

Learning with symbols within neural networks: Motivating examples

Neuro-symbolic modeling



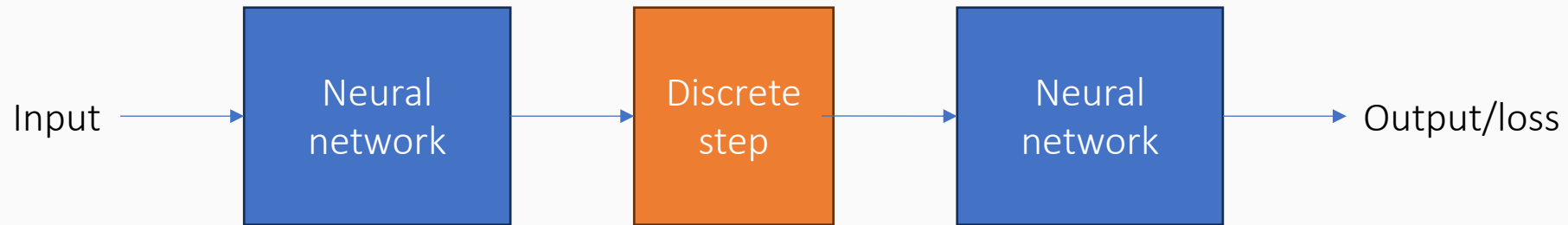
This lecture

- Motivating examples
- The straight-through estimator
- The Gumbel trick
- REINFORCE

(others if time permits)

Not all these approaches
are always applicable

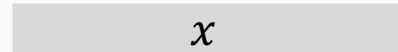
Neural networks containing discrete elements



