

Универзитет Црне Горе  
Природно-математички факултет

---

Математички софтверски пакети

# Примјена Active Inference-а на проблем Cliff Walking

**Студент:**  
Никола Поповић 1/24 Б

**Ментор:**  
мр Никола Пижурица

мај 2025.

# Садржај

Увод	1
<b>1 Опис проблема и окружења</b>	<b>2</b>
<b>2 Имплементација Active Inference агента</b>	<b>3</b>
2.1 Стандардни модел . . . . .	3
2.1.1 Учитавање библиотека. Визуелизација окружења и увјерења. .	3
2.1.2 Иницијализација окружења. Дефинисање основних параметара и структуре стања . . . . .	6
2.1.3 Генеративни модел . . . . .	7
2.1.4 Функције за инференцију и рачунање очекиване слободне енергије	10
2.1.5 Понашање агента . . . . .	16
2.2 Модел са шумом у матрицама A и D . . . . .	17
2.3 Модел са клизавом површином (is_slippery=True) . . . . .	21

## Увод

У овом раду представљена је имплементација агента заснованог на принципу Active Inference у окружењу Cliff Walking, доступном у оквиру библиотеке Gymnasium. Циљ агента је да научи како да безбједно дође до циљног поља, избегавајући опасне провалије, ослањајући се искључиво на свој генеративни модел свијета и принцип минимизације очекиване слободне енергије.

У оквиру рада реализована је имплементација агента који не учи структуру свијета кроз интеракцију, већ почиње са унапријед дефинисаним генеративним моделом: матрицама  $A$  (вјероватноћа запажања у зависности од скривених стања),  $B$  (матрица транзиција стања),  $C$  (преференције) и  $D$  (почетна увјерења). Агент планира низ акција (политике) и бира ону која минимизује очекивану слободну енергију.

Поред основне верзије агента, дата је имплементација у још два сценарија:

1. када је додат шум у матрицама  $A$  и  $D$ , и
2. када окружење постаје стохастичко (укључивањем `is_slippery` параметра).

# 1 Опис проблема и окружења

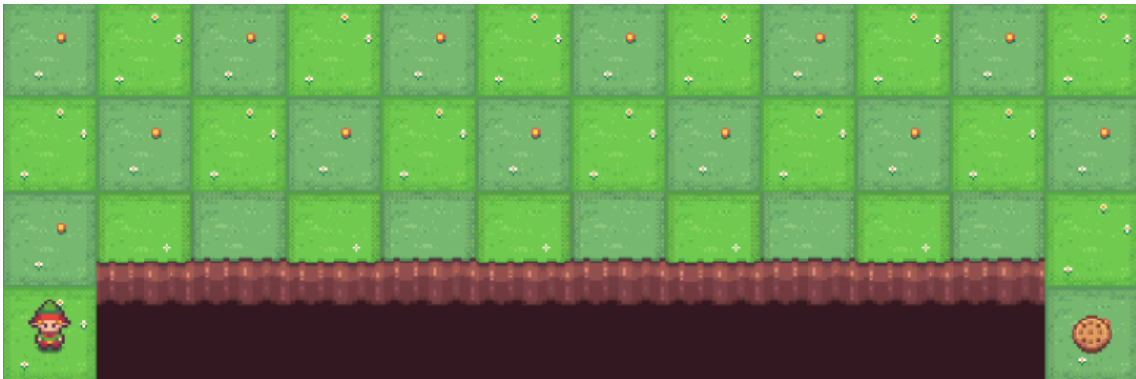
Окружење представља мрежу величине  $4 \times 12$ , у којој агент почиње у позицији  $[3, 0]$ , односно у доњем лијевом углу. Циљ агента је да достигне позицију  $[3, 11]$ , која се налази у доњем десном углу мреже. Када агент достигне циљ, епизода се завршава.

Дуж реда  $[3, 1]$  до  $[3, 10]$  налази се провалија. Ако агент крочи на неко од ових поља, одмах се враћа на почетну позицију. Агент извршава акције док не достигне циљ, али уколико направи грешку и упадне у провалију, мора поново да крене од почетка.

Простор акција агента састоји се од четири могућа покрета: помјерања горе, десно, доље или лијево, које су означене бројевима од 0 до 3. При томе, постоји могућност да окружење буде клизаво (is\_slippery режим, подразумевано искључен), што значи да агент понекад може да се помјери у неком другом смјеру, који се налази нормално у односу на онај који је наумио.

Сваки корак агента носи награду од  $-1$ , док је казна за корак у провалију значајно већа и износи  $-100$ . Епизода се завршава у тренутку када агент стигне на циљну позицију.

Простор опсервација садржи  $3 \times 12 + 1 = 37$  могућих стања. Агент се не може налазити на провалији, нити на циљну позицију јер долазак на циљ прекида епизоду. Тренутна позиција агента представљена је као један цио број, израчунат по формули: тренутни\_ред  $\times$  број\_колона  $+$  тренутна\_колона, гдје индексирање реда и колона креће од 0. На примјер, почетна позиција  $[3, 0]$  одговара индексу  $3 \times 12 + 0 = 36$ .



Слика 1.1: Почетно стање агента у Cliff Walking окружењу

## 2 Имплементација Active Inference агента

### 2.1 Стандардни модел

У овој фази рада фокусираћемо се на класичан модел окружења за Active Inference агента, где су матрице вјероватноћа запажања у зависности од скривених стања и почетних уверења детерменистичке, односно без шума. То значи да матрица  $A$ , која описује вероватноћу посматрања датог стања, и матрица  $D$ , која представља агентову почетну претпоставку о свом положају, садрже само сигурне и јасне вриједности без случајних варијација.

Поред тога, параметар клизавог окружења `is_slippery` је (подразумијевано) у овом моделу искључен, што значи да се агент помјера строго у правцу изабране акције.

