



# Non-Parametric Up-and-Down Experimentation Revisited<sup>1</sup>

Michael N. Katehakis

*Rutgers Business School*

18TH ANNUAL APPLIED PROBABILITY DAY:  
IN HONOR AND MEMORY OF CYRUS DERMAN

THE CENTER FOR APPLIED PROBABILITY AT COLUMBIA UNIVERSITY  
December 2, 2011

---

<sup>1</sup>Joint work with Cyrus Derman

# The Basic Problem

Derman C. (1957). Non-Parametric Up-and-down Experimentation *The Annals of Mathematical Statistics* , 28(3), pp. 795-798.

- Let  $Y(x)$  be a random variable such that:

$$Y(x) = \begin{cases} 1 & P(Y(x) = 1) = F(x) \\ 0 & P(Y(x) = 0) = 1 - F(x) \end{cases}$$

where  $F(x)$  is an **unknown** distribution function.

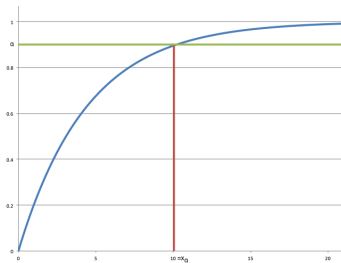
- **Objective:** given  $\alpha$  estimate the  $\alpha$ - quantile of  $F(x)$ ,

$$x_\alpha = F^{-1}(\alpha)$$

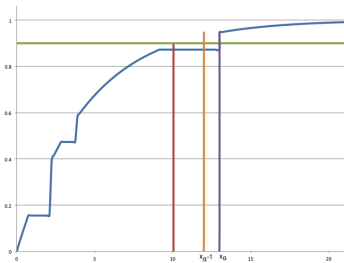
with observations distributed like  $Y(x)$  where the choice of  $x$  is under control.

- **Special Case:** Median:  $x_{0.50} = L_{0.50}$ ,  $\alpha = .50$ .

# Possible Cases



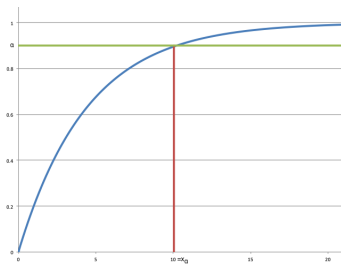
$$F(x_\alpha) = \alpha$$



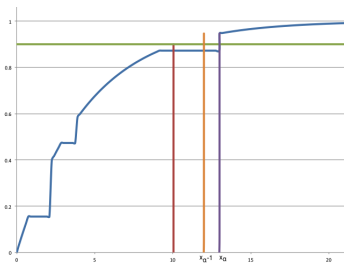
$$x_\alpha = \inf\{x \in \mathbf{N} : F(x) \geq \alpha\}$$

$$\underline{x}_\alpha = \sup\{x \in \mathbf{N} : F(x) \leq \alpha\}$$

# Possible Cases



$$F(x_\alpha) = \alpha$$



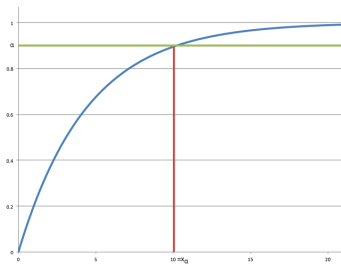
$$x_\alpha = \inf\{x \in \mathbb{N} : F(x) \geq \alpha\}$$

$$\underline{x}_\alpha = \sup\{x \in \mathbb{N} : F(x) \leq \alpha\}$$

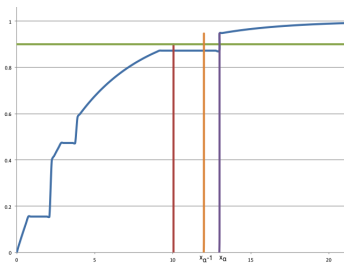
## Objective:

Start with  $x_0$ , observe  $Y(x_0) = 1, 0$ ,  $\Pr(Y(x_0) = 1) = F(x_0)$ .  
choose  $X_1$ , observe  $Y(X_1) = 1, 0$ ,  $\Pr(Y(X_1) = 1) = F(X_1)$ .

# Possible Cases



$$F(x_{\alpha}) = \alpha$$



$$x_{\alpha} = \inf\{x \in \mathbf{N} : F(x) \geq \alpha\}$$

$$\underline{x}_{\alpha} = \sup\{x \in \mathbf{N} : F(x) < \alpha\}$$

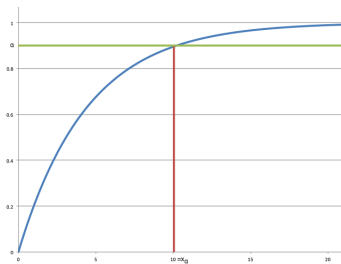
## Objective:

Start with  $x_0$ , observe  $Y(x_0) = 1, 0$ ,  $\Pr(Y(x_0) = 1) = F(x_0)$ .

