

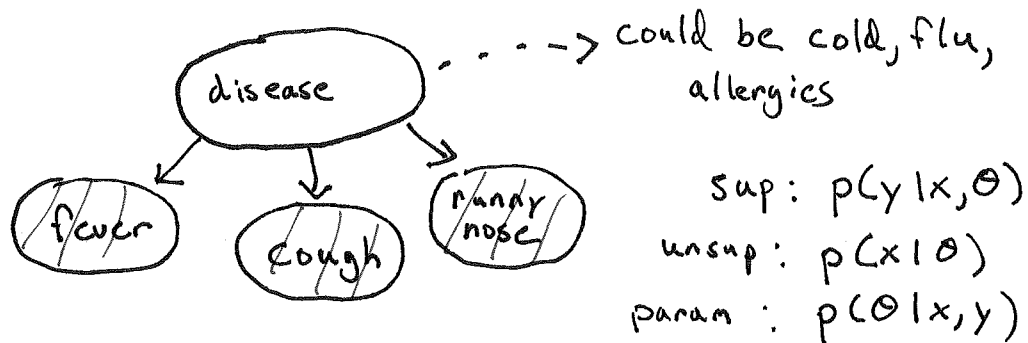
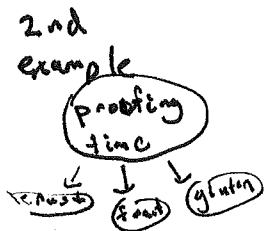
□ What is machine learning? AI meets data/stats

Expert Knowledge
sophisticated reasoning
and search

learn, train data
simpler models

ML combines these

□ We're going to focus on one particular area (take 181 for breadth): generative models of the world



① disease generates symptoms (but then we infer diseases from symptoms)

② this picture is a "graphical model" shows generative dependencies

③ we could also learn the params from data (e.g. what the diseases are) } supervised vs. unsupervised

④ Write Bayes rule:

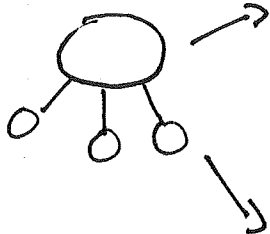
$$p(\text{disease} | \text{data}) \propto p(\text{fever} | \text{disease}) p(\text{cough} | \text{disease}) \cdots p(\text{disease})$$

$$\propto \underbrace{\prod p(\text{symptom} | \text{disease})}_{\substack{\text{likelihood} \\ p(\text{data} | \text{model})}} \underbrace{p(\text{disease})}_{\text{prior}}$$

"prob of data"

- in some ways like rule-based systems, need models
- note Bayes rule can be used with some symptoms missing (missing data). By modeling everything we can answer any question (drawbacks & benefits) "reasoning under uncertainty"

(4)



Models: Simple discrete models, Gaussian models, Markov random fields, GLM, Factor analysis, HMMs, Latent Dirichlet Allocation, Deep models

Inference & Learning: Belief Propagation, Junction Tree, Variational, MCMC, LP Relaxations

We'll alternate between models and methods.

□ Format of class / course

- Class: first 10 min. reading check: be on time, have a web-enabled device
- Lecture & turn to neighbor
- Course: 5 psets, first today! midterm
Math & code reviews → section: stanchet training
- Final Project! (UAI/ICML/KDD) Turn to your partner and talk about research / why you're in this class
- piazza! sign-up.

□ Remind Class

- HW1 release today
- HVO due Friday
- tryout quiz
- Help with piazza

□ Discrete Models:

Take on countable values $\{0, 1\}$, $\{\text{cold, flu, asthma}\}$

Today: some of the simplest discrete models,
how to manipulate, "tactics" Murphy 3.3

Running Example: $p(\text{heads}) = \theta$

□ Prior #1: Three manufacturing processes: $\theta = .4, .5, .6$

→ Expert knowledge told us this fact. v.p. .1 .8 .1
Mixture models coming later. (empirical dist.)

for now: $p(\theta) = .1 \delta(\theta=.4) + .8 \delta(\theta=.5) + .1 \delta(\theta=.6)$

Likelihood: $\text{Bin}(N, N_1, \theta) = \underbrace{\binom{N}{N_1}}_{\substack{\text{\# of ways} \\ \text{to get } N_1 \text{ heads} \\ \text{in } N; \text{ constant} \\ \text{w.r.t. } \theta}} \underbrace{\theta^{N_1} (1-\theta)^{N-N_1}}_{\theta^x (1-\theta)^x \text{ per coin}}$

Suppose we observe N_0, N_1, \dots how do we
perform inference? $\Rightarrow p(\theta | x)$

(N_0 = tails)

Discrete 2-2

Several Options:

□ Maximum Likelihood [MLE]

$$\theta_{MLE} = \underset{\theta}{\operatorname{argmax}} p(N_0, N_1 | \theta) = \underset{\theta}{\operatorname{argmax}} \log p(N_0, N_1 | \theta)$$

$$= \log \binom{N_1 + N_0}{N_1} + N_1 \log \theta + N_0 \log(1 - \theta)$$

doesn't depend
on θ ...

in fact we can
ignore, pretend
one sequence
of heads

take derivatives wrt θ

$$\frac{d(\cdot)}{d\theta} = \frac{N_1}{\theta} + \frac{N_0}{1-\theta} (-1) = 0$$

$$\theta = \frac{N_1}{N_0 + N_1}$$

N = surprises!

(?) What if you needed to predict whether a coin was going to be heads? would you guess w.p. θ ? (No partial credit) No!

• Proportion right if you always guess H : θ

• " " " " guess w.p. θ : $\theta^2 + (1-\theta)^2 \Rightarrow 1 - 2\theta + 2\theta^2$

If we plug-in .6 : .6 vs. .52

\Rightarrow Predicting \neq Decisions!

□ MAP / full posterior [discrete_coins.m] $p(\theta | x) \propto p(x | \theta) p(\theta)$

$$p(\theta = .4 | N_0, N_1) \propto \underbrace{\binom{N_0 + N_1}{N_1}}_{\text{can ignore}} \underbrace{(.4)^{N_1} (1-.4)^{N_0}}_{\text{unnormalized likelihood}} \underbrace{(.1)}_{\text{prior}}$$

full posterior
MAP
 $\underset{\theta}{\operatorname{argmax}} p(\theta | N)$
 \Downarrow
not MLE!

$$p(\theta = .45 | N_0, N_1) = 0 \quad [\text{Big effect}]$$

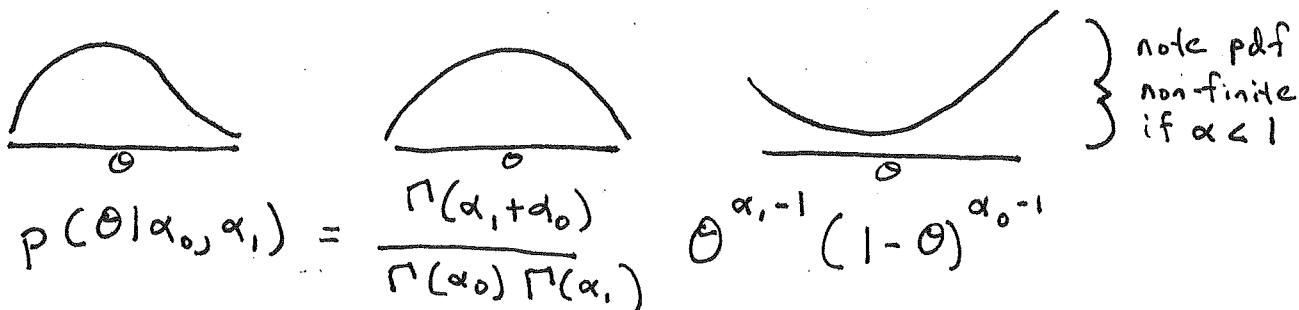
$$p(\theta = .51) p(\theta = .6 | 1) \rightarrow \text{same as above}$$

normalize after

(?) when will $\text{Map} = \text{MLE}$

□ Prior #2

Beta distribution:



$$p(\theta | \alpha_0, \alpha_1) = \frac{\Gamma(\alpha_0 + \alpha_1)}{\Gamma(\alpha_0) \Gamma(\alpha_1)}$$

$$\theta^{\alpha_1 - 1} (1 - \theta)^{\alpha_0 - 1}$$

here, α_0, α_1 are
pseudocounts

plot_beta

$$\begin{aligned} p(\theta | N_0, N_1) &= \frac{p(N_0, N_1 | \theta) p(\theta)}{p(N_0, N_1)} \quad \left\{ \begin{array}{l} \text{constant} \end{array} \right. \\ &= \frac{\Gamma(\alpha_0 + \alpha_1)}{\Gamma(\alpha_0) \Gamma(\alpha_1)} \theta^{\alpha_1 - 1} (1 - \theta)^{\alpha_0 - 1} \theta^{N_1} (1 - \theta)^{N_0} \quad (\text{constant}) \\ &= \theta^{N_1 + \alpha_1 - 1} (1 - \theta)^{N_0 + \alpha_0 - 1} \end{aligned}$$

only part with θ . need to integrate to 1.

- oh but it has form of beta, so:

$$= \frac{\Gamma(N_1 + N_0 + \alpha_0 + \alpha_1)}{\Gamma(N_1 + \alpha_1) \Gamma(N_0 + \alpha_0)} \theta^{N_1 + \alpha_1 - 1} (1 - \theta)^{N_0 + \alpha_0 - 1}$$

(?) A \curvearrowright distribution corresponds to a "sparse" prior that things are either 0 or 1. How does this change if we had:

→ 10 counts to place

→ 8 heads, 1 tails, where would the 10th go

No more sparsity once we have evidence →
no "noise" model

□ Predictive Distribution

$$\begin{aligned}
 p(x|N_0, N_1) &= \int p(x|N_0, N_1, \theta) p(\theta|N_0, N_1) d\theta \\
 &= \int p(x|\theta) p(\theta|N_0, N_1) d\theta = \int \theta p(\theta|N_0, N_1) d\theta \\
 &= \mathbb{E}_{\theta \sim p(\theta|N_0, N_1)} [\theta] = \frac{\alpha_1 + N_1}{\alpha_0 + \alpha_1 + N_0 + N_1} \Rightarrow \text{mean of Beta dist.}
 \end{aligned}$$

note: doesn't care about ↻ or ↺

□ Marginal Likelihood

$p(\text{data}|\alpha)$ = prob of model

$$\begin{aligned}
 p(N_0, N_1) &= \int_{\theta} p(x_1, \dots, x_{N_1}|\theta) p(\theta) d\theta \Rightarrow \\
 &= \int_{\theta} \frac{\Gamma(\alpha_1 + \alpha_0)}{\Gamma(\alpha_0)\Gamma(\alpha_1)} \theta^{\alpha_1 - 1} (1-\theta)^{\alpha_0 - 1} d\theta = \frac{\Gamma(\alpha_1 + N_1) \Gamma(\alpha_0 + N_0)}{\Gamma(\alpha_1 + N_1 + \alpha_0 + N_0)}
 \end{aligned}$$

□ Extensions: What can we do with coins?

