Everything You Wanted to Know About the Loss Surface but Were Afraid to Ask – The Talk

Boris Hanin

Texas A&M

August 21, 2018

Boris Hanin Loss Surface

Plan

<ロ> <同> <同> < 同> < 同>

æ

Linear models:

- Baldi-Hornik (1989)
- Kawaguchi-Lu (2016)

æ

Э

P

Linear models:

- Baldi-Hornik (1989)
- Kawaguchi-Lu (2016)
- One hidden layer:
 - Ge-Lee-Ma (2016)
 - Mei-Montanari-Nguyen (2018)
 - Venturi-Bandeira-Bruna (2018)

Linear models:

- Baldi-Hornik (1989)
- Kawaguchi-Lu (2016)
- One hidden layer:
 - Ge-Lee-Ma (2016)
 - Mei-Montanari-Nguyen (2018)
 - Venturi-Bandeira-Bruna (2018)
- More than one hidden layer:
 - Gori-Tesi (1992) + Nguyen-Hein (2017)

But First ... Turns out Skip Connections are Good



Figure 1: The loss surfaces of ResNet-56 with/without skip connections. The vertical axis is logarithmic to show dynamic range. The proposed filter normalization scheme is used to enable comparisons of sharpness/flatness between the two figures.

Figure: From Visualizing the Loss Landscape of Neural Nets by Li, Xu, Taylor, Studer, Goldstein

But First ... Turns out Skip Connections are Good



(a) 110 layers, no skip connections

(b) DenseNet, 121 layers

Figure 6: (left) The loss surfaces of ResNet-110-noshort, without skip connections. (right) DenseNet, the current state-of-the-art network for CIFAR-10.

Figure: From Visualizing the Loss Landscape of Neural Nets by Li, Xu, Taylor, Studer, Goldstein

Boris Hanin Loss Surface

æ

Suppose

$$X = (x_j, 1 \leq j \leq N), \quad Y = (y_j, 1 \leq j \leq N), \ x_j \in \mathbb{R}^n, \ y_j \in \mathbb{R}^m.$$

æ

Suppose

$$X = (x_j, 1 \leq j \leq N), \quad Y = (y_j, 1 \leq j \leq N), \ x_j \in \mathbb{R}^n, \ y_j \in \mathbb{R}^m,$$

• Define empirical (co)-variances:

$$\Sigma_{YX} = \sum_{j} y_j x_j^T, \qquad \Sigma_{XX} = \sum_{j} x_j x_j^T, \qquad \Sigma_{YY} = \sum_{j} y_j y_j^T.$$

Suppose

$$X = (x_j, 1 \leq j \leq N), \quad Y = (y_j, 1 \leq j \leq N), \ x_j \in \mathbb{R}^n, \ y_j \in \mathbb{R}^m$$

• Define empirical (co)-variances:

$$\Sigma_{YX} = \sum_{j} y_j x_j^T, \qquad \Sigma_{XX} = \sum_{j} x_j x_j^T, \qquad \Sigma_{YY} = \sum_{j} y_j y_j^T.$$

• Classical least squares regression

$$A^* = \operatorname{argmin} \mathcal{L}(A) = \|AX - Y\|_F^2 = \sum_j \|Ax_j - y_j\|^2$$

is a convex problem and can be solved by GD or analytically by

$$A^* = \Sigma_{YX} \Sigma_{XX}^{-1}.$$

Boris Hanin Loss Surface

æ

• Q. What about rank-constrained least squares:

$$A_k^* = \operatorname{argmin} \mathcal{L}(A) = \|AX - Y\|_F^2, \quad \operatorname{rank}(A) \le k < n?$$

• Q. What about rank-constrained least squares:

$$A_k^* = \operatorname{argmin} \mathcal{L}(A) = \|AX - Y\|_F^2$$
, $\operatorname{rank}(A) \le k < n$?

• A. A_k^* is top k principal components of A^* .

