



→ Simple dist → complex dist

eg. VAE

Generative  
Process

$z \sim p(z)$   
latent  
variable

$x \sim p_\theta(x|z)$   
likelihood,  
decoder

$p(x)$  → complex  
eg. image, audio

encoder  
Posterior  
 $p_\theta(z|x) = \frac{p_\theta(x|z) p(z)}{p_\theta(x)}$   
prior  
eg.  $\mathcal{N}(0, I)$   
learn  $\theta$   
decoder given  $z$  what image?  
likelihood  
prior  
evidence

$$p_\theta(x) = \int p_\theta(x|z) p(z) dz$$

is intractable

approx posterior

$q_\phi(z|x)$   
need to learn  
encoder

Select  $q$  → normalizing constant  
sample

Goal: Minimize

$$KL(q_\phi(z|x) \parallel p_\theta(z|x)) = E_{q_\phi} \left[ \log \frac{q_\phi(z|x)}{p_\theta(z|x)} \right]$$

ELBO (Evidence lower bound)

(Variational Inference)

intractable

$$\log p_\theta(x) = \text{ELBO} + KL$$

(log likelihood)

$$KL \geq 0$$

$$\log p_\theta(x) \geq \text{ELBO}$$

ELBO:

$$\underbrace{\mathcal{L}(x; \theta, \phi)}_{\text{minimize}} = E_{q_\phi(z|x)} \left[ \log p_\theta\left(\frac{x}{2}\right) \right] - KL(q_\phi(z|x) \parallel p(z))$$

Interlude

$$D_{KL}(q_{\phi} || p_{\theta}) = E_{q_{\phi}} \left[ \log \frac{q_{\phi}(z|x)}{p_{\theta}(z|x)} \right]$$

Applies to both discrete and cont, use for distillation, generative models etc.  
 The less likely an event, the more surprising it is, so surprise is  $1/P(x)$

90% chance no surprise quiz! You can chill.  
Surprise quiz!

No quiz? Expected!

less likely  $\rightarrow$  more surprising

$$I(x) = \log \frac{1}{P(x)}$$

mon quiz	tues quiz	wed quiz
-------------	--------------	-------------

3x surprise

Ind Probabilities multiply  $\rightarrow$  Surprise adds

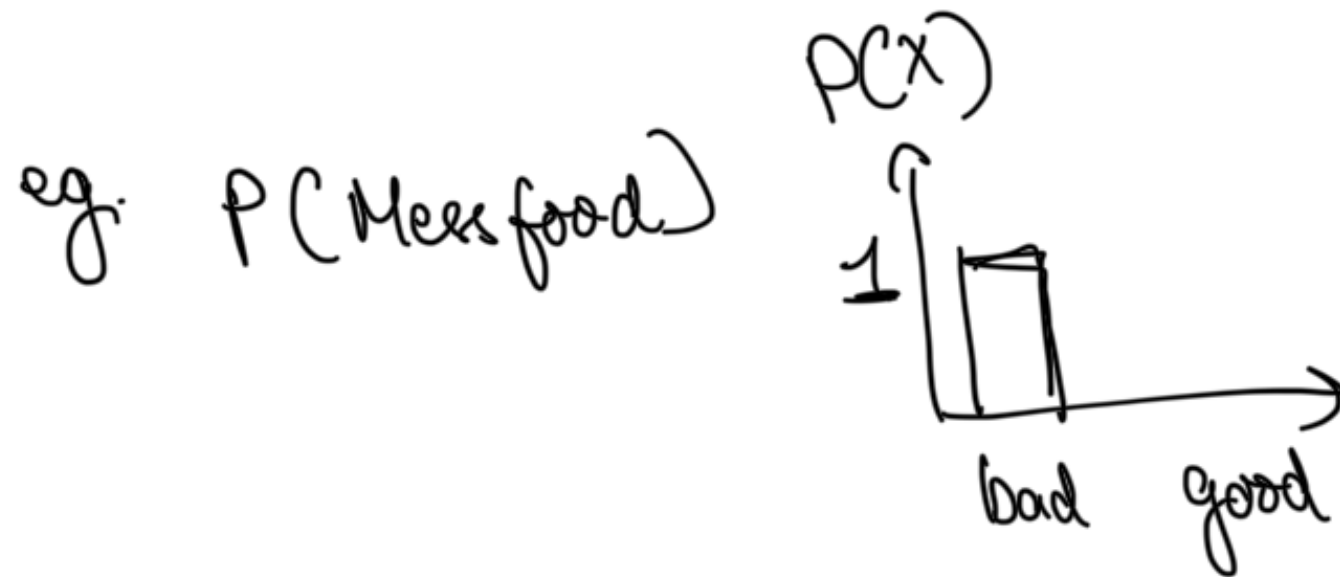
$$\text{Entropy} = E[I] = H(P) = 90\% \log \frac{1}{90\%} + 10\% \log \frac{1}{10\%}$$

