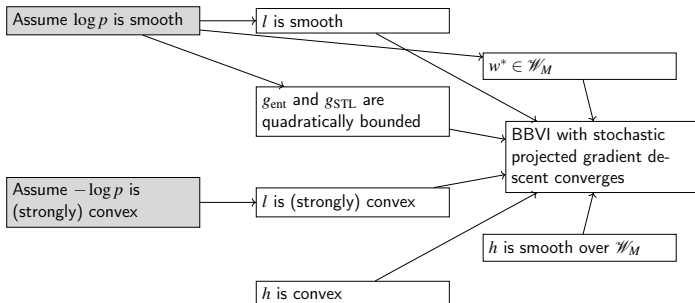# Convergence Guarantees for Variational Inference

Justin Domke, University of Massachusetts Amherst

these slides:   t.ly/sICHy or people.cs.umass.edu/domke/convergence.pdf

# Outline

**Inference**: Given $p(z,x)$ and observed data $x$, approximate $p(z|x)$

**Inference**: Given $p(z,x)$ and observed data $x$, approximate $p(z|x)$

**Variational inference**: ...by choosing some family $q_w(z)$ and minimizing $KL(q_w(z)\|p(z|x))$

**Inference**: Given $p(z,x)$ and observed data $x$, approximate $p(z|x)$

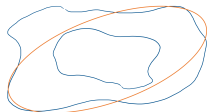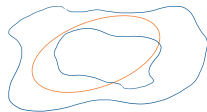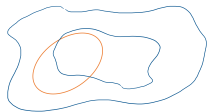**Variational inference**: ...by choosing some family $q_w(z)$ and minimizing $KL(q_w(z)\|p(z|x))$

**Black box variational inference**: ...while only evaluating $\log p(z,x)$ or $\nabla_z \log p(z,x)$.

# Black box VI in practice

- Let $q_w(z)$ be the set of dense Gaussians
- Inialize $w$ somehow.
- Repeat:
    - Get stochastic estimate $g$ of $\nabla_w KL\left(q_w(z)\|p(z|x)\right)$.
    - Take gradient step: $w \leftarrow w - \gamma g$.

# Black box VI in practice

- Let $q_w(z)$ be the set of dense Gaussians
- Inialize $w$ somehow.
- Repeat:
  - Get stochastic estimate $g$ of $\nabla_w KL(q_w(z) \| p(z|x))$.
  - Take gradient step: $w \leftarrow w - \gamma g$.

# Black box VI in practice

- Let $q_w(z)$ be the set of dense Gaussians
- Inialize $w$ somehow.
- Repeat:
  - ▸ Get stochastic estimate $g$ of $\nabla_w KL\left(q_w(z)\|p(z|x)\right)$.
  - ▸ Take gradient step: $w \leftarrow w - \gamma g$.



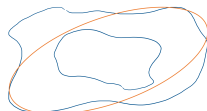Easy to find $g$ via autodiff, seems to work well in practice.
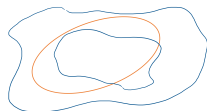
# Black box VI in practice

- Let $q_w(z)$ be the set of dense Gaussians
- Inialize $w$ somehow.
- Repeat:
  - Get stochastic estimate $g$ of $\nabla_w KL(q_w(z) \| p(z|x))$.
  - Take gradient step: $w \leftarrow w - \gamma g$.



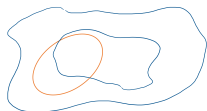Easy to find $g$ via autodiff, seems to work well in practice.

**This talk**: But can we prove anything?

Results on `fires` with *exact* gradients, initialized with $\Sigma = C_0 C_0^\top$.

# Bad news

Can we guarantee anything? If $p(z,x)$ could be *anything*, then no.

# Bad news

Can we guarantee anything? If $p(z,x)$ could be *anything*, then no.

Best we can hope for: If $p$ is "nice" then BBVI optimization is "nice".

# How $p$ might be nice

Plausible properties for $f(z) = -\log p(z,x)$:

# How $p$ might be nice

Plausible properties for $f(z) = -\log p(z, x)$:

- Convex ($\nabla_z^2 f(z) \succeq 0$)
- Strongly convex ($\nabla_z^2 f(z) \succeq cI$)
- Smooth ($\nabla_z^2 f(z) \preceq MI$)

# How $p$ might be nice

Plausible properties for $f(z) = -\log p(z, x)$:

- Convex ($\nabla_z^2 f(z) \succeq 0$)
- Strongly convex ($\nabla_z^2 f(z) \succeq cI$)
- Smooth ($\nabla_z^2 f(z) \preceq MI$)

| $p(z, x)$ | convex | strongly covex | smooth |
|---|:---:|:---:|:---:|
| Gaussian | ✓ | ✓ | ✓ |
| Bayesian linear regression | ✓ | ✓ | ✓ |
| Bayesian logistic regression | ✓ | ✓ | ✓ |
| Heirarchical logistic regression | ✓ | × | ✓ |

## But is optimization nice?

$$\min_w F(w) := \underbrace{\mathop{\mathbb{E}}_{q_w(z)} \left[ -\log p(z,x) \right]}_{\text{"energy" } l(w)} + \underbrace{\mathop{\mathbb{E}}_{q_w(z)} \log q_w(z)}_{\text{"neg-entropy" } h(w)}$$

## But is optimization nice?

$$\min_w F(w) := \underbrace{\mathbb{E}_{q_w(z)} \left[ -\log p(z,x) \right]}_{\text{"energy" } l(w)} + \underbrace{\mathbb{E}_{q_w(z)} \log q_w(z)}_{\text{"neg-entropy" } h(w)}$$

Assume henceforth that $q_w(z) = \mathcal{N}(z|m, CC^\top), \quad w = (m, C)$.

## But is optimization nice?

$$\min_w F(w) := \underbrace{\mathop{\mathbb{E}}_{q_w(z)} \left[ -\log p(z,x) \right]}_{\text{"energy" } l(w)} + \underbrace{\mathop{\mathbb{E}}_{q_w(z)} \log q_w(z)}_{\text{"neg-entropy" } h(w)}$$

Assume henceforth that $q_w(z) = \mathcal{N}(z|m, CC^\top), \quad w = (m, C)$.

