Density estimation

Computing, and avoiding, partition functions

Roadmap:

- Motivation: density estimation
- Understanding annealing/tempering
- NADE

lain Murray

School of Informatics, University of Edinburgh

Includes work with Ruslan Salakhutdinov and Hugo Larochelle

Probabilistic model ${\cal H}$

Predict new images: $P(\mathbf{x} | \mathcal{H})$





High density can be boring

$P(\mathbf{x} \,|\, \mathcal{H})$





Image reconstruction

Observation model: $P(\mathbf{y} | \mathbf{x})$

Underlying image:

$P(\mathbf{x} \mid \mathbf{y}) \propto P(\mathbf{y} \mid \mathbf{x}) P(\mathbf{x})$

(e.g., Zoran and Weiss, 2011; Lucas Theis's work)

Roadmap

— Unsupervised learning and $P(\mathbf{x} | \mathcal{H})$

- Evaluating $P(\mathbf{x} \mid \mathcal{H})$

Salakhutdinov and Murray (2008) Murray and Salakhutdinov (2009) Wallach, Murray, Salakhutdinov & Mimno (2009)

NADE: "density estimation put first" Larochelle and Murray (2011)

Restricted Boltzmann Machines

$$P(\mathbf{x} | \theta) = \sum_{\mathbf{h}} P(\mathbf{x}, \mathbf{h} | \theta) = \frac{1}{\mathcal{Z}(\theta)} \underbrace{\sum_{\mathbf{h}} \exp[\cdots]}_{f(\mathbf{x}; \theta), \text{ tractable}}$$

\mathcal{Z} a normalizer

Annealing / Tempering

e.g. $P(\mathbf{x};\beta) \propto P^*(\mathbf{x})^\beta \pi(\mathbf{x})^{(1-\beta)}$ $\beta = 0$ $\beta = 0.01$ $\beta = 0.1$ $\beta = 0.25$ $\beta = 0.5$ $\beta = 1$

$1/\beta$ = "temperature"

Annealed Importance Sampling

$$\mathcal{P}(X) = \frac{P^*(\mathbf{x}_K)}{\mathcal{Z}} \prod_{k=1}^K \widetilde{T}_k(\mathbf{x}_{k-1}; \mathbf{x}_k), \qquad \mathcal{Q}(X) = \pi(\mathbf{x}_0) \prod_{k=1}^K T_k(\mathbf{x}_k; \mathbf{x}_{k-1})$$

Standard importance sampling of $\mathcal{P}(X) = \frac{\mathcal{P}^{*}(X)}{\mathcal{Z}}$ with $\mathcal{Q}(X)$

Annealed Importance Sampling

Parallel tempering

Standard MCMC: Transitions + swap proposals on joint:

$$P(X) = \prod_{\beta} P(X;\beta)$$

- larger system
- \bullet information from low β diffuses up by slow random walk

Tempered transitions

Proposal: swap order, final point \check{x}_0 putatively $\sim P(x)$

Acceptance probability:

$$\min\left[1, \ \frac{P_{\beta_1}(\hat{x}_0)}{P(\hat{x}_0)} \cdots \frac{P_{\beta_K}(\hat{x}_{K-1})}{P_{\beta_{K-1}}(\hat{x}_0)} \frac{P_{\beta_{K-1}}(\check{x}_{K-1})}{P_{\beta_K}(\check{x}_{K-1})} \cdots \frac{P(\check{x}_0)}{P_{\beta_1}(\check{x}_0)}\right]$$

Whirlwind tour of annealing / tempering

Must be able to get anywhere in distribution

Methods to use generally for hardest problems.

An experiment

Take 60,000 binarized MNIST digits, like these:

- Train an RBM using CD (and then find \mathcal{Z})
- Train a mixture of multivariate Bernoullis with EM

Compare samples and test-set log-probs

A comparison

Samples from:

- mixture of Bernoullis, -143 nats/test digit
- \bullet RBM, -106 nats/test digit

Which is which?