• Q. What about rank-constrained least squares:

$$A_k^* = \operatorname{argmin} \mathcal{L}(A) = \|AX - Y\|_F^2$$
, $\operatorname{rank}(A) \le k < n$?

- **A.** A_k^* is top k principal components of A^* .
- **Q.** Note that if A is $n \times m$ then

$$\operatorname{rank}(A) \leq k \Rightarrow A = BC, B - n \times k, C - k \times m.$$

• Q. What about rank-constrained least squares:

$$A_k^* = \operatorname{argmin} \mathcal{L}(A) = \|AX - Y\|_F^2$$
, $\operatorname{rank}(A) \le k < n$?

- **A.** A_k^* is top k principal components of A^* .
- **Q.** Note that if A is $n \times m$ then

$$\operatorname{rank}(A) \leq k \Rightarrow A = BC, B - n \times k, C - k \times m.$$

• RC least squares \leftrightarrow linear net with one layer and L^2 loss:

$$(A^*, B^*) = \operatorname{argmin} \mathcal{L}(A, B) = \|ABX - Y\|_F^2$$

up to $A \mapsto AC, B \mapsto C^{-1}B, \quad C \in GL_k.$

• Q. What about rank-constrained least squares:

$$A_k^* = \operatorname{argmin} \mathcal{L}(A) = \|AX - Y\|_F^2$$
, $\operatorname{rank}(A) \le k < n$?

- **A.** A_k^* is top k principal components of A^* .
- **Q.** Note that if A is $n \times m$ then

$$\operatorname{rank}(A) \leq k \quad \Rightarrow \quad A = BC, \quad B - n \times k, \quad C - k \times m.$$

• RC least squares \leftrightarrow linear net with one layer and L^2 loss:

$$(A^*, B^*) = \operatorname{argmin} \mathcal{L}(A, B) = \|ABX - Y\|_F^2$$

up to $A \mapsto AC, B \mapsto C^{-1}B, \quad C \in GL_k.$

• Key. $\mathcal{L}(A, B)$ is not convex, so unclear if can solve by GD.

Boris Hanin Loss Surface

æ

- Suppose Σ_{XX} invertible and $\Sigma = \Sigma_{YX} \Sigma_{XX}^{-1} \Sigma_{XY}$ has full rank.
- **Thm.** All local minima of $\mathcal{L}(A, B)$ are global minima. Thus, GD should be fine.

- Suppose Σ_{XX} invertible and $\Sigma = \Sigma_{YX} \Sigma_{XX}^{-1} \Sigma_{XY}$ has full rank.
- **Thm.** All local minima of $\mathcal{L}(A, B)$ are global minima. Thus, GD should be fine.
- Thm. The saddle points of $\mathcal{L}(A, B)$ correspond to

$$AB = \operatorname{proj}_k(A^*),$$

where proj_k is projection onto some k principal component directions of A^* .

- Suppose Σ_{XX} invertible and $\Sigma = \Sigma_{YX} \Sigma_{XX}^{-1} \Sigma_{XY}$ has full rank.
- **Thm.** All local minima of $\mathcal{L}(A, B)$ are global minima. Thus, GD should be fine.
- Thm. The saddle points of $\mathcal{L}(A, B)$ correspond to

$$AB = \operatorname{proj}_k(A^*),$$

where proj_k is projection onto some k principal component directions of A^* .

• The proofs are by explicit computation.

Proof Idea Baldi-Hornik (1989)

• Write $\mathcal{L}(A, B) = \operatorname{vec} (ABX - Y)^T \operatorname{vec} (ABX - Y)$.

э

< □ > < □ >

Proof Idea Baldi-Hornik (1989)

• Write $\mathcal{L}(A, B) = \operatorname{vec} (ABX - Y)^T \operatorname{vec} (ABX - Y)$.

Use identity

$$\operatorname{vec}(ABC) = (C^T \otimes A) \operatorname{vec}(B)$$

to differentiate with respect to vec(A), vec(B).

