

---

# MHCflurry Documentation

*Release 1.6.0*

Timothy O'Donnell

Mar 23, 2020



# CONTENTS

<b>1</b>	<b>Introduction and setup</b>	<b>1</b>
1.1	Installation (pip)	1
1.2	Using conda	2
<b>2</b>	<b>Command-line tutorial</b>	<b>3</b>
2.1	Downloading models	3
2.2	Generating predictions	3
2.3	Scanning protein sequences for predicted MHC I ligands	4
2.4	Fitting your own models	5
2.5	Environment variables	7
<b>3</b>	<b>Python library tutorial</b>	<b>9</b>
3.1	Loading a predictor	9
3.2	Predicting for individual peptides	9
3.3	Scanning protein sequences	11
3.4	Lower level interfaces	12
<b>4</b>	<b>Command-line reference</b>	<b>13</b>
4.1	mhcflurry-predict	13
4.2	mhcflurry-predict-scan	15
4.3	mhcflurry-downloads	17
4.3.1	mhcflurry-downloads fetch	17
4.3.2	mhcflurry-downloads info	18
4.3.3	mhcflurry-downloads path	18
4.3.4	mhcflurry-downloads url	18
4.4	mhcflurry-class1-train-allele-specific-models	18
4.5	mhcflurry-class1-select-allele-specific-models	19
4.6	mhcflurry-class1-train-pan-allele-models	21
4.7	mhcflurry-class1-select-pan-allele-models	22
4.8	mhcflurry-class1-train-processing-models	23
4.9	mhcflurry-class1-select-processing-models	25
4.10	mhcflurry-class1-train-presentation-models	26
<b>5</b>	<b>API Documentation</b>	<b>27</b>
5.1	Submodules	48
5.2	mhcflurry.allele_encoding module	48
5.3	mhcflurry.amino_acid module	49
5.4	mhcflurry.calibrate_percentile_ranks_command module	50
5.5	mhcflurry.class1_affinity_predictor module	51
5.6	mhcflurry.class1_neural_network module	57

5.7	mhcflurry.class1_presentation_predictor module . . . . .	63
5.8	mhcflurry.class1_processing_neural_network module . . . . .	68
5.9	mhcflurry.class1_processing_predictor module . . . . .	70
5.10	mhcflurry.cluster_parallelism module . . . . .	73
5.11	mhcflurry.common module . . . . .	74
5.12	mhcflurry.custom_loss module . . . . .	77
5.13	mhcflurry.data_dependent_weights_initialization module . . . . .	79
5.14	mhcflurry.downloads module . . . . .	79
5.15	mhcflurry.downloads_command module . . . . .	81
5.16	mhcflurry.encodable_sequences module . . . . .	83
5.17	mhcflurry.ensemble_centrality module . . . . .	86
5.18	mhcflurry.fasta module . . . . .	86
5.19	mhcflurry.flanking_encoding module . . . . .	86
5.20	mhcflurry.hyperparameters module . . . . .	88
5.21	mhcflurry.local_parallelism module . . . . .	88
5.22	mhcflurry.percent_rank_transform module . . . . .	90
5.23	mhcflurry.predict_command module . . . . .	91
5.24	mhcflurry.predict_scan_command module . . . . .	92
5.25	mhcflurry.random_negative_peptides module . . . . .	92
5.26	mhcflurry.regression_target module . . . . .	94
5.27	mhcflurry.scoring module . . . . .	94
5.28	mhcflurry.select_allele_specific_models_command module . . . . .	95
5.29	mhcflurry.select_pan_allele_models_command module . . . . .	96
5.30	mhcflurry.select_processing_models_command module . . . . .	97
5.31	mhcflurry.testing_utils module . . . . .	98
5.32	mhcflurry.train_allele_specific_models_command module . . . . .	98
5.33	mhcflurry.train_pan_allele_models_command module . . . . .	98
5.34	mhcflurry.train_presentation_models_command module . . . . .	100
5.35	mhcflurry.train_processing_models_command module . . . . .	100
5.36	mhcflurry.version module . . . . .	101
<b>Python Module Index</b>		<b>103</b>
<b>Index</b>		<b>105</b>

## INTRODUCTION AND SETUP

MHCflurry is an open source package for peptide/MHC I binding affinity prediction. It aims to provide competitive accuracy with a fast and documented implementation.

You can download pre-trained MHCflurry models fit to mass spec-identified MHC I ligands and peptide/MHC affinity measurements deposited in IEDB (plus a few other sources) or train a MHCflurry predictor on your own data.

Starting in version 1.6.0, the default MHCflurry binding affinity predictors are “pan-allele” models that support most sequenced MHC I alleles across humans and a few other species (about 14,000 alleles in total). This version also introduces two experimental predictors, an “antigen processing” predictor that attempts to model MHC allele-independent effects such as proteosomal cleavage and a “presentation” predictor that integrates processing predictions with binding affinity predictions to give a composite “presentation score.” Both models are trained on mass spec-identified MHC ligands.

MHCflurry supports Python 3.4+. It uses the [keras](#) neural network library via either the Tensorflow or Theano backends. GPUs may optionally be used for a modest speed improvement.

If you find MHCflurry useful in your research please cite:

T. J. O'Donnell, et al., “MHCflurry: Open-Source Class I MHC Binding Affinity Prediction,” *Cell Systems*, 2018. [https://www.cell.com/cell-systems/fulltext/S2405-4712\(18\)30232-1](https://www.cell.com/cell-systems/fulltext/S2405-4712(18)30232-1).

If you have questions or encounter problems, please file an issue at the MHCflurry github repo: <https://github.com/openvax/mhcflurry>

### 1.1 Installation (pip)

Install the package:

```
$ pip install mhcflurry
```

Then download our datasets and trained models:

```
$ mhcflurry-downloads fetch
```

From a checkout you can run the unit tests with:

```
$ pip install nose
$ nosetests .
```

## 1.2 Using conda

You can alternatively get up and running with a `conda` environment as follows. Some users have reported that this can avoid problems installing tensorflow.

```
$ conda create -q -n mhcflurry-env python=3.6 'tensorflow<2.0.0'
$ source activate mhcflurry-env
```

Then continue as above:

```
$ pip install mhcflurry
$ mhcflurry-downloads fetch
```

## COMMAND-LINE TUTORIAL

### 2.1 Downloading models

Most users will use pre-trained MHCflurry models that we release. These models are distributed separately from the pip package and may be downloaded with the *mhcflurry-downloads* tool:

```
$ mhcflurry-downloads fetch models_class1_presentation
```

Files downloaded with *mhcflurry-downloads* are stored in a platform-specific directory. To get the path to downloaded data, you can use:

```
$ mhcflurry-downloads path models_class1_presentation  
/Users/tim/Library/Application Support/mhcflurry/4/1.6.0/models_class1_presentation/
```

We also release a number of other “downloads,” such as curated training data and some experimental models. To see what’s available and what you have downloaded, run `mhcflurry-downloads info`.

Most users will only need `models_class1_presentation`, however, as the presentation predictor includes a peptide / MHC I binding affinity (BA) predictor as well as an antigen processing (AP) predictor.

---

**Note:** The code we use for *generating* the downloads is in the `downloads_generation` directory in the repository (<https://github.com/openvax/mhcflurry/tree/master/downloads-generation>)

---

### 2.2 Generating predictions

The *mhcflurry-predict* command generates predictions for individual peptides (see the next section for how to scan protein sequences for epitopes). By default it will use the pre-trained models you downloaded above. Other models can be used by specifying the `--models` argument.

Running:

```
$ mhcflurry-predict  
  --alleles HLA-A0201 HLA-A0301  
  --peptides SIINFELK SIINFELD SIINFELQ  
  --out /tmp/predictions.csv  
Predicting processing.  
Predicting affinities.  
Wrote: /tmp/predictions.csv
```

results in a file like this:

```
$ cat /tmp/predictions.csv
allele,peptide,mhcflurry_affinity,mhcflurry_affinity_percentile,mhcflurry_processing_
↪score,mhcflurry_presentation_score
HLA-A0201,SIINFEDL,12906.786866783777,6.566374999999997,0.10147304146084934,0.
↪012502709574418581
HLA-A0201,SIINFEDK,34491.778276696794,44.612750000000001,0.015448467253008857,0.
↪00387464531700359
HLA-A0201,SIINFEDQ,31476.142209308568,30.217625,0.0359264743165113,0.
↪004518140004135804
HLA-A0301,SIINFEDL,32253.5874797359,20.37875,0.10147304146084934,0.0055786327707357695
HLA-A0301,SIINFEDK,37075.86269121571,58.8585,0.015448467253008857,0.
↪0036345557828772576
HLA-A0301,SIINFEDQ,33258.4276356094,24.029125,0.0359264743165113,0.004303226784859033
```

The binding affinity predictions are given as affinities (KD) in nM in the `mhcflurry_affinity` column. Lower values indicate stronger binders. A commonly-used threshold for peptides with a reasonable chance of being immunogenic is 500 nM.

The `mhcflurry_affinity_percentile` gives the percentile of the affinity prediction among a large number of random peptides tested on that allele (range 0 - 100). Lower is stronger. Two percent is a commonly-used threshold.

The last two columns give the antigen processing and presentation scores, respectively. These range from 0 to 1 with higher values indicating more favorable processing or presentation.

---

**Note:** The processing predictor is experimental. It models allele-independent effects that influence whether a peptide will be detected in a mass spec experiment. The presentation score is a simple logistic regression model that combines the (log) binding affinity prediction with the processing score to give a composite prediction. The resulting prediction may be useful for prioritizing potential epitopes, but no thresholds have been established for what constitutes a “high enough” presentation score.

---

In most cases you’ll want to specify the input as a CSV file instead of passing peptides and alleles as commandline arguments. If you’re relying on the processing or presentation scores, you may also want to pass the upstream and downstream sequences of the peptides from their source proteins for potentially more accurate cleavage prediction. See the [mhcflurry-predict](#) docs.

## 2.3 Scanning protein sequences for predicted MHC I ligands

Starting in version 1.6.0, MHCflurry supports scanning proteins for MHC-binding peptides using the `mhcflurry-predict-scan` command.

We’ll generate predictions across `example.fasta`, a FASTA file with two short sequences:

```
>protein1
MSSSSTPVCNPGNCGV
>protein2
MVENKRLLEGMEMIFGVIPGA
```

Here’s the `mhcflurry-predict-scan` invocation to scan the proteins for binders to either of two MHC I genotypes (using a 100 nM threshold):

```
$ mhcflurry-predict-scan
  example.fasta
  --alleles
```

(continues on next page)



(continued from previous page)

```

HLA-A*02:01,HLA-A*03:01,HLA-B*57:01,HLA-B*45:01,HLA-C*02:02,HLA-C*07:02
HLA-A*01:01,HLA-A*02:06,HLA-B*44:02,HLA-B*07:02,HLA-C*01:02,HLA-C*03:01
--results-filtered affinity
--threshold-affinity 100
Guessed input file format: fasta
Read input fasta with 2 sequences
sequence_id      sequence
0      protein1  MSSSSTPVCNPGPGNCQV
1      protein2  MVENKRLLEGMEMIFGQVIPGA
Predicting processing.
Predicting affinities.
sequence_name,pos,peptide,n_flank,c_flank,sample_name,affinity,best_allele,affinity_
↪percentile,processing_score,presentation_score
protein2,5,RLLEGMEMI,MVENK,FGQVIPGA,genotype_00,15.506647240532306,HLA-A*02:01,0.
↪06800000000000002,0.9045842811465263,0.9899999181261543
protein2,5,RLLEGMEMI,MVENK,FGQVIPGA,genotype_01,17.618050003789662,HLA-A*02:06,0.
↪09587500000000006,0.9045842811465263,0.9888061128232571
protein2,12,MIFGQVIPGA,MVENKRLLEGME,,genotype_01,27.74054476375278,HLA-A*02:06,0.
↪2527499999999999,0.9025756493210793,0.9831738588268119
protein2,4,KRLLEGMEM,MVEN,IFGQVIPGA,genotype_00,29.368389637130125,HLA-C*07:02,0.
↪09150000000000004,0.49434878677129745,0.9171519649768841
protein2,12,MIFGQVIPGA,MVENKRLLEGME,,genotype_00,35.083004993507345,HLA-A*02:01,0.
↪34662499999999985,0.9025756493210793,0.9793299707138564
protein1,8,CPNPGPGNCQV,MSSSSTPV,,genotype_01,44.167261275374614,HLA-B*07:02,0.236125,0.
↪16842720657587051,0.6796971953334389
protein2,10,MEMIFGQVI,MVENKRLLEG,PGA,genotype_00,48.90184726168596,HLA-B*45:01,0.
↪180375,0.40690354630351067,0.8325127154379427
protein2,10,MEMIFGQVI,MVENKRLLEG,PGA,genotype_01,77.43363203601268,HLA-B*44:02,0.
↪3624999999999998,0.40690354630351067,0.7673069170016388
protein2,4,KRLLEGMEMIF,MVEN,GQVIPGA,genotype_00,88.06509237659311,HLA-C*07:02,0.
↪8623749999999999,0.46066924184560776,0.7842639917310846

```

See the [mhcflurry-predict-scan](#) docs for more options.

## 2.4 Fitting your own models

If you have your own data and want to fit your own MHCflurry models, you have a few options. If you have data for only one or a few MHC I alleles, the best approach is to use the [mhcflurry-class1-train-allele-specific-models](#) command to fit an “allele-specific” predictor, in which separate neural networks are used for each allele.

To call [mhcflurry-class1-train-allele-specific-models](#) you’ll need some training data. The data we use for our released predictors can be downloaded with [mhcflurry-downloads](#):

```
$ mhcflurry-downloads fetch data_curated
```

It looks like this:

```

$ bzcat "$(mhcflurry-downloads path data_curated)/curated_training_data.csv.bz2" |
↪head -n 3
allele,peptide,measurement_value,measurement_inequality,measurement_type,measurement_
↪kind,measurement_source,original_allele
BoLA-1*21:01,AENDTLVSV,7817.0,=,quantitative,affinity,Barlow - purified MHC/
↪competitive/fluorescence,BoLA-1*02101
BoLA-1*21:01,NQFNGGCLLV,1086.0,=,quantitative,affinity,Barlow - purified MHC/direct/
↪fluorescence,BoLA-1*02101

```

Here's an example invocation to fit a predictor:

```
$ mhcflurry-class1-train-allele-specific-models \
  --data curated_training_data.csv.bz2 \
  --hyperparameters hyperparameters.yaml \
  --min-measurements-per-allele 75 \
  --out-models-dir models
```

The `hyperparameters.yaml` file gives the list of neural network architectures to train models for. Here's an example specifying a single architecture:

```
- activation: tanh
  dense_layer_l1_regularization: 0.0
  dropout_probability: 0.0
  early_stopping: true
  layer_sizes: [8]
  locally_connected_layers: []
  loss: custom:mse_with_inequalities
  max_epochs: 500
  minibatch_size: 128
  n_models: 4
  output_activation: sigmoid
  patience: 20
  peptide_amino_acid_encoding: BLOSUM62
  random_negative_affinity_max: 50000.0
  random_negative_affinity_min: 20000.0
  random_negative_constant: 25
  random_negative_rate: 0.0
  validation_split: 0.1
```

The available hyperparameters for binding predictors are defined in [Class1NeuralNetwork](#). To see exactly how these are used you will need to read the source code.

**Note:** MHCflurry predictors are serialized to disk as many files in a directory. The model training command above will write the models to the output directory specified by the `--out-models-dir` argument. This directory has files like:

```
info.txt
manifest.csv
model_selection.csv.bz2
model_selection_data.csv.bz2
...
weights_PATR-B*24:01-6-da42511a2164a8fe.npz
weights_PATR-B*24:01-7-4e3f5ded9cc1d851.npz
weights_PATR-B*24:01-8-491d1b4d85da0dc4.npz
weights_PATR-B*24:01-9-8f295e814502ffa1.npz
```

The `manifest.csv` file gives metadata for all the models used in the predictor. There will be a `weights_...` file for each model giving its weights (the parameters for the neural network). The `percent_ranks.csv` stores a histogram of model predictions for each allele over a large number of random peptides. It is used for generating the percent ranks at prediction time.

To fit pan-allele models like the ones released with MHCflurry, you can use a similar tool, [mhcflurry-class1-train-pan-allele-models](#). You'll probably also want to take a look at the scripts used to generate the production models, which are available in the `downloads-generation` directory in the MHCflurry repository. See the scripts in the `models_class1_pan` subdirectory to see how the fitting and model selection was done for models currently distributed with MHCflurry.

**Note:** The production MHCflurry models were fit using a cluster with several dozen GPUs over a period of about two days. If you model select over fewer architectures, however, it should be possible to fit a predictor using less resources.

---

## 2.5 Environment variables

MHCflurry behavior can be modified using these environment variables:

**MHCFLURRY\_DEFAULT\_CLASS1\_MODELS** Path to models directory. If you call `Class1AffinityPredictor.load()` with no arguments, the models specified in this environment variable will be used. If this environment variable is undefined, the downloaded models for the current MHCflurry release are used.

**MHCFLURRY\_OPTIMIZATION\_LEVEL** The pan-allele models can be somewhat slow. As an optimization, when this variable is greater than 0 (default is 1), we “stitch” the pan-allele models in the ensemble into one large tensorflow graph. In our experiments it gives about a 30% speed improvement. It has no effect on allele-specific models. Set this variable to 0 to disable this behavior. This may be helpful if you are running out of memory using the pan-allele models.

**MHCFLURRY\_DEFAULT\_PREDICT\_BATCH\_SIZE** For large prediction tasks, it can be helpful to increase the prediction batch size, which is set by this environment variable (default is 4096). This affects both allele-specific and pan-allele predictors. It can have large effects on performance. Alternatively, if you are running out of memory, you can try decreasing the batch size.



## PYTHON LIBRARY TUTORIAL

The MHCflurry Python API exposes additional options and features beyond those supported by the commandline tools and can be more convenient for interactive analyses and bioinformatic pipelines. This tutorial gives a basic overview of the most important functionality. See the [API Documentation](#) for further details.

### 3.1 Loading a predictor

Most prediction tasks can be performed using the `Class1PresentationPredictor` class, which provides a programmatic API to the functionality in the `mhcflurry-predict` and `mhcflurry-predict-scan` commands.

Instances of `Class1PresentationPredictor` wrap a `Class1AffinityPredictor` to generate binding affinity predictions and a `Class1ProcessingPredictor` to generate antigen processing predictions. The presentation score is computed using a logistic regression model over binding affinity and processing predictions.

Use the `load` static method to load a trained predictor from disk. With no arguments this method will load the predictor released with MHCflurry (see [Downloading models](#)). If you pass a path to a models directory, then it will load that predictor instead.

```
>>> from mhcflurry import Class1PresentationPredictor
>>> predictor = Class1PresentationPredictor.load()
>>> predictor.supported_alleles[:5]
['Atbe-B*01:01', 'Atbe-E*03:01', 'Atbe-G*03:01', 'Atbe-G*03:02', 'Atbe-G*06:01']
```

### 3.2 Predicting for individual peptides

To generate predictions for individual peptides, we can use the `predict` method of the `Class1PresentationPredictor`, loaded above. This method returns a `pandas.DataFrame` with binding affinity, processing, and presentation predictions:

```
>>> predictor.predict(
...     peptides=["SIINFEKL", "NLVPMVATV"],
...     alleles=["HLA-A0201", "HLA-A0301"],
...     verbose=0)
  peptide  peptide_num sample_name  affinity best_allele  processing_score  presentation_score
↪-----
0  SIINFEKL           0    sample1  12906.786173   HLA-A0201           0.101473           0.012503
↪
1  NLVPMVATV          1    sample1   15.038358   HLA-A0201           0.676289           0.975463
↪
```

Here, the list of alleles is taken to be an individual's MHC I genotype (i.e. up to 6 alleles), and the strongest binder across alleles for each peptide is reported.

**Note:** MHCflurry normalizes allele names using the `mhcnames` package. Names like HLA-A\*02:01 or A\*02:01 will be normalized to HLA-A\*02:01, so most naming conventions can be used with methods such as `predict`.

If you have multiple sample genotypes, you can pass a dict, where the keys are arbitrary sample names:

```
>>> predictor.predict(
...     peptides=["KSEYMTSWFY", "NLVPMVATV"],
...     alleles={
...         "sample1": ["A0201", "A0301", "B0702", "B4402", "C0201", "C0702"],
...         "sample2": ["A0101", "A0206", "B5701", "C0202"],
...     },
...     verbose=0)
  peptide peptide_num sample_name      affinity best_allele  processing_score  presentation_score
0 KSEYMTSWFY           0      sample1  16737.745268      A0301           0.381632           0.026550
1 NLVPMVATV           1      sample1   15.038358      A0201           0.676289           0.975463
2 KSEYMTSWFY           0      sample2   62.540779      A0101           0.381632           0.796731
3 NLVPMVATV           1      sample2   15.765500      A0206           0.676289           0.974439
```

Here the strongest binder for each sample / peptide pair is returned.

Many users will focus on the binding affinity predictions, as the processing and presentation predictions are experimental. If you do use the latter scores, however, when available you should provide the upstream (N-flank) and downstream (C-flank) sequences from the source proteins of the peptides for a small boost in accuracy. To do so, specify the `n_flank` and `c_flank` arguments, which give the flanking sequences for the corresponding peptides:

```
>>> predictor.predict(
...     peptides=["KSEYMTSWFY", "NLVPMVATV"],
...     n_flanks=["NNNNNNN", "SSSSSSSS"],
...     c_flanks=["CCCCCCCC", "YYYYAAAA"],
...     alleles={
...         "sample1": ["A0201", "A0301", "B0702", "B4402", "C0201", "C0702"],
...         "sample2": ["A0101", "A0206", "B5701", "C0202"],
...     },
...     verbose=0)
  peptide  n_flank  c_flank peptide_num sample_name      affinity best_allele  processing_score  presentation_score
0 KSEYMTSWFY NNNNNNN CCCCCCCC           0      sample1  16737.745268      A0301           0.605816           0.056190
1 NLVPMVATV SSSSSSSS  YYYYAAAA           1      sample1   15.038358      A0201           0.824994           0.986719
2 KSEYMTSWFY NNNNNNN CCCCCCCC           0      sample2   62.540779      A0101           0.605816           0.897493
3 NLVPMVATV SSSSSSSS  YYYYAAAA           1      sample2   15.765500      A0206           0.824994           0.986155
```

### 3.3 Scanning protein sequences

The `predict_sequences` method supports scanning protein sequences for MHC ligands. Here's an example to identify all peptides with a predicted binding affinity of 500 nM or tighter to any allele across two sample genotypes and two short peptide sequences.

```
>>> predictor.predict_sequences(
...     sequences={
...         'protein1': "MDSKGSSQKGSRLLLLLLVSNLL",
...         'protein2': "SSLPTPEDKEQAQQTHH",
...     },
...     alleles={
...         "sample1": ["A0201", "A0301", "B0702"],
...         "sample2": ["A0101", "C0202"],
...     },
...     result="filtered",
...     comparison_quantity="affinity",
...     filter_value=500,
...     verbose=0)
sequence_name pos peptide n_flank c_flank sample_name affinity_u
↪best_allele affinity_percentile processing_score presentation_score
0 protein1 13 LLLLVVSNL MDSKGSSQKGSRL L sample1 38.206225 ↪
↪ A0201 0.380125 0.017644 0.571060
1 protein1 14 LLLVVSNNL MDSKGSSQKGSRL sample1 42.243472 ↪
↪ A0201 0.420250 0.090984 0.619213
2 protein1 5 SSQKGSRL MDSKG LLLVVSNNL sample2 66.749223 ↪
↪ C0202 0.803375 0.383608 0.774468
3 protein1 6 SQKGSRLLL MDSKGS LLVVSNNL sample2 178.033467 ↪
↪ C0202 1.820000 0.275019 0.482206
4 protein1 13 LLLLVVSNLL MDSKGSSQKGSRL sample1 202.208167 ↪
↪ A0201 1.112500 0.058782 0.261320
5 protein1 12 LLLLLVVSNNL MDSKGSSQKGSRL sample1 202.506582 ↪
↪ A0201 1.112500 0.010025 0.225648
6 protein2 0 SSLPTPEDK EQAQQTHH sample1 335.529377 ↪
↪ A0301 1.011750 0.010443 0.156798
7 protein2 0 SSLPTPEDK EQAQQTHH sample2 353.451759 ↪
↪ C0202 2.674250 0.010443 0.150753
8 protein1 8 KGSRLLLLL MDSKGSSQ VVSNNL sample2 410.327286 ↪
↪ C0202 2.887000 0.121374 0.194081
9 protein1 5 SSQKGSRL MDSKG LLLLVVSNLL sample2 477.285937 ↪
↪ C0202 3.107375 0.111982 0.168572
```

When using `predict_sequences`, the flanking sequences for each peptide are automatically included in the processing and presentation predictions.

See the documentation for `Class1PresentationPredictor` for other useful methods.

### 3.4 Lower level interfaces

The `Class1PresentationPredictor` delegates to a `Class1AffinityPredictor` instance for binding affinity predictions. If all you need are binding affinities, you can use this instance directly.

Here's an example:

```
>>> from mhcfurry import ClassAffinityPredictor
>>> predictor = ClassAffinityPredictor.load()
>>> predictor.predict_to_dataframe(allele="HLA-A0201", peptides=["SIINFEKL", "SIINFEQL", "SIINFEKL", "SIINFEQL"])
```

	peptide	allele	prediction	prediction_low	prediction_high	prediction_percentile
0	SIINFEKL	HLA-A0201	12906.786173	8829.460289	18029.923061	6.566375
1	SIINFEQL	HLA-A0201	13025.300796	9050.056312	18338.004869	6.623625

The `prediction_low` and `prediction_high` fields give the 5-95 percentile predictions across the models in the ensemble. This detailed information is not available through the higher-level *Class1PresentationPredictor* interface.

Under the hood, `Class1AffinityPredictor` itself delegates to an ensemble of of `Class1NeuralNetwork` instances, which implement the neural network models used for prediction. To fit your own affinity prediction models, call `fit`.

You can similarly use `Class1ProcessingPredictor` directly for antigen processing prediction, and there is a low-level `Class1ProcessingNeuralNetwork` with a `fit` method.

See the API documentation of these classes for details.



## COMMAND-LINE REFERENCE

See also the [tutorial](#).

### 4.1 mhcflurry-predict

Run MHCflurry predictor on specified peptides.

By default, the presentation predictor is used, and predictions for MHC I binding affinity, antigen processing, and the composite presentation score are returned. If you just want binding affinity predictions, pass `--affinity-only`.

Examples:

Write a CSV file containing the contents of INPUT.csv plus additional columns giving MHCflurry predictions:

```
$ mhcflurry-predict INPUT.csv --out RESULT.csv
```

The input CSV file is expected to contain columns “allele”, “peptide”, and, optionally, “n\_flank”, and “c\_flank”.

If `--out` is not specified, results are written to stdout.

You can also run on alleles and peptides specified on the commandline, in which case predictions are written for *all combinations* of alleles and peptides:

```
$ mhcflurry-predict --alleles HLA-A*02:01 H-2Kb --peptides SIINFEKL DENDREKLLL
```

Instead of individual alleles (in a CSV or on the command line), you can also give a comma separated list of alleles giving a sample genotype. In this case, the tightest binding affinity across the alleles for the sample will be returned. For example:

```
$ mhcflurry-predict --peptides SIINFEKL DENDREKLLL --alleles HLA-A*02:01,HLA-A*03:01,HLA-B*57:01,HLA-B*45:01,HLA-C*02:01,HLA-C*07:02 HLA-A*01:01,HLA-A*02:06,HLA-B*44:02,HLA-B*07:02,HLA-C*01:01,HLA-C*03:01
```

will give the tightest predicted affinities across alleles for each of the two genotypes specified for each peptide.

```
usage: mhcflurry-predict [-h] [--list-supported-alleles]
                        [--list-supported-peptide-lengths] [--version]
                        [--alleles ALLELE [ALLELE ...]]
                        [--peptides PEPTIDE [PEPTIDE ...]]
                        [--allele-column NAME] [--peptide-column NAME]
                        [--n-flank-column NAME] [--c-flank-column NAME]
                        [--no-throw] [--out OUTPUT.csv]
                        [--prediction-column-prefix NAME]
                        [--output-delimiter CHAR] [--no-affinity-percentile]
                        [--always-include-best-allele] [--models DIR]
                        [--affinity-only] [--no-flanking]
                        [INPUT.csv]
```

**input.csv**  
Input CSV

**-h, --help**  
Show this help message and exit

**--list-supported-alleles**  
Prints the list of supported alleles and exits

**--list-supported-peptide-lengths**  
Prints the list of supported peptide lengths and exits

**--version**  
show program's version number and exit

**--alleles <allele>**  
Alleles to predict (exclusive with passing an input CSV)

**--peptides <peptide>**  
Peptides to predict (exclusive with passing an input CSV)

**--allele-column <name>**  
Input column name for alleles. Default: 'allele'

**--peptide-column <name>**  
Input column name for peptides. Default: 'peptide'

**--n-flank-column <name>**  
Column giving N-terminal flanking sequence. Default: 'n\_flank'

**--c-flank-column <name>**  
Column giving C-terminal flanking sequence. Default: 'c\_flank'

**--no-throw**  
Return NaNs for unsupported alleles or peptides instead of raising

**--out <output.csv>**  
Output CSV

**--prediction-column-prefix <name>**  
Prefix for output column names. Default: 'mhcflurry\_'

**--output-delimiter <char>**  
Delimiter character for results. Default: ','

**--no-affinity-percentile**  
Do not include affinity percentile rank

**--always-include-best-allele**  
Always include the best\_allele column even when it is identical to the allele column (i.e. all queries are monoallelic).

**--models <dir>**  
Directory containing models. Either a binding affinity predictor or a presentation predictor can be used. Default: /Users/tim/Library/Application Support/mhcflurry/4/1.6.0/models\_class1\_presentation/models

**--affinity-only**  
Affinity prediction only (no antigen processing or presentation)

**--no-flanking**  
Do not use flanking sequence information even when available

## 4.2 mhcflurry-predict-scan

Scan protein sequences using the MHCflurry presentation predictor.

By default, sub-sequences (peptides) with affinity percentile ranks less than 2.0 are returned. You can also specify `--results-all` to return predictions for all peptides, or `--results-best` to return the top peptide for each sequence.

Examples:

Scan a set of sequences in a FASTA file for binders to any alleles in a MHC I genotype:

```
$ mhcflurry-predict-scan test/data/example.fasta --alleles HLA-A*02:01,HLA-A*03:01,HLA-B*57:01,HLA-B*45:01,HLA-C*02:01,HLA-C*07:02
```

Instead of a FASTA, you can also pass a CSV that has “sequence\_id” and “sequence” columns.

You can also specify multiple MHC I genotypes to scan as space-separated arguments to the `--alleles` option:

```
$ mhcflurry-predict-scan test/data/example.fasta --alleles HLA-A*02:01,HLA-A*03:01,HLA-B*57:01,HLA-B*45:01,HLA-C*02:02,HLA-C*07:02 HLA-A*01:01,HLA-A*02:06,HLA-B*44:02,HLA-B*07:02,HLA-C*01:02,HLA-C*03:01
```

If `--out` is not specified, results are written to standard out.

You can also specify sequences on the commandline:

```
mhcflurry-predict-scan --sequences MGYINVFAFPFTIYSLLLCRMNSRNYIAQVDVVNFNLT --alleles HLA-A*02:01,HLA-A*03:01,HLA-B*57:01,HLA-B*45:01,HLA-C*02:02,HLA-C*07:02
```

```
usage: mhcflurry-predict-scan [-h] [--list-supported-alleles]
                             [--list-supported-peptide-lengths] [--version]
                             [--input-format {guess, csv, fasta}]
                             [--alleles ALLELE [ALLELE ...]]
                             [--sequences SEQ [SEQ ...]]
                             [--sequence-id-column NAME]
                             [--sequence-column NAME] [--no-throw]
                             [--peptide-lengths PEPTIDE_LENGTHS [PEPTIDE_LENGTHS ...]]
                             [--results-all]
                             [--results-best {presentation_score, processing_score,
↪ affinity, affinity_percentile}]
                             [--results-filtered {presentation_score, processing_
↪ score, affinity, affinity_percentile}]
                             [--threshold-presentation-score THRESHOLD_PRESENTATION_
↪ SCORE]
                             [--threshold-processing-score THRESHOLD_PROCESSING_
↪ SCORE]
                             [--threshold-affinity THRESHOLD_AFFINITY]
                             [--threshold-affinity-percentile THRESHOLD_AFFINITY_
↪ PERCENTILE]
                             [--out OUTPUT.csv] [--output-delimiter CHAR]
                             [--no-affinity-percentile] [--models DIR]
                             [--no-flanking]
                             [INPUT]
```

### input

Input CSV or FASTA

### -h, --help

Show this help message and exit

**--list-supported-alleles**  
Print the list of supported alleles and exits

**--list-supported-peptide-lengths**  
Print the list of supported peptide lengths and exits

**--version**  
show program's version number and exit

**--input-format** {guess, csv, fasta}  
Format of input file. By default, it is guessed from the file extension.

**--alleles** <allele>  
Alleles to predict

**--sequences** <seq>  
Sequences to predict (exclusive with passing an input file)

**--sequence-id-column** <name>  
Input CSV column name for sequence IDs. Default: 'sequence\_id'

**--sequence-column** <name>  
Input CSV column name for sequences. Default: 'sequence'

**--no-throw**  
Return NaNs for unsupported alleles or peptides instead of raising

**--peptide-lengths** <peptide\_lengths>  
Peptide lengths to consider. Default: [8, 9, 10, 11].

**--results-all**  
Return results for all peptides regardless of affinity, etc.

**--results-best** {presentation\_score, processing\_score, affinity, affinity\_percentile}  
Take the top result for each sequence according to the specified predicted quantity

**--results-filtered** {presentation\_score, processing\_score, affinity, affinity\_percentile}  
Filter results by the specified quantity.

