

## Week 1

13. 08. 2018

### Introduction to Bayesian methods and Conjugate priors

#### (1) Think Bayesian

Principle 1: Use prior knowledge

Principle 2: Choose answer that explains observations the most.

Principle 3: Avoid making extra assumptions

#### (2) Review of probability

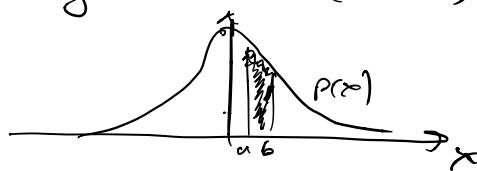
Random variables  $\xrightarrow{\text{discrete}}$   $\xrightarrow{\text{continuous}}$

#### Discrete: Probability Mass Function (PMF)

$$P(X) = \begin{cases} 0.2 & X=1 \\ 0.5 & X=3 \\ 0.3 & X=7 \\ 0 & \text{otherwise} \end{cases}$$

#### Continuous: Probability Density Function (PDF)

$$P(X \in [a, b]) = \int_a^b p(x) dx$$



## Independence

$$X \perp Y \Leftrightarrow P(X, Y) = P(X)P(Y)$$

joint                      marginals

## Conditional probability

$$P(X|Y) = \frac{P(X, Y)}{P(Y)} \begin{matrix} \text{- joint} \\ \text{- marginal} \end{matrix}$$

conditional

## Chain rule

$$P(X, Y) = P(X|Y)P(Y)$$

$$P(X, Y, Z) = P(X|Y, Z) \cdot P(Y|Z) \cdot P(Z)$$

$$P(X_1, \dots, X_n) = \prod_{i=1}^n P(X_i | X_1, \dots, X_{i-1})$$

## Sum rule

$$P(X) = \int_{\mathbb{R}} p(x, y) dy$$

## Bayes theorem

$\Theta$  - parameters     $X$  - observations

$$P(\theta|X) = \frac{P(X|\theta)}{P(X)} = \frac{P(X|\theta)P(\theta)}{\int P(X)}$$

posteriorium      evidence      likelihood      prior /

### ③ Bayesian approach to statistics

frequentist

objective

$\theta$  is fixed

$X$  is random

$|X| \gg |\theta|$

max likelihood

$$\hat{\theta} = \arg \max_{\theta} P(X|\theta)$$

Classification

$$P(\theta | X_{tr}, y_{tr}) = \frac{P(y_{tr} | X_{tr}, \theta) P(\theta)}{P(y_{tr} | X_{tr})}$$

↑ training

$$P(y_{ts} | X_{ts}, X_{tr}, y_{tr}) = \int P(y_{ts} | X_{ts}, \theta) P(\theta | X_{tr}, y_{tr}) d\theta$$

↑ prediction

Regularization

Bayesian

subjective

$\theta$  is random

$X$  is fixed

for any  $|X|$

bayes theorem:

$$P(\theta | X) = \frac{P(X | \theta) P(\theta)}{P(X)}$$

$$P(\theta|X) = \frac{P(X|\theta) P(\theta)}{P(X)} \leftarrow \text{regularizer}$$

## On-line learning

$$P_{ic}(\theta) = P(\theta|x_k) = \frac{P(x_k|\theta) P_{r-1}(\theta)}{\left[ P(x_k) \right]}$$

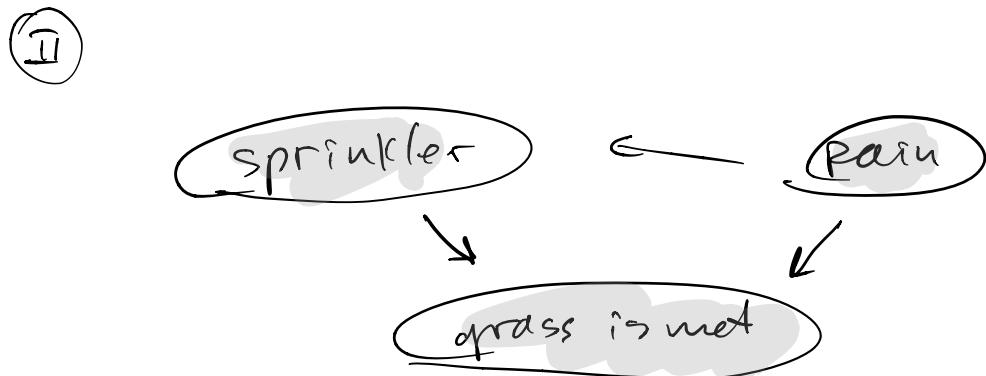
new prior                      posterior                      likelihood                      prior

## ④ How to define a model

### Bayesian network

Nodes: random variables

Edges: direct impact

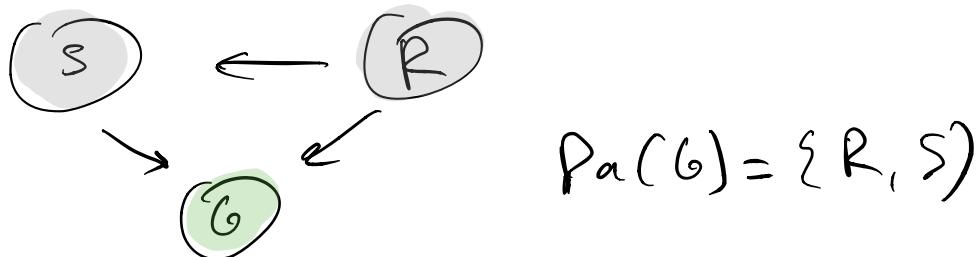


## Probabilistic model from BN

model: joint probability over all variables

$$P(X_1, \dots, X_n) = \prod_{k=1}^n P(X_k | \text{Pa}(X_k))$$

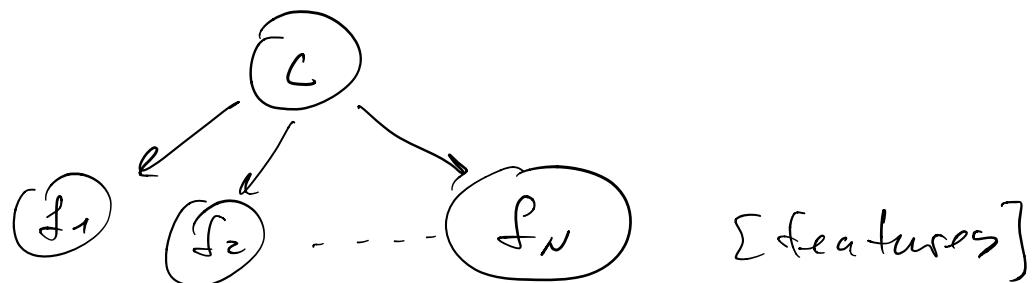
↑  
parent



$$P(S, R, G) = P(G|S, R) \cdot P(S|R) \cdot P(R)$$

probability model.

## Naïve Bayes classifier



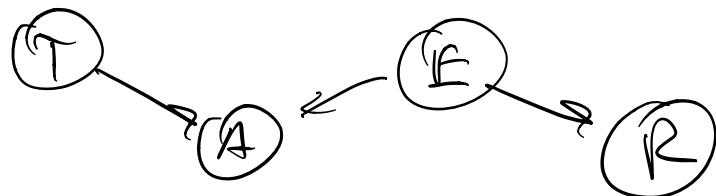
$$P(c, f_1, \dots, f_n) = P(c) \prod_{i=1}^n P(f_i | c)$$

↑  
naïve



## ⑤ Example: thief & alarm

Model:



$$P(t, a, e, r) = P(t) \cdot P(e) \cdot P(a|t, e) \cdot P(r|a)$$

Distributions

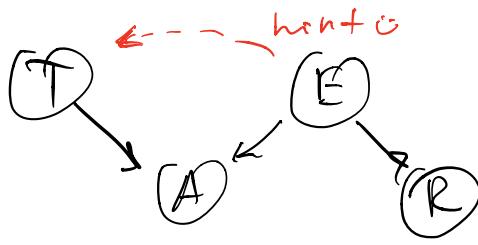
Priors

$P(t=1)$	$10^{-3}$
$P(e=1)$	$10^{-2}$

$P(a=1 t, e)$	$e=0$	$e=1$
$t=0$	0	$1/10$
$t=1$	1	1

$P(r a)$	$a=0$	$a=1$
$a=0$	0	$\pi_1$
$a=1$	$\pi_2$	1

$$\begin{matrix} T=1 & T \\ T=0 & \bar{T} \end{matrix}$$



Priors	
$P(T)$	$10^{-3}$
$P(E)$	$10^{-2}$

$P(A T, E)$		$\bar{E}$	$E$
$\bar{T}$	$T$	0	$1/10$
1	2		1

$P(R E)$	
$\bar{E}$	$E$
	$1/2$

$$\begin{aligned}
 P(T|A) &= \frac{P(T, A)}{P(A)} = \frac{P(T, A, E) + P(T, A, \bar{E})}{P(A, E) + P(A, \bar{E})} = \\
 &= \frac{\cancel{P(T, A, E)} + \cancel{P(T, A, \bar{E})}}{\cancel{P(T, A, E)} + \cancel{P(T, A, \bar{E})} + \cancel{P(\bar{T}, A, E)} + \cancel{P(\bar{T}, A, \bar{E})}} \quad (=)
 \end{aligned}$$

$$\begin{aligned}
 P(T, A, E) &= P(A|T, E) \cdot P(T, E) = \\
 &= P(A|T, E) \cdot P(T|E) \cdot P(E) = \\
 &= P(A|T, E) \cdot \underset{10^{-3}}{P(T)} \cdot \underset{10^{-2}}{P(E)} = 10^{-5}
 \end{aligned}$$

$$\begin{aligned}
 P(\bar{T}, A, \bar{E}) &= P(A|\bar{T}, \bar{E}) \cdot P(\bar{T}, \bar{E}) = \underset{(1-10^{-3})}{P(A|\bar{T}, \bar{E})} \underset{(1-10^{-2})}{P(\bar{T}, \bar{E})} \\
 &= P(A|\bar{T}, \bar{E}) \cdot P(\bar{T}|\bar{E}) P(\bar{E}) = P(A|\bar{T}, \bar{E}) \cdot P(\bar{T}) P(\bar{E}) = 0
 \end{aligned}$$

$$\begin{aligned}
 P(T, A, \bar{E}) &= P(A|T, \bar{E}) \cdot P(T, \bar{E}) = \\
 &= P(A|T, \bar{E}) \cdot P(T|\bar{E}) P(\bar{E}) =
 \end{aligned}$$

$$= P(A|T, \bar{E}) \cdot P(T) \cdot P(\bar{E}) = 1 \cdot 10^{-3} (1 - 10^{-2}) = \\ = 10^{-3} - 10^{-5} = 0.001 - 0.00001 = 99 \cdot 10^{-6}$$

$\approx 50\%$ .

$$P(T|A, R) = \frac{P(T, A, R)}{P(A, R)} = \\ = \frac{P(T, A, R, \bar{E}) + P(T, A, R, E)}{P(A, R, T) + P(A, R, \bar{T})} = \\ = \frac{P(T, A, R, \bar{E}) + P(T, A, R, E)}{P(A, R, T, \bar{E}) + P(A, R, T, E) + P(A, R, \bar{T}, \bar{E}) + P(A, R, \bar{T}, E)}$$

$$P(T, A, R, \bar{E}) = P(A|R, T, \bar{E}) \cdot P(R, T, \bar{E}) = \\ = P(A|R, T, \bar{E}) \cdot P(R|T, \bar{E}) \cdot P(T, \bar{E}) = \\ = P(A|R, T, \bar{E}) \cdot P(R|T, \bar{E}) P(T|\bar{E}) P(\bar{E})$$

$\approx 1\%$ .

## ⑥ Example: linear regression

### Univariate normal

$$\mathcal{N}(x|\mu, \sigma^2) = \frac{1}{\sqrt{2\pi\sigma^2}} e^{-\frac{(x-\mu)^2}{2\sigma^2}}$$

$$\mathbb{E} x = \mu$$

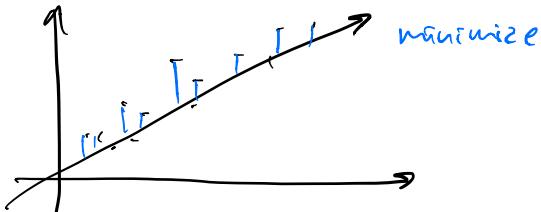
$$\text{Var}[x] = \sigma^2$$

### Multivariate normal

$$\mathcal{N}(x|\mu, \Sigma) = \frac{1}{\sqrt{|2\pi\Sigma|}} \exp\left\{-\frac{1}{2} (x-\mu)^T \Sigma^{-1} (x-\mu)\right\}$$

$$\mathbb{E} x = \mu \quad \text{Cov}[x] = \Sigma$$

### Linear regression

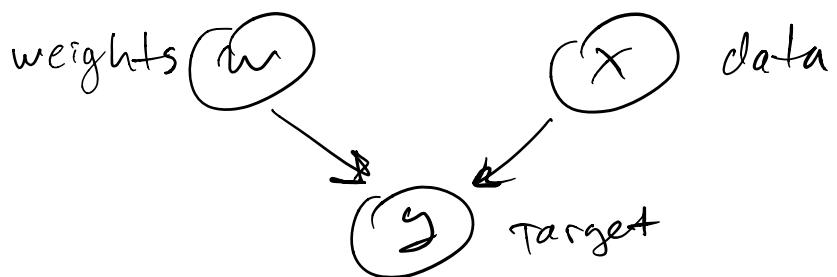


least squares problem

$$L(w) = \sum_{i=1}^N (w^T x_i - y_i)^2 = \|w^T X - y\|^2 \rightarrow \min_w$$

$$\hat{w} = \underset{w}{\operatorname{arg\,min}} L(w)$$

### Model



$$P(w, y | x) = P(y | x, w) P(w)$$

$$P(y | w, x) = \mathcal{N}(y | w^\top x, \sigma^2 I)$$

$$P(w) = \mathcal{N}(w | 0, j^2 I)$$

$$P(w | y, x) = \frac{P(y, w | x)}{P(y | x)} \rightarrow \max_w$$

$$P(y, w | x) = P(y | x, w) \cdot P(w) \rightarrow \max_w$$

$$\log P(y, w | x) = \log [P(y | x, w) \cdot P(w)]$$

$$\log \{ P(y | x, w) \cdot P(w) \} =$$

$$= \log P(y | x, w) + \log P(w) =$$

$$= \log C_1 \cdot \exp \left( -\frac{1}{2} (y - w^\top x)^\top \{ \sigma^2 I \}^{-1} (y - w^\top x) \right) +$$

$$+ \log C_2 \exp \left( -\frac{1}{2} w^\top \{ j^2 I \}^{-1} w^\top \right) =$$

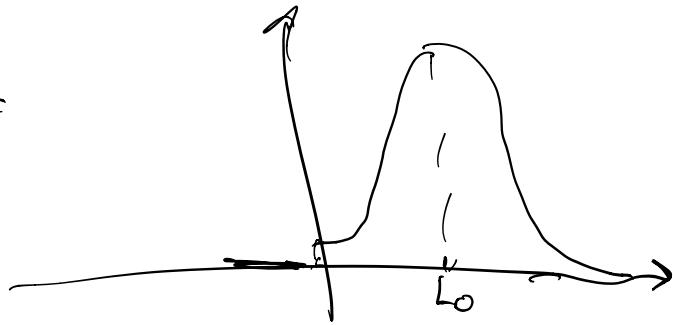
$$= -\frac{C_1}{2C^2} (y - w^\top x)^\top (y - w^\top x) - \frac{C_2}{2\gamma^2} w^\top w =$$

margin

$$= -\frac{1}{2C^2} \|y - w^\top x\|^2 - \frac{1}{2\gamma^2} \|w\|^2 \quad \left. \right\} - 1/2C^2$$

$$\min_w \quad \|y - w^\top x\|^2 + \underbrace{\frac{C^2}{\gamma^2} \|w\|^2}_{L_2}$$

$$P(O, L, L_0, T, t) = \\ = P(O | L, L_0, T, t)$$



$$P(O | L) \cdot P(L | L_0, \text{temp}) \cdot p(\text{temp} | L_0, T) \cdot P(L_0 | PC)$$

$$P(L, L_0, T)$$

(2) choose correct statements

$$p(a|b) = \sum p(a, b, c) / c = (\rightarrow)$$

$$= \int \frac{p(a, b, c)}{p(b, c)} dc = \frac{p(a|b, c) p(b|c) p(c)}{}$$

$$p(a|b) = \frac{p(a, b)}{p(b)}$$

$$\int \frac{p(a, b, c)}{p(b)} dc$$

$$p(a|b) = \sum p(a, c|b) dc$$

$$\frac{p(a, b)}{p(b)} \quad \int \frac{p(a, c, b)}{p(b)} dc \quad (\checkmark)$$

$$p(a|b) \quad \int p(a|b,c) p(c) dc$$

$$\frac{p(a,b)}{p(b)} \quad \int \frac{p(a,b,c)}{p(b,c)} \cdot p(c) dc$$

$$\underbrace{\int \frac{p(a,b,c)}{p(b)}}_{} = \int \frac{p(a|b,c)}{p(b)p(c)} p(c) dc$$

$$p(a,b|c) \quad p(a|b,c) \quad p(b|c)$$

$$\frac{p(a,b,c)}{p(c)} \quad \frac{p(a,b,c)}{p(b,c)} \cdot \frac{p(b,c)}{p(c)}$$

③

$$\frac{p(a|b,c)}{\underbrace{p(a,b,c)}_{p(b,c)p(c)}} \quad \frac{p(a|b)}{p(b)} \quad \frac{p(a|c)}{p(c)}$$

$$\int p(a,b,c) dc \quad \int p(a,b) p(a,c) dc$$

$$p(a|bc) = p(a|b) p(a|c)$$

1

$$p(a|b,c) = \frac{p(b|a,c) \cdot p(a|c)}{\int p(b|a',c) p(a'|c) da'}$$

$$p(a|b,c) = \frac{p(a,b,c)}{p(b,c)} = p(b|a,c)$$

$$p(a|b,c) = \frac{p(a,b,c)}{p(b,c)} = \frac{p(b|a,c) \cdot p(a|c)p(c)}{p(b,c)} =$$

$$= \frac{p(b|a,c) p(a|c) p(c)}{\int p(a',b,c) da'} =$$

$$= \frac{p(b|a,c) p(a|c) p(c)}{\int p(b|a',c) p(a'|c) p(c) da'}$$

$$p(a|b) = \frac{p(a,c|b)}{p(c|a,b)}$$

$$\frac{p(a,b)}{p(b)} = \frac{\cancel{p(a,c,b)}}{p(b)} \cdot \frac{p(a,b)}{\cancel{p(c,a,b)}}$$

$$p(a|b)p(b) + p(a|\bar{b})p(\bar{b})$$

$$\sum_{\forall b'} p(a|b')p(b') \neq$$

$$= \sum_{\forall b'} p(a,b) = p(a)$$

$$p(a|b) \neq p(a|b)$$

$$\sum_{\forall b'} p(a|b') = \sum_{\forall b'} \frac{p(a,b')}{p(b')}$$

## Conjugate priors

### 1. Analytical inference

Posterior distribution

$$P(\theta|X) = \frac{P(X|\theta) P(\theta)}{P(X)} \quad \begin{matrix} \leftarrow \text{prior} \\ \leftarrow \text{evidence} \\ \text{likelihood} \end{matrix}$$

what is  $P(X)$ ?

