

T2.6 Teoría de la decisión Bayesiana

Índice

- 1 Conceptos básicos
 - 1.1 Conceptos básicos
 - 1.2 Conceptos (un poco menos) básicos
- 2 Problemas de clasificación
 - 2.1 Problemas de clasificación: pérdida 01
 - 2.2 Matrices de confusión multiclase
- 3 Problemas de regresión
- 4 Problemas de predicción probabilística

1 Conceptos básicos

Inferencia Bayesiana: cálculo de la **posterior** $p(H \mid \mathbf{x})$ mediante la regla de Bayes actualizar nuestras creencias sobre cantidades ocultas H a partir de datos \mathbf{x}

Teoría de la decisión Bayesiana: usa la inferencia para decidir cuál es la mejor de las posibles **acciones** a realizar

1.1 Conceptos básicos

Agente: debe escoger una acción de un conjunto de acciones posibles, \mathcal{A}

Estado de la naturaleza: $h \in \mathcal{H}$, condiciona los costes y beneficios que se derivan de tomar cada acción posible

Función de pérdida: indica el coste incurrido al tomar la acción $a \in \mathcal{A}$ cuando el estado de la naturaleza es $h \in \mathcal{H}$

$$\ell(h, a)$$

Riesgo (pérdida) esperado a posteriori: de a tras observar \mathbf{x}

$$R(a \mid \mathbf{x}) = \mathbb{E}_{p(h|\mathbf{x})}[\ell(h, a)] = \sum_{h \in \mathcal{H}} \ell(h, a) p(h \mid \mathbf{x})$$

Política óptima o estimador de Bayes: obtiene una acción de mínimo riesgo por cada observación posible

$$\pi^*(\mathbf{x}) = \operatorname{argmin}_{a \in \mathcal{A}} R(a \mid \mathbf{x})$$

1.2 Conceptos (un poco menos) básicos

Función de utilidad: deseabilidad de cada acción posible en cada posible estado, esto es, riesgo con signo cambiado

$$U(h, a) = -\ell(h, a)$$

Principio de utilidad esperada máxima: estimador de Bayes expresado en términos de utilidad

$$\pi^*(\mathbf{x}) = \operatorname{argmax}_{a \in \mathcal{A}} \mathbb{E}_h[U(h, a)]$$

Sensibilidad al riesgo: asumimos que el agente es **neutral**, esto es, insensible al riesgo, pero podría no ser así; por ejemplo, nos da igual obtener 50 EUR con seguridad, o con 50% de probabilidad de 0 y 100 EUR

2 Problemas de clasificación

2.1 Problemas de clasificación: pérdida 01

Estados de la naturaleza y acciones: etiquetas de clase, $\mathcal{H} = \mathcal{Y} = \{1, \dots, C\}$ y $\mathcal{A} = \mathcal{Y}$

Pérdida 01 para dos clases: $\mathcal{Y} = \{0, 1\}$

$$\ell_{01}(y^*, \hat{y}) = \left[\begin{array}{c|cc} & \hat{y} = 0 & \hat{y} = 1 \\ \hline y^* = 0 & 0 & 1 \\ y^* = 1 & 1 & 0 \end{array} \right] = \mathbb{I}(y^* \neq \hat{y})$$

Pérdida esperada a posteriori: la probabilidad de error a posteriori es uno menos la de acertar a posteriori

$$R(\hat{y} \mid \mathbf{x}) = \sum_y \ell_{01}(y, \hat{y}) p(y \mid \mathbf{x}) = \sum_{y \neq \hat{y}} p(y \mid \mathbf{x}) = 1 - p(\hat{y} \mid \mathbf{x})$$

Estimador de Bayes: estimador máximo a posteriori (MAP), esto es, la etiqueta más probable o **moda** de la probabilidad a posteriori

$$\pi(\mathbf{x}) = \operatorname{argmax}_{y \in \mathcal{Y}} p(y \mid \mathbf{x})$$

2.2 Matrices de confusión multiclase

Datos: conjunto de pares etiqueta real-predicha, $\mathcal{D} = \{(y_m, \hat{y}_m)\}_{m=1}^M$, obtenidos al clasificar M muestras (de test)

