

## Model Architecture

### 1. Environment

**State :** allele sequence; peptide sequence

**Action :**

1. stop or not: determine whether or not to stop the edit of sequence. (stop edits when stop == 1)
2. position: determine the position of edits (**replace one amino acid with another one**).
3. amino acid: determine the type of amino acid at the predicted position.

### 2. Policy Network $\pi(s)$

#### 2.1 Representation Learning

**Amino acid embedding:** Bi-LSTM

For each amino acid, we will represent it with the forward and backward embeddings.

$$h_i^* = \overrightarrow{f}(h_{i-1}^*, x_i); b_i = \overleftarrow{f}(b_{i+1}, x_i)$$

$$h_i = \text{Concat}(h_i^*, b_i)$$

For each allele sequence, we can select one of model architecture from CNN/RNN/MLP/Transformer to embed the allele sequence as  $a$ . (use MLP and CNN for parameter tuning.)

#### 2.2 Action Prediction

##### Stop prediction

Predict whether stopping the edit of sequence or not using the embeddings of peptide sequences and alleles

$$\text{Softmax}(m_1(h_n, a))$$

##### Amino acid position selection

Given a peptide sequence of length  $n$ , we can learn its embedding  $p_t = (h_0, \dots, h_{n-1}) \in R^{n \times k}$ .  
**(add** position encoding here)

$$f_i = \text{Softmax}(m_2(h_i, a, pos_i))$$

##### Amino acid type prediction

predict the type of added amino acid as below,

$$f_2 = \text{Softmax}(m_3(h_i, a))$$

$m_3$  will output an vector of length 20.

The probabilities of illegal actions should be masked.

### 3. Model training

Use PPO to optimize the model.

Sampling Method (used to sample the initial state):

All peptides are sampled from intermediate or negative binding peptides in IEDB dataset.

Some basic settings:

- feature type: blosum + onehot + deep
- latent dimension: (40, 20)
- value network and policy network: (80, 40)
- architecture for allele sequence: (CNN, FC)
- entropy coefficients: (0.0, 0.001, 0.01, 0.1)
- Hidden dimension for allele: (320)
- latent dimension for allele: (80)
- latent dimension for peptide (Bi-LSTM): (40)

#### Value Network $V(s)$

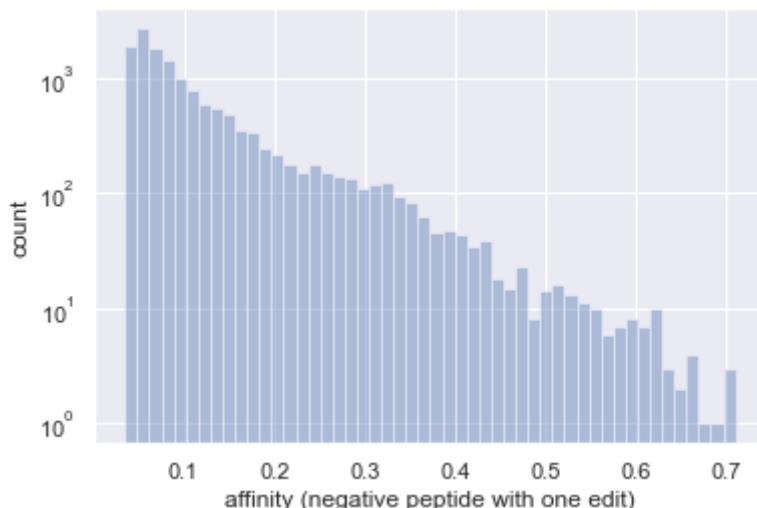
MLPs: Given a peptide sequence, predict the estimation of its future reward.

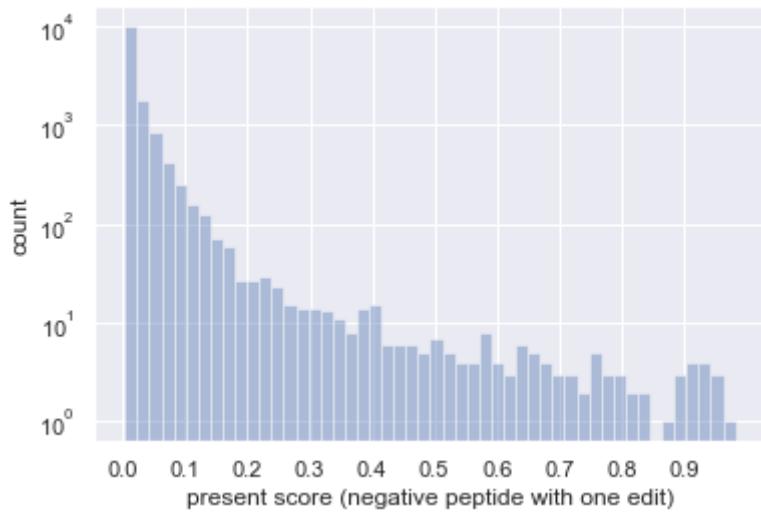
### 4. Reward Design

Analysis about the presentation scores predicted by MHCflurry:

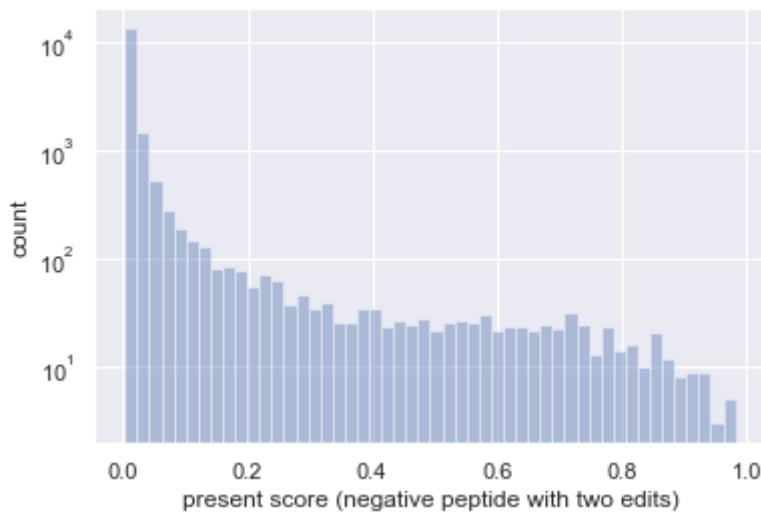
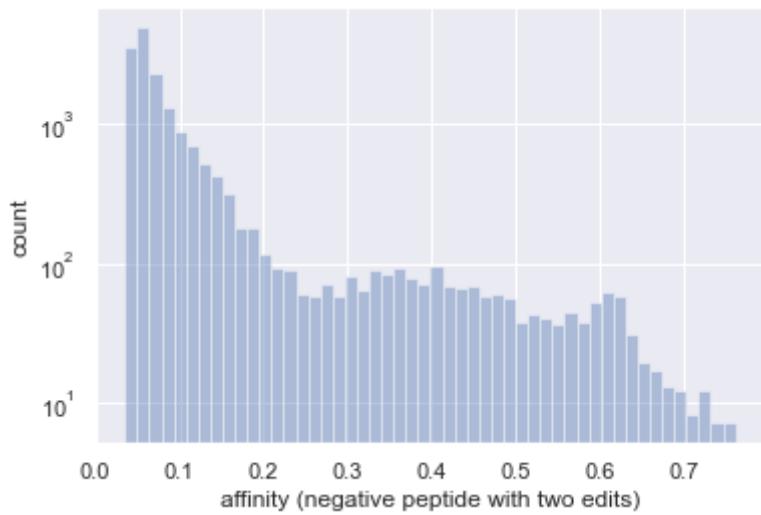
#### 4.1 Negative Sample

We sampled 100 peptides with binding affinity > 5000, and kept only samples with predicted presentation scores < 0.1. Then, we edited these peptides by replacing one amino acid and calculated the predicted affinities and presentation scores of new sequences.





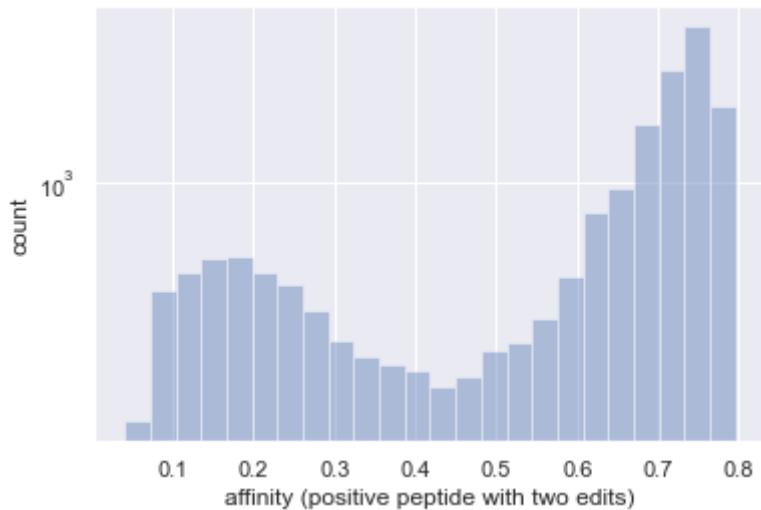
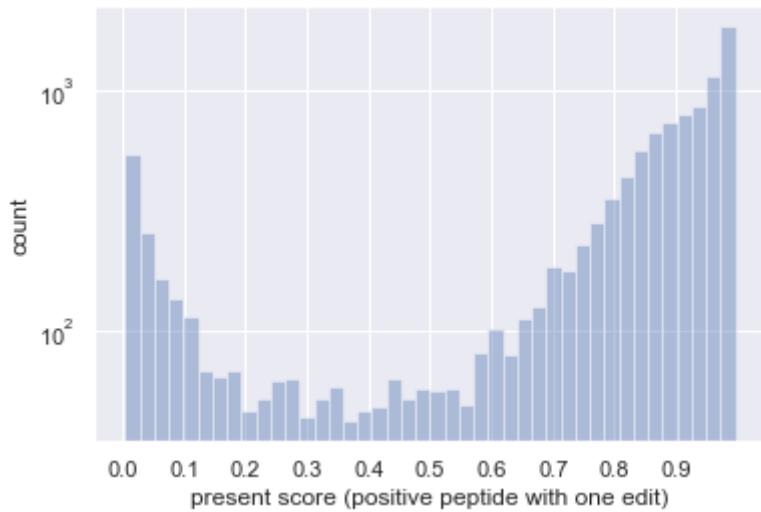
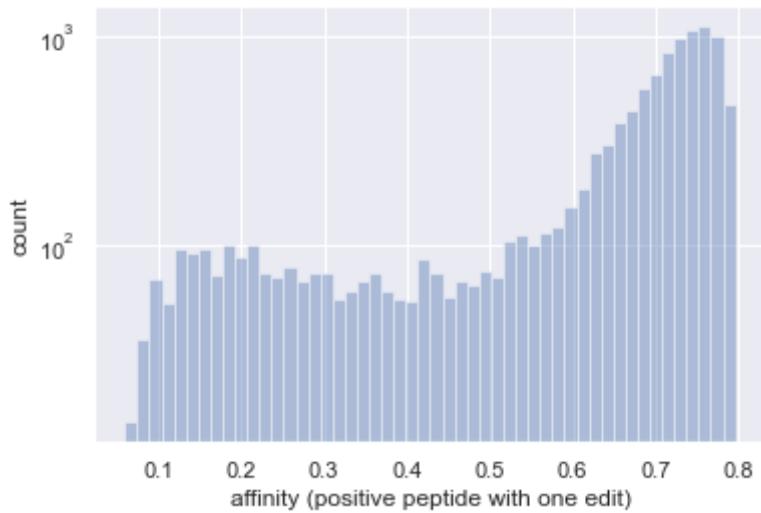
We sampled 100 peptides from the edited sequences, and edited these peptides again. Then, We calculated the predicted affinities and presentation scores of new sequences.

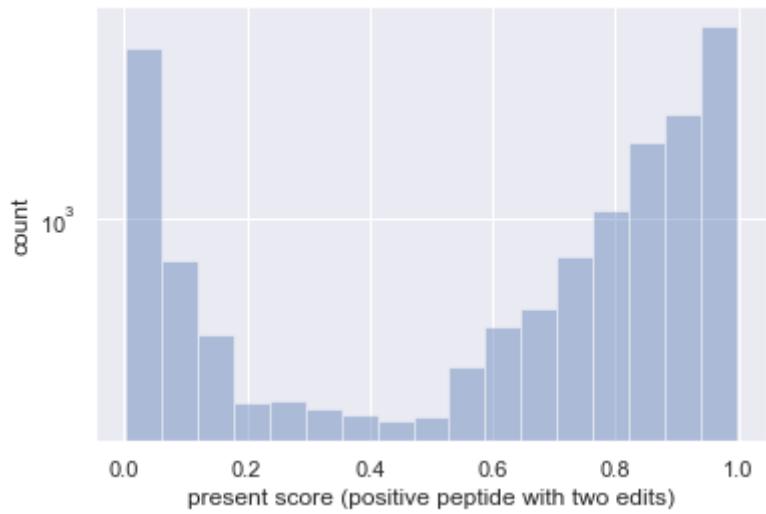


Conclusion: If we start from negative peptides, it is very likely that we can hardly receive any rewards (always 0 - 0.05).

#### 4.2 Positive Sample

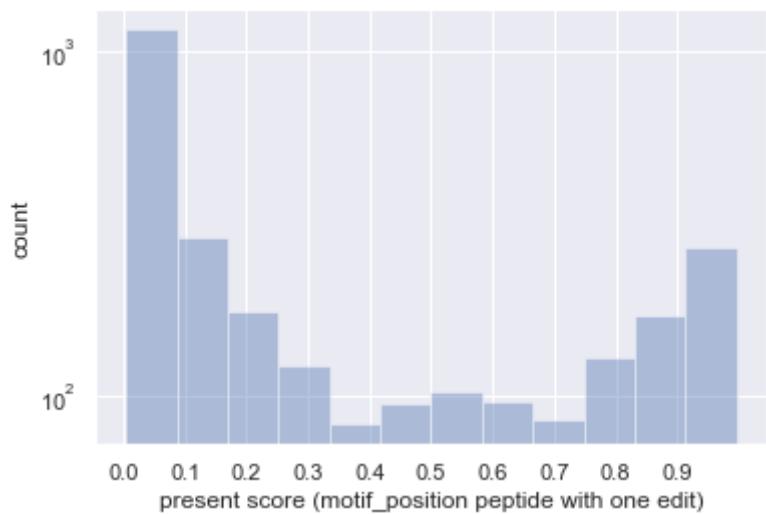
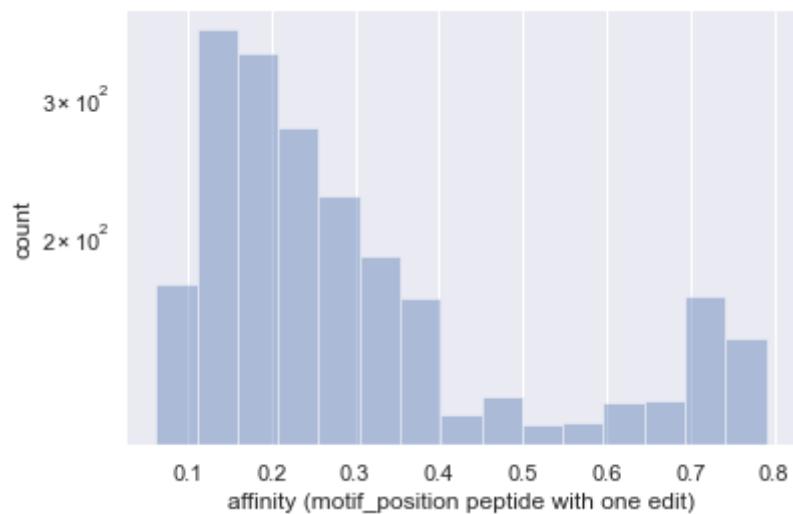
We sampled 100 peptides of length 9 binding to the allele HLA-A\*02:01 with binding affinity < 100, and kept only samples with predicted presentation scores > 0.8. Then, I edited these peptides by replacing one amino acid and calculated the predicted affinities and presentation scores of new sequences.

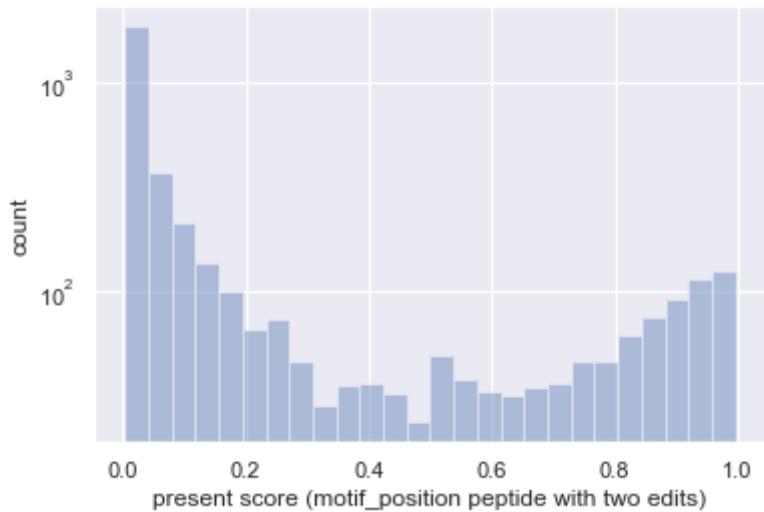
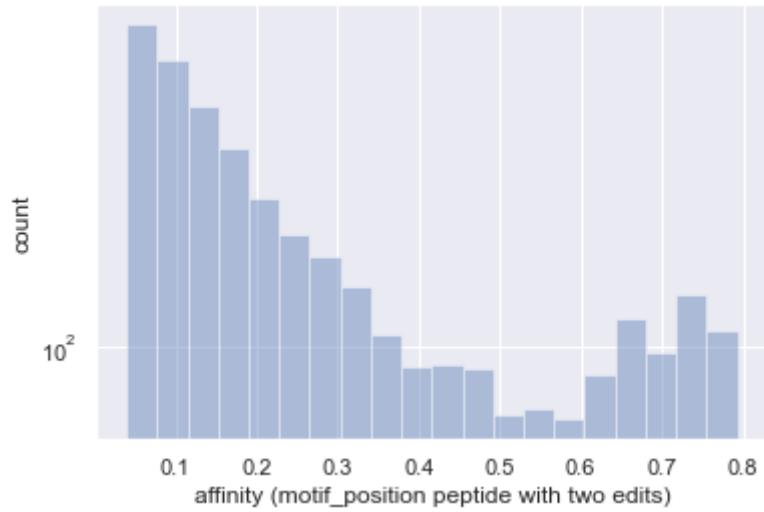




## Motif Position

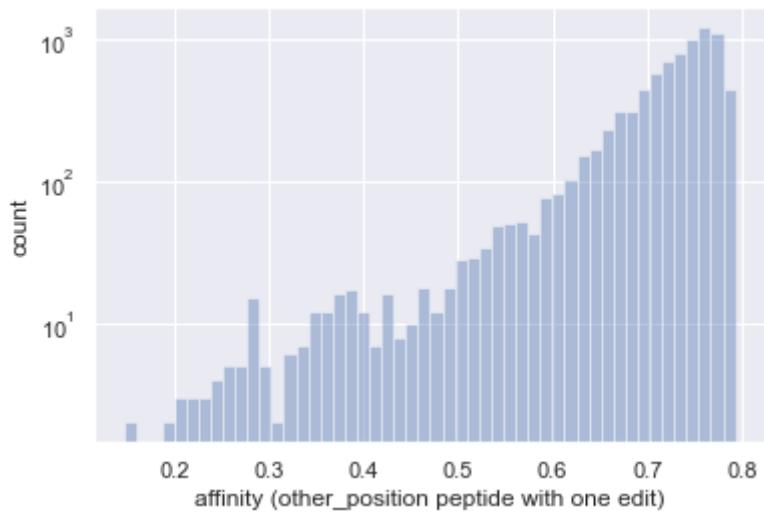
If we only edit the amino acids at the motif positions in peptides,

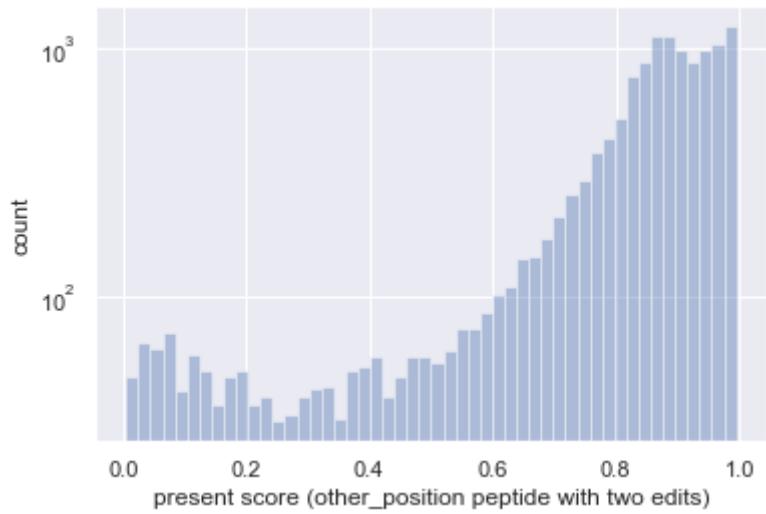
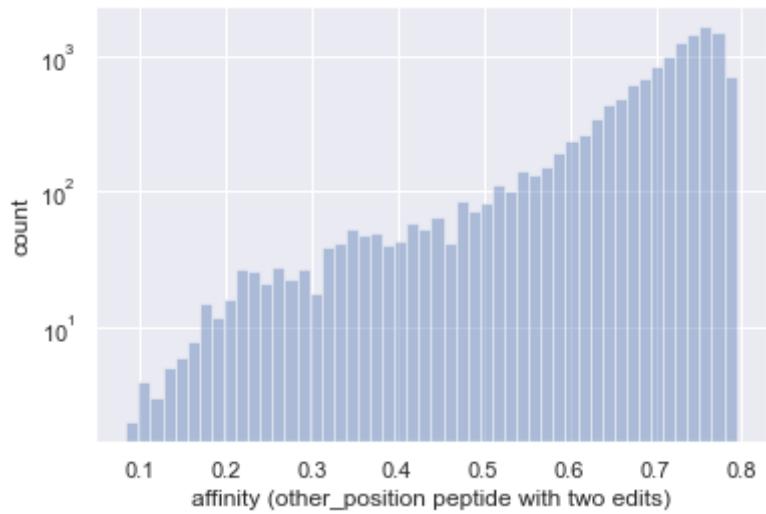
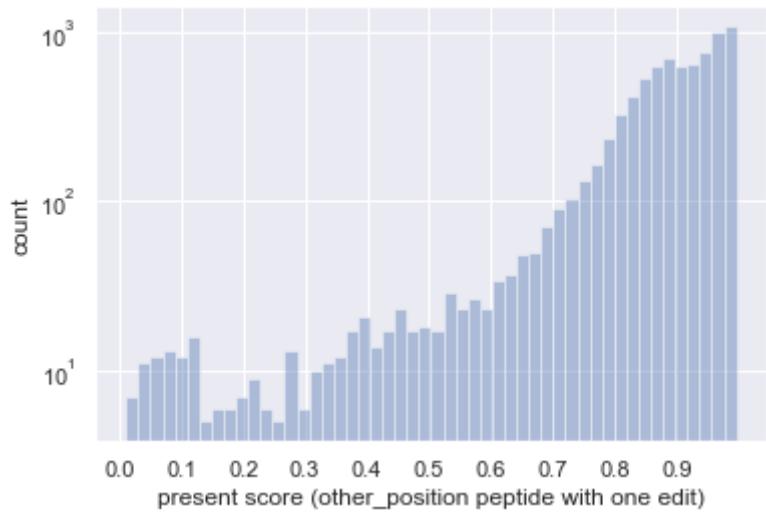




### Other position

If we edit the amino acids at the non-motif positions in peptides,





### 4.3 Reward Design

- Maze game and Mountain Car

For each step that the agent does not reach the goal, the environment returns a negative reward (penalty) to avoid long path.

**Note:** The end of edit process cannot be determined or stopped by policy network; otherwise, the policy network can always tend to stop the edit process as soon as possible. That means, we need to set a goal so that the environment can stop the edit process by either reaching the goal or reaching the maximum steps.

[Our reward design]:

Intermediate reward: -0.3

Final reward:  $10 * \text{presentation score}$

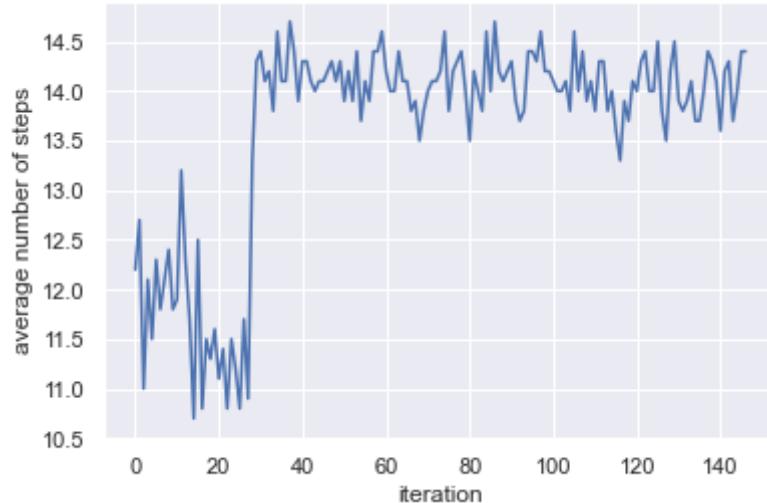
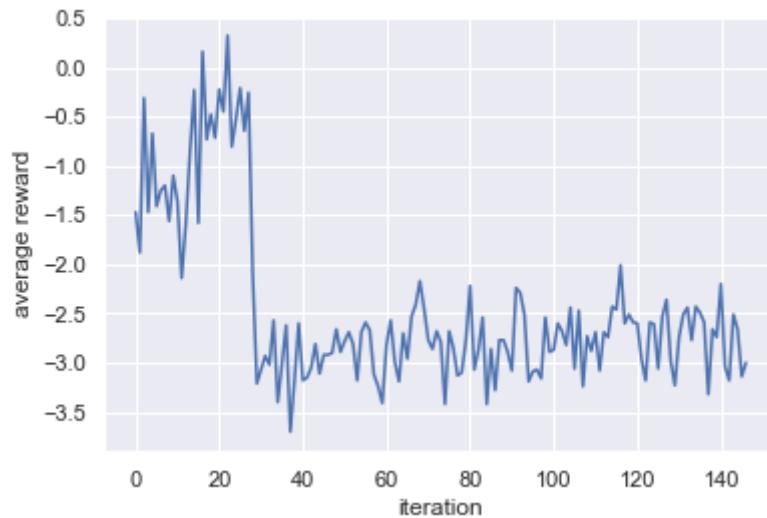
Maximum steps: 15

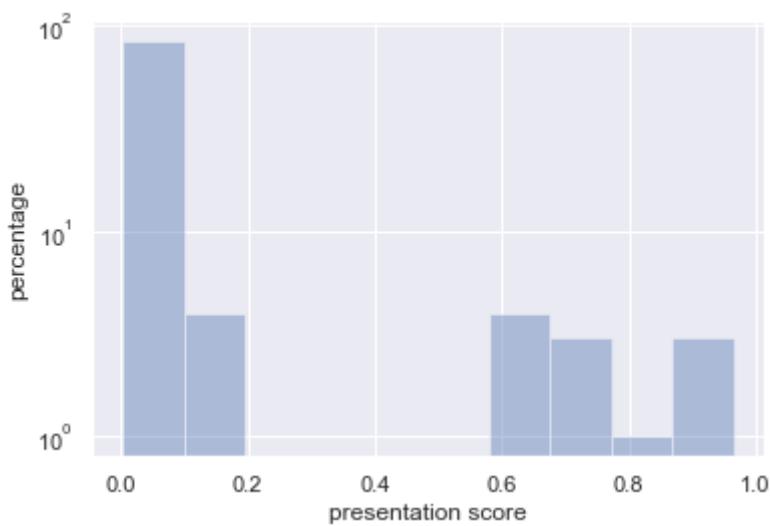
Goal: 0.6

Initial good samples: peptides with presentation score from 0.2 to 0.5, with the assumption that these peptides are easier to be optimized and thus provide more signals than random peptides.

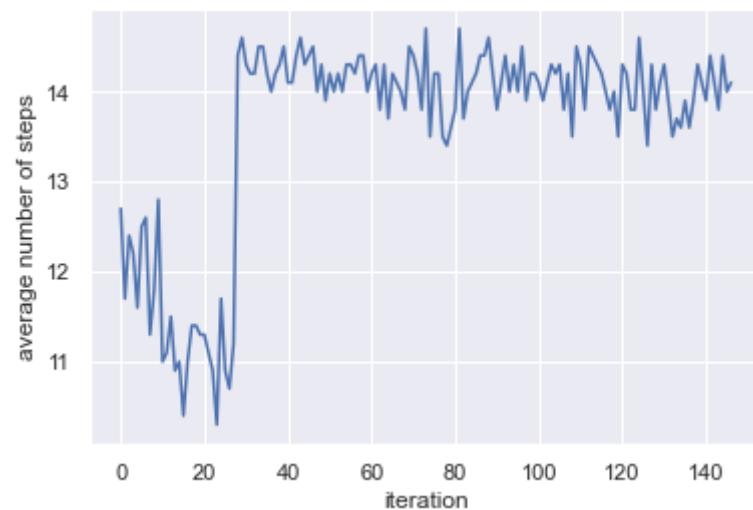
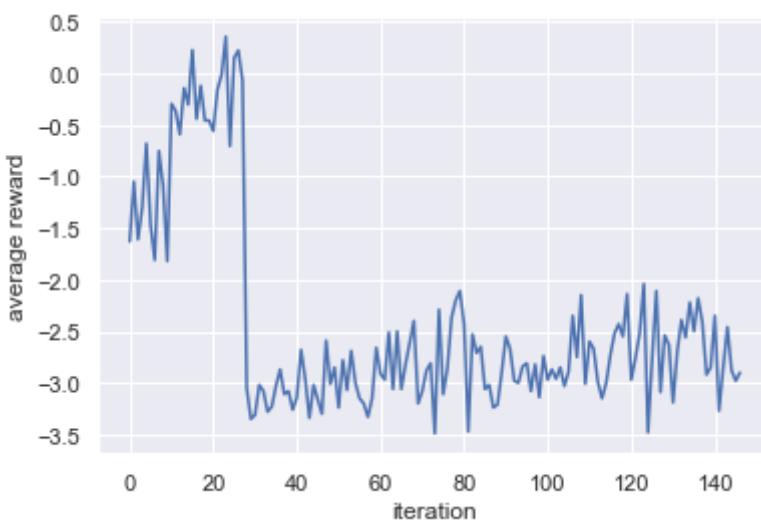
Results:

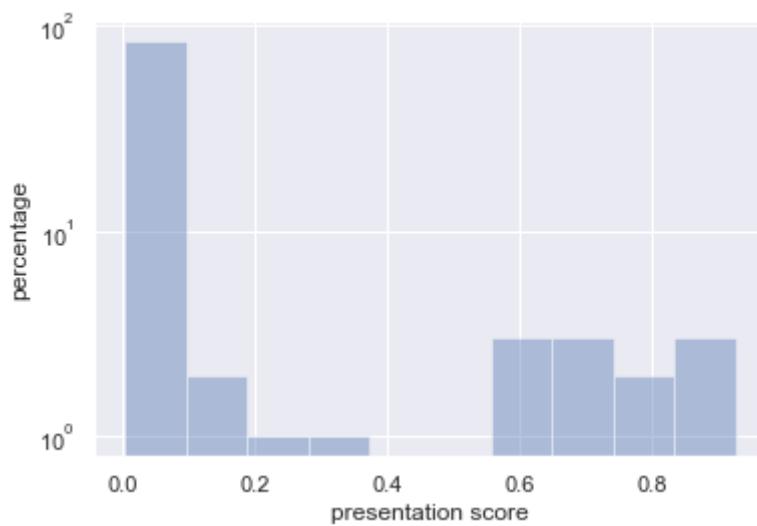
With CNN:





With FC:





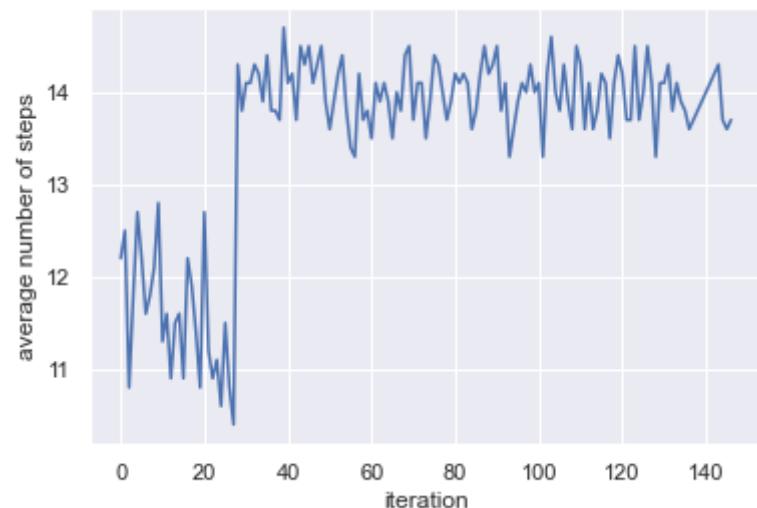
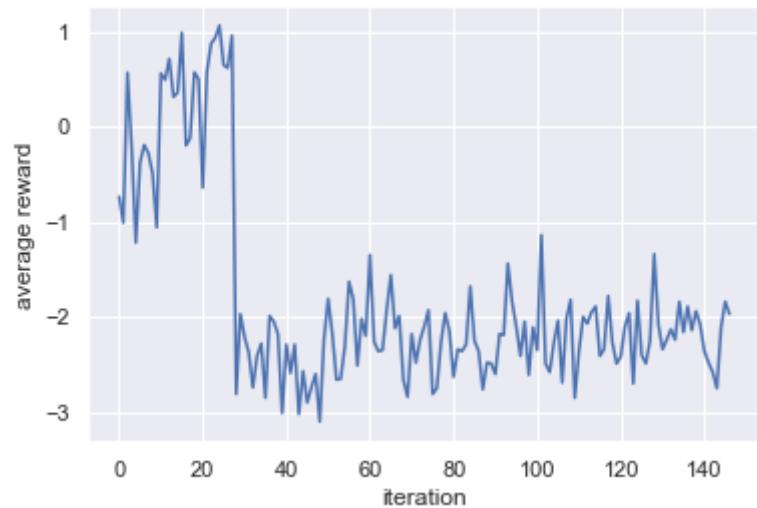
[Our reward design]:

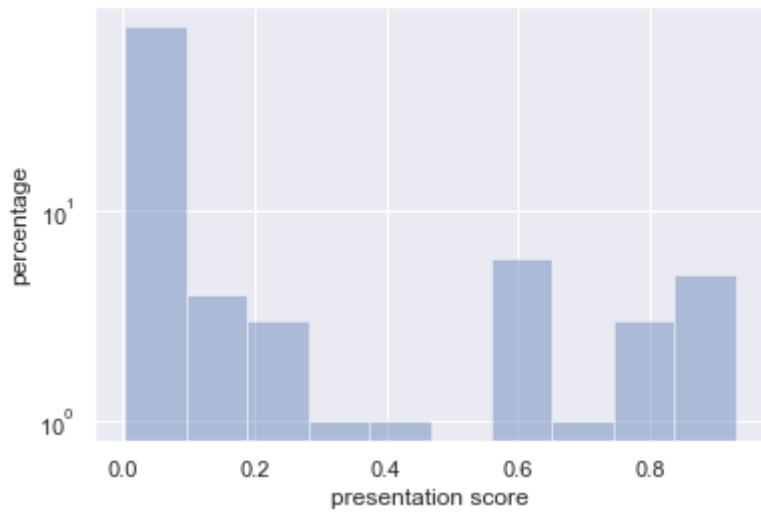
Intermediate reward:  $-1 + \text{presentation score}$

Final reward:  $10 * \text{presentation score}$

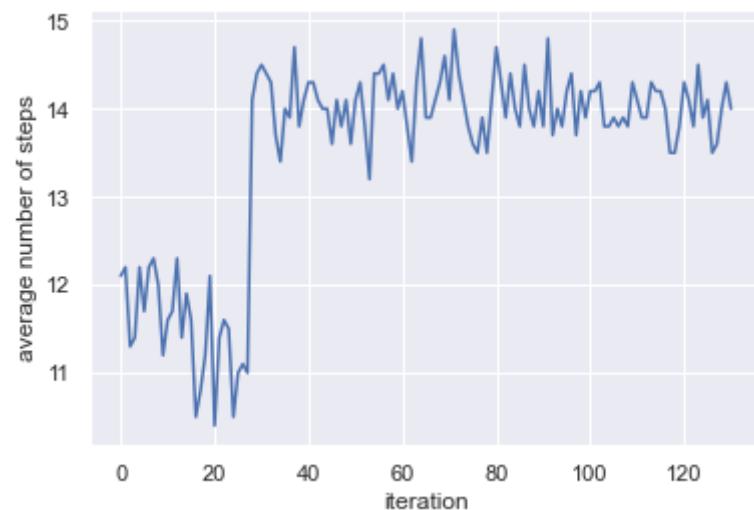
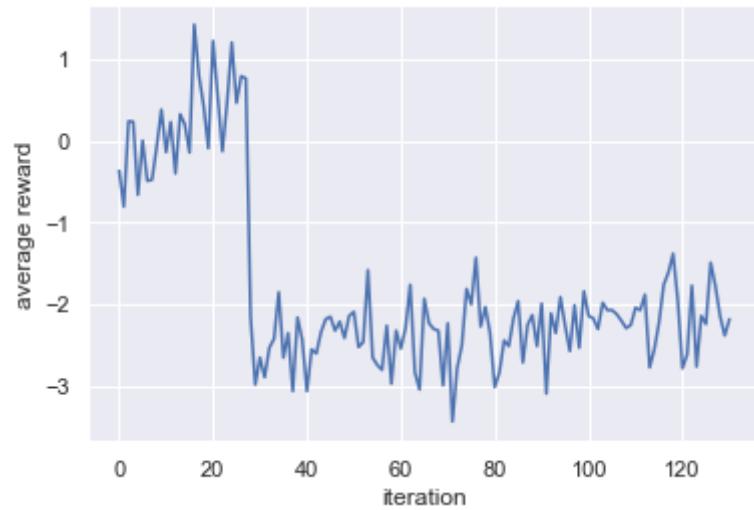
Maximum steps: 15

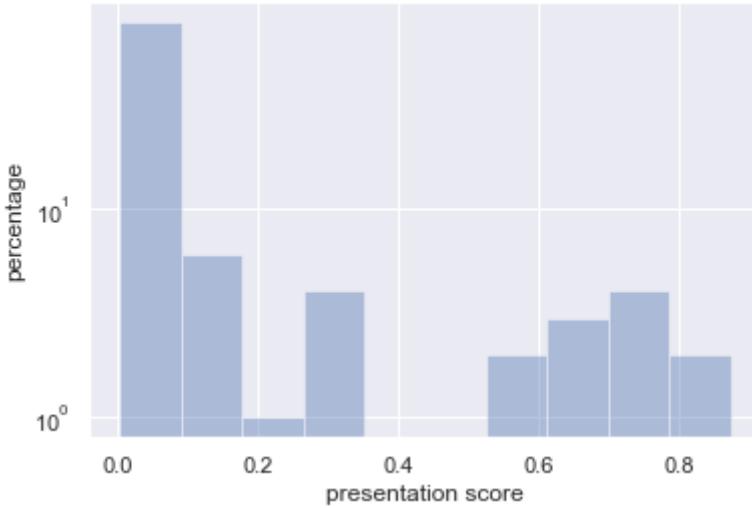
With CNN:





With FC:





- Molecule modification

For each step, the environment returns a property score as a reward.

The environment defines the maximum number of steps  $T$

The agent can choose "no modification" action to keep the molecules unchanged until the step  $T$ .

At step  $t$ , the reward will be discounted by  $\beta^{T-t}$  and  $\beta < 1$ , so that the final reward will have the largest contribution to the cumulative reward.

- Other directions

- choose the maximum reward

$$\text{finite horizon undiscounted return: } R(\tau) = \sum_{t=0}^T r_t$$

$$\text{infinite horizon discounted return: } R(\tau) = \sum_{t=0}^T \gamma^t r_t$$

The corresponding Bellman equation for  $Q^\pi(s, a)$  with the expected return defined as  $R(\tau) = \sum_{t=0}^{\infty} \gamma^t r_t$  is

$$Q^\pi(s_t, a_t) = \mathbb{E}_{\substack{s_{t+1} \sim P(\cdot | s_t, a_t) \\ a_{t+1} \sim \pi(\cdot | s_{t+1})}} [r(s_t, a_t) + \gamma Q^\pi(s_{t+1}, a_{t+1})]$$

optimize maximum reward achieved in a trajectory,

$$R(\tau) = \max_{t \geq 0} \gamma^t r_t$$

$\gamma < 1$  allows the agent to prefer rewards sooner rather than later.

The action-value function for a stationary policy  $\pi$ , denoted  $Q_{\max}^\pi : \mathcal{S} \times \mathcal{A} \rightarrow \mathbb{R}$ , represents the expected maximum of discounted rewards along the trajectories induced by the policy in the MDP:

$$Q_{\max}^\pi(s_t, a_t) \triangleq \mathbb{E} \left[ \max_{t' \geq t} (\gamma^{t'-t} r(s_{t'}, a_{t'})) \mid (s_t, a_t), \pi \right], \forall (s_t, a_t) \in \mathcal{S} \times \mathcal{A} \quad (2)$$

- New settings (discussion on 3/3)

action: 3 different actions (no modification/ modification; position; type)

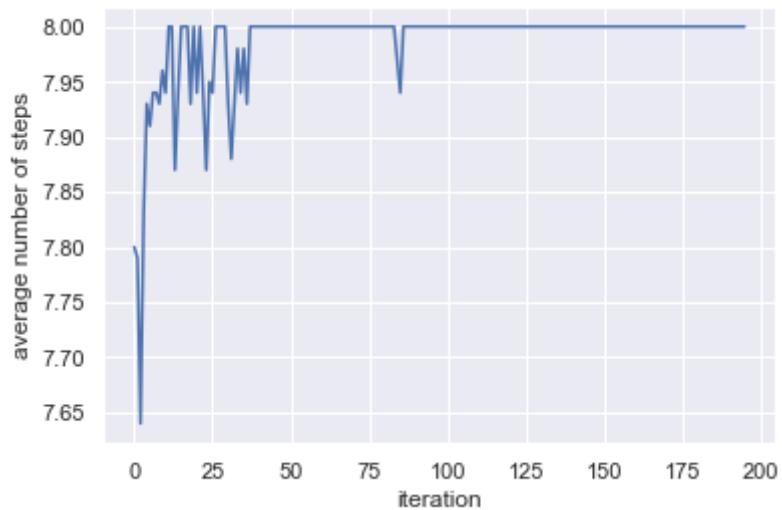
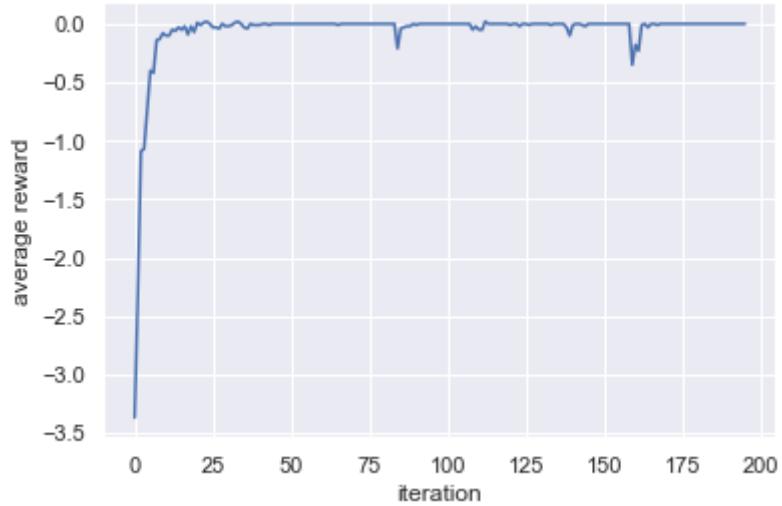
intermediate reward  $r_t$ : +1 (increase larger than  $\delta$ ); -1 (increase smaller than  $\delta$ ); 0 (no modification); ( $\delta = 0.1$ )

stop: reach the goal / maximum steps

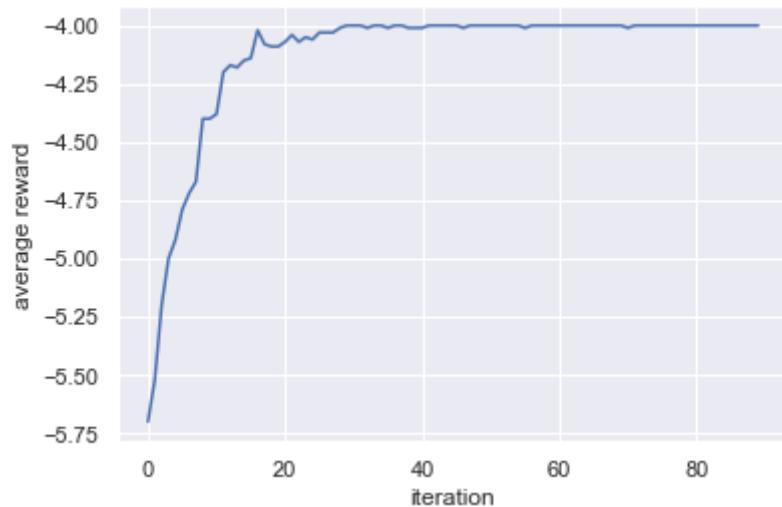
use this setting for pre-training, and then use the setting with only final rewards

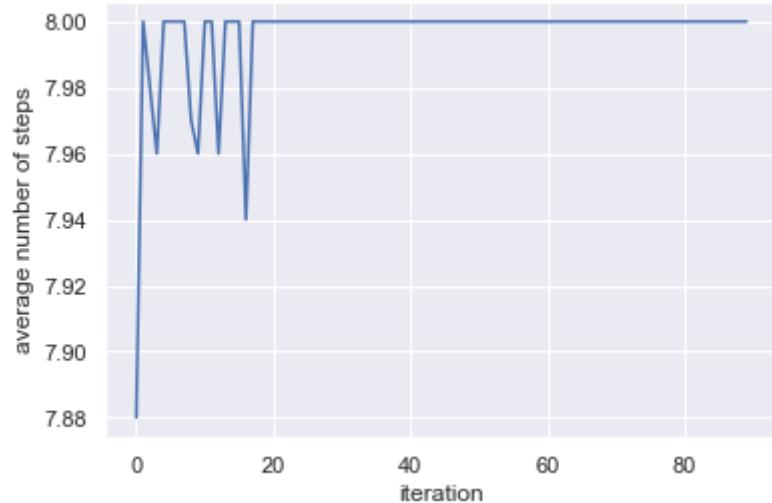
- Results with (no final reward + no penalty for "no modification" action)

The policy stuck into the "no modification" actions....

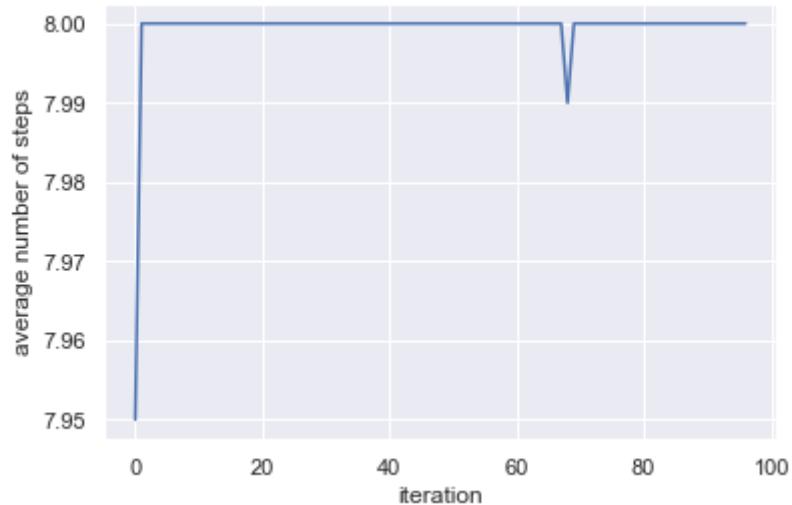
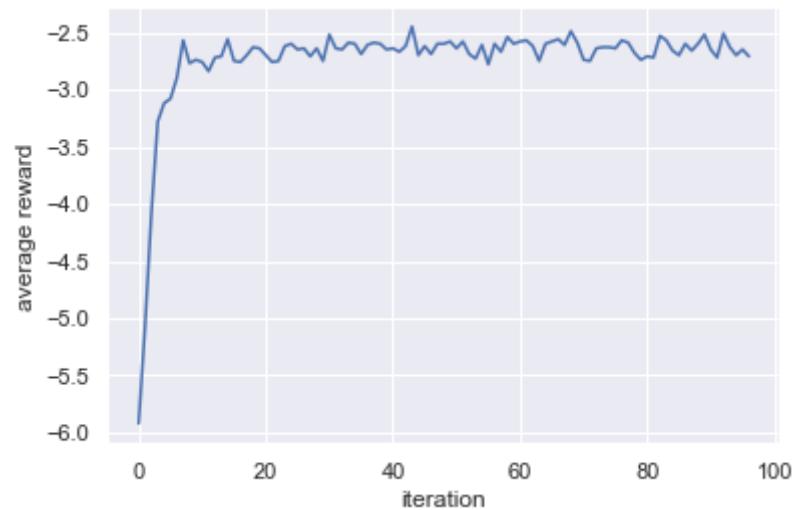


- o Results with (no final reward + penalty for "no modification action") (penalty = -0.5)

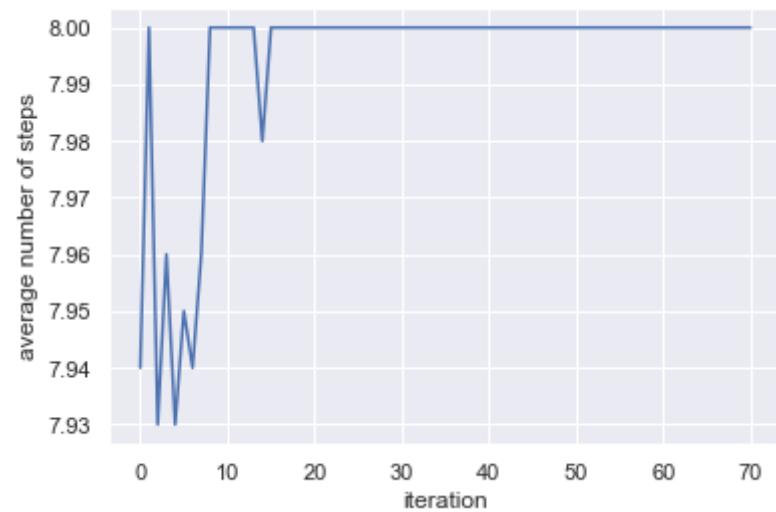
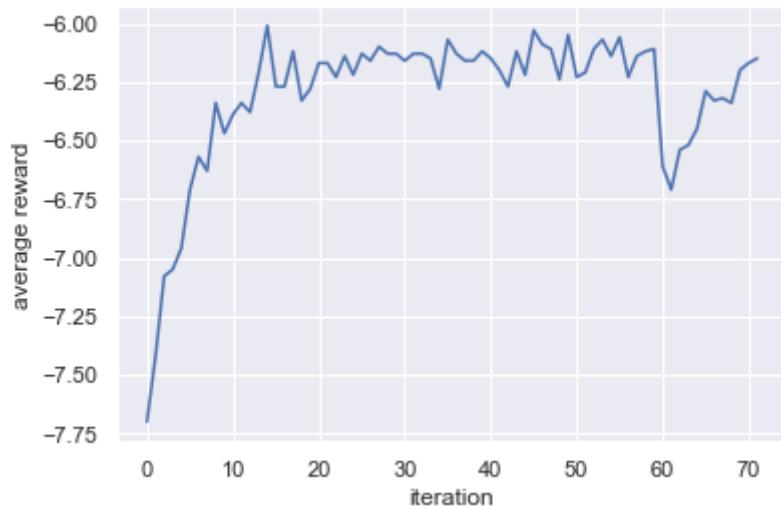




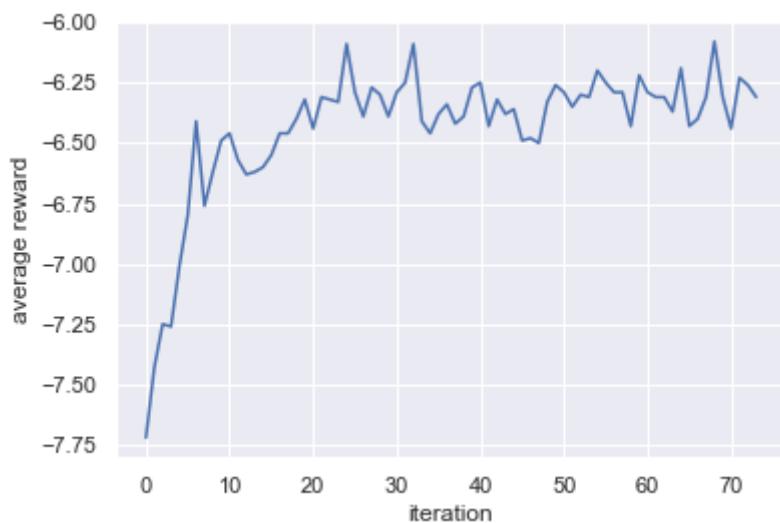
- Results with (final reward score \* 10 + no penalty for "no modification action")

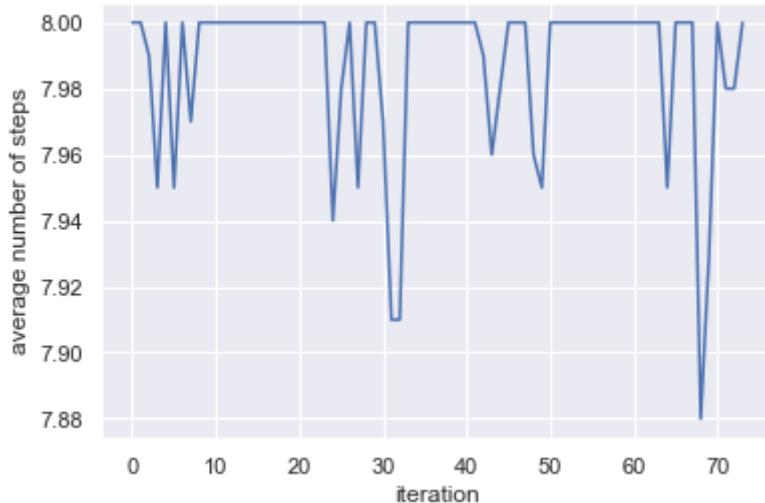


- Results with (final reward score \* 10 + penalty for "no modification action") (penalty = -0.5)



with large coefficients for entropy loss to encourage exploration:





- Pretrain policy network. (generate sequences for pretraining?)
- Edit motif position.
- remove "no modification" action from the action space and remove final reward.
- schedule exploration. use a large exploration ratio at the beginning.

## Update (4-14)

Add prior distribution into the sampling.

Given an allele, find positive peptides at different length from the IEDB dataset.

If the number of positive peptides at length  $l$  for allele  $a$  is greater than 10, then we derive the distribution of amino acids at each position; otherwise, we find the most similar allele (calculated from BLOSUM62 matrix) with positive peptides greater than 10 and directly use the distribution of that most similar allele.

After deriving the distribution of amino acids at each position, we can use the entropy of amino acid distribution to demonstrate whether or not those positive peptides have preference on some specific amino acids at the position. In other words, we can use the entropy to identify the motif position.

Therefore, with such a distribution, we can get the motif-position probability distribution  $p_p(pos|a)$  using the entropy, and get the amino acid distribution  $p_m(amino|a, pos)$  at a specific position.

At the first 50k steps, we directly combine the derived prior distribution with the calculated distribution from policy network.

$$pos \sim \text{Categorical}(p_p(pos|a, l) + p_\theta(pos|p, a))$$

$$amino \sim \text{Categorical}(p_m(amino|a, l, pos) + p_\theta(amino|p, a, pos))$$

(use these two distributions separately)

### Reward Design:

action: 3 different actions (no modification/ modification; position; type)

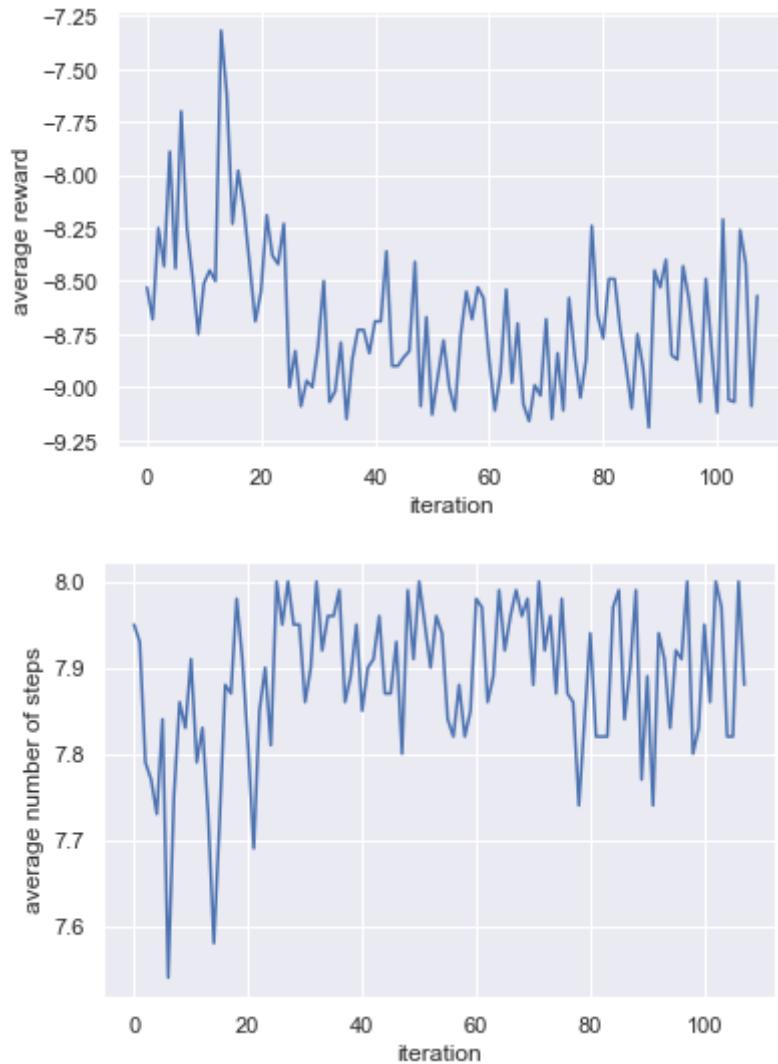
intermediate reward  $r_t$ : +1 (increase larger than  $\delta$ ); -1 (increase smaller than  $\delta$ ); -0.9 (no modification); ( $\delta = 0.05$ )

Final reward: present score \* 10

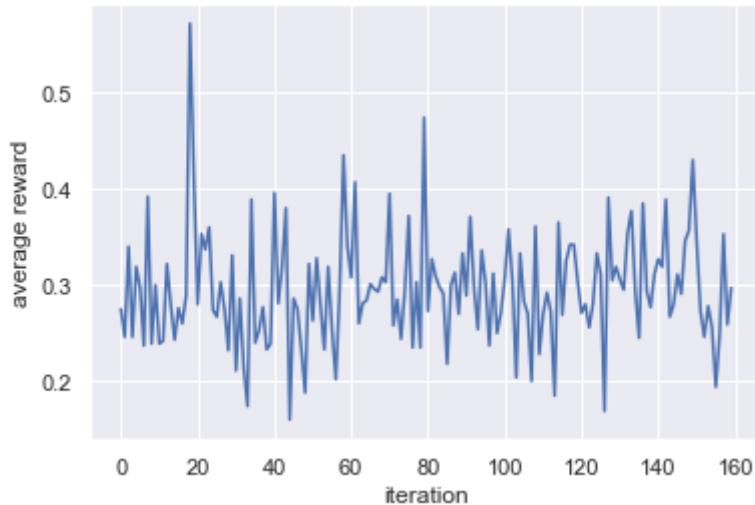
Max steps: 8

Goal: 0.6

50K step ~ 25 iteration with prior distribution.



Another setting from molecule modification (this setting doesn't perform very well....):  $\gamma^{T-i}$



### Conclusion:

The usage of prior distribution seems to be incorrect.

- either use the prior distribution or use the distribution from policy network to sample trajectories.
- use the prior distribution for more iterations.

## Update (4-21)

Objective:

Free generation without given alleles.

Given the Covid peptides, we find the mutated peptides for the given alleles?

Search space can be large even with the hard constraint (amino acid difference) or the soft constraint (Blosum62 matrix distance).

Plan:

1. Train the model for more iterations.
2. Update the prior distribution during training.

- The usage of prior distribution:

At each step, we sample the actions from the prior distribution with probability  $\beta$ , and from the distribution predicted by policy networks with probability  $(1 - \beta)$ .

We set the initial  $\beta = 0.5$  and decrease the value of  $\beta$  by 0.05 every 20k steps.

Each iteration has 2,048 steps.  $\beta = 0$  at iteration 100.

- A potential problem with the prior distributions is that:

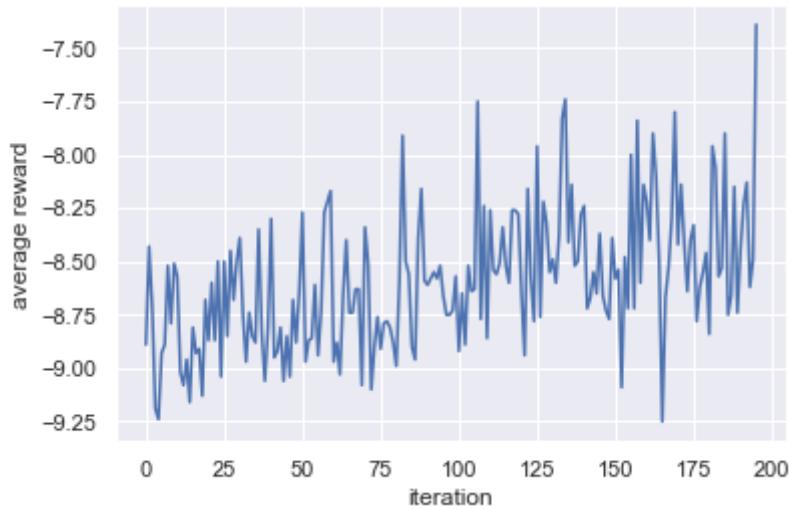
The derived prior distributions from the dataset can be very sparse.

For allele "BoLA-2\*12:01", the positive peptide sequences at length 11 in IEDB dataset are very similar:

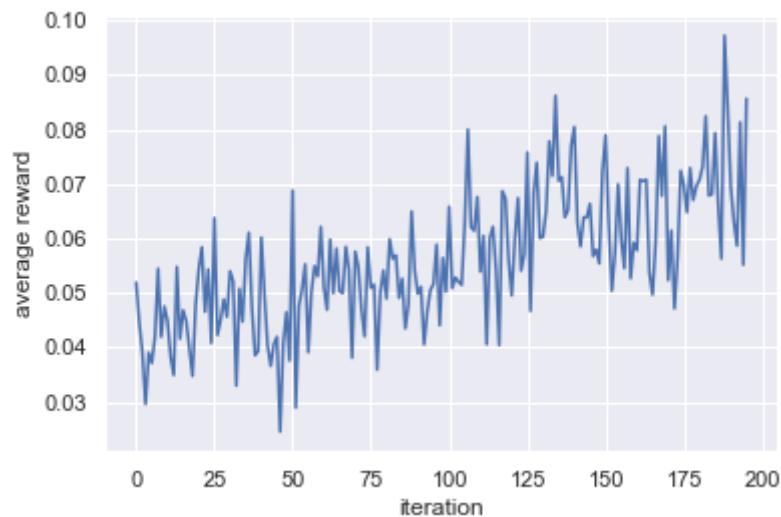
'KAAHGMGKVGK', 'KASHGMGKVGK', 'KSAHGMGKVGK', 'KSKHGMGKVGK',  
'KSSAGMGKVGK', 'KSSHAMGKVGK', 'KSSHGAGKVGK', 'KSSHGMAKVGK',  
'KSSHGMGAVGK', 'KSSHGMGEVGK', 'KSSHGMGKAGK', 'KSSHGMGKVAK'

- Setting 1 (+1, -1, -0.9)

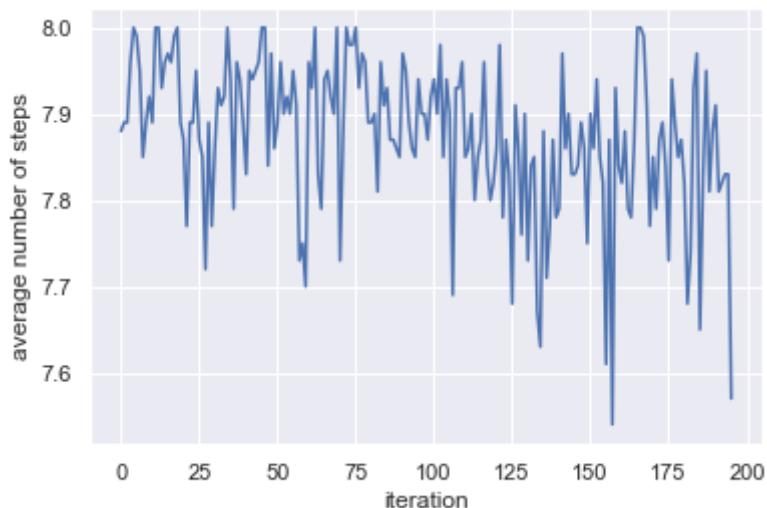
Cumulative reward (return):



Final reward /10 (presentation score):

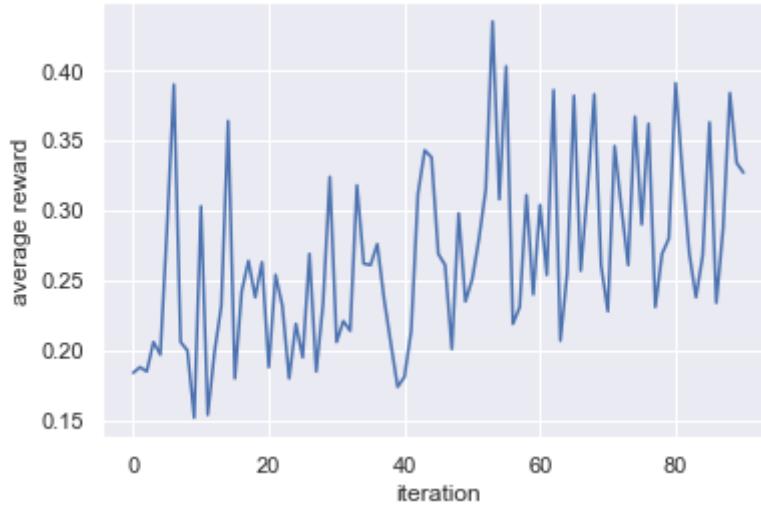


Average length:

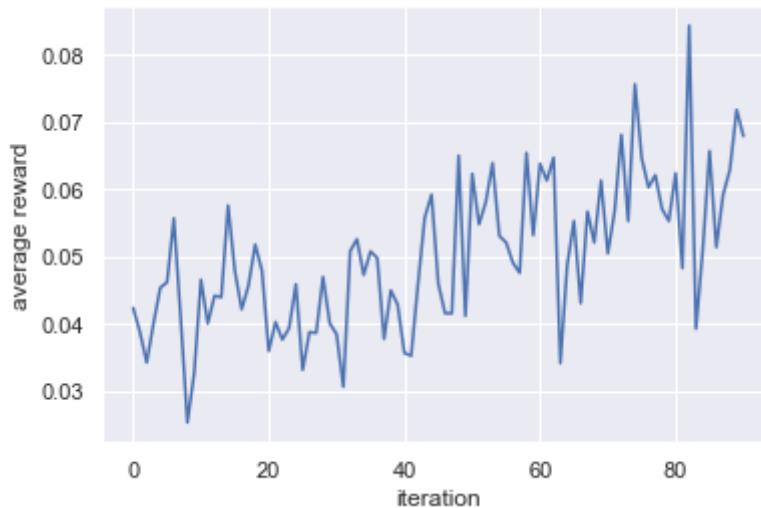


- Setting 2 ( $r_t = s * \gamma^{T-t}$ )

Cumulative reward:



Final reward:



## Update (5-5)

Next step:

1. Use the setting without intermediate reward.
2. Change the reward of "no modification" action (-0.5, 0, 0.5?)
3. Use the setting with intermediate reward first and then

Discussion:

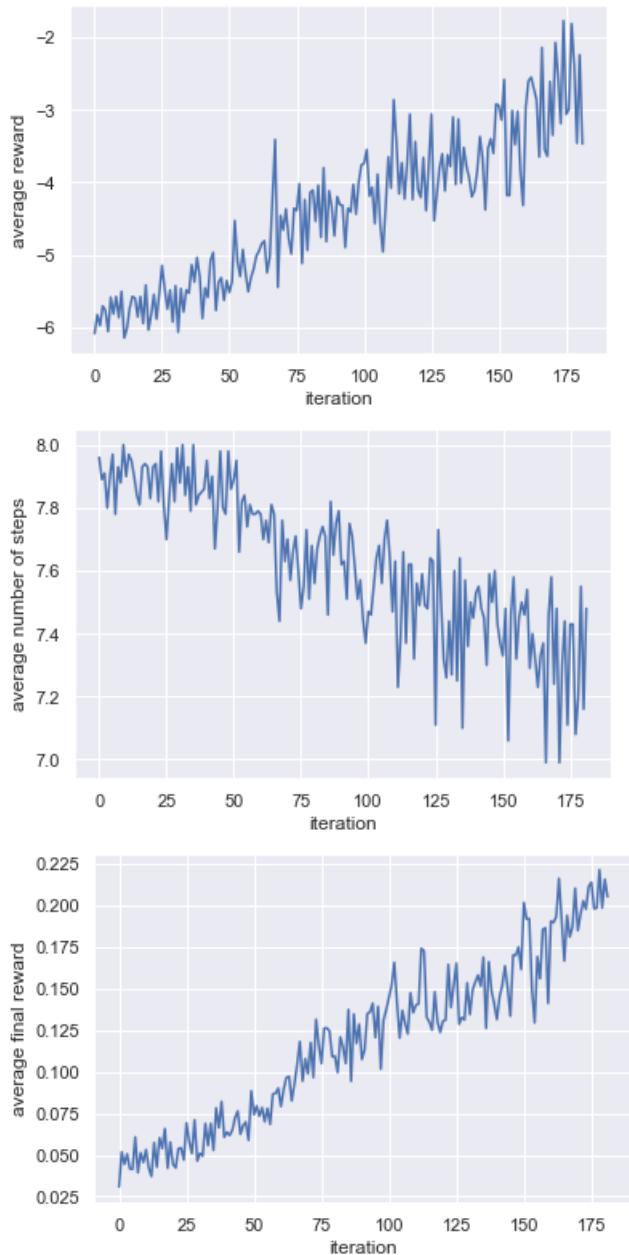
1. Given any allele, we want to generate a database of peptides that can trigger immune response.
2. Given any allele, we want to generate a database of peptides that can target the virus genome.

Title:

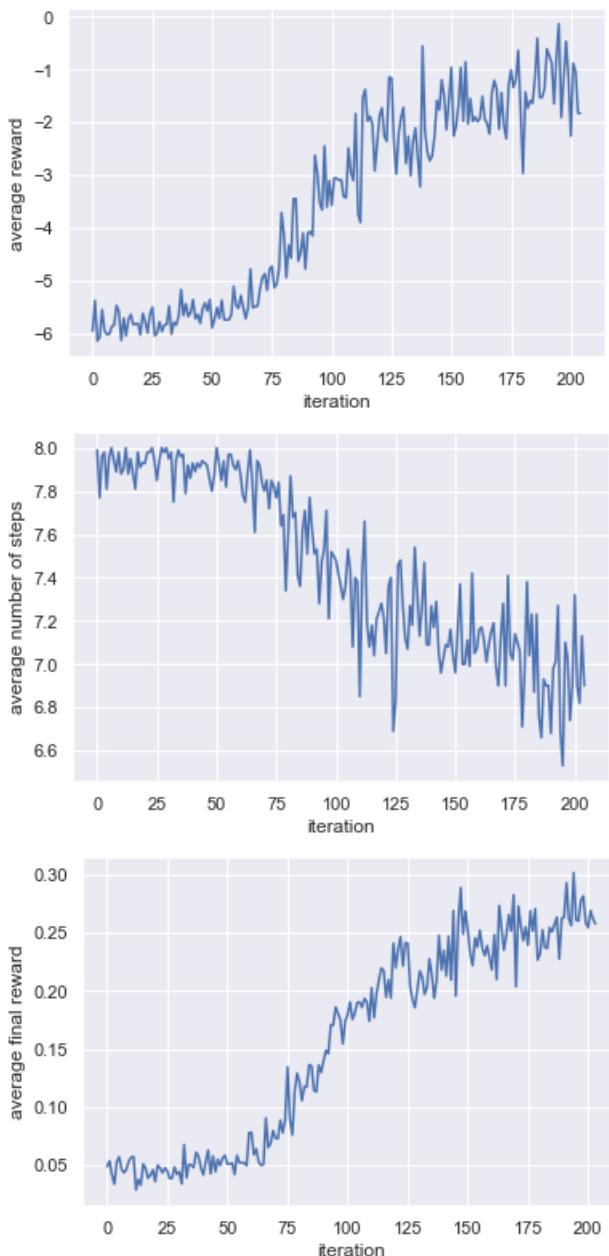
peptide mutation policy with combinatorial constraint (i.e., mutate the peptide sequences with the constraint of amino acid difference).

- Setting 1 (+1, -1, -0.9)

- Without warm start (dictionary)



- With warm start



Example:

```

terminal: False; action: 0,1,19; old_peptide: VTYLALLAAF; new_peptides: VYYLALLAAF;
allele: YFAMYEEESAAHTDVDTLYIIYRDYTWAARAYTWY; rewards: -1.0000; process_score:
0.0331; present_score: 0.0047; affinity: 0.0688;

terminal: False; action: 0,0,19; old_peptide: VYYLALLAAF; new_peptides: YYYLALLAAF;
allele: YFAMYEEESAAHTDVDTLYIIYRDYTWAARAYTWY; rewards: -1.0000; process_score:
0.0447; present_score: 0.0048; affinity: 0.0656;

terminal: False; action: 0,4,15; old_peptide: YYYLALLAAF; new_peptides: YYYLPLLAAF;
allele: YFAMYEEESAAHTDVDTLYIIYRDYTWAARAYTWY; rewards: -1.0000; process_score:
0.1511; present_score: 0.0074; affinity: 0.0712;

terminal: False; action: 0,3,15; old_peptide: YYYLPLLAAF; new_peptides: YYYPPPLLAAF;
allele: YFAMYEEESAAHTDVDTLYIIYRDYTWAARAYTWY; rewards: 1.0000; process_score:
0.7602; present_score: 0.1436; affinity: 0.1628;

terminal: False; action: 0,5,15; old_peptide: YYYPPPLLAAF; new_peptides: YYYPPPLAAF;
allele: YFAMYEEESAAHTDVDTLYIIYRDYTWAARAYTWY; rewards: 1.0000; process_score:
0.9098; present_score: 0.3412; affinity: 0.2197;

```

```

terminal: False; action: 0,6,15; old_peptide: YYYPPPLAAF; new_peptides: YYYPPPPAAF;
allele: YFAMYEEESAHTDVDTLYIYRDYTWAARAYTWY; rewards: 1.0000; process_score:
0.8753; present_score: 0.5267; affinity: 0.3041;

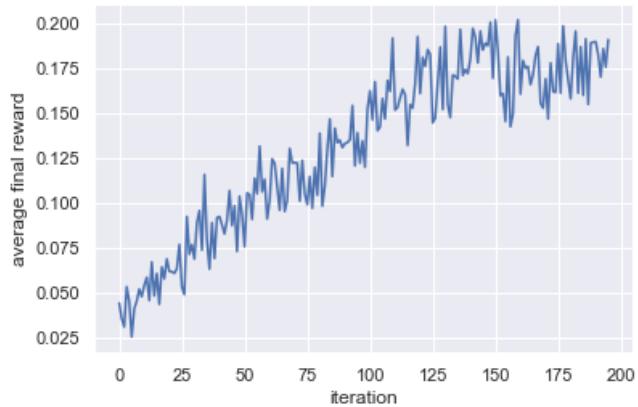
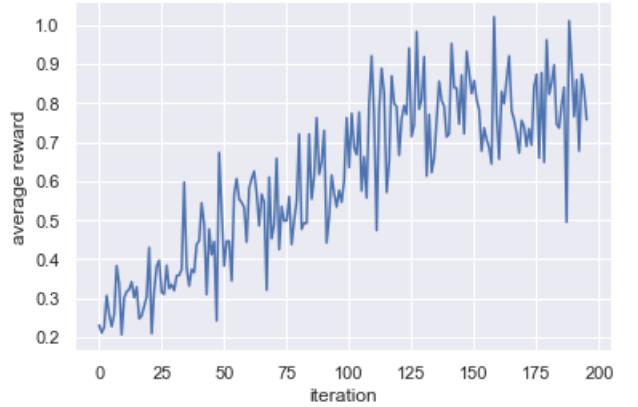
terminal: False; action: 0,7,15; old_peptide: YYYPPPPAAF; new_peptides: YYYPPPPPAF;
allele: YFAMYEEESAHTDVDTLYIYRDYTWAARAYTWY; rewards: -1.0000; process_score:
0.8238; present_score: 0.4826; affinity: 0.3046;

terminal: True; action: 0,2,15; old_peptide: YYYPPPPPAF; new_peptides: YYPPPPPAF;
allele: YFAMYEEESAHTDVDTLYIYRDYTWAARAYTWY; rewards: 3.5450; process_score:
0.4992; present_score: 0.3545; affinity: 0.3635;

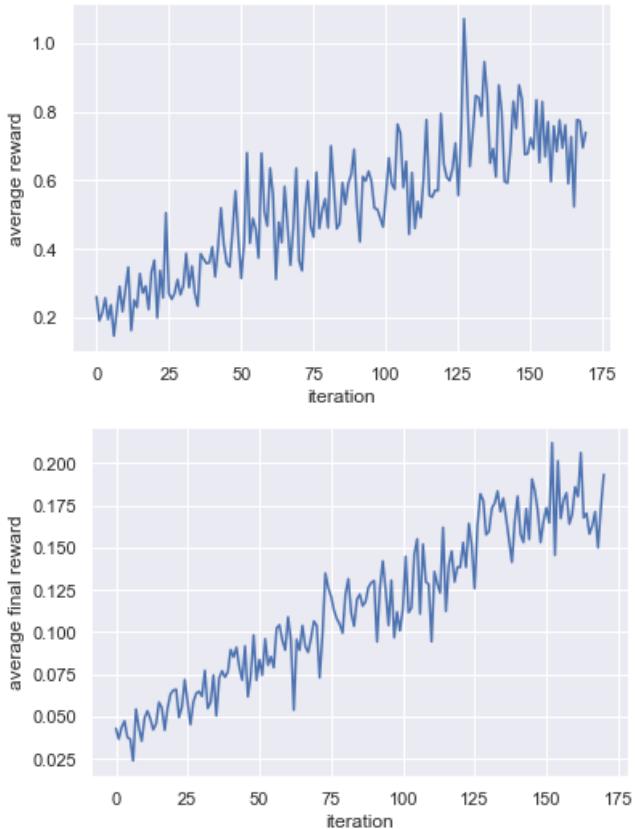
```

- Setting 2 ( $r_t = s * \gamma^{T-t}$ )

- Without warm start (dictionary)



- With warm start



terminal: False; action: 0,8,20; old\_peptide: GRMEKKTWK; new\_peptides: GRMEKKTWV; allele: YYTKYREISTNTYESNLYWRYNLYTWAEALAYLWY; rewards: 0.0058; process\_score: 0.2096; present\_score: 0.0121; affinity: 0.0990;

terminal: False; action: 0,7,20; old\_peptide: GRMEKKTWK; new\_peptides: GRMEKKTVV; allele: YYTKYREISTNTYESNLYWRYNLYTWAEALAYLWY; rewards: 0.0191; process\_score: 0.3005; present\_score: 0.0359; affinity: 0.1745;

terminal: False; action: 0,6,9; old\_peptide: GRMEKKTVV; new\_peptides: GRMEKKHVV; allele: YYTKYREISTNTYESNLYWRYNLYTWAEALAYLWY; rewards: 0.0131; process\_score: 0.1694; present\_score: 0.0222; affinity: 0.1715;

terminal: False; action: 0,1,20; old\_peptide: GRMEKKHVV; new\_peptides: GVMEKKHVV; allele: YYTKYREISTNTYESNLYWRYNLYTWAEALAYLWY; rewards: 0.0352; process\_score: 0.2614; present\_score: 0.0536; affinity: 0.2274;

terminal: False; action: 0,4,9; old\_peptide: GVMEKKHVV; new\_peptides: GVMEHKHVV; allele: YYTKYREISTNTYESNLYWRYNLYTWAEALAYLWY; rewards: 0.0510; process\_score: 0.2005; present\_score: 0.0699; affinity: 0.2748;

terminal: False; action: 0,4,20; old\_peptide: GVMEHKHVV; new\_peptides: GVMEVKHVV; allele: YYTKYREISTNTYESNLYWRYNLYTWAEALAYLWY; rewards: 0.1047; process\_score: 0.4881; present\_score: 0.1292; affinity: 0.2428;

terminal: False; action: 0,5,20; old\_peptide: GVMEVKHVV; new\_peptides: GVMEVHVV; allele: YYTKYREISTNTYESNLYWRYNLYTWAEALAYLWY; rewards: 0.5094; process\_score: 0.9276; present\_score: 0.5660; affinity: 0.3016;

terminal: True; action: 0,0,9; old\_peptide: GVMEVHVV; new\_peptides: HVMEVHVV; allele: YYTKYREISTNTYESNLYWRYNLYTWAEALAYLWY; rewards: 0.6848; process\_score: 0.9049; present\_score: 0.6848; affinity: 0.3578;

## Update (5-12)

Suggestions:

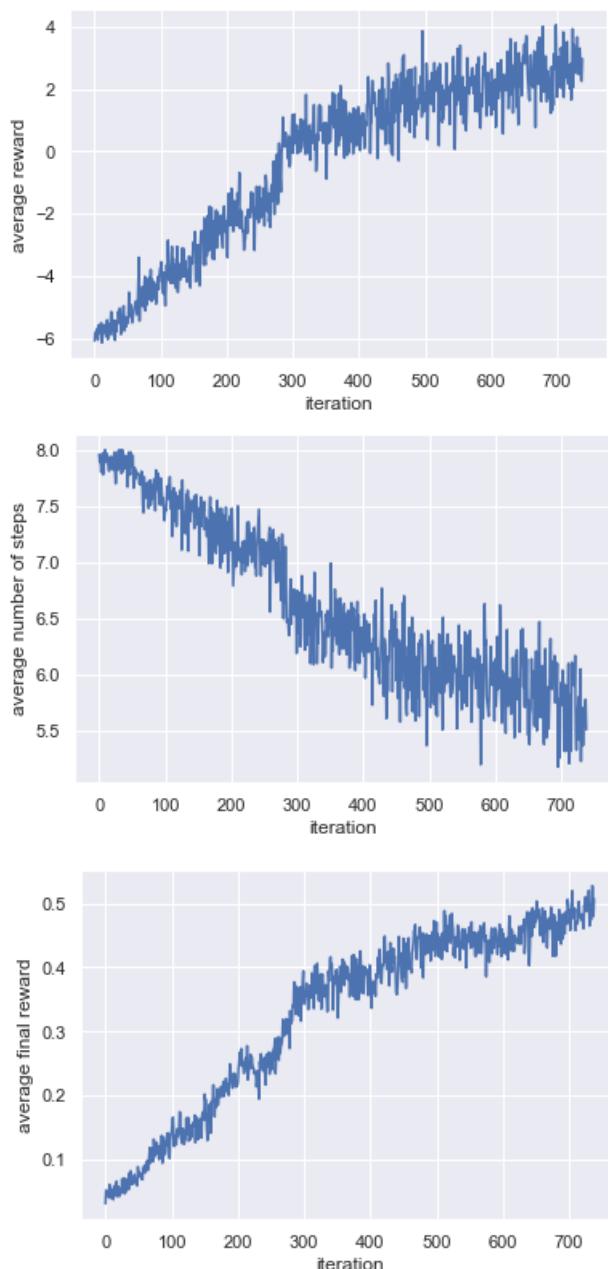
1. Use a larger episode length
2. Set the threshold as 0.75

### 1. Result Analysis with longer running time

- Setting 1 (+1, -1, -0.9) without warm start (i.e., prior distribution dictionary)

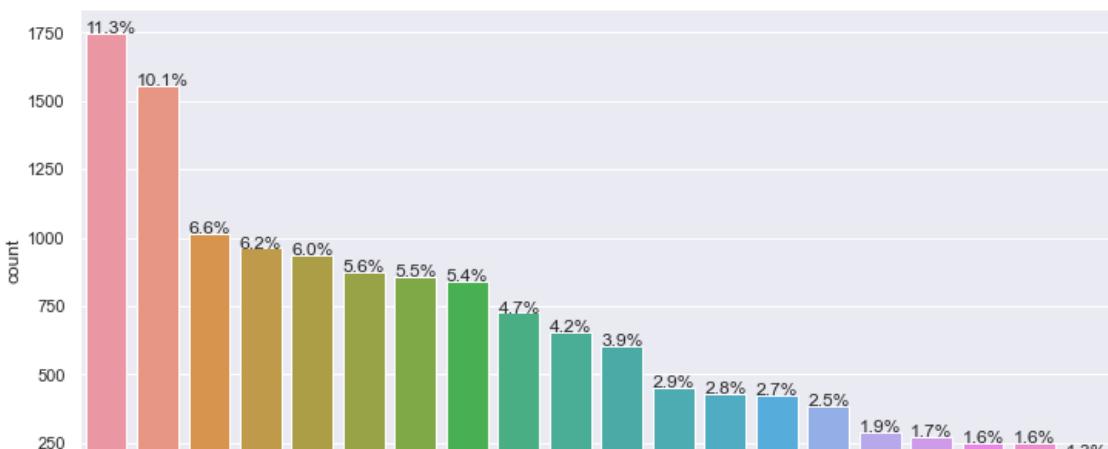
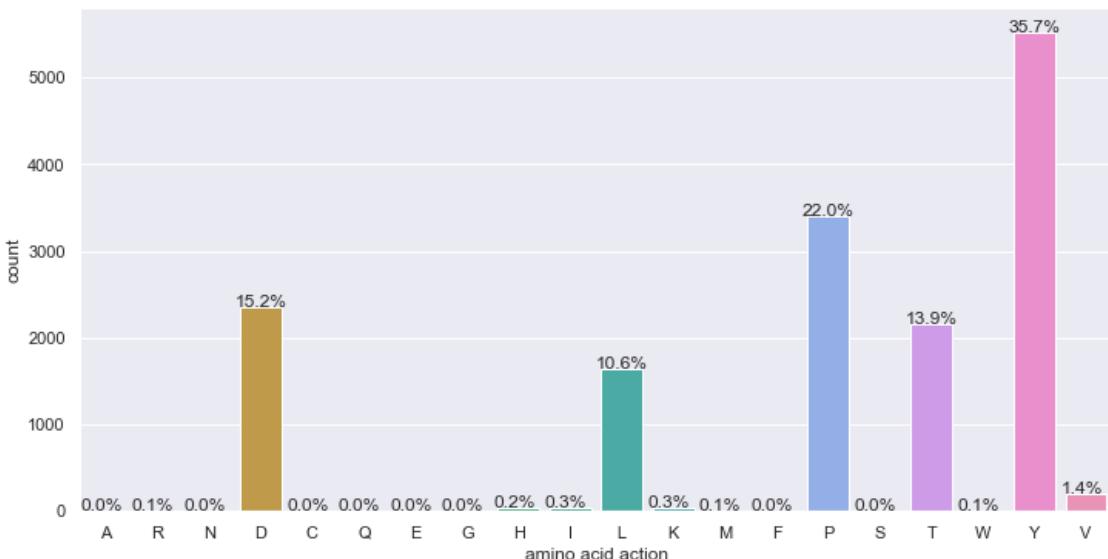
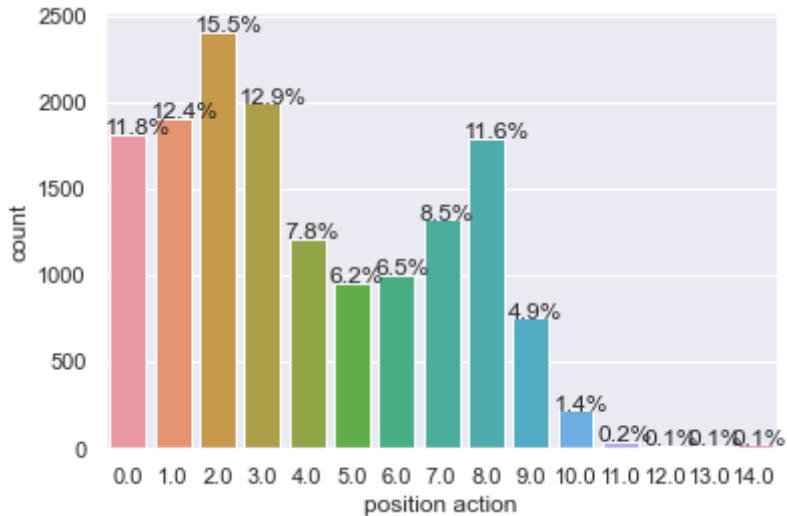
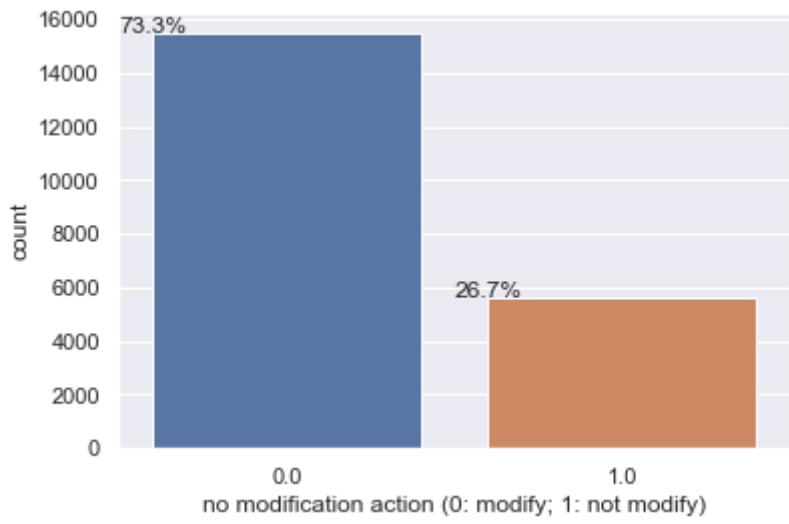
Among all the samples (i.e., the initial state of each episode/trajectory), 80% of them are sampled from the IEDB dataset; 20% from randomly generated peptides.