How stochastic optimization guarantees usually work:

1. Prove that gradient has **bounded noise** (either $\mathbb{E}\|g\|_2^2 \leq b$ or $\mathbb{V}[g] \leq b$)
2. Prove that objective is **convex** or **strongly convex**
3. Prove that objective is Lipschitz **smooth**.

## But is optimization nice?

$$\min_w F(w) := \underbrace{\mathbb{E}_{q_w(z)} \left[ -\log p(z,x) \right]}_{\text{"energy" } l(w)} + \underbrace{\mathbb{E}_{q_w(z)} \log q_w(z)}_{\text{"neg-entropy" } h(w)}$$

Assume henceforth that $q_w(z) = \mathcal{N}(z|m, CC^\top), \quad w = (m, C)$.

How stochastic optimization guarantees usually work:

1. Prove that gradient has **bounded noise** (either $\mathbb{E}\|g\|_2^2 \le b$ or $\mathbb{V}[g] \le b$)
2. Prove that objective is **convex** or **strongly convex**
3. Prove that objective is Lipschitz **smooth**.

**Trouble**: If $p(z|x) = \mathcal{N}(z|0, I)$, then 1 and 3 are false!

# Table of properties

$$F(w) := \underbrace{\mathop{\mathbb{E}}_{q_w(z)}[-\log p(z,x)]}_{\text{"energy" } l(w)} + \underbrace{\mathop{\mathbb{E}}_{q_w(z)}\log q_w(z)}_{\text{"neg-entropy" } h(w)}$$

| Condition on $-\log p(z,x)$ | Consequence |
| --- | --- |
| none | |
| | |
| convex | |
| $c$-strongly convex | |
| | |
| $M$-smooth | |

# Outline

# Neg-entropy

**Theorem**

$h(w)$ is convex, but not strongly convex and not smooth.

# Neg-entropy

## Theorem

$h(w)$ is convex, but not strongly convex and not smooth.

## Proof.

$h(w) = -\log|\det C| + \frac{d}{2}\log(2\pi e)$

$\square$

# Neg-entropy

**Theorem**

$h(w)$ *is convex, but not strongly convex and not smooth.*

**Proof.**

$h(w) = -\log|\det C| + \frac{d}{2}\log(2\pi e) = -\sum_i \log \sigma_i(C) + \text{const.}$

$\square$

Blows up when singular values of $C$ are small.

# Table of properties

$$F(w) := \underbrace{\underset{q_w(z)}{\mathbb{E}}[-\log p(z,x)]}_{\text{"energy" } l(w)} + \underbrace{\underset{q_w(z)}{\mathbb{E}} \log q_w(z)}_{\text{"neg-entropy" } h(w)}$$

| Condition on $-\log p(z,x)$ | Consequence |
|---|---|
| none | $h(w)$ convex (when $C$ symmetric or triangular) |
|  | $h(w)$ *not* strongly convex, *not* smooth |
| convex | |
| $c$-strongly convex | |
| $M$-smooth | |

# Outline

# (Strong) convexity

**Theorem**

*If $-\log p(z,x)$ is convex, then $l(w)$ is also convex.*

**Theorem**

*If $-\log p(z,x)$ is $c$-strongly convex, then $l(w)$ is also $c$-strongly convex.*

# (Strong) convexity

**Theorem**

*If $-\log p(z,x)$ is convex, then $l(w)$ is also convex.*

**Theorem**

*If $-\log p(z,x)$ is $c$-strongly convex, then $l(w)$ is also $c$-strongly convex.*

**Proof.**

Easy.
(Convexity result due to Titsias and Lázaro-gredilla (2014))
(Strong convexity result (D., 2019) generalizes Challis and Barber (2013)) $\qquad\square$

# Smoothness

**Theorem**

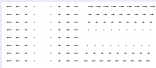*If $\log p(z,x)$ is M-smooth, then $l(w)$ is also M-smooth.*

# Smoothness

## Theorem

If $\log p(z,x)$ is M-smooth, then $l(w)$ is also M-smooth.

## Proof.

Define inner-product space + Bessel's inequality + various exact calculations.

$$\frac{dt_w(u)}{dC_{11}} \qquad \frac{dt_w(u)}{dC_{12}}$$

$$\frac{dt_w(u)}{dC_{21}} \qquad \frac{dt_w(u)}{dC_{22}}$$

$$\frac{dt_w(u)}{dm_1} \qquad \frac{dt_w(u)}{dm_2}$$

**Lemma 2.** $\langle a, b \rangle_s = \mathbb{E}_{u \sim s} a(u)^\top b(u)$ is a valid inner-product on squared-integrable $a: \mathbb{R}^d \to \mathbb{R}^k$.

*Proof.* The space of square integrable functions is $\{a : \mathbb{R}^d \to \mathbb{R}^k \mid \mathbb{E}_{u \sim s} a_i(u)^2 \le \infty \ \forall i \in \{1,\dots,k\}\}$. Since each component $a_i(u)$ and $b_i(u)$ is square-integrable with respect to $s(u)$ we know (by Cauchy-Schwarz) that $\mathbb{E}_{u \sim s} a_i(u) b_i(u) \le \sqrt{\mathbb{E}_{u \sim s} a_i(u)^2} \sqrt{\mathbb{E}_{u \sim s} b_i(u)}$ is finite and real. Therefore, we have by linearity of expectation that

$$\sum_{i=1}^k \mathbb{E}_{u \sim s} a_i(u) b_i(u) = \mathbb{E}_{u \sim s} a(u)^\top b(u)$$
$$= \langle a, b \rangle_s$$

is finite and real for all $a, b \in V_s$. To show that $(V_s, \langle \cdot, \cdot \rangle_s)$ is a valid inner-product space, it is easy to establish all the necessary properties of the inner-product, namely for all $a, b, c \in V_s$,

$\langle a, b \rangle = \langle b, a \rangle$

$\langle \theta a, b \rangle = \theta \langle a, b \rangle$ for $\theta \in \mathbb{R}$

$\langle a + b, c \rangle = \langle a, c \rangle + \langle b, c \rangle$

$\langle a, a \rangle \ge 0$

$\langle a, a \rangle = 0 \Leftrightarrow a = 0$. (Where $0(x)$ is a function that always returns a vector of $k$ zeros.) □

**Lemma 3.** Let $a_i(u) = \frac{d}{du_i} t_w(u)$. This is independent of $w$ and $\frac{dl(w)}{dw} = \langle a_i, \nabla f \circ t_w \rangle_s$.