Avg Surprise

= bits

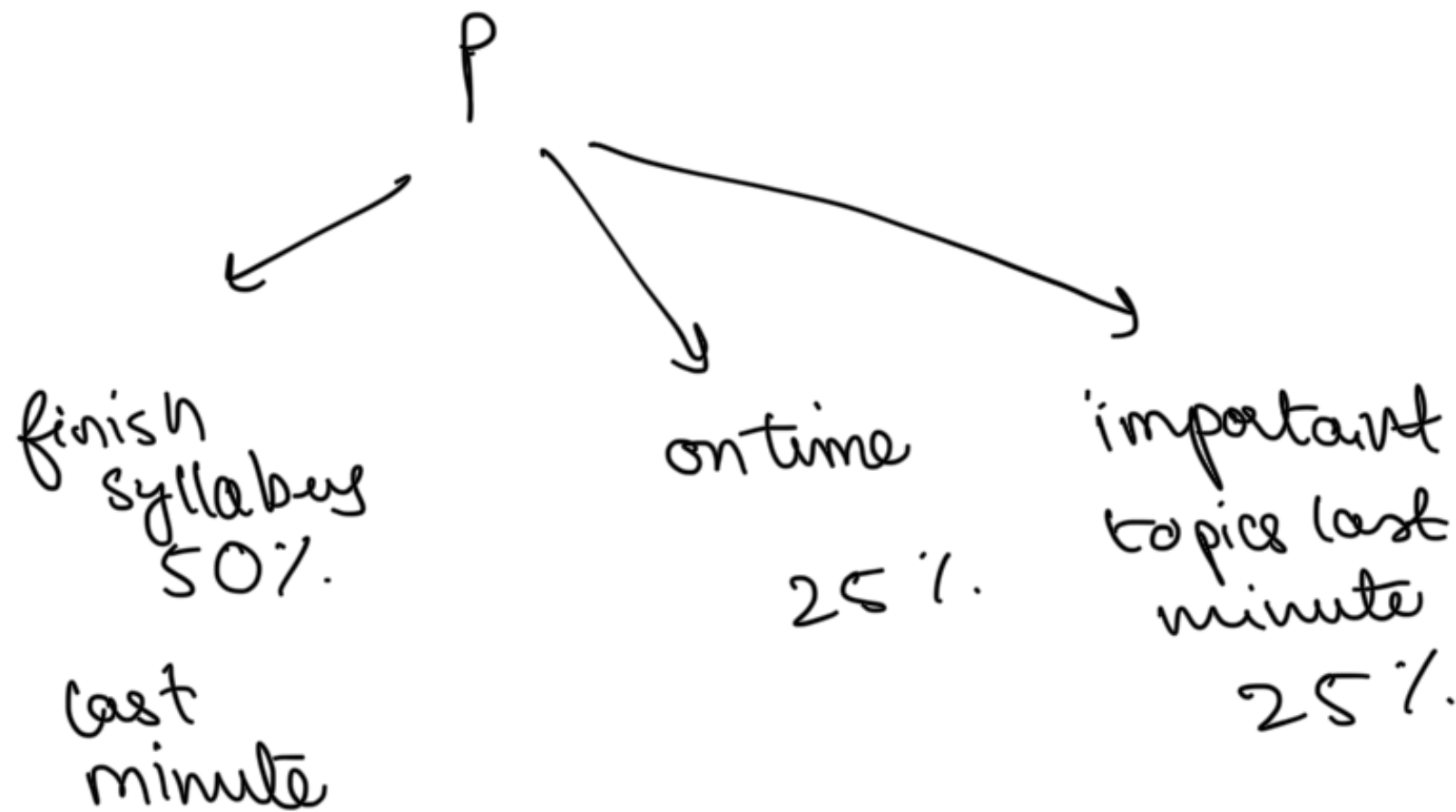
encode  
outcome

$$= \sum_i p(x_i) \log_2 \frac{1}{p(x_i)}$$



$$H(P) = 1 \log \frac{1}{1} = 0$$

= 0 bits to encode this  
information  
= no surprise



$$H(P) = \frac{1}{2} \log \frac{1}{(1/2)} + \frac{1}{4} \log \frac{1}{(1/4)} + \frac{1}{4} \log \frac{1}{(1/4)}$$

$$= 1.5$$

Encoding:

0

10

11

Say  
... optimal

suboptimal  
encoding

10  
assume  
25%

0  
(50%)

11  
(25%)

wrong idea  
code optimized  
for wrong model

Cross Entropy  
 $H(p, q)$   
↑     ↑  
real   assumption

