

Vladislav Belavin, Maxim Borisyak

# Variational Optimization



2021



Yandex



EPFL

S<sup>3</sup>T  
Schaffhausen  
Institute of  
Technology

# Variational bound



# Variational bound

Variational Optimization replaces problem:

$$f(\theta) \rightarrow \min_{\theta};$$

with:

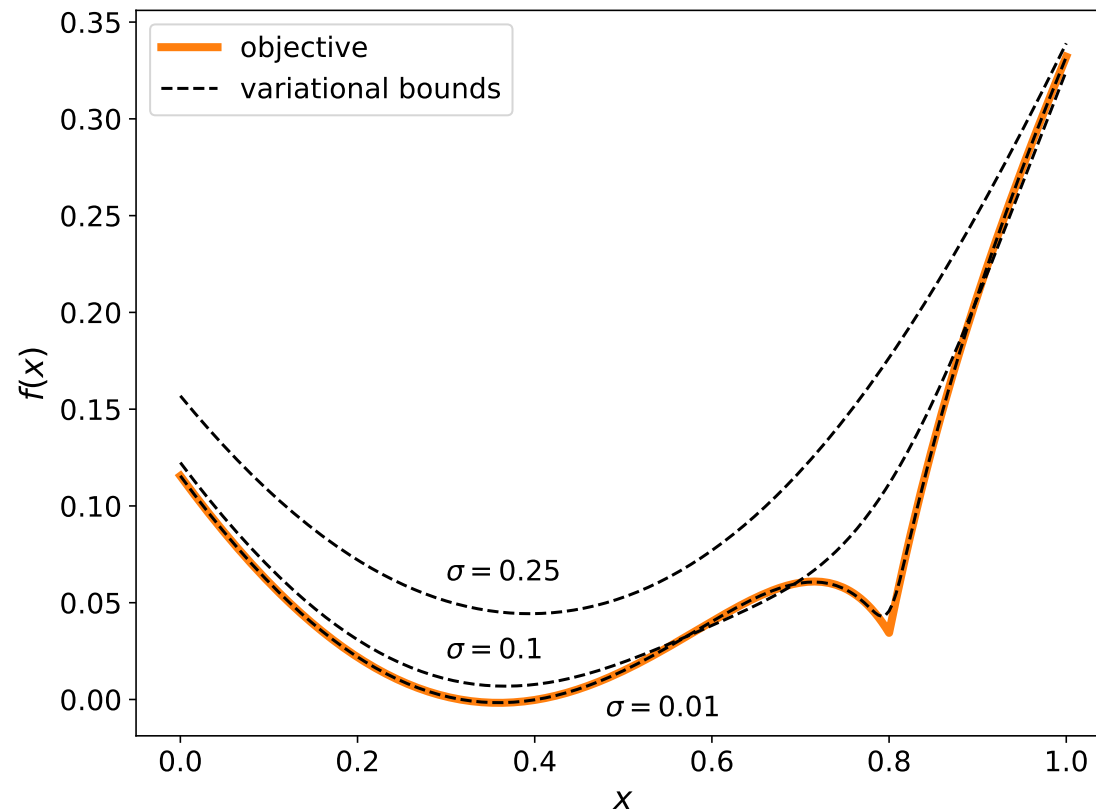
$$J(\psi) = \mathbb{E}_{\theta \sim P(\cdot | \psi)} f(\theta) \rightarrow \min_{\psi}$$

where:

- ▶  $J(\psi)$  - variational bound;
- ▶  $P(\cdot | \psi)$  - search distribution.

This variational bound is not the only one, nevertheless, it is the most common in Variational Optimization.

