

Variational Inference for HLP

29.11.2018.

• github

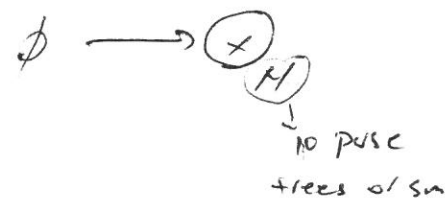
- learn a distribution over observed and unobserved data
- supervised - with known probability (mass/density)

$$X \sim \text{Cat}(\pi_1, \dots, \pi_k) \rightarrow k \text{ genders } k \text{ classes}$$

or $X \sim \mathcal{N}(\mu, \sigma^2) \rightarrow \text{height of population}$

estimate params that assign
maximum likelihood

- for every problem we have side info & observation
a sentence ϕ its translation x



- have HMs predict parameters of our probabilistic model and proceed to estimate params θ of the HM

$$x|\phi \sim \text{Cat}(\pi_{\theta}(\phi)) \quad \text{or} \quad x|\phi \sim \mathcal{N}(\mu_{\theta}(\phi), \sigma_{\theta}(\phi)^2)$$

- HMs - as task-driven feature extraction

- they don't generate the data (they don't output sentiment), but parametrise distributions that by assumption govern data (they output probability per sentiment labels). \rightarrow prediction is done by us, not learning & statistics or math.

$$\mathcal{L}(\theta|x^{(1:H)}) = \log \prod_{s=1}^H P(x^{(s)}|\theta) = \sum_{s=1}^H \log P(x^{(s)}|\theta)$$

\rightarrow neg-log-likelihood fn gives us something to chase

\rightarrow given the data, we see this params
 \rightarrow gradient-based optimisation

- for large N , computing the gradient is inconvenient

$$\sum_{s=1}^N \left(\frac{1}{N} \right) H \nabla_{\theta} \log p(x^{(s)} | \theta)$$

uniform distribution
 H

$$= \sum_{s=1}^N U(s | 1/N) H \nabla_{\theta} \log$$

$$= \mathbb{E}_{S \sim U(1/N)} [H \nabla_{\theta} \log p(x^{(s)} | \theta)]$$

S selects data point uniformly at random
 → expected gradient → don't compute \mathbb{E} N times,
 but use Monte Carlo to get a sample & get
 fewer N s

$$\stackrel{MC}{\approx} \frac{1}{M} \sum_{m=1}^M H \nabla_{\theta} \log$$

$$S_i \sim U(1/N)$$

get unbiased gradient estimation.

• DL in HCP

- MLE → which loss to optimise (eg. neg log-likelihood)
- automatic differentiation (backprop)
- stochastic opt powered by backprop

• when do we have intractable likelihood?

→ latent variables

$$p(x | \theta) = \sum_{c=1}^K \text{Cat}(c | \pi_1, \dots, \pi_K) N(x | \mu_{\theta}(c), \Sigma_{\theta}(c)^2)$$

inconvenient: discrete latent variable, continuous observations

SLIDES ▽ impossible: continuous - - - , discrete - - -

• exact gradient is intractable

$$\frac{\text{joint}}{\text{marginal}} = \text{conditional} \quad (18/23) \quad \downarrow \text{do math}$$

get expected gradient
?

do ~~me~~ sampling? we don't know
the posterior $p(z|x, \theta) \rightarrow$ hard to
estimate. How?

• Why latent variable modeling then??

- better handle on stat. assumptions
- organise a massive collection of data
- learn from unlabelled data, semi supervised learning
- induce discrete rep.
- uncertainty quantification

• Deep Gen. Models \rightarrow probabilistic models parametrised by HM

II. Variational Inference: the Basics

• Generative Model

- Joint Distribution of x and z random variables

\downarrow
observed

\downarrow
latent
unobserved

$$\begin{aligned} p(x, z) &= p(x)p(z|x) \\ &\approx p(z)p(x|z) \\ \text{assumption} &= p(x) \cdot p(z) \end{aligned}$$

- $p(x|z)$ is the likelihood
- $p(z)$ is the prior over z

$$\underbrace{p(z|x)}_{\text{posterior}} = \frac{p(x|z) \cdot p(z)}{\underbrace{p(x)}_{\text{marginal likelihood / evidence}}}$$

- We want to compute the posterior over latent variables $p(z|x)$. Must compute

$$p(x) = \int p(x, z) dz$$

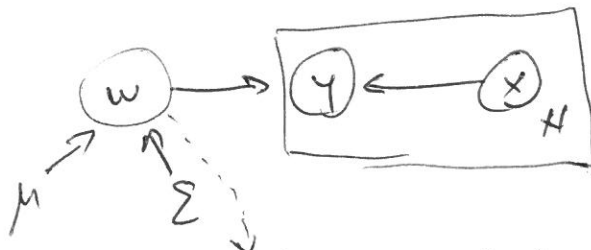
which is often intractable \Rightarrow approximate inference

accept that we can't solve this.