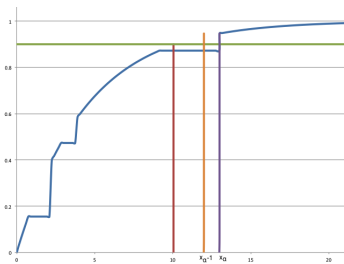
choose  $X_1$ , observe  $Y(X_1) = 1, 0$ ,  $\Pr(Y(X_1) = 1) = F(X_1)$ .

$\dots$ ,  $\dots$

# Possible Cases



$$F(x_\alpha) = \alpha$$



$$x_\alpha = \inf\{x \in \mathbb{N} : F(x) \geq \alpha\}$$

$$\underline{x}_\alpha = \sup\{x \in \mathbb{N} : F(x) \leq \alpha\}$$

## Objective:

Start with  $x_0$ , observe  $Y(x_0) = 1, 0$ ,  $\Pr(Y(x_0) = 1) = F(x_0)$ .

choose  $X_1$ , observe  $Y(X_1) = 1, 0$ ,  $\Pr(Y(X_1) = 1) = F(X_1)$ .

$\dots$ ,  $\dots$

choose  $X_n$ , such that  $\hat{x}_\alpha(x_0, X_1, \dots, X_n) \xrightarrow{a.s.} x_\alpha$ .

## News vendor Model

- When  $x$  items are ordered in a period we observe if there is :  
no shortage  $Y(x) = 1$  or shortage  $Y(x) = 0$ .
- The optimal order quantity  $x_\alpha$  is determined by a quantile requirement:  $x_\alpha = \inf\{x \in \mathbf{N} : F(x) \geq \alpha = (p - c)/p\}$
- $F$  is the demand distribution.

## News vendor Model

- When  $x$  items are ordered in a period we observe if there is :  
no shortage  $Y(x) = 1$  or shortage  $Y(x) = 0$ .
- The optimal order quantity  $x_\alpha$  is determined by a quantile requirement:  $x_\alpha = \inf\{x \in \mathbf{N} : F(x) \geq \alpha = (p - c)/p\}$
- $F$  is the demand distribution.

## Drug Testing

- Dosage is  $x$
- Observe a success  $Y(x) = 1$  or failure  $Y(x) = 0$ .
- The optimal dosage  $x_\alpha$  is determined by a quantile requirement:  
 $x_\alpha = \inf\{x \in \mathbf{N} : F(x) \geq \alpha\}$
- $F$  is the quantal response function.



## News vendor Model

- When  $x$  items are ordered in a period we observe if there is :  
no shortage  $Y(x) = 1$  or shortage  $Y(x) = 0$ .
- The optimal order quantity  $x_\alpha$  is determined by a quantile requirement:  $x_\alpha = \inf\{x \in \mathbf{N} : F(x) \geq \alpha = (p - c)/p\}$
- $F$  is the demand distribution.

## Drug Testing

- Dosage is  $x$
- Observe a success  $Y(x) = 1$  or failure  $Y(x) = 0$ .
- The optimal dosage  $x_\alpha$  is determined by a quantile requirement:  
 $x_\alpha = \inf\{x \in \mathbf{N} : F(x) \geq \alpha\}$
- $F$  is the quantal response function.

## Educational Testing

- Difficulty level of a test question is  $x$
- Observe a correct  $Y(x) = 1$  or wrong  $Y(x) = 0$  answer.
- The student's ability is the  $x_\alpha$  for which  $F(x_\alpha) = \alpha$

## Newsvendor Model

- When  $x$  items are ordered in a period we observe if there is :  
no shortage  $Y(x) = 1$  or shortage  $Y(x) = 0$ .
- The optimal order quantity  $x_\alpha$  is determined by a quantile requirement:  $x_\alpha = \inf\{x \in \mathbf{N} : F(x) \geq \alpha = (p - c)/p\}$
- $F$  is the demand distribution.

## Drug Testing

- Dosage is  $x$
- Observe a success  $Y(x) = 1$  or failure  $Y(x) = 0$ .
- The optimal dosage  $x_\alpha$  is determined by a quantile requirement:  
 $x_\alpha = \inf\{x \in \mathbf{N} : F(x) \geq \alpha\}$
- $F$  is the quantal response function.

