

# 強化学習：Q 学習

2017/03/23

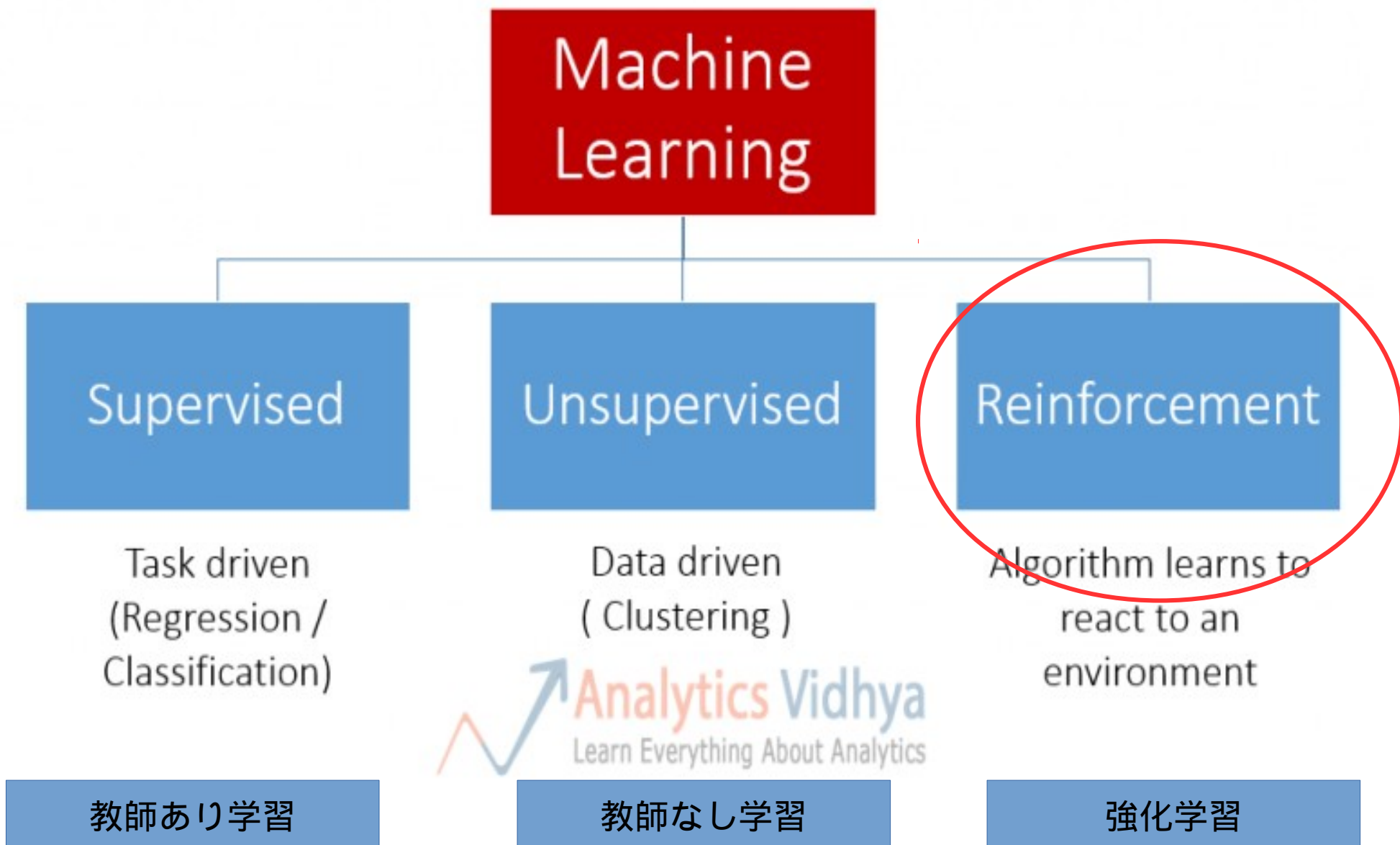
1426071 スプラトマン ジョシュア

# 学習とは？

- 人工知能の一部
- 人間の学習能力を機械（コンピュータ）で真似する
- あるデータに対して学習し、そこからパターンを見つけ出す
- <https://youtu.be/qv6UV0Q>

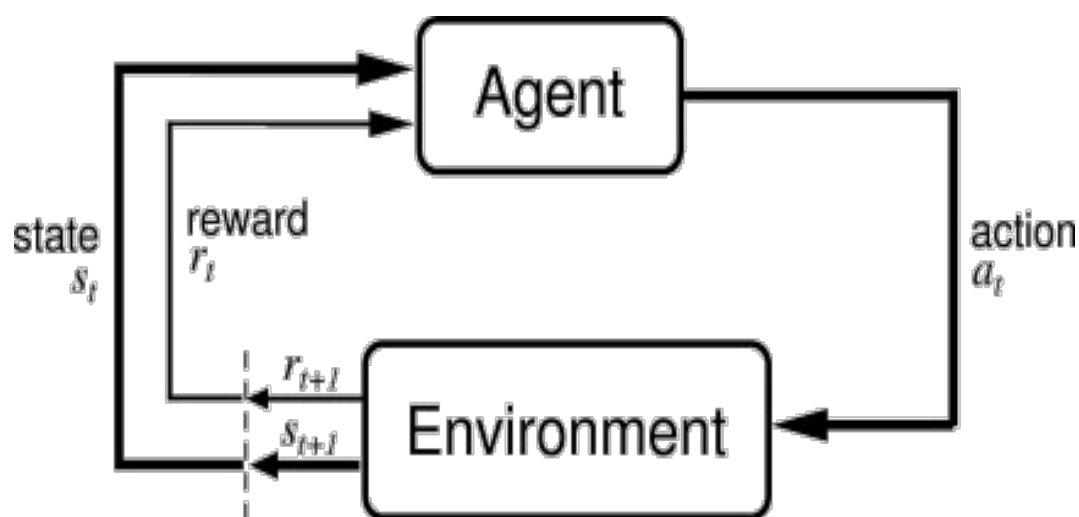


# Types of Machine Learning



# 強化学習 (RL) とは？

- 状態に対してどんな行動をする学習タイプ
- 最大の報酬を探す
- 赤ちゃんが歩きたい時考えよう！



# 強化学習 vs

## 教師あり学習 (SL)

- 教師データを比較しながら学習
- チェスを考える
  - SL: あやゆる戦略を事前対策する
  - RL: 状態に対して行動を決める

## • 教師なし学習 (UL)

- データのパターンを探す
- 好きな雑誌を進める AI を考える
  - UL: 過去読んだ似たような雑誌をお進め
  - RL: たまたま新しい雑誌とか出し、フィードバックを受けながら好みを探す