Let's see some examples

Example 1: An autoencoder

Input



Example 1: An autoencoder

Latent representation

z

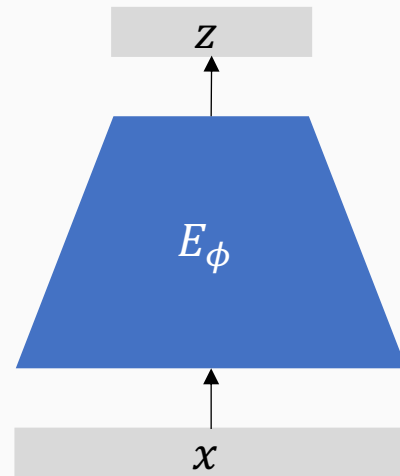
Encoder

E_ϕ

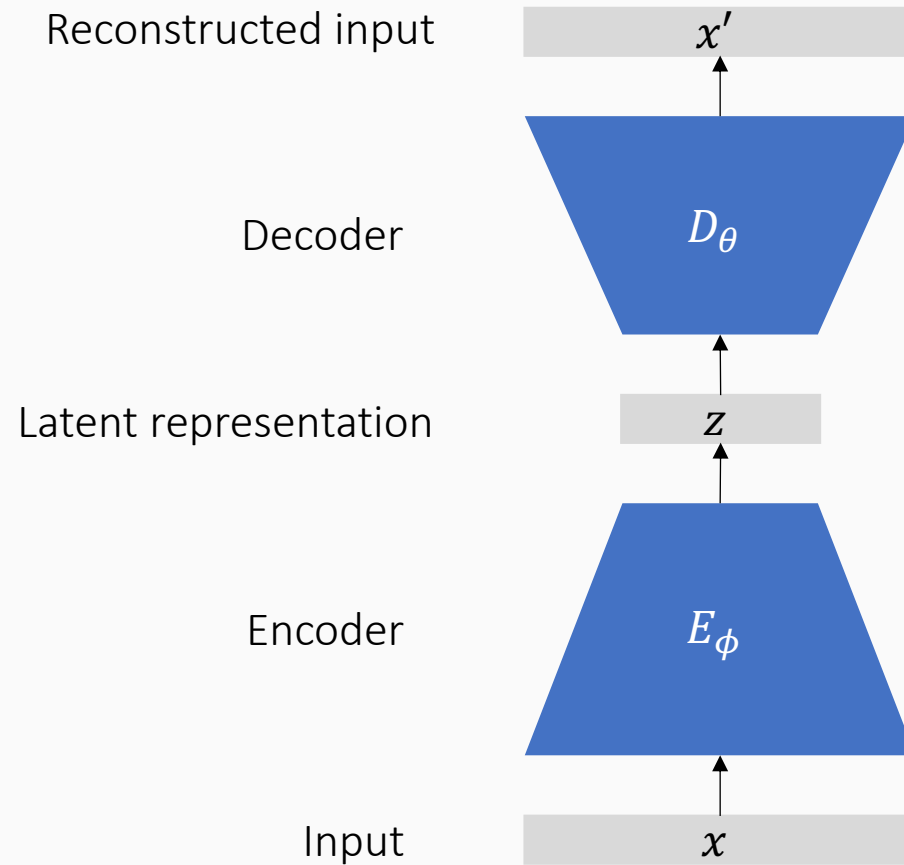
Input

x

This could be any neural network architecture, parameterized by ϕ

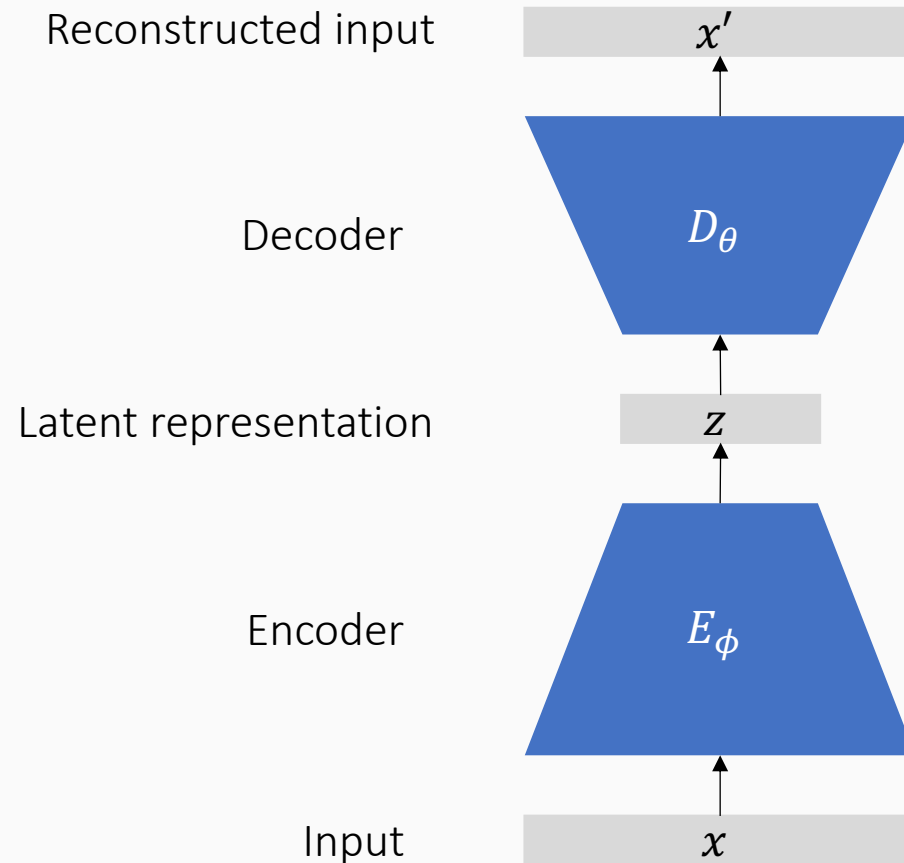


Example 1: An autoencoder



This could be any neural network architecture, parameterized by θ

Example 1: An autoencoder

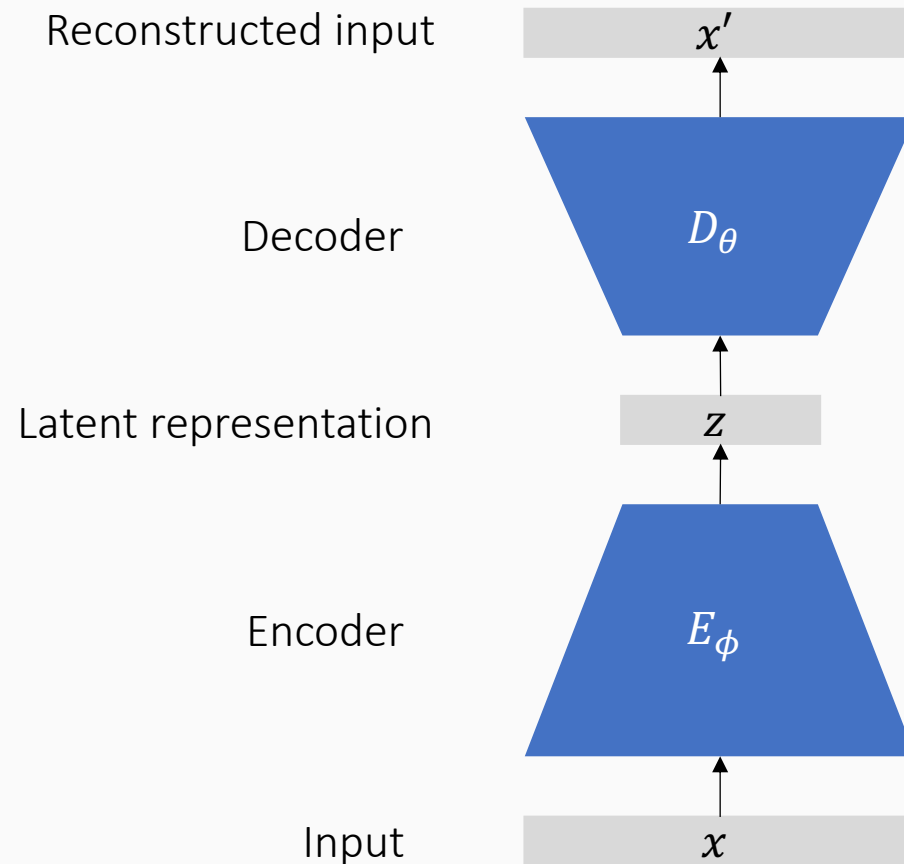


The goal of learning

To ensure that the reconstructed input x' is close to the original input x

If the encoder and decoder are differentiable and the latent representation is a real vector, we have an end-to-end differentiable function