## Educational Testing

- Difficulty level of a test question is  $x$
- Observe a correct  $Y(x) = 1$  or wrong  $Y(x) = 0$  answer.
- The student's ability is the  $x_\alpha$  for which  $F(x_\alpha) = \alpha$

## Manufacturing

## News vendor Model

- When  $x$  items are ordered in a period we observe if there is :  
no shortage  $Y(x) = 1$  or shortage  $Y(x) = 0$ .
- The optimal order quantity  $x_\alpha$  is determined by a quantile requirement:  $x_\alpha = \inf\{x \in \mathbf{N} : F(x) \geq \alpha = (p - c)/p\}$
- $F$  is the demand distribution.

## Drug Testing

- Dosage is  $x$
- Observe a success  $Y(x) = 1$  or failure  $Y(x) = 0$ .
- The optimal dosage  $x_\alpha$  is determined by a quantile requirement:  
 $x_\alpha = \inf\{x \in \mathbf{N} : F(x) \geq \alpha\}$
- $F$  is the quantal response function.

## Educational Testing

- Difficulty level of a test question is  $x$
- Observe a correct  $Y(x) = 1$  or wrong  $Y(x) = 0$  answer.
- The student's ability is the  $x_\alpha$  for which  $F(x_\alpha) = \alpha$

## Manufacturing

## The Sensitivity of an Explosive

# Derman's Up & Down Method

- Grid or experimental range of  $x$  to a set of numbers of the form

$$b + hn \quad (-\infty < b < \infty, h > 0, n = 0, \pm 1, \dots).$$

For convenience one can assume  $b = 0, h = 1$ .

- Procedure:

Start with  $x_0$  (init. guess)  $y_0 = Y(x_0)$  is observed  
where  $P(Y(x_0) = 1) = F(x_0) = 1 - P(Y(x_0) = 0)$ .

Given  $x_0, y(x_0), \dots, x_{n-1}, y(x_{n-1})$  for  $n \geq 0$ , define:

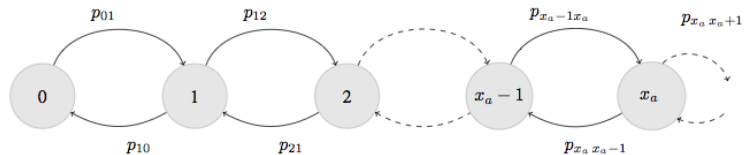
$$x_{n+1} = \begin{cases} x_n + 1 & \text{if } \begin{cases} y(x_n) = 0 & \text{with probability: } 1 \\ y(x_n) = 1 & \text{with probability: } 1 - \frac{1}{2\alpha} \end{cases} \\ x_n - 1 & \text{if } y(x_n) = 1 \quad \text{with probability: } \frac{1}{2\alpha} \end{cases}$$

where w.l.o.g.  $\alpha > 1/2$

- The estimate  $\hat{x}_{\alpha,n}$  of  $x_\alpha$  based on  $n$  observations is

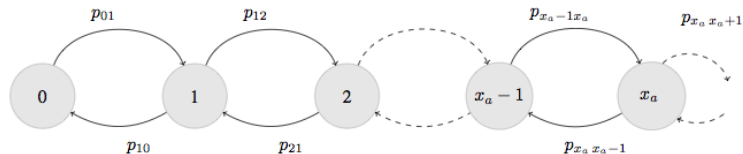
$$\hat{x}_{\alpha,n} = \begin{cases} \text{the most frequent value of } x\text{'s,} & \text{if unique,} \\ \text{the arithmetic average of the most frequent levels} & \text{otherwise.} \end{cases}$$

# Derman's Main Result



$\{X_n, \}_{n \geq 1}$  process:  $p_{i, i+1} = 1 - F(i)/2\alpha$

# Derman's Main Result

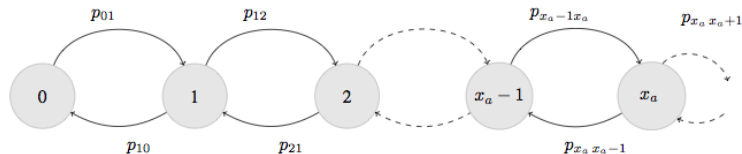


$\{X_n, \}_{n \geq 1}$  process:  $p_{i, i+1} = 1 - F(i)/2\alpha$

Condition A: If  $F(x)$  is strictly increasing for  $x \in [x_\alpha - 1, x_\alpha]$  then:

$$\Pr(\max(|\overline{\lim}_n \hat{x}_{\alpha, n} - x_\alpha|, |\underline{\lim}_n \hat{x}_{\alpha, n} - x_\alpha|) < 1) = 1$$

