

## Билет 1

(Стр. 3)

### Понятие точечной оценки параметра. Состоятельности и несмещенности оценок

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

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

Пусть  $F(x)$  - функция распр.  $\xi_i$ , предположим, она зав. от  $\theta$   
Опр. Оценкой пар.  $\theta$  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.  $\xi_i$ , хотим найти величину, от которой она зависит.

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

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

Пример:  $\xi_1, \dots, \xi_n$  - iid,  $E\xi_i = ?$

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

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

## Билет 2

### Метод максимального правдоподобия (ММП). для оценки 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.,  $\theta$  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  $\xi$  от  $\theta$ ,  $x_i$  - измерения



### Билет 3

сир. 2

#### Интервальный анализ. Доверительный интервал

нельзя точечной оценкой  $\theta^*$  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  $\theta$  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  $u_\alpha$ )  
 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 $\chi^2$ к распр. $\chi^2$ (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  $(z_i, z_{i+1})$   
 $= n P_i$ ,  $Y_{n,i} = \sum_{j=1}^n \frac{(n_{i,j} - n P_i)^2}{n P_i}$  - величина эмпирической дисперсии  
 и  $Y_{n,i} > c$  где  $c$  - заданное значение по ф-ции распр. данная наблюд. значения и гипотеза против. гипотез.  
 можно добиться  $P(Y_{n,i} > c) = \alpha$ , где  $\alpha$  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.  $\chi^2$  с  $n$  степенями свободы  
 $R_n^2 \sim \chi_n^2$



Теорема  $\forall x \in \mathbb{R}, \forall n \rightarrow \infty P(\sum_{i=1}^n \xi_i < x) \rightarrow P(x_{n-1}^z, x), \text{ где}$  (сир. 3)

с.в.  $\xi_{n-1}^z$ , имеет  $\chi^2$ -распр. с  $n-1$  ст. своб.

Опр. Статистический критерий: Множество событий,

если  $\xi_{n,i} \geq c$

Множество не противоречит наоб., если произошло проис. в полном смысле.

### Билет 5

#### Выбор из двух гипотез. Ошибки первого и второго рода

Пусть  $\xi_1, \dots, \xi_n \sim N(a, \sigma)$ , i.i.d,  $\sigma$  - известно,  $a$  - неизвест.

2 гипотезы:  $M_0: a = a_0, M_1: a = a_1, a_0 < a_1, M_0$  - осн. гипотеза,  $M_1$  - конкурирующая.

$\varphi$ -функция распредел. величины  $\xi$  для разных гипотез различны,  $P_0(P_1)$  - вер-ность события, выпад. 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  
 $M_0 \quad M_1$

Вероятность

$M_0 \quad M_0 + \quad M_0 - \quad \alpha$

$M_1 \quad M_0 - \quad M_0 + =$   
 $\beta$

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,  $c = a_0 + \frac{\sigma}{\sqrt{n}} \cdot \gamma$ , т.е.  $\gamma = \frac{\sqrt{n}(\bar{\xi} - a_0)}{\sigma}$  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  $\alpha, \beta$  в завис. от  
степени непереносимости equivoc. с миним. ошибкой решения

### Билет 6

#### Вероятностное пространство $\Omega$ - дан распределением

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

$\varphi$ -функция  $p(x|\eta) := \frac{p(x, \eta)}{p(\eta)}$ , по-усл. вер-ности

Формула полной вер-ности  $p(x) = \int p(x|\eta) p(\eta) d\eta$  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., гамма-распредел., Г-распредел.

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

Теорема о лин. байесовском классификаторе: Если для дискр. масс. множества распредел. имеют едм. вид:  $p(x|y) = h(x)g_y(\theta_y)e^{\eta_y(\theta_y)x}$ , сред. признаков есть постоянная по веро-ти разд. пов-ост минимизация  $(w, x) = \ln \frac{1}{1-z}$   $p(y=1|x) = \sigma(w, x)$ , где  $\sigma(z) = \frac{1}{1+e^{-z}}$  - лог. ф-ция.

