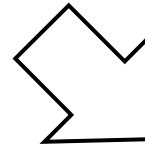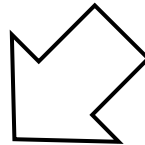# 時間整合的マルコフ決定過程

## 恐神 貴行

IBM東京基礎研究所

森村哲郎, Mikael Onsjoe（IBM東京基礎研究所）との共同研究

# Today's route selection is limited

Shortest path algorithm
finds an optimal path

## For car navigation

Prescriptive analytics for drivers

## For traffic simulation
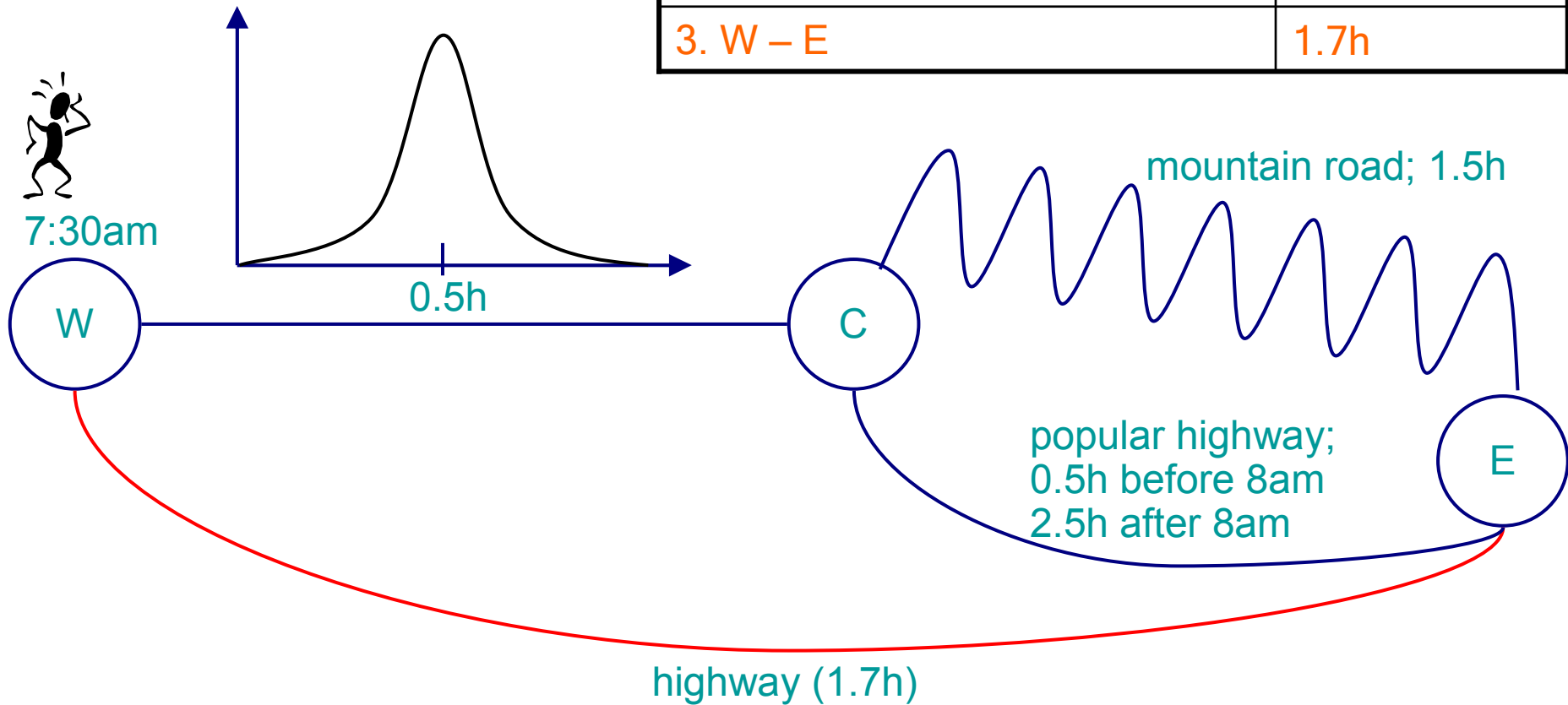
Descriptive model of drivers

cf. sophisticated drivers select routes dynamically depending on latest conditions

# Outline

- **Limitations of traditional route selection**

- Difficulties in selecting objective functions

- Time-consistent Markov decision processes

- Effectiveness of time-consistent Markov decision processes

# Consider an example with three paths

| Path | Expected time |
|------|---------------|
| 1. W – C – mountain road – E | 2h |
| 2. W – C – popular highway – E | 2h |
| 3. W – E | 1.7h |

