

- $\chi$  (here  $\chi$  are bold letters) : A vector
- All vectors are assumed to be col. vectors
- $v$  (here  $v$  are bold Roman letters e.g.  $M$ ) denote matrices
- $(w_1, w_2, \dots, w_m)$ : A row vector with  $m$  elts

→ Generally col. vecn:  $w = (w_1, \dots, w_m)^T$

$[a, b]$ : Closed interval from  $a$  to  $b$   
inclusive of  $\underline{a}$  &  $\underline{b}$ .

$(a, b)$ : Open interval; excludes  $\underline{a}'$  &  $\underline{b}'$ .

$[a, b)$ :  $\underline{a}$  inclusive &  $\underline{b}'$  excluded.

$I_M$ :  $(M \times M)$  identity matrix.

$E_x [f(x, y)]$ : Expectation of a fun  $f(x, y)$  w.r.t. to a R.V.  $x'$

Conditional Expectation:  $E_x [f(x) | z]$   
If distib  $y | x'$  is conditioned on  $z$

Variance:  $\text{Var}[f(x)]$

Covariance:  $\text{cov}[x, y]$

$$\text{cov}[x, y] = \text{cov}[y, x]$$

$x = (x_1, x_2, \dots, x_D)^T$ :

$n$  values  $(x_1, x_2, \dots, x_n)$  of

a  $D$ -dimensional vector

$$x = (x_1, \dots, x_D)^T$$

$x = (x_1, \dots, x_n)^T$ : A training set

containing  $n$  observations of ' $x$ '.

$$\textcircled{O} \quad \text{Obj}(x, w) = w_0 + \sum_{j=1}^m w_j x^j + w_2 x^2 + \dots + w_m x^m$$
$$= \sum_{j=0}^m w_j \cdot x^j$$

↓

A polynomial func of  $y$  with order

$w = \text{vec}(w)$  of polynomial coeff's  $w_0, w_1, \dots, w_N$ .

Coeff. values will be determined by fitting the polynomial to the training data.

This can be done by minimizing an "Error func" that measures the misfit b/w the func  $y(x, w)$ , for any given value of  $\underline{w}$  & the training set data points.

$$E(w) = \frac{1}{2} \sum_{n=1}^N \{y(x_n, w) - t_n\}^2$$

## Regularization:

→ A technique to control the over fitting phenomenon

→ It adds a penalty ~~term~~ term to the error func in order to discourage the coeffs from reaching

large values.

$$\tilde{E}(\omega) = \frac{1}{2} \sum_{n=1}^N \{ y(x_n, \omega) - t_n \}^2 + \frac{\lambda}{2} \|\omega\|^2$$

$$\text{where } \|\omega\|^2 = \omega^T \cdot \omega$$

$$= \omega_0^2 + \omega_1^2 + \dots + \omega_N^2$$

The coeff.  $\lambda$  governs the relative importance of the regularization term compared with the SSE term.

Joint probability: The prob. that ' $x$ ' will take the value ' $x_i$ ' & ' $y$ ' will take the value ' $y_j$ ' is written

as

$$P(x=x_i, y=y_j)$$

and is called the Joint prob.

$X = x_i$  &  $Y = y_j$ .

Given by the no. of points falling in the cell  $(i,j)$  as a fraction of the total no. of points:

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

VVIMP:  $X, Y$ : 2 R.V.'s

$\rightarrow X$  takes the values  $\{x_i\}_{i=1,2,\dots,M}$

$\rightarrow Y$  takes the values  $\{y_j\}_{j=1,2,\dots,L}$

$\rightarrow n_{ij}$ :  $n_{ij}$  denotes the no. of points in the corresponding cell of the array.

$i \rightarrow$  column;  $j$ : row( $y_j$ )  
( $y_j$ )

⑥ Relationship b/w Conditional prob. &

$$p(y|x) = \frac{p(x|y) \cdot p(y)}{p(x)}$$

(Bayes theorem)

Dr. in Bayes theorem is a

→ normalization constant

→ Regd. to ensure that the sum of the conditional prob. on LHS over all values of  $y$  equals One.

