

Deep Bayes (HSE +)

$$\text{conditional} = \frac{\text{joint}}{\text{marginal}} = \frac{p(x, y)}{p(y)}$$

$$p(x, y, z) = p(x|y, z) p(y|z) p(z)$$

$$p(y) = \int p(x, y) dx \quad \hookrightarrow \text{one dim cond. distn}$$

estimate $p(y|x)$ from $p(x, y, z)$ w/o knowing z .

$$p(y|x) = \frac{p(y, x)}{p(x)} = \frac{\int p(x, y, z) dz}{\int p(x, y, z) dy dz} \quad (\text{In theory})$$

Bayes theorem

$$p(y|x) = \frac{p(x, y)}{p(x)} = \frac{p(x|y)p(y)}{p(x)} = \frac{p(x|y)p(y)}{\int p(x|y)p(y)dy}$$

$$\text{Posterior} = \frac{\text{Likelihood} \times \text{Prior}}{\text{Evidence}}$$

Statistical inference

data $x = (x_1, \dots, x_n) \Leftrightarrow p(x|\theta)$, θ ?

2 methods

Frequentist / Classical :

MLE.

$$\begin{aligned}\hat{\theta}_{ML} &= \arg \max \rho(x|\theta) = \arg \max \prod_{i=1}^n p(x_i|\theta) \\ &= \arg \max \sum_i \log p(x_i|\theta)\end{aligned}$$

Bayesian:

encode uncertainty about θ in a prior $p(\theta)$

then

$$p(\theta|x) = \frac{\prod_{i=1}^n p(x_i|\theta) p(\theta)}{\int \prod_{i=1}^n p(x_i|\theta) p(\theta) d\theta}$$

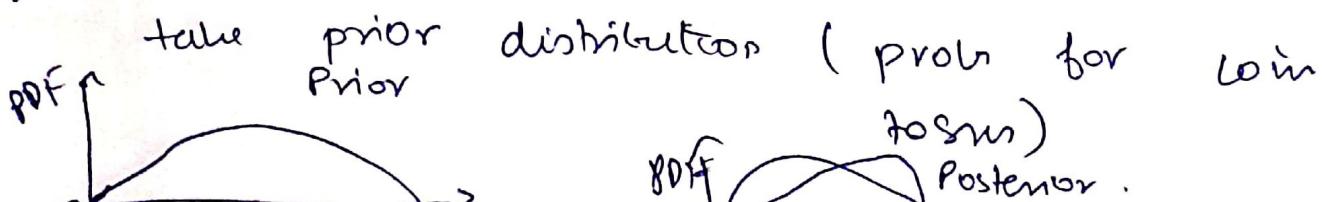
Coin tossing ex:

2 tosses. Outcome H. H

estimate prob θ of landing H.

Frequentist: $\hat{\theta}_{ML} = 1$ (Both are heads)

Bayesian:



In Freq. random variabn is not predictable.
Simply collect statistics on data.

In Bayesian, random variable is outcome of a deterministic proc, though some of the factors are known to observer.

Almost all data is random from Bayesian point of view.

Q1) Radioactive decay, quantum mechanics not bayesian

In MLE, no of data \gg dim of θ .

In Bayesian, $\forall n$.

Freq framework is a limit case of Bayesian.

$$\lim_{n/d \rightarrow \infty} p(\theta | x_1, \dots) = \delta(\theta - \theta_m)$$

→ centered around θ_m

There is no contradiction b/w classical & Bayesian.

Advantages of Bayesian

- prior info
- prior is a form of regularisation
- point estimate of θ posterior contains info about uncertainty.

Probabilistic ML model

$x \rightarrow$ observed data (features)

$y \rightarrow$ hidden / latent

$\theta \rightarrow$ model parameters

Discriminative model

$$P(y, \theta | x)$$

$$\approx P(y|x, \theta)P(\theta)$$

usually prior
distn θ is as
indip. of x .

generative

$$p(x, y, \theta) \approx P(x, y | \theta)P(\theta)$$

can generate
new $P(x, y)$

Training By models

given $(x_r, y_r) \in$ model $P(y, \theta | x)$

training: Bayesian inf over θ

$$P(\underline{\theta}) \propto P(\theta | x_r, y_r) = \frac{P(y_r | x_r, \theta)P(\theta)}{\int P(y_r | x_r, \theta)P(\theta) d\theta}$$

for testing:

$$P(y | x_r, x_r, y_r) = \int P(y | x, \theta)P(\theta | x_r, y_r) d\theta$$

↓
weighting across the
learned distribution

The lower integrals may be intractable.

Conjugate distribution

$p(x)$ & $p(x|y)$ are conjugate iff $p(y|x)$ belongs to same parametric family as $p(y)$. Then it can be integrated.

e.g. for coin toss.

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

Bern $p(x|\theta) = x(1-\theta)^{1-x}$.

now $p(\theta)$ the prior should also belong to parametric family of $p(x|\theta)$, which is bernoulli.

What distn is conjugate to bernoulli?

Beta distri matches requirements

$$\text{Beta}(\theta|a,b) = \frac{1}{B(a,b)} \theta^{a-1} (1-\theta)^{b-1}$$

and this can be tweaked to make it bell shaped to model coin toss.

$$p(\theta) = C\theta^{\alpha}(1-\theta)^{\beta}$$

$$p(\theta|x) = C' p(x|\theta)p(\theta) \propto C' \theta^{\alpha} (1-\theta)^{1-x} \theta^{\alpha-1} (1-\theta)^{b-1} / B(a,b)$$

$$= C'' \theta^{\alpha} (1-\theta)^{\beta}$$

∴ these two are conjugate.

$$P(\theta|x) = \frac{1}{Z} P(x|\theta) P(\theta) = \frac{1}{Z} \prod_i P(x_i|\theta) p(\theta)$$

$$= \frac{1}{Z} \pi^\theta \cdot \theta^{\sum x_i} (1-\theta)^{\sum 1-x_i} \frac{1}{B(a,b)} \theta^{a-1} (1-\theta)^{b-1}$$

$$= \frac{1}{Z'} \theta^{\sum x_i + a - 1} (1-\theta)^{\sum 1-x_i + b - 1}$$

new parameters $\stackrel{?}{=} \text{Beta}(\theta|a', b')$

$$a' = a + \sum x_i \quad b' = b + n - \sum x_i$$

Conjugate distributions

Likelihood $p(x y)$		Conjugate prior $p(y)$
gaussian	y	gaussian
gaussian	μ	gamma
gaussian	σ^{-2}	gaussian-gamma
Multivariate gaussian	(μ, σ^{-2})	wishart
Bernoulli	Σ^{-1}	
Multinomial	P	Beta
Poisson	P	Dirichlet
Uniform	λ	Gamma

$x_i | \theta) p(\theta)$
what to do if there is no conjugacy?

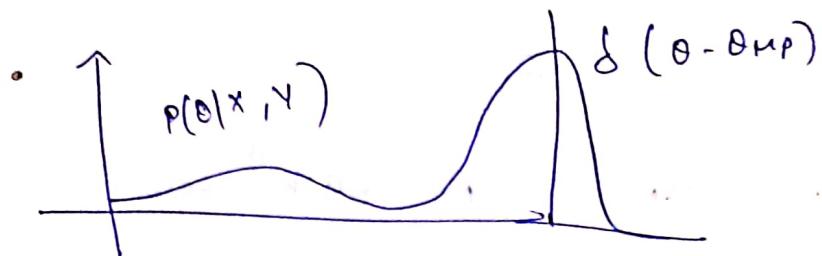
- Max posterior estimation

$$\hat{\theta}_{MP} = \arg \max \theta p(\theta | x_{tr}, y_{tr}) = \arg \max \theta p(y_b | x_{tr}, \theta) p(\theta)$$

~~This~~ cov L₂ Reg corresponds to this

In testing:

$$P(y | x, x_{tr}, y_{tr}) = P(y | x, \hat{\theta}_{MP})$$



Bayesian Reasoning

1) dark mark

stays 20% if the dir.

stays 100% if he is alive

survival chanu ~~0.001~~, 1%.

what is P of being alive if mark is present?

$x \in (0,1)$ 1 if he is ~~not~~ alive

$y \in (0,1)$ 1 if mark is visible

$$P(y=1 | x=1) = 1 \quad P(y=1 | x=0) = 0.2 \quad P(x=1) = \frac{1}{10}$$

$$P(x_2=1 | y_2=1) = ?$$

$$P(x_2=1 | y_2=1) = \frac{P(y_2=1 | x_2=1) P(x_2=1)}{\sum P(y_2=1 | x_2=j) P(x_2=j)}$$
$$= \frac{1 \times \frac{1}{101}}{\frac{1 \times \frac{1}{101}}{1} + \frac{1}{5} \times \frac{100}{101}}$$
$$= \frac{1}{21}$$

2) $x_2 \in \{x_1, x_2, \dots, x_N\}$, N independent rolls

$$N_k = \sum_{n=1}^N I(x_n = k) - \text{counts}$$

$$P(x|\theta) = \prod_{k=1}^K \theta_k^{N_k} - \text{multinomial likelihood}$$

MLE

$$\theta_{ML} = \arg \max \log P(x|\theta)$$

Refer [github](#).[Bookmark]

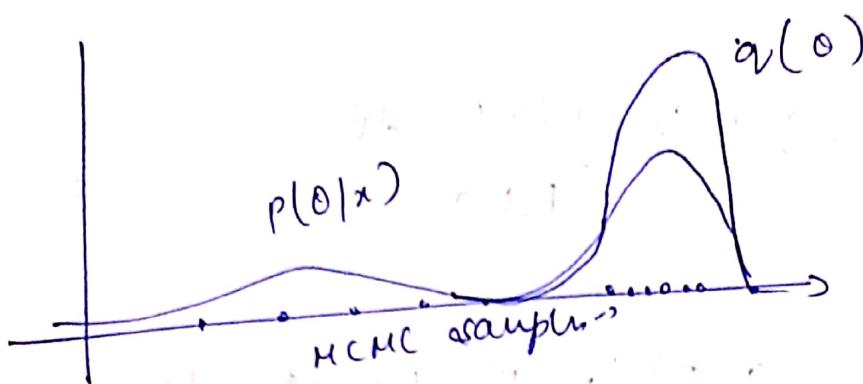
Approximate inference.

Variational inference (VI)

- Approx $P(\theta|x) \propto q(\theta) \epsilon \Phi$
- Biased
- Fast & $\epsilon_{\text{var}}^{\text{max}}$ scalable.

Monte Carlo
Markov chain MCMC

- Samples from unnormalised $p(\theta|x)$
- unbiased
- lot of samples needed



Dissimilarity measure for VI

- can be done using KL divergence

$$F(q) := \text{KL}(q(\theta)||p(\theta|x)) \xrightarrow[q(\theta) \in Q]{\min}$$

$$\text{KL}(q||p) = \int q(\theta) \log \frac{q(\theta)}{p(\theta|x)} d\theta$$

Properties

$$\text{KL}(q||p) \geq 0$$

$$\text{KL}(q||p) = 0 \Leftrightarrow q = p$$

$$\text{KL}(q||p) \neq \text{KL}(p||q)$$

To compute KL div, $P(\theta|x)$ is needed.
 Though if it is obtained, then KL div is not needed.

Some math

$$\log P(x) = \int q(\theta) \log (P(x)) d\theta \quad (\text{Identity})$$

$$= \int q(\theta) \log \frac{P(x, \theta)}{P(\theta|x)} d\theta$$

$$= \int q(\theta) \frac{\log P(x, \theta) q(\theta)}{P(\theta|x) q(\theta)} d\theta$$

$$= \int q(\theta) \frac{\log P(x, \theta)}{q(\theta|x) q(\theta)} d\theta + \int q(\theta) \log \frac{q(\theta)}{\frac{P(x, \theta)}{P(\theta|x) q(\theta)}} d\theta$$

$$\stackrel{\log P(x)}{\downarrow} \quad \stackrel{\text{doesn't depend}}{\downarrow} \quad \stackrel{L(\theta | q(\theta))}{\downarrow} + \text{KL}(q(\theta) || P(\theta|x))$$

Evidence lower bound (ELBO) KL div
depends on q

Instead of mini KL div, we can maximise ELBO in above eq

In ELBO, everything is computable
 No true posterior estimate is needed

$$\begin{aligned}
 L(\theta|o) &= \int q(\theta) \log p(o|\theta) d\theta + \int q(\theta) \log \frac{p(\theta)}{q(\theta)} d\theta \\
 &= \underbrace{\mathbb{E}_{q(\theta)} \log p(o|\theta)}_{\text{encourages MLE}} - \underbrace{\text{KL}(q(\theta)||p(\theta))}_{\substack{\text{to maximize } L, \\ \text{minimize this,} \\ \text{so make it 0}}} \\
 &\quad (\text{converges to prior})
 \end{aligned}$$

If L is maximised over all distributions, then ~~KL~~ we will arrive at true posterior distn $p(\theta|x)$. (not actually true since it is intractable)

how to perform optimisation wrt distn?

