

# Probability & Statistics

Machine learning, 2021

Mansoor Rezghi

Department of Computer science, TMU

# References

- K. Murphy Machine Learning: A Probabilistic Perspective, MIT Press, 2012.(KM)-chap2
- C. Bishop, Pattern Recognition and Machine Learning, Springer, 2006.(CB)-chap2

# Probability: Discrete random variable

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

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

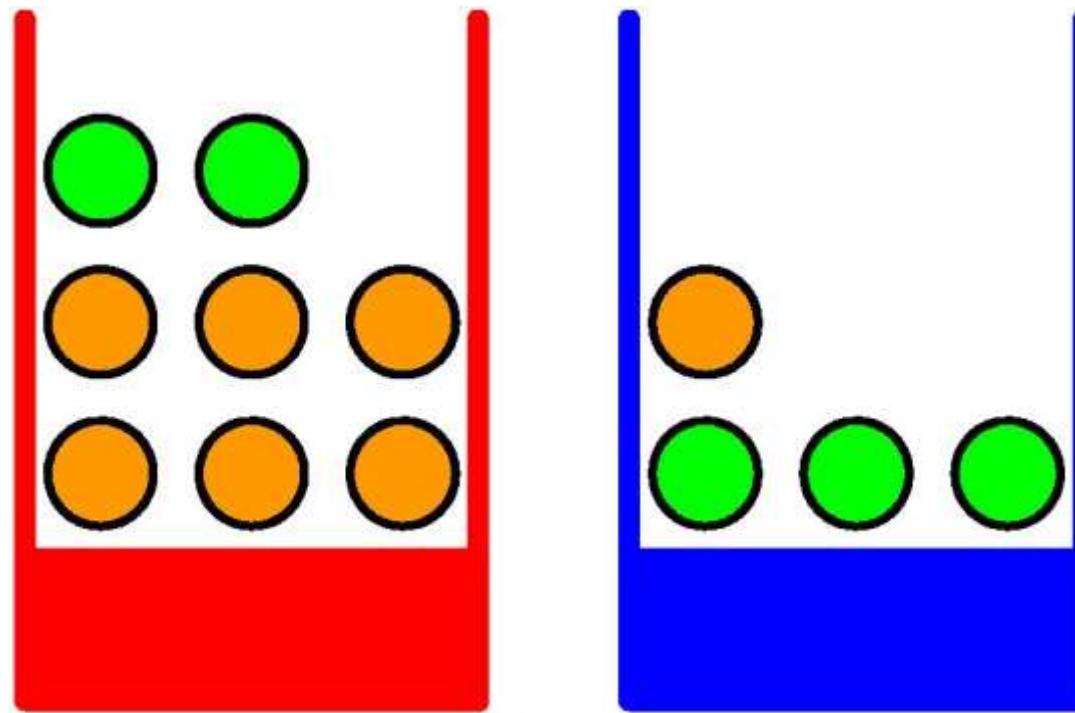
Conditional Probability  $p(A|B) = \frac{p(A, B)}{p(B)}$  if  $p(B) > 0$

Bayes rule

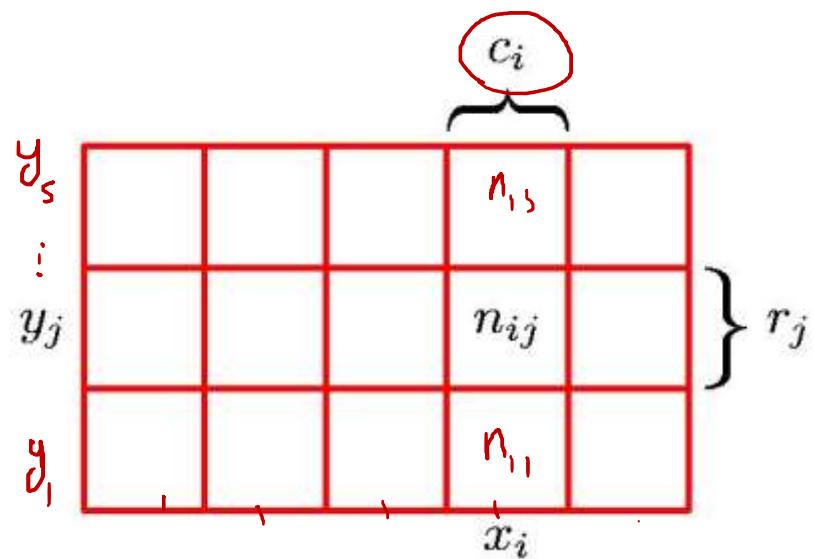
$$p(X = x|Y = y) = \frac{p(X = x, Y = y)}{p(Y = y)} = \frac{p(X = x)p(Y = y|X = x)}{\sum_{x'} p(X = x')p(Y = y|X = x')}$$

# Probability Theory

Apples and Oranges



# Probability Theory



Joint Probability

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

$$X = \{x_1, \dots, x_n\} \quad n_X = N$$
$$Y = \{y_1, \dots, y_m\} \quad n_Y = M$$

$$X = x_i, \quad P(X = x_i, Y = y_j) =$$

$$Y = y_j$$

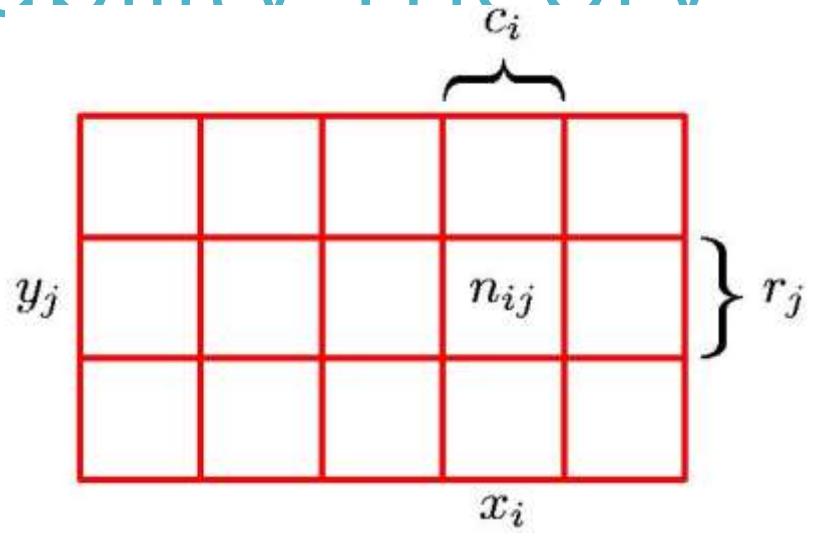
Marginal Probability

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

Conditional Probability

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

# Probability Theory



Sum Rule

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

Product Rule

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

$$p(x, y) = p(y|x) p(x) \quad \checkmark$$

# The Rules of Probability

Sum Rule

$$p(X) = \sum_{Y \in \mathcal{Y}} p(X, Y)$$

Product Rule

$$p(X, Y) = p(Y|X)p(X)$$
 ✓

# Bayes Theorem

$$P(X|Y)$$

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

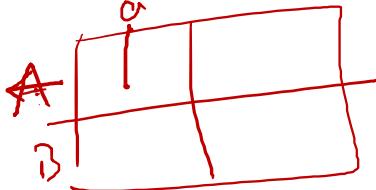
$$\underbrace{p(X)}_{Y} = \sum_Y \underbrace{p(X|Y)p(Y)}_{Y}$$