### Maximum a posteriori

$$\theta_{\text{MAP}} = \arg \max_{\theta} P(\theta|X)$$

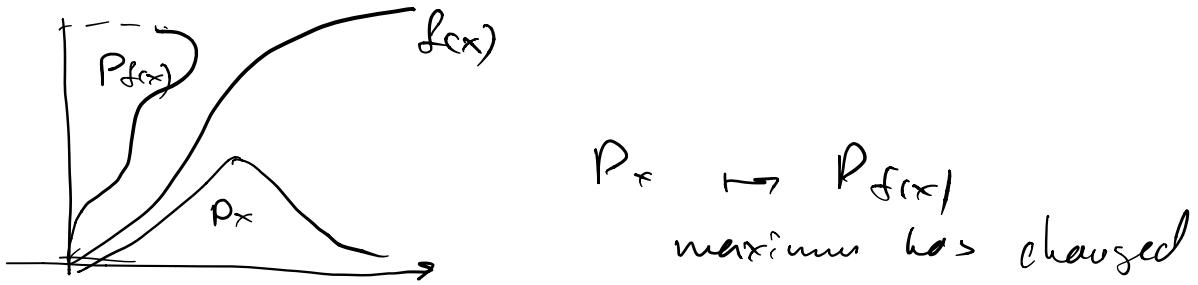
$$\theta_{\text{MAP}} = \arg \max_{\theta} \frac{P(X|\theta) P(\theta)}{P(X)}$$

$$\theta_{\text{MAP}} = \arg \max_{\theta} P(X|\theta) P(\theta)$$

optimization problem

MAP : problems

Not invariant to reparametrization

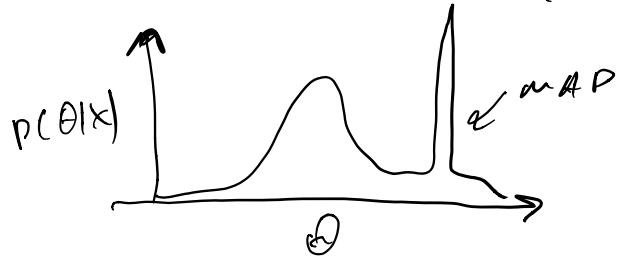


Can't use as prior

$$P_{k+1}(\theta) = \frac{P(x_k | \theta) P_{k-1}(\theta)}{P(x_k)}$$

$$P_k(\theta) = \frac{P(x_k | \theta) \delta(\theta - \theta_{mp})}{P(x_k)} = \delta(\theta - \theta_{mp})$$

MAP is a solution to  $L(\theta) = \sum [\theta \neq \theta^*] \min_{\theta}$



Summary:

Pros

- Easy to compute

Cons:

- Not invariant to reparametrization
- Can't use as a prior
- finds untypical point
- Can't compute credible intervals

## Conjugate distributions

$$P(\theta|X) = \frac{P(X|\theta) P(\theta)}{P(X)} \leftarrow \text{our own choice}$$

↓  
 fixed  
by  
model
 
 ↓  
 fixed by data

$P(\theta)$  is conjugate to  $P(X|\theta)$ :

$$P(\theta|X) = \frac{P(X|\theta) P(\theta)}{P(X)} \leftarrow \mathcal{A}(v)$$

↓  
 $\mathcal{A}(v')$

$\mathcal{A}$  - same family of distributions

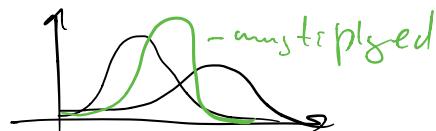
### Example

$$P(X|\theta) = N(X|\theta, \sigma^2) \quad A(v) - ?$$

$$P(\theta|X) = \frac{P(X|\theta) P(\theta)}{P(X)} \leftarrow N(X|\theta, \sigma^2)$$

Two gaussians

$$P(X_1) \sim N(\mu_1, \sigma_1^2) \quad P(X_2) \sim N(\mu_2, \sigma_2^2)$$



$$N(X|\theta, \sigma^2) \quad \leftarrow N(\theta|\mu, \sigma^2)$$

solution:

$$N(\theta|a, b^2) \quad P(\theta|X) \rightarrow \frac{P(X|\theta) P(\theta)}{P(X)}$$

$$p(\theta|x) = \frac{p(x|\theta)p(\theta)}{p(x)} = \frac{N(x|\theta, 1) N(\theta|0, 1)}{p(x)}$$

$$p(\theta|x) \sim e^{-\frac{1}{2}(x-\theta)^2} \cdot e^{-\frac{1}{2}\theta^2} \sim e^{-(\theta - \frac{x^2}{2})^2}$$

$$p(\theta|x) = N(\theta | \frac{x}{2}, \frac{1}{2})$$

Examples: Normal, precision

3. Distributions: Gamma

$$f(j|a, b) = \frac{b^a}{r(a)} j^{a-1} \cdot e^{-bj} \quad j, a, b > 0$$

$$\mathbb{E}[j] = a/b \quad \text{Mode}[j] = \frac{a-1}{b}$$

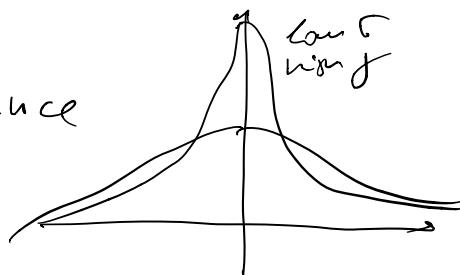
$$\text{Var}[j] = a/b^2$$

example: you ran  $\sqrt{a} \approx 100m$  a dog  
 $\frac{a}{b} = 5000 \quad \frac{a}{b^2} \approx 100$  expectation std

4. Example: Normal, precision

Precision  $j = \frac{1}{\sigma^2}$  - variance

$$N(x|\mu, \sigma^2) = \frac{1}{\sqrt{2\pi\sigma^2}} e^{-\frac{(x-\mu)^2}{2\sigma^2}}$$



$$N(x|\mu, j^2) = \frac{1}{\sqrt{2\pi}} e^{-j \frac{(x-\mu)^2}{2}}$$

$$N(x|\mu, j^2) \sim j^{\frac{1}{2}} e^{-bj} \quad \text{wrong}$$

$$P(j) \sim j^{\frac{1}{2}} e^{-bj}$$

$$P(j|x) = \frac{P(x|j) P(j)}{P(x)} \sim j e^{-j(b + \frac{(x-\mu)^2}{2})}$$

$$P(j) \sim j^{a-1} e^{-bj}$$

$$P(j) = \Gamma(j|a, b) \sim j^{a-1} b^{-1} j$$

$$P(j|x) \sim P(x|j) P(j) \sim \left( j^{\frac{1}{2}} e^{-j \frac{(x-\mu)^2}{2}} \right) \cdot \left( j^{a-1} e^{-bj} \right)$$

$$P(j|x) \sim j^{\frac{1}{2}+a-1} e^{-j(b + \frac{(x-\mu)^2}{2})}$$

$$P(j|x) = \Gamma(a + \frac{1}{2}, b + \frac{(x-\mu)^2}{2})$$

Example: Bernoulli

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

$$x \in [0, 1] \quad a, b > 0$$

$$B(a, b) = \frac{\Gamma(a+b)}{\Gamma(a)\Gamma(b)}$$

$$\mathbb{E}[x] = \frac{\alpha}{\alpha + b} \quad \text{Mode}[x] = \frac{\alpha - 1}{\alpha + b - 2}$$

$$\text{Var}[x] = \frac{\alpha b}{(\alpha + b)^2 (\alpha + b - 1)}$$

Beta prior

$$p(x|\theta) = \theta^{N_s} (1-\theta)^{N_o}$$

$$p(\theta) = B(\theta|a, b) \sim \theta^{a-1} (1-\theta)^{b-1}$$

$$p(\theta|x) \sim p(x|\theta) p(\theta)$$

$$p(\theta|x) \sim \theta^{N_s} (1-\theta)^{N_o} \cdot \theta^{a-1} (1-\theta)^{b-1}$$

$$p(\theta|x) \sim \theta^{N_s+a-1} (1-\theta)^{N_o+b-1}$$

$$p(\theta|x) = B(N_s+a, N_o+b)$$

2. Summary

$$p(\theta|x) = \frac{p(x|\theta)p(\theta)}{p(x)} \quad \begin{matrix} \text{choose prior} \\ \text{evidence} \end{matrix}$$

- pros:
- Exact posterior
  - Easy for on-line learning
- e.g.  $p(\theta|x) = B(N_s+a, N_o+b)$

cons:

- conjugate prior may be inadequate.

## Week 2

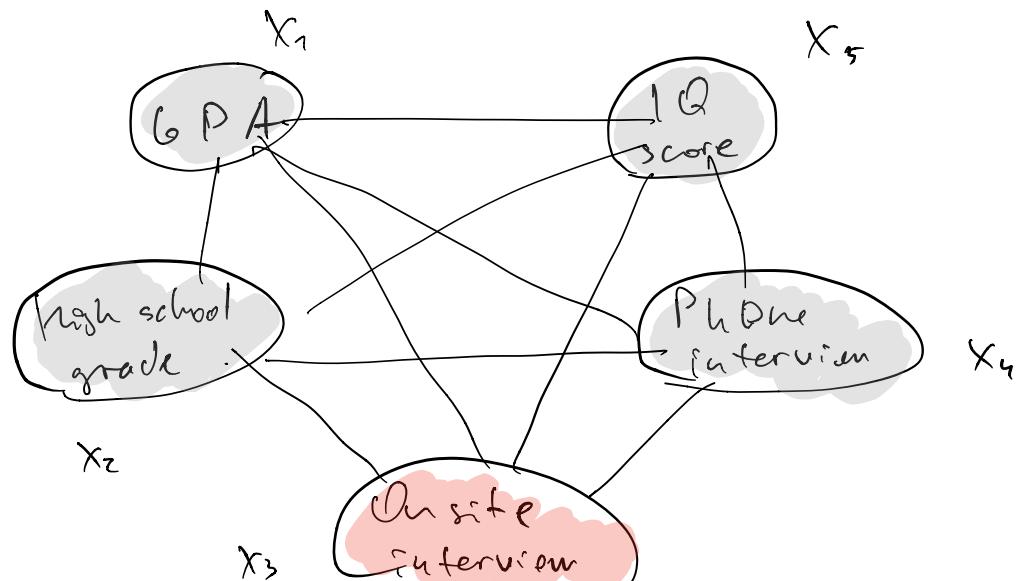
23.08.2018

### Expectation-Maximization algorithm

#### Latent Variable Models

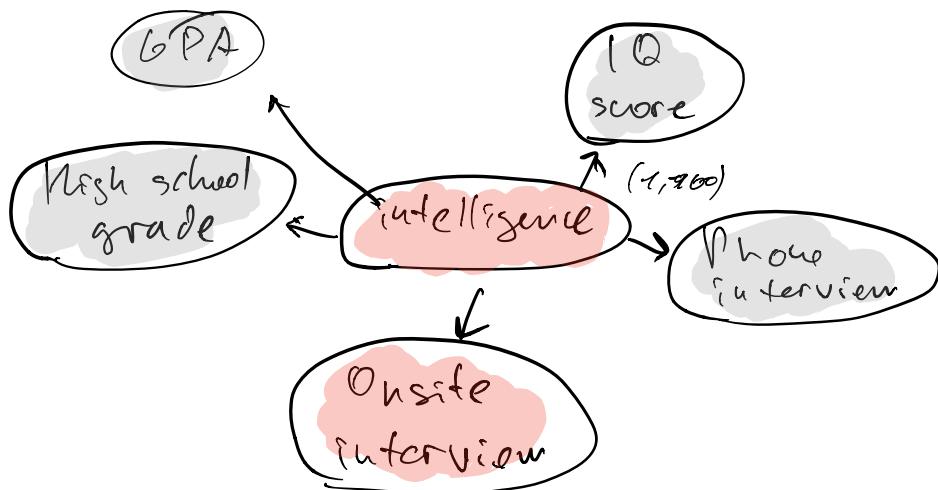
Latent (hidden) variable is a variable that you never observe

probabilistic model for hiring



$$P(X_1, X_2, X_3, X_4, X_5) = \frac{\exp(-w^T x)}{Z_{\text{in practical}}}$$

let's introduce new variable: intelligence



$$\begin{aligned}
 p(x_1, x_2, x_3, x_4, x_5) &= \sum_{I=1}^{100} p(x_1, x_2, x_3, x_4, x_5 | I) p(I) = \\
 &= \sum_{I=1}^{100} p(x_1 | I) \dots p(x_5 | I) p(I)
 \end{aligned}$$

### Latent variable

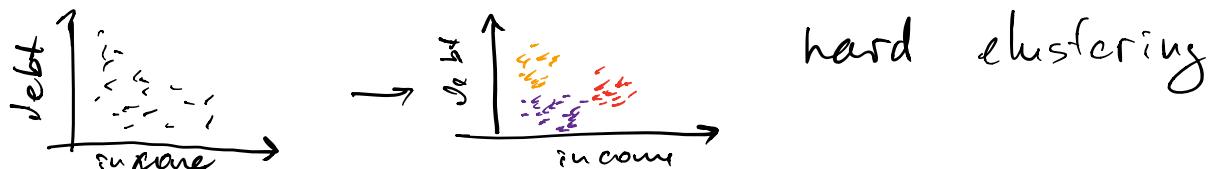
Pros:

- Simpler models (less edges)
- Fewer parameters
- Are sometimes meaningful

Cons:

- harder to work with

### Probabilistic clustering



soft clustering

hard clustering

$$p(\text{cluster id } k | \mathbf{x})$$

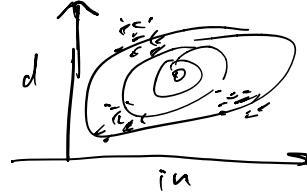
$$\text{cluster id } k = f(\mathbf{x})$$

## Gaussian Mixture Model

probabilistic model of data

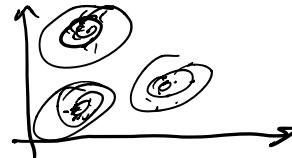
$$p(\mathbf{x} | \theta) = N(\mathbf{x} | \mu, \Sigma)$$

$$\theta = \{\mu, \Sigma\}$$



$$p(\mathbf{x} | \theta) = \pi_1 \cdot N(\mathbf{x} | \mu_1, \Sigma_1) + \pi_2 \cdot N(\mathbf{x} | \mu_2, \Sigma_2) + \pi_3 \cdot N(\mathbf{x} | \mu_3, \Sigma_3), \quad \sum \pi_i = 1$$

$$\theta = \{\pi_1, \pi_2, \pi_3, \mu_1, \mu_2, \mu_3, \Sigma_1, \Sigma_2, \Sigma_3\}$$



## Training GMM

$$\max_{\theta} p(\mathbf{x} | \theta) = \prod_{i=1}^N p(\mathbf{x}_i | \theta) = \left( \text{dataset has } n \text{ points} \right)$$

$$= \prod_{i=1}^N (\pi_1 \cdot N(\mathbf{x}_i | \mu_1, \Sigma_1) + \dots)$$

$$\text{subject to } \sum \pi_i = 1, \forall i \quad \pi_i \geq 0$$

$$\Sigma_{kk} > 0; \left( \begin{array}{l} \text{SGD can't solve problem} \\ \text{with this constraint} \end{array} \right)$$

### Summary

- GMM is a flexible probability distribution
- It is hard to fit (train) with SGD

### Introducing latent variable

$$p(x|\theta) = \pi_1 \cdot N(x|\mu_1, \Sigma_1) + \pi_2 \cdot N(x|\mu_2, \Sigma_2) + \pi_3 \cdot N(x|\mu_3, \Sigma_3)$$



$$p(t=c | \theta) = \pi_c$$

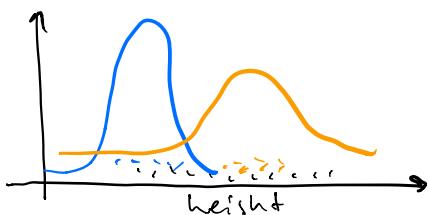
$$p(x | t=c, \theta) = N(x | \mu_c, \Sigma_c)$$

$$p(x | \theta) = \sum_{c=1}^3 p(x | t=c, \theta) \cdot p(t=c | \theta)$$

$$p(x | \theta) = \int p(x | t, \theta) p(t | \theta) dt$$

### Expectation Maximization

Dataset:  $\{x_1, \dots, x_n\}$



How to estimate parameter  $\theta$ ?

If sources + are known, e.g.:

$$p(x | t=1, \theta) = N(x | \mu_1, \sigma_1^2)$$

$$\mu_1 = \frac{\sum_{\text{blue}: x_i}}{\# \text{ of blue points}} \quad \sigma_1^2 = \frac{\sum_{\text{blue}:} (x_i - \mu_1)^2}{\# \text{ of blue points}}$$

If we know only:  $f_i: p(f_i | x_i, \theta)$

$$\mu_1 = \frac{\sum_i p(f_i=1 | x_i, \theta) x_i}{\sum_i p(f_i=1 | x_i, \theta)}$$

$$\sigma_1^2 = \frac{\sum_i p(f_i=1 | x_i, \theta) (x_i - \mu_1)^2}{\sum_i p(f_i=1 | x_i, \theta)}$$

But, we don't know the sources

Given:  $p(x | t=1, \theta) = N(-2, 1)$

Find:  $p(t=1 | x, \theta)$

$$p(t=1 | x, \theta) = \frac{p(x | t=1, \theta) p(t=1 | \theta)}{\sum}$$

## Chicken and egg problem

- Need Gaussian parameters to estimate sources
- Need sources to estimate Gaussian parameters

## EM algorithm

1. Start with 2 randomly placed Gaussians parameters  $\Theta$
2. Until convergence repeat:
  - a) for each point compute  $p(c|x_i, \Theta)$ : does  $x_i$  look like it came from cluster  $c$ ?
  - b) update Gaussian parameters  $\Theta$  to fit points assigned to them.

## Example of GMM training

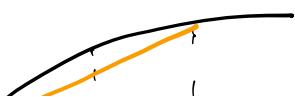
....

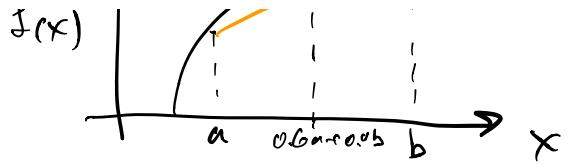
## Expectation Maximization algorithm

### Jensen's inequality & Kullback-Leibler divergence

## General form of EM

concave functions





Def:  $f(x)$  is concave if

$$\forall a, b, \lambda : f(\lambda a + (1-\lambda)b) \geq \underline{\lambda f(a) + (1-\lambda)f(b)} \quad 0 \leq \lambda \leq 1$$

### Jensen's inequality

$$f(\lambda a_1 + (1-\lambda)b) \geq \lambda f(a_1) + (1-\lambda)f(b) \quad 0 \leq \lambda \leq 1$$

$$f(\lambda_1 a_1 + \lambda_2 a_2 + \lambda_3 a_3) \geq \lambda_1 f(a_1) + \lambda_2 f(a_2) + \lambda_3 f(a_3)$$

$$\lambda_1 + \lambda_2 + \lambda_3 = 1 \quad \lambda_i \geq 0$$

Let's call  $\lambda_i$ -probabilities

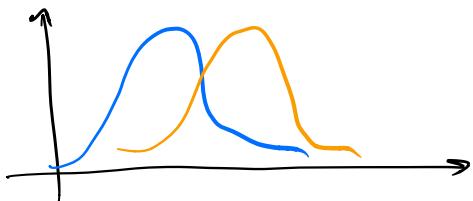
$$p(t = a_1) = \lambda_1$$

$$p(t = a_2) = \lambda_2$$

$$p(t = a_3) = \lambda_3$$

$$f(E_{p(a)} t) \geq E_{p(a)} f(t) : f\text{-concave}$$

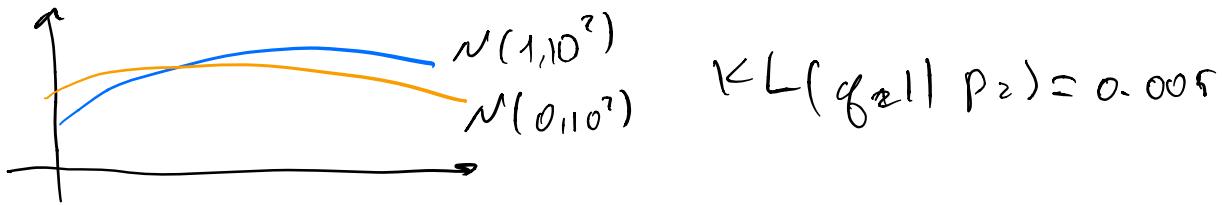
### Kullback-Leibler divergence



$$N(0, 1)$$

$$N(1, 1)$$

$$KL(q_1 \| p_1) = 0.5$$



$$KL(q \parallel p) = \int q(x) \log \frac{q(x)}{p(x)} dx$$

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

$$1. \quad KL(q_0 \parallel p) \neq KL(p \parallel q)$$

$$2. \quad KL(q \parallel q) = 0$$

$$3. \quad KL(q \parallel p) \geq 0$$

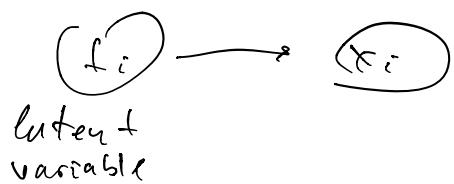
Proof:  $-KL(q \parallel p) = -\mathbb{E}_q \left[ \log \frac{q}{p} \right] =$

$$= \mathbb{E}_q \left[ -\log \frac{q}{p} \right] = \mathbb{E}_q \left[ \log \frac{p}{q} \right] \leq$$

$$\leq \log \mathbb{E}_q \left( \frac{p}{q} \right) = \log \int q(x) \frac{p(x)}{q(x)} dx = \log \int p(x) dx = 0$$

## Expectation-Maximization algorithm

General form EM



$$p(x_i | \theta) = \sum_{c=1}^s p(x_i | f_i = c, \theta) p(f_i = c | \theta)$$

$$\max_{\theta} p(X|\theta)$$

$$\max_{\theta} \left[ \log p(X|\theta) \right] = \log \prod_{i=1}^N p(x_i|\theta) =$$

↑  
assume  
independence

$$= \sum_{i=1}^N \log p(x_i|\theta)$$

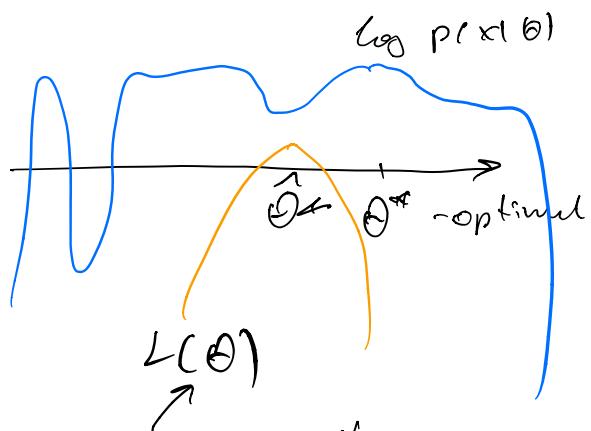
$$\log p(X|\theta) = \sum_{i=1}^N \log p(x_i|\theta) =$$

$$= \sum_{i=1}^N \log \left( \sum_{c=1}^3 p(x_i, f_i=c | \theta) \right)$$

V/ Jensen's inequality

lower bound  $L(\theta)$

example:



$$\log p(X|\theta) = \sum_{i=1}^N \log p(x_i|\theta) = \sum_{i=1}^N \log \sum_{c=1}^3 p(x_i, f_i=c | \theta) p(f_i=c | \theta)$$

$$= \sum_{i=1}^N \log \sum_{c=1}^3 \underbrace{q_{\theta}(f_i=c)}_{q_{\theta-1}(f_i=c)} \underbrace{p(x_i, f_i=c | \theta)}_{\sum}$$

