

① Approximate means:

Suppose X_1, \dots, X_n are $\overbrace{\text{indep. and identically dist.}}^{\text{iid}}$ random variables.

$X_i \in \{-1, 1\}, \quad \mathbb{P}(X_i=1) = p$

$\mathbb{E}[X_i] = 1\mathbb{P}(X_i=1) + (-1)\mathbb{P}(X_i=-1) = 2p-1$

$P \approx \frac{1}{n} \sum_{i=1}^n \mathbb{1}\{X_i=1\}$ $1-p \approx \frac{1}{n} \sum_{i=1}^n \mathbb{1}\{X_i=-1\}$

\uparrow id. func. $\begin{cases} = 1, \text{ if } X_i=1 \\ = 0, \text{ ow} \end{cases}$

$\mathbb{E}[X_i] \approx 1 \left(\frac{1}{n} \sum_{i=1}^n \mathbb{1}\{X_i=1\} \right) + (-1) \left(\frac{1}{n} \sum_{i=1}^n \mathbb{1}\{X_i=-1\} \right)$

$= \frac{1}{n} \left[1 (\text{no. of } X_i\text{'s}=1) + (-1) (\text{no. of } X_i\text{'s}=-1) \right]$

$= \frac{1}{n} \sum_{i=1}^n X_i$

We can generalize this!

X_1, \dots, X_n iid with PDF (prob. density func) $f(x)$.

$\mathbb{E}[X_i] = \int_{-\infty}^{\infty} x f(x) dx \approx \frac{1}{n} \sum_{i=1}^n X_i$

$\mathbb{E}[g(X_i)] = \int_{-\infty}^{\infty} g(x) f(x) dx \approx \frac{1}{n} \sum_{i=1}^n g(X_i)$ for any func. $g: \mathbb{R} \rightarrow \mathbb{R}$

② Decomposing SE:

$L(w) = \int_{-\infty}^{\infty} (\hat{f}_m(x) - f(x))^2 dx = \underbrace{\int_{-\infty}^{\infty} \hat{f}_m(x)^2 dx - 2 \int_{-\infty}^{\infty} \hat{f}_m(x) f(x) dx}_{\text{approximate this!}} + \int_{-\infty}^{\infty} f(x)^2 dx$

\nwarrow unknown constant

③ [Rudemo 1982] Leave-one-out:

$\int_{-\infty}^{\infty} \hat{f}_m(x) f(x) dx \approx \frac{1}{m} \sum_{i=1}^m \hat{f}_m(X_i)$ where X_1, \dots, X_m are iid $f(x)$

Let $\hat{f}_{m,-i}$ be the histogram with samples $X_1, \dots, X_{i-1}, X_{i+1}, \dots, X_m$ (leave X_i out) ($\hat{f}_{m,-i} \approx \hat{f}_m$)

$$\int_{-\infty}^{\infty} \hat{f}_m(x) f(x) dx \approx \frac{1}{m} \sum_{i=1}^m \hat{f}_{m, \tau_i}(X_i) \leftarrow \text{less biased!}$$



④ Algebra...

$$\int_{-\infty}^{\infty} \hat{f}_m(x)^2 dx - 2 \int_{-\infty}^{\infty} \hat{f}_m(x) f(x) dx$$

$$= \sum_{k=1}^n w \left(\frac{\hat{p}_k}{w} \right)^2 - 2 \frac{1}{m} \sum_{i=1}^m \hat{f}_{m, \tau_i}(X_i)$$

If X_i is in bin k , $\frac{m \hat{p}_k - 1}{w(m-1)} = \hat{f}_{m, \tau_i}(X_i)$

$$= \frac{1}{w} \sum_{k=1}^n \hat{p}_k^2 - 2 \frac{1}{m} \sum_{k=1}^n \underbrace{m \hat{p}_k}_{\text{no. of } X_i \text{'s in bin } k} \cdot \frac{m \hat{p}_k - 1}{w(m-1)}$$

$$= \frac{1}{w} \sum_{k=1}^n \hat{p}_k^2 - \frac{2m}{w(m-1)} \sum_{k=1}^n \hat{p}_k^2 + \frac{2}{w(m-1)} \left(\sum_{k=1}^n \hat{p}_k \right) \rightarrow = 1$$

$$= \underbrace{\left[\frac{2}{w(m-1)} - \frac{(m+1)}{w(m-1)} \sum_{k=1}^n \hat{p}_k^2 \right]}_{J(w)}$$

$$L(w) \approx J(w) + \underbrace{\int_{-\infty}^{\infty} f(x)^2 dx}_{\text{constant}}$$

PROBABILITY

Aug 30, 2024
Anuran Makur

o Probability Space: $(\Omega, \mathcal{F}, \mathbb{P})$

1) Sample space Ω - set of outcomes of random experiment
e.g: $\{H, T\}, \{1, 2, 3, 4, 5, 6\}, [0, 1], \mathbb{R}$

2) σ -algebra \mathcal{F} - collection of subsets/events of Ω to which we will assign prob.s

- $\emptyset \in \mathcal{F}$
- $A \in \mathcal{F} \Rightarrow A^c \in \mathcal{F}$ "complement"
- $A_1, A_2, A_3, \dots \in \mathcal{F} \Rightarrow \bigcup_{i=1}^{\infty} A_i \in \mathcal{F}$

3) Probability measure $\mathbb{P}: \mathcal{F} \rightarrow [0, 1]$ (σ -additive set func.)

- $\mathbb{P}(\Omega) = 1$ [normalization]
 - $A_1, A_2, A_3, \dots \in \mathcal{F}$ s.t. $A_i \cap A_j = \emptyset$ (disjoint)
- $$\mathbb{P}\left(\bigcup_{i=1}^{\infty} A_i\right) = \sum_{i=1}^{\infty} \mathbb{P}(A_i)$$

[σ -additivity]

} Kolmogorov axioms

Example: (Coin toss) $\Omega = \{H, T\}$, $\mathcal{F} = \{\emptyset, \{H\}, \{T\}, \{H, T\}\}$
 $\mathbb{P}(\emptyset) = 0$, $\mathbb{P}(H) = \frac{1}{2}$, $\mathbb{P}(T) = \frac{1}{2}$, $\mathbb{P}(\Omega) = 1$

o Properties:

1) Monotonicity: $A \subseteq B \Rightarrow \mathbb{P}(A) \leq \mathbb{P}(B)$



2) Inclusion-Exclusion: $\mathbb{P}(A \cup B) = \mathbb{P}(A) + \mathbb{P}(B) - \mathbb{P}(A \cap B)$



3) Union Bound (Boole): $\mathbb{P}\left(\bigcup_{i=1}^{\infty} A_i\right) \leq \sum_{i=1}^{\infty} \mathbb{P}(A_i)$

↑ need not be disjoint

"is a subset of"

4) Continuity: $A_1 \subseteq A_2 \subseteq A_3 \subseteq \dots \Rightarrow \lim_{n \rightarrow \infty} \mathbb{P}(A_n) = \mathbb{P}\left(\bigcup_{i=1}^{\infty} A_i\right)$
 $A_1 \supseteq A_2 \supseteq A_3 \supseteq \dots \Rightarrow \lim_{n \rightarrow \infty} \mathbb{P}(A_n) = \mathbb{P}\left(\bigcap_{i=1}^{\infty} A_i\right)$

} equivalent to σ -additivity axiom

o Conditional Probability: If A is an event with $\mathbb{P}(A) > 0$, then the cond. prob. of B given A is $\mathbb{P}(B|A) \triangleq \frac{\mathbb{P}(B \cap A)}{\mathbb{P}(A)}$.

$\rightarrow (A, \{A \cap B: B \in \mathcal{F}\}, \mathbb{P}(\cdot|A))$ is a prob. space!

$$\mathbb{P}(B|A) = \frac{\mathbb{P}(B \cap A)}{\mathbb{P}(A)} = \frac{\mathbb{P}(A|B) \mathbb{P}(B)}{\mathbb{P}(A)} \quad [\text{Bayes Rule}]$$

◦ Independence: (occurrence of event A does not affect event B)

A and B are independent events if $\underbrace{P(A|B)} = P(A)$.

$$\frac{P(A \cap B)}{P(B)}$$

$$P(A \cap B) = P(A)P(B)$$

← $P(A)$ or $P(B)$ could be 0!

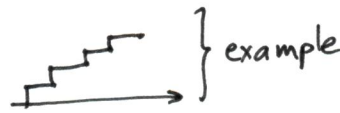
◦ Random Variable: A random variable is a function $X: \Omega \rightarrow \mathbb{R}$ such that for every $c \in \mathbb{R}$, $\{\omega \in \Omega : X(\omega) \leq c\} \in \mathcal{F}$.