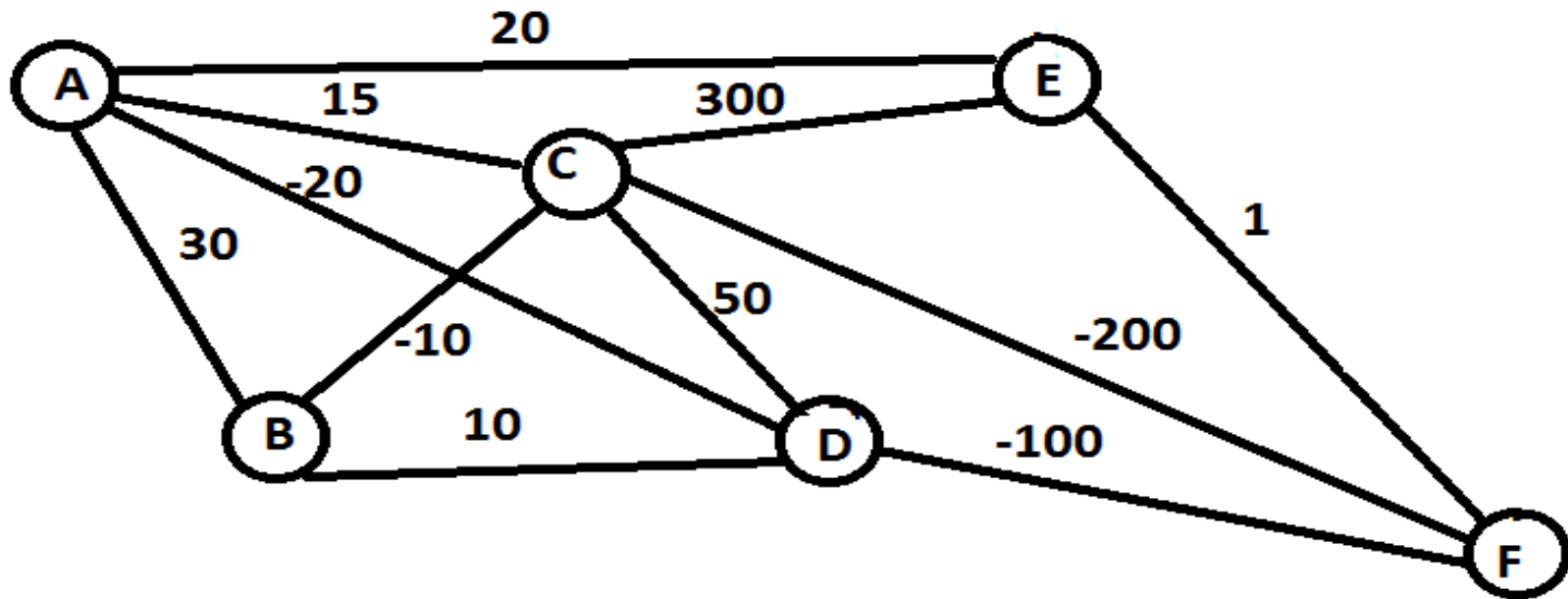
# Basic of 強化学習

- 搾取 vs 探検どっち？ 適当はダメ
- markov モデル：
  - 状態 :  $S$
  - 行動 :  $A$
  - 報酬 :  $R$
  - ポリシー :  $\pi$  ( 状態 + 行動 )
  - 値 :  $V$



