

Билет 1

(Стр. 3)

Понятие точечной оценки параметра. Состоятельности и несмещенности оценок

Точечная оценка Пусть x_1, \dots, x_n - набор измерений какой-то величины. $\xi_1, \xi_2, \dots, \xi_n$ - набор iid. a - неизв. параметр.

ЦПТ: ξ_1, \dots, ξ_n - iid, $M_{\xi_i} = a$, $D_{\xi_i} = \sigma^2 \Rightarrow$ при $n \rightarrow \infty$ $\left(\frac{\xi_1 + \dots + \xi_n - na}{\sigma \sqrt{n}} < x \right) \rightarrow$

$\rightarrow N(0,1)$ или $\xi_1 + \xi_n - na \xrightarrow{\text{по распр.}} N(0,1)$

Пусть $F(x)$ - функция распр. ξ_i , предположим, она зав. от θ
Опр. Оценкой пар. θ наз. любая функция от (x_1, \dots, x_n) (при этом хотим, чтобы этой функ. были как можно ближе к измерениям)

Измеряя не знаем распр. ξ_i , хотим найти величину, от которой она зависит.

Опр. θ_n^* - несмещенная оценка θ , если $E\theta_n^* = \theta$

Опр. θ_n^* - состоятельная оценка, если $\forall \varepsilon > 0 \ P(|\theta_n^* - \theta| < \varepsilon) \xrightarrow{n \rightarrow \infty} 1$

Пример: ξ_1, \dots, ξ_n - iid, $E\xi_i = ?$

$\theta_n^* = \frac{\xi_1 + \dots + \xi_n}{n}$; $\theta_n^* = \xi_i$ - несмещенная оценка

Опр. θ_n^* - эффективна в мн. множестве оценок, если она обладает в нем минимальной дисперсией

Билет 2

Метод максимального правдоподобия (ММП). по оценке параметров

Основной способ получения (точечных) оценок по выборке - метод макс. правдоподобия

Опр. Ф-ция правдоподобия $L = L(x_1, \dots, x_n) = p(x_1, \dots, x_n, \theta)$ совм. вер. выборки

Если iid: $L = p(x_1, \theta) p(x_2, \theta) \dots p(x_n, \theta)$ x_1, \dots, x_n - фикс., θ могут быть различными для разных x_i

Опр. Оценка максимального правдоподобия $\theta_{м.п.} = \arg \max_{\theta} L$

Часто удобнее искать $\max \ln L(\theta)$, которая совм. с $\max L(\theta)$ в силу монотонности логарифма

$$\frac{\partial L}{\partial \theta} = 0 \Rightarrow \max \frac{\partial \ln L}{\partial \theta} = 0$$

Пример Бернулли $x_i \in \{0, 1\}$

$$L(x_1, \dots, x_n, p) = \prod_{i=1}^n p^{x_i} (1-p)^{1-x_i}$$

$$\ln L = \sum_{i=1}^n x_i \ln p + (1-x_i) \ln(1-p)$$

$$\frac{\partial \ln L}{\partial p} = \sum_{i=1}^n \frac{x_i}{p} - \sum_{i=1}^n (1-x_i) \frac{1}{1-p}$$

На самом деле $L = P_{\xi_1}(\theta) P_{\xi_2}(\theta) \dots P_{\xi_n}(\theta) (x_1, \dots, x_n)$, где $\xi_i(\theta)$ описывает зависимость случайной величины ξ от θ , x_i - измерения

Билет 3

сир. 2

Интервальный анализ. Доверительный интервал

имею точечной оценкой θ^* можно указать интервал $(\underline{\theta}, \bar{\theta})$. $\underline{\theta}(x_1, \dots, x_n), \bar{\theta}(x_1, \dots, x_n) : P(\underline{\theta}(x_1, \dots, x_n) < \theta < \bar{\theta}(x_1, \dots, x_n)) = 1 - 2\alpha$
 $(\underline{\theta}, \bar{\theta})$ - доверительный интервал для θ с дов. вероятностью $1 - 2\alpha$.