*Proof.* Now, we can write $l(w)$ as

$$l(w) = \mathbb{E}_{z \sim q_w} f(z) = \mathbb{E}_{u \sim s} f(t_w(u)).$$

Since $t_w(u) = Cu + m$ is an affine function, it's easy to see that both $\frac{d}{du_i} t_w(u)$ and $\frac{d}{du_i} t_w(u)$ are independent of $w$. Therefore, the gradient of $l(w)$ can be written as

$$\nabla_w l(w) = \nabla_w \mathbb{E}_{u \sim s} f(t_w(u))$$
$$= \mathbb{E}_{u \sim s} \nabla_w t_w(u)^\top \nabla f(t_w(u)).$$
$$= \langle a_i, \nabla f \circ t_w \rangle_s.$$
□

**Lemma 4.** If $s$ is standardized, then the functions $\{a_i\}$ are orthonormal in $\langle \cdot, \cdot \rangle_s$.

*Proof.* It is easy to calculate that

$$\frac{d}{dm_i} t_w(u) = e_i$$
$$\frac{d}{dC_{ij}} t_w(u) = e_i u_j,$$

where $e_i$ is the indicator vector in the $i$-th component. Therefore, we have that

$$\mathbb{E}_{u \sim s} \left( \frac{d}{dm_i} t_w(u) \right)^\top \left( \frac{d}{dm_j} t_w(u) \right)$$
$$= \mathbb{E}_{u \sim s} e_i^\top e_j$$
$$= I[i = j]$$

$$\mathbb{E}_{u \sim s} \left( \frac{d}{dm_i} t_w(u) \right)^\top \left( \frac{d}{dm_j} t_w(u) \right)$$
$$= \mathbb{E}_{u \sim s} u_j e_i^\top e_k$$
$$= I[i = k] \mathbb{E}_{u \sim s} u_j$$
$$= 0$$
(since zero mean)

$$\mathbb{E}_{u \sim s} \left( \frac{d}{dC_{ij}} t_w(u) \right)^\top \left( \frac{d}{dC_{kl}} t_w(u) \right)$$
$$= \mathbb{E}_{u \sim s} u_j u_l e_i^\top e_k$$
$$= I[i = k] \mathbb{E}_{u \sim s} u_j u_l$$
$$= I[i = k] I[j = l]$$
(since unit variance and zero mean)

These three identities are equivalent to stating that $\{a_i\}$ are orthonormal in $\langle \cdot, \cdot \rangle_s$. □

**Lemma 5.** If $s$ is standardized, then $\mathbb{E}_{u \sim s} \|t_w(u) - t_v(u)\|_2^2 = \|w - v\|_2^2$.

*Proof.* Let $\Delta m$ and $\Delta S$ denote the difference of the $m$ and $S$ parts of $w$, respectively. We want to calculate

$$\mathbb{E}_{u \sim s} \|t_w(u) - t_v(u)\|_2^2$$
$$= \mathbb{E}_{u \sim s} \|\Delta Cu + \Delta m\|_2^2$$
$$= \mathbb{E}_{u \sim s} \left[ \|(\Delta C)u\|_2^2 + 2\Delta m^\top \Delta Cu + \|\Delta m\|_2^2 \right].$$

It is easy to see that the expectation of the middle term is zero, and the last is a constant. The expectation of the first term is

$$\mathbb{E}_{u \sim s} \|(\Delta C)u\|_2^2 = \mathbb{E}_{u \sim s} u^\top (\Delta C)^\top (\Delta C) u$$
$$= \mathbb{E}_{u \sim s} \mathrm{tr}\left( u^\top (\Delta C)^\top (\Delta C) u \right)$$
$$= \mathbb{E}_{u \sim s} \mathrm{tr}\left( (\Delta C)^\top (\Delta C) u u^\top \right)$$
$$= \mathrm{tr}\left( (\Delta C)^\top (\Delta C) \right) = \|\nabla C\|_F^2$$
(since zero mean and unit variance)

Putting this together gives that

$$\mathbb{E}_{u \sim s} \|t_w(u) - t_v(u)\|_2^2 = \|\Delta C\|_F^2 + \|\Delta m\|_2^2$$
$$= \|w - v\|_2^2.$$
□

*Proof of Thm. 1.* Take two parameter vectors, $w$ and $v$. Apply Lem. 3 to each component of the gradients $\nabla l(w)$ and $\nabla l(v)$ to get that

$$\|\nabla l(w) - \nabla l(v)\|_2^2$$
$$= \sum_i \left( \langle a_i, \nabla f \circ t_w \rangle_s - \langle a_i, \nabla f \circ t_v \rangle_s \right)^2$$
$$= \sum_i \langle a_i, \nabla f \circ t_w - \nabla f \circ t_v \rangle_s^2.$$

Lem. 4 showed that the functions $\{a_i\}$ are orthonormal in the inner-product $\langle \cdot, \cdot \rangle_s$. Thus, by Bessel's inequality,

$$\|\nabla l(w) - \nabla l(v)\|_2^2 \le \|\nabla f \circ t_w - \nabla f \circ t_v\|_s^2 \quad (5)$$
$$= \mathbb{E}_{u \sim s} \|\nabla f(t_w(u)) - \nabla f(t_v(u))\|_2^2$$

where $\|\cdot\|_s$ denotes the norm corresponding to $\langle \cdot, \cdot \rangle_s$. Now apply the smoothness of $f$ to get that

$$\|\nabla l(w) - \nabla l(v)\|_2^2 \le M^2 \mathbb{E}_{u \sim s} \|t_w(u) - t_v(u)\|_2^2 \quad (6)$$
$$= M^2 \|w - v\|_2^2, \quad (7)$$

where the last equality follows from Lem. 5. □

□

# Table of properties

$$F(w) := \underbrace{\mathbb{E}_{q_w(z)}[-\log p(z,x)]}_{\text{"energy" } l(w)} + \underbrace{\mathbb{E}_{q_w(z)}\log q_w(z)}_{\text{"neg-entropy" } h(w)}$$

| Condition on $-\log p(z,x)$ | Consequence |
|---|---|
| none | $h(w)$ convex (when $C$ symmetric or triangular) |
| | $h(w)$ *not* strongly convex, *not* smooth |
| convex | $l(w)$ convex |
| $c$-strongly convex | $l(w)$ $c$-strongly convex |
| $M$-smooth | $l(w)$ $M$-smooth |

