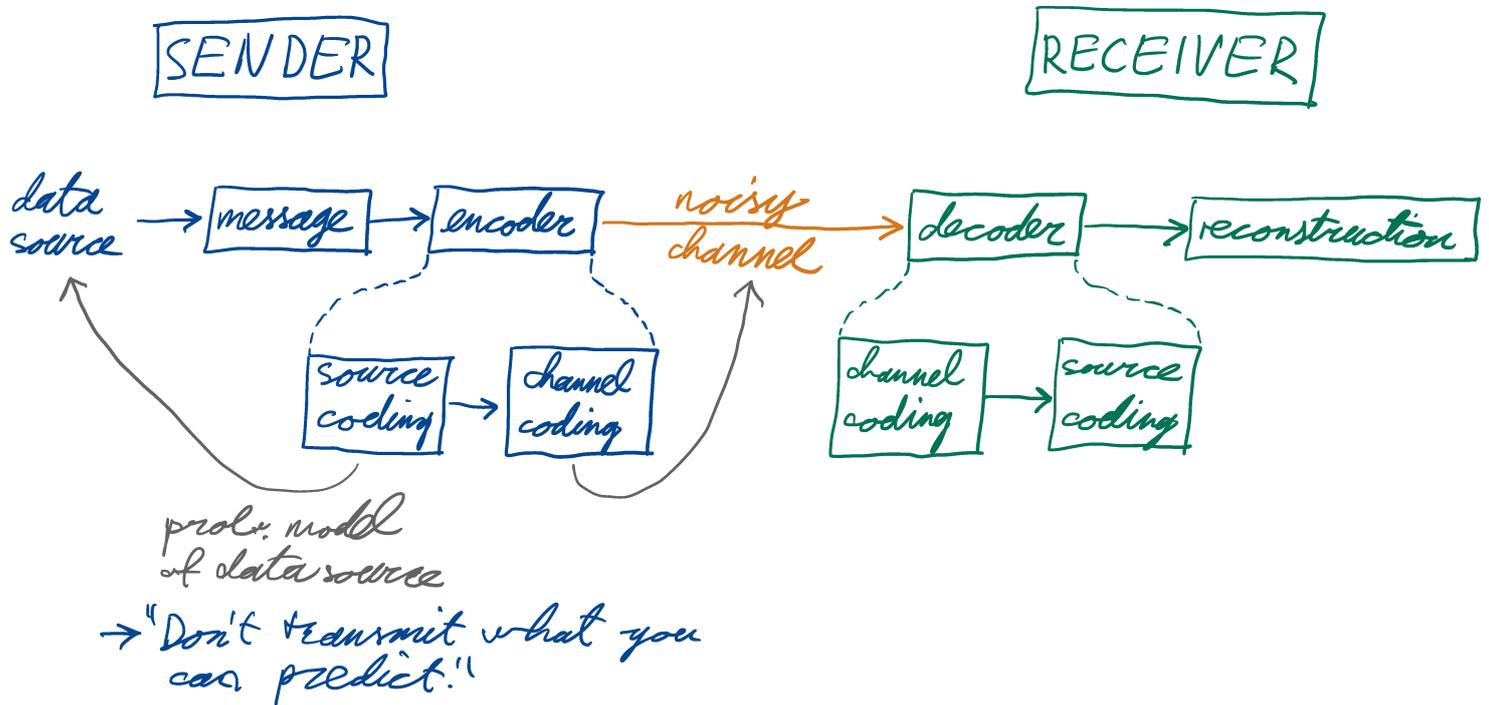


# Probabilistic Models of Data Sources

Reminder: the big picture



Qualitatively: better prob. models  $\Rightarrow$  better compression performance

## Quantifying the Modeling Error: The Kullback-Leibler Divergence

- consider general lossless compression setup:  
(i.e. no longer to symbol codes)
  - data source generates messages  $x$  with probs.  $P_{data}(x)$
  - Def: bit rate  $R(x) :=$  total no. of bits in compressed rep. of  $x$  for a given lossless compression method
- ↑  
prob. dist. over entire messages (not just single symbols)

