# bayesian inference: principles and practice
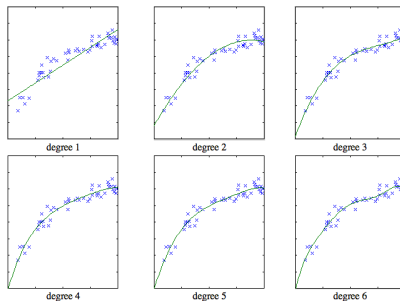
jake hofman

http://jakehofman.com

july 9, 2009

principles
practice

background
bayes' theorem
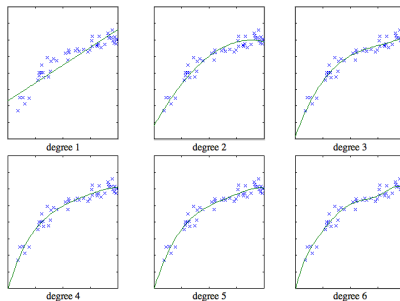bayesian probability
bayesian inference

## motivation

- would like models that:
    - provide predictive and explanatory power
    - are complex enough to describe observed phenomena
    - are simple enough to generalize to future observations

principles
practice

background
bayes' theorem
bayesian probability
bayesian inference

# motivation

- would like models that:
    - provide predictive and explanatory power
    - are complex enough to describe observed phenomena
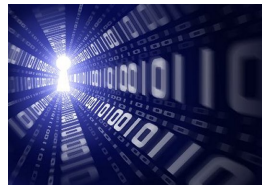    - are simple enough to generalize to future observations



- claim: bayesian inference provides a systematic framework to infer such models from observed data

principles
practice

background
bayes' theorem
bayesian probability
bayesian inference

## motivation

- principles behind bayesian interpretation of probability and bayesian inference are well established (bayes, laplace, etc., 18th century)
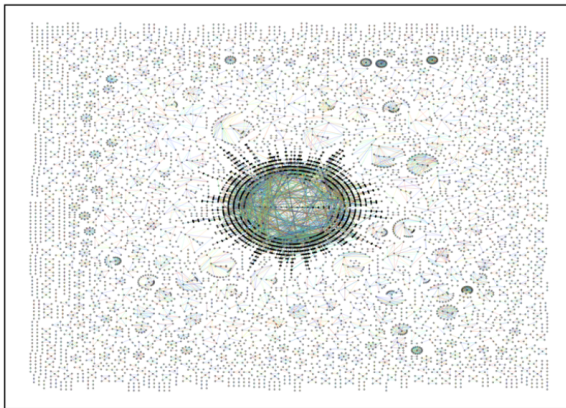


$+$



- recent advances in mathematical techniques and computational resources have enabled successful applications of these principles to real-world problems

principles
practice

background
bayes' theorem
bayesian probability
bayesian inference

# motivation: a bayesian approach to network modularity

principles
practice

background
bayes' theorem
bayesian probability
bayesian inference

## outline

1. **principles (what we'd like to do)**
   - background: joint, marginal, and conditional probabilities
   - bayes' theorem: inverting conditional probabilities
   - bayesian probability: unknowns as random variables
   - bayesian inference: bayesian probability $+$ bayes' theorem

2. practice (what we're able to do)
   - monte carlo methods: representative samples
   - variational methods: bound optimization
   - references

principles
practice

background
bayes' theorem
bayesian probability
bayesian inference

## joint, marginal, and conditional probabilities

### joint distribution

$p_{XY}(X = x, Y = y)$: probability $X = x$ *and* $Y = y$

### conditional distribution

$p_{X|Y}(X = x | Y = y)$: probability $X = x$ *given* $Y = y$

### marginal distribution

$p_X(X)$: probability $X = x$ (*regardless of* $Y$)

principles
practice

background
bayes' theorem
bayesian probability
bayesian inference

## sum and product rules

#### sum rule

sum out settings of irrelevant variables:

$$p(x) = \sum_{y \in \Omega_Y} p(x, y) \tag{1}$$

#### product rule

the joint as the product of the conditional and marginal:

$$\begin{aligned}
p(x, y) &= p(x|y) \, p(y) \tag{2} \\
&= p(y|x) \, p(x) \tag{3}
\end{aligned}$$

principles
practice

background
bayes' theorem
bayesian probability
bayesian inference

# outline

1. ### principles (what we'd like to do)
   - background: joint, marginal, and conditional probabilities
   - **bayes' theorem: inverting conditional probabilities**
   - bayesian probability: unknowns as random variables
   - bayesian inference: bayesian probability $+$ bayes' theorem

2. practice (what we're able to do)
   - monte carlo methods: representative samples
   - variational methods: bound optimization
   - references

principles
practice

background
**bayes' theorem**
bayesian probability
bayesian inference

# inverting conditional probabilities

equate far right- and left-hand sides of product rule

$$p(y|x)\, p(x) = p(x, y) = p(x|y)\, p(y) \tag{4}$$

and divide:

### bayes' theorem (bayes and price 1763)

the probability of $Y$ given $X$ from the probability of $X$ given $Y$:

$$p(y|x) = \frac{p(x|y)\, p(y)}{p(x)} \tag{5}$$

where $p(x) = \sum_{y \in \Omega_Y} p(x|y)\, p(y)$ is the normalization constant

principles
practice

background
bayes' theorem
bayesian probability
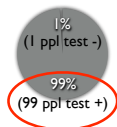bayesian inference

# example: diagnoses a la bayes

- population 10,000
- 1% has (rare) disease
- test is 99% (relatively) effective, i.e.
  - given a patient is sick, 99% test positive
  - given a patient is healthy, 99% test negative

principles
practice

background
**bayes' theorem**
bayesian probability
bayesian inference

# example: diagnoses a la bayes

- population 10,000
- 1% has (rare) disease
- test is 99% (relatively) effective, i.e.
  - given a patient is sick, 99% test positive
  - given a patient is healthy, 99% test negative
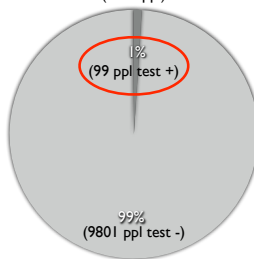- given positive test, what is probability the patient is sick?[1]

---

[1]follows wiggins (2006)

principles
practice

background
bayes' theorem
bayesian probability
bayesian inference

## example: diagnoses a la bayes



sick population
(100 ppl)

1%
(1 ppl test -)

99%
(99 ppl test +)

healthy population
(9900 ppl)

1%
(99 ppl test +)

99%
(9801 ppl test -)

- 99 sick patients test positive, 99 healthy patients test positive

principles
practice

background
bayes' theorem
bayesian probability
bayesian inference

# example: diagnoses a la bayes



sick population
(100 ppl)

1%
(1 ppl test -)

99%
(99 ppl test +)

healthy population
(9900 ppl)

1%
(99 ppl test +)

99%
(9801 ppl test -)

- 99 sick patients test positive, 99 healthy patients test positive
- given positive test, 50% probability that patient is sick

principles
practice

background
**bayes' theorem**
bayesian probability
bayesian inference

## example: diagnoses a la bayes

- know probability of testing positive/negative given sick/healthy
- use bayes' theorem to "invert" to probability of sick/healthy given positive/negative test

$$p\left(sick|test\ +\right) = \frac{\overbrace{p\left(test\ +|sick\right)}^{99/100}\overbrace{p\left(sick\right)}^{1/100}}{\underbrace{p\left(test\ +\right)}_{198/100^2}} = \frac{99}{198} = \frac{1}{2} \quad (6)$$

principles
practice

background
**bayes' theorem**
bayesian probability
bayesian inference

## example: diagnoses a la bayes

- know probability of testing positive/negative given sick/healthy
- use bayes' theorem to "invert" to probability of sick/healthy given positive/negative test

$$p\left(sick|test\,+\right) = \frac{\overbrace{p\left(test\,+|sick\right)}^{99/100}\overbrace{p\left(sick\right)}^{1/100}}{\underbrace{p\left(test\,+\right)}_{198/100^2}} = \frac{99}{198} = \frac{1}{2} \quad (6)$$