#### 2.1.1 Учитавање библиотека. Визуелизација окружења и уверења.

Испод је приказан Python код за учитавање потребних библиотека:

```
1 import gymnasium as gym
2 import numpy as np
3 import matplotlib.pyplot as plt
4 import seaborn as sns
5 import itertools
6 from pymdp import utils
7 from pymdp.maths import softmax, spm_log_single as log_stable
8 from pymdp.control import construct_policies
```

Код 2.1.1: Учитавање библиотека

За боље разумијевање процеса учења и понашања агента, имплементирани су функције за визуелизацију. Оне су преузете из `Active inference from scratch` и димензије су прилагођене нашем проблему.

```
1 def plot_likelihood(matrix, xlabel=None, ylabel=None,
2                     title_str="Likelihood distribution (A)":
3     """
4     Plots a 2-D likelihood matrix as a square heatmap with a controlled
5     plot size.
6     """
7     if not np.isclose(matrix.sum(axis=0), 1.0).all():
8         raise ValueError("Distribution not column-normalized!
9                             Please normalize")
```

```

9             (ensure matrix.sum(axis=0) == 1.0 for all
10             columns)")
11
12     if xlabels is None:
13         xlabels = list(range(matrix.shape[1]))
14     if ylabels is None:
15         ylabels = list(range(matrix.shape[0]))
16
17     # Adjust the figure size manually for more controlled sizing
18     max_size = 8 # Maximum figure size (in inches)
19     size = max(matrix.shape) # Get the larger of (rows, cols)
20
21     # Use a fixed figure size to prevent the plot from being too big
22     # Scale the figure size based on the matrix shape but constrain the
23     # size
24     fig_size = (min(size, max_size), min(size, max_size))
25
26     plt.figure(figsize=fig_size) # Set the figure size
27
28     sns.heatmap(matrix, xticklabels=xlabels, yticklabels=ylabels,
29                 cmap='gray', cbar=False, vmin=0.0, vmax=1.0, square=True)
30
31     plt.title(title_str)
32
33     plt.show()
34
35 def plot_grid(grid_locations, num_x=4, num_y=12):
36     """
37     Plots the spatial coordinates of GridWorld as a heatmap, with each
38     (X, Y) coordinate labeled with its linear index (its 'state id').
39     """
40     grid_heatmap = np.zeros((num_x, num_y))
41
42     for linear_idx, location in enumerate(grid_locations):
43         y, x = location
44         grid_heatmap[y, x] = linear_idx
45
46     sns.set(font_scale=1.5)
47
48     sns.heatmap(grid_heatmap, annot=True, cbar = False,
49                 fmt='.0f', cmap='crest')
50
51     plt.show()
52 def plot_point_on_grid(state_vector, grid_locations):

```

```

52     """
53     Plots the current location of the agent on the grid world
54     """
55
56     state_index = np.where(state_vector)[0][0]
57
58     y, x = grid_locations[state_index]
59
60     grid_heatmap = np.zeros((4,12))
61
62     grid_heatmap[y,x] = 1.0
63
64     sns.heatmap(grid_heatmap, cbar = False, fmt='.0f')
65
66     plt.show()
67
68
69 def plot_beliefs(belief_dist, title_str="Belief Distribution"):
70     """
71     Plot a categorical distribution or belief distribution, stored in
72     the 1-D numpy vector 'belief_dist'.
73     """
74
75     if not np.isclose(belief_dist.sum(), 1.0):
76         raise ValueError("Distribution not normalized!
77                             Please normalize")
78
79     plt.figure(figsize=(12, 4))
80
81     plt.grid(zorder=0)
82
83     plt.bar(range(belief_dist.shape[0]), belief_dist, color='r',
84             zorder=3)
85
86     plt.xticks(range(belief_dist.shape[0]))
87
88     plt.title(title_str)
89
90     plt.show()

```

Код 2.1.2: Функције за цртање графика

### 2.1.2 Иницијализација окружења. Дефинисање основних параметара и структуре стања

У овом дијелу креирамо окружење „CliffWalking-v0“ из библиотеке Gymnasium са режимом приказа који враћа RGB низ пиксела. Функција `display_frame` служи за приказивање тренутног стања окружења као слике, гдје се искључују осе и додаје наслов ради боље визуализације, а команда `plt.show(block = True)` осигурава да се прозор са сликом приказује интерактивно и да се програм не настави док корисник не затвори тај прозор.

```
1 env = gym.make("CliffWalking-v0", render_mode="rgb_array")
2 obs, info = env.reset()
3
4 def display_frame(frame, title_str=""):
5     plt.figure(figsize=(10,4))
6     plt.imshow(frame)
7     plt.axis("off")
8     plt.title(title_str)
9     plt.show(block=True)
```

Код 2.1.3: Иницијализација и функција за приказ окружења

Овдје је дат код који описује структуру мреже стања (grid world), број стања и опажања, као и скуп акција које агент може предузети. Ове основне компоненте представљају темељ за изградњу модела и даљу симулацију.

```
1 grid_locations = list(itertools.product(range(4), range(12)))
2
3 n_states = len(grid_locations)
4 n_observations = len(grid_locations)
5
6 actions = ["UP", "RIGHT", "DOWN", "LEFT"]
```

