

Kernel theories of networks and their use in neuroscience

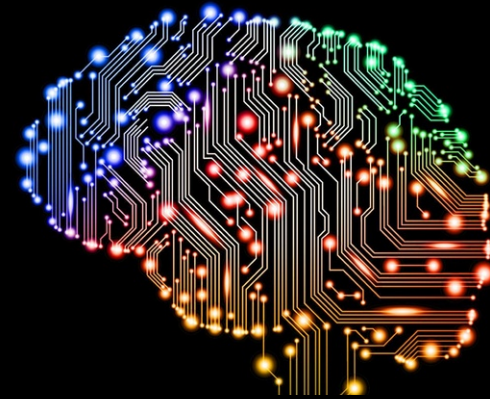
Kameron Decker Harris

University of Washington
Paul G. Allen School of Comp Science & Engineering, Biology

Western Washington University
Computer Science



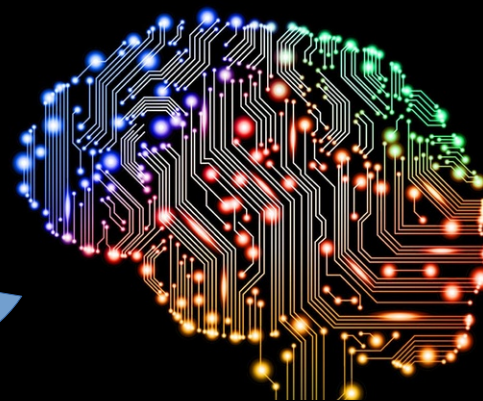
Neuro



CS



Neuro



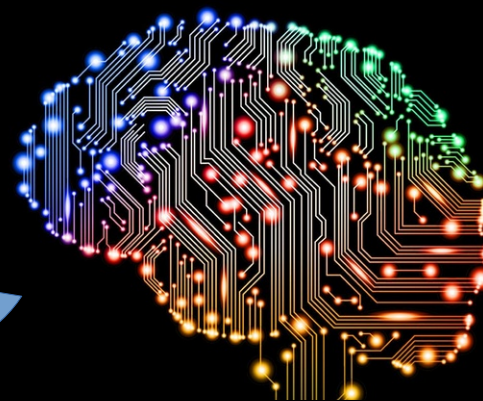
CS



Neuro



Big data management



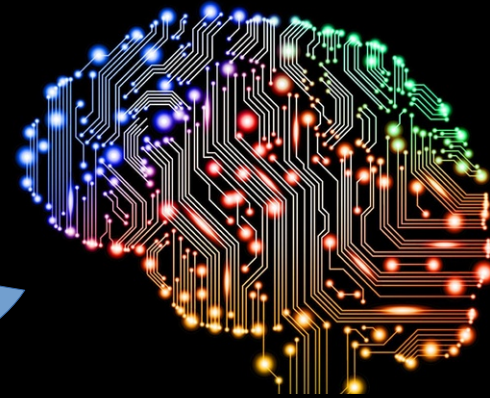
CS



Neuro



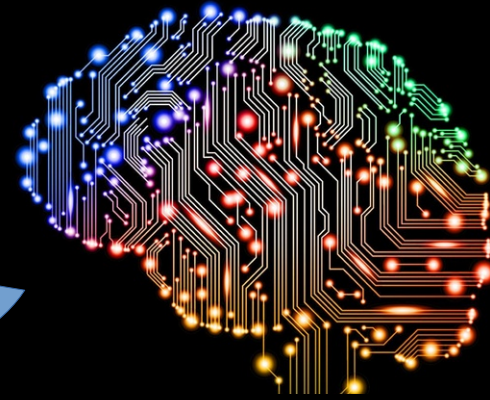
Big data management
Analysis of experiments



CS



Neuro

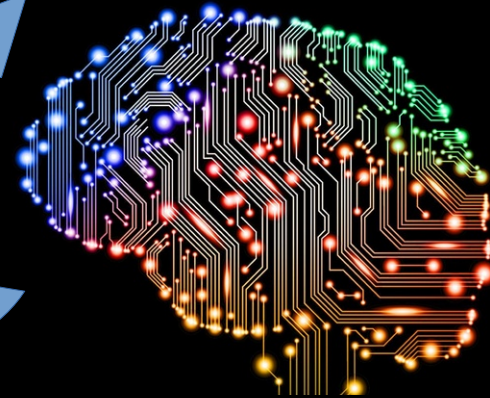
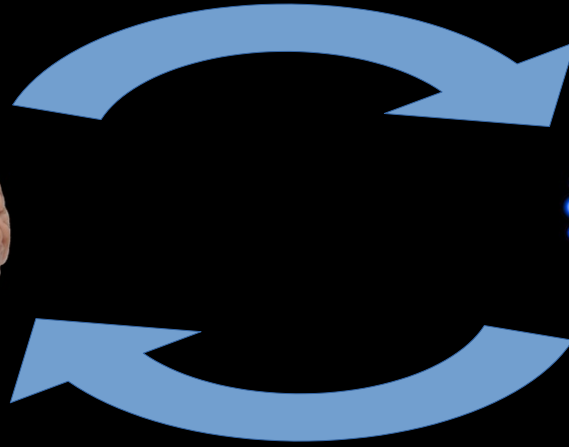


