

# Course Description

Standard supervised classification setup:

- Data distribution:  $(x, y) \in \mathcal{X} \times \mathcal{Y} \subseteq \mathbb{R}^d \times \mathbb{R}$   
 $(x, y) \sim \nu$ ,  $\nu$ : probability measure in  $\mathcal{X} \times \mathcal{Y}$ .
- Loss function  $l(\hat{y}, y)$ ; convex with respect to  $\hat{y}$ .  
 e.g.  $l(\hat{y}, y) = \|\hat{y} - y\|^2$ ;  $l(\hat{y}, y) = \log(1 + e^{-\hat{y} \cdot y})$
- Model  $\hat{y} = \phi(x; \theta)$ ;  $\theta \in \mathcal{T}$
- Empirical Risk Minimization (ERM) / Structural Risk Min.

Empirical loss/Risk:  $(x_1, y_1) \dots (x_n, y_n) \stackrel{i.i.d.}{\sim} \nu$

$$\hat{\nu} = \frac{1}{n} \sum_{i=1}^n \delta_{(x_i, y_i)}$$

empirical measure

$$\hat{R}(\theta) = \mathbb{E}_{\hat{\nu}} [l(\phi(x, \theta), y)] + R(\theta)$$

regularization

$$= \frac{1}{n} \sum_{i=1}^n l(\phi(x_i, \theta), y_i) + R(\theta)$$

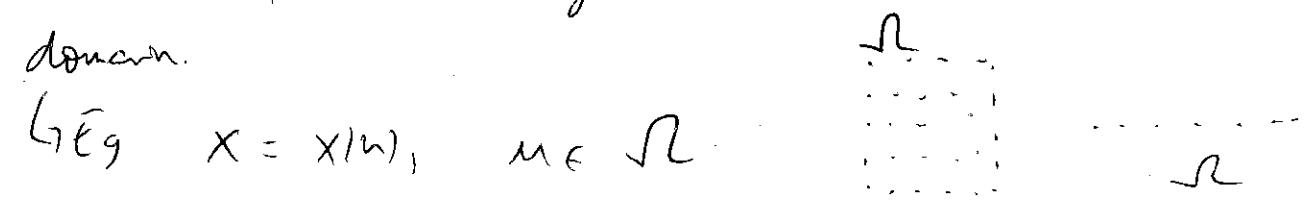
- Optimization:  $\min_{\theta} \hat{R}(\theta)$

Many questions:

- choice of model? (approximation) when is
- choice of optimization? deep learning
- generalization? (statistics) specific?

## Part I: Geometry of data

- input  $x$  in a high-dimensional space -  $x \in \mathbb{R}^d$  d.s.s.1
- But in many applications,  $x$  may be itself modeled as a function defined over a low-dimensional domain.



- Impact of this extra structure in  $\Omega$ ?
- In particular, assuming no noise in targets,  $y = f(x)$ , how to build representations  $\Phi(x; \Omega)$  that are stable to known transformations  $T$  that st  $|f(Tx) - f(x)| \approx 0$ ?
- Domain  $\Omega$  can be a group (lots of global symmetries) (torus,  $\mathbb{R}^d$ )
- a manifold
- a graph

- Q: general framework? Applications to inverse problems.
- Density estimation: choice of metric? [from MLE to OT].

## Part II: Geometry of Optimization and Learning

- Previous questions do not concern any learning/optimization.
- Basics of convex optimization and Nesterov acceleration.
- Non-convex optimization (or how to escape saddle points eff).
- Continuous-time analysis (stochastic and deterministic).

~~over-parametrized~~ region,  
role of batch normalization  
positive / negative results.

Asymptotic behavior of "over-parametrized" neural networks dynamics.

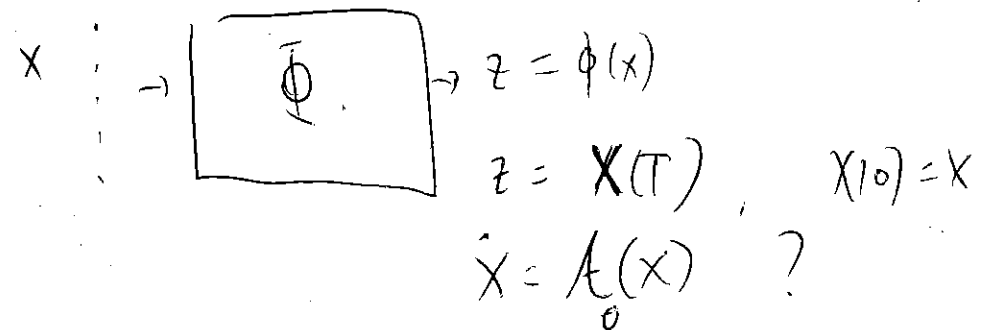
- Reproducing Kernel Hilbert Spaces (RKHS).
- Mean field and optimal transport.

Implicit bias of gradient descent.

Lecture 1: The curse of dimensionality.

Parallel curricula:

DFL 1: The Neural ODE [Chen et al. '18]



DFL 2: Exploration in RL.

- Unifying count-based Exploration and Intrinsic Motivation [Bellemare et al '16]

Lecture 1: The curse of dimensionality

↳ Generally refers to an exponential scaling of some cost / budget with respect to the input dimensionality of the problem.

• RKHS Hypothesis class  $\mathcal{F}$ : a function space  
 $\mathcal{F} = \{ f: \mathcal{X} \rightarrow \mathbb{R} \mid \mathcal{X} = \mathbb{R}^d \}$

observe  $n$  points  $(x_1, f(x_1)) \dots (x_n, f(x_n))$ .  
 $x_i \sim \text{iid}$ .

$f^*$  unknown, assumed in  $\mathcal{F}$ .

Q: How does the risk of our best estimator  $f$  grow with  $n, d$ ?

1st case:  $\mathcal{F}$  contains linear functions

$f(x) = \langle x, \theta \rangle + w, w \sim \mathcal{N}(0, \sigma^2)$

$f(x) = \langle \theta, x \rangle \quad L(\theta) = \frac{1}{2n} \sum_{i=1}^n (w_i + \langle x_i, \theta \rangle - \langle x_i, \theta^* \rangle)^2$

$\theta \sim \mathcal{N}(0, \sigma^2 (X^T X)^{-1}) \quad L(\theta) = \frac{1}{2n} \sum_{i=1}^n (w_i + \langle x_i, \theta^* - \theta \rangle)^2$

$\nabla L(\theta) = \frac{1}{n} \sum x_i (w_i + \langle x_i, \theta^* - \theta \rangle)$

$\hat{\theta} = \frac{1}{\sum x}^{-1} \frac{1}{\sum x y}$        $\sum_x = \frac{1}{n} \sum x_i x_i^T$

$\|\hat{\theta} - \theta^*\| \sim \sqrt{\frac{1}{n}}$

2nd case:  $\mathcal{F}$  contains Lipschitz functions.

$$\|f(x) - f(x')\| \leq \beta \|x - x'\|$$

$$\sup \| \hat{f}(x) - f^*(x) \| \leq \epsilon ?$$

As  $n \sim \epsilon^{-d}$

## Lecture 2: Geometric Stability in Euclidean Domain (Scattering Transform)

→ Linearization.



$f(x), x \in X \rightarrow$  high-dim space.

$\hat{f}(x) = \sigma(a^T \phi(x) + b)$  is a good approximation

If we think about  $f(x)$  as encoding a two-class classif, the change of variables  $x \mapsto \Phi(x)$  linearizes the boundary.

↳ in particular,  $a^T (\phi(x) - \phi(x')) = 0$

$$\Rightarrow f(x) = f(x')$$

(level sets of  $f$  mapped to hyperplanes by  $\Phi$ )  
(intra-class variability becomes flat)  
How to construct such  $\Phi$ ?

→ Invariance & Symmetry.  $x \in \Omega$  (input domain)

∃ a global symmetry  $\phi$ 's operator  $\varphi \in \text{Aut}(\Omega)$  st  
 $f(x) = f(\varphi(x)) \quad \forall x.$

∴ They can be absorbed by  $\phi$  as

$$\left\{ \begin{array}{l} \text{Invariants: } \phi(\varphi(x)) = \phi(x) \quad \forall x. \\ \text{Equivariants: } \phi(\varphi(x)) = \tilde{\varphi} \phi(x) \quad \forall x \end{array} \right.$$

Q: symmetries in image recognition problems?

Translations:  $x \in L_2(\mathbb{R}^2)$ ,  $\boxed{3} \rightarrow \boxed{3}$

$\{ \varphi_v ; v \in \mathbb{R}^2 \}$ ,  $\varphi_v(x)(u) = x(u-v)$

Dilations:  $\{ \varphi_s ; s \in \mathbb{R}_+ \}$ ,  $\varphi_s(x)(u) = s^{-1}x(s^{-1}u)$

Rotations  $\{ \varphi_\theta ; \theta \in [0, 2\pi) \}$ ,  $\varphi_\theta(x)(u) = x(R_\theta u)$

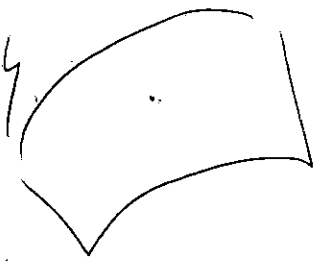
These are subgroups of the group of rigid motions Affine Group  $\text{Aff}(\mathbb{R}^2)$ , we can associate it with the space of affine transforms

$$\begin{pmatrix} u_1 \\ u_2 \end{pmatrix} \mapsto \begin{pmatrix} v_1 \\ v_2 \end{pmatrix} + \begin{pmatrix} a_1 & a_2 \\ a_3 & a_4 \end{pmatrix} \begin{pmatrix} u_1 \\ u_2 \end{pmatrix}$$

- ↳ Not-commutative in general.
- ↳ In analysis, this group contains the Heisenberg group, with linear part of the form  $\begin{pmatrix} 1 & a \\ 0 & 1 \end{pmatrix}$
- ↳ Lie group; ~~symplectic manifold~~ (manifolds with symplectic)

Invariant Representations:

The action of a transformation group  $G$  on  $x$  is an orbit:  $G \cdot x = \{ \varphi_g(x) ; g \in G \}$



Recall  $\text{Aff}(\mathbb{R}^2)$  like  $\text{U}(1)$   
 P. General approach to obtain Group Invariant Represent?

One-parameter unitary groups

Def. A one-parameter unitary group  $\{ \varphi_t \in \text{Aut}(\mathcal{H}) \}_{t \in \mathbb{R}}$

(i)  $\forall t, s \quad \varphi_{s+t} = \varphi_s \varphi_t$

(ii)  $\lim_{s \rightarrow t} \| \varphi_s - \varphi_t \| = 0$ , (iii)  $\| \varphi_t x \| = \| x \| \quad \forall t, x$

In particular, these are Abelian groups (why?)  
 (e.g. Rotation / Translation)

Theorem [Stone, 30s]  $\mathcal{H}$  Hilbert space

$\left\{ \begin{array}{l} \text{self-adjoint operators on } \mathcal{H} \\ \mathcal{H} \end{array} \right\} \xleftrightarrow{1:1} \left\{ \begin{array}{l} \text{one-parameter} \\ \text{unitary groups of} \\ \text{Aut}(\mathcal{H}) \end{array} \right\}$

In particular; given  $\{ \varphi_t \}_{t \in \mathbb{R}}$ ,  $\exists$  A self-adjoint  $A$  st  $\forall t, \varphi_t = e^{itA}$   
 $(Au, w) = (u, Aw)$

Recall the Fourier Transform: for  $x \in L^2(\mathbb{R})$  we define

$$\hat{x}(\beta) = \int x(u) e^{-i\beta u} du$$

- Main properties:
- [Linearity]:  $z = \alpha x + \beta y \Rightarrow \hat{z} = \alpha \hat{x} + \beta \hat{y}$
  - [Parseval]:  $\| \hat{x} \| = \| x \|$ ,  $\langle x, y \rangle = \langle \hat{x}, \hat{y} \rangle$
  - [IFT]:  $x(u) = \int \hat{x}(\beta) e^{i\beta u} d\beta$
  - [Transl]:  $\hat{y} = \varphi_v x(u)$ ;  $\hat{y}(\beta) = e^{i\beta v} \hat{x}(\beta)$
  - [Dil]:  $y = x(su)$ ;  $\hat{y}(\beta) = s^{-1} \hat{x}(s^{-1}\beta)$

Stone thm and Fourier: (Translations are diagonal operators in the Fourier domain. Stone thm formalizes)  
 Fourier simultaneously diagonalises all translations. (possible because they commute)  
 Stone thm generalises this notion. "Fourier transform" of  $A = V^* \text{diag}(\lambda_1, \dots, \lambda_n) V$

→ larger Abelian groups also ok

$$(G = G_1 \times \dots \times G_p)$$

one-parameter groups

→ Generalized to non-abelian case by Peter-Weyl theorem using representation theory.

→ Q: How to obtain invariants in that case?

$A = V^* \text{diag}(\lambda_1 - \lambda_d) V$   
 $y = Vx$   
 $y_{1k}' = e^{i\lambda_{1k}t} y_{1k}$   
 $y_{1k} = \langle x, v_{1k} \rangle$

•  $\phi(x) = |Vx|$  is thus  $G$ -invariant:

$$\forall x, t \quad \phi(\varphi_t(x)) = \phi(x)$$

$$A = V^* \Lambda V \Rightarrow e^{itA} = V^* \text{diag}(e^{it\lambda_1}, \dots, e^{it\lambda_d}) V$$

$$V\varphi_t x = V e^{itA} x = V V^* \text{diag}(e^{it\lambda_j}) V x$$

$$\Rightarrow |V\varphi_t x| = |Vx| \quad \forall x, t$$

• Thus, in the commutative case, a single layer is sufficient to obtain invariance.

• But, even in the commutative case, is this good enough?

↳ Symmetry is a very strict criterion. How to relax?

↳ Symmetry groups in images are low-dimensional; Not that helpful!

→ Deformations:  $x \in L^1(\mathbb{R}^m)$   $\varphi: \mathbb{R}^m \rightarrow \mathbb{R}^m$  diffeomorphism

$$x_\tau(u) := x|_{u-\tau(u)} \quad \varphi(u) = u - \tau(u)$$

$(\|\nabla\tau\|_\infty < 1)$

→ Geometric Stability:  $\forall x, \tau$ , deformation cost.

$$|f(x) - f(x_\tau)| \leq \|\tau\|$$

E.g.  $\|\tau\| = \sup_u \|\nabla\tau(u)\|$ : measures the stretching of the domain.

Stable representations:  $\forall x, \tau$ ,  $\|\phi(x) - \phi(x_\tau)\| \leq C\|\tau\|$

since  $\hat{f}(x) = \langle \phi(x), \beta \rangle$ , it follows that

$$|\hat{f}(x) - \hat{f}(x_\tau)| \leq \|\beta\| \|\phi(x) - \phi(x_\tau)\| \leq C\|\beta\|\|\tau\|$$

→ Are the Fourier Invariants stable to deformations?

↳ We know  $\phi(\varphi_t x) = \phi(x) \quad \forall t, x$ , but what happens for transformations  $\varphi \notin G$ ?

→ In some applications, we are interested in local invariance instead.

