

# Machine Learning

## Deterministic Approximative Inference

---

Carl Henrik Ek - [carlhenrik.ek@bristol.ac.uk](mailto:carlhenrik.ek@bristol.ac.uk)

December 2, 2019

<http://www.carlhenrik.com>

$$p(y) = \sum_{n=1}^N p(y, x_i)$$

- Number of atoms in the universe

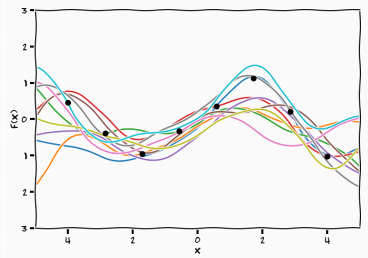
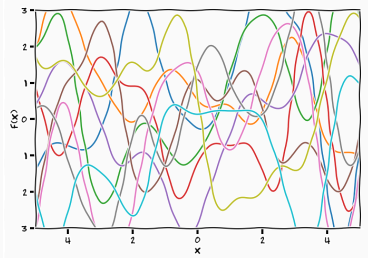
$$10^{80} \approx (2^{\frac{10}{3}})^{80} \approx 2^{267}$$

- Age of the universe in seconds

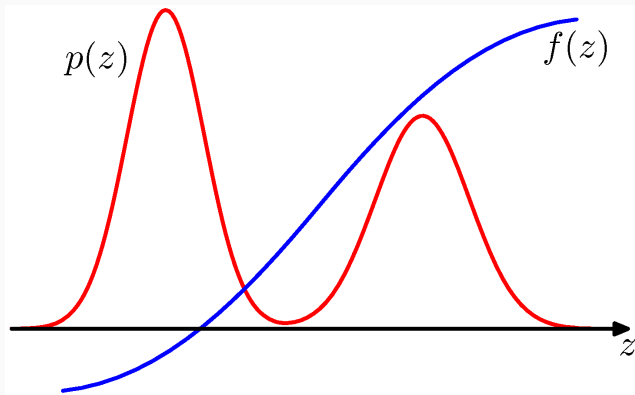
$$4.35 \cdot 10^{17} \approx 2^{59}$$



# Big Number



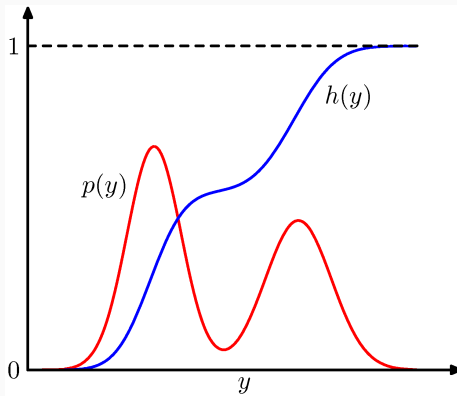
$$p(y|x) = \int p(y|f)p(f|x)df$$



$$\mathbb{E}_{p(z)}[f] = \int f(z)p(z)dz \approx \frac{1}{L} \sum_{l=1}^L f(z^{(l)})$$

