

Derived by PMF/CPF
PMF

Mass
 $X \sim F(x)$

$$\rightarrow P(X=x)$$



Density
 $X \sim F(x)$

$$F(x) \rightarrow P(X \leq x)$$



- in domain
- $\forall x \in X, P(x) \geq 0$
- $\int P(x) dx = 1$

Joint

$$P(x, y) \rightarrow P(x=x, y=y)$$

Marginal

$$P(x) \rightarrow P(x=x, y=\text{arbitrary}) = \sum_y P(x, y=y)$$

or

$$\int P(x, y) dy$$

Conditional

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

Chain Rule immediately follows:

Do Bayes Rule here
↓

$$P(a, b, c, \dots) = \underbrace{P(a, b, c)}_{P(a|b, c)} \cdot \underbrace{P(b, c)}_{P(b|c)} \cdot P(c)$$

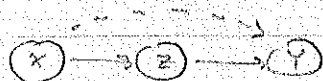
$$P(a|b, c) \cdot P(b|c) \cdot P(c)$$

Independence:

Share no information $\xrightarrow{P(a)}$ $a \perp b$

$$P(a, b) = P(a|b) \cdot P(b)$$

CS - Show no information given you know z.



$X \perp Y | Z$

$$P(X, Y | Z) = \underbrace{P(X, Y, Z)}_{P(Y, Z)} \cdot \underbrace{P(Y, Z)}_{P(Z)}$$

$$P(X | Y, Z) \cdot P(Y | Z)$$

$$\text{no info gained} \rightarrow P(X | Z)$$

Examples:

0.05	0.1	0.2
0.3	0.15	0.2

0.35
0.65

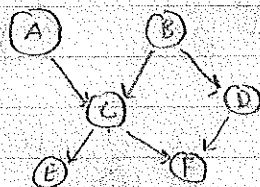
0.5	0.25	0.57
0.46	0.33	0.31

0.35	0.25	0.4
------	------	-----

5	10	20
30		20

0.15	0.4	0.5
0.25	0.6	0.5

Introduce
disks

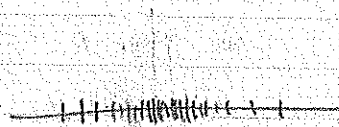
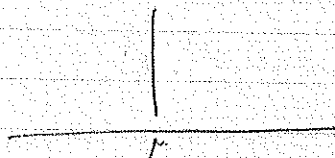


$$P(F | C, D) \cdot P(E | C) \cdot P(D | B) \cdot P(C | A, B) \cdot P(B) \cdot P(A)$$

See Bishop 85

Dirac Delta

Empirical Dist.



Mixture: $Y \sim \text{Bernoulli}(P)$

$$(X|Y=y) \sim N(y\mu_1 + (1-y)\mu_2, 1)$$

$$Y=0 \rightarrow P$$

$$Y=1 \rightarrow (1-P)$$

$$\text{CDF of } X: P(X \leq x) = \sum_{y=0}^1 P(X \leq x | Y=y) P(Y=y)$$

$$F_X(x) =$$

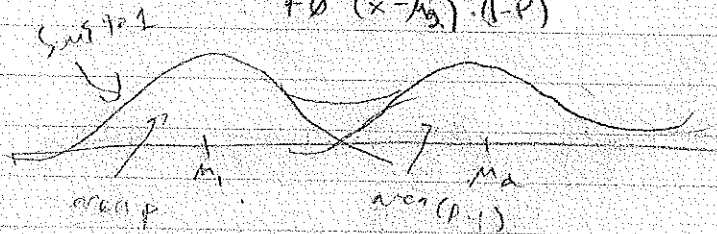
$$= P(X \leq x | Y=0) P(Y=0) + P(X \leq x | Y=1) P(Y=1)$$

$$= P(X=\mu_1 \leq x-\mu_1 | Y=0) P(Y=0) + P(X=\mu_2 \leq x-\mu_2 | Y=1) P(Y=1)$$

$$= \Phi(x-\mu_1) P(Y=0) + \Phi(x-\mu_2) P(Y=1)$$

$$\frac{d}{dx} F_X(x) = \phi(x-\mu_1) P$$

$$+ \phi(x-\mu_2) (1-P)$$



Y can be latent variable

CLT: Let X_1, X_2, \dots, X_n i.i.d. with mean μ and var $\sigma^2 < \infty$

$$\sqrt{n}(\bar{X} - \mu) \Rightarrow N(0, \sigma^2)$$

LLN: $\bar{X}_n \xrightarrow{P} \mu$

$$\bar{X}_n \xrightarrow{P} \mu$$

$$\text{Var}(f(x)) = \mathbb{E}((f(x) - \mathbb{E}(f(x)))^2)$$

$$\text{Cov}(f(x), g(y)) = \mathbb{E}((f(x) - \mathbb{E}(f(x)))(g(y) - \mathbb{E}(g(y))))$$

$$\mathbb{E}(x) = \frac{1}{n} \sum_{i=1}^n x_i \text{ or } \int x \cdot p(x) dx$$

$$\frac{1}{n} \sum_{i=1}^n (x_i - \mu)^2$$

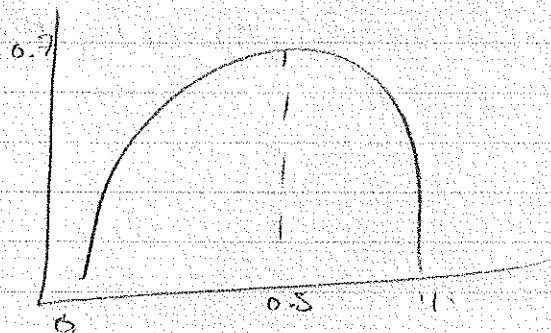
$$\frac{1}{n} \sum_{i=1}^n (x_i - \mu_x)(y_i - \mu_y)$$

Useful functions

Sigmoid $\sigma(x) = \frac{1}{1+e^{-x}}$ Proline Bern,
Saturates

$$\text{Softplus } \phi(x) = \log(1 + \exp(x))$$

Softened version of ReLU



$$H(x) = \text{Exp}[I(x)] = -\text{Exp}[\text{Log } P(x)]$$

Sunny	cloudy	raining	foggy
00	01	10	11

need 2 bits \rightarrow motivates $\log_2(\frac{1}{p})$

If: Sunny more likely than hazy

$\frac{1}{2}$	$\frac{1}{4}$	$\frac{1}{8}$	$\frac{1}{8}$
0	10	110	111

$$\log(2) = 1$$

$$\log(4) = 2$$

$$\log(8) = 3$$

$$\cdot \frac{1}{2}$$

$$\cdot \frac{1}{4}$$

$$\cdot \frac{1}{8} = \frac{1}{8}$$

$$\frac{1}{2}$$

$$\frac{1}{2}$$

$$\frac{3}{8}$$

$$= 1\frac{3}{8}$$

or

$$\frac{1}{2} \cdot 2 = \frac{1}{2} \cdot 4 = 2$$

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

$$-\sum p \ln\left(\frac{p}{q}\right) \geq -\sum p \left(\frac{q}{p} - 1\right)$$

$$\text{by } \ln x \leq x - 1$$

$$-\sum p \left(\frac{q}{p} - 1\right) = -\sum q + \sum p = 0$$

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

how much it cost to code Q
under P 's coding scheme.

$$H(P, Q) = H(P) + D_{KL}(P \parallel Q)$$

minimizing is equiv.