

Probability & Statistics

Machine learning, 2021

Mansoor Rezghi

Department of Computer science, TMU

Ref:KM, CB

References

- K. Murphy Machine Learning: A Probabilistic Perspective, MIT Press, 2012.(KM)-chapter 2
- C. Bishop, Pattern Recognition and Machine Learning, Springer, 2006.(CB)- chapter 2
- I. Goodfellow, Deep Learning, MIT Press, 2015, chapter 3

Probability: Discrete random variable

Discrete random variable X , which can take on any value from a finite or countable infinite set χ .

$X=x$ event

$P(X=x)$ or $P(x)$

$0 \leq P(x) \leq 1$

$$\sum_{x \in \chi} p(x) = 1$$

Continues random variable X :
Uncertain continues quantity

$$p(x \in (a, b)) = \int_a^b p(x) dx.$$

$$\begin{aligned} p(x) &\geq 0 \\ \int_{-\infty}^{\infty} p(x) dx &= 1. \end{aligned}$$

Basic rules

$$p(A \vee B) = p(A) + p(B) - p(A \wedge B)$$

Joint probability

$$p(A, B) = p(A \wedge B) = p(A|B)p(B)$$

Marginal distribution

$$p(A) = \sum_b p(A, B) = \sum_b p(A|B = b)p(B = b)$$

Conditional Probability

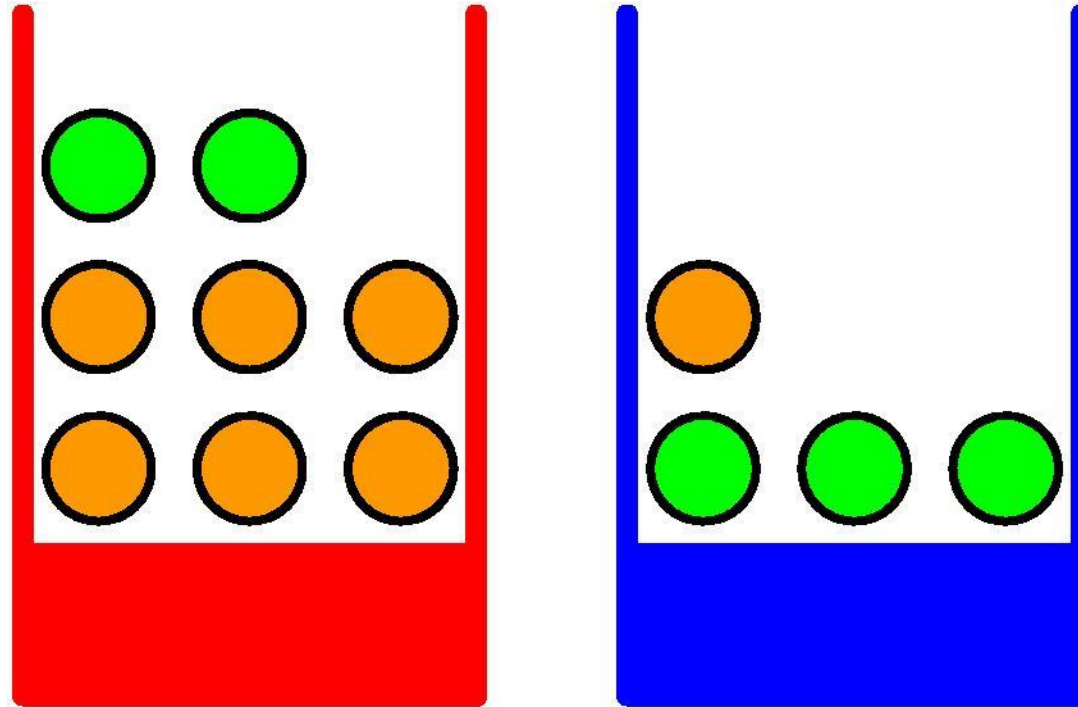
$$p(A|B) = \frac{p(A, B)}{p(B)} \text{ if } p(B) > 0$$

Bayes rule

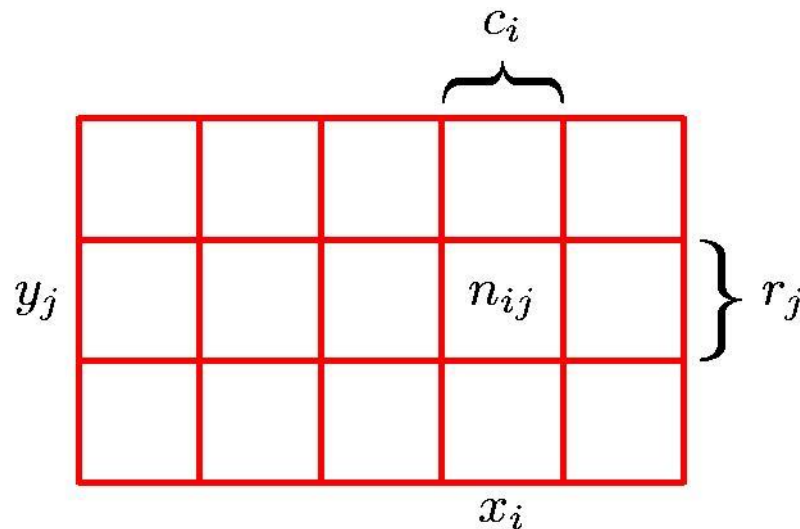
$$p(X = x|Y = y) = \frac{p(X = x, Y = y)}{p(Y = y)} = \frac{p(X = x)p(Y = y|X = x)}{\sum_{x'} p(X = x')p(Y = y|X = x')}$$