- ① Mean field approx MFA

② Parametric approx

① $q(\theta) = \prod q_j(\theta_j)$, $\theta = [\theta_1, \theta_2, \dots, \theta_m]$

q is factored as product of m multipliers

② $q(\theta) = q(\theta|x)$

① MFA

all θ_i 's are assumed ~~indep~~ independent of each other

(Simpler approx)

Using Block coordinate ascent:

At each step, fix all factors $\{q_i(\theta_i)\}_{i \neq j}$,
except one and optimize w.r.t to it.

$$L(q(\theta)) \rightarrow \max_{q_j(\theta_j)}$$

do it for all $j \in m$ till convergence.

$$L(q(\theta)) = E \log p(x, \theta) - E \log q(\theta)$$

$$\Rightarrow \quad \quad \quad - \sum E_{q_k(\theta_k)} \log q_k \theta_k$$

$$= E_{q_j(\theta_j)} [E_{q_{j \neq j}} \log p(x, \theta)] - E_{q_j(\theta_j)} \log q_j(\theta_j)$$

Consider

$$\{r_j(\theta_j) = \frac{1}{z_j} \exp(E_{q_{j \neq j}} \log p(x, \theta))\} \quad \text{--- ①}$$

$r_j \rightarrow \text{PDL}$

$$\Rightarrow E_{q_j(\theta_j)} \log \frac{r_j(\theta_j)}{q_j(\theta_j)} + C = -KL(q_j(\theta_j) || r_j(\theta_j)) + C$$

for -ve KL,

optimal $q_j = r_j$

$$q_j(\theta) = r_j(\theta) = \frac{1}{z_j} \exp(E_{q_{j \neq j}} \log p(x, \theta))$$

Update each q_j ,
convergence. This
converge. repeat till
is guaranteed to

• conditions to complete ①

$$p(x, \theta) = p(x | \theta) p(\theta) \quad \theta = [\theta_1, \dots, \theta_m]$$

• conditional conjugacy of likelihood and prior on each θ_j conditioned on all other $\{\theta_i\}_{i \neq j}$

$$p(\theta_j | \theta_{i \neq j}) \in A(\alpha), \quad p(x | \theta_j, \theta_{i \neq j}) \in B(\theta_j) \rightarrow$$

• the prior w.r.t. each

$$p(\theta_j | \theta_{i \neq j}) \in A(\alpha)$$

subgroup θ_j given all other $\theta_{i \neq j}$ fixed,

and MLE fn as a fn of θ_j given all other

$\theta_{i \neq j}$ fixed, form conjugate pair of distri.

Then posterior distri $i|j$ belongs to same
parametric family as prior, this should

hold true for each subgroup. Then if

this holds true, then update equation ①
can be done in closed form

In practice,

$$\log q_j(\theta_j) = \sum_{x_{i|j}} \log p(x, \theta) + \text{const}$$

Inference methods: Summary

Full Bayesian inference: $p(\theta|x)$

MAP

$$: p(\theta|x) \propto \delta(\theta - \theta_{\text{MAP}})$$

Mean field variational Inference: $p(\theta|x) \approx q(\theta) = \prod_i q_i(\theta_i)$

Parametric variational inference: $p(\theta|x) \approx q(\theta) = q(\theta)$

Latent Variable Models and

EM

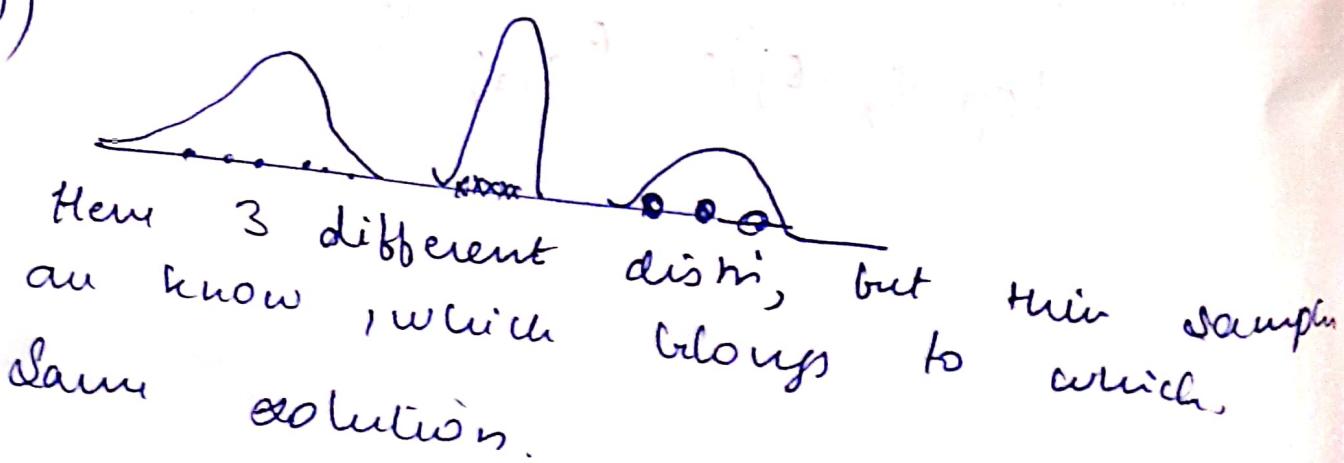
LVM eg:

i) $x_i \sim N(x_i | \mu, \sigma^2) = \frac{1}{\sqrt{2\pi\sigma^2}} \exp\left(-\frac{(x-\mu)^2}{2\sigma^2}\right)$

estimate $\mu \& \sigma$

Soln: MLE eq for $\mu \& \sigma$.

ii)



iii)

now we don't know which one was derived from which gaussian
- we can try to fit on gaussian.



but not the best choice because there is date mining b/w many factors etc

- To solve this, establish a LVM.

for each x_i , establish an additional latent variable z_i which denotes index of gaussian IP from which x_i was derived.
if data is assumed IID (independent identical distri)

$$p(x, z | \theta) = \prod p(x_i | z_i, \theta) p(z_i | \theta)$$

$$= \prod p(x_i | z_i, \theta) p(z_i | \theta)$$

$$= \prod \prod_{i=1}^n N(x_i | \mu_{z_i}, \sigma_{z_i}^2)$$

$\pi_j = p(z_j = j)$ are prior of j^{th} gaussian

if both X & z are known,

$$\theta_{ML} = \arg \max_{\theta} p(x, z | \theta)$$

$$\theta_{ML} = \arg \max_{\theta} \log p(x, z | \theta)$$

$$= \arg \max_{\theta} \log \prod_{i=1}^n N(x_i | \mu_{z_i}, \sigma_{z_i}^2)$$

But Z is not known.

Then we need to maxi wrt Θ_z the log
of incomplete likelihood.

$$\log P(x|\theta) = \int q(z) \log \frac{P(x, z|\theta)}{q(z)} dz \xrightarrow{\text{ELBO}}$$
$$+ \int q(z) \log \frac{q(z)}{q(z|x, \theta)} dz \xrightarrow{\text{②}}$$
$$\downarrow$$
$$KL(\text{non neg})$$
$$= L(q, \theta) + KL(q||P) \geq L(q, \theta)$$

Instead of optimising log of incomplete likelihood $\log P(x|\theta)$, we optimise

ELBO wrt $\Theta \& q(z)$.

The block-coordinate algorithm is known as EM algorithm

Variational lower bound definition

fn $g(\xi, x)$ a VLB of a fn $f(x)$ iff

- for all $\xi \in \mathbb{R}^n$, it follows $f(x) \geq g(\xi, x)$,
- for any x_0 , there exists $\xi(x_0)$ such that $f(x_0) = g(\xi(x_0), x_0)$

If such a VLB exists, then $f(n) \rightarrow \max_x$

we can perform block-coordinate update
EM algorithm

$$L(q, \theta) = \int q(z) \log \frac{P(x, z | \theta)}{q(z)} dz \rightarrow \max_{q, \theta}$$

Start with θ_0 & iterate

E-step:

$$q(z) = \arg \max_q L(q, \theta_0) = \arg \min_q KL(q || p) \\ = p(z | x, \theta_0)$$

due to decompositio, of ④, max q 1st is min q

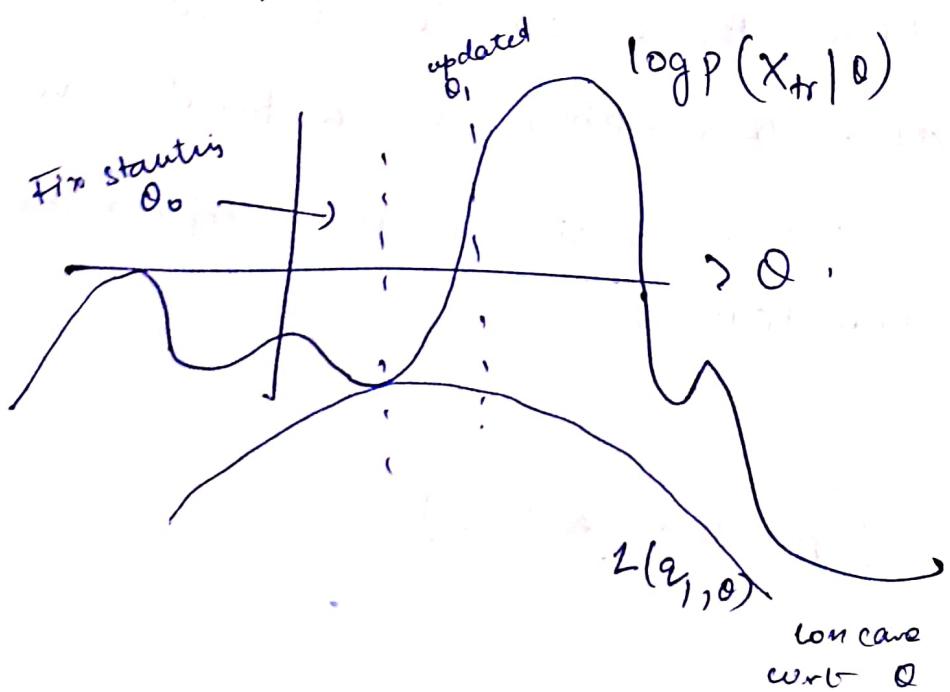
M-step

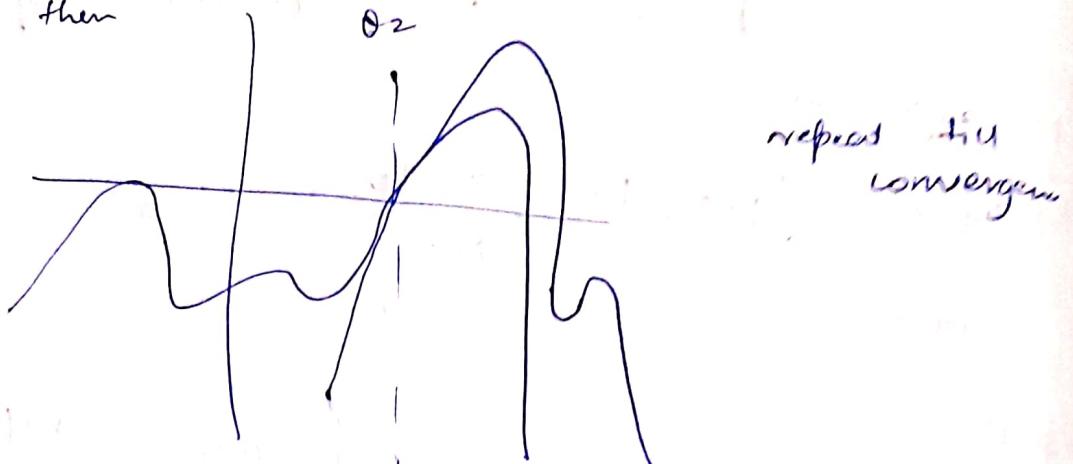
fix q .

$$\theta_* = \arg \max_{\theta} L(q, \theta) = \arg \max_{\theta} \sum_z \log p(x, z | \theta)$$

log of likelihood is usually concave fn.

maximisation of concave fn, is a convex optim problem





Converges to local maximum.

Benefits of EM

- E & M steps can be done in closed form (eg: mix of gaussians)
 - allows to build complicated models using mix of simple distributions
 - If true posterior is intractable, then closed $q(z)$ can be solved by optimisation
 - Allows to process missing data by treating them as latent variables
- discrete ^(latent) distn of variables is a finite mix of dists
- continuous latent. variable is representation having
- ef: PCA, ICA
- [problem set 2]

Stochastic inference and variational VAE

continuous latent variable

$$p(x_i | \theta) = \int p(x_i, z_i | \theta) dz_i = \int p(x_i | z_i, \theta) p(z_i | \theta) dz_i$$

$$q(z_i) = p(z_i | x_i, \theta) = \frac{p(x_i | z_i, \theta) p(z_i | \theta)}{\int p(x_i | z_i, \theta) p(z_i | \theta) dz_i}$$

e.g.: PCA

$$x \in \mathbb{R}^D, z \in \mathbb{R}^d, D \gg d$$

$$p(x, z | \theta) = \prod p(x_i | z_i, \theta) p(z_i | \theta)$$

assume standard gaussian as latent.

$$\approx \prod N(\alpha_i | v z_i + \mu, \sigma^2 I) N(z_i | \mu, \sigma^2)$$

$\theta \rightarrow D \times d$ matrix v , D dim μ & σ^2 .

Can use EM to find $\arg \max_{\theta} p(x_t, f_t | \theta)$

Advantages of prob PCA:

- closed form soln has $O(D^2)$, EM has $O(n D d)$
- can process missing parts in x and present parts in z
- can determine d if $p(\theta)$ is established
- can be generalized to more general models such as mixture of PCA

Mixture of PCA

- Two types of latent mode variables:
discrete $t \in \{1, \dots, k\}$ & $z \in \mathbb{R}^d$

t is the index of corresponding subspace.

$$p(x, z, t | \theta) = \prod p(x_i | t_i, z_i, \theta) p(t_i | \theta)$$

$$= \prod \pi_i N(x_i | \nu_t, z_i + \mu_t, \sigma_t^2 I) N(z_i | 0, I)$$

E

$$q(z, t) = \prod q_{t_i}(z_i, t_i) = \prod p(z_i, t_i | x_i, \theta)$$

$$= \prod \frac{N(x_i | \nu_t, z_i + \mu_t, \sigma_t^2 I) N(z_i | 0, I) \pi_{t_i}}{\sum_{t=1}^k N(x_i | \nu_t, z_i + \mu_t, \sigma_t^2 I) N(z_i | 0, I) \pi_t}$$

M

$$E_{z, t} \log p(x, z, t | \theta)$$

$$= \sum_{i=1}^n E_{t_i, z_i} (\log p(x_i | t_i, z_i, \theta) + \log p(z_i | 0))$$

$$+ \log p(t_i | \theta)) \rightarrow \max_{\theta}$$

non-linear generalisation by PCA \rightarrow VAE

$$P(x, z | \theta) \propto \prod_{i=1}^n p(x_i | z_i; \theta) p(z_i)$$

$$\approx \prod_{i=1}^n \left(\prod_{j=1}^D N(x_{ij} | \mu_j(z_i), \sigma_j^2(z_i)) \right) N(z_i | \mu, \Sigma)$$

↓ multi-modal prior

here μ, σ is substituted by some non linear fn of z . μ_j , σ_j are \circ/ρ of θ , a NN, parametrized by θ .