It takes 60 hours to get the results below.



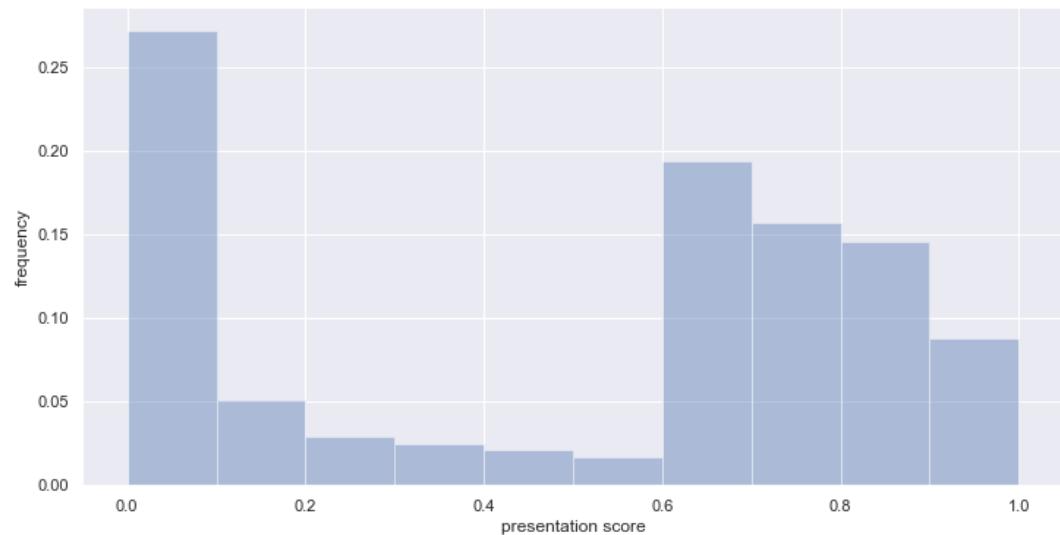
Each iteration has 2,112 steps. In each step, a specific given state will be updated by the policy network. Within 739 iterations, we have 238,343 samples in total. In order to analyze the trained policy network, we select the corresponding optimized peptides (i.e., the peptides in the terminal state) of the last 3,750 (about 5%) samples in the last 10 iterations for further analysis.

- Analysis 1: action analysis

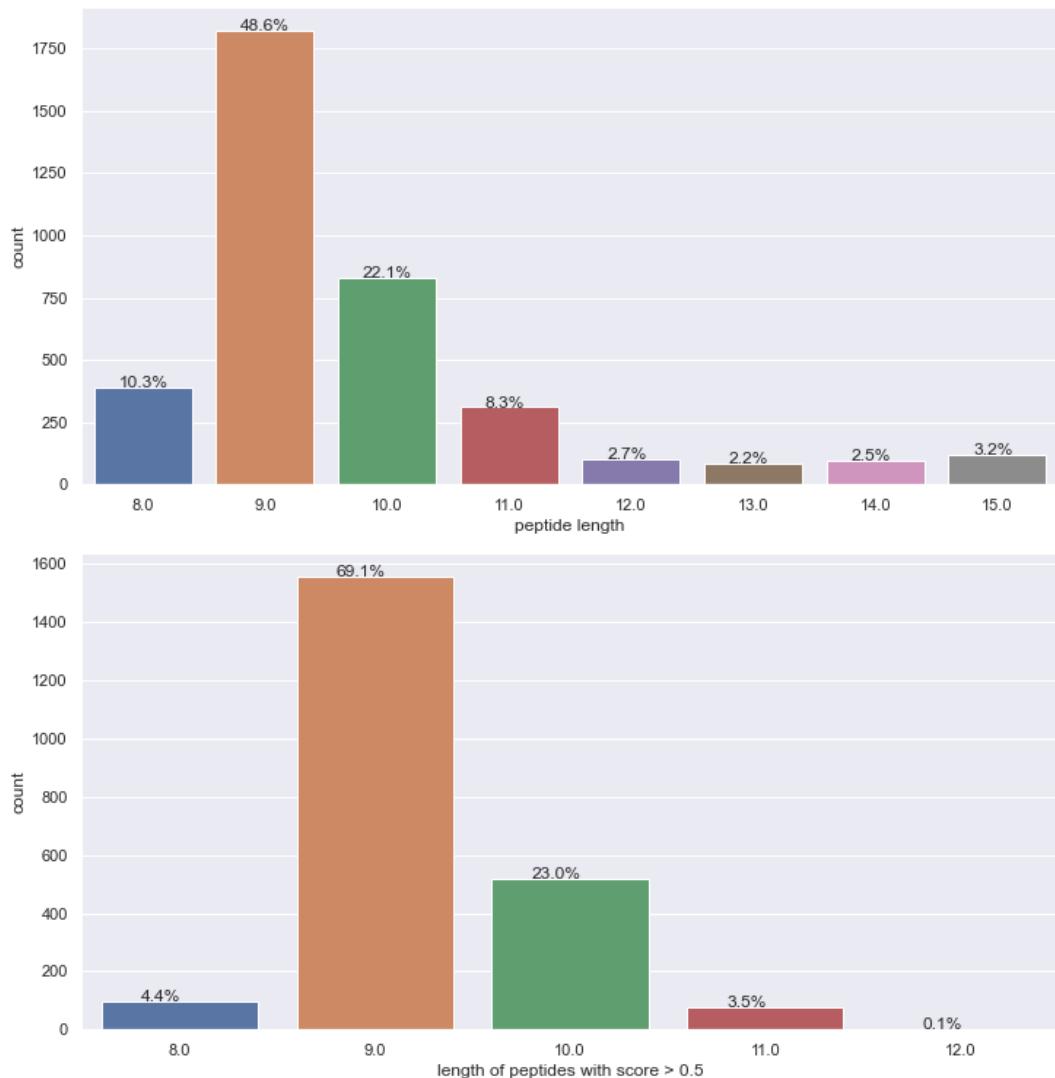




- Analysis 2: Reward distribution (60.16% optimized peptides have the presentation scores > 0.5)



- Analysis 3: peptide length distribution



Examples:

- Episodes with many "no modification" actions

terminal: False; action: 0,7,11; old\_peptide: WGCCPEVS; new\_peptides: WGCCPEVL; allele: YYAGYREKYRQTDVSNLYIRYDYYTWAESWAYLWY; rewards: -1.0000; process\_score: 0.0054; present\_score: 0.0059; affinity: 0.0984;  
terminal: False; action: 1,0,1; old\_peptide: WGCCPEVL; new\_peptides: WGCCPEVL; allele: YYAGYREKYRQTDVSNLYIRYDYYTWAESWAYLWY; rewards: -0.9000; process\_score: 0.0054; present\_score: 0.0059; affinity: 0.0984;  
terminal: False; action: 1,0,1; old\_peptide: WGCCPEVL; new\_peptides: WGCCPEVL; allele: YYAGYREKYRQTDVSNLYIRYDYYTWAESWAYLWY; rewards: -0.9000; process\_score: 0.0054; present\_score: 0.0059; affinity: 0.0984;  
terminal: False; action: 1,0,1; old\_peptide: WGCCPEVL; new\_peptides: WGCCPEVL; allele: YYAGYREKYRQTDVSNLYIRYDYYTWAESWAYLWY; rewards: -0.9000; process\_score: 0.0054; present\_score: 0.0059; affinity: 0.0984;  
terminal: False; action: 1,0,1; old\_peptide: WGCCPEVL; new\_peptides: WGCCPEVL; allele: YYAGYREKYRQTDVSNLYIRYDYYTWAESWAYLWY; rewards: -0.9000; process\_score: 0.0054; present\_score: 0.0059; affinity: 0.0984;  
terminal: False; action: 0,2,15; old\_peptide: WGCCPEVL; new\_peptides: WGPCPEVL; allele: YYAGYREKYRQTDVSNLYIRYDYYTWAESWAYLWY; rewards: -1.0000; process\_score: 0.0208; present\_score: 0.0148; affinity: 0.1822;  
terminal: False; action: 1,0,1; old\_peptide: WGPCPEVL; new\_peptides: WGPCPEVL; allele: YYAGYREKYRQTDVSNLYIRYDYYTWAESWAYLWY; rewards: -0.9000; process\_score: 0.0208; present\_score: 0.0148; affinity: 0.1822;  
terminal: True; action: 1,0,1; old\_peptide: WGPCPEVL; new\_peptides: WGPCPEVL; allele: YYAGYREKYRQTDVSNLYIRYDYYTWAESWAYLWY; rewards: 0.1481; process\_score: 0.0208; present\_score: 0.0148; affinity: 0.1822;

---

terminal: False; action: 1,0,1; old\_peptide: PVINVQDL; new\_peptides: PVINVQDL; allele: YFAMYGEKVAHTHVDTLVRYQYYTWAFLAYTWY; rewards: -0.9000; process\_score: 0.0246; present\_score: 0.0039; affinity: 0.0524;  
terminal: False; action: 1,0,1; old\_peptide: PVINVQDL; new\_peptides: PVINVQDL; allele: YFAMYGEKVAHTHVDTLVRYQYYTWAFLAYTWY; rewards: -0.9000; process\_score: 0.0246; present\_score: 0.0039; affinity: 0.0524;  
terminal: False; action: 0,6,17; old\_peptide: PVINVQDL; new\_peptides: PVINVQTL; allele: YFAMYGEKVAHTHVDTLVRYQYYTWAFLAYTWY; rewards: -1.0000; process\_score: 0.1038; present\_score: 0.0055; affinity: 0.0589;  
terminal: False; action: 1,0,1; old\_peptide: PVINVQTL; new\_peptides: PVINVQTL; allele: YFAMYGEKVAHTHVDTLVRYQYYTWAFLAYTWY; rewards: -0.9000; process\_score: 0.1038; present\_score: 0.0055; affinity: 0.0589;  
terminal: False; action: 0,1,15; old\_peptide: PVINVQTL; new\_peptides: PPINVQTL; allele: YFAMYGEKVAHTHVDTLVRYQYYTWAFLAYTWY; rewards: -1.0000; process\_score: 0.0141; present\_score: 0.0038; affinity: 0.0546;  
terminal: False; action: 0,0,19; old\_peptide: PPINVQTL; new\_peptides: YPINVQTL; allele: YFAMYGEKVAHTHVDTLVRYQYYTWAFLAYTWY; rewards: -1.0000; process\_score: 0.4524; present\_score: 0.0213; affinity: 0.0720;  
terminal: False; action: 0,3,4; old\_peptide: YPINVQTL; new\_peptides: YPIDVQTL; allele: YFAMYGEKVAHTHVDTLVRYQYYTWAFLAYTWY; rewards: -1.0000; process\_score: 0.3719; present\_score: 0.0181; affinity: 0.0836;  
terminal: True; action: 0,5,17; old\_peptide: YPIDVQTL; new\_peptides: YPIDVTTL; allele: YFAMYGEKVAHTHVDTLVRYQYYTWAFLAYTWY; rewards: 0.1499; process\_score: 0.3121; present\_score: 0.0150; affinity: 0.0853;

## 2. Examples with large final rewards

terminal: False; action: 0,0,19; old\_peptide: KEEESEGAIW; new\_peptides: YEEESEGAIW; allele: YYAEYRNIYAQTDESNLYIRYDSYTAWVLAYLWY; rewards: -1.0000; process\_score: 0.3170; present\_score: 0.0973; affinity: 0.2699;  
terminal: False; action: 0,2,19; old\_peptide: YEEESEGAIW; new\_peptides: YEYESEGAIW; allele: YYAEYRNIYAQTDESNLYIRYDSYTAWVLAYLWY; rewards: 1.0000; process\_score: 0.7595; present\_score: 0.5092; affinity: 0.3364;  
terminal: False; action: 0,9,11; old\_peptide: YEYESEGAIW; new\_peptides: YEYESEGAIL; allele: YYAEYRNIYAQTDESNLYIRYDSYTAWVLAYLWY; rewards: -1.0000; process\_score: 0.6002; present\_score: 0.4058; affinity: 0.3502;  
terminal: False; action: 0,1,19; old\_peptide: YEYESEGAIL; new\_peptides: YYYESEGAIL; allele: YYAEYRNIYAQTDESNLYIRYDSYTAWVLAYLWY; rewards: 1.0000; process\_score: 0.6715; present\_score: 0.4866; affinity: 0.3574;  
terminal: True; action: 0,3,4; old\_peptide: YYYESEGAIL; new\_peptides: YYYDSEGAIL; allele: YYAEYRNIYAQTDESNLYIRYDSYTAWVLAYLWY; rewards: 7.7289; process\_score: 0.7591; present\_score: 0.7729; affinity: 0.4496;

---

terminal: False; action: 0,8,11; old\_peptide: STAGVNNG; new\_peptides: STAGVNVL; allele: YYSKYRNIYAQTDESNLYLSYDYYTWAERAYEWY; rewards: -1.0000; process\_score: 0.1842; present\_score: 0.0533; affinity: 0.2529;  
terminal: False; action: 0,3,15; old\_peptide: STAGVNVL; new\_peptides: STLPGVNVL; allele: YYSKYRNIYAQTDESNLYLSYDYYTWAERAYEWY; rewards: 1.0000; process\_score: 0.3234; present\_score: 0.2261; affinity: 0.3626;  
terminal: True; action: 0,2,4; old\_peptide: STLPGVNVL; new\_peptides: STDPGVNVL; allele: YYSKYRNIYAQTDESNLYLSYDYYTWAERAYEWY; rewards: 7.3298; process\_score: 0.7001; present\_score: 0.7330; affinity: 0.4490;

## 3. Examples with final presentation score < 0.5

Problem: how to encourage examples like this to select as many "no modification" actions as possible?

terminal: False; action: 0,4,15; old\_peptide: TLVQYMDDI; new\_peptides: TLVQPMDDI; allele: YHTKYREISTNTYESNLWRYNLHTWAELAYLWY; rewards: -1.0000; process\_score: 0.0111; present\_score: 0.0044; affinity: 0.0687;  
terminal: False; action: 0,6,17; old\_peptide: TLVQPMDDI; new\_peptides: TLVQPMTDI; allele: YHTKYREISTNTYESNLWRYNLHTWAELAYLWY; rewards: -1.0000; process\_score: 0.0426; present\_score: 0.0052; affinity: 0.0740;  
terminal: False; action: 0,7,17; old\_peptide: TLVQPMTDI; new\_peptides: TLVQPMTTI; allele: YHTKYREISTNTYESNLWRYNLHTWAELAYLWY; rewards: -1.0000; process\_score: 0.2484; present\_score: 0.0133; affinity: 0.0952;  
terminal: False; action: 0,3,15; old\_peptide: TLVQPMTTI; new\_peptides: TLVPPMTTI; allele: YHTKYREISTNTYESNLWRYNLHTWAELAYLWY; rewards: -1.0000; process\_score: 0.1736; present\_score: 0.0123; affinity: 0.1131;  
terminal: False; action: 0,2,4; old\_peptide: TLVPPMTTI; new\_peptides: TLDPPMTTI; allele: YHTKYREISTNTYESNLWRYNLHTWAELAYLWY; rewards: -1.0000; process\_score: 0.5161; present\_score: 0.0623; affinity: 0.1569;  
terminal: False; action: 0,5,15; old\_peptide: TLDPPMTTI; new\_peptides: TLDPPPTTI; allele: YHTKYREISTNTYESNLWRYNLHTWAELAYLWY; rewards: 1.0000; process\_score: 0.6149; present\_score: 0.1562; affinity: 0.2211;  
terminal: False; action: **0,1,19**; old\_peptide: TLDPPPTTI; new\_peptides: TYDPPPTTI; allele:

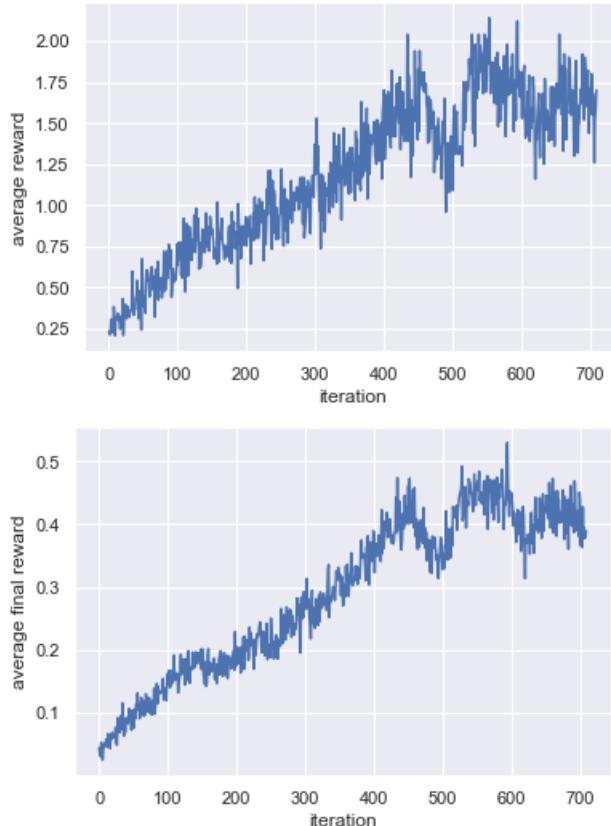
```

YHTKYREISTNTYESNLYWRYNLHTWAELAYLWY; rewards: -1.0000; process_score: 0.5693;
present_score: 0.1088; affinity: 0.1969;
terminal: True; action: 0,1,4; old_peptide: TYDPPPPTI; new_peptides: TDDPPPPTI; allele:
YHTKYREISTNTYESNLYWRYNLHTWAELAYLWY; rewards: 3.1405; process_score: 0.1116;
present_score: 0.3140; affinity: 0.4767;

```

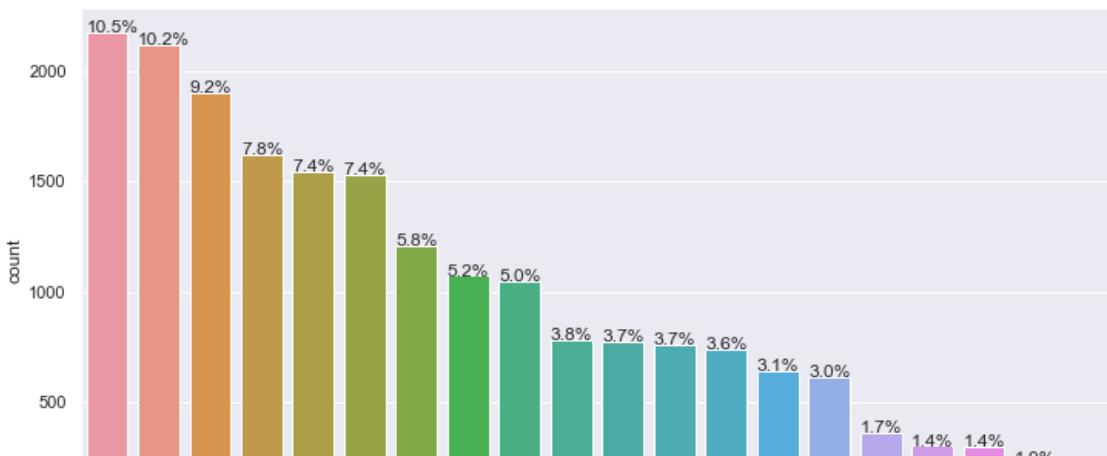
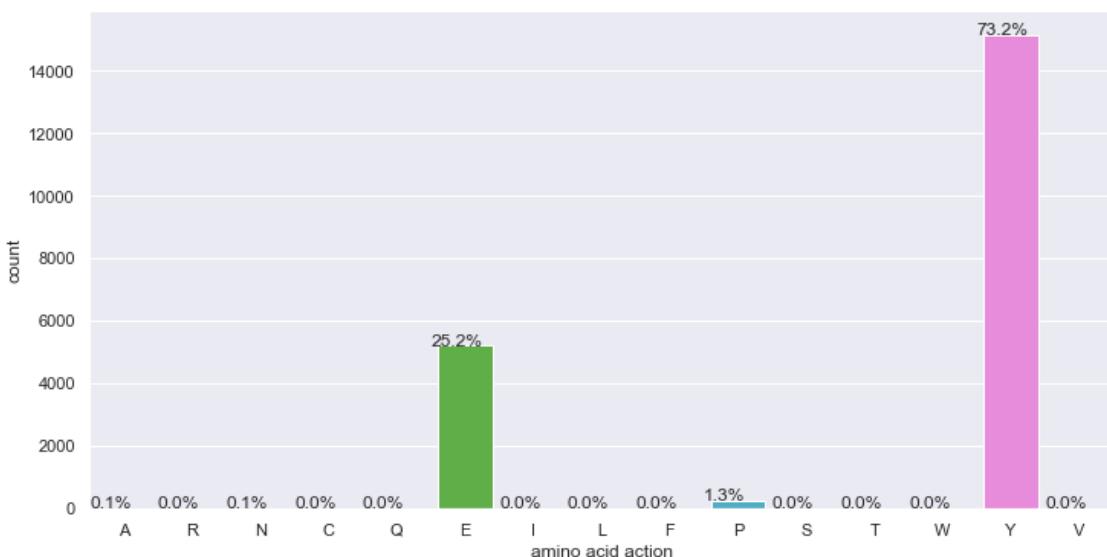
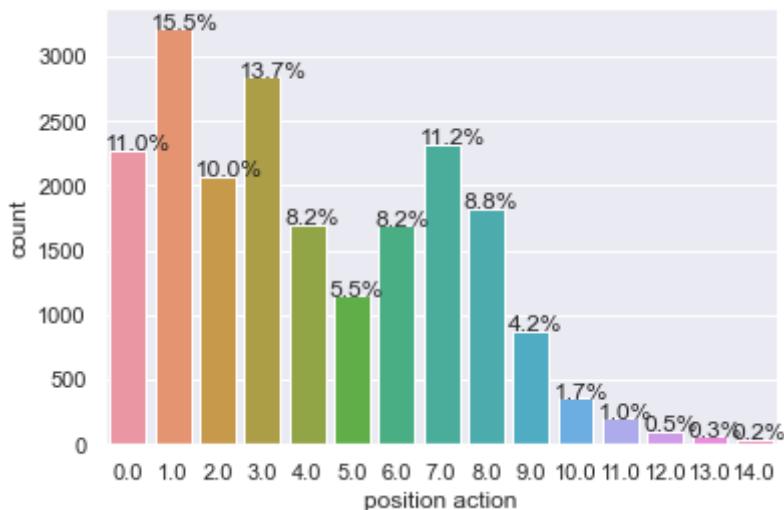
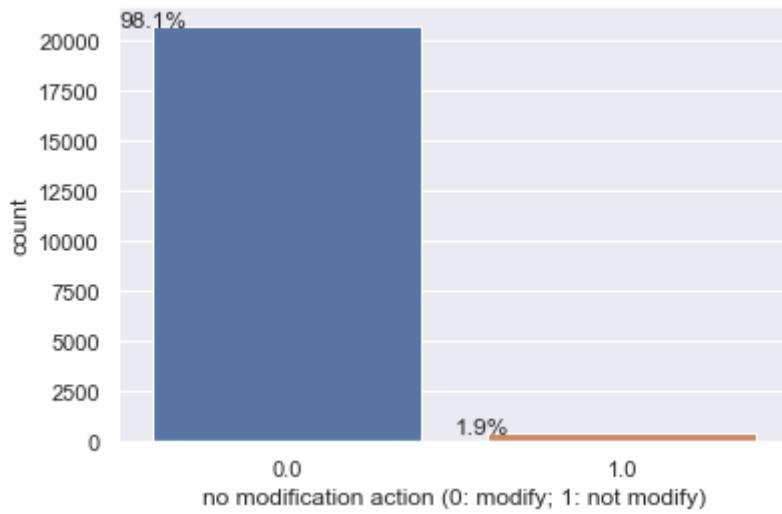
- Setting 2 ( $r_t = s * \gamma^{T-t}$ ) without warm start (dictionary)

Among all the samples (i.e., the initial state of each episode/trajectory), 80% of them are sampled from the IEDB dataset; 20% from randomly generated peptides.



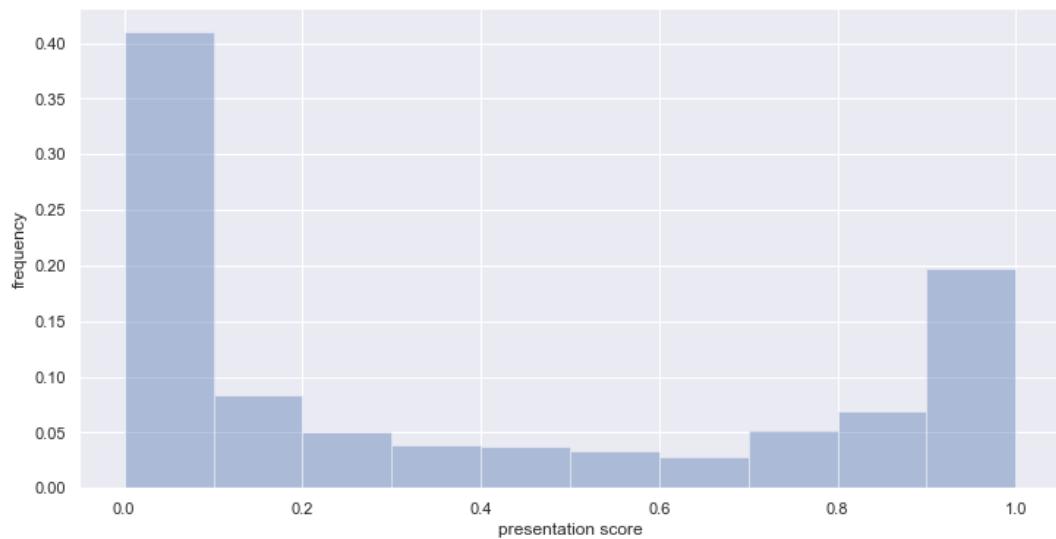
Each iteration has 2,112 steps. In each step, a specific given state will be updated by the policy network. Within 710 iterations, we have 187,440 samples in total. In order to analyze the trained policy network, we select the corresponding optimized peptides (i.e., the peptides in the terminal state) of the last 2,640 (about 1.4%) samples in the last 10 iterations for further analysis.

- Analysis 1: action analysis

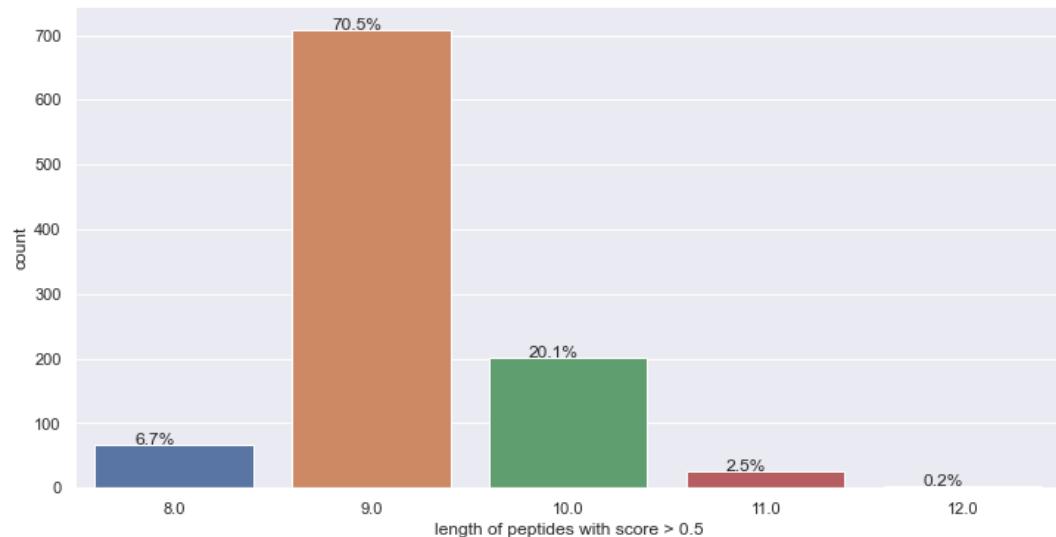
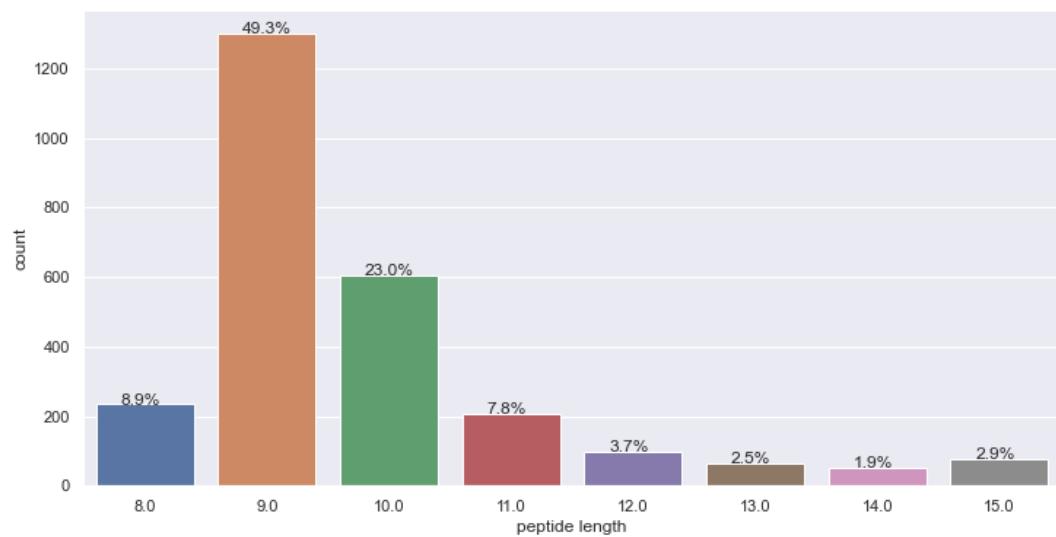




- Analysis 2: Reward distribution (38.07% optimized peptides have the presentation scores > 0.5)

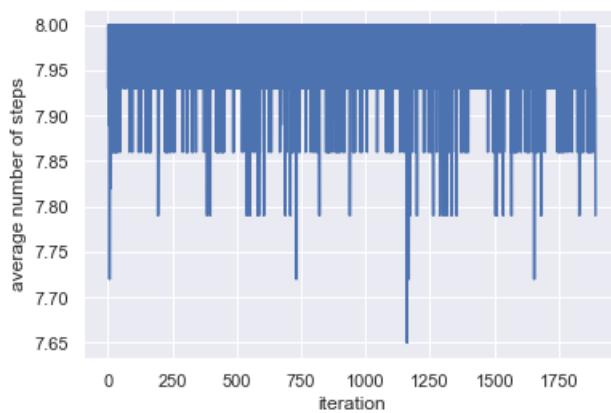
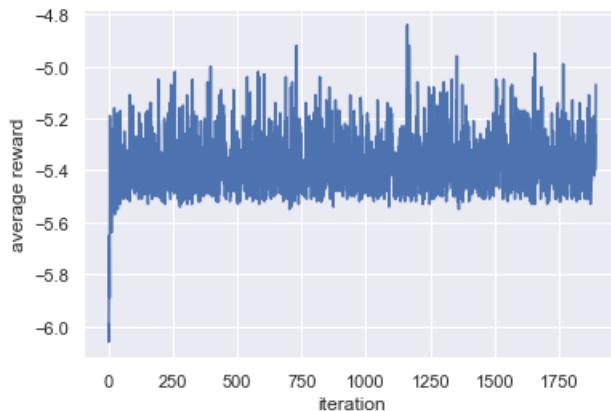


- Analysis 3: peptide length distribution

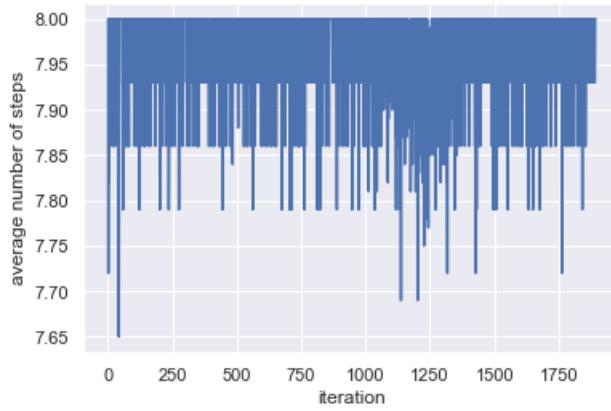
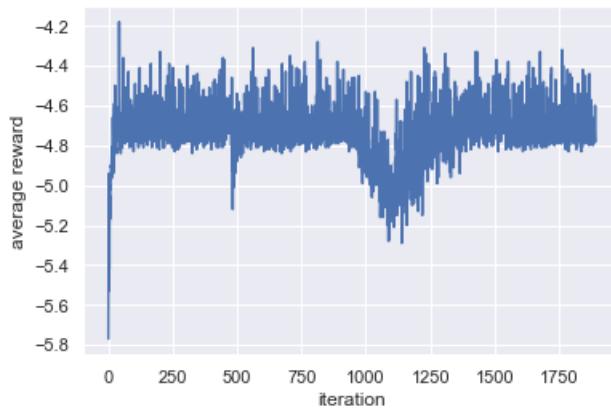


## 2. Result Analysis with smaller penalty (i.e., greater reward) for "no modification" action.

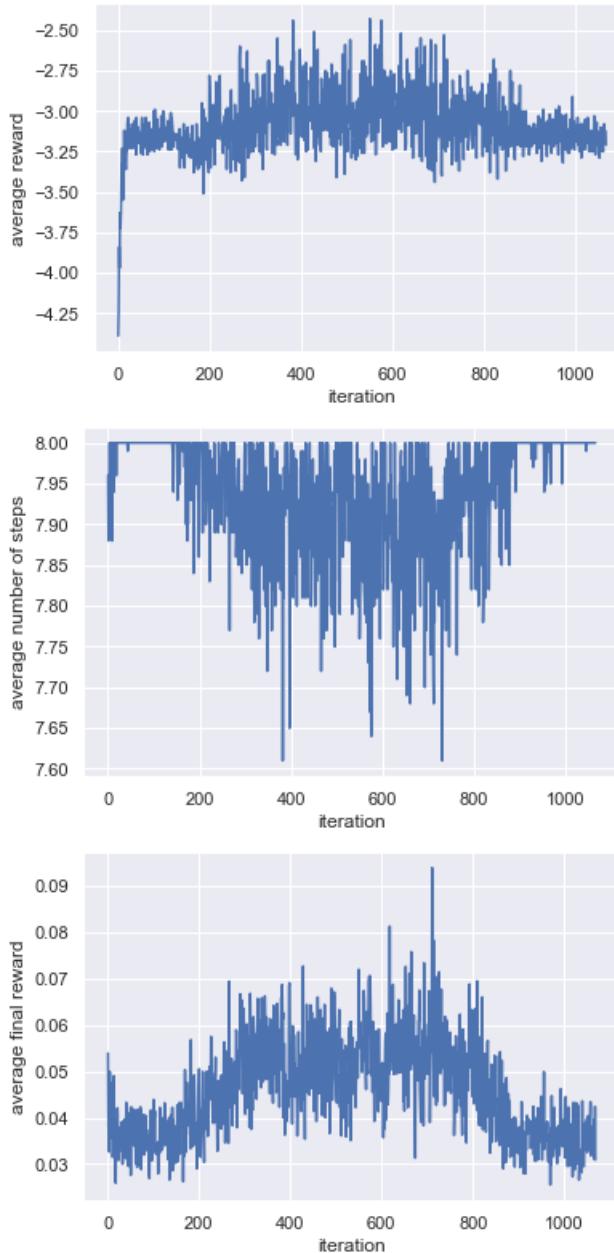
- Setting 1 (+1, -1, -0.8) without warm start



- Setting 1 (+1, -1, -0.7) without warm start



- Setting 1 (+1, -1, -0.5) without warm start



### 3. Result Analysis with new action spaces

Question:

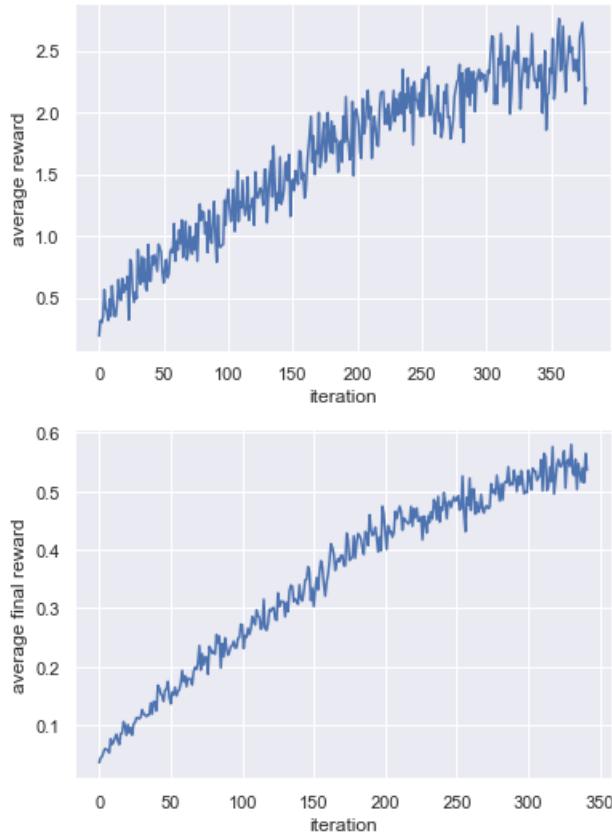
The reason why the training of model fails with penalty greater than -0.9 could be that, with current action space, if we select "no modification" action, our model cannot update the network for position prediction and amino acid prediction. With small penalty for "no modification" action, the training of our model can get stuck at the "no modification" action and don't update the other two networks.

Therefore, based on this assumption, I wonder whether combining the "no modification" action with the "position" action can help alleviate this issue. With this new action space,

$$R = \sum \gamma^t r_t$$

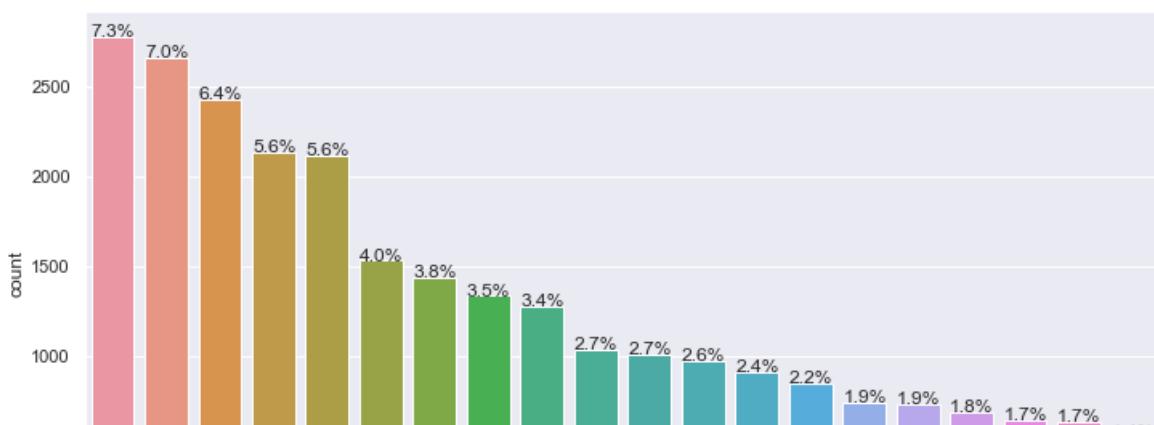
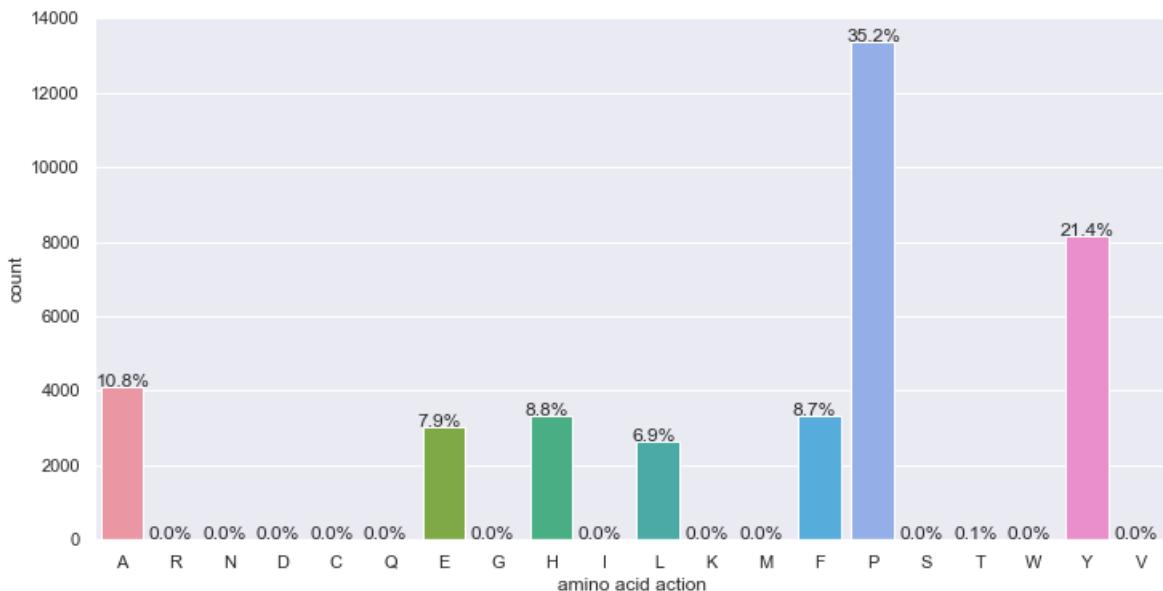
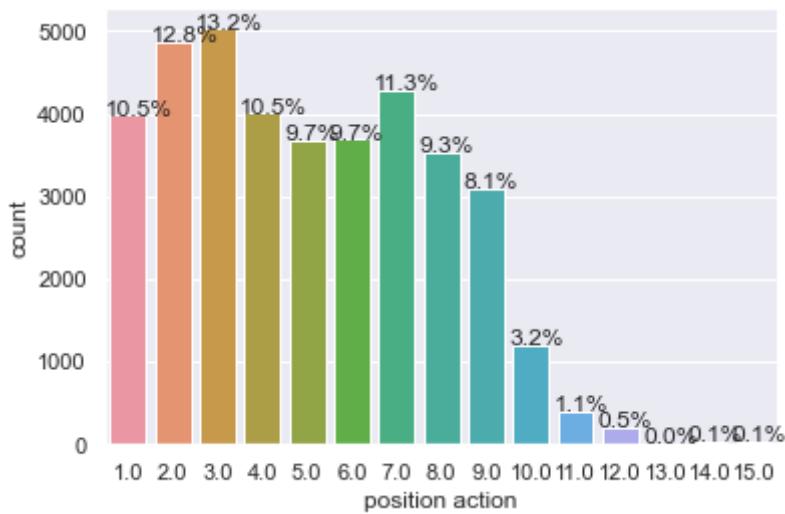
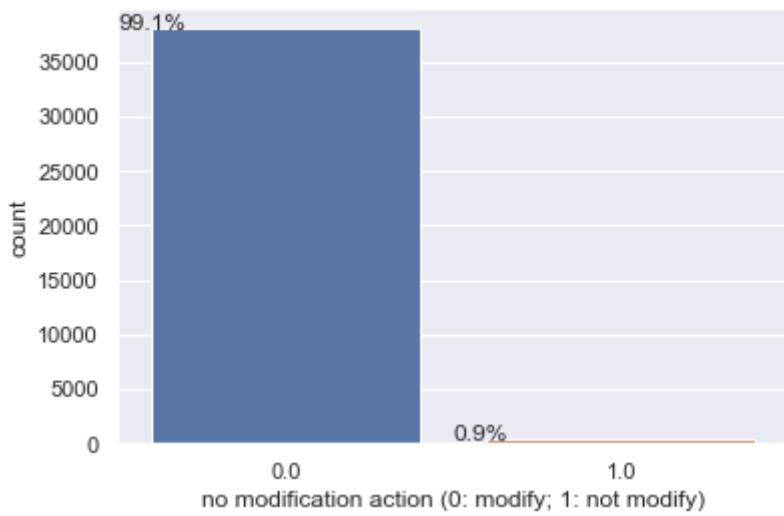
- Setting 2 ( $r_t = s_t * \alpha^{T-t} + \beta$ ) with two actions (position & no modification, amino acid).

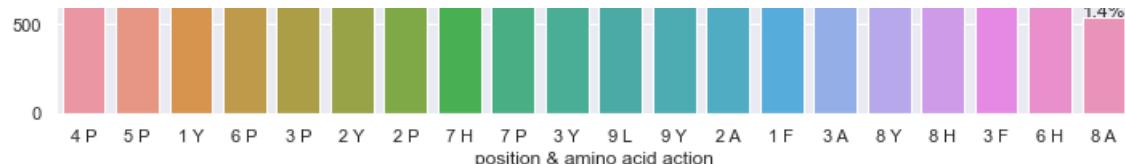
Among all the samples (i.e., the initial state of each episode/trajectory), 80% of them are sampled from the IEDB dataset; 20% from randomly generated peptides.



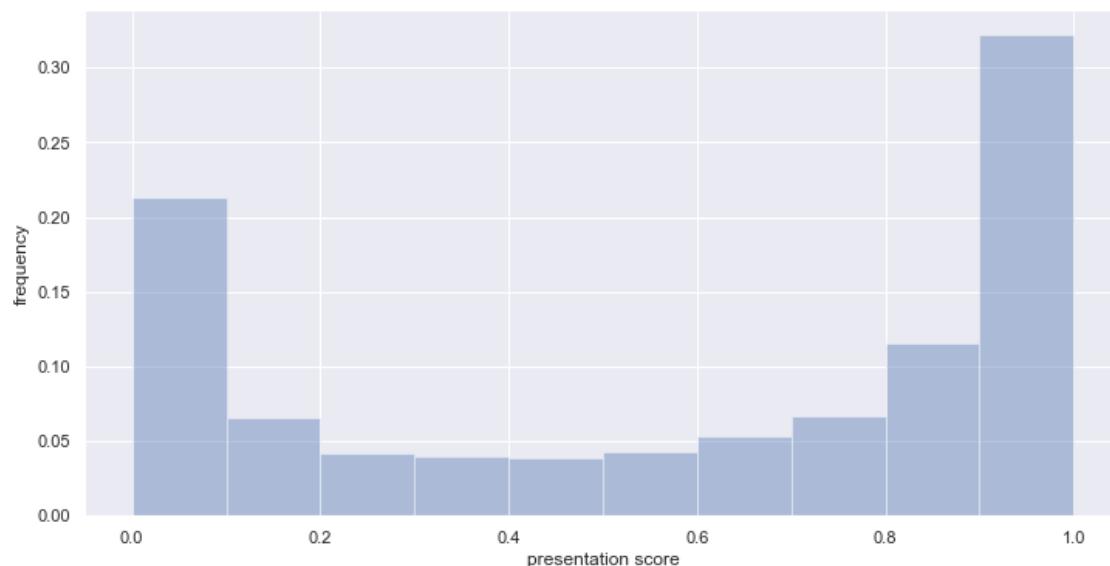
Each iteration has 3,840 steps. In each step, a specific given state will be updated by the policy network. Within 378 iterations, we have 181,440 samples in total. In order to analyze the trained policy network, we select the corresponding optimized peptides (i.e., the peptides in the terminal state) of the last 4,800 (about 2.6%) samples in the last 10 iterations for further analysis.

- Analysis 1: action analysis

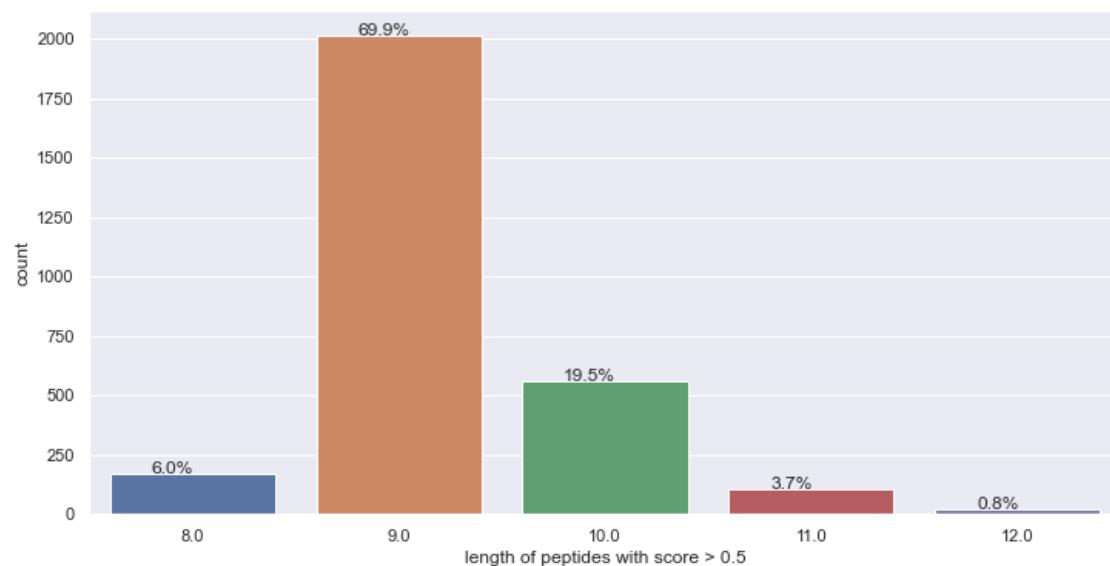


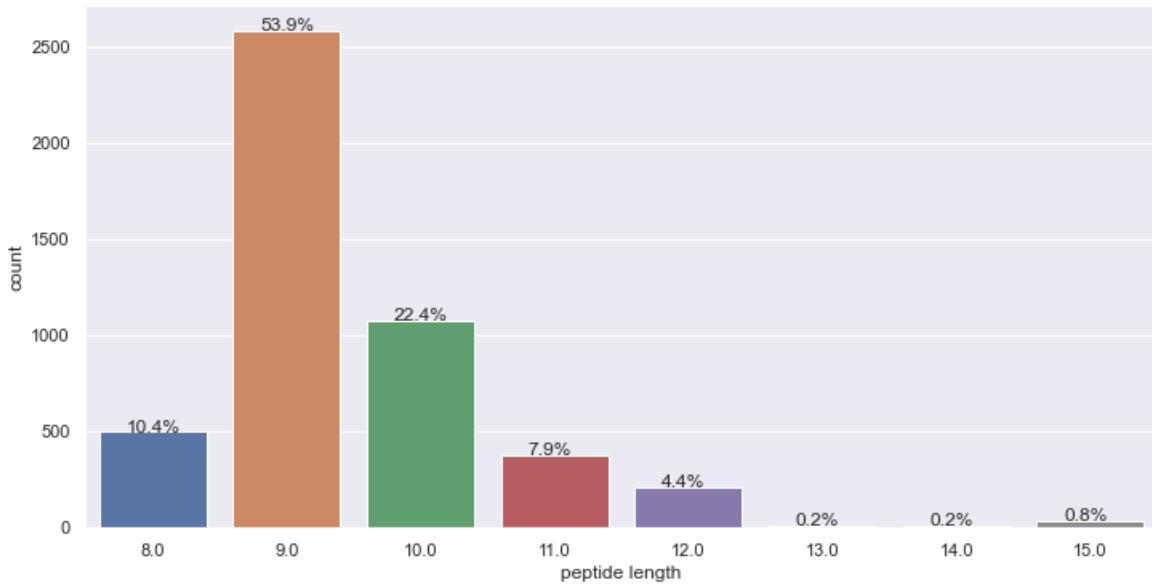


- Analysis 2: Reward distribution (60.08% optimized peptides have the presentation scores > 0.5)



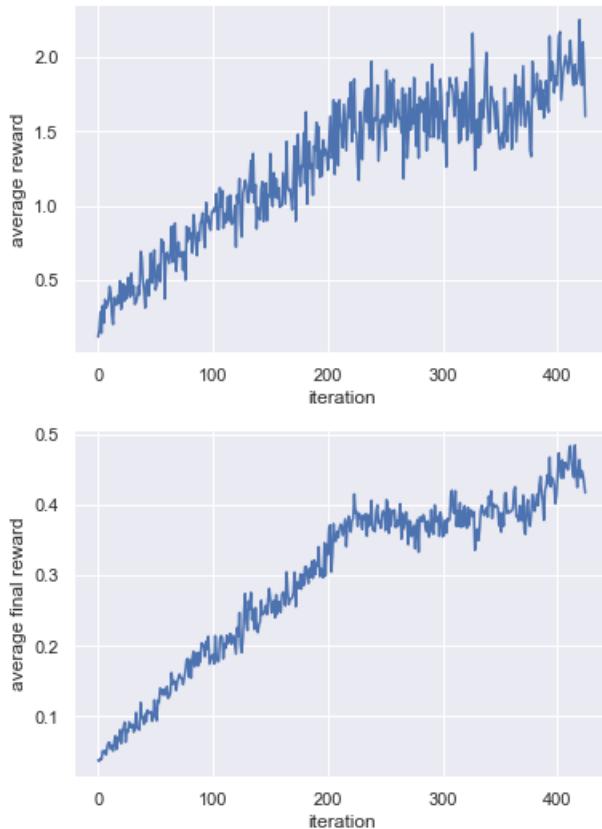
- Analysis 3: peptide length distribution





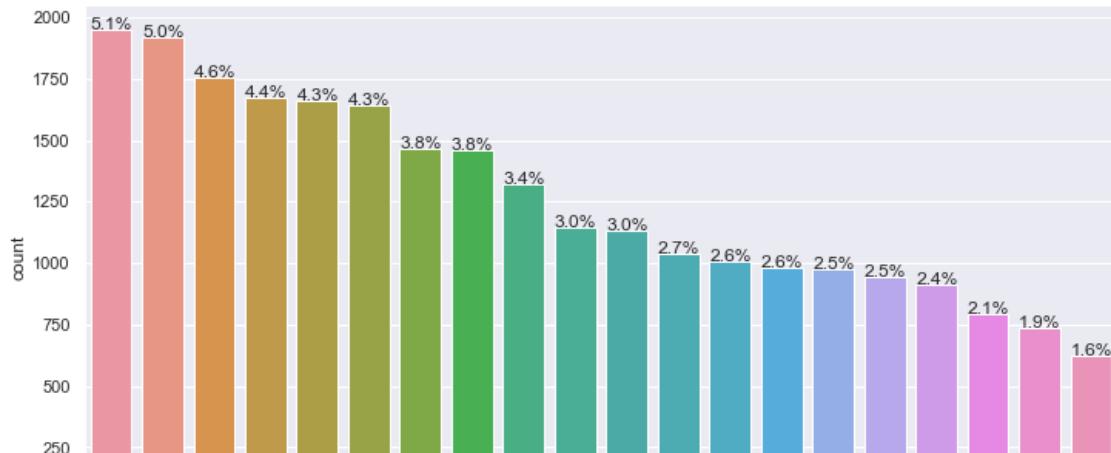
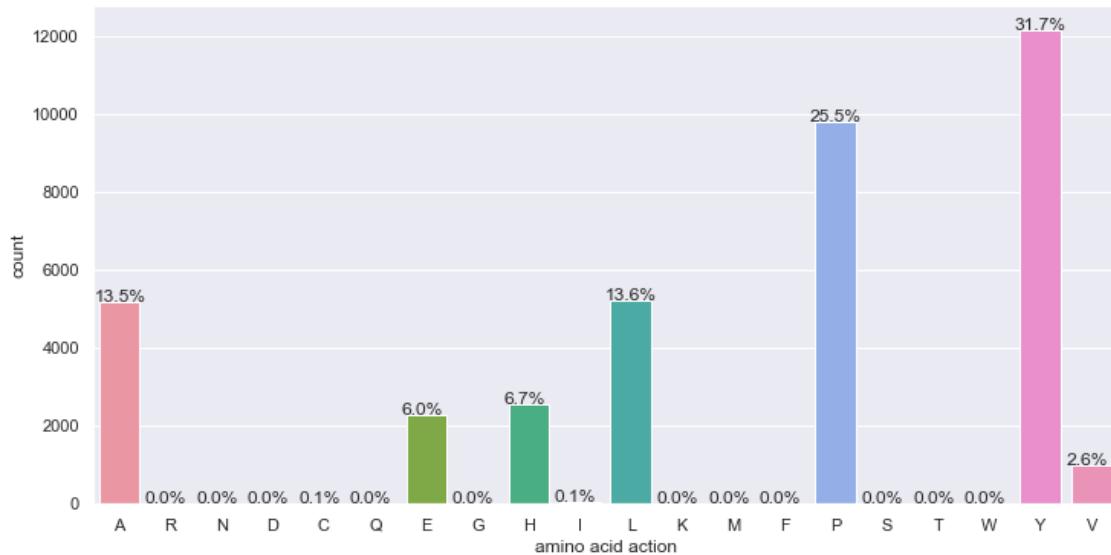
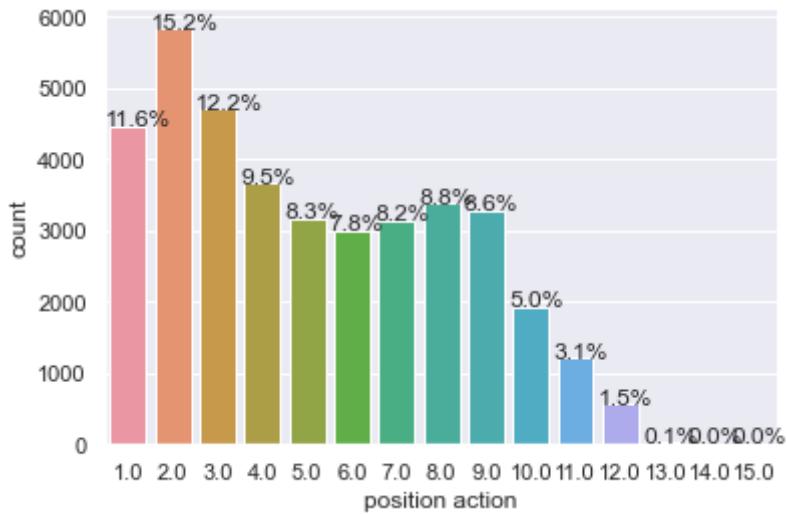
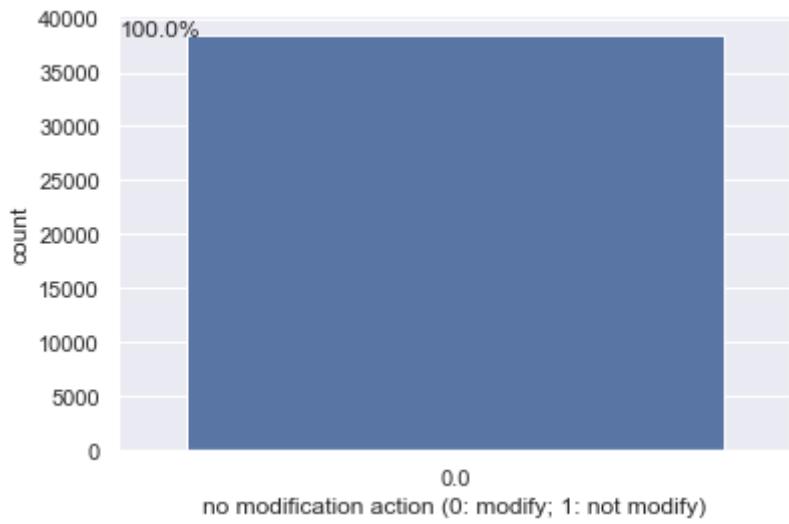
#### 4. Result Analysis with different sampling strategy.

- Setting 2 with 50% IEDB samples and 50% randomly generated peptides



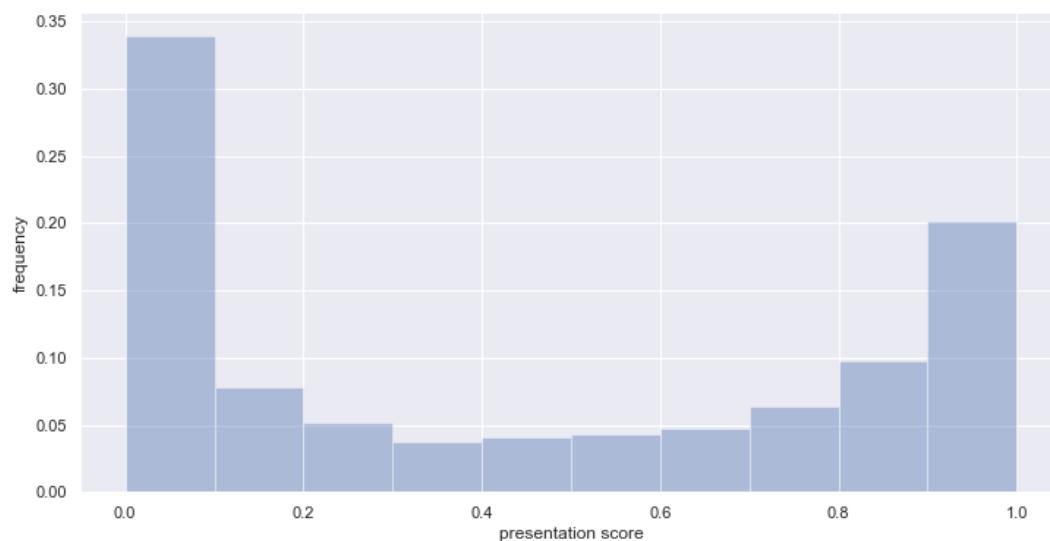
Each iteration has 3,840 steps. In each step, a specific given state will be updated by the policy network. Within 425 iterations, we have 204,000 samples in total. In order to analyze the trained policy network, we select the corresponding optimized peptides (i.e., the peptides in the terminal state) of the last 4,800 (about 2.4%) samples in the last 10 iterations for further analysis.

- Analysis 1: action analysis

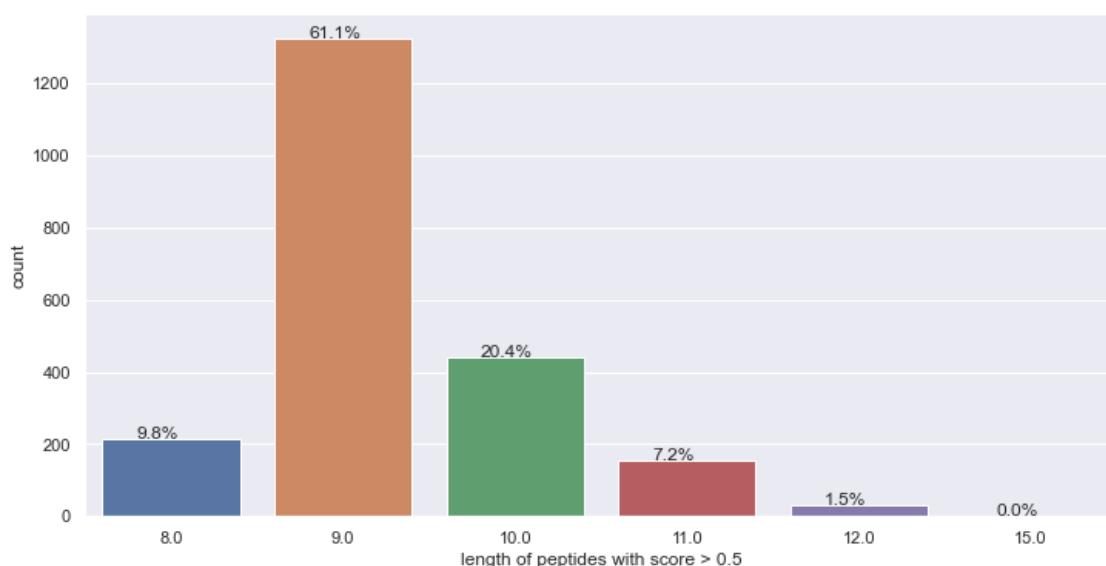
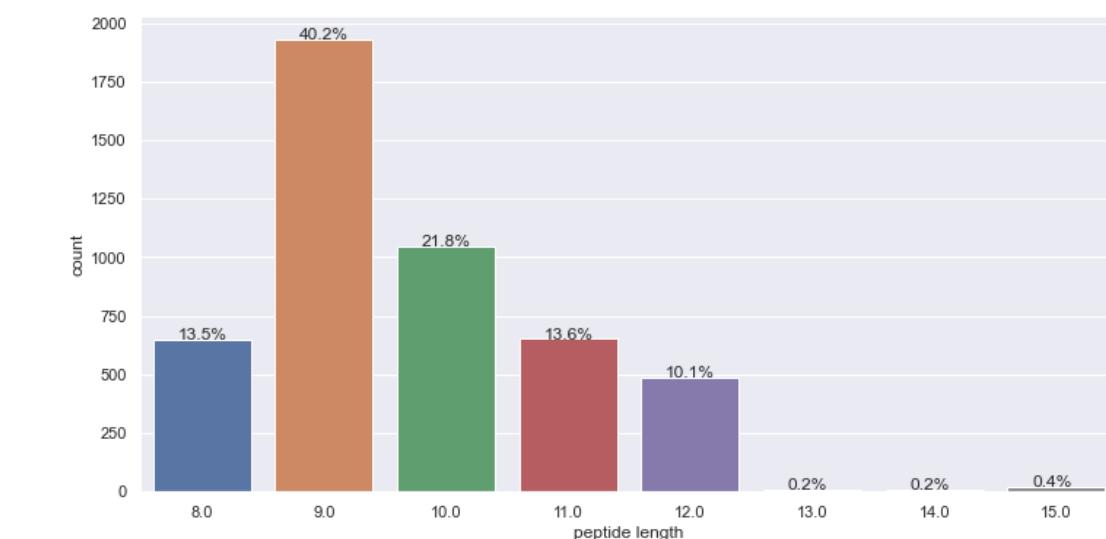




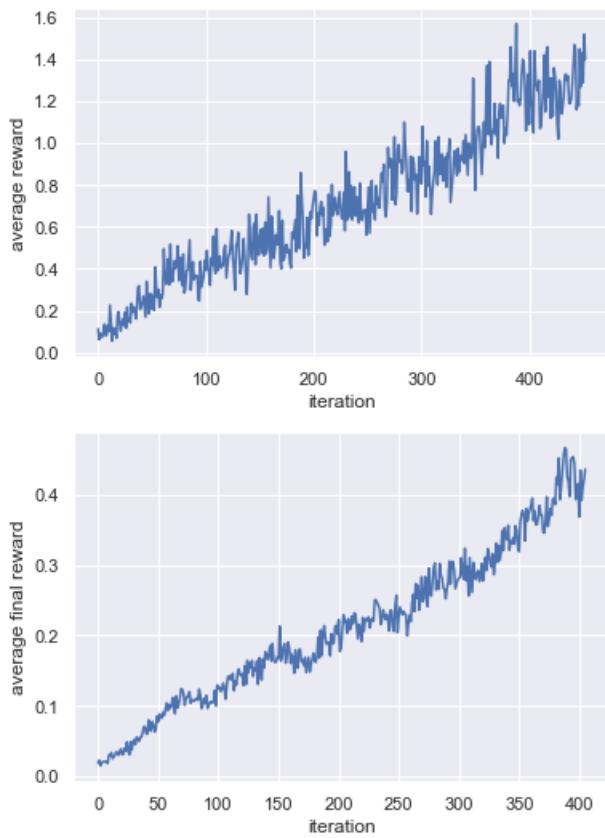
- Analysis 2: Reward distribution (45.27% optimized peptides have the presentation scores > 0.5)



- Analysis 3: peptide length distribution



- Setting 2 with all randomly generated peptides



## Update (5-19)

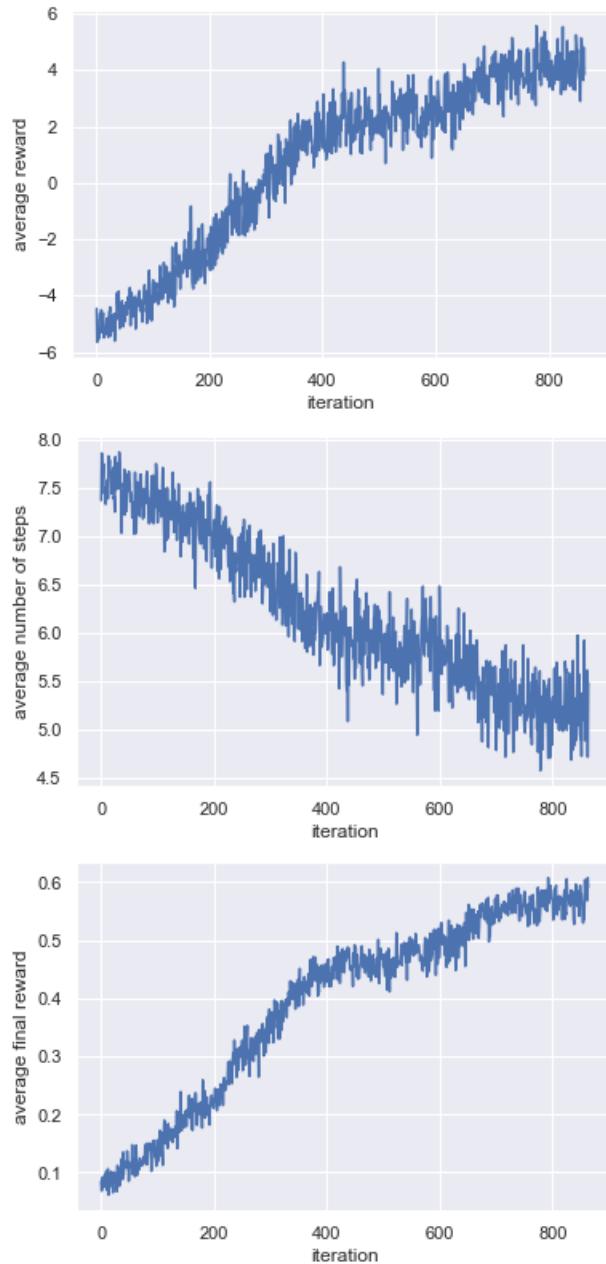
Suggestions:

- 1) Using a large entropy coefficient at the beginning to encourage exploration and then decrease the entropy coefficient. To avoid the learned policies being conserved at certain amino acids, it can be helpful to improve the exploration using a large coefficient.
- 2) Using prioritized experience replay (i.e., importance sampling) to increase the exploitation of good episodes (i.e., episodes with good presentation scores and limited modification).
- 3) For constrained generation (i.e., generating positive peptides from Covid-19 peptides), comparing the performance of the model with prioritized experience replay and the model with less maximum steps (change maximum number of steps to 3 or 4).
- 4) Improving the performance of models under the setting with only final rewards, through warming up (i.e., using the intermediate rewards at the beginning of training and then removing the intermediate rewards) or using the prioritized replay. In addition, trying to change the "no modification" action to "termination" action so that shorter episodes with large representation scores can receive higher cumulative rewards than longer episodes. Note that "termination" action can lead to less exploration and thus worse performance. Therefore, we need to limit the usage of "termination" action at the beginning of training.
- 5)

## 1. Result Analysis with longer running time

- Setting 1 (+1, -1, -0.9) (stop criteria = 0.6; sample rate = 80% from IEDB and 20% from random)

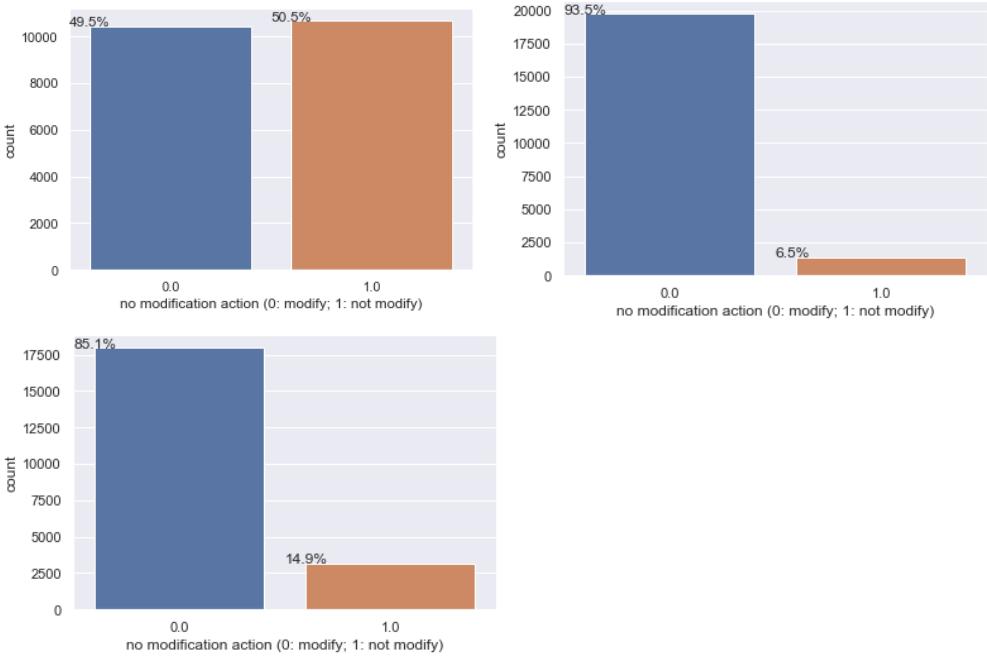
It takes 72 hours to get the results below....



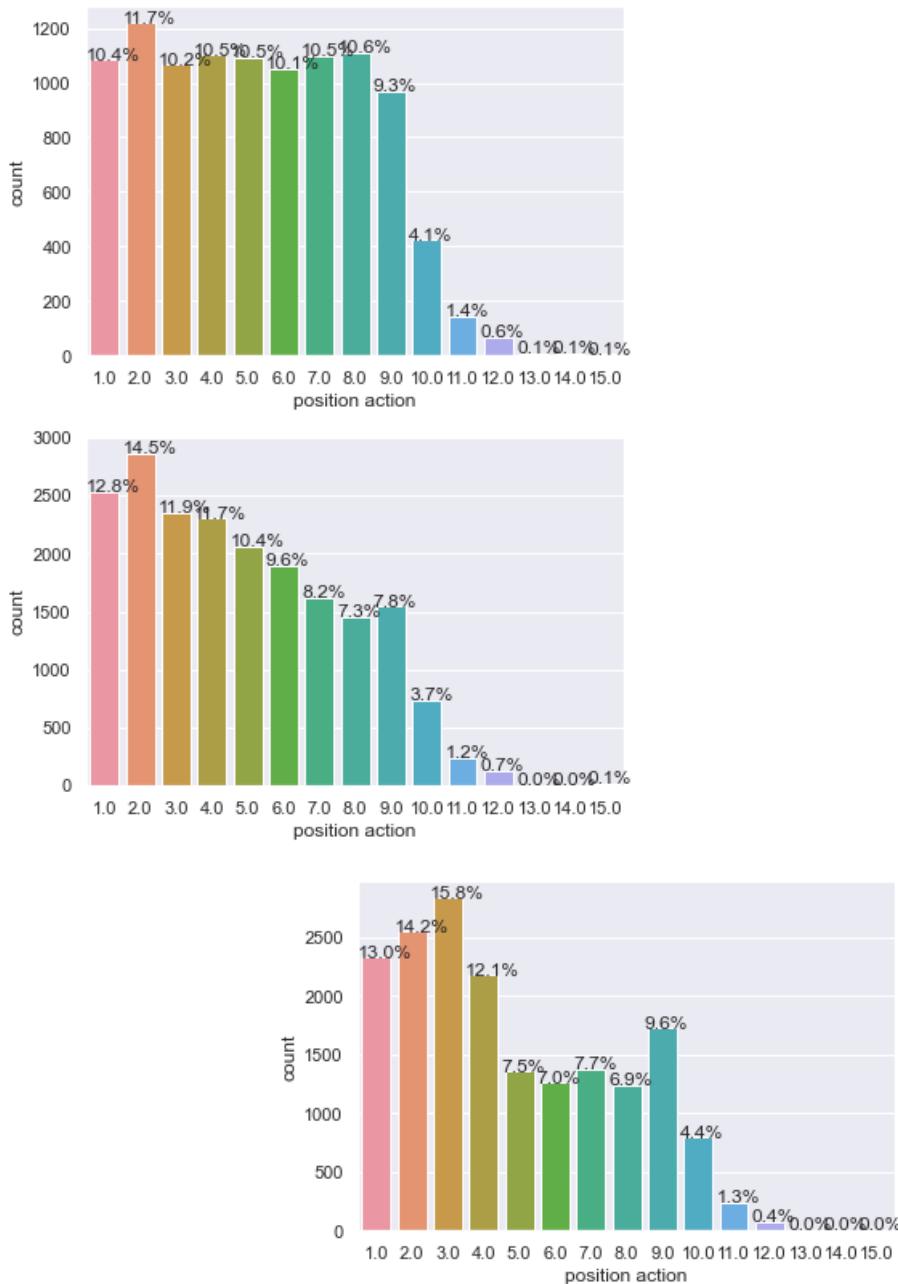
### Action analysis

#### 1. no modification action

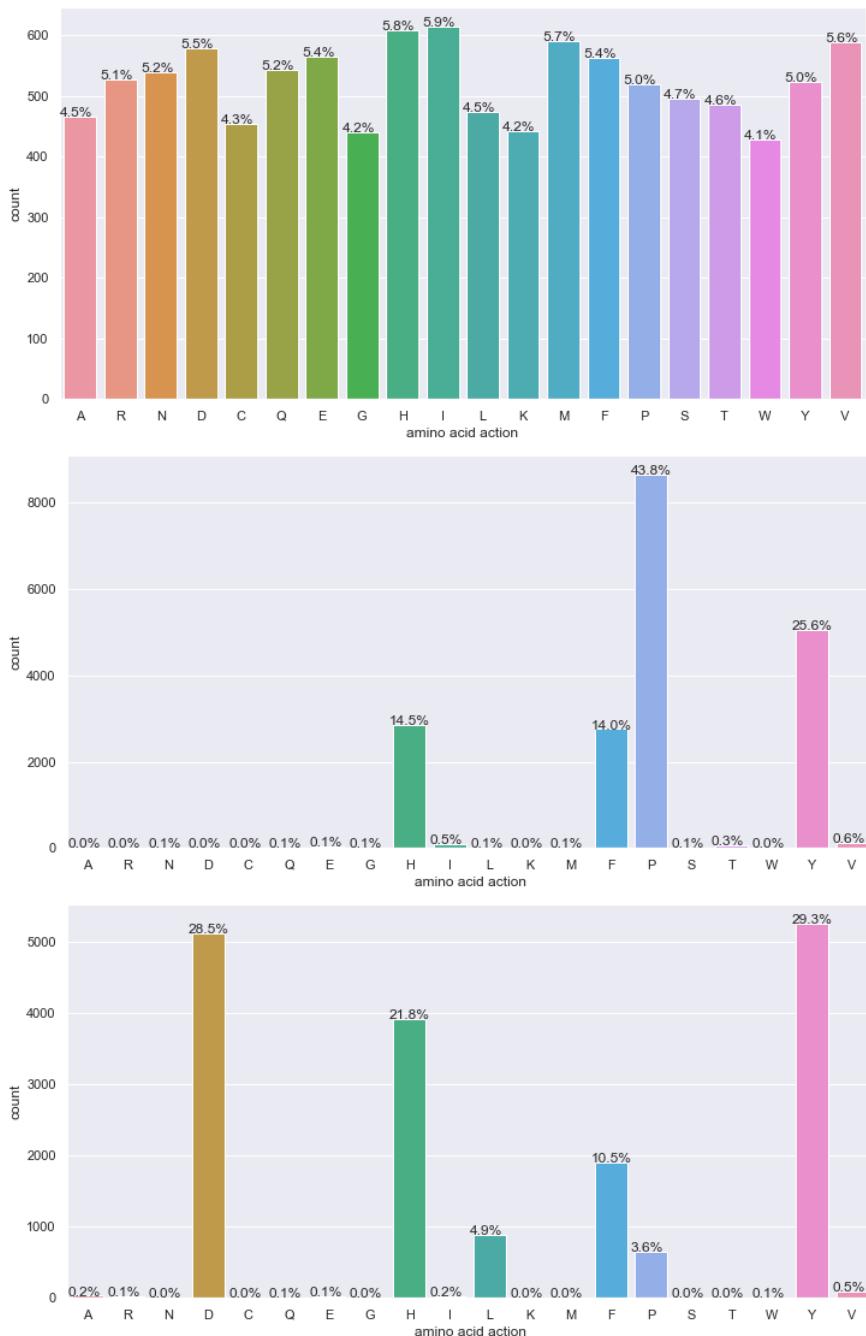
(first 10 iterations) (mid 10 iterations) (last 10 iterations)



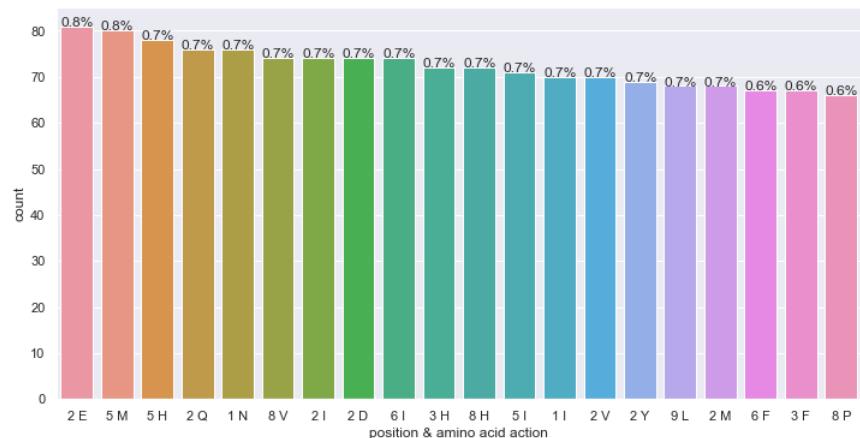
## 2. position select action

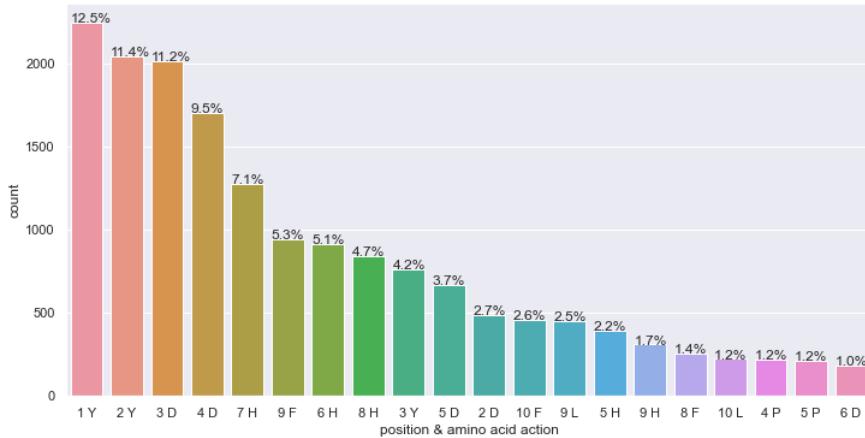
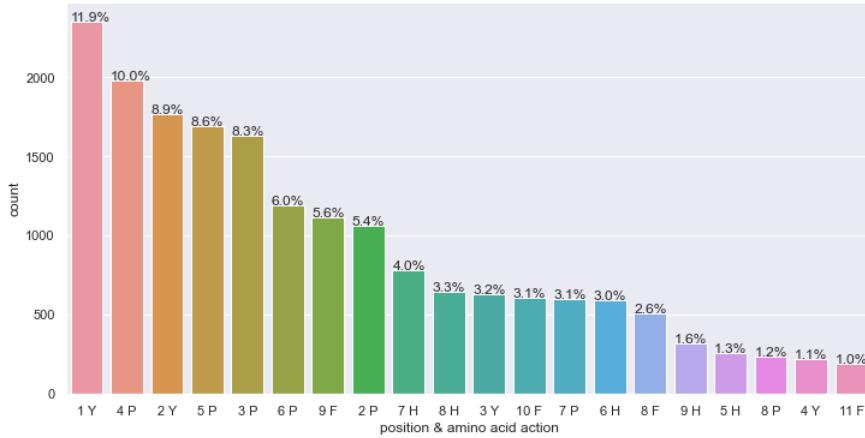


### 3. amino acid action

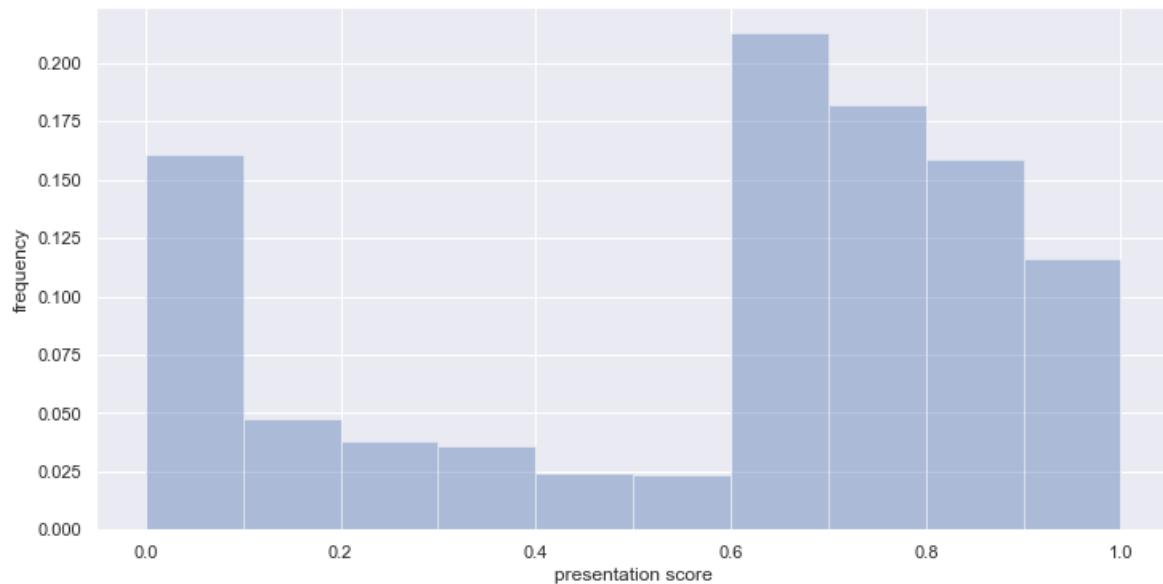


### 4. pos & amino action

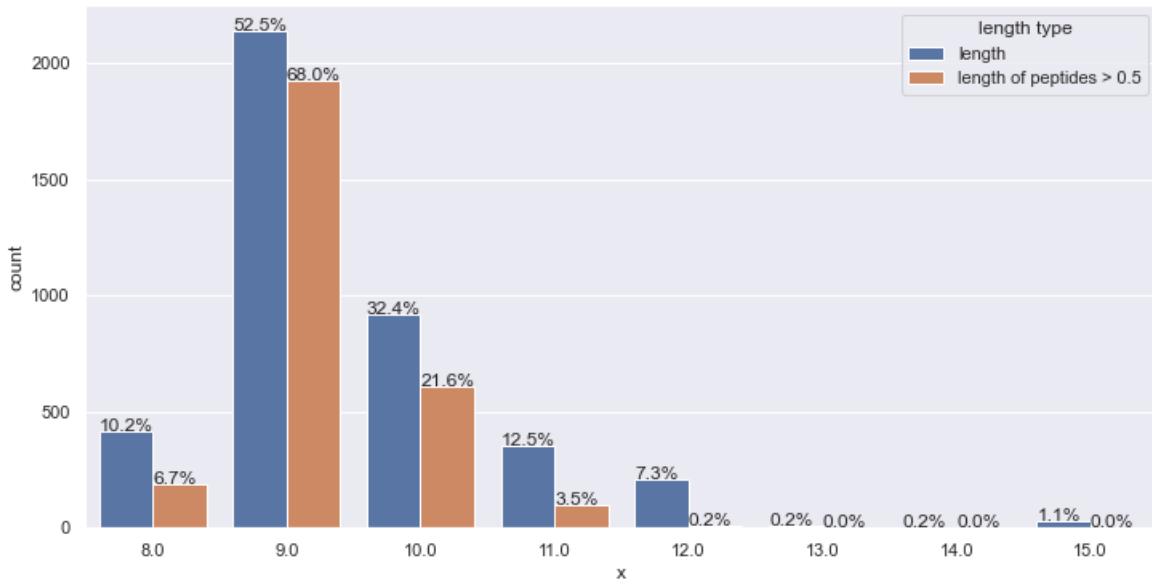




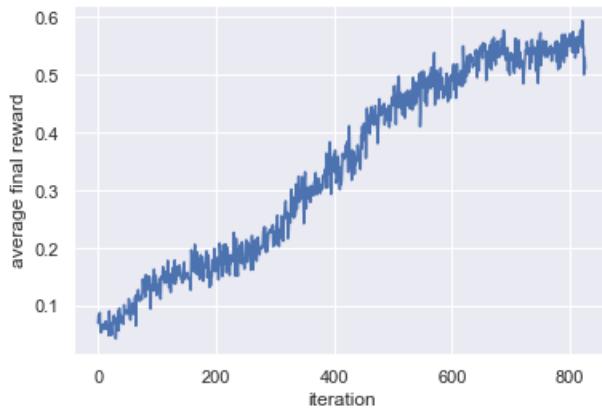
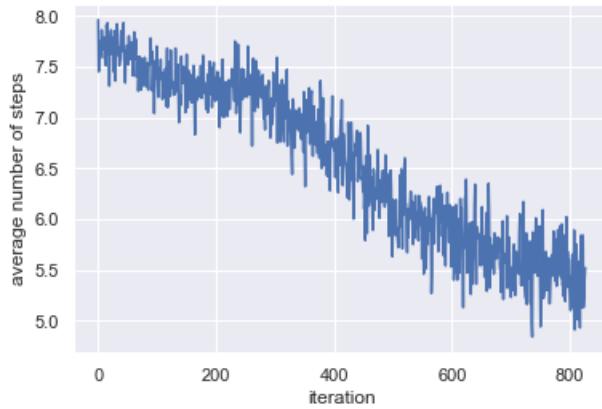
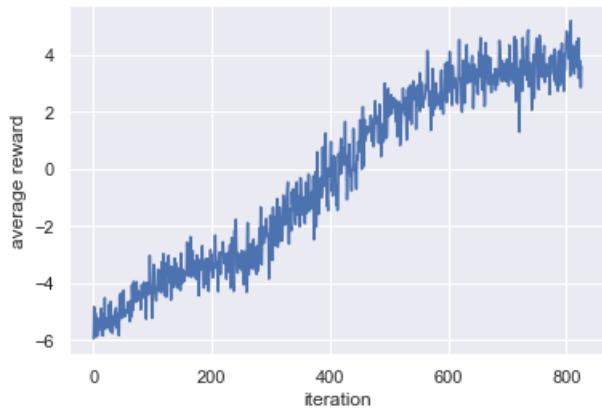
Presentation score analysis: (69.4% optimized peptides have presentation scores > 0.5)