Probability Theory

Apples and Oranges



Probability Theory



Joint Probability

$$p(X = x_i, Y = y_j) = \frac{n_{ij}}{N}$$

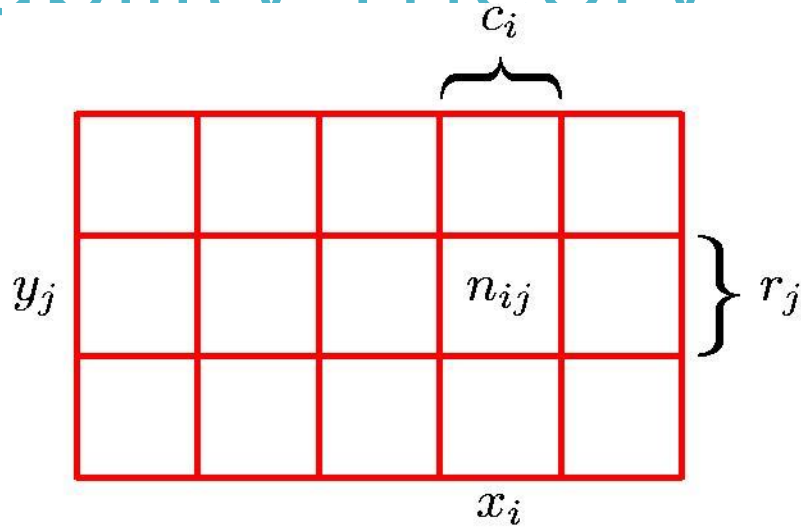
Marginal Probability

$$p(X = x_i) = \frac{c_i}{N}.$$

Conditional Probability

$$p(Y = y_j | X = x_i) = \frac{n_{ij}}{c_i}$$

Probability Theory



Sum Rule

$$\begin{aligned} p(X = x_i) &= \frac{c_i}{N} = \frac{1}{N} \sum_{j=1}^L n_{ij} \\ &= \sum_{j=1}^L p(X = x_i, Y = y_j) \end{aligned}$$

Product Rule

$$\begin{aligned} p(X = x_i, Y = y_j) &= \frac{n_{ij}}{N} = \frac{n_{ij}}{c_i} \cdot \frac{c_i}{N} \\ &= p(Y = y_j | X = x_i) p(X = x_i) \end{aligned}$$

The Rules of Probability

Sum Rule

$$p(X) = \sum_Y p(X, Y)$$

Product Rule

$$p(X, Y) = p(Y|X)p(X)$$

Bayes Theorem

$$p(Y|X) = \frac{p(X|Y)p(Y)}{p(X)}$$

$$p(X) = \sum_Y p(X|Y)p(Y)$$

posterior \propto likelihood \times prior

Bayes rule: Example

Y: you have cancer or NOT, Y=1: You have cancer	X: Results of a test about your cancer X=1: test shows your cancer
--	---

the test has a **sensitivity** of 80%, which means, if you have cancer, the test will be positive with probability 0.8.

$$p(x = 1|y = 1) = 0.8$$

if $p(y = 1) = 0.004$ and

$p(x = 1 y = 0) = 0.1$	false positive or false alarm.
------------------------	--------------------------------



$$\begin{aligned} p(y = 1|x = 1) &= \frac{p(x = 1|y = 1)p(y = 1)}{p(x = 1|y = 1)p(y = 1) + p(x = 1|y = 0)p(y = 0)} \\ &= \frac{0.8 \times 0.004}{0.8 \times 0.004 + 0.1 \times 0.996} = 0.031 \end{aligned}$$

Independence

Independence and conditional independence

We say X and Y are **unconditionally independent** or **marginally independent**, denoted $X \perp Y$, if we can represent the joint as the product of the two marginals (see Figure 2.2), i.e.,

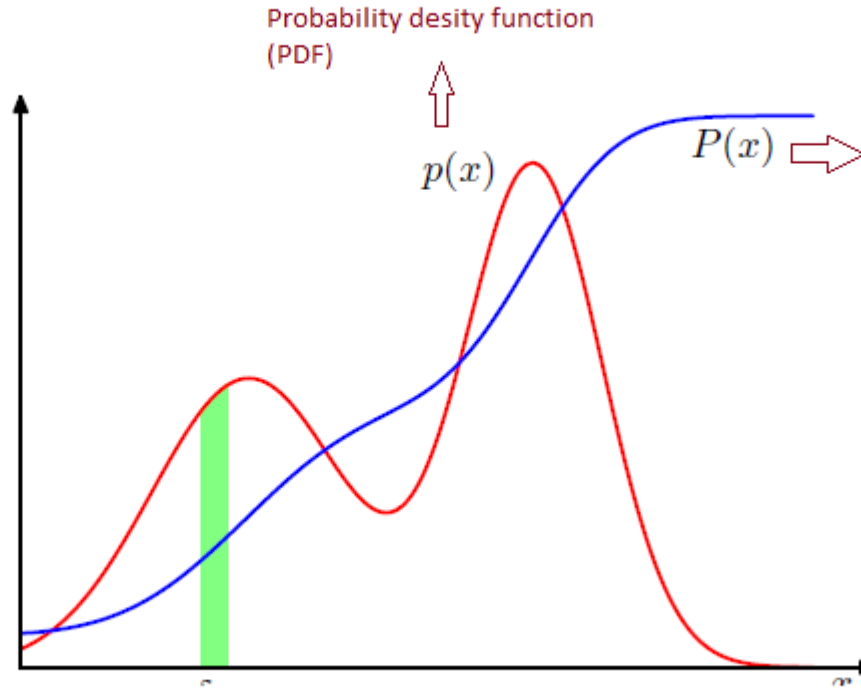
$$X \perp Y \iff p(X, Y) = p(X)p(Y) \quad (2.14)$$