Proof Idea Baldi-Hornik (1989)

• Write $\mathcal{L}(A, B) = \operatorname{vec} (ABX - Y)^T \operatorname{vec} (ABX - Y)$.

Use identity

$$\operatorname{vec}(ABC) = (C^T \otimes A) \operatorname{vec}(B)$$

to differentiate with respect to vec(A), vec(B).

• Check that at critical point W = (A, B), we must have

$$P_A \Sigma P_A = P_A \Sigma = P_A \Sigma,$$

where $P_A = \operatorname{proj}_{\operatorname{col}(A)}$.

Boris Hanin Loss Surface

æ

Э

• Want to understand local minima of

$$\mathcal{L}(W_1,\ldots,W_d) = \|W_d\cdots W_1X - Y\|_F^2,$$

when W_1, \ldots, W_d have fixed shape.

э

• Want to understand local minima of

$$\mathcal{L}(W_1,\ldots,W_d) = \|W_d\cdots W_1X - Y\|_F^2,$$

when W_1, \ldots, W_d have fixed shape.

• Thm. All local minima of \mathcal{L} are global minima.

• Want to understand local minima of

$$\mathcal{L}(W_1,\ldots,W_d) = \|W_d\cdots W_1X - Y\|_F^2,$$

when W_1, \ldots, W_d have fixed shape.

- Thm. All local minima of \mathcal{L} are global minima.
- Idea of Proof:
 - Severy local minimum $\{W_j\}$ can be perturbed to local minimum $\{\widehat{W}_j\}$ with same value of loss and all \widehat{W}_j of full rank.

• Want to understand local minima of

$$\mathcal{L}(W_1,\ldots,W_d) = \|W_d\cdots W_1X - Y\|_F^2,$$

when W_1, \ldots, W_d have fixed shape.

- Thm. All local minima of \mathcal{L} are global minima.
- Idea of Proof:
 - Every local minimum $\{W_j\}$ can be perturbed to local minimum $\{\widehat{W}_j\}$ with same value of loss and all \widehat{W}_j of full rank.
 - Perturbations of W_j's give all rank ≤ min {layer widths} pertrubations of ∏_j W_j.

• Want to understand local minima of

$$\mathcal{L}(W_1,\ldots,W_d) = \|W_d\cdots W_1X - Y\|_F^2,$$

when W_1, \ldots, W_d have fixed shape.

- Thm. All local minima of \mathcal{L} are global minima.
- Idea of Proof:
 - Severy local minimum $\{W_j\}$ can be perturbed to local minimum $\{\widehat{W}_j\}$ with same value of loss and all \widehat{W}_j of full rank.
 - Perturbations of W_j's give all rank ≤ min {layer widths} pertrubations of ∏_j W_j.
 - So Thus, $\left\{\prod_{j \neq \min} \widehat{W_j}, \widehat{W}_{\min}\right\}$ is local minimum of Baldi-Hornik problem. Hence, global minimum.

Loss for One Layer, Gaussian Inputs: Ge-Lee-Ma (2016)

Loss for One Layer, Gaussian Inputs: Ge-Lee-Ma (2016)

Input/Output Distribution. Inputs x_j ~ N(0, Id_n) i.i.d. and outputs y_i produced by net with one hidden layer:

$$y_j = a^{*T} \sigma(B^*x) = \sum_{k=1}^n a_k^* \sigma(\langle b_k^*, x_j \rangle), \quad ||b_k^*|| = ||a^*|| = 1.$$

Loss for One Layer, Gaussian Inputs: Ge-Lee-Ma (2016)

Input/Output Distribution. Inputs x_j ~ N(0, Id_n) i.i.d. and outputs y_j produced by net with one hidden layer:

$$y_j = a^{*T} \sigma(B^*x) = \sum_{k=1}^n a_k^* \sigma(\langle b_k^*, x_j \rangle), \quad ||b_k^*|| = ||a^*|| = 1.$$

