

① Shannon info / entropy

- should I email you that class is canceled?

- $X = 1$ class canceled
 $= 0$ not canceled

- $P(X=1) = \frac{1}{1024}$ (rarely canceled)

- what is the "information content" of that event?

- $h(X=1) = \log_2 \frac{1}{P(X=1)} = 10 \text{ bits}$ (information)

- $h(X=0) = \log_2 \frac{1}{P(X=0)} \approx 0.0015 \text{ bits}$

- Not much info in $X=0$, so only email if $X=1$

- information = "surprise"

- should you email me to ask if class is canceled?

- How much info do you expect to gain?

- $\mathbb{E}[h(X)] = \sum_x P(X) \log_2 \frac{1}{P(X)} \equiv H(X)$ (entropy)

- in this case, the binary entropy function

$$H_2(p) \triangleq p \log_2 \frac{1}{p} + (1-p) \log_2 \frac{1}{1-p} \\ \equiv H(X), \text{ for } X \sim \text{Bernoulli}(p).$$

- $H_2\left(\frac{1}{1024}\right) = \frac{1}{1024} \times 10 + \frac{1023}{1024} \times 0.0015 \approx 0.01 \text{ bits}$

- so, not much expected info from asking.

- $\frac{\partial}{\partial p} H_2(p) = 0 \iff p = 1-p = 0.5$

- maximum expected info = maximum uncertainty

- entropy = "uncertainty"

- $H_2(0.5) = 1 \text{ bit}$; binary variable is worth at most 1 bit

- more generally, $X \in \{1, \dots, K\}$, $P = P_1 \dots P_K$

$$H(X) \equiv H(P) = \sum_x P_x \log_2 \frac{1}{P_x}$$

$$\max_p H(p) = \log_2 K \text{ bits} = \text{"raw bit content of } X" \equiv H_0(X)$$

② "Game of Submarine" (Mackay Chap 4.1)

• sub is somewhere in 64-cell grid

• $P(Z=k) = \frac{1}{64}$

• $P(\text{guess correctly on 1st try}) \equiv P(X_1=1) = \frac{1}{64}$

• info gained: $h(X_1=1) = h(Z=k) = \log_2 64 = 6 \text{ bits}$

• if we miss: $h(X_1=0) = h(Z \neq k) = \log_2 \frac{64}{63} \approx 0.023 \text{ bits}$

• $P(\text{guess correctly on 2nd try}) \equiv P(X_2=1 | X_1=0) = \frac{1}{63}$

• total info gained after 32 misses

$$h(X_1=0 \dots X_{32}=0) = \log_2 \frac{64}{63} + \log_2 \frac{63}{62} + \dots + \log_2 \frac{33}{32}$$

$$= \log_2 \frac{64}{32} = \log_2 2 = 1 \text{ bit}$$

• total info gained after hit on 33rd attempt

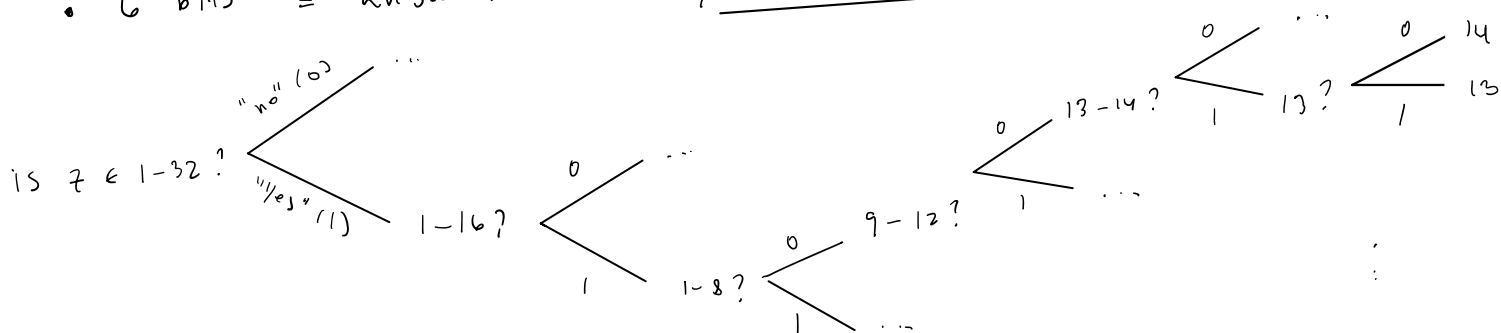
$$h(X_1=0, \dots, X_{32}=0, X_{33}=1) = 1 + \log_2 32 = \underline{6 \text{ bits}}$$

• What about hitting after 48 misses?

$$= \log_2 \frac{64}{32} + \dots + \log_2 \frac{17}{16} + \log_2 16 = 6 \text{ bits}$$

• 7 is always worth 6 bits. why?