# Outline

# Challenge: Non-smooth objective

$$F(w) := \underbrace{\underset{q_w(z)}{\mathbb{E}}\left[-\log p(z,x)\right]}_{\text{"energy" } l(w)} + \underbrace{\underset{q_w(z)}{\mathbb{E}}\log q_w(z)}_{\text{"neg-entropy" } h(w)}$$

**Problem**: $h$ is not smooth. So $F$ is (probably) not smooth.

**Gradient descent** (need $l+h$ smooth)

$$
\begin{aligned}
w' &= w - \gamma(\nabla l(w) + \nabla h(w)) \\
&= \underset{v}{\operatorname{argmin}} \underbrace{l(w) + h(w) + \langle \nabla l(w) + \nabla h(w),\ v - w \rangle}_{\text{local affine approximation of } l(v)+h(v)} + \underbrace{\frac{1}{2\gamma} \|v - w\|_2^2}_{\text{penalty term}}
\end{aligned}
$$

**Gradient descent** (need $l + h$ smooth)

$$
\begin{aligned}
w' &= w - \gamma(\nabla l(w) + \nabla h(w)) \\
&= \underset{v}{\operatorname{argmin}} \underbrace{l(w) + h(w) + \langle \nabla l(w) + \nabla h(w),\ v - w \rangle}_{\text{local affine approximation of } l(v) + h(v)} + \underbrace{\frac{1}{2\gamma} \|v - w\|_2^2}_{\text{penalty term}}
\end{aligned}
$$

**Proximal gradient descent**: (only need $l$ smooth)

$$
\begin{aligned}
w' &= \underset{v}{\operatorname{argmin}} \underbrace{l(w) + \langle \nabla l(w),\ v - w \rangle}_{\text{local affine approximation of } l} + \underbrace{h(v)}_{\text{exact } h} + \underbrace{\frac{1}{2\gamma} \|v - w\|_2^2}_{\text{penalty term}} \\
&= \underset{\gamma h}{\operatorname{prox}} [w - \gamma \nabla l(w)]
\end{aligned}
$$

**Gradient descent** (need $l + h$ smooth)

$$
\begin{aligned}
w' &= w - \gamma(\nabla l(w) + \nabla h(w)) \\
&= \operatorname*{argmin}_{v} \underbrace{l(w) + h(w) + \langle \nabla l(w) + \nabla h(w),\ v - w \rangle}_{\text{local affine approximation of } l(v) + h(v)} + \underbrace{\frac{1}{2\gamma} \|v - w\|_2^2}_{\text{penalty term}}
\end{aligned}
$$

**Proximal gradient descent**: (only need $l$ smooth)

$$
\begin{aligned}
w' &= \operatorname*{argmin}_{v} \underbrace{l(w) + \langle \nabla l(w),\ v - w \rangle}_{\text{local affine approximation of } l} + \underbrace{h(v)}_{\text{exact } h} + \underbrace{\frac{1}{2\gamma} \|v - w\|_2^2}_{\text{penalty term}} \\
&= \operatorname*{prox}_{\gamma h} [w - \gamma \nabla l(w)]
\end{aligned}
$$

Computing $\operatorname{prox}_{\gamma h}[w] = \operatorname{argmin}_v h(v) + \frac{1}{2\gamma} \|w - v\|_2^2$ is easy when $C$ is triangular.

**Gradient descent** (need $l + h$ smooth)

$$
\begin{aligned}
w' &= w - \gamma(\nabla l(w) + \nabla h(w)) \\
&= \operatorname*{argmin}_{v} \underbrace{l(w) + h(w) + \langle \nabla l(w) + \nabla h(w),\ v - w \rangle}_{\text{local affine approximation of } l(v) + h(v)} + \underbrace{\frac{1}{2\gamma} \|v - w\|_2^2}_{\text{penalty term}}
\end{aligned}
$$

**Proximal gradient descent**: (only need $l$ smooth)

$$
\begin{aligned}
w' &= \operatorname*{argmin}_{v} \underbrace{l(w) + \langle \nabla l(w),\ v - w \rangle}_{\text{local affine approximation of } l} + \underbrace{h(v)}_{\text{exact } h} + \underbrace{\frac{1}{2\gamma} \|v - w\|_2^2}_{\text{penalty term}} \\
&= \operatorname*{prox}_{\gamma h} [w - \gamma \nabla l(w)]
\end{aligned}
$$

Computing $\operatorname{prox}_{\gamma h}[w] = \operatorname{argmin}_v h(v) + \frac{1}{2\gamma} \|w - v\|_2^2$ is easy when $C$ is triangular.

**Standard theory**: Converges if $l$ is (strongly) convex and smooth.

# Outline

# Solution guarantees

$$F(w) := \underbrace{\mathop{\mathbb{E}}_{q_w(z)}\left[-\log p(z,x)\right]}_{\text{"energy" } l(w)} + \underbrace{\mathop{\mathbb{E}}_{q_w(z)}\log q_w(z)}_{\text{"neg-entropy" } h(w)}$$

Hmmmm...

# Solution guarantees

$$F(w) := \underbrace{\mathbb{E}_{q_w(z)}\left[-\log p(z,x)\right]}_{\text{"energy" } l(w)} + \underbrace{\mathbb{E}_{q_w(z)}\log q_w(z)}_{\text{"neg-entropy" } h(w)}$$

Hmmmm...