The "probability dist." of X is characterized by the probabilities of $\underbrace{\{\omega \in \Omega : X(\omega) \leq c\}}_{\{X \leq c\} \text{ (notation)}}$.

Cumulative Distribution Function (CDF): For rand. var. X , its CDF is $F_X: \mathbb{R} \rightarrow [0, 1]$, $F_X(x) \triangleq P(X \leq x)$. ← All rand. vars have CDF!

3 "pure" rand. vars: A pure rand. variable X and its CDF F has one of 3 forms:

① Discrete rand. var.: X takes countably many values with > 0 probability and F is discontinuous with jumps, & piecewise constant.



② Continuous rand. var.: X takes uncountably many values (e.g. $[0, 1]$) and F absolutely continuous.

↙ (Think: continuously differentiable)
 \exists a derivative $f = F'$ s.t. $F(x) = \int^x f(t) dt$.

We refer to f as the probability density function (PDF).

③ Singular rand. var.: F is continuous but not abs. cont. } not for exam!
 eg: Cantor function (Devil's staircase)

Decomposition Thm: The probability dist. of a general rand. var., P_X , can be uniquely decomposed into a mixture of the 3 pure components:

$$P_X = \lambda_1 P_{\text{disc.}} + \lambda_2 P_{\text{cont.}} + \lambda_3 P_{\text{sing.}}$$

$$\lambda_1, \lambda_2, \lambda_3 \geq 0, \quad \lambda_1 + \lambda_2 + \lambda_3 = 1 \quad \left. \vphantom{\lambda_1, \lambda_2, \lambda_3} \right\} \text{convex weights}$$

o Continuous rand. variables:

Prop: A function $f: \mathbb{R} \rightarrow \mathbb{R}$ is a PDF of some rand. var. iff

① $f(x) \geq 0$ for all $x \in \mathbb{R}$,

② $\int_{-\infty}^{\infty} f(x) dx = 1$.

o CDF: Prop: A function $F: \mathbb{R} \rightarrow [0, 1]$ is a CDF of some rand. var. iff

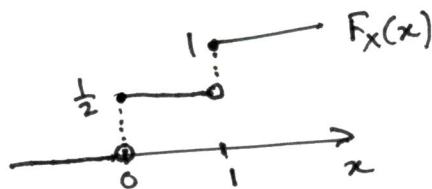
① $\lim_{x \rightarrow -\infty} F(x) = 0$,

② $\lim_{x \rightarrow \infty} F(x) = 1$,

③ If $x \leq y$, then $F(x) \leq F(y)$,

④ F is right-continuous.

Example $X \in \{0, 1\}$, $X \sim \text{Ber}(\frac{1}{2})$



o Discrete rand. var: X takes values in a set $\mathcal{X} \subseteq \mathbb{R}$

↑ countable

Prop: A function $f: \mathcal{X} \rightarrow [0, 1]$ is a PMF iff

① $f(x) \geq 0$ for all $x \in \mathcal{X}$,

② $\sum_{x \in \mathcal{X}} f(x) = 1$.

• Expected Value:

① Discrete: $X \in \mathbb{Z}$ $\mathbb{E}[X] = \sum_{x \in \mathbb{Z}} x f_x(x)$, $\mathbb{E}[g(x)] = \sum_{x \in \mathbb{Z}} g(x) f_x(x)$

↙ range(x)
↙ PMF
↙ Function g(·)
↙ PMF

② Continuous: $\mathbb{E}[X] = \int_{-\infty}^{\infty} x f_x(x) dx$, $\mathbb{E}[g(x)] = \int_{-\infty}^{\infty} g(x) f_x(x) dx$

↑ PDF
↑ PDF

$g(x) = x^2$ $\mathbb{E}[g(x)] = \mathbb{E}[X^2] = \int_{-\infty}^{\infty} x^2 f_x(x) dx$

$g(x) = x$ $\mathbb{E}[g(x)] = \mathbb{E}[X] = \int_{-\infty}^{\infty} x f_x(x) dx$

Linearity: For two r.v.s X, Y and two constants $a, b \in \mathbb{R}$,
 $\mathbb{E}[aX + bY] = a \mathbb{E}[X] + b \mathbb{E}[Y]$.

Variance: $\text{var}(x) \triangleq \mathbb{E}[(x - \mathbb{E}[x])^2] \geq 0$

$$= \mathbb{E}[x^2 - 2x \mathbb{E}[x] + \mathbb{E}[x]^2]$$

$$= \mathbb{E}[x^2] - 2 \mathbb{E}[x] \mathbb{E}[x] + \mathbb{E}[x]^2$$

$$= \boxed{\mathbb{E}[x^2] - \mathbb{E}[x]^2} \geq 0$$

• Probability plots:

① P-P plot: ← probability

Given two CDFs F and G , their P-P plot is the graph $\{(F(x), G(x)) : x \in \mathbb{R}\}$.

② Q-Q plot: ← quantile ↑ hard to "grid" this

Given two CDFs F and G , and their generalized inverses F^{-1} and G^{-1} , their Q-Q plot is the graph $\{(F^{-1}(q), G^{-1}(q)) : q \in [0, 1]\}$.
 ↑ can "grid" this!

* Generalized inverse: Given CDF $F: \mathbb{R} \rightarrow [0, 1]$, its gen. inverse is $F^{-1}: [0, 1] \rightarrow \mathbb{R} \cup \{\pm\infty\}$, $F^{-1}(q) \triangleq \inf\{x \in \mathbb{R} : F(x) \geq q\}$.
 ↑ "infimum" [min if set is closed]
 eg: $\inf(0, \infty) = 0$ but min doesn't exist

• Example: $X \stackrel{\text{"is dist. as"}}{\sim} \text{Bernoulli}(p)$
 $E[X] = 0(1-p) + 1(p) = \underline{p}$
 $\text{var}(X) = E[X^2] - E[X]^2$
 $= (0^2(1-p) + 1^2(p)) - p^2$
 $= p - p^2$
 $= \underline{p(1-p)}$

Example: $X \sim \text{Exp}(\lambda)$
 $E[X] = \int_0^{\infty} x \lambda e^{-\lambda x} dx = \frac{1}{\lambda} \int_0^{\infty} u e^{-u} du$
 $= \frac{1}{\lambda} \left(\lim_{u \rightarrow \infty} \frac{u}{e^u} = 0 + \int_0^{\infty} e^{-u} du \right)$ [integrate by parts.]
 $= \frac{1}{\lambda}$
 $\text{var}(X) = \frac{1}{\lambda^2}$ ← Exercise

Example: $X \sim \text{Binomial}(n, p)$
 X_1, X_2, \dots, X_n iid $\text{Ber}(p)$ variables
 $X = X_1 + X_2 + \dots + X_n$
 $P(X = k) = \binom{n}{k} p^k (1-p)^{n-k}$
 $E[X] = E[X_1 + \dots + X_n]$
 $= E[X_1] + \dots + E[X_n]$
 $= \underline{np}$

Binomial coefficient
 $\binom{n}{k} = \frac{n(n-1)(n-2)\dots(n-k+1)(n-k)\dots 2 \cdot 1}{k!(n-k)\dots 2 \cdot 1}$
 $= \frac{n!}{k!(n-k)!}$
 $\text{var}(X) = \underline{np(1-p)}$ ← Exercise

① Properties of mean and variance:

1) For rand. var.s X, Y and constants $a, b, c \in \mathbb{R}$,

$$\mathbb{E}[aX + bY + c] = a\mathbb{E}[X] + b\mathbb{E}[Y] + c.$$

2) For rand. var.s X, Y that are independent, and constants $a, b, c \in \mathbb{R}$,

$$\text{var}(aX + bY + c) = a^2 \text{var}(X) + b^2 \text{var}(Y).$$

② Unbiased estimators of mean and variance:

Suppose X_1, \dots, X_n are iid rand. var.s with $\mathbb{E}[X] = \mu$ and variance $\text{var}(X) = \sigma^2$.

$$\bar{X} = \frac{1}{n} \sum_{i=1}^n X_i$$

 ↑ sample mean

$$s^2 = \frac{1}{n-1} \sum_{i=1}^n (X_i - \bar{X})^2$$

 ↑ sample variance
 ↖ Bessel's correction

$$\mathbb{E}[\bar{X}] = \mathbb{E}\left[\frac{1}{n} \sum_{i=1}^n X_i\right] = \frac{1}{n} \sum_{i=1}^n \underbrace{\mathbb{E}[X_i]}_{\mu} = \mu \quad [\text{unbiased}]$$

\bar{X} is an unbiased estimator of μ .

