

BLEU Score [Bilingual Evaluation Understudy]

→ A method for Automatic Evaluation of Machine Translation.

→ BLEU is a score for comparing a candidate translation of text to one or more reference translations (human translation).

→ It can also be used to evaluate text generated for a suite of NLP tasks (like language generation, image caption generation, text summarization, speech recognition).

→ BLEU is a metric for evaluating a generated sentence (after translation) to a reference sentence.

→ BLEU score can be used to evaluate the performance of a Machine Translation system particularly when there are multiple equally good answers.

→ BLEU metric ranges from 0 to 1.

(range) Score = 1.0 - A perfect match

Score = 0.0 - A perfect mismatch results

→ BLEU score was developed for evaluating the predictions made by automatic machine translation system, but it is not perfect evaluation method.

* Benefits of BLEU score

- 1) Quick & inexpensive to calculate
- 2) Easy to understand
- 3) Language independent
- 4) It correlates highly with human evaluation
- 5) Widely adopted method.

→ Machine translation evaluation system require

- i) A numerical translation closeness metric
- ii) A good quality human reference translations

→ BLEU score is based on "text string matches".

OR
⇒ It is a score that quantifies how good machine translation is by computing a



Similarity score based on N-gram
precisely

(N gram) Precision - look at each word in the candidate sentence & assign it a score of 1 if it shows up in any of the reference sentence otherwise 0.

Ex NT output - she cat the cat on the mat
Reference 1: There is a cat on the mat

Reference 2: The cat is on the mat

N gram: - Sequence of N consecutive words
Like

i) Unigram - one word Ex -> hello

ii) bigram - two word Ex -> how are

iii) trigram - three word sequence of words
Ex -> hello how are

iv) 4-gram - four word seq. of words
Ex -> hello how are you

The main task of BLEU implements is to compare N-gram of the candidate with two N-grams of the reference translation & count the number of matches

→ Here the nos. of matches are position independent



Calculation

Behaviour

→ If there are more nos. of matches b/w candidate & reference translation, then better is the translation

→ Higher the BLEU score better is the translation

→ More reference human translation -> more results in better & accurate scores.

Remaining Next Day

UNIT-4



~~★~~ Unigram Precision

Unique Unigram	Count / Score
the	3
Cat	2
on	1
mat	1
<hr/>	
Total counts for the unique unigram in the candidate sentence	= 7 (It should be b/w 0 & 1)

Candidate Sentence

the cat the cat on the mat
 1 2 3 4 5 6 7

Nos. of unigram in candidate sentence = 7

$$\text{Unigram precision} = \frac{7}{7} = 1.0 \text{ (normalized)}$$

7 \rightarrow total nos. of words in C.S.

ii) Bigram Precision

Candidate Sentence

the cat the cat on the mat

Bigram - the cat, ~~the~~ cat on, the mat = 5
 1 2 3 4 5

Uniqe Pigum	2
the cat	1
Cat the	1
cat on	1
on the	1
<u>the mat</u>	<u>4</u>
- total cost	6

$$\text{Brgum precision} = \frac{5}{6} = 0.833$$

A Daaback

A drawback
Chinese precision is very easy to be
over-confident about the quality of
machine translation.

Example of poor MT output with high precision

Candidate Sentence : the the the the the the the
 ngram = $\frac{7}{7} = 1.0$
 precision

→ To overcome the unigram precision problem, clipped count & modified precision were proposed

~~A~~ Modified N-gram precision

Candidate - the cat the cat on the mat
Reference 1 - there is a cat on the mat
Reference 2 - the cat is on the mat

- ↳ Unique modified 0-gram precision

Reference Unigram	Count	Clipped Count
the	3	2
Cat	2	1
on	1	1
mat	1	1
total	7	5
clipped		
Count		

$C_{max} < \text{Count}$
 Clipped Count
 $= C_{max}$

Column modified precision = $\frac{5}{7} = 0.714$

	Unique Bigram	Count	Clipped Count	Bigram NP
Release the cat		2	1	
Release the		1	0	
the cat		1	0	
				$= 4$

$$\frac{4}{9} = 0.444\ldots$$

iii) Modified N-gram precision

Candidate Sentence: the the the the the the
Reference Sentence: there is a cat on the mat
the cat is on the mat