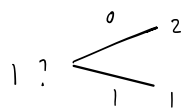
• 6 bits = answers to 6 yes-no questions



• Binary encoding of Z

Z	code
2	111111
...	...
13	110011
14	110010
...	...
64	000000

← 6-bit fixed length codes



③ Lossy compression of fixed length codes

- should not be surprising that z can be encoded into 6 bits
- $H_0(z) = \log_2 64 = 6$ bits
- in this case $H(z) = H_0(z)$ (max. uncertainty)
- can we do better if z is more predictable? $\rightarrow H(z) < H_0(z)$
- Historical context: Shannon, Bell Labs, WW2
- Encoding/decoding messages x
- "Alphabet" $x \in A_x$
- Error tolerance δ
- say $A_x = \{1 \dots 64\}$ and $p(x \in \underbrace{\{1 \dots 32\}}_{\triangleq S_\delta}) \geq 1 - \delta$
- Then only code for S_δ :

x	code
1	1111
2	1110
\vdots	\vdots
32	0000
33	\times
\vdots	\vdots
64	\times

4-bit codes

no code available

- raw bit content: $H_\delta(x) < H_0(x)$ compression
 $\log_2 32 < \log_2 64$

- error prob $\delta \rightarrow$ lossy compression
- Can we do better? yes!

- send/encode blocks $X^N = (x_1 \dots x_N)$

- say $A_x = \{1, 2, \dots, 26\}$, $p(x = x) = p_x$

- $x_i \stackrel{\text{iid}}{\sim} p(x) \rightarrow H(x_1 \dots x_N) = NH(x)$ (derive this?)

\hookrightarrow this additivity is why Shannon defined info as $\log_2 \frac{1}{p(x)}$

- X^N can be encoded in $NH_0(x)$ bits with no loss

- Source coding theorem:

X^N can be encoded in $NH(x)$ bits
 with negligible loss as $N \rightarrow \infty$

$$p(x)_{\text{typical}} = \prod_i p(x_i) \approx \prod_{x \in A} p_x^{(Np_x)}$$

$$h(x) = \log_2 \dots \approx \sum_x N p_x \log_2 p_x$$

• Define: $N_k \triangleq \sum_{i=1}^N \mathbb{1}(X_i = k)$ $= N H(p)$

• In a typical sequence X^N , there will be

$$N_k \approx N p_k$$

• $P(X^N) \approx \prod_{k \in \mathcal{A}_X} p_k^{(N p_k)}$ (typically)

• $h(X^N) \approx N \sum_k p_k \log_2 \frac{1}{p_k} = \underline{N H(X)}$

• info content of a typical X^N is $\approx N H(X)$

• X^N is ϵ -typical if $\left| \frac{1}{N} \sum_{i=1}^N h(X_i) - H(X) \right| < \epsilon$

• typical set:

$$\mathcal{A}_{\epsilon, N} \triangleq \left\{ X^N : X^N \text{ is } \epsilon\text{-typical} \right\}$$

• How likely is $X^N \in \mathcal{A}_{\epsilon, N}$? $P(X^N \in \mathcal{A}_{\epsilon, N})$?

• $R_i \triangleq h(X_i) = \log_2 \frac{1}{p(X_i)}$

• $R_i \stackrel{\text{iid}}{\sim} P(R)$, $i=1 \dots n$

• $\mathbb{E}[R] = \mathbb{E}[h(X)] = H(X)$

• $\lim_{N \rightarrow \infty} P\left(\left|\frac{1}{N} \sum_{i=1}^N R_i - \mathbb{E}[R]\right| < \epsilon\right) \rightarrow 1$
(Law of Large Numbers)

• " X^N is almost surely in the typical set for large N "

• How big is $\mathcal{A}_{\epsilon, N}$?