$$= \sum p(x) \log_2 \frac{1}{q(x)}$$

true     approx  
↓       ↓

$$D_{KL}(P \parallel Q) = H(P, Q) - H(P)$$

$$= \left(\frac{1}{2}\right) \log_2 \left(\frac{1}{1/4}\right) + \frac{1}{4} \log_2 \left(\frac{1}{1/2}\right) + \frac{1}{4} \log_2 \left(\frac{1}{1/4}\right)$$

Self Information      $I(x) = \log \frac{1}{P(x)}$

$$= 1 + \frac{1}{4} + \frac{1}{2}$$

$$= 1.75$$

Cross Entropy : Avg surprise under wrong belief



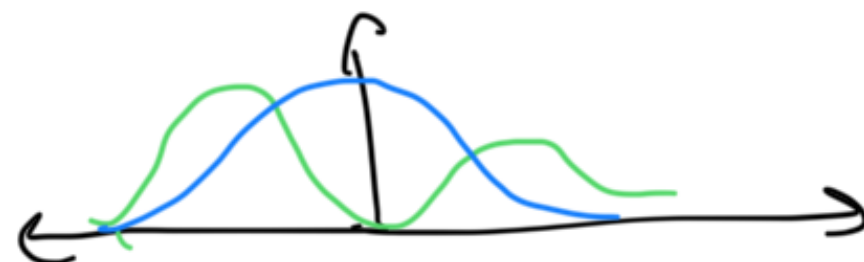
so on average 0.25 bits extra

forward KL

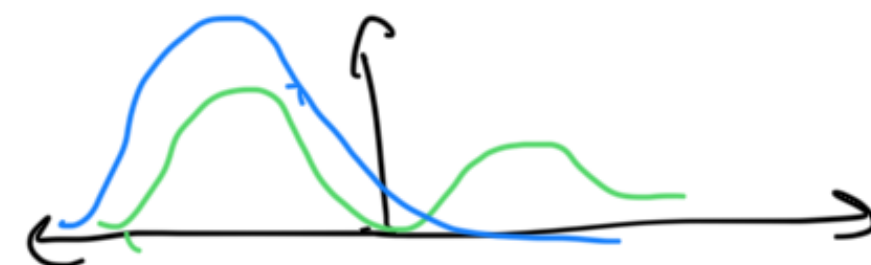
$D_{KL}(P||Q) \rightarrow$  mode covering