7:30am

0.5h

W

C

E

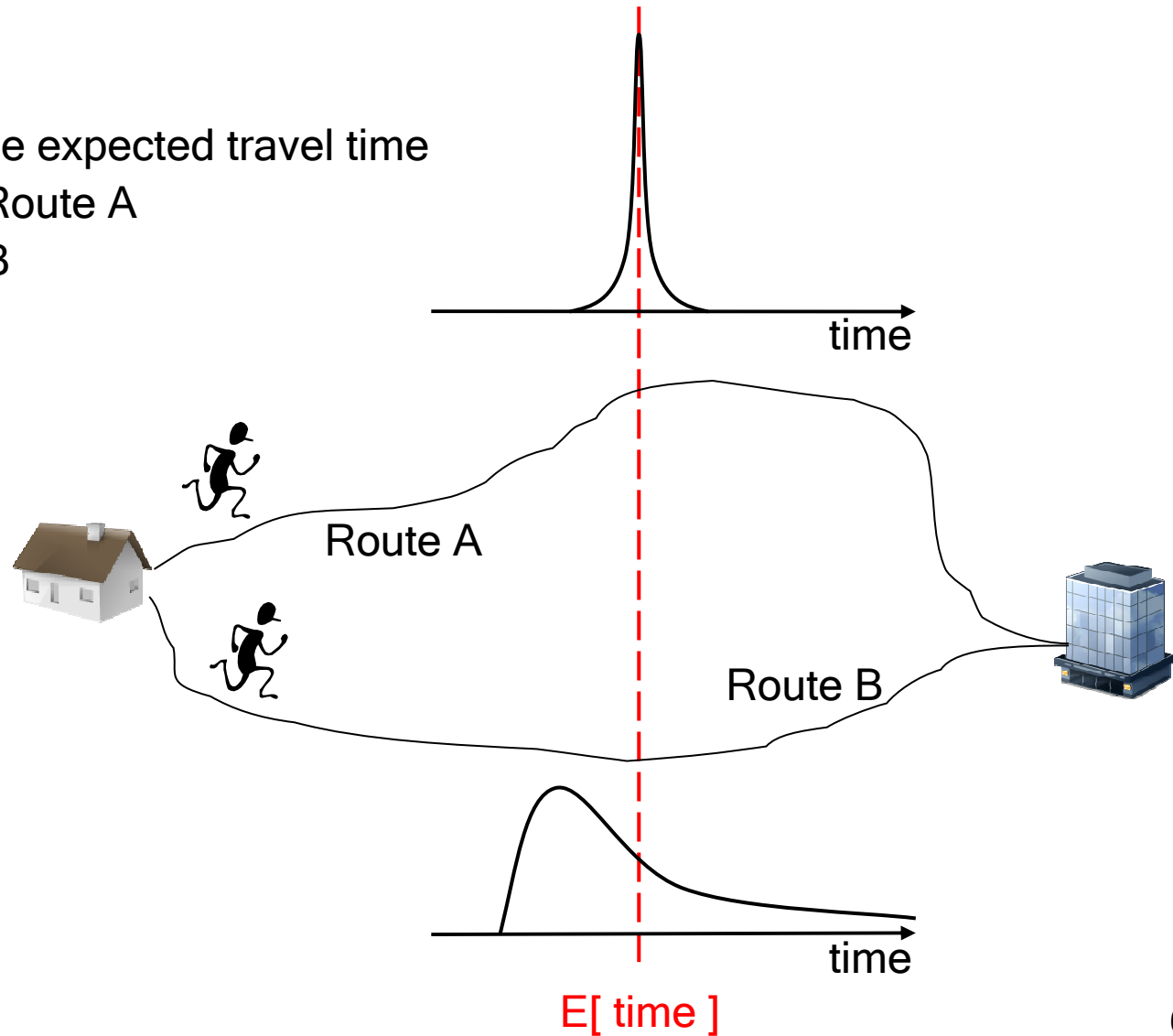mountain road; 1.5h

popular highway;
0.5h before 8am
2.5h after 8am

highway (1.7h)

# No path is better than a dynamic strategy

Expected time with the dynamic strategy is
0.5h + 0.5 x 0.5h + 0.5 x 1.5h = 1.5h

after 8am (50%)

mountain road; 1.5h

7:30am

W

0.5h

C

before 8am (50%)

popular highway;
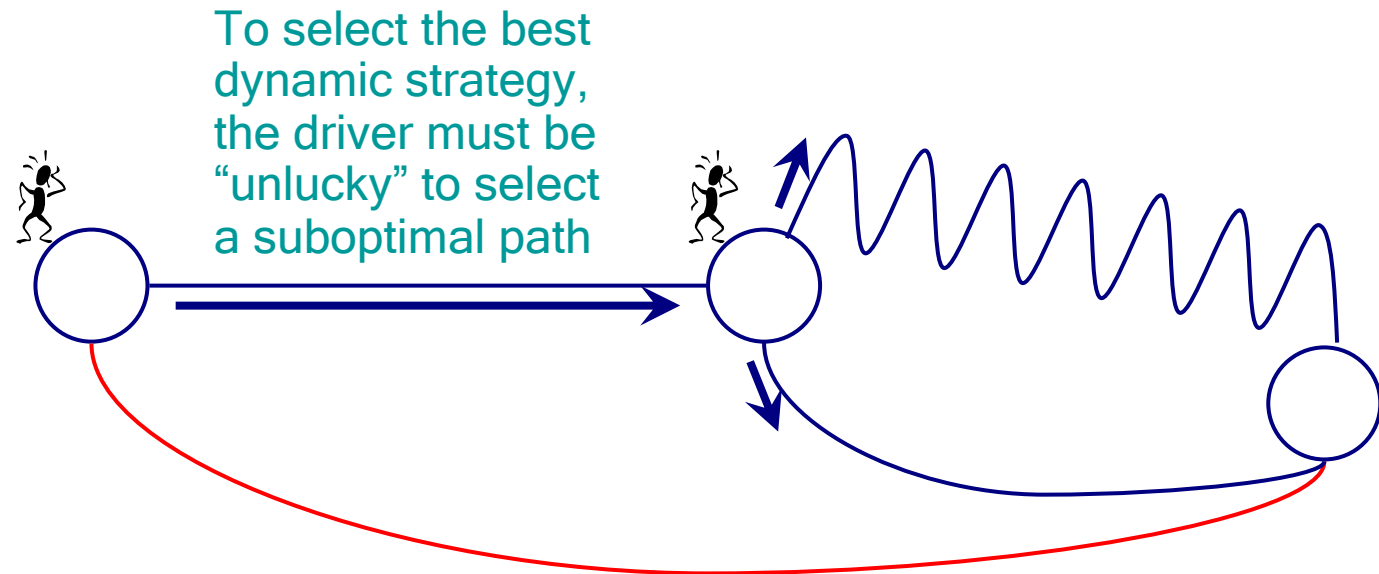0.5h before 8am
2.5h after 8am