Пример: Пусть $\bar{x} = \frac{x_1 + \dots + x_n}{n}$, $x_1, \dots, x_n \sim N(a, \sigma)$ iid, $E\bar{x} = a$, $D\bar{x} = \frac{\sigma^2}{n}$

Центральной предельной теореме $\frac{\bar{x} - a}{\sigma/\sqrt{n}} \sim N(0, 1)$
 где $\alpha = \frac{1}{\sqrt{2\pi}} \int_{-\infty}^{-u_\alpha} e^{-x^2/2} dx$ (функция такова u_α)
 тогда дов. интервал с дов. вер-ностью $1 - 2\alpha$: $(\bar{x} - u_\alpha \frac{\sigma}{\sqrt{n}}, \bar{x} + u_\alpha \frac{\sigma}{\sqrt{n}})$

Если для 2 indep. значений пар. $\underline{\theta}_i, \bar{\theta}_i$ ($i=1, \dots, n$)
 Если для 2 indep. значений пар. $\underline{\theta}_i, \bar{\theta}_i$ ($i=1, \dots, n$)
 $P(\underline{\theta}(x_1, \dots, x_n) < \theta < \bar{\theta}(x_1, \dots, x_n)) \xrightarrow{n \rightarrow \infty} 1 - 2\alpha$, то такой интервал наз. асимптотическим.

Билет 4

Подход к проверке гипотезы о заданном распределении. Статистические примеры. Теорема о сходимости χ^2 к распр. χ^2 (без доказательства)

Статистическая гипотеза - некоторое предположение о законе и хар-ках indep. или совокупностей
 Проверка на осн. анализа наблюдаемых выборок. При этом статистические проверки не доказывают, что гипотеза верна.
 Пусть x_1, \dots, x_n iid, $F(x)$ - непрерывная д.ф.р. заданного ф-ции распределения x .
 Разобьем ось на промежутки: $(-\infty, z_1), (z_1, z_2), \dots, (z_{j-1}, z_j), (z_j, \infty)$ (1)
 Пусть F ф-ция распр. x , то $P_1 = P(x_i \in (-\infty, z_1)) = F(z_1)$
 $P_2 = P(x_i \in (z_1, z_2)) = F(z_2) - F(z_1)$
 \vdots
 $P_{j-1} = P(x_i \in (z_{j-1}, z_j)) = F(z_j) - F(z_{j-1})$
 $P_j = P(x_i \in (z_j, \infty)) = 1 - F(z_j)$

спуст. вел. x_i сводятся к независимым ед. образам: i -ое испытание - попадание x_i в некоторый интервал из (1)
 Пусть $n_i = n \cdot P_i$ (x_1, \dots, x_n) - число раз, попавших в (z_i, z_{i+1})
 $= n P_i$, $\chi^2_{n, j} = \sum_{i=1}^j \frac{(n_i - n P_i)^2}{n P_i}$ - величина эмпирического χ^2
 и $\chi^2_{n, j} > c$ где c - заданное значение по ф-ции распр. данная таблица значений и гипотеза принята, иначе нет.
 можно добиться $P(\chi^2_{n, j} > c) = \alpha$, где α малое
 Пусть спуст. вел. x_1, \dots, x_n распр. по закону $N(0, 1)$ и iid. Тогда пред. с.в. $R_n^2 = x_1^2 + \dots + x_n^2$ наз. распр. χ^2 с n степенями свободы $R_n^2 \sim \chi_n^2$

Теорема $\forall x \in \mathbb{R}, \forall n \rightarrow \infty P(\sum_{i=1}^n \xi_i < x) \rightarrow P(x, \infty)$ (сип. 3)
 с.в. ξ_i , имеют χ^2 -распр. с $k-1$ ст. своб.

Опр. Статистический критерий: Множество событий, есть $\xi_{n,i} \geq c$
 Множество не противоречит наоб., если произошло проис.
 противоположно событие.