# Derman's Main Result



$$\{X_n, \}_{n \geq 1} \text{ process: } p_{i, i+1} = 1 - F(i)/2\alpha$$

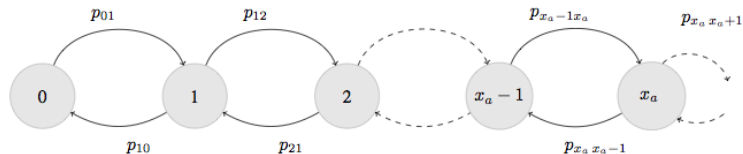
Condition A: If  $F(x)$  is strictly increasing for  $x \in [x_\alpha - 1, x_\alpha]$  then:

$$\Pr(\max(|\overline{\lim}_n \hat{x}_{\alpha, n} - x_\alpha|, |\underline{\lim}_n \hat{x}_{\alpha, n} - x_\alpha|) < 1) = 1$$

Main Tool:  $\pi(x) = \lim_n \Pr(X_n = x | x_0)$

$$\pi(0) \leq \pi(1) \leq \dots \leq \pi(x_\alpha - 1) < \pi(x_\alpha) \geq \pi(x_\alpha + 1) \geq \pi(x_\alpha + 2) \geq \dots$$

# Derman's U-D Revisited

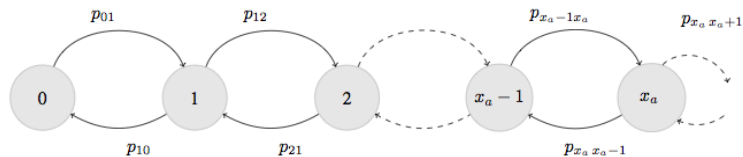


$\{X_n, \}_{n \geq 1}$  process:  $p_{i, i+1} = 1 - F(i)/2\alpha$

- What is an efficient *Stopping Criterion* ?
- What is the *Error Probability* ?
- What if *Condition A* does not hold ?

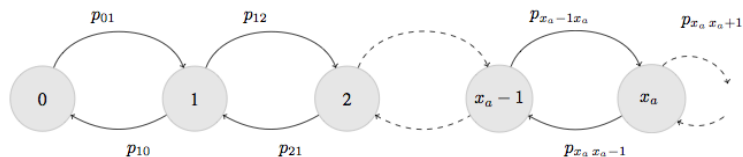


# Derman's U-D Revisited - Answers



$\{X_n, \}_{n \geq 1}$  process:  $p_{i, i+1} = 1 - F(i)/2\alpha$

# Derman's U-D Revisited - Answers

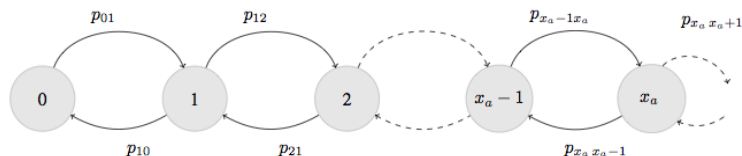


$$\{X_n, \}_{n \geq 1} \text{ process: } p_{i, i+1} = 1 - F(i)/2\alpha$$

Observation 1 if  $x_\alpha$  is on the grid

$$p_{00} = 1 \geq p_{12} \geq \dots \geq p_{x_\alpha - 1, x_\alpha} \geq p_{x_\alpha, x_\alpha + 1} = 1/2 \geq p_{x_\alpha + 1, x_\alpha + 2} \geq \dots$$

# Derman's U-D Revisited - Answers



$$\{X_n, \}_{n \geq 1} \text{ process: } p_{i, i+1} = 1 - F(i)/2\alpha$$

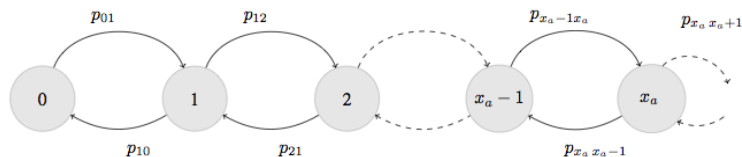
Observation 1 if  $x_\alpha$  is on the grid

$$p_{00} = 1 \geq p_{12} \geq \dots \geq p_{x_\alpha - 1, x_\alpha} \geq p_{x_\alpha, x_\alpha + 1} = 1/2 \geq p_{x_\alpha + 1, x_\alpha + 2} \geq \dots$$

Observation 2 if  $x_\alpha$  is not on the grid:

$$p_{00} = 1 \geq p_{12} \geq \dots \geq p_{x_\alpha - 1, x_\alpha} > 1/2 \geq p_{x_\alpha, x_\alpha + 1} \geq p_{x_\alpha + 1, x_\alpha + 2} \geq \dots$$