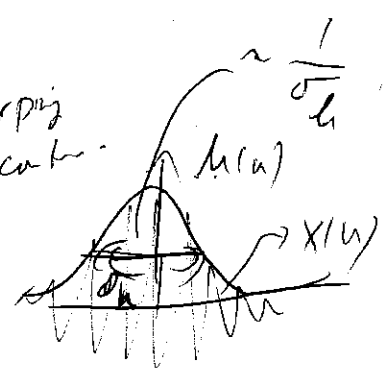
→ We consider the local deformation cost

$$\|\tau\| = \underbrace{2^{-\beta}}_{\text{displacement} \leq 2^\beta} \|\tau\|_\infty + \underbrace{\|\nabla\tau\|_\infty}_{\text{warping index cost}}$$

→ Shallow (Fourier) invariants are unstable.

$h(u)$  low pass window  $x(u) = h(u) e^{i\beta u}$

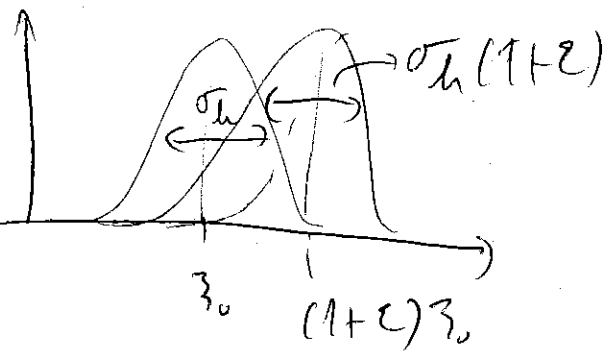
→ Deformation of the form  $\varphi_\epsilon x(u) = x((1+\epsilon)u)$ ,  $\epsilon \ll 1$



$$(\hat{\varphi}_\tau x)(z) = (1+\varepsilon)^{-1} \hat{x}((1+\varepsilon)^{-1}z)$$

$$\hat{x}(z) = \hat{h}(z-z_0)$$

$$(\hat{\varphi}_\tau x)(z) = (1+\varepsilon)^{-1} \hat{h}((1+\varepsilon)^{-1}(z-z_0))$$



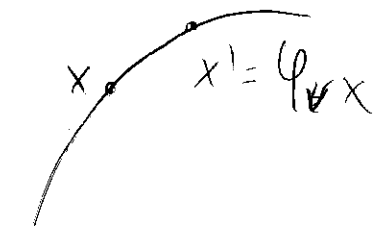
If  $(1+\varepsilon)z_0 - z_0 = \varepsilon z_0 \gg \sigma_h(2+\varepsilon)$

then supports are disjoint,

$$\| |\hat{x}| - |\hat{x}_2| \| \sim \|x\|$$

→ Moreover,  $|\hat{x}|$  loses information in phase correlations.

→ Fix? Local Invariants.



$$\| \phi(x) - \phi(\varphi_v x) \| \leq C 2^{-d} \|v\|$$

$$\forall v, \|x\|=1, \frac{\| \phi(x) - \phi(\varphi_v x) \|}{\|v\|} \approx 2^{-d}$$

↳ smooth along the orbits  $\rightarrow$  Haar measure on  $G$

$$\Phi(x) = 2^{-d} \int_v \phi(2^{-d}v) \varphi_v x \, dv, \quad \left( \int \phi(v) \, dv = 1 \right)$$

$$= \int \phi_2(v) x(u-v) \, dv = x \cdot \phi_2(u), \quad \phi_2 = 2^{-d} \phi(2^{-d}v)$$

and in fact we have smoothness beyond the orbit, on general deformations:

Prop:  $\phi(x) = x \cdot \phi_2$  satisfies  $\forall \|x\|=1 \in L^2, \forall \tau$

$$\| \phi(x) - \phi(\varphi_\tau x) \| \leq C \|\tau\|$$

↳ Proof: Schur lemma for integral operators

$$|K f(u)| = \int |k(u, u')| f(u') \, du'$$

$$\|K\| \leq \max \left( \sup_u \int |k(u, u')| \, du', \sup_{u'} \int |k(u, u')| \, du \right)$$

$$\phi(x) = x \cdot \phi_2; \quad \phi(x)(u) = \int x(v) \phi_2(u-v) \, dv$$

$$\varphi_\tau(x)(u) = x(u-\tau(u))$$

$$\phi(x_\tau)(u) = \int x(v-\tau(v)) \phi_2(u-v) \, dv = \int x(v) \phi_2(u-v+\tau(v)) \cdot |I-\tau'(v)| \, dv$$

$$\Rightarrow (\phi(x) - \phi(\varphi_\tau x))(u) = \int x(v) [\phi_2(u-v) - \phi_2(u-v+\tau(v)) \cdot |I-\tau'(v)|] \, dv$$

Q: other stable linear operators?

$$\forall v, \phi(x) = \phi(\varphi_v x) \Rightarrow \phi(x) = \frac{1}{\mu(G)} \int \phi(\varphi_v x) \, dv$$

$$\Rightarrow \phi(x) = \phi \left( \frac{1}{\mu(G)} \int \varphi_v x \, dv \right) = \phi \left( \frac{1}{\mu(G)} \int x(u) \, dv \right)$$

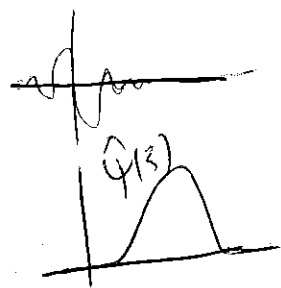
# Lecture 3: The Scattering Transform

→ We just saw that the only linear invariant is essentially the average. This loses all the relevant information.

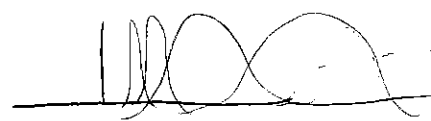
Corollary: High frequency representation must involve a non-linearity. → recall we want (i) stability (ii) "richness" (ability to detect)

Wavelets  $\psi$  bandpass  $L^1(\mathbb{R}^d)$   $\int \psi(u) du = 0$

$\hat{\psi}(0) = 0$  We may ask for more vanishing moments, equivalently  $\int u^j \hat{\psi}(u) du = 0 \forall j \leq m$ . (a single vanishing moment)



In order to extract all the frequencies, we need



wavelets at different scales:  $\psi_j(u) = 2^{-j} \psi(2^{-j}u)$

$$\psi_{j,0}(u) = 2^{-j} \psi(2^{-j}u) \quad (1D) \quad j \in \mathbb{Z}$$

$$\psi_{j,10}(u) = 2^{-j} \psi(2^{-j}u) \quad (2D)$$

Littlewood-Paley transform: let  $x \in L^2(\Omega)$ , then

$$Wx = \left\{ \begin{array}{l} x * \psi_{j,0} \quad (j \in \mathbb{Z}) \\ x * \phi_j \quad k \in \mathbb{Z} \end{array} \right\} \in L^2(\Omega \times \mathbb{Z} \times \mathbb{Z})$$

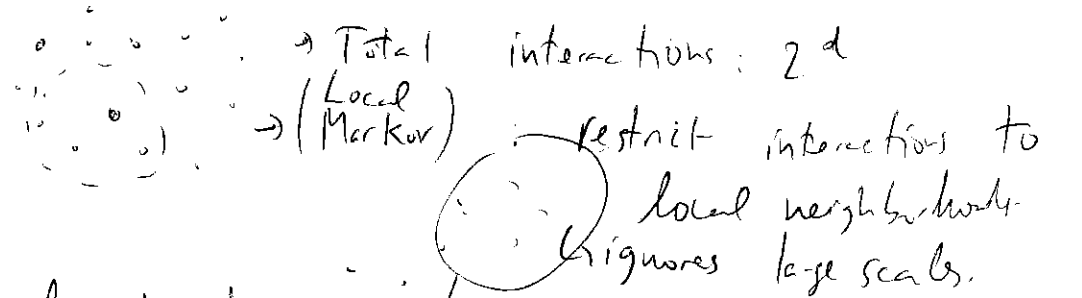
Thm (Littlewood-Paley): if  $\exists \delta > 0$  such that

$$\forall \epsilon, \quad 1 - \delta \leq |\hat{\phi}(\xi)|^2 + \frac{1}{\epsilon} \sum_k |\hat{\psi}(2^k \xi)|^2 \leq 1,$$

then  $\forall x \in L^2, \quad (1 - \delta) \|x\|_{L^2(\Omega)}^2 \leq \|Wx\|_{L^2(\Omega \times \mathbb{Z})}^2 \leq \|x\|_{L^2(\Omega)}^2$

- Link this to additive stability.
- wavelets used extensively in harmonic analysis to characterize local regularity (eg through Besov-spaces).
- used successfully in non-parametric statistics to obtain optimal rates in denoising via the shrinkage estimator (Donoho).

→ wavelets and scale separation: suppose  $d$  particles interacting.



→ Multiscale interaction. each particle: from  $d$  to  $\log d$  interactions.

↳ Related to Fast Multipole Method (Greengard & Rokhlin)

→  $\psi$ : Can we then simply concatenate  $W$  with an averaging? No, since they have 0 mean. We can create a new invariant out of each wavelet sub-band with a non-linearity that restores an informative average.

For example using a complex modulus:  $\int |x * \psi_k| du$

$Wx$  is translation equivalent:  $\tilde{x}(u) = x(u - u_0)$   
 $W_k \tilde{x}(u) = W_k x(u - u_0)$

$\psi$ : stable to deformations?

As before, let  $\psi_{\tau} x(u) := x(u - \tau(u)) \forall x$ .

Theorem [Mallat '12]: There exists  $C > 0$  such that  $\forall j \in \mathbb{Z}$

$$\| [W_\tau, \Psi_\tau] \| \leq C \left[ \|\nabla \tau\|_\infty \max(1, \log \frac{\|\Delta \tau\|_\infty}{\|\nabla \tau\|_\infty}) + \|\Delta \tau\|_\infty \right]$$

↳ Remark: For the compass, we had local stability/invariance  
 $\| \phi - \phi \circ \Psi_\tau \| \leq C \cdot \|\tau\|$ , here it is equivalence.

↳ Key tool behind this result: Cotlar-Stein near orthogonal  
Lemma: A family of bounded linear operators between two Hilbert spaces  $E, F$ ;  $T_j: E \rightarrow F$ .

$\{T_j\}$  is almost orthogonal if

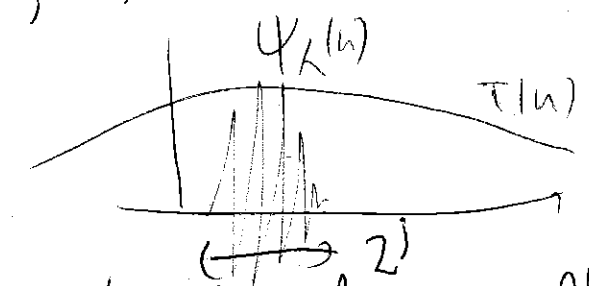
$$A = \sup_j \sum_k \sqrt{a_{jk}} < +\infty, \quad B = \sup_k \sum_j \sqrt{b_{jk}} < +\infty$$

$$a_{jk} := \| T_j T_k^* \|_E, \quad b_{jk} := \| T_j^* T_k \|_F$$

If  $\{T_j\}$  are almost orthogonal, then  $\sum_j T_j$  converges,

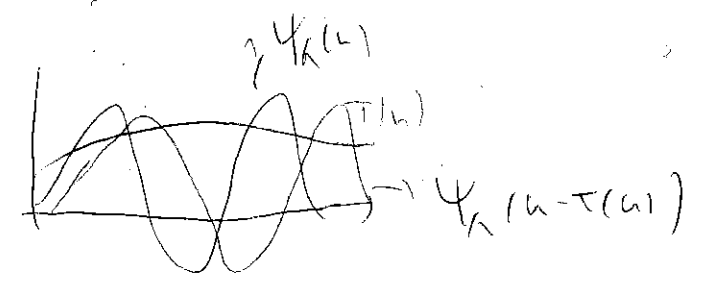
$$\text{and } \| \sum_j T_j \| \leq \sqrt{AB}$$

Intuition:



(\*) at small scales,  $\Psi_\lambda$  has small support, and for  $u, v$  within support of  $\Psi_\lambda$ , and because  $\tau$  is smooth,  
 $|\tau(u) - \tau(v)| \sim 2^j \|\nabla \tau\|_\infty$ .

$$\text{Thus } |(\Psi_\tau x) \circ \Psi_\lambda(u) - (x \circ \Psi_\lambda)(u - \tau(u))| \sim \|\nabla \tau\|_\infty$$

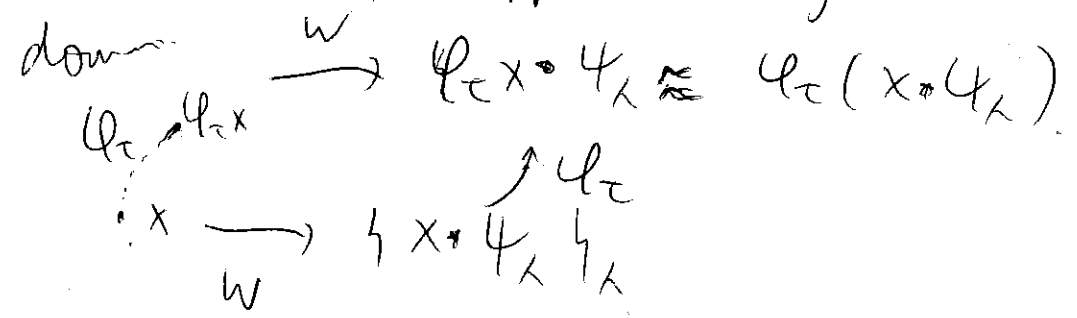


(i) at large scales,  $\Psi_\lambda$  is itself smooth, thus  
 $|(\Psi_\tau(x \circ \Psi_\lambda)) - (\Psi_\tau x) \circ \Psi_\lambda| \sim \|\nabla \tau\|_\infty$ .

(similar as our result on local stability for  $\phi$ )

(iii) Finally, different scales interact weakly in the sense of the Cotlar-Stein Lemma.

↳ This result says that deformations in the input domain are approximately mapped to deformations in the wavelet domain.



so we can hope to extract stable measurements as we did in the input with the average. But we need a non-linear operator in between (why?).

→ characterizing stable non-linearity.  $M: C^2(\Omega) \rightarrow C^2(\Omega)$ .  
 • Preserve additive stability from Littlewood-Paley.

$$\| Mx - Mx' \| \leq \| x - x' \|$$

• Preserve geometric stability from wavelet commutator: it is sufficient



Theorem (B'12): If  $M$  is non-expansive in  $L^2$  such that  $M\varphi_\varepsilon = \varphi_\varepsilon M$  for  $\varphi_\varepsilon$  diffeo, then  $M$  is point-wise:

$$Mx(w) = f(x(w)) \quad M: L^2 \rightarrow L^2 \text{ but } f: \mathbb{R} \rightarrow \mathbb{R}$$

↳ Idea of the proof: Look at the isotropy group of  $x \in L^2(\Omega)$   $G(x) = \{ \phi \in \text{Diff}(\Omega); x = x \circ \phi \}$  (symmetries of the level sets).

observation:  $G(x) \subseteq G(Mx) \quad \forall x$ .

