

Tarea 3

1. Considere un problema de regresión lineal en el que queremos ponderar de forma distinta cada ejemplo de entrenamiento, tal y como vimos en clase con el método de regresión ponderada localmente. Específicamente, queremos minimizar

$$J(\underline{\theta}) = \frac{1}{2} \sum_{i=1}^m w^{(i)} (\underline{\theta}^T \underline{x}^{(i)} - y^{(i)})^2$$

- (a) Demuestre que para el caso general $J(\underline{\theta})$ se puede reescribir como

$$J(\underline{\theta}) = \frac{1}{2} (\underline{X}\underline{\theta} - \underline{y})^T W (\underline{X}\underline{\theta} - \underline{y})$$

para una matriz diagonal W apropiada donde \underline{X} es la matriz de diseño e \underline{y} el vector de salidas, tal y como lo definimos en clase.

Solución:

Para esta parte lo que se hará es desarrollar la expresión de $J(\underline{\theta})$ con \underline{X} y W como:

$$\underline{X} = \begin{bmatrix} X_{1,1} & X_{1,2} & \cdots & X_{1,n} \\ X_{2,1} & X_{2,2} & \cdots & X_{2,n} \\ \vdots & \vdots & \ddots & \vdots \\ X_{n,1} & X_{n,2} & \cdots & X_{n,n} \end{bmatrix}$$

donde, cada fila de la matriz representa un vector de valores de entrada.

$$W = \begin{bmatrix} W_{1,1} & 0 & \cdots & 0 \\ 0 & W_{2,2} & \cdots & 0 \\ \vdots & \vdots & \ddots & \vdots \\ 0 & 0 & \cdots & W_{n,n} \end{bmatrix}$$

Para el desarrollo $W_{1,1} = W_1, W_{2,2} = W_2, \dots, W_{n,n} = W_n$.

Desarrollando:

$$J(\underline{\theta}) = \frac{1}{2} \left(\begin{bmatrix} X_{1,:} \\ X_{2,:} \\ \vdots \\ X_{n,:} \end{bmatrix} \begin{bmatrix} \theta_1 \\ \theta_2 \\ \vdots \\ \theta_n \end{bmatrix} - \begin{bmatrix} y_1 \\ y_2 \\ \vdots \\ y_n \end{bmatrix} \right)^T \begin{bmatrix} W_1 & 0 & \cdots & 0 \\ 0 & W_2 & \cdots & 0 \\ \vdots & \vdots & \ddots & \vdots \\ 0 & 0 & \cdots & W_n \end{bmatrix} \left(\begin{bmatrix} X_{1,:} \\ X_{2,:} \\ \vdots \\ X_{n,:} \end{bmatrix} \begin{bmatrix} \theta_1 \\ \theta_2 \\ \vdots \\ \theta_n \end{bmatrix} - \begin{bmatrix} y_1 \\ y_2 \\ \vdots \\ y_n \end{bmatrix} \right)$$

$$J(\underline{\theta}) = \frac{1}{2} \left(\begin{bmatrix} X_{1,:} \cdot \underline{\theta} - y_1 \\ X_{2,:} \cdot \underline{\theta} - y_2 \\ \vdots \\ X_{n,:} \cdot \underline{\theta} - y_n \end{bmatrix} \right)^T \begin{bmatrix} W_1 & 0 & \cdots & 0 \\ 0 & W_2 & \cdots & 0 \\ \vdots & \vdots & \ddots & \vdots \\ 0 & 0 & \cdots & W_n \end{bmatrix} \left(\begin{bmatrix} X_{1,:} \cdot \underline{\theta} - y_1 \\ X_{2,:} \cdot \underline{\theta} - y_2 \\ \vdots \\ X_{n,:} \cdot \underline{\theta} - y_n \end{bmatrix} \right)$$

$$J(\underline{\theta}) = \frac{1}{2} \left(\begin{bmatrix} W_1(X_{1,:} \cdot \underline{\theta} - y_1) \\ W_2(X_{2,:} \cdot \underline{\theta} - y_2) \\ \vdots \\ W_n(X_{n,:} \cdot \underline{\theta} - y_n) \end{bmatrix} \right)^T \left(\begin{bmatrix} X_{1,:} \cdot \underline{\theta} - y_1 \\ X_{2,:} \cdot \underline{\theta} - y_2 \\ \vdots \\ X_{n,:} \cdot \underline{\theta} - y_n \end{bmatrix} \right)$$

$$J(\underline{\theta}) = \frac{1}{2} [W_1(X_{1,:} \cdot \underline{\theta} - y_1)^2 + W_2(X_{2,:} \cdot \underline{\theta} - y_2)^2 + \dots + W_n(X_{n,:} \cdot \underline{\theta} - y_n)^2]$$

$$\therefore J(\underline{\theta}) = \frac{1}{2} \sum_{i=1}^n W_i(X_{i,:} \cdot \underline{\theta} - y_i)^2$$

