# Homework 3: Bayesian ML and Multiclass

## Bayesian logistic regression with Gaussian Prior



Ans. 1

Ans 2.

②　let $w \sim N(0, \Sigma)$

MAP estimate for $w$ after observing data $D$

$$\arg\max_{w} p(w|D) = \arg\max_{w} p(D|w) \cdot p(w)$$

$$= \arg\max_{w} \left[ \log\left(p(D|w)\right) + \log\left(p(w)\right) \right] \quad \text{[taking log as it is strictly inc. f}^n$$
$$\text{\& makes compute easy]}$$

$$\underbrace{\qquad}_{\substack{\log \text{ likelihood function} \\ \text{for log. reg.}}}$$

$$\arg\max_{w} \left[ \sum_{i=1}^{n} -\frac{1}{2}\left( y^{(i)}\log\left(1+\exp(-x^{(i)T}w)\right) + (1-y^{(i)})\left(1+\exp(x^{(i)T}w)\right)\right) + \log(p(w)) \right]$$

$$= \arg\max_{w} \left( -NLL_D(w) + \log|2\pi\Sigma|^{-\frac{1}{2}} - \frac{1}{2}w^T\Sigma^{-1}w \right)$$

$$\underbrace{\qquad}_{\substack{\downarrow \text{can be removed as it is} \\ \text{not dep. on } w}}$$

$$= \arg\min_{w} \left( NLL_D(w) + \frac{1}{2}w^T\Sigma^{-1}w \right) \equiv \arg\min_{w} \left( \frac{1}{n}NLL_D(w) + \frac{1}{2n}w^T\Sigma^{-1}w \right)$$
$$\underbrace{\qquad}_{①}$$

$L_2$ regularized logistic regression

$$= \arg\min_{w} \left( \frac{1}{n}\sum_{i=1}^{n} \log\left(1+\exp(-y^{(i)}x^{(i)T}w)\right) + \lambda\|w\|_2^2 \right)$$

$$\underline{\downarrow\text{can be rewritten as}}$$

$$\frac{1}{2n}\left[ \sum_{i=1}^{n} y^{(i)}\log\left(1+\exp(-x^{(i)T}w)\right) + (1-y^{(i)})\log\left(1+\exp(+x^{(i)T}w)\right) \right]$$

$$= \frac{1}{n}NLL_D(w)$$

$$\Rightarrow \arg\min_{w} \left( \frac{1}{n}NLL_D(w) + \lambda\|w\|_2^2 \right)$$
$$\underbrace{\qquad}_{②}$$

Equating eq'ls (2) and (1) we see that

$$\frac{1}{2n}\left(w^T \Sigma^{-1} w\right) = \lambda \|w\|_2^2$$

$$= \lambda w^T w$$

$$\Rightarrow \boxed{\Sigma = \frac{1}{2n\lambda} I}$$

please turn over

Ans 3.

(3) for ERM to be the mode of the posterior dist.

such that $\Sigma = I$.

$$\Rightarrow \quad I = \frac{1}{2n\lambda} I$$

$$\boxed{\Rightarrow \quad \lambda = \frac{1}{2n}}$$

please turn over

/

# Coin Flipping with partial observability

Coin flipping with partial observability.

$$p(z = H \mid \theta_1) = \theta_1 \quad \text{// biased coin}.$$

someone reports $z$ to us as $x$
there is a $\frac{1}{2}$ chance $x$ is incorrect if $z = H$.
This is denoted by

$$p(x = H \mid z = H, \theta_2) = \theta_2$$

and

$$p(x = T \mid z = T) = 1.$$

/

Ans 4.

④ To prove

$$p(x=H \mid \theta_1, \theta_2) = \theta_1 \theta_2. \qquad ①$$

Proof :

$$p(x=H \mid \theta_1, \theta_2) = \Big[ \big( p(x=H \mid z=H, \theta_1, \theta_2) \cdot p(z=H \mid \theta_1) \big)$$
$$+ \underbrace{p(x=H \mid z=T, \theta_1, \theta_2) \cdot p(z=T \mid \theta_1)}_{②} \Big]$$

Evaluating ①

$$\Longrightarrow \quad \underline{p(x=H \mid z=H, \theta_1, \theta_2) \cdot p(z=H \mid \theta_1)}$$