- most "work" in calculating denominator (normalization)

principles
practice

background
bayes' theorem
**bayesian probability**
bayesian inference

## outline

1. principles (what we'd like to do)
   - background: joint, marginal, and conditional probabilities
   - bayes' theorem: inverting conditional probabilities
   - bayesian probability: unknowns as random variables
   - bayesian inference: bayesian probability + bayes' theorem

2. practice (what we're able to do)
   - monte carlo methods: representative samples
   - variational methods: bound optimization
   - references

principles
practice

background
bayes' theorem
**bayesian probability**
bayesian inference

# interpretations of probabilities
(just enough philosophy)

- frequentists: limit of relative frequency of events for large number of trials
- bayesians: measure of a state of knowledge, quantifying degrees of belief (jaynes 2003)

principles
practice

background
bayes' theorem
**bayesian probability**
bayesian inference

# interpretations of probabilities
(just enough philosophy)

- frequentists: limit of relative frequency of events for large number of trials
- bayesians: measure of a state of knowledge, quantifying degrees of belief (jaynes 2003)
- key difference: bayesians permit assignment of probabilities to unknown/unobservable hypotheses (frequentists do not)

principles    background
practice    bayes' theorem
**bayesian probability**
bayesian inference

# interpretations of probabilities
(just enough philosophy)

- e.g., inferring model parameters $\Theta$ from observed data $\mathcal{D}$:
  - frequentist approach: calculate parameter setting that maximizes likelihood of data (point estimate),

  $$\widehat{\Theta} = \underset{\Theta}{\operatorname{argmax}}\, p(\mathcal{D}|\Theta) \tag{7}$$

  - bayesian approach: calculate distribution over parameter settings given data,

  $$p(\Theta|\mathcal{D}) = ? \tag{8}$$

principles
practice

background
bayes' theorem
bayesian probability
**bayesian inference**

## outline

1. principles (what we'd like to do)
   - background: joint, marginal, and conditional probabilities
   - bayes' theorem: inverting conditional probabilities
   - bayesian probability: unknowns as random variables
   - bayesian inference: bayesian probability $+$ bayes' theorem

2. practice (what we're able to do)
   - monte carlo methods: representative samples
   - variational methods: bound optimization
   - references

principles
practice

background
bayes' theorem
bayesian probability
bayesian inference

## bayesian probability $+$ bayes' theorem

- bayesian inference:
    - treat unknown quantities as random variables
    - use bayes' theorem to systematically update prior knowledge in the presence of observed data

$$\overbrace{p(\Theta|\mathcal{D})}^{posterior} = \frac{\overbrace{p(\mathcal{D}|\Theta)}^{likelihood}\overbrace{p(\Theta)}^{prior}}{\underbrace{p(\mathcal{D})}_{evidence}} \qquad (9)$$

principles
practice

background
bayes' theorem
bayesian probability
bayesian inference

## example: coin flipping

- observe independent coin flips (bernoulli trials)
- infer distribution over coin bias