Код 2.1.4: Дефинисање стања, опажања и акција



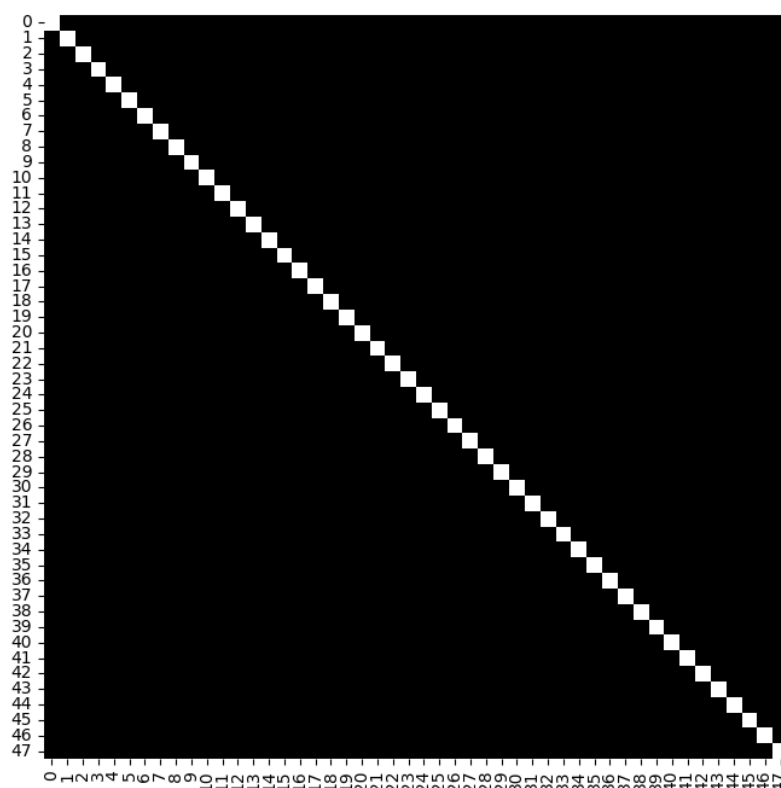
### 2.1.3 Генеративни модел

У наставку формулишемо и описујемо основне компоненте генеративног модела које карактеришу понашање окружења и унутрашња уверења агента.

Матрица  $A$  представља условне вероватноће опсервација у зависности од тренутног стања агента. Када је матрица  $A$  јединична, то значи да свака опсервација јасно и директно одговара једном стању, без икакве неодређености или шума. Другим ријечима, ако се агент налази у стању  $s_i$ , он са сигурношћу 1 опсервира стање  $o_i$ . Оваква поставка поједностављује модел, јер нема нејасноћа у посматрањима — свака позиција је прецизно идентификована.

```
1 A = np.eye(n_observations, n_states)
2
3 plot_likelihood(A)
```

Код 2.1.5: Дефинисање матрице  $A$  и њен графички приказ



Слика 2.1: График вјероватноћа запажања у зависности од скривених стања

Матрица транзиције В дефинише како се стања мијењају као резултат извршене акције. За сваку акцију и тренутно стање, ова матрица садржи вјероватноће преласка у следеће стање. Ако акција води ван граница дефинисане мреже, позиција остаје непромијењена. Посебан случај представљају провалије на дну мреже, гдје је агент уколико ступи у ту зону враћен на почетну позицију. Важно је напоменути да су у овом моделу преласци потпуно детерминистички: ако агент одабере неку акцију, она ће засигурно бити извршена, осим ако не води ван оквира мреже.

```
1 def create_B_matrix():
2     B = np.zeros( (len(grid_locations), len(grid_locations),
3                   len(actions)) )
4
5     for action_id, action_label in enumerate(actions):
6
7         for curr_state, grid_location in enumerate(grid_locations):
8
9             y, x = grid_location
10
11             if (y == 3 and 1 <= x <= 10):
12                 x = 0
13                 y = 3
14
15             if action_label == "UP":
16                 next_y = y - 1 if y > 0 else y
17                 next_x = x
18             elif action_label == "DOWN":
19                 next_y = y + 1 if y < 3 else y
20                 next_x = x
21             elif action_label == "LEFT":
22                 next_x = x - 1 if x > 0 else x
23                 next_y = y
24             elif action_label == "RIGHT":
25                 next_x = x + 1 if x < 11 else x
26                 next_y = y
27
28             new_location = (next_y, next_x)
29             next_state = grid_locations.index(new_location)
30             B[next_state, curr_state, action_id] = 1.0
31     return B
32
33
34 B = create_B_matrix()
```

Код 2.1.6: Дефинисање матрице В

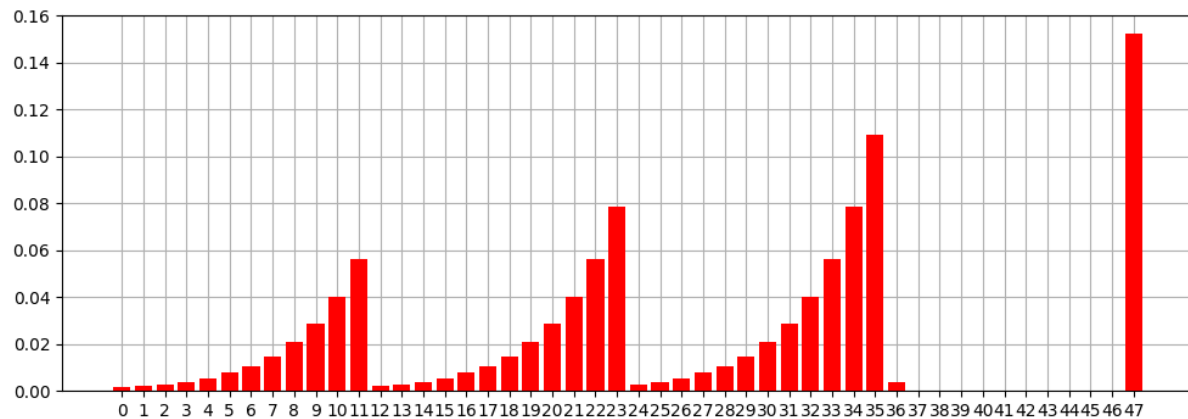
Вектор преференција  $C$  одражава жељена стања, односно циљеве агента. Вриједности су засноване на негативној скалираној Менхетн удаљености од циљне позиције — што је неко стање даље од циља, то је мање пожељно. Конкретно, удаљеност се дијели са 3 да би се ублажила разлика између стања. Провалије имају изузетно ниске вриједности, чиме се избјегавају као непожељне. Након иницијализације,  $C$  се нормализује коришћењем softmax функције, добијајући вјероватосну дистрибуцију коју агент користи приликом планирања.

```

1 C = np.full(n_states, -1.0)
2
3 for idx in range(n_states):
4     i, j = grid_locations[idx]
5     distance = abs(3 - i) + abs(11 - j)
6     c_value = -distance / 3
7
8     if i == 3 and 1 <= j <= 10:
9         c_value = -100
10    C[idx] = c_value
11
12 C = softmax(C)
13
14
15 plot_beliefs(C)