Peptide length distribution

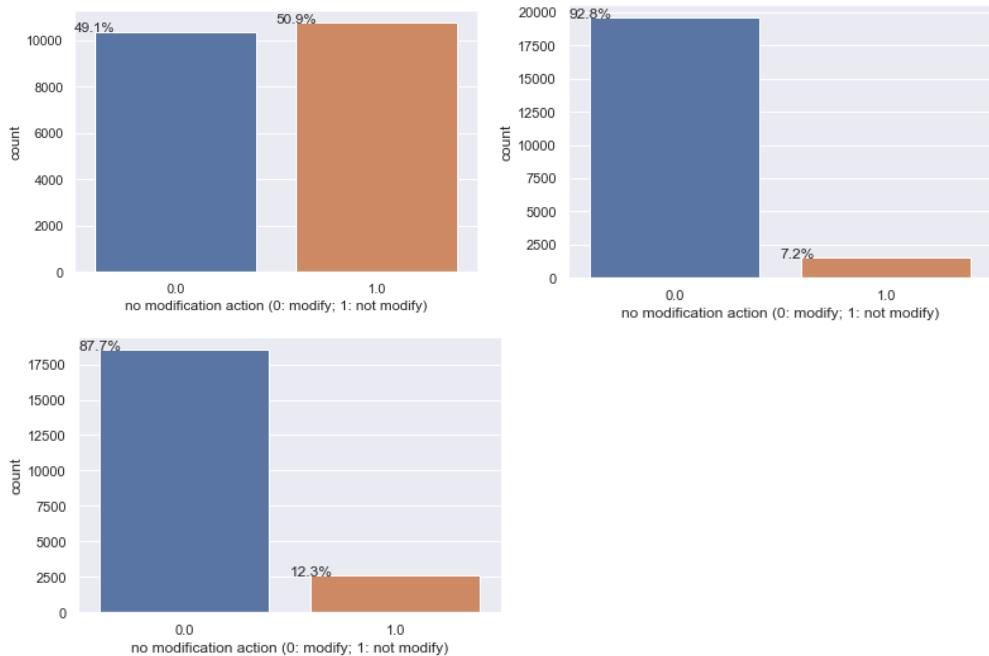


- Setting 1 (+1, -1, -0.9) (stop criteria = 0.6; sample rate = 50% from IEDB and 50% from random)

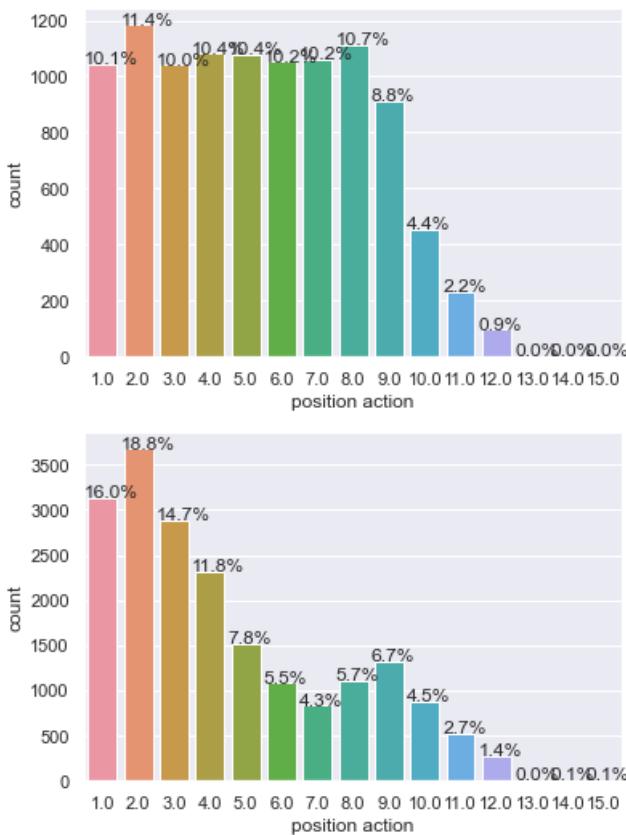


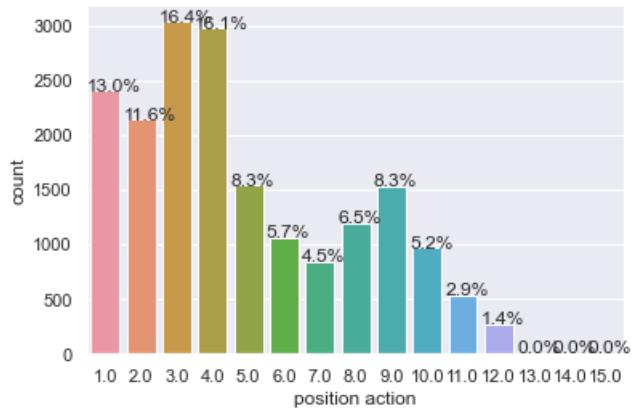
1. no modification action

(first 10 iterations)  
 iterations)

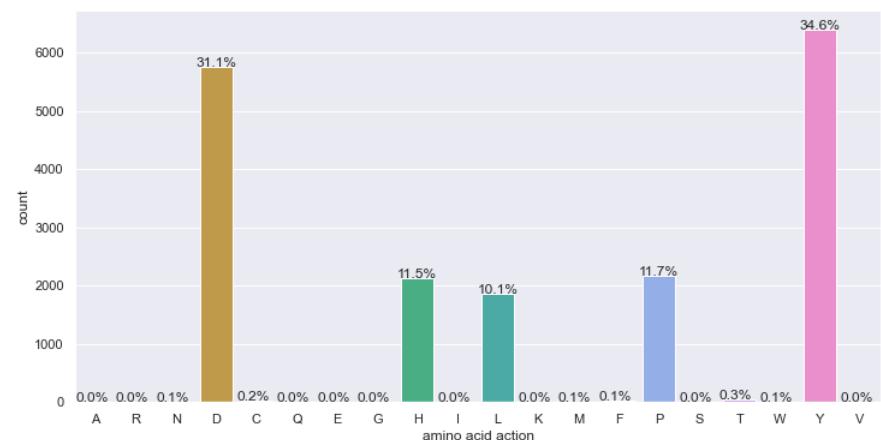
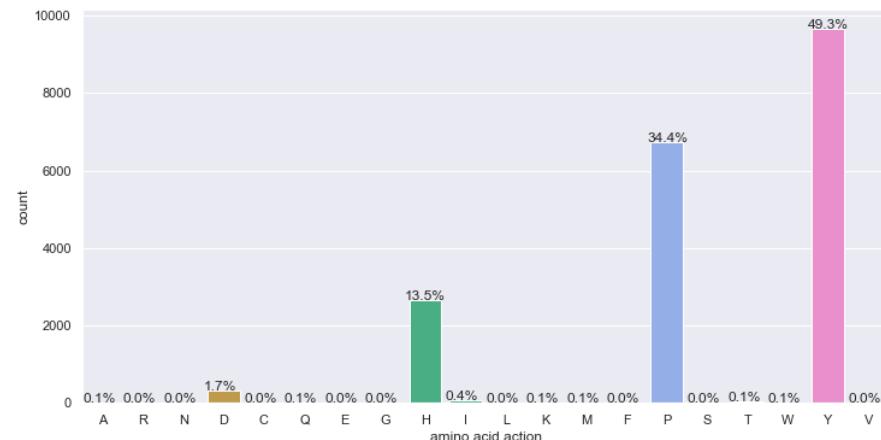
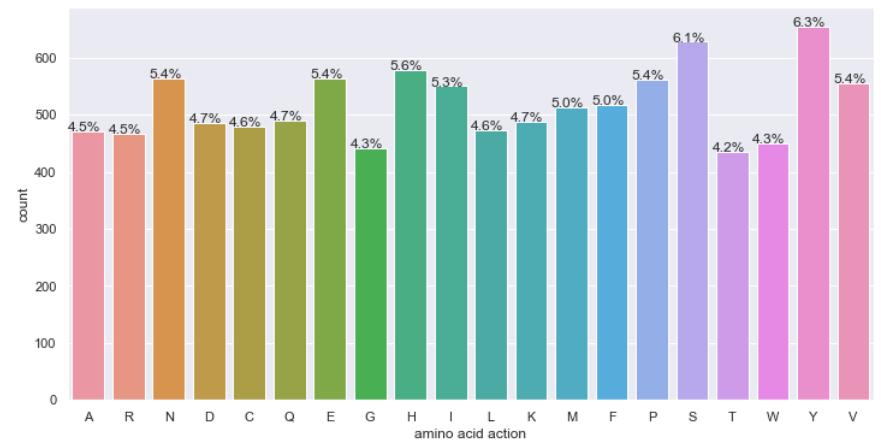


## 2. position select action

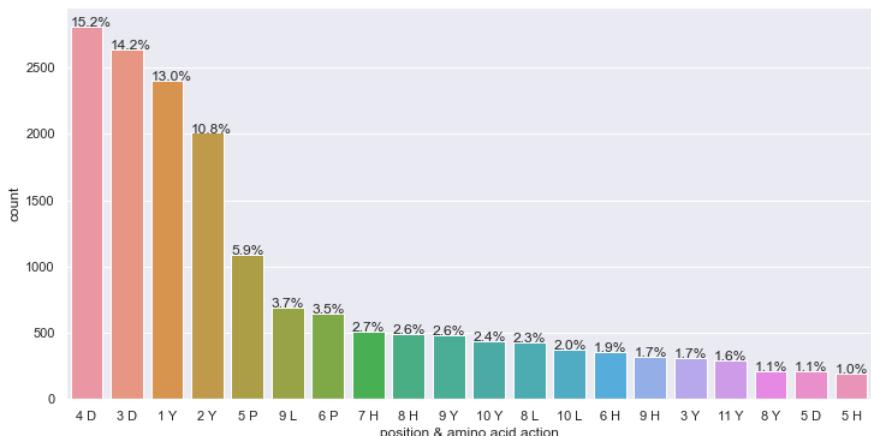
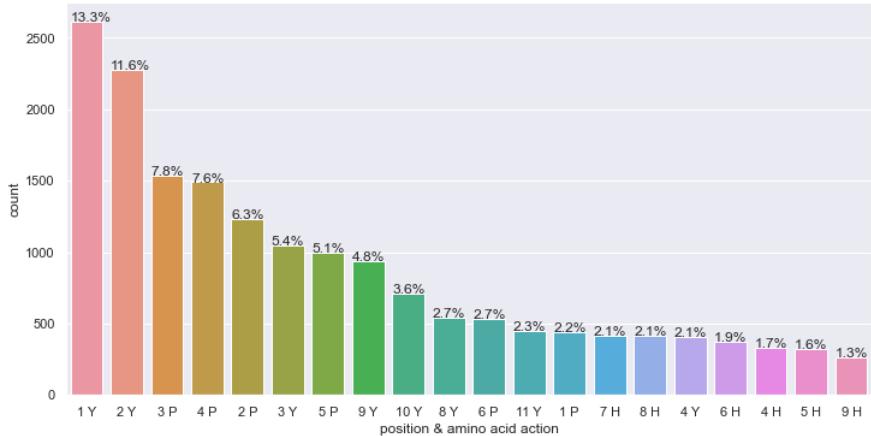
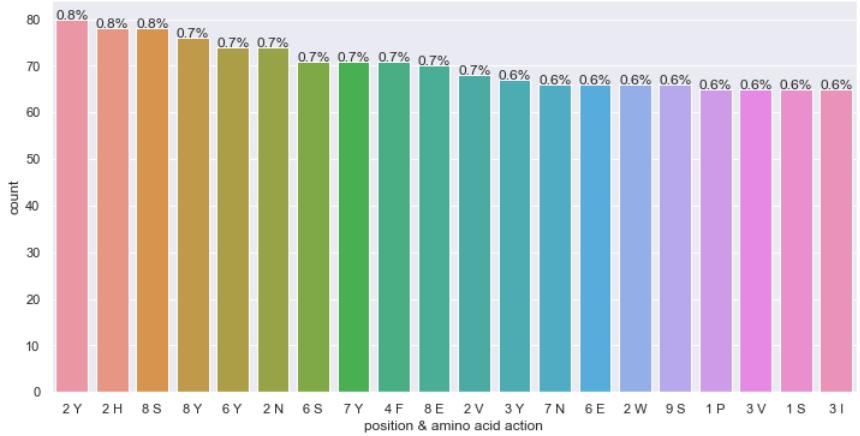




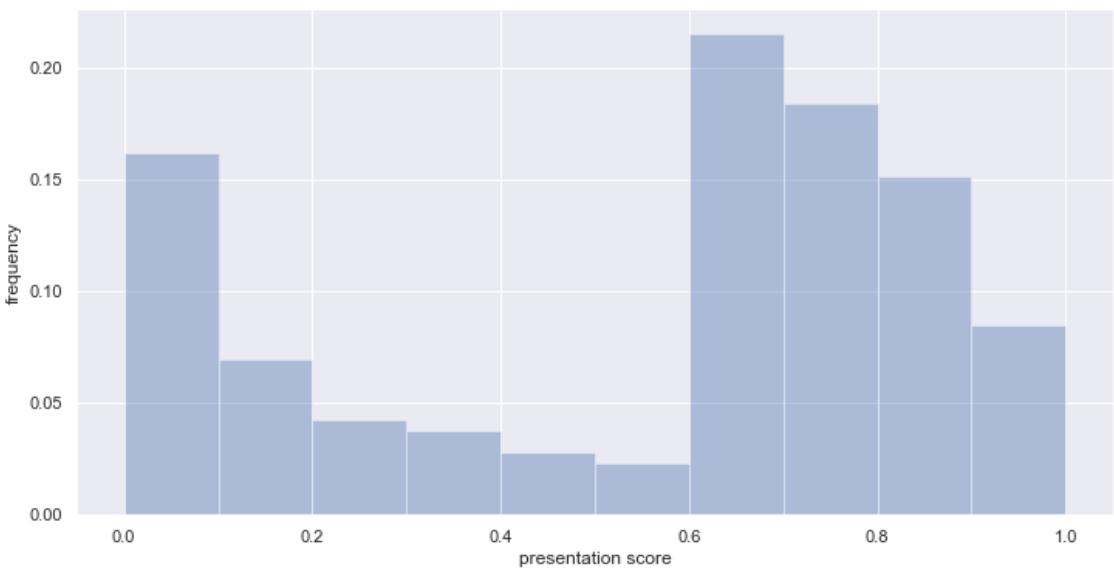
### 3. amino acid action

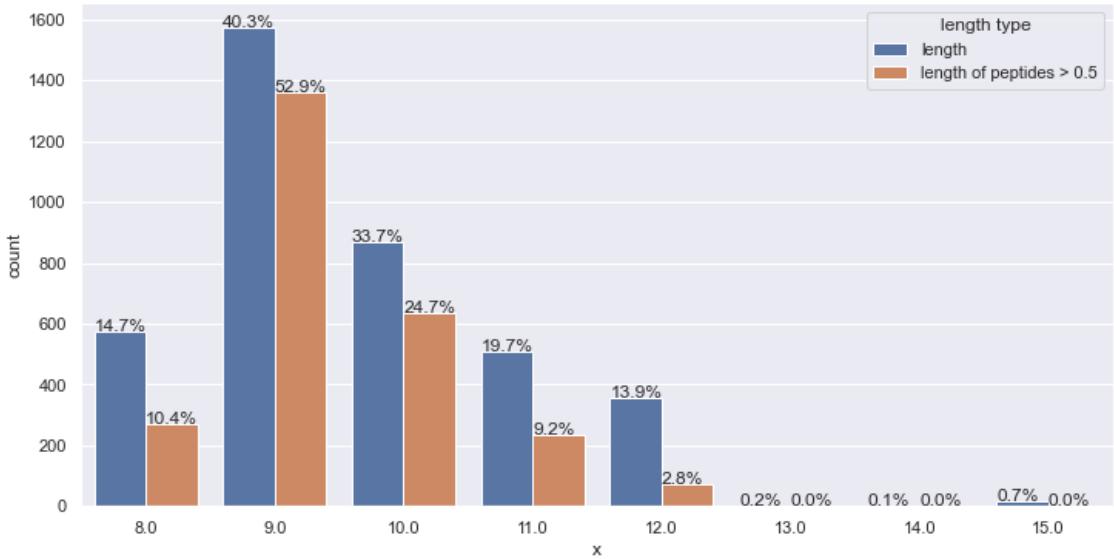


### 4. position & amino action

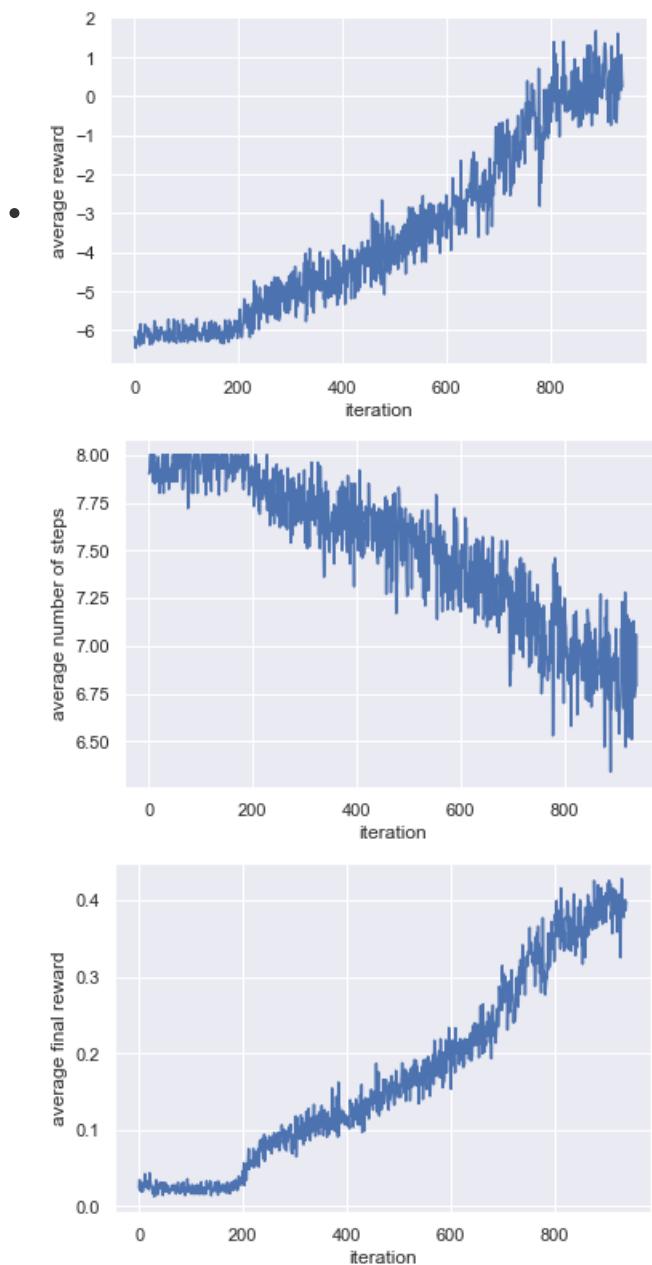


presentation score analysis: (65.9% optimized peptides have presentation scores > 0.5)





- Setting 1 (+1, -1, -0.9) (stop criteria = 0.6; sample rate = 100% from random)

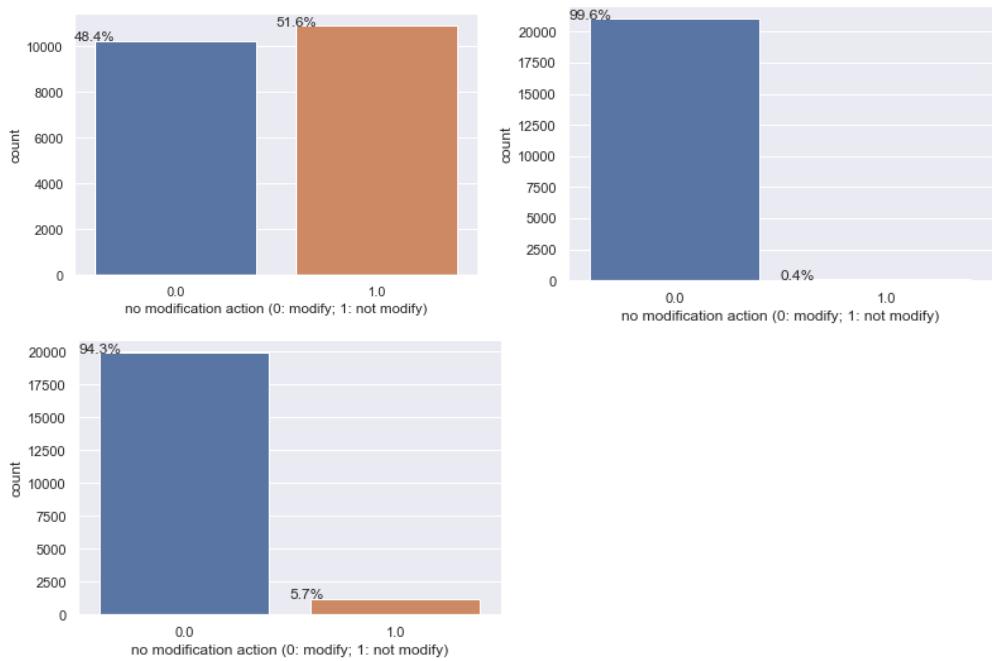


1. no modification action

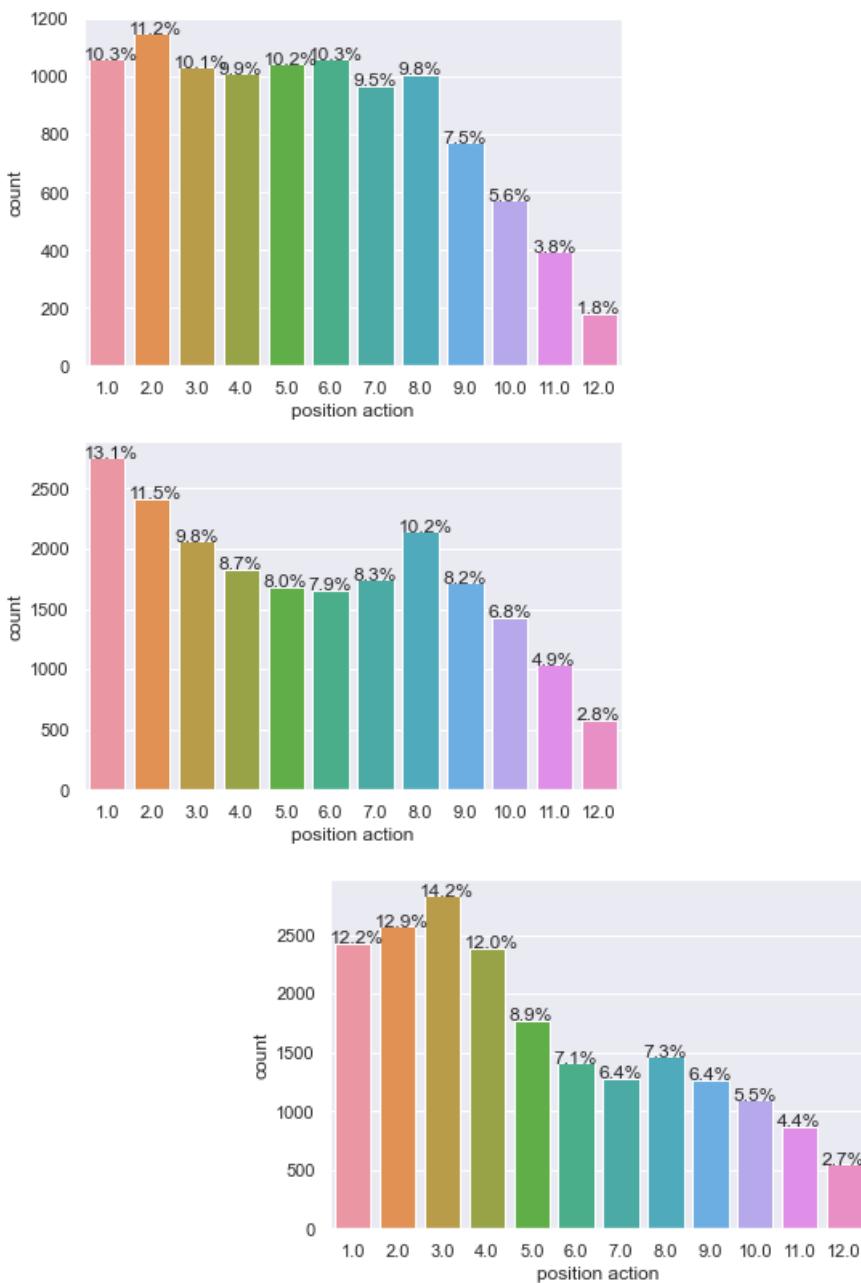
(first 10 iterations)  
iterations)

(mid 10 iterations)

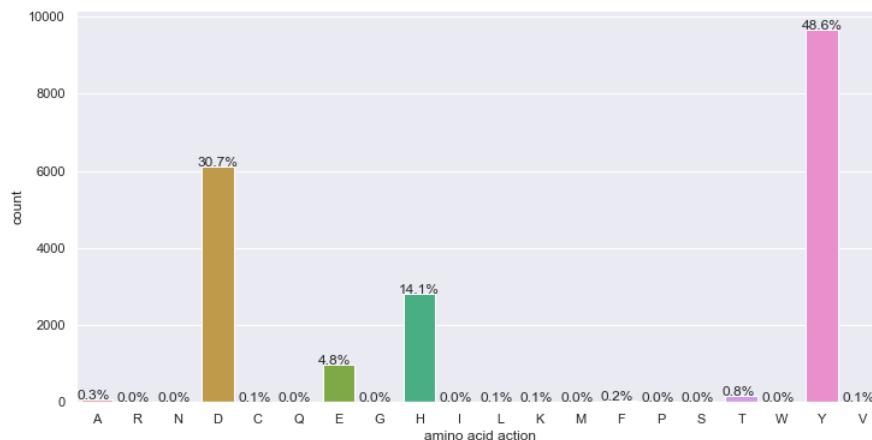
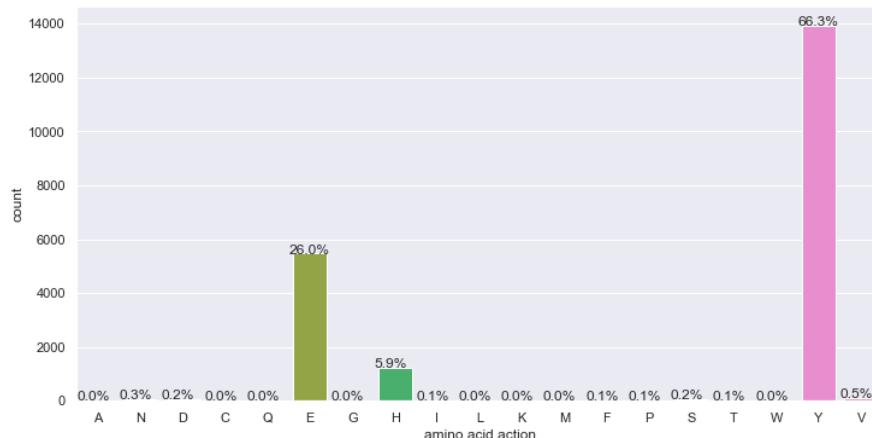
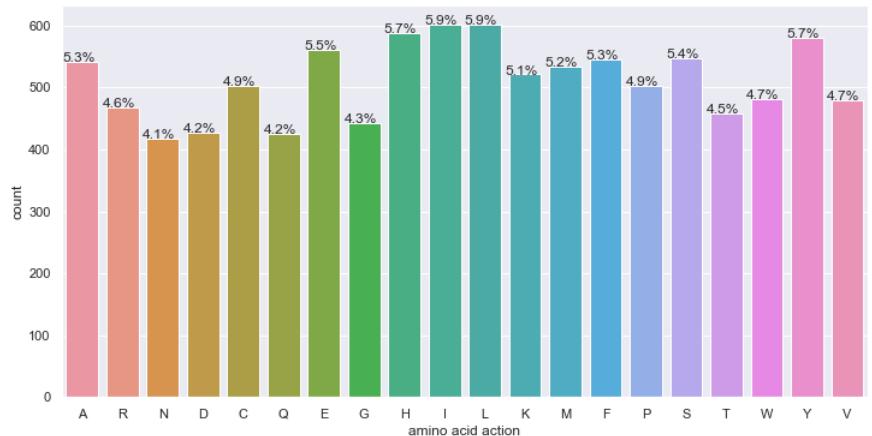
(last 10



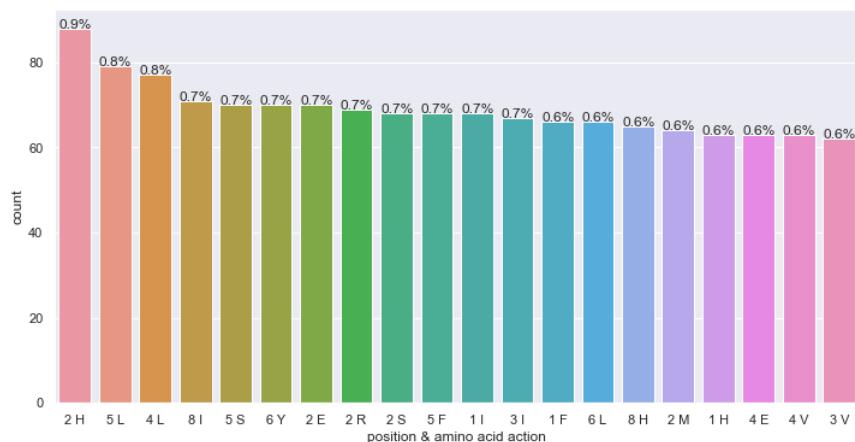
## 2. position select action

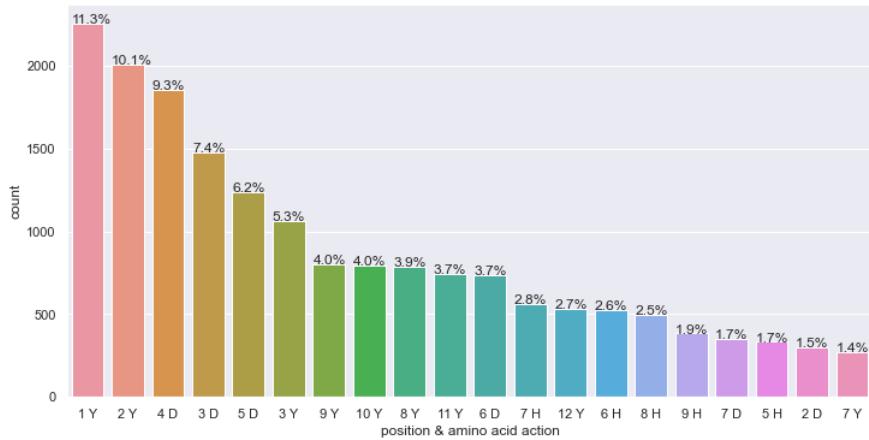
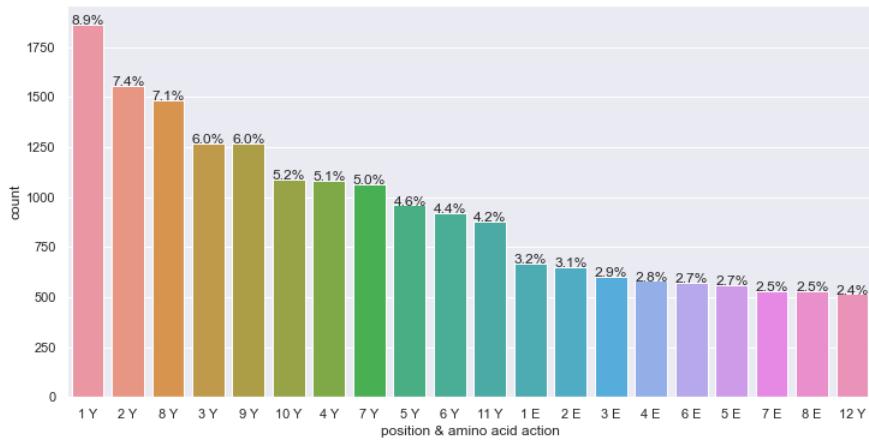


## 3. amino acid action

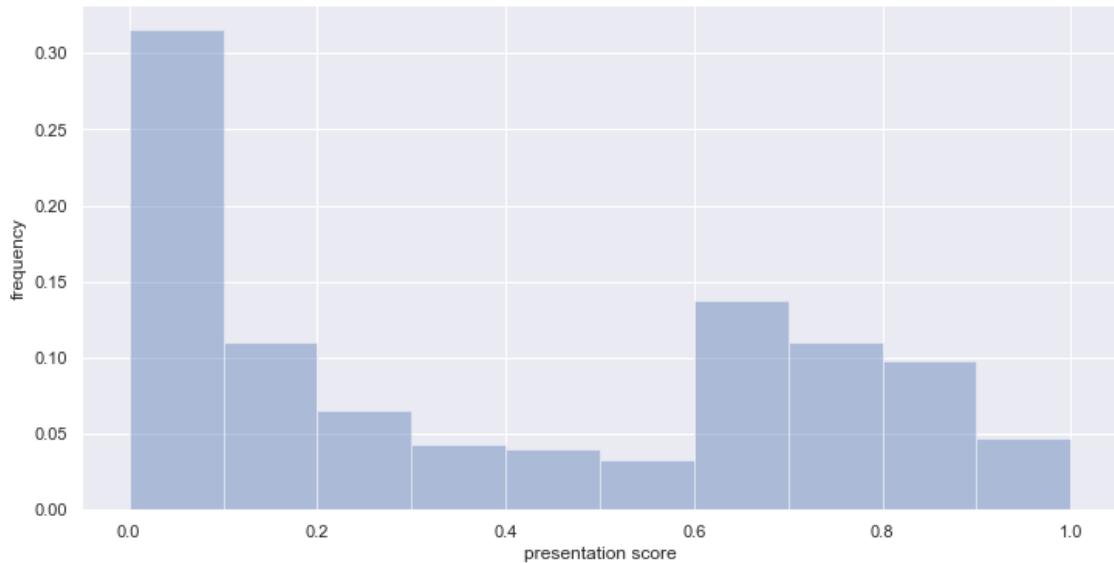


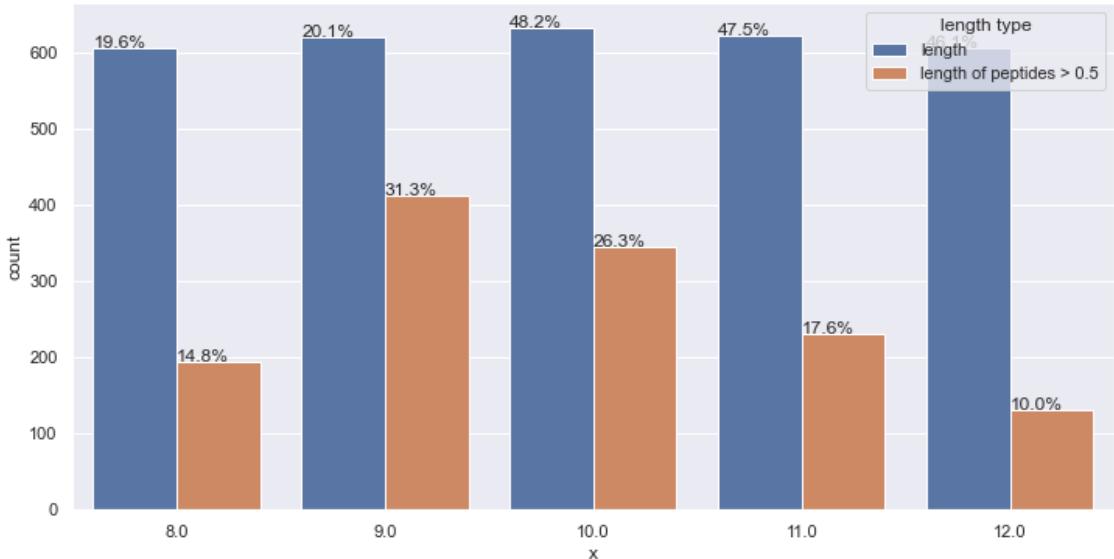
#### 4. pos & amino action





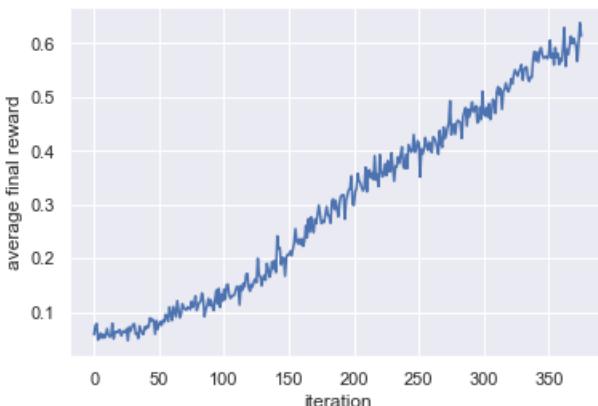
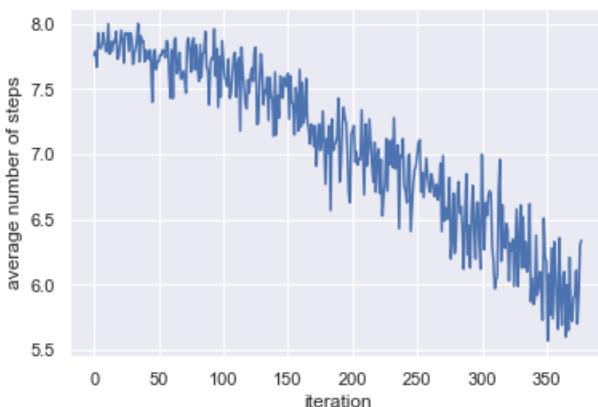
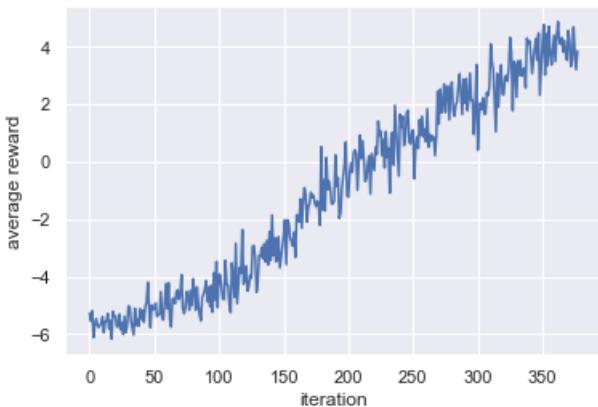
presentation score analysis: (42.5% optimized peptides have presentation scores > 0.5)



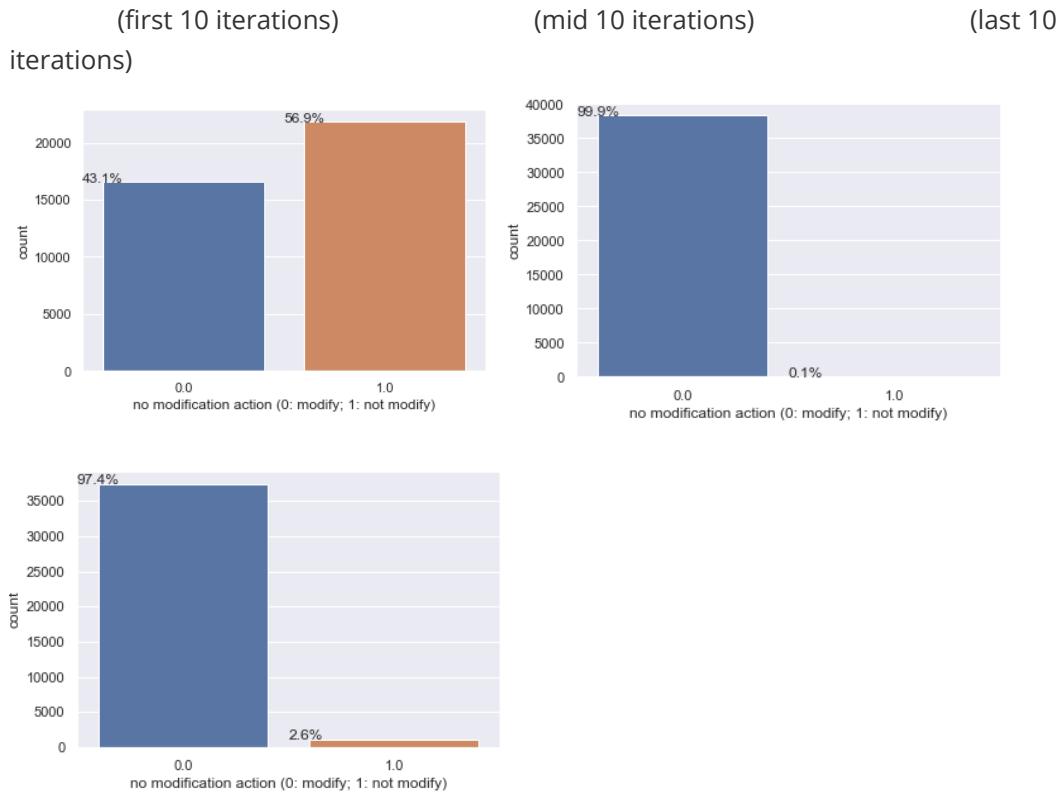


## 2. Result Analysis with different stopping criteria

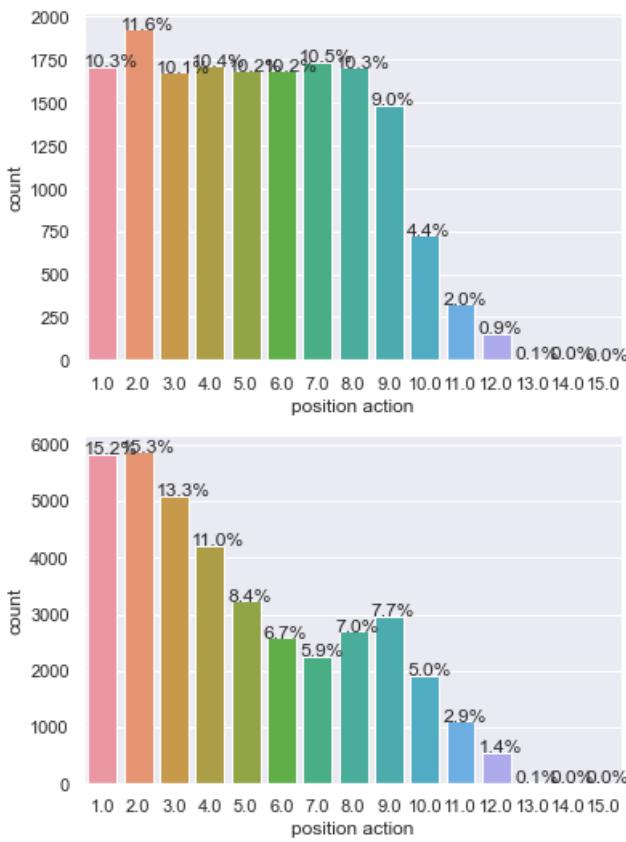
- Setting 1 (+1, -1, -0.9) (stop criteria = 0.75; sample rate = 50% from IEDB and 50% from random)

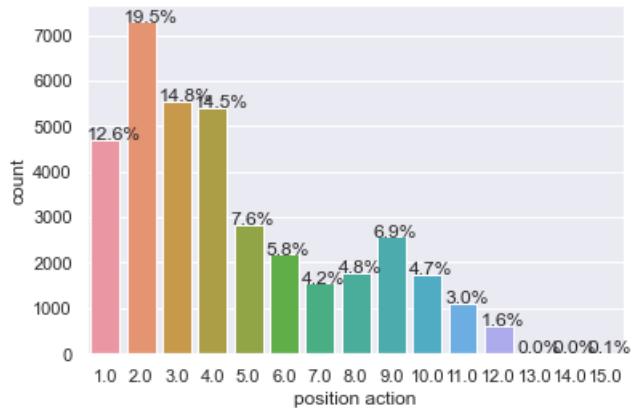


## 1. no modification action

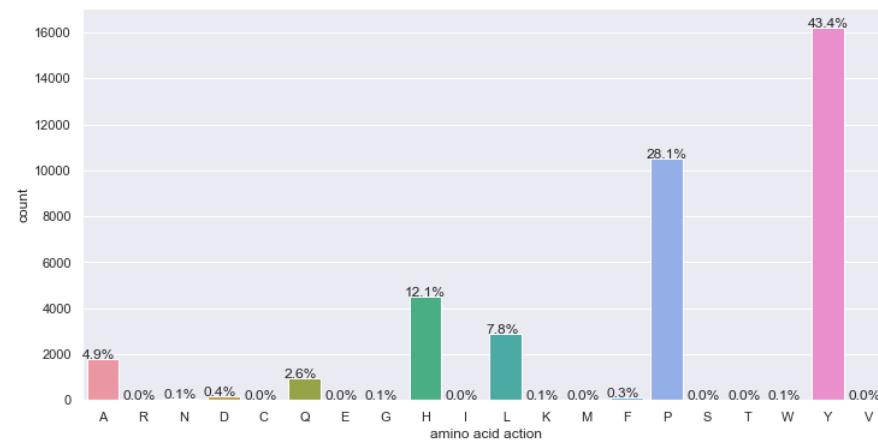
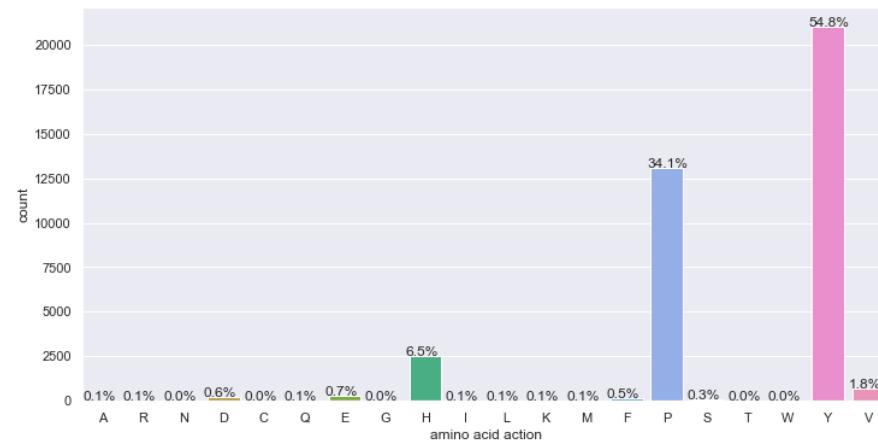
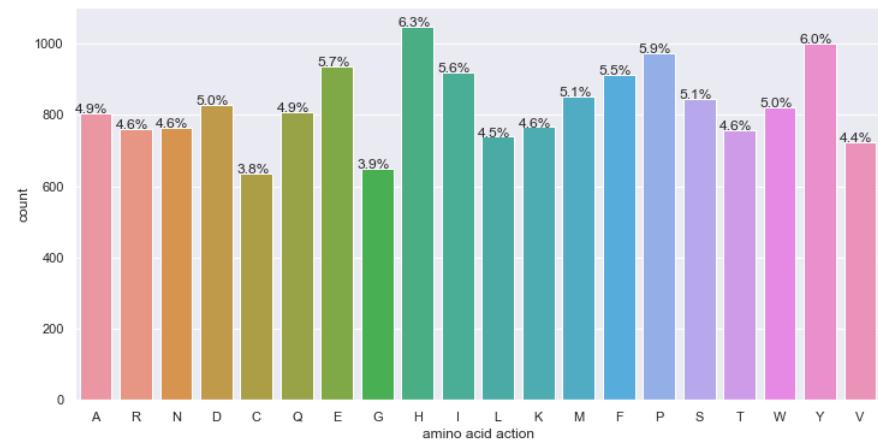


## 2. position select action

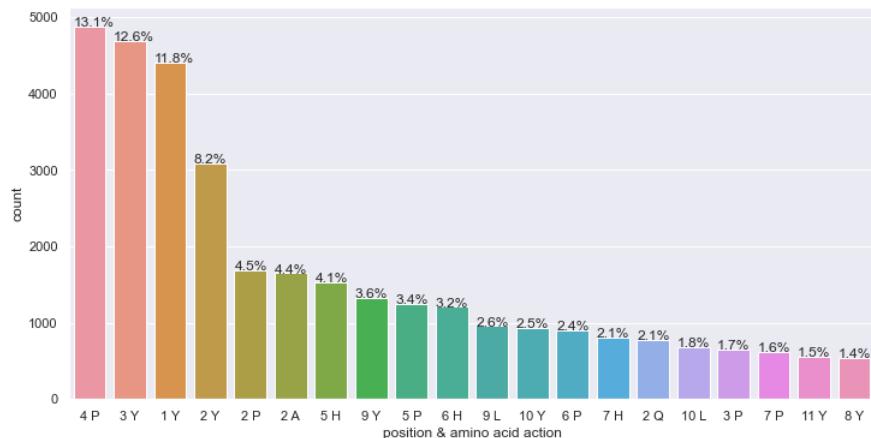
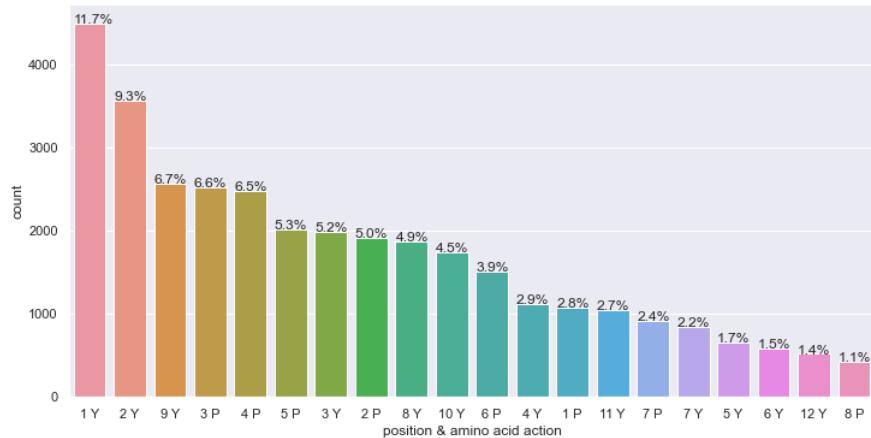
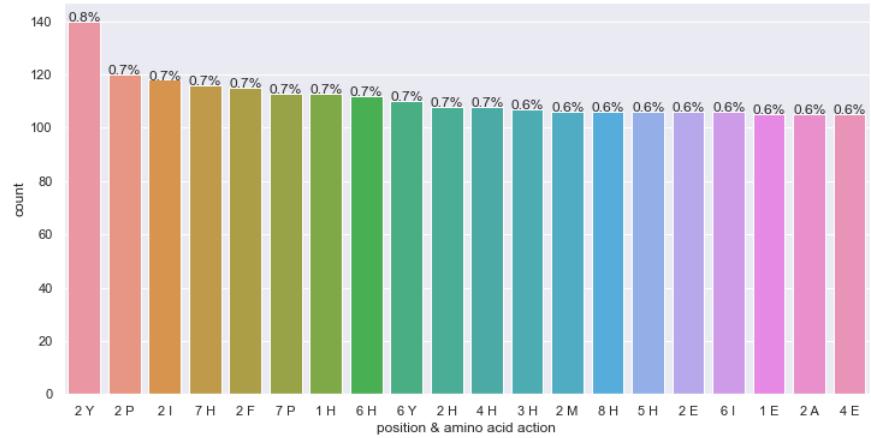




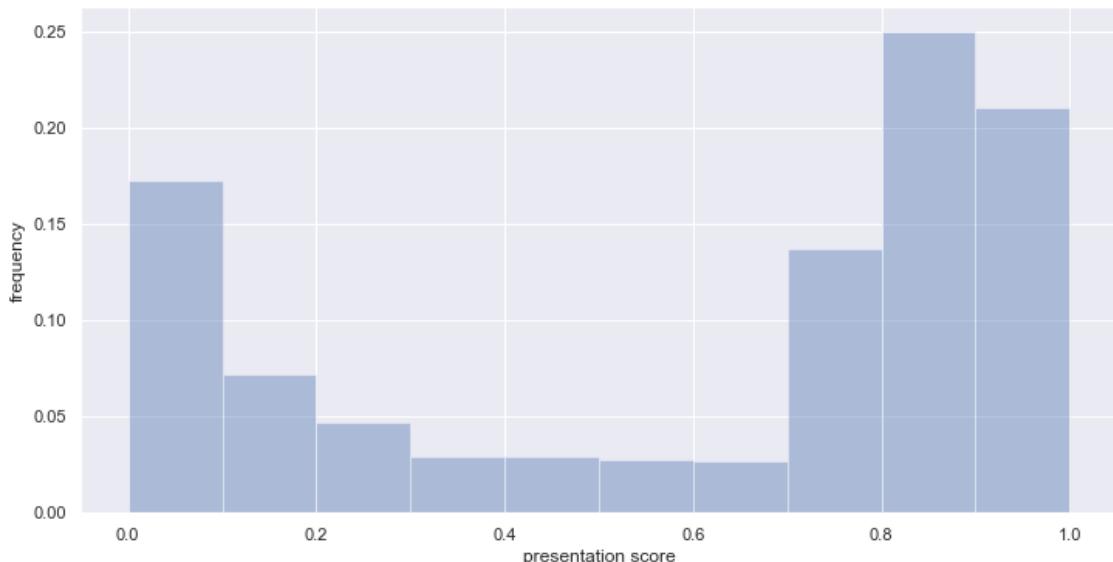
### 3. amino acid action

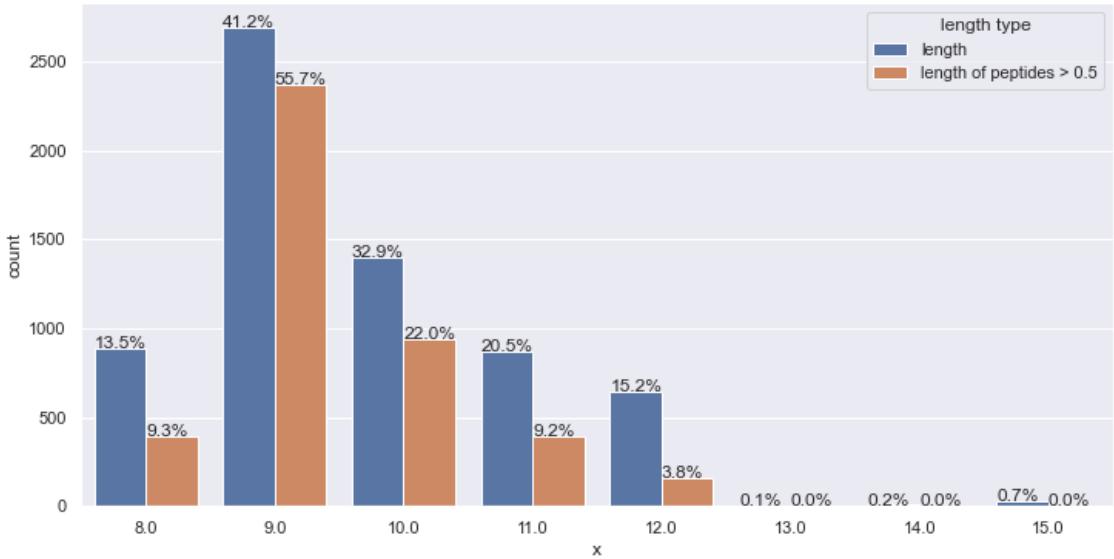


### 4. pos & amino action

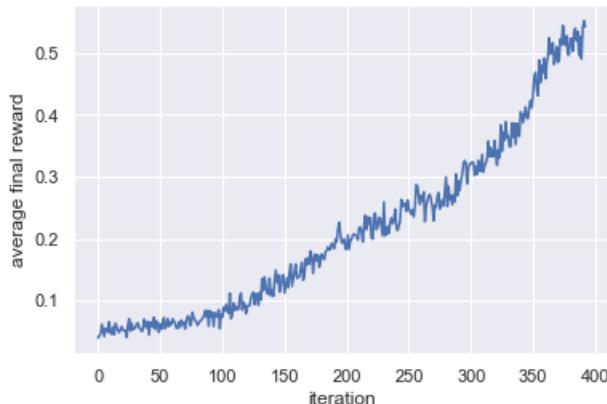
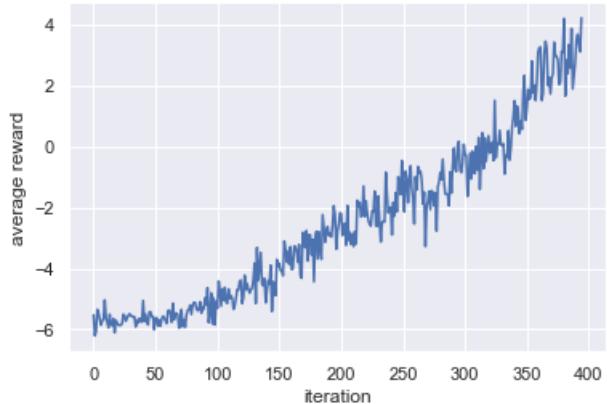


presentation score analysis: (65.1% optimized peptides have presentation scores > 0.5)

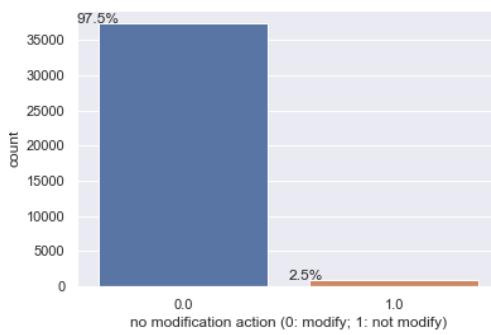
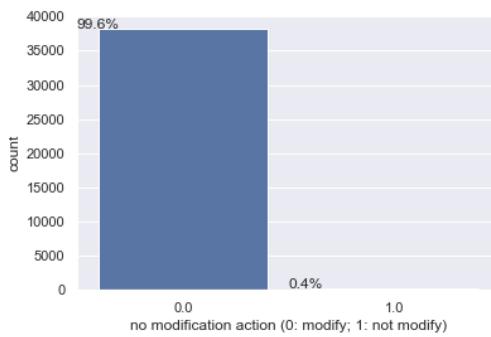
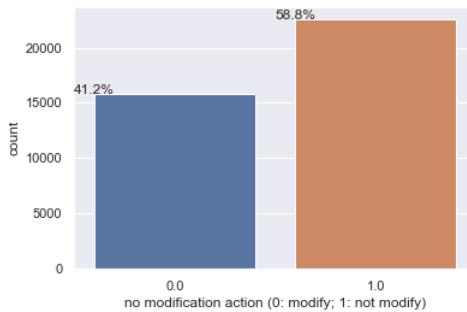




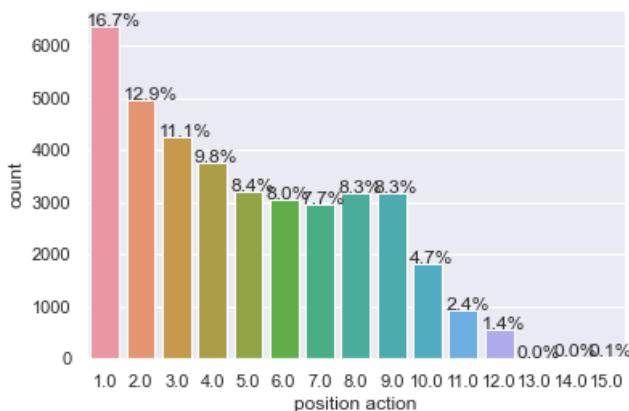
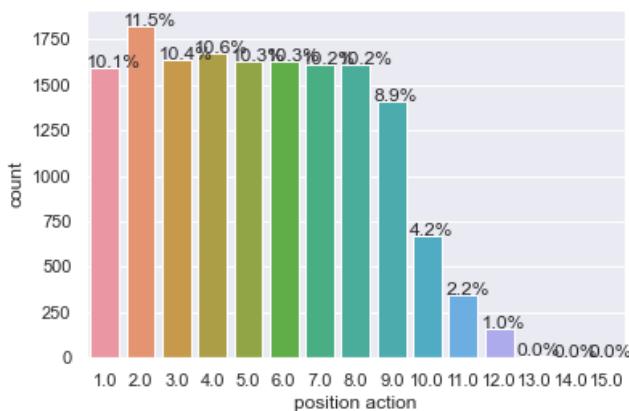
- Setting 1 (+1, -1, -0.9) (stop criteria = 1.0; sample rate = 50% from IEDB and 50% from random)

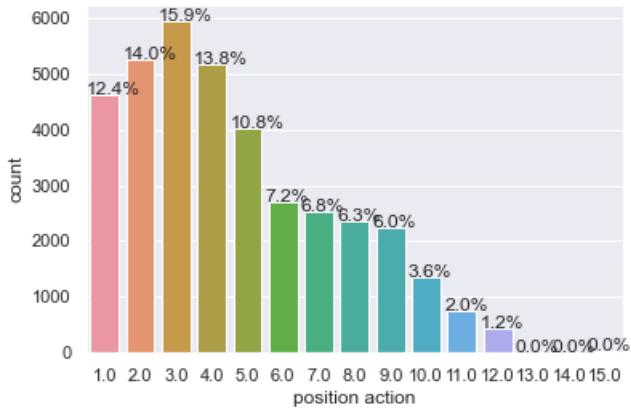


1. no modification action

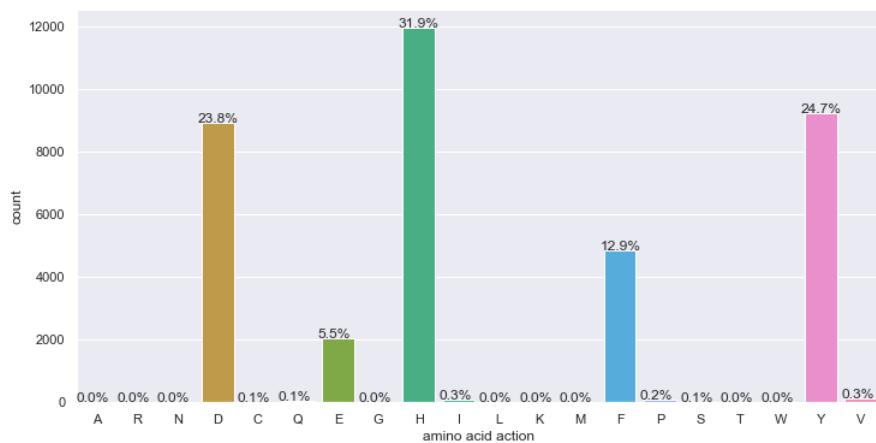
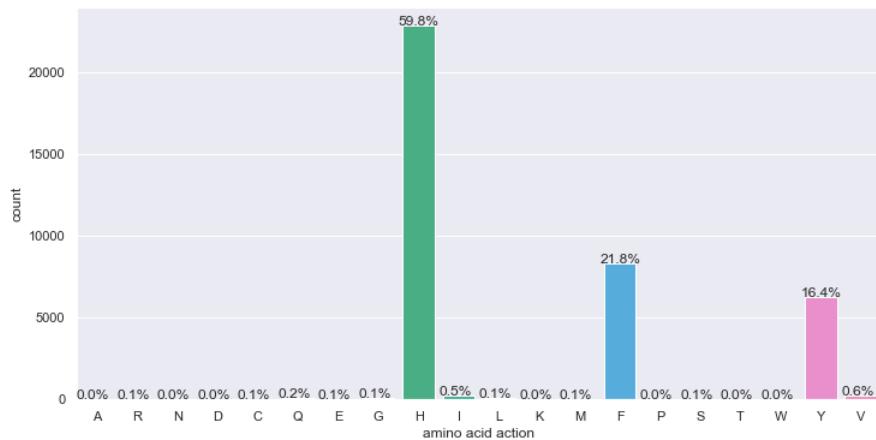
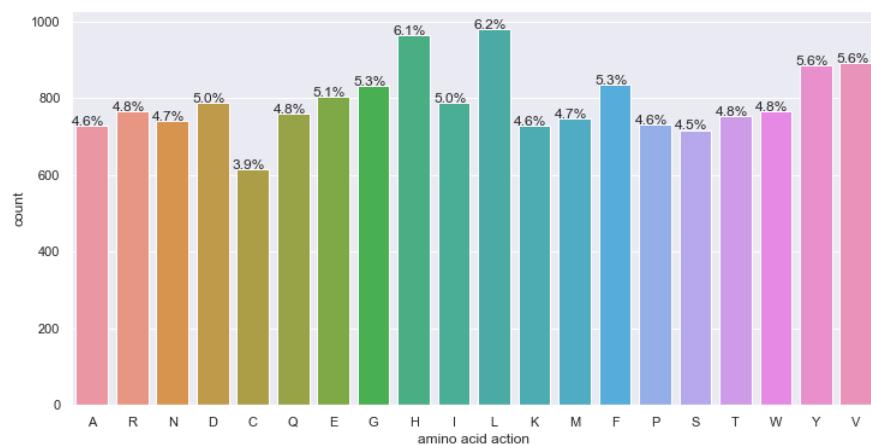


## 2. position select action

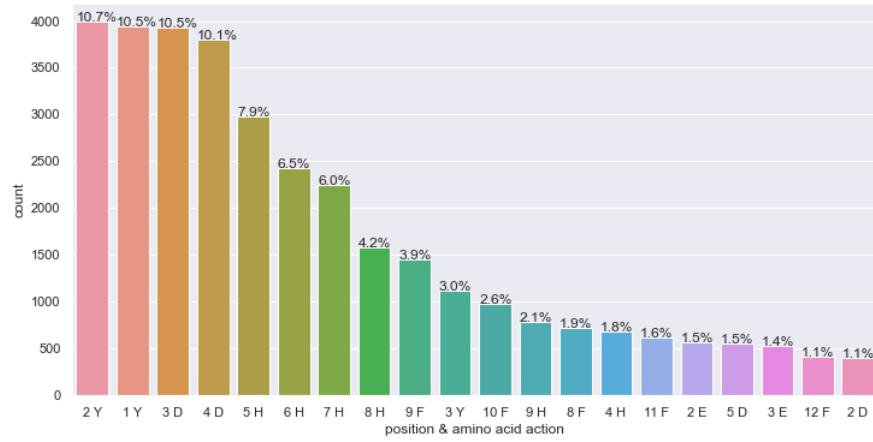
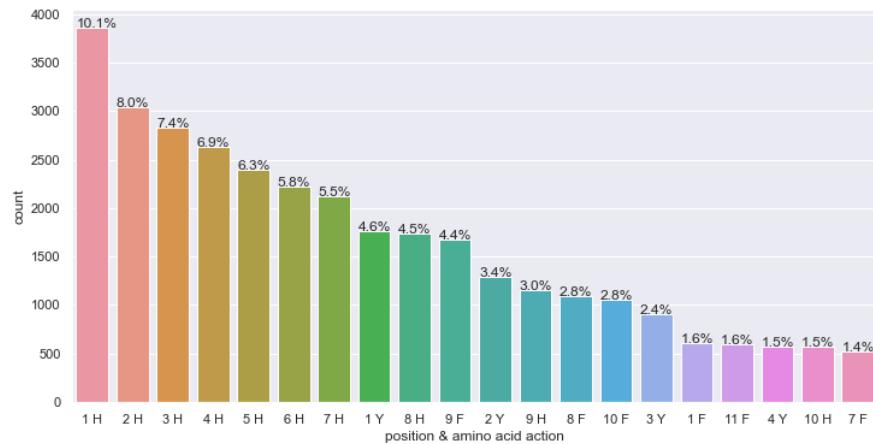
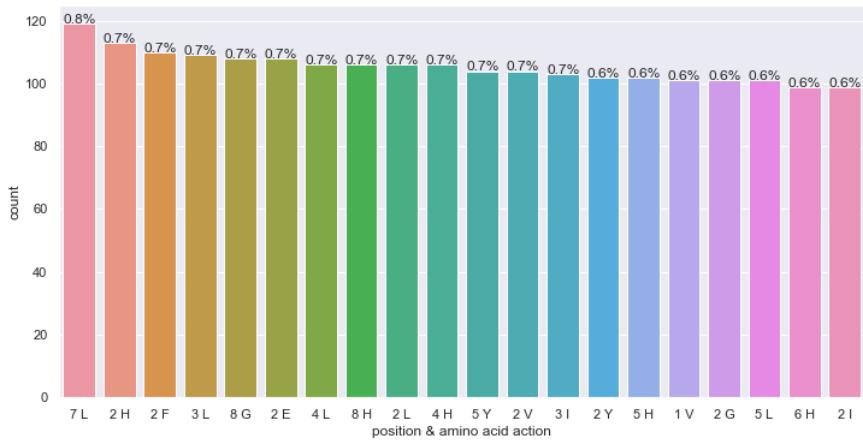




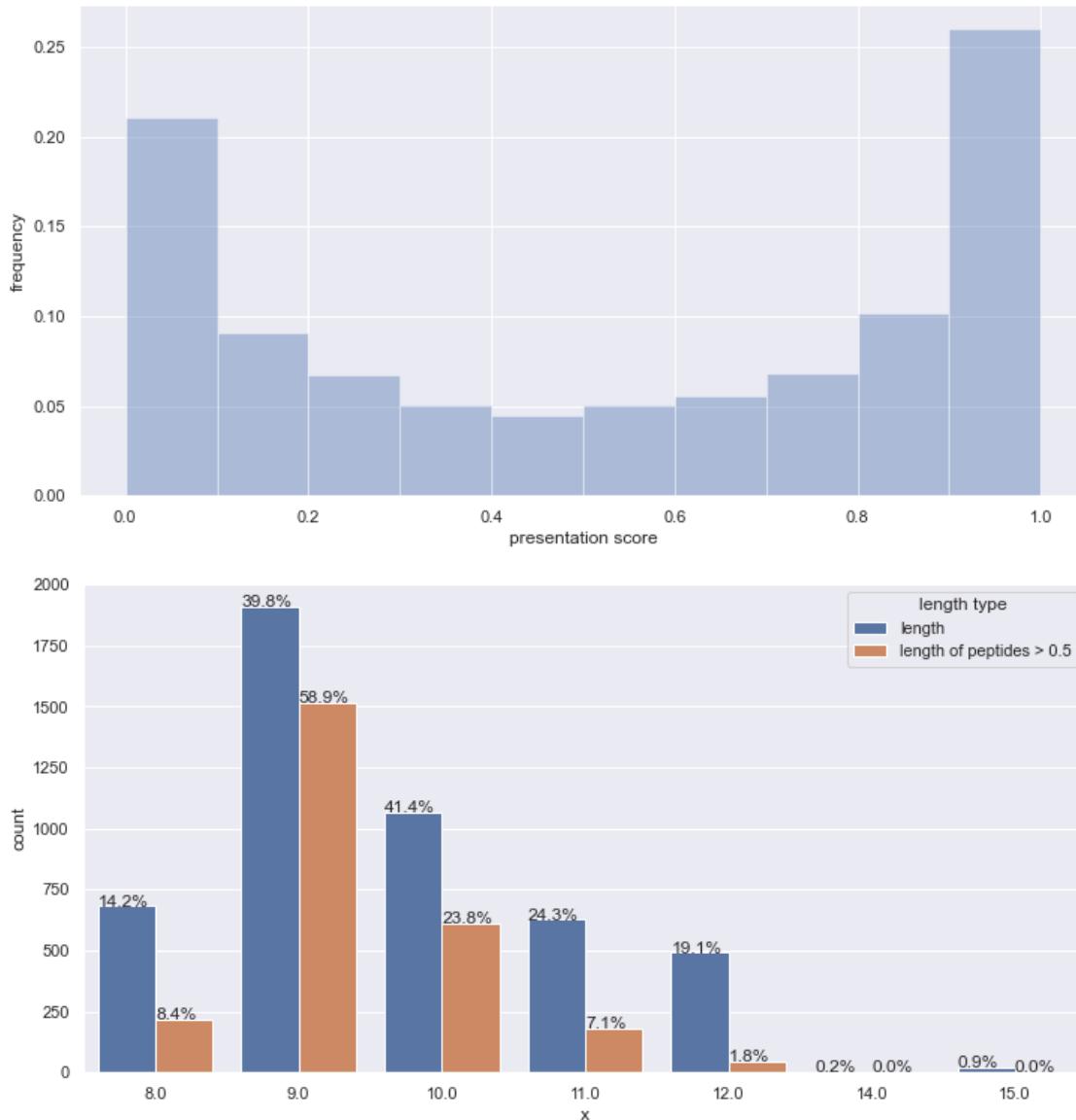
### 3. amino acid action



#### 4. pos & amino action

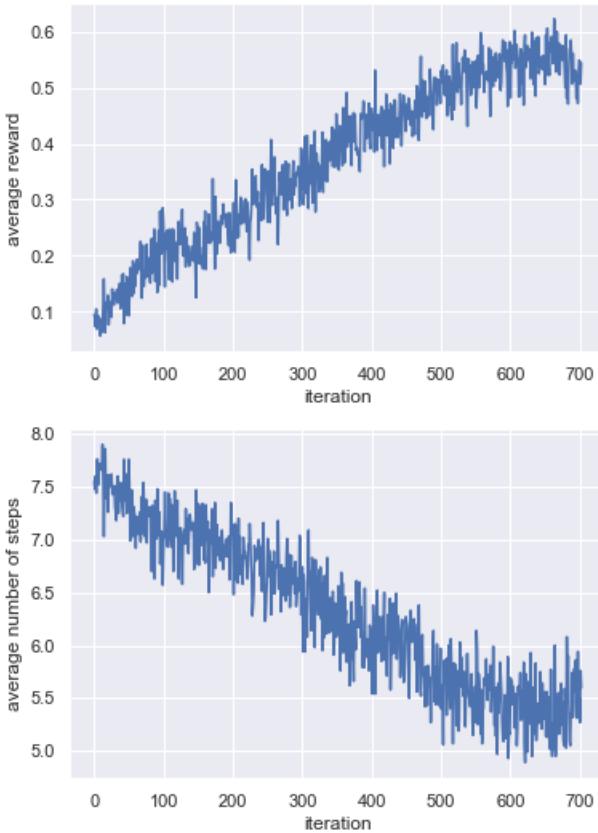


presentation score analysis: (53.6% optimized peptides have presentation scores > 0.5)

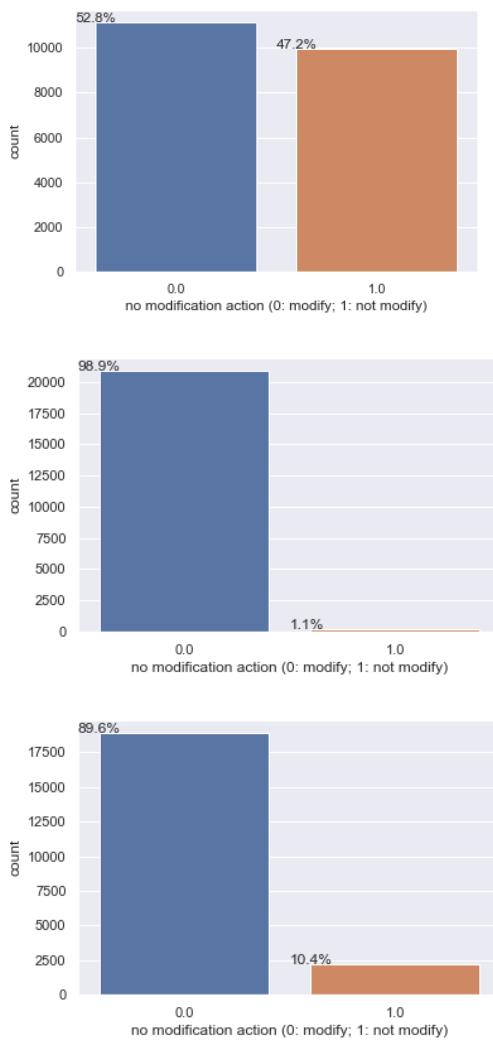


### 3. Result Analysis with only final rewards

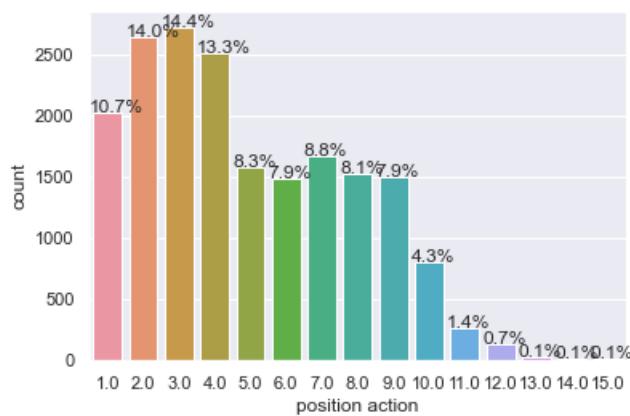
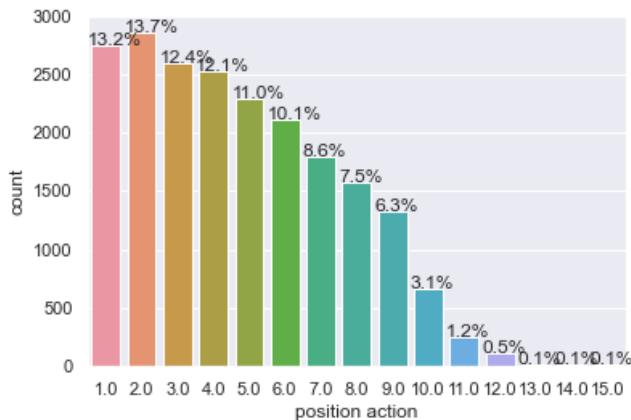
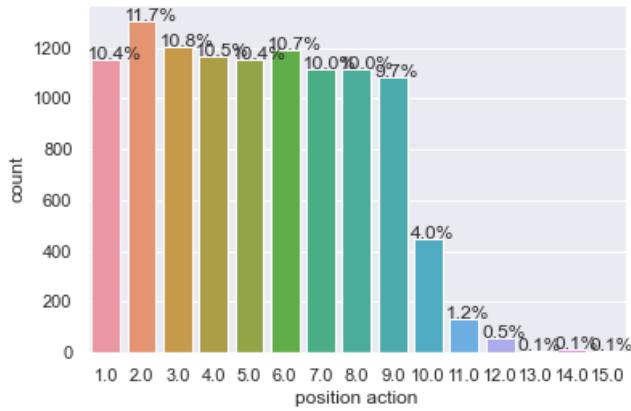
- Setting 1 (stop criteria = 0.6; sample rate = 80% from IEDB and 20% from random)



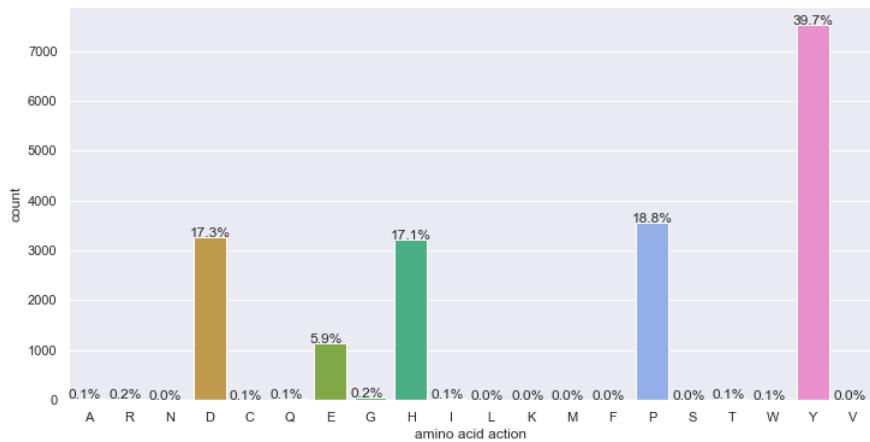
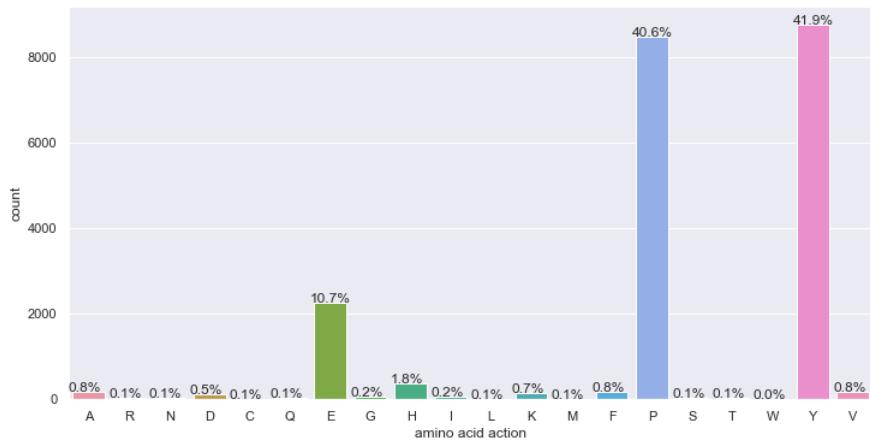
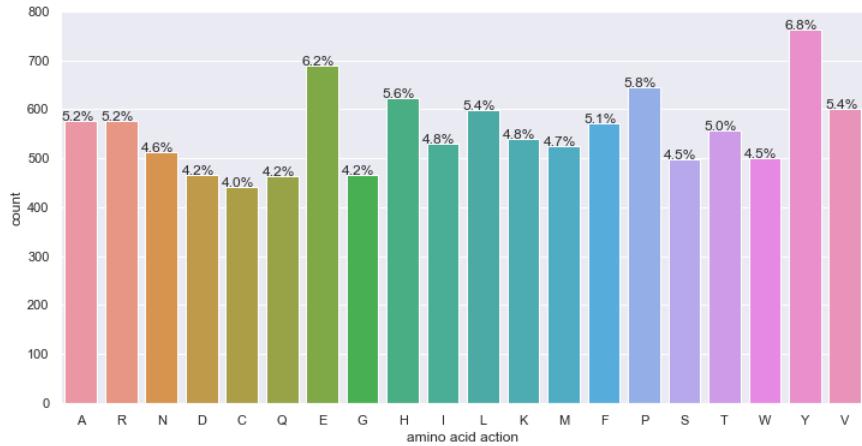
1. no modification action



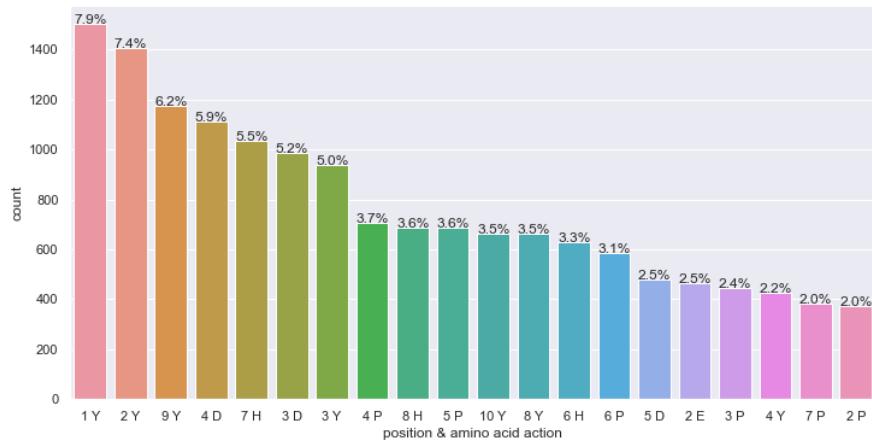
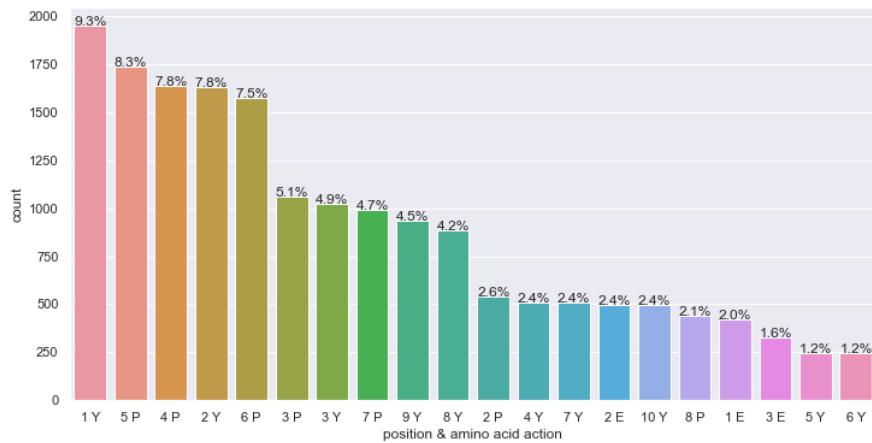
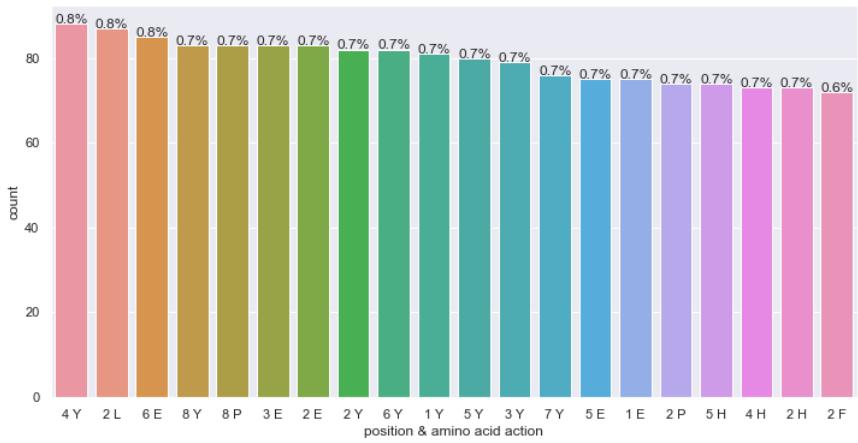
2. position select action



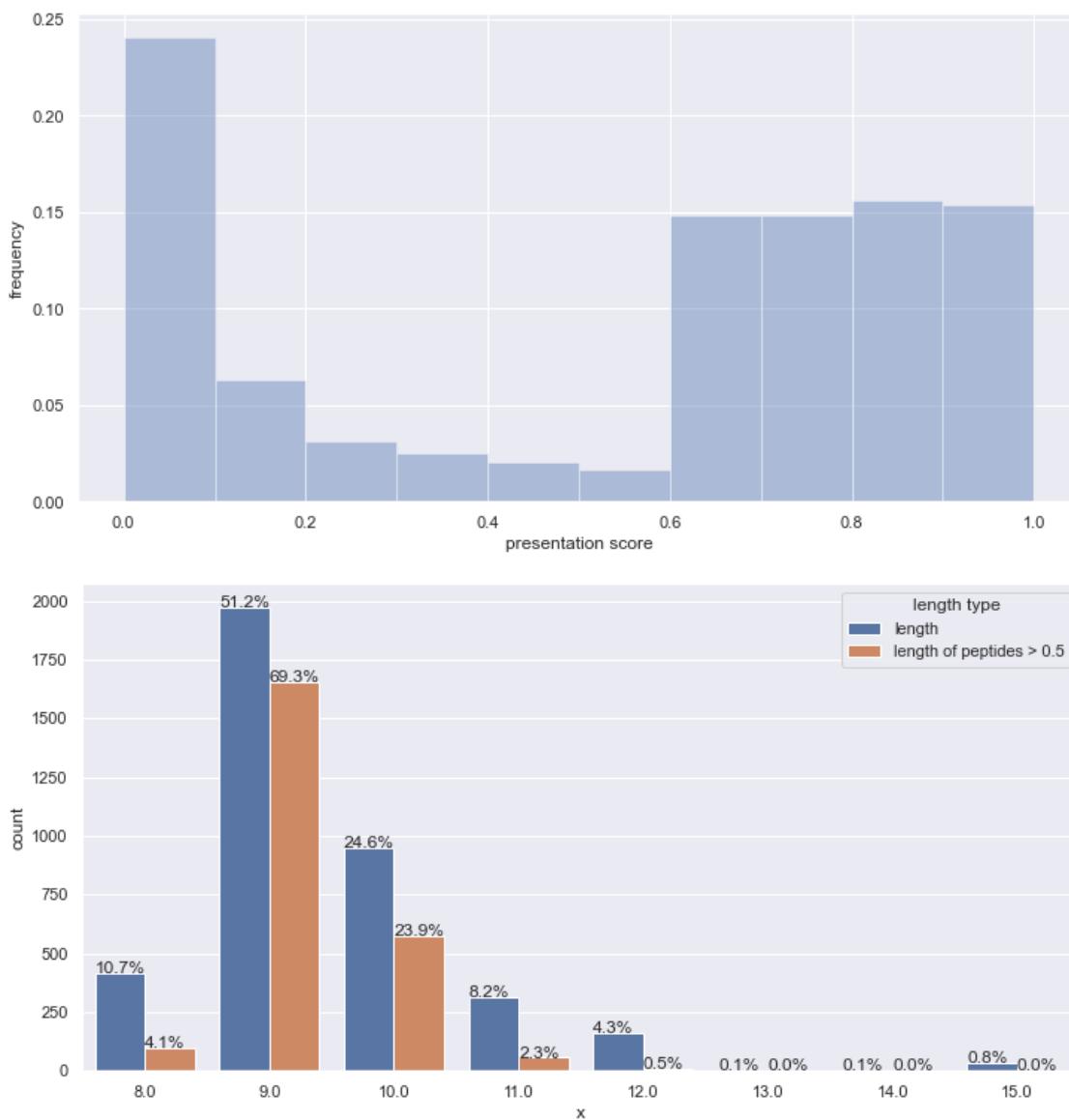
### 3. amino acid action



#### 4. pos & amino action



presentation score analysis: (53.6% optimized peptides have presentation scores > 0.5)



- Setting 1 (stop criteria = 0.6; sample rate = 80% from IEDB and 20% from random)

## How to group human peptides into buckets using LSH?

human peptides: 197,850,667

Below are the first 3 peptides in the file:

MAFSAEDVL

AFSAEDVLK

FSAEDVLKE

.....