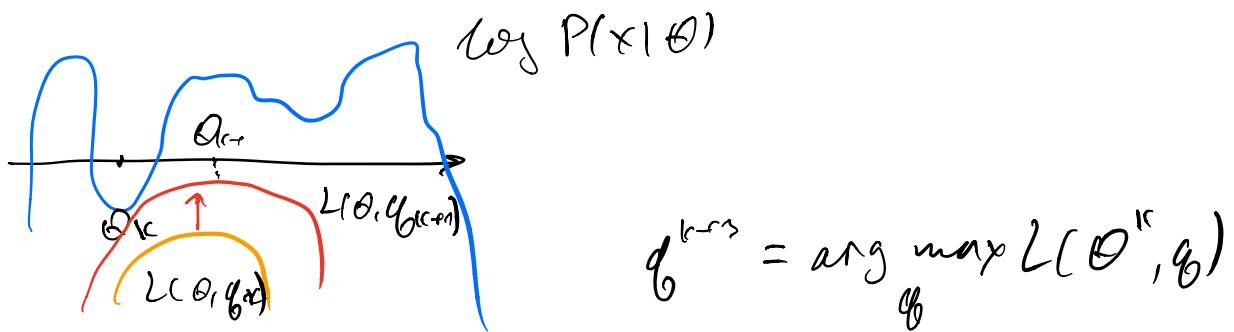
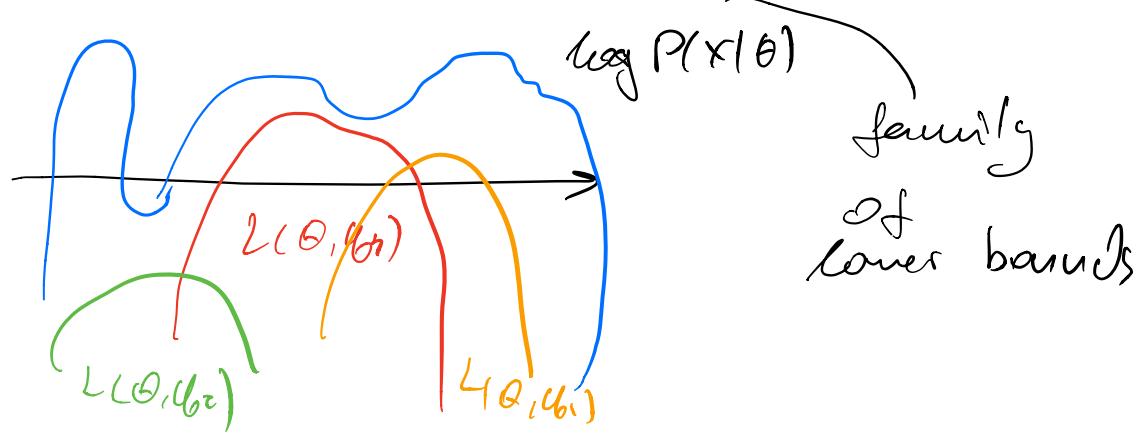
$i=1$

$L(\theta)$

$$\log \left( \sum_c a_c v_c \right) \geq \sum_c a_c \log(v_c)$$

$$\geq \sum_{i=1}^N \sum_{c=1}^S \theta^{(t_i=c)} \log \frac{p(x_i, t_i=c | \theta)}{q_\theta(t_i=c)}$$

$$\log p(x|\theta) \geq L(\theta, q) \text{ for any } q$$



### Summary of Expectation Maximization

$$\log P(x|\theta) \geq L(\theta, q) \text{ for any } q$$

(variational  
lower bound

### E-step

$$q^{(k+1)} = \arg \max_{q_b} L(\Theta^k, q_b)$$

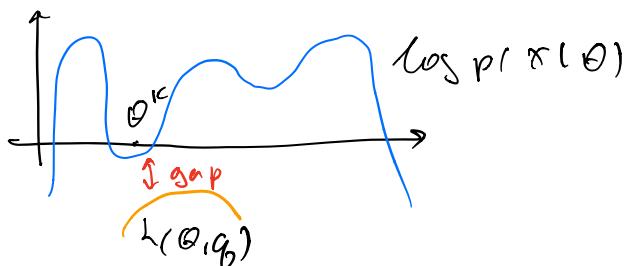
### M-step

$$\Theta^{(k+1)} = \arg \max_{\Theta} L(\Theta, q^{(k+1)})$$

### E-step details

$$\log P(X|\Theta) \geq L(\Theta, q_b)$$

E-step:  $\max_{q_b} L(\Theta^k, q_b)$



$$GAP = \log P(X|\Theta) - L(\Theta, q_b) =$$

$$= \sum_{i=1}^N \log p(x_i|\Theta) - \sum_{i=1}^N \sum_{c=1}^3 q_b(f_i=c) \log \frac{p(x_i, f_i=c|\Theta)}{q_b(f_i=c)} =$$

$$= \sum_{i=1}^N \left( \overbrace{\log p(X|\Theta)}^{} + \sum_{c=1}^3 q_b(f_i=c) - \sum_{c=1}^3 q_b(f_i=c) \log \frac{p(x_i, f_i=c|\Theta)}{q_b(f_i=c)} \right) =$$

$$\begin{aligned}
&= \sum_{i=1}^N \sum_{c=1}^3 q(f_i=c) \left( \log P(x_i | \theta) - \log \frac{P(x_i, f_i=c | \theta)}{q(f_i=c)} \right) = \\
&= \sum_{i=1}^N \sum_{c=1}^3 q(f_i=c) \left[ \log \frac{P(x_i | \theta) q(f_i=c)}{P(x_i, f_i=c | \theta)} \right] = \\
&= \sum_{i=1}^N \underbrace{\sum_{c=1}^3 q(f_i=c) \left\{ \log \frac{q(f_i=c)}{p(f_i=c | x_i, \theta)} \right\}}_{KL(q(f_i) || p(f_i=c | x_i, \theta))} \\
&\quad \text{GAP} = \sum_{i=1}^N \underbrace{KL(q(f_i) || p(f_i=c | x_i, \theta))}_{\min_q} \geq 0
\end{aligned}$$

### M-step details

$$\begin{aligned}
L(\theta, q) &= \sum_i \sum_c q(f_i=c) \log \frac{P(x_i, f_i=c | \theta)}{q(f_i=c)} = \\
&= \sum_i \sum_c q(f_i=c) \log P(x_i, f_i=c | \theta) - \\
&\quad - \sum_i \sum_c q(f_i=c) \log q(f_i=c) \quad \text{constant w.r.t. } \theta
\end{aligned}$$

$$\Leftrightarrow \mathbb{E}_q \log p(X, T | \theta) \sim \text{const}$$

(Usually) concave function w.r.t  $\Theta$ , easy to optimize

### Expectation-Maximization algorithm

For  $k=1, \dots$

#### E-step

$$q^{(k+1)} = \arg \min KL[q^{(k)} \| p(T|X, \Theta^k)]$$

$\Leftrightarrow$

$$q^{(k+1)}(f_i) = p(f_i | x_i, \Theta^k)$$

#### M-step

$$\Theta^{(k+1)} = \arg \max_{\Theta} E_{q^{(k+1)}} \log p(X, T | \Theta)$$

#### Convergence guarantees

$$\begin{aligned} \log p(X | \Theta^{(k+1)}) &\geq L(\Theta^{(k+1)}, q^{(k+1)}) \geq \\ &\geq L(\Theta^k, q^{(k+1)}) = \log p(X | \Theta^k) \end{aligned}$$

$\Downarrow$

$$\log p(X | \Theta^{(k+1)}) \geq \log p(X | \Theta^k)$$

- On each iteration EM doesn't decrease the objective

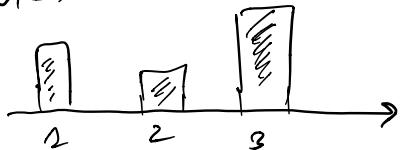
- Guaranteed to converge to a local minimum

maximum (or saddle point)

Example: EM for discrete mixture

E - Step

$$N_1 = 30 \quad N_2 = 20 \quad N_3 = 50$$



EM example

you want to fit distribution to the data

$$P(x_i) = \gamma \cdot P_1(x_i) + (1-\gamma) P_2(x_i)$$

$P(x_i | \theta)$   
every answer  
means

	1	2	3
$t=1$	$\alpha$	$1-\alpha$	0
$t=2$	0	$1-\beta$	$\beta$

Find  $\gamma, \alpha, \beta$  with EM

$$\alpha_0 = \beta_0 = \alpha_0 = 0.5$$



$$P(t_i=1) = \gamma$$

$$P(t_i=2) = 1-\gamma$$

$$P(x_i | t_i=2) = P_2(x_i)$$

E step

$$q_t(f_i=c) = P(f_i=c | x_i)$$

$$p(f_i=2 | x_i=1) = \frac{p(x_i=1 | f_i=1) p(f_i=1)}{p(x_i=1 | f_i=1) p(f_i=1) + p(x_i=1 | f_i=2) p(f_i=2)} =$$

$$= \frac{\lambda \cdot j}{\lambda \cdot j + \alpha \cdot (1-j)} = 1$$

$$p(f_i=1 | x_i=3) = \frac{p(x_i=3 | f_i=1) p(f_i=1)}{p(x_i=3 | f_i=1) p(f_i=1) + p(x_i=3 | f_i=2) p(f_i=2)}$$

$$\approx \overbrace{\dots}^0 = 0$$

$$p(f_i=1 | x_i=2) = \frac{(1-\beta)\gamma}{(1-\lambda)\gamma + (1-\beta)(1-\lambda)} = 0.5$$

M-step

E-step result

$$q_0(f_i=1) = p(f_i=1 | x_i) = \begin{cases} 1, & x_i=1 \\ 0.5, & x_i=2 \\ 0, & x_i=3 \end{cases}$$

$$q_0(f_i=2) = 1 - q_0(f_i=1)$$

$$\max_{\lambda, \beta, \gamma} \sum_{i=1}^N \mathbb{E}_{q_0(f_i)} \log p(x_i | f_i) p(f_i) =$$

$$= \sum_{i=1}^n q_i(f_i=2) \cdot \log p(x_i | f_i=2) \stackrel{j}{p(f_i=2)} +$$

$$+ \sum_{i=1}^n q_i(f_i=2) \cdot \log p(x_i | f_i=2) \stackrel{(1-j)}{p(f_i=2)} =$$

$$\begin{aligned}
 &= 30 \cdot 1 \underbrace{\cdot p(f_i=1 | x_i=1)}_1 \log d \cdot j + \\
 &\quad + 20 \cdot 0.5 \cdot \log(1-d) \cdot j + \quad \left. \right\} \text{first component} \\
 &\quad + 60 \cdot 0 \cdot \log(0) \cdot j + \\
 &\quad + 30 \cdot p(f_i=2 | x_i=1) \stackrel{=0}{=} \dots + \quad \left. \right\} \text{second component} \\
 &\quad + 20 \cdot 0.5 \cdot \log(1-\beta) \cdot (1-j) + \\
 &\quad + 60 \cdot 1 \cdot \log(\beta) \cdot (1-j) \quad \left. \right\} = \\
 &= 30 \cdot \log d \cdot j + 10 \cdot \log(1-d) \cdot j + \text{const}(d)
 \end{aligned}$$

$$\nabla = 30 \frac{1}{j} \cdot j - 10 \frac{1}{1-j} \cdot j = 0$$

$$\frac{30}{2} = \frac{10}{1-d} \Rightarrow 30 - 50d = 10d \Rightarrow d = \frac{30}{60}$$

$$\beta = \frac{6}{7} \quad j = \frac{4}{7}$$

## Summary of Expectation Maximization

. Method for training Latent Variable Models



. Handles missing data

. Sequence of simple task instead of one hard

. Guarantees to converge

. Helps with complicated parameter constraints

$$\sum_c > 0$$

. Numerous extensions

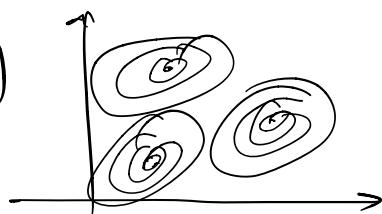
## Applications and examples

### General EM for GMM

GMM.

$$p(x|\theta) = \sum_{i=1}^k \pi_i N(x|\mu_i, \Sigma_i)$$

$$\theta = \{\pi_i, \mu_i, \Sigma_i\}$$



## E-Step

EM: for each point compute

$$q(\ell_i | \cdot) = p(\ell_i | x_i, \theta)$$

GMM: for each point compute

$$p(\ell_i | x_i, \theta)$$

## M-Step

EM: Update parameters to maximize

$$\max_{\theta} \mathbb{E}_q \log p(x, \tau | \theta)$$

GMM: Update Gaussian parameters to fit points assigned to them

$$\mu_2 = \frac{\sum_i p(\ell_i=2 | x_i, \theta) x_i}{\sum_i p(\ell_i=2 | x_i, \theta)}$$

proof:

## M-Step for GMM

$$\max_{\theta} \mathbb{E}_q \log p(x, \tau | \theta) =$$

$$= \max_{\theta} \underbrace{\sum_{i=1}^n \mathbb{E}_{q(\ell_i)} \log p(x_i, \ell_i | \theta)}_{\approx}$$

$$\sum_{i=1}^n \sum_{c=1}^3 q(\ell_i=c) \log [p(x_i | \ell_i, \theta) p(\ell_i, \theta)] =$$

$$= \sum_{i=1}^n \sum_{c=1}^3 q(f_i=c) \log \left[ \frac{1}{2} \exp\left(-\frac{(x_i - \mu_c)^2}{2\sigma_c^2}\right) \cdot \pi_c \right] =$$

$$= \sum_{i=1}^n \sum_{c=1}^3 q(f_i=c) \left[ \log \frac{\pi_c}{2} - \frac{(x_i - \mu_c)^2}{2\sigma_c^2} \right]$$

$$\frac{\partial}{\partial \mu_1} = \sum_{i=1}^n q(f_i=1) \cdot \frac{(x_i - \mu_1)}{\sigma_1^2} = 0 \quad | \quad \sigma_1^2$$

$$\sum_{i=1}^n q(f_i=1) (x_i - \mu_1) = 0$$

$$\sum q(f_i=1) x_i - \sum_{i=1}^n q(f_i=1) \mu_1 = 0$$

$$\mu_1 = \frac{\sum q(f_i=1) x_i}{\sum q(f_i=1)}$$

same for  $\sigma_1^2$

$$\sigma_1^2 = \frac{\sum_i (x_i - \mu_1)^2 \cdot q(f_i=1)}{\sum_i q(f_i=1)}$$

$$\pi_c \geq 0 \quad \text{same with constraint}$$

$$\sum \pi_c = 1$$

$$\pi_c = \frac{\sum_{i=1}^N q_i (t; c=c)}{N}$$

## K-means from probabilistic perspective

1. Randomly initialize parameters  $\Theta = \{\mu_1, \dots, \mu_K\}$
2. Until convergence repeat:

a) For each point compute closest centroid

$$c_i = \arg \min_c \|x_i - \mu_c\|^2$$

b) Update centroids

$$\mu_c = \frac{\sum_{i: c_i=c} x_i}{\# \{i: c_i=c\}}$$

From GMM to K-means

. Fix covariances to be identical  $\Sigma_c = I$

. Fix weight to be uniform

$$\pi_c = \frac{1}{\# \text{ of Gaussian}}$$

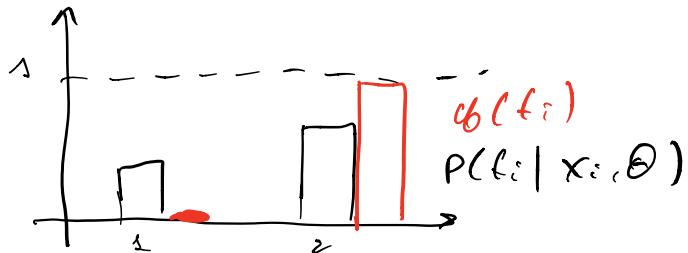
$$p(x_i | t_i=c, \Theta) = \frac{1}{Z} \exp (-0.5 \|x_i - \mu_c\|^2)$$

E-Step:

$$q^{k+1} = \arg \min_{q \in Q} KL\{q(T) \| p(T|X, \Theta^k)\}$$

where  $\mathcal{Q}$  is the set of delta-functions

example:



$$KL(q \parallel p) = \mathbb{E}_q \log \frac{q}{p} = \\ = 0 \cdot \dots + 1 \cdot \log \frac{1}{0.6} = 0.52$$

E-step:

$$q^{(k+1)}(f_i) = \begin{cases} 1 & \text{if } f_i = c_i \\ 0 & \text{otherwise} \end{cases}$$

$$c_i = \arg \max_c p(f_i = c | x_i, \theta)$$

$$p(f_i | x_i, \theta) = \frac{1}{Z} p(x_i | f_i, \theta) p(f_i | \theta) = \\ = \frac{1}{Z} \exp(-0.5 \|x_i - \mu_c\|^2) \pi_c$$

$$c_i = \arg \max_c p(f_i = c | x_i, \theta) = \arg \min_c \|x_i - \mu_c\|^2$$

K-means M-Step

$$\begin{aligned}
 & \max_{\theta} \mathbb{E}_{q^*} \log(p(x, t | \theta)) = \\
 &= \max_{\theta} \sum_{i=1}^N \underbrace{\mathbb{E}_{q^*(t_i)} \log(p(x_i, t_i | \theta))}_{q^*(t_i=c)} \\
 &= \sum_{i=1}^N \sum_{c=1}^3 q^*(t_i=c) \log p(x_i | t_i, \theta) p(t_i | \theta) =
 \end{aligned}$$

$$\mu_c = \frac{\sum_{i=1}^N q^*(t_i=c) \cdot x_i}{\sum_{i=1}^N q^*(t_i=c)} = \frac{\sum_{i: c_i=c} x_i}{\sum_{i: c_i=c} 1}$$

$$q^*(t_i) = \begin{cases} 1, & t_i = c \\ 0, & t_i \neq c \end{cases}$$

$$\theta^{(k+1)} = \arg \max_{\theta} \mathbb{E}_{q^{k+1}} \log p(x, t | \theta)$$

$$\mu_c^{k+1} = \frac{\sum_{i: c_i=c} x_i}{\sum_{i: c_i=c} 1}$$

Summary:

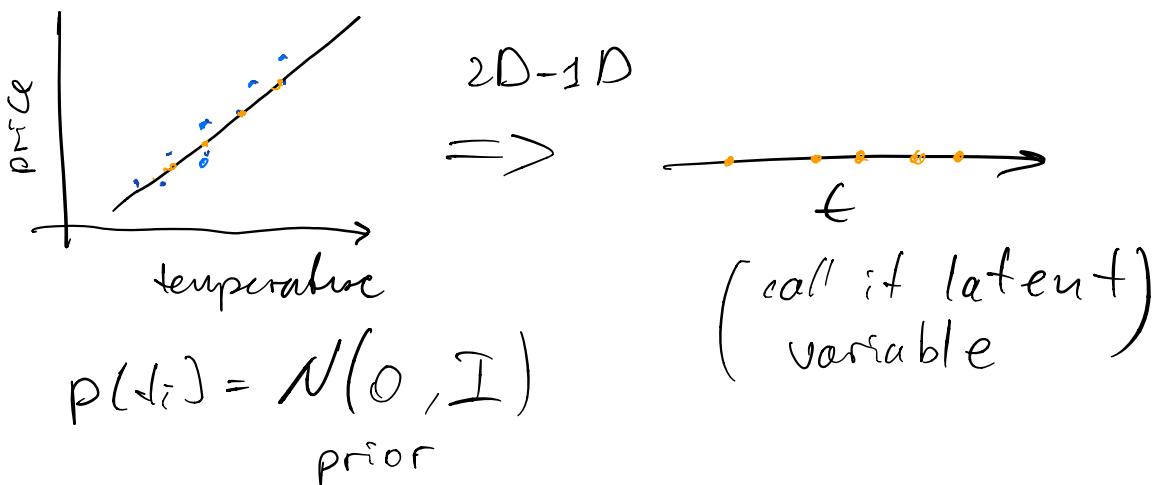
K-means is actually EM for GMM but

- With fixed covariance matrices:

$$\Sigma_c = \mathbb{I}$$

- Simplified E-step (approximate  $p(f_i | x_i, \Theta)$  with delta function)

## Probabilistic PCA



$$x_i = (s, t) \cdot f_i + (z, u)$$

$$x_i = W f_i + b$$

$$x_i = W f_i + b + \varepsilon_i, \quad \varepsilon \sim N(0, \Sigma)$$



$$p(f_i) = N(0, I)$$

$$p(x_i | f_i, \Theta) = N(W f_i + b, \Sigma)$$

$$\max_{\theta} p(x|\theta) = \prod_{i=1}^N p(x_i|\theta) =$$

$$= \prod_i \underbrace{\int p(x_i|f_i, \theta) p(f_i) df_i}_{\text{conjugacy } N(\mu_i, \Sigma_i)}$$

## EM for Probabilistic PCA

### E-step

$$q_\theta(f_i) = p(f_i|x_i, \theta) = \frac{p(x_i|f_i, \theta) p(f_i)}{Z} =$$

$$= N(\hat{\mu}_i, \hat{\Sigma}_i)$$

### M-step

$$\max_{\theta} E_{q(f_i)} \sum_i \log p(x_i|f_i, \theta) p(f_i) =$$

$$= \sum_i E_{q(f_i)} \log \left( \frac{1}{Z} \exp(-) \exp(-) \right) =$$

$$= \sum_i \log \frac{1}{Z} + \sum_i E_{q(f_i)} \log \left( \exp(-) \exp(-) \right)$$

$$\sum_i \mathbb{E}_{q(x_i)} \log \left[ \exp \left( -\frac{(x_i - w_i^\top - b)^2}{2C^2} \right) \exp \left( -\frac{f_i^2}{\sigma^2} \right) \right]$$

