

پردازش گفتار

مروری بر آمار و احتمال

هادی ویسی

h.veisi@ut.ac.ir

دانشگاه تهران - دانشکده علوم و فنون نوین



فهرست

- معرفی
- احتمال
- متغیرهای تصادفی
- میانگین و واریانس
- قانون اعداد بزرگ
- توابع توزیع
- نظریه تخمین
 - کمینه میانگین مربعات خطا (MMSE)
 - تخمین بیشینه شباهت (MLE)
 - تخمین بیز (Bayesian)



معرفی

- یک کلمه خاص (مانند "گفتار") را ۱۰ بار مختلف ضبط کنید
 - در هیچ دو حالتی فایل های ضبط شده دیجیتالی آنها دقیقاً یکسان نیست!
 - فرآیند تولید گفتار دارای ذات تصادفی است
- نقش پیرنگ تصادف و عدم قطعیت در پردازش زبان گفتاری
 - فرموله کردن مسائل پردازش گفتار در یک چارچوب احتمالاتی
 - رایج ترین روش ها و الگوریتم های پردازش گفتار آماری هستند



احتمال ...

○ بیان میزان اطمینان از خروجی وقایعی (مشاهده‌هایی) که قطعی نیستند

○ تعاریف

- فضای نمونه (Sample Space): مجموعه‌ای از تمام خروجی‌های ممکن $S =$
- رویداد (Event): زیرمجموعه‌ای از فضای نمونه $A =$
- احتمال (Probability) یک رویداد: فراوانی نسبی رخداد آن رویداد با فرض تکرار این فرایند به تعداد دفعات زیاد تحت شرایط مشابه $P(A) =$

$$P(A) = \frac{N_A}{N_S}$$

تعداد مشاهده‌هایی که خروجی آن‌ها متعلق به رویداد A است

تعداد کل تمام مشاهده‌ها

$$0 \leq P(A) \leq 1 \text{ for all } A$$



احتمال ...

○ افراز (Partition)

• اگر تعداد n رویداد A_1, A_2, \dots, A_n داشته باشیم که $\bigcup_{i=1}^n A_i = S$

$$P(A_1 \cup A_2 \cup \dots \cup A_n) = \sum_{i=1}^n P(A_i) = 1$$

• آنگاه داریم

○ احتمال توأم (Joint Probability)

• برای دو رویداد A و B که به طور همزمان اتفاق می افتند

$$P(AB) = \frac{N_{AB}}{N_S}$$

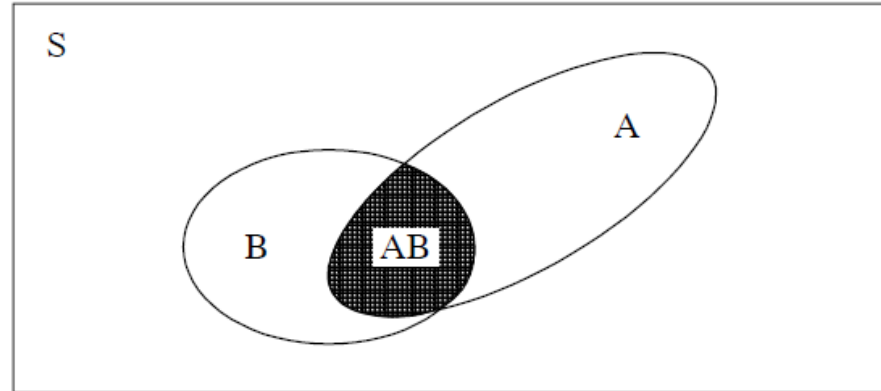


احتمال ...

○ احتمال شرطی (Conditional Probability)

- رخ دادن رویداد A با دانستن اینکه رویداد دیگری مانند B رخ داده است

$$P(A|B) = \frac{P(AB)}{P(B)} = \frac{N_{AB}/N_S}{N_B/N_S}$$



$$\hat{W} = \arg \max_w P(W | X) = \arg \max_w \frac{P(W)P(X | W)}{P(X)}$$

کلمه بیان شده

سیگنال گفتاری

- در بازشناسی گفتار



احتمال ...

○ قاعده زنجیری (Chain Rule)

$$P(A|B) = \frac{P(AB)}{P(B)} = \frac{N_{AB}/N_S}{N_B/N_S} \quad \longrightarrow \quad P(AB) = P(A|B)P(B) = P(B|A)P(A)$$

- حالت ساده
- می‌تواند احتمال توأم چندین رویداد را بر حسب ضرب چند احتمال شرطی مشخص کند
- استفاده برای تجزیه یک مسئله احتمالاتی توأم پیچیده به زنجیره‌ای از احتمال‌های شرطی

$$P(A_1 A_2 \cdots A_n) = P(A_n | A_1 \cdots A_{n-1}) \cdots P(A_2 | A_1) P(A_1) \quad \bullet \text{ حالت کلی}$$

○ مستقل (Independent) بودن

- رخ دادن یک رویداد هیچ ارتباط و تأثیری بر رخ دادن دیگر ندارد.
- احتمال شرطی: برابر با احتمال غیر شرطی است. $P(A|B) = P(A)$
- احتمال توأم: برابر با حاصل ضرب دو احتمال است. $P(AB) = P(A) P(B)$

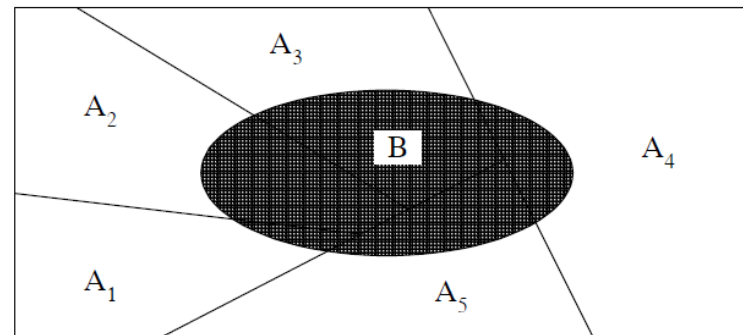


احتمال ...

افراز (Partition) رویداد B

- اگر تعداد n رویداد A_1, A_2, \dots, A_n یک افراز از S باشد و B یک رویداد در S باشد
- آنگاه رویدادهای BA_1, BA_2, \dots, BA_n یک افراز از B را شکل می‌دهد

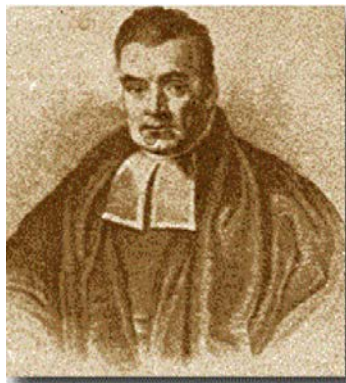
$$B = A_1B \cup A_2B \cup \dots \cup A_nB$$



$$P(B) = \sum_{k=1}^n P(A_k B)$$

- چون رویدادهای BA_1, BA_2, \dots, BA_n مجزا هستند
- احتمال رویداد B از حاصل جمع احتمال‌های توأم محاسبه می‌شود

احتمال حاشیه‌ای (Marginal Probability) رویداد B



Bayes, Thomas (1763) An essay towards solving a problem in the doctrine of chances. *Philosophical Transactions of the Royal Society of London*, 53:370-418

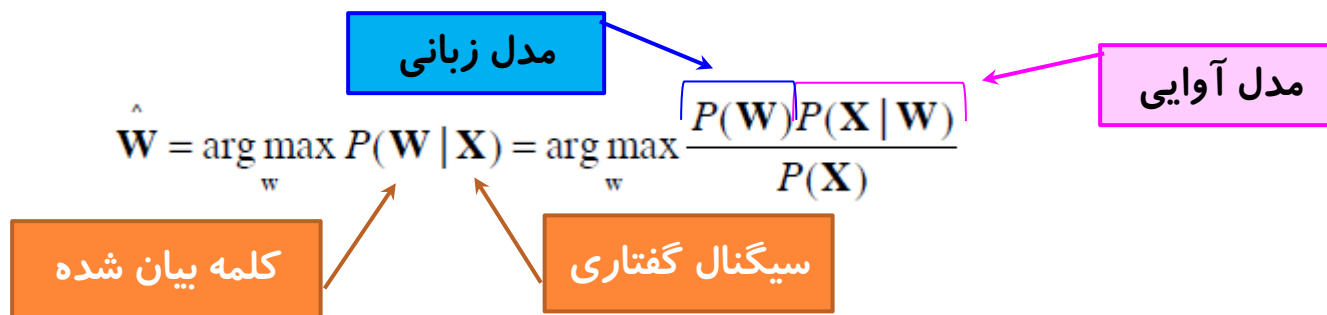
○ قانون بیز (Bayes' Rule)

$$P(B) = \sum_{k=1}^n P(A_k)P(B | A_k)$$

- با توجه به قاعده زنجیری

$$P(A_i | B) = \frac{P(A_i B)}{P(B)} = \frac{P(B | A_i)P(A_i)}{\sum_{k=1}^n P(B | A_k)P(A_k)}$$

- این قانون مبنای بازشناسی الگو (مانند بازشناسی گفتار) است





متغیرهای تصادفی ...

○ متغیر تصادفی (Random Variable)

- متغیر X که بیانگر یک کمیت عددی در یک فضای نمونه است
- عناصر یک فضای نمونه را می‌توان شماره‌گذاری کرد و با آن شماره‌ها به آن‌ها ارجاع کرد.
- تابعی که هر خروجی ممکن s در فضای نمونه S را به یک عدد حقیقی $X(s)$ نگاشت می‌کند.
- یک رویداد به صورت مجموعه‌ای از $\{s\}$ نشان داده می‌شود که $\{s \mid X(s)=x\}$

○ مثال: پرتاب سکه

- فضای نمونه $S = \{\text{شیر، خط}\}$
 - متغیر تصادفی X
- $$X(s) = \begin{cases} 1 & \text{if } s = \text{شیر} \\ 0 & \text{if } s = \text{خط} \end{cases}$$

○ احتمال اینکه $X=x$ باشد

$$P(X = x) = P(s \mid X(s) = x)$$

- در مثال سکه $P(X = 1) = P(s \mid X(s) = 1) = P(s = \text{شیر}) = 0.5$



متغیرهای تصادفی ...

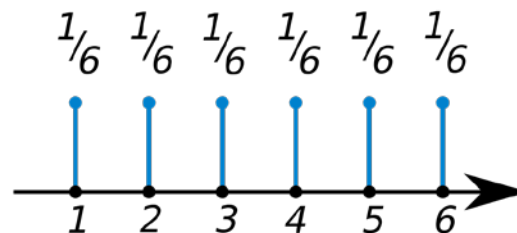
○ متغیر تصادفی گسسته (Discrete) ...

- فقط تعداد متناهی n از مقادیر مختلف را می‌گیرد (دارای توزیع گسسته)
 - مثال: پرتاب تاس (فقط ۶ حالت دارد)

• تابع احتمال (Probability Function)

$$p_X(x) = P(X = x)$$

- یا تابع جرم احتمال (Probability Mass Function)
- برای هر عدد حقیقی x ، بیانگر میزان احتمال متغیر تصادفی گسسته است



- مثال: پرتاب تاس
- متغیرهای تصادفی: 1 تا 6

- حاصل جمع جرم احتمال در تمام مقادیر متغیر تصادفی برابر با یک است

$$\sum_{k=1}^n p(x_k) = \sum_{k=1}^n P(X = x_k) = 1$$



متغیرهای تصادفی ...

○ متغیر تصادفی گسسته (Discrete)

• احتمال حاشیه‌ای $P_X(x_i) = P(X = x_i) = \sum_{k=1}^m P(X = x_i, Y = y_k) = \sum_{k=1}^m P(X = x_i | Y = y_k) P(Y = y_k)$

• قانون زنجیری $P(X_1 = x_1, \dots, X_n = x_n) = P(X_n = x_n | X_1 = x_1, \dots, X_{n-1} = x_{n-1}) \cdots P(X_2 = x_2 | X_1 = x_1) P(X_1 = x_1)$

• قانون بیز $P(X = x_i | Y = y) = \frac{P(X = x_i, Y = y)}{P(Y = y)} = \frac{P(Y = y | X = x_i) P(X = x_i)}{\sum_{k=1}^n P(Y = y | X = x_k) P(X = x_k)}$

• اگر متغیرهای تصادفی X و Y به لحاظ آماری از هم مستقل باشند

$$P(X = x_i, Y = y_j) = P(X = x_i) P(Y = y_j) = p_X(x_i) p_Y(y_j) \quad \forall \text{ all } i \text{ and } j$$



متغیرهای تصادفی ...

○ متغیر تصادفی پیوسته (Continuous) ...