conditionally independent

$$X \perp Y | Z \iff p(X, Y | Z) = p(X | Z)p(Y | Z)$$

Continues Random Variable

The concept of probability for discrete variables can be extended to that of a probability density $p(x)$ over a continuous variable x and is such that the probability of x lying in the interval $(x, x + \delta x)$ is given by $p(x)\delta x$ for $\delta x \rightarrow 0$. The probability density can be expressed as the derivative of a cumulative distribution function $P(x)$.



cumulative
distribution function (CDF)

$$P(z) = \int_{-\infty}^z p(x) dx$$

Important concepts

Expectations & Variance

$$\mathbb{E}[f] = \sum_x p(x) f(x)$$

$$\mathbb{E}[f] = \int p(x) f(x) \, dx$$

Conditional Expectation
(discrete)

$$\mathbb{E}_x[f|y] = \sum_x p(x|y) f(x)$$


Variance

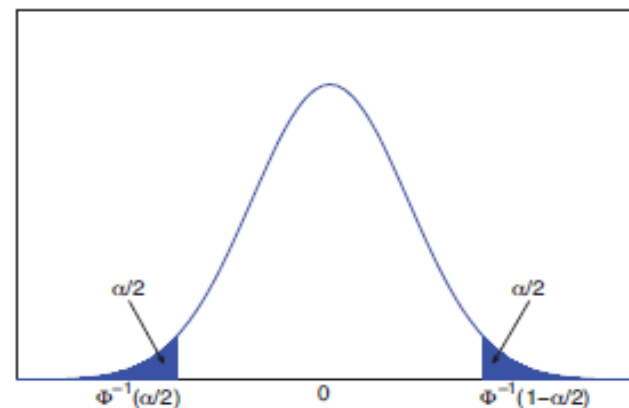
$$\text{var}[f] = \mathbb{E} \left[(f(x) - \mathbb{E}[f(x)])^2 \right] = \mathbb{E}[f(x)^2] - \mathbb{E}[f(x)]^2$$

Quintile

α **quantile** of F is $F^{-1}(\alpha)$ is the value of x_α such that $P(X \leq x_\alpha) = \alpha$

$F^{-1}(0.5)$ is the **median**

$F^{-1}(0.25)$ and $F^{-1}(0.75)$ are the lower and upper **quartiles**



Discrete Distributions

The binomial distribution

Suppose we toss a coin n times. Let $X \in \{0, \dots, n\}$ be the number of heads. If the probability of heads is θ , then we say X has a **binomial** distribution, written as $X \sim \text{Bin}(n, \theta)$. The pmf is given by

$$\text{Bin}(k|n, \theta) \triangleq \binom{n}{k} \theta^k (1 - \theta)^{n-k}$$

$$\text{mean} = \theta, \quad \text{var} = n\theta(1 - \theta)$$

The Bernoulli distribution

$$X \in \{0, 1\}$$

$$X \sim \text{Ber}(\theta),$$

$$\text{Ber}(x|\theta) = \theta^{\mathbb{I}(x=1)} (1 - \theta)^{\mathbb{I}(x=0)}$$

The multinomial distribution

let $\mathbf{x} = (x_1, \dots, x_K)$ be a random vector, where x_j is the number of times side j of the die occurs. Then \mathbf{x} has the following pmf:

$$\text{Mu}(\mathbf{x}|n, \boldsymbol{\theta}) \triangleq \binom{n}{x_1 \dots x_K} \prod_{j=1}^K \theta_j^{x_j} \quad n = \sum_{k=1}^K x_k$$

The Poisson distribution

We say that $X \in \{0, 1, 2, \dots\}$ has a **Poisson** distribution with parameter $\lambda > 0$, written $X \sim \text{Poi}(\lambda)$, if its pmf is

$$\text{Poi}(x|\lambda) = e^{-\lambda} \frac{\lambda^x}{x!}$$

The first term is just the normalization constant, required to ensure the distribution sums to 1.

The Poisson distribution is often used as a model for counts of rare events like radioactive decay and traffic accidents.

Continues distributions

Gaussian distribution

$$\mathcal{N}(x|\mu, \sigma^2) \triangleq \frac{1}{\sqrt{2\pi\sigma^2}} e^{-\frac{1}{2\sigma^2}(x-\mu)^2}$$



$$\mu = \mathbb{E}[X] \quad \sigma^2 = \text{var}[X]$$

Precision: $\lambda = 1/\sigma^2$

Advantages

It has two parameters which are easy to interpret mean and variance.