- $h(w) = -\log|\det C| + \text{const.}$ is smooth except when singular values of $C$ are small.

# Solution guarantees

$$F(w) := \underbrace{\mathop{\mathbb{E}}_{q_w(z)}[-\log p(z,x)]}_{\text{"energy" } l(w)} + \underbrace{\mathop{\mathbb{E}}_{q_w(z)}\log q_w(z)}_{\text{"neg-entropy" } h(w)}$$

Hmmmm...

- $h(w) = -\log|\det C| + \text{const.}$ is smooth except when singular values of $C$ are small.
- $h(w)$ also becomes really *large* when the singular values of $C$ are small.

# Solution guarantees

$$F(w) := \underbrace{\underset{q_w(z)}{\mathbb{E}}\,[-\log p(z,x)]}_{\text{"energy" } l(w)} + \underbrace{\underset{q_w(z)}{\mathbb{E}}\,\log q_w(z)}_{\text{"neg-entropy" } h(w)}$$

Hmmmm...

- $h(w) = -\log|\det C| + \text{const.}$ is smooth except when singular values of $C$ are small.
- $h(w)$ also becomes really *large* when the singular values of $C$ are small.
- Maybe the singular values of $C$ can't be too small *at the solution*?

## Solution guarantees

$$F(w) := \underbrace{\mathbb{E}_{q_w(z)} [-\log p(z,x)]}_{\text{"energy" } l(w)} + \underbrace{\mathbb{E}_{q_w(z)} \log q_w(z)}_{\text{"neg-entropy" } h(w)}$$

Hmmmm...

- $h(w) = -\log|\det C| + \text{const.}$ is smooth except when singular values of $C$ are small.
- $h(w)$ also becomes really *large* when the singular values of $C$ are small.
- Maybe the singular values of $C$ can't be too small *at the solution*?
- And maybe we can exploit that somehow?

$$\mathscr{W}_M := \left\{ (m, C) \,\middle|\, \sigma_{\min}(C) \geq \frac{1}{\sqrt{M}} \right\}$$

$$\mathscr{W}_M := \left\{ (m, C) \mid \sigma_{\min}(C) \geq \frac{1}{\sqrt{M}} \right\}$$

### Theorem

*If $\log p(z, x)$ is M-smooth and $w^*$ minimizes $l(w) + h(w)$, then $w^* \in \mathscr{W}_M$. (D. 2020, Thm. 7)*

$$\mathscr{W}_M := \left\{ (m,C) | \sigma_{\min}(C) \geq \frac{1}{\sqrt{M}} \right\}$$

### Theorem

*If $\log p(z,x)$ is M-smooth and $w^*$ minimizes $l(w) + h(w)$, then $w^* \in \mathscr{W}_M$. (D. 2020, Thm. 7)*

### Lemma

*$h(w)$ is M-smooth over $\mathscr{W}_{\mathscr{M}}$. (D. 2020, Lemma 12)*

$$\mathscr{W}_M := \left\{ (m,C) | \sigma_{\min}(C) \geq \frac{1}{\sqrt{M}} \right\}$$

### Theorem

*If $\log p(z,x)$ is M-smooth and $w^*$ minimizes $l(w) + h(w)$, then $w^* \in \mathscr{W}_M$. (D. 2020, Thm. 7)*

### Lemma

*$h(w)$ is M-smooth over $\mathscr{W}_{\mathscr{M}}$. (D. 2020, Lemma 12)*

**Projected gradient descent**:

$$w' = \text{proj}_{\mathscr{W}_M}[w - \gamma(\nabla l(w) + \nabla h(w))]$$

$\text{proj}_{\mathscr{W}_M}[w] = \text{argmin}_{w' \in \mathscr{W}_M} \|w - w'\|_2$

$$\mathscr{W}_M := \left\{ (m,C) \mid \sigma_{\min}(C) \geq \frac{1}{\sqrt{M}} \right\}$$

### Theorem

*If $\log p(z,x)$ is M-smooth and $w^*$ minimizes $l(w) + h(w)$, then $w^* \in \mathscr{W}_M$. (D. 2020, Thm. 7)*

### Lemma

*$h(w)$ is M-smooth over $\mathscr{W}_{\mathscr{M}}$. (D. 2020, Lemma 12)*

**Projected gradient descent**:

$$w' = \text{proj}_{\mathscr{W}_M}[w - \gamma(\nabla l(w) + \nabla h(w))]$$

$\text{proj}_{\mathscr{W}_M}[w] = \text{argmin}_{w' \in \mathscr{W}_M} \|w - w'\|_2$ is easy to compute but requires an SVD of $C$.

$$\mathscr{W}_M := \left\{ (m,C) | \sigma_{\min}(C) \geq \frac{1}{\sqrt{M}} \right\}$$

### Theorem
*If $\log p(z,x)$ is M-smooth and $w^*$ minimizes $l(w) + h(w)$, then $w^* \in \mathscr{W}_M$. (D. 2020, Thm. 7)*

### Lemma
*$h(w)$ is M-smooth over $\mathscr{W}_{\mathscr{M}}$. (D. 2020, Lemma 12)*

**Projected gradient descent**:

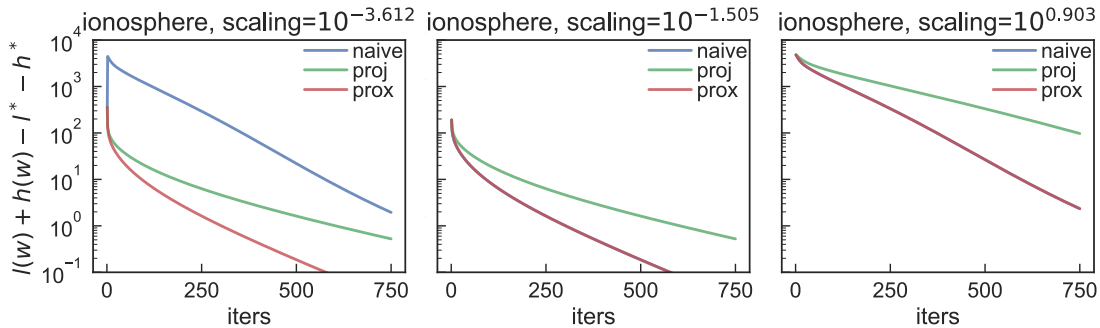$$w' = \text{proj}_{\mathscr{W}_M}[w - \gamma(\nabla l(w) + \nabla h(w))]$$

$\text{proj}_{\mathscr{W}_M}[w] = \text{argmin}_{w' \in \mathscr{W}_M} \|w - w'\|_2$ is easy to compute but requires an SVD of $C$.