$a f_i^2 + C f_i + d$

## Probabilistic formulation of PCA

- Allows for missing data
- Straightforward iterative scheme for large dimensionalities
- Can do mixture of PPCA
- Hyperparameter tuning (number of components or choose between diagonal and full covariance)

## Week 3

### Variational Inference & Latent Dirichlet Allocation

#### Variational Inference

##### Why approximate inference?

$$P^*(z) = P(z|X) = \frac{P(X|z)P(z)}{P(X)}$$

- Easy for conjugate priors
- Hard otherwise

Example:  $P(x|z) = \mathcal{N}(x|\mu(z), \Sigma(z))$  [VAE]

↑  
neural networks

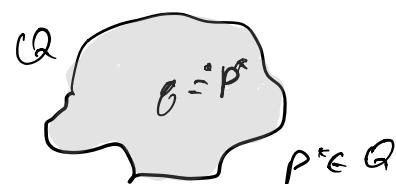
#### VI

1. Select a family of distributions  $Q$   
example:  $\mathcal{N}(\mu, (\Sigma_{ij}))$

2. Find best approximation  $q(z)$  of  $P^*(z)$

$$KL[q(z) || P^*(z)] \rightarrow \min_{q \in Q}$$

#### choice of variational family



### Unnormalized distribution

$$p^*(z) = p(z|x) = \frac{p(x|z)p(z)}{p(x)} = \frac{\hat{p}(z)}{Z}$$

optimization

$$KL\left\{q_{\theta}(z) \parallel \frac{\hat{p}(z)}{Z}\right\} = \int q_{\theta}(z) \log \frac{q_{\theta}(z)}{\hat{p}(z)/Z} dz =$$

$$= \int q_{\theta}(z) \log \frac{q_{\theta}(z)}{\hat{p}(z)} dz + \int q_{\theta}(z) \log Z dz =$$

$$= KL\{q_{\theta}(z) \parallel \hat{p}(z)\} + \log Z$$

$$KL\{q_{\theta}(z) \parallel \hat{p}(z)\} \rightarrow \max_{\theta \in Q}$$

### Mean field approximation

1. Select a family of distributions  $\mathcal{Q}$

$$\mathcal{Q} = \{q_{\theta} \mid q_{\theta}(z) = \prod_{i=1}^d q_{\theta_i}(z_i)\} \quad \begin{matrix} \text{distributions} \\ \text{can be factorized} \end{matrix}$$

2. find best approximation  $q_{\theta}(z)$  of  $p^*(z)$

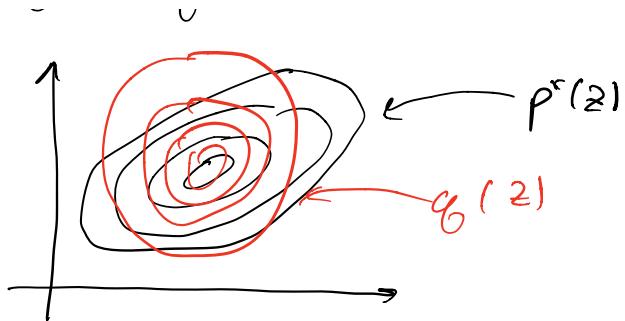
$$KL\{q_{\theta}(z) \parallel p^*(z)\} \rightarrow \min_{\theta \in Q}$$

Example:

$$p^*(z_1, z_2) \approx q_{\theta_1}(z_1) q_{\theta_2}(z_2)$$

$$p^*(z_1, z_2) = N(0, \Sigma)$$

$$q_{\theta_1}(z_1) q_{\theta_2}(z_2) = N(0, \begin{pmatrix} \Sigma_1 & 0 \\ 0 & \Sigma_2 \end{pmatrix})$$



### Optimization

$$KL(q \parallel p^*) = KL\left\{\prod_{i=1}^d q_i \parallel p^*\right\} \rightarrow \min_{q_1, q_2, \dots, q_d}$$

coordinate descend:

$$1. KL(q \parallel p^*) \rightarrow \min_{q_1}$$

$$2. KL(q \parallel p^*) \rightarrow \min_{q_2}$$

:

One step:

$$KL(q \parallel p^*) \rightarrow \min_{q_1, q_2, \dots, q_d}$$

$$KL(q \parallel p^*) = KL\left\{\prod_{i=1}^d q_i \parallel p^*\right\} =$$

$$= \int \prod_{i=1}^d q_i(z_i) \log \frac{\prod_{i=1}^d q_i(z_i)}{p^*(z)} dz =$$

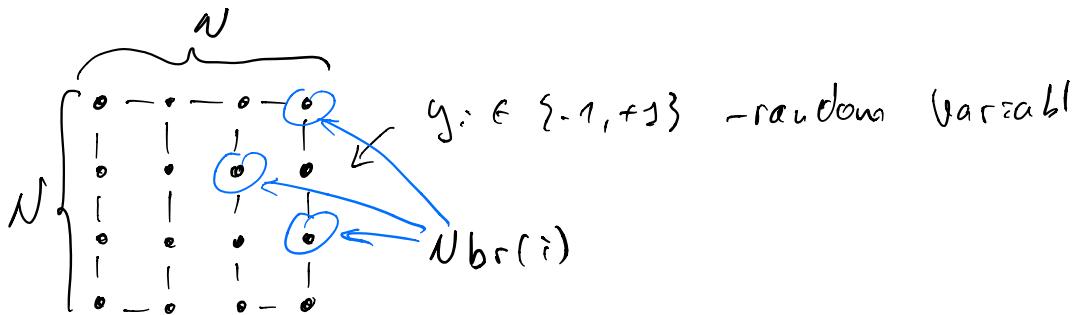
$$= \sum_{i=1}^d \int \prod_{j=1}^d q_j \cdot \log q_i dz - \int \prod_{j=1}^d q_j \cdot \log p^* dz =$$

$$\begin{aligned}
&= \underbrace{\sum_{j=1}^d q_j \log q_{j|k} dz}_{- \sum_{j=1}^d q_j \log p^* dz} + \sum_{i \neq k} \underbrace{\sum_{j=1}^d q_j \log q_{i|k} dz}_{\text{const w.r.t. } k} - \\
&\quad \left( \sum_{j=1}^d q_j \log p^* dz \right) \Rightarrow \\
&\quad \int q_{k|k} \log q_{k|k} \left[ \sum_{j \neq k} q_j dz_{j|k} \right] dz_{k|k} = \\
&= \int q_{k|k} \log q_{k|k} dz_{k|k} \\
&\quad \sum_{i \neq k} \int q_i \log q_i dz_i \\
&\quad \text{const} \\
\Leftrightarrow & \int q_{k|k} \log q_{k|k} dz_{k|k} - \int q_{k|k} \left[ \sum_{j \neq k} q_j \log p^* dz_{j|k} \right] dz_{k|k} = \\
&= \int q_{k|k} \left\{ \log q_{k|k} - \underbrace{\left[ \sum_{j \neq k} q_j \log p^* dz_{j|k} \right]}_{h(z_k)} \right\} dz_{k|k} - \text{const} = \\
&\quad h(z_k) = \mathbb{E}_{q_{k|k}} \log p^* \\
&f(z_k) = \frac{e^{h(z_k)}}{\int e^{h(z_k)} dz_k} - \text{our new distributions}
\end{aligned}$$

$$\begin{aligned}
&= \int q_{k|k} \log \frac{q_{k|k}}{f} dz_{k|k} + \text{const} \rightarrow \min \\
&K \perp [q_{k|k} || f] \rightarrow \min
\end{aligned}$$

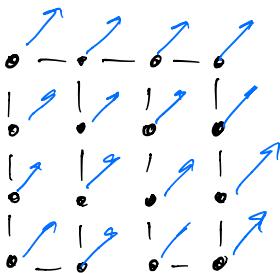
$$q_{ik} = t \Rightarrow \log q_{ik} = \mathbb{E}_{\theta_k} \log p^* + \text{const}$$

Example: Ising model

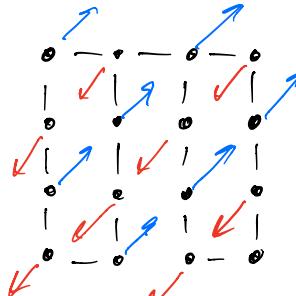


$$p(s) \propto \exp \left( \frac{1}{Z} \cdot \sum_i \sum_{j \in Nbr(i)} y_i y_j + \sum_i b_i y_i \right)$$

$\phi(s)$



Ferromagnetic  
 $J > 0$



Anti-ferromagnetic  
 $J < 0$

normalization constant

$$P(s) = \frac{1}{Z} \phi(s) \quad Z = \sum_s \phi(s) \quad 2^{N^2} \text{ terms}$$

$$P(s) \approx q(s) = \prod_i q_i(y_i)$$

$$P(y) \propto \exp\left(\frac{1}{2} \cdot J \cdot \sum_i \sum_{j \in \text{nbr}(y_i)} y_i y_j + \sum_i b_i y_i\right)$$

$Nbr(y_u)$

$f_k(y_k) - ?$

$$\begin{aligned}
 \log q_k &= \mathbb{E}_{q_k} \log P + \text{const} = \\
 &= \mathbb{E}_{q_k} \left[ \frac{1}{2} J \cdot \sum_i \sum_{j \in Nbr(y_i)} y_i y_j + \sum_i b_i y_i \right] - \text{const} = \\
 &= \mathbb{E}_{q_k} \left\{ J \cdot \sum_{j \in Nbr(y_k)} y_k y_j + b_k y_k \right\} + \text{const} = \\
 &= J \cdot \sum_{j \in Nbr(y_k)} y_k \underbrace{\mathbb{E}_{q_k} y_j}_{\mu_j} + b_k y_k - \text{const} = \\
 &= y_k \cdot \underbrace{\left( J \cdot \sum_{j \in Nbr(y_k)} \mu_j + b_k \right)}_M + \text{const}
 \end{aligned}$$

$$q_k = C \cdot \exp(y_k \cdot M)$$

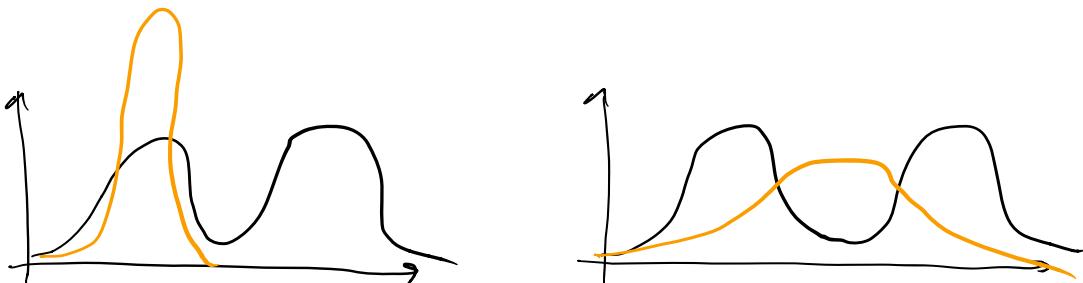
$$q_k(1) + q_k(-1) = 1$$

$$C \cdot e^M + C \cdot e^{-M} = 1 \quad C = \frac{1}{e^M + e^{-M}}$$

$$q_{\text{enc}}(+1) = \frac{e^m}{e^m + e^{-m}} = \frac{1}{1 + e^{-2m}} = \sigma(2M)$$

$$\mu_m = 1 \cdot q(+1) - (-1) q(-1) = q(+1) - q(-1) = \\ = \frac{e^m - e^{-m}}{e^m + e^{-m}} = \tanh(M)$$

### Optimization solutions



x capture statistics v

v mode has high probability x

$$KL[q||p^*] = \int q(z) \log \frac{q(z)}{p^*(z)} dz = +\infty$$

$\circlearrowleft$

### Variational EM & Review

$$\log p(x|\theta) \geq L(\theta, g) = \mathbb{E}_{g(t)} \log \frac{p(x, t | \theta)}{g(t)} dt$$

↑  
marginal  
likelihood

↑  
variational  
lower  
bound

↓  
max

E-step

$$L(\theta, q) \rightarrow \max_q \Leftrightarrow KL[q_\theta(T) \parallel p(T|X, \theta)] \rightarrow \min_q$$

M-step

$$L(\theta, q) \rightarrow \max_{\theta} \Leftrightarrow \mathbb{E}_{q_\theta(T)} \log p(X, T | \theta) \rightarrow \max_{\theta}$$

E-step

$$KL[q_\theta(T) \parallel p(T|X, \theta)] \rightarrow \min_q$$

full posterior

$$q_\theta(T) = p(T|X, \theta)$$

variational inference

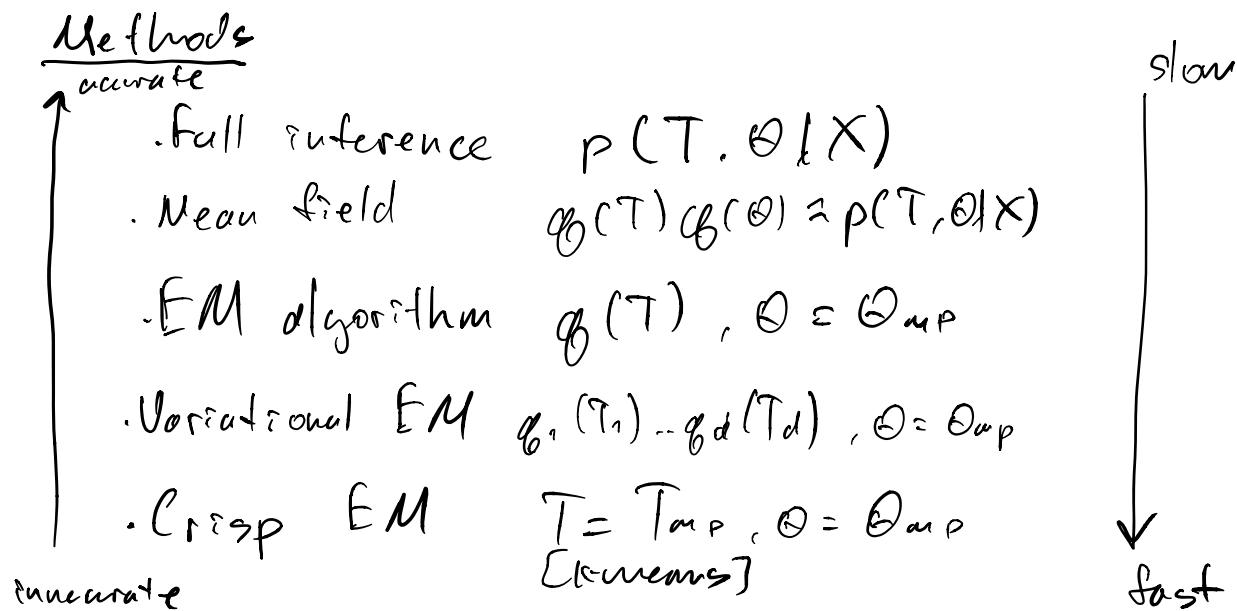
$$KL[q_\theta(T) \parallel p(T|X, \theta)] \rightarrow \min_{q \in Q}$$

Model

Known:  $X$  data

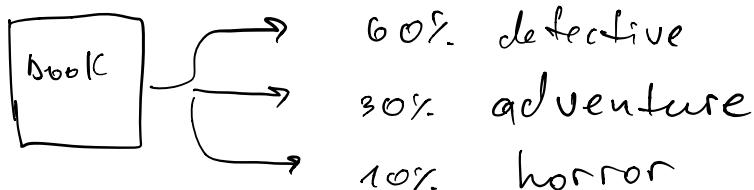
Unknown:  $\theta$  parameters

Unknown:  $T$  latent variables



## Latent Dirichlet Allocation

### Topic Modeling



Document is a distribution over topics

sport	economy	politics	(topics)
20% football	24% money	10% president	
10% hockey	9% dollar	4% usa	

Football player from USA has

salary in dollars

Topic is a distribution over words

### Similarity

Book1 → 60% detective  
30% adventure =  $\begin{pmatrix} 0.6 \\ 0.3 \\ 0.1 \end{pmatrix} = a$   
10% horror

Book2 → 62% detective  
33% adventure =  $\begin{pmatrix} 0.62 \\ 0.33 \\ 0.05 \end{pmatrix} = b$   
5% horror

### Goals

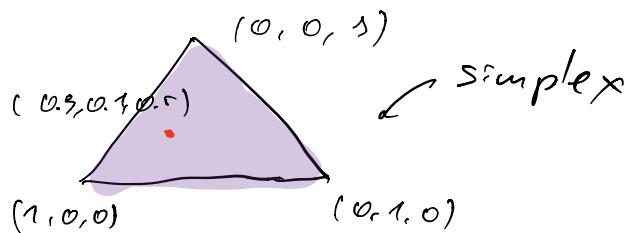
- ① Construct topics
- ② Assign topics to texts

### Dirichlet distribution

$$\text{Dir}(\theta | \alpha) = \frac{1}{B(\alpha)} \prod_{k=1}^K \theta_k^{\alpha_k - 1}$$

$$\sum \theta_k = 1$$

$$\theta_k > 0$$



## Statistics

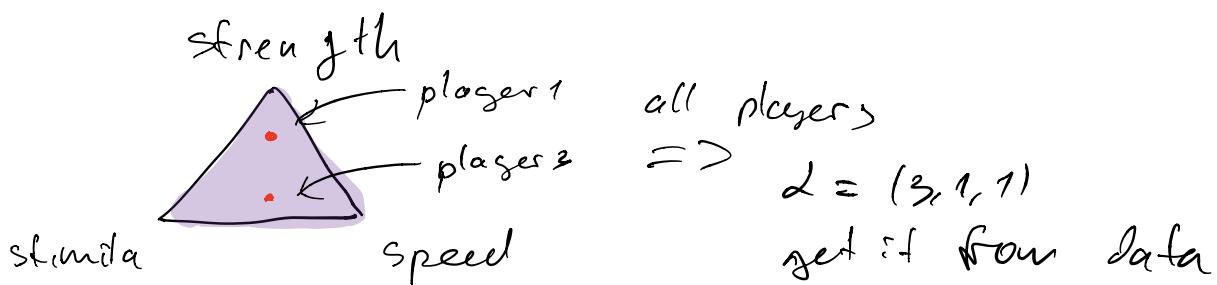
$$\mathbb{E} \theta_i = \frac{d_i}{d_0}$$

$$\text{cov}(\theta_i, \theta_j) = \frac{d_i d_0 [i=j] - d_i d_j}{d_0^2 (d_0 + 1)}$$

$$d_0 = \sum_{i=1}^k d_i$$

## Example

massively multiplayer online  
role-playing game (MMORPG)



## Conjugate prior

$P(\theta)$  is conjugate to  $P(X|\theta)$

$$P(\theta|X) = \frac{P(X|\theta) P(\theta)}{P(X)} \propto A(v)$$

$\int A(v)$

## Multinomial likelihood

$$P(X|\theta) = \frac{u!}{x_1! \dots x_K!} \theta_1^{x_1} \dots \theta_K^{x_K}$$

$$P(\theta) = \text{Dir}(\theta|\alpha) = \frac{1}{B(\alpha)} \prod_{k=1}^K \theta_k^{\alpha_k - 1}$$

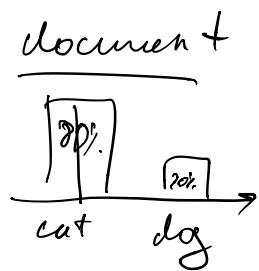
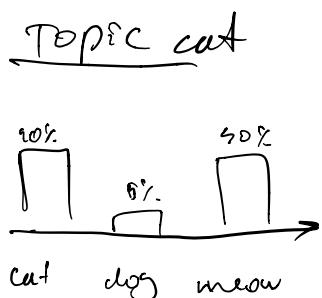
$$p(\theta|X) \propto \prod_{k=1}^K \theta_k^{\alpha_k + x_k - 1}$$

$$p(\theta|X) = \text{Dir}(\theta | (\underbrace{\dots}_{\alpha_k}, \underbrace{\dots}_{x_k}))$$