Билет 5

Выбор из двух гипотез. Ошибки первого и второго рода

Пусть $\xi_1, \dots, \xi_n \sim N(a, \sigma)$, i.i.d, σ - известно, a - неизвест.
 2 гипотезы: $M_0: a = a_0, M_1: a = a_1, a_0 < a_1, M_0$ - осн. гипотеза, M_1 - альтернативная.

φ -функция распредел. величины ξ для разных гипотез различны, $P_0(P_1)$ - вероятности событий, вытекл. при гипотезе $M_0(M_1)$

Выберем оценку $\bar{\xi} = \frac{\xi_1 + \dots + \xi_n}{n}$ (она несмещ. и состоятельная)

Критерий: Выберем $a_0 < c < a_1$. Если $\bar{\xi} > c$, примем M_1 , иначе M_0

Опр. Ошибка первого рода - вероятность принять M_1 , когда верна M_0
 $\alpha = P_0(\bar{\xi} > c)$

Опр. Ошибка второго рода - вероятность принять M_0 , когда верна M_1
 $\beta = P_1(\bar{\xi} < c)$

Верная гипотеза
 $M_0 \quad M_1$

	M_0	$M_0 +$	$M_0 -$	α
принимать	M_1	$M_0 +$	$M_0 -$	

Пусть $\alpha = P_0(\bar{\xi} > c)$ задано. Тогда $\alpha = P_0(\bar{\xi} - a_0 > (c - a_0) \frac{\sqrt{n}}{\sigma}) = P_0(\frac{\bar{\xi} - a_0}{\sigma/\sqrt{n}} > (c - a_0) \frac{\sqrt{n}}{\sigma})$
 Следовательно, $c = a_0 + \frac{\sigma}{\sqrt{n}} \cdot \Phi^{-1}(\alpha)$, т.е. $\Phi^{-1}(\alpha) = \frac{(c - a_0) \sqrt{n}}{\sigma}$ распр. норм. и станд.
 Аналог. для ошибки второго рода $\beta = P_1(\bar{\xi} < c)$
 Получим $\alpha + \beta = \frac{a_1 - a_0}{\sigma} \sqrt{n} \Rightarrow$ можно выбрать α, β в завис. от
 степени непереносимости equivoc. с миним. ошибкой решения

Билет 6

Вероятностное пространство Ω - дан распределением

Опр. Пусть с.в. ξ, η независимы, то с.в. (ξ, η) имеет распр. заданное совместной плотностью $p(x, y) = p_{\xi, \eta}(x, y)$ тогда $p(x, y) = p(x|\eta) p(\eta) = p(\eta|x) p(x)$

φ -функция $p(x|\eta) = \frac{p(x, y)}{p(\eta)}$, по-усл. вероятности

Формула полной вероятности $p(x) = \int p(x|\eta) p(\eta) d\eta$ или $p(x) = \sum_{y \in Y} p(x|y) p(y)$

Теорема Байеса $p(x|\eta) = \frac{p(\eta|x) p(\eta)}{p(x)} = \frac{p(\eta|x) p(\eta)}{\int p(\eta|x) p(\eta) d\eta}$

Предположим: Пусть известно едм. множество \mathcal{P} распредел. $p(x, y)$ на $X \times Y$, где $X \subset \mathbb{R}^n$, $Y = \{0, 1\}$, зад. ор-ую номер $L(a(x), y)$

Опр. Среднее вел. номер для апоримна $a(x)$:

$$R(a) = \iint L(a(x), y) dP(x, y) = \iint L(a(x), y) P(x, y) dx dy$$

Хотим найти глуду суммарной массорпимации (по нас. $X \times Y$ ор. сивенн X, Y)

Задача Найти такой апоримн $a^*(x)$, где $a^*(x) = \arg \min_n R(a)$

Будем называть модель a^* оптимальной, R^* - значение мин. значения среднего вел.

Байес

Лин. Б.К. Лин и квадрат.

Дискр.