**Standard theory**: converges if $l + h$ is (strongly) convex and smooth.

# Table of properties

$$F(w) := \underbrace{\mathbb{E}_{q_w(z)}\left[-\log p(z,x)\right]}_{\text{"energy" } l(w)} + \underbrace{\mathbb{E}_{q_w(z)}\log q_w(z)}_{\text{"neg-entropy" } h(w)}$$

| Condition on $-\log p(z,x)$ | Consequence |
|---|---|
| none | $h(w)$ convex (when $C$ symmetric or triangular) |
| | $h(w)$ *not* strongly convex, *not* smooth |
| | <span style="color:red">$h(w)$ is $M$-smooth over $\mathscr{W}_M$</span> |
| convex | $l(w)$ convex |
| $c$-strongly convex | $l(w)$ $c$-strongly convex |
| | |
| $M$-smooth | $l(w)$ $M$-smooth |
| | <span style="color:red">$w^* \in \mathscr{W}_M$</span> |

Bayesian logistic regression. ("Exact" gradients by reducing evaluation of 1-D integral, precomputed using numerical quadrature.)

# Outline

# Summary so far

BBVI with proximal or projected gradient descent converges, assuming:

1. $-\log p(z,x)$ is smooth
2. $-\log p(z,x)$ is (strongly) convex
3. You can compute the exact gradient.

# Summary so far

BBVI with proximal or projected gradient descent converges, assuming:

1. $-\log p(z, x)$ is smooth $\longleftarrow$ Sometimes true
2. $-\log p(z, x)$ is (strongly) convex $\longleftarrow$ Sometimes true
3. You can compute the exact gradient. $\longleftarrow$ Almost never true

## Estimating gradients

Can "reparameterize" using $t_w(u) = Cu + m$:

$$l(w) = - \mathop{\mathbb{E}}_{q_w(z)} \log p(z,x) = - \mathop{\mathbb{E}}_{\mathcal{N}(u|0,I)} \log p(t_w(u),x).$$

## Estimating gradients

Can "reparameterize" using $t_w(u) = Cu + m$:

$$l(w) = - \mathop{\mathbb{E}}_{q_w(z)} \log p(z,x) = - \mathop{\mathbb{E}}_{\mathcal{N}(u|0,I)} \log p(t_w(u),x).$$

### Definition

Typical gradient estimator (for $\nabla l(w)$):

$$g_{\text{energy}} = -\nabla_w \log p\left(t_w(u),x\right)$$

## Estimating gradients

Can "reparameterize" using $t_w(u) = Cu + m$:

$$l(w) = - \mathop{\mathbb{E}}_{q_w(z)} \log p(z,x) = - \mathop{\mathbb{E}}_{\mathcal{N}(u|0,I)} \log p(t_w(u),x).$$

### Definition

Typical gradient estimator (for $\nabla l(w)$):

$$g_{\text{energy}} = -\nabla_w \log p\left(t_w(u),x\right)$$

### Definition

Other gradient estimators (for $\nabla l(w) + \nabla h(w)$):

$$
\begin{aligned}
g_{\text{ent}} &= -\nabla_w \log p\left(t_w(u),x\right) + \nabla_w h(w) \\
g_{\text{STL}} &= -\nabla_w \log p\left(t_w(u),x\right) + \left[\nabla_w \log q_v(t_w(u))\right]_{v=w}
\end{aligned}
$$

## Quadratic bounds

Stochastic optimization proofs often assume $\mathbb{E}\|g\|_2^2$ (or $\mathbb{V}[g]$) is *uniformly* bounded. Not true for us!

# Quadratic bounds

Stochastic optimization proofs often assume $\mathbb{E}\|g\|_2^2$ (or $\mathbb{V}[g]$) is *uniformly* bounded. Not true for us!

### Definition

A gradient estimator $g$ for $\nabla\phi$ is **quadratically bounded** with parameters $(a, b, w^*)$ if $\mathbb{E}[g] = \nabla\phi(w)$ and

$$\mathbb{E}\|g\|_2^2 \le a\|w - w^*\|_2^2 + b.$$

# Quadratic bounds

Stochastic optimization proofs often assume $\mathbb{E}\|g\|_2^2$ (or $\mathbb{V}[g]$) is *uniformly* bounded. Not true for us!

## Definition

A gradient estimator $g$ for $\nabla\phi$ is **quadratically bounded** with parameters $(a, b, w^*)$ if $\mathbb{E}[g] = \nabla\phi(w)$ and

$$\mathbb{E}\|g\|_2^2 \leq a\|w - w^*\|_2^2 + b.$$

## Theorem

*If $\log p(z, x)$ is M-smooth, then $g_{\text{energy}}$, $g_{\text{ent}}$, and $g_{\text{STL}}$ are all quadratically bounded (D., 2019, D., Garrigos, and Gower, 2023)*

# Table of properties

$$F(w) := \underbrace{\mathbb{E}_{q_w(z)}[-\log p(z,x)]}_{\text{"energy" } l(w)} + \underbrace{\mathbb{E}_{q_w(z)}\log q_w(z)}_{\text{"neg-entropy" } h(w)}$$

| Condition on $-\log p(z,x)$ | Consequence |
|---|---|
| none | $h(w)$ convex (when $C$ symmetric or triangular) |
| | $h(w)$ *not* strongly convex, *not* smooth |
| | $h(w)$ is $M$-smooth over $\mathcal{W}_M$ |
| convex | $l(w)$ convex |
| $c$-strongly convex | $l(w)$ $c$-strongly convex |
| | |
| $M$-smooth | $l(w)$ $M$-smooth |
| | $w^* \in \mathcal{W}_M$ |
| | gradient estimators quadratically bounded |

# Outline

# An optimization "hole"

We have:

- Varying noise (quadratically bounded).
- Composite non-smooth objective.
- Objective is smooth inside of $\mathscr{W}_M$, but not *locally* smooth.

