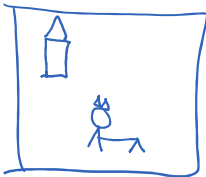


today: • latent variable SVMstruct - CCP  
• deep learning

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 part configurations

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

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

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

⊗ CRF (p(y|x)) <sup>generalize</sup> → hidden CRF p(y, z|x)

similar to latent variable modeling with graph model

ML → EM (expectation-maximization)

↳ analog for latent SVMstruct is CCP

latent SVMstruct

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

generalize structured hinge loss:

$w(w)$ : convex function of w

$J(x, y, w) \triangleq \max_{\tilde{y}, \tilde{z}} \langle w, \phi(x, \tilde{y}, \tilde{z}) \rangle + Q(y, (\tilde{y}, \tilde{z}))$    
 $\rightarrow \max_{z' \in Z} \langle w, \phi(x, y, z') \rangle \geq Q(y, (\tilde{y}, \tilde{z}))$

$$d(x, y, w) = \max_{\tilde{y}, \tilde{z}} \langle w, \varphi(x, y, \tilde{z}) \rangle + \ell(y, (\tilde{y}, \tilde{z})) \quad \underbrace{- \max_{z' \in Z} \langle w, \varphi(x, y, z') \rangle}_{\triangleq v(w)} \geq \ell(y, (y, z))$$

here  $f(x, y, w) = u(w) - v(w)$  where  $u$  &  $v$  are convex fct. of  $w$

"difference of convex functions"

$\Leftrightarrow$  CCCP procedure is to approximate minimize this

CCCP procedure:

- linearize  $v(w)$  at  $w_t$  to get an upper bound
- $w_{t+1}$  is obtained by minimizing this upper bound
- repeat

$\hookrightarrow$  a majorization-minimization procedure (EM is another example)

$$\left[ \begin{array}{l} f_t(w) = u(w) - [v(w_t) + \langle \nabla v(w_t), w - w_t \rangle] \geq f(w) \quad \forall w \\ \text{(or subgradient)} \quad \text{and } f_t(w_t) = f(w_t) \\ w_{t+1} = \arg \min_w f_t(w) \end{array} \right.$$

properties of procedure: • like EM, descent procedure i.e.  $f(w_{t+1}) \leq f(w_t)$

$$f(w_t) = f_t(w_t) \geq f_t(w_{t+1}) \stackrel{\text{upper bound}}{\geq} f(w_{t+1}) //$$

• local linear convergence to a stationary point [see NIPS opt 2012 paper for latent SVM struct]

\* for SVM struct:  $v(w) = \max_z \langle w, \varphi(x, y, z) \rangle \triangleq \arg \max_z \langle w, \varphi(x, y, z) \rangle$

$$\partial v(w_t) = \varphi(x, y, \hat{z}(x, y, w_t))$$

$$\Rightarrow f_t(w) = \max_{(\tilde{y}, \tilde{z})} \langle w, \varphi(x, \tilde{y}, \tilde{z}) \rangle + \ell(y, (\tilde{y}, \tilde{z})) - \langle w, \varphi(x, y, \hat{z}) \rangle + \text{const}$$

$\rightarrow$  like SVM struct objective

CCCP algorithm for latent SVM struct:

repeat:  $\bullet$  fill in  $\hat{z}_t^{(i)}$  for all ground truth  $y^{(i)}$  using  $w_t$

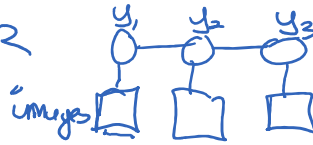
repeat:   
 • fill in  $\hat{z}_t^{(i)}$  for all ground truth  $y_t^{(i)}$  using  $w_t$    
 • solve a standard SVM struct to get  $w_{t+1}$    
 • repeat

Deep Learning

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

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

example: OCR



so you  $\ell_t(x_t, y_t) = \begin{pmatrix} 0 \\ x_t \\ 0 \end{pmatrix} \leftarrow y_t^{th}$

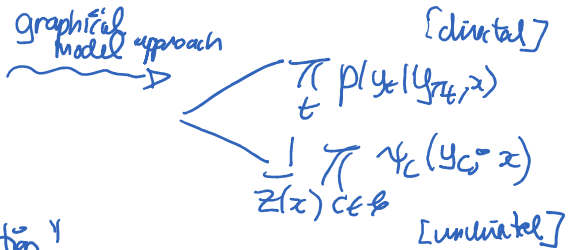
instead  $\ell_t(x_t, y_t) = \begin{pmatrix} 0 \\ NN_{\Theta}(x_t) \\ 0 \end{pmatrix} \leftarrow y_t$

example: [Vu et al. ICCV 2015]   
 "context-aware CNNs for person head detection"   
 learned on images e.g.

