

Computing Wasserstein barycenters
using gradient descent algorithms.

joint with Sibo Chewi, Tyler Staune & Austin Stromme
(MIT)

1. Averages

$X_1, \dots, X_n \stackrel{iid}{\sim} P$ over \mathbb{R} ,

$$\hat{\theta}_n = \bar{X}_n = \frac{1}{n} \sum_{i=1}^n X_i \quad \theta = \mathbb{E}[X_1]$$

• LLN: $\hat{\theta}_n \rightarrow \theta$

• Quadratic risk: $\mathbb{E}[(\bar{X}_n - \mathbb{E}[X_1])^2] = \frac{\text{Var}(P)}{n}$

is a standard statistical technique

→ Regression, Maximum likelihood, empirical risk minimization,

Note that $\underset{m}{\operatorname{argmin}} \frac{1}{n} \sum_{i=1}^n (X_i - m)^2$ } both metric
& $\underset{m}{\operatorname{argmin}} \mathbb{E}[(X_1 - m)^2]$ } quantities

$$\bar{X}_n = \underset{m}{\operatorname{argmin}} \frac{1}{n} \sum_{i=1}^n d^2(X_i, m)$$

$$E[X] = \operatorname{argmin}_m \int d^2(x, m) dP(x)$$

↳

2. Wasserstein barycenters

This talk: $(\mathbb{R}, |\cdot|) \rightsquigarrow (\mathcal{P}_2(\mathbb{R}^d), W_2) =: \mathcal{W}_2$
 Wasserstein space.

Captures interesting geometric features
 in graphics and Machine learning

Setup: • P probability measure over \mathcal{W}_2

• $b^* = \operatorname{argmin}_{b \in \mathcal{W}_2} \int W_2^2(b, \mu) dP(\mu)$

↳ true barycenter

• Observe $\mu_1, \dots, \mu_n \stackrel{\text{iid}}{\sim} P$ (images, datasets, ...)

$$\hat{b} = \operatorname{argmin}_{b \in \mathcal{W}} \frac{1}{n} \sum_{i=1}^n W_2^2(b, \mu_i)$$

↳ empirical barycenter

Wasserstein barycenters (Agueh & Carlier 2011)

Questions: ∇ Statistical: $E W_2^2(\hat{b}, b^*) \lesssim \frac{\sigma^2}{n}$?

∇ Computational: - How to compute \hat{b}

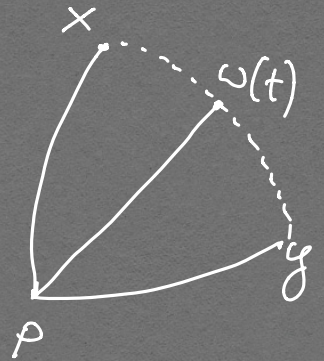
$$\text{at } \mathbb{E} W_2^2(\tilde{b}, b^*) \leq \frac{\sigma^2}{n}$$

Why different from Euclidean case?

Main answer: CURVATURE

Def: $\mathcal{M} = (\mathcal{M}, d)$ geodesic space
 $\text{curv}(\mathcal{M}) \geq 0$ if ...

if $d(x, y) = \|x - y\|$ = Hilbert norm



$$d^2(p, \omega(t)) = (1-t)d^2(p, x) + td^2(p, y) - t(1-t)d^2(x, y)$$

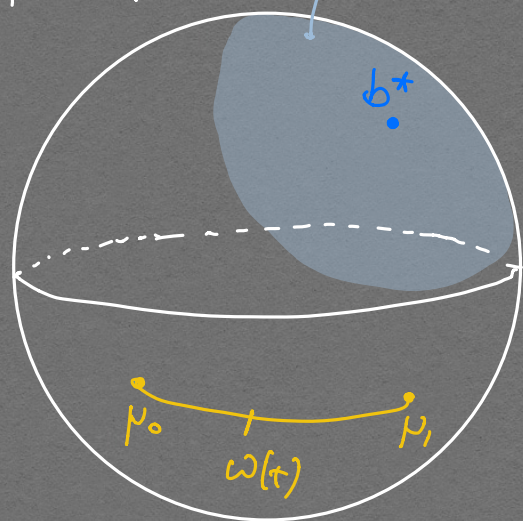
$$\forall p, x, y$$

geodesic $\omega(t) = (1-t)x + ty$

Claim [AGS]:

1. \mathcal{W}_2 is a geodesic space with geodesic between p_0 and p_1 , given by the law of

$$(1-t)x_0 + tx_1, \quad \begin{matrix} \leftarrow & \text{optimally coupled} & \rightarrow \\ x_0 \sim p_0, & x_1 \sim p_1 \end{matrix}$$



2. $\text{Curv}(\mathcal{W}_2) \geq 0$

3. Statistical rates

- Ahidar-Coutrix, Le Gouic, Paris ('18)

$$\mathbb{E} W_2^2(\hat{b}, b^*) \lesssim n^{-1/d}$$

(Also: first rates on general Alexandrov spaces)

- Le Gouic, Paris, R., Strömme ('19) [LPRS]

if $T_{b^* \rightarrow \mu}$ is gradient of $\begin{cases} \beta\text{-smooth} \\ \alpha\text{-str. convex} \end{cases}$, $\forall \mu \in \text{supp}(P)$, $\beta - \alpha < 1$

then

$$\mathbb{E} W_2^2(\hat{b}, b^*) \lesssim \frac{\sigma^2}{n}$$

where $\sigma^2 = \int W_2^2(b, b^*) dP(b)$

(+ similar results for general Alexandrov spaces)

- Koshnin, Spokoyny, Suvorikova

$$\mathbb{E} W_2^2(\hat{b}, b^*) \lesssim \frac{\sigma^2}{n}$$

if P is supported on Gaussians (+ CLT)

Remarks: the Gaussian case is interesting.

- "totally geodesic"

$N(\mu_0, \Sigma_0)$

$N(\mu, \Sigma)$

↖ all Gaussians

- Viable alternative to naive discretization.

[Cuturi & Doucet '14] (First order algorithm to compute Wasserstein barycenters)

$\mu_1, \dots, \mu_n \rightarrow$ discretize \rightarrow algorithm.
sample \swarrow \searrow histogram



\hookrightarrow Statistical modeling \rightarrow summarize a distribution in a few key parameters (e.g.: mean, variance)

4 - Gradient descent over the Wasserstein space

Want to minimize:

$$b \mapsto F(b) = \frac{1}{2} \int W_2^2(b, \mu) dP(\mu) \text{ over } \mathcal{W}_2$$

(e.g. when $P = P_n = \frac{1}{n} \sum \delta_{\mu_i}$)

\rightarrow gradient descent.

Otto calculus: (see Ambrosio,igli, Savaré '08)

$$f: \mathbb{R}^d \rightarrow \mathbb{R}$$

$$x \in \mathbb{R}^d$$

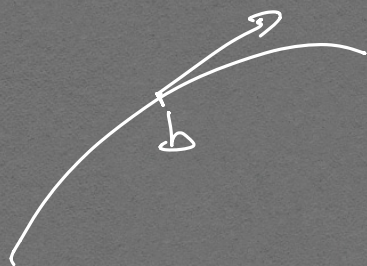
$$v \in T_x \mathbb{R}^d = \mathbb{R}^d$$

$$D_v f(x) = \lim_{h \rightarrow 0} \frac{f(x + hv) - f(x)}{h}$$

gradient $\nabla f(x) \in (\mathcal{T}_x \mathbb{R}^d)^\perp = \mathbb{R}^d$ is st.

$$\langle \nabla f(x), v \rangle = D_v f(x) \quad \forall v \in \mathcal{T}_x \mathbb{R}^d$$

Tangent space of \mathcal{W}_2 at b ?



$T_{b \rightarrow p}$ - id (displacement map)
 $\nabla \varphi, \varphi_{conv}$ is a tangent vector.

$$\mathcal{T}_b \mathcal{W}_2 = \left\{ \lambda (\nabla \varphi - id), \lambda > 0 \right\}_{\varphi_{conv}}^{L^2}$$