principles
practice

background
bayes' theorem
bayesian probability
**bayesian inference**

## example: coin flipping

prior $p(\Theta)$ over coin bias before observing flips



prior over coin fairness
$\alpha_0 = 2, \beta_0 = 2$

principles
practice

background
bayes' theorem
bayesian probability
**bayesian inference**

## example: coin flipping

observe flips: HTHHHTTHHHH

principles
practice

background
bayes' theorem
bayesian probability
**bayesian inference**

## example: coin flipping

update posterior $p(\Theta|\mathcal{D})$ using bayes' theorem

principles
practice

background
bayes' theorem
bayesian probability
**bayesian inference**

## example: coin flipping

observe flips: HHHHHHHHHHHHTHHHHHHHHHH
HHHHHHHHHHHHHHHHHHHHHHH
HHHHHHTHHHHHHHHHHHHHHHH
HHHHHHHHHHHHHHHHHHHHHHH HHHHHHHHT

principles
practice

background
bayes' theorem
bayesian probability
**bayesian inference**

## example: coin flipping

update posterior $p(\Theta|\mathcal{D})$ using bayes' theorem

principles
practice

background
bayes' theorem
bayesian probability
**bayesian inference**

## quantities of interest

- bayesian inference maintains full posterior distributions over unknowns

- many quantities of interest require expectations under these posteriors, e.g. posterior mean and predictive distribution:

$$\bar{\Theta} = \mathbb{E}_{p(\Theta|\mathcal{D})}[\Theta] = \int d\Theta \ \Theta \ p(\Theta|\mathcal{D}) \tag{10}$$

$$p(x|\mathcal{D}) = \mathbb{E}_{p(\Theta|\mathcal{D})}[p(x|\Theta, \mathcal{D})] = \int d\Theta \ p(x|\Theta, \mathcal{D}) \ p(\Theta|\mathcal{D}) \tag{11}$$

principles
practice

background
bayes' theorem
bayesian probability
bayesian inference

## quantities of interest

- bayesian inference maintains full posterior distributions over unknowns

- many quantities of interest require expectations under these posteriors, e.g. posterior mean and predictive distribution:

$$\bar{\Theta} = \mathbb{E}_{p(\Theta|\mathcal{D})}\left[\Theta\right] = \int d\Theta \; \Theta \; p(\Theta|\mathcal{D}) \tag{10}$$

$$p\left(x|\mathcal{D}\right) = \mathbb{E}_{p(\Theta|\mathcal{D})}\left[p\left(x|\Theta, \mathcal{D}\right)\right] = \int d\Theta \; p\left(x|\Theta, \mathcal{D}\right) p(\Theta|\mathcal{D}) \tag{11}$$

- often can't compute posterior (normalization), let alone expectations with respect to it $\rightarrow$ approximation methods

# outline

1. principles (what we'd like to do)
   - background: joint, marginal, and conditional probabilities
   - bayes' theorem: inverting conditional probabilities
   - bayesian probability: unknowns as random variables
   - bayesian inference: bayesian probability + bayes' theorem

2. practice (what we're able to do)
   - monte carlo methods: representative samples
   - variational methods: bound optimization
   - references

## representative samples

- general approach: approximate intractable expectations via sum over representative samples[2]

$$\Phi = \mathbb{E}_{p(x)}[\phi(x)] = \int dx \underbrace{\phi(x)}_{\textit{arbitrary function}} \underbrace{p(x)}_{\textit{target density}} \quad (12)$$

---

[2]follows mackay (2003), including stolen images

## representative samples

- general approach: approximate intractable expectations via sum over representative samples[2]

$$\Phi = \mathbb{E}_{p(x)}\left[\phi(x)\right] = \int dx \underbrace{\phi(x)}_{\text{arbitrary function}} \underbrace{p(x)}_{\text{target density}} \quad (12)$$

$$\Downarrow$$

$$\widehat{\Phi} = \frac{1}{R}\sum_{r=1}^{R}\phi(x^{(r)}) \quad (13)$$

---

[2]follows mackay (2003), including stolen images

## representative samples

- general approach: approximate intractable expectations via sum over representative samples[2]

$$\Phi = \mathbb{E}_{p(x)}\left[\phi(x)\right] = \int dx \underbrace{\phi(x)}_{arbitrary\ function} \underbrace{p(x)}_{target\ density} \qquad (12)$$

$$\Downarrow$$

$$\widehat{\Phi} = \frac{1}{R} \sum_{r=1}^{R} \phi(x^{(r)}) \qquad (13)$$

- shifts problem to finding "good" samples

---

[2]follows mackay (2003), including stolen images

## representative samples

- further complication: in general we can only evaluate the target density to within a multiplicative (normalization) constant, i.e.

