

Чтобы r была распределена по нормальному закону (так удобно считать квантили), с r производится преобразование Фишера:

$$Z = \frac{1}{2} \ln \frac{1+r}{1-r}$$

$$M(Z) = \frac{1}{2} \ln \frac{1+p}{1-p} ; D(Z) = \frac{1}{n-3}$$

Тогда

$$U = \sqrt{n-3} \left(Z - \frac{1}{2} \ln \frac{1+p}{1-p} \right) \sim N(0, 1)$$

распределена по стандартному нормальному закону.

Она с вероятностью $1-\varepsilon$ будет заключена в пределах $\pm U_{1-\varepsilon/2}$. Если решить получен-

ное неравенство

$$P \left(-U_{1-\varepsilon/2} \leq \sqrt{n-3} \left(Z - \frac{1}{2} \ln \frac{1+p}{1-p} \right) \leq U_{1-\varepsilon/2} \right) = 1-\varepsilon,$$

то получаются строгие и обратные формулы для вычисления r и проверки U -критерия.

Этот доверительный интервал соответствует к отклонению от нормальности и годится только для нормально распределенных двумерных случайных величин.

Билет № 15. Проверка статистических гипотез. Вероятность ошибки 1-го и 2-го рода. Уровень значимости критерия и мощность критерия.

Пусть есть некоторая случайная выборка, распределенная по неизвестному закону:

$$x_1; x_2; \dots; x_n = \xi; \xi \in F(x).$$

И пусть есть некая характеристика этого закона θ . (Мат. ожид., дисперсия или что-нибудь еще).

Мы предполагаем, что она равна некоторому числу θ_0 : вводим нулевую гипотезу:

$$H_0: \theta = \theta_0.$$

И альтернативную ей гипотезу H_1 :

$$H_1: \theta \neq \theta_0, \text{ простая альтернатива}$$

$$H_1: \theta \neq \theta_0, \text{ двусторонняя альтернатива}$$

$H_0: \theta \neq \theta_0$
 $\theta > \theta_0$
 $\theta < \theta_0$ } симметричные альтернативы

$H_1: \theta > \theta_0$
 $\theta < \theta_0$ } односторонние альтернативы

Выбрав нулевую гипотезу, мы берем новую выборку и вводим для нее некоторую функцию - статистику критерия $\varphi = \varphi(x_1; x_2; \dots; x_n)$.

Если H_0 верна, то эта функция будет иметь одно распределение, если не верна - то другое. В области значений этой ф-ии выбирается некоторая область W (критическая область), такая, что вероятность попадания значения ф-ии φ в эту область не превосходит заданного малого значения α - уровня значимости критерия. Т.е. если верна H_0 , то мы не будем выходить за критическую область, а умеем мы все время попадать за нее, значит, H_0 не верна, а верна H_1 .

Ошибки.

• Ошибка 1-го рода: верна H_0 , но мы ее отвергаем (наша гипотеза верна, но мы говорим, что не верна).

• Ошибка 2-го рода: верна H_1 , но мы ее отвергаем (наша гипотеза не верна, но мы говорим, что она верна).

Вероятность ошибки 1-го рода α и есть уровень значимости критерия (вероятность попасть в критическую область при верной H_0).

Вероятность ошибки 2-го рода β .

$1-\beta$ - вероятность правильно принять H_1 , когда верна H_1 . $1-\beta$ - мощность критерия.

Чем меньше вероятность ошибки 2-го рода, тем больше $1-\beta$, тем более мощный критерий.

Пример: проверка гипотез о математических ожиданиях.

Пусть $Z \sim N(\mu; \sigma^2)$. $H_0: \mu = \mu_0$; и есть выборка x_1, x_2, \dots, x_n
 $H_1: \mu = \mu_1$; значений случайной величины Z .

Но если $Z \sim N(\mu; \sigma^2)$, то $\frac{Z - \mu}{\sigma} \sim N(0; 1)$ - стандартное распределение.

