TODAY

ENTROPY

- Yesterday: We claimed that $H(x)$ = expected number of bits needed to convey $X$ to reciver who knows $P_x$.

- Was a lie:

  - Problem: $H(x) + H(y) \neq H(x,y)$ when

    $X, Y$ independent

    auording to definition above

  - Example: $X = 0$    w.p.    .99

    $= 1$    w.p.    .01

  - Still need one bit to convey $X$

  - But $H(X_1 \dots X_{100})$ where $X_1 \dots X_{100}$ i.i.d.

    $\sim P_x$

    is much less (Why?)

# Correct "Operational" Definition

many
copies
of X

· $H(X) \triangleq$ amortized expected # bits to convey $\cancel{\cancel{X}} X$

ie, $X_1, \ldots, X_n$ where $X_1 \ldots X_n$ i.i.d. $\sim \overset{P}{\cancel{\cancel{X}}}$
   convey

amortized $= \frac{1}{n}$ (# bits)  ;  take $\lim\limits_{n \to \infty}$.

· Let $E_n, D_n$ be functions s.t.

$$E_n: \Sigma^n \to \{0,1\}^*$$

$$D_n: \{0,1\}^* \to \Sigma^n$$

$$\forall \underline{x} \in \Sigma^n \quad \cancel{E} \quad D_n(E_n(\underline{x})) = \underline{x} \quad ; \quad E_n \text{ } \underline{\text{prefix-free}}$$

Then $H(X) \triangleq \lim\limits_{n \to \infty} \frac{1}{n} \left\{ \underset{\underline{X} = (X_1 \ldots X_n) \sim P_X^n}{\mathbb{E}} \left[ |E_n(\underline{x})| \right] \right\}$

· Prefix-Free : $\cancel{E} \; \forall \; \underline{x}, \underline{y} \quad E(\underline{x}) \; \underline{\text{not}} \; \text{prefix of}$

$$E(\underline{y}).$$

·

" Can we compute $H(X)$ given $P_X$ ? " (in finite time)

( -studied in Inf. Theory as "Single-letter Characterization")

Yes ... as we will see ...

$$ \underline{\hspace{3cm}} \quad X \quad \underline{\hspace{3cm}} $$

• Suppose $X \sim \text{Bern}(p)$       (Think $p$ small)

$$\left[ \begin{array}{l} X = 0 \quad \text{w.p. } 1-p \\ \phantom{X} = 1 \quad \text{w.p. } p \end{array} \right]$$

## Potential Compression of $X_1, \ldots, X_n$

$E(X_1 \ldots X_n) = \quad (R ; i)$

$R = \Sigma X_i$ ;    $i = $ index of $X_1 .. X_n$

among $\binom{n}{R}$ strings of length $n$ with $R$ ones.

Length of compression $= \underbrace{\log n}_{\text{to convey } R} + \log \binom{n}{R}$

But how large is $R$?

"Chernoff Bounds": $\quad E[X_i] = p \qquad E[\sum X_i] = np$  (4)

$$\Pr\left[\,|\sum X_i - np| \geq \lambda \cdot \sqrt{n}\,\right] \leq 2 \cdot e^{-\frac{\lambda^2}{2}}$$

• So, ... ~~comp.~~ with very high prob. $\quad R \approx pn$

$$\left[\begin{array}{l} R \geq (p-\epsilon)n \\ k \leq (p+\epsilon)n \\ \text{w.p.} \geq 1 - 2^{-\epsilon^2 n} \end{array}\right]$$

• $\log \binom{n}{pn} \overset{?}{\approx} 2^{h(p) \cdot n}$

$$h(p) = -p\log_2 p - (1-p)\log(1-p)$$

Exercise: Prove this using Stirling's approx.

$$n! \approx \frac{1}{\sqrt{2\pi n}}\left(\frac{n}{e}\right)^n$$

So conclusion:
$\forall \epsilon$ for sufficient large $n$

$$H(X) \leq \frac{1}{n}\left[\,O(\log n) + (h(p)+\epsilon)\cdot n\,\right]$$

$$\Rightarrow \quad H(X) \leq h(p) . \quad \left[\text{taking limits on } \epsilon \,\&\, n\right].$$

**Theorem** : $H(X) = h(p)$ . $\left[ \text{for} \quad X = \text{Bern}(p) \right]$

**Proof** : Already seen $H(X) \leq h(p)$.

**Converse** : ① Whp $k \in \left( (p-\epsilon)n, (p+\epsilon)n \right) \approx pn$

② Even if $k$ known to receiver

need to distinguish between $\binom{n}{pn}$ possibilities

~~By~~ "~~Prefix~~" So strings of length $\leq \log \binom{n}{pn} - t$

Only occur w.p. $\leq \frac{1}{2} \cdot 2^{-t}$

Conclude that w.p. $\geq (1 - 2^{-t})$ (

Encoding length $\geq \log \binom{n}{pn} - t$ in

any valid encoding.