- optimal expected bit rate <sup>(lowest theoret. possible)</sup>

$$E_{x \sim p_{data}} [R_{opt}(x)] = H(p_{data}) + \epsilon$$

consider entire set possible messages as alphabet  $\rightarrow x$  is a single symbol  
 $\rightarrow L_{opt} = H(p_{data}) + \epsilon$

$E_{p_{data}} [-\log p_{data}] < 1 \text{ bit}$   
 (typically irrelevant)

- reminder: to reach optimal expected bit rate, a compression method has to satisfy  $R_{opt}(x) = -\log p_{data}(x) + \epsilon \forall x$

Problem: in practice, we don't know  $p_{data}$

$\rightarrow$  distinguish

$p_{data}$

vs.

$p_{model}$

↑  
true dist. of the data

$\hookrightarrow$  we don't know this

$\hookrightarrow$  but we may have samples from  $p_{data}$

$\Rightarrow$  we can evaluate empirical averages, to estimate true expectation values

↑  
(for now:) assume we can explicitly evaluate

$p_{model}(x) \forall x$

$\Rightarrow$  optimal lossless compression code:  $R(x) = -\log p_{model}(x) \forall x$

$\Rightarrow$  expected bit rate

$$E_{x \sim p_{data}} [R(x)] = E_{x \sim p_{data}} [-\log p_{model}(x)] = H(p_{data}, p_{model})$$

"cross entropy"

$\rightarrow$  we have to minimize  $H(p_{data}, p_{model})$  over parameters of  $p_{model}$ .

⇒ Overhead due to  $p_{\text{model}} = p_{\text{data}}$ :

## Kullback - Leibler Divergence

$$D_{\text{KL}}(p_{\text{data}} \parallel p_{\text{model}}) = \underbrace{H(p_{\text{data}}, p_{\text{model}})}_{\text{actual bit rate}} - \underbrace{H(p_{\text{data}})}_{\text{theoret. lower bound}}$$

$$= \mathbb{E}_{x \sim p_{\text{data}}} \left[ \log \frac{p_{\text{data}}(x)}{p_{\text{model}}(x)} \right]$$

- Problem 3.1(c): prove that  $D_{\text{KL}}(p \parallel q) \geq 0$  ("Gibbs Theorem")

So far:  $p_{\text{model}}(\underline{x}) = (p(k))^{\prod_{i=1}^k} p(x_i)$

$\uparrow$   
 $x^*$

$\underbrace{\hspace{10em}}$   
prob. of a single symbol

→ assumed that symbols are "i.i.d.":

independent & identically distributed

difficult part: we haven't been able to model correlations between symbols

↖ easy to drop this restriction:

$$p_{\text{model}}(\underline{x}) = p(k) \prod_{i=1}^k p_i(x_i)$$

→ prefix codes still work

# Interlude: Probability Theory & Random Variables

Goal: efficiently model correlations between parts of the message

- sample space  $\Omega$ : (abstract) space of "states of the world"
  - event  $E \subset \Omega$ : "event  $E$  occurs" = "the world is in a state  $\omega \in E$ "
  - probability measure:  $P: \Sigma \rightarrow [0, 1]$ 
    - $\hookrightarrow P(\Omega) = 1$
    - $\hookrightarrow P(\emptyset) = 0$
    - $\hookrightarrow P(\bigcup_{i=1}^{\infty} E_i) = \sum_{i=1}^{\infty} P(E_i)$  if  $\{E_i\}$  are pairwise disjoint
    - $\hookrightarrow P(\bigcup_{i=1}^k E_i) = \sum_{i=1}^k P(E_i)$
- $\Sigma$ -algebra on  $\Omega$   
(set of subsets of  $\Omega$ )

Remark: for continuous states,  $P(\{\omega\})$  for a single  $\omega \in \Omega$  typically doesn't make much sense.

e.g.:  $\Omega \in \mathbb{R}$ ,  $\omega \in \Omega$  is the arrival time of a bus