~~Байесовский классификатор~~
~~Теорема об оптимальности байесовского~~
~~классификатора~~

Опр. Распределение признака. Значимые семейства распределений, если м-ностр. распредел. можно лишь записать в следующем виде: $p(x|\theta) = h(x)g(\theta)e^{\eta(\theta)T(x)}$

Пример: нормальное, логнорм., бинамиальное, Г-распр.

Предн. $T(x)=y$, $p(x|y) = h(x)g_y(\theta_y)e^{\eta_y(\theta_y)x}$

Теорема о лин. байесовском классификаторе: Если для дискр. масс. множеств распредел. имеют едн. вид: $p(x|y) = h(x)g_y(\theta_y)e^{\eta_y(\theta_y)x}$, среди признаков есть постоянное по верн. 1) разд. пов-ост минимизация $(w, x) = \ln \frac{1}{1-z}$ $p(y|x) = \sigma(\langle w, x \rangle)$, где $\sigma(z) = \frac{1}{1+e^{-z}}$ - лог. ф-ция.

Д-во: $\frac{p(y=+1|x)}{p(y=-1|x)} = \frac{p(+1|y=+1)p(y=+1)}{p(+1|y=-1)p(y=-1)} = \frac{p(y=+1)h(x)g_+(\theta_+)e^{\eta_+(\theta_+)x}}{p(y=-1)h(x)g_-(\theta_-)e^{\eta_-(\theta_-)x}} =$
 $= \frac{p(y=+1)g_+(\theta_+)e^{\eta_+(\theta_+)x}}{p(y=-1)g_-(\theta_-)e^{\eta_-(\theta_-)x}} =$

Сделаем предн. о const. Введем θ и получим $e^{\langle w, x \rangle}$
 Из логнорм. вел. и то, что $p(y=+1|x) + p(y=-1|x) = 1 \Rightarrow$
 $\Rightarrow p(y|x) = \sigma(\langle w, x \rangle)$, где $\sigma(z) = \frac{1}{1+e^{-z}}$

Для дискр. масс. разд. поверхностей оптимального байесовского классификатора имеем вид $\frac{p(y=+1|x)}{p(y=-1|x)} - \frac{1}{1} = e^{\langle w, x \rangle} - \frac{1}{1} = 0$

Классификатор Quadratic Discriminant Analysis QDA

• предн., что найд. из каждого класса норм. распредел.
 • предн., что каждый класс имеет свое сов-е метр.

Классификатор Linear Discriminant Analysis LDA

• -1- единичную
 • -1- ~~сво~~ метр. метр.
QDA $a(x) = \frac{1}{2} x^T A x + (w, x) - b = 0$, где $A = \Sigma_0^{-1} - \Sigma_1^{-1}$
 μ_y - вектор метр. в классе y $w = \mu_1^T \Sigma_1^{-1} - \mu_0^T \Sigma_0^{-1}$

Z_y - вектор. матриц. размер. x $b = \ln \frac{p_1}{p_0} + \frac{1}{2} \ln \frac{\det \Sigma_0}{\det \Sigma_1} - \mu_1^T \Sigma_1^{-1} (\mu_1 + \mu_0) + \Sigma_0^{-1} \mu_0$

LDA: $a(x) = (w, x) - b = 0$, где $w = (\mu_1 - \mu_0)^T \Sigma^{-1}$

$b = \ln \frac{p_1}{p_0} - \frac{1}{2} (\mu_1 - \mu_0)^T \Sigma^{-1} (\mu_0 + \mu_1)$

Билет 9

Байесовский максимизатор.

