# Statistics

## ML / Stats - Regression Analysis

- **Coefficients** $\beta_i$
  - Holding all other variables constant, 1 unit change of $x_i$ averagely leads to $\beta_i$ units of change in $\hat{y}$
- **P value**
  - Null Hypothesis: $\beta_i = 0 \rightarrow$ 假设这个variable 对 $Y$ 没有影响
  - If P-value $\leq 0.05$ (0.01) $\rightarrow$ Reject Null and favor that $\beta_i \not\equiv 0$
- **Feature Importance**
  - Coefficient $\beta_i$ 值的大小 代表该variable对 $Y$ 的影响程度大小
    - **Note: 讨论importance时, variables都需要standardize**
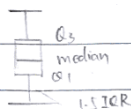  - Statistical Significance: 确认该variable 的 p-value $< 0.05$

## Stats - A/B Testing

- **P Value**: How extreme the observed value is under the null hypothesis
- Law of Large Numbers (LLM)
  - $\bar{Y}_n$ converges in probability to $\mu_Y \equiv E[Y]$ when sample size $n$ is large
  - https://builtin.com/data-science/law-of-large-numbers
- **Central Limit Theorem** (CLT)
  - https://www.scribbr.com/statistics/central-limit-theorem/
  - CLT allows us to study $\bar{Y}_n$ even if sample data are not normally distributed

## Stats - How to choose a test

- 如果是 **Numerical** & **sample mean** 有足够意义
  - 如果 population $\mu, \sigma$ 已知 $\rightarrow$ Z-test
  - 如果 population $\mu, \sigma$ 未知 $\rightarrow$ **T-test**
    - 如果treatment 对于个体差异非常大 $\rightarrow$ paired t-test
      - **https://math.stackexchange.com/questions/1732771/paired-t-test-vs-welchs-t-test**
    - 如果个体差异小
      - https://www.statology.org/paired-vs-unpaired-t-test/
      - 组间差距大 (variance 不同, heteroskedasticity) $\rightarrow$ **Welch's t-test**
      - 组间差距小 (variance 相同, homoscedasticity) $\rightarrow$ unpaired (Student's) t-test
    - 如果有超过2组 $\rightarrow$ ANOVA https://www.qualtrics.com/experience-management/research/anova/
- 如果是 **Categorical** $\rightarrow \chi^2$ test https://www.investopedia.com/terms/c/chi-square-statistic.asp
- 如果 sample mean 不是讨论对象
  - 用KS Test 来 compare distribution https://towardsdatascience.com/kolmogorov-smirnov-test-84c92fb4158d

Variable: Numeric (0,1,2) Continuous. Categorical: Nominal, Ordinal

Simple Random Sample. Stratefied $\Big\{$ under grad    diff.
    grad.

cluster Sampling: gym 1,2,3

$Q_3$
median
$Q_1$
1.5 IQR

Bias: Bad Sample Frame Bias. Convenience Sample Bias. Volunteer Bias

Disjoint: No common outcomes. $P(A \text{ or } B) = P(A) + P(B) (- P(A \cap B))$

$\Big\{$ Indep: $P(A|B) = P(A)$ / $P(A \cap B) = P(A) \cdot P(B)$    Conditional $P(A|B) = \frac{P(A \cap B)}{P(B)}$

Bayes: $P(A|B) = \frac{P(B|A) \cdot P(A)}{P(B)} = \frac{P(B|A) \cdot P(A)}{P(B|A) \cdot P(A) + P(B|A^c) \cdot P(A^c)}$

Random Variable.

$Var(X) = E[X^2] - E^2[X]$

Discrete Random Variable: 有 infinitely many outcomes

$SD(aX) = |a| \cdot SD(X)$

$Var(X) = \sum (x - \mu)^2 \cdot p(x)$.    $Var(X \pm c) = Var(X)$.    $Var(aX) = a^2 Var(X)$

① First success. $X \sim Geom(p \text{ success})$. $p(x) = (1-p)^{x-1} \cdot p$. $E(x) = \frac{1}{p}$. $SD(x) = \frac{\sqrt{1-p}}{p}$

② k success in n trials. $X \sim Binom(n, p)$ $p(x) = c(n,k) q^{n-k} \cdot p^k$.

$E(x) = np$. $SD(x) = \sqrt{npq} = \sqrt{np(1-p)}$

x: # of times an event occurs when average rate of occurance is $\lambda$.

③ Average. $p(个例)$. $X \sim poisson(\lambda)$. $p(x) = \frac{\lambda^x \cdot e^{-\lambda}}{x!}$, $x \in N^*$. $0! = 1$

$E(x) = \lambda$. $SD = \sqrt{\lambda}$.

④ $k^{th}$ success on $n^{th}$ trial. $X \sim NegBinom(k, p)$.