All the human peptides and virus peptides are extracted by sliding a 9-mer window along the human genome or Cov-19 virus.

Note that in our experiments, the distance we care about should be hamming distance instead of edit distance. Therefore, instead of simply counting the occurrences of  $k$ -mers, we need to take the position of  $k$ -mers into consideration.

However, 3 amino acid differences on hamming distance can significantly change the  $k$ -mers in the peptides. For example, "XAFSXEDVX" has 3 amino acid difference from the human peptide "MAFSAEDVL". However, these two peptides have only two common 3-mers among all the seven 3-mers, and 4 common 2-mers among all the 8 2-mers. With 3 amino acid differences, for the

human peptide "MAFSAEDVL", the above example "**XAFSXEDVX**" or "**XAFXAEDVX**" will change the common k-mers most significantly while the example "**XXXSAEDVL**" will lead to the least amount of change on the common k-mers.

For contiguous 3-mers,

For discontiguous *k*-mers,

Inspired by the above analysis, the question is **how we can find a hashing function so that it can guarantee that peptides with less than 3 amino acid differences can be grouped into the same buckets?**

"**XXXSAEDVL**"

"**XXFXAEDVL**"

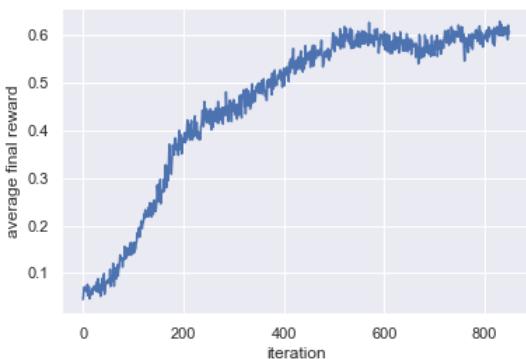
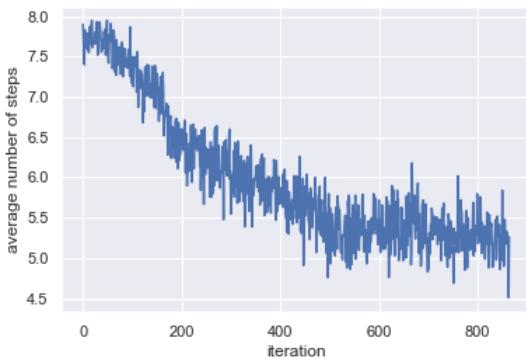
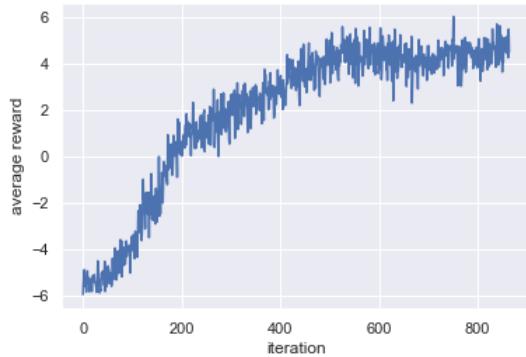
"**XXFSXEDVL**"

## **Update (5-26)**

### **1. Result Analysis with longer running time and different stopping criteria**

**Experiment: stop criteria = 0.6**

Setting 1 (sample rate = 50% from IEDB and 50% from random)

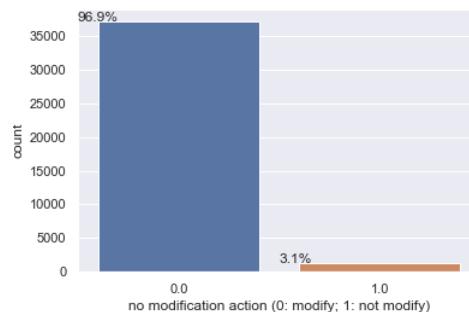
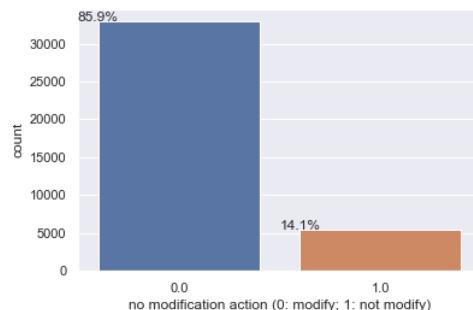
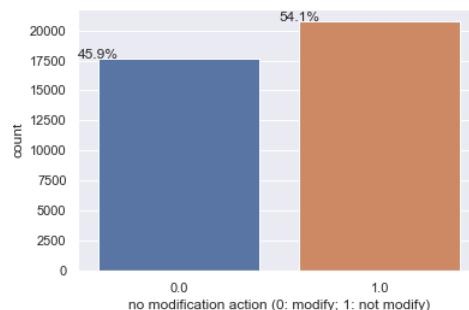


### 1. "No modification" action analysis

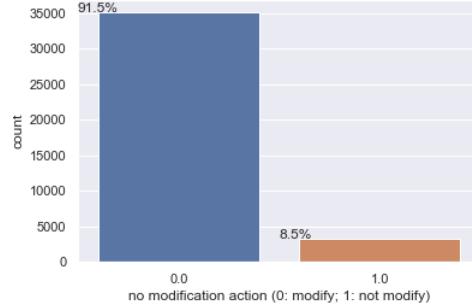
(first 10 iterations)  
the middle)

(10 iterations at 1/4)

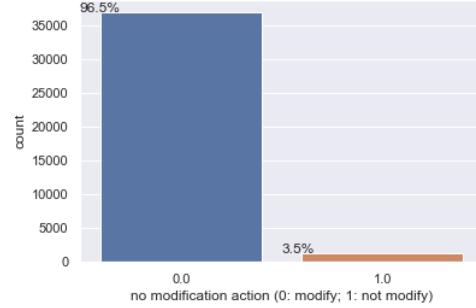
(10 iterations in



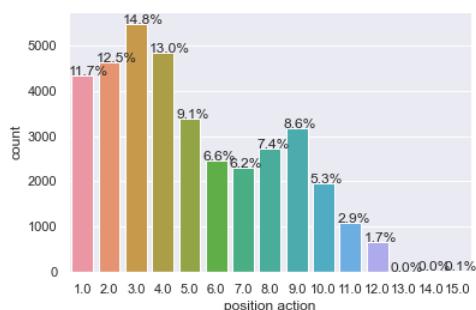
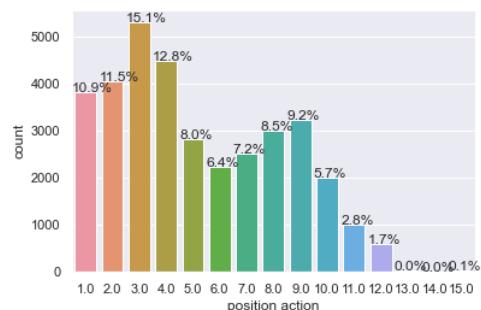
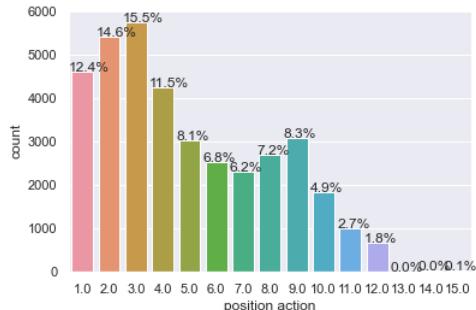
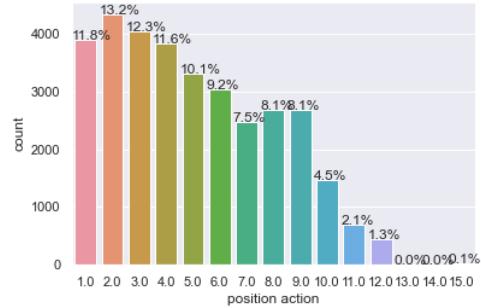
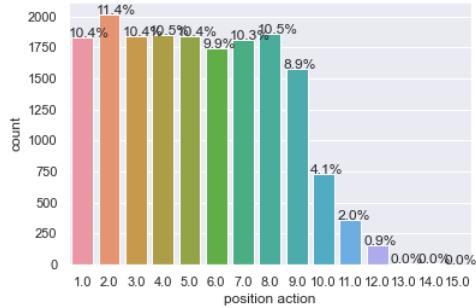
(10 iterations at 3/4)



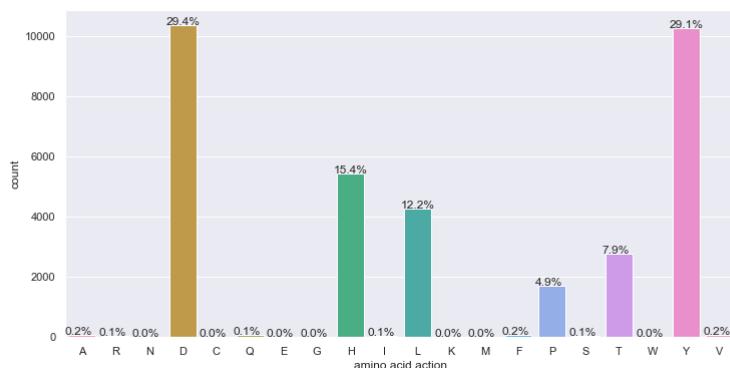
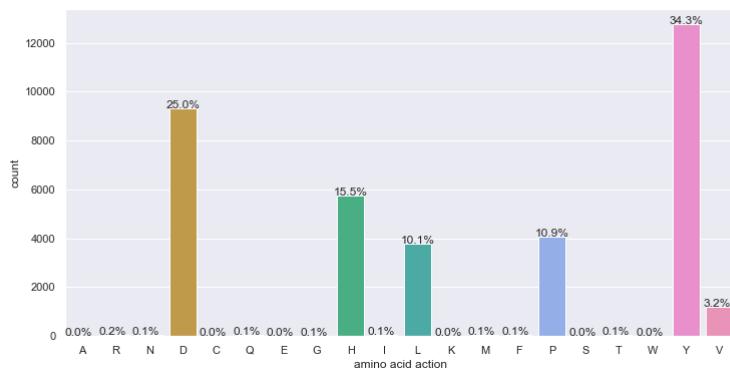
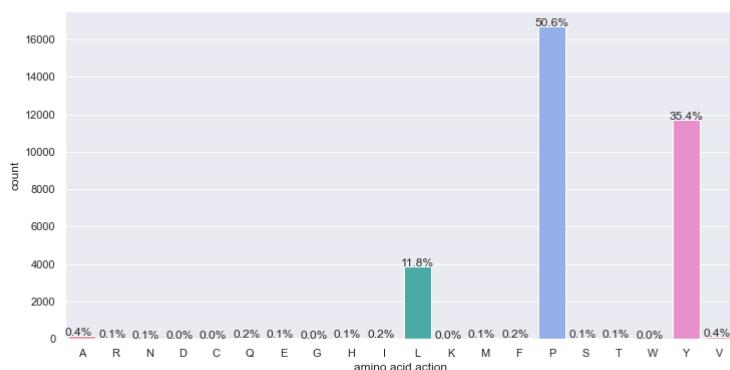
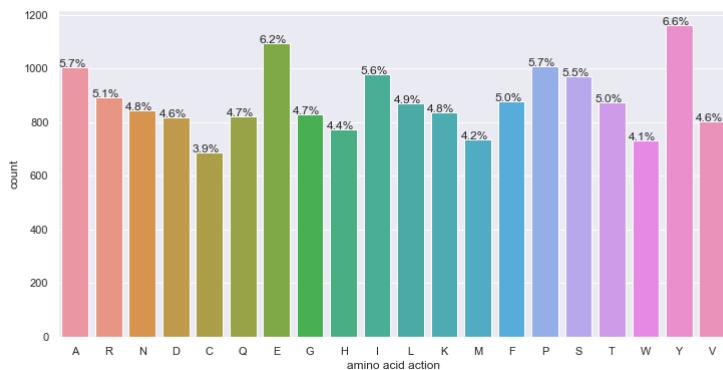
(last 10 iterations)

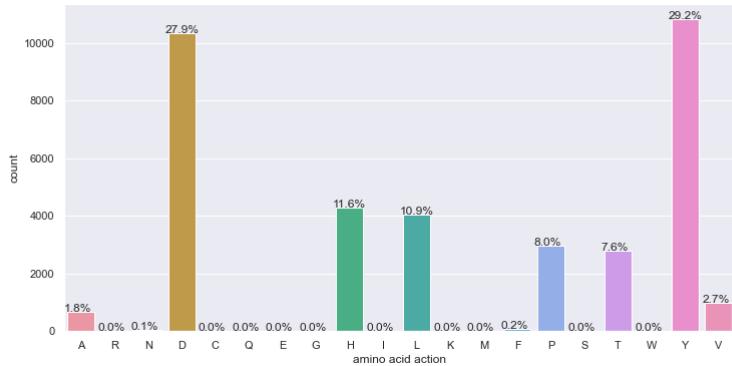


## 2. Position action analysis

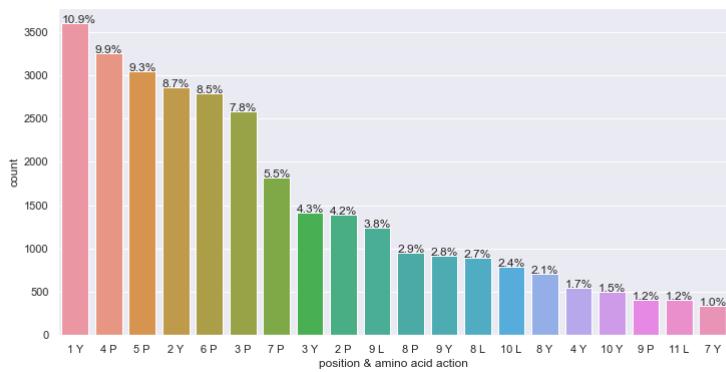
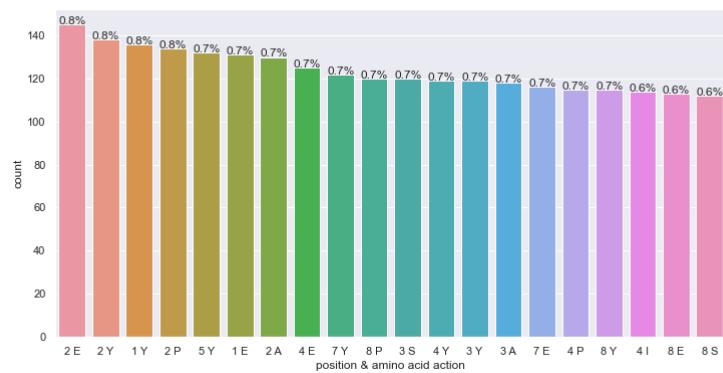


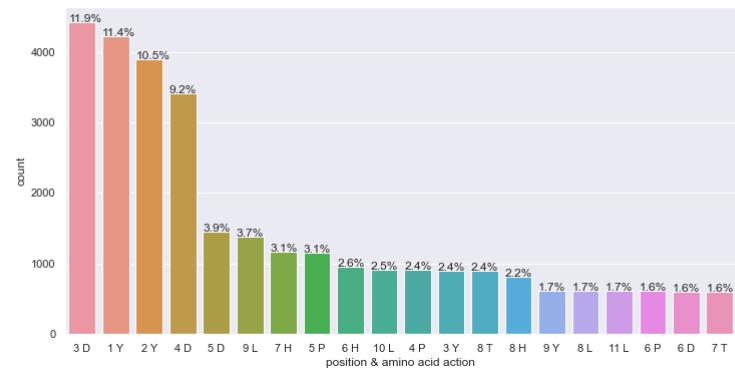
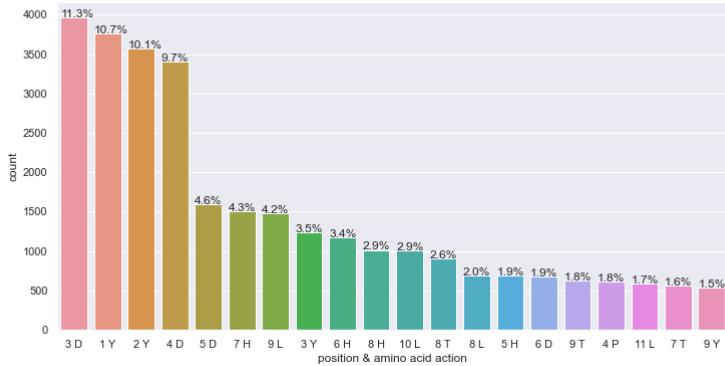
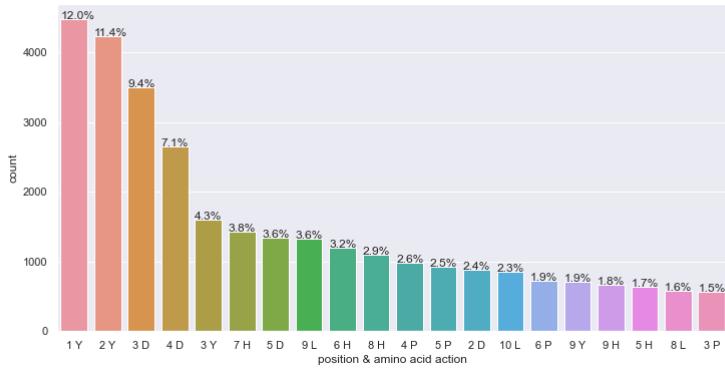
## 3. Amino action analysis



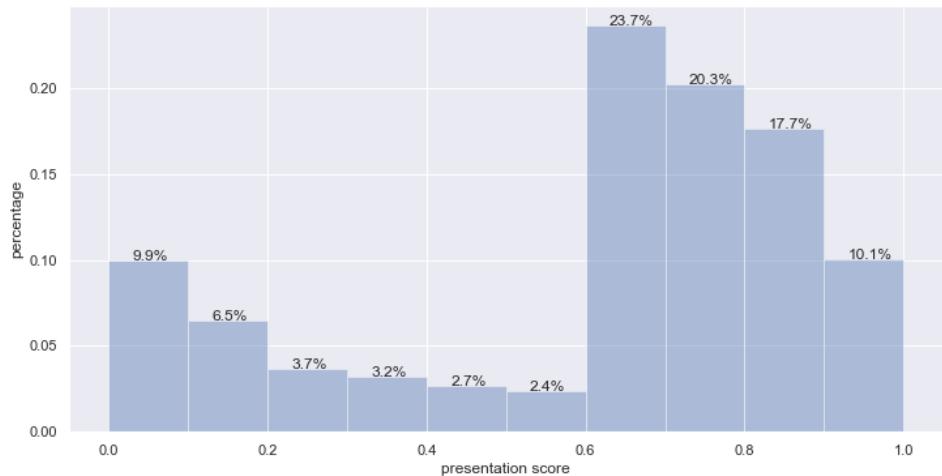


#### 4. Position & Amino action analysis

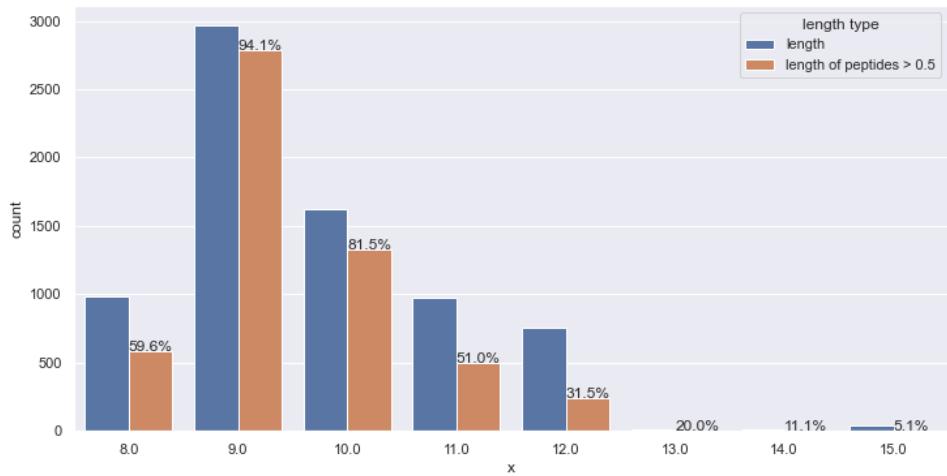




## 5. Reward analysis

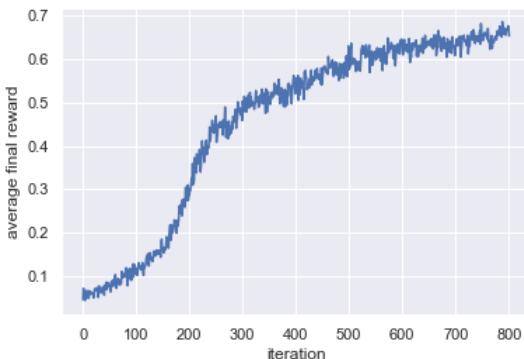
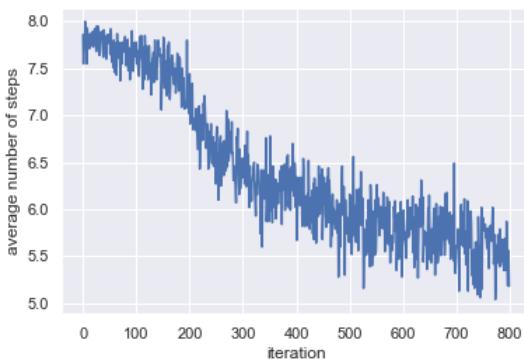
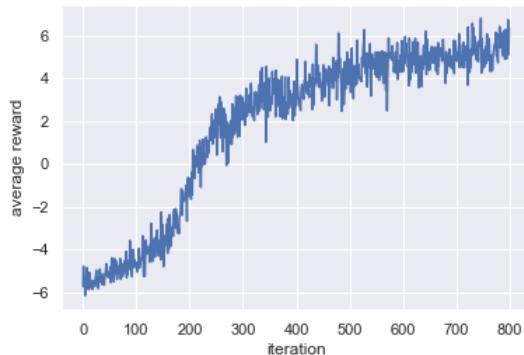


## 6. Length analysis



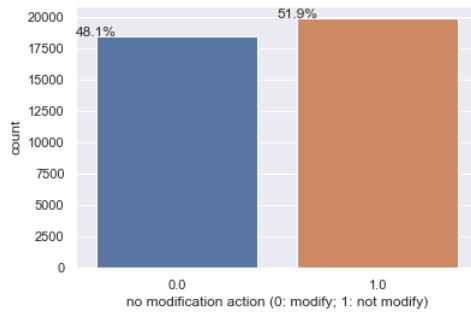
### Experiment: stop criteria = 0.75

Setting 1 (sample rate = 50% from IEDB and 50% from random)

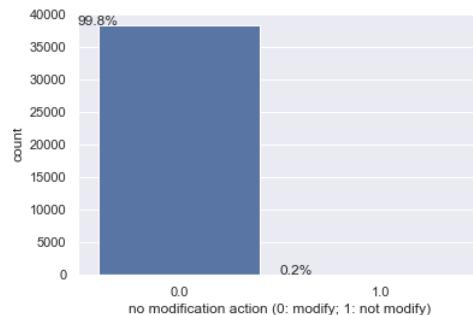


1. "No modification" action analysis

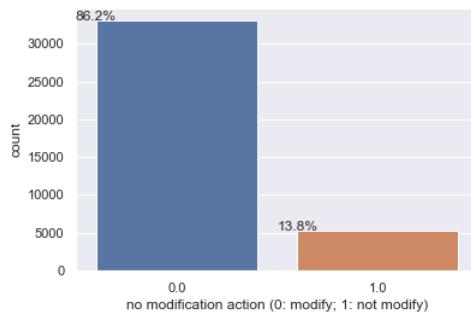
(first 10 iterations)  
the middle)



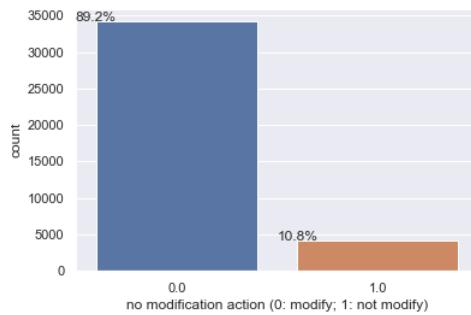
(10 iterations at 1/4)



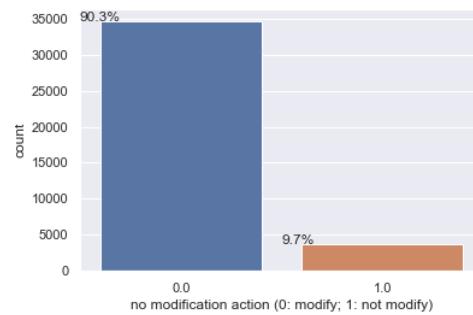
(10 iterations in  
the middle)



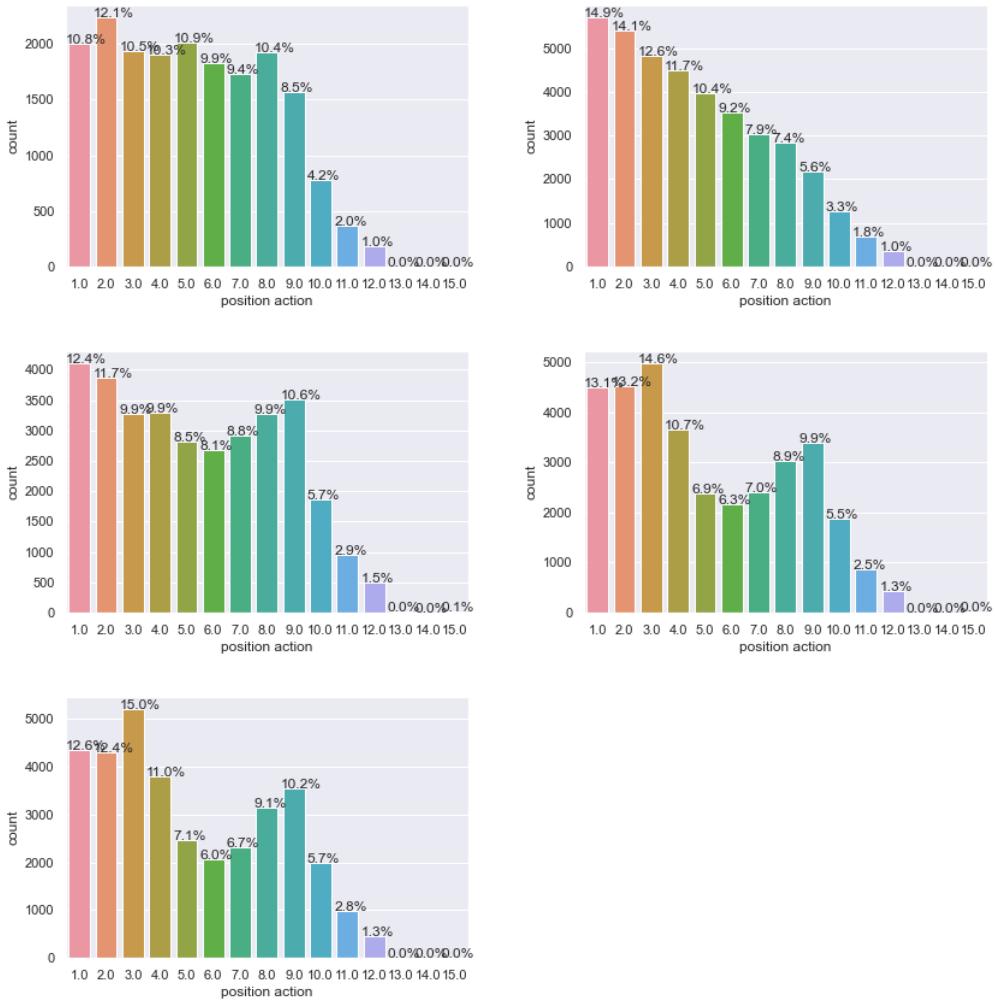
(10 iterations at 3/4)



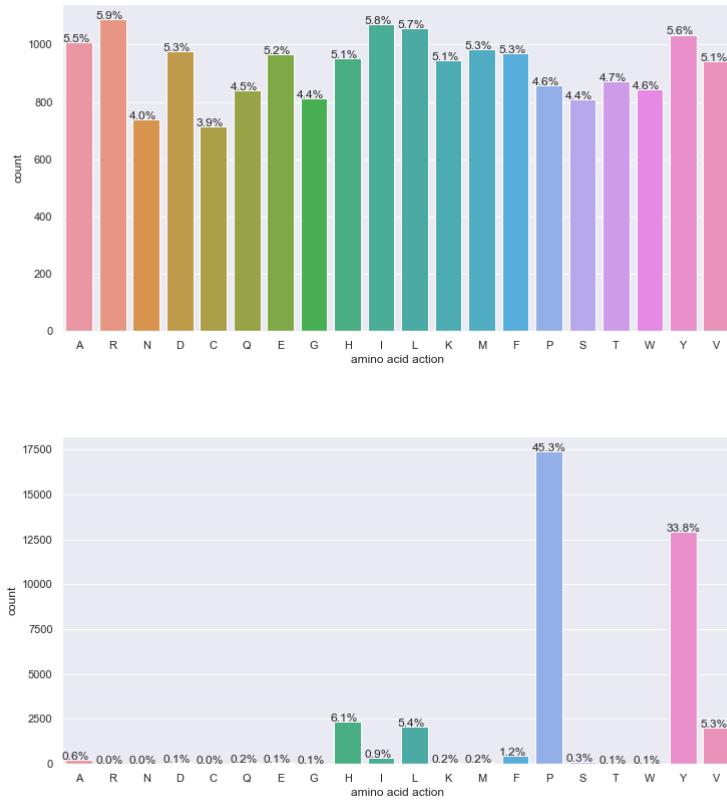
(last 10 iterations)

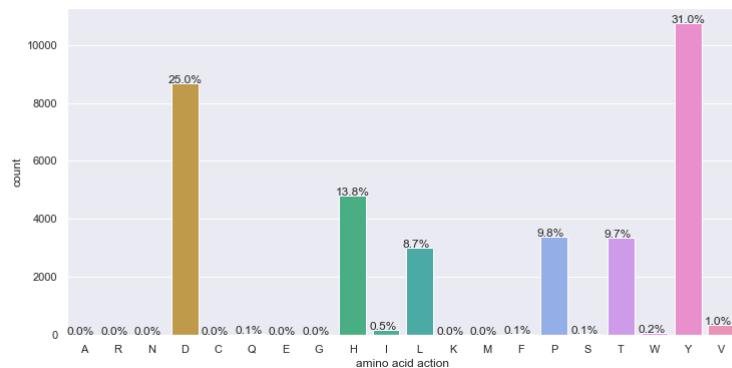
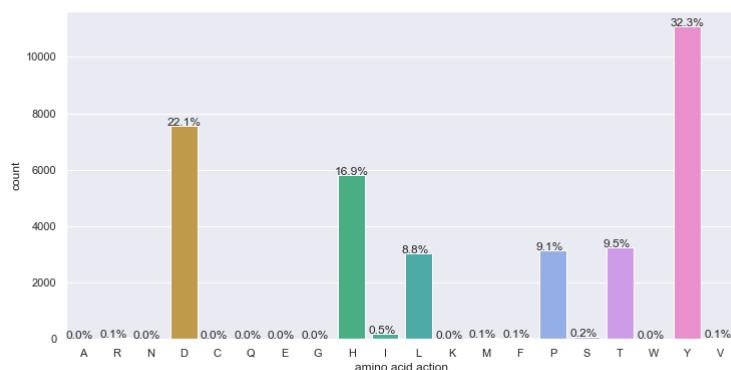
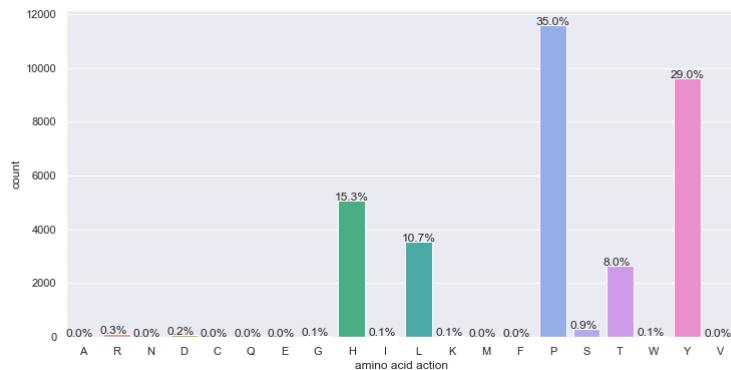


## 2. Position action analysis

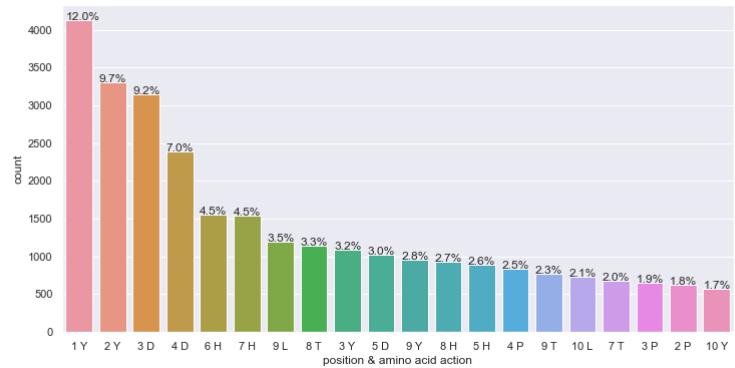
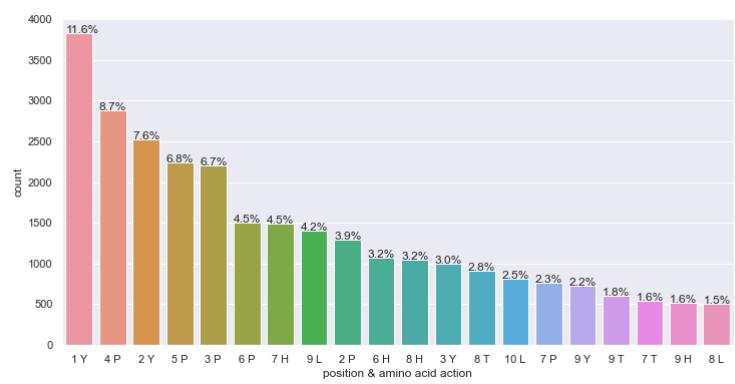
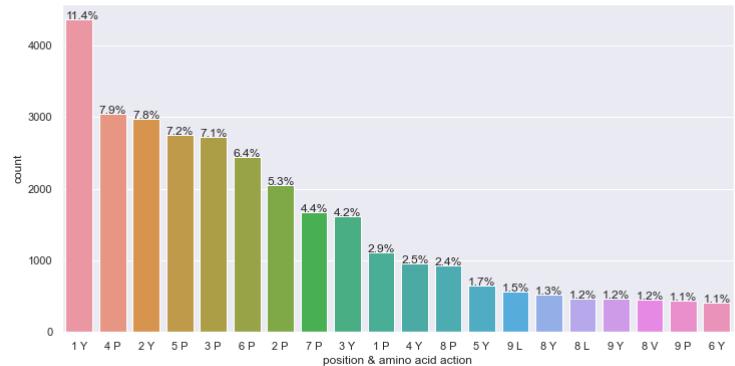
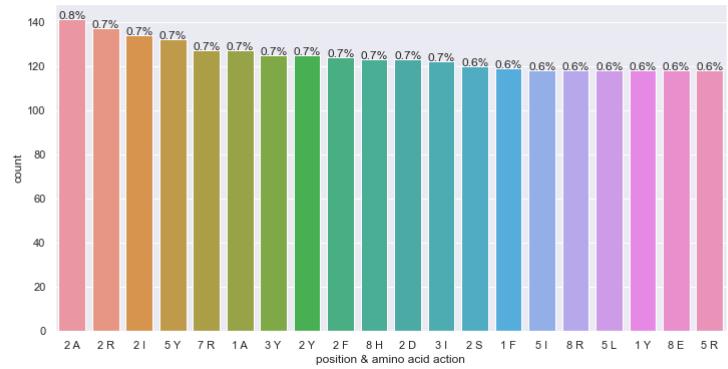


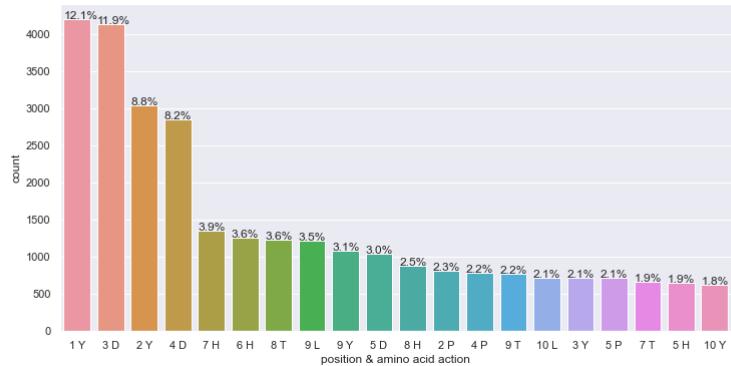
### 3. Amino action analysis



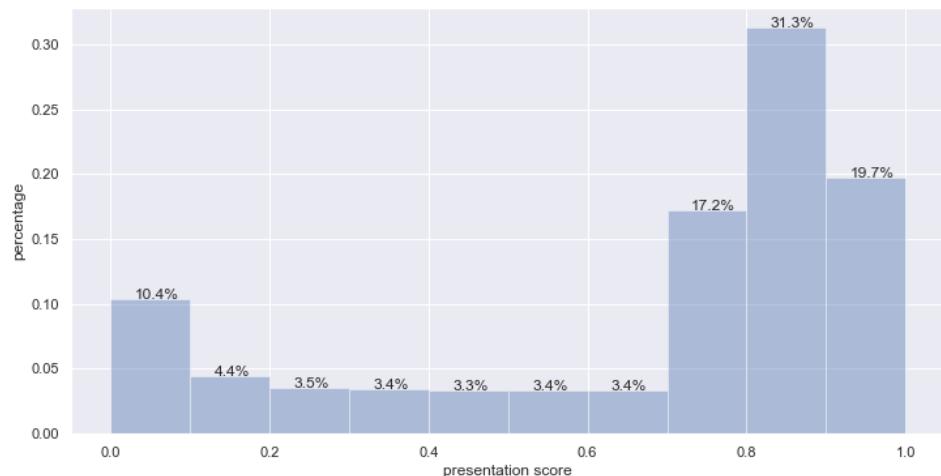


#### 4. Position & Amino action analysis

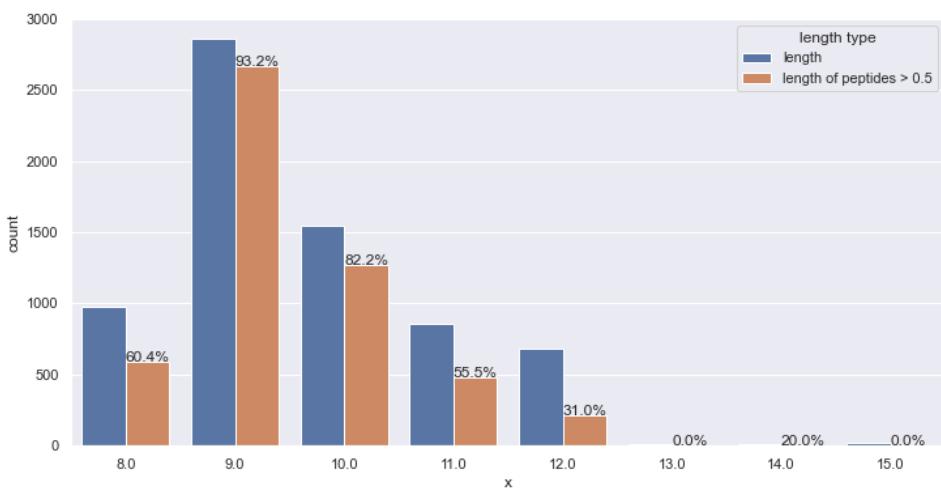




## 5. Reward analysis

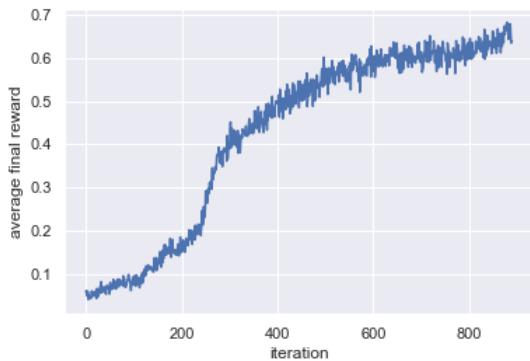
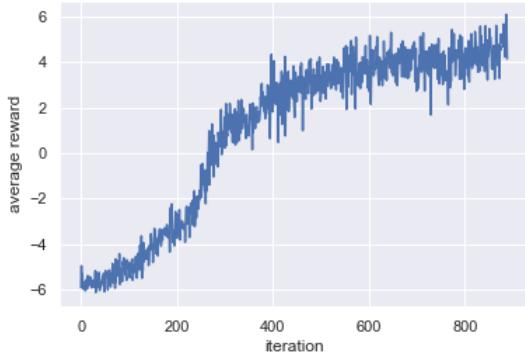


## 6. Length analysis



**Experiment: stop criteria = 1.0**

Setting 1 (sample rate = 50% from IEDB and 50% from random)

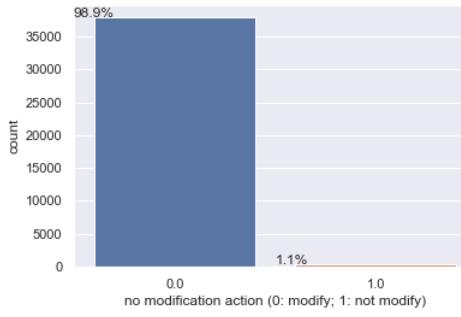
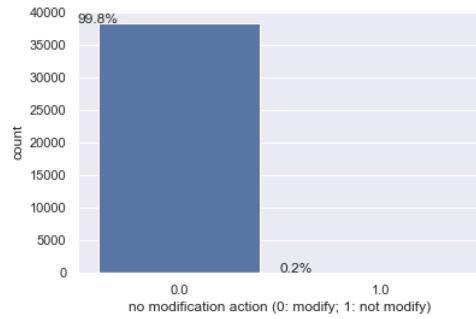
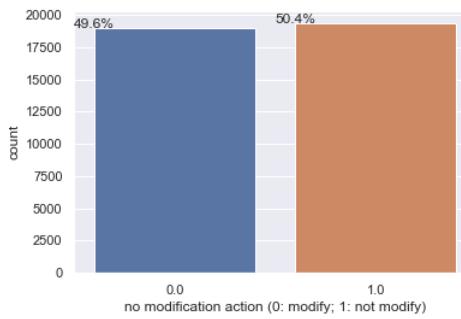


### 1. "No modification" action analysis

(first 10 iterations)  
the middle)

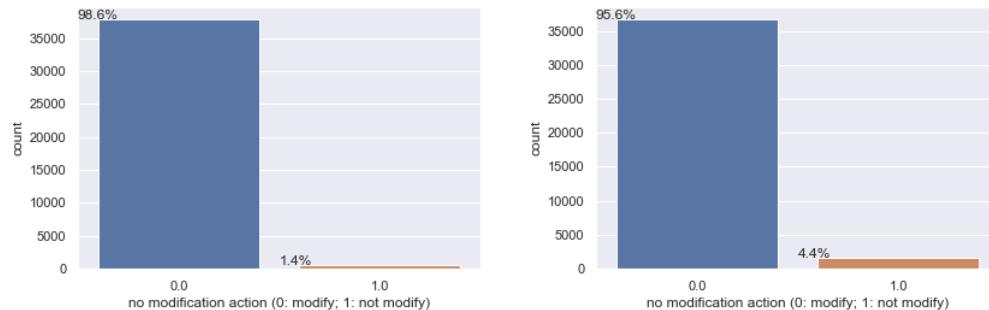
(10 iterations at 1/4)

(10 iterations in

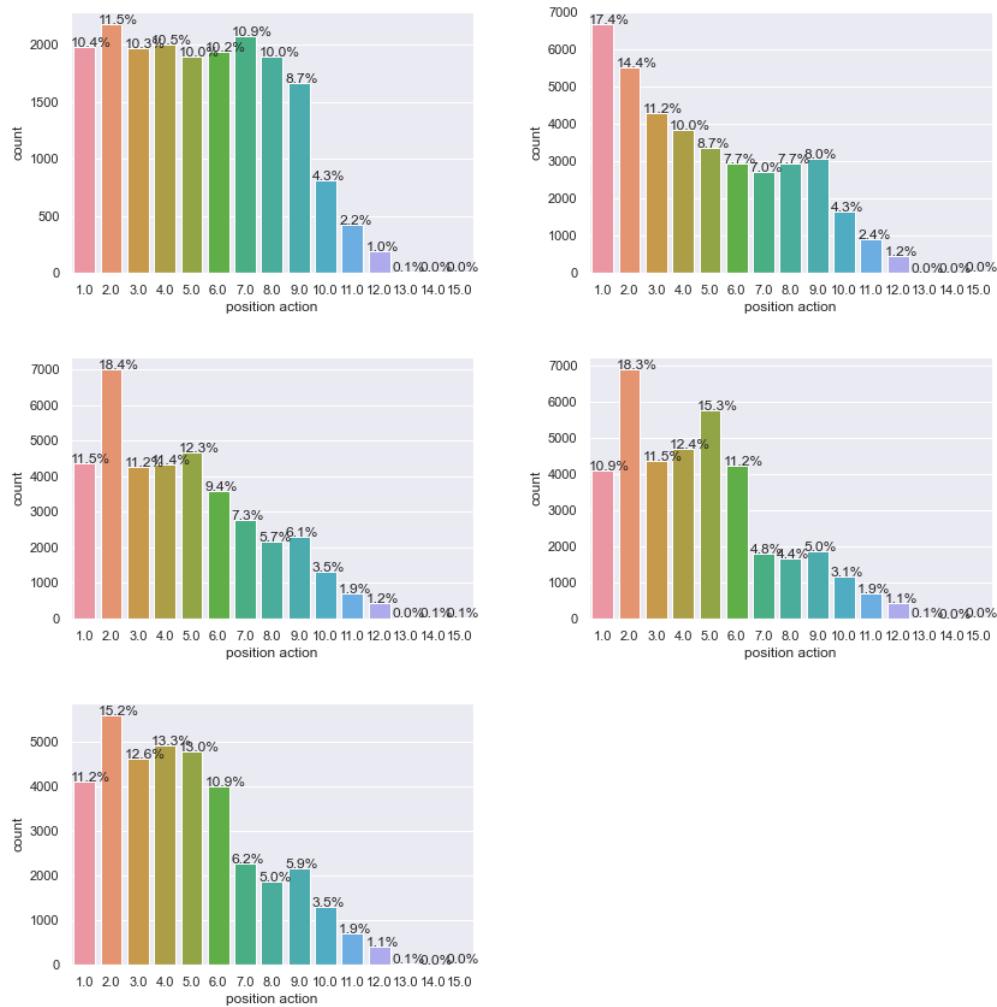


(10 iterations at 3/4)

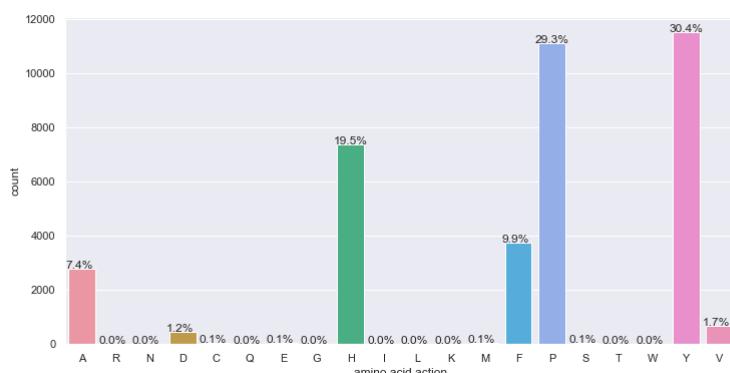
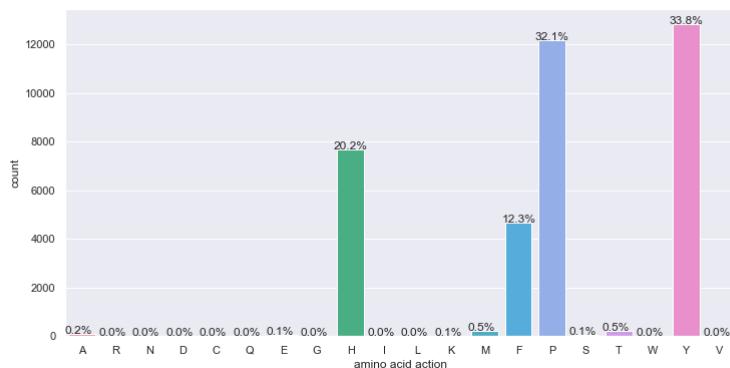
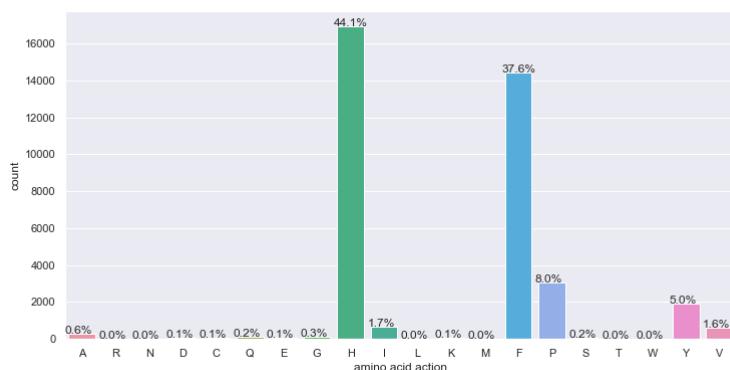
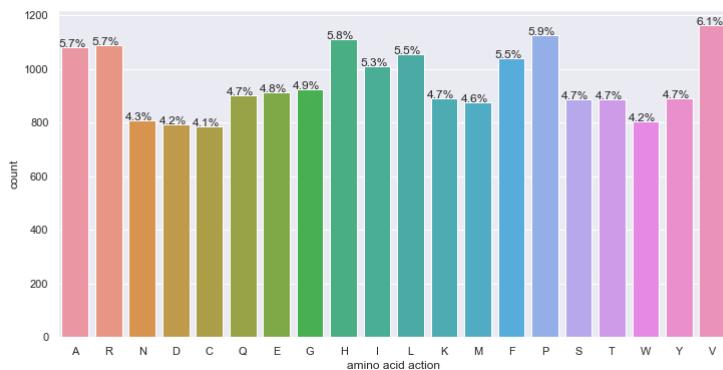
(last 10 iterations)

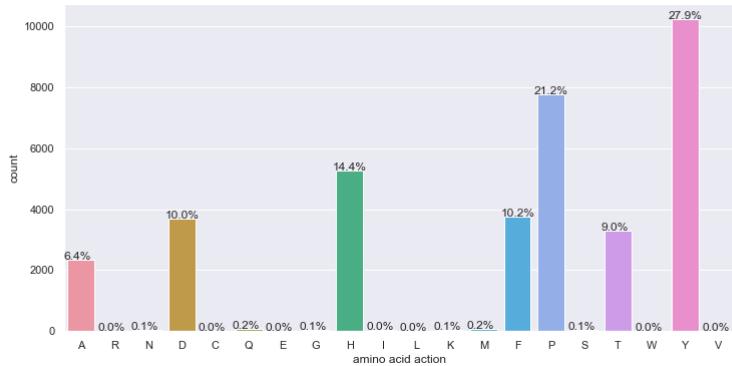


## 2. Position action analysis

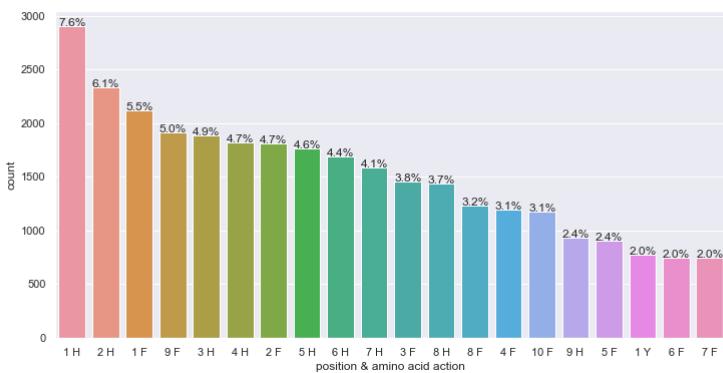
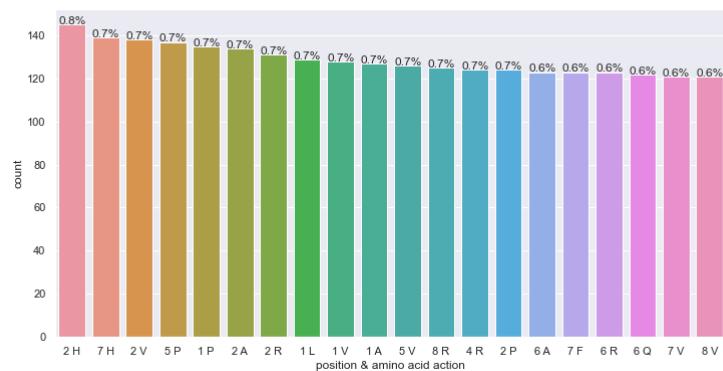


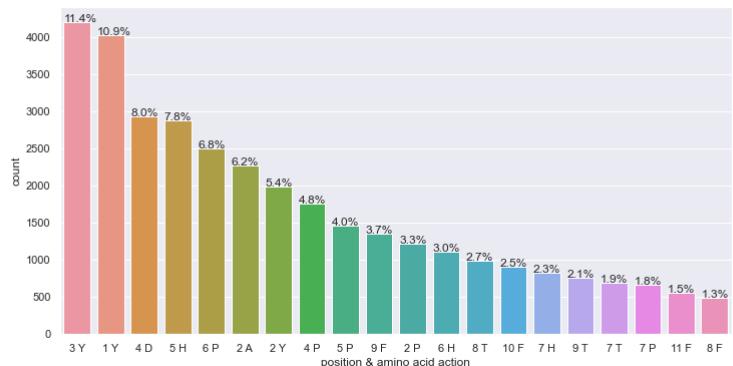
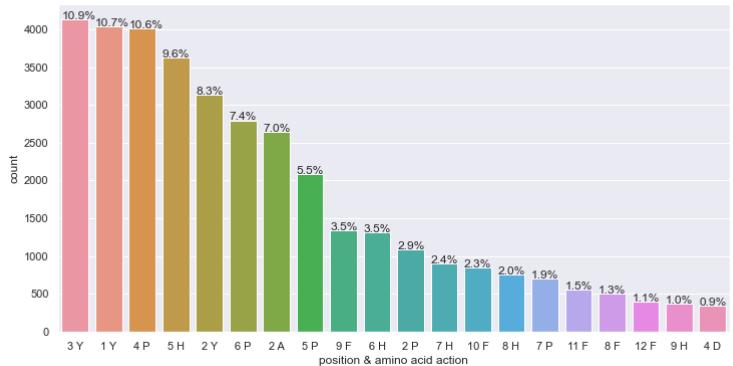
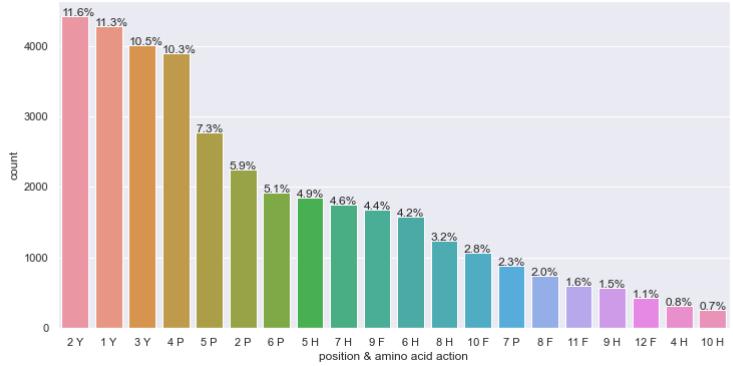
## 3. Amino action analysis



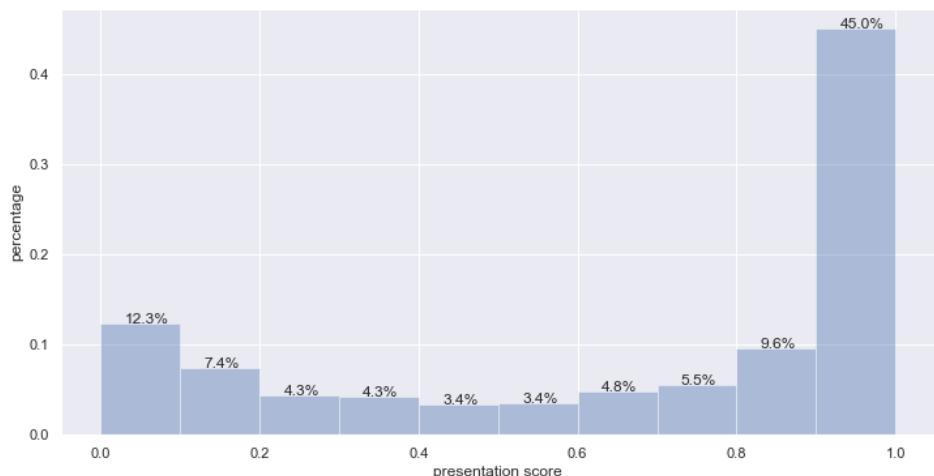


#### 4. Position & Amino action analysis

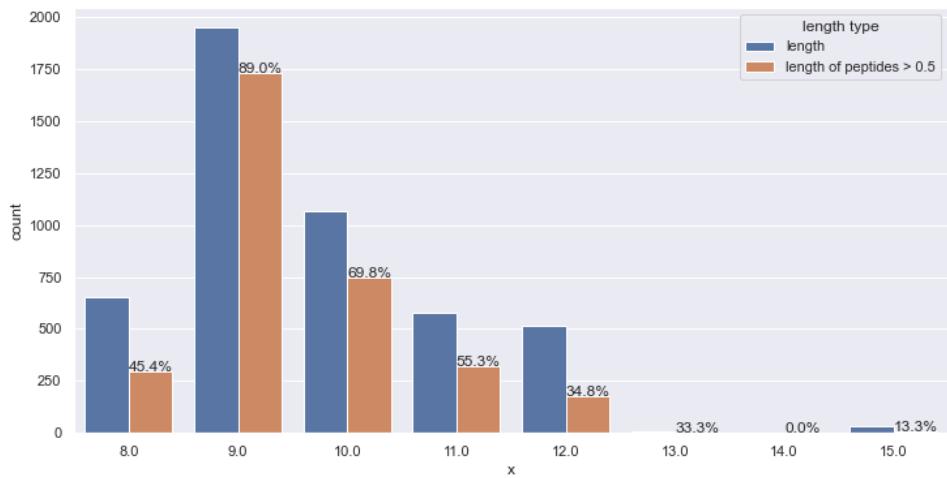




## 5. Reward analysis



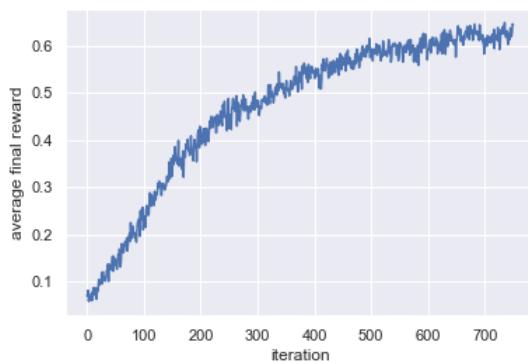
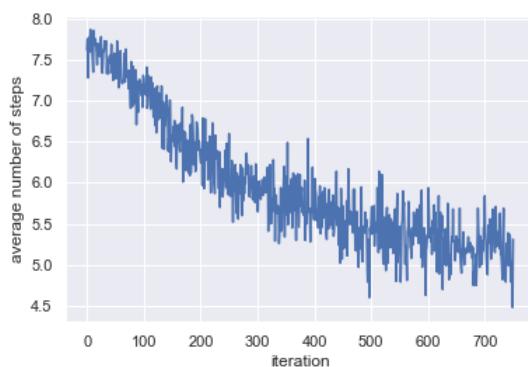
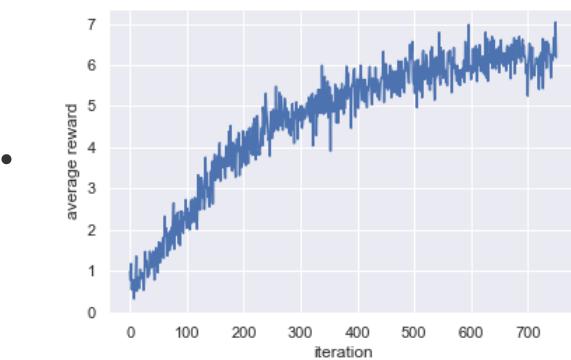
## 6. Length analysis



## 2. Result Analysis with only final rewards

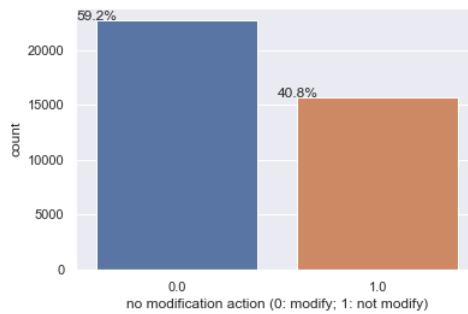
**Experiment: stop criteria = 0.6**

Setting 1 (only final rewards; sample rate = 50% from IEDB and 50% from random)

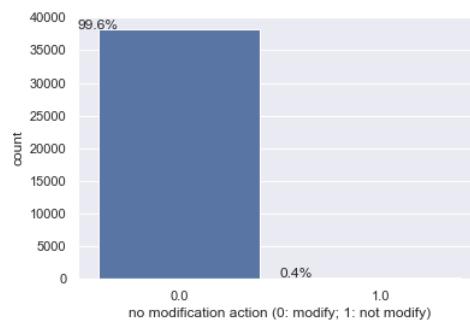


1. "No modification" action analysis

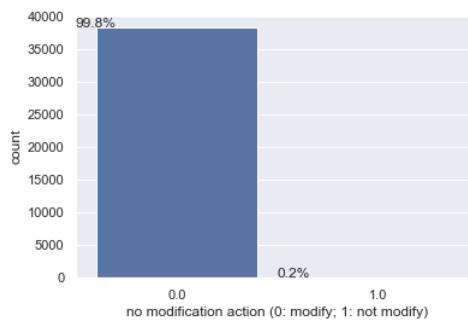
(first 10 iterations)  
the middle)



(10 iterations at 1/4)

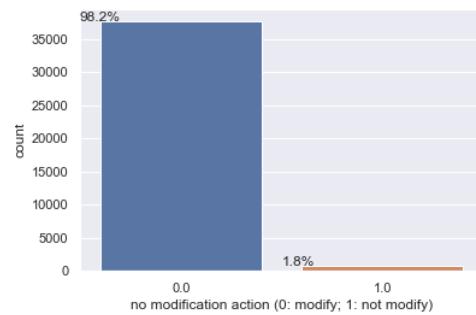
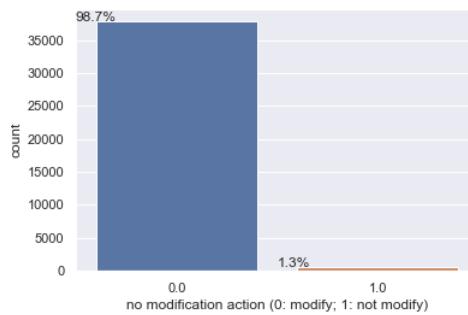


(10 iterations in

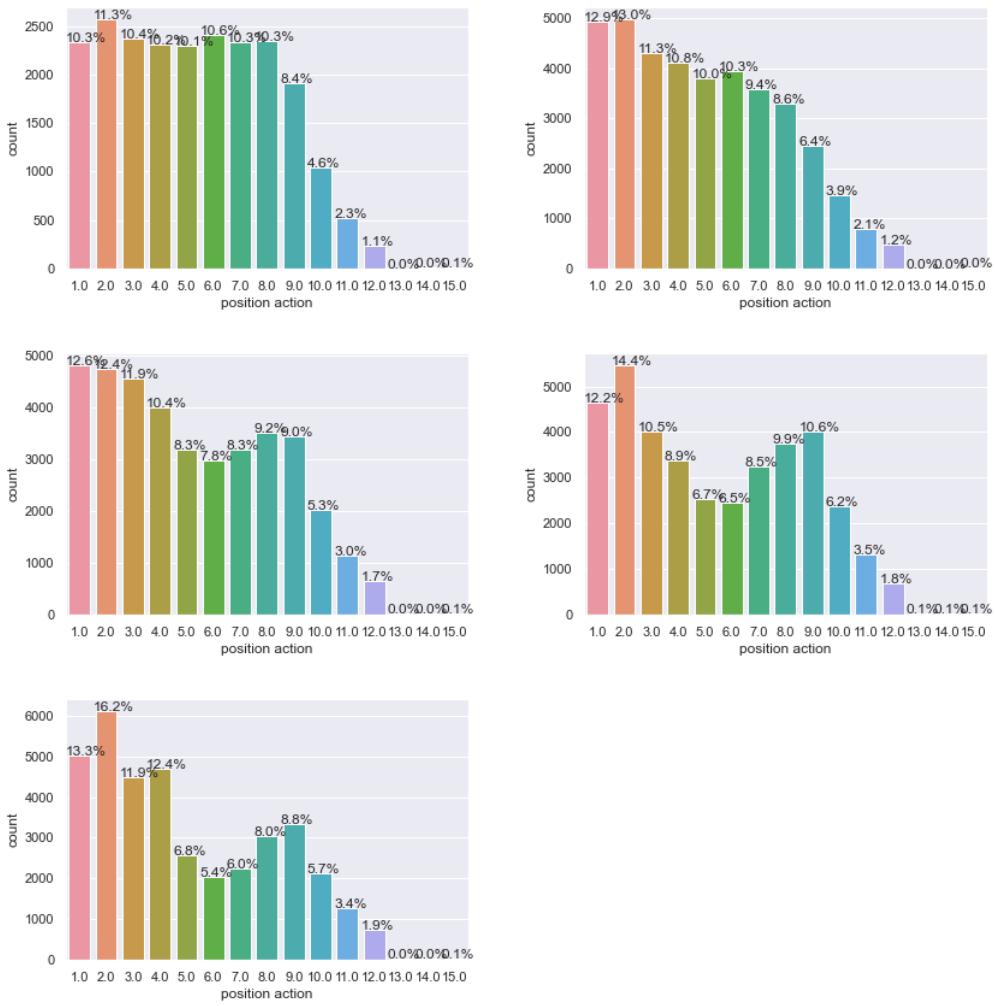


(10 iterations at 3/4)

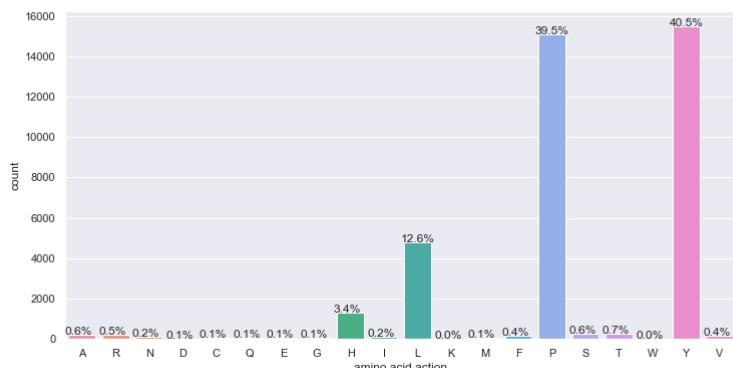
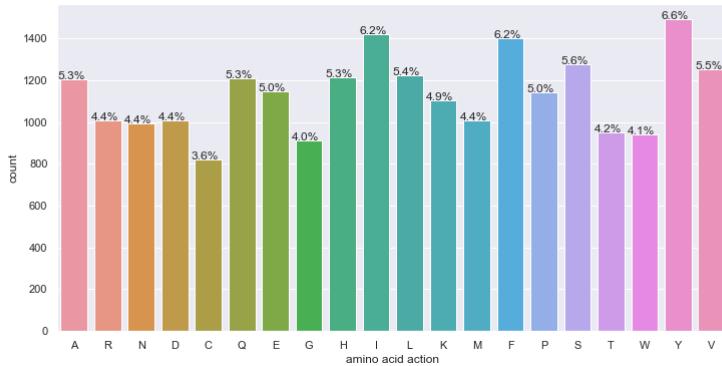
(last 10 iterations)

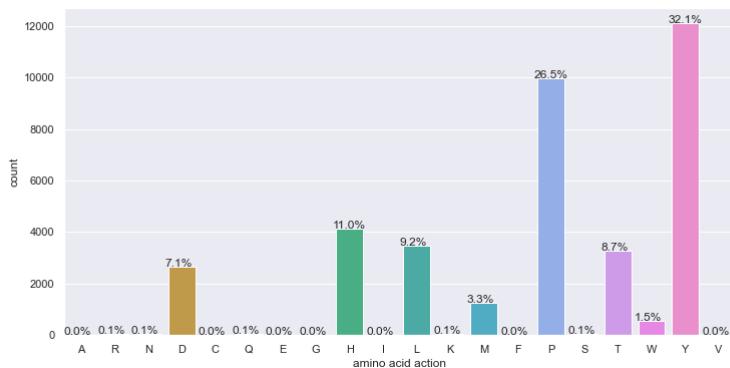
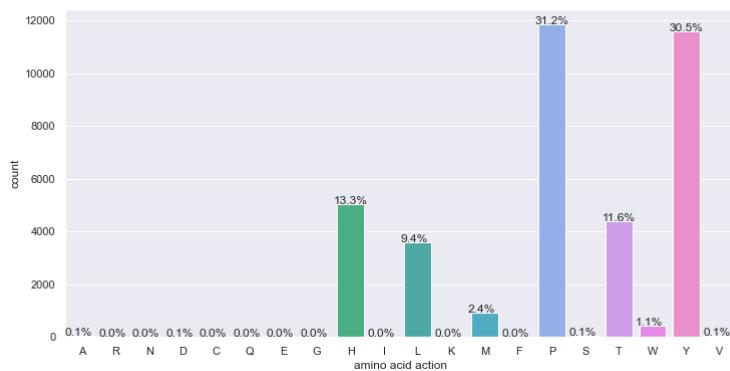
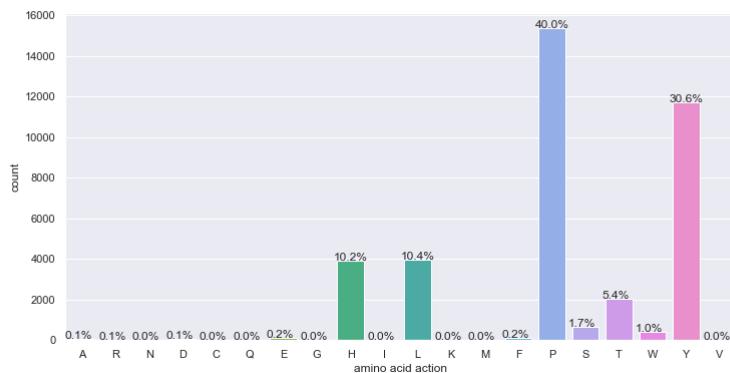


## 2. Position action analysis

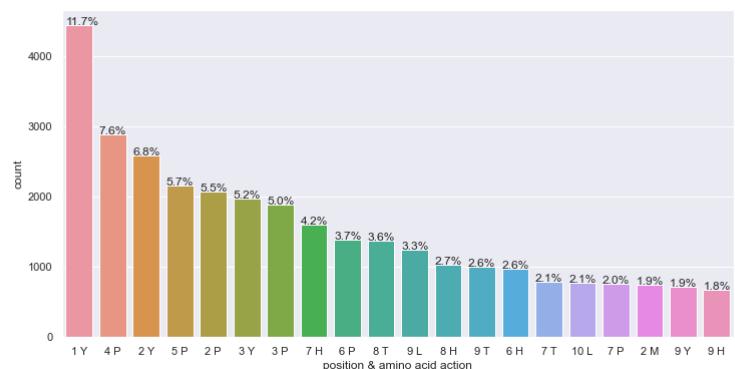
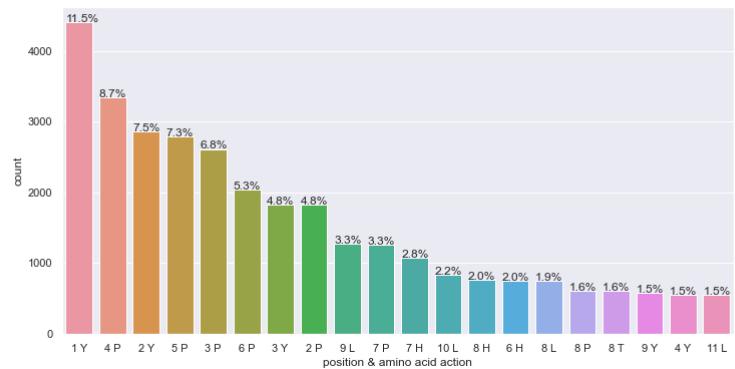
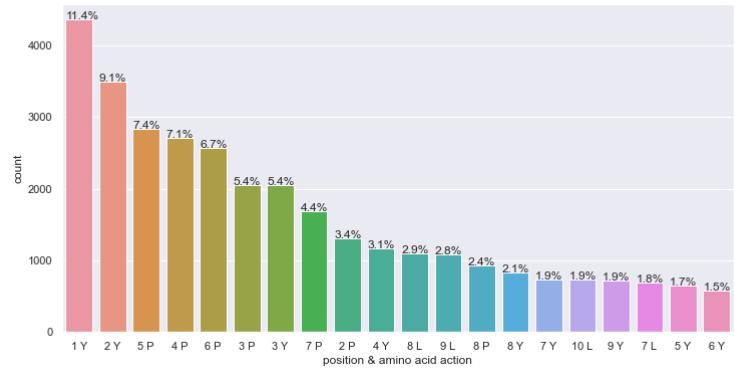
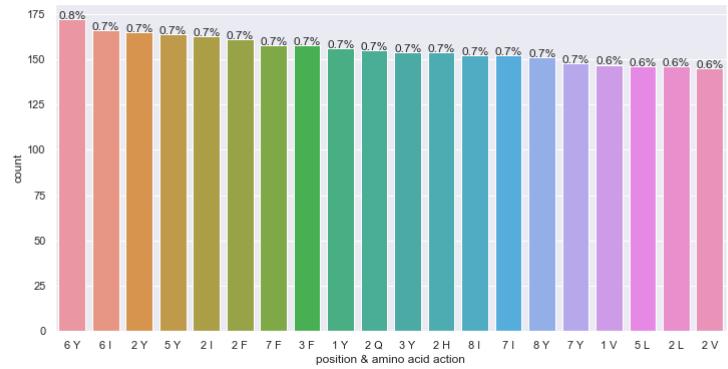


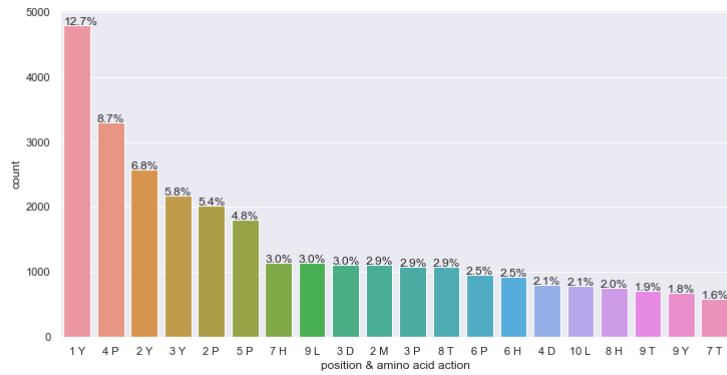
### 3. Amino action analysis



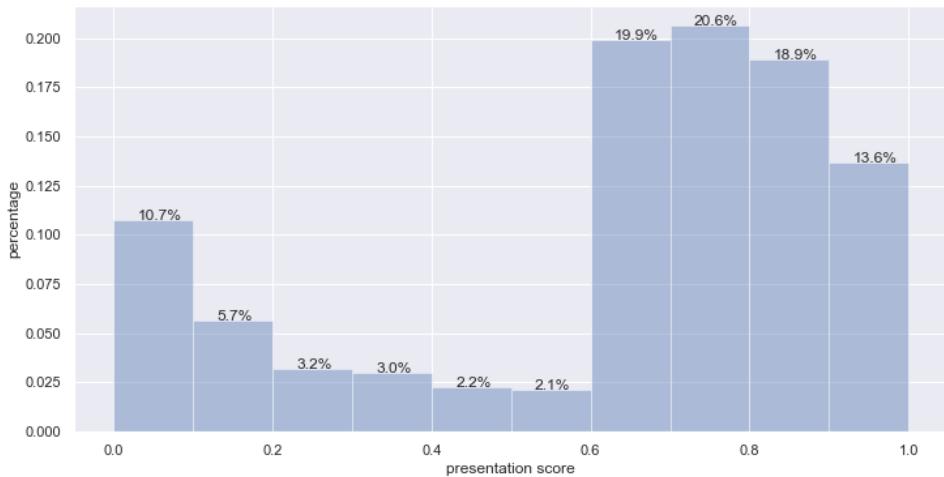


#### 4. Position & Amino action analysis

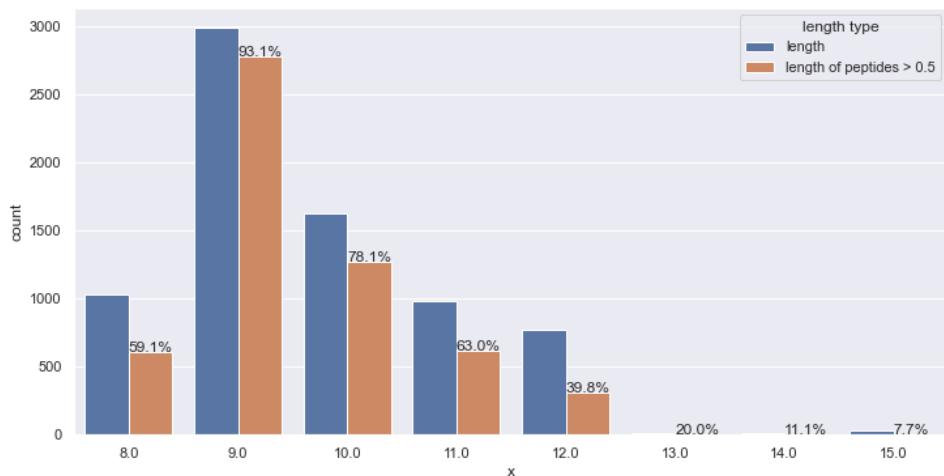




## 5. Reward analysis

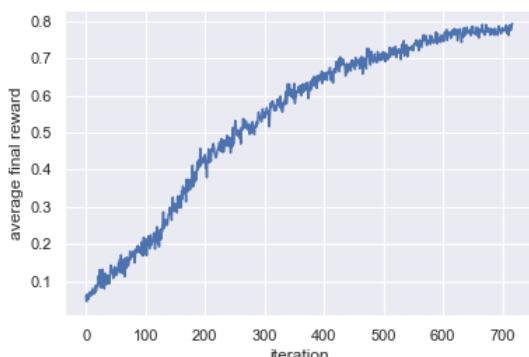
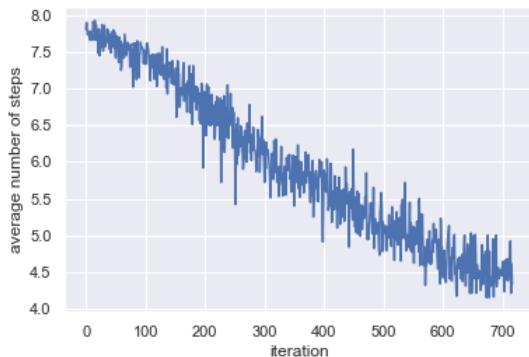
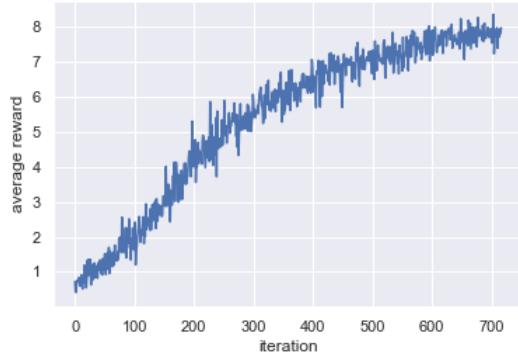


## 6. Length analysis



**Experiment: stop criteria = 0.75**

Setting 1 (only final rewards; sample rate = 50% from IEKB and 50% from random)

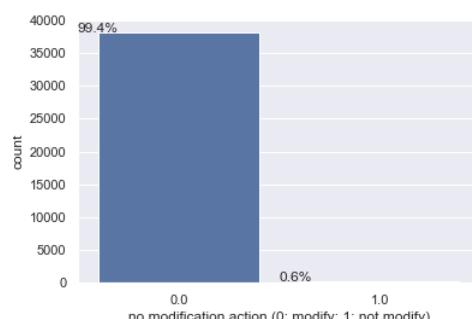
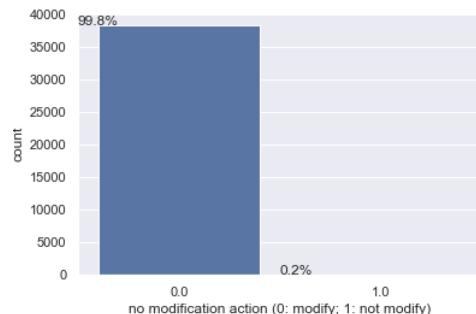
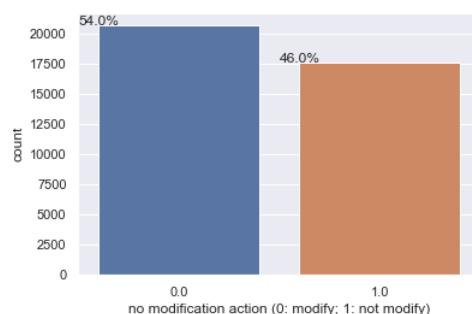


### 1. "No modification" action analysis

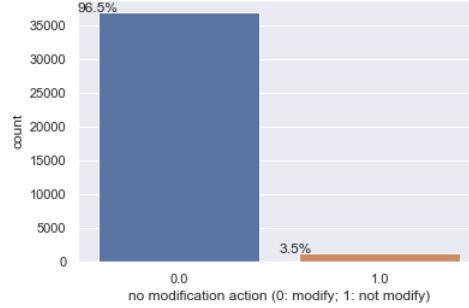
(first 10 iterations)  
the middle)

(10 iterations at 1/4)

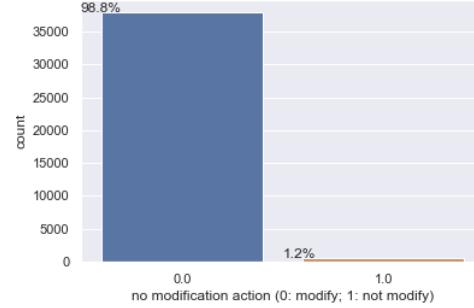
(10 iterations in



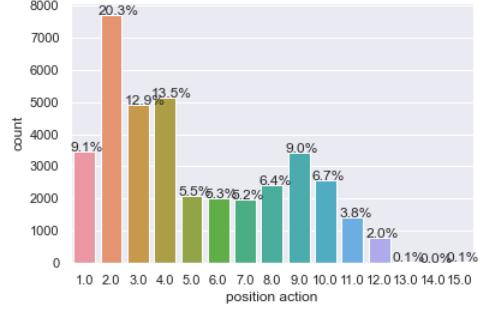
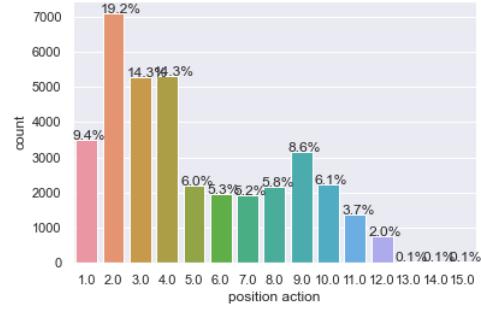
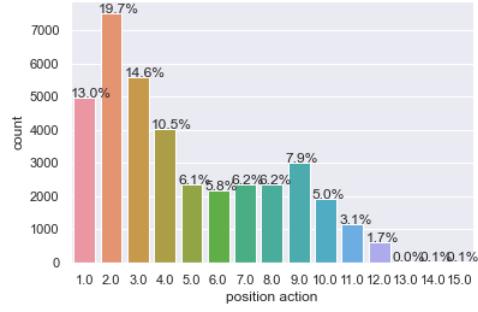
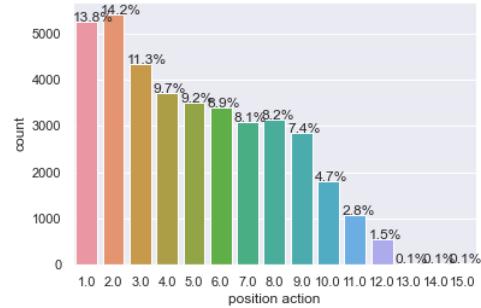
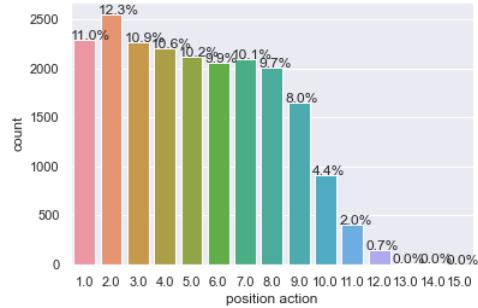
(10 iterations at 3/4)



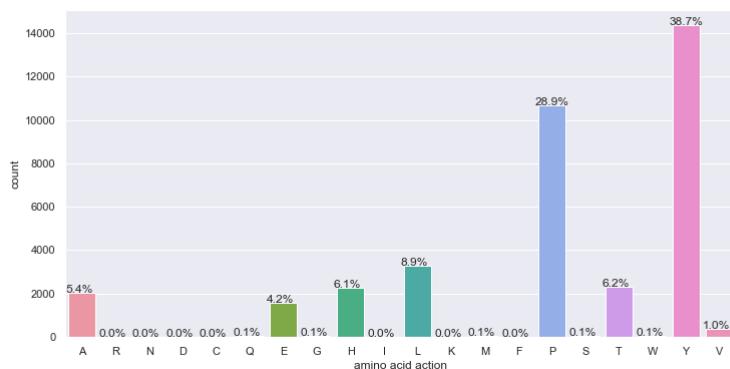
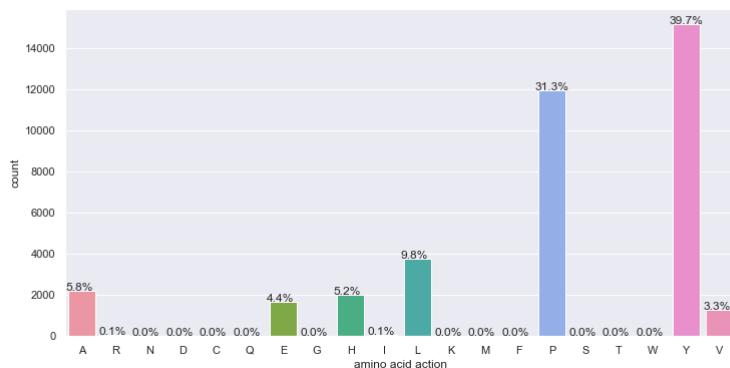
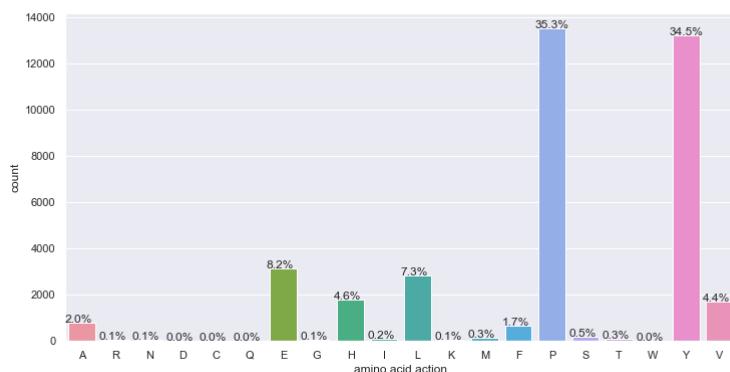
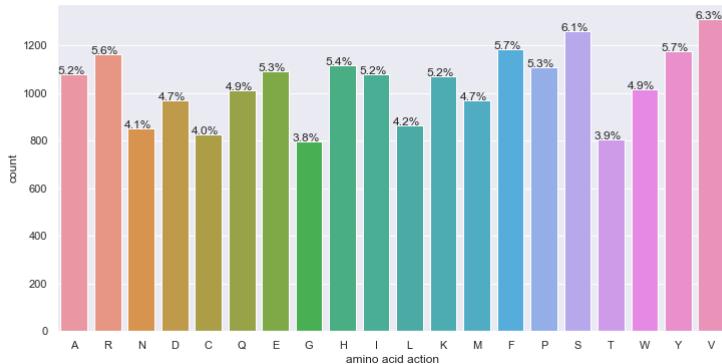
(last 10 iterations)