Теорема об оптимальности байесовского максимизатора

Види ф-ции потерь для задачи бм. максимизации:

1) cross-entropy $L = p(y=1) \log(a(x)=1) + p(y=0) \log(a(x)=0) = -\log \text{loss}$

2) минимизатор $\begin{cases} 1, a(x)=y \\ 0, a(x) \neq y \end{cases}$

Рассм. 2) ф-цию потерь $L(a(x), y) = [a(x) \neq y]$

Средний риск $R(a) = \iint L(a(x), y) p(x, y) dx dy = \int \sum_y [a(x) \neq y] p(x|y) P(y) dy$

$= \int \sum_y (1 - [a(x)=y]) p(x|y) P(y) dx = \int \sum_y p(x|y) P(y) dx - \int \sum_y [a(x)=y] p(x|y) P(y) dx$

Получим $R(a) = \arg \max_a \int \sum_y [a(x)=y] p(x|y) P(y) dx = \arg \max_y p(x|y) P(y)$

Пусть ф-ция потерь $L(a(x), y) = \lambda_y [a(x) \neq y]$, где $\lambda_y \geq 0$

Теорема Минимум средних потерь при ф-ции потерь $L(a(x), y)$

достигается байесовским максимизатором $a(x) = \arg \max_y \lambda_y p(y|x)$

Следствие: Оптимальное решение максимизации при одинаковых потерях для всех классов.

Билет 6

Неивный байесовский максимизатор

Предп. Все признаки явл. незав. е.в. $p(x|y) = \prod p(x_i|y)$

Неивный байесовский максимизатор: $a(x) = \arg \max_{y \in Y} p(y) \prod p(x_i|y)$

Если предполож, что $p(x_i|y) = \frac{1}{\sqrt{2\pi}\sigma_y} e^{-\frac{(x_i - \mu_y)^2}{2\sigma_y^2}}$ тогда максимизатор

будет экв. неивному байесовскому логистическому максимизатору.

Параметры байесовского максимизатора $(P(y), \mu, \sigma)$ оцениваются

по обуча. мн-ву.

Замечание: работаем на малых выборках (нет корр. матриц)

Недостатки: размер. множества и угл.

однажды очень много обуча. выборки, т.е. реальн. размер.

Основные подходы:

• вычисл. мощность растет по входн. данным.

• идеальн. предп. о независимости (ф. распредел.) и по

данным несправедливо хороши

• нек. методы, методы вычисл. мощности,

Методы оценки вероятности
методов распр. Оценка Парзена-Розенблатта
(одног. случай)

Задача: По выборке $x^1 = (x_i, y_i)_{i=1}^n$, оценить m -ую вероятности $p^1(x)$ (без введ. нар. моментов)

Оцен. случай: Придана x , может применять мет. метод Гус-тени. Гус-тени. оценка $\hat{p}(x) = \frac{1}{n} \sum_{i=1}^n I(x=x_i)$

Непр. случай (одном.) По сур. m -ти $p(x) = \lim_{h \rightarrow 0} \frac{1}{2h} P(x-h, x+h)$

Выберем $x, h \in \mathbb{R}$ и положим h : $\hat{p}_h(x) = \frac{1}{2h} \sum_{i=1}^n I(x-h < x_i < x+h)$ $\hat{p}_h(x) = \frac{1}{2h} \sum_{i=1}^n I(\frac{x-x_i}{h} \in (-1, 1))$. Интервальная оценка $\hat{p}_h(x)$ от h жем. m -ого вер. гом. объектов выборки

Оцен. Парзена-Розенблатта

Интервальная оценка плотности по сур. m -ти h : $\hat{p}_h(x) = \frac{1}{2h} \sum_{i=1}^n k(\frac{x-x_i}{h})$
где $k(z) = \frac{1}{2} \exp(-|z|)$, $k(z) = \frac{1}{2} \exp(-|z|)$, $k(z) = \frac{1}{2} \exp(-|z|)$, $k(z) = \frac{1}{2} \exp(-|z|)$
коэф. $h > 0$, неопр. p -чис

Теорема ($x \in \mathbb{R}$)

Пусть $k(z)$ непрерывна, $k(z) \geq 0$, $\int_{-\infty}^{\infty} k(z) dz = 1$, $\lim_{h \rightarrow 0} h = 0$ и $\lim_{h \rightarrow 0} \frac{1}{h} = \infty$
Тогда 1) $\hat{p}_h(x) \rightarrow p(x)$ для $x \in X$
2) скорость сходимости порядка $O(h^2)$