• X^N is ϵ -typical if:

$$\begin{aligned} H(X) - \epsilon &\leq \frac{1}{N} \log_2 \frac{1}{P(X^N)} \leq H(X) + \epsilon \\ \Leftrightarrow \frac{-N(H(X) + \epsilon)}{2} &\leq P(X^N) \leq \frac{-N(H(X) - \epsilon)}{2} \end{aligned}$$

- so, if all $x^N \sim P(x^N)$ are in $A_{\epsilon, N}$ and $P(x^N) \approx 2^{-NH(x)}$
then $|A_{\epsilon, N}| \approx 2^{NH(x)}$

- this is the Asymptotic equipartition principle (AEP)

- this is much smaller than $2^{NH_0(x)} = |A_x|^N$
(raw bit context)

- e.g.: $|A_x| = 2$ (binary), $P(X=1) = 0.4$

$$\frac{|A_{\epsilon, N}|}{|A_x|^N} \approx 2^{N(H(x) - H_0(x))} = 2^{N(0.97 - 1.0)} = 2^{-0.03N}$$

- so only 2^{-30} fraction of $N=3000$ sequences are typical

- conclusion: a fixed-length code for x^N can be
 $\log_2 |A_{\epsilon, N}| = NH(x)$ bits with negligible loss of info

④ Lossless variable-length symbol codes

- $\mathcal{X} = \{a, b, c, d\}$

- symbol code c_i for $i \in \mathcal{X}$, e.g. $c_a = 11$

i	c_i	ℓ_i	p_i
a	11	2	1/4
b	10	2	1/4
c	00	2	1/4
d	01	2	1/4

$abcd \rightarrow \underbrace{11}_{c_a} \underbrace{10}_{c_b} \underbrace{00}_{c_c} \underbrace{01}_{c_d}$

- must be prefix codes to be uniquely decodable

- $\mathbb{E}[\ell] = \sum_i p_i \ell_i$ "expected code length"

- for unequal probabilities we can do better, e.g.

i	c_i	ℓ_i	p_i
a	0	1	1/2
b	100	3	1/6
c	101	3	1/6
d	111	3	1/6

$$\mathbb{E}[\ell] = \frac{1}{2} + \frac{1}{6} \cdot 3 \cdot 3 = \frac{1}{2} + \frac{9}{6} = \frac{1}{2} + \frac{3}{2} = \frac{3}{1} = 3$$

bits

(vs. 2 bits)

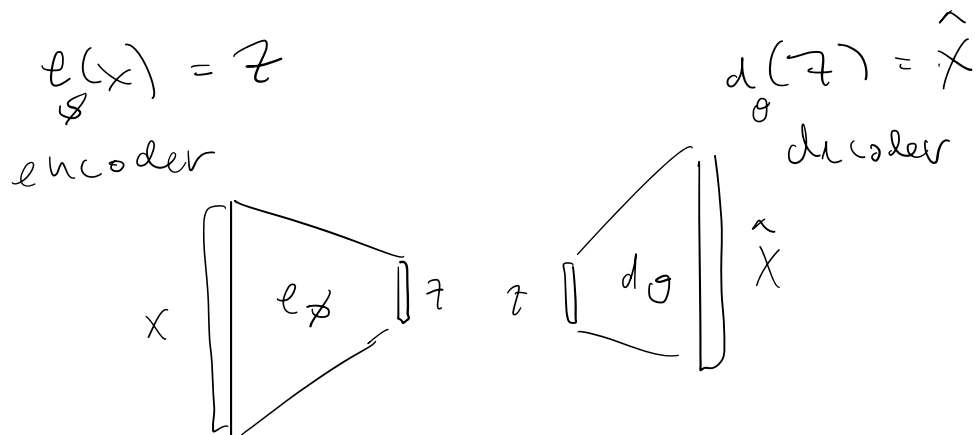
- Theorem: $\mathbb{E}[\ell] \geq H(X)$

$\mathbb{E}[\ell] = H(X)$ only if $\ell_i = \log_2 1/p_i$

- So again: $NH(X)$ bits (on average) to encode X^N

⑤ Learning/modeling as compression

- How do we learn an optimal coding scheme?
- Source coding theorem only says one exists
- Example: autoencoders