E

highway (1.7h)

# Expectation is obviously limited in representing drivers' preference under risk

- Two routes have same expected travel time
- Some drivers prefer Route A
- Others prefer Route B



Route A

Route B

time

time

E[ time ]

# We study models for selecting dynamic strategies

- Interpretation of the dynamic strategy with a model of path selection is convoluted

To select the best dynamic strategy, the driver must be "unlucky" to select a suboptimal path



- Want to select optimal dynamic strategies with respect to a broad class of objective functions
  - personalized recommendation of dynamic strategies
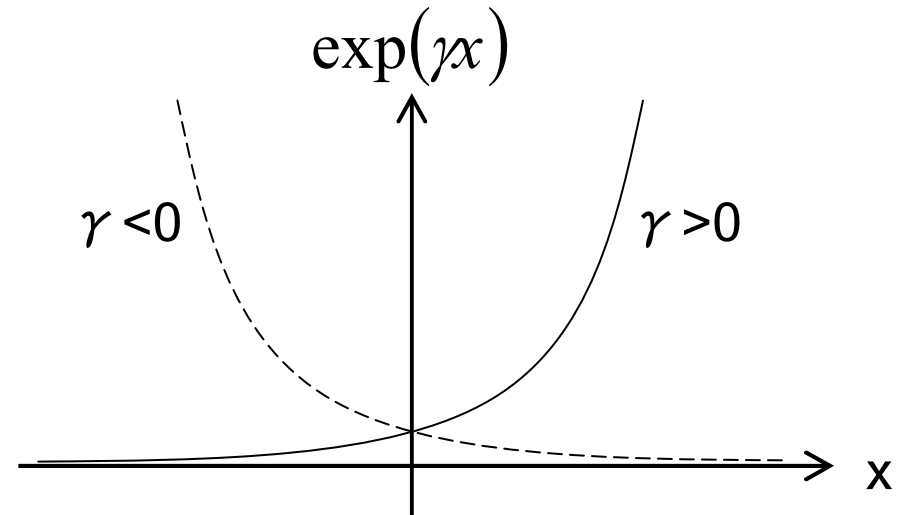  - realistic traffic simulation

# Outline

- Limitations of traditional route selection

- Difficulties in selecting objective functions

- Time-consistent Markov decision processes

- Effectiveness of time-consistent Markov decision processes

# Expected exponential utility is the standard objective of risk-sensitive Markov decision processes



$$\min. \ \mathrm{E}\!\left[\exp(\gamma X)\right] \quad \Leftrightarrow \quad \min. \ \mathrm{ERM}_\gamma[X] \equiv \frac{1}{\gamma}\ln\mathrm{E}\!\left[\exp(\gamma X)\right]$$

(works only for $\gamma > 0$)

- $\gamma > 0 \Rightarrow$ risk-averse
- $\gamma < 0 \Rightarrow$ risk-seeking

Minimization of *expected exponential utility* is essentially equivalent to minimization of *entropic risk measure*

# Which route would you take?



| Probability | 0.9 (normal) | 0.1 (busy) |
|---|---|---|
| Route P | 10 min. | Unif[20,80] min. |
| Route Q | Unif[0,20] min. | 50 min. |

Route P

Route Q

90%

10%

Travel time (min.)

90%

10%

Travel time (min.)

# Route P is never optimal with respect to any entropic risk measure

- Some drivers choose Route P
- Others choose Route Q
- They are all rational (e.g., E[P] = E[Q])

Route P

Route Q



$$\mathrm{ERM}_\gamma[T_Q] \le \mathrm{ERM}_\gamma[T_P], \forall \gamma$$

New

11

# Expected utility is the standard objective function for decision making under risk

- Choose a dynamic strategy such that

$$\mathrm{E}[u(T)]$$

is minimized
  - T: travel time
  - u: utility function

cf. Expected utility theory
(von Neumann & Morgenstern 1944)

- Entropic risk measure is a particular expected utility
  - u(x) = exp( $\gamma$ x)

# For every path, does there exist a utility such that the path is optimal with respect to the expected utility?

L

M

H

A

B

| | time | probability |
|---|---|---|
| $T_L$ | 20 | 1.0 |
| | 10 | 0.3 |
| $T_M$ | 20 | 0.5 |
| | 30 | 0.2 |
| | 10 | 0.6 |
| $T_H$ | 30 | 0.4 |

# Expected utility is limited in representing driver's preference



For any utility function, $u$:

$$\mathrm{E}[u(T_M)] = 0.5\,\mathrm{E}[u(T_L)] + 0.5\,\mathrm{E}[u(T_H)]$$

$\Rightarrow$ We can have only

$$\mathrm{E}[u(T_L)] \leq \mathrm{E}[u(T_M)] \leq \mathrm{E}[u(T_H)]$$

or

$$\mathrm{E}[u(T_H)] \leq \mathrm{E}[u(T_M)] \leq \mathrm{E}[u(T_L)]$$

Never choose M with expected utility

|       | time | probability |
|-------|------|-------------|
| $T_L$ | 20   | 1.0         |
|       | 10   | 0.3         |
| $T_M$ | 20   | 0.5         |
|       | 30   | 0.2         |
|       | 10   | 0.6         |
| $T_H$ | 30   | 0.4         |

# Conditional tail expectation is a popular risk measure in finance

- Choose a dynamic strategy such that

$$\text{CTE}_{\alpha}[T] \equiv \frac{(1-\beta)\text{E}[T \mid T > Q_{\alpha}] + (\beta - \alpha)Q_{\alpha}}{1-\alpha}$$

is minimized

  - T: travel time

$$Q_{\alpha} \equiv \min\{t \mid \text{Pr}(T \leq t) \geq \alpha\}$$

$$\beta \equiv \text{Pr}(T \leq V_{\alpha})$$

- When T is continuous,

$$\text{CTE}_{\alpha}[T] \equiv \text{E}[T \mid T > Q_{\alpha}]$$

CTE$_{0.75}$[T] = 1.5

25%

$Q_{0.25}$=1     3

15

# Choose either the road to B or that to B' when we leave A to reach C

Normal (90%):
60 min.

