

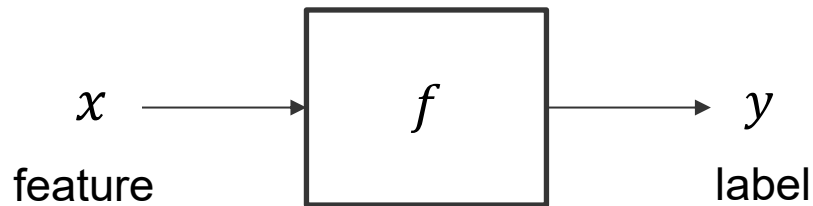
Some Recent Advances in Bandit Theory

Chen-Yu Wei 魏振宇

Postdoc @ UC Berkeley

Learning to Make Decisions

Machine Learning \approx Looking for a Function



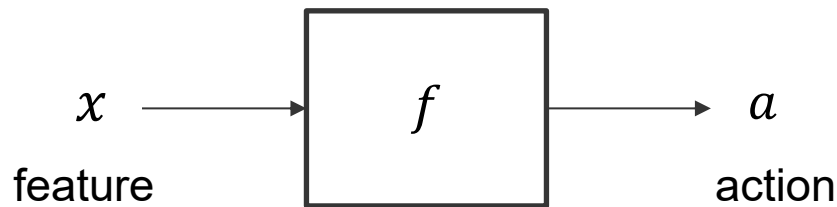
$$f \left(\text{image of a cat} \right) = \text{Cat} \quad (\text{classification})$$

$$f \left(\text{氣溫、濕度、風速...} \right) = 1000\text{mm 降雨量} \quad (\text{regression})$$

Decision Making Problem

Problems where the **learner's behavior** affects the **feedback**.

Learning to Make Decisions

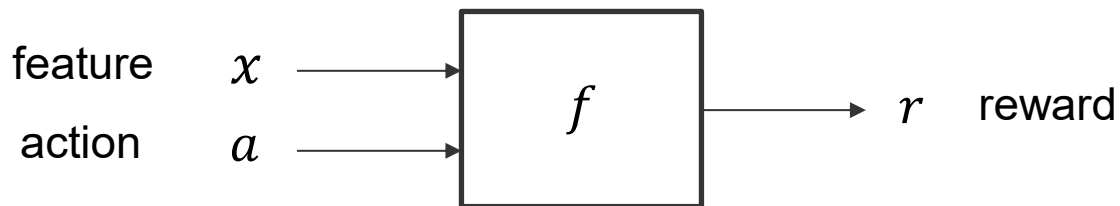


$$f \left(\text{movie watching history} \right) = \text{recommendation}$$

The equation shows a function f applied to "movie watching history" (represented by a screenshot of a Netflix movie grid) to produce a "recommendation" (represented by a movie poster for the film "Parasite").

Reduction from **decision-making** to **classification**

Learning to Make Decisions



f ( , ) = 7.5

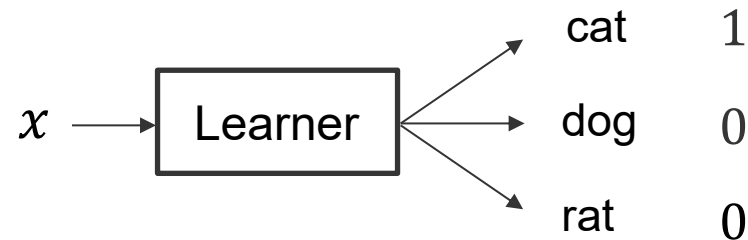
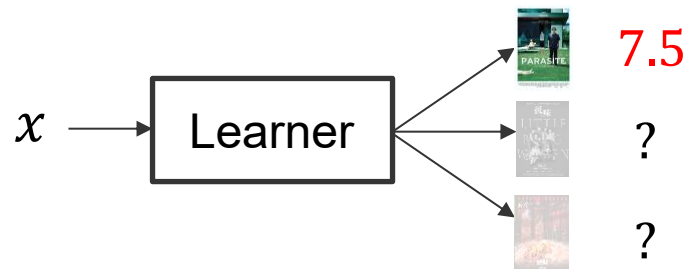
movie watching history , recommendation score

The equation is $f(\text{movie watching history}, \text{recommendation}) = 7.5$. The 'movie watching history' is represented by a screenshot of a Netflix movie grid. The 'recommendation' is represented by a movie poster for the film 'Parasite'. The result is '7.5 score'.