posterior  $\propto$  likelihood  $\times$  prior

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

11

# Bayes rule: Example



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

$$p(x=1|y=1) = 0.8 \quad \checkmark$$

If  $p(y=1) = 0.004$

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

$$P(A|B)$$

$$P(B|A)$$

...  
Be  
Y = 0  
Y = 1

X = 0  
X = 1

→  $P(X=1|Y=1)$

$$P(x=1|y=1) \quad \checkmark$$

$$P(y=1|x=1) \quad ?$$



$$p(y=1|x=1) = \frac{p(x=1|y=1)p(y=1)}{p(x=1|y=1)p(y=1) + p(x=1|y=0)p(y=0)}$$

$$= \frac{0.8 \times 0.004}{0.8 \times 0.004 + 0.1 \times 0.996} = 0.031$$

$$P(y=0|x=0) \quad ? \quad ?$$

# 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.,

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

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

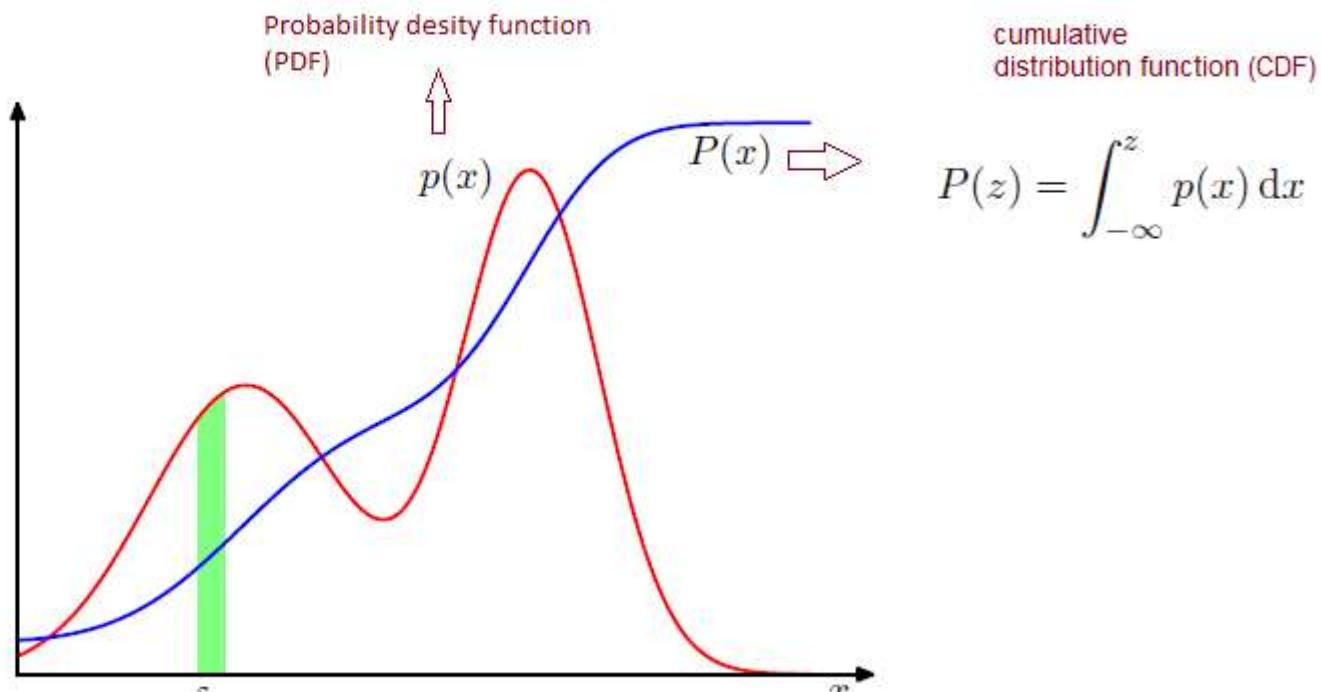
conditionally independent

$$\frac{P(X)}{P(X|Y)}$$

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

# 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  $P(x)$ .



$$P(u) = \left\{ \frac{1}{6}, \frac{1}{6} \right\}$$

# Important concepts

$\chi$   
 $f(n)$

18 18 19 19 18

$x = 18, 18, 19$

~~18 18~~

$$\frac{2 \times 18 + 18 + 2 \times 19}{5}$$

$$= 18 \times \frac{2}{5} + \frac{1}{5} \times 18 + 19 \times \frac{2}{5}$$

$$P(x) = \frac{1}{5}, \frac{1}{5}, \frac{2}{5}$$

$$E(x) = \sum x P(x)$$

$$E(x) = \sum x p(x)$$

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

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

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

$$\begin{array}{r} 17 \\ 18 \\ \hline 17 \end{array} \quad \begin{array}{r} 17 \\ 18 \\ \hline 17 \end{array} \quad \begin{array}{r} 17 \\ 18 \\ \hline 19 \end{array}$$

Variance

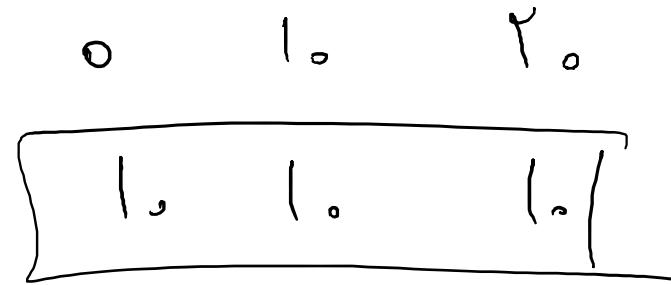
$$\text{var}[f] = \mathbb{E} \left[ (f(x) - \mathbb{E}[f(x)])^2 \right] = \mathbb{E}[f(x)^2] - \mathbb{E}[f(x)]^2$$

$$\text{Var}(x) = E \left( (x - E(x))^2 \right)$$

$$x \rightarrow p(n) \rightarrow E(x) = \mu$$

$$(x - \mu)^r \rightarrow P(n)$$

$$\begin{aligned} E((x - \mu)^r) &= \sum (x - \mu)^r p(n) \\ &= \sum (x - E(x))^r P(x) \end{aligned}$$



$P(\cdot)$

# Quintile

$$P(X \leq x_\alpha) = \alpha$$

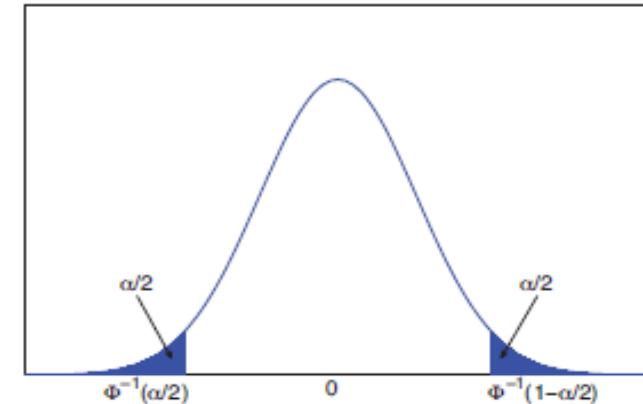
$x_\alpha$        $x_{.25}$        $x_{.5}$        $x_{.75}$

$$P(X \leq 4) = \sum_{m < 4} P(x)$$

$\alpha$  quantile of  $F = F^{-1}(\alpha)$  is the value of  $x_\alpha$  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

$$P(X=1) = \theta$$

$$P(X=0) = 1 - \theta$$

$$P(X=x) = \theta^x (1-\theta)^{1-x}$$

## 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  $\theta$ , then we say  $X$  has a **binomial** distribution, written as  $X \sim \text{Bin}(n, \theta)$ . The pmf is given by