# Variational bound: example



# Properties

$$J(\psi) = \mathbb{E}_{\theta \sim P(\cdot | \psi)} f(\theta) \rightarrow \min_{\psi}$$

- ▶ upper bound:

$$\forall \psi : J(\psi) \geq \min f(\theta);$$

- ▶ if  $P(\cdot | \psi)$  is allowed to (nearly) collapse into delta function, then

$$P(\cdot | \psi^*) \approx \delta(\theta^*).$$

# Gradient of the variational bound

$$\begin{aligned}\frac{\partial}{\partial \psi} J(\psi) &= \frac{\partial}{\partial \psi} \mathbb{E}_{\theta \sim P(\cdot | \psi)} f(\theta) = \\ &= \frac{\partial}{\partial \psi} \int_{\theta} d\theta f(\theta) P(\theta | \psi) = \\ &= \int_{\theta} d\theta f(\theta) \frac{\partial}{\partial \psi} P(\theta | \psi) = \\ &= \int_{\theta} d\theta f(\theta) P(\theta | \psi) \frac{\partial}{\partial \psi} \log P(\theta | \psi) = \\ &= \mathbb{E}_{\theta \sim P(\cdot | \psi)} f(\theta) \frac{\partial}{\partial \psi} \log P(\theta | \psi)\end{aligned}$$

# Gradient of the variational bound

$$\nabla_{\psi} J(\psi) = \mathbb{E}_{\theta \sim P(\cdot | \psi)} f(\theta) \nabla_{\psi} \log P(\theta | \psi)$$

$\nabla_{\psi} J(\psi)$  does not depend on  $\nabla_{\theta} f(\theta)$

# Gradient of the variational bound

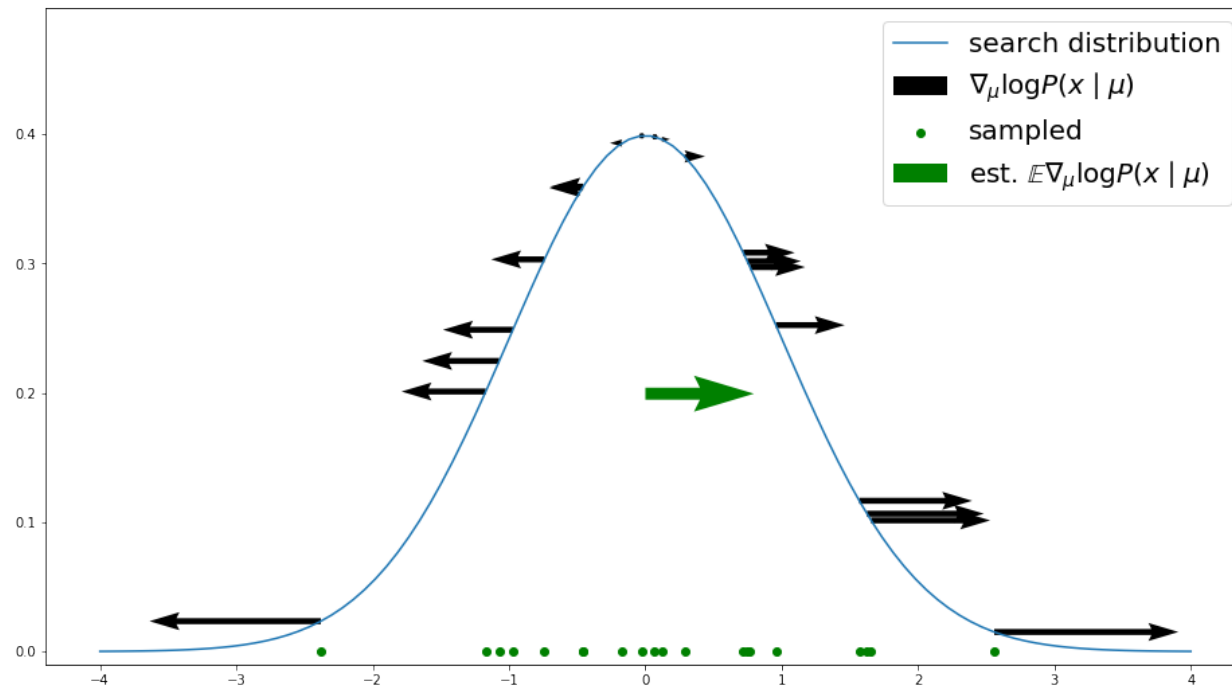


Figure 1: Gradient of  $\mathcal{N}(x; \mu, \sigma)$  w.r.t.  $\mu$ , i.e.  $\nabla_{\mu} \log \mathcal{N}(x; \mu, \sigma)$