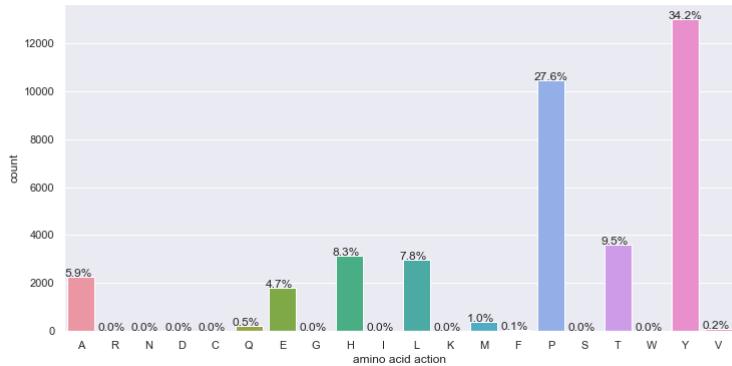


## 2. Position action analysis

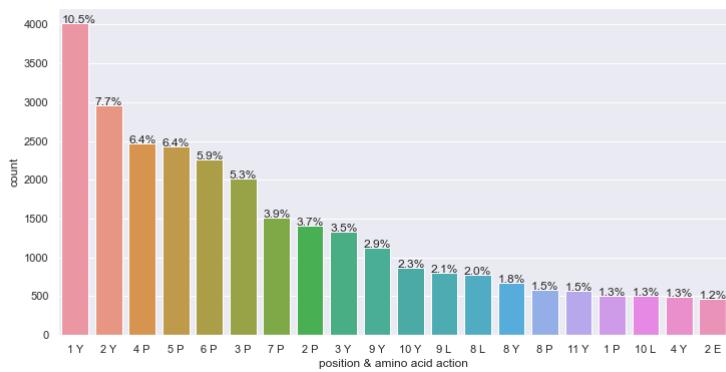
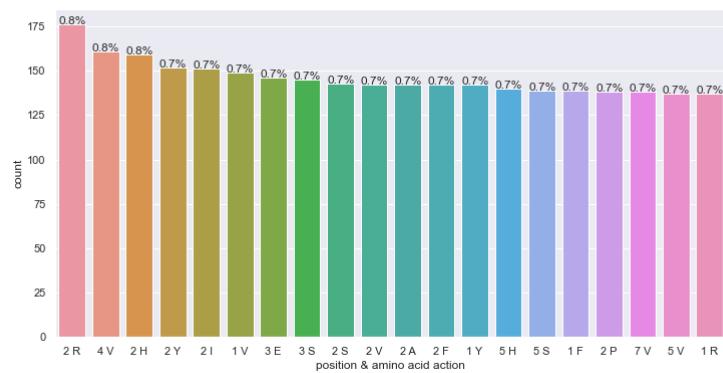


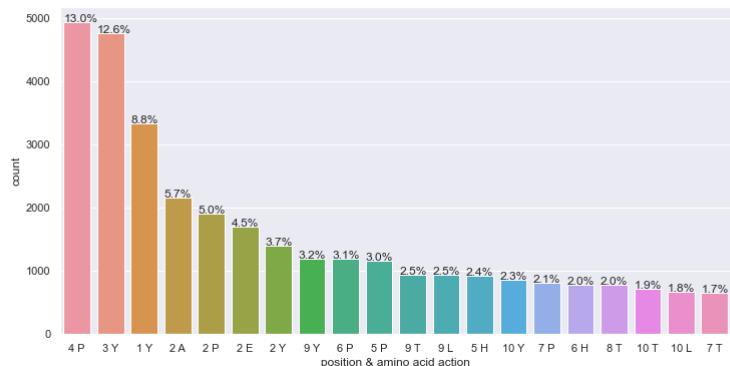
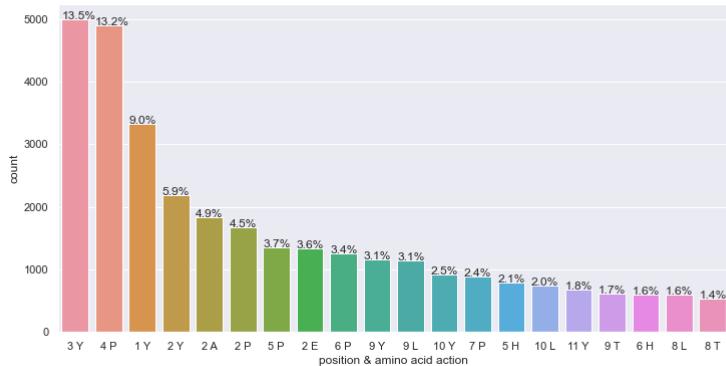
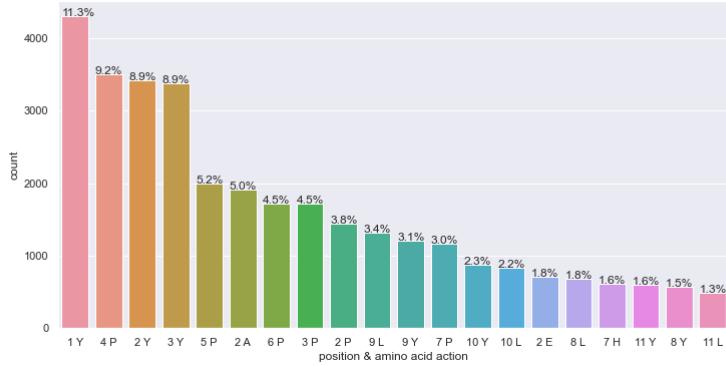
## 3. Amino action analysis



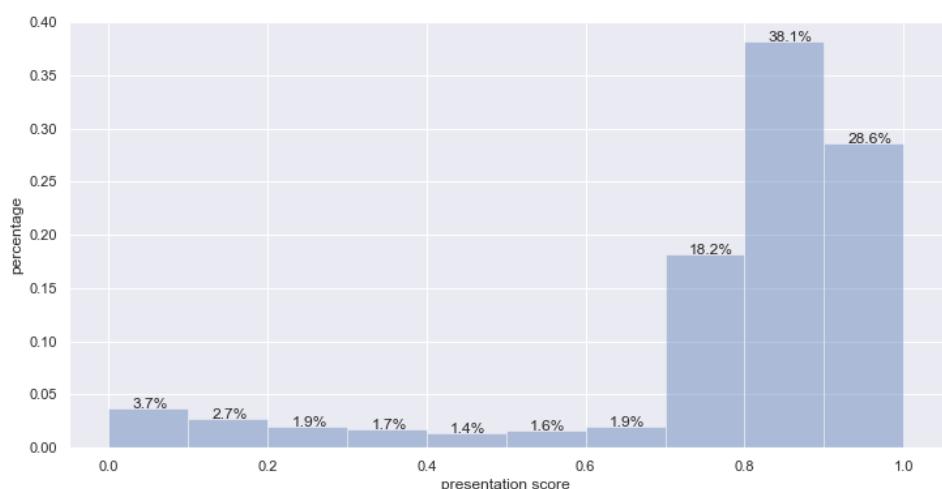


#### 4. Position & Amino action analysis

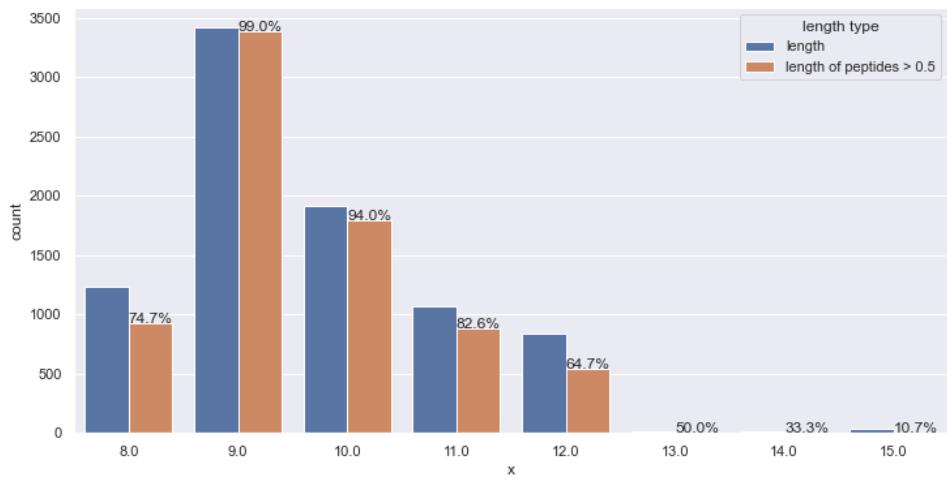




## 5. Reward analysis

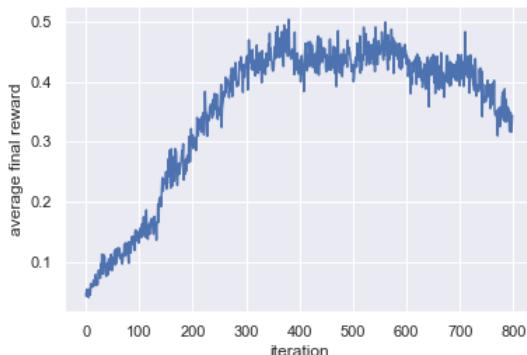
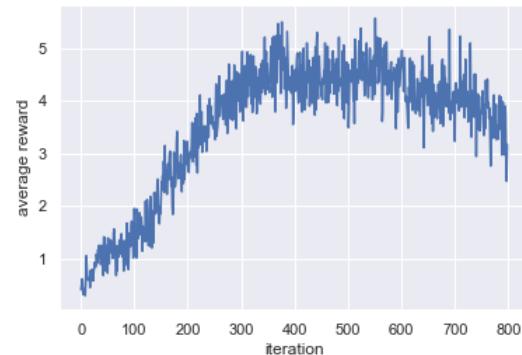


## 6. Length analysis



### Experiment: stop criteria = 1.0

Setting 1 (only final rewards; sample rate = 50% from IEDB and 50% from random)

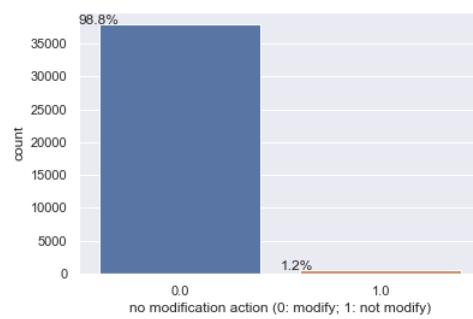
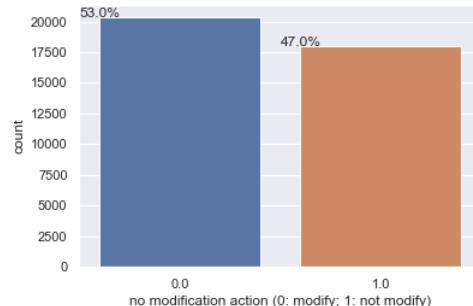


#### 1. "No modification" action analysis

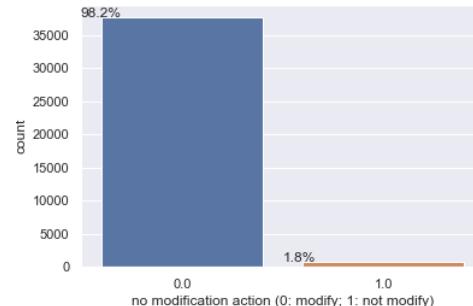
(first 10 iterations)  
the middle)

(10 iterations at 1/4)

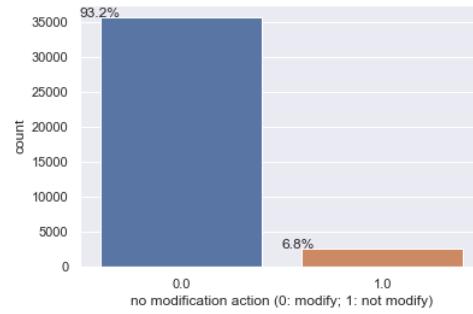
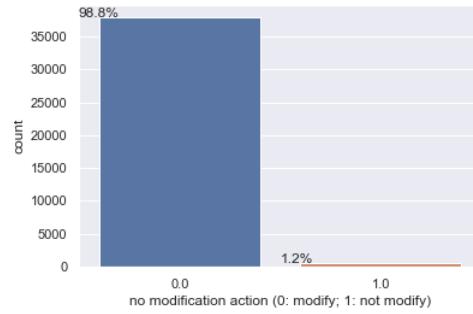
(10 iterations in



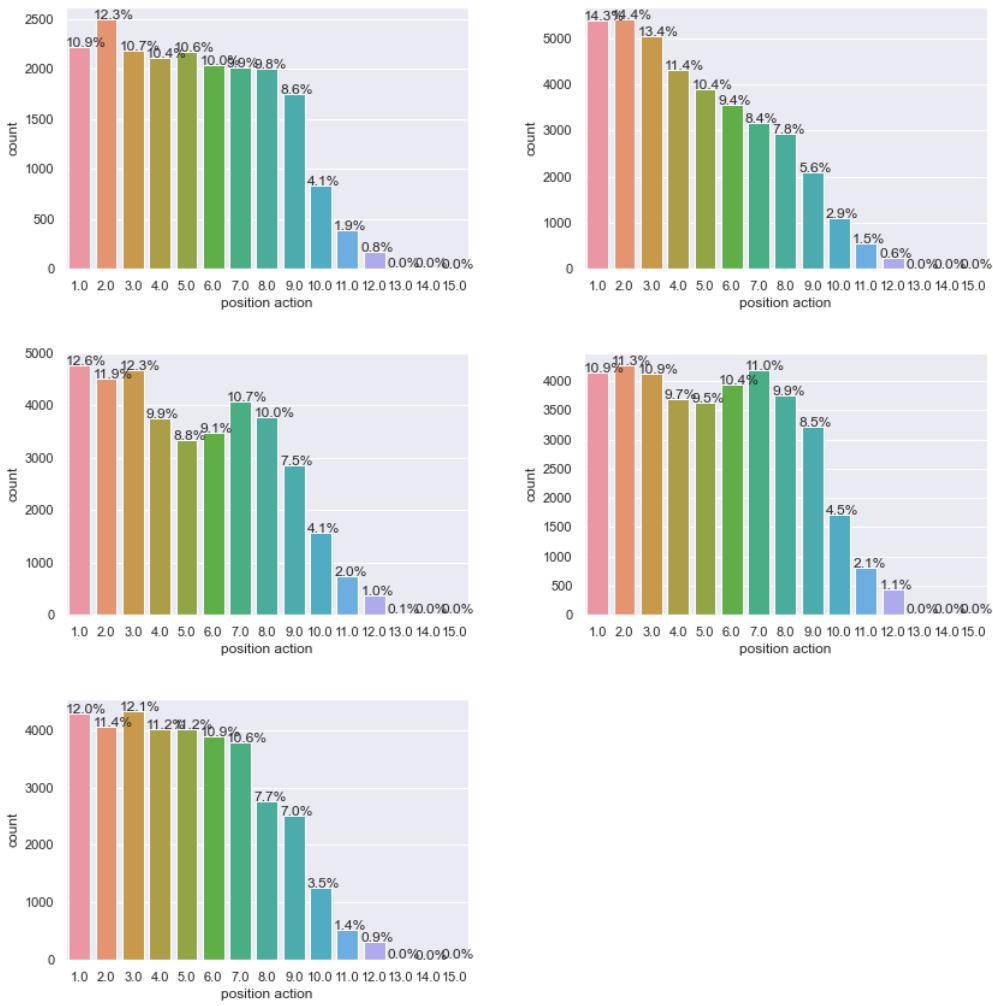
(10 iterations at 3/4)



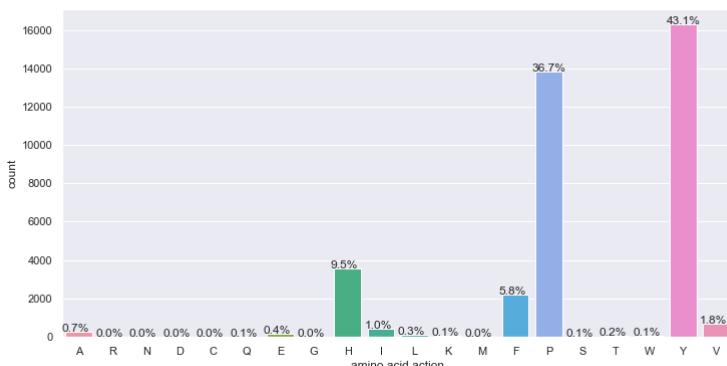
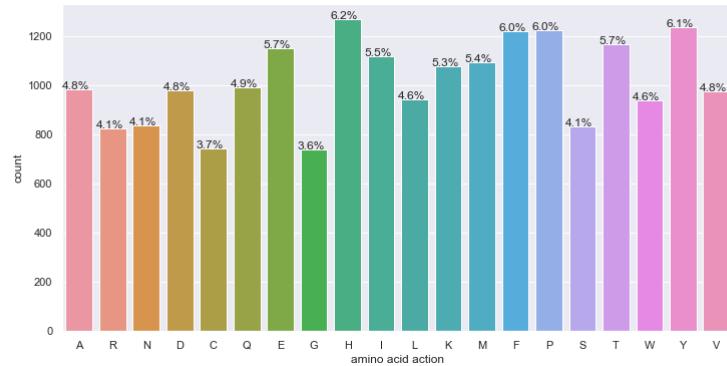
(last 10 iterations)

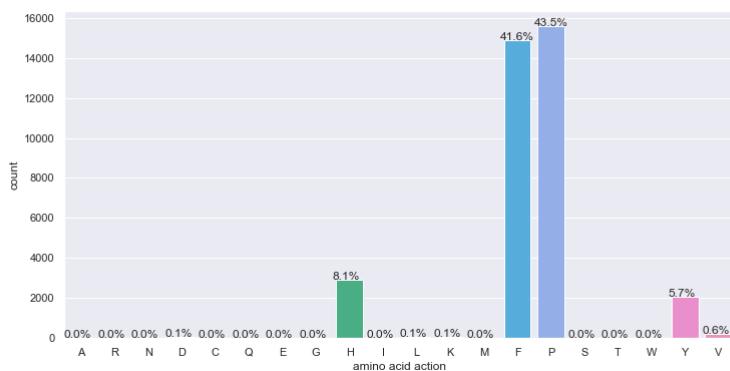
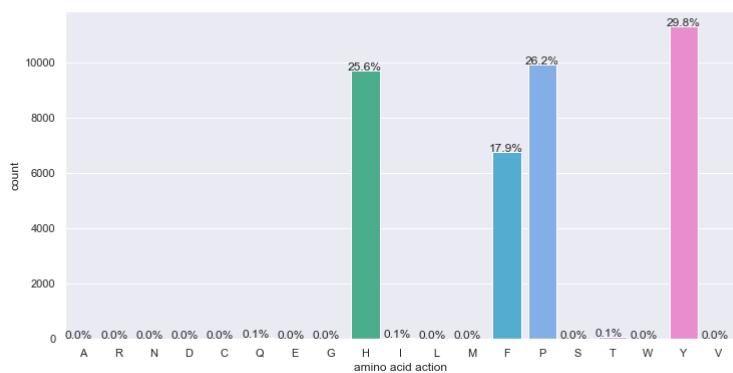
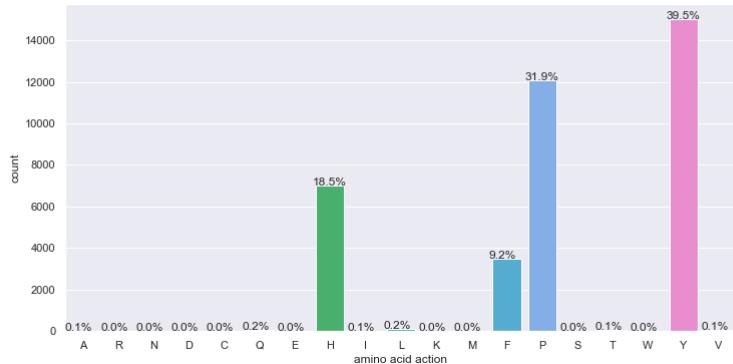


## 2. Position action analysis

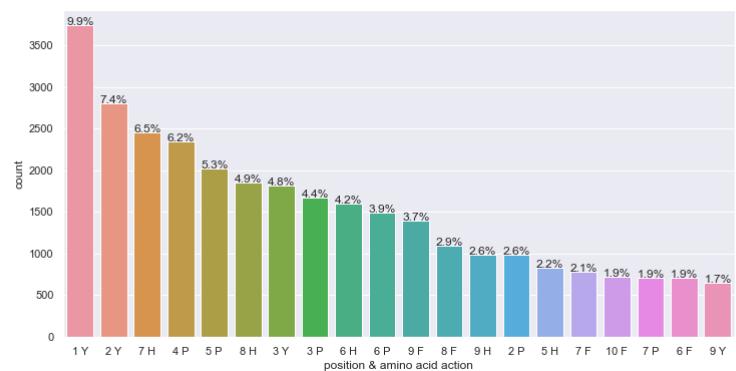
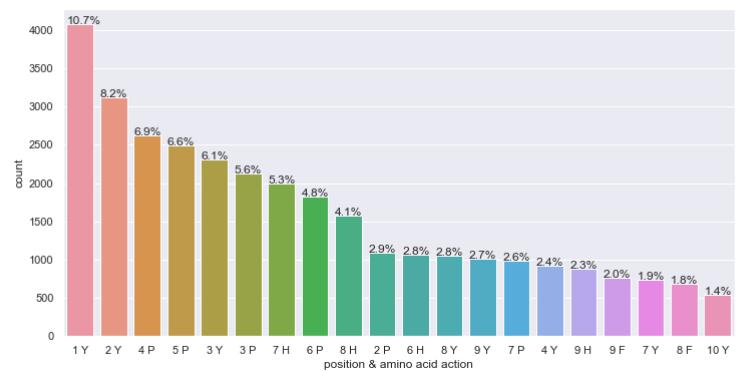
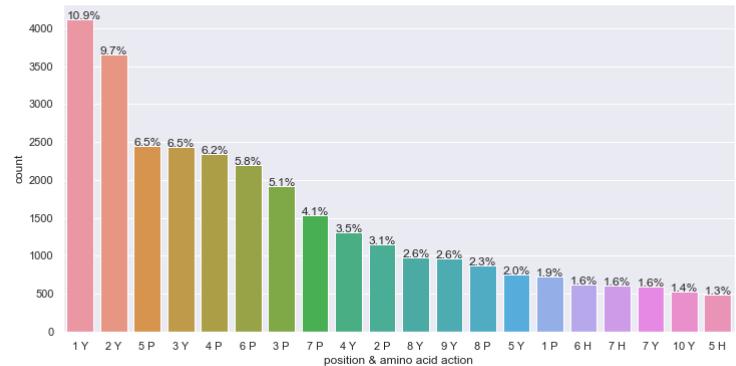
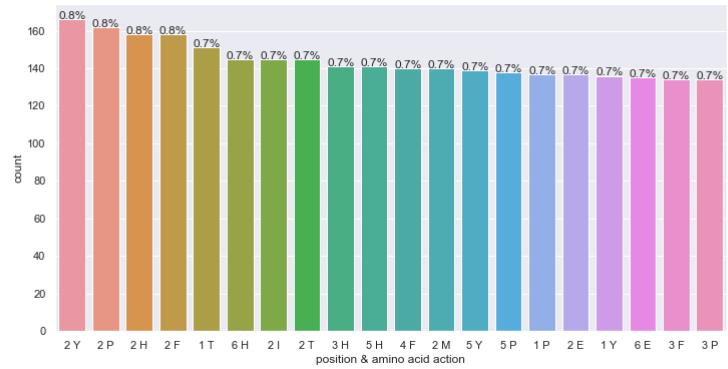


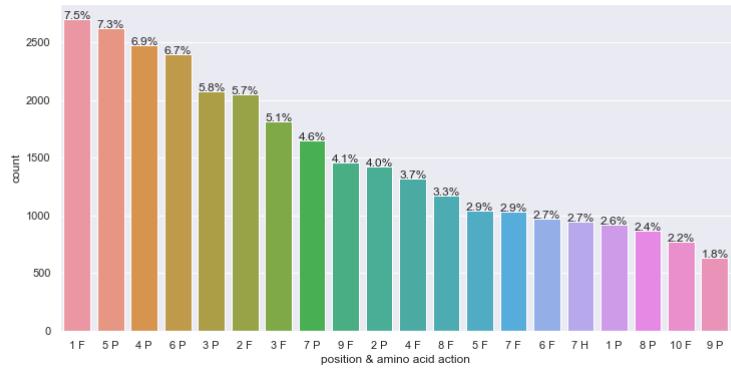
### 3. Amino action analysis



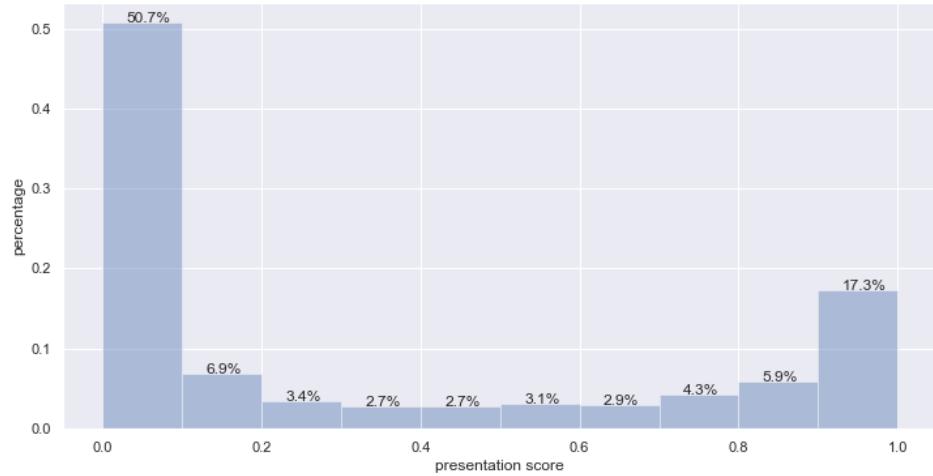


#### 4. Position & Amino action analysis

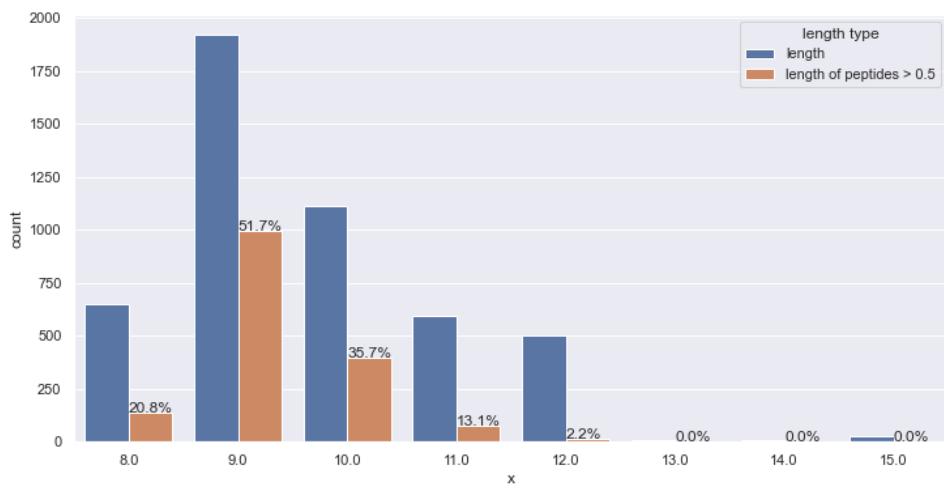




## 5. Reward analysis



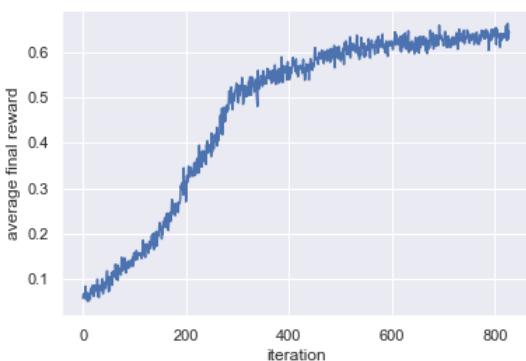
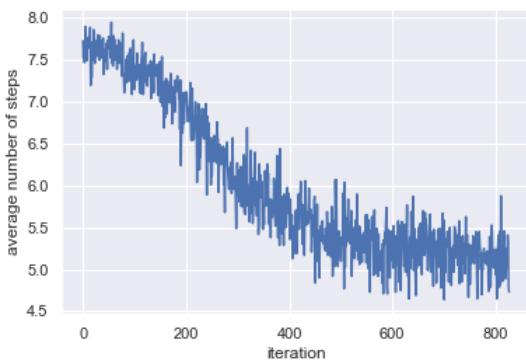
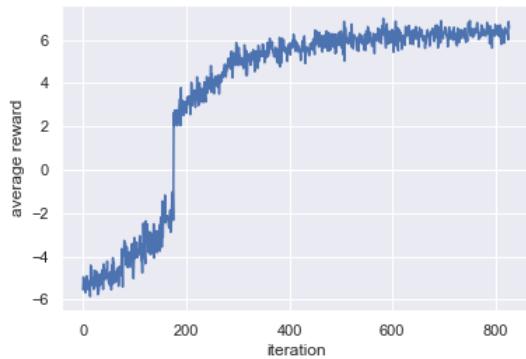
## 6. Length analysis



## 3. Result Analysis with intermediate reward at first and then only final rewards

**Experiment: stop criteria = 0.6**

Setting 1 (sample rate = 50% from IEDB and 50% from random)



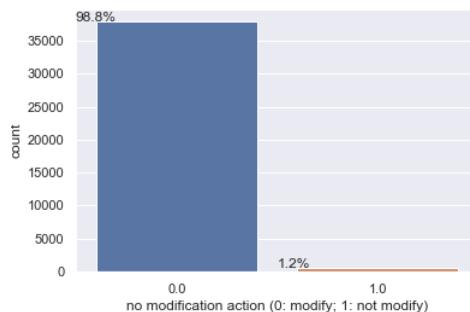
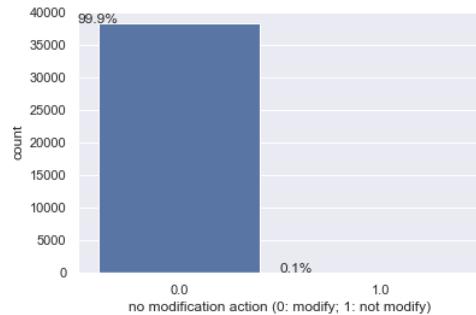
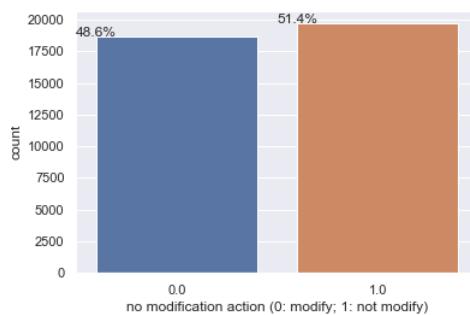
### 1. "No modification" action analysis

(first 10 iterations)

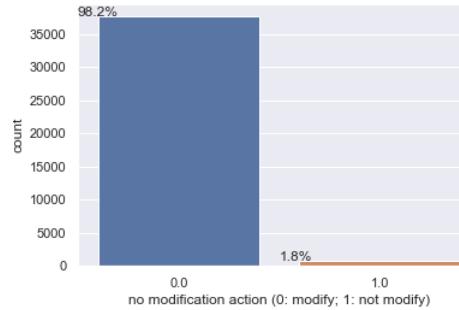
the middle)

(10 iterations at 1/4)

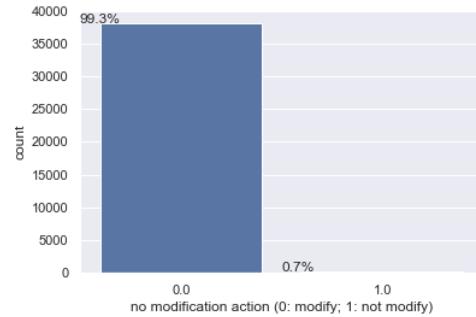
(10 iterations in



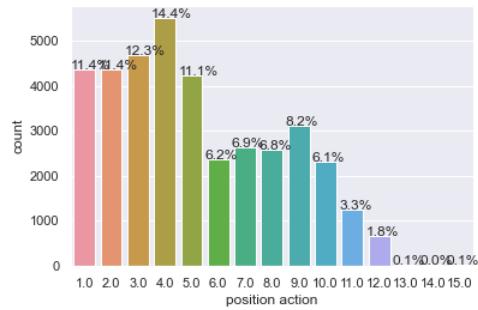
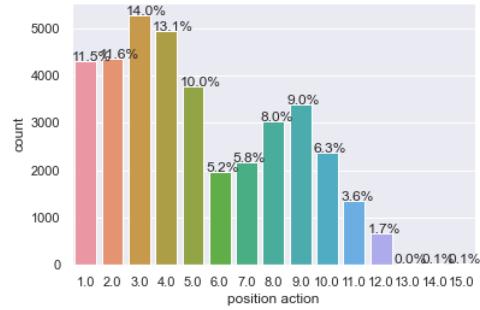
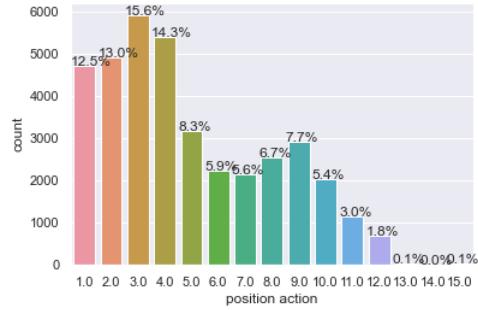
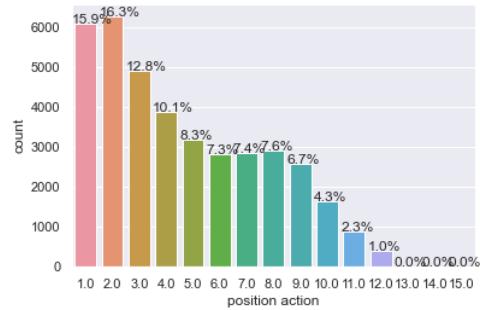
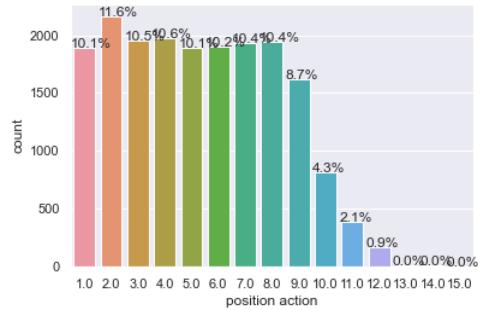
(10 iterations at 3/4)



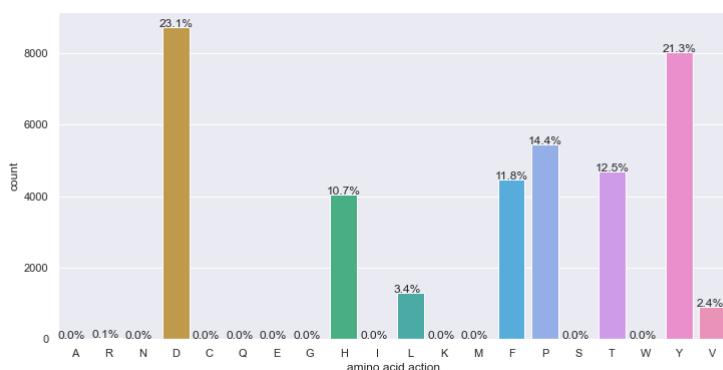
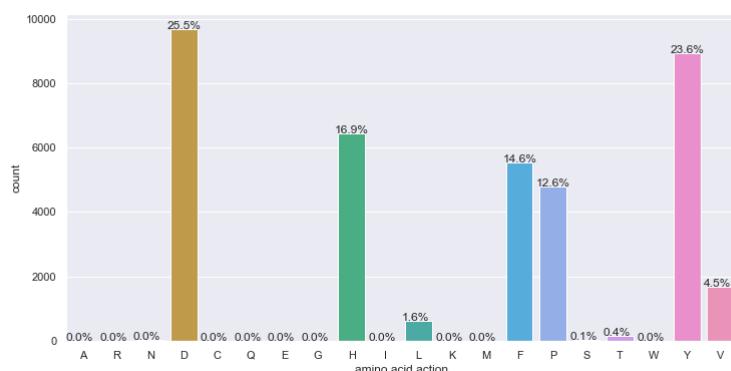
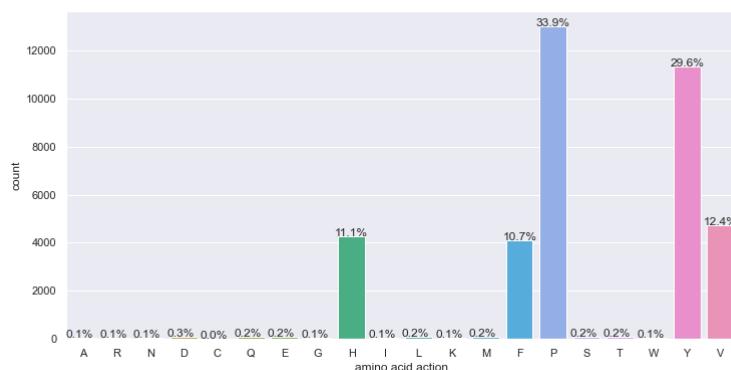
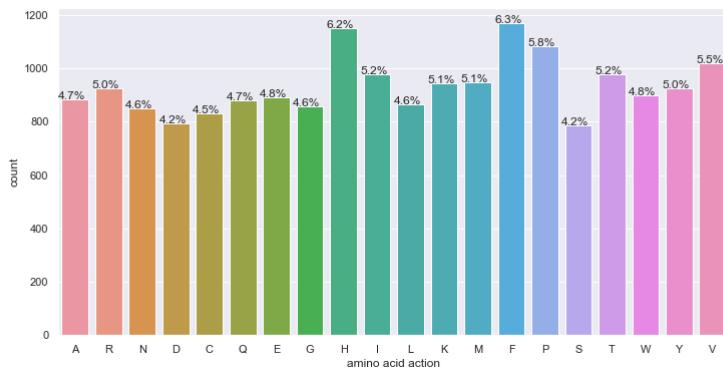
(last 10 iterations)

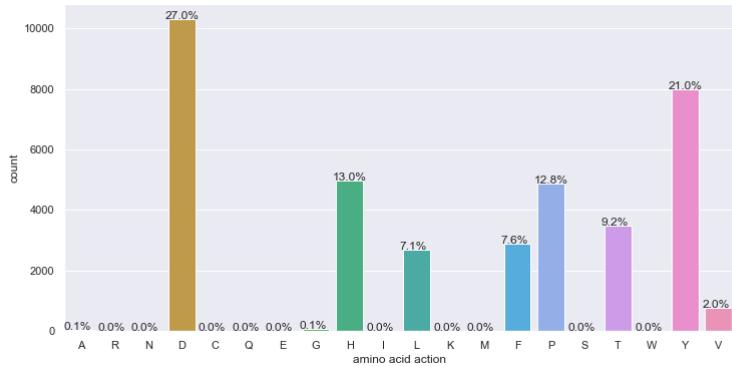


## 2. Position action analysis

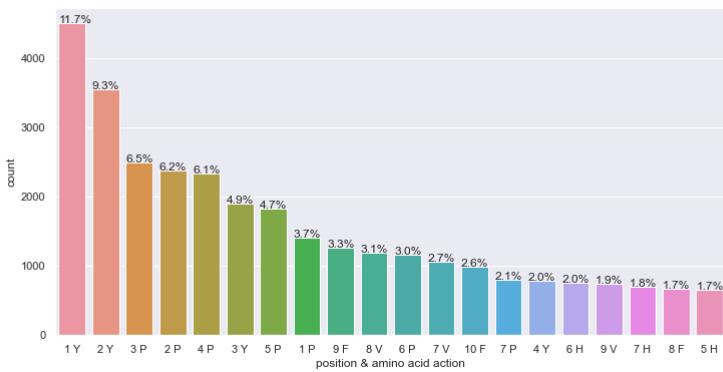
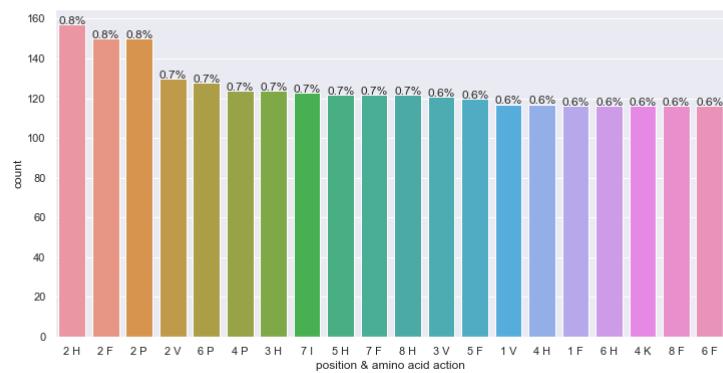


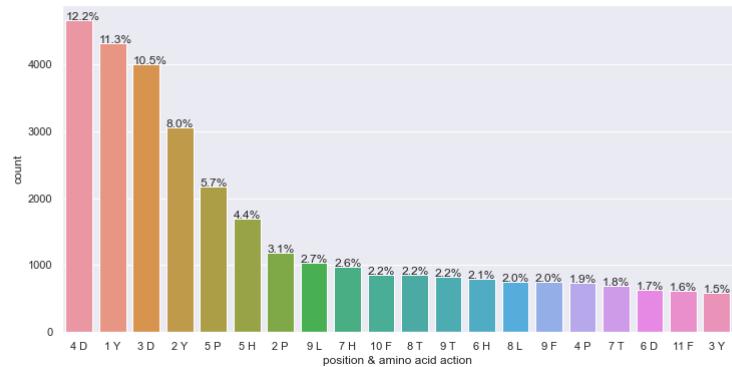
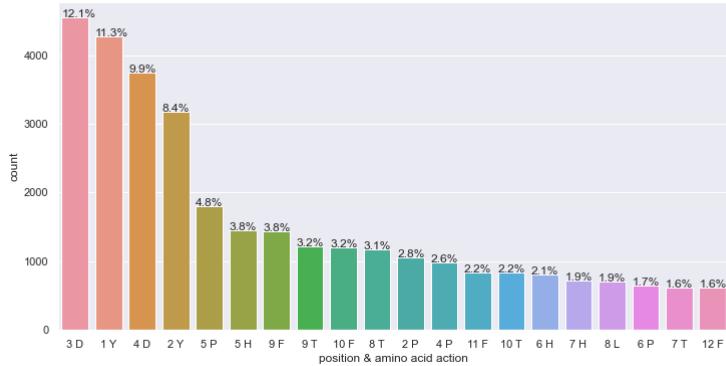
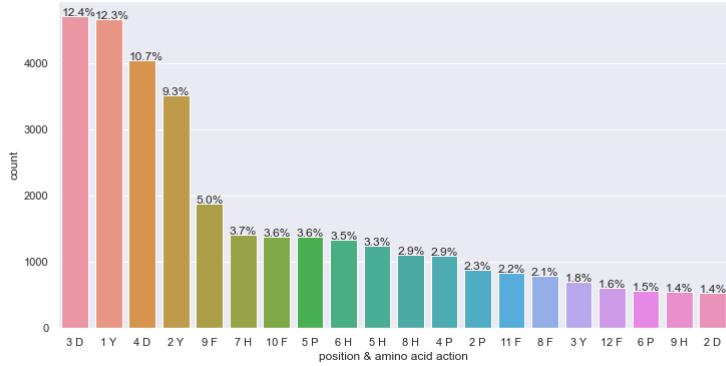
## 3. Amino action analysis



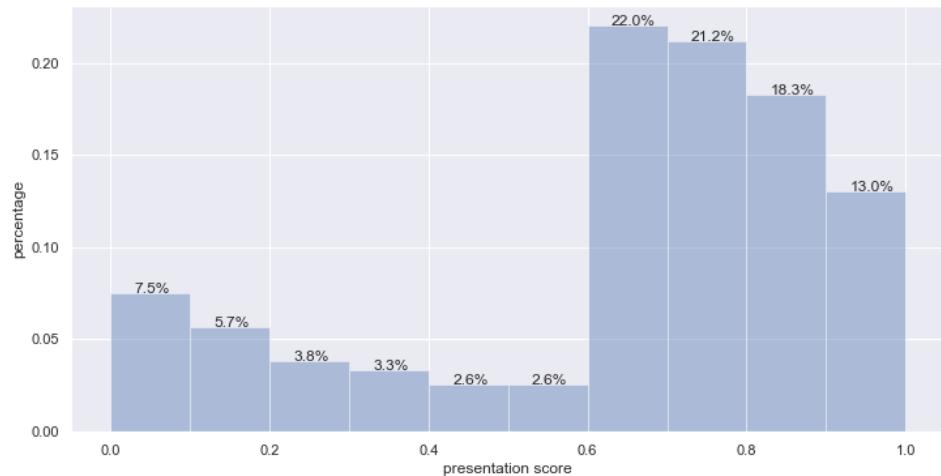


#### 4. Position & Amino action analysis

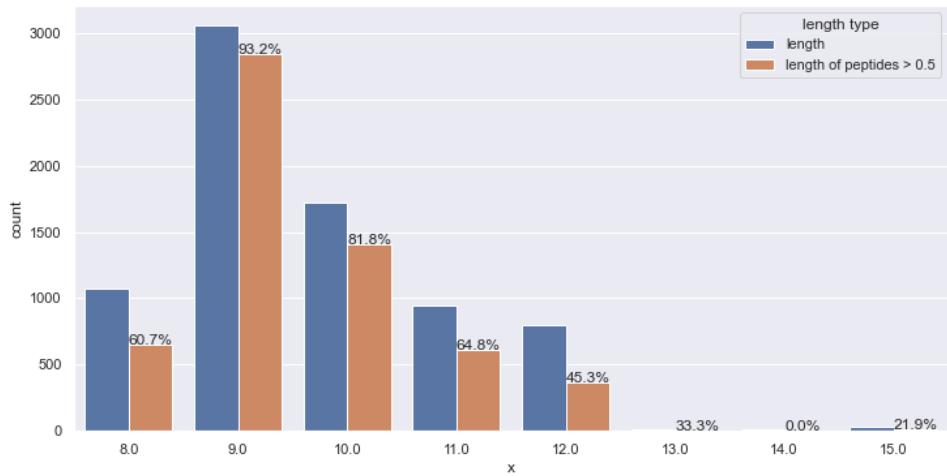




## 5. Reward analysis

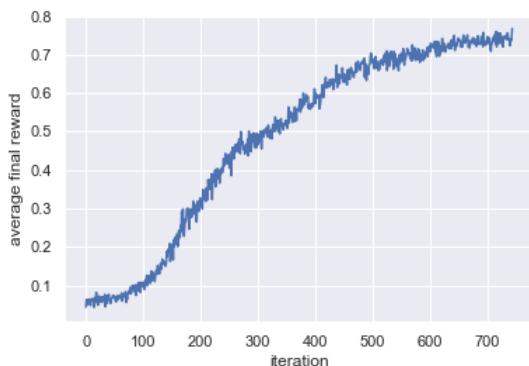
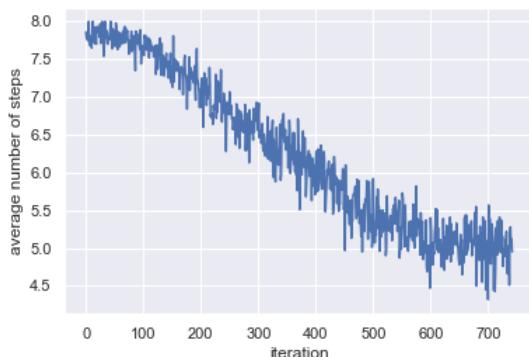
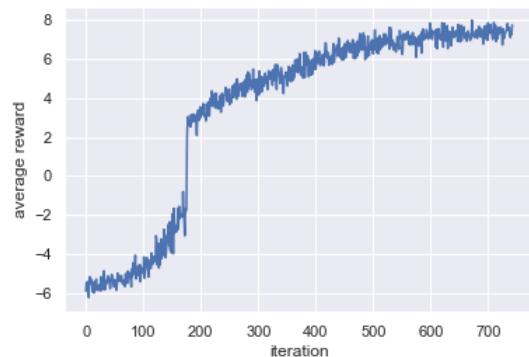


## 6. Length analysis



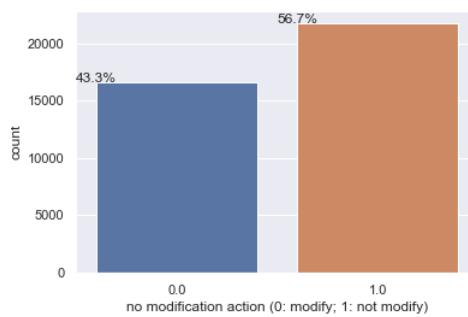
### Experiment: stop criteria = 0.75

Setting 1 (sample rate = 50% from IEDB and 50% from random)

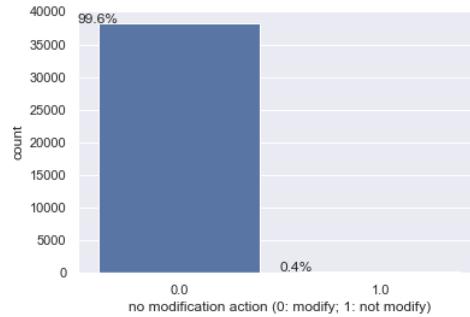


1. "No modification" action analysis

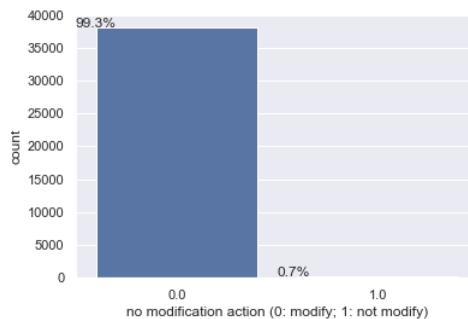
(first 10 iterations)  
the middle)



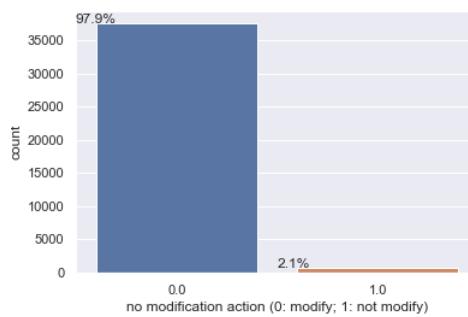
(10 iterations at 1/4)



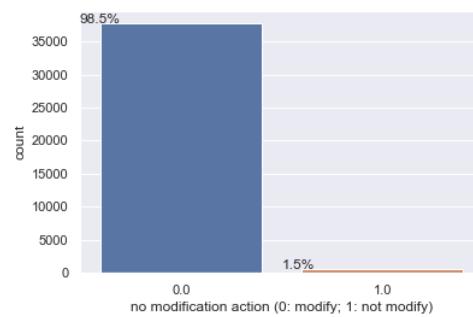
(10 iterations in



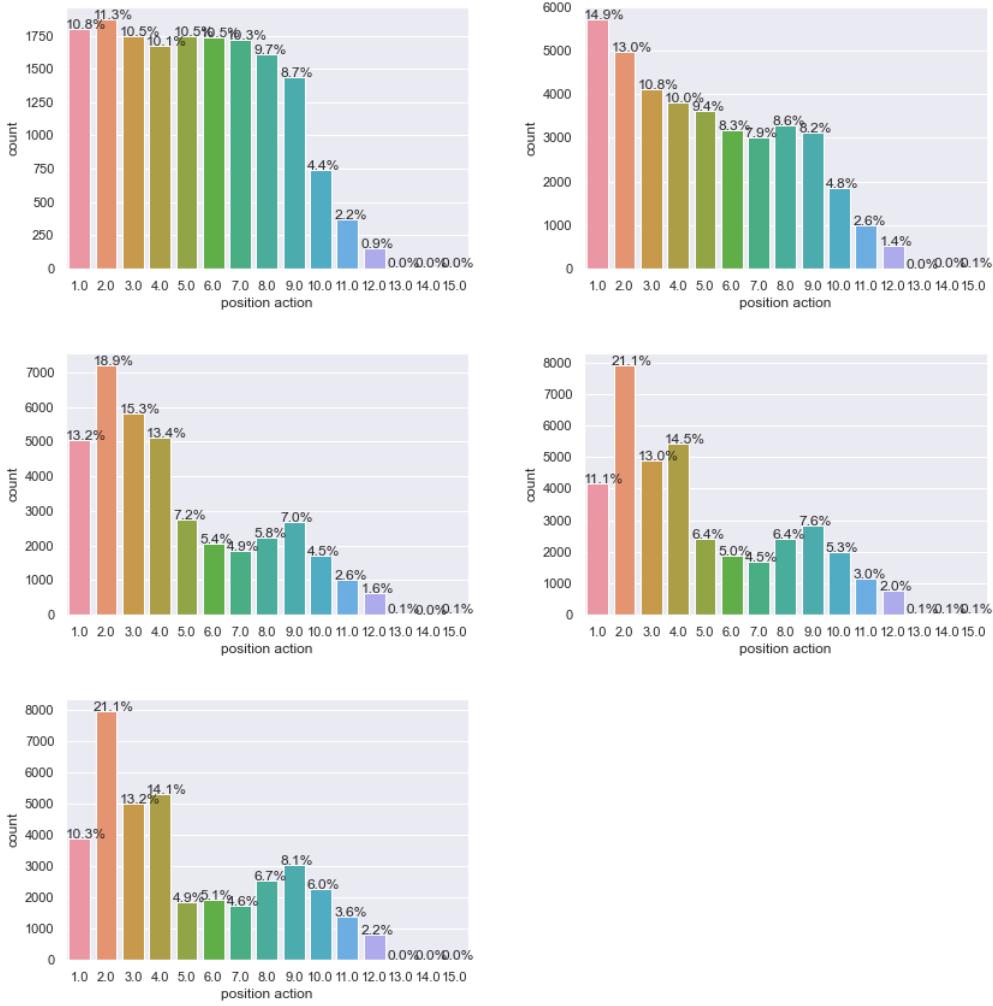
(10 iterations at 3/4)



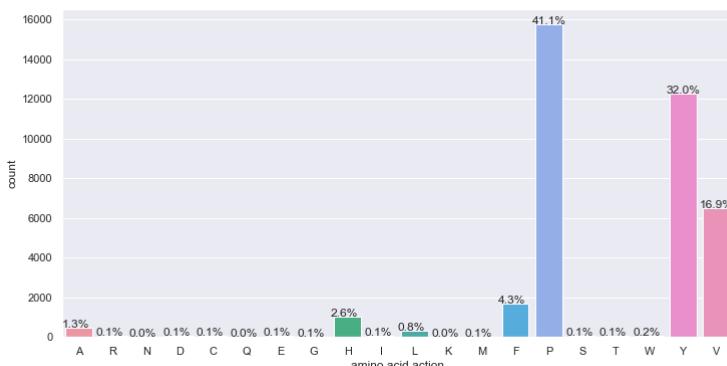
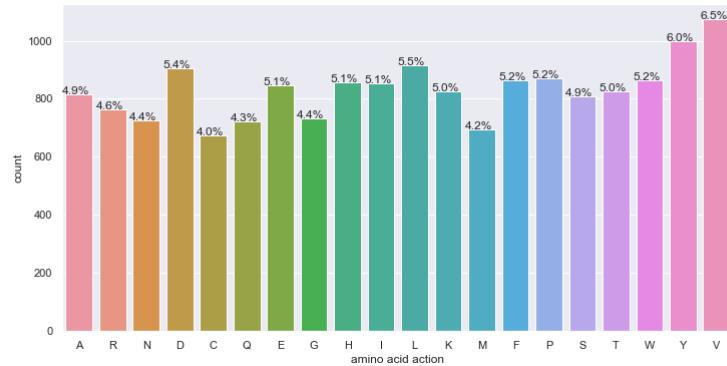
(last 10 iterations)

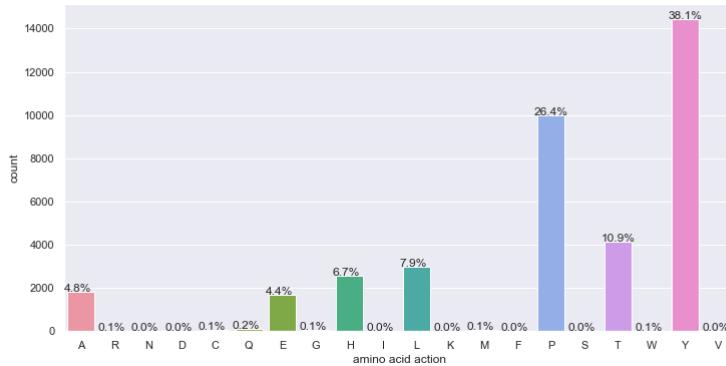
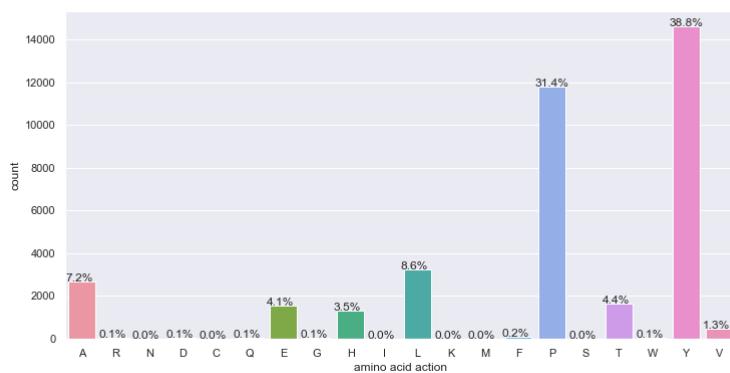
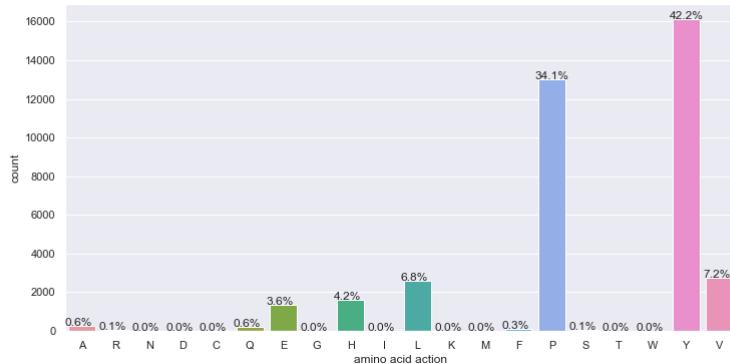


## 2. Position action analysis

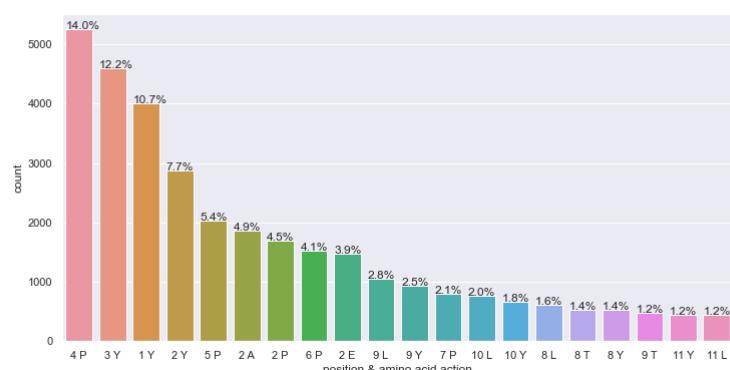
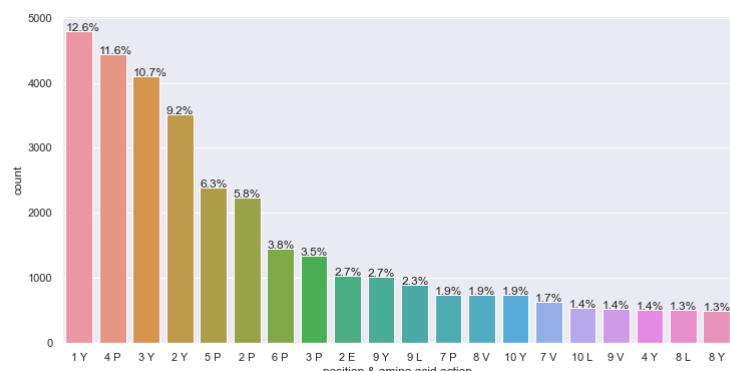
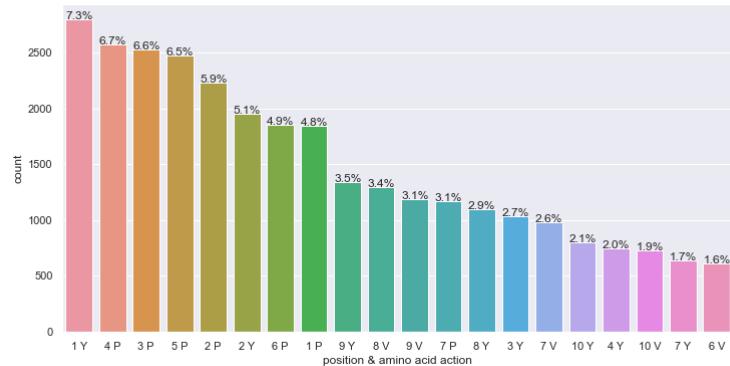
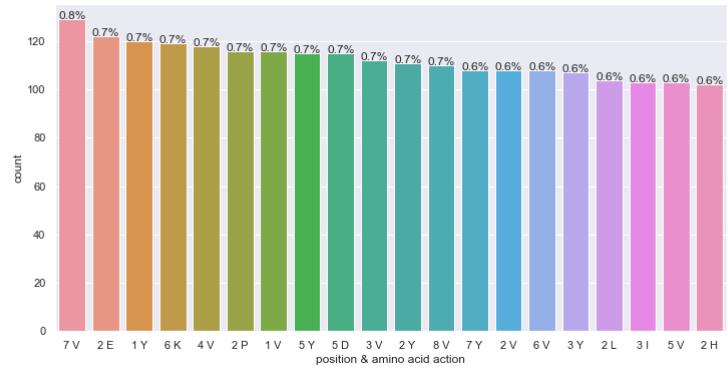


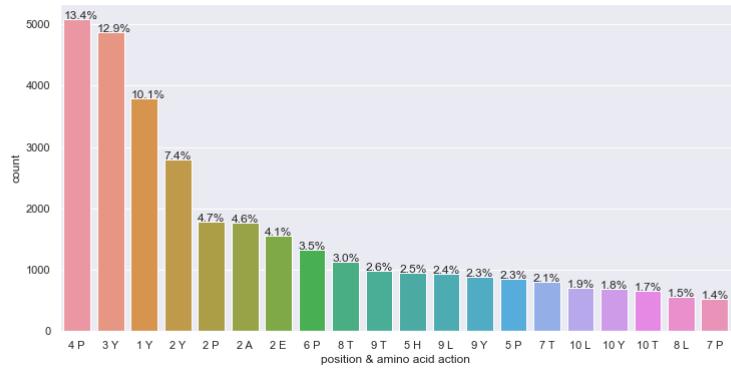
### 3. Amino action analysis



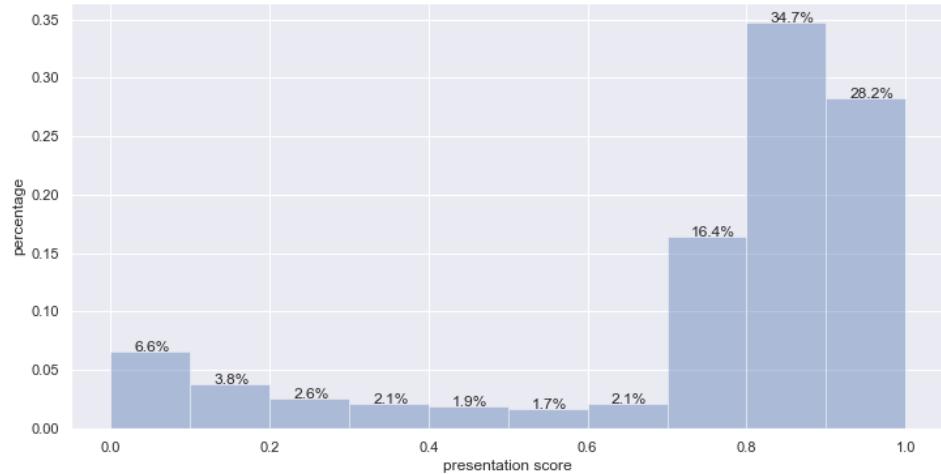


#### 4. Position & Amino action analysis

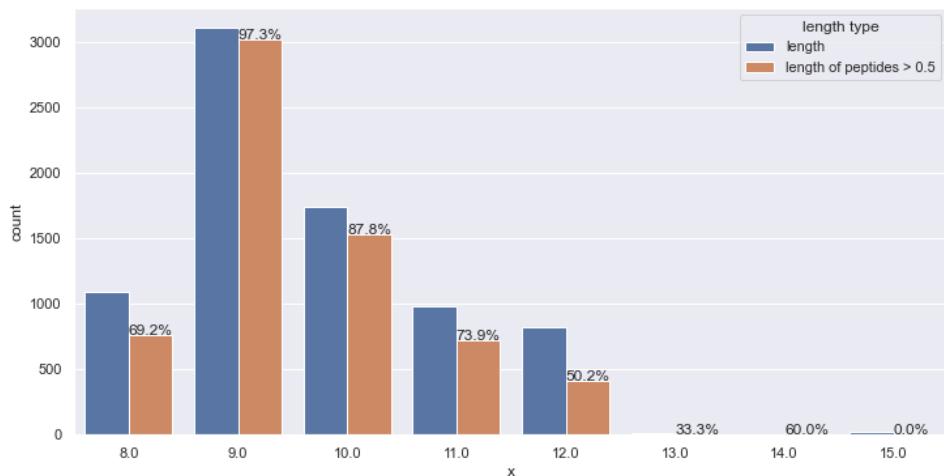




## 5. Reward analysis

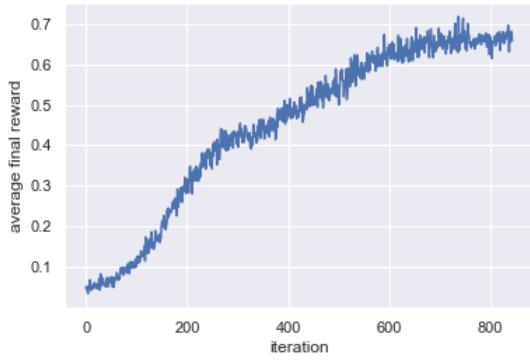
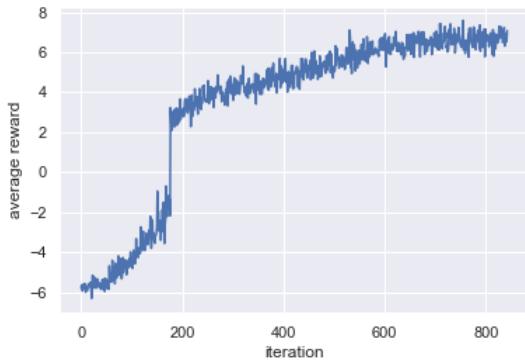


## 6. Length analysis



**Experiment: stop criteria = 1.0**

Setting 1 (sample rate = 50% from IEDB and 50% from random)

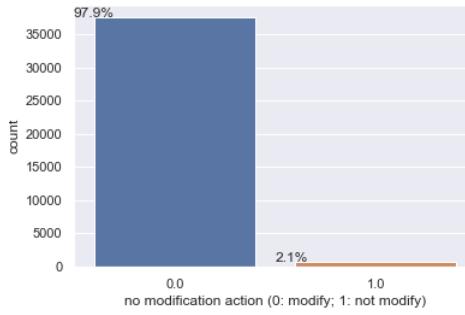
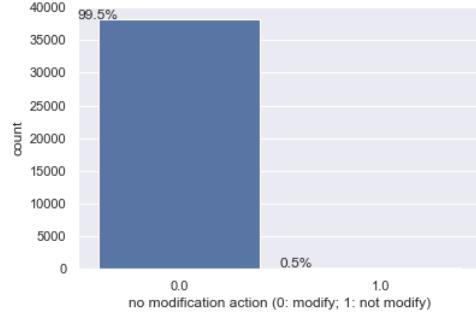
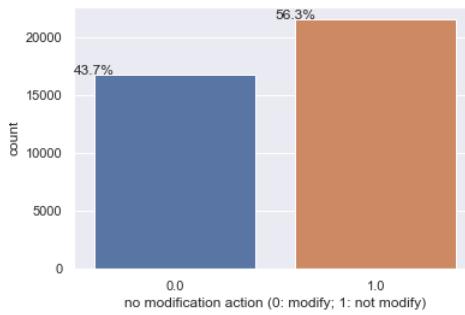


### 1. "No modification" action analysis

(first 10 iterations)  
the middle)

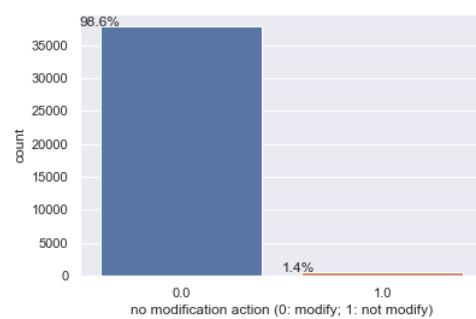
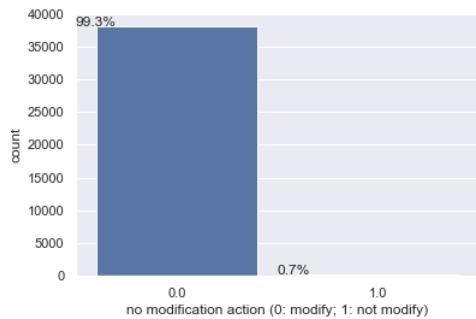
(10 iterations at 1/4)

(10 iterations in

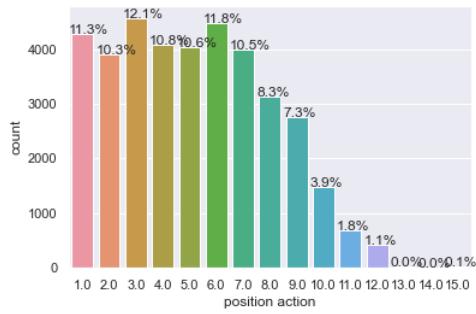
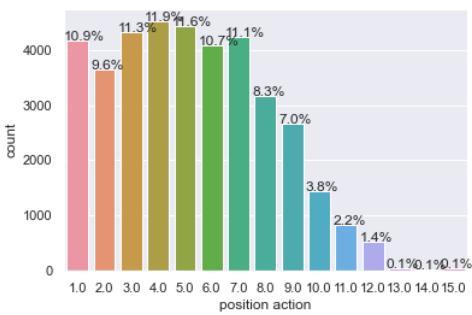
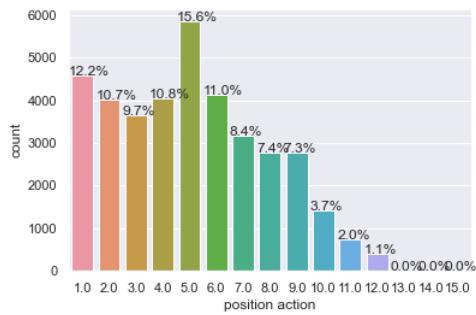
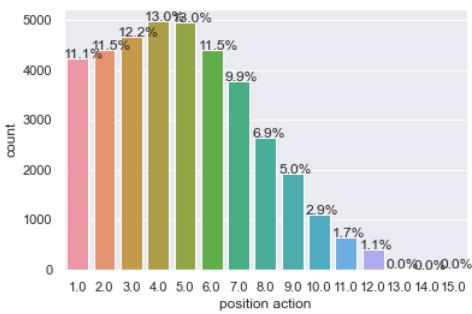
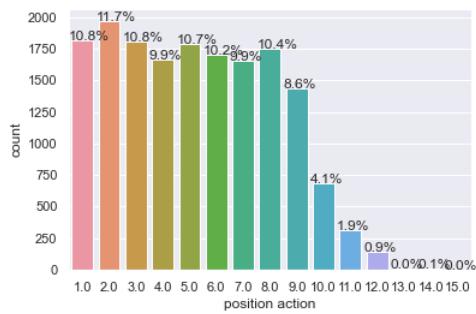


(10 iterations at 3/4)

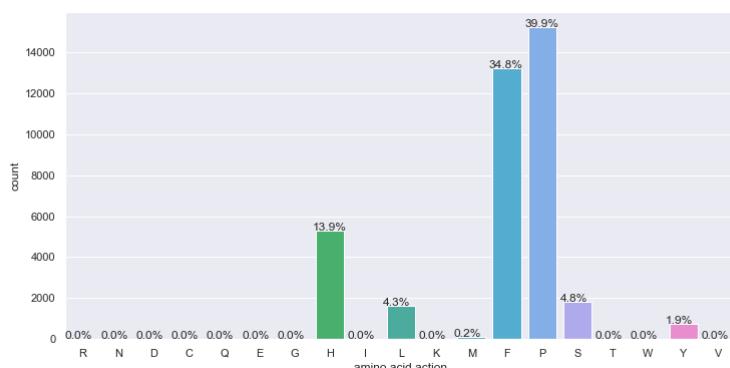
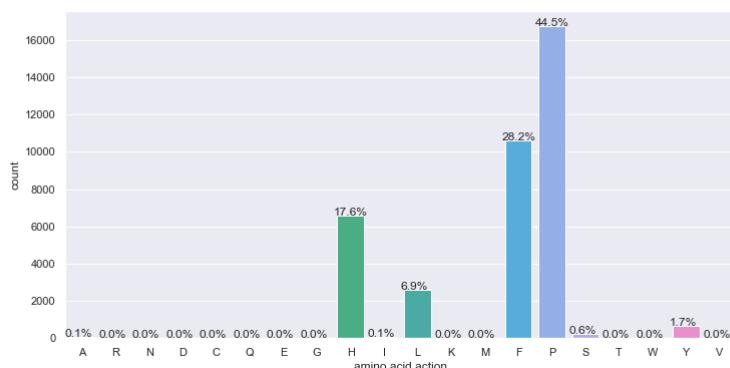
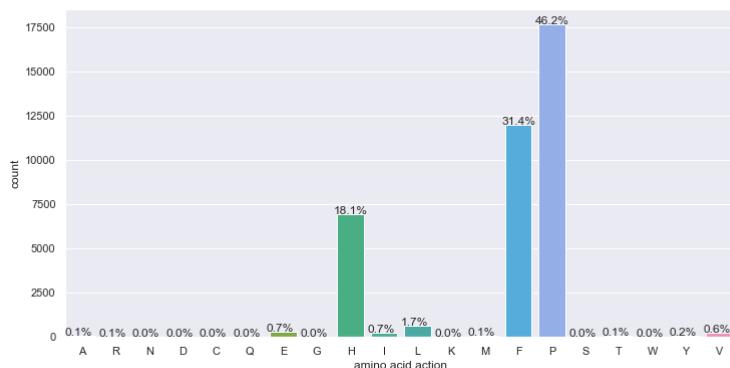
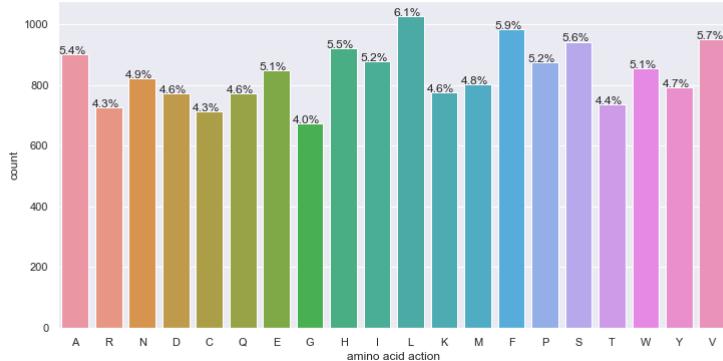
(last 10 iterations)

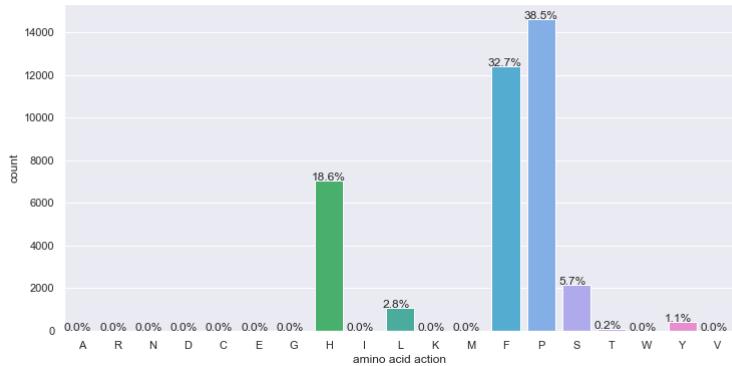


## 2. Position action analysis

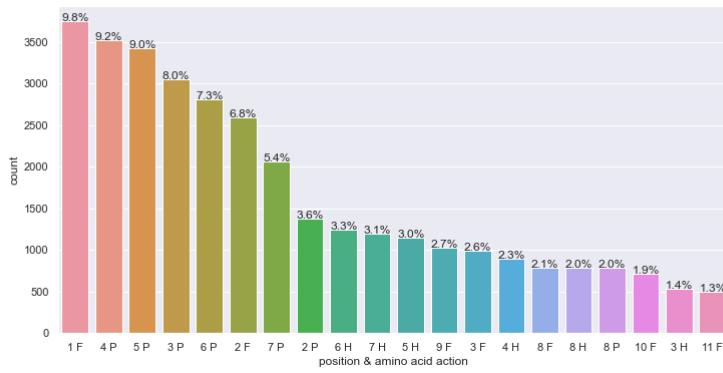
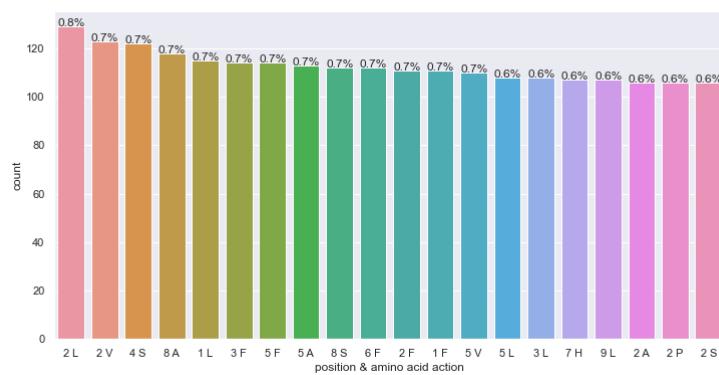


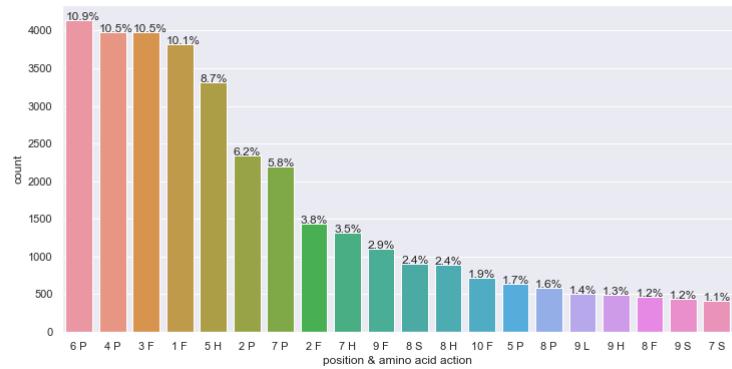
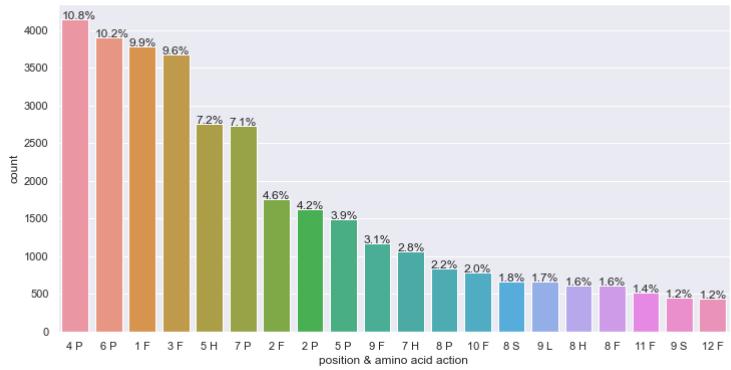
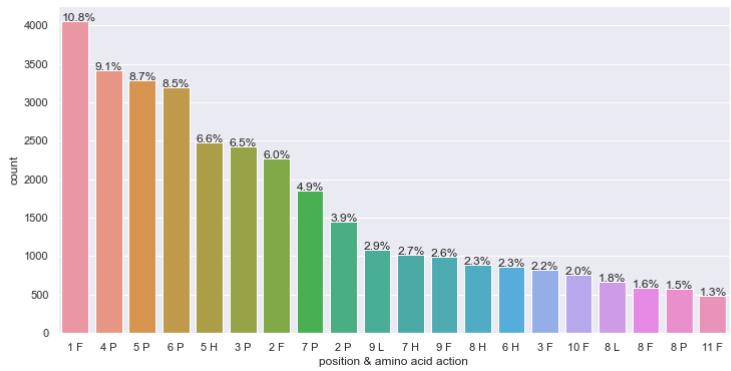
## 3. Amino action analysis



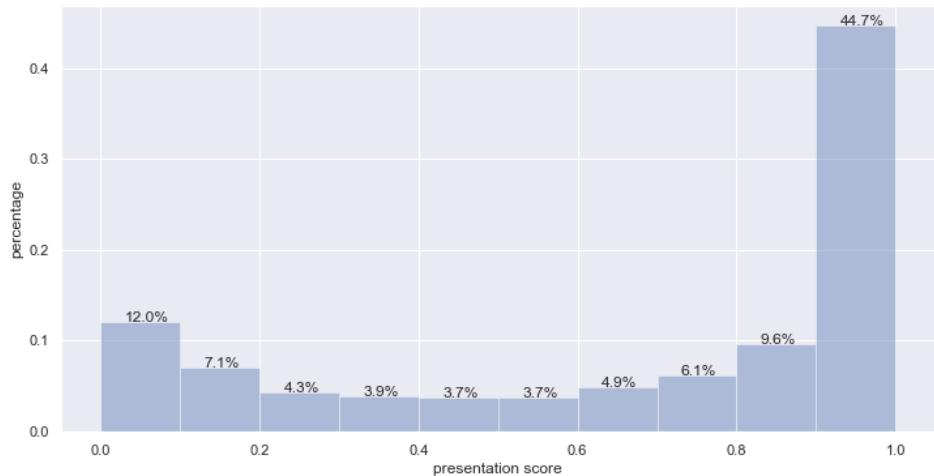


#### 4. Position & Amino action analysis

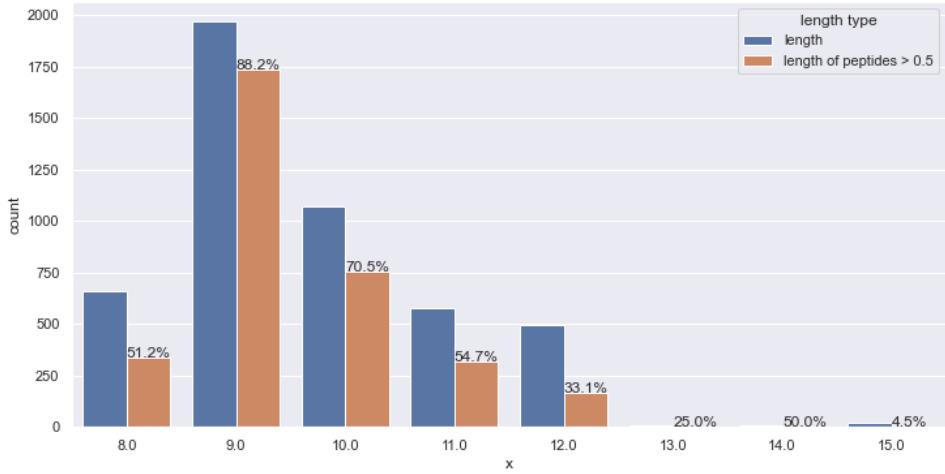




## 5. Reward analysis



## 6. Length analysis



## 4. Prioritized Experience Replay

In the original paper, they used Temporal Difference error to estimate the priority of each transition  $(S_{t-1}, A_{t-1}, S_t, R_t)$ .

$$\text{TD error: } \delta_t = R_t + \gamma \max_a(Q(S_t, a)) - Q(S_{t-1}, A_{t-1})$$

$$Q = R_t + \gamma * (R_{t+1} + \dots + \gamma^{T-t-1} R_T)$$

$$\text{Priority: } p_i = |\delta_i| + \epsilon \rightarrow P(i) = \frac{p_i^\alpha}{\sum_k p_k^\alpha}$$

"The exponent  $\alpha$  determines how much prioritization is used, with  $\alpha = 0$  corresponding to the uniform case."

$$\text{Correction with importance-sampling weights: } w_i = \left( \frac{1}{N} \cdot \frac{1}{P(i)} \right)^\beta$$

With  $P(i)$ , transitions are sampled with non-uniform probability. By setting  $\beta = 1$  at the end of training, the above weight can fully compensate for the non-uniform probabilities. At the beginning of training, the value of  $\beta$  should be small.

$\max(w_i)$  is used to normalize the weight  $w_i$  to be smaller than 1.

In our experiments, instead of using TD error to calculate the priority of single transition, we can use the hamming distance and the final presentation scores to calculate the priorities of the entire episodes. The episode with smaller hamming distance and higher presentation scores will have higher priority.

### How to design the priority function?

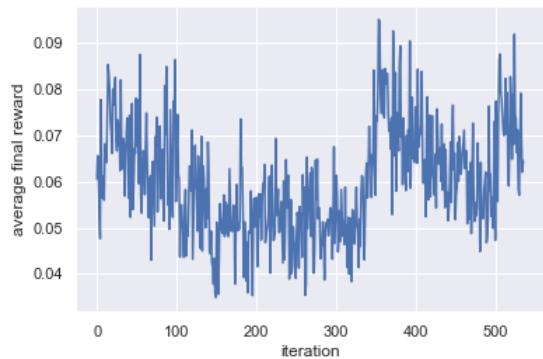
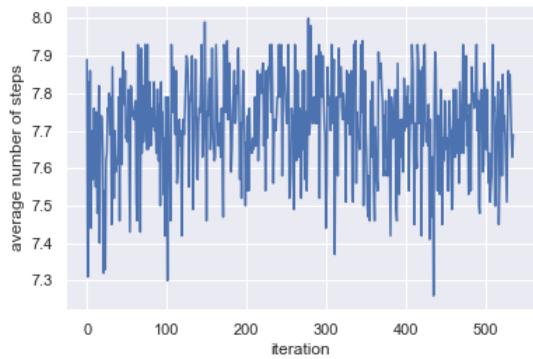
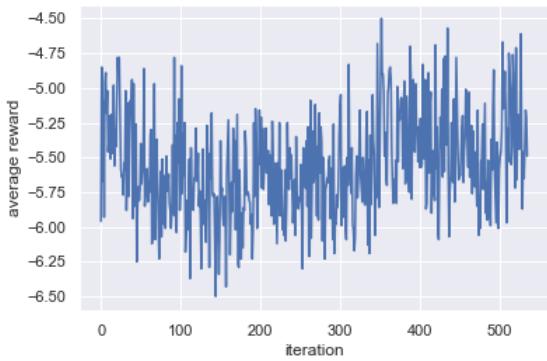
$$p_i = (d_{\max} - d + 1)^\phi \cdot s_T \rightarrow P(i) = \frac{p_i^\alpha}{\sum_k p_k^\alpha}$$

change  $s_T$  to  $(s_T - s_0)$

where  $s_T$  is the final presentation score;  $d$  is the hamming distance between the initial state and the termination state;  $d_{\max}$  is the maximum number of steps;  $\phi$  and  $\alpha$  are two hyperparameters with values smaller than 1. Currently, both  $\phi$  and  $\alpha$  are set to be 0.5.

Note that  $(d_{\max} - d + 1)$  ensures that the value is greater than 0.