→ many coins, correlated θ ? a model of binary data

→ many-sided coins? a model of categorical data

$$p(\theta|\alpha) = \frac{\Gamma(\sum \alpha_k)}{\prod_k \Gamma(\alpha_k)} \prod_k \theta_k^{\alpha_k - 1} \quad \text{"Dirichlet"}$$

$$p(\text{data}|\theta) = \frac{(\sum x_k)!}{\prod_k x_k!} \prod_k \theta_k^{x_k} \quad \text{"multinomial"}$$

□ Think of a language with distribution over words

the	.15
some	.35
Squid	.05
bow-tie	.05
octopus	.05

Change the model:

$$p(x|\alpha) = \int \text{Mult}(x|\theta) p(\theta|\alpha) d\theta$$

$$p(\theta|x, \alpha) = \frac{1}{Z} \prod \theta_k^{x_k + \alpha_k - 1}$$

$$\frac{(\sum_k x_k + \alpha_k - 1)!}{\prod (x_k + \alpha_k)} \quad \leftarrow \quad \frac{\Gamma(\sum x_k + \alpha_k)}{\prod \Gamma(x_k + \alpha_k)}$$

(?) How to "train"? (inference over parameters)

Coming weeks

Δ inv orthogonal matrix
 Δ eigen decomp $\Sigma = U \Delta U^T$
 Δ pos. def $\Delta > 0$

Gaussians
3-1

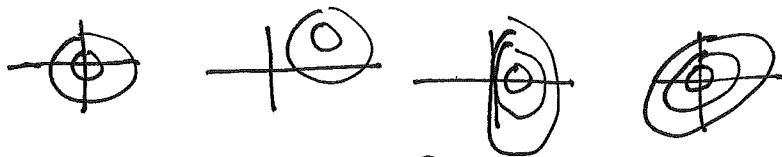
HVO proves important properties

Project Ideas

Review Gaussians

$$N(x | \mu, \Sigma) = \underbrace{(2\pi)^{-\frac{D}{2}} |\Sigma|^{-\frac{1}{2}}}_Z \exp \left\{ \underbrace{-\frac{1}{2} (x - \mu)^T \Sigma^{-1} (x - \mu)}_{\text{important part quadratic form}} \right\}$$

(?) What do μ, Σ look like for:



Decompose: $(x - \mu)^T \Sigma^{-1} (x - \mu) = (x - \mu)^T \left(\sum_d \frac{1}{\lambda_d} u_d u_d^T \right) (x - \mu)$

$= \sum_d \frac{1}{\lambda_d} \underbrace{[(x - \mu)^T u_d]}_{y_d \text{ scalar}} \underbrace{[u_d^T (x - \mu)]}_{\text{normal}}$

\rightarrow gives ellipse for each direction d and elongation λ

(?) suppose $u_1 = [1 \ 0]$ $\lambda_1 = 3$

$u_2 = [0 \ 1]$ $\lambda_2 = 0.1$

(?) Tiled gaussians
 $\mu \rightarrow \tilde{\mu}$
 $\Sigma \rightarrow \tilde{\Sigma}$

what do we get? u_1 picks out 1st dim

Manipulating: Stretch, rotate, shift [Change of variables]

let $y = Ax + b$ $x \sim N(0, I)$ $p(y)$

recall: $y = f(x)$

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

$p(x) dx \rightarrow p(f^{-1}(y)) \left| \frac{d}{dy} f^{-1}(y) \right| dy$
 $p(y) = p(x) \left| \frac{dx}{dy} \right|$

here: $x = A^{-1}(y - b)$

$\left| \frac{dx}{dy} \right| = |A^{-1}|$

So we get

$$\begin{aligned}
 p(y) &= (2\pi)^{-\frac{D}{2}} |A|^{-1} \exp\left\{-\frac{1}{2} x^T x\right\} \leftarrow \text{initial} \\
 &= \dots \dots \dots \left\{-\frac{1}{2} (A^{-1}(y-b))^T (A^{-1}(y-b))\right\} \\
 &\qquad \qquad \qquad \left\{-\frac{1}{2} (y-b)^T \underset{\substack{\downarrow \\ \mu_y = b}}{A^{-T} A^{-1}} (y-b)\right\} \\
 &\qquad \qquad \qquad \qquad \qquad \qquad \qquad \qquad \searrow \Sigma_y = A^T A
 \end{aligned}$$

But wait, easier for Gaussians.

$$\mathbb{E}[y] = \mathbb{E}(Ax+b) = A \mathbb{E}(x) + b = b$$

$$\text{Cov}[y] = \mathbb{E}(\cancel{(y-b)}^T \cancel{(y-b)}) = A^T A \text{ (exercise)}$$

} true in general.
but completely
define gaussians
(first 2 moments)

□ Now write $\Sigma = U^T \Lambda U$ and $\mu = b$

$$y = U \Lambda^{\frac{1}{2}} x + b$$

$\uparrow \quad \quad \uparrow \quad \quad \nwarrow$
 rotate scale shift

plot
bigauss 2

□ Detour: High-dim Gaussians

$$x \sim N(0, \frac{1}{D} I)$$

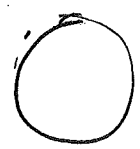
(?) What is expected length $\|x\|^2$?

$$\mathbb{E}[\|x\|^2] = \mathbb{E}\left[\sum_a x_a^2\right] = D \sigma^2 = D \frac{1}{D} = 1$$

(?) What is the variance? $\mathbb{E}[x^4] = 3\sigma^4$

$$\text{var}[\|x\|^2] \Rightarrow \mathbb{E}[x^4] - \mathbb{E}[x^2]^2 = 3\sigma^4 - \sigma^4 = 2\sigma^4 = \frac{D^2}{D^4} = \frac{2}{D}$$

(?) implications?



□ Key Formulas for MVN

$$\mu = \begin{bmatrix} \mu_1 \\ \mu_2 \end{bmatrix} \quad \Sigma = \begin{bmatrix} \Sigma_{11} & \Sigma_{12} \\ \Sigma_{21} & \Sigma_{22} \end{bmatrix}$$

$$p(x_1, x_2) = \mathcal{N}(\begin{smallmatrix} x_1, x_2 \\ \mu, \Sigma \end{smallmatrix})$$

Marginalization: $p(x_1) = \int_{x_2} p(x_1, x_2)$

$$\begin{aligned} p(x_1) &= \int_{x_2} \frac{1}{Z} \exp \left\{ (x_1 - \mu_1)^T \Sigma_{11}^{-1} (x_1 - \mu_1) + 2(x_1 - \mu_1)^T \Sigma_{12}^{-1} (x_2 - \mu_2) \right. \\ &\quad \left. + (x_2 - \mu_2)^T \Sigma_{22}^{-1} (x_2 - \mu_2) \right\} \\ &= p(x_1) \int_{x_2} p(x_2 | x_1) = \mathcal{N}(x_1 | \mu_1, \Sigma_{11}) \end{aligned}$$

Conditionals: $p(x_1 | x_2)$

$$p(x_1 | x_2) = \mathcal{N}(\mu_{x_1|x_2}, \Sigma_{x_1|x_2}) = \mathcal{N}(\mu_{1|2}, \Sigma_{11|2})$$

$$\mu_{1|2} = \mu_1 + \Sigma_{12} \Sigma_{22}^{-1} (x_2 - \mu_2) \Leftarrow \text{residual [later]}$$

$$\Sigma_{11|2} = \Sigma_{11} - \Sigma_{12} \Sigma_{22}^{-1} \Sigma_{21} \quad (\text{proof in Murphy})$$

Information Form

Instead of Σ and μ , use Σ^{-1} and $\Sigma^{-1}\mu$

This format makes conditioning trivial (Σ_{11}^{-1})
but marginals more complicated.

□ MLE

$$\operatorname{argmax}_{\mu, \Sigma} \sum_{n=1}^N \left\{ -\frac{1}{2} \ln |\Sigma| - \frac{1}{2} (x_n - \mu)^T \Sigma^{-1} (x_n - \mu) \right\}$$

$$\frac{d}{d\mu} : \sum_{n=1}^N (x_n - \mu) = 0 \Rightarrow \mu = \frac{1}{N} \sum_{n=1}^N x_n \quad (\text{sample mean})$$

$$\frac{d}{d\Sigma} \quad [\text{Leave as exercise}] \quad \frac{\partial}{\partial A} \ln |A| = A^{-1}$$

$$\frac{\partial}{\partial A} \operatorname{tr} \{ B A \} = B^T$$

$$\operatorname{tr}(ABC) = \operatorname{tr}(CAD) = \operatorname{tr}(BCA)$$

$$\Sigma = \frac{1}{N} \sum_{n=1}^N x_n x_n^T \quad (\text{sample covariance})$$

□ Conjugate Priors [just mean]

□ Predictive

□ Regression

inputs x , outputs y
for now \mathbb{R} in \mathbb{R}

□ Gaussian Noise Model

$$p(y|x, \theta) = \mathcal{N}(y | w^T x, \sigma^2)$$

Write down likelihood:

$$\begin{aligned} \ell(\theta) &= \log p(D|\theta) \\ &= \sum_n \log p(y_n | x_n, \theta) \\ &= \sum_n \log \left(\frac{1}{\sqrt{2\pi\sigma^2}} \exp \left\{ -\frac{1}{2\sigma^2} (y_n - w^T x_n)^2 \right\} \right) \\ &= \underbrace{-\frac{1}{2\sigma^2} \sum_n (y_n - w^T x_n)^2}_{\substack{\text{scaling} \\ \text{residual sum} \\ \text{of squares}}} - \underbrace{\frac{N}{2} \log(2\pi\sigma^2)}_{\text{constant}} \end{aligned}$$

Note how gaussian \Leftrightarrow squared loss

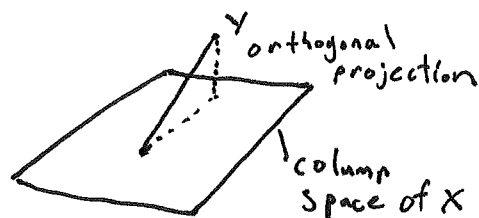
MLE:

$$(y - Xw)^T (y - Xw) = w^T X^T X w - 2 w^T X^T y$$

$$\begin{aligned} \frac{d}{dw} (2X^T X w - 2X^T y) &= 0 \\ w &= (X^T X)^{-1} X^T y \end{aligned}$$

□ Geometry of MLE

$$y = \begin{bmatrix} y \\ y \end{bmatrix} \quad x = \begin{bmatrix} x & x \end{bmatrix}$$



□ Being Bayesian

$$Y = W^T X + \epsilon$$

\uparrow observed \uparrow random \nwarrow fixed \swarrow random

$p(w)$? \rightarrow in HVL, consider $N(w|0, \tau^2)$

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

\uparrow assume 0 centered \nwarrow assume known

$$\propto \exp \left\{ -\frac{1}{2\sigma^2} \|y - Xw\|_2^2 \right\}$$

(?) Why can we ignore z ?
 \hookrightarrow only care about w .

