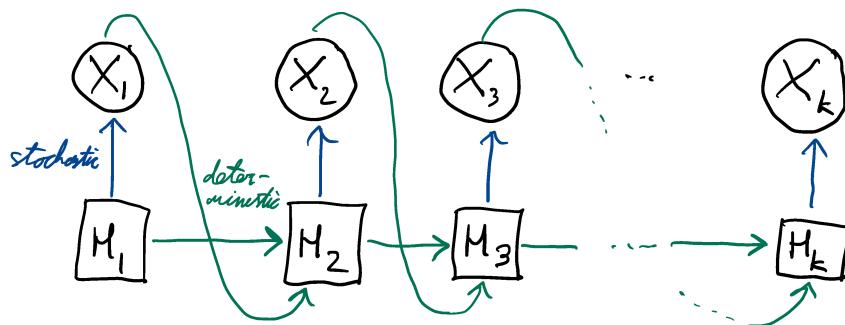


Bits-Back Coding With Latent Variable Models

Last video: • overview of probability theory & random vars.

- modeling error: KL-Divergence
→ need to model correlations! $H_p(X) + H_p(Y) \geq H_p((X,Y))$ Problem 4.2 & 6.2
- autoregressive models



- ☺ can model long (-ish) range correlations
- ☺ compact to store
- ☹ struggle with very long correlations
- ☹ not parallelizable

This video: • latent variable models
• Bayesian inference
• Bits-back coding

Latent Variable Models

Consider these hypothesized news headlines.

Parliament Votes on New Labor Bill.

} topic
"politics"

Labor Union Votes to Extend Strikes.

Soccer Player Scores First Goal Since Joining New Team.

Guest Team is Leading by One Goal.

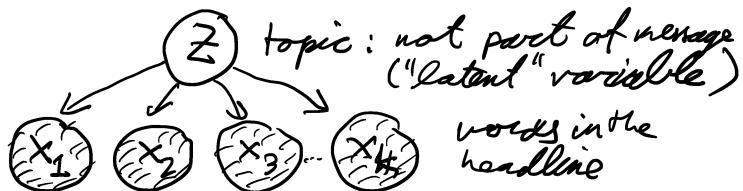
} topic
"sports"

Observation: words within a headline appear to be correlated:

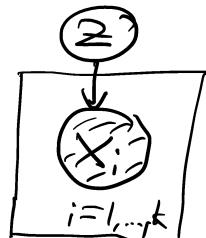
consider two positions $i \neq j$; claim: X_i & X_j are not stat. indep
 ↑
 words at these positions

$$\text{e.g.: } P(X_i = \text{"Goal"}, X_j = \text{"Team"}) > P(X_i = \text{"Goal"}) P(X_j = \text{"Team"})$$

Model of the "Generative Process"



Abbreviated depiction



These pictures denote a joint prob. dist. that factorizes as follows:

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

("Topic Model", e.g. LDA: Blei & Ng 2003, before ingesters Pitcock et al. 2000)

⇒ marginal dist of the message \underline{X} :

$$P(\underline{X}) = \sum_z P(\underline{X}, Z=z) = \sum_z \left(P(Z=z) \prod_{i=1}^k P(X_i|Z=z) \right)$$

Claim: this model can capture correlations like

$$P(X_i = \text{"Goal"}, X_j = \text{"Team"}) > P(X_i = \text{"Goal"}) P(X_j = \text{"Team"})$$

→ proof exercise

Data Compression With Latent Variables Models

$$P(\underline{X}, \underline{Z}) = P(\underline{Z}) P(\underline{X} | \underline{Z})$$

message ↑ latent var ↑

ideally we would like to compress \underline{X} with this model

→ Problem: we don't know value of \underline{Z}

Problem Set 5: implement & compare 3 compression methods for l.v.m.

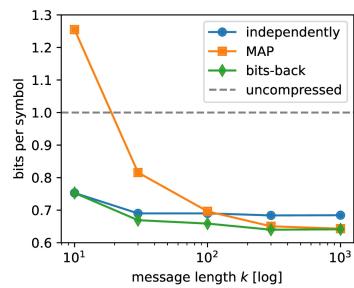
→ Problem 5.2: treat X_i as independent

→ ignore correlations

$$\rightarrow E_p[R(\underline{x})] = \sum_{i=1}^n H_p(X_i) \geq H_p(\underline{X})$$

→ Problem 5.3: MAP-estimate

$$\rightarrow E_p[R(\underline{x})] = -\log P(\underline{Z}=\underline{z}^*) + H_p(\underline{X}|\underline{Z}=\underline{z}^*)$$



