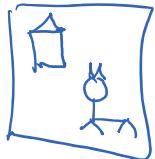


today:

- 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 et al., TPAMI 2010]

"deformable part models" for object recognition
 $\hookrightarrow Z$ 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 $\operatorname{argmax}_{y \in \mathcal{Y}, z \in \mathcal{Z}} s(x, y, z; w) \triangleq h_w(x)$

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

similar to latent variable modeling
 with PGM
 \rightarrow marginalize Z out

ML \rightarrow EM (expectation-maximization)

SVMstruct
 (surrogate loss
 minimization)

analog for latent SVMstruct is CCCP
 (majorization-minimization
 procedure \rightarrow like EM)

Latent SVMstruct

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

generalized structural hinge loss:

$$S(x, y; w) \triangleq \max_{\tilde{y}, \tilde{z}} \langle w, \phi(x, \tilde{y}, \tilde{z}) \rangle + l(y, (\tilde{y}, \tilde{z})) - \max_{z' \in \mathcal{Z}} \langle w, \phi(x, y, z') \rangle \geq l(y, h_w(x))$$



$$f(x, y; \omega) \triangleq \max_{\tilde{y}, \tilde{z}} \langle \omega, \ell(x, \tilde{y}, \tilde{z}) \rangle + l(y, (\tilde{y}, \tilde{z}))$$

$\stackrel{\triangle}{=} u(\omega)$ $\left(\max_{z' \in Z} \langle \omega, \ell(x, y, z') \rangle \right) \geq l(y, h_\omega(x))$
 $\stackrel{\triangle}{=} v(\omega)$ $f(\tilde{y}, \tilde{z})$

here $f(x, y; \omega) = u(\omega) - v(\omega)$ where $u \& v$ are convex fct. of ω

"difference of convex functions" (dc)

\hookrightarrow CCCP procedure is to approx. minimize this

CCCP procedure:

- minimize $v(\omega)$ at w_t to get an upper bound
- w_{t+1} is obtained by minimizing upper bound
- repeat \hookrightarrow a majorization-minimization procedure
(EM is another example)

$$\begin{cases} f_t(\omega) = u(\omega) - [v(w_t) + \langle \nabla v(w_t), \omega - w_t \rangle] \geq f(\omega) \quad \forall \omega \\ w_{t+1} = \underset{\omega}{\operatorname{argmin}} \ f_t(\omega) \end{cases}$$

↑ or subgradient and $f_t(w_t) \approx f(w_t)$

properties of procedure:

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

- local linear convergence to a stationarity pt.
for latent SVMstruct

[see NeurIPS OPT 2012 paper]

* CCCP for SVMstruct :

$$v(\omega) = \max_{\tilde{z}} \langle \omega, \ell(x, y, \tilde{z}) \rangle \stackrel{\triangle}{=} \underset{\tilde{z}}{\operatorname{argmax}} \langle \omega_t, \ell(x, y, \tilde{z}) \rangle$$

$$\nabla v(\omega_t) = \ell(x, y, \hat{z}(x, y; \omega_t))$$

(project idea:
do amortization
with $NN_{\ell(x, y)}$)

$$\Rightarrow f_t(\omega) = \max_{(\tilde{y}, \tilde{z})} \langle \omega, \ell(x, \tilde{y}, \tilde{z}) \rangle + l(y, (\tilde{y}, \tilde{z})) - \langle \omega, \ell(x, y, \hat{z}) \rangle + c.t.$$

\leadsto like SVMstruct obj.

