

MASARYKOVA UNIVERZITA  
INNOVATION LECTURES (I.N.N.O.I.E.C.) www.muni.cz

**Binding and Kinetics for Experimental Biologists**  
Lecture 2  
**Evolutionary Computing: Initial Estimate Problem**

Petr Kuzmič, Ph.D.  
BioKin, Ltd.  
WATERTOWN, MASSACHUSETTS, U.S.A.

Tento projekt je spolufinancován Evropským sociálním fondem a státním rozpočtem České republiky.

EVROPSKÁ UNIE **esf** MINISTERSTVO ŠKOLSTVÍ, Mládeže a Tělovýchovy OP Vzdělávání pro konkurenceschopnost UNIVERZITA PAVLA PRÁGUE

INVESTICE DO ROZVOJE VZDĚLÁVÁNÍ

BioKin I.L.B.

**Lecture outline**

- The problem:**  
Fitting nonlinear data usually requires an **initial estimate** of model parameters. This initial estimate must be **close enough** to the "true" values.
- The solution:**  
Use a data-fitting method that **does not depend** on initial estimates.
- An implementation:**  
The **Differential Evolution** algorithm (Price *et al.*, 2005).
- An example:**  
Kinetics of **forked DNA** binding to the protein-protein complex formed by DNA-polymerase sliding clamp (gp45) and clamp loader (gp44/62).

BKEB Lec 2: Evolutionary Computing 2

**The ultimate goal of analyzing kinetic / binding data**

SELECT AMONG POSSIBLE **MOLECULAR MECHANISMS**

$E + S \rightleftharpoons E \cdot S \rightarrow E + P$   
 $E + I \rightleftharpoons E \cdot I$

mechanism A  
mechanism B  
mechanism C

A VARIETY OF POSSIBLE MECHANISMS

computer

signal  
concentration

EXPERIMENTAL DATA

Select most plausible model

BKEB Lec 2: Evolutionary Computing 3

**Most models in natural sciences are nonlinear**

LINEAR VS. NONLINEAR MODELS

**Linear**  
 $y = A + k \cdot x$

**Nonlinear**  
 $y = A [1 - \exp(-k \cdot x)]$

$y = 2 + 5x$   
 $y = 2 [1 - \exp(-5x)]$

BKEB Lec 2: Evolutionary Computing 4

**We need initial estimates of model parameters**

NONLINEAR MODELS REQUIRE INITIAL ESTIMATES OF PARAMETERS

$E + S \xrightleftharpoons[k_{-1}]{k_{+1}} E \cdot S \xrightarrow{k_{+2}} E + P$   
 $E + I \xrightleftharpoons[k_{-3}]{k_{+3}} E \cdot I$

A GIVEN MODEL

ESTIMATED PARAMETERS  
 $k_{+1} \quad k_{-1}$   
 $k_{+2} \quad k_{-3}$

signal  
concentration

EXPERIMENTAL DATA

computer

REFINED PARAMETERS  
 $k_{+1} \quad k_{-1}$   
 $k_{+2} \quad k_{-3}$

**The Initial Estimate Problem:**

- Estimated parameters must be "close enough".
- How can we guess them?
- How can we be sure that they are "close enough"?

BKEB Lec 2: Evolutionary Computing 5

**The crux of the problem: Finding *global* minima**

SUM OF SQUARED DEVIATIONS  $\sum(\text{data} - \text{model})^2$   
data - model = "residual"

MODEL PARAMETER

- Least-squares fitting **only** goes "downhill"
- How do we know where to start?**

BKEB Lec 2: Evolutionary Computing 6

### Charles Darwin to the rescue

BIOLOGICAL EVOLUTION IMITATED IN "DE"

Charles Darwin (1809-1882)

**Differential Evolution**  
A Practical Approach to Global Optimization  
Kenneth V. Price - Rainer M. Storn - Jozsef A. Lampinen  
ISBN-10: 3540209506  
Springer

BKEB Lec 2: Evolutionary Computing 7

### Specialized numerical software: DynaFit

CHAPTER TEN

**DYNAFIT—A SOFTWARE PACKAGE FOR ENZYMOLOGY**  
Petr Kuzmíč

DOWNLOAD <http://www.biokin.com/dynafit>

DynaFit implements the **Differential Evolution** algorithm for global sum-of-squares minimization.