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

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

$$E(X) = \sum n p(n)$$

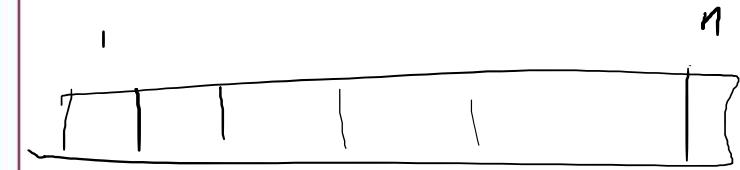
$$P(n)$$

## The Bernoulli distribution

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

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

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



$$k = 1$$

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

$$E(n) = 1 \times \theta + 0 \times (1-\theta) = \theta$$

$$E(n) = \sum n_i p(n_i)$$

$$X_1 = \begin{pmatrix} 1 \\ 0 \\ 0 \\ 0 \\ 0 \\ 0 \end{pmatrix} \quad X_2 = \begin{pmatrix} 0 \\ 1 \\ 0 \\ 0 \\ 0 \\ 0 \end{pmatrix}$$

$$X = \begin{pmatrix} n_1 \\ n_2 \\ n_3 \\ n_4 \\ n_5 \\ n_6 \end{pmatrix}$$

$n \leq n_i \leq 8$   
 $\sum n_i = n$

### The multinomial distribution

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

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

### 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

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

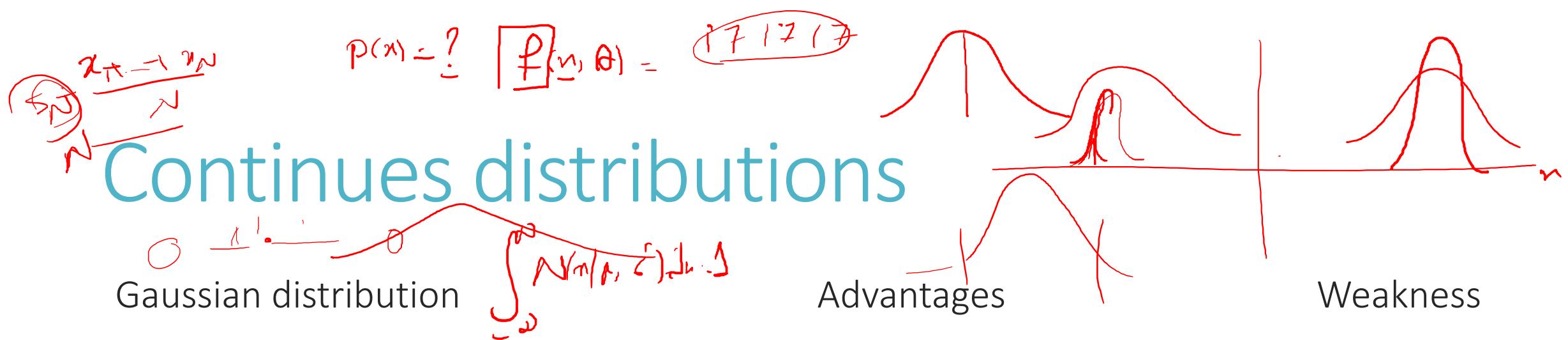
The first term is just the normalization constant, required to ensure the distribution sums to 1.

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

$$X = \begin{pmatrix} 4 \\ 9 \\ 0 \\ 0 \\ 0 \\ 7 \end{pmatrix}$$

$$\boldsymbol{\theta} = \begin{pmatrix} \theta_1 \\ \theta_2 \\ \theta_3 \\ \theta_4 \end{pmatrix}$$

$$\frac{n!}{x_1! \dots x_4!} \cdot \frac{1}{6} \cdot \frac{1}{6} \cdot \frac{1}{6}$$



$\mathcal{N}(x|\mu, \sigma^2)$

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

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

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

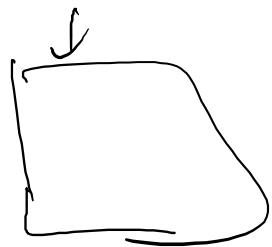
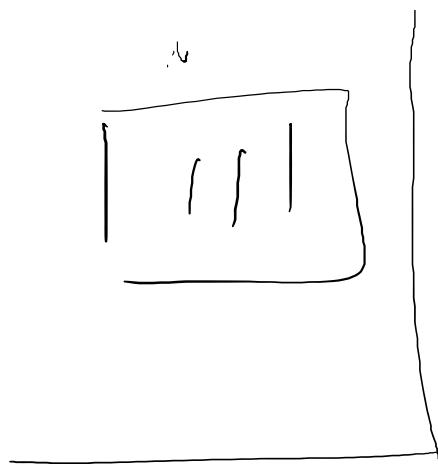
It has two parameters which are easy to interpret  
mean and variance.

$S_N = \sum x_1 + \dots + x_N$

Central limit Theorem  
sums of independent random variables have an approximately Gaussian distribution.

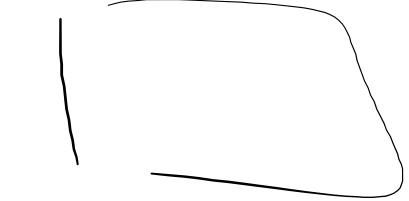
Maximum Entropy  
Gaussian distribution makes the least number of assumptions

Sensitive to outlier



-

-

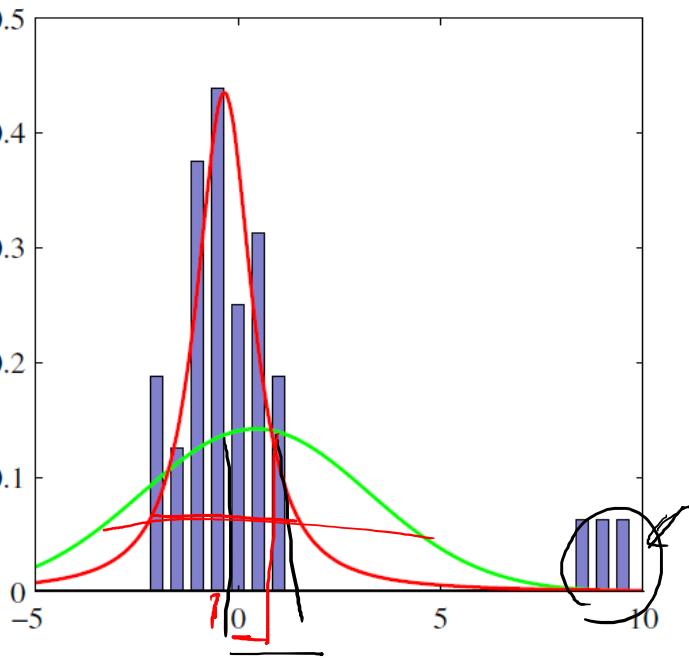
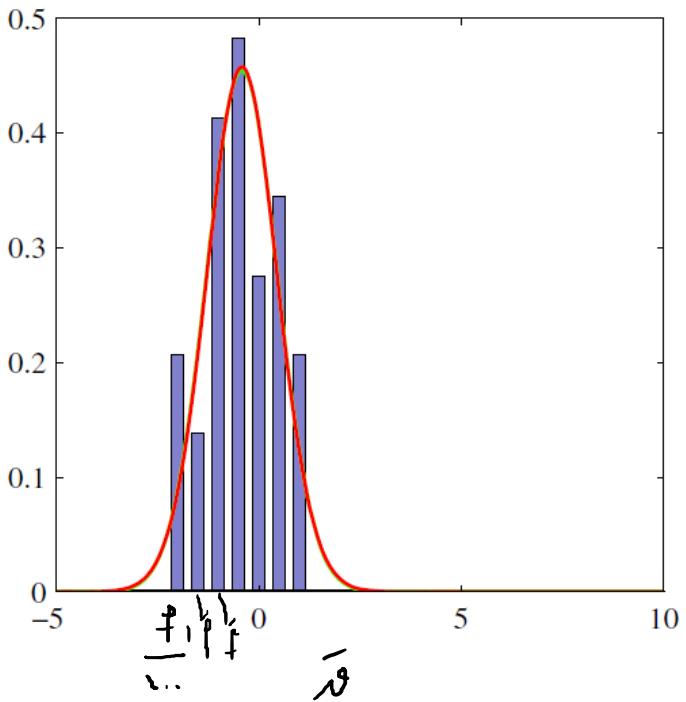


# Gaussian and outlier

17 D

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  $\mu$  and precision  $\lambda$