- $\hat{\theta}, \hat{\theta} \leftarrow \underset{\theta, \theta}{\operatorname{argmin}} \frac{1}{n} \sum_{i=1}^n \ell(x_i, d_{\theta}(e_{\theta}(x_i)))$
for some loss $\ell(\dots)$
- "bits back" and VAEs (save for later)

⑥ Info, entropy, etc. with dependent variables

- Is the sub in this search cell?

$$z=1 \quad \text{yes} \quad \pi(z=1) = \pi$$

$$z=0 \quad \text{no} \quad \pi(z=0) = 1-\pi$$

- We can take magnetometer measurements x

- should we? Do we expect to gain info about z ?

- Assume we know



$$\pi^* \equiv \pi^*(z=1 | x=x) = \frac{\pi P(x | z=1)}{\pi P(x | z=1) + (1-\pi) P(x | z=0)}$$

- Did x reduce our uncertainty?

$$H_2(\pi) \equiv H(z) \quad \text{vs.} \quad H_2(\pi^*) \equiv H(z | x=x)$$

- Not necessarily... e.g., if $\pi = 0.9$, $\pi^* = 0.8$

- Uncertainty about z might increase for a given x

- What about in expectation? the random variable X ,
not a value $x \rightarrow$

$$\mathbb{E}_{X \sim p(X)} [H(Z|X=x)] \equiv H(Z|X)$$

"conditional entropy"

- Fact: $H(Z|X) \leq H(Z)$
- You can show this using Jensen's inequality

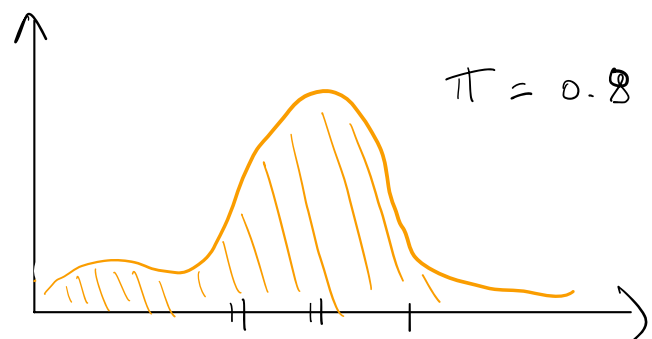
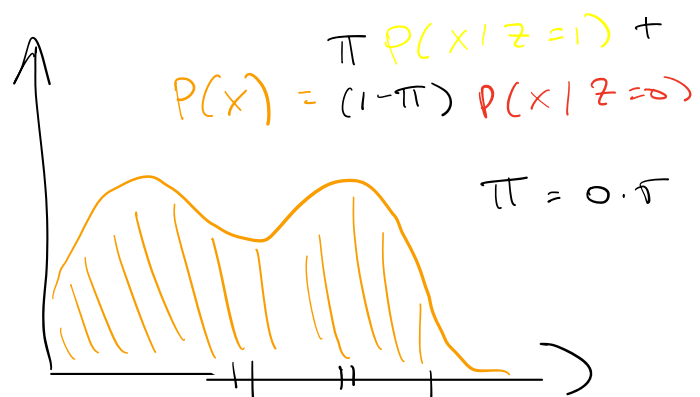
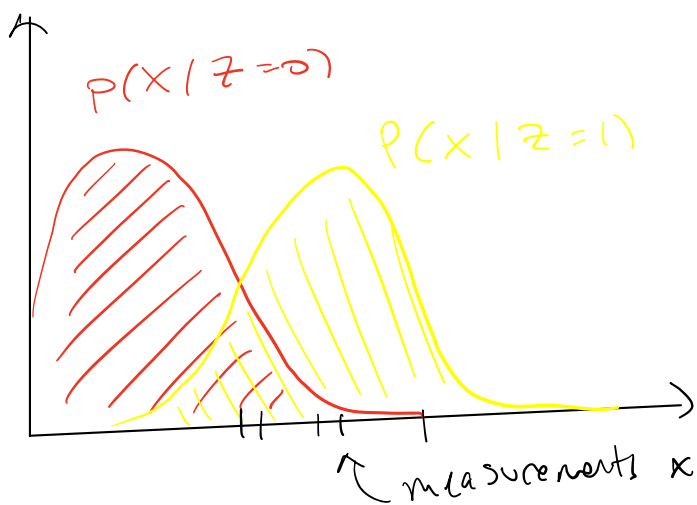
$$\mathbb{E}_x [H(p(Z|x))] \leq H(\mathbb{E}_x [p(Z|x)])$$