Example 1: An autoencoder

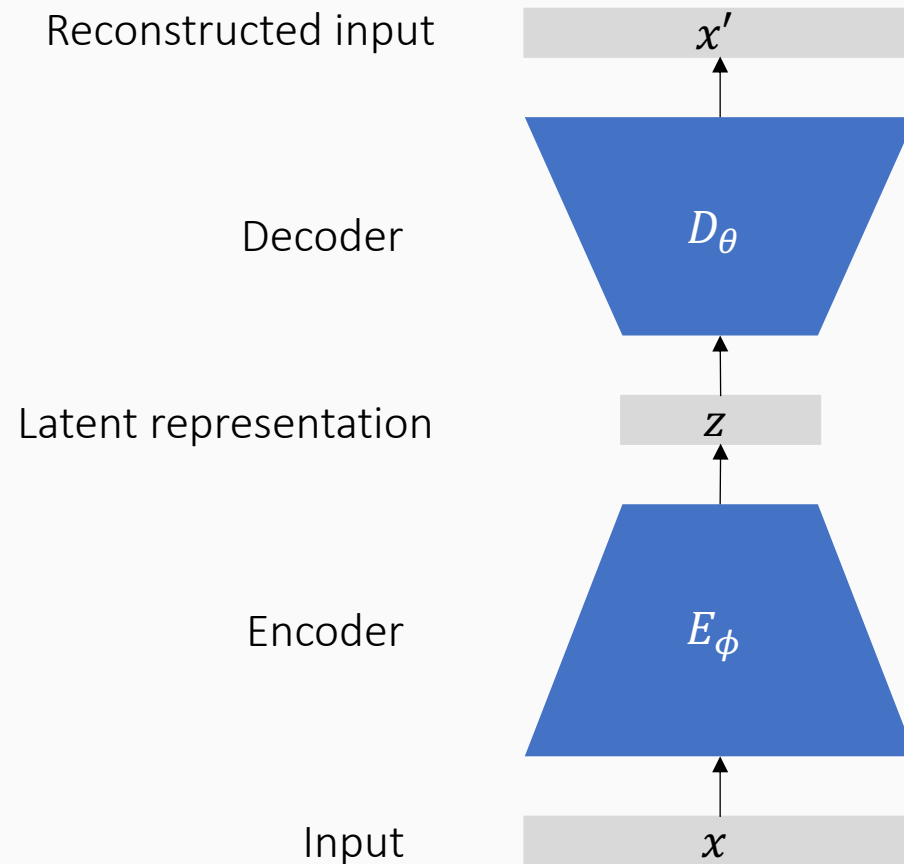


The goal of learning

To ensure that the reconstructed input x' is close to the original input x

The latent representation z can be seen as an *encoding* of the input that captures enough information to be able to reconstruct the input

Example 1: An autoencoder



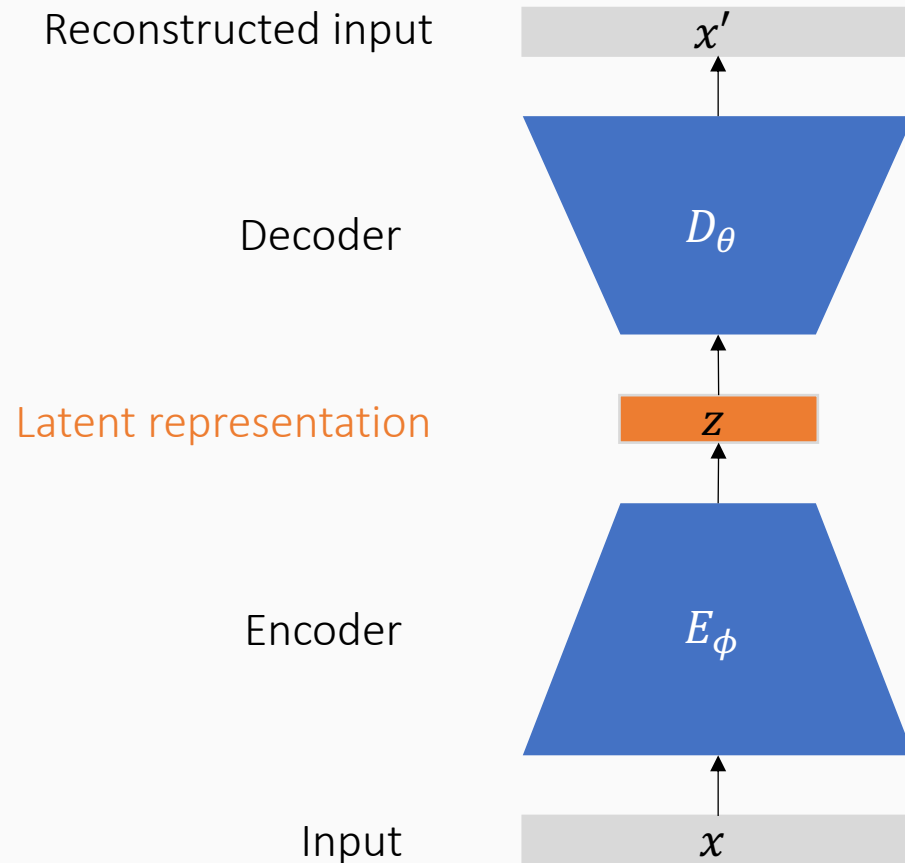
The goal of learning

To ensure that the reconstructed input x' is close to the original input x

The latent representation z can be seen as an *encoding* of the input that captures enough information to be able to reconstruct the input