Able to handle out layer

# Laplace distribution



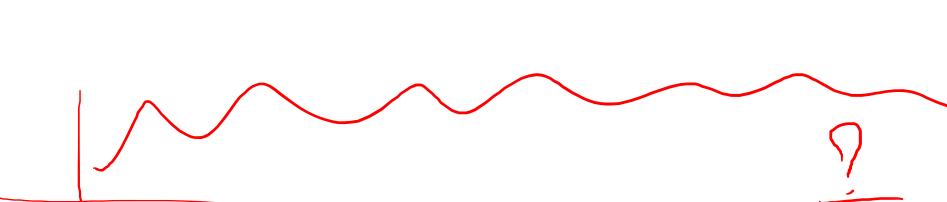
$$\text{Lap}(x|\mu, b) \triangleq \frac{1}{2b} \exp\left(-\frac{|x - \mu|}{b}\right)$$

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

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

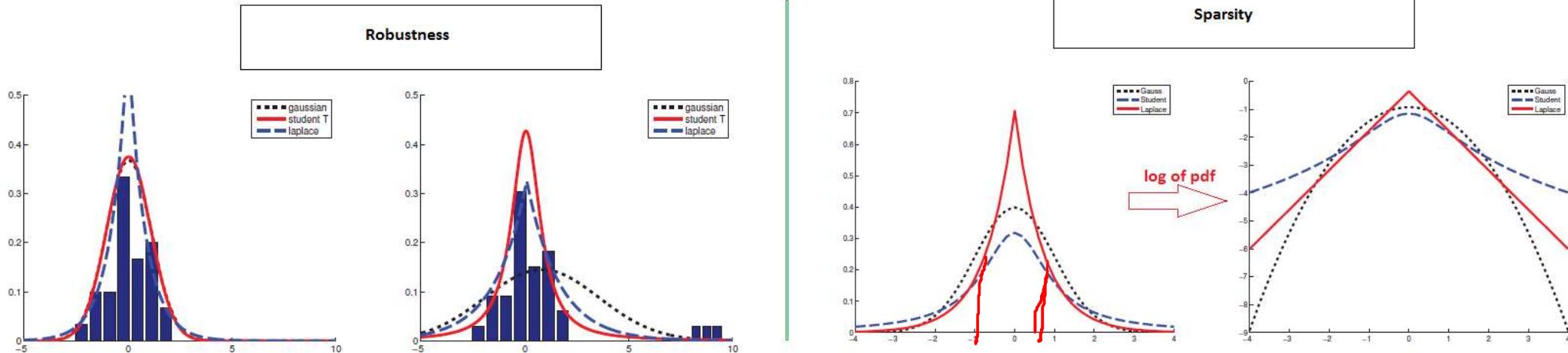
**Robust undr outlier**

**Sparse distribution**



# Laplace : Robustness and sparsity





Laplace distribution, which is always log-concave

$$x = \begin{pmatrix} x_1 \\ \vdots \\ x_n \end{pmatrix} \quad x \in R$$

$$x \in \mathbb{R}^d \quad x \in R^d$$

$P(X=x)$

$P(\mathbf{X}=x) =$

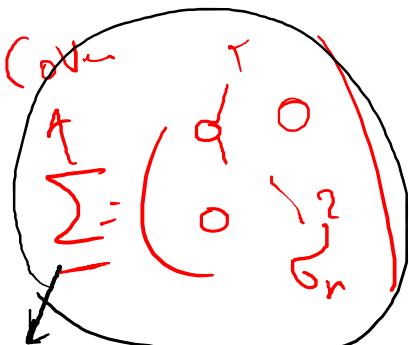
$$x = \begin{pmatrix} x_1 \\ \vdots \\ x_d \end{pmatrix} \quad P(X) = P(x_1, x_2, \dots, x_n) = \prod_{i=1}^n P(x_i)$$

$$x_i \sim N(n | \mu_i, \sigma_i^2)$$

$$N(x | \mu_i, \sigma_i^2) = \frac{1}{\sqrt{2\pi\sigma_i^2}} e^{-\frac{(x_i - \mu_i)^2}{2\sigma_i^2}}$$

$$= \frac{1}{(2\pi)^{\frac{d}{2}} \prod_{i=1}^d \sigma_i^2} e^{-(x - \mu)^T \Sigma^{-1} (x - \mu)}$$

$$= \frac{1}{(2\pi)^{\frac{d}{2}} |\Sigma|^{\frac{1}{2}}} e$$



$$P(X) = \prod_{i=1}^d N(x_i | \mu_i, \sigma_i^2) ?$$

$$= \left( \prod_{i=1}^d \frac{1}{\sqrt{2\pi\sigma_i^2}} \right) e^{\sum_i -\frac{(x_i - \mu_i)^2}{2\sigma_i^2}}$$

$$x = \begin{pmatrix} x_1 \\ \vdots \\ x_n \end{pmatrix} \quad \mu = \begin{pmatrix} \mu_1 \\ \vdots \\ \mu_n \end{pmatrix}$$

## Multivariate

$$X \perp Y \quad \text{or} \quad X \perp\!\!\! \perp Y \quad P(X, Y) = P(X)P(Y)$$

$P(X, Y)$

# Multivariate

$$E[n] = \sum n P(n)$$

$$X = \begin{pmatrix} X_1 \\ \vdots \\ X_n \end{pmatrix}$$

$$\text{Cov}(X) = (\text{Cov}(X_i, X_j))_{1 \leq i, j \leq n}$$

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.

$$\underline{P(X, Y)} =$$

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

$$X = \begin{pmatrix} X_1 \\ \vdots \\ X_d \end{pmatrix}$$

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

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

# Correlation

$$\text{Cov}(X, Y) = \frac{1}{\sqrt{\sigma_x^2 \sigma_y^2}} E \left( \frac{(X - \mu_x)(Y - \mu_y)}{\sqrt{\sigma_x^2 \sigma_y^2}} \right)$$

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

➡

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

A **correlation matrix** has the form

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

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

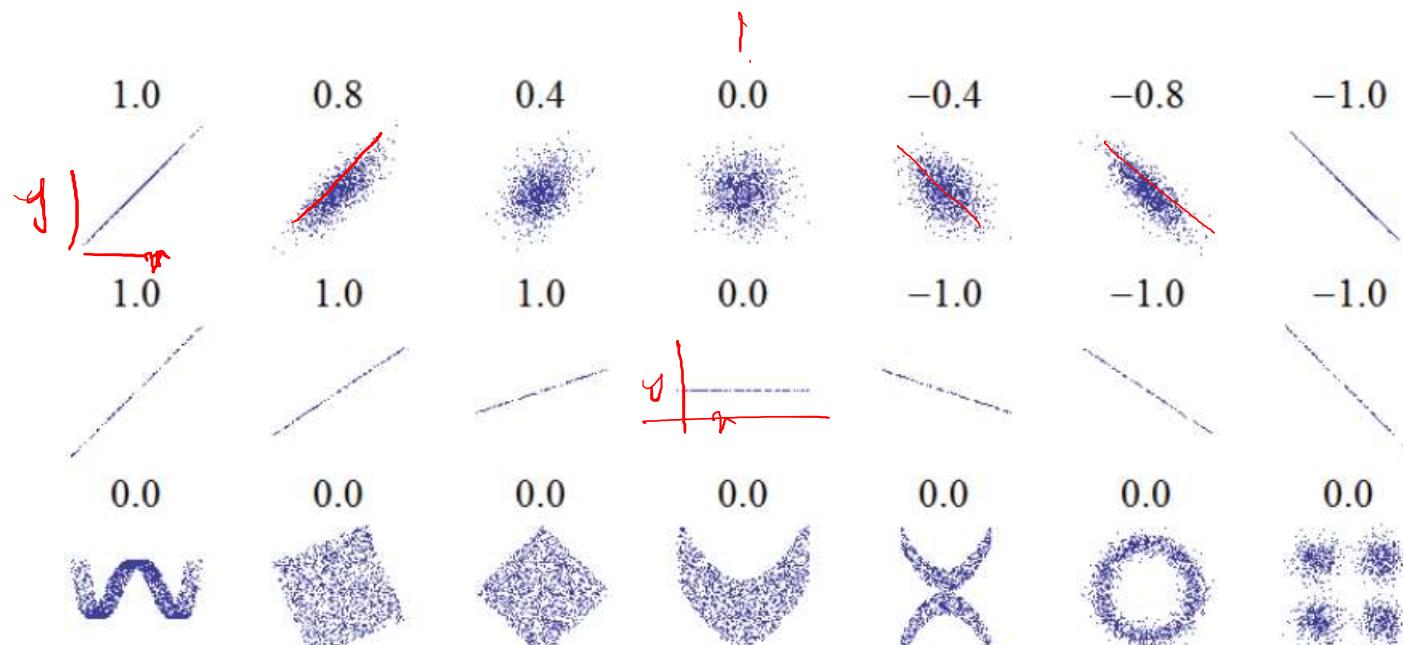
$$\text{corr}[X, Y] = 1 \Leftrightarrow Y = aX + b \quad a < 0$$

