

# Gibbs Sampling

Stephen F. Altschul

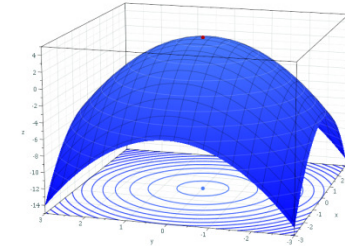
National Center for Biotechnology Information  
National Library of Medicine  
National Institutes of Health

# Optimization in High-Dimensional Space

## Smooth and simple landscapes

Relatively easy to find optimum.

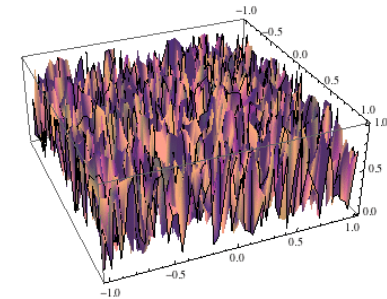
Algorithms: Newton's method; gradient descent.



## Random landscapes

Finding optimal solution intractable.

Algorithms: Brute force enumeration.



## Rough but correlated landscapes

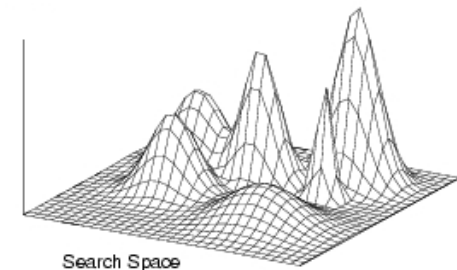
Difficult to find provably optimum solution.

Fairly effective heuristic methods available.

Algorithms: Simulated annealing; Gibbs sampling.

Success depends on details of landscape.

Difficulties: Local optima.



# Local Multiple Alignment

## Simple Version of Problem

Input:  $N$  sequences; pattern width  $W$ .

Problem: Find “highest-scoring”, ungapped local multiple alignment, involving one segment of length  $W$  from each sequence.

Search space:  $L^N$

One may score the local alignment in various ways. Here, we will use *BILD* scores.

### Gibbs sampling:

Given an alignment, one may easily derive a profile or scoring matrix.

Given a profile, one may easily calculate its likelihoods, implied by various segments within a sequence.

Gibbs sampling alternates between generating profiles from given alignments, and sampling alignment positions based on given profile, until “convergence”.

# 1. Initialization

Choose random length- $W$  segments from within the input sequences:

MTQPSKTTKLTKEVDRLISDYQTKQDEQAQETL**VRVYTNLVDMLAKKY**SKGKSFHEDLRQVGMIGLLGAIKRYD  
PVVGKSFEAFIPTIIGEIKRFLRDKTWSVHVPRRIKELGPRIKMAVDQLTTETQ RSPKVEEIAEFLDVSEEEVL  
ETMEMGKSYQALSVDHSIEADSDGSTVTILDIVGSQEDGYERNQQQLMLQSVLHVLSDREKQIIDLTYYIQNKSQK  
ETGDILGISQMHVSRLQRKAVKKLREALIEDPSMELM

MPPLFVMNNEILMHLRALKKTKKDVSLHDPIGQDKEGNEISLIDVLKSENEVDVIDTIQLNMELEKVKQYIDILDD  
REKEVIVGRFGLDLKKEKTQREIAKELGISRSYVSRIEKR**ALMKMFHEFYRAEKE**KRKKAKGK

MELRDLDLNLVVFNQLLVDRRVSITAENLGLTQPAVSNALKRLRTSLQDPLFVRTHQGMEPTPYAAHLAEPVTS  
AMHALRNALQHHESEFDPLTSEFTFTLAM**TDIGEIFYMPRLMDV**LAHQAPNCVISTVRDSSMSLMQALQNGTVDLA  
VGLLPNLQTGFFQRRLLQNHVCLCRKDHVPVTRPLTLERFCSYGHVRVIAAGTGHGEVDTYMTRVGIRRDIRE  
VPHFAAVGHILQRTDLLATVPIRLADCCVEPFGLSALPHPVVLPEIAINMFWHAKYHKDLANIWLRQLMFDLFTD

MNAYTVSRLALDAGVSVHIVRDYLLRGLLRPVACTTGGYGLFDDAALQRLCFVRAAFEAGIGLGALARLCRALDA  
ANCDETAQ**QAVLRQFVERREAL**ANLEVQLAAMPTAPAQHAESLP