# Gradient of the variational bound

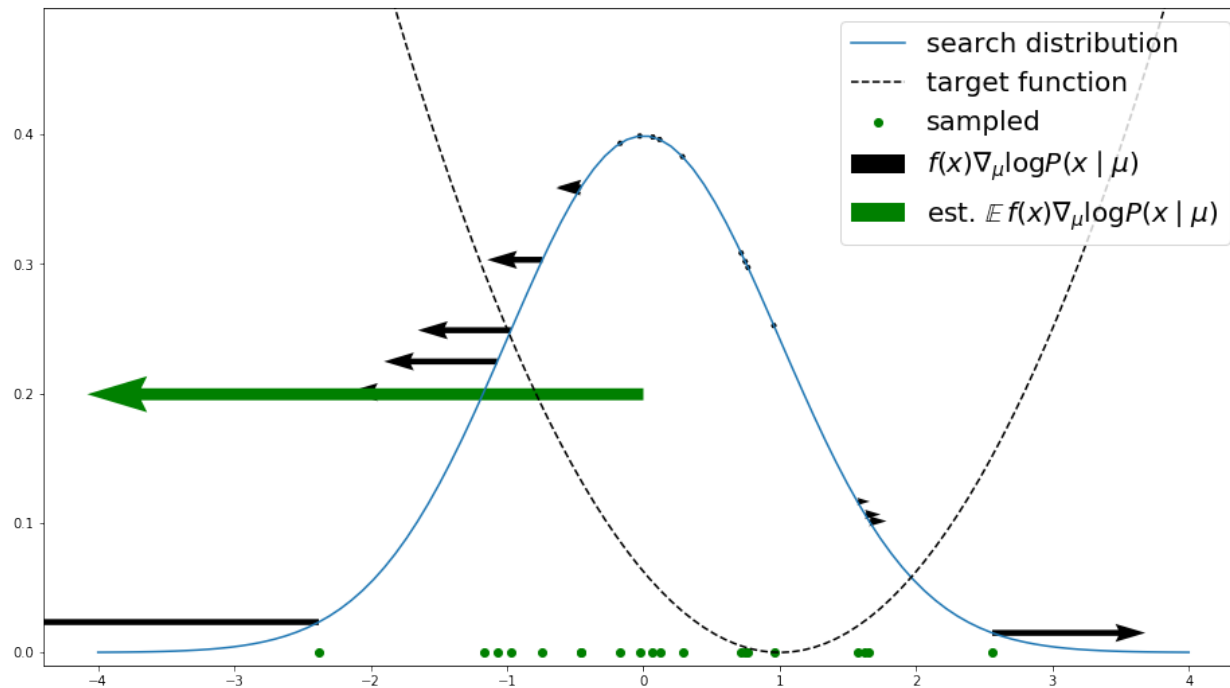
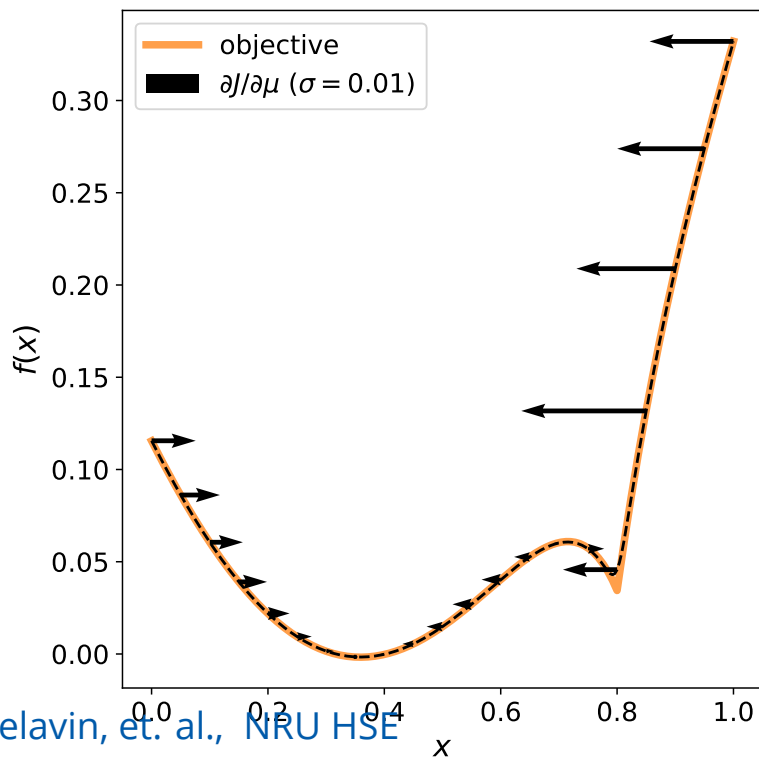


