

Lecture 19 - scribbles - deep learning

Tuesday, March 27, 2018 14:21

- proximal gradient method & Catalyst
- latent SVM struct + CCCP
- deep learning - RNN
- learning to search - EARNN

proximal gradient method

↳ generalization of projected gradient method to other non-smooth functions

Composite framework : $F(w) \triangleq f(w) + \Omega(w)$ where f is convex & L -smooth
 Ω is convex but not nec. smooth

constrained opt. $\Omega(w) = \delta_M(w) \triangleq \begin{cases} 0 & \text{if } w \in M \\ +\infty & \text{otherwise} \end{cases}$
indicator of M

ℓ_1 -regularization $\Omega(w) = \|w\|_1$

proximal gradient update:

$$w_{t+1} = \arg \min_w \underbrace{f(w_t) + \langle \nabla f(w_t), w - w_t \rangle + \frac{1}{2\delta_t} \|w - w_t\|^2}_{\triangleq B_t(w)} + \Omega(w)$$

if $\delta_t \leq \frac{1}{L}$ then $f(w) \leq B_t(w) \forall w$

we can rewrite $B_t(w) = \frac{1}{2\delta_t} \|w - [w_t - \delta_t \nabla f(w_t)]\|^2 + \text{cst.}$ (by completing the square)

\Rightarrow if $\Omega(w) = \delta_M(w)$; we get projected gradient alg.

↓ \Rightarrow if $\Omega(w) = \delta_M(w)$; we get projected gradient alg.

$$w_{t+1} = \text{Prox}_{\gamma_t}^{\Omega} (w_t - \gamma_t \nabla f(w_t))$$

"proximal operator" $\text{prox}_{\gamma}^{\Omega}(z) \triangleq \underset{w}{\text{argmin}} \left\{ \Omega(w) + \frac{1}{2\gamma} \|w - z\|^2 \right\}$

like for projection, prox operator is non-expansive (i.e. 1-Lipschitz)

$$\text{i.e. } \|\text{prox}_{\gamma}^{\Omega}(w) - \text{prox}_{\gamma}^{\Omega}(w')\|_2 \leq \|w - w'\|_2$$

\Rightarrow convergence rates of prox gradient method are same as unconstrained gradient descent

e.g. if f is μ -strongly convex; linear rate with $\rho = \frac{\mu}{\mu + L}$

* to be useful, need $\text{prox}_{\gamma}^{\Omega}$ to be efficiently computable

$$\text{prox}_{\gamma}^{\|w\|_1}(z) = \underset{w}{\text{argmin}} \|w\|_1 + \frac{1}{2\gamma} \|w - z\|_2^2$$

$$\text{"soft-thresholding"} = \begin{cases} \text{sgn}(z) [|z| - \gamma] & \text{if } |z| \geq \gamma \\ 0 & \text{o.w.} \end{cases}$$

FISTA \rightarrow accelerated prox gradient method

Catalyst algorithm

[Lin, Maillard & Harchaoui NIPS 2015]

"meta-algorithm" = outer loop which uses a linearly convergent alg. in the inner loop to get overall acceleration(?)

main idea: use the accelerated proximal point algorithm
with approximation inner loop of prox operator

proximal pt. alg.: is proximal gradient with $f=0$
 $\hookrightarrow w_{t+1} = \text{prox}_{\frac{1}{2\gamma}}^F(w_t)$

Catalyst alg.: (for μ -strongly convex F)

let $q \triangleq \frac{\mu}{\mu+1}$ (γ is an algorithmic parameter)

repeat:

$$w_{t+1} \approx \underset{w}{\text{argmin}} \underbrace{F(w) + \frac{1}{2\gamma} \|w - z_t\|^2}_{\triangleq G_t(w)} \quad \text{s.t.} \quad G_t(w_{t+1}) - \min_w G_t(w) \leq \epsilon_t$$

using an inner loop algorithm [e.g. SAGA or AFW]

$$z_{t+1} = w_{t+1} + \beta_{t+1} (w_{t+1} - w_t)$$

"extrapolation" \nearrow like a "momentum"

[accelerated Nesterov]
trick piece

β_{t+1} is found using fancy equation so that everything works

solve for α_{t+1} in eq. $\alpha_{t+1}^2 = (1 - \alpha_{t+1})\alpha_t^2 + q\alpha_{t+1}$

(pick $\alpha_{t+1} \in]0, 1[$)

$$\beta_{t+1} \triangleq \frac{\alpha_t(1-\alpha_t)}{\alpha_t^2 + \alpha_{t+1}}$$

catalyst trick: use $\gamma \in \mathbb{R}$
s.t. overall # of inner loop calls
give an overall acceleration

$$\tilde{\mu} \approx \mu + \frac{1}{\gamma}$$

with clever analysis of warm starting

acceleration result: if inner loop alg. has convergence $\exp(-\frac{\mu}{L} t)$

strong convexity
of inner loop problem

then with correct constants:

linear rate: $\beta = \frac{1}{k} \leadsto \approx \frac{1}{\sqrt{k}}$ for strongly convex case

$\frac{1}{k} \leadsto \frac{1}{k^2}$ for convex case

result e.g. get accelerated SAGA
" SUBG
" AFW
etc...