Reduction from **decision-making** to **regression**

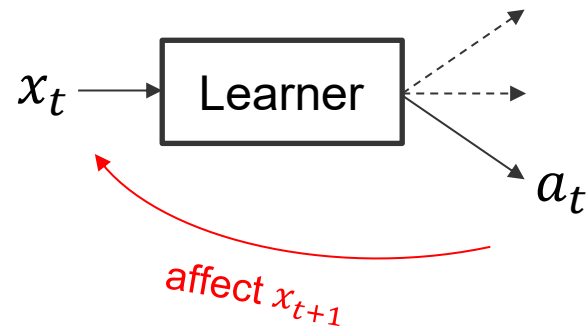
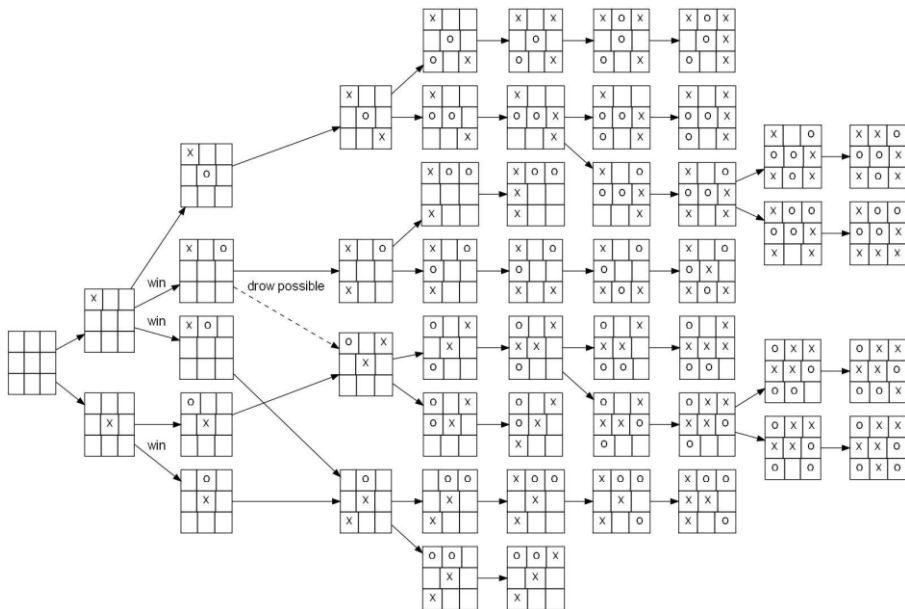
Challenges in Decision Problems (1/2)

1. Partial feedback → Need to **explore**

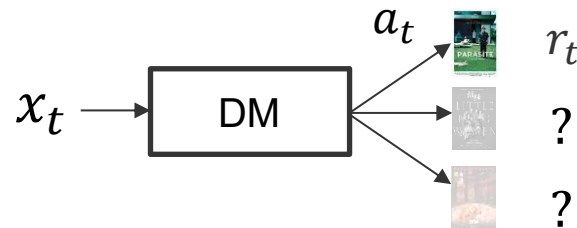


Challenges in Decision Problems (2/2)

2. Delayed feedback / long-term dependency → Need to do **credit assignment**



Today, we will consider problems with partial feedback but no long-term dependency.