**--threshold-presentation-score** <threshold\_presentation\_score>  
Threshold if filtering by presentation score. Default: 0.7

**--threshold-processing-score** <threshold\_processing\_score>  
Threshold if filtering by processing score. Default: 0.5

**--threshold-affinity** <threshold\_affinity>  
Threshold if filtering by affinity. Default: 500

**--threshold-affinity-percentile** <threshold\_affinity\_percentile>  
Threshold if filtering by affinity percentile. Default: 2.0

**--out** <output.csv>  
Output CSV

**--output-delimiter** <char>  
Delimiter character for results. Default: ','

**--no-affinity-percentile**  
Do not include affinity percentile rank

**--models** <dir>  
Directory containing presentation models. Default: /Users/tim/Library/Application Support/mhcflurry/4/1.6.0/models\_class1\_presentation/models

**--no-flanking**

Do not use flanking sequence information in predictions

## 4.3 mhcflurry-downloads

Download MHCflurry released datasets and trained models.

Examples

**Fetch the default downloads:** \$ mhcflurry-downloads fetch**Fetch a specific download:** \$ mhcflurry-downloads fetch models\_class1\_pan**Get the path to a download:** \$ mhcflurry-downloads path models\_class1\_pan**Get the URL of a download:** \$ mhcflurry-downloads url models\_class1\_pan**Summarize available and fetched downloads:** \$ mhcflurry-downloads info

```
usage: mhcflurry-downloads [-h] [--quiet] [--verbose]
                           {fetch,info,path,url} ...
```

**-h, --help**

show this help message and exit

**--quiet**

Output less

**--verbose, -v**

Output more

### 4.3.1 mhcflurry-downloads fetch

```
usage: mhcflurry-downloads fetch [-h] [--keep] [--release RELEASE]
                                  [--already-downloaded-dir DIR]
                                  [DOWNLOAD [DOWNLOAD ...]]
```

**download**

Items to download

**-h, --help**

show this help message and exit

**--keep**

Don't delete archives after they are extracted

**--release <release>**

Release to download. Default: 1.6.0

**--already-downloaded-dir <dir>**

Don't download files, get them from DIR

### 4.3.2 mhcflurry-downloads info

```
usage: mhcflurry-downloads info [-h]
```

**-h, --help**  
show this help message and exit

### 4.3.3 mhcflurry-downloads path

```
usage: mhcflurry-downloads path [-h] [download_name]
```

**download\_name**

**-h, --help**  
show this help message and exit

### 4.3.4 mhcflurry-downloads url

```
usage: mhcflurry-downloads url [-h] [download_name]
```

**download\_name**

**-h, --help**  
show this help message and exit

## 4.4 mhcflurry-class1-train-allele-specific-models

```
usage:  
Train Class1 single allele models.
```

**-h, --help**  
show this help message and exit

**--data** <file.csv>  
Training data CSV. Expected columns: allele, peptide, measurement\_value

**--out-models-dir** <dir>  
Directory to write models and manifest

**--hyperparameters** <file.json>  
JSON or YAML of hyperparameters

**--allele** <allele>  
Alleles to train models for. If not specified, all alleles with enough measurements will be used.

**--min-measurements-per-allele** <n>  
Train models for alleles with >=N measurements.

**--held-out-fraction-reciprocal** <n>  
Hold out 1/N fraction of data (for e.g. subsequent model selection. For example, specify 5 to hold out 20 percent of the data.

**--held-out-fraction-seed** <n>  
Seed for randomizing which measurements are held out. Only matters when --held-out-fraction is specified. Default: 0.

**--ignore-inequalities**  
Do not use affinity value inequalities even when present in data

**--n-models** <n>  
Ensemble size, i.e. how many models to train for each architecture. If specified here it overrides any 'n\_models' specified in the hyperparameters.

**--max-epochs** <n>  
Max training epochs. If specified here it overrides any 'max\_epochs' specified in the hyperparameters.

**--allele-sequences** <file.csv>  
Allele sequences file. Used for computing allele similarity matrix.

**--save-interval** <n>  
Write models to disk every N seconds. Only affects parallel runs; serial runs write each model to disk as it is trained.

**--verbosity** <verbosity>  
Keras verbosity. Default: 0

**--num-jobs** <n>  
Number of local processes to parallelize training over. Set to 0 for serial run. Default: 0.

**--backend** {tensorflow-gpu,tensorflow-cpu,tensorflow-default}  
Keras backend. If not specified will use system default.

**--gpus** <n>  
Number of GPUs to attempt to parallelize across. Requires running in parallel.

**--max-workers-per-gpu** <n>  
Maximum number of workers to assign to a GPU. Additional tasks will run on CPU.

**--max-tasks-per-worker** <n>  
Restart workers after N tasks. Workaround for tensorflow memory leaks. Requires Python >=3.2.

**--worker-log-dir** <worker\_log\_dir>  
Write worker stdout and stderr logs to given directory.

## 4.5 mhcflurry-class1-select-allele-specific-models

```
usage:
Model select class1 single allele models.
```

**-h, --help**  
show this help message and exit

**--data** <file.csv>  
Model selection data CSV. Expected columns: allele, peptide, measurement\_value

**--exclude-data** <file.csv>  
Data to EXCLUDE from model selection. Useful to specify the original training data used

**--models-dir** <dir>  
Directory to read models

**--out-models-dir** <dir>  
Directory to write selected models

**--out-unselected-predictions** <file.csv>  
Write predictions for validation data using unselected predictor to FILE.csv

**--unselected-accuracy-scorer** <scorer>  
**--unselected-accuracy-scorer-num-samples** <unselected\_accuracy\_scorer\_num\_samples>  
**--unselected-accuracy-percentile-threshold** <x>  
**--allele** <allele>  
Alleles to select models for. If not specified, all alleles with enough measurements will be used.  
**--combined-min-models** <n>  
Min number of models to select per allele when using combined selector  
**--combined-max-models** <n>  
Max number of models to select per allele when using combined selector  
**--combined-min-contribution-percent** <x>  
Use only model selectors that can contribute at least X % to the total score. Default: 1.0  
**--mass-spec-min-measurements** <n>  
Min number of measurements required for an allele to use mass-spec model selection  
**--mass-spec-min-models** <n>  
Min number of models to select per allele when using mass-spec selector  
**--mass-spec-max-models** <n>  
Max number of models to select per allele when using mass-spec selector  
**--mse-min-measurements** <n>  
Min number of measurements required for an allele to use MSE model selection  
**--mse-min-models** <n>  
Min number of models to select per allele when using MSE selector  
**--mse-max-models** <n>  
Max number of models to select per allele when using MSE selector  
**--scoring** <scoring>  
Scoring procedures to use in order  
**--consensus-min-models** <n>  
Min number of models to select per allele when using consensus selector  
**--consensus-max-models** <n>  
Max number of models to select per allele when using consensus selector  
**--consensus-num-peptides-per-length** <consensus\_num\_peptides\_per\_length>  
Num peptides per length to use for consensus scoring  
**--mass-spec-regex** <regex>  
Regular expression for mass-spec data. Runs on measurement\_source col. Default: mass[- ]spec.  
**--verbosity** <verbosity>  
Keras verbosity. Default: 0  
**--num-jobs** <n>  
Number of local processes to parallelize training over. Set to 0 for serial run. Default: 0.  
**--backend** {tensorflow-gpu,tensorflow-cpu,tensorflow-default}  
Keras backend. If not specified will use system default.  
**--gpus** <n>  
Number of GPUs to attempt to parallelize across. Requires running in parallel.  
**--max-workers-per-gpu** <n>  
Maximum number of workers to assign to a GPU. Additional tasks will run on CPU.



**--max-tasks-per-worker** <n>  
Restart workers after N tasks. Workaround for tensorflow memory leaks. Requires Python >=3.2.

**--worker-log-dir** <worker\_log\_dir>  
Write worker stdout and stderr logs to given directory.

## 4.6 mhcflurry-class1-train-pan-allele-models

```
usage:
Train Class1 pan-allele models.
```

**-h, --help**  
show this help message and exit

**--data** <file.csv>  
Training data CSV. Expected columns: allele, peptide, measurement\_value

**--pretrain-data** <file.csv>  
Pre-training data CSV. Expected columns: allele, peptide, measurement\_value

**--out-models-dir** <dir>  
Directory to write models and manifest

**--hyperparameters** <file.json>  
JSON or YAML of hyperparameters

**--held-out-measurements-per-allele-fraction-and-max** <x>  
Fraction of measurements per allele to hold out, and maximum number

**--ignore-inequalities**  
Do not use affinity value inequalities even when present in data

**--num-folds** <n>  
Number of training folds.

**--num-replicates** <n>  
Number of replicates per (architecture, fold) pair to train.

**--max-epochs** <n>  
Max training epochs. If specified here it overrides any 'max\_epochs' specified in the hyperparameters.

**--allele-sequences** <file.csv>  
Allele sequences file.

**--verbosity** <verbosity>  
Keras verbosity. Default: 0

**--debug**  
Launch python debugger on error

**--continue-incomplete**  
Continue training models from an incomplete training run. If this is specified then the only required argument is --out-models-dir

**--only-initialize**  
Do not actually train models. The initialized run can be continued later with --continue-incomplete.

**--num-jobs** <n>  
Number of local processes to parallelize training over. Set to 0 for serial run. Default: 0.

**--backend** {tensorflow-gpu,tensorflow-cpu,tensorflow-default}  
Keras backend. If not specified will use system default.

**--gpus** <n>  
Number of GPUs to attempt to parallelize across. Requires running in parallel.

**--max-workers-per-gpu** <n>  
Maximum number of workers to assign to a GPU. Additional tasks will run on CPU.

**--max-tasks-per-worker** <n>  
Restart workers after N tasks. Workaround for tensorflow memory leaks. Requires Python >=3.2.

**--worker-log-dir** <worker\_log\_dir>  
Write worker stdout and stderr logs to given directory.

**--cluster-parallelism**

**--cluster-submit-command** <cluster\_submit\_command>  
Default: sh

**--cluster-results-workdir** <cluster\_results\_workdir>  
Default: ./cluster-workdir

**--additional-complete-file** <additional\_complete\_file>  
Additional file to monitor for job completion. Default: STDERR

**--cluster-script-prefix-path** <cluster\_script\_prefix\_path>

**--cluster-max-retries** <cluster\_max\_retries>  
How many times to rerun failing jobs. Default: 3

## 4.7 mhcflurry-class1-select-pan-allele-models

```
usage:
Model select class1 pan-allele models.

APPROACH: For each training fold, we select at least min and at most max models
(where min and max are set by the --{min/max}-models-per-fold argument) using a
step-up (forward) selection procedure. The final ensemble is the union of all
selected models across all folds.
```

**-h, --help**  
show this help message and exit

**--data** <file.csv>  
Model selection data CSV. Expected columns: allele, peptide, measurement\_value

**--models-dir** <dir>  
Directory to read models

**--out-models-dir** <dir>  
Directory to write selected models

**--min-models-per-fold** <n>  
Min number of models to select per fold

**--max-models-per-fold** <n>  
Max number of models to select per fold

**--mass-spec-regex** <regex>  
Regular expression for mass-spec data. Runs on measurement\_source col.Default: mass[- ]spec.

**--verbosity** <verbosity>  
Keras verbosity. Default: 0

**--num-jobs** <n>  
Number of local processes to parallelize training over. Set to 0 for serial run. Default: 0.

**--backend** {tensorflow-gpu,tensorflow-cpu,tensorflow-default}  
Keras backend. If not specified will use system default.

**--gpus** <n>  
Number of GPUs to attempt to parallelize across. Requires running in parallel.

**--max-workers-per-gpu** <n>  
Maximum number of workers to assign to a GPU. Additional tasks will run on CPU.

**--max-tasks-per-worker** <n>  
Restart workers after N tasks. Workaround for tensorflow memory leaks. Requires Python >=3.2.

**--worker-log-dir** <worker\_log\_dir>  
Write worker stdout and stderr logs to given directory.

**--cluster-parallelism**

**--cluster-submit-command** <cluster\_submit\_command>  
Default: sh

**--cluster-results-workdir** <cluster\_results\_workdir>  
Default: ./cluster-workdir

**--additional-complete-file** <additional\_complete\_file>  
Additional file to monitor for job completion. Default: STDERR

**--cluster-script-prefix-path** <cluster\_script\_prefix\_path>

**--cluster-max-retries** <cluster\_max\_retries>  
How many times to rerun failing jobs. Default: 3

## 4.8 mhcflurry-class1-train-processing-models

```
usage:
Train Class1 processing models.
```

**-h, --help**  
show this help message and exit

**--data** <file.csv>  
Training data CSV. Expected columns: peptide, n\_flank, c\_flank, hit

**--out-models-dir** <dir>  
Directory to write models and manifest

**--hyperparameters** <file.json>  
JSON or YAML of hyperparameters

**--held-out-samples** <n>  
Number of experiments to hold out per fold

**--num-folds** <n>  
Number of training folds.

**--num-replicates** <n>  
Number of replicates per (architecture, fold) pair to train.

**--max-epochs** <n>  
Max training epochs. If specified here it overrides any 'max\_epochs' specified in the hyperparameters.

**--verbosity** <verbosity>  
Keras verbosity. Default: 0

**--debug**  
Launch python debugger on error

**--continue-incomplete**  
Continue training models from an incomplete training run. If this is specified then the only required argument is `--out-models-dir`

**--only-initialize**  
Do not actually train models. The initialized run can be continued later with `--continue-incomplete`.

**--num-jobs** <n>  
Number of local processes to parallelize training over. Set to 0 for serial run. Default: 0.

**--backend** {tensorflow-gpu,tensorflow-cpu,tensorflow-default}  
Keras backend. If not specified will use system default.

**--gpus** <n>  
Number of GPUs to attempt to parallelize across. Requires running in parallel.

**--max-workers-per-gpu** <n>  
Maximum number of workers to assign to a GPU. Additional tasks will run on CPU.

**--max-tasks-per-worker** <n>  
Restart workers after N tasks. Workaround for tensorflow memory leaks. Requires Python >=3.2.

**--worker-log-dir** <worker\_log\_dir>  
Write worker stdout and stderr logs to given directory.

**--cluster-parallelism**

**--cluster-submit-command** <cluster\_submit\_command>  
Default: sh

**--cluster-results-workdir** <cluster\_results\_workdir>  
Default: ./cluster-workdir

**--additional-complete-file** <additional\_complete\_file>  
Additional file to monitor for job completion. Default: STDERR

**--cluster-script-prefix-path** <cluster\_script\_prefix\_path>

**--cluster-max-retries** <cluster\_max\_retries>  
How many times to rerun failing jobs. Default: 3

## 4.9 mhcflurry-class1-select-processing-models

```
usage:
Model select antigen processing models.

APPROACH: For each training fold, we select at least min and at most max models
(where min and max are set by the --{min/max}-models-per-fold argument) using a
step-up (forward) selection procedure. The final ensemble is the union of all
selected models across all folds. AUC is used as the metric.
```

**-h, --help**  
show this help message and exit

**--data <file.csv>**  
Model selection data CSV. Expected columns: peptide, hit, fold\_0, ..., fold\_N

**--models-dir <dir>**  
Directory to read models

**--out-models-dir <dir>**  
Directory to write selected models

**--min-models-per-fold <n>**  
Min number of models to select per fold

**--max-models-per-fold <n>**  
Max number of models to select per fold

**--verbosity <verbosity>**  
Keras verbosity. Default: 0

**--num-jobs <n>**  
Number of local processes to parallelize training over. Set to 0 for serial run. Default: 0.

**--backend {tensorflow-gpu,tensorflow-cpu,tensorflow-default}**  
Keras backend. If not specified will use system default.

**--gpus <n>**  
Number of GPUs to attempt to parallelize across. Requires running in parallel.

**--max-workers-per-gpu <n>**  
Maximum number of workers to assign to a GPU. Additional tasks will run on CPU.

**--max-tasks-per-worker <n>**  
Restart workers after N tasks. Workaround for tensorflow memory leaks. Requires Python >=3.2.

**--worker-log-dir <worker\_log\_dir>**  
Write worker stdout and stderr logs to given directory.

**--cluster-parallelism**

**--cluster-submit-command <cluster\_submit\_command>**  
Default: sh

**--cluster-results-workdir <cluster\_results\_workdir>**  
Default: ./cluster-workdir

**--additional-complete-file <additional\_complete\_file>**  
Additional file to monitor for job completion. Default: STDERR

**--cluster-script-prefix-path <cluster\_script\_prefix\_path>**

**--cluster-max-retries** <cluster\_max\_retries>  
How many times to rerun failing jobs. Default: 3

## 4.10 mhcflurry-class1-train-presentation-models

usage:  
Train Class1 presentation models.

**-h, --help**  
show this help message and exit

**--data** <file.csv>  
Training data CSV. Expected columns: peptide, n\_flank, c\_flank, hit

**--out-models-dir** <dir>  
Directory to write models and manifest

**--affinity-predictor** <dir>  
Affinity predictor models dir

**--processing-predictor-with-flanks** <dir>  
Processing predictor with flanks

**--processing-predictor-without-flanks** <dir>  
Processing predictor without flanks

**--verbosity** <verbosity>  
Default: 1

**--debug**  
Launch python debugger on error

**--hla-column** <hla\_column>  
Column in data giving space-separated MHC I alleles

**--target-column** <target\_column>  
Column in data giving hit (1) vs decoy (0)

## API DOCUMENTATION

Class I MHC ligand prediction package

```
class mhcflurry.Class1AffinityPredictor (allele_to_allele_specific_models=None,  
                                         class1_pan_allele_models=None, al-  
                                         lele_to_sequence=None, manifest_df=None,  
                                         allele_to_percent_rank_transform=None, meta-  
                                         data_dataframes=None)
```

Bases: object

High-level interface for peptide/MHC I binding affinity prediction.

This class manages low-level *Class1NeuralNetwork* instances, each of which wraps a single Keras network. The purpose of *Class1AffinityPredictor* is to implement ensembles, handling of multiple alleles, and predictor loading and saving. It also provides a place to keep track of metadata like prediction histograms for percentile rank calibration.

#### Parameters

**allele\_to\_allele\_specific\_models** [dict of string -> list of *Class1NeuralNetwork*] Ensemble of single-allele models to use for each allele.

**class1\_pan\_allele\_models** [list of *Class1NeuralNetwork*] Ensemble of pan-allele models.

**allele\_to\_sequence** [dict of string -> string] MHC allele name to fixed-length amino acid sequence (sometimes referred to as the pseudosequence). Required only if class1\_pan\_allele\_models is specified.

**manifest\_df** [pandas.DataFrame, optional] Must have columns: model\_name, allele, config\_json, model. Only required if you want to update an existing serialization of a Class1AffinityPredictor. Otherwise this dataframe will be generated automatically based on the supplied models.

**allele\_to\_percent\_rank\_transform** [dict of string -> PercentRankTransform, optional] PercentRankTransform instances to use for each allele

**metadata\_dataframes** [dict of string -> pandas.DataFrame, optional] Optional additional dataframes to write to the models dir when save() is called. Useful for tracking provenance.

#### property manifest\_df

A pandas.DataFrame describing the models included in this predictor.

Based on: - self.class1\_pan\_allele\_models - self.allele\_to\_allele\_specific\_models

#### Returns

**pandas.DataFrame**

**clear\_cache** (*self*)

Clear values cached based on the neural networks in this predictor.

**Users should call this after mutating any of the following:**

- self.class1\_pan\_allele\_models
- self.allele\_to\_allele\_specific\_models
- self.allele\_to\_sequence

Methods that mutate these instance variables will call this method on their own if needed.

**property neural\_networks**

List of the neural networks in the ensemble.

**Returns**

list of *Class1NeuralNetwork*

**classmethod merge** (*predictors*)

Merge the ensembles of two or more *Class1AffinityPredictor* instances.

Note: the resulting merged predictor will NOT have calibrated percentile ranks. Call *calibrate\_percentile\_ranks* on it if these are needed.

**Parameters**

**predictors** [sequence of *Class1AffinityPredictor*]

**Returns**

*Class1AffinityPredictor* instance

**merge\_in\_place** (*self, others*)

Add the models present in other predictors into the current predictor.

**Parameters**

**others** [list of *Class1AffinityPredictor*] Other predictors to merge into the current predictor.

**Returns**

list of string [names of newly added models]

**property supported\_alleles**

Alleles for which predictions can be made.

**Returns**

list of string

**property supported\_peptide\_lengths**

(minimum, maximum) lengths of peptides supported by *all models*, inclusive.

**Returns**

(int, int) tuple

**check\_consistency** (*self*)

Verify that self.manifest\_df is consistent with: - self.class1\_pan\_allele\_models - self.allele\_to\_allele\_specific\_models

Currently only checks for agreement on the total number of models.

Throws AssertionError if inconsistent.



**save** (*self*, *models\_dir*, *model\_names\_to\_write=None*, *write\_metadata=True*)

Serialize the predictor to a directory on disk. If the directory does not exist it will be created.

The serialization format consists of a file called “manifest.csv” with the configurations of each `Class1NeuralNetwork`, along with per-network files giving the model weights. If there are pan-allele predictors in the ensemble, the allele sequences are also stored in the directory. There is also a small file “index.txt” with basic metadata: when the models were trained, by whom, on what host.

#### Parameters

**models\_dir** [string] Path to directory. It will be created if it doesn’t exist.

**model\_names\_to\_write** [list of string, optional] Only write the weights for the specified models. Useful for incremental updates during training.

**write\_metadata** [boolean, optional] Whether to write optional metadata

**static load** (*models\_dir=None*, *max\_models=None*, *optimization\_level=None*)

Deserialize a predictor from a directory on disk.

#### Parameters

**models\_dir** [string] Path to directory. If unspecified the default downloaded models are used.

**max\_models** [int, optional] Maximum number of `Class1NeuralNetwork` instances to load

**optimization\_level** [int] If >0, model optimization will be attempted. Defaults to value of environment variable MHCFLURRY\_OPTIMIZATION\_LEVEL.

#### Returns

`Class1AffinityPredictor` instance

**optimize** (*self*, *warn=True*)

EXPERIMENTAL: Optimize the predictor for faster predictions.

Currently the only optimization implemented is to merge multiple pan- allele predictors at the tensorflow level.

The optimization is performed in-place, mutating the instance.

#### Returns

**bool** Whether optimization was performed

**static model\_name** (*allele*, *num*)

Generate a model name

#### Parameters

**allele** [string]

**num** [int]

#### Returns

**string**

**static weights\_path** (*models\_dir*, *model\_name*)

Generate the path to the weights file for a model

#### Parameters

**models\_dir** [string]

**model\_name** [string]

**Returns**

string

**property master\_allele\_encoding**

An AlleleEncoding containing the universe of alleles specified by self.allele\_to\_sequence.

**Returns**

AlleleEncoding

**fit\_allele\_specific\_predictors** (*self*, *n\_models*, *architecture\_hyperparameters\_list*,  
*allele*, *peptides*, *affinities*, *inequalities=None*,  
*train\_rounds=None*, *models\_dir\_for\_save=None*,  
*verbose=0*, *progress\_preamble="*,  
*progress\_print\_interval=5.0*)

Fit one or more allele specific predictors for a single allele using one or more neural network architectures.

The new predictors are saved in the Class1AffinityPredictor instance and will be used on subsequent calls to *predict*.

**Parameters**

**n\_models** [int] Number of neural networks to fit

**architecture\_hyperparameters\_list** [list of dict] List of hyperparameter sets.

**allele** [string]

**peptides** [EncodableSequences or list of string]

**affinities** [list of float] nM affinities

**inequalities** [list of string, each element one of ">", "<", or "="] See *Class1NeuralNetwork.fit* for details.

**train\_rounds** [sequence of int] Each training point *i* will be used on training rounds *r* for which *train\_rounds*[*i*] > *r*, *r* >= 0.

**models\_dir\_for\_save** [string, optional] If specified, the Class1AffinityPredictor is (incrementally) written to the given models dir after each neural network is fit.

**verbose** [int] Keras verbosity

**progress\_preamble** [string] Optional string of information to include in each progress update

**progress\_print\_interval** [float] How often (in seconds) to print progress. Set to None to disable.

**Returns**

list of *Class1NeuralNetwork*

**fit\_class1\_pan\_allele\_models** (*self*, *n\_models*, *architecture\_hyperparameters*, *alleles*, *peptides*, *affinities*, *inequalities*, *models\_dir\_for\_save=None*, *verbose=1*, *progress\_preamble="*, *progress\_print\_interval=5.0*)

Fit one or more pan-allele predictors using a single neural network architecture.

The new predictors are saved in the Class1AffinityPredictor instance and will be used on subsequent calls to *predict*.

**Parameters**

**n\_models** [int] Number of neural networks to fit

**architecture\_hyperparameters** [dict]

**alleles** [list of string] Allele names (not sequences) corresponding to each peptide

**peptides** [EncodableSequences or list of string]

**affinities** [list of float] nM affinities

**inequalities** [list of string, each element one of “>”, “<”, or “=”] See `Class1NeuralNetwork.fit` for details.

**models\_dir\_for\_save** [string, optional] If specified, the `Class1AffinityPredictor` is (incrementally) written to the given models dir after each neural network is fit.

**verbose** [int] Keras verbosity

**progress\_preamble** [string] Optional string of information to include in each progress update

**progress\_print\_interval** [float] How often (in seconds) to print progress. Set to `None` to disable.

#### Returns

list of *Class1NeuralNetwork*

**add\_pan\_allele\_model** (*self*, *model*, *models\_dir\_for\_save=None*)

Add a pan-allele model to the ensemble and optionally do an incremental save.

#### Parameters

**model** [Class1NeuralNetwork]

**models\_dir\_for\_save** [string] Directory to save resulting ensemble to

**percentile\_ranks** (*self*, *affinities*, *allele=None*, *alleles=None*, *throw=True*)

Return percentile ranks for the given ic50 affinities and alleles.

The ‘allele’ and ‘alleles’ argument are as in the *predict* method. Specify one of these.

#### Parameters

**affinities** [sequence of float] nM affinities

**allele** [string]

**alleles** [sequence of string]

**throw** [boolean] If True, a `ValueError` will be raised in the case of unsupported alleles. If False, a warning will be logged and `NaN` will be returned for those percentile ranks.

#### Returns

`numpy.array` of float

**predict** (*self*, *peptides*, *alleles=None*, *allele=None*, *throw=True*, *centrality\_measure='mean'*, *model\_kwargs={}*)

Predict nM binding affinities.

If multiple predictors are available for an allele, the predictions are the geometric means of the individual model (nM) predictions.

One of ‘allele’ or ‘alleles’ must be specified. If ‘allele’ is specified all predictions will be for the given allele. If ‘alleles’ is specified it must be the same length as ‘peptides’ and give the allele corresponding to each peptide.

#### Parameters

**peptides** [EncodableSequences or list of string]

**alleles** [list of string]

**allele** [string]

**throw** [boolean] If True, a ValueError will be raised in the case of unsupported alleles or peptide lengths. If False, a warning will be logged and the predictions for the unsupported alleles or peptides will be NaN.

**centrality\_measure** [string or callable] Measure of central tendency to use to combine predictions in the ensemble. Options include: mean, median, robust\_mean.

**model\_kwargs** [dict] Additional keyword arguments to pass to Class1NeuralNetwork.predict

#### Returns

**numpy.array of predictions**

```
predict_to_dataframe (self, peptides, alleles=None, allele=None, throw=True,
                        include_individual_model_predictions=False,
                        include_percentile_ranks=True, include_confidence_intervals=True,
                        centrality_measure='mean', model_kwargs={})
```

Predict nM binding affinities. Gives more detailed output than *predict* method, including 5-95% prediction intervals.

If multiple predictors are available for an allele, the predictions are the geometric means of the individual model predictions.

One of ‘allele’ or ‘alleles’ must be specified. If ‘allele’ is specified all predictions will be for the given allele. If ‘alleles’ is specified it must be the same length as ‘peptides’ and give the allele corresponding to each peptide.

#### Parameters

**peptides** [EncodableSequences or list of string]

**alleles** [list of string]

**allele** [string]

**throw** [boolean] If True, a ValueError will be raised in the case of unsupported alleles or peptide lengths. If False, a warning will be logged and the predictions for the unsupported alleles or peptides will be NaN.

**include\_individual\_model\_predictions** [boolean] If True, the predictions of each individual model are included as columns in the result DataFrame.

**include\_percentile\_ranks** [boolean, default True] If True, a “prediction\_percentile” column will be included giving the percentile ranks. If no percentile rank info is available, this will be ignored with a warning.

**centrality\_measure** [string or callable] Measure of central tendency to use to combine predictions in the ensemble. Options include: mean, median, robust\_mean.

**model\_kwargs** [dict] Additional keyword arguments to pass to Class1NeuralNetwork.predict

#### Returns

**pandas.DataFrame of predictions**

**calibrate\_percentile\_ranks** (*self*, *peptides=None*, *num\_peptides\_per\_length=100000*, *alleles=None*, *bins=None*, *motif\_summary=False*, *summary\_top\_peptide\_fractions=[0.001]*, *verbose=False*, *model\_kwargs={}*)

Compute the cumulative distribution of ic50 values for a set of alleles over a large universe of random peptides, to enable taking quantiles of this distribution later.

#### Parameters

**peptides** [sequence of string or EncodableSequences, optional] Peptides to use

**num\_peptides\_per\_length** [int, optional] If peptides argument is not specified, then num\_peptides\_per\_length peptides are randomly sampled from a uniform distribution for each supported length

**alleles** [sequence of string, optional] Alleles to perform calibration for. If not specified all supported alleles will be calibrated.

**bins** [object] Anything that can be passed to numpy.histogram's "bins" argument can be used here, i.e. either an integer or a sequence giving bin edges. This is in ic50 space.

**motif\_summary** [bool] If True, the length distribution and per-position amino acid frequencies are also calculated for the top x fraction of tightest-binding peptides, where each value of x is given in the summary\_top\_peptide\_fractions list.

**summary\_top\_peptide\_fractions** [list of float] Only used if motif\_summary is True

**verbose** [boolean] Whether to print status updates to stdout

**model\_kwargs** [dict] Additional low-level Class1NeuralNetwork.predict() kwargs.

#### Returns

dict of string -> pandas.DataFrame

If motif\_summary is True, this will have keys "frequency\_matrices" and "length\_distributions". Otherwise it will be empty.

**model\_select** (*self*, *score\_function*, *alleles=None*, *min\_models=1*, *max\_models=10000*)

Perform model selection using a user-specified scoring function.

This works only with allele-specific models, not pan-allele models.

Model selection is done using a "step up" variable selection procedure, in which models are repeatedly added to an ensemble until the score stops improving.

#### Parameters

**score\_function** [Class1AffinityPredictor -> float function] Scoring function

**alleles** [list of string, optional] If not specified, model selection is performed for all alleles.

**min\_models** [int, optional] Min models to select per allele

**max\_models** [int, optional] Max models to select per allele

#### Returns

Class1AffinityPredictor [predictor containing the selected models]

**class** mhcflurry.Class1NeuralNetwork (\*\*hyperparameters)

Bases: object

Low level class I predictor consisting of a single neural network.

Both single allele and pan-allele prediction are supported.

Users will generally use `Class1AffinityPredictor`, which gives a higher-level interface and supports ensembles.

**`network_hyperparameter_defaults`** = `<mhcflurry.hyperparameters.HyperparameterDefaults object>`  
Hyperparameters (and their default values) that affect the neural network architecture.

**`compile_hyperparameter_defaults`** = `<mhcflurry.hyperparameters.HyperparameterDefaults object>`  
Loss and optimizer hyperparameters.

**`fit_hyperparameter_defaults`** = `<mhcflurry.hyperparameters.HyperparameterDefaults object>`  
Hyperparameters for neural network training.

**`early_stopping_hyperparameter_defaults`** = `<mhcflurry.hyperparameters.HyperparameterDefaults object>`  
Hyperparameters for early stopping.

**`miscellaneous_hyperparameter_defaults`** = `<mhcflurry.hyperparameters.HyperparameterDefaults object>`  
Miscellaneous hyperparameters. These parameters are not used by this class but may be interpreted by other code.

**`hyperparameter_defaults`** = `<mhcflurry.hyperparameters.HyperparameterDefaults object>`  
Combined set of all supported hyperparameters and their default values.

**`hyperparameter_renames`** = `{'embedding_init_method': None, 'embedding_input_dim': None}`

**`classmethod apply_hyperparameter_renames`** (*hyperparameters*)  
Handle hyperparameter renames.

#### Parameters

**`hyperparameters`** [dict]

#### Returns

**`dict`** [updated hyperparameters]

**`KERAS_MODELS_CACHE`** = `{}`

Process-wide keras model cache, a map from: architecture JSON string to (Keras model, existing network weights)

**`classmethod clear_model_cache`** ()

Clear the Keras model cache.