For instance we choose the complex modulus (Reals is ok for real wavelets too).

$$x \rightarrow \{ |x \circ \psi_k| \psi_k \} \rightarrow \{ \| |x \circ \psi_k| \circ \psi_k^{-1} \|_{k, k'} \} \rightarrow \dots$$

$$\downarrow \qquad \qquad \qquad \downarrow$$

$$\int |x(w)| \psi(w) dw \quad \int |x \circ \psi_k| \psi(w) dw \quad \int \| |x \circ \psi_k| \circ \psi_k^{-1} \|_{k, k'} dw$$

(with  $dw$ )

Even time we extract an invariant, we loose information

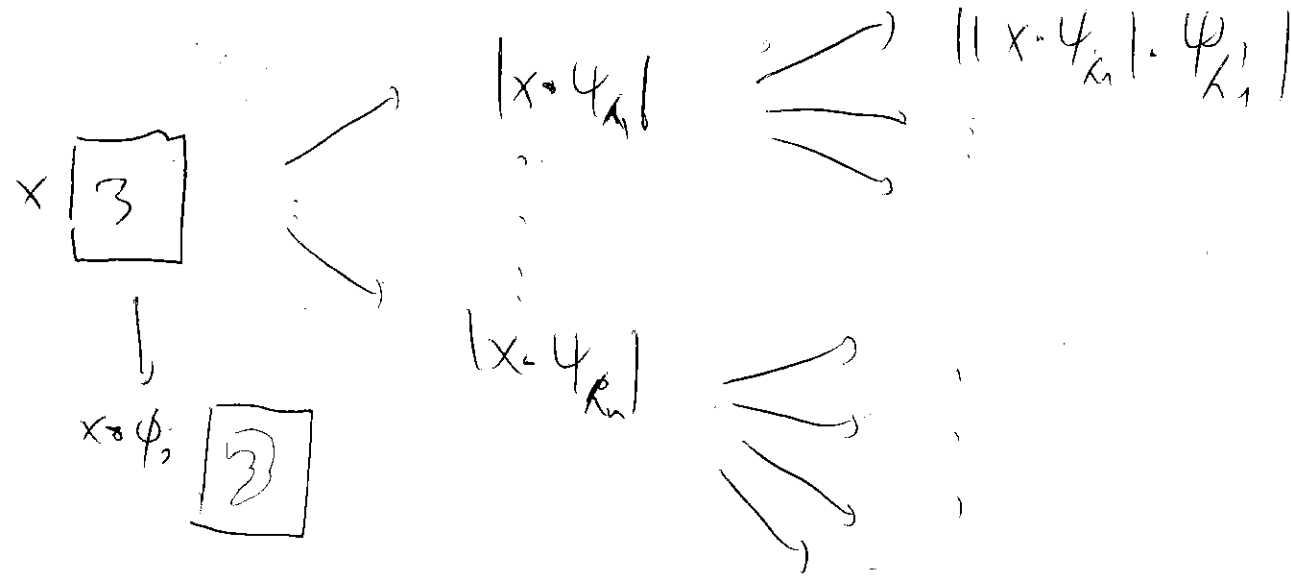
Q: How to systematically recover the info that is lost?

A new wavelet decomposition captures lost high-frequencies!

$$\int \| |x \circ \psi_k| \circ \psi_k^{-1} \|_{k, k'} (w) dw$$

Def: For each "path"  $p = (k_1, \dots, k_m)$  with  $k_i \in \Lambda_3$  and  $x$ ,

the scattering transform of  $x$  is



→ Scattering metric and Energy conservation.

For a given set of paths  $\Gamma$ , the Euclidean norm defined by scattering coefficients is

$$\| S_\Gamma(x) \|^2 := \sum_{p \in \Gamma} \| S_p(x) \|^2 \quad x \in L^2(\Omega)$$

is it well-defined? / let  $\Lambda_3^m$  denote all paths of length  $m$ , and denote by

$$\bar{U}_3 x = \{ x \circ \psi_3, |x \circ \psi_k| \psi_k \} \quad \text{the "propagator" operator}$$

→ It is non-expansive thanks to the Littlewood-Paley property:

$$\| \bar{U}_3 x - \bar{U}_3 x' \|^2 = \| x \circ \psi_3 - x' \circ \psi_3 \|^2 + \sum_k \| |x \circ \psi_k| - |x' \circ \psi_k| \|^2 \leq \| x - x' \|^2$$

(\*)

But observe that scattering coeffs are constructed by cascading

$$\bar{U}_3 : \text{if } U(p)x := \| \dots |x \circ \psi_{k_1}| \circ \dots \psi_{k_m} \| \quad p = (k_1, \dots, k_m)$$

$$U(p) \bar{U}_3 x \leq \dots \leq \| x \circ \psi_3 \|^2 + \sum_k \| |x \circ \psi_k| \|^2$$

$\mathcal{P}_j$ : set of all paths.  $\mathcal{P}_j = \bigcup_m \Lambda_j^m$

$\Lambda_j^m$ : paths of length  $m$ .

(Induction over  $m$ .  $\|S_j x - S_j x'\| \leq \|x - x'\|$ )

Apply (\*) to each  $|x \cdot \psi_k|, |x' \cdot \psi_k|$ .

### Lecture 4: Scattering Extensions, Graph Neural Networks.

↳ Consider now unitary wavelet decompositions:

$$(*) \quad \|x\|^2 = \|x \cdot \phi_j\|^2 + \sum \|x \cdot \psi_k\|^2$$

Thm: For appropriate wavelets,  $\|x\|^2 = \|S_j x\|^2$

By applying (\*) on each output  $x_p := |x \cdot \psi_{k_1}| \dots |x \cdot \psi_{k_m}|$

we have

$$\|x\|^2 = \sum_{|p|=m} \|S_j^m x\|^2 + \sum_{|p|=m} \| |x \cdot \psi_{k_1}| \dots |x \cdot \psi_{k_m}| \|^2$$

The result amounts to showing that

$$\lim_{m \rightarrow \infty} \sum_{|p|=m} \| |x \cdot \psi_{k_1}| \dots |x \cdot \psi_{k_m}| \|^2 = 0$$

• Energy moves progressively towards low frequencies

• Decay is Exponential for band-limited signals

[Ewaldspurger '17, Czaja '16]

(Ewaldspurger '17, Czaja '16)

### Geometric Stability of Euclidean Scattering

recall  $\|x\| := \sqrt{\|Ux\|_0^2 + \|Vx\|_0^2 + \|Hx\|_0^2}$  (max(L, D, H))

$$\text{recall } \|\tau\| := 2^{-j} \|\tau\|_0 + \|\nabla \tau\|_0 + \|H \tau\|_0$$

Thm. (Mallet '10) There exists  $C > 0$  st for all  $x \in C^2(\mathbb{R}^d)$

and all  $m$ , the  $m$ -th order scattering representation

$$\text{satisfies } \|S_j x - S_j x_\tau\| \leq C \cdot m \|x\| \|\tau\|$$

Proof architecture:

• Denote  $Ax = x \cdot \phi_j$ ,  $Wx = \{x \cdot \psi_k\}$ ,  $Mx = |x|$ .

• We know that

$$\|A - A\psi_\tau\| \leq C \|\tau\|$$

$$\|W\psi_\tau - \psi_\tau W\| \leq C \|\tau\|$$

$$M\psi_\tau = \psi_\tau M \quad \left( \begin{array}{l} \| [W, \psi_\tau] \| \leq C \|\tau\| \\ [M, \psi_\tau] = 0 \end{array} \right)$$

• By def, the operator  $S_j$  is

$$\{ A, A M W, A M W M W, \dots \} \quad (m \text{ times})$$

$$\|S_j - S_j \psi_\tau\|^2 = \|A - A\psi_\tau\|^2 + \|A M W - A M W \psi_\tau\|^2 + \dots + \|A(MW)^{\otimes m} - A(MW)^m \psi_\tau\|^2$$

$$\|AU^k - AU^k \psi_\tau\| \leq \|AU^k - AU^{k-1} \psi_\tau U\| +$$

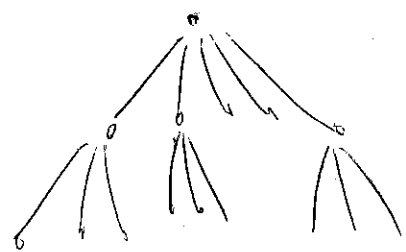
$$\|AU^{k-1} \psi_\tau U - AU^k \psi_\tau\|$$

$$\leq \|AU^{(k)} - AU^{(k-1)}\varphi_\varepsilon\| + \|AU^{(k-1)}[U, \varphi_\varepsilon]\|$$

$$\leq \|AU^{(k-1)} - AU^{(k-1)}\varphi_\varepsilon\| + C\|\varepsilon\|$$

$$\leq K C \|\varepsilon\| + \|A - A\varphi_\varepsilon\| \leq \tilde{C} K \|\varepsilon\| \quad \square$$

### Limitations of separable scattering



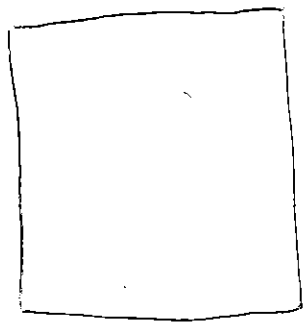
→ This is the specific convolutional architecture given by  $Sx = \{ \mathcal{S}_p x; |p| \leq m \}$ .

- ↳ No feature dimensionality reduction: each new layer increases the number of feature maps.
- ↳ Each wavelet band is assumed to be deformed independently.
- ↳ No adaptivity.

### Joint vs Separable Invariance

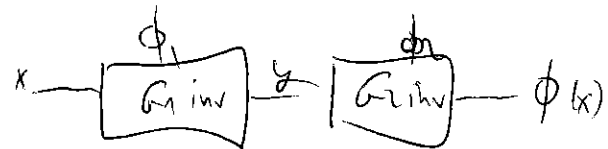
$$G = S^1 \times S^1$$

Each copy of  $S^1$  acts on a different coordinate (horiz vs vertical).



$$(p^1, p^2, u) = x(u_1 - a, u_2) \cdot \varphi^{\sim} x(u_1, u_2) = x(u_1, u_2) U_2 \text{ of}$$

$$G = G_1 \times G_2$$



$y = \phi_1(x)$  such that  $\phi_1(\varphi^1 x) = \phi_1(x) \quad \forall \varphi^1 \in G_1$

$\phi_2(y)$  such that  $\phi_2(\varphi^2 x) = \phi_2(x) \quad \forall \varphi^2 \in G_2$

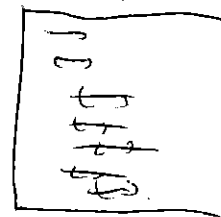
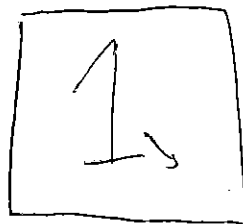
so if  $\phi_1$  is  $\left\{ \begin{array}{l} G_1 \text{ invariant} \\ G_2 \text{ equivariant} \end{array} \right\}$  then  $\phi = \phi_2 \circ \phi_1$  is  $G$ -invariant.

$\phi_2$  is  $G_2$  invariant

$$\begin{aligned} \phi_2 \phi_1(\varphi^1 \varphi^2 x) &= \phi_2 \phi_1(\varphi^2 x) = \phi_2 \varphi^2 \phi_1(x) = \\ &= \phi_2 \phi_1(x) = \phi(x). \end{aligned}$$

So we can compose invariances by using separate group structures.

However, this factorization does not take into account the joint action of  $G$  along  $(u_1, u_2)$ .



Ex: Pseudo-translation group.  $G = \mathbb{R}^2 \times S^1$   
( $u, \alpha$ )

$$g \in G, g = (u, \alpha)$$

$$(g \cdot x)(u) = x(R_\alpha T_u u)$$

Non-commutative:

$$g \cdot g' = R_\alpha \cdot T_\nu R_\alpha T_\nu \neq R_\alpha T_\nu R_\alpha T_\nu$$

This group acts on 1st-layer features similarly as translation acts on input images:

$$f(g \cdot x * \psi_k) = X_1(R_\alpha T_\nu u; \beta) \alpha + \theta$$

"  
(via)

separable spatial convolution

↳ so we can replace convolution with roto-transl. convolution:

def Let  $G$  be a compact group equipped with the Haar measure  $\mu$ , acting on  $\Omega$ , and  $h \in C(G)$

$$x *_{G} h(u) = \int_G x(\psi_g u) \cdot h(g) d\mu(g)$$

↳ If  $x = x(u, \beta, \theta)$  and  $G$  are roto-translations, then convolution recombines the channels (not separable across orientation).

↳ Scattering with roto-translation wavelets:  
[Sifre & Mallat '13] [Oyallon & Mallat '15]

Proposals: Give feedback this week.  
Inverse Curriculum of Aaron leave of absence.

### → Applications of scattering.

- In computer vision: [B & Mallat '11] MNIST/Texture. [Sifre & Mallat '13] Texture. [Oyallon & Mallat '15] CIFAR-10/100.
- Speech / Recognition: [B. et al '16] Super-Resolution. Music / Audio Processing: [Anden & Mallat '14] [Anden & Lostanien '16].
- Multifractal Analysis: [B et al '14, '15].
- Unsupervised Quantum Chemistry: [Eickenberg & Girn, '16].
- Unsupervised Learning [Angulo & Mallat '18] GAN [B & Mallat '18] Geometric Probability Theory.

### Lecture 5:

#### From Scattering to CNNs.

The essential stability properties of scattering stem from the commutation properties  $\| [W, L_\tau] \| \leq \| \tau \|$  and  $\| A - A_\tau \| \leq \| \tau \|$ .

$W = \{ \psi_k \}$  contains a family of differential operators  $D_k X = X * \psi_k$  which nearly commute with deformation.

A CNN can be written as  $X_{k'}^{(l+1)} = f \left[ \sum_k D_k X^{(l)} \cdot \theta_{k,k'} \right]$  We learn linear combinations of stable differential

From Euclidean to non-Euclidean stability

↳ so far, we measured geometric stability in terms of diffeomorphisms  $\phi: \Omega \rightarrow \Omega$ , in the sense that  $\phi(u)$  "close" to a rigid translation should satisfy

$$\|\phi(x) - \phi(\phi x)\| \approx 0.$$

↳ Q: If  $\Omega$  is not Euclidean / not even continuous? we can think of a diffeo as a change of metric.

$$(\Omega, d_\mu) \xrightarrow{\phi} (\Omega, \phi(d_\mu))$$

$$\langle x, x' \rangle = \int x(u) x'(u) \mu(du)$$

$$\langle x, x' \rangle_\phi = \int x(\phi(u)) \cdot x'(\phi(u)) \mu(du) =$$

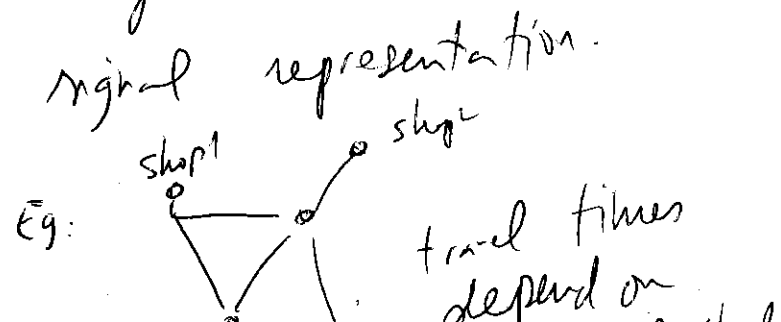
$$= \int x(u) \cdot x'(u) |\mathbf{I} - \nabla \phi(u)| \mu(du)$$

↳  $d_\mu(u) = |\mathbf{I} - \nabla \phi(u)| d_\mu(u)$  is a change of metric.