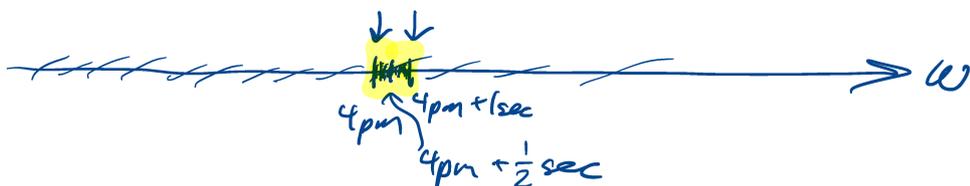
Q: with what probability that the bus arrives exactly at 4pm today

$$P(\{\text{"4pm"}\}) = 0$$

$$P([\text{"4pm"}, \text{"4pm + 1sec"}])$$

can be  $> 0$

expressed as a real number according to some standard



• Random variable:  $X: \Omega \rightarrow \mathbb{R}$

Example: "Simplified Game of Monopoly"

$$\hookrightarrow \Omega = \{ (a, b) : a, b \in \{1, 2, 3\} \}$$

value of red die  $\nearrow$   
value of blue die  $\nwarrow$

$$\hookrightarrow \Sigma = \mathcal{P}(\Omega) = 2^\Omega := \{ \text{all subsets of } \Omega \}$$

$$\hookrightarrow P(E) = \frac{|E|}{|\Omega|} = \frac{|E|}{9}$$

$\hookrightarrow$  Random variables:

• value of red die:  $X_r(a, b) = a$

• value of blue die:  $X_b(a, b) = b$

• sum of red + blue die:  $X_s(a, b) = a + b \in \{2, 3, 4, 5, 6\}$

Properties of a single random variable

$\hookrightarrow$  expectation value (of discrete r.v.)

$$\begin{aligned} \mathbb{E}_P[X] &= \sum_{\omega \in \Omega} P(\{\omega\}) X(\omega) \\ &= \sum_{x \in X(\Omega)} \underbrace{P(X^{-1}(x))}_{{= \{ \omega \in \Omega : X(\omega) = x \}}} x \end{aligned}$$

$$\mathbb{E}_P[X_r] = \mathbb{E}_P[X_b] = 2$$

$$\mathbb{E}_P[X_s] = 4$$

$\hookrightarrow$  expectation value (of a continuous r.v.)  $\swarrow$  probability density fct.

$$\mathbb{E}_P[X] = \int X(\omega) dP(\omega) = \int_{-\infty}^{\infty} x p(x) dx$$

integration measure  $\swarrow$  in this course  $\swarrow$   $p(x) \geq 0 \forall x$   
 $\sum_{-\infty}^{\infty} p(x) dx = 1$   
 $p(x)$  can be  $> 1$

$\hookrightarrow$  prob. dist. of a r.v.:  $P(X=x) = P(\{\omega \in \Omega : X(\omega) = x\})$

•  $P(X): \mathbb{R} \mapsto [0, 1], x \mapsto P(X=x)$   
"the fct.  $P(X=\cdot)$ "

## Properties of two r.v.s

↳ joint probability distribution  $X$  &  $Y$ :

$$P(X=x, Y=y) = P(\{\omega \in \Omega : X(\omega)=x \wedge Y(\omega)=y\})$$

$$P(X, Y): \mathbb{R} \times \mathbb{R} \mapsto [0, 1], \quad (x, y) \mapsto P(X=x, Y=y)$$

↳ Def: 2 r.v.s  $X$  &  $Y$  are (statistically) independent iff:

$$P(X, Y) = P(X) P(Y) \quad \text{"marginal distribution"}$$

$$\text{(i.e.: } P(X=x, Y=y) = P(X=x) P(Y=y) \quad \forall x, y)$$

- $X_r, X_b$  are independent
- $X_r, X_s$  are not independent:

$$\text{eg.: } P(X_r=1, X_s=3) = P(\{(1, 2)\}) = \frac{1}{9}$$

$$\text{but } P(X_r=1) P(X_s=3) = \frac{1}{3} \times \underbrace{P(\{(1, 2), (2, 1)\})}_{2/9} = \frac{2}{27} \neq \frac{1}{9}$$