Put a prior: $p(w|w_0, V_0) \propto \exp \left\{ -\frac{1}{2} (w - w_0)^T V_0^{-1} (w - w_0) \right\}$

$$p(w|x, y, \mu, \sigma^2) = \underbrace{N}_{\text{prior}} \underbrace{N}_{\text{likelihood}}$$

$$= N(w | w_N, V_N)$$

After algebra:

$$w_N = V_N V_0^{-1} w_0 + \frac{1}{\sigma^2} V_N X^T y = V_N \left(V_0^{-1} w_0 + \frac{X^T}{\sigma^2} y \right)$$

\uparrow original mean

$$V_N = \left(V_0^{-1} + \frac{X^T X}{\sigma^2} \right)^{-1}$$

\rightarrow original information

Note that $V_0 \rightarrow \infty$
 becomes MLE $(X^T X)^{-1} X^T y$

□ Posterior Predictive

$$\int N(y|x^T w, \sigma^2) N(w|w_N, V_N) dw = N(w_N^T x, \sigma^2 + x^T V_N x)$$

\uparrow noise \uparrow variance on x

Variance depends on x !

$$y = x^T w + \epsilon$$

\uparrow sum of gaussians
 \uparrow fixed variance

Versus $p(y|x, w_{MAP}, \sigma^2)$

Show Demo
Samples

Before: Regression $x \rightarrow y \in \mathbb{R}$

Now $y \in \{0, 1\}$ (Later multiclass)

□ Naive Bayes

$y \sim \text{cat}(\pi)$ $x_j \sim \text{L}(y)$ Generative model of x, y

$$p(x|y=c, \theta) = \prod_{j=1}^D p(x_j|y=c, \theta_{jc})$$

naïve conditional
independence
assumption

1) Multivariate Bernoulli Naive Bayes

$$x_j \sim \text{Ber}(\mu_{jc}) \text{ if } y=c$$

Features are binary-variables

2) Categorical NB

$$x_j \sim \text{Cat}(\underline{\mu}_{jc}) \text{ if } y=c \quad [\text{recall } p(\underline{x}|\underline{\mu}) = \prod_d \mu_d^{x_d}]$$

3) Gaussian NB

$$x \sim \mathcal{N}(\underline{\mu}_c, \Sigma_{\text{diag}}^c) \text{ if } y=c$$

MLE

$$\begin{aligned} p(x_{ij}, y_i | \theta) &= p(y_i | \pi) \prod_j p(x_{ij} | \theta_j) \\ &= \prod_c \pi_c^{(y_i=c)} \prod_{j,c} p(x_{ij} | \theta_{jc})^{(y_i=c)} \end{aligned}$$

$$\underset{\mu, \pi}{\text{argmax}} \log p(x_{ij}, y_i | \theta) = \sum_{c=1}^C N_c \log \pi_c + \sum_j \sum_c \sum_{i: y_i=c} \log p(x_{ij} | \theta_{jc})$$

$$\hat{\pi}_c = \frac{N_c}{N} \quad (\text{class counts})$$

MLE for distribution

Use a factored Prior, assume parameter independence.

$$p(\theta) = p(\pi) \prod_j \prod_c p(\theta_{jc})$$

(?) Which priors should we use here?

π parameters for Cat - Dirichlet

θ parameters for class-conditional, so - Beta or Dirichlet or MVN

Multivariate Bernoulli case

$$\pi \sim \text{Dir}(\alpha)$$

$$\theta_{jc} \sim \text{Beta}(\beta_0, \beta_1)$$

(often set to 1)
assume 1 pseudo-count
for features and class

Recall that ^{posterior} updates have a simple form

$$p(\pi | D) = \text{Dir}(N_1 + \alpha_1, \dots, N_C + \alpha_C)$$

$$p(\theta_{jc} | D) = \text{Beta}((N_c - N_{jc}) + \beta_0, N_{jc} + \beta_1)$$

Exercise: Compute Naïve Bayes for gaussian with prior.

Posterior Predictive

$$p(y | x, D) \propto \underbrace{p(y | D)}_{\int d\pi} \underbrace{\prod p(x_j | y, D)}_{\int d\theta_{jc}}$$

Integrate over params

By earlier class we know this gives the mean of the posterior

$$\bar{\pi}_c = \frac{N_c + \alpha_c}{N + \sum_c \alpha_c}$$

$$\bar{\theta}_{jc} = \frac{N_{jc} + \beta_1}{N_c + \beta_1 + \beta_0}$$

$$p(y=c | x, D) \propto \bar{\pi}_c \prod_{j=1}^n (\bar{\theta}_{jc})^{(x_j=1)} (1 - \bar{\theta}_{jc})^{(x_j=0)}$$

Exponential Form

$$p(y=c|x,D) \propto \pi_c \prod_j \theta_{jc}^{(x_j=1)} \theta_{j\bar{c}}^{(x_j=0)}$$

Take exp of log

$$\propto \exp \left\{ \log \pi_c + \sum_j x_j \log \theta_{jc} + (1-x_j) \log \theta_{j\bar{c}} \right\}$$

$$= \exp \{ w_c^T x + b_c \} \quad \text{where} \quad b = \log \pi_c + \sum_j \log \{ 1 - \theta_{j\bar{c}} \}$$

$$w_{jc} = \log \frac{\theta_{jc}}{1 - \theta_{j\bar{c}}}$$

This shows that NB is a transformation of linear functions of x . The formulation is known as the log-odds of the data.

conversely if we have w we can recover θ

$$\theta = \text{sign}(u) = \sigma(u) = \frac{1}{1 + e^{-u}} \quad (\text{exercise})$$

This is the sigmoid function

Finally to normalize the probability distribution

$$p(y=c|x,D) = \frac{\exp(w_c^T x + b_c)}{\sum_{c'} \exp(w_{c'}^T x + b_{c'})} = \text{softmax} \left(\begin{bmatrix} w_1^T x + b_1 \\ \vdots \\ w_C^T x + b_C \end{bmatrix} \right)_c$$

$$\text{where } \text{softmax}(z)_c = \frac{\exp(z_c)}{\sum_{c'} \exp(z_{c'})}$$

Last two classes linear regression and classification.

Both linear at heart, but different outputs.

Today Exponential Families as a unifying concept.

- Central concept behind many core distributions: normal, bernoulli, categorical, gamma, etc.
 - Provides basis for conjugacy in Bayesian reasoning
 - Central tool for graphical models and variational inf.
 - We will use to^{re} derive logistic regression etc.
- Warning: notation changes here a bit. Be careful w.r.t previous sections.

Exponential Family

$$p(x|\theta) = \frac{1}{Z(\theta)} h(x) \exp \{ \theta^T \phi(x) \}$$

$$= h(x) \exp \{ \theta^T \phi(x) - A(\theta) \}$$

$$A(\theta) = \log(Z(\theta)) = \log \int \exp \{ \theta^T \phi(x) \} dx$$

- θ - the natural parameters (function of "parameters" θ etc.)
- Z, A - the (log) partition function
- $\phi(x)$ - the sufficient statistics (informally features)
- $h(x)$ - scaling (not really important)

A representation is overcomplete if there is

Bernoulli

Exp Families

6-2

$$\text{Ber}(x|m) = \binom{n}{x} (1-m)^{n-x} = \exp \{ x \log m + (1-x) \log (1-m) \}$$
$$= \exp \left\{ x \log \frac{m}{1-m} + \log (1-m) \right\}$$

$$A = -\log(1-m) = -\log(1-\sigma(\theta)) = \theta + \log(1+e^{-\theta}) \quad \text{exercise}$$

$$\phi(x) = x$$

$$\theta = \log \frac{m}{1-m}$$

$$m = \sigma(\theta) = \frac{1}{1+e^{-\theta}} \quad \psi \text{ mapping mean function and inverse}$$

$$h = 1$$

Alternatively, overcomplete representation:

$$\exp \left\{ \begin{bmatrix} x \\ 1-x \end{bmatrix} \begin{bmatrix} \log m \\ \log (1-m) \end{bmatrix} \right\}$$

Univariate Gaussian

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

$$= \exp \left[-\frac{1}{2\sigma^2} x^2 + \frac{x\mu}{\sigma^2} - \frac{1}{2\sigma^2} \mu^2 \right]$$

$$\phi(x) = \begin{bmatrix} x^2 & x \end{bmatrix}$$

$$\theta = \begin{bmatrix} -\frac{1}{2\sigma^2} & \frac{\mu}{\sigma^2} \end{bmatrix}$$

$$\eta = \sqrt{2\pi} \sigma \exp \frac{\mu^2}{2\sigma^2}$$

$$\mu = \theta_1 \sigma^2$$

$$= -\frac{1}{2\sigma^2} \theta_1$$

$$\sigma^2 = \theta_2 = -\frac{1}{2\theta_2} \quad \psi \text{ mean function}$$

$$A = \log \sqrt{2\pi} + \log \sigma + \frac{\mu^2}{2\sigma^2}$$
$$= \frac{1}{2} \log \sqrt{2\pi} + \frac{1}{2} \log -2\theta_2 + \frac{\theta_1^2}{-4\theta_2}$$

Key properties

- Derivatives of the log-partition function A are the cumulants of the distribution, e.g. $\mathbb{E}[X]$, $\text{var}[X]$, etc.

$$\begin{aligned}\frac{dA}{d\eta} &= \frac{d}{d\eta} \left(\log \int \exp(\eta^T \phi) h(x) dx \right) \\ &= \frac{\int \phi \exp(\eta^T \phi) h(x) dx}{\exp(A(\theta))} = \int \phi(x) p(x) dx = \mathbb{E}[\phi(x)]\end{aligned}$$

Exercise: Show proof for variance and multivariate case

Bernoulli: $A = \theta + \log(1 + e^{-\theta})$

$$\frac{dA}{d\theta} = 1 - \frac{e^{-\theta}}{1 + e^{-\theta}} = \frac{1}{1 + e^{-\theta}} = \sigma(\theta) = \mu$$

Gaussian

$$A = -\frac{1}{2} \log 2\pi - \frac{1}{2} \log(-2\theta_2) - \frac{\theta_1^2}{4\theta_2}$$

$$\frac{dA}{d\theta_1} = -\frac{\theta_1}{2\theta_2} = \frac{\mu/\sigma^2}{1/\sigma^2} = \underline{\mu}$$

Fun

MLE

$$\ell(\theta) = \log p(D|\theta) = \underbrace{\theta^T \phi(D)}_{\text{features of data}} - \underbrace{N A(\theta)}_{\text{data points}}$$

$$\frac{d\ell(\theta)}{d\theta} = \phi(D) - N \frac{d}{d\theta} A(\theta) = \phi(D) - N \mathbb{E}[\phi(D)]$$

Can show concave, so sufficient that

$$\frac{\phi(D)}{N} = \mathbb{E}[\phi(D)]$$

Moment matching.