$\quad \hookrightarrow$ since it is known that $z=H$ for this case

$$\therefore \quad p(x=H \mid z=H, \theta_1, \theta_2) = p(x=H \mid z=H, \theta_2) = \theta_2$$

and

$$p(z=H \mid \theta_1) = \theta_1$$

$$\therefore \quad ① = \theta_1 \theta_2$$

Evaluating ②

$$p(x=H \mid z=T, \theta_1, \theta_2) \cdot p(z=T \mid \theta_1)$$

$\quad \hookrightarrow$ since $x$ always correctly reports $z$ if $z=T$

$$\therefore \quad p(x=H \mid z=T, \theta_1, \theta_2) = 0$$

and

$$p(z=T \mid \theta_1) = (1-\theta_1)$$

$\Big\}$ $\therefore$ ② = 0

$$\therefore \quad \boxed{p(x=H \mid \theta_1, \theta_2) = ① + ② = \theta_1 \theta_2.} \quad /\!/ \text{ Hence proved}$$

/

Ans 5.

⑤ assuming all events are independent, we get-

$$\chi_{D_x}(\theta_1, \theta_2) = \prod_{i=1}^{N_x} p(x \mid \theta_1, \theta_2)$$

→ incorrectly reported head + all tails correctly reported.

$$\boxed{\chi_{D_x}(\theta_1, \theta_2) = (\theta_1 \theta_2)^{n_h} (1 - \theta_1 \theta_2)^{n_t}}$$

└ correctly reported H

/

## Ans 6.

⑥ Estimating $\theta_1$ and $\theta_2$ through MLE is not usable as say for a given distribution, we get

$$\chi_{D_\alpha}(\theta_1, \theta_2) = (\theta_1 \theta_2)^{n_h} (1 - \theta_1 \theta_2)^{n_t}$$

if we try MLE, we get

$$\arg\max_{\theta_1, \theta_2} \log \chi_{D_\alpha}(\theta_1, \theta_2)$$

$$= \arg\max_{\theta_1, \theta_2} n_h \log \theta_1 \theta_2 + n_t \log (1 - \theta_1 \theta_2)$$

$$\frac{\partial \log \chi(\theta_1, \theta_2)}{\partial \theta_1} = 0 = \frac{n_h}{\theta_1 \theta_2} \cdot \theta_2 + \frac{n_t}{1 - \theta_1 \theta_2} \cdot (-\theta_2) \qquad \left.\begin{array}{l} \theta_1, \theta_2 \in (0, 1] \\ \text{[for this to be differentiable]} \end{array}\right.$$

$$\Rightarrow \theta_2 (n_h - n_h \theta_1 \theta_2) - n_t \theta_1 \theta_2^2 = 0$$

$$\Rightarrow n_h \theta_2 - n_h \theta_1 \theta_2^2 - n_t \theta_1 \theta_2^2 = 0$$

$$\Rightarrow (\theta_2)(n_h - \theta_1 \theta_2 (n_h + n_t)) = 0$$

$$\theta_1 \theta_2 = \frac{n_h}{n_h + n_t} \qquad\qquad\qquad \text{▷ identical}$$

$$\frac{\partial \log l(\theta_1, \theta_2)}{\partial \theta_2} = 0 = \frac{n_h}{\theta_1 \theta_2} \theta_1 + \frac{n_t}{1 - \theta_1 \theta_2} \cdot (-\theta_1)$$

$$\Rightarrow (\theta_1)(n_h - (n_h + n_t)\theta_1 \theta_2) = 0$$

$$\Rightarrow \theta_1 \theta_2 = \frac{n_h}{n_h + n_t}$$

∴ we have 2 variables 1 eqⁿ non trivial solⁿ only possible.

Hence since we get non trivial solⁿs for $\theta_1$ and $\theta_2$ we cannot estimate $\theta_1$ and $\theta_2$ using MLE.

/

Ans 7.