$p(x) = c(n-1, k-1) p^k \cdot (1-p)^{n-k}$. $E(x) = \frac{k}{p}$. $SD(x) = \frac{\sqrt{k(1-p)}}{p}$

Continuous Random. Variable: pdf. ① $f(x) \geq 0$ for all x. ② $\int_{-\infty}^{\infty} f(x) dx = 1$

$\mu = E[x] = \int_{-\infty}^{\infty} x \cdot f(x) dx$. $Var(x) = \int (x - \mu)^2 f(x) dx = E[x^2] - \mu^2$

$E[x^2] = \int_{-\infty}^{\infty} x^2 f(x) dx$.

Hypo:

$p(\text{type I error}) = \alpha$. (reject $H_0$ when $H_0$ is true)    Type II: favor $H_0$ when $H_0$ false.

Mean
One sample T: T-dist. Inference: Indep. (Random + < 10% population)
    $\Big\{$ Nearly normal pop. dist.

Uniform Dist. $f(x) = \begin{cases} \frac{1}{b-a}, \\ 0, \end{cases}$    $E[x] = \frac{a+b}{2}$. $SD(x) = \frac{b-a}{\sqrt{12}}$.

Exponential Dist. $f(x) = \lambda e^{-\lambda x}$. $X \sim Exp(\lambda)$. $E(x) = \frac{1}{\lambda}$, $SD(x) = \frac{1}{\lambda}$.

且 $P(X \geq s+t | X \geq s) = P(X \geq t)$. 无记忆

Regression Inference . $SE(b_1) = \dfrac{Se \ (residual)}{S_x \cdot \sqrt{n-1}}$ . hist: $t_{n-2}$ .     Estimate std Error   t·value

$\phantom{xxx}$ Intercept

cI : estimate $\pm t^*_{n-2} \cdot SE(b_1)$     $\square$     k     A     k/A .

$\phantom{xxxxxxx}\underset{b_1}{\Vert}$

$\phantom{xxxxx}$ $H_0: \beta = 0$ . $H_A: \beta \neq 0 . \longrightarrow$ an association .

__Two sample t-test. (Means).__

$\phantom{xx}$ ① Paired $\underline{t\text{-test}}$ . $df = n-1$ . cI : $\bar{d} \pm t^*_{df} \cdot SE_{\bar{d}}$ . $T = \dfrac{\bar{d} - 0}{SE_{\bar{d}}}$ .

$\phantom{xxxx}$ check on differences : Indep of each other ( Random, < 10% pop ) . Diff. nearly normal .

$\phantom{xx}$ ② Unpair two sample t-test $\quad df = min(n_1-1, n_2-1)$ . cI : $(\bar{x}_1 - \bar{x}_2) \pm t^*_{df} SE_{\bar{x}_1 - \bar{x}_2}$ .

$\phantom{xxxx}$ $SE = \sqrt{\dfrac{S_1^2}{n_1} + \dfrac{S_2^2}{n_2}}$ $\qquad T = \dfrac{(\bar{x}_1 - \bar{x}_2) - 0}{SE}$

$\phantom{xxx}$ Inference : Indep. ( Random, < 10% ) . Nearly normal population dist .

$\phantom{xxxxxxx}$ Indep. of two samples .     rejecting fence   $H_0$ dist

$\phantom{xxxxxxxxxxx}$ true dist.

__Power :__ $\qquad$ p ( reject $H_0$ | $H_0$ is false ) .     $\nearrow$ power.

__Means for 3 more samples :__ ANOVA. $\quad H_0: \mu_1 = \mu_2 = \cdots \mu_k$ . $H_A:$ some $\mu_i$ is diff.     in each group.

$\phantom{xx}$ F - dist: $\quad$ Inference : Data are indep. within / across groups ; data in each group nearly normal ; spread roughly

$\phantom{xxxxxxxxxxxxxxxxxxxxxxxxxxxxxxxxxxxxxxxxxxxxxxxxxxxxxxxxxxxxxx}$ equal .

$\phantom{xxxx}$ check : Indep ( random, < 10% ) . side-by-side boxplot $\leftarrow$      DF   Sum Sq  Mean Sq  F·value  P value

$\phantom{xxx}$ mean square across groups

$\phantom{xxx}$ $F = \dfrac{MSG}{MSE} = \dfrac{\frac{1}{k-1} \Sigma n_j (\bar{x}_j - \bar{x})^2}{\frac{1}{n-k} \Sigma (n_1-1)^2 S_j^2}$ $\quad$ k 组数据 $\quad$ year (k-1) $\quad$ A $\quad$ A/(k-1) $\quad \dfrac{A/(k-1)}{B/(n-k)}$ $\quad \cdots$