```

Код 2.1.7: Дефинисање матрице  $C$  и њен графички приказ

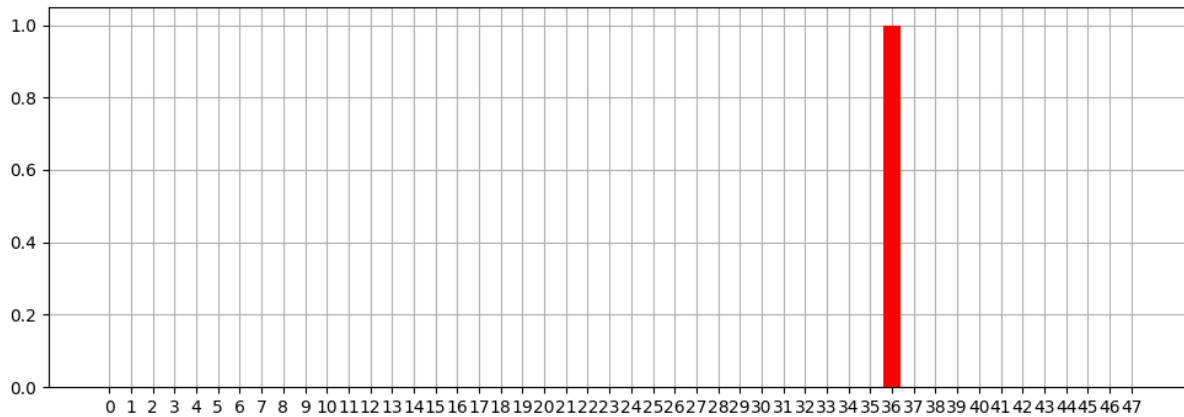


Слика 2.2: График пожељних стања

Почетна уверења  $D$  представљају почетно стање у којем агент верује да се налази, што је у овом случају локација  $[3, 0]$ . То је моделовано као јединични вектор (one-hot) који има вриједност 1 на почетном стању и 0 на свим осталим.

```
1 D = utils.onehot(grid_locations.index((3,0)), n_states)
2
3 plot_beliefs(D)
```

Код 2.1.8: Дефинисање матрице  $D$  и њен графички приказ



Слика 2.3: График вјероватноћа почетних стања

#### 2.1.4 Функције за инференцију и рачунање очекиване слободне енергије

Након дефинисања генеративног модела, следећи корак јесте доношење закључака о скривеним стањима на основу опсервација и претходних увјерења, као и предвиђање последица различитих потенцијалних акција. На овом мјесту представљене су двије суштинске компоненте Active Inference-а: ажурирање увјерења, којом агент изводи претпоставке о свом тренутном стању, и планирање путем очекиване слободне енергије, која омогућава агенту да процијени и изабере најпожељније акције на основу предвиђених исхода. Ове компоненте заједно чине основу интелигентног понашања заснованог на принципу минимизације слободне енергије.

Ове функције су преузете из Active Inference from Scratch и прилагођене нашем задатку. Тамо се може прочитати нешто више о њима.

Функција `infer_states` ажурира увјерење агента о скривеним стањима на основу нове опсервације и претходног увјерења, користећи Бајесово правило у логаритамском простору. Конкретно, логаритам условне вероватноће запажања у зависности од скривеног стања и логаритам дистрибуције претходног увјерења сабирају се елемент по елемент. Тако добијени збир се затим нормализује примјеном `softmax` функције.

```
1 def infer_states(observation_index, A, prior):
2     log_likelihood = log_stable(A[observation_index,:])
3
4     log_prior = log_stable(prior)
5
6     qs = softmax(log_likelihood + log_prior)
7
8     return qs
```

Код 2.1.9: Функција `infer_states`

Сада ћемо имплементирати функцију која предвиђа која ће скривена стања вероватно настати један корак у будућност, ако се изврши одређена акција.

```
1 def get_expected_states(B, qs_current, action):
2     qs_u = B[:, :, action].dot(qs_current)
3
4     return qs_u
```

Код 2.1.10: Функција `get_expected_states`

Наредна функција предвиђа која ће опажања највероватније бити добијена након што агент уђе у одређена стања.

```
1 def get_expected_observations(A, qs_u):
2     qo_u = A.dot(qs_u)
3
4     return qo_u
```

Код 2.1.11: Функција `get_expected_observations`

Даћемо сада имплементацију функција `entropy` и `kl_divergence`. Прва израчунава ентропију (мјеру неизвјесности) за свако скривено стање, на основу вјероватноћа различитих опажања у том стању. Друга израчунава Кулбак-Лајблерову дивергенцију између двије вјероватносне дистрибуције: `qo_u`, очекивана опажања (шта агент

очекује да види) и C, агентове претходне преференције (шта жели да види). Она мјери колико се очекивана опажања разликују од онога што агент жели.

```
1 def entropy(A):
2     H_A = - (A * log_stable(A)).sum(axis=0)
3
4     return H_A
5
6
7 def kl_divergence(qo_u, C):
8     return (log_stable(qo_u) - log_stable(C)).dot(qo_u)
```

Код 2.1.12: Функције entropy и kl\_divergence

Следећа функција представља централни механизам планирања у оквиру Active Inference-a. Она израчунава очекивану слободну енергију (G) за сваки могући низ акција, односно политику. За сваку политику симулира како би се увјерење о скривеним стањима и очекиваним опажањима развијало током предвиђеног временског хоризонта.

Очекивана слободна енергија се састоји од два главна члана:

1. Очекивана ентропија опажања — мјера неизвјесности коју агент предвиђа да ће искусити.
2. Кулбак-Лајблерова дивергенција између предвиђених и пожељних опажања — мјера тога колико политика доводи до опажања која су у складу са преференцијама агента.