$$\mathbb{E}\left[\sum_{i=1}^n (X_i - \bar{X})^2\right] = \mathbb{E}\left[\sum_{i=1}^n X_i^2 - 2X_i\bar{X} + \bar{X}^2\right] = \mathbb{E}\left[\sum_{i=1}^n X_i^2 - 2\bar{X} \sum_{i=1}^n X_i + \sum_{i=1}^n \bar{X}^2\right]$$

$$= \mathbb{E}\left[\sum_{i=1}^n X_i^2 - 2n\bar{X}^2 + n\bar{X}^2\right]$$

$$= \sum_{i=1}^n \mathbb{E}[X_i^2] - n\mathbb{E}[\bar{X}^2]$$

$$= n\mathbb{E}[X^2] - n\mathbb{E}[\bar{X}^2]$$

$$= n(\sigma^2 + \mu^2) - n\left(\frac{\text{var}(\bar{X})}{\sigma^2/n} + \underbrace{\mathbb{E}[\bar{X}^2]}_{\mu^2}\right)$$

$$= n\sigma^2 + n\mu^2 - n\left(\frac{\sigma^2}{n}\right) - n\mu^2$$

$$= (n-1)\sigma^2$$

$$\underbrace{\mathbb{E}\left[\frac{1}{n-1} \sum_{i=1}^n (X_i - \bar{X})^2\right]}_{s^2} = \sigma^2$$

$$\begin{aligned} \text{var}(X) &= \mathbb{E}[X^2] - \mathbb{E}[X]^2 \\ \text{var}(\bar{X}) &= \text{var}\left(\frac{1}{n} \sum_{i=1}^n X_i\right) \\ &= \frac{1}{n^2} \sum_{i=1}^n \text{var}(X_i) \\ &= \frac{1}{n^2} n\sigma^2 \\ &= \frac{\sigma^2}{n} \end{aligned}$$

③ Tail bounds:

• Markov's inequality: For any rand. var. $X \geq 0$ and any constant $a > 0$,

$$\mathbb{P}(X \geq a) \leq \frac{\mathbb{E}[X]}{a}.$$

Pf: (continuous)

$$\begin{aligned} \mathbb{E}[X] &= \int_0^{\infty} x f_X(x) dx = \underbrace{\int_0^a x f_X(x) dx}_{\geq 0} + \int_a^{\infty} x f_X(x) dx \geq \int_a^{\infty} x f_X(x) dx \\ &\geq \int_a^{\infty} a f_X(x) dx = a \int_a^{\infty} f_X(x) dx = a \mathbb{P}(X \geq a). \quad \square \end{aligned}$$

• Chebyshev's inequality: For any rand. var. X and any constant $a > 0$,

$$\mathbb{P}(|X - \mathbb{E}[X]| \geq a) \leq \frac{\text{var}(X)}{a^2}.$$

Pf: Let $Y = (X - \mathbb{E}[X])^2 \geq 0$. Then, by Markov,

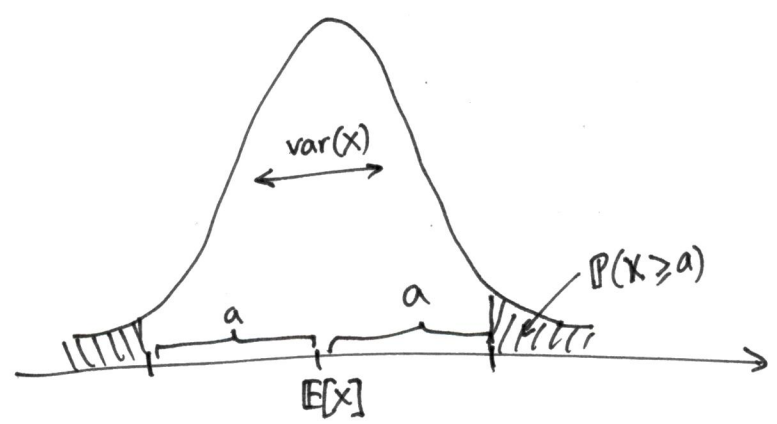
$$\begin{aligned} \mathbb{P}(Y \geq a^2) &\leq \frac{\mathbb{E}[Y]}{a^2} \leftarrow \mathbb{E}[Y] = \mathbb{E}[(X - \mathbb{E}[X])^2] = \text{var}(X) \\ &= \frac{\text{var}(X)}{a^2} \end{aligned}$$

Argue this!

$$\{Y \geq a^2\} = \{(X - \mathbb{E}[X])^2 \geq a^2\} \stackrel{*}{=} \{|X - \mathbb{E}[X]| \geq a\}$$

$$\Rightarrow \mathbb{P}(|X - \mathbb{E}[X]| \geq a) = \mathbb{P}(Y \geq a^2) \leq \frac{\text{var}(X)}{a^2}. \quad \square$$

Picture:



④ • Law of Large Numbers: (weak)

For iid random variables X_1, \dots, X_n with mean $E[X] = \mu$ and variance $\text{var}(X) = \sigma^2$,

we have:

$$\forall \varepsilon > 0, \lim_{n \rightarrow \infty} P\left(\left| \underbrace{\frac{1}{n} \sum_{i=1}^n X_i}_{\bar{X}} - \underbrace{\mu}_{E[X]} \right| \geq \varepsilon\right) = 0.$$

Pf: By Chebyshev, for any $\varepsilon > 0$,

$$0 \leq P\left(\left| \underbrace{\frac{1}{n} \sum_{i=1}^n X_i}_{\bar{X}} - \underbrace{\mu}_{E[X]} \right| \geq \varepsilon\right) \leq \frac{\text{var}(\bar{X})}{\varepsilon^2} = \frac{\sigma^2}{n\varepsilon^2}$$

$$\lim_{n \rightarrow \infty} P(|\bar{X} - E[\bar{X}]| \geq \varepsilon) = 0. \quad \square$$

⑤ • Central Limit Theorem: (Lindeberg-Lévy)

For iid random variables X_1, \dots, X_n with $E[X] = \mu$ and $\text{var}(X) = \sigma^2$,

$$\forall x \in \mathbb{R}, \lim_{n \rightarrow \infty} P\left(\underbrace{\frac{1}{\sqrt{n}\sigma} \sum_{i=1}^n (X_i - \mu)}_{\text{mean 0 and variance 1}} \leq x\right) = \underbrace{\int_{-\infty}^x \frac{1}{\sqrt{2\pi}} e^{-t^2/2} dt}_{\text{standard Normal CDF with mean 0 and variance 1}}.$$

Slides statement:

$$\frac{\sigma}{\sqrt{n}} \left(\frac{1}{\sqrt{n}\sigma} \sum_{i=1}^n (X_i - \mu) \right) \stackrel{n \rightarrow \infty}{\approx} N\left(0, \frac{\sigma^2}{n}\right)$$

$$\rightarrow \mu + \left(\frac{1}{n} \sum_{i=1}^n (X_i - \mu) \right) \approx N\left(\mu, \frac{\sigma^2}{n}\right)$$

$$\rightarrow \boxed{\frac{1}{n} \sum_{i=1}^n X_i \approx N\left(\mu, \frac{\sigma^2}{n}\right)} \quad (\text{in slides})$$

Induced Distributions:

① Change-of-Variables: (linear case)

Let X be a continuous random var. with PDF f_x and CDF F_x .

For any $a \neq 0$ and $b \in \mathbb{R}$, let $Y = \underline{aX + b}$ be another continuous rand. var. with PDF f_y and CDF F_y .

Can we write f_y in terms of f_x ?

$$F_y(y) = \mathbb{P}(Y \leq y) = \mathbb{P}(aX + b \leq y) = \begin{cases} \mathbb{P}\left(X \leq \frac{y-b}{a}\right), & a > 0 \\ \mathbb{P}\left(X \geq \frac{y-b}{a}\right), & a < 0 \end{cases}$$

$$= \begin{cases} F_x\left(\frac{y-b}{a}\right), & a > 0 \\ 1 - F_x\left(\frac{y-b}{a}\right), & a < 0 \end{cases}$$

$$\begin{array}{c} f_y(y) \\ \uparrow \\ \text{PDF} \end{array} = \frac{d}{dy} F_y(y) = \begin{cases} f_x\left(\frac{y-b}{a}\right) \frac{1}{a}, & a > 0 \\ -f_x\left(\frac{y-b}{a}\right) \frac{1}{a}, & a < 0 \end{cases} = \frac{1}{|a|} f_x\left(\frac{y-b}{a}\right)$$

[Chain rule of calculus]

$$\boxed{f_y(y) = \frac{1}{|a|} f_x\left(\frac{y-b}{a}\right)} \quad \text{for all } y \in \mathbb{R}$$

② Convolution: (discrete)