$$E(r_t | \pi, s_t)$$

# Basic of 強化学習



状態：？

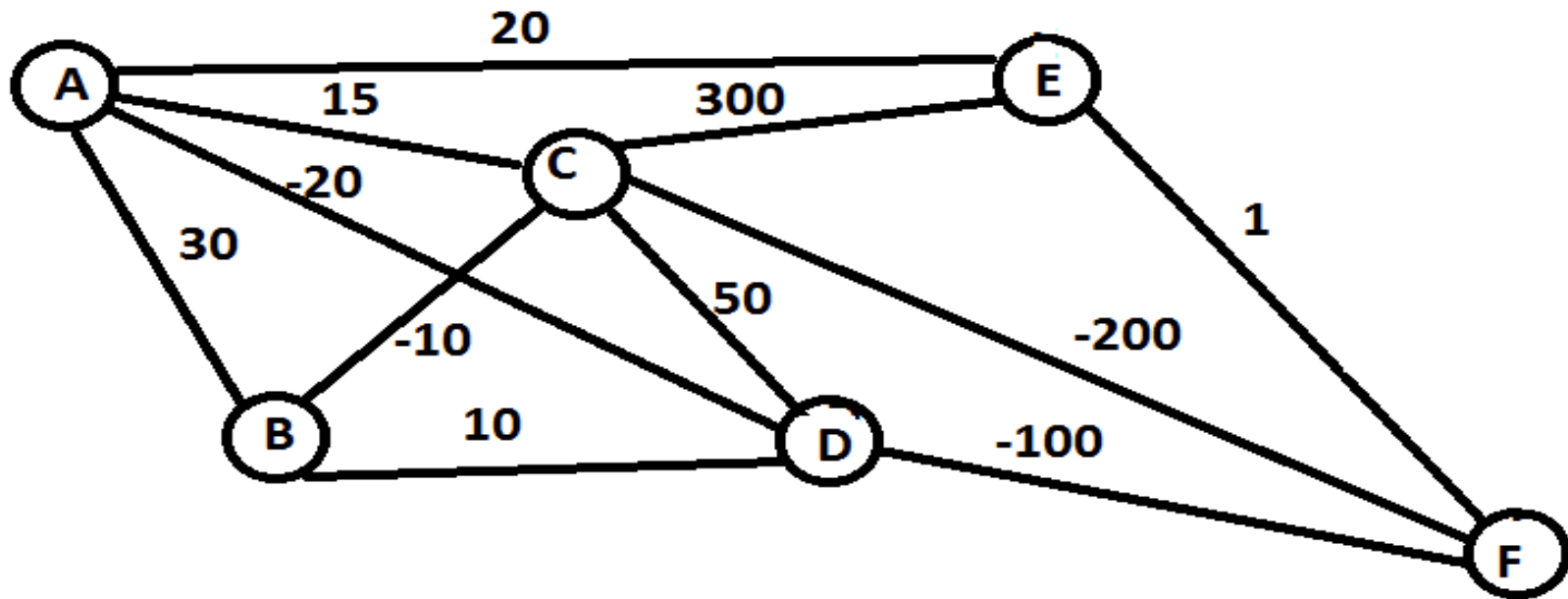
報酬：？

値：

行動：？

ポリシ：？

# Basic of 強化学習

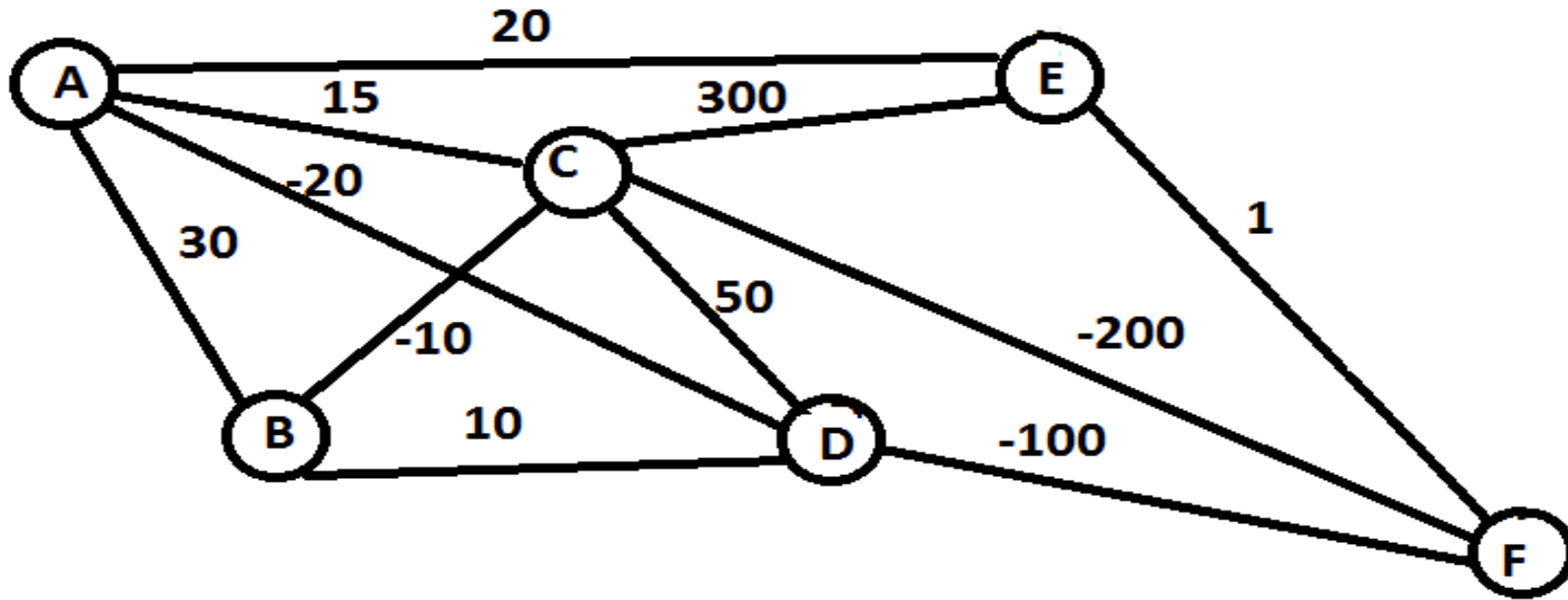


状態 : A,C,D    報酬 : 20,15    値 : 20+1,etc

行動 : A->C    ポリシ : A-E-F, C-B, C-D-F

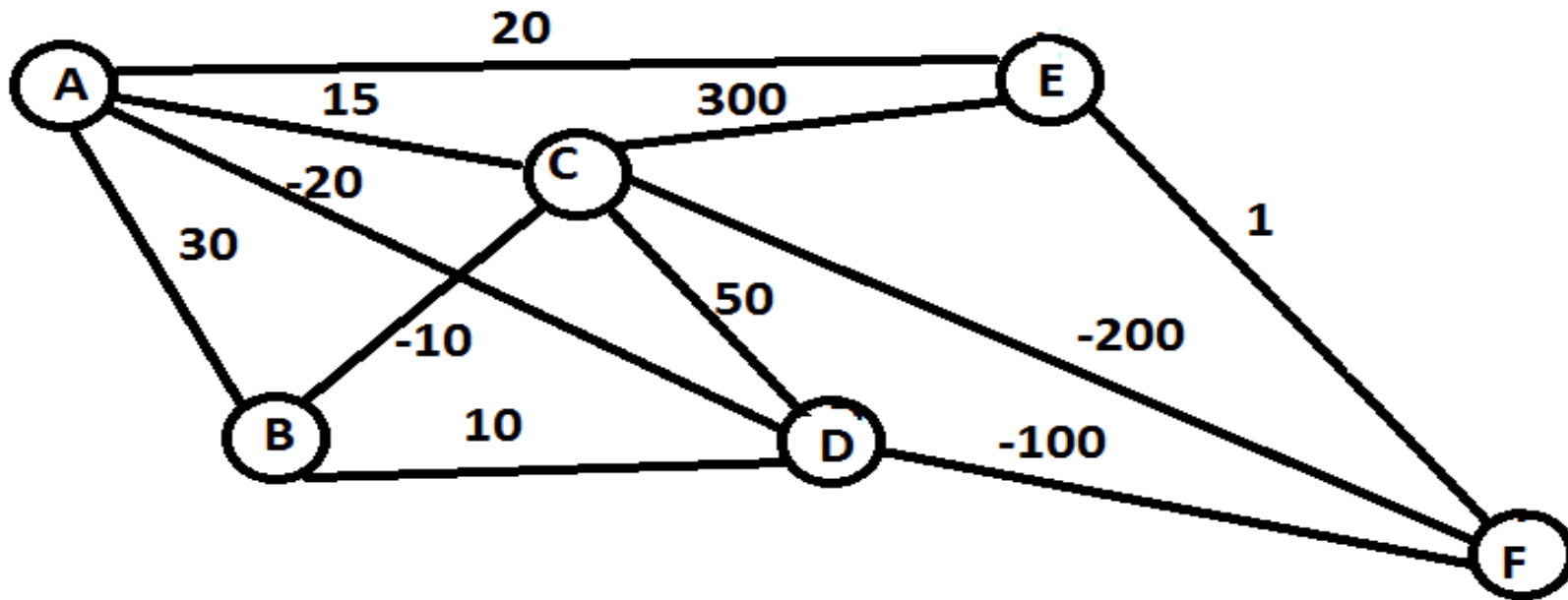


# Basic of 強化学習



- スタート : A    ゴール : F    最短経路 ?