P.T.O.

② Gaussian Distrib.: one of the most imp dists

- Used for continuous var.s
- Also called "Normal Distrib."
- Let  $x$  = a single real valued var.

Formula

$$N(x|\mu, \sigma^2) = \frac{1}{(2\pi\sigma^2)^{1/2}} \cdot \exp \left\{ -\frac{1}{2\sigma^2} (x-\mu)^2 \right\}$$

Governed by 2 parameters:

→  $\mu$  = Mean

→  $\sigma^2$  = variance

→  $\sqrt{\sigma^2} = \sigma$  = Std. deviation

→  $\beta = \frac{1}{\sigma^2}$  = Precision

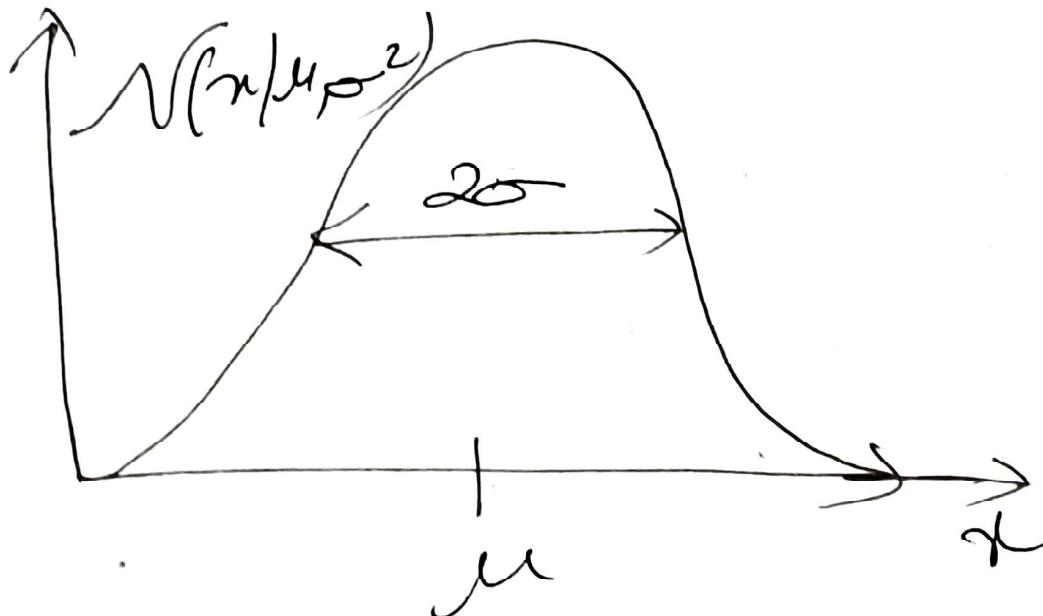
The Gaussian density satisfies:

$$N(x|\mu, \sigma^2) > 0$$

It is straightforward to show  
that the Gaussian is normalized.  
so that

$$\int_{-\infty}^{\infty} N(x|\mu, \sigma^2) dx = 1$$

Plot of univariate Gaussian



we can readily find expectations of func of 'x' under (1).

In particular, the avg. value of 'x' is given by:

$$\begin{aligned} E[x] &= \int_{-\infty}^{\infty} N(x|\mu, \sigma^2) x dx \\ &= \mu \end{aligned}$$

$\mu$  = Represents the avg. value of 'x'  
under the distrib., it is  
referred to as the mean.

likg for the 2nd order moment:

$$\begin{aligned} E[x^2] &= \int_{-\infty}^{\infty} N(x|\mu, \sigma^2) x^2 dx \\ &= \mu^2 + \sigma^2 \end{aligned}$$

