

Pauta Control 1 - Máquinas de Aprendizaje Computacional (INF-578)

Semestre II 2017 - Viernes 13.10.17

Respuestas:

1. Falso, ya que al aumentar la cardinalidad del espacio de hipótesis, hay menor garantía de haber encontrado la hipótesis óptima que modele la distribución conjunta subyacente usando el conjunto de entrenamiento. Esto ocurre, porque en un espacio más grande se puede encontrar una mayor cantidad de hipótesis con error 0 en el conjunto de entrenamiento.
- 2.

$$\begin{aligned}\frac{1}{2}(\mathbf{Y} - \tilde{\mathbf{X}}\beta)^T(\mathbf{Y} - \tilde{\mathbf{X}}\beta) &= \frac{1}{2}(\mathbf{Y}^T\mathbf{Y} - 2\mathbf{Y}^T\tilde{\mathbf{X}}\beta + \beta^T\tilde{\mathbf{X}}^T\tilde{\mathbf{X}}\beta) \\ \frac{1}{2}(\mathbf{Y} - \tilde{\mathbf{X}}\beta)^T(\mathbf{Y} - \tilde{\mathbf{X}}\beta) &= \frac{1}{2}(\mathbf{Y}^T\mathbf{Y} - 2\mathbf{Y}^T(\mathbf{X} + \mathbf{E})\beta + \beta^T(\mathbf{X} + \mathbf{E})^T(\mathbf{X} + \mathbf{E})\beta) \\ &= \frac{1}{2}(\mathbf{Y}^T\mathbf{Y} - 2\mathbf{Y}^T\mathbf{X}\beta - 2\mathbf{Y}^T\mathbf{E}\beta + \beta^T(\mathbf{X}^T\mathbf{X} + \mathbf{E}^T\mathbf{X} + \mathbf{X}^T\mathbf{E} + \mathbf{E}^T\mathbf{E})\beta) \\ &= \frac{1}{2}(\mathbf{Y}^T\mathbf{Y} - 2\mathbf{Y}^T\mathbf{X}\beta - 2\mathbf{Y}^T\mathbf{E}\beta + \beta^T\mathbf{X}^T\mathbf{X}\beta + 2\beta^T\mathbf{E}^T\mathbf{X}\beta + \beta^T\mathbf{E}^T\mathbf{E}\beta)\end{aligned}$$

Entonces, por la linealidad del valor esperado,

$$\begin{aligned}\mathbb{E}_\epsilon \frac{1}{2}(\mathbf{Y} - \beta\tilde{\mathbf{X}})^T(\mathbf{Y} - \beta\tilde{\mathbf{X}}) &= \frac{1}{2}(\mathbf{Y}^T\mathbf{Y} - 2\mathbf{Y}^T\beta\mathbf{X} - 2\mathbf{Y}^T\beta\mathbb{E}_\epsilon(\mathbf{E}) + \beta^T\mathbf{X}^T\mathbf{X}\beta \\ &\quad + \beta^T\mathbb{E}_\epsilon(\mathbf{E})^T\mathbf{X}\beta + \beta^T\mathbf{X}^T\mathbb{E}_\epsilon(\mathbf{E})\beta + \beta^T\mathbb{E}_\epsilon(\mathbf{E}^T\mathbf{E})\beta) \\ &= \frac{1}{2}(\mathbf{Y}^T\mathbf{Y} - 2\mathbf{Y}^T\beta\mathbf{X} - 2\mathbf{Y}^T\beta\mathbf{0} + \beta^T\mathbf{X}^T\mathbf{X}\beta \\ &\quad + \beta^T\mathbf{0}^T\mathbf{X}\beta + \beta^T\mathbf{X}^T\mathbf{0}\beta + \sigma^2\beta^T\mathbf{I}\beta) \\ &= \frac{1}{2}(\mathbf{Y} - \beta\mathbf{X})^T(\mathbf{Y} - \beta\mathbf{X}) + \frac{\sigma^2}{2}\|\beta\|^2\end{aligned}$$

que es equivalente a Ridge con $\lambda = \frac{\sigma^2}{2}$.