Бизнес 11

Постановка задачи регрессии
Аналитическое решение задачи минимизации
перемен. Связь с методом наименьших квадратов

Постановка задачи:

Пусть X - объекты из \mathbb{R}^n , Y - значения из \mathbb{R}
 $x^i = (x_i, y_i)_{i=1}^n$ - обучающая выборка, $y_i = y(x_i)$, где $y: X \rightarrow \mathbb{R}$ - непрерывная
функ. зависимость. $a(x) = f(x, w)$ - модель, $w \in \mathbb{R}^p$ - вектор парам. модели

Метод наименьших квадратов (в обш. виде)

$Q(w, x^i) = \sum_{i=1}^n (x_i, (f(x_i, w) - y_i))^2 \rightarrow \min$, где x_i - все, что есть в выборке x ,
 $Q(w, x^i)$ - сум. квадратов

Предположим, что $a(x) = f(x, w) = w_0 + w_1 x_1 + \dots + w_n x_n$, где $w = (w_0, \dots, w_n)^T \in \mathbb{R}^{n+1}$
параметры модели

Другая форма: $a(x) = \frac{w^T \cdot x}{\|x\|}$, где $x = (1, x_1, \dots, x_n)^T \in \mathbb{R}^{n+1}$

Метод наим. квадратов: $L(w, x_{train}) = MSE(w, x_{train}) = \frac{1}{n} \sum_{i=1}^n (w^T x_i - y_i)^2$

Найдем $\hat{w} = \arg \min_w (L(w, x_{train}))$

Теорема: Решением задачи $\arg\min (\sum (w^T x_i - y_i)^2)$ (стр. 7)

сбн. $\hat{w} = (X^T X)^{-1} X^T y$, где $X = \{x_i\}$ - матрица признаков, $y = (y_1, \dots, y_n)$

Д-во: $\|Xw - y\|^2 \rightarrow \min_w$

$$\frac{\partial}{\partial w} \|Xw - y\|^2 = 0 \Rightarrow \frac{\partial}{\partial w} \|Xw - y\|^2 = \frac{\partial}{\partial w} (Xw - y)^T (Xw - y) = \frac{\partial}{\partial w} (Xw)^T Xw - y^T Xw + y^T y = \frac{\partial}{\partial w} w^T (X^T X) w - y^T Xw + y^T y = 2 \frac{\partial}{\partial w} (X^T X) w - y^T X = 0$$

Опр. Пусть $w = (w_1, \dots, w_n)$, $z = z(w_1, \dots, w_n)$. Тогда $\frac{\partial z}{\partial w} = (\frac{\partial z}{\partial w_1}, \dots, \frac{\partial z}{\partial w_n})^T$

Лемма: $\frac{\partial}{\partial x} x^T a = a$

Лемма: $\frac{\partial}{\partial x} x^T A x = (A + A^T)x$

$$\frac{\partial}{\partial w} w^T (X^T X) w = 2 \frac{\partial}{\partial w} (X^T y)^T w = 2 X^T y w - 2 X^T y = 0 \Rightarrow w = (X^T X)^{-1} X^T y$$

Билет 12

Регрессия и метод наименьших квадратов (МНК). Теорема об универсальности

Решение МНК и МРП

Модель данных с некоррелированными аддитивными шумами:
 $y(x_i) = f(x_i; w) + \varepsilon_i$, $\varepsilon_i \sim N(0, \sigma_i^2)$, $i = 1, \dots, n$

Метод наименьших квадратов:

$$L(\varepsilon_1, \dots, \varepsilon_n | w) = \prod_{i=1}^n \frac{1}{\sigma_i \sqrt{2\pi}} e^{-\frac{\varepsilon_i^2}{2\sigma_i^2}} \rightarrow \max$$

$$-\log L(\varepsilon_1, \dots, \varepsilon_n | w) = \text{const } w + \frac{1}{2} \sum_{i=1}^n \frac{1}{\sigma_i^2} (f(x_i; w) - y_i)^2 \rightarrow \min_w$$

- метод наим. кв.