↑
 this is a subspace of $L^2(\mathbb{R}^d)$

$$\forall U, V \in \mathcal{T}_b \mathcal{W}_2, \langle V, U \rangle_b = \mathbb{E}_b \langle V(x), U(x) \rangle_{\mathbb{R}^d}$$

$$\exp_b : \mathcal{T}_b \mathcal{W}_2 \rightarrow \mathcal{W}_2$$

$$V \rightarrow (id + V) \# b$$

Gradient of the barycenter functional $b \mapsto F(b)$

$$V \in \mathcal{T}_b \mathcal{W}_2, \quad V = T_{b \rightarrow b'} - id$$

$$D_V F(b) = \lim_{h \rightarrow 0} \frac{F(id + h(T_{b \rightarrow b'} - id) \# b) - F(b)}{h}$$

$$\begin{aligned}
&= \lim_{h \rightarrow 0} \frac{1}{2h} \left\{ \underbrace{W_2^2(h(T_{b \rightarrow b'} - \text{id})_{\#} b, \mu) - W_2^2(b, \mu)}_{\text{opt. between } (b, \mu)} \right\} dP(\mu) \\
&\rightarrow \leq \int \|X + h(T_{b \rightarrow b'}(X) - X) - Y\|^2 d\mathcal{P}(X, Y) \\
&= \int \|X - Y\|^2 + 2h \langle T_{b \rightarrow b'}(X) - X, X - Y \rangle d\mathcal{P} + o(h^2) \\
&= W_2^2(X, Y) + 2h \langle T_{b \rightarrow b'} - \text{id}, \text{id} - T_{b \rightarrow \mu} \rangle_{L^2(b)}^{+d(h^2)}
\end{aligned}$$

$$\begin{aligned}
D_1 F(b) &\leq \left\langle T_{b \rightarrow b'} - \text{id}, \text{id} - \int T_{b \rightarrow \mu} dP(\mu) \right\rangle_{L^2(b)} \\
&= \underbrace{\langle T_{b \rightarrow b'} - \text{id}, \text{id} - \int T_{b \rightarrow \mu} dP(\mu) \rangle}_V
\end{aligned}$$

$$\boxed{\nabla F(b) = \text{id} - \int T_{b \rightarrow \mu} dP(\mu)}$$

Gradient descent:

$$\begin{aligned}
b_{t+1} &= [\text{id} - \eta_t \nabla F(b_t)]_{\#} b_t \\
&= \left[(1 - \eta_t) \text{id} + \eta_t \underbrace{\int T_{b_t \rightarrow \mu} dP(\mu)}_{\frac{1}{n} \sum_{i=1}^n T_{b_t \rightarrow \mu_i}} \right]_{\#} b_t
\end{aligned}$$



SGD: $p_1, \dots, p_n \stackrel{iid}{\sim} \mathcal{P}$

$$b_{t+1} = [(1 - \eta_t) \text{id} + \eta_t T_{b_t \rightarrow p_t}] \# b_t$$

Rates of convergence?

5- Optimization redux

$f: \mathbb{R}^d \rightarrow \mathbb{R}$ convex

$$\frac{\alpha}{2} \|x - y\|^2 \leq f(y) - f(x) - \langle \nabla f(x), y - x \rangle \leq \frac{\beta}{2} \|x - y\|^2$$

α -str. cvx

β -strongly smooth

GD: $x_{t+1} = x_t - \eta_t \nabla f(x_t)$

↙ minimizer

$$f(x_{t+1}) - f(x^*) = f(x_t - \eta_t \nabla f(x_t)) - f(x^*)$$

$$\leq f(x_t) - f(x^*) - \eta_t \|\nabla f(x_t)\|^2 + \frac{\eta_t^2 \beta}{2} \|\nabla f(x_t)\|^2$$

$$\boxed{\eta = \frac{1}{\beta}} = f(x_t) - f(x^*) - \frac{1}{2\beta} \|\nabla f(x_t)\|^2$$

ASSUME Polyak - Łojasiewicz (PL) inequality:

$$f(x) - f(x^*) \leq C_{PL} \|\nabla f(x)\|^2$$

$$\leq \left(1 - \frac{1}{2\beta C_{PL}}\right) (f(x_t) - f(x^*))$$

$$\leq \left(1 - \frac{1}{2\beta C_{PL}}\right)^{t+1} [f(x_0) - f(x^*)]$$

if α -strongly convex (at x^*), we get

$$\|x_t - x^*\|^2 \leq \frac{2}{\alpha} \left(1 - \frac{1}{2\beta C_{PL}}\right)^{t+1} [f(x_0) - f(x^*)]$$