Matriz de confusión para C clases: $\mathbf{M} = [M_{y,\hat{y}}]$ con $M_{y,\hat{y}} = \sum_m \mathbb{I}(y_m = y) \mathbb{I}(\hat{y}_m = \hat{y})$

y	$\hat{1}$	$\hat{2}$	\dots	\hat{C}	Suma fila
1	$M_{1,\hat{1}}$	$M_{1,\hat{2}}$	\dots	$M_{1,\hat{C}}$	$M_{1,:}$
2	$M_{2,\hat{1}}$	$M_{2,\hat{2}}$	\dots	$M_{2,\hat{C}}$	$M_{2,:}$
\vdots	\vdots	\vdots	\vdots	\vdots	\vdots

y	$\hat{1}$	$\hat{2}$	\dots	\hat{C}	Suma fila
C	$M_{C,\hat{1}}$	$M_{C,\hat{2}}$	\dots	$M_{C,\hat{C}}$	$M_{C,:}$
Suma:	$M_{:, \hat{1}}$	$M_{:, \hat{2}}$	\dots	$M_{:, \hat{C}}$	M

Normalización por filas: estimación empírica de $p(\hat{y} | y)$

Normalización por columnas: estimación empírica de $p(y | \hat{y})$

Normalización por filas y columnas: estimación empírica de $p(y, \hat{y})$

Análisis de una clase específica: se reduce a matriz binaria considerando el resto de clases como una única clase (negativa)

3 Problemas de regresión

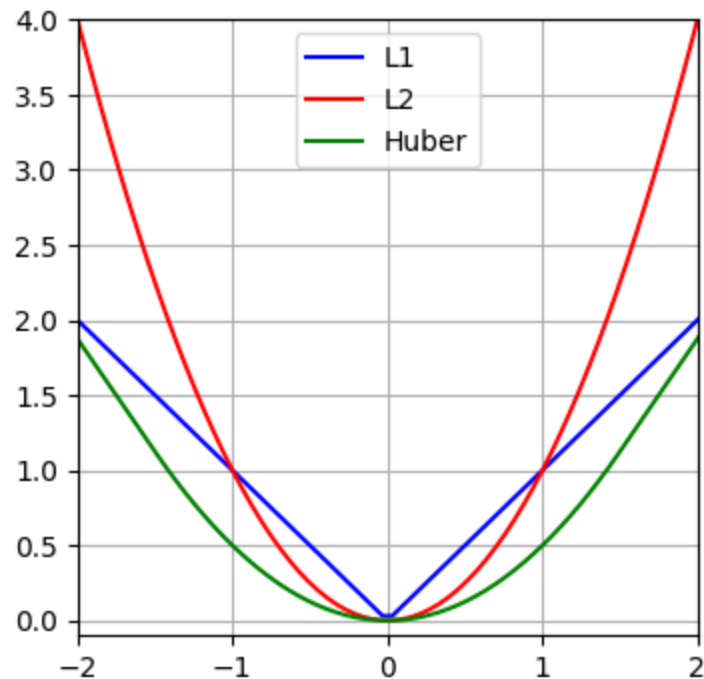
Estados de la naturaleza y acciones: reales, $\mathcal{H} = \mathcal{A} = \mathcal{Y} = \mathbb{R}$

Pérdidas L2 (ℓ_2 , cuadrática o error cuadrático), L1 (ℓ_1) y Huber: Huber combina L1 y L2 de acuerdo con un parámetro $\delta \geq 0$

$$\ell_2(h, a) = (h - a)^2 \quad \ell_1(h, a) = |h - a| \quad \ell_\delta(h, a) = \begin{cases} \frac{(h-a)^2}{2} & \text{si } |h - a| \leq \delta \\ \delta|h - a| - \frac{\delta^2}{2} & \text{si } |h - a| > \delta \end{cases}$$

Comparación gráfica: L2, L1 y Huber en función de la desviación de la verdad, $h - a$

```
In [2]: import numpy as np; import matplotlib.pyplot as plt
e = np.linspace(-3.0, 3.0, 100) # e = h - a (error)
L1 = abs(e); L2 = np.square(e); delta = 1.5; i = abs(e) <= delta
Huber = (abs(e)<=delta) * 0.5*L2 + (abs(e)>delta) * delta*(L1-delta/2);
plt.figure(figsize=(4, 4)); plt.xlim((-2, 2)); plt.ylim((-0.1, 4)); plt.grid();
plt.plot(e, L1, 'b'); plt.plot(e, L2, 'r'); plt.plot(e, Huber, 'g'); plt.legend(['L1', 'L2', 'Huber']);
```