Unique N-gram modified precision

Unique N-gram	Count	Cropped Count	Rel/Max
the	7	2	<u>2</u>

$$\text{Precision} = \frac{2}{7}$$

A Problem in Modified N-gram precision

→ BLEU scores tend to favor short translations, which can produce very high precision scores, even using modified precision

Ex

Candidate Sentence - the cat
Reference Sentence - the cat is on the mat

$$\text{Modified unigram precision} = \frac{|H|}{2} = 1$$

highest

A Solution is to Penalty Penalty

$$BP = \begin{cases} 1 & \text{if } c > r \\ c^{(1-x/2)} & \text{if } c \leq r \end{cases}$$

$r \rightarrow$ length of reference sentence
 $c \rightarrow$ length of candidate sentence

For above example
 $C=2, r=6$
 as $C < r$ $BP = c^{(1-6/2)} = c^{(-3)} = c^{-2}$

$$\text{Bleu score} = BP \times \exp \left[\frac{1}{N} \sum_{n=1}^N \ln p_n \right]$$

\exp exponential
 $\sum_{n=1}^N \ln p_n$ Geometric Mean of all N-gram precision

$P_n \rightarrow$ Modified precision for n-gram
 if using bigram, unigram, digram, N=3

→ BLEU only utilizes a BP for shorter sentence

→ the standard & longer sentences are not penalized properly

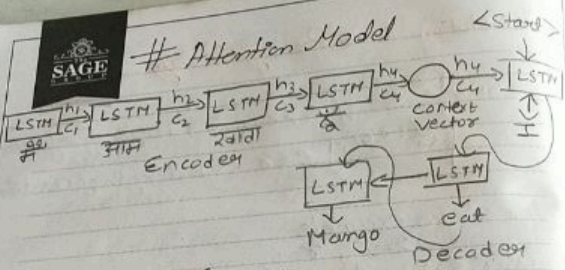


Fig → Translation using Encoder-Decoder

* Context Vector (h_4, c_4) → Contains good Summary of Input Sequence (fixed size)

- For long sentences (input) context vector may not have all necessary information.
- To solve this problem, attention mechanism is used with Encoder-Decoder, RNN/LSTM

Attention Model

→ Attention means directing your focus at something & taking greater notice

OR Concentrating on one or a few things while ignoring others.


- The attention mechanism is an improvement over the encoder-decoder based on neural based machine translation system in NLP
- In encoder-decoder RNN/LSTM
 → Encoder process the entire input sentence & encode it into a context vector which is the last hidden state of LSTM/RNN

Ex मैं आम खाता हूँ
 t_1 मैं [1 0 0 0] - I
 t_2 खाता हूँ [0 0 0 0] - eat
 t_3 आम [0 1 0 0] - Mango

α_{ij} : i → importance of word
 j → timestamp

	मैं	आम	खाता	हूँ
t_1	1	0	0	0
α_{11}	α_{21}	α_{31}	α_{41}	
t_2	0	0	0.5	0.5
α_{12}	α_{22}	α_{32}	α_{42}	

α_{12} → This word of the sequence at time stamp 2



 t_3 $\begin{matrix} \alpha_1^3 & \alpha_2^3 & \alpha_3^3 & \alpha_4^3 \end{matrix}$

→ It is expected that context vector contains good summary of input sequence & information of the whole sentence. If the intermediate states of the encoder are ignored, & the final state is supposed to be the initial hidden state of the decoder.

→ Decoder units produce the words in a sentence one after another.

* Drawback of Encoder-decoder Machine Translation system

→ If the encoder makes a bad summary the translation will also be bad.

→ Usually encoder creates a bad summary when it tries to understand a large sentence.

→ RNNs cannot remember long sentence & sequences due to vanishing/exploding

gradient problem. It can remember the words which it has just seen.

→ So the performance of the encoder-decoder network degrades rapidly as the length of the input sentence increases.

→ LSTM is supposed to do better than RNN, but it became forgetful in specific cases.

→ But in both RNN & LSTM there is no mechanism to give more importance to some of the input words compared to others while translating the sentence.

→ So when model generates a sentence (machine translation) it searches for a set of positions in the encoder hidden states where the most relevant information is available. This idea is called **Attention**.

* Attention is proposed as a solution to the limitation of encoder-decoder model encoding the input sequence to one fixed length vector from which to decode each output time step. This issue is believed to be more of a problem when decoding long sequences.