- (b) Si todos los pesos $w^{(i)}$ son iguales a 1, entonces en clase vimos que la ecuación normal es simplemente

$$X^T X \underline{\theta} = X^T \underline{y}$$

y el valor de $\underline{\theta}$ que minimiza $J(\underline{\theta})$ está dado por $(X^T X)^{-1} X^T \underline{y}$.

Encuentre el gradiente $\nabla_{\underline{\theta}} J(\underline{\theta})$ e igualelo a cero para encontrar una versión generalizada de las ecuaciones normales en este contexto con ponderación. Encuentre una forma cerrada del valor de $\underline{\theta}$ que minimiza a $J(\underline{\theta})$, en función de X , W e \underline{y} .

Solución:

El mínimo se encuentra buscando $\nabla_{\underline{\theta}} J(\underline{\theta}) = 0$

$$\nabla_{\underline{\theta}} J(\underline{\theta}) = \frac{1}{2} \nabla_{\underline{\theta}} (X \underline{\theta} - \underline{y})^T W (X \underline{\theta} - \underline{y}) = 0$$

$$\nabla_{\underline{\theta}} J(\underline{\theta}) = \nabla_{\underline{\theta}} \frac{1}{2} (\underline{\theta}^T X^T W X \underline{\theta} - \underline{\theta}^T X^T W \underline{y} - \underline{y}^T W X \underline{\theta} + \underline{y}^T W \underline{y}) = 0$$

Puesto que la expresión entre paréntesis es un escalar real

$$\nabla_{\underline{\theta}} J(\underline{\theta}) = \nabla_{\underline{\theta}} \frac{1}{2} \text{tr}(\underline{\theta}^T X^T W X \underline{\theta} - \underline{\theta}^T X^T W \underline{y} - \underline{y}^T W X \underline{\theta} + \underline{y}^T W \underline{y})$$

$$\nabla_{\underline{\theta}} J(\underline{\theta}) = \nabla_{\underline{\theta}} \frac{1}{2} (\text{tr}(\underline{\theta}^T X^T W X \underline{\theta}) - 2 \text{tr}(\underline{y}^T W X \underline{\theta}))$$

y utilizando $\nabla_{A^T} A B A^T C = B^T A^T C^T + B A^T C$, con $A = \underline{\theta}^T$, $B = X X^T$ y $C = W$

$$\nabla_{\underline{\theta}} J(\underline{\theta}) = \nabla_{\underline{\theta}} \frac{1}{2} (\text{tr}(\underline{\theta}^T X^T W X \underline{\theta}) - 2 \text{tr}(\underline{y}^T W X \underline{\theta}))$$

$$\nabla_{\underline{\theta}} J(\underline{\theta}) = \frac{1}{2} (X^T W X \underline{\theta} + X^T W X \underline{\theta} - 2 X^T W \underline{y})$$

$$\nabla_{\underline{\theta}} J(\underline{\theta}) = X^T W X \underline{\theta} - X^T W \underline{y}$$

Recordando que para minimizar $\underline{\theta}$ el gradiente $\Delta_{\underline{\theta}} J(\underline{\theta}) = \underline{0}$

$$X^T W X \underline{\theta} = X^T W \underline{y}$$

$$\therefore \underline{\theta} = (X^T W X)^{-1} X^T W \underline{y}$$

- (c) Suponga que tenemos un conjunto de entrenamiento $\{(\underline{x}^{(i)}, y^{(i)}); i = 1 \dots m\}$ de m ejemplos independientes, pero en el que los $y^{(i)}$ se observaron con varianzas distintas. Específicamente, suponga que

$$p(y^{(i)}, \underline{x}^{(i)}; \underline{\theta}) = \frac{1}{\sqrt{2\pi}\sigma^{(i)}} \exp\left(-\frac{(y^{(i)} - \underline{\theta}^T \underline{x}^{(i)})^2}{2(\sigma^{(i)})^2}\right)$$

o en otras palabras, $y^{(i)}$ tiene media $\underline{\theta}^T \underline{x}^{(i)}$ y varianza $(\sigma^{(i)})^2$, donde las $\sigma^{(i)}$ son constantes fijas, conocidas. Demuestre que encontrar el estimado de máxima verosimilitud de $\underline{\theta}$ se reduce a resolver un problema de regresión lineal ponderada. Establezca claramente que los $w^{(i)}$ se calculan en términos de $\sigma^{(i)}$.