А поскольку выборочное среднее имеет математическое ожидание исходной случайной величины, и дисперсию, равную дисперсии исходной случайной величины, по-

деленной на \sqrt{n} , то мы заменим ξ на \bar{x} , а σ/\sqrt{n} и получаем стандартизованное выборочное среднее:

$$u = \frac{\bar{x} - \mu_0}{\sigma/\sqrt{n}} \quad \begin{aligned} D(u) &= 1 \\ M(u) &= 0 \end{aligned}$$

А т.к. выборочное среднее - оценка для математического ожидания, то мы его выбираем в качестве статистики критерия.

H_0 : $\mu = \mu_0$; распределение U стандартно.

H_1 : $\mu = \mu_1$; распределение отличается от стандартного.

В качестве статистики критерия выбирается величина, характеризующая степень отклонения от нулевой гипотезы.

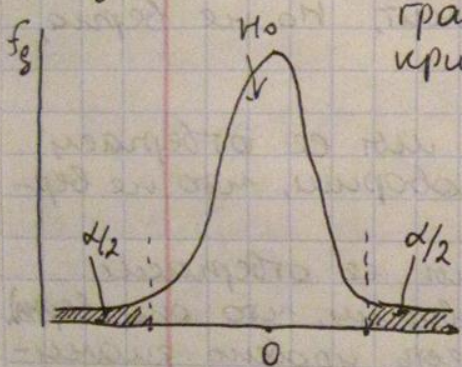


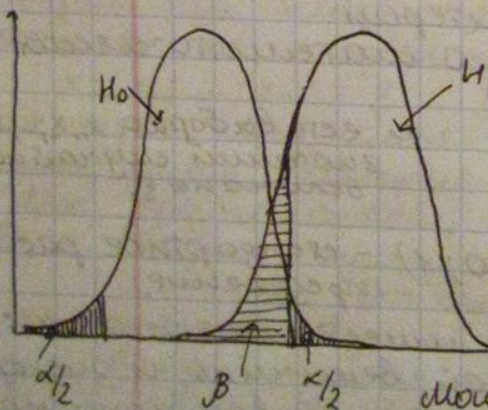
График функции плотности статистики критерия: при верной H_0 $\mu = \mu_0$.

Если мы получаем значение выборочного среднего и такие, что они попадают в критическую область, то есть повод задуматься, что H_0 не верна, а верна H_1 . Но в то же время H_0 может быть и верна с вероятностью α , а мы ее отвергли,

потому что u попало в критическую область. Поэтому α - вероятность ошибки 1-го рода.

Если верна H_1 , то распределение статистики U отличается от нормального на некоторую величину Δ :

$$\Delta = \frac{\mu_1 - \mu_0}{\sigma/\sqrt{n}}; \quad U \sim N(\Delta; 1)$$



Чем больше Δ , тем меньше вероятность ошибки 2-го рода β .

Чем меньше дисперсия, тем меньше вероятность β .

Чем больше α , тем меньше β .

Поскольку ошибки 1-го рода опаснее, то мы в первую очередь стараемся за α и стараемся его уменьшить. А β - как получится.

Мощность критерия - способность обнаружить истинное отклонение от нулевой гипотезы.

Ошибка 1 рода - ошибка ложного обнаружения несуществования эффекта
Ошибка 2 рода - ошибка ложного необнаружения существования эффекта.

Бишет в 16. Одновыборочный t-критерий.

Проверка гипотезы о равенстве заданному числу математического ожидания нормально распределенной случайной величины с неизвестной дисперсией.

В предыдущем случае нам было известна дисперсия, и мы брали в качестве статистики критерия стандартизованное выборочное среднее $u = \frac{\bar{x} - \mu_0}{\sigma/\sqrt{n}}$

Оно у нас было распределено по стандартному закону. Если дисперсия неизвестна, то вместо σ придется брать ее оценку S , но $s^2 = \frac{\sum_{i=1}^n (x_i - \bar{x})^2}{n-1}$, поэтому $u = \frac{\bar{x} - \mu_0}{s/\sqrt{n}}$ будет

иметь t-распределение с $(n-1)$ степенями свободы:

$$t = \frac{\bar{x} - \mu_0}{s/\sqrt{n}} \quad (\text{если верна } H_0).$$

Тогда критическая область для проверки $H_0: \mu = \mu_0$ при альтернативе $H_1: \mu \neq \mu_0$ будет состоять из двух бесконечных полуинтервалов

$$(-\infty; t_{n-1; \frac{\alpha}{2}}] \cup [t_{n-1; 1-\frac{\alpha}{2}}; \infty),$$

где $t_{n-1; \frac{\alpha}{2}}$ и $t_{n-1; 1-\frac{\alpha}{2}}$ - квантили t-распределения с $(n-1)$ степенями свободы порядков $\frac{\alpha}{2}$ и $1-\frac{\alpha}{2}$ (а в силу симметричности t-распределения $t_{\frac{\alpha}{2}} = -t_{1-\frac{\alpha}{2}}$).