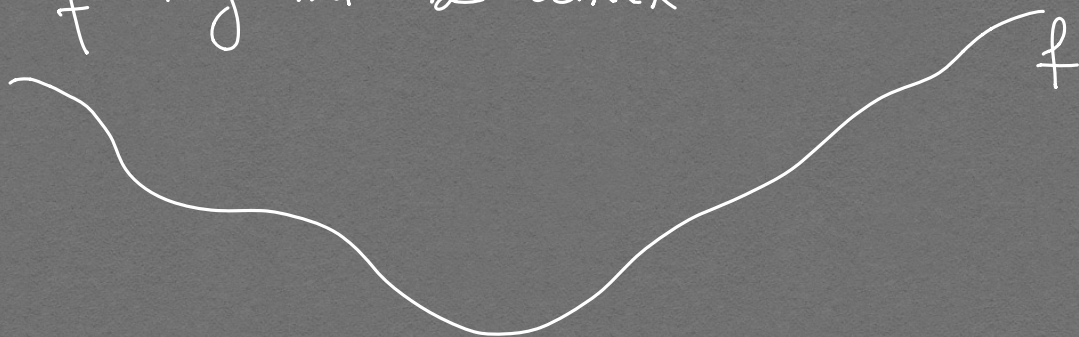
Recap: we used:

A. $f(y) - f(x) \leq \langle \nabla f(x), y-x \rangle + \frac{\beta}{2} \|y-x\|^2$

B. $\frac{\alpha}{2} \|x-x^*\|^2 \leq f(x) - f(x^*)$ at x^*

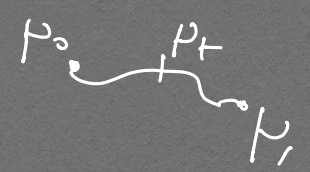
C. $f(x) - f(x^*) \leq C_{PL} \|\nabla f(x)\|^2$

↳ Note: f may not be convex



6. Back to Wasserstein

$G: \mathcal{W}_2 \rightarrow \mathbb{R}$ is geodesically convex if
 \forall geodesic $\mu_t \in \mathcal{W}_2$ from μ_0 to μ_1 , it holds

$$G(\mu_t) \leq (1-t)G(\mu_0) + tG(\mu_1)$$


(can also define smoothness, strong convexity).

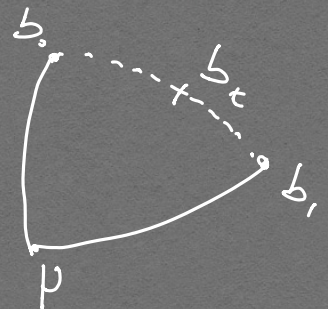
Example: if π is log-concave, then,

$b \mapsto \text{KL}(b \parallel \pi)$ is g -convex.

Fact: F is NOT geodesically convex

But: $A. F$ is 1 -smooth

$$\text{Curv}(\mathcal{W}_2) \geq 0$$



$$\int \frac{W_2^2(b_t, p) - W_2^2(b_0, p)}{2t} dP(p)$$

$$\geq \int (\cancel{1-t}) W_2^2(b_0, p) + \cancel{t} W_2^2(b_1, p) - \frac{\cancel{t}(1-t)}{2} W_2^2(b_0, b_1) dP(p)$$

$$= F(b_1) - F(b_0) - \frac{(1-t)}{2} W_2^2(b_0, b_1)$$

$t \rightarrow 0$:

$$\langle \nabla F(b_0), T_{b_0 \rightarrow b_1} - id \rangle_b \geq F(b_1) - F(b_0) - \frac{1}{2} W_2^2(b_0, b_1)$$

Rearranging:

$$F(b_1) - F(b_0) \leq \langle \nabla F(b_0), T_{b_0 \rightarrow b_1} - id \rangle_b + \frac{1}{2} W_2^2(b_0, b_1)$$

In fact 1-smoothness of d^2 is a characterization of $\text{Curv} \geq 0$.

$$B. \quad F(b) - F(b^*) \geq C W_2^2(b, b^*)$$

↳ Variance inequality.