Solución:

La máxima verosimilitud se basa en elegir $\underline{\theta}$ de modo que los datos sean lo más probables posibles.

Para esto, se utiliza el logaritmo natural (\ln), es decir, verosimilitud logarítmica.

$$\ell(\underline{\theta}) = \ln L(\underline{\theta})$$

$$\ell(\underline{\theta}) = \ln \prod_{i=1}^m \frac{1}{\sqrt{2\pi}\sigma^{(i)}} \exp\left(-\frac{(y^{(i)} - \underline{\theta}^T \underline{x}^{(i)})^2}{2(\sigma^{(i)})^2}\right)$$

$$\ell(\underline{\theta}) = \sum_{i=1}^m \ln \frac{1}{\sqrt{2\pi}\sigma^{(i)}} \exp\left(-\frac{(y^{(i)} - \underline{\theta}^T \underline{x}^{(i)})^2}{2(\sigma^{(i)})^2}\right)$$

$$\ell(\underline{\theta}) = \ln \frac{1}{\sqrt{2\pi}\sigma^{(i)}} - \sum_{i=1}^m \frac{1}{\sigma^{(i)2}} \frac{1}{2} (y^{(i)} - \underline{\theta}^T \underline{x})^2$$

Donde el primer término es una constante, sacando el $\frac{1}{2}$ por distributividad y haciendo a $w^{(i)} = \frac{1}{\sigma^{(i)2}}$ se puede observar que el segundo término tiene la forma de la función de error en una RL ponderada por lo que se puede concluir que para maximizar la verosimilitud se requiere minimizar "la función de error" al igual que la misma RL.

2. Visualización de los datos

- (a) Use las ecuaciones normales para implementar la regresión lineal (no ponderada) $y = \underline{\theta}^T \underline{x}$ en el primer ejemplo de entrenamiento (esto es, la primera fila que no es el encabezado). En una figura, grafique tanto los datos crudos como la línea recta resultante del ajuste. Indique el $\underline{\theta}$ óptimo resultante de la regresión lineal

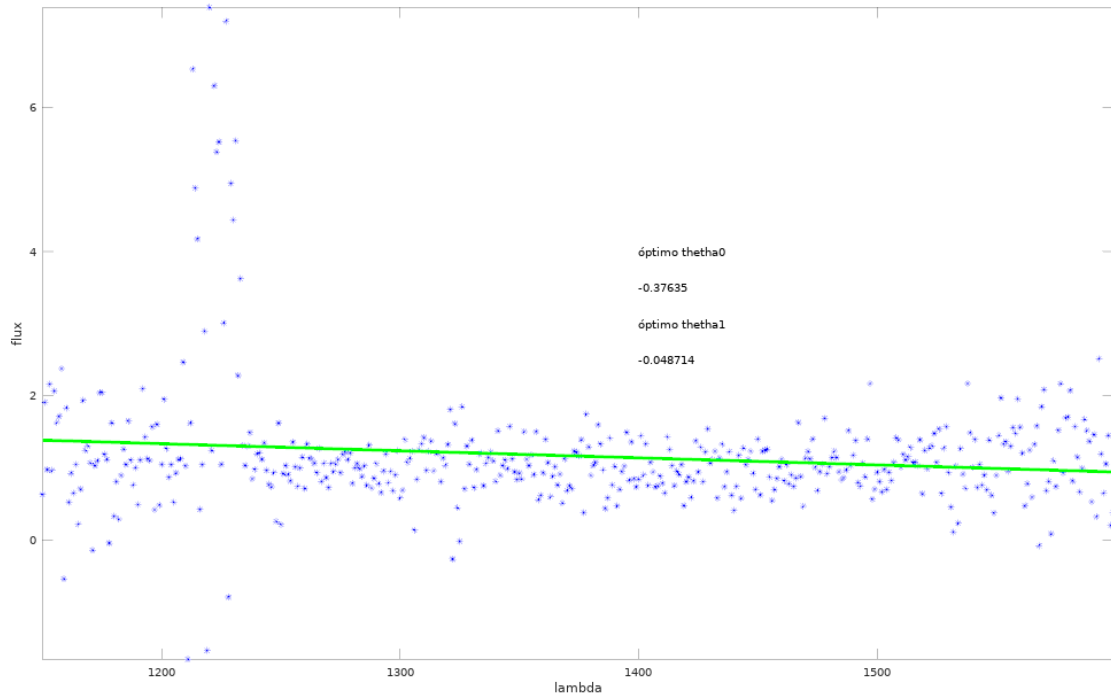


Figura 1: Regresión lineal no ponderada.