**`classmethod borrow_cached_network`** (*network\_json*, *network\_weights*)

Return a keras Model with the specified architecture and weights. As an optimization, when possible this will reuse architectures from a process-wide cache.

The returned object is “borrowed” in the sense that its weights can change later after subsequent calls to this method from other objects.

If you’re using this from a parallel implementation you’ll need to hold a lock while using the returned object.

#### Parameters

**`network_json`** [string of JSON]

**`network_weights`** [list of numpy.array]

#### Returns

**`keras.models.Model`**

**`network`** (*self*, *borrow=False*)

Return the keras model associated with this predictor.

#### Parameters

**borrow** [bool] Whether to return a cached model if possible. See `borrow_cached_network` for details

**Returns**

**keras.models.Model**

**update\_network\_description** (*self*)

Update `self.network_json` and `self.network_weights` properties based on this instances's neural network.

**static keras\_network\_cache\_key** (*network\_json*)

Given a Keras JSON description of a neural network, return a key that uniquely defines this network. Networks that share the same key should have compatible weights matrices and give the same prediction outputs when their weights are the same.

**Parameters**

**network\_json** [string]

**Returns**

**string**

**get\_config** (*self*)

serialize to a dict all attributes except model weights

**Returns**

**dict**

**classmethod from\_config** (*config*, *weights=None*, *weights\_loader=None*)

deserialize from a dict returned by `get_config`).

**Parameters**

**config** [dict]

**weights** [list of array, optional] Network weights to restore

**weights\_loader** [callable, optional] Function to call (no arguments) to load weights when needed

**Returns**

**Class1NeuralNetwork**

**load\_weights** (*self*)

Load weights by evaluating `self.network_weights_loader`, if needed.

After calling this, `self.network_weights_loader` will be `None` and `self.network_weights` will be the weights list, if available.

**get\_weights** (*self*)

Get the network weights

**Returns**

**list of numpy.array giving weights for each layer or None if there is no**

**network**

**peptides\_to\_network\_input** (*self*, *peptides*)

Encode peptides to the fixed-length encoding expected by the neural network (which depends on the architecture).

**Parameters**

**peptides** [EncodableSequences or list of string]

**Returns**

**numpy.array**

**property supported\_peptide\_lengths**

(minimum, maximum) lengths of peptides supported, inclusive.

**Returns**

**(int, int) tuple**

**allele\_encoding\_to\_network\_input** (*self, allele\_encoding*)

Encode alleles to the fixed-length encoding expected by the neural network (which depends on the architecture).

**Parameters**

**allele\_encoding** [AlleleEncoding]

**Returns**

**(numpy.array, numpy.array)**

**Indices and allele representations.**

**static data\_dependent\_weights\_initialization** (*network, x\_dict=None, method='lsuv', verbose=1*)

Data dependent weights initialization.

**Parameters**

**network** [keras.Model]

**x\_dict** [dict of string -> numpy.ndarray] Training data as would be passed keras.Model.fit().

**method** [string] Initialization method. Currently only “lsuv” is supported.

**verbose** [int] Status updates printed to stdout if verbose > 0

**fit\_generator** (*self, generator, validation\_peptide\_encoding, validation\_affinities, validation\_allele\_encoding=None, validation\_inequalities=None, validation\_output\_indices=None, steps\_per\_epoch=10, epochs=1000, min\_epochs=0, patience=10, min\_delta=0.0, verbose=1, progress\_callback=None, progress\_preamble="", progress\_print\_interval=5.0*)

Fit using a generator. Does not support many of the features of fit(), such as random negative peptides.

Fitting proceeds until early stopping is hit, using the peptides, affinities, etc. given by the parameters starting with “**validation\_**”.

This is used for pre-training pan-allele models using data synthesized by the allele-specific models.

**Parameters**

**generator** [generator yielding (alleles, peptides, affinities) tuples] where alleles and peptides are lists of strings, and affinities is list of floats.

**validation\_peptide\_encoding** [EncodableSequences]

**validation\_affinities** [list of float]

**validation\_allele\_encoding** [AlleleEncoding]

**validation\_inequalities** [list of string]

**validation\_output\_indices** [list of int]



**steps\_per\_epoch** [int]  
**epochs** [int]  
**min\_epochs** [int]  
**patience** [int]  
**min\_delta** [float]  
**verbose** [int]  
**progress\_callback** [thunk]  
**progress\_preamble** [string]  
**progress\_print\_interval** [float]

**fit** (*self*, *peptides*, *affinities*, *allele\_encoding=None*, *inequalities=None*, *output\_indices=None*, *sample\_weights=None*, *shuffle\_permutation=None*, *verbose=1*, *progress\_callback=None*, *progress\_preamble=""*, *progress\_print\_interval=5.0*)  
 Fit the neural network.

#### Parameters

**peptides** [EncodableSequences or list of string]  
**affinities** [list of float] nM affinities. Must be same length of as peptides.  
**allele\_encoding** [AlleleEncoding] If not specified, the model will be a single-allele predictor.  
**inequalities** [list of string, each element one of ">", "<", or "=".] Inequalities to use for fitting. Same length as affinities. Each element must be one of ">", "<", or "=" . For example, a ">" will train on  $y_{pred} > y_{true}$  for that element in the training set. Requires using a custom losses that support inequalities (e.g. `mse_with_inequalities`). If None all inequalities are taken to be "=".

**output\_indices** [list of int] For multi-output models only. Same length as affinities. Indicates the index of the output (starting from 0) for each training example.

**sample\_weights** [list of float] If not specified, all samples (including random negatives added during training) will have equal weight. If specified, the random negatives will be assigned weight=1.0.

**shuffle\_permutation** [list of int] Permutation (integer list) of same length as peptides and affinities If None, then a random permutation will be generated.

**verbose** [int] Keras verbosity level

**progress\_callback** [function] No-argument function to call after each epoch.

**progress\_preamble** [string] Optional string of information to include in each progress update

**progress\_print\_interval** [float] How often (in seconds) to print progress update. Set to None to disable.

**predict** (*self*, *peptides*, *allele\_encoding=None*, *batch\_size=4096*, *output\_index=0*)  
 Predict affinities.

If peptides are specified as EncodableSequences, then the predictions will be cached for this predictor as long as the EncodableSequences object remains in memory. The cache is keyed in the object identity of the EncodableSequences, not the sequences themselves. The cache is used only for allele-specific models (i.e. when `allele_encoding` is None).

**Parameters**

**peptides** [EncodableSequences or list of string]

**allele\_encoding** [AlleleEncoding, optional] Only required when this model is a pan-allele model

**batch\_size** [int] batch\_size passed to Keras

**output\_index** [int or None] For multi-output models. Gives the output index to return. If set to None, then all outputs are returned as a samples x outputs matrix.

**Returns**

**numpy.array of nM affinity predictions**

**classmethod merge** (*models*, *merge\_method*='average')

Merge multiple models at the tensorflow (or other backend) level.

Only certain neural network architectures support merging. Others will result in a NotImplementedError.

**Parameters**

**models** [list of Class1NeuralNetwork] instances to merge

**merge\_method** [string, one of “average”, “sum”, or “concatenate”] How to merge the predictions of the different models

**Returns**

**Class1NeuralNetwork** The merged neural network

**make\_network** (*self*, *peptide\_encoding*, *allele\_amino\_acid\_encoding*, *allele\_dense\_layer\_sizes*, *peptide\_dense\_layer\_sizes*, *peptide\_allele\_merge\_method*, *peptide\_allele\_merge\_activation*, *layer\_sizes*, *dense\_layer\_l1\_regularization*, *dense\_layer\_l2\_regularization*, *activation*, *init*, *output\_activation*, *dropout\_probability*, *batch\_normalization*, *locally\_connected\_layers*, *topology*, *num\_outputs*=1, *allele\_representations*=None)

Helper function to make a keras network for class 1 affinity prediction.

**clear\_allele\_representations** (*self*)

Set allele representations to an empty array. Useful before saving to save a smaller version of the model.

**set\_allele\_representations** (*self*, *allele\_representations*, *force\_surgery*=False)

Set the allele representations in use by this model. This means mutating the weights for the allele input embedding layer.

Rationale: instead of passing in the allele sequence for each data point during model training or prediction (which is expensive in terms of memory usage), we pass in an allele index between 0 and n-1 where n is the number of alleles in some universe of possible alleles. This index is used in the model to lookup the corresponding allele sequence. This function sets the lookup table.

See also: AlleleEncoding.allele\_representations()

**Parameters**

**allele\_representations** [numpy.ndarray of shape (a, l, m)]

where **a** is the total number of alleles, **l** is the allele sequence length, **m** is the length of the vectors used to represent amino acids

**class** mhcflurry.**Class1ProcessingPredictor** (*models*, *manifest\_df*=None, *meta-data\_dataframes*=None)

Bases: object

User-facing interface to antigen processing prediction.

Delegates to an ensemble of `Class1ProcessingNeuralNetwork` instances.

Instantiate a new `Class1ProcessingPredictor`

Users will generally call `load()` to restore a saved predictor rather than using this constructor.

#### Parameters

**models** [list of `Class1ProcessingNeuralNetwork`] Neural networks in the ensemble.

**manifest\_df** [`pandas.DataFrame`] Manifest dataframe. If not specified a new one will be created when needed.

**metadata\_dataframes** [dict of string -> `pandas.DataFrame`] Arbitrary metadata associated with this predictor

#### **property sequence\_lengths**

Supported maximum sequence lengths.

Passing a peptide greater than the maximum supported length results in an error.

Passing an N- or C-flank sequence greater than the maximum supported length results in some part of it being ignored.

#### Returns

**dict of string -> int**

Keys are “peptide”, “n\_flank”, “c\_flank”. Values give the maximum supported sequence length.

**add\_models** (*self*, *models*)

Add models to the ensemble (in-place).

#### Parameters

**models** [list of `Class1ProcessingNeuralNetwork`]

#### Returns

**list of string**

Names of the new models.

#### **property manifest\_df**

A `pandas.DataFrame` describing the models included in this predictor.

#### Returns

**pandas.DataFrame**

**static model\_name** (*num*)

Generate a model name

#### Returns

**string**

**static weights\_path** (*models\_dir*, *model\_name*)

Generate the path to the weights file for a model

#### Parameters

**models\_dir** [string]

**model\_name** [string]

#### Returns

**string**

**predict** (*self*, *peptides*, *n\_flanks=None*, *c\_flanks=None*, *batch\_size=4096*)  
Predict antigen processing.

**Parameters**

**peptides** [list of string] Peptide sequences  
**n\_flanks** [list of string] Upstream sequence before each peptide  
**c\_flanks** [list of string] Downstream sequence after each peptide  
**batch\_size** [int] Prediction keras batch size.

**Returns**

**numpy.array**  
**Processing scores. Range is 0-1, higher indicates more favorable processing.**

**predict\_to\_dataframe** (*self*, *peptides*, *n\_flanks=None*, *c\_flanks=None*, *batch\_size=4096*)  
Predict antigen processing.

See [predict](#) method for parameter descriptions.

**Returns**

**pandas.DataFrame**  
**Processing predictions are in the “score” column. Also includes peptides and flanking sequences.**

**predict\_to\_dataframe\_encoded** (*self*, *sequences*, *batch\_size=4096*)  
Predict antigen processing.

See [predict](#) method for more information.

**Parameters**

**sequences** [FlankingEncoding]  
**batch\_size** [int]

**Returns**

**pandas.DataFrame**

**check\_consistency** (*self*)  
Verify that self.manifest\_df is consistent with instance variables.

Currently only checks for agreement on the total number of models.

Throws AssertionError if inconsistent.

**save** (*self*, *models\_dir*, *model\_names\_to\_write=None*, *write\_metadata=True*)  
Serialize the predictor to a directory on disk. If the directory does not exist it will be created.

The serialization format consists of a file called “manifest.csv” with the configurations of each Class1ProcessingNeuralNetwork, along with per-network files giving the model weights.

**Parameters**

**models\_dir** [string] Path to directory. It will be created if it doesn’t exist.

**classmethod load** (*models\_dir=None, max\_models=None*)

Deserialize a predictor from a directory on disk.

#### Parameters

**models\_dir** [string] Path to directory. If unspecified the default downloaded models are used.

**max\_models** [int, optional] Maximum number of models to load

#### Returns

*Class1ProcessingPredictor* instance

**class** mhcflurry.**Class1ProcessingNeuralNetwork** (*\*\*hyperparameters*)

Bases: object

A neural network for antigen processing prediction

**network\_hyperparameter\_defaults** = <mhcflurry.hyperparameters.HyperparameterDefaults object>  
Hyperparameters (and their default values) that affect the neural network architecture.

**fit\_hyperparameter\_defaults** = <mhcflurry.hyperparameters.HyperparameterDefaults object>  
Hyperparameters for neural network training.

**early\_stopping\_hyperparameter\_defaults** = <mhcflurry.hyperparameters.HyperparameterDefaults object>  
Hyperparameters for early stopping.

**compile\_hyperparameter\_defaults** = <mhcflurry.hyperparameters.HyperparameterDefaults object>  
Loss and optimizer hyperparameters. Any values supported by keras may be used.

**auxiliary\_input\_hyperparameter\_defaults** = <mhcflurry.hyperparameters.HyperparameterDefaults object>  
Allele feature hyperparameters.

**hyperparameter\_defaults** = <mhcflurry.hyperparameters.HyperparameterDefaults object>

**property sequence\_lengths**  
Supported maximum sequence lengths

#### Returns

dict of string -> int

Keys are “peptide”, “n\_flank”, “c\_flank”. Values give the maximum supported sequence length.

**network** (*self*)

Return the keras model associated with this network.

**update\_network\_description** (*self*)

Update self.network\_json and self.network\_weights properties based on this instances’s neural network.

**fit** (*self, sequences, targets, sample\_weights=None, shuffle\_permutation=None, verbose=1, progress\_callback=None, progress\_preamble="", progress\_print\_interval=5.0*)  
Fit the neural network.

#### Parameters

**sequences** [FlankingEncoding] Peptides and upstream/downstream flanking sequences

**targets** [list of float] 1 indicates hit, 0 indicates decoy

**sample\_weights** [list of float] If not specified all samples have equal weight.

**shuffle\_permutation** [list of int] Permutation (integer list) of same length as peptides and affinities If None, then a random permutation will be generated.

**verbose** [int] Keras verbosity level

**progress\_callback** [function] No-argument function to call after each epoch.

**progress\_preamble** [string] Optional string of information to include in each progress update

**progress\_print\_interval** [float] How often (in seconds) to print progress update. Set to None to disable.

**predict** (*self*, *peptides*, *n\_flanks=None*, *c\_flanks=None*, *batch\_size=4096*)  
Predict antigen processing.

#### Parameters

**peptides** [list of string] Peptide sequences

**n\_flanks** [list of string] Upstream sequence before each peptide

**c\_flanks** [list of string] Downstream sequence after each peptide

**batch\_size** [int] Prediction keras batch size.

#### Returns

**numpy.array**

Processing scores. Range is 0-1, higher indicates more favorable processing.

**predict\_encoded** (*self*, *sequences*, *batch\_size=4096*)  
Predict antigen processing.

#### Parameters

**sequences** [FlankingEncoding] Peptides and flanking sequences

**batch\_size** [int] Prediction keras batch size.

#### Returns

**numpy.array**

**network\_input** (*self*, *sequences*)  
Encode peptides to the fixed-length encoding expected by the neural network (which depends on the architecture).

#### Parameters

**sequences** [FlankingEncoding] Peptides and flanking sequences

#### Returns

**numpy.array**

**make\_network** (*self*, *amino\_acid\_encoding*, *peptide\_max\_length*, *n\_flank\_length*, *c\_flank\_length*, *flanking\_averages*, *convolutional\_filters*, *convolutional\_kernel\_size*, *convolutional\_activation*, *convolutional\_kernel\_l1\_l2*, *dropout\_rate*, *post\_convolutional\_dense\_layer\_sizes*)  
Helper function to make a keras network given hyperparameters.

**get\_weights** (*self*)  
Get the network weights

#### Returns

list of numpy.array giving weights for each layer or None if there is no

**network**

**get\_config** (*self*)  
serialize to a dict all attributes except model weights

**Returns**

**dict**

**classmethod from\_config** (*config*, *weights=None*)  
deserialize from a dict returned by get\_config().

**Parameters**

**config** [dict]

**weights** [list of array, optional] Network weights to restore

**Returns****Class1ProcessingNeuralNetwork**

```
class mhcflurry.Class1PresentationPredictor (affinity_predictor=None, process-  
                                           ing_predictor_with_flanks=None, pro-  
                                           cessing_predictor_without_flanks=None,  
                                           weights_dataframe=None, meta-  
                                           data_dataframes=None)
```

Bases: object

A logistic regression model over predicted binding affinity (BA) and antigen processing (AP) score.

Instances of this class delegate to Class1AffinityPredictor and Class1ProcessingPredictor instances to generate BA and AP predictions. These predictions are combined using a logistic regression model to give a “presentation score” prediction.

Most users will call the *load* static method to get an instance of this class, then call the *predict* method to generate predictions.

**model\_inputs** = ['affinity\_score', 'processing\_score']

**property supported\_alleles**

List of alleles supported by the underlying Class1AffinityPredictor

**property supported\_peptide\_lengths**

(min, max) of supported peptide lengths, inclusive.

**predict\_affinity** (*self*, *peptides*, *alleles*, *sample\_names=None*, *include\_affinity\_percentile=True*,  
 *verbose=1*, *throw=True*)

Predict binding affinities across samples (each corresponding to up to six MHC I alleles).

Two modes are supported: each peptide can be evaluated for binding to any of the alleles in any sample (this is what happens when *sample\_names* is None), or the *i*'th peptide can be evaluated for binding the alleles of the sample given by the *i*'th entry in *sample\_names*.

For example, if we don't specify *sample\_names*, then predictions are taken for all combinations of samples and peptides, for a result size of num peptides \* num samples:

```
>>> predictor = Class1PresentationPredictor.load()
>>> predictor.predict_affinity(
...     peptides=["SIINFEKL", "PEPTIDE"],
...     alleles={
...         "sample1": ["A0201", "A0301", "B0702"],
...         "sample2": ["A0101", "C0202"],
...     },
```

(continues on next page)

(continued from previous page)

```

...     verbose=0)
    peptide peptide_num sample_name      affinity best_allele
0  SIINFEKL           0    sample1  12906.787792      A0201
1   PEPTIDE           1    sample1  36827.681130      B0702
2  SIINFEKL           0    sample2   3588.413748      C0202
3   PEPTIDE           1    sample2  34362.109211      C0202

```

In contrast, here we specify `sample_names`, so peptide is evaluated for binding the alleles in the corresponding sample, for a result size equal to the number of peptides:

```

>>> predictor.predict_affinity(
...     peptides=["SIINFEKL", "PEPTIDE"],
...     alleles={
...         "sample1": ["A0201", "A0301", "B0702"],
...         "sample2": ["A0101", "C0202"],
...     },
...     sample_names=["sample2", "sample1"],
...     verbose=0)
    peptide peptide_num sample_name      affinity best_allele
0  SIINFEKL           0    sample2   3588.412141      C0202
1   PEPTIDE           1    sample1  36827.682779      B0702

```

### Parameters

**peptides** [list of string] Peptide sequences

**alleles** [dict of string -> list of string] Keys are sample names, values are the alleles (genotype) for that sample

**sample\_names** [list of string [same length as peptides]] Sample names corresponding to each peptide. If None, then predictions are generated for all sample genotypes across all peptides.

**include\_affinity\_percentile** [bool] Whether to include affinity percentile ranks

**verbose** [int] Set to 0 for quiet.

**throw** [verbose] Whether to throw exception (vs. just log a warning) on invalid peptides, etc.

### Returns

**pandas.DataFrame** [predictions]

**predict\_processing** (*self*, *peptides*, *n\_flanks=None*, *c\_flanks=None*, *verbose=1*)

Predict antigen processing scores for individual peptides, optionally including flanking sequences for better cleavage prediction.

### Parameters

**peptides** [list of string]

**n\_flanks** [list of string [same length as peptides]]

**c\_flanks** [list of string [same length as peptides]]

**verbose** [int]

### Returns

**numpy.array** [Antigen processing scores for each peptide]



**fit** (*self*, *targets*, *peptides*, *sample\_names*, *alleles*, *n\_flanks*=None, *c\_flanks*=None, *verbose*=1)

Fit the presentation score logistic regression model.

#### Parameters

**targets** [list of int/float] 1 indicates hit, 0 indicates decoy

**peptides** [list of string [same length as targets]]

**sample\_names** [list of string [same length as targets]]

**alleles** [dict of string -> list of string] Keys are sample names, values are the alleles for that sample

**n\_flanks** [list of string [same length as targets]]

**c\_flanks** [list of string [same length as targets]]

**verbose** [int]

**get\_model** (*self*, *name*=None)

Load or instantiate a new logistic regression model. Private helper method.

#### Parameters

**name** [string] If None (the default), an un-fit LR model is returned. Otherwise the weights are loaded for the specified model.

#### Returns

**sklearn.linear\_model.LogisticRegression**

**predict** (*self*, *peptides*, *alleles*, *sample\_names*=None, *n\_flanks*=None, *c\_flanks*=None, *include\_affinity\_percentile*=False, *verbose*=1, *throw*=True)

Predict presentation scores across a set of peptides.

Presentation scores combine predictions for MHC I binding affinity and antigen processing.

This method returns a pandas.DataFrame giving presentation scores plus the binding affinity and processing predictions and other intermediate results.

Example:

```
>>> predictor = Class1PresentationPredictor.load()
>>> predictor.predict(
...     peptides=["SIINFEKL", "PEPTIDE"],
...     n_flanks=["NNN", "SNS"],
...     c_flanks=["CCC", "CNC"],
...     alleles={
...         "sample1": ["A0201", "A0301", "B0702"],
...         "sample2": ["A0101", "C0202"],
...     },
...     verbose=0)
  peptide n_flank c_flank  peptide_num sample_name  affinity best_
↪allele  processing_score presentation_score
0  SIINFEKL    NNN    CCC           0    sample1  12906.787792
↪A0201          0.802466          0.140365
1  PEPTIDE    SNS    CNC           1    sample1  36827.681130
↪B0702          0.105260          0.004059
2  SIINFEKL    NNN    CCC           0    sample2  3588.413748
↪C0202          0.802466          0.338647
3  PEPTIDE    SNS    CNC           1    sample2  34362.109211
↪C0202          0.105260          0.004317
```

You can also specify `sample_names`, in which case peptide is evaluated for binding the alleles in the corresponding sample only. See [predict\\_affinity](#) for an examples.

### Parameters

**peptides** [list of string] Peptide sequences

**alleles** [list of string or dict of string -> list of string] If you are predicting for a single sample, pass a list of strings (up to 6) indicating the genotype. If you are predicting across multiple samples, pass a dict where the keys are (arbitrary) sample names and the values are the alleles to predict for that sample.

**sample\_names** [list of string [same length as peptides]] If you are passing a dict for 'alleles', you can use this argument to specify which peptides go with which samples. If it is None, then predictions will be performed for each peptide across all samples.

**n\_flanks** [list of string [same length as peptides]] Upstream sequences before the peptide. Sequences of any length can be given and a suffix of the size supported by the model will be used.

**c\_flanks** [list of string [same length as peptides]] Downstream sequences after the peptide. Sequences of any length can be given and a prefix of the size supported by the model will be used.

**include\_affinity\_percentile** [bool] Whether to include affinity percentile ranks

**verbose** [int] Set to 0 for quiet.

**throw** [verbose] Whether to throw exception (vs. just log a warning) on invalid peptides, etc.

### Returns

**pandas.DataFrame**

**Presentation scores and intermediate results.**

**predict\_sequences** (*self*, *sequences*, *alleles*, *result='best'*, *comparison\_quantity='presentation\_score'*, *filter\_value=None*, *peptide\_lengths=(8, 9, 10, 11)*, *use\_flanks=True*, *include\_affinity\_percentile=True*, *verbose=1*, *throw=True*)

Predict presentation across protein sequences.

Example:

```
>>> predictor = Class1PresentationPredictor.load()
>>> predictor.predict_sequences(
...     sequences={
...         'protein1': "MDSKGSSQKGSRLLLLLVVSNNLL",
...         'protein2': "SSLPTPEDKEQAQQTHH",
...     },
...     alleles={
...         "sample1": ["A0201", "A0301", "B0702"],
...         "sample2": ["A0101", "C0202"],
...     },
...     result="filtered",
...     comparison_quantity="affinity",
...     filter_value=500,
...     verbose=0)
  sequence_name pos      peptide      n_flank      c_flank sample_name
↪affinity best_allele  affinity_percentile  processing_score  presentation_
↪score
```

(continues on next page)

(continued from previous page)

0	protein1	13	LLLLVVSNL	MDSKGSSQKGSRL	L	sample1	38.
↪206225	A0201			0.380125	0.017644	0.	
↪571060							
1	protein1	14	LLLLVSNLL	MDSKGSSQKGSRL		sample1	42.
↪243472	A0201			0.420250	0.090984	0.	
↪619213							
2	protein1	5	SSQKGSRL	MDSKG	LLLLVSNLL	sample2	66.
↪749223	C0202			0.803375	0.383608	0.	
↪774468							
3	protein1	6	SQKGSRL	MDSKGS	LLVVSNNL	sample2	178.
↪033474	C0202			1.820000	0.275019	0.	
↪482206							
4	protein1	13	LLLLVVSNNL	MDSKGSSQKGSRL		sample1	202.
↪208167	A0201			1.112500	0.058782	0.	
↪261320							
5	protein1	12	LLLLLVSNL	MDSKGSSQKGSR	L	sample1	202.
↪506582	A0201			1.112500	0.010025	0.	
↪225648							
6	protein2	0	SSLPTPEDK		EQAQQTHH	sample1	335.
↪529377	A0301			1.011750	0.010443	0.	
↪156798							
7	protein2	0	SSLPTPEDK		EQAQQTHH	sample2	353.
↪451759	C0202			2.674250	0.010443	0.	
↪150753							
8	protein1	8	KGSRL	MDSKGSSQ	VVSNNL	sample2	410.
↪327286	C0202			2.887000	0.121374	0.	
↪194081							
9	protein1	5	SSQKGSRL	MDSKG	LLLLVVSNNL	sample2	477.
↪285954	C0202			3.107375	0.111982	0.	
↪168572							

## Parameters

**sequences** [str, list of string, or string -> string dict] Protein sequences. If a dict is given, the keys are arbitrary ( e.g. protein names), and the values are the amino acid sequences.

**alleles** [list of string, list of list of string, or dict of string -> list of string] MHC I alleles. Can be: (1) a string (a single allele), (2) a list of strings (a single genotype), (3) a list of list of strings (multiple genotypes, where the total number of genotypes must equal the number of sequences), or (4) a dict giving multiple genotypes, which will each be run over the sequences.

**result** [string] Specify ‘best’ to return the strongest peptide for each sequence, ‘all’ to return predictions for all peptides, or ‘filtered’ to return predictions where the comparison\_quantity is stronger (i.e (<) for affinity, (>) for scores) than filter\_value.

**comparison\_quantity** [string] One of “presentation\_score”, “processing\_score”, “affinity”, or “affinity\_percentile”. Prediction to use to rank (if result is “best”) or filter (if result is “filtered”) results.

**filter\_value** [float] Threshold value to use, only relevant when result is “filtered”. If comparison\_quantity is “affinity”, then all results less than (i.e. tighter than) the specified nM affinity are retained. If it’s “presentation\_score” or “processing\_score” then results greater than the indicated filter\_value are retained.

**peptide\_lengths** [list of int] Peptide lengths to predict for.

**use\_flanks** [bool] Whether to include flanking sequences when running the AP predictor

(for better cleavage prediction).

**include\_affinity\_percentile** [bool] Whether to include affinity percentile ranks in output.

**verbose** [int] Set to 0 for quiet mode.

**throw** [boolean] Whether to throw exceptions (vs. log warnings) on invalid inputs.

#### Returns

**pandas.DataFrame with columns:** peptide, n\_flank, c\_flank, sequence\_name, affinity, best\_allele, processing\_score, presentation\_score

**save** (*self*, *models\_dir*)

Save the predictor to a directory on disk. If the directory does not exist it will be created.

The wrapped Class1AffinityPredictor and Class1ProcessingPredictor instances are included in the saved data.

#### Parameters

**models\_dir** [string] Path to directory. It will be created if it doesn't exist.

**classmethod load** (*models\_dir=None*, *max\_models=None*)

Deserialize a predictor from a directory on disk.

This will also load the wrapped Class1AffinityPredictor and Class1ProcessingPredictor instances.

#### Parameters

**models\_dir** [string] Path to directory. If unspecified the default downloaded models are used.

**max\_models** [int, optional] Maximum number of affinity and processing (counted separately) models to load

#### Returns

*Class1PresentationPredictor* instance

## 5.1 Submodules

## 5.2 mhcflurry.allele\_encoding module

```
class mhcflurry.allele_encoding.AlleleEncoding (alleles=None, allele_to_sequence=None, row_from=None)
```

Bases: object

A place to cache encodings for a sequence of alleles.

We frequently work with alleles by integer indices, for example as inputs to neural networks. This class is used to map allele names to integer indices in a consistent way by keeping track of the universe of alleles under use, i.e. a distinction is made between the universe of supported alleles (what's in `allele_to_sequence`) and the actual set of alleles used for some task (what's in `alleles`).

#### Parameters

**alleles** [list of string] Allele names. If any allele is None instead of string, it will be mapped to the special index value -1.

**allele\_to\_sequence** [dict of str -> str] Allele name to amino acid sequence

**borrow\_from** [AlleleEncoding, optional] If specified, do not specify allele\_to\_sequence. The sequences from the provided instance are used. This guarantees that the mappings from allele to index and from allele to sequence are the same between the instances.

**compact** (*self*)

Return a new AlleleEncoding in which the universe of supported alleles is only the alleles actually used.

**Returns**

AlleleEncoding

**allele\_representations** (*self*, *encoding\_name*)

Encode the universe of supported allele sequences to a matrix.

**Parameters**

**encoding\_name** [string] How to represent amino acids. Valid names are “BLOSUM62” or “one-hot”. See `amino_acid.ENCODING_DATA_FRAMES`.

**Returns**

**numpy.array of shape** (num alleles in universe, sequence length, vector size)

**where vector size is usually 21 (20 amino acids + X character)**

**fixed\_length\_vector\_encoded\_sequences** (*self*, *encoding\_name*)

Encode allele sequences (not the universe of alleles) to a matrix.

**Parameters**

**encoding\_name** [string] How to represent amino acids. Valid names are “BLOSUM62” or “one-hot”. See `amino_acid.ENCODING_DATA_FRAMES`.

**Returns**

**numpy.array with shape:** (num alleles, sequence length, vector size)

**where vector size is usually 21 (20 amino acids + X character)**

## 5.3 mhcflurry.amino\_acid module

Functions for encoding fixed length sequences of amino acids into various vector representations, such as one-hot and BLOSUM62.

`mhcflurry.amino_acid.available_vector_encodings()`

Return list of supported amino acid vector encodings.

**Returns**

list of string

`mhcflurry.amino_acid.vector_encoding_length(name)`

Return the length of the given vector encoding.

**Parameters**

**name** [string]

**Returns**

int

`mhcflurry.amino_acid.index_encoding(sequences, letter_to_index_dict)`

Encode a sequence of same-length strings to a matrix of integers of the same shape. The map from characters to integers is given by `letter_to_index_dict`.

Given a sequence of `n` strings all of length `k`, return a `k * n` array where the `(i, j)`th element is `letter_to_index_dict[sequence[i][j]]`.

**Parameters**

**sequences** [list of length `n` of strings of length `k`]

**letter\_to\_index\_dict** [dict]

**Returns**

**numpy.array of integers with shape `(k, n)`**

`mhcflurry.amino_acid.fixed_vectors_encoding(index_encoded_sequences, letter_to_vector_df)`

Given a `n x k` matrix of integers such as that returned by `index_encoding()` and a dataframe mapping each index to an arbitrary vector, return a `n * k * m` array where the `(i, j)`th element is `letter_to_vector_df.iloc[sequence[i][j]]`.

The dataframe index and columns names are ignored here; the indexing is done entirely by integer position in the dataframe.

**Parameters**

**index\_encoded\_sequences** [`n x k` array of integers]

**letter\_to\_vector\_df** [pandas.DataFrame of shape `(alphabet_size, m)`]

**Returns**

**numpy.array of integers with shape `(n, k, m)`**

## 5.4 mhcflurry.calibrate\_percentile\_ranks\_command module

Calibrate percentile ranks for models. Runs in-place.

`mhcflurry.calibrate_percentile_ranks_command.run(argv=['-b', 'latex', '-v', '-d', '_build/doctrees', '.', '_build/latex'])`

`mhcflurry.calibrate_percentile_ranks_command.do_calibrate_percentile_ranks(alleles, constant_data={})`

`mhcflurry.calibrate_percentile_ranks_command.calibrate_percentile_ranks(allele, pre-dic-tor, pep-tides=None, mo-tif_summary=False, sum-mary_top_peptide_fraction, ver-bose=False, model_kwargs={})`

## 5.5 mhcflurry.class1\_affinity\_predictor module

```
class mhcflurry.class1_affinity_predictor.Class1AffinityPredictor (allele_to_allele_specific_models=None,  
                                                                class1_pan_allele_models=None,  
                                                                al-  
                                                                lele_to_sequence=None,  
                                                                mani-  
                                                                fest_df=None,  
                                                                al-  
                                                                lele_to_percent_rank_transform=None,  
                                                                meta-  
                                                                data_dataframes=None)
```

Bases: `object`

High-level interface for peptide/MHC I binding affinity prediction.