If  $\mathcal{F} = L^2(X)$  is our data input space: square-integrable functions defined over a metric space  $X$ .

$$\Phi: L^2(X) \rightarrow \mathbb{R}^K \quad \text{signal representation.}$$

$$X = (\Omega, d) \quad \text{distance.}$$



$$X^{(t=0)}, X^{(t=1)}$$

$\text{dist}(X^{(0)}, X^{(1)})$  measures how close the metrics are between  $X^{(0)}$

Stability to metric changes:

$$\left\| \frac{\Phi}{X} - \frac{\Phi}{X'} \right\| \lesssim \text{dist}(X, X')$$

Q: How to define meaningful distances? Graphs are flexible data structures to describe discrete metric domains.

→ Diffusion Wavelets on Graphs.

$$G = (V, E) \quad V \text{ set of nodes}$$

$E$ : set of edges  $(i,j), (i,j) \in V$ .

we focus mostly on undirected <sup>weighted</sup> graphs:

Let  $|V|=n$ , and  $W \in \mathbb{R}^{n \times n}$  be a symmetric weighted adjacency matrix.  $D = \text{diag}(W\mathbf{1})$  degree matrix.

↳ we define the <sup>symmetric</sup> diffusion operator  $A := D^{-1/2} W D^{-1/2}$

$$\tilde{A} = D^{-1} W \text{ is a Markov Chain,}$$

$\tilde{A}$  and  $A$  are similar  $\Rightarrow$  they share the same eigenvalues

[Coifman, Lafon, Nadler, '06]

def: Diffusion distance at time  $s$  between two nodes is  $d_G^s(x, x') = \|A^s \delta_x - A^s \delta_{x'}\|$

def: Diffusion distance between  $G$  and  $G'$  at time  $s$  is defined as

$$d(G, G') = \inf_{\Pi \in \Pi_n} \|A^{2s} - \Pi^+(A')^{2s} \Pi\|$$

As  $s$  increases, distance is weaker (why?)  
 Our distance results in a stronger topology than another usual distance defined over metric space, the Gromov-Hausdorff distance:

$$d_{GH}(G, G') = \inf_{\Pi \in \Pi_n} \sup_{x, x'} |d_G^s(x, x') - d_{G'}^s(\Pi(x), \Pi(x'))|$$

Here we used  $|V| = |V'| = n$ , but this can be easily extended to varying size using transportation plans rather than permutations.

### Graph Diffusion Scattering

We saw before that scattering transforms in the Euclidean domain are stable to deformations (= metric changes).

Can we extend that to non-Euclidean settings?

↳ We have defined the equivalent of a deformation (= diffusion metric).

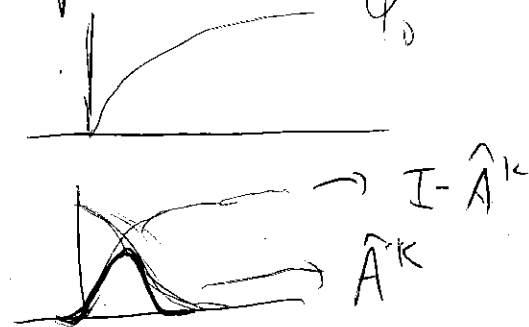
↳ We need the equivalent of wavelets.

In fact, we can also use diffusions on the graph to define a wavelet decomposition.

def: Diffusion wavelets are obtained by using powers of the adjacency:

$$\psi_0 = I - A, \quad \psi_j = A^{2^{j-1}} (I - A^{2^{j-1}}) \quad j > 0.$$

↳  $A$  is a low-pass filter: we are building multiscale filter bank by combining diffusions at different fine-scales.



Def: Diffusion Scattering Transform:

$$\Psi: L^2(G) \rightarrow (L^2(G))^{\mathbb{J}}$$

$$x \mapsto \{\psi_j(x)\}_j$$

$$\phi(x) = \{A^j x, A^j \psi(\Psi x), A^j \psi(W \psi(W x)), \dots\}.$$

Thm: [Game, Ribeiro, B. '19]: The graph scattering representation is stable to metric perturbations, as measured by graph diffusion metrics:

$$\|\Phi_G(x) - \Phi_{G'}(x)\| \leq \text{Ind}(G, G') \|x\|$$

GNN: Instead of predefining how to combine neighbors, the Laplacian diffusion  $A = D^{-1}L$ , we replace it by a data-driven diffusion.

## Part II: Optimization and Geometry in DL

• Prototypical problem in ~~ML~~ ML is of the form

$$\min_{x \in X} f(x)$$

(exceptions: saddle point problems, e.g. adversarial training)

• Characteristic:  $X \subseteq \mathbb{R}^d$ ,  $d$  huge.

• Worse case performance: searching for minimum is hard - exponential in  $d$ . How to overcome this?

First important class of functions that overcome curse of dimensionality: convex functions.

Assume for now  $X = \mathbb{R}^d$ .

Def:  $f$  convex if  $\forall x, y, t$ ,  $f(tx + (1-t)y) \leq tf(x) + (1-t)f(y)$ .

If  $f$  is differentiable, then

$$\lim_{t \rightarrow 0} \frac{f(tx + (1-t)y) - f(y)}{t} \leq f(x) - f(y)$$

$$\langle \nabla f(y), x - y \rangle \leq f(x) - f(y)$$

In general, we can extend this notion to non-differentiable functions, using the notion of sub-gradients.

def let  $f: \mathbb{R}^d \rightarrow \mathbb{R}$ .  $g \in \mathbb{R}^d$  is a subgradient of  $f$  at  $x$  if  $\forall y$ ,  $f(x) - f(y) \leq \langle g, x - y \rangle$ .

fact: If  $f$  is convex, then  $\partial f(x) = \{g; g \text{ is a subgradient at } x\}$  is non-empty. If  $f$  is differentiable and convex, then  $\nabla f(x) \in \partial f(x)$ .

### Gradient Descent for smooth functions

$\rightarrow$   $f$  continuously differentiable function is  $\beta$ -smooth if

$$\nabla f \text{ is } \beta\text{-Lipschitz: } \|\nabla f(x) - \nabla f(y)\| \leq \beta \|x - y\|$$

$\rightarrow$  if  $f$  is twice differentiable,  $\nabla^2 f \leq \beta I$ .

Thm:  $f$  convex and  $\beta$ -smooth in  $\mathbb{R}^d$ . Then gradient descent

$$x_{t+1} = x_t - \eta \nabla f(x_t) \text{ with } \eta = \beta^{-1} \text{ satisfies}$$

$$f(x_t) - \min_x f(x) \leq \frac{2\beta \|x_1 - x^*\|^2}{t-1}$$

We need two structural lemmas first.

Lemma:  $f$   $\beta$ -smooth. Then  $\forall x, y$

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

$$f(x) - f(y) = \int_0^1 \langle \nabla f(y + t(x-y)), x - y \rangle dt$$

$$|f(x) - f(y) - \langle \nabla f(y), x - y \rangle| = \left| \int_0^1 \langle \nabla f(y + t(x-y)) - \nabla f(y), x - y \rangle dt \right|$$

$$\leq \int_0^1 \|x - y\| \|\nabla f(y + t(x-y)) - \nabla f(y)\| dt$$

$$\leq \int_0^1 \|x - y\|^2 t \beta dt = \frac{\beta}{2} \|x - y\|^2. \quad \square$$

↳ If  $f$  is convex, since  $\nabla f(x) \in \partial f(x)$ , we have

$$0 \leq f(x) - f(y) - \langle \nabla f(y), x - y \rangle \leq \frac{\beta}{2} \|x - y\|^2$$

Choosing  $y = x - \beta^{-1} \nabla f(x)$ , we obtain  $\frac{\beta}{2} \|x - y\|^2$

~~$f(x - \beta^{-1} \nabla f(x)) - f(x) + \langle \nabla f(x), \beta^{-1} \nabla f(x) \rangle \leq \frac{\beta}{2} \|x - y\|^2$~~

$$f(x - \beta^{-1} \nabla f(x)) - f(x) + \langle \nabla f(x), \beta^{-1} \nabla f(x) \rangle \leq \frac{\beta}{2} \|x - y\|^2$$

$$f(x - \beta^{-1} \nabla f(x)) - f(x) + \beta^{-1} \|\nabla f(x)\|^2 \leq \frac{\beta^{-1}}{2} \|\nabla f(x)\|^2$$

$$f(x - \beta^{-1} \nabla f(x)) - f(x) \leq -\frac{\beta^{-1}}{2} \|\nabla f(x)\|^2$$

Lemma: (co-wercivity). For  $f$  convex and  $\beta$ -smooth, we have

$$f(x) - f(y) \leq \langle \nabla f(x), x - y \rangle - \frac{1}{2\beta} \|\nabla f(x) - \nabla f(y)\|^2$$

$$z = y - \frac{1}{\beta} (\nabla f(y) - \nabla f(x)).$$

$$f(x) - f(y) = f(x) - f(z) + f(z) - f(y)$$

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

$$= \langle \nabla f(x), x - y \rangle + \langle \nabla f(x) - \nabla f(y), y - z \rangle + \frac{1}{2\beta} \|\nabla f(x) - \nabla f(y)\|^2$$

$$= \langle \nabla f(x), x - y \rangle - \frac{1}{2\beta} \|\nabla f(x) - \nabla f(y)\|^2 \quad \square$$

Proof of thm: Using one-step improvement we have

$$f(x_{s+1}) - f(x_s) \leq -\frac{1}{2\beta} \|\nabla f(x_s)\|^2$$

$$\delta_s = f(x_s) - f(x^*) \rightarrow \delta_{s+1} \leq \delta_s - \frac{1}{2\beta} \|\nabla f(x_s)\|^2$$

using convexity

$$\delta_s \leq \langle \nabla f(x_s), x_s - x^* \rangle \leq \|x_s - x^*\| \|\nabla f(x_s)\|$$

Assume for now that  $\|x_s - x^*\|$  is decreasing  $s$ .

$$\text{Then } \|\nabla f(x_s)\| \geq \frac{\delta_s}{\|x_s - x^*\|} \geq \frac{\delta_s}{\|x_1 - x^*\|}$$

$$\Rightarrow \delta_{s+1} \leq \delta_s - \frac{\delta_s^2}{2\beta \|x_1 - x^*\|^2} = w \delta_s$$

$$\Rightarrow w \delta_s^2 + \delta_{s+1} \leq \delta_s \Rightarrow \frac{w \delta_s}{\delta_{s+1}} + \frac{1}{\delta_s} \leq \frac{1}{\delta_{s+1}} \Rightarrow \frac{1}{\delta_{s+1}} - \frac{1}{\delta_s} \geq w$$

$$\Rightarrow \frac{1}{\delta_{s+1}} \geq \frac{1}{\delta_s} + w$$



using previous lemma, we have

$$\frac{1}{\beta} \|\nabla f(x) - \nabla f(y)\|^2 \leq \langle \nabla f(x) - \nabla f(y), x-y \rangle$$

$$\begin{aligned} \text{Then } \|x_{s+1} - x^*\|^2 &= \|x_s - \beta^{-1} \nabla f(x_s) - x^*\|^2 \\ &= \|x_s - x^*\|^2 - \frac{2}{\beta} \langle \nabla f(x_s), x_s - x^* \rangle + \frac{1}{\beta^2} \|\nabla f(x_s) - \nabla f(x^*)\|^2 \\ &\leq \|x_s - x^*\|^2 - \frac{2}{\beta^2} \|\nabla f(x_s)\|^2 + \frac{1}{\beta^2} \|\nabla f(x_s)\|^2 \\ &\leq \|x_s - x^*\|^2 \quad \square \end{aligned}$$

↳ Rate  $1/t$  is not optimal in the class of  $\beta$ -smooth convex functions (we will see this later).

But first, let's push smoothness a bit further:

def:  $f$  is  $\alpha$ -strongly convex if

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

↳ if  $f$  is twice differentiable,  $\nabla^2 f(x) \succeq \alpha I$ .

↳ Strong convexity + smoothness means we can sandwich  $f(y)$  between two quadratic forms.

$$q_x^-(y) = f(x) + \langle \nabla f(x), y-x \rangle + \frac{\alpha}{2} \|x-y\|^2$$

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

$$q_x^-(y) \leq f(y) \leq q_x^+(y)$$

→ denote  $k = \frac{\beta}{\alpha} (\geq 1)$  the "condition" number of  $f$ .

lemma:  $f$   $\beta$ -smooth,  $\alpha$ -strongly convex. Then

$$\langle \nabla f(x) - \nabla f(y), x-y \rangle \geq \frac{\alpha\beta}{\alpha+\beta} \|x-y\|^2 + \frac{1}{\beta+\alpha} \|\nabla f(x) - \nabla f(y)\|^2$$

Recall that if  $f$  is convex and  $\beta$ -smooth, then

$$\langle \nabla f(x) - \nabla f(y), x-y \rangle \geq \frac{1}{\beta} \|\nabla f(x) - \nabla f(y)\|^2$$

If  $f$  is  $\beta$ -smooth,  $\alpha$ -strongly convex, then  $\varphi(x) := f(x) - \frac{\alpha}{2} \|x\|^2$

is convex and  $(\beta-\alpha)$ -smooth.

$$\Rightarrow \langle \nabla \varphi(x) - \nabla \varphi(y), x-y \rangle \geq \frac{1}{\beta-\alpha} \|\nabla \varphi(x) - \nabla \varphi(y)\|^2$$

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

$$\langle \nabla f(x) - \nabla f(y), x-y \rangle \geq \alpha \|x-y\|^2 + \frac{1}{\beta-\alpha} \left[ \|\nabla f(x) - \nabla f(y)\|^2 + \alpha \|x-y\|^2 - 2\alpha \langle \nabla f(x) - \nabla f(y), x-y \rangle \right]$$

Theorem:  $f$   $\beta$ -smooth and  $\alpha$ -strongly convex;  $\gamma = \frac{2}{\alpha+\beta}$  gives

$$f(x_{t+1}) - f(x^*) \leq \frac{\beta}{2} \exp\left(-\frac{4t}{k+1}\right) \|x_1 - x^*\|^2$$

Proof: By smoothness,  $f(x_t) - f(x^*) \leq \frac{\beta}{2} \|x_t - x^*\|^2$

$$\begin{aligned} \|x_{t+1} - x^*\|^2 &= \|x_t - \gamma \nabla f(x_t) - x^*\|^2 \\ &= \|x_t - x^*\|^2 - 2\gamma \langle \nabla f(x_t), x_t - x^* \rangle + \gamma^2 \|\nabla f(x_t)\|^2 \end{aligned}$$

$$\leq \left(1 - 2\frac{\gamma\alpha\beta}{\alpha+\beta}\right) \|x_t - x^*\|^2 + \left(\gamma^2 - 2\frac{\gamma}{\beta+\alpha}\right) \|\nabla f(x_t)\|^2$$

$$\leq (k-1)^2 \|x_t - x^*\|^2 \leq \dots$$

Oracle lower bounds:

↳ Smooth case: there exists  $\beta$ -smooth convex function  $f$  st for any black-box procedure with  $x_{t+1} \in \text{span}\{g_1, \dots, g_t\}$ ,

$$\min_{S \subseteq T} f(x_S) - f(x^*) \geq \frac{3}{32} \beta \frac{\|x_1 - x^*\|^2}{(t+1)^2} \quad (t < \frac{n}{2})$$

→ Remark: function is a Dirichlet Energy (quadratic).