II) "end-to-end" training: structured prediction energy networks (SPENs)

III) recurrent neural networks (RNN)

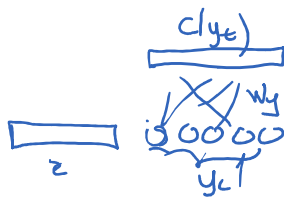
motivation:  $p(y|x) = \prod_{t \leq T} p(y_t | y_{1:t-1}, x)$    
 (chain rule)



RNN  $\rightarrow$  "structured parametrization" of  $p(y_t | y_{1:t-1}, x)$    
 using NN

with no cond. indep. assumption

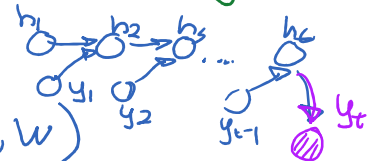
$\hookrightarrow$  usually lose exact drooping



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

$h_t = f(f(\dots f(y_1, \dots, y_t, \dots), x, y_{t-1}, w))$

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



$p(y_t | y_{1:t-1}, x)$    
 given by a (deep) ANN architecture

Standard learning: use M.L.

Standard learning: use M.L.

$$\text{i.e. } \min_{W, \tilde{W}} -\frac{1}{n} \sum_{i=1}^n \underbrace{\log p(y^{(i)} | x^{(i)})}_{\sum_t \log p(y_t^{(i)} | y_{1:t-1}^{(i)}, x^{(i)})}$$

output of a deep NN

...structure  
 "teacher forcing"  
 ↓  
 exposure problem  
 i.e. do not know  
 $P(\hat{y}_t | \text{unseen } x)$   
 prefix  
 $y_{1:t-1}$

14h41

for ML, do SGD objective

gradient of  $\log p(y_t^{(i)} | y_{1:t-1}^{(i)}, x^{(i)}; W, \tilde{W})$   
 → use backpropagation

decoding:  $\text{argmax}_{y \in \mathcal{Y}} \sum_t \log p(y_t | y_{1:t-1}, x)$  → NP hard!

→ need approximation

- greedy decoding  $\hat{y}_t = \text{argmax}_{y_t \in \mathcal{Y}_t} p(y_t | y_{1:t-1}, x)$
- beam search "greedy decoding with memory of size  $k$ "  
 "beam"

beam search: construct  $\hat{y}_1, \dots, \hat{y}_T$

beam of size  $L$  (memory)

- at step  $t$ , you have  $L$  candidate solution prefixes  $\hat{y}_{1:t}^{(1)}$  to  $\hat{y}_{1:t}^{(L)}$

- expand possible next choice:  $|\mathcal{Y}_{t+1}| \cdot L$

- score them (e.g.  $\log p(y_{t+1} | \hat{y}_{1:t}^{(k)}) + \log p(\hat{y}_{1:t}^{(k)})$ )

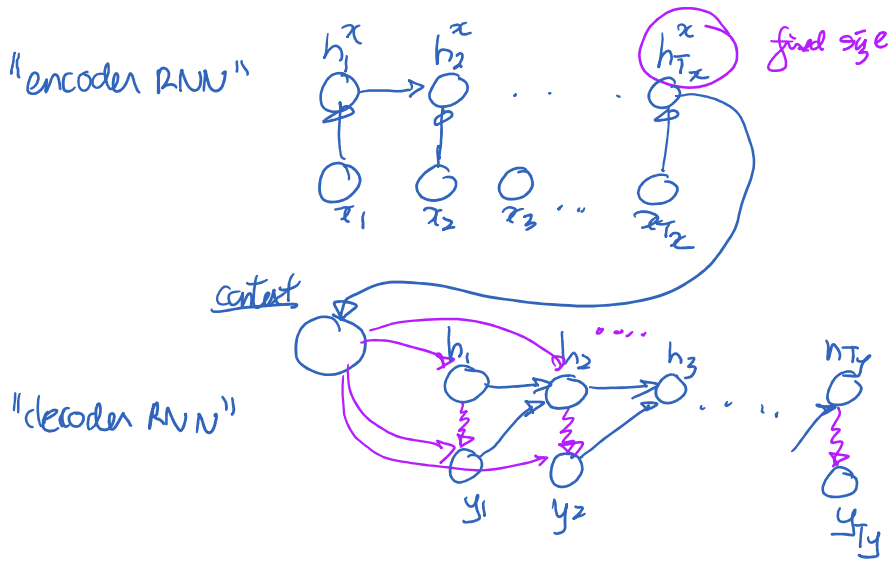
then keep top  $L$  candidates as  $\{\hat{y}_{1:t+1}^{(k)}\}_{k=1}^L$

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

seq2seq data encoder/decoder architecture

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

when  $x$  has variable length



Issues:

a) variable length output?

→ end-of-sequence sequence character

b) long input sequence  $x$ ?

problem: need to summarize input sentence in <sup>context</sup> fixed length

solution: "attention mechanism"

c) vanishing gradient?

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