

معیارهای مرکزیت داده‌ها

مد (mode) داده با بیشترین تکرار

میانه (median) داده وسطی. (برای تعداد فرد، میانگین دو داده وسط)

میانگین (mean) $\bar{x} = \frac{\sum_{i=1}^n x_i}{n}$

معیارهای پراکندگی داده‌ها

دامنه تغییرات $R = \max - \min$

میانگین قدر مطلق انحراف از میانگین $MAD = \frac{1}{n} \sum_{i=1}^n |x_i - \bar{x}|$

میانگین مربع انحراف از میانگین $MSD = \frac{1}{n} \sum_{i=1}^n (x_i - \bar{x})^2 = \sum_{i=1}^m (x_i - \bar{x})^2 \frac{f_i}{n}$

جزر میانگین مربع انحراف از میانگین $RMSD = \sqrt{\frac{1}{n} \sum_{i=1}^n (x_i - \bar{x})^2}$

واریانس $\sigma^2 = \frac{1}{n-1} \sum_{i=1}^n (x_i - \bar{x})^2$

انحراف معیار $\sigma = \sqrt{\frac{1}{n-1} \sum_{i=1}^n (x_i - \bar{x})^2}$

احتمال

اصول احتمال کولموگوروف: $\sum P_i = 1$ $P \in [0, 1]$
جمع احتمال: $P(A \cup B) = P(A) + P(B) - P(A \cap B)$
احتمال متمم: $P(\bar{A}) = 1 - P(A)$
استقلال:

• وقوع یکی از پیشامدها روی دیگری تاثیر ندارد

• $P(A \cap B) = P(A)P(B)$

• $P(A|B) = P(A)$ or $P(B|A) = P(B)$

• $P(A|B) = P(A|\bar{B})$ or $P(B|A) = P(B|\bar{A})$

احتمال شرطی:

• $P(A|B) = \frac{P(A \cap B)}{P(B)}$

• $P(A \cap B) = P(B) \times P(A|B) = P(A) \times P(B|A)$

قانون بیز: $P(A_i|B) = \frac{P(B|A_i)P(A_i)}{\sum_{i=1}^n P(A_i)P(B|A_i)}$

متغیر تصادفی گسسته

تابع جرمی احتمال (pmf) $f(x) = P(X = x)$

$$f(x) \geq 0 \quad \forall x \in X, \quad \sum_x f(x) = 1$$

تابع توزیع جمعی (cdf) $F(x) = P(X \leq x) = \sum_{k \leq x} f(k)$

• میانگین/امید ریاضی $\mathbb{E}[X] = \sum_x x f(x)$

• واریانس $Var(X) = \sigma^2 = \sum_x (x - \mu)^2 f(x) = \mathbb{E}[X^2] - (\mathbb{E}[X])^2$

متغیر تصادفی پیوسته

تابع چگالی احتمال (pdf) $P(a \leq X \leq b) = \int_a^b f(x) dx$

$$f(x) \geq 0, \quad \forall x, \quad \int_{-\infty}^{\infty} f(x) dx = 1$$

تابع توزیع جمعی (cdf) $F(x) = P(X \leq x) = \int_{-\infty}^x f(t) dt$

• میانگین/امید ریاضی $\mathbb{E}[X] = \int_{-\infty}^{\infty} X f(X) dx$

• واریانس $Var(X) = \sigma^2 = \int_{-\infty}^{\infty} (x - \mu)^2 f(x) dx = \mathbb{E}[X^2] - (\mathbb{E}[X])^2$

دو متغیر تصادفی

- $\mathbb{E}[aX + bY + c] = a\mathbb{E}[X] + b\mathbb{E}[Y] + c$
- $Var(aX + bY + c) = a^2 Var(X) + b^2 Var(Y) + 2ab Cov(X, Y)$
- $Cov(X, Y) = \mathbb{E}[(X - \mathbb{E}[X])(Y - \mathbb{E}[Y])] = \mathbb{E}[XY] - \mathbb{E}[X]\mathbb{E}[Y]$