Оба ова члана се сабирају за сваки корак унутар политике, а њихов збир представља укупну очекивану слободну енергију политике. Нижа вриједност G означава пожељнију политику — ону која води ка предвидљивијим (мање неодређеним) и пожељнијим опажањима,

```
1 def calculate_G_policies(A, B, C, qs_current, policies):
2
3     G = np.zeros(len(policies)) # initialize the vector of expected free
4     # energies, one per policy
5     H_A = entropy(A) # can calculate the entropy of the A
6     # matrix beforehand, since it'll be the same for all policies
7
8     for policy_id, policy in enumerate(policies): # loop over policies -
9     # policy_id will be the linear index of the policy (0, 1, 2, ...) and
10    # 'policy' will be a column vector where 'policy[t,0]' indexes the
```

```

7      action entailed by that policy at time 't'
8
9      t_horizon = policy.shape[0] # temporal depth of the policy
10
11      G_pi = 0.0 # initialize expected free energy for this policy
12
13      for t in range(t_horizon): # loop over temporal depth of the policy
14
15          action = policy[t,0] # action entailed by this particular policy,
16          at time 't'
17
18          # get the past predictive posterior - which is either your
19          current posterior at the current time (not the policy time) or the
20          predictive posterior entailed by this policy, one timestep ago (in
21          policy time)
22          if t == 0:
23              qs_prev = qs_current
24          else:
25              qs_prev = qs_pi_t
26
27          qs_pi_t = get_expected_states(B, qs_prev, action) # expected
28          states, under the action entailed by the policy at this particular
29          time
30          qo_pi_t = get_expected_observations(A, qs_pi_t) # expected
31          observations, under the action entailed by the policy at this
32          particular time
33
34          kld = kl_divergence(qo_pi_t, C) # Kullback-Leibler divergence
35          between expected observations and the prior preferences C
36
37          G_pi_t = H_A.dot(qs_pi_t) + kld # predicted uncertainty +
38          predicted divergence, for this policy & timepoint
39
40          G_pi += G_pi_t # accumulate the expected free energy for each
41          timepoint into the overall EFE for the policy
42
43      G[policy_id] += G_pi
44
45  return G

```

Код 2.1.13: Функција calculate\_G\_policies

У наставку је дата функција која израчунава вјероватноће избора сваке појединачне акције на основу вјероватноћа додељених свакој политици ( $Q_\pi$ ). Пошто се одлука о акцији доноси у садашњем тренутку, узима се у обзир само прва акција

у свакој политици. Умјесто сабирања, користи се максимална вриједност  $Q_\pi$  за сваку акцију (односно највећа вјероватноћа неке политике која ту акцију предвиђа као прву), што значи да се као изабрана узима акција са почетка политике која има највећу вјероватноћу. На крају, резултујући вектор се нормализује како би представљао валидну дистрибуцију вероватноћа над акцијама.

```

1 def compute_prob_actions(actions, policies, Q_pi):
2     P_u = np.zeros(len(actions)) # initialize the vector of probabilities
   of each action
3
4     for policy_id, policy in enumerate(policies):
5         P_u[int(policy[0,0])] = max(P_u[int(policy[0,0])], Q_pi[policy_id])
6
7     P_u = utils.norm_dist(P_u) # normalize the action probabilities
8
9     return P_u

```

Код 2.1.14: Функција compute\_prob\_actions

Функција active\_inference\_with\_planning представља главну петљу агента који користи приступ Active Inference са планирањем да би доносио одлуке и интераговао са окружењем током задатог броја временских корака. На почетку, агент поставља своје почетно увјерење о скривеном стању користећи матрицу  $D$ , ресетује окружење и добија прво опажање. Такође конструише све могуће политике – низове акција одређене дужине – које ће користити за планирање.

У сваком кораку, агент најпре опсервира своје тренутно стање у окружењу. Ако је достигао циљну локацију (тј. координату 47), прекида епизоду и враћа коначно увјерење и акумулирану награду. У супротном, врши инференцију над скривеним стањем на основу запажања и матрице  $A$ , користећи Бајесово правило. Ово увјерење представља његово тренутно разумијевање гдје се налази у простору скривених стања.

На основу овог увјерења, агент симулира сваки могући низ акција (свака политика) и израчунава њихову очекивану слободну енергију, користећи функцију calculate\_G\_policies. Ова вриједност укључује предвиђену неизвјесност у опажањима као и разлику између предвиђених и пожељних опажања. Затим, примјењује softmax функцију над негативним вредностима слободне енергије да би добио вероватноће избора сваке политике ( $Q_\pi$ ), и на основу тога израчунава маргиналну вероватноћу избора сваке акције ( $P_u$ ). Из  $P_u$  се затим бира акција са највећом вероватноћом.



Након избора акције, агент предвиђа расподелу скривених стања која би могла настати као последица те акције, користећи транзициону матрицу  $B$  и своје текуће увјерење. Креће се у окружењу, добија ново опажање и награду, и наставља циклус све док не истекне вријеме, или се епизода не заврши услед постизања циља. На крају, функција враћа коначно увјерење о скривеном стању и укупну награду прикупљену током епизоде.

```

1 def active_inference_with_planning(A, B, C, D, n_actions, env,
  policy_len=2, T=5):
2
3     """ Initialize prior, first observation, and policies """
4     prior = D
5     obs, _ = env.reset()
6     cumulative_reward = 0
7
8     policies = construct_policies([n_states], [n_actions], policy_len=
  policy_len)
9
10
11     for t in range(T):
12         print(f'\nTime {t}: Agent observes itself in location: {obs} (
  Grid: {grid_locations[obs]})')
13
14         if obs == 47:
15             print(f"Agent reached the target at {grid_locations[obs]}!
  Stopping.")
16             return qs_current, cumulative_reward # Stop and return
17
18         # Inference (use the integer observation directly)
19         qs_current = infer_states(obs, A, prior)
20
21         # Calculate expected free energy
22         G = calculate_G_policies(A, B, C, qs_current, policies)
23         Q_pi = softmax(-G)
24         P_u = compute_prob_actions(actions, policies, Q_pi)
25         chosen_action = np.argmax(P_u)
26
27         # Update prior
28         prior = B[:, :, chosen_action].dot(qs_current)
29
30         next_obs, reward, terminated, truncated, info = env.step(
  chosen_action)
31
32         cumulative_reward += reward