It takes z_i as i/p.

allows us to construct very non linear low-dim manifolds of latent sub-spaces of \mathbb{R}^d .

Applications: more conjugacy due to N.N

$$\int p(x, z | \theta) dz \propto \prod_{i=1}^n \int p(x_i | z_i; \theta) p(z_i) dz_i$$

not optimisable

if we use EM, even it is intractable for true posterior distri

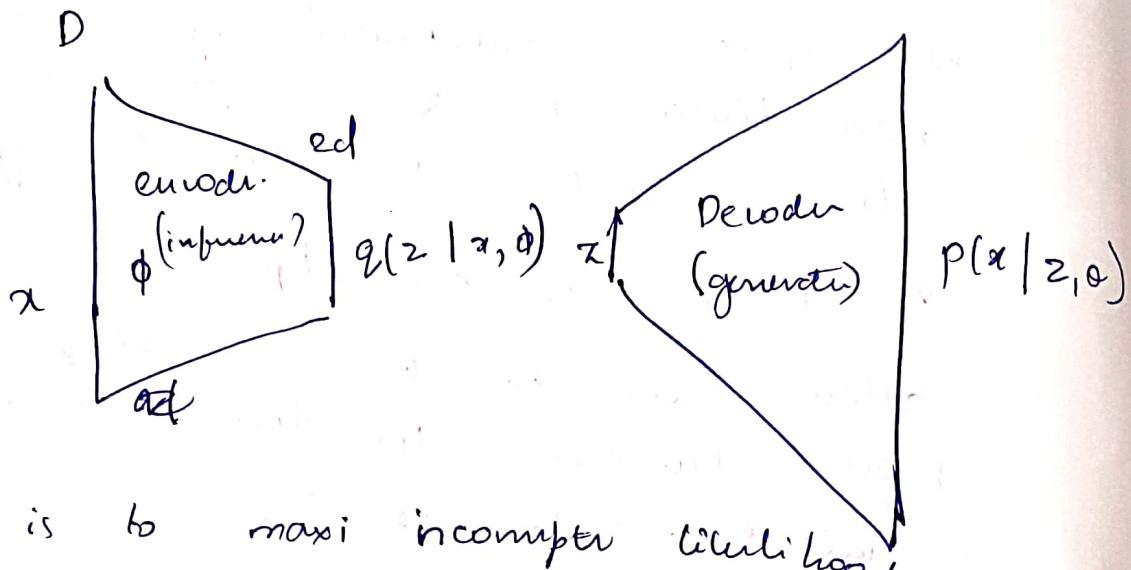
using variational bayes true posterior can be approximated.

$$q(z_i|x_i, \phi) \approx p(z_i|x_i, \theta)$$

$$= \prod N(z_{ij} | \mu_j(x_i), \sigma_j^2(x_i))$$

this μ, Σ, σ is given by

another neural n/w. with weights γ ,



goal is to maximize in-computer likelihood.

$\mathcal{E}/$

$$q(z|x, \phi) = \arg \min_{\phi} \text{KL}(q(z|x, \phi) || p(z|x, \theta))$$

This is equivalent to optimising ELBO,

$$L(\phi, \theta) = \int q(z|x, \phi) \log \frac{p(x|z, \theta) p(z)}{q(z|x, \phi)} dz$$

Block of cord-opti

$\rightarrow \max_{\phi, \theta}$

wrt encoder $\phi \rightarrow \mathcal{E}$

Block cord opti

wrt decoder $\theta \rightarrow \mathcal{M}$

- Training data is assumed to be large
- integral in ELBO still intractable.

Soln:

Compute stochastic gradients using mini batch
& monte carlo estimation.

Practical Session: VAE

(notebook tasks)

1. loss fn for VAEs in pytorch.
2. DRAW.

Discrete variable modals

Why discrete?

- Easier to interpret
- Allow the model to make a discrete choice
 - hard attn.
 - binary masks where to look at
 - classifies modified images
- To get certain properties for discrete pred.
 - GANs for text

Relax objective over discrete random variables

an objective over continuous random samples
& w/ reparametrisation trick.

Keep discrete in test case

for reparametrisation, $f(z)$ should be differentiable

Gumbel Max Trick

$Z \sim \text{categorical}(\pi_1, \dots, \pi_k)$

Can be computed as

$$Z \stackrel{d}{=} \underset{k}{\operatorname{argmax}} \frac{\pi_k}{\pi_k + \epsilon}, \quad \epsilon \sim \text{Exp}(1)$$

↓

This will have categorical distribution.

• apply $-\log$ to make it computationally stable

$$Z \stackrel{d}{=} \underset{k}{\operatorname{argmax}} [\log \pi_k - \log \bar{\pi}_k] \quad \bar{\pi} \sim \text{Exp}(1)$$

argmax is not differentiable.

approx argmax with softmax

softmax with Temp to relax SL.

$$\text{softmax}_T(x)_j := \frac{\exp(x_j/T)}{\sum \exp(x_i/T)}$$

temp allows to control stepness of softmax

at $T=0$,

argmax = softmax,

at $T=\infty$, softmax leads to uniform normed

assum distin

Z is a one-hot vector \in

replace with a continuous relaxation \tilde{Z}

$$\tilde{Z}(\gamma, \pi) := \text{softmax}(\log \pi_1 + \gamma_1 + \dots + \log \pi_k + \gamma_k)$$

Then can reparametrize and gradient estimator.

choosing temp

- bias / variance trade off
 - small temp leads to high var, low bias
 - large temp leads to low var, high bias
- grid search.
- Toy example \rightarrow implementation

going for conti distir for training & discrete distir for testing, then is the mismatch

Using Reinfor with reduce in Variance:

Control variance: $\mathbb{E}_{q(z)} f(z) = \mathbb{E}_{q(z)} b(z)$

convinient b(z) with tractable expectation

$$\mu^* = \mathbb{E}_{q(z)} b(z)$$

$$\mathbb{E}_{q(z)} f(z) = \mathbb{E}_{q(z)} [(f(z) - b(z)) + \mu]$$

↓

the control variance might be lower var estimate

if $f(z) \neq b(z)$ are

positively correlated

- Unbiased
- lower bound on B
- convenient if $b(z)$ is zero-mean
- Can take several samples to reduce variance further.

$b(z)$ is called baseline fn in Reinforce choosing baseline:

- constant

$$b(z) = c$$

$$(f(z) - c) \frac{\partial}{\partial \phi} \log q_\phi(z) + \underbrace{\frac{\partial}{\partial \phi} E_{q_\phi(z)}}_{=c}$$

formula for
optimal constant
such that var is
minimal

$$c_2 = \frac{\sum \text{cov}[f(z)] \frac{\partial}{\partial \phi_d} [\log q_\phi(z) \frac{\partial}{\partial \phi_d} \log q_\phi(z)]}{\sum \text{Var}[\frac{\partial}{\partial \phi_d} \log q_\phi(z)]}$$

- can be estimated by moving average,
but variance might go higher.
- $b(z)$ as a first-order Taylor series expansion of $f(z)$ at some point μ .

$$b(z) = f(\mu) + \frac{\partial f}{\partial z}(\mu)^T (z - \mu)$$

$$\text{H:z } M(\phi) = E_{q_\phi(z)}^2$$

- Backpropagates through mean, then fine tunes inaccuracies with reinforce
- taylor series approx is expensive

- variance minimisation
 - minimising $\text{grad } \varphi$: var.
 - $\text{var}[g(z, \phi)] = E[g(z, \phi)]^2 - (Eg(z, \phi))^2$
 - In general minimising variance leads to increase in bias
- more complicated gradient estimator - REBAR
 - uses gumbel-relaxed f as baseline
 - Practical

Fairness in ML

- Bias is everywhere
 - uncontrolled bias will lead to unfairness
- Sources of unfairness

- Bias from data
- Bias from algorithm

Training data

- Sampling bias
not representative of the overall target population
- Selective label
only observe outcome of one side of decision

- proxy labels

e.g. want to predict who will commit crime but have only data on who is arrested.

From algo

- when there is class imbalance (could be due to sampling error)

- feedback effect

model at time $t+1$ uses data of time t

- In imSitu situation recognition dataset, 33% more women than men in data, this amplifies to 68% by the algo.

- In Adult Income dataset (UCI), male income is high while 11% of female has high income. The skewness ratio is amplified from 3:1 to 5:1.

Enforcing fairness

Fairness definition

- Discrimination in the law.

- treating ppl from diff groups diffly
(direct design)

- seemingly treat all groups equally but consequences are different
(indirect)

- To remove direct, remove info on sensitive grouping.

- Indirect: enforce equality of outcomes, enforcing equality on outcomes?

- same no of positive predictions
(statistical demographic parity)

- this produces lower accuracy.
Claim we are enforcing against biased data \rightarrow forcing misclassification

- True positive rate outcome should be same.

- enforce true positive rate same for all groups.

- Equalised odds

- True positive rate & false positive rate should be same for all claim.

Statistical fairness notions

- Statistical parity ($\hat{y} \perp s$): $\Pr(\hat{y}_{=1}|s_{=0}) = \Pr(\hat{y}_{=1}|s_{=1})$
 independent

prediction \hat{y}

sensitive group s

\hat{y} independent of s

- Equalised odds ($\hat{y} \perp s \mid M$)

$$\Pr(\hat{y} = y|s=0, M) = \Pr(\hat{y} = y|s=1, M)$$

Conditional independence.

$\hat{y} \perp s$ given true label info.

forcing TP & TN to be equal.

- Predictive parity:

$$\Pr(y_{=1}|s_{=0}, \hat{y}_{=1}) = \Pr(y_{=1}|s_{=1}, \hat{y}_{=1})$$

Reverse the conditional independence.

Not possible to enforce all 3.

By Bayes Rule.

Equalised odds & predictive parity can be enforced at same time if priors are equal (perfect data) or if $TPR_2 = 1$, $FPR_2 = 0$ (Perfect classifier)

(S,1) equalised odds and statistical parity
at once only in case of perfect dodant

Problem Statement

total applicants = 20k.

50% blue group

50% green group,

acceptance rate 50%.

entrance test:

80% who passed will graduate

10% who failed could graduate

60% of blue will pass the test

40% of green will pass the test

confusion tables:

statistical panels

		blue	
		Accepted	not
graduates	graduates	4000 (80%)	1200
	doesn't	1000 (20%)	3800
		5000	

		green	
		Accepted	not
graduates	graduates	3300	600 (10%)
	doesn't	1700	4500 (90%)
		5000	

10% of qualified blue rejected while

10% of unqualified green accepted

equality of opportunity

	Accepted	not	<u>bias</u>
graduate	4440	760	
doesn't	1110	3690	
	<u>SSD</u>		<u>green</u>
graduate	3245	555	
doesn't	1205	4995	

SSY. of blue is 44.5% of green, TRL = 88.4%, for both

Predictive parity

60: 40 accept band on test only.

	accepted	not
graduate	4800	400
doesn't	1200	3600
	<u>6000</u>	
graduate	3200	600
doesn't	800	5400
	<u>4000</u>	

This can reinforce bias
 (what if blue is well educated etc)
 (if more represented in future, other approaches)

- Only group effect is observed, not wrt each individual
- Can't use data for individual farmers.
eg: other sources like expert data from other sources to tackle this
- need to talk with domain expert to decide on which action

Algorithmic fairness methods

- pre processing
- post processing
- in processing (with learning)

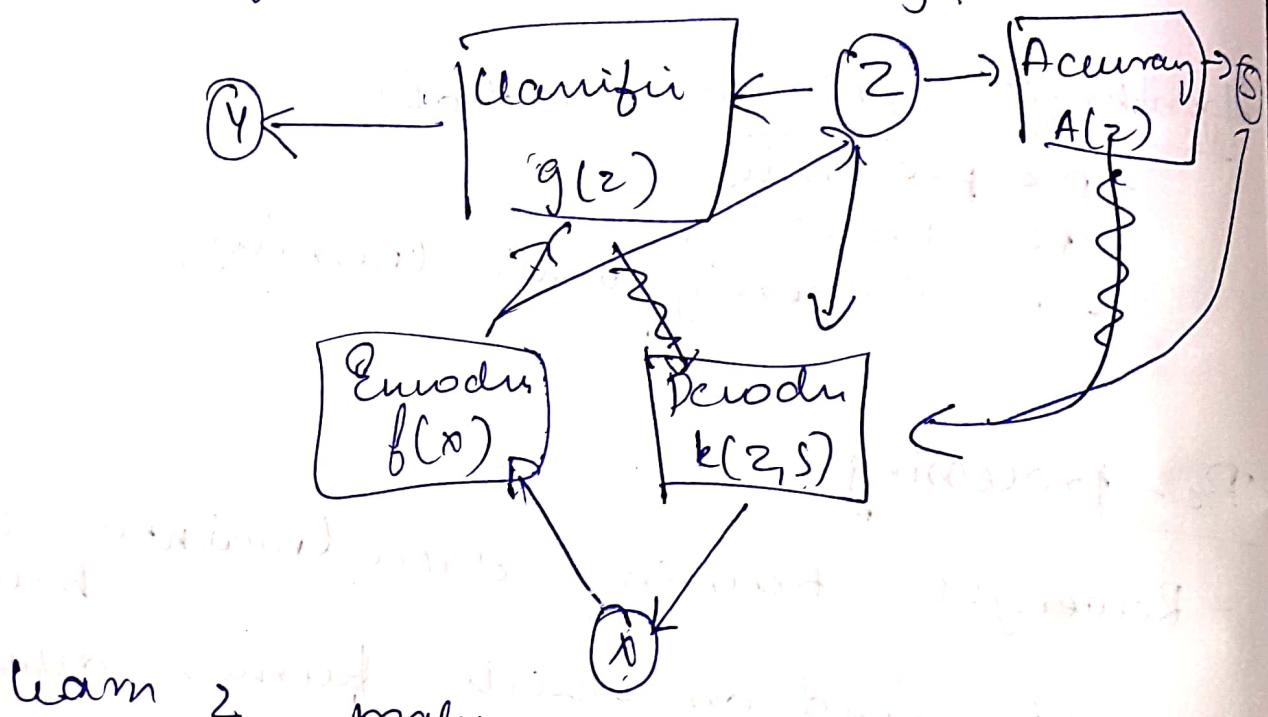
Pre-processing

- Reweight training data (maintain class balance)
weights based on which fairness definition to be followed.

$$w(s=0, y_{21}) = \frac{Pr(y_{21}) Pr(s=0)}{Pr(y_{21}, s=0)}$$

, for ↓ lesser data class,
larger weight, so
more times fed to
model.