CCCP alg. for latent SVM struct

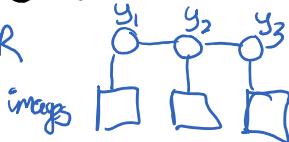
repeat :
 ? fill in $\hat{z}_t^{(i)}$ for all ground truth $y_t^{(i)}$ using w_t
 • solve a standard SVMstruct to get w_{t+1} ,
 • repeat

Deep learning:

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

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

examples: OCR



can learn

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

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

$$\text{instead } \varphi_t(x_t, y_t) = \begin{pmatrix} 0 \\ \text{NN}_t(x_t) \\ 0 \end{pmatrix} \xrightarrow{\text{pre-trained on image! e.g.}} \text{ } y_t^{\text{th}} \text{ position}$$

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

14h28



III) recurrent neural networks (RNN)

$$\text{motivation: } p(y|x) = \prod_{t \leq T} p(y_t | y_{1:t-1}, x) \quad \begin{matrix} \text{chain rule} \\ \text{graphical modelling approach} \end{matrix}$$

$$\prod_t p(y_t | y_{1:t-1}, x) \quad \begin{matrix} \text{Directed} \\ \text{Graphical modelling approach} \end{matrix}$$

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

using NN

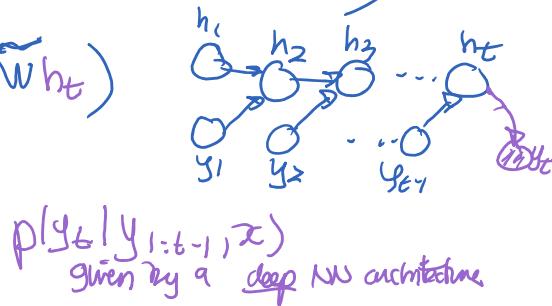
$\xrightarrow{\text{usually loose exact decoding}}$

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

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

$$\text{define } p(y_t | y_{1:t-1}, x) \text{ or } \exp(c \tilde{h}_t^T W h_t)$$

e.g.
word embedding
in LSTM and
can be fine tuned



Standard Learning: Using ML

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

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

output of a deep NN

"teacher forcing"

for ML, do SGD on objective

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

→ use back propagation

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

decoding: augment $\sum_{y \in \Sigma} \log p(y_t | y_{1:t-1}, x) \rightarrow \text{NP hard?}$

→ need approximation

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

beam search

"greedy decoding with
memory if size L"

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 solutions prefers $y_{1:t}^{(1)}$
 $y_{1:t}^{(L)}$

- expand possible next choices $|Y_{t+1}| \cdot L$

score them (e.g. $\log p(y_{t+1} | \hat{y}_{1:t}^{(0)}, x) + \log p(\hat{y}_{1:t}^{(0)}, x)$)

then beam to L candidates as $\hat{y}_{1:t+1}^{(0)} \dots \hat{y}_{1:t+1}^L$.

score $\tau_{t+1} = \log p(y_{t+1}|y_{1:t}, x) + \log p(y_{1:t}, x)$

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

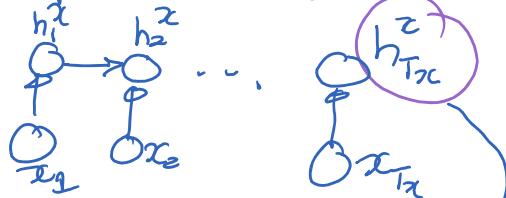
vs.

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

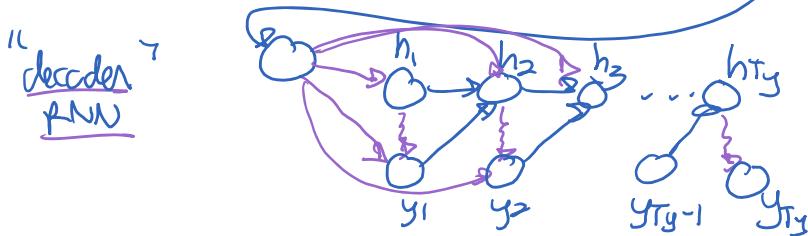
Seq2seq aka encoder/decoder architecture

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

"encoder RNN"



→ fixed size representation



issues:

a) variable length output? → end-of-sequence character

b) how to handle long 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

→ fixed by "transformer architecture"

→ GPT & al.