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

Сделаем предп. о const. Введем  $\theta$  и получим  $e^{(w, x)}$   
из полученного впр. и того, что  $p(y=+1|x) + p(y=-1|x) = 1 \Rightarrow$   
 $\Rightarrow p(y|x) = \sigma(w, x)$ , где  $\sigma(z) = \frac{1}{1+e^{-z}}$

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

### Классификатор Quadratic Discriminant Analysis QDA

• предп., что подл. из каждого класса норм. распредел.  
• предп., что каждый класс имеет свое covар. метр.

### Классификатор Linear Discriminant Analysis LDA

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



$Z_y$  - вектор. матрица, 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  $L = p(y=1) \log(a(x)=1) + p(y=0) \log(a(x)=0)$  ~~cross~~  $\logloss$

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

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. Априорная вероятность максимизации

## 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 предположим, 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 эквивалентен байесовскому логарифмическому максимизатору.

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

по обуча. мн-ву.

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

Недостатки: размер выборки не изв. одно есть только обуча. выборка, т.е. реальн. размер.

### Основные подходы

- вычисл. мощность растет по входн. данным.
- идея предп. о независимости (ф. разн.) и по данным максимизации пар-ров
- нек. методы, методы вычисл. мощности



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

Задача: По выборке  $x^1 = (x_i)_{i=1}^n$ , оценить  $m$ -ую производную  $p^m(x)$  (без введ. пар. функции)

Однор. случай: При этом  $x_i$  может принимать  $m$ -ую производную  $p(x)$ . Тогда  $p(x) = \frac{1}{n} \sum_{i=1}^n \delta(x - x_i)$

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

Выберем  $x, h \in \mathbb{R}$  и положим  $h$ :  $\hat{p}_h(x) = \frac{1}{2h} \sum_{i=1}^n I((x-x_i) < h) = \frac{1}{2h} \sum_{i=1}^n I(\frac{x-x_i}{h} < 1)$ . Тогда  $\hat{p}_h(x)$  есть  $h$  раз  $m$ -ой вер. гон. объектов выборки

Оценка Парзена-Розенблатта

Тогда  $\hat{p}_h(x) = \frac{1}{2h} \sum_{i=1}^n k(\frac{x-x_i}{h})$  где  $k(z) = \frac{1}{2} \int_{-1}^1 I(z-t) I(z+t) dt$  где  $k(z)$  - одно, удовл. требованиям:  $k(z) \geq 0$ ,  $\int_{-\infty}^{\infty} k(z) dz = 1$ ,  $k(z) \rightarrow 0$  при  $|z| \rightarrow \infty$

Теорема ( $x \in \mathbb{R}$ )

Пусть  $k(z)$  - одно, удовл.  $\int_{-\infty}^{\infty} k^2(z) dz < \infty$ ;  $k(z) \rightarrow 0$  при  $|z| \rightarrow \infty$ . Тогда 1)  $\hat{p}_h(x) \rightarrow p(x)$  для  $x \in \mathbb{R}$   
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  $\mathbb{R}$ . Пусть  $x^1 = (x_i)_{i=1}^n$  - обучающая выборка,  $y_i = y(x_i)$ , где  $y: x \rightarrow y$  - неизвестная зависимость.  $a(x) = f(x, w)$  - модель,  $w \in \mathbb{R}^p$  - вектор парам. модели

Метод наименьших квадратов (в обш. виде)

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

$Q(w, x^1)$  - сум. квадратов

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

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 гауссовыми 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:

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

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

- метод наим. кв.

Теорема: При предпос. выше (про регул. 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 непрерывных  $a(x)$

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

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

(Сур. 8)

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

$$L_1, L_2 \text{ - регуляризаторы: } L(w, X_{\text{train}}) = \text{MSE}(w, X_{\text{train}}) + \frac{\alpha}{2} \sum_{i=0}^n w_i^2 =$$

$L_2$  - регуляризатор:  $L(w, X_{\text{train}}) = \text{MSE}(w, X_{\text{train}}) + \frac{\alpha}{2} \sum_{i=0}^n w_i^2$  -  $p$ -число потерь  
 $= \frac{1}{2} \sum_i (w^T x_i - y_i)^2 + \frac{\alpha}{2} \sum_{i=0}^n w_i^2$  -  $p$ -число потерь

Задача: найти  $\hat{w} = \arg \min_w (L(w, X_{\text{train}}))$   
Метрика: Решением задачи  $\arg \min_w (\sum_i (w^T x_i - y_i)^2 + \alpha \sum_{i=0}^n w_i^2)$  экв.