- دارای مقادیر پیوسته (و در نتیجه توزیع پیوسته) است
- مثال: قد افراد یک کشور، مقدار دامنه سیگنال گفتار

- اگر تابع غیرمنفی f وجود داشته باشد که روی مقادیر حقیقی تعریف شده و برای بازه A

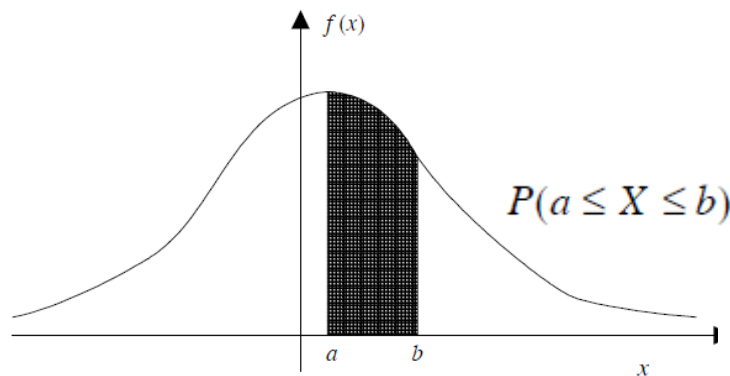
$$P(X \in A) = \int_A f_X(x) dx$$

$$f(x) \geq 0 \text{ for } -\infty \leq x \leq \infty$$

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

- f_x : تابع توزیع احتمال (pdf: Probability Density Function)

- احتمال در یک بازه معنی دارد
- احتمال در یک نقطه برابر با صفر است





متغیرهای تصادفی ...

○ متغیر تصادفی پیوسته (Continuous) ...

• شرایط تابع pdf

$$f(x) \geq 0 \text{ for } -\infty \leq x \leq \infty$$

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

$$f_X(x) = \int_{-\infty}^{\infty} f_{X,Y}(x, y) dy = \int_{-\infty}^{\infty} f_{X|Y}(x | y) f_Y(y) dy$$

• احتمال حاشیه‌ای

$$f_{X_1, \dots, X_n}(x_1, \dots, x_n) = f_{X_n | X_1, \dots, X_{n-1}}(x_n | x_1, \dots, x_{n-1}) \cdots f_{X_2 | X_1}(x_2 | x_1) f_{X_1}(x_1)$$

• قانون زنجیری

$$f_{X|Y}(x | y) = \frac{f_{X,Y}(x, y)}{f_Y(y)} = \frac{f_{Y|X}(y | x) f_X(x)}{\int_{-\infty}^{\infty} f_{Y|X}(y | x) f_X(x) dx}$$

• قانون بیز



متغیرهای تصادفی ...

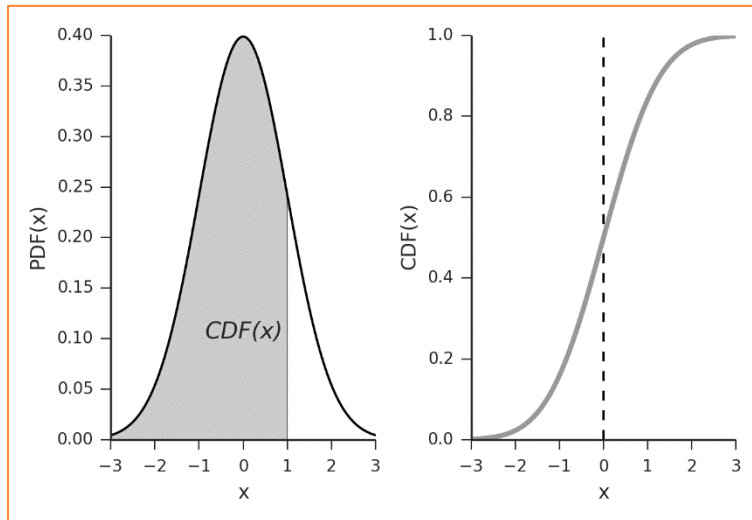
تابع توزیع (Distribution Function)

- یا تابع توزیع تجمعی (CDF: Cumulative Distribution Function)
- بیانگر [جمع] احتمال‌های مقادیر کوچک‌تر از x

$$F(x) = P(X \leq x) \text{ for } -\infty \leq x \leq \infty$$

مقدار حقیقی

- برای متغیر تصادفی گسسته یا پیوسته



- برای حالت پیوسته داریم

$$F(x) = \int_{-\infty}^x f_X(x) dx$$

$$f_X(x) = \frac{dF(x)}{dx}$$



میانگین و واریانس ...

○ امید ریاضی (Expectation) یا میانگین (Mean)

$$E(X) = \sum_x xf(x) \quad \bullet \text{ برای متغیر تصادفی گسسته } X$$

$$E(X) = \int_{-\infty}^{\infty} xf(x)dx \quad \bullet \text{ برای متغیر تصادفی پیوسته } X$$

$$E[g(X)] = \int_{-\infty}^{\infty} g(x)f(x)dx \quad \bullet \text{ برای تابعی از متغیر } X$$

○ مرکز جرم توزیع احتمال

○ امید ریاضی یک عملگر خطی است؛ دارای ویژگی‌های جمع‌پذیری و همگنی

○ حتی در صورت مستقل نبودن X_i ها

$$E(a_1X_1 + \dots + a_nX_n + b) = a_1E(X_1) + \dots + a_nE(X_n) + b$$

مقدار ثابت



میانگین و واریانس ...

○ امید ریاضی (Expectation) یا میانگین (Mean)

• مثال ۱: انداختن یک تاس

○ متغیر تصادفی با شش مقدار 1, 2, ..., 6 با احتمال برابر 1/6 برای هر کدام

$$E(X) = \sum_x xf(x) = \frac{1}{6} \times 1 + \frac{1}{6} \times 2 + \frac{1}{6} \times 3 + \frac{1}{6} \times 4 + \frac{1}{6} \times 5 + \frac{1}{6} \times 6 = 3.5$$

• مثال ۲: محاسبه معدل درسی یک دانشجو

○ متغیر تصادفی (X)؟

○ نمره درس

○ احتمال (تابع توزیع)؟

○ متناسب با تعداد واحد درسها (تعداد واحد درس تقسیم بر کل واحد)

تعداد واحد	نمره	درس
۲	۱۸	آشناسی
۴	۱۲	برنامه نویسی
۳	۱۵	ریاضیات
۱	۲۰	روش تحقیق

$$E(x) = \frac{2}{10} \times 18 + \frac{4}{10} \times 12 + \frac{3}{10} \times 15 + \frac{1}{10} \times 20 = 14.9$$



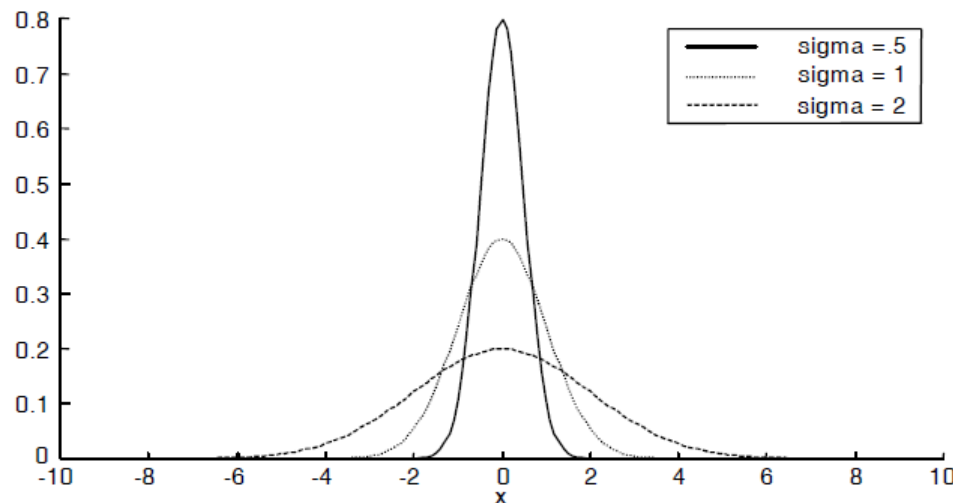
میانگین و واریانس ...

○ واریانس (Variance) ...

$$\text{Var}(X) = \sigma^2 = E[(X - \mu)^2]$$

• میانگین متغیر X : $\mu = E(X)$

- مجذور غیرمنفی واریانس = انحراف معیار (Standard Deviation) σ
- بیانگر میزان پراکندگی یا انتشار توزیع در اطراف میانگین
- مقدار کوچک واریانس = توزیع فشرده احتمال در اطراف میانگین
- مقدار بزرگ واریانس = توزیع احتمال در اطراف میانگین دارای پراکندگی بیشتری است





میانگین و واریانس ...

○ واریانس (Variance)

• ممان (گشتاور) k ام X = امید ریاضی $E(X^k)$

○ برای هر متغیر تصادفی X و هر عدد صحیح مثبت k

• برای واریانس داریم

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

• واریانس اختلاف بین ممان دوم و مربع ممان اول است

• ویژگی‌های واریانس

○ جمع‌پذیری: در صورت مستقل بودن متغیر تصادفی X و Y

$$Var(X + Y) = Var(X) + Var(Y)$$

○ همگنی: برقرار نیست

○ واریانس مقدار ثابت صفر است

$$Var(aX) = a^2 Var(X)$$

$$Var(a_1 X_1 + \dots + a_n X_n + b) = a_1^2 Var(X_1) + \dots + a_n^2 Var(X_n)$$



میانگین و واریانس ...

○ امید ریاضی شرطی (Conditional Expectation)

$$E_{Y|X}(Y | X = x) = \sum_y y f_{Y|X}(y | x)$$

- برای متغیرهای گسسته X و Y
- امید ریاضی شرطی Y : تابعی از X

$$E_{Y|X}(Y | X = x) = \int_{-\infty}^{\infty} y f_{Y|X}(y | x) dy$$

- برای متغیرهای پیوسته X و Y

- خود $E(Y | X)$ یک متغیر تصادفی است
- تابعی از متغیر تصادفی X است

$$E_X [E_{Y|X}(Y | X)] = E_{X,Y}(Y)$$

- فرض کنید X و Y یک توزیع توأم پیوسته دارند و $g(X, Y)$ تابعی از X و Y است

$$E_{Y|X} [g(X, Y) | X = x] = \int_{-\infty}^{\infty} g(x, y) f_{Y|X}(y | x) dy$$

$$E_X \{E_{Y|X} [g(X, Y) | X]\} = E_{X,Y} [g(X, Y)]$$



میانگین و واریانس

○ میانه (Median)

- احتمال کل را به دو قسمت مساوی تقسیم می‌کند

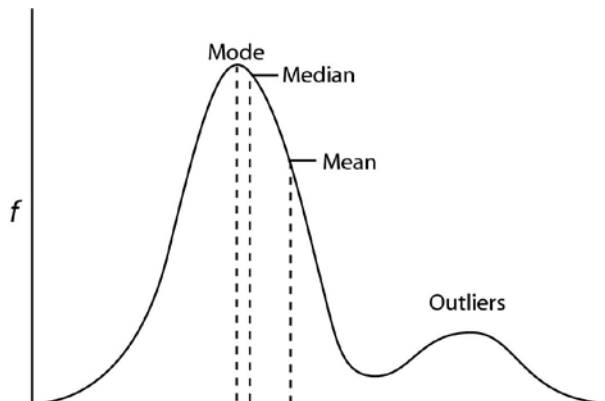
$$P(X \leq m) \geq 1/2 \text{ and } P(X \geq m) \geq 1/2$$

- میانه توزیع متغیر X نقطه‌ای مانند m است

○ احتمال سمت چپ m و احتمال سمت راست m دقیقاً 0.5 است

○ نما (Mode)

- جایی که تابع توزیع دارای بیشترین مقدار خود است
- یک توزیع می‌تواند بیش از یک میانه داشته باشد





قانون اعداد بزرگ ...

○ قانون اعداد بزرگ (Law of Large Numbers)

- میانگین نمونه (Sample Mean) و واریانس نمونه (Sample Variance)
- مقدار میانگین و واریانس تعدادی از نمونه‌های حاصل از یک آزمایش آماری است (آنچه ما در عمل محاسبه می‌کنیم)

○ فرض کنید یک توزیع با میانگین μ و واریانس σ^2 داریم

- متغیرهای تصادفی X_1, X_2, \dots, X_n از این توزیع تولید می‌شوند
- متغیرهای تصادفی i.i.d: - Independent Identically Distributed
 - مستقل و با توزیع یکسان
 - هر یک دارای μ و واریانس σ^2
- میانگین حسابی n نمونه
 - همان میانگین نمونه

$$\bar{X}_n = \frac{1}{n}(X_1 + \dots + X_n)$$

خودش متغیری
تصادفی است



قانون اعداد بزرگ

○ داریم

$$E(\bar{X}_n) = \mu$$

• میانگین "میانگین نمونه"

○ میانگین "میانگین نمونه" برابر با میانگین توزیع است

$$\text{Var}(\bar{X}_n) = \frac{\sigma^2}{n}$$

• واریانس "میانگین نمونه"

○ واریانس "میانگین نمونه" برابر با $1/n$ واریانس توزیع است
○ توزیع "میانگین نمونه" به نسبت توزیع اصلی در اطراف میانگین متمرکزتر است

○ قانون اعداد بزرگ

• بیان می کند "میانگین نمونه" به میانگین توزیع نزدیک می شود

○ وقتی اندازه نمونه (n) بزرگ باشد

$$\lim_{n \rightarrow \infty} P(|\bar{X}_n - \mu| < \varepsilon) = 1 \text{ for any given number } \varepsilon > 0$$



کواریانس و همبستگی ...