Let X and Y be discrete rand. vars with values in \mathbb{Z} and PMFs f_x and f_y , respectively. Assume X and Y are independent. Let $Z = X + Y$ be a disc. rand. var. with PMF f_z . Can we write f_z in terms of f_x and f_y ?

$$\begin{aligned}
 f_z(k) &= \mathbb{P}(Z=k) = \mathbb{P}(\exists j \in \mathbb{Z}, X=j \text{ and } Y=k-j) \\
 &= \sum_{j=-\infty}^{\infty} \mathbb{P}(X=j \text{ and } Y=k-j) = \sum_{j=-\infty}^{\infty} \mathbb{P}(X=j) \mathbb{P}(Y=k-j) \\
 &= \sum_{j=-\infty}^{\infty} f_x(j) f_y(k-j)
 \end{aligned}$$

Define convolution: $(f_x \star f_y)(k) \stackrel{\text{"convolve"}}{=} \sum_{j=-\infty}^{\infty} f_x(j) f_y(k-j)$

$f_z = f_x \star f_y$

Note: If X and Y were cont. r.v.s with PDFs f_x, f_y .
 Let $Z = X + Y$ have PDF f_z .
 $f_z(t) = \int_{-\infty}^{\infty} f_x(u) f_y(t-u) du$

③ Fact: If X and Y are independent Gaussians with means μ_x and μ_y , and variances σ_x^2 and σ_y^2 , respectively,
 then for any $a, b, c \in \mathbb{R}$, $aX + bY + c$ is Gaussian
 with mean $a\mu_x + b\mu_y + c$ and variance $a^2\sigma_x^2 + b^2\sigma_y^2$.

↓
 can be generalized.

Two-Sample Test & Confidence Interval:

- $\bar{x}_0 \sim N(\mu_0, \frac{\sigma_0^2}{n_0})$ [CLT] (Sample 1) , $\bar{x}_1 \sim N(\mu_1, \frac{\sigma_1^2}{n_1})$ [CLT] (Sample 2) (\bar{x}_0 and \bar{x}_1 are independent)

Statistic

$$\bar{x} = \bar{x}_0 - \bar{x}_1 \sim N(\underbrace{\mu_0 - \mu_1}_{=0}, \frac{\sigma_0^2}{n_0} + \frac{\sigma_1^2}{n_1})$$

$$\Rightarrow \bar{x} \sim N(0, \frac{\sigma_0^2}{n_0} + \frac{\sigma_1^2}{n_1})$$

This is used in two-sample z-test.

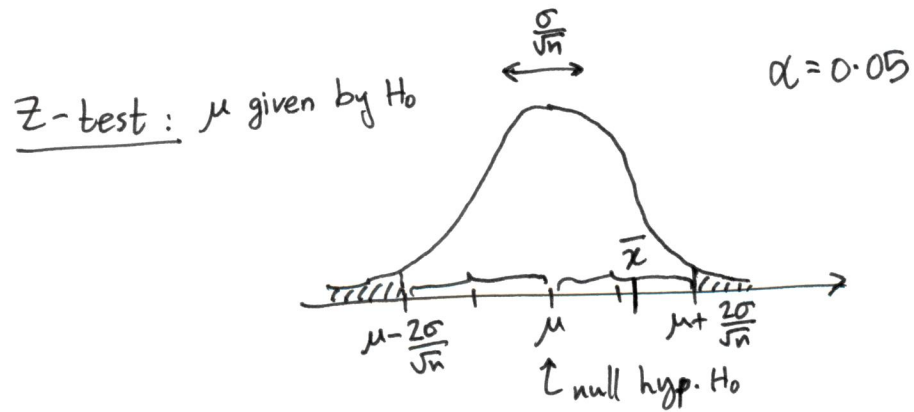
Recall: $\text{var}(aX + bY) = a^2 \text{var}(X) + b^2 \text{var}(Y)$

- A confidence level $c \in [0, 1]$ refers to the following prob:
 confidence interval

$$c = P(\mu \in (\bar{x} - u, \bar{x} + u)) = P(|\bar{x} - \mu| < u)$$

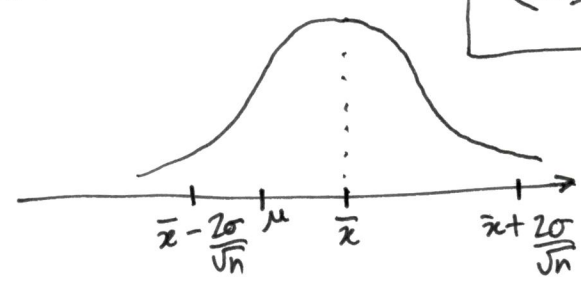
What is this? $u = z_c \frac{\sigma}{\sqrt{n}}$ (slides)

Relation to z-test



Does \bar{x} live in $(\mu - \frac{2\sigma}{\sqrt{n}}, \mu + \frac{2\sigma}{\sqrt{n}})$?

Confidence int.: μ unknown



Note:
 $\bar{x} \in (\mu - \frac{2\sigma}{\sqrt{n}}, \mu + \frac{2\sigma}{\sqrt{n}})$
 $\Leftrightarrow \mu \in (\bar{x} - \frac{2\sigma}{\sqrt{n}}, \bar{x} + \frac{2\sigma}{\sqrt{n}})$

T-DISTRIBUTION:

Sep 20, 2024
Anuran Makur

- Cauchy Distribution: (Ex. t-dist with $\nu=1$)

$$X \sim \text{Cauchy} \quad f_X(x) = \frac{1}{\pi(1+x^2)} \rightarrow \mathbb{E}[X], \text{var}(x) \text{ undefined!}$$

↑ PDF

$$\left. \begin{aligned} \int_{-\infty}^{\infty} \frac{1}{1+x^2} dx &= \int_{-\frac{\pi}{2}}^{\frac{\pi}{2}} \frac{1}{1+\tan^2\theta} (1+\tan^2\theta) d\theta = \int_{-\frac{\pi}{2}}^{\frac{\pi}{2}} 1 d\theta = \underline{\underline{\pi}} \\ &\quad \uparrow x = \tan\theta \\ &\quad \quad = \frac{\sin\theta}{\cos\theta} \\ \frac{dx}{d\theta} &= \frac{\cos\theta}{\cos^2\theta} + \frac{\sin\theta(-1)(-\sin\theta)}{\cos^2\theta} = 1 + \tan^2\theta \end{aligned} \right\} \text{normalization}$$

$$\mathbb{E}[X] = \int_{-\infty}^{\infty} x f_X(x) dx = \int_{-\infty}^{\infty} \frac{x}{\pi(1+x^2)} dx = \frac{1}{2\pi} \int_{-\infty}^{\infty} \frac{2x}{1+x^2} dx = \frac{1}{2\pi} \left[\log(1+x^2) \right]_{-\infty}^{\infty} = \underline{\underline{\infty - \infty}}$$

-
- t-confidence interval: $\left(\bar{x} - t_c \frac{s}{\sqrt{n}}, \bar{x} + t_c \frac{s}{\sqrt{n}} \right)$
 - sample mean
 - t_c - value
 - sample std. dev
 - no. of samples

Binomial Coefficients:

- Recall $X \sim \text{Binomial}(n, p)$
 \uparrow no. of trials (pos. integer) \uparrow success prob. $\in [0, 1]$

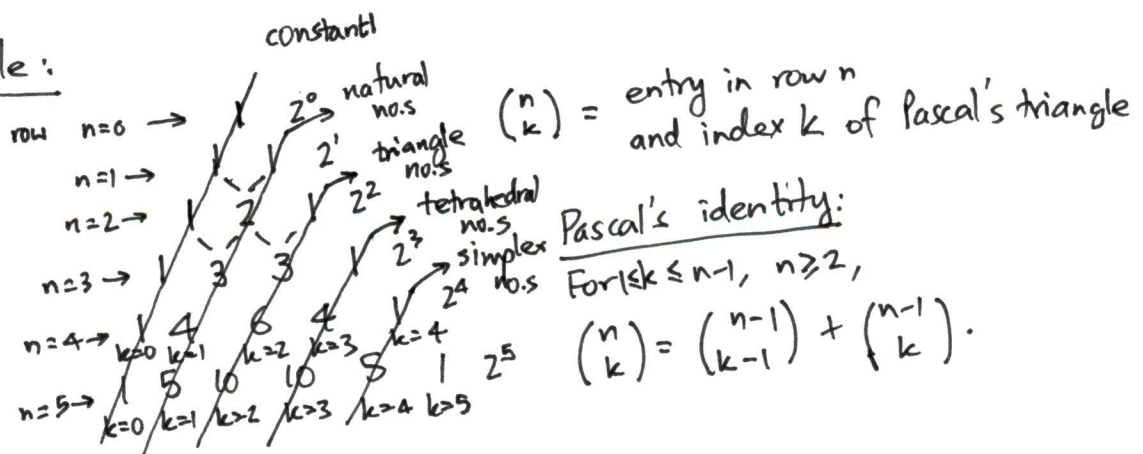
PMF: $\mathbb{P}(X=k) = \binom{n}{k} p^k (1-p)^{n-k}$, $k=0, 1, 2, \dots, n-1, n$
 \uparrow binomial coefficient

$$\sum_{k=0}^n \mathbb{P}(X=k) = \sum_{k=0}^n \binom{n}{k} p^k (1-p)^{n-k} \stackrel{\text{binomial thm}}{=} (p+1-p)^n = 1$$

