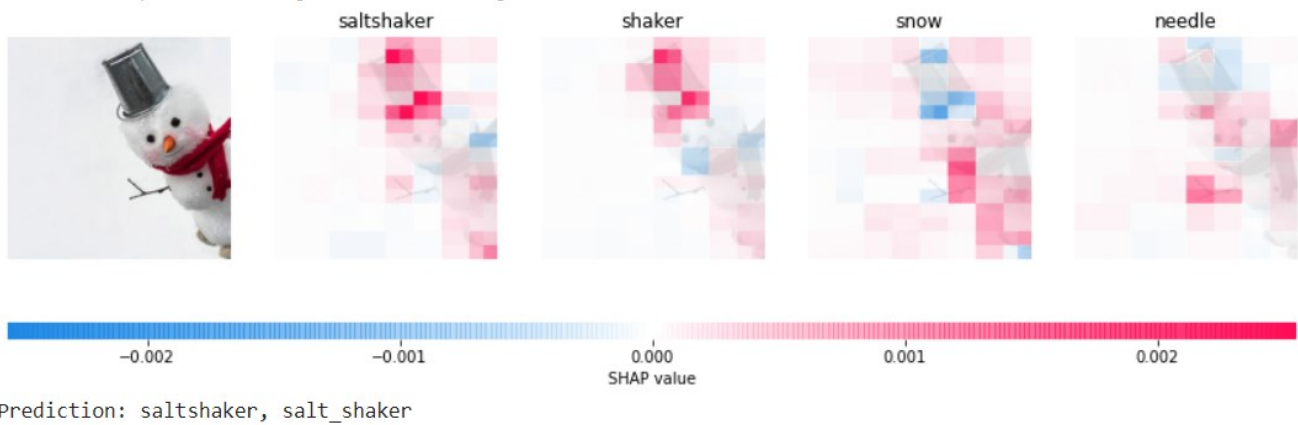


Q13 (a) As it could be seen from the corresponding image, the faces are the most important part. The pixels on Obama's face in the class related to the president is red. Also, Bush's face on ex-president class seems to be the most important part as the pixels of that area are red.

(b) Probably, the curves of the mask under the person's eyes and her eyebrows form a shape that looks like a Domino mask. Also, the shape of the lower part of her face has not changed after wearing the mask. The color of this mask might be another reason.



(c) As the body of the snowman is white and its hat is silver, the model misidentified it as a saltshaker. The pixels of the hat are red, so they played an important role and the hat is like a saltshaker head.

- (d) The model suffers from gender bias, and has detected the Chancellor as a first lady. It is probably due to the fact that the number of the women in power is by far smaller than men, and the model has been trained on a dataset that was imbalanced in this manner.
- (e) when we want to train on a dataset, not only we should pay attention to the class balance, but we should also consider the diversity of examples in each group. For example, in the picture of our instructor which is labeled as eccentric and geeky, the glasses are important features. It means that in the data set, these classes were full of pictures of people with glasses, so it became a feature to determine the labels. Another example could be the number of each gender in different job labels, so the model won't misclassify based on the gender of each person.

Q2: (a) i - True

The learnt prediction function is: $\text{sign}(w^T \phi(x_i))$
 $= \text{sign}(\sum y_i \alpha_i \phi(x_i) \phi(x))$

So if α_i for i -th instance is zero, its contribution to the sum is 0

ii - False

Slack variable s larger than zero means that there are some points on the other side of the $+1$ (or -1) plane. So, there is no guarantee that we don't have a point that will even penetrate the decision boundary and in that case, it would be misclassified.

iii - True

In this case, we will choose the closest points of the groups to be the support vectors. The dataset is linearly separable, so this distance exists and there won't be other points in between.

iv - False

As we might have some misclassified points, it should be $\min_i \left| \frac{y_i w^T \phi(x_i)}{\|w\|} \right|$

Proof: $x \rightarrow$ nearest point

$\bar{x} \rightarrow$ on the separation line

$$\begin{aligned} \Rightarrow x = \bar{x} + \lambda w \quad \left. \begin{array}{l} \Rightarrow w^T x = 0 + \lambda w^T w \\ w^T \bar{x} = 0 \end{array} \right\} & \Rightarrow \lambda = \frac{w^T x}{w^T w} \Rightarrow |\lambda| = \frac{|w^T x|}{\|w\|^2} \end{aligned}$$

$$|\lambda| \cdot \|w\| = \frac{|w^T x|}{\|w\|} \xrightarrow{\min_{x \in \phi(x)}} \min_i \frac{|w^T \phi(x_i)|}{\|w\|} = \min_i \left| y_i \frac{w^T \phi(x_i)}{\|w\|} \right|$$

V- False

$$\frac{1}{\|w\|} \frac{x - p_{(x)}}{\|w\|} = \min_i \frac{1}{\|w\|}$$

In this case, the margins change and some points might get on the other side of these lines. So, we will have some new ξ 's and this will change the value of our function $\frac{w^T w}{4} + \sum_{i=1}^n C \xi_i$. So, the w might change and so will the distance we calculated in the previous part.