○ کواریانس متغیرهای تصادفی X و Y

$$Cov(X, Y) = E[(X - \mu_X)(Y - \mu_Y)] = Cov(Y, X)$$



$$Cov(X, Y) = E(XY) - E(X)E(Y)$$

○ ضریب همبستگی (Correlation Coefficient)

$$\rho_{XY} = \frac{Cov(X, Y)}{\sigma_X \sigma_Y}$$

• مقدار در بازه -1 و 1 $-1 \leq \rho(X, Y) \leq 1$

• بیانگر همبستگی خطی (linear dependency) بین دو متغیر X و Y

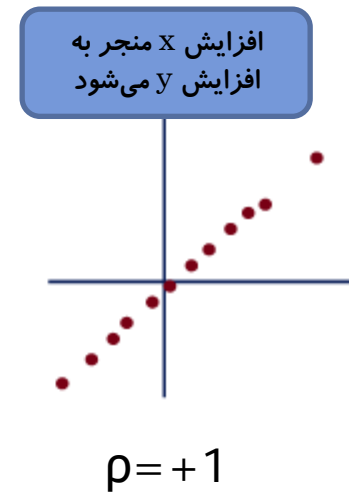
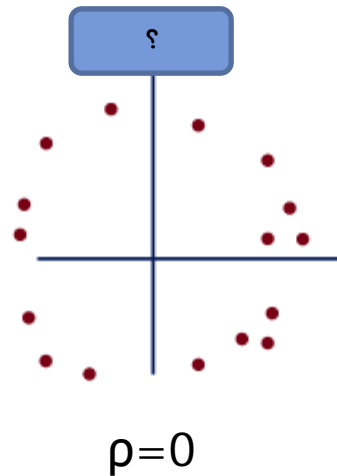
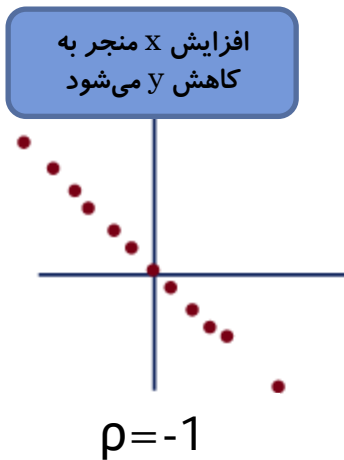
• دو متغیر تصادفی X و Y متعامد (Orthogonal) هستند اگر $E(XY) = 0$



کوواریانس و همبستگی ...

ضریب همبستگی (Correlation Coefficient)

- اگر $\rho > 0$ باشد، X و Y همبستگی مثبت دارند
- اگر $\rho < 0$ باشد، همبستگی منفی دارند
- اگر $\rho = 0$ باشد، همبسته (Correlated) نیستند





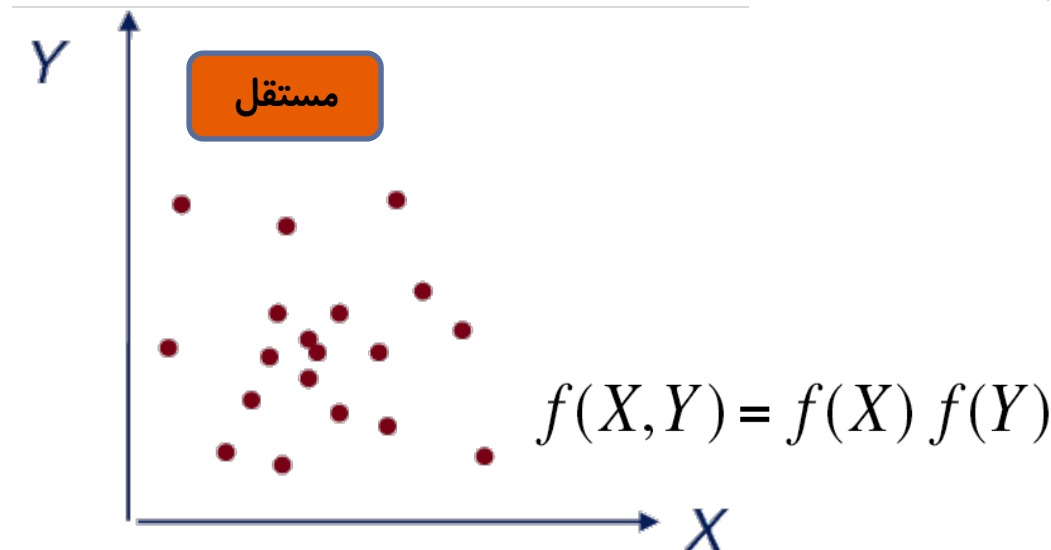
کوواریانس و همبستگی ...

○ متغیرهای تصادفی مستقل (Independent)

$$\text{Cov}(X, Y) = \rho_{XY} = 0$$

• اگر داشته باشیم

- آیا ناهمبسته بودن (uncorrelated)، استقلال (independence) را نتیجه می‌دهد؟
 - مستقل بودن، ناهمبسته بودن را نیز نتیجه می‌دهد، اما برعکس آن درست نیست (به غیر از توزیع نرمال)





کوواریانس و همبستگی

○ چند قضیه

• اگر رابطه متغیرهای X و Y به صورت $Y=aX+b$ باشد، آنگاه

○ اگر $a>0$ باشد، $\rho_{XY} = +1$

○ اگر $a<0$ باشد، $\rho_{XY} = -1$

• برای هر دو متغیر X و Y داریم

$$Var(X + Y) = Var(X) + Var(Y) + 2Cov(X, Y)$$

• اگر n متغیر تصادفی X_1, X_2, \dots, X_n داشته باشیم، آنگاه

$$Var\left(\sum_{i=1}^n X_i\right) = \sum_{i=1}^n Var(X_i) + 2\sum_{i=1}^n \sum_{j=1}^{i-1} Cov(X_i, X_j)$$



بردارهای تصادفی ...

○ بردار تصادفی

- وقتی یک متغیر تصادفی یک بردار باشد و نه یک عدد

$$\mathbf{X} = (X_1, \dots, X_n)$$



$$\mathbf{X} = \begin{bmatrix} X_1 \\ \vdots \\ X_n \end{bmatrix}$$

بردارى با n مؤلفه

• بردار میانگین

- یک بردار n بعدی که مؤلفه‌های آن امید ریاضی‌های تک تک مؤلفه‌های \mathbf{X} است

$$E(\mathbf{X}) = \begin{bmatrix} E(X_1) \\ \vdots \\ E(X_n) \end{bmatrix}$$

• ماتریس کوواریانس

- مؤلفه‌های قطر اصلی ماتریس کوواریانس = واریانس‌های هر کدام از X_i ها

$$\text{Cov}(X_i, X_j) = \text{Cov}(X_j, X_i) \quad \text{○ کوواریانس متقارن است}$$

$$\text{Cov}(\mathbf{X}) = E \left[[\mathbf{X} - E(\mathbf{X})][\mathbf{X} - E(\mathbf{X})]^t \right] = \begin{bmatrix} \text{Cov}(X_1, X_1) & \dots & \text{Cov}(X_1, X_n) \\ \vdots & & \vdots \\ \text{Cov}(X_n, X_1) & \dots & \text{Cov}(X_n, X_n) \end{bmatrix}$$



پردازش‌های تصادفی

○ رابطه خطی پردازش‌های تصادفی

- بردار n بعدی X
- بردار m بعدی Y

$$Y = AX + B$$

یک ماتریس $m \times n$

یک بردار m بعدی

- میانگین و کواریانس Y بر حسب میانگین و کواریانس X

$$E(Y) = AE(X) + B$$

○ کاربرد در تبدیل ویژگی‌ها (کاهش بعد)

$$Cov(Y) = ACov(X)A^t$$

○ کاربرد در تطبیق (adapt) پارامترهای مدل‌های آوایی و زبانی

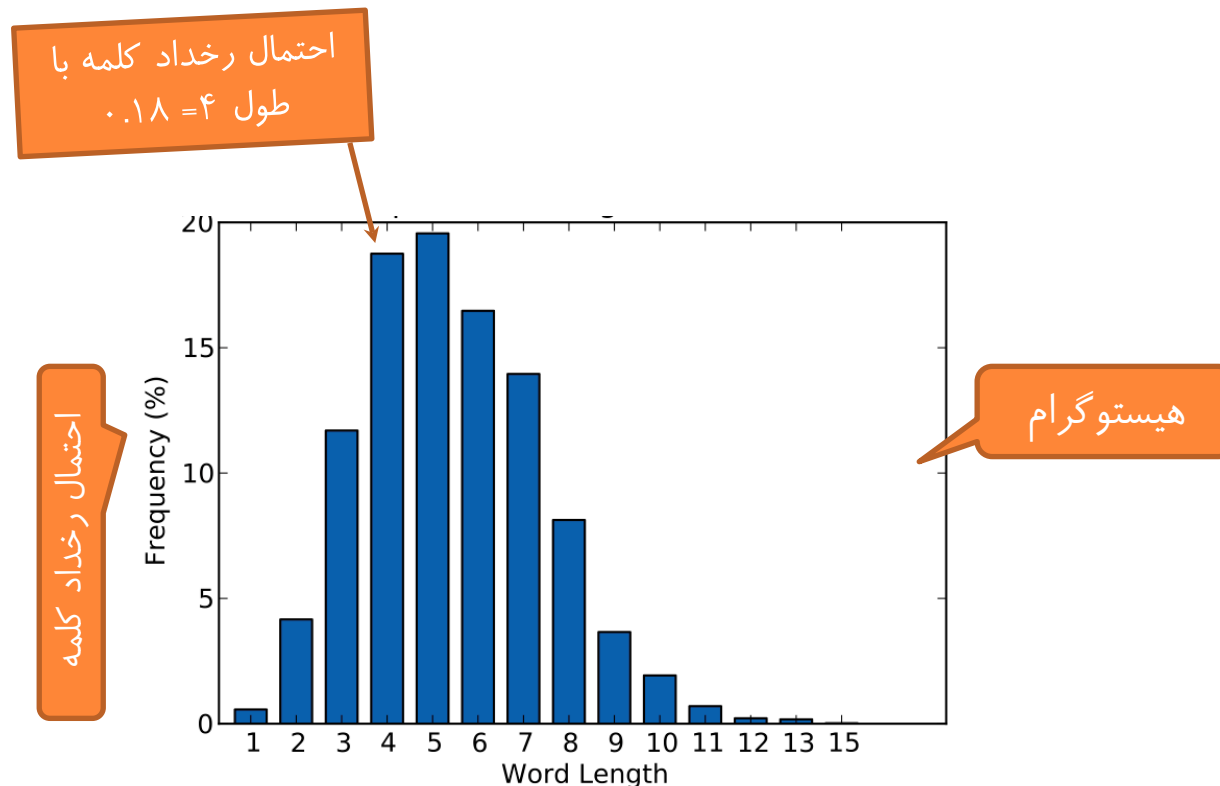
ترانزفاده



توابع توزیع ...

○ مثال: پردازش متن (طول کلمات) ...

- یک لغتنامه داریم
- می‌خواهیم تعداد کلمات با طول ۱ حرف، با طول ۲ حرف، ...، طول ۲۰ حرف را بشماریم



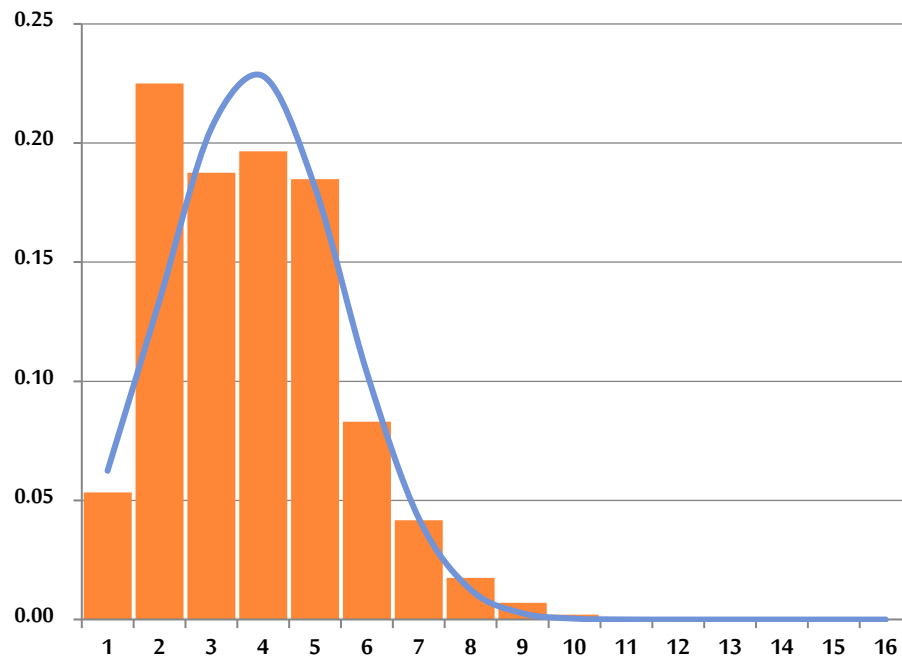


توابع توزیع ...