تابع احتمال شرطی

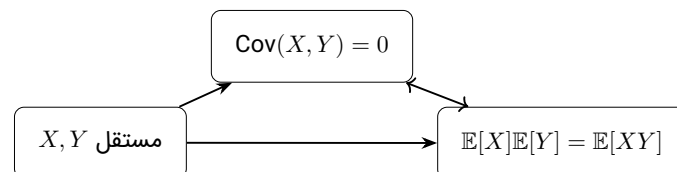
$$p(x, y) = \frac{p(x, y)}{p(y)} \quad \mathbb{E}[X|Y = y] = \sum_x x p(x, y)$$

استقلال

$$Pr(X = x, Y = y) = Pr(X = x).Pr(Y = y)$$

کوواریانس معیاری برای ارتباط دو متغیر تصادفی

$X, Y : Cov(X, Y) > 0$ هم‌راستا $X, Y : Cov(X, Y) < 0$ غیرهم‌راستا



جمع و میانگین متغیر تصادفی‌های مستقل
جمع

$$Y = a_1X_1 + a_2X_2 + \dots + a_nX_n$$

$$\mathbb{E}[Y] = \sum_{i=1}^n a_i \mathbb{E}[X_i] \quad \text{Var}(Y) = \sum_{i=1}^n a_i^2 \text{Var}(X_i)$$

میانگین

$$\bar{X} = \frac{X_1 + X_2 + \dots + X_n}{n}$$

$$\mathbb{E}[\bar{X}] = \frac{1}{n} \sum_{i=1}^n \mathbb{E}[X_i] \quad \text{Var}(\bar{X}) = \frac{1}{n^2} \sum_{i=1}^n \text{Var}(X_i)$$

اعمال تابع

$$\mathbb{E}[g(X)] = \sum_x g(x)p(x) \quad \mathbb{E}[g(X, Y)] = \sum_x g(x, y)p(x, y)$$

بیشینه و کمینه متغیرهای تصادفی مستقل

* برای متغیر تصادفی X_i ، cdf را با $F_{X_i}(x)$ و pdf را با $f_{X_i}(x)$ نشان می‌دهیم*
بیشینه ($V = \max\{X_1, X_2, \dots, X_n\}$)

$$F_V(v) = P(V \leq v) = P(X_1 \leq v, X_2 \leq v, \dots, X_n \leq v)$$

$$= P(X_1 \leq v)P(X_2 \leq v) \dots P(X_n \leq v) = \prod_{i=1}^n F_{X_i}(v)$$

$$f_V(v) = F'_V(v) = \frac{d}{dv} F_V(v)$$

کمینه ($U = \min\{X_1, X_2, \dots, X_n\}$)

$$F_U(u) = P(U \leq u) = 1 - P(U > u) = 1 - P(X_1 > u, X_2 > u, \dots, X_n > u)$$

$$= 1 - P(X_1 > u)P(X_2 > u) \dots P(X_n > u)$$

$$= 1 - \prod_{i=1}^n [1 - F_{X_i}(u)]$$

$$f_U(u) = \frac{d}{du} F_U(u)$$

توزیع t - student

فقط برای نمونه‌های با اندازه کوچک که σ نامعلوم است استفاده می‌شود.
همچنین جامعه باید نرمال باشد.

توزیع‌های پیوسته

توزیع یکنواخت (Uniform): $X \sim U(a, b)$

میانگین:

واریانس:

$$\mu = \mathbb{E}(X) = \frac{a+b}{2}$$

$$\sigma^2 = \text{Var}(X) = \frac{(b-a)^2}{12}$$

cdf

$$F_X(x) = \begin{cases} 0 & x < a \\ \frac{x-a}{b-a} & a \leq x \leq b \\ 1 & x > b \end{cases}$$

pdf

$$f_X(x) = \begin{cases} \frac{1}{b-a} & a \leq x \leq b \\ 0 & \text{otherwise} \end{cases}$$

توزیع نمایی (Exponential): $X \sim \text{Exp}(\lambda)$

میانگین:

واریانس:

$$\mu = E(X) = \frac{1}{\lambda}$$

$$\sigma^2 = \text{Var}(X) = \frac{1}{\lambda^2}$$

cdf

$$F_X(x) = \begin{cases} 1 - e^{-\lambda x} & x \geq 0 \\ 0 & x < 0 \end{cases}$$

pdf

$$f_X(x) = \begin{cases} \lambda e^{-\lambda x} & x \geq 0 \\ 0 & x < 0 \end{cases}$$

توزیع نرمال

توزیع نرمال $X \sim N(\mu, \sigma^2)$

$$p(x) = \frac{1}{\sqrt{2\pi}\sigma} e^{-\frac{(x-\mu)^2}{2\sigma^2}}$$

نرمال استاندارد $Z \sim N(0, 1)$ و $z = \frac{X-\mu}{\sigma}$

$$p(z) = \frac{1}{\sqrt{2\pi}} e^{-z^2/2}$$

• تقریباً 68% داده‌ها در فاصله σ میانگین اند.

• تقریباً 95% داده‌ها در فاصله 2σ میانگین اند.

• تقریباً 99.7% داده‌ها در فاصله 3σ میانگین اند.

قضیه حد مرکزی (CLT)

اگر n نمونه تصادفی از یک جامعه/توزیع با میانگین μ و واریانس σ^2 باشند، برای n به اندازه بزرگ ($n \geq 20$) آنگاه:

$$\bar{X} = \frac{X_1 + X_2 + \dots + X_n}{n} \sim N(\mu, \frac{\sigma^2}{n})$$

این قضیه در مورد جمع نمونه‌ها نیز برقرار است:

$$T = X_1 + X_2 + \dots + X_n = n\bar{X} \quad T \sim N(n\mu, n\sigma^2)$$

متغیر تصادفی برنولی و دوجمله‌ای

برنولی

فضای نمونه دو حالت.

$$Pr(\text{Success}) = p, \quad Pr(\text{Failure}) = 1 - p$$

توزیع برنولی:

$$X \sim \text{Bernoulli}(p)$$

$$\text{pmf: } P(X = x) = p^x(1-p)^{1-x}, \quad x = 0, 1$$

$$\mu = \mathbb{E}[X] = p$$

$$\sigma^2 = \text{Var}(X) = p(1-p)$$

میانگین:

واریانس:

دوجمله‌ای

متغیر تصادفی دوجمله‌ای X ، تعداد موفقیت‌ها در n تلاش مستقل که همه موفقیت‌ها، احتمال برابر p دارند.

توزیع دوجمله‌ای:

$$X \sim \text{Bin}(n, p)$$

$$\text{pmf} \quad P(X = x) = \binom{n}{x} p^x (1-p)^{n-x}$$

$$\text{cdf} \quad P(X \leq x) = \sum_{i=0}^x \binom{n}{i} p^i (1-p)^{n-i}$$

$$\mu = \mathbb{E}[X] = np$$

$$\sigma^2 = \text{Var}(X) = np(1-p)$$

میانگین:

واریانس:

تخمین نرمال توزیع دوجمله‌ای

اگر $X \sim \text{Bin}(n, p)$ و با تصحیح پیوستگی، داشته باشیم $np \geq 5$ و $n(1-p) \geq 5$ آنگاه :

$$X \sim N(np, np(1-p)), \text{ approx.}$$

نمونه‌گیری

• با جاگذاری: نمونه مستقل

• بدون جاگذاری: نمونه وابسته. واریانس از با جاگذاری کمتر است. با نمونه کمتری می‌توان μ را تخمین زد

• نمونه‌گیری نسبتاً کوچک از جامعه بزرگ: فرقی بین با/بدون جاگذاری نیست

$$\mu_{\bar{X}} = \mu \quad \sigma_{\bar{X}} = \frac{\sigma}{\sqrt{n}}$$

بدون جاگذاری

$$\text{Var}(\bar{X}) = \frac{\sigma^2}{n} \frac{m-n}{m-1} \xrightarrow{m \rightarrow \infty} \text{Var}(\bar{X}) = \frac{\sigma^2}{n}$$

