now, consider $s(x, y, z; \omega) = \langle \omega, \ell(x, y, z) \rangle$

as before, could predict with argument $\max_{y \in \mathcal{Y}, z \in \mathcal{Z}} s(x, y, z; \omega)$

* CRF $(p(y|x)) \xrightarrow{\text{generalize}} \text{hidden CRF } p(y, z|x)$

similar to latent variable modeling
with graph. model

ML $\xrightarrow{\text{ML}} \text{expectation-maximization (EM)}$

\hookrightarrow analog for latent SVMstruct is CCCP

Latent SVMstruct

$$l(y, (\tilde{y}, \tilde{z}))$$

generalization of structured hinge loss:

$$l(x, y, \omega) \triangleq \max_{\tilde{y}, \tilde{z}} \langle \omega, \ell(x, \tilde{y}, \tilde{z}) \rangle + \Omega(y, (\tilde{y}, \tilde{z})) - \max_{z \in \mathcal{Z}} \langle \omega, \ell(x, y, z) \rangle \geq l(y, h_\omega(x))$$



here $\Omega(\gamma, y, \omega) = u(\omega) - v(\omega)$ where u, v are convex fcts of ω

"difference of convex functions"

↳ CCCP procedure is used
to approx. minimize this

CCCP procedure

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

→ a majorization-minimization procedure
(EM is another example)

$$g_t(w) = u(w) - [v(w_t) + \langle \nabla v(w_t), w - w_t \rangle]$$

↑
(or subgradient)

$$w_{t+1} = \underset{w}{\operatorname{argmin}} g_t(w)$$

properties of this procedure:

- like EM, descent procedure ie. $g(w_{t+1}) \leq g(w_t)$
$$g(w_t) = g_t(w_t) \geq g_t(w_{t+1}) = g_{t+1}(w_{t+1})$$
- local linear convergence to a stationary point [see NIPS OPT 2012 paper]
for latent SVMstruct

* for SVMstruct

$$v(w) = \max_u \langle w, \varphi(\gamma_i u, z) \rangle$$

$$V(w) = \max_z \langle w, \psi(x, y, z) \rangle$$

$$\partial V(w_t) = \psi(x, y, \hat{z}(x, w_t)) \xrightarrow{\text{argmax } z} \psi(x, y, \hat{z})$$

$$\Rightarrow \mathcal{L}_t(x, y; w) = \max_{(\hat{y}, \hat{z})} \langle w, \psi(x, \hat{y}, \hat{z}) \rangle - \langle w, \psi(x, y, \hat{z}) \rangle + \text{cst.}$$

~ like SUMstruct objective

CCCP for latent SUMstruct : repeat:

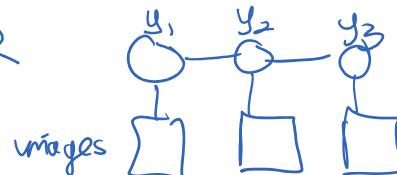
- fill in $\hat{z}_t^{(i)}$ for all ground truth $y^{(i)}$ using w_t
- solve standard SUMstruct to get w_{t+1}
- repeat

Deep learning

go from $\langle w, \psi(x, y) \rangle$ to $\langle w, \psi(x, y, \Theta) \rangle$ can learn

I) plug in "deep learning" features in a structured predict. model

example: OCR



this does: $\psi_t(x_t, y_t) = \begin{pmatrix} 0 \\ x_t \\ 0 \end{pmatrix} \odot y_t$

instead: $\psi_t(x_t, y_t) = \begin{pmatrix} 0 \\ \text{NN}_G(x_t) \\ 0 \end{pmatrix} \odot y_t$

example

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

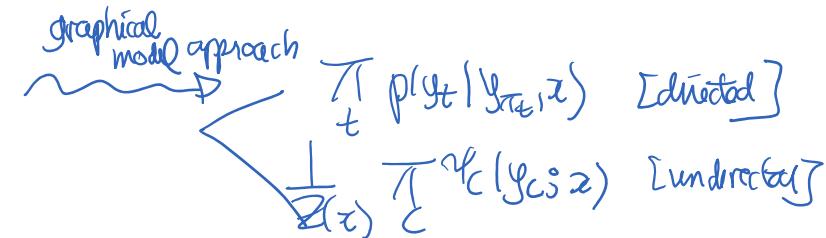
Learned on Imagenet e.g.

instead : $\psi_t(x_t, y_t) = \text{NN}_{\phi}(x_t) \rightarrow y_t$

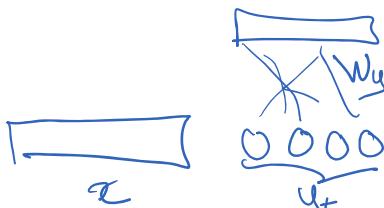
- II) "end-to-end" training : structured prediction energy network (SPEN)
- III) recurrent neural networks (RNN)

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

chain rule



RNN \rightarrow "structured parameterization" of $p(y_t | y_{1:t-1}, x)$ with (in general) no card. indep assumptions

