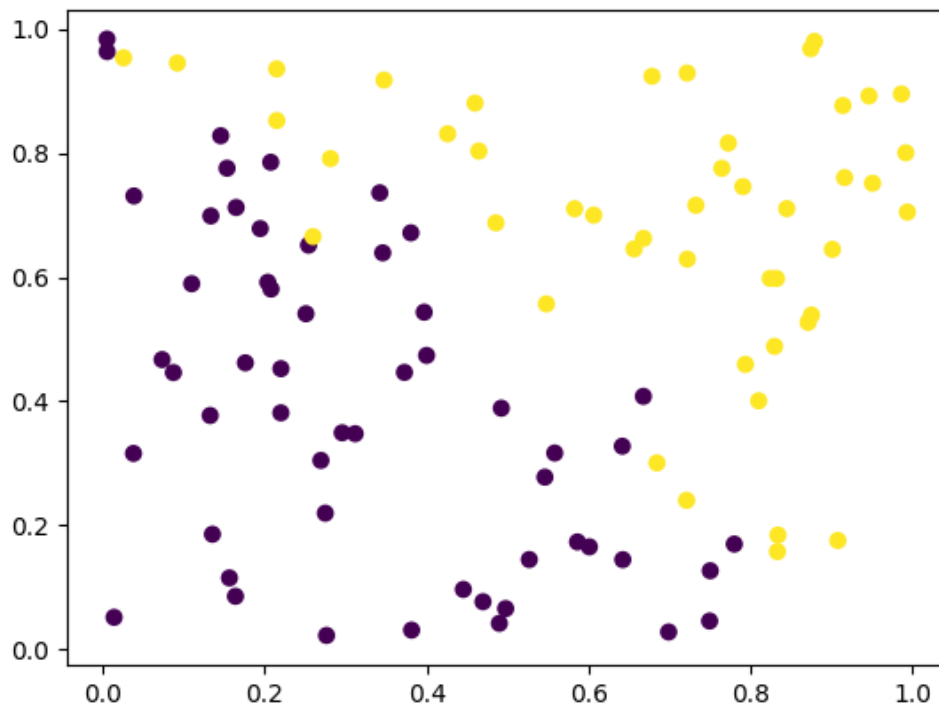
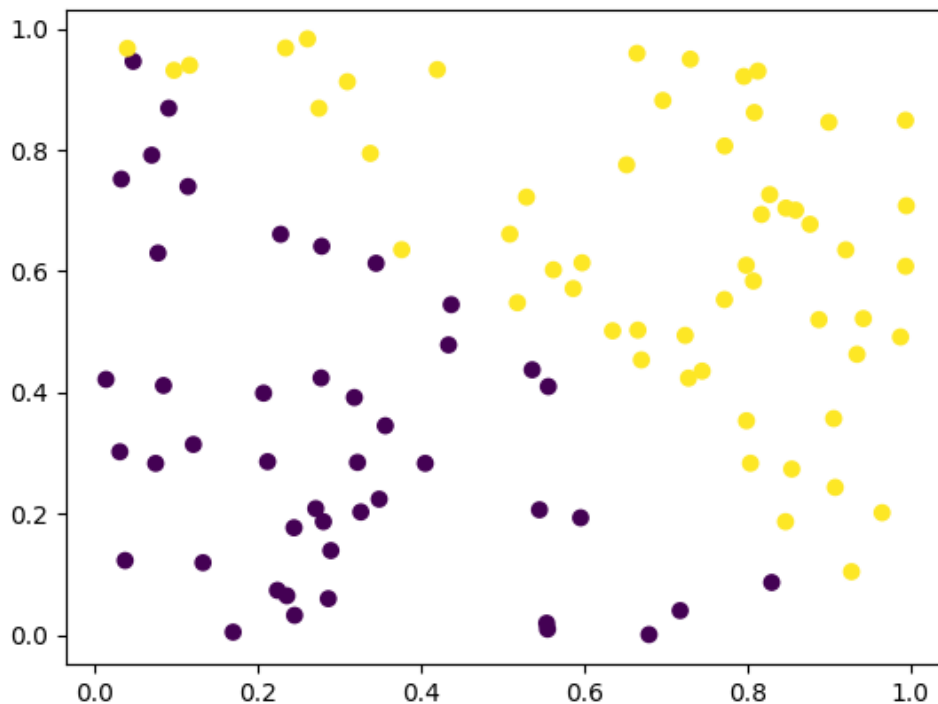