Kuzmíc (2009) *Meth. Enzymol.*, **467**, 247-280

BKEB Lec 2: Evolutionary Computing 8

### Biological metaphor: "Gene, allele"

BIOLOGY	COMPUTER
<p><b>gene</b></p> <p>...AAGTCG...GTAACCG...</p> <p>"keratin"</p> <p>four-letter alphabet variable length</p>	<p>• sequence of bits representing a number</p> <p>...01110011000001101101110011...</p> <p>"K<sub>M</sub>" "K<sub>cat</sub>"</p> <p>• two letter alphabet • fixed length (16 or 32 bits)</p>

BKEB Lec 2: Evolutionary Computing 9

### "Chromosome, genotype, phenotype"

BIOLOGY	COMPUTER
<p><b>genotype</b></p> <p>...AAGTCGGT CdGAAGTCGGT TA...</p> <p>keratin oncoprotein</p> <p><b>phenotype</b></p>	<p>• particular combination of all model parameters</p> <p>0110101101 01111001101 001111101101</p> <p>V<sub>max</sub> = 1.23 K<sub>M</sub> = 4.56 K<sub>is</sub> = 78.9</p> <p>full set of parameters</p> $v = \frac{V_{max} [S] / K_M}{1 + [S] / K_M + [S]^2 / K_M K_{is}}$

BKEB Lec 2: Evolutionary Computing 10

### "Organism, fitness"

BIOLOGY	COMPUTER
<p><b>genotype</b></p> <p>...AAGTCGGT CdGAAGTCGGT TA...</p> <p>keratin oncoprotein</p> <p><b>FITNESS:</b> "agreement" with the environment</p>	<p>• <b>FITNESS:</b> agreement between the data and the model</p> <p>V<sub>max</sub> = 1.3 K<sub>M</sub> = 9.1 K<sub>is</sub> = 137.8</p>

BKEB Lec 2: Evolutionary Computing 11

### "Population"

BIOLOGY	COMPUTER
<p>high fitness</p> <p>medium fitness</p> <p>low fitness</p>	<p>V<sub>max</sub> K<sub>M</sub> K<sub>is</sub></p> <p>V<sub>max</sub> K<sub>M</sub> K<sub>is</sub></p>

BKEB Lec 2: Evolutionary Computing 12

### DE Population size in DynaFit

number of population members **per optimized model parameter**

number of population members **per order of magnitude**

```

(DifferentialEvolution)
PopulationSizeFixed      = 0
PopulationSizeMinimal    = 300
PopulationSizeParameter  = 5
PopulationSizePerOrderOfMag = 3
MinimumGenerationsPerParameter = 5
MaximumGenerationsPerParameter = 100
MaximumEvolution        = 4
MinimumEvolution        = 1
RandomSeed               = 1234
  
```

BioKin  
BKEB Lec 2: Evolutionary Computing 13

### "Sexual reproduction, crossover"

BIOLOGY	COMPUTER
	random crossover point
	mother: 01101011001111001101 00011111011
	father: 0110101101001111001101 11100011011
	"sexual mating" probability $p_{cross}$
	child: 011010110110 01111001101 11100011011
	$V_{max}$ $K_M$ $K_S$

BioKin  
BKEB Lec 2: Evolutionary Computing 14

### "Mutation, genetic diversity"

BIOLOGY	COMPUTER
	father: 01101011011 00111100110 11100011011
	$V_{max}$ $K_M$ $K_S$
	mutation
	mutant father: 1100111011 00101101010 11001011001
	$V_{max}^{(*)}$ $K_M^{(*)}$ $K_S^{(*)}$

BioKin  
BKEB Lec 2: Evolutionary Computing 15

### "Mutation, genetic diversity"

THE "DIFFERENTIAL" IN DIFFERENTIAL EVOLUTION ALGORITHM - STEP 1

Compute difference between two randomly chosen "auntie" phenotypes

aunt #1: 01101011011 00111100110 11100011011  
 $V_{max}^{(1)}$   $K_M^{(1)}$   $K_S^{(1)}$

aunt #2: 1100111011 00101101010 11001011001  
 $V_{max}^{(2)}$   $K_M^{(2)}$   $K_S^{(2)}$

subtract

aunt #2 minus aunt #1: 1100111011 00101101010 11001011001  
 $V_{max}^{(2)} - V_{max}^{(1)}$   $K_M^{(2)} - K_M^{(1)}$   $K_S^{(2)} - K_S^{(1)}$