⑦ Assuming all events are independent, we can write the likelihood as

$$\lambda(\theta_1, \theta_2) = p(D_r, D_c \mid \theta_1, \theta_2) = p(D_r \mid \theta_1, \theta_2) \cdot p(D_c \mid \theta_1, \theta_2)$$

$$= (\theta_1 \theta_2)^{n_h} (1 - \theta_1 \theta_2)^{n_t} (\theta_1)^{c_h} (1 - \theta_1)^{c_t}$$

$$\arg\max_{\theta_1, \theta_2} \left(\log \lambda(\theta_1, \theta_2)\right) =$$

$$l(\theta_1, \theta_2) = n_h \log(\theta_1 \theta_2) + n_t \log(1 - \theta_1 \theta_2) + c_h \log \theta_1 + c_t \log(1 - \theta_1)$$

Differentiating wrt $\theta_1$ we get :-

$$\frac{\partial l}{\partial \theta_1} = \frac{n_h}{\theta_1 \theta_2} \theta_2 + \frac{n_t (-\theta_2)}{1 - \theta_1 \theta_2} + \frac{c_h}{\theta_1} + \frac{c_t}{1 - \theta_1}(-1) = 0$$

$$= \frac{\theta_2}{(\theta_1 \theta_2)(1 - \theta_1 \theta_2)} \left(n_h - \theta_1 \theta_2 (n_h + n_t)\right) + \frac{1}{\theta_1 (1 - \theta_1)} (c_h - \theta_1 (c_h + c_t)) = 0$$

$$\Rightarrow 0 = \frac{\left(n_h - (\theta_1 \theta_2)(n_h + n_t)\right)}{\theta_1 (1 - \theta_1 \theta_2)} + \frac{(c_h - \theta_1 (c_h + c_t))}{\theta_1 (1 - \theta_1)} \qquad —①$$

Differentiating wrt $\theta_2$ :-

$$\frac{\partial l}{\partial \theta_2} = 0 = \frac{\left(n_h - (\theta_1 \theta_2)(n_h + n_t)\right)}{\theta_2 (1 - \theta_1 \theta_2)} \qquad —②$$

solving these eqⁿs of ② we get

$$\boxed{\theta_1 = \frac{c_h}{c_h + c_t}}$$  ———— non trivial solution.

and

$$\boxed{\theta_2 = \frac{n_h}{n_h + n_t} \cdot \frac{c_h + c_t}{c_h}}$$

∴ we can estimate $\theta_1, \theta_2$ using MLE.

please turn over

## Ans 8.

⑧ $\quad\quad\quad\quad\quad \theta_1 \sim Beta\,(h, t)$

$$\Rightarrow p(\theta_1) = \theta_1^{h-1}\,(1-\theta_1)^{t-1}$$

prior is provided to mitigate the overfitting in $D_c$

$$\Rightarrow \text{posterior of } \theta_1 \text{ based on clean results } D_c$$

$$p(\theta_1 | D_c) \propto p(D_c | \theta_1) \cdot p(\theta_1)$$

$$\Rightarrow p(\theta_1 | D_c) = K \cdot \theta_1^{c_h}(1-\theta_1)^{c_t} \cdot \theta_1^{h-1}(1-\theta_1)^{t-1}$$

$$= K\,\theta_1^{c_h+h-1}\,(1-\theta_1)^{c_t+t-1}$$

$$p(\theta_1 | D_c) = Beta\,(c_h+h,\ c_t+t)$$

$\therefore \arg\max\limits_{\theta_1} p(\theta_1 | D_c) = \boxed{\dfrac{c_h+h-1}{c_h+h+c_t+t-2} = \theta_1\,_{MAP}}$

$\therefore \boxed{\theta_2 = \dfrac{n_h}{n_h+n_t} \cdot \dfrac{c_h+h+c_t+t-2}{c_h+h-1}}$ $\quad \hookrightarrow$ formula on slides.

$\hookrightarrow$ [ from eqn ② of problem 7 ]

*please turn over*

$1$

# ℓ2-regularized empirical risk function for multiclass hinge loss

Ans 9.

⑨  To prove : $J(w)$ is convex function of $w$.

Proof:

Let $f_i(w) = \Delta(y_i, y) + \left( \langle w, \psi(x_i, y) - \psi(x_i, y_i) \rangle \right)$

$\langle w, \psi(x_i, y) \rangle = \langle w, \psi(x_i, y) \rangle + 0$  } affine function

$\Delta(y_i, y) \in \mathbb{R}$

$\langle w, \psi(x_i, y_i) \rangle \in \mathbb{R}$.

∴ $f_i(w)$ is an affine function of $y$ on $w$ and thus convex.

