

Probability Theory Review

Random variable

The r.v. X takes a value in \mathcal{X} (set of possible values) according to the outcome of an event.

Eg. Event: flip a coin

$$\mathcal{X} = \{H, T\}$$

X has value H when coin is heads, & T when it's tails.

Associated with r.v. X is a distribution $P(X=x)$ that describes the frequency of the r.v. value events.

Examples: Discrete r.v.

indicator variable

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

didn't happen

it happened

probability mass function (pmf)

$P(X=x)$ = probability of x occurring

$$\sum_{x \in \mathcal{X}} P(X=x) = 1$$

$$0 \leq P(X=x) \leq 1, \forall x \in \mathcal{X}$$

Notation: $P(X=x) = p(x)$
 $= P_X(x)$
 value
 r.v.

of people in a room

$$\mathcal{X} = \mathbb{Z}_+$$

non-negative integers

Continuous r.v.

Sensor reading

$$\mathcal{X} = \mathbb{R}$$

probability density function (pdf)

$p(x)$ = likelihood of x

$$P(a \leq x \leq b) = \int_a^b p(x) dx$$

$$\int p(x) dx = 1$$

$$0 \leq p(x) \forall x \in \mathcal{X}$$

Example Distributions

• Bernoulli (coin)

$$\mathcal{X} = \{0, 1\}, \pi = \text{probability that 1 occurs.}$$

$$\begin{cases} P_X(1) = \pi \\ P_X(0) = 1 - \pi \end{cases}$$

$$p(x) = \pi^x (1 - \pi)^{1-x}$$

$$x=0 \Rightarrow \pi^0 (1-\pi)^1 = 1-\pi$$

$$x=1 \Rightarrow \pi^1 (1-\pi)^0 = \pi$$

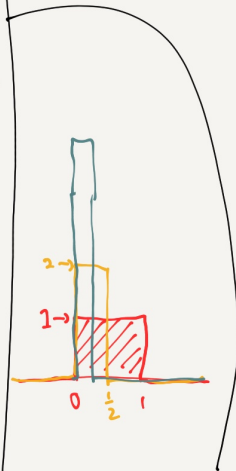
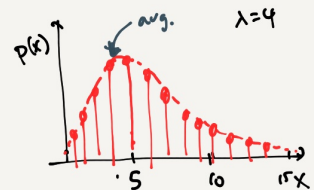
• Poisson (# of arrivals in a fixed time period)

$$\mathcal{X} = \mathbb{Z}_+ = \{0, 1, 2, \dots\}, \lambda = \text{average arrival rate } (\lambda \geq 0)$$

$$p(x) = \frac{1}{x!} e^{-\lambda} \lambda^x$$

x factorial

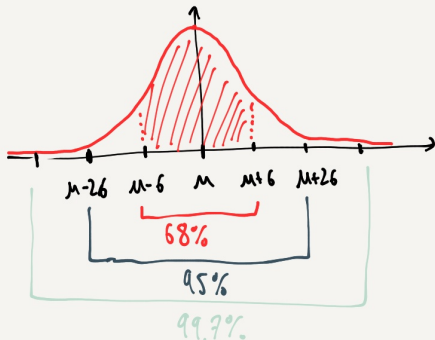
$$x \cdot (x-1) \cdot (x-2) \dots \cdot 1$$



Normal (Gaussian)

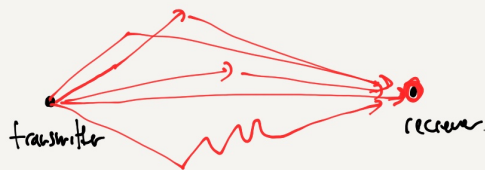
$X = \mathbb{R}$, μ = mean σ^2 = variance > 0
 σ = standard deviation

$$p(x) = \frac{1}{\sqrt{2\pi\sigma^2}} e^{-\frac{1}{2\sigma^2}(x-\mu)^2}$$

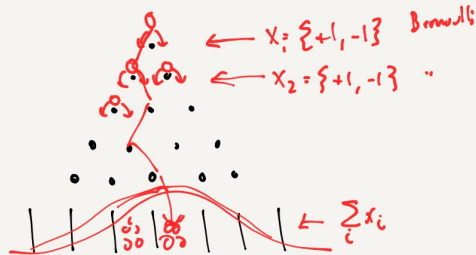


Central Limit Theorem (CLT)

Sum of N r.v. \rightarrow Gaussian distribution for large N



LLN



Joint distributions

distribution over more than 1 r.v.

probability/likelihood of $X=x$ & $Y=y$:

$$P(X=x, Y=y) = P(x, y)$$

Example: $X = \{0, 1\}$
 $Y = \{0, 1\}$

Marginal distribution

distribution over one variable in the joint.

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

\rightarrow summing over other r.v.

$$p(x) = \int_{y \in \mathcal{Y}} p(x, y) dy$$

\rightarrow integrating out other r.v.