# Basic of 強化学習



- スタート : A    ゴール : F    最短経路 : A-D-F
- A->(B,C,D,E)    A-D が一番コスト低い → -20
- D->(B,C,F)    D-F が一番コスト低い → 値 : -120

# おめでとう！

- さっき強化学習アルゴリズムをやりました！
- epsilon greedy 法
  - 搾取 vs 探検どっち？



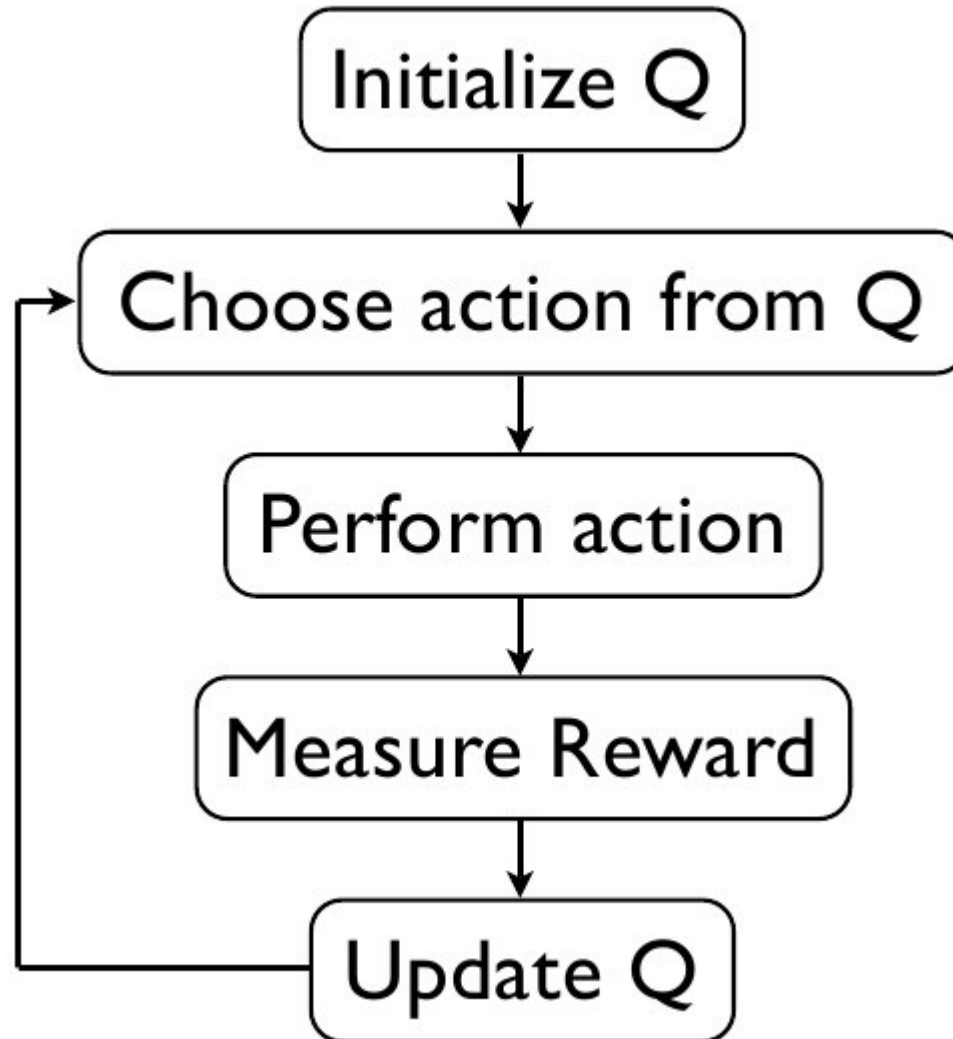
# 他のやりかた？

- Value Based 最適の値を重心として進む
- Policy Based 最適のポリシを重心として進む
- ちなみに epsilon greedy 法は？

# Q 学習とは？

- 強化学習のアルゴリズムの一つ
- Policy Based
- <https://youtu.be/JS1fImajAEE>
- <https://youtu.be/OisMuSW2F2w>
- AlphaGo! Neural Network で Q 学習！

# Q 学習のながれ ( 訓練 )



# Q 学習のながれ

- 今回はシンプルな物を作ります
- Neural Network ではなく QTable (nxn) 行列を使います
- QTable を学習することが目的
- 学習した QTable を使って実装する