Тогда если t попадет в интервал $[t_{n-1; \frac{\alpha}{2}}; t_{n-1; 1-\frac{\alpha}{2}}]$, то

H_0 согласуется с экспериментальными данными.

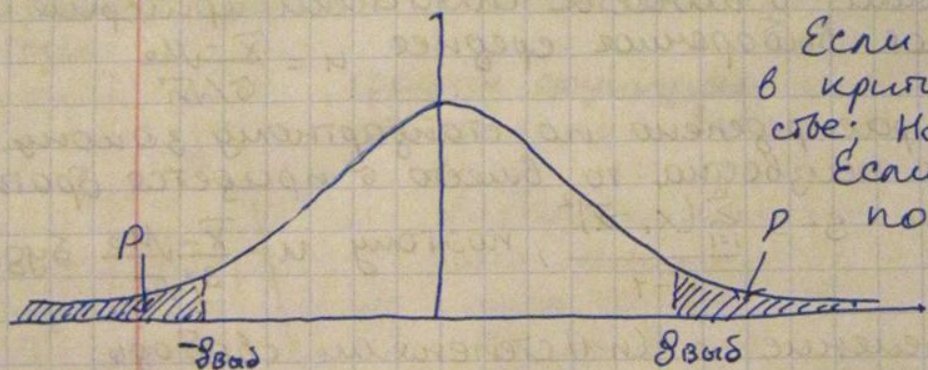
Для случаев, когда статистика выборки попала близко к границе и очень легко может через нее перейти, если повaryировать данными, а значит, резко возрастает вероятность ошибки 2 рода, существует p-значение.

P -значение - вероятность того, что статистика критерия g по модулю превзойдет $g_{\text{крит.}}$, вычисленное по выборке:

$$P = P \{ |g| > g_{\text{крит.}} \}$$

Т.е. мы считаем g для некоторой выборки, и вероятность того, что g для любой другой выборки будет больше, чем для этой, будет P -значением.

Чем меньше P , тем дальше "заехали" g для нашей выборки, тем меньше шансов его превзойти и тем больше шансов у нашей выборки вообще выйти за границу в критическое множество:



Если $P < \alpha$, то мы в критическом множестве. Но не верно.
 Если $P > \alpha$, то мы не попали в критическое множество. Но верно.
 Т.е. P -значение показывает, как близко мы от границы.

Применение t -критерия требует нормальности исходной случайной величины, но его можно использовать и при умеренных отклонениях от нормальности.

Билет № 17. Двухвыборочный t -критерий для независимых и связанных выборок.

Проверка гипотезы о равенстве математических ожиданий двух нормально распределенных случайных величин (двухвыборочный t -критерий).

1) Независимые выборки

Пусть есть 2 случайные выборки x_1, x_2, \dots, x_n и y_1, y_2, \dots, y_n значений двух независимых нормально распределенных случайных величин $\xi \sim N(M(\xi), D(\xi))$ и $\eta \sim N(M(\eta), D(\eta))$.

Гипотеза $H_0: M(\xi) = M(\eta)$.

$H_1: M(\xi) \neq M(\eta)$.

1) Известно, что $D(x) = D(y) = \sigma^2$, значение σ^2 неизвестно. Тогда для σ^2 можно получить объединенную несмещенную оценку:

$$s^2 = \frac{\sum_{i=1}^n (x_i - \bar{x})^2 + \sum_{j=1}^m (y_j - \bar{y})^2}{n+m-2}$$

Тогда s^2/n и s^2/m будут несмещенными оценками для дисперсии выборочных средних \bar{x} и \bar{y} , а сумма $\frac{s^2}{n} + \frac{s^2}{m}$ будет несмещенной оценкой для дисперсии разности средних $\bar{x} - \bar{y}$. Тогда статистика

$$t = \frac{\bar{x} - \bar{y}}{\sqrt{\frac{s^2}{n} + \frac{s^2}{m}}}$$

имеет t -распределение с $n+m-2$ степенями свободы.