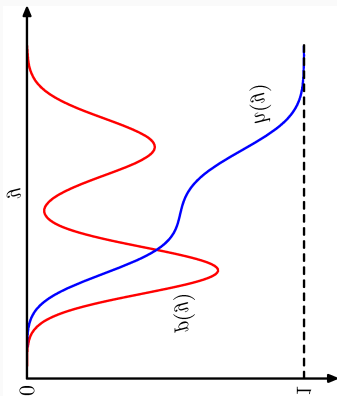
$$z^{(l)} \sim p(z)$$

# Basic Probabilities



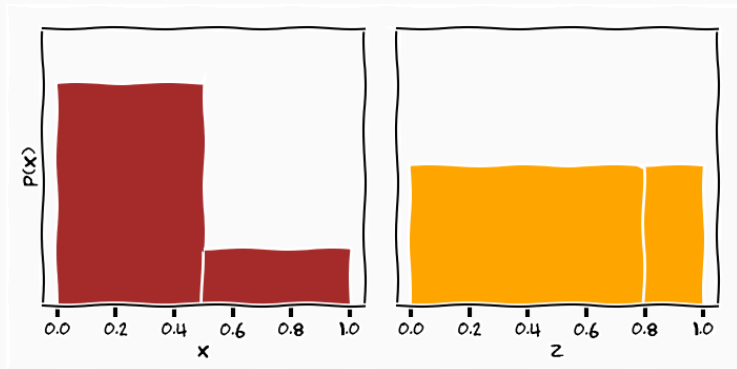
$$z = h(y) = \int_{-\infty}^y p(y) dy$$

# Basic Probabilities



$$y = h^{-1}(z)$$

# Change of Variables





$$\int_{y \in \mathcal{Y}} p(y) dy = \int_{x \in \mathcal{X}} p(x) dx$$

$$\int_{y \in \mathcal{Y}} p(y) dy = \int_{x \in \mathcal{X}} p(x) dx$$
$$x = x(y)$$

# Change of Variables

$$\int_{y \in \mathcal{Y}} p(y) dy = \int_{x \in \mathcal{X}} p(x) dx$$

$$x = x(y)$$

$$\int_{y \in \mathcal{Y}} p(y) dy = \int_{y \in \mathcal{Y}} p(x(y)) \frac{dx}{dy} dy$$

# Change of Variables

$$\int_{y \in \mathcal{Y}} p(y) dy = \int_{x \in \mathcal{X}} p(x) dx$$

$$x = x(y)$$

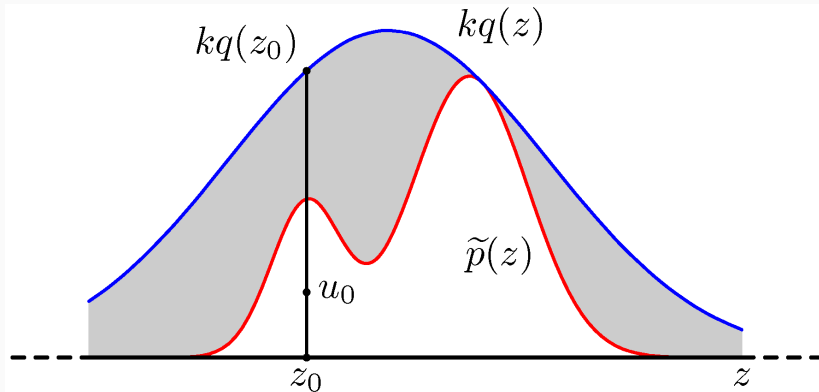
$$\int_{y \in \mathcal{Y}} p(y) dy = \int_{y \in \mathcal{Y}} p(x(y)) \frac{dx}{dy} dy$$

$$p(y) = p(x(y)) \frac{dx}{dy}$$

# Starting point

- We can sample random numbers from the uniform distribution
- We can use the indefinite integral to transform the uniform to any distribution
- Want to use these distributions as proxies

# Rejection Sampling

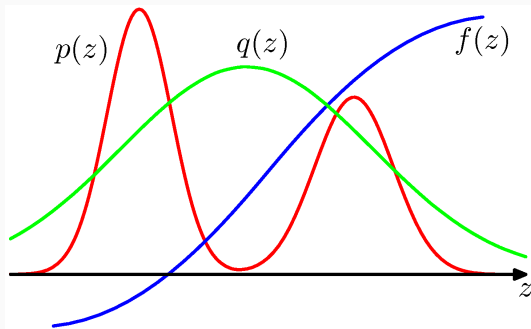


$$\mathbb{E}_{p(\mathbf{z})}[f] = \int f(\mathbf{z})p(\mathbf{z})d\mathbf{z} = \int f(\mathbf{z})\frac{p(\mathbf{z})}{q(\mathbf{z})}q(\mathbf{z})d\mathbf{z} = \mathbb{E}_{q(\mathbf{z})}\left[f(\mathbf{z})\frac{p(\mathbf{z})}{q(\mathbf{z})}\right]$$

$$\approx \frac{1}{L} \sum_{l=1}^L \frac{p(\mathbf{z}^{(l)})}{q(\mathbf{z}^{(l)})} \cdot f(\mathbf{z}^{(l)})$$