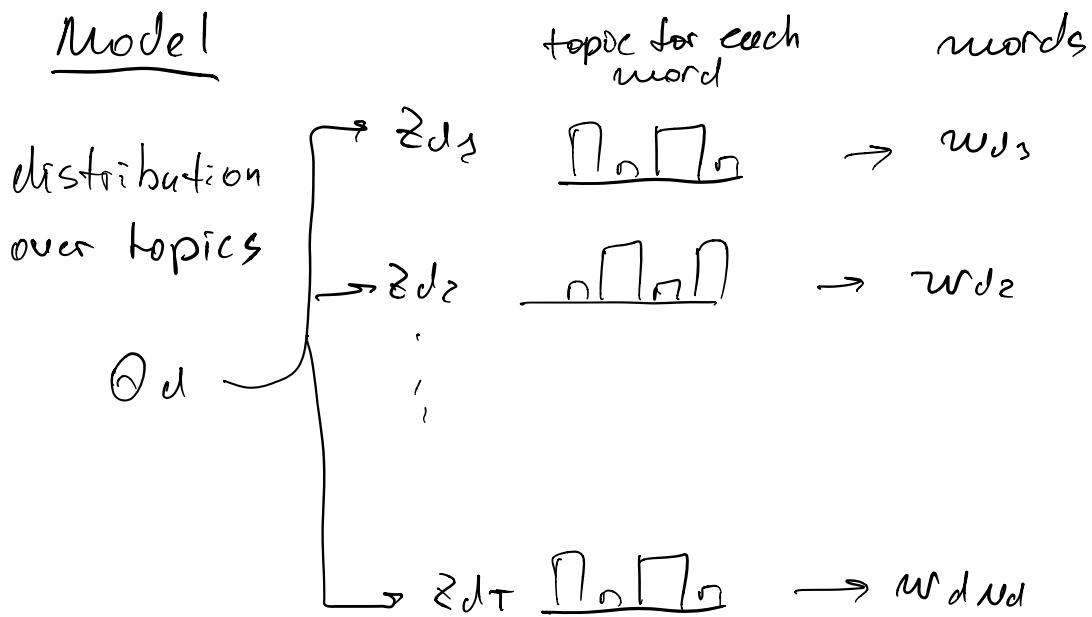
## Latent Dirichlet Allocation

Document is a distribution over topics

Topic is a distribution over words

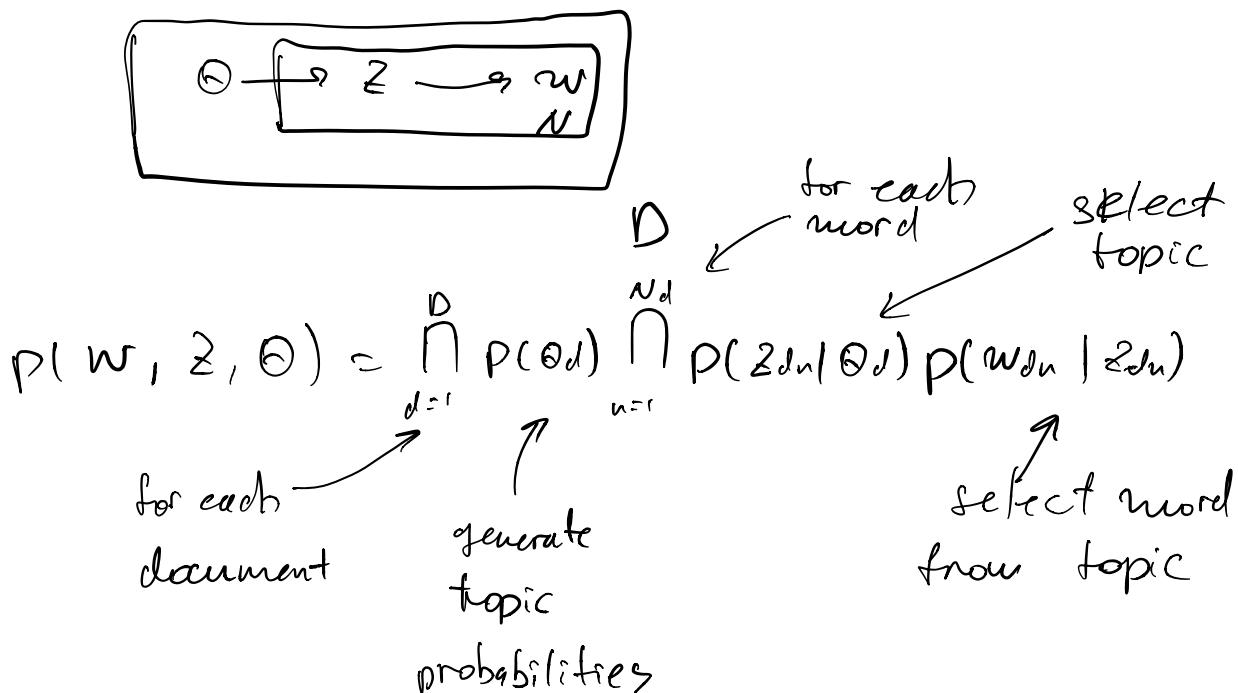


## Model



$$z_{dn} \in \{1, \dots, T\} \quad w_{dn} \in \{1, \dots, V\}$$

## LDA Model



$$p(W, Z, \Theta) = \prod_{d=1}^D p(\Theta_d) \prod_{n=1}^{N_d} p(z_{dn} | \Theta_d) p(w_{dn} | z_{dn})$$

$$p(\Theta_d) \sim \text{Dir}(\alpha)$$

$$p(z_{dn} | \Theta_d) = \Theta_{dz_{dn}}$$

$$p(w_{dn} | z_{dn}) = \phi_{z_{dn} w_n}$$

constraints:

$\phi_{wm} \geq 0$

$\sum_w \phi_{wm} = 1$

Known: W data

Unknown:  $\phi$  parameters, distribution over words for each topic

Unknown:  $Z$  latent variables, topic of each word

Unknown:  $\Theta$  latent variables, distribution over topics for each document

LDA: E-step, theta

$$p(\Theta, Z, w) = \prod_{d=1}^D p(\Theta_d) \prod_{n=1}^{N_d} p(z_{dn} | \Theta_d) p(w_{dn} | z_{dn}) =$$

$$= \prod_{d=1}^D p(\Theta_d) \prod_{n=1}^{N_d} \Theta_{dz_{dn}} \cdot \phi_{z_{dn} w_n}$$

$$P(W|t) \xrightarrow{\max} \phi$$

$$\text{E-step: } KL(q(\theta)q(z) \| P(\theta, z|w))$$

$\downarrow$   
min

$$q(\theta), q(z)$$

M-step

$$\mathbb{E}_{q(\theta)q(z)} \log P(\theta, z|w) \xrightarrow{\max} \phi$$

$$\log q(\theta) = \mathbb{E}_{q(z)} \log P(\theta, z|w) + \text{const} \quad (\textcircled{=})$$

$$P(\theta, z|w) = \frac{P(\theta, z, w)}{p(w)} \text{ does not depend on } q(z)$$

$$\textcircled{=} \mathbb{E}_{q(z)} \log P(\theta, z|w) + \text{const} \quad (\textcircled{=})$$

$$\log p(\theta, z, w) =$$

$$= \sum_{d=1}^D \left[ \sum_{t=1}^T (\delta_{dt}) \log \theta_{dt} + \underbrace{\sum_{n=1}^{N_d} \sum_{t=1}^T [z_{dn=t}] (\log \phi_{dt} + \log \phi_{dn})}_{\text{constant w.r.t. } \theta} \right] - \text{const}$$

$$\textcircled{2} \quad \mathbb{E}_{q(\mathbf{z})} \left( \sum_{d=1}^D \left[ \sum_{t=1}^T (\alpha_t - 1) \log \theta_{dt} + \sum_{n=1}^{N_d} \sum_{t=1}^T [z_{dn}=t] \log \theta_{dt} \right] \right) \text{const} =$$

$$= \sum_{d=1}^D \left[ \sum_{t=1}^T (\alpha_t - 1) \log \theta_{dt} + \sum_{n=1}^{N_d} \sum_{t=1}^T \mathbb{E}_{\substack{q(\mathbf{z}) \\ z_{dn}}} [z_{dn}=t] \cdot \log \theta_{dt} \right] + \text{const}$$

*$f_{dn}$*

$$= \sum_{d=1}^D \sum_{t=1}^T \left[ (\alpha_t - 1) + \sum_{n=1}^{N_d} f_{dn} \right] \cdot \log \theta_{dt} + \text{const}$$

$$q(\boldsymbol{\theta}) = \prod_{d=1}^D \prod_{t=1}^T \theta_{dt}^{[\alpha_t + \sum_n f_{dn} - 1]} \cdot \text{const}$$

$$q(\boldsymbol{\theta}) = \prod_{d=1}^D q_d(\theta_d), \quad q_d(\theta_d) = \text{Dir}(\theta_d | d + \sum f_{dn})$$

LDA : E-step, 2

$$\log q(\mathbf{z}) = \mathbb{E}_{q(\mathbf{z})} \log p(\boldsymbol{\theta}, \mathbf{z}, \mathbf{w}) + \text{const} =$$

$$p(\boldsymbol{\theta}, \mathbf{z} | \mathbf{w}) = \frac{p(\boldsymbol{\theta}, \mathbf{z}, \mathbf{w})}{p(\mathbf{w})} \quad \begin{matrix} \text{does not depend} \\ \text{on } q(\mathbf{z}) \end{matrix}$$

$$= \mathbb{E}_{q(\mathbf{z})} \left[ \sum_d \left[ \sum_t (\alpha_t - 1) \log \theta_{dt} + \sum_n \sum_t [z_{dn}=t] (\log \theta_{dt} + \log \phi_{dn}) \right] \right]$$

*does not depend on  $\mathbf{z}$*

+ const =

$$= \mathbb{E}_{q(\theta)} \sum_d \sum_n \sum_+ [z_{dn} = +] (\log \theta_{dt} + \log q_{+m_n}) + \text{const} =$$

$$= \sum_d \sum_n \sum_+ [z_{dn} = +] \left[ \mathbb{E}_{q(\theta)} \log \theta_{dt} + \log q_{+m_n} \right] + \text{const}$$

$$q_\theta(z) = \prod_d^D \prod_n^{N_d} q_\theta(z_{dn}) + \text{const}$$

$$q_\theta(z_{dn} = +) = \frac{q_{+m_n} \cdot \exp \left( \mathbb{E}_{q(\theta_{dt})} \log \theta_{dt} \right)}{\sum_{t=1}^T q_{+m_n} \cdot \exp \left( \mathbb{E}_{q(\theta_{dt})} \log \theta_{dt} \right)} = f_{dn}(t)$$

$$f_{dn} = \mathbb{E}_{q_\theta(z_{dn})} [z_{dn} = +]$$

### LDA : M-step & Prediction

$$\mathbb{E}_{q(\theta) q(T)} \log P(\theta, z, w) \rightarrow \max_{\phi}$$

$$\begin{aligned} \mathbb{E}_{q(\theta) q(T)} \log P(\theta, z, w) &= \\ &= \mathbb{E}_{q(\theta) q(z)} \left[ \sum_d \left[ \sum_+ (x_{d+1}) \cancel{\log \theta_{dt}} + \sum_n \sum_+ [z_{dn} = +] \cancel{(\log \theta_{dt} + \log q_{+m_n})} \right] \right] = \end{aligned}$$

$$= \mathbb{E}_{\theta(0)\theta(2)} \left[ \sum_d \sum_n \sum_{+}^{+} \{z_{dn}=+ \} \cdot \log \varphi + w_{dn} \right]_{\text{exact}}$$

max



$$\checkmark \varphi + w > 0 \quad \forall t, w$$

$$\sum_{w=1}^V \varphi + w = 1 \quad \forall t \quad (\text{using Lagrangian})$$

$$L = \mathbb{E}_{\theta(0)\theta(2)} \sum_d \sum_n \sum_{+}^{+} \{z_{dn}=+ \} \log \varphi + w_{dn} +$$

$$+ \sum_{t=1}^T \lambda_+ \left( \sum_w \varphi + w - 1 \right) =$$

$$= \sum_d \sum_n \sum_{+}^{+} \varphi_{dn}^+ \log (\varphi + w_{dn}) + \sum_t \lambda_+ \left( \sum_w \varphi + w - 1 \right)$$

$$\frac{\partial L}{\partial \varphi + w} = \sum_d \sum_n \varphi_{dn}^+ \frac{\frac{\partial}{\partial \varphi + w_{dn}} \cdot \{w_{dn}=w\}}{\varphi + w_{dn}} + \lambda_+ = 0$$

$$\Rightarrow \varphi + w = \frac{\sum_d \sum_n \varphi_{dn}^+ \cdot \{w_{dn}=w\}}{-\lambda_+}$$

$$\sum_w q_{t+w} = \sum_w \frac{\sum_d \sum_n f_{dn}^+ [w_{dn}=w]}{\lambda +} = 1$$

$$\lambda + = \sum_w \sum_d \sum_n f_{dn}^+ [w_{dn}=w]$$

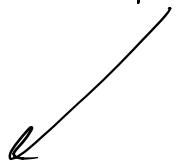
$$q_{t+w} = \frac{\sum_d \sum_n f_{dn}^+ [w_{dn}=w]}{\sum_w \sum_d \sum_n f_{dn}^+ [w_{dn}=w]}$$

predictions

$p(\Theta_{d^*}, z_{d^*} | w, \phi)$  - approximate

$d^*$ -new document

$w$ - train data



$$KL(q(\Theta_{d^*})q(z_{d^*}) || p(\Theta_{d^*}, z_{d^*} | w, \phi)) \rightarrow \min$$

$q(\Theta_{d^*})$   
 $q(z_{d^*})$

Extensions of LDA

$$p(w, z, \theta) = \prod_{d=1}^D p(\theta_d) \prod_{n=1}^{w_d} p(z_{dn} | \theta_d) p(w_{dn} | z_{dn})$$

$$p(\theta_d) \sim \text{Dir}(\alpha)$$

$\alpha \uparrow \Rightarrow$  More topics for each document

$\alpha \downarrow \Rightarrow$  Less topics for each document

$\alpha$  can be selected as  $p(w | \alpha) \rightarrow \max_{\alpha}$

### Sparcity of topics

Sparse prior on  $\phi$

$$p(w, z, \theta, \phi) = \prod_{t=1}^T p(\phi_t) p(w, z, \theta | \phi)$$

$$p(\phi_t) \sim \text{Dir}(\beta)$$

### Topics correlation

logistic normal distribution

$$p(\theta_d) \sim \mathcal{P}(\mathcal{N}(\mu, \Sigma))$$

### Dynamic Topic Model

$$p(B_{t+1}^{\uparrow} | B_t^{\uparrow}) \sim \mathcal{N}(B_t^{\uparrow}, \sigma^2 I)$$

$$\phi_{t+1}^{\uparrow} = \text{softmax}[B_{t+1}^{\uparrow}]$$

## Summary

- Many topics are interpretable
- Works well with rare words
- Fast even for huge text collections
- Multicore & distributed implementations
- Many features can be added with extensions



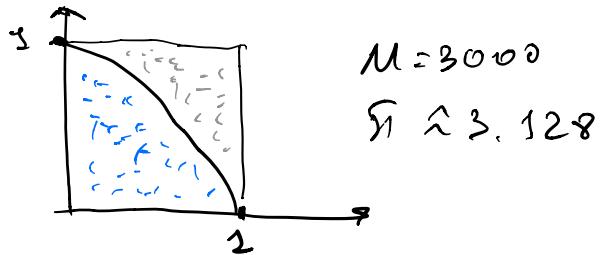
## Week 4

# Markov chain Monte Carlo

## Monte Carlo estimation

- MCMC - silver bullet of probabilistic modeling

Estimate expected value by sampling



$$M = 3000$$

$$\pi \approx 3.128$$

$$\hat{E}_q = E[x^2 + y^2 \leq 1] \approx \frac{1}{M} \sum_{s=1}^M [x_s^2 + y_s^2 \leq 1]$$

$x_s, y_s \sim U(0, 1)$

$E_{p(x)} f(x)$  - general form

$$E_{p(x)} f(x) \approx \frac{1}{M} \sum_{s=1}^M f(x_s), \quad x_s \sim p(x)$$

Why do we need to estimate expected values?

- Fall Bayesian inference (week 1)

$$p(y | x, Y_{\text{train}}, X_{\text{train}}) =$$

$$= \int p(g|x, w) p(w|X_{\text{train}}, Y_{\text{train}}) dw \quad (\textcircled{E})$$

$p(g|x, w)$  - neural network,  $w$ -weight of nn

$$\textcircled{E} \mathbb{E}_{p(w|X_{\text{train}}, Y_{\text{train}})} p(g|x, w)$$

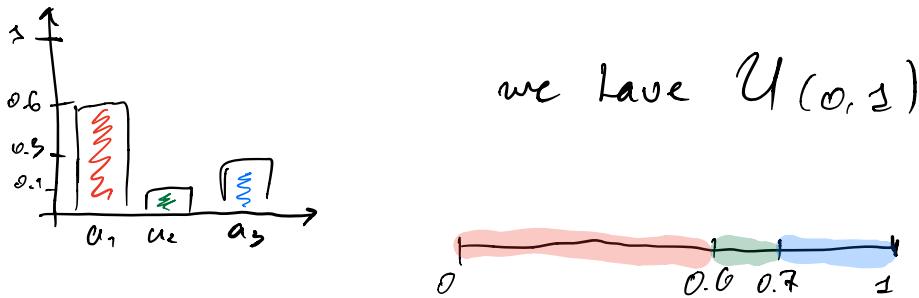
$$p(w|Y_{\text{train}}, X_{\text{train}}) = \frac{p(Y_{\text{train}}, X_{\text{train}}, w) p(w)}{\mathcal{Z}}$$

### • M-step of EM-algorithm (week 2)

$$\max_Q \mathbb{E}_q \log p(X, T|Q)$$

Sampling from s-d distributions

s-d sampling (discrete)



1-d discrete distributions with finite number of values are easy  
At least then number of values is  $< 10000$

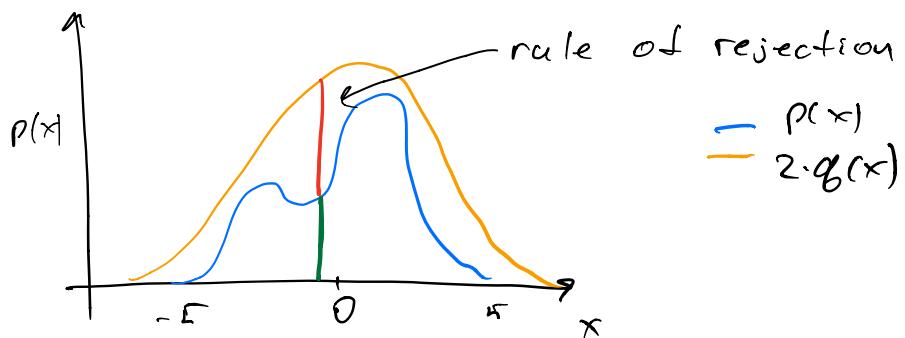
## continuous sampling

Sampling from Gaussian distribution

1 idea: central limit theorem

$$z = \sum_{i=1}^{12} x_i - 6, \quad x_i \sim U(0, 1)$$

$$p(z) \approx N(0, 1)$$



$$\hat{x} \sim q(x) \quad g \sim U[0, 2q(\hat{x})]$$

Accept  $\hat{x}$  with probability  $\frac{p(x)}{2q(x)}$ : if  $g \leq p(x)$

How is it effective?

$$p(x) \leq M q(x)$$

Accept  $\frac{1}{M}$  points on average

$$\frac{\hat{p}(x)}{2} \leq M q(x) \Rightarrow \hat{p}(x) \leq \underbrace{2M}_{\text{new } M} q(x)$$

works if you can upper bound your unnormalized distribution

## Summary of rejection sampling

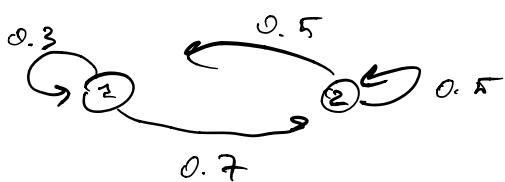
pros:

- Works for most distributions  
(even unnormalized)

cons:

- If  $q$  and  $p$  are too different ( $M$  is large), rejects most of the points
- $M$  is large for d-dimensional distributions

## Markov Chains



$$T(L \rightarrow L) = 0.3$$

$$T(L \rightarrow R) = 0.7$$

$$T(R \rightarrow L) = 0.5$$

$$T(R \rightarrow R) = 0.5$$

T-transition probabilities

	1	2
$x^1$	1	0
$x^2$	0.3	0.7
$x^3$	$0.3^2 + 0.7 \cdot 0.5$ 0.49	$0.3 \cdot 0.7 + 0.7 \cdot 0.5$ 0.56
$\vdots$		
$x^n$		

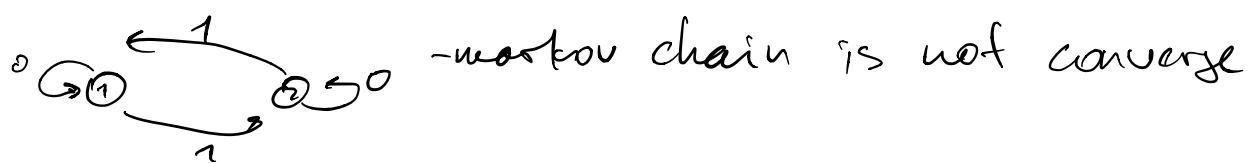
$$P(x^3) = P(x^3 | x^2=L) P(x^2=L) +$$

$$+ P(x^3 | x^2=R) P(x^2=R)$$

rule of sum prob

$$P(L) \approx 0.42 \quad P(R) \approx 0.58$$

- We want to sample from  $p(x)$
  - Build a Markov chain that converge to  $p(x)$
  - Start from any  $x^0$
  - For  $k = 0, 1, \dots$
- $$x^{k+1} \sim T(x^k \rightarrow x^{k+1})$$
- Eventually  $x^k$  will look like sample from  $p(x)$



### Definition

A distribution  $\pi$  is called stationary if

$$\pi(x') = \sum_x T(x \rightarrow x') \pi(x)$$

### Theorem:

If  $T(x \rightarrow x') > 0$  for all  $x, x'$  then exist unique  $\pi$ :

$$\pi(x) = \sum_x T(x \rightarrow x') \pi(x)$$

and Markov chain converges to  $\pi$  from any starting point

## Gibbs sampling

$$p(x_1, x_2, x_3) = \frac{\hat{P}(x_1, x_2, x_3)}{Z}$$

start with  $(x_1^0, x_2^0, x_3^0)$  e.g.  $(0, 0, 0)$

$$x_1^1 \sim p(x_1 | x_2 = x_2^0, x_3 = x_3^0) = \frac{\hat{P}(x_1, x_2^0, x_3^0)}{Z_1}$$

$$x_2^1 \sim p(x_2 | x_1 = x_1^1, x_3 = x_3^0)$$

$$x_3^1 \sim p(x_3 | x_1 = x_1^1, x_2 = x_2^1)$$

For  $k = 0, 1, 2, \dots$

$$x_1^{k+1} \sim p(x_1 | x_2 = x_2^k, x_3 = x_3^k)$$

$$x_2^{k+1} \sim p(x_2 | x_1 = x_1^{k+1}, x_3 = x_3^k)$$

$$x_3^{k+1} \sim p(x_3 | x_1 = x_1^{k+1}, x_2 = x_2^{k+1})$$

## Proof of Gibbs Sampling

$$p(x', y', z') = \sum_{x, y, z} q_0(x, y, z \rightarrow x', y', z') p(x, y, z)$$