⟹ $\max\limits_{y \in Y} f_y(w) \equiv \max\{f_1(w), f_2(w) \cdots f_k(w)\}$

is also convex with domain
$dom\, f_1 \wedge dom\, f_2 \wedge \cdots \wedge dom\, f_k$

→ (Rosenberg's notes 3.2.4)

1

$\Rightarrow \qquad \frac{1}{n} \sum\limits_{i=1}^{n} \max\limits_{y \in Y} f_y(w) \quad$ is convex $\Big]$ (Rosenberg's notes 3.2.1)

Now $\quad \lambda \|w\|_2^2 \quad$ is convex as $\lambda \in \mathbb{R}$ and $\lambda > 0$ and $\|w\|_2^2$ is convex

$\qquad\qquad\qquad\qquad\qquad\qquad\qquad\qquad$ ↳ (Rosenberg's notes 3.1.3)

$\therefore \quad \lambda \|w\|_2^2 + \frac{1}{n} \sum\limits_{i=1}^{n} \max\limits_{y \in Y} f_y(w) \quad$ is convex [sum of 2 convex $f$'s is convex]

$\Longrightarrow \quad \lambda \|w\|_2^2 + \frac{1}{n} \sum\limits_{i=1}^{n} \max\limits_{y \in Y} \Big[ \Delta(y_i, y) + \langle w, \psi(x_i, y) - \psi(x_i, y_i) \rangle \Big]$

is a convex function.

~~Hence~~ Proved.

## Ans 10.

⑩ $J(w) = \lambda \|w\|_2^2 + \frac{1}{n} \sum_{i=1}^{n} \max_{y \in Y} [\Delta(y_i^o, y) + \langle w, \psi(x_i^o, y) - \psi(x_i, y_i^o) \rangle]$

Let $R(w) = \lambda \|w\|_2^2$ has subgradient

$$g_R(w) = 2\lambda w \in \partial R(w)$$

Let $H(w) = \frac{1}{n} \sum_{i=1}^{n} \max_{y \in Y} [\Delta(y_i, y) + \langle w, \psi(x_i, y) - \psi(x_i, y_i) \rangle]$

for each $i$,

$$\hat{y_i^o} = \arg\max_{y \in Y} [\Delta(y_i, y) + \langle w, \psi(x_i, y) - \psi(x_i, y_i) \rangle]$$

↳ $\hat{y_i^o}$ is the $y$ that maximizes above exp.

from Assignment 2 we know that

given $f(x) = \max f_k(x)$ $\qquad (k = \{1, 2 \cdots, n\})$

if $k$ is any index for which

$$f_k(x) = f(x)$$

then if $g(x) \in \partial f_k(x)$

$$\Rightarrow \quad g(x) \in \partial f(x)$$

Let $f_y(w) = \Delta(y_i, y) + \langle w, \psi(x_i, y) - \psi(x_i, y_i) \rangle$

∴ subgradient of $f_y(w)$ at $w$ can be given by subgradient of $f_{\hat{y_i}}(w)$

where $\hat{y_i} = \arg\max_{y \in Y} f_y(w)$

one subgradient of $f_{\hat{y_i^o}}(w)$ wrt $w$ is $\Rightarrow$

$$g_{f_{\hat{y_i^o}}}(w) = \psi(x_i, \hat{y_i}) - \psi(x_i, y_i) \in \partial(f(w))$$

∴ combining all components of $H(w)$ we get

$$g_H(w) = \frac{1}{n} \sum_{i=1}^{n} (\psi(x_i, \hat{y_i}) - \psi(x_i, y_i^o))$$

we   can   write

$$g_J(w) = g_R(w) + g_H(w)$$

$$\Rightarrow \boxed{g_J(w) = 2\lambda w + \frac{1}{n}\sum_{i=1}^{n}\left(\psi(x_i, \hat{y}_i) - \psi(x_i, y_i)\right)}$$

$$\left\{\text{where } \hat{y}_i = \arg\max_{y \in Y}\left[\Delta(y_i, y) + \langle w, \psi(x_i, y) - \psi(x_i, y_i)\rangle\right]\right\}$$

$$\boxed{g_J(w) \in \partial C J(w)}$$

please turn over

/

Ans 11.

