

Reward Variance in Markov Chains

Tom Verhoeff

Technische Universiteit Eindhoven

Faculty of Math. & Computing Science
Software Construction

T.Verhoeff@TUE.NL

<http://www.win.tue.nl/~wstomv/>

Spider on Cube

A spider walks randomly on the faces of a cube:



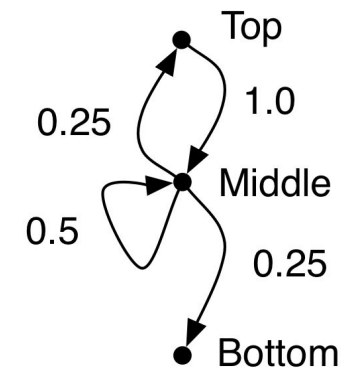
What is the **expected time** for the spider to get off the cube?

What is the corresponding **variance**?

© 2004, Tom Verhoeff

Reward Variance-2

Markov Chain



Finite state machine with transition **probabilities** and **rewards**

© 2004, Tom Verhoeff

Reward Variance-3

Expected Reward

Transition **probabilities** $p_{s,t}$ and **rewards** $r_{s,t}$ for states $s, t \in \Omega$

Lift to walks $w = s_0 s_1 \dots s_n$

Probability of walk w : $P.w = p_{s_0.s_1} * \dots * p_{s_{n-1}.s_n}$

Reward of walk w : $R.w = r_{s_0.s_1} + \dots + r_{s_{n-1}.s_n}$

Walks from s to absorption in A : $W.s = \bigcup_{t \in \Omega} \{stv \mid tv \in W.t\}$ ($s \notin A$)

Expected reward from s to absorption in A :

$$\mathcal{E}_{W.s}[R] = \sum_{w \in W.s} P.w * R.w$$

© 2004, Tom Verhoeff

Reward Variance-4

Expected Walk Length for Spider

Obtain equations by **generalizing** and **conditioning**.

Expected walk lengths $\mu_s = \mathcal{E}_{W.s}[R]$ from face s to bottom:

$$\begin{aligned}\mu_T &= 1 + \mu_M \\ \mu_M &= 0.25(1 + \mu_T) + 0.5(1 + \mu_M) + 0.25(1 + \mu_B) \\ \mu_B &= 0\end{aligned}$$

Solution:

$$\begin{aligned}\mu_T &= 6 \\ \mu_M &= 5 \\ \mu_B &= 0\end{aligned}$$

© 2004, Tom Verhoeff

Reward Variance-5

Conditioning on First Transition to A

For $s \notin A$:

$$\begin{aligned}
 & \mathcal{E}_{W.s}[X] \\
 = & \{ \text{definition of } \mathcal{E} \} \\
 & \sum_{w \in W.s} P.w * X.w \\
 = & \{ \text{write } w = sv \text{ for } v \in W.t, \text{ because } s \notin A \} \\
 & \sum_{t \in \Omega} \sum_{v \in W.t} P.sv * X.sv \\
 = & \{ \text{recurrence for walk probability: } P.sv = p.s.t * P.v \text{ for } v \in W.t \} \\
 & \sum_{t \in \Omega} \sum_{v \in W.t} p.s.t * P.v * X.sv \\
 = & \{ \text{distribute } p.s.t * \text{ outside } \sum_v \text{ (} p.s.t \text{ does not depend on } v \text{)} \} \\
 & \sum_{t \in \Omega} p.s.t * \sum_{v \in W.t} P.v * X.sv \\
 = & \{ \text{definition of } \mathcal{E} \} \\
 & \mathcal{E}_{t \in \Omega} [\mathcal{E}_{v \in W.t} [X.sv]]
 \end{aligned}$$

Equations for Expected Reward until Absorption

For $s \notin A$:

$$\begin{aligned}
 & \mathcal{E}_{W.s}[R] \\
 = & \{ \text{conditioning on first state } t \text{ after state } s, \text{ using } s \notin A \} \\
 & \mathcal{E}_{t \in \Omega} [\mathcal{E}_{v \in W.t} [R.sv]] \\
 = & \{ \text{recurrence for walk reward: } R.sv = r.s.t + R.v \text{ for } v \in W.t \} \\
 & \mathcal{E}_{t \in \Omega} [\mathcal{E}_{v \in W.t} [r.s.t + R.v]] \\
 = & \{ \text{linearity of expectation (} r.s.t \text{ is independent of } v \text{)} \} \\
 & \mathcal{E}_{t \in \Omega} [r.s.t + \mathcal{E}_{v \in W.t} [R.v]] \\
 = & \{ \text{simplify notation} \} \\
 & \mathcal{E}_{t \in \Omega} [r.s.t + \mathcal{E}_{W.t}[R]]
 \end{aligned}$$