• Thm. If $||a|| = ||b_k|| = 1$, then L^2 population risk $\mathbb{E} \left[||y(x, a, B) - y(x, a^*, B^*)||^2 \right]$

is a sum of rank *m* tensor-norms:

$$\sum_{m\geq 0} \widehat{\sigma}_m^2 \left\| \sum_{k=1}^n a_k b_k^{\otimes m} - \sum_{k=1}^n a_k^* (b_k^*)^{\otimes m} \right\|_F^2,$$

where $\widehat{\sigma}_m$ are the Hermite coefficients of σ :

$$\sigma(t) = \sum_{m \ge 0} \frac{\widehat{\sigma}_m}{m!} H_m(t).$$

æ

Orthogonality:

$$\int_{\mathbb{R}} H_n(x) H_m(x) e^{-x^2/2} \frac{dx}{\sqrt{2\pi}} = \delta_{n=m} n!.$$

æ

æ

Orthogonality:

$$\int_{\mathbb{R}} H_n(x) H_m(x) e^{-x^2/2} \frac{dx}{\sqrt{2\pi}} = \delta_{n=m} n!.$$

Sum-product formula for Hermite polynomials:

$$H_n(x \cdot y) = \sum_{|p|=n} \binom{n}{p} \prod_{j=1}^d H_{p_j}(x_j) y^{p_j}, \qquad x, y \in \mathbb{R}^d, \ \|y\| = 1,$$

Orthogonality:

$$\int_{\mathbb{R}} H_n(x) H_m(x) e^{-x^2/2} \frac{dx}{\sqrt{2\pi}} = \delta_{n=m} n!.$$

Sum-product formula for Hermite polynomials:

$$H_n(x \cdot y) = \sum_{|p|=n} \binom{n}{p} \prod_{j=1}^d H_{p_j}(x_j) y^{p_j}, \qquad x, y \in \mathbb{R}^d, \ \|y\| = 1,$$

where the sum is over all multi-indices p.

Wick formula

$$\mathbb{E}\left[\frac{H_k(v^T x)}{k!}\frac{H_j(w^T x)}{j!}\right] = \delta_{j=k} \quad \frac{1}{k!} \langle v, w \rangle^k.$$

Boris Hanin Loss Surface

Input/Output Distribution. Inputs x_j ~ ℙ i.i.d. and outputs y_j produced by net with one hidden layer:

$$y_j = \frac{1}{N} \sum_{i=1}^N \sigma_*(x; \theta_i) = \frac{1}{N} \sum_{i=1}^N a_i \sigma \left(\langle x, w_i \rangle + b_i \right).$$

Input/Output Distribution. Inputs x_j ~ ℙ i.i.d. and outputs y_j produced by net with one hidden layer:

$$y_j = \frac{1}{N} \sum_{i=1}^N \sigma_*(x;\theta_i) = \frac{1}{N} \sum_{i=1}^N a_i \sigma \left(\langle x, w_i \rangle + b_i \right).$$

• $L^2 \log R_N(\theta) = \frac{1}{2} \mathbb{E}[||y_N - \hat{y}||^2]$ is pair interaction in potential: $R(\rho) := \int V(\theta) d \rho^{(N)}(\theta) + \frac{1}{2} \int U(\theta_1, \theta_2) d\rho^{(N)}(\theta_1) d\rho^{(N)}(\theta_2),$

Input/Output Distribution. Inputs x_j ~ ℙ i.i.d. and outputs y_j produced by net with one hidden layer:

$$y_j = \frac{1}{N} \sum_{i=1}^N \sigma_*(x;\theta_i) = \frac{1}{N} \sum_{i=1}^N a_i \sigma \left(\langle x, w_i \rangle + b_i \right).$$