$\phantom{xxxxx}$ mean square error $\phantom{xxxxxxx}$ 每组 n↑ . $\quad$ Residuals (n-k) $\quad$ B $\quad$ B/(n-k)

$\phantom{xxxxxxxxxxxxxxxxxxxxx}$ margin of error

__Proportions.__ $\quad$ ① one sample , $\quad$ cI : $\hat{p} \pm z^* \cdot SE$ . $\quad$ hypo 看 $N(p, \sqrt{\frac{pq}{n}})$ . $\quad SE = \sqrt{\dfrac{p \cdot q}{n}}$

Margin of Error: $\quad$ Inference : Indep ( random, < 10% ) . At Least 10 success / failures

$z^* \cdot SE.$ $\qquad$ ② two proportions. $\quad$ cI : $(\hat{p}_1 - \hat{p}_2) \pm z^* \cdot SE_{\hat{p}_1 - \hat{p}_2}$ . $\quad SE_{\hat{p}_1 - \hat{p}_2} = \sqrt{\dfrac{p_1 q_1}{n_1} + \dfrac{p_2 q_2}{n_2}}$

$= z^* \cdot \sqrt{\frac{pq}{n}}$ $\qquad$ Inference 同上 , Indep. between two samples . $\quad Z = \dfrac{(\hat{p}_1 - \hat{p}_2) - 0}{SE_{pool}}$ $\quad$ ( 用 $SE_{pool}$ ) .

$\phantom{xxxx}$ $H_0: p_1 = p_2 \Longrightarrow$ hypo : $\hat{p}_{pooled} = \dfrac{success1 + success2}{n_1 + n_2}$ , $SE_{pooled} = \sqrt{\dfrac{\hat{p}_{pool} \hat{q}}{n_1} + \dfrac{\hat{p}_{pool} \hat{q}_{pool}}{n_2}}$

__Counts.__ $\quad$ ① $\chi^2$ : $H_0:$ $\cdots$ matches population , $\quad H_A:$ some difference in $\cdots$

Goodness -of -Fit . $\quad \chi^2 = \dfrac{(O_1 - E_1)^2}{E_1} + \cdots + \dfrac{(O_k - E_k)^2}{E_k}$ , k categories , $\chi^2_{k-1}$ . $\quad df = k-1$

生日、基因型 $\quad$ Inference : ① one-dimensional table of counts $\quad$ ② counts in cells are indep $\quad$ ③ Expected counts ≥ 5

Test for indep : ② 2 more populations. $\quad H_0:$ no association , $\quad H_A:$ some an association .

$\phantom{xxxx}$ expected value : $\dfrac{(row \ total) \cdot (col'n \ total)}{table \ total}$ $\qquad df = (row -1)(col'n -1)$ $\quad$ Inference 同上

__Regression.__

$\phantom{xxx}$ Scatter plot : Form ( lin. curve ) . Direction ( pos. neg ) Strength ( weak. strong ) $\quad$ Outliers

correlation coeff : R : Inference : ① quantitative $\quad$ ② straight enough $\quad$ ③ No outliers $\qquad r = \dfrac{\Sigma (x - \bar{x})(y - \bar{y})}{(n-1) \cdot S_x \cdot S_y}$

$\phantom{xxxx}$ error : $\Sigma (residuals)^2 = \Sigma (obs - exp)^2$ $\quad$ $b_1 = r \cdot \dfrac{S_y}{S_x}$ . $\quad b_0 = \bar{y} - b_1 \bar{x}$ . $\quad$ ④ Indep data .

$\phantom{xx}$ Inference : ① graph roughly linear $\quad$ ② hist of residuals nearly normal $\quad$ ③ constant variability around line

$\phantom{xxx}$ $r^2$ proportion of the variation in y-variable explained by variation in x-variable .

## Stats - 异动归因

- 问题类型: 某metric突然变化，如何分析？
- 首先确认问题**真实性**
  - 读取data，存储新data 是否出错？(e.g. 新上线商品分类错误、计量错误)
  - 问题原本是否有，突然的提升 是否statistically significant ?
- 罗列因素
  - **外部因素**: PEST
    - Policy, Economics, Social, Technology
  - **内部因素**: index拆解
    - 公式分解
    - 维度考虑，e.g. gender, age, operation system, vertical
    - 比较前后的**distribution是否相同**
      - $E[Y] = \sum E[Y_i|X_i] \cdot P(X_i)$，$Y_i$ 变化 & $P(X_i)$分布变化 都有影响
      - **Simpson's Paradox**: https://towardsdatascience.com/simpsons-paradox-and-interpreting-data-6a0443516765
      - 解决方案
        - **Segregate the data in groups** or **aggregate data together**
        - https://towardsdatascience.com/simpsons-paradox-how-to-prove-two-opposite-arguments-using-one-dataset-1c9c917f5ff9