↳ smooth & strongly convex case:  $\exists \beta$ -smooth and  $\alpha$ -strongly convex  $f$  st for any  $t$ ,

$$f(x_t) - f(x^*) \geq \frac{\alpha}{2} \left( \frac{\sqrt{k}-1}{\sqrt{k+1}} \right)^{2(t+1)} \|x_1 - x^*\|^2$$

Gap: How to close it?

→ Quadratic case:  $f(x) = \frac{1}{2} (x - x^*)^T H (x - x^*)$

$\nabla^2 f(x) = H$ , so we assume  $\alpha I \leq H \leq \beta I$ .

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

$$= x_t - \eta H(x_t - x^*) = [I - \eta H] x_t + \eta \underbrace{H x^*}_b$$

Suppose  $x_0 = 0$ ;  $x_{t+1} = \left( \sum_{k=0}^t [I - \eta H]^k \right) \eta b$

$$\alpha I \leq H \leq \beta I \Rightarrow \cancel{I} (1 - \eta\beta) I \leq I - \eta H \leq (1 - \eta\alpha) I$$

$\|I - \eta H\| < 1$   
if  $\eta < \beta^{-1} \Rightarrow (\eta H)^{-1} = [I - (I - \eta H)]^{-1} = \sum_{k=0}^{\infty} [I - \eta H]^k$

$$x_{t+1} = H_t^{-1} \eta b = \left[ (\eta H)^{t+1} + H_t - (\eta H)^t \right] \eta b$$

$$= x_0 + (H_t - (\eta H)^t) \eta b$$

$$\|H_t - (\eta H)^t\| = O(\|I - \eta H\|^t) = O\left(\left(1 - \frac{\alpha}{\beta}\right)^t\right)$$

$$= O\left(\left(1 - \frac{\alpha}{\beta}\right)^t\right)$$

→ Discrete vs continuous time optimization.

Consider  $x_{k+1} = x_k - \gamma \nabla f(x_k)$

Q: What happens when  $\gamma \rightarrow 0$ ? Do trajectories "pile-up"?

Assume  $\nabla f$  is bounded and Lipschitz (smooth).

Introduce the Ansatz  $x_{k\delta} \approx X(k\delta)$ , where  $X(t)$  is a smooth curve defined for  $t \geq 0 \in \mathbb{R}_+$ . Let  $\frac{t}{\delta} = k$ .

$$X(t) \Leftrightarrow x_{t/\delta} = x_k$$

$$X(t+\delta) \Leftrightarrow x_{k+1}$$

$$\frac{x_{k+1} - x_k}{\delta} = -\nabla f(x_k)$$

$$\downarrow$$

$$\frac{X(t+\delta) - X(t)}{\delta} = -\nabla f(x_t)$$

$$\downarrow$$

$$\dot{X}(t) = -\nabla f(X_t)$$

↳ Gradient flow.

Gradient Descent is a Forward-Euler discretization

We can get a first understanding of optimization algorithms by first looking at their behavior for  $t \rightarrow \infty$  continuous time.

Ex.  $f(x) = \frac{1}{2} x^T H x$ ,  $H$  <sup>(symmetric)</sup> pd. :  $H \succ \alpha I$ .

$$\dot{x}(t) = -H x(t) \Rightarrow x(t) = e^{-tH} x(0)$$

$$\Rightarrow \|x(t)\| \leq \|x(0)\| e^{-t\alpha}$$

Ex:  $f(x)$  convex, smooth.

$$\dot{x}(t) = -\nabla f(x(t))$$

Recall the convergence rate

we obtained:  $f(x_k) \leq \frac{\beta C \|x_k - x^*\|^2}{k} = \frac{C \|x_k - x^*\|^2}{k/\beta}$

Can we use continuous-time to derive it much <sup>st</sup> <sub>f</sub> <sup>quicker?</sup>

Consider the ~~approx~~ function

$$L(t) = t \cdot [f(x(t)) - f(x^*)] + \frac{1}{2} \|x(t) - x^*\|^2$$

$$\begin{aligned} \dot{L}(t) &= f(x(t)) - f(x^*) + t \langle \nabla f(x(t)), -\nabla f(x(t)) \rangle \\ &\quad + \langle \nabla f(x(t)), x(t) - x^* \rangle \end{aligned}$$

$$= \underbrace{-t \|\nabla f(x_t)\|^2}_{\leq 0} + \underbrace{f(x(t)) - f(x^*) + \langle \nabla f(x_t), x^* - x_t \rangle}_{\text{by convexity}}$$

$$L(t) \leq L(0) = \frac{1}{2} \|x(0) - x^*\|^2$$

$$t [f(x(t)) - f(x^*)] \leq \frac{\|x(0) - x^*\|^2}{2t}$$

Q: How good is gradient descent for the class of (smooth) convex functions?

A: A reasonable black-box procedure is a mapping from the past history to the next point:

$$(x_1, g_1, x_2, g_2, \dots, x_t, g_t) \mapsto x_{t+1}$$

we consider the "linear" case, where we ask

$$x_{t+1} \in \text{Span}(g_1, \dots, g_t)$$

Q: How to close this gap?

A: Through the notion of acceleration.

Consider  $x_{k+1} = y_k - \frac{1}{k+1} \nabla f(y_k)$

$$y_{k+1} = x_{k+1} + \frac{k}{k+3} (x_{k+1} - x_k)$$

Theorem (Nesterov '83):

$$f(x_t) - f^* \leq \frac{2 \|x_0 - x^*\|^2}{(t+1)^2 \alpha^2}$$

→ Acceleration for  $\alpha$ -strongly convex,  $\beta$ -smooth functions: just need to slightly modify scheme:

$$x_{k+1} = y_k - \beta^{-1} \nabla f(y_k)$$

$$y_{k+1} = x_{k+1} + \frac{1 - \sqrt{\alpha/\beta}}{1 + \sqrt{\alpha/\beta}} (x_{k+1} - x_k)$$

Theorem (Master):  $f(x_k) - f(x^*) \leq \beta \left(1 - \sqrt{\frac{\alpha}{\beta}}\right)^k \|x_0 - x^*\|^2$   
 ( $\sim (1 - \sqrt{k^{-1}})^k$  vs  $(1 - k^{-1})^k$  from GD).

↳ continuous-time interpretation?

↳ What is the corresponding ODE?

$$\ddot{X} + \frac{3}{\tau} \dot{X} + \nabla f = 0$$

↳ Theorem: [Su, Boyd, Candès, '15] The ODE has unique global solution in  $X \in C^2((0, \infty), \mathbb{R}^d) \cap C^1([0, \infty), \mathbb{R}^d)$ .

Moreover, as  $\eta \rightarrow 0$ , Nesterov's scheme converges to the ODE: for all fixed  $T > 0$ ,

$$\lim_{\eta \rightarrow 0} \sup_{k \leq T/\eta} \|x_k - X(k\sqrt{\eta})\| = 0$$

Simple consequence: we show easily that we get  $1/\epsilon^2$  convergence rate.

Theorem (S, B, C): Consider now

$$L(t) = t^2 (f(x(t)) - f^*) + 2 \left\| x + \frac{t}{2} \dot{x} - x^* \right\|^2$$

$$\dot{L}(t) = 2t (f(x(t)) - f^*) + t^2 \langle \nabla f, \dot{x} \rangle + 4 \left\langle x + \frac{t}{2} \dot{x} - x^*, \frac{3}{2} \dot{x} + \frac{t}{2} \ddot{x} \right\rangle$$

since  $\ddot{x} + \frac{3}{\tau} \dot{x} = -\nabla f$ ,

$$\frac{3}{2} \dot{x} + \frac{t}{2} \ddot{x} = -\frac{t}{2} \nabla f$$

$$\begin{aligned} \dot{L}(t) &= 2t (f(x(t)) - f^*) + t^2 \langle \nabla f, \dot{x} \rangle - 2t \langle \nabla f, x - x^* + \frac{t}{2} \dot{x} \rangle \\ &= 2t (f(x(t)) - f^*) - 2t \langle \nabla f, x - x^* \rangle \leq 0. \quad \square \end{aligned}$$

↳ Interpretation in the strongly convex case: Chebyshev polynomials. Recall we were approximating the inverse matrix

$(\eta H)^{-1}$  with a finite power series  $\sum_{j=0}^k [I - \eta H]^j$

$$\left\| \frac{(\eta H)^{-1}}{\lambda} - \sum_{j=0}^k [I - \eta H]^j \right\| \rightarrow 0 \text{ as } k \rightarrow \infty, \text{ with error } o((1 - \kappa)^k)$$

Q: What is the best polynomial of degree  $k$  to approximate  $\lambda^{-1}$ ?  $\leftarrow \min \|T_k(A) - \lambda^{-1}\|$

Lemma: (Chebyshev polynomial): There exist a polynomial  $q_{Tic}$  of degree  $O(\sqrt{(\frac{\beta}{\alpha}) \log(\frac{1}{\epsilon})})$  such that

$$q_{Tic}(0) = 1 \quad \text{and} \quad |q_{Tic}(x)| \leq \epsilon \quad \forall x \in [\alpha, \beta].$$

$q_{Tic}$  is a Chebyshev polynomial, computed recursively from  $q_{Tic-1}, q_{Tic-2} \Rightarrow$  Nesterov scheme.

$\rightarrow$  Further explain Chebyshev.

$$\min \|A^{-1} - q_{Tic}(A)\| \cong \min \|I - \underbrace{A q_{Tic}(A)}_{P_K(A)}\|$$

$\cdot P_{Tic}(A)$  commutes with  $A$ ;

$\lambda$  eigenvalue of  $A \Leftrightarrow P_{Tic}(\lambda) = 1 - \lambda q_{Tic}(\lambda)$  eigenvalue of  $P_{Tic}(A)$ .

$\cdot$  Eigenvalues of  $A$  are in  $\lambda \in [\alpha, \beta]$ .

$\cdot \min_{x \in (\alpha, \beta)} |P_{Tic}(x)|$  st  $P_{Tic}(0) = 1 \rightarrow$  Chebyshev polynomials.

$\hookrightarrow$  polynomial  $P_{Tic}$  of degree  $O(\sqrt{(\frac{\beta}{\alpha}) \log(\frac{1}{\epsilon})})$  st  $P_{Tic}(0) = 1$  and  $|P_{Tic}(x)| \leq \epsilon$  for  $x \in (\alpha, \beta)$ .

$\rightarrow$  Today: Further on discrete vs continuous time.  
Stochastic Gradient Descent.

$\hookrightarrow$  Previously: take a discrete scheme and derive a continuous-time limit.

$\hookrightarrow$  Q: Can we take opposite route?

Given convex function  $g$ , we need two steps:

1) Construct an ODE such that its solutions  $x(t)$  satisfy  $g(x(t)) - g(x^*) \leq t^{-p}$

2) Discretize the ODE such that convergence rate is preserved.

A note on discretization of ODEs.

$\frac{dy}{dt} = f(t, y)$  integrate from  $t_n$  to  $t_{n+1} = t_n + h$

$$y(t_{n+1}) - y(t_n) = \int_{t_n}^{t_{n+1}} f(t, y(t)) dt$$

Two basic strategies:

(i)  $\int_{t_n}^{t_{n+1}} f(t, y(t)) dt \approx h f(t_n, y(t_n))$ , yields

$$y(t_{n+1}) = y(t_n) + h f(t_n, y(t_n)) \quad \text{Explicit Euler Scheme.}$$

(ii)  $\int_{t_n}^{t_{n+1}} f(t, y(t)) dt \approx h f(t_{n+1}, y(t_{n+1}))$  thus

$$y_{n+1} = y_n + h f(t_{n+1}, y_{n+1}) \quad \text{Implicit Euler Scheme.}$$

# The Bregman Lagrangian [Wisibono et al. '16]

↳ we assume  $g$  ~~is~~  $\min_{\theta \in \Theta} g(\theta)$  admits a unique minimiser  $\theta^*$ .

↳ let  $h$  be another convex function.

df: A Bregman divergence in  $\Theta$  is

$$D_h(\theta, \eta) = h(\theta) - h(\eta) - \langle \nabla h(\eta), \theta - \eta \rangle$$

non-negative since  $h$  is convex,  $D_h(\theta, \theta) = 0$ ; but not symm.

Locally like a Hessian metric:

$$D_h(\theta, \eta) = \frac{1}{2} (\theta - \eta)^T \nabla^2 h(\eta) (\theta - \eta) + o(\|\theta - \eta\|)$$

df: Bregman Lagrangian is

$$\mathcal{L}(X, V, t) := e^{\alpha(t) + \gamma(t)} (D_h(X + e^{-\alpha(t)} V, X) - e^{\beta(t)} g(X))$$

$X$  position

$V$  velocity

$t$  time.

with  $\dot{\beta}(t) \leq e^{\alpha(t)}$   
 $\dot{\gamma}(t) = e^{\alpha(t)}$

$X_t \in \Theta, t \geq 0$

Action on the path is  $J(X) = \int_{t \geq 0} \mathcal{L}(X_t, \dot{X}_t, t) dt$

Optimization system: ~~minimize~~ curves minimize

satisfy the Euler-Lagrange equation:

$$\frac{d}{dt} \frac{\partial \mathcal{L}}{\partial V} = \frac{\partial \mathcal{L}}{\partial X}; \text{ by plugging the Bregman Lagrangian}$$

it becomes

$$\frac{d}{dt} \nabla h(X_t + e^{-\alpha(t)} \dot{X}_t) = -e^{\alpha(t) + \beta(t)} \nabla g(X_t)$$

Q: Do these solutions minimize  $g$ ? How fast?

Theorem: [Wisibono et al. '17] The solutions ~~of~~ to the Euler-Lagrange satisfy  $g(X_t) - g(\theta_0) \leq \mathcal{O}(e^{-\beta(t)})$

↳ Consider the Lyapunov function

$$E_t = D_h(\theta^*, X_t + e^{-\alpha(t)} \dot{X}_t) + e^{\beta(t)} (g(X_t) - g(\theta_0))$$

↳ ~~can~~ show as before that  $\dot{E}(t) \leq 0$ .

Optimal rate is achieved with the relation  $\beta(t) = e^{\alpha(t)}$

Q: How to discretize these ODEs while preserving rate?

$$\alpha(t) = \log p - \log t, \quad \beta(t) = p \log t + \log C, \quad \gamma(t) = p \log t$$

Resulting Euler-Lagrange becomes  $p > 0, C > 0$ .

$$\ddot{X}_t + \frac{p+1}{t} \dot{X}_t + Cp^2 t^{p-2} [\nabla^2 h(X_t + t p^{-1} \dot{X}_t)]^{-1} \nabla g(X_t) = 0$$

This is equivalent to a decoupled system of 1st order eqs:

$$\begin{cases} \dot{z}_t = X_t + \frac{t}{p} \dot{X}_t \\ (*) \quad \frac{d}{dt} \nabla h(z_t) = -C p t^{p-1} \nabla g(X_t) \end{cases}$$

A "naive" idea is to use a standard discretization scheme for (\*); eg forward-backward Euler.

$$\begin{cases} \eta_k = \underset{z}{\operatorname{argmin}} \{ C p k^{p-1} \langle \nabla g(\theta_k), z \rangle + \delta^{-p} D_h(z, \eta_{k-1}) \} \\ \theta_{k+1} = \frac{p}{k} \eta_k + \frac{k-p}{k} \theta_k \end{cases}$$

It is not stable, does not preserve the convergence rate.