$$\text{Corr}(X, Y) = \frac{\text{Cov}(X, Y)}{\sqrt{\text{Var}(X)} \sqrt{\text{Var}(Y)}} = \frac{\text{Cov}(X, Y)}{\sqrt{\text{Var}(X)}} \frac{\text{Cov}(Y, Y)}{\sqrt{\text{Var}(Y)}} = \frac{\text{Cov}(X, Y)}{\sqrt{\text{Var}(X)}} = \frac{\sum_{i=1}^n (X - \mu_X)(Y - \mu_Y)P(X_i, Y_i)}{\sqrt{\text{Var}(X)}}$$

$-1 \leq \text{Corr}(X, Y) \leq 1$

$\text{Corr}(X, Y) = 1 \quad Y = aX + b \quad a > 0$

$\text{Corr}(X, Y) = -1 \quad Y = aX + b \quad a < 0$



**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](http://en.wikipedia.org/wiki/File:Correlation_examples.png)

$$\text{Cov}(X, Y) = \sum_{x,y} (x - \mu_x)(y - \mu_y) p(x, y) / P(X)P(Y) = \left[ \sum_n (n - \mu_x) p(n) \right] \left[ \sum_n (y - \mu_y) p(y) \right]$$

## Dependence Vs Correlation

$$\frac{\sum n p(n) - \mu_x \sum p(n)}{\mu_x - \mu_n} = 0$$

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



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

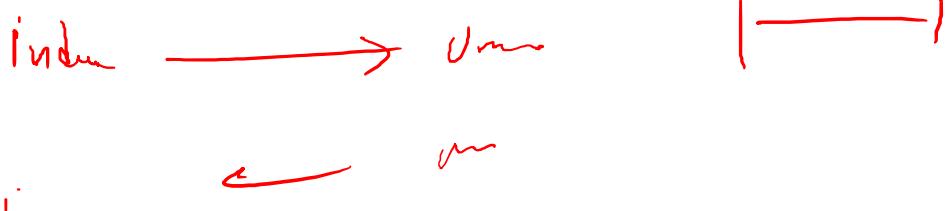
so they are uncorrelated.

**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.*



$$X = \begin{pmatrix} x_1 \\ x_2 \\ \vdots \\ x_d \end{pmatrix}$$

# Multivariate Distributions

## The multivariate Gaussian

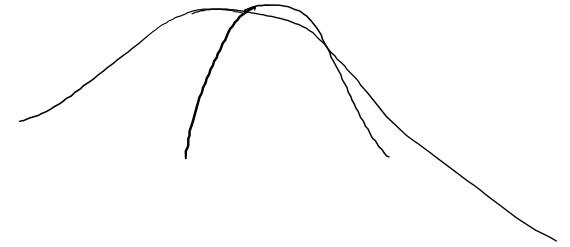
The **multivariate Gaussian** or **multivariate normal (MVN)** is the most widely used joint probability density function for continuous variables. We discuss MVNs in detail in Chapter 4; here we just give some definitions and plots.

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).

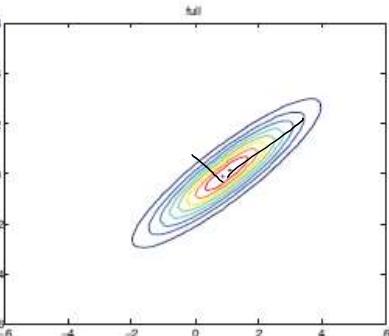
Figure 2.13 plots some MVN densities in 2d for three different kinds of covariance matrices. A full covariance matrix has  $D(D + 1)/2$  parameters (we divide by 2 since  $\boldsymbol{\Sigma}$  is symmetric). A diagonal covariance matrix has  $D$  parameters, and has 0s in the off-diagonal terms. A **spherical**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



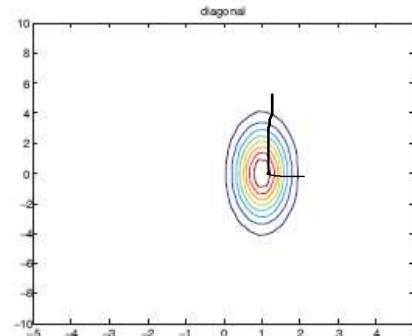
$$\Sigma \in \mathbb{R}^{n \times n}$$

$$\Sigma = \begin{pmatrix} \sigma_1 & 0 \\ 0 & \sigma_2 \end{pmatrix}$$

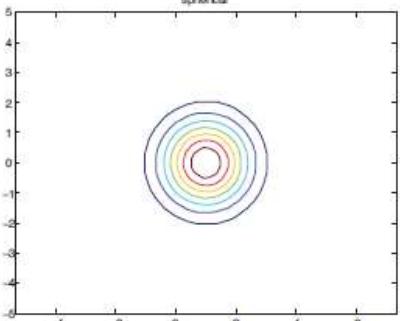
$$\Sigma = \sigma^2 I$$



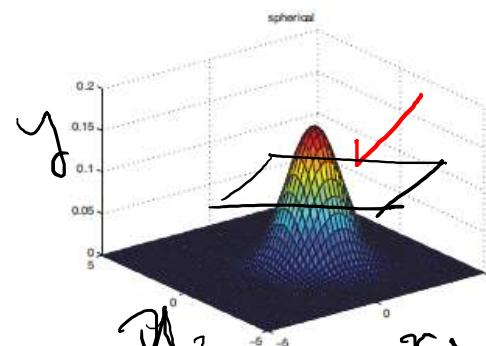
(a)



(b)



(c)



(d)

~~ax<sub>1</sub><sup>2</sup> + bx<sub>2</sub><sup>2</sup>~~ +  $\boxed{c x_1 x_2} = ?$

$$y = f(x_1, x_2)$$

$$f(x_1, x_2) = 2$$

$$\begin{pmatrix} x_1 \\ x_2 \end{pmatrix}$$

$$(x_1, x_2)$$

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).

$$\begin{aligned}
 x \in \mathbb{R}^2 & \quad x = \begin{pmatrix} x_1 \\ x_2 \end{pmatrix} \\
 & - (x_1 - \mu_1)^T \begin{pmatrix} \alpha & \beta \\ \gamma & \varepsilon \end{pmatrix} \begin{pmatrix} x_1 - \mu_1 \\ x_2 - \mu_2 \end{pmatrix} = f' \\
 \text{or } e & - \left[ \begin{pmatrix} x_1 - \mu_1 & x_2 - \mu_2 \end{pmatrix} \begin{pmatrix} \alpha & \beta \\ \gamma & \varepsilon \end{pmatrix} \begin{pmatrix} x_1 - \mu_1 \\ x_2 - \mu_2 \end{pmatrix} \right] = f'' 
 \end{aligned}$$