• به $(m-n)/(m-1)$ ضریب کاهش گویند

• اگر $n = 1$ ، ضریب کاهش برابر ۱

• اگر $n \gg m$ می‌توان نمونه‌ها را مستقل گرفت. (فرقی بین با/بدون جاگذاری نیست)

نمونه‌گیری از جامعه برنولی

$$P = \frac{\sum X_i}{n} \quad P \sim \mathcal{N}\left(\pi, \frac{\pi(1-\pi)}{n}\right)$$

تفاضل میانگین دو جامعه

با فرض برابری واریانس دو جامعه، تخمین‌گر واریانس جامعه $\mu_1 - \mu_2$

$$s_p^2 = \frac{1}{n_1 + n_2 - 2} \left[\sum_{i=1}^{n_1} (x_{i1} - \bar{x}_1)^2 + \sum_{i=1}^{n_2} (x_{i2} - \bar{x}_2)^2 \right]$$

بازه اطمینان برای $\mu_1 - \mu_2$:

$$(\bar{x}_1 - \bar{x}_2) \pm t_{\frac{\alpha}{2}, (n_1+n_2-2)} s_p \sqrt{\frac{1}{n_1} + \frac{1}{n_2}}$$

جامعه برنولی

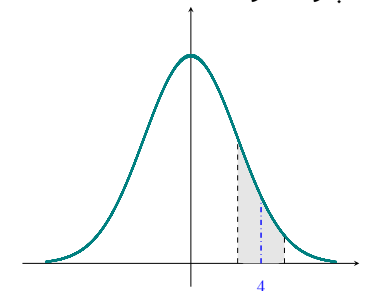
برای n بزرگ:

$$\pi_1 - \pi_2 = P_1 - P_2 \pm z_{\frac{\alpha}{2}} \sqrt{\frac{P_1(1-P_1)}{n_1} + \frac{P_2(1-P_2)}{n_2}}$$

اصلاح پیوستگی

اگر Y پیوسته و X گسسته باشد:

- $P(X > 4) = P(X \geq 5) = P(Y \geq 4.5)$
- $P(X \geq 4) = P(Y \geq 3.5)$
- $P(X < 4) = P(X \leq 3) = P(Y \leq 3.5)$
- $P(X \leq 4) = P(Y \leq 4.5)$
- $P(X = 4) = P(3.5 \leq Y \leq 4.5)$



بازه اطمینان

بازه اطمینان $\alpha\%$:

$$\bar{x} \pm z_{\frac{\alpha}{2}} \frac{\sigma}{\sqrt{n}} \quad \sigma^2 \text{ معلوم:}$$

$$\bar{x} \pm t_{\frac{\alpha}{2}, (n-1)} \frac{s}{\sqrt{n}} \quad \sigma^2 \text{ نامعلوم:}$$

تخمین گر نقطه‌ای

محک

• اوربیبی

$$B_{\hat{\theta}} = \mathbb{E}[\hat{\theta}] - \theta$$

اوربیبی صفر، تخمین گر خوب

• کارایی

کارایی $\hat{\theta}$ نسبت به $\hat{\theta}$:

$$\frac{Var(\hat{\theta})}{Var(\hat{\theta})}$$

برای اوربیب‌ها:

$$\frac{\mathbb{E}[\hat{\theta} - \theta]^2}{\mathbb{E}[\hat{\theta} - \theta]^2}$$

تخمین گر کارا، بهترین تخمین گر

• سازگاری

$$\lim_{n \rightarrow \infty} \mathbb{E}[(\hat{\theta} - \theta)^2] = 0 \iff B_{\hat{\theta}}^2 \rightarrow 0, \quad Var(\hat{\theta}) \rightarrow 0$$

در $n \rightarrow \infty$ ، سازگاری بهترین معیار و در n کوچک، اوربیبی.

انتخاب

• MME عملی که قرار است روی جامعه انجام شود، روی نمونه انجام بده.