• L^2 loss $R_N(\theta) = \frac{1}{2}\mathbb{E}[||y_N - \hat{y}||^2]$ is pair interaction in potential:

$$\mathsf{R}(
ho) := \int \mathsf{V}(heta) d \
ho^{(\mathsf{N})}(heta) \ + \ rac{1}{2} \int \mathit{U}(heta_1, heta_2) \ d \
ho^{(\mathsf{N})}(heta_1) d \
ho^{(\mathsf{N})}(heta_2),$$

plus a constant where $d
ho^{(N)}(heta) = \sum_{i=1}^N \delta_{ heta_i}$ and

 $V(\theta) = -\mathbb{E}\left[y\sigma_*(x,\theta)\right], \quad U(\theta_1,\theta_2) = \mathbb{E}\left[\sigma_*(x,\theta_1)\sigma_*(x,\theta_2)\right].$

Input/Output Distribution. Inputs x_j ~ ℙ i.i.d. and outputs y_j produced by net with one hidden layer:

$$y_j = \frac{1}{N} \sum_{i=1}^N \sigma_*(x;\theta_i) = \frac{1}{N} \sum_{i=1}^N a_i \sigma \left(\langle x, w_i \rangle + b_i \right).$$

• L^2 loss $R_N(\theta) = \frac{1}{2}\mathbb{E}[||y_N - \widehat{y}||^2]$ is pair interaction in potential:

$$R(\rho) := \int V(\theta) d \rho^{(N)}(\theta) + \frac{1}{2} \int U(\theta_1, \theta_2) d \rho^{(N)}(\theta_1) d \rho^{(N)}(\theta_2),$$

plus a constant where $d
ho^{(N)}(heta) = \sum_{i=1}^N \delta_{ heta_i}$ and

$$V(heta) = -\mathbb{E}\left[y\sigma_*(x, heta)
ight], \quad U(heta_1, heta_2) = \mathbb{E}\left[\sigma_*(x, heta_1)\sigma_*(x, heta_2)
ight].$$

• **Q.** What does SGD looks like in ρ -space?

Boris Hanin Loss Surface

• **Thm.** Consider step sizes $s_k = \epsilon \xi(k\epsilon)$ and write $\rho_k^{(N)}$ for empirical measure after k GD steps. Then

$$\rho_{t/\epsilon}^{(N)} \implies \rho_t$$

as $N \to \infty$ and $\epsilon \to 0$, where ρ_t evolves under gradient flow for Wasserstein metric

$$egin{aligned} &\partial_t\,
ho_t = \xi(t)\,
abla_ heta\,(
ho_t\,
abla_ heta\Psi(heta,
ho_t)) \ &\Psi(heta,
ho) &:= V(heta) + \int U(heta, heta')d
ho(heta'). \end{aligned}$$

In particular, have a good approximation when $k = t/\epsilon$ with $\epsilon, N^{-1} \ll 1/{\rm input}$ -dim.

L^2 -loss in Deep and Wide Networks: Nguyen-Hein (2017)



L^2 -loss in Deep and Wide Networks: Nguyen-Hein (2017)

- σ smooth, analytic, bounded
- N distinct training samples with L^2 -loss

L²-loss in Deep and Wide Networks: Nguyen-Hein (2017)

- σ smooth, analytic, bounded
- N distinct training samples with L^2 -loss
- Thm. Fix a critical point. Suppose
 - ∃ layer k with $n_k \ge N 1$ neurons with Hessian of loss restricted to parameters in layers $\ge k + 1$ is non-degenerate
 - 2 weight matrices in layers $\geq k + 1$ all have full column rank.

Then this critical point is a global minimum.

• Rmk. Second condition requires "pyramidal" structure.

Boris Hanin Loss Surface

æ

• **Thm.** (Gori-Tesi 1992) If inputs are linearly independent, layer widths monotonically decreasing, then any crit at which the weight matrices have full rank is a global min.