$$p(x) = \frac{p^*(x)}{Z} \qquad (14)$$

  and $p^*(x^{(r)})$ can be evaluated with $Z$ unknown

## sampling methods

- monte carlo methods
    - uniform sampling
    - importance sampling
    - rejection sampling
    - . . .

- markov chain monte carlo (mcmc) methods
    - metropolis-hastings
    - gibbs sampling
    - . . .

## uniform sampling

- sample uniformly from state space of all $x$ values
- evaluate non-normalized density $p^*(x^{(r)})$ at each $x^{(r)}$
- approximate normalization constant as

$$Z_R = \sum_{r=1}^{R} p^*(x^{(r)}) \tag{15}$$

- estimate expectation as

$$\widehat{\Phi} = \sum_{r=1}^{R} \phi(x^{(r)}) \frac{p^*(x^{(r)})}{Z_R} \tag{16}$$

## uniform sampling

- sample uniformly from state space of all $x$ values
- evaluate non-normalized density $p^*(x^{(r)})$ at each $x^{(r)}$
- approximate normalization constant as

$$Z_R = \sum_{r=1}^{R} p^*(x^{(r)}) \qquad (15)$$

- estimate expectation as

$$\widehat{\Phi} = \sum_{r=1}^{R} \phi(x^{(r)}) \frac{p^*(x^{(r)})}{Z_R} \qquad (16)$$

- requires prohibitively large number of samples in high dimensions with concentrated density

## importance sampling

- modify uniform sampling by introducing a sampler density $q(x) = \frac{q^*(x)}{Z_Q}$
- choose $q(x)$ simple enough that $q^*(x)$ can be sampled from, with hope that $q^*(x)$ is a reasonable approximation to $p^*(x)$

## importance sampling

- adjust estimator by weighting "importance" of each sample

$$\widehat{\Phi} = \frac{\sum_{r=1}^{R} w_r \phi(x^{(r)})}{\sum_R w_r} \qquad (17)$$

where

$$w_r = \frac{p^*(x^{(r)})}{q^*(x^{(r)})} \qquad (18)$$

# importance sampling

- adjust estimator by weighting "importance" of each sample

$$\widehat{\Phi} = \frac{\sum_{r=1}^{R} w_r \phi(x^{(r)})}{\sum_R w_r} \tag{17}$$

where

$$w_r = \frac{p^*(x^{(r)})}{q^*(x^{(r)})} \tag{18}$$

- difficult to choose "good" $q^*(x)$ as well as estimate reliability of estimator

## rejection sampling

- similar to importance sampling, but proposal density strictly bounds target density, i.e.

$$cq^*(x) > p^*(x), \qquad (19)$$

for some known value $c$ and all $x$



(a)

$cQ^*(x)$

$P^*(x)$

$x$

## rejection sampling

- generate sample $x$ from $q^*(x)$
- generate uniformly random number $u$ from $[0, cq^*(x)]$
- add $x$ to set $\{x^{(r)}\}$ if $p^*(x) > u$
- estimate expectation as $\widehat{\Phi} = \frac{1}{R} \sum_{r=1}^{R} \phi(x^{(r)})$



(b)

$cQ^*(x)$

$P^*(x)$

$u$

$x$      $x$

# rejection sampling

- generate sample $x$ from $q^*(x)$
- generate uniformly random number $u$ from $[0, cq^*(x)]$
- add $x$ to set $\{x^{(r)}\}$ if $p^*(x) > u$
- estimate expectation as $\widehat{\Phi} = \frac{1}{R} \sum_{r=1}^{R} \phi(x^{(r)})$