Критическая область будет состоять из двух бесконечных полуинтервалов

$$(-\infty; t_{n+m-2; \frac{\alpha}{2}}] \cup [t_{n+m-2; 1-\frac{\alpha}{2}}; \infty), \text{ где}$$

$t_{n+m-2; \alpha/2}$ и $t_{n+m-2; 1-\alpha/2}$ - квантили t -распределения порядков $\alpha/2$ и $1-\alpha/2$ со степенями свободы $n+m-2$.

2) Если $D(x) \neq D(y)$, то для каждой из дисперсий вычисляется своя оценка:

$$s_x^2 = \frac{\sum_{i=1}^n (x_i - \bar{x})^2}{n-1} \quad \text{и} \quad s_y^2 = \frac{\sum_{j=1}^m (y_j - \bar{y})^2}{m-1}$$

Тогда статистика критерия

$$t = \frac{\bar{x} - \bar{y}}{\sqrt{\frac{s_x^2}{n} + \frac{s_y^2}{m}}}$$

имеет t -распределение с числом степеней свободы, равным целой части от $1/k$, где

$$k = \frac{\left(\frac{s_x^2/n}{s_x^2/n + s_y^2/m} \right)^2}{n-1} + \frac{\left(\frac{s_y^2/m}{s_x^2/n + s_y^2/m} \right)^2}{m-1}$$

2) Связанные выборки.

Пусть $x_1; x_2; \dots; x_n$ и $y_1; y_2; \dots; y_n$ - связанные случайные выборки из нормальных распределений $\xi \in N(M(\xi); D(\xi))$ и $\eta \in N(M(\eta); D(\eta))$.

$$H_0: M(x) = M(y);$$

$$H_1: M(x) \neq M(y).$$

Статистика критерия: $t = \frac{(\bar{x} - \bar{y}) \sqrt{n}}{\sqrt{S_x^2 + S_y^2 - 2S_{xy}}}$, где

$$S_{xy} = \frac{1}{n-1} \sum_{i=1}^n (x_i - \bar{x})(y_i - \bar{y}).$$

t имеет t -распределение с $(n-1)$ степенями свободы; критическая область состоит из двух интервалов $(-\infty; -t_{n-1; 1-\alpha/2})$ и $(t_{n-1; 1-\alpha/2}; \infty)$, где

$t_{n-1; 1-\alpha/2}$ - квантиль t -распределения порядка $1-\alpha/2$ с числом степеней свободы $n-1$.

Билет № 18. Двухвыборочный F-критерий.

Проверка гипотез о равенстве дисперсий двух независимых нормально распределенных случайных величин.

Пусть есть 2 случайные выборки $x_1; x_2; \dots; x_n$ и $y_1; y_2; \dots; y_m$ значений двух независимых нормально распределенных случайных величин

$$x \sim N(M(x); D(x)); \quad y \sim N(M(y); D(y));$$

$$H_0 = D(x) = D(y)$$

$$H_1 = D(x) \neq D(y).$$

Статистика критерия - отношение наименьших оценок дисперсий этих случайных величин:

$$F = \frac{S_x^2}{S_y^2} \text{ имеет } F\text{-распределение с } (n-1) \text{ и } (m-1) \text{ степенями свободы}$$

(n и m - объемы двух выборок).

Критическая область состоит из двух интервалов:

$$[0; F_{n-1; m-1; \alpha/2}] \text{ и } [F_{n-1; m-1; 1-\alpha/2}; \infty), \text{ где}$$

$F_{n-1; m-1; \alpha/2}$ и $F_{n-1; m-1; 1-\alpha/2}$ - квантили порядка $\alpha/2$ и $1-\alpha/2$ F-распределения с $n-1$ и $m-1$ степенями свободы.

В отличие от t -критерия, F-критерий чувствителен к отклонениям исходных случайных величин от нормальности.