# Derman's U-D Revisited - Answers



$$\{X_n, \}_{n \geq 1} \text{ process: } p_{i, i+1} = 1 - F(i)/2\alpha$$

Observation 1 if  $x_\alpha$  is on the grid

$$p_{00} = 1 \geq p_{12} \geq \dots \geq p_{x_\alpha - 1, x_\alpha} \geq p_{x_\alpha, x_\alpha + 1} = 1/2 \geq p_{x_\alpha + 1, x_\alpha + 2} \geq \dots$$

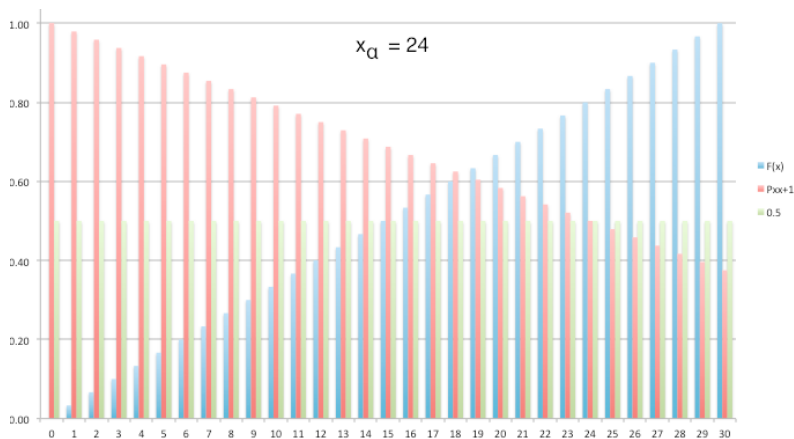
Observation 2 if  $x_\alpha$  is not on the grid:

$$p_{00} = 1 \geq p_{12} \geq \dots \geq p_{x_\alpha - 1, x_\alpha} > 1/2 \geq p_{x_\alpha, x_\alpha + 1} \geq p_{x_\alpha + 1, x_\alpha + 2} \geq \dots$$

Compare with :

$$\pi(0) \leq \pi(1) \leq \dots \leq \pi(x_\alpha - 1) \leq \pi(x_\alpha) \geq \pi(x_\alpha + 1) \geq \pi(x_\alpha + 2) \geq \dots$$

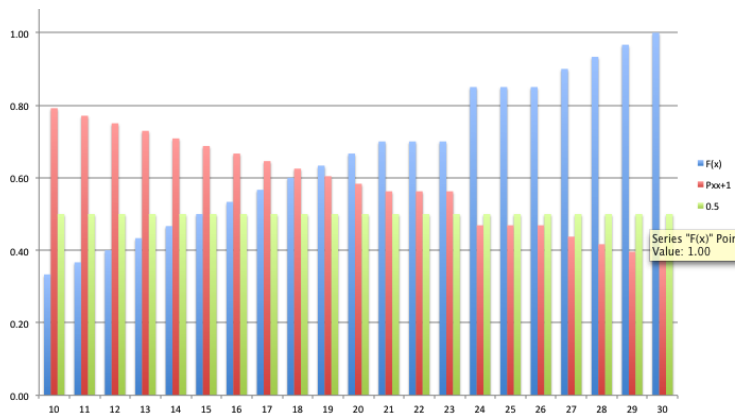
# Derman's U-D Revisited -Continued



$\alpha = .80$ ,  $x_\alpha = 24$  **on the grid**:  $F(24) = .8$

$$p_{23,24} > 1/2 = p_{24,25} > p_{25,26} < 1/2$$

## Derman's U-D Revisited -Continued



$\alpha = .80$ ,  $x_\alpha = 24$  **not on the grid**:  $F(23) < 0.8 < F(24)$

$$p_{23,24} > 1/2 > p_{24,25} > p_{25,26} > \dots$$

# Derman and Katehakis (2010-2011) New Results

For simplicity consider the case  $x_\alpha$  is on the grid.