- So conditioning always reduces uncertainty

- But by how much?

$$H(Z) - H(Z|X) = I(Z; X) \text{ "mutual information"}$$

$$= H(X) - H(X|Z) \text{ (symmetric)}$$



$$= H(X) - \left[\pi H(X|Z=1) + (1-\pi) H(X|Z=0) \right]$$

- Mutual info is high when

- $H(Z)$ or $H(X)$ is high

- $H(X|Z)$ or $H(Z|X)$ is small at likely values of Z or X

- Another view:

$$I(X; Z) = \sum_x \sum_z p(x, z) \log_2 \frac{p(x, z)}{p(x) p(z)}$$

$$= KL \left(p(x, z) \parallel p(x) p(z) \right)$$

$$= 0 \iff X \perp Z$$

- "if $X \perp Z$, then X does not reduce uncertainty about Z "

- Yet another view:

$$I(X; Z) = H(X) + H(Z) - H(X, Z)$$

- compare to covariance:

$$\text{cov}(X, Y) = E[XY] - E[X]E[Y]$$

⑦ Which cell to search?

- Sub $z \in \{1 \dots k\}$, $\pi(z=k) = \pi_k$
- Search cell k and find it : $y_k = 1$

$$P(y_k = 1 \mid z = k) = 1$$

$$P(y_k = 1 \mid z \neq k) = 0$$

- $\arg\max_k H(z) - H(z \mid y_k)$

$$= \arg\max_k H(y_k) - \underbrace{H(y_k \mid z)}_{=0}$$

$$P(y_k = 1) = \sum_{j=1}^k \pi_j \underbrace{P(y_k = 1 \mid z = j)}_{=0 \text{ if } j \neq k, 1 \text{ if } j = k} = \pi_k$$

$$\left. \right\} = \arg\max_k H_2(\pi_k)$$

$$= \arg\max_k \pi_k$$

⑧ SEPs

$$P(y_k = 1 \mid z = k) = q_k \leq 1$$

$$P(y_k = 1 \mid z \neq k) = 1$$

$$P(y_k = 1) = \pi_k q_k$$

- $H(y_k | z) = \pi_k H_2(q_k)$

- $\arg \max_k H_2(\pi_k q_k) - \pi_k H_2(q_k)$

- $\max_{\text{possible}} I(z; y_k)$ is $\pi_k = 0.5, q_k = 1.0$

- e.g. $\pi_k = \frac{1}{8}, q_k = 1 \rightarrow H_2(\frac{1}{8}) - 0 \approx \underline{0.5435}$

vs. $\pi_{k'} = \frac{1}{2}, q_{k'} = \frac{1}{2} \rightarrow H_2(\frac{1}{4}) - \frac{1}{2} H_2(\frac{1}{2}) \approx \underline{0.311}$

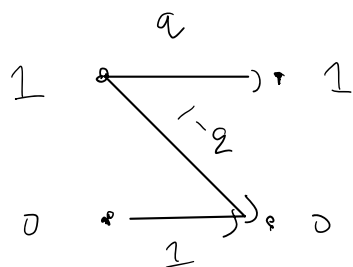
↳ more expected information from $\pi_k = \frac{1}{8}$.

