

Regularización y aprendizaje en línea

Lección 16

Dr. Pablo Alvarado Moya

CE5506 Introducción al reconocimiento de patrones
Área de Ingeniería en Computadores
Tecnológico de Costa Rica

II Semestre, 2019

Contenido

- 1 Enfoques de estimación de parámetros
 - Enfoque frecuentista
 - Enfoque Bayesiano
- 2 Funciones de error equivalentes
- 3 Aprendizaje en línea

Enfoque frecuentista

- Retomemos como ejemplo regresión lineal/logística
- Elegimos parámetros que maximizan la verosimilitud

$$\max_{\underline{\theta}} \prod_i p(y^{(i)} | \underline{\mathbf{x}}^{(i)}; \underline{\theta})$$

- Aquí, $\underline{\theta}$ **no** es una variable aleatoria, sino una constante desconocida
- Nuestra tarea es encontrar métodos para estimar ese vector constante desconocido $\underline{\theta}$
- Este es el enfoque usual en estadística **frecuentista**.

Enfoque Bayesiano

(1)

- En el enfoque estadístico **Bayesiano** $\underline{\theta}$ también es una variable aleatoria de valor desconocido
- Se especificaría una distribución **a priori** $p(\underline{\theta})$ que expresa las creencias previas sobre los parámetros (**antes** de ver cualquier dato)
- Por ejemplo: $p(\underline{\theta}) \sim \mathcal{N}(\underline{\mu}, \underline{\Sigma})$

Enfoque Bayesiano

(2)

- Dado el conjunto de entrenamiento $\mathcal{S} = \{(\underline{\mathbf{x}}^{(i)}, y^{(i)}) \mid i = 1 \dots m\}$ se puede calcular la distribución *a posteriori* de los parámetros:

$$\begin{aligned} p(\underline{\theta}|\mathcal{S}) &= \frac{p(\mathcal{S}|\underline{\theta})p(\underline{\theta})}{p(\mathcal{S})} \\ &= \frac{(\prod_{i=1}^m p(y^{(i)}|\underline{\mathbf{x}}^{(i)}, \underline{\theta})) p(\underline{\theta})}{\int_{\underline{\theta}} (\prod_{i=1}^m p(y^{(i)}|\underline{\mathbf{x}}^{(i)}, \underline{\theta})p(\underline{\theta})) d\underline{\theta}} \end{aligned}$$

(después de ver los datos de entrenamiento)

Eligiendo la verosimilitud

- A $p(y^{(i)}|\underline{\mathbf{x}}^{(i)},\underline{\boldsymbol{\theta}})$ lo determina el modelo que se usa para el aprendizaje
- Con regresión logística Bayesiana, puede elegirse:

$$p\left(y^{(i)}|\underline{\mathbf{x}}^{(i)},\underline{\boldsymbol{\theta}}\right) = h_{\underline{\boldsymbol{\theta}}}\left(\underline{\mathbf{x}}^{(i)}\right)^{y^{(i)}}\left(1 - h_{\underline{\boldsymbol{\theta}}}\left(\underline{\mathbf{x}}^{(i)}\right)\right)^{1-y^{(i)}}$$

con

$$h_{\underline{\boldsymbol{\theta}}}(\underline{\mathbf{x}}^{(i)}) = \frac{1}{1 + \exp(-\underline{\boldsymbol{\theta}}^T \underline{\mathbf{x}}^{(i)})}$$

Predicción

- Cuando nos piden predecir y para un nuevo \underline{x} calculamos la distribución *a posteriori*

$$p(y|\underline{x}, \mathcal{S}) = \int_{\underline{\theta}} p(y|\underline{x}, \underline{\theta}) p(\underline{\theta}|\mathcal{S}) d\underline{\theta}$$

- Si queremos predecir el valor esperado de y dado \underline{x} , entonces:

$$E[y|\underline{x}, \mathcal{S}] = \int_y y p(y|\underline{x}, \mathcal{S}) dy$$

- Estos pasos son una predicción completamente Bayesiana
- Calcular la distribución *a posteriori* $p(\underline{\theta}|\mathcal{S})$ es en general difícil, por la necesidad de integrar en $\underline{\theta}$ (usualmente de muchas dimensiones), y no se puede hacer de forma cerrada.

Aproximación de distribución *a posteriori*

(1)

- Por lo general aproximaremos la distribución *a posteriori* de $\underline{\theta}$
- Una opción usual es reemplazar la distribución anterior con una estimación con un solo punto.
- El estimado MAP (*maximum a posteriori*) de $\underline{\theta}$ es

$$\hat{\underline{\theta}}_{\text{MAP}} = \arg \max_{\underline{\theta}} p(\underline{\theta} | \mathcal{S}) = \arg \max_{\underline{\theta}} \left(\prod_{i=1}^m p(y^{(i)} | \mathbf{x}^{(i)}, \underline{\theta}) \right) p(\underline{\theta})$$

y para predecir se usa $h_{\hat{\underline{\theta}}_{\text{MAP}}}(\mathbf{x}) = g(\hat{\underline{\theta}}_{\text{MAP}}^T \mathbf{x})$

- Esto es muy parecido a máxima verosimilitud, excepto por el término *a priori* $p(\underline{\theta})$

Aproximación de distribución *a posteriori*

(2)

- En la práctica se usa $\underline{\theta} \sim \mathcal{N}(0, \tau^2 \mathbf{I})$
(por ejemplo, clasificación de texto con 50 000 parámetros)
- Con eso los parámetros $\hat{\underline{\theta}}_{\text{MAP}}$ usualmente tienen menor norma que los elegidos con ML, lo que los hace menos susceptibles a sobreajuste.
- Esto se usa con éxito en clasificación de textos