Теорема: При предпос. выше (сро. логич. шум) решение МНК и МРП совп., причем веса обратно пропорц. дисперсии шума $w_i \propto \sigma_i^{-2}$

Билет 13

Теорема о минимальности величины
среднего потерь для регрессии с ф-цией
ошибки MSE.

Пусть известно совм. распр. $p(x, y)$ на $X \times Y$ (x - вход, y - ответ)

Пусть задана ф-ция потерь $L(a(x), y)$

Опр. Среднее значение потерь для аннотированного $a(x)$

$$R(a) = \iint L(a(x), y) dP(x, y) = \iint L(a(x), y) p(x, y) dx dy$$

Задача: Найти $a^*(x)$: $a^*(x) = \arg\min_a R(a)$ где a^* - оптимальная модель.

R - значение минимального среднего риска.

Теорема: Если $L(a(x), y) = (a(x) - y)^2$, то мин. среднее значение потерь

получается при $a^* = E(y|x) \leftarrow y_{\text{ср.}}$ минимиз.

$$D-во. R(a) = \iint L(a(x), y) p(x, y) dy dx = \iint (a(x) - y)^2 p(x, y) dy dx =$$

$$= \int \int (a(x) - y)^2 p(y|x) dy p(x) dx = \int E((y - a(x))^2 | x) p(x) dx$$

$$\text{Лемма: } E((y - a(x))^2 | x) = E((y - E(y|x))^2 | x) + E((a(x) - E(y|x))^2 | x)$$

$$R(a) = \int E((y - a(x))^2 | x) p(x) dx = \int E((y - E(y|x))^2 | x) p(x) dx +$$

$$+ \int E((a(x) - E(y|x))^2 | x) p(x) dx \geq \int E((y - E(y|x))^2 | x) p(x) dx$$

Билет 14

(Сир. 8)

Решение задачи предельной регрессии
 L_1, L_2 , Elastic-Net регуляризации. Вероятностный
смысл регуляризации (ср. 10 стр.)

L_2 -регуляризация: $L(w, X_{train}) = MSE(w, X_{train}) + \frac{\alpha}{2} \sum_{i=1}^n w_i^2$ - p -число потерь
 $= \frac{1}{2} \sum_i (w^T x_i - y_i)^2 + \frac{\alpha}{2} \sum_{i=1}^n w_i^2$

Задача: найти $\hat{w} = \arg\min_w (L(w, X_{train}))$
 Минимум: Решением задачи $\arg\min (\sum_i (w^T x_i - y_i)^2 + \alpha \sum_{i=1}^n w_i^2)$ экв.

Решение: $\hat{w} = (X^T X + \alpha I_{n+1})^{-1} X^T y$, где $x_i = x_i^T, y = (y_1, \dots, y_n)^T, I_{n+1}$ - ед. матрица.

Лемма: $\frac{\partial}{\partial w} x^T x = 2x$
 $\|Xw - y\|^2 + \alpha \|w\|^2 \rightarrow \min_w$

Необх. усл. минимума: $\frac{\partial}{\partial w} ((Xw - y)^T (Xw - y) + \alpha w^T w) = 2X^T Xw - 2X^T y + 2\alpha w = 0$

Свойства: • регуляризатор не даёт парам модели быть слишком большими

• Чем правее рег. экстр. большую свободу. Свойства

• Более жесткие и выдерживают, чем нестрогие при

• Больше переобучения, чем нестрогие при

Берем $\alpha = \frac{1}{T^2}$, T - среднее значение. Свойства

L_1 -регуляризация: $L(w, X_{train}) = MSE(w, X_{train}) + \alpha \sum_{i=1}^n |w_i|$ - p -число потерь

Задача: найти $\hat{w} = \arg\min_w (L(w, X_{train}))$

С-ва: Обеспечивает выбор признаков, тем самым. реш.