Linear equations with unknowns $\mu_s = \mathcal{E}_{W.s}[R]$:

$$\mu_s = \sum_{t \in \Omega} p.s.t * (r.s.t + \mu_t)$$

Variance

Definition:

$$\mathcal{V}[X] = \mathcal{E} [(X - \mathcal{E}[X])^2]$$

Often computed via second moment:

$$\begin{aligned}
 \mathcal{V}[X] &= \mathcal{E} [(X - \mathcal{E}[X])^2] \\
 &= \mathcal{E} [X^2 - 2X\mathcal{E}[X] + \mathcal{E}^2[X]] \\
 &= \mathcal{E} [X^2] - 2\mathcal{E}[X]\mathcal{E}[X] + \mathcal{E}^2[X] \\
 &= \mathcal{E} [X^2] - \mathcal{E}^2[X]
 \end{aligned}$$

Numerical disadvantage: loss of accuracy through **cancellation**

Basic Properties of Variance

A constant offset does not affect the variance:

$$\mathcal{V}[c + X] = \mathcal{E} [(c + X - \mathcal{E}[c + X])^2] = \mathcal{E} [(X - \mathcal{E}[X])^2] = \mathcal{V}[X]$$

Via second moment:

$$\mathcal{V}[c + X] = \mathcal{E} [(c + X)^2] - (c + \mathcal{E}[X])^2$$

Combine:

$$\mathcal{E}[(c + X)^2] = (c + \mathcal{E}[X])^2 + \mathcal{V}[X] \quad (1)$$

Reward Variance on Walks until Absorption

For $s \notin A$, we calculate

$$\begin{aligned}
 & \mathcal{V}_{W.s}[R] \\
 = & \{ \text{definition of } \mathcal{V} \} \\
 & \mathcal{E}_{W.s} \left[(R - \mathcal{E}_{W.s}[R])^2 \right] \\
 = & \{ \text{conditioning on first state } t \text{ after state } s, \text{ using } s \notin A \} \\
 & \mathcal{E}_{t \in \Omega} \left[\mathcal{E}_{v \in W.t} \left[(R.s.v - \mathcal{E}_{W.s}[R])^2 \right] \right] \\
 = & \{ \text{recurrence for reward: } R.s.v = r.s.t + R.v \text{ for } v \in W.t \} \\
 & \mathcal{E}_{t \in \Omega} \left[\mathcal{E}_{v \in W.t} \left[(r.s.t + R.v - \mathcal{E}_{W.s}[R])^2 \right] \right] \\
 = & \{ (1), \text{ using that } r.s.t - \mathcal{E}_{W.s}[R] \text{ does not depend on } v \} \\
 & \mathcal{E}_{t \in \Omega} \left[(r.s.t + \mathcal{E}_{W.t}[R] - \mathcal{E}_{W.s}[R])^2 + \mathcal{V}_{W.t}[R] \right]
 \end{aligned}$$

Equations for Reward Variance

System of linear equations with unknowns $\sigma_s^2 = \mathcal{V}_{W.s}[R]$ for $s \in \Omega$,

involving $\mu_s = \mathcal{E}_{W.s}[R]$ as parameters:

$$\sigma_s^2 = \sum_{t \in \Omega} p.s.t * \left((r.s.t + \mu_t - \mu_s)^2 + \sigma_t^2 \right)$$

Variance in Walk Length for Spider

Variance in walk length $\sigma_s^2 = \mathcal{V}_{W.s}[R]$ from face s to bottom:

$$\begin{aligned}
 \sigma_T^2 &= (1 + \mu_M - \mu_T)^2 + \sigma_M^2 \\
 \sigma_M^2 &= 0.25 \left((1 + \mu_T - \mu_M)^2 + \sigma_T^2 \right) + \\
 & \quad 0.5 \left((1 + \mu_M - \mu_M)^2 + \sigma_M^2 \right) + \\
 & \quad 0.25 \left((1 + \mu_B - \mu_M)^2 + \sigma_B^2 \right) \\
 \sigma_B^2 &= 0
 \end{aligned}$$

Solution:

$$\begin{aligned}
 \sigma_T^2 &= 22 \\
 \sigma_T &\approx 4.69
 \end{aligned}$$

Conclusion

- Importance of the variance
- Simple direct formula to calculate variance in reward until absorption in Markov chain
- Numerically attractive
- Applied to score variance of optimal strategy for Yahtzee
Acyclic Markov chain with $\approx 10^9$ states