•  
•  
•

## 2. Remove one segment from alignment

Select a sequence “X” at random from among the input sequences, and remove its segment from the multiple alignment:

**VRVYTNLVDMLAKKY**

**ALMKMFHEFYRAEKE**

**TDIGEIFYFMPRLMDV**

**LAVLRQFVERRREAL**

**PSPLYPWMRSQFGKC**

**DDTAIRTVLNQALSR**

**QWERG DSEPTGKNLF**

**YHHIKKEKSPKGKSS**

**RIESALLNKIAMLGT**

### 3. Construct a profile from the remaining alignment

Multiple alignment:

VRVYTNLVDMLAKKY  
ALMKMFHEFYRAEKE  
TDIGEIFYFMPRLMDV  
LAVLRQFVERRREAL

⋮



Log-odds scores:

<b>A:</b>	1	1	0	-1	0	
<b>C:</b>	-2	0	-2	-1	-1	
<b>D:</b>	0	2	-3	0	1	...
<b>E:</b>	1	1	-2	0	0	
⋮				⋮		

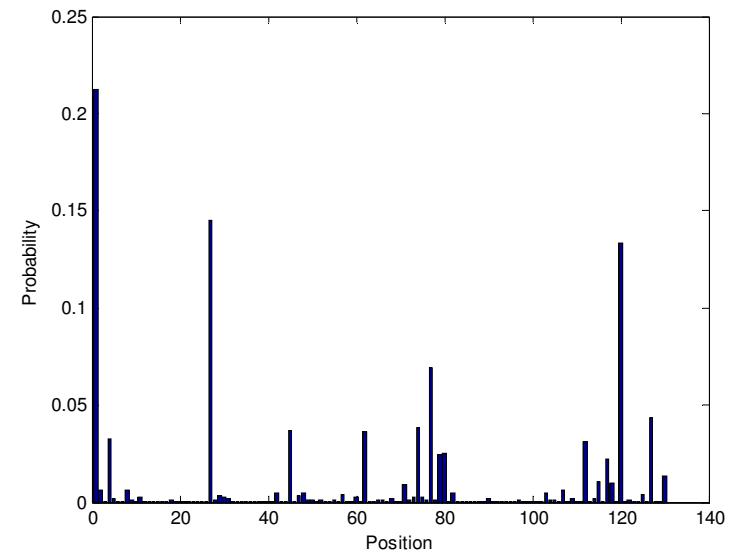
## 4. Calculate relative likelihoods at all positions, and sample

Sequence  $X$ :



Normalized likelihoods:

(  $2^{\text{bit score}}$  )



Sample a random position from sequence  $X$ , weighted by normalized likelihoods.

Add the segment at this position to the multiple local alignment.

If this new alignment is better than any so far seen, remember it.

If there has been no improvement in the last  $I$  iterations, stop.

Otherwise, return to step 2, and remove a new segment.