Recall for Bernoulli: $\frac{1}{N} \sum_{i=1}^N I(x_i=1)$

Gaussian sample mean
sample variance

Generalized Linear Models

Exponential families make it easy to generalize linear regression and classification.

$$p(y|x, w) = h(y) \exp(\underbrace{\phi(y)^T \theta(x, w)}_{\text{conditional}} - A(\theta))$$

↳ break from Murphy

univariate

$$\theta = \underbrace{\psi}_{\text{transform}}(\underbrace{g^{-1}}_{\text{linear}}(w^T x))$$

f - response function

Linear regression

Select exponential family as gaussian

Set ψ and g^{-1} to identity

Estimate w (closed form)

New: Logistic Regression

Select exp. family as Bernoulli

Set ψ to ~~identity~~ _{logit} and $g^{-1}(w^T x) = \text{sigmoid}(w^T x)$

$$\mu = \sigma(w^T x)$$

$$\theta = \log \frac{\mu}{1-\mu}$$

Basic idea

use linear model to estimate μ , ~~then use g^{-1} to match range~~ with squashing g^{-1} to match range. then use ψ to map to natural parameters.

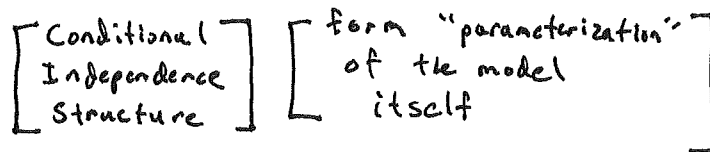
Fitting Models

Exponential Families

6-5

Graphical Models

- Core tool for rest of semester
- Separate out



- Will provide key "modularity" for doing inference

□ High-Level: When does a joint distribution simplify?

- Always use chain-rule $p(A, B, C) = p(A|B, C)p(B|C)p(C)$

But does it factor more? e.g.

$$p(A|B)p(B|C)p(C)$$



Formally, directed GM or Bayes net.

- Graph $G = (V, E)$ with $(s, t) \in E$ $s \neq t$

- $pa(x)$ parents
of x

- Each node corresponds to a random variable.

- Each edge " " to a conditioning decision

- Graph is topologically ordered because of chain rule

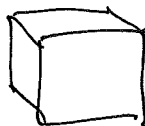
- Nodes that are shaded indicate observed RVs.

Discrete BGMs

- Each node associated with a sample space set

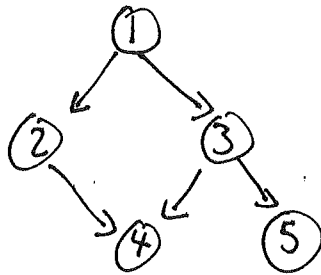
- Local conditional probabilities defined by a CPT

- CPT size of $p(x_i | x_1, \dots, x_{i-1}) = O(2^{|pa(x_i)|})$

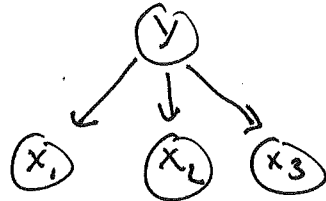


Examples

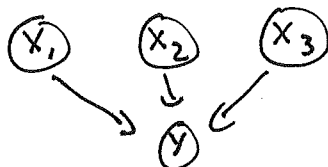
5 CPTs



(?) Write out DGM for naive bayes



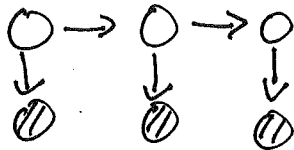
LR



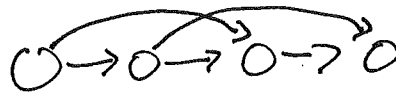
hidden

much worse. simplicity in parameterization

HMM



Auto regressive / Markov chain

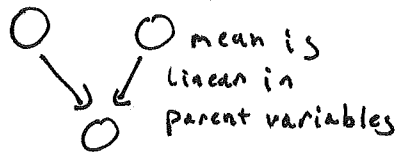


Factorial HMM

□ Gaussian Directed Models

Special Case: much easier

$$p(x_i | \text{pa}(x_i)) = \mathcal{N}(x_i | \mu_i + w_i^T (\text{pa}(x_i) - \mu), \sigma_i^2)$$



$$x_i = \mu_i + \sum_j w_{ij} (x_j - \underbrace{\mu_j}_{\text{constant for simplicity}}) + \sigma_i z_i \quad \forall i, z_i \sim \mathcal{N}(0, 1)$$

Can derive global mean $\underline{\mu}$ as (μ_1, \dots, μ_n)

Let $S = \text{diag}(\underline{\sigma})$ local ~~var~~ standard dev.

$$(x - \underline{\mu}) = W(x - \underline{\mu}) + SZ \quad [\text{matrix-vector form}]$$

$$S_z = (I - W)(x - \underline{\mu})$$

$$x - \underline{\mu} = (I - W)^{-1} S_z \rightarrow u$$

$$\Sigma = \text{cov}[x - \underline{\mu}] = \text{cov}[USZ] = US \text{cov}[Z] S U^T = US^2 U^T$$

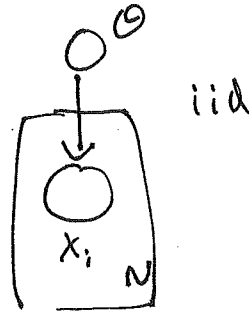
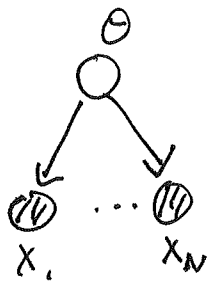
↓
I

(?) Why is $(I - W)$ invertible?

$$\text{GGM} = \mathcal{N}(\underline{\mu}, US^2 U^T)$$

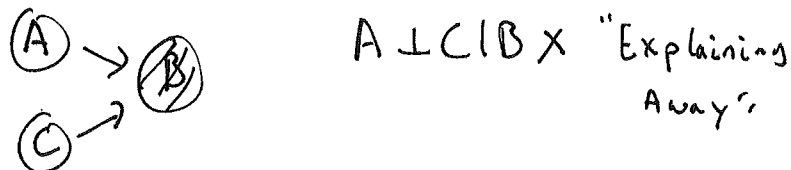
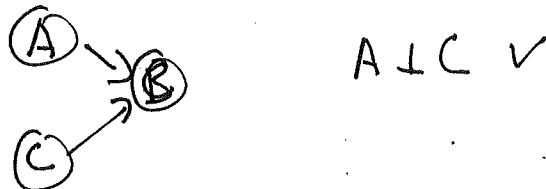
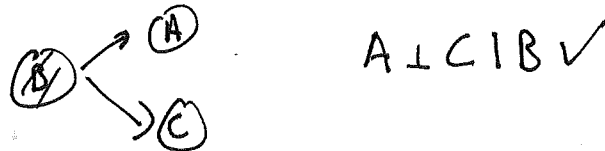
↪ invariant weights
S is local variances

Properties of Bayes Nets: Plates



D-separation and conditional independence

(?)



Undirected Graphical Models 8-1

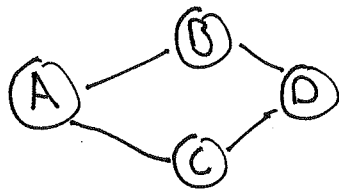
Last class: Directed graphical models. At

- Attempt to describe the conditioning relationships
- can directly use to find local conditional

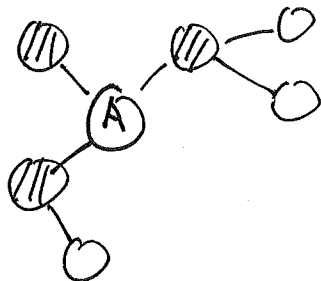
To day: Undirected Graphical models (Markov Random Fields)

- Simpler conditional independence rules
- Describes a different class of distributions
- (personal bias) Often more useful

High-Level: Independence properties

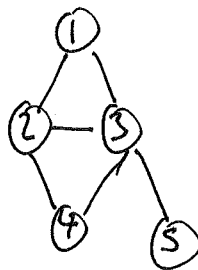
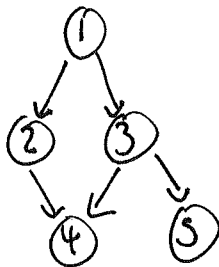


$A \perp D \mid S$ if ~~any~~ no path between A and D that does not cross through S

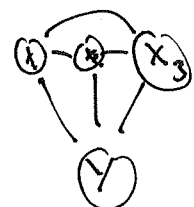
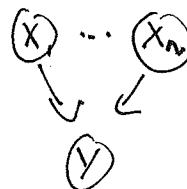
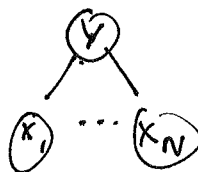
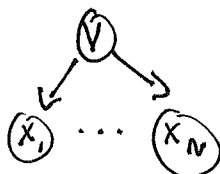


Fundamental property: A is conditionally independent from rest of graph conditioned on its markov blanket (neighbors in G)

Conversion from directed

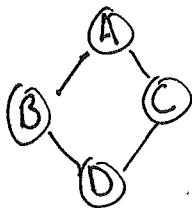
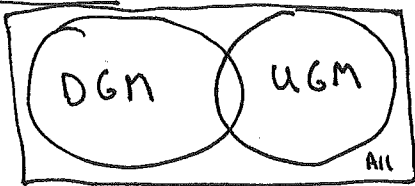


"marry parents"

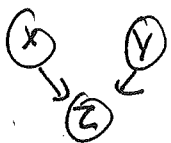
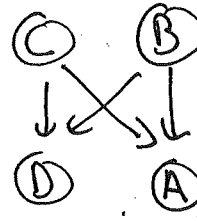
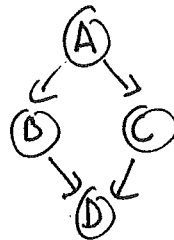


clear why
↓ difficult.

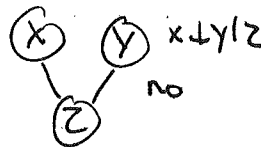
Corner Cases



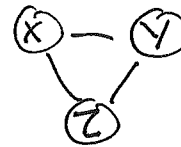
$A \perp D | B, C$
 $B \perp C | A, D$



$X \perp Y | Z$
 $X \perp Y$



$X \perp Y | Z$
no

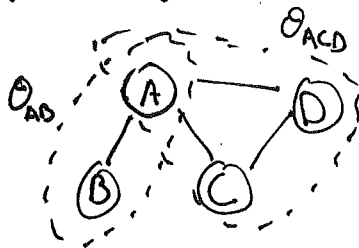


$X \perp Y$

(X)

MRF Parameterization

$$p(x_1, \dots, x_N) = h(x) \exp \{ \theta^T \phi(x) - A(\theta) \} \quad (\text{exponential function family})$$

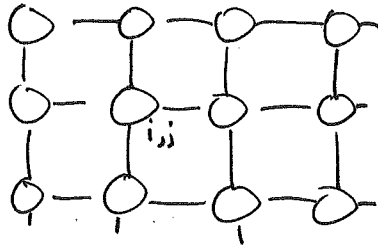


one log-potential θ associated with each clique in the graph

$$\exp \{ \theta_{AB}(x_A, x_B) + \theta_{ACD}(x_A, x_D, x_C) - A(\theta) \}$$

- Note: unlike DGM there are no ^{easy} local probabilities
- $\theta_{AB}(x_A, x_B)$ is a local energy, but is unnormalized, compare to CPT
- To compute $p(x_1, \dots, x_n)$ need log partition function
- ~~For~~ In general computing A is NP-complete sum or integral over all structures.