Busy (10%):
120 min.

10 min.

B

30 min.

A

B'

10 min.

C

Normal (90%):
50 min. w.p. 99%
100 min. w.p. 1%

Busy (10%):
100 min. w.p. 99%
200 min. w.p. 1%

B-C is normal iff B'-C is normal

# A-B'-C has smaller risk than A-B-C with respect to $CTE_{0.99}$



Normal (90%):
60 min.

Busy (10%):
120 min.

Start at 6am

10 min.

B

30 min.

B'

10 min.

C

Normal (90%):
50 min. w.p. 99%
100 min. w.p. 1%

Busy (10%):
100 min. w.p. 99%
200 min. w.p. 1%

|  | $CTE_{0.99}$ |
|---|---|
| A-B-C | 130 min. |
| A-B'-C | $\dfrac{210 \times 0.001 + 110 \times 0.009}{0.01}$ = 120 min. |

B-C is normal iff B'-C is normal

17

# If traffic conditions are normal at B', B'-B-C appears to have smaller risk than B'-C with respect to $CTE_{0.99}$



Normal:
60 min.

Start at 6am

10 min.

B

30 min.

A

C

B'

10 min.

Normal:
50 min. w.p. 99%
100 min. w.p. 1%

|  | $CTE_{0.99}$ |
|---|---|
| A-B'-C | 110 |
| A-B'-B-C | 100 |

B-C is normal iff B'-C is normal

# If traffic conditions are busy at B', B'-B-C appears to have smaller risk than B'-C with respect to $CTE_{0.99}$

Busy:
120 min.

Start at 6am

10 min.

B

30 min.

A

C

B'

10 min.

|  | $CTE_{0.99}$ |
|---|---|
| A-B'-C | 210 |
| A-B'-B-C | 160 |

Busy:
100 min. w.p. 99%
200 min. w.p. 1%

B-C is normal iff B'-C is normal

# Following "optimal" directions, we end up in taking a poor route surely

Normal (90%):
60 min.

Busy (10%):
120 min.

10 min.

Start at 6am

A

B

30 min.

C

B'

10 min.

Normal (90%):
50 min. w.p. 99%
100 min. w.p. 1%

Busy (10%):
100 min. w.p. 99%
200 min. w.p. 1%

B-C is normal iff B'-C is normal

# Outline

- Limitations of traditional route selection

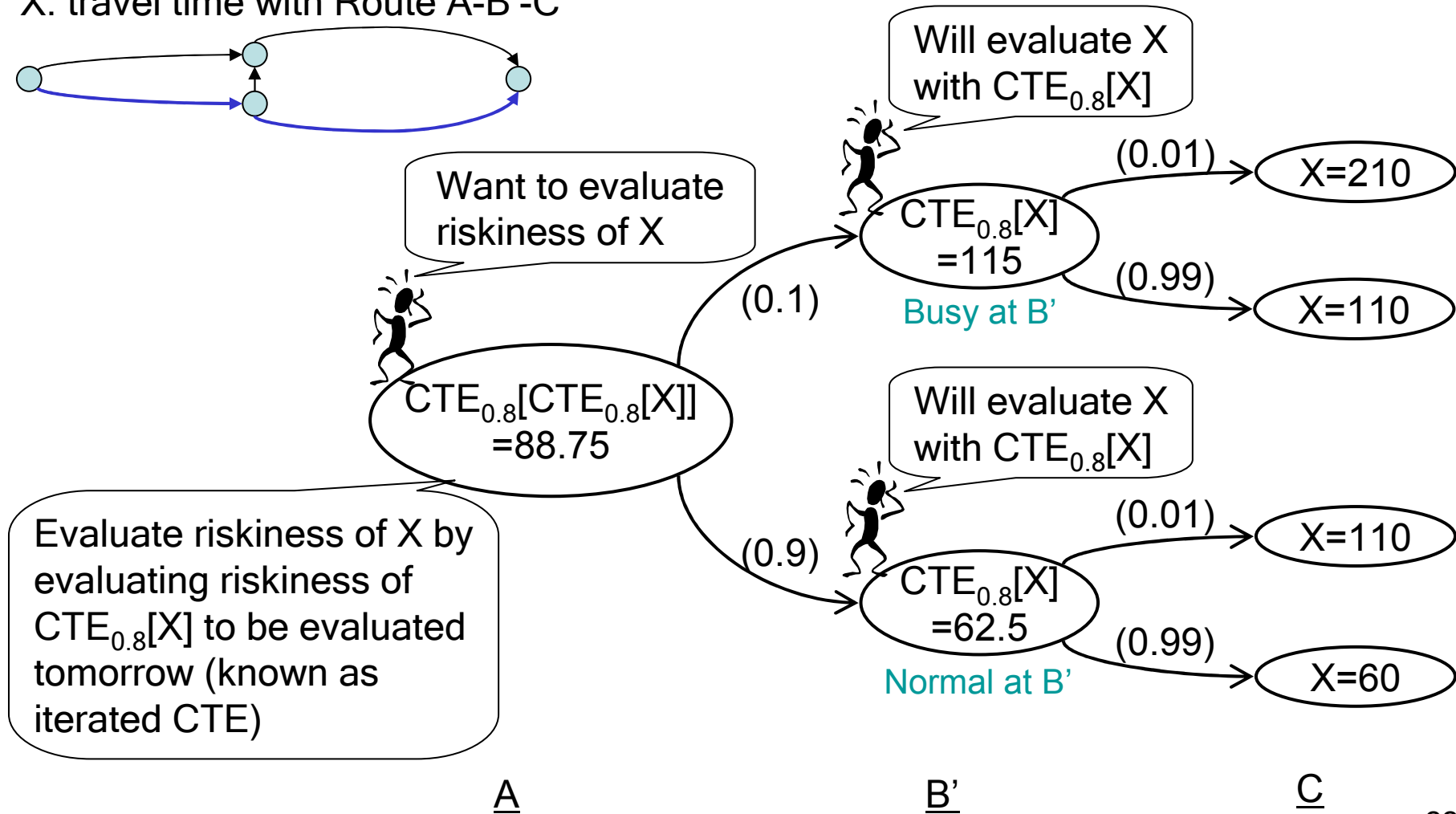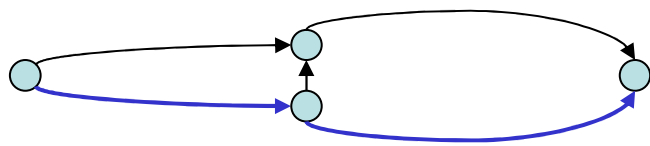- Difficulties in selecting objective functions

