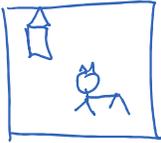


today: • latent variable SVMstruct - CCCP
 • deep learning - RNN

latent variables

motivation: semantic segmentation



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

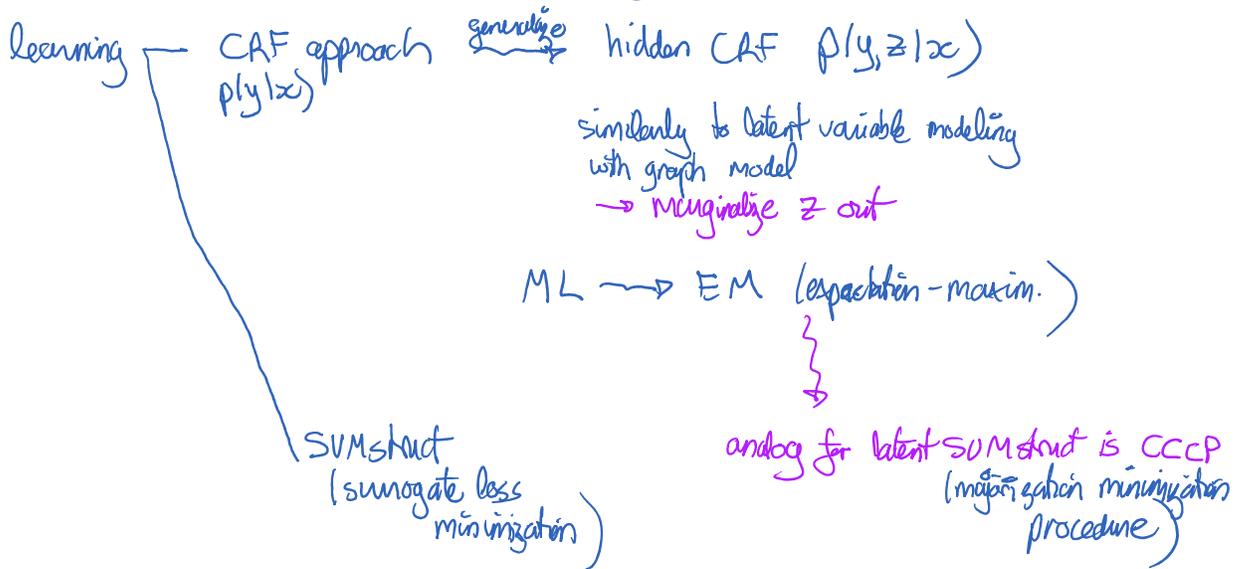
perhaps only have class labels $\rightarrow y$

also: [Felzenszwalb & al IFAAMI 2010]
 "deformable part 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 $\text{argmax}_{y \in \mathcal{Y}, z \in \mathcal{Z}} s(x, y, z; w)$



latent SVMstruct

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

generalize structured hinge loss

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

(best score for ground truth)
(best score for ground truth)

$$\underbrace{\tilde{y}, \tilde{z}}_{\triangleq u(w)} \quad \underbrace{z' \in Z}_{\triangleq v(w)} \quad \left(\begin{matrix} \tilde{y}, \tilde{z} \\ z' \end{matrix} \right)$$

here $f(x, y; w) = u(w) - v(w)$ 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
- w_{t+1} is obtained by minimizing this upper bound
- repeat \rightarrow a majorization-minimization procedure (EM is another example)

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

properties of procedure:

- like EM, descent procedure i.e. $f(w_{t+1}) \leq f(w_t)$
- $f(w_t) = f_t(w_t) \geq f_t(w_{t+1}) \geq f(w_{t+1})$ (upper bound)

- local linear convergence to a stationary pt. for latent SVM struct [see NIPS OPT 2012 paper]

* CCCP for SVM struct:

$$v(w) = \max_{z'} \langle w, \phi(x, y, z') \rangle$$

$$\triangleq \operatorname{argmax}_{z'} \langle w_t, \phi(x, y, z') \rangle$$

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

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

\leadsto like SVM struct objective

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:
- solve in z_t for all ground truth y^* using w_t
 - solve a standard SVM struct to get w_{t+1}
 - repeat

Deep Learning

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

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

example: OCR

so far $\phi_t(x_t, y_t) = \begin{pmatrix} 0 \\ x_t \\ 0 \end{pmatrix} \leftarrow y_t^{\text{th}}$

instead $\phi_t(x_t, y_t) = \begin{pmatrix} 0 \\ \text{NN}_{\theta}(x_t) \\ 0 \end{pmatrix} \leftarrow y_t^{\text{th}}$

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

pre-trained on images e.g.