### Multivariate Student *t* distribution

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

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

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

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  $\nu$  is, the fatter the tails. As  $\nu \rightarrow \infty$ , the

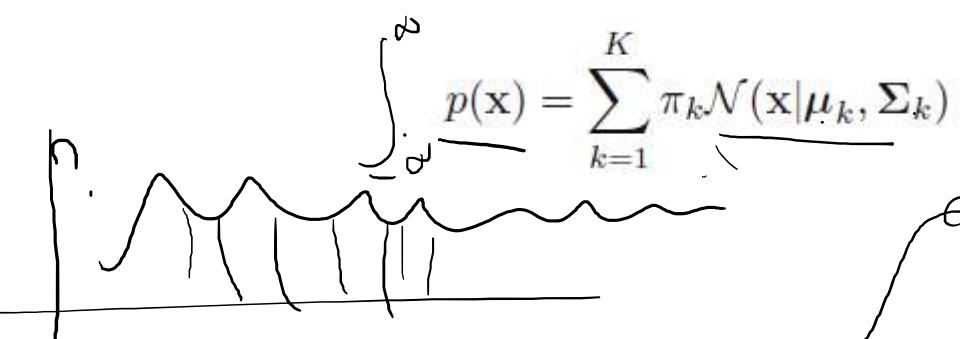
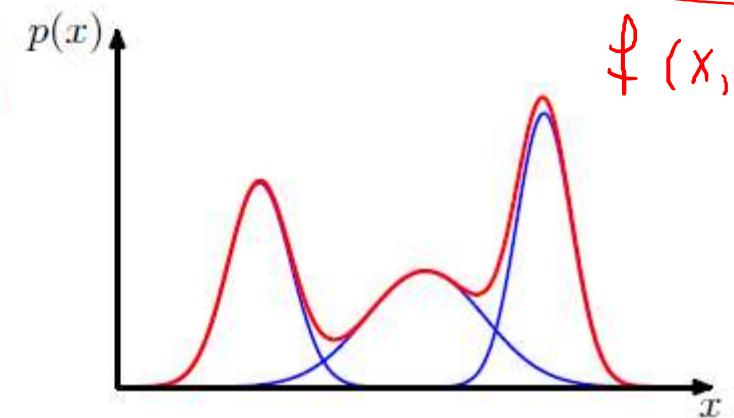
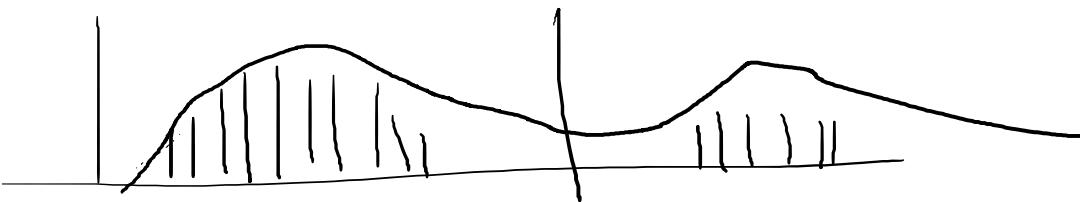
# Mixture Models

$$P(x, \theta_1, \theta_2) = \pi_1 P(m, \theta_1) + \pi_2 f(x, \theta_2)$$

$\underbrace{\phantom{P(x, \theta_1, \theta_2) = \pi_1 P(m, \theta_1) + \pi_2 f(x, \theta_2)}}$   
 $\pi_1, \pi_2$

## 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.



$$\sum_{k=1}^K \pi_k = 1.$$

$\downarrow$



$$p(x) = \sum_{k=1}^K \pi_k \mathcal{N}(x | \mu_k, \Sigma_k)$$

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

$$\sum_{k=1}^K \pi_k = 1.$$

$0 \leq \pi_k \leq 1.$

$$1 - \int_{-\infty}^{\infty} f(x, \theta_1, \theta_2, \pi, z) dx$$

$$= \pi_1 \left\{ \int_{-\infty}^{\infty} f(m, \theta_1) dm + \frac{1}{\pi_1} \int_{-\infty}^{\infty} f(z, \theta_2) dz \right\}$$

# Transformation of Random Variables

# Transformation

Linear

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



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

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

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  $\mu$  and variance  $\sigma^2$ . 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

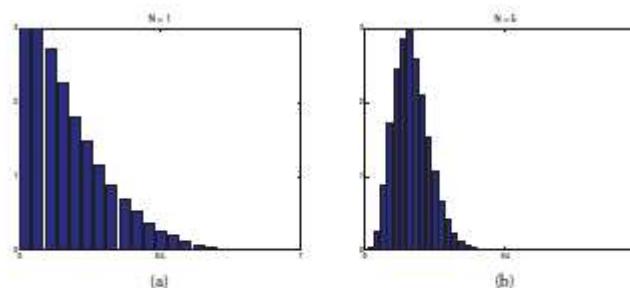
$$p(S_N = s) = \frac{1}{\sqrt{2\pi N\sigma^2}} \exp\left(-\frac{(s - N\mu)^2}{2N\sigma^2}\right) \quad (2.96)$$

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.

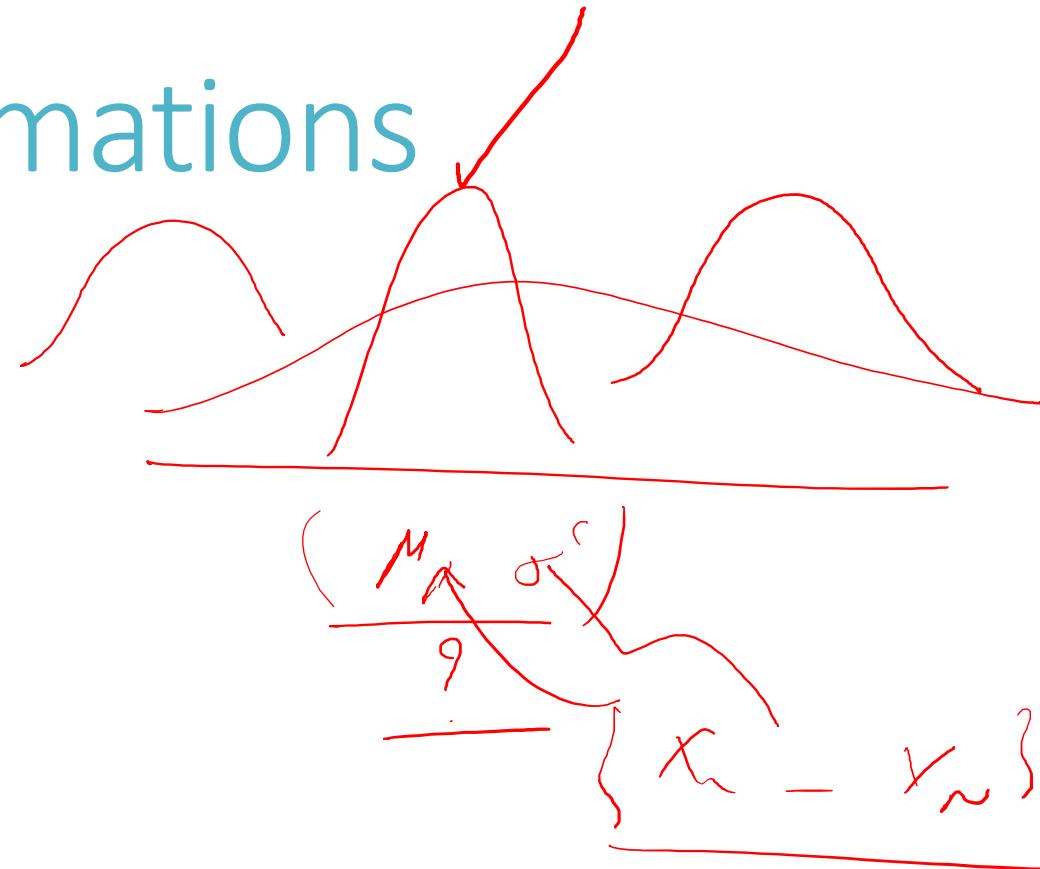


**Figure 2.17** The central limit theorem in pictures. We plot a histogram of  $\frac{1}{N} \sum_{i=1}^N x_{ij}$ , where  $x_{ij} \sim \text{Beta}(1, 5)$ , for  $j = 1 : 10000$ . As  $N \rightarrow \infty$ , the distribution tends towards a Gaussian. (a)  $N = 1$ . (b)  $N = 5$ . Based on Figure 2.6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

