# Some Theory Behind Algorithms for Stochastic Optimization

Zelda Zabinsky

University of Washington
Industrial and Systems Engineering

May 24, 2010
NSF Workshop on Simulation Optimization

# Overview

o Problem formulation

o Theoretical performance of stochastic adaptive search methods

o Algorithms based on Hit-and-Run to approximate theoretical performance

o Incorporate random sampling and noisy objective functions

# What is Stochastic Optimization?

o Randomness in algorithm AND/OR in function evaluation

o Related terms:
  - Simulation optimization
  - Optimization via simulation
  - Random search methods
  - Stochastic approximation
  - Stochastic programming
  - Design of experiments
  - Response surface optimization

# Problem Formulation

o Minimize $f(x)$ subject to $x$ in $S$

o $x$:   $n$ variables, continuous and/or discrete

o $f(x)$:   objective function, could be black-box, ill-structured, noisy

o S:  feasible set, nonlinear constraints, or membership oracle

o Assume an optimum $x*$ exists, with $y*=f(x*)$

# Example Problem Formulations

o Maximize expected value
  subject to standard deviation < b

o Minimize standard deviation
  subject to expected value > t

o Minimize CVaR (conditional value at risk)

o Minimize sum of least squares from data
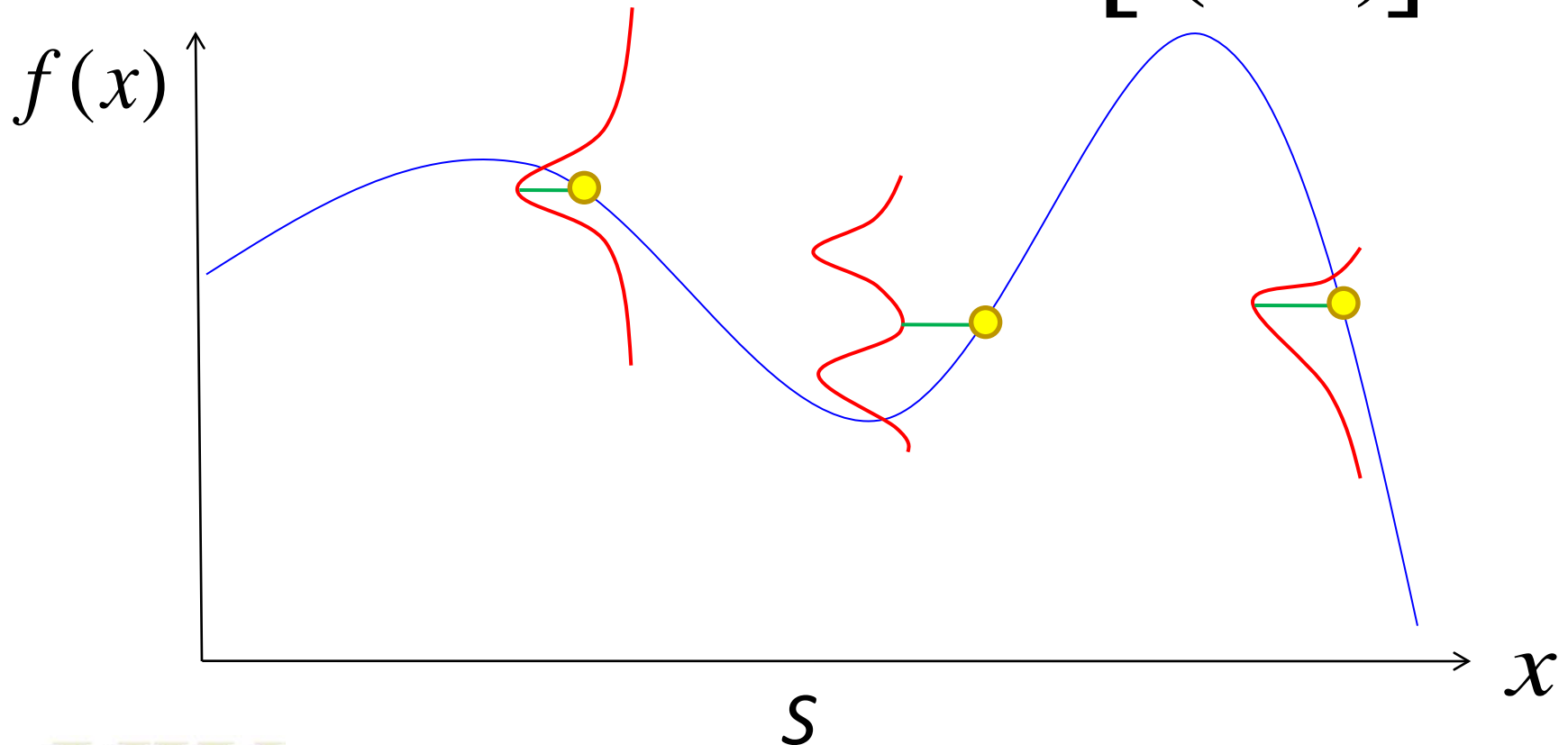
o Maximize probability of satisfying noisy
  constraints

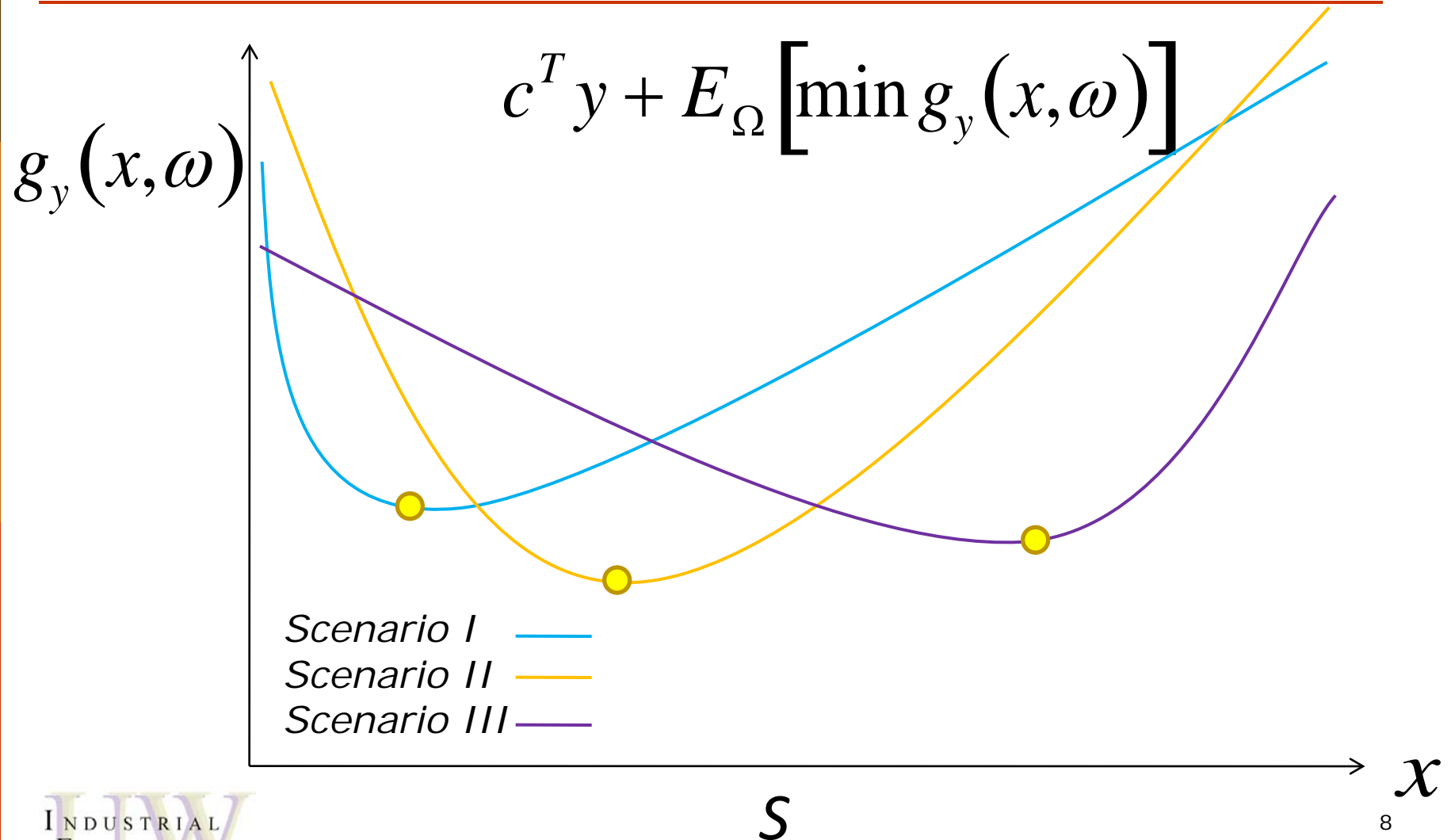INDUSTRIAL
ENGINEERING

# Approximate or Estimate *f(x)* ?

o Approximate a complicated function:

- Taylor series expansion
- Finite element analysis
- Computational fluid dynamics

o Estimate a noisy function with:

- Replications
- Length of discrete-event simulation run

# Noisy Objective Function

$$f(x) = E_\Omega\big[g(x, \omega_x)\big]$$

$f(x)$

$s$

$x$

# Scenario-based Recourse Function

$$c^T y + E_{\Omega}\left[\min g_y(x,\omega)\right]$$

$g_y(x,\omega)$



Scenario I
Scenario II
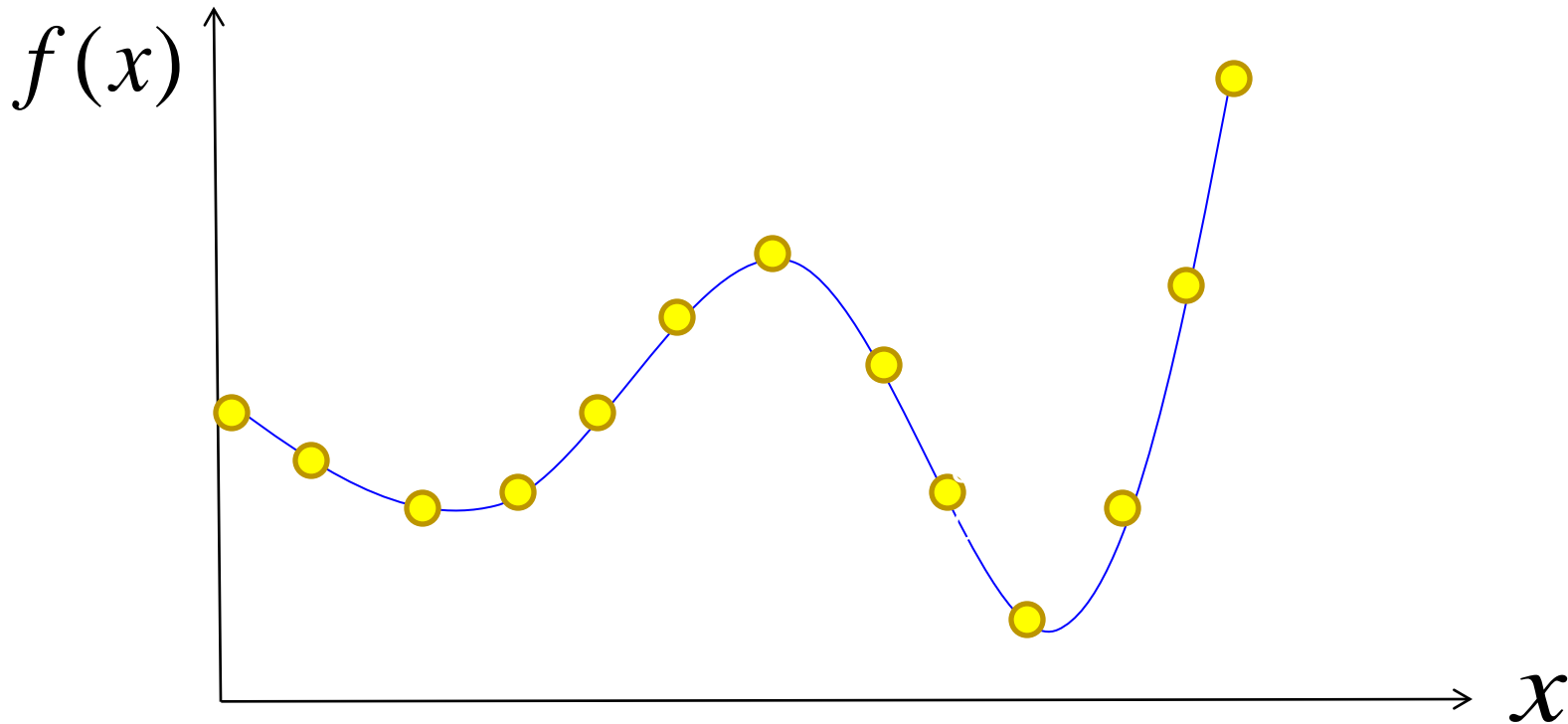Scenario III

$x$

$S$

INDUSTRIAL
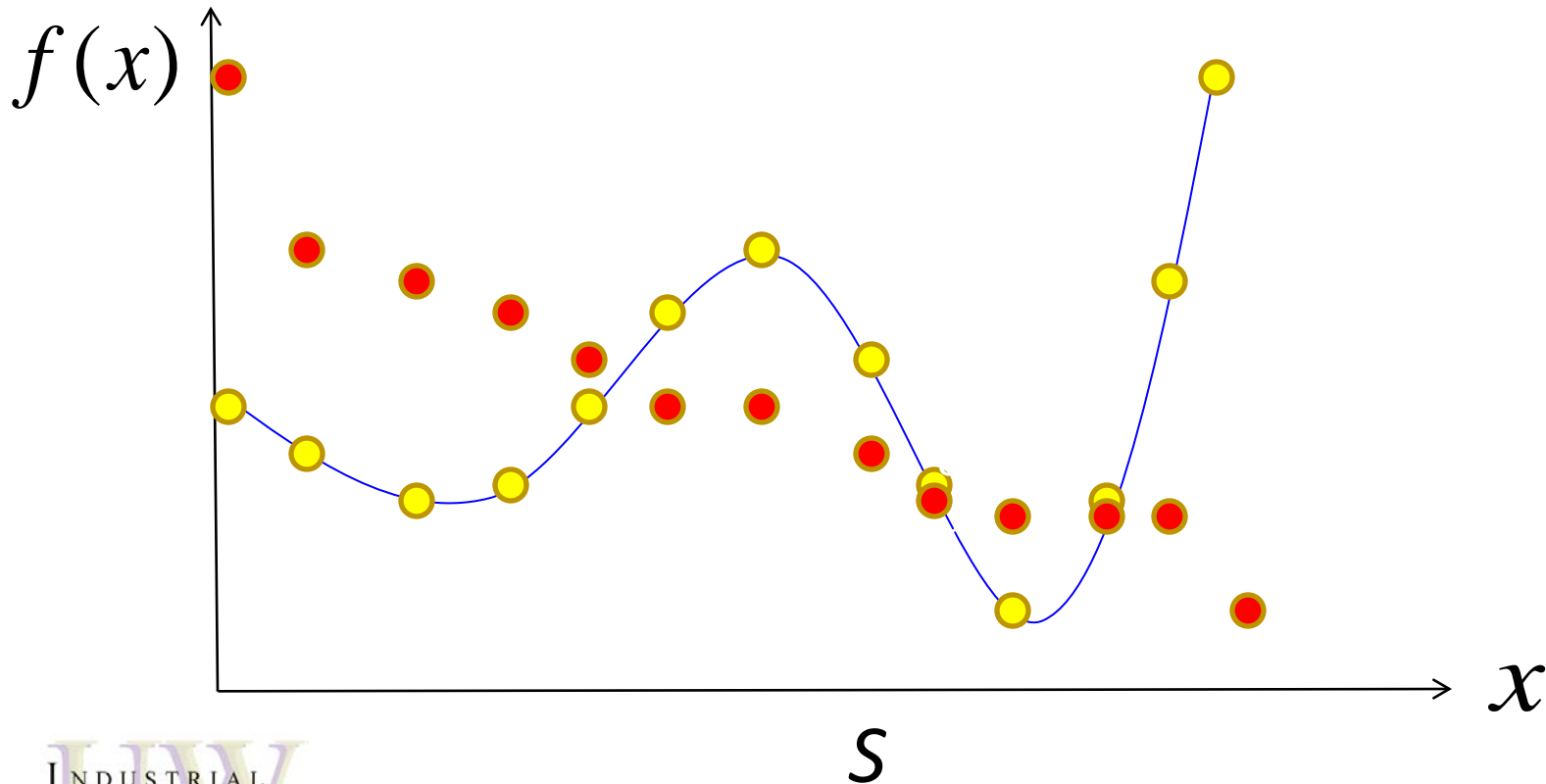ENGINEERING

# Local versus Global Optima

- o "Local" optima are relative to the neighborhood and algorithm

$f(x)$

$x$

$S$

# Local versus Global Optima

o "Local" optima are relative to the neighborhood and algorithm

# Research Question:
# What Do We Really Want?

o Do we really just want the optimum?

o What about sensitivity?

o Do we want to approximate the entire surface?

o Multi-criteria?

o Role of objective function and constraints?

o Where does randomness appear?

INDUSTRIAL
ENGINEERING

# How can we solve…?

IDEAL Algorithm:

o Optimizes any function quickly and accurately

o Provides information on how "good" the solution is

o Handles black-box and/or noisy functions, with continuous and/or discrete variables

o Is easy to implement and use
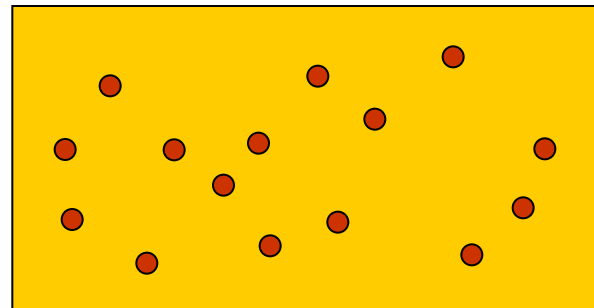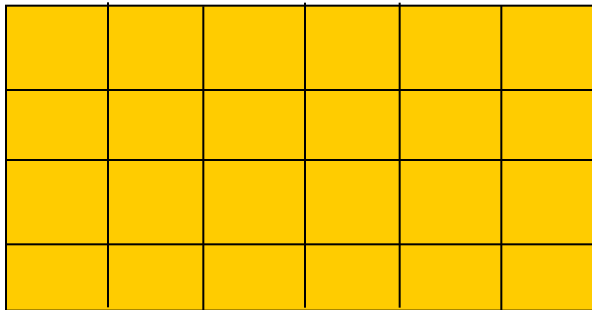
# Theoretical Performance of Stochastic Adaptive Search

o What kind of performance can we hope for?

o Global optimization problems are NP-hard

o Tradeoff between accuracy and computation

o Sacrifice guarantee of optimality for speed in finding a "good" solution

o Three theoretical constructs:

- Pure adaptive search (PAS)

- Hesitant adaptive search (HAS)

- Annealing adaptive search (AAS)

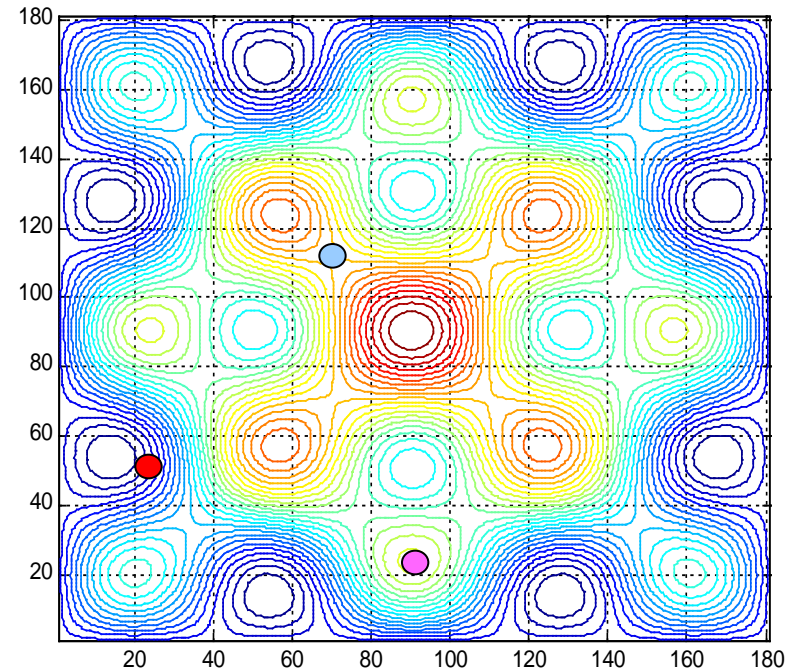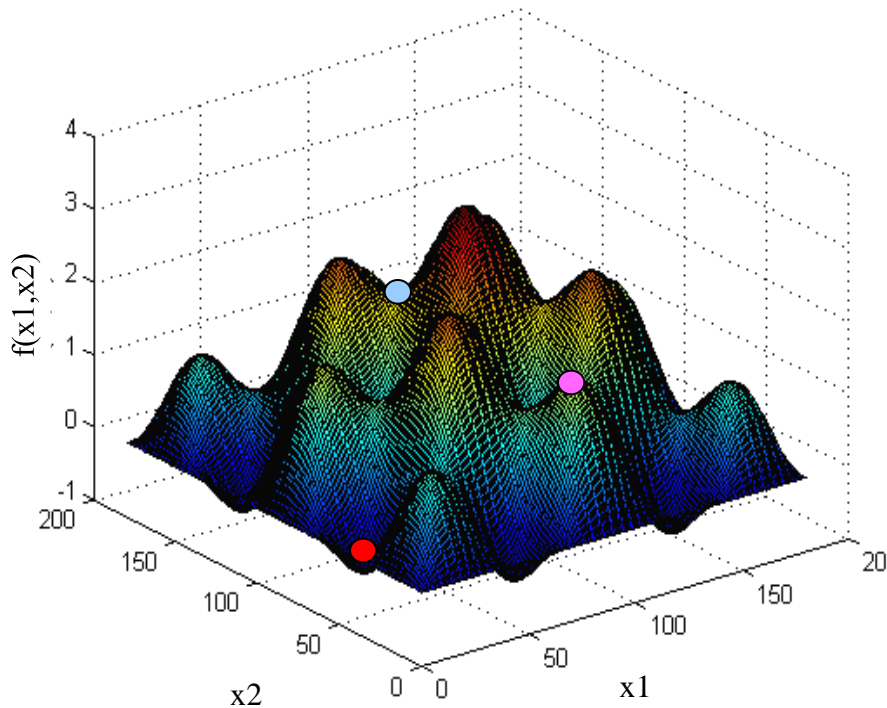INDUSTRIAL
ENGINEERING

13

# Performance of Two Simple Methods

- **Grid Search**:  Number of grid points is $O((L/\varepsilon)^n)$, where $L$ is the Lipschitz constant, $n$ is the dimension, and $\varepsilon$ is distance to the optimum

- **Pure Random Search**:  Expected number of points is $O(1/p(y^*+\varepsilon))$, where $p(y^*+\varepsilon)$ is the probability of sampling within $\varepsilon$ of the optimum $y^*$

- Complexity of both is exponential in dimension

# Pure Adaptive Search (PAS)

o PAS:   chooses points uniformly distributed in improving level sets

INDUSTRIAL
ENGINEERING

# Bounds on Expected Number of Iterations

o  PAS (continuous):

$$E[N(y^*+\varepsilon)] \leq 1 + \ln (1/p(y^*+\varepsilon))$$

where $p(y^*+\varepsilon)$ is the probability of PRS sampling within $\varepsilon$ of the global optimum $y^*$

o  PAS (finite):

$$E[N(y^*)] \leq 1 + \ln (1/p_1)$$

where $p_1$ is the probability of PRS sampling the global optimum

[Zabinsky and Smith, 1992]
[Zabinsky, Wood, Steel and Baritompa, 1995]

INDUSTRIAL
ENGINEERING

# Pure Adaptive Search

o **Theoretically, PAS is LINEAR in dimension**

o **Theorem:**

For any global optimization problem in $n$ dimensions, with Lipschitz constant at most $L$, and convex feasible region with diameter at most $D$, the expected number of PAS points to get within $\varepsilon$ of the global optimum is:

$$E[N(y^* + \varepsilon)] \leq 1 + n \ln(LD / \varepsilon)$$
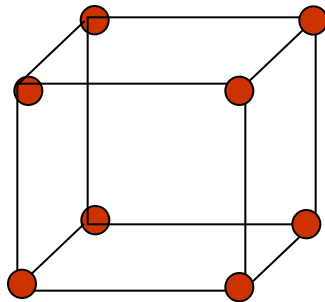
[Zabinsky and Smith, 1992]

INDUSTRIAL
ENGINEERING

# Finite PAS

o **Analogous LINEARITY result**

o **Theorem:**

For an *n* dimensional lattice { *1,...,k*}$^n$, with distinct objective function values, the expected number of points for PAS, sampling uniformly, to first reach the global optimum is:

$$E[N(y^*)] < 2 + n\ ln(k)$$

[Zabinsky, Wood, Steel and Baritompa, 1995]

# Hesitant Adaptive Search (HAS)

o What if we sample improving level sets with "bettering" probability *b(y)* and "hesitate" with probability *1-b(y)* ?

$$E[N(y^*+\varepsilon)] = \int_{y^*+\varepsilon}^{\infty} \frac{d\rho(t)}{b(t)\,p(t)}$$

where $\rho(t)$ is the underlying sampling distribution and $p(t)$ is the probability of sampling *t* or better

[Bulger and Wood, 1998]

INDUSTRIAL
ENGINEERING

19

# General HAS

o For a mixed discrete and continuous global optimization problem, the expected value of $N(y^*+\varepsilon)$, the variance, and the complete distribution can be expressed using the sampling distribution $\rho(t)$ and bettering probabilities $b(y)$

[Wood, Zabinsky and Kristinsdottir, 2001]

INDUSTRIAL
ENGINEERING

# Annealing Adaptive Search (AAS)

o What if we sample from the original feasible region each iteration, but change distributions?

o Generate points over the whole domain using a Boltzmann distribution parameterized by temperature *T*

- Boltzmann distribution becomes more concentrated around the global optima as the temperature decreases
- Temperature is determined by a cooling schedule

o The *record values* of AAS are dominated by PAS and thus LINEAR in dimension

[Romeijn and Smith, 1994]