Central limit Theorem

sums of independent random variables have an approximately Gaussian distribution,

Maximum Entropy

Gaussian distribution makes the least number of assumptions

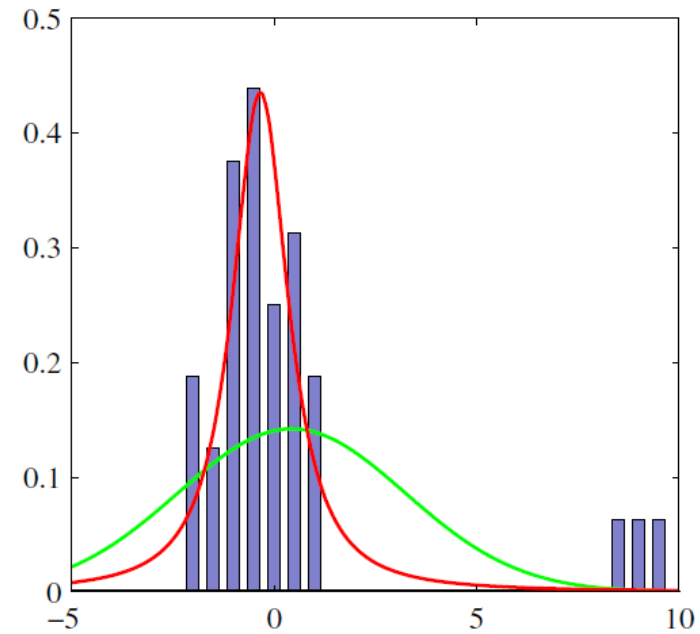
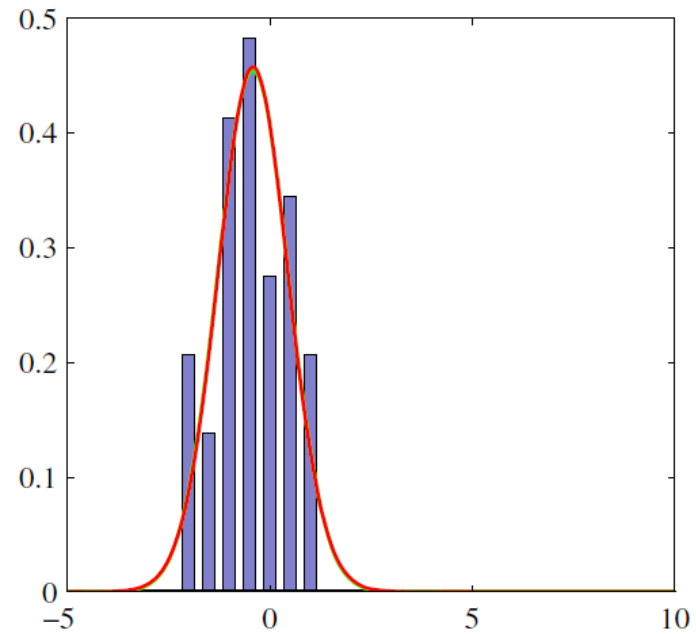
Weakness

Sensitive to outlier

Gaussian and outlier

Red: T-student

Green:
Gaussian



T-student

$$\mathcal{T}(x|\mu, \sigma^2, \nu) \propto \left[1 + \frac{1}{\nu} \left(\frac{x - \mu}{\sigma} \right)^2 \right]^{-\left(\frac{\nu+1}{2}\right)}$$

↑ degree of freedom
↓ Scale Parameter

$$\text{mean} = \mu, \text{mode} = \mu, \text{var} = \frac{\nu\sigma^2}{(\nu-2)}$$

$\nu = 1$, the t-distribution reduces to the *Cauchy* distribution.

$\nu \rightarrow \infty$ the t-distribution $\text{St}(x|\mu, \lambda, \nu)$

becomes a Gaussian $\mathcal{N}(x|\mu, \lambda^{-1})$

with mean μ and precision λ

Able to handle out layer

Laplace distribution

$$\text{Lap}(x|\mu, b) \triangleq \frac{1}{2b} \exp\left(-\frac{|x - \mu|}{b}\right)$$

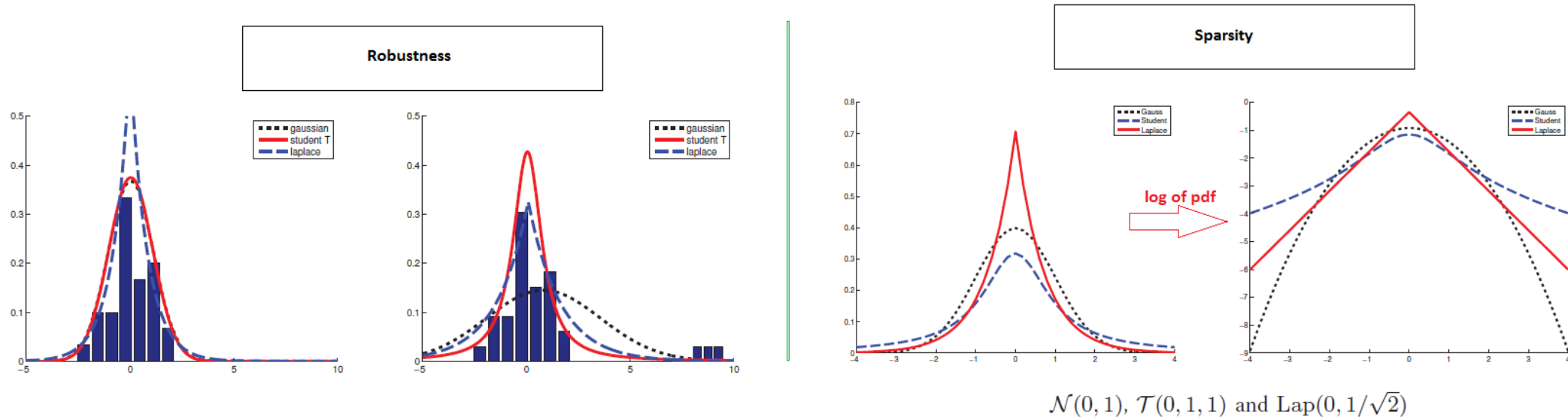
Here μ is a location parameter and $b > 0$ is a scale parameter.