This class manages low-level `Class1NeuralNetwork` instances, each of which wraps a single Keras network. The purpose of `Class1AffinityPredictor` is to implement ensembles, handling of multiple alleles, and predictor loading and saving. It also provides a place to keep track of metadata like prediction histograms for percentile rank calibration.

### Parameters

**allele\_to\_allele\_specific\_models** [dict of string -> list of `Class1NeuralNetwork`] Ensemble of single-allele models to use for each allele.

**class1\_pan\_allele\_models** [list of `Class1NeuralNetwork`] Ensemble of pan-allele models.

**allele\_to\_sequence** [dict of string -> string] MHC allele name to fixed-length amino acid sequence (sometimes referred to as the pseudosequence). Required only if `class1_pan_allele_models` is specified.

**manifest\_df** [`pandas.DataFrame`, optional] Must have columns: `model_name`, `allele`, `config_json`, `model`. Only required if you want to update an existing serialization of a `Class1AffinityPredictor`. Otherwise this dataframe will be generated automatically based on the supplied models.

**allele\_to\_percent\_rank\_transform** [dict of string -> `PercentRankTransform`, optional] `PercentRankTransform` instances to use for each allele

**metadata\_dataframes** [dict of string -> `pandas.DataFrame`, optional] Optional additional dataframes to write to the models dir when `save()` is called. Useful for tracking provenance.

### property `manifest_df`

A `pandas.DataFrame` describing the models included in this predictor.

Based on: - `self.class1_pan_allele_models` - `self.allele_to_allele_specific_models`

### Returns

**`pandas.DataFrame`**

### `clear_cache` (*self*)

Clear values cached based on the neural networks in this predictor.

Users should call this after mutating any of the following:

- `self.class1_pan_allele_models`
- `self.allele_to_allele_specific_models`

- `self.allele_to_sequence`

Methods that mutate these instance variables will call this method on their own if needed.

**property `neural_networks`**

List of the neural networks in the ensemble.

**Returns**

**list of `Class1NeuralNetwork`**

**classmethod `merge` (*predictors*)**

Merge the ensembles of two or more `Class1AffinityPredictor` instances.

Note: the resulting merged predictor will NOT have calibrated percentile ranks. Call `calibrate_percentile_ranks` on it if these are needed.

**Parameters**

**predictors** [sequence of `Class1AffinityPredictor`]

**Returns**

**`Class1AffinityPredictor` instance**

**`merge_in_place` (*self*, *others*)**

Add the models present in other predictors into the current predictor.

**Parameters**

**others** [list of `Class1AffinityPredictor`] Other predictors to merge into the current predictor.

**Returns**

**list of string** [names of newly added models]

**property `supported_alleles`**

Alleles for which predictions can be made.

**Returns**

**list of string**

**property `supported_peptide_lengths`**

(minimum, maximum) lengths of peptides supported by *all models*, inclusive.

**Returns**

**(int, int) tuple**

**`check_consistency` (*self*)**

Verify that `self.manifest_df` is consistent with: - `self.class1_pan_allele_models` - `self.allele_to_allele_specific_models`

Currently only checks for agreement on the total number of models.

Throws `AssertionError` if inconsistent.

**`save` (*self*, *models\_dir*, *model\_names\_to\_write=None*, *write\_metadata=True*)**

Serialize the predictor to a directory on disk. If the directory does not exist it will be created.

The serialization format consists of a file called “manifest.csv” with the configurations of each `Class1NeuralNetwork`, along with per-network files giving the model weights. If there are pan-allele predictors in the ensemble, the allele sequences are also stored in the directory. There is also a small file “index.txt” with basic metadata: when the models were trained, by whom, on what host.

**Parameters**



**models\_dir** [string] Path to directory. It will be created if it doesn't exist.

**model\_names\_to\_write** [list of string, optional] Only write the weights for the specified models. Useful for incremental updates during training.

**write\_metadata** [boolean, optional] Whether to write optional metadata

**static load** (*models\_dir=None, max\_models=None, optimization\_level=None*)

Deserialize a predictor from a directory on disk.

#### Parameters

**models\_dir** [string] Path to directory. If unspecified the default downloaded models are used.

**max\_models** [int, optional] Maximum number of `Class1NeuralNetwork` instances to load

**optimization\_level** [int] If >0, model optimization will be attempted. Defaults to value of environment variable `MHCFLURRY_OPTIMIZATION_LEVEL`.

#### Returns

*Class1AffinityPredictor* instance

**optimize** (*self, warn=True*)

EXPERIMENTAL: Optimize the predictor for faster predictions.

Currently the only optimization implemented is to merge multiple pan-allele predictors at the tensorflow level.

The optimization is performed in-place, mutating the instance.

#### Returns

**bool** Whether optimization was performed

**static model\_name** (*allele, num*)

Generate a model name

#### Parameters

**allele** [string]

**num** [int]

#### Returns

**string**

**static weights\_path** (*models\_dir, model\_name*)

Generate the path to the weights file for a model

#### Parameters

**models\_dir** [string]

**model\_name** [string]

#### Returns

**string**

**property master\_allele\_encoding**

An `AlleleEncoding` containing the universe of alleles specified by `self.allele_to_sequence`.

#### Returns

### AlleleEncoding

```
fit_allele_specific_predictors (self, n_models, architecture_hyperparameters_list,  
                                allele, peptides, affinities, inequalities=None,  
                                train_rounds=None, models_dir_for_save=None,  
                                verbose=0, progress_preamble="",  
                                progress_print_interval=5.0)
```

Fit one or more allele specific predictors for a single allele using one or more neural network architectures.

The new predictors are saved in the `Class1AffinityPredictor` instance and will be used on subsequent calls to `predict`.

#### Parameters

**n\_models** [int] Number of neural networks to fit

**architecture\_hyperparameters\_list** [list of dict] List of hyperparameter sets.

**allele** [string]

**peptides** [`EncodableSequences` or list of string]

**affinities** [list of float] nM affinities

**inequalities** [list of string, each element one of ">", "<", or "="] See `Class1NeuralNetwork.fit` for details.

**train\_rounds** [sequence of int] Each training point *i* will be used on training rounds *r* for which `train_rounds[i] > r`, `r >= 0`.

**models\_dir\_for\_save** [string, optional] If specified, the `Class1AffinityPredictor` is (incrementally) written to the given models dir after each neural network is fit.

**verbose** [int] Keras verbosity

**progress\_preamble** [string] Optional string of information to include in each progress update

**progress\_print\_interval** [float] How often (in seconds) to print progress. Set to `None` to disable.

#### Returns

list of `Class1NeuralNetwork`

```
fit_class1_pan_allele_models (self, n_models, architecture_hyperparameters, alleles, pep-  
                                tides, affinities, inequalities, models_dir_for_save=None, ver-  
                                bose=1, progress_preamble="", progress_print_interval=5.0)
```

Fit one or more pan-allele predictors using a single neural network architecture.

The new predictors are saved in the `Class1AffinityPredictor` instance and will be used on subsequent calls to `predict`.

#### Parameters

**n\_models** [int] Number of neural networks to fit

**architecture\_hyperparameters** [dict]

**alleles** [list of string] Allele names (not sequences) corresponding to each peptide

**peptides** [`EncodableSequences` or list of string]

**affinities** [list of float] nM affinities

**inequalities** [list of string, each element one of ">", "<", or "="] See `Class1NeuralNetwork.fit` for details.

**models\_dir\_for\_save** [string, optional] If specified, the Class1AffinityPredictor is (incrementally) written to the given models dir after each neural network is fit.

**verbose** [int] Keras verbosity

**progress\_preamble** [string] Optional string of information to include in each progress update

**progress\_print\_interval** [float] How often (in seconds) to print progress. Set to None to disable.

### Returns

**list of Class1NeuralNetwork**

**add\_pan\_allele\_model** (*self*, *model*, *models\_dir\_for\_save=None*)

Add a pan-allele model to the ensemble and optionally do an incremental save.

### Parameters

**model** [Class1NeuralNetwork]

**models\_dir\_for\_save** [string] Directory to save resulting ensemble to

**percentile\_ranks** (*self*, *affinities*, *allele=None*, *alleles=None*, *throw=True*)

Return percentile ranks for the given ic50 affinities and alleles.

The ‘allele’ and ‘alleles’ argument are as in the [predict](#) method. Specify one of these.

### Parameters

**affinities** [sequence of float] nM affinities

**allele** [string]

**alleles** [sequence of string]

**throw** [boolean] If True, a ValueError will be raised in the case of unsupported alleles. If False, a warning will be logged and NaN will be returned for those percentile ranks.

### Returns

**numpy.array of float**

**predict** (*self*, *peptides*, *alleles=None*, *allele=None*, *throw=True*, *centrality\_measure='mean'*, *model\_kwargs={}*)

Predict nM binding affinities.

If multiple predictors are available for an allele, the predictions are the geometric means of the individual model (nM) predictions.

One of ‘allele’ or ‘alleles’ must be specified. If ‘allele’ is specified all predictions will be for the given allele. If ‘alleles’ is specified it must be the same length as ‘peptides’ and give the allele corresponding to each peptide.

### Parameters

**peptides** [EncodableSequences or list of string]

**alleles** [list of string]

**allele** [string]

**throw** [boolean] If True, a ValueError will be raised in the case of unsupported alleles or peptide lengths. If False, a warning will be logged and the predictions for the unsupported alleles or peptides will be NaN.

**centrality\_measure** [string or callable] Measure of central tendency to use to combine predictions in the ensemble. Options include: mean, median, robust\_mean.

**model\_kwargs** [dict] Additional keyword arguments to pass to `Class1NeuralNetwork.predict`

#### Returns

**numpy.array of predictions**

```
predict_to_dataframe (self, peptides, alleles=None, allele=None, throw=True,
                        include_individual_model_predictions=False, include_percentile_ranks=True,
                        include_confidence_intervals=True, centrality_measure='mean', model_kwargs={})
```

Predict nM binding affinities. Gives more detailed output than `predict` method, including 5-95% prediction intervals.

If multiple predictors are available for an allele, the predictions are the geometric means of the individual model predictions.

One of ‘allele’ or ‘alleles’ must be specified. If ‘allele’ is specified all predictions will be for the given allele. If ‘alleles’ is specified it must be the same length as ‘peptides’ and give the allele corresponding to each peptide.

#### Parameters

**peptides** [EncodableSequences or list of string]

**alleles** [list of string]

**allele** [string]

**throw** [boolean] If True, a `ValueError` will be raised in the case of unsupported alleles or peptide lengths. If False, a warning will be logged and the predictions for the unsupported alleles or peptides will be NaN.

**include\_individual\_model\_predictions** [boolean] If True, the predictions of each individual model are included as columns in the result `DataFrame`.

**include\_percentile\_ranks** [boolean, default True] If True, a “prediction\_percentile” column will be included giving the percentile ranks. If no percentile rank info is available, this will be ignored with a warning.

**centrality\_measure** [string or callable] Measure of central tendency to use to combine predictions in the ensemble. Options include: mean, median, robust\_mean.

**model\_kwargs** [dict] Additional keyword arguments to pass to `Class1NeuralNetwork.predict`

#### Returns

**pandas.DataFrame of predictions**

```
calibrate_percentile_ranks (self, peptides=None, num_peptides_per_length=100000,
                             alleles=None, bins=None, motif_summary=False, summary_top_peptide_fractions=[0.001],
                             verbose=False, model_kwargs={})
```

Compute the cumulative distribution of `ic50` values for a set of alleles over a large universe of random peptides, to enable taking quantiles of this distribution later.

#### Parameters

**peptides** [sequence of string or EncodableSequences, optional] Peptides to use

**num\_peptides\_per\_length** [int, optional] If peptides argument is not specified, then num\_peptides\_per\_length peptides are randomly sampled from a uniform distribution for each supported length

**alleles** [sequence of string, optional] Alleles to perform calibration for. If not specified all supported alleles will be calibrated.

**bins** [object] Anything that can be passed to numpy.histogram’s “bins” argument can be used here, i.e. either an integer or a sequence giving bin edges. This is in ic50 space.

**motif\_summary** [bool] If True, the length distribution and per-position amino acid frequencies are also calculated for the top x fraction of tightest-binding peptides, where each value of x is given in the summary\_top\_peptide\_fractions list.

**summary\_top\_peptide\_fractions** [list of float] Only used if motif\_summary is True

**verbose** [boolean] Whether to print status updates to stdout

**model\_kwargs** [dict] Additional low-level Class1NeuralNetwork.predict() kwargs.

#### Returns

**dict of string -> pandas.DataFrame**

If motif\_summary is True, this will have keys “frequency\_matrices” and “length\_distributions”. Otherwise it will be empty.

**model\_select** (*self*, *score\_function*, *alleles=None*, *min\_models=1*, *max\_models=10000*)

Perform model selection using a user-specified scoring function.

This works only with allele-specific models, not pan-allele models.

Model selection is done using a “step up” variable selection procedure, in which models are repeatedly added to an ensemble until the score stops improving.

#### Parameters

**score\_function** [Class1AffinityPredictor -> float function] Scoring function

**alleles** [list of string, optional] If not specified, model selection is performed for all alleles.

**min\_models** [int, optional] Min models to select per allele

**max\_models** [int, optional] Max models to select per allele

#### Returns

**Class1AffinityPredictor** [predictor containing the selected models]

## 5.6 mhcflurry.class1\_neural\_network module

**class** mhcflurry.class1\_neural\_network.**Class1NeuralNetwork** (\*\*hyperparameters)

Bases: object

Low level class I predictor consisting of a single neural network.

Both single allele and pan-allele prediction are supported.

Users will generally use Class1AffinityPredictor, which gives a higher-level interface and supports ensembles.

**network\_hyperparameter\_defaults** = <mhcflurry.hyperparameters.HyperparameterDefaults object>  
Hyperparameters (and their default values) that affect the neural network architecture.

**compile\_hyperparameter\_defaults** = <mhcflurry.hyperparameters.HyperparameterDefaults object>  
Loss and optimizer hyperparameters.

**fit\_hyperparameter\_defaults** = <mhcflurry.hyperparameters.HyperparameterDefaults object>  
Hyperparameters for neural network training.

**early\_stopping\_hyperparameter\_defaults** = <mhcflurry.hyperparameters.HyperparameterDefaults object>  
Hyperparameters for early stopping.

**miscellaneous\_hyperparameter\_defaults** = <mhcflurry.hyperparameters.HyperparameterDefaults object>  
Miscellaneous hyperparameters. These parameters are not used by this class but may be interpreted by other code.

**hyperparameter\_defaults** = <mhcflurry.hyperparameters.HyperparameterDefaults object>  
Combined set of all supported hyperparameters and their default values.

**hyperparameter\_renames** = {'embedding\_init\_method': None, 'embedding\_input\_dim': None}

**classmethod apply\_hyperparameter\_renames** (*hyperparameters*)  
Handle hyperparameter renames.

#### Parameters

**hyperparameters** [dict]

#### Returns

**dict** [updated hyperparameters]

**KERAS\_MODELS\_CACHE** = {}

Process-wide keras model cache, a map from: architecture JSON string to (Keras model, existing network weights)

**classmethod clear\_model\_cache** ()  
Clear the Keras model cache.

**classmethod borrow\_cached\_network** (*network\_json*, *network\_weights*)  
Return a keras Model with the specified architecture and weights. As an optimization, when possible this will reuse architectures from a process-wide cache.

The returned object is “borrowed” in the sense that its weights can change later after subsequent calls to this method from other objects.

If you’re using this from a parallel implementation you’ll need to hold a lock while using the returned object.

#### Parameters

**network\_json** [string of JSON]

**network\_weights** [list of numpy.array]

#### Returns

**keras.models.Model**

**network** (*self*, *borrow=False*)  
Return the keras model associated with this predictor.

#### Parameters

**borrow** [bool] Whether to return a cached model if possible. See `borrow_cached_network` for details

#### Returns

**keras.models.Model****update\_network\_description** (*self*)

Update self.network\_json and self.network\_weights properties based on this instances's neural network.

**static keras\_network\_cache\_key** (*network\_json*)

Given a Keras JSON description of a neural network, return a key that uniquely defines this network. Networks that share the same key should have compatible weights matrices and give the same prediction outputs when their weights are the same.

**Parameters**

**network\_json** [string]

**Returns**

string

**get\_config** (*self*)

serialize to a dict all attributes except model weights

**Returns**

dict

**classmethod from\_config** (*config*, *weights=None*, *weights\_loader=None*)

deserialize from a dict returned by get\_config().

**Parameters**

**config** [dict]

**weights** [list of array, optional] Network weights to restore

**weights\_loader** [callable, optional] Function to call (no arguments) to load weights when needed

**Returns**

**Class1NeuralNetwork**

**load\_weights** (*self*)

Load weights by evaluating self.network\_weights\_loader, if needed.

After calling this, self.network\_weights\_loader will be None and self.network\_weights will be the weights list, if available.

**get\_weights** (*self*)

Get the network weights

**Returns**

list of numpy.array giving weights for each layer or None if there is no network

**peptides\_to\_network\_input** (*self*, *peptides*)

Encode peptides to the fixed-length encoding expected by the neural network (which depends on the architecture).

**Parameters**

**peptides** [EncodableSequences or list of string]

**Returns**

numpy.array

**property supported\_peptide\_lengths**

(minimum, maximum) lengths of peptides supported, inclusive.

**Returns**

(int, int) tuple

**allele\_encoding\_to\_network\_input** (*self, allele\_encoding*)

Encode alleles to the fixed-length encoding expected by the neural network (which depends on the architecture).

**Parameters**

**allele\_encoding** [AlleleEncoding]

**Returns**

(numpy.array, numpy.array)

**Indices and allele representations.**

**static data\_dependent\_weights\_initialization** (*network, x\_dict=None, method='lsuv', verbose=1*)

Data dependent weights initialization.

**Parameters**

**network** [keras.Model]

**x\_dict** [dict of string -> numpy.ndarray] Training data as would be passed keras.Model.fit().

**method** [string] Initialization method. Currently only “lsuv” is supported.

**verbose** [int] Status updates printed to stdout if verbose > 0

**fit\_generator** (*self, generator, validation\_peptide\_encoding, validation\_affinities, validation\_allele\_encoding=None, validation\_inequalities=None, validation\_output\_indices=None, steps\_per\_epoch=10, epochs=1000, min\_epochs=0, patience=10, min\_delta=0.0, verbose=1, progress\_callback=None, progress\_preamble="", progress\_print\_interval=5.0*)

Fit using a generator. Does not support many of the features of fit(), such as random negative peptides.

Fitting proceeds until early stopping is hit, using the peptides, affinities, etc. given by the parameters starting with “**validation\_**”.

This is used for pre-training pan-allele models using data synthesized by the allele-specific models.

**Parameters**

**generator** [generator yielding (alleles, peptides, affinities) tuples] where alleles and peptides are lists of strings, and affinities is list of floats.

**validation\_peptide\_encoding** [EncodableSequences]

**validation\_affinities** [list of float]

**validation\_allele\_encoding** [AlleleEncoding]

**validation\_inequalities** [list of string]

**validation\_output\_indices** [list of int]

**steps\_per\_epoch** [int]

**epochs** [int]

**min\_epochs** [int]



**patience** [int]  
**min\_delta** [float]  
**verbose** [int]  
**progress\_callback** [thunk]  
**progress\_preamble** [string]  
**progress\_print\_interval** [float]  
**fit** (*self*, *peptides*, *affinities*, *allele\_encoding=None*, *inequalities=None*, *output\_indices=None*, *sample\_weights=None*, *shuffle\_permutation=None*, *verbose=1*, *progress\_callback=None*, *progress\_preamble=""*, *progress\_print\_interval=5.0*)  
 Fit the neural network.

#### Parameters

**peptides** [EncodableSequences or list of string]  
**affinities** [list of float] nM affinities. Must be same length of as peptides.  
**allele\_encoding** [AlleleEncoding] If not specified, the model will be a single-allele predictor.  
**inequalities** [list of string, each element one of ">", "<", or "=".] Inequalities to use for fitting. Same length as affinities. Each element must be one of ">", "<", or "=".. For example, a ">" will train on  $y_{pred} > y_{true}$  for that element in the training set. Requires using a custom losses that support inequalities (e.g. `mse_with_inequalities`). If None all inequalities are taken to be "=".  
**output\_indices** [list of int] For multi-output models only. Same length as affinities. Indicates the index of the output (starting from 0) for each training example.  
**sample\_weights** [list of float] If not specified, all samples (including random negatives added during training) will have equal weight. If specified, the random negatives will be assigned weight=1.0.  
**shuffle\_permutation** [list of int] Permutation (integer list) of same length as peptides and affinities If None, then a random permutation will be generated.  
**verbose** [int] Keras verbosity level  
**progress\_callback** [function] No-argument function to call after each epoch.  
**progress\_preamble** [string] Optional string of information to include in each progress update  
**progress\_print\_interval** [float] How often (in seconds) to print progress update. Set to None to disable.  
**predict** (*self*, *peptides*, *allele\_encoding=None*, *batch\_size=4096*, *output\_index=0*)  
 Predict affinities.

If peptides are specified as EncodableSequences, then the predictions will be cached for this predictor as long as the EncodableSequences object remains in memory. The cache is keyed in the object identity of the EncodableSequences, not the sequences themselves. The cache is used only for allele-specific models (i.e. when `allele_encoding` is None).

#### Parameters

**peptides** [EncodableSequences or list of string]  
**allele\_encoding** [AlleleEncoding, optional] Only required when this model is a pan-allele model

**batch\_size** [int] batch\_size passed to Keras

**output\_index** [int or None] For multi-output models. Gives the output index to return. If set to None, then all outputs are returned as a samples x outputs matrix.

#### Returns

**numpy.array of nM affinity predictions**

**classmethod merge** (*models*, *merge\_method*='average')

Merge multiple models at the tensorflow (or other backend) level.

Only certain neural network architectures support merging. Others will result in a `NotImplementedError`.

#### Parameters

**models** [list of `Class1NeuralNetwork`] instances to merge

**merge\_method** [string, one of “average”, “sum”, or “concatenate”] How to merge the predictions of the different models

#### Returns

**Class1NeuralNetwork** The merged neural network

**make\_network** (*self*, *peptide\_encoding*, *allele\_amino\_acid\_encoding*, *allele\_dense\_layer\_sizes*, *peptide\_dense\_layer\_sizes*, *peptide\_allele\_merge\_method*, *peptide\_allele\_merge\_activation*, *layer\_sizes*, *dense\_layer\_l1\_regularization*, *dense\_layer\_l2\_regularization*, *activation*, *init*, *output\_activation*, *dropout\_probability*, *batch\_normalization*, *locally\_connected\_layers*, *topology*, *num\_outputs*=1, *allele\_representations*=None)

Helper function to make a keras network for class 1 affinity prediction.

**clear\_allele\_representations** (*self*)

Set allele representations to an empty array. Useful before saving to save a smaller version of the model.

**set\_allele\_representations** (*self*, *allele\_representations*, *force\_surgery*=False)

Set the allele representations in use by this model. This means mutating the weights for the allele input embedding layer.

Rationale: instead of passing in the allele sequence for each data point during model training or prediction (which is expensive in terms of memory usage), we pass in an allele index between 0 and n-1 where n is the number of alleles in some universe of possible alleles. This index is used in the model to lookup the corresponding allele sequence. This function sets the lookup table.

See also: `AlleleEncoding.allele_representations()`

#### Parameters

**allele\_representations** [numpy.ndarray of shape (a, l, m)]

where **a** is the total number of alleles, **l** is the allele sequence length, **m** is the length of the vectors used to represent amino acids

## 5.7 mhcflurry.class1\_presentation\_predictor module

**class** mhcflurry.class1\_presentation\_predictor.**Class1PresentationPredictor** (*affinity\_predictor=None*, *processing\_predictor=None*, *ing\_predictor\_with\_flan=None*, *pro-cess-ing\_predictor\_without\_f=None*, *weights\_dataframe=None*, *meta-data\_dataframes=None*)

Bases: object

A logistic regression model over predicted binding affinity (BA) and antigen processing (AP) score.

Instances of this class delegate to Class1AffinityPredictor and Class1ProcessingPredictor instances to generate BA and AP predictions. These predictions are combined using a logistic regression model to give a “presentation score” prediction.

Most users will call the `load` static method to get an instance of this class, then call the `predict` method to generate predictions.

**model\_inputs** = ['affinity\_score', 'processing\_score']

**property supported\_alleles**

List of alleles supported by the underlying Class1AffinityPredictor

**property supported\_peptide\_lengths**

(min, max) of supported peptide lengths, inclusive.

**predict\_affinity** (*self*, *peptides*, *alleles*, *sample\_names=None*, *include\_affinity\_percentile=True*, *verbose=1*, *throw=True*)

Predict binding affinities across samples (each corresponding to up to six MHC I alleles).

Two modes are supported: each peptide can be evaluated for binding to any of the alleles in any sample (this is what happens when *sample\_names* is *None*), or the *i*’th peptide can be evaluated for binding the alleles of the sample given by the *i*’th entry in *sample\_names*.

For example, if we don’t specify *sample\_names*, then predictions are taken for all combinations of samples and peptides, for a result size of *num peptides* \* *num samples*:

```
>>> predictor = Class1PresentationPredictor.load()
>>> predictor.predict_affinity(
...     peptides=["SIINFPEKL", "PEPTIDE"],
...     alleles={
...         "sample1": ["A0201", "A0301", "B0702"],
...         "sample2": ["A0101", "C0202"],
...     },
...     verbose=0)
  peptide  peptide_num  sample_name  affinity  best_allele
0  SIINFPEKL           0      sample1  12906.787792    A0201
1    PEPTIDE           1      sample1  36827.681130    B0702
2  SIINFPEKL           0      sample2   3588.413748    C0202
3    PEPTIDE           1      sample2  34362.109211    C0202
```

In contrast, here we specify *sample\_names*, so peptide is evaluated for binding the alleles in the corresponding sample, for a result size equal to the number of peptides:

```
>>> predictor.predict_affinity(
...     peptides=["SIINFEKL", "PEPTIDE"],
...     alleles={
...         "sample1": ["A0201", "A0301", "B0702"],
...         "sample2": ["A0101", "C0202"],
...     },
...     sample_names=["sample2", "sample1"],
...     verbose=0)
  peptide  peptide_num sample_name      affinity best_allele
0 SIINFEKL           0    sample2  3588.412141      C0202
1  PEPTIDE           1    sample1  36827.682779      B0702
```

### Parameters

**peptides** [list of string] Peptide sequences

**alleles** [dict of string -> list of string] Keys are sample names, values are the alleles (genotype) for that sample

**sample\_names** [list of string [same length as peptides]] Sample names corresponding to each peptide. If None, then predictions are generated for all sample genotypes across all peptides.

**include\_affinity\_percentile** [bool] Whether to include affinity percentile ranks

**verbose** [int] Set to 0 for quiet.

**throw** [verbose] Whether to throw exception (vs. just log a warning) on invalid peptides, etc.

### Returns

**pandas.DataFrame** [predictions]

**predict\_processing** (*self*, *peptides*, *n\_flanks=None*, *c\_flanks=None*, *verbose=1*)

Predict antigen processing scores for individual peptides, optionally including flanking sequences for better cleavage prediction.

### Parameters

**peptides** [list of string]

**n\_flanks** [list of string [same length as peptides]]

**c\_flanks** [list of string [same length as peptides]]

**verbose** [int]

### Returns

**numpy.array** [Antigen processing scores for each peptide]

**fit** (*self*, *targets*, *peptides*, *sample\_names*, *alleles*, *n\_flanks=None*, *c\_flanks=None*, *verbose=1*)

Fit the presentation score logistic regression model.

### Parameters

**targets** [list of int/float] 1 indicates hit, 0 indicates decoy

**peptides** [list of string [same length as targets]]

**sample\_names** [list of string [same length as targets]]

**alleles** [dict of string -> list of string] Keys are sample names, values are the alleles for that sample

**n\_flanks** [list of string [same length as targets]]

**c\_flanks** [list of string [same length as targets]]

**verbose** [int]

**get\_model** (*self*, *name=None*)

Load or instantiate a new logistic regression model. Private helper method.

#### Parameters

**name** [string] If None (the default), an un-fit LR model is returned. Otherwise the weights are loaded for the specified model.

#### Returns

**sklearn.linear\_model.LogisticRegression**

**predict** (*self*, *peptides*, *alleles*, *sample\_names=None*, *n\_flanks=None*, *c\_flanks=None*, *include\_affinity\_percentile=False*, *verbose=1*, *throw=True*)

Predict presentation scores across a set of peptides.

Presentation scores combine predictions for MHC I binding affinity and antigen processing.

This method returns a pandas.DataFrame giving presentation scores plus the binding affinity and processing predictions and other intermediate results.

Example:

```
>>> predictor = Class1PresentationPredictor.load()
>>> predictor.predict(
...     peptides=["SIINFEKL", "PEPTIDE"],
...     n_flanks=["NNN", "SNS"],
...     c_flanks=["CCC", "CNC"],
...     alleles={
...         "sample1": ["A0201", "A0301", "B0702"],
...         "sample2": ["A0101", "C0202"],
...     },
...     verbose=0)
  peptide n_flank c_flank  peptide_num sample_name  affinity best_
↪ allele processing_score presentation_score
0  SIINFEKL    NNN    CCC           0    sample1  12906.787792
↪ A0201         0.802466         0.140365
1  PEPTIDE    SNS    CNC           1    sample1  36827.681130
↪ B0702         0.105260         0.004059
2  SIINFEKL    NNN    CCC           0    sample2  3588.413748
↪ C0202         0.802466         0.338647
3  PEPTIDE    SNS    CNC           1    sample2  34362.109211
↪ C0202         0.105260         0.004317
```

You can also specify *sample\_names*, in which case peptide is evaluated for binding the alleles in the corresponding sample only. See [predict\\_affinity](#) for an examples.

#### Parameters

**peptides** [list of string] Peptide sequences

**alleles** [list of string or dict of string -> list of string] If you are predicting for a single sample, pass a list of strings (up to 6) indicating the genotype. If you are predicting across multiple

samples, pass a dict where the keys are (arbitrary) sample names and the values are the alleles to predict for that sample.

**sample\_names** [list of string [same length as peptides]] If you are passing a dict for 'alleles', you can use this argument to specify which peptides go with which samples. If it is None, then predictions will be performed for each peptide across all samples.

**n\_flanks** [list of string [same length as peptides]] Upstream sequences before the peptide. Sequences of any length can be given and a suffix of the size supported by the model will be used.

**c\_flanks** [list of string [same length as peptides]] Downstream sequences after the peptide. Sequences of any length can be given and a prefix of the size supported by the model will be used.

**include\_affinity\_percentile** [bool] Whether to include affinity percentile ranks

**verbose** [int] Set to 0 for quiet.

**throw** [verbose] Whether to throw exception (vs. just log a warning) on invalid peptides, etc.

## Returns

**pandas.DataFrame**

**Presentation scores and intermediate results.**

**predict\_sequences** (*self*, *sequences*, *alleles*, *result='best'*, *comparison\_quantity='presentation\_score'*, *filter\_value=None*, *peptide\_lengths=(8, 9, 10, 11)*, *use\_flanks=True*, *include\_affinity\_percentile=True*, *verbose=1*, *throw=True*)

Predict presentation across protein sequences.

Example:

```
>>> predictor = ClassIPresentationPredictor.load()
>>> predictor.predict_sequences(
...     sequences={
...         'protein1': "MDSKGSSQKGSRLLLLLLVSNLL",
...         'protein2': "SSLPTPEDKEQAQQTHH",
...     },
...     alleles={
...         "sample1": ["A0201", "A0301", "B0702"],
...         "sample2": ["A0101", "C0202"],
...     },
...     result="filtered",
...     comparison_quantity="affinity",
...     filter_value=500,
...     verbose=0)
sequence_name pos      peptide      n_flank      c_flank sample_name
↪affinity best_allele  affinity_percentile  processing_score  presentation_
↪score
0      protein1    13      LLLLVVSNL      MDSKGSSQKGSRL      L      sample1    38.
↪206225      A0201      0.380125      0.017644      0.
↪571060
1      protein1    14      LLLVVSNNL      MDSKGSSQKGSRL      L      sample1    42.
↪243472      A0201      0.420250      0.090984      0.
↪619213
2      protein1     5      SSQKGSRL      MDSKG      LLLVVSNNL      sample2    66.
↪749223      C0202      0.803375      0.383608      0.
↪774468
```

(continues on next page)

(continued from previous page)

3	protein1	6	SQKGSRLLL	MDSKGS	LLVVSNNL	sample2	178.
↪033474		C0202		1.820000	0.275019		0.
↪482206							
4	protein1	13	LLLLVVSNNL	MDSKGSSQKGSRL		sample1	202.
↪208167		A0201		1.112500	0.058782		0.
↪261320							
5	protein1	12	LLLLLVVSNL	MDSKGSSQKGSR	L	sample1	202.
↪506582		A0201		1.112500	0.010025		0.
↪225648							
6	protein2	0	SSLPTPEDK		EQAQQTHH	sample1	335.
↪529377		A0301		1.011750	0.010443		0.
↪156798							
7	protein2	0	SSLPTPEDK		EQAQQTHH	sample2	353.
↪451759		C0202		2.674250	0.010443		0.
↪150753							
8	protein1	8	KGSRLLLLL	MDSKGSSQ	VVSNNL	sample2	410.
↪327286		C0202		2.887000	0.121374		0.
↪194081							
9	protein1	5	SSQKGSRL	MDSKG	LLLLVVSNNL	sample2	477.
↪285954		C0202		3.107375	0.111982		0.
↪168572							

## Parameters

**sequences** [str, list of string, or string -> string dict] Protein sequences. If a dict is given, the keys are arbitrary ( e.g. protein names), and the values are the amino acid sequences.