○ مثال: پردازش متن (طول کلمات) - برای فارسی

• روی پیکره کوچک

○ متوسط طول کلمات: ۳.۸

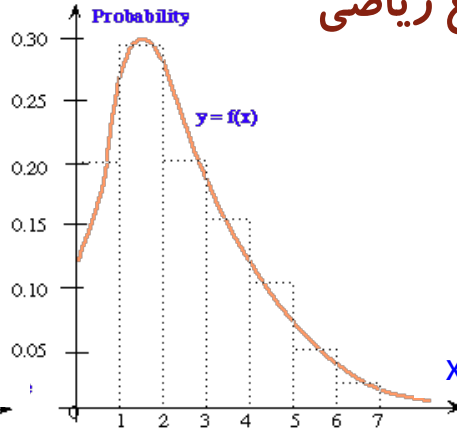
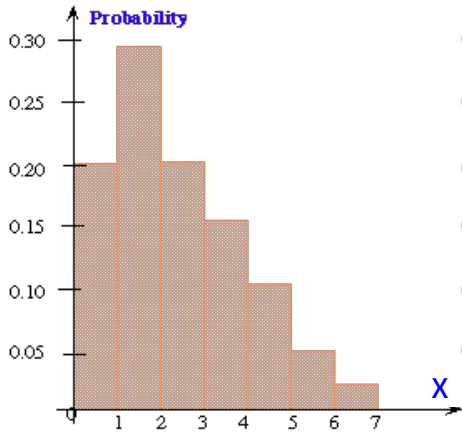




توابع توزیع ...

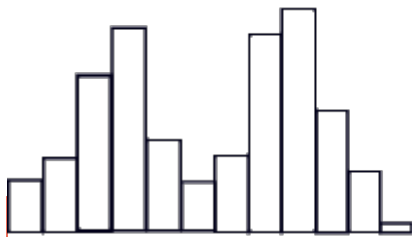
○ از هیستوگرام به تابع توزیع

- نمایش شکل توزیع احتمال‌ها با یک تابع ریاضی

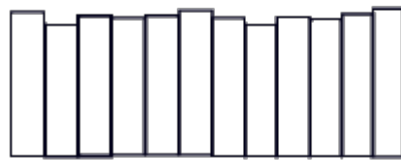


- می‌تواند شکل‌های مختلفی داشته باشد

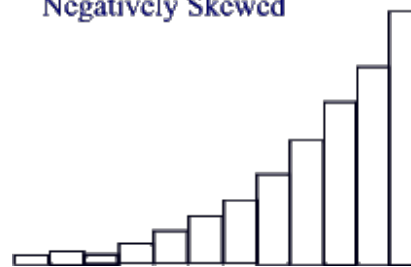
Bi-Modal Distribution



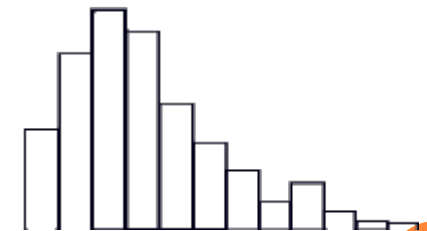
Unitary Distribution



Negatively Skewed



- Positively Skewed





توابع توزیع ...

○ توزیع یکنواخت (Uniform Distribution)

- تابع احتمال یا تابع توزیع احتمال، یک تابع ثابت است
- احتمال وقوع همه مقادیر یکسان است
- مثال: احتمال انتخاب هر نقطه در یک بازه عددی مشخص

$$P(X = x_i) = \frac{1}{n} \quad 1 \leq i \leq n$$

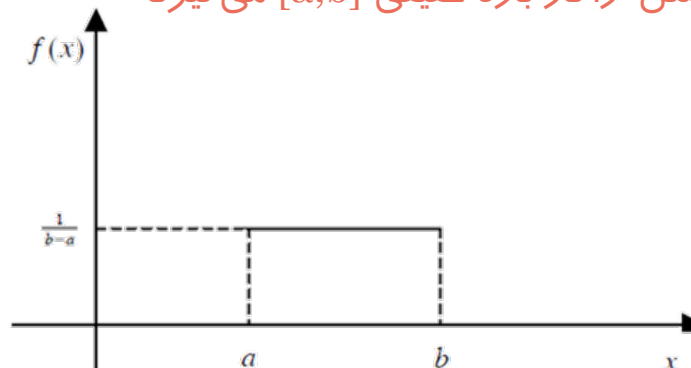
- برای متغیر گسسته

○ X فقط مقادیر ممکن از ۱ تا n را می‌گیرد

$$f(x) = \frac{1}{b-a} \quad a \leq x \leq b$$

- برای متغیر پیوسته

○ X فقط مقادیر ممکن را در بازه حقیقی $[a, b]$ می‌گیرد





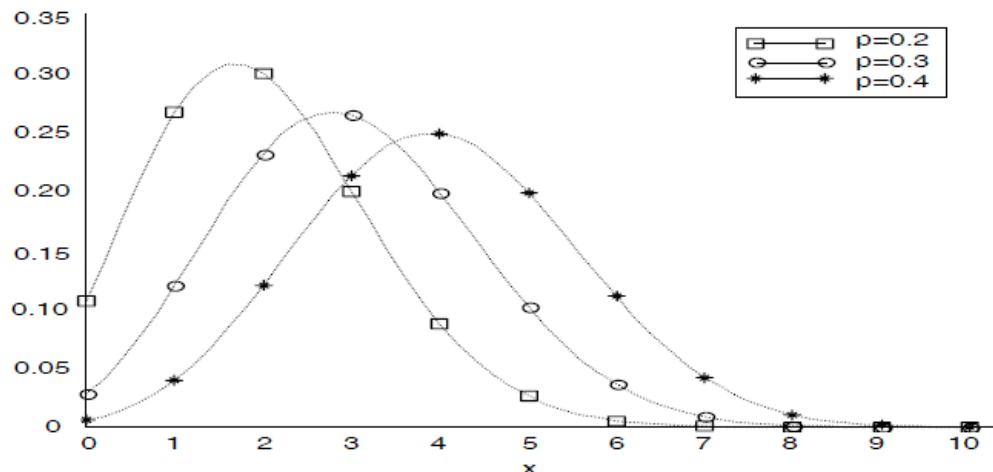
توابع توزیع ...

توزیع دو جمله‌ای (Binomial Distribution) ...

- برای توصیف رویدادهایی با تصمیم‌گیری دودویی
 - آزمایش برنولی: شکست (۰) یا پیروزی (۱)
- مثال: انداختن یک سکه (دو حالت شیر یا خط)
 - احتمال آمدن شیر p و احتمال خط $1-p$
- سکه را n بار بیندازیم (n بار تکرار مستقل آزمایش برنولی) و تعداد شیرهای مشاهده شده (مجموع پیروزی‌ها) را با X نشان دهیم، متغیر تصادفی X دارای تابع احتمال دو جمله‌ای است

$$P(X = x) = f(x | n, p) = \binom{n}{x} p^x (1-p)^{n-x}$$

تعداد x بار از n بار شیر بیاید





توابع توزیع ...

توزیع دوجمله‌ای (Binomial Distribution)

$$P(X = x) = f(x | n, p) = \binom{n}{x} p^x (1-p)^{n-x}$$

$$E(X) = np$$

$$Var(X) = np(1-p)$$

- میانگین و واریانس توزیع دوجمله‌ای

- مثال: اگر یک تیرانداز با احتمال 0.7 تیری را به هدف بزند و X تعداد تیرهای به هدف خورده در 5 شلیک باشد

$$p(k) = \binom{n}{k} (0.7)^k (0.3)^{5-k}$$

- تابع توزیع

$$P(X = 3) = \binom{5}{3} (0.7)^3 (0.3)^2 = 0.3087$$

- احتمال اینکه دقیقاً 3 تیر به هدف بزند

$$P(X \leq 2) = p(0) + p(1) + p(2) = 0.16308$$

- احتمال اینکه حداکثر 2 تیر به هدف بزند

- مثال: امتحان چهار گزینه‌ای

○ پیروزی = گزینه درست (0.25)، شکست = سایر گزینه‌ها (0.75)



توابع توزیع ...

○ توزیع هندسی (Geometric Distributions) ...

- بیانگر تعداد آزمایش‌های (زمان) برنولی تا رسیدن به اولین پیروزی در توزیع دوجمله‌ای
- مثال: انداختن سکه – تعداد آزمایش‌هایی که تکرار می‌شود تا یک بار خط بیاید
- مثال: تعداد بارهای انداختن تاس تا آمدن شش در منج!
- مثال: چند مرتبه آزمایش کنیم تا یک رمز عبور ۸ کاراکتری یک کامپیوتر را حدس بزنیم
- X = متغیر تصادفی زمان (تعداد انداختن‌ها) قبل از آمدن اولین خط
 - احتمال آمدن شیر = p و احتمال آمدن خط = $1-p$

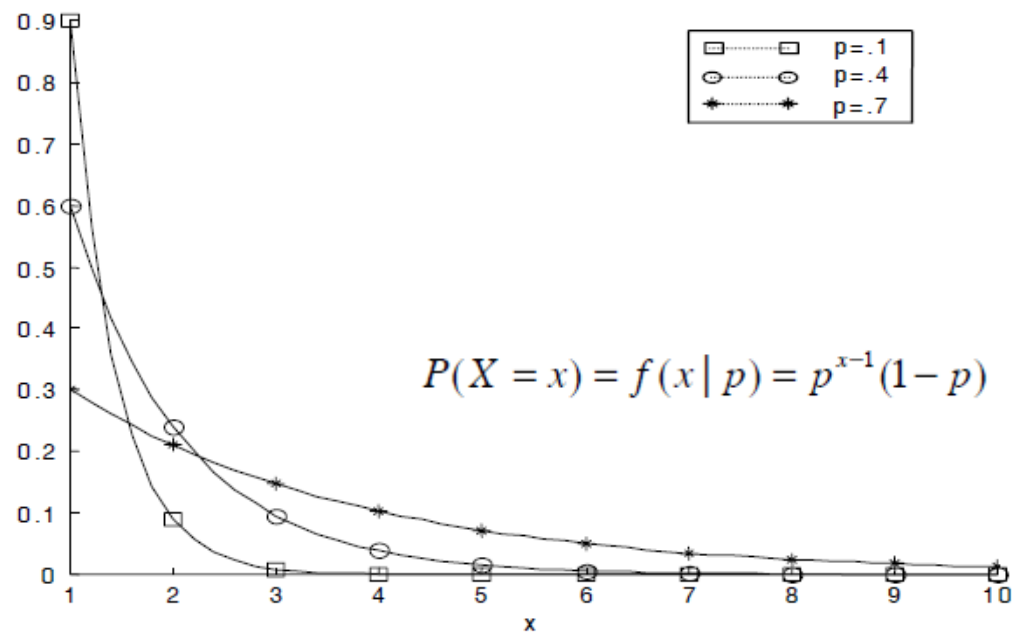
$$P(X = x) = f(x | p) = p^{x-1}(1-p) \quad x = 1, 2, \dots \quad \text{and } 0 < p < 1$$

$$\text{Var}(X) = \frac{1}{(1-p)^2} \quad E(X) = \frac{1}{1-p} \quad \bullet \text{ میانگین و واریانس}$$



توابع توزیع ...

توزیع هندسی (Geometric Distributions)



• مثال: توزیع طول حالت (State Duration) در مدل مخفی مارکوف (HMM)

○ بیانگر مدت زمانی است که در یک حالت خاص باقی می‌مانیم



توابع توزیع ...

○ توزیع چندجمله‌ای (Multinomial Distribution) ...

- حالت کلی‌تر توزیع دوجمله‌ای با k حالت (به جای دو حالت)
- مثال: کیسه‌ای حاوی توپ‌هایی با k رنگ مختلف داریم
 - نسبت توپ‌های رنگ i برابر با p_i است
 - فرض کنید n توپ به طور تصادفی از کیسه انتخاب شده‌اند
 - فرض کنید X_i بیانگر تعداد توپ‌های انتخاب شده با رنگ i باشد
 - آنگاه بردار تصادفی $\mathbf{X}=(X_1, \dots, X_k)$ دارای توزیع چندجمله‌ای با پارامترهای n و $\mathbf{p}=(p_1, \dots, p_k)$ است

$$P(\mathbf{X} = \mathbf{x}) = f(\mathbf{x} | n, \mathbf{p}) = \begin{cases} \frac{n!}{x_1! \dots x_k!} p_1^{x_1} \dots p_k^{x_k} & \text{where } x_i \geq 0 \ \forall i = 1, \dots, k \\ & \text{and } x_1 + \dots + x_k = n \\ 0 & \text{otherwise} \end{cases}$$

- توزیع‌های چندجمله‌ای معمولاً با آزمون مربعات-کای به کار می‌رود
 - از پرکاربردترین روش‌های آزمون فرضیه‌های نیکویی برازش (Goodness-of-Fit)



توابع توزیع ...

