

Today: • latent variable SVMstruct - CCCP
 • deep learning - RNN

Latent variables

motivation: semantic segmentation



segmentation is expensive $\rightarrow z$ "latent variable"

perhaps only have class labels $\rightarrow y$

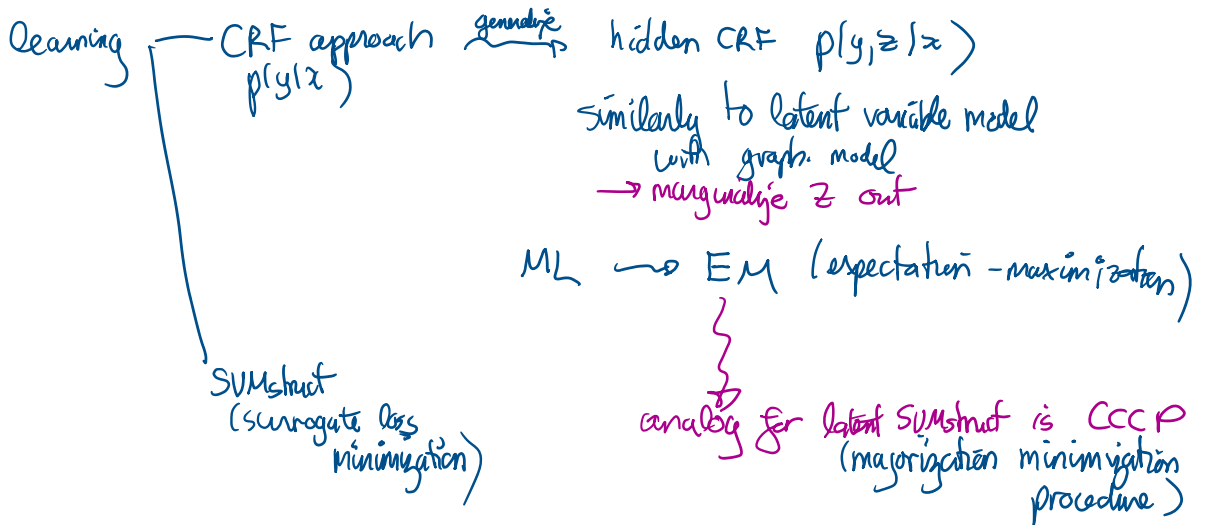
also: [Felzenszwalb & al TPAMI 2010
 "deformable parts models" for object recognition

$\rightarrow z$ there was an object part configuration

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 $\arg \max_{y \in Y, z \in Z} s(x, y, z; w)$



Latent SVMstruct

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

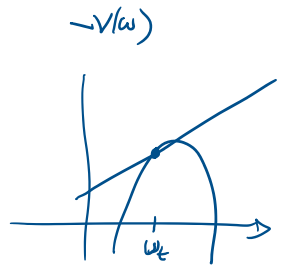
generalized structured hinge loss

$$s(x, y, w) \Leftrightarrow \underbrace{\max_{\tilde{y}, \tilde{z}} \langle w, \phi(x, \tilde{y}, \tilde{z}) \rangle + l(y, (\tilde{y}, \tilde{z}))}_{\hat{=} u(w)} - \underbrace{\max_{z' \in Z} \langle w, \phi(x, y, z') \rangle}_{\hat{=} v(w)} \geq l(y, h_w(x))$$

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

"difference of convex functions"

↳ CCCP procedure is to approx. minimize this



CCCP procedure:

- linearize $v(w)$ at w_t to get an upper bound on $-v(w)$
- w_{t+1} is obtained by minimizing the upperbound $u(w)$ + linearize piece (w)
- repeat \rightarrow a majorization-minimization procedure (EM is another example)

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

properties of procedure:

- like EM, descent procedure i.e. $\mathcal{J}(w_{t+1}) \leq \mathcal{J}(w_t)$
- $\mathcal{J}(w_t) = \mathcal{J}_t(w_t) \geq \mathcal{J}_t(w_{t+1}) \Rightarrow \mathcal{J}(w_{t+1})$
 \uparrow upper bound

• local linear convergence to a station. pt. for latent SVM struct

[see NIPS OPT workshop 2012 paper]

* CCCP for SVM struct:

$$\begin{aligned} v(w) &= \max_{z'} \langle w, \ell(x, y, z') \rangle \rightarrow \equiv \operatorname{argmax}_{z'} \langle w_t, \ell(x, y, z') \rangle \\ \partial v(w_t) &= \ell(x, y, \hat{z}(x, y, w_t)) \\ \Rightarrow \mathcal{J}_t(w) &= \max_{(\tilde{y}, \tilde{z})} \langle w, \ell(x, \tilde{y}, \tilde{z}) \rangle + \ell(y, \tilde{y}, \tilde{z}) - \langle w, \ell(x, y, \hat{z}(x, y, w_t)) \rangle + \text{const.} \\ &\rightarrow \text{like SVM struct objective} \end{aligned}$$

CCCP algorithm for latent SVM struct:

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

10/16

Deep learning

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

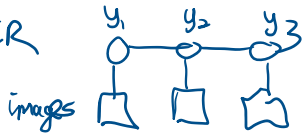
can learn

\rightarrow now for "deep learning" problems in a structured input-output

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

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

example: OCR



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

so for $\psi_t(x_t, y_t) = \begin{pmatrix} 0 \\ x_t \\ 0 \end{pmatrix} \cdot y_t^T$

instead $\psi_t(x_t, y_t) = \begin{pmatrix} 0 \\ \text{NN}(x_t) \\ 0 \end{pmatrix} \cdot y_t^T$

pre-trained on images e.g.

