

Lecture 7 - Jan 31

Crash Course in Probability + Statistics

Reading

Data Mining and Machine Learning

2: Numerical Attributes

Upcoming Deadlines

Homework 2 (Feb 10)

A random experiment results in one of a number of outcomes from a sample space Ω

Subsets of the sample space $E \subseteq \Omega$ are events

A probability measure is a function $P: \{\text{events}\} \rightarrow [0, 1]$ where

$$\textcircled{1} \quad P(\Omega) = 1$$

\textcircled{2} If A_1, A_2, \dots are disjoint events,

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

Properties of probability measures

$$P(A^c) = 1 - P(A)$$

A does not occur

probability of complement

$$P(A \cup B) = P(A) + P(B) - P(A \cap B)$$

principle of inclusion-exclusion

$$P(A \cup B) = P(A) + P(B) - P(A \cap B)$$

Principle of
Inclusion-exclusion

Conditional Probability

The probability event A occurs given event B has occurred is

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

Conditional probability

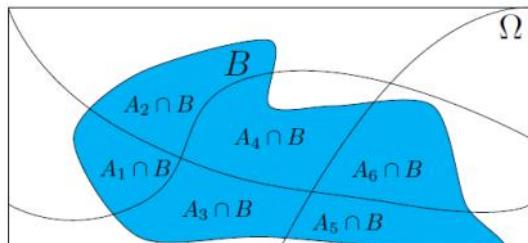
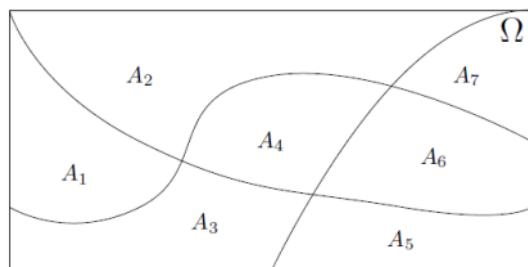
$P(A)$ - prior probability $P(A|B)$ - posterior probability

The Law of Total Probability

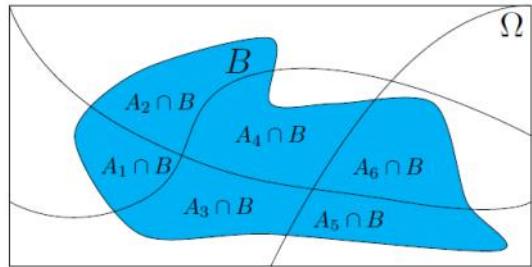
If events A_1, A_2, \dots partition the sample space Ω , then

$$P(B) = \sum_{i=1}^{\infty} P(A_i \cap B)$$

$$= \sum_{i=1}^{\infty} P(B|A_i)P(A_i)$$



$$= \sum_{i=1}^{\infty} P(B|A_i)P(A_i)$$



Bayes' Theorem

If $P(A) > 0$ and $P(B) > 0$,

$$P(A|B) = \frac{P(B|A)P(A)}{P(B)}$$

Version 2 If A_1, A_2, \dots partition Ω ,

$$P(A_i|B) = \frac{P(B|A_i)P(A_i)}{P(B)}$$

$$P(A_i|B) = \frac{P(B|A_i)P(A_i)}{\sum_{j=1}^{\infty} P(B|A_j)P(A_j)}$$

Events A and B are independent if $P(A \cap B) = P(A)P(B)$
 or .. $P(A|B) = P(A)$.. , $P(B|A) = P(B)$

Events A and B are independent if they occur

↳ this implies $P(A|B) = \frac{P(A)P(B)}{P(B)} = P(A)$, i.e. event B's occurrence has no impact on A.

A random variable is a measurable function $X: \Omega \rightarrow \mathbb{R}$ mapping outcomes to numbers.

Random variables valued in countable sets are discrete

Random variables valued in intervals (e.g. \mathbb{R}) are continuous

A probability distribution describes how likely a R.V. or set of R.V.s is to take different values.

In the discrete case, a probability mass function (PMF) f fully describes a probability distribution

$$f(x) = P(X=x)$$

for each x in the range of X .