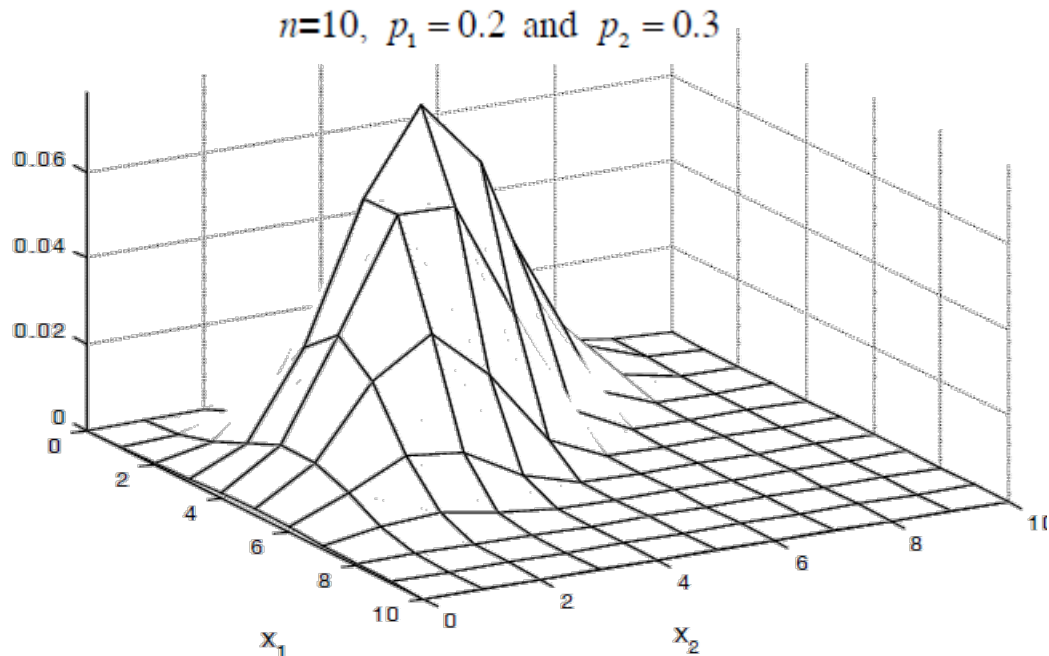
توزیع چندجمله‌ای (Multinomial Distribution)

- میانگین، واریانس و کواریانس

$$E(X_i) = np_i$$

$$Var(X_i) = np_i(1 - p_i) \quad \forall i = 1, \dots, k$$

$$Cov(X_i, X_j) = -np_i p_j$$

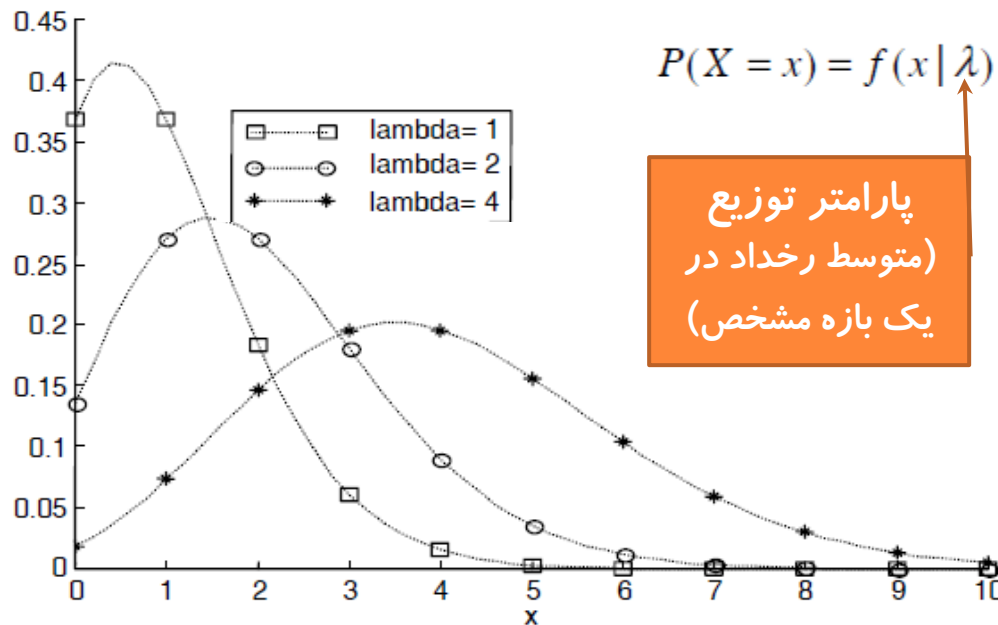




توابع توزیع ...

توزیع پواسون (Poisson Distribution)

- از توزیع‌های گسسته متداول
- $X =$ تعداد کل رخداد‌های یک پدیده در طول یک مدت زمان یا یک منطقه مکانی ثابت
- تعداد تماس‌های تلفنی دریافت شده توسط یک مرکز مخابراتی در یک مدت زمان ثابت



$$P(X = x) = f(x | \lambda) = \begin{cases} \frac{e^{-\lambda} \lambda^x}{x!} & \text{for } x=0,1,2,\dots \\ 0 & \text{otherwise} \end{cases}$$

پارامتر توزیع
(متوسط رخداد در
یک بازه مشخص)

- میانگین و واریانس

$$E(X) = Var(X) = \lambda$$

- کاربرد در پردازش گفتار

• برای مدل‌سازی دیرش (duration) یک واج



توابع توزیع ...

توزیع گاما (Gamma Distribution) ...

• متغیر تصادفی پیوسته X

پارامتر شکل (shape)

$$f(x | \alpha, \beta) = \begin{cases} \frac{\beta^\alpha}{\Gamma(\alpha)} x^{\alpha-1} e^{-\beta x} & x > 0 \\ 0 & x \leq 0 \end{cases}$$

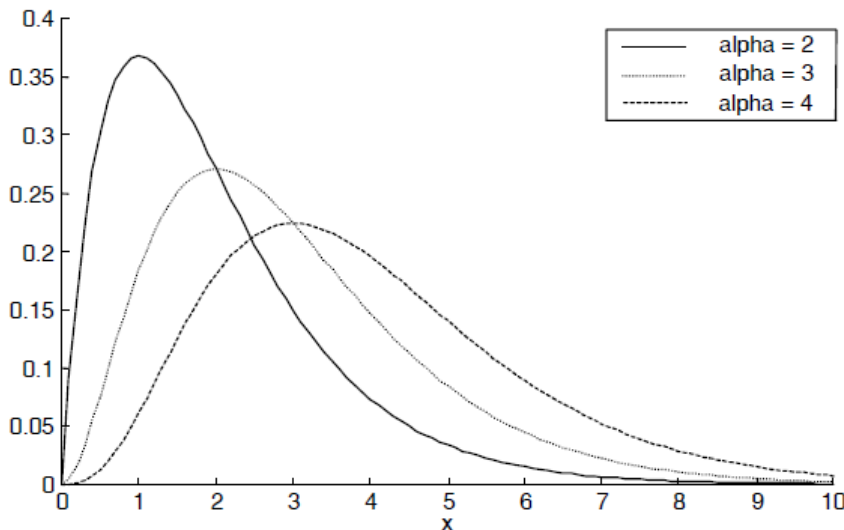
پارامتر مقیاس (scale)

تابع گاما

$$\Gamma(\alpha) = \int_0^{\infty} x^{\alpha-1} e^{-x} dx$$

$$\Gamma(n) = \begin{cases} (n-1)! & n = 2, 3, \dots \\ 1 & n = 1 \end{cases}$$

تابع فاکتوریل برای n حقیقی



• میانگین و واریانس

$$E(X) = \frac{\alpha}{\beta} \text{ and } Var(X) = \frac{\alpha}{\beta^2}$$



توابع توزیع ...

○ توزیع گاما (Gamma Distribution)

- اگر متغیرهای تصادفی X_1, \dots, X_k مستقل باشند،
و هر متغیر تصادفی X_i یک توزیع گاما با پارامترهای α_i و β داشته باشد،
آنگاه مجموع $X_1 + \dots + X_k$ نیز یک توزیع گاما با پارامترهای $\alpha_1 + \dots + \alpha_k$ و β دارد

• در پردازش گفتار

- سیگنال گفتار نویزی دارای توزیع گاما است (دامنه طیف)
- فرکانس‌های بالای سیگنال گفتار در حوزه DCT و بخش‌های حقیقی و موهومی طیف



توابع توزیع ...

توزیع نمایی (Exponential Distribution)

$$f(x|\beta) = \begin{cases} \beta e^{-\beta x} & x > 0 \\ 0 & x \leq 0 \end{cases}$$

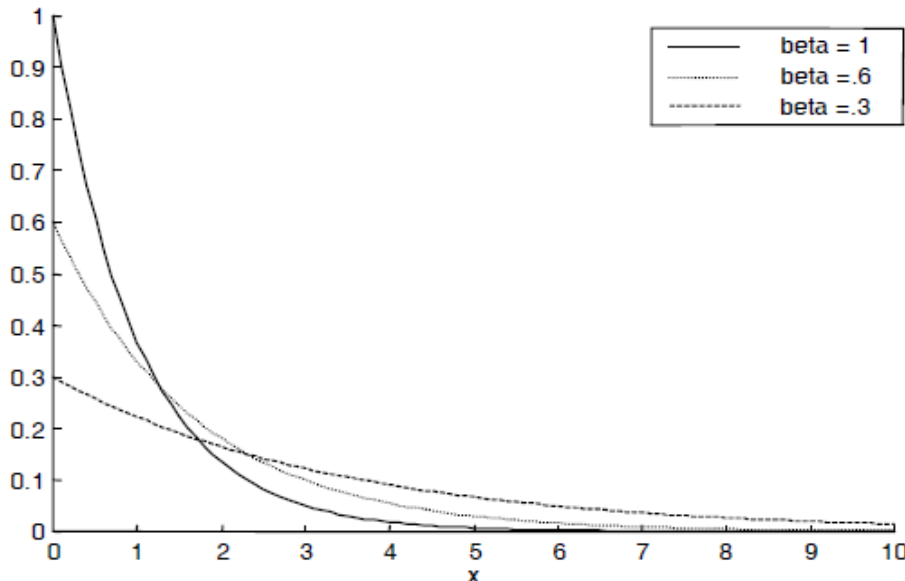
• حالت خاصی از توزیع گاما که $\alpha=1$

• تخمین زدن مدت زمان لازم برای رخداد یک پیشامد خاص

○ زمان ورود مشتریها

○ طول عمر یک وسیله

○ زمان بین دو رویداد در فرایند پواسن



• میانگین و واریانس

$$E(X) = \frac{1}{\beta} \quad \text{and} \quad \text{Var}(X) = \frac{1}{\beta^2}$$



توابع توزیع ...

توزیع لاپلاس (Laplace Distribution)

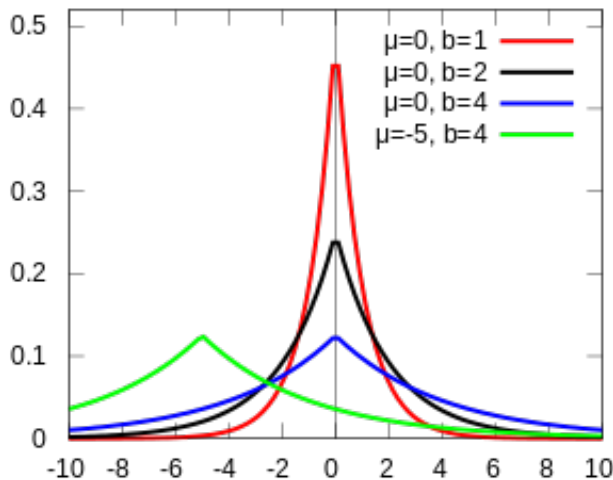
- برای متغیر پیوسته

- ترکیب دو توزیع گاما

پارامتر مقیاس (scale)

پارامتر مکان (location)

$$f(x|\mu, b) = \frac{1}{2b} \exp\left(-\frac{|x - \mu|}{b}\right) = \frac{1}{2b} \begin{cases} \exp\left(-\frac{\mu - x}{b}\right) & \text{if } x < \mu \\ \exp\left(-\frac{x - \mu}{b}\right) & \text{if } x \geq \mu \end{cases} = \frac{1}{2b} \exp\left(-\frac{|x - \mu|}{b}\right)$$



- میانگین = میانه = نما = μ

- واریانس = $2b^2$

- کاربرد در پردازش گفتار

○ سیگنال گفتار تمیز دارای توزیع لاپلاس است (دامنه طیف)

○ فرکانس‌های پایین سیگنال گفتار در حوزه DCT و بخش‌های حقیقی و موهومی طیف



توابع توزیع ...

○ توزیع گاوسی (Gaussian Distribution) ...