- From reweighting to resampling.
 - Sampling data with replacement according to weights
- Learn a new, fair representation of data.
 - using adversarial learning.



Learn z , make z is informative
wrt y , but cannot predict s
di limitations

- latent embeddings still has some info on s if a diff classifier is used on z .

In processing

- specify fairness metric as constraint
- optimise for accuracy under this condition
- Then low accuracy

Constraint:

loss for $L(\theta)$

normal: $\min_{\theta} L(\theta)$

statistical family:

subject to $P(\hat{y}_{21}|s_{20}) \geq P(\hat{y}_{21}|s_{21})$.

- Not Convex
- Need to find a better way.
- Return a set of models that trade-off b/w fairness & accuracy
- Likelihood with dependency on s
 $S: P(y_{21}|x, \theta, s)$.

Introduce latent target label \tilde{y}

$$Pr(y_{21}|x, \theta, s) \geq \sum_{\tilde{y}} Pr(y_{21}|\tilde{y}, s) Pr(\tilde{y}|x, \theta, s)$$

Post-processing.

Using S at prediction time,
use 2 diff. thresholds

Interpretability in fairness

a fair rep. $T_w(x)$ has sensitive reward
but is dependent on x

- Separate fair & unfair loss
- use generative training, adversarial learning to generate data points with pre-specified sensitive attributes

GRANS

Parameterisation

Sample from two distribution
 $p(x) \in q(x)$ taken.

Parameterising model

- Parameter density function

$$q_\theta(z) = N(z; \theta, I)$$

- Define distribution $q(z)$ implicitly
 $z \sim N(0, I)$

$$g_\theta(z) \sim q_\theta(z)$$

mapping $z \in x$ + simple sampling

Images.

Noise $\sim \mathcal{N}(0, \sigma)$



Classifier

$$f_{\phi}(x) = P_{\phi}(y=1|x)$$

Loss:

BCE

$$L(\phi, \theta) = E_{P(x)}[-\log f_{\phi}(x)] + E_{P(x)}[-\log(1 - f_{\phi}(G_{\theta}(z)))]$$

update classifier

$$\phi^* = \arg \max L(\phi, \theta)$$

update generator

$$\theta^{\text{new}} = \theta^{\text{old}} + \frac{\partial L(\phi^*, \theta^{\text{old}})}{\partial \theta}$$

repeating
learning scheme

1. update guide

- $P(x)$
- $\frac{P(x) + q(x)}{P(x) + q(x)}$
- $\frac{P(x)}{q(x)}$
- $D(p||q)$

2. update generator

- Move $q(\mathbf{z}|\mathbf{x})$ closer to $p(\mathbf{z})$.

3. Repeat.

Prescribed (eg VAE)

- direct access to both prior, likelihood, conditional distributions.

(implicit (GAN)) :

- evaluate Σ , sample prior $p(\mathbf{z})$, posterior can only be sampled. Can approx $q(\mathbf{z})$, $q(\mathbf{z}|\mathbf{x})$ from sampler.

for classification

$$\text{Instead of } \min \mathcal{L}(\phi, \theta)$$

$$\max -\mathcal{L}(\phi, \theta) = -\log 4 + 2D_{JS}(p||q_{\phi})$$

Estimate distance $D(p||q_{\phi})$
 $P(\mathbf{x})$ and $q(\mathbf{x})$

update generator

- same
- minimize distance.

This becomes a min max problem w/r/t
classifying a generator

F-divergence & w-distance is taken.

F-divergence:

for $P(\pi) \in \Phi(\pi)$.

$$D_f(P||Q) = \int_{\pi} f\left(\frac{P(x)}{Q(x)}\right) Q(x) dx$$

different forms of f gives different divergences.

$$f(t) = t \log(t) \rightarrow \text{KLD}$$

$$D_f(P||Q) = \text{KL}(P||Q)$$

$$f(t) = -\log(t) \rightarrow \text{Reverse KLD}$$

$$D_f(Q||P) = \text{KL}(Q||P)$$

$$f(u) = \frac{1}{2}|t-1| \rightarrow \text{total variation}$$

$$D_f(P||Q) = \frac{1}{2} \int_{\pi} |P(x) - Q(x)| dx$$

w-distance

cost of transporting x to $y \rightarrow c(x, y)$

$$\text{eg: } c(x, y) = \|x - y\|$$

optimal transport

$$T(P, Q) = \inf_{\Gamma \in \mathcal{P}(x \sim P, y \sim Q)} E_{(x,y) \sim \Gamma} [c(x, y)]$$

↓
Set of all joint distri
of (x, y)
with $P \in \mathcal{D}$.

↓
expected cost

OT dual:

- Primal

$$T(P, Q) \geq \inf_{E \in \mathcal{C}} E_{x \sim P, y \sim Q} [c(x, y)]$$

- Dual (ω_1 mimic)

$$T(P, Q) = \omega_1(P, Q) = \sup_{\|f\|_L \leq 1} E_{x \sim P} f(x) - E_{x \sim Q} f(x)$$

Liebnitz

constant ≤ 1

b/w continuity
and smoothness.

- ① ↴ can be enforced by
clip weights of n/w 2 methods
hope ops of n/w for
similar derivati points do not
much

② gradient penalty.

Evaluate expectation over set of \tilde{x} .

$$\lambda E_{\tilde{x} \sim P(\tilde{x})}$$

$$\left(\| \nabla_x D(\tilde{x}) \|_2 - 1 \right)^2$$

$\tilde{x} \rightarrow$ convex combination b/w real & fake data points.

enforce gradient of the discriminator to be almost everywhere.

Optimal Transport vs f -D

$$z \sim U[0, 1]$$

$$\nu \in (0, 2) \text{ sfm at } 0$$

$$\theta \in (0, 2) \text{ sfm at } \theta$$

$$W_1(P, Q) \geq 0 \rightarrow \text{gives conti signal}$$

$$JS(P||Q) = \begin{cases} \log(2) & \theta \neq 0 \\ 0 & \theta = 0 \end{cases} \rightarrow \text{grad} \approx 0$$

$$KL(P||Q) = \begin{cases} \infty & \theta \neq 0 \\ 0 & \theta = 0 \end{cases} \quad \begin{matrix} \text{don't care} \\ \text{due to } \infty \end{matrix}$$

Practical tricks to make GAN training easier
(Always use them)

• Spectral Normalisation

- grad penalty - is a strict constraint
- introduces significant bias

- easier way to introduce L-1 conti

- idea is now is composition of some fns, if we can make each of these L-1 conti, then we will be a L-1 conti fn. So linear fn can be made L-1 ifn. It is good with some non-linearity.
- Can conv be made L-1 conti?

$$\begin{bmatrix} x_1 & x_2 & x_3 \\ x_4 & x_5 & x_6 \\ x_7 & x_8 & x_9 \end{bmatrix} \xrightarrow{\begin{pmatrix} a & b \\ c & d \\ e & f \\ g & h & i \end{pmatrix}} \begin{bmatrix} w_1 & w_2 \\ w_3 & w_4 \\ w \end{bmatrix} \xrightarrow{x_0 w^2}$$

$$\begin{bmatrix} w_1 & w_2 \\ w_3 & w_4 \end{bmatrix} \xrightarrow{x_0 w^2} \begin{pmatrix} a & b \\ c & d \end{pmatrix}$$

↓
linear o/p.

$$x_0 w^2 = \begin{bmatrix} x_1 w_1 + x_2 w_2 + x_4 w_3 + x_5 w_4 + x_7 w_1 + \dots \\ x_4 w_1 + \dots \end{bmatrix}$$

↓
linear o/p

reformat it into multi of some matrix \mathbf{A} & some vector

$$\begin{bmatrix} w_1 & w_2 & 0 & w_3 & w_4 & 0 & 0 & 0 & 0 \\ 0 & w_1 & w_2 & 0 & w_3 & w_4 & 0 & 0 & 0 \\ 0 & 0 & 0 & w_1 & w_2 & 0 & w_3 & w_4 & 0 \\ 0 & 0 & 0 & 0 & w_1 & w_2 & 0 & 0 & w_3 \end{bmatrix}$$

\mathbf{Aw}

$$\begin{bmatrix} x_1 \\ x_2 \\ \vdots \\ x_9 \end{bmatrix} \xrightarrow{\mathbf{A}} \mathbf{x}$$

\therefore it is a row 2 'A ω ' $\xrightarrow{\text{reshape this from } Q \times D \text{ to } (2 \times 2)}$

L-1 conti of a linear operator

$$\|f(x_1) - f(x_2)\|_2 \leq L \|x_1 - x_2\|_2$$

(can take any norm)

$$\|Ax_1 - Ax_2\|_2 \leq L \|x_1 - x_2\|_2$$

$$\frac{\|Ax\|_2}{\|x\|_2} \leq L$$

$$\sigma_{\max} = \sup_{\|x\|_2=1} \frac{\|Ax\|_2}{\|x\|_2} \leq L$$

Spectral norm of matrix
equal to max singular value of matrix

$$L = \sup_x \frac{\|\frac{A}{\sigma_{\max}}x\|_2}{\|x\|_2} \leq L$$

Singular values:

using SVD.

$$A = U \Sigma V^T \rightarrow \text{right singular}$$

left singular

(orthogonal)

$$\sigma = u^T A v$$

u & v can be used to find σ_{\max} .

to find σ_{\max} :
power iteration

assume EVD for simplicity

$$A = Q \Lambda Q^{-1}$$

$$A = Q \Lambda Q^{-1}$$

let x_0 be init a vander with unit norm

Run the iterations:

$$1. \quad x_{i+1} = Ax_i$$

$$2. \quad x_{i+1} = \underline{x_{i+1}}$$

$$\|x_{i+1}\|_2$$

we converge to EV, q corresponding to λ_m

$$\sigma_{\max} = q^T A q.$$

Similarly σ_{\min} is computed.

All ops are diff wrt A , & w

This is done one step at each grid step
at gen / discr. Keep the previous q and
use them to get better approx.

Theory vs practice

Instead of Aw , a reshaped vector, w

$$\text{Out} \text{ Hout} \text{ Wout} \times \text{Cin} \text{ Hin} \text{ Win}$$

$\text{Cout} \times \text{Cin} K_h K_w$

resulting σ_{\max} is a lower bound to true σ_m
using this power iteration method.