- Can't compute the posterior when:

(1) functional form of the posterior is unknown

- bayesian log-linear POS tagger



don't know distr. \rightarrow assume it's Gaussian.

(2) Very hard computation

Factorial HMMs

- complexity of inference $O(L^{\#chains} \cdot T)$

EDA - assume posterior is
categorical distr. (per word)
Dirichlet distr. (for docs)

\rightarrow Rule of Thumb \rightarrow assume the posterior is in the same family as the prior.

Variational Inference

$p(z|x)$ is not computable

\Rightarrow approximate by auxiliary distr. $q(z)$

Requirement: choose $q(z)$ as close as possible to $p(z|x)$

KL divergence

$$KL(q(z) \parallel p(z|x)) = \mathbb{E}_{q(z)} \left[\log \left(\frac{q(z)}{p(z|x)} \right) \right]$$

assume support of q
is included in support
of p .

continuous $\int q(z) \log(\dots) dz$

discrete $\sum_z q(z) \log$

- properties: $KL \geq 0$ (we want it 0)

$$-KL \leq 0$$

VI derivation I.

$$\log p(x) = \log \int p(x, z) dz = \dots = \log \left(\mathbb{E}_{q(z)} \left[\frac{p(x, z)}{q(z)} \right] \right) \geq \underbrace{\mathbb{E}_{q(z)} [\log(1)]}_{\text{negative}}$$

Found a lower bound for
log-likelihood:

$$\begin{aligned} \log p(x) &\geq \mathbb{E}_{q(z)} \left[\log \left(\frac{p(z|x) p(x)}{q(z)} \right) \right] \\ &= \int q(z) \log \left(\frac{p(z|x) p(x)}{q(z)} \right) dz + \log p(x) \\ &= -KL(q(z) \parallel p(z|x)) + \log p(x) \end{aligned}$$

The ~~gap~~ introduced gap
by the assumption between
the actual log-likelihood
and approx. is exactly



• VI derivation II.

- minimize KL

$$\min_{q(z)} KL = \max_{q(z)} -KL$$

$$\max_{q(z)} \int q(z) \log \frac{p(z|x)}{q(z)} dz = \dots = \max_{q(z)} \mathbb{E}_{q(z)} [\log(p(x, z))] + H(q(z))$$

key: ignore $\log p(x)$ because it's constant

ELBO
evidence lower bound

fit the data

regularizer
basically ∇

so we don't go far from 0.

• Performing VI:

(1) Maximize regularized expected log-density

(2) Optimize generative model:

$$\max_{p(x, z)} \mathbb{E}_{q(z)} [\log(p(x, z))] + H(q(z))$$

constant, we found it

- Do (1) & (2) iteratively

- EM Algorithm - fix one thing

• Mean Field Inference

- how to choose q ? From a good family, must be tractable.

- make all latent variables independent under $q(z)$

- approx. posterior $q(s, z) = \prod_{t=1}^T q(s_t) q(z_t)$

poor approx. bcz we assume everything is independent.

Make same latent variables independent

• posterior inference is intractable coz marginal likelihood $p(x)$

\Rightarrow VI approx. (start with mean field approx)

David Blei

III. Deep Generative Models

$$p(x|z, \theta) = p(z) p(x|z, \theta)$$

↓
NN params

marginal likelihood $p(x|\theta)$ is intractable.

• Wake-sleep Algorithm

- generalize latent variables to NN

(1) a generation network to model the data (θ)

(2) an inference (recognition) network to model the latent variable (λ param.)

- binary hidden units, "hard EM" fashion

(1) Wake phase

- update gen. params θ

(2) Sleep phase

- update λ , but can't get proper gradient \rightarrow must use fictional data \rightarrow we can go in the wrong direction

- its supervised learning with made up labels (by inference n.)

- pros: - conditionally independent layer-wise updates
- amortised inference

cons: - trained on different objectives

- λ are updated on fake data \tilde{x}

- just gets worse \therefore

• Turn to VI

$$\log p(x|\theta) \geq \text{ELBO} \quad \# \text{math on slides}$$

$$\arg \max_{\theta, \lambda} \underbrace{E_{q(z|x, \lambda)} [\log p(x|z, \theta)]}_{\text{approx. by sampling}} - \underbrace{KL(q(z|x, \lambda) \| p(z))}_{\text{easy, analytical true}}$$

Generator Network ~~Gener~~ Gradient

$$\frac{\partial}{\partial \theta} \{ \mathbb{E}_z [\log p(x|z, \theta)] - \text{KL} \}$$

constant

$$= \mathbb{E}_{q(z|x, \lambda)} \left[\frac{\partial}{\partial \theta} \log p(x|z, \theta) \right]$$

$$\stackrel{MC}{\approx} \frac{1}{S} \sum_{i=1}^S \frac{\partial}{\partial \theta}$$

we can do it easily with backprop

Inference Network

- KL part is analytical computation
- focus again on:

$$\frac{\partial}{\partial \lambda} \mathbb{E}_{q(z|x, \lambda)} [\log p(x|z, \theta)]$$

→ gradient is stuck, can't compute ⇒ reparameterisation trick

- $h: z \rightarrow \epsilon$,
 ϵ doesn't depend on λ .