# Why Does the Algorithm Work?

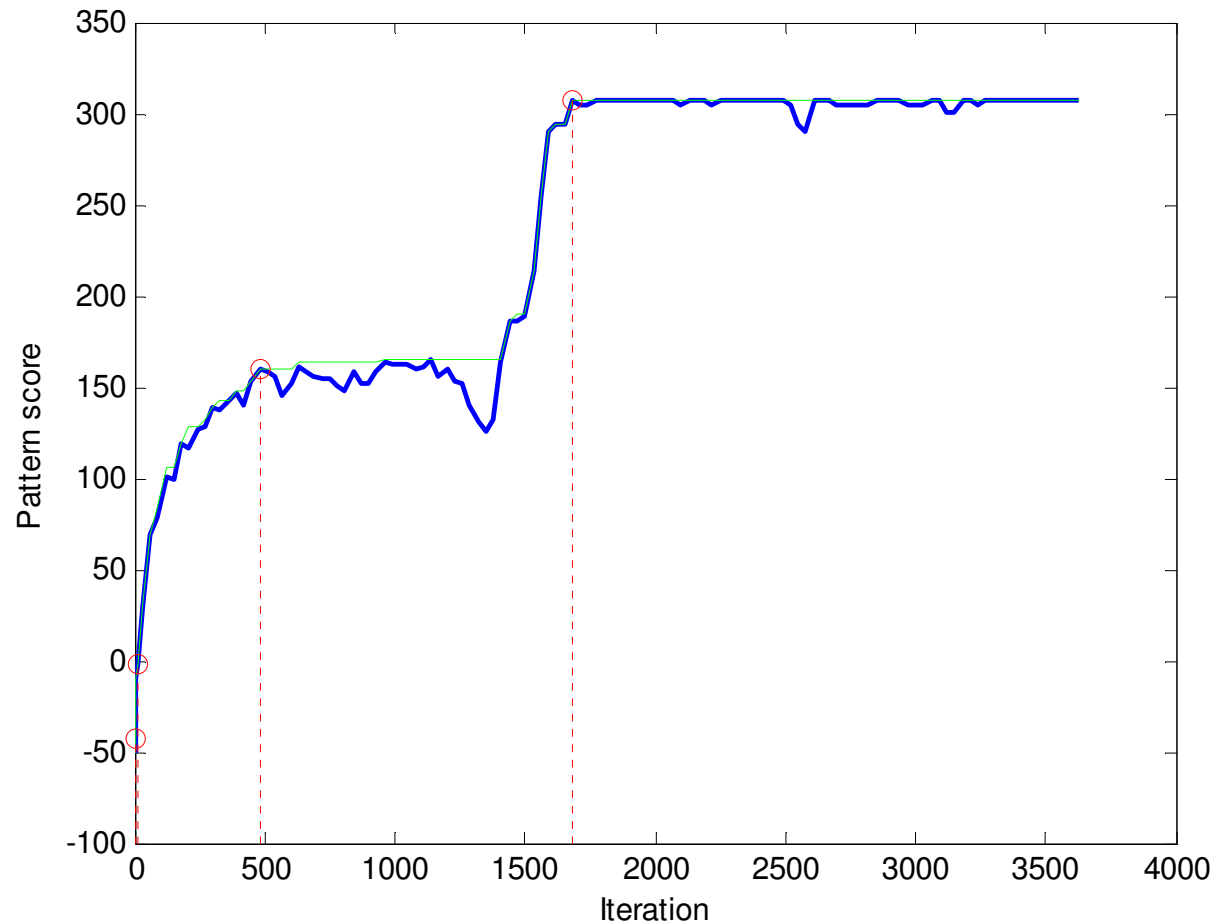
When no common pattern is represented in the multiple alignment, the positions in sequence  $X$  to be sampled have roughly equal likelihoods, so the algorithm performs a random walk through the solution space.

Once a single segment is chosen that is similar to segments found in most or all sequences, these other segments are slightly favored, and a second related segment may well be sampled.

As more related segments are found, the process accelerates, converging on a locally optimal solution. If there are no other good local optima, this solution has a good chance of being the global optimum.



# Behavior of the Objective Function



Input

30 HTH sequences

Pattern width 15

Search space

$\sim 10^{68}$

Time

< 1.0 seconds

# The Evolving Multiple Alignment

MTQPSKTTKLTDEV  
MPPLFVMNNEILMHL  
VVFNQLLVDRRVST  
WFQNRMRKWKKENKT  
SGTGKELVARALHDY  
RIRYRRKNLKHTQRS  
ALDAGVSVHIVRDYL  
QLNGQDVNDLYELVL  
LEIYHHIKKEKSPKG  
SQISRWKRDWIPKFS  
GSVAVLKDEEGKEM  
TINADGSVYAEVVKP  
EIVTAGALKYQENAY  
QLLLRRMEAINESLH  
DLSGKMPNLRQQMMR  
GGLDSYIRAANAWPM  
TRLAWPGNVRQLENT  
ETAATMKDVALKAKV  
PRSASHYLLSDQKSR  
YHNEQKERQAIEQLI  
RLLQLSQGQAVKGNQ  
TRPTEKQYETLENQL  
SNSLKAAPVELRQWL  
AFVKFNCAALPDNLL  
EQLNEREKQIMELRF  
EDKISGTSERPGLK  
TIHQPKDSLGETAFN  
FIGGEDEPGKADIRE  
ARQQEVFDLIRDHIS  
EDEELAEAKKVAHL