Widely studied with several algorithms, models, etc

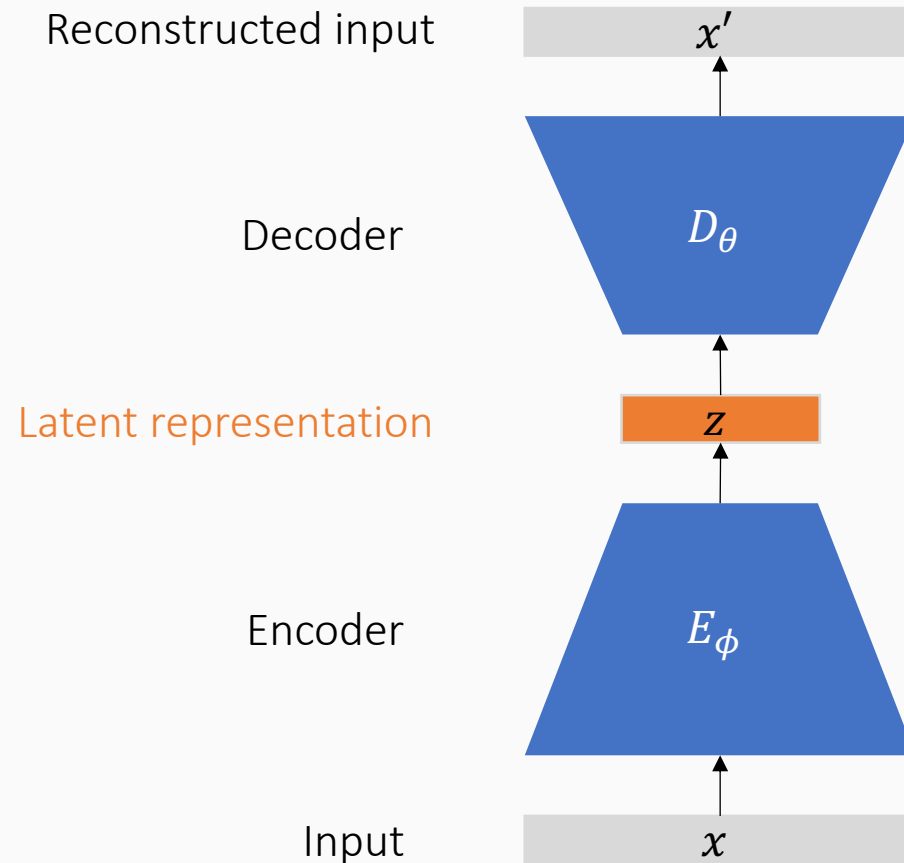
Example 1: An autoencoder



Suppose we want the latent representation to be **discrete**? E.g. a binary vector that is sampled from the sigmoids produced by the encoder

What might some advantages of such a discrete representation be?

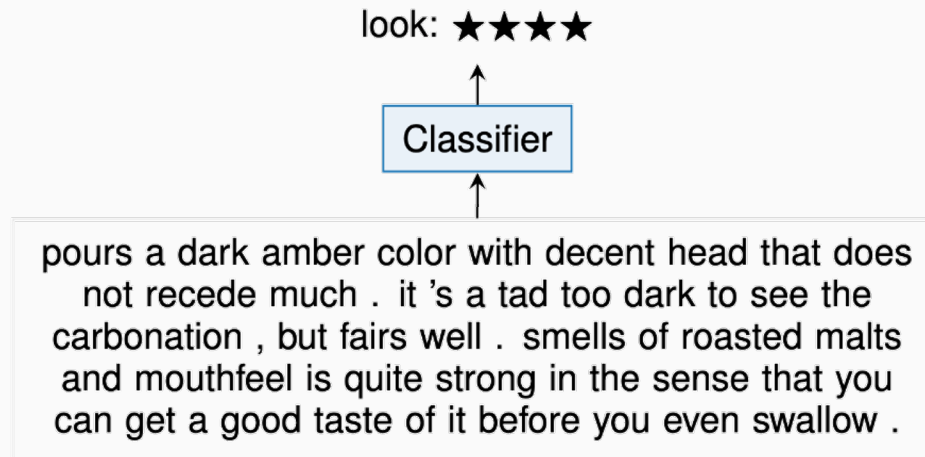
Example 1: An autoencoder



Suppose we want the latent representation to be **discrete**? E.g. a binary vector that is sampled from the sigmoids produced by the encoder

How does this affect learning?

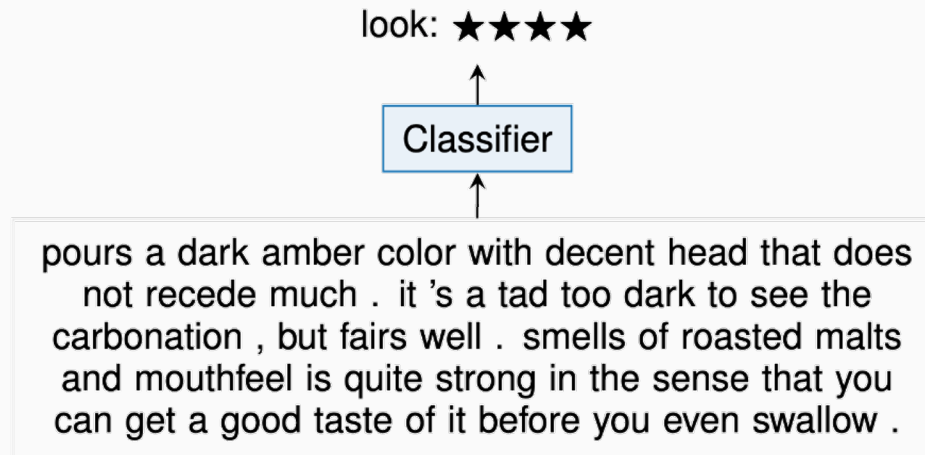
Example 2: Text classification and rationales