Deep learning:

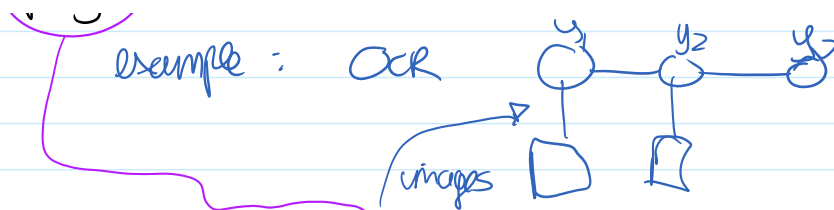
go from $\langle w, \ell(x; y) \rangle$ to $\langle w, \ell(x; y; \Theta) \rangle$

can learn features

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

example: OCR





this loss: $\ell_t(x_t, y_t) = x_t$ → learned on unigram e.g.
 instead $\ell_t(x_t, y_t) = \text{NN}_{\theta}(x_t)$

II) "end-to-end" training: structured energy prediction network (SEPN)

III) recurrent neural networks (RNN)

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

graphical model approach

$\prod_{t \leq T} p(y_t | y_{1:t-1}, x)$ [directed]

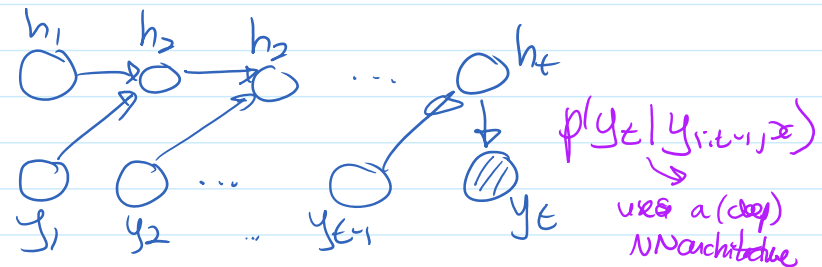
$\frac{1}{Z(x)} \prod_{t \leq T} \psi_t(y_t; x)$ [undirected]

RNN → "structured parameterization" of $p(y_t | y_{1:t-1}, x)$

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

$$h_t = f(f(f(\dots(y_0), \dots), x, y_t, w))$$

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



learning: use maximum likelihood i.e. $\min_{w, \tilde{w}} \sum_{i=1}^n \log p(y^{(i)} | x^{(i)})$

"teacher forcing"

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

exposure problem

do SGD on this
chain rule \rightarrow backpropagation

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

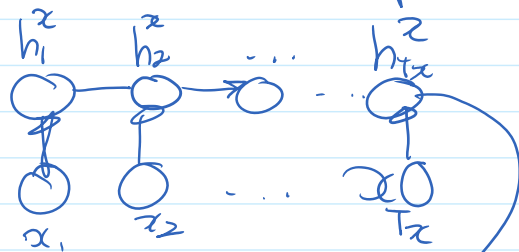
\rightarrow need approximation

greedy decoding $\hat{y}_t = \arg \max_{y \in \mathcal{Y}_t} p(y | \hat{y}_{1:t-1}, x)$

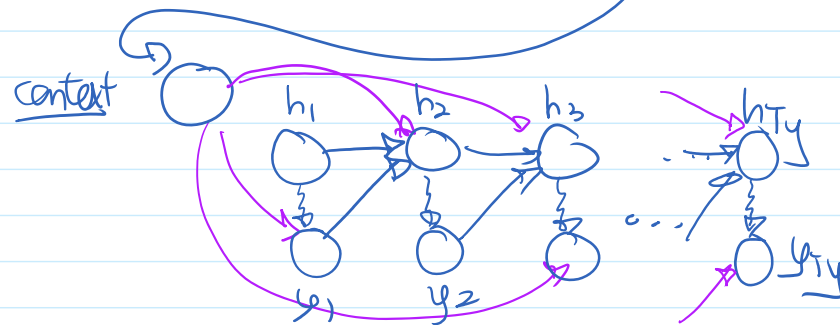
beam search "greedy with memory of size k "
"beam"

Seq2seq / encoder / decoder

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



"encoder RNN"



"decoder RNN"

issues:

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

b) long input sequence x ?

problem: needs to be summarized in fixed length context vector

solution: "attention mechanism"

c) vanishing gradient?