- Setting 1 (stop criteria = 0.6; sample rate = 50% from IEDB and 50% from random)

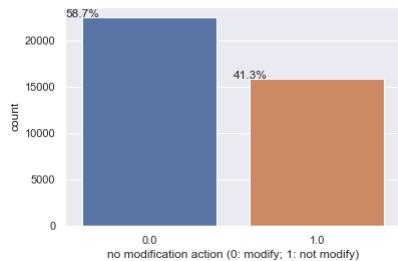
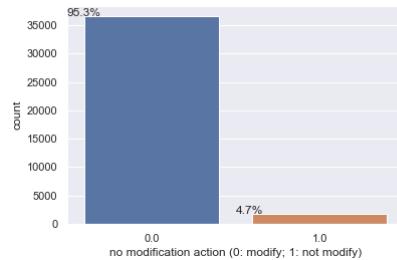
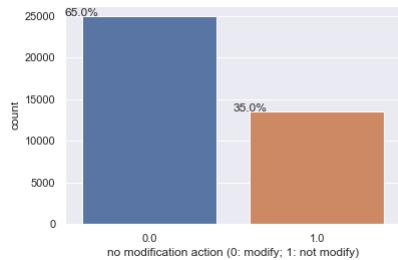


### 1. "No modification" action analysis

(first 10 iterations)  
the middle)

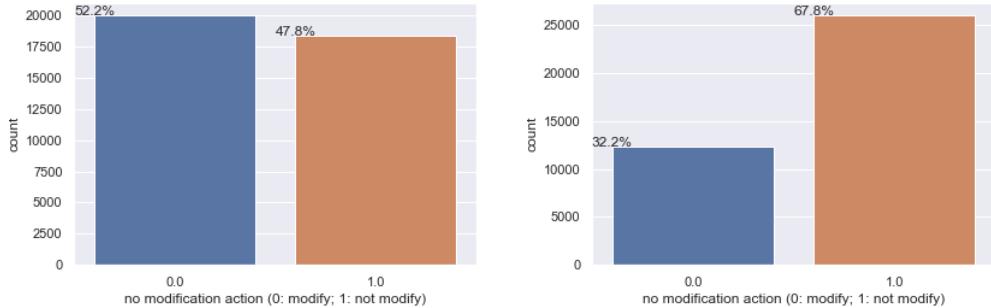
(10 iterations at 1/4)

(10 iterations in

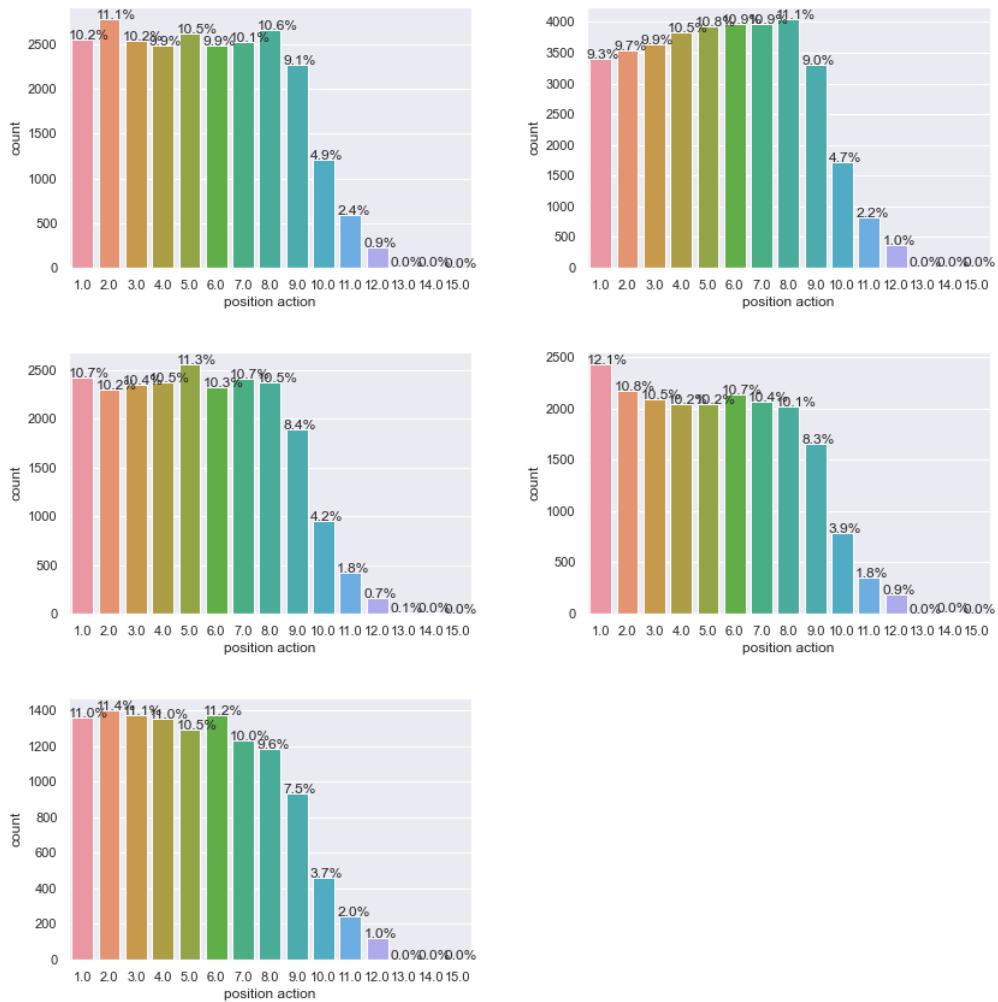


(10 iterations at 3/4)

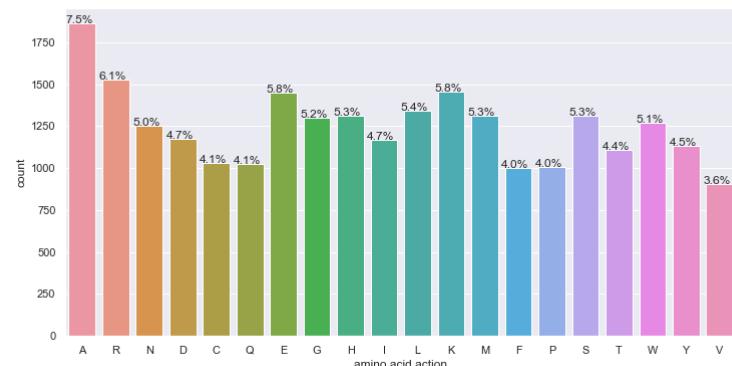
(last 10 iterations)

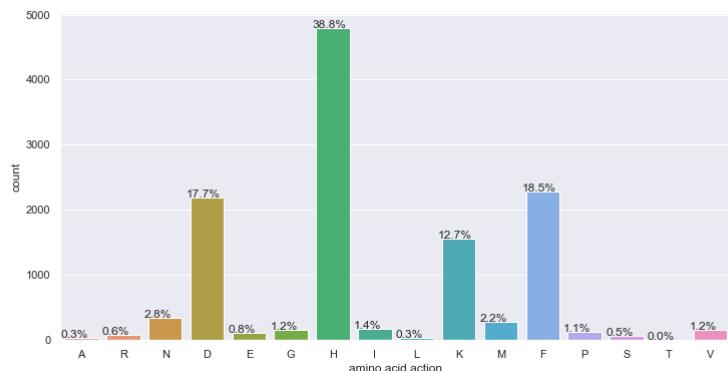
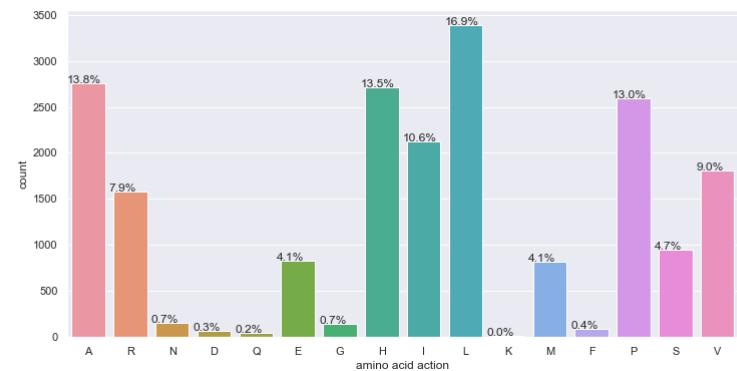
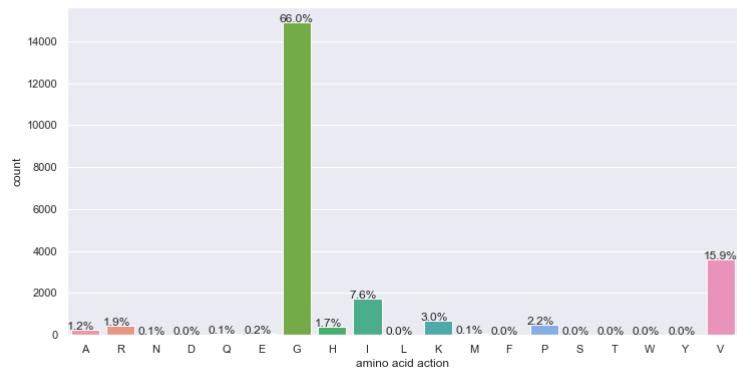
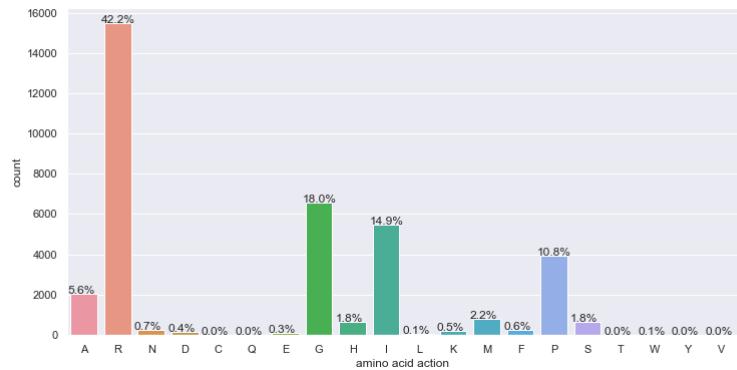


## 2. Position action analysis

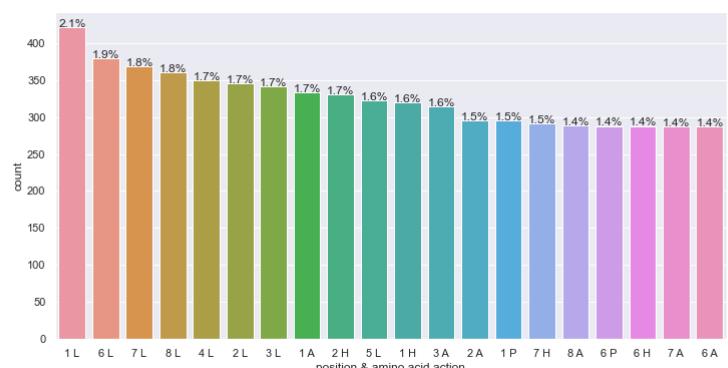
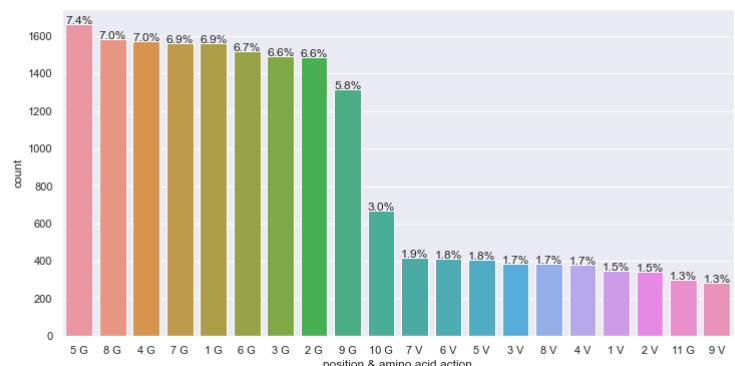
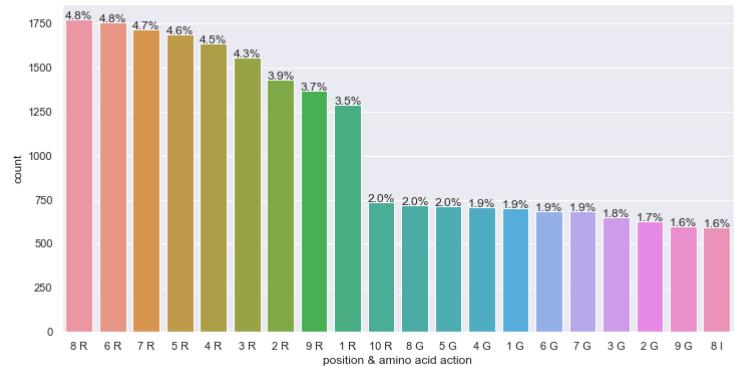
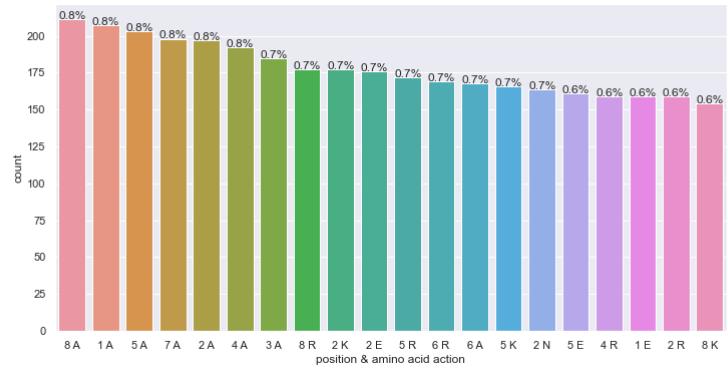


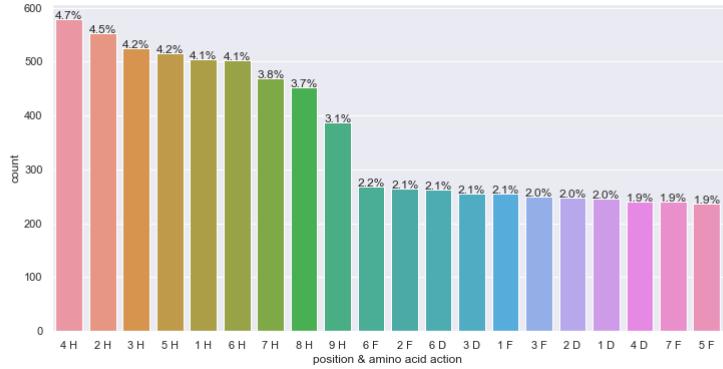
## 3. Amino action analysis



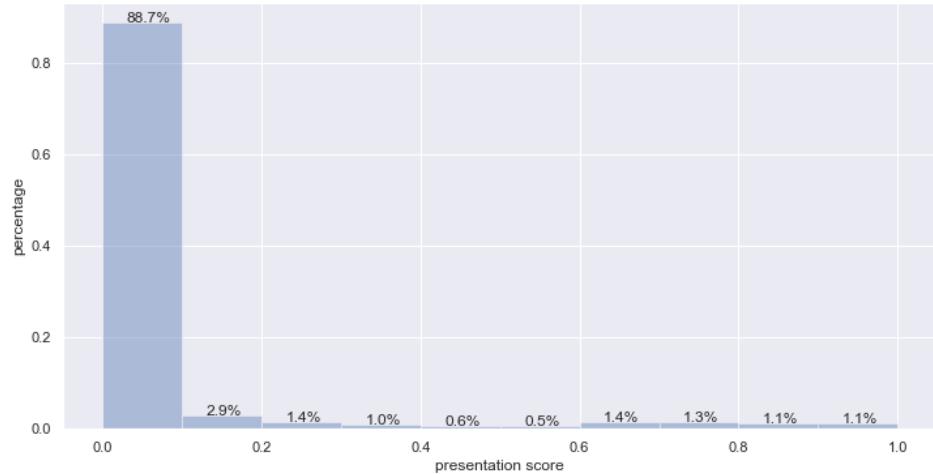


#### 4. Position & Amino action analysis

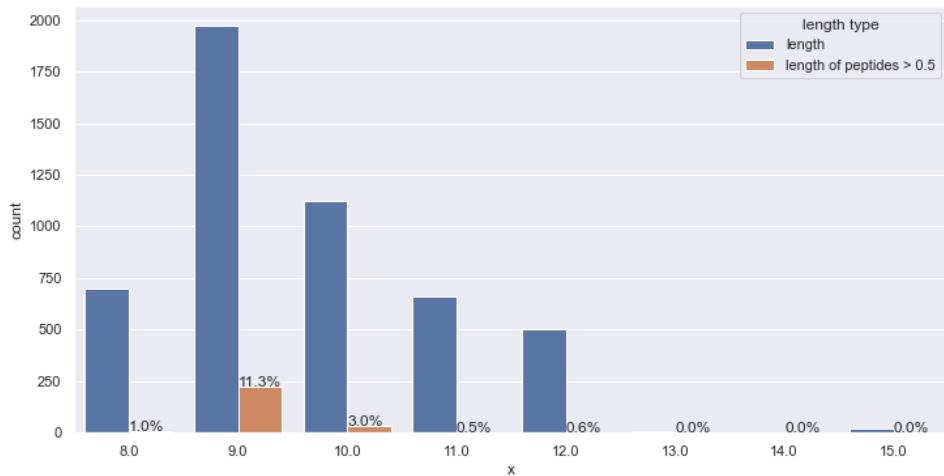




## 5. Reward analysis



## 6. Length analysis



## 5. Using undiscounted factors for "no modification" actions at the end of episodes

The returns for "modification" actions on non-terminal states  $S_{t-1}$ :  $y = R_t + \gamma V(S_{t+1})$ , that is the sum of current reward  $R_t$  and the estimated future rewards  $V(S_{t+1})$ .

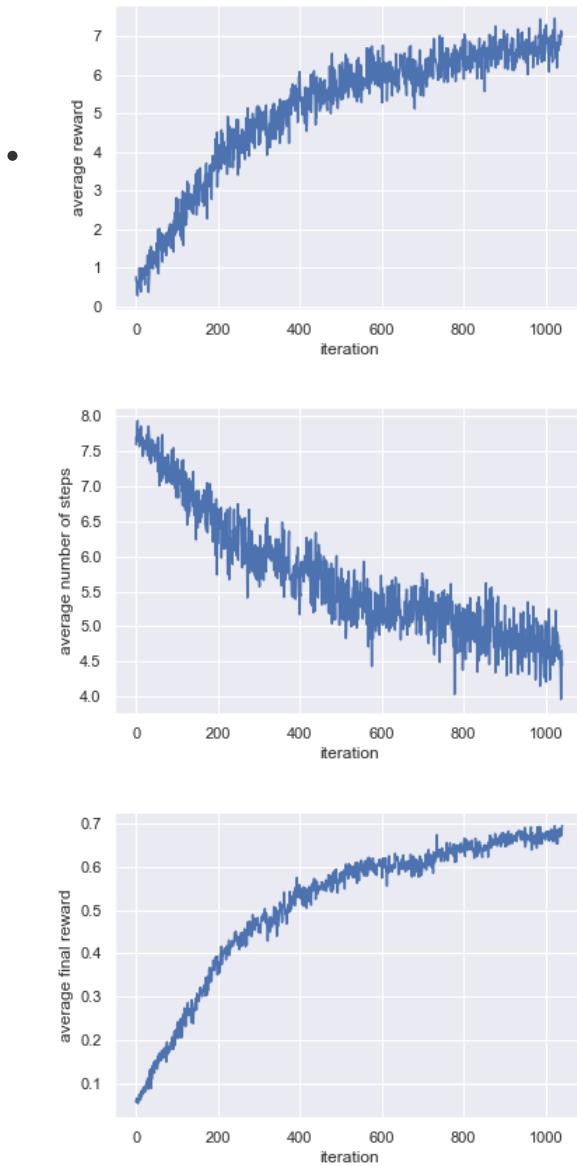
If states are terminal states, the return  $y = R_T$ .

This return can be used to calculate the advantage function or TD error and update the value function.

Instead of using the return function  $y = R_t + \gamma V(S_{t+1})$  for those "no modification" transitions at the end of episodes, we can use  $y = R_T$  to calculate their returns. In other words, all these "no modification" transitions are viewed as the same termination state.

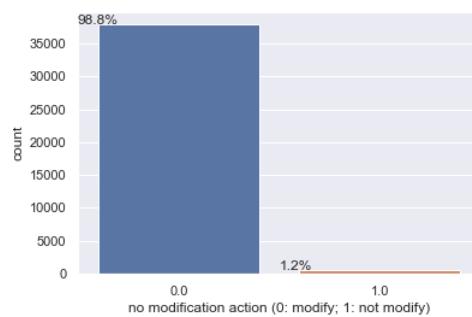
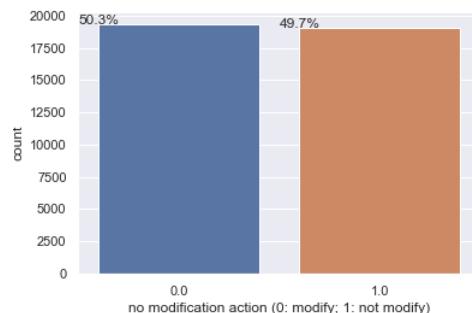
### **Experiment: stop criteria = 0.6**

Setting 1 (only final rewards; sample rate = 50% from IEDB and 50% from random)

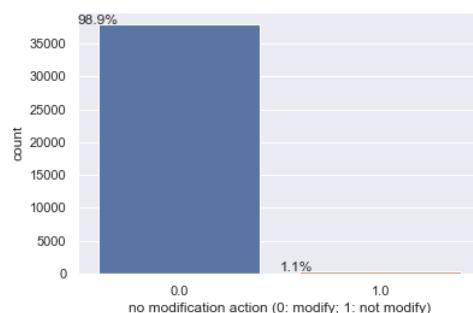


#### 1. "No modification" action analysis

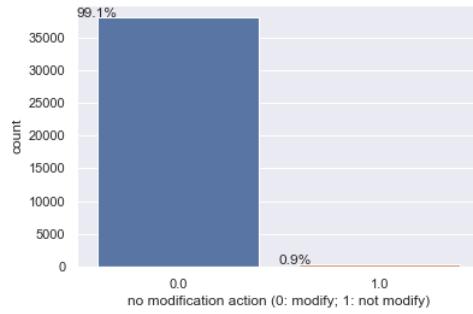
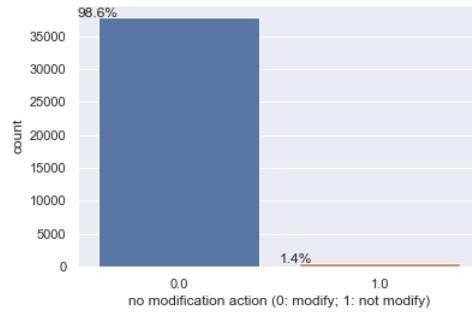
- |                                      |                        |                                  |
|--------------------------------------|------------------------|----------------------------------|
| (first 10 iterations)<br>the middle) | (10 iterations at 1/4) | (10 iterations in<br>the middle) |
|--------------------------------------|------------------------|----------------------------------|



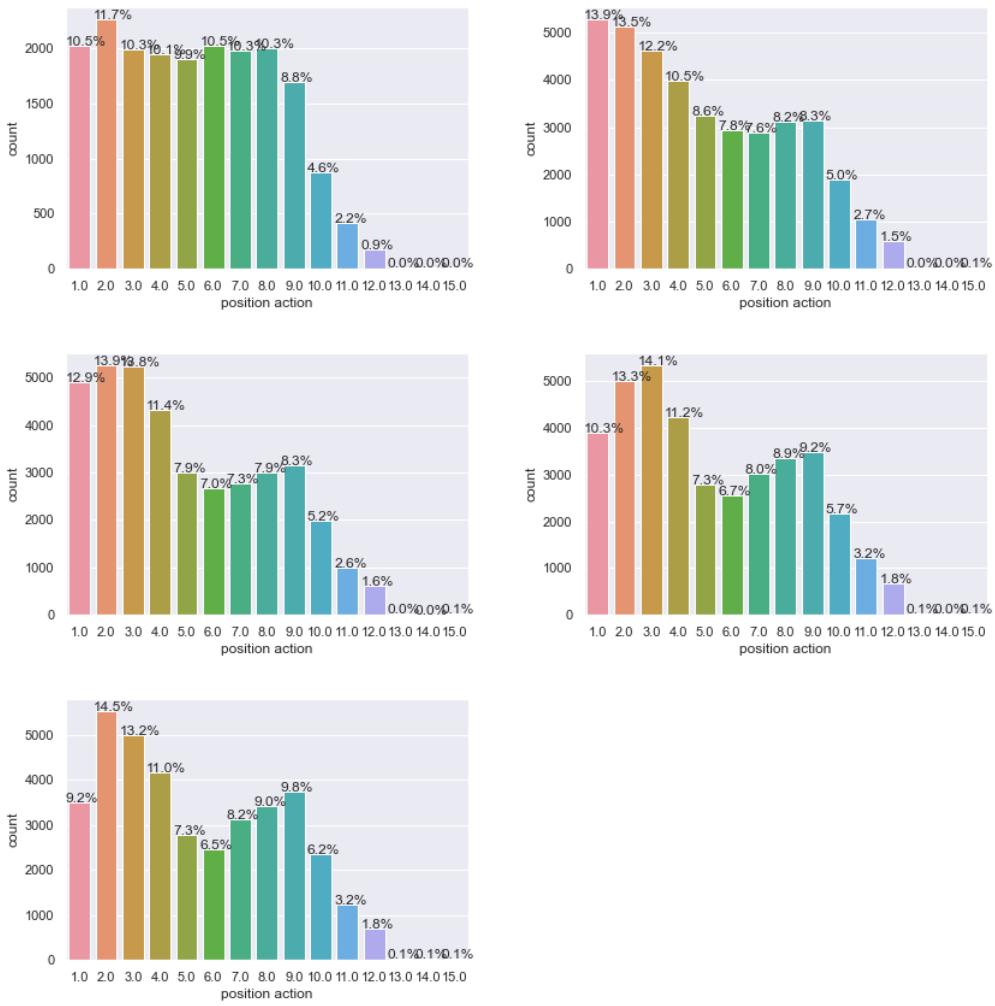
(10 iterations at 3/4)



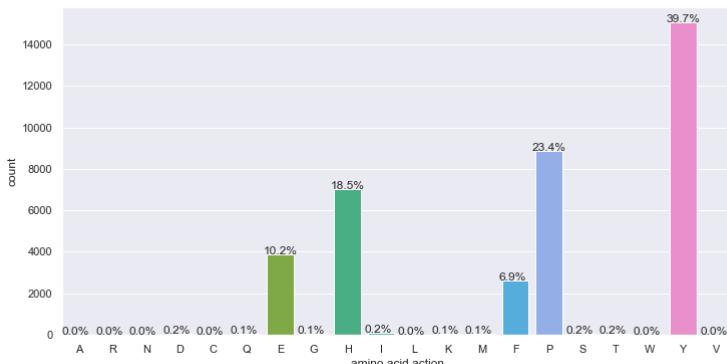
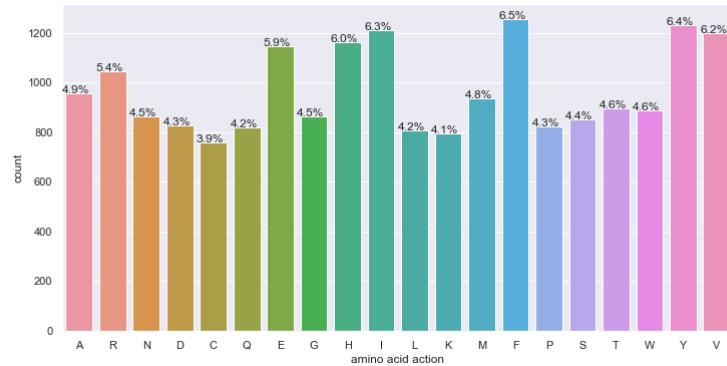
(last 10 iterations)

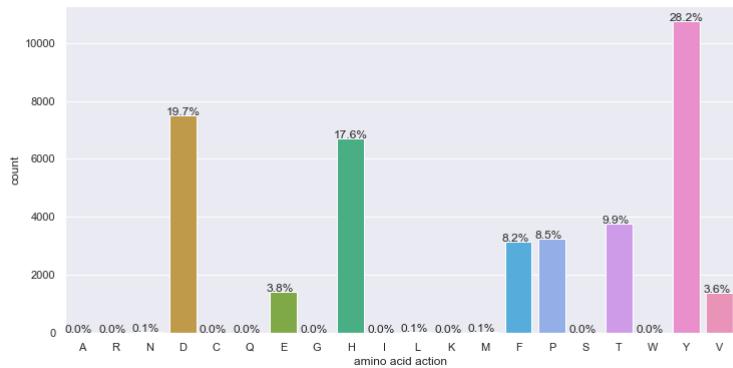
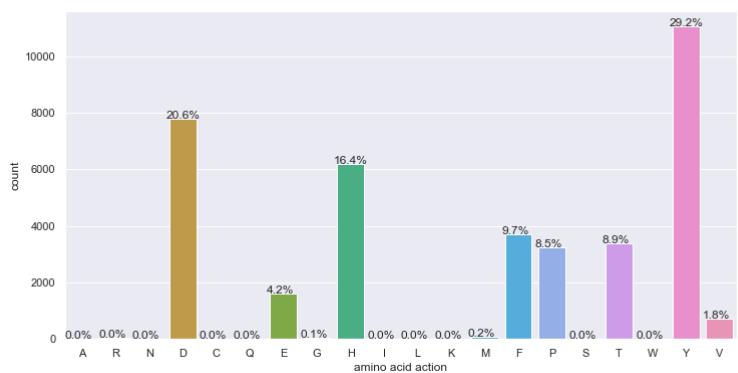
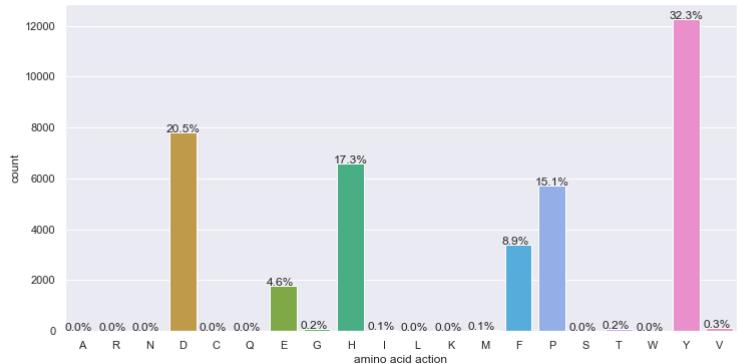


## 2. Position action analysis

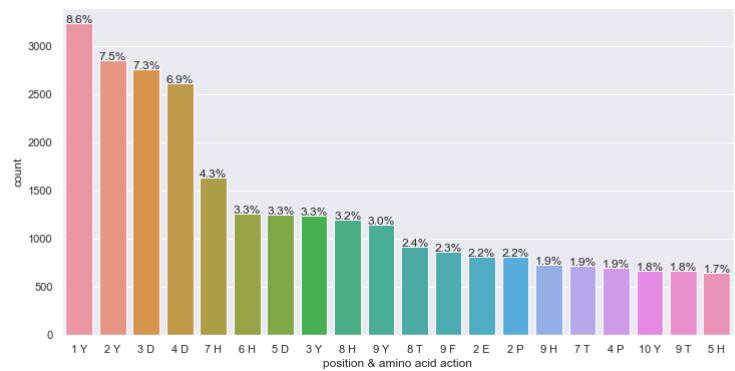
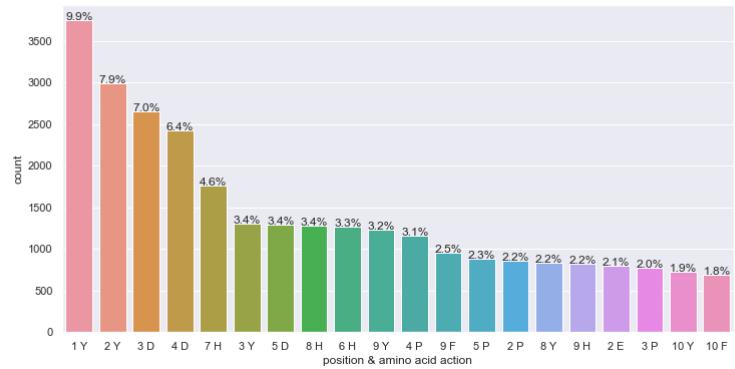
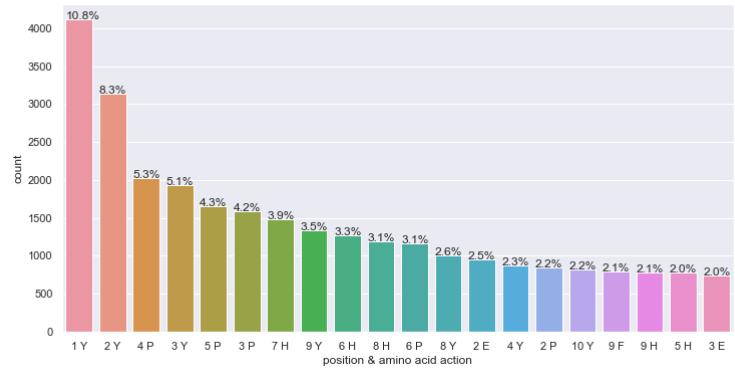
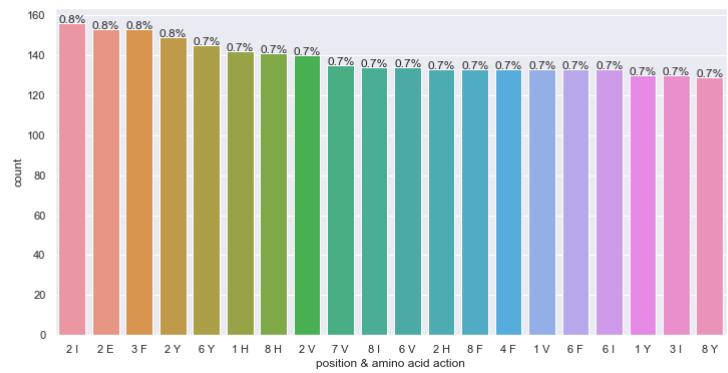


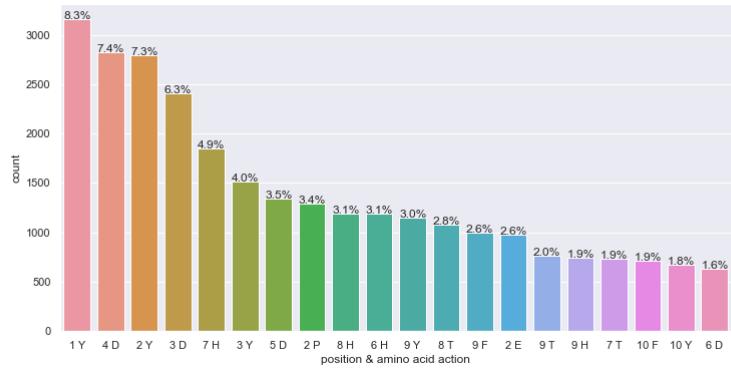
### 3. Amino action analysis



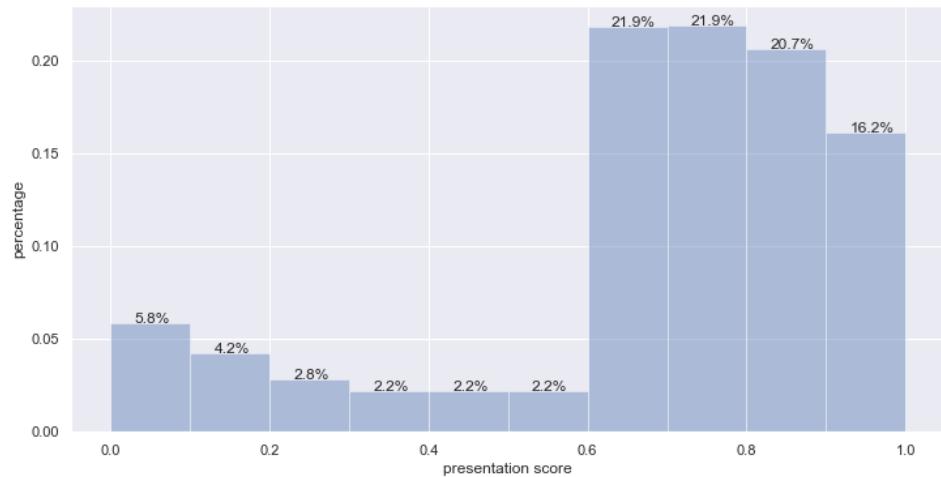


#### 4. Position & Amino action analysis

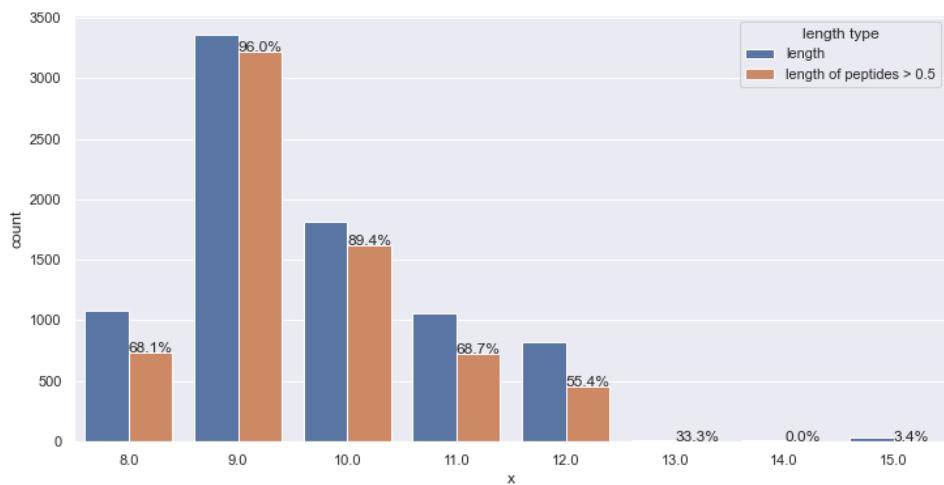




## 5. Reward analysis

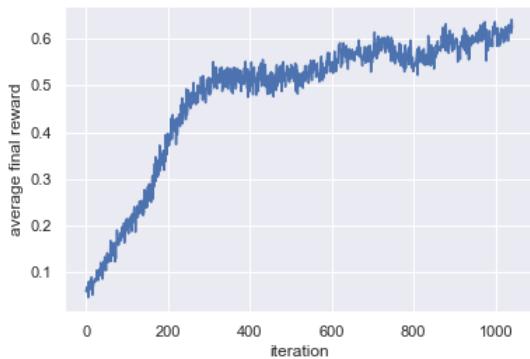
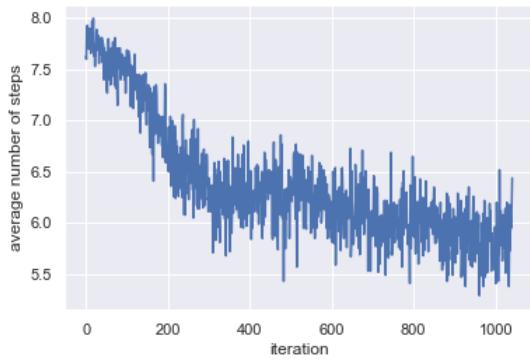
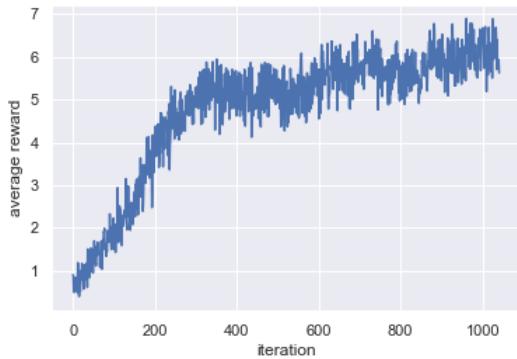


## 6. Length analysis



**Experiment: stop criteria = 0.75**

Setting 1 (only final rewards; sample rate = 50% from IEDB and 50% from random)



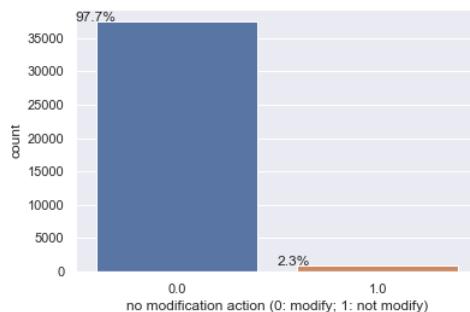
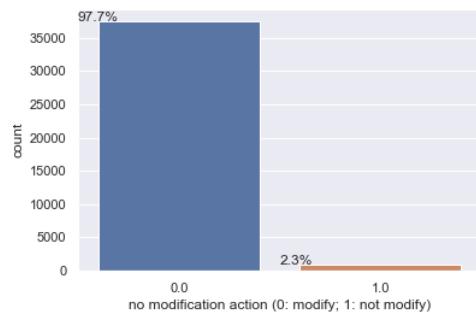
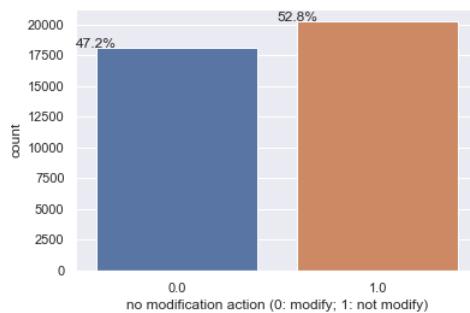
### 1. "No modification" action analysis

(first 10 iterations)

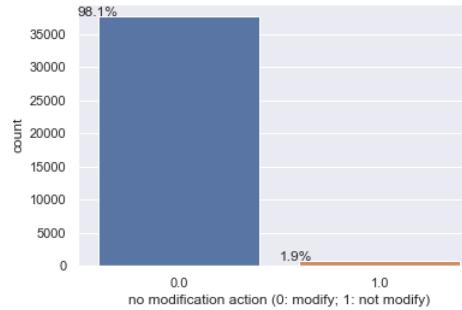
the middle)

(10 iterations at 1/4)

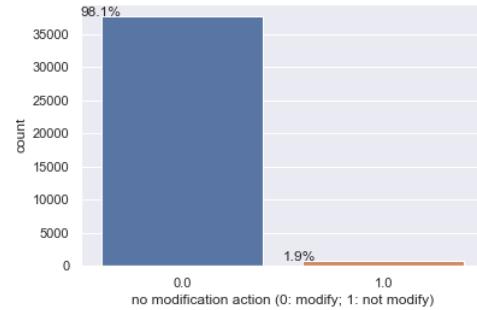
(10 iterations in



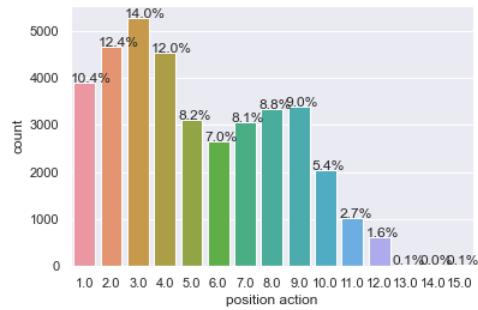
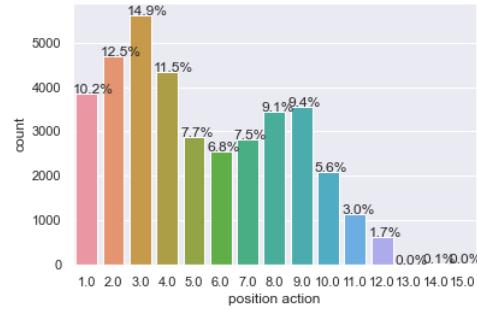
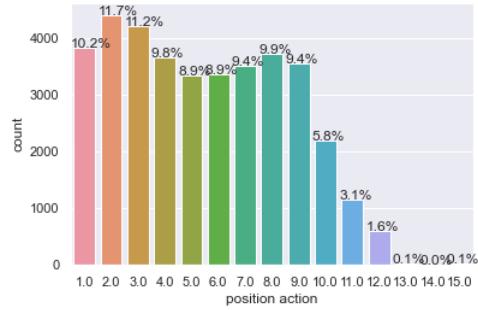
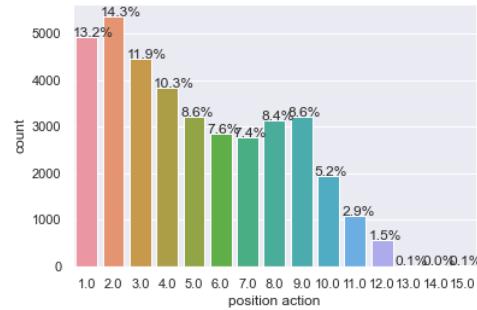
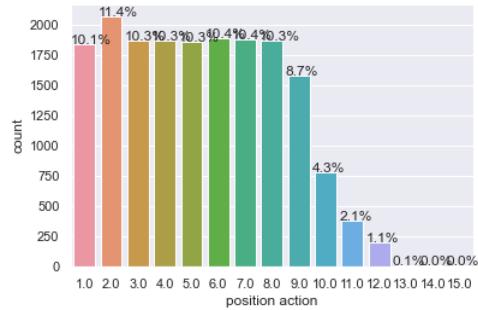
(10 iterations at 3/4)



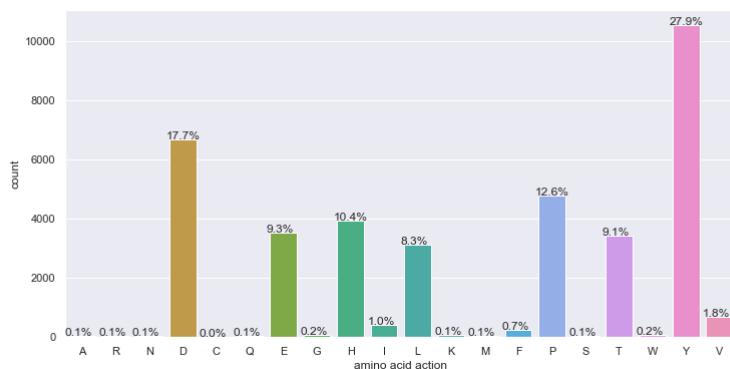
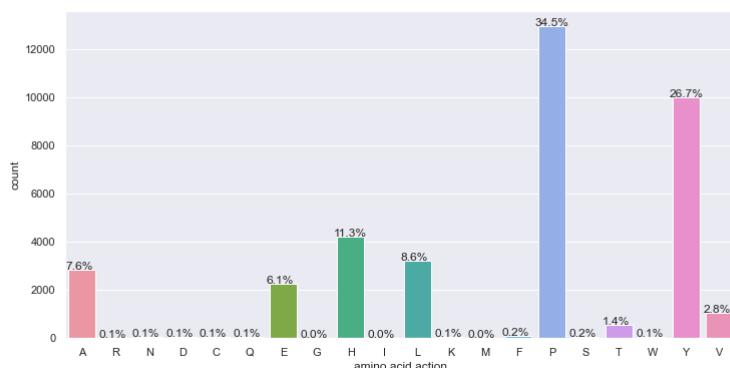
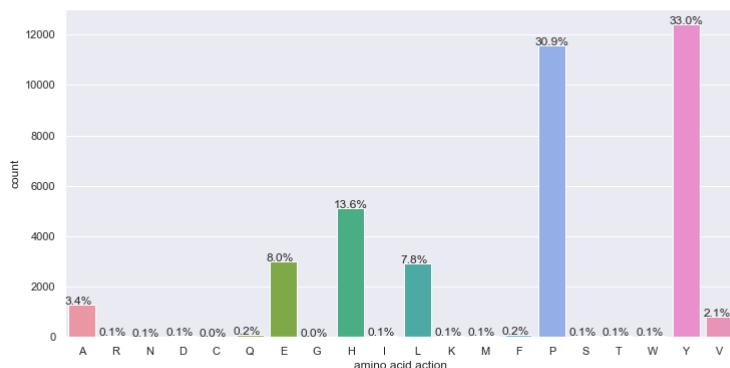
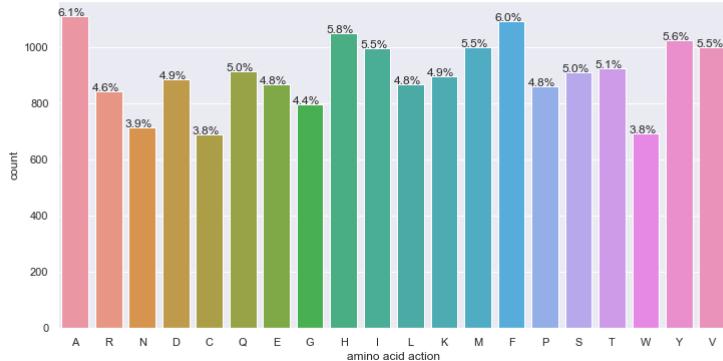
(last 10 iterations)

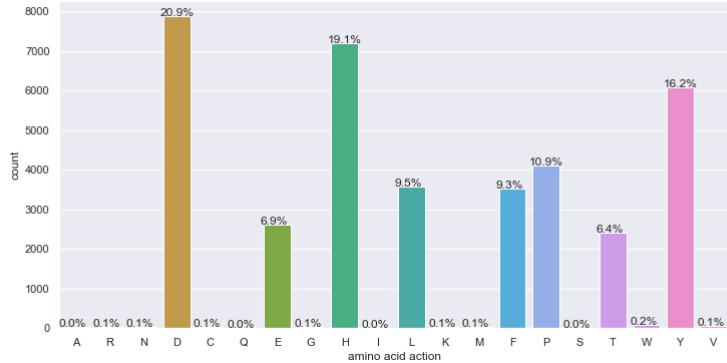


## 2. Position action analysis

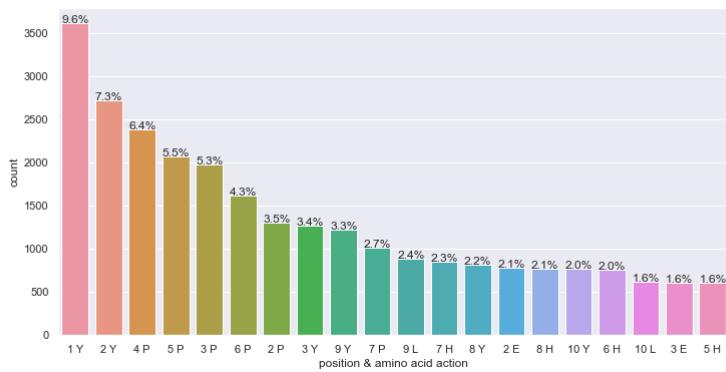
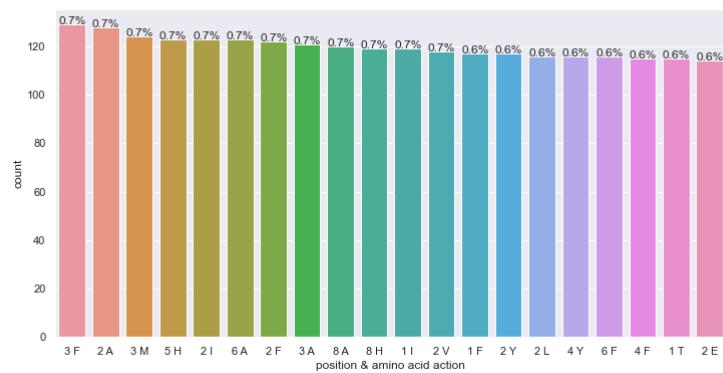


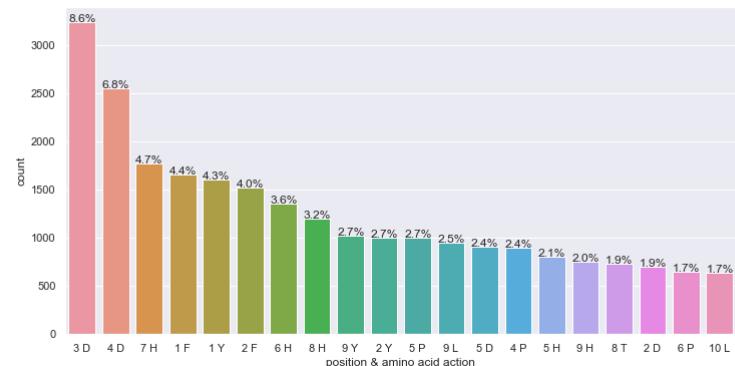
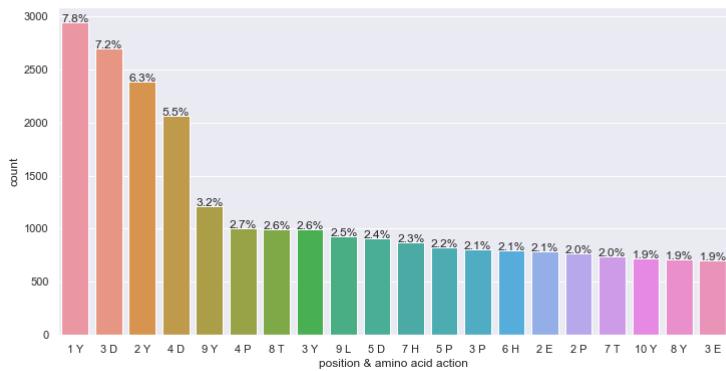
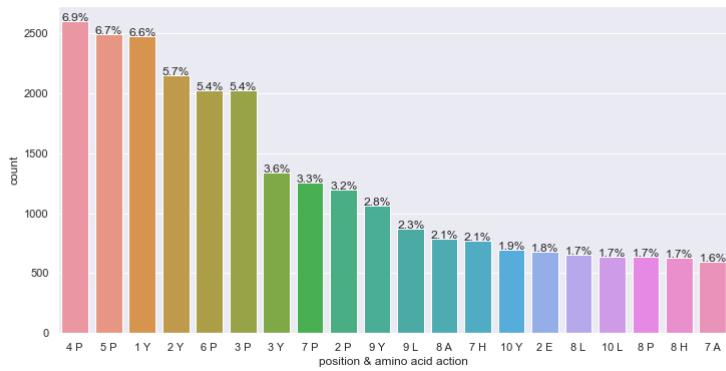
## 3. Amino action analysis



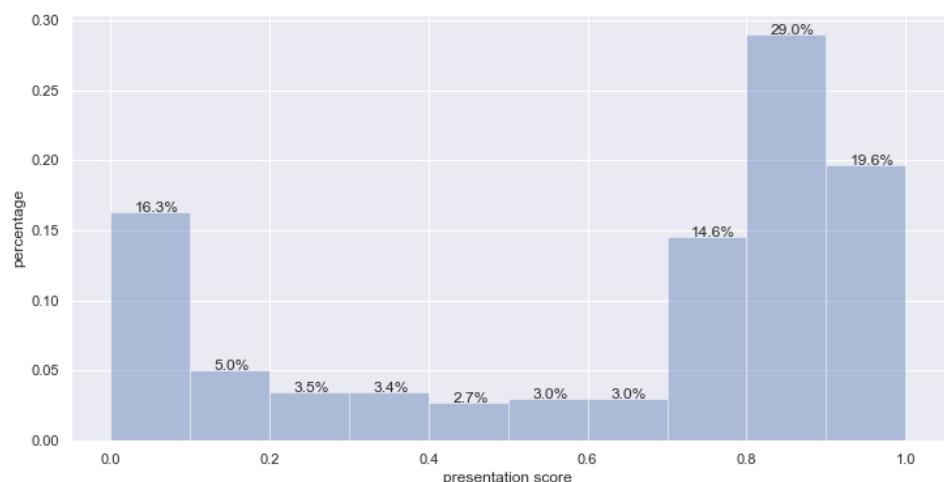


#### 4. Position & Amino action analysis

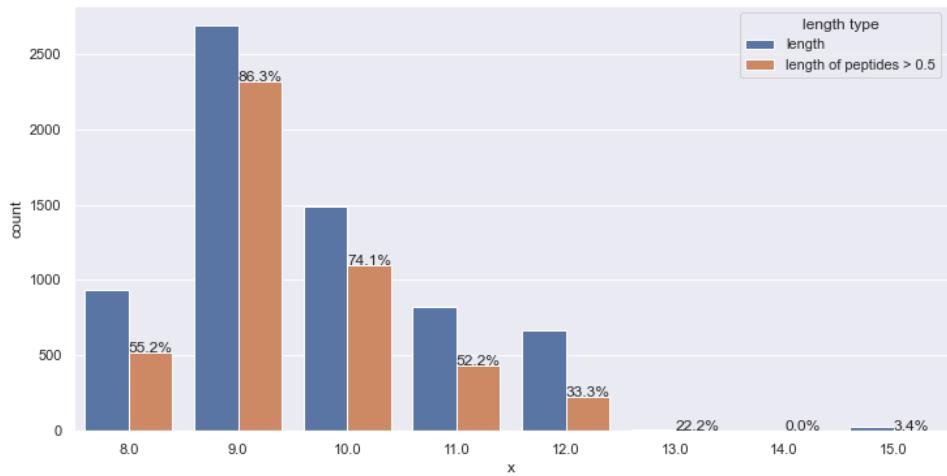




## 5. Reward analysis

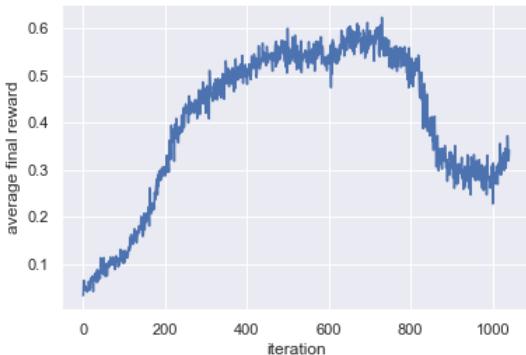
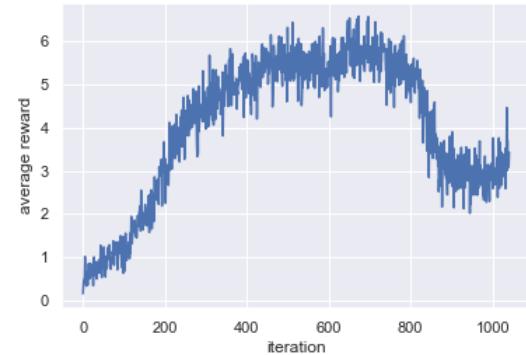


## 6. Length analysis



### Experiment: stop criteria = 1.0

Setting 1 (only final rewards; sample rate = 50% from IEDB and 50% from random)

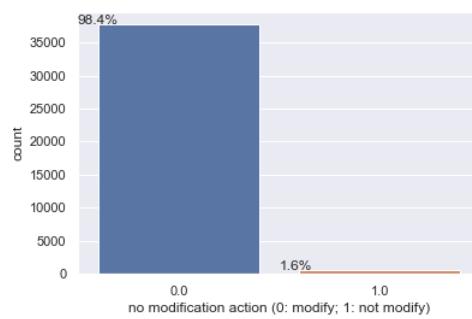
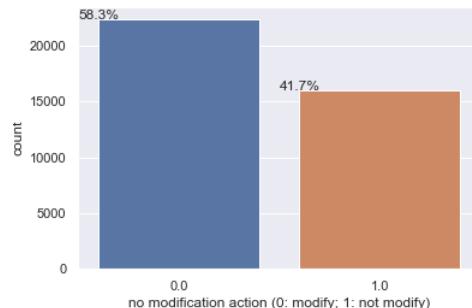


#### 1. "No modification" action analysis

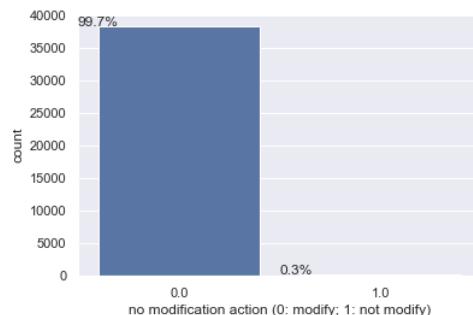
(first 10 iterations)  
the middle)

(10 iterations at 1/4)

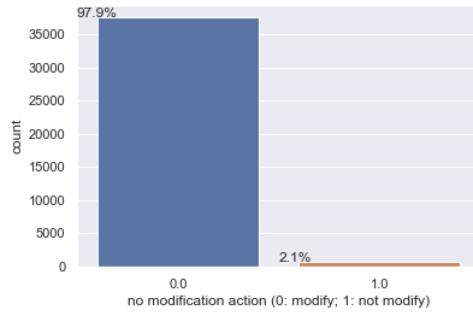
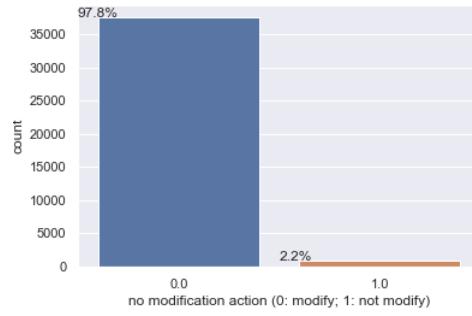
(10 iterations in



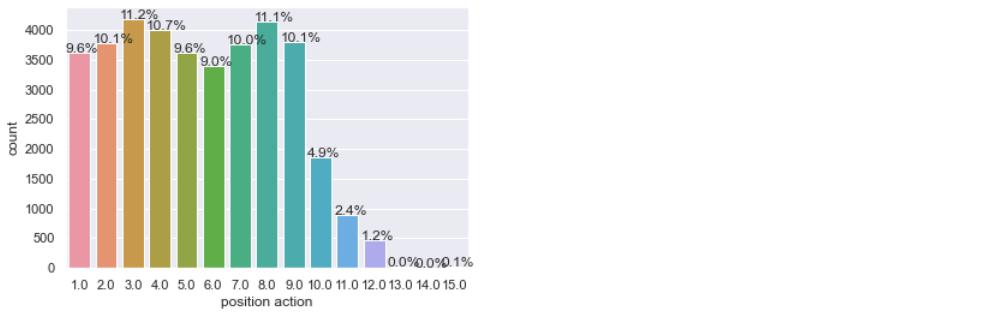
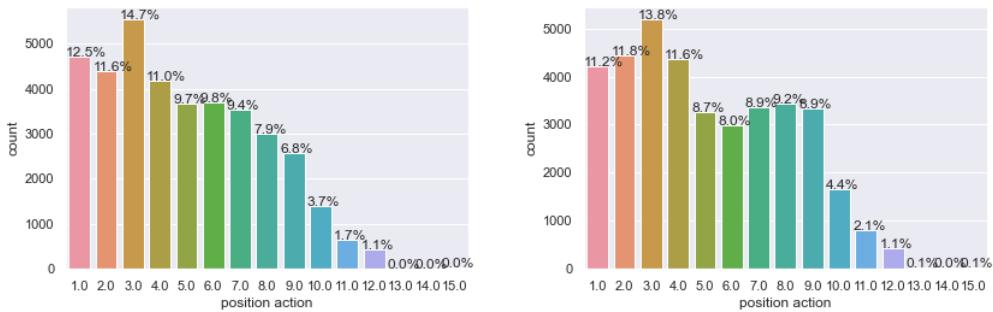
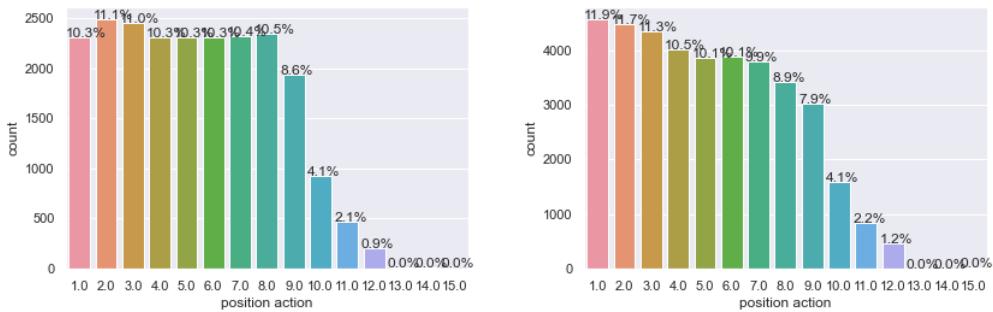
(10 iterations at 3/4)



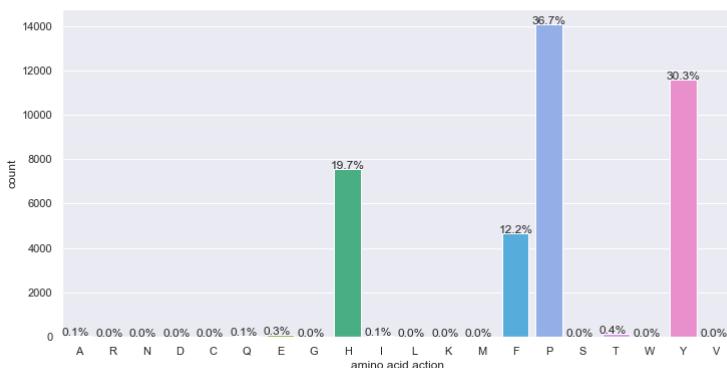
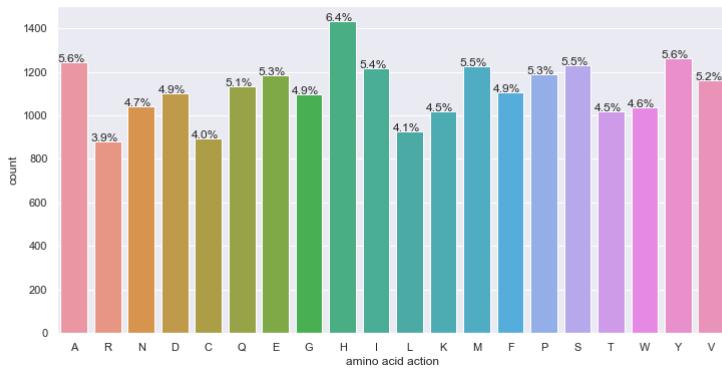
(last 10 iterations)

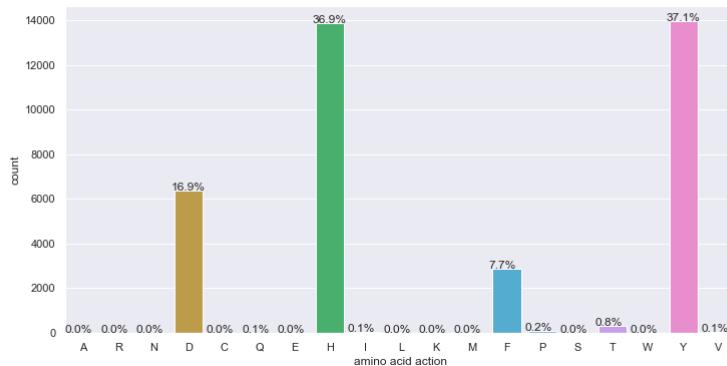
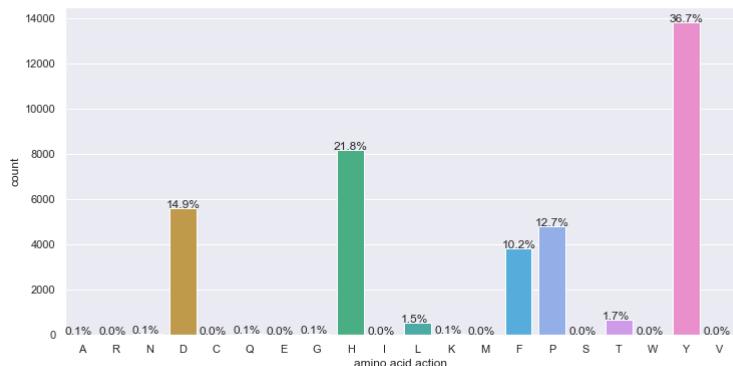
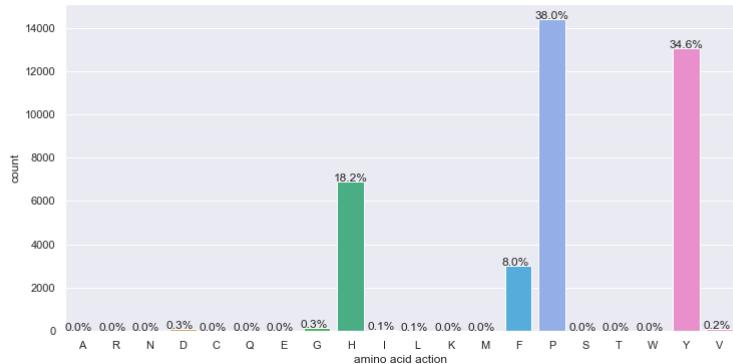


## 2. Position action analysis

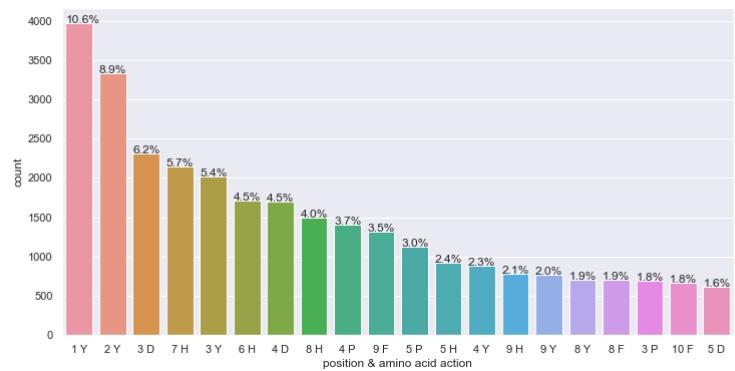
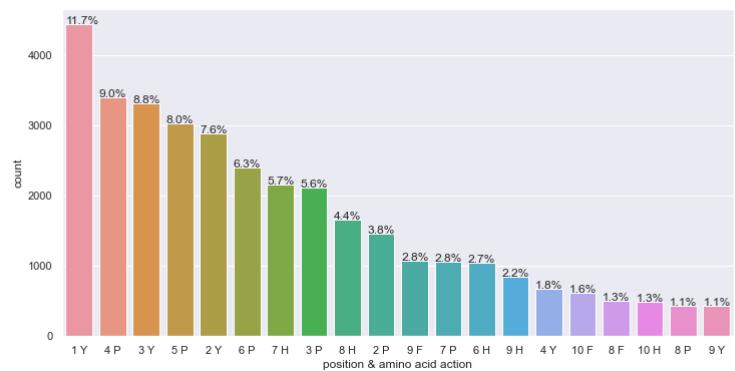
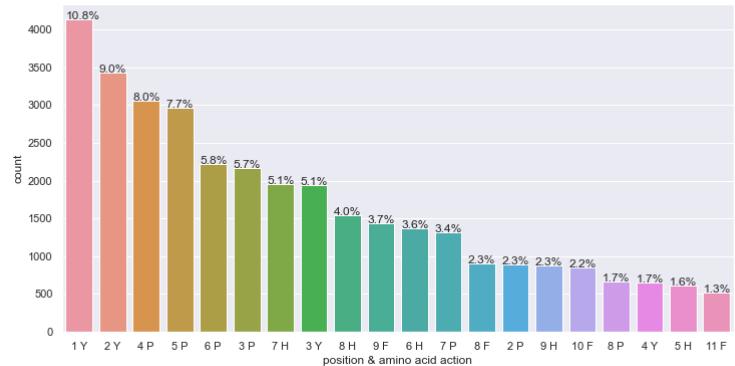
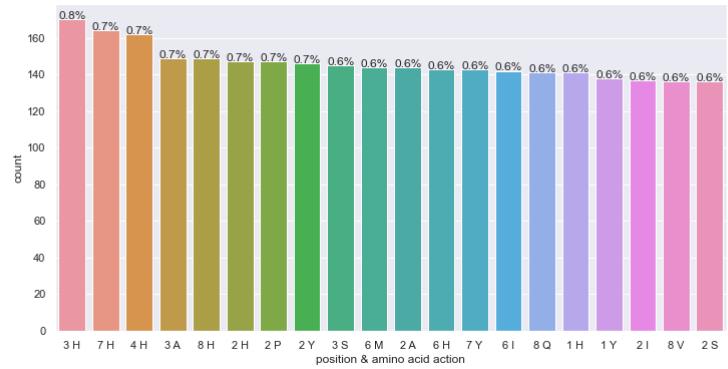


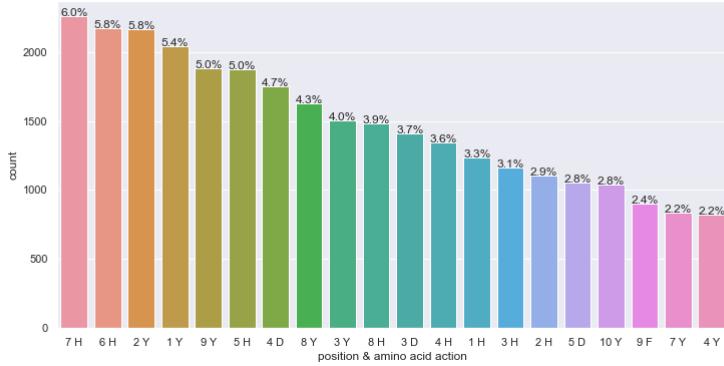
### 3. Amino action analysis



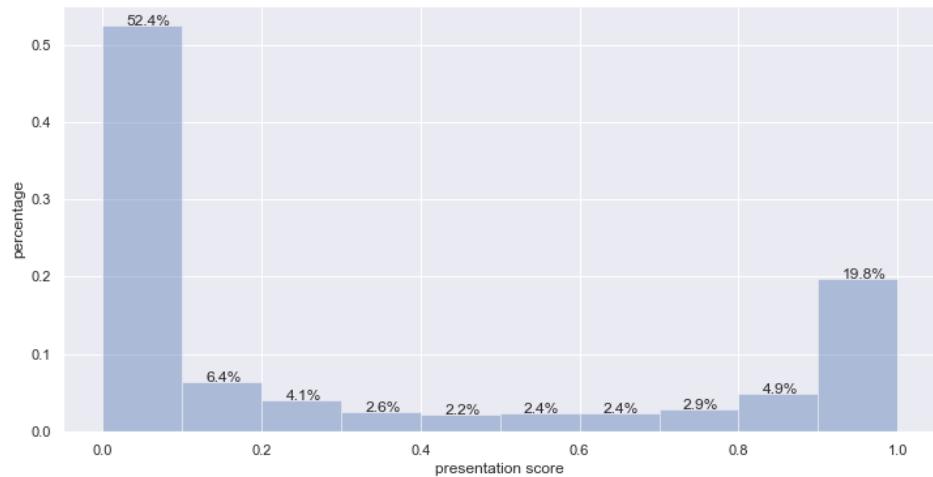


#### 4. Position & Amino action analysis

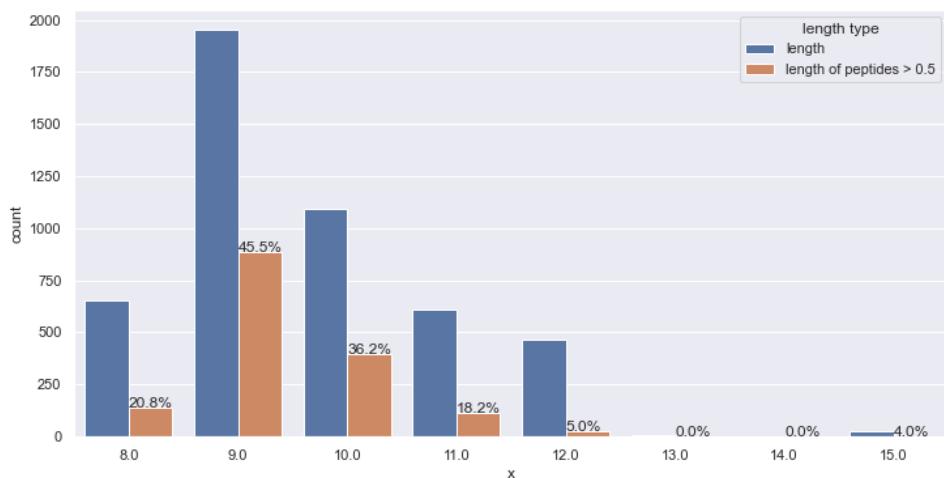




## 5. Reward analysis



## 6. Length analysis



## Update (6-2)

Next step:

- 1) add the constraint that the generated peptides are dissimilar with human peptides. Use as a component of the reward function or as a hard constraint or as a **post-processing** step.
- 2) write the paper
- 3) remove "no modification" action

Future plan:

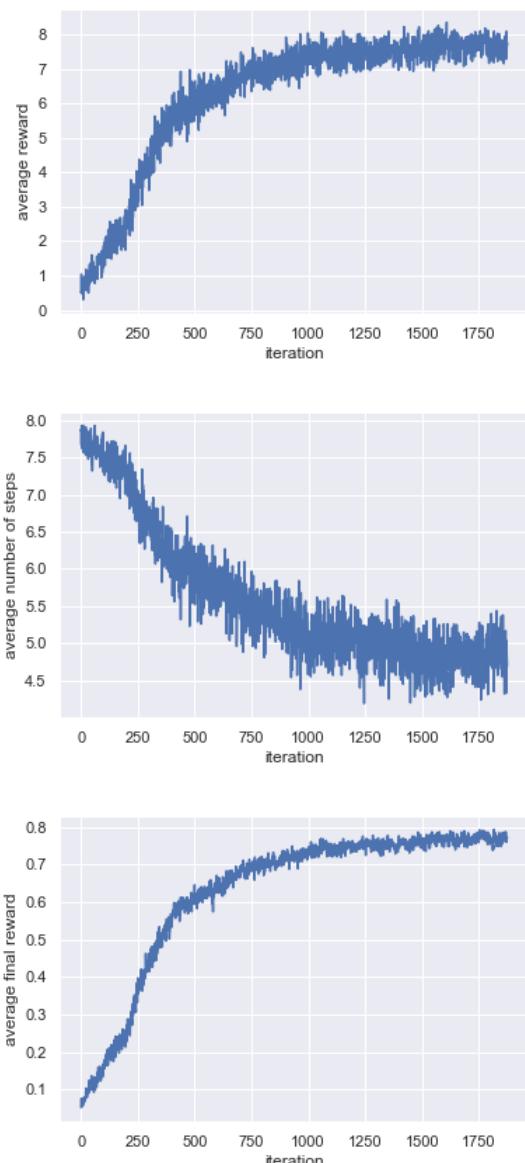
(1) Physics Guided DRL for TCR-pMHC Modeling in the Small Positive Sample Setting

(2) Combine generative classifier with DRL to optimize peptide

## 1. Result Analysis with only final rewards

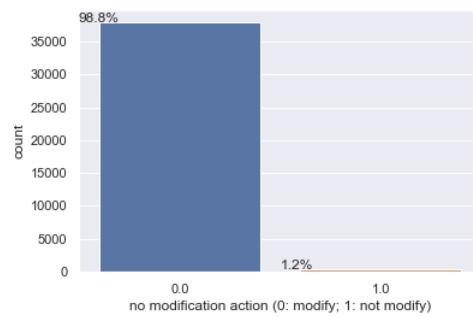
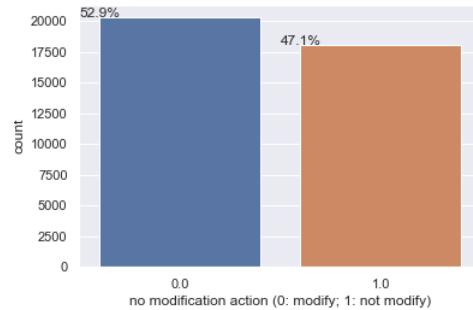
**Experiment: stop criteria = 0.75**

Setting 1 (only final rewards; sample rate = 50% from IEDB and 50% from random)

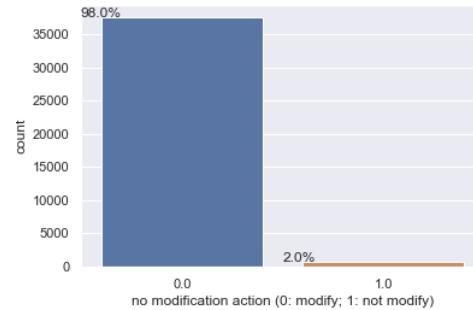


1. "No modification" action analysis

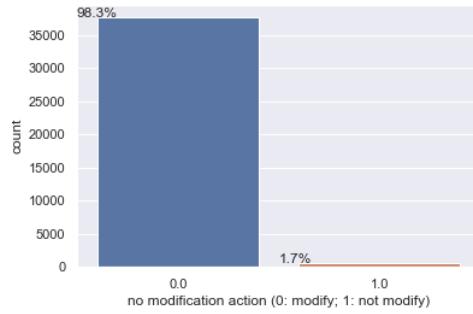
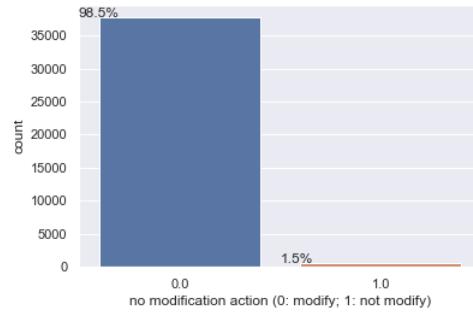
(first 10 iterations)  
(10 iterations at 1/4)  
the middle)  
(10 iterations in



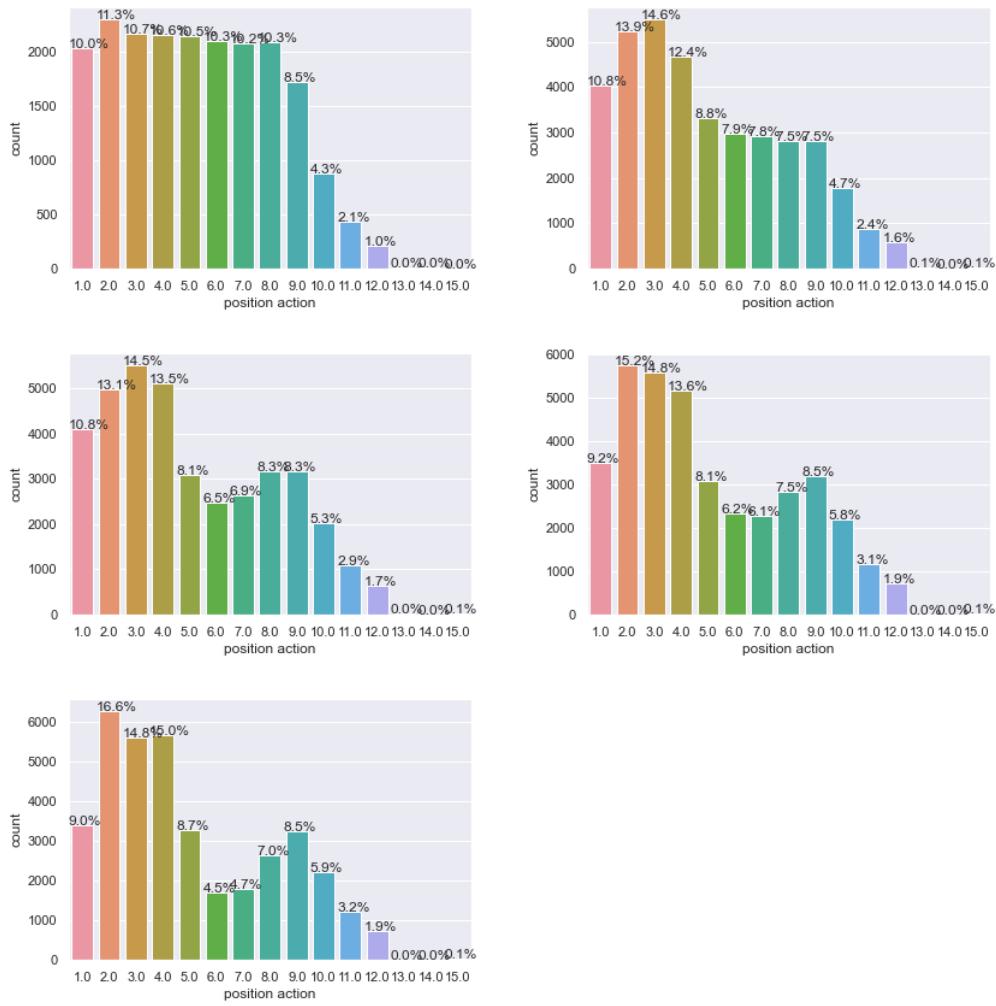
(10 iterations at 3/4)



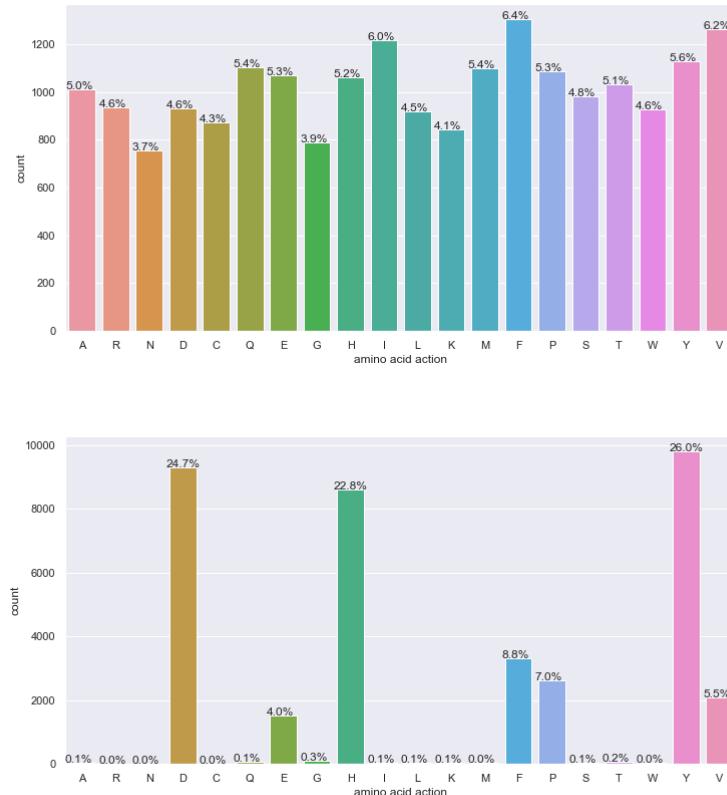
(last 10 iterations)

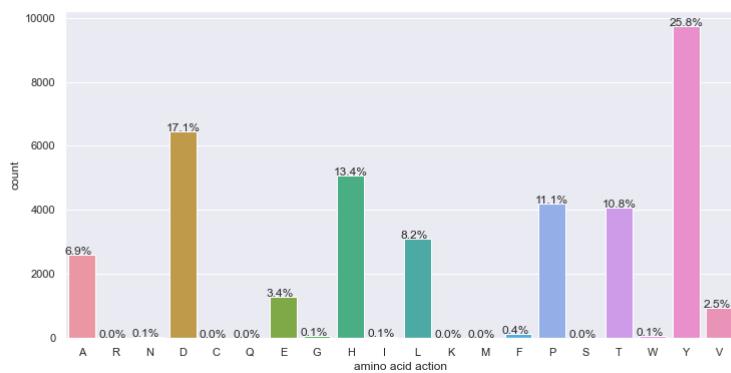
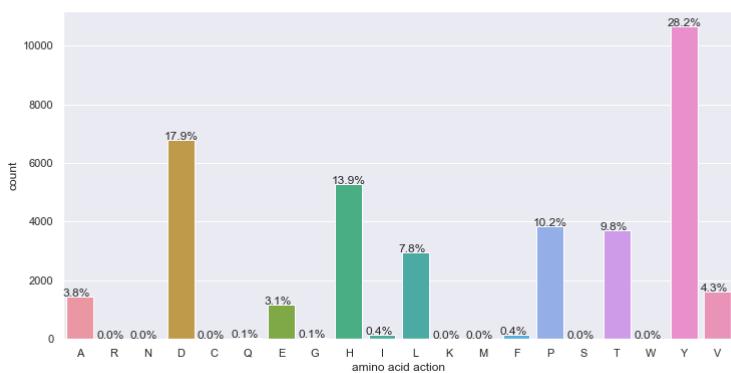
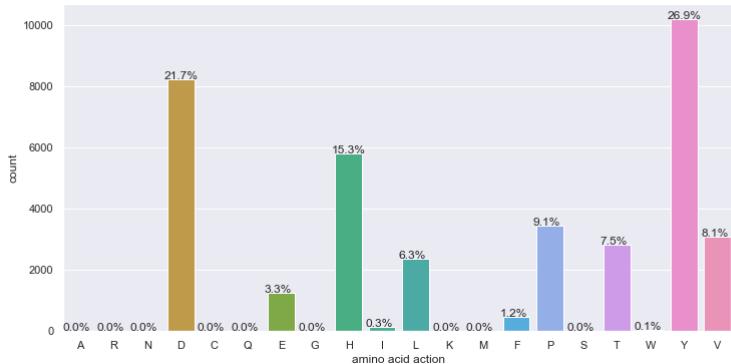


## 2. Position action analysis

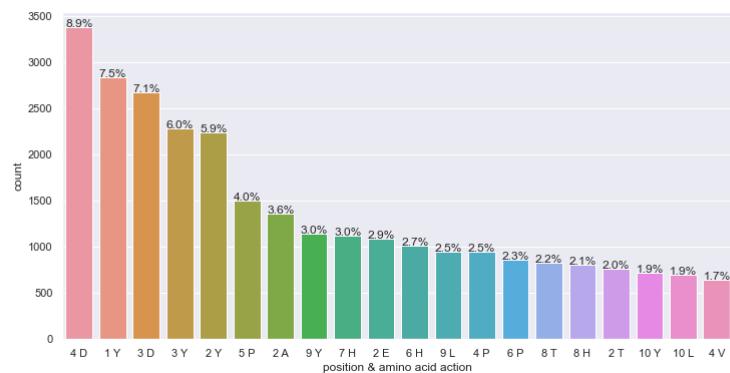
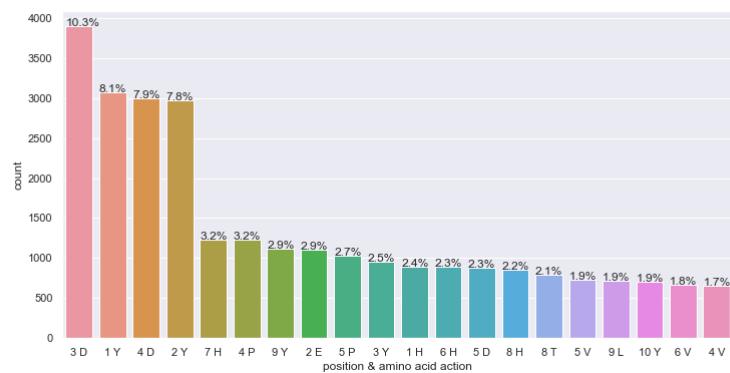
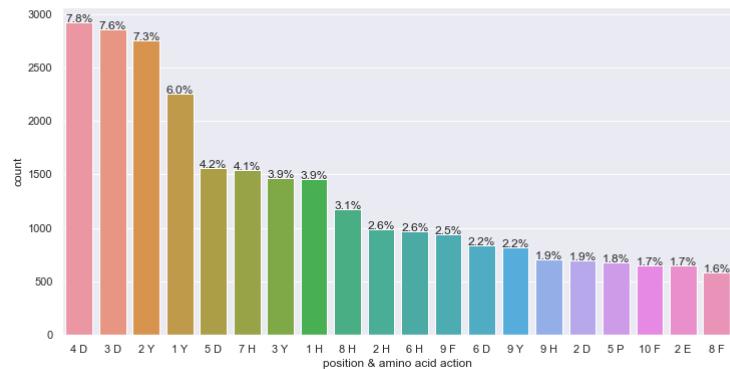
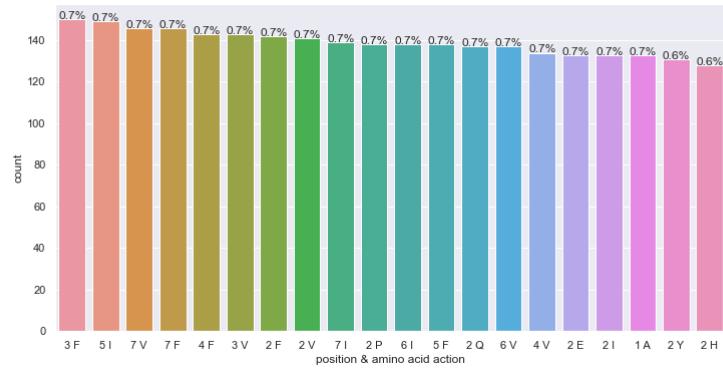


