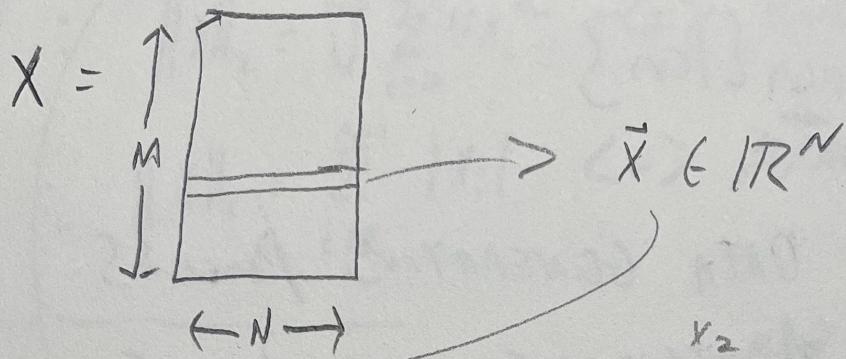


FEATURE REPRESENTATION

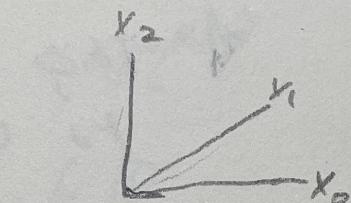
$$X \in \mathbb{R}^{M \times N}$$

M = Number of Observations

N = Number of Features



D VECTOR SPACE:

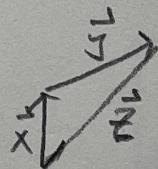


$$x = [0.2, 0.5, 1.0]$$

Norm: $f(\cdot) : \mathbb{R}^N \rightarrow \mathbb{R}$

$f(\cdot)$ MUST SATISFY:

$$\textcircled{1} \quad f(\vec{x}) = 0 \text{ iff } \vec{x} = \vec{\phi}$$

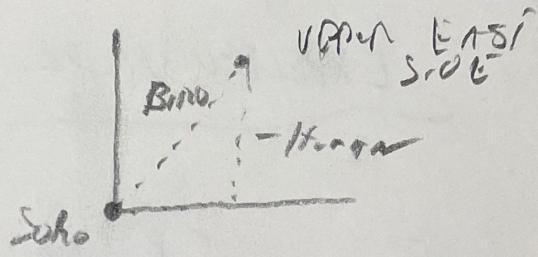


$$\textcircled{2} \quad f(\vec{x} + \vec{y}) \leq f(\vec{x}) + f(\vec{y})$$

$$\textcircled{3} \quad f(\alpha \vec{x}) = |\alpha| f(\vec{x}) \quad \alpha \in \mathbb{R}$$

L_p Norms

$$\|x\|_p = \left[\sum_{i=0}^{n-1} |x_i|^p \right]^{\frac{1}{p}}$$



$$L_2 : \|x\|_2 = \sqrt{\sum_{i=0}^{n-1} x_i^2} \iff \text{EUCLIDEAN}$$

$$L_1 : \|x\|_1 = \sum_{i=0}^{n-1} |x_i| \iff \text{MANHATTAN}$$

FROBENIUS NORM $A \in \mathbb{R}^{M \times N}$

$$\|A\|_F = \sqrt{\sum_{i=0}^{M-1} \sum_{j=0}^{N-1} A_{ij}^2} \in \mathbb{R}$$

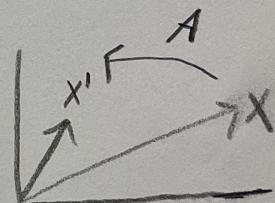
INNER PRODUCT ("DOT" Product)

$$\langle x, y \rangle = x \cdot y = \sum_{i=0}^{n-1} x_i y_i \in \mathbb{R}$$

TENSOR (LINEAR TRANSFORMER)

$$x' = Ax \quad x, x' \in \mathbb{R}^N$$

$$A \in \mathbb{R}^{N \times N}$$



$$\begin{bmatrix} A_{00} & \dots & A_{0N} \\ \vdots & & \vdots \\ A_{m0} & \dots & A_{mN} \end{bmatrix} \begin{bmatrix} x_0 \\ \vdots \\ x_N \end{bmatrix} = \begin{bmatrix} \sum_{i=0}^N A_{0i} x_i \\ \vdots \\ \sum_{i=0}^N A_{Ni} x_i \end{bmatrix}$$

Basis of words (BOW)

$$x = [0, 0, 0, \dots, 1, 0, 0, \dots, 2, \dots]_N$$

↳ Sparse Representation

Singular Value Decomposition (SVD)

$$X \in \mathbb{R}^{M \times N}$$

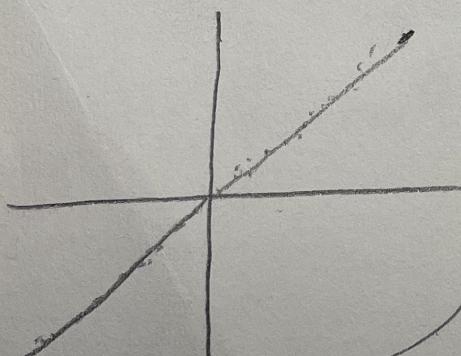
U, V are ORTHOGONAL

$$U^{-1} = U^T$$

$$X = U \sum_{M \times N} V^T \quad M \times M \quad M \times N \quad N \times N$$

$$V^{-1} = V^T$$

$$\Sigma = \begin{bmatrix} \sigma_1 & & & & 0 \\ & \ddots & & & 0 \\ 0 & & \ddots & & \sigma_n \\ \hline & & & & \text{Null Space} \end{bmatrix}$$