- **Thm.** (Gori-Tesi 1992) If inputs are linearly independent, layer widths monotonically decreasing, then any crit at which the weight matrices have full rank is a global min.
- Idea. Backprop relates derivative of loss with respect to activations at successive layers:

$$\frac{\partial \mathcal{L}}{\partial \operatorname{Act}^{(j)}} = \frac{\partial \mathcal{L}}{\partial \operatorname{Act}^{(j+1)}} \underbrace{\mathcal{D}^{(j)} \mathcal{W}^{(j)}}_{Z^{(j)}},$$

where

$$D^{(j)} = \operatorname{Diag}\left(\phi'\left(\operatorname{preact}_{\beta}^{(j)}\right), \ \beta = 1, \ldots, n_j\right).$$

- **Thm.** (Gori-Tesi 1992) If inputs are linearly independent, layer widths monotonically decreasing, then any crit at which the weight matrices have full rank is a global min.
- Idea. Backprop relates derivative of loss with respect to activations at successive layers:

$$\frac{\partial \mathcal{L}}{\partial \operatorname{Act}^{(j)}} = \frac{\partial \mathcal{L}}{\partial \operatorname{Act}^{(j+1)}} \underbrace{\mathcal{D}^{(j)} \mathcal{W}^{(j)}}_{Z^{(j)}},$$

where

$$D^{(j)} = \operatorname{Diag}\left(\phi'\left(\operatorname{preact}_{\beta}^{(j)}\right), \ \beta = 1, \dots, n_j\right).$$

• If $\phi' \neq 0$ and $W^{(j)}$ has full rank, then can do "forward prop"

$$\frac{\partial \mathcal{L}}{\partial \operatorname{Act}^{(j+1)}} = \frac{\partial \mathcal{L}}{\partial \operatorname{Act}^{(j)}} Z^{(j) T} \left[Z^{(j)} Z^{(j) T} \right]^{-1}$$

- **Thm.** (Gori-Tesi 1992) If inputs are linearly independent, layer widths monotonically decreasing, then any crit at which the weight matrices have full rank is a global min.
- Idea. Backprop relates derivative of loss with respect to activations at successive layers:

$$\frac{\partial \mathcal{L}}{\partial \operatorname{Act}^{(j)}} = \frac{\partial \mathcal{L}}{\partial \operatorname{Act}^{(j+1)}} \underbrace{\mathcal{D}^{(j)} \mathcal{W}^{(j)}}_{Z^{(j)}},$$

where

$$D^{(j)} = \operatorname{Diag}\left(\phi'\left(\operatorname{preact}_{\beta}^{(j)}
ight), \ \beta = 1, \dots, n_j
ight).$$

• If $\phi' \neq 0$ and $W^{(j)}$ has full rank, then can do "forward prop"

$$\frac{\partial \mathcal{L}}{\partial \operatorname{Act}^{(j+1)}} = \frac{\partial \mathcal{L}}{\partial \operatorname{Act}^{(j)}} Z^{(j) T} \left[Z^{(j)} Z^{(j) T} \right]^{-1}$$

• Now just transplant this condition to some intermediate layer.

Boris Hanin Loss Surface

• Consider one layer network with no bias, L^2 loss, and hidden layer of width p

$$\Phi(x,\theta) = U\rho(Wx), \qquad L(\theta) := \mathbb{E}\left[\|\Phi(X,\theta) - Y\|^2 \right].$$

• Consider one layer network with no bias, L^2 loss, and hidden layer of width p

$$\Phi(x,\theta) = U\rho(Wx), \qquad L(\theta) := \mathbb{E}\left[\|\Phi(X,\theta) - Y\|^2\right].$$

• Consider span of possible features

$$V := \operatorname{Span} \left\{ \psi_{w}, \ w \in \mathbb{R}^{n} \right\}, \qquad \psi_{w} := \rho\left(\langle w, \cdot \rangle \right)$$

• Consider one layer network with no bias, L^2 loss, and hidden layer of width p

$$\Phi(x,\theta) = U\rho(Wx), \qquad L(\theta) := \mathbb{E}\left[\|\Phi(X,\theta) - Y\|^2 \right].$$

• Consider span of possible features

$$V := \operatorname{Span} \{ \psi_{w}, w \in \mathbb{R}^{n} \}, \qquad \psi_{w} := \rho(\langle w, \cdot \rangle)$$

Thm. If p ≥ dim(V), then there are no bad local minima. If p ≥ 2 dim(V), then all local=global minima are connected.