- LSTM
- gated recurrent unit (GRU)

Learning to search (L2S)

$$h_w: X \rightarrow Y$$

$$h_w(x) = \operatorname{argmax}_{y \in Y} s(x, y; w)$$

special case: learning to do greedy search

split y in ordered # of decisions

$$(y_1, \dots, y_T)$$

learn a classifier $\pi(\text{feature}(\hat{y}_1, \dots, \hat{y}_{t-1}, x)) = \hat{y}_t$
↑ classifier "decoding policy"

L2S framework's

from $(x^{(i)}, y^{(i)})_{i=1}^n$ and $l(\cdot, \cdot)$

learn a good classifier/policy π_w s.t. $\hat{y}_t = \pi_w(\hat{y}_1, \dots, \hat{y}_{t-1}, x)$

$h_w(x) \triangleq (\hat{y}_1, \dots, \hat{y}_t)$ "greedy decoder"

s.t. $l(y^{(i)}, h_w(x^{(i)}))$ is good

SEARN Hal Daume's phd thesis

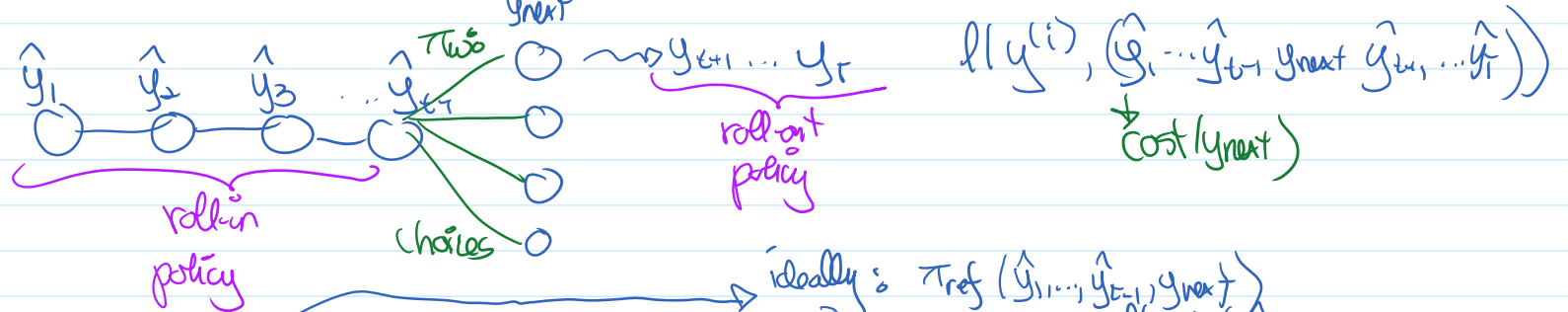
s.t. $l(y^{(i)}, h_w(x^{(i)}))$ is good

→ central idea: "reduction": where reduces structural prediction learning to problem of cost-sensitive classification learning for π_w

method: generate training data for classifier π_w i.e. $(\hat{y}_{\text{context}}^{(i)}, x^{(i)}, \text{cost}(y_{\text{next}}))$
 $\hat{y}_{\text{context}}^{(i)}$ is a prefix sequence

"roll-in" policy → determines how get $\hat{y}_{1:t-1}$ context

"roll-out" policy → " " how get y_{t+1}, \dots, y_T "completion" (target)



ideally: $\pi_{\text{ref}}(\hat{y}_1, \dots, \hat{y}_{t-1}, y_{\text{next}}) = \arg \min_{y_{\text{completion}}} l(y^{(i)}, (y_{1:t-1}, y_{\text{next}}, y_{\text{completion}}))$

e.g. π_{ref} for Hamming loss → copy ground truth

in practice: heuristic approximate it

	roll-in	roll-out	
roll-in	reference	mixture $\frac{1}{2}\pi_{\text{ref}} + \frac{1}{2}\pi_w$	learned π_w
reference policy		← inconsistent →	
learned	consistent not locally optimal	consistency "locally optimal"	R.L.

"teacher forcing" →

LOLS ICML 2015

"locally optimal learning to search"

→ learning better than your teacher

i.e. when heuristic is not optimal

SEARNN: apply LS to RNN training i.e. $\pi_w(y_{1:t-1}, x) \rightarrow$ RNN cell

$p(y_{\text{next}} | y_{1:t-1}, x)$ of a RNN

if use $-\log p(\hat{y}_{\text{target}} | y_{1:t-1}, x)$ as cost "target learning"

and $\hat{y}_{\text{target}} =$ ground truth

then LS with ref roll-in is standard MLE

- * address exposure bias using learned roll-in
- make use of structured loss $\ell(\dots)$ vs. MLE