Suppose we have a text classification task

Given some text of a review, we need to assign it a rating for some aspect (here “look”)

Example 2: Text classification and rationales

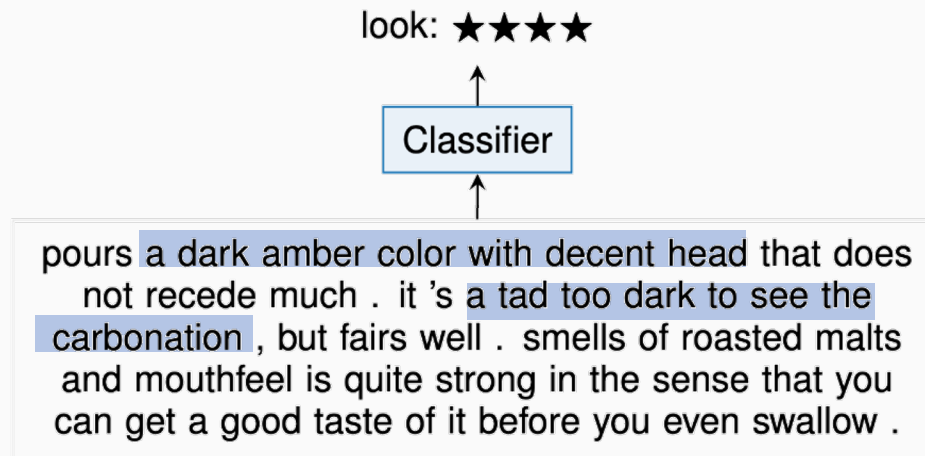


Suppose we have a text classification task

Given some text of a review, we need to assign it a rating for some aspect (here “look”)

We can ask: Why did the classifier predict a four star rating for this example?

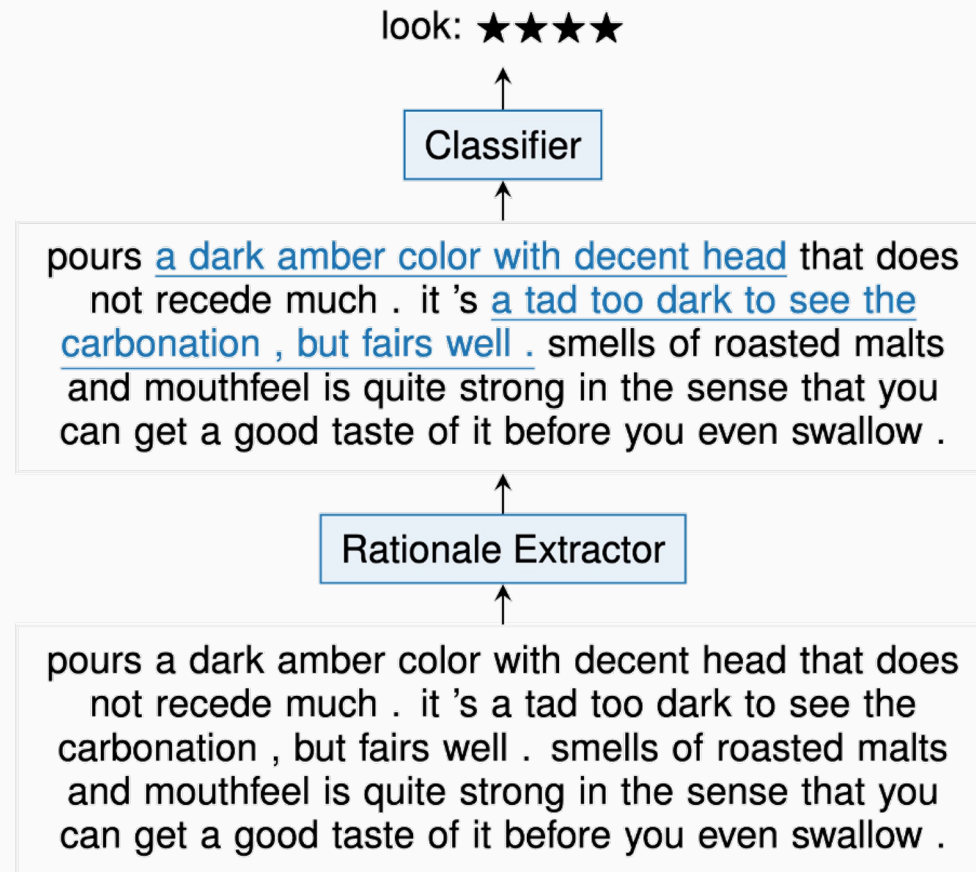
Example 2: Text classification and rationales