Canonical Example



$$\exp \left\{ \sum_{ij} \theta_{ij}^{\uparrow} (x_{ij}, x_{i-1,j}) + \theta_{ij}^{\rightarrow} (x_{ij}, x_{i,j+1}) \dots \right\}$$

- E \rightarrow

partition $\log \sum_x e^{-E(x)}$: super-intractable!

What next?

$p(x_1, \dots, x_N)$ - likelihood of data

$p(x_i)$ - marginals

$\arg \max_{x_{1:N}} p(x_{1:N})$ - arg max / MAP



Gaussian

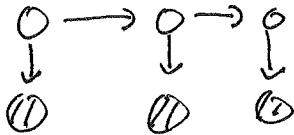
Given "information" form Σ^{-1} can read MRF off of

$\Sigma_{ij}^{-1} \neq 0$ implies $x_i \leftrightarrow x_j$ edge

$$\Sigma^{-1} = (I - W)(S^2)^{-1}(I - W)^T$$

$$\Sigma = U S^2 U^T$$

Recall

HMM

□ Gaussian

$$X_t = X_{t-1} + \epsilon \quad \epsilon \sim N(0, \sigma_\epsilon^2)$$

$$Z_t = X_t + \epsilon' \quad \epsilon' \sim N(0, \sigma_n^2)$$

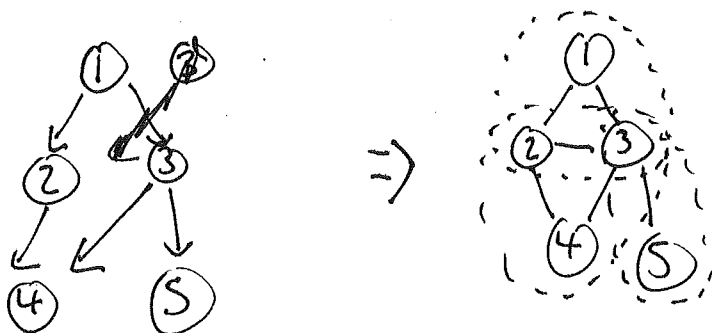
$$\text{Joint: } \prod_t p(x_t | x_{t-1}) p(z_t | x_t)$$

Time series GM

$$8\frac{1}{2}-2$$

This Class:

Exact marginal inference in undirected discrete GMs.



$$p(x_1) p(x_2|x_1) p(x_3|x_1) p(x_4|x_2, x_3) p(x_5|x_3) =$$

$$\exp \{ \theta_{123}(x_1, x_2, x_3) + \theta_{234}(x_2, x_3, x_4) + \theta_{35} - A(\theta) \}$$

where $\theta_{123} = \log p(x_1) + \log p(x_2|x_1) + \log p(x_3|x_1)$

$$\theta_{234} = \log p(x_4|x_2, x_3)$$

$$\theta_{35} = \log p(x_5|x_3) \quad [\text{may lose some CI information}]$$

Hammersley-Clifford: A positive distribution satisfies

CI properties of a graph iff it can be represented as

$$p(y|\theta) = \exp \{ \sum_{C \in \mathcal{C}} \theta_C(y_C) - A(\theta) \} \quad \text{or } \frac{1}{Z(\theta)} \prod_C \psi_C(y_C|\theta_C)$$

Where \mathcal{C} is the set of cliques in the graph. ^{Murphy}

Therefore for simplicity we will consider UGM to start with.

Ex: Conditional random field.

A CRF is a conditional ~~UGM~~ UGM

$$P(y_1 \dots y_n | x) = \exp \left\{ \sum_c \theta_c(y_c | x) - A(\theta) \right\}$$

They are heavily used for labeling style problems.



Linear chain CRF

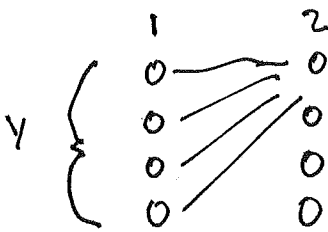
$$P(y_1 \dots y_n | x) = \exp \left\{ \sum_{i=1}^{n-1} \theta_{i,i+1}(y_i, y_{i+1} | x) - A(\theta) \right\}$$

use ψ instead

$$A(\theta) = \log \sum_{y'} \sum_{i=1}^{n-1} \theta_{i,i+1}(y'_i, y'_{i+1}) = \log \sum_{y'_{12}} \exp \{ \theta_{12}(y'_{12}) \} \sum_{y_3} \exp \{ \theta_{23}(y_{23}) \}$$

exponential size set

Dynamic Programming.



$$\alpha_{ij} = \sum_{y_1, \dots, y_{i-1}} \exp \left\{ \sum_{i=1}^{I-1} \theta_{i,i+1}(y_i, y_{i+1}) \right\}$$

$$p(y_t | x) \text{ bel}_t^-(y_t) = \frac{\exp \left\{ \sum_{t-1 \rightarrow t} m_{t-1 \rightarrow t}^-(y_t) - A_t \right\}}{\sum_{y_t} \exp \left\{ \sum_{t-1 \rightarrow t} m_{t-1 \rightarrow t}^-(y_t) - A_t \right\}}$$

$$m_{t-1 \rightarrow t}^-(y_t) = \sum_{y_{t-1}} \theta_{t-1,t}(y_{t-1}, y_t) \text{ bel}_{t-1}^-(y_{t-1})$$

$$m_{t+1 \rightarrow t}^+(y_t) = \sum_{y_{t+1}} \theta_{t,t+1}(y_t, y_{t+1}) \text{ bel}_{t+1}^+(y_{t+1})$$

$$m_{t+1 \rightarrow t}^+(y_t) = \sum_{y_{t+1}} \exp \{ \theta_{t,t+1}(y_t, y_{t+1}) \} m_{t+1 \rightarrow t}^+(y_{t+1})$$

$$\text{bel}_t^-(y_t) = \frac{1}{Z_t} m_{t-1 \rightarrow t}^-(y_t)$$

$$m_{t-1 \rightarrow t}^-(y_t) = \sum_{y_{t-1}} \psi_{t-1}(y_{t-1}, y_t) \text{ bel}_{t-1}^-(y_{t-1})$$

forward

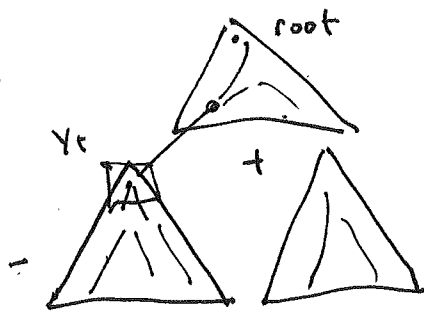
$$\text{bel}_t(y_t) \propto b_t(y_t) m_{t+1 \rightarrow t}^+(y_t) \propto p(y_t | x)$$

$$m_{t+1 \rightarrow t}^+(y_t) = \sum_{y_{t+1}} \psi_{t+1}(y_t, y_{t+1}) m_{t+1 \rightarrow t+1}^+(y_{t+1})$$

backward

General Case: Tree

Inference
9-3



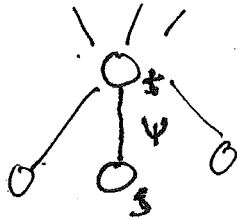
- multiple nodes below

$$bel_t^-(y_t) = \frac{1}{Z} \prod_{c \in ch(t)} m_{c \rightarrow t}^-(y_t)$$

$$m_{c \rightarrow t}^-(y_t) = \sum_{y_c} \psi_{ct}(y_c, y_t) bel_c^-(y_c)$$

$$bel_s(y_s) \propto bel_s^-(y_s) \prod_{c \in pa(s)} m_{c \rightarrow s}^+(y_s)$$

~~$$m_{s \rightarrow t}^+(y_t) = \sum_{y_s} \psi_{st}(y_s, y_t) bel_s^-(y_s)$$~~



$$m_{t \rightarrow s}^+(y_s) = \sum_{y_t} \psi_{st}(y_s, y_t) \frac{bel_t^-(y_t)}{m_{s \rightarrow t}^-(y_t)}$$

around the sides
saved

$$= \sum_{y_t} \psi_{st}(y_s, y_t) \prod_{c \in ch(t)} m_{c \rightarrow t}^-(y_t) \prod_{p \in pa(t)} m_{p \rightarrow t}^+(y_t)$$

Algorithm

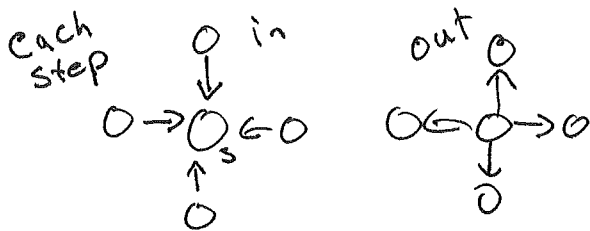
BP implementation

9-3.5
Inference

We assumed bottom-up and top down.

However can implement in parallel.

Assume $bel_s^o(x_s) = \text{unif.}$



$$\text{step 1: } bel_s(x_s) \propto \prod_{t \in \text{nbr}(s)} m_{t \rightarrow s}(x_s)$$

(Compute beliefs)

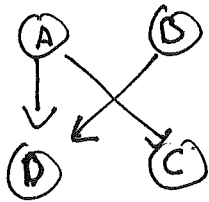
$$\text{step 2: } m_{s \rightarrow t}(x_t) = \sum_{x_s} (\psi_{st}(x_s, x_t) \prod_{u \in \text{nbr}(t)} m_{u \rightarrow s}(x_s))$$

(marginalize self-out)

Variable Elimination

Inference

Consider non-trees



$$\begin{aligned}
 p(D) &= \sum_{A,B,C} p(A,B,C,D) = \sum_{A,B,C} p(D|A) p(C|A) p(A) p(B) \\
 &= \sum \psi(D,A,B) \psi(A,C) \psi(A) \psi(B) \\
 &= \sum_{A,B} \psi(D,A,B) \psi(B) \sum_C \psi(A,C) \psi(A) \\
 &\quad \quad \quad p(A) \\
 &= \sum_A \psi(A) \sum_{B,D} \psi(D,A,B) \psi(B) \\
 &\quad \quad \quad \psi(A) p(A) p(A)
 \end{aligned}$$

Computational complexity, exponential in largest factor

Parallel version of algorithm

Sofar all methods have been exact.

However this is only for a special case of models.