INDUSTRIAL
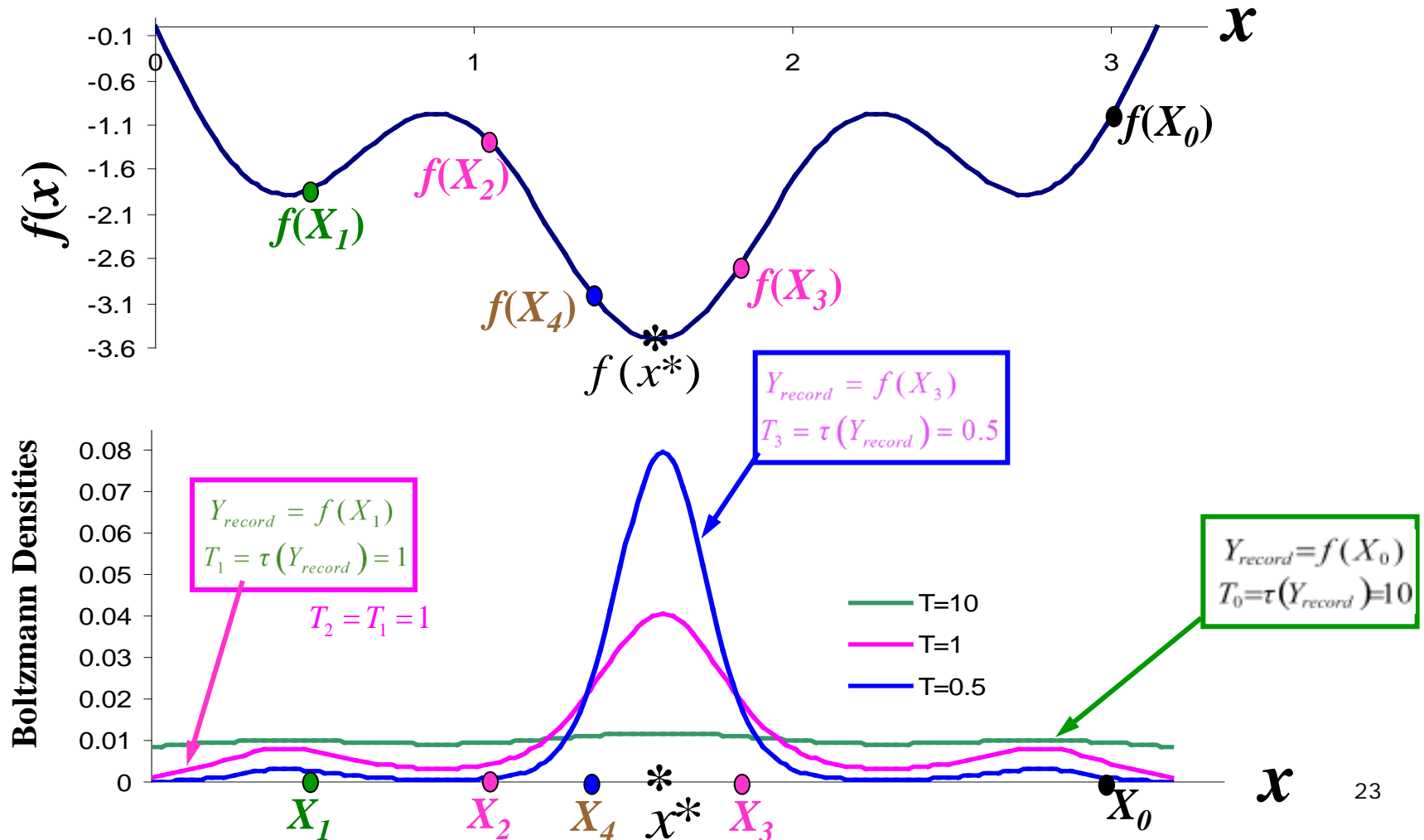ENGINEERING

# Performance of Annealing Adaptive Search

o The expected number of *sample points* of AAS is bounded by HAS with a specific *b(y)*

o Select the next temperature so that the probability of generating an improvement under that Boltzmann distribution is at least 1-α , i.e.,

$$P\left(Y_{R(k)+1}^{AAS} < y \mid Y_{R(k)}^{AAS} = y\right) \geq 1 - \alpha$$

o Then the *expected number of AAS sample points* is LINEAR in dimension

[Shen, Kiatsupaibul, Zabinsky and Smith, 2007]

22

# AAS with Adaptive Cooling Schedule

# Research Areas

o **Develop theoretical analysis of PAS, HAS, AAS for noisy or approximate functions**

- Model approximation or estimation error
- Characterize impact of error on performance

o **Use theory to develop algorithms**

- Approximate sampling from improving sets (as PAS) or Boltzmann distributions (as AAS)
- Use HAS, with $\rho(t)$ and $b(y)$, to quantify and balance accuracy and efficiency

INDUSTRIAL
ENGINEERING

# Random Search Algorithms

o **Instance-based methods**

- Sequential random search
- Multi-start and population-based algorithms

o Model-based methods

- Importance sampling
- Cross-entropy    [Rubinstein and Kroese, 2004]
- Model reference adaptive search [Hu, Fu and Marcus, 2007]

[Zlochin, Birattari, Meuleau and Dorigo, 2004]

INDUSTRIAL
ENGINEERING

# Sequential Random Search

o Stochastic approximation [Robbins and Monro, 1951]

o Step-size algorithms [Rastrigin, 1960] [Solis and Wets, 1981]

o Simulated annealing
   [Romeijn and Smith, 1994],  [Alrafaei and Andradottir, 1999]

o Tabu search [Glover and Kochenberger, 2003]

o Nested partition [Shi and Olafsson, 2000]

o COMPASS [Hong and Nelson, 2006]

o View these algorithms as Markov chains with
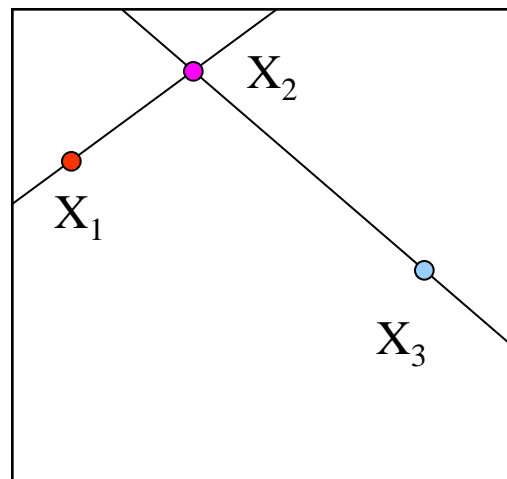  • Candidate point generators
  • Update procedures

# Use Hit-and-Run to Approximate AAS

o Hit-and-Run is a Markov chain Monte Carlo (MCMC) sampler

- converges to a uniform distribution

  [Smith, 1984]

- in polynomial time $O(n^3)$
  [Lovász, 1999]

- can approximate any arbitrary distribution by using a filter

o The difficulty of implementing AAS is to generate points directly from a family of Boltzmann distributions
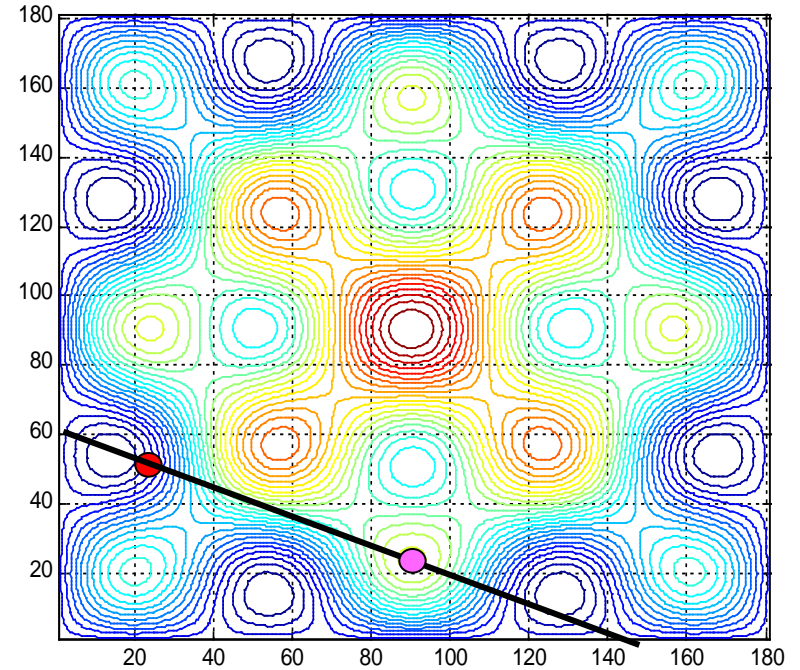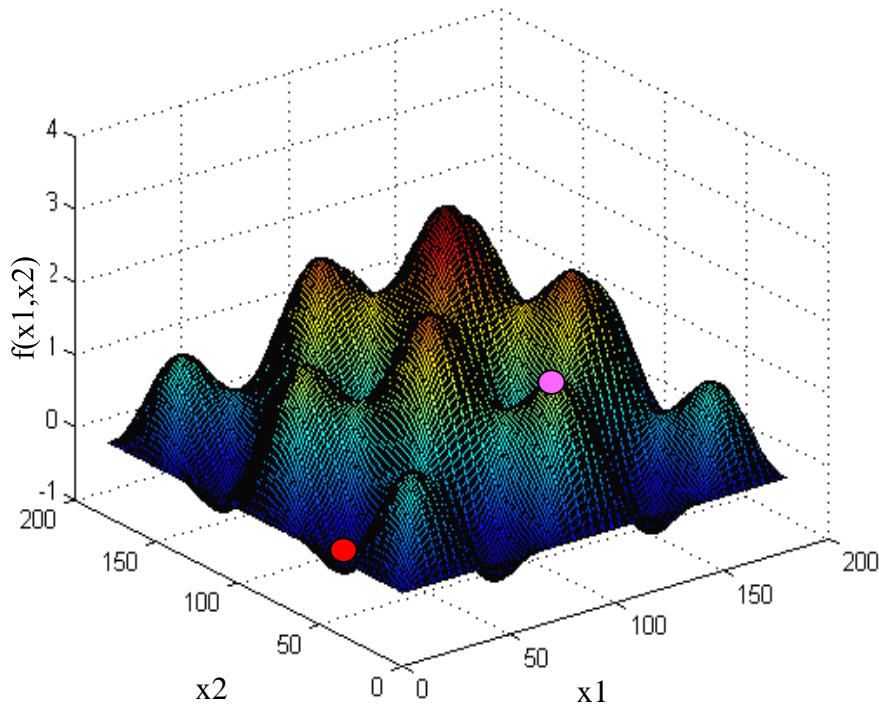
INDUSTRIAL
ENGINEERING

# Hit-and-Run

o Hit-and-Run generates a random direction (uniformly distributed on a hypersphere) and a random point (uniformly distributed on the line)
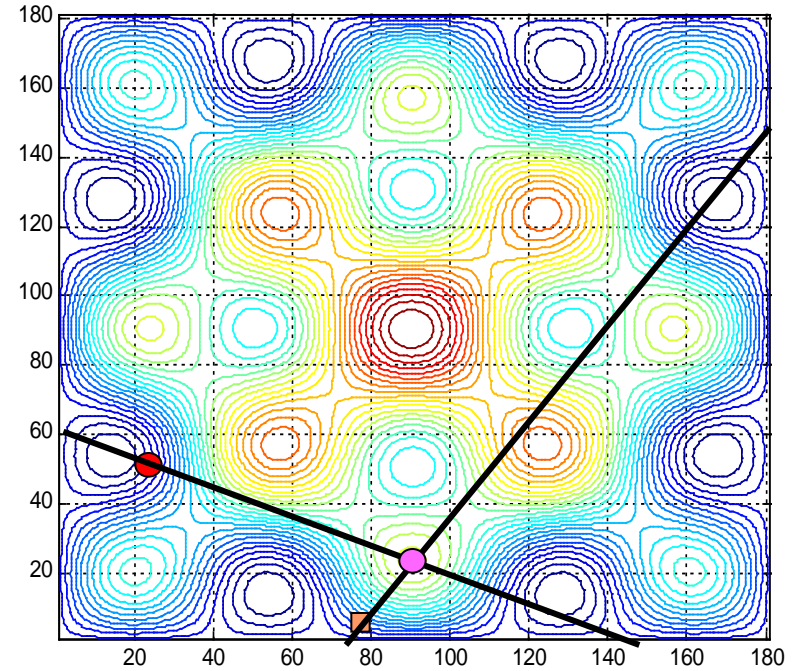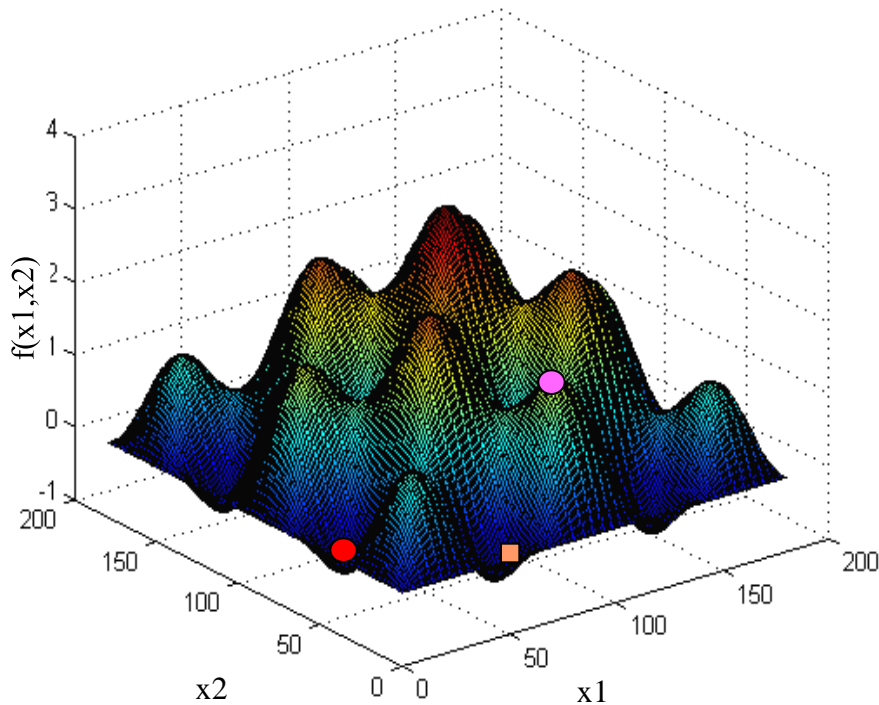
# Improving Hit-and-Run