# An optimization "hole"

We have:

- Varying noise (quadratically bounded).
- Composite non-smooth objective.
- Objective is smooth inside of $\mathscr{W}_M$, but not *locally* smooth.

Questions:

- Does proximal gradient descent work with quadratically bounded noise?
- Does projected gradient descent work with quadratically bounded noise?

# New optimization theory

Does stochastic proximal gradient descent work with quadratically bounded noise?

# New optimization theory

Does stochastic proximal gradient descent work with quadratically bounded noise?

### Theorem

*Yes. Converges at a $1/T$ rate if objective is smooth and strongly convex, or $1/\sqrt{T}$ if smooth and merely convex. (D., Gairrigos, and Gower, 2023, Thms. 7+8)*

# New optimization theory

Does stochastic proximal gradient descent work with quadratically bounded noise?

### Theorem

*Yes. Converges at a $1/T$ rate if objective is smooth and strongly convex, or $1/\sqrt{T}$ if smooth and merely convex. (D., Gairrigos, and Gower, 2023, Thms. 7+8)*

Does stochastic projected gradient descent work with quadratically bounded noise?

# New optimization theory

Does stochastic proximal gradient descent work with quadratically bounded noise?

### Theorem

*Yes. Converges at a $1/T$ rate if objective is smooth and strongly convex, or $1/\sqrt{T}$ if smooth and merely convex. (D., Gairrigos, and Gower, 2023, Thms. 7+8)*

Does stochastic projected gradient descent work with quadratically bounded noise?

### Theorem

*Yes. Converges at a $1/T$ rate if objective is smooth and strongly convex, or $1/\sqrt{T}$ if smooth and merely convex. (D., Gairrigos, and Gower, 2023, Thms. 10+11)*
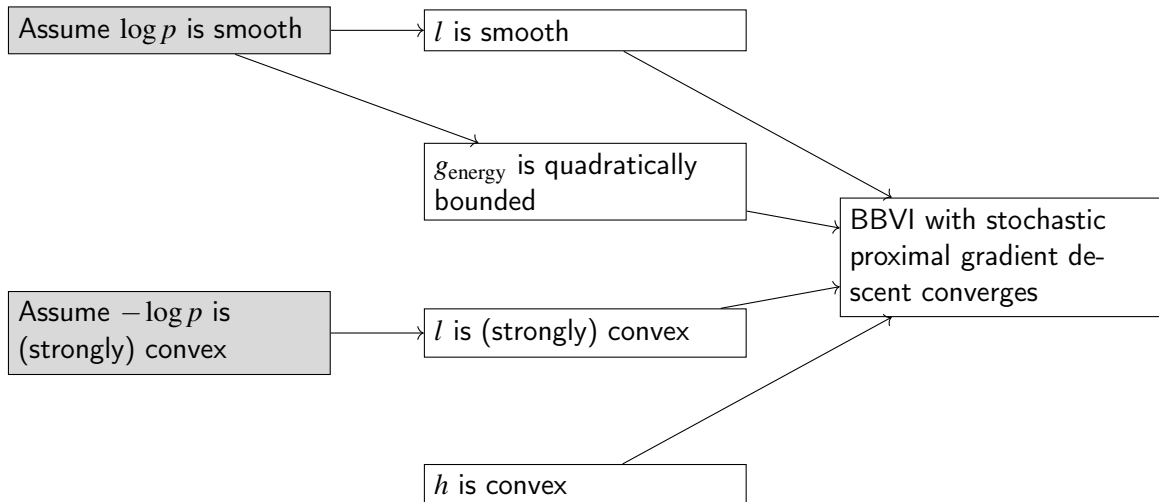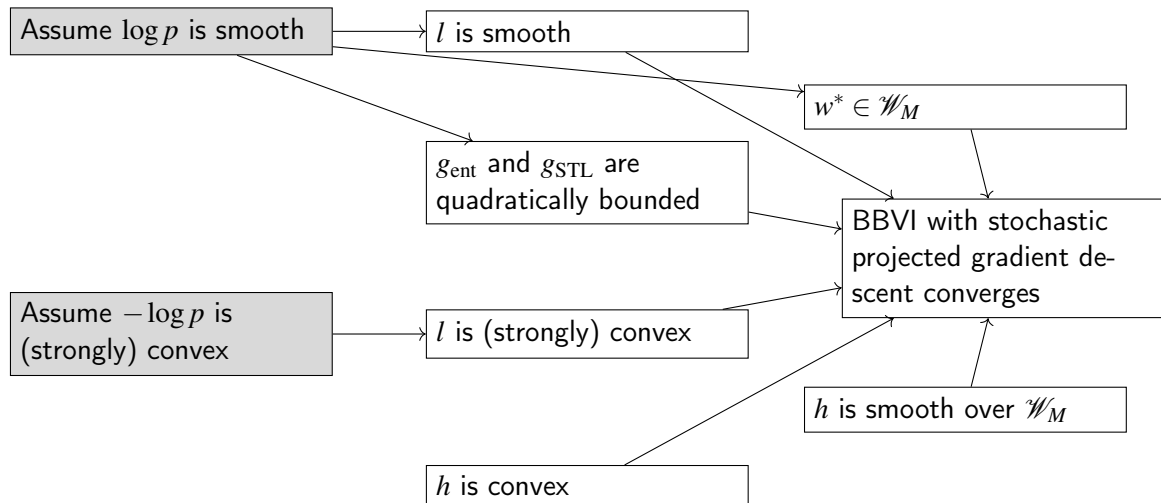
Putting the pieces together (proximal gradient descent)

# Putting the pieces together (proximal gradient descent)

Putting the pieces together (projected gradient descent)