σ_{\max} is equal only for a linear layer.

this outperforms true spectral norm

game theory view on GRANS

In theory.

$$\begin{aligned} \psi^* &: \arg \max_{\psi} L(\theta^{old}, \psi) \rightarrow \text{optimal discr} \\ \theta^{new} &= \theta^{old} - \alpha \nabla_{\theta} L(\theta^{old}, \psi^*) \rightarrow \text{gen single step} \\ \theta^{old} &= \theta^{new} \end{aligned}$$

\uparrow
 $D(P||Q)$

- For most Div, optimal value will be constant.
(zero grad)
- In practice, we don't get exact Div, we estimate it which has non zero grad

In practice:

$$\begin{aligned} \psi^{new} &\in \psi^{old} + \alpha \nabla_{\psi} L(\theta^{old}, \psi^{old}) \\ \theta^{new} &= \theta^{old} - \alpha \nabla_{\theta} L(\theta^{old}, \psi^{new}) \quad \rightarrow \\ \theta^{old} &= \theta^{new}, \quad \psi^{old} = \psi^{new} \quad \text{alternating G-D} \end{aligned}$$

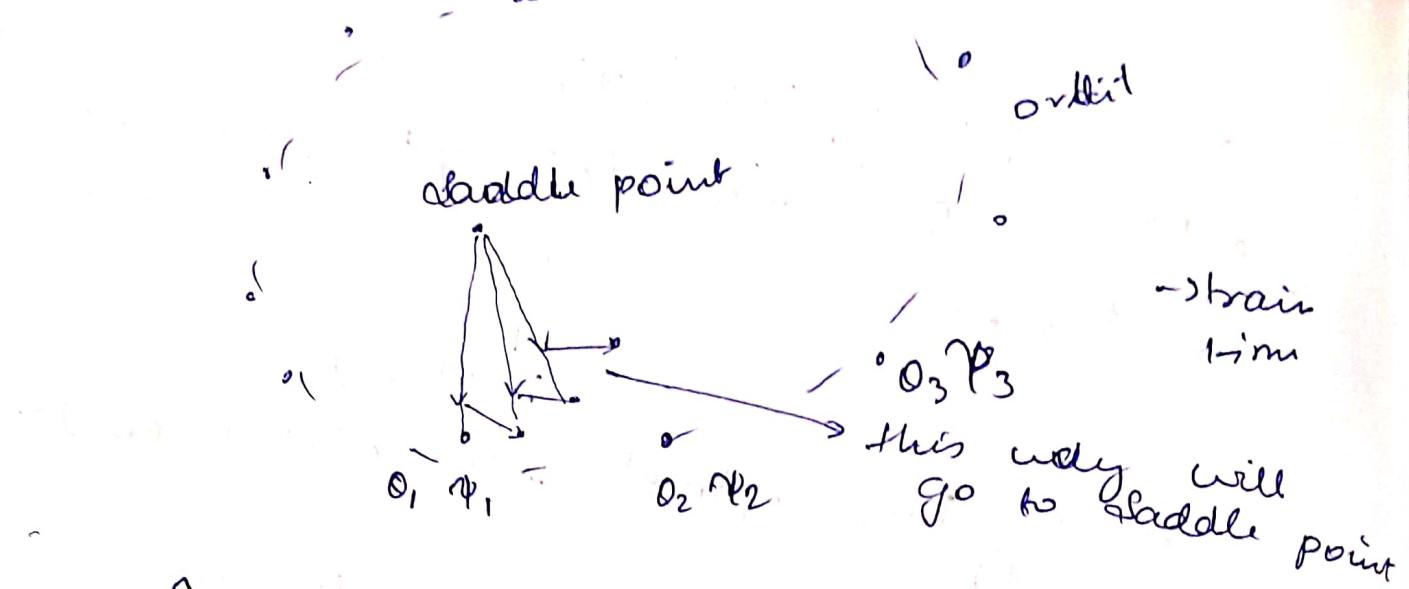
\hookrightarrow minimize lower bound on div

simplifying loss

$$\begin{aligned} \min_{\theta} \max_{\psi} E_{x \sim \delta_{\theta}} f(\psi_x) + E_{z \sim \delta_{\theta}} f(\psi_z) \\ \min_{\theta} \max_{\psi} f(\theta) + f(\psi_{\theta}) \end{aligned}$$

$$\approx \min_{\theta} \max_{\psi} f(\psi_{\theta})$$

Exponential moving average.



$\theta_{i+1} = (1-\alpha) \theta_i + \alpha \theta_i + \dots \rightarrow \text{EMA of weights averaging.}$

average all points

$$\text{Or } \frac{1}{N-n} \sum_{i=n+1}^N \theta_i$$

When to start averaging?

- have to compute the orbit
- and keep going cyclic
- but can't track this (tricky)

Soln
stop at some epoch, collect weights & take avg,

Quality Measurements

- Best metrics are hand crafted & problem specific.
(Human evaluation for e.g.)
- Metrics for Images

Inception Score (IS)

- Evaluates objectiveness & class distn
- Assumption - single object in image
- take pretrained n/w, gen samples, run n/w on samples to get prob on each class on pretrained data
- Each prediction should have low entropy (high prob for single class)
- For collection of images, marginal distn should be high.

$$\log IS = \mathbb{E}_x KL(p(y|x) || p(y))$$

$$= \sum_x p(x) \sum_y p(y|x) \log \frac{p(y)}{p(y|x)}$$

$$= \underbrace{\sum_y \log p(y) \sum_x p(y|x) p(x)}_{\text{entropy}}$$

= entropy

$$- \underbrace{\sum_{xy} p(y|x) p(x) \log p(y|x)}_{H(y|x)}$$

$$H(y|x)$$

$$= H(y) - H(y|x)$$

Frechet Inception Distance (FID)

- Measures quality of image.
- Compares activations of pre-trained convnet b/w real & generated.

$$FID = \| \mu_r - \mu_g \|^2 + \text{Tr} (\Sigma_r + \Sigma_g - 2(\Sigma_r \Sigma_g)^{-1})$$

Compare only means & pairwise correlations

Applications

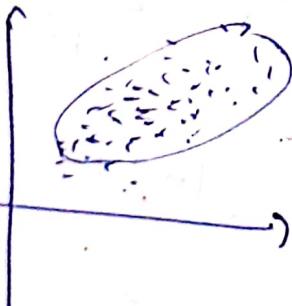
- 2k inequality img. generations (StyleGAN)
- extend datasets
- Super resolution, text2speech etc.

Normalising flows

Probabilistic view on generative modeling

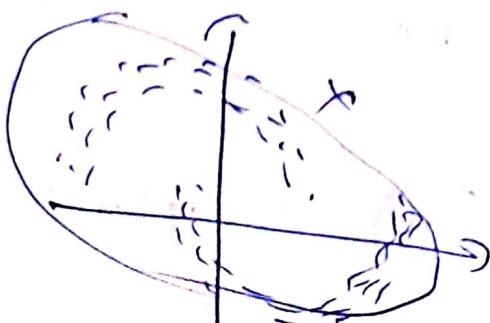
- learning dists on data

$$P_{\text{model}}(x) \approx P_{\text{data}}(x)$$



→ can use gaussian dist

$$n(x | \mu, \sigma^2)$$



→ can't use gaussian

more applications

- Probability estimate (out of distn estim distn)
- Representation learning (semi-supervised)

grad based training of probabilistic GRM

samples $x = (x_1, x_2, \dots, x_n)$ $x_i \sim P_{\text{data}}(x)$
from distn

parametric model $P_\theta(x)$

- assume model is suitable for grad opt by θ
- prob constraints (valid pdf output)
- desirable to have a tractable sampling procedure.

Train using MLE

$$\log P_\theta(x) = \sum \log P_\theta(x_i) \rightarrow \max_{\theta}$$

solved using S.G. ascent, Adam etc,

eg:

$$P_\theta(x) = N(x|\mu, \Sigma) \quad \theta = (\mu, \Sigma)$$

$$\text{MLE: } \sum_i \log N(x_i|\mu, \Sigma) \rightarrow \max_{\mu, \Sigma}$$

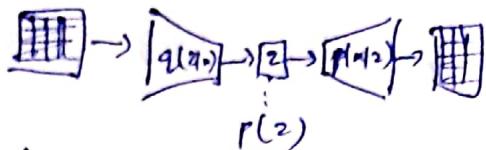
update μ, Σ

$$\mu \in [2 \ 2]$$

$$\Sigma \in [0.4 \ 0.2]$$

Choosing model for complex data.

Variational AEs
(semi-implicit
model)



$$\log p_\theta(x), \log E_{p(r|z)} p_\theta(x|r)$$

expressive and tractable
pdf

$$p_\theta(x) = \text{Normalising flow}$$

- expressive parametric model
- Tractable & exact $p_{\theta(x)}$
- fast & tractable sampling

→ unbiased estimate

Variational inference

$$\log p_\theta(x) \geq L(x, \theta) \Rightarrow \max_{\theta}$$

→ an unbiased estimate,
apply non linear $f_\theta(\log)$
then it won't remain
unbiased.

Normalising flow NF

change of variable formula

1d - case

$$\text{given } p(z) \quad ?$$

given $p(z)$

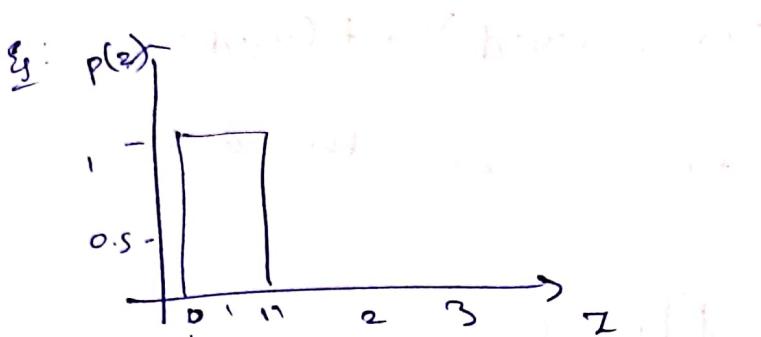
$$x = f(z)$$

$$f(z) = ?$$

$$p(x)|dx| = p(z)|dz|$$

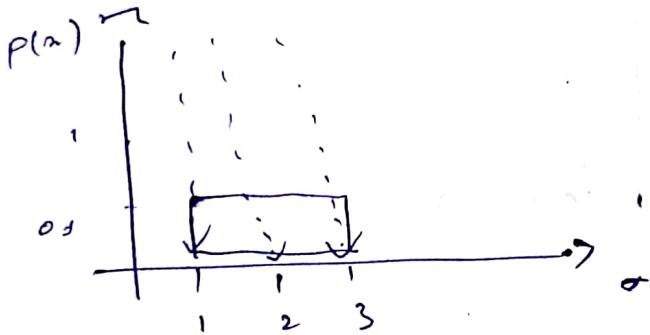
$$p(x) = p(z) \left| \frac{dz}{dx} \right|$$

$$= p(f(x)) \left| \frac{df(x)}{dx} \right|$$



$$x = 2z+1 \quad \frac{dx}{dz} = 2 \quad \frac{db(x)}{dx} = 0.5$$

$$z = \left(\frac{x-1}{2} \right)$$



prob goes down

N-d case

$$p(x) = p(f(x)) \left| \det \frac{\partial f(x)}{\partial x^i} \right|$$

Jacobian.

NF Concepts



f_0 should be reversible

$$p(x) = p(f(x)) \left| \det \frac{\partial f(x)}{\partial x^i} \right|$$

Coupling layers.

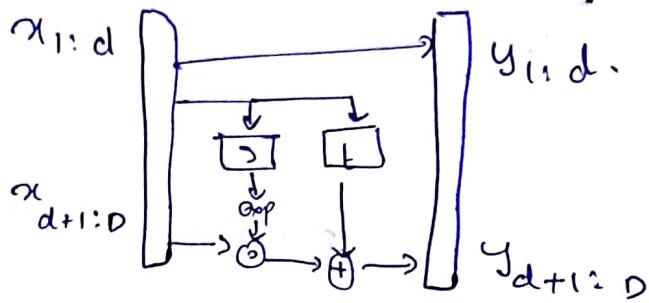
Input α , output y .

Forward:

$$y_{1:d} = x_{1:d}$$

$$y_{d+1:D} = \alpha_{d+1:D} \odot \exp(s(x_{1:d})) + t(x_{1:d})$$

s and t are conditioned on the d component



Inverse:

$$x_{1:d} = y_{1:d}$$

$$\alpha_{d+1:D} = (y_{d+1:D} - t(y_{1:d})) \odot \exp(-s(y_{1:d}))$$

\therefore it is reversible

Jacobian:

$$\frac{\partial y}{\partial x^T} = \begin{bmatrix} I_d & 0 \\ \frac{\partial y_{d+1:D}}{\partial x^T_{1:d}} & \text{diag}(\exp[s(x_{1:d})]) \end{bmatrix}$$

This is a lower Δ matrix.

\therefore det is product of diag elements

$$\det \left| \frac{\partial y}{\partial x^T} \right| = \exp \sum s(x_{1:d})_j$$

- can't stack coupling layers directly
 - so add permutations in between coupling layers
- model.

$$f_0(x) : \mathbb{R}^D \rightarrow \mathbb{R}^D$$

$$f_0(x) = f^n \circ \dots \circ f^1(x)$$

(coupling layers)

$$\log p_0(x) = \log p(f_0(x)) + \log |\det \frac{\partial f_0(x)}{\partial x^T}|$$

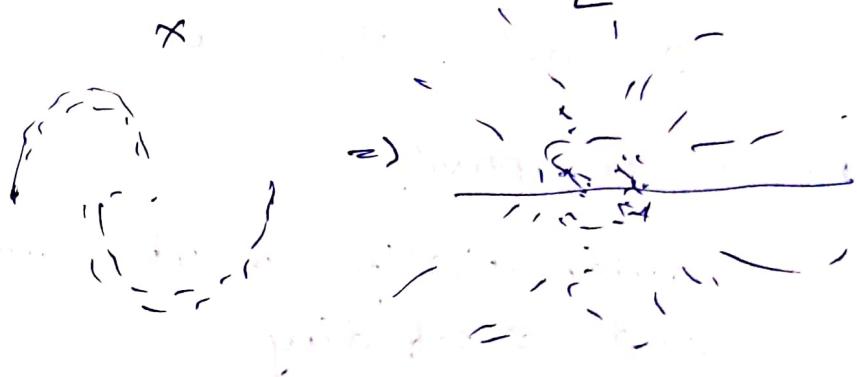
$$= \log p(f(x)) + \sum_i^n \log \left| \det \frac{\partial f^i}{\partial x^i} \right|$$

$$\log p_0(x) = \sum_i \log p(x_i) \rightarrow \max_{\theta}$$

Inference

$$x \sim p_x$$

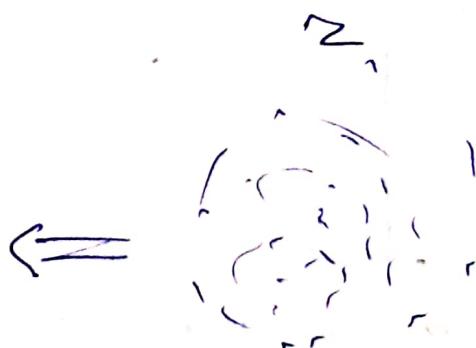
$$z = f(x)$$



generation

$$z \sim p_z$$

$$x = f^{-1}(z)$$



each coupling layer perform domain scaling and shuffling.