**alleles** [list of string, list of list of string, or dict of string -> list of string] MHC I alleles. Can be: (1) a string (a single allele), (2) a list of strings (a single genotype), (3) a list of list of strings (multiple genotypes, where the total number of genotypes must equal the number of sequences), or (4) a dict giving multiple genotypes, which will each be run over the sequences.

**result** [string] Specify ‘best’ to return the strongest peptide for each sequence, ‘all’ to return predictions for all peptides, or ‘filtered’ to return predictions where the comparison\_quantity is stronger (i.e (<) for affinity, (>) for scores) than filter\_value.

**comparison\_quantity** [string] One of “presentation\_score”, “processing\_score”, “affinity”, or “affinity\_percentile”. Prediction to use to rank (if result is “best”) or filter (if result is “filtered”) results.

**filter\_value** [float] Threshold value to use, only relevant when result is “filtered”. If comparison\_quantity is “affinity”, then all results less than (i.e. tighter than) the specified nM affinity are retained. If it’s “presentation\_score” or “processing\_score” then results greater than the indicated filter\_value are retained.

**peptide\_lengths** [list of int] Peptide lengths to predict for.

**use\_flanks** [bool] Whether to include flanking sequences when running the AP predictor (for better cleavage prediction).

**include\_affinity\_percentile** [bool] Whether to include affinity percentile ranks in output.

**verbose** [int] Set to 0 for quiet mode.

**throw** [boolean] Whether to throw exceptions (vs. log warnings) on invalid inputs.

## Returns

**pandas.DataFrame with columns:** peptide, n\_flank, c\_flank, sequence\_name, affinity, best\_allele, processing\_score, presentation\_score

**save** (*self*, *models\_dir*)

Save the predictor to a directory on disk. If the directory does not exist it will be created.

The wrapped Class1AffinityPredictor and Class1ProcessingPredictor instances are included in the saved data.

#### Parameters

**models\_dir** [string] Path to directory. It will be created if it doesn't exist.

**classmethod load** (*models\_dir=None*, *max\_models=None*)

Deserialize a predictor from a directory on disk.

This will also load the wrapped Class1AffinityPredictor and Class1ProcessingPredictor instances.

#### Parameters

**models\_dir** [string] Path to directory. If unspecified the default downloaded models are used.

**max\_models** [int, optional] Maximum number of affinity and processing (counted separately) models to load

#### Returns

*Class1PresentationPredictor* instance

## 5.8 mhcflurry.class1\_processing\_neural\_network module

Antigen processing neural network implementation

**class** mhcflurry.class1\_processing\_neural\_network.Class1ProcessingNeuralNetwork (\*\*hyperparameter

Bases: object

A neural network for antigen processing prediction

**network\_hyperparameter\_defaults** = <mhcflurry.hyperparameters.HyperparameterDefaults ob  
Hyperparameters (and their default values) that affect the neural network architecture.

**fit\_hyperparameter\_defaults** = <mhcflurry.hyperparameters.HyperparameterDefaults object  
Hyperparameters for neural network training.

**early\_stopping\_hyperparameter\_defaults** = <mhcflurry.hyperparameters.HyperparameterDefa  
Hyperparameters for early stopping.

**compile\_hyperparameter\_defaults** = <mhcflurry.hyperparameters.HyperparameterDefaults ob  
Loss and optimizer hyperparameters. Any values supported by keras may be used.

**auxiliary\_input\_hyperparameter\_defaults** = <mhcflurry.hyperparameters.HyperparameterDef  
Allele feature hyperparameters.

**hyperparameter\_defaults** = <mhcflurry.hyperparameters.HyperparameterDefaults object>

**property sequence\_lengths**  
Supported maximum sequence lengths

#### Returns

dict of string -> int



Keys are “peptide”, “n\_flank”, “c\_flank”. Values give the maximum supported sequence length.

**network** (*self*)

Return the keras model associated with this network.

**update\_network\_description** (*self*)

Update self.network\_json and self.network\_weights properties based on this instances’s neural network.

**fit** (*self*, *sequences*, *targets*, *sample\_weights=None*, *shuffle\_permutation=None*, *verbose=1*, *progress\_callback=None*, *progress\_preamble=""*, *progress\_print\_interval=5.0*)  
Fit the neural network.

#### Parameters

**sequences** [FlankingEncoding] Peptides and upstream/downstream flanking sequences

**targets** [list of float] 1 indicates hit, 0 indicates decoy

**sample\_weights** [list of float] If not specified all samples have equal weight.

**shuffle\_permutation** [list of int] Permutation (integer list) of same length as peptides and affinities If None, then a random permutation will be generated.

**verbose** [int] Keras verbosity level

**progress\_callback** [function] No-argument function to call after each epoch.

**progress\_preamble** [string] Optional string of information to include in each progress update

**progress\_print\_interval** [float] How often (in seconds) to print progress update. Set to None to disable.

**predict** (*self*, *peptides*, *n\_flanks=None*, *c\_flanks=None*, *batch\_size=4096*)

Predict antigen processing.

#### Parameters

**peptides** [list of string] Peptide sequences

**n\_flanks** [list of string] Upstream sequence before each peptide

**c\_flanks** [list of string] Downstream sequence after each peptide

**batch\_size** [int] Prediction keras batch size.

#### Returns

**numpy.array**

Processing scores. Range is 0-1, higher indicates more favorable processing.

**predict\_encoded** (*self*, *sequences*, *batch\_size=4096*)

Predict antigen processing.

#### Parameters

**sequences** [FlankingEncoding] Peptides and flanking sequences

**batch\_size** [int] Prediction keras batch size.

#### Returns

**numpy.array**

**network\_input** (*self*, *sequences*)

Encode peptides to the fixed-length encoding expected by the neural network (which depends on the architecture).

**Parameters**

**sequences** [FlankingEncoding] Peptides and flanking sequences

**Returns**

**numpy.array**

**make\_network** (*self*, *amino\_acid\_encoding*, *peptide\_max\_length*, *n\_flank\_length*, *c\_flank\_length*, *flanking\_averages*, *convolutional\_filters*, *convolutional\_kernel\_size*, *convolutional\_activation*, *convolutional\_kernel\_l1\_l2*, *dropout\_rate*, *post\_convolutional\_dense\_layer\_sizes*)

Helper function to make a keras network given hyperparameters.

**get\_weights** (*self*)

Get the network weights

**Returns**

**list of numpy.array giving weights for each layer or None if there is no network**

**get\_config** (*self*)

serialize to a dict all attributes except model weights

**Returns**

**dict**

**classmethod from\_config** (*config*, *weights=None*)

deserialize from a dict returned by get\_config().

**Parameters**

**config** [dict]

**weights** [list of array, optional] Network weights to restore

**Returns**

**Class1ProcessingNeuralNetwork**

## 5.9 mhcfurry.class1\_processing\_predictor module

```
class mhcfurry.class1_processing_predictor.Class1ProcessingPredictor (models,  
                                                                mani-  
                                                                fest_df=None,  
                                                                meta-  
                                                                data_dataframes=None)
```

Bases: object

User-facing interface to antigen processing prediction.

Delegates to an ensemble of Class1ProcessingNeuralNetwork instances.

Instantiate a new Class1ProcessingPredictor

Users will generally call load() to restore a saved predictor rather than using this constructor.

**Parameters**

**models** [list of Class1ProcessingNeuralNetwork] Neural networks in the ensemble.

**manifest\_df** [pandas.DataFrame] Manifest dataframe. If not specified a new one will be created when needed.

**metadata\_dataframes** [dict of string -> pandas.DataFrame] Arbitrary metadata associated with this predictor

**property sequence\_lengths**

Supported maximum sequence lengths.

Passing a peptide greater than the maximum supported length results in an error.

Passing an N- or C-flank sequence greater than the maximum supported length results in some part of it being ignored.

**Returns**

**dict of string -> int**

**Keys are “peptide”, “n\_flank”, “c\_flank”. Values give the maximum supported sequence length.**

**add\_models** (*self*, *models*)

Add models to the ensemble (in-place).

**Parameters**

**models** [list of Class1ProcessingNeuralNetwork]

**Returns**

**list of string**

**Names of the new models.**

**property manifest\_df**

A pandas.DataFrame describing the models included in this predictor.

**Returns**

**pandas.DataFrame**

**static model\_name** (*num*)

Generate a model name

**Returns**

**string**

**static weights\_path** (*models\_dir*, *model\_name*)

Generate the path to the weights file for a model

**Parameters**

**models\_dir** [string]

**model\_name** [string]

**Returns**

**string**

**predict** (*self*, *peptides*, *n\_flanks=None*, *c\_flanks=None*, *batch\_size=4096*)

Predict antigen processing.

**Parameters**

**peptides** [list of string] Peptide sequences  
**n\_flanks** [list of string] Upstream sequence before each peptide  
**c\_flanks** [list of string] Downstream sequence after each peptide  
**batch\_size** [int] Prediction keras batch size.

**Returns**

**numpy.array**  
**Processing scores. Range is 0-1, higher indicates more favorable processing.**

**predict\_to\_dataframe** (*self*, *peptides*, *n\_flanks=None*, *c\_flanks=None*, *batch\_size=4096*)  
Predict antigen processing.

See [predict](#) method for parameter descriptions.

**Returns**

**pandas.DataFrame**  
**Processing predictions are in the “score” column. Also includes peptides and flanking sequences.**

**predict\_to\_dataframe\_encoded** (*self*, *sequences*, *batch\_size=4096*)  
Predict antigen processing.

See [predict](#) method for more information.

**Parameters**

**sequences** [FlankingEncoding]  
**batch\_size** [int]

**Returns**

**pandas.DataFrame**

**check\_consistency** (*self*)  
Verify that self.manifest\_df is consistent with instance variables.

Currently only checks for agreement on the total number of models.

Throws AssertionError if inconsistent.

**save** (*self*, *models\_dir*, *model\_names\_to\_write=None*, *write\_metadata=True*)  
Serialize the predictor to a directory on disk. If the directory does not exist it will be created.

The serialization format consists of a file called “manifest.csv” with the configurations of each Class1ProcessingNeuralNetwork, along with per-network files giving the model weights.

**Parameters**

**models\_dir** [string] Path to directory. It will be created if it doesn’t exist.

**classmethod load** (*models\_dir=None*, *max\_models=None*)  
Deserialize a predictor from a directory on disk.

**Parameters**

**models\_dir** [string] Path to directory. If unspecified the default downloaded models are used.

**max\_models** [int, optional] Maximum number of models to load

#### Returns

*Class1ProcessingPredictor* instance

## 5.10 mhcflurry.cluster\_parallelism module

Simple, relatively naive parallel map implementation for HPC clusters.

Used for training MHCflurry models.

`mhcflurry.cluster_parallelism.add_cluster_parallelism_args(parser)`

Add commandline arguments controlling cluster parallelism to an argparse ArgumentParser.

#### Parameters

**parser** [argparse.ArgumentParser]

`mhcflurry.cluster_parallelism.cluster_results_from_args(args, work_function, work_items, constant_data=None, input_serialization_method='pickle', result_serialization_method='pickle', clear_constant_data=False)`

Parallel map configurable using commandline arguments. See the `cluster_results()` function for docs.

The `args` parameter should be an `argparse.Namespace` from an `argparse` parser generated using the `add_cluster_parallelism_args()` function.

#### Parameters

**args**

**work\_function**

**work\_items**

**constant\_data**

**result\_serialization\_method**

**clear\_constant\_data**

#### Returns

**generator**

`mhcflurry.cluster_parallelism.cluster_results(work_function, work_items, constant_data=None, submit_command='sh', results_workdir='./cluster-workdir', additional_complete_file=None, script_prefix_path=None, input_serialization_method='pickle', result_serialization_method='pickle', max_retries=3, clear_constant_data=False)`

Parallel map on an HPC cluster.

Returns [work\_function(item) for item in work\_items] where each invocation of work\_function is performed as a separate HPC cluster job. Order is preserved.

Optionally, “constant data” can be specified, which will be passed to each work\_function() invocation as a keyword argument called constant\_data. This data is serialized once and all workers read it from the same source, which is more efficient than serializing it separately for each worker.

Each worker’s input is serialized to a shared NFS directory and the submit\_command is used to launch a job to process that input. The shared filesystem is polled occasionally to watch for results, which are fed back to the user.

#### Parameters

**work\_function** [A -> B]

**work\_items** [list of A]

**constant\_data** [object]

**submit\_command** [string] For running on LSF, we use “bsub” here.

**results\_workdir** [string] Path to NFS shared directory where inputs and results can be written

**script\_prefix\_path** [string] Path to script that will be invoked to run each worker. A line calling the \_mhcflurry-cluster-worker-entry-point command will be appended to the contents of this file.

**result\_serialization\_method** [string, one of “pickle” or “save\_predictor”] The “save\_predictor” works only when the return type of work\_function is Class1AffinityPredictor

**max\_retries** [int] How many times to attempt to re-launch a failed worker

**clear\_constant\_data** [bool] If True, the constant data dict is cleared on the launching host after it is serialized to disk.

#### Returns

**generator of B**

`mhcflurry.cluster_parallelism.worker_entry_point` (argv=['-b', 'latex', '-v', '-d', '\_build/doctrees', '.\_', '\_build/latex'])

Entry point for the worker command.

#### Parameters

**argv** [list of string]

## 5.11 mhcflurry.common module

`mhcflurry.common.set_keras_backend` (backend=None, num\_threads=None, gpu\_device\_nums=None)

Configure Keras backend to use GPU or CPU. Only tensorflow is supported.

#### Parameters

**backend** [string, optional] one of ‘tensorflow-default’, ‘tensorflow-cpu’, ‘tensorflow-gpu’

**gpu\_device\_nums** [list of int, optional] GPU devices to potentially use

**num\_threads** [int, optional] Tensorflow threads to use

`mhcflurry.common.configure_logging` (verbose=False)

Configure logging module using defaults.

**Parameters**

**verbose** [boolean] If true, output will be at level DEBUG, otherwise, INFO.

`mhcflurry.common.amino_acid_distribution` (*peptides*, *smoothing*=0.0)

Compute the fraction of each amino acid across a collection of peptides.

**Parameters**

**peptides** [list of string]

**smoothing** [float, optional] Small number (e.g. 0.01) to add to all amino acid fractions. The higher the number the more uniform the distribution.

**Returns**

**pandas.Series indexed by amino acids**

`mhcflurry.common.random_peptides` (*num*, *length*=9, *distribution*=None)

Generate random peptides (kmers).

**Parameters**

**num** [int] Number of peptides to return

**length** [int] Length of each peptide

**distribution** [pandas.Series] Maps 1-letter amino acid abbreviations to probabilities. If not specified a uniform distribution is used.

**Returns**

**list of string**

`mhcflurry.common.positional_frequency_matrix` (*peptides*)

Given a set of peptides, calculate a length x amino acids frequency matrix.

**Parameters**

**peptides** [list of string] All of same length

**Returns**

**pandas.DataFrame** Index is position, columns are amino acids

`mhcflurry.common.save_weights` (*weights\_list*, *filename*)

Save model weights to the given filename using numpy's ".npz" format.

**Parameters**

**weights\_list** [list of numpy array]

**filename** [string]

`mhcflurry.common.load_weights` (*filename*)

Restore model weights from the given filename, which should have been created with *save\_weights*.

**Parameters**

**filename** [string]

**Returns**

**list of array**

```
class mhcflurry.common.NumpyJSONEncoder(*, skipkeys=False, ensure_ascii=True,
                                         check_circular=True, allow_nan=True,
                                         sort_keys=False, indent=None, separators=None,
                                         default=None)
```

Bases: `json.encoder.JSONEncoder`

JSON encoder (used with json module) that can handle numpy arrays.

Constructor for JSONEncoder, with sensible defaults.

If `skipkeys` is false, then it is a `TypeError` to attempt encoding of keys that are not `str`, `int`, `float` or `None`. If `skipkeys` is `True`, such items are simply skipped.

If `ensure_ascii` is true, the output is guaranteed to be `str` objects with all incoming non-ASCII characters escaped. If `ensure_ascii` is false, the output can contain non-ASCII characters.

If `check_circular` is true, then lists, dicts, and custom encoded objects will be checked for circular references during encoding to prevent an infinite recursion (which would cause an `OverflowError`). Otherwise, no such check takes place.

If `allow_nan` is true, then `NaN`, `Infinity`, and `-Infinity` will be encoded as such. This behavior is not JSON specification compliant, but is consistent with most JavaScript based encoders and decoders. Otherwise, it will be a `ValueError` to encode such floats.

If `sort_keys` is true, then the output of dictionaries will be sorted by key; this is useful for regression tests to ensure that JSON serializations can be compared on a day-to-day basis.

If `indent` is a non-negative integer, then JSON array elements and object members will be pretty-printed with that indent level. An indent level of 0 will only insert newlines. `None` is the most compact representation.

If specified, `separators` should be an (item\_separator, key\_separator) tuple. The default is `(' ', ': ')` if `indent` is `None` and `(';', ': ')` otherwise. To get the most compact JSON representation, you should specify `(';', ':')` to eliminate whitespace.

If specified, `default` is a function that gets called for objects that can't otherwise be serialized. It should return a JSON encodable version of the object or raise a `TypeError`.

**default** (*self*, *obj*)

Implement this method in a subclass such that it returns a serializable object for `o`, or calls the base implementation (to raise a `TypeError`).

For example, to support arbitrary iterators, you could implement `default` like this:

```
def default(self, o):
    try:
        iterable = iter(o)
    except TypeError:
        pass
    else:
        return list(iterable)
    # Let the base class default method raise the TypeError
    return JSONEncoder.default(self, o)
```



## 5.12 mhcfllurry.custom\_loss module

Custom loss functions.

For losses supporting inequalities, each training data point is associated with one of (=), (<), or (>). For e.g. (>) inequalities, penalization is applied only if the prediction is less than the given value.

`mhcfllurry.custom_loss.get_loss(name)`

Get a custom\_loss.Loss instance by name.

### Parameters

**name** [string]

### Returns

`custom_loss.Loss`

**class** `mhcfllurry.custom_loss.Loss` (*name=None*)

Bases: `object`

Thin wrapper to keep track of neural network loss functions, which could be custom or baked into Keras.

Each subclass or instance should define these properties/methods: - `name` : string - `loss` : string or function

This is what gets passed to `keras.fit()`

- **encode\_y** [numpy.ndarray -> numpy.ndarray] Transformation to apply to regression target before fitting

**loss** (*self, y\_true, y\_pred*)

**get\_keras\_loss** (*self, reduction='sum\_over\_batch\_size'*)

**class** `mhcfllurry.custom_loss.StandardKerasLoss` (*loss\_name='mse'*)

Bases: `mhcfllurry.custom_loss.Loss`

A loss function supported by Keras, such as MSE.

**supports\_inequalities** = **False**

**supports\_multiple\_outputs** = **False**

**static encode\_y** (*y*)

**class** `mhcfllurry.custom_loss.TransformPredictionsLossWrapper` (*loss,*  
*y\_pred\_transform=None*)

Bases: `mhcfllurry.custom_loss.Loss`

Wrapper that applies an arbitrary transform to `y_pred` before calling an underlying loss function.

The `y_pred_transform` function should be a tensor -> tensor function.

**encode\_y** (*self, \*args, \*\*kwargs*)

**loss** (*self, y\_true, y\_pred*)

**class** `mhcfllurry.custom_loss.MSEWithInequalities` (*name=None*)

Bases: `mhcfllurry.custom_loss.Loss`

Supports training a regression model on data that includes inequalities (e.g.  $x < 100$ ). Mean square error is used as the loss for elements with an (=) inequality. For elements with e.g. a ( $> 0.5$ ) inequality, then the loss for that element is  $(y - 0.5)^2$  (standard MSE) if  $y < 500$  and 0 otherwise.

This loss assumes that the normal range for `y_true` and `y_pred` is 0 - 1. As a hack, the implementation uses other intervals for `y_pred` to encode the inequality information.

`y_true` is interpreted as follows:

**between 0 - 1** Regular MSE loss is used. Penalty  $(y_{\text{pred}} - y_{\text{true}})^2$  is applied if `y_pred` is greater or less than `y_true`.

**between 2 - 3:** Treated as a “>” inequality. Penalty  $(y_{\text{pred}} - (y_{\text{true}} - 2))^2$  is applied only if `y_pred` is less than `y_true - 2`.

**between 4 - 5:** Treated as a “<” inequality. Penalty  $(y_{\text{pred}} - (y_{\text{true}} - 4))^2$  is applied only if `y_pred` is greater than `y_true - 4`.

```
name = 'mse_with_inequalities'
```

```
supports_inequalities = True
```

```
supports_multiple_outputs = False
```

```
static encode_y(y, inequalities=None)
```

```
loss(self, y_true, y_pred)
```

```
class mhcflurry.custom_loss.MSEWithInequalitiesAndMultipleOutputs (name=None)
```

Bases: `mhcflurry.custom_loss.Loss`

Loss supporting inequalities and multiple outputs.

This loss assumes that the normal range for `y_true` and `y_pred` is 0 - 1. As a hack, the implementation uses other intervals for `y_pred` to encode the inequality and output-index information.

Inequalities are encoded into the regression target as in the `MSEWithInequalities` loss.

Multiple outputs are encoded by mapping each regression target `x` (after transforming for inequalities) using the rule `x -> x + i * 10` where `i` is the output index.

The reason for explicitly encoding multiple outputs this way (rather than just making the regression target a matrix instead of a vector) is that in our use cases we frequently have missing data in the regression target. This encoding gives a simple way to penalize only on (data point, output index) pairs that have labels.

```
name = 'mse_with_inequalities_and_multiple_outputs'
```

```
supports_inequalities = True
```

```
supports_multiple_outputs = True
```

```
static encode_y(y, inequalities=None, output_indices=None)
```

```
loss(self, y_true, y_pred)
```

```
class mhcflurry.custom_loss.MultiallelicMassSpecLoss (delta=0.2, multiplier=1.0)
```

Bases: `mhcflurry.custom_loss.Loss`

```
name = 'multiallelic_mass_spec_loss'
```

```
supports_inequalities = True
```

```
supports_multiple_outputs = False
```

```
static encode_y(y)
```

```
loss(self, y_true, y_pred)
```

```
mhcflurry.custom_loss.check_shape(name, arr, expected_shape)
```

Raise `ValueError` if `arr.shape != expected_shape`.

#### Parameters

**name** [string] Included in error message to aid debugging

`mhcflurry.custom_loss.cls`  
alias of `mhcflurry.custom_loss.MultiallelicMassSpecLoss`

### 5.13 mhcflurry.data\_dependent\_weights\_initialization module

### Layer-sequential unit-variance initialization for neural networks.

**See:** Mishkin and Matas, “All you need is a good init”. 2016. <https://arxiv.org/abs/1511.06422>

```
mhcflurry.data_dependent_weights_initialization.svd_orthonormal(shape)
mhcflurry.data_dependent_weights_initialization.get_activations(model, layer, X_batch)
mhcflurry.data_dependent_weights_initialization.lsuv_init(model, batch, verbose=True, margin=0.1, max_iter=100)
```

Initialize neural network weights using layer-sequential unit-variance initialization.

**See:** Mishkin and Matas, “All you need is a good init”. 2016. <https://arxiv.org/abs/1511.06422>

## Parameters

**model** [keras.Model]  
**batch** [dict] Training data, as would be passed keras.Model.fit()  
**verbose** [boolean] Whether to print progress to stdout  
**margin** [float]  
**max\_iter** [int]

## Returns

**keras.Model** Same as what was passed in.

## 5.14 mhcflurry.downloads module

Manage local downloaded data.

```

mhcflurry.downloads.get_downloads_dir()
    Return the path to local downloaded data

```

`mhcflurry.downloads.get_current_release()`  
Return the current downloaded data release

```
mhcflurry.downloads.get_downloads_metadata()
```

Return the contents of downloads.yml as a dict

`mhcflurry.downloads.get_default_class1_models_dir` (*test\_exists=True*)  
Return the absolute path to the default class1 models dir.

If environment variable `MHCFLURRY_DEFAULT_CLASS1_MODELS` is set to an absolute path, return that path. If it's set to a relative path (i.e. does not start with `/`) then return that path taken to be relative to the `mhcflurry` downloads dir.

If environment variable `MHCFLURRY_DEFAULT_CLASS1_MODELS` is NOT set, then return the path to downloaded models in the “models\_class1” download.

**Parameters**

**test\_exists** [boolean, optional] Whether to raise an exception of the path does not exist

**Returns**

**string** [absolute path]

```
mhcflurry.downloads.get_default_class1_presentation_models_dir(test_exists=True)
```

Return the absolute path to the default class1 presentation models dir.

See `get_default_class1_models_dir`.

If environment variable `MHCFLURRY_DEFAULT_CLASS1_PRESENTATION_MODELS` is set to an absolute path, return that path. If it’s set to a relative path (does not start with /) then return that path taken to be relative to the mhcflurry downloads dir.

**Parameters**

**test\_exists** [boolean, optional] Whether to raise an exception of the path does not exist

**Returns**

**string** [absolute path]

```
mhcflurry.downloads.get_default_class1_processing_models_dir(test_exists=True)
```

Return the absolute path to the default class1 processing models dir.

See `get_default_class1_models_dir`.

If environment variable `MHCFLURRY_DEFAULT_CLASS1_PROCESSING_MODELS` is set to an absolute path, return that path. If it’s set to a relative path (does not start with /) then return that path taken to be relative to the mhcflurry downloads dir.

**Parameters**

**test\_exists** [boolean, optional] Whether to raise an exception of the path does not exist

**Returns**

**string** [absolute path]

```
mhcflurry.downloads.get_current_release_downloads()
```

Return a dict of all available downloads in the current release.

The dict keys are the names of the downloads. The values are a dict with two entries:

**downloaded** [bool] Whether the download is currently available locally

**metadata** [dict] Info about the download from downloads.yml such as URL

**up\_to\_date** [bool or None] Whether the download URL(s) match what was used to download the current data.  
This is None if it cannot be determined.

```
mhcflurry.downloads.get_path(download_name, filename="", test_exists=True)
```

Get the local path to a file in a MHCflurry download

**Parameters**

**download\_name** [string]

**filename** [string] Relative path within the download to the file of interest

**test\_exists** [boolean] If True (default) throw an error telling the user how to download the data if the file does not exist

**Returns****string giving local absolute path**`mhcflurry.downloads.configure()`

Setup various global variables based on environment variables.

## 5.15 mhcflurry.downloads\_command module

Download MHCflurry released datasets and trained models.

Examples

**Fetch the default downloads:** `$ mhcflurry-downloads fetch`**Fetch a specific download:** `$ mhcflurry-downloads fetch models_class1_pan`**Get the path to a download:** `$ mhcflurry-downloads path models_class1_pan`**Get the URL of a download:** `$ mhcflurry-downloads url models_class1_pan`**Summarize available and fetched downloads:** `$ mhcflurry-downloads info``mhcflurry.downloads_command.run(argv=['-b', 'latex', '-v', '-d', '_build/doctrees', '.', '_build/latex'])``mhcflurry.downloads_command.mkdir_p(path)`Make directories as needed, similar to `mkdir -p` in a shell.From: <http://stackoverflow.com/questions/600268/mkdir-p-functionality-in-python>`mhcflurry.downloads_command.yes_no(booleant)`

```
class mhcflurry.downloads_command.TqdmUpTo(
    iterable=None, desc=None, total=None,
    leave=True, file=None, ncols=None, mininterval=0.1,
    maxinterval=10.0, miniters=None,
    ascii=None, disable=False, unit='it',
    unit_scale=False, dynamic_ncols=False,
    smoothing=0.3, bar_format=None,
    initial=0, position=None, postfix=None,
    unit_divisor=1000, gui=False, **kwargs)
```

Bases: `tqdm._tqdm.tqdm`Provides `update_to(n)` which uses `tqdm.update(delta_n)`.**Parameters****iterable** [iterable, optional] Iterable to decorate with a progressbar. Leave blank to manually manage the updates.**desc** [str, optional] Prefix for the progressbar.**total** [int, optional] The number of expected iterations. If unspecified, `len(iterable)` is used if possible. As a last resort, only basic progress statistics are displayed (no ETA, no progressbar). If `gui` is `True` and this parameter needs subsequent updating, specify an initial arbitrary large positive integer, e.g. `int(9e9)`.**leave** [bool, optional] If [default: `True`], keeps all traces of the progressbar upon termination of iteration.**file** [`io.TextIOWrapper` or `io.StringIO`, optional] Specifies where to output the progress messages (default: `sys.stderr`). Uses `file.write(str)` and `file.flush()` methods.

**ncols** [int, optional] The width of the entire output message. If specified, dynamically resizes the progressbar to stay within this bound. If unspecified, attempts to use environment width. The fallback is a meter width of 10 and no limit for the counter and statistics. If 0, will not print any meter (only stats).

**mininterval** [float, optional] Minimum progress display update interval, in seconds [default: 0.1].

**maxinterval** [float, optional] Maximum progress display update interval, in seconds [default: 10]. Automatically adjusts `miniters` to correspond to `mininterval` after long display update lag. Only works if `dynamic_miniters` or `monitor` thread is enabled.

**miniters** [int, optional] Minimum progress display update interval, in iterations. If 0 and `dynamic_miniters`, will automatically adjust to equal `mininterval` (more CPU efficient, good for tight loops). If > 0, will skip display of specified number of iterations. Tweak this and `mininterval` to get very efficient loops. If your progress is erratic with both fast and slow iterations (network, skipping items, etc) you should set `miniters=1`.

**ascii** [bool, optional] If unspecified or False, use unicode (smooth blocks) to fill the meter. The fallback is to use ASCII characters 1-9 #.

**disable** [bool, optional] Whether to disable the entire progressbar wrapper [default: False]. If set to None, disable on non-TTY.

**unit** [str, optional] String that will be used to define the unit of each iteration [default: it].

**unit\_scale** [bool or int or float, optional] If 1 or True, the number of iterations will be reduced/scaled automatically and a metric prefix following the International System of Units standard will be added (kilo, mega, etc.) [default: False]. If any other non-zero number, will scale `total` and `n`.

**dynamic\_ncols** [bool, optional] If set, constantly alters `ncols` to the environment (allowing for window resizes) [default: False].

**smoothing** [float, optional] Exponential moving average smoothing factor for speed estimates (ignored in GUI mode). Ranges from 0 (average speed) to 1 (current/instantaneous speed) [default: 0.3].

**bar\_format** [str, optional] Specify a custom bar string formatting. May impact performance. [default: '{l\_bar}{bar}{r\_bar}'], where `l_bar`='{desc}: {percentage:3.0f}%' and `r_bar`='{n\_fmt}/{total\_fmt} [{elapsed}<{remaining}], '{rate\_fmt}{postfix}']

**Possible vars:** `l_bar`, `bar`, `r_bar`, `n`, `n_fmt`, `total`, `total_fmt`, `percentage`, `rate`, `rate_fmt`, `rate_noinv`, `rate_noinv_fmt`, `rate_inv`, `rate_inv_fmt`, `elapsed`, `remaining`, `desc`, `postfix`.

Note that a trailing ": " is automatically removed after {desc} if the latter is empty.

**initial** [int, optional] The initial counter value. Useful when restarting a progress bar [default: 0].

**position** [int, optional] Specify the line offset to print this bar (starting from 0) Automatic if unspecified. Useful to manage multiple bars at once (eg, from threads).

**postfix** [dict or \*, optional] Specify additional stats to display at the end of the bar. Calls `set_postfix(**postfix)` if possible (dict).

**unit\_divisor** [float, optional] [default: 1000], ignored unless `unit_scale` is True.

**gui** [bool, optional] WARNING: internal parameter - do not use. Use `tqdm_gui(...)` instead. If set, will attempt to use matplotlib animations for a graphical output [default: False].

**Returns****out** [decorated iterator.]**update\_to** (*self*, *b=1*, *bsize=1*, *tsize=None*)**b** [int, optional] Number of blocks transferred so far [default: 1].**bsize** [int, optional] Size of each block (in tqdm units) [default: 1].**tsize** [int, optional] Total size (in tqdm units). If [default: None] remains unchanged.`mhcflurry.downloads_command.fetch_subcommand(args)``mhcflurry.downloads_command.info_subcommand(args)``mhcflurry.downloads_command.path_subcommand(args)`

Print the local path to a download

`mhcflurry.downloads_command.url_subcommand(args)`

Print the URL(s) for a download

## 5.16 mhcflurry.encodable\_sequences module

Class for encoding variable-length peptides to fixed-size numerical matrices

**exception** `mhcflurry.encodable_sequences.EncodingError` (*message*, *supported\_peptide\_lengths*)Bases: `ValueError`

Exception raised when peptides cannot be encoded

**class** `mhcflurry.encodable_sequences.EncodableSequences` (*sequences*)Bases: `object`

Class for encoding variable-length peptides to fixed-size numerical matrices

This class caches various encodings of a list of sequences.