A simple protocol (Contextual Bandit)

For $t = 1, 2, \dots$

Environment generates a **context** x_t (user profile, watching list)

Decision-maker takes an **action** a_t (a movie)

Environment reveals a **reward** r_t (user scoring)

(r_t depends on x_t and a_t)

(a_t does NOT affect x_{t+1})

Supervised
Learning

(full feedback)

<

**Contextual
Bandits**

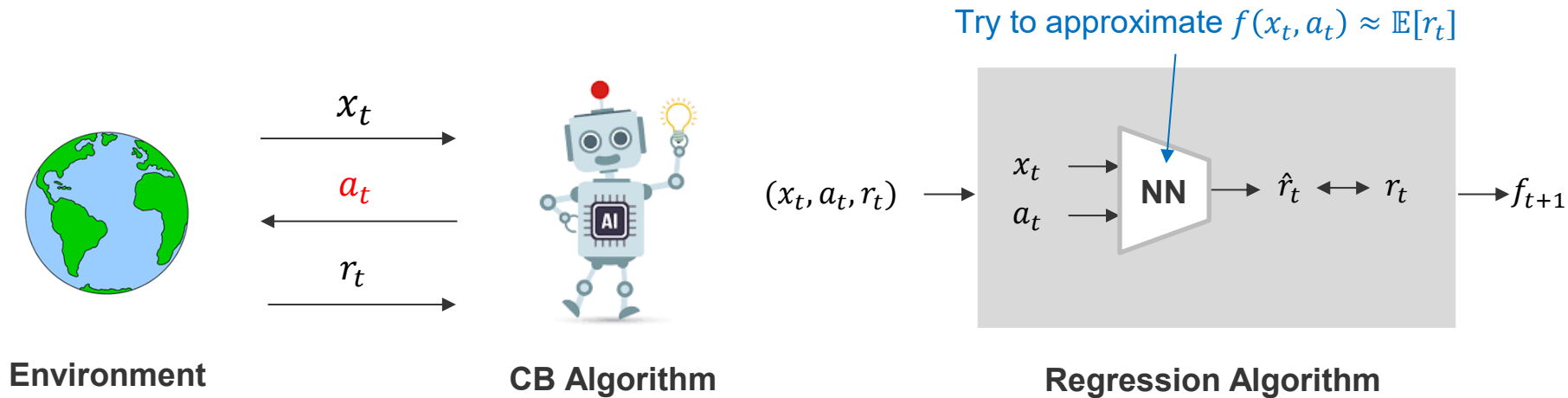
(partial feedback)

<<

Reinforcement
Learning

(partial + delayed feedback)

Reduction to Regression



Greedy: $a_t = \operatorname{argmax}_a f_t(x_t, a)$

Failure of the Greedy Strategy

Adding Exploration

- ϵ -greedy

$$a_t = \begin{cases} \operatorname{argmax}_a f_t(x_t, a) & \text{w. p. } 1 - \epsilon \\ \text{uniformly random} & \text{w. p. } \epsilon \end{cases}$$

- Boltzmann exploration

$$a_t \sim p_t(a) = \frac{\exp(\gamma \cdot f_t(x_t, a))}{\sum_{a'} \exp(\gamma \cdot f_t(x_t, a'))}$$

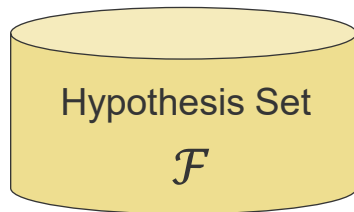
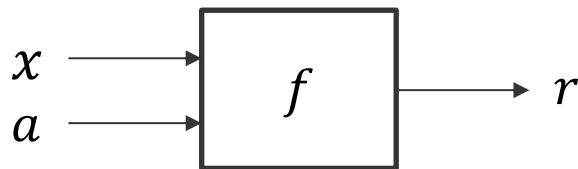
李宏毅老師 “DRL: Q-learning”

How good are they?

Today's mission: show an *optimal* exploration strategy for CB.