o IHR: choose a random direction and a random point, accept only improving points
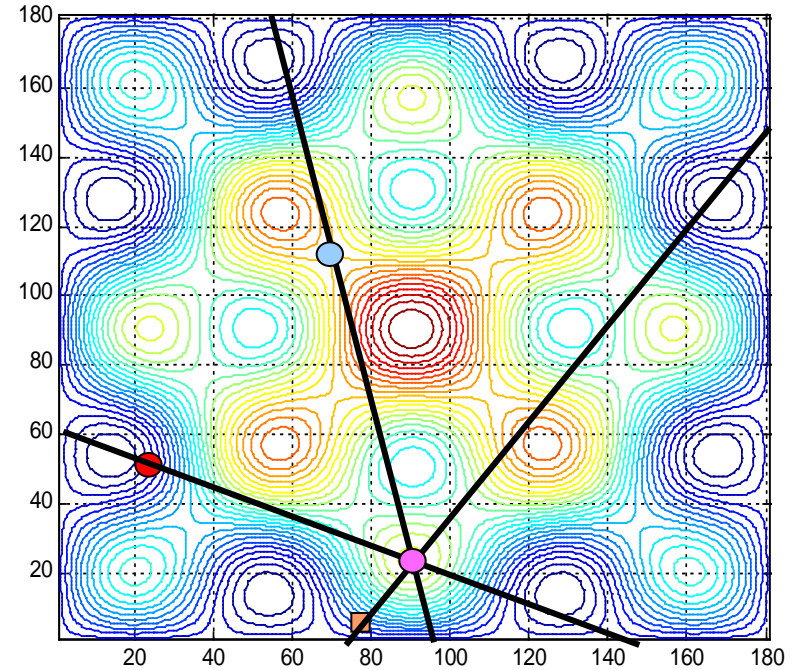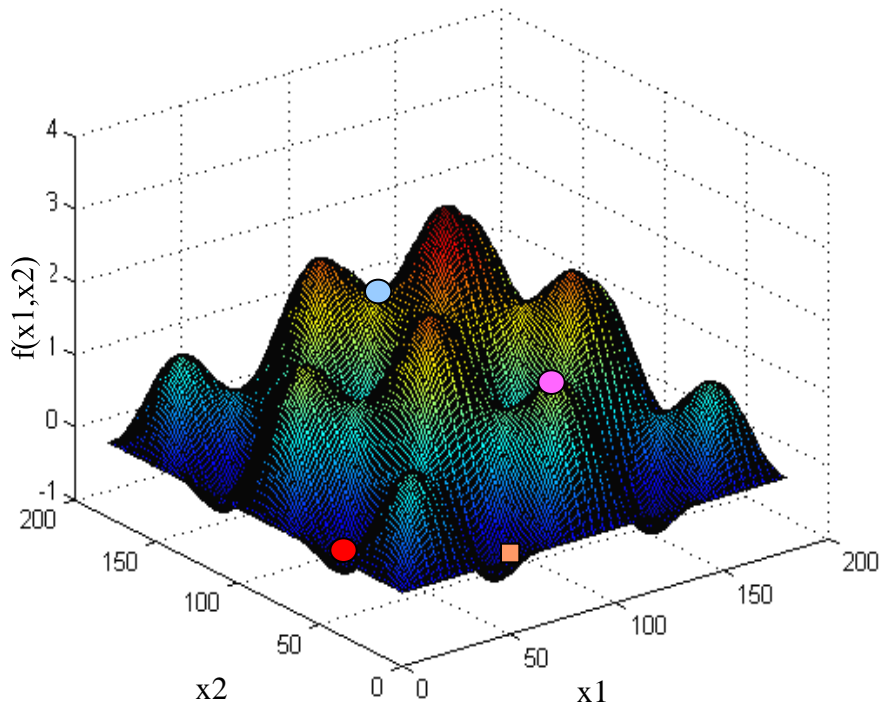
# Improving Hit-and-Run

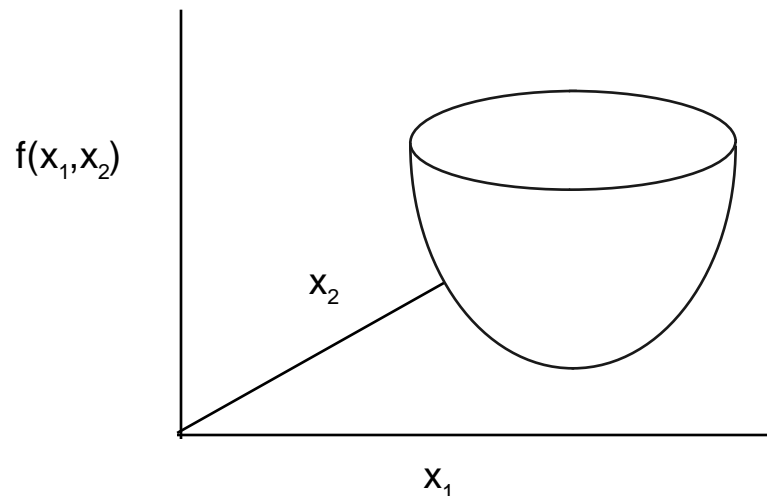o   IHR:  choose a random direction and a random point, accept only improving points

INDUSTRIAL
ENGINEERING

# Improving Hit-and-Run

o IHR: choose a random direction and a random point, accept only improving points

# Is IHR Efficient in Dimension?

o Theorem:
   For any elliptical program in *n* dimensions, the expected number of function evaluations for IHR is: $O(n^{5/2})$

   [Zabinsky, Smith, McDonald, Romeijn and Kaufman, 1993]



$f(x_1,x_2)$

$x_2$

$x_1$

# Is IHR Efficient in Dimension?

o  Theorem:
   For any elliptical program in *n* dimensions, the
   expected number of function evaluations for
   IHR is:      $O(n^{5/2})$
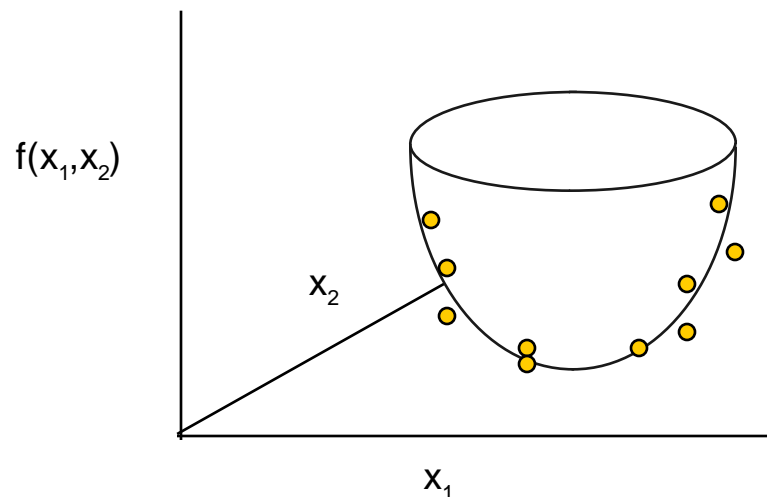
[Zabinsky, Smith, McDonald, Romeijn and Kaufman, 1993]

$f(x_1, x_2)$

$x_2$

$x_1$

# Use Hit-and-Run to Approximate Annealing Adaptive Search

o Hide-and-Seek: add a probabilistic Metropolis acceptance-rejection criterion to Hit-and-Run to approximate the Boltzmann distribution
  [Romeijn and Smith, 1994]

o Converges in probability with almost any cooling schedule driving temperature to zero

o AAS Adaptive Cooling Schedule:
  - Temperature values according to AAS to maintain $1-\alpha$ probability of improvement
  - Update temperature when record values are obtained     [Shen, Kiatsupaibul, Zabinsky and Smith, 2007]

INDUSTRIAL
ENGINEERING

# Research Possibilities:

o How long should we execute Hit-and-Run at a fixed temperature?

o What is the benefit of sequential temperatures (warm starts) on convergence rate?

o Hit-and-Run has fast convergence on "well-rounded" sets; how can we modify transition kernel in general?

o Incorporate new Hit-and-Run on mixed integer/continuous sets

  • Discrete hit-and-run
    [Baumert, Ghate, Kiatsupaibul, Shen, Smith and Zabinsky, 2009]

  • Pattern hit-and-run
    [Mete, Shen, Zabinsky, Kiatsupaibul and Smith, 2010]

INDUSTRIAL ENGINEERING

# Simulated Annealing with Multi-start: When to Stop or Restart a Run?

o Use HAS to model progress of a heuristic random search algorithm and estimate associated parameters

o Dynamic Multi-start Sequential Search
  - If current run appears "stuck" according to HAS analysis, stop and restart
  - Estimate probability of achieving $y^*+\varepsilon$ based on observed values and estimated parameters
  - If probability is high enough, terminate

[Zabinsky, Bulger and Khompatraporn, 2010]

INDUSTRIAL
ENGINEERING

# Meta-control of Interacting-Particle Algorithm

o **Interacting-Particle Algorithm**

- Combines simulated annealing and population based algorithms

- Uses statistical physics and Feynman-Kac formulas to develop selection probabilities

[Del Moral, *Feynman-Kac Formulae: Genealogical and Interacting Particle Systems with Applications*, 2004]

o Meta-control approach to dynamically heat and cool temperature
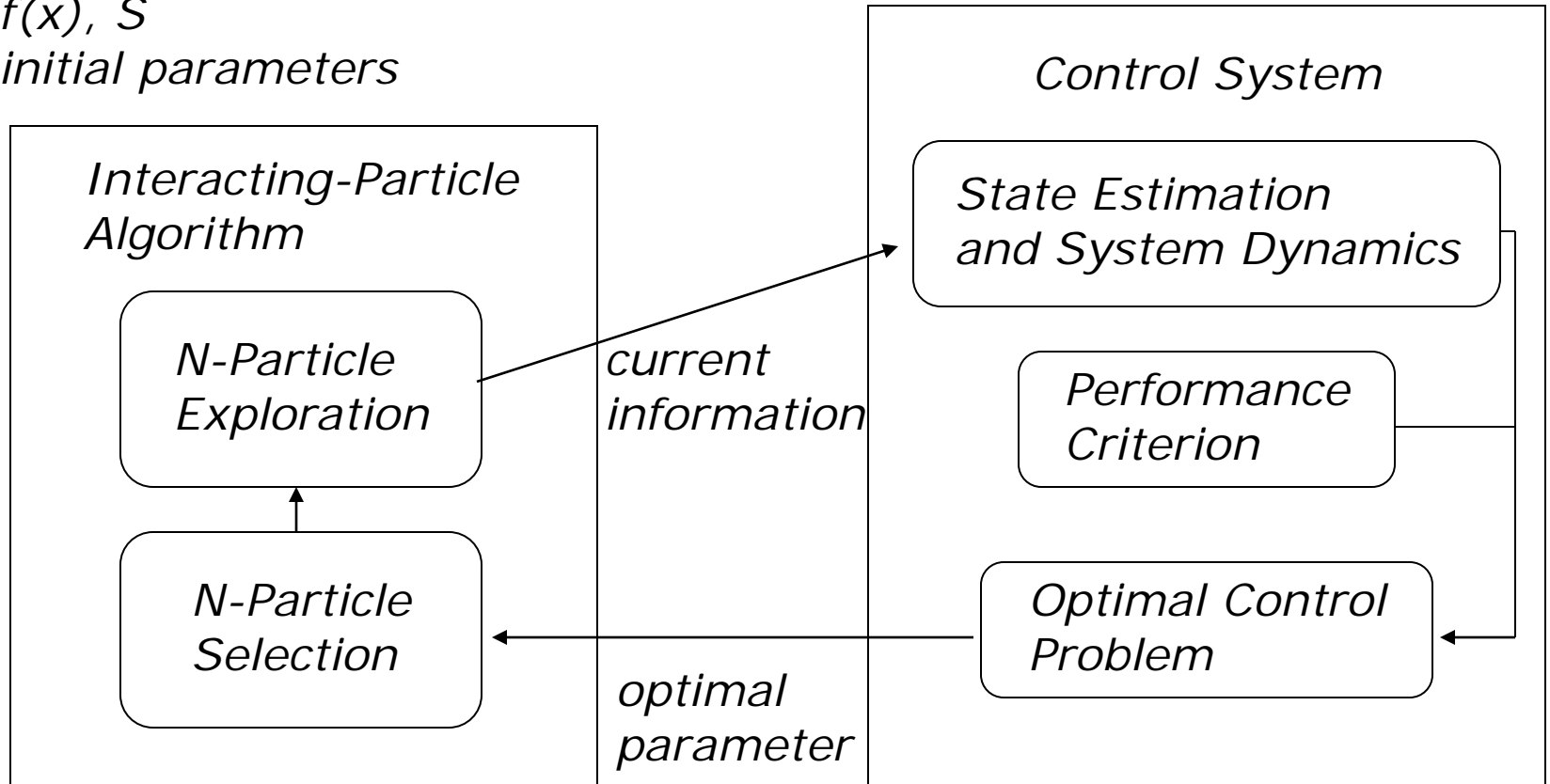
[Kohn, Zabinsky and Brayman, 2006]
[Molvalioglu, Zabinsky and Kohn, 2009]

INDUSTRIAL
ENGINEERING

# Meta-control Approach

f(x), S
initial parameters

Control System

Interacting-Particle Algorithm

State Estimation
and System Dynamics

N-Particle
Exploration

current
information

Performance
Criterion

N-Particle
Selection

optimal
parameter

Optimal Control
Problem

# Research Possibilities

o Combine theoretical analyses with MCMC and meta-control to:

- Control the exploration transition probabilities
- Obtain stopping criterion and quality of solution
- Relate interacting particles to cloning/splitting

o Combine theoretical analyses and meta-control with model-based approach

# Another Research Area: Quantum Global Optimization

o **Grover's Adaptive Search can implement PAS on a quantum computer**
[Baritompa, Bulger and Wood, 2005]

o **Apply research on quantum control theory to global optimization**

- [Gardiner, *Handbook of Stochastic Methods for Physics, Chemistry and the Natural Sciences*, 2004]
- [Del Moral, *Feynman-Kac Formulae: Genealogical and Interacting Particle Systems with Applications*, 2004]
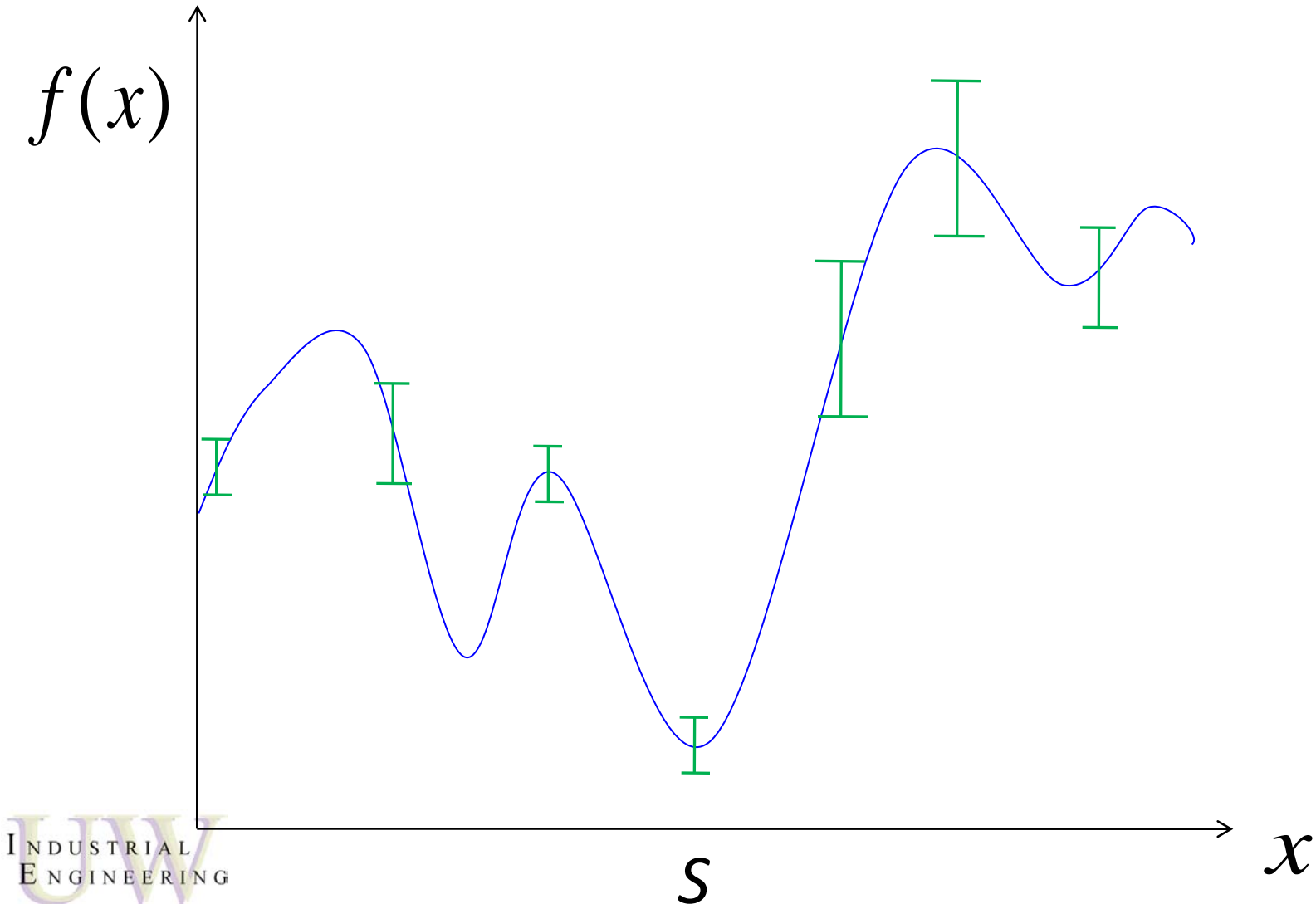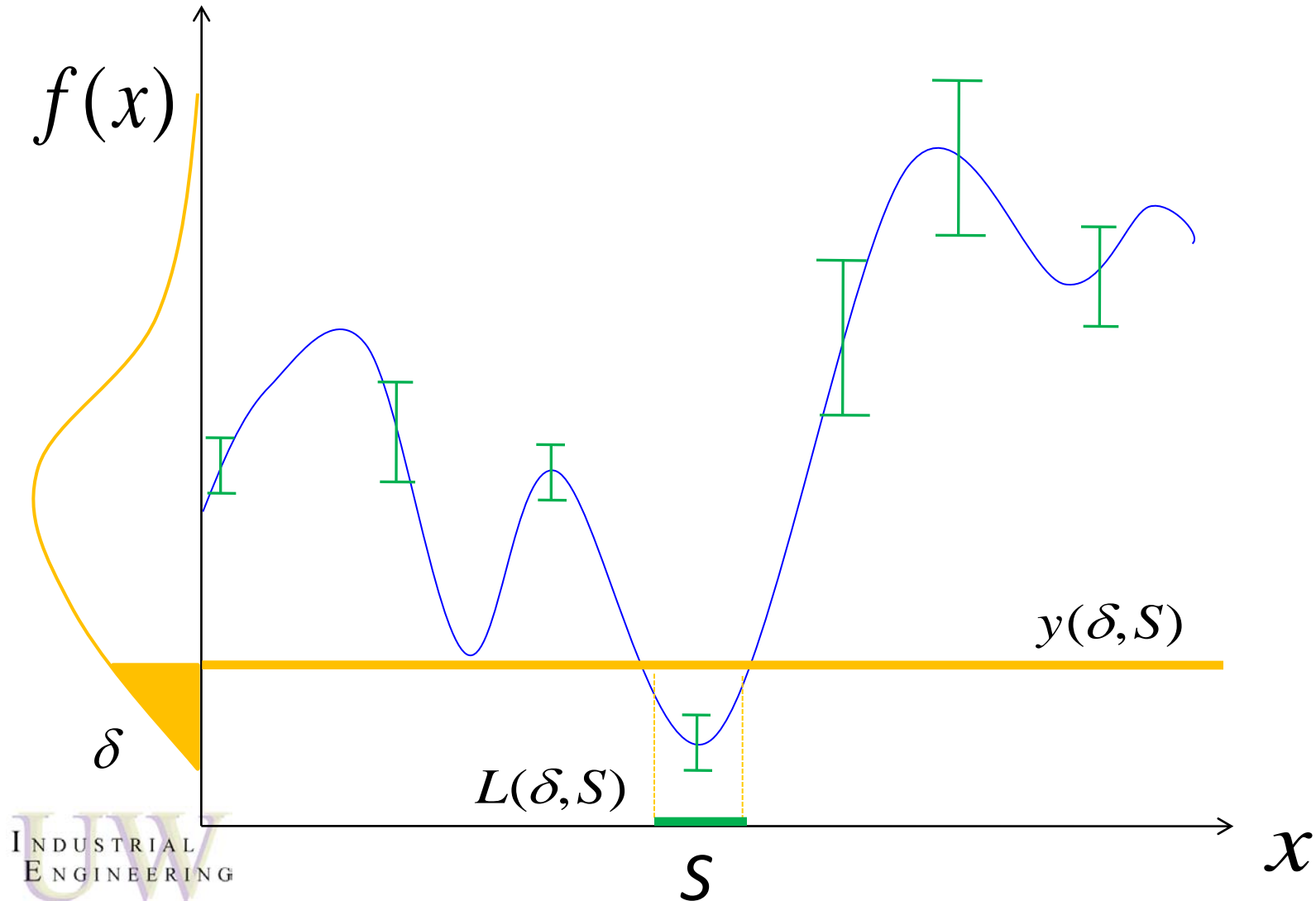
# Optimization of Noisy Functions

o Use random sampling to explore the feasible region and estimate the objective function with replications

o Recognize two sources of noise:

- Randomness in the sampling distribution
- Randomness in the objective function

o Adaptively adjust the number of samples and the number of replications

INDUSTRIAL
ENGINEERING

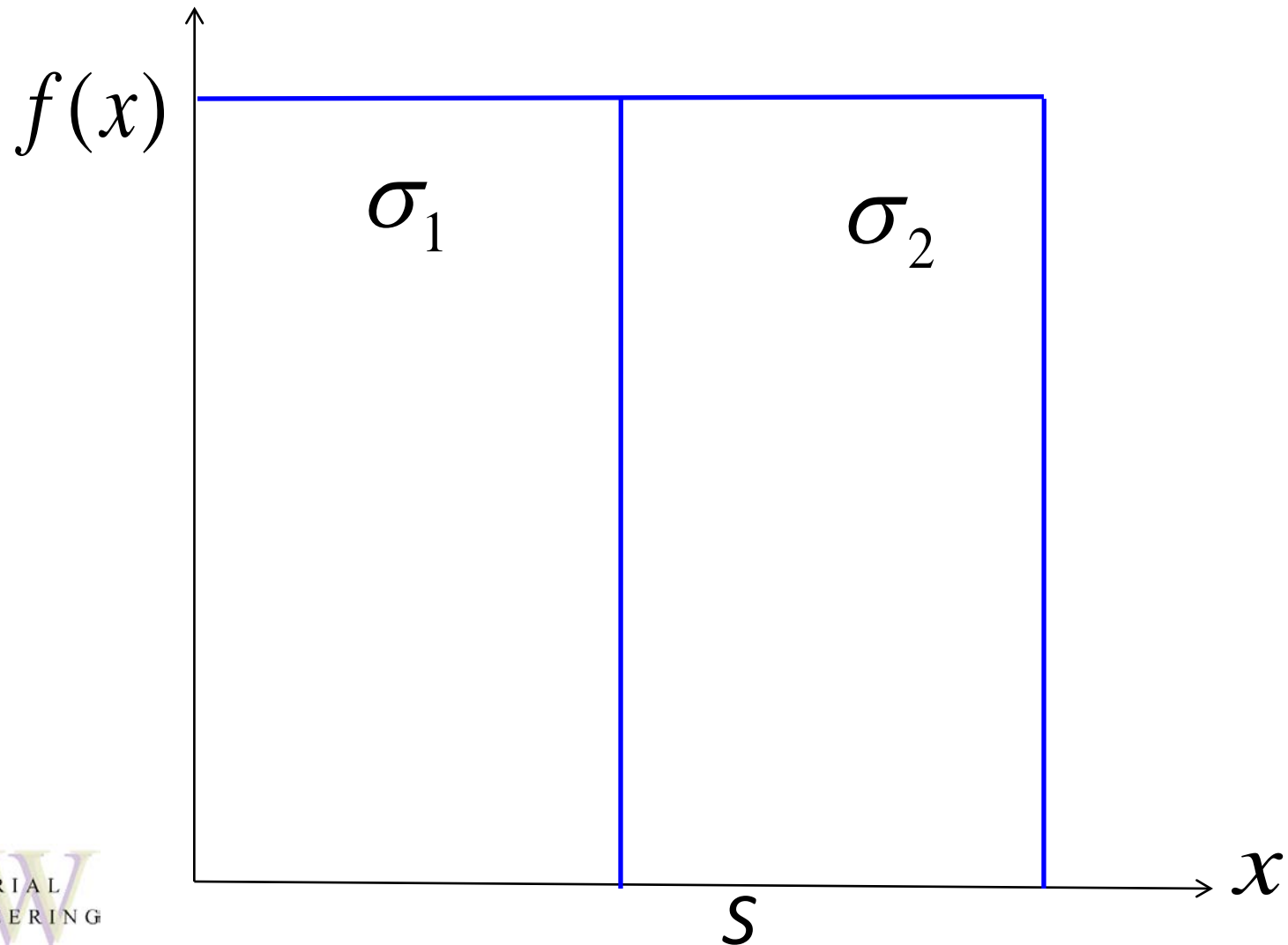# Noisy Objective Function



$f(x)$

$s$

$x$

INDUSTRIAL
ENGINEERING

# Noisy Objective Function



$f(x)$

$y(\delta, S)$
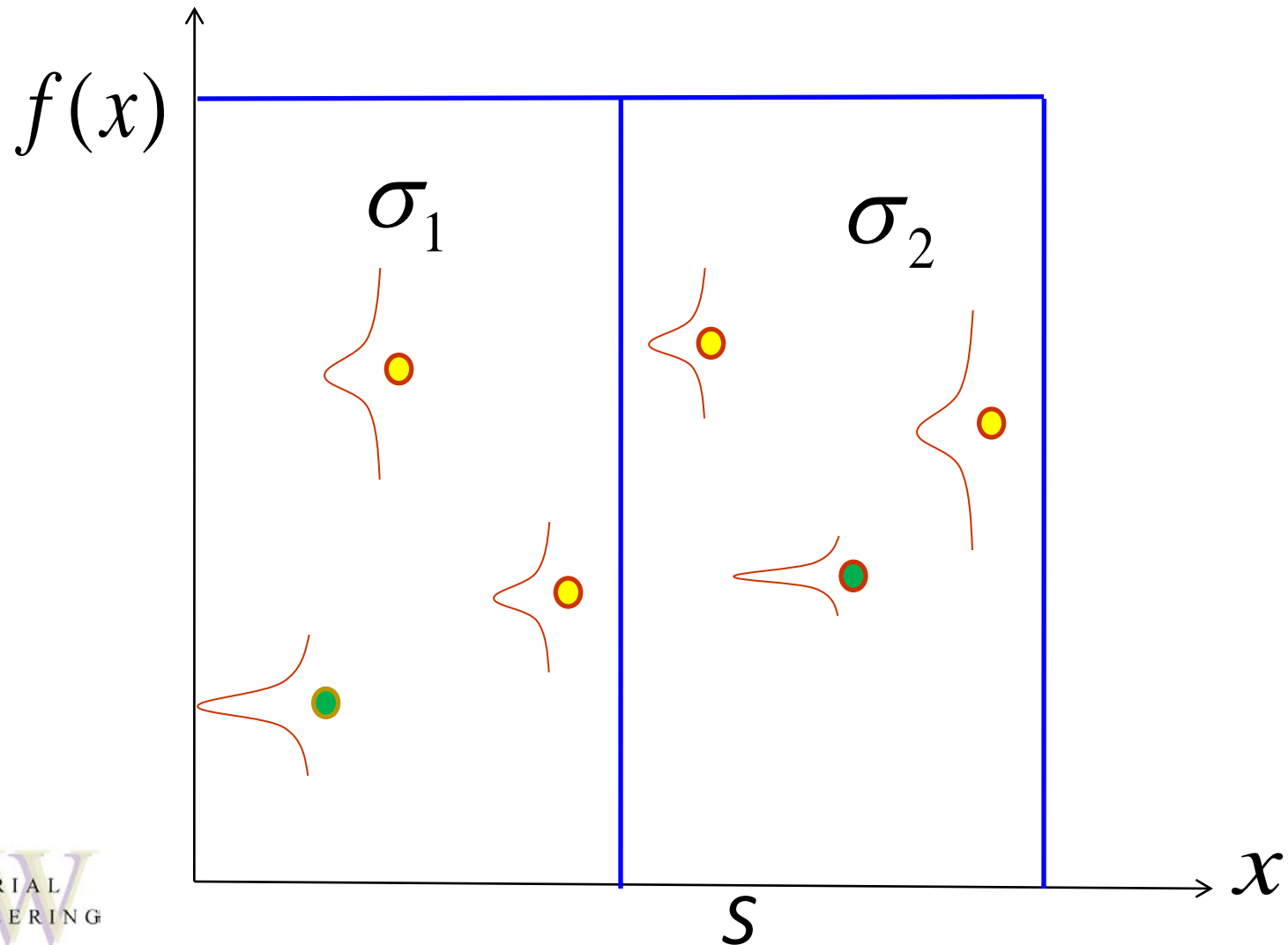
$\delta$

$L(\delta, S)$

$S$

$x$

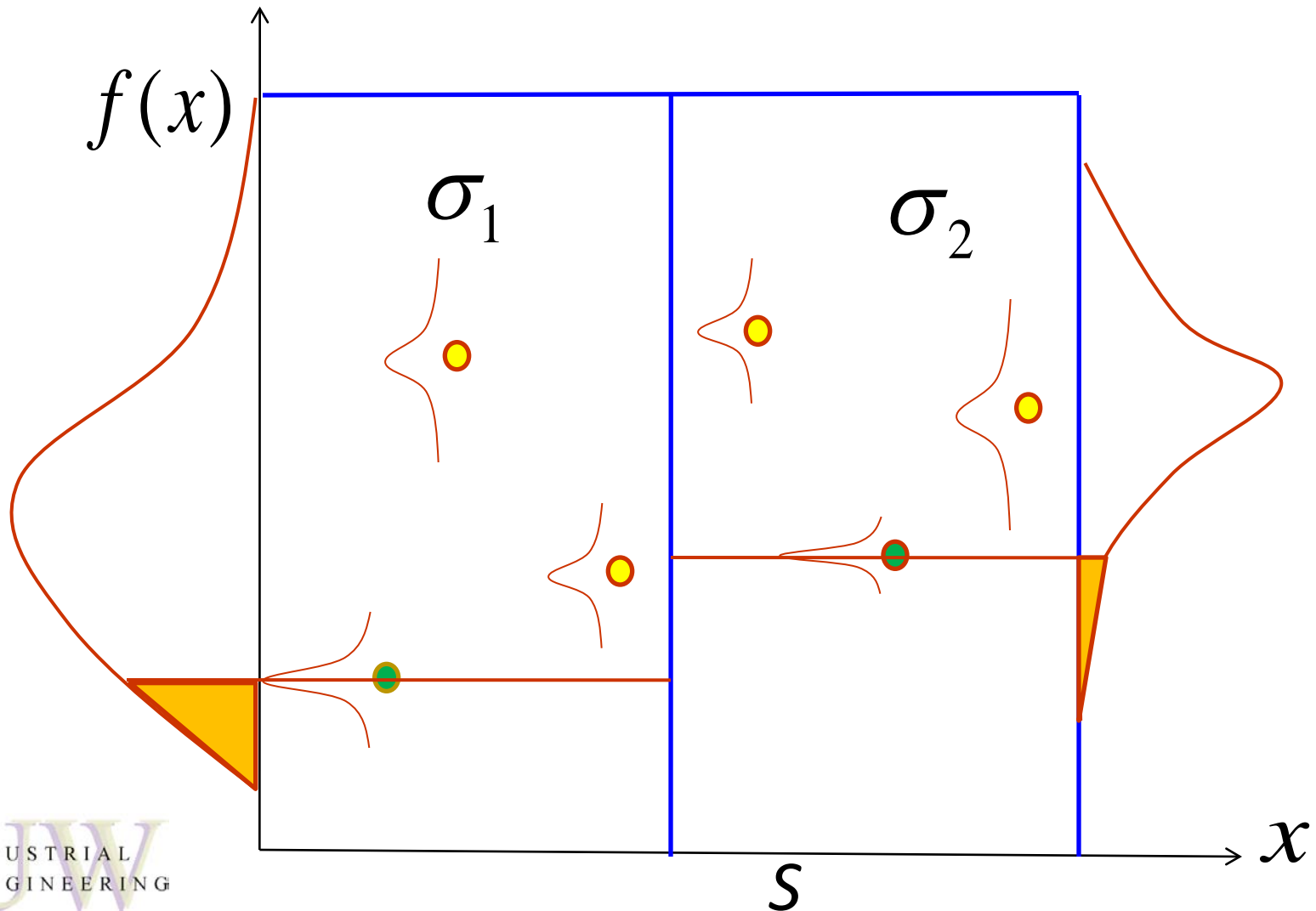INDUSTRIAL
ENGINEERING

# Probabilistic Branch-and-Bound (PBnB)

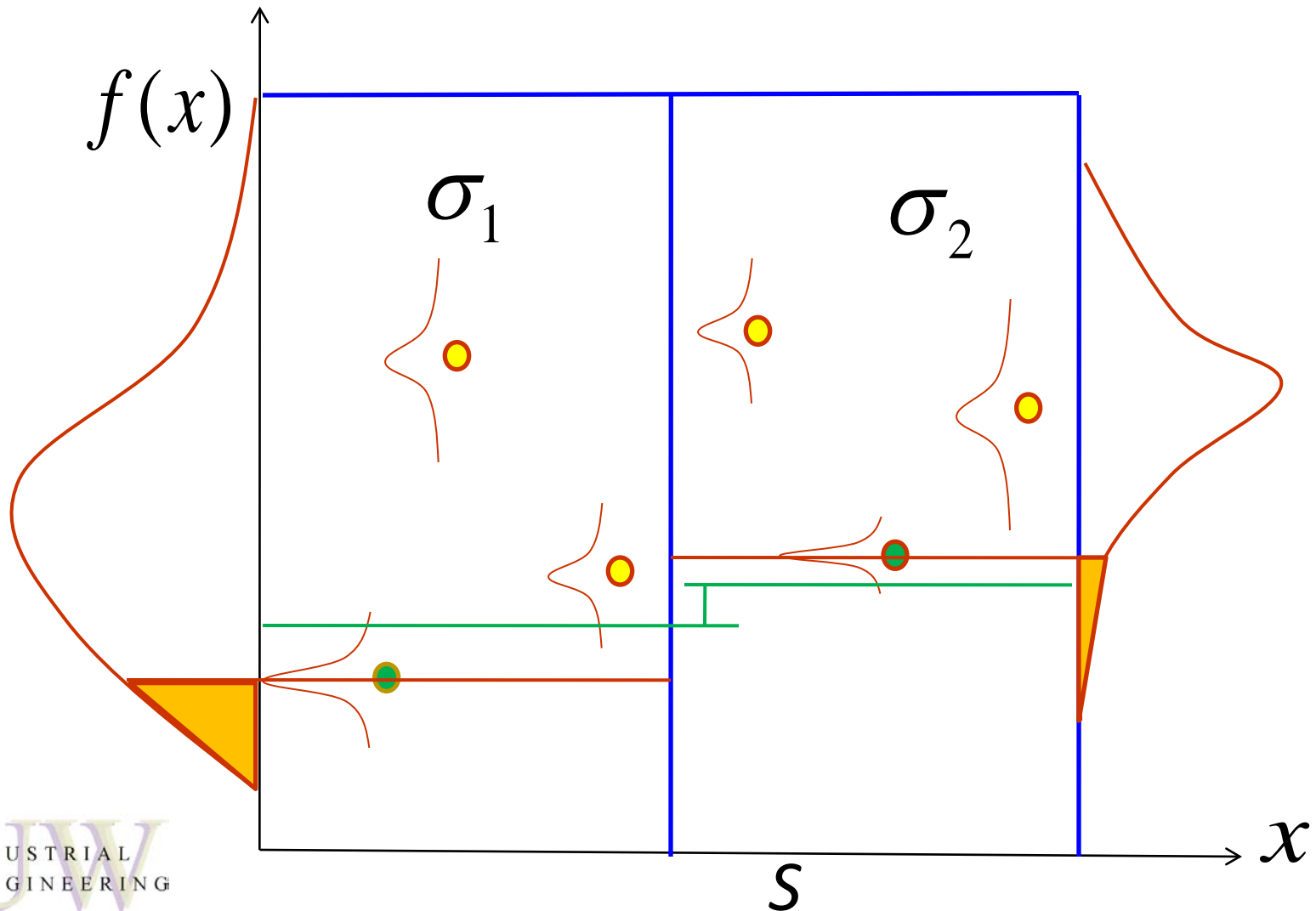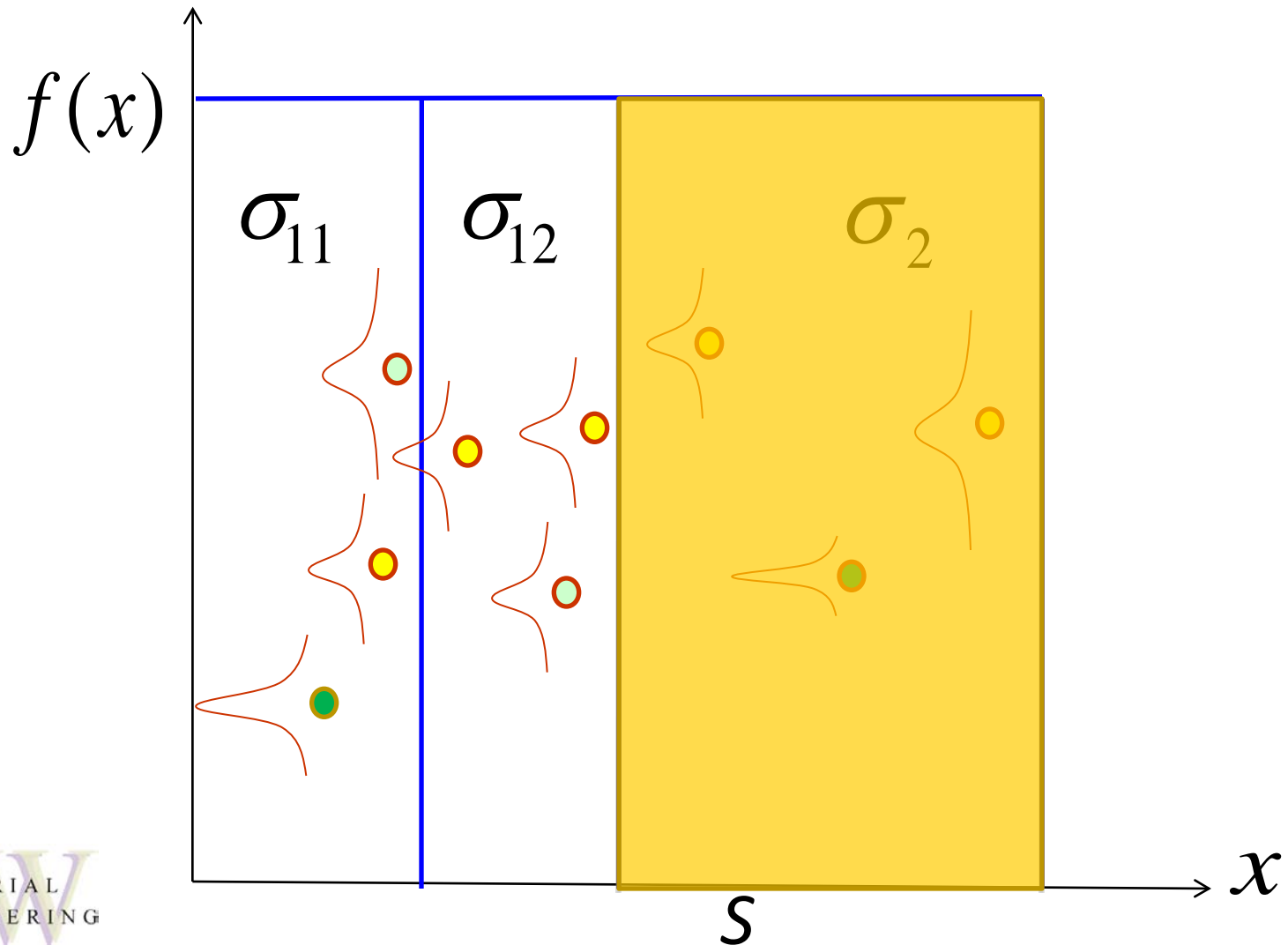# Sample $N^*$ Uniform Random Points with $R^*$ Replications

