

## \* Prerequisites

In this assignment you will implement the Naive Bayes Classifier. Before starting this assignment, make sure you understand the concepts discussed in the videos in Week 2 about Naive Bayes. You can also find it useful to read Chapter 1 of the textbook.

Also, make sure that you are familiar with the `numpy.ndarray` class of python's `numpy` library and that you are able to answer the following questions:

Let's assume `a` is a numpy array.

- What is an array's shape (e.g., what is the meaning of `a.shape`)?
- What is numpy's reshaping operation? How much computational overhead would it induce?
- What is numpy's transpose operation, and how is it different from reshaping? Does it cause computation overhead?
- What is the meaning of the commands `a.reshape(-1, 1)` and `a.reshape(-1)`?
- What would happen to the variable `a` after we call `b = a.reshape(-1)`? Does any of the attributes of `a` change?
- How do assignments in python and numpy work in general?
  - Does the `b=a` statement use copying by value? Or is it copying by reference?
  - Would the answer to the previous question change depending on whether `a` is a numpy array or a scalar value?

You can answer all of these questions by

- Reading numpy's documentation from <https://numpy.org/doc/stable/>.
- Making trials using dummy variables.

## \*Assignment Summary

The UC Irvine machine learning data repository hosts a famous dataset, the Pima Indians dataset, on whether a patient has diabetes originally owned by the National Institute of Diabetes and Digestive and Kidney Diseases and donated by Vincent Sigillito. You can find it at <https://www.kaggle.com/ucml/pima-indians-diabetes-database/data>. This data has a set of attributes of patients, and a categorical variable telling whether the patient is diabetic or not. For several attributes in this data set, a value of 0 may indicate a missing value of the variable. It has a total of 768 data-points.

- Part 1-A)** First, you will build a simple naive Bayes classifier to classify this data set. We will use 20% of the data for evaluation and the other 80% for training.

You should use a normal distribution to model each of the class-conditional distributions.

Report the accuracy of the classifier on the 20% evaluation data, where accuracy is the number of correct predictions as a fraction of total predictions.

- Part 1-B)** Next, you will adjust your code so that, for attributes 3 (Diastolic blood pressure), 4 (Triceps skin fold thickness), 6 (Body mass index), and 8 (Age), it regards a value of 0 as a missing value when estimating the class-conditional distributions, and the posterior.

Report the accuracy of the classifier on the 20% that was held out for evaluation.

- Part 1-C)** Last, you will have some experience with SVMlight, an off-the-shelf implementation of Support Vector Machines or SVMs. For now, you don't need to understand much about SVM's we will explore them in more depth in the following exercises. You will install SVMlight, which you can find at <http://svmlight.joachims.org>, to train and evaluate an SVM to classify this data.

You should NOT substitute NA values for zeros for attributes 3, 4, 6, and 8.

Report the accuracy of the classifier on the held out 20%

## 0. Data

### 0.1 Description

The UC Irvine's Machine Learning Data Repository Department hosts a Kaggle Competition with famous collection of data on whether a patient has diabetes (the Pima Indians dataset), originally owned by the National Institute of Diabetes and Digestive and Kidney Diseases and donated by Vincent Sigillito.

You can find this data at <https://www.kaggle.com/ucml/pima-indians-diabetes-database/data>. The Kaggle website offers valuable visualizations of the original data dimensions in its dashboard. It is quite insightful to take the time and make sense of the data using their dashboard before applying any method to the data.

### 0.2 Information Summary

- Input/Output:** This data has a set of attributes of patients, and a categorical variable telling whether the patient is diabetic or not.
- Missing Data:** For several attributes in this data set, a value of 0 may indicate a missing value of the variable.
- Final Goal:** We want to build a classifier that can predict whether a patient has diabetes or not. To do this we will train multiple kinds of models, and will be handing the missing data with different approaches for each method (i.e. some methods will ignore their existence, while others may do something about the missing data).