## Conditional Probability Distribution

↳ for events: "conditional prob. of event  $E_2$  given event  $E_1$ "

$$P(E_2 | E_1) = \frac{P(E_1 \cap E_2)}{P(E_1)}$$

$$\left( \Rightarrow \underbrace{P(\neg E_2 | E_1)}_{\neg E_2 = \Omega \setminus E_2} = \frac{P(E_2 \setminus E_1)}{P(E_1)} \Rightarrow P(E_2 | E_1) + P(\neg E_2 | E_1) = 1 \right)$$

↳ for r.v.s:  $P(Y | X) = \frac{P(X, Y)}{P(Y)}$

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

"What is the prob. of  $Y$  being  $y$  if I already know that  $X=x$ ?"

→ if  $X, Y$  are indep:  $P(X, Y) = P(X) P(Y)$

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

↑  
For indep. v.v.s

Important: writing  $P(Y|X)$  does not imply causality, i.e. it does not mean that  $X$  is the cause of  $Y$ .

→ even if  $X$  is the cause of  $Y$ , we can still calculate:

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

"Bayesian inference"

Chain rule of probabilities: (follows directly from def. of cond. probs)

$$P(X, Y) = P(X) P(Y|X) = P(Y) P(X|Y)$$

$$P(X, Y, Z) = P(X) P(Y|X) P(Z|X, Y) = \dots$$

Back to source coding

Problem 3.2: You'll implement a compression method for natural language; it will exploit correlations between symbols (chars).

message:  $\underline{X} = (x_1, x_2, x_3, \dots, x_k)$

$$P(\underline{X}) = P(x_1) P(x_2|x_1) P(x_3|x_1, x_2) \dots P(x_k|x_1, x_2, \dots, x_{k-1})$$

↑  
chain rule (always correct)

Issue: this general (exact) factorization of the joint distribution is not computationally feasible:

$P(X_k | X_1, X_2, \dots, X_{k-1})$  is an extremely complicated fct.

→ need ways to:

capture relevant correlations

while still

maintaining  
 • compact model representation  
 • computational efficiency

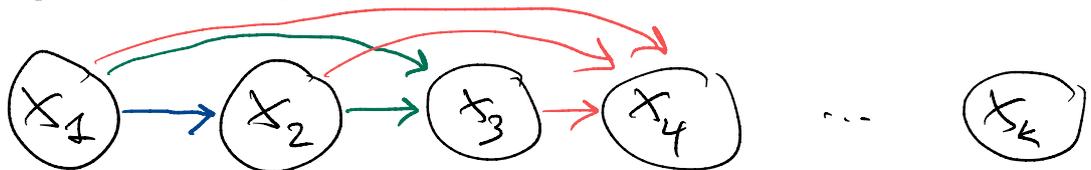
→ general strategy: enforce conditional independence:

for r.v.s  $X, Y, Z$ :  $Y$  &  $Z$  are cond. indep given  $X$  iff:

$$P(Z | X, Y) = P(Z | X)$$

(exercise:  $\Leftrightarrow P(Y, Z | X) = P(Y | X) P(Z | X)$ )

• general chain rule:



$$P(\underline{X}) = P(X_1) P(X_2 | X_1) P(X_3 | X_1, X_2) P(X_4 | X_1, X_2, X_3) \dots$$