F

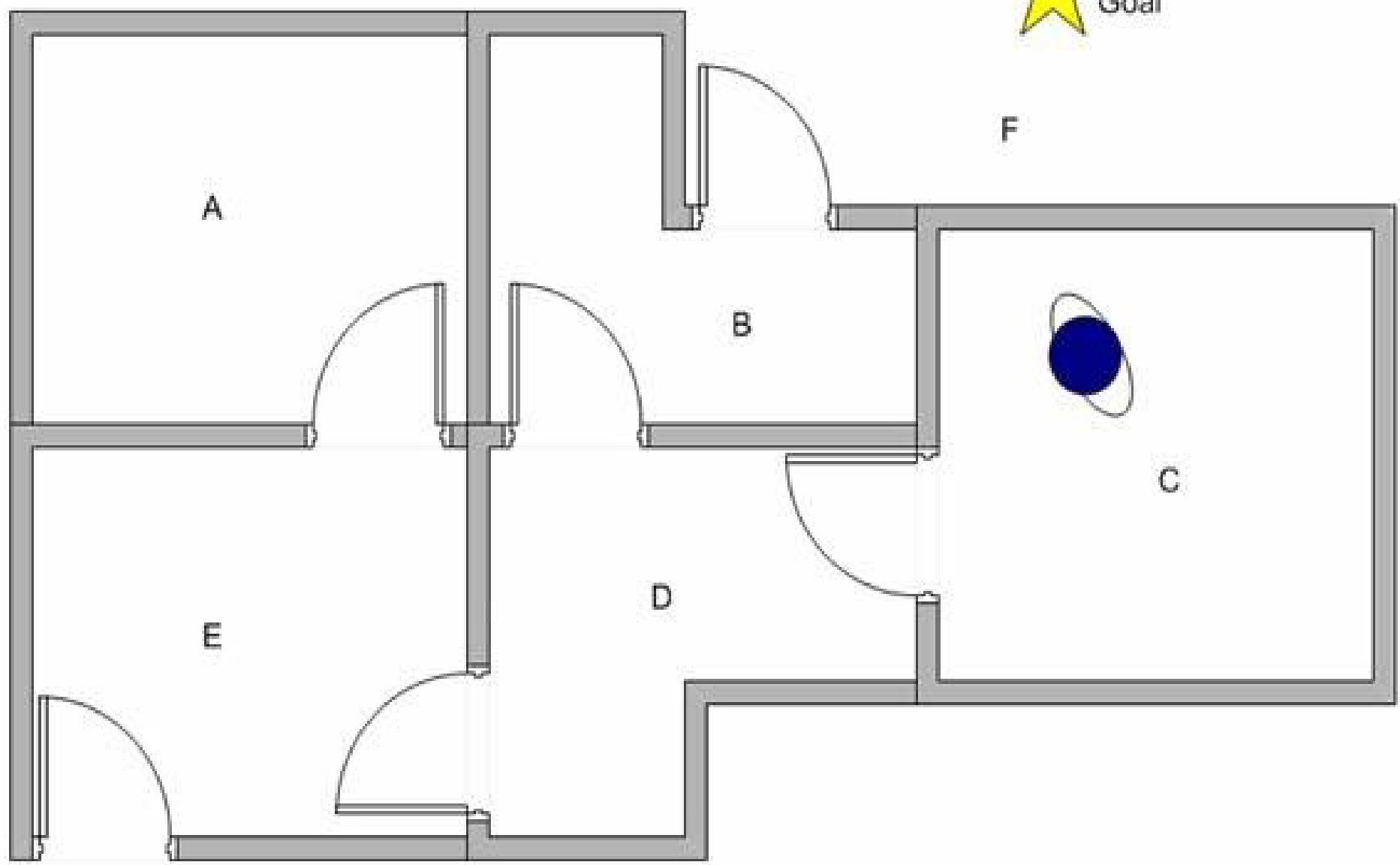
A

B

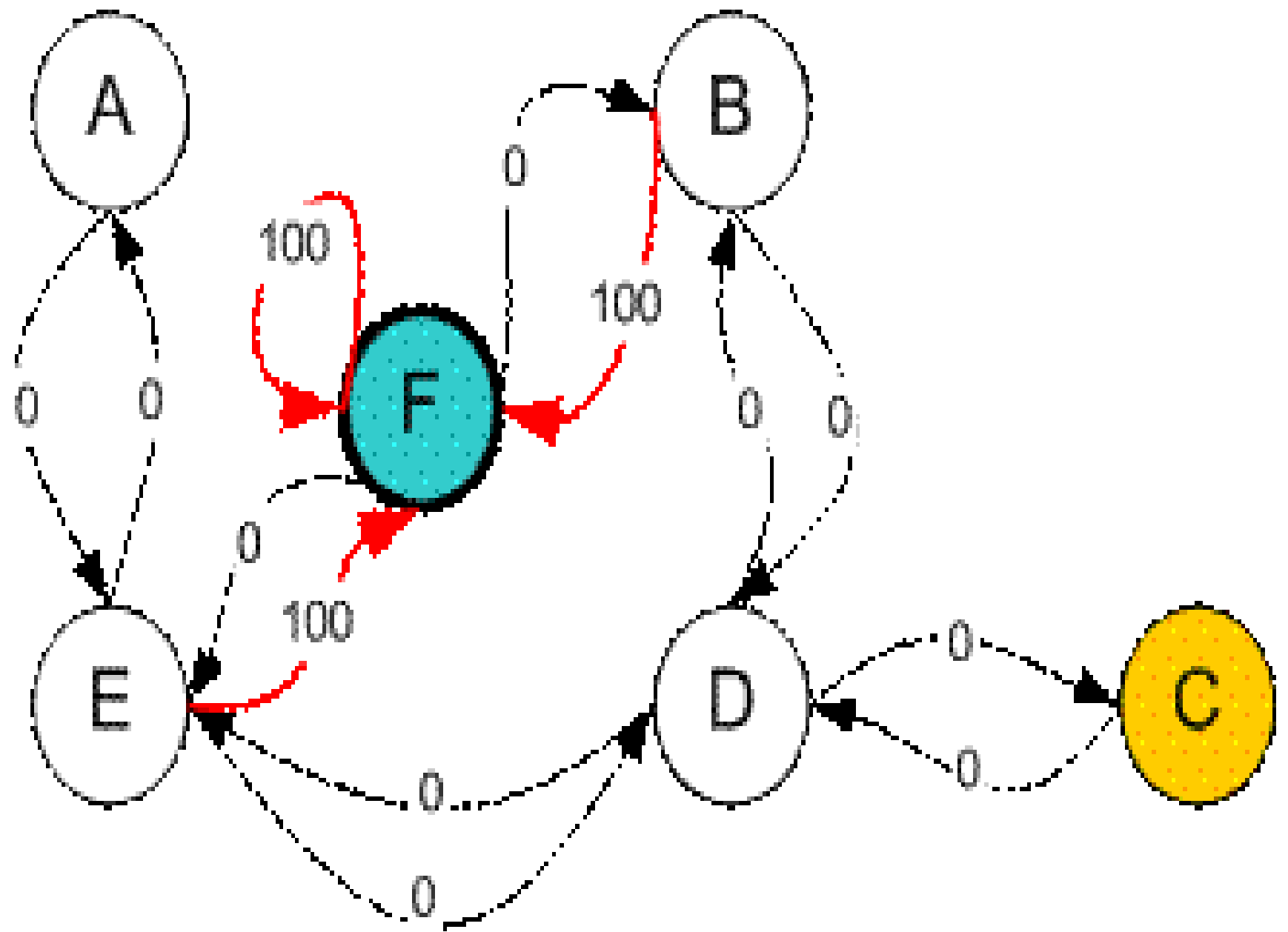
C

D

E





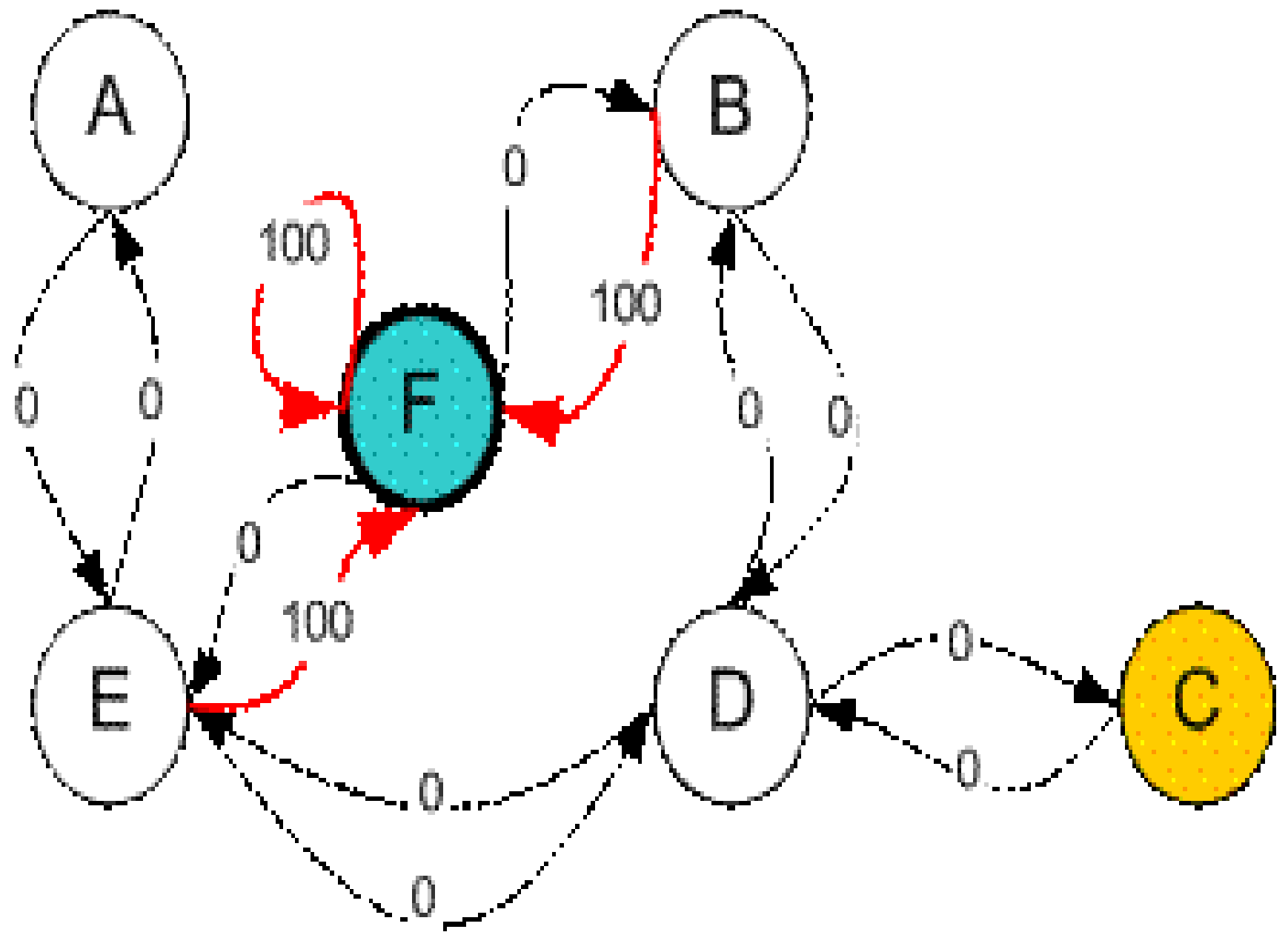


# markov モデル

- 状態  $S$ : ?
- 行動  $A$ : ?
- 報酬  $R$ : ?
- ポリシ  $\pi$ : ?
- 値  $V$ : ?

# markov モデル

- 状態  $S$ : A,B,C,D,E,F
- 行動  $A$ : A-E, B-D, B-F, C-D, D-E, E-F, F-F
- 報酬  $R$ : 0,100 <- 行列で表す
- ポリシ  $\pi$ : QTable <- これを学習する
- 値  $V$ : QTable からの値