- Recall: $\binom{n}{k} = \frac{n!}{k!(n-k)!}$ \leftarrow no. of ways to choose k elements from a set of n elements
 \uparrow Factorial: $k! = k(k-1)(k-2)\dots 3 \cdot 2 \cdot 1$, $0! = 1$

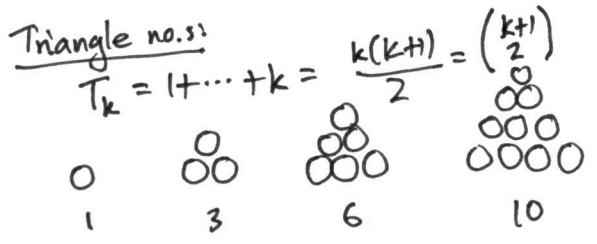
from hereon, not for exam

Pascal's Triangle:

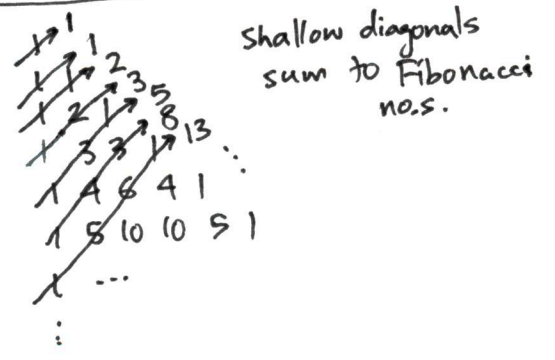


Pascal's identity:
 For $0 \leq k \leq n-1$, $n \geq 2$,
 $\binom{n}{k} = \binom{n-1}{k-1} + \binom{n-1}{k}$

Powers of 2: $\sum_{k=0}^n \binom{n}{k} = (1+1)^n = 2^n$
 \uparrow binomial thm

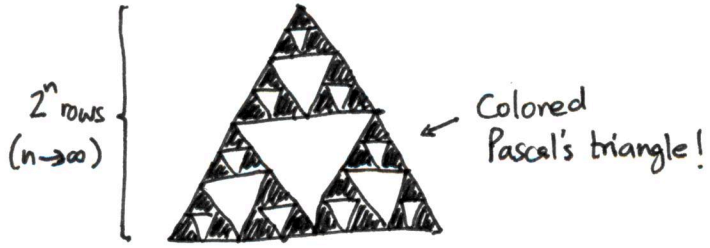


Shallow Diagonals:



Sierpiński Triangle:

- Color all odd numbers black and even numbers white.



- Famous example of fractal!

LINEAR REGRESSION:

Oct 14, 2024
Anuran Makur

① Gradient:

Given a differentiable function $f: \mathbb{R}^d \rightarrow \mathbb{R}$, its gradient is the vector field $\nabla f: \mathbb{R}^d \rightarrow \mathbb{R}^d$ given by:

$$\forall x \in \mathbb{R}^d, \quad \nabla f(x) = \begin{bmatrix} \frac{\partial f}{\partial x_1}(x) \\ \frac{\partial f}{\partial x_2}(x) \\ \vdots \\ \frac{\partial f}{\partial x_d}(x) \end{bmatrix}.$$

② Quadratic Form:

Given fixed $A \in \mathbb{R}^{d \times d}$, define $f(x) = \overbrace{x^T A x}^{\text{real no.}}$ for $x \in \mathbb{R}^d$.

Prop: For $x \in \mathbb{R}^d$, $\nabla f(x) = (A + A^T)x$.

Pf: Observe $f(x) = \sum_{i=1}^d x_i [Ax]_i = \sum_{i=1}^d x_i \sum_{j=1}^d A_{ij} x_j = \sum_{i=1}^d \sum_{j=1}^d x_i x_j A_{ij}$.

$$\frac{\partial f}{\partial x_k}(x) = \frac{\partial}{\partial x_k} \left(x_k^2 A_{kk} + \sum_{j \neq k} x_k x_j A_{kj} + \sum_{i \neq k} x_i x_k A_{ik} \right)$$

$$= 2x_k A_{kk} + \sum_{j \neq k} x_j A_{kj} + \sum_{i \neq k} x_i A_{ik}$$

$$= \sum_{i=1}^d x_i \underbrace{A_{ik}}_{[A^T]_{ki}} + \sum_{j=1}^d x_j \underbrace{A_{kj}}_{[Ax]_k}$$

$$= [A^T x]_k + [Ax]_k$$

$$= [(A^T + A)x]_k \quad \blacksquare$$

③ Linear Form:

Given fixed $b \in \mathbb{R}^d$, define $f(x) = b^T x$ for $x \in \mathbb{R}^d$.

Prop: For $x \in \mathbb{R}^d$, $\nabla f(x) = b$.

Pf: $\frac{\partial f}{\partial x_k}(x) = \frac{\partial}{\partial x_k} \sum_{i=1}^d b_i x_i = b_k$. \(\blacksquare\)

④ Normal Equation:

Recall: $y \in \mathbb{R}^N$ ← no. of data samples
 ↖ target vector
 $X \in \mathbb{R}^{N \times (M+1)}$ ← no. of explanatory variables
 ↖ feature matrix

} given dataset

Mean-squared Error: $E(\beta) = \frac{1}{N} \|y - X\beta\|_2^2$ for $\beta \in \mathbb{R}^{M+1}$
 ↖ minimize
 ↖ regression coefficients
 ↖ ℓ^2 norm

$$\begin{aligned} \Rightarrow E(\beta) &= \frac{1}{N} (y - X\beta)^T (y - X\beta) \\ &= \frac{1}{N} (y^T - \beta^T X^T) (y - X\beta) = \frac{1}{N} [y^T y - y^T X\beta - \beta^T X^T y + \beta^T X^T X\beta] \\ &= \frac{1}{N} [y^T y - 2(X^T y)^T \beta + \beta^T (X^T X)\beta] \end{aligned}$$

↖ $-(X^T y)^T \beta$

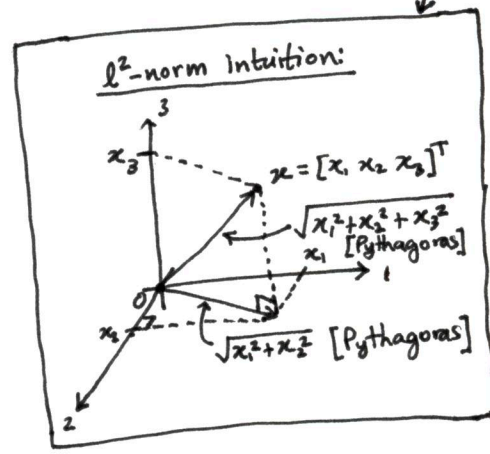
$$\begin{aligned} \Rightarrow \nabla E(\beta) &= -\frac{2}{N} X^T y + \frac{1}{N} (X^T X + X^T X)\beta \\ &= \frac{2}{N} (X^T X\beta - X^T y) \end{aligned}$$

To min. $E(\beta)$, we use stationarity condition:

$$\nabla E(\beta) = 0 \leftarrow \text{zero vector}$$

$$\Leftrightarrow \frac{2}{N} (X^T X\beta - X^T y) = 0$$

$$\Leftrightarrow \boxed{X^T X\beta = X^T y} \left. \vphantom{\boxed{X^T X\beta = X^T y}} \right\} \text{normal equations}$$



Note: $\nabla^2 E(\beta) = \frac{2}{N} X^T X$
 is positive semidefinite
 \Rightarrow Stationary point is global minimum

⑤ Solutions of Normal Eq:

• Normal Equation always has solns because $\text{row}(X) = \text{col}(X^T X)$.
 ↖ rowspace ↖ col-space

Let $y = Xr + v$, for $r \in \mathbb{R}^{M+1}$ and $v \in \text{left-null}(X)$.
 ↖ left-nullspace(X)
 ↖ $(X^T v = 0)$
 ↖ unique decomposition of \mathbb{R}^N into $\text{col}(X)$ and its \perp -complement.

$$\Rightarrow X^T y = X^T (Xr + v) = X^T Xr + X^T v = X^T Xr$$

\Rightarrow Normal eq. has solution!