### 0.3 Loading

```
In [1]: import pandas as pd
import numpy as np
import seaborn as sns
import matplotlib.pyplot as plt

from anl_utils import test_case_checker

In [2]: df = pd.read_csv('../BasicClassification-lib/diabetes.csv')
df.head()
```

	Pregnancies	Glucose	BloodPressure	SkinThickness	Insulin	BMI	DiabetesPedigreeFunction	Age	Outcome
0	6	148	72	35	0	33.6		0.627	50
1	1	85	66	29	0	26.6		0.351	31
2	8	183	64	0	0	23.3		0.672	32
3	1	89	66	23	94	28.1		0.167	21
4	0	137	40	35	168	43.1		0.248	33

```
In [3]: # Let's generate the split ourselves.
np.random.seed(2345)
rand_unifs = np.random.uniform(0,1,size=df.shape[0])
division_thresh = np.percentile(rand_unifs, 80)
train_indicator = rand_unifs < division_thresh
eval_indicator = rand_unifs >= division_thresh

In [4]: train_df = df[train_indicator].reset_index(drop=True)
train_features = train_df.loc[:, train_df.columns != 'Outcome'].values
train_labels = train_df['Outcome'].values
eval_df = df[~train_indicator].reset_index(drop=True)
eval_features = eval_df.loc[:, eval_df.columns != 'Outcome'].values
eval_labels = eval_df['Outcome'].values
eval_df.head()
```

	Pregnancies	Glucose	BloodPressure	SkinThickness	Insulin	BMI	DiabetesPedigreeFunction	Age	Outcome
0	1	85	66	29	0	26.6		0.351	31
1	8	183	64	0	0	23.3		0.672	32
2	1	89	66	23	94	28.1		0.167	21
3	0	137	40	35	168	43.1		0.248	33
4	5	116	74	0	0	25.6		0.201	30

```
In [5]: eval_df = df[eval_indicator].reset_index(drop=True)
eval_features = eval_df.loc[:, eval_df.columns != 'Outcome'].values
eval_labels = eval_df['Outcome'].values
eval_df.head()
```

	Pregnancies	Glucose	BloodPressure	SkinThickness	Insulin	BMI	DiabetesPedigreeFunction	Age	Outcome
0	6	148	72	35	0	33.6		0.627	50
1	3	78	50	32	88	31.0		0.248	26
2	10	168	74	0	0	38.0		0.537	34
3	0	118	84	47	230	45.8		0.551	31
4	7	107	74	0	0	29.6		0.254	31

```
In [6]: train_features.shape, train_labels.shape, eval_features.shape, eval_labels.shape

Out[6]: ((614, 8), (614,), (154, 8), (154,))
```

### 0.2 Pre-processing The Data

Some of the columns exhibit missing values. We will use a Naive Bayes Classifier later that will treat such missing values in a special way. To be specific, for attribute 3 (Diastolic blood pressure), attribute 4 (Triceps skin fold thickness), attribute 6 (Body mass index), and attribute 8 (Age), we should regard a value of 0 as a missing value.

Therefore, we will be creating the `train_features_with_nans` and `eval_features_with_nans` numpy arrays to be just like their `train_features` and `eval_features` counter-parts, but with the zero-values in such columns replaced with nans.

```
In [7]: train_df_with_nans = train_df.copy(deep=True)
eval_df_with_nans = eval_df.copy(deep=True)
for col_with_nans in ['BloodPressure', 'SkinThickness', 'BMI', 'Age']:
    train_df_with_nans[col_with_nans] = train_df_with_nans[col_with_nans].replace(0, np.nan)
    eval_df_with_nans[col_with_nans] = eval_df_with_nans[col_with_nans].replace(0, np.nan)
train_features_with_nans = train_df_with_nans.loc[:, train_df_with_nans.columns != 'Outcome']
eval_features_with_nans = eval_df_with_nans.loc[:, eval_df_with_nans.columns != 'Outcome']

In [8]: print('Here are the training rows with at least one missing values.')
print('')
print('You can see that such incomplete data points constitute a substantial part of ')
print('')
nan_training_data = train_df_with_nans[train_df_with_nans.isna().any(axis=1)]
nan_training_data
```

Here are the training rows with at least one missing values.

You can see that such incomplete data points constitute a substantial part of the data.