- We can **Stop** the procedure:
  - $\tau_\alpha^1 = \inf\{k : \hat{p}_{k,k+1}(x_0, \dots, X_k) \in (1/2 - \epsilon, 1/2 + \epsilon)\}$
  - $\tau_\alpha^2 = \inf\{k : \hat{p}_{k,k+1}(x_0, \dots, X_k) \in (1/2 - \epsilon, 1/2 + \epsilon) \ \& \ V_k = \max_{k'}\{V_{k'}\} \}$
  - $\hat{x}_a = X_{\tau_\alpha^i}$
- Have  $\Pr(\tau_\alpha^i > u) = c_F^i(e^{-u} + \epsilon_F^i(n, u))$   
where  $\epsilon_F^i(n, u) \rightarrow 0$  ( $n \rightarrow \infty$ )
- We can modify the procedure by taking a second sample at  $\tau_\alpha$   
(Two Stage)
- Working on using techniques of Adaptive MDPs to obtain:

$$R_N^\pi \geq R_N^{\pi^{DK}} = M_{DK}(P)\log N + o(\log N)$$

# Derman and Katehakis (2010-2011) New Results

For simplicity consider the case  $x_\alpha$  is on the grid.

- We can **Stop** the procedure:
  - $\tau_\alpha^1 = \inf\{k : \hat{p}_{k,k+1}(x_0, \dots, X_k) \in (1/2 - \epsilon, 1/2 + \epsilon)\}$
  - $\tau_\alpha^2 = \inf\{k : \hat{p}_{k,k+1}(x_0, \dots, X_k) \in (1/2 - \epsilon, 1/2 + \epsilon) \ \& \ V_k = \max_{k'}\{V_{k'}\} \}$
  - $\hat{x}_a = X_{\tau_\alpha^i}$
- Have  $\Pr(\tau_\alpha^i > u) = c_F^i(e^{-u} + \epsilon_F^i(n, u))$   
where  $\epsilon_F^i(n, u) \rightarrow 0$  ( $n \rightarrow \infty$ )
- We can modify the procedure by taking a second sample at  $\tau_\alpha$   
(Two Stage)
- Working on using techniques of Adaptive MDPs to obtain:

$$R_N^\pi \geq R_N^{\pi^{DK}} = M_{DK}(P)\log N + o(\log N)$$

Similar results hold in the case  $x_\alpha$  is not on the grid or  
Condition A does not hold.



# Background

## Non Parametric “Up & Down” Methods

- Anderson, T., McCarthy, P., Tukey, J., (1946). Staircase method of sensitivity testing. *Naval Ordinance Report 65-46*, Statistical Research Group, Princeton University, Princeton, NJ.
- Dixon, W.J. and Mood, A.M. (1948). A method for obtaining and analyzing sensitive data. *Journal of the American Statistical Association* 43, pp. 109-126. *The validity of their procedure depends on the assumption that  $F(x)$  is normal.*
- Derman C. (1957). Non-Parametric Up-and-down Experimentation
- Wetherill, G.B. (1963). Sequential estimation of quantal response curves.

# Background

## Non Parametric “Up & Down” Methods

- Anderson, T., McCarthy, P., Tukey, J., (1946). Staircase method of sensitivity testing. *Naval Ordnance Report 65-46*, Statistical Research Group, Princeton University, Princeton, NJ.
- Dixon, W.J. and Mood, A.M. (1948). A method for obtaining and analyzing sensitive data. *Journal of the American Statistical Association* 43, pp. 109 -126. *The validity of their procedure depends on the assumption that  $F(x)$  is normal.*
- Derman C. (1957). Non-Parametric Up-and-down Experimentation
- Wetherill, G.B. (1963). Sequential estimation of quantal response curves.
- Wetherill, G.B., Glazebrook, K.D., (1986). *Sequential Methods in Statistics*, Chapman & Hall, London.
- Durham, D.S., Flournoy N. and Rosenberger W. F (1997). A Random Walk Rule for Phase I Clinical Trials
- Bortot P. and Giovagnolia A. (2005). Up-and-down experiments of first and second order
- Ivanova A(2006). Dose-Finding in Oncology-Nonparametric Methods,
- Pollak, R., Palazotto, A., Nicholas, T. (2006). A simulation-based investigation of the stair- case method for fatigue strength testing.
- Baldi Antognini A. and . Giovagnoli A (2010). Compound Optimal Allocation for Individual and Collective Ethics in Binary Clinical Trials,

## Background: “Stochastic Approximation” and Related Methods

- Robbins, H. and Monro, S. (1951). A stochastic approximation method, *The Annals of Mathematical Statistics*, 22(3), pp. 400–407.
- Derman C. and J. Sacks (1959). On Dvoretzky's Stochastic Approximation Theorem, *The Annals of Mathematical Statistics*, 30(2), pp. 601-606.
- Frederic M. Lord (1971) Tailored Testing, An Application of Stochastic Approximation *Journal of the American Statistical Association* Vol. 66, No. 336, pp. 707-711.

*The Robbins and Monro SA scheme can be used for estimating any quantile and it imposes no parametric assumptions on  $F(x)$ .*

- The method does assume, however, that the range of possible experimental values of  $x$  is the real line.
- It has slow convergence rate.