→ Problem 5.4: bits-back coding

$$\rightarrow E_p[R_{\text{nat}}(\underline{x})] = H_p(\underline{X})$$

Naive approach: MAP-estimate

$$P(\underline{X}, \underline{Z}) = P(\underline{Z}) P(\underline{X} | \underline{Z})$$

→ Idea: • encode some value \underline{z} for \underline{Z} using $P(\underline{Z})$ & transmit it

• then encode \underline{X} using $P(\underline{X} | \underline{Z}=\underline{z})$ ($= \prod_{i=1}^n P(X_i | Z=z_i)$)

⇒ decoder can • decode \underline{z} using $P(\underline{Z})$

• decode \underline{X} using $P(\underline{X} | \underline{Z}=\underline{z})$

$$\begin{aligned} \text{bit rate: } R^{(z)}(\underline{x}) &= -\log P(\underline{Z}=\underline{z}) - \log P(\underline{X}=\underline{x} | \underline{Z}=\underline{z}) \\ &= -\log P(\underline{X}=\underline{x}, \underline{Z}=\underline{z}) \end{aligned}$$

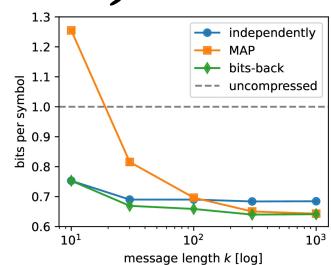
$$\text{choose } \underline{z}^* := \arg \min_{\underline{z}} R^{(z)}(\underline{x}) = \arg \max_{\underline{z}} P(\underline{X}=\underline{x}, \underline{Z}=\underline{z})$$

↑
"maximum a-posteriori" (MAP)
estimate of \underline{Z}

overhead over theoret. bound:

$$R^{(z^*)}(\underline{x}) - (-\log P(\underline{X}=\underline{x})) = -\log P(\underline{X}=\underline{x}, \underline{Z}=\underline{z}^*) + \log P(\underline{X}=\underline{x})$$

$$= -\log P(\underline{Z}=\underline{z}^* | \underline{X}=\underline{x}) \stackrel{> 0}{\leftarrow} \text{posterior distribution}$$



Bayesian Inference

• model: $P(\underline{X}, \underline{Z}) = P(\underline{Z}) P(\underline{X} | \underline{Z})$

↳ know \underline{X} , don't know \underline{Z} (\Rightarrow MAP-est. method has an overhead)

↳ But: knowing \underline{X} typically reveals some information about \underline{Z}

Parliament Votes on New Labor Bill.

} topic
"politics"

Labor Union Votes to Extend Strikes.

Soccer Player Scores First Goal Since Joining New Team.

} topic
"sports"

Guest Team is Leading by One Goal.

→ However: there can still be some ambiguity about \underline{Z} (even after you know \underline{X})

Parliament Votes on Aid for Community Sports Teams. (*)

→ can only make prob. statements about \underline{Z}

$$\rightarrow \text{posterior} \quad P(\underline{Z} | \underline{X} = \underline{x}) = \frac{P(\underline{Z}) P(\underline{X} = \underline{x} | \underline{Z})}{\underbrace{P(\underline{X} = \underline{x})}_{\text{posterior}}} = \frac{P(\underline{Z}) P(\underline{X} = \underline{x} | \underline{Z})}{\sum_{z'} P(z=z') P(\underline{X} = \underline{x} | \underline{Z} = z')}$$

Remarks:

- in principle posterior distib. is known once you know $P(\underline{X}, \underline{Z})$ & \underline{X}
- in practice, however, calculating the posterior is often prohibitively expensive (\Rightarrow Lecture 7: approximate Bayesian inference)

Understanding the overhead of MAP-est. method

→ we could encode $(*)$ in two different ways

(a) $\underline{Z} = \text{"politics"}$, then we use $P(\underline{X} | \underline{Z} = \text{"politics"})$

two different compressed bit strings of same vector
→ wasted

(b) $\underline{Z} = \text{"sports"}$, then we use $P(\underline{X} | \underline{Z} = \text{"sports"})$

→ remember: overhead = $-\log P(\underline{Z} = z^* | \underline{X} = \underline{x})$

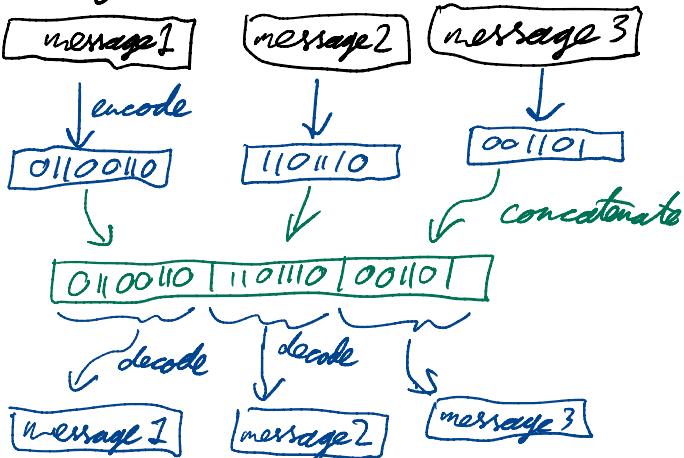
Bits-Back Coding

(Wallace 1990, Hinton & Camp 1993,
practical: BB-ANS: Townsend et al. 2019
lossy: Yang et al. 2020)