(want to prove)

$$\begin{aligned} & \sum_{x, y, z} p(x' | y=y', z=z') p(y' | x=x', z=z') \cdot \\ & \quad \cdot p(z' | x=x', y=y') \cdot p(x, y, z) = \end{aligned}$$

$$= p(z'|x=x', g=g') \sum_{x,y,z} p(x|g=g, z=z) p(g|x=x', z=z) p(x, g, z) =$$

$$= p(z'|x=x', g=g') \sum_{y,z} p(x|g=g, z=z) p(g|x=x', z=z) \underbrace{\sum_x p(x, g, z)}_{p(y, z)} =$$

$$= p(z'|x=x', g=g') \sum_{y,z} p(x|g=g, z=z) p(g|x=x', z=z) p(y, z) =$$

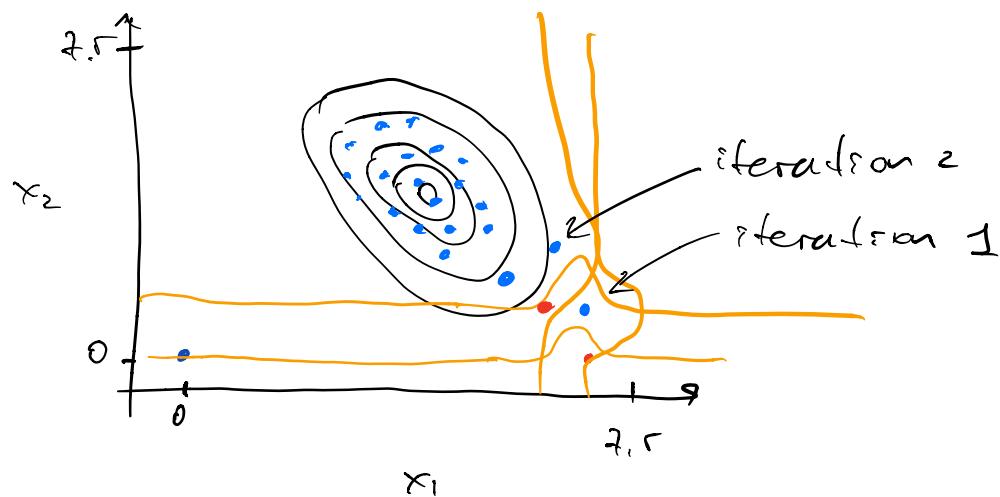
$$= p(z'|x=x', g=g') \sum_z p(g|x=x', z=z) \cdot \sum_y p(x|g=g, z=z) p(y, z) =$$

$$= p(z'|x=x', g=g') \sum_z p(g|x=x', z=z) p(x', z) =$$

$$= p(z'|x=x', g=g') \sum_z p(y', x', z) =$$

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

Example of Gibbs sampling



## Summary

pros:

- reduce multidimensional sampling to sequence of 1d sampling
- a few lines of code

cons:

- highly correlated samples
- slow convergence (mixing)
- not parallel

## Metropolis - Hastings

Apply rejection sampling to Markov Chains

For  $k = 1, 2, \dots$

- sample  $x'$  from a wrong  $Q(x^k \rightarrow x')$
- Accept proposal  $x'$  with probability  $A(x^k \rightarrow x')$
- Otherwise stay at  $x^k$ :  $x^{k+1} = x^k$

$$T(x \rightarrow x') = Q(x \rightarrow x') A(x \rightarrow x') \text{ for all } x \neq x'$$

$$\begin{aligned} T(x \rightarrow x) &= Q(x \rightarrow x) A(x \rightarrow x) + \\ &+ \sum_{x' \neq x} Q(x \rightarrow x') (1 - A(x \rightarrow x')) \end{aligned}$$

How to choose A:

$$\pi(x') = \sum_x \pi(x) T(x \rightarrow x')$$

### Detailed Balance

$$\pi(x') = \sum_x \pi(x) T(x \rightarrow x') \quad \text{stationary distribution}$$

If:

$$\pi(x) \overleftarrow{T}(x \rightarrow x') = \pi(x') \overleftarrow{T}(x' \rightarrow x)$$

detailed  
balance  
question

Then:

$$\pi(x') = \sum_x \pi(x) \overleftarrow{T}(x \rightarrow x')$$

Proof:

$$\begin{aligned} \sum_x \pi(x) \overleftarrow{T}(x \rightarrow x') &= \sum_x \pi(x') \overleftarrow{T}(x' \rightarrow x) \\ \pi(x') \underbrace{\sum_x \overleftarrow{T}(x' \rightarrow x)}_1 &= \pi(x') \end{aligned}$$

How to choose A:

$$\pi(x) \overrightarrow{T}(x \rightarrow x') = \pi(x') \overrightarrow{T}(x' \rightarrow x)$$

## Metropolis-Hastings: choosing the critic

choosing a critic:

$$\pi(x) \underbrace{Q(x \rightarrow x')}_{\tilde{\pi}(x \rightarrow x')} \underbrace{A(x \rightarrow x')}_{\tilde{A}(x \rightarrow x')} = \pi(x') \underbrace{Q(x' \rightarrow x)}_{\tilde{\pi}(x' \rightarrow x)} \underbrace{A(x' \rightarrow x)}_{\tilde{A}(x' \rightarrow x)}$$

$$\frac{A(x \rightarrow x')}{A(x' \rightarrow x)} = \frac{\pi(x') Q(x' \rightarrow x)}{\pi(x) Q(x \rightarrow x')} = p < 1$$

assumption

$$A(x \rightarrow x') = p$$

$$A(x' \rightarrow x) = 1$$

$$\frac{A(x \rightarrow x')}{A(x' \rightarrow x)} = \frac{\pi(x') Q(x' \rightarrow x)}{\pi(x) Q(x \rightarrow x')} = p > 1$$

$$A(x \rightarrow x') = 1$$

$$A(x' \rightarrow x) = 1/p$$

↓

$$A(x \rightarrow x') = \min \left\{ 1, \frac{\hat{\pi}(x') Q(x' \rightarrow x)}{\hat{\pi}(x) Q(x \rightarrow x')} \right\}$$

we don't care  
about normalization const

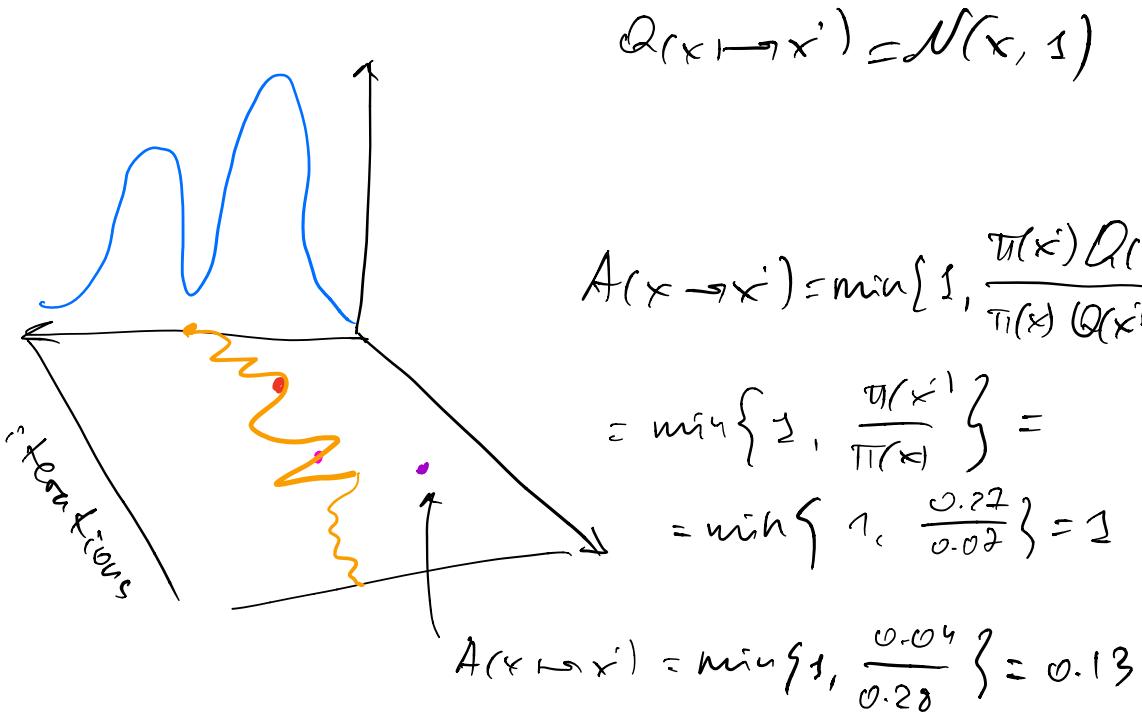
Choice of  $Q$

$$Q(x \rightarrow x') > 0$$

opposing forces:

- $Q$  should spread out, to improve mixing and reduce correlation
- But then acceptance probability is often low

## Example of Metropolis-Hastings



$Q = N(x, 0.02)$  - very slow convergence  
stay at the same region

$Q = N(x, 100)$  - very fast convergence  
stay at same place

## Metropolis-Hastings as correction scheme

Recall Gibbs sampling  
lets make it parallel

$$x_1^{(k+1)} \sim p(x_1 | x_2 = x_2^k, x_3 = x_3^k)$$

$$x_2^{(k+1)} \sim p(x_2 | x_1 = x_1^k, x_3 = x_3^k)$$

$$x_3^{(k+1)} \sim p(x_3 | x_1 = x_1^k, x_2 = x_2^k)$$

It's wrong now, but can correct  
with Metropolis Hastings!

### Summary

Rejection sampling applied to Markov Chains

Pros:

- You can choose among family of Markov chains
- Works for unnormalized densities
- Easy to implement

Cons:

- Samples are still correlated
- Have to choose among family of Markov chains

### Markov Chain Monte Carlo summary

Two MCMC approaches:

- Gibbs sampling - reducing multidimensional sampling to a sequence of 1d
- Metropolis Hastings - rejected sampling for markov chains

$\underset{\text{Monte Carlo}}{MC}$  → approximate via sampling  
 sampling via MC (markov chain)

### Monte Carlo vs Variational Inference

#### Monte Carlo

$$\mathbb{E}_{p(x)} f(x) \approx \frac{1}{M} \sum_{s=1}^M f(x_s), \quad x_s \sim p(x)$$

unbiased estimation (larger  $M \Rightarrow$  better accuracy)

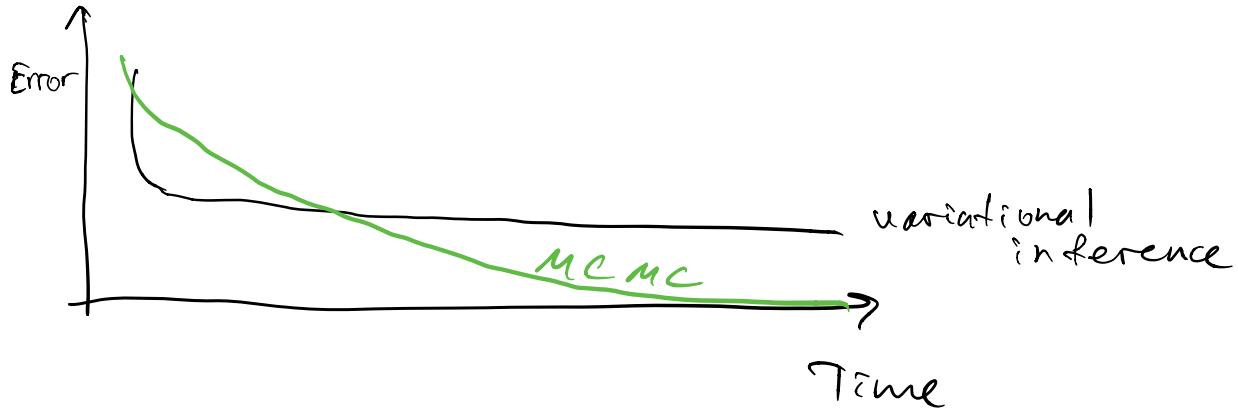
$$\mathbb{E}_{p(x)} \frac{1}{M} \sum_{s=1}^M f(x_s) = \mathbb{E}_{p(x)} f(x)$$

#### Variational Inference

$$p(x) \approx q(x)$$

$$\mathbb{E}_{p(x)} f(x) \approx \mathbb{E}_{q(x)} f(x)$$

#### Schematic illustration



### Methods

- best
- full inference  $p(T, \Theta | X)$
  - mean field  $q_T(T) q_\Theta(\Theta) \approx p(T, \Theta | X)$
  - **MCMC**  $T_s, \Theta_s \sim p(T, \Theta | X)$
  - EM algorithm  $q_T(T), \Theta = \Theta_{\text{up}}$
  - Variational EM  $q_1(T_1) \dots q_d(T_d), \Theta = \Theta_{\text{up}}$
  - **MCMC EM**  $T_s \sim p(T | \Theta, X), \Theta = \Theta_{\text{up}}$
- worst

### Summary MCMC

Pros:

- easy to implement
- easy to parallelize
- unbiased estimates (more  $\Rightarrow$  accuracy)

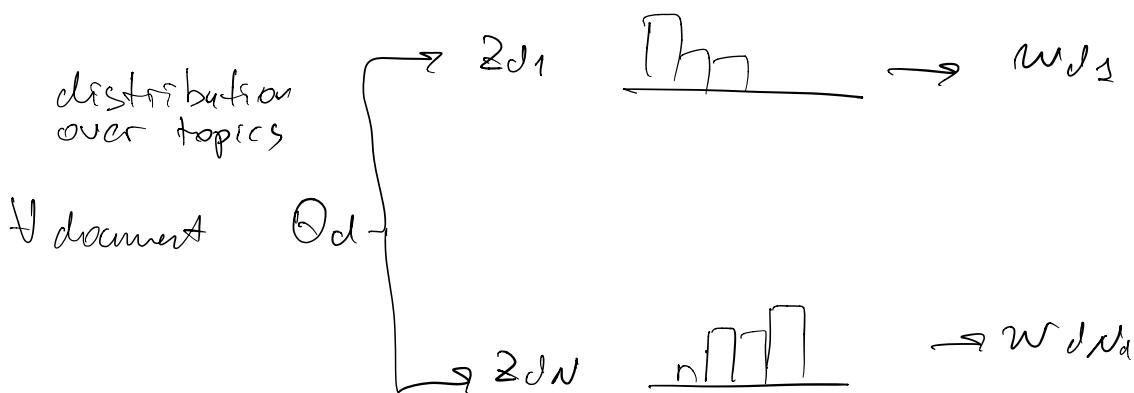
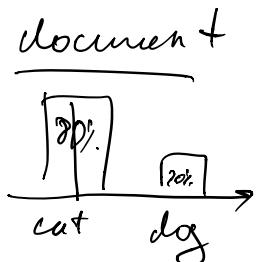
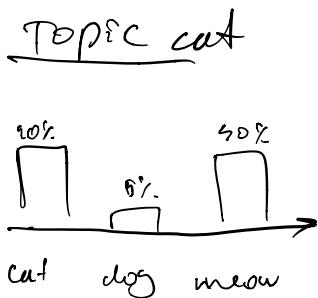
Cons:

- usually slower than alternatives

## UCMC to LDA

Document is a distribution over topics

Topic is a distribution over words



$$z_{dn} \in \{1, \dots, T\} \quad w_{dn} \in \{1, \dots, V\}$$

Mean field for LDA

$$\max_{\phi} p(w|\phi)$$

E-step  $p(z, \Theta | w, \phi) \approx q(z, \Theta)$

M-step  $\max_{\phi} \mathbb{E}_q \log p(z, \Theta, w | \phi)$

## MCMC for LDA

Known:  $W$  data

Unknown:  $\phi$  parameters, distribution over words for each topic

Unknown: 2 latent variables, topic of each word

Unknown:  $\theta$  latent variables, distribution over topics for each document  
↳ let's go full Bayesian

Known:  $W$  data

Unknown:  $\phi$  ~~latent variables~~, distribution over words for each topic

Unknown: 2 latent variables, topic of each word

Unknown:  $\theta$  latent variables, distribution over topics for each

$$p(\phi, \theta, z | w) \sim \{\text{Gibbs Sampling}\}$$

Init:  $\phi^0, \theta^0, z^0$

$$\left\{ \phi_3^t \sim p(\phi_3 | \phi_2^0, \phi_3^1, \dots, \theta^0, z^0, w) \right.$$

$$\phi_2^1 \sim p(\phi_2 | \phi_1^1, \phi_3^0, \dots, \Theta^0, Z^0, w)$$

$$\phi_i^1 \sim p(\phi_i | \phi_1^1, \dots, \phi_{i-1}^1, \phi_{i+1}^0, \dots, \Theta^0, Z^0, w)$$

$$\Theta_i^1 \sim p(\Theta_i | \Phi^1, \Theta_1^1, \dots, \Theta_{i-1}^1, \Theta_{i+1}^0, \dots, Z^0, w)$$

$$Z_i^1 \sim p(Z_i | \Phi^1, \Theta^1, Z_1^1, \dots, Z_{i-1}^1, Z_{i+1}^0, \dots, w)$$

for  $k=3, 2, \dots$

### Collapsed Gibbs for LDA

Model

$$\begin{aligned}
 p(\Theta_d) &= \text{Dir}(\beta) & p(\Phi) &= \text{Dir}(\alpha) \\
 p(Z_d | \Theta_d) &= \Theta_d z_{dn} & p(W_{dn} | Z_{dn}, \Phi) &= \Phi_{z_{dn}} w_{dn}
 \end{aligned}$$

conjugate  
can compute analytically

$$p(\Theta | Z)$$

$$p(\Theta) = \int p(Z | \Theta) p(\Theta) d\Theta =$$

$$= \frac{p(Z | \Theta) p(\Theta)}{p(\Theta | Z)}$$

conjugate

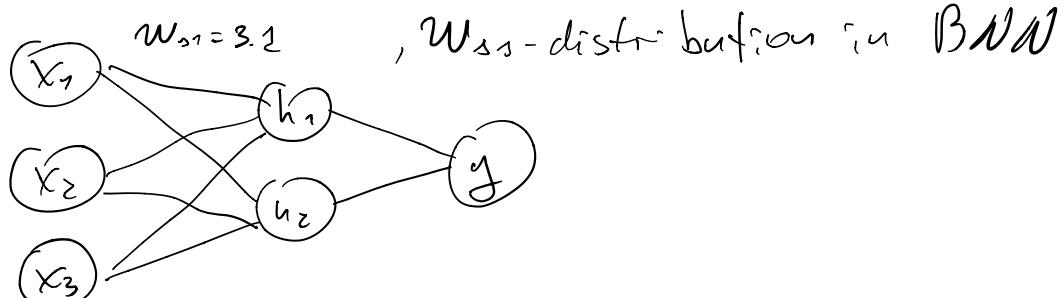
$$p(\phi | z, w) \\ p(w | z) = \frac{p(w | z, \phi) p(\phi)}{p(\phi | z, w)}$$

$$p(z | w) = \frac{p(w | z) p(z)}{c}$$

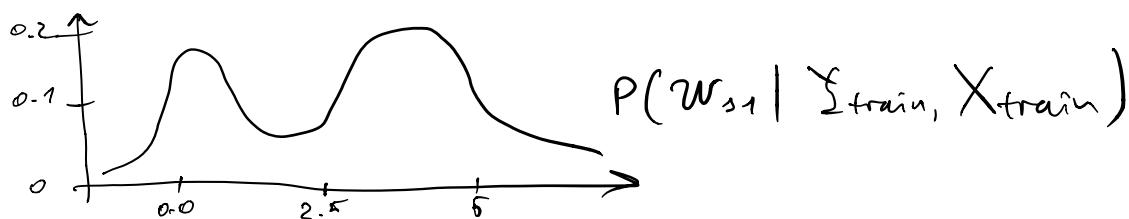
$P(z | w)$  ~ {Gibbs Sampling}

$$p(\phi | w) = \int p(\phi | w, z) p(z | w) dz = \\ = \mathbb{E}_{p(z | w)} P(\phi | w, z) \approx \\ \approx p(\phi | w, \hat{z})$$

## Bayesian Neural Networks



BdNN uses distributions instead numbers



$$\begin{aligned}
 p(y|x, \mathcal{Y}_{\text{train}}, X_{\text{train}}) &= \\
 &= \int \underbrace{p(y|x, w)}_{N \text{ output}} p(w|\mathcal{Y}_{\text{train}}, X_{\text{train}}) dw = \\
 &= \mathbb{E}_{p(w|\mathcal{Y}_{\text{train}}, X_{\text{train}})} p(y|x, w) \\
 p(w|\mathcal{Y}_{\text{train}}, X_{\text{train}}) &\sim \{\text{Gibbs}\} \\
 p(w|\mathcal{Y}_{\text{train}}, X_{\text{train}}) &= \frac{p(\mathcal{Y}_{\text{train}}|X_{\text{train}}, w) p(w)}{Z}
 \end{aligned}$$

problem: depends on the whole dataset!