## Background: Optimal Adaptive Policies for MABs

- Gittins, J.C. and Jones, D.M. (1979). "A dynamic allocation index for the discounted multiarmed bandit problem",
- Katehakis M. N. and A. F. Veinott Jr. (1987). "The Multi-Armed Bandit problem: decomposition and computation",

# Background: Optimal Adaptive Policies for MABs

- Gittins, J.C. and Jones, D.M. (1979). “A dynamic allocation index for the discounted multiarmed bandit problem”,
- Katehakis M. N. and A. F. Veinott Jr. (1987). “The Multi-Armed Bandit problem: decomposition and computation”,
- Lai T.L. and H. E. Robbins (1985) “Asymptotically Efficient Adaptive Allocation Rules”,

$$R_N^\pi = \max\{\mu_1, \dots, \mu_n\}N - V_N^\pi$$

$$R_N^\pi \geq R_N^{\pi^{LR}} = M_{LR}(P)\log N + o(\log N)$$

- Note:

$$R_N^{\pi^0} = M_{LR}(P)\log N + o(\log N) = o(N^a) \quad \forall a > 0$$

- Even better:

$$\overline{\lim}_N R_N^{\pi^0} / R_N^\pi \leq 1$$

- Katehakis M. N. and H. E. Robbins (1995). “Sequential choice from several populations”,

# Background: Optimal Adaptive Policies for MDPs

- Burnetas, A.N. and M. N. Katehakis (1996) "Optimal Adaptive Policies for Sequential Allocation Problems",
- Burnetas, A.N. and M. N. Katehakis (1997) "Optimal Adaptive Policies for Markov Decision Processes",

$$R_N^\pi = V_N - V_N^\pi$$
$$R_N^\pi \geq R_N^{\pi^{BK}} = M_{BK}(P)\log N + o(\log N)$$

# Background: Optimal Adaptive Policies for MDPs

- Burnetas, A.N. and M. N. Katehakis (1996) "Optimal Adaptive Policies for Sequential Allocation Problems",
- Burnetas, A.N. and M. N. Katehakis (1997) "Optimal Adaptive Policies for Markov Decision Processes",

$$R_N^\pi = V_N - V_N^\pi$$
$$R_N^\pi \geq R_N^{\pi^{BK}} = M_{BK}(P)\log N + o(\log N)$$

## Near-optimal Regret Bounds for Reinforcement Learning

- Auer, P. and R. Ortner (2007) Logarithmic online regret bounds for undiscounted reinforcement learning.

$$R_N^\pi \geq R_N^{\pi^{AO}} = M_{AO}(P)\log N + o(\log N)$$

- Tewari A. and P. Bartlett (2008). "Optimistic linear programming gives logarithmic regret for irreducible MDPs",

$$R_N^\pi \geq R_N^{\pi^{TB}} = M_{TB}(P)\log N + o(\log N)$$

- Auer, P. and Jaksch, T. and Ortner, R. (2009). "Near-optimal regret bounds for reinforcement learning",

# Background: Optimal Adaptive Policies for MDPs

- Burnetas, A.N. and M. N. Katehakis (1996) "Optimal Adaptive Policies for Sequential Allocation Problems",
- Burnetas, A.N. and M. N. Katehakis (1997) "Optimal Adaptive Policies for Markov Decision Processes",

$$R_N^\pi = V_N - V_N^\pi$$
$$R_N^\pi \geq R_N^{\pi^{BK}} = M_{BK}(P)\log N + o(\log N)$$

## Near-optimal Regret Bounds for Reinforcement Learning

- Auer, P. and R. Ortner (2007) Logarithmic online regret bounds for undiscounted reinforcement learning.

$$R_N^\pi \geq R_N^{\pi^{AO}} = M_{AO}(P)\log N + o(\log N)$$

- Tewari A. and P. Bartlett (2008). "Optimistic linear programming gives logarithmic regret for irreducible MDPs",

$$R_N^\pi \geq R_N^{\pi^{TB}} = M_{TB}(P)\log N + o(\log N)$$

- Auer, P. and Jaksch, T. and Ortner, R. (2009). "Near-optimal regret bounds for reinforcement learning",

$$M_{AO}(P) > M_{TB}(P) > M_{BK}(P) \forall P$$



# Background: Optimal Adaptive Policies for MDPs

- Burnetas, A.N. and M. N. Katehakis (1996) "Optimal Adaptive Policies for Sequential Allocation Problems",
- Burnetas, A.N. and M. N. Katehakis (1997) "Optimal Adaptive Policies for Markov Decision Processes",

$$R_N^\pi = V_N - V_N^\pi$$
$$R_N^\pi \geq R_N^{\pi^{BK}} = M_{BK}(P)\log N + o(\log N)$$