1)

a) La diferencia más notable a la hora de correr el código y entrenar el modelo es que en el caso del dataset A la convergencia llega rápidamente y en cambio cuando se entrena sobre el conjunto B este tarda demasiado en converger.

b) Lo realizado para intentar comprender esta diferencia fue graficar mediante gráficos de puntos ambos datasets. Al observar los gráficos lo más notable al momento de comparar es el hecho de que en el conjunto B podríamos clasificar casi perfectamente los datos con un correcto borde de decisión lineal. A continuación los gráficos de puntos:





Si observamos las figuras con detalle podríamos notar que en el dataset “b” a diferencia del dataset “a” se podría separar/clasificar casi perfectamente los datos con un correcto y justo borde de decisión lineal. Esta diferencia podría ser la causa por la cual la convergencia tarda tanto en el conjunto “b”.

c) Lo que determina el valor de la tasa de aprendizaje es que tan rápido puede convergerse hacia los valores óptimos de los parámetros del modelo. Lo que se establece es que un LR que es muy chico el modelo puede converger lento y tomarse más iteraciones para llegar a una solución óptima y por otra parte un LR grande puede hacer que el modelo no converja, ósea los parámetros nunca llegan a valores óptimos. En el caso para el trabajos se probó estableciendo la variable de “learnig rate” tanto con valores altos, bajos y otros varios y no fue solución para que el dataset “b”.

En lo que es escalado de datos, es importante ya que esto ayuda a que los modelos converjan más rápido y además de manera más estable. Si las características no están en la misma escala, algunas de estas pueden dominar sobre otras en términos de contribución a la función de coste, lo

que puede hacer que el proceso de optimización sea más lento. Se ha probado con escalar las Xs previo a realizar el entrenamiento pero tampoco tuvo resultados positivos en la velocidad de convergencia.

Al aplicar un término de regularización a la función de costo se penalizan los valores grandes de los coeficientes del modelo. Lo que se hizo fue aplicar dicho término a la función de costo como un valor lambda según vimos en clase y probando con varios valores se logró que el dataset “b” converja y a su vez que esto no afecte gravemente en el conjunto “a”.

d) No, las SVM en realidad son menos propensas a sufrir problemas de convergencia ya que son menos sensibles a la elección de los hiperparámetros o a la calidad de los datos (en comparación de la regresión logística). La función de pérdida de Hinge, utilizada por las SVM, están diseñadas para maximizar el margen entre las clases, lo que significa que se enfoca en los datos cercanos al borde de decisión que están clasificados incorrectamente.

2)

a) El desarrollo matemático está plasmado en la hoja. Partimos de la función de costo que ya conocemos de antes de haberla utilizado en Análisis multivariado y a lo largo de la cursada actual. Lo que hacemos con ella es buscar el gradiente que como ya sabemos es derivar a la función de costo con respecto a θ . Una vez obtenido el gradiente pasamos a realizar la distributiva de la sumatoria y nos quedan dos términos, por un lado tenemos a la sumatoria de la sigmoide y por otro a la sumatoria negativa de las probabilidades estimadas. Por último pasamos hacia el otro lado de la igualdad el segundo término para tenerlo positivo. Entonces tenemos la siguiente igualdad:

$$\sum_{i=1}^m h_{\theta}(x^{(i)}) = \sum_{i=1}^m y^{(i)}$$

Si tenemos en cuenta lo detallado en el enunciado y teoría previamente vista podemos decir por un lado que $h_{\theta}(x^{(i)}) = P(y^{(i)} = 1|x, \theta)$ y por

otro que $\sum_{i=1}^m y^{(i)} = \sum_{i=1}^m I(y^{(i)} = 1)$ y con esto damos por demostrada la propiedad.

Handwritten mathematical derivation on grid paper:

$$J(\theta) = -\frac{1}{m} \sum_{i=1}^m \left[y^{(i)} \log(h_{\theta}(x^{(i)})) + (1 - y^{(i)}) \log(1 - h_{\theta}(x^{(i)})) \right]$$

$$\frac{dJ(\theta)}{d\theta} = -\frac{1}{m} \sum_{i=1}^m (h_{\theta}(x^{(i)}) - y^{(i)}) = 0$$

$$\sum_{i=1}^m h_{\theta}(x^{(i)}) - \sum_{i=1}^m y^{(i)} = 0$$

$$\sum_{i=1}^m h_{\theta}(x^{(i)}) = \sum_{i=1}^m y^{(i)}$$

A green arrow points from the boxed equation to the following line:

$$\sum_{i=1}^m h_{\theta}(x^{(i)}) = P(y^{(i)} = 1 | x, \theta)$$

$$\sum_{i=1}^m y^{(i)} = \sum_{i=1}^m I(y^{(i)} = 1)$$

b) Esto no es necesariamente así ya que podemos plantear y destacar la diferencia entre que un modelo perfectamente calibrado y un modelo perfectamente preciso. Por el lado de la calibración según dice el ejercicio tenemos que cuando las probabilidades generadas por un modelo coinciden con la observación empírica estamos frente a un modelo bien calibrado. Es decir si el modelo predice una probabilidad del 70% alrededor del 70% de las muestras con esa probabilidad deberían clasificarse como positivas. Y por el lado de la precisión lo que sabemos es que tan bien el modelo clasifica cada muestra, por ejemplo si estamos ante una precisión perfecta el modelo clasifica correctamente todas las muestras sin cometer error alguno.

La relación que establecemos es que un modelo puede estar perfectamente calibrado sin necesariamente tener una precisión perfecta, es decir el modelo aún puede cometer errores de clasificación según el umbral de decisión que se elija, más allá de que las probabilidades predichas son precisas en términos de la probabilidad real. Y por otro lado, un modelo puede ser perfectamente preciso pero no estar correctamente calibrado, en

este caso el modelo clasificará todas las muestras correctamente pero sus probabilidades predichas no reflejan correctamente las probabilidades reales.

c) Según lo investigado la regularización L2 en Reg. Log. Es un término adicional que se le agrega a la función de costo. En este caso la función de costo se compone por dos partes, el término log-loss, que es la pérdida logarítmica, y el término de regularización L2.

$$J(\theta) = -\frac{1}{m} \sum_{i=1}^m \left[y^{(i)} \log(h_{\theta}(x^{(i)})) + (1 - y^{(i)}) \log(1 - h_{\theta}(x^{(i)})) \right] + \frac{\lambda}{2m} \sum_{j=1}^n \theta_j^2$$

Al incluir esta modificación, la regularización L2 puede afectar de varias maneras a la calibración. Puede aportar mayor estabilidad ya que L2 tiende a reducir la varianza del modelo al no tener en cuenta a los coeficientes grandes. También mejora cuando estamos frente a casos de overfitting, es decir cuando el modelo tiende a ajustarse demasiado a los datos de entrenamiento L2 ayuda a que el modelo se confíe en sus predicciones. Cabe destacar que estas mejoras de L2 están atadas a la elección del parámetro lambda ya que si este es grande el modelo puede volverse más conservador y perder algo de precisión y en el caso contrario donde es pequeño el modelo puede tender a sobre ajustar (hacer overfitting) y perder calibración.

3)

a)

3) a)

TEOREMA DE BAYES $= P(\theta|x) = \frac{P(y|x, \theta) \cdot P(\theta)}{P(y|x)}$

$\theta_{\text{MAP}} = \underset{\theta}{\text{arg max}} \frac{P(y|x, \theta) \cdot P(\theta)}{P(y|x)}$

$\theta_{\text{MAP}} = \underset{\theta}{\text{arg max}} P(y|x, \theta) \cdot P(\theta) \rightarrow P(\theta) = P(\theta|x)$

$\theta_{\text{MAP}} = \underset{\theta}{\text{arg max}} P(y|x, \theta) \cdot P(\theta|x)$

$\underset{\theta}{\text{arg max}} P(\theta|x, y) = \underset{\theta}{\text{arg max}} P(y|x, \theta) \cdot P(\theta|x)$

Lo que hacemos aquí es partir del teorema de Bayes, luego ignoramos o apartamos el denominador ya que no contiene ni depende de theta. Luego aprovechando lo detallado en la consigna, remplazamos al prior por la probabilidad de theta dado X.

b)

3) b)

$\theta_{\text{MAP}} = P(y|x, \theta) \cdot P(\theta) \rightarrow P(\theta) = \exp\left(-\frac{\|\theta\|_2^2}{2\gamma^2}\right)$

$\theta_{\text{MAP}} = P(y|x, \theta) \cdot \exp\left(-\frac{1}{2\gamma^2} \|\theta\|_2^2\right)$

$\ln(\theta_{\text{MAP}}) = \ln(P(y|x, \theta)) - \frac{1}{2\gamma^2} \|\theta\|_2^2$