Variance of  $x$ :

$$\text{var}[x] = E[x^2] - E[x]^2$$

$$= (\mu^2 + \sigma^2) - (\mu)^2$$

$$= \sigma^2$$

$\hat{=}$  Variance Parameter

The PDF defined over a D-dimensional vector ' $x$ ' of continuous variables is given by:

$$p(x|\mu, \Sigma) \cdot \frac{1}{(2\pi)^{D/2}}$$

$$= \frac{1}{(2\pi)^{D/2}} \cdot \frac{1}{|\Sigma|^{1/2}} \cdot$$

$$\exp \left\{ -\frac{1}{2} (x-\mu)^T \Sigma^{-1} (x-\mu) \right\}$$

$\mu$  = a D-dimensional vector  
called the mean

$\Sigma = D \times D$  matrix called  
covariance

$$|\Sigma| = \text{Determinant of } \Sigma$$

we have a data set of observations:

$$\mathbf{x} = (x_1, \dots, x_N)^T$$

represents  $N$  observations of the scalar variable  $x$ .

e.g.:  $\underline{\mathbf{x}}$  is a vector of the form:

	$f_1$	$\vdots$	$f_D$
$n_1$			
$n_2$			
$\vdots$			
$n_D$			

(1)  $D$ -Dimensions (or)  
features

(2)  $x_1, \dots, x_N$   
are no. of observations  
of  $x$ .

Assumptions) The observations are drawn independently from a  $\text{GD}$  where  $\mu$  &  $\sigma^2$  are constant;

& we would like to determine these parameters from the data set.

u.i.i.d: data points that are drawn independently from the same distib.

→ (1) Unconditional prob:

Marginal prob: The prob. of an event irrespective of the outcome

of another variable.  $P(A)$   
ex: A card drawn is red ( $P(\text{red}) = \frac{2}{4} = 0.5$ )

Conditional prob: The prob. of one event occurring in the presence of a second event.

Joint prob: prob. of 2 events occurring simultaneously.

P(A): prob. of event A (or)  $\downarrow P(x=A)$   
a R.V.

P(n): 'n' is a R.V;  
 $P(n)$  assigns a prob. to all  
 $\downarrow$   
values of 'n'.  
at n(c)

probability =  $\frac{\text{No. of desired outcomes}}{\text{Total no. of possible outcomes}}$

Note: Sum of probs. of all outcomes  
would be equal to one.

{prob. of a certain outcome = 1.0}

Complement: prob. of an event not  
occurring.

Denoted as  $\boxed{1 - P(A)}$

(A) (1) The prob. of a row of data is the joint prob. across each I/P var.

(2) Probs. of a specific value  $y$  one I/P var. is the marginal prob. across the values of other I/P vars.

(3) The "predictive model" ~~is~~ itself is an estimate of the condn. prob. of an O/P given an I/P  $\underline{\underline{eg}}$

Joint prob. Distrib: The joint prob. of 2 or more R.V.'s.  $\leftrightarrow$

$$\text{eg: } \boxed{P(A \text{ and } B)} \\ = P(A \text{ given } B) + P(B)$$

## (Q) Fundamental Rule of prob.

(Q1)

Product Rule (Q1) Chain Rule

= calculation of Joint prob.

Joint prob. is symmetrical i.e.

$$\boxed{P(A \text{ and } B) = P(B \text{ and } A)}$$

$$\begin{aligned} P(A \text{ and } B) &= P(A \text{ given } B) * P(B) \\ &= P(B \text{ given } A) * P(A) \end{aligned}$$

$$P(A \text{ given } B) = \frac{P(A \text{ and } B)}{P(B)}$$

$$\text{Joint prob: } P(A \text{ and } B) = P(A) * P(B)$$

Note: The "marginal prob." for an event  
fr an independent R.V. is simply  
the prob. of the event.

Cond. prob:  $P(A \text{ given } B) = P(A)$

$\Rightarrow$  prob. of B has no effect  
on A'.