- یا توزیع نرمال (Normal Distribution)

- مهم‌ترین توزیع احتمال

- متغیرهای تصادفی مطالعه شده در آزمایش‌های مختلف فیزیکی (از جمله سیگنال‌های گفتاری) دارای توزیع‌هایی هستند که تقریباً گاوسی است

- محاسبات آن آسان است (به ویژه در تخمین‌ها)

- قضیه حد مرکزی (Central Limit Theorem)

میانگین

واریانس

$$f(x | \mu, \sigma^2) = N(\mu, \sigma^2) = \frac{1}{\sqrt{2\pi}\sigma} \exp\left[-\frac{(x-\mu)^2}{2\sigma^2}\right]$$

- برای یک متغیر تصادفی پیوسته

- تابع گاوسی حول میانگین تقارنی است

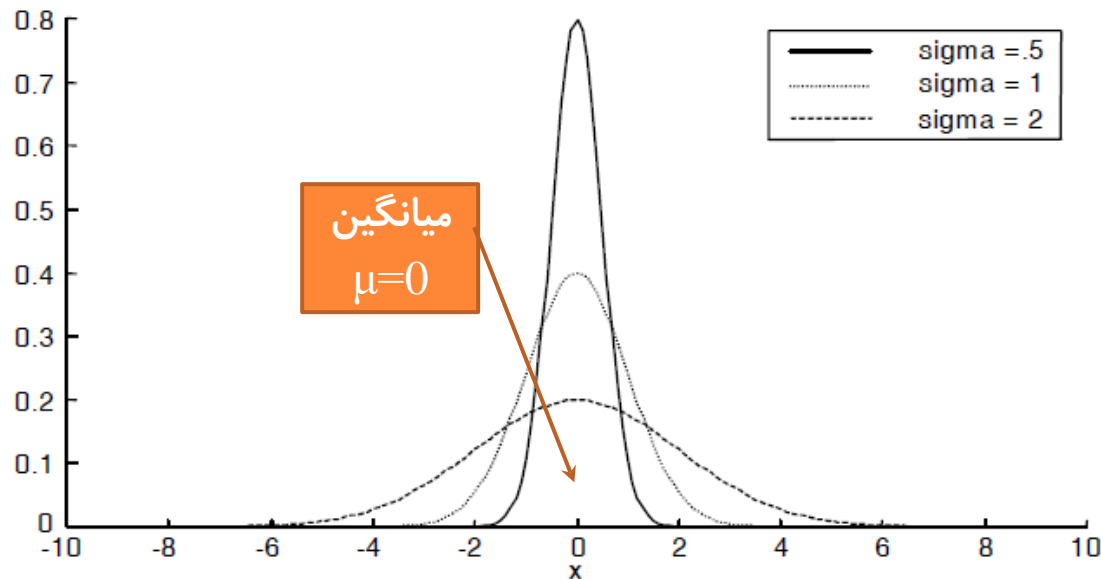
- میانگین، میانه (نقطه تقارن) و نمای (بیشینه مقدار) توزیع یک نقطه است



توابع توزیع ...

توزیع گاوسی (Gaussian Distribution) ...

$$f(x|\mu, \sigma^2) = N(\mu, \sigma^2) = \frac{1}{\sqrt{2\pi}\sigma} \exp\left[-\frac{(x-\mu)^2}{2\sigma^2}\right]$$



کاهش واریانس = تراکم بیشتر حول میانگین



توابع توزیع ...

○ توزیع گاوسی (Gaussian Distribution) ...

- اگر متغیر تصادفی X یک توزیع گاوسی با میانگین μ و واریانس σ^2 باشد
 آنگاه هر تابع خطی $Y=aX+b$ نیز یک توزیع گاوسی دارد
 Y یک توزیع گاوسی با میانگین $a\mu+b$ و واریانس $a^2\sigma^2$ دارد
- حالت کلی: مجموع $X_1 + \dots + X_n$ از متغیرهای تصادفی مستقل X_1, \dots, X_n نیز، که در آن هر متغیر تصادفی X_i یک توزیع گاوسی دارد، یک توزیع گاوسی است
- توزیع گاوسی استاندارد یا توزیع گاوسی $N(0,1) =$
 - میانگین صفر و واریانس یک
 - رفتار توزیع گاوسی را می‌توان فقط با استفاده از توزیع گاوسی استاندارد توضیح داد
 - تبدیل خطی توزیع گاوسی یک توزیع گاوسی است
 - اگر متغیر تصادفی X یک توزیع گاوسی با میانگین μ و واریانس σ^2 باشد، می‌توان نشان داد

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



توابع توزیع ...

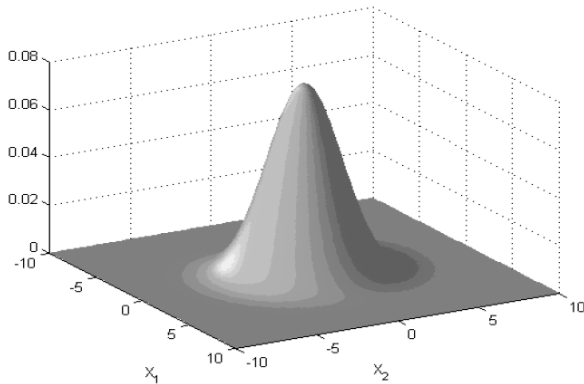
توزیع گاوسی چندمتغیره (Multivariate)

- برای بردار تصادفی n بعدی

n = تعداد ابعاد بردار x (متغیرها)

$$f(\mathbf{X} = \mathbf{x} | \boldsymbol{\mu}, \boldsymbol{\Sigma}) = N(\mathbf{x}; \boldsymbol{\mu}, \boldsymbol{\Sigma}) = \frac{1}{(2\pi)^{n/2} |\boldsymbol{\Sigma}|^{1/2}} \exp\left[-\frac{1}{2} (\mathbf{x} - \boldsymbol{\mu})^t \boldsymbol{\Sigma}^{-1} (\mathbf{x} - \boldsymbol{\mu})\right]$$

$|\cdot|$ = دترمینان
 t = ترانهاده
 -1 = معکوس



برای دو بعد

• میانگین $\boldsymbol{\mu} = E(\mathbf{x})$

• کواریانس $\boldsymbol{\Sigma} = E[(\mathbf{x} - \boldsymbol{\mu})(\mathbf{x} - \boldsymbol{\mu})^t]$

- متقارن و مثبت-معین (positive definite) (دترمینان مثبت)
- عناصر قطر اصلی σ_{ii} = واریانس متغیر متناسب x_i
- عناصر غیرقطر اصلی σ_{ij} = کواریانس متغیرهای x_i و x_j

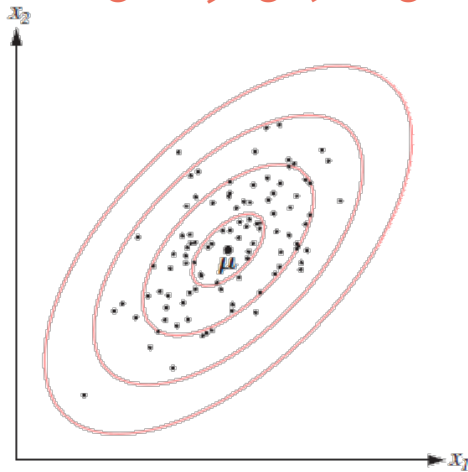
$$\sigma_{ij}^2 = E[(x_i - \mu_i)(x_j - \mu_j)]$$



توابع توزیع ...

○ شکل داده‌های دارای توزیع نرمال

- نمونه‌های داده که از توزیع نرمال پیروی می‌کنند، داخل یک خوشه قرار می‌گیرد
- مرکز خوشه = میانگین، شکل خوشه (بیضی شکل) = تعیین شده توسط واریانس (ماتریس کواریانس)



- بردار ویژه ماتریس کواریانس = محورهای اصلی بیضی
- مقادیر ویژه ماتریس کواریانس = طول محورهای اصلی بیضی

○ فاصله ماهالونوبیس (Mahalanobis distance)

- فاصله بین یک مجموعه داده معین (با پارامترهای میانگین و واریانس) و یک نمونه داده
- در نظر گرفتن وابستگی بین داده‌ها (متفاوت با فاصله اقلیدسی)
- تغییرناپذیر با مقیاس (scale invariance)

$$(\mathbf{x} - \boldsymbol{\mu})^T \boldsymbol{\Sigma}^{-1} (\mathbf{x} - \boldsymbol{\mu})$$



توابع توزیع ...

○ قضیه حد مرکزی (Central Limit Theorem) ...

- n متغیر تصادفی X_1, \dots, X_n که i.i.d هستند (مستقل و با توزیع یکسان) داریم
- توزیع این متغیرها دارای میانگین μ و واریانس σ^2 است

$$Y_n = \frac{n(\bar{X}_n - \mu)}{\sqrt{n\sigma^2}} \sim N(0,1)$$

میانگین نمونه‌ای متغیرها

- با افزایش n به سمت بی‌نهایت، داریم
- متغیر تصادفی Y دارای توزیع گاوسی استاندارد است

- متغیر تصادفی میانگین نمونه‌ای دارای توزیع گاوسی با میانگین μ و واریانس σ^2/n است



توابع توزیع ...

○ قضیه حد مرکزی (Central Limit Theorem)

- توسعه برای حالتی که توزیع‌ها یکسان نیستند (Liapounov 1901)

○ متغیرهای تصادفی X_1, \dots, X_n مستقل هستند و $E(|X_i - \mu_i|^3) < \infty$

○ آنگاه متغیر زیر دارای توزیع گاوسی است $Y_n = \left(\sum_{i=1}^n X_i - \sum_{i=1}^n \mu_i \right) / \left(\sum_{i=1}^n \sigma_i^2 \right)^{1/2}$

○ مجموع متغیرهای تصادفی X_1, \dots, X_n دارای توزیع گاوسی با میانگین $\sum_{i=1}^n \mu_i$ و واریانس $\left(\sum_{i=1}^n \sigma_i^2 \right)^{1/2}$ است

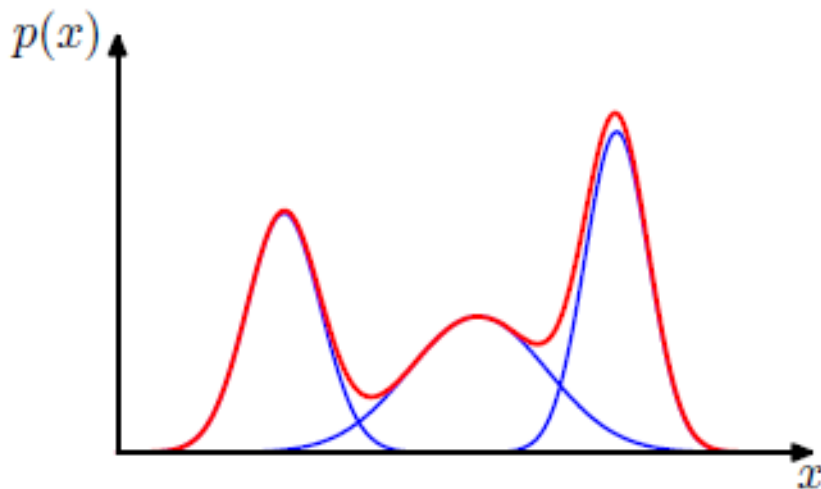
حاصل جمع تعداد زیادی متغیر تصادفی مستقل، صرف نظر از توزیع‌های اصلی هر یک از آن‌ها، با بزرگ شدن تعداد متغیرهای تصادفی، دارای توزیع گاوسی است.



توابع توزیع

○ مدل‌های مخلوط (Mixture Model)

- ترکیب (خطی) چند مدل با همدیگر
- مدل کردن توزیع‌های پیچیده با بیشینه‌های محلی چند گانه
- برای توزیع نرمال (گوسی): مدل مخلوط گاوسی (GMM: Gaussian Mixture Model)
 - از پرکاربردترین روش‌های مدل‌سازی



ضریب مخلوط

$$f(\mathbf{x}) = \sum_{k=1}^K c_k N_k(\mathbf{x}; \boldsymbol{\mu}_k, \boldsymbol{\Sigma}_k)$$

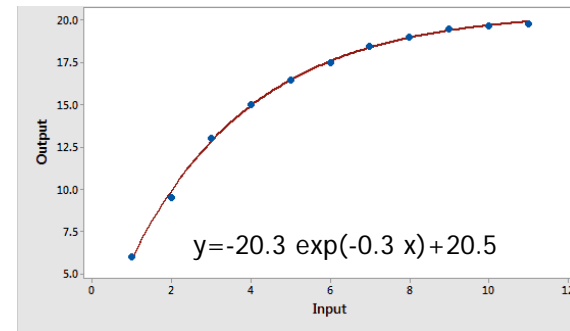
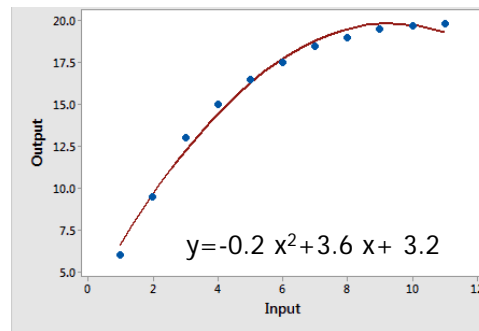
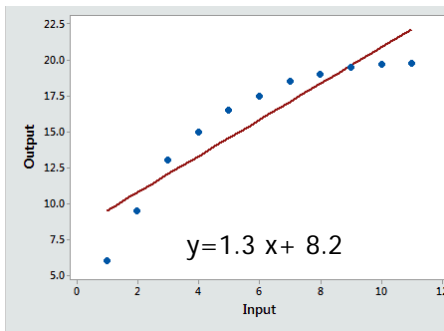
$$c_k \geq 0 \text{ and } \sum_{k=1}^K c_k = 1$$



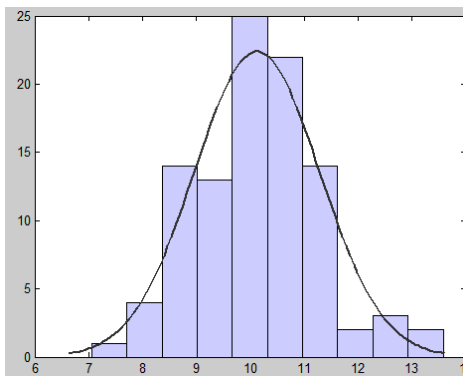
نظریه تخمین ...