↳ As it turns out, in order to recover matching convergence, we need to ~~recover~~ <sup>have</sup> some form of co-coercivity property.

$$\theta_{k+1} = \frac{p}{k+p} \eta_k + \frac{k}{k+p} \zeta_k, \quad \boxed{\text{"Nesterov estimate sequence"}}$$

$$\eta_k = \underset{z}{\operatorname{argmin}} \{ C p k^{p-1} \langle \nabla g(\zeta_k), z \rangle + \delta^{-p} D_h(z, \eta_{k-1}) \}$$

$$\zeta_k \text{ s.t. } \langle \nabla g(\zeta_k), \theta_k - \zeta_k \rangle \geq \frac{1}{\mu} \delta^{p/(p-1)} \|\nabla g(\zeta_k)\|^{p/p}$$

↳ This can be used to derive high-order <sup>(explicit)</sup> methods.

→ Some other related works } Symplectic Optim.  
Bettencourt, Wilson, Jordan  
Lyapunov analysis  
(Wilson, Recht, Jordan)

### Stochastic Gradient Descent

→ As before, consider the minimisation of function  $g$  defined in  $\mathbb{R}^d$ .

→ However, we no longer have access to  $\nabla g(\theta)$  directly; only given access to unbiased estimates  $\nabla g_n(\theta_n)$ .

Ex:  $g_n$  is the loss for a single data-point:

$$g_n(\theta) = \ell(\phi(x_n; \theta), y_n)$$

$$g(\theta) = \mathbb{E}_{(x,y) \sim \mathcal{D}} \ell(\phi(x; \theta), y)$$

### Stochastic Approximation [Robbins & Munro '51]

General setting: find the zeros of a function  $h: \mathbb{R}^d \rightarrow \mathbb{R}^d$  from random observations at certain points

→ If  $h = \nabla g$  we link to stochastic optimization.

Robbins & Munros  $\theta_n = \theta_{n-1} - \alpha_n [h(\theta_{n-1}) + \varepsilon_n]$

Questions: } ① conditions for convergence?  
② ...

Motivating example: estimate the mean from samples.

→ Starting from  $\theta_0 = 0$ , we get data  $x_n \in \mathbb{R}^d$

$$\theta_n = (1 - \gamma_n) \theta_{n-1} + \gamma_n x_n = \theta_{n-1} - \gamma_n (\theta_{n-1} - x_n)$$

$$\left( \begin{array}{l} \text{Ex: } \gamma_n = \frac{1}{n} \rightarrow \theta_n = \frac{1}{n} \sum_{k \leq n} x_k \\ \gamma_n = \frac{2}{n+1}, \theta_n = \frac{2}{n(n+1)} \sum_{k \leq n} k x_k \end{array} \right)$$

→ If  $x_n$  are iid with  $\mathbb{E} x_n = x$ ,  $\mathbb{E} \|x_n - x\|^2 = \sigma^2$

$$\theta_n - x = \prod_{k \leq n} (1 - \gamma_k) (\theta_0 - x) + \sum_{i \leq n} \prod_{k \geq i} (1 - \gamma_k) \gamma_i (x_i - x)$$

$$\mathbb{E} \|\theta_n - x\|^2 = \prod_{k \leq n} (1 - \gamma_k)^2 \|\theta_0 - x\|^2 + \sigma^2 \sum_{i \leq n} \gamma_i^2 \prod_{k \geq i} (1 - \gamma_k)^2$$

error has two contributions:

↳ ~~Approx~~ Initial conditions:  $\prod_{k \leq n} (1 - \gamma_k) \rightarrow 0$  as  $n \rightarrow \infty$ .

~~we~~ we want to forget them.

↳ Robustness to noise:  $\sum_{i \leq n} \gamma_i^2 \prod_{k \geq i} (1 - \gamma_k)^2 \rightarrow 0$  as  $n \rightarrow \infty$

→ If  $\gamma_n \rightarrow 0$ ,  $\log \prod_{k \leq n} (1 - \gamma_k) \approx - \sum_{k \leq n} \gamma_k$ , so

$\sum_{k \leq n} \gamma_k$  should diverge to  $\infty$ .  $\ominus$  in order to

Ex: if  $\gamma_n = C n^{-\alpha}$ ,

$$\alpha = 1 \quad \sum_{i \leq n} i^{-1} = \log n + C' + O(n^{-1})$$

$$\alpha > 1 \quad \sum_{i \leq n} i^{-\alpha} = C' + O(n^{1-\alpha}) \rightarrow \text{init conditions not forgotten.}$$

$$\alpha < 0 \quad \sum_{i \leq n} i^{-\alpha} = C' n^{1-\alpha} + O(1).$$

Noise Term: assume  $\gamma_n$  non increasing and  $\gamma_n \leq \frac{1}{\mu}$ .

$$\text{Then } \forall m \leq n, \sum_{k \leq n} \gamma_k^2 \prod_{k \geq i} (1 - \gamma_k)^2 = \sum_{i \leq n} \gamma_i^2 \prod_{k \geq i} (1 - 2\gamma_k)$$

Suppose  $\mu < 2$ . Then

$$\sum_i \gamma_i^2 \prod_{k \geq i} (1 - 2\gamma_k) \leq \sum_i \gamma_i^2 \prod_{k \geq i} (1 - \mu \gamma_k) =$$

$$= \sum_{i=1}^m \prod_{k \geq i} (1 - \mu \gamma_k) + \sum_{i=m+1}^n \prod_{k \geq i} (1 - \mu \gamma_k)$$

$$\leq \prod_{i=m+1}^n (1 - \mu \gamma_i) \sum_{k=1}^n \gamma_k^2 + \gamma_m \sum_{i=m+1}^n \gamma_i \prod_{k \geq i} (1 - \mu \gamma_k)$$

$$\stackrel{\text{concavity}}{\leq} \exp\left(-\mu \sum_{i=m+1}^n \gamma_i\right) \sum_{k=1}^n \gamma_k^2 + \frac{\gamma_m}{\mu} \left( O(n) \right).$$

↳ so we ~~need~~ require  $\gamma_n$  to go to 0.

Ex:  $\gamma_n = C/n \rightarrow$  noise term converges in  $O(1/n)$



## Convergence of stochastic optimization

→ When  $g$  is convex, we want to study how

$\theta_n \rightarrow \theta_*$  (convergence in iterates) or

$g(\theta_n) \rightarrow g(\theta_*)$  (convergence in value).

→ Several convergence criteria in the stochastic case:

• A.S. convergence:  $\text{Prob}(g(\theta_n) \rightarrow g(\theta_*)) = 1$

• Convergence in proba:  $\forall \epsilon > 0, \text{Prob}(|g(\theta_n) - g(\theta_*)| \geq \epsilon) \rightarrow 0$

• Convergence in Moments:  $E|g(\theta_n) - g(\theta_*)|^r \rightarrow 0$

→ R & M asymptotic normality [Fabian '68]

$$\gamma_n = Cn^{-1}$$

$$E(\theta_n - \theta_*)(\theta_n - \theta_*)^T \approx n^{-2} CA (\theta_0 - \theta_*)(\theta_0 - \theta_*)^T$$

$$+ n^{-1} C^2 (2CA - I)^{-1} \Sigma$$

$$A = D^2 g(\theta_*) \quad \Sigma = E(\epsilon_n \epsilon_n^T) \quad \text{sto convergence} \\ (\nabla g_n = \nabla g + \epsilon_n)$$

⇒  $2C \lambda_{\min}(A) \gg 1$  for convergence.

$C$  too small: no convergence (due to memory of init. cond)

$C$  too large: bias variance

## Polyak-Ruppert Averaging

• Previous R&M algo ~~is~~ suffers from sensitivity of step-size, and dependence on unknown conditioning of the problem.

• We modify bias-variance tradeoff by considering an averaging over the iterates:

$$\bar{\theta}_n = \frac{1}{n} \sum_{k \in S_n} \theta_k \quad \left( \bar{\theta}_n = \left(1 - \frac{1}{n}\right) \bar{\theta}_{n-1} + n^{-1} \theta_n \right)$$

Theorem: (Cesaro's) Suppose  $\theta_n \rightarrow \theta_*$  with rate  $\|\theta_n - \theta_*\|$

Then  $\bar{\theta}_n \rightarrow \theta_*$  with rate  $\bar{\alpha}_n \leq \frac{1}{n} \sum_{k \in S_n} \alpha_k$

~~convergence~~ <sup>one can show that</sup> if  $\sum \alpha_n < +\infty$  then rate is always  $\frac{1}{n}$ . We lose convergence speed, but we gain robustness.

Several convergence rates for ~~the~~ different assumptions of  $g$  and choices of step-size:

Global minimax rates [Nemirovsky '83, Agarwal '12, Bach '17]

Strongly convex case:  $O((\alpha n)^{-1})$   
 Non-strongly convex:  $O(n^{-1/2})$  both attained with averaged stochastic gradient descent

Take-home: In smooth problems,  $\delta_n \sim n^{-1/2}$  + averaging gives adaptivity to strong convexity [Bach, Moulines '11]

Q: constant step-size?

Setup with least squares optim:

$$g(\theta) = \mathbb{E}[(Y - \langle X, \theta \rangle)^2], \theta \in \mathbb{R}^d$$

• Generic Covariance  $H = \mathbb{E}[XX^T]$  (no assumption pd!)

• constant step-size  $\gamma = 1/4R^2$ , assuming  $\|X\| \leq R$ .

Theorem (bach, Moulines): Averaged stochastic GD satisfies:

$$\mathbb{E} g(\bar{\theta}_n) - g(\theta_*) \leq \frac{C}{n}$$

Matches statistical lower bound (Tsybakov '03).

→ The SGD step is

$$\theta_n = \theta_{n-1} - \gamma X_n (\langle X_n, \theta_{n-1} \rangle - y_n)$$

$$\theta_n | \theta_{n-1} = [I - \gamma X_n X_n^T] \theta_{n-1} + \gamma \underbrace{y_n X_n}_{z_n}$$

Homogeneous Markov chain

convergence to a stationary distribution  $\pi_\gamma$ .

Expectation is  $\bar{\theta}_\gamma = \int \theta \pi_\gamma(d\theta)$

• In the LS setting, it turns out that  $\bar{\theta}_\gamma$  is independent of  $\gamma$ , and satisfies

$$\bar{\theta} = (\mathbb{E}[XX^T] + \mathbb{E}[X^T y]) \text{ pseudo-inverse}$$

However,  $\theta_n$  does NOT converge to  $\theta_*$ ; it oscillates around it, with oscillations of order  $\sqrt{\gamma}$ .

How to produce a convergent sequence? Average!

Ergodic Theorem:  $\bar{\theta}_n \rightarrow \bar{\theta}_\gamma$  at rate  $O(1/\sqrt{n})$ .  
(CLT).

$$\theta_n = \theta_{n-1} - \gamma \nabla g_n(\theta_{n-1})$$

Stationary dis  $\pi_\gamma$

$$\int \theta \pi_\gamma(d\theta) = \int \theta \pi_\gamma(d\theta) - \gamma \int \nabla g(\theta) \pi(d\theta)$$

$$\Rightarrow \int \nabla g(\theta) \pi(d\theta) = 0$$

$\bar{\theta}_\gamma = \theta_*$  and  $\bar{\theta}_n = \frac{1}{n} \sum \theta_k$  converges to  $\bar{\theta}_\gamma$ .

~~... slightly stronger assumptions~~  
~~...  $\mathbb{E}(\theta_n - \theta_*)^2$~~

For non-quadratic objectives, what happens?

→  $\theta_n = \theta_{n-1} - \gamma \nabla g_n(\theta_{n-1})$   
 is also a Homogeneous Markov chain

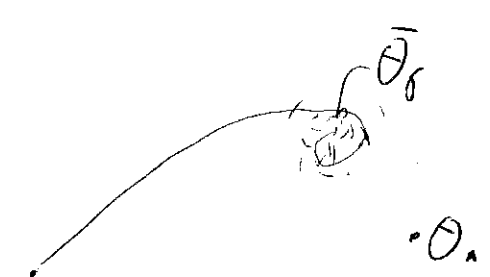
→ Its stationary distribution  $\pi_\gamma$  satisfies

$$\int \nabla g(\theta) \pi_\gamma(d\theta) = 0 \quad (g = \mathbb{E}g_n)$$

+ However, its mean  $\bar{\theta}_\gamma = \int \theta \pi_\gamma(d\theta)$  does not minimize  $g$ :

$$\nabla g(\bar{\theta}_\gamma) = \nabla g(\int \theta \pi_\gamma(d\theta)) \neq \int \nabla g(\theta) \pi_\gamma(d\theta) = 0$$

→ Thus,  $\theta_n$  oscillates around the "wrong" stationary point.  $\bar{\theta}_\gamma \neq \theta_*$



Moreover,  $\|\theta_n - \bar{\theta}_\gamma\| = O(\gamma)$

$$\bar{\theta}_n - \theta_* = \underbrace{\bar{\theta}_n - \bar{\theta}_\gamma}_{\text{stochastic error}} + \underbrace{\bar{\theta}_\gamma - \theta_*}_{\text{deterministic}}$$

→ [Bach, Durvasula, Dinevari '17] If  $g$  is strongly convex, for small  $\gamma$  we have

$$\bar{\theta}_n = \theta_* + \gamma \Delta + \mathcal{O}(\gamma^2)$$

with  $\|\Delta\| \leq C\gamma^2$ , and

$$\mathbb{E}[\bar{\theta}_{1c}^{(1)} - \theta_*] = \frac{A(\theta_*, \gamma)}{1c} + \gamma \Delta + \mathcal{O}(\gamma^2)$$

→ Consider now two chains  $(\bar{\theta}_{1c}^{(2\gamma)})_{1c}$ ,  $(\bar{\theta}_{1c}^{(\gamma)})_{1c}$  associated with  $\gamma$  and  $2\gamma$  respectively. Then we have that

$2\bar{\theta}_{1c}^{(2\gamma)} - \bar{\theta}_{1c}^{(\gamma)}$  satisfies

$$\mathbb{E}[2\bar{\theta}_{1c}^{(2\gamma)} - \bar{\theta}_{1c}^{(\gamma)} - \theta_*] = \frac{2A(\theta_*, 2\gamma) - A(\theta_*, \gamma)}{1c} + \alpha\gamma$$

Non-convex and SGD [Bottou, Curtis, Wouda '16]

→ How about non-convex case? (Optim Methods for Large-Scale ML)  
 we still assume smoothness.