- Time-consistent Markov decision processes

- Effectiveness of time-consistent Markov decision processes

# We define a time-consistent MDP as the MDP whose objective is to minimize an iterated risk measure



X: travel time with Route A-B'-C

Will evaluate X with $CTE_{0.8}[X]$

Want to evaluate riskiness of X

$CTE_{0.8}[X] = 115$

Busy at B'

(0.01) → X=210

(0.99) → X=110

(0.1)

$CTE_{0.8}[CTE_{0.8}[X]] = 88.75$

Will evaluate X with $CTE_{0.8}[X]$

Evaluate riskiness of X by evaluating riskiness of $CTE_{0.8}[X]$ to be evaluated tomorrow (known as iterated CTE)

(0.9)

$CTE_{0.8}[X] = 62.5$

Normal at B'

(0.01) → X=110

(0.99) → X=60

A          B'          C

# Formally, an iterated risk measure is a dynamic risk measure having a recursive structure

- $(\Omega, F, P)$: Filtered probability space
  - $F_0 \subseteq F_1 \subseteq \ldots \subseteq F_N = F$
- Y: $F$-measurable random variable

- We say that $\rho$ is an iterated risk measure if
  - $\rho_N[Y] = Y$
  - $\rho_n[Y] = r_n[\rho_{n+1}[Y]]$
  - $r_n$: conditional risk measure mapping $F_{n+1}$-measurable random variable to $F_n$-measurable random variable