Observaciones que se derivan de la comparación gráfica:

- L1 penaliza **linealmente** las desviaciones de la verdad
- L2 penaliza **cuadráticamente** las desviaciones de la verdad, por lo que es **más sensible a outliers** que L1
- Huber representa un compromiso entre L1 y L2

Pérdida L2 esperada a posteriori:

$$R(a | \mathbf{x}) = \mathbb{E}[(h - a)^2 | \mathbf{x}] = \mathbb{E}[h^2 | \mathbf{x}] - 2a\mathbb{E}[h | \mathbf{x}] + a^2$$

Regresor de Bayes L2 o minimum mean squared error (MMSE): media a posteriori

$$\frac{\partial}{\partial a} R(a | \mathbf{x}) = -2\mathbb{E}[h | \mathbf{x}] + 2a = 0 \quad \Rightarrow \quad \pi(\mathbf{x}) = \mathbb{E}[h | \mathbf{x}] = \int h p(h | \mathbf{x}) dh$$

Pérdida L1 esperada a posteriori:

$$R(a | \mathbf{x}) = \mathbb{E}[|h - a| | \mathbf{x}] = \int |h - a| p(h | \mathbf{x}) dh = \int_{-\infty}^a (a - h) p(h | \mathbf{x}) dh + \int_a^{\infty} (h - a) p(h | \mathbf{x}) dh$$

Regresor de Bayes L1: mediana a posteriori

$$a : P(h < a | \mathbf{x}) = P(h \geq a | \mathbf{x}) = 0.5$$

Pérdidas para \mathbb{R}^D : las pérdidas usuales para \mathbb{R} pueden extenderse fácilmente a \mathbb{R}^D y usarse para calcular los parámetros óptimos que debe devolver un estimador, la acción óptima que debe realizar un robot, etc.

4 Problemas de predicción probabilística

Estados de la naturaleza y acciones: distribuciones de probabilidad; $h = p(Y | \mathbf{x})$ y buscamos una $a = q(Y | \mathbf{x})$ que minimice $\mathbb{E}[\ell(p, q)]$ para un \mathbf{x} dado

Función de pérdida: divergencia de Kullback-Leibler (KL), en función de la entropía de p , $\mathbb{H}(p)$, y la entropía cruzada entre p y q , $\mathbb{H}(p, q)$

$$\mathbb{KL}(p||q) = -\mathbb{H}(p) + \mathbb{H}(p, q)$$

La minimización de KL equivale a minizar la entropía cruzada:

$$\begin{aligned} q^*(Y | \mathbf{x}) &= \operatorname{argmin}_q \mathbb{KL}(p(Y||\mathbf{x}), q(Y | \mathbf{x})) \\ &= \operatorname{argmin}_q -\mathbb{H}(p) + \mathbb{H}(p(Y | \mathbf{x}), q(Y | \mathbf{x})) \\ &= \operatorname{argmin}_q \mathbb{H}(p(Y | \mathbf{x}), q(Y | \mathbf{x})) \\ &= \operatorname{argmin}_q - \sum_{y \in \mathcal{Y}} p(y | \mathbf{x}) \log q(y | \mathbf{x}) \end{aligned}$$

En clasificación equivale a usar la log-pérdida: si h es one-hot, $h = p(Y | \mathbf{x}) = \delta(Y = c)$

$$\mathbb{H}(\delta(Y = c), q) = - \sum_{y \in \mathcal{Y}} \delta(y = c) \log q(y | \mathbf{x}) = -\log q(c | \mathbf{x})$$

Regla de puntuación propia: pérdida $\ell(p, q)$ cuya minimización en q converge a p

Puntuación de Brier: regla propia menos sensible a eventos raros que la entropía cruzada

$$\ell(p, q) = \frac{1}{C} \sum_{c=1}^C (q(y = c \mid \mathbf{x}) - p(y = c \mid \mathbf{x}))^2$$