- Sample from proposal distribution and re-weight samples
- Accepts all samples

# Importance Sampling



$$r_l = \frac{p(z^{(l)})}{q(z^{(l)})}$$



## Metropolis Sampling

1. start with state  $z^{(0)}$

## Metropolis Sampling

1. start with state  $\mathbf{z}^{(0)}$
2. sample from conditional proposal distribution  $q(\mathbf{z}^*|\mathbf{z}^{(0)})$

## Metropolis Sampling

1. start with state  $\mathbf{z}^{(0)}$
2. sample from conditional proposal distribution  $q(\mathbf{z}^*|\mathbf{z}^{(0)})$
3. compute acceptance probability

$$A(\mathbf{z}^*, \mathbf{z}^{(0)}) = \min \left( 1, \frac{\tilde{p}(\mathbf{z}^*)}{\tilde{p}(\mathbf{z}^{(0)})} \right)$$

## Metropolis Sampling

1. start with state  $\mathbf{z}^{(0)}$
2. sample from conditional proposal distribution  $q(\mathbf{z}^*|\mathbf{z}^{(0)})$
3. compute acceptance probability

$$A(\mathbf{z}^*, \mathbf{z}^{(0)}) = \min \left( 1, \frac{\tilde{p}(\mathbf{z}^*)}{\tilde{p}(\mathbf{z}^{(0)})} \right)$$

4. Draw uniform random number  $u \sim \text{Uniform}(0, 1)$

## Metropolis Sampling

1. start with state  $\mathbf{z}^{(0)}$
2. sample from conditional proposal distribution  $q(\mathbf{z}^*|\mathbf{z}^{(0)})$
3. compute acceptance probability

$$A(\mathbf{z}^*, \mathbf{z}^{(0)}) = \min \left( 1, \frac{\tilde{p}(\mathbf{z}^*)}{\tilde{p}(\mathbf{z}^{(0)})} \right)$$

4. Draw uniform random number  $u \sim \text{Uniform}(0, 1)$ 
  - if  $A(\mathbf{z}^*, \mathbf{z}^{(0)}) > u \rightarrow \mathbf{z}^{(1)} = \mathbf{z}^*$

## Metropolis Sampling

1. start with state  $\mathbf{z}^{(0)}$
2. sample from conditional proposal distribution  $q(\mathbf{z}^*|\mathbf{z}^{(0)})$
3. compute acceptance probability

$$A(\mathbf{z}^*, \mathbf{z}^{(0)}) = \min \left( 1, \frac{\tilde{p}(\mathbf{z}^*)}{\tilde{p}(\mathbf{z}^{(0)})} \right)$$

4. Draw uniform random number  $u \sim \text{Uniform}(0, 1)$ 
  - if  $A(\mathbf{z}^*, \mathbf{z}^{(0)}) > u \rightarrow \mathbf{z}^{(1)} = \mathbf{z}^*$
  - otherwise reject  $\mathbf{z}^*$  and start over

1. Initialise  $\mathbf{z}^{(0)}$

1. Initialise  $\mathbf{z}^{(0)}$
2. Pick single variable  $z_i \in \mathbf{z}$



1. Initialise  $\mathbf{z}^{(0)}$
2. Pick single variable  $z_i \in \mathbf{z}$
3. Formulate posterior  $p(z_i | \mathbf{z}_{\neg i})$

1. Initialise  $\mathbf{z}^{(0)}$
2. Pick single variable  $z_i \in \mathbf{z}$
3. Formulate posterior  $p(z_i | \mathbf{z}_{\neg i})$
4. Sample from posterior

$$z_i^{(1)} \sim p(z_i | \mathbf{z}_{\neg i})$$

1. Initialise  $\mathbf{z}^{(0)}$
2. Pick single variable  $z_i \in \mathbf{z}$
3. Formulate posterior  $p(z_i | \mathbf{z}_{\neg i})$
4. Sample from posterior

$$z_i^{(1)} \sim p(z_i | \mathbf{z}_{\neg i})$$

5. cycle through variables

- We can sample from distributions whose explicit form we know

- We can sample from distributions whose explicit form we know
- Formulate different proxies using these distributions to the general case