$D_{KL}(Q||P) \rightarrow$  mode seeking  
(reverse)

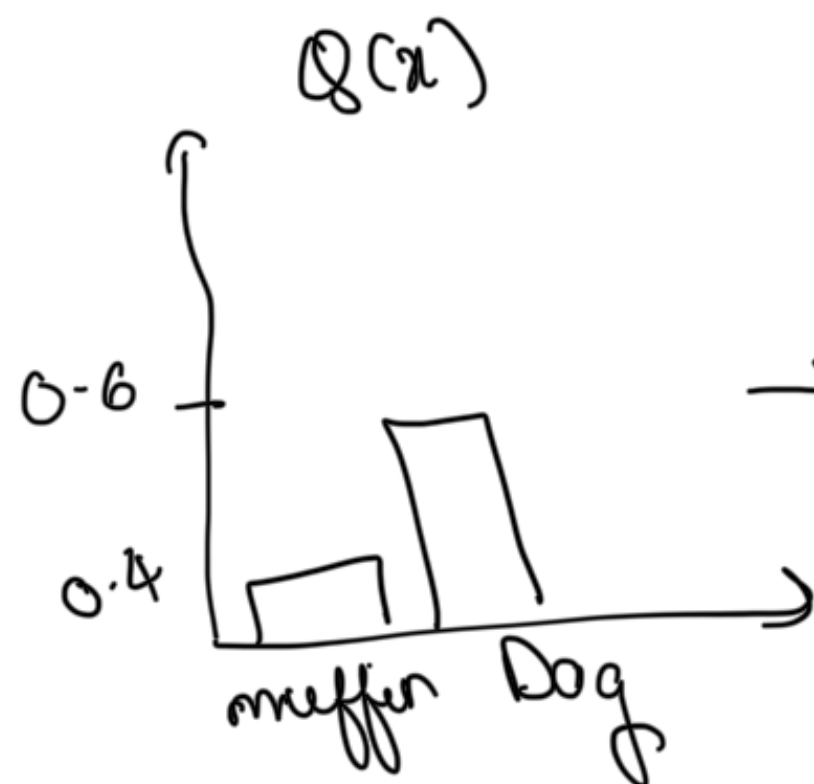
$$E \left[ Q(x) \log \frac{Q(x)}{P(x)} \right]$$