### Langevin Monte Carlo

Gibbs and Metropolis Hastings can't do mini-batches

Say we want to sample from  $p(w|D)$

Start from  $w^0$

For  $k=1, \dots$

$$w^{k+1} = \underbrace{w^k + \varepsilon \nabla \log p(w^k|D)}_{\zeta^k \sim N(0, 2\varepsilon I)} =$$

gradient ascent

$$= w^k + \epsilon \triangleright \left( \log p(w^k) + \sum_{i=1}^N \log p(y_i | x_i, w^k) \right) + \underbrace{\frac{1}{2} \|w^k\|^2}_{\text{weight decay}}$$

usual cross entropy

scheme:

initialize weight  $w^0$

- Do say 100 iterations with usual SGD, but add Gaussian noise  $\zeta^k \sim N(0, 2\epsilon I)$  to each update
- After 100 epochs decide if Markov Chain converged and start collecting weight values
- For a new object predict compare average prediction of CNN with weight  $w^{100}, w^{101}, \dots, w^{200}$   
 $\uparrow$  better
- Train another CNN to mimic the ensemble



Week 5

20.10.2018

## Variational Autoencoder

### Scaling Variational Inference & Unbiased estimates

#### Scale Variational Inference

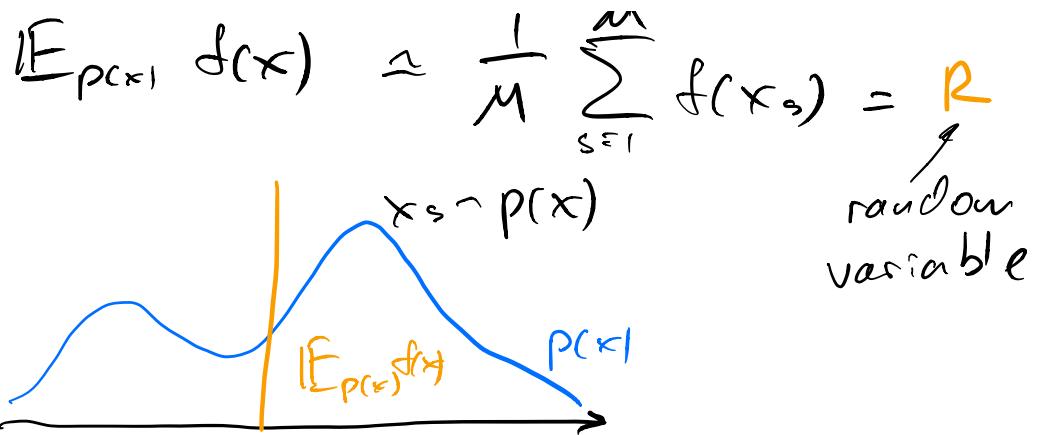
People used to think that Bayes is for small datasets

- Too slow for big data
- Not very beneficial any way

Things changed when Bayes met Deep learning

- Understand why and how combine Deep learning and Bayesian methods
- Learn how to synthesize images with VAE
- Learn state-of-the-art Bayesian neural networks and their applications

#### Unbiased estimates



$$\mathbb{E}_{p(x)} R = \mathbb{E}_{p(x)} f(x)$$

not unbiased

$$\log(\mathbb{E}_{p(x)} f(x)) \stackrel{?}{\approx} \log\left(\frac{1}{M} \sum_{s=1}^M f(x_s)\right) = G$$

$x_s \sim p(x)$

$$\mathbb{E}_{p(x)} G \neq \log(\mathbb{E}_{p(x)} f(x))$$

- Estimator called unbiased if its expected value equals to thing it estimates.

- This is unbiased estimator:

$$\mathbb{E}_{p(x)} f(x) \approx \frac{1}{M} \sum_{s=1}^M f(x_s) = R$$

Modeling a distribution of images

Let's start with fitting  $p(x)$  into a dataset.

But why do we need it?

- Generate new data.
- Detect anomalies and outliers.
- Work with missing data.
- Represent your data in a nice way (e.g. model p(molecule) to search for drugs)

### How to model $p(x)$

- $\log \hat{p}(x) = CNN(x)$

$$p(x) = \exp(CNN(x))$$

infeasible  $\rightarrow \mathcal{Z}$  - billions of images.

- Use the chain rule

$x_1$	$x_2$	$x_3$
$x_a$	$x_r$	$x_b$
$x_z$	$x_s$	$x_g$

image

pixel recurrent neural networks

$$p(x_1, \dots, x_d) = p(x_1) p(x_2|x_1) \dots \\ \dots p(x_d|x_1, \dots, x_{d-1})$$

$$p(x_k | x_1, \dots, x_{k-1}) = RNN(x_1, \dots, x_{k-1})$$

cool, but slow to generate

- $p(x_1, \dots, x_d) = p(x_1) \dots p(x_d)$

too restrictive

- Mixture of several Gaussians (GMM)

still too restrictive

- Mixture of infinitely many Gaussians

$$p(x) = \int p(x|t) p(t) dt$$

Using CNNs with a mixture  
of Gaussians

$$p(x) = \int p(x|t) p(t) dt$$

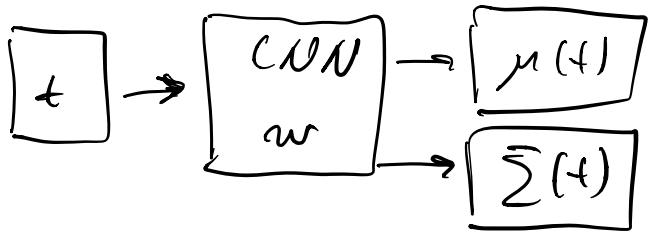
prior:  $p(t) = N(0, I)$

likelihood:  $p(x|t) = N(\mu(t), \Sigma(t))$

if:  $\mu(t) = Wt + b$ ,  $\Sigma(t) = \Sigma_0$  get PPCA  
(metk 2)

But if  $x$  is image, why not

$$\mu(t) = CNN_z(t) \quad \Sigma(t) = CNN_z(t)$$



$$p(x|w) = \int p(x|t, w) p(t) dt$$

$$p(t) = N(0, I)$$

$$p(x|t, w) = N(\mu(t, w), \Sigma(t, w))$$

diag( $\Sigma^2(t, w)$ )

Scaling up Expectation Maximization

$$\max_w p(x|w) = \int p(x|T, w) p(T) dt$$

Latent Variable model - use EM!

$$\log p(x|w) \geq L(w, q) - E \text{ step}$$

$$\underset{w, q}{\text{maximize}} \quad L(w, q) \quad - M \text{ step}$$

But E-step is intractable: Need to compute  $p(T|X, w)$

MLMC - ?

$$\mathbb{E}_g \log p(x, T | w) \approx \frac{1}{M} \sum_{s=1}^M \log p(x, T_s | w)$$

$$T_s \sim q(T)$$

Variational EM!

$$\log p(x | w) \geq L(w, g)$$

$$\underset{w, g}{\text{maximize}} \quad L(w, g)$$

$$\text{subject to: } q_i(x_i) = \hat{q}_i(f_{i1}) \dots \hat{q}_i(f_{im})$$

but again intractable.

Scaling up Variational EM

$$\underset{w, g}{\text{maximize}} \quad L(w, g)$$

$$\text{subject to: } q_i(x_i) = \hat{q}_i(f_{i1}) \dots \hat{q}_i(f_{im})$$



$$\underset{w, q_1, \dots, q_M}{\text{maximize}} \quad L(w, q_1, \dots, q_M)$$

$$\text{subject to: } q_i(x_i) = \hat{q}_i(f_{i1}) \dots \hat{q}_i(f_{im})$$

$\Downarrow$  another approximation

maximize  $L(w, \phi_1, \dots \phi_N)$

$w, m_1, \dots m_N$   
 $s_1, \dots s_N$

subject to:  $g_i(t_i) = N(m_i, \text{diag}(s_i^2))$

But this may ~100 parameters

for each training object

And is not clear what is used for test objects

$\Downarrow$

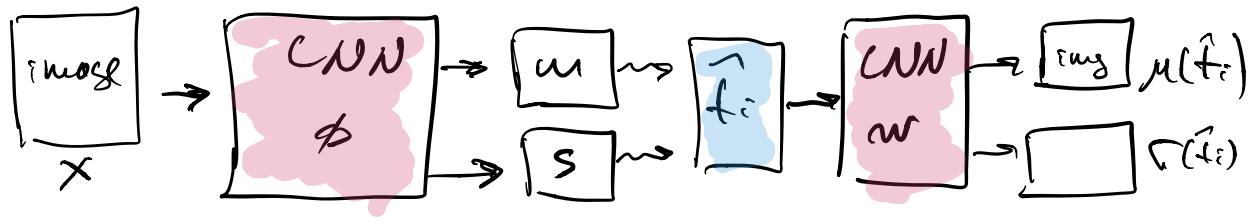
maximize  $L(w, \phi_1, \dots \phi_N)$   
 $w, \phi$

subject to  $g_i(t_i) = N(m(x_i, \phi), \text{diag}(s^2(x_i, \phi)))$

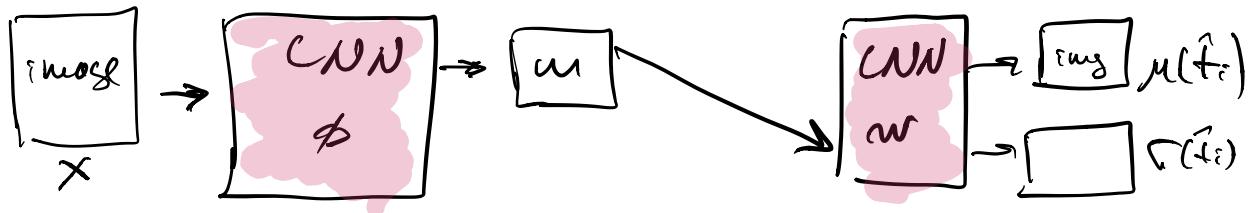


$$L(w, \phi_1, \dots \phi_N) = \sum_i I_{g_i} \log \frac{p(x_i | t_i, w) p(t_i)}{g_i(t_i)}$$

$$\hat{t}_i \sim N(m(x_i, \phi), \text{diag}(s^2(x_i, \phi)))$$



if  $s(x)=0$  then  $\hat{f}_i = m(x_i, \phi)$ : usual  
autoencoder



$$\max_{w, \phi} \sum_i \mathbb{E}_{q_{\phi}(f_i)} \log \frac{p(x_i | f_i, w) p(f_i)}{q_{\phi}(f_i)} =$$

$$= \sum_i \mathbb{E}_{q_{\phi}(f_i)} \log p(x_i | f_i, w) + \underbrace{\mathbb{E}_{q_{\phi}(f_i)} \log \frac{p(f_i)}{q_{\phi}(f_i)}}_{-KL(q_{\phi}(f_i) || p(f_i))}$$

$$- \|x_i - \mu(f_i)\|^2 + \text{const}$$

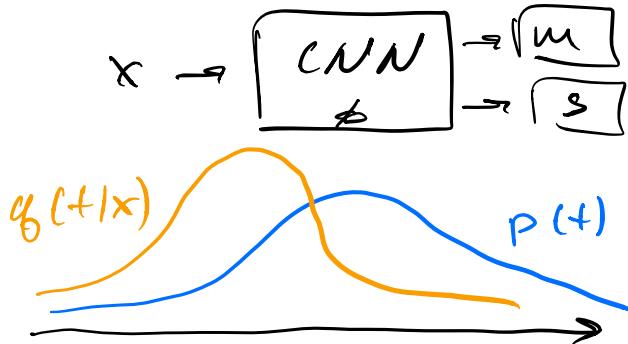
if  $s(x_i) = 1$  for simplicity

approximate posterior reconstruction loss

$$q_{\phi}(f_i) \approx p(f_i | x_i, w)$$

regularization

## Detecting outliers

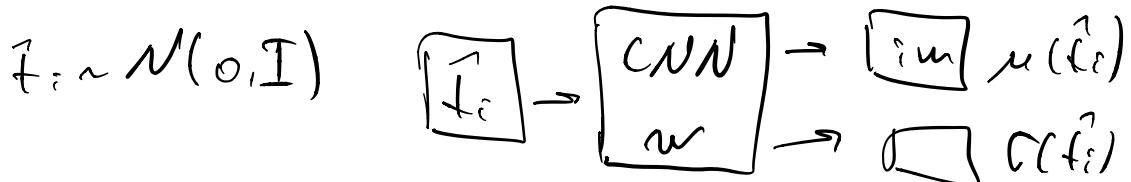


$$KL(q(t|x) \parallel p(t)) = 0.54$$

but for fraud  $KL(q(t|x) \parallel p(t)) = 6.25$   
for example

## Generating new samples

$$p(x|w) = \int p(x|t, w)p(t)dt$$



## Gradient of decoders

$$\max_{w, \phi} \sum_i \mathbb{E}_{q(t_i)} \log p(x_i | t_i, w) - KL(q_\phi(t_i) \parallel p(t_i))$$

easy and analytical

$$KL(q_{\phi}(t_i) \parallel p(t_i)) = \sum_j (-\log q_j(t_i) + \frac{\sigma_j^2(t_i) + \mu_j^2(t_i) - 1}{2})$$

$$f(w, \phi) = \sum_i \mathbb{E}_{q_{\phi}(t_i)} \log p(x_i | t_i, w)$$

$$q_{\phi}(t_i) = q(t_i | x_i, \phi) = \mathcal{N}(\mu_i, \text{diag}(S_i))$$

$$f(w, \phi) = \sum_i \mathbb{E}_{q(t_i | x_i, \phi)} \log p(x_i | t_i, w)$$

$$\nabla_w f(w, \phi) = \nabla_w \sum_i \mathbb{E}_{q(t_i | x_i, \phi)} \log p(x_i | t_i, w) =$$

$$= \nabla_w \sum_i \int q(t_i | x_i, \phi) \underbrace{\log p(x_i | t_i, w)}_{\text{smooth function}} dt_i =$$

$$= \sum_i \int \nabla_w [q(t_i | x_i, \phi) \cdot \log p(x_i | t_i, w)] dt_i =$$

(if this is smooth function)

$$= \sum_i \int q(t_i | x_i, \phi) \nabla_w \log p(x_i | t_i, w) dt_i =$$

$$= \sum_i \mathbb{E}_{q(t_i | x_i, \phi)} \nabla_w \log p(x_i | t_i, w) \approx$$

$\approx \sum_i \nabla_w \log p(x_i | \tilde{t}_i, w)$  [approximate by sampling]  
 $\tilde{t}_i \sim q(t_i | x_i, \phi)$

$$\nabla_w f(w, \phi) \approx \sum_{i=1}^n \nabla_w \log p(x_i | \hat{t}_i, w) \approx$$

$\hat{t}_i \sim g(t_i | x_i, \phi)$  gradient of standard NN  
 $t_i$  - sample

$$\approx \frac{1}{n} \sum_{s=1}^n \nabla_w \log p(x_{is} | \hat{t}_{is}, w)$$

Stochastic gradient of standard NN  
 $i_s \sim U\{1, \dots, N\}$

### Log derivative trick

$$\begin{aligned}\nabla_\phi f(w, \phi) &= \nabla_\phi \sum_i \mathbb{E}_{g(t_i | x_i, \phi)} \log p(x_i | t_i, w) = \\ &= \sum_i \int \nabla_\phi g(t_i | x_i, \phi) \log p(x_i | t_i, w) dt_i = \\ &\quad (\text{only for smooth functions}) \\ &= \sum_i \int \frac{g(t_i | x_i, \phi)}{g(t_i | x_i, \phi)} \nabla_\phi g(t_i | x_i, \phi) \log p(x_i | t_i, w) dt_i \quad \text{①}\end{aligned}$$

$$\nabla \log g(\phi) = \frac{\nabla g(\phi)}{g(\phi)}$$

$$\begin{aligned}& \textcircled{1} \quad \sum_i \int g(t_i | x_i, \phi) \nabla \log g(t_i | x_i, \phi) \log p(x_i | t_i, w) dt_i = \\ &= \sum_i \mathbb{E}_{g(t_i | x_i, \phi)} \nabla \log g(t_i | x_i, \phi) \log p(x_i | t_i, w) \quad \text{②}\end{aligned}$$

(log derivative trick)

approx  
with  
samples

## Reparameterization trick

$$\nabla_{\phi} f(w, \phi) = \sum_i \nabla_{\phi} \mathbb{E}_{q(x_i | x_i, \phi)} \log p(x_i | x_i, w) =$$

$$= \sum_i \mathbb{E}_{q(x_i | x_i, \phi)} \nabla_{\phi} \log q(x_i | x_i, \phi) \log p(x_i | x_i, w) dt_i$$

$\log p(x_i | x_i, w)$  - very big negative number, because net doesn't know new images, so variance of estimation  $\nabla_{\phi} f(w, \phi)$  will be huge.

$$t_i \sim q(f_i | x_i, \phi) = \mathcal{N}(m_i, \text{diag}(S_i^2))$$

$$t_i = \varepsilon_i \odot S_i + m_i = g(\varepsilon_i, x_i, \phi)$$

$$\varepsilon_i \sim p(\varepsilon_i) = \mathcal{N}(0, I) \quad \leftarrow \text{deterministic function}$$

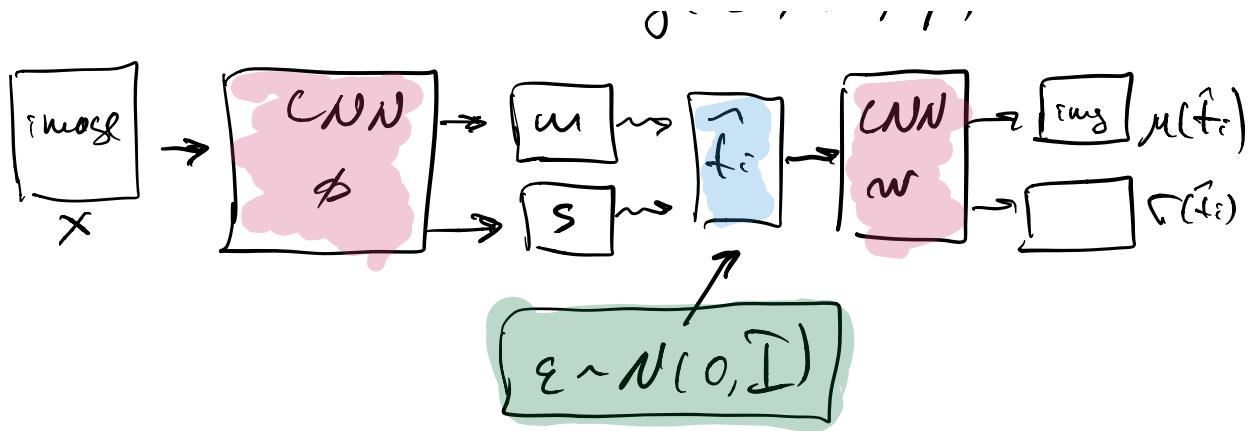
$$\textcircled{1} \sum_i \nabla_{\phi} \mathbb{E}_{p(\varepsilon_i)} \log p(x_i | g(\varepsilon_i, x_i, \phi), w) =$$

$$= \sum_i \int p(\varepsilon_i) \nabla_{\phi} \log p(x_i | g(\varepsilon_i, x_i, \phi), w) d\varepsilon_i =$$

$$= \sum_i \mathbb{E}_{p(\varepsilon_i)} \nabla_{\phi} \log p(x_i | g(\varepsilon_i, x_i, \phi), w)$$

$$p(\varepsilon_i) = \mathcal{N}(0, I)$$

$$t_i = \varepsilon_i \odot S_i + m_i = g(\varepsilon_i, x_i, \phi)$$



## Variational Autoencoder summary

- Infinite mixture of Gaussians
- To learn: EM + approximate  $q$  with Gaussians + stochastic variational inference
- Like plain autoencoder, but with noise and KL regularization
- Generates nice images

# Variational Dropout

## Learning with priors

### Ideal Bayesian learning

- $(X, Y)$  and probabilistic classifier  
 $p(y|X, w)$

- define a reasonable prior  $p(w)$

- Training stage

$$p(w|X, Y) = \frac{p(Y|X, w)p(w)}{\int p(Y|X, w)p(w)dw}$$

↑  
posterior

usually intractable

- Test stage

$$p(y^*|\hat{x}, X, Y) = \int p(y^*|\hat{x}, w)p(w|X, Y)dw$$

## Variational Inference

Inference  $p(w|X, Y)$  - origin problem

Optimization

$$g(w|\phi) = \underset{\phi}{\operatorname{arg\,min}} D(g(w|\phi), p(w|X, Y))$$

variational distribution family

distance ↴ measure between 2 distributions

$$\text{If } D(g(w|\phi), p(w|X, Y)) =$$

$$= KL[g(w|\phi) || p(w|X, Y)]$$

Then optimization problem is equivalent to maximizing ELBO (evidence lower bound)

$$L(\phi) = \int g(w|\phi) \log \frac{P(Y|X, w)p(w)}{g(w|\phi)} dw$$

↓

max  $\neq$

steel retractable  
but ↴

Use properties of ELBO

$$L(\phi) = \int g(w|\phi) \log \frac{P(Y|X, w)p(w)}{g(w|\phi)} dw$$

- We can use mini-batching for optimization since

$$\log P(Y|X, w) = \sum_{i=1}^n \log p(y_i|x_i, w)$$

- We can use MC estimates for computing stochastic gradient - especially effective when reparameterization trick is applicable



$$L(\phi) = \int r(\varepsilon) \log \frac{p(x|\mathbf{z}, w(\varepsilon, \phi)) p(w(\varepsilon, \phi))}{q(w(\varepsilon, \phi) | \phi)} d\varepsilon$$

- The richer is variational family  $\{q(w|\phi)\}_{\phi}$  the better we approximate true posterior
- Can split ELBO into data term and KL-term

$$L(\phi) = \underbrace{\int q(w|\phi) \log p(x|z, w) dw}_{\text{data term}} - \underbrace{- \text{KL}(q(w|\phi) || p(w))}_{\text{regularizer}}$$

## Dropout as Bayesian procedure

- Popular regularization technique
- Injects noise to the weights or activations at each iteration during training
- The magnitude of noise is defined by dropout rate

## Gaussian drop out

We inject multiplicative Gaussian noise with mean  $s$  and variance  $\lambda$   
 $w_{ij} = \theta_{ij} \epsilon_{ij}, \quad \epsilon_{ij} \sim \mathcal{N}(\mu_{ij} | s, \lambda)$

Then we compute stochastic gradient of log-likelihood given noised weight  $W$

$$\begin{aligned} \text{stoch grad}_\theta \log p(\mathbf{y} | \mathbf{x}, \mathbf{w}) &= \\ &= \text{stoch grad}_\theta \log p(\mathbf{y} | \mathbf{x}, \theta \hat{\epsilon}) \\ &\quad \hat{\epsilon} \sim \mathcal{N}(\mu | s, \lambda \mathbf{I}) \end{aligned}$$

But we obtain exactly the same stochastic gradient if we consider

$$\text{stoch grad}_\theta \int \mathcal{N}(\mathbf{w} | \theta, \lambda \theta^2) \log p(\mathbf{y} | \mathbf{x}, \mathbf{w}) d\mathbf{w}$$

where  
 $\mathcal{N}(\mathbf{w} | \theta, \lambda \theta^2) = \prod_{ij} \mathcal{N}(w_{ij} | \theta_{ij}, \lambda \theta_{ij}^2)$

$$\text{stoch grad}_\theta \int \mathcal{N}(\epsilon | s, \lambda) \log p(\mathbf{y} | \mathbf{x}, \mathbf{w}) d\mathbf{w}$$

= stochastic gradient  $\log p(Y|X, W)$   
 where  $\hat{\epsilon} \sim \mathcal{N}(\epsilon | 0, 2I)$

Gaussian dropout optimizes

$$\int q(W|\Theta, \lambda) \log P(Y|X, W) dW \xrightarrow{\Theta} \max$$

Looks similar to the data term  
 in ELBO with variational distribution

$$q(W|\Theta, \lambda) = \mathcal{N}(W|\Theta, \lambda\Theta^2) = \prod_{ij} \mathcal{N}(w_{ij} | \theta_{ij}, \lambda\theta_{ij}^2)$$

But where is KL-term?