(b)

i - The left one belongs to second order kernel (it has a convex shape)

The middle one belongs to the linear kernel.

The right one belongs to the fifth order kernel.

As it could be seen from the pictures, using a higher order polynomial lead to having more support vectors and it won't perform well on new data points.

ii - The left one $\rightarrow C = 200$ Increasing C means more penalty on ξ 's and
The right one $\rightarrow C = 2$ hence, less number of support vectors. (which is the figure A case)
A small C leads to overfitting. It captures all datapoints and won't generalize to test points.

iii - The RBF kernel formula is $\exp\left(-\frac{\|x-z\|^2}{2\sigma^2}\right) = \exp(-\gamma \|x-z\|^2)$; $\gamma = \frac{1}{2\sigma^2}$

The parameter γ defines the influence of a datapoint ; If it is low, then σ^2 will be high and the influence will reach far distances. So, it is proportional to the inverse of the radius of influence of samples. So we have:

The left figure $\rightarrow \gamma = 0.0025$

The middle figure $\rightarrow \gamma = 400$

The right figure $\rightarrow \gamma = 1$

A large gamma leads to overfitting, as the points should be closer together.

$$(C) \quad i - y_i f(x_i) = y_i (w^T \phi(x_i)) = y_i \sum_j \alpha_j y_j e^{-\gamma \|x-x_j\|^2} = y_i \sum_P \alpha_j y_j e^{-\gamma \|x-x_j\|^2} + y_i \sum_D \alpha_j y_j e^{-\gamma \|x-x_j\|^2}$$

$$= + \sum_P \alpha_j e^{-\gamma \|x-x_j\|^2} - \sum_D \alpha_j e^{-\gamma \|x-x_j\|^2}$$

We also have:

$$\begin{cases} \|x-x_j\| \leq s_1 & \text{for } x \in P \rightarrow e^{-\gamma \|x-x_j\|^2} \geq e^{-\gamma s_1^2} \\ \|x-x_j\| \geq s_2 & \text{for } x \in D \rightarrow e^{-\gamma \|x-x_j\|^2} \leq e^{-\gamma s_2^2} \\ -y_j^2 \geq -1 \end{cases}$$

$$\Rightarrow y_i f(x_i) - y_j^2 \geq \sum_{j \in P} \alpha_j e^{-\gamma s_1^2} - \sum_{j \in D} \alpha_j e^{-\gamma s_2^2} - 1 \quad \checkmark$$

ii - From the right handside of the above equation, we have:

$$\sum_{j \in P} \alpha_j e^{-\gamma s_1^2} - \sum_{j \in D} \alpha_j e^{-\gamma s_2^2} - 1 + \left(\sum_{j \in P} \alpha_j e^{-\gamma s_2^2} - \sum_{j \in P} \alpha_j e^{-\gamma s_2^2} \right)$$

$$= \sum_{j \in P} \alpha_j (e^{-\gamma s_1^2} - e^{-\gamma s_2^2}) + \left(\sum_{j \in P} \alpha_j - \sum_{j \in D} \alpha_j \right) e^{-\gamma s_2^2} - 1$$

$$\textcircled{1} \text{ we have } \sum_{i=1}^m \alpha_i y_i = 0 \rightarrow \sum_P \alpha_j y_j + \sum_D \alpha_j y_j = 0 \Rightarrow \sum_P \alpha_j - \sum_D \alpha_j = 0$$

\downarrow $+1 (-1)$ \downarrow $-1 (1)$

$$\textcircled{2} 0 \leq \alpha_j \leq C \rightarrow \sum_P \alpha_j \leq |P| \cdot C \leq mC \rightarrow -\sum_D \alpha_j \geq -mC \quad (*)$$

$$\textcircled{3} \text{ for convex functions we have: } f(x) - f(y) \geq f'(y)(x-y) \rightarrow \begin{cases} x := s_1^2 \rightarrow f(x) = e^{-rx} \\ y := s_2^2 \rightarrow f(y) = e^{-ry} \rightarrow f'(y) = -re^{-ry} \end{cases}$$

$$\Rightarrow e^{-rx} - e^{-ry} \geq -re^{-ry}(x-y) \Rightarrow e^{-rs_1^2} - e^{-rs_2^2} \geq -re^{-rs_2^2}(s_1^2 - s_2^2)$$

$$\Rightarrow y_i (f(x_i) - y_i) \geq \sum_P \alpha_j (e^{-rs_1^2} - e^{-rs_2^2}) - 1 \geq \sum_P \alpha_j (-re^{-rs_2^2}(s_1^2 - s_2^2)) - 1$$