(11) stochastic gradient descent will be based on a single Training e.g. $(x_i, y_i)$

$$\gamma\left(J(w, x_i^o, y_i^o)\right) = \left(2\lambda w + \left(\psi(x_i^o, \hat{y}_i) - \psi(x_i, y_i^o)\right)\right)$$

please turn over

/

Ans 12.

(12) minibatch gradient will be based on avg of samples in the mini-batch.

$$\partial\left(J(w, (x_i, y_i), (x_{i+1}, y_{i+1}) \cdots (x_{i+m-1}, y_{i+m-1}))\right)$$

$$\|$$

$$\left(2\lambda w + \frac{1}{m} \sum_{t=i}^{i+m-1} \left(\psi(x_t, \hat{y}_t) - \psi(x_t, y_t)\right)\right)$$

/

## Hinge Loss is a specialized case of Generalized Hinge Loss

Let $\mathcal{Y} = \{-1, 1\}$. Let $\Delta(y, \hat{y}) = \mathbb{1}y \neq \hat{y}$. If $g(x)$ is the score function in our binary classification setting, then define our compatibility function as

$$
\begin{aligned}
h(x, 1) &= g(x)/2 \\
h(x, -1) &= -g(x)/2.
\end{aligned}
$$

Show that for this choice of $h$, the multiclass hinge loss reduces to hinge loss:

$$
\ell\left(h, (x, y)\right) = \max_{y' \in \mathcal{Y}} \left[\Delta\left(y, y'\right)\right) + h(x, y') - h(x, y)\right] = \max\{0, 1 - yg(x)\}
$$

In this problem we will work on a simple three-class classification example. The data is generated and plotted for you in the skeleton code.

### Hinge loss is a Specialized Case of Generalized Hinge loss

$$\Delta(y, \hat{y}) = 1_{y \neq \hat{y}} = \begin{cases} 0, & y = \hat{y} \\ 1, & y \neq \hat{y} \end{cases}$$

Compatibility function

$$h(x, 1) = g(x)/2 \quad ; \quad h(x, -1) = -g(x)/2.$$

generalized multiclass hinge loss

$$\ell(h, (x, y)) = \max_{y' \in Y} \left[ \Delta(y, \hat{y}) + h(x, y') - h(x, y) \right]$$

Case 1: correct classification $(y' = y)$

$$\Delta(y, y') = 0$$
$$h(x, y) = h(x, y')$$
$$\implies \Delta(y, y') + h(x, y') - h(x, y) = 0$$

Case 2: incorrect classification $(\hat{y} \neq y)$

$$\implies \Delta(y, \hat{y}) = 1$$

Case 2a: $y = 1, \hat{y} = -1$

$$\implies \Delta(y, \hat{y}) + h(x, y') - h(x, y)$$
$$= 1 + \frac{-g(x)}{2} - \frac{g(x)}{2}$$
$$= 1 - g(x)$$

Case 2b: $y = -1, \hat{y} = 1$

$$\implies \Delta(y, \hat{y}) + h(x, y') - h(x, y)$$
$$= 1 + g(x).$$

We can combine cases 2a and 2b to get a gen. expr. for case 2 as follows:

Case 2:
$$= 1 - y \cdot g(x)$$

$\therefore$

$$l(h, (x, y)) = \max_{y \in Y} [\Delta(y, y') + h(x, y') - h(x, y)]$$

$$= \max(0, 1 - y \cdot g(x)) \quad \rbrace \text{ taking the max of case 1 and case 2}$$

Hence Proved.

please turn over

/

# One-vs-All (also known as One-vs-Rest)

Ans 13.

- 
```python
def fit(self, X, y=None):
    #Your code goes here
    for i in range(self.n_classes):
        class_i_labels = np.where(y == i, 1, 0)
        self.estimators[i].fit(X, class_i_labels)
    self.fitted = True
    return self
```

- 
```python
def decision_function(self, X):
    if not self.fitted:
        raise RuntimeError("You must train classifer \
            before predicting data.")

    if not hasattr(self.estimators[0], "decision_function"):
        raise AttributeError(
            "Base estimator doesn't have a\
                decision_function attribute.")
    scores = np.empty((X.shape[0], self.n_classes))
    for i in range(self.n_classes):
        class_i_predicted_labels =
self.estimators[i].decision_function(X)
        scores[:, i] = class_i_predicted_labels
    return scores
```

- 
```python
def predict(self, X):
    """
    Predict the class with the highest score.
    @param X: array-like, shape = [n_samples,n_features] input data
    @returns array-like, shape = [n_samples,] the predicted classes
for each input
    """
    scores = self.decision_function(X)
    return np.argmax(scores, axis = 1)
```