### 3 possible simplifications

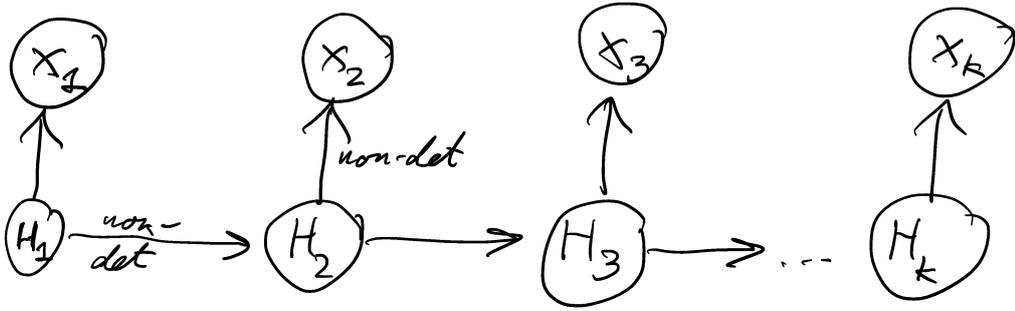
(a) if  $X_i$  are generated in sequence by some memoryless process → Markov Process

→ assumes cond. indep. of  $X_i$  & all  $X_j$  with  $j < i-1$  given  $X_{i-1}$



$$P(\underline{X}) = \prod_{i=1}^k P(X_i | X_{i-1})$$

## b) Hidden Markov Model

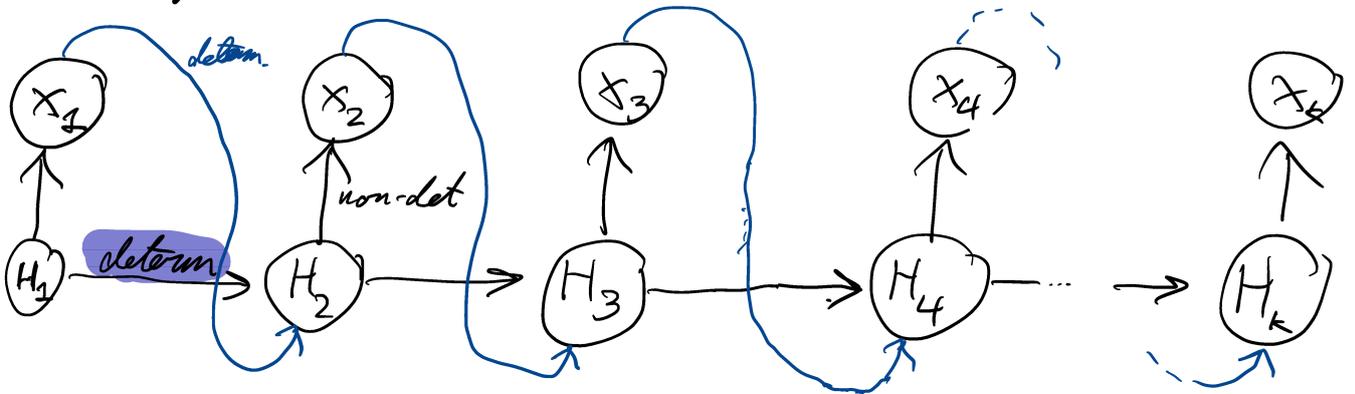


(exercise: this can capture long-range correlations

→ in a model like this  $X_1$  &  $X_3$  do not have to be cond. indep given  $X_2$ )

→ difficult for compression (e.g. using bits back coding)

## c) Autoregressive Model

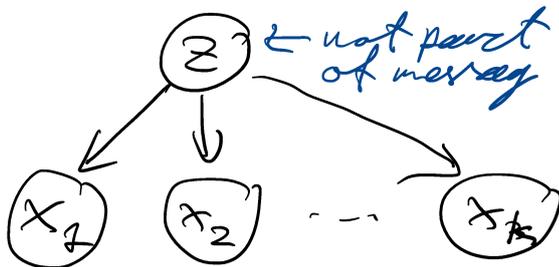


$H_{i+1}$  is a determ fun. of  $H_i, X_i$

😊 can capture long-range correlations  
(i.e.  $X_2, X_4$  are not cond. indep given  $X_3$ )

😞 hard to parallelize

→ next video: latent variable models



😊 can capture correlations between  $X_i$ 's  
😊 can be parallelized  
😊 how to use this for compression → bits-back coding