Función de error equivalente

- Surge la pregunta: si utilizamos MAP con la estrategia anterior, ¿cuál es la función de error que estamos minimizando?
- Resulta que la nueva función de error para regresión lineal es:

$$J(\underline{\theta}) = \sum_{i=1}^m \left| y^{(i)} - \underline{\theta}^T \underline{\mathbf{x}}^{(i)} \right|^2 + \lambda \|\underline{\theta}\|_2^2$$

- El término de regularización $\lambda \|\underline{\theta}\|_2^2$ fuerza pequeños valores de los parámetros, lo que se puede interpretar como una simplificación del modelo
- Es usual dejar a θ_0 fuera de la regularización
- Las ecuaciones normales son en este caso:

$$\underline{\theta} = (\mathbf{X}^T \mathbf{X} + \lambda \mathbf{I})^{-1} \mathbf{X}^T \mathbf{y}$$

Regresión logística con regularización

- En el caso de regresión logística, regularizamos con

$$\begin{aligned} J(\underline{\theta}) &= \ell(\underline{\theta}) - \lambda \|\underline{\theta}\|_2^2 \\ &= \left(\sum_{i=1}^m y^{(i)} \ln(h_{\underline{\theta}}(\underline{\mathbf{x}}^{(i)})) + (1 - y^{(i)}) \ln(1 - h_{\underline{\theta}}(\underline{\mathbf{x}}^{(i)})) \right) - \lambda \|\underline{\theta}\|_2^2 \end{aligned}$$

con gradiente

$$\nabla_{\underline{\theta}} J(\underline{\theta}) = \left[\sum_{i=1}^{\infty} \left(y^{(i)} - h_{\underline{\theta}}(\underline{\mathbf{x}}^{(i)}) \right) \underline{\mathbf{x}}^{(i)} \right] - 2\lambda \underline{\theta}$$

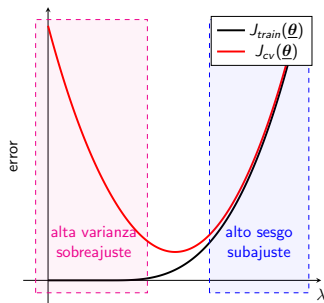
donde $\lambda \|\underline{\theta}\|_2^2$ se resta puesto que aquí **maximizamos** $J(\underline{\theta})$.

Ejemplos

- 1 Regresión polinomial: `nregress_norm.m`
- 2 Regresión logística: `rlr.m`

Selección de λ

- En general optimizamos con términos de regularización $\lambda \|\underline{\theta}\|_2^2$ (excluyendo a θ_0).
- La estimación del error para análisis se hace **sin** ese término.



Tipos de aprendizaje

- Hemos trabajado hasta ahora con aprendizaje **fuera de línea** (*off-line*)
- En particular trabajamos con aprendizaje **por lotes** y aprendizaje **secuencial**
- En estos métodos entrenamos primero la hipótesis h con un conjunto de datos y probamos en un conjunto aparte de prueba
- Existe además el aprendizaje **en línea**
- En estos métodos el algoritmo debe predecir continuamente mientras aprende.

Definición de aprendizaje en línea

- Dado un conjunto de datos

$$(\underline{\mathbf{x}}^{(1)}, y^{(1)}), (\underline{\mathbf{x}}^{(2)}, y^{(2)}), \dots, (\underline{\mathbf{x}}^{(m)}, y^{(m)})$$

que se producen en exactamente ese orden.

- El algoritmo primero recibe $\underline{\mathbf{x}}^{(1)}$, y predice $\hat{y}^{(1)}$
- Después de eso se le revela al algoritmo el valor verdadero $y^{(1)}$
- Con $y^{(1)}$ el algoritmo aprende
- Ahora se le presenta al algoritmo $\underline{\mathbf{x}}^{(2)}$, y predice $\hat{y}^{(2)}$.
- Con $y^{(2)}$ algoritmo aprende.
- Proceso se repite hasta m -ésima iteración
- Nos interesa el **error en-línea total**:

$$\varepsilon = \sum_{i=1}^m 1 \left\{ \hat{y}^{(i)} \neq y^{(i)} \right\}$$

Perceptron en-línea

- Se podría utilizar cualquier algoritmo de los anteriores usando todos los datos previos
- Aprendizaje secuencial se adapta casi directamente al caso:
- Por ejemplo, el perceptron

Inicialice $\underline{\theta} = \underline{0}$, $i = 1$

repeat

$$\left| \begin{array}{l} \underline{\theta} \leftarrow \underline{\theta} + \alpha(y^{(i)} - h(\underline{x}^{(i)}))\underline{x}^{(i)} \\ i \leftarrow i + 1 \end{array} \right.$$

until *convergencia*

- Se ha demostrado que esto siempre converge para datos separables con un margen geométrico γ .
- Ver notas [cs229-notes6.pdf](#) de Andrew Ng

Resumen

- 1 Enfoques de estimación de parámetros
 - Enfoque frecuentista
 - Enfoque Bayesiano
- 2 Funciones de error equivalentes
- 3 Aprendizaje en línea

Este documento ha sido elaborado con software libre incluyendo \LaTeX , Beamer, GNUPlot, GNU/Octave, XFig, Inkscape, GNU-Make y Subversion en GNU/Linux



Este trabajo se encuentra bajo una Licencia Creative Commons Atribución-NoComercial-LicenciarIgual 3.0 Unported. Para ver una copia de esta Licencia, visite <http://creativecommons.org/licenses/by-nc-sa/3.0/> o envíe una carta a Creative Commons, 444 Castro Street, Suite 900, Mountain View, California, 94041, USA.

© 2017–2019 Pablo Alvarado-Moya Área de Ingeniería en Computadores Instituto Tecnológico de Costa Rica