	Pregnancies	Glucose	BloodPressure	SkinThickness	Insulin	BMI	DiabetesPedigreeFunction	Age	Outcome
1	8	183	64	0	0	23.3		0.672	32
4	5	116	74	0	0	25.6		0.201	30
5	10	168	74	0	0	38.0		0.537	34
7	8	118	84	47	230	45.8		0.551	31
8	4	110	92.0	NaN	0	37.6		0.191	30
...	...	...	...	...	...	...	...	...	...
598	6	162	62.0	NaN	0	24.3		0.178	50
599	4	136	70.0	NaN	0	31.2		0.182	22
605	1	106	76.0	NaN	0	37.5		0.197	26
606	6	190	92.0	NaN	0	35.5		0.278	66
612	1	126	60.0	NaN	0	30.1		0.349	47

186 rows x 9 columns

## 1. Part 1 (Building a simple Naive Bayes Classifier)

Consider a single sample  $(\mathbf{x}, y)$ , where the feature vector is denoted with  $\mathbf{x}$ , and the label is denoted with  $y$ . We will also denote the  $j^{th}$  feature of  $\mathbf{x}$  with  $x^{(j)}$ .

According to the textbook, the Naive Bayes Classifier uses the following decision rule:

\*Choose  $y$  such that

$$\left[ \log p(y) + \sum_j \log p(x^{(j)}|y) \right]$$

is the largest"

However, we first need to define the probabilistic models of the prior  $p(y)$  and the class-conditional feature distributions  $p(x^{(j)}|y)$  using the training data.

- Modeling the prior  $p(y)$ : We fit a Bernoulli distribution to the Outcome variable of `train_df`.
- Modelling the class-conditional feature distributions  $p(x^{(j)}|y)$ : We fit Gaussian distributions, and infer the Gaussian mean and variance parameters from `train_df`.

### Task 1

Write a function `log_prior` that takes a numpy array `train_labels` as input, and outputs the following vector as a column numpy array (i.e. with shape  $(2, 1)$ ).

$$\log p_y = \begin{bmatrix} \log p(y=0) \\ \log p(y=1) \end{bmatrix}$$

Try and avoid the utilization of loops as much as possible. No loops are necessary.

**Hint:** Make sure all the array shapes are what you need and expect. You can reshape any numpy array without any tangible computational overhead.

```
In [9]: #Method for to get the output column numpy array with the shape (2,1).

#Parameter:
#train_labels : numpy.ndarray
#A numpy array

#Returns:
#A column numpy array with required shape (2,1)

def log_prior(train_labels):

    # For Getting the numbr of elements in input numpy arrays
    no_train = train_labels.size

    # to get the first one that is get log p(y = 0)
    log_p_y0 = np.log(train_labels[train_labels == 0].size / no_train)

    # For Second one that is to get log p(y = 1)
    log_p_y1 = np.log(train_labels[train_labels == 1].size / no_train)

    # Creating the matrix and reshaping it to (2, 1)
    log_py = np.array([log_p_y0, log_p_y1]).reshape((2, 1))

    # Assert thr shapes of given matrix which is identical to 2x1
    assert log_py.shape == (2, 1)

    # now Return matrix
    return log_py

In [10]: # Performing sanity checks on your implementation
some_labels = np.array([0, 1, 0, 1, 0, 0, 1, 1, 0, 0, 1, 1, 1, 0, 0, 1, 1, 1, 1])
some_log_py = log_prior(some_labels)
assert np.array_equal(some_log_py.round(3), np.array([[0.916], [-0.511]]))

# Checking against the pre-computed test database
test_results = test_case_checker(log_prior, task_id=1)
assert test_results['passed'], test_results['message']

In [11]: # This cell is left empty as a separator. You can leave this cell as it is, and you s

In [12]: log_py = log_prior(train_labels)
log_py

Out[12]: array([[0.41610786],
                [-1.0776608]])
```

### Task 2

Write a function `cc_mean_ignore_missing` that takes the numpy arrays `train_features` and `train_labels` as input, and outputs the following matrix with the shape  $(8, 2)$ , where 8 is the number of features.

$$\mu_y = \begin{bmatrix} E[x^{(0)}|y=0] & E[x^{(0)}|y=1] \\ E[x^{(1)}|y=0] & E[x^{(1)}|y=1] \\ \vdots & \vdots \\ E[x^{(7)}|y=0] & E[x^{(7)}|y=1] \end{bmatrix}$$

Some points regarding this task:

- The `train_features` numpy array has a shape of  $(N, 8)$  where  $N$  is the number of training data points, and 8 is the number of the features.
- The `train_labels` numpy array has a shape of  $(N, 1)$ .
- You can assume that `train_features` has no missing elements in this task.
- Try and avoid the utilization of loops as much as possible. No loops are necessary.

```
In [13]: def cc_mean_ignore_missing(train_features, train_labels):
    N, d = train_features.shape
    # your code here

    x_when_zero = train_features[np.where(train_labels == 0)[0], :].copy()
    mean_zero = np.mean(x_when_zero, axis=0).reshape(-1, 1)

    x_when_one = train_features[np.where(train_labels == 1)[0], :].copy()
    mean_one = np.mean(x_when_one, axis=0).reshape(-1, 1)

    mu_y = np.hstack((mean_zero, mean_one))

    assert mu_y.shape == (d, 2)

    return mu_y

In [14]: # Performing sanity checks on your implementation
some_feats = np.array([[ 1., 85., 66., 29., 0., 26.6, 0.4, 31. ],
                        [ 8., 183., 64., 0., 0., 23.3, 0.7, 32. ],
                        [ 1., 89., 66., 23., 94., 28.1, 0.2, 21. ],
                        [ 0., 137., 40., 35., 168., 43.1, 2.3, 33. ],
                        [ 5., 116., 74., 0., 0., 25.6, 0.2, 30. ]])
some_labels = np.array([0, 1, 0, 1, 0])
some_mu_y = cc_mean_ignore_missing(some_feats, some_labels)
assert np.array_equal(some_mu_y.round(2), np.array([[ 2.33, 4. ],
                                                    [ 96.67, 160. ],
                                                    [ 68.67, 52. ],
                                                    [ 17.33, 17.5 ],
                                                    [ 31.33, 84. ],
                                                    [ 26.77, 33.2 ],
                                                    [ 0.27, 1.5 ],
                                                    [ 27.33, 32.5 ]]))

# Checking against the pre-computed test database
test_results = test_case_checker(cc_mean_ignore_missing, task_id=2)
assert test_results['passed'], test_results['message']

In [15]: mu_y = cc_mean_ignore_missing(train_features, train_labels)

Out[15]: array([[ 3.4681975, 4.31860029],
                [109.49753086, 142.30143541],
                [ 68.77037037, 70.66028708],
                [19.51398025, 21.97129187],
                [ 66.25679012, 100.55980861],
                [30.31703704, 35.1492823 ],
                [ 4.42825926, 0.55279904],
                [31.57283951, 37.39712919]])

Task 3
```

Write a function `cc_std_ignore_missing` that takes the numpy arrays `train_features` and `train_labels` as input, and outputs the following matrix with the shape  $(8, 2)$ , where 8 is the number of features.

$$\sigma_y = \begin{bmatrix} \text{std}[x^{(0)}|y=0] & \text{std}[x^{(0)}|y=1] \\ \text{std}[x^{(1)}|y=0] & \text{std}[x^{(1)}|y=1] \\ \vdots & \vdots \\ \text{std}[x^{(7)}|y=0] & \text{std}[x^{(7)}|y=1] \end{bmatrix}$$

Some points regarding this task:

- The `train_features` numpy array has a shape of  $(N, 8)$  where  $N$  is the number of training data points, and 8 is the number of the features.
- The `train_labels` numpy array has a shape of  $(N, 1)$ .
- You can assume that `train_features` has no missing elements in this task.
- Try and avoid the utilization of loops as much as possible. No loops are necessary.

```
In [16]: def cc_std_ignore_missing(train_features, train_labels):
    N, d = train_features.shape
    #code start

    #code end
    raise NotImplementedError

    assert sigma_y.shape == (d, 2)

    return sigma_y

In [17]: # This cell is left empty as a separator. You can leave this cell as it is, and you s

In [18]: # Performing sanity checks on your implementation
some_feats = np.array([[ 1., 85., 66., 29., 0., 26.6, 0.4, 31. ],
                        [ 8., 183., 64., 0., 0., 23.3, 0.7, 32. ],
                        [ 1., 89., 66., 23., 94., 28.1, 0.2, 21. ],
                        [ 0., 137., 40., 35., 168., 43.1, 2.3, 33. ],
                        [ 5., 116., 74., 0., 0., 25.6, 0.2, 30. ]])
some_labels = np.array([0, 1, 0, 1, 0])
some_std_y = cc_std_ignore_missing(some_feats, some_labels)
assert np.array_equal(some_std_y.round(3), np.array([[ 1.886, 4. ],
                                                    [13.768, 23. ],
                                                    [ 9.771, 12. ],
                                                    [ 46.499, 17.5 ],
                                                    [44.312, 84. ],
                                                    [10.27, 9.9 ],
                                                    [ 0.094, 0.8 ],
                                                    [ 4.497, 0.5 ]]))

# Checking against the pre-computed test database
test_results = test_case_checker(cc_std_ignore_missing, task_id=3)
assert test_results['passed'], test_results['message']

-----
NotImplementedError                                Traceback (most recent call last)
~/ipykernel_60/607607561673.py in <module>
      7
      8 some_labels = np.array([0, 1, 0, 1, 0])
----> 9 some_mu_y = cc_std_ignore_missing(some_feats, some_labels)
      10
      11 assert np.array_equal(some_std_y.round(3), np.array([[ 1.886, 4. ],
~/tmp/ipykernel_60/1028123466.py in cc_std_ignore_missing(train_features, train_labels)
      5
      6 #code end
----> 7 raise NotImplementedError
      8
      9 assert sigma_y.shape == (d, 2)
NotImplementedError:

In [24]: # This cell is left empty as a separator. You can leave this cell as it is, and you s

In [25]: sigma_y = cc_std_ignore_missing(train_features, train_labels)
sigma_y

-----
NotImplementedError                                Traceback (most recent call last)
~/ipykernel_60/181470216.py in <module>
----> 1 sigma_y = cc_std_ignore_missing(train_features, train_labels)
      2
      3 #code end
----> 4 raise NotImplementedError
      5
      6 #code end
      7
      8 assert sigma_y.shape == (d, 2)
NotImplementedError:

Task 4
```

Write a function `log_prob` that takes the numpy arrays `train_features`,  $\mu_y$ ,  $\sigma_y$  and  $\log p_y$  as input, and outputs the following matrix with the shape  $(N, 2)$

$$\log p_{x,y} = \begin{bmatrix} \log p(y=0) + \sum_{j=1}^7 \log p(x^{(j)}|y=0) & \log p(y=1) + \sum_{j=1}^7 \log p(x^{(j)}|y=1) \\ \log p(y=0) + \sum_{j=1}^7 \log p(x^{(j)}|y=0) & \log p(y=1) + \sum_{j=1}^7 \log p(x^{(j)}|y=1) \\ \vdots & \vdots \\ \log p(y=0) + \sum_{j=1}^7 \log p(x^{(N)}|y=0) & \log p(y=1) + \sum_{j=1}^7 \log p(x^{(N)}|y=1) \end{bmatrix}$$

where

- $N$  is the number of training data points.
- $x_i$  is the  $i^{th}$  training data point.

Try and avoid the utilization of loops as much as possible. No loops are necessary.

**Hint:** Remember that we are modelling  $p(x^{(j)}|y)$  with a Gaussian whose parameters are defined inside  $\mu_y$  and  $\sigma_y$ . Write the Gaussian PDF expression and take its natural log on paper, then implement it.