Figure 2: Gradient of variational bound with  $\mathcal{N}(x; \mu, \sigma)$  as search distribution w.r.t.  $\mu$ , i.e.  $f(x)\nabla_{\mu} \log \mathcal{N}(x; \mu, \sigma)$

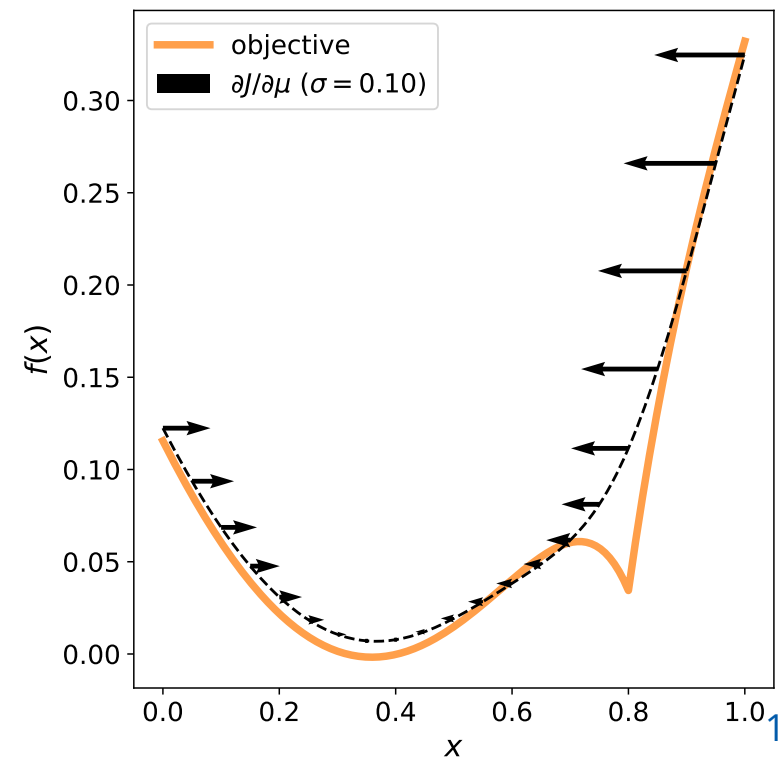
# Gradient of the variational bound

Trade-off: higher  $\sigma$  – smooth gradients, ability to avoid local minima, lower  $\sigma$  – closer to original objective, but more prone to noise.