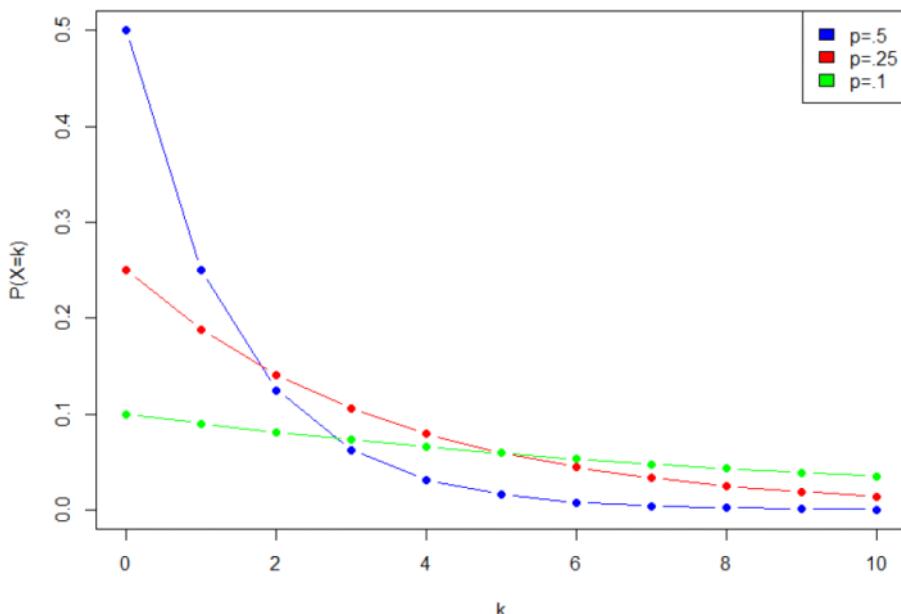
Example: Geometric R.V.

Let X be the number of failures before the first success in a sequence of independent binary experiments, each with

a sequence of independent binary experiments, each with probability P of success.

$$P(X=k) = f(k) = \begin{cases} (1-p)^k p, & \text{if } k=0, 1, 2, \dots \\ 0, & \text{else} \end{cases}$$

Geometric Probability Mass Functions



In the continuous case, a probability density function (pdf) f fully describes a probability distribution

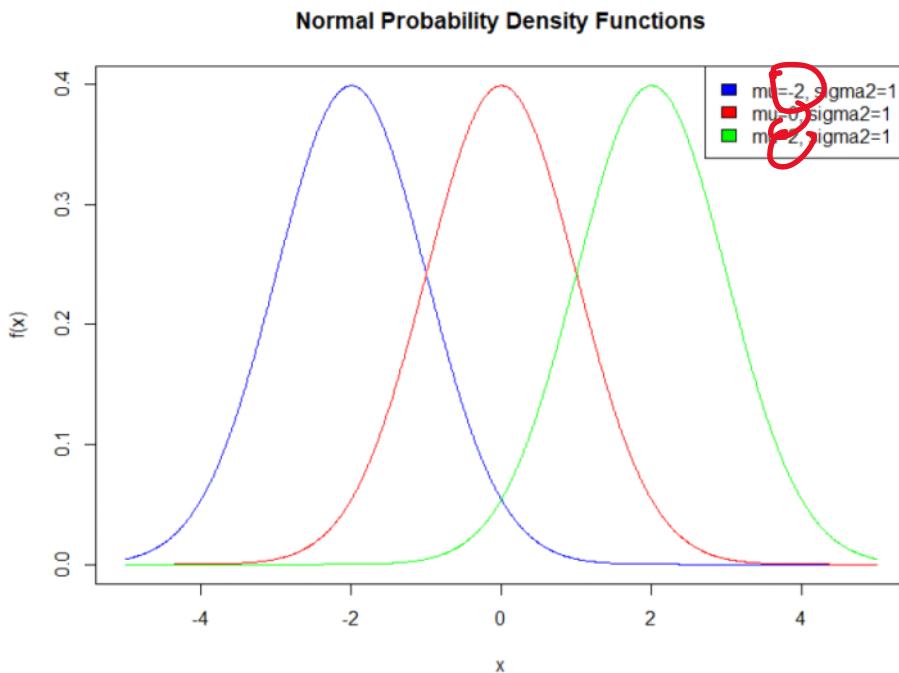
$$\int_E f(x) dx = P(X \in E)$$

for a set $E \subseteq \mathbb{R}$.