II) "end-to-end" training

structured prediction energy networks (SPEENs)



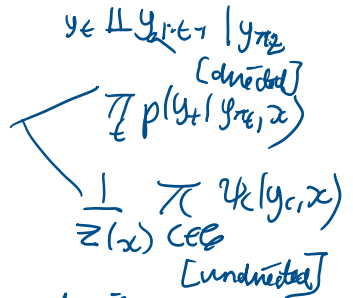
$\psi(x, y, \xi) = \langle w, \psi(x, y, \xi) \rangle$

III) recurrent neural networks (RNN)

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

chain rule

graphical modeling approach



RNN \rightarrow "structured parameterization" with no cond. indep. assumptions

of $p(y_t | y_{1:t-1}, x)$

using a NN

\hookrightarrow usually loose exact decoding

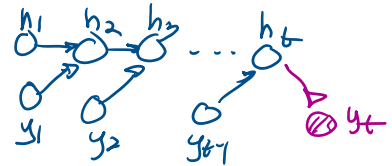
$h_{t+1} \triangleq f(h_t, x, y_t, w)$

$h_t = f(f(\dots f(f(h_1, x, y_1, w), x, y_2) \dots), x, y_{t-1}, w)$

define

$p(y_t | y_{1:t-1}, x) \propto \exp(c(y_t)^T \tilde{w} h_t)$

eg. "code" for y_t



$p(y_t | y_{1:t-1}, x)$

\hookrightarrow given by a deep NN architecture

e.g. word embedding in NMT and can be fine-tuned

standard learning: using MCL

standard training: using MCL

$$\text{f.e. min}_{W, \tilde{W}} \frac{1}{n} \sum_{i=1}^n \log p(y^{(i)} | x^{(i)})$$

$$\sum_t \log p(y_t^{(i)} | y_{1:t-1}^{(i)}, x^{(i)})$$

output of a deep NN

"teacher forcing"

yield
"exposure problem"
i.e. don't know
 $p(y_t | \text{unseen}(y_{1:t-1}))$

for ML, do SGD on obj.

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

→ use back propagation

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

→ need approximation

greedy decoding $\hat{y}_t = \text{argmax}_{y_t \in \mathcal{Y}_t} p(y_t | \hat{y}_1, \dots, \hat{y}_{t-1}, x)$

beam search
"greedy decoding with memory of size K"
↳ size of beam

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

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 $|\mathcal{Y}_{t+1}| \cdot L$

score them (e.g. $\log p(y_{t+1}^{(i)} | y_{1:t}^{(i)}, x) + \log p(y_{1:t}^{(i)}, x)$)

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

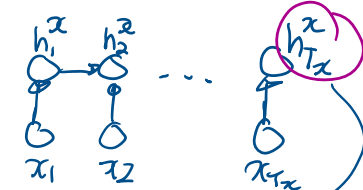
vs.

viterbi alg. which does "backtracking" to correct past mistakes

seq 2 seq a.k.a. encoder/decoder architecture

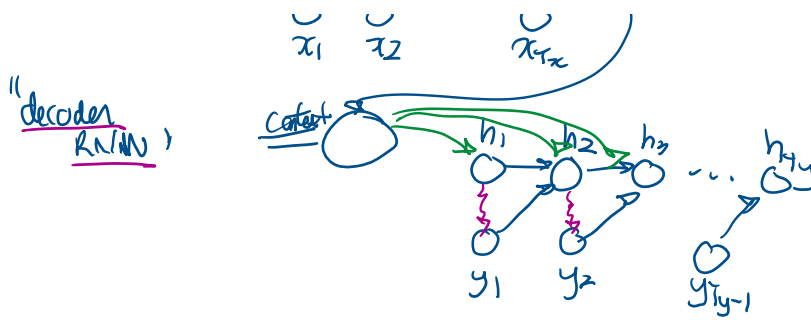
↳ useful way to get $p(y_t | y_{1:t-1}, x)$ for a RNN when x has variable length

"encoder RNN"



→ fixed size representation

"Decoder"



Issues:

a) variable length output? \rightarrow end-of-sequence symbol

b) how to handle long input sequence x ?

problem: need to summarize input sentence in one context vector of fixed length

solution: "attention mechanism"

c) vanishing gradient?

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

d) sequential processing "transformer" architecture

\rightarrow fixed by "transformer" architecture

"attention is all you need"

\rightarrow GPT & al.

"self-attention")

\hookrightarrow quadratic cost
 \Rightarrow need tricks to handle long context / sequence