

Lecture 6

Sampling/MCMC/Graph Cuts/Linear Programming

Note Title

1/24/2010

Brief introduction to other techniques for optimization.

First, mention dynamic programming & Warwights TRW.

Linear Programming →

What is the limit of a Bethe/Convex Free energy as $T \rightarrow 0$?

$$\# \min_{\{b_{ij}, \lambda_i\}} \sum_{ij} b_{ij}(x_i, x_j) \psi_{ij}(x_i, x_j) + \sum_i b_i(x_i) \phi_i(x_i),$$

subject to linear constraints

$$\sum_{x_j} b_{ij}(x_i, x_j) = b_i(x_i)$$

$$\sum_{x_i} b_i(x_i) = 1.$$

This is an example of linear programming (LP).
See wikipedia.

$$\text{maximize } \underline{C}^T \underline{x} \quad \text{subject to } A \underline{x} \leq \underline{b}$$

There are courses taught on this at UCLA
eg Prof. Vandenberghe

There are software packages available online
Caution - general purpose packages probably won't exploit the special structure of computer vision problems → there may be faster algorithms

Note: the minimum of # usually does not correspond to the minimum

$$\# \{ \hat{x}_i \} = \arg \min_{\{x_i\}} \left\{ \sum_{ij} \psi_{ij}(x_i, x_j) + \sum_i \phi_i(x_i) \right\}$$

Current research investigates the relationship of BP/TRW for small T and linear programming and the solution of the optimization problem (#).

→ e.g. dual methods like BP/TRW may be more effective than standard linear programming methods for some problems.

(2) Stochastic Sampling & MCMC methods

Also taught at UCLA. Introduction Stat 202C.
More advanced classes. \rightarrow Stat 231b (2nd)

Basic idea -
if we can sample from a distribution $P(x)$
to obtain random independent samples $x_i \sim P(x)$
then we can estimate properties of interest by
$$\int_x f(x) P(x) dx \approx \frac{1}{N} \sum_{i=1}^N f(x_i)$$

How many samples N required?

Key Result \rightarrow the estimate $\frac{1}{N} \sum_{i=1}^N f(x_i)$ is a
random variable, its expectation is $\int_x f(x) P(x) dx$

its variance is independent of the dimension
of the space of x ! Decreases like $1/N$.

Claim \rightarrow big advantage if x lies in a high-dim
space.
compared to alternative ways to approximate $\int_x f(x) P(x) dx$

In particular, we can estimate quantities such as
 $\int_x x P(x) dx$ which can be used for inference

How to get independent samples from a
distribution $P(x)$?

For simple distributions \rightarrow e.g. $P(x) = \frac{1}{\sqrt{2\pi}\sigma} e^{-\frac{(x-\mu)^2}{2\sigma^2}}$
1-D Gaussian

this can be done directly. (Random no. generators)

Many techniques -

- importance sampling

\sim sample from $q(x)$ to get $x_1 \dots x_m$

\sim weight samples by $w_i = \frac{P(x_i)}{q(x_i)}$
and set $\tilde{w}_i = w_i / \sum w_j$

\sim resample from $\{x_1 \dots x_m\}$ by probs \tilde{w}_i

$q(x)$ easy to
sample from.

(3) In general, we cannot sample directly from probability distribution \rightarrow certainly not from MRFs

Instead, need a procedure / algorithm that results in a sample from $P(x)$.

Markov Chain Monte Carlo (MCMC)

Define a transition kernel $K(x|x')$

$$\sum_x K(x|x') = 1, \forall x'$$
$$K(x|x') > 0$$

Condition $\sum_{x'} K(x|x') P(x') = P(x)$

sufficient - detail balance $K(x|x') \frac{P(x')}{P(x)} = K(x'|x) P(x)$ fixed point

require for any states x, y , there exist a trajectory (x_1, \dots, x_N) st. $x_1 = x, x_N = y$
& $K(x_i | x_{i-1}) > 0$.
(you can get to y from x)

MCMC: initialize x_0 at random,

Sample x_1 from $K(x_1|x_0)$
 x_2 " $K(x_2|x_1)$
 \vdots
 x_N " $K(x_N|x_{N-1})$

then x_N will be a random sample from $P(x)$ for suff. large N ?

How quick does this converge? How large N ?

Treat $K(x'|x)$ as a matrix

Calculate its eigenvalues & eigenvectors

MCMC converges like $e^{-N\lambda_2}$

where λ_2 is the second largest (mod) eigenvalue of $K(x'|x)$

(largest eigenvalue is 1)

Convergence is exponentially fast - but usually impossible to calculate λ_2

How do you know you have a sample? Apply set of tests.

(4) How to find a kernel $K(x|x')$?
 Two important kernels that obey detailed balance.

Gibbs sampler:

$$K(x|x') = \sum_r p(r) K_r(x|x')$$

$p(r)$ is a dist. on graph nodes (or groups of nodes)

$$K(x|x') = P(x_r | x'_{N(r)}) \delta_{x'_{N(r)}, x_{N(r)}}$$

$N(r)$ all the nodes except for r .
 Conditionally distributed

Note: conditional distribution is independent of the normalization constant Z of the distribution.

$$P(x) = \frac{1}{Z} e^{-\sum_{ij} \psi_{ij}(x_i, x_j) - \sum_i \phi_i(x_i)}$$

The Z is hard to compute in general and is one reason why "direct sampling" is hard & MCMC required.

$$P(x_i | x'_{N(i)}) = \frac{e^{-\sum_j \psi_{ij}(x_i, x_j) - \phi_i(x_i)}}{\sum_{x_i} e^{-\sum_j \psi_{ij}(x_i, x_j) - \phi_i(x_i)}}$$

usually easy to do.