Assumption 1:  $F: \mathbb{R}^d \rightarrow \mathbb{R}$  is continuously differentiable, and  $\nabla F$  is  $\beta$ -Lipschitz continuous.

$$F(x) \leq F(y) + \langle \nabla F(y), x-y \rangle + \frac{1}{2} \beta \|x-y\|^2$$

Generic SGD algorithm: stochastic gradient.

$$\theta_{1c+1} = \theta_{1c} - \eta_{1c} \nabla g(\theta_{1c}, \zeta_{1c})$$

Lemma:  $\mathbb{E}_{\zeta_{1c}}[F(\theta_{1c+1})] - F(\theta_{1c}) \leq -\eta_{1c} \langle \nabla F(\theta_{1c}), \mathbb{E}_{\zeta_{1c}} \nabla g(\theta_{1c}, \zeta_{1c}) \rangle + \frac{1}{2} \eta_{1c}^2 \beta \mathbb{E}_{\zeta_{1c}}[\|\nabla g(\theta_{1c}, \zeta_{1c})\|^2]$

$$F(\theta_{1c+1}) - F(\theta_{1c}) \leq \langle \nabla F(\theta_{1c}), \theta_{1c+1} - \theta_{1c} \rangle + \frac{1}{2} \beta \|\theta_{1c} - \theta_{1c+1}\|^2$$

$$\leq -\eta_{1c} \langle \nabla F(\theta_{1c}), \nabla g(\theta_{1c}, \zeta_{1c}) \rangle + \frac{1}{2} \eta_{1c}^2 \beta \|\nabla g(\theta_{1c}, \zeta_{1c})\|^2$$

→ So, if  $g(\theta_{1c}, z_{1c})$  is an unbiased estimator of

$\nabla F(\theta_{1c})$ , we obtain

$$\mathbb{E}_{z_{1c}} [F(\theta_{1c+1})] - F(\theta_{1c}) \leq -\eta_{1c} \|\nabla F(\theta_{1c})\|^2 + \frac{1}{2} \eta_{1c}^2 \beta (\mathbb{E} \|g(\theta_{1c}, z_{1c})\|^2)$$

→ Assumption: Denote by  $V_{z_{1c}}(g(\theta_{1c}, z_{1c})) = \mathbb{E} \|g\|^2 - \|\mathbb{E} g\|^2$  the variance of the stochastic gradient.

→ Assume  $V_{z_{1c}}[g] \leq M + M_V \|\nabla F\|^2$

for some scalars  $M, M_V$ .

and  $\mathbb{E}_{z_{1c}} g = \nabla F(\theta_{1c})$ .

→ Lemma:  $\mathbb{E}_{z_{1c}} [F(\theta_{1c+1})] - F(\theta_{1c}) \leq -(1 - \frac{1}{2} \eta_{1c} \beta (M+1)) \eta_{1c} \|\nabla F(\theta_{1c})\|^2 + \frac{1}{2} \beta \eta_{1c}^2 M$

□

Theorem (Non-convex Objective, Fixed Step-size).

Suppose  $F$  is bounded below, and  $\eta$  fixed and

$$0 < \eta < \frac{1}{\beta(1+M)}$$

Then  $\mathbb{E} \left[ \frac{1}{K} \sum_{k=1}^K \|\nabla F(\theta_{1k})\|^2 \right] \leq \frac{\eta \beta M + 2(F(\theta_1) - F_{\min})}{K \eta}$

→  $\eta \beta M$ .

Proof:  $\mathbb{E} [F(\theta_{1c+1})] - \mathbb{E} [F(\theta_{1c})] \leq$

$$\leq -\left(1 - \frac{1}{2} \eta \beta (1+M)\right) \eta \mathbb{E} \|\nabla F(\theta_{1c})\|^2 + \frac{1}{2} \eta^2 \beta M$$

$$\leq -\frac{1}{2} \eta \mathbb{E} \|\nabla F(\theta_{1c})\|^2 + \frac{1}{2} \eta^2 \beta M$$

$$F_{\min} - F(\theta_1) \leq \mathbb{E} [F(\theta_{K+1})] - F(\theta_1) \leq -\frac{1}{2} \eta \sum_{k=1}^K \mathbb{E} \|\nabla F(\theta_{1k})\|^2 + \frac{K}{2} \eta^2 \beta M$$

$$\Rightarrow \mathbb{E} \left[ \frac{1}{K} \sum_{k=1}^K \|\nabla F(\theta_{1k})\|^2 \right] \leq \frac{2(F(\theta_1) - F_{\min})}{K \eta} + \eta \beta M \quad \square$$

↳  $M=0$ ? No luck  $\Rightarrow \|\nabla F(\theta_{1c})\|_2 \rightarrow 0$

→  $\eta$  is in a tradeoff }  $\eta$  small  $\rightarrow$  average norm will be smaller, but it will take longer to converge!

→ Decreasing step-size? Suppose  $\eta_{1c}$  satisfy

$$\sum_{k=1}^{\infty} \eta_{1k} = +\infty, \quad \sum_{k=1}^{\infty} \eta_{1k}^2 < +\infty$$

Theorem: Let  $A_{1K} := \sum_{k \leq K} \eta_{1k}$ . Then, under previous

assumptions,  $\mathbb{E} \left[ \frac{1}{A_{1K}} \sum_{k \leq K} \eta_{1k} \|\nabla F(\theta_{1k})\|^2 \right] \rightarrow 0 \quad (K \rightarrow \infty)$ .

Proof: We know by hypothesis that  $\eta_{1c} \rightarrow 0$ , so wlog

$$\eta_{1c} \beta (1+M) \leq 1 \quad \forall c.$$

As before, we have

$$E[F(\theta_{k+1})] - E[F(\theta_k)] \leq -\frac{1}{2} \eta_k E[\|\nabla F(\theta_k)\|^2] + \frac{1}{2} \eta_k^2 \beta M$$

Summing again, we have

$$F_{\text{inf}} - E[F(\theta_1)] \leq -\frac{1}{2} \sum_{k=1}^K \eta_k E[\|\nabla F(\theta_k)\|^2] + \frac{1}{2} \beta M \sum_{k=1}^K \eta_k^2$$

$$\Rightarrow \sum_{k=1}^K \eta_k E[\|\nabla F(\theta_k)\|^2] \leq 2(F(\theta_1) - F_{\text{inf}}) + \beta M \sum_{k=1}^K \eta_k^2$$

$$\Rightarrow \frac{1}{A_k} \sum_{k=1}^K \eta_k E[\|\nabla F(\theta_k)\|^2] \leq \frac{C}{A_k} \rightarrow 0.$$

Corollary If we further assume that  $F$  is twice diff

and  $\theta \mapsto \|\nabla^2 F(\theta)\|$  is Lipschitz, then

$$\lim_{k \rightarrow \infty} E[\|\nabla F(\theta_k)\|^2] = 0.$$

$\hookrightarrow \epsilon$ -approximate first-order critical points

$$\{\theta; \|\nabla F(\theta)\| \leq \epsilon\}$$

Q: How fast can gradient descent reach first-order critical points?

Theorem (Nesterov '98):  $g$   $\beta$ -smooth.  $G \subset \mathbb{D}$  with step size  $\gamma = \beta^{-1}$  requires  $\beta \frac{g(\theta_1) - g^*}{\epsilon^2}$  iterations to reach a  $\epsilon$ -1st order stationary point of  $g$ .

Proof (simplified for continuous-time).

$$\dot{X}(t) = -\nabla g(X(t)) \quad g \text{ bounded below.}$$

$$g(X(T)) - g(X(0)) = \int_0^T \langle \nabla g(X(t)), \dot{X}(t) \rangle dt$$

$$= -\int_0^T \|\nabla g(X(t))\|^2 dt$$

$$\Rightarrow \int_0^T \underbrace{\|\nabla g(X(t))\|^2}_{h(t)} dt \leq K \quad \forall T.$$

$$\text{so } h(t) = O(t^{-1}).$$

Remarks: • No curse of dimensionality in this bound.

• Here we reach a 1st order stationary point.

In the case of convex functions, 1st order stationary is sufficient — not anymore.

$\rightarrow$  Classification of critical points.

def:  $\theta^*$  is a strict saddle point if  $\nabla g(\theta^*) = 0$  and

→ A local minima satisfies  $\lambda_{\min}(\nabla^2 g(\theta^*)) \geq 0$   
(not necessarily strict).

→ Every critical point  $(\nabla g(\theta^*) = 0)$  is an equilibrium point of gradient. Which are stable equilibria?

↳ Intuition: strict saddles are unstable.

Questions: A reasonable strategy to avoid unstable equilibria is to add noise (SGD). How  
Is noise necessary? sufficient?

→ Worst case analysis: GD can provably converge to saddle points with appropriately (bad) init [Kesten '09].  
But what about "generic" init?

→ Example: quadratic case: Consider the non-convex quadratic function  $g(\theta) = \frac{1}{2} \theta^T H \theta$ ,  $H = \text{diag}(k_1, \dots, k_d)$ ,

$$k_1, \dots, k_s > 0 \quad k_{s+1}, \dots, k_d < 0$$

↳ A single critical point  $\theta^* = 0$ , which is a strict saddle.

↳ Gradient Descent initialised at  $\theta_0$  is

$$\theta_{k+1} = \theta_k - \gamma \nabla g(\theta_k) = \sum_{i=1}^d (1 - \gamma k_i)^{k+1} \theta_{0,i} e_i$$

Then  $\|\theta_k\| \approx \dots$  A reasonable  $\dots$

→ If  $\theta_0$  is sampled from a distribution which is absolutely continuous w.r.t ambient dim, then GD avoids this strict saddle with probability 1.

→ How general is that?

Stable Manifold Theorem: assume  $g$  twice differentiable.

• A discrete-time optim algorithm is a mapping

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

$$\text{ex: GD} \quad \varphi(x) = x - \gamma \nabla g(x)$$

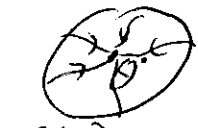
• Iterate  $k$  is obtained as  $\theta_k = \varphi^k(\theta_0)$ .

•  $\mathcal{X}^* = \{ \text{strict saddle points of } g \}$ .

• Definition (Global Stable Set): The global stable set of strict saddles is

$$S_\varphi = \{ \theta_0; \lim_{k \rightarrow \infty} \varphi^k(\theta_0) \in \mathcal{X}^* \}$$

$S_\varphi$  thus contains the initial values that eventually will land in a strict saddle.

→ In a neighborhood of a critical point  $\theta^*$ , how to capture the local attractive set, i.e.  the points around  $\theta^*$  that will be pushed to  $\theta^*$ ?

$$\text{Spec} \{ \dots \} \quad \theta^T \nabla^2 g(\theta^*) \theta > 0$$

- If this local attractor set has zero measure within a small neighborhood of  $\theta^*$ , then with probability 1 an init close to  $\theta^*$  will leave this neighborhood.
- Q: How to go from local to global stability? A priori, it is not immediate: escaping a saddle point does not guarantee we won't fall into another one later.

Def: Given a diffeomorphism  $\varphi$ , an unstable fixed point set is  $A_\varphi^* = \{ \theta; \varphi(\theta) = \theta; \mu_{\text{max}} \mu_{\text{min}} \|D\varphi(\theta)\| > 1 \}$

Theorem [Stable Manifold, Smale] [67]: Let  $\varphi$  be  $C^1$  (diffeo) mapping  $X \rightarrow X$  and  $\text{Det}(D\varphi(\theta)) \neq 0 \ \forall \theta \in X$ .

Then the set of initial points that converge to an unstable fixed point has measure zero:

$$\mu(\{ \theta_0; \lim_{k \rightarrow \infty} \varphi^k(\theta_0) \in A_\varphi^* \}) = 0.$$

Corollary: If  $X^* \subseteq A_\varphi$ , then  $\mu(S_\varphi) = 0$ .

Theorem [Gradient Descent Avoids Strict Saddles, Lee et al '16]

Assume  $g$  is  $\beta$ -smooth and  $\gamma < \beta^{-1}$ . Then if  $\theta^*$  is a strict saddle, then  $\Pr(\lim_{k \rightarrow \infty} \theta_k = \theta^*) = 0$

Proof:  $\rightarrow$  Every strict saddle of  $g$  is an unstable fixed point of  $\varphi(x) = x - \gamma \nabla g(x)$ .

$$D\varphi(x) = I - \gamma D^2 g(x)$$

$$\theta^* \text{ strict saddle} \Rightarrow \lambda_{\min}(D^2 g(\theta^*)) < 0 \Rightarrow$$

$$\Rightarrow \lambda_{\max}(I - \gamma D^2 g(\theta^*)) > 1.$$

$\rightarrow$  If  $\gamma < \beta^{-1}$ , then  $\det(D\varphi(\theta)) \neq 0 \ \forall \theta$ ,  $\lambda_{\min}(I - \gamma D^2 g(x)) > 0 \ \forall x$ .  $\square$

Extension

$\rightarrow$  Prove that  $\varphi(x)$  is a diffeo. it is invertible.

•  $\varphi$  injective: suppose  $x, y$  st  $\varphi(x) = \varphi(y)$

$$\Rightarrow x - y = \gamma (\nabla g(x) - \nabla g(y))$$

$$\|x - y\| = \gamma \|\nabla g(x) - \nabla g(y)\| \leq \gamma \beta \|x - y\|$$

but  $\gamma \beta < 1$ .

•  $\varphi$  surjective: Consider the implicit function

$$\phi(y) = \underset{x}{\text{argmin}} \frac{1}{2} \|x - y\|^2 - \gamma g(x) = \underset{x}{\text{min}} F_y(x)$$

since  $\gamma < \beta^{-1}$ ,  $F_y$  is strongly convex wrt  $x$ ,  $\Rightarrow$  unique min.

$$\text{Using KKT we have } \phi(y) - y - \gamma \nabla g(\phi(y)) = 0$$

$$y = \phi(y) - \gamma \nabla g(\phi(y)) = \varphi(\phi(y))$$

→ Extensions to proximal method, mirror descent, coordinate descent.

→ This result establishes that GD escapes strict saddles. Does this imply it converges to local minimizers?

Required additional properties

(i) Set of strict saddles cannot be too "large" countably finite / isolated saddles are ok.

(ii)  $\lim_k \theta_k$  exists.

Two sufficient conditions are

(i) Isolated critical points and compact sublevel sets.

(ii) Local Łojasiewicz inequality:  $\exists m, a, \epsilon$  st  $\|\nabla g(\theta)\| \geq m |g(\theta) - g(\theta^*)|^a$ ,  $a < 1$ .

for  $\theta \in \mathcal{N}(\theta^*)$ ;  $g(\theta^*) < g(\theta) < g(\theta^*) + \epsilon$   
↳ ensures length traveled by iterates of GD is finite.

→ This result shows that the tiniest amount of randomness (at init) is sufficient in mild assumptions to avoid saddles. What is the benefit of extra noise in the dynamics?

## Escaping from saddle points

→ We are interested in finding  $\epsilon$ -second order stationary points:  $D^2g$  Lipschitz, they are defined as

$$\|\nabla g(\theta^*)\| \leq \epsilon \quad \text{and} \quad \lambda_{\min}(D^2g(\theta^*)) \geq -\sqrt{\epsilon} \rho$$

with  $\rho = \text{Lip}(D^2g)$ .

→ Stochastic GD finds  $\epsilon$ -second-order critical points.

Theorem [Be'15] If  $g$  satisfies the strict saddle, then noisy GD converges to a local minimum in polynomial time.

"Vanilla" SGD:  $\theta_{k+1} = \theta_k - \gamma(\nabla g(\theta_k) + \epsilon_{k+1})$ ,

Time is dimension-dependent  $\tau \sim N^{10, I}$ .

(Several improvements in subsequent papers).