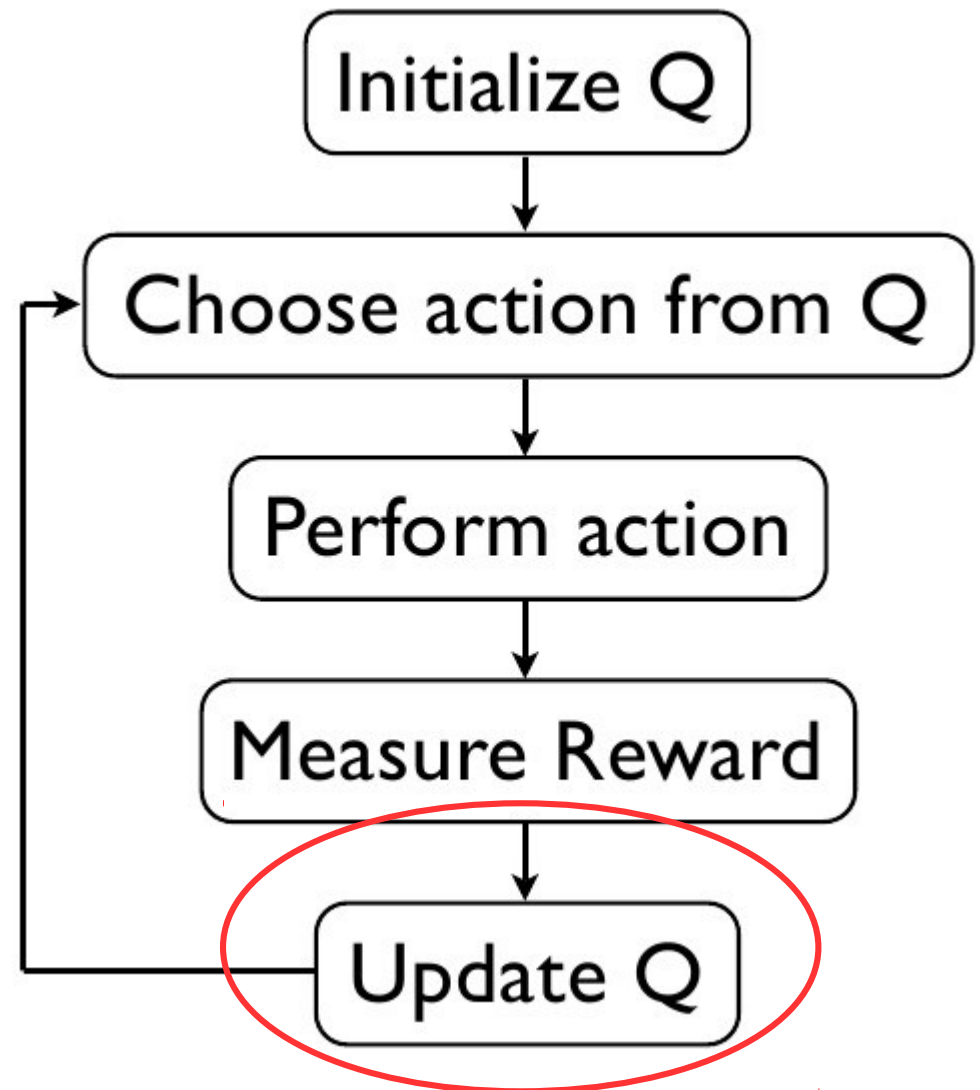
# QTable と R 行列

$$\mathbf{R} = \begin{array}{c|cccccc} state \backslash action & A & B & C & D & E & F \\ \hline A & - & - & - & - & 0 & - \\ B & - & - & - & 0 & - & 100 \\ C & - & - & - & 0 & - & - \\ D & - & 0 & 0 & - & 0 & - \\ E & 0 & - & - & 0 & - & 100 \\ F & - & 0 & - & - & 0 & 100 \end{array} \quad \mathbf{Q} = \begin{array}{c|cccccc} & A & B & C & D & E & F \\ \hline A & 0 & 0 & 0 & 0 & 0 & 0 \\ B & 0 & 0 & 0 & 0 & 0 & 0 \\ C & 0 & 0 & 0 & 0 & 0 & 0 \\ D & 0 & 0 & 0 & 0 & 0 & 0 \\ E & 0 & 0 & 0 & 0 & 0 & 0 \\ F & 0 & 0 & 0 & 0 & 0 & 0 \end{array}$$

# Update 式

$$Q(state, action) = R(state, action) + \gamma \cdot \text{Max}[Q(next\ state, all\ actions)]$$

- $\gamma$  は学習パラメータ (0~1)
- 報酬がすぐか遅れるか



# やってみよう！

$$\mathbf{Q} = \begin{matrix} & \begin{matrix} A & B & C & D & E & F \end{matrix} \\ \begin{matrix} A \\ B \\ C \\ D \\ E \\ F \end{matrix} & \begin{bmatrix} 0 & 0 & 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 & 0 & 0 \end{bmatrix} \end{matrix}$$

- B にいる場合
- ランダムに F に移動
- セッション終了

$$\mathbf{Q}(\text{state}, \text{action}) = \mathbf{R}(\text{state}, \text{action}) + \gamma \cdot \text{Max}[\mathbf{Q}(\text{next state}, \text{all actions})]$$

$$\mathbf{Q}(B, F) = \mathbf{R}(B, F) + 0.8 \cdot \text{Max}\{\mathbf{Q}(F, B), \mathbf{Q}(F, E), \mathbf{Q}(F, F)\} = 100 + 0.8 \cdot 0 = 100$$

$$\mathbf{Q} = \begin{matrix} & \begin{matrix} A & B & C & D & E & F \end{matrix} \\ \begin{matrix} A \\ B \\ C \\ D \\ E \\ F \end{matrix} & \begin{bmatrix} 0 & 0 & 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 & 0 & 100 \\ 0 & 0 & 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 & 0 & 0 \end{bmatrix} \end{matrix}$$

# やってみよう！

$$\mathbf{Q} = \begin{matrix} & \begin{matrix} A & B & C & D & E & F \end{matrix} \\ \begin{matrix} A \\ B \\ C \\ D \\ E \\ F \end{matrix} & \begin{bmatrix} 0 & 0 & 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 & 0 & 100 \\ 0 & 0 & 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 & 0 & 0 \end{bmatrix} \end{matrix}$$

- Dにいる場合
- ランダムに B に移動
- ゴールまで続く

$$\mathbf{Q}(\text{state}, \text{action}) = \mathbf{R}(\text{state}, \text{action}) + \gamma \cdot \text{Max}[\mathbf{Q}(\text{next state}, \text{all actions})]$$

$$\mathbf{Q}(D, B) = \mathbf{R}(D, B) + 0.8 \cdot \text{Max}\{\mathbf{Q}(B, D), \mathbf{Q}(B, F)\} = 0 + 0.8 \cdot \text{Max}\{0, 100\} = 80$$

$$\mathbf{Q} = \begin{matrix} & \begin{matrix} A & B & C & D & E & F \end{matrix} \\ \begin{matrix} A \\ B \\ C \\ D \\ E \\ F \end{matrix} & \begin{bmatrix} 0 & 0 & 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 & 0 & 100 \\ 0 & 0 & 0 & 0 & 0 & 0 \\ 0 & 80 & 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 & 0 & 0 \end{bmatrix} \end{matrix}$$