• Normal equation $X^T X \beta = X^T y$ has unique solution β^* when:

equivalent $\left\{ \begin{array}{l} - X^T X \text{ is invertible/non-singular} \\ - X \text{ has linearly indep. cols.} \\ - X \text{ full-rank} \end{array} \right\} \beta^* = (X^T X)^{-1} X^T y$

• What if $X^T X$ is not invertible? Choose solution β^* with min l^2 -norm. \hookrightarrow prevents "overfitting"

⑥ Moore-Penrose Pseudoinverse: [not in exam]

• For any $A \in \mathbb{R}^{m \times n}$, its singular value decomposition is: unique!

$$A = U D V^T \leftarrow V \in \mathbb{R}^{n \times n}, V^T V = I \text{ (orthogonal)}$$

$U \in \mathbb{R}^{m \times m}$
 $U^T U = I \text{ (orthogonal)}$

$D \in \mathbb{R}^{m \times n}$ diagonal
 $D_{ii} \geq 0, D_{ij} = 0 \text{ for } i \neq j$

• For any $D \in \mathbb{R}^{m \times n}$ diagonal, its Moore-Penrose pseudoinverse is:

$$D^\dagger \in \mathbb{R}^{n \times m} \text{ diagonal}$$

$$[D^\dagger]_{ii} = \begin{cases} \frac{1}{D_{ii}}, & D_{ii} \neq 0 \\ 0, & \text{otherwise} \end{cases}, [D^\dagger]_{ij} = 0 \text{ for } i \neq j$$

• For any $A \in \mathbb{R}^{m \times n}$, its Moore-Penrose pseudoinverse is:

$$A^\dagger = V D^\dagger U^T \in \mathbb{R}^{n \times m}$$

• $\beta^* = X^\dagger y$ is the unique min. l^2 -norm solution to normal equation $X^T X \beta = X^T y$.
← "dagger"

- Facts:
- 1) If X is invertible, $X^\dagger = X^{-1}$.
 - 2) If $X^T X$ is invertible, $X^\dagger = (X^T X)^{-1} X^T$. [see unique sol'n case]

• When $X^T X \beta = X^T y$ has ∞ sol's, they are of the form:

$$\beta = X^\dagger y + v, \quad v \in \text{null}(X)$$

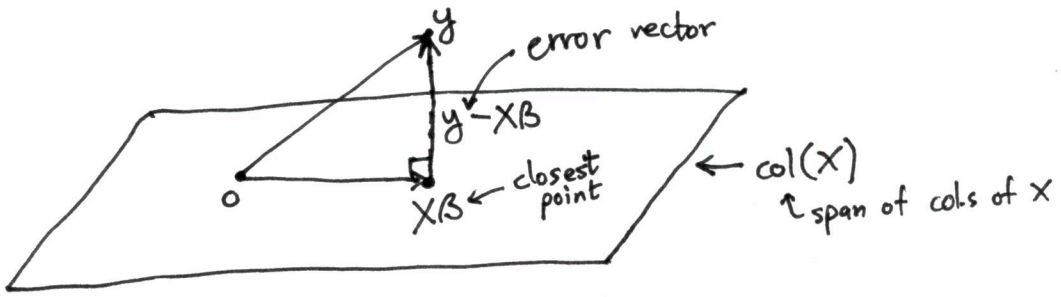
\uparrow nullspace

• For ridge regression, the stationary cond. is $(X^T X + \lambda I) \beta = X^T y$, where $\lambda > 0$.
regularization parameter

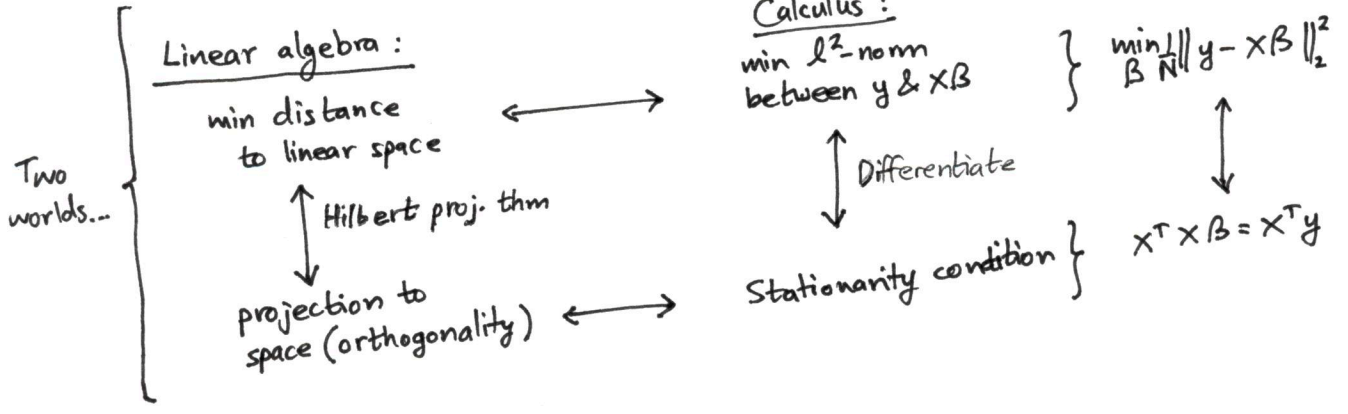
\Rightarrow Solution: $\beta^* = (X^T X + \lambda I)^{-1} X^T y$ ← always exists & unique!
 Moore-Penrose pseudoinv. is limit of ridge sol'n! $\rightarrow \lim_{\lambda \rightarrow 0^+} (X^T X + \lambda I)^{-1} X^T = X^\dagger$

⑦ Geometry of linear regression:

$\min_B \frac{1}{N} \|y - XB\|_2^2 \leftarrow$ Given y , find a vector $XB \in \text{col}(X)$ that is closest to y !



• error is \perp to every col. of X !
 $\rightarrow X^T(\underbrace{y - XB}_{\text{error vector}}) = \mathbf{0} \Leftrightarrow X^T y = X^T X B \leftarrow$ "normal" eqn!
 \uparrow zero vector



CLUSTERING:

Nov 1, 2024
Anuran Makur

• k-means clustering:

Given $\{x_1, \dots, x_N\} \subseteq \mathbb{R}^d$, find clusters S_1, \dots, S_k such that:

training data $\min_{S_1, \dots, S_k} \underbrace{\sum_{i=1}^k \sum_{x \in S_i} \|x - \mu_i\|_2^2}_{f(S_1, \dots, S_k, \mu_1, \dots, \mu_k)}$ [Objective Value], where $\mu_i = \frac{1}{|S_i|} \sum_{x \in S_i} x$.
 centroid

- Assignment Step:

Fix μ_i 's. Then, the best assignments to min. $f(S_1, \dots, S_k, \mu_1, \dots, \mu_k)$ are:

For x_j , choose cluster $\arg \min_i \|x_j - \mu_i\|_2^2$.

↑ break ties deterministically

Assignment will reduce value of Objective func. or keep it the same.

- Update Step:

Fix S_i 's. Then, find best μ_i 's:

$$\min_{\mu_1, \dots, \mu_k} \sum_{i=1}^k \sum_{x \in S_i} \|x - \mu_i\|_2^2 = \sum_{i=1}^k \underbrace{\min_{\mu_i} \sum_{x \in S_i} \|x - \mu_i\|_2^2}_{[*]}$$

$$[*] \min_{\mu_i = (\mu_{i1}, \dots, \mu_{id})} \frac{1}{|S_i|} \sum_{x \in S_i} \sum_{j=1}^d (x_j - \mu_{ij})^2$$

↑
constant

$$= \sum_{j=1}^d \min_{\mu_{ij}} \underbrace{\frac{1}{|S_i|} \sum_{x \in S_i} (x_j - \mu_{ij})^2}_{E[\cdot]}$$

So, it suffices to solve the following problem. For rand. var. $X \in \mathbb{R}$,

what is $\min_{c \in \mathbb{R}} E[(X-c)^2]$?

Prop: $\min_{c \in \mathbb{R}} E[(X-c)^2] = \text{var}(X)$, $\arg \min_{c \in \mathbb{R}} E[(X-c)^2] = E[X]$

Pf: $E[(X-c)^2] = E[(X - E[X] + E[X] - c)^2] = \underbrace{E[(X - E[X])^2]}_{\text{var}(X)} + \underbrace{2E[(X - E[X])(E[X] - c)]}_{=0} + \underbrace{(E[X] - c)^2}_{\geq 0}$

$\geq \text{var}(X)$ with eq. iff $c = E[X]$. □

Hence, Update will reduce value of Objective func. or keep it the same.