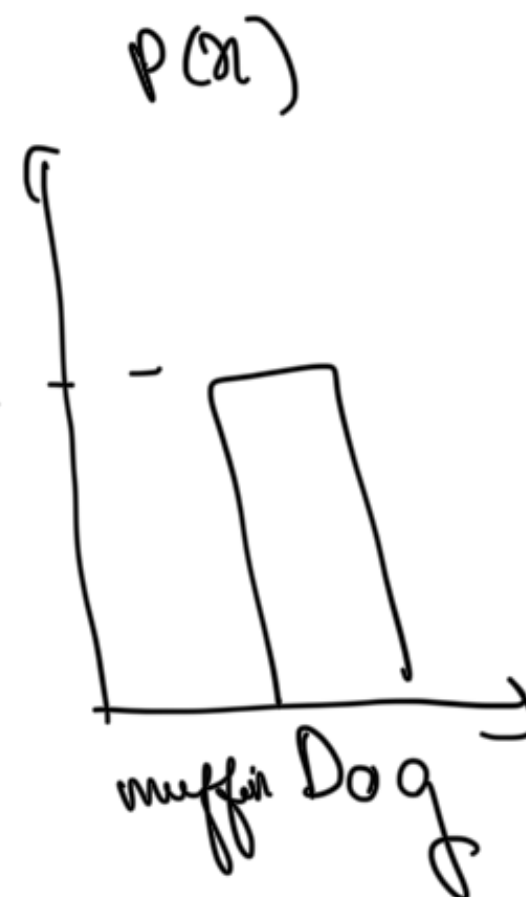
picks one peak



Penalized for assigning mass where P has none



$\rightarrow D_{KL} \leftarrow 1$



$$\min D(P||Q) = E \left[ p(x) \log \frac{p(x)}{q(x)} \right]$$



ie minimize

$$- \sum P(x_i) \log Q(x_i)$$

Problem

Vocab size LCM

$\approx 200K$  tokens

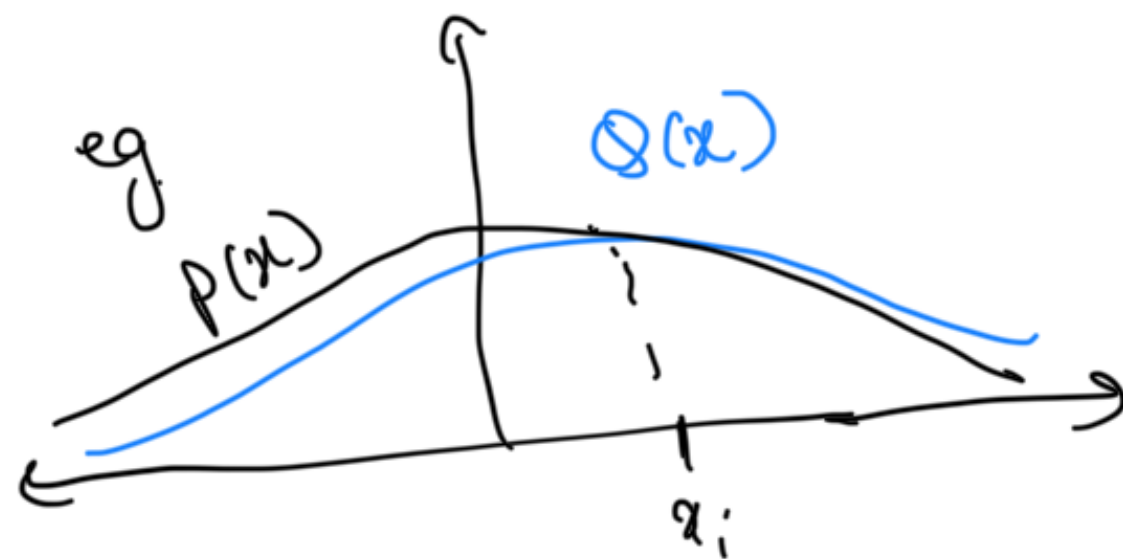
no closed form solution!

$$E_{x \sim P} \left[ \log \frac{P(x)}{Q(x)} \right]$$

Monte Carlo

Sample  $x_1, \dots, x_N \sim P(x)$

$$\hat{J}_{MC} = \frac{1}{N} \sum_{k=1}^N \log \frac{P(x_k)}{Q(x_k)} \quad \left. \vphantom{\sum_{k=1}^N} \right\} \text{Unbiased}$$



lot of variance

negative if  $Q(x_k) > p(x_k)$

$$= \frac{1}{N} \sum \frac{1}{2} ( )^2$$

lower variance  
if  $p, Q$  are close enough

Bias deviation  
from true value

Control Variates

simple est  $\rightarrow$  add term  $ST$   $E[\text{term}] = 0$

$f(x_1)$   $f(x_2)$   $f(x_3)$   $f(x_4)$

$$= \frac{1}{N} \sum \left[ \frac{p(x_k)}{Q(x_k)} + \lambda \left( r(x_k) - \mathbb{E}_{a \sim p} [r(a)] \right) \right]$$