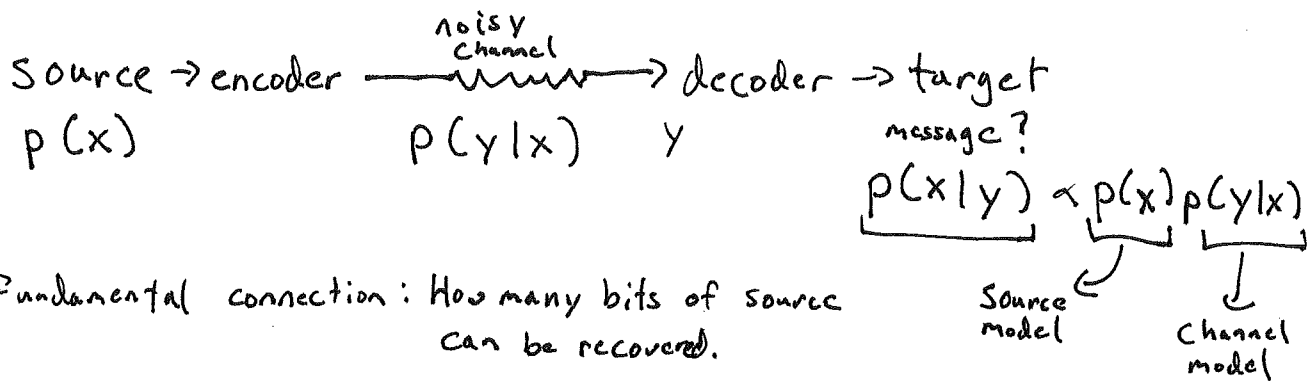
Most models of interest will be approximations, the focus of 2nd part of class.

→ Before we do that, we need some more fundamentals

Weirdly, these will come from information theory.

Whole textbooks written on this connection (Cover and Thomas 2006
MacKay 2003)

Information Theory



Entropy

$$H(p) = H(X) \triangleq - \sum_{k=1}^K p(X=k) \log_2 p(X=k) = \mathbb{E}(\log_2 p(X=k))$$

- measure of the uncertainty of the distribution
- unit of measure is "bits" (or nats if \ln)
- ^{Avg} number of bits required to encode the distribution

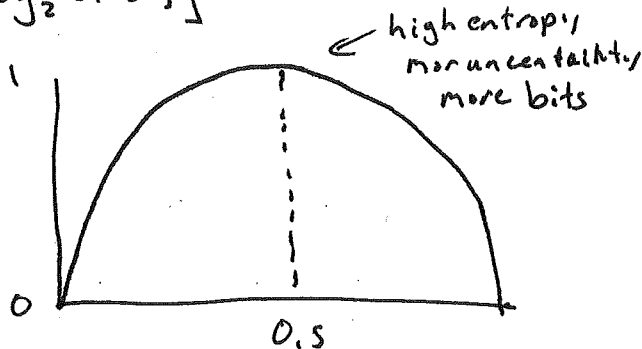
$\exp(H(X)) = \text{perplexity}$ effective uncertainty of distribution

Shannon Game: example guess next word

[

For binary case

$$H(X) = - \left[\underbrace{\theta \log \frac{\theta}{1-\theta}}_{\text{log odds}} + \log_2(1-\theta) \right]$$



Cross Entropy

$$H(p, q) = \sum_k p(X=k) \log q(X=k) = \mathbb{E}_p[\log q(X=k)]$$

p - true distribution

q - our distribution

- number of avg. bits required when true dist is p , but we use q .

Example: language model

- Use q to estimate next word, but true dist is p .

$$\mathbb{E}_p[\log q(X=k)] \approx \sum_{k \in \mathcal{V}} p(X=k) \log q(X=k) = H(p, q)$$

- in fact MLE of categorical classifier

$$\log p(y|x) = \sum_n \log p(y_n | x_n) = \sum_n \underbrace{\sum_{y' \in \mathcal{Y}} \Pi(y' = y_n)}_{\text{Delta dist. over } y's} \underbrace{\log p(y_n | x_n)}_{\substack{q \text{ dist as above} \\ \text{discrete}}}$$

In deep learning, the most common ^{discrete} loss is "Cross-entropy" loss,

KL divergence relative entropy

Information Theory
9.2-3

- The most common way to compare distances.

$$KL(p||q) \triangleq \sum_k p_k \log \frac{p_k}{q_k} = \mathbb{E}_p \left[\log \frac{p_k}{q_k} \right] = \underbrace{-H(p)}_{\text{bits to encode } p} + \underbrace{H(p, q)}_{\text{bits to encode } p \text{ with } q}$$

- extra number of bits needed with q .

$$\operatorname{argmin}_q KL(p||q) = H(p, q) \Rightarrow \text{MLE estimate of } q \text{ with obs } p$$

$$\operatorname{argmin}_p KL(p||q) = \underbrace{-H(p)}_{\text{max entropy } p \text{ that can be encoded w/ } q} + H(p, q) \Rightarrow \text{max entropy } p \text{ that can be encoded w/ } q$$

Theorem: $KL(p||q) \geq 0$ minus (full support)

$$\begin{aligned} -KL(p||q) &= \mathbb{E}_p \left[\log \frac{q_k}{p_k} \right] \leq \log \mathbb{E}_p \left[\frac{q_k}{p_k} \right] = \log \sum q_k(x) \\ &= \log 1 = 0 \end{aligned}$$

Jensen's Inequality: $f(\mathbb{E}[x]) \leq \mathbb{E}(f(x))$ if f is convex

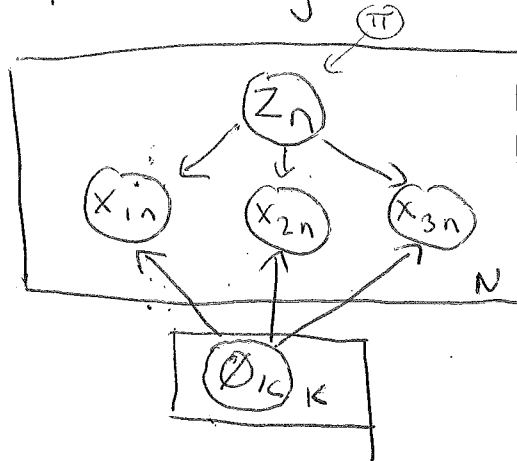
We will use for log

Equality when $p=q$.

□ Until now: Supervised models

input: $x \in \mathbb{R}^d$, output $y \in \mathbb{R}, \{0,1\}, \text{GLM}$

□ Now unsupervised setting: Latent variables z that are unseen.

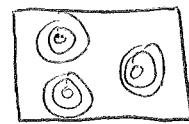


z_n controls which ϕ_k generates the data

Specific instantiations:

1) ϕ are means of gaussians (GMM)

2) ϕ are separate multinomials (mixtures of multinomials)



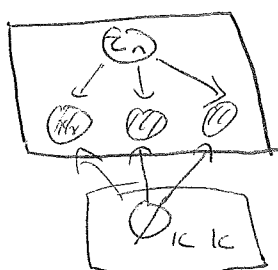
□ Why is this challenging? Graphical models are highly connected.

$$\{x_n\}, \pi, \{z_n\}, \{\phi_k\}$$

$$p(\{x_n\}, \{z_n\}) = \prod_n \prod_k (\pi_{1k} p(x_n | \phi_{1k}))^{z_{n1k}} \quad \text{Complete Data likelihood}$$

$$p(\{x_n\}) = \sum_{\{z_n\}} \prod_n \prod_k \pi_{1k} p(x_n | \phi_{1k})^{z_{n1k}} \quad \text{"}$$

But $p(\{x_n\}, \{z_n\})$ looks okay...



$$z_n \neq \phi_k | x_n$$

$$z_n \neq z_i | x_n \quad \text{!}$$

□ Let's write with logs.

$$\log p(x_n, z_n | \pi, \phi) =$$

$$\sum_n \sum_k z_{nk} \ln \pi_k + z_{nk} \ln p(x_n | \phi_k)$$

□ Expectation-Maximization \rightarrow local coordinate ascent

Goal: Maximize expected complete log-likelihood

1) Suppose we have distribution over z_n .

$$q_{nk} = p(z_n = k)$$

$$\mathbb{E}_q [\ln p(x_n, z_n | \pi, \phi)] = \sum_n \sum_k q_{nk} \ln \pi_k + q_{nk} \ln p(x_n | \phi_k)$$

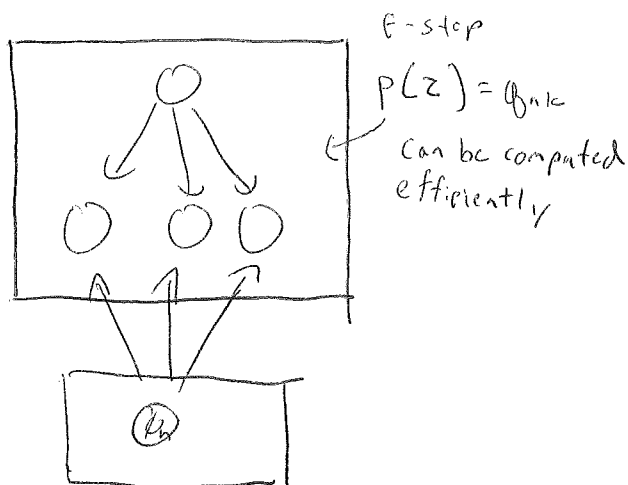
2) E-step: improve approx of q_{nk} given π, ϕ_k estimates

$$q_{nk} \leftarrow \frac{\pi_k p(x_n | \phi_k)}{\sum_k \pi_k p(x_n | \phi_k)} \quad \text{Formula for } p(z_n)$$

3) M-step: improve $\{\phi\}, \pi$ assuming q_{nk} is data.

$$\pi_k \propto \sum_n q_{nk} \quad (\text{categorical})$$

Ⓢ ϕ 's are model specific.



□ Justification of EM

$$\begin{aligned}
 p(\{x_n\}) &= \prod_n \sum_{z_n} p(x_n, z_n | \pi, \phi_k) \\
 \log &= \sum_n \log \sum_k p(x_n, z_n = k | \pi, \phi_k) \\
 &= \sum_n \log \sum_k \underbrace{q(z_n = k) \frac{p(x_n, z_n = k)}{q(z_n = k)}}_{\text{jensen flip log exp}} \\
 &= \sum_n \log \mathbb{E}_{k \sim q} \left[\frac{p(x_n, z_n = k)}{q(z_n = k)} \right] \\
 &\geq \sum_n \mathbb{E}_{k \sim q} \log \dots \\
 &= \sum_n \underbrace{\mathbb{E}_{k \sim q} \log p(x_n, z_n | \dots)}_{\substack{\text{Expected complete data} \\ \text{likelihood}}} + H[q(z_n = k)] \\
 &\quad \text{Always a lower bound we maximize}
 \end{aligned}$$

Now put entropy back

Can show

$$\begin{aligned}
 &= \sum_n \underbrace{KL[q(z) \parallel p(z | x, \pi, \phi)]}_{\substack{\text{alternating opt, here with} \\ \text{marginalization}}} + \text{const} \\
 &\quad q = p(z | x, \pi, \phi)
 \end{aligned}$$

□ EM solves the problem by splitting into two parts

→ local variables z_n

→ global variables ϕ, π

Common strategy

Last week: Exact inference rarely possible

Today: Beginning of unit on approximate inference

▣ Variational Inference: idea

if finding $p(z, \theta | D)$ is too hard find $q(z, \theta)$

that is close: $q^* = \operatorname{argmin}_{q \in \text{Easy}} d(q, p)$

- if q^* is easy, can compute marginals on q^*

- p can be any distribution.