A better fitted RBM

RBM samples Training set examples

Test log-prob now 20 nats better (-86 nats/digit)

Dependent latent variables

"Deep Belief Net"

Lateral connections

Directed model

 $P(\mathbf{x}) = \frac{1}{\mathcal{Z}} \sum_{\mathbf{h}} P^*(\mathbf{x}, \mathbf{h}), \text{ not available}$

Chib-style estimates

Bayes Rule:

$$P(\mathbf{h} \,|\, \mathbf{x}) = \frac{P(\mathbf{h}, \mathbf{x})}{P(\mathbf{x})}$$

For any special state h^* :

$$P(\mathbf{x}) = \frac{P(\mathbf{h}^*, \mathbf{x})}{P(\mathbf{h}^* | \mathbf{x})} \leftarrow \mathsf{Estimate}$$

Murray and Salakhutdinov (2009)

Variational approach

$$\log P(\mathbf{x}) = \log \sum_{\mathbf{h}} \frac{1}{\mathcal{Z}} P^*(\mathbf{x}, \mathbf{h})$$

$$\geq \sum_{\mathbf{h}} Q(\mathbf{h}) \log P^*(\mathbf{x}, \mathbf{h}) - \log \mathcal{Z} + \mathcal{H}[Q(\mathbf{h})]$$

Results MNIST

Results Natural Scenes

$P(\mathbf{x} \,|\, \mathcal{H})$ taught me

— RBM: state-of-the-art for binary dists

— Deep nets only very slightly better on MNIST

— Some Gaussian RBMs are really bad... ...and going deep won't help

— Most topic model $P(\mathbf{x} \mid \mathcal{H})$ ests wrong

Roadmap

— Unsupervised learning and $P(\mathbf{x} | \mathcal{H})$

Evaluating $P(\mathbf{x} | \mathcal{H})$ Salakhutdinov and Murray (2008) Murray and Salakhutdinov (2009) Wallach, Murray, Salakhutdinov & Mimno (2009)

NADE: "density estimation put first"

Larochelle and Murray (2011)

Decompose into scalars

$P(\mathbf{x}) = P(x_1) P(x_2 | x_1) P(x_3 | x_1, x_2) \dots$

 $= P(x_k | \mathbf{x}_{< k})$ k

FVSBN: Fully Visible Sigmoid Belief Net

Logistic regression for conditionals

FVSBNs beat mixtures, but not RBMs

Approximate RBM

$P(x_k | \mathbf{x}_{< k})$ from MCMC

or mean field

(Requires fitted RBM. Creates new model.)

One Mean Field step

 $\hat{x}_k = \sigma \left(b_k^x + W_{k,.} \mathbf{h}^{(k)} \right)$ $\mathbf{h}^{(k)} = \sigma \left(\mathbf{b}^h + W_{..< k}^\top \mathbf{x}_{< k} \right)$

NADE Neural Autoregressive Distribution Estimator

Fit as new model

$P(\mathbf{x} | \mathcal{H}) = \prod_k \hat{p}(x_k)$

Tractable, $\mathcal{O}(DH)$

NADE results

Model	ADULT	CONNECT-4	DNA	MUSHROOMS	NIPS-0-12	OCR-LETTERS	RCV1	WEB
MoB	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00
	± 0.10	± 0.04	± 0.53	± 0.10	± 1.12	± 0.32	± 0.11	± 0.23
RBM	4.18	0.75	1.29	-0.69	12.65	-2.49	-1.29	0.78
	± 0.06	± 0.02	± 0.48	± 0.09	± 1.07	± 0.30	± 0.11	± 0.20
RBM	4.15	-1.72	1.45	-0.69	11.25	0.99	-0.04	0.02
mult.	± 0.06	± 0.03	± 0.40	± 0.05	± 1.06	± 0.29	± 0.11	± 0.21
RBForest	4.12	0.59	1.39	0.04	12.61	3.78	0.56	-0.15
	± 0.06	± 0.02	± 0.49	± 0.07	± 1.07	± 0.28	± 0.11	± 0.21
FVSBN	7.27	11.02	14.55	4.19	13.14	1.26	-2.24	0.81
	\pm 0.04	± 0.01	\pm 0.50	± 0.05	± 0.98	± 0.23	± 0.11	± 0.20
NADE	7.25	11.42	13.38	4.65	16.94	13.34	0.93	1.77
	\pm 0.05	\pm 0.01	± 0.57	\pm 0.04	\pm 1.11	\pm 0.21	\pm 0.11	\pm 0.20
Normalization	-20.44	-23.41	-98.19	-14.46	-290.02	-40.56	-47.59	-30.16

★ Little variation when changing input ordering: DNA = +/- 0.05 MUSHROOMS = +/- 0.045 NIPS-0-12 = +/- 0.15