→ The attention model allows an RNN to pay attention to specific parts of the input that is considered as being important which improves the performance of the resulting model in practice.

→ Instead of sending last hidden state output of decoder, we could take a weighted average of the corresponding word representation & feed it to the decoder.

→ For ex, at timestamp 2, we can take a weighted average of the representation of words to get output 'eat' from the decoder.

By doing this we are not overloading the decoder with irrelevant information (like # 31111)

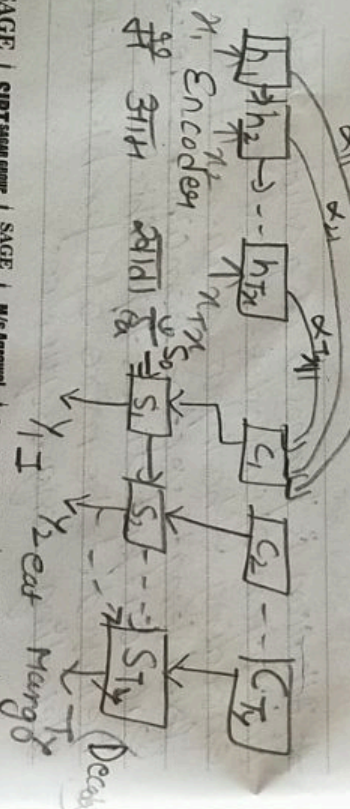


Fig → Encoder-Decoder Model with Attention Mechanism

Prediction of T word

$T=2$

$$C_2 = \alpha_{12} \cdot h_1 + \alpha_{22} \cdot h_2 + \alpha_{32} \cdot h_3 + \alpha_{42} \cdot h_4$$

$$\alpha_{12} = a(s_1, h_1) \quad \alpha_{12} = \frac{e^{\theta_{12}}}{e^{\theta_{12}} + e^{\theta_{22}} + e^{\theta_{32}} + e^{\theta_{42}}}$$

$$\alpha_{22} = a(s_1, h_2) \quad \alpha_{22} = \frac{e^{\theta_{22}}}{e^{\theta_{12}} + e^{\theta_{22}} + e^{\theta_{32}} + e^{\theta_{42}}}$$

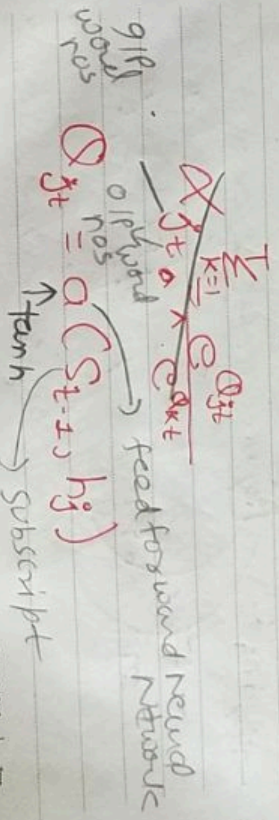
$$\alpha_{32} = a(s_1, h_3) \quad \dots$$

$$C_t = \sum_{j=1}^T \alpha_{jt} \cdot h_j$$

Context vector

h_j = encoder hidden state at j^{th} timestamp

α_{jt} - Probability of focusing on the j^{th} word to produce the word.



$$X_{jt} = \frac{\exp(Q_{jt})}{\sum_{k=1}^K \exp(Q_{kt})}$$

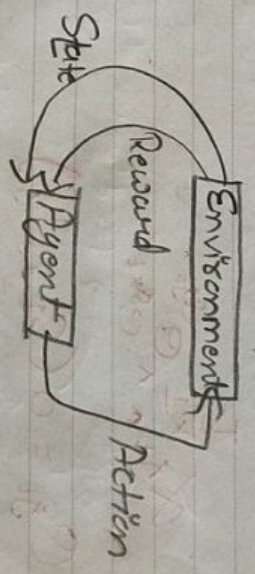
→ The importance of the j^{th} input word for deciding the t^{th} output word can be calculated by

$$Q_{jt} = f_{HTT}(S_{t-1}, h_j) \text{ or } a$$

A/f_{HTT} = feed forward Neural Network

The alignment model scores (Q) how well back encoded input (h) matches the current output of the decoder (S)

