

Problemas Convexos

Carlos Malanche

20 de marzo de 2018

1. Máxima verosimilitud (*Maximum Likelihood*)

Podemos ver el problema de una manera alterna y un tanto complicada, es una manera de ver las cosas basados en probabilidad.

Como ya lo habíamos dicho antes, podemos suponer que nuestra información fue generada por una parte determinística más una parte de ruido

$$y = f(\mathbf{x}) + \epsilon \quad (1)$$

Con una restricción más, supongamos que la función determinística es una función lineal, tal y como lo hicimos con mínimos cuadrados

$$y = \mathbf{x}^T \mathbf{c} + \epsilon \quad (2)$$

Dado un \mathbf{c} fijo, podemos preguntarnos cuál es la probabilidad de haber obtenido las mediciones y .

$$P(y|\mathbf{x}, \mathbf{c}) = \mathcal{N}(y|\mathbf{x}^T \mathbf{c}, \sigma^2) = \frac{1}{\sqrt{2\pi\sigma^2}} e^{-\frac{(y-\mathbf{x}^T \mathbf{c})^2}{2\sigma^2}} \quad (3)$$

En el caso de \mathbf{y} un vector, la expresión es un tanto más complicada

$$P(\mathbf{y}|\mathbf{X}, \mathbf{c}) = \mathcal{N}(\mathbf{y}|\mathbf{X}\mathbf{c}, \Sigma) = \frac{1}{\sqrt{(2\pi)^m \det(\Sigma)}} \exp\left[-\frac{1}{2}(\mathbf{y} - \mathbf{X}\mathbf{c})^T \Sigma^{-1}(\mathbf{y} - \mathbf{X}\mathbf{c})\right] \quad (4)$$

La idea de *maximum likelihood* es encontrar el vector de parámetros \mathbf{c} que maximiza la probabilidad de haber obtenido el vector de *outcomes* \mathbf{y} .

Para resolver esto: ¿La probabilidad de obtener \mathbf{y} dado \mathbf{X} como función de \mathbf{c} qué forma tiene? ¿Es un problema convexo? ¡No! y tampoco es cóncavo. Primero vamos a arreglar lo primero. Hay una forma cuadrática en el argumento de la exponencial, la cual ya sabemos será convexa (excepto por el signo menos...). Vamos a intentar sacarla aplicando un logaritmo a la probabilidad.

$$\log(P(\mathbf{y}|\mathbf{X}, \mathbf{c})) = -\frac{1}{2}(\mathbf{y} - \mathbf{X}\mathbf{c})^T \Sigma^{-1}(\mathbf{y} - \mathbf{X}\mathbf{c}) + \text{cnst} \quad (5)$$

En donde hemos decidido abreviar esa constante

$$\text{cnst} = -\frac{1}{2}m\log((2\pi)^m \det(\Sigma))$$

Basta con multiplicar toda la expresión por un menos y deshacernos de la constante (no va a modificar el problema) para movernos a un marco muy similar al que estabamos utilizando anteriormente (un problema de minimizar en lugar de maximizar).

$$\arg \min_{\mathbf{c}} \left\{ \frac{1}{2}(\mathbf{y} - \mathbf{X}\mathbf{c})^T \Sigma^{-1}(\mathbf{y} - \mathbf{X}\mathbf{c}) \right\} \quad (6)$$

El problema es casi idéntico al de mínimos cuadrados! En realidad, es una generalización, pues basta con hacer $\Sigma = n\mathbb{I}$ para recuperar los mínimos cuadrados. Hacer esto le quitaría el sentido estadístico a nuestra interpretación así que no lo haremos.

En lugar de ello, todo se simplifica cuando la matriz de covarianza es diagonal, pues implica que cada una de las medidas fue independiente (lo cual es algo relativamente seguro de hacer en la mayoría de los casos). Incluso

podemos asumir que la desviación estándar σ es la misma para cada uno de los parámetros, con lo que obtenemos casi la función de pérdidas cuadrática

$$\arg \min_{\underline{c}} \left\{ \frac{1}{2} (\mathbf{y} - \mathbf{X}\underline{c})^T \Sigma^{-1} (\mathbf{y} - \mathbf{X}\underline{c}) \right\} \quad (7)$$

2. Regresión regularizada

Recordemos que el objetivo es encontrar un vector de parámetros (los cuales definen una función) que minimice el error de predicción sobre una serie de datos conocida. Si seguimos denotando con $J_{\hat{f}}(\underline{c})$ la función de costo adaptada a nuestra función estimadora \hat{f} , definiremos pues la función objetivo como

$$\mathcal{L}(\underline{c}) := J_{\hat{f}}(\underline{c}) \quad (8)$$

se ve un poco redundante pero tiene un propósito inventar esta nueva función. Lo que buscamos es minimizar la función objetivo

$$\arg \min_{\underline{c}} \mathcal{L}(\underline{c}) \quad (9)$$

Hemos visto que cuando J es la función de pérdidas cuadrática (FPC) la solución es analítica y está dada por las ecuaciones normales. Sin embargo, el problema puede volverse inestable dependiendo del rango de la matriz de observaciones \mathbf{X} (columnas casi linealmente independientes convierten en casi-singular a la matriz, y su inversión es inestable). Llegan a salvarnos los **regularizadores**, que son funciones que se añaden a la función objetivo para poner imponer condiciones.