مساله

- تعدادی نمونه داده داریم، می‌خواهیم آنها را با یک تابع (خطی/غیرخطی) مدل کنیم
 - مثال: وزن افراد بر حسب قد آنها، سیگنال تمیز بر حسب سیگنال نویزی



- تعدادی نمونه داده داریم، می‌خواهیم تابع توزیع احتمال آنها را بدست آوریم
 - مثال: توزیع گاوسی به طول کلمات در یک زبان



$$\frac{1}{\sqrt{2\pi}\sigma} \exp\left[-\frac{(x-\mu)^2}{2\sigma^2}\right]$$



نظریه تخمین ...

○ نظریه تخمین (Estimation theory)

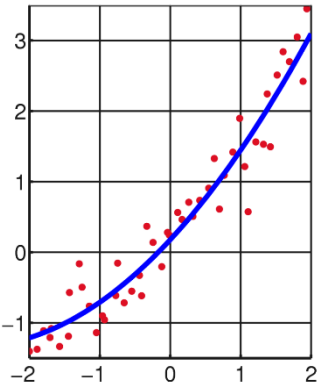
- در مدل‌سازی آماری یک تابع توزیع (مانند گاوسی) برای داده‌ها فرض می‌کنیم
- و باید از روی داده‌های آموزش، پارامترهای آن توزیع (مانند میانگین و واریانس) را تخمین بزنیم
- متغیرهای تصادفی X_1, \dots, X_n را که i.i.d هستند (مستقل و با توزیع یکسان) داریم
- هدف تخمین پارامترهای Φ
- تابع تخمین‌گر $\theta(X_1, \dots, X_n)$

○ روش‌های تخمین

- کمینه میانگین مربعات خطا (MMSE: Minimum Mean Square Error)
- تخمین بیشینه شباهت (MLE: Maximum-Likelihood Estimation)
- تخمین بیز (Bayesian Estimation)



نظریه تخمین ...



○ کمینه میانگین مربعات خطا (MMSE) ...

- کمینه کردن امید ریاضی مربعات خطای بین مقدار واقعی و مقدار تخمین زده شده

$$E(Y - \hat{Y})^2 = E(Y - g(X))^2$$

- فرض کنید هدف ما تخمین مقدار Y با داشتن X باشد، یعنی $\hat{Y} = g(X)$
- که $g(X)$ تابعی بر حسب پارامترهای Φ است، یعنی $g(X, \Phi)$

- و با داشتن پارامترهای Φ ، تابع $g()$ به صورت کامل مشخص می‌شود
- پس: هدف تخمین پارامترهای Φ است

$$\hat{\Phi}_{MMSE} = \arg \min_{\Phi} [E[(Y - g(X, \Phi))^2]]$$

• تخمین LSE: Least Square Error

- در عمل به جای تابع توزیع توأم X و Y ، نمونه‌هایی از x_i و y_i معادل داریم

$$\Phi_{LSE} = \arg \min_{\Phi} \sum_{i=1}^n [y_i - g(x_i, \Phi)]^2$$

- قانون اعداد بزرگ: وقتی تعداد نمونه‌ها به بی‌نهایت میل می‌کند، LSE و MMSE برابر می‌شوند



نظریه تخمین ...

○ کمینه میانگین مربعات خطا (MMSE): برای تابع ثابت ...

• تابع ثابت $\hat{Y} = g(x) = c$

○ پارامتر $c =$

• هدف کمینه کردن خطاست $E(Y - \hat{Y})^2 = E(Y - c)^2$

○ مشتق گرفتن و برابر صفر قرار دادن $c_{MMSE} = E(Y)$

• خطای مجذور میانگین کمینه دقیقاً برابر با واریانس Y است

• تخمین LSE

○ میانگین نمونه‌ای

$$\min \sum_{i=1}^n [y_i - c]^2$$

$$c_{LSE} = \frac{1}{n} \sum_{i=1}^n y_i$$



نظریه تخمین ...

○ کمینه میانگین مربعات خطا (MMSE): برای تابع خطی ...

• تابع خطی $\hat{Y} = g(x) = ax + b$

○ پارامترها: a و b

$$e(a, b) = E(Y - \hat{Y})^2 = E(Y - ax - b)^2$$

$$\frac{\partial e}{\partial a} = 0, \text{ and } \frac{\partial e}{\partial b} = 0$$



$$a = \frac{\text{cov}(X, Y)}{\text{Var}(X)} = \rho_{XY} \frac{\sigma_Y}{\sigma_X}$$

$$b = E(Y) - \rho_{XY} \frac{\sigma_Y}{\sigma_X} E(X)$$

$$\hat{\mathbf{Y}} = \mathbf{X}\mathbf{A} \text{ or } \begin{pmatrix} y_1 \\ y_2 \\ \vdots \\ y_n \end{pmatrix} = \begin{pmatrix} 1 & x_1^1 & \cdots & x_1^d \\ 1 & x_2^1 & \cdots & x_2^d \\ \vdots & \vdots & & \vdots \\ 1 & x_n^1 & \cdots & x_n^d \end{pmatrix} \begin{pmatrix} a_0 \\ a_1 \\ \vdots \\ a_d \end{pmatrix}$$

• برای تخمین LSE

○ فرض: بردار \mathbf{x} دارای d بعد است و n نمونه داریم

$$e(\mathbf{A}) = \|\hat{\mathbf{Y}} - \mathbf{Y}\|^2 = \sum_{i=1}^n (\mathbf{A}^t \mathbf{x}_i - y_i)^2 \quad \longrightarrow \quad \nabla e(\mathbf{A}) = \sum_{i=1}^n 2(\mathbf{A}^t \mathbf{x}_i - y_i) \mathbf{x}_i = 2\mathbf{X}^t (\mathbf{X}\mathbf{A} - \mathbf{Y})$$

شبه معکوس

$$\mathbf{X}^t \mathbf{X} \mathbf{A} = \mathbf{X}^t \mathbf{Y} \quad \longrightarrow \quad \mathbf{A}_{LSE} = (\mathbf{X}^t \mathbf{X})^{-1} \mathbf{X}^t \mathbf{Y}$$

نظریه تخمین ...

کمینه میانگین مربعات خطا (MMSE): برای تابع غیرخطی

• تابع غیرخطی

$$\min_{g(\bullet) \in G_{nl}} E[Y - g(X)]^2$$

$$E_X [E_{Y|X}(Y | X)] = E_{X,Y}(Y)$$

$$E_{X,Y} [Y - g(X)]^2 = E_X \left\{ E_{Y|X} \left[[Y - g(X)]^2 \mid X = x \right] \right\}$$

$$= \int_{-\infty}^{\infty} E_{Y|X} \left[[Y - g(X)]^2 \mid X = x \right] f_X(x) dx$$

$$= \int_{-\infty}^{\infty} E_{Y|X} \left[[Y - g(x)]^2 \mid X = x \right] f_X(x) dx$$

مقدار مثبت، کمینه کردن انتگرال معادل کمینه کردن این مقدار است

$$E(Y - \hat{Y})^2 = E(Y - c)^2$$

$$c_{MMSE} = E(Y)$$

$$\min_{g(x) \in R} E_{Y|X} \left[[Y - g(x)]^2 \mid X = x \right] \Rightarrow \hat{Y} = g_{MMSE}(X) = E_{Y|X}(Y | X) = \int_{-\infty}^{\infty} y f_Y(y | X = x) dy$$

تخمین طیف گفتار تمیز

$$\hat{S}(\omega_k) = E[S(\omega_k) | \mathbf{Y}] = \frac{P_{ys}(\omega_k)}{P_{yy}(\omega_k)} = \frac{P_{ss}(\omega_k)}{P_{ss}(\omega_k) + P_{dd}(\omega_k)} Y(\omega_k)$$

• در بهسازی گفتار

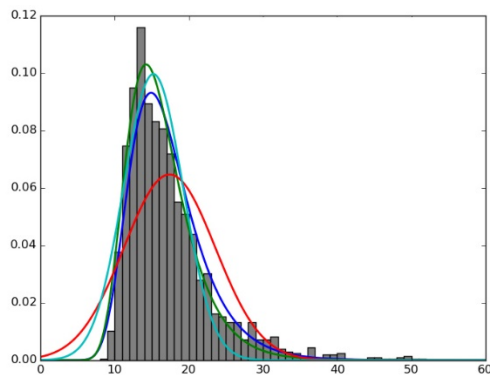
سیگنال گفتار نویزی

توان طیف گفتار تمیز

توان طیف نویز



نظریه تخمین ...



○ تخمین بیشینه شباهت (MLE) ...

• پرکاربردترین روش تخمین پارامتری

• تخمین توزیع n نمونه داده i.i.d به صورت $p(\mathbf{x} | \Phi) = X_1, \dots, X_n$

• فرض: پارامترهای Φ دارای مقادیر ثابت، اما نامشخص هستند

○ اگر تابع توزیع گاوسی باشد، $\Phi = \{\mu, \Sigma\}$

• تخمین پارامترهای توزیع به نحوی که احتمال بدست آوردن نمونه داده‌ها از روی این توزیع بیشینه باشد

تابع درست‌نمایی

$$p_n(\mathbf{x} | \Phi) = \prod_{k=1}^n p(x_k | \Phi)$$

$$\Phi_{MLE} = \underset{\Phi}{\operatorname{argmax}} p_n(\mathbf{x} | \Phi)$$

• چون متغیرهای تصادفی مستقل هستند

• هدف: بیشینه کردن تابع درست‌نمایی

• لگاریتم شباهت (Log-Likelihood)

○ عدم تغییر مساله (تابع یکنوای صعودی)

$$l(\Phi) = \log p_n(\mathbf{x} | \Phi) = \sum_{k=1}^n \log p(x_k | \Phi) \quad \text{○ ساده کردن محاسبات و فرمول‌ها (تبدیل ضرب به جمع)}$$



نظریه تخمین ...

○ تخمین بیشینه شباهت (MLE) ...

• بیشینه کردن تابع درست‌نمایی (یا لگاریتم آن) با گرادیان

○ مشتق گرفتن بر حسب پارامترها و برابر صفر قرار دادن

$$\nabla_{\Phi} = \begin{bmatrix} \frac{\partial}{\partial \Phi_1} \\ \vdots \\ \frac{\partial}{\partial \Phi_k} \end{bmatrix}$$

$$\nabla_{\Phi} l(\Phi) = \sum_{k=1}^n \nabla_{\Phi} \log p(x_k | \Phi) = 0$$

$$p(x | \Phi) = \frac{1}{\sqrt{2\pi}\sigma} \exp\left[-\frac{(x-\mu)^2}{2\sigma^2}\right]$$

• مثال ۱: توزیع گاوسی تک متغیره

○ لگاریتم تابع درست‌نمایی

$$\log p_n(\mathbf{x} | \Phi) = \sum_{k=1}^n \log p(x_k | \Phi) = \sum_{k=1}^n \log\left(\frac{1}{\sqrt{2\pi}\sigma} \exp\left[-\frac{(x_k - \mu)^2}{2\sigma^2}\right]\right) = -\frac{n}{2} \log(2\pi\sigma^2) - \frac{1}{2\sigma^2} \sum_{k=1}^n (x_k - \mu)^2$$

$$\frac{\partial}{\partial \mu} \log p_n(x | \Phi) = \sum_{k=1}^n \frac{1}{\sigma^2} (x_k - \mu) = 0$$

$$\frac{\partial}{\partial \sigma^2} \log p_n(x | \Phi) = -\frac{n}{2\sigma^2} + \sum_{k=1}^n \frac{(x_k - \mu)^2}{2\sigma^4} = 0$$



$$\mu_{MLE} = \frac{1}{n} \sum_{k=1}^n x_k = E(x)$$

$$\sigma_{MLE}^2 = \frac{1}{n} \sum_{k=1}^n (x_k - \mu_{MLE})^2 = E[(x - \mu_{MLE})^2]$$

همان میانگین و واریانس نمونه‌ای



نظریه تخمین ...

○ تخمین بیشینه شباهت (MLE)

- مثال ۲: توزیع گاوسی چندمتغیره

$$p(\mathbf{x} | \Phi) = \frac{1}{(2\pi)^{d/2} |\Sigma|^{1/2}} \exp\left[-\frac{1}{2}(\mathbf{x} - \boldsymbol{\mu})^t \Sigma^{-1}(\mathbf{x} - \boldsymbol{\mu})\right]$$



$$\hat{\boldsymbol{\mu}}_{MLE} = \frac{1}{n} \sum_{k=1}^n \mathbf{x}_k$$

$$\hat{\Sigma}_{MLE} = \frac{1}{n} \sum_{k=1}^n (\mathbf{x}_k - \hat{\boldsymbol{\mu}}_{MLE})(\mathbf{x}_k - \hat{\boldsymbol{\mu}}_{MLE})^t = E[(\mathbf{x}_k - \hat{\boldsymbol{\mu}}_{MLE})(\mathbf{x}_k - \hat{\boldsymbol{\mu}}_{MLE})^t]$$

$$\mathcal{E} \left[\frac{1}{n} \sum_{i=1}^n (x_i - \hat{\mu})^2 \right] = \frac{n-1}{n} \sigma^2 \neq \sigma^2$$

- تخمین ML برای واریانس بایاس شده است

○ امید ریاضی واریانس تخمینی با واریانس واقعی برابر نیست

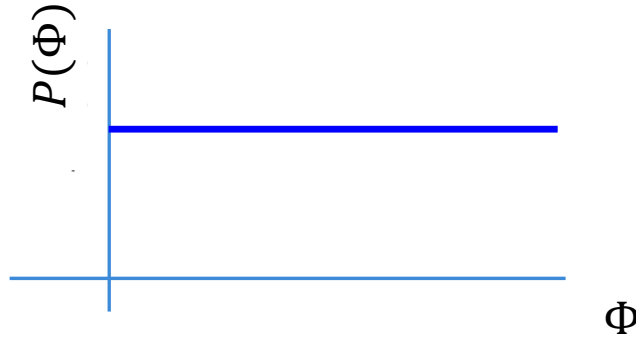
○ با میل کردن n به سمت بی نهایت اثر بایاس کم می شود

$$\hat{\sigma}_{\text{unbiased}}^2 = \frac{1}{n-1} \sum_{k=1}^n (x_k - \hat{\mu})^2$$

- تخمین غیربایاس



نظریه تخمین ...