① $\text{Rank}(X) = \text{Number of Non-Zero } \sigma_i^{\text{'}} \text{ s} := K$

② Truncated SVD

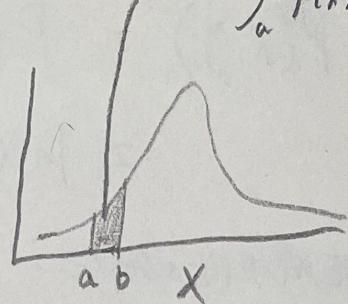
$$X = U' \Sigma' V'^T \quad M \times N \quad M \times K \quad K \times K \quad K \times N$$

$$③ q \in \mathbb{R}^N \Rightarrow \hat{q}' = q V' \quad 1 \times K \quad 1 \times N \quad N \times K$$

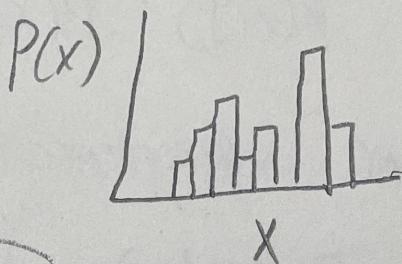
RANDOM VARIABLE (RV)

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

CONTINUOUS : PDF $f(x)$



DISCRETE : PMF



(D) PMF must satisfy

$$\textcircled{1} \quad \forall x \in X \quad 0 \leq P(x) \leq 1$$

$$\textcircled{2} \quad \sum_{x \in X} P(x) = 1$$

(Y) $\sigma_{\text{softmax}}(x) = \frac{e^x}{\sum_{x' \in X} e^{x'}}$

Joint, Marginal, Conditional O/S/T.

Joint: $P(X, Y) = \frac{1}{|X|} \boxed{}$

MARGINAL: $P(X) = \sum_{y \in Y} P(X, y) \Rightarrow P(Y) = \sum_{x \in X} P(y, x)$

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

Product Rule:

$$\begin{aligned} P(x,y) &= P(y|x)P(x) \\ &= P(x|y)P(y) \end{aligned}$$

INDEPENDENCE:

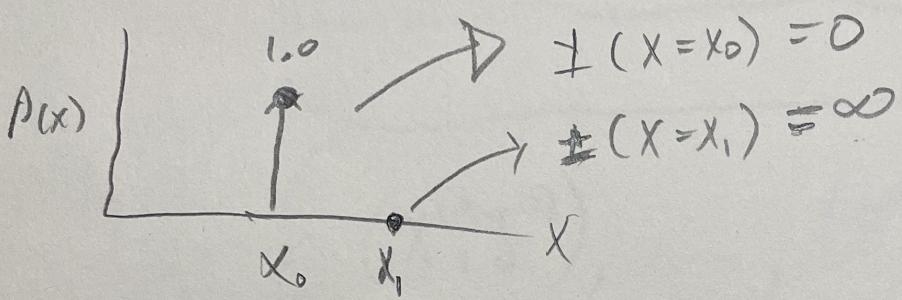
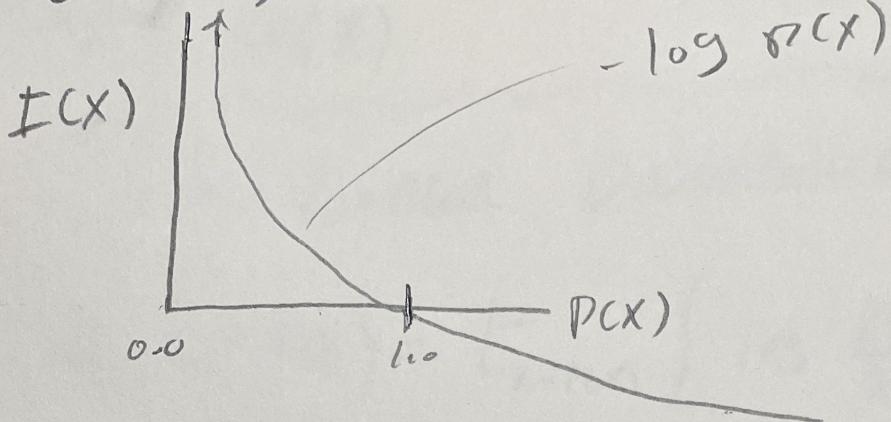
$$P(x,y) = P(x)P(y)$$