Смысл α : α - коэффициент регуляризации. Свойства

В формуле $L(w, X_{train}) = MSE(w, X_{train}) + \alpha \sum_{i=1}^n |w_i|$

Elastic Net: $L(w, X_{train}) = MSE(w, X_{train}) + \alpha \sum_{i=1}^n |w_i| + (1-\alpha) \frac{\alpha}{2} \sum_{i=1}^n w_i^2$ - p -число потерь

Задача: найти $\hat{w} = \arg\min_w (L(w, X_{train}))$

С-ва: Чем сильнее регуляризация, тем сильнее Ridge и Lasso рег.

Билет 15

Непараметрическая регрессия
Формула Кадарова-Вансена (ср. 12)

Пошаговая зад. регр. и регуляризации

• X - обучающая выборка (размер n), y - целевой (размер n), на X p -мерная

• $x_i \in \mathbb{R}^p, y_i \in \mathbb{R}$ - обучающие выборки, $y_i = y(x_i)$

• $a(x) = f(x, w)$ - непараметрическая модель зависимости от \mathbb{R}^p -вектор-параметров

• Метод наименьших квадратов: $Q(w, X^p) = \sum_i (f(x_i, w) - y_i)^2 \rightarrow \min_w$ где α - вес, степень важности обучающих данных;
 Недостатком: надо иметь хорошую начальную модель $f(x, w)$
 Будем приближ. модель $\text{const } f(x, w) = w$ в окр-ти $x \in X$

$Q(w, x^e) = \sum (x_i(x)(w - y_i)^2) \rightarrow \min_w$, где $x_i(x) = k\left(\frac{p(x, x_i)}{h}\right)$ - весовая функция
 x_i - обучающие, $k(z)$ - ядро, широкое, узкое, гауссовое, h - ширина или
 масштабирование

Формула Н-В: $a_n(x, X^e) = \frac{\sum_{i=1}^e y_i x_i(x)}{\sum_{i=1}^e x_i(x)} = \frac{\sum_{i=1}^e y_i k\left(\frac{p(x, x_i)}{h}\right)}{\sum_{i=1}^e k\left(\frac{p(x, x_i)}{h}\right)}$

Теорема Пусть непрерывны $y(x)$.

1) Выборка $X^e = (x_i, y_i)_{i=1}^e$ произв. из распр. $p(x, y)$

2) Ядро $k(z)$ вып. и $\int_{-\infty}^{\infty} k(z) dz < \infty$, $\lim_{z \rightarrow \infty} z k(z) = 0$

3) Зависит масса $E(y|x)$ и имеет верн. плотность:

$$E(y^2|x) = \int y^2 p(y|x) dy < \infty \text{ при } \forall x \in X$$

4) Коэф. h_e убывает, но не быстро: $\lim_{e \rightarrow \infty} h_e = 0$, $\lim_{e \rightarrow \infty} h_e' = \infty$

Тогда имеет место экв. по вер-ти: $a_{h_e}(x, X^e) \xrightarrow{p(x)} E(y|x)$
 в любой $x \in X$, где $E(y|x)$, $p(x)$, $D(x|y)$ вып. и $p(x) > 0$

Виды ядер: $\Pi(z) = [|z| \leq 1]$ - прямоугор. $E(z) = (1-z^2) [|z| \leq 1]$ - квадратичн.
 $T(z) = (1-|z|) [|z| \leq 1]$ - треугольн. $Q(z) = (1-z^2)^2 [|z| \leq 1]$ - кубическое
 $G(z) = e^{-z^2}$ - гауссовское

$G(z)$ - медленн. затухающее, широкое или узкое зависит
 от параметра h

$T(z)$ - вып. мин. шир., не вып. в центре ступенчатое T . выборки
 $\Pi(z)$ - вып. макс. шир., выбор ядра можно выбирать на основе