- We can sample from distributions whose explicit form we know
- Formulate different proxies using these distributions to the general case
  - Rejection sampling

- We can sample from distributions whose explicit form we know
- Formulate different proxies using these distributions to the general case
  - Rejection sampling
  - Importance sampling

- We can sample from distributions whose explicit form we know
- Formulate different proxies using these distributions to the general case
  - Rejection sampling
  - Importance sampling
- Create a random walk to encourage exploration



- We can sample from distributions whose explicit form we know
- Formulate different proxies using these distributions to the general case
  - Rejection sampling
  - Importance sampling
- Create a random walk to encourage exploration
  - Gibbs sampling

# Deterministic Approximations

---

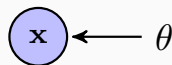
- Stochastic inference
  - approximate expectation with sum
  - works in the limit
  - hard to know how well we are doing
  - usually slow

- Stochastic inference
  - approximate expectation with sum
  - works in the limit
  - hard to know how well we are doing
  - usually slow
- Idea
  - *can we reformulate inference as optimisation?*

$$p(\mathbf{Y})$$

- Given some observed data  $\mathbf{Y}$
- Find a probabilistic model such that the probability of the data is maximised
- Idea: find an approximate model  $q$  that we can integrate

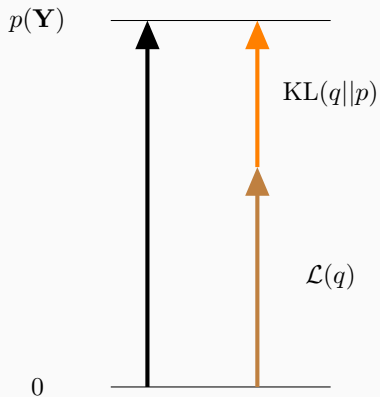
## Lower Bound



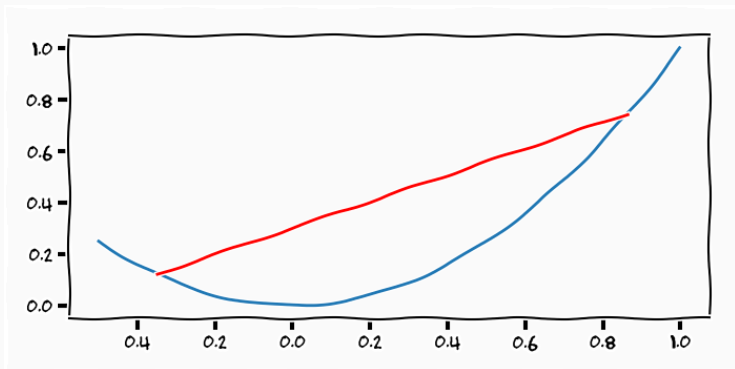
$$p(y) = \int_x p(y|x)p(x) = \frac{p(y|x)p(x)}{p(x|y)}$$

$$q_\theta(x) \approx p(x|y)$$

# Deterministic Approximation



# Jensen Inequality



## Convex Function

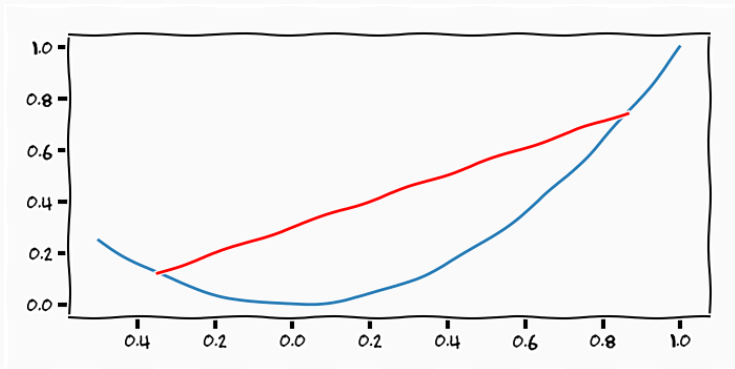
$$\lambda f(x_0) + (1 - \lambda)f(x_1) \geq f(\lambda x_0 + (1 - \lambda)x_1)$$