Reinforcement Learning

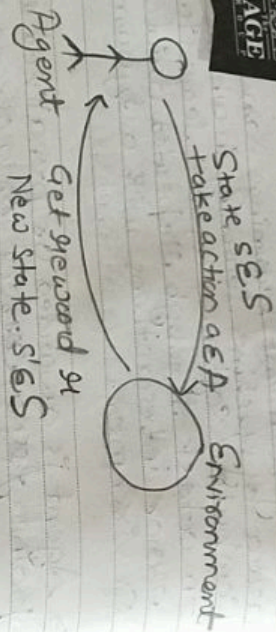


→ Reinforcement learning is an area of machine learning in which computer/software agent learns to perform a task through repeated trial & error interactions with a dynamic environment

→ Agent is awarded or penalized with a point for a correct or a wrong answer
 → On the basis of the positive reward points gained the model trains itself.

→ The agent is acting in an environment. The agent can stay in one of many states (S) of the environment & choose to take one of many actions (A) to switch from one state to another. Which state the agent will arrive is decided by transition probabilities between states (P).

→ Once an action is taken, the environment delivers a reward (r) or feedback



→ In RL agent, interacts with the environment trying to take smart actions to maximize cumulative rewards

* Goal of RL

Find an optimal policy that maximize total cumulative rewards of the agent

The interaction b/w the agent & the environment involves a sequence of actions & observed rewards in time $t=1, 2, \dots, T$.

→ During the process the agent accumulates the knowledge about the environment learns the optimal policy & makes decisions on which action to take next so as to efficiently learn the best policy

→ RL used in video games, computer games etc.

* Terms used in Reinforcement learning

1) Agent → An entity that takes action

2) **Action (A)** → It is the move that an agent makes in a given state in the environment

→ Usually there is a list of discrete possible actions that agent chooses (move left/right, up, down, jump etc).

A → Set of all actions
 a → a particular action

3) **State (S)** → A state is used to represent current situation of the environment from the agent's view (like obstacles, enemies etc).

S → Set of all the states
 s → a particular state.

4) **Reward (R)** → R reward is the immediate feedback given by environment to the agent for its action in the given state.

marco touches a coin - the coin's point
marco touches enemy - he dies

R → total or cumulative reward
r → immediate rewards received after action a...

5) Environment :-
The physical world in which the agent interacts.

→ Environment takes the agent's current state & action as input & returns agent reward & its next state as output.

6) Policy (π) :-
Method to map agent's state to actions or the strategy that the agent employs to determine its action based on the current state.

→ The agent's policy $\pi(s)$ provide the guideline on what is the optimal action to take in a certain state with the goal to maximize the total rewards.

7) Value - Future reward that an agent could receive by taking an action in a

particular state

{ total amount of reward the agent expect
cumulative reward = winning the game }

Ex PacMan game
Agent → Pacman

Environment - grid world
States - location of Pacman in the grid world
Reward - eating food
Punishment (Killed by ghost) - loses the game.

Markov Decision Process (MDP)

→ MDP is a mathematical framework that can solve most Reinforcement learning problems with discrete actions.

→ With MDP, an agent can assure at an optimal policy for maximum rewards over time.

→ Any reinforcement learning task composed of a set of states, actions & rewards that follows the Markov property would be considered as Markov Decision Process

→ To make a decision about what action to take at time $t+1$, we only have to consider the environment at time t & no more (Markov property)

* Goal:

The aim of MDP is to train an agent to find a policy that will return the maximum cumulative rewards from taking a series of actions in one or more states

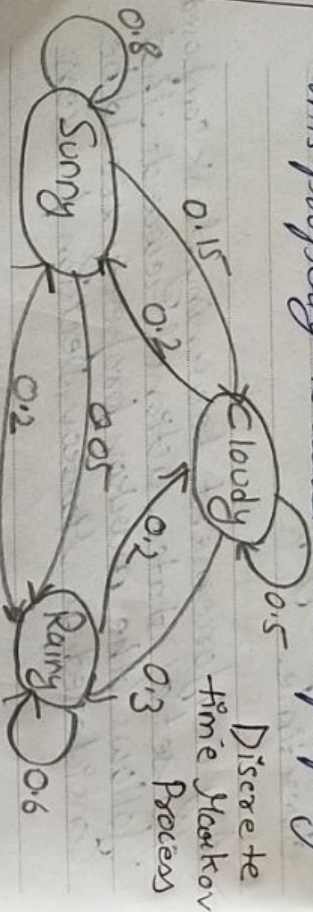
$$\Pi(GS) : S \rightarrow A$$

* Markov Process:-

Given entire history of states (S_0, S_1, \dots, S_t) the future state (S_{t+1}) depends only on present state (S_t)

$$P[S_{t+1} / S_0, S_1, \dots, S_t] = P[S_{t+1} / S_t]$$

This property is called Markov property



Let state S_t be weather on any given day. Tomorrow's weather depends on today's weather