**Important Note:** Do not use third-party and non-standard implementations for computing  $\log p(x_i^{(j)}|y)$ . Using functions that we are modelling  $p(x^{(j)}|y)$  with a Gaussian whose parameters are defined inside  $\mu_y$  and  $\sigma_y$ . Write the Gaussian PDF expression and take its natural log on paper, then implement it. The other PDF values can easily become extremely small numbers that cannot be represented using floating point standards and thus would be stored as zero. Taking the log of a zero value will throw an error. On the other hand, it is unnecessary to compute and store  $p(x^{(j)}|y)$  and the features. This latter approach is numerically stable, and can be applied when the PDF values are much smaller than could be stored using the common standards.

```
In [28]: def log_prob(train_features, mu_y, sigma_y, log_py):
    N, d = train_features.shape
    # my code start

    pdf_y0 = (np.log(1) - (0.5*(np.log(2) + np.log(np.pi) + np.log(sigma_y[0,0])) +
    log_p_y0 = np.array(log_py[0] + np.sum(pdf_y0, axis=1))
    pdf_y1 = (np.log(1) - (0.5*(np.log(2) + np.log(np.pi) + np.log(sigma_y[1,1])) +
    log_p_y1 = np.array(log_py[1] + np.sum(pdf_y1, axis=1))
    log_p_x_y = np.c_[log_p_y0, log_p_y1]
    print(log_p_x_y)
    # my code end
    raise NotImplementedError

    assert log_p_x_y.shape == (N, 2)

    return log_p_x_y

In [29]: # Performing sanity checks on your implementation
some_feats = np.array([[ 1., 85., 66., 29., 0., 26.6, 0.4, 31. ],
                        [ 8., 183., 64., 0., 0., 23.3, 0.7, 32. ],
                        [ 1., 89., 66., 23., 94., 28.1, 0.2, 21. ],
                        [ 0., 137., 40., 35., 168., 43.1, 2.3, 33. ],
                        [ 5., 116., 74., 0., 0., 25.6, 0.2, 30. ]])
some_labels = np.array([0, 1, 0, 1, 0])
some_mu_y = cc_mean_ignore_missing(some_feats, some_labels)
some_log_py = cc_std_ignore_missing(some_feats, some_labels)
some_log_py_y = log_prior(some_labels)
some_log_p_x_y = log_prob(some_feats, some_mu_y, some_std_y, some_log_py_y)
assert np.array_equal(some_log_p_x_y.round(3), np.array([[ -20.822, -36.606],
                                                        [-60.879, -27.944],
                                                        [-21.774, -295.68 ],
                                                        [-417.359, -27.944],
                                                        [-23.2, -42.6 ]]))

# Checking against the pre-computed test database
test_results = test_case_checker(log_prob, task_id=4)
assert test_results['passed'], test_results['message']

-----
NotImplementedError                                Traceback (most recent call last)
~/ipykernel_60/800464219.py in <module>
      9
      10 some_std_y = cc_std_ignore_missing(some_feats, some_labels)
----> 11 some_log_py = log_prior(some_labels)
      12
~/tmp/ipykernel_60/1028123466.py in cc_std_ignore_missing(train_features, train_labels)
      5
      6 #code end
----> 7 raise NotImplementedError
      8
      9 assert sigma_y.shape == (d, 2)
NotImplementedError:

In [ ]: # This cell is left empty as a separator. You can leave this cell as it is, and you s

In [ ]: log_p_x_y = log_prob(train_features, mu_y, sigma_y, log_py)
log_p_x_y

1.1. Writing the Simple Naive Bayes Classifier
```

```
In [ ]: class NBClassifier():
    def __init__(self, train_features, train_labels):
        self.train_features = train_features
        self.train_labels = train_labels
        self.log_py = log_prior(train_labels)
        self.mu_y = self.get_cc_means()
        self.sigma_y = self.get_cc_std()

    def get_cc_means(self):
        mu_y = cc_mean_ignore_missing(self.train_features, self.train_labels)
        return mu_y

    def get_cc_std(self):
        sigma_y = cc_std_ignore_missing(self.train_features, self.train_labels)
        return sigma_y

    def predict(self, features):
        log_p_x_y = log_prob(features, self.mu_y, self.sigma_y, self.log_py)
        return log_p_x_y.argmax(axis=1)

In [ ]: diabetes_classifier = NBClassifier(train_features, train_labels)
train_pred = diabetes_classifier.predict(train_labels)
eval_pred = diabetes_classifier.predict(eval_features)

In [ ]: train_acc = (train_pred==train_labels).mean()
eval_acc = (eval_pred==eval_labels).mean()
print(f'The training data accuracy of your trained model is {train_acc}')
print(f'The evaluation data accuracy of your trained model is {eval_acc}')
```

## 1.2 Running an off-the-shelf implementation of Naive-Bayes For Comparison

```
In [ ]: from sklearn.naive_bayes import GaussianNB
gnb = GaussianNB().fit(train_features, train_labels)
train_pred_sk = gnb.predict(train_features)
eval_pred_sk = gnb.predict(eval_features)
print(f'The training data accuracy of your trained model is {(train_pred_sk == train_labels).mean()}')
print(f'The evaluation data accuracy of your trained model is {(eval_pred_sk == eval_labels).mean()}')
```

## Part 2 (Building a Naive Bayes Classifier Considering Missing Entries)

In this part, we will modify some of the parameter inference functions of the Naive Bayes classifier to make it able to ignore the NA entries when inferring the Gaussian mean and stds.