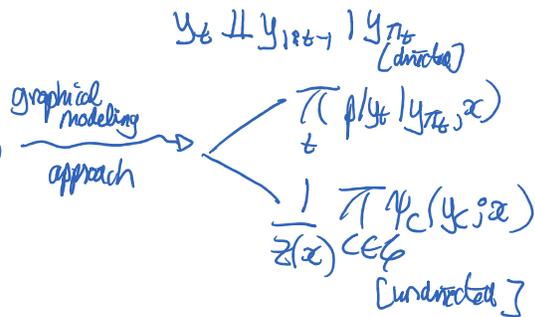
II) "end-to-end" training structured prediction energy networks (SPENs)



III) recurrent neural networks (RNN)

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

chain rule



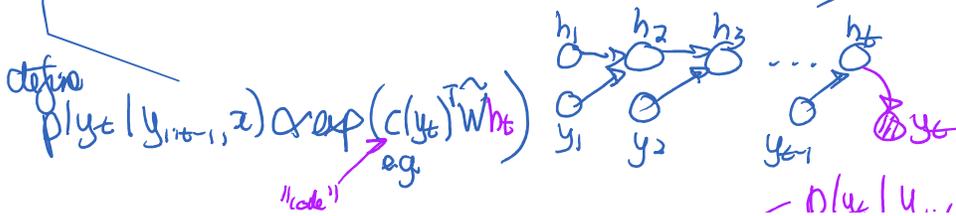
RNN \rightarrow "structured parametrization" of $p(y_t | y_{1:t-1}, x)$ using NN

with no cond. indep. assumptions

\hookrightarrow usually lose exact decoding

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

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



$p(y_t | y_{1:t-1}, x)$ is approx $(C(y_t | w^{(t)}) \bar{y}_1 \bar{y}_2 \dots \bar{y}_{t-1})$
 "code" for y_t
 e.g. word embedding in NMT and can be fine tuned
 given by a deep NN architecture
 $p(y_t | y_{1:t-1}, x)$

Standard Learning: using ML

ie. $\min_{W, \tilde{W}} \frac{1}{n} \sum_{i=1}^n \log p(y^{(i)} | x^{(i)})$

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

"teacher forcing"

↓
 "exposure problem"
 ie. don't know $p(\tilde{y}_t | \text{unseen } x, y_{1:t-1})$

14h43

for ML, do SGD derivative

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

decoding: $\arg \max_{y \in \mathcal{Y}} \sum_t \log p(y_t | y_{1:t-1}, x) \rightarrow$ NP hard?

need approximation

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

• expand possible next choice $|\mathcal{S}_{t+1}| \cdot L$

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

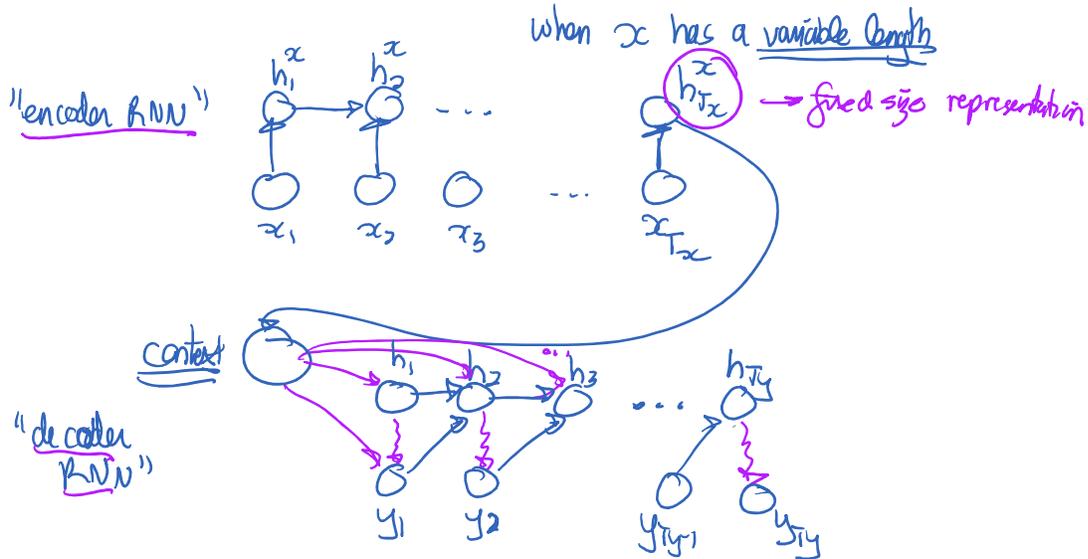
then keep top L candidates as $\sum_{i=1}^L y_{1:t+1}^{(i)}$

vs.

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

Seq2seq a.k.a. encoder/decoder architecture

↳ useful way to get $p(y_{1:T} | x)$ for a RNN



issues:

- variable length output? → end-of-sequence characters
- 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 gradients?

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