$$x \in [x_{min}, x_{max}]$$

$$\lambda \in [0, 1]$$



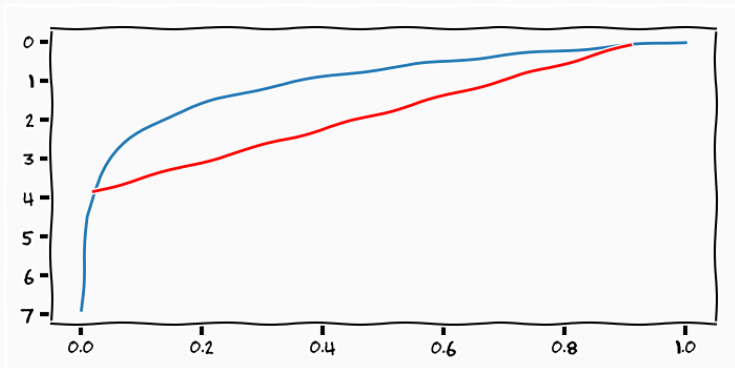
# Jensen Inequality



$$\mathbb{E}[f(x)] \geq f(\mathbb{E}[x])$$

$$\int f(x)p(x)dx \geq f\left(\int xp(x)dx\right)$$

# Jensen Inequality in Variational Bayes



$$\int \log(x)p(x)dx \leq \log\left(\int xp(x)dx\right)$$

*moving the log inside the the integral is a lower-bound on the integral*

$$p(y)$$

$$\log p(y)$$

$$\log p(y) = \log p(y) + \int \log \frac{p(x|y)}{p(x|y)} dx$$

$$\begin{aligned}\log p(y) &= \log p(y) + \int \log \frac{p(x|y)}{p(x|y)} dx \\ &= \int q(x) \log p(y) dx + \int q(x) \log \frac{p(x|y)}{p(x|y)} dx\end{aligned}$$

$$\begin{aligned}\log p(y) &= \log p(y) + \int \log \frac{p(x|y)}{p(x|y)} dx \\&= \int q(x) \log p(y) dx + \int q(x) \log \frac{p(x|y)}{p(x|y)} dx \\&= \int q(x) \log \frac{p(x|y)p(y)}{p(x|y)} dx\end{aligned}$$

$$\begin{aligned}\log p(y) &= \log p(y) + \int \log \frac{p(x|y)}{p(x|y)} dx \\&= \int q(x) \log p(y) dx + \int q(x) \log \frac{p(x|y)}{p(x|y)} dx \\&= \int q(x) \log \frac{p(x|y)p(y)}{p(x|y)} dx \\&= \int q(x) \log \frac{q(x)}{q(x)} dx + \int q(x) \log p(x, y) dx + \int q(x) \log \frac{1}{p(x|y)} dx\end{aligned}$$



$$\begin{aligned}\log p(y) &= \log p(y) + \int \log \frac{p(x|y)}{p(x|y)} dx \\&= \int q(x) \log p(y) dx + \int q(x) \log \frac{p(x|y)}{p(x|y)} dx \\&= \int q(x) \log \frac{p(x|y)p(y)}{p(x|y)} dx \\&= \int q(x) \log \frac{q(x)}{q(x)} dx + \int q(x) \log p(x, y) dx + \int q(x) \log \frac{1}{p(x|y)} dx \\&= \int q(x) \log \frac{1}{q(x)} dx + \int q(x) \log p(x, y) dx + \int q(x) \log \frac{q(x)}{p(x|y)} dx\end{aligned}$$