$$\mathcal{L}(\underline{c}) := J_{\hat{f}}(\underline{c}) + \Omega(\underline{c}) \quad (10)$$

Esto es todo un campo de investigación, pues no basta con pasarle condiciones a la función objetivo a diestra y siniestra, no tiene caso hacerlo si no sabemos como resolver el problema. Un ejemplo de un *regularizador* es esta modificación a la función indicadora

$$\Omega(\underline{c}) = 1_S := \begin{cases} 0, & \text{si } \underline{c} \notin S \\ \infty, & \text{o.c.} \end{cases} \quad (11)$$

Basta con que la función de costo no sea infinita para una $\underline{c} \in S$ para que el vector solución \underline{c} no pueda estar fuera de S .

Por el momento nos vamos a enfocar en regularizadores con un comportamiento un poco menos agresivo (en general, derivables) que de todos modos nos pueden ayudar, en particular a que la matriz de las ecuaciones normales no quede mal condicionada.

2.1. Regularización de Tikhonov

Definamos los regularizadores de Tikhonov como aquellas función objetivo en donde el regularizador Ω tiene la forma

$$\Omega(\underline{c}) := \|\mathbf{\Gamma}\underline{c}\|^2 \quad (12)$$

Donde $\mathbf{\Gamma}$ es una matriz que podemos elegir para sacar características lineales del vector de parámetros. Busquen cuál es el resultado de encontrar el gradiente del problema con regularización de Tikhonov, es muy sencillo.

2.2. Regularización de arista

Caso particular, tomaremos

$$\Omega(\underline{c}) := \frac{\gamma}{2n} \|\underline{c}\|^2 \quad (13)$$

es decir, $\mathbf{\Gamma} = (\gamma/2n)^{1/2} \mathbb{I}$. Vamos a repetir el ejercicio de la clase pasada, y vamos a tratar de encontrar un punto óptimo para este caso

$$\nabla \mathcal{L}(\underline{c}) = \nabla \left(\frac{1}{2n} \|\mathbf{y} - \mathbf{X}\mathbf{c}\|^2 \right) + \frac{\gamma}{2n} \nabla \|\underline{c}\|^2 \quad (14)$$

$$= \frac{1}{n} \mathbf{X}^T (\mathbf{X}\mathbf{c} - \mathbf{y}) + \frac{\gamma}{n} \mathbf{c} \quad (15)$$

Ya sólo igualamos a cero y resolvemos para \mathbf{c}

$$\mathbf{c} = (\mathbf{X}^T \mathbf{X} + \gamma \mathbb{I})^{-1} \mathbf{X}^T \mathbf{y} \quad (16)$$

Para los atentos, el efecto de sumarle ese múltiplo de la identidad a una matriz semidefinida positiva es el de levantar los eigenvalores, con lo que podemos obligar a la matriz completa a ser invertible

Teorema. Sea \mathbf{G} una matriz de Gram generada por la matriz \mathbf{X} . Entonces los valores propios de $(\mathbf{G} + \gamma \mathbb{I})$ se elevan por γ

Demostración. Consideremos la diagonalización de la matriz de Gram, dada por las matriz de vectores propios \mathbf{V} y la matriz diagonal de valores propios \mathbf{P} .

$$\mathbf{G} = \mathbf{V}^T \mathbf{P} \mathbf{V}$$

Con ella, escribimos de nuevo nuestra matriz por invertir

$$\mathbf{X}^T \mathbf{X} + \gamma \mathbb{I} = \mathbf{V}^T \mathbf{P} \mathbf{V} + \gamma \mathbb{I} = \mathbf{V}^T \mathbf{P} \mathbf{V} + \gamma \mathbf{V}^T \mathbf{V} = \mathbf{V}^T (\mathbf{P} + \gamma \mathbb{I}) \mathbf{V} \quad (17)$$

Nota: Esta demostración se auxilió de saber que las matrices de vectores propios de una matriz de Gram son ortonormales, lo cual se puede ver con la descomposición de valores singulares de su matriz generadora. \square

Casi por diversión, de una manera alternativa podemos extraer los valores propios usando el coeficiente de *Rayleigh*

$$\frac{\mathbf{v}^T (\mathbf{X}^T \mathbf{X} + \gamma \mathbb{I}) \mathbf{v}}{\mathbf{v}^T \mathbf{v}} = \frac{\mathbf{v}^T \mathbf{X}^T \mathbf{X} \mathbf{v}}{\mathbf{v}^T \mathbf{v}} + \gamma \frac{\mathbf{v}^T \mathbf{v}}{\mathbf{v}^T \mathbf{v}} = \frac{\mathbf{v}^T \mathbf{X}^T \mathbf{X} \mathbf{v}}{\mathbf{v}^T \mathbf{v}} + \gamma$$

Añadimos el conocimiento de que esa matriz de Gram es semidefinida positiva, lo cual quiere decir que los valores propios son iguales o mayores a cero, con lo que podemos concluir que los valores propios de la matriz por invertir **son al menos** γ .