$$\text{mean} = \mu, \text{ mode} = \mu, \text{ var} = 2b^2$$

Robust undr outlier

Sparse distribution

Laplace : Robustness and sparsity



Laplace distribution, which is always log-concave

Multivariate

Multivariate


A **joint probability distribution** has the form $p(x_1, \dots, x_D)$ for a set of $D > 1$ variables, and models the (stochastic) relationships between the variables.

$$\text{cov}[X, Y] \triangleq \mathbb{E}[(X - \mathbb{E}[X])(Y - \mathbb{E}[Y])] = \mathbb{E}[XY] - \mathbb{E}[X]\mathbb{E}[Y]$$

If \mathbf{x} is a d -dimensional random vector, its **covariance matrix** is defined to be the following symmetric, positive definite matrix:

$$\begin{aligned} \text{cov}[\mathbf{x}] &\triangleq \mathbb{E}[(\mathbf{x} - \mathbb{E}[\mathbf{x}])(\mathbf{x} - \mathbb{E}[\mathbf{x}])^T] \\ &= \begin{pmatrix} \text{var}[X_1] & \text{cov}[X_1, X_2] & \cdots & \text{cov}[X_1, X_d] \\ \text{cov}[X_2, X_1] & \text{var}[X_2] & \cdots & \text{cov}[X_2, X_d] \\ \vdots & \vdots & \ddots & \vdots \\ \text{cov}[X_d, X_1] & \text{cov}[X_d, X_2] & \cdots & \text{var}[X_d] \end{pmatrix} \end{aligned}$$

Correlation

$$\text{corr}[X, Y] \triangleq \frac{\text{cov}[X, Y]}{\sqrt{\text{var}[X] \text{var}[Y]}}$$


A **correlation matrix** has the form

$$\mathbf{R} = \begin{pmatrix} \text{corr}[X_1, X_1] & \text{corr}[X_1, X_2] & \cdots & \text{corr}[X_1, X_d] \\ \vdots & \vdots & \ddots & \vdots \\ \text{corr}[X_d, X_1] & \text{corr}[X_d, X_2] & \cdots & \text{corr}[X_d, X_d] \end{pmatrix}$$

$$-1 \leq \text{corr}[X, Y] \leq 1$$

$$\text{corr}[X, Y] = 1 \text{ if and only if } Y = aX + b \quad \mathbf{a > 0}$$

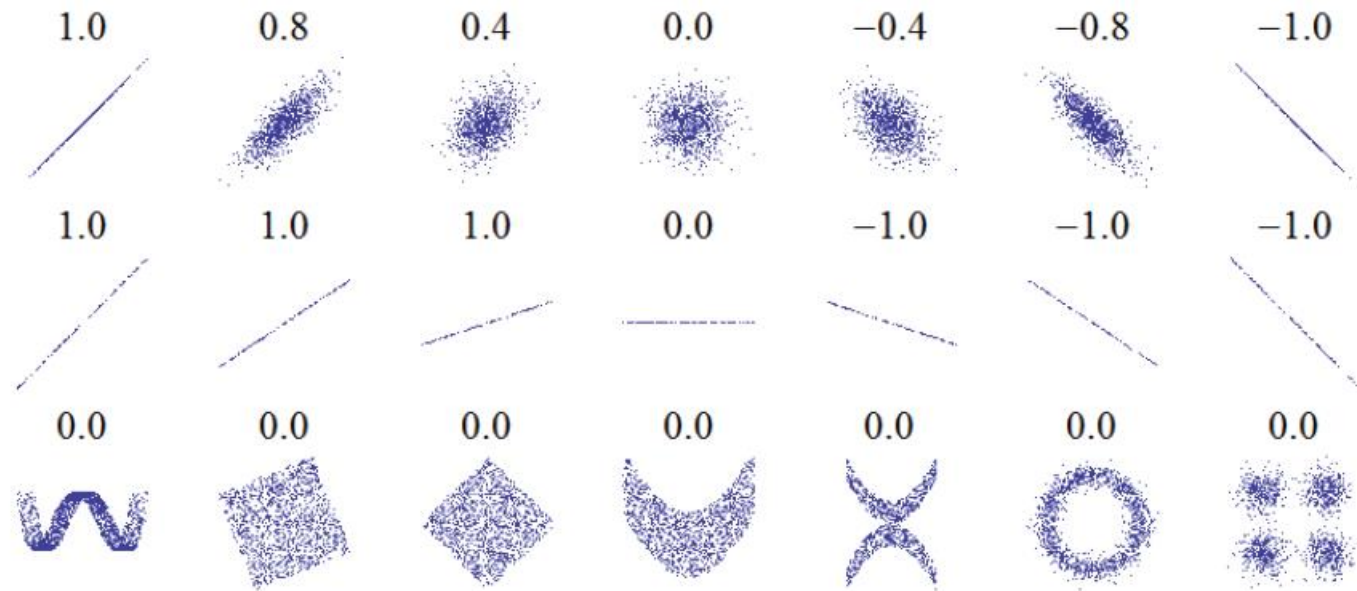


Figure 2.12 Several sets of (x, y) points, with the correlation coefficient of x and y for each set. Note that the correlation reflects the noisiness and direction of a linear relationship (top row), but not the slope of that relationship (middle), nor many aspects of nonlinear relationships (bottom). N.B.: the figure in the center has a slope of 0 but in that case the correlation coefficient is undefined because the variance of Y is zero. Source: http://en.wikipedia.org/wiki/File:Correlation_examples.png

Dependence Vs Correlation

If X and Y are independent, meaning $p(X, Y) = p(X)p(Y)$



$\text{cov}[X, Y] = 0$, and hence $\text{corr}[X, Y] = 0$

so they are uncorrelated.