Ans 14.

- 
```
#Here we test the OneVsAllClassifier
from sklearn import svm
svm_estimator = svm.LinearSVC(loss='hinge', fit_intercept=False,
C=200)
clf_onevsall = OneVsAllClassifier(svm_estimator, n_classes=3)
clf_onevsall.fit(X,y)

for i in range(3) :
    print("Coeffs %d"%i)
    print(clf_onevsall.estimators[i].coef_) #Will fail if you
haven't implemented fit yet

# create a mesh to plot in
h = .02  # step size in the mesh
x_min, x_max = min(X[:,0])-3,max(X[:,0])+3
y_min, y_max = min(X[:,1])-3,max(X[:,1])+3
xx, yy = np.meshgrid(np.arange(x_min, x_max, h),
                     np.arange(y_min, y_max, h))
mesh_input = np.c_[xx.ravel(), yy.ravel()]

Z = clf_onevsall.predict(mesh_input)
Z = Z.reshape(xx.shape)
plt.contourf(xx, yy, Z, cmap=plt.cm.coolwarm, alpha=0.8)
# Plot also the training points
plt.scatter(X[:, 0], X[:, 1], c=y, cmap=plt.cm.coolwarm)


from sklearn import metrics
metrics.confusion_matrix(y, clf_onevsall.predict(X))
```
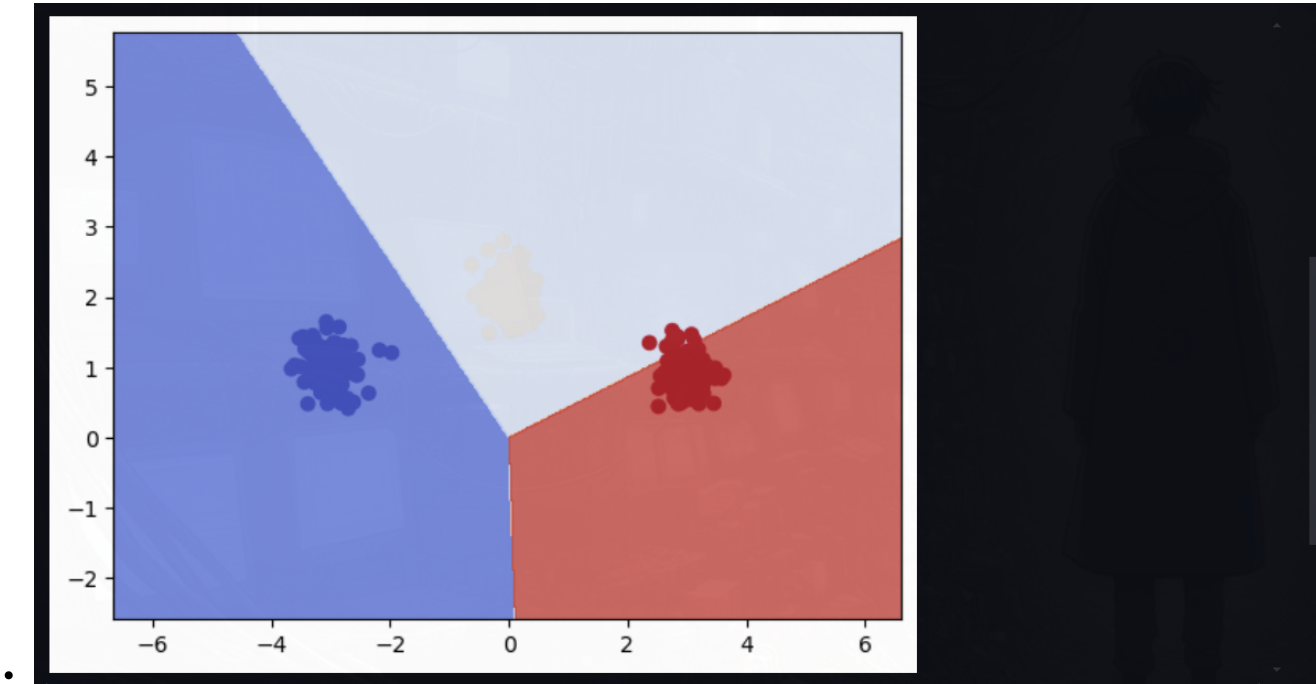
- 
```
Coeffs 0
[[-1.05853334 -0.90294603]]
Coeffs 1
[[0.42121645 0.27171776]]
Coeffs 2
[[ 0.89164752 -0.82601734]]
/home/arjun-prasad/anaconda3/envs/deep_learning/lib/python3.12/site-packages/sklearn/svm/_base
  warnings.warn(
```

- 
```
array([[100,   0,   0],
       [  0, 100,   0],
       [  0,  11,  89]])
```