Some Theory

Formalization



Assumption: $\exists f^* \in \mathcal{F}$, s.t. $\mathbb{E}[r] = f^*(x, a)$

Optimal policy: choose $\operatorname{argmax}_a f^*(x_t, a)$

\therefore Performance Gap (a.k.a. **Regret**)

$$= \underbrace{\sum_{t=1}^T \max_a f^*(x_t, a)}_{\text{Optimal policy's expected total reward}} - \underbrace{\sum_{t=1}^T f^*(x_t, a_t)}_{\text{Learner's expected total reward}}$$

Optimal policy's
expected total reward

Learner's
expected total reward

The goal of the learner:
Minimize Regret

Strategy and Its Interpretation

$$p_t = \min_p \max_{f \in \mathcal{F}} \left\{ \underbrace{\max_{a'} f(x_t, a) - \mathbb{E}_{a \sim p}[f(x_t, a)]}_{(1)} - \underbrace{\gamma \mathbb{E}_{a \sim p} \left[(f(x_t, a) - \underbrace{f_t(x_t, a)}_{\substack{\text{output from the} \\ \text{regression algorithm}}})^2 \right]}_{(2)} \right\}$$

Diagram illustrating the strategy and its interpretation. The equation shows the minimization of regret and maximization of information gain simultaneously, with components labeled 1, 2, 3, and 4.

(1) = **regret** in round t (supposed that $f^* = f$)

(2) = **information gain** in round t (supposed that $f^* = f$)

(3) Minimize regret and maximize information gain simultaneously

(4) Consider the worst-case f^*

Foster and Rakhlin. Beyond UCB: Optimal and Efficient Contextual Bandits with Regression Oracles. ICML 2020.

Simplification

It suffices to use the following:

(the exact solution of the min-max program assuming $\mathcal{F}(x_t, \cdot) = \mathbb{R}^A$)

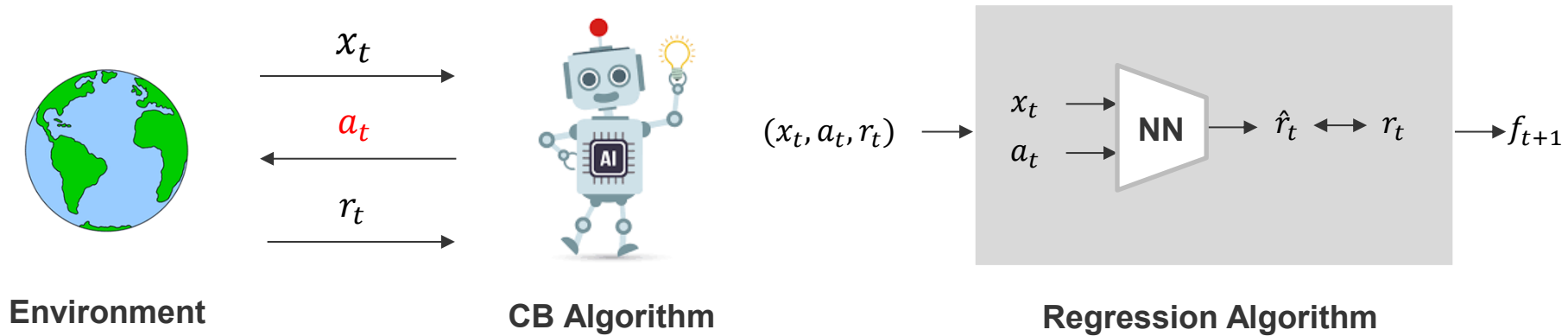
$$p_t(a) = \frac{1}{\lambda + \gamma \left(\max_{a'} f_t(x_t, a') - f_t(x_t, a) \right)}$$

(Inverse-Gap Weighting)

sub-optimality gap of action a predicted by f_t

normalization factor making $\sum_a p_t(a) = 1$

Reduction to Regression



$$\text{IGW: } a_t \sim p_t(a) = \frac{1}{\lambda + \gamma \left(\max_{a'} f_t(x_t, a') - f_t(x_t, a) \right)}$$

Theorem by Foster and Rakhlin (2020)

With Inverse-Gap Weighting,

$$\text{Regret}_{\text{CB}} \leq \frac{AT}{\gamma} + \underbrace{\gamma \sum_{t=1}^T (f_t(x_t, a_t) - f^*(x_t, a_t))^2}_{\substack{\text{Performance of the regression algorithm} \\ \text{Can be } O(\log |\mathcal{F}|)}} \leq O\left(\sqrt{AT \log |\mathcal{F}|}\right)$$