CS

Big data management
Analysis of experiments
Artificial neural networks



Neuro



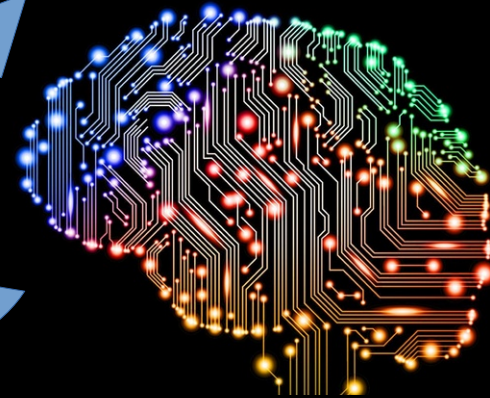
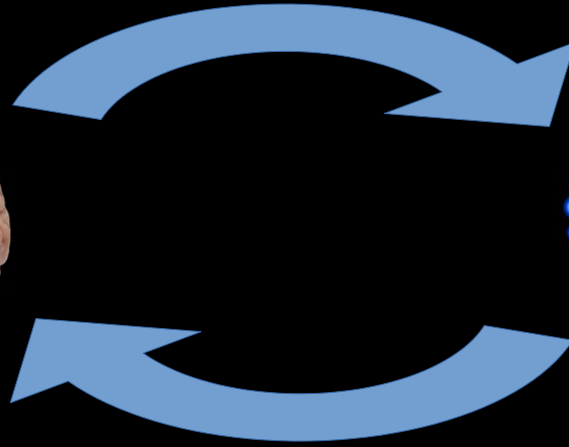
CS

Big data management
Analysis of experiments
Artificial neural networks

Brain-inspired algorithms



Neuro



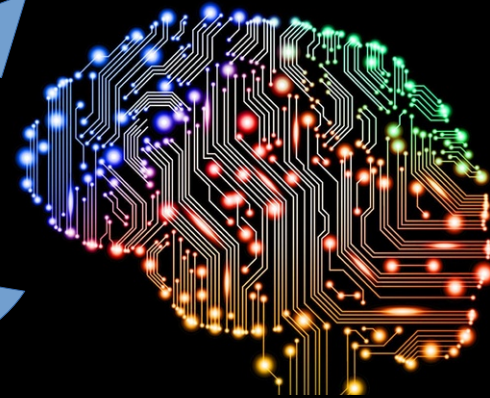
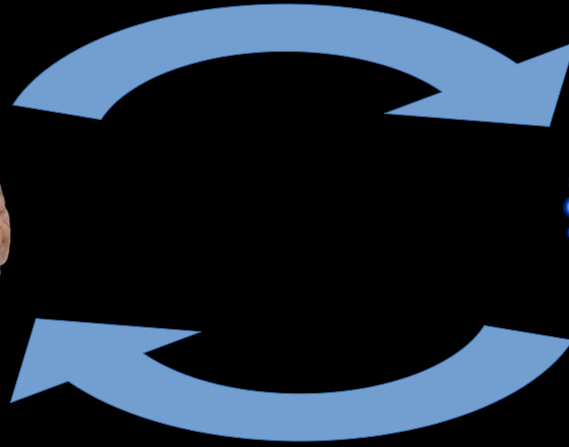
CS

Big data management
Analysis of experiments
Artificial neural networks

Brain-inspired algorithms
Emphasis on dynamics



Neuro



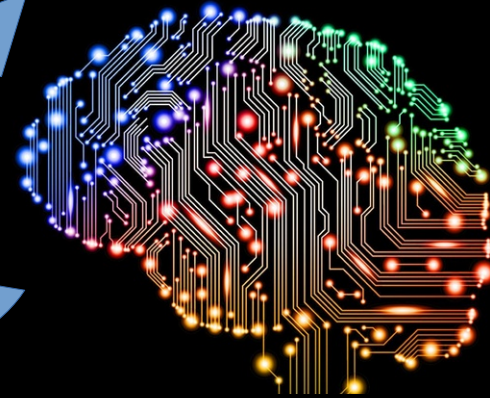
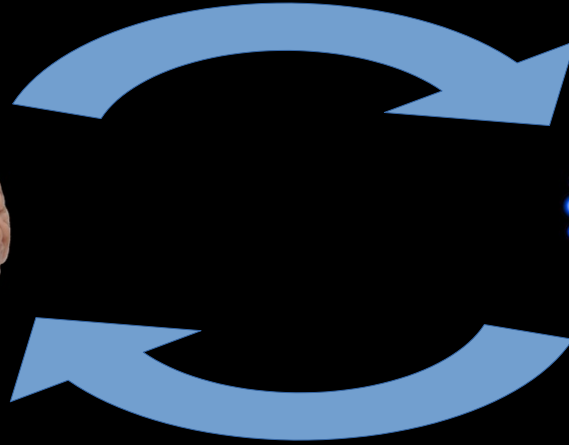
CS

Big data management
Analysis of experiments
Artificial neural networks

Brain-inspired algorithms
Emphasis on dynamics
Alternate models of computation



Neuro



CS

Big data management
Analysis of experiments
Artificial neural networks

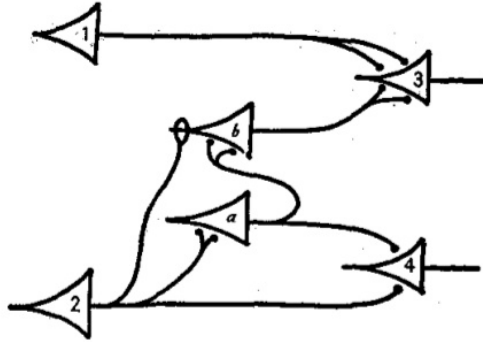
Overview

- Context: structure, randomness, & neuroscience
- Kernel theory:
 - Mathematical framework
 - Network sparsity \rightarrow additive functions
 - Tuning curves \rightarrow basis change
- Conclusions

Overview

- **Context: structure, randomness, & neuroscience**
- Kernel theory:
 - Mathematical framework
 - Network sparsity \rightarrow additive functions
 - Tuning curves \rightarrow basis change
- Conclusions

Structure and function