(b)

- $c$ often prohibitively large for poor choice of $q^*(x)$ or high dimensions

## metropolis-hastings

- use a local proposal density $q(x'; x^{(t)})$, depending on current state $x^{(t)}$



- construct markov chain through state space, converging to target density
- note: proposal density needn't closely approximate target density

## metropolis-hastings

- at time $t$, generate tenative state $x'$ from $q(x'; x^{(t)})$
- evaluate

$$a = \frac{p^*(x^{(t)})}{p^*(x')} \frac{q(x^{(t)}; x')}{q(x'; x^{(t)})} \qquad (20)$$

- if $a \geq 1$, accept the new state; else accept the new state with probability $a$
- if new state is rejected, set $x^{(t+1)} = x^{(t)}$

## metropolis-hastings

- at time $t$, generate tenative state $x'$ from $q(x'; x^{(t)})$
- evaluate

$$a = \frac{p^*(x^{(t)})}{p^*(x')} \frac{q(x^{(t)}; x')}{q(x'; x^{(t)})} \qquad (20)$$

- if $a \geq 1$, accept the new state; else accept the new state with probability $a$
- if new state is rejected, set $x^{(t+1)} = x^{(t)}$
- effective for high dimensional problems, but difficult to assess "convergence" of markov chain[3]

---

[3]see neal (1993)

## gibbs sampling

- metropolis method where proposal density is chosen as conditional distribution, i.e.

$$q(x_i'; x^{(t)}) = p\left(x_i | \{x_j^{(t)}\}_{j \neq i}\right) \qquad (21)$$

- useful when joint density factorizes, as in sparse graphical model[4]

---

[4]see wainwright & jordan 2008

# gibbs sampling

- metropolis method where proposal density is chosen as conditional distribution, i.e.

$$q(x_i'; x^{(t)}) = p\left(x_i | \{x_j^{(t)}\}_{j \neq i}\right) \tag{21}$$

- useful when joint density factorizes, as in sparse graphical model[4]

- similar difficulties to metropolis, but no concerns about adjustable parameters

---

[4]see wainwright & jordan 2008

## outline

1. principles (what we'd like to do)
   - background: joint, marginal, and conditional probabilities
   - bayes' theorem: inverting conditional probabilities
   - bayesian probability: unknowns as random variables
   - bayesian inference: bayesian probability $+$ bayes' theorem

2. practice (what we're able to do)
   - monte carlo methods: representative samples
   - variational methods: bound optimization
   - references

# bound optimization

- general approach: replace integration with optimization
- construct auxiliary function upper-bounded by log-evidence, maximize auxiliary function



---

## variational bayes

- bound log of expected value by expected value of log using jensen's inequality[6]:

$$
\begin{aligned}
-\ln p(\mathcal{D}) &= -\ln \int d\Theta \; p(\mathcal{D}|\Theta)p(\Theta) \\
&= -\ln \int d\Theta \; \frac{p(\mathcal{D}|\Theta)p(\Theta)}{q(\Theta)}q(\Theta) \\
&\leq -\int d\Theta \; \ln\left[\frac{p(\mathcal{D}|\Theta)p(\Theta)}{q(\Theta)}\right]q(\Theta)
\end{aligned}
$$



- for sufficiently simple (i.e. factorized) approximating distribution $q(\Theta)$, right-hand side can be easily evaluated and optimized

---

[6]image from feynman (1972)

## variational bayes

- iterative coordinate ascent algorithm provides controlled analytic approxmations to posterior and evidence
- approximate posterior $q(\Theta)$ minimizes kullback-leibler distance to true posterior
- resulting deterministic algorithm is often fast and scalable

# variational bayes

- iterative coordinate ascent algorithm provides controlled analytic approxmations to posterior and evidence

- approximate posterior $q(\Theta)$ minimizes kullback-leibler distance to true posterior