"marginalization"

conditional distribution

distribution of one r.v. when the value of another r.v. is known (given)

$$P(X=x | Y=y) = \frac{P(X=x, Y=y)}{P(Y=y)}$$

\uparrow r.v. \uparrow "given" \uparrow "known value"

$$p(x, y) = p(x|y)p(y)$$

Statistical independence

distribution of a r.v. does not change when given the value of another r.v.

$$\Rightarrow \textcircled{1} X \perp\!\!\!\perp Y \text{ iff } p(x|y) = p(x)$$

$$\textcircled{2} X \perp\!\!\!\perp Y \text{ iff } p(x, y) = p(x)p(y) \leftarrow \text{joint is product of marginals.}$$

$p(x, y)$	$Y=0$	$Y=1$	$P(x)$
$X=0$	0.08	0.12	0.2
$X=1$	0.32	0.48	0.8
$P(y)$	0.4	0.6	

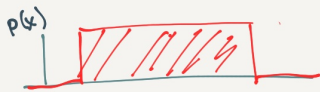
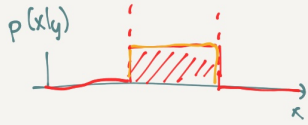
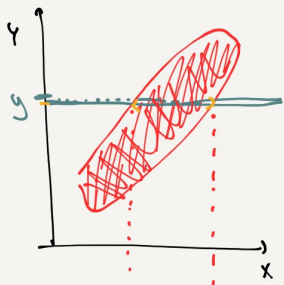
why is it called marginalization? "margin of table"

e.g.

$$P(X=0 | Y=0) = \frac{0.08}{0.4} = 0.2$$

$$P(X=1 | Y=0) = \frac{0.32}{0.4} = 0.8$$

in our example $X \perp\!\!\!\perp Y$.



Bayes' Rule

$$\left. \begin{aligned} p(x, y) &= p(x|y)p(y) \\ p(x, y) &= p(y|x)p(x) \end{aligned} \right\} p(y|x)p(x) = p(x|y)p(y)$$

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

$$\begin{aligned} p(x) &= \int p(x, y) dy \\ &= \int p(x|y)p(y) dy \end{aligned}$$

given only $p(x|y)$ & $p(y)$,
we can "invert" the conditioning
to get $p(y|x)$

$$\Rightarrow p(y|x) = \frac{p(x|y)p(y)}{\int p(x|y)p(y) dy}$$

Y = class label (digit)

X = feature value

$p(X|Y)$ = distribution of features for class Y

$p(Y)$ = prob. of class Y .

\Rightarrow Bayes' Rule $\Rightarrow p(Y|X)$
↑ ↑
class feature

Expectations

Suppose we have a function $f(x)$ & r.v. X .

On average, what is the value of $f(X)$?

$$E_X[f(x)] = \sum_{x \in \mathcal{X}} f(x)p(x) \leftarrow \text{weighted average of function values.}$$

$$E_X[f(x)] = \int f(x)p(x)dx$$

• mean: $E[X] = \int x p(x) dx = \mu_x$

• variance $\text{Var}(x) = E[(X - E(x))^2] = \sigma_x^2$
 $= E[X^2] - (E[X])^2$

• covariance: $\text{cov}(x, y) = E_{xy}[(X - E(x))(Y - E(y))]$
 $= \int (x - \mu_x)(y - \mu_y) p(x, y) dx dy$
 $= E_{xy}(x, y) - E(x) \cdot E(y) = \sigma_{xy}$

Conditional Expectation

$$E_{x|y}[x] = \int x p(x|y) dx = \text{function of } y$$

known

$$E_{x|y}[f(x)] = \int f(x) p(x|y) dx$$

Brief Linear Algebra Review

column vector: $x \in \mathbb{R}^d$ $x = \begin{bmatrix} x_1 \\ \vdots \\ x_d \end{bmatrix}$ (usually we are column-centric)

matrix: $A \in \mathbb{R}^{m \times n}$ $A = \begin{bmatrix} a_{11} & \dots & a_{1n} \\ \vdots & \ddots & \vdots \\ a_{m1} & \dots & a_{mn} \end{bmatrix} = \begin{bmatrix} | & \dots & | \\ a_1 & \dots & a_n \\ | & \dots & | \end{bmatrix}$
of rows # of columns
 a_i is the i th column of A .

inner product: $x^T y = \sum_i x_i y_i$ (similarity b/w vectors x & y)

$$\text{length (norm)}: \|x\| = \sqrt{x^T x} = \sqrt{\sum_i x_i^2}$$

$$\text{distance: } d(x, y) = \|x - y\|$$

outer product: $xy^T = \begin{bmatrix} x_1 y_1 & \dots & x_1 y_n \\ \vdots & \ddots & \vdots \\ x_n y_1 & \dots & x_n y_n \end{bmatrix}$ ← matrix of all pairwise products of elements in x & y .