We will use $d(q, p) = \text{KL}(q \| p) = \int q(z, \theta) \log \frac{q(z, \theta)}{p(z, \theta)}$

Rel to EM Recall with EM:

$$\begin{aligned} \log p(x) &= \log \int_{\theta} p(x, \theta) d\theta = \log \int_{\theta} \underbrace{q(\theta)} \frac{p(x, \theta)}{\underbrace{q(\theta)}} d\theta = \log \mathbb{E}_{\theta \sim q} \frac{p(x, \theta)}{q(\theta)} \\ &\geq \mathbb{E}_{\theta \sim q} \log \frac{p(x, \theta)}{q(\theta)} \quad [\text{Jensen's inequality}] \end{aligned}$$

Now consider gap in likelihood

$$\begin{aligned} \log p(x) - \mathbb{E}_q \log \frac{p(x, \theta)}{q(\theta)} &= \cancel{\int_{\theta} q(\theta) \log p(x)} - \cancel{\int_{\theta} q(\theta) \log \frac{p(x, \theta)}{q(\theta)}} \\ &\quad \mathbb{E}_q \log(p(x)) - \log \frac{p(x, \theta)}{q(\theta)} \\ &= \mathbb{E}_q \left[\log \frac{q(\theta)}{p(\theta|x)} \right] = \text{KL}(q, p) \end{aligned}$$

Directly minimizing KL is equivalent to minimizing lower bound coordinate ascent versus opt.

Solving

- This is an optimization problem any method is okay
- Today Mean-field

1) Select $q(z) = \prod_i q_i(z_i)$, utilize to approximate $p(z)$

2) Recall goal is to maximize lower-bound $KL(q, p)$, we do this by fitting each q_i individually.

$$\begin{aligned}
 q_i^*(z_i) &\leftarrow \operatorname{argmin}_{q_i} KL(q \| p) \\
 &= \operatorname{argmin}_{q_i} \int_Z \left(\prod_j q_j(z_j) \right) \log \left(\frac{\prod_j q_j(z_j)}{p(z)} \right) \\
 &= \operatorname{argmin}_{q_i} -H(q_i) - \int_Z \left(\prod_j q_j(z_j) \right) \log(p(z)) + \dots \\
 &= \operatorname{argmin}_{q_i} -H(q_i) - \int_{z_i} q_i(z_i) \underbrace{\int_{z_j \neq z_i} \prod_{j \neq i} q_j(z_j) \log p(z)}_{\mathbb{E}_{j \neq i} [\log p(z)]} \quad \text{(other entropy terms)} \\
 &= \operatorname{argmin}_{q_i} KL(q_i \| \tilde{p}) \quad \text{Call this } \log \tilde{p}(z_i)
 \end{aligned}$$

$$q_i \leftarrow \tilde{p}; \quad q_i(z_i) \propto \exp \left\{ \mathbb{E}_{j \neq i} [\log p(z)] \right\}$$

$KL(p \| q)$ version \mathbb{E}_P does same thing

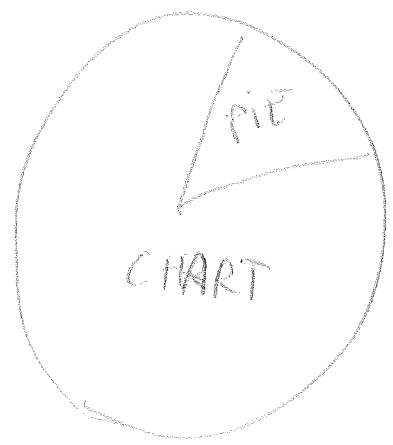
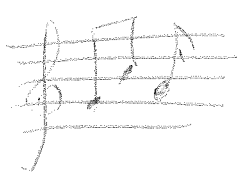


$$0.32 \pm 0.46 \cdot SE$$

$$\sqrt{\sin(\cos(\frac{3}{x}))^5}$$

$$\int_1^2 y^2$$

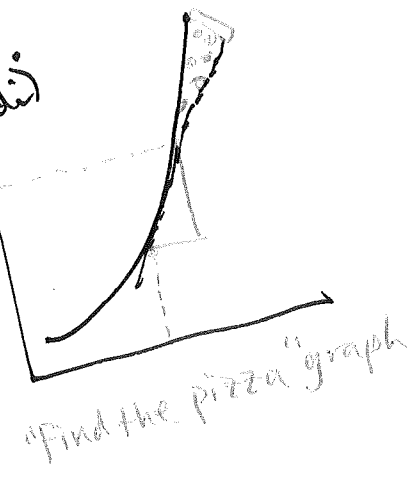
put the entropy back
now take it out again
put it in!
mix it around
do me holey jokey
DESPACITO
hey baby



$$P\left[\frac{1}{x}\right]$$

$$x = \frac{-b \pm \sqrt{4ac^2}}{2a}$$

$$\sqrt{\text{Var}(\hat{S}(t))} = \hat{S}(t)^2 \sum_{i \in \mathcal{I}_t} \frac{d_i}{t_i} \frac{d_i}{d_i}$$



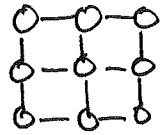
□ Example

Variational

12-3

Ising model

$$p(x) \propto \exp \{ \theta_v^T x + x^T \theta_e^T x \} \quad \mu_i = \mathbb{E}[x_i] = p(x_i=1) + p(x_i=0)(0)$$



updates $\log q_i(x_i) \propto \mathbb{E}_{j \neq i} [\theta_i x_i + \sum_{j \in N(i)} \theta_{ij} x_i x_j] + \text{const}$

$$\ln(\mu_i) = \log q_i(x_i) \propto \theta_i x_i + \mathbb{E}[\sum \theta_{ij} x_i x_j] = \theta_i x_i + \sum \theta_{ij} \mu_j$$

$$\mu_i = \frac{1}{1 + \exp \{ -(\theta_i + \sum \theta_{ij} \mu_j) \}} = \sigma(-(\theta_i + \sum \theta_{ij} \mu_j))$$

VLB $\log Z \geq H(q, p) + H$

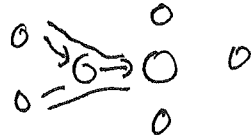
Bayes GMM

Loopy belief propagation

$$1) m_{s \rightarrow t}(x_t) = \sum_{x_s} (\psi_{st}(x_s, x_t) \prod_{u \in \text{nbr}(s), u \neq t} m_{u \rightarrow s}(x_s)) = \sum \exp \left\{ \theta_{st}(x_s, x_t) + \sum m_{s \rightarrow t} \right\}$$

$$2) bc_s(x_s) \propto \prod_{t \in \text{nbr}(s)} m_{t \rightarrow s}(x_s)$$

$q_s =$

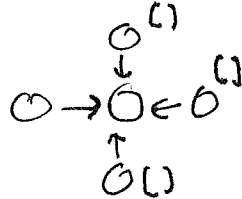
mean field

~~$$q_s(x_s) \propto \prod_{t \in \text{nbr}(s)} \sum_{x_t} \psi_{st}(x_s, x_t) q_t(x_t)$$~~

$$\log q_s(x_s) \propto \mathbb{E} \left[\sum_{j \in \text{nbr}_s} \theta_{sj}(x_s, x_j) \right] + \text{const.}$$

$$\sum_{j \in \text{nbr}_s} \sum_{x_j} q(x_j) \theta(x_s, x_j)$$

$$q_s(x_s) \propto \prod_{j \in \text{nbr}_s} \exp \left(\sum_{x_j} q(x_j) \theta(x_s, x_j) \right)$$



$$q_s(x_s) \propto \exp \left\{ \sum_{j \in N_s} \mathbb{E} \theta(x_s, x_j) \right\}$$

Adv. Variational

13-2

Max product BP

MAP - Relaxations

14-1

Linear programming relaxation

MAP / Relaxations

14-2

MAP/Relaxations

14-3

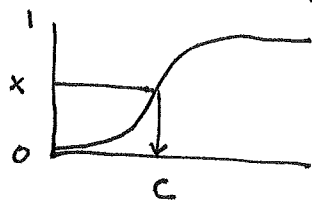
□ How to approx a distribution?

- exact: good
- variational: fuzzy "marginals"
- Monte Carlo: lots of samples, "Hard assignments"

□ Start at the beginning. Drawing samples.

Assume we have $x \sim \text{unif}(0,1)$

If we know cdf $f(x) = p(y \leq x)$ then $f^{-1}(x) = c$ is a sample

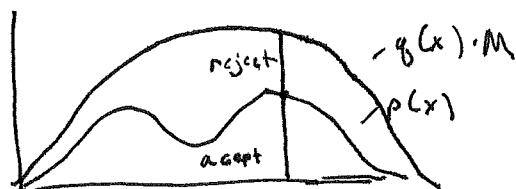


Only really works with univariate, need CDF

□ Gaussian Samples (Box mueller trick)

□ Rejection Sampling

- Assume can evaluate $p(x)$ or $\tilde{p}(x)$ but not sample (no partition)



"throw a dart at Mq , see if it lands in p "

Algorithm:
 $x_n \sim q(x)$
 $u \sim U(0,1)$
 if $u < \frac{p(x_n)}{Mq(x_n)}$

Where $Mq(x) > p(x) \forall x$
 (2) or $\tilde{M}q(x) > \tilde{p}(x)$
 $\tilde{M} = MZ$

□ Proof: $p(x < x_0 | x \text{ accepted}) = \frac{p(x < x_0, x \text{ accepted})}{p(x \text{ accepted})}$

Monte Carlo 15-2

$$= \frac{\int_{-\infty}^{x_0} \int_0^1 q(x) \cdot \frac{\tilde{p}(x)}{M q(x)} dx du}{\int_{-\infty}^{\infty} \int_0^1 q(x) \cdot \frac{\tilde{p}(x)}{M q(x)} dx du} = \frac{\frac{1}{M} \int_{-\infty}^{x_0} \tilde{p}(x) dx}{\frac{1}{M} \int_{-\infty}^{\infty} \tilde{p}(x) dx}$$

$$= \int_{-\infty}^{x_0} p(x) dx \quad \text{and note that } p(x \text{ accepted}) \approx \frac{1}{M}$$

$$\frac{1}{M} \int_{-\infty}^{\infty} \tilde{p}(x) dx = \frac{1}{M} \int_{-\infty}^{\infty} p(x) dx = \frac{1}{M}$$

□ Examples

1) Bayes: let $q(\theta)$ be prior $p(\theta)$ to get samples from posterior $p(\theta|x)$

$$\tilde{p}(\theta|x) = p(D|\theta)p(\theta)$$

$$q(\theta) = p(\theta)$$

$$M = p(D|\hat{\theta}) = \text{MLE}$$

$$\Rightarrow \frac{\tilde{p}(\theta)}{M q(\theta)} = \frac{p(D|\theta)}{p(D|\hat{\theta})} \leftarrow \text{posterior} \in \text{MLE}$$

Retain on high-likelihood. Prior controls sampling.