2.3. Regularización *Lasso*

Otro caso divertido es la regularización de *Lasso* (por sus siglas *least absolute shrinkage and selection operator*), en donde la función de regularización es la norma en L_1 , es decir

$$\Omega(\underline{c}) = \|\underline{c}\|_{L_1}$$

léase, la suma de los parámetros que definen al modelo. Lo curioso de este regularizador es que favorece soluciones con muchas entradas en cero. Piensen en un caso simple de dos dimensiones, y se podrán imaginar por qué pasa esto.

3. Problemas convexos

Nosotros agarramos y minimizamos el problema asumiendo que era convexo pero no lo probamos. Pues bien, definamos lo que es un problema convexo

Definición. Una función $f : U \rightarrow \mathbb{R}$ definida en un conjunto convexo y abierto es convexa si para cualesquiera dos puntos $\mathbf{s}_1, \mathbf{s}_2 \in U \subset \mathbb{R}^m$ y $\lambda \in [0, 1]$ se cumple que

$$f(\lambda \mathbf{s}_1 + (1 - \lambda) \mathbf{s}_2) \leq \lambda f(\mathbf{s}_1) + (1 - \lambda) f(\mathbf{s}_2)$$

La definición es muy intuitiva, pues nos dice que el resultado de evaluar f en cualquiera de los puntos sobre la recta que une un par de puntos $\underline{s}_2, \underline{s}_2 \in S$ (donde S es el dominio de la función) es inferior al resultado de evaluar el mismo punto en la línea que une $f(\underline{s}_1)$ con $f(\underline{s}_2)$.

Para el caso de mínimos cuadrados, se tendría que verificar que la siguiente desigualdad es cierta

$$\frac{1}{2n} \|\mathbf{X}(\lambda \mathbf{s}_1 + (1 - \lambda) \mathbf{s}_2) - \mathbf{y}\|^2 \leq \lambda \frac{1}{2n} \|\mathbf{X} \mathbf{s}_1 - \mathbf{y}\|^2 + (1 - \lambda) \frac{1}{2n} \|\mathbf{X} \mathbf{s}_2 - \mathbf{y}\|^2 \quad (18)$$

Ese ejercicio queda al lector. Vamos a intentar llegar a una definición alterna.

$$\begin{aligned} f(\lambda \mathbf{s}_1 + (1 - \lambda) \mathbf{s}_2) &\leq \lambda f(\mathbf{s}_1) + (1 - \lambda) f(\mathbf{s}_2) \\ f(\lambda(\mathbf{s}_1 - \mathbf{s}_2) + \mathbf{s}_2) &\leq \lambda(f(\mathbf{s}_1) - f(\mathbf{s}_2)) + f(\mathbf{s}_2) \\ \frac{f(\lambda(\mathbf{s}_1 - \mathbf{s}_2) + \mathbf{s}_2) - f(\mathbf{s}_2)}{\lambda} &\leq f(\mathbf{s}_1) - f(\mathbf{s}_2) \end{aligned}$$

Sin pérdida de generalidad, vamos a tomar el límite cuando λ tienda a cero (ya suficiente libertad tenemos con permitir a \underline{s}_1 y a \underline{s}_2 ser cualesquiera dos puntos en el dominio)

$$\begin{aligned} \lim_{\lambda \rightarrow 0} \frac{f(\lambda(\mathbf{s}_1 - \mathbf{s}_2) + \mathbf{s}_2) - f(\mathbf{s}_2)}{\lambda} &\leq f(\mathbf{s}_1) - f(\mathbf{s}_2) \\ (\mathbf{s}_1 - \mathbf{s}_2)^T \nabla f(\mathbf{s}_2) &\leq f(\mathbf{s}_1) - f(\mathbf{s}_2) \end{aligned}$$

Donde hemos hecho uso de una de las definiciones de derivada direccional (una que no se encuentra normalizada). Haciendo un último despeje, tenemos esta definición alterna para funciones diferenciables:

Definición. Una función $f : \mathbb{R}^m \rightarrow \mathbb{R}$ diferenciable es convexa si para cualesquiera dos puntos $\mathbf{s}_1, \mathbf{s}_2 \in S \subset \mathbb{R}^m$ se cumple que

$$f(\mathbf{s}_1) - f(\mathbf{s}_2) \geq (\mathbf{s}_1 - \mathbf{s}_2)^T \nabla f(\mathbf{s}_2)$$