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

Лемма  $\frac{\partial}{\partial x} x^T x = 2x$   
 $\|Xw - y\|^2 + \alpha \|w\|^2 \rightarrow \min_w$   
 Необход. 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: • регуляризаторы не дают парам модели быть слишком большими  
 • Чем правее регр. гиперпл. - тем больше обобщ. способность  
 • Более жесткие и вдумчивые, чем более правее модель

Вер. смысл  $\alpha$ :  $\alpha = \frac{1}{T^2}$ ,  $T$  - среднее значение сигнатур регр. на  $w$   
 $L_1$  - регуляризатор:  $L(w, X_{\text{train}}) = \text{MSE}(w, X_{\text{train}}) + \alpha \sum_{i=0}^n |w_i| = \sum_i (w^T x_i - y_i)^2 + \alpha \sum_{i=0}^n |w_i|$

Задача: найти  $\hat{w} = \arg \min_w (L(w, X_{\text{train}}))$   
С-ва: Обеспечиваем выбор признаков, тем сильнее регу.

Смысл  $\alpha$ :  $\alpha$  - обр. пропорц. среднеквадр. ошибки.  $p(w) = \frac{1}{T} e^{-\frac{\alpha w^2}{2T}}$   
 В смысле сигнатур это - регр. Лемма:  $p(w) = \frac{1}{T} e^{-\frac{\alpha w^2}{2T}}$

Elastic Net:  $L(w, X_{\text{train}}) = \text{MSE}(w, X_{\text{train}}) + 2\alpha \sum_{i=0}^n |w_i| + (1-2)\frac{\alpha}{2} \sum_{i=0}^n w_i^2$  -  $p$ -число потерь  
 $= \sum_i (w^T x_i - y_i)^2 + 2\alpha \sum_{i=0}^n |w_i| + (1-2)\frac{\alpha}{2} \sum_{i=0}^n w_i^2$

Задача: найти  $\hat{w} = \arg \min_w (L(w, X_{\text{train}}))$   
С-ва: Чем сильнее регуляризатор, тем сильнее сдвиг Ridge и Lasso регр.

## Билет 15

Непараметрическая регрессия  
Формула Кадаров-Вансена (ср. регр.)

Пошаговая зад. регр. и регуляризации  
 •  $X$  - обучающая (голова  $\mathbb{R}^n$ ),  $y$  - целевые (голова  $\mathbb{R}$ ), на  $X$   $p$ -мерная  
 •  $X^L = (x_i, y_i)_{i=1}^L$  - обучающая выборка,  $y_i = y(x_i)$

•  $a(x) = f(x, w)$  - непараметрическая модель зависимости из  $\mathbb{R}^p$  - вектор-параметров  
 • Метод наименьших квадратов:  $Q(w, X^L) = \sum_i (f(x_i, w) - y_i)^2 \rightarrow \min_w$  идеал

Все, что есть важности обучения;  
Недостатки: надо иметь хорошую парам. модель  $f(x, w)$   
 Будем прибр. 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  $x_i(x) = k\left(\frac{p(x, x_i)}{h}\right)$  - вес образца  $x_i$  относительно  $x$ ,  $k(z)$  - ядро, широкое, узкое, гауссовское,  $h$  - ширина или стандартное отклонение

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. уз. разл.  $p(x, y)$

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

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

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

4) Плот.  $h_e$  убывает, но не может:  $\lim_{e \rightarrow \infty} h_e = 0$ ,  $\lim_{e \rightarrow \infty} h_e = \infty$

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

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

$G(z)$  - медленная аппроксимация, широкое или узкое отклонение  
 Вывод на основе широты.

$T(z)$  - узк. или шир., не вып. в случае отклонения  $T$ . выборки  
 $\Pi(z)$  - узк. или шир., выбор ядра можно вывести на основе.