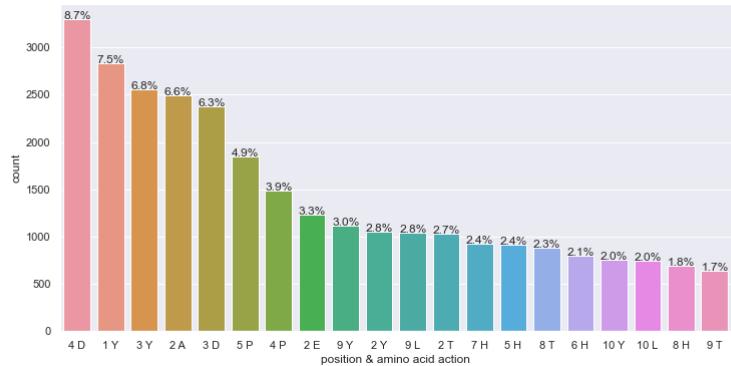
### 3. Amino action analysis



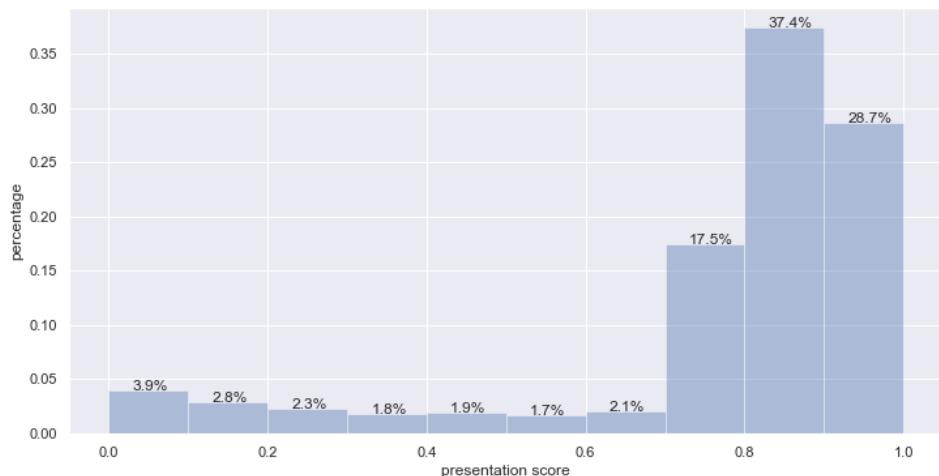


#### 4. Position & Amino action analysis

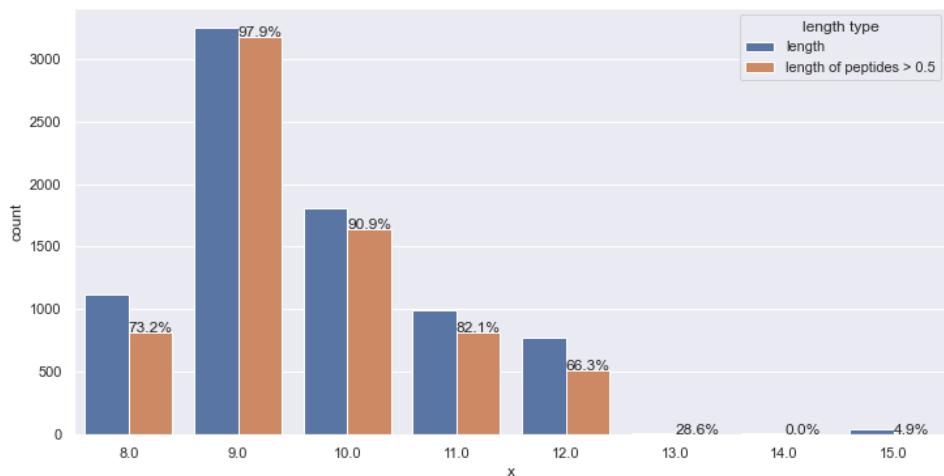




## 5. Reward analysis

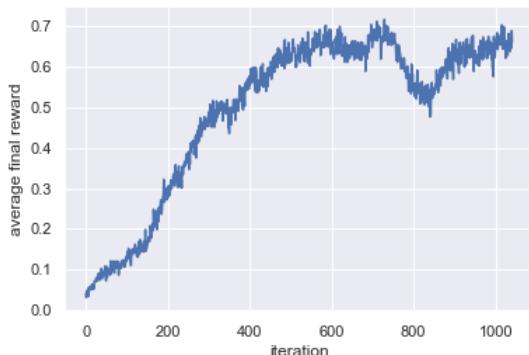
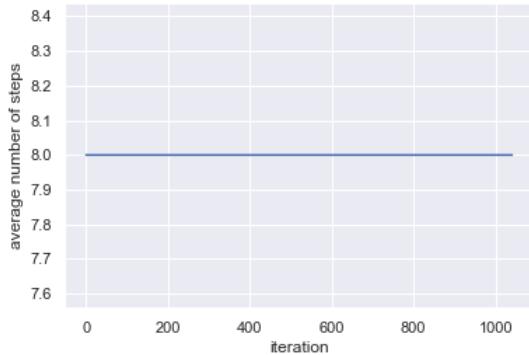
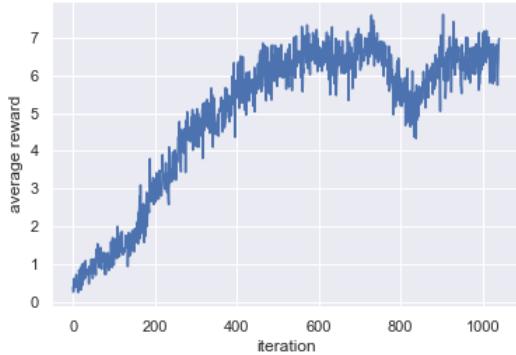


## 6. Length analysis



**Experiment: stop criteria = 1.0**

Setting 1 (only final rewards; sample rate = 50% from IEKB and 50% from random)

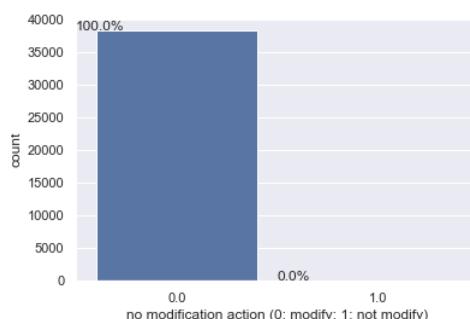
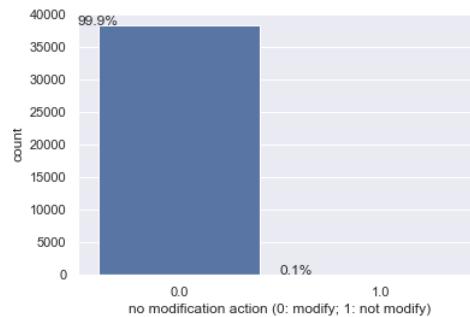
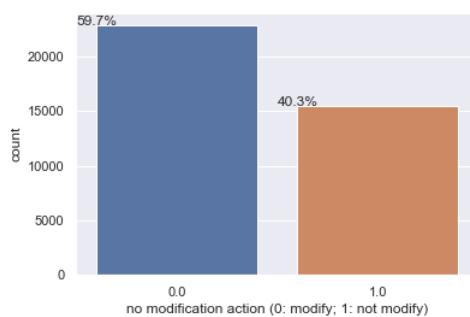


### 1. "No modification" action analysis

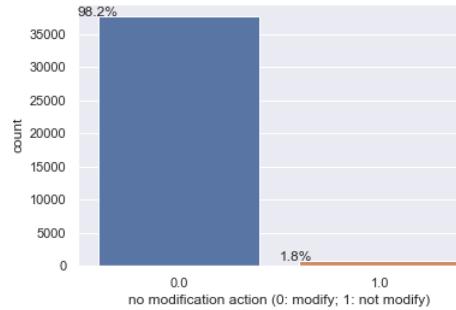
(first 10 iterations)  
the middle)

(10 iterations at 1/4)

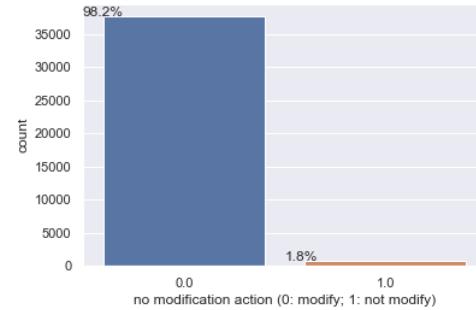
(10 iterations in



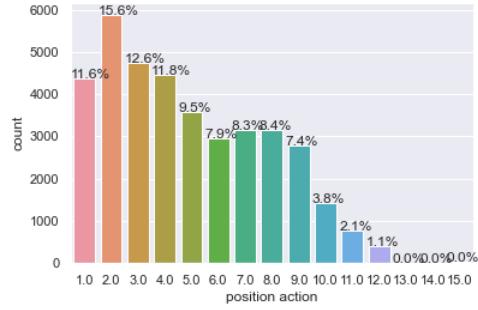
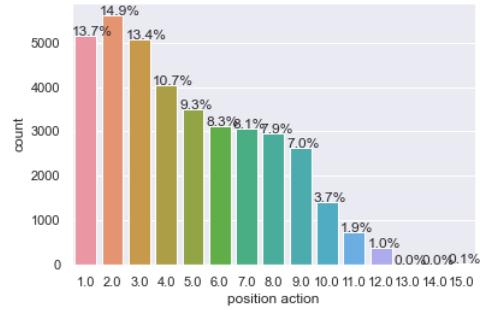
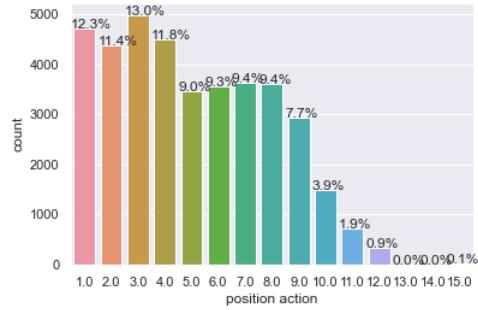
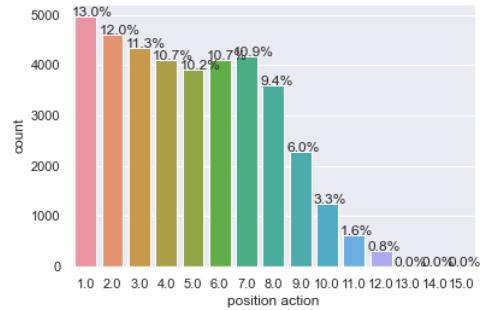
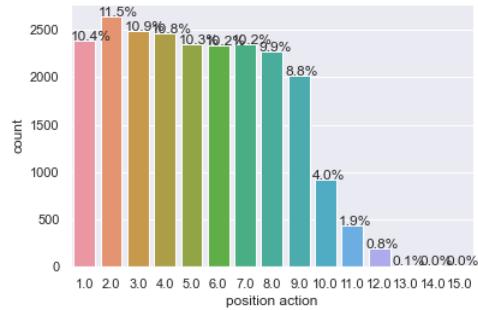
(10 iterations at 3/4)



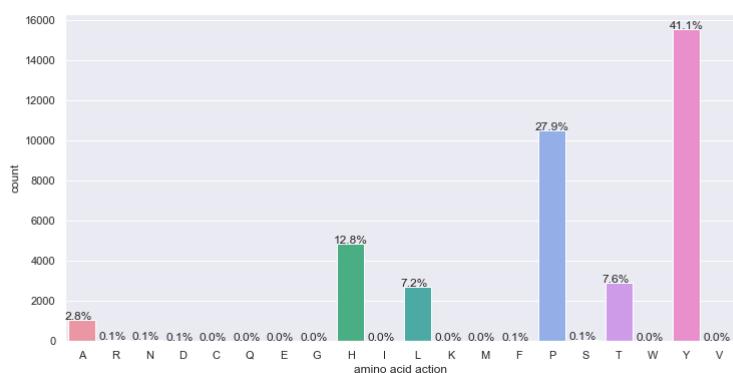
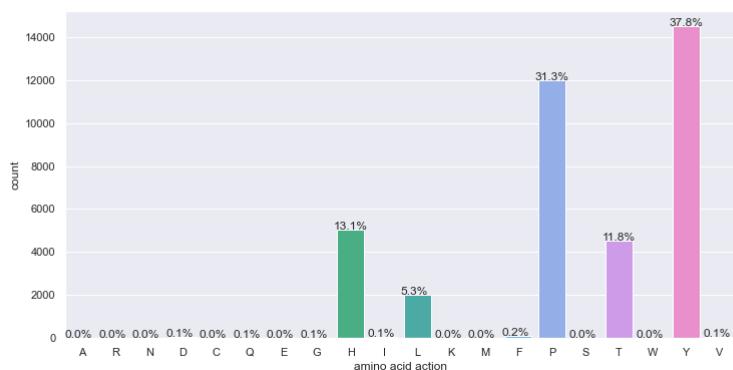
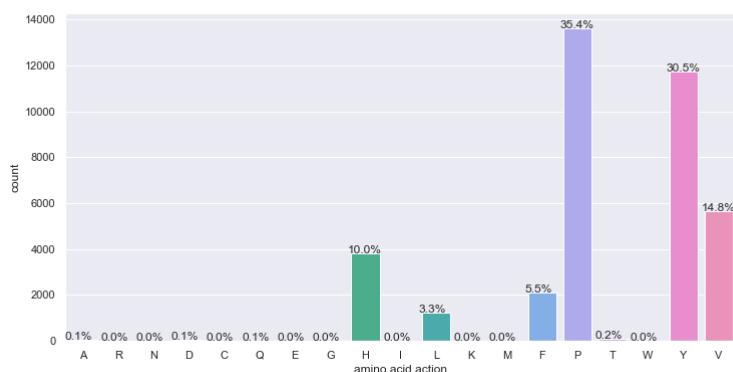
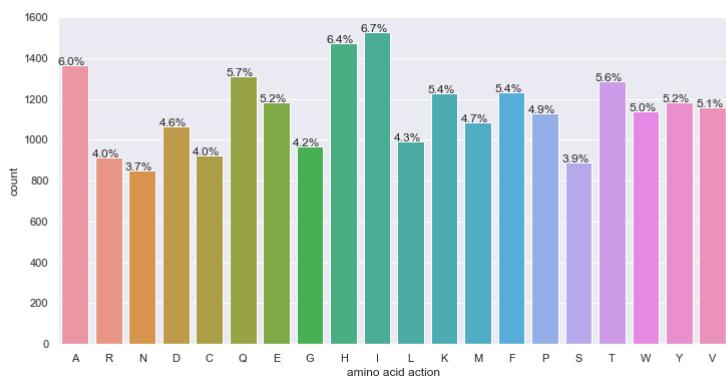
(last 10 iterations)

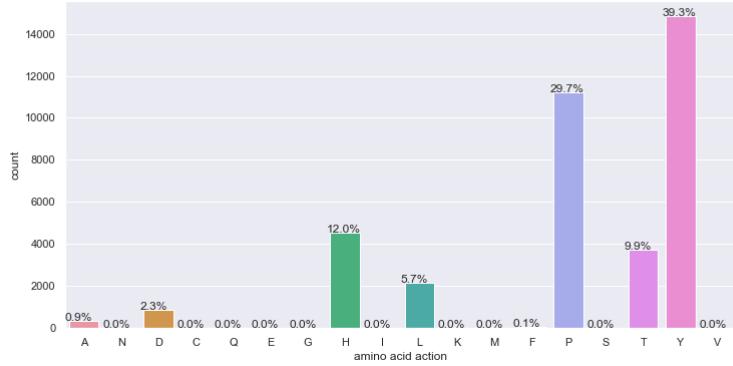


## 2. Position action analysis

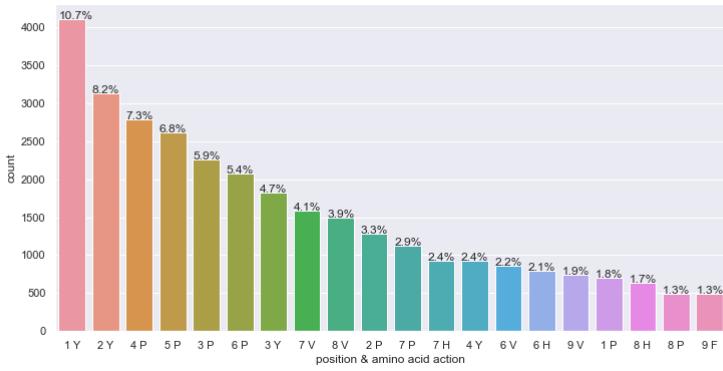
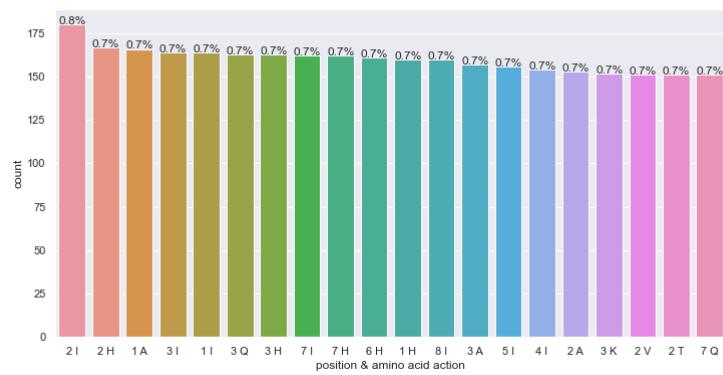


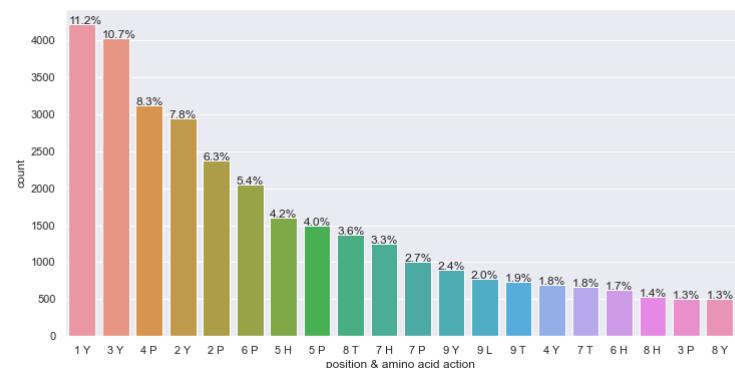
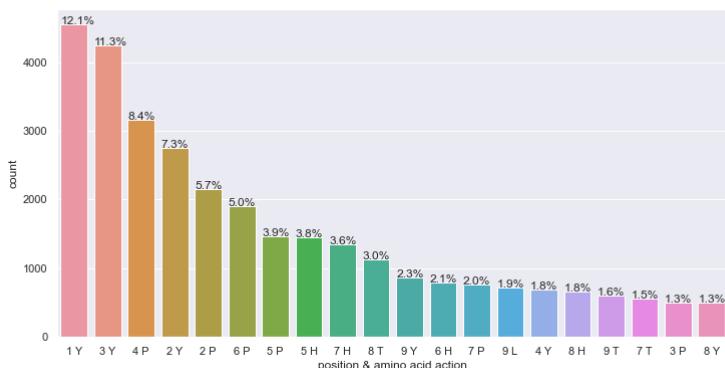
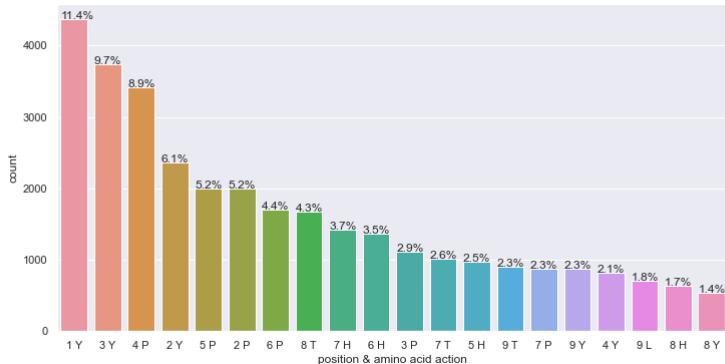
## 3. Amino action analysis



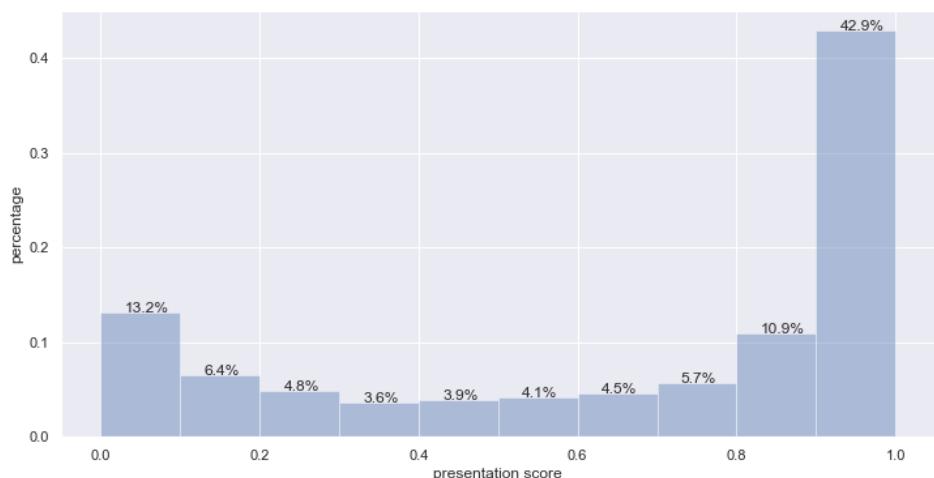


#### 4. Position & Amino action analysis

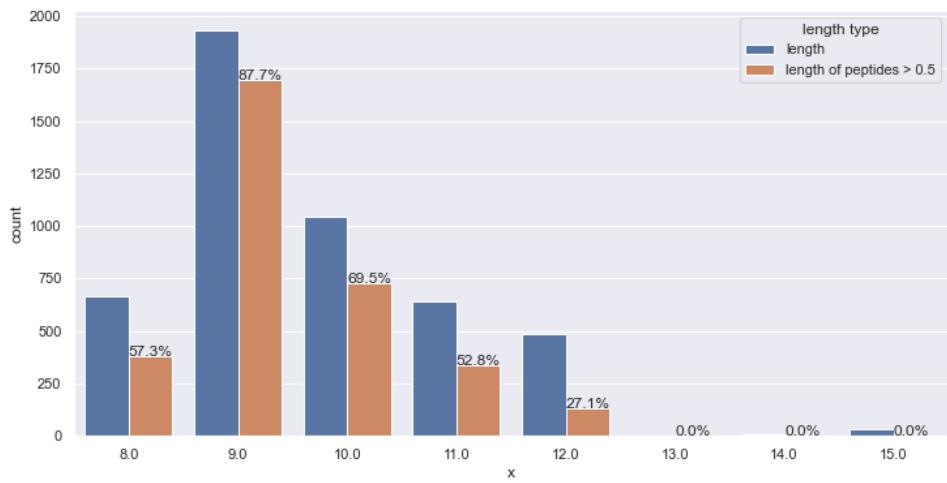




## 5. Reward analysis



## 6. Length analysis



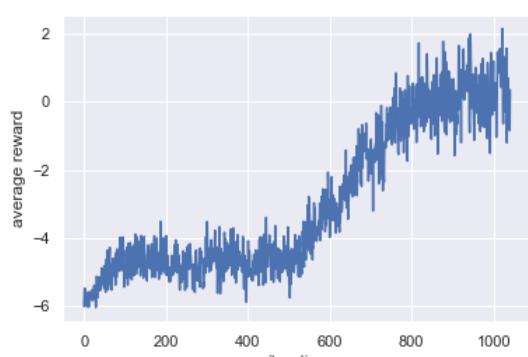
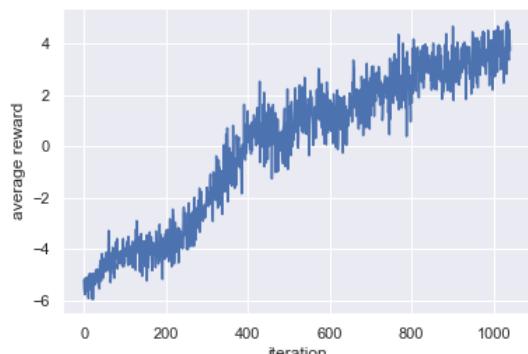
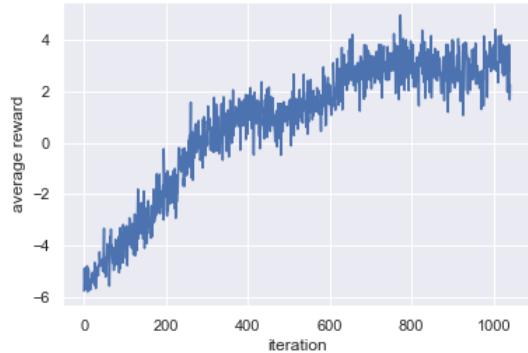
## 2. Prioritized Experience Replay

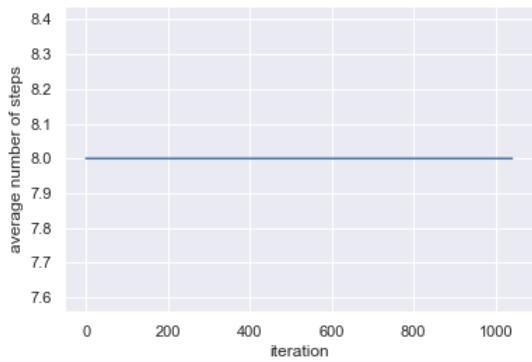
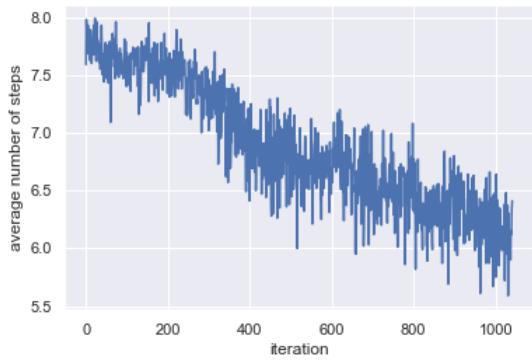
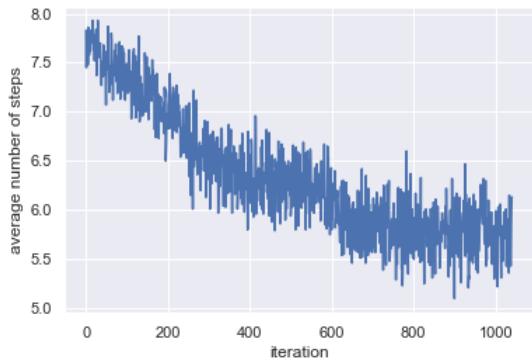
**Experiment: immediate reward**

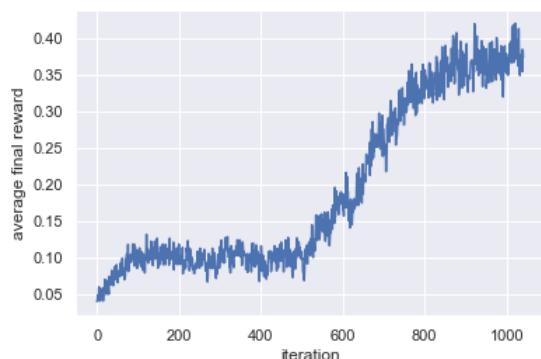
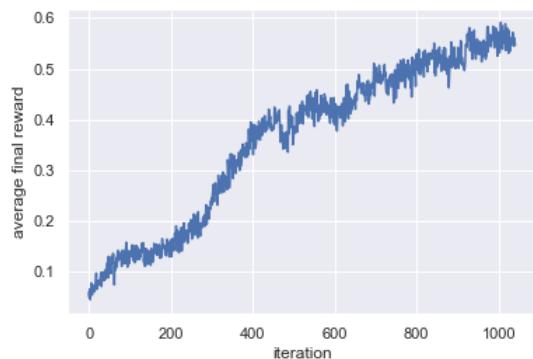
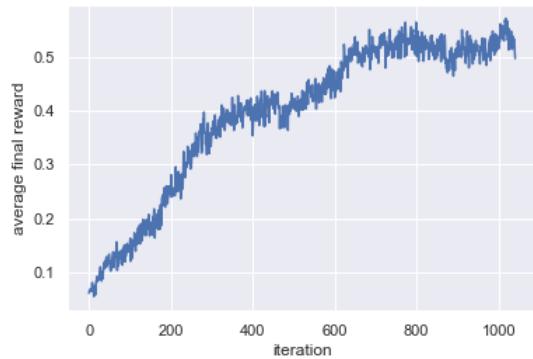
stop criteria = 0.6

stop criteria = 1.0

stop criteria = 0.75

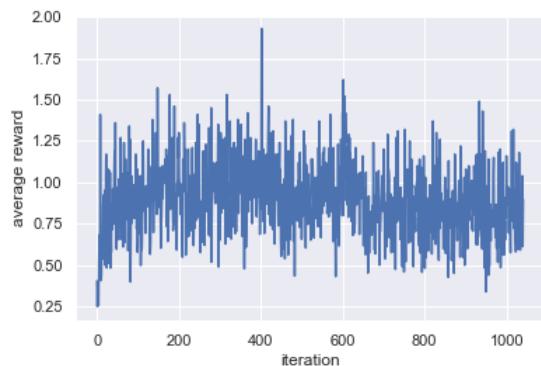
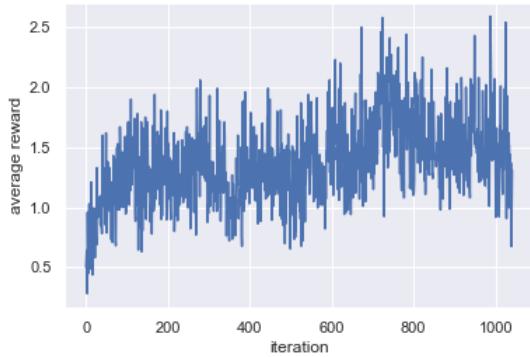
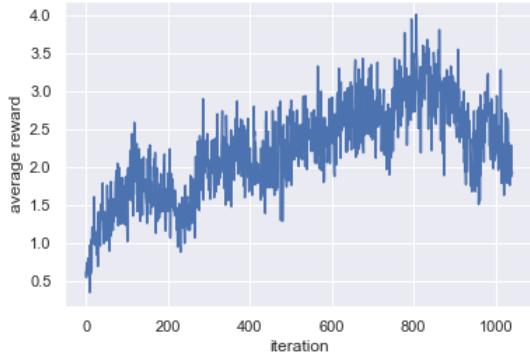


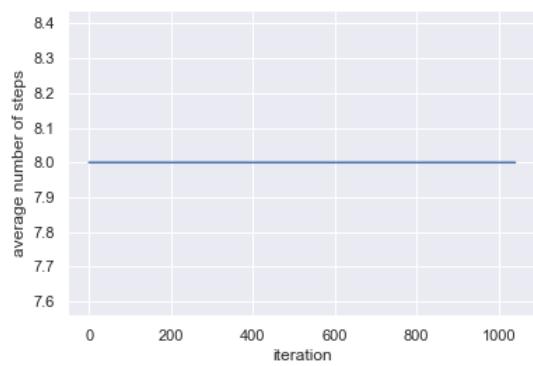
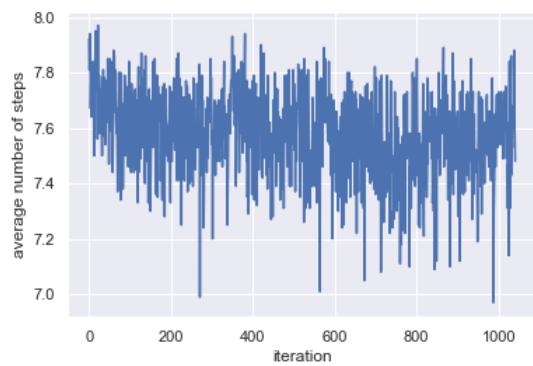
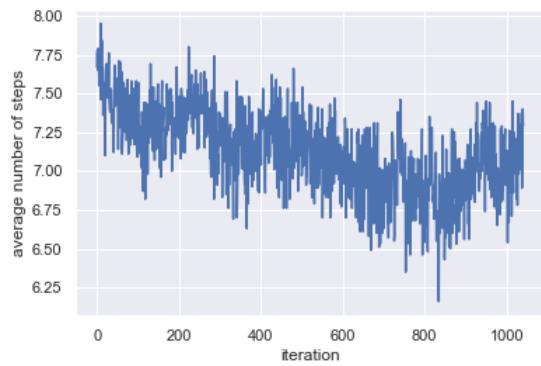


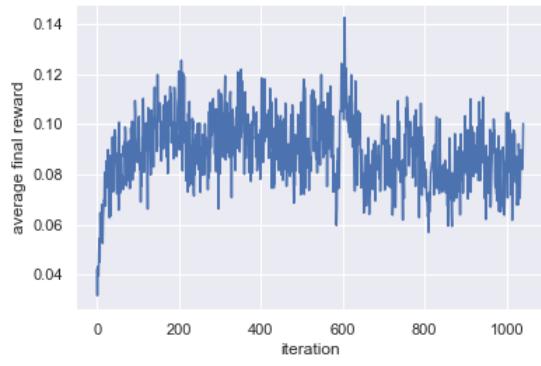
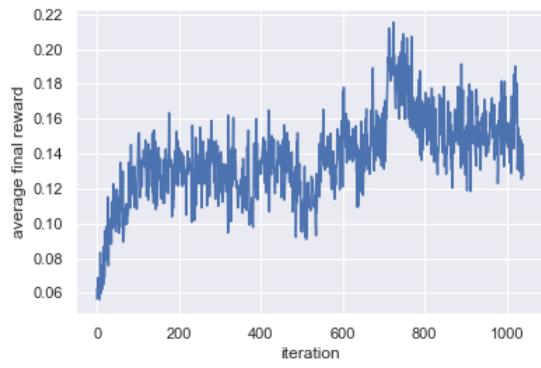
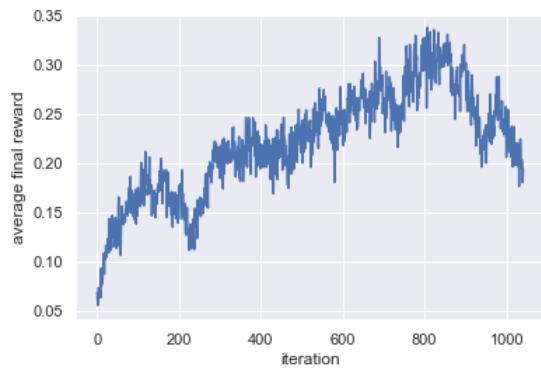


## Experiment: only final reward

stop criteria = 0.6 stop criteria = 0.75  
stop criteria = 1.0



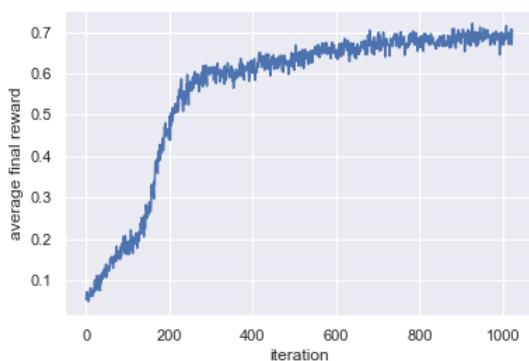
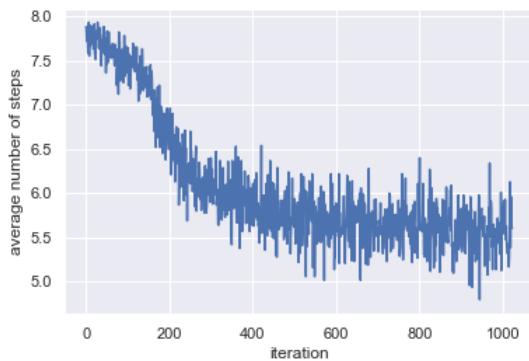
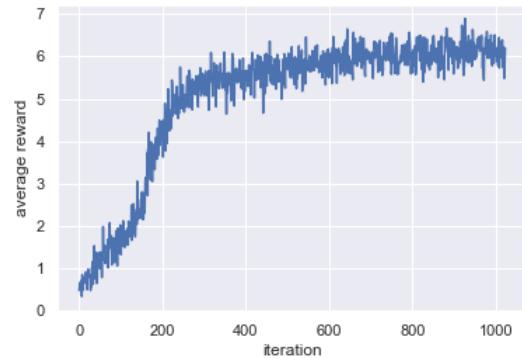




### 3. Positive rewards

**Experiment: stop criteria = 0.75**

Setting 1 (only final rewards; sample rate = 50% from IEDB and 50% from random)

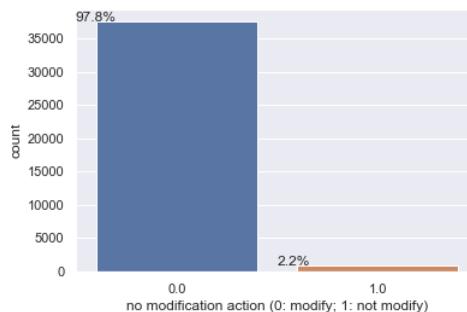
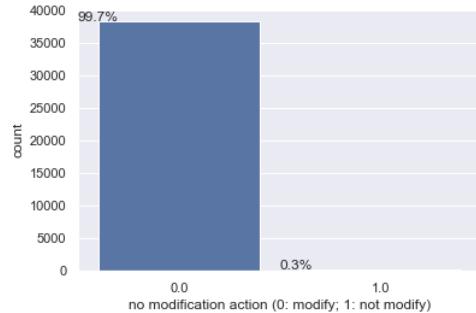
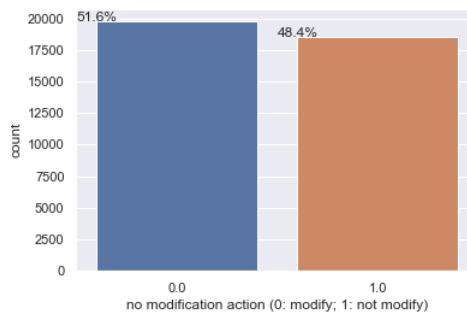


### 1. "No modification" action analysis

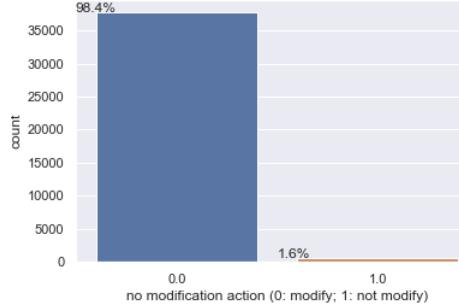
(first 10 iterations)  
the middle)

(10 iterations at 1/4)

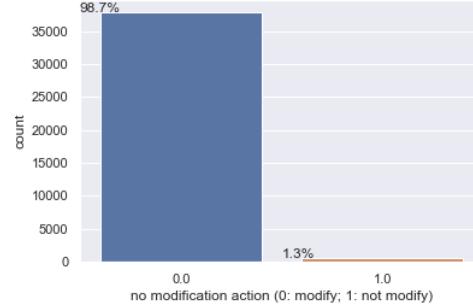
(10 iterations in



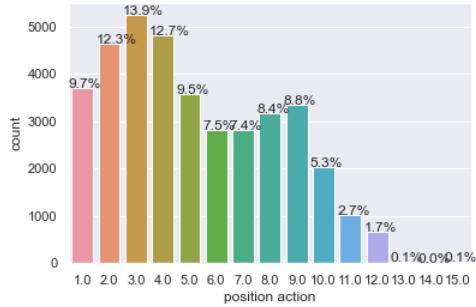
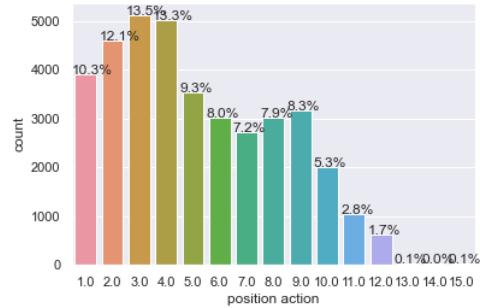
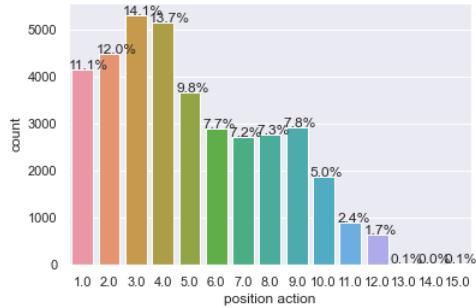
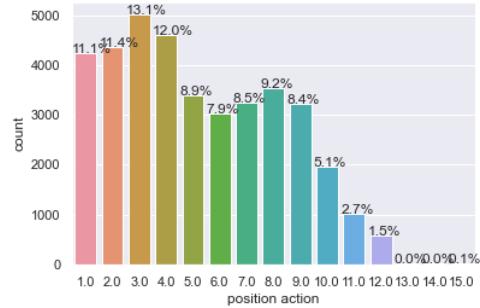
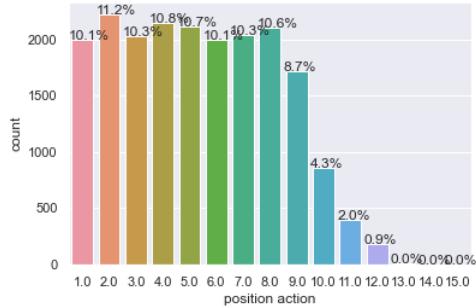
(10 iterations at 3/4)



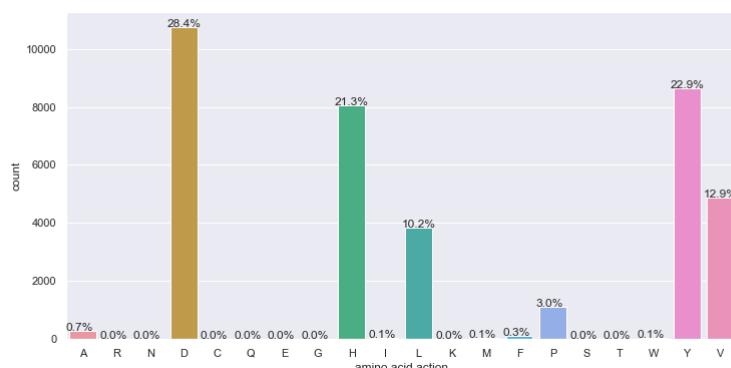
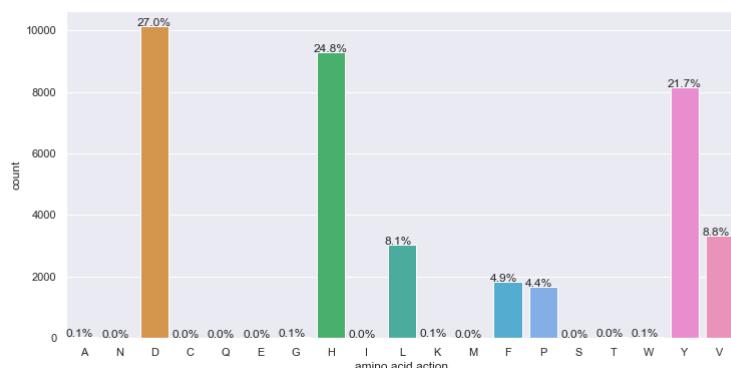
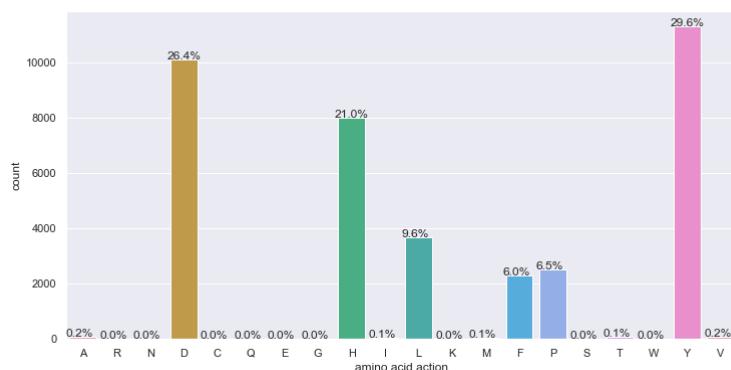
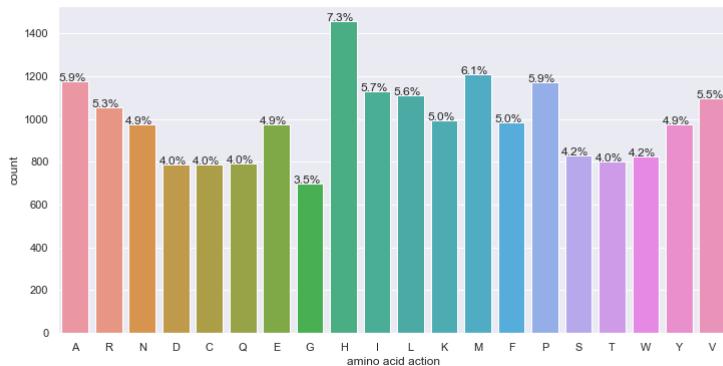
(last 10 iterations)

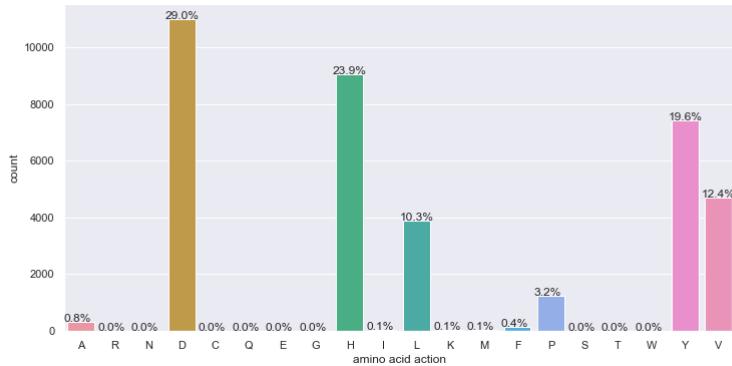


## 2. Position action analysis

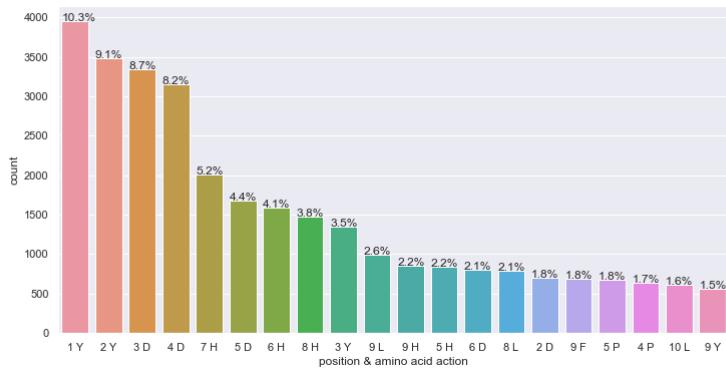
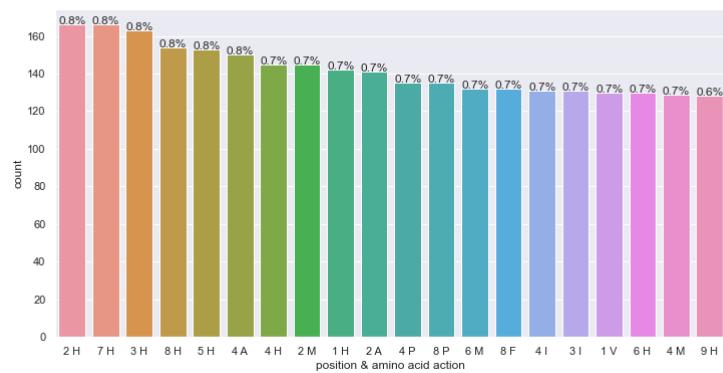


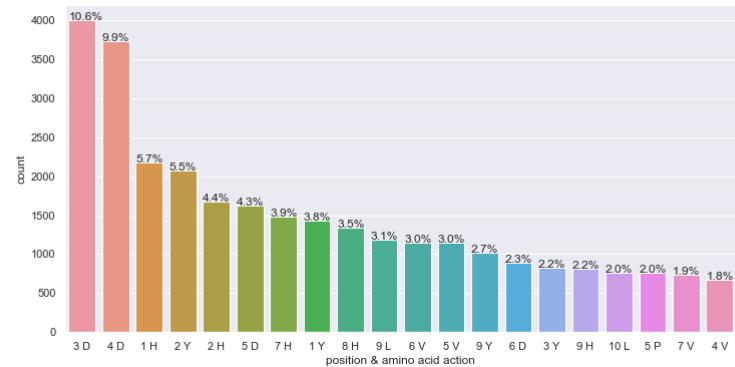
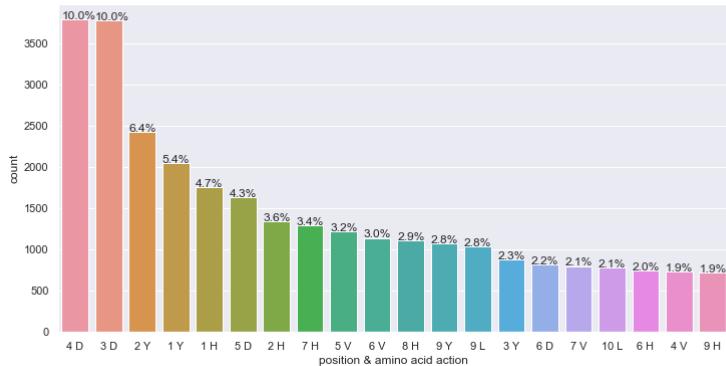
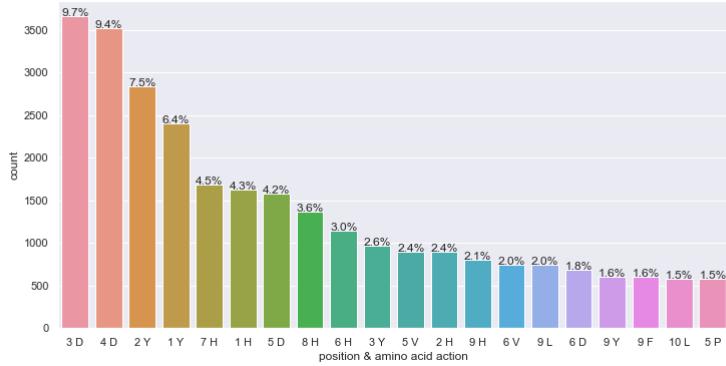
## 3. Amino action analysis



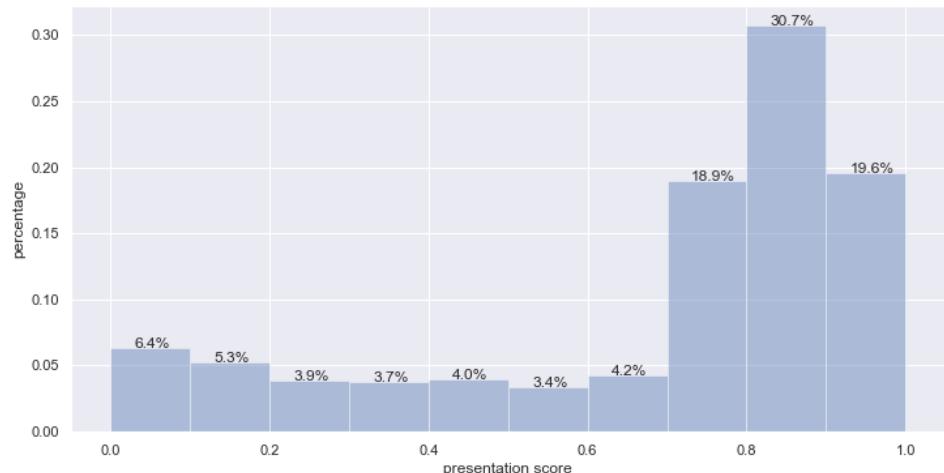


#### 4. Position & Amino action analysis

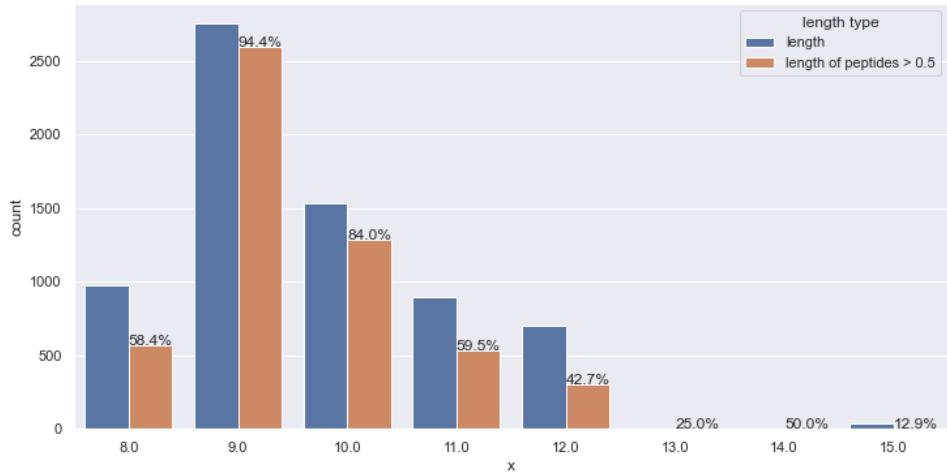




## 5. Reward analysis

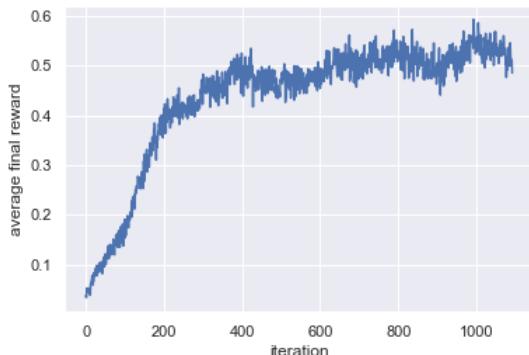
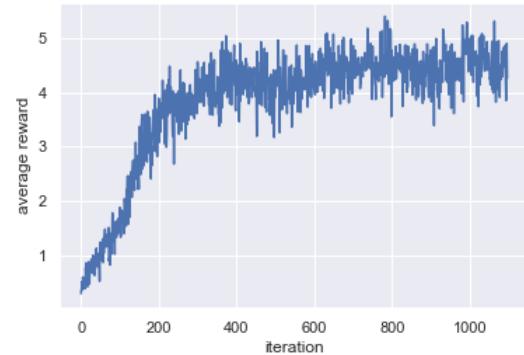


## 6. Length analysis



### Experiment: stop criteria = 1.0

Setting 1 (only final rewards; sample rate = 50% from IEDB and 50% from random)

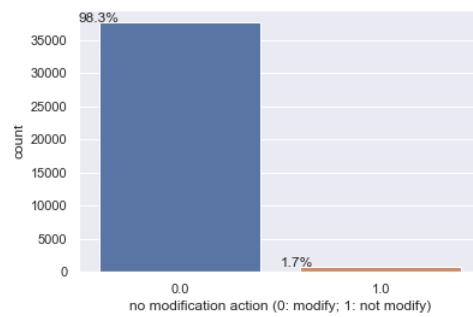
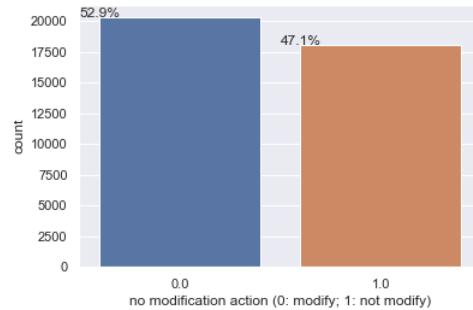


#### 1. "No modification" action analysis

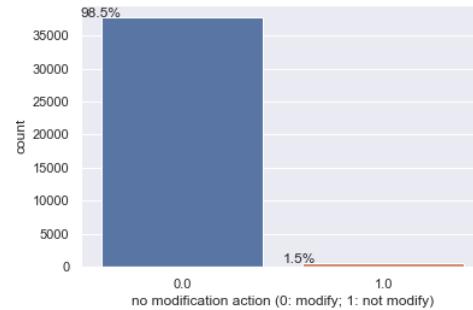
(first 10 iterations)  
the middle)

(10 iterations at 1/4)

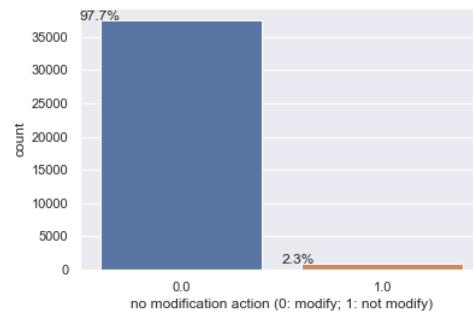
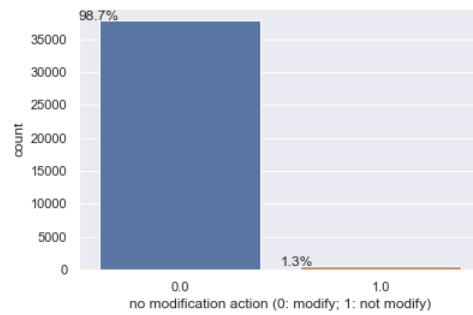
(10 iterations in



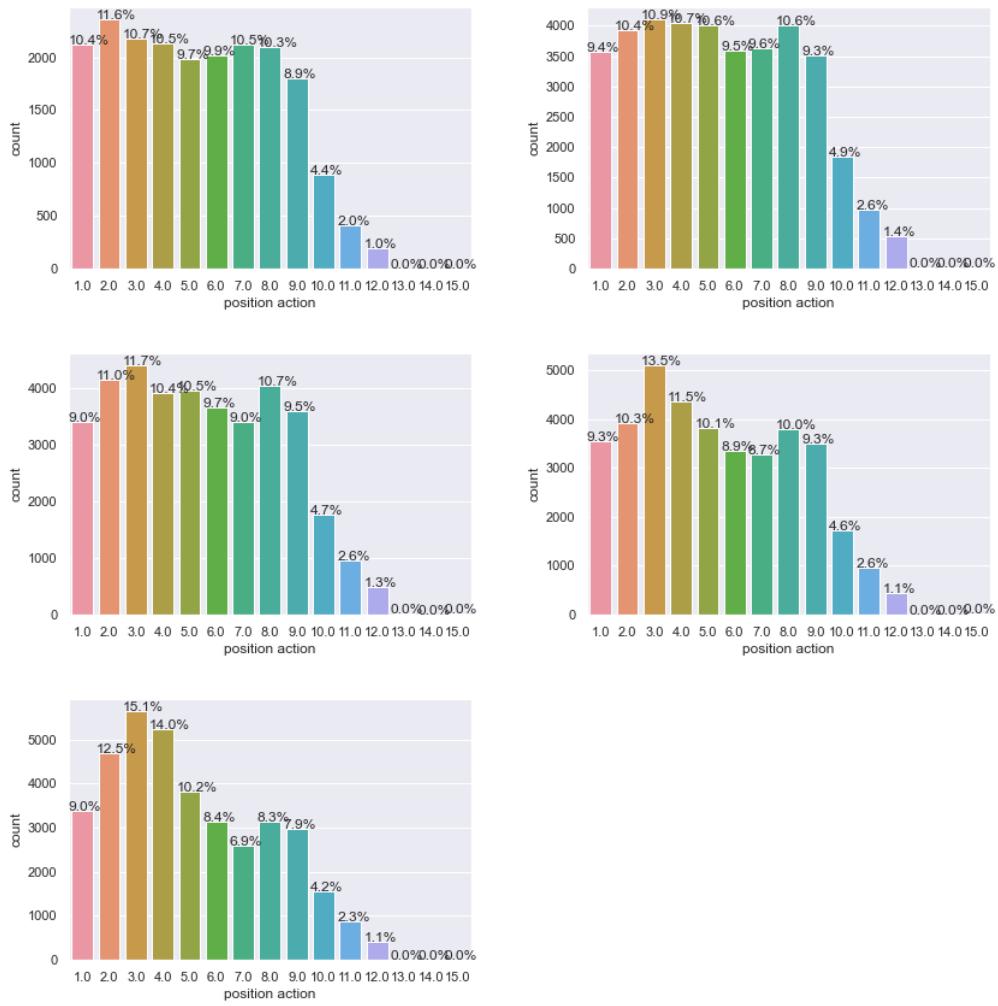
(10 iterations at 3/4)



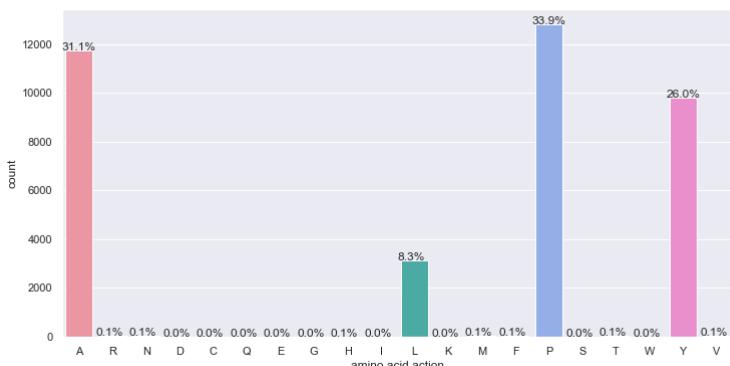
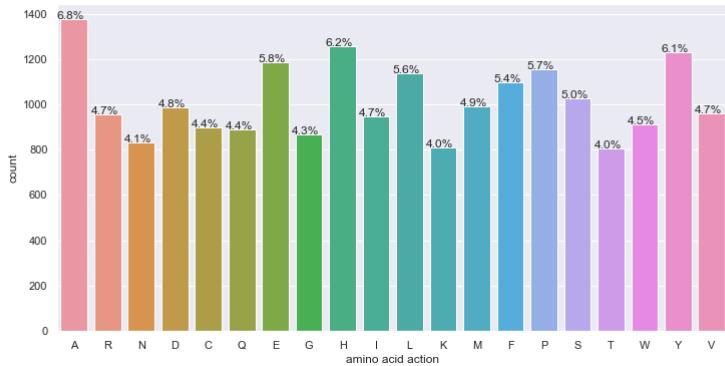
(last 10 iterations)

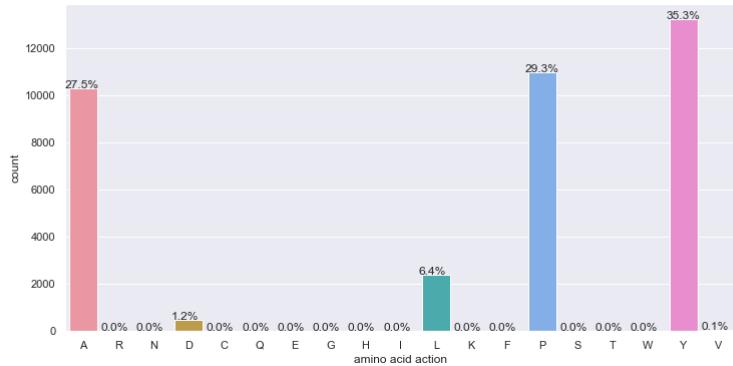
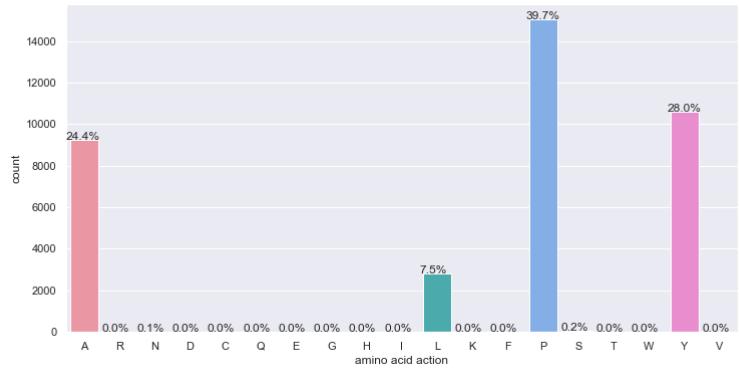
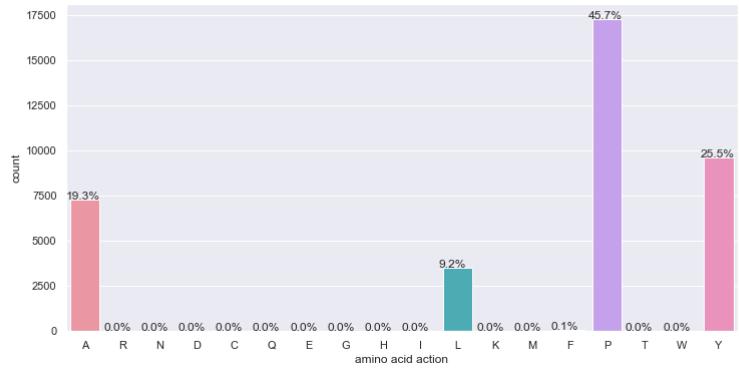


## 2. Position action analysis

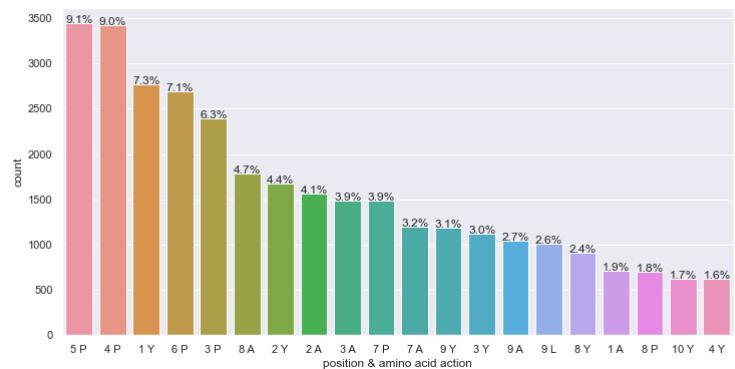
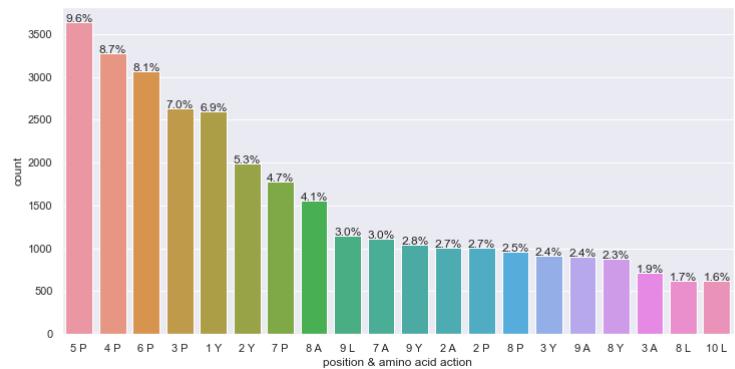
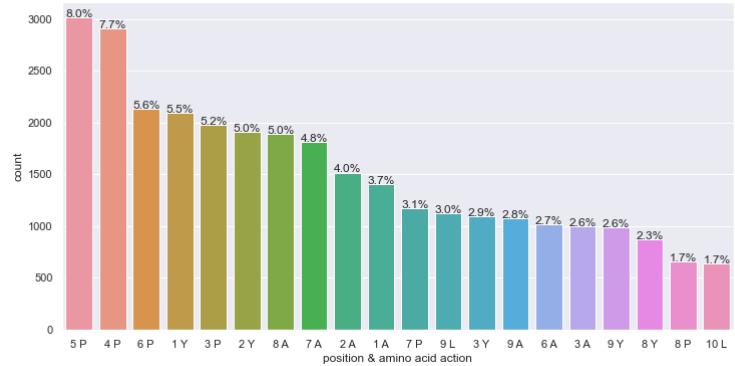
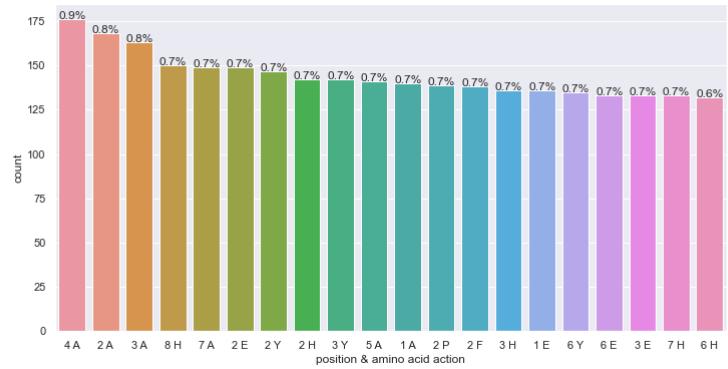


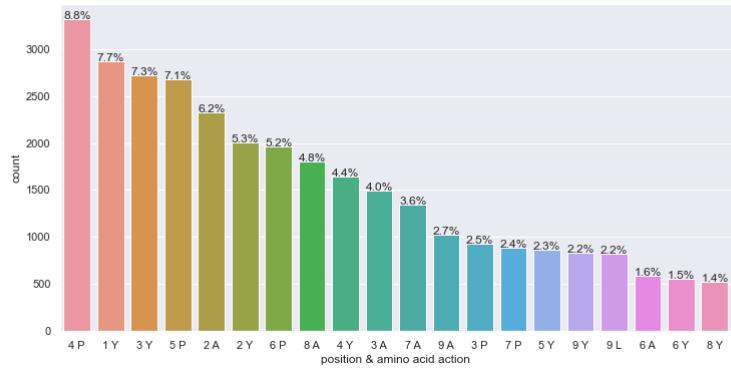
### 3. Amino action analysis



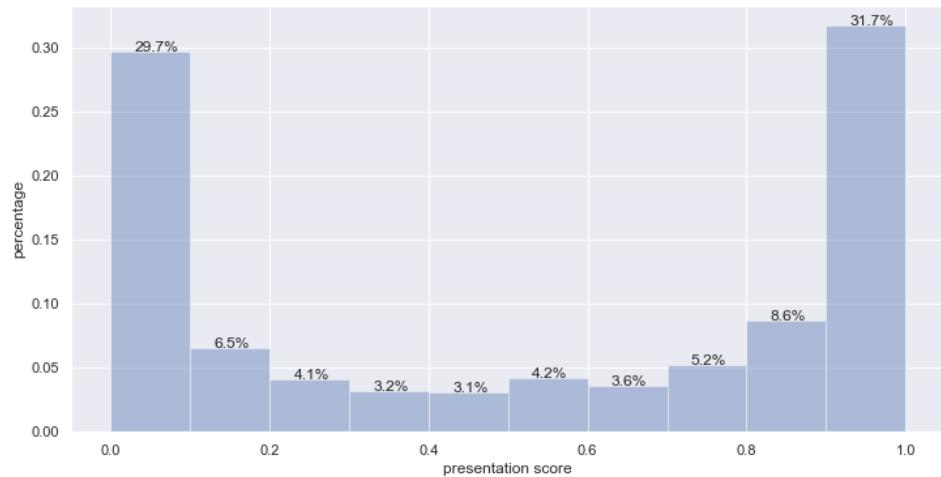


#### 4. Position & Amino action analysis

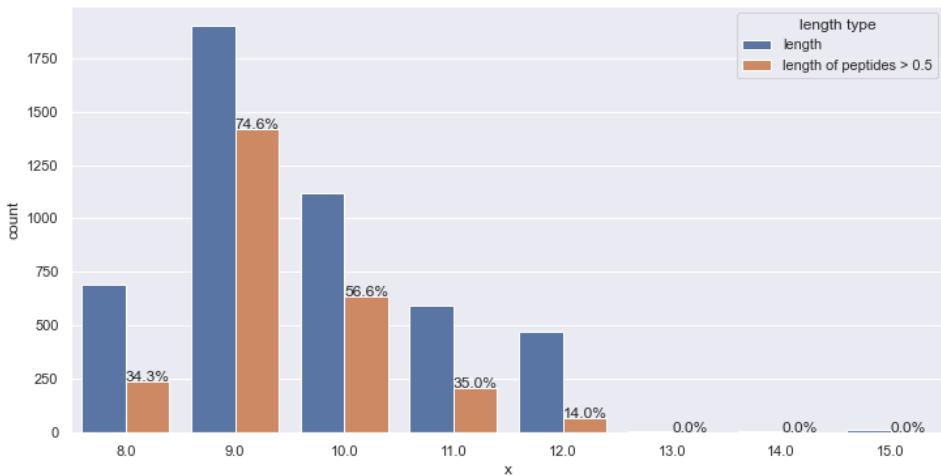




## 5. Reward analysis



## 6. Length analysis



## Update (6-9)

Trick:

Add the number of steps as a feature to predict the actions  $a \sim \pi(a|s, t)$  and the future rewards  $Q(s, a)$ . This trick is from the paper below <sup>1</sup>.

Next step:

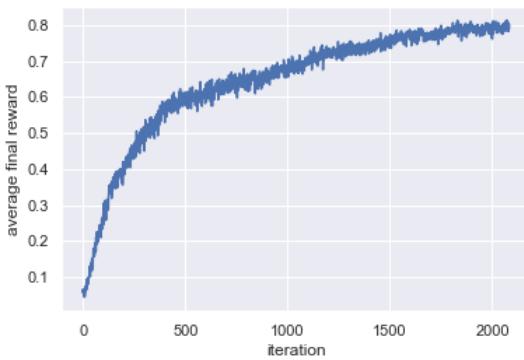
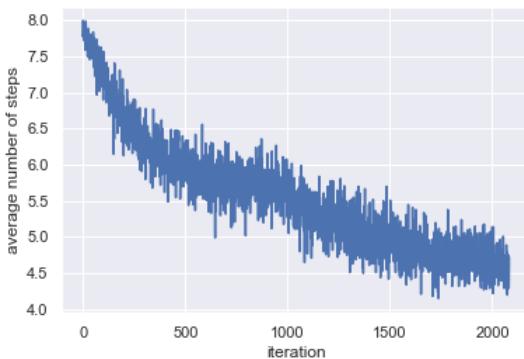
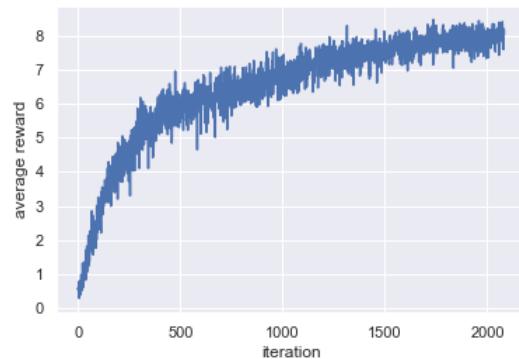
1. Add human peptide constraint when the data is available.
  2. Compare the baselines (i.e., MCTS, motif, random) with the learned policy.
  3. Write the paper
  4. find alleles such as HLA-C with little training data, generate peptides for  
find alleles such as HLA-A with

## Idea:

- ### 1. Combine MCTS with generative classifier

## Experiment: stop criteria = 0.75

Setting 1 (only final rewards; sample rate = 50% from IEDB and 50% from random)

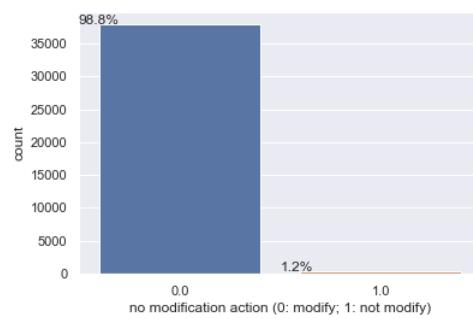
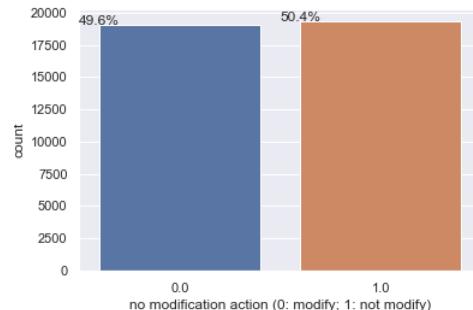


- ## 1. "No modification" action analysis

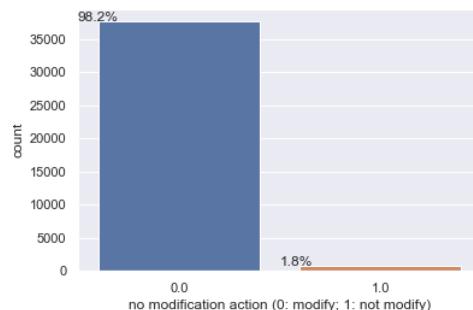
(first 10 iterations)  
the middle)

(10 iterations at 1/4)

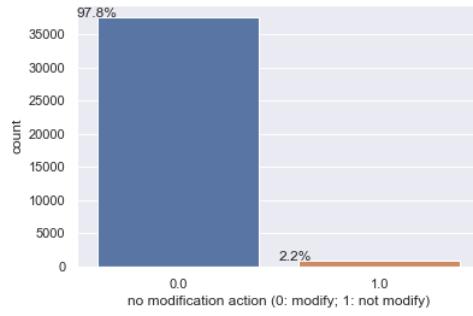
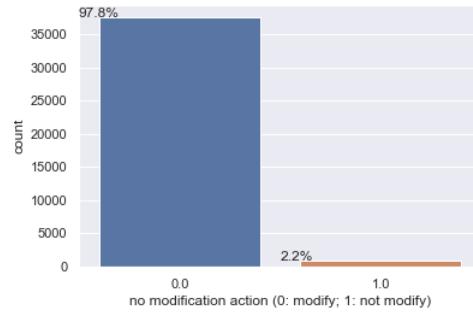
(10 iterations in



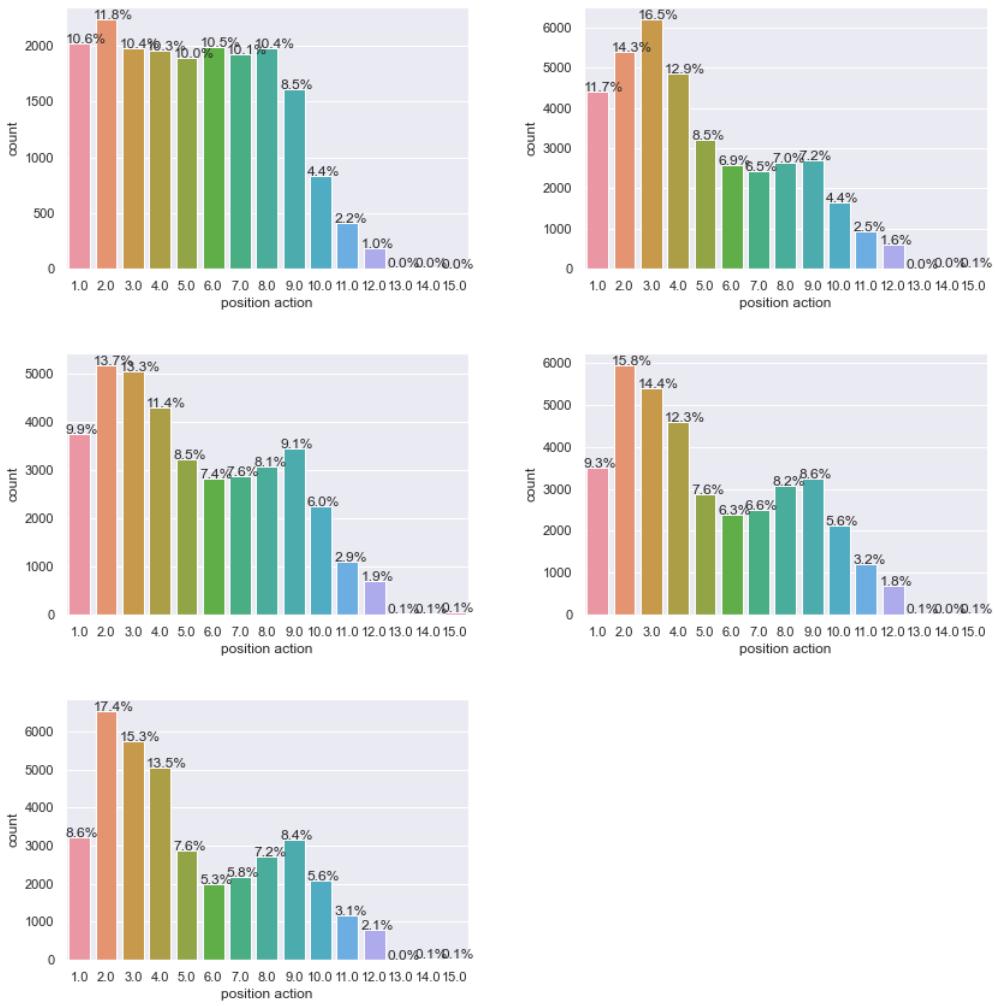
(10 iterations at 3/4)



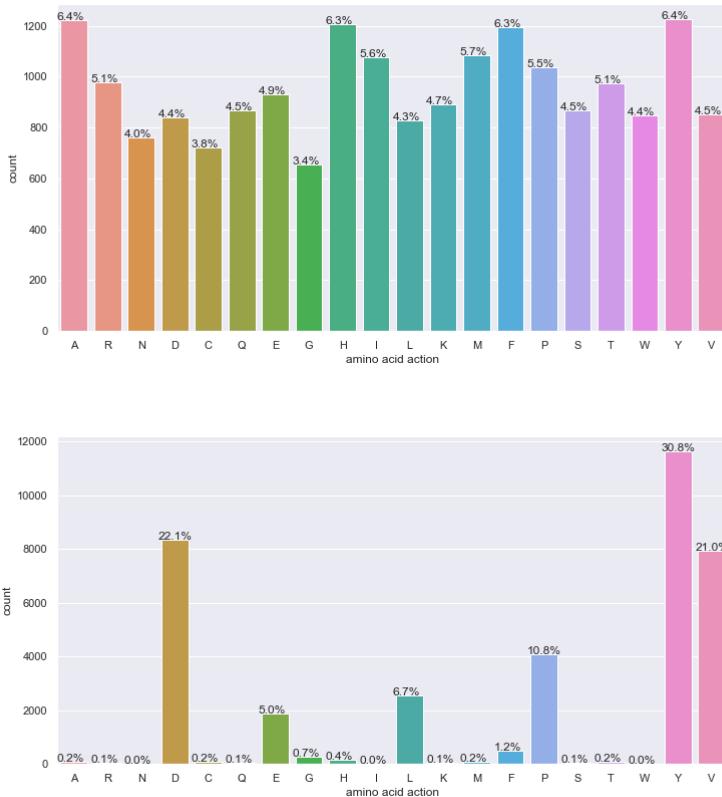
(last 10 iterations)

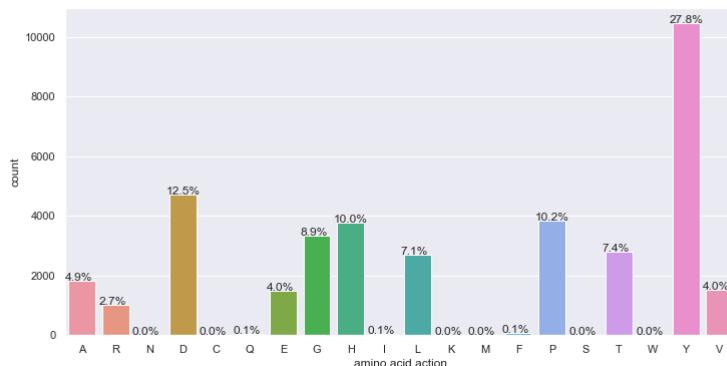
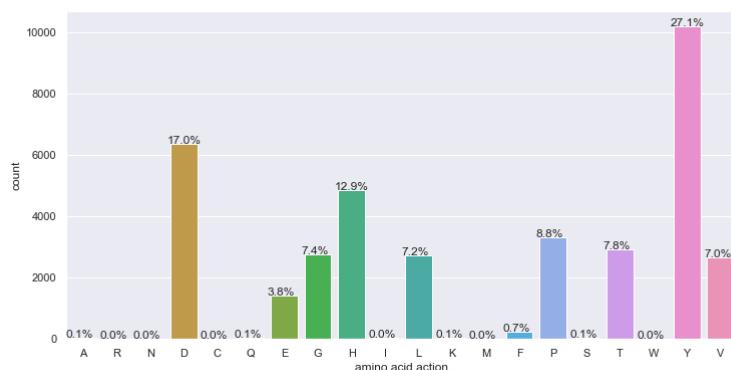
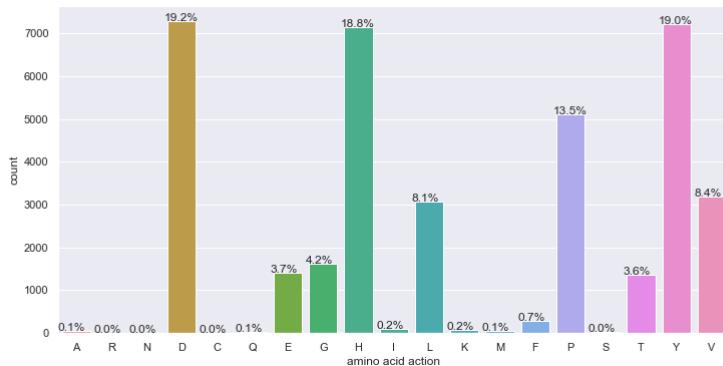


## 2. Position action analysis

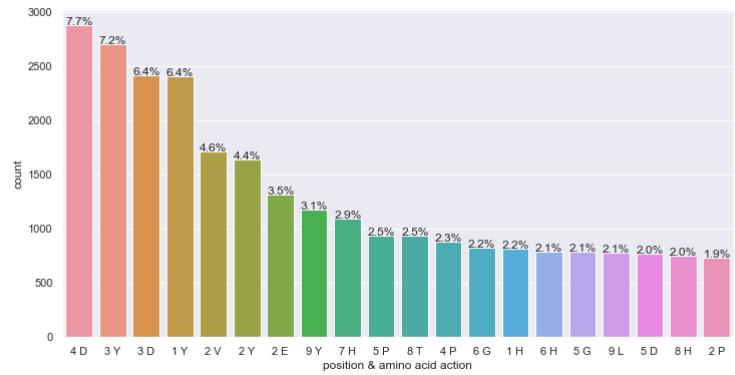
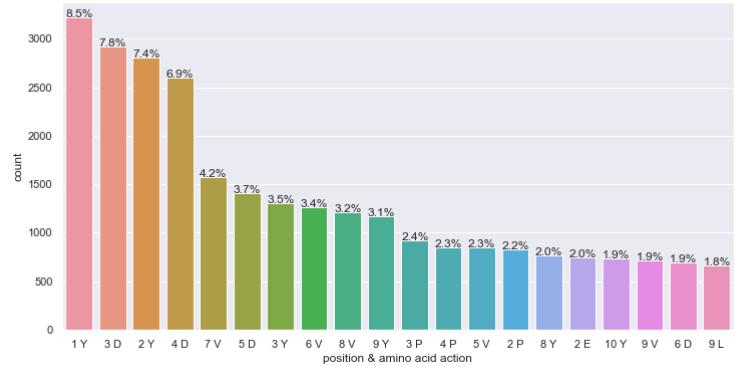
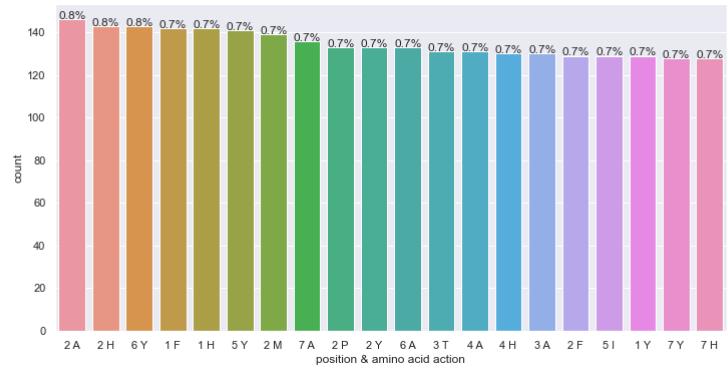


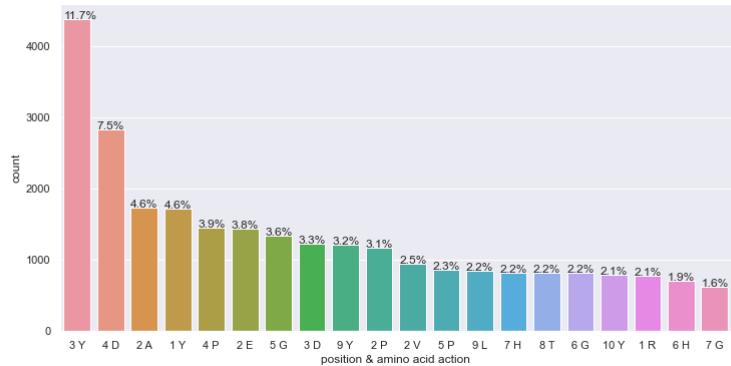
### 3. Amino action analysis



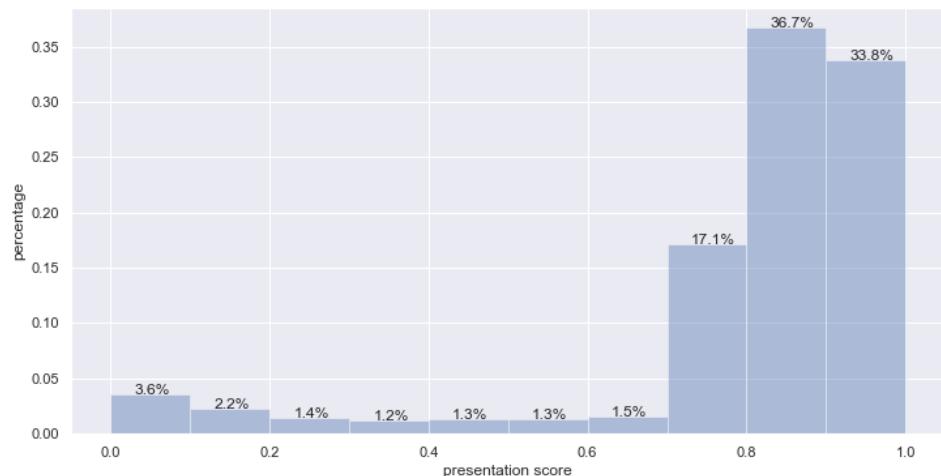


#### 4. Position & Amino action analysis

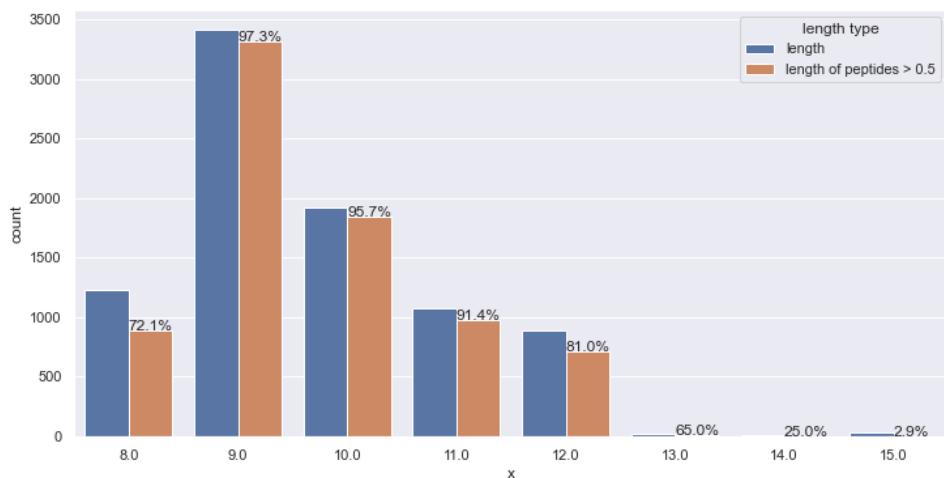




## 5. Reward analysis

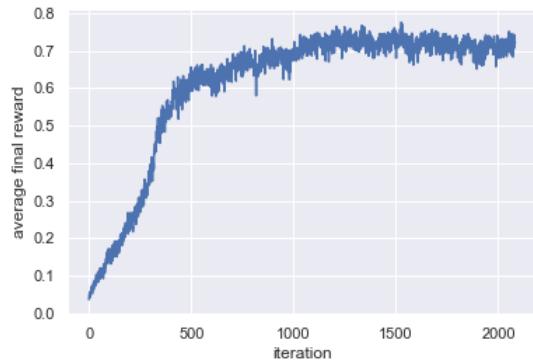
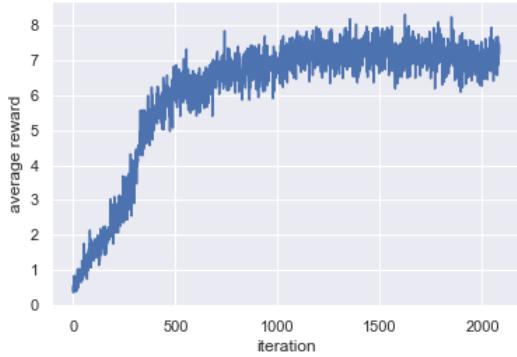