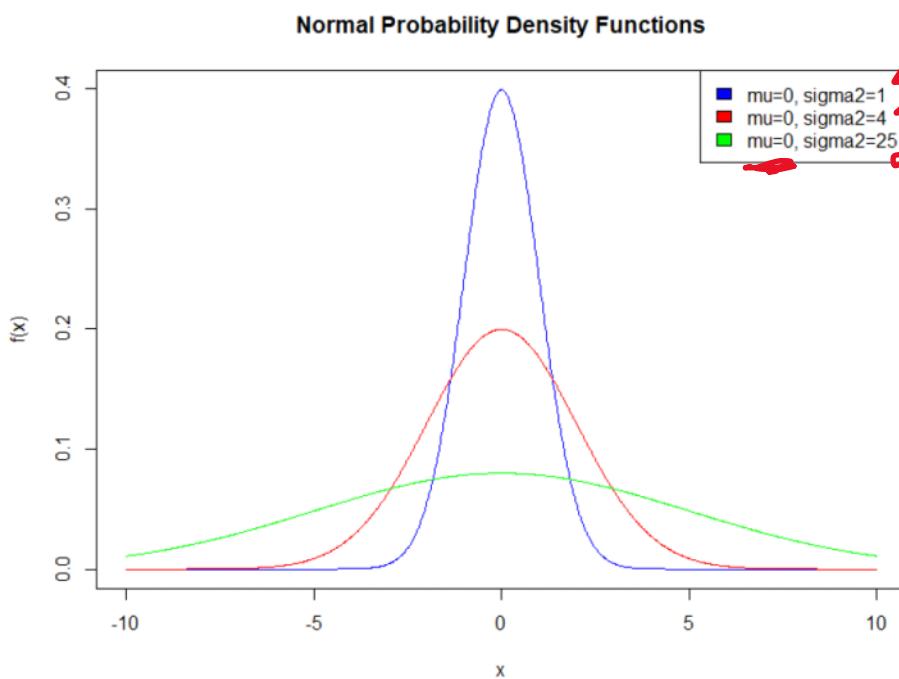
Example : Normal (Gaussian) RV

X has a normal distribution with parameters $\mu \in \mathbb{R}$, $\sigma > 0$
 if its pdf is

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



μ is the center
 of the distribution



σ controls how dispersed the probability is
 large $\sigma \rightarrow$ more "spread out"
 small $\sigma \rightarrow$ more localized around μ

If $\mu=0$, $\sigma=1$, X is a standard normal R.V.

The expected value of a R.V. X with range R is

$$E[X] = \sum_{x \in R} x P(X=x) = \sum_{x \in R} xf(x) \quad (\text{discrete})$$

$$E[X] = \int_R x f(x) dx \quad (\text{continuous})$$

This is the mean value X takes

Properties of expectation

① If $g(x)$ is a R.V., $E[g(x)] = \begin{cases} \sum_{x \in R} g(x)f(x), & X \text{ is discrete} \\ \int_R g(x)f(x) dx, & X \text{ is continuous} \end{cases}$

② $E[aX+bY] = aE[X] + bE[Y]$

③ If X, Y are independent, $E[XY] = E[X]E[Y]$

The variance of X is $\text{Var}(X) = E[(X - E[X])^2]$
SSE between X & its mean

↳ measures how "spread out" the pdf of X is

Properties of Variance

$$\textcircled{1} \quad \text{Var}(X) = E[X^2] - E[X]^2$$

$$\textcircled{2} \quad \text{Var}(aX + bY) = a^2\text{Var}(X) + b^2\text{Var}(Y) \quad \underline{\text{if}} \quad X, Y \text{ are independent}$$

If X_1, X_2, \dots, X_n are discrete R.V.s, their joint PMF is

$$f(x_1, \dots, x_n) = P(X_1=x_1, \dots, X_n=x_n)$$

If they are continuous, the joint pdf is f such that

$$\int_E f(x_1, \dots, x_n) dx_1 \cdots dx_n = P((X_1, \dots, X_n) \in E)$$

for any $E \subseteq \mathbb{R}^n$.