I'm lost

Gaussian Transformation for

- we can write any ~~standard~~ Gaussian with $\mathcal{N}(0, 1)$

$$h(z, \lambda) = \frac{z - \mu(\lambda)}{\sigma(\lambda)} = \epsilon \sim \mathcal{N}(0, 1)$$

replace with standard Gaussian so $\frac{\partial}{\partial \lambda}$ can go in:

$$\int \mathcal{N}(z | \mu, \sigma^2) \log p(x|z) dz$$

$$\Downarrow$$

$$\int \mathcal{N}(\epsilon | 0, 1) \log p(x | h^{-1}(\epsilon, \lambda), \theta) d\epsilon$$

standard Gauss.
shift

- now we ~~have~~ can sample

• Unigram Document Model

- slides

$$h = \text{relu}(w_1 z + b_1)$$

$$f(z, \theta) = \text{softmax}(w_2 h + b_2)$$

dim. of vocabulary (probabilities)

$$\theta = \{w_1, b_1, w_2, b_2\}$$

Gen.
net

Inference model:

average emb.
of doc.

$$s = \sum_{i=1}^N E_{x_i}$$

$$h = \text{relu}(M_1 s + c_1)$$

$$\mu(x_i^N, \lambda) = M_2 h + c_2$$

$$\sigma(x_i^N, \lambda) = \text{softplus}(M_3 h + c_3)$$

$$\lambda = \{E, M_1, M_2, M_3, c_1, c_2, c_3\}$$

- VAE trains both networks with the same objective
→ solves wake-sleep paper

IV. DGM: Discrete Latent Variable

- reparametrisation gradient

$$\dots - \frac{\partial}{\partial \theta} \text{KL}(\underbrace{q(z|\lambda) \| p(z|\theta)})$$

depends on θ , not constant anymore.

However, it's easy and stuff exist that solves it easily.

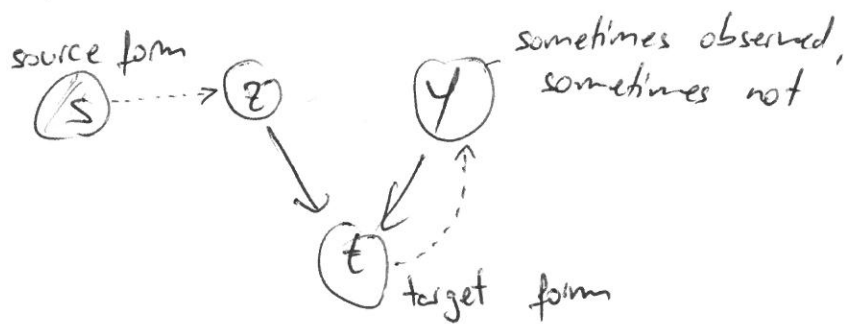
- Gotta do more manipulation to do it for discrete variables
(in latent space)
depends on the problem
you're working on

- compute the Jacobian to transform the variables
→ transformation will expand/shrink the space, must control with $|J|$

- Example - combine cont. & disc. latent variables
→ partially observed data

- Morphological Reinflection task

- plays → played
- z - lemma (real vector)
- y - morph. information - gender of a word etc. (discrete vector)



simple posterior approximation

What we know?

IV.

- DGMs
- objective - lower bound on log-likelihood \rightarrow can't be computed exactly \rightarrow MC
- MC is not differentiable \rightarrow score function estimator \rightarrow reparameterised gradient
- Gamma distribution - distr. of Gaussian σ^2
- Dirichlet - distr. over categorical distributions
- retractable ex. If we have Gaussian \rightarrow reparametrisation
- ADVI - Automatic Differentiation variational inference

VI. Normalising Flows

- map transform complex distributions to simpler ones (NFs)
- we need a more flexible environment

\rightarrow not just for known distributions ~~we know~~
ones we know

the reparametrisation for

\downarrow
express the density of
variable Y in terms of a var X .
Assume that a differentiable, invertible
mapping $h: X \rightarrow Y$ exists.

- \Rightarrow Let's ~~learn~~ the transformation h (or its inverse)
learn

- our data x has unknown cont. density $p(x)$. (art handcraft
likelihood \rightarrow word embeddings)