### Task 5

Write a function `cc_mean_consider_missing` that

- has exactly the same input and output types as the `cc_mean_ignore_missing` function,
  - and has similar functionality to `cc_mean_ignore_missing` except that it can handle and ignore the NA entries when computing the class conditional means.
- You can borrow most of the code from your `cc_mean_ignore_missing` implementation, but you should make it compatible with the existence of NA values in the features.
- Try and avoid the utilization of loops as much as possible. No loops are necessary.

- Hint:** You may find the `np.isnan` function useful.

```
In [ ]: def cc_mean_consider_missing(train_features_with_nans, train_labels):
    N, d = train_features_with_nans.shape
    # your code here
    raise NotImplementedError

In [ ]: # Performing sanity checks on your implementation
some_feats = np.array([[ 1., 85., 66., 29., 0., 26.6, 0.4, 31. ],
                        [ 8., 183., 64., 0., 0., 23.3, 0.7, 32. ],
                        [ 1., 89., 66., 23., 94., 28.1, 0.2, 21. ],
                        [ 0., 137., 40., 35., 168., 43.1, 2.3, 33. ],
                        [ 5., 116., 74., 0., 0., 25.6, 0.2, 30. ]])
for i, n in [(0,0), (1,1), (2,3), (3,4), (4, 2)]:
    some_feats[i,:] = np.nan
    some_mu_y = cc_mean_consider_missing(some_feats, some_labels)
assert np.array_equal(some_mu_y.round(2), np.array([[ 3., 4. ],
                                                    [ 96.67, 137. ],
                                                    [ 66., 52. ],
                                                    [ 14.5, 17.5 ],
                                                    [44.31, 0. ],
                                                    [ 10.27, 9.9 ],
                                                    [ 0.09, 0.8 ],
                                                    [ 4.5, 0.5 ]]))

# Checking against the pre-computed test database
test_results = test_case_checker(cc_mean_consider_missing, task_id=5)
assert test_results['passed'], test_results['message']

In [ ]: # This cell is left empty as a separator. You can leave this cell as it is, and you s

In [ ]: mu_y = cc_mean_consider_missing(train_features_with_nans, train_labels)
mu_y

Task 6
```



1. We have to export the training data to a special format called `svmlight/libsvm`. This can be done using `scikit-learn`.

2. We have to run the `svm_learn` program to learn the model and then store it.

3. We have to import the model back to python.

### 3.1 Exporting the training data to libsvm format

```
In [ ]: from sklearn.datasets import dump_svmlight_file
dump_svmlight_file(train_features, 2*train_labels-1, 'training_feats.data',
                  zero_based=False, comment=None, query_id=None, multilabel=False)
```

### 3.2 Training SVMlight

```
In [ ]: !chmod +x ../BasicClassification-lib/svmlight/svm_learn
from subprocess import Popen, PIPE
process = Popen(['../BasicClassification-lib/svmlight/svm_learn', './training_feats.d
stdout, stderr = process.communicate()
print(stdout.decode("utf-8"))
```

### 3.3 Importing the SVM Model

```
In [ ]: from svmweight import get_svmlight_weights
svm_weights, thresh = get_svmlight_weights('svm_model.txt', printOutput=False)

def svmlight_classifier(train_features):
    return (train_features @ svm_weights - thresh).reshape(-1) >= 0.
```

```
In [ ]: train_pred = svmlight_classifier(train_features)
eval_pred = svmlight_classifier(eval_features)
```

```
In [ ]: train_acc = (train_pred==train_labels).mean()
eval_acc = (eval_pred==eval_labels).mean()
print(f"The training data accuracy of your trained model is {train_acc}")
print(f"The evaluation data accuracy of your trained model is {eval_acc}")
```

```
In [ ]: # Cleaning up after our work is done
!rm -rf svm_model.txt training_feats.data
```