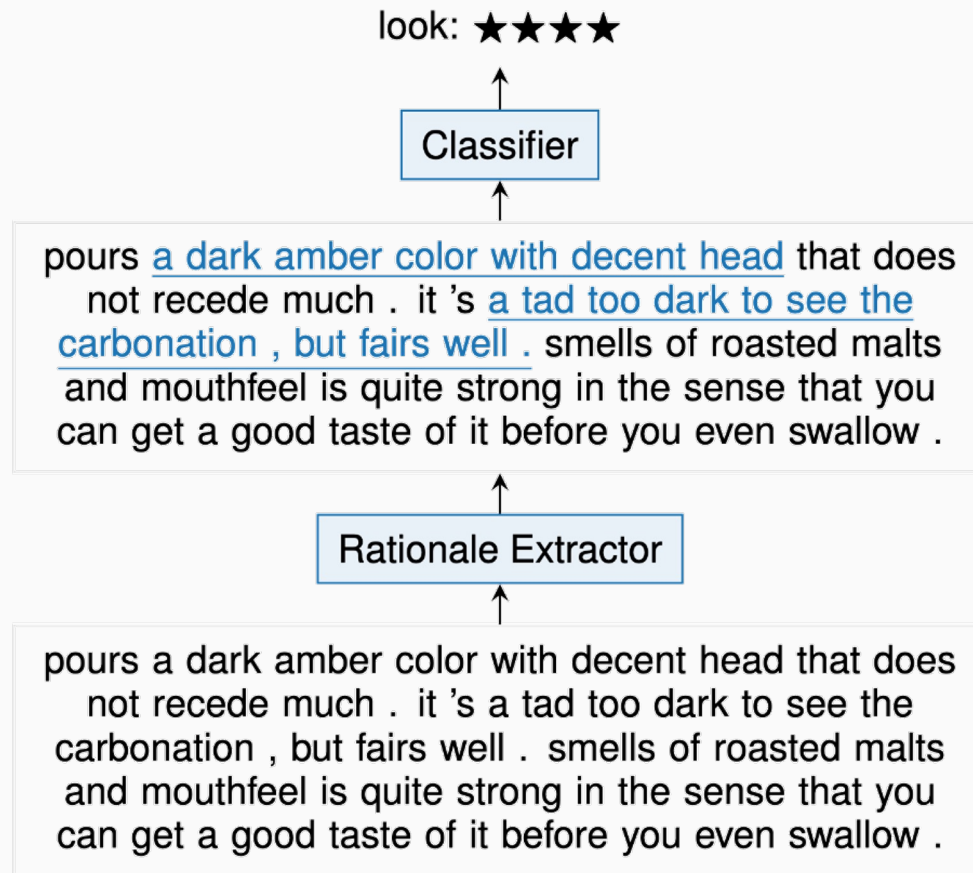
Why did the classifier predict a four star rating for this example?

One answer: Because of the highlighted words

Example 2: Text classification and rationales



Example 2: Text classification and rationales



An alternative approach:

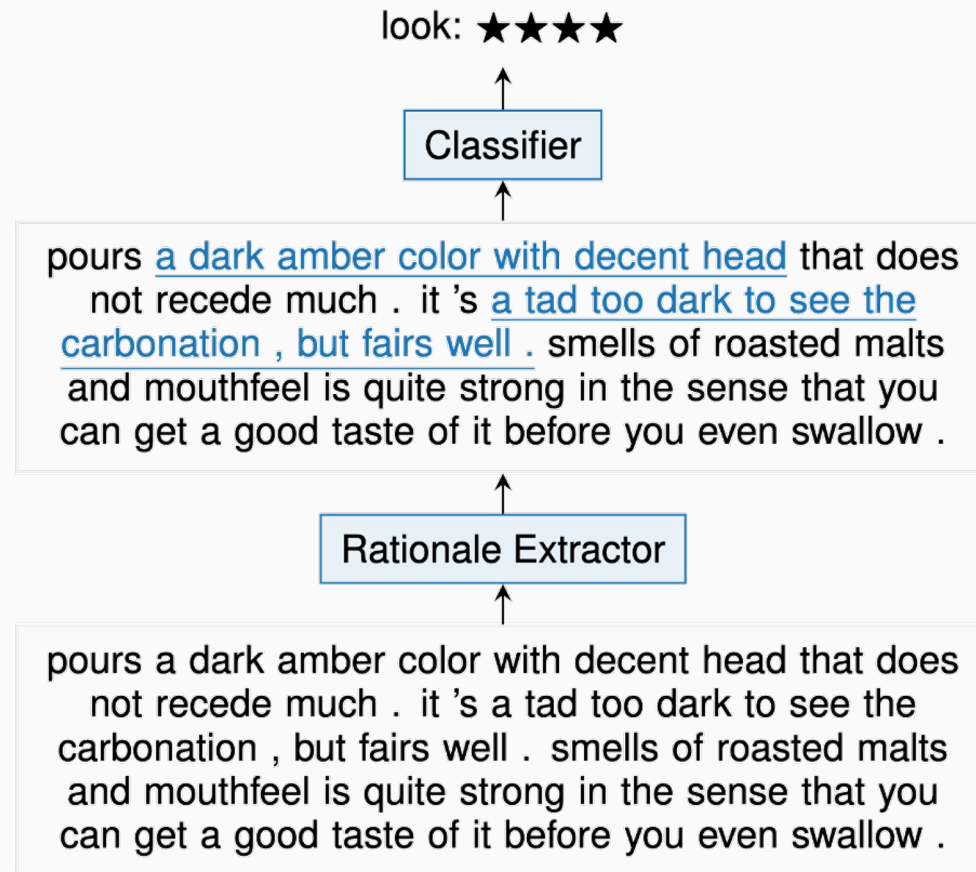
1. Identify the rationale (the highlighted words)

$$Z_i | x \sim \text{Bernoulli}(g_i(x, \phi))$$

For each token, this represents whether the token is relevant or not

Each Z_i can be seen as a Boolean proposition

Example 2: Text classification and rationales



An alternative approach:

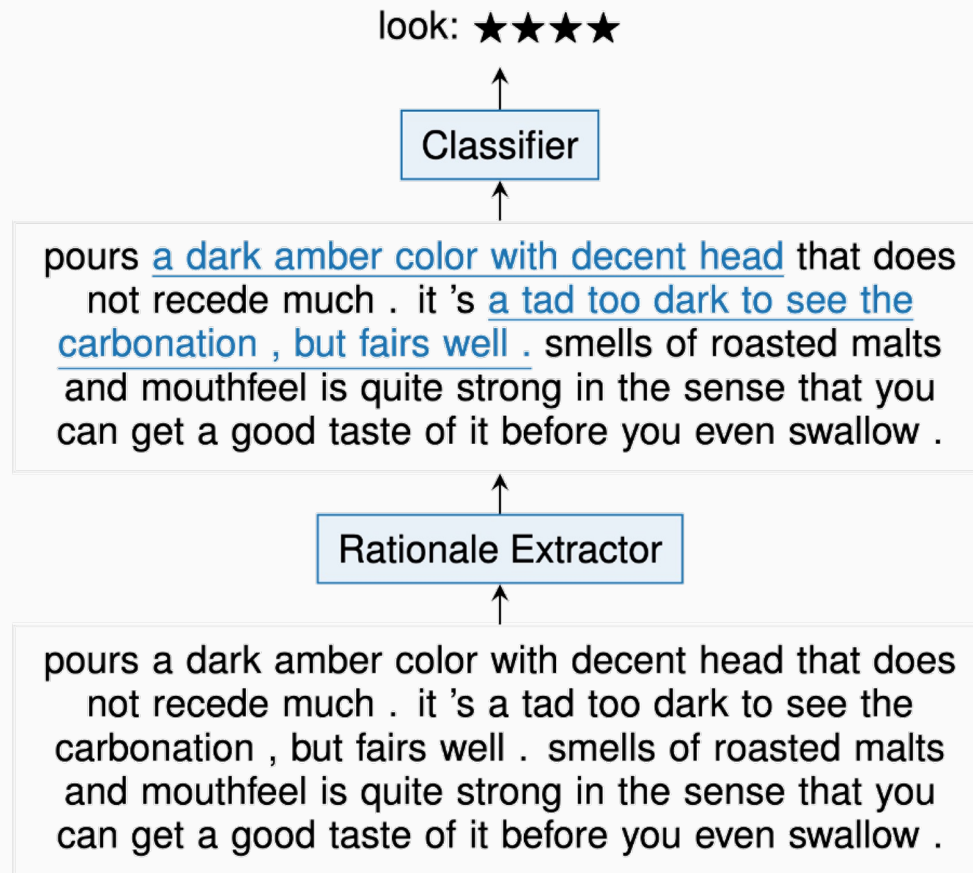
1. Identify the rationale (the highlighted words)
2. Use only the highlighted words as input to the classifier

$$Z_i | x \sim \text{Bernoulli}(g_i(x, \phi))$$

$$Y | x \sim \text{Categorical}(f(x \odot Z, \theta))$$

Mask out irrelevant words

Example 2: Text classification and rationales



An alternative approach:

1. Identify the rationale (the highlighted words)
2. Use only the highlighted words as input to the classifier

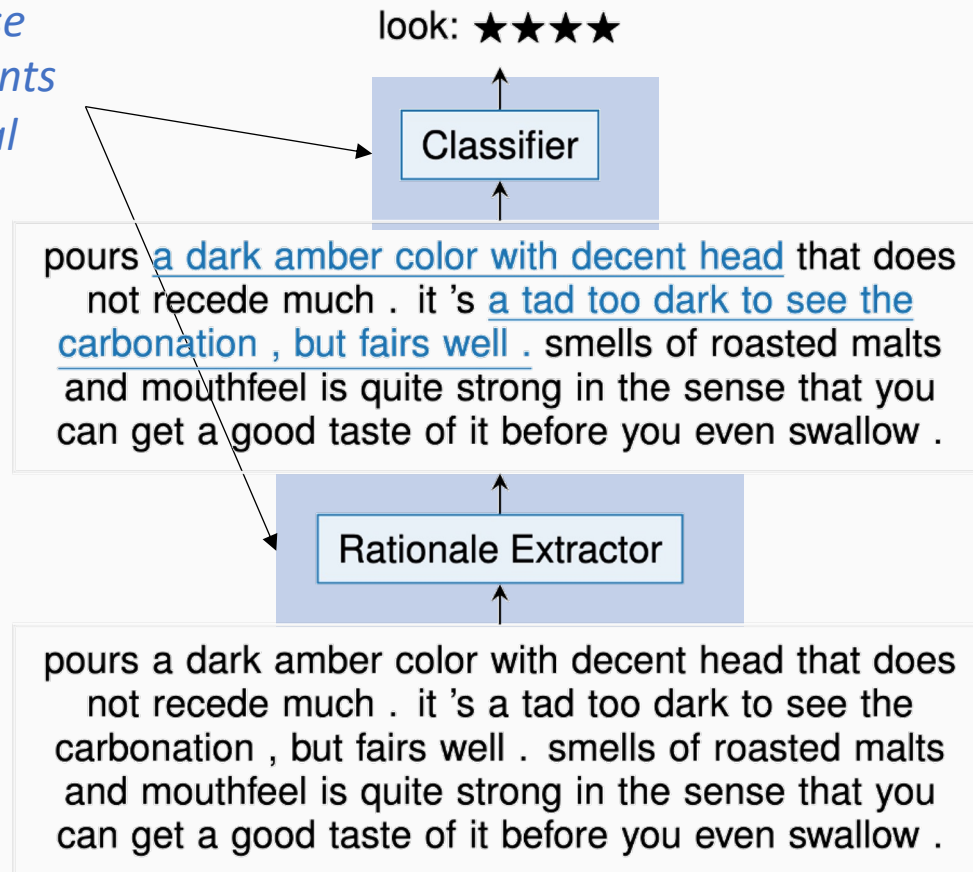
$$Z_i | x \sim \text{Bernoulli}(g_i(x, \phi))$$