	Model	Log. Like.	7	7	3	į		З,	2		Ø.	2	7
Intractable	*	0	- C	5		(VI	6	Z	1	Ł	8	4	0
	MoB [*]	-137.64	Å.	E	9	3	ÿ	*4 2 T	1	S	3	6	4
	RBM (CDI)*	\approx -125.53		1	Ç	2	Ģ	ι,	t, the	7	4	0	
		\approx -105.50	ß	2	q	5	5	12	ij	3 .	3	17	F
	RBM (CD3)*		47		Č,	1.5		9	7	4	4	1	
		~ 96.21				~		1	<u> </u>	<u> </u>	1	<u>ار</u> بور.	
		\approx -80.34	6	4	**>	¥.	2	\checkmark	مقعا	/	5	-7	é
	FVSBN	-97.45	5		1		1	2	5	×	2	6	6
				B		$ _{L_{1}}$	0	}	6			60	
	NADE	-88.86		2	6	Ċ,	4.1	2	a	G,	G	A	$\dot{\circ}$
										4		<u> </u>	$\mathbf{\nabla}$

When should we learn $P(\mathbf{x} | \mathcal{H})$?

Monte Carlo methods

Autoregressive models (see also Lucas Theis)

A longer talk on NADE: http://videolectures.net/aistats2011_larochelle_neural/

Appendix slides

Markov chain estimation

Stationary condition for Markov chain:

$$P(\mathbf{h}^*|\mathbf{x}) = \sum_{\mathbf{h}} T(\mathbf{h}^* \leftarrow \mathbf{h}) P(\mathbf{h}|\mathbf{x})$$
$$\approx \left[\frac{1}{S} \sum_{s=1}^{S} T(\mathbf{h}^* \leftarrow \mathbf{h}^{(s)}), \quad \mathbf{h}^{(s)} \sim \mathcal{P}(H) \right] = \hat{p}$$

 $\mathcal{P}(H)$ draws a sequence from an equilibrium Markov chain:

$$\underbrace{\mathbf{h}^{(1)}}_{T} \xrightarrow{\mathbf{h}^{(2)}}_{T} \underbrace{\mathbf{h}^{(3)}}_{T} \underbrace{\mathbf{h}^{(4)}}_{T} \underbrace{\mathbf{h}^{(5)}}_{T} \underbrace{\mathbf{h}^{(5)}}_{T}$$

Bias in answer

$$P(\mathbf{x}) = \frac{P(\mathbf{h}^*, \mathbf{x})}{P(\mathbf{h}^* | \mathbf{x})} = \frac{P(\mathbf{h}^*, \mathbf{x})}{\mathbb{E}[\hat{p}]} \le \mathbb{E}\left[\frac{P(\mathbf{h}^*, \mathbf{x})}{\hat{p}}\right]$$

Idea: bias Markov chain by starting at \mathbf{h}^*

$$rac{1}{S}\sum_{s=1}^{S}T(\mathbf{h}^*\!\leftarrow\!\mathbf{h}^{(s)})$$
 will often overestimate $P(\mathbf{h}^*|\mathbf{x})$

New estimator

We actually need a slightly more complicated \mathcal{Q} :

$$\mathbf{h}^{(1)} \leftarrow \widetilde{T} \qquad \mathbf{h}^{(2)} \leftarrow \widetilde{T} \qquad \mathbf{h}^{(3)} \leftarrow \mathbf{h}^{(4)} \qquad T \leftarrow \mathbf{h}^{(5)}$$

$$\widetilde{T} \qquad \widetilde{T} \qquad \widetilde{T} \qquad \widetilde{T} \qquad \mathbf{h}^{(5)}$$

$$\widetilde{T} \qquad \widetilde{T} \qquad \mathbf{h}^{(5)} \qquad \widetilde{T} \qquad \mathbf{h}^{(5)}$$

$$\widetilde{T} \qquad \mathbf{h}^{(5)} \qquad \mathbf{h}^{($$

$$\mathbb{E}_{\mathcal{Q}(H)}\left[1\left/\frac{1}{S}\sum_{s=1}^{S}T(\mathbf{h}^{*}\leftarrow\mathbf{h}^{(s)})\right]=\frac{1}{P(\mathbf{h}^{*}|\mathbf{x})}$$

 $\hat{P}(\mathbf{x}) = \frac{P(\mathbf{x}, \mathbf{h}^*)}{\hat{P}(\mathbf{h}^* | \mathbf{x})} \text{ unbiased } \Rightarrow \text{ stochastic lower bound on } \log P(\mathbf{x})$