1 iteration

SKTTKLTDEVDRLI  
FVMNNEILMHLRALK  
QLLVDRRVSITAENL  
RYLTRRRRIEIAHAL  
KELVARALHDYGRRR  
RRKNLKHTQRSLAKA  
GVSVHIVRDYLLRGL  
QDVNDLYELVLA EVE  
HHIKKEKSPKGKSSI  
RWKRDWIPKFSMLLA  
VLIKDEEGKEMILSY  
QTKTAKDLGVYQSAI  
AGALKYQENAYRQAA  
RRMEAINESLHPPMD  
QDMILLSSKKNAEER  
SYIRAANAWPML SAD  
RLARHFLQIAARELG  
TMKDVALKAKVSTAT  
LVEEKRRRAAKLAATL  
QKERQAIEQLIRHRC  
AMLVANDQMALGAMR  
KNKRALLDALAIEML  
KAAPVELRQWLEEV L  
FNCAALPDNLL ESEL  
EREKQIMELRFGLVG  
SGTSERPGLKKLLR  
PKDSLGETAFNM LLD  
EDEPGKADIREVAF A  
EVFDLIRDHISQTGM  
LAELAKKVAHLLTKE

10 iterations

GISQMHVSRLQRKAV  
GISRSYVSRIEKRAL  
TVRDSSMSLMQALQN  
GVPQQQQQQQQPSQ  
KLDQAALERLKQHRW  
PESQDTQLAEMRAR  
VLRQFVERRREALAN  
PLRDSVKQALKNYFA  
FIMESNLTKVEQHTL  
GVDKSQISRWKRDWI  
RIAQTLLNLAKQPD A  
GVYQSAINKAIHAGR  
GISDAAVSQWKEVIP  
LLEQLLLRRMEAIN E  
NLRQQMMRLMSGEIK  
RVRQLEKNAMKKLRA  
MLPD SWATLLQWAD  
KVSQATRNRVEKAAR  
LLSDQKSRLVEEKRR  
KERQAIEQLIRHRC A  
ALADSLMQLARQVSR  
VLEDQEHQVAK EERE  
YSAAMAEQRHQEWLR  
LSR ATEASKTLQEVL  
GISQSYISRLEKRII  
MERELIVERTKAGLE  
FEPESGYRAMQQILS  
FSSSSGYELAKQMLA  
HISQTGMPPTRAEIA  
GINESQISRWKGDFI

480 iterations

ETGDI LGISQMHVSR  
EIAKELGISRSYVSR  
ITAENLGLTQPAVSN  
EIAHALCLTERQIKI  
RAADLLGLNRNTLRK  
SLAKALKISHVSVSQ  
RAAFEAGICLGALAR  
RAALMGINRGT LRK  
EVAKKCGITPLQVRV  
KTAEAVGVDKSQISR  
EIQQIVGCSRET VGR  
KTAKDLGVYQSAINK  
AVAKALGISDAAVSQ  
SVAQHVCLSPSRLSH  
DIGNYLGLTVETISR  
ELADRYGVSAERVQR  
EAARLLGWGRNTLTR  
DVALKAKVSTATVSR  
DAAALLGVSEMTIRR  
DVARLAGVSVATVSR  
DVAEYAGVSYQTVSR  
KLAQKLGVEQPTLYW  
ELKNELGAGIATITR  
KAARLLGMTPRQVAY  
DVADMMGISQSYISR  
KVAIIYDVGVSTLYK  
DVAKRANVSTTTVSH  
DIAIEAGVSLATVSR  
EIAQRLGFRSPNAAE  
KVADALGINESQISR