# Use Order Statistics to Assess Range Distribution
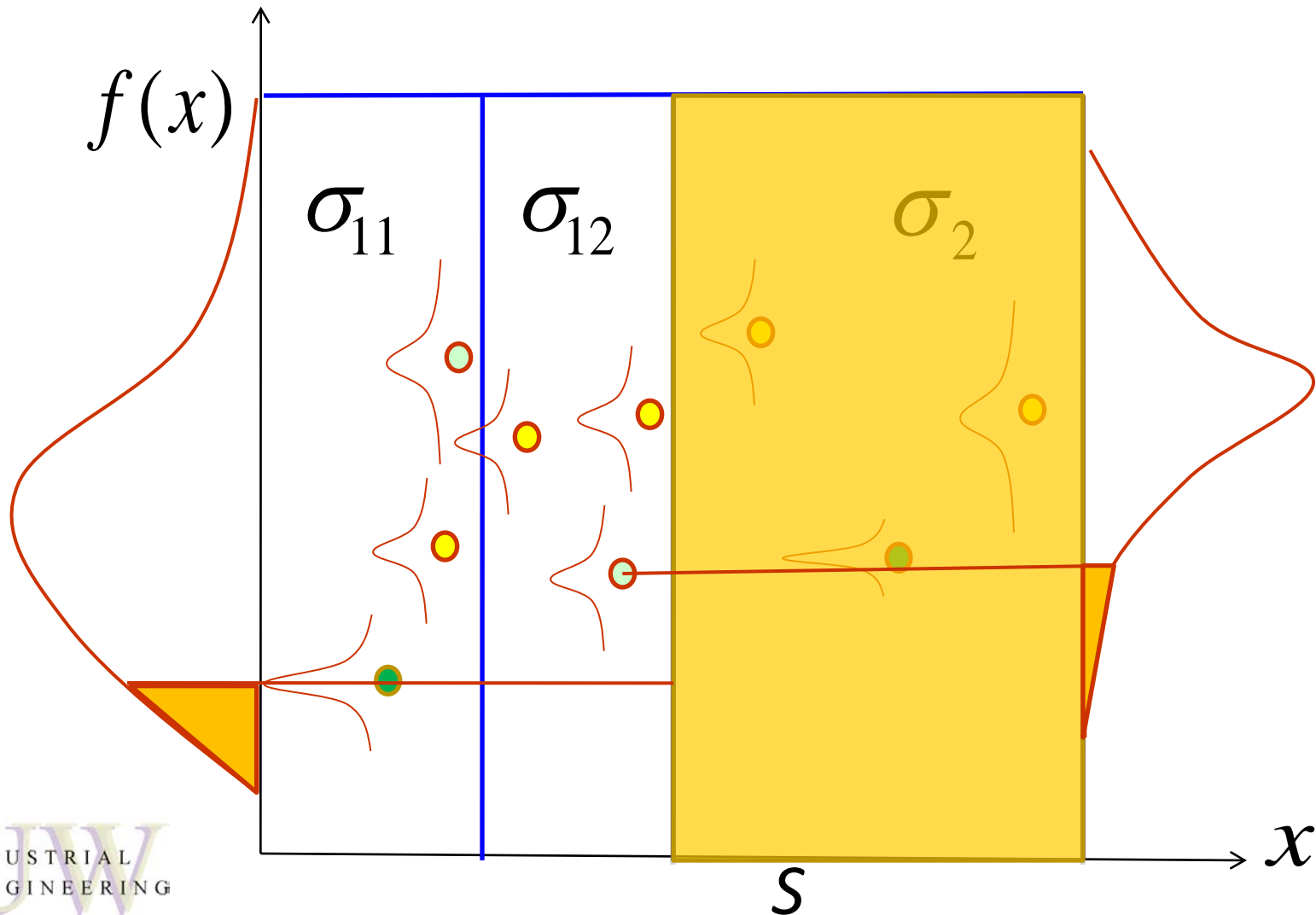
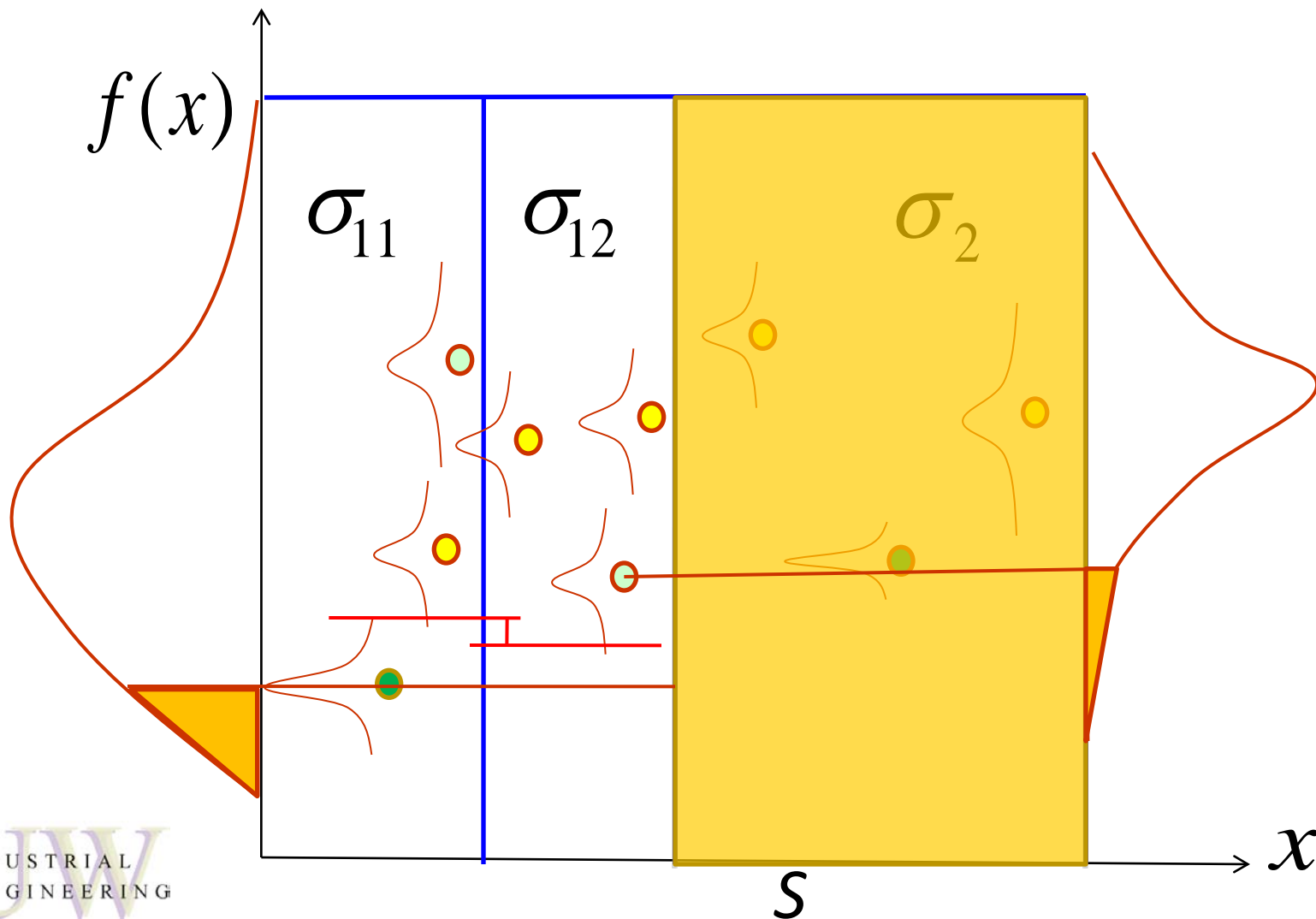# Prune, if Statistically Confident

# Subdivide & Sample Additional Points
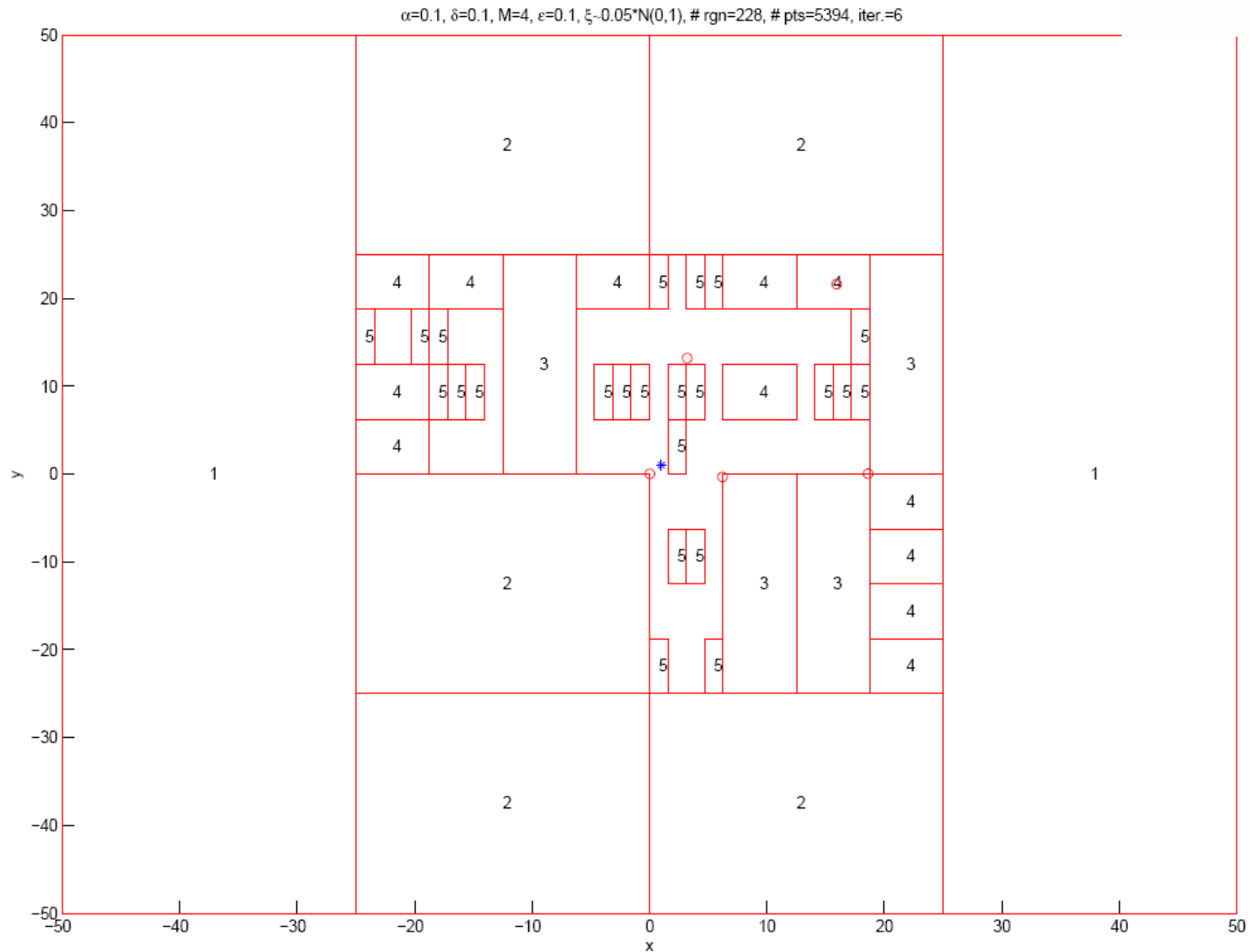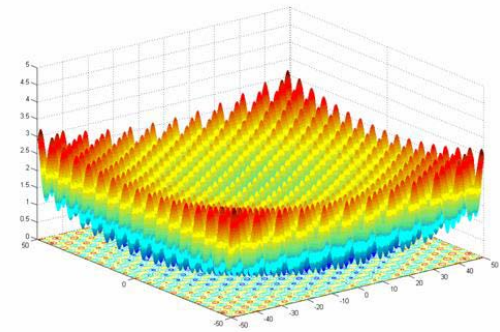
# Reassess Range Distribution

# If No Pruning, Then Continue …

# PBnB:
# Numerical Example

# Research Areas

o Develop theory to tradeoff accuracy with computational effort

o Use theory to develop algorithms that give insight into original problem

- Global optima and sensitivity
- Shape of the function or range distribution

o Use interdisciplinary approaches to incorporate feedback

# Summary

o Theoretical analysis of PAS, HAS, AAS motivates random search algorithms

o Hit-and-Run is a MCMC method to approximate theoretical performance

o Meta-control theory allows adaptation based on observations

o Probabilistic Branch-and-Bound incorporates noisy function evaluations and sampling noise into analysis

# Additional Slides for Details on Interacting Particle Algorithm

# Interacting-Particle Algorithm

o Simulated Annealing: Markov chain Monte Carlo method for approximating a sequence of Boltzmann distributions

$$\eta_t(dx) = \frac{e^{-f(x)/T_t}}{\int_S e^{-f(y)/T_t} \, dy} \, dx$$

o Population-based Algorithms: simulate a distribution (e.g. Feynman-Kac annealing model)such that

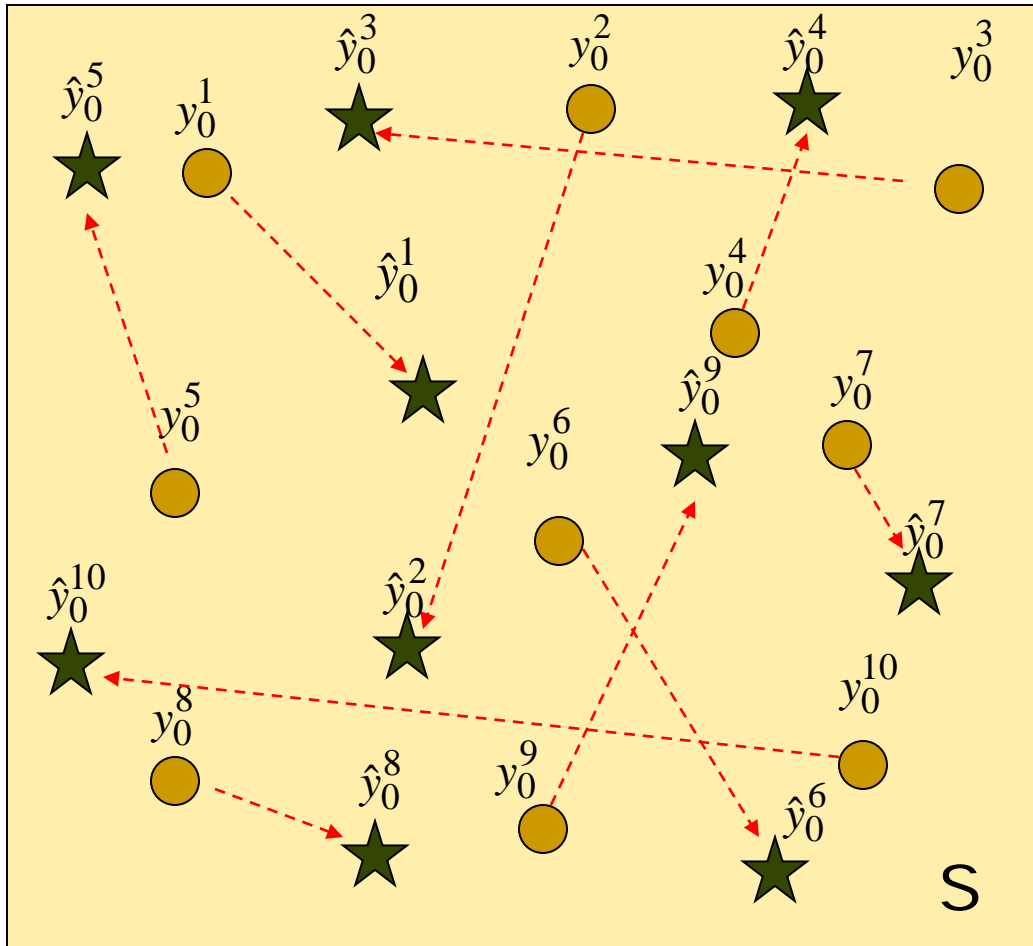$$E_{\eta_t}(f) \to y^* \ \text{ as } \ t \to \infty$$