Библия № 19. Критерий согласия χ^2 и Колмогорова-Смирнова.

Критерий согласия - критерий для проверки согласия между распределением выборочных значений и заданным теоретическим распределением.

Пусть есть выборка x_1, x_2, \dots, x_n значений случайной величины ξ с неизвестной φ -ей распределением $F(x)$

$$H_0: F(x) = F_0(x)$$

$$H_1: F(x) \neq F_0(x)$$

$F_0(x)$ - некоторое заданное распределение.

Если $F_0(x)$ задано полностью, то мы проверим простую нулевую гипотезу, если о $F_0(x)$ нам известно только то, что она принадлежит к некоторому виду ($F_0(x)$ - нормальное распределение, или F -распределение, или t -распределение), но больше мы о $F_0(x)$ ничего не знаем (ни ее дисперсия, ни мат. ожидание, и т.д.) - то это сложная нулевая гипотеза.

1) Критерий согласия χ^2 .

$F(x)$ может быть и непрерывным, и дискретным.

1) Простая нулевая гипотеза

Область изменения значения выборки разбивается на k интервалов так, чтобы число наблюдений n_i , попавших в i -ый интервал, было не менее 10.

n_i - эмпирическое число наблюдений i -ого интервала.

np_i - теоретическое число наблюдений i -ого интервала, т.е. какими оно должно быть, если бы H_0 была верна, и выборка принадлежала к $F_0(x)$.

Здесь n - объем выборки, p_i - вероятность попадания в i -ый интервал, рассчитанная исходя из известного распределения $F_0(x)$.

Тогда статистика критерия

$$\chi^2 = \sum_{i=1}^k \frac{(n_i - np_i)^2}{np_i}$$

Если H_0 верна, то распределение этой статистики приближается χ^2 -распределением с $k-1$ степенью свободы. Критическое множество состоит из одного полуинтервала $[\chi^2_{k-1; 1-\alpha}; +\infty)$, где $\chi^2_{k-1; 1-\alpha}$ - квантиль χ^2 -распределения с числом степеней свободы $k-1$ порядка $1-\alpha$.

2) Сложная нулевая гипотеза.

Мы разбиваем выборку на k интервалов, оцениваем выборочные характеристики и цуних вычисляем оценки вероятностей попадания в тот или иной интервал.

n_i - эмпирическое число наблюдений в i -м интервале

$n\hat{p}_i$ - теоретическое число наблюдений в i -м интервале, где n - объем выборки, \hat{p}_i - оценка вероятности попадания в i -ый интервал

Тогда
$$\chi^2 = \sum_{i=1}^k \frac{(n_i - n\hat{p}_i)^2}{n\hat{p}_i}$$

и если H_0 верна, то χ^2 при $n \rightarrow \infty$ распределена асимптотически (т.е. стремится к распределению как χ^2 с числом степеней свободы $k-r-1$, где r - число параметров $F_0(x)$).

Критий: $[\chi^2_{k-r-1; 1-\alpha}; +\infty)$, где

$\chi^2_{k-r-1; 1-\alpha}$ - квантиль χ^2 -распределения с числом степеней свободы $k-r-1$ порядка $1-\alpha$.

Критерии Колмогорова и Смирнова.

Применяются для проверки соответствия непрерывного распределения $F(x)$ заданному $F_0(x)$.

$H_0: F(x) = F_0(x)$

$H_1: F(x) \neq F_0(x)$

Статистика Колмогорова - мера близости эмпирической ф-ии распределения $\hat{F}(x)$ и теоретической $F_0(x)$:

$$D_n = \sup_x |\hat{F}(x) - F_0(x)|$$

Чем ближе друг к другу $\hat{F}(x)$ и $F_0(x)$, тем меньше верхняя граница их разницы.

Статистика Смирнова:

$H_0: F(x) = F_0(x)$

$H_1^+: F(x) > F_0(x)$

$$D_n^+ = \sup_x |\hat{F}(x) - F_0(x)|$$

Критическое множество для проверки H_0 против H_1^+ :

$$[D_n; 1-\alpha; +\infty)$$

Для проверки H_0 против H_1^+ : $[D_n^+; 1-\alpha; +\infty)$

где $D_{n, 1-\alpha}$; $D_{n, 1-\alpha}^+$ - критические значения статистик D_n и D_n^+ .