- resulting deterministic algorithm is often fast and scalable

- complexity of approximation often limited (to, e.g., mean-field theory, assuming weak interaction between unknowns)

- iterative algorithm requires restarts, no guarantees on quality of approximation

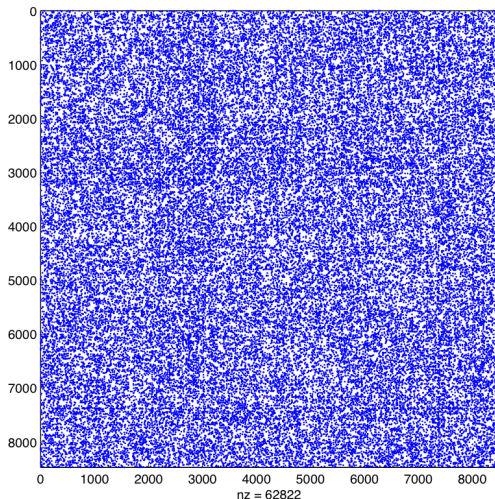# example: a bayesian approach to network modularity

# example: a bayesian approach to network modularity
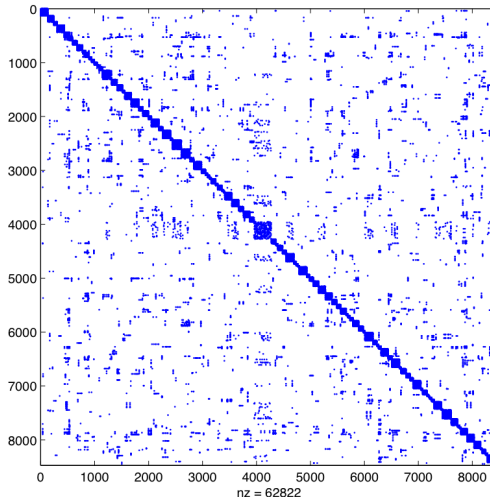
nodes: authors, edges: co-authored papers



can we infer (community) structure in the giant component?

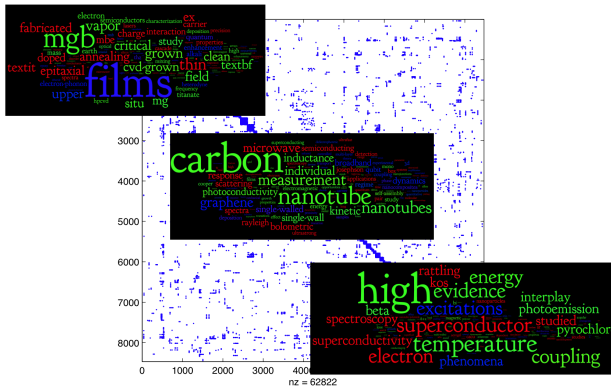# example: a bayesian approach to network modularity

# example: a bayesian approach to network modularity

## example: a bayesian approach to network modularity

inferred topological communities correspond to sub-disciplines

## outline

(1) principles (what we'd like to do)
- background: joint, marginal, and conditional probabilities
- bayes' theorem: inverting conditional probabilities
- bayesian probability: unknowns as random variables
- bayesian inference: bayesian probability + bayes' theorem

(2) practice (what we're able to do)
- monte carlo methods: representative samples
- variational methods: bound optimization
- references

- "information theory, inference, and learning algorithms", mackay (2003)
- "pattern recognition and machine learning", bishop (2006)
- "bayesian data analysis", gelman, et. al. (2003)
- "probabilistic inference using markov chain monte carlo methods", neal (1993)
- "graphical models, exponential families, and variational inference", wainwright & jordan (2006)
- "probability theory: the logic of science", jaynes (2003)
- "what is bayes' theorem ...", wiggins (2006)
- bayesian inference view on cran
- variational-bayes.org
- variational bayesian inference for network modularity