Sturm ('03):

$$\text{Curv}(\mathcal{M}) \leq 0 \Leftrightarrow \int d^2(b, x) - d^2(b^*, x) dP(x) \geq 1 \cdot d^2(b, b^*)$$

$\forall P$

Th) CMRS

Assume that $T_{b^* \rightarrow p} = \nabla \varphi_p$ for some $\alpha(p)$ strict potential φ_p

\Rightarrow

$$F(b) - F(b^*) \geq C_{\text{var}} W_2^2(b, b^*)$$

with $C_{\text{var}} = \int \alpha(p) dP(p)$

Rk: strict improvement over LPRS (no smoothness)

C. Integrated PL inequality

Kantorovich duality yields:

$$\frac{1}{2} W_2^2(b, \mu) = \int \left(\frac{\|\cdot\|^2}{2} - \varphi_{\mu \rightarrow b} \right) d\mu + \int \left(\frac{\|\cdot\|^2}{2} - \varphi_{b \rightarrow \mu} \right) db$$

$$\forall b' \quad \frac{1}{2} W_2^2(b^*, \mu) \geq \int \left(\frac{\|\cdot\|^2}{2} - \varphi_{\mu \rightarrow b} \right) d\mu + \int \left(\frac{\|\cdot\|^2}{2} - \varphi_{b \rightarrow \mu} \right) db^*$$

$$\Rightarrow F(b) - F(b^*) \leq \underbrace{\int \left(\frac{\|\cdot\|^2}{2} - \int \varphi_{\mu \rightarrow b} dP(\mu) \right) d(b - b^*)}_f$$

take $X \sim b$, $Y \sim b^*$ optimally coupled

$$\mathbb{E}[f(X) - f(Y)] = \int_0^1 \mathbb{E}[\langle \nabla f(tY + (1-t)X), Y - X \rangle] dt$$

$$\leq \sqrt{\int_0^1 \mathbb{E} \|\nabla f(X_t)\|^2 dt} \sqrt{\mathbb{E} \|Y - X\|^2}$$

$$= W_2(b, b^*) \left(\int_0^1 \|\nabla f\|_{L^2(b_t)}^2 dt \right)^{\frac{1}{2}}$$



$$\leq \sqrt{\frac{F(b) - F(b^*)}{\text{Cvar}}} \left(\int_0^1 \|\nabla f\|_{L^2(b_t)}^2 dt \right)^{\frac{1}{2}}$$

$$\Rightarrow F(b) - F(b^*) \leq \int_0^1 \|\nabla f\|_{L^2(b_t)}^2 dt$$

to obtain a real PL, need to control

$$\|\nabla f\|_{L^2(b_t)}^2 \leq \|\nabla f\|_{L^2(b)}^2 \quad \forall t \in [0,1]$$

would be sufficient to have $\frac{db_t}{db} \leq c$

main issue: No control on iterates of GD.

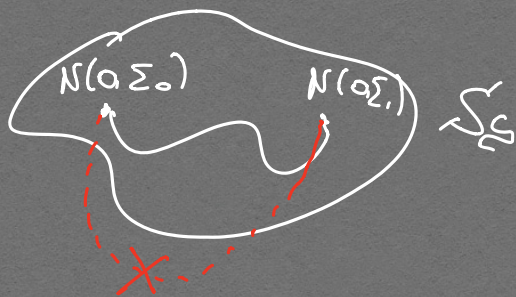
7. Gaussian distributions

$$\text{supp}(P) \subset \mathcal{S}_\xi$$

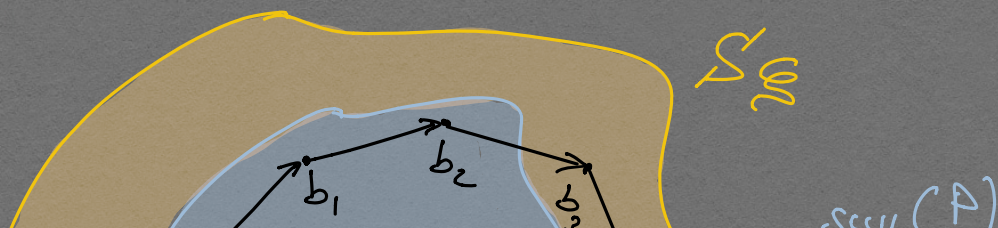
$$\mathcal{S}_\xi := \{ \mathcal{N}(0, \Sigma) : \|\Sigma\|_{\text{op}} \leq 1, \det \Sigma \geq 2\xi \}$$