The expected value of some R.V. $g(X_1, \dots, X_n)$ is

$$E[g(X_1, \dots, X_n)] = \begin{cases} \sum_{x_1} \cdots \sum_{x_n} g(x_1, \dots, x_n) f(x_1, \dots, x_n), & \text{if discrete} \\ \int_R \cdots \int_R g(x_1, \dots, x_n) f(x_1, \dots, x_n) dx_1 \cdots dx_n, & \text{if continuous} \end{cases}$$

The covariance between $X_j + X_k$ is

$$\text{cov}(X_j, X_k) = E[(X_j - E[X_j])(X_k - E[X_k])] = \sigma_{jk}$$

Covariance quantities how dependent $X_j + X_k$ are, but its value is dependent on the scale of $X_j + X_k$...

Correlation is a normalized version valued in $[0,1]$.

The Covariance matrix for RVs X_1, \dots, X_d is

$$\Sigma = \begin{pmatrix} \sigma_1^2 & \sigma_{12} & \cdots & \sigma_{1d} \\ \sigma_{21} & \sigma_2^2 & \cdots & \sigma_{2d} \\ \vdots & \vdots & \ddots & \vdots \\ \sigma_{d1}, \sigma_{d2}, \dots, \sigma_d^2 \end{pmatrix} \quad \leftarrow \text{Square, symmetric matrix}$$

The correlation of $X_j + X_k$ is

$$\text{Corr}(X_j, X_k) = \frac{\text{Cov}(X_j, X_k)}{\sigma_j \sigma_k}$$

Multivariate Normal (Gaussian) Distribution

A random vector $X = (X_1, \dots, X_d)$ has a MV. normal distribution with mean vector $\mu \in \mathbb{R}^d$ and covariance matrix Σ if its

Mahalanobis distance from x to μ takes into account variance & cov. mat.

with mean vector $\mu \in \mathbb{R}^d$ and covariance matrix Σ , the joint pdf is

$$f(x|\mu, \Sigma) = \frac{1}{(2\pi)^{d/2} \sqrt{|\Sigma|}} \exp\left(-\frac{(x-\mu)^T \Sigma^{-1} (x-\mu)}{2}\right)$$

Mahalanobis distance from x to μ
 taking into account variance & cov. mat.

↑
 vector

X has a standard MV. normal distribution if $\mu = \begin{pmatrix} 0 \\ \vdots \\ 0 \end{pmatrix}$ and

$$\Sigma = I_d = \begin{pmatrix} 1 & 0 & \cdots & 0 \\ 0 & 1 & \cdots & 0 \\ \vdots & \vdots & \ddots & \vdots \\ 0 & 0 & \cdots & 1 \end{pmatrix} \quad \rightarrow \text{here, each } X_i \text{ is an independent standard normal RV}$$

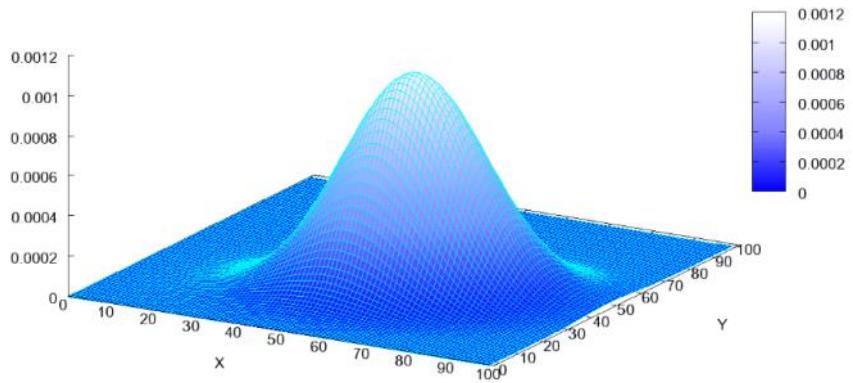
↳ $f(x|0, I) = \frac{1}{(2\pi)^{d/2}} \exp\left(-\frac{(x-\mu)^T (x-\mu)}{2}\right) = \frac{1}{(2\pi)^{d/2}} \exp\left(-\frac{\|x-\mu\|^2}{2}\right)$