```

```

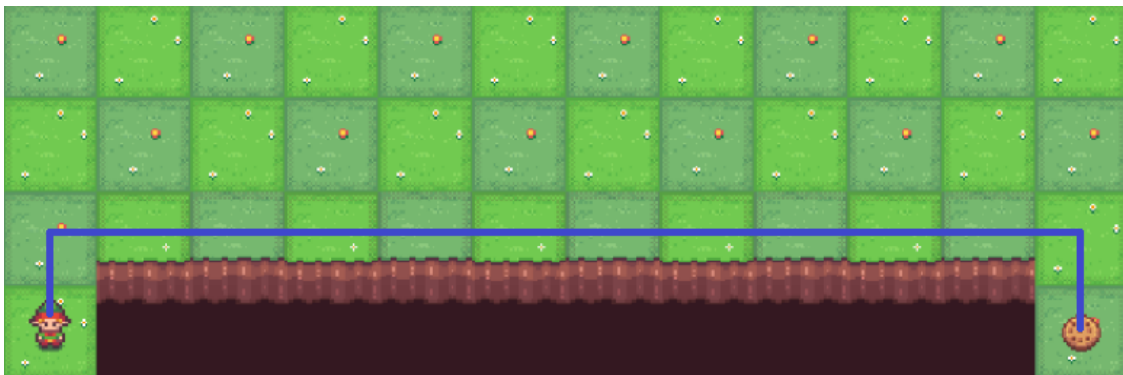
33
34
35     print(f"Chosen action: {actions[chosen_action]} | Reward: {
reward}")
36     frame = env.render() # returns an RGB array since env was
created with render_mode="rgb_array"
37     display_frame(frame, title_str=f"Time {t + 1}: Action {actions[
chosen_action]}")
38
39     if terminated or truncated:
40         print(f"Episode terminated! Total reward: {
cumulative_reward}")
41         return qs_current, cumulative_reward
42     obs = next_obs
43
44     print(f"\nFinal cumulative reward after {T} timesteps: {
cumulative_reward}")
45
46     return qs_current, cumulative_reward

```

Код 2.1.15: Функција active\_inference\_with\_planning

### 2.1.5 Понашање агента

Након извршене имплементације, посматрајмо понашање агента у симулираном окружењу. Циљ агента је да стигне до задате локације, при чему се очекује да прелази што краћи пут и избегава провалије. На доњој слици приказана је оптимална путања од почетне позиције до циљне локације. За тај пут му је потребно да направи 13 корака.



Слика 2.4: Оптимална путања агента

У оквиру тестирања, агент разматра политике дужине 7 корака унапријед (`policy_len = 7`), док је максимално трајање епизоде ограничено на 20 временских корака ( $T = 20$ ).

Резултати симулације показују да агент, вођен принципом Active Inference, доследно прати ову руту. Ово указује да агент успјешно планира низ акција које минимизују очекивану слободну енергију, чиме постиже и функционално и рационално понашање у датом окружењу.

## 2.2 Модел са шумом у матрицама A и D

У овој варијанти модела, уводи се додатна неодређеност у опажањима и почетно увјерење агента, како би се симулирали услови ограничене или зашумљене перцепције. У циљу моделирања ситуације у којој агент нема потпуно поуздана сензорна опажања и у којој његова почетна увјерења о локацији нису прецизна, модификован је основни модел у неколико кључних аспеката.

Најприје, матрица опажања A проширена је тако да не одражава више детерминистичку везу између стања и опажања. Конкретно, ако опсервира да је у нултом реду, постоји једнака вјероватноћа да је у опаженом или у стању испод. Слично томе, када агент опази да се налази у стању у првом или другом реду мреже, модел му даје једнаку вјероватноћу да се заправо налази у опаженом стању или у оном изнад. Посебан случај је стање 36 (почетно стање), гдје се такође укључује неизвјесност према горњем стању. Ове вјероватноће симулиране су тако што су одговарајуће вриједности у матрици A постављене на једнаке бројеве, а затим је матрица нормализована по колонама. На тај начин, опажања агента више нису поуздан одраз стварног стања, већ садрже шум.

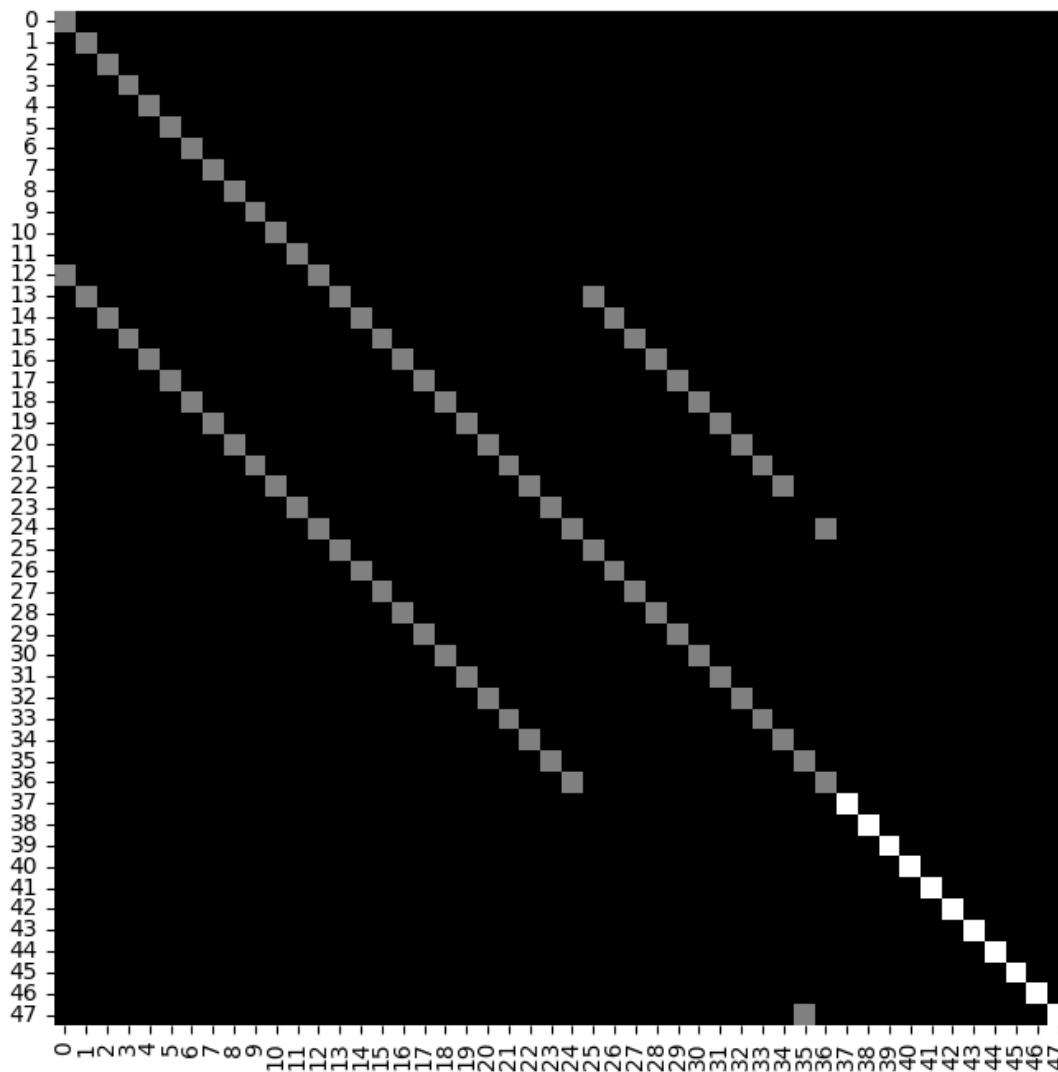
```
1 A = np.eye(n_observations, n_states)
2
3 for i in range(n_observations):
4     if 0 <= i <= 11:
5         A[i+12, i] = 1
6     elif 12 <= i <= 35:
7         if 24 < i < 35:
8             A[i-12, i] = 1
9         else:
10            A[i+12, i] = 1
11     elif i == 36:
12         A[i-12, i] = 1
13
```

```

14 A = A / A.sum(axis=0, keepdims=True)
15
16
17 plot_likelihood(A)