## The "posterior" term

$$KL(q(x)||q(x|y)) = \int q(x) \log \frac{q(x)}{p(x|y)} dx$$

## The "posterior" term

$$\begin{aligned} KL(q(x)||q(x|y)) &= \int q(x) \log \frac{q(x)}{p(x|y)} dx \\ &= - \int q(x) \log \frac{p(x|y)}{q(x)} dx \end{aligned}$$

## The "posterior" term

$$\begin{aligned} KL(q(x)||p(x|y)) &= \int q(x) \log \frac{q(x)}{p(x|y)} dx \\ &= - \int q(x) \log \frac{p(x|y)}{q(x)} dx \\ &\geq -\log \int p(x|y) dx = -\log 1 = 0 \end{aligned}$$

## The "posterior" term

$$KL(q(x)||p(x|y)) = \int q(x) \log \frac{q(x)}{p(x|y)} dx$$

## The "posterior" term

$$KL(q(x)||p(x|y)) = \int q(x) \log \frac{q(x)}{p(x|y)} dx$$

= {Lets assume that  $q(x) = p(x|y)$ }

## The "posterior" term

$$\begin{aligned} KL(q(x)||p(x|y)) &= \int q(x) \log \frac{q(x)}{p(x|y)} dx \\ &= \{ \text{Lets assume that } q(x) = p(x|y) \} \\ &= \int p(x|y) \log \underbrace{\frac{p(x|y)}{p(x|y)}}_{=1} dx \end{aligned}$$

## The "posterior" term

$$\begin{aligned} KL(q(x)||p(x|y)) &= \int q(x) \log \frac{q(x)}{p(x|y)} dx \\ &= \{\text{Lets assume that } q(x) = p(x|y)\} \\ &= \int p(x|y) \log \underbrace{\frac{p(x|y)}{p(x|y)}}_{=1} dx \\ &= 0 \end{aligned}$$



$$KL(q(x)||p(x|y)) = \int q(x) \log \frac{q(x)}{p(x|y)} dx$$

- Measure of divergence between distributions
- Not a metric (not symmetric)
- $KL(q(x)||p(x|y)) = 0 \Leftrightarrow q(x) = p(x|y)$
- $KL(q(x)||p(x|y)) \geq 0$

## The "other terms"

$$\int q(x) \log \frac{1}{q(x)} dx + \int q(x) \log p(x, y) dx =$$

## The "other terms"

$$\begin{aligned} \int q(x) \log \frac{1}{q(x)} dx + \int q(x) \log p(x, y) dx &= \\ = \int q(x) \log \frac{p(x, y)}{q(x)} dx \end{aligned}$$

## The "other terms"

$$\begin{aligned} & \int q(x) \log \frac{1}{q(x)} dx + \int q(x) \log p(x, y) dx = \\ &= \int q(x) \log \frac{p(x, y)}{q(x)} dx \\ &= \{ \text{Lets assume that } q(x) = p(x|y) \} \end{aligned}$$

## The "other terms"

$$\begin{aligned} & \int q(x) \log \frac{1}{q(x)} dx + \int q(x) \log p(x, y) dx = \\ &= \int q(x) \log \frac{p(x, y)}{q(x)} dx \\ &= \{ \text{Lets assume that } q(x) = p(x|y) \} \\ &= \int p(x|y) \log \frac{p(x, y)}{p(x|y)} dx \end{aligned}$$

## The "other terms"

$$\begin{aligned} & \int q(x) \log \frac{1}{q(x)} dx + \int q(x) \log p(x, y) dx = \\ &= \int q(x) \log \frac{p(x, y)}{q(x)} dx \\ &= \{ \text{Let's assume that } q(x) = p(x|y) \} \\ &= \int p(x|y) \log \frac{p(x, y)}{p(x|y)} dx = \int p(x|y) \log \frac{p(x|y)p(y)}{p(x|y)} dx \end{aligned}$$