# Parameter Estimation

# Parameter estimations

- Maximum likelihood
- Maximum a posteriori
- Expected Maximization



# ML,MAP,EM

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

$x_i$ : samples

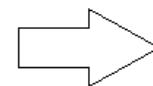
$X=(X_1, \dots, X_n)$ : iid

MAP: Prior & posterior

EM: usually Mixture models

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

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



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

$p_{\text{data}}(\mathbf{x})$

Unknown

$\mathbb{X} = \{\mathbf{x}^{(1)}, \dots, \mathbf{x}^{(m)}\}$  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

$$\boldsymbol{\theta}_{\text{ML}} = \arg \max_{\boldsymbol{\theta}} p_{\text{model}}(\mathbb{X}; \boldsymbol{\theta})$$

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



ML

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

# Maximum likelihood

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

We assume that

$$\begin{aligned} x^t &\sim p(x|\theta) \\ x^t &\text{ are independent.} \end{aligned}$$

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

---

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

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

~~param f(x)~~

### ML for Bernoulli

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

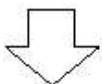
iid sample  $\mathcal{X} = \{x^t\}_{t=1}^N$ , where  $x^t \in \{0,1\}$

$$x^* = \underset{x}{\arg \max} f(x)$$

~~f: incres~~

$$x^* = \underset{x}{\arg \max} g(f(x))$$

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



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

## Gaussian (Normal) Density

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

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

---

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

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



## Shanon Entropy

Information  $\propto \frac{1}{\text{Possibility}}$

$$\text{Information} \propto \frac{1}{P(x)}$$

**Additivity:** For independent events  $Z_1$  and  $Z_2$

$$\mathbf{Inf}(Z_1 + Z_2) \propto \mathbf{Inf}(Z_1) + \mathbf{Inf}(Z_2)$$

$$I(X) = -\log P(x)$$

Covers two mwntioed properties

Shannon entropy  $H(x) = \mathbb{E}_{x \sim P}[I(x)] = -\mathbb{E}_{x \sim P}[\log P(x)].$

$$H(x) \propto -\frac{1}{P(n)}, \quad H(x,y) = H(n) + H(y) \leftarrow H(n) = \log \left( \frac{1}{P(x)} \right) = -\log P(x)$$

$x, y$  in data

The **entropy** of a random variable  $X$  with distribution  $p$ , denoted by  $H(X)$  or sometimes  $H(p)$ , is a measure of its uncertainty. In particular, for a discrete variable with  $K$  states, it is defined by

In order to co

$$H(X) \triangleq -\sum_{k=1}^K p(X=k) \log_2 p(X=k)$$

$$\log_2 0 = \infty$$

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

$$H[x] = -8 \times \frac{1}{8} \log_2 \frac{1}{8} = 3 \text{ bits.}$$

usage of length 3 bits.

$$E[R_{x_1}] = \sum p_{x_1} r_{x_1}$$

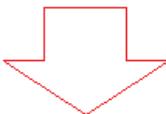
$$H(x) = E[-\log p(x)] = -\sum p(x) \log p(x) \quad \text{Entropy}$$

**Event**  $\{a, b, c, d, e, f, g, h\}$

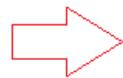


$$H[x] = -\frac{1}{2} \log_2 \frac{1}{2} - \frac{1}{4} \log_2 \frac{1}{4} - \frac{1}{8} \log_2 \frac{1}{8} - \frac{1}{16} \log_2 \frac{1}{16} - \frac{4}{64} \log_2 \frac{1}{64}$$

**Probability**  $\frac{1}{2}, \frac{1}{4}, \frac{1}{8}, \frac{1}{16}, \frac{1}{64}, \frac{1}{64}, \frac{1}{64}, \frac{1}{64}$



**Coding**  
0, 10, 110, 1110, 111100, 111101, 111110, 111111



$$\text{average code length} = \frac{1}{2} \times 1 + \frac{1}{4} \times 2 + \frac{1}{8} \times 3 + \frac{1}{16} \times 4 + 4 \times \frac{1}{64}$$

$$x = \{0, 1\}$$

$$P(x) = \frac{1}{2} \quad \frac{1}{2}$$

$$\cancel{(1 \log 0)} \cancel{x \frac{1}{2}} - \cancel{\log(1) \frac{1}{2}} = \frac{-1}{2} \times \cancel{\log 0} = \cancel{\frac{1}{2} \log 0} \approx \infty$$

$$x = \{0, 1\}$$

$$P(x) = 1 \quad 0$$

$$-\infty \cancel{x \log 1} + 1 \log 0 = -\cancel{\log 0} = \infty$$

$$y = \ln x$$

$$x = e^y$$

$$-\sum p(x_i) \log p(x_i)$$

$$-\left( 1 \times \cancel{\log 1} + 0 \times \cancel{\log 0} \right) = 0$$

$$= -\left[ \frac{1}{2} \times \log \frac{1}{2} + \frac{1}{2} \times \log \left( \frac{1}{2} \right) \right]$$

$$= -\log \frac{1}{2} = \log 2$$

$$\log_2 x_1$$

$$\begin{matrix} p(x) \\ q(x) \end{matrix}$$

$$D_{KL}(P(x), q(x))$$

$$-\int p(x) \ln p(x) dx$$

$$-\int p(x) \ln q(x) dx$$

Σ 7

## Kullback-Leibler (KL) divergence

$$q(x)$$

Coding by inexact dist



$$\begin{aligned} KL(p\|q) &= -\underbrace{\int p(x) \ln q(x) dx}_{\text{Coding by inexact dist}} - \left( -\int p(x) \ln p(x) dx \right) \\ &= -\int p(x) \ln \left\{ \frac{q(x)}{p(x)} \right\} dx. \end{aligned}$$



Coding by exact dist

p(x): True unknown distribution

$$\begin{matrix} q(x) \\ P(x) \end{matrix} \rightarrow \ln q(x)$$

q(x): An approximation of p(x)

Extra information for decoding

$$KL(p\|q) \neq KL(q\|p).$$

$$KL(p\|q) \geq 0$$

P

$$\begin{array}{cccc} x_1 & x_2 & x_3 & x_4 \\ \frac{1}{4} & \frac{1}{4} & \frac{1}{4} & \frac{1}{4} \end{array} \quad \left| \quad \begin{array}{cccc} Q & x_1 & x_2 & x_3 & x_4 \\ \frac{1}{4} & \frac{1}{8} & \frac{1}{8} & \frac{1}{8} & \frac{1}{8} \end{array} \right.$$

$$\begin{array}{cccc} x_1 & x_2 & x_3 & x_4 \\ \frac{1}{2} & \frac{1}{2} & 0 & \frac{1}{4} \end{array}$$

$$E(f_{\pi(x)}) = \sum f_{\pi(x)} P(x) - \log p_{\pi(x)}$$

Kullback-Leibler (KL) divergence: The similarity of two distributions

$$\checkmark D_{KL}(P||Q) = \mathbb{E}_{x \sim P} \left[ \log \frac{P(x)}{Q(x)} \right] = \underline{-H(P)}$$

$$D_{KL}(P||Q) \neq D_{KL}(Q||P)$$

cross-entropy

$$H(P, Q) = H(P) + D_{KL}(P||Q) = -\mathbb{E}_{x \sim P} \log Q(x)$$

$$E_{x \sim P} \left( -\log p_{\pi(x)} \right)$$

$$\operatorname{argmin}_Q H(P, Q) = \operatorname{argmin} D_{KL}(P||Q)$$

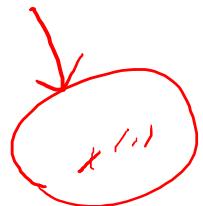
$$\boxed{- \int P(x) \log P(x) - \int P(x) \log Q(x)}$$