- Iteration: (Assignment + Update)

- An iteration produces new clustering \Leftrightarrow Objective func. strictly decreases.
- As there are $\leq k^N$ cluster assignments, and each iter. produces new clusterings, k-means converges in $O(k^N)$ steps. [Stop when 2 iters. have same Obj. Value]

Summary:

- k-means converges in finite time
- $O(k^N)$ iteration complexity in worst-case, but fast in practice
- Only converges to local optima
 - ↳ multiple init.s
 - ↳ kmeans++

EVALUATION METRICS:

Nov 8, 2024
Anuran Makur

• F1-score / Dice-Sørensen coefficient:

$$\frac{1}{F1} \triangleq \frac{1}{2} \left(\frac{1}{\text{Precision}} + \frac{1}{\text{Recall}} \right)$$

$$\Rightarrow 0 \leq F1 \leq 1$$

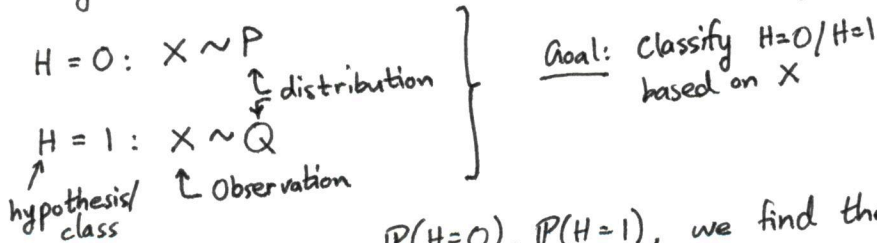
$$F1 = 0 \text{ iff Precision} = 0 \text{ or Recall} = 0$$

$$F1 = 1 \text{ iff Precision} = \text{Recall} = 1$$

Ex: $F1 \approx 0.947...$ for example in slides!

• ROC Curve:

- Theory of Neyman-Pearson hypothesis testing [not to be confused with Fisher/significance testing]



- When we do not know $P(H=0), P(H=1)$, we find the best tradeoff between detection and false-alarm prob. (1-specificity/Type I error)

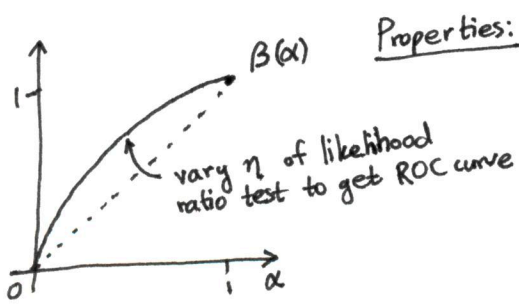
$$\forall \alpha \in [0, 1], B(\alpha) \triangleq \max_{\hat{H}(\cdot)} P(\hat{H}(X) = 1 | H = 1)$$

\uparrow ROC (receiver operating characteristic / Neyman-Pearson func.)
 \uparrow false-alarm prob.
 \uparrow max. over all decision rules/classifiers based on X

AUROC: Area under ROC

$$\text{AUROC} = \int_0^1 B(\alpha) d\alpha \in [0, 1]$$

$$= \frac{\text{Gini coeff.} + 1}{2}$$



Properties:

- 1) $B(0) = 0, B(1) = 1$
- 2) B is non-decreasing and concave
- 3) $B(\alpha) \geq \alpha$
- 4) Neyman-Pearson Lemma: (informal)

For each α , $B(\alpha)$ achieved by likelihood ratio test:

$$\frac{P(x)}{Q(x)} \begin{matrix} \hat{H}(x) = 0 \\ \geq \eta \\ \hat{H}(x) = 1 \end{matrix} \leftarrow \eta \text{ depends on } \alpha$$

- Note: In ML, we do not have (P, Q)
 \Rightarrow Build ROC curve for classifier, eg, logistic regression, by varying "threshold" on "soft" output.

NAIVE BAYES SUPPLEMENT: (NB)

Nov 11, 2024
Anuran Makur

• Multinomial Distribution:

Fix a PMF (parameters) (p_1, p_2, \dots, p_k) over $\{1, \dots, k\}$.
 $\left\{ \begin{array}{l} p_i \geq 0, \\ \sum_{i=1}^k p_i = 1 \end{array} \right.$

Sample n iid points from (p_1, \dots, p_k) . Let $X_i =$ no. of samples equal to i .

The multinomial dist. is the (joint) PMF of (X_1, \dots, X_k) :

$$\mathbb{P}(X_1 = x_1, X_2 = x_2, \dots, X_k = x_k) = \underbrace{\binom{n}{x_1, x_2, \dots, x_k}}_{\text{multinomial coeff.}} p_1^{x_1} p_2^{x_2} \dots p_k^{x_k} \quad \text{for all } x_1, \dots, x_k \geq 0 \text{ and } \sum_{i=1}^k x_i = n$$
$$\binom{n}{x_1, \dots, x_k} \triangleq \frac{n!}{x_1! x_2! \dots x_k!}$$

- Note: $k=2$ is binomial.

- Multinomial Thm: $(a_1 + a_2 + \dots + a_k)^n = \sum_{\substack{x_1, \dots, x_k \geq 0 \\ x_1 + \dots + x_k = n}} \binom{n}{x_1, \dots, x_k} a_1^{x_1} a_2^{x_2} \dots a_k^{x_k}$

- To estimate p_i , take $\frac{X_i}{n}$.

• Three Types of NB:

In each case, the data looks like $(\underbrace{X_1, \dots, X_k}_{k\text{-dim feature}}, \underbrace{Y}_{\text{label}})$.

[* Naive assumptions: X_1, \dots, X_k are indep. given Y *]

① Gaussian NB: Each X_i is Gaussian given Y (diff. param's over i and $Y=y$)
→ Learn $O(k)$ param's (w/o NB, learn $O(k^2)$ param's)

② Bernoulli NB: Each X_i is Bernoulli given Y (diff. param's over i and $Y=y$)
→ Learn $O(k)$ param's (w/o NB, learn $O(2^k)$ param's)

③ Multinomial NB: (X_1, \dots, X_k) is multinomial given Y (diff param's over $Y=y$)
→ Learn $O(k)$ param's [no explicit cond. indep. assumed]

LAPLACE SMOOTHING:

Nov 11, 2024
Anuran Makur

★ Sunrise Problem: What is the prob. of Sun rising tomorrow?

★ Laplace's answer: "Rule of Succession"

- Assume Sun rises indep'tly with prob. p each day.
- Assume p is unknown: $p \sim \text{Uniform}([0,1]) \leftarrow f(p)=1$ for $p \in [0,1]$ (PDF)
- Known: Sun rose on days $1, \dots, N$. ($N < +\infty$)

Want $\rightarrow P(\text{Sun rises on day } N+1 \mid \text{Sun rose on days } 1, \dots, N)$

$$\stackrel{\substack{\text{def.} \\ \text{of cond.} \\ \text{prob.}}}{=} \frac{P(\text{Sun rises on days } 1, \dots, N+1)}{P(\text{Sun rises on days } 1, \dots, N)} \stackrel{[*]}{=} \frac{(\frac{1}{N+2})}{(\frac{1}{N+1})} = \boxed{\frac{N+1}{N+2}} \dots$$

$$[*] P(\text{Sun rises on days } 1, \dots, N) \stackrel{\substack{\text{Law of total} \\ \text{probability}}}{=} \int_0^1 \underbrace{P(\text{Sun rises on days } 1, \dots, N \mid p)}_{= p^N} \underbrace{f(p)}_{= 1} dp$$

$$= \int_0^1 p^N dp$$

$$= \underline{\underline{\frac{1}{N+1}}}$$

★ Concept:

Naive answer	$\frac{N}{N+0}$ <small>Sun rises</small>	$\frac{0}{N+0}$ <small>Sun does not rise</small>	\rightarrow	$\frac{N}{N+0} = 1$	\nwarrow estimates of prob.
Laplace's answer	$\frac{N+1}{N+1}$ <small>Sun rises</small>	$\frac{0+1}{N+1}$ <small>Sun does not rise</small>	\rightarrow	$\frac{N+1}{N+1+1} = \boxed{\frac{N+1}{N+2}}$	\nwarrow (Laplace smoothing)

★ Lidstone Smoothing:

- Add $\alpha \in [0, \infty)$ to each bin!
- \uparrow not an integer

GRADIENT DESCENT:

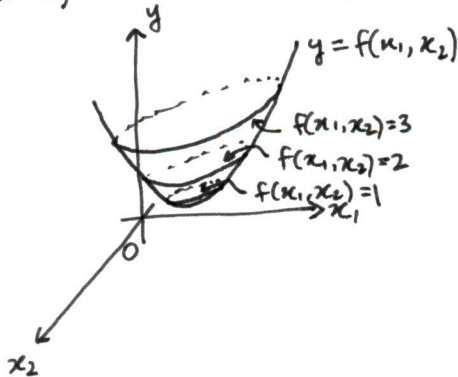
Nov 20, 2024
Anuran Makur

1) Greedy Algorithm:

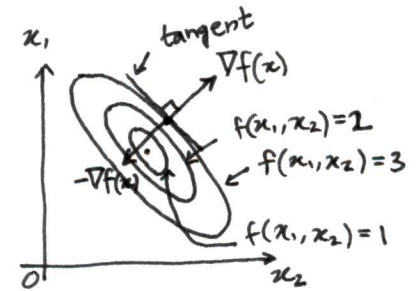
Given a function $f: \mathbb{R}^d \rightarrow \mathbb{R}$, $f(x_1, \dots, x_d)$,

- Problem: $\min_{x \in \mathbb{R}^d} f(x)$.