But $\text{corr}[X, Y] = 0$.

Clearly Y is dependent on X

For example, let $X \sim U(-1, 1)$ and $Y = X^2$.

uncorrelated does not imply independent.

Multivariate Distributions

The multivariate Gaussian

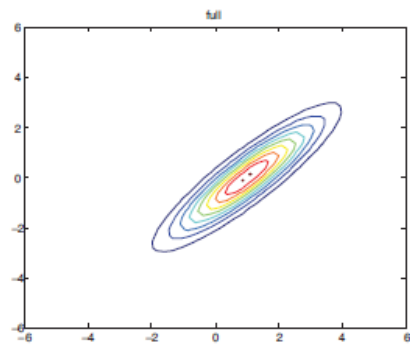
The **multivariate Gaussian** or **multivariate normal** (MVN) is the most widely used joint probability density function for continuous variables. We discuss MVNs in detail in Chapter 4; here we just give some definitions and plots.

The pdf of the MVN in D dimensions is defined by the following:

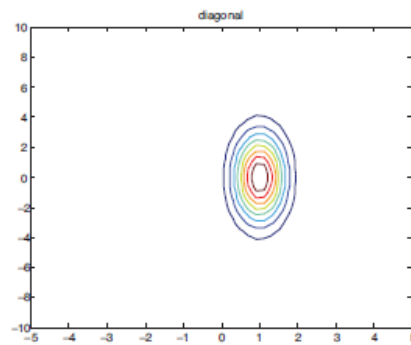
$$\mathcal{N}(\mathbf{x}|\boldsymbol{\mu}, \boldsymbol{\Sigma}) \triangleq \frac{1}{(2\pi)^{D/2}|\boldsymbol{\Sigma}|^{1/2}} \exp \left[-\frac{1}{2}(\mathbf{x} - \boldsymbol{\mu})^T \boldsymbol{\Sigma}^{-1}(\mathbf{x} - \boldsymbol{\mu}) \right] \quad (2.70)$$

where $\boldsymbol{\mu} = \mathbb{E}[\mathbf{x}] \in \mathbb{R}^D$ is the mean vector, and $\boldsymbol{\Sigma} = \text{cov}[\mathbf{x}]$ is the $D \times D$ covariance matrix. Sometimes we will work in terms of the **precision matrix** or **concentration matrix** instead. This is just the inverse covariance matrix, $\boldsymbol{\Lambda} = \boldsymbol{\Sigma}^{-1}$. The normalization constant $(2\pi)^{-D/2}|\boldsymbol{\Lambda}|^{1/2}$ just ensures that the pdf integrates to 1 (see Exercise 4.5).

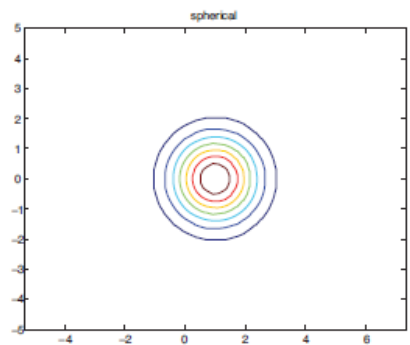
Figure 2.13 plots some MVN densities in 2d for three different kinds of covariance matrices. A full covariance matrix has $D(D+1)/2$ parameters (we divide by 2 since $\boldsymbol{\Sigma}$ is symmetric). A diagonal covariance matrix has D parameters, and has 0s in the off-diagonal terms. A **spherical** or **isotropic** covariance, $\boldsymbol{\Sigma} = \sigma^2 \mathbf{I}_D$, has one free parameter.



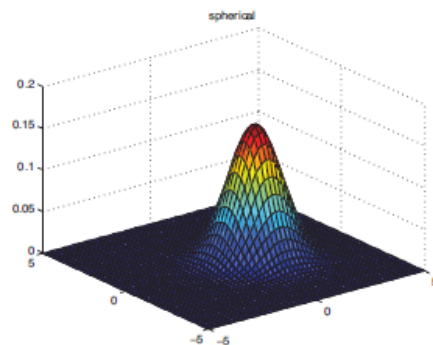
(a)



(b)



(c)



(d)

Figure 2.13 We show the level sets for 2d Gaussians.

(a) A full covariance matrix has elliptical contours.

(b) A diagonal covariance matrix is an **axis aligned** ellipse.

(c) A spherical covariance matrix has a circular shape.

(d) Surface plot for the spherical Gaussian in (c).