$$\text{we also have } rs_2^2 \geq 0 \Rightarrow -e^{-rs_2^2} \geq -1 \rightarrow y_i (f(x_i) - y_i) \geq \sum_P \alpha_j (-r(s_1^2 - s_2^2)) - 1$$

$$\textcircled{*} \rightarrow y_i (f(x_i) - y_i) \geq -mC r(s_1^2 - s_2^2) - 1 = Cm r(s_2^2 - s_1^2) - 1 \quad \checkmark$$

$$\text{iii - For } y_i = 1 \rightarrow \begin{cases} f(x_i) > 0 & \Leftrightarrow \text{correctly classified} \\ \underline{f(x_i) > 1} & \Leftrightarrow \text{not misclassified outside the margin} \end{cases} \quad \textcircled{I}$$

$$\text{For } y_i = -1 \Rightarrow \begin{cases} f(x_i) < 0 & \Leftrightarrow \text{correctly classified} \\ \underline{f(x_i) < -1} & \Leftrightarrow \text{not misclassified outside the margin} \end{cases} \quad \textcircled{II}$$

$$\left. \begin{array}{l} \textcircled{I} \rightarrow y_i f(x_i) > 1 \\ \textcircled{II} \rightarrow y_i f(x_i) < -1 \end{array} \right\} \Rightarrow \text{The condition for not being misclassified outside of the margin is } y_i f(x_i) > 1$$

$$\Rightarrow y_i f(x_i) - y_i^2 > 1 - 1 = 0$$

$$\Rightarrow Cm r(s_2^2 - s_1^2) - 1 > 0 \xrightarrow{s_1 > s_2} Cr < \frac{1}{m(s_2^2 - s_1^2)}$$

$$\rightarrow Cr < \frac{1}{m(s_1^2 - s_2^2)}$$

Q4: a) i) Here, we define a loss function $J(x; w)$ such that maximizing this function leads to a weak performance by the classifier. However, we also don't want to change the input hugely; so we should constraint the perturbation.

By using a first-order approximation, we can generalize this loss function:

$$J(x; w) = J(x_0 + r; w) \approx J(x_0; w) + \nabla_x J(x_0; w)^T r$$

To formulate our optimization problem, we should find this r (which is $\leq \eta$) so it will maximize $J(x; w)$

$$\underset{r}{\text{maximize}} \quad J(x_0; w) + \nabla_x J(x_0; w)^T r \quad \text{s.t.} \quad \|r\|_\infty \leq \eta$$

ii) The above equation can be regarded as:

$$\underset{r}{\text{minimize}} \quad -J(x_0; w) - \nabla_x J(x_0; w)^T r \quad \text{s.t.} \quad \|r\|_\infty \leq \eta$$

we define: $w_0 = \nabla_x J(x_0; w)$ and $w^T r = -\nabla_x J(x_0; w)^T r$

$$\Rightarrow \underset{r}{\text{minimize}} \quad w^T r - w_0 \quad \text{s.t.} \quad \|r\|_\infty \leq \eta$$

Now, we use the Holder's inequality: $\langle w, r \rangle = |w^T r| \leq \|w\|_p \|r\|_q$ $\left\{ \begin{array}{l} p=1, q=\infty \\ \frac{1}{p} + \frac{1}{q} = 1 \end{array} \right\} \rightarrow |w^T r| \leq \|w\|_1 \|r\|_\infty$

$$\Rightarrow -\|w\|_1 \|r\|_\infty \leq w^T r \leq \|w\|_1 \|r\|_\infty$$

As we have $\|r\|_\infty \leq \eta \rightarrow w^T r \geq -\eta \|w\|_1$ if we choose $r = -\eta \text{sign}(w)$, we will reach the lower bound $(-\eta \|w\|_1)$

$$w^T r = -\eta w^T \text{sign}(w) = -\eta \sum_i w_i \text{sign}(w_i) = -\eta \sum_i |w_i| = -\eta \|w\|_1$$

\Rightarrow So, the answer to our problem is $r = -\eta \text{sign}(-\nabla_x J(x_0; w)) = \eta \text{sign}(\nabla_x J(x_0; w))$

$$\Rightarrow x = x_0 + r = x_0 + \eta \text{sign}(\nabla_x J(x_0; w))$$

b) $z = w_1 x + b$ $y = w_2 z = w_2 w_1 x + w_2 b$ $J = \sum_i (\hat{y}_i - y_i)^2$

$$\begin{aligned} \text{FGSM attack} \rightarrow x &= x_0 + \eta \text{sign}(\nabla_x (\sum_i (\hat{y}_i - w_2 w_1 x_i - w_2 b)^2)) \\ &= x_0 + \eta \text{sign}(2(-w_1 w_2)(\hat{y}_i - w_2 w_1 x_i - w_2 b)) \end{aligned}$$