Choose a suitable γ

Exploration Strategies

- ϵ -greedy

$$a_t = \begin{cases} \operatorname{argmax}_a f_t(x_t, a) & \text{w. p. } 1 - \epsilon \\ \text{uniformly random} & \text{w. p. } \epsilon \end{cases}$$

- Boltzmann exploration

Cesa-Bianchi, Gentile, Lugosi, Neu.
Boltzmann exploration done right. NeurIPS 2017.

$$a_t \sim p_t(a) \propto \exp(\gamma \cdot f_t(x_t, a)) \propto \exp(-\gamma \cdot \text{Gap}(a))$$

- Inverse-gap weighting (optimal reduction from Contextual Bandit to Regression)

$$a_t \sim p_t(a) = \frac{1}{\lambda + \gamma \cdot \text{Gap}(a)}$$

Foster and Rakhlin. Beyond UCB: Optimal and Efficient
Contextual Bandits with Regression Oracles. ICML 2020.

Milestones in the Theory of Contextual Bandits

Abe and Long. Associative reinforcement learning using linear probabilistic concepts. 1999.

Langford and Zhang. The epoch-greedy algorithm for multi-armed bandits with side information. 2007.
(coin down the name *contextual bandit*)

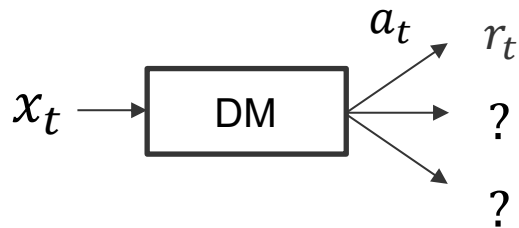
Reduction to classification:

- Dudik, et al. Efficient optimal learning for contextual bandits. 2011.
(optimal regret)
- Agarwal et al. Taming the monster: a fast and simple algorithm for contextual bandits. 2014.
(computationally efficient + optimal regret)

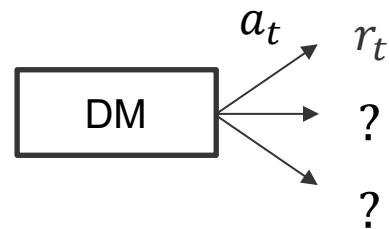
Reduction to regression:

- Agarwal et al. Contextual bandit learning with predictable rewards. 2012.
(optimal regret)
- Foster and Rakhlin. Beyond UCB: optimal and efficient contextual bandits with regression oracles. 2020.
(computationally efficient + optimal regret)

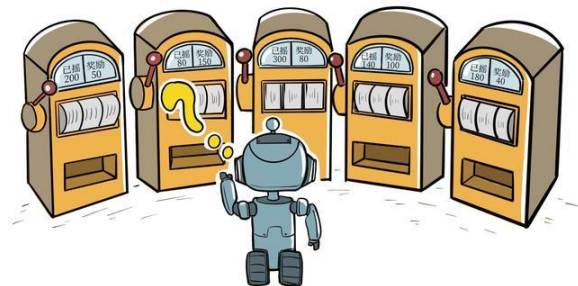
Robust and Adaptive Bandit Algorithms

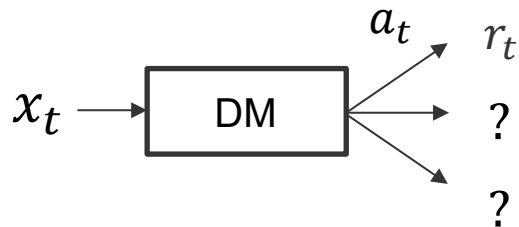


Contextual Bandits

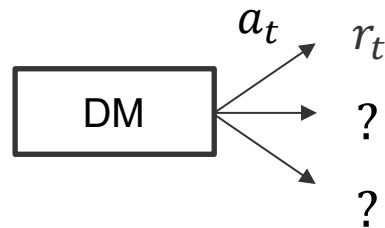


Multi-Armed Bandits (多臂吃角子老虎機)





Contextual Bandits



Multi-Armed Bandits (多臂吃角子老虎機)