## The "other terms"

$$\begin{aligned} & \int q(x) \log \frac{1}{q(x)} dx + \int q(x) \log p(x, y) dx = \\ &= \int q(x) \log \frac{p(x, y)}{q(x)} dx \\ &= \{ \text{Lets assume that } q(x) = p(x|y) \} \\ &= \int p(x|y) \log \frac{p(x, y)}{p(x|y)} dx = \int p(x|y) \log \frac{p(x|y)p(y)}{p(x|y)} dx \\ &= \int p(x|y) \log \underbrace{\frac{p(x|y)}{p(x|y)}}_{=1} dx + \int p(x|y) \log p(y) dx \end{aligned}$$

## The "other terms"

$$\begin{aligned} & \int q(x) \log \frac{1}{q(x)} dx + \int q(x) \log p(x, y) dx = \\ &= \int q(x) \log \frac{p(x, y)}{q(x)} dx \\ &= \{ \text{Lets assume that } q(x) = p(x|y) \} \\ &= \int p(x|y) \log \frac{p(x, y)}{p(x|y)} dx = \int p(x|y) \log \frac{p(x|y)p(y)}{p(x|y)} dx \\ &= \int p(x|y) \log \underbrace{\frac{p(x|y)}{p(x|y)}}_{=1} dx + \int p(x|y) \log p(y) dx \\ &= \underbrace{\int p(x|y) dx}_{=1} \log p(y) \end{aligned}$$



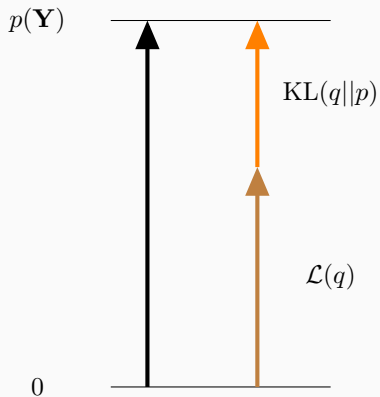
## The "other terms"

$$\begin{aligned} & \int q(x) \log \frac{1}{q(x)} dx + \int q(x) \log p(x, y) dx = \\ &= \int q(x) \log \frac{p(x, y)}{q(x)} dx \\ &= \{ \text{Let's assume that } q(x) = p(x|y) \} \\ &= \int p(x|y) \log \frac{p(x, y)}{p(x|y)} dx = \int p(x|y) \log \frac{p(x|y)p(y)}{p(x|y)} dx \\ &= \int p(x|y) \log \underbrace{\frac{p(x|y)}{p(x|y)}}_{=1} dx + \int p(x|y) \log p(y) dx \\ &= \underbrace{\int p(x|y) dx}_{=1} \log p(y) = \log p(y) \end{aligned}$$

$$\begin{aligned}\log p(y) &= \int q(x) \log \frac{1}{q(x)} dx + \int q(x) \log p(x, y) dx + \int q(x) \log \frac{q(x)}{p(x|y)} dx \\ &\geq - \int q(x) \log q(x) dx + \int q(x) \log p(x, y) dx\end{aligned}$$

- The Evidence Lower BOnd
- Tight if  $q(x) = p(x|y)$

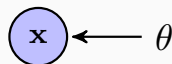
# Deterministic Approximation



$$\begin{aligned}\log p(y) &\geq - \int q(x) \log q(x) dx + \int q(x) \log p(x, y) dx \\ &= \mathbb{E}_{q(x)} [\log p(x, y)] - H(q(x)) = \mathcal{L}(q(x))\end{aligned}$$

- if we maximise the ELBO we,
  - find an approximate posterior
  - lower bound the marginal likelihood
- *maximising*  $p(y)$  is learning
- finding  $q(x) \approx p(x|y)$  is prediction

## Lower Bound



$$p(y) = \int_x p(y|x)p(x) = \frac{p(y|x)p(x)}{p(x|y)}$$

$$q_\theta(x) \approx p(x|y)$$

# Why is this useful?

Why is this a sensible thing to do?