• MLE

$$Pr(X_i = x_i; \theta) = \prod_{i=1}^n F(X_i = x_i; \theta) := L(\theta)$$

$L(\theta)$ را بیشینه کن

رگرسیون +

سری‌های زمانی:

بررسی روند تغییر یک پارامتر در طول بازه زمانی.
روش‌های تصحیح داده‌ها:

• متغیر کمکی dummy با مقادیر گسسته

• تصحیح فصلی با میانگین‌گیری. مناسب داده‌های periodic

قضیه گاوس-مارکوف: در کلاس تخمین‌گرهای خطی نااوربیب برای α, β ، کمترین مجموع مربعات، کمترین واریانس را دارد.

آزمون فرض

• H_0 : باور عمومی. چیزی که می‌خواهیم در برابر آن بجنگیم

• H_1 : مشاهده ما. ادعای ما. چیزی که تلاش می‌کنیم اثبات کنیم.

در آزمون فرض، صرفاً با مقایسه دو فرضیه، درستی یا نادرستی H_0 را بررسی می‌کنیم.
انواع خطا

• خطای نوع اول α : رد H_0 به شرط درستی H_0 یا به عبارتی: $Pr(H_0 \times | H_0)$

• خطای نوع دوم β : پذیرفتن H_0 به شرط نادرستی H_0 یا به عبارتی: $Pr(H_0 \checkmark | H_1)$

H_1 در خطای نوع اول دخیل نیست.

خطای نوع اول، علاوه بر H_0 به اندازه نمونه نیز بستگی دارد.

خطای اول و دوم، رفتار الاکلنگی (trade-off) دارند.

α خطای مهمتری از β است. مثال قاضی

با افزایش اندازه نمونه، α کاهش می‌یابد.

تابع قدرت: مقدار $1 - \beta$ را قدرت تست و برای فرضیات مرکب، تابع قدرت گوییم.
دوست داریم H_1 را بپذیریم. پس دوست داریم α کم و قدرت تست بالا داشته باشیم.

p-value:

کمترین احتمال خطای نوع اول که آماره آزمون موجب رد فرض صفر شود، به شرط آنکه فرض صفر صحیح باشد: $p - value = Pr(X > x | H_0)$

H_0 قابل پذیرش \iff بازه اطمینان شامل H_0 باشد

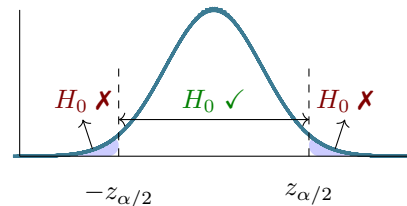
$\alpha > p - value \Rightarrow \text{reject } H_0$

بازه اطمینان، مجموعه‌ای از فرض‌های قابل قبول است.

Rejection Rules:

Consider test statistic z , and significance value α .

- Lower-tail test: Reject H_0 if $z \leq z_{\alpha}$
- Upper-tail test: Reject H_0 if $z \geq z_{\alpha}$
- Two-tail test: Reject H_0 if $|z| \geq z_{\frac{\alpha}{2}}$



برای مقایسه r میانگین. مساله رگرسیون خطی با متغیرهای تمامی dummy

تعداد جامعه‌ها r

تعداد نمونه‌های هر جامعه n

انحراف معیار نمونه i -ام s_i^2

میانگین نمونه i -ام $\bar{x}_i = \frac{1}{n} \sum x_i$

میانگین \bar{X}_i ها $\bar{\bar{x}} = \frac{1}{r} \sum \bar{x}_i$

واریانس \bar{X}_i ها $s_{\bar{X}_i}^2 = \frac{1}{r-1} \sum (\bar{x}_i - \bar{\bar{x}})^2$

میانگین واریانس‌ها $s_p^2 = \frac{1}{r} \sum s_i^2$

$$F = \frac{ns_{\bar{x}}^2}{s_p^2} = \frac{\text{Explained Var}}{\text{Unexplained Var}}$$

فرضیات:

- برای هر جامعه، متغیر مورد نظر، توزیع نرمال دارد.
- واریانس متغیر مورد نظر، برای همه جوامع برابر σ^2 است.
- مشاهدات مستقلند