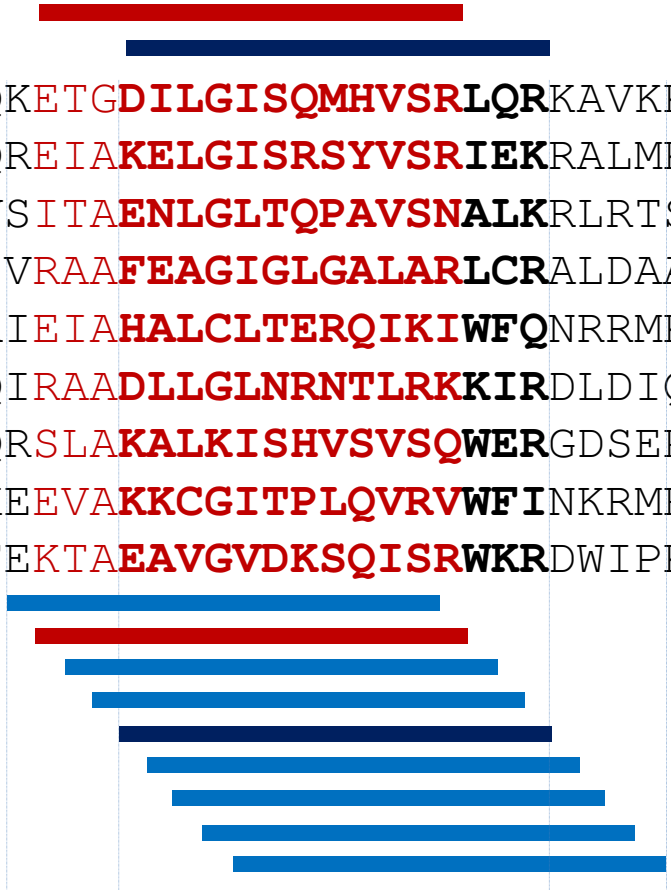
1680 iterations

# Phase Shifts

The Gibbs sampling algorithm may easily converge on a local optimum that is a “phase-shifted” version of the global optimum. **Why?**

Optimal solution:

Solution found:



SQKETG**DILGISQMHVSR**LQ**R**KAVKKL  
TQRE**IAKELGISRSYVSRIE**KRALMKM  
RVS**ITAENLGLTQPAVSNAL**KRLRTSL  
CFV**RAAFEAGIGLGALARLC**RALDAAN  
... RR**IEIAHALCLTERQIKIWFQ**NRRMKW ...  
NQI**RAADLLGLNRNTLRKKIR**DLDIQV  
TQR**SLAKALKISHVSVSQWER**GDSEPT  
EKE**EVAKKCGITPLQVRVWF**INKRMRS  
GTE**KTAEAVGVDKSQISR**WKRDWIPKF

One remedy is to add a separate “phase-shift sampling step”.

No segments are removed, but likelihoods are calculated for the current alignment and several phase-shifted alternatives. These alignments are then sampled among.

This can be understood as changing the topology, of definition of distance, on the underlying “alignment space.”

# Pattern Width

## How does one choose pattern width?

Choosing  $W$  too small discards available information for locating a pattern, while choosing  $W$  too large adds unnecessary noise. The Gibbs sampling algorithm, however, should be fairly robust to deviations that are not too far from the optimal  $W$ .

## What is a reasonable criterion for optimal pattern width?

It can be difficult to compare multiple alignment scores directly for different choices of  $W$ , especially when all column scores are positive. One criterion for selecting  $W$  is the *Minimum Description Length Principle*.

For ungapped local multiple alignments, this is equivalent to optimizing the *BILD* score along a single high-dimensional diagonal, which can be achieved using a variation of the Smith-Waterman algorithm.

Employing the criterion of optimal BILD score,  $W$  may be modified dynamically, within a Gibbs sampling program.

Grunwald, P.D. (2007) *The Minimum Description Length Principle*. MIT Press, Cambridge, MA.

Altschul, S.F., *et al.* (2010) "The construction and use of log-odds substitution scores for multiple sequence alignment." *PLoS Comput. Biol.* **6**:e1000852.

# Close Sequences

If two input sequences are too similar to one another, they can cause each other to “stick” during the sampling stage. In other words, even when they are misaligned, the current position in one sequence will cause the equivalent position in the other sequence to be selected, and vice versa.

## Possible remedies

One may remove extra copies of sequences that are too similar to one another from the input set, and add them back in at a later stage. Paradoxically, this suggests that the most distantly related sequences should be aligned first.

Alternatively, one may employ a strategy analogous to the “realignment stage” in MUSCLE. The relative alignment of a set of closely related sequences can be fixed. Then segments from these sequences can be removed in tandem from the multiple alignment, and new segments (in their previously-fixed relative alignment) sampled in one pass.

# Several Generalizations of the Problem

Some sequences may be missing the pattern.

Some sequences may have multiple copies of the pattern.

The sequences may contain multiple distinct patterns, either consistently ordered or in arbitrary order.

The best alignment between the consensus pattern and its occurrences within the sequences may contain gaps.