In practice this is used only for peptides. To encode MHC allele sequences, see `AlleleEncoding`.**unknown\_character** = 'X'**classmethod** `create` (*sequences*)Factory that returns an `EncodableSequences` given a list of strings. As a convenience, you can also pass it an `EncodableSequences` instance, in which case the object is returned unchanged.**variable\_length\_to\_fixed\_length\_categorical** (*self*, *alignment\_method='pad\_middle'*, *left\_edge=4*, *right\_edge=4*, *max\_length=15*)

Encode variable-length sequences to a fixed-size index-encoded (integer) matrix.

See `sequences_to_fixed_length_index_encoded_array` for details.**Parameters****alignment\_method** [string] One of “pad\_middle” or “left\_pad\_right\_pad”**left\_edge** [int, size of fixed-position left side] Only relevant for pad\_middle alignment method**right\_edge** [int, size of the fixed-position right side] Only relevant for pad\_middle alignment method**max\_length** [maximum supported peptide length]

**Returns**

**numpy.array of integers with shape (num sequences, encoded length)**

For `pad_middle`, the encoded length is `max_length`. For `left_pad_right_pad`, it's `3 * max_length`.

```
variable_length_to_fixed_length_vector_encoding(self,          vec-  
                                              tor_encoding_name,    align-  
                                              ment_method='pad_middle',  
                                              left_edge=4,      right_edge=4,  
                                              max_length=15, trim=False, al-  
                                              low_unsupported_amino_acids=False)
```

Encode variable-length sequences to a fixed-size matrix. Amino acids are encoded as specified by the `vector_encoding_name` argument.

See [sequences\\_to\\_fixed\\_length\\_index\\_encoded\\_array](#) for details.

See also: `variable_length_to_fixed_length_categorical`.

**Parameters**

**vector\_encoding\_name** [string] How to represent amino acids. One of “BLO-SUM62”, “one-hot”, etc. Full list of supported vector encodings is given by `available_vector_encodings()`.

**alignment\_method** [string] One of “pad\_middle” or “left\_pad\_right\_pad”

**left\_edge** [int] Size of fixed-position left side. Only relevant for `pad_middle` alignment method

**right\_edge** [int] Size of the fixed-position right side. Only relevant for `pad_middle` alignment method

**max\_length** [int] Maximum supported peptide length

**trim** [bool] If True, longer sequences will be trimmed to fit the maximum supported length. Not supported for all alignment methods.

**allow\_unsupported\_amino\_acids** [bool] If True, non-canonical amino acids will be replaced with the X character before encoding.

**Returns**

**numpy.array with shape (num sequences, encoded length, m)**

**where**

- `m` is the vector encoding length (usually 21).
- encoded length is `max_length` if `alignment_method` is `pad_middle`; `3 * max_length` if it's `left_pad_right_pad`.

```
classmethod sequences_to_fixed_length_index_encoded_array(sequences, align-  
                                                         ment_method='pad_middle',  
                                                         left_edge=4,  
                                                         right_edge=4,  
                                                         max_length=15,  
                                                         trim=False, al-  
                                                         low_unsupported_amino_acids=False)
```

Encode variable-length sequences to a fixed-size index-encoded (integer) matrix.

How variable length sequences get mapped to fixed length is set by the “`alignment_method`” argument. Supported alignment methods are:



**pad\_middle** Encoding designed for preserving the anchor positions of class I peptides. This is what is used in allele-specific models.

Each string must be of length at least `left_edge + right_edge` and at most `max_length`. The first `left_edge` characters in the input always map to the first `left_edge` characters in the output. Similarly for the last `right_edge` characters. The middle characters are filled in based on the length, with the X character filling in the blanks.

Example:

AAAACDDDD -> AAAAXXXCXXXDDDD

**left\_pad\_centered\_right\_pad** Encoding that makes no assumptions on anchor positions but is 3x larger than `pad_middle`, since it duplicates the peptide (left aligned + centered + right aligned). This is what is used for the pan-allele models.

Example:

AAAACDDDD -> AAAACDDDDXXXXXXXXXXAAAACDDDDXXXXXXXXXXAAAACDDDD

**left\_pad\_right\_pad** Same as `left_pad_centered_right_pad` but only includes left- and right-padded peptide.

Example:

AAAACDDDD -> AAAACDDDDXXXXXXXXXXXXAAAACDDDD

### Parameters

**sequences** [list of string]

**alignment\_method** [string] One of “`pad_middle`” or “`left_pad_right_pad`”

**left\_edge** [int] Size of fixed-position left side. Only relevant for `pad_middle` alignment method

**right\_edge** [int] Size of the fixed-position right side. Only relevant for `pad_middle` alignment method

**max\_length** [int] maximum supported peptide length

**trim** [bool] If True, longer sequences will be trimmed to fit the maximum supported length. Not supported for all alignment methods.

**allow\_unsupported\_amino\_acids** [bool] If True, non-canonical amino acids will be replaced with the X character before encoding.

### Returns

**numpy.array of integers with shape (num sequences, encoded length)**

For `pad_middle`, the encoded length is `max_length`. For `left_pad_right_pad`,

it's `2 * max_length`. For `left_pad_centered_right_pad`, it's

`3 * max_length`.

## 5.17 mhcflurry.ensemble\_centrality module

Measures of centrality (e.g. mean) used to combine predictions across an ensemble. The input to these functions are log affinities, and they are expected to return a centrality measure also in log-space.

`mhcflurry.ensemble_centrality.robust_mean(log_values)`

Mean of values falling within the 25-75 percentiles.

### Parameters

**log\_values** [2-d numpy.array] Center is computed along the second axis (i.e. per row).

### Returns

**center** [numpy.array of length log\_values.shape[1]]

## 5.18 mhcflurry.fasta module

Adapted from pyensembl, [github.com/openvax/pyensembl](https://github.com/openvax/pyensembl) Original implementation by Alex Rubinsteyn.

The worse sin in bioinformatics is to write your own FASTA parser. We're doing this to avoid adding another dependency to MHCflurry, however.

`mhcflurry.fasta.read_fasta_to_dataframe(filename)`

**class** `mhcflurry.fasta.FastaParser`

Bases: object

FastaParser object consumes lines of a FASTA file incrementally.

**iterate\_over\_file** (*self, fasta\_path*)

Generator that yields identifiers paired with sequences.

**static open\_file** (*fasta\_path*)

Open either a text file or compressed gzip file as a stream of bytes.

## 5.19 mhcflurry.flanking\_encoding module

Class for encoding variable-length flanking and peptides to fixed-size numerical matrices

**class** `mhcflurry.flanking_encoding.EncodingResult` (*array, peptide\_lengths*)

Bases: tuple

Create new instance of EncodingResult(array, peptide\_lengths)

**property array**

Alias for field number 0

**property peptide\_lengths**

Alias for field number 1

**class** `mhcflurry.flanking_encoding.FlankingEncoding` (*peptides, n\_flanks, c\_flanks*)

Bases: object

Encode peptides and optionally their N- and C-flanking sequences into fixed size numerical matrices. Similar to EncodableSequences but with support for flanking sequences and the encoding scheme used by the processing predictor.

Instances of this class have an immutable list of peptides with flanking sequences. Encodings are cached in the instances for faster performance when the same set of peptides needs to be encoded more than once.

Constructor. Sequences of any lengths can be passed.

#### Parameters

**peptides** [list of string] Peptide sequences

**n\_flanks** [list of string [same length as peptides]] Upstream sequences

**c\_flanks** [list of string [same length as peptides]] Downstream sequences

**unknown\_character** = 'X'

**vector\_encode** (*self*, *vector\_encoding\_name*, *peptide\_max\_length*, *n\_flank\_length*, *c\_flank\_length*, *allow\_unsupported\_amino\_acids=True*)

Encode variable-length sequences to a fixed-size matrix.

#### Parameters

**vector\_encoding\_name** [string] How to represent amino acids. One of “BLOSUM62”, “one-hot”, etc. See `amino_acid.available_vector_encodings()`.

**peptide\_max\_length** [int] Maximum supported peptide length.

**n\_flank\_length** [int] Maximum supported N-flank length

**c\_flank\_length** [int] Maximum supported C-flank length

**allow\_unsupported\_amino\_acids** [bool] If True, non-canonical amino acids will be replaced with the X character before encoding.

#### Returns

**numpy.array with shape (num sequences, length, m)**

where

- num sequences is number of peptides, i.e. `len(self)`
- length is `peptide_max_length + n_flank_length + c_flank_length`
- m is the vector encoding length (usually 21).

**static encode** (*vector\_encoding\_name*, *df*, *peptide\_max\_length*, *n\_flank\_length*, *c\_flank\_length*, *allow\_unsupported\_amino\_acids=False*)

Encode variable-length sequences to a fixed-size matrix.

Helper function. Users should use `vector_encode`.

#### Parameters

**vector\_encoding\_name** [string]

**df** [pandas.DataFrame]

**peptide\_max\_length** [int]

**n\_flank\_length** [int]

**c\_flank\_length** [int]

**allow\_unsupported\_amino\_acids** [bool]

#### Returns

**numpy.array**

## 5.20 mhcflurry.hyperparameters module

Hyperparameter (neural network options) management

**class** mhcflurry.hyperparameters.HyperparameterDefaults (\*\*defaults)

Bases: object

Class for managing hyperparameters. Thin wrapper around a dict.

Instances of this class are a specification of the hyperparameters *supported* by a model and their defaults. The particular hyperparameter settings to be used, for example, to train a model are kept in plain dicts.

**extend** (self, other)

Return a new HyperparameterDefaults instance containing the hyperparameters from the current instance combined with those from other.

It is an error if self and other have any hyperparameters in common.

**with\_defaults** (self, obj)

Given a dict of hyperparameter settings, return a dict containing those settings augmented by the defaults for any keys missing from the dict.

**subselect** (self, obj)

Filter a dict of hyperparameter settings to only those keys defined in this HyperparameterDefaults .

**check\_valid\_keys** (self, obj)

Given a dict of hyperparameter settings, throw an exception if any keys are not defined in this HyperparameterDefaults instance.

**models\_grid** (self, \*\*kwargs)

Make a grid of models by taking the cartesian product of all specified model parameter lists.

### Parameters

The valid kwarg parameters are the entries of this HyperparameterDefaults instance. Each parameter must be a list giving the values to search across.

### Returns

list of dict giving the parameters for each model. The length of the list is the product of the lengths of the input lists.

## 5.21 mhcflurry.local\_parallelism module

Infrastructure for “local” parallelism, i.e. multiprocessing parallelism on one compute node.

mhcflurry.local\_parallelism.add\_local\_parallelism\_args (parser)

Add local parallelism arguments to the given argparse.ArgumentParser.

### Parameters

**parser** [argparse.ArgumentParser]

mhcflurry.local\_parallelism.worker\_pool\_with\_gpu\_assignments\_from\_args (args)

Create a multiprocessing.Pool where each worker uses its own GPU.

Uses commandline arguments. See [worker\\_pool\\_with\\_gpu\\_assignments](#).

### Parameters

**args** [argparse.ArgumentParser]

#### Returns

**multiprocessing.Pool**

`mhcflurry.local_parallelism.worker_pool_with_gpu_assignments` (*num\_jobs*,  
*num\_gpus=0*,  
*backend=None*,  
*max\_workers\_per\_gpu=1*,  
*max\_tasks\_per\_worker=None*,  
*worker\_log\_dir=None*)

Create a multiprocessing.Pool where each worker uses its own GPU.

#### Parameters

**num\_jobs** [int] Number of worker processes.

**num\_gpus** [int]

**backend** [string]

**max\_workers\_per\_gpu** [int]

**max\_tasks\_per\_worker** [int]

**worker\_log\_dir** [string]

#### Returns

**multiprocessing.Pool**

`mhcflurry.local_parallelism.make_worker_pool` (*processes=None*, *initializer=None*, *initializer\_kwargs\_per\_process=None*,  
*max\_tasks\_per\_worker=None*)

Convenience wrapper to create a multiprocessing.Pool.

This function adds support for per-worker initializer arguments, which are not natively supported by the multiprocessing module. The motivation for this feature is to support allocating each worker to a (different) GPU.

**IMPLEMENTATION NOTE:** The per-worker initializer arguments are implemented using a Queue. Each worker reads its arguments from this queue when it starts. When it terminates, it adds its initializer arguments back to the queue, so a future process can initialize itself using these arguments.

There is one issue with this approach, however. If a worker crashes, it never repopulates the queue of initializer arguments. This will prevent any future worker from re-using those arguments. To deal with this issue we add a second ‘backup queue’. This queue always contains the full set of initializer arguments: whenever a worker reads from it, it always pushes the pop’d args back to the end of the queue immediately. If the primary arg queue is ever empty, then workers will read from this backup queue.

#### Parameters

**processes** [int] Number of workers. Default: num CPUs.

**initializer** [function, optional] Init function to call in each worker

**initializer\_kwargs\_per\_process** [list of dict, optional] Arguments to pass to initializer function for each worker. Length of list must equal the number of workers.

**max\_tasks\_per\_worker** [int, optional] Restart workers after this many tasks. Requires Python >=3.2.

#### Returns

**multiprocessing.Pool**

```
mhcflurry.local_parallelism.worker_init_entry_point(init_function, arg_queue=None,  
                                                    backup_arg_queue=None)
```

```
mhcflurry.local_parallelism.worker_init(keras_backend=None, gpu_device_nums=None,  
                                       worker_log_dir=None)
```

**exception** mhcflurry.local\_parallelism.WrapException

Bases: Exception

Add traceback info to exception so exceptions raised in worker processes can still show traceback info when re-raised in the parent.

```
mhcflurry.local_parallelism.call_wrapped(function, *args, **kwargs)
```

Run function on args and kwargs and return result, wrapping any exception raised in a WrapException.

**Parameters**

**function** [arbitrary function]

Any other arguments provided are passed to the function.

**Returns**

object

```
mhcflurry.local_parallelism.call_wrapped_kwargs(function, kwargs)
```

Invoke function on given kwargs and return result, wrapping any exception raised in a WrapException.

**Parameters**

**function** [arbitrary function]

**kwargs** [dict]

**Returns**

object

result of calling function(\*\*kwargs)

## 5.22 mhcflurry.percent\_rank\_transform module

Class for transforming arbitrary values into percent ranks given a distribution.

**class** mhcflurry.percent\_rank\_transform.PercentRankTransform

Bases: object

Transform arbitrary values into percent ranks.

**fit** (*self*, *values*, *bins*)

Fit the transform using the given values (in our case ic50s).

**Parameters**

**values** [ic50 values]

**bins** [bins for the cumulative distribution function] Anything that can be passed to numpy.histogram’s “bins” argument can be used here.

**transform** (*self*, *values*)

Return percent ranks (range [0, 100]) for the given values.

**to\_series** (*self*)

Serialize the fit to a pandas.Series.

The index on the series gives the bin edges and the values give the CDF.

**Returns**

**pandas.Series**

**static from\_series** (*series*)

Deserialize a PercentRankTransform the given pandas.Series, as returned by *to\_series()*.

**Parameters**

**series** [pandas.Series]

**Returns**

**PercentRankTransform**

## 5.23 mhcflurry.predict\_command module

Run MHCflurry predictor on specified peptides.

By default, the presentation predictor is used, and predictions for MHC I binding affinity, antigen processing, and the composite presentation score are returned. If you just want binding affinity predictions, pass *-affinity-only*.

Examples:

Write a CSV file containing the contents of INPUT.csv plus additional columns giving MHCflurry predictions:

```
$ mhcflurry-predict INPUT.csv -out RESULT.csv
```

The input CSV file is expected to contain columns “allele”, “peptide”, and, optionally, “n\_flank”, and “c\_flank”.

If *--out* is not specified, results are written to stdout.

You can also run on alleles and peptides specified on the commandline, in which case predictions are written for *all combinations* of alleles and peptides:

```
$ mhcflurry-predict -alleles HLA-A*02:01 H-2Kb -peptides SIINFEKL DENDREKLLL
```

Instead of individual alleles (in a CSV or on the command line), you can also give a comma separated list of alleles giving a sample genotype. In this case, the tightest binding affinity across the alleles for the sample will be returned. For example:

```
$ mhcflurry-predict -peptides SIINFEKL DENDREKLLL -alleles HLA-A*02:01,HLA-A*03:01,HLA-B*57:01,HLA-B*45:01,HLA-C*02:01,HLA-C*07:02 HLA-A*01:01,HLA-A*02:06,HLA-B*44:02,HLA-B*07:02,HLA-C*01:01,HLA-C*03:01
```

will give the tightest predicted affinities across alleles for each of the two genotypes specified for each peptide.

```
mhcflurry.predict_command.run(argv=['-b', 'latex', '-v', '-d', '_build/doctrees', '.', '_build/latex'])
```

## 5.24 mhcflurry.predict\_scan\_command module

Scan protein sequences using the MHCflurry presentation predictor.

By default, sub-sequences (peptides) with affinity percentile ranks less than 2.0 are returned. You can also specify `--results-all` to return predictions for all peptides, or `--results-best` to return the top peptide for each sequence.

Examples:

Scan a set of sequences in a FASTA file for binders to any alleles in a MHC I genotype:

```
$ mhcflurry-predict-scan test/data/example.fasta --alleles HLA-A*02:01,HLA-A*03:01,HLA-B*57:01,HLA-B*45:01,HLA-C*02:01,HLA-C*07:02
```

Instead of a FASTA, you can also pass a CSV that has “sequence\_id” and “sequence” columns.

You can also specify multiple MHC I genotypes to scan as space-separated arguments to the `--alleles` option:

```
$ mhcflurry-predict-scan test/data/example.fasta --alleles HLA-A*02:01,HLA-A*03:01,HLA-B*57:01,HLA-B*45:01,HLA-C*02:02,HLA-C*07:02 HLA-A*01:01,HLA-A*02:06,HLA-B*44:02,HLA-B*07:02,HLA-C*01:02,HLA-C*03:01
```

If `--out` is not specified, results are written to standard out.

You can also specify sequences on the commandline:

```
mhcflurry-predict-scan --sequences MGYINVFAFPFTIYSLLLCRMNSRNYIAQVDVVNFNLT --alleles HLA-A*02:01,HLA-A*03:01,HLA-B*57:01,HLA-B*45:01,HLA-C*02:02,HLA-C*07:02
```

```
mhcflurry.predict_scan_command.run(argv=['-b', 'latex', '-v', '-d', '_build/doctrees', '', '_build/latex'])
```

## 5.25 mhcflurry.random\_negative\_peptides module

```
class mhcflurry.random_negative_peptides.RandomNegativePeptides (**hyperparameters)
    Bases: object
```

Generate random negative (peptide, allele) pairs. These are used during model training, where they are resampled at each epoch.

```
hyperparameter_defaults = <mhcflurry.hyperparameters.HyperparameterDefaults object>
    Hyperparameters for random negative peptides.
```

**Number of random negatives will be:** `random_negative_rate * (num measurements) + random_negative_constant`

where the exact meaning of (num measurements) depends on the particular `random_negative_method` in use.

If `random_negative_match_distribution` is True, then the amino acid frequencies of the training data peptides are used to generate the random peptides.

**Valid values for `random_negative_method` are:**

“**by\_length**”: used for allele-specific prediction. See description in [`RandomNegativePeptides.plan\_by\_length`](#) method.

“**by\_allele**”: used for pan-allele prediction. See [`RandomNegativePeptides.plan\_by\_allele`](#) method.

“**by\_allele\_equalize\_nonbinders**”: used for pan-allele prediction. See [`RandomNegativePeptides.plan\_by\_allele\_equalize\_nonbinders`](#) method.



“recommended”: the default. Use `by_length` if the predictor is allele-specific and `by_allele` if it’s pan-allele.

**plan** (*self*, *peptides*, *affinities*, *alleles=None*, *inequalities=None*)

Calculate the number of random negatives for each allele and peptide length. Call this once after instantiating the object.

#### Parameters

**peptides** [list of string]

**affinities** [list of float]

**alleles** [list of string, optional]

**inequalities** [list of string (“>”, “<”, or “=”), optional]

#### Returns

**pandas.DataFrame** indicating number of random negatives for each length and allele.

**plan\_by\_length** (*self*, *df\_all*, *df\_binders=None*, *df\_nonbinders=None*)

Generate a random negative plan using the “by\_length” policy.

Parameters are as in the `plan` method. No return value.

Used for allele-specific predictors. Does not work well for pan-allele.

Different numbers of random negatives per length. Alleles are sampled proportionally to the number of times they are used in the training data.

**plan\_by\_allele** (*self*, *df\_all*, *df\_binders=None*, *df\_nonbinders=None*)

Generate a random negative plan using the “by\_allele” policy.

Parameters are as in the `plan` method. No return value.

For each allele, a particular number of random negatives are used for all lengths. Across alleles, the number of random negatives varies; within an allele, the number of random negatives for each length is a constant

**plan\_by\_allele\_equalize\_nonbinders** (*self*, *df\_all*, *df\_binders*, *df\_nonbinders*)

Generate a random negative plan using the “by\_allele\_equalize\_nonbinders” policy.

Parameters are as in the `plan` method. No return value.

Requires that the `random_negative_binder_threshold` hyperparameter is set.

In a first step, the number of random negatives selected by the “by\_allele” method are added (see `plan_by_allele`). Then, the total number of non-binders are calculated for each allele and length. This total includes non-binder measurements in the training data plus the random negative peptides added in the first step. In a second step, additional random negative peptides are added so that for each allele, all peptide lengths have the same total number of non-binders.

**get\_alleles** (*self*)

Get the list of alleles corresponding to each random negative peptide as returned by `get_peptides`. This does NOT change and can be safely called once and reused.

#### Returns

**list of string**

**get\_peptides** (*self*)

Get the list of random negative peptides. This will be different each time the method is called.

#### Returns

list of string

**get\_total\_count** (*self*)

Total number of planned random negative peptides.

**Returns**

int

## 5.26 mhcflurry.regression\_target module

`mhcflurry.regression_target.from_ic50` (*ic50*, *max\_ic50=50000.0*)

Convert ic50s to regression targets in the range [0.0, 1.0].

**Parameters**

**ic50** [numpy.array of float]

**Returns**

numpy.array of float

`mhcflurry.regression_target.to_ic50` (*x*, *max\_ic50=50000.0*)

Convert regression targets in the range [0.0, 1.0] to ic50s in the range [0, 50000.0].

**Parameters**

**x** [numpy.array of float]

**Returns**

numpy.array of float

## 5.27 mhcflurry.scoring module

Measures of prediction accuracy

`mhcflurry.scoring.make_scores` (*ic50\_y*, *ic50\_y\_pred*, *sample\_weight=None*, *threshold\_nm=500*,  
*max\_ic50=50000*)

Calculate AUC, F1, and Kendall Tau scores.

**Parameters**

**ic50\_y** [float list] true IC50s (i.e. affinities)

**ic50\_y\_pred** [float list] predicted IC50s

**sample\_weight** [float list [optional]]

**threshold\_nm** [float [optional]]

**max\_ic50** [float [optional]]

**Returns**

dict with entries “auc”, “f1”, “tau”

## 5.28 mhcflurry.select\_allele\_specific\_models\_command module

Model select class1 single allele models.

```
mhcflurry.select_allele_specific_models_command.run(argv=['-b', 'latex', '-v', '-d', '_build/doctrees', '', '_build/latex'])
```

```
class mhcflurry.select_allele_specific_models_command.ScrambledPredictor(predictor)
    Bases: object
```

```
    predict(self, peptides, allele)
```

```
mhcflurry.select_allele_specific_models_command.model_select(allele)
```

```
mhcflurry.select_allele_specific_models_command.cache_encoding(predictor, peptides)
```

```
class mhcflurry.select_allele_specific_models_command.ScoreFunction(function,
                                                                    sum-
                                                                    mary=None)
    Bases: object
```

Thin wrapper over a score function (Class1AffinityPredictor -> float). Used to keep a summary string associated with the function.

```
class mhcflurry.select_allele_specific_models_command.CombinedModelSelector(model_selectors,
                                                                    weights=None,
                                                                    min_contribution_per_
                                                                    allele=None)
    Bases: object
```

Model selector that computes a weighted average over other model selectors.

```
    usable_for_allele(self, allele)
```

```
    plan_summary(self, allele)
```

```
    score_function(self, allele, dry_run=False)
```

```
class mhcflurry.select_allele_specific_models_command.ConsensusModelSelector(predictor,
                                                                    num_peptides_per_allele,
                                                                    multi-
                                                                    ply_score_by_value)
    Bases: object
```

Model selector that scores sub-ensembles based on their Kendall tau consistency with the full ensemble over a set of random peptides.

```
    usable_for_allele(self, allele)
```

```
    max_absolute_value(self, allele)
```

```
    plan_summary(self, allele)
```

```
    score_function(self, allele)
```

```
class mhcflurry.select_allele_specific_models_command.MSEModelSelector(df,
                                                                    pre-
                                                                    dic-
                                                                    tor,
                                                                    min_measurements=1,
                                                                    mul-
                                                                    ti-
                                                                    ply_score_by_data_size=True)
```

Bases: object

Model selector that uses mean-squared error to score models. Inequalities are supported.

**usable\_for\_allele** (*self*, *allele*)

**max\_absolute\_value** (*self*, *allele*)

**plan\_summary** (*self*, *allele*)

**score\_function** (*self*, *allele*)

```
class mhcflurry.select_allele_specific_models_command.MassSpecModelSelector(df,
                                                                    pre-
                                                                    dic-
                                                                    tor,
                                                                    de-
                                                                    coys_per_length=0,
                                                                    min_measurements=1,
                                                                    mul-
                                                                    ti-
                                                                    ply_score_by_data_size=True)
```

Bases: object

Model selector that uses PPV of differentiating decoys from hits from mass-spec experiments.

**static ppv** (*y\_true*, *predictions*)

**usable\_for\_allele** (*self*, *allele*)

**max\_absolute\_value** (*self*, *allele*)

**plan\_summary** (*self*, *allele*)

**score\_function** (*self*, *allele*)

## 5.29 mhcflurry.select\_pan\_allele\_models\_command module

Model select class1 pan-allele models.

APPROACH: For each training fold, we select at least min and at most max models (where min and max are set by the `--min/max-models-per-fold` argument) using a step-up (forward) selection procedure. The final ensemble is the union of all selected models across all folds.

```
mhcflurry.select_pan_allele_models_command.mse(predictions, actual, in-
                                                equalities=None, affin-
                                                ties_are_already_01_transformed=False)
```

Mean squared error of predictions vs. actual

### Parameters

**predictions** [list of float]

**actual** [list of float]

**inequalities** [list of string (“>”, “<”, or “=”)]

**affinities\_are\_already\_01\_transformed** [boolean] Predictions and actual are taken to be nanomolar affinities if `affinities_are_already_01_transformed` is False, otherwise 0-1 values.

#### Returns

**float**

```
mhcflurry.select_pan_allele_models_command.run(argv=['-b', 'latex', '-v', '-d',
                                                    '_build/doctrees', ':', '_build/latex'])
```

```
mhcflurry.select_pan_allele_models_command.do_model_select_task(item, constant_data={})
```

```
mhcflurry.select_pan_allele_models_command.model_select(fold_num, models,
                                                         min_models, max_models,
                                                         constant_data={})
```

Model select for a fold.

#### Parameters

**fold\_num** [int]

**models** [list of Class1NeuralNetwork]

**min\_models** [int]

**max\_models** [int]

**constant\_data** [dict]

#### Returns

dict with keys ‘fold\_num’, ‘selected\_indices’, ‘summary’

## 5.30 mhcflurry.select\_processing\_models\_command module

Model select antigen processing models.

APPROACH: For each training fold, we select at least min and at most max models (where min and max are set by the `--min/max-models-per-fold` argument) using a step-up (forward) selection procedure. The final ensemble is the union of all selected models across all folds. AUC is used as the metric.

```
mhcflurry.select_processing_models_command.run(argv=['-b', 'latex', '-v', '-d',
                                                    '_build/doctrees', ':', '_build/latex'])
```

```
mhcflurry.select_processing_models_command.do_model_select_task(item, constant_data={})
```

```
mhcflurry.select_processing_models_command.model_select(fold_num, models,
                                                         min_models, max_models,
                                                         constant_data={})
```

Model select for a fold.

#### Parameters

**fold\_num** [int]

**models** [list of Class1NeuralNetwork]

**min\_models** [int]

**max\_models** [int]

**constant\_data** [dict]

**Returns**

dict with keys 'fold\_num', 'selected\_indices', 'summary'

## 5.31 mhcflurry.testing\_utils module

Utilities used in MHCflurry unit tests.

`mhcflurry.testing_utils.startup()`

Configure Keras backend for running unit tests.

`mhcflurry.testing_utils.cleanup()`

Clear tensorflow session and other process-wide resources.

## 5.32 mhcflurry.train\_allele\_specific\_models\_command module

Train ClassI single allele models.

`mhcflurry.train_allele_specific_models_command.run(argv=['-b', 'latex', '-v', '-d', '_build/doctrees', '._build/latex'])`

`mhcflurry.train_allele_specific_models_command.alleles_by_similarity(allele)`

`mhcflurry.train_allele_specific_models_command.train_model(n_models, allele_num, n_alleles, hyperparameter_set_num, num_hyperparameter_sets, allele, hyperparameters, verbose, progress_print_interval, predictor, save_to)`

`mhcflurry.train_allele_specific_models_command.subselect_df_held_out(df, reciprocal_held_out_fraction=10, seed=0)`

## 5.33 mhcflurry.train\_pan\_allele\_models\_command module

Train ClassI pan-allele models.

`mhcflurry.train_pan_allele_models_command.assign_folds(df, num_folds, held_out_fraction, held_out_max)`

Split training data into multiple test/train pairs, which we refer to as folds. Note that a given data point may be assigned to multiple test or train sets; these folds are NOT a non-overlapping partition as used in cross validation.

A fold is defined by a boolean value for each data point, indicating whether it is included in the training data for that fold. If it's not in the training data, then it's in the test data.

Folds are balanced in terms of allele content.

**Parameters**

**df** [pandas.DataFrame] training data

**num\_folds** [int]

**held\_out\_fraction** [float] Fraction of data to hold out as test data in each fold

**held\_out\_max** For a given allele, do not hold out more than held\_out\_max number of data points in any fold.

**Returns**

**pandas.DataFrame** index is same as df.index, columns are “fold\_0”, ... “fold\_N” giving whether the data point is in the training data for the fold

`mhcflurry.train_pan_allele_models_command.pretrain_data_iterator` (*filename*,  
*master\_allele\_encoding*,  
*peptides\_per\_chunk=1024*)

Step through a CSV file giving predictions for a large number of peptides (rows) and alleles (columns).

**Parameters**

**filename** [string]

**master\_allele\_encoding** [AlleleEncoding]

**peptides\_per\_chunk** [int]

**Returns**

**Generator of (AlleleEncoding, EncodableSequences, float affinities) tuples**

```
mhcflurry.train_pan_allele_models_command.run (argv=['-b', 'latex', '-v', '-d',
'_build/doctrees', '.', '_build/latex'])
mhcflurry.train_pan_allele_models_command.main (args)
mhcflurry.train_pan_allele_models_command.initialize_training (args)
mhcflurry.train_pan_allele_models_command.train_models (args)
mhcflurry.train_pan_allele_models_command.train_model (work_item_name,
work_item_num,
num_work_items, architecture_num, num_architectures,
fold_num, num_folds, replicate_num, num_replicates,
hyperparameters, pretrain_data_filename, verbose,
progress_print_interval, predictor, save_to, constant_data={})
```

## 5.34 mhcflurry.train\_presentation\_models\_command module

Train Class1 presentation models.

```
mhcflurry.train_presentation_models_command.run(argv=['-b', 'latex', '-v', '-d',
                                                    '_build/doctrees', '.', '_build/latex'])
mhcflurry.train_presentation_models_command.main(args)
```

## 5.35 mhcflurry.train\_processing\_models\_command module

Train Class1 processing models.

```
mhcflurry.train_processing_models_command.assign_folds(df, num_folds,
                                                    held_out_samples)
```

Split training data into multiple test/train pairs, which we refer to as folds. Note that a given data point may be assigned to multiple test or train sets; these folds are NOT a non-overlapping partition as used in cross validation.

A fold is defined by a boolean value for each data point, indicating whether it is included in the training data for that fold. If it's not in the training data, then it's in the test data.