```

Код 2.2.1: Модификована матрица A



Слика 2.5: График вјероватноће запажања у зависности од скривених стања са шумом

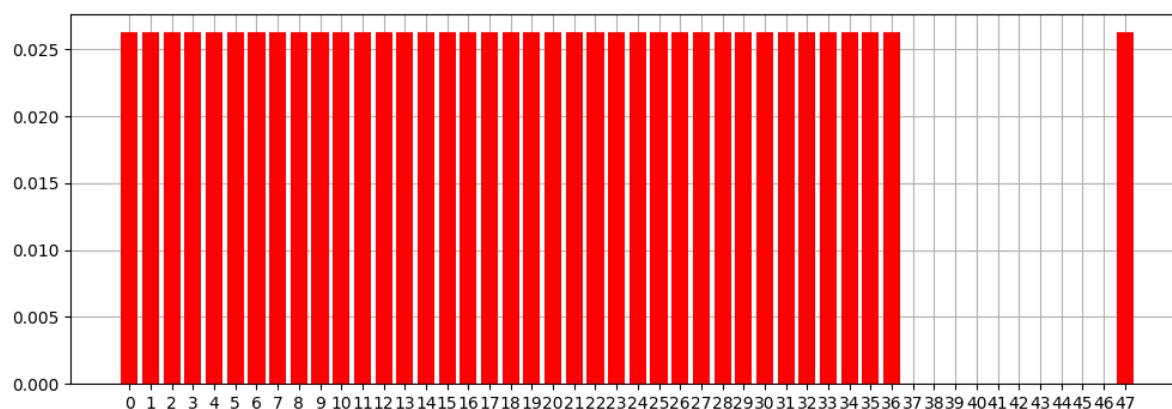
Поред тога, измјењена је и почетна распоdjела увjерења, односно матрица D. Вриједности у овој дистрибуцији су подешене тако да стања која представљају литица (индекси 37–46) имају вjероватноћу једнаку нули, након чега је цијела распоdjела нормализована. На овај начин агент на самом почетку нема никакву претпоставку о томе у ком се конкретном стању налази, већ креће са равномjерном неизвјесношћу над преосталим стањима. Оваква иницијализација подстиче агента да кроз кретање и опажања постепено изгради слику о својој позицији у окружењу.

```

1 D = np.ones(n_states)
2
3 D[37:47] = 0
4
5 D = D / D.sum()
6
7 plot_beliefs(D)

```

Код 2.2.2: Модификована матрица D



Слика 2.6: График вjероватноће почетних стања са равномjерном распоdjелом над дозвољеним позицијама

Да би се агент подстакао на истраживање окружења у условима повећане неизвјесности, извршене су двије кључне измјене у моделу.

У функцији `compute_prob_actions`, начин обраде вjероватноћа појединачних акција промијењен је са избора политике са највећом вjероватноћом (`max`) на сабирање доприноса свих политика (`+`). Ова промјена значи да ће се вjероватноћа одређене

акције повећавати пропорционално томе колико се често та акција појављује као прва у високовјероватним политикама.

```
1 def compute_prob_actions(actions, policies, Q_pi):
2     P_u = np.zeros(len(actions)) # initialize the vector of probabilities
      of each action
3
4     for policy_id, policy in enumerate(policies):
5         P_u[int(policy[0,0])] += Q_pi[policy_id] # get the marginal
      probability for the given action, entailed by this policy at the
      first timestep
6
7     P_u = utils.norm_dist(P_u) # normalize the action probabilities
8
9     return P_u
```

Код 2.2.3: Измијењена compute\_prob\_actions функција

На крају, избор акције више није базиран на детерминистичком критеријуму (argmax), већ се користи стохастички избор преко utils\_sample. Овај приступ омогућава да се и акције са нешто мањом вјероватноћом повремено бирају, што директно подржава истраживање. Тако се боље истиче значај направљених промјена у условима неизвјесности агента о његовом тренутном положају и окружењу.

```
1 def active_inference_with_planning(...):
2     # dio funkcije prije izbora akcije...
3
4     # Izbor akcije: ranije je birana akcija sa najvecom vjerovatnocom
5
6     # chosen_action = np.argmax(P_u)
7
8     # Sada se koristi uzorkovanje iz raspodjele, da bi se podstaklo
      istrazivanje:
9     chosen_action = utils.sample(P_u)
10
11     # ostatak funkcije...
```

Код 2.2.4: active\_inference\_with\_planning

Током тестирања коришћена је дужина политике (policy\_len) једнака 7 и дужина епизоде (T) је постављена на 30 корака.