CONDITIONAL INDEPENDENCE:

$$P(x,y|z) = P(x|z)P(y|z)$$

INFORMATION THEORY

- Claude Shannon

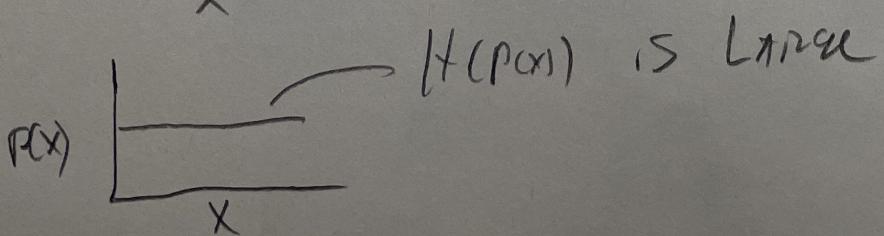
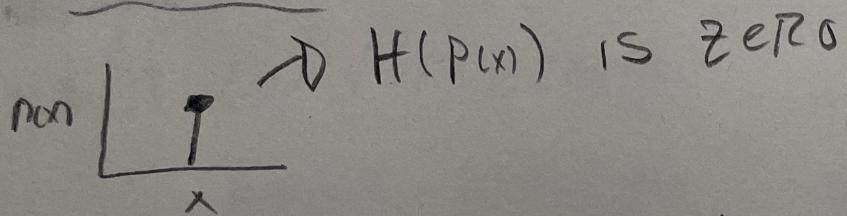


4) SELF INFORMATION: $I(x) = -\log P(x)$

"The" $\rightarrow P(X = "The")$ is fairly high, $I(x)$ is low

"Russia" $\rightarrow P(X = "Russia")$ is lower $I(x)$ is high

ENTROPY:



$$H(P) = E_{x \sim P(x)} [I(x)]$$

$$= - \sum_{x \in X} P(x) \log P(x)$$

Divergence Between two Distr. Btwn

$P(x), Q(x)$

KULLBACK-LEIBLER Divergence

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

MAXIMUM LIKELIHOOD

coin FLIP:

R.V. : $y \in Y$ $Y = \{0, 1\}$

$$\text{PMF} : P(y) = \begin{cases} \theta & y=0 \\ 1-\theta & y=1 \end{cases} \quad \text{"Bernoulli"}$$

EXPERIMENT consists of M -TRAILS

$$D = \{y^{(1)}, \dots, y^{(M)}\}$$

Likelihood is the joint distribution over our DATA:

$$- P(y^{(1)}, \dots, y^{(M)}) \rightarrow \text{Joint O. S. B.} \sim$$

↳ for EACH trial, i , we are drawing $y^{(i)}$ from $\underline{P(y_i; \theta)}$

↳ EACH SAMPLE THAT WE DRAW, i.e., EACH TRIAL OUTCOME, IS INDEPENDENT FROM ALL OTHER $M-1$ TRIALS

↳ THIS ALLOWS US TO REWRITE THE JOINT D.S.

AS:

$$P(y^{(1)}, \dots, y^{(M)}) = \underline{P(y_1; \theta)} \cdot P(y_2; \theta) \cdot \dots \cdot P(y_M; \theta)$$

MLE Continued . . .

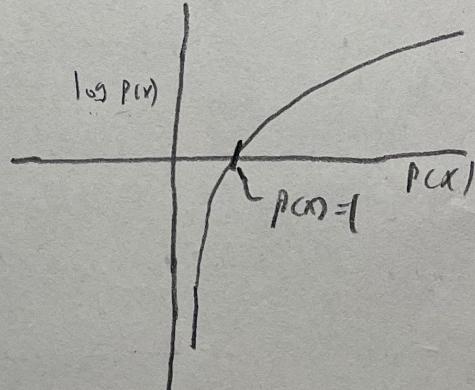
$$P(y^1, \dots, y^m) = \prod_{i=1}^m P(y^{(i)}; \theta)$$

(*) Learning Problem is at Carle to
Find the "Best" θ)

(*) $\hat{\theta}$ is the most likely θ
given on set of observations
in $D \rightarrow D = \{y^{(1)}, \dots, y^{(n)}\}$

$$(D) \hat{\theta} = \operatorname{ARGMAX}_{\theta} \prod_{i=1}^m P(y^{(i)}; \theta)$$

(*) Numerical constraint:



$$(D) \hat{\theta} = \operatorname{ARGMAX}_{\theta} \sum_{i=1}^m \log P(y^{(i)}; \theta)$$

(*) EQUIVALENT TO THE
ABOVE BECAUSE OF THE FACT THAT
 $P(x)$ AND $\log P(x)$ - change monotonically!

$$(D) \boxed{\hat{\theta} = \operatorname{ARGMIN}_{\theta} \sum_{i=1}^m -\log P(y^{(i)}; \theta)}$$

BAYES Rule:

↳ PRODUCT Rule of Probabilities

$$\begin{aligned} P(X,Y) &= P(X|Y)P(Y) \\ &= P(Y|X)P(X) \end{aligned}$$

↳ BAYES Rule:

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

$$\text{Posterior} = \frac{\text{Likelihood} \cdot \text{Prior}}{\text{EVIDENCE}}$$