Ex: ($d=2$)



contours/
level sets



- The gradient of $f: \mathbb{R}^d \rightarrow \mathbb{R}$ is $\nabla f: \mathbb{R}^d \rightarrow \mathbb{R}^d$, $\nabla f(x) \triangleq \left[\frac{\partial f}{\partial x_1}(x) \quad \frac{\partial f}{\partial x_2}(x) \quad \dots \quad \frac{\partial f}{\partial x_d}(x) \right]^T$.

- Facts:
 - ∇f points in direction of steepest ascent.
 - $-\nabla f$ points in direction of steepest descent.

- Recall: (Directional derivative)

For a unit vector $u \in \mathbb{R}^d$, $\|u\|_2 = 1$, $(D_u f)(x) \triangleq \lim_{\delta \rightarrow 0} \frac{f(x + \delta u) - f(x)}{\delta} = \nabla f(x)^T u$

$$-\|\nabla f(x)\|_2 \leq \underbrace{\nabla f(x)^T u}_{\text{directional deriv.}} \leq \|\nabla f(x)\|_2 \underbrace{\|u\|_2}_{=1} = \|\nabla f(x)\|_2 \quad \left[\text{Cauchy-Schwarz inequality} \right]$$

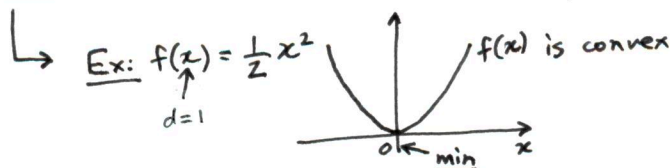
\uparrow eq. iff $u = \frac{-\nabla f(x)}{\|\nabla f(x)\|_2}$ \uparrow eq. iff $u = \frac{\nabla f(x)}{\|\nabla f(x)\|_2}$

Hence, $u \propto -\nabla f(x)$ is direction of steepest descent at x .

- Gradient Descent (GD):
 - 1) Initialize at any arb. $x^0 \in \mathbb{R}^d$.
 - 2) $x^{t+1} = x^t - \eta \nabla f(x^t)$ [greedy!]

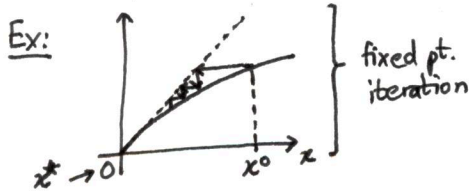
\downarrow $\eta > 0$ step-size (small)

- Facts:
 - GD converges to local min. or saddle point (under cond.s). [Why? See next page.]
 - If f is (strongly) convex, then GD converges to global min.



2) Fixed Point: [not in exam]

• Banach's fixed point thm: Given any $g: \mathbb{R}^d \rightarrow \mathbb{R}^d$ such that g is contractive, i.e., $\forall x, y \in \mathbb{R}^d, \|g(x) - g(y)\|_2 \leq L \|x - y\|_2$ with $L \in (0, 1)$, there is a unique $x^* \in \mathbb{R}^d$ such that:
 Lipschitz continuous



- 1) $g(x^*) = x^*$, [fixed point] iteration initialization
- 2) For any $y \in \mathbb{R}^d, x^{t+1} = g(x^t), x^0 = y,$
 $\lim_{t \rightarrow \infty} x^t = x^*$. [fixed point iteration]
 ↑ conv. exp. fast!

• Consider $g(x) = x - \eta \nabla f(x)$. Then, $x^* = g(x^*) = x^* - \eta \nabla f(x^*) \Leftrightarrow \nabla f(x^*) = 0$.
 fixed point of g zero vector
 stationary point of f

So, fixed pt iteration \Leftrightarrow GD algorithm:

$$x^{t+1} = g(x^t) = x^t - \eta \nabla f(x^t)$$

$$\Rightarrow \lim_{t \rightarrow \infty} x^t = x^*$$

Hence, GD converges to stat. pt. under cond.s.

3) Gradient Flow: [not in exam]

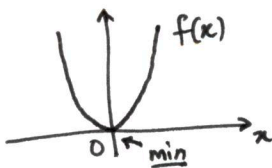
$x^{t+1} = x^t - \eta \nabla f(x^t)$ ← GD iteration
 discrete interval 1

• Change discrete time interval to η : $x^{t+\eta} = x^t - \eta \nabla f(x^t)$

ordinary differential equation $\Rightarrow \frac{dx}{dt} \approx \frac{x^{t+\eta} - x^t}{\eta} = -\nabla f(x^t)$
 approx. of derivative!

• ODE: $\frac{dx}{dt} = -\nabla f(x)$ [Gradient flow]

Ex: $f(x) = \frac{1}{2}x^2$ ($d=1$)



$x(0) = 1$ [Initialize] linear ODE
 $\frac{dx}{dt} = -f'(x) = -x \Rightarrow x(t) = e^{-t}$
 exp. fast convergence!
 $\lim_{t \rightarrow \infty} x(t) = 0$ ← This is the argmin!

• Fact: (informal) eg: Lyapunov stability
 If ODE is "stable", then its solution $x(t)$ converges to $\lim_{t \rightarrow \infty} x(t) = x^*$, where x^* is an equilibrium point (i.e., $\nabla f(x^*) = 0$). (Indeed, $x(t) = x^*$ is a steady-state solution to the ODE.)

RELATION BETWEEN SIGMOID & TANH:

Nov 22, 2024
Anuran Makur

- Sigmoid/Logistic function: $\sigma: \mathbb{R} \rightarrow [0, 1]$,
$$\sigma(x) \triangleq \frac{1}{1+e^{-x}}$$

- Hyperbolic Tangent function: $\tanh: \mathbb{R} \rightarrow [-1, 1]$,

$$\tanh(x) \triangleq \frac{e^x - e^{-x}}{e^x + e^{-x}} = \frac{1 - e^{-2x}}{1 + e^{-2x}} = \frac{1}{1 + e^{-2x}} - \left(1 - \frac{1}{1 + e^{-2x}}\right) = \frac{2}{1 + e^{-2x}} - 1$$

$$\Rightarrow \boxed{\tanh(x) = 2\sigma(2x) - 1} \leftarrow \text{rescaled sigmoid}$$

- Recall: $\tan(x) = \frac{\sin(x)}{\cos(x)}$

$$= \frac{1}{i} \left(\frac{e^{ix} - e^{-ix}}{e^{ix} + e^{-ix}} \right)$$

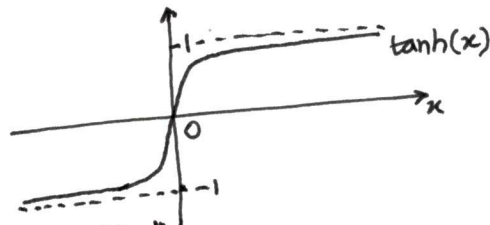
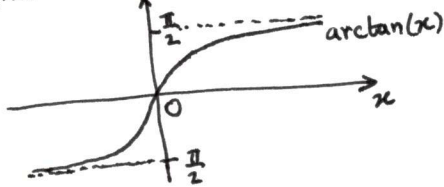
$$\Rightarrow \boxed{\tanh(x) = -i \tan(ix)} \leftarrow \text{relation to } \tan(x)$$

$i = \sqrt{-1}$ (imaginary unit)
 $e^{ix} = \cos(x) + i \sin(x)$ [Euler's formula]

$e^{-ix} = \cos(x) - i \sin(x)$ \leftarrow conjugate

$$\begin{cases} \cos(x) = \frac{e^{ix} + e^{-ix}}{2} \\ \sin(x) = \frac{e^{ix} - e^{-ix}}{2i} \end{cases}$$

- Recall: [not in exam]



Relationship? See " Gudermannian function".
 \hookrightarrow relates hyperbolic & circular angles