В случае сложной нулевой гипотезы

$$\hat{D}_n = \sup_x |F(x) - F_0(x; \hat{\theta}_1, \dots, \hat{\theta}_r)|,$$

где $\hat{\theta}_1, \dots, \hat{\theta}_r$ - оценки неизвестных параметров. Если для простой нулевой гипотезы распределение статистик D_n и D_n^+ при справедливости H_0 не зависит от типа $F_0(x)$, то в случае сложной нулевой гипотезы при верной H_0 распределение D_n и D_n^+ уже зависит от конкретного вида распределения $F_0(x)$.

Бишет №20. Критерий знаков и ранговых знаков.

Непараметрические критерии - не требуют знания вида исходного распределения $F(x)$ за исключением предположения о непрерывности $F(x)$.

Используются, если выборка не принадлежит к нормальному распределению.

Ранг - порядковый номер наблюдения при их упорядочении по возрастанию.

Одновыборочные непараметрические критерии - для проверки гипотезы о равенстве медианы заданному значению.

Пусть есть выборка x_1, x_2, \dots, x_n значений случайной величины Z с неизвестной непрерывной ф-ей распределения $F(x; M_e)$, где M_e - неизвестная медиана.

1) Критерий знаков

$$H_0: M_e = M_{e0}$$

$$H_1: M_e \neq M_{e0}$$

Статистика критерия:

n^+ - число положительных

разностей $x_i - M_{e0}$; $i = 1, \dots, n$.

при верной H_0 $P(x_i > M_{e0}) = P(x_i < M_{e0}) = 1/2$.

Статистика n^+ - дискретная случайная величина, распределенная по биномиальному закону,

$p = 1/2$.

Билет № 19. Проверка гипотезы о равенстве заданному числу коэффициента корреляции

Пусть $(x_1, y_1); (x_2, y_2); \dots; (x_n, y_n)$ - случайная выборка пар значений двумерной случайной величины (ξ, η) , имеющей двумерное нормальное распределение

$$H_0: \rho = \rho_0 \quad H_1: \rho \neq \rho_0.$$

Коэффициент корреляции $\rho = \rho_0$ - заданное число.

Для этого можно использовать статистику:

$$U = \frac{\sqrt{n-3}}{2} \left(\ln \frac{1+r}{1-r} - \ln \frac{1+\rho_0}{1-\rho_0} \right).$$

При верной H_0 и при большой n распределение ее приближается к стандартному нормальному

$$\text{Интервал: } (-\infty; U_{\alpha/2}] \cup [U_{1-\alpha/2}; \infty).$$

Если проверяется $H_0: \rho = 0$, то это эквивалентно проверке гипотезы о независимости ξ и η .

Тогда

$$U = \frac{\sqrt{n-3}}{2} \left(\ln \frac{1+r}{1-r} \right).$$

Билет № 23. Классификация методов многомерного статистического анализа.

Данные, описываемые любым числом переменных - многомерные данные.

Представляются в виде матрицы, строки - наблюдения, столбцы - переменные.

Независимые переменные - факторы.

Анализ данных



Анализ зависимостей

анализ факторов

анализ взаимосвязей между переменными

анализ структуры многомерных данных

переменные



количественные

качественные

y-зависимые переменные, x-независимые

Анализ факторов

количественные переменные

Факторный анализ
(н.р. метод главных компонент)

качественные переменные

Кластерный анализ
(агломеративно-иерархический метод - разбиение и систематизация всей совокупности наблюдений).

Анализ зависимостей

Y-количественные

Y-качественные

X-колич.

X-качеств.

X-колич.

X-качеств.

Регрессионный анализ - поиск функциональной зависимости Y от X.

Дисперсионный анализ
Установление связи между X и Y.

Дискриминантный анализ - получение правила, позволяющего на основе X предсказывать Y.

Сегментационный анализ
Последовательное разбиение совокупности наблюдений Y.

Билет №24. Регрессионный анализ.

Пусть есть матрица наблюдений и переменных, и есть переменные Y, X_1, X_2, \dots, X_m .