$$Y | x \sim \text{Categorical}(f(x \odot Z, \theta))$$

Construct a categorical distribution over the remaining tokens

Example 2: Text classification and rationales

Both these components are neural networks



An alternative approach:

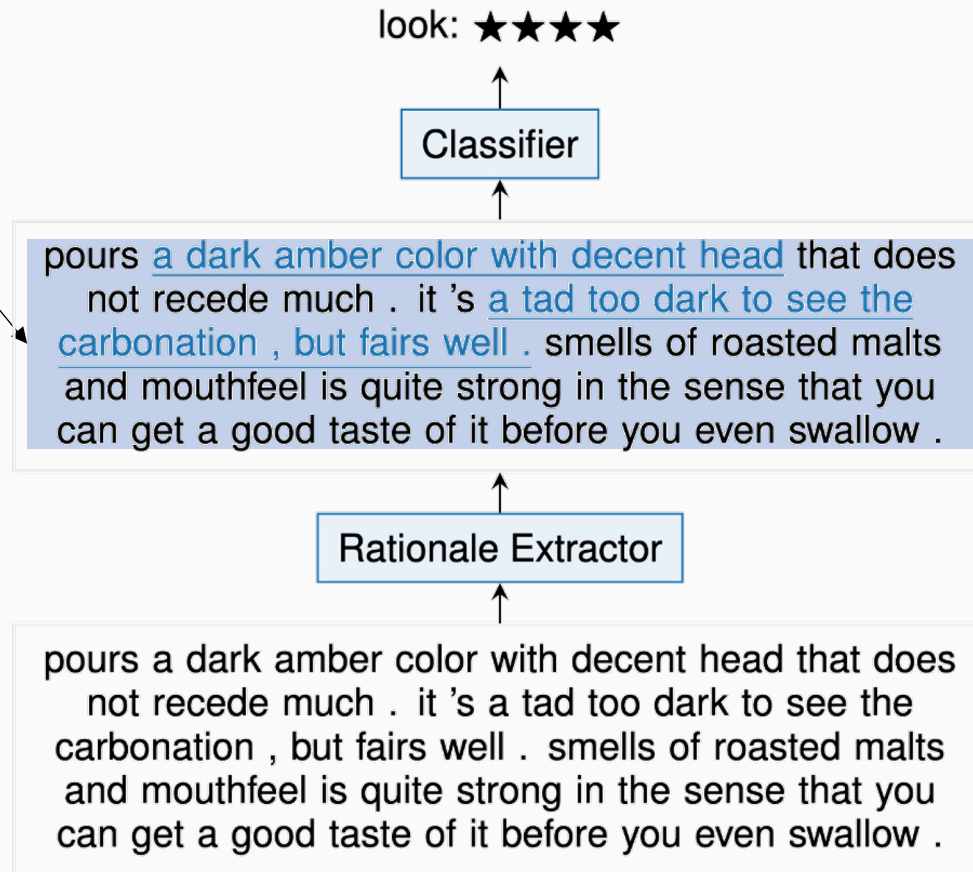
1. Identify the rationale (the highlighted words)
2. Use only the highlighted words as input to the classifier

$$Z_i | x \sim \text{Bernoulli}(g_i(x, \phi))$$

$$Y | x \sim \text{Categorical}(f(x \odot Z, \theta))$$

Example 2: Text classification and rationales

But this step is discrete.



An alternative approach:

1. Identify the rationale (the highlighted words)
2. Use only the highlighted words as input to the classifier

$$Z_i | x \sim \text{Bernoulli}(g_i(x, \phi))$$

$$Y | x \sim \text{Categorical}(f(x \odot Z, \theta))$$

Other examples: Suppose we have ...

...one neural component to permute a collection of inputs (e.g. sort it), and another network operating on the permutations

...one network identifying relevant parts of an image, and a different one using only the relevant components to make its prediction

...a black box program whose inputs are produced by a neural network, and whose outputs are consumed by another one

In all cases: How should we train the two networks?

Some approaches for addressing such problems

- The straight-through estimator
- The Gumbel trick
- Policy gradient and other RL approaches
- Differentiable structured layers

Not all these approaches
are always applicable