One step transition probability matrix P

S_t	S_{t+1}		
	Sunny	Cloudy	Rainy
Sunny	0.8	0.15	0.05
Cloudy	0.2	0.5	0.3
Rainy	0.2	0.2	0.6
	-1	-1	-1

→ At each time step, the process is in some state S_t

→ The agent may choose any action A_t available in state S_t

→ The process responds at the next time step by moving into state S_{t+1} & giving the agent a corresponding reward R_t, A_t, S_{t+1}

→ The state S_t captures all relevant information from the history $(S_0, S_1, \dots, S_{t-1})$. Once the

state S_t is known the history may be thrown away.

→ As MDP satisfies Markov property the probability of transitioning from state S_t to S_{t+1} depends only on S_t & a_t

$$\begin{aligned} P[S_{t+1} / S_0, \dots, S_t, a_0, a_1, \dots, a_t] \\ = P[S_{t+1} / S_t, a_t] \\ = P[S_{t+1} / S_t] \end{aligned}$$

* State Transition Matrix

For a Markov state s & successor state s' , the state transition probability is defined by

$$P_{ss'} = P[S_{t+1} = s' / S_t = s]$$

State transition Matrix P defines transition probabilities from all states S to all successor states S' , where each row of the matrix sums to 1

$$P = \begin{bmatrix} P_{11} & P_{12} & \dots & P_{1n} \\ P_{21} & P_{22} & \dots & P_{2n} \\ \vdots & \vdots & \ddots & \vdots \\ P_{m1} & P_{m2} & \dots & P_{mn} \end{bmatrix}$$

* Markov Decision Process -

→ That means an agent is supposed to decide the best action to select based on his current state & when this step is repeated, the problem is known as MDP

→ The MDP process is used to model an environment for an agent to learn

→ It is mathematical representation of the environment

* Parts of MDP model - $\langle S, A, T, R, \gamma \rangle$

i) States (S) :- A state is used to represent the current situation of the environment from the agent's view.

S - set of possible states

ii) Action $A(s)$:- $A \rightarrow$ set of possible actions. An action is how an agent gets from one state to another state s'

iii) Transition function $T(s, a, s')$ or $P(s' / s, a)$
The transition function is the probability that an agent gets to the next

state (s') given that he takes an action (a) from his current state (s)

OR This is the probability of getting to state s' given that we take action a in state s

iv) Rewards

This function rewards an agent for taking the right actions & punishes (with negative rewards) for wrong actions.

$R(s, s')$ - the reward for moving from state s to s'
or $R(s, a)$

The purpose MDP is to find a policy $\pi(s)$ this policy specifies the action to be taken by the agent when in a state s

Discount factor

v) γ - discount factor

Help us to evaluate the expected reward relative to the advantage or disadvantage of each state. For ex-1 Given the state gain in assuming 9 take

The best possible action now & at each subsequent step what long term reward can we expect

Bellman Equations:-

Firstly we know that goal of RL is given the current state we are in, choose the optimal action which will maximize the long-term expected reward provided by the environment

\Rightarrow Bellman Equation helps to solve MDP or we can say it helps in finding optimal policy & value function.

i) Value Function

\rightarrow Each state is associated with a value function $V(s)$ predicting the expected amount of future rewards we are able to receive if this state by acting the corresponding policy.