$\begin{pmatrix} y_1 & x_{11} & x_{12} & \dots & x_{1m} \\ y_2 & x_{21} & x_{22} & \dots & x_{2m} \\ \dots & \dots & \dots & \dots & \dots \\ y_n & x_{n1} & x_{n2} & \dots & x_{nm} \end{pmatrix}$ - матрица экспериментальных данных

В регрессионном анализе рассматривается связь между переменной Y (зависимой) и переменными X_1, X_2, \dots, X_n - независимыми, эта связь описывается моделью:

$$Y = F(X_1, X_2, \dots, X_m; \beta_0, \beta_1, \dots, \beta_k) + \varepsilon,$$

где $\beta_0, \beta_1, \dots, \beta_k$ - неизвестные параметры - коэффициенты регрессии;
 ε - случайное отклонение Y от F (ошибка измерения).
У каждого результата измерения X_i есть своя ошибка измерения, но ошибки X_i несоизмеримы с ошибками Y, поэтому ошибками X_i мы пренебрегаем.
Ортогональная регрессия - ошибки X и Y соизмеримы.
Виды регрессии.

1 независимая переменная - простая линейная регрессия

$$y_i = \beta_0 + \beta_1 x_{i1} + \varepsilon_i$$

Много независимых переменных - множественная линейная регрессия:

$$y_i = \beta_0 + \beta_1 x_{i1} + \beta_2 x_{i2} + \dots + \beta_m x_{im} + \varepsilon_i$$

ε - независимые случайные величины, распределенные по нормальному закону:

$$\varepsilon_i \sim N(0; \sigma^2)$$

Если модель по независимым переменным нелинейна, а по параметрам линейна, то она все равно остается линейной. Н-р, полиномиальная:

$$y_i = \beta_0 + \beta_1 x_{i1} + \dots + \beta_k x_{i1}^k + \varepsilon_i$$

Она сводится к линейной путем введения новых независимых переменных.

Экспоненциальная модель вида

$$y = a \varepsilon e^{bx}$$

нелинейна по параметрам, но может быть приведена к линейной логарифмированием.

Но

$$y = a_1 \varepsilon_1 e^{b_1 x} + a_2 \varepsilon_2 e^{b_2 x}$$

к линейной уже свести нельзя.

$$y_i = \beta_0 (1 - e^{-\beta_1 x_{i1}})^{\beta_2} \quad \text{- S-образная кривая.}$$

Нахождение параметров регрессии.

1. Оценивание параметров
2. Нахождение доверительных интервалов для параметров
3. Проверка гипотез на эти параметры.

Т.е. мы ищем не функцию от параметров, а подбираем какую-либо модель и проверяем ее применимость.

Простой линейной регрессионной анализ
Строим двумерную диаграмму рассеивания и прямую, такую, чтобы расстояния от точек, соответствующих наблюдениям, до прямой, были наименьшими.

Получаем уравнение прямой:

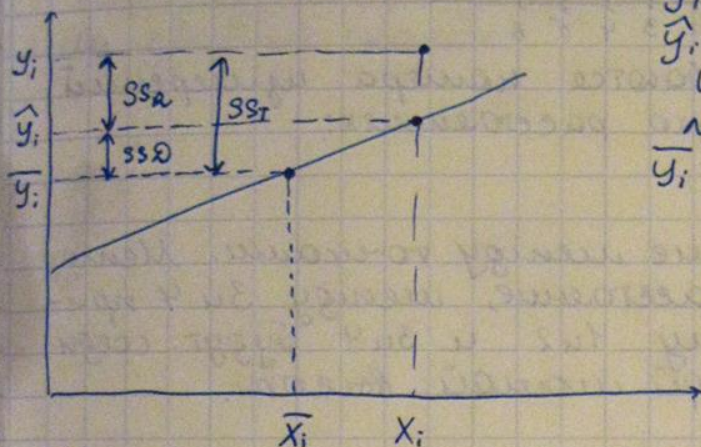
$$y_i = \beta_0 + \beta_1 x_i + \varepsilon_i$$

β_0 - точка пересечения линии с Oy

β_1 - tg угла наклона.