V. Belavin, et. al., NRU HSE

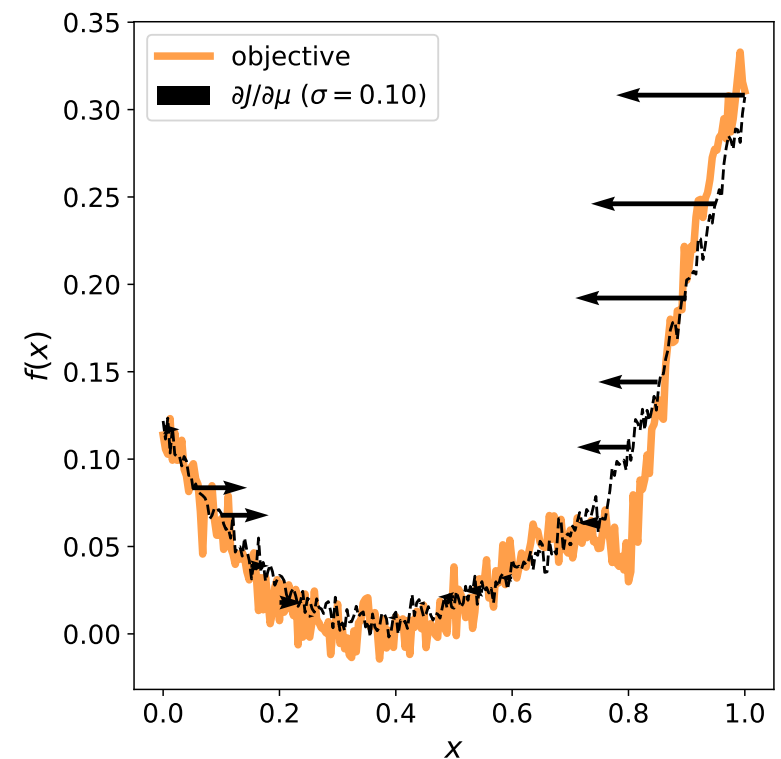
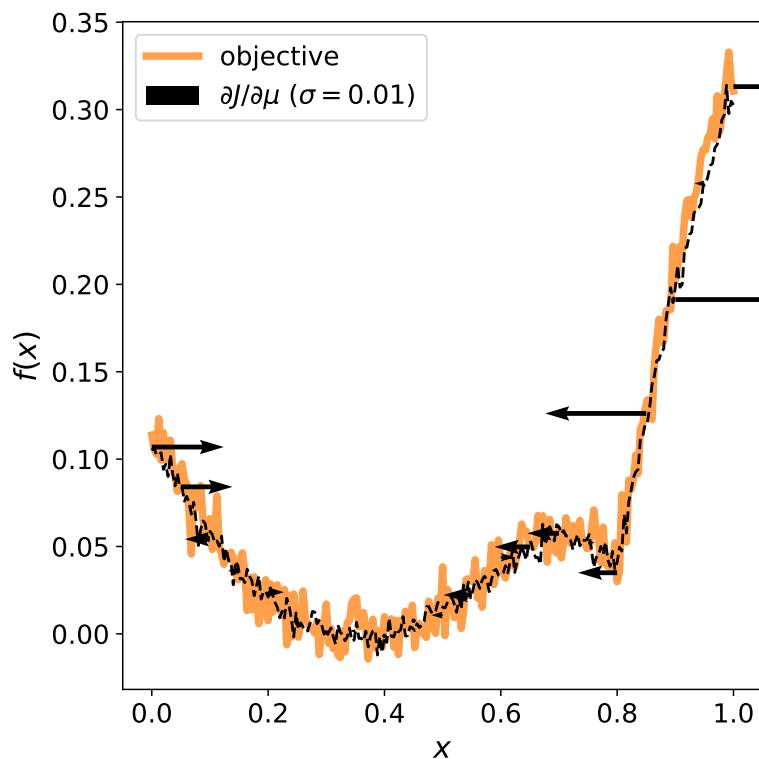
July 24, 2020