If $d=2$, the multivariate normal distribution simplifies to

$$f(x, y) = \frac{1}{2\pi\sigma_1\sigma_2\sqrt{1-\rho^2}} e^{-\frac{1}{2(1-\rho^2)} \left[\frac{(x-\mu_1)^2}{\sigma_1^2} - \frac{2\rho(x-\mu_1)(y-\mu_2)}{\sigma_1\sigma_2} + \frac{(y-\mu_2)^2}{\sigma_2^2} \right]}$$

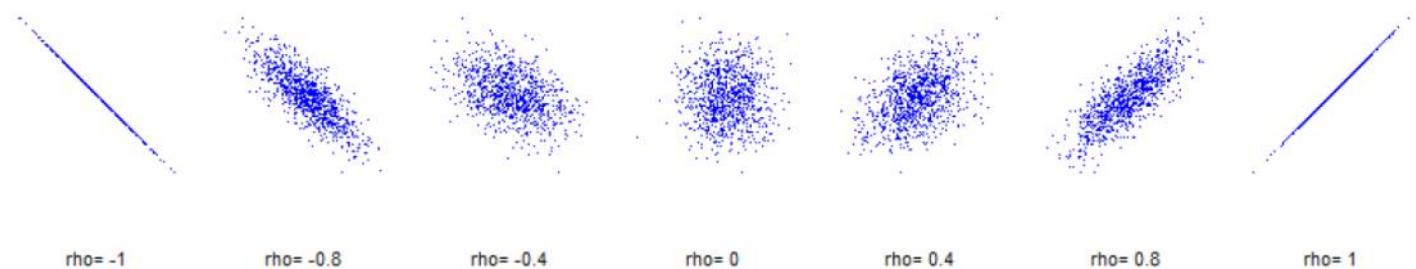
If $\rho = 0$,

Multivariate Normal Distribution



BivariateNormal

If ρ varies



A random sample is R.V.s X_1, \dots, X_n that are independent and identically distributed (i.i.d.)

↳ a sample comes from some (usually unknown) population with some fixed but usually unknown distribution. Parameters of this distribution are population parameters

A sample statistic is a R.V. $g(\underbrace{X_1, \dots, X_n}_{\text{sample}})$ for some given g .

Univariate statistics

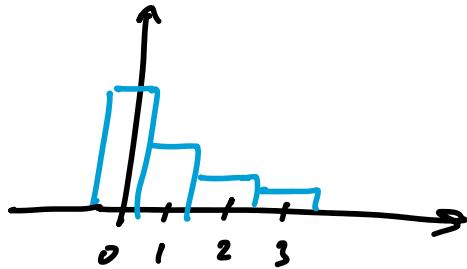
Let $D = \begin{pmatrix} x_1 \\ \vdots \\ x_n \end{pmatrix}$ be values of the R.V.s in a sample.

We treat it as a point in \mathbb{R}^n .

Let X be a R.V. with the population's distribution.

The empirical PMF of X is

$$\hat{f}(x) = \frac{1}{n} \sum_{i=1}^n \mathbf{1}_{\{x_i = x\}}$$



$\mathbf{1}_A = \begin{cases} 1, & \text{if } A \text{ is true} \\ 0, & \text{else} \end{cases}$ is an indicator function

The sample mean of X is $\hat{\mu} = \frac{1}{n} \sum_{i=1}^n x_i$

The centered feature vector of X is $\bar{X} = \begin{pmatrix} x_1 - \hat{\mu} \\ \vdots \\ x_n - \hat{\mu} \end{pmatrix}$

The sample mode of X is $\text{mode}(x) = \arg \max_x \hat{f}(x)$

The sample variance of X is $\hat{\sigma}^2 = \frac{1}{n} \sum_{i=1}^n (x_i - \hat{\mu})^2 = \frac{1}{n} \|\bar{X}\|^2$

↳ note $\frac{1}{n} \sum_{i=1}^n (x_i - \hat{\mu})^2 = \frac{1}{n} \bar{X}^T \bar{X} = \frac{1}{n} \|\bar{X}\|^2$

The standard score (z-score) of x_i is $z_i = \frac{x_i - \hat{\mu}}{\hat{\sigma}}$

Multivariate Statistics

$$D = \begin{pmatrix} x_{11} & x_{12} & \cdots & x_{1d} \\ \vdots & \vdots & \ddots & \vdots \\ x_{n1} & x_{n2} & \cdots & x_{nd} \end{pmatrix} \quad \begin{matrix} \leftarrow x_1^T \\ \leftarrow x_n^T \\ \uparrow \\ x_1 \\ \uparrow \\ x_2 \\ \uparrow \\ x_d \end{matrix}$$

Below, we give multivariate versions of sample statistics

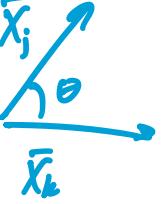
The sample mean vector is $\hat{\mu} = \frac{1}{n} \sum_{i=1}^n x_i = \begin{pmatrix} \hat{\mu}_1 \\ \vdots \\ \hat{\mu}_n \end{pmatrix}$

The sample covariance between $X_j + X_k$ is $\hat{G}_{jk} = \frac{1}{n} \sum_{i=1}^n (x_{ij} - \hat{\mu}_j)(x_{ik} - \hat{\mu}_k)$
 $= \underline{\underline{\frac{1}{n} \bar{X}_j^T \bar{X}_k}}$

↳ notice $\sum (x_{ij} - \hat{\mu}_j)(x_{ik} - \hat{\mu}_k) = \bar{X}_j^T \bar{X}_k$

The sample correlation between $X_j + X_k$ is $\frac{\hat{G}_{jk}}{\hat{\sigma}_j \hat{\sigma}_k} = \frac{\bar{X}_j^T \bar{X}_k}{\|\bar{X}_j\| \cdot \|\bar{X}_k\|}$ 

$$\frac{\hat{G}_{jk}}{\hat{\sigma}_j \hat{\sigma}_k} = \frac{\frac{1}{n} \bar{X}_j^T \bar{X}_k}{\sqrt{\frac{1}{n} \bar{X}_j^T \bar{X}_j} \sqrt{\frac{1}{n} \bar{X}_k^T \bar{X}_k}} = \frac{\cancel{\frac{1}{n}} \bar{X}_j^T \bar{X}_k}{\cancel{\frac{1}{n}} \sqrt{\bar{X}_j^T \bar{X}_j} \sqrt{\bar{X}_k^T \bar{X}_k}} = \frac{\bar{X}_j^T \bar{X}_k}{\|\bar{X}_j\| \|\bar{X}_k\|} = \cos \theta$$

where  this is called the cosine similarity between the vectors in some contexts.

The sample covariance matrix is $\hat{\Sigma} = \begin{pmatrix} \hat{G}_1^2, \hat{G}_{12}, \dots, \hat{G}_{1d} \\ \hat{G}_{21}, \hat{G}_2^2, \dots, \hat{G}_{2d} \\ \vdots & \vdots & \ddots & \vdots \\ \hat{G}_{d1}, \hat{G}_{d2}, \dots, \hat{G}_d^2 \end{pmatrix}$

by the formulas above, we see $\hat{\Sigma} = \frac{1}{n} \begin{pmatrix} \bar{X}_1^T \bar{X}_1, \bar{X}_1^T \bar{X}_2, \dots, \bar{X}_1^T \bar{X}_d \\ \bar{X}_2^T \bar{X}_1, \bar{X}_2^T \bar{X}_2, \dots, \bar{X}_2^T \bar{X}_d \\ \vdots & \vdots & \ddots & \vdots \\ \bar{X}_d^T \bar{X}_1, \bar{X}_d^T \bar{X}_2, \dots, \bar{X}_d^T \bar{X}_d \end{pmatrix}$

$$= \frac{1}{n} (\bar{D}^T \bar{D})$$

where $\bar{D} = D - \mathbf{1} \cdot \hat{\mu} = \begin{pmatrix} \bar{x}_1^T - \hat{\mu}^T \\ \vdots \\ \bar{x}_d^T - \hat{\mu}^T \end{pmatrix}$ ↑ sample scatter matrix

The sample total variance is $\text{tr}(\hat{\Sigma}) = \hat{G}_1^2 + \dots + \hat{G}_d^2$

The sample generalized variance is $|\hat{\Sigma}| = \det(\hat{\Sigma})$