# Recursive definition implies dynamic programming finds the optimal policy for time-consistent MDP

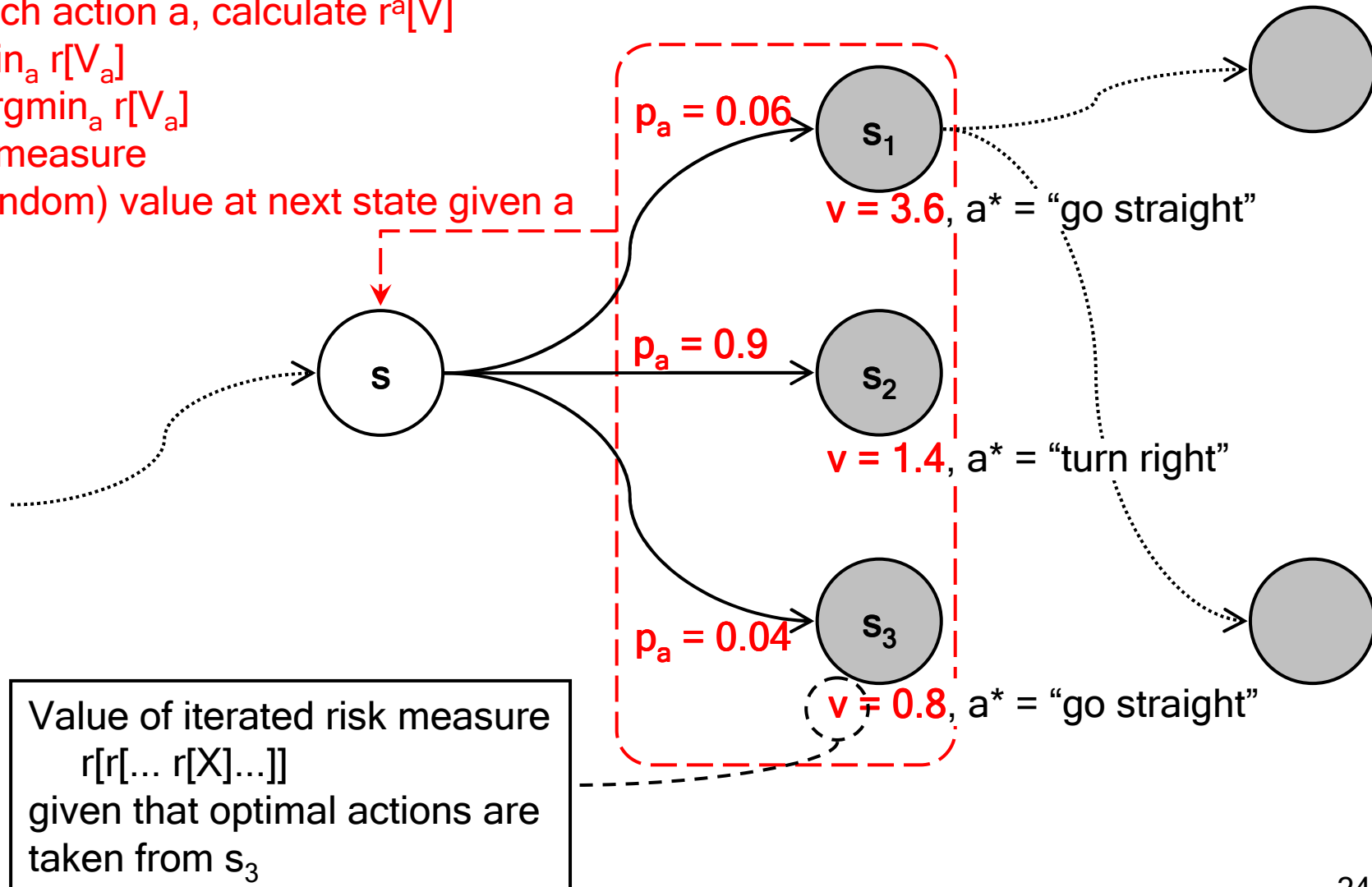For each action a, calculate $r^a[V]$

$v = \min_a r[V_a]$

$a^* = \text{argmin}_a r[V_a]$

r: risk measure

$V_a$: (random) value at next state given a

$p_a = 0.06$

$p_a = 0.9$

$p_a = 0.04$

s

$s_1$

$v = 3.6$, $a^* =$ "go straight"

$s_2$

$v = 1.4$, $a^* =$ "turn right"

$s_3$

$v = 0.8$, $a^* =$ "go straight"

Value of iterated risk measure
  $r[r[... r[X]...]]$
given that optimal actions are
taken from $s_3$

New

24

# More precisely, risk measures must be monotonic

| Properties of a risk measure, r | Optimal policy for MDPs with respect to the corresponding iterated risk measure |
|---|---|
| Monotonic:<br><br>$$X \leq Y \Rightarrow r(X) \leq r(Y)$$ | Can be found with dynamic programming<br>Need augmented states<br>    state := (state, accumulated cost) |
| Monotonic &<br><br>Translation invariant:<br><br>$$r(X + c) \leq r(X) + c$$ | No need for augmented states<br>Cannot discount future cost |
| Monotonic & Translation invariant &<br>Positive homogeneite:<br><br>$$r(aX) \leq ar(X)$$ | No need for augmented states<br>Can discount future cost |

**New**

# Dynamic programming with monotonic and translation-invariant iterated risk measures

- Markov decision process
  - $S_n$: State at time n (random variable, $F_n$-measurable)
  - $A_n$: Action at time n (random variable , $F_n$-measurable)
  - $C_n$: Cost between time n and time n+1, depending on $S_n$, $A_n$, $S_{n+1}$ (random variable, $F_{n+1}$-measurable)
  - $\mathbf{S}_n$: State space at time n (set)
  - $\mathbf{A}(s)$: Action space from state s (set)
  - $\Pi$ : Set of candidate policies (set)

- Find $\pi$ that minimizes $\rho_n\left[\sum_{\ell=0}^{N-1} C_\ell \mid S_n = s, \pi\right]$ or equivalently $\rho_n\left[\sum_{\ell=n}^{N-1} C_\ell \mid S_n = s, \pi\right]$

  for every $s \in \mathbf{S}_n$, n=0,…,N-1

$$V_n^*(s) \equiv \min_{\pi \in \Pi} \rho_n\left[\sum_{\ell=n}^{N-1} C_\ell \mid S_n = s, \pi\right]$$
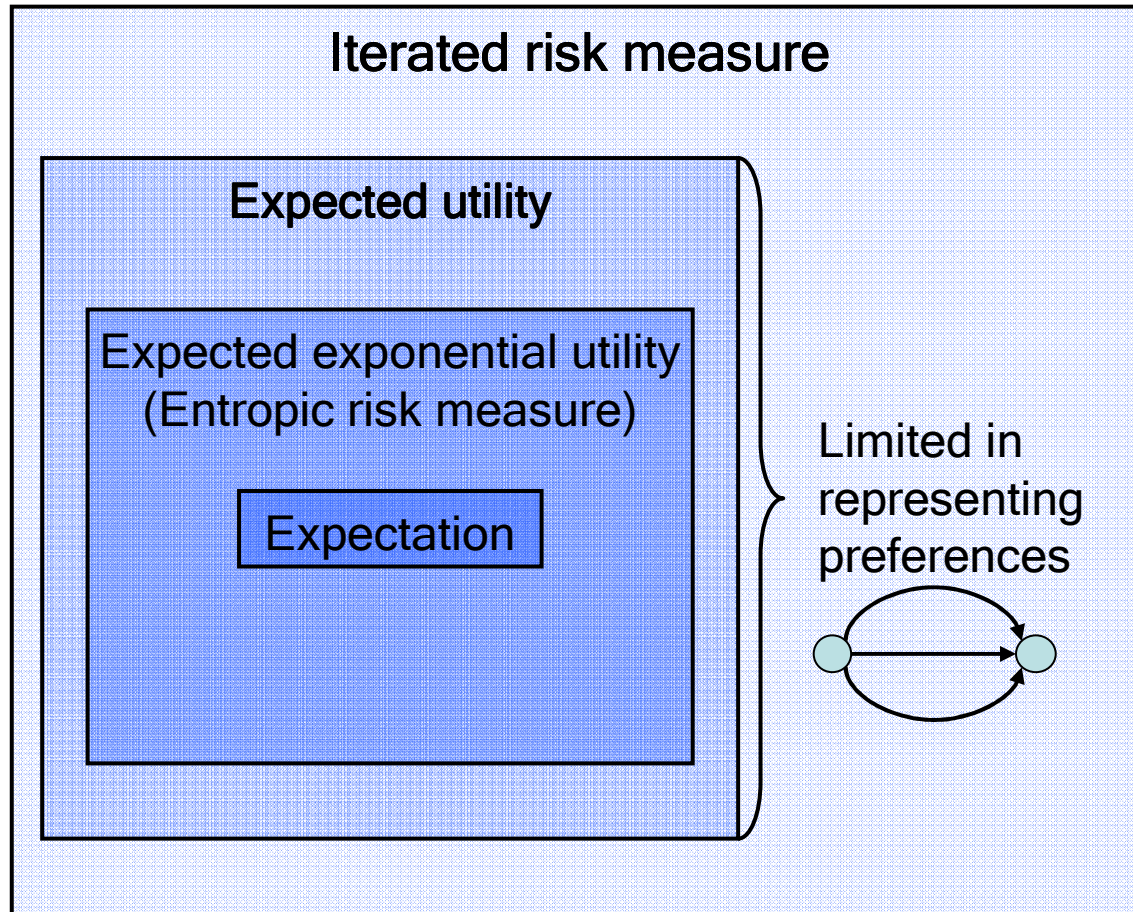
by translation-invariance
$\rho(X+c) = \rho(X) + c$