Measuring quality

Continuous data \rightarrow avg test log likelihood

- Normalising flow can be applied to variational inference (e.g. in VAE after encoder)

Gaussian Processes and Bayesian optimization

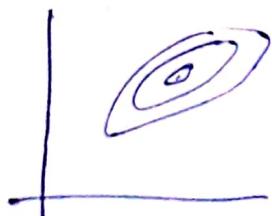
Applications

- Real valued regression
- Classification
- Ordinal regression

Model complexity:

- automatically select model complexity and overfitting.

gaussian distri



divide gaussian vector f into 2 sub vectors

$$p(f_1, f_2) \sim \mathcal{N}(f_1, f_2 | \mu, \Sigma)$$

Joint

also divide μ & Σ into μ_1, μ_2 & $\Sigma_{11}, \Sigma_{12}, \Sigma_{21}, \Sigma_{22}$

$$\begin{pmatrix} b_1 \\ b_2 \end{pmatrix} \sim N \left(\begin{bmatrix} \mu_1 \\ \mu_2 \end{bmatrix}, \begin{bmatrix} \Sigma_{11} & \Sigma_{12} \\ \Sigma_{21} & \Sigma_{22} \end{bmatrix} \right)$$

$$p(b_1) \rightarrow \int p(b_1, b_2) db_2 = N(b_1 | \mu_1, \Sigma_{11})$$

gaussian process model

- assume data is generated by a conti. stochastic process.

- assume covariance matrix is from a kernel function

$$K_2 \{k_{ij}\} = \{K\{x_i, x_j\}\}$$

- squared exponential kernel

$$K(x, x') = \sigma_b^2 \exp \left\{ - \sum_{i=1}^d \frac{(x_i - x'_i)^2}{2\tau^2} \right\}$$

a GP need not use the gaussian kernel.

- more general ex:

$$K(x, x') = \sigma_b^2 \exp \left\{ - \frac{|x - x'|}{r} \right\} + \sigma_n^2 + \sigma_e^2 \delta(x, x')$$

σ_b^2 - marginal func variance

σ_n^2 - variance of bias

σ_e^2 - noise variance

r - lengthscale

$\propto > 1$ - roughness

- Probabilistic prior assumptions in form of covariances
- new covariances can be constructed from old

GP regression

$$\text{data } S_m = \{x, y\} = \{(x_i, y_i)\}_{i=1}^m$$

Model

$$y_i = f(x_i) + \epsilon_i$$

↳ gaussian white noise
 $\mu = 0, \sigma^2$

$$f \sim GP(0, K)$$

- $f(x_i)$ s are generated by noisy version is observed.

$$\text{prior } p(f) = N(f|0, K)$$

model noise or likelihood.

$$p(y|f) = N(y|f, \sigma^2 I_m)$$

marginal likelihood

$$p(y) = \int p(y|f) p(f) df = N(0, K + \sigma^2 I_m)$$

Prediction

$$y_* = b_* + \epsilon_* \quad b_* = f(x_*)$$

$$p(y_*|f) = N(0, [K_*^\top K_*]^{-1})$$

$$k_* = \{K(x_*, x_i)\}_i^m \quad K_{**} = K(x_*, x_*)$$

$$P(f_*|y) \approx N(f_* | \mu_*, \sigma_*^2)$$

where $\mu_* = K_*^\top [K + \sigma^2 I_m]^{-1} y$; $\sigma_*^2 = K_{**} - R_*^\top [K + \sigma^2 I_m]^{-1} R_*$

Even in a GP with no parameter, there can be hyperparameter σ^2 : noise variance \in params of cov fn.

Learning GP model

- ability to choose hyperparam and cov directly from data
- minimizes negative log likelihood $L(\theta)$ wrt cov fn params \in noise θ .

$$\therefore p(y) \sim \mathcal{N}(0, K + \sigma^2 I_m)$$

then (marginal log-likelihood)

$$L = -\log p(y|\theta)$$

$$= \underbrace{\frac{1}{2} \log \det C(\theta)}_{\text{constant}} + \underbrace{\frac{1}{2} y^\top C^{-1}(\theta) y}_{\text{data-fit}} + \frac{m}{2} \log \det$$

Reg

$$C \geq K + \sigma^2 I_m$$

- minimisation of $L(\theta)$ is a. mon convex optim problem
- gradients

$$\frac{\partial L}{\partial \theta} = \frac{1}{2} \text{tr} \left(C^{-1} \frac{\partial C}{\partial \theta_i} \right) - \frac{1}{2} y^T C^{-1} \frac{\partial C}{\partial \theta_i} C^{-1} y$$

\downarrow

$O(m^3)$ $C \rightarrow mxm$

Cost can be approximated using random fourier features, special quadrature band formulas. So overall complexity reduces to linear(almost?)

Bayesian optim

global optim:

Consider $f: X \rightarrow \mathbb{R}$ $X \subseteq \mathbb{R}^d$

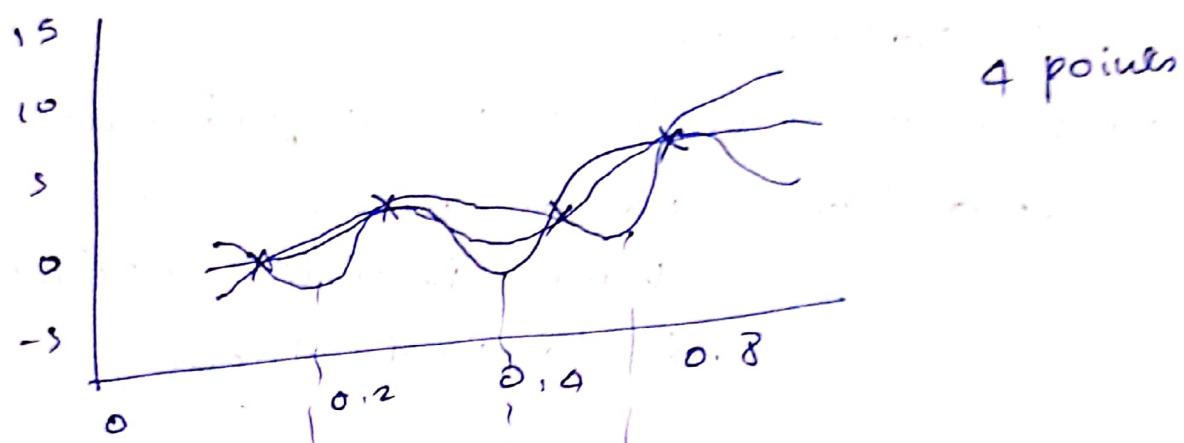
$$x_{\min} = \arg \min_{x \in X} f(x)$$

- f is multimodal
- so random search can't be done
- grads can't be computed numerically
- Real world problems often deal with this, e.g: optim of struct of aircraft, lot of params, single sim might take days.
- e.g: optim geometry of a street to mini time consuming, physical dims needed on each optim.

Ej: Modern DNNs.

- params:
 - No. of layers
 - Units per layers
 - Reg. Coeffs
 - lr, etc.
- might take a lot of time
- how to get best struct?

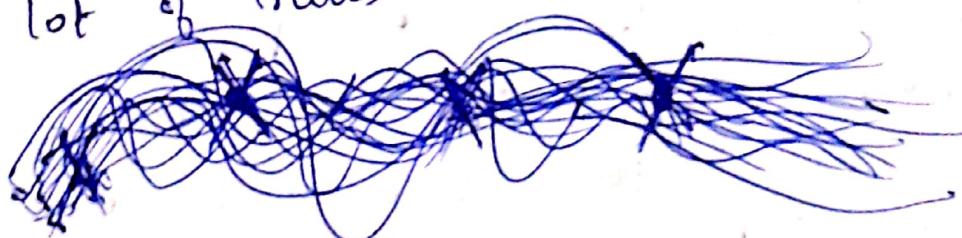
Ej: design of grafts to avoid in heart surgery



→ 3 diff minimas
- for 3 $f(x)$

the more $f(n)$ we optimise, the more locations of minima found.

after lot of runs



→ minima concentrate

→ gaussian process

① is the posterior given the minimum constraint
(but we don't consider global minima?)

- prior assumption on f .
- sample values of f .
- calculate posterior distn of position of minimum
- sample around the locations of minimum
- add these new observations & iterate
- In this, we need to support exploration and exploitation
 - exploration \rightarrow variance is large (big uncertainty)
 - exploitative \rightarrow where mean is low.
(which indicates optimum is located there)
- consider some heuristic criterion to sample points
 - On gaussian processes, if the $f(x)$ is generated by some gaussian process then if it follows the process criterion mentioned below, the process will converge to a minimum faster, no. of steps to global minimum. Speed of conv is lower than G.D.
- take $f(x)$
- sample from $f(x)$
- Using this, construct GP to approx $f(x)$

- Using one of the criterion, select the next point which will evaluate the $f(x)$. Repeat.

Criterions to select candidate:

{ own defined para
fun

$$\leftarrow \alpha_{LCB}(x) = \mu_x(x) + \{ \cdot \sigma_x(x) \}$$

↓
- ve sign because maximum μ

- have to find real minima

- $\sigma(x)$ is the uncertainty estimate to add the error in prediction
- can be considered as a confidence interval for ^{predictor} position of min value.

[CnP: upper(lower) Confidence Band]

→ Expected Improvement

- prob criterion

- estimates avg value of improv over existing obs.

$$\epsilon(x) = y_{best} - \mu_x(x), y_{best} = \min_{i=1, \dots, m} y_i,$$

↓
have to mini this

$$\alpha_{EI}(x) = \int \max(0, y_{best} - y_x) p(y_x|x) dy_x ?$$

we don't care about -ve

- maximum prob of improvement

$$r(x) = \frac{u(x) - y_{\text{best}}}{\sigma(x)} \rightarrow \text{normalised } u(x)$$

$$\alpha_{MP1}(x) = P(f(x) < y_{\text{best}}) = \phi(gV(x))$$

Bayesian

Instead of optim initial fn,
construct acquisition fn

$$x_{t+1} = \arg \max_{x \in X} \alpha(x|s_t)$$

$\alpha(x) \rightarrow$ not so expensive

grads are typically available

- no need to find global optim for many cans, just need better dir to getting to a local opt fn

Tuning hyper

param of SVM

n samples 2500

n feature = 43

n interparams f = 18

n redundant f = 5

- optimise wrt $g(x) = (c, \gamma)$

c - penalisation

γ - kernel width

- target fn $L(x)$ is AUC ; 3 fold cv

- put random values of C, σ^2 , then estimate AUC
- Construct gaussian process by finding minima obtained
- optimise on a criterion (Expected Improv) to obtain point for next step. Repeat

About GP?

- Simple to use
- not many parameters
- can estimate uncertainty

Deepbayes (continued)

Deep Gaussian Processes

linear models

linear combination of basis fun

$$f(x) = w^T \phi(x)$$

$$\phi(x) = (\phi_1(x), \dots, \phi_d(x))^T$$

- In regression, obs' are modelled as

$$p(y|w, x, \sigma) = N(\phi_w(x), \sigma^2 I)$$

$$\phi_w = \phi(x) = \begin{bmatrix} \phi_1(x_1) & \dots & \phi_d(x_1) \\ \vdots & \ddots & \vdots \\ \phi_1(x_n) & \dots & \phi_d(x_n) \end{bmatrix}$$

$\phi_w \rightarrow$ mean of latent fun.

$$x = (x_1, \dots, x_n)^T$$

$$y = (y_1, \dots, y_n)^T$$

$$w = (w_1, \dots, w_d)^T$$

quadratic loss \Rightarrow maximising likelihood

given a $p(w)$ prior, we learn $p(w|y, x)$

$$p(w|y, x) = \frac{p(y|x, w)p(w)}{\int p(y|x, w)p(w)dw}$$

$$\text{let } p(w) \sim N(0, \sigma^2)$$

- prior of f is again gaussian

$$p(f) = N(0, K)$$

$$K = \text{Cov}(f) = E[\phi \omega \omega^T \phi^T] = \phi \Sigma \phi^T$$

we can solve either by

$$p(y|\omega, x, \lambda) = N(\phi\omega, \lambda I) \quad p(\omega) = N(0, \sigma^2)$$

or

$$p(y|\omega, x, \lambda) = N(f, \lambda I) \quad p(f) = N(0, \sigma^2)$$

(prior in terms of f)

- Bayesian linear regression is a solved problem.
we can analytically compute posterior.

- But, we need to specify basis function

- Now we can work with functions
such as cov K , instead of parameters.
Using this, we can choose K from
inputs and evaluate ϕ instead of
other way around. K can be
selected such that ϕ is infinite dim

$$\text{Eg: } k(x_i, x_j) = \exp\left(-\frac{\|x_i - x_j\|^2}{2}\right)$$

\exp doesn't allow disentanglement
of relation b/w x_i & x_j .

Taylor expansion is needed.
This turns out to be
scalar product of things that
go to infinite dim.