# Interacting-Particle Algorithm

- **Initialization:** Sample the initial locations $y_0^i$ for particle $i=1,...,N$ from the distribution $\eta_0$

For $t=0,1,....,$

- **N-particle exploration:** Move particle $i=1,...,N$ to location $\hat{y}_t^i$ with probability distribution $E(y_t^i, d\hat{y}_t^i)$

- **Temperature Parameter Update:**
$$T_t = (1+\varepsilon_t)T_{t-1}, \qquad \varepsilon_t \geq -1$$

- **N-particle selection:** Set the particles' locations $y_{t+1}^k, i=1,...,N$

$$y_{t+1}^k = \hat{y}_t^i \text{ with probability} \qquad \frac{e^{-f(\hat{y}_t^i)/T_t}}{\displaystyle\sum_{j=1}^{N} e^{-f(\hat{y}_t^j)/T_t}}$$

# Illustration of Interacting-Particle Algorithm



**Initialization:** Sample N=10 points uniformly on S

**N-particle exploration:**

Move particles from ⬤ to ★ using Markov kernel E

# Illustration of Interacting-Particle Algorithm


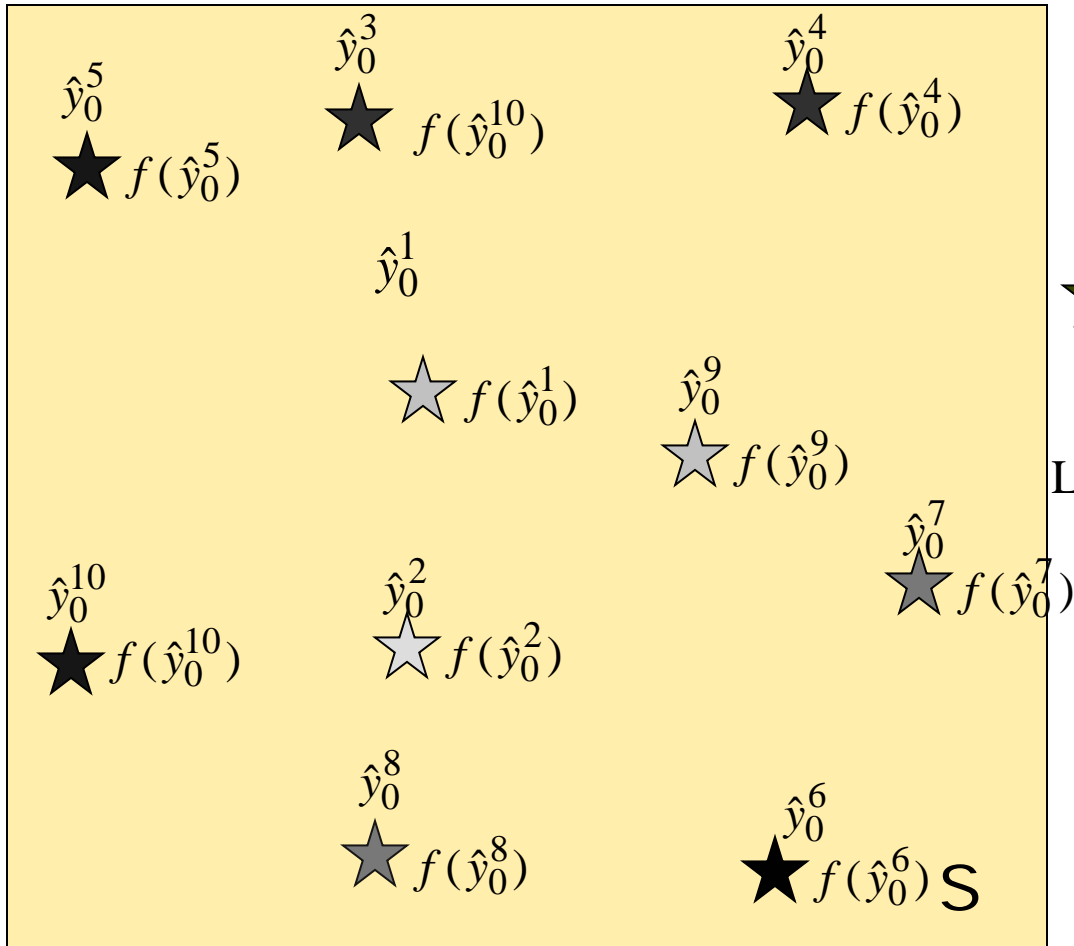
**Initialization:** Sample N=10 points uniformly on S

**N-particle exploration:**

Move particles from ⬤ to

⭐ using Markov kernel E

Evaluate $f(.)$

Lower $f(.)$ ⟷ Higher $f(.)$

$\hat{y}_0^3$

$\hat{y}_0^4$ $f(\hat{y}_0^4)$

$\hat{y}_0^5$

$f(\hat{y}_0^5)$

$f(\hat{y}_0^{10})$

$\hat{y}_0^1$

$f(\hat{y}_0^1)$

$\hat{y}_0^9$

$f(\hat{y}_0^9)$

$\hat{y}_0^7$

$f(\hat{y}_0^7)$

$\hat{y}_0^{10}$

$f(\hat{y}_0^{10})$

$\hat{y}_0^2$

$f(\hat{y}_0^2)$

$\hat{y}_0^8$

$f(\hat{y}_0^8)$

$\hat{y}_0^6$

$f(\hat{y}_0^6)$ S

# Illustration of Interacting-Particle Algorithm



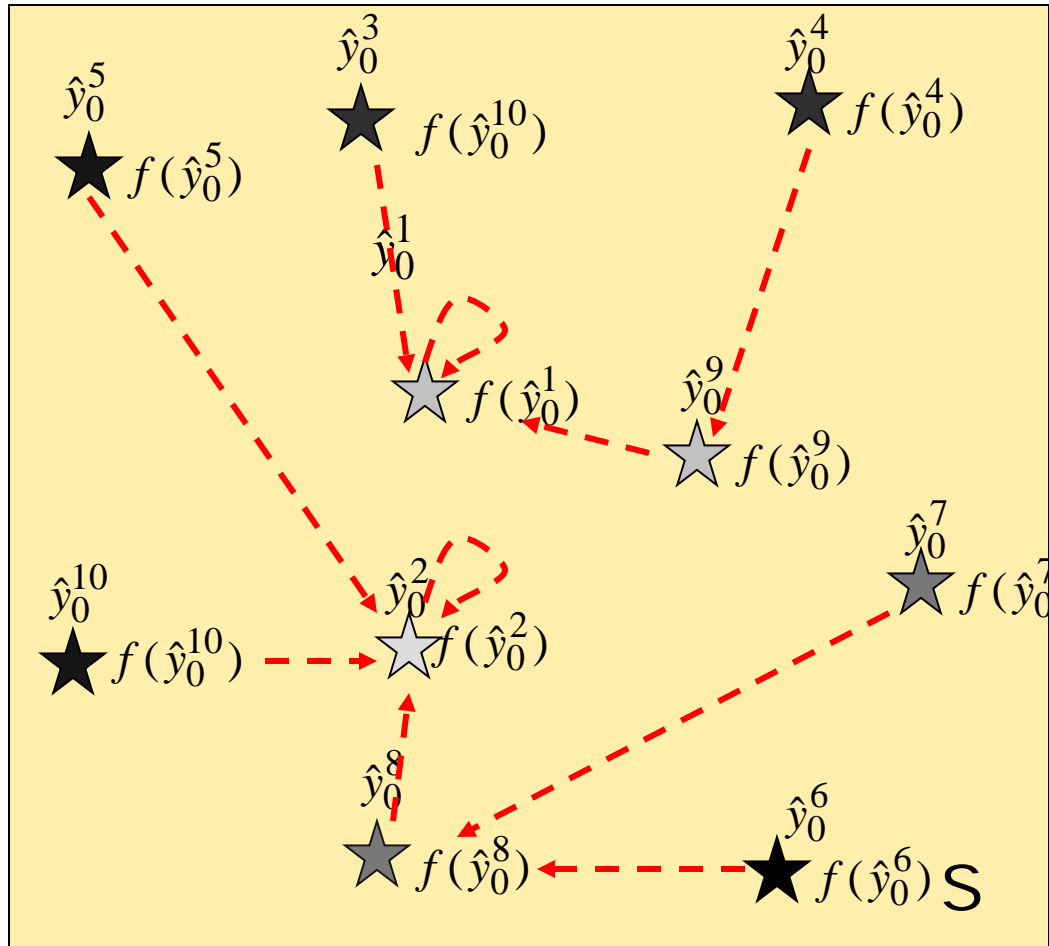**Initialization:** Sample N=10 points uniformly on S

**N-particle exploration:**

Move particles from ⬤ to

⭐ using Markov kernel E

Evaluate $f(.)$

Lower $f(.)$ ⟷ Higher $f(.)$

**N-particle selection:**

Move particles from ⭐ to

⬤ according to their objective function values

# Illustration of Interacting-Particle Algorithm



**Initialization:** Sample N=10 points uniformly on S

**N-particle exploration:**

Move particles from ⬤ to

⭐ using Markov kernel E

Evaluate $f(.)$

Lower $f(.)$ ⟵⟶ Higher $f(.)$
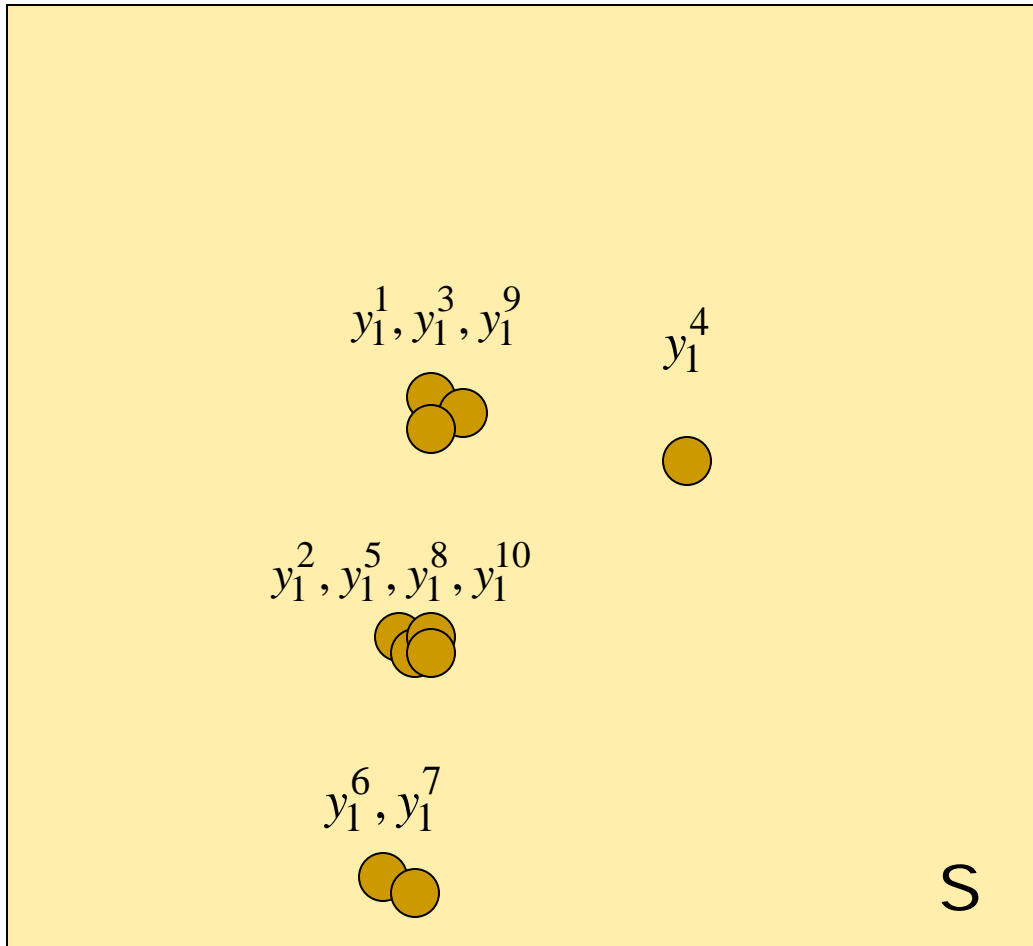
**N-particle selection:**

Move particles from ⭐ to

⬤ according to their objective function values

$y_1^1, y_1^3, y_1^9$

$y_1^4$

$y_1^2, y_1^5, y_1^8, y_1^{10}$

$y_1^6, y_1^7$

S

# Multi-start or Population-based Algorithms

o **Multi-start and clustering algorithms**
[Rinnooy Kan and Timmer, 1987]   [Locatelli and Schoen, 1999]

o **Genetic algorithms** [Davis, 1991]

o **Evolutionary programming** [Bäck, Fogel and Michalewicz, 1997]

o **Particle swarm optimization** [Kennedy, Eberhart and Shi, 2001]

o **Interacting particle algorithm** [del Moral, 2004]

INDUSTRIAL
ENGINEERING