BioKin  
BKEB Lec 2: Evolutionary Computing 16

### "Mutation, genetic diversity"

THE "DIFFERENTIAL" IN DIFFERENTIAL EVOLUTION ALGORITHM - STEP 2

Add **weighted** difference between two "uncle" phenotypes to "father"

father: 01101011011 00111100110 11100011011  
 $V_{max}$   $K_M$   $K_S$

add a fraction of

aunt #2 minus aunt #1: 1100111011 00101101010 11001011001  
 $V_{max}^{(2)} - V_{max}^{(1)}$   $K_M^{(2)} - K_M^{(1)}$   $K_S^{(2)} - K_S^{(1)}$

mutant father: 1100111011 00101101010 11001011001  
 $V_{max}^{(*)}$   $K_M^{(*)}$   $K_S^{(*)}$

BioKin  
BKEB Lec 2: Evolutionary Computing 17

### "Mutation, genetic diversity"

THE "DIFFERENTIAL" IN DIFFERENTIAL EVOLUTION ALGORITHM

EXAMPLE: Michaelis-Menten equation  $v = V_{max} \frac{[S]}{[S] + K_M}$

"father"  $\rightarrow K_M^* = K_M + F \times (K_M^{(1)} - K_M^{(2)})$

"aunt 1"  $\rightarrow K_M^{(1)}$

"aunt 2"  $\rightarrow K_M^{(2)}$

"mutant father"  $\rightarrow$  weight (fraction) mutation rate

BioKin  
BKEB Lec 2: Evolutionary Computing 18

### DE "undocumented" settings in DynaFit

File Edit View Help  
Input Output

```

(DifferentialEvolution) = n
CombineGenerations = 0
ReplaceStragglersPercent = y
Constrained = y
Strategy = 2
Weight = 0.8
Crossover = 1
Jitter = 0.01
Distribution = uniform
AddUserEstimate = n
Scaling = logarithmic
NormalDeviation = 0.5
ExponentialLambda = 2
ReportFrequency = 1
TestParameterRange = y
TestParameterRangeAll = y
TestParameterRangeFull = n
StopParameterRange = 0.01
TestCostFunctionRange = y
StopCostFunctionRange = 0.000001
TestCostFunctionChange = y
StopCostFunctionChange = 0.000001
TestCostFunctionChangeCount = 10
  
```

six different mutation strategies

fractional difference used in mutations

$$K_M^* = K_M + F \times (K_M^{(1)} - K_M^{(2)})$$

probability that "child" inherits "father's" genes, not "mother's" genes

These DE tuning constants are "undocumented" in the DynaFit distribution.

BioKin 1.1.0 BKEB Lec 2: Evolutionary Computing 19

### "Selection"

BIOLOGY	COMPUTER
<p>high fitness</p> <p>more likely to breed</p>	<p>low sum of squares</p> <p>0110101101100111001111011</p> <p><math>V_{max}</math> <math>K_M</math> <math>K_{15}</math></p> <p>more likely to be carried to the next generation</p>
<p>low fitness</p> <p>less likely to breed</p>	<p>high sum of squares</p> <p>00000000001111111111110000000000</p> <p><math>V_{max}</math> <math>K_M</math> <math>K_{15}</math></p> <p>less likely to be carried to the next generation</p>

BioKin 1.1.0 BKEB Lec 2: Evolutionary Computing 20

### Basic Differential Evolution Algorithm - Summary

- 1 Randomly create the initial population (size  $N$ )
- Repeat until almost all population members have very high fitness:
  - 2 Evaluate fitness: sum of squares for all population members
  - 3 Mutation: random gene modification (mutate *father*, weight  $F$ )
  - 4 Sexual reproduction: random crossover with probability  $P_{cross}$
  - 5 Natural selection: keep *child* in gene pool if more fit than *mother*

BioKin 1.1.0 BKEB Lec 2: Evolutionary Computing 21

### Example: DNA + clamp / clamp loader complex

DETERMINE ASSOCIATION AND DISSOCIATION RATE CONSTANT IN AN  $A + B \rightleftharpoons AB$  SYSTEM

Schematic representation of the clamp loading onto forked-DNA substrate. In the forked DNA-substrate, the primer carries a Cy3 fluorescent donor and the gp45 clamp contains an acceptor Cy5 dye.

see Lecture 1 for details

Courtesy of Senthil Perumal, Penn State University (Steven Benkovic lab)