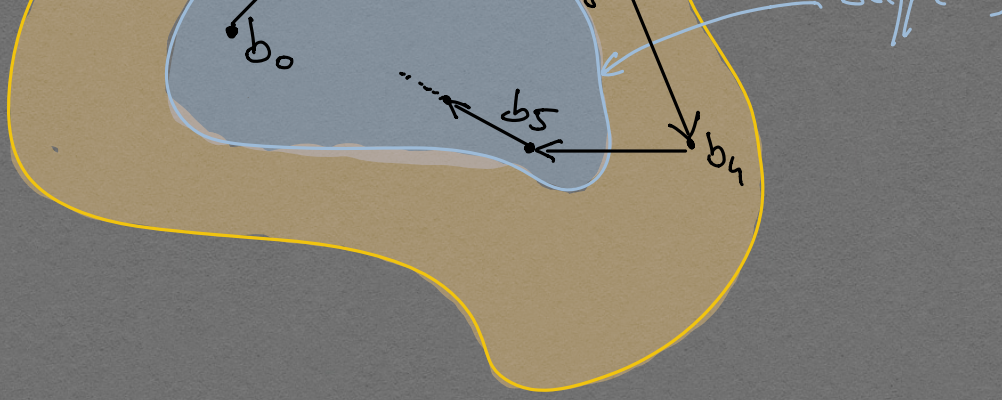
density upper bound.

Claim: \mathcal{S}_ξ is a geodesically convex set



Recall: GD & SGD push along geodesics





A. (1-smoothness) is always true

B. Variance inequality:

$$F(b) - F(b^*) \geq \frac{\sigma}{2} W_2^2(b, b^*)$$

C.

The CMRS 20

$\forall b \in \mathcal{S}_\xi$:

$$F(b) - F(b^*) \leq \frac{1}{2\xi^2} \|\nabla F(b)\|_{L^2(b)}^2$$

(PL ineq)

Concretely

Algorithm 1 Bures-Wasserstein GD

```

1: procedure BURES-GD( $\Sigma_0, Q, T$ )
2:   for  $t = 1, \dots, T$  do
3:      $S_t \leftarrow \int \Sigma_{t-1}^{-1/2} \{ \Sigma_{t-1}^{1/2} \Sigma(\mu) \Sigma_{t-1}^{1/2} \}^{1/2} \Sigma_{t-1}^{-1/2} dQ(\mu)$ 
4:      $\Sigma_t \leftarrow S_t \Sigma_{t-1} S_t$ 
5:   end for
6:   return  $\Sigma_T$ 
7: end procedure

```

Algorithm 2 Bures-Wasserstein SGD

```

1: procedure BURES-SGD( $\Sigma_0, (\eta_t)_{t=1}^T, (K_t)_{t=1}^T$ )
2:   for  $t = 1, \dots, T$  do
3:      $\hat{S}_t \leftarrow \Sigma_{t-1}^{-1/2} \{ \Sigma_{t-1}^{1/2} K_t \Sigma_{t-1}^{1/2} \}^{1/2} \Sigma_{t-1}^{-1/2}$ 
4:      $\Sigma_t \leftarrow ((1-\eta_t)I_D + \eta_t \hat{S}_t) \Sigma_{t-1} ((1-\eta_t)I_D + \eta_t \hat{S}_t)$ 
5:   end for
6:   return  $\Sigma_T$ 
7: end procedure

```

↑ Studied in Alvarez-Esteban, del Barrio & Cuesta-Albertos '16
as fixed point algorithm ($\nabla F(b^*) = 0$). No guarantees)

Main results

Run either GD on $F(b) = \frac{1}{n} \sum_{i=1}^n W_2^2(b, \mu_i)$
or SGD on μ_1, \dots, μ_n

Th1 GD with stepsize 1

$$W_2^2(b_T, \hat{b}_n) \lesssim \frac{(1 - \xi^2)^T}{\xi}$$

↓ + Koshnin et al.

$$\mathbb{E} W_2^2(b_T, b^*) \lesssim n^{-1/2}$$

for $T \geq C \log n$

Th2 SGD with $\eta_t \sim \frac{1}{t}$

$$\mathbb{E} W_2^2(b_n, b^*) \lesssim \frac{\text{Var}(P)}{n \xi^5}$$

Some open questions

1. Better dependence on ξ ?
2. Averaging
3. Other functionals?
4. Beyond Gaussians.