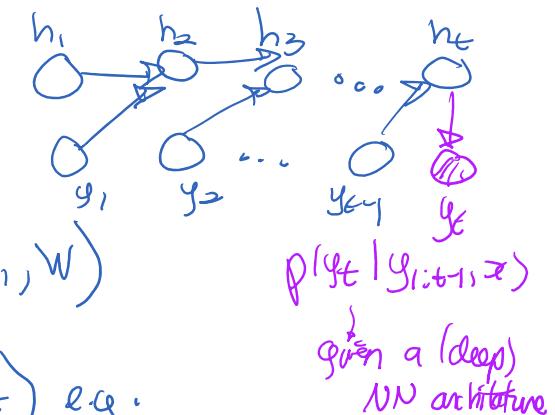


using a NN

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

$$h_t = f(f(f(\dots(y_0), \dots), x, y_{t-1}, w))$$

define $p(y_t | y_{1:t-1}, x) \propto \exp(c(y_t)^T \tilde{W} h_t)$ e.g.



Standard learning : use maximum likelihood

i.e. $\min_{W, \tilde{W}} - \frac{1}{n} \sum_i \log p(y^{(i)} | x^{(i)})$

$\sum \log p(u^{(i)} | v_{1:n-1}^{(i)}, z^{(i)})$

"teacher's forcing"

14h⁴⁰

do SGD on this

$$\sum_i \log p(y^{(i)} | y_{1:t}, x^{(i)})$$

exposure problem

chain rule \rightarrow backpropagation

i.e. do not know
 $p(y | \text{unseen}_y, x)$

decoding: $\underset{y_t}{\operatorname{argmax}} \sum_i \log p(y^{(i)} | y_{1:t}, x) \rightarrow \text{NP hard!}$

\rightsquigarrow need approximation

greedy decoding

$$\hat{y}_t = \underset{y_t \in \Sigma_t}{\operatorname{argmax}} p(y | \hat{y}_{1:t}, x)$$

beam search

"greedy decoding with
memory size K ",
"beam"

beam search: construct y_1, \dots, y_t

beam of size L (memory)

• at step t , you have L candidate solution prefixes $y_{1:t}^{(1)}, y_{1:t}^{(L)}$

• expand possible next choice: $L \cdot |\Sigma_{t+1}|$

• score them (e.g. $p(y_{t+1} | y_{1:t}^{(l)})$)

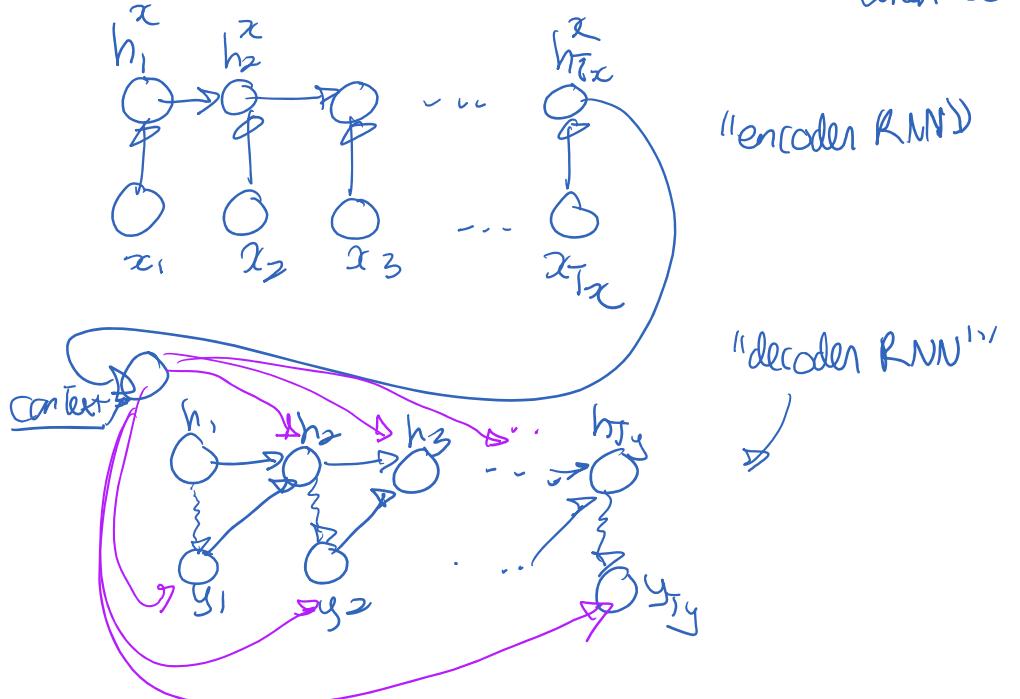
then keep top L candidates as $\{y_{1:t+1}^{(l)}\}_{l=1}^L$

vs. Viterbi alg. which does "backtracking" to correct past mistakes

Seq2seq / encoder/decoder architecture

↳ useful way to get $p(y_t | y_{1:t-1}, x)$ for a RNN

when x can be variable length



issues:

a) Variable length output?
→ end-of-sequence special character

b) Long input sequence?

problem: need to be summarized in fixed length context vector

solution: "attention mechanism"

c) vanishing gradient?

- LSTM
- gated recurrent unit (GRU)
- etc ...