For $t = 1, 2, \dots$

Decision-maker takes an **action** a_t

Environment reveals a **reward** r_t

Two Reward Generation Processes (Figures from Wouter Koolen)

For $t = 1, 2, \dots$

Environment decides **rewards** $r_t(a)$ for all actions a

Decision-maker takes an **action** a_t

Environment reveals a **reward** $r_t(a_t)$



fixed over time

$\mathbb{E}[r_t(a)] = f^*(a)$
(stochastic world)



$r_t(a)$ are arbitrarily decided
(adversarial world)

$$\text{Regret} = \max_a \sum_{t=1}^T r_t(a) - \sum_{t=1}^T r_t(a_t)$$

Why consider adversarial worlds?

(Robustness)

李彥寰老師“預測、學習與賽局”



Development of Stochastic and Adversarial MAB



Thompson 1933

Lai and Robbins 1985

Agrawal 1995

Auer et al. 2002 $O\left(\frac{A \log T}{\Delta}\right)$ **optimal**



Auer et al. 2002

Audibert and Bubeck, 2009 $O(\sqrt{AT})$ **optimal**

Δ = difference between the expected reward of the best and the second-best action

Bubeck and Slivkins 2012:

Is there a **single algorithm** with optimal regret in both worlds?

The **best-of-both-world** problem, “**robust and adaptive**” learning

Development of Stochastic and Adversarial MAB

Regret



Strategy

$\text{Gap}(a)$: current estimation for the sub-optimality gap of action a

Auer et al., 2002
(classic algorithm for adversarial MAB)

$$\sqrt{AT \log A}$$

$$\sqrt{AT \log A}$$

$$p(a) = \text{softmax}(a) \propto e^{-\gamma \cdot \text{Gap}(a)}$$

Seldin and Slivkins, 2014

$$\frac{A \log^2 T}{\Delta} + \frac{A^2}{\Delta^3}$$

$$\sqrt{AT \log A}$$

$$p(a) \propto \text{softmax}(a) + \min \left\{ \frac{1}{\sqrt{At}}, \frac{1}{t \text{Gap}(a)^2} \right\}$$

W and Luo, 2018

$$\frac{A \log T}{\Delta}$$

$$\sqrt{AT \log T}$$

$$p(a) = \frac{1}{\lambda + \gamma_t \text{Gap}(a)}$$

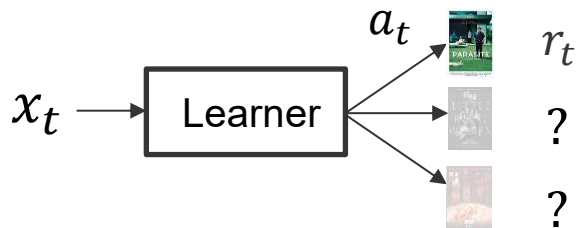
Zimmert and Seldin, 2019

$$\frac{A \log T}{\Delta}$$

$$\sqrt{AT}$$

$$p(a) = \frac{1}{\left(\lambda + \sqrt{t} \text{Gap}(a) \right)^2}$$

Summary

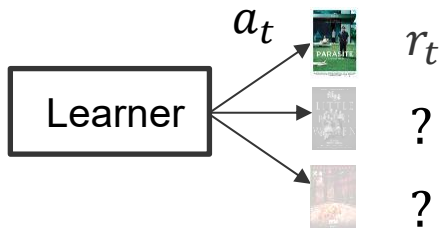


Contextual Bandit

Assumption: $\mathbb{E}[r_t] = f^*(x_t, a_t)$

Can solve it by **regression** + **exploration strategy**
(**Inverse Gap Weighting**)

$$\text{Regret} \leq \frac{AT}{\gamma} + \gamma \cdot \text{regression_error}$$



Multi-Armed Bandit

Either adversarial or stochastic ($\mathbb{E}[r_t] = f^*(a_t)$)

Can get robustness and adaptivity (best-of-both-world)
by **Inverse Squared Gap Weighting**

$$\text{Regret} \leq \sqrt{AT} \text{ (adversarial), or } \frac{A \log T}{\Delta} \text{ (stochastic)}$$