$\ln(\theta_{\text{MAP}}) = \ln(P(y|x, \theta)) - \lambda \|\theta\|_2^2$

$\theta_{\text{MAP}} = \underset{\theta}{\text{ARG MIN}} (-\ln[P(y|x, \theta)] + \lambda \|\theta\|_2^2)$

Partiendo de lo hecho en el inciso (a), continuamos tomando al prior como:

$\exp(-\frac{\|\theta\|^2}{2\eta^2})$ y remplazamos en la formula. Lo que hacemos a

continuación es aplicar LN y luego separamos en suma aplicando la

propiedad logarítmica. Ahora por último si tomamos $-\frac{1}{2\eta^2} = -\lambda$

obtendremos la expresión que queremos buscar que la misma puede ser escrita o reformulada como se ve en la hoja.

c)

Handwritten mathematical derivation on grid paper:

c)

$$P(\theta | y, x) = P(y | x, \theta) \cdot P(\theta)$$

$$P(y | x, \theta) = \frac{1}{\sqrt{2\pi\sigma^2}} \exp\left(-\frac{(y - \theta^T \cdot x)^2}{2\sigma^2}\right)$$

$$P(\theta) = \frac{1}{\sqrt{2\pi\eta^2}} \exp\left(-\frac{1}{2\eta^2} \cdot \theta^T \cdot \theta\right)$$

$$P(\theta | y, x) = \frac{1}{\sqrt{2\pi\sigma^2}} \exp\left(-\frac{(y - \theta^T \cdot x)^2}{2\sigma^2}\right) \cdot \frac{1}{\sqrt{2\pi\eta^2}} \exp\left(-\frac{1}{2\eta^2} \cdot \theta^T \cdot \theta\right)$$

$$\log(P(\theta | y, x)) = \frac{-(y - \theta^T \cdot x)^2}{2\sigma^2} - \frac{\theta^T \cdot \theta}{2\eta^2} - \underbrace{\frac{1}{2} \ln(2\pi\sigma^2) - \frac{D}{2} \ln(2\pi\eta^2)}_{\text{CONSTANTE}}$$

$$\nabla L_m(P(\theta|y, x)) = \nabla \left(\frac{(y - \theta^T x)^2}{2\sigma^2} - \frac{\theta^T \cdot \theta}{2\lambda^2} \right)$$

$$\nabla \left(\frac{(y - \theta^T x)^2}{2\sigma^2} \right) = \frac{1}{2\sigma^2} \nabla (y - \theta^T x)^2$$

$$= \frac{1}{2\sigma^2} \cdot 2x(y - \theta^T x)$$

$$= \frac{x(y - \theta^T x)}{\sigma^2}$$

$$\nabla \left(\frac{\theta^T \cdot \theta}{2\lambda^2} \right) = \frac{1}{2\lambda^2} \nabla (\theta^T \cdot \theta)$$

$$= \frac{1}{2\lambda^2} \cdot 2\theta$$

$$= \frac{\theta}{\lambda^2}$$

$$\nabla L_m(P(\theta|y, x)) = \frac{x(y - \theta^T x)}{\sigma^2} - \frac{\theta}{\lambda^2} = 0$$

$$\theta_{MAP} = x(y - \theta^T x) = 0$$

Investigando en internet se encontró que podemos comenzar por remplazar los valores del prior y de la probabilidad de “y” dado X y theta. Al tener esos valores lo que hacemos es remplazar en la fórmula de $P(\theta|x, y)$. Ahora para encontrar la función cerrada nos conviene utilizar el logaritmo natural ya que estamos en una distribución gaussiana. Aplicamos y luego de obtener el desarrollo ignoramos todos los términos constantes. Lo último que queda por hacer es derivar lo obtenido. Separamos en términos y derivamos. Luego de obtener la derivada despejamos Theta y obtenemos lo resaltado en la hoja.

d) ----

4)

b) Para calcular el factor de compresión se utilizaron los bits de cada imagen. Para calcular los bits, se obtuvo la resolución de la imagen y se lo multiplica por la cantidad de bits por pixel (y por canal de color), es decir, $24(8 \times 3)$. Luego, los bits de la imagen comprimida es la resolución por el logaritmo en base 2 de 16 (lo que calcula los bits por pixel teniendo en cuenta esa cantidad de colores), es decir, 4 bits por pixel. Al obtener esta cuenta, se divide la cantidad de bits de la imagen original por la cantidad de bits de la imagen comprimida, y se llega a que el factor de compresión al disminuir la cantidad de colores a 16 en total es de 6.