## Near-optimal Regret Bounds for Reinforcement Learning

- Auer, P. and R. Ortner (2007) Logarithmic online regret bounds for undiscounted reinforcement learning.

$$R_N^\pi \geq R_N^{\pi^{AO}} = M_{AO}(P)\log N + o(\log N)$$

- Tewari A. and P. Bartlett (2008). "Optimistic linear programming gives logarithmic regret for irreducible MDPs",

$$R_N^\pi \geq R_N^{\pi^{TB}} = M_{TB}(P)\log N + o(\log N)$$

- Auer, P. and Jaksch, T. and Ortner, R. (2009). "Near-optimal regret bounds for reinforcement learning",

$$M_{AO}(P) > M_{TB}(P) > M_{BK}(P) \forall P$$

# Background: Optimal Adaptive Policies for MDPs

- Burnetas, A.N. and M. N. Katehakis (1996) "Optimal Adaptive Policies for Sequential Allocation Problems",
- Burnetas, A.N. and M. N. Katehakis (1997) "Optimal Adaptive Policies for Markov Decision Processes",

$$R_N^\pi = V_N - V_N^\pi$$
$$R_N^\pi \geq R_N^{\pi^{BK}} = M_{BK}(P)\log N + o(\log N)$$

## Near-optimal Regret Bounds for Reinforcement Learning

- Auer, P. and R. Ortner (2007) Logarithmic online regret bounds for undiscounted reinforcement learning.

$$R_N^\pi \geq R_N^{\pi^{AO}} = M_{AO}(P)\log N + o(\log N)$$

- Tewari A. and P. Bartlett (2008). "Optimistic linear programming gives logarithmic regret for irreducible MDPs",

$$R_N^\pi \geq R_N^{\pi^{TB}} = M_{TB}(P)\log N + o(\log N)$$

- Auer, P. and Jaksch, T. and Ortner, R. (2009). "Near-optimal regret bounds for reinforcement learning",

$$M_{AO}(P) > M_{TB}(P) > M_{BK}(P) \forall P$$

# The Robbins Monro Method for $x_\alpha$

Compute  $x_\alpha$  such that

$$M(x_\alpha) = \alpha$$

- Deterministic case  $M(x)$  is **known** solution by:

$$x_n = x_{n-1} + a_n[\alpha - M(x_{n-1})],$$

- Stochastic case  $M(x)$  is **unknown**

$$E(Y(x)) = M(x) = \alpha. \quad \forall x$$

Solution by:

$$x_n = x_{n-1} + a_n[\alpha - y_{n-1}],$$

The estimate  $\hat{x}_{\alpha,n}$  of  $x_\alpha$  based on  $n$  observations is  $\hat{x}_{\alpha,n} = x_n$ ,  
*Under regularity conditions*<sup>2</sup>

---

<sup>2</sup>E.g.,  $M$  is non-decreasing, there exists a solution  $M(x_\alpha) = \alpha$ ,  $\exists \frac{M(x)}{dx} > 0$  at  $x_\alpha$ , and  $\sum_{n=0}^{\infty} a_n = \infty$ ,  $\sum_{n=0}^{\infty} a_n^2 < \infty$

# The Dixon and Mood Method for $LD_{50}$ ( $x_{0.50}$ )

- Grid or experimental range<sup>3</sup> of  $x$  to a set of numbers of the form

$$b + hn \quad (-\infty < b < \infty, h > 0, n = 0, \pm 1, \dots).$$

For convenience one can assume  $b = 0, h = 1$ .

- Data: Treatments are administered sequentially at dosage:  $X_i$ , as follows:  
Start with  $x_0$  (arbitrary guess)  $y_0 = Y(x_0)$  is observed where  
 $P(Y(x_0) = 1) = F(x_0) = 1 - P(Y(x_0) = 0)$ .  
Given  $x_0, y(x_0), \dots, x_k, y(x_{n-1})$  for  $n \geq 0$ , define recursively

$$x_n = \begin{cases} x_{n-1} + 1 & \text{if } y(x_k) = 0, \\ x_{n-1} - 1 & \text{if } y(x_k) = 1. \end{cases}$$

- The estimate  $\hat{x}_\alpha$  of  $x_\alpha$  based on  $n$  observations is

$$\hat{x}_{0.50} = \frac{1}{n+1} \sum_{j=1}^{n+1} x_j,$$

---

<sup>3</sup>Natural limitations such as when  $x$  is obtained by a counting procedure - limitations on the precision of measuring instruments.