$$r(x) = \frac{Q(x)}{P(x)}$$

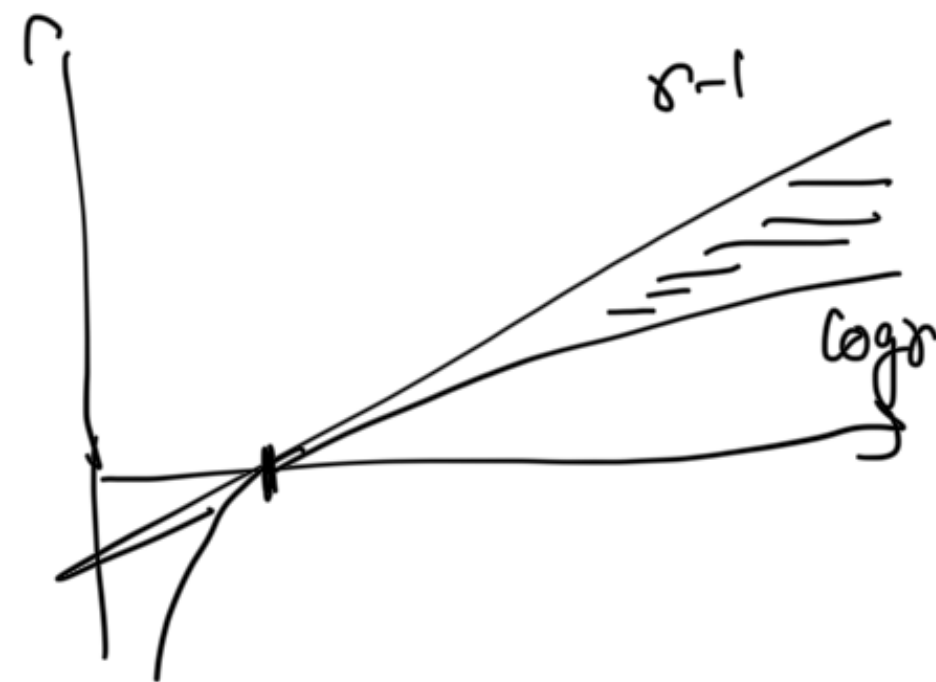
$$\begin{aligned} \mathbb{E} r(x) &= \int p(x) r(x) dx \\ &= 1 \end{aligned}$$

now  $\lambda$ ?

$$= \frac{1}{N} \sum \left[ -\log r(x_k) + \lambda (r - 1) \right]$$

$\lambda = 1$ , estimator remains non negative

$$= \frac{1}{N} \sum \left[ r(x_k) - 1 - \log r(x_k) \right]$$



ELBO

$$D_{KL}(q_\phi || p_\theta) = E_{q_\phi} \left[ \log \frac{q_\phi(z|x)}{p_\theta(z|x)} \right]$$

$$= E_{q_\phi} [\log q_\phi(z|x)] - E_{q_\phi} [\log p_\theta(z|x)]$$

$$= E_{q_\phi} [\log q_\phi(z|x)] - E_{q_\phi} [\log p(z, x)] + \log p_\theta(x)$$

$$\int q_\phi(z|x) dz$$

$$= E_{q_\phi} [\log q_\phi(z|x)] - E [\log p_\theta(z, x)] + \log p_\theta(x)$$

marginal log  
likelihood

component

$$\log p_{\theta}(x) = -E_{q_{\phi}}[\log q_{\phi}(z|x)] + E_{q_{\phi}}[\log p_{\theta}(z,x)]$$

can't compute

log evident

$$+ D_{KL}(q_{\phi} || p_{\theta})$$

component 2  
intractable

always positive

$\geq 0$

$$\log p_{\theta}(x)$$

$\geq$



Evident lower bound



maximize this  
to minimize

$$D_{KL}(q_{\phi} || p_{\theta})$$

$$ELBO = -E_{q_{\phi}}[\log q_{\phi}(z|x)] + E_{q_{\phi}}[\log p_{\theta}(z,x)]$$

$$= E_{q_{\phi}} \left[ \log p_{\theta} \left( \overset{\text{orig}}{\frac{x}{z}} \right) \right] - E_{q_{\phi}} \left[ \log \frac{q_{\phi}(z|x)}{p_{\theta}(z)} \right]$$

latent space

Expected  
Reconstruction  
error

learnable

approx posterior  $q$

prior

we know

KL divergence b/w

We could also pick  $q_{\phi}(z) \rightarrow$  not conditioned on  $x$