– Ryan Adams<sup>1</sup>

---

<sup>1</sup>Talking Machines Season 2, Episode 5

# Why is this useful?

## Why is this a sensible thing to do?

- If we can't formulate the joint distribution there isn't much we can do

– Ryan Adams<sup>1</sup>

---

<sup>1</sup>[Talking Machines Season 2, Episode 5](#)

## Why is this a sensible thing to do?

- If we can't formulate the joint distribution there isn't much we can do
- Taking the expectation of a log is usually easier than the expectation

– Ryan Adams<sup>1</sup>

---

<sup>1</sup>[Talking Machines Season 2, Episode 5](#)



# Why is this useful?

## Why is this a sensible thing to do?

- If we can't formulate the joint distribution there isn't much we can do
- Taking the expectation of a log is usually easier than the expectation
- We are allowed to choose the distribution to take the expectation over

– Ryan Adams<sup>1</sup>

---

<sup>1</sup>Talking Machines Season 2, Episode 5

## How to choose Q?

$$\mathcal{L}(q(x)) = \mathbb{E}_{q(x)} [\log p(x, y)] - H(q(x))$$

- We have to be able to compute an expectation over the joint distribution
- The second term should be trivial

$$q(\mathbf{X}) = \prod_i q_i(\mathbf{x}_i)$$

$$\mathcal{L}(q_j) = \mathcal{L}_j(q_j) + \mathcal{L}_{\neg j}(q_{\neg j}),$$

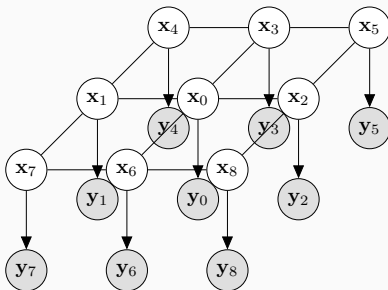
- Model originating if Physics
- We model marginals rather than the full distribution
- We can update each distribution in turn and cycle

1. Formulate joint distribution over data and latent parameters

1. Formulate joint distribution over data and latent parameters
2. Formulate fully factorised approximative posterior over latent variables

1. Formulate joint distribution over data and latent parameters
2. Formulate fully factorised approximative posterior over latent variables
3. Fit marginal approximation by making bound tight

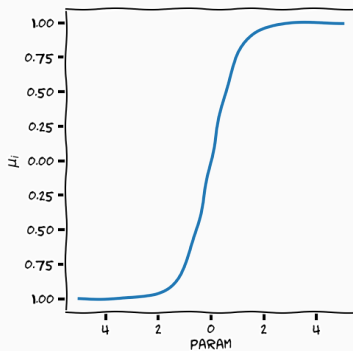
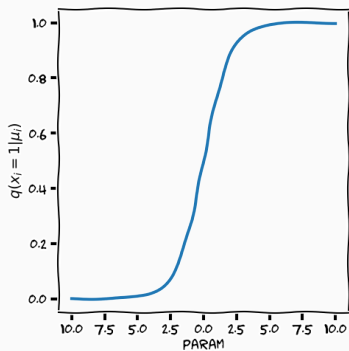
1. Formulate joint distribution over data and latent parameters
2. Formulate fully factorised approximative posterior over latent variables
3. Fit marginal approximation by making bound tight
4. Iterate through variables



$$q(\mathbf{x}, \boldsymbol{\mu}) = \prod_i^N q(x_i, \mu_i)$$

$$\mu_i = \mathbb{E}[x_i]$$





## Summary

---

- Variational methods can be **very** efficient
  - really fun to work with
- Can be made black-box [2]
- Will never be correct
- Provides us with approximative posterior for predictions

eof

## References

---



Christopher M. Bishop.

***Pattern Recognition and Machine Learning (Information Science and Statistics).***

Springer-Verlag New York, Inc., Secaucus, NJ, USA, 2006.



Rajesh Ranganath, Sean Gerrish, and David Blei.

**Black Box Variational Inference.**

In Samuel Kaski and Jukka Corander, editors, *Proceedings of the Seventeenth International Conference on Artificial Intelligence and Statistics*, volume 33 of *Proceedings of Machine Learning Research*, pages 814–822, Reykjavik, Iceland, 22–25 Apr 2014. PMLR.