What if KL-term depends only on  $\lambda$ ?

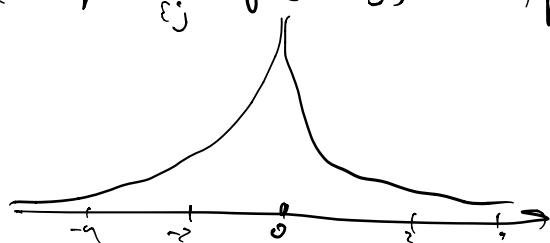
$$KL[q(W|\Theta, \lambda) || p(W)] = R(\lambda)$$

Then Gaussian dropout is exactly  
 variational Bayesian inference!

Log-uniform prior

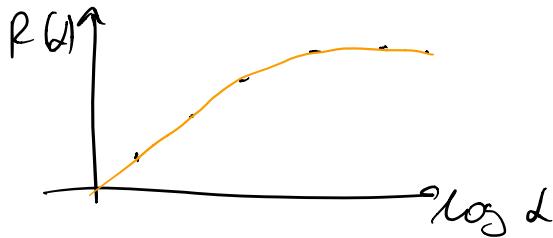
Such prior  $p(W)$  exists and is  
 very interpretable

$$p(W) = \prod_{ij} p(w_{ij}) \quad , p(w_{ij}) \propto \frac{1}{|w_{ij}|}$$



## KL-term

- It can be shown that with such choice of prior KL-term depends only on  $d$
- Although KL is intractable it can be approximated well a differentiable function



## Conclusion

- Dropout prevents overfitting
- It injects noise during training
- Gaussian dropout corresponds to variational Bayesian inference with log-uniform prior
- Now we may construct various generalizations of dropout

## Sparse variational dropout

### Variational dropout

- Gaussian dropout with fixed noise variance is equivalent to optimizing ELBO

$$L(\theta) = \int q(w|\theta, \lambda) \log p(y|x, w) dw - R(\lambda)$$

$$\max_{\lambda}$$

where  $R(\lambda)$  stands for KL-term

- Why not to optimize it w.r.t. noise variance  $\lambda$  as well?

- This will only improve our variational approximation



$$L(\theta) = \int q(w|\theta, \lambda) \log p(y|x, w) dw - R(\lambda)$$

$$\max_{\theta, \lambda}$$

- Note it was not possible until we came with Bayesian interpretation of Gaussian dropout

Our variational approximation for the true posterior

$$q(w|\theta, \lambda) = \prod_j N(w_j | \theta_j, \lambda \theta_j^2)$$

$\Downarrow$  individual drop out rates

$$q_f(w_i | \Theta, \alpha) = \prod_{ij} N(w_{ij} | \theta_{ij}, \alpha_{ij} \theta_{ij}^2)$$

### Sparseification property

Remember that maximum of  $R(\alpha)$  is achieved at regularity  
we also show that

$$\alpha_{ij} \rightarrow \infty \Rightarrow \theta_{ij} \rightarrow 0 \text{ and } \alpha_{ij} \theta_{ij}^2 \rightarrow 0$$

This means that

$$\lim_{\alpha_{ij} \rightarrow \infty} q(w_{ij} | \theta_{ij}, \alpha_{ij}) = \delta(0)$$

we may remove the corresponding weight!

### Sparse variational dropout

- Assign log-uniform prior over the weight  $p(w)$
- Fix variational family of distributions  $q_\phi(w | \Theta, A) = \prod_{ij} N(w_{ij} | \theta_{ij}, \alpha_{ij} \theta_{ij}^2)$
- Perform variational inference

$$L(\theta, A) = \int q(w|\theta, A) \log P(y|x, w) dW - R(A)$$

$\max_{\theta, A}$

- Remove all weights whose  $|w_{ij}| \gg 1$

Up to 33.5% of the weights become irrelevant without accuracy drop!

### Conclusions

- Bayesian dropout allows to remove high redundancy of modern DNN
- Variational Bayesian inference is highly scalable procedure that allows to optimize millions of variational parameters
- One of many examples of successful combination of Bayesian methods and deep learning

Quiz: Categorical Reparameterization

with Gumbel - Softmax

3, 4, 8

read the paper



Week 6

5. 11. 2018

## Gaussian processes & Bayesian optimization

Non parametric methods

Parametric methods

① Define parametric model

$$p(g | X, \Theta)$$

② Find best parameters using MAP estimation:

$$p(\Theta | g, X) \rightarrow \max_{\Theta}$$

Parametric methods: fixed number of parameters

Non-parametric: Number of parameters depends on dataset size

K-nearest neighbors

$$g = \frac{1}{k} \sum_{i=1}^k y_i$$

Neumann - Watson

$$g(x) = \sum_{i=1}^n w_i(x) y_i ; \quad w_i(x) = \frac{k(x, x_i)}{\sum_{j=1}^n k(x, x_j)}$$

## Some kernels

$$K(x_1, x_2) = e^{-\frac{1}{2\sigma^2} \|x_1 - x_2\|^2}$$

$$K(x_1, x_2) = \prod \left[ \begin{cases} \|x_1 - x_2\| < h \end{cases} \right]$$

P:

- limited complexity
- fast inference
- slow learning

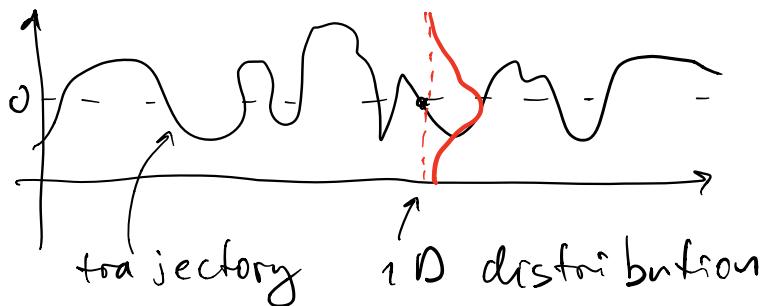
non P:

- arbitrary complex
- Need to process all data for prediction
- Learning: remember all data

## Gaussian processes

### Random process

For any  $x \in \mathbb{R}^d$  assign random variable  $f(x)$



### Gaussian process

Random process  $f$  is Gaussian, if for any finite number points, their joint distribution is normal

$f_n \in \mathcal{M}$ ,  $\forall x_1, x_2, \dots, x_n \in \mathbb{R}^d$

$$\begin{pmatrix} f(x_1) \\ \vdots \\ f(x_n) \end{pmatrix} \sim N$$

Parameters:

$$\mathbb{E} f(x) = m(x)$$

$$\text{Cov}\{f(x_1), f(x_2)\} = K(x_1, x_2)$$

finally:

$$\begin{pmatrix} f(x_1) \\ \vdots \\ f(x_n) \end{pmatrix} \sim N\left(\begin{pmatrix} m(x_1) \\ \vdots \\ m(x_n) \end{pmatrix}, \begin{pmatrix} K(x_1, x_1) & \dots & K(x_1, x_n) \\ \vdots & \ddots & \vdots \\ K(x_n, x_1) & \dots & K(x_n, x_n) \end{pmatrix}\right)$$

### Stationary process

Random process  $f$  is stationary, if its finite-dimensional distributions depend only on relative position of the points

Gaussian process is stationary, if

$$m(x) = \text{const}$$

$$K(x_1, x_2) = \hat{K}(x_1 - x_2)$$

Variance:  $\text{Var}[f(x)] = \hat{K}(0)$



## kernel

$$\text{RBF} : \hat{K}(x_1, -x_2) = C^2 \exp\left(-\frac{(x_1 - x_2)^2}{2\ell^2}\right)$$

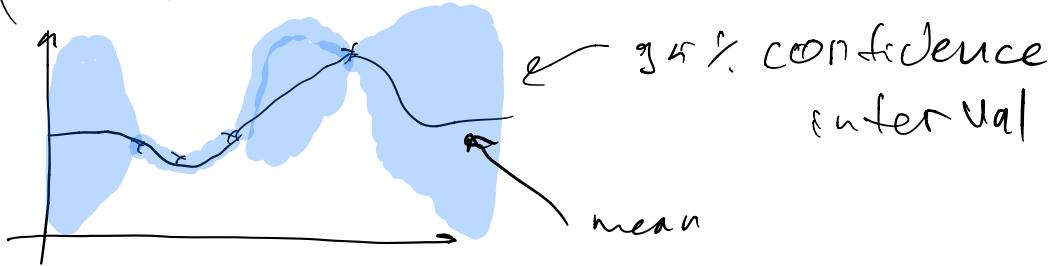
$$\text{Rational Quadratic} : \hat{K}(x_1, -x_2) = C^2 \left(1 + \frac{(x_1 - x_2)^2}{2\alpha\ell^2}\right)^{-\alpha}$$

$$\text{White noise} : \hat{K}(x_1, -x_2) = C^2 \delta(x_1, -x_2)$$

## GP for machine learning

### Task

$$\begin{pmatrix} f(x_1) \\ \vdots \\ f(x_n) \end{pmatrix} \Rightarrow f(x) \approx ? \Rightarrow p(f(x) | f(x_1), \dots, f(x_n)) = ?$$



### Prediction

$f(x)$  is a Gaussian Process with stationary prior,  $m(x) = 0$

$$p(f(x) | f(x_1), \dots, f(x_n)) = \frac{P(f(x), f(x_1), \dots, f(x_n))}{P(f(x_1), \dots, f(x_n))} =$$

$$= \frac{N(f(x), f(x_1), \dots, f(x_n) | 0, \bar{C})}{N(f(x_1), \dots, f(x_n) | 0, C)} \quad (\textcircled{=})$$

$$C = \begin{pmatrix} K(0) & \dots & K(x_c - x_n) \\ \vdots & \ddots & \vdots \\ K(x_n - x_c) & \dots & K(0) \end{pmatrix}$$

$$\bar{C} = \begin{pmatrix} K(0) & K^\top \\ K & C \end{pmatrix} \quad K = \begin{pmatrix} K(x - x_1) \\ \vdots \\ K(x - x_n) \end{pmatrix}$$

$$(\textcircled{\textcircled{}}) N(f(x) | \mu, \sigma^2)$$

$$\mu = K^\top C^{-1} f$$

$$\sigma^2 = K(0) - K^\top C^{-1} K$$

$$f = \begin{pmatrix} f(x_1) \\ \vdots \\ f(x_n) \end{pmatrix}$$

### Preprocessing

- Far from data :  $E[f(x)] = 0, \text{Var}[f(x)] = K(0)$
- Remove trend, sea sohatly
- Subtract mean and normalize

### Derivation of main formula

$$P(f(x) | f(x_1), \dots, f(x_n)) = \frac{P(f(x), f(x_1), \dots, f(x_n))}{P(f(x_1), \dots, f(x_n))}$$

$$\begin{aligned}
&= \frac{\mathcal{N}(\{f^*\} | 0, \begin{bmatrix} K(0) & K^T \\ K & C \end{bmatrix})}{\mathcal{N}(f | 0, C)} \quad \text{and} \\
&\propto \exp\left(-\frac{1}{2} \{f^*\}^T \begin{bmatrix} K(0) & K^T \\ K & C \end{bmatrix}^{-1} \begin{bmatrix} f^* \\ f \end{bmatrix} + f^T C^{-1} f\right) \\
&= \exp\left(-\frac{1}{2} \{f^*\}^T \begin{bmatrix} K(0) & K^T \\ K & C \end{bmatrix}^{-1} \begin{bmatrix} f^* \\ f \end{bmatrix} + f^T C^{-1} f\right) = \\
&\quad \boxed{\begin{bmatrix} d & b^T \\ b & A \end{bmatrix}} \\
&= \exp\left(-\frac{1}{2} \left[ d(f^*)^2 + 2b^T f \cdot f^* + \text{const}\right]\right) \quad \text{and} \\
&\propto \exp\left(-\frac{1}{2d} \left\{ f^* + \frac{b^T f}{d}\right\}^2\right) = \\
&\quad \boxed{\mu = -\frac{b^T f}{d} \quad C^2 = d^{-1}}
\end{aligned}$$

$$= \mathcal{N}(f^* | \mu, C^2)$$

Let's find  $b, d$

$$\begin{bmatrix} K(0) & K^T \\ K & C \end{bmatrix} \begin{bmatrix} d & b^T \\ b & A \end{bmatrix} = \begin{bmatrix} I & 0 \\ 0 & I \end{bmatrix}$$

$$\underbrace{|K(0)|d + K^T b}_{=1} \Rightarrow |K(0)|d - K^T C^{-1} K d = 1$$

$$\begin{cases} \mathbf{K}(0) \mathbf{b}^\top + \mathbf{K}^\top \mathbf{A} = \mathbf{0} \\ \underline{\mathbf{d}^\top \mathbf{K} + \mathbf{C}^\top \mathbf{b} = 0} \Rightarrow \mathbf{b} = -\mathbf{C}^{-1} \mathbf{K} \mathbf{d} \\ \mathbf{K} \mathbf{b}^\top + \mathbf{C} \mathbf{A} = \mathbf{I} \end{cases}$$

$\mathbf{d} = \frac{1}{\mathbf{I}(0) - \mathbf{K}^\top \mathbf{C}^{-1} \mathbf{K}}$

$\Rightarrow \mathbf{b} = -\frac{\mathbf{C}^{-1} \mathbf{K}}{\mathbf{I}(0) - \mathbf{K}^\top \mathbf{C}^{-1} \mathbf{K}}$

$$\Rightarrow \mathbf{f}^2 = \mathbf{I}(0) - \mathbf{K}^\top \mathbf{C}^{-1} \mathbf{K}$$

$$\mu = -\frac{\mathbf{b}^\top \mathbf{f}}{\mathbf{d}^\top \mathbf{f}} =$$

## Nuance of GP

### Noisy observations

$\hat{f}(x) = f(x) + \varepsilon$  - independent gaussian noise

$$\varepsilon \sim \mathcal{N}(0, s^2)$$

$$m(x) = 0$$

$$k(x_i - x_j) = k(x_i - x_j) + s^2 I_{\{x_i = x_j\}}$$

## Kernel parameters

$$\hat{k}(x_i, -x_i) = \gamma^2 \exp\left(-\frac{(x_i - x_i)^2}{2\ell^2}\right) + s^2 \prod_{j \neq i} \{x_i = x_j\}$$

$$p(f(x_1), f(x_2), \dots, f(x_n) | \gamma^2, \ell, s^2) \rightarrow \max_{\gamma^2, \ell, s^2}$$

$$N(f(x_1), f(x_2), \dots, f(x_n) | \varphi, C) \rightarrow \max_{\gamma^2, \ell, s^2}$$

Optimize with gradient ascent

## Classification

$$y \in \{-1, 1\}$$

Latent process:  $f(x)$

$$\text{Class probabilities: } p(y|f) = \frac{1}{1 + \exp(-yf)}$$

## Training:

- Approximate latent process from data

$$p^*(f(x)) = p(f(x) | y_1(x_1), \dots, y_n(x_n))$$

- Compute predictions

$$\pi(x) = \int p(y(x) | f(x)) \cdot p^*(f(x)) d f(x)$$

## Inducing inputs

$$\mu = k^T C^{-1} f \quad \mathcal{O}(n^3)$$

$$C^2 = I(n) - k^T C^{-1} k$$

Idea:

- replace dataset with small number of points (like sum)
- Predict using those points

Speed:

- Precomputing  $\mathcal{O}(m^2n)$
- Mean  $\mathcal{O}(m)$
- Variance  $\mathcal{O}(m^2)$

## Bayesian optimization

## Black-box optimization

$$f(x) \rightarrow \max_x$$

Gradient is known:

- Gradient descent with restarts

Gradient is unknown:

- Numerically estimate gradient
- Grid search / Random search

Cases:

- $x$  geographic coordinates,  $f(x)$  - amount of oil  
     $\uparrow$  sample =  $10^6$  \$
- $x$  hyperparameters of NN,  $f(x)$  - objective  
     $\uparrow$  sample = 10 hours

- $x$  drug,  $f(x)$  - effectiveness against disease  
• sample = 2 month, \$10k, life of a rat

Goal: Optimize with minimum number of trials

Surrogate model  $\hat{f} \approx f$

- Approximates true function
- Cheap to evaluate

Acquisition function:  $\mu(x)$

- Estimate profit for optimization
- Uses surrogate model

Surrogate model

$$\hat{f} \approx f$$

Should model arbitrary complex functions  
 $\Rightarrow$  Nonparametric method

Profitable to estimate uncertainty

$\Rightarrow$  Gaussian process

Acquisition function

Exploration:

Search in regions with high uncertainty

Exploitation:

Search in regions with high estimated value

### Maximum probability of improvement (MPI)

Current best value:  $f^*$

$$\mu(x) = \text{IP}(\hat{f}(x) \geq f^* + \varepsilon) = \Phi\left(\frac{\mathbb{E}\hat{f}(x) - f^* - \varepsilon}{\text{Var}[\hat{f}(x)]}\right)$$

Works well if value of maximum is known

### Upper confidence bound (UCB)

$$\mu(x) = \mathbb{E}\hat{f}(x) + \sqrt{\text{Var}[\hat{f}(x)]}$$

### Expected improvement (EI)

$$\begin{aligned}\mu(x) &= \mathbb{E} \max(f(x) - f^*, 0) = \\ &= \text{Var}[\hat{f}(x)] \cdot [z \cdot \Phi(z) + \phi(z)]\end{aligned}$$

$$z = \frac{\mathbb{E}\hat{f}(x) - \mu(x)}{\sqrt{\text{Var}[\hat{f}(x)]}}$$

most widely used

Example:

Start with few points  
while not converged:

1. Train GP
2. Find maximum of  $\mu(x)$   
using e.g. gradient ascent
3. Evaluate function at  
maximum of  $\mu(x)$

## Random search vs gaussian processes

RS

- Parallelizable
- Needs many more points for high dimensions

Any function

GP

- Hard to parallelize experiments
- Requires less points on average

Each evaluation is expensive

Train BM << Train NN

## Applications of Bayesian optimization

### Hyperparameter tuning

Network parameters

- Number of layers
- Layer sizes
- Dropout on/off

- Batch normalization on/off

Training parameters

- Learning rate
- Momentum

Usually finds better optima than when tuned by hand

Honest comparison with other methods in research

### Discrete and continuous variables

- Treat discrete variables as continuous when fitting process
  - Maximize  $\mu(x)$  for each possible value of discrete variables
  - Multi-armed bandit: all variables are discrete
-

Final Task

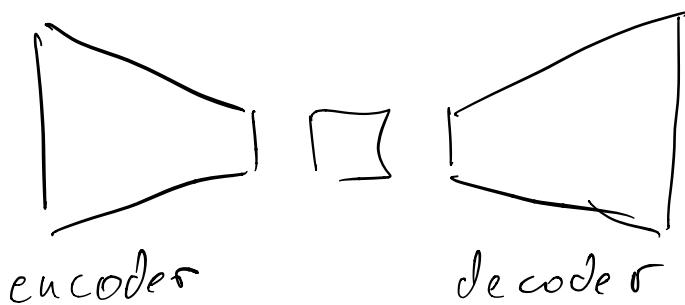
28.12.2018

6.01.2019

- ①  $2^5 \cdot$  samples from VAE

vae

13 poems



- ② Search procedure.