# Putting the pieces together (projected gradient descent)

# Putting the pieces together

## Theorem

*If $-\log p(z,x)$ is M-smooth and (strongly) convex, then stochastic proximal gradient descent using the $g_{\mathrm{energy}}$ estimator with a dense Gaussian variational family with triangular $C$ with an appropriate stepsize sequence converges to the optimum of the ELBO at a $1/\sqrt{T}$ $(1/T)$ rate. (D., Gairrigos, and Gower, 2023, Cor. 12)*

# Putting the pieces together

## Theorem

*If $-\log p(z,x)$ is M-smooth and (strongly) convex, then stochastic proximal gradient descent using the $g_{\text{energy}}$ estimator with a dense Gaussian variational family with triangular $C$ with an appropriate stepsize sequence converges to the optimum of the ELBO at a $1/\sqrt{T}$ $(1/T)$ rate. (D., Gairrigos, and Gower, 2023, Cor. 12)*

## Theorem

*If $-\log p(z,x)$ is M-smooth and (strongly) convex, then stochastic projected gradient descent (projecting onto $\mathcal{W}_M$) using either the $g_{\text{STL}}$ or $g_{\text{ent}}$ estimators with a dense Gaussian variational family with symmetric $C$ with an appropriate stepsize sequence converges to the optimum of the ELBO at a $1/\sqrt{T}$ $(1/T)$ rate. (D., Gairrigos, and Gower, 2023, Cor. 13)*

# Outline

# Related work

- Kim et al. (2023) give a similar $1/T$ rate for proximal SGD using $g_{\text{energy}}$ with smoothness and strong convexity.

## Related work

- Kim et al. (2023) give a similar $1/T$ rate for proximal SGD using $g_{\text{energy}}$ with smoothness and strong convexity.

- Xu and Campbell (2023) give a $1/\sqrt{T}$ rate for projected-SGD using $g_{\text{ent}}$ with a particular rescaling which is *asymptotic* in the number of observations (☺) and *local* (☺) but does not require convexity (☺).

# Related work

- Kim et al. (2023) give a similar $1/T$ rate for proximal SGD using $g_{\text{energy}}$ with smoothness and strong convexity.

- Xu and Campbell (2023) give a $1/\sqrt{T}$ rate for projected-SGD using $g_{\text{ent}}$ with a particular rescaling which is *asymptotic* in the number of observations (☺) and *local* (☺) but does not require convexity (☺).

- Lambert et al. (2022) give a $1/T$ rate for a VI-like SGD algorithm from a discretization of a Wasserstein gradient flow with smoothness+strong convexity. Diao et al. (2023) give a related proximal with a $1/T$ rate or $1/\sqrt{T}$ with just convexity. These require the Hessian of the log-posterior (☺) but are very beautiful (☺).

# Open questions

- Why does *regular* SGD seem to work so well?

## Open questions

- Why does *regular* SGD seem to work so well?
- Guarantees with Adam instead of SGD?

# Open questions

- Why does *regular* SGD seem to work so well?

- Guarantees with Adam instead of SGD?

- Guarantees without assuming we know smoothness/strong convexity constants?

# Open questions

- Why does *regular* SGD seem to work so well?

- Guarantees with Adam instead of SGD?

- Guarantees without assuming we know smoothness/strong convexity constants?

- Guarantees without assuming smoothness or (strong) convexity at all?

## Open questions

- Why does *regular* SGD seem to work so well?
- Guarantees with Adam instead of SGD?
- Guarantees without assuming we know smoothness/strong convexity constants?
- Guarantees without assuming smoothness or (strong) convexity at all?
- Guarantees with more general variational families (e.g. normalizing flows)?

## Open questions

- Why does *regular* SGD seem to work so well?
- Guarantees with Adam instead of SGD?
- Guarantees without assuming we know smoothness/strong convexity constants?
- Guarantees without assuming smoothness or (strong) convexity at all?
- Guarantees with more general variational families (e.g. normalizing flows)?
- Is this "inference research" or "optimization research"?

## Open questions

- Why does *regular* SGD seem to work so well?

- Guarantees with Adam instead of SGD?

- Guarantees without assuming we know smoothness/strong convexity constants?

- Guarantees without assuming smoothness or (strong) convexity at all?

- Guarantees with more general variational families (e.g. normalizing flows)?

- Is this "inference research" or "optimization research"?

# Thank you!

these slides:  `t.ly/sICHy` or `people.cs.umass.edu/domke/convergence.pdf`

# Citations

- D. *Provable gradient variance guarantees for black-box variational inference*. NeurIPS 2019.
- D. *Provable smoothness guarantees for black-box variational inference.* ICML 2020.
- D., Gairrigos, and Gower. *Provable convergence guarantees for black-box variational inference.* NeurIPS 2023.
- Kim, Oh, Wu, Ma, and Gardner. *On the convergence and scale parameterizations of black-box variational inference.* NeurIPS 2023.
- Xu and Campbell. *The computational asymptotics of gaussian variational inference and the laplace approximation.* Stat Comput, (32), 2023.
- Lambert, Chewi, Bach, Bonnabel, and Rigollet. *Variational inference via Wasserstein gradient flows.* NeurIPS 2022.
- Diao, Balasubramanian, Chewi, and Salim. *Forward- backward Gaussian variational inference via JKO in the Bures-Wasserstein space.* ICML 2023.

$$F(w) := \underbrace{\mathop{\mathbb{E}}_{q_w(z)}[-\log p(z,x)]}_{\text{"energy" } l(w)} + \underbrace{\mathop{\mathbb{E}}_{q_w(z)} \log q_w(z)}_{\text{"neg-entropy" } h(w)}$$

| Condition on $-\log p(z,x)$ | Consequence |
|---|---|
| none | $h(w)$ convex (when $C$ symmetric or triangular) |
| | $h(w)$ *not* strongly convex, *not* smooth |
| | $h(w)$ is $M$-smooth over $\mathscr{W}_M$ |
| convex | $l(w)$ convex |
| $c$-strongly convex | $l(w)$ $c$-strongly convex |
| | $\|C\|_F^2 + \|m - z^*\|_2^2 \leq \frac{d}{c}$ at solution |
| $M$-smooth | $l(w)$ $M$-smooth |
| | $w^* \in \mathscr{W}_M$ |
| | gradient estimators quadratically bounded |