Нахождение расстояний между прямой и точками (метод наименьших квадратов)

Пусть у нас не ортогональная регрессия, где расстояния измеряются по перпендикуляру.



y_i - наблюдение

\hat{y}_i - соответствующее ему значение на прямой регрессионной линии, зависит от параметра β

\bar{y} - среднее наблюдение

Тогда

$$SS_R = \sum_{i=1}^n (y_i - \hat{y}_i)^2 \text{ должно быть минимальным.}$$

Другие параметры:

$$SS_D = \sum_{i=1}^n (\hat{y}_i - \bar{y})^2$$

$$SS_T = \sum_{i=1}^n (y_i - \bar{y})^2$$

$$SS_T = SS_D + SS_R$$

Если $SS_D = 0$, то прямая

идет по линии $y = \bar{y}$. т.е. y не зависит от x . Поэтому это не очень хорошая ситуация.

$R^2 = \frac{SS_D}{SS_T}$ - коэффициент детерминации (определяет качество предсказания).

$R^2 = 0$ плохое предсказание

$SS_D = 0$

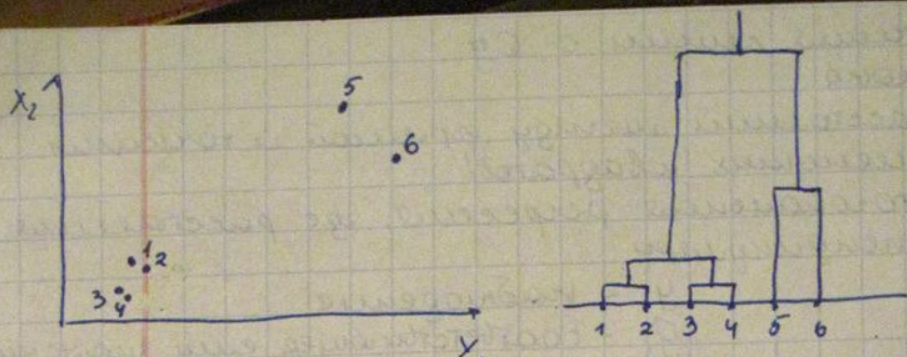
$R^2 = 1$ $SS_R = 0$ хорошее предсказание.

Смысл R^2 - какой процент общего разброса данных объяснен с помощью линейной регрессии.

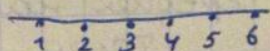
Бинет №25. Кластерный анализ.

Пусть есть n переменных и некоторое число наблюдений. Кластерный анализ оценивает, какие из этих наблюдений близки друг к другу, а какие - далеки. находит взаимосвязи и классифицирует данные.

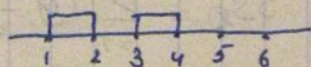
Если 2 переменных, то можно построить двумерную диаграмму рассеяния, а в многомерном случае строится дендрограмма.



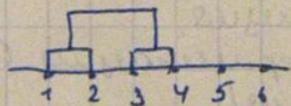
На дендрограмме указываются номера измерений на одинаковых друг от друга расстояниях:



Потом измеряется расстояние между точками. Между 1 и 2 самое маленькое расстояние, между 3 и 4 примерно такое же. Поэтому 1 и 2 и 3 и 4 будут соединяться столбиками самой малой высоты.



Затем измеряется расстояние между группами 1-2 и 3-4. Оно уже побольше, и столбик между ними будет выше.



5 и 6 вообще далеко друг от друга - между ними высокий столбик. Ну а самый высокий будет между 1234 и 56.

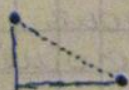
Оценка расстояний между точками:

• Евклидово расстояние

$$P_{1,2} = \sqrt{(x_{11} - x_{12})^2 + (x_{21} - x_{22})^2}$$

Но у разных переменных разные единицы измерения, и складывать различия в возрасте и в цвете глаз как-то неправильно.

• Манхеттенское расстояние



сумма катетов прямоугольного треугольника, гипотенуза которого соединяет точки

Оценка расстояния от точки до совокупности точек

- Метод средней связи - расстояние до точек 1 и 2 считается как среднее между расстоянием до точки 1 и точки 2.
 - Метод ближайшего соседа - расстояние до группы точек считается как расстояние до ближайшей из этих точек
- Для разбиения на группы достаточным считается пробной или двойной скачок.

об
и
ри-

и
ри-

и
ри-

и