3. Con predictores ortogonales y asumiendo que la matrix \mathbf{X} está centrada tenemos que

$$\begin{aligned}\mathbf{X}^T\mathbf{X} &= \mathbf{\Sigma} = \text{diag}(\sigma_1^2, \sigma_2^2, \dots, \sigma_d^2) \\ \mathbf{\Sigma}^{-1} &= \text{diag}(\sigma_1^{-2}, \sigma_2^{-2}, \dots, \sigma_d^{-2}),\end{aligned}$$

La solución de mínimos cuadrados es

$$\hat{\beta} = (\mathbf{X}^T\mathbf{X})^{-1}\mathbf{X}^T\mathbf{Y} = \mathbf{\Sigma}^{-1}\mathbf{X}^T\mathbf{Y}$$

Si denotamos por $\mathbf{X}_{(j)}$ al vector correspondiente a la j -ésima columna de \mathbf{X} tenemos que el j -ésimo coeficiente de $\hat{\beta}$ es

$$\hat{\beta}_j = \frac{1}{\sigma_j^2}\mathbf{X}_{(j)}^T\mathbf{Y},$$

y su z -score viene dado por

$$z_j = \frac{\hat{\beta}_j}{\hat{\sigma}/\sigma_j} = \frac{\hat{\beta}_j \sigma_j}{\hat{\sigma}},$$

Si dos coeficientes $\hat{\beta}_j$ y $\hat{\beta}_i$ tienen el mismo peso

$$\frac{1}{\sigma_j^2} \mathbf{X}_{(j)}^T \mathbf{Y} = \frac{1}{\sigma_i^2} \mathbf{X}_{(i)}^T \mathbf{Y}$$

De estos dos coeficientes, aquel con mayor z -score es aquél con mayor σ_j (o, equivalentemente, mayor σ_j^2). Ahora, la solución de Ridge es

$$\tilde{\beta} = (\mathbf{X}^T \mathbf{X} + \lambda \mathbf{I})^{-1} \mathbf{X}^T \mathbf{Y}$$

cuyo j -ésimo coeficiente es

$$\tilde{\beta}_j = \frac{1}{\sigma_j^2 + \lambda} \mathbf{X}_{(j)}^T \mathbf{Y}$$

Tenemos que demostrar que si $\sigma_j^2 > \sigma_i^2$,

$$\frac{1}{\sigma_j^2 + \lambda} \mathbf{X}_{(j)}^T \mathbf{Y} > \frac{1}{\sigma_i^2 + \lambda} \mathbf{X}_{(i)}^T \mathbf{Y}$$

o equivalentemente, que

$$1 = \frac{\frac{1}{\sigma_j^2} \mathbf{X}_{(j)}^T \mathbf{Y}}{\frac{1}{\sigma_i^2} \mathbf{X}_{(i)}^T \mathbf{Y}} > \frac{\sigma_i^2 (\sigma_j^2 + \lambda)}{\sigma_j^2 (\sigma_i^2 + \lambda)}$$

En efecto, para λ estrictamente positivo,

$$\begin{aligned} \sigma_j^2 > \sigma_i^2 &\Rightarrow \lambda \sigma_j^2 > \lambda \sigma_i^2 \Rightarrow \sigma_j^2 \sigma_i^2 + \lambda \sigma_j^2 > \sigma_j^2 \sigma_i^2 + \lambda \sigma_i^2 \\ &\Leftrightarrow \sigma_j^2 (\sigma_i^2 + \lambda) > \sigma_i^2 (\sigma_j^2 + \lambda) \\ &\Leftrightarrow \frac{\sigma_i^2 (\sigma_j^2 + \lambda)}{\sigma_j^2 (\sigma_i^2 + \lambda)} < 1 \end{aligned}$$