### Parameters

**df** [pandas.DataFrame] training data

**num\_folds** [int]

**held\_out\_samples** [int]

### Returns

**pandas.DataFrame** index is same as df.index, columns are "fold\_0", ... "fold\_N" giving whether the data point is in the training data for the fold

```
mhcflurry.train_processing_models_command.run(argv=['-b', 'latex', '-v', '-d',
                                                    '_build/doctrees', '.', '_build/latex'])
```

```
mhcflurry.train_processing_models_command.main(args)
```

```
mhcflurry.train_processing_models_command.initialize_training(args)
```

```
mhcflurry.train_processing_models_command.train_models(args)
```

```
mhcflurry.train_processing_models_command.train_model(work_item_name,
                                                    work_item_num,
                                                    num_work_items, architecture_num, num_architectures,
                                                    fold_num, num_folds, replicate_num, num_replicates,
                                                    hyperparameters, verbose,
                                                    progress_print_interval,
                                                    predictor, save_to, constant_data={})
```



## 5.36 mhcflurry.version module



## PYTHON MODULE INDEX

### m

`mhcflurry`, 27

`mhcflurry.allele_encoding`, 48

`mhcflurry.amino_acid`, 49

`mhcflurry.calibrate_percentile_ranks_command`, 50

`mhcflurry.class1_affinity_predictor`, 51

`mhcflurry.class1_neural_network`, 57

`mhcflurry.class1_presentation_predictor`, 63

`mhcflurry.class1_processing_neural_network`, 68

`mhcflurry.class1_processing_predictor`, 70

`mhcflurry.cluster_parallelism`, 73

`mhcflurry.common`, 74

`mhcflurry.custom_loss`, 77

`mhcflurry.data_dependent_weights_initialization`, 79

`mhcflurry.downloads`, 79

`mhcflurry.downloads_command`, 81

`mhcflurry.encodable_sequences`, 83

`mhcflurry.ensemble_centrality`, 86

`mhcflurry.fasta`, 86

`mhcflurry.flanking_encoding`, 86

`mhcflurry.hyperparameters`, 88

`mhcflurry.local_parallelism`, 88

`mhcflurry.percent_rank_transform`, 90

`mhcflurry.predict_command`, 91

`mhcflurry.predict_scan_command`, 92

`mhcflurry.random_negative_peptides`, 92

`mhcflurry.regression_target`, 94

`mhcflurry.scoring`, 94

`mhcflurry.select_allele_specific_models_command`, 95

`mhcflurry.select_pan_allele_models_command`, 96

`mhcflurry.select_processing_models_command`, 97

`mhcflurry.testing_utils`, 98

`mhcflurry.train_allele_specific_models_command`, 98

`mhcflurry.train_pan_allele_models_command`, 98

`mhcflurry.train_presentation_models_command`, 100

`mhcflurry.train_processing_models_command`, 100

`mhcflurry.version`, 101



## INDEX

### Symbols

--additional-complete-file  
    <additional\_complete\_file>  
    mhcflurry-class1-select-pan-allele-models command line option, 23  
    mhcflurry-class1-select-processing-models command line option, 25  
    mhcflurry-class1-train-pan-allele-models command line option, 22  
    mhcflurry-class1-train-processing-models command line option, 24

--affinity-only  
    mhcflurry-predict command line option, 14

--affinity-predictor <dir>  
    mhcflurry-class1-train-presentation-models command line option, 26

--allele <allele>  
    mhcflurry-class1-select-allele-specific-models command line option, 20  
    mhcflurry-class1-train-allele-specific-models command line option, 18

--allele-column <name>  
    mhcflurry-predict command line option, 14

--allele-sequences <file.csv>  
    mhcflurry-class1-train-allele-specific-models command line option, 19  
    mhcflurry-class1-train-pan-allele-models command line option, 21

--alleles <allele>  
    mhcflurry-predict command line option, 14  
    mhcflurry-predict-scan command line option, 16

--already-downloaded-dir <dir>  
    mhcflurry-downloads-fetch command line option, 17

--always-include-best-allele  
    mhcflurry-predict command line option, 14

--backend {tensorflow-gpu,tensorflow-cpu,tensorflow-default,  
    mhcflurry-class1-select-allele-specific-models command line option, 20  
    mhcflurry-class1-select-pan-allele-models command line option, 23  
    mhcflurry-class1-select-processing-models command line option, 25  
    mhcflurry-class1-train-allele-specific-models command line option, 19  
    mhcflurry-class1-train-pan-allele-models command line option, 21  
    mhcflurry-class1-train-processing-models command line option, 24

--c-flank-column <name>  
    mhcflurry-predict command line option, 14

--cluster-max-retries  
    <cluster\_max\_retries>  
    mhcflurry-class1-select-pan-allele-models command line option, 23  
    mhcflurry-class1-select-processing-models command line option, 25  
    mhcflurry-class1-train-pan-allele-models command line option, 22  
    mhcflurry-class1-train-processing-models command line option, 24

--cluster-parallelism  
    mhcflurry-class1-select-pan-allele-models command line option, 23  
    mhcflurry-class1-select-processing-models command line option, 25  
    mhcflurry-class1-train-pan-allele-models command line option, 22  
    mhcflurry-class1-train-processing-models command line option, 24

--cluster-results-workdir  
    <cluster\_results\_workdir>  
    mhcflurry-class1-select-pan-allele-models command line option, 23  
    mhcflurry-class1-select-processing-models command line option, 25  
    mhcflurry-class1-train-pan-allele-models command line option, 22

```

    mhcflurry-class1-train-processing-models mhcflurry-class1-train-allele-specific-models
        command line option, 24                                command line option, 18
--cluster-script-prefix-path                mhcflurry-class1-train-pan-allele-models
    <cluster_script_prefix_path>            command line option, 21
    mhcflurry-class1-select-pan-allele-models mhcflurry-class1-train-presentation-models
        command line option, 23                                command line option, 26
    mhcflurry-class1-select-processing-models mhcflurry-class1-train-processing-models
        command line option, 25                                command line option, 23
    mhcflurry-class1-train-pan-allele-models mhcflurry-class1-train-pan-allele-models
        command line option, 22                                command line option, 21
    mhcflurry-class1-train-processing-models mhcflurry-class1-train-presentation-models
        command line option, 24                                command line option, 26
--cluster-submit-command                    mhcflurry-class1-train-processing-models
    <cluster_submit_command>                command line option, 24
    mhcflurry-class1-select-pan-allele-models mhcflurry-class1-select-allele-specific-models
        command line option, 23                                command line option, 19
    mhcflurry-class1-select-processing-models mhcflurry-class1-select-allele-specific-models
        command line option, 25                                command line option, 19
    mhcflurry-class1-train-pan-allele-models mhcflurry-class1-select-allele-specific-models
        command line option, 22                                command line option, 20
    mhcflurry-class1-train-processing-models mhcflurry-class1-select-pan-allele-models
        command line option, 24                                command line option, 23
--combined-max-models <n>                  mhcflurry-class1-select-processing-models
    mhcflurry-class1-select-allele-specific-models command line option, 25
        command line option, 20
--combined-min-contribution-percent         mhcflurry-class1-train-allele-specific-models
    <x>                                     command line option, 19
    mhcflurry-class1-select-allele-specific-models mhcflurry-class1-train-pan-allele-models
        command line option, 20                                command line option, 22
--combined-min-models <n>                  mhcflurry-class1-train-processing-models
    mhcflurry-class1-select-allele-specific-models command line option, 24
        command line option, 20
--consensus-max-models <n>                  mhcflurry-class1-train-allele-specific-models
    mhcflurry-class1-select-allele-specific-models command line option, 18
        command line option, 20
--consensus-min-models <n>                  mhcflurry-class1-train-allele-specific-models
    mhcflurry-class1-select-allele-specific-models command line option, 18
        command line option, 20
--consensus-num-peptides-per-length        <x>
    <consensus_num_peptides_per_length> mhcflurry-class1-train-pan-allele-models
        command line option, 21
    mhcflurry-class1-select-allele-specific-models mhcflurry-class1-select-pan-allele-models
        command line option, 20                                command line option, 22
--continue-incomplete                      mhcflurry-class1-train-processing-models
    mhcflurry-class1-train-pan-allele-models command line option, 23
        command line option, 21
    mhcflurry-class1-train-processing-models mhcflurry-class1-select-allele-specific-models
        command line option, 24                                command line option, 19
--data <file.csv>                          mhcflurry-class1-select-pan-allele-models
    mhcflurry-class1-select-allele-specific-models command line option, 22
        command line option, 19
    mhcflurry-class1-select-pan-allele-models mhcflurry-class1-select-processing-models
        command line option, 25                                command line option, 25
    mhcflurry-class1-select-pan-allele-models mhcflurry-class1-train-allele-specific-models
        command line option, 22                                command line option, 18
    mhcflurry-class1-select-processing-models mhcflurry-class1-train-pan-allele-models
        command line option, 25                                command line option, 25

```

command line option, [21](#)  
 mhcflurry-class1-train-presentation-models command line option, [20](#)  
 command line option, [26](#)  
 mhcflurry-class1-train-processing-models mhcflurry-class1-select-allele-specific-models  
 command line option, [23](#)  
 mhcflurry-downloads command line  
 option, [17](#)  
 mhcflurry-downloads-fetch command  
 line option, [17](#)  
 mhcflurry-downloads-info command  
 line option, [18](#)  
 mhcflurry-downloads-path command  
 line option, [18](#)  
 mhcflurry-downloads-url command  
 line option, [18](#)  
 mhcflurry-predict command line  
 option, [14](#)  
 mhcflurry-predict-scan command  
 line option, [15](#)  
 --hla-column <hla\_column>  
 mhcflurry-class1-train-presentation-models  
 command line option, [26](#)  
 --hyperparameters <file.json>  
 mhcflurry-class1-train-allele-specific-models  
 command line option, [18](#)  
 mhcflurry-class1-train-pan-allele-models  
 command line option, [21](#)  
 mhcflurry-class1-train-processing-models  
 command line option, [23](#)  
 --ignore-inequalities  
 mhcflurry-class1-train-allele-specific-models  
 command line option, [19](#)  
 mhcflurry-class1-train-pan-allele-models  
 command line option, [21](#)  
 --input-format {guess,csv,fasta}  
 mhcflurry-predict-scan command  
 line option, [16](#)  
 --keep  
 mhcflurry-downloads-fetch command  
 line option, [17](#)  
 --list-supported-alleles  
 mhcflurry-predict command line  
 option, [14](#)  
 mhcflurry-predict-scan command  
 line option, [15](#)  
 --list-supported-peptide-lengths  
 mhcflurry-predict command line  
 option, [14](#)  
 mhcflurry-predict-scan command  
 line option, [16](#)  
 --mass-spec-max-models <n>  
 mhcflurry-class1-select-allele-specific-models  
 command line option, [20](#)  
 --mass-spec-min-measurements <n>  
 mhcflurry-class1-select-allele-specific-models  
 command line option, [20](#)  
 --mass-spec-min-models <n>  
 mhcflurry-class1-select-allele-specific-models  
 command line option, [20](#)  
 --mass-spec-regex <regex>  
 mhcflurry-class1-select-allele-specific-models  
 command line option, [20](#)  
 mhcflurry-class1-select-pan-allele-models  
 command line option, [22](#)  
 --max-epochs <n>  
 mhcflurry-class1-train-allele-specific-models  
 command line option, [19](#)  
 mhcflurry-class1-train-pan-allele-models  
 command line option, [21](#)  
 mhcflurry-class1-train-processing-models  
 command line option, [24](#)  
 --max-models-per-fold <n>  
 mhcflurry-class1-select-pan-allele-models  
 command line option, [22](#)  
 mhcflurry-class1-select-processing-models  
 command line option, [25](#)  
 --max-tasks-per-worker <n>  
 mhcflurry-class1-select-allele-specific-models  
 command line option, [20](#)  
 mhcflurry-class1-select-pan-allele-models  
 command line option, [23](#)  
 mhcflurry-class1-select-processing-models  
 command line option, [25](#)  
 mhcflurry-class1-train-allele-specific-models  
 command line option, [19](#)  
 mhcflurry-class1-train-pan-allele-models  
 command line option, [22](#)  
 mhcflurry-class1-train-processing-models  
 command line option, [24](#)  
 --max-workers-per-gpu <n>  
 mhcflurry-class1-select-allele-specific-models  
 command line option, [20](#)  
 mhcflurry-class1-select-pan-allele-models  
 command line option, [23](#)  
 mhcflurry-class1-select-processing-models  
 command line option, [25](#)  
 mhcflurry-class1-train-allele-specific-models  
 command line option, [19](#)  
 mhcflurry-class1-train-pan-allele-models  
 command line option, [22](#)  
 mhcflurry-class1-train-processing-models  
 command line option, [24](#)  
 --min-measurements-per-allele <n>  
 mhcflurry-class1-train-allele-specific-models  
 command line option, [18](#)  
 --min-models-per-fold <n>  
 mhcflurry-class1-select-pan-allele-models  
 command line option, [22](#)

```

mhcflurry-class1-select-processing-models mhcflurry-class1-select-processing-models
    command line option, 25
--models <dir>
    mhcflurry-predict command line
        option, 14
    mhcflurry-predict-scan command
        line option, 16
--models-dir <dir>
    mhcflurry-class1-select-allele-specific-models
        command line option, 19
    mhcflurry-class1-select-pan-allele-models
        command line option, 22
    mhcflurry-class1-select-processing-models
        command line option, 23
--mse-max-models <n>
    mhcflurry-class1-select-allele-specific-models
        command line option, 20
--mse-min-measurements <n>
    mhcflurry-class1-select-allele-specific-models
        command line option, 20
--mse-min-models <n>
    mhcflurry-class1-select-allele-specific-models
        command line option, 20
--n-flank-column <name>
    mhcflurry-predict command line
        option, 14
--n-models <n>
    mhcflurry-class1-train-allele-specific-models
        command line option, 19
--no-affinity-percentile
    mhcflurry-predict command line
        option, 14
    mhcflurry-predict-scan command
        line option, 16
--no-flanking
    mhcflurry-predict command line
        option, 14
    mhcflurry-predict-scan command
        line option, 16
--no-throw
    mhcflurry-predict command line
        option, 14
    mhcflurry-predict-scan command
        line option, 16
--num-folds <n>
    mhcflurry-class1-train-pan-allele-models
        command line option, 21
    mhcflurry-class1-train-processing-models
        command line option, 23
--num-jobs <n>
    mhcflurry-class1-select-allele-specific-models
        command line option, 20
    mhcflurry-class1-select-pan-allele-models
        command line option, 23
    mhcflurry-class1-select-processing-models
        command line option, 25
    mhcflurry-class1-train-allele-specific-models
        command line option, 19
    mhcflurry-class1-train-pan-allele-models
        command line option, 21
    mhcflurry-class1-train-processing-models
        command line option, 24
    --only-initialize
        mhcflurry-class1-train-pan-allele-models
            command line option, 21
        mhcflurry-class1-train-processing-models
            command line option, 24
    --out-delimiter <char>
        mhcflurry-predict command line
            option, 14
        mhcflurry-predict-scan command
            line option, 16
    --out-models-dir <dir>
        mhcflurry-class1-select-allele-specific-models
            command line option, 19
        mhcflurry-class1-select-pan-allele-models
            command line option, 22
        mhcflurry-class1-select-processing-models
            command line option, 25
        mhcflurry-class1-train-allele-specific-models
            command line option, 18
        mhcflurry-class1-train-pan-allele-models
            command line option, 21
        mhcflurry-class1-train-presentation-models
            command line option, 26
        mhcflurry-class1-train-processing-models
            command line option, 23
    --out-unselected-predictions
        <file.csv>
        mhcflurry-class1-select-allele-specific-models
            command line option, 19
    --output-delimiter <char>
        mhcflurry-predict command line
            option, 14
        mhcflurry-predict-scan command
            line option, 16
    --peptide-column <name>
        mhcflurry-predict command line
            option, 14
    --peptide-lengths <peptide_lengths>
        mhcflurry-predict-scan command
            line option, 16
    --peptides <peptide>

```



```

    mhcflurry-predict command line
        option, 14
--prediction-column-prefix <name>
    mhcflurry-predict command line
        option, 14
--pretrain-data <file.csv>
    mhcflurry-class1-train-pan-allele-models mhcflurry-predict-scan command
        command line option, 21
--processing-predictor-with-flanks
    <dir>
    mhcflurry-class1-train-presentation-models mhcflurry-predict-scan command
        command line option, 26
--processing-predictor-without-flanks
    <dir>
    mhcflurry-class1-train-presentation-models mhcflurry-class1-select-allele-specific-models
        command line option, 26
--quiet
    mhcflurry-downloads command line
        option, 17
--release <release>
    mhcflurry-downloads-fetch command
        line option, 17
--results-all
    mhcflurry-predict-scan command
        line option, 16
--results-best {presentation_score, processing_score, affinity, affinity_percentile}
    mhcflurry-predict-scan command
        line option, 16
--results-filtered
    {presentation_score, processing_score, affinity, affinity_percentile}
    mhcflurry-predict-scan command
        line option, 16
--save-interval <n>
    mhcflurry-class1-train-allele-specific-models mhcflurry-class1-select-allele-specific-models
        command line option, 19
--scoring <scoring>
    mhcflurry-class1-select-allele-specific-models mhcflurry-class1-train-pan-allele-models
        command line option, 20
--sequence-column <name>
    mhcflurry-predict-scan command
        line option, 16
--sequence-id-column <name>
    mhcflurry-predict-scan command
        line option, 16
--sequences <seq>
    mhcflurry-predict-scan command
        line option, 16
--target-column <target_column>
    mhcflurry-class1-train-presentation-models mhcflurry-class1-select-allele-specific-models
        command line option, 26
--threshold-affinity
    <threshold_affinity>
    mhcflurry-predict-scan command
        line option, 16
--threshold-affinity-percentile
    <threshold_affinity_percentile>
    mhcflurry-predict-scan command
        line option, 16
--threshold-presentation-score
    <threshold_presentation_score>
    mhcflurry-predict-scan command
        line option, 16
--threshold-processing-score
    <threshold_processing_score>
    mhcflurry-predict-scan command
        line option, 16
--unselected-accuracy-percentile-threshold
    <x>
    mhcflurry-class1-select-allele-specific-models
        command line option, 20
--unselected-accuracy-scorer <scorer>
    mhcflurry-class1-select-allele-specific-models
        command line option, 19
--unselected-accuracy-scorer-num-samples
    <unselected_accuracy_scorer_num_samples>
    mhcflurry-class1-select-allele-specific-models
        command line option, 20
--verbose
    mhcflurry-downloads command line
        option, 17
--verbosity <verbosity>
    mhcflurry-class1-select-allele-specific-models
        command line option, 20
--threshold-affinity-percentile
    <threshold_affinity_percentile>
    mhcflurry-class1-select-allele-specific-models
        command line option, 23
--threshold-presentation-score
    <threshold_presentation_score>
    mhcflurry-class1-select-processing-models
        command line option, 25
--threshold-processing-score
    <threshold_processing_score>
    mhcflurry-class1-train-allele-specific-models
        command line option, 19
--unselected-accuracy-percentile-threshold
    <x>
    mhcflurry-class1-train-pan-allele-models
        command line option, 21
--unselected-accuracy-scorer <scorer>
    mhcflurry-class1-train-presentation-models
        command line option, 26
--unselected-accuracy-scorer-num-samples
    <unselected_accuracy_scorer_num_samples>
    mhcflurry-class1-train-processing-models
        command line option, 24
--version
    mhcflurry-predict command line
        option, 14
    mhcflurry-predict-scan command
        line option, 16
--worker-log-dir <worker_log_dir>
    mhcflurry-class1-select-allele-specific-models
        command line option, 21
    mhcflurry-class1-select-pan-allele-models
        command line option, 23
    mhcflurry-class1-select-processing-models
        command line option, 25
    mhcflurry-class1-train-allele-specific-models

```

command line option, 19

mhcflurry-class1-train-pan-allele-models command line option, 22

mhcflurry-class1-train-processing-models command line option, 24

-h

mhcflurry-class1-select-allele-specific-models command line option, 19

mhcflurry-class1-select-pan-allele-models command line option, 22

mhcflurry-class1-select-processing-models command line option, 25

mhcflurry-class1-train-allele-specific-models command line option, 18

mhcflurry-class1-train-pan-allele-models command line option, 21

mhcflurry-class1-train-presentation-models command line option, 26

mhcflurry-class1-train-processing-models command line option, 23

mhcflurry-downloads command line option, 17

mhcflurry-downloads-fetch command line option, 17

mhcflurry-downloads-info command line option, 18

mhcflurry-downloads-path command line option, 18

mhcflurry-downloads-url command line option, 18

mhcflurry-predict command line option, 14

mhcflurry-predict-scan command line option, 15

-v

mhcflurry-downloads command line option, 17

**A**

add\_cluster\_parallelism\_args() (in module *mhcflurry.cluster\_parallelism*), 73

add\_local\_parallelism\_args() (in module *mhcflurry.local\_parallelism*), 88

add\_models() (*mhcflurry.class1\_processing\_predictor.Class1ProcessingPredictor* method), 71

add\_models() (*mhcflurry.Class1ProcessingPredictor* method), 39

add\_pan\_allele\_model() (*mhcflurry.class1\_affinity\_predictor.Class1AffinityPredictor* method), 55

add\_pan\_allele\_model() (*mhcflurry.Class1AffinityPredictor* method), 31

allele\_encoding\_to\_network\_input() (*mhcflurry.class1\_neural\_network.Class1NeuralNetwork* method), 60

allele\_encoding\_to\_network\_input() (*mhcflurry.Class1NeuralNetwork* method), 36

allele\_representations() (*mhcflurry.allele\_encoding.AleleEncoding* method), 49

AlleleEncoding (class in *mhcflurry.allele\_encoding*), 48

alleles\_by\_similarity() (in module *mhcflurry.train\_allele\_specific\_models\_command*), 98

amino\_acid\_distribution() (in module *mhcflurry.common*), 75

apply\_hyperparameter\_renames() (*mhcflurry.class1\_neural\_network.Class1NeuralNetwork* class method), 58

apply\_hyperparameter\_renames() (*mhcflurry.Class1NeuralNetwork* class method), 34

array() (*mhcflurry.flanking\_encoding.EncodingResult* property), 86

assign\_folds() (in module *mhcflurry.train\_pan\_allele\_models\_command*), 98

assign\_folds() (in module *mhcflurry.train\_processing\_models\_command*), 100

auxiliary\_input\_hyperparameter\_defaults (*mhcflurry.class1\_processing\_neural\_network.Class1ProcessingNeuralNetwork* attribute), 68

auxiliary\_input\_hyperparameter\_defaults (*mhcflurry.Class1ProcessingNeuralNetwork* attribute), 41

available\_vector\_encodings() (in module *mhcflurry.amino\_acid*), 49

**B**

borrow\_cached\_network() (*mhcflurry.class1\_neural\_network.Class1NeuralNetwork* class method), 58

borrow\_cached\_network() (*mhcflurry.Class1NeuralNetwork* class method), 34

**C**

cache\_encoding() (in module *mhcflurry.select\_allele\_specific\_models\_command*), 95

calibrate\_percentile\_ranks() (in module *mhcflurry.calibrate\_percentile\_ranks\_command*), 50

calibrate\_percentile\_ranks() (*mhcflurry.class1\_affinity\_predictor.Class1AffinityPredictor* method), 56

`calibrate_percentile_ranks()` (*mhcflurry.Class1AffinityPredictor* method), 32  
`call_wrapped()` (in module *mhcflurry.local\_parallelism*), 90  
`call_wrapped_kwargs()` (in module *mhcflurry.local\_parallelism*), 90  
`check_consistency()` (*mhcflurry.class1\_affinity\_predictor.Class1AffinityPredictor* method), 52  
`check_consistency()` (*mhcflurry.class1\_processing\_predictor.Class1ProcessingNeuralNetwork* method), 72  
`check_consistency()` (*mhcflurry.Class1AffinityPredictor* method), 28  
`check_consistency()` (*mhcflurry.Class1ProcessingPredictor* method), 40  
`check_shape()` (in module *mhcflurry.custom\_loss*), 78  
`check_valid_keys()` (*mhcflurry.hyperparameters.HyperparameterDefaults* method), 88  
`Class1AffinityPredictor` (class in *mhcflurry*), 27  
`Class1AffinityPredictor` (class in *mhcflurry.class1\_affinity\_predictor*), 51  
`Class1NeuralNetwork` (class in *mhcflurry*), 33  
`Class1NeuralNetwork` (class in *mhcflurry.class1\_neural\_network*), 57  
`Class1PresentationPredictor` (class in *mhcflurry*), 43  
`Class1PresentationPredictor` (class in *mhcflurry.class1\_presentation\_predictor*), 63  
`Class1ProcessingNeuralNetwork` (class in *mhcflurry*), 41  
`Class1ProcessingNeuralNetwork` (class in *mhcflurry.class1\_processing\_neural\_network*), 68  
`Class1ProcessingPredictor` (class in *mhcflurry*), 38  
`Class1ProcessingPredictor` (class in *mhcflurry.class1\_processing\_predictor*), 70  
`cleanup()` (in module *mhcflurry.testing\_utils*), 98  
`clear_allele_representations()` (*mhcflurry.class1\_neural\_network.Class1NeuralNetwork* method), 62  
`clear_allele_representations()` (*mhcflurry.Class1NeuralNetwork* method), 38  
`clear_cache()` (*mhcflurry.class1\_affinity\_predictor.Class1AffinityPredictor* method), 51  
`clear_cache()` (*mhcflurry.Class1AffinityPredictor* method), 27  
`clear_model_cache()` (*mhcflurry.class1\_neural\_network.Class1NeuralNetwork* class method), 58  
`clear_model_cache()` (*mhcflurry.Class1NeuralNetwork* class method), 34  
`cls` (in module *mhcflurry.custom\_loss*), 79  
`cluster_results()` (in module *mhcflurry.cluster\_parallelism*), 73  
`cluster_results_from_args()` (in module *mhcflurry.cluster\_parallelism*), 73  
`CombinedModelSelector` (class in *mhcflurry.select\_allele\_specific\_models\_command*), 95  
`compact()` (*mhcflurry.allele\_encoding.AlleleEncoding* method), 49  
`compile_hyperparameter_defaults` (*mhcflurry.class1\_neural\_network.Class1NeuralNetwork* attribute), 57  
`compile_hyperparameter_defaults` (*mhcflurry.class1\_processing\_neural\_network.Class1ProcessingNeuralNetwork* attribute), 68  
`compile_hyperparameter_defaults` (*mhcflurry.Class1NeuralNetwork* attribute), 34  
`compile_hyperparameter_defaults` (*mhcflurry.Class1ProcessingNeuralNetwork* attribute), 41  
`configure()` (in module *mhcflurry.downloads*), 81  
`configure_logging()` (in module *mhcflurry.common*), 74  
`ConsensusModelSelector` (class in *mhcflurry.select\_allele\_specific\_models\_command*), 95  
`create()` (*mhcflurry.encodable\_sequences.EncodableSequences* class method), 83

## D

`data_dependent_weights_initialization()` (*mhcflurry.class1\_neural\_network.Class1NeuralNetwork* static method), 60  
`data_dependent_weights_initialization()` (*mhcflurry.Class1NeuralNetwork* static method), 36  
`default()` (*mhcflurry.common.NumpyJSONEncoder* method), 76  
`do_calibrate_percentile_ranks()` (in module *mhcflurry.calibrate\_percentile\_ranks\_command*), 50  
`do_model_select_task()` (in module *mhcflurry.select\_pan\_allele\_models\_command*), 97  
`do_model_select_task()` (in module *mhcflurry.select\_processing\_models\_command*), 97

download  
     mhcflurry-downloads-fetch command  
         line option, 17  
download\_name  
     mhcflurry-downloads-path command  
         line option, 18  
     mhcflurry-downloads-url command  
         line option, 18

## E

early\_stopping\_hyperparameter\_defaults  
     (mhcflurry.class1\_neural\_network.Class1NeuralNetwork  
     attribute), 58  
early\_stopping\_hyperparameter\_defaults  
     (mhcflurry.class1\_processing\_neural\_network.Class1ProcessingNeuralNetwork  
     attribute), 68  
early\_stopping\_hyperparameter\_defaults  
     (mhcflurry.Class1NeuralNetwork attribute), 34  
early\_stopping\_hyperparameter\_defaults  
     (mhcflurry.Class1ProcessingNeuralNetwork  
     attribute), 41  
EncodableSequences (class in  
     mhcflurry.encodable\_sequences), 83  
encode() (mhcflurry.flanking\_encoding.FlankingEncoding  
     static method), 87  
encode\_y() (mhcflurry.custom\_loss.MSEWithInequalities  
     static method), 78  
encode\_y() (mhcflurry.custom\_loss.MSEWithInequalitiesAndMultiAllelicMassSpecLoss  
     static method), 78  
encode\_y() (mhcflurry.custom\_loss.MultiallelicMassSpecLoss  
     static method), 78  
encode\_y() (mhcflurry.custom\_loss.StandardKerasLoss  
     static method), 77  
encode\_y() (mhcflurry.custom\_loss.TransformPredictionsLossWrapper  
     method), 77  
EncodingError, 83  
EncodingResult (class in  
     mhcflurry.flanking\_encoding), 86  
extend() (mhcflurry.hyperparameters.HyperparameterDefaults  
     method), 88

## F

FastaParser (class in mhcflurry.fasta), 86  
fetch\_subcommand() (in  
     mhcflurry.downloads\_command), 83  
fit() (mhcflurry.class1\_neural\_network.Class1NeuralNetwork  
     method), 61  
fit() (mhcflurry.class1\_presentation\_predictor.Class1PresentationPredictor  
     method), 64  
fit() (mhcflurry.class1\_processing\_neural\_network.Class1ProcessingNeuralNetwork  
     method), 69  
fit() (mhcflurry.Class1NeuralNetwork method), 37  
fit() (mhcflurry.Class1PresentationPredictor method),  
     44  
fit() (mhcflurry.Class1ProcessingNeuralNetwork  
     method), 41  
fit() (mhcflurry.percent\_rank\_transform.PercentRankTransform  
     method), 90  
fit\_allele\_specific\_predictors()  
     (mhcflurry.class1\_affinity\_predictor.Class1AffinityPredictor  
     method), 54  
fit\_allele\_specific\_predictors()  
     (mhcflurry.Class1AffinityPredictor method), 30  
fit\_class1\_pan\_allele\_models()  
     (mhcflurry.class1\_affinity\_predictor.Class1AffinityPredictor  
     method), 54  
fit\_class1\_pan\_allele\_models()  
     (mhcflurry.Class1AffinityPredictor method), 30  
fit\_class1\_pan\_allele\_models()  
     (mhcflurry.class1\_neural\_network.Class1NeuralNetwork  
     method), 60  
fit\_generator() (mhcflurry.Class1NeuralNetwork  
     method), 36  
fit\_hyperparameter\_defaults  
     (mhcflurry.class1\_neural\_network.Class1NeuralNetwork  
     attribute), 58  
fit\_hyperparameter\_defaults  
     (mhcflurry.class1\_processing\_neural\_network.Class1ProcessingNeuralNetwork  
     attribute), 68  
fit\_hyperparameter\_defaults  
     (mhcflurry.Class1NeuralNetwork attribute), 34  
fit\_hyperparameter\_defaults  
     (mhcflurry.Class1ProcessingNeuralNetwork  
     attribute), 41  
fit\_length\_vector\_encoded\_sequences()  
     (mhcflurry.allele\_encoding.AlleleEncoding  
     method), 49  
fixed\_vectors\_encoding() (in  
     module  
     mhcflurry.amino\_acid), 50  
FlankingEncoding (class in  
     mhcflurry.flanking\_encoding), 86  
from\_config() (mhcflurry.class1\_neural\_network.Class1NeuralNetwork  
     class method), 59  
from\_config() (mhcflurry.class1\_processing\_neural\_network.Class1ProcessingNeuralNetwork  
     class method), 70  
from\_config() (mhcflurry.Class1NeuralNetwork  
     class method), 35  
from\_config() (mhcflurry.Class1ProcessingNeuralNetwork  
     class method), 43  
from\_ic50() (in module mhcflurry.regression\_target),  
     94  
from\_series() (mhcflurry.percent\_rank\_transform.PercentRankTransform  
     static method), 91

## G

get\_activations() (in  
     module  
     mhcflurry.data\_dependent\_weights\_initialization),  
     79