# #7 17 17 17 18

# ML and Cross Entropy Equivalence

$$\frac{17+17+18+19}{4} = \frac{2 \times 17 + 18 + 19}{4}$$

$$= 17 \times \frac{1}{2} + 18 \times \frac{1}{4} + 19 \times \frac{1}{4}$$



$$\begin{array}{c} 17 \ 17 \quad 18 \ 19 \\ \hline \frac{1}{2} \quad \frac{1}{4} \quad \frac{1}{4} \end{array}$$

$$\theta_{\text{ML}} = \arg \max_{\theta} \sum_{i=1}^m \log p_{\text{model}}(\mathbf{x}^{(i)}; \theta) \rightarrow \theta_{\text{ML}} = \arg \max_{\theta} \mathbb{E}_{\mathbf{x} \sim \hat{p}_{\text{data}}} [\log p_{\text{model}}(\mathbf{x}; \theta)] = \arg \min_{\theta} -\mathbb{E}_{\mathbf{x} \sim \hat{p}_{\text{data}}} [\log p_{\text{model}}(\mathbf{x})] = \arg \min_{\theta} H(P, Q) \text{ P.mle.}$$

$$D_{\text{KL}}(\hat{p}_{\text{data}} \| p_{\text{model}}) = \mathbb{E}_{\mathbf{x} \sim \hat{p}_{\text{data}}} [\log \hat{p}_{\text{data}}(\mathbf{x}) - \log p_{\text{model}}(\mathbf{x})]$$

$H(P_{\text{data}}, P_{\text{model}})$

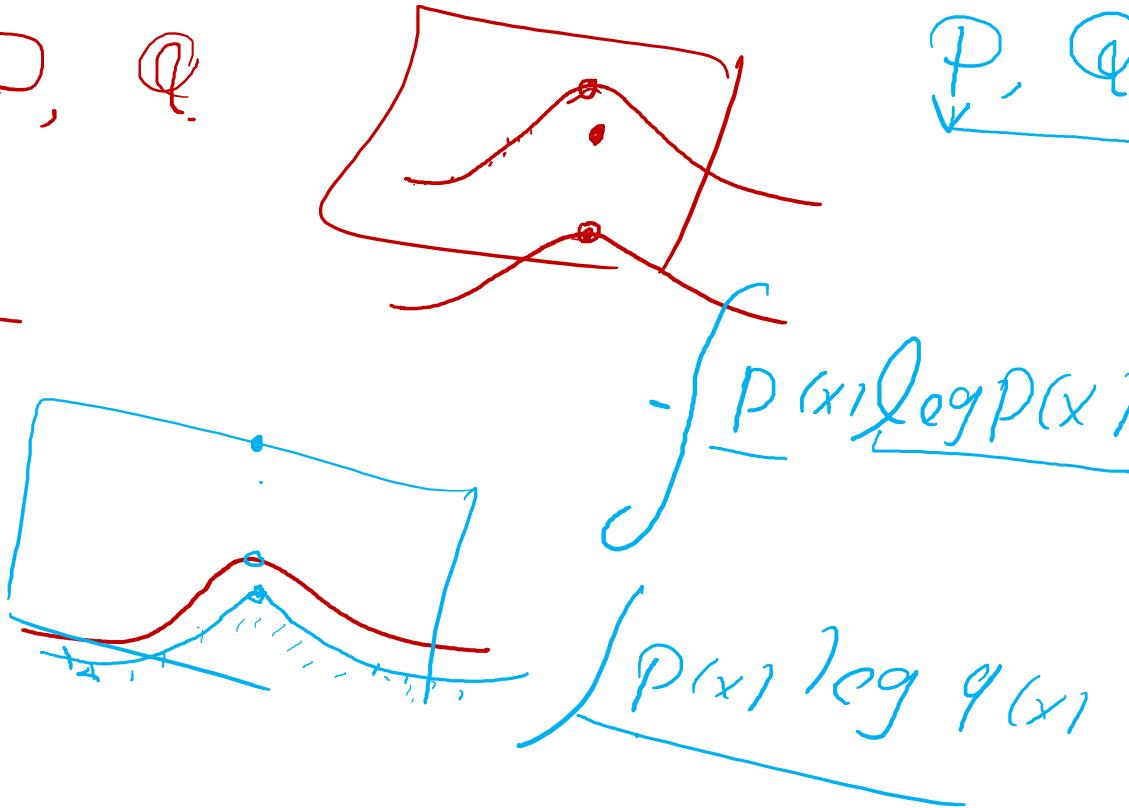
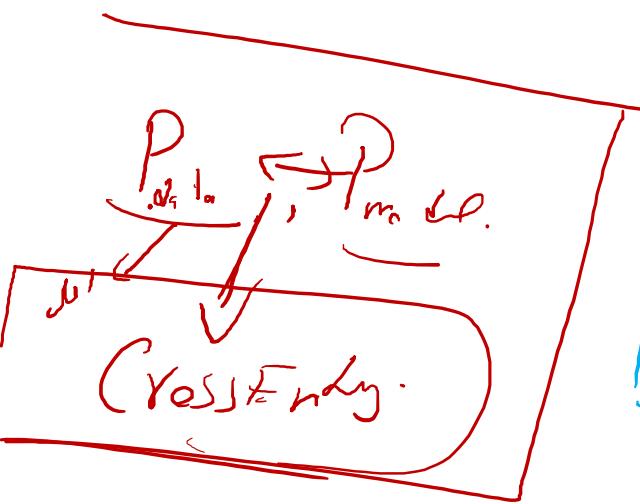
One way to interpret maximum likelihood estimation is to view it as minimizing the dissimilarity between the empirical distribution  $\hat{p}_{\text{data}}$  defined by the training set and the model distribution, with the degree of dissimilarity between the two measured by the KL divergence.

$$= \mathbb{E}_{\mathbf{x} \sim P_{\text{data}}} [-\log Q(\mathbf{x})] = H(P, Q)$$

P, q

P, Q

P, Q



161

# Mutual Information

$$KL \left( p(x,y) \parallel p(x)p(y) \right)$$

$$\begin{aligned} I[x, y] &\equiv KL(p(x, y) \| p(x)p(y)) \\ &= - \iint p(x, y) \ln \left( \frac{p(x)p(y)}{p(x, y)} \right) dx dy \end{aligned}$$

$$\underbrace{Q(x, y)}_{\text{Def}} = p(x)p(y)$$

$$\begin{aligned} &- \iint p(x, y) \ln p(x, y) + \iint \cancel{[p(x)p(y)]} \ln p(x)p(y) \\ &\quad - \iint p(x, y) \ln \frac{p(x)p(y)}{p(x, y)} dx dy \end{aligned}$$

$$I[x, y] = H[x] - H[x|y] = H[y] - H[y|x]$$

# Conditional log-Likelihood

If  $\mathbf{X}$  represents all our inputs and  $\mathbf{Y}$  all our observed targets, then the conditional maximum likelihood estimator is

$$\boldsymbol{\theta}_{\text{ML}} = \arg \max_{\boldsymbol{\theta}} P(\mathbf{Y} \mid \mathbf{X}; \boldsymbol{\theta}).$$



If the examples are assumed to be i.i.d., then this can be decomposed into

$$\boldsymbol{\theta}_{\text{ML}} = \arg \max_{\boldsymbol{\theta}} \sum_{i=1}^m \log P(\mathbf{y}^{(i)} \mid \mathbf{x}^{(i)}; \boldsymbol{\theta}).$$