- There are other kernels satisfying this property.
- In Bayesian LR,

prior is σ , likelihood is $\ln L$.

posterior must be σ .

over ω : - posterior $p(\omega | x, y, \lambda) \sim N(\mu_\omega, \Sigma_\omega)$

$$\text{Cov } \Sigma_\omega = \left(\frac{1}{\lambda} \phi^T \phi + S^{-1} \right)^{-1}$$

mean

$$\mu_\omega = \frac{1}{\lambda} \Sigma_\omega \phi^T y$$

Predictions

$$p(y_{*} | x, y, x_*, \lambda) \sim N(\psi(x_*)^T \mu_\omega,$$

can be rewritten using $K = x + \psi(x_*)^T \Sigma_\omega \psi(x_*)$

over f :

$$p(f | x, y, x) \sim N(\mu_f, \Sigma_f)$$

$$\Sigma_f = K - K(K+I)^{-1} K$$

$$\mu_f = K(K+I)^{-1} y$$

$$p(y_* | x, y, x_*, \lambda) \sim N(K_x^T (K+I)^{-1} y, \\ \lambda + K_x^T (K+I)^{-1} K_x)$$

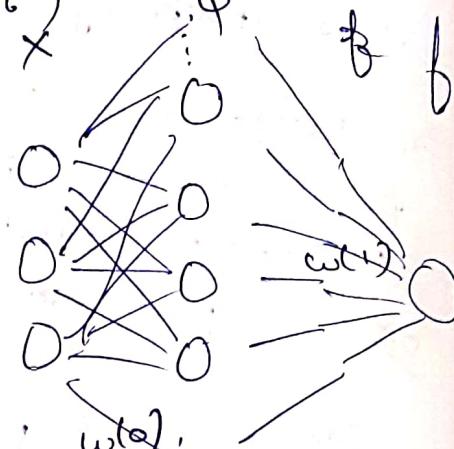
Both are equivalent.

$\psi(\cdot)$ costs $O(D^2)$ mem, $O(D^3)$ time

$K(\cdot, \cdot)$ costs $O(N^2)$ mem, $O(N^3)$ time

gaussian processes as infinitely wide
Shallow nn (1996)

$$\begin{aligned} \{\omega^{(0)}\}_{ij} &\sim \mathcal{N}(0, \alpha_0) \\ \{\omega'\}_{ij} &\sim \mathcal{N}(0, \alpha_1) \end{aligned}$$



- Central limit theorem implies that f is gaussian

- f has zero mean

$$\begin{aligned} \text{cov}(f) &= \mathbb{E}_P(\omega^{(0)}, \omega^{(1)}) [\phi(x\omega^{(0)}) \omega^{(1)} \phi(x\omega^{(0)})^T] \\ &\geq \alpha_1 \mathbb{E}_P(\omega^{(0)}) [\phi(x\omega^{(0)}) \phi(x\omega^{(0)})^T] \end{aligned}$$

- Some choices of ϕ leads to analytic expression of known kernels.

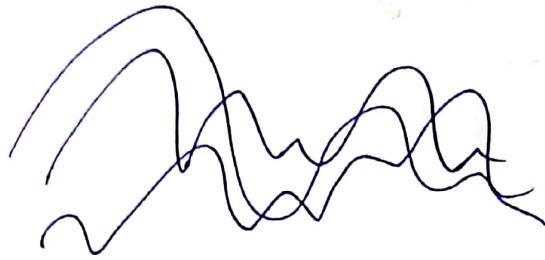
Challenges & Limitations on GP

- Kernel design $k(\cdot; \theta)$ (before it was born (can); have parameters for more flexibility)
- GPs can be too expensive
- GPs might not be tractable (say for classification)
↓
need approx.

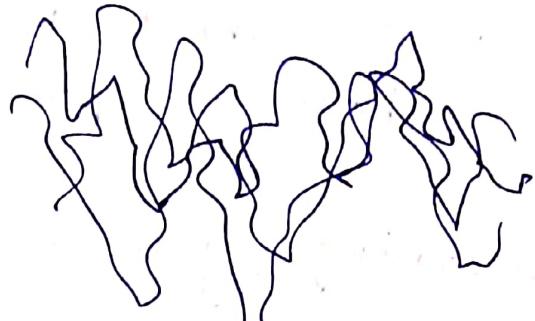
Eg:

$$p(y|x, \theta) = \int p(y|f) p(f|x, \theta) df$$

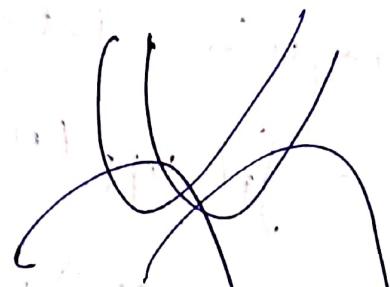
Can't compute (no conjugacy) \approx sum other distn gaussian



gaussian



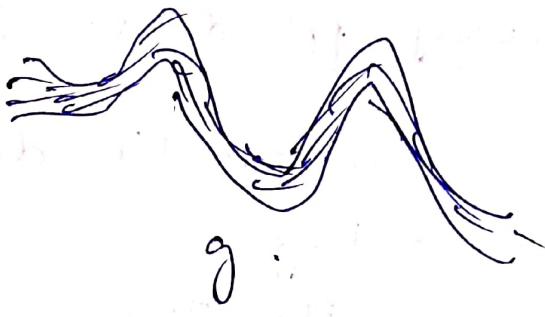
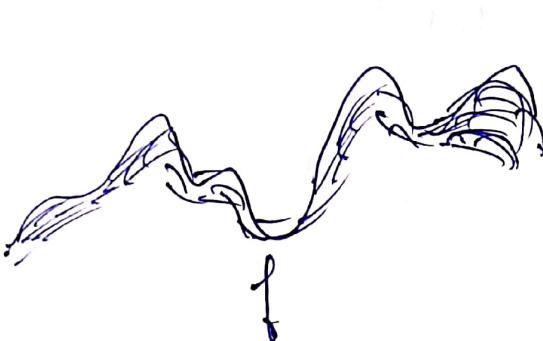
Matérn



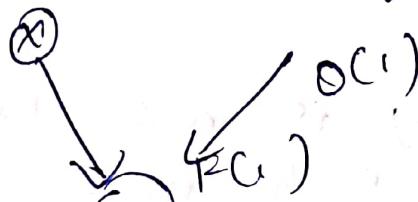
polynomial

Deep GP for large Rep Power

- composition of processes



(x)



$\xrightarrow{O(1)}$



$\xrightarrow{F(2)}$

$\xrightarrow{O(2)}$

y

- increase the rep power of GP
- also complicate the model

Pathologies of deep gaussian process

- give input to later layers to help better representation.

deep learning x gaussian process

- Inference requires integral Calc.

$$p(y|x, \theta) = \int p(y|f^{(N_h)}, \theta^{(N_h)}) x$$

$$p(f^{(N_h)} | f^{(N_h-1)}, \theta^{(N_h-1)}) x$$

$$\dots \times p(f^{(1)} | x, \theta^{(0)}) df^{(N)}_{df}$$

why ppl do DL & not GP?

- automatic diff.
- Reg
- GPUs, TPUs
- mini batch
- Applications specific reps (CNNs)
- In GPUs, mini batching is really hard to do, need to invert K to mini batch.

Main ingredients for Inference for DGP's

posterior $p(f|x, y, \theta) \approx q(\theta)$

(can be non gaussian as likelihood
might not be gaussian)
(so approx it with a gaussian)

other approx which deal
approx to gaussian

- Laplace approx
- Expectation propagation
- Variational Bayes
- MC or L

- Now, Scale model to lot of datapoints

Sparse GP \Rightarrow Nyström approx

introduce M new latent
variables u at location z .

$$K \approx K_x z^T z z^T K_{zz}$$

Fully Independent Training
Conditionals.

diagonal of K is not correct, so
- partially independent training
Conditionals ($P(z|c)$ is done)

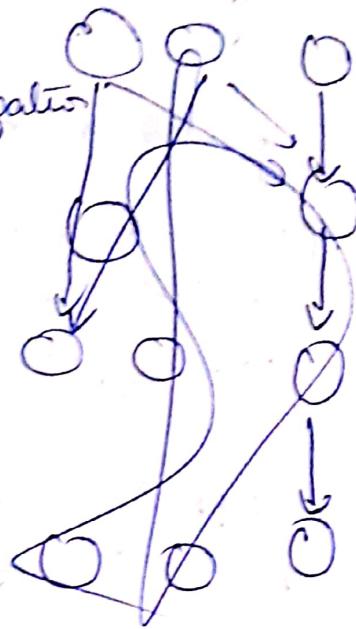
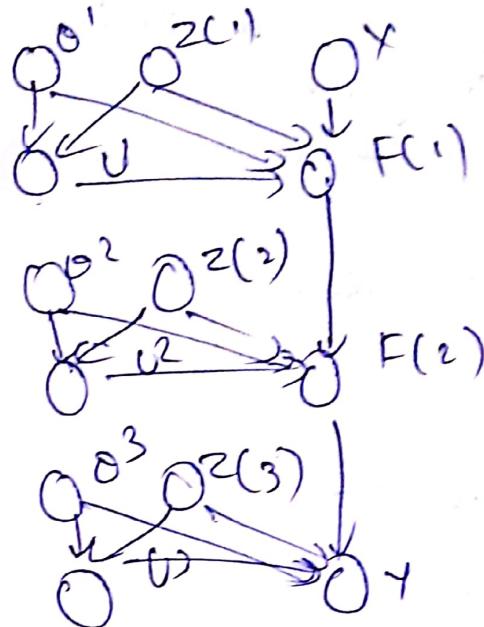
- Variationally sparse GP
- introduces distri over f and u
 $q(f, u) \approx p(f|u) q(u)$

- Random feature expansion

- Kronecker (Toeplitz) Tensor Structured LFs.

Induction

- Inference for DGPs
scalar expectation propagation

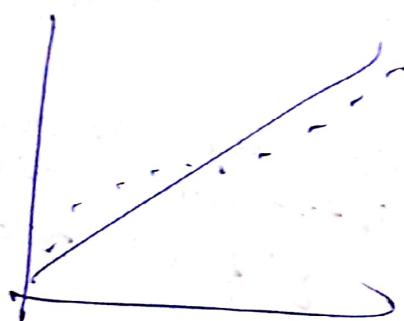


- new random variables z at every layer.
- hyperparameters for kernels & latent vars

- based on Fourier transform due to shift invariance
- DGPs with random features become DNNs
- can be learnt using stochastic variational

Calibration as a means of Quantification

of Uncertainty



\rightarrow for deep CNNs

Post Calibration fixes it.

what is the principled way of developing classifiers, which by construction calibrated

- gaussian process + CNN.

CNN layers \rightarrow DGP \rightarrow OLP

- Bayesian CNNs are calibrated

- Need to be bayesian w/ filters

- have to regularize filters

- Eg: Using Monte Carlo Dropout

another method

- Replace linear with Random Feature Expansions

- some new trends mentioned.

Adaptive skip grams

Representation learning

- dependent on features

- K-means might give similar colour code not similar semantic code.

Rep h for words

- goal is to learn close embeddings for semantically similar words.

- Skipgram.

- learn ilp embed & olp embed

- take current word x and take context words y (words surrounding x)

- $P(y_i|x) = \prod_j P(y_{ij}|x)$
- assume they are independent
- based on distributional hypothesis:
similar words appear on similar context.
- $\Theta \leftarrow \Theta + \eta \nabla \log P(y_i|x_i)$
then pick most in Θ pre-
convergence
- learns rich embeddings
- sparse gradients
- very efficient parallel training.

Is it good enough?

- for some words, only one meaning was learnt.
- for some words, diff meanings might get mixed up in one vector.

Why is it happening?

- word ambiguity

Soln:

- latent variable model

training using EM

Observed $D = \{(x_i, y_i)\}_1^N$

Latent $Z = \{z_i\}_1^N$, each meaning in local context, local

Prm: $\Pi_z = \{\pi_w\}_1^V$, meaning probability, global

$\Theta = \{in_{w,k}\}_{w=1, k=1}^{V, K} \cup \{out_w\}_1^V$ (embeddings)

target $\log p(D|\theta)$

$$\approx \log \int p(\pi|\alpha) \int p(z|\pi) \prod_{i=1}^N p(y_i|z_i, \theta) dz_i \rightarrow \max_{\theta}$$

↓
intractable.

So consider a fully factored posterior approx

$$q(z, \pi) \approx \pi q(\pi_w) \prod_i q(z_i) \propto p(z, \pi | D, \theta)$$

then introduce variational bound.

Then,

E-sup

- assume current approx of $q(z)$

- solve $a_k(z) = \mathbb{E}_{q^*}(z_i) = \arg \max_{q(z)} \lambda(q(\pi), q(z))$

- Perform word dem. disambiguation

M-sup

- stochastic update of global params

- Due to conjugacy $q(\pi_w) \sim \text{Dirichlet}(\gamma_{w_1}, \dots, \gamma_{w_k})$

- VLB has a finite-dim numeric parameterization.

$$\lambda^*(\gamma, \theta) = \mathbb{E}_q \left[\underbrace{\log p(\pi|k) - \log(\pi|v)}_{\rightarrow k \sim \text{Dirichlet}} + \right]$$

$$\sum \sum q^*(z_i, k) \log p(y_i | z_i, k, \theta)$$

$\underbrace{\quad \quad \quad}_{\text{weighted Stig gram pred}}$

$$+ H(q^*(z))$$

$\underbrace{\quad \quad \quad}_{\text{constant}}$

SVI for Multi-meaning Skipgram model

*Stochastic
variational
inference*

- Sample a context: $(x, y) \sim D$

- Disambiguate word meaning

$$q(z) \propto \exp(E_{\pi(x)} \log(z)) + \log P(y|z, \theta)$$

- update word embeddings

$$\theta \leftarrow \theta + \eta \nabla \sum_{k=1}^K q(z=k) \log P(y|x, z=k, \theta)$$

- update params of $q(\pi_x | \psi_n)$

$$\gamma_{x,k} \leftarrow \gamma_{x,k} + \underbrace{\eta (q(z=k) \pi_x - \gamma_{x,k})}_{\text{natural grad.}}$$

Choosing number of meanings.