Агент се не понаша увијек најоптималније, јер због имплементираног истраживања понекад бира акције које нису најкраћи пут до циља. Међутим, и поред тога, он

увијек успијева да стигне до циља у релативно малом броју корака. Такође, агент вјешто избјегава литице и никада не упада у њих. Иако често није сигуран да ли се налази у опсервираном пољу или у пољу изнад њега, оптималне путање у оба случаја се у великој мјери поклапају, па ова неизвјесност не доводи до озбиљних проблема у кретању. У неким другим сценаријима, са другачије дефинисаним матрицама А и D, ова врста перцептивне неодређености би могла више утицати на понашање агента.

## 2.3 Модел са клизавом површином (`is_slippery=True`)

У овој варијанти модела уводи се додатна неодређеност у кретању агента кроз окружење, с циљем симулирања услова кретања по клизавој подлози. Агент се са вјероватноћом од  $1/3$  креће у намјераваном правцу, док постоји једнака вјероватноћа од по  $1/3$  да ће завршити у једном од два правца која су нормална у односу на онај који је изабрао. Ова промјена уводи неизвјесност у саму динамику окружења, захтијевајући од агента да прилагоди процес планирања и доношења одлука у условима у којима више нема потпуну контролу над сопственим кретањем.

У односу на стандардни модел, једина промјена у овој варијанти односи се на дефинисање транзиционе матрице B, која описује вјероватноће преласка из једног скривеног стања у друго као резултат одређене акције. Док је у основном моделу прелаз био детерминистички, овдје је уведена стохастичка компонента која симулира клизаву подлогу: агент се са вјероватноћом од  $1/3$  креће у намјераваном правцу, а са по  $1/3$  у један од два правца која су нормална на жељени.

На тај начин, матрица B је прилагођена тако да, у зависности од тренутне позиције и намјераване акције, расподјели вјероватноће на три могуће наредне позиције. Испод је приказан код којим је дефинисана нова транзициона матрица:

```
1 def create_B_matrix():
2     B = np.zeros( (len(grid_locations), len(grid_locations), len(actions))
3                 )
4     for action_id, action_label in enumerate(actions):
5
6         for curr_state, grid_location in enumerate(grid_locations):
7
8             y, x = grid_location
9
10            if (y == 3 and 1 <= x <= 10):
11                x = 0
12                y = 3
```

```

13
14     if action_label == "UP":
15         next_y = y - 1 if y > 0 else y
16         next_x = x
17         next_x1 = max(x - 1, 0)
18         next_y1 = y
19         next_x2 = min(x + 1, 11)
20         next_y2 = y
21     elif action_label == "DOWN":
22         next_y = y + 1 if y < 3 else y
23         next_x = x
24         next_x1 = max(x - 1, 0)
25         next_y1 = y
26         next_x2 = min(x + 1, 11)
27         next_y2 = y
28     elif action_label == "LEFT":
29         next_x = x - 1 if x > 0 else x
30         next_y = y
31         next_x1 = x
32         next_y1 = max(y - 1, 0)
33         next_x2 = x
34         next_y2 = min(y + 1, 3)
35     elif action_label == "RIGHT":
36         next_x = x + 1 if x < 11 else x
37         next_y = y
38         next_x1 = x
39         next_y1 = max(y - 1, 0)
40         next_x2 = x
41         next_y2 = min(y + 1, 3)
42
43
44     new_location = (next_y, next_x)
45     new_location1 = (next_y1, next_x1)
46     new_location2 = (next_y2, next_x2)
47     next_state = grid_locations.index(new_location)
48     next_state1 = grid_locations.index(new_location1)
49     next_state2 = grid_locations.index(new_location2)
50
51     if next_state != next_state1 and next_state != next_state2:
52         B[next_state, curr_state, action_id] = 1/3
53         B[next_state1, curr_state, action_id] = 1/3
54         B[next_state2, curr_state, action_id] = 1/3
55     elif next_state == next_state1 and next_state != next_state2:
56         B[next_state, curr_state, action_id] = 2/3
57         B[next_state2, curr_state, action_id] = 1/3

```



```

58     elif next_state != next_state1 and next_state == next_state2:
59         B[next_state, curr_state, action_id] = 2/3
60         B[next_state1, curr_state, action_id] = 1/3
61     else:
62         B[next_state, curr_state, action_id] = 1.0
63
64
65     return B
66
67 B = create_B_matrix()

```

Код 2.3.1: Дефинисање матрице B у моделу са клизавом подлогом

Током тестирања коришћена је дужина политике `policy_len = 2`, што се показало као погодан избор, будући да агент не може прецизно да предвиди позицију у даљој будућности, па краће политике боље одражавају значај актуелних одлука. Максимална дужина епизоде постављена је на  $T = 60$ .

Подсјетимо се [овдје](#) изгледа окружења. Када се агент налази у почетном стању (стање 36), он бира акцију LEFT. Оваква одлука минимизира ризик од пада у литицу (која се налази десно), а омогућава могућност да се агент омакне у горњи ред, док у најгорем случају остаје у истом положају. Након што се позиционира у реду 2, бира акцију UP, чиме избјегава опасну зону испод, а има могућност да се помјери и бочно. Када дође у ред 1, почиње да се креће RIGHT кроз мрежу, настављајући тако и у случају да се омакне у ред 0. Овај образац кретања – вертикално уздизање ка безбједнијој зони, а затим хоризонтално кретање ка циљу – омогућава агенту да избјегне литице и да постепено напредује.

По доласку у последњу колону, агент испитује своју позицију. Ако се налази у стању 11 или 23 (нулти и први ред), користи акцију DOWN да би дошао до стања 35. Уколико је већ у стању 35, бира акцију RIGHT, ослањајући се на могућност да се омакне у циљно стање 47 које се налази испод.

Због клизаве површине и сталне могућности непредвиђеног помјерања, агент веома тешко стиже до циља — али његово понашање и даље одражава оптималну адаптацију у датом окружењу.