[get\\_alleles\(\)](#) ([mhcflurry.random\\_negative\\_peptides.RandomNegativePeptides](#) defaults method), 93  
[get\\_config\(\)](#) ([mhcflurry.class1\\_neural\\_network.Class1NeuralNetwork](#) attribute), 34  
[get\\_config\(\)](#) ([mhcflurry.class1\\_neural\\_network.Class1NeuralNetwork](#) hyperparameter\_defaults method), 59  
[get\\_config\(\)](#) ([mhcflurry.class1\\_processing\\_neural\\_network.Class1ProcessingNeuralNetwork](#) hyperparameter\_defaults method), 70  
[get\\_config\(\)](#) ([mhcflurry.Class1NeuralNetwork](#) attribute), 35  
[get\\_config\(\)](#) ([mhcflurry.Class1ProcessingNeuralNetwork](#) hyperparameter\_renames method), 43  
[get\\_current\\_release\(\)](#) (in module [mhcflurry.downloads](#)), 79  
[get\\_current\\_release\\_downloads\(\)](#) (in module [mhcflurry.downloads](#)), 80  
[get\\_default\\_class1\\_models\\_dir\(\)](#) (in module [mhcflurry.downloads](#)), 79  
[get\\_default\\_class1\\_presentation\\_models\\_dir\(\)](#) (in module [mhcflurry.downloads](#)), 80  
[get\\_default\\_class1\\_processing\\_models\\_dir\(\)](#) (in module [mhcflurry.downloads](#)), 80  
[get\\_downloads\\_dir\(\)](#) (in module [mhcflurry.downloads](#)), 79  
[get\\_downloads\\_metadata\(\)](#) (in module [mhcflurry.downloads](#)), 79  
[get\\_keras\\_loss\(\)](#) ([mhcflurry.custom\\_loss.Loss](#) method), 77  
[get\\_loss\(\)](#) (in module [mhcflurry.custom\\_loss](#)), 77  
[get\\_model\(\)](#) ([mhcflurry.class1\\_presentation\\_predictor.Class1PresentationPredictor](#) method), 65  
[get\\_model\(\)](#) ([mhcflurry.Class1PresentationPredictor](#) method), 45  
[get\\_path\(\)](#) (in module [mhcflurry.downloads](#)), 80  
[get\\_peptides\(\)](#) ([mhcflurry.random\\_negative\\_peptides.RandomNegativePeptides](#) method), 93  
[get\\_total\\_count\(\)](#) ([mhcflurry.random\\_negative\\_peptides.RandomNegativePeptides](#) method), 94  
[get\\_weights\(\)](#) ([mhcflurry.class1\\_neural\\_network.Class1NeuralNetwork](#) method), 59  
[get\\_weights\(\)](#) ([mhcflurry.class1\\_processing\\_neural\\_network.Class1ProcessingNeuralNetwork](#) method), 70  
[get\\_weights\(\)](#) ([mhcflurry.Class1NeuralNetwork](#) method), 35  
[get\\_weights\(\)](#) ([mhcflurry.Class1ProcessingNeuralNetwork](#) method), 42

## H

[hyperparameter\\_defaults](#) ([mhcflurry.class1\\_neural\\_network.Class1NeuralNetwork](#) attribute), 58

## L

[load\(\)](#) ([mhcflurry.class1\\_affinity\\_predictor.Class1AffinityPredictor](#) static method), 59  
[load\(\)](#) ([mhcflurry.class1\\_presentation\\_predictor.Class1PresentationPredictor](#) class method), 68



`load()` (`mhcflurry.class1_processing_predictor.Class1ProcessingPredictor` `mhcflurry.select_allele_specific_models_command`),  
 class method), 72 96  
`load()` (`mhcflurry.Class1AffinityPredictor` static `master_allele_encoding()`  
 method), 29 (`mhcflurry.class1_affinity_predictor.Class1AffinityPredictor`  
`load()` (`mhcflurry.Class1PresentationPredictor` class `property`), 53  
 method), 48 `master_allele_encoding()`  
`load()` (`mhcflurry.Class1ProcessingPredictor` class `mhcflurry.Class1AffinityPredictor` `property`),  
 method), 40 30  
`load_weights()` (in module `mhcflurry.common`), 75 `max_absolute_value()`  
`load_weights()` (`mhcflurry.class1_neural_network.Class1NeuralNetwork` `mhcflurry.select_allele_specific_models_command.ConsensusModelSelector`  
 method), 59 method), 95  
`load_weights()` (`mhcflurry.Class1NeuralNetwork` `max_absolute_value()`  
 method), 35 (`mhcflurry.select_allele_specific_models_command.MassSpecModelSelector`  
 method), 96  
`Loss` (class in `mhcflurry.custom_loss`), 77 `max_absolute_value()`  
`loss()` (`mhcflurry.custom_loss.Loss` method), 77 `max_absolute_value()`  
`loss()` (`mhcflurry.custom_loss.MSEWithInequalities` `mhcflurry.select_allele_specific_models_command.MSEModelSelector`  
 method), 78 method), 96  
`loss()` (`mhcflurry.custom_loss.MSEWithInequalitiesAndMultipleQueries` `mhcflurry.class1_affinity_predictor.Class1AffinityPredictor`  
 method), 78 class method), 52  
`loss()` (`mhcflurry.custom_loss.MultiallelicMassSpecLoss` `merge()` (`mhcflurry.class1_neural_network.Class1NeuralNetwork`  
 method), 78 class method), 62  
`loss()` (`mhcflurry.custom_loss.TransformPredictionsLossWrapper` `mhcflurry.Class1AffinityPredictor` class  
 method), 77 method), 28  
`lsuv_init()` (in module `merge()` (`mhcflurry.Class1NeuralNetwork` class  
`mhcflurry.data_dependent_weights_initialization`), method), 38  
 79 `merge_in_place()` (`mhcflurry.class1_affinity_predictor.Class1AffinityPredictor`  
 method), 52  
`merge_in_place()` (`mhcflurry.Class1AffinityPredictor`  
 method), 52  
**M**  
`main()` (in module `mhcflurry.train_pan_allele_models_command`), method), 28  
 99 `mhcflurry` (module), 27  
`main()` (in module `mhcflurry.train_presentation_models_command`), `mhcflurry.allele_encoding` (module), 48  
 100 `mhcflurry.amino_acid` (module), 49  
`main()` (in module `mhcflurry.train_processing_models_command`), `mhcflurry.calibrate_percentile_ranks_command`  
 100 (module), 50  
`make_network()` (`mhcflurry.class1_neural_network.Class1NeuralNetwork` `mhcflurry.class1_affinity_predictor`  
 method), 62 (module), 51  
`make_network()` (`mhcflurry.class1_processing_neural_network.Class1ProcessingNeuralNetwork` (module),  
 method), 70 57  
`make_network()` (`mhcflurry.Class1NeuralNetwork` `mhcflurry.class1_presentation_predictor`  
 method), 38 (module), 63  
`make_network()` (`mhcflurry.Class1ProcessingNeuralNetwork` `mhcflurry.class1_processing_neural_network`  
 method), 42 (module), 68  
`make_scores()` (in module `mhcflurry.scoring`), 94 `mhcflurry.class1_processing_predictor`  
`make_worker_pool()` (in module `mhcflurry` (module), 70  
`mhcflurry.local_parallelism`), 89 `mhcflurry.cluster_parallelism` (module), 73  
`manifest_df()` (`mhcflurry.class1_affinity_predictor.Class1AffinityPredictor` `mhcflurry.common` (module), 74  
 property), 51 `mhcflurry.custom_loss` (module), 77  
`manifest_df()` (`mhcflurry.class1_processing_predictor.Class1ProcessingPredictor` `mhcflurry.data_dependent_weights_initialization`  
 property), 71 (module), 79  
`manifest_df()` (`mhcflurry.Class1AffinityPredictor` `mhcflurry.downloads` (module), 79  
 property), 27 `mhcflurry.downloads_command` (module), 81  
`manifest_df()` (`mhcflurry.Class1ProcessingPredictor` `mhcflurry.encodable_sequences` (module), 83  
 property), 39 `mhcflurry.ensemble_centrality` (module), 86  
`MassSpecModelSelector` (class in `mhcflurry.fasta` (module), 86

---

```

mhcflurry.flanking_encoding (module), 86
mhcflurry.hyperparameters (module), 88
mhcflurry.local_parallelism (module), 88
mhcflurry.percent_rank_transform (module), 90
mhcflurry.predict_command (module), 91
mhcflurry.predict_scan_command (module), 92
mhcflurry.random_negative_peptides (module), 92
mhcflurry.regression_target (module), 94
mhcflurry.scoring (module), 94
mhcflurry.select_allele_specific_models_command (module), 95
mhcflurry.select_pan_allele_models_command (module), 96
mhcflurry.select_processing_models_command (module), 97
mhcflurry.testing_utils (module), 98
mhcflurry.train_allele_specific_models_command (module), 98
mhcflurry.train_pan_allele_models_command (module), 98
mhcflurry.train_presentation_models_command (module), 100
mhcflurry.train_processing_models_command (module), 100
mhcflurry.version (module), 101
mhcflurry-class1-select-allele-specific-models <cluster_results_workdir>, 23
  command line option
  --allele <allele>, 20
  --backend {tensorflow-gpu,tensorflow-cpu,tensorflow-default}, 23
  --combined-max-models <n>, 20
  --combined-min-contribution-percent <x>, 20
  --combined-min-models <n>, 20
  --consensus-max-models <n>, 20
  --consensus-min-models <n>, 20
  --consensus-num-peptides-per-length <consensus_num_peptides_per_length>, 20
  --data <file.csv>, 19
  --exclude-data <file.csv>, 19
  --gpus <n>, 20
  --help, 19
  --mass-spec-max-models <n>, 20
  --mass-spec-min-measurements <n>, 20
  --mass-spec-min-models <n>, 20
  --mass-spec-regex <regex>, 20
  --max-tasks-per-worker <n>, 20
  --max-workers-per-gpu <n>, 20
  --models-dir <dir>, 19
  --mse-max-models <n>, 20
  --mse-min-measurements <n>, 20
  --mse-min-models <n>, 20
  --num-jobs <n>, 20
  --out-models-dir <dir>, 19
  --out-unselected-predictions <file.csv>, 19
  --scoring <scoring>, 20
  --unselected-accuracy-percentile-threshold <x>, 20
  --unselected-accuracy-scorer <scorer>, 19
  --unselected-accuracy-scorer-num-samples <unselected_accuracy_scorer_num_samples>, 20
  --verbosity <verbosity>, 20
  --worker-log-dir <worker_log_dir>, 21
  -h, 19
mhcflurry-class1-select-pan-allele-models command line option
  --additional-complete-file <additional_complete_file>, 23
  --backend {tensorflow-gpu,tensorflow-cpu,tensorflow-default}, 23
  --cluster-max-retries <cluster_max_retries>, 23
  --cluster-parallelism, 23
  --cluster-results-workdir <cluster_results_workdir>, 23
  --cluster-script-prefix-path <cluster_script_prefix_path>, 23
  --cluster-submit-command <cluster_submit_command>, 23
  --data <file.csv>, 22
  --gpus <n>, 23
  --help, 22
  --mass-spec-regex <regex>, 22
  --max-models-per-fold <n>, 22
  --max-tasks-per-worker <n>, 23
  --max-workers-per-gpu <n>, 23
  --min-models-per-fold <n>, 22
  --models-dir <dir>, 22
  --num-jobs <n>, 23
  --out-models-dir <dir>, 22
  --verbosity <verbosity>, 23
  --worker-log-dir <worker_log_dir>, 23
  -h, 22
mhcflurry-class1-select-processing-models command line option
  --additional-complete-file <additional_complete_file>, 25
  --backend {tensorflow-gpu,tensorflow-cpu,tensorflow-default}, 25

```

```
25
--cluster-max-retries
    <cluster_max_retries>, 25
--cluster-parallelism, 25
--cluster-results-workdir
    <cluster_results_workdir>, 25
--cluster-script-prefix-path
    <cluster_script_prefix_path>,
    25
--cluster-submit-command
    <cluster_submit_command>, 25
--data <file.csv>, 25
--gpus <n>, 25
--help, 25
--max-models-per-fold <n>, 25
--max-tasks-per-worker <n>, 25
--max-workers-per-gpu <n>, 25
--min-models-per-fold <n>, 25
--models-dir <dir>, 25
--num-jobs <n>, 25
--out-models-dir <dir>, 25
--verbosity <verbosity>, 25
--worker-log-dir <worker_log_dir>,
    25
-h, 25
mhcflurry-class1-train-allele-specific-models
    command line option
--allele <allele>, 18
--allele-sequences <file.csv>, 19
--backend {tensorflow-gpu,tensorflow-cpu,tensorflow-gpu,
    19
--data <file.csv>, 18
--gpus <n>, 19
--held-out-fraction-reciprocal <n>,
    18
--held-out-fraction-seed <n>, 18
--help, 18
--hyperparameters <file.json>, 18
--ignore-inequalities, 19
--max-epochs <n>, 19
--max-tasks-per-worker <n>, 19
--max-workers-per-gpu <n>, 19
--min-measurements-per-allele <n>,
    18
--n-models <n>, 19
--num-jobs <n>, 19
--out-models-dir <dir>, 18
--save-interval <n>, 19
--verbosity <verbosity>, 19
--worker-log-dir <worker_log_dir>,
    19
-h, 18
mhcflurry-class1-train-pan-allele-models
    command line option
--additional-complete-file
    <additional_complete_file>, 22
--allele-sequences <file.csv>, 21
--backend {tensorflow-gpu,tensorflow-cpu,tensor
    21
--cluster-max-retries
    <cluster_max_retries>, 22
--cluster-parallelism, 22
--cluster-results-workdir
    <cluster_results_workdir>, 22
--cluster-script-prefix-path
    <cluster_script_prefix_path>,
    22
--cluster-submit-command
    <cluster_submit_command>, 22
--continue-incomplete, 21
--data <file.csv>, 21
--debug, 21
--gpus <n>, 22
--held-out-measurements-per-allele-fraction-and
    <x>, 21
--help, 21
--hyperparameters <file.json>, 21
--ignore-inequalities, 21
--max-epochs <n>, 21
--max-tasks-per-worker <n>, 22
--max-workers-per-gpu <n>, 22
--num-folds <n>, 21
--num-jobs <n>, 21
--only-initialize, 21
--out-models-dir <dir>, 21
--pretrain-data <file.csv>, 21
--verbosity <verbosity>, 21
--worker-log-dir <worker_log_dir>,
    22
-h, 21
mhcflurry-class1-train-presentation-models
    command line option
--affinity-predictor <dir>, 26
--data <file.csv>, 26
--debug, 26
--help, 26
--hla-column <hla_column>, 26
--out-models-dir <dir>, 26
--processing-predictor-with-flanks
    <dir>, 26
--processing-predictor-without-flanks
    <dir>, 26
--target-column <target_column>, 26
--verbosity <verbosity>, 26
-h, 26
mhcflurry-class1-train-processing-models
    command line option
```



```

--additional-complete-file
    <additional_complete_file>, 24
--backend {tensorflow-gpu,tensorflow-cpu,tensorflow-gpu},
    24
--cluster-max-retries
    <cluster_max_retries>, 24
--cluster-parallelism, 24
--cluster-results-workdir
    <cluster_results_workdir>, 24
--cluster-script-prefix-path
    <cluster_script_prefix_path>,
    24
--cluster-submit-command
    <cluster_submit_command>, 24
--continue-incomplete, 24
--data <file.csv>, 23
--debug, 24
--gpus <n>, 24
--held-out-samples <n>, 23
--help, 23
--hyperparameters <file.json>, 23
--max-epochs <n>, 24
--max-tasks-per-worker <n>, 24
--max-workers-per-gpu <n>, 24
--num-folds <n>, 23
--num-jobs <n>, 24
--num-replicates <n>, 23
--only-initialize, 24
--out-models-dir <dir>, 23
--verbosity <verbosity>, 24
--worker-log-dir <worker_log_dir>,
    24
-h, 23
mhcflurry-downloads command line
    option
--help, 17
--quiet, 17
--verbose, 17
-h, 17
-v, 17
mhcflurry-downloads-fetch command line
    option
--already-downloaded-dir <dir>, 17
--help, 17
--keep, 17
--release <release>, 17
-h, 17
download, 17
mhcflurry-downloads-info command line
    option
--help, 18
-h, 18
mhcflurry-downloads-path command line
    option
--help, 18
-h, 18
download_name, 18
mhcflurry-downloads-url command line
    option
--help, 18
-h, 18
download_name, 18
mhcflurry-predict command line option
--affinity-only, 14
--allele-column <name>, 14
--alleles <allele>, 14
--always-include-best-allele, 14
--c-flank-column <name>, 14
--help, 14
--list-supported-alleles, 14
--list-supported-peptide-lengths, 14
--models <dir>, 14
--n-flank-column <name>, 14
--no-affinity-percentile, 14
--no-flanking, 14
--no-throw, 14
--out <output.csv>, 14
--output-delimiter <char>, 14
--peptide-column <name>, 14
--peptides <peptide>, 14
--prediction-column-prefix <name>,
    14
--version, 14
-h, 14
input.csv, 13
mhcflurry-predict-scan command line
    option
--alleles <allele>, 16
--help, 15
--input-format {guess,csv,fasta}, 16
--list-supported-alleles, 15
--list-supported-peptide-lengths, 16
--models <dir>, 16
--no-affinity-percentile, 16
--no-flanking, 16
--no-throw, 16
--out <output.csv>, 16
--output-delimiter <char>, 16
--peptide-lengths
    <peptide_lengths>, 16
--results-all, 16
--results-best
    {presentation_score,processing_score,affinit
    16
--results-filtered
    {presentation_score,processing_score,affinit
    16
--sequence-column <name>, 16

```

```

--sequence-id-column <name>, 16
--sequences <seq>, 16
--threshold-affinity
    <threshold_affinity>, 16
--threshold-affinity-percentile
    <threshold_affinity_percentile>, 16
--threshold-presentation-score
    <threshold_presentation_score>, 16
--threshold-processing-score
    <threshold_processing_score>, 16
--version, 16
-h, 15
input, 15
miscellaneous_hyperparameter_defaults
    (mhcflurry.class1_neural_network.Class1NeuralNetwork
    attribute), 58
miscellaneous_hyperparameter_defaults
    (mhcflurry.Class1NeuralNetwork attribute), 34
mkdir_p() (in module
    mhcflurry.downloads_command), 81
model_inputs (mhcflurry.class1_presentation_predictor.Class1PresentationPredictor
    attribute), 63
model_inputs (mhcflurry.Class1PresentationPredictor
    attribute), 43
model_name() (mhcflurry.class1_affinity_predictor.Class1AffinityPredictor
    static method), 53
model_name() (mhcflurry.class1_processing_predictor.Class1ProcessingPredictor
    static method), 71
model_name() (mhcflurry.Class1AffinityPredictor
    static method), 29
model_name() (mhcflurry.Class1ProcessingPredictor
    static method), 39
model_select() (in module
    mhcflurry.select_allele_specific_models_command), 95
model_select() (in module
    mhcflurry.select_pan_allele_models_command), 97
model_select() (in module
    mhcflurry.select_processing_models_command), 97
model_select() (mhcflurry.class1_affinity_predictor.Class1AffinityPredictor
    method), 57
model_select() (mhcflurry.Class1AffinityPredictor
    method), 33
models_grid() (mhcflurry.hyperparameters.HyperparameterDefaults
    method), 88
mse() (in module mhcflurry.select_pan_allele_models_command), 96
MSEModelSelector (class in
    mhcflurry.select_allele_specific_models_command), 95
MSEWithInequalities (class in
    mhcflurry.custom_loss), 77
MSEWithInequalitiesAndMultipleOutputs
    (class in mhcflurry.custom_loss), 78
MultiallelicMassSpecLoss (class in
    mhcflurry.custom_loss), 78

N
name (mhcflurry.custom_loss.MSEWithInequalities at-
    tribute), 78
name (mhcflurry.custom_loss.MSEWithInequalitiesAndMultipleOutputs
    attribute), 78
name (mhcflurry.custom_loss.MultiallelicMassSpecLoss
    attribute), 78
network() (mhcflurry.class1_neural_network.Class1NeuralNetwork
    method), 58
network() (mhcflurry.class1_processing_neural_network.Class1ProcessingNeuralNetwork
    method), 69
network() (mhcflurry.Class1NeuralNetwork method), 34
network() (mhcflurry.Class1ProcessingNeuralNetwork
    method), 41
network_hyperparameter_defaults
    (mhcflurry.class1_neural_network.Class1NeuralNetwork
    attribute), 57
network_hyperparameter_defaults
    (mhcflurry.class1_processing_neural_network.Class1ProcessingNeuralNetwork
    attribute), 68
network_hyperparameter_defaults
    (mhcflurry.Class1NeuralNetwork attribute), 34
network_hyperparameter_defaults
    (mhcflurry.Class1ProcessingNeuralNetwork
    attribute), 41
network_input() (mhcflurry.class1_processing_neural_network.Class1ProcessingNeuralNetwork
    method), 69
network_input() (mhcflurry.Class1ProcessingNeuralNetwork
    method), 42
neural_networks()
    (mhcflurry.class1_affinity_predictor.Class1AffinityPredictor
    property), 52
neural_networks()
    (mhcflurry.Class1AffinityPredictor property), 28
NumPyJSONEncoder (class in mhcflurry.common), 75

O
open_file() (mhcflurry.fasta.FastaParser static
    method), 86
optimize() (mhcflurry.class1_affinity_predictor.Class1AffinityPredictor
    method), 53
optimize() (mhcflurry.Class1AffinityPredictor
    method), 29

```

## P

`path_subcommand()` (in module `mhcflurry.downloads_command`), 83  
`peptide_lengths()` (`mhcflurry.flanking_encoding.EncodingResult` property), 86  
`peptides_to_network_input()` (`mhcflurry.class1_neural_network.Class1NeuralNetwork` method), 59  
`peptides_to_network_input()` (`mhcflurry.Class1NeuralNetwork` method), 35  
`percentile_ranks()` (`mhcflurry.class1_affinity_predictor.Class1AffinityPredictor` method), 55  
`percentile_ranks()` (`mhcflurry.Class1AffinityPredictor` method), 31  
`PercentRankTransform` (class in `mhcflurry.percent_rank_transform`), 90  
`plan()` (`mhcflurry.random_negative_peptides.RandomNegativePeptides` method), 93  
`plan_by_allele()` (`mhcflurry.random_negative_peptides.RandomNegativePeptides` method), 93  
`plan_by_allele_equalize_nonbinders()` (`mhcflurry.random_negative_peptides.RandomNegativePeptides` method), 93  
`plan_by_length()` (`mhcflurry.random_negative_peptides.RandomNegativePeptides` method), 93  
`plan_summary()` (`mhcflurry.select_allele_specific_models_command.CombinedModelSelector` method), 95  
`plan_summary()` (`mhcflurry.select_allele_specific_models_command.ConsensusModelSelector` method), 95  
`plan_summary()` (`mhcflurry.select_allele_specific_models_command.MassSpecModelSelector` method), 96  
`plan_summary()` (`mhcflurry.select_allele_specific_models_command.MSEModelSelector` method), 96  
`positional_frequency_matrix()` (in module `mhcflurry.common`), 75  
`ppv()` (`mhcflurry.select_allele_specific_models_command.MassSpecModelSelector` static method), 96  
`predict()` (`mhcflurry.class1_affinity_predictor.Class1AffinityPredictor` method), 55  
`predict()` (`mhcflurry.class1_neural_network.Class1NeuralNetwork` method), 61  
`predict()` (`mhcflurry.class1_presentation_predictor.Class1PresentationPredictor` method), 65  
`predict()` (`mhcflurry.class1_processing_neural_network.Class1ProcessingNeuralNetwork` method), 69  
`predict()` (`mhcflurry.class1_processing_predictor.Class1ProcessingPredictor` method), 71  
`predict()` (`mhcflurry.Class1AffinityPredictor` method), 31  
`predict()` (`mhcflurry.Class1NeuralNetwork` method), 37  
`predict()` (`mhcflurry.Class1PresentationPredictor` method), 45  
`predict()` (`mhcflurry.Class1ProcessingNeuralNetwork` method), 42  
`predict()` (`mhcflurry.Class1ProcessingPredictor` method), 40  
`predict()` (`mhcflurry.select_allele_specific_models_command.Scramble` method), 95  
`predict_affinity()` (`mhcflurry.class1_presentation_predictor.Class1PresentationPredictor` method), 63  
`predict_affinity()` (`mhcflurry.Class1PresentationPredictor` method), 43  
`predict_encoded()` (`mhcflurry.class1_processing_neural_network.Class1ProcessingNeuralNetwork` method), 69  
`predict_encoded()` (`mhcflurry.Class1ProcessingNeuralNetwork` method), 42  
`predict_processing()` (`mhcflurry.class1_presentation_predictor.Class1PresentationPredictor` method), 64  
`predict_processing()` (`mhcflurry.Class1PresentationPredictor` method), 44  
`predict_sequences()` (`mhcflurry.class1_presentation_predictor.Class1PresentationPredictor` method), 66  
`predict_sequences()` (`mhcflurry.Class1PresentationPredictor` method), 46  
`predict_to_dataframe()` (`mhcflurry.class1_processing_predictor.Class1ProcessingPredictor` method), 72  
`predict_to_dataframe()` (`mhcflurry.Class1AffinityPredictor` method), 32  
`predict_to_dataframe()` (`mhcflurry.Class1ProcessingPredictor` method), 40  
`predict_to_dataframe_encoded()` (`mhcflurry.class1_processing_predictor.Class1ProcessingPredictor` method), 72  
`predict_to_dataframe_encoded()` (`mhcflurry.Class1ProcessingPredictor` method), 40  
`pretrain_data_iterator()` (in module `mhcflurry.train_pan_allele_models_command`), 99

## R

`random_peptides()` (in module `mhcflurry.common`), 75  
`RandomNegativePeptides` (class in `mhcflurry.random_negative_peptides`), 92  
`read_fasta_to_dataframe()` (in module `mhcflurry.fasta`), 86  
`robust_mean()` (in module `mhcflurry.ensemble_centrality`), 86  
`run()` (in module `mhcflurry.calibrate_percentile_ranks_command`), 50  
`run()` (in module `mhcflurry.downloads_command`), 81  
`run()` (in module `mhcflurry.predict_command`), 91  
`run()` (in module `mhcflurry.predict_scan_command`), 92  
`run()` (in module `mhcflurry.select_allele_specific_models_command`), 95  
`run()` (in module `mhcflurry.select_pan_allele_models_command`), 97  
`run()` (in module `mhcflurry.select_processing_models_command`), 97  
`run()` (in module `mhcflurry.train_allele_specific_models_command`), 98  
`run()` (in module `mhcflurry.train_pan_allele_models_command`), 99  
`run()` (in module `mhcflurry.train_presentation_models_command`), 100  
`run()` (in module `mhcflurry.train_processing_models_command`), 100

## S

`save()` (`mhcflurry.class1_affinity_predictor.Class1AffinityPredictor` method), 52  
`save()` (`mhcflurry.class1_presentation_predictor.Class1PresentationPredictor` method), 68  
`save()` (`mhcflurry.class1_processing_predictor.Class1ProcessingPredictor` method), 72  
`save()` (`mhcflurry.Class1AffinityPredictor` method), 28  
`save()` (`mhcflurry.Class1PresentationPredictor` method), 48  
`save()` (`mhcflurry.Class1ProcessingPredictor` method), 40  
`save_weights()` (in module `mhcflurry.common`), 75  
`score_function()` (`mhcflurry.select_allele_specific_models_command.ConsensusModelSelector` method), 95  
`score_function()` (`mhcflurry.select_allele_specific_models_command.MassSpecModelSelector` method), 95  
`score_function()` (`mhcflurry.select_allele_specific_models_command.MSEModelSelector` method), 96  
`score_function()` (`mhcflurry.select_allele_specific_models_command.MSEModelSelector` method), 96  
`ScoreFunction` (class in `mhcflurry.select_allele_specific_models_command`), 95  
`ScrambledPredictor` (class in `mhcflurry.select_allele_specific_models_command`), 95  
`sequence_lengths()` (`mhcflurry.class1_processing_neural_network.Class1ProcessingNeuralNetwork` property), 68  
`sequence_lengths()` (`mhcflurry.class1_processing_predictor.Class1ProcessingPredictor` property), 71  
`sequence_lengths()` (`mhcflurry.Class1ProcessingNeuralNetwork` property), 41  
`sequence_lengths()` (`mhcflurry.Class1ProcessingPredictor` property), 39  
`sequences_to_fixed_length_index_encoded_array()` (`mhcflurry.encodable_sequences.EncodableSequences` class method), 84  
`set_allele_representations()` (`mhcflurry.class1_neural_network.Class1NeuralNetwork` method), 62  
`set_allele_representations()` (`mhcflurry.Class1NeuralNetwork` method), 38  
`set_keras_backend()` (in module `mhcflurry.common`), 74  
`StandardKerasLoss` (class in `mhcflurry.custom_loss`), 77  
`startup()` (in module `mhcflurry.testing_utils`), 98  
`subselect()` (`mhcflurry.hyperparameters.HyperparameterDefaults` method), 88  
`subselect_df_held_out()` (in module `mhcflurry.train_allele_specific_models_command`), 98  
`supported_alleles()` (`mhcflurry.class1_affinity_predictor.Class1AffinityPredictor` property), 52  
`supported_alleles()` (`mhcflurry.class1_presentation_predictor.Class1PresentationPredictor` property), 63  
`supported_alleles()` (`mhcflurry.Class1AffinityPredictor` property), 28  
`supported_alleles()` (`mhcflurry.Class1PresentationPredictor` property), 43  
`supported_peptide_lengths()` (`mhcflurry.class1_processing_predictor.Class1ProcessingPredictor` property), 52  
`supported_peptide_lengths()` (`mhcflurry.class1_neural_network.Class1NeuralNetwork` property), 59  
`supported_peptide_lengths()` (`mhcflurry.class1_presentation_predictor.Class1PresentationPredictor` property), 63

property), 63

supported\_peptide\_lengths() (mhcflurry.ClassIAffinityPredictor property), 28

supported\_peptide\_lengths() (mhcflurry.ClassINeuralNetwork property), 36

supported\_peptide\_lengths() (mhcflurry.ClassIPresentationPredictor property), 43

supports\_inequalities (mhcflurry.custom\_loss.MSEWithInequalities attribute), 78

supports\_inequalities (mhcflurry.custom\_loss.MSEWithInequalitiesAndMultipleOutputs attribute), 78

supports\_inequalities (mhcflurry.custom\_loss.MultiallelicMassSpecLoss attribute), 78

supports\_inequalities (mhcflurry.custom\_loss.StandardKerasLoss attribute), 77

supports\_multiple\_outputs (mhcflurry.custom\_loss.MSEWithInequalities attribute), 78

supports\_multiple\_outputs (mhcflurry.custom\_loss.MSEWithInequalitiesAndMultipleOutputs attribute), 78

supports\_multiple\_outputs (mhcflurry.custom\_loss.MultiallelicMassSpecLoss attribute), 78

supports\_multiple\_outputs (mhcflurry.custom\_loss.StandardKerasLoss attribute), 77

svd\_orthonormal() (in module mhcflurry.data\_dependent\_weights\_initialization), 79

## T

to\_ic50() (in module mhcflurry.regression\_target), 94

to\_series() (mhcflurry.percent\_rank\_transform.PercentRankTransform method), 90

TqdmUpTo (class in mhcflurry.downloads\_command), 81

train\_model() (in module mhcflurry.train\_allele\_specific\_models\_command), 98

train\_model() (in module mhcflurry.train\_pan\_allele\_models\_command), 99

train\_model() (in module mhcflurry.train\_processing\_models\_command), 100

train\_models() (in module mhcflurry.train\_pan\_allele\_models\_command), 99

train\_models() (in module mhcflurry.train\_processing\_models\_command), 100

transform() (mhcflurry.percent\_rank\_transform.PercentRankTransform method), 90

TransformPredictionsLossWrapper (class in mhcflurry.custom\_loss), 77

## U

unknown\_character (mhcflurry.encodable\_sequences.EncodableSequences attribute), 83

unknown\_character (mhcflurry.flanking\_encoding.FlankingEncoding attribute), 87

update\_network\_description() (mhcflurry.class1\_neural\_network.Class1NeuralNetwork method), 59

update\_network\_description() (mhcflurry.class1\_processing\_neural\_network.Class1ProcessingNeuralNetwork method), 69

update\_network\_description() (mhcflurry.Class1NeuralNetwork method), 35

update\_network\_description() (mhcflurry.Class1ProcessingNeuralNetwork method), 41

update\_to() (mhcflurry.downloads\_command.TqdmUpTo method), 83

url\_subcommand() (in module mhcflurry.downloads\_command), 83

usable\_for\_allele() (mhcflurry.select\_allele\_specific\_models\_command.CombinedModelSequences method), 95

usable\_for\_allele() (mhcflurry.select\_allele\_specific\_models\_command.ConsensusModelSequences method), 95

usable\_for\_allele() (mhcflurry.select\_allele\_specific\_models\_command.MassSpecModelSequences method), 96

usable\_for\_allele() (mhcflurry.select\_allele\_specific\_models\_command.MSEModelSequences method), 96

## V

variable\_length\_to\_fixed\_length\_categorical() (mhcflurry.encodable\_sequences.EncodableSequences method), 83

variable\_length\_to\_fixed\_length\_vector\_encoding() (mhcflurry.encodable\_sequences.EncodableSequences method), 84

vector\_encode() (mhcflurry.flanking\_encoding.FlankingEncoding method), 87

`vector_encoding_length()` (in module  
`mhcflurry.amino_acid`), 49

## W

`weights_path()` (`mhcflurry.classI_affinity_predictor.ClassIAffinityPredictor`  
static method), 53

`weights_path()` (`mhcflurry.classI_processing_predictor.ClassIProcessingPredictor`  
static method), 71

`weights_path()` (`mhcflurry.ClassIAffinityPredictor`  
static method), 29

`weights_path()` (`mhcflurry.ClassIProcessingPredictor`  
static method), 39

`with_defaults()` (`mhcflurry.hyperparameters.HyperparameterDefaults`  
method), 88

`worker_entry_point()` (in module  
`mhcflurry.cluster_parallelism`), 74

`worker_init()` (in module  
`mhcflurry.local_parallelism`), 90

`worker_init_entry_point()` (in module  
`mhcflurry.local_parallelism`), 89

`worker_pool_with_gpu_assignments()` (in  
module `mhcflurry.local_parallelism`), 89

`worker_pool_with_gpu_assignments_from_args()`  
(in module `mhcflurry.local_parallelism`), 88

`WrapException`, 90

## Y

`yes_no()` (in module `mhcflurry.downloads_command`),  
81