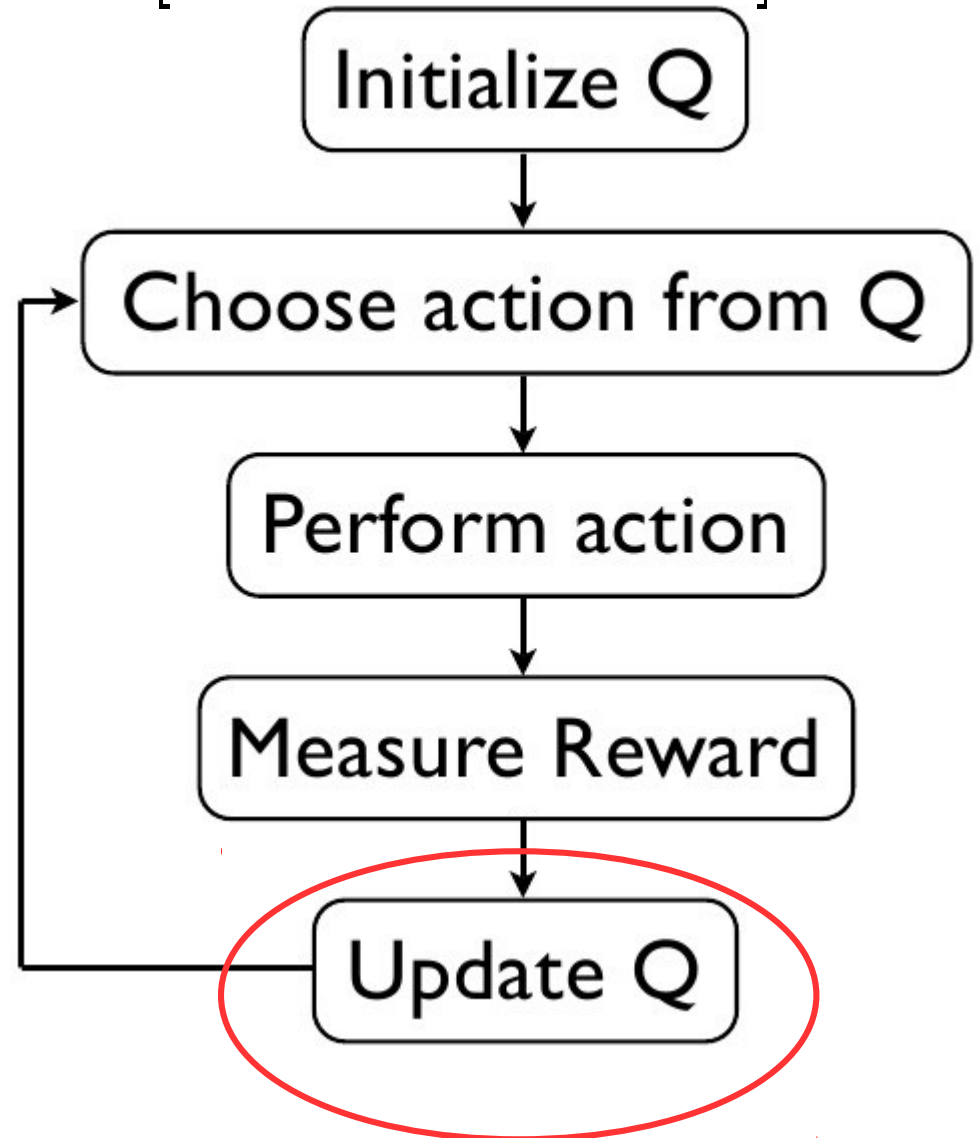
続けると。。。。

$$\mathbf{Q} = \begin{array}{c|cccccc} \textit{state} \backslash \textit{action} & A & B & C & D & E & F \\ \hline A & - & - & - & - & 400 & - \\ B & - & - & - & 320 & - & 500 \\ C & - & - & - & 320 & - & - \\ D & - & 400 & 256 & - & 400 & - \\ E & 320 & - & - & 320 & - & 500 \\ F & - & 400 & - & - & 400 & 500 \end{array}$$

# quiz1

$$Q(state, action) = R(state, action) + \gamma \cdot \text{Max}[Q(next\ state, all\ actions)]$$

- quiz1.py
- この手順を書いてみる
- 難しいなら quiz11.py
  - 式を書いてみる

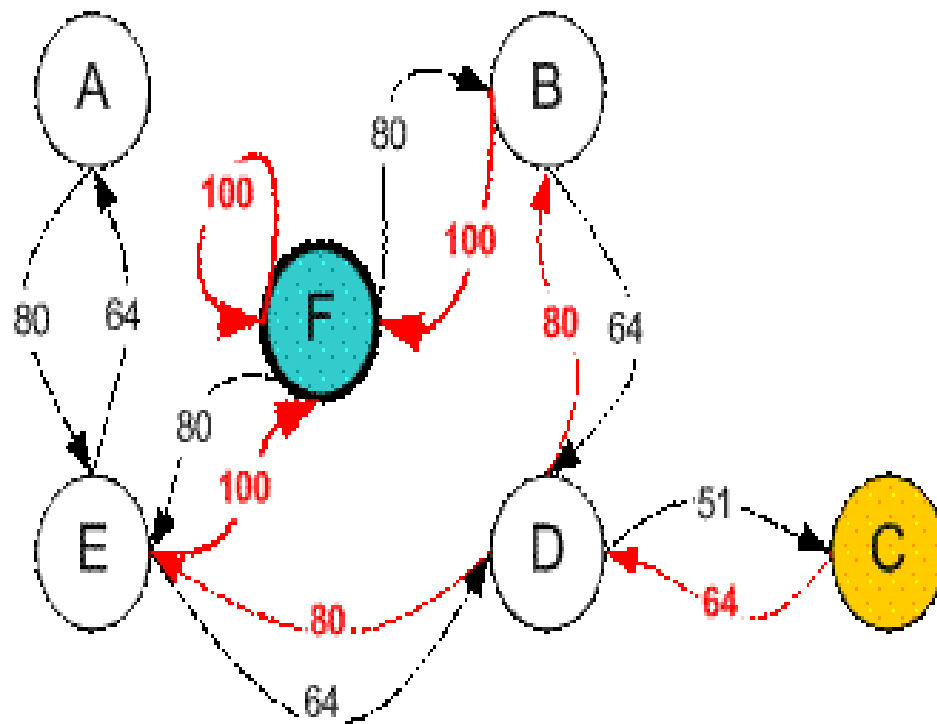


# QTable を実装する

$$\mathbf{Q} = \begin{array}{c|cccccc} \textit{state} \backslash \textit{action} & A & B & C & D & E & F \\ \hline A & - & - & - & - & 400 & - \\ B & - & - & - & 320 & - & 500 \\ C & - & - & - & 320 & - & - \\ D & - & 400 & 256 & - & 400 & - \\ E & 320 & - & - & 320 & - & 500 \\ F & - & 400 & - & - & 400 & 500 \end{array}$$

$$\hat{\mathbf{Q}} = \begin{array}{c|cccccc} \textit{state} \backslash \textit{action} & A & B & C & D & E & F \\ \hline A & - & - & - & - & 80 & - \\ B & - & - & - & 64 & - & 100 \\ C & - & - & - & 64 & - & - \\ D & - & 80 & 51 & - & 80 & - \\ E & 64 & - & - & 64 & - & 100 \\ F & - & 80 & - & - & 80 & 100 \end{array}$$

# QTable を実装する



$\hat{Q} =$

<i>state \ action</i>	<i>A</i>	<i>B</i>	<i>C</i>	<i>D</i>	<i>E</i>	<i>F</i>
<i>A</i>	–	–	–	–	80	–
<i>B</i>	–	–	–	64	–	100
<i>C</i>	–	–	–	64	–	–
<i>D</i>	–	80	51	–	80	–
<i>E</i>	64	–	–	64	–	100
<i>F</i>	–	80	–	–	80	100

# quiz2

- quiz2.py
- QTable を使う実装プログラムを作ってみる。

# quiz3

- quiz3.py
- 報酬 R を書いてみよう。