10/27

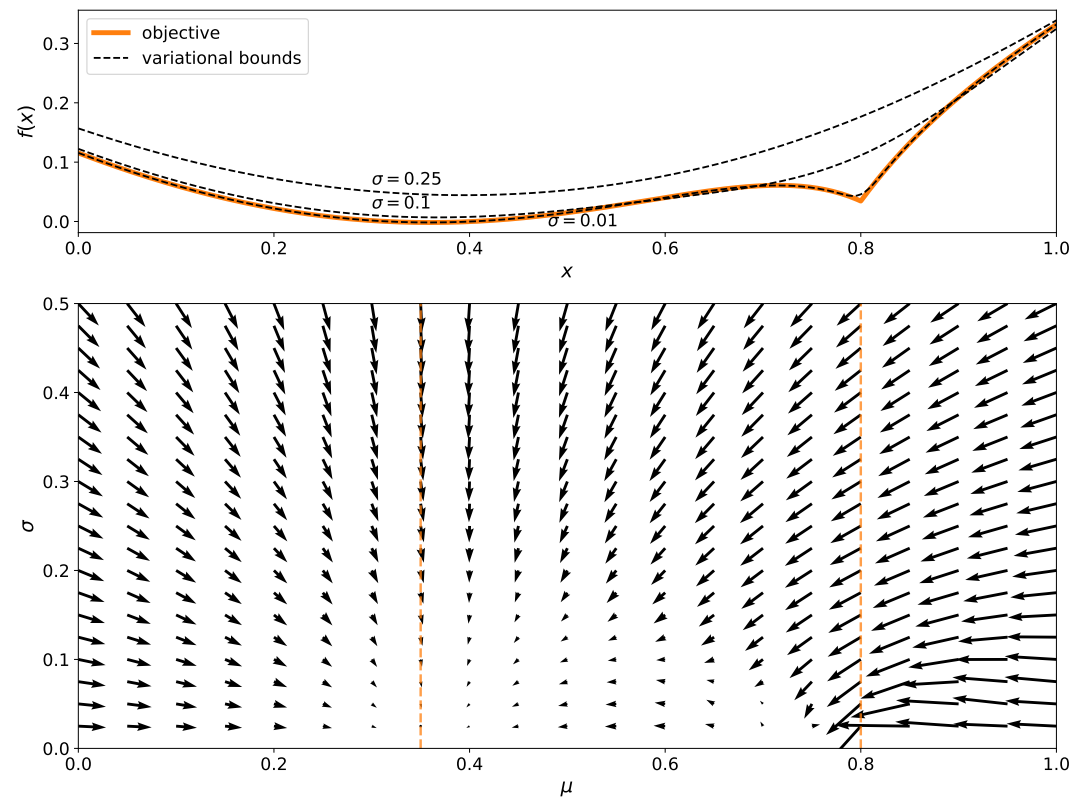
# Gradient of the variational bound

Higher  $\sigma$  also negates sampling noise effects.



# Gradient of the variational bound

In practice,  $\sigma$  usually collapses from any initialization.



# Variational optimization



# SGD-VO

---

**Algorithm 1** SGD-VO

---

```
1: initialize  $P(\cdot \mid \psi)$ 
2: while not converged do
3:   sample  $\theta$  from  $P(\cdot \mid \psi)$ ;
4:    $\nabla_{\psi} J(\psi) \leftarrow f(\theta) \nabla_{\psi} \log P(\theta \mid \psi)$ ;
5:    $\psi \leftarrow \psi - \gamma \nabla_{\psi} J(\psi)$ ;
6: end while
```

---

# Variational Optimization

- ▶ allows usage of stochastic gradient methods for black-box problems:
  - VO is much slower in contrast to using analytical gradient;
- ▶ search distribution is chosen to be simple:
  - e.g. normal distribution;
- ▶ dimensionality of the problem can be retained:
  - at least 1 additional parameter to allow the search distribution to collapse ( $\sigma$ );
  - $2 \cdot n$  parameters for a normal distribution with diagonal covariance;
  - $O(n^2)$  for a full covariance matrix.

# Discrete variables





# Variational optimization

Variational optimization could be straightforwardly used to estimate gradient with the well know formula.

$$\nabla_{\psi} J(\psi) = \mathbb{E}_{\theta \sim P(\cdot | \psi)} f(\theta) \nabla_{\psi} \log P(\theta | \psi)$$

For example, imagine the simplest problem, where  $\frac{df}{dx}$  do not exist:

$$f(x) = \begin{cases} 0.45, & x = 0 \\ 0.53, & x = 1 \end{cases}$$

# Variational optimization

$$f(x) = \begin{cases} 0.45, & x = 0 \\ 0.53, & x = 1 \end{cases}$$

As a search distribution we will choose  $p(x|\psi) = \text{Bernoulli}(x|\psi)$ :

$$J(\psi) = \mathbb{E}_{x \sim P(x|\psi)} f(x) = 0.53\psi + 0.45(1 - \psi) = 0.08\psi + 0.45$$

Voila! Even though  $f(x)$  can not be differentiated, we are still able to compute gradient of variational bound ( $\mathbb{E}_{x \sim P(x|\psi)} f(x)$ ) w.r.t.  $\psi$  and optimize upper bound.