Example → simple MRF, with $x_i \in \{0, 1\}$

$$\psi_{ij}(x_i, x_j) = \sum_j \psi_{ij} x_i x_j$$

$$\phi_i(x_i) = \sum_i x_i \phi_i$$

Gibbs: choose node i at random, choose x_i by sampling from $P(x_i | x_{-i})$

$$P(x_i | x_{-i}) = \frac{1}{1 + e^{x_i (\sum_j \psi_{ij} x_j + \phi_i)}}$$

Note: Gibbs sampling can be approximated by MFT.

$$b_i^{t+1} = \frac{1}{1 + \exp\{2 \sum_j \psi_{ij} b_j^t + \phi_i\}} \quad \text{(DIA)}$$

also Gibbs sampling can be approximated by Belief Propagation (Koller & Friedman, Jordan, Tenenbaum 2005)

Gibbs is easy to implement

Probably not the fastest MCMC (judged by empirical comparison)

Metropolis Hastings:

$$K(x|x') = q(x|x') \min\left\{1, \frac{p(x) q(x'|x)}{p(x') q(x|x')}\right\}$$

$q(x|x')$ proposal probability for $x \neq x'$

(5)

M-H has two steps:

- (1) sample from proposal $q(x|x')$
- (2) accept or reject proposal with prob $\min\left\{1, \frac{p(x)q(x'|x)}{p(x')q(x|x')}\right\}$ $x \neq x'$

Note: (1) $q(x|x')$ can be anything (provided it enables you to get to all parts of the space)

(2) M-H only depends on the ratio $p(x)/p(x')$ which eliminates the normalization term!

Connection to steepest descent.

$$\rightarrow \text{suppose } q(x'|x) = q(x|x') = 1$$

$$p(x) = \frac{1}{Z} e^{-E(x)}$$

then proposal $x' \rightarrow x$ is accepted with certainty, if $E(x) \leq E(x')$ and with prob. $e^{E(x') - E(x)}$, if $E(x) > E(x')$

Note: MCMC does not converge to a fixed state. Only to a probability distribution.

\rightarrow so prob. that M-H can go 'uphill' and escape local minima.

Note: add temperature (annealing)

$$E(x) \rightarrow E(x)/T$$

can go uphill more easily if T is large
" " with difficulty if T is small.

Note: the convergence rate of an Metropolis-Hastings depends on the proposal probability $q(x|x')$.

\rightarrow how to choose $q(x|x')$? This requires knowledge about the problem domain.

Example in Vision \rightarrow data-driven MCMC.

(Turk, Zhu, Paoli 2002)

Extensions: • Suppose $\{K_\mu(x|x') : \mu \in \Delta\}$ are

transition kernels that obey detailed balance,

then $\sum_{\mu \in \Delta} \alpha_\mu K_\mu(x|x')$, $\alpha_\mu \geq 0$, $\sum \alpha_\mu = 1$ obeys detailed balance.

so, all these kernels can be used.

- kernels can be used to create and destroy nodes
- kernels can be used to group regions
- kernels can do anything

(6)

Max-Flow / Min Cut.

CS Algorithms

Check: Cormen, Leiserson, Rivest.

Max-Flow: Source-to-Sink.

(wiki)

Graph $G = (V, E)$ edges $(u, v) \in E$ have capacity $c(u, v) \geq 0$

Flow

$$f: V \times V \rightarrow \mathbb{R}$$

(1) capacity: $\forall u, v \in V, f(u, v) \leq c(u, v)$

(2) skew-symmetry: $\forall u, v \in V, f(u, v) = -f(v, u)$

(3) conservation: $\forall u \in V \setminus \{s, t\}, \sum_{v \in V} f(u, v) = 0$.

Water flowing down pipes (some pipes allow backward flow).

Total Flow $|f| = \sum_{v \in V} f(s, v)$

Algorithm

eg Ford-Fulkerson enable us to compute the maximal flow.

Theorem: the maximal flow corresponds to the min-cut.

cut (S, T) partition node V into $V = S \cup T$ with $S \cap T = \emptyset$
 $s \in S, t \in T$.

net flow across cut is $f(S, T) = \sum_{u \in S, v \in T} f(u, v)$

Theorem: Certain classes of binary energy minimization problems can be converted into min-cut (max flow) - hence solved in polynomial time.

$$E(x) = \sum_{i, j} \psi_{ij} x_i x_j + \sum_i \phi_i x_i \quad x_i \in \{0, 1\}$$

convert to energy function that depends only on terms for which x_i & x_j take different values

eg. $x_i x_j \rightarrow -x_i(1-x_j) + x_i$

source node $x_s = 0$, sink node $x_t = 1$.

(require condition $\psi_{ij} \leq 0, \forall i, j$)

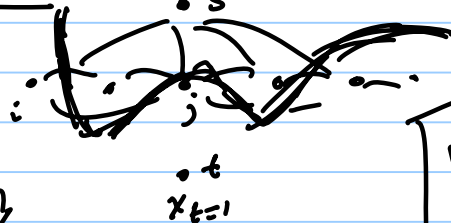
An assignment of variables

$$x_i \rightarrow 0, 1$$

puts the graph nodes into two sets

$$S = \{i : x_i = 0\}$$

$$T = \{i : x_i = 1\}$$



Many variations and extensions!

The energy contributions only occur if at

edges $i - j$ st. $x_i \neq x_j$, or $s - i$ st. $x_i \neq 0$ or $i - t$ st. $x_i \neq 1$

The min cut minimizes the energy