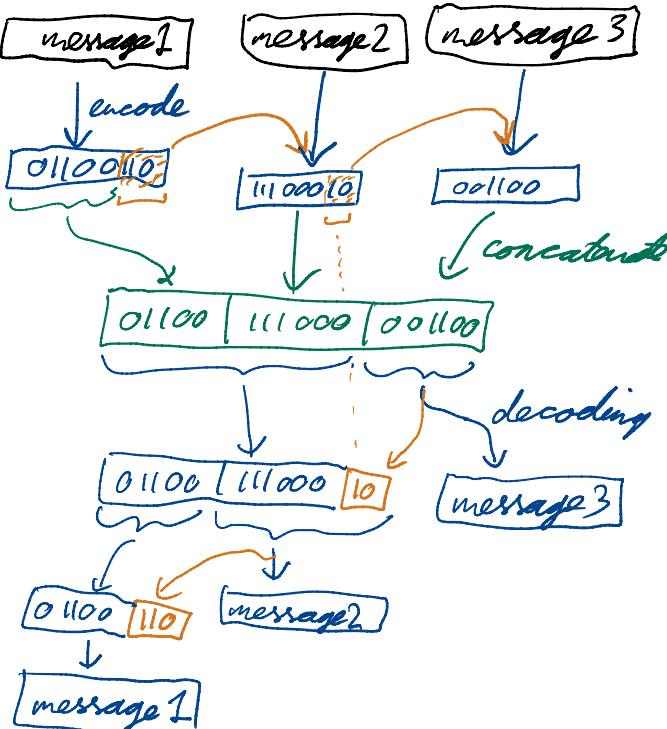
→ Idea: "piggyback" some additional message into the choice of z

Setup: communicate multiple messages (e.g. multiple image patches) over a single channel

→ usually:



→ bits-back: (operates as a stack, i.e. "last in first out")



Algorithm: "Bits-Back Coding"
subroutine encode (\underline{x} , compressed, P):

$z \leftarrow \underline{\text{decode}}$ from compressed
using $P(z | \underline{x} = \underline{x})$,
encode \underline{x} using $P(\underline{x} | z = z)$ onto compressed
encode z using $P(z)$ onto compressed
return compressed

subroutine decode (compressed, P):

$z \leftarrow \underline{\text{decode}}$ from compressed
using $P(z)$

$\underline{x} \leftarrow \underline{\text{decode}}$ from compressed
using $P(\underline{x} | z = z)$

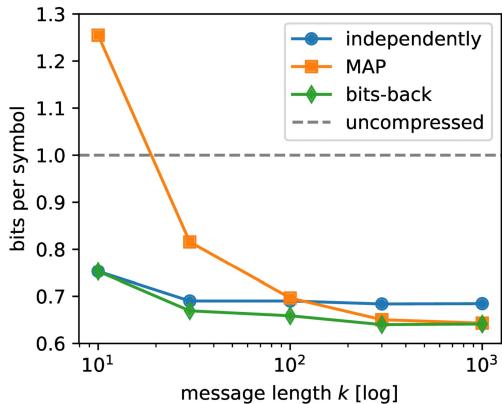
encode z onto compressed using
 $P(z | \underline{x} = \underline{x})$

return (\underline{x} , compressed)

net bit rate of bits-back coding:

$$\begin{aligned} R_{\text{net}}(\underline{x}) &= -\log P(\underline{X}=\underline{x} | \underline{Z}=\underline{z}) - \log P(\underline{Z}=\underline{z}) - (-\log P(\underline{Z}=\underline{z} | \underline{X}=\underline{x})) \\ &= -\log \frac{P(\underline{X}=\underline{x}, \underline{Z}=\underline{z})}{P(\underline{X}=\underline{x}, \underline{Z}=\underline{z})} = -\log P(\underline{X}=\underline{x}) \end{aligned}$$

⇒ bits-back coding is optimal (net).



Next steps:

- How do we encode-decode fractional numbers of bits with stuck semantics.
- What if we don't know the exact posterior?
→ Lecture 7: approximate Bayesian inference
- How can we efficiently train deep latent variable models?
→ Lecture 7 & subsequent:
 - variational expectation maximization
 - deep generative models