## 6. Length analysis



## Experiment: no stopping criteria

Setting 1 (only final rewards; sample rate = 50% from IEDB and 50% from random)

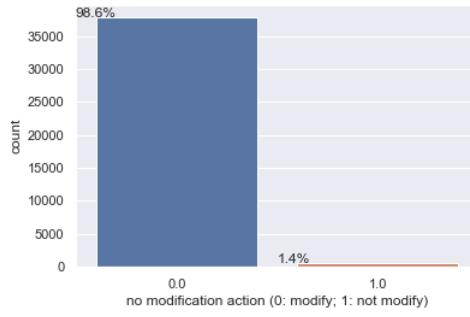
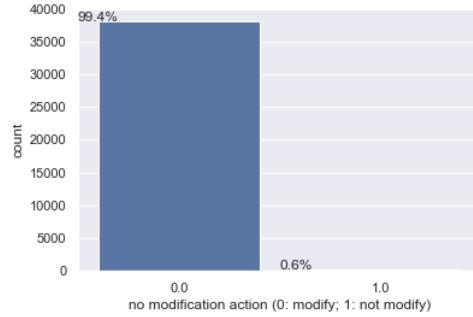
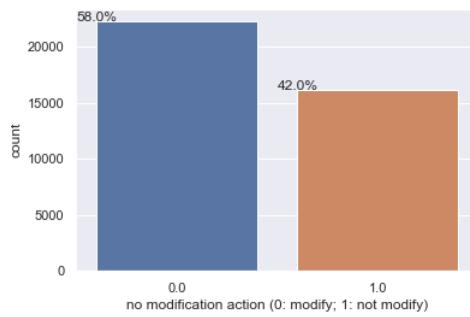


### 1. "No modification" action analysis

(first 10 iterations)  
the middle)

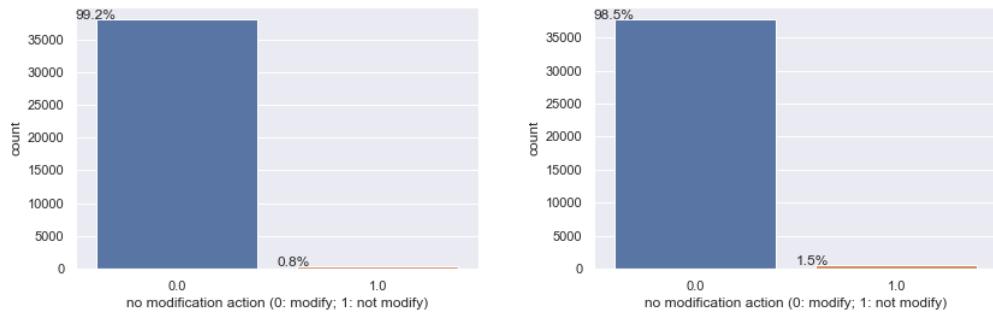
(10 iterations at 1/4)

(10 iterations in

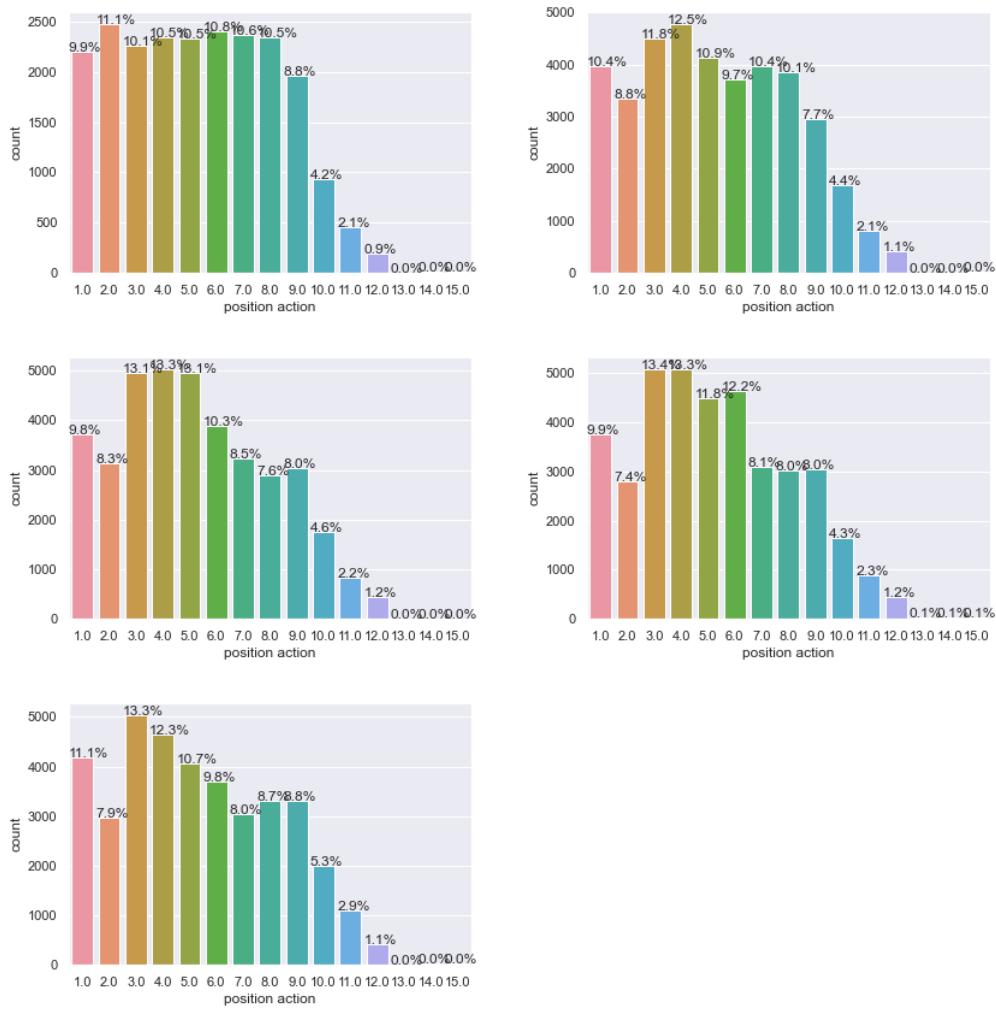


(10 iterations at 3/4)

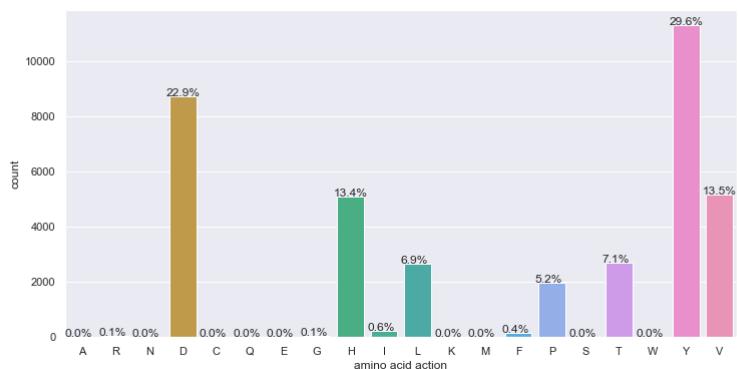
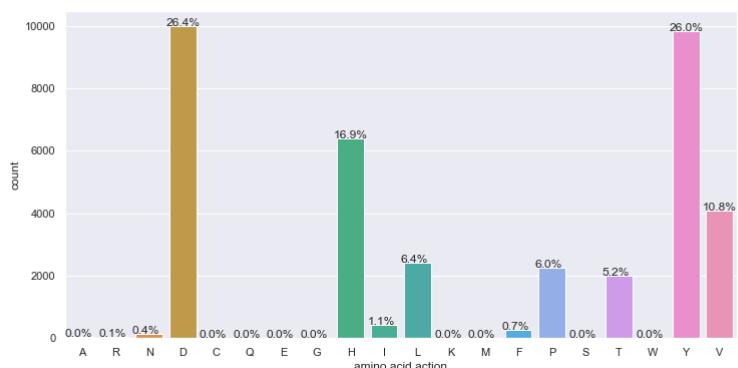
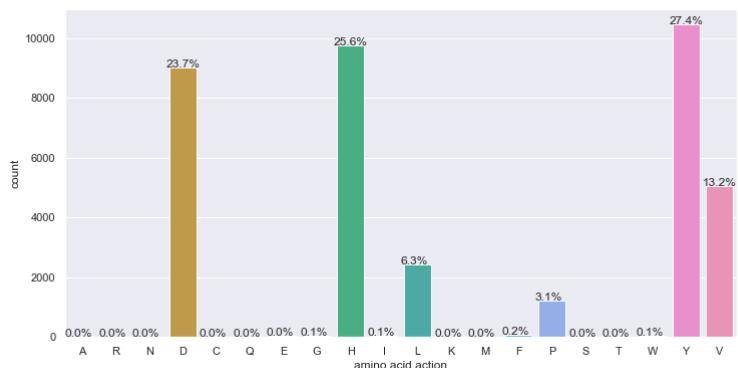
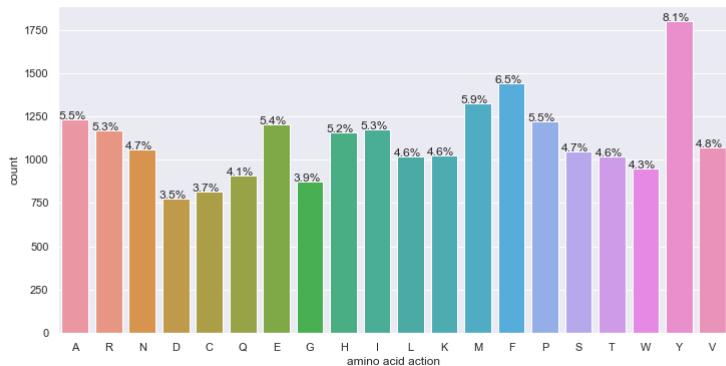
(last 10 iterations)

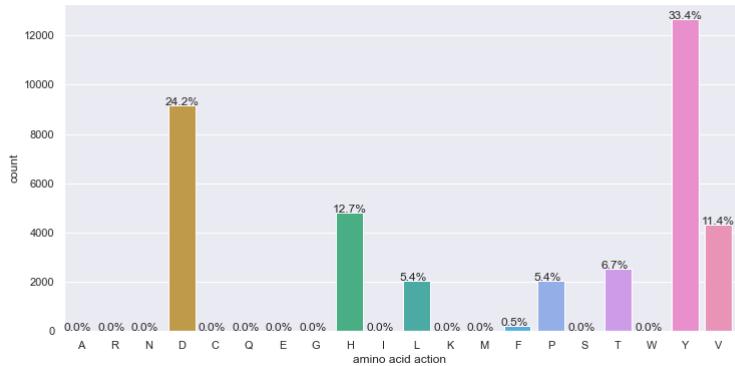


## 2. Position action analysis

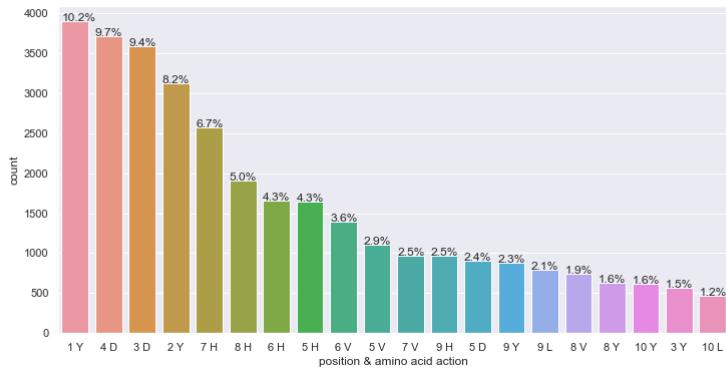
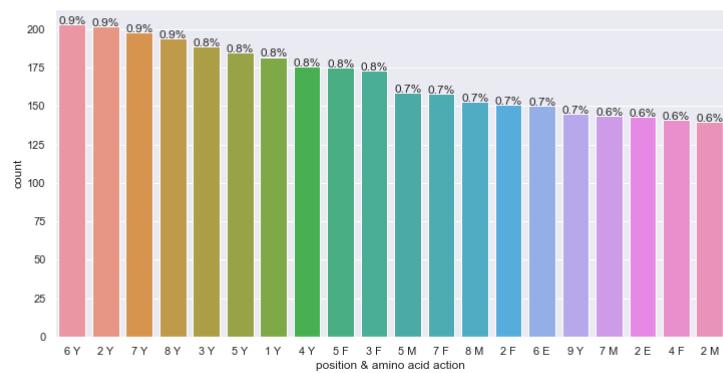


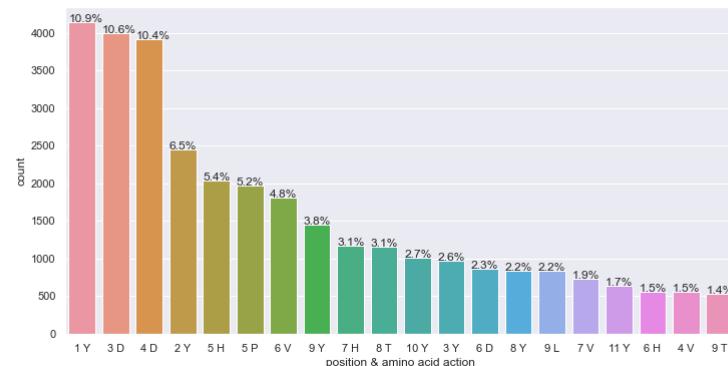
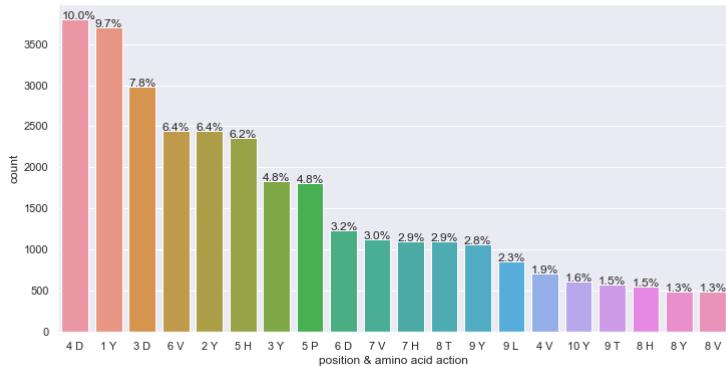
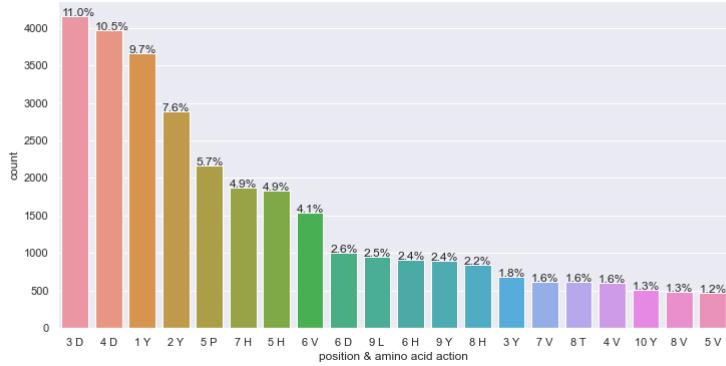
## 3. Amino action analysis



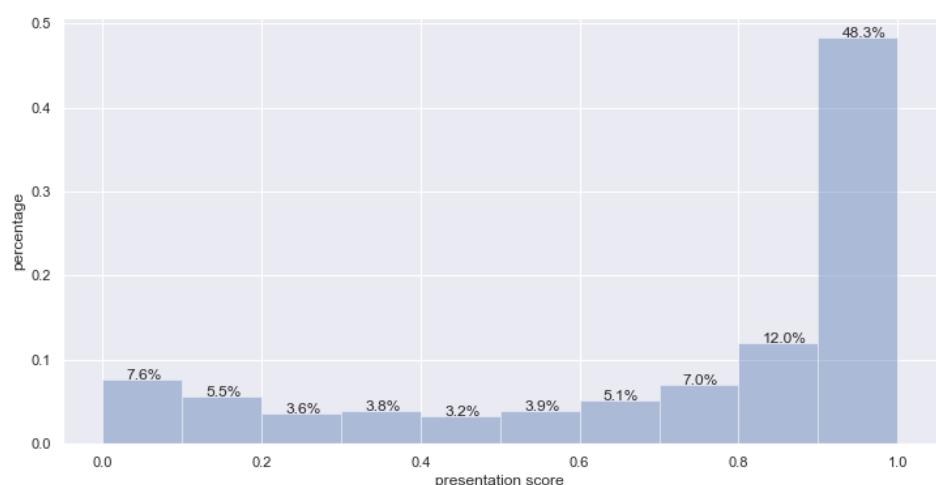


#### 4. Position & Amino action analysis

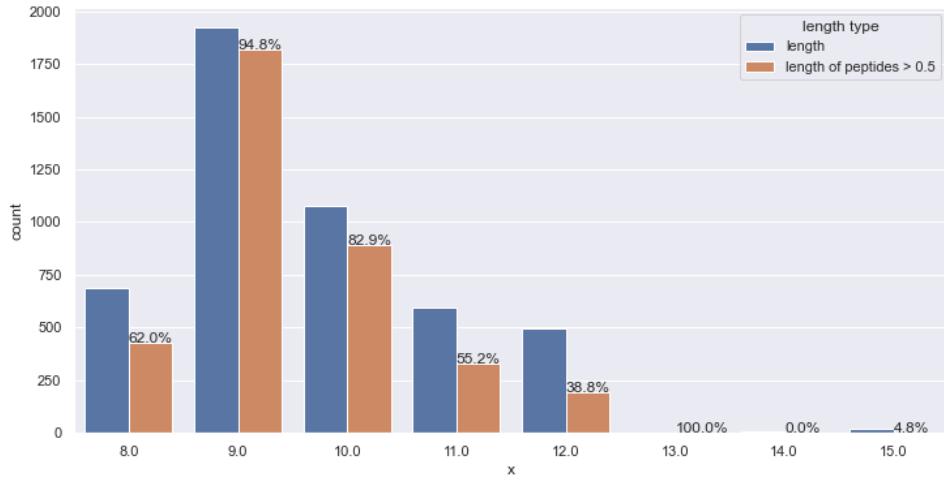




## 5. Reward analysis



## 6. Length analysis



## Baseline 1: MCTS

important hyperparameters:

1. number of rollout for each allele
2. a criteria used to output good peptides
3. a tradeoff between exploitation and exploration within the UCB1 formula.

Algorithm detail:

- Each state is a sequence of amino acids. The root state is the empty sequence.
- For each state, we can pick one amino acid and get the next state by appending it to the sequence of this state.
  - We pick the amino acid that has the maximum value with the UCB1 formula

$$a_k = \arg \max_a \frac{\sum_i r_i}{N(s,a)} + c_{puct} * \sqrt{\frac{2N(s)}{1+N(s,a)}}$$

where  $c_{puct}$  is the tradeoff between exploitation and exploration;  $N(s, a)$  is the visit count of state-action pair and  $N(s)$  is the visit count of state  $s$ ;  $r_i$  is the maximum representation score of  $i$ -th rollout going through the state-action pair.

- The new state will be added to the tree.
- For each state  $s$  with sequence of length  $l \geq 8$ , we use *MHCflurry* to evaluate the presentation score of that sequence  $p(s)$ .
- The terminal state has the sequence of length  $l = 15$ .
- After reaching the terminal state, we backpropagate the presentation score  $r_i^{t-1} = \max(r_i^t, p(s_i^{t-1}))$  in which  $r_i^T = p(s_i^T)$ .

## Use the motif position to reorder the sequence

Question:

1. How to specify the number of rollouts?

In the RL model with stop criteria = 0.75, we have 3,688 alleles in total. Each allele is associated with about 400 episodes.

It could be unfair to set the number of rollouts to be 400. But I guess 5000 rollouts could be a good choice?

## 2. How to compare the MCTS and other baselines with the RL model?

Compare the space of generated peptides, e.g., t-sne plot.

For each allele, run baselines and our RL model for 1000 rounds.

Some metrics: number of optimized peptides (> 0.75); percentage of unique peptides; average pairwise edit distances between optimized peptides; overlap with training data;

*t-sne* plot

Some preliminary results:

(10,000 rollouts for each allele,  $c_{puct} = 0.5$ )

time cost for each allele: about 3 hours

```
Threshold = 0.5:  
Number of peptides, allele name  
1325      HLA-B*57:01  
896       HLA-A*29:02  
610       HLA-A*24:02  
544       HLA-B*35:01  
534       HLA-B*44:02  
464       HLA-A*02:03  
435       HLA-A*02:01  
432       HLA-A*31:01  
420       HLA-A*02:07  
364       HLA-A*01:01  
362       HLA-A*03:01  
282       HLA-A*68:02  
175       HLA-B*51:01  
164       HLA-B*54:01
```

Threshold = 0.75:

```
554      HLA-B57:01  
412      HLA-A29:02  
280      HLA-A24:02  
274      HLA-B35:01  
266      HLA-B44:02  
258      HLA-A02:03  
223      HLA-A02:01  
182      HLA-A02:07  
177      HLA-A31:01  
153      HLA-A03:01  
138      HLA-A01:01  
127      HLA-A68:02
```

## Baseline 2: Position Weight Matrix (Motif)

Similar to the prior distribution we derived in the previous experiments.

## Baseline 3: Random Policy

## Update (6-23)

- An example of comparison among methods with different experimental settings:

Model	Max. Penalized logP	Model
GVAE + BO (Kusner et al., 2017) <sup>1</sup>	$2.87 \pm 0.06$	VAE
SD-VAE (Dai et al., 2018) <sup>1</sup>	$3.50 \pm 0.44$	VAE
CVAE + BO (Gómez-Bombarelli et al., 2018) <sup>2</sup>	$4.85 \pm 0.17$	VAE
ORGAN (Guimaraes et al., 2017) <sup>1</sup>	$3.52 \pm 0.08$	GAN
JT-VAE (Jin et al., 2018a) <sup>1</sup>	$4.90 \pm 0.33$	VAE
ChemTS (Yang et al., 2017)	$5.6 \pm 0.5$	RNN
GCPN (You et al., 2018) <sup>1</sup>	$7.87 \pm 0.07$	PN + GAN
Random SELFIES	$6.19 \pm 0.63$	Random Search
GB-GA (Jensen, 2019) <sup>3</sup>	$7.4 \pm 0.9$	GA
GB-GA (Jensen, 2019) <sup>4</sup>	$15.76 \pm 5.71$	GA
GA (here)	$12.61 \pm 0.81$	GA
GA + D (here)	$13.31 \pm 0.63$	GA + DNN
(GA + D(t)) (here) <sup>5</sup>	$20.72 \pm 3.14$	GA + DNN

<sup>1</sup> average of three best molecules after performing Bayesian optimization on the latent representation of a trained VAE

<sup>2</sup> two best molecules of a single run

<sup>3</sup> averaged over 10 runs with 20 molecules per generation with a molecular weight (excl. hydrogen) smaller than  $39.15 \pm 3.50$  g/mol run for 50 generations

<sup>4</sup> averaged over 10 runs with 500 molecules up to 81 characters per generation and 100 generations

<sup>5</sup> averaged over 5 runs with 1000 generations each

- The baselines used in the paper "model-based RL for biological sequence design"
- RegEvolution:** Local search based on regularized evolution (Real et al., 2019), which has performed well on other black-box optimization tasks and can be seen as an instance of directed evolution.
- DbAs:** Cross-entropy optimization using variational autoencoders (Brookes & Listgarten, 2018).
- FBGAN:** Cross entropy optimization using generative adversarial networks (Gupta & Zou, 2018).
- Bayesopt GP:** Bayesian optimization using a Gaussian process regressor and activation maximization as acquisition function solver.
- Bayesopt ENN** Bayesian optimization using an ensemble of neural network regressors and activation maximization as acquisition function solver.

*activation maximization*<sup>2</sup> simply means back propagation. Specifically, the discrete biological one-hot sequence vectors are preprocessed with continuous relaxation by adding a latent variable  $z$ . This

## 1. Comparison between MCTS and RL model

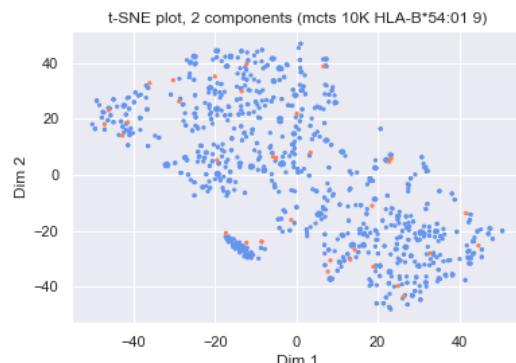
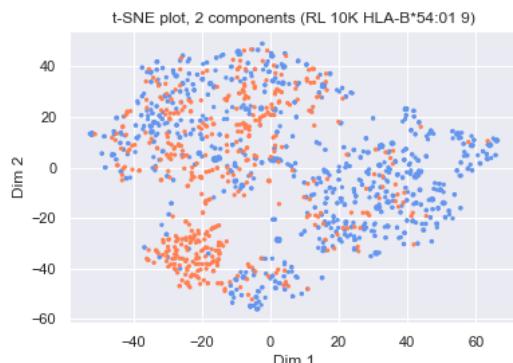
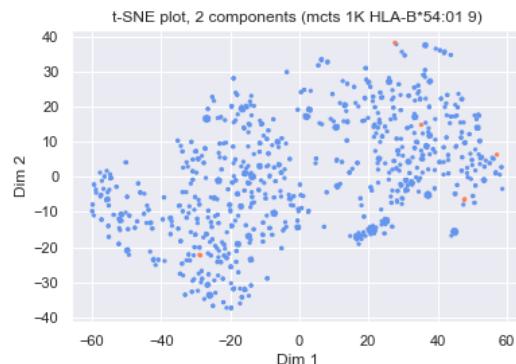
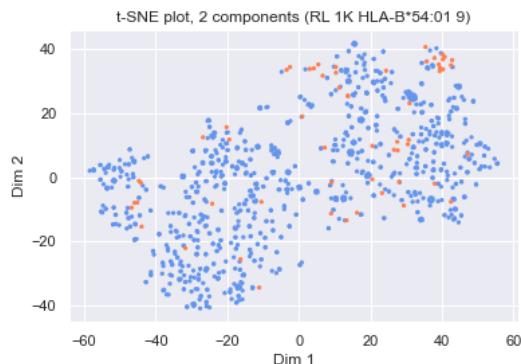
- calculate the distance between peptides of the same length.
- compare the overlapping peptides generated by the RL and generated by the MCTS.
- T-sne plot
- Use the same representation with the MHCflurry and MLP and fully connected layer. *Do not use the gradients to update the padding part.*

Advantages of the RL approach:

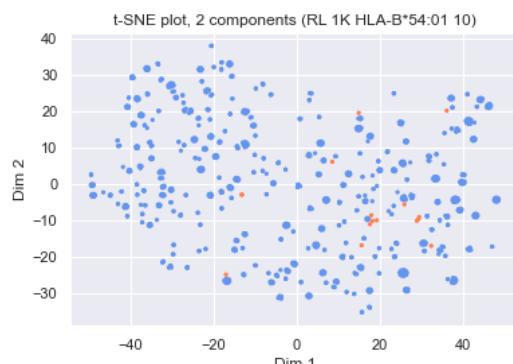
- Good performance
- Wider space
- Interpretability

### 1.1 t-sne plot

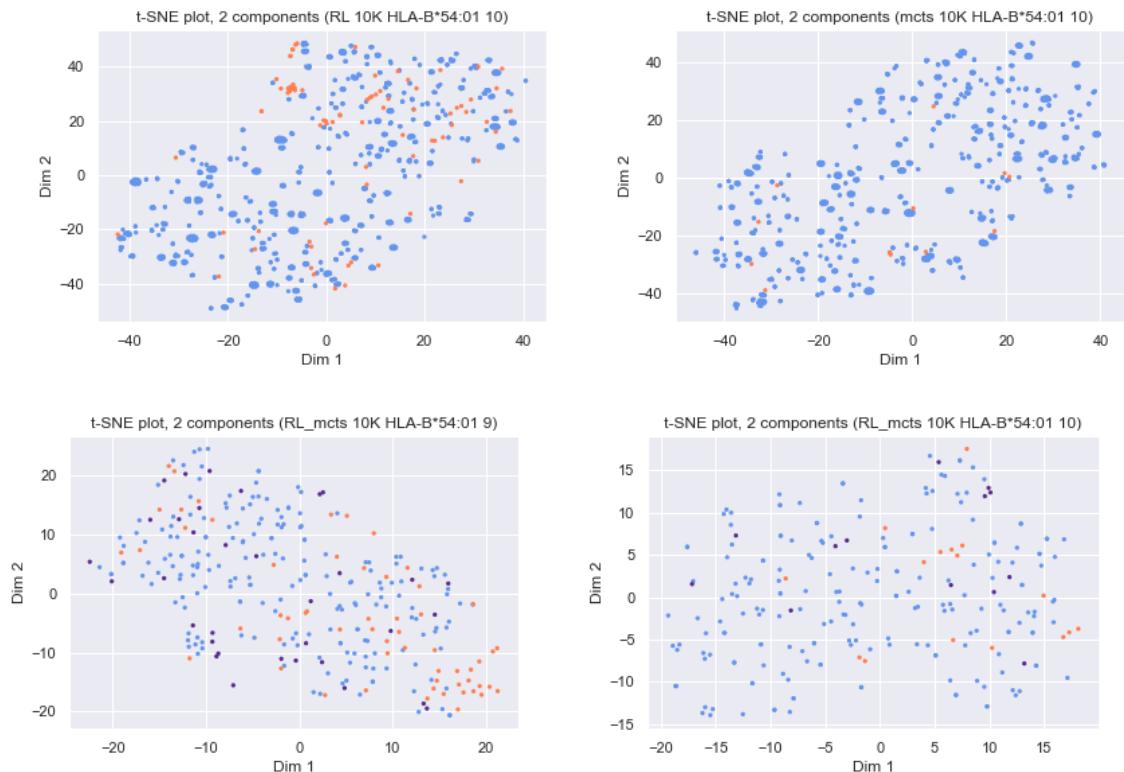
- Allele: HLA-B\*54:01 ; Length = 9



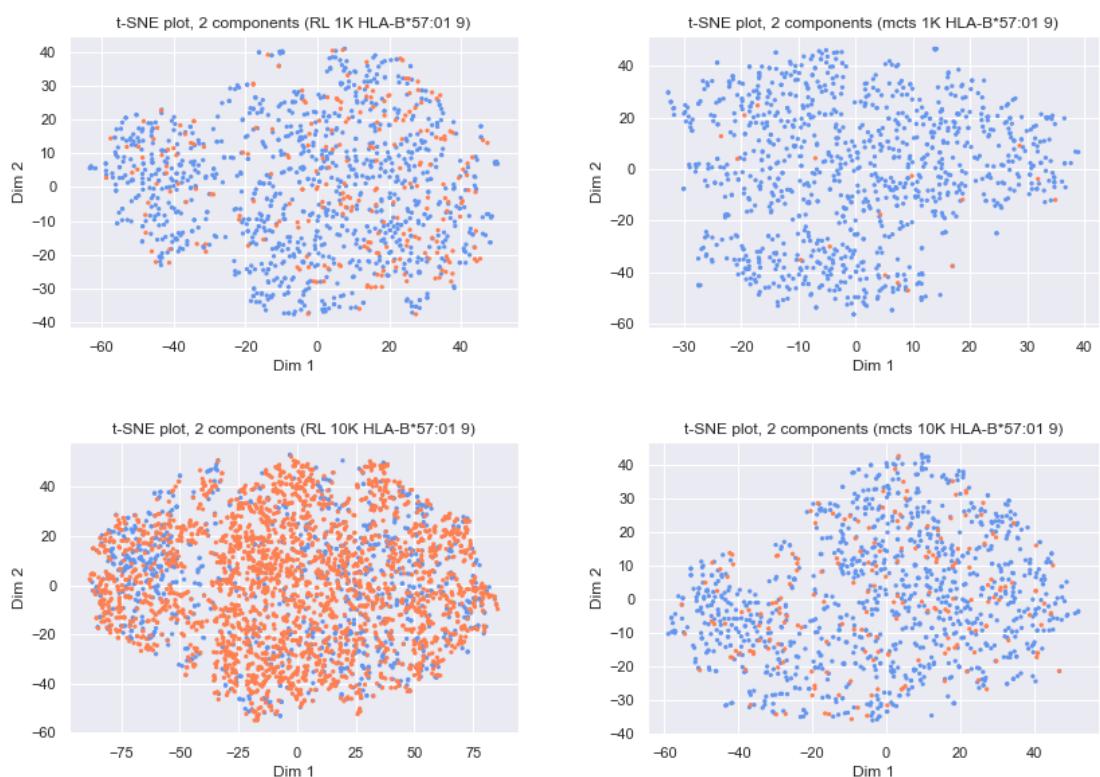
- Allele: HLA-B\*54:01 ; Length = 10



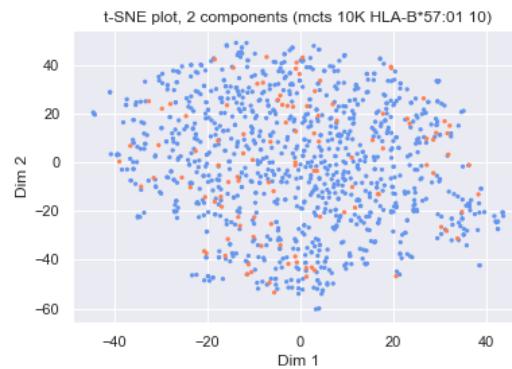
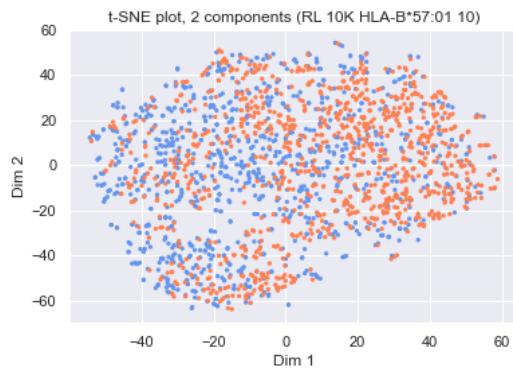
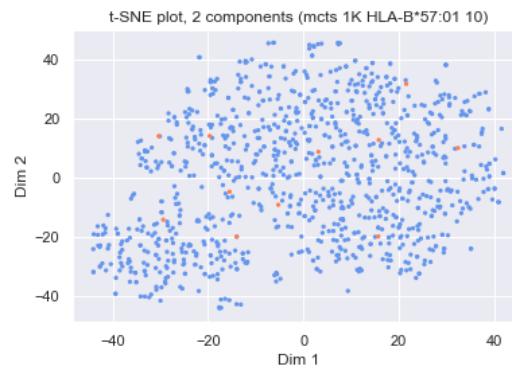
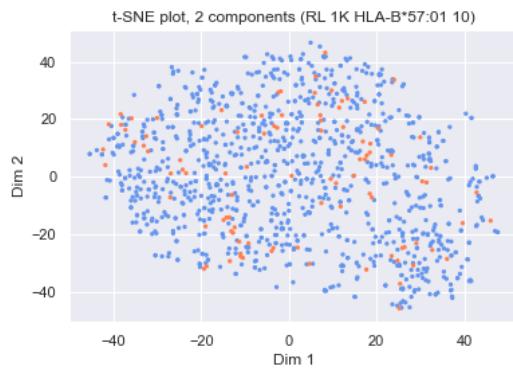
None



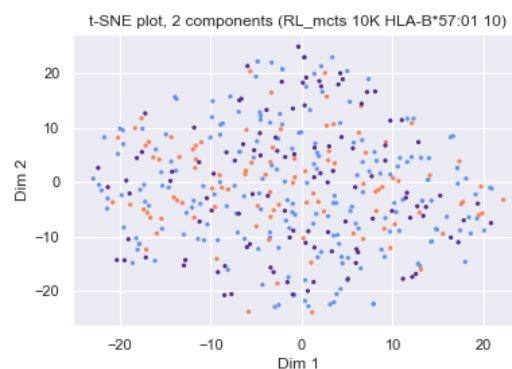
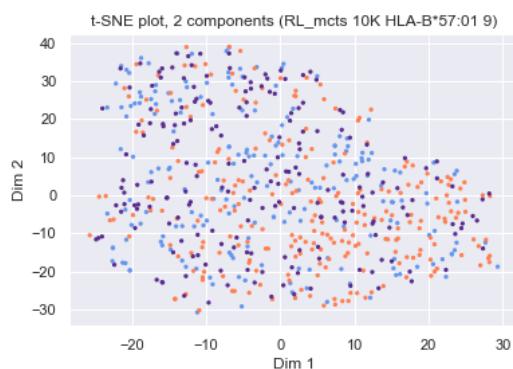
- Allele: HLA-B\*57:01 ; Length = 9



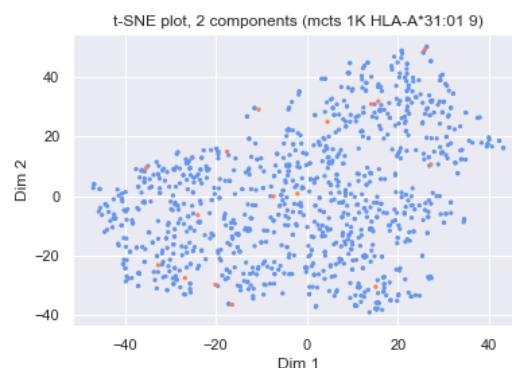
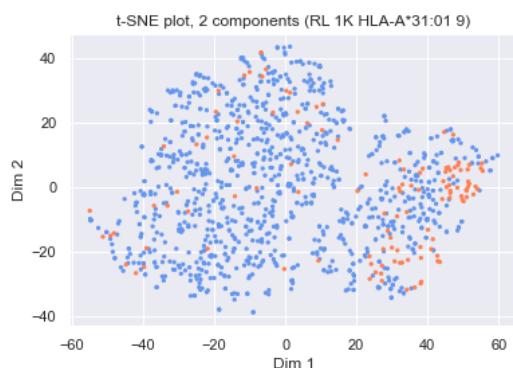
- Allele: HLA-B\*57:01 ; Length = 10

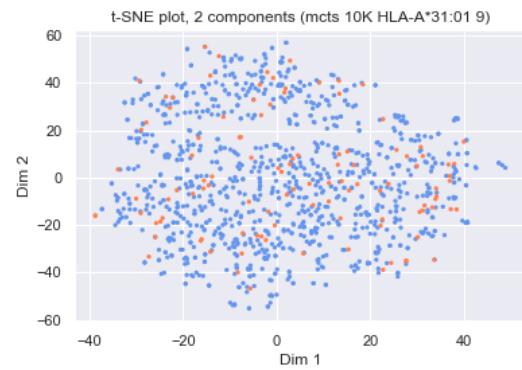
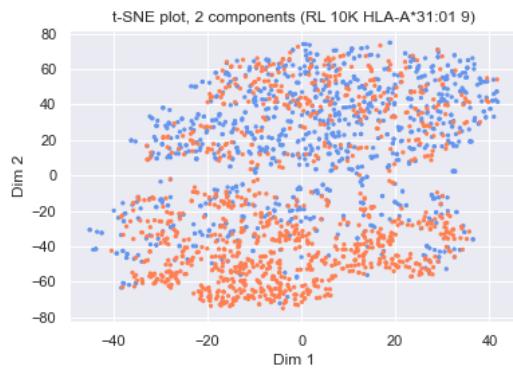


### RL 1K vs MCTS 10K

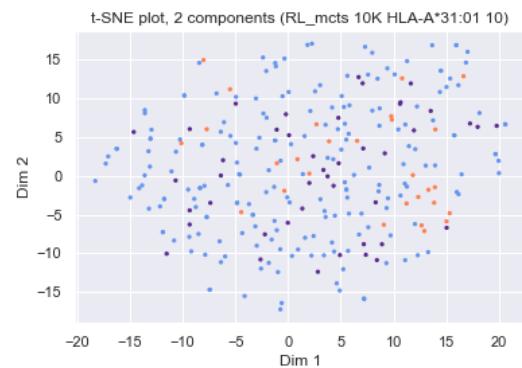
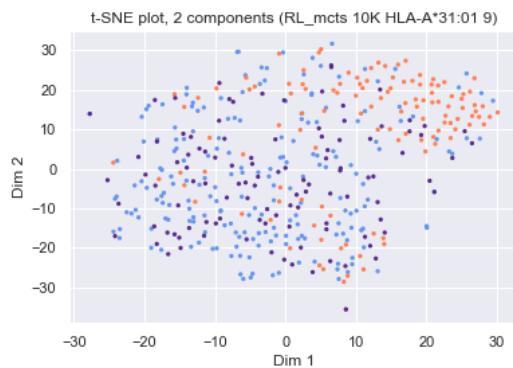
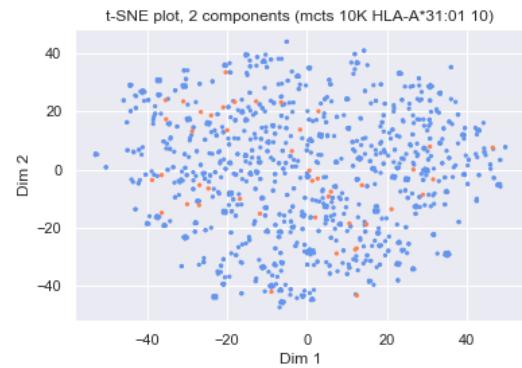
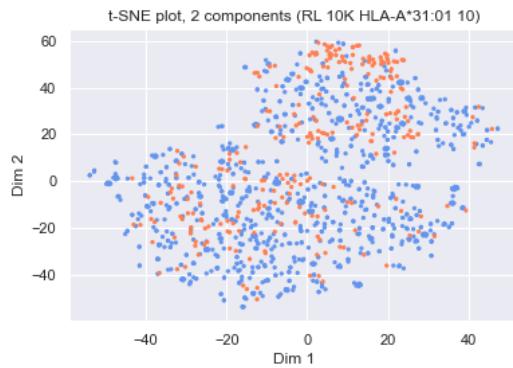
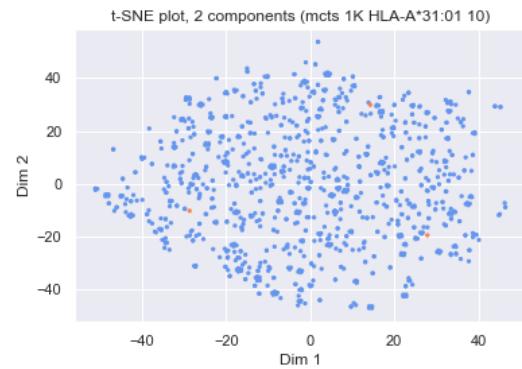
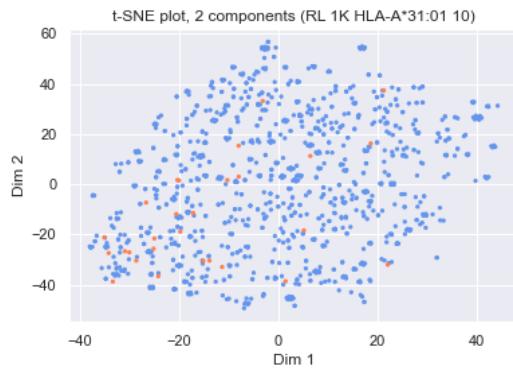


- Allele: HLA-A\*31:01 ; Length = 9



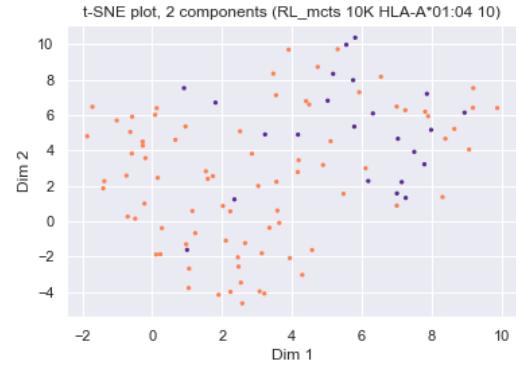
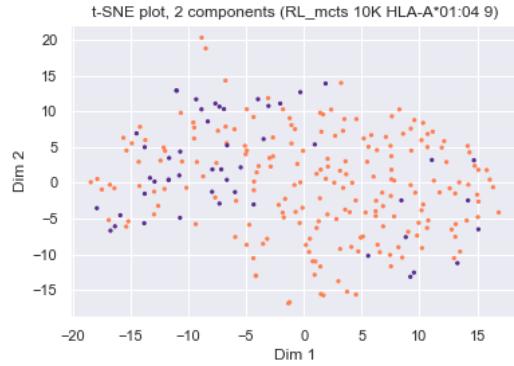


- Allele: HLA-A\*31:01 ; Length = 10



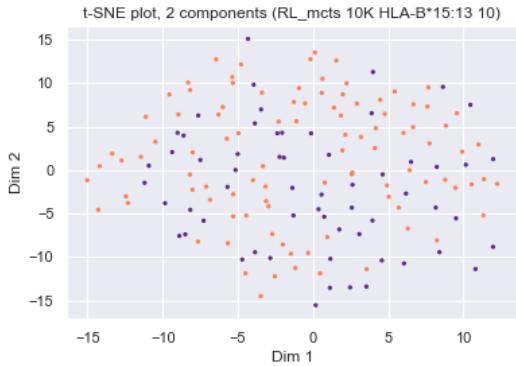
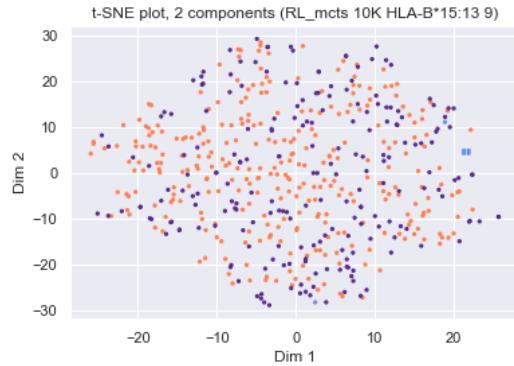
- Allele: HLA-A\*01:04 ; Length = 9 & 10

## RL 1K (Orange) vs MCTS 10K (Purple)



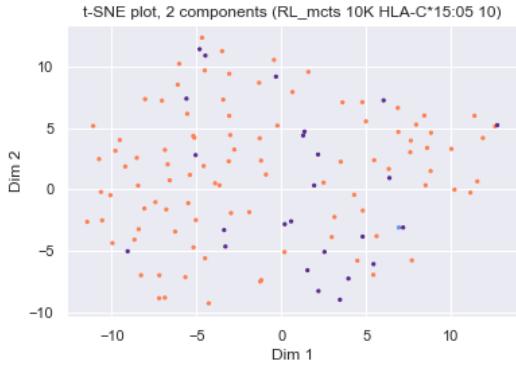
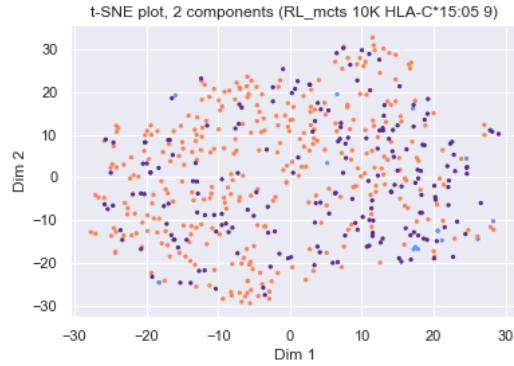
- Allele: HLA-B\*15:13 ; Length = 9 & 10

## RL 1K vs MCTS 10K



- Allele: HLA-C\*15:05 ; Length = 9 & 10

## RL 1K vs MCTS 10K



## 2. Preliminary results of bayesian optimization

The current reconstruct accuracy is about 95%.

The current results of Bayesian optimization with 100 initial points and 1000 iterations are not good.

Possible ways to optimize the results:

1. improve the reconstruct accuracy to 99%
2. try to tune the parameter  $\kappa$  in UCB formula ( $UCB(x) = \mu(x) + \kappa * \sigma(x)$ )
3. improve the number of iterations or initial points so that the GP within BO can converge.

Note that this BO is an allele-specific model. Remember that GP in BO need to model the predicted presentation score  $f(s|a, p)$  between allele and peptide. If we can only maximize peptide sequences, then the input of GP should be only peptide sequence. Therefore, GP actually models the correlation between peptide sequence  $p$  and its presentation score with a given allele  $a$ , that is  $p(s|a, p) = f^a(s|p)$ . Because  $a$  is not an input of GP model

- Generate 500 peptides for each HLA-A allele in the table using RL model;
- Extract 500 peptides (including positive peptides and negative peptides) with binding affinities for the largest allele in IEDB dataset.
- Build a google site

## Update (7-7)

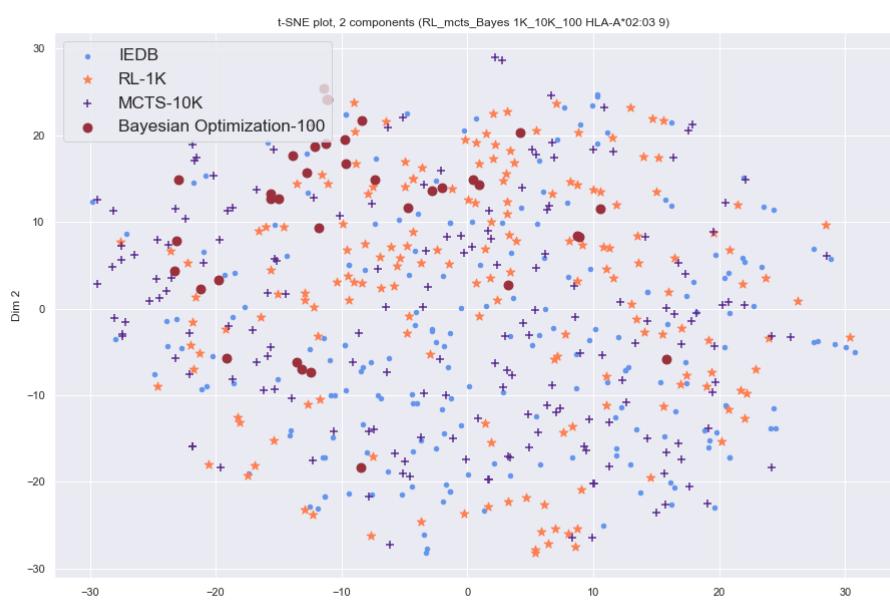
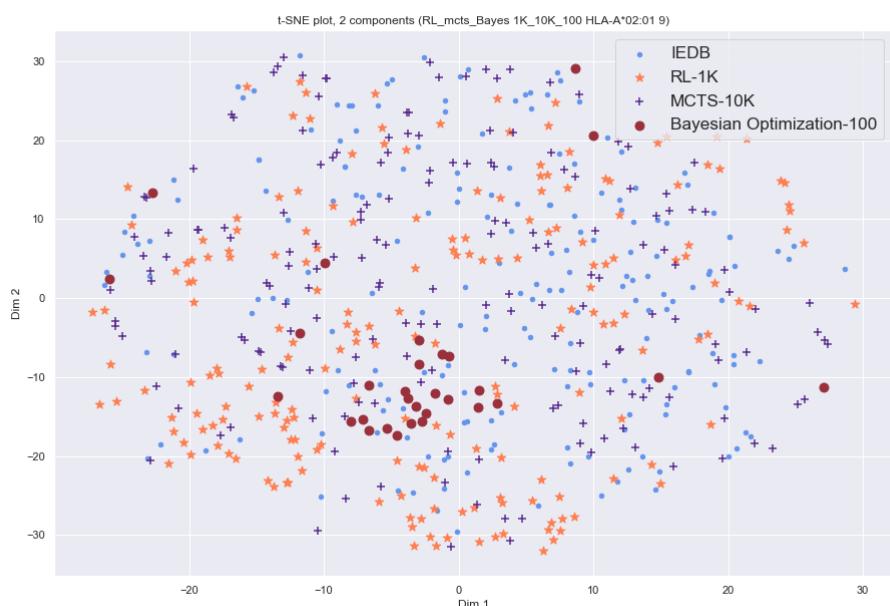
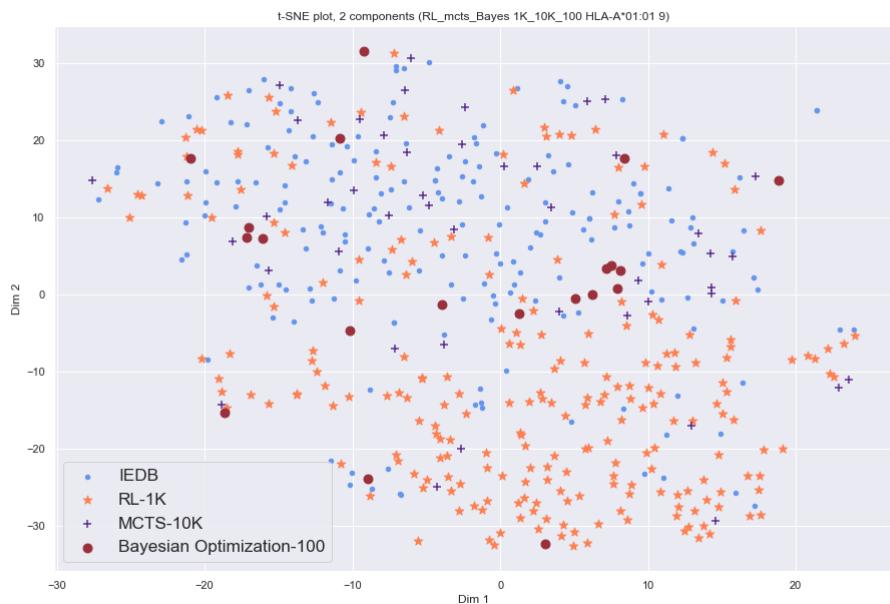
### 1. Result of Bayesian Optimization

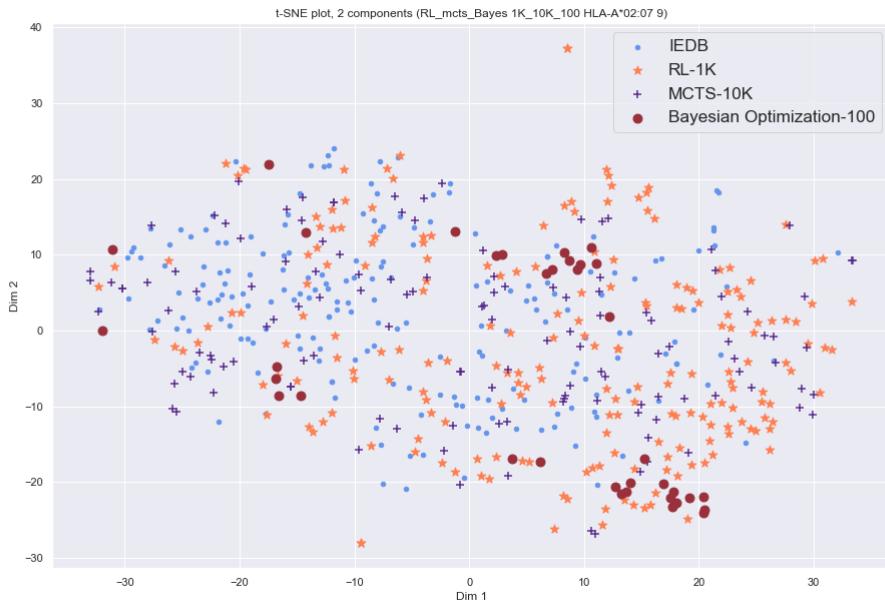
- The previous way I used to run the Bayesian Optimization is incorrect... I ran BO for one time and each BO would optimize the peptides for 1000 steps. However, by checking other papers using the BO, I realized that I should run BO for multiple times and each BO should optimize the peptides for only 5 - 10 steps. In each run, the gaussian process will have different random initial points, and thus the BO could have different optimized peptides.
- This method is very time-consuming. The algorithm has been ran in a paralleled way on 15 cores and still can cost more than 5 hours with 500 initial points (used to train a Gaussian Process model) and only 100 rollouts (i.e., times).
- The number of initial points can significantly influence the running time of this algorithm. The model with 500 initial points achieves significantly better results than that with 100 initial points. However, this model with 500 initial points cost 8 hours to finish the optimization on 23 alleles with 100 rollouts while the model with 100 initial points cost only 2.4 hours.

I just found another BO package and will test whether that one is faster than the one I used.

### Visualization:







## Update (7-12)

### 1. Result of backpropagation

- o initialize the peptide latent embedding by encoding a random peptide.
- o update the peptide latent embedding with the gradients for 50 iterations.
- o output all the latent embeddings with the predicted scores greater than 0.75.
- o decode these latent embeddings to real peptides.
- o check the presentation scores of these peptides.

Observation:

(1) For different alleles, the backpropagation method can find very different number of unique peptides.

This could be due to that some random peptides are optimized to the same good peptide.

That means Backpropagation cannot find diverse peptides for some alleles.

(2) This method is very efficient. It took about 1.5 hours to get the results with 10K rollouts.

### 2. Update on Bayesian Optimization

The package BoTorch is faster than the previous package I used. Now, it took about 70-80 mins to get the results with 100 rollouts. With the previous package, BO with the same parameters can run more than 8 hours.

With BoTorch, we can save the trained GP model and restore it when training on the new sampled data. The results seem to be better.

It can cost more than 13 hours to get the results with 1000 rollouts.

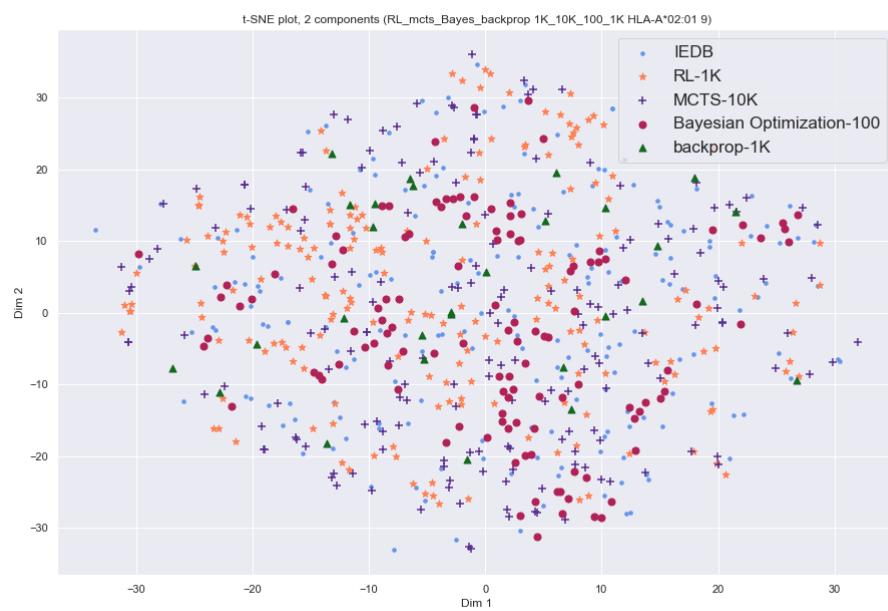
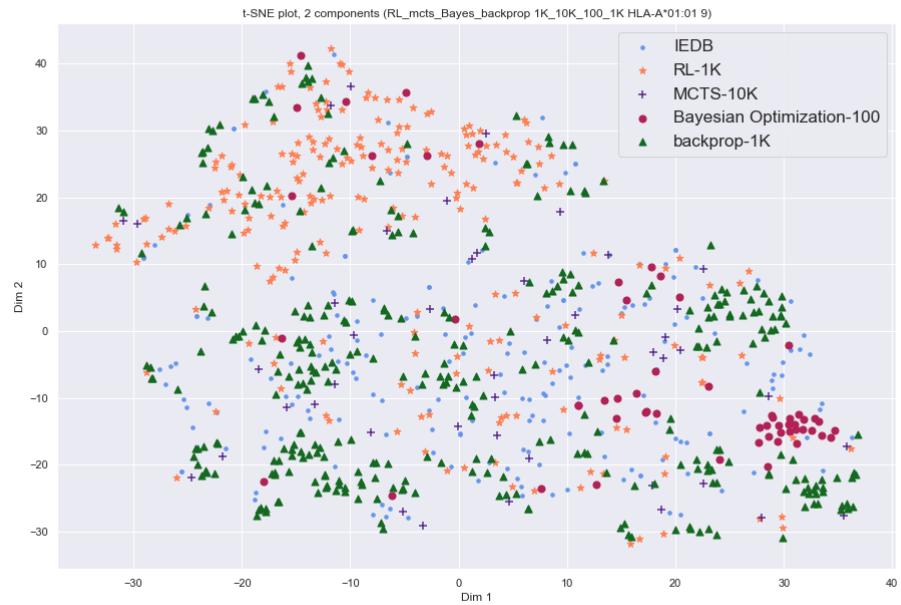
Observation:

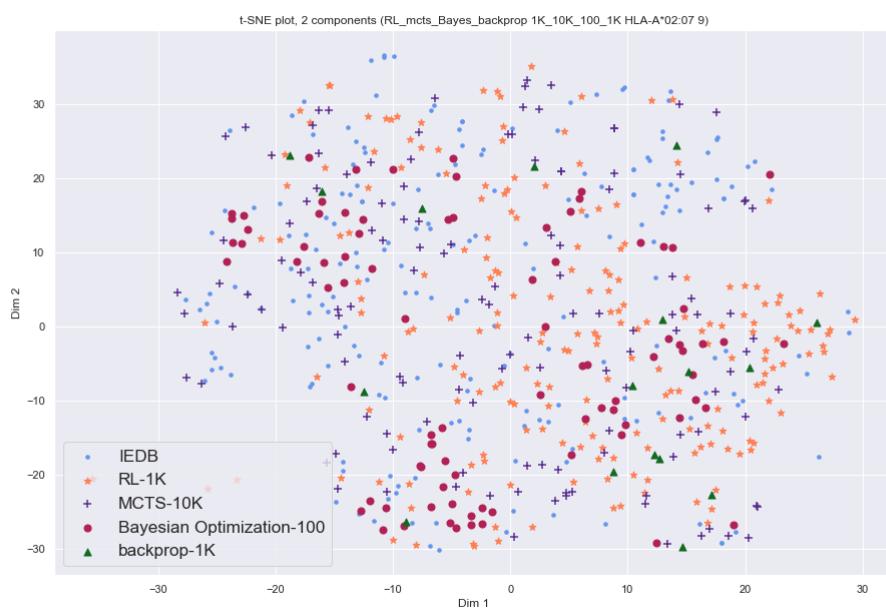
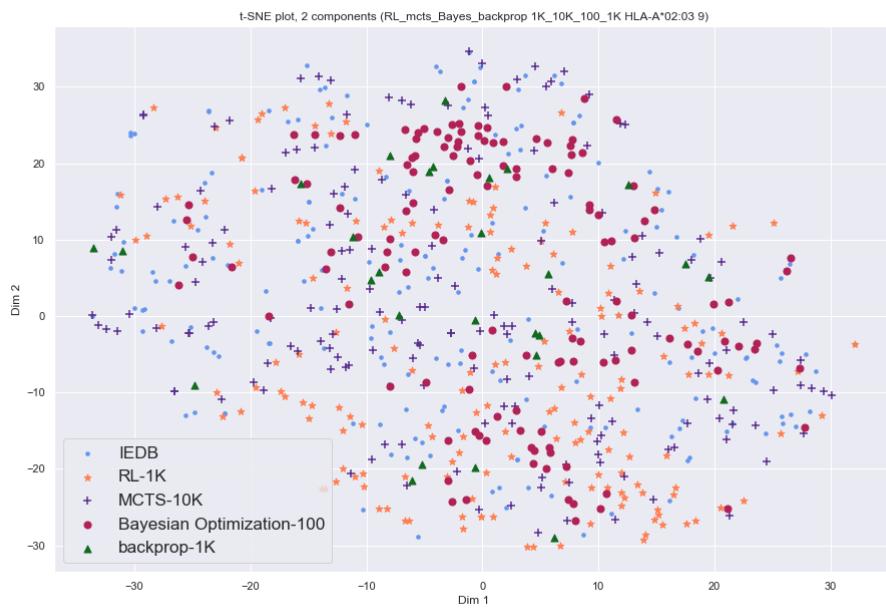
(1) Same with backpropagation, this method can find very different number of unique peptides.

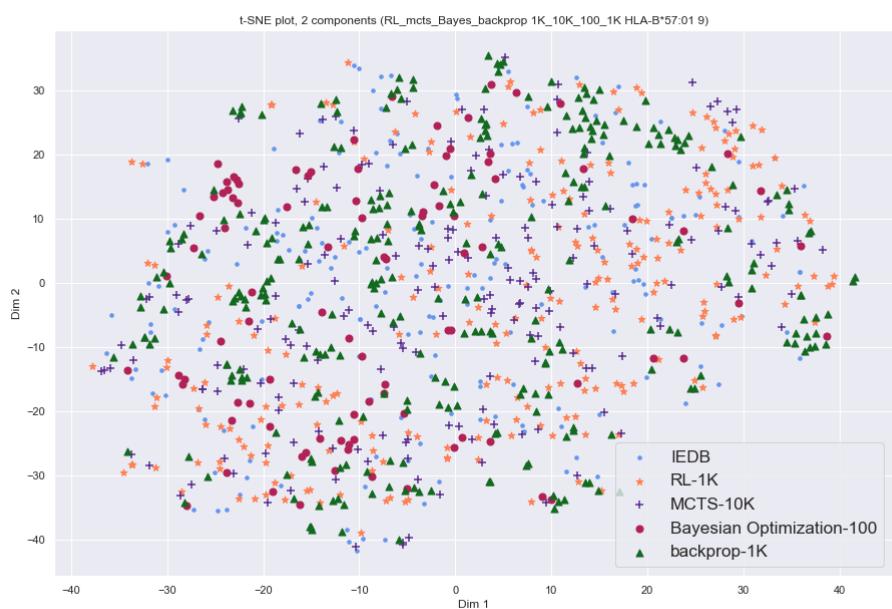
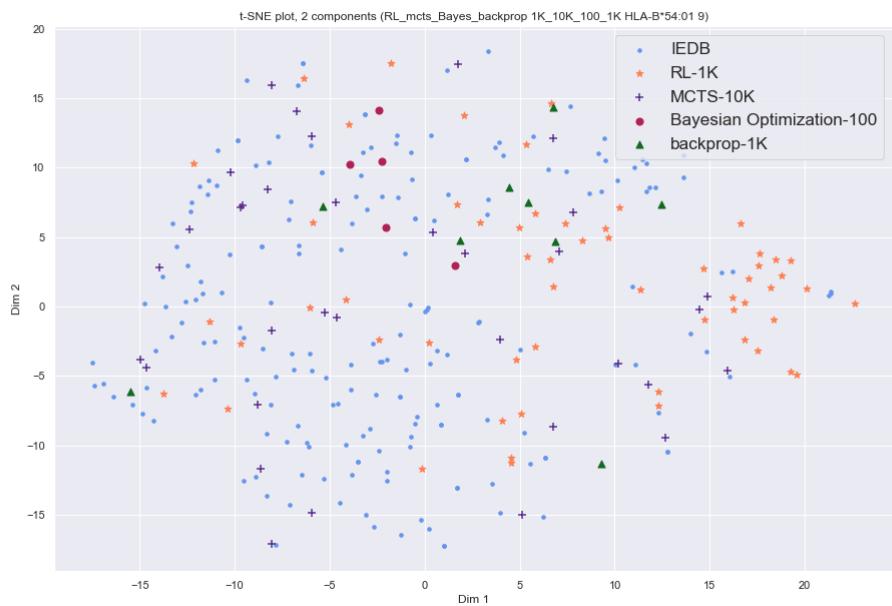
#### Verification:

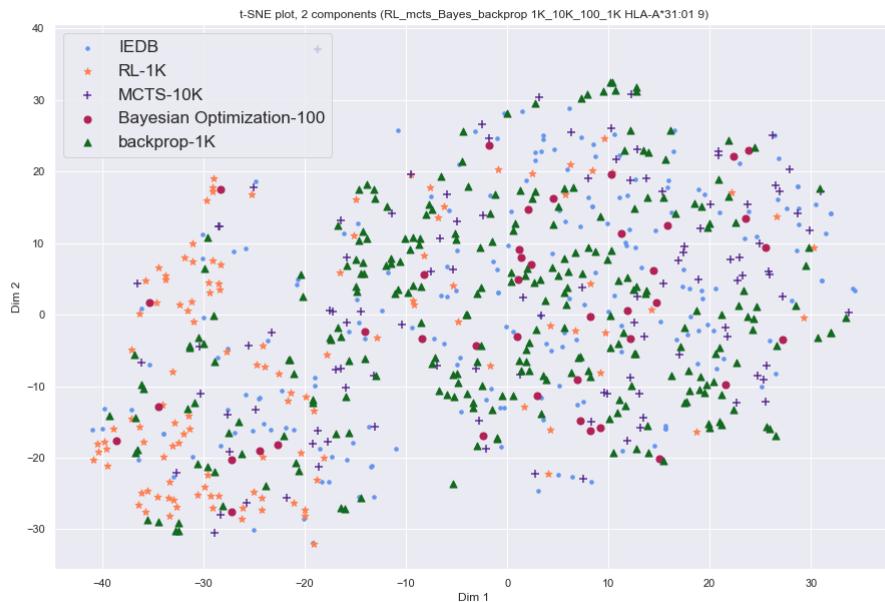
- 1) Check whether the peptides generated by Backpropagation and BO are similar or not.
- 2) Check whether the VAE performs bad on reconstruction of positive peptides of HLA-B5401 and HLA-B5101.

#### Visualization:









### Metrics:

Quanti metrics

- Diversity/Standard Deviation (of top 100?)
- Memory (leave running time to the)
- count of qualified peptides
- average presentation scores

Current metric - pairwise hamming distance doesn't make sense.

## Update (7-23)

### 1. Random Methods

With the batch size 128, this random method can generate 300 - 700 peptides per second.

Among the generated peptides, 0.07% ~ 0.9% peptides are qualified. With 0.07%, this method can generate 1 qualified peptides with 2.1 s ~ 4.8 s. That means, even with the smallest positive rate, this random method is faster than all the other methods.

### 2. Motif Methods

For alleles with the data, this motif-based method can generate many positive peptides.

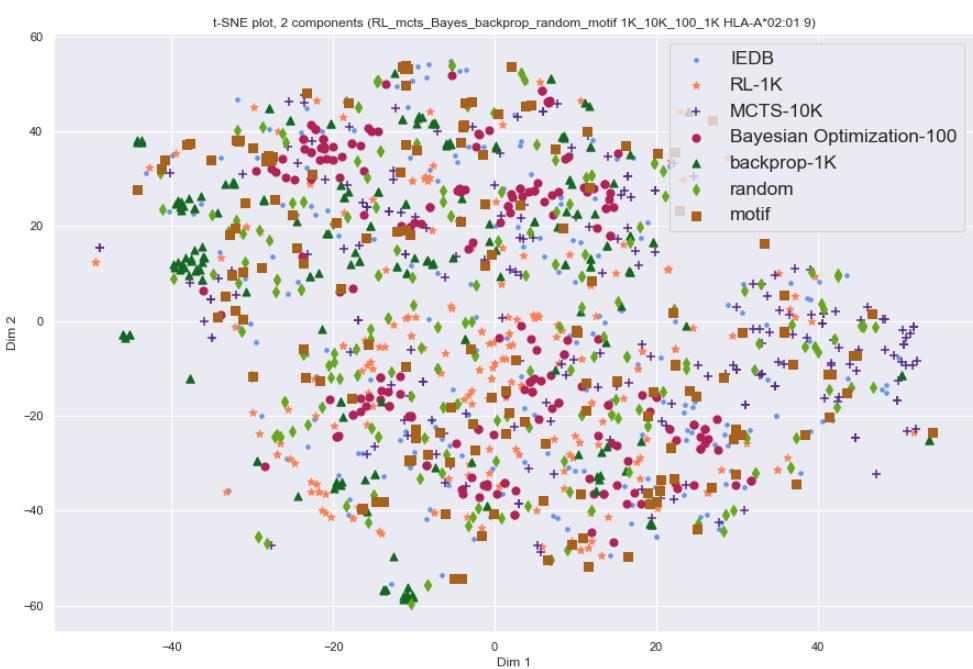
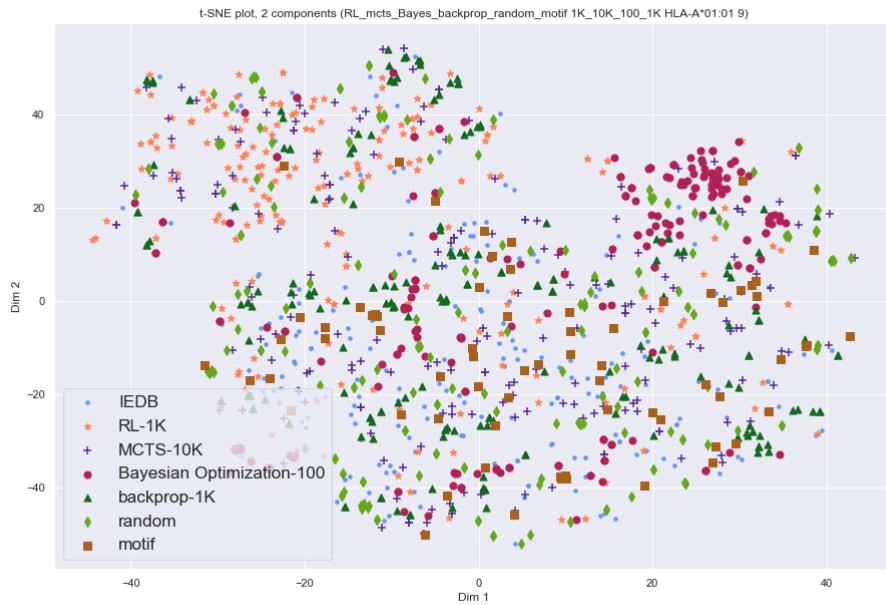
For alleles without the data, sometimes it can generate qualified peptides as many as the random methods, but sometimes it can not. It depends on the quality of the motif from other alleles.

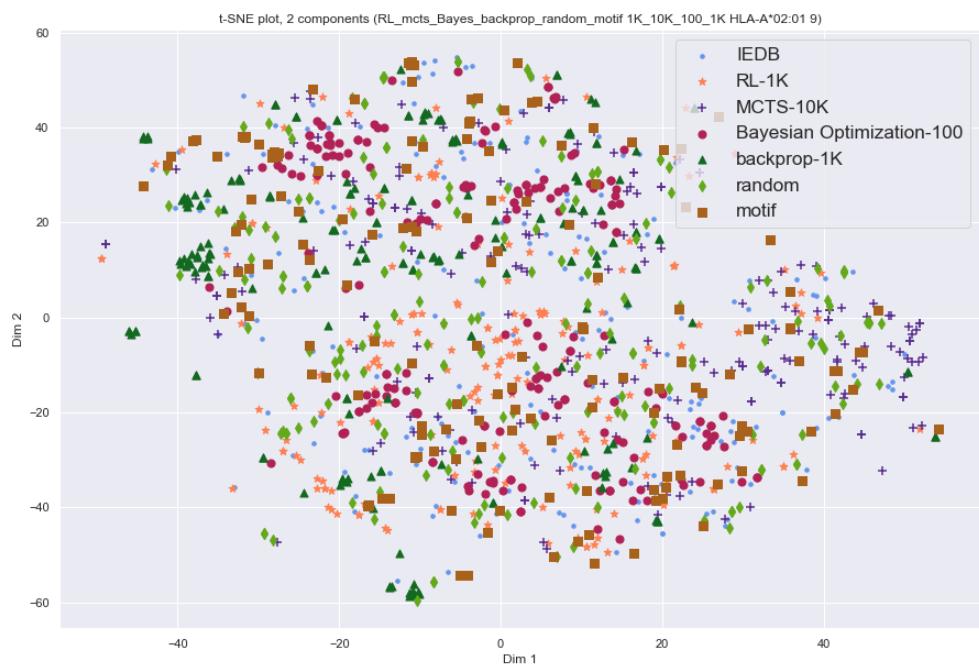
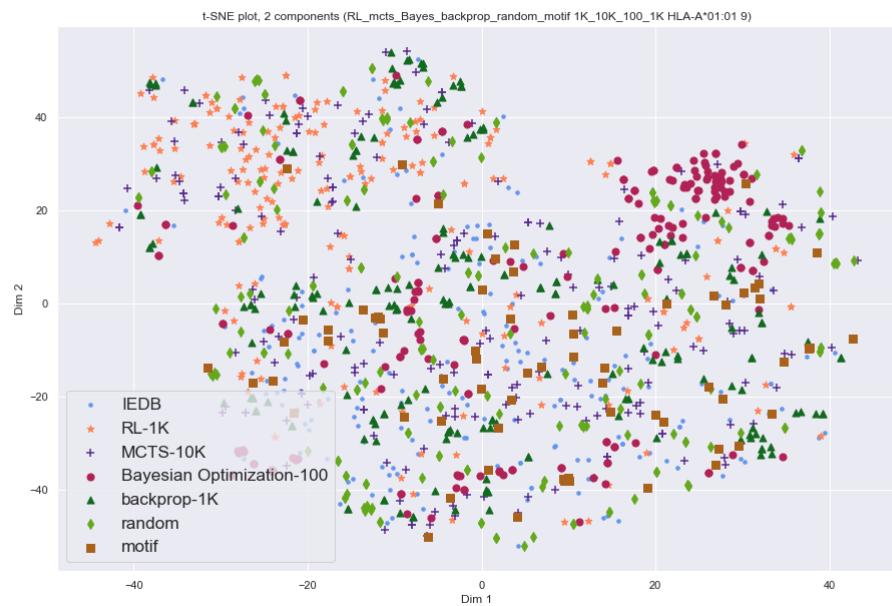
In a conclusion, compared with other baselines and the RL method, random generation and motif-based generation are the most efficient way to do free generation. The high efficiency of these two methods come from the oracle (MHCflurry). Among other baseline methods, at each step, the oracle will be called to evaluate only one peptide, while these two methods asked the oracle to predict 128 peptides.

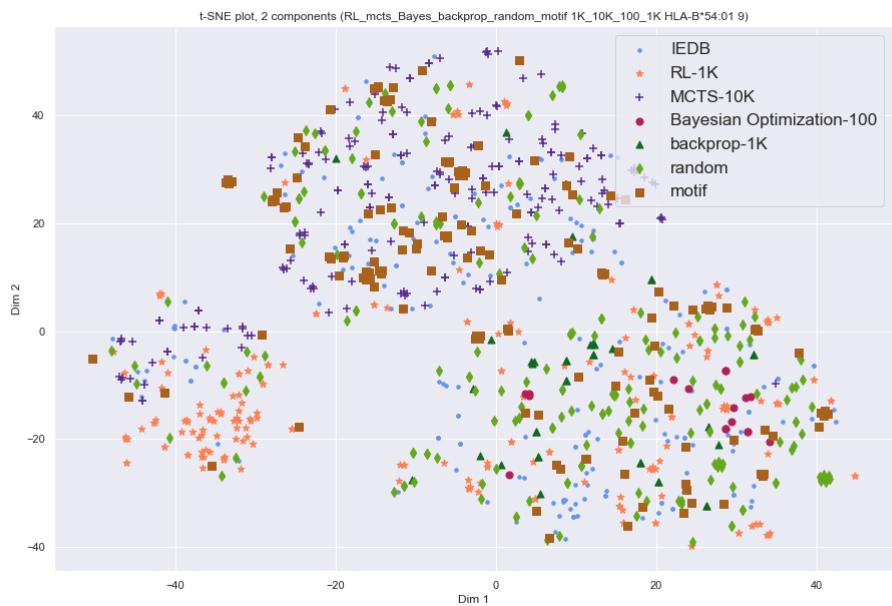
Conclusion: If the objective is to generate sequences without any constraints, and the oracle itself is efficient enough, and the positive rate of random data is not like 0.00[....]01%, using brutal force is a good choice.

Visualization:

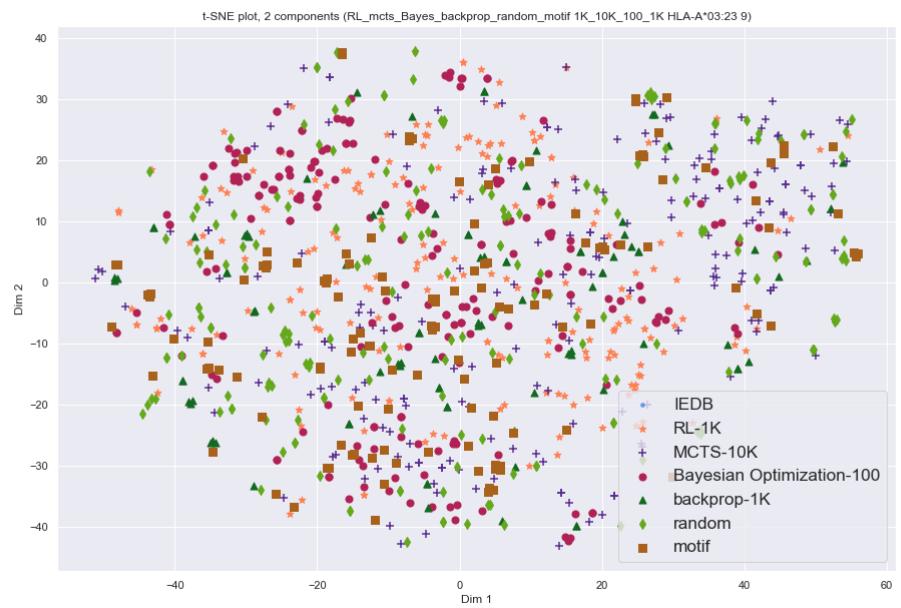
- Alleles with the Data

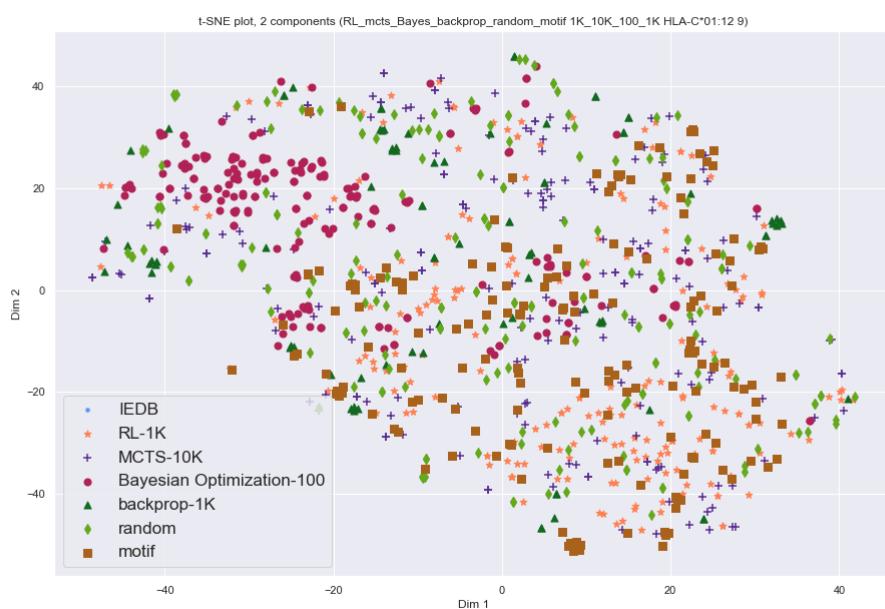
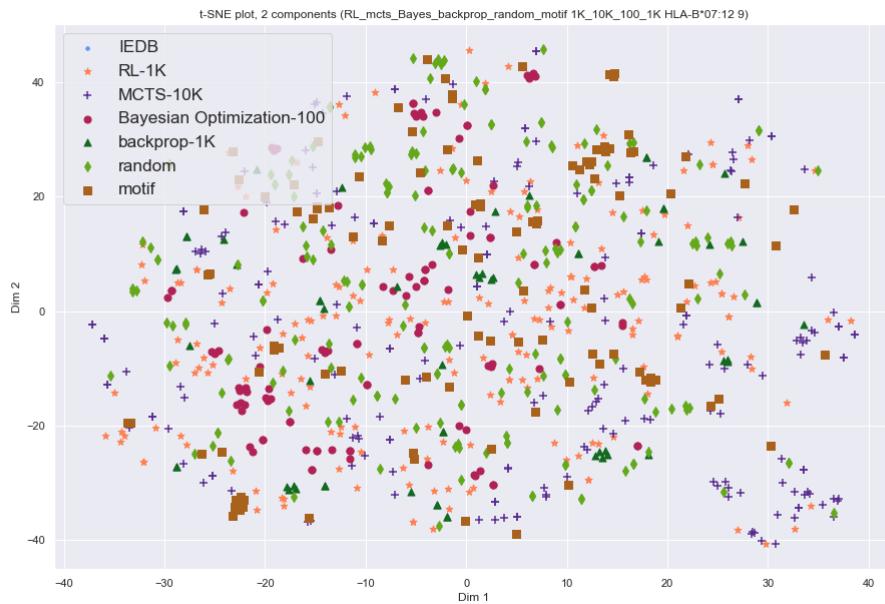






- Alleles without the Data





How can we demonstrate the performance?

1. We can let all the models to generate the same number of peptides, and count the percentage of positive peptides.

The data of MHC class 2 is more complicated since the length of peptides ranges from 12 to 19.

## Update (7-26)

Objective: use **physical model** as reward function to generate peptides for alleles **without data**.

- Why "random" cannot be the solution for the project with physical model ?

Since we need to use physical model, sample efficiency could be a problem. In this case, we should not use random to generate millions of peptides and evaluate them to get hundreds of positive peptides, because this approach assumes that our oracle (i.e., MHCflurry or Physical model) is very efficient so that millions of peptides can be evaluated.

- Why do we need the RL? Can we directly use "random" or any other baselines on the updated classifier with the augmented data from physical model?

If we directly use the physical model to generate a batch of pMHC data and then update the classifier, we do not need RL to optimize on that classifier.

However, by using the physical model in an interactive way, RL could have advantages below:

- (1) achieve better exploration in an efficient way;
- (2) infer motif from the learned policy;
- (3) learn the peptide generation problem with physical model in an active learning way.
- (4) directly learn the signals from physical model.

- How can we solve this "sample efficiency" problem?

Though we can still use the data from physical model to update MHCflurry and get more accurate predictions, generating data for 10k alleles using physical model is unpractical.

Therefore, if we want to use physical model, no matter with RL or other baselines (random / backpropagation), we have to make sure that the testing pairs for physical model are very important and can help reduce **uncertainty**.

- How can we quantify the uncertainty?

The uncertainty of this peptide generation problem comes from two parts: (1) alleles without data; (2) unvisited peptides.

The most simplest way is to use the uncertainty from the ensemble of networks. But the uncertainty estimation could be improved, for example, we can formulate this problem in an hierarchical way.

- (1) solution for alleles without data;

Inspired by the uncertainty estimation method for active learning <sup>3</sup>, I think maybe we can maintain a graph with all the alleles.

Each allele represents a node in the graph and has the number of peptides as an attribute. Similar alleles will be connected with the edges. The weight of the edges could be determined by the number of common positive peptides or the similarity on latent embeddings.

(For RL, the alleles with high uncertainty could be selected with higher probability during the training.)

- (2) solution for unvisited peptides; I don't have a specific idea.....

- How can we implement the uncertainty with RL? any related work?

I just found a paper that could be related to our objective. The paper name is "MUSBO: Model-based Uncertainty Regularized and Sample Efficient Batch Optimization for Deployment Constrained Reinforcement Learning". They used an ensemble of networks to calculate the uncertainty of each state-action pair. The state-action pair with higher uncertainty will contribute less to the update of agent.

We can use similar framework to calculate the uncertainty of the predicted reward. The state-action pair with higher uncertainty will be evaluated using physical model.

- If we do not have the physical model available, how can we prepare for the AAAI submission?

(1) change the current model to model-based. use an ensemble of reward function models to calculate the uncertainty.

(2) select those pairs with the highest uncertainty for the evaluation of MHCflurry. To make sure the simulation is accurate, we can only use the alleles with the data for training and testing.

(3) limit the calls of MHCflurry for all the baselines (e.g., 10k calls)? Some baselines are allele-specific..... Maybe we can use only pan-allele methods such as random methods and backpropagation methods.

We can even conduct an ablation study to evaluate whether the performance of backpropagation method can increase with the selected uncertain pMHC data from RL.

- 
1. Pardo, Fabio, et al. "Time limits in reinforcement learning." *International Conference on Machine Learning*. PMLR, 2018.  [↗](#)
  2. Killoran, Nathan, et al. "Generating and designing DNA with deep generative models." *arXiv preprint arXiv:1712.06148* (2017).  [↗](#)
  3. Multi-Class Active Learning by Uncertainty Sampling with Diversity Maximization  [↗](#)