Multivariate Student t distribution

A more robust alternative to the MVN is the **multivariate Student t** distribution, whose pdf is given by

$$\mathcal{T}(\mathbf{x}|\boldsymbol{\mu}, \boldsymbol{\Sigma}, \nu) = \frac{\Gamma(\nu/2 + D/2)}{\Gamma(\nu/2)} \frac{|\boldsymbol{\Sigma}|^{-1/2}}{\nu^{D/2} \pi^{D/2}} \times \left[1 + \frac{1}{\nu} (\mathbf{x} - \boldsymbol{\mu})^T \boldsymbol{\Sigma}^{-1} (\mathbf{x} - \boldsymbol{\mu}) \right]^{-\left(\frac{\nu+D}{2}\right)} \quad (2.71)$$

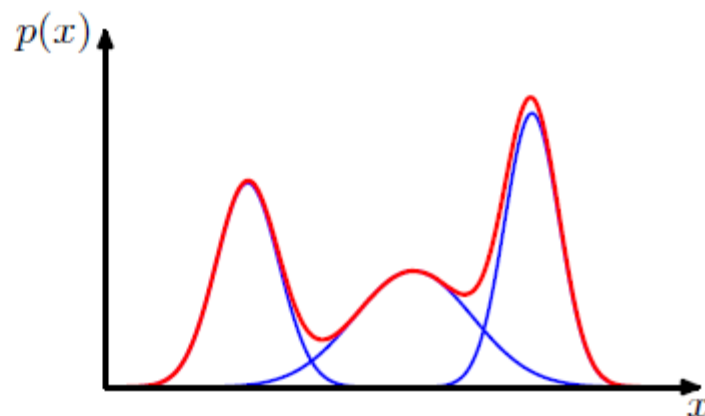
$$= \frac{\Gamma(\nu/2 + D/2)}{\Gamma(\nu/2)} |\pi \mathbf{V}|^{-1/2} \times [1 + (\mathbf{x} - \boldsymbol{\mu})^T \mathbf{V}^{-1} (\mathbf{x} - \boldsymbol{\mu})]^{-\left(\frac{\nu+D}{2}\right)} \quad (2.72)$$

where $\boldsymbol{\Sigma}$ is called the scale matrix (since it is not exactly the covariance matrix) and $\mathbf{V} = \nu \boldsymbol{\Sigma}$. This has fatter tails than a Gaussian. The smaller ν is, the fatter the tails. As $\nu \rightarrow \infty$, the

Mixture Models

Mixtures of Gaussians

Example of a Gaussian mixture distribution in one dimension showing three Gaussians (each scaled by a coefficient) in blue and their sum in red.



$$p(\mathbf{x}) = \sum_{k=1}^K \pi_k \mathcal{N}(\mathbf{x} | \boldsymbol{\mu}_k, \boldsymbol{\Sigma}_k) \quad \sum_{k=1}^K \pi_k = 1. \quad 0 \leq \pi_k \leq 1.$$

Transformation of Random Variables

Transformation

Linear

$$\mathbf{y} = f(\mathbf{x}) = \mathbf{Ax} + \mathbf{b}$$



$$\mathbb{E}[\mathbf{y}] = \mathbb{E}[\mathbf{Ax} + \mathbf{b}] = \mathbf{A}\boldsymbol{\mu} + \mathbf{b}$$

$$\text{cov}[\mathbf{y}] = \text{cov}[\mathbf{Ax} + \mathbf{b}] = \mathbf{A}\boldsymbol{\Sigma}\mathbf{A}^T$$

General

$$p_y(y) = \sum_{x: f(x)=y} p_x(x)$$

$$p_y(\mathbf{y}) = p_x(\mathbf{x}) \left| \det \left(\frac{\partial \mathbf{x}}{\partial \mathbf{y}} \right) \right|$$

Example

For example, suppose $X \sim U(-1, 1)$, and $Y = X^2$. Then $p_y(y) = \frac{1}{2}y^{-\frac{1}{2}}$.

As a simple example, consider transforming a density from Cartesian coordinates $\mathbf{x} = (x_1, x_2)$ to polar coordinates $\mathbf{y} = (r, \theta)$, where $x_1 = r \cos \theta$ and $x_2 = r \sin \theta$. Then

$$\mathbf{J}_{\mathbf{y} \rightarrow \mathbf{x}} = \begin{pmatrix} \frac{\partial x_1}{\partial r} & \frac{\partial x_1}{\partial \theta} \\ \frac{\partial x_2}{\partial r} & \frac{\partial x_2}{\partial \theta} \end{pmatrix} = \begin{pmatrix} \cos \theta & -r \sin \theta \\ \sin \theta & r \cos \theta \end{pmatrix}$$

and

$$|\det \mathbf{J}| = |r \cos^2 \theta + r \sin^2 \theta| = |r|$$

Hence

$$\begin{aligned} p_{\mathbf{y}}(\mathbf{y}) &= p_{\mathbf{x}}(\mathbf{x}) |\det \mathbf{J}| \\ p_{r, \theta}(r, \theta) &= p_{x_1, x_2}(x_1, x_2) r = p_{x_1, x_2}(r \cos \theta, r \sin \theta) r \end{aligned}$$

Central limit theorem

Now consider N random variables with pdf's (not necessarily Gaussian) $p(x_i)$, each with mean μ and variance σ^2 . We assume each variable is **independent and identically distributed** or **iid** for short. Let $S_N = \sum_{i=1}^N X_i$ be the sum of the rv's. This is a simple but widely used transformation of rv's. One can show that, as N increases, the distribution of this sum approaches

$$p(S_N = s) = \frac{1}{\sqrt{2\pi N\sigma^2}} \exp\left(-\frac{(s - N\mu)^2}{2N\sigma^2}\right) \quad (2.96)$$

Hence the distribution of the quantity

$$Z_N \triangleq \frac{S_N - N\mu}{\sigma\sqrt{N}} = \frac{\bar{X} - \mu}{\sigma/\sqrt{N}} \quad (2.97)$$

converges to the standard normal, where $\bar{X} = \frac{1}{N} \sum_{i=1}^N x_i$ is the sample mean. This is called the **central limit theorem**. See e.g., (Jaynes 2003, p222) or (Rice 1995, p169) for a proof.