# Discrete variables in Deep Learning

Neural networks with discrete random variables are a powerful technique for representing processes encountered in language modeling, attention mechanisms, and robotics control. Discrete representations are often more interpretable. However, networks with discrete variables are hard to train.

**Cons of VO:** Monte-Carlo estimate of  $\nabla_{\psi} \mathbb{E}_{x \sim P(x|\psi)} f(x)$  has high variance and, consequently, slow convergence. The variance scales linearly with the number of dimensions of the sample vector\*, making it especially challenging to use for categorical distributions.

\*Source: <https://arxiv.org/abs/1401.4082>

# Gumbel Max Trick

Given probabilities  $\{\pi_i\}_{i=1}^n$ ,  $\sum \pi_i = 1$ , we want to be able to sample from this discrete distribution.

**Proposition.** Following procedure:

$$z = \arg \max \{\log \pi_i + G_i\}, \quad G_i = -\log(-\log u_i), \quad u_i \sim U[0, 1],$$

generates samples from discrete probability distribution defined by  $\{\pi_i\}_{i=1}^n$ .

Yet we still can't backprop through  $z$  w.r.t.  $\pi$ , but it's only a first step.

# Gumbel Softmax Trick

To be able to backprop through discrete sample we will relax our problem with the softmax operation:

$$y_i = \frac{\exp((\log \pi_i + \mathbf{G}_i)/\tau)}{\sum_j \exp((\log \pi_j + \mathbf{G}_j)/\tau)}, \quad \mathbf{G}_i = -\log(-\log u_i), \quad u_i \sim U[0, 1],$$

# Gumbel Softmax Trick

During forward pass:

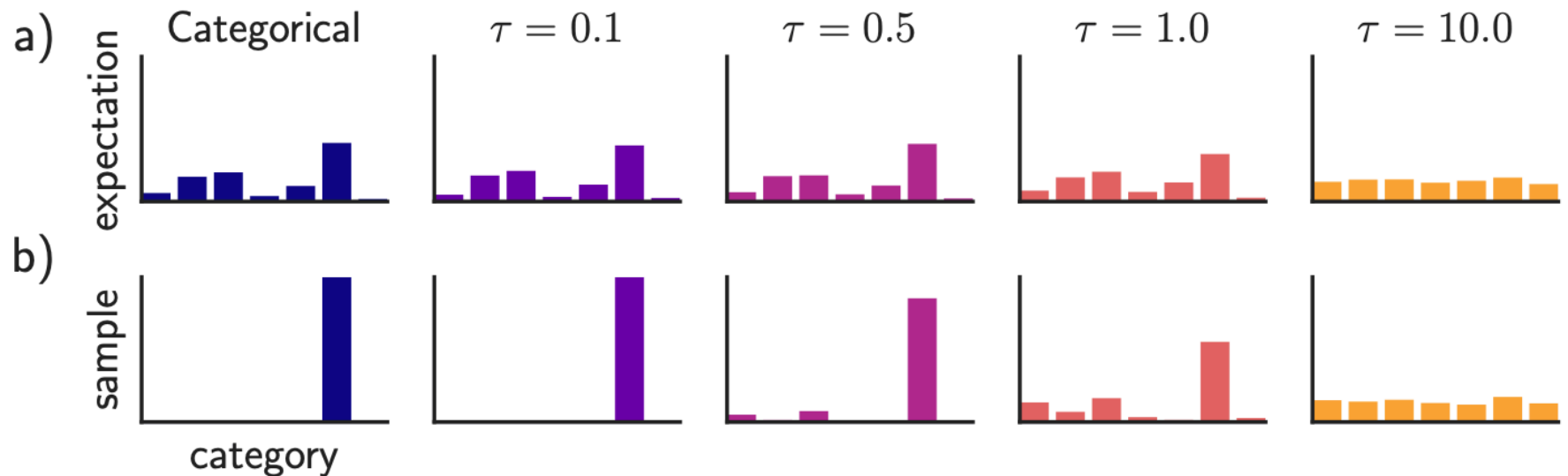
$$z = \text{one-hot}(\arg \max\{y_i\}),$$