- Puzzle: $q_k = q_{k'} = 0.86$ (equal!)

$$\pi_k = 0.455 < \pi_{k'} = 0.5$$

$$\begin{array}{ccc} \swarrow & !! & \searrow \\ I(z; y_k) \approx 0.685 & > & I(z; y_{k'}) \approx 0.678 \end{array}$$

- Hint: Mackay example 9.11, "z-channels"



9

"Bits back"

- optimal encoding of x is in $\log_2 \frac{1}{p(x)}$ bits
- what if $p(x)$ is intractable?
- However $p(x|z)$ is $p(x) = \sum_z p(x, z)$
- Transmit x and z
- code x according to $p(x|z)$, z according to $p(z)$
- Which z ? $\arg \max_z p(z|x)$?
- That would minimize the code length of x but not z
- Instead: sample $z \sim p(z|x)$

naive cost: $\log_2 \frac{1}{p(x|z)} + \log_2 \frac{1}{p(z)}$

- However! The receiver can decode the random bits and gain them back:

$$\log_2 \frac{1}{p(x|z)} + \log_2 \frac{1}{p(z)} - \log_2 \frac{1}{p(z|x)}$$

"bits back"

$$= \log_2 \frac{1}{p(x)} \leftarrow \text{optimal rate!}$$

- This is the "bits back" argument

(10)

Minimum Description Length (MDL)

- "Bits back" motivates latent variable models
- But which model?
- $P(X; M_1)$ vs. $P(X; M_2)$
- MDL principle: minimize the joint length of the data and the model

$$\min_M L(X; M) + L(M) \\ \quad \quad \quad \underbrace{\hspace{1cm}} \\ \quad \quad \quad = h(X; M)$$

- This can be understood as defining a prior over models:

$$P(M) = 2^{-L(M)}$$

and doing MAP:

$$\hookrightarrow \max_M \log P(X; M) + P(M)$$

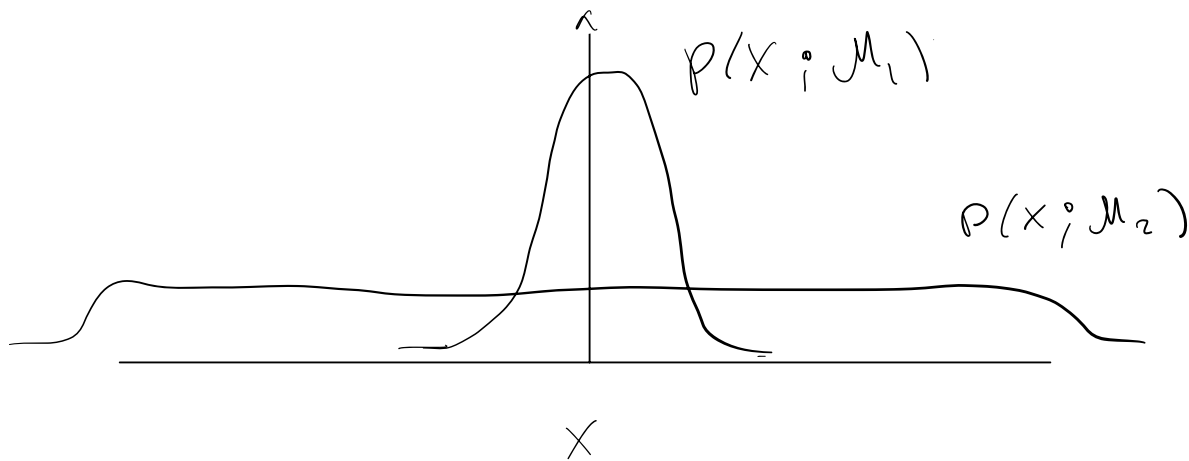
- for two models, the Bayes factor:

$$\frac{P(M_1 | X)}{P(M_2 | X)} = \frac{P(X; M_1)}{P(X; M_2)} \frac{P(M_1)}{P(M_2)}$$

↑
evidence

⑪ Occam's Razor and model evidence

- Mackay argues that model evidence alone captures the principle of Occam's Razor.
- Say M_1 is simpler than M_2
- M_2 can fit to more data but flatter spreads probability across many possible x :



⑫ Recent updates to this perspective

- Fong and Holmes (2020):

"Leave p -out" CV score:

$$S_{CV}(x_{1:n}; p) = \frac{1}{\binom{n}{p}} \sum_{t=1}^{\binom{n}{p}} \frac{1}{p} \sum_{j=1}^p S(\tilde{x}_j^t | x_{1:n-p}^t)$$

Theorem:

$$\log P(x_{1:n}) = \sum_{p=1}^n S_{CV}(x_{1:n}; p)$$

$$S(\tilde{x}_j^t | x_{1:n-p}^t) = \log P(\tilde{x}_j | x_{1:n-p}; M)$$