\rightarrow It is equal to expected total reward for an agent starting from state s

$$V(s) = E[R/s_t = s]$$

OR $V(s) = E[R_{t+1} + \gamma V(s_{t+1}) / s_t = s]$

OR $V_{\pi}(s) = E[R / s_t = s]$

π - Policy, γ - discount factor

OR $V_{\pi}(s) = E[\sum_{t=0}^{\infty} \gamma^t r_{t+1} / s_0 = s]$

→ The value function estimates value of a state

→ It depends on the policy by which the agent picks action to perform

→ The state with high V is more valuable than the state with low V

γ → discount factor - a factor that is multiplied to future reward to get its present value

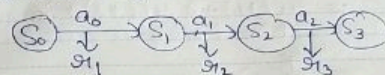
Usually $0 \leq \gamma \leq 1$

Present reward = $\gamma^k \times$ future reward after k steps

Ex → 10 points after 2 timesteps ⇒

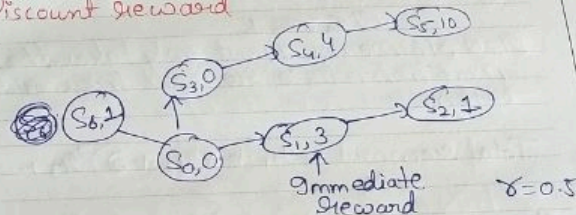
Present reward = $(0.5)^2 \times 10$

At time $t=0$



Present Value of total reward at $t=0$ for $S_0 = r_1 + \gamma r_2 + \gamma^2 r_3 + \gamma^3 r_4 + \dots$

Discount reward



$S_0 \rightarrow S_1 \rightarrow S_2$ total reward = $3 + 1 = 4$

$S_0 \rightarrow S_3 \rightarrow S_4 \rightarrow S_5$ total reward = $0 + 4 + 10 = 14$

(Discount reward) Cumulative reward of state

$S_0 = 0 + 0.5 \times 4 + 0.5^2 \times 10$

So, the discounted return at time step t is given by

total discount reward
 $= r_1(t+1) + \gamma r_1(t+2) + \gamma^2 r_1(t+3) + \dots$
 $= \sum_{i=0}^{\infty} \gamma^i r_{i+1}$

Goal of agent is to maximize expected (discounted) return

$0 \leq \gamma \leq 1$
 $\gamma=0$, it means agent only takes the most immediate reward into consideration

Total reward $= r_1 + \gamma r_{12} + \gamma^2 r_{13} + \dots$
 $= r_1$

If $\gamma=1$, the return is without discount

Total reward $= r_1 + \gamma r_{12} + \gamma^2 r_{13} + \dots$
 $= r_1 + r_{12} + r_{13} + \dots$

For large value of γ , the agent takes care of more distant future to take more rewards.

Value function $V_{\pi}(s) = E \left[\sum_{i=0}^{\infty} \gamma^i r_{i+1} \mid s \right]$

Markov decision process objective is to find optimal policy where all states have max value function

$V^*(s) = \max_{\pi} V_{\pi}(s)$

★ Bellman Equation

→ It helps to solve Markov decision process or we can say it helps in finding optimal policy & value function

For Deterministic Environment
 Bellman equation is given by

π -policy $V_{\pi}(s) = \max_a \{ R(s,a) + \gamma V_{\pi}(s') \}$

Bellman equation answers the following question

Given the state I am in, assuming I take the best possible action now & at each subsequent step, what long term reward can I expect?
OR

What is the value of the state?

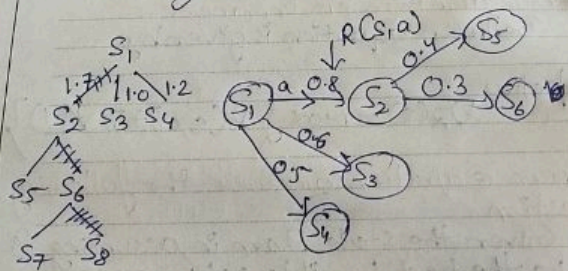
where

$V(s)$ = value for being in a certain state s

$V(s')$ - is the value of being in the next state after taking an action a

γ → discount factor

\max_a → Here we want Maximization of action (which maximizes the value) so that agent is in the optimal state



$$V(s_1) = \max \begin{cases} 0.8 + \gamma V(s_2) \\ 0.6 + \gamma V(s_3) \\ 0.5 + \gamma V(s_4) \end{cases}$$

$$V(s_2) = \max \begin{cases} 0.4 + \gamma V(s_5) \\ 0.3 + \gamma V(s_6) \end{cases}$$

Bellman equation refers to a set of equations that decompose the value function into immediate reward plus the discounted future values

For ^{non} deterministic / Stochastic environment

Bellman equations

$$V(s) = \max_a \left\{ R(s, a) + \gamma \sum_{s'} P(s, a, s') V(s') \right\}$$