matrix-vector multiplication

$$\textcircled{1} y = Ax = \begin{bmatrix} 1 & \dots & 1 \\ a_1 & \dots & a_n \\ 1 & \dots & 1 \end{bmatrix} \begin{bmatrix} x_1 \\ \vdots \\ x_n \end{bmatrix} = \sum_i x_i a_i \in \mathbb{R}^m$$

$A \in \mathbb{R}^{m \times n}$ $x \in \mathbb{R}^n$
coefficient vector
linear combo of the columns of A w/ entries in x as the coeff.

$$\textcircled{2} y = A^T x = \begin{bmatrix} -a_1^T \\ \vdots \\ -a_m^T \end{bmatrix} \begin{bmatrix} 1 \\ x \\ 1 \end{bmatrix} = \begin{bmatrix} a_1^T x \\ \vdots \\ a_m^T x \end{bmatrix} \in \mathbb{R}^m$$

$A \in \mathbb{R}^{d \times m}$ $x \in \mathbb{R}^d$
inner products b/w columns of A & vector x .

Matrix-matrix multiplication

① $AB = A[b_1 \dots b_n] = [Ab_1 \dots Ab_n] \in \mathbb{R}^{m \times n}$
 $\uparrow \quad \uparrow$
 $\mathbb{R}^{m \times d} \quad \mathbb{R}^{d \times n}$
"A multiplied w/ each column of B"

② $A^T B = \begin{bmatrix} -a_1^T \\ \vdots \\ -a_m^T \end{bmatrix} \begin{bmatrix} b_1 \dots b_n \\ \vdots \\ b_n \end{bmatrix} = [a_i^T b_j]_{i,j} \in \mathbb{R}^{m \times n}$
 $\uparrow \quad \uparrow$
 $A \in \mathbb{R}^{d \times m} \quad \mathbb{R}^{d \times n}$
"all the inner products btwn columns of A & B"

③ $AB^T = \begin{bmatrix} 1 & \dots & 1 \\ a_1 & \dots & a_d \\ \vdots & & \vdots \\ 1 & \dots & 1 \end{bmatrix} \begin{bmatrix} -b_1^T \\ \vdots \\ -b_d^T \end{bmatrix} = \sum_{i=1}^d a_i b_i^T \in \mathbb{R}^{m \times n}$
 $\uparrow \quad \uparrow$
 $\mathbb{R}^{m \times d} \quad \mathbb{R}^{d \times n}$
"sum of outer-products of columns of A & B"

Vector r.v.

$$X = \begin{bmatrix} x_1 \\ \vdots \\ x_d \end{bmatrix}, \quad \mathcal{X} = \mathbb{R}^d$$

Notation

$$p(x_1, x_2, x_3, \dots, x_d) = p(x) \quad , \quad \int p(x) dx = 1$$

\uparrow vector \downarrow

$$\int \dots \int p(x_1, x_2, \dots) dx_1 dx_2 \dots = 1$$

Mean vector

$$\mu = \begin{bmatrix} \mu_1 \\ \vdots \\ \mu_d \end{bmatrix} = E[X] = \int_{\mathbb{R}^d} x p(x) dx \in \mathbb{R}^d$$

\uparrow vector \uparrow number

Covariance matrix

$$\text{cov}(X) = E[(X - E[X])(X - E[X])^T] = E[X X^T] - E[X] E[X]^T$$

2d e.g. $= E \begin{bmatrix} (x_1 - \mu_1)^2 & (x_1 - \mu_1)(x_2 - \mu_2) \\ (x_1 - \mu_1)(x_2 - \mu_2) & (x_2 - \mu_2)^2 \end{bmatrix}$

$$= \begin{bmatrix} \sigma_1^2 & \sigma_{12}^2 \\ \sigma_{12}^2 & \sigma_2^2 \end{bmatrix}$$

"covariances btwn dim. of r.v."
"variances of each dim. of r.v."

Multivariate Gaussian

$$\mathcal{X} = \mathbb{R}^d$$

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

$$\text{cov matrix } \Sigma \in \mathbb{S}_{++}^d$$

$\leftarrow d \times d$
positive
definite
(symmetric)
matrix

$$p(x) = \frac{1}{(2\pi)^{d/2} |\Sigma|^{1/2}} e^{-\frac{1}{2} \|x - \mu\|_{\Sigma}^2} = \mathcal{N}(x | \mu, \Sigma)$$

Mahalanobis distance:

$$\|x - \mu\|_{\Sigma}^2 = (x - \mu)^T \Sigma^{-1} (x - \mu) \rightarrow \frac{(x - \mu)^2}{6^2}$$

determinant = $|\Sigma|$ = "volume of Gaussian"

Special cases: Σ is a diagonal matrix = $\begin{bmatrix} \sigma_1^2 & 0 \\ 0 & \ddots & \sigma_d^2 \end{bmatrix}$

$$p(x) = \prod_{i=1}^d \mathcal{N}(x_i | \mu_i, \sigma_i^2) = \text{product of univariate Gaussians.}$$

\uparrow joint \uparrow product \uparrow marginals

i.e. d independent univariate Gaussians on each dim.

(TBC in tutorial)