$$V_N^*(s) = 0 \qquad\qquad \forall s \in \mathbf{S}_N$$
$$V_n^*(s) = \min_{a \in \mathbf{A}(s)} r_n\left[C_n + V_{n+1}^*(S_{n+1}) \mid S_n = s, A_n = a\right] \quad \forall s \in \mathbf{S}_n, n = 0,..., N-1$$
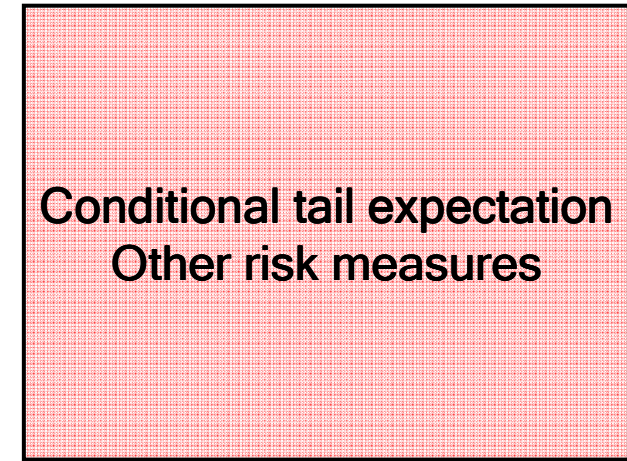
**New**

26

# Outline

- Limitations of traditional route selection

- Difficulties in selecting objective functions

- Time-consistent Markov decision processes

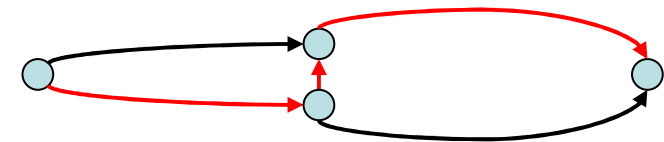- Effectiveness of time-consistent Markov decision processes

# Iterated risk measures overcome limitations of expected utility and other risk measures



**Iterated risk measure**

Expected utility

Expected exponential utility
(Entropic risk measure)