$\Rightarrow \quad H(X) \geq \underset{\substack{\lim \\ n \to \infty}}{2^{\,\cancel{H(p)n}}} \quad \# \quad \dfrac{h(p) \cdot n - \epsilon n - t}{n}$

$\geq \quad h(p)$ .

What about non-Bernouli $X$?

- $X \in \{1 \dots k\} = \Omega$.

  $P_1 \dots P_k$ ~~denote~~ $P_i \triangleq Pr[X = i]$

- typical string $\underline{X} \in \Omega^n$ has $\quad P_1 n \qquad 1$'s

  $\qquad\qquad\qquad\qquad\qquad\qquad\qquad P_2 n \qquad 2$'s

  $\qquad\qquad\qquad\qquad\qquad\qquad\qquad \vdots$

  $\qquad\qquad\qquad\qquad\qquad\qquad\qquad P_n n \qquad k$'s.

- All strings with these #'s $1$'s $, \dots k$'s are

  $\qquad\qquad\qquad$ equally likely.

$\Rightarrow$ Compression length

$$\approx \log \binom{n}{P_1 n \; P_2 n \, \dots \, P_k n} \pm o(n).$$

$$\approx \; \cancel{h(P_x) \cdot n} \qquad \cancel{to} \quad h(P_x) \cdot n \pm o(n)$$

$$h(P_x) \triangleq \sum_{i=1}^{R} P_i \log \frac{1}{P_i} \; .$$

<u>Theorem</u> : $H(X) = \sum_{w \in \Omega} P_x(w) \cdot \log \frac{1}{P_x(w)}$

# Conditional Entropy

$$H(Y|X) \overset{=}{\phantom{.}} \underset{\omega \sim P_X}{E}\left[H(Y|X=\omega)\right]$$

$$H(Y|X) = \sum_{\omega \in \Omega_X} P_X(\omega) \cdot H(Y|X=\omega)$$

$$= \sum_{\omega \in \Omega_X} \sum_{r \in \Omega_Y} P_X(\omega) \cdot P_{Y|X=\omega}(r) \cdot \log \frac{1}{P_{Y|X=\omega}(r)}$$

$$= \sum_{(\omega, r)} P_{XY}(\omega, r) \cdot \lg \frac{P_X(\omega)}{P_{XY}(\omega, r)}$$

$Y \in \Omega_Y \qquad X \in \Omega_X$

**Exercise :** Verify $\qquad H(X,Y) = H(X) + H(Y|X)$.

**Exercise :** Show that $X^{\#} \in \Omega^n$ can be $\underset{\wedge}{\overset{\text{always}}{\text{compressed}}}$

~~deterministically~~ to length $\overset{\leq}{\phantom{.}} H(X) \cdot n + \epsilon n$ so that

decompression error $\leq \exp\left(-f(\epsilon, R) \cdot n\right)$

where $\quad f(\epsilon, R) \to > 0$ for every $\epsilon > 0$ & $R$.

(more on this next.)

- So far we have talked about expected length of compression. Can we get "always"?

$$[\text{must allow error - why?}]$$

$$[\text{what if we allow error} \to 0].$$

- Asymptotic Equipartition Principle :

- $\forall$ distribution $P_X$ on finite $\Omega$,

- $\forall \epsilon$ for sufficiently large $n$

$$\exists \ S \subseteq \Omega^n \quad \text{s.t.}$$

• $\forall \ \underline{\omega} \in S$ , $\Pr\left[X = \underline{\omega}\right] \in \left[\frac{1}{|S|^{1+\epsilon}}, \frac{1}{|S|^{1-\epsilon}}\right]$

  $\underbrace{\hspace{4cm}}_{\text{nearly uniform on } S}$

• $\Pr\left[X \notin S\right] \leq \epsilon$.

$$[\text{\# every distribution looks like uniform}]$$

- Exercise : Prove $|S| \approx 2^{H(X) \cdot n}$
- Exercise : Derive AEP
• derive : "always with small error" compression from above ⊠