During backward pass:

$$\frac{dy}{d\pi} \rightarrow \frac{dz}{d\pi}, \quad \tau \rightarrow 0$$

# Gumbel Softmax Trick

Temperature annealing of gumbel softmax:



# Gradient estimation for discrete variables overview

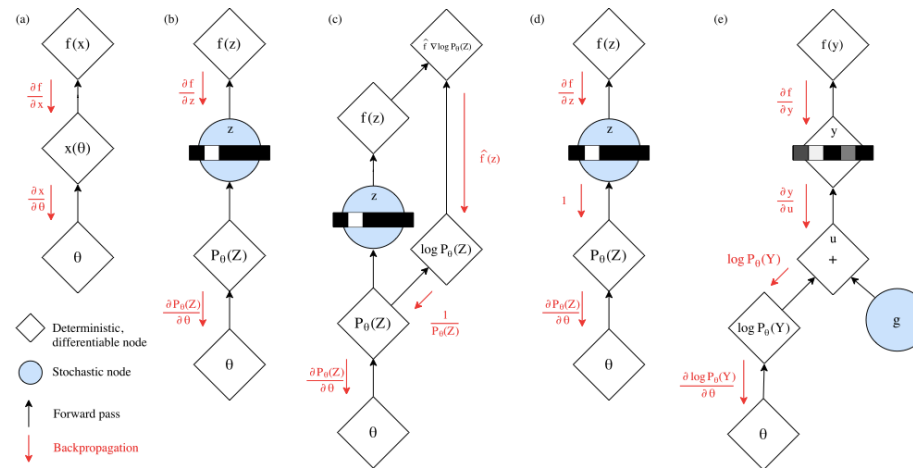


Figure 3: (1) Differentiable linear layer, (2) Non-differentiable discrete layer, (3) variational gradient, (4) Straight-Through estimator, (5) Gumbel-Softmax estimator

Source: <https://arxiv.org/abs/1401.4082>



# Conclusion

## Variational optimization:

- ▶ good for optimization of non-differentiable functions;
- ▶ local optimization algorithm:
  - suffer less from the curse of dimensionality than global algorithms;
  - might converge into local optima;
- ▶ sensitive to the choice of (hyper)parameters.

## Discrete variables:

- ▶ variational optimization could be used;
  - high variance => slow convergence;
- ▶ relaxation methods, like Gumbel Softmax Trick, worth trying.

# Quiz

You want to estimate gradient of the variational bound for the function  $f(x) = \sin x$  with normal search distribution  $p(x|\mu, \sigma^2) = \mathcal{N}(x; \mu, \sigma^2)$ . Choose the correct sample MC estimate  $x_{\text{sample}} \sim \mathcal{N}(x; \mu, \sigma^2)$  of the gradient w.r.t.  $\mu$ :

1.  $\cos(x_{\text{sample}}) \cdot p(x_{\text{sample}}|\mu, \sigma^2)$
2.  $\sin(x_{\text{sample}}) \cdot \frac{x_{\text{sample}} - \mu}{\sigma^2}$
3.  $\cos(x_{\text{sample}}) \cdot \frac{x_{\text{sample}} - \mu}{\sigma^2}$
4.  $\sin(x_{\text{sample}}) \cdot p(x_{\text{sample}}|\mu, \sigma^2)$

# Thank you for your attention!

Vladislav Belavin

 vbelavin@hse.ru

 SchattenGenie

 hse\_lambda