(?) Why is this okay? MLE ensures $\frac{\tilde{p}(\theta)}{M q(\theta)} \leq 1$

2) Gaussians:

$$p(x) = \mathcal{N}(0, \sigma_p^2 I)$$

$$q(x) = \mathcal{N}(0, \sigma_q^2 I)$$

$$\sigma_q^2 > \sigma_p^2$$



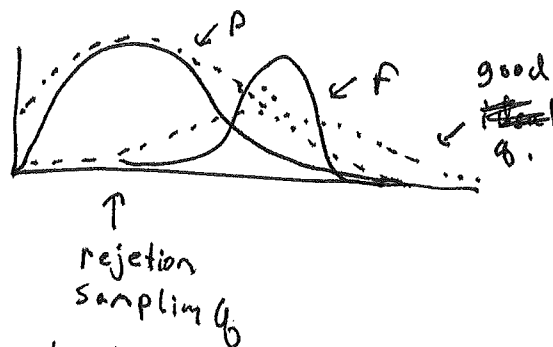
heights $\frac{\left(\frac{1}{\sqrt{2\pi}}\right)^D \left(\frac{1}{\sigma_p}\right)^D \exp\{0\}}{\left(\frac{1}{\sqrt{2\pi}}\right)^D \left(\frac{1}{\sigma_q}\right)^D} = \left(\frac{\sigma_q}{\sigma_p}\right)^D \leftarrow \text{dimensionality!}$

□ Importance Sampling

- In practice, we are usually sampling to compute expectations $\mathbb{E}_p[f(x)] = \int p(x) f(x) dx$

- If we have access to f , it wastes time when $p(x) \uparrow f(x)$

Example:



Recall variational ~~method~~ approach.

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

Here we introduce q into expectation

$$\int q(x) \frac{p(x)}{q(x)} f(x) dx = \mathbb{E}_q \left[f(x) \frac{p(x)}{q(x)} \right] \approx \frac{1}{N} \sum_{\substack{x \sim q \\ \text{samples from } q}} f(x) \underbrace{\frac{p(x)}{q(x)}}_{\text{weight } w(x)}$$

We can show optimal $q^* = \frac{|f(x)| p(x)}{\int |f(x')| p(x') dx}$

Not sure if we need this...

□ What about high-dimensional x_1, \dots, x_N
(Next class)

Monte Carlo
IS-4

Last class: Importance sampling $\int f(x)p(x)dx$

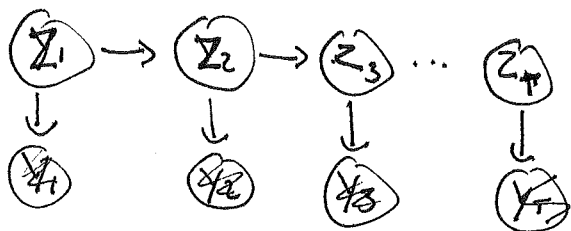
utilize function $q(x)$ to draw samples and then reweight.

Today: X_1, \dots, X_D many variables compute sequentially

First note: $p(x) = \int f(x)p(x)dx$ when $f(x) = \delta_{x'}(x)$

~~can use~~

Return to time series models



Assume we have seen y_1, \dots, y_t , want to estimate $p(z_1, \dots, z_t | y_1, \dots, y_t)$ by sampling

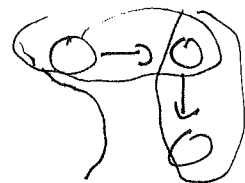
$$p(z_1, \dots, z_t | y_1, \dots, y_t) = \sum_{z'} p(z'_1, \dots, z'_t | y_1, \dots, y_t) \delta_{z'}(z) \approx \sum_s \underbrace{\hat{w}_t^s}_{\text{weight}} \underbrace{\delta_{z^s}(z_{1:t})}_{\substack{\text{sample same} \\ \text{as } z}}$$

Again need a proposal $q(z_{1:t} | y_{1:t})$ as before

$$\hat{w}_t^s = \frac{p(z_{1:t} | y_{1:t})}{q(z_{1:t} | y_{1:t})}$$

Now let's exploit the conditional independence structure

$$\begin{aligned} p(z_{1:t} | y_{1:t}) &\propto p(z_{1:t} | y_{1:t-1}) p(y_t | z_{1:t}, y_{1:t-1}) \\ &\propto p(y_t | z_t) p(z_t | z_{t-1}) p(z_{1:t-1} | y_{t-1}) \end{aligned}$$



$$q(z_{1:t} | y_{1:t}) = q(z_t | z_{t-1}, y_{1:t}) q(z_{1:t-1} | y_{1:t-1}) \quad \text{no indep. needed}$$

$$\hat{w}_t^s = \hat{w}_{t-1}^s \frac{p(y_t | z_t^s) p(z_t^s | z_{t-1}^s)}{q(z_t^s | z_{1:t-1}^s, y_{1:t})}$$

\downarrow or Markov
 $q(z_t^s | z_{t-1}^s, y_{1:t})$

$\sim \begin{matrix} \text{prob} \\ \text{non markov} \end{matrix}$

 \uparrow has evidence

D Algorithm

- sample "particles" $z_t^s \sim q(z_t | z_{t-1}^s, y_t)$
- weight particles as $\tilde{w}_t^s = \frac{p(y_t | z_t^s) p(z_t^s | z_{t-1}^s)}{q(z_t | z_{t-1}^s, y_t)} \tilde{w}_{t-1}^s$
- Normalize to compute "filter" $p(z_t | y_1 \dots y_t) = \sum_s \hat{w}_t^s \delta(z_t)$
particles

□ Issue: Sampling in a very high dimensional space

Use approximation of coverage of distribution.

$$\hat{S}_{\text{eff}} = \frac{1}{\sum_s (\tilde{w}_t^s)^2} \quad \text{i.e. how much of the posterior is captured? many of the samples are being used.}$$

two solutions

1) Resample:

Each time step compute $p(z_t | y_1 \dots y_t)$

If $\hat{S}_{\text{eff}} \leq \text{cutoff}$, then resample from $p(z_t | y_1 \dots y_t)$
and start with $\tilde{w}_t^s = \frac{1}{S}$ (uniform weighting)

2)

$$q(z_t | z_{t-1}^s, y_t) = p(z_t | z_{t-1}^s, y_t)$$

not

$$q(z_t | \dots) = p(z_t | z_{t-1}^s)$$



most kept

0 → 0

most lost

Application: Linear Gaussian

Show updates for

Particle Filtering
16-23

□ Monte-Carlo Principle

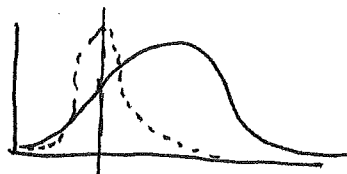
$$\int f(x)p(x)dx \approx \frac{1}{N} \sum_n f(x_n) \quad x_n \sim p(x)$$

unbiased, variance $1/N$

How to get $x_n \sim p(x)$?

- exactly (invert, univariate)
- rejection (perfect samples)
- importance (factor in $f(x)$), sample from q .

Philosophical,



Sample here. This sample is very good.

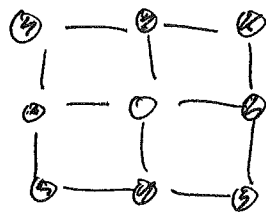
Why forget about it? No memory in rej./importance

□ Markov Chain Monte-Carlo

Tradeoff - correlated samples, lose independence, but more exploration in high prob. regions.

Example: Gibbs Sampling

Idea: assume we have x_1, \dots, x_D , sample x_d



$$p(x_d | x_1, \dots, x_{d-1}, x_{d+1}, \dots, x_D) \propto \exp\{x_d^T \theta \dots\}$$

Sample each variable in turn

Comparison to mean field.

- Similar update process, fix markov blanket
- Mean field: Compute expected value, using expectations of neighbors
- Gibbs: Compute hard assignment by sampling

Gibbs can often be easier to compute.

Markov Chain

MCMC

17-2



• Transition distribution $T(x'/x)$

• Finite example: Transition matrix \mathbb{R}''

- Start with ~~π_0~~ π_0 initial dist.

- Apply transition distribution t times

$$\pi_t = T^t \pi_0$$

- fundamental theorem. Will converge

$$\pi = T\pi$$

- equilibrium point π is "stationary" distribution

- By definition this is eigenvector with $\lambda = 1$

- 2nd largest eigenvalue gives rate.

~~to~~ dominance of rank 1 matrix

• Few more requirements of MCMC

- irreducible/ergodic: if $\pi(x) > 0$ must be able to reach x from x' in finite steps

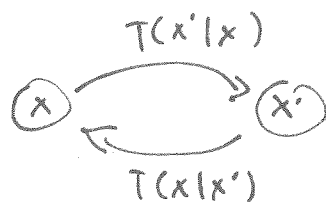
- aperiodic: x only accessible at times

- invariance: generalize from matrix

$$\pi(x') = \int T(x'/x) \pi(x) dx$$

Question: Can we pick T to make $\pi(x) = p(x)$

Detailed Balance : Sufficient Condition



Suppose we don't know
direction / reversible
chain / orthogonal matrix

$$\pi(x) T(x'|x) = \pi(x') T(x|x')$$

↳ implies that

$$\int \pi(x) T(x'|x) dx = \pi(x')$$

□ Properties: $E\left[\frac{1}{N} \sum_n f(x_n)\right] = \frac{1}{N} \sum_n E[f(x_n)] = \frac{1}{N} \sum_n \int \pi(x) f(x) dx = E[f(x)]$
↳ $x \sim \text{MCMC chain}$ (at stationary dist. asymptotically)

$$\text{Variance: } \text{var}\left(\frac{1}{N} \sum_n f(x_n)\right) = \frac{1}{N^2} \text{var}\left(\sum_n f(x_n)\right)$$

$$= \frac{1}{N^2} \left[\underbrace{\sum_n \text{var}(f(x_n))}_{\text{good / as expected}} + 2 \sum_n \sum_{n' > n} \underbrace{\text{cov}(f(x_n), f(x_{n'}))}_{\text{big hit if chain is highly correlated}} \right]$$

good / as expected
 $\frac{1}{N}$

big hit if chain
is highly correlated
(trade off, easy samples
versus hard uncorrelated)

□ Metropolis-Hastings

- Define random walk $q(x'|x)$ "proposal dist"
- Reject if $q(x'|x)$ takes us outside distribution

Formally

1) draw x' from $q(x'|x)$

2) accept w.p. $\min(1, \frac{p(x')}{p(x)} \frac{q(x|x')}{q(x'|x)})$

← $\frac{1}{2}$ cancels

$$\pi \leftarrow \mathbb{E}[f(x)] = \int f(x) \frac{p(x')}{p(x)} dx = \int q(x'|x) \frac{p(x')}{p(x) q(x'|x)} (Book)$$

Proof: reverse

$$x' \neq x$$

$$\begin{aligned} & p(x) q(x'|x) \min\left(1, \frac{p(x')}{p(x)} \frac{q(x|x')}{q(x'|x)}\right) \\ &= \underbrace{\frac{p(x)}{p(x)} \frac{q(x'|x)}{q(x'|x)}}_{\text{cancels}} \underbrace{\min(p(x)q(x'|x), p(x')q(x|x'))}_{\text{symmetric}} \end{aligned}$$

Stationary if $\underline{p(x)}$

Mixing Considerations

(?) What if $x \rightarrow x_{\text{map}}$ w.p. 1 but $x_{\text{map}} \rightarrow \text{unif}$?

Show pictures