# Multiclass SVM

Ans 15.

- ```python
  def featureMap(X,y,num_classes) :
      #The following line handles X being a 1d-array or a 2d-array
      num_samples, num_inFeatures = (1,X.shape[0]) if \
          len(X.shape) == 1 else (X.shape[0],X.shape[1])
      #your code goes here, and replaces following return
      if X.ndim == 1:
          X.reshape(1, -1)

      featureMap = np.zeros((num_samples, num_classes*num_inFeatures))
      loc = num_inFeatures*y
      for i in range(num_samples):
          featureMap[i, loc : loc + num_inFeatures] = X[i]

      return featureMap
  ```

Ans 16.

- 
```python
    def sgd(X, y, num_outFeatures, subgd, eta = 0.1, T = 10000):
        num_samples = X.shape[0]
        #your code goes here and replaces following return statement
        w = np.zeros(num_outFeatures)
        for iter in range(T):
            idx = np.random.randint(0, num_samples)
            grad = subgd(X[idx], y[idx], w)
            w = w - (eta*grad)
        return w
```

Ans 17.

- 
```python
def subgradient(self,x,y,w):
    '''
    Computes the subgradient at a given data point x,y
    @param x: sample input
    @param y: sample class
    @param w: parameter vector
    @return returns subgradient vector at given x,y,w
    '''
    #Your code goes here and replaces the following return statement
    subgrad_R = 2*self.lam*w
    y_hat = -1

    f_y_w_max = -np.inf
    for yi in range(self.num_classes):
        if yi == y:
            continue # no point in evaluating this
        f_yi_w = self.Delta(yi, y) + (w @ (self.Psi(x, yi).flatten() - \
            self.Psi(x, y).flatten()))
        if f_yi_w > f_y_w_max:
            f_y_w_max = f_yi_w
            y_hat = yi
    return subgrad_R + (self.Psi(x, y_hat).flatten() - self.Psi(x, y).flatten())
```

- 
```python
def decision_function(self, X):
    if not self.fitted:
        raise RuntimeError("You must train classifer before predicting data.")

    #Your code goes here and replaces following return statement
    num_samples = X.shape[0]
    class_scores = np.zeros((num_samples, self.num_classes))
    for yi in range(self.num_classes):
        class_scores[:, yi] = self.Psi(X, yi) @ self.coef_
    return class_scores
```

- 
```python
def predict(self, X):
    #Your code goes here and replaces following return statement
    class_scores = self.decision_function(X)
    return np.argmax(class_scores, axis = 1)
```

please turn over

Ans 18.

- 
```
from skeleton_code import zeroOne, featureMap, sgd, MulticlassSVM
#the following code tests the MulticlassSVM and sgd
#will fail if MulticlassSVM is not implemented yet
est = MulticlassSVM(6,lam=1)
est.fit(X,y)
print("w:")
print(est.coef_)
Z = est.predict(mesh_input)
Z = Z.reshape(xx.shape)
plt.contourf(xx, yy, Z, cmap=plt.cm.coolwarm, alpha=0.8)
# Plot also the training points
plt.scatter(X[:, 0], X[:, 1], c=y, cmap=plt.cm.coolwarm)


from sklearn import metrics
metrics.confusion_matrix(y, est.predict(X))
```

- 
```
w:
[-0.41923524 -0.41923524 -0.06249277 -0.06249277  0.48172801  0.48172801]

array([[100,   0,   0],
       [  0,   0, 100],
       [  0,   0, 100]])
```

-