Expectation

Limited in representing preferences

**Conditional tail expectation
Other risk measures**
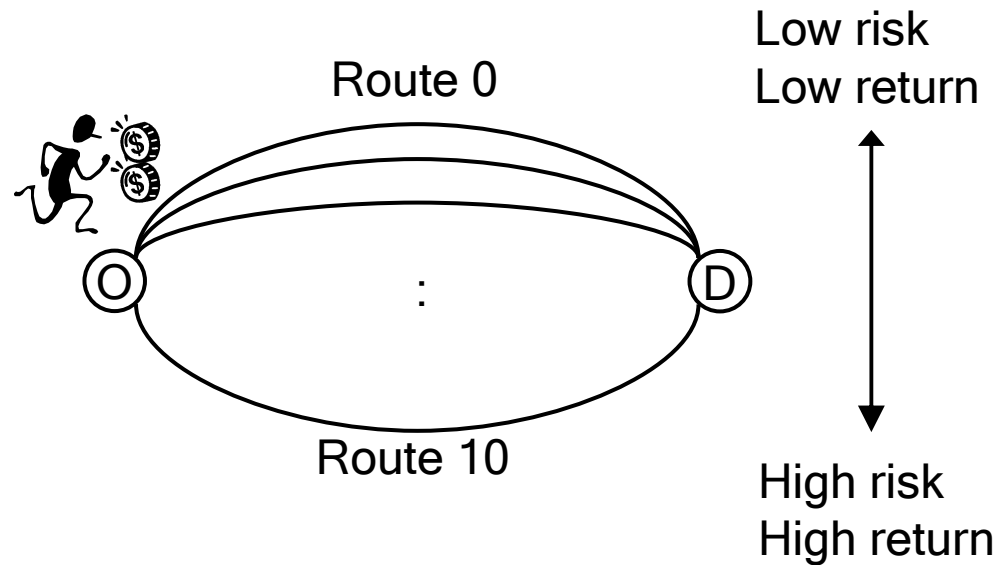
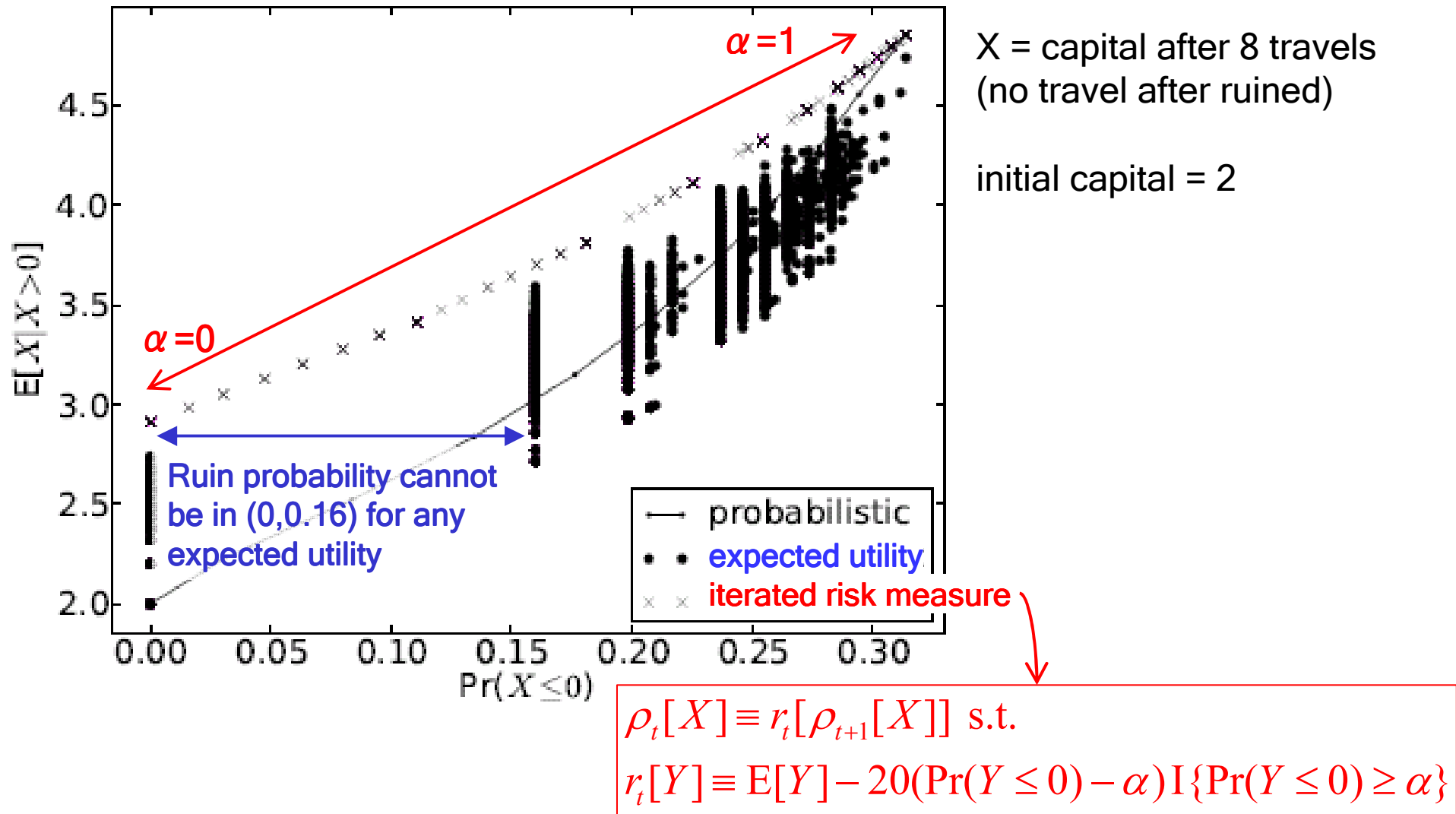Inconsistent decisions

Consistent decisions

# A driver takes only extreme routes if his decisions follow expected utility

**For any utility function, <span style="color:red">Route 0 (equivalent to doing nothing) or Route 10 (riskiest) is most preferable</span>**

Route 0

Low risk
Low return

O    :    D

High risk
High return

Route 10

| Route | 0 | 1 | | | ... | i | | | ... | 9 | | | 10 | |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| Net gain | 0 | -1 | 0 | +1 | ... | -1 | 0 | +1 | ... | -1 | 0 | +1 | -1 | +1 |
| Probability | 1 | 0.04 | 0.9 | 0.06 | ... | 0.04i | 1 – 0.1i | 0.06i | ... | 0.36 | 0.1 | 0.56 | 0.4 | 0.6 |

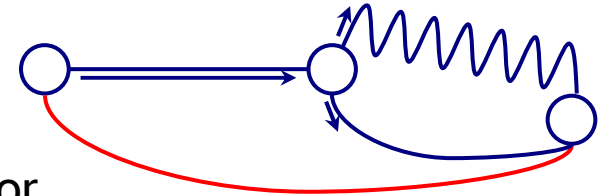# Iterated risk measures can represent the preference that cannot be represented with any expected utility



X = capital after 8 travels (no travel after ruined)

initial capital = 2

$$\rho_t[X] \equiv r_t[\rho_{t+1}[X]] \text{ s.t.}$$
$$r_t[Y] \equiv \mathrm{E}[Y] - 20(\Pr(Y \leq 0) - \alpha)\,\mathrm{I}\{\Pr(Y \leq 0) \geq \alpha\}$$

New

# Takeaways

- "Shortest path" is limited
  - Route selection follows a dynamic strategy
- Expected utility can only represent limited preferences for
  - personalized recommendation of dynamic strategies
  - realistic traffic simulation
- Traditional risk measures lead to inconsistent decisions
  - Inconsistent decision maker can surely lose infinite capital against rational decision maker
- Time-consistent MDP is defined with iterated risk measures
  - can represent broad preferences with consistent decisions
  - optimal policy found with dynamic programming

**Non-standard MDP**

**Standard MDP**

**Standard risk-sensitive MDP**

**Time-consistent MDP**

Inconsistent decisions

Consistent decisions

# 時間整合的マルコフ決定過程

## 恐神 貴行

## IBM東京基礎研究所

森村哲郎, Mikael Onsjoe （IBM東京基礎研究所）との共同研究

References:
T. Osogami, "Iterated risk measures for risk-sensitive Markov decision processes with discounted cost", in *Proceedings of UAI 2011.*
T. Osogami and T. Morimura, *Time-consistency of optimization problems*, manuscript, 2011.
T. Osogami and M. Onsjoe, *Overcoming a limitation of expected utility*, manuscript 2011.