→ We have seen that both GD and SGD escape strict saddles.

↳ Dependency on input dim? cursed?

↳ Does more improve or decrease efficiency?

Theorem [Jin et al '17] Noisy Gradient descent finds  $\epsilon$ -second-order stationary point whp with  $\tilde{O}\left(\frac{\beta(g(\theta_0) - g^*)}{\epsilon^2}\right)$  (up to polylog factors).

→ Adding noise is sufficient to escape saddles.

→ Up to log factors, this matches GD rate for 1st-order



→ It turns out that noise is also necessary.  
 Even with uniform init, one can construct smooth functions for which GD requires exponential time  
~~Theorem~~ (Du et al '17).

### High-dimensional Energy Landscapes

- What are good models for non-convex, yet tractable high-dimensional energy landscapes?
- Models from statistical physics are very rich mathematical structure.

Eg Spherical spin glass:

$$E(\theta) = \frac{1}{N^{(k-1)/2}} \sum_{k_1, \dots, k_p=1}^N z_{k_1, \dots, k_p} \theta_{k_1} \dots \theta_{k_p}$$

$$z_p \sim N(0, \sigma^2), \quad \|\theta\|^2 = 1$$

Hamiltonian of the spherical  $k$ -spin spin glass

Kac-Rice Formula

$$\lim_{N \rightarrow \infty} \frac{1}{N} \log \mathbb{E} \text{Crt}_{N, K}(u) = \theta_{K, p}(u)$$



### Reproducing kernels and Kernel Dynamics

Goal: Infinite-dimensional least squares.  
~~hold / fitting~~

→ Generic learning in an infinite-dimensional space:

$$\min_{f \in \mathcal{F}} \|f^* - f\|^2. \quad \text{TKR / problem}$$

→ Since we measure error in an empirical way, we need some form of regularization in the fitting phase. What is the richest structure available?

Hilbert space → kernels.

→ Def: Given a set of objects  $X$ , a psd kernel is a symmetric function  $K: X \times X \rightarrow \mathbb{R}$  st for all finite sequences  $x_i \in X$  and  $\alpha_i \in \mathbb{R}$ ,  $i=1, \dots, n$ ,

$$\sum_{i,j} \alpha_i \alpha_j K(x_i, x_j) \geq 0.$$

Ex:  $K(x, x') = \langle \phi(x), \phi(x') \rangle$  for arbitrary mappings:  
 $\phi: X \rightarrow \mathcal{F}$ .  $\mathcal{F}$  Hilbert space.

→ Def

Theorem (Aronszajn '50):  $K$  is a psd Kernel if and only if there exists  $\mathcal{F}$  Hilbert space and a mapping  $\phi: \mathcal{X} \rightarrow \mathcal{F}$  st  $\forall (x, x') \quad K(x, x') = \langle \phi(x), \phi(x') \rangle_{\mathcal{F}}$ .

$\mathcal{F}$ : feature space;  $\phi(x)$  = feature map of  $x$ .

Examples:

→ Linear Kernel:  $K(x, y) = x^T y$ . ( $\phi(x) = x$ )

→ Polynomial Kernel:  $K(x, y) = (1 + x^T y)^d$  ( $\phi(x)$  = monomials)

↳ If  $K_1$  and  $K_2$  are kernels, then

$\left. \begin{array}{l} \rightarrow \alpha K_1 + \beta K_2 \text{ is a Kernel. } (\alpha, \beta > 0) \\ \cdot K_1, K_2(x, x') := K_1(x, y) \cdot K_2(x, y) \end{array} \right\} \begin{array}{l} \text{are} \\ \text{also} \\ \text{kernels.} \end{array}$

→ Gaussian Kernel  $K(x, y) = \exp(-\alpha \|x - y\|^2)$   
 $\phi(x) = ?$  finite-dimensional?

### Reproducing Kernel Hilbert Space

In Aronszajn Theorem, we have the representation

$$K(x, x') = \langle \phi(x), \phi(x') \rangle_{\mathcal{F}}$$

↳ We instantiate  $\mathcal{F}$  as a function space over  $\mathcal{X}$ .

$$\phi(x) = K(\cdot, x) \quad \mathcal{F}: \mathcal{X} \rightarrow \mathbb{R}$$

$$x' \mapsto K(x', x)$$

function evaluation:  $f(x) = \langle f, \phi(x) \rangle \quad f \in \mathcal{F}$

Reproducing property:  $K(x, y) = \langle K(\cdot, x), K(\cdot, y) \rangle \rightarrow$  "kernel Trick"

Evaluating of  $f$  at a point  $x$  is an inner product in the feature space.

Equivalently, the evaluation functional

$L_x: f \mapsto f(x) \quad \forall f \in \mathcal{H}$  is bounded:

$\exists M > 0$  st  $|L_x(f)| = |f(x)| \leq M \|f\|_{\mathcal{H}} \quad \forall f$ .

### Regularization and representer theorem

Suppose data  $x_i \in \mathcal{X}$ , labels  $y_i \in \mathcal{Y} \quad i=1, \dots, n$   
 kernel  $k$  with RKHS  $\mathcal{F}$ .

Learn by We consider

$$\min_{f \in \mathcal{F}} \sum_{i=1}^n \ell(y_i, f^T \phi(x_i)) + \frac{\lambda}{2} \|f\|^2$$

Representer Theorem (Kimeldorf & Wahba, '90) The minimiser is

of the form  $f = \sum_{j=1}^n \alpha_j \Phi(x_j) = \sum_{j=1}^n \alpha_j K(x_j, \cdot)$

Proof: etc.

## Kernel Ridge Regression

Data  $x_1, \dots, x_n \in X$   $y_1, \dots, y_n \in \mathbb{R}$ , psd kernel  $k$ .

\* Infinite-dimensional Least Squares

$$\min_{f \in \mathcal{F}} \frac{1}{n} \sum_{i=1}^n (y_i - f(x_i))^2 + \lambda \|f\|_{\mathcal{F}}^2$$

Option 1: Use previous representer theorem

$$f = \sum_{i=1}^n \alpha_i K(\cdot, x_i), \text{ obtained with}$$

$$\min_{\alpha \in \mathbb{R}^n} \frac{1}{n} \sum_{i=1}^n (y_i - (K\alpha)_i)^2 + \lambda \alpha^T K \alpha$$

$$\alpha^* = (K + n\lambda I)^{-1} y + \varepsilon, \quad K\varepsilon = 0$$

$f^*$  is unique.

Option 2: Suppose finite-dimensional feature space

$$\mathcal{F} \subset \mathbb{R}^d, \quad \Phi \in \mathbb{R}^{n \times d}$$

$$\min_{w \in \mathbb{R}^d} \frac{1}{n} \|y - \Phi w\|^2 + \lambda \|w\|^2$$

$$\rightarrow w = (\Phi^T \Phi + n\lambda I)^{-1} \Phi^T y$$

$$w = \underbrace{\Phi^T (\Phi \Phi^T + n\lambda I)^{-1}}_{\Phi^T \alpha} y = \Phi^T (K + n\lambda I)^{-1} y = \Phi^T \alpha$$

$\rightarrow$  Primal / Dual perspective: linear systems of different sizes. ( $n$ : # of data-points,  $d$ : # of features)  $\rightarrow$  dual perspective is more natural for generic kernels, where  $d = \infty$ !).

## Kernels and Neural Networks

$\rightarrow$  we have seen that kernels are dot products in generic feature spaces and that adding kernels is also a kernel.

$\rightarrow$  let  $\varphi(x, v)$   $x \in X$ ,  $v \in V$  be a "neuron".

we consider now a fixed probability measure  $\tau$  on  $V$ , assumed to be compact. ( $\varphi_v: X \rightarrow \mathbb{R}$  measurable.) Let

$$\mathcal{F} = \left\{ f : f(x) = \int_V p(v) \varphi_v(x) d\tau(v) \quad \forall x \in X \right\}$$

$p: V \rightarrow \mathbb{R}$  is square- $\tau$ -integrable

Prop (Bach '16):  $\mathcal{F}$  is an RKHS, with kernel

$$K(x, y) = \int_V \varphi_v(x) \varphi_v(y) d\tau(v)$$

Proof: Consider  $T: L_2(\tau) \rightarrow F$  defined by

$$(Tf)(x) := \int_V p(v) \varphi_v(x) d\tau(v). \text{ Let } K \text{ be the}$$

null space of  $T$ .  $T|_{K^\perp}$  is a bijection from

$K^\perp$  to  $F$ . We define  $\langle f, g \rangle_F := \int_V (T^{-1}f)(v) (T^{-1}g)(v) d\tau(v)$ .

•  $K(\cdot, y) \in F$  for all  $y \in X$ .

$$\begin{aligned} \langle f, K(\cdot, y) \rangle &= \int_V (T^{-1}f)(v) \cdot p(v) d\tau(v) = \\ &= \int_V (T^{-1}f)(v) \cdot \varphi_v(y) d\tau(v) = T(V^{-1}f)(y) = f(y) \end{aligned}$$

↳ Reproducing Property.

↳ Kernels as expectations are well-suited to approximations by Monte-Carlo sampling.

→ Consider  $v_1, \dots, v_m$  iid samples from  $\tau$ , and

$$\text{let } \hat{K}(x, y) = \frac{1}{m} \sum_{i=1}^m \varphi_{v_i}(x) \varphi_{v_i}(y).$$

→ What is Kernel regression in this setting?

Our predictors have the form

$$\frac{1}{m} \sum_{i=1}^m \alpha_i \varphi_{v_i}(x) \text{ adjusting } \alpha_i: \text{ feature learning}$$

As  $m \rightarrow \infty$ ,  $\hat{K} \rightarrow K$  (Rahimi & Recht '07)

Rate is essentially the Monte-Carlo Rate:  $\frac{1}{\sqrt{m}}$

Kernel Dynamics and the NTK

↳ Consider a generic parametric non-linear model

$$f(x; \theta) \quad (\text{eg } f(x; \theta) = \sum_{i=1}^n \alpha_i \varphi(x, \beta_i))$$

↳ We evaluate this model with the convex loss

$$R(f(x; \theta)) ; (\text{eg } R(f(x; \theta)) = \|f^* - f(x; \theta)\|^2)$$

$$\dot{\theta}_t = \nabla_{\theta} f(x, \theta_t)^\top R'(f(\cdot, \theta_t)) , \text{ this}$$

$$\dot{f}(\cdot, \theta_t) = \underbrace{\nabla_{\theta} f(\cdot, \theta_t)^\top \nabla_{\theta} f(\cdot, \theta_t)}_{K(t)}$$

↳ This update is the same as kernel regression,

$$\text{using the kernel } K_t(x, x') = \nabla f(x; \theta(t))^\top \nabla f(x'; \theta(t))$$

Fact: If  $\theta(t)$  does not move too much,

$K_t \approx K_0$ , thus the general GD dynamics will

be well approximated with Kernel dynamics.

↳ The assumption that  $K_t \approx K_0$  is equivalent to the assumption the model is linear: our parametric model thus becomes

$$T_f(x; \theta) = f(x; \theta_0) + (\theta - \theta_0) \cdot \nabla_{\theta} f(x; \theta_0)$$

↳ Up to an affine offset, this model is equivalent to a kernel method of the form  $K(x, x') = \nabla_{\theta} f(x; \theta_0)^T \nabla_{\theta} f(x'; \theta_0)$

### ↳ Lazy Training Regime

→ Present the lazy regime using linear approximation in function space.

Q: What is the corresponding linear model?

It is defined by the Kernel Tangent Kernel:

$$K(x, x') = \nabla_{\theta} f(x; \theta_0)^T \nabla_{\theta} f(x'; \theta_0)$$

→ Scaled objective function and comparison with kernel dynamics.

Theorem by Christ & Bach.

→ Analysis in terms of  $f(x; \theta) = \alpha(n) \sum_{i=1}^n \phi(x; \theta_i)$

The quantity that controls the validity of

linear approximation is roughly  $\frac{\|D\bar{f}_{\theta_0}\|^2}{L_{D\bar{f}}}$

If  $f_n \rightarrow f$  then this ratio converges to a constant. If we scale by something larger, it explodes.

$$\mathbb{E} \|D\bar{f}_{\theta_0}\|^2 = \alpha(n)^2 \mathbb{E} \|D\phi_{\theta_0}\|^2$$

$$L_{D\bar{f}} = \alpha(n) L_{D\phi}; \quad \mathbb{E} \|\nabla f(\theta_0)\|^2 = n \alpha(n) \mathbb{E} \|\phi\|^2$$

$$\mathbb{E} \| \bar{f}_n(\theta_*) - \bar{f}(\theta_*) \| \leq \max(1, \alpha(n)\sqrt{n})$$

Q: What happens in the regime  $\alpha(n) = 1/n$ ?

In this regime, we develop the particle ~~model~~ interacting systems:

$$\bar{f}_n(x; \theta) = \frac{1}{n} \sum_{i=1}^n \phi(x; \theta_i) \quad \theta = (\theta_1, \dots, \theta_n)$$

$\theta_i \in D$

$$L(\theta) = \frac{1}{2} \|\bar{f}_n(\theta) - f^*\|^2$$

$$= C_{f^*} - \langle \bar{f}_n(\theta), f^* \rangle + \frac{1}{2} \|\bar{f}_n(\theta)\|^2$$

$$= C_{f^*} - \frac{1}{n} \sum_{i=1}^n \langle \phi(\theta_i), f^* \rangle + \frac{1}{2n^2} \sum_{i,j=1}^n \langle \phi(\theta_i), \phi(\theta_j) \rangle$$

Let  $\hat{\mu}_n = \frac{1}{n} \sum_{i=1}^n \delta_{\theta_i}$  empirical measure

$$\bar{f}_n(\theta) = \int \phi(\theta) \mu_n(d\theta)$$

$$L(\theta) = \left( f - \int F(\theta) \mu_n(d\theta) + \iint K(\theta, \theta') \mu_n(d\theta) \mu_n(d\theta') \right)$$

$$F(\theta) = \langle \phi(\theta), f^\circ \rangle, \quad K(\theta, \theta') = \langle \phi(\theta), \phi(\theta') \rangle,$$

$$\dot{\theta}_i = -\nabla_{\theta_i} L(\theta) = -\frac{1}{n} \langle f^\circ, \nabla \phi(\theta_i) \rangle + \frac{1}{n^2} \sum_{\theta, \theta'} \langle \nabla \phi(\theta), \phi(\theta') \rangle$$

→ We have  $\frac{n}{2} \nabla_{\theta_i} L(\theta) = \nabla V(\theta_i; \mu_n)$  with

$$V(\theta) = -F(\theta) + \int K(\theta, \theta') \mu_n(d\theta').$$

↳ The gradient of each neuron corresponds to ~~feel~~ a velocity field evaluated at each location.

↳ Describe the evolution of a measure ~~as~~ subject to a velocity field?

Continuity equation:

$$\partial_t \mu_t = \operatorname{div}(\nabla V \cdot \mu_t)$$

$$\partial_t \left( \int f \mu_t \right) = - \int \langle \nabla f, \nabla V \rangle \mu_t \quad \forall f \in C_c^1$$

→ Important properties/questions:

(i) Convergence of this PDE?

(ii) Behavior of finite-particle dynamics around mean-field? Fluctuation analysis.

(iii) How is generalization affected under this regime?

↳ now features move at maximal speed.

(iv) Extension to deeper architectures?