آزمون فرض:

$$H_0 : \mu_1 = \mu_2 = \dots = \mu_r$$

$$H_1 : \mu_i \neq \mu_j \text{ for } i \neq j$$

بازه اطمینان همزمان

$$\sum C_i \mu_i = \sum C_i \bar{X}_i \pm \sqrt{F_{0.05, sp}} \sqrt{\frac{(r-1)(\sum C_i^2)}{n}}, \quad \sum C_i = 0$$

The ANOVA Table:

F-ratio	Mean Square	Sum of Squares	df	Variation of Source
$\frac{MSSr}{MSSu}$	$\frac{SSr}{r-1} = MSSr$	SSr	$r-1$	Explained
	$\frac{SSu}{r(n-1)} = MSSu$	SSu	$r(n-1)$	Unexplained
		SST	$nr-1$	Total

Test Statistic:

$$v_1 = df(SSr) = r-1$$

$$v_2 = df(SSu) = r(n-1)$$

$$F_{obs} = \frac{MSSr}{MSSu} \sim F_{v_1, v_2}$$

اگر $F_{obs} \geq F_{\alpha, v_1, v_2}$ آنگاه H_0 رد می‌شود.

$$\text{reject } H_0 \iff F > 1$$

اگر F نزدیک ۱ باشد، H_0 را می‌پذیریم.

کمترین مجموع مربعات:

$$\operatorname{argmin}_{\alpha, \beta} \sum_{i=1}^n (y_i - (bx_i + a))^2 \Rightarrow a = \bar{Y} \quad b = \frac{\sum Y_i x_i}{\sum x_i^2}$$

مدلسازی:

$$Y_i = \alpha + \beta x_i + \varepsilon_i \sim \hat{Y}_i = a + bx_i \quad d.f = n-2$$

$$\begin{array}{c|c|c|c|c|c} \mathbb{E}[Y_i] & \alpha + \beta x_i & \mathbb{E}[\varepsilon_i] & 0 & \mathbb{E}[a] & \alpha \\ \hline \operatorname{Var}(Y_i) & \sigma^2 & \operatorname{Var}(\varepsilon_i) & \sigma^2 & \operatorname{Var}(a) & \frac{\sigma^2}{n} \\ & & & & \operatorname{Var}(b) & \frac{\sigma^2}{\sum x_i^2} \end{array}$$

Y_i ها متغیر تصادفی مستقلند

بازه اطمینان:

$$\beta = b \pm t_{0.025}^{(n-2)} \frac{s}{\sqrt{\sum x_i^2}} \quad \alpha = a \pm t_{0.025}^{(n-2)} \frac{s}{\sqrt{n}}$$

بازه اطمینان پیش‌بینی:

$$Y_0 = \hat{Y}_0 \pm t_{0.025}^{(n-2)} s \sqrt{\frac{1}{n} + \frac{x_0^2}{\sum_{i=1}^n x_i^2} + 1}$$

مدلسازی:

$$Y_i = \alpha + \beta x_i + \gamma z_i + \varepsilon_i \sim \hat{Y}_i = a + bx_i + cz_i \quad d.f = n-3$$

کمترین مجموع مربعات:

$$\operatorname{argmin}_{\alpha, \beta, \gamma} \sum_{i=1}^n (y_i - (a + bx_i + cz_i))^2 \Rightarrow a = \bar{Y} \quad \begin{cases} \sum Y_i x_i = b \sum x_i^2 + c \sum x_i z_i \\ \sum Y_i z_i = b \sum x_i z_i + c \sum z_i^2 \end{cases}$$

همخطی چندگانه (multicollinearity): یک بعد را از دست می‌دهیم! وقتی یکی از پارامترها ضریب دیگری باشد یا ارتباط خطی داشته باشند اتفاق می‌افتد.

• رگرسیون استاندارد: رگرسیون با متغیرهای عددی

• ANOVA: رگرسیون با متغیرهای dummy

• ANOCOVA: رگرسیون استاندارد + ANOVA