BioKin 1.1.0 BKEB Lec 2: Evolutionary Computing 22

### Example: DynaFit script for Differential Evolution

INSERT A SINGLE LINE IN THE [TASK] SECTION

```

DynaFit: fit_004.txt
File Edit View Help
Input Output
[[task]]
task = fit
data = progress
algorithm = differential-evolution
[mechanism]
DNA + Clamp.Loader <=> Complex : kon koff
[constants]
kon = 1 ?
koff = 1 ?
[responses]
Complex = 1 ? (0.01 ... 100)
[data]
file ./courses/bksh/lab1/stb/data/d1-edit.txt
offset 0.3 ? (0.2 ... 0.4)
  
```

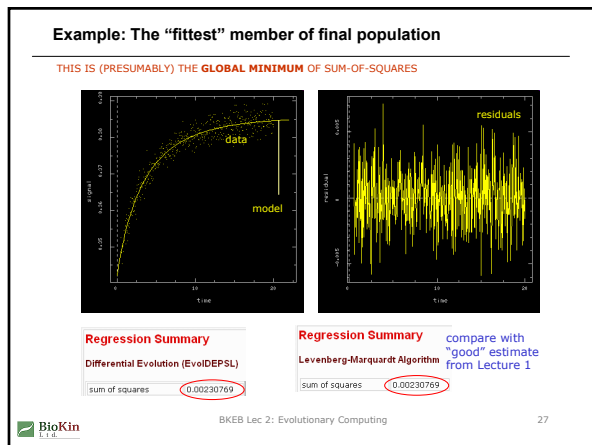
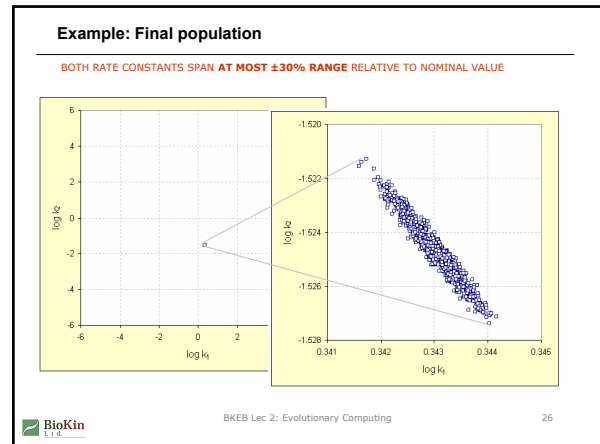
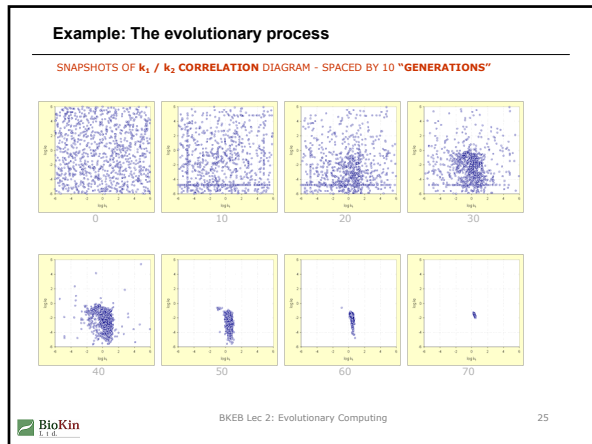
constraints !

BioKin 1.1.0 BKEB Lec 2: Evolutionary Computing 23

### Example: Initial population

BOTH RATE CONSTANTS SPAN TWELVE ORDERS OF MAGNITUDE

BioKin 1.1.0 BKEB Lec 2: Evolutionary Computing 24



### Example: Comparison of DE and regular data fitting

DIFFERENTIAL EVOLUTION (DE) FOUND THE SAME FIT AS THE "GOOD" ESTIMATE

	initial estimate	sum of squares	relative sum of sq.	"best-fit" constants
lecture 1	"good" $k_1 = 1$ $k_2 = 1$	0.002308	1.00	$k_1 = 2.2 \pm 0.5$ $k_2 = 0.030 \pm 0.015$
	"bad" $k_1 = 100$ $k_2 = 0.01$	0.002354	1.02	$k_1 = 0.2 \pm 3.4$ $k_2 = 0.2 \pm 0.6$
1000 random estimates	$k_1 = 10^{-6} - 10^{+6}$ $k_2 = 10^{-6} - 10^{+6}$	0.002308	1.00	$k_1 = 2.2 \pm 0.5$ $k_2 = 0.030 \pm 0.015$

BioKin  
BKEB Lec 2: Evolutionary Computing 28

### Significant disadvantage of DE: very slow

DYNAFIT CAN TAKE MULTIPLE DAYS TO RUN A COMPLEX PROBLEM

DynaFit 4.065 on DNA / clamp / clamp loader example:

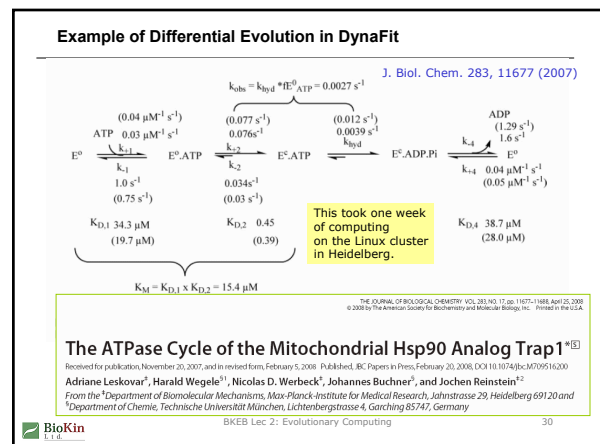
algorithm	computation time	relative time
Levenberg-Marquardt with two restarts	0.88 sec	1
Differential Evolution with four restarts (population size: 1000)	12 min 31 sec	853

1 second  
1 minute  
10 minutes

→

15 minutes  
15 hours  
6 days

BioKin  
BKEB Lec 2: Evolutionary Computing 29



### Example: Systematic scan of many initial estimates

CAREFUL! THIS IS FASTER THAN DIFFERENTIAL EVOLUTION BUT DOES NOT ALWAYS WORK

```

DynaFit - fit-007.txt
File Edit View Help
Input Output

[task]
task = estimate
data = progress

[mechanism]
DNA + Clamp.Loader <=> Complex : kon koff

[constants]
kon = [0.001, 0.01, 0.1, 1, 10, 100, 1000] ?
koff = [0.001, 0.01, 0.1, 1, 10, 100, 1000] ?

[concentrations]
DNA = 0.1
Clamp.Loader = 0.1
    
```

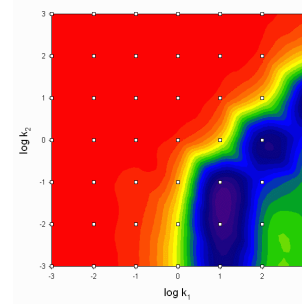
#### ALGORITHM

1. generate all possible combinations of rate constants
2. compute initial sum of squares for each combination
3. rank combinations by initial sum of squares
4. select the best **N** combinations
5. perform a full fit for those **N**
6. rank results again

7 × 7 = 49 combinations of  $k_{on}$  and  $k_{off}$

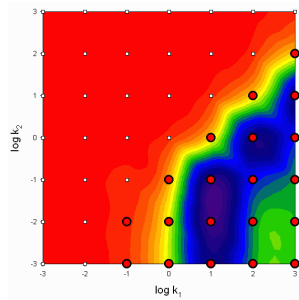
### Example: Systematic scan – Phase 1

AFTER EVALUATING THE INITIAL SUM OF SQUARES FOR ALL 49 COMBINATIONS OF  $k_1$  and  $k_2$



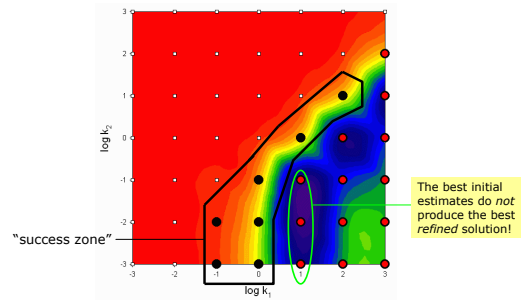
### Example: Systematic scan – Phase 2

AFTER RANKING THE INITIAL ESTIMATES AND SELECTING 20 BEST ONES BY SUM OF SQUARES



### Example: Systematic scan – Phase 3

AFTER PERFORMING FULL REFINEMENT FOR 20 BEST ESTIMATES OUT OF 49 TRIED



### Summary and conclusions

1. Finding good-enough initial estimates is a very difficult problem.
2. One should use system-specific information as much as possible. This includes using the literature and/or general principles for "intelligent" guesses.
3. Always use the "Try" method in DynaFit to display the initial fit. Make sure that the initial estimate is at least approximately correct.
4. The Differential Evolution algorithm almost always helps. However, it can be excruciatingly slow (running typically for multiple hours).
5. The systematic scan (task = estimate) sometimes helps. However, the "best" initial estimates almost never produce the desired solution!
6. DynaFit is not a "silver bullet": You must still use your brain a lot.