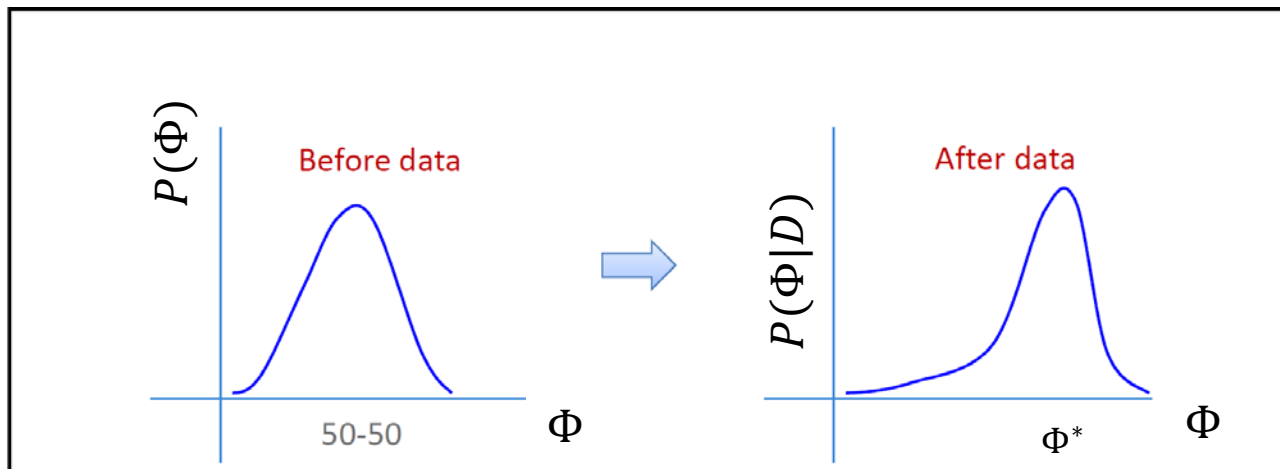


○ در تخمین ML

- فرض = ثابت بودن پارامترهای توزیع
- توزیع $p(\Phi)$ یکنواخت است

○ تخمین بیز (Bayesian)

- پارامتر Φ یک متغیر تصادفی است





نظریه تخمین ...

○ تخمین بیز (Bayesian Estimation) ...

- پارامتر Φ یک متغیر تصادفی و نامشخص است
 - در تخمین ML این پارامتر تصادفی نیست و ثابت است
 - تصادفی بودن پارامتر Φ به معنی وجود احتمال پیشین توزیع $p(\Phi)$ برای آن است
- شکل تابع توزیع $p(\mathbf{x} | \Phi)$ مشخص است (مثلاً توزیع نرمال)
- مجموعه n نمونه داده $\mathcal{D} = \{\mathbf{x}_1, \mathbf{x}_2, \dots, \mathbf{x}_n\}$ i.i.d دارای توزیع $p(\mathbf{x} | \Phi)$
 - داده‌ها حاوی اطلاعاتی از پارامتر Φ

احتمال پسین: احتمال پارامتر پس از مشاهده داده‌ها

$$p(\Phi | \mathbf{x}) = \frac{p(\mathbf{x} | \Phi)p(\Phi)}{p(\mathbf{x})} \propto p(\mathbf{x} | \Phi)p(\Phi) \quad \bullet \text{ با توجه به قانون بیز}$$

$$\text{posterior} = \frac{\text{likelihood} \times \text{prior}}{\text{evidence}}$$

$$p(x | D) = \int p(x | \Phi)p(\Phi | D)d\Phi \quad \bullet \text{ هدف}$$



نظریه تخمین ...

○ تخمین بیز (Bayesian Estimation) ...

$$p(x | D) = \int p(x | \Phi) p(\Phi | D) d\Phi$$

- نخستین گام در تخمین بیز: محاسبه $p(\Phi | D)$

$$p(\Phi | D) = \frac{1}{\alpha} p(D | \Phi) p(\Phi)$$

- که α ثابت نرمال کننده است (مستقل از Φ)
- حذف نمی شود
- در عمل می توان مقدار ثابتی را به عنوان تخمین آن قرار داد

- با فرض مستقل بودن نمونه داده ها

$$p(D | \Phi) = \prod_{k=1}^n p(x_k | \Phi)$$



نظریه تخمین ...

○ تخمین بیز (Bayesian Estimation) ...

- مثال: داده‌ها با توزیع گاوسی، تک متغیره، واریانس معلوم σ^2 و میانگین نامعلوم Φ

$$p(\mathbf{x} | \Phi) = \frac{1}{(2\pi)^{n/2} \sigma^n} \exp\left[-\frac{1}{2} \sum_{i=1}^n \left(\frac{x_i - \Phi}{\sigma}\right)^2\right] \propto \exp\left[-\frac{1}{2} \sum_{i=1}^n \left(\frac{x_i - \Phi}{\sigma}\right)^2\right]$$

$$p(\mathbf{x} | \Phi) \propto \exp\left[-\frac{n}{2\sigma^2} (\Phi - \bar{x}_n)^2\right] \exp\left[-\frac{1}{2\sigma^2} \sum_{i=1}^n (x_i - \bar{x}_n)^2\right]$$

$$\sum_{i=1}^n (x_i - \Phi)^2 = n(\Phi - \bar{x}_n)^2 + \sum_{i=1}^n (x_i - \bar{x}_n)^2$$

$$\bar{x}_n = \frac{1}{n} \sum_{i=1}^n x_i$$

- فرض: $p(\Phi)$ نیز توزیع گاوسی با میانگین μ و واریانس v^2 دارد

○ فرض گاوسی بودن توزیع بودن اولیه ضروری نیست و می‌تواند هر توزیع دیگری برای آن فرض شود

$$p(\Phi) = \frac{1}{(2\pi)^{1/2} v} \exp\left[-\frac{1}{2} \left(\frac{\Phi - \mu}{v}\right)^2\right] \propto \exp\left[-\frac{1}{2} \left(\frac{\Phi - \mu}{v}\right)^2\right]$$

$$p(\Phi | \mathbf{x}) \propto \exp\left\{-\frac{1}{2} \left[\frac{n}{\sigma^2} (\Phi - \bar{x}_n)^2 + \frac{1}{v^2} (\Phi - \mu)^2 \right]\right\}$$

- بنابراین



نظریه تخمین ...

○ تخمین بیز (Bayesian Estimation)

- مثال: داده‌ها با توزیع گاوسی، تک متغیره، واریانس معلوم σ^2 و میانگین نامعلوم Φ

$$p(\Phi | \mathbf{x}) \propto \exp \left\{ -\frac{1}{2} \left[\frac{n}{\sigma^2} (\Phi - \bar{x}_n)^2 + \frac{1}{v^2} (\Phi - \mu)^2 \right] \right\} \quad \rho = \frac{\sigma^2 \mu + n v^2 \bar{x}_n}{\sigma^2 + n v^2} \quad \tau^2 = \frac{\sigma^2 v^2}{\sigma^2 + n v^2}$$

$$p(\Phi | \mathbf{x}) \propto \exp \left\{ -\frac{1}{2} \left[\frac{1}{\tau^2} (\Phi - \rho)^2 + \frac{n}{\sigma^2 + n v^2} (\bar{x}_n - \mu)^2 \right] \right\}$$

مقدار ثابت

$$p(\Phi | \mathbf{x}) = \frac{1}{\sqrt{2\pi}\tau} \exp \left[\frac{-1}{2\tau^2} (\Phi - \rho)^2 \right]$$

- توزیع گاوسی با میانگین ρ و واریانس τ^2

○ ترکیبی از اطلاعات اولیه (μ, v^2) و اطلاعات نمونه‌ها (جمع وزن دار میانگین)

○ با افزایش تعداد نمونه‌ها ($n \rightarrow \infty$), داریم:

○ واریانس به صفر میل می‌کند، عدم قطعیت تخمین میانگین کم می‌شود (واریانس = میزان عدم قطعیت میانگین تخمینی است)

○ میانگین، به میانگین نمونه‌ای نمونه‌ها میل می‌کند



نظریه تخمین ...

○ تخمین بیشینه احتمال پسین (MAP: maximum a posteriori) ...

$$p(\Phi | \mathbf{x}) = \frac{p(\mathbf{x} | \Phi)p(\Phi)}{p(\mathbf{x})} \propto p(\mathbf{x} | \Phi)p(\Phi)$$

• هدف: بیشینه کردن $p(\Phi | \mathbf{x})$

○ در ML هدف بیشینه کردن $p(\mathbf{x} | \Phi)$

○ پارامترها متغیرهای تصادفی با توزیع پیشین $p(\Phi)$ هستند

○ متداولترین تخمین گر بیزی

$$\Phi_{MAP} = \theta_{MAP}(\mathbf{x}) = \underset{\Phi}{\operatorname{argmax}} p(\Phi | \mathbf{x}) = \underset{\Phi}{\operatorname{argmax}} p(\mathbf{x} | \Phi)p(\Phi)$$

$$\Phi_{MAP} = \underset{\Phi}{\operatorname{argmax}} \log p(\mathbf{x} | \Phi) + \log p(\Phi)$$

$$\frac{\partial \log p(\mathbf{x} | \Phi)}{\partial \Phi} + \frac{\partial \log p(\Phi)}{\partial \Phi} = 0$$

$$\left. \frac{\partial \log p(\mathbf{x} | \Phi)}{\partial \Phi} \right|_{\Phi=\Phi_{MAP}} = - \left. \frac{\partial \log p(\Phi)}{\partial \Phi} \right|_{\Phi=\Phi_{MAP}}$$

• تخمین ML و MAP یکسان هستند وقتی توزیع پیشین $p(\Phi)$ یکنواخت باشد



نظریه تخمین

○ تخمین بیشینه احتمال پسین (MAP: maximum a posteriori)

- مثال (قبلی): داده‌ها با توزیع گاوسی، تک متغیره، واریانس معلوم σ^2 و میانگین نامعلوم Φ
- توزیع پیشین Φ : گاوسی با میانگین μ و واریانس ν^2

$$p(\Phi | \mathbf{x}) = \frac{1}{\sqrt{2\pi\tau}} \exp\left[-\frac{1}{2\tau^2}(\Phi - \rho)^2\right]$$

○ مشتق‌گیری نسبت به Φ

○ متوسط وزن‌دار میانگین نمونه‌ها و میانگین قبلی

$$\Phi_{MAP} = \rho = \frac{\sigma^2 \mu + n\nu^2 \bar{x}_n}{\sigma^2 + n\nu^2}$$

میانگین نمونه‌ها تعداد نمونه‌ها

• کاربرد در تطبیق (adaptation)

- آموزش صدای یک کاربر جدید به سیستم بازشناسی گفتار
- آموزش پارامترهای مدل با پایگاه داده‌های مستقل از گوینده (با چندین گوینده) = توزیع پیشین
- تطبیق با محاسبه میانگین نمونه‌های یک گوینده خاص و استفاده از رابطه بالا



مقایسه تخمین‌گرها

تخمین‌گر بیشینه احتمال پسین (MAP)	تخمین‌گر بیشینه شباهت (ML)	تخمین‌گر بیز
تخمین‌گر نقطه	تخمین‌گر نقطه	تخمین‌گر توزیع
$p(\mathbf{x} \mathcal{D}) = p(\mathbf{x} \hat{\theta})$	$p(\mathbf{x} \mathcal{D}) = p(\mathbf{x} \hat{\theta})$	$p(\mathbf{x} \mathcal{D}) = \int p(\mathbf{x} \theta)p(\theta \mathcal{D})d\theta$
تخمین نقطه	تخمین نقطه	تخمین توزیع
$\hat{\theta} = \arg \max_{\theta} \ln p(\mathcal{D} \theta)p(\theta)$	$\hat{\theta} = \arg \max_{\theta} \ln p(\mathcal{D} \theta)$	$p(\theta \mathcal{D}) = \frac{1}{\alpha} p(\mathcal{D} \theta)p(\theta)$
1. استفاده از اطلاعات پیشین پارامتر توزیع	1. تفسیر ساده‌تر (نقطه‌ای) 2. محاسبات کم‌تر	1. استفاده بیشتر از اطلاعات 2. کارایی بهتر در صورت عدم سازگاری بین توزیع فرض شده و توزیع واقعی 3. در نظر گرفتن بایاس واریانس