- Fix K for each word

(not very good)

- Use a heuristic, based on word freq
(bit non-trivial)

- external knowledge, speech tagging

Given a distri of points, how many clusters?

traditional model selection

- define a finite mixture model

$$\pi(\alpha) \sim \text{Dir}(\alpha_1, \dots, \alpha_K)$$

$$\phi_k \sim H, k = 1, \dots, K$$

(cluster
params)

$H \rightarrow$ prior distri over cluster params

$$\mu, \Sigma \sim P_H(\mu, \Sigma)$$

$$z_i | \Pi \sim \text{Categorical}(\pi)$$

$$\alpha | z_i, \phi \sim p(\alpha | \phi_i)$$

(for each data point,
choose a cluster
from categorical
distri)

- For each K measure $\log p(\mathbf{x}_{\text{val}} | \boldsymbol{\theta}_{\text{min}})$
- This is computationally inefficient.
- Not very elegant (have to make sure val data comes from same distn)

Model selection with Dirichlet priors

- Define a discrete mixture model

$$\mathbf{b}_k = \sum_{i=1}^{\infty} \pi_k \delta(\phi - \phi_{ik}) \quad \begin{matrix} \text{infinite number} \\ \text{of cluster param.} \\ \text{cat with } i \text{ weight.} \end{matrix}$$

$$p(\mathbf{x}|\mathbf{b}) = \sum_i \pi_k p(\mathbf{x}, \mathbf{b}_k)$$

- choose a nonparametric prior.
- $G|H, \alpha \sim DP(H, \alpha) \xrightarrow{\text{posterior over cluster param}} \text{scalar value.}$
- \hookrightarrow Dirichlet priors

- Select the model one-shot

$$p(\mathbf{x}|\mathbf{X}) = \underbrace{\int p(\mathbf{b}|\mathbf{x}) p(\mathbf{x}|\mathbf{b}) d\mathbf{b}}_{\text{model}}$$

$$p(\mathbf{b}|\mathbf{x}) \propto \underbrace{p(\mathbf{b}) p(\mathbf{x}|\mathbf{b})}_{\text{model selection}}$$

What is a good non parametric prior?

- should allow potential infinite number of clusters.
- for a finite dataset, then no of clusters $K < n$.
- Model complexity to grow with more data

⇒ two representations of Dir. Priors.

- Chinese Restaurant Process
- Stick breaking process

Chinn Restaurant Probs

- imagine you are first customer & restaurant is empty. Take a table.

- next customer has a choice, to sit with you on same table or pick a new table.

- so at every time step, new customer comes and has a choice

$$p(z_i = k | z_{\leq i}) \propto \begin{cases} n_k, & k \leq K \\ \alpha, & k = K+1 \end{cases}$$

$$p(z) \propto \prod_i p(z_i | z_{\leq i})$$

- treat them z as cluster assignments.

- each table is a cluster

- for each non empty table, sample cluster params from prior π_1 .

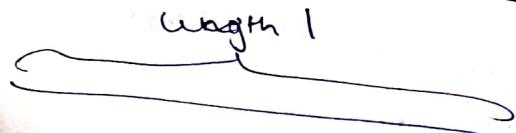
- Then for generativity, sample from corresponding mixture parameter component

- So here, we need to replace old prior over word over mixing prob with the CRP prior.

- But state is not anymore IEP: Because z_i are all conditioned on each other.

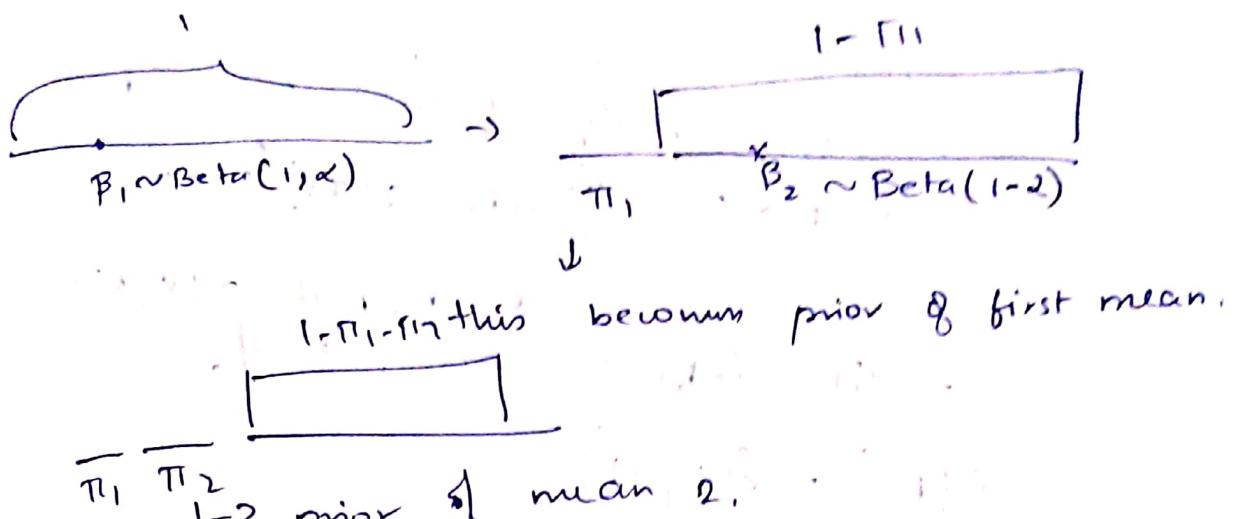
- So not convenient for SVI to train.

Stick Breaking probs.



On beginning, then break whole stick, then divide where to break it by ...

breaking proportion from Beta distri. $\beta(1, \kappa)$.
 $\beta \in (0, 1)$



Repeat this infinitely.

$$\beta_k \sim \text{Beta}(1, \alpha) \quad k \rightarrow \infty$$

$$\phi_n \sim H, \quad k = 1, 2, \dots, \infty$$

π_k prior prob. of cluster.

$$\pi_k = \beta_k \prod_{t \neq k} (1 - \beta_t), \approx \beta_k (1 - \sum_{t \neq k} \pi_t)$$

$z_i | \pi \sim \text{Categorical}(\pi)$

$$x_i | z_i \sim p(x_i | \phi_k)$$

$$G(z) = \sum_{k=1}^{\infty} \pi_k \delta(\phi - \phi_k) \sim DP(H, \alpha) \quad z \sim CRP(\alpha)$$

In CRP, cluster assignment was sequential.

In SBP, it is done independently.

Adaptive - skipogram

- based on infinite mixture.

- prior over infinite many means

- automatic model selection.

- control of means granularly.

to compute it practically:

Variational inference for SBP

- finite dim var approx

- assume fully factored var approx

$$q(\beta, \phi, z) = \prod_{i=1}^{\infty} q(\beta_i) q(\phi_i) q(z_i)$$

$$\approx p(\beta, \phi, z | x, k)$$

- still infinitely many β, ϕ .

$$q(\beta_T) \approx \delta(\beta_T - 1)$$

$$q(\beta_k) \approx \text{Beta}(1/\alpha), q(\phi_k) \approx p_H(\phi_k), k \geq T$$

prior.

We don't care after T

- for all $k > T$, posterior over β_k is

the same as prior

and for $k > T$, posterior over ϕ is
same as prior

- Due to prior conjugacy:

$$q(\beta_k) \approx \text{Beta}(\beta_k | a_k, b_k)$$

lower bound decomposes

$$L(q) = \mathbb{E}_q[\log p(x | \beta, \pi, z) - \sum_{k=1}^T KL(q(\beta_k) || p(\beta_k)) - \sum_{k=1}^T KL(q(\phi_k) || p(\phi_k))]$$

$$\mathbb{E}_q = \sum_{k=1}^T \mathbb{E}_{\beta_k, \phi_k}$$

$$- \sum_{k=1}^T KL(q(\beta_k) q(\phi_k) || p(\beta_k) p(\phi_k))$$

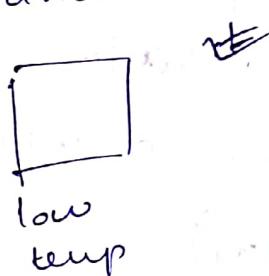
$$- KL(q(z) || p(z | \beta))$$

Markov Chain Monte Carlo

- problem with ELBO maximisation, $q(\theta|z)$ could be a poor approx for the true posterior.
- There are cases where good approx are necessary.

Ising model

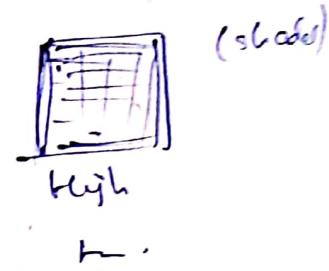
- model from statistical physics that try to describe magnetic properties of solids depending on temp. In a short interval of temperature, magnet loses its properties (phase transition effect)
- Variational inference for Ising model variational approx



low
temp

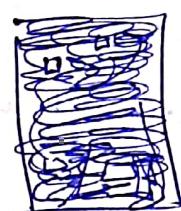


Critical
temp



High
temp

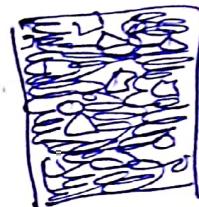
from true posterior



low temp



Critical
Temp



High
temp

Variational proxy is too poor for this model

- If variational proxy is more complex, it's not efficient to estimate

MCMC

- able to model exactly and to estimate needed statistics using this
- going to construct a sample from the distribution.

given $p(x) \rightarrow \frac{1}{Z} \tilde{p}(x)$, $Z = \int \tilde{p}(x) dx$.

MCMC $\rightarrow (x_1, x_2, \dots, x_m) \sim p(x)$ using only $\tilde{p}(x)$

given these samples,

$$\mathbb{E}_{p(x)} f(x) \approx \int f(x) p(x) dx \underset{\text{Statistics}}{\approx} \frac{1}{m} \sum_{i=1}^m f(x_i)$$

- drawing samples using some
Specially construct markov chains,
Transition probabilities $q(x|y)$.

- generate x_1 from a initial dist'n $p_0(x)$
- x_2 from $q(x|x_1)$
- x_3 from $q(x|x_2)$

x_1, x_2, \dots are not independent, but
can still be used to estimate $\mathbb{E}_{p(x)} f(x)$

- first sample can be really far from true
- In practice, first set of samples are dropped.

How to converge?

Markov chain properties

- invariant distribution

given a sample from true distri., transitioning using a transition prob., it should give sample from our distribution.
In practice, this is ensured by sampler.

- ergodic

initial $p_0(x)$.

$p_i(x)$, i'm sup of a MC.

Then

$$p_i(n) \rightarrow \pi_i(x) \quad , \quad i \rightarrow \infty . \quad \forall p_0(x)$$

gibbs sampling

generally $p(n) \propto \prod p(x_i)$ $x_i(x_1, \dots, x_n)$

On this, consider iD conditionals

$$x_i^{\text{new}} \sim p(x_i | x_1, \dots, x_{i-1}, x_{i+1}, \dots, x_n) \quad (\text{first cord of new point})$$

$$x_2^{\text{new}} \sim p(x_2 | x_1^{\text{new}}, x_3, \dots, x_n) \quad (\text{second cord of new point})$$

$$x_k^{\text{new}} \sim p(x_k | x_1^{\text{new}}, x_2^{\text{new}}, \dots, x_{k-1}^{\text{new}})$$

$$x^{\text{new}} \approx \{x_1^{\text{new}}, x_2^{\text{new}}, \dots, x_n^{\text{new}}\}$$

can be proved about its invariance.

- if all $p(x_i | x_{-i}) > 0$, then ergodic

gibbs sampling : conclusions

- for discrete, conti, very general
- no params for tuning
- inefficient for high dims.

Metropolis - Hastings Sampling

goal is to generate $p(x) = \frac{1}{Z} \tilde{p}(x)$

- introduce proposal distri $r(x|y)$.
(this can be anything)

- to make it invariant, add a rejection step.

\Rightarrow sample trial $\sim r(x|x_{\text{old}})$

- accept or reject

$$A \approx \min \left(1, \frac{\tilde{p}(x_{\text{trial}}) r(x_{\text{old}}|x_{\text{trial}})}{\tilde{p}(x_{\text{old}}) r(x_{\text{trial}}|x_{\text{old}})} \right)$$

$\rightarrow x_{\text{new}} = x_{\text{trial}}$ with prob A

$x_{\text{new}} = x_{\text{old}}$ with prob 1-A

Metropolis Sampling

proposal is symmetric

$$r(x|y) \approx r(y|x)$$

$$A \approx \min \left(1, \frac{\tilde{p}(x_{\text{trial}})}{\tilde{p}(x_{\text{old}})} \right)$$

- generally $x_{\text{trial}} \sim r(x|x_{\text{old}})$

- if $\tilde{p}(x_{\text{trial}}) \geq \tilde{p}(x_{\text{old}})$ $x_{\text{new}} = x_{\text{trial}}$

else $x_{\text{new}} = x_{\text{trial}}$ with prob $\frac{\tilde{p}(x_{\text{trial}})}{\tilde{p}(x_{\text{old}})}$

Hamiltonian Monte Carlo

- particular case of Metropolis-Hastings but with good proposal distri.
- In MH, τ doesn't come out from desired distri

Hamiltonian equation

consider simple physical setting,

$x \rightarrow$ particle

$p \rightarrow$ momentum

$$p = m \frac{dx}{dt}$$

$U(x) \rightarrow$ potential E

$$K(p) = \frac{p^T p}{m} - \text{kinetic } E$$

$$H(x, p) = U(x) + K(p)$$

Hamiltonian

Hamiltonian eqs:

$$\frac{dx_i}{dt} = \frac{\partial H}{\partial p_i}$$

$$\frac{\partial p_i}{\partial t} = \frac{\partial H}{\partial x_i}$$