• Consider one layer network with no bias, L^2 loss, and hidden layer of width p

$$\Phi(x,\theta) = U\rho(Wx), \qquad L(\theta) := \mathbb{E}\left[\|\Phi(X,\theta) - Y\|^2 \right].$$

• Consider span of possible features

$$V := \operatorname{Span} \left\{ \psi_{w}, \ w \in \mathbb{R}^{n} \right\}, \qquad \psi_{w} := \rho\left(\langle w, \cdot \rangle \right)$$

- Thm. If p ≥ dim(V), then there are no bad local minima. If p ≥ 2 dim(V), then all local=global minima are connected.
- Ex. $\rho = polynomial$, linear

Boris Hanin Loss Surface

• Output $\Phi(x, \theta) \in V$.

- Output $\Phi(x,\theta) \in V$.
- V is reproducing kernel Hilbert space (RKHS):

- Output $\Phi(x,\theta) \in V$.
- V is reproducing kernel Hilbert space (RKHS):
 - Choose basis $\left\{\psi_{w_j}\right\}$ and define $\langle\cdot,\cdot\rangle$ with $\left\{\psi_{w_j}\right\}_j = \mathsf{ONB}$

- Output $\Phi(x,\theta) \in V$.
- V is reproducing kernel Hilbert space (RKHS):
 - Choose basis $\{\psi_{w_j}\}$ and define $\langle\cdot,\cdot\rangle$ with $\{\psi_{w_j}\}_j = \mathsf{ONB}$
 - Define reproducing kernel $K_V(x, y) := \sum_j \psi_{w_j}(x) \psi_{w_j}(y)$.

- Output $\Phi(x,\theta) \in V$.
- V is reproducing kernel Hilbert space (RKHS):
 - Choose basis $\{\psi_{w_j}\}$ and define $\langle\cdot,\cdot\rangle$ with $\{\psi_{w_j}\}_j = \mathsf{ONB}$
 - Define reproducing kernel $K_V(x,y) := \sum_j \psi_{w_j}(x) \psi_{w_j}(y)$.
 - $\phi(x) := K_V(x, \cdot) \in V$ is kernel of evaluation at x

- Output $\Phi(x,\theta) \in V$.
- V is reproducing kernel Hilbert space (RKHS):
 - Choose basis $\{\psi_{w_j}\}$ and define $\langle\cdot,\cdot\rangle$ with $\{\psi_{w_j}\}_j = \mathsf{ONB}$
 - Define reproducing kernel $K_V(x, y) := \sum_j \psi_{w_j}(x) \psi_{w_j}(y)$.
 - $\phi(x) := K_V(x, \cdot) \in V$ is kernel of evaluation at x
 - Then $\psi_w(x) = \langle \psi_w, \phi(x) \rangle$

- Output $\Phi(x, \theta) \in V$.
- V is reproducing kernel Hilbert space (RKHS):
 - Choose basis $\{\psi_{w_j}\}$ and define $\langle\cdot,\cdot\rangle$ with $\{\psi_{w_j}\}_j = \mathsf{ONB}$
 - Define reproducing kernel $K_V(x, y) := \sum_j \psi_{w_j}(x) \psi_{w_j}(y)$.
 - $\phi(x) := K_V(x, \cdot) \in V$ is kernel of evaluation at x
 - Then $\psi_w(x) = \langle \psi_w, \phi(x) \rangle$
- Obtain $\Phi = \langle U\psi(W), \phi(x) \rangle$

「同 ト イ ヨ ト イ ヨ ト ― ヨ