

Vector Valued Rvs:

Here the Range set of the fun^c(RV) is \mathbb{R}^d , where "d" is a scalar

$$X: \Omega \rightarrow \mathbb{R}^d$$

Joint distribution:

Let Ω be a sample space

Define 2 fun's, x_1, x_2

$$\begin{aligned} x_1: \Omega &\rightarrow \mathbb{R} \\ x_2: \Omega &\rightarrow \mathbb{R} \end{aligned} \quad] \quad \text{multiple RV on same } \Omega$$

Corresponding P_{x_1}, P_{x_2} be the respective DF.

define joint Prob distributions as

$$P_{x_1, x_2}(a, b) = \text{IP}[E: \text{Intersection of inverse images of } (-\infty, a] \text{ & } (-\infty, b)]$$

under x_1, x_2 respectively]

The above idea can be extended to "d" scalar RV.

Define Conditional Probabilities / Distribution s.

if $A \& B \in \mathcal{F}$

$$P[A|B] = \frac{P[A \cap B]}{P[B]}$$

Suppose $X \& Y$ are 2 RVs defined on same Ω

(Range can be diff)

$$\begin{aligned} X: \Omega &\rightarrow \mathbb{R}^d \\ Y: \Omega &\rightarrow \mathbb{R}^k \end{aligned} \quad d \neq k$$

$$P_{X|Y}(x|y) = \frac{P_{XY}}{P_Y}$$

Joint
Marginal

Marginal distribution : If $X \& Y$ are 2 RVs

marginal of X is defined as.

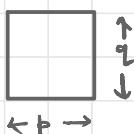
$$P_X = \int_y P_{XY} dy$$

→ Assuming that integration is possible.

$$P_Y = \int_x P_{XY} dx$$

→ depending on the dim it can be a multidimensional integral

Example : Consider an image of dim $p \times q$



Can be viewed as a \mathbb{R}^{pq} dim Vector.

- * In general any datapoint is a "d" dim vector.
- * every datapoint is an element in the range space of a random variable
- * The distribution function indicates ② quantifies the likelihood of observing a datapoint under \mathbf{x} .
- * Distribution func completely specifies the sample space.
- * Typically a "label" is taken to be another additional RV, defined on the same sample space.
- * In this scenario, we have 2 RVs, \mathbf{x} & y

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

$$y \in \mathbb{R}^k$$

- * The convention is as follows :

"Learning" starts with a dataset

$$D = \{(x_1, y_1), (x_2, y_2), \dots, (x_n, y_n)\} \sim \text{iid } P_{xy}$$

x : Feature ② Data Space

y : Target ② Label space.

* iid : Independent & Identically distributed.

$$P[A \cap B] = P[A] P[B] : \text{statistically indpt.}$$

* in a vector valued RV, we are not assuming independence among the scalars of that vector RV.

* our assumption is 2 data points $\in \mathbb{R}^d$ are independent.

* Independence is across data points & not data dimensions.

* All problems in machine learning can be seen as
Given D n unknown P_x , estimate P_x ② Conditionals on P_x
↳ Sample from it.

Lecture 4: 14th Jan 2026

Recall we have given a dataset

$$D = \{(x_i, y_i)\}_{i=1}^n \sim \text{iid } P_{xy} : \text{Unknown.}$$

$$x_i \in \mathbb{R}^d$$

$$y_i \in \mathbb{R}^k \text{ or } \{1, 2, \dots, K\}$$

Typically X : input / features / data

Y : Label / output

X & Y are random variables on a Sample Space

The Fundamental problem of ML :

Given $D \sim$ Unknown distribution.

↳ Estimate the distribution.

↳ Sample from the distribution.

Let's consider the Distribution Estimation.

Given $D = \{(x_i, y_i)\}_{i=1}^n \sim \text{iid } P_{xy}$

estimate the distribution (or) the moments of the distribution.

Examples :

1) Estimate $IP_{Y|X}$

Classification $Y \in \{1, 2, \dots\}$

Regression $Y \in \mathbb{R}^d$

(Bounding box)

2) Estimate

IP_x

IP_y

IP_{XY}

$P_{x|y}$

* Estimating joint (a) marginal, does not mean its generative.

dummy \leftarrow \rightarrow fixed
 $(x | y=y)$

* In all conditional $IP_{x|y}$

Conditional distribution is for a specific "y"

Probability density func:

Given a CRV with a dist func IP_x ,

density func $f_x : X \rightarrow \mathbb{R}^+$

$$P_x(x) = \int_{-\infty}^x p_x(x) dx.$$

- * $p_x(x)$ does not correspond to probability
 - * Not all RV has density func
 - * Our assumption is all distribution func that we consider has well behaved density func
 - + we will estimating density from now on.
-

The challenge with ML

- * Given David form Unknown p_x , estimate p_x
- * challenge is P_x is completely unknown
- * Given the samples from P_x can we estimate P_x .

Question :

- * How to estimate a density func given its samples.
- * Consider we have a dataset D ,

$$D = \{x_i\}_{i=1}^N \sim \text{iid } p_x$$

Mathematically how do we solve it?

① Assume a parametric functional form on ϕ_x .

& let's call it a ϕ_θ

$$\text{eg: a) } \phi_\theta^a(x) = w^T x + b$$

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

$$w \in \mathbb{R}^d$$

$$b \in \mathbb{R}$$

$$\text{b) } \phi_\theta^b(x) = N(x; \mu, \Sigma)$$

$$\theta = \{\mu, \Sigma\}$$

$$\mu \in \mathbb{R}^d$$

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

These are different model choices

- * This assumption on ϕ_x is a huge step.
- * We also have universal func approx.
- * Once we make an assumption on parametric assumption

on ϕ_θ ,

our algorithm can say is that parametric form which one to use.

② Define ③ Compute a distance metric between ϕ_x & ϕ_θ

We need say how good our assumption is.

* Let D_m denote the "distance metric" b/w

$$\hat{p}_x \text{ & } p_\theta$$

$$* D_m(\hat{p}_x || p_\theta) : \hat{p}_x \times p_\theta \rightarrow \mathbb{R}^+$$

③ Find ② Estimate the parameters θ by solving the following optimization problem.

$$\Theta^* = \underset{\theta}{\operatorname{argmin}} D_m(\hat{p}_x || p_\theta)$$

This is the Training of
model.