- (b) Implemente la regresión lineal ponderada localmente en el primer ejemplo de entrenamiento. Use las ecuaciones normales que usted derivó en el punto 1.2. En una figura aparte, grafique tanto los datos crudos como la curva suave resultante de su regresión. Cuando evalúe $h(\cdot)$ en un punto \underline{x} use los pesos

$$w^{(i)} = \exp\left(-\frac{\|\underline{x} - \underline{x}^{(i)}\|^2}{2\tau^2}\right)$$

con el parámetro de ancho de banda $\tau = 5$

- (c) Repita el punto 2.2 cuatro veces más con $\tau = 1, 10, 100$ y 1000 . Grafique las curvas resultantes en una misma figura. Indique en una frase corta que ocurre a la curva de regresión conforme τ crece.

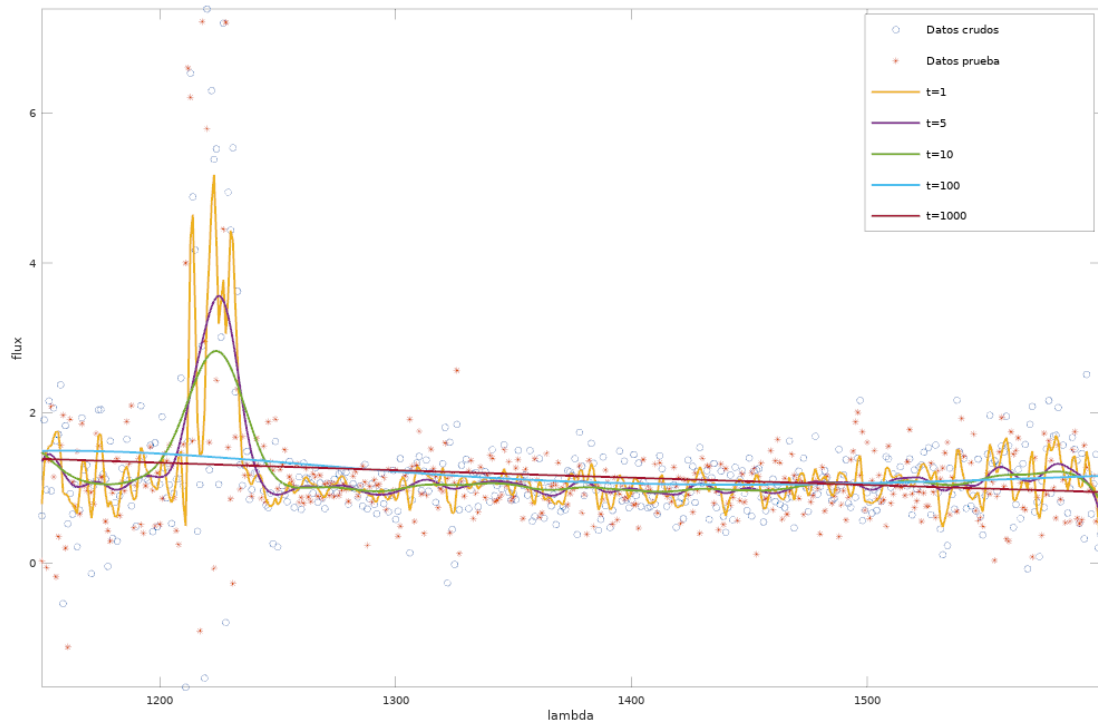


Figura 2: Regresión lineal ponderada para distintos valores de tau.

En la figura se puede notar que conforme el valor de tau disminuye, el modelo se aproxima a los datos de entrenamiento, mientras que si aumenta, la hipótesis se aproxima a la hipótesis de regresión lineal, por lo tanto, a valores altos de tau se genera subajuste y a valores pequeños provocan un sobreajuste de la hipótesis.