In Figure 2.17 we give an example in which we compute the mean of rv's drawn from a beta distribution. We see that the sampling distribution of the mean value rapidly converges to a Gaussian distribution.

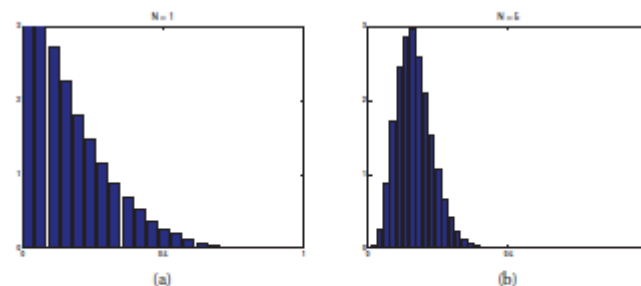


Figure 2.17 The central limit theorem in pictures. We plot a histogram of $\frac{1}{N} \sum_{i=1}^N x_{ij}$, where $x_{ij} \sim \text{Beta}(1,5)$, for $j = 1 : 10000$. As $N \rightarrow \infty$, the distribution tends towards a Gaussian. (a) $N = 1$. (b) $N = 5$. Based on Figure 2.6 of (Bishop 2006a). Figure generated by `centralLimitDemo`.

Parameter Estimation

Parameter estimations

- Maximum likelihood
- Maximum a posteriori
- Expected Maximization

ML, MAP, EM

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

xi: samples

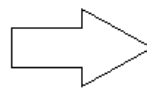
X=(X1,...,Xn): iid

MAP: Prior & posterior

EM: usually Mixture models

$$\theta_{MLE} = \arg \max_{\theta} P(X|\theta) = \arg \max_{\theta} \prod_i P(x_i|\theta)$$

$$P(\theta|X) = \frac{P(X|\theta)P(\theta)}{P(X)} \\ \propto P(X|\theta)P(\theta)$$



$$\theta_{MAP} = \arg \max_{\theta} P(X|\theta)P(\theta) \\ = \arg \max_{\theta} \log P(X|\theta)P(\theta)$$

$\mathbb{X} = \{\mathbf{x}^{(1)}, \dots, \mathbf{x}^{(m)}\}$ Samples drawn **independently** from the true but **Unknown** data generating distributions $p_{\text{data}}(\mathbf{x})$

$p_{\text{model}}(\mathbf{x}; \boldsymbol{\theta})$: Parametric Model

$$p_{\text{model}}(\mathbb{X}; \boldsymbol{\theta}) = \prod_{i=1}^m p_{\text{model}}(\mathbf{x}^{(i)}; \boldsymbol{\theta})$$

Maximum Likelihood

$$\begin{aligned}\boldsymbol{\theta}_{\text{ML}} &= \arg \max_{\boldsymbol{\theta}} p_{\text{model}}(\mathbb{X}; \boldsymbol{\theta}) \\ &= \arg \max_{\boldsymbol{\theta}} \prod_{i=1}^m p_{\text{model}}(\mathbf{x}^{(i)}; \boldsymbol{\theta})\end{aligned}$$



$$\boldsymbol{\theta}_{\text{ML}} = \arg \max_{\boldsymbol{\theta}} \sum_{i=1}^m \log p_{\text{model}}(\mathbf{x}^{(i)}; \boldsymbol{\theta})$$

Maximum likelihood

$$\mathcal{X} = \{x^t\}_{t=1}^N$$

We assume that

$x^t \sim p(x \theta)$ x^t are independent,
--

We want to find θ that makes sampling x^t from $p(x|\theta)$ as likely as possible

$$\underset{\theta}{\text{Max}} \quad l(\theta|\mathcal{X}) \equiv p(\mathcal{X}|\theta) = \prod_{t=1}^N p(x^t|\theta)$$

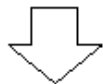
$$\underset{\theta}{\text{Max}} \quad \mathcal{L}(\theta|\mathcal{X}) \equiv \log l(\theta|\mathcal{X}) = \sum_{t=1}^N \log p(x^t|\theta)$$

ML for Bernoulli

$$P(x) = p^x(1-p)^{1-x}, x \in \{0, 1\}$$

iid sample $\mathcal{X} = \{x^t\}_{t=1}^N$, where $x^t \in \{0, 1\}$

$$\begin{aligned}\mathcal{L}(p|\mathcal{X}) &= \log \prod_{t=1}^N p^{(x^t)}(1-p)^{(1-x^t)} \\ &= \sum_t x^t \log p + \left(N - \sum_t x^t\right) \log(1-p)\end{aligned}$$



$$\hat{p} = \frac{\sum_t x^t}{N}$$

Gaussian (Normal) Density

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

$$\mathcal{X} = \{x^t\}_{t=1}^N \text{ with } x^t \sim \mathcal{N}(\mu, \sigma^2)$$

$$\mathcal{L}(\mu, \sigma | \mathcal{X}) = -\frac{N}{2} \log(2\pi) - N \log \sigma - \frac{\sum_t (x^t - \mu)^2}{2\sigma^2}$$

$$\begin{aligned} m &= \frac{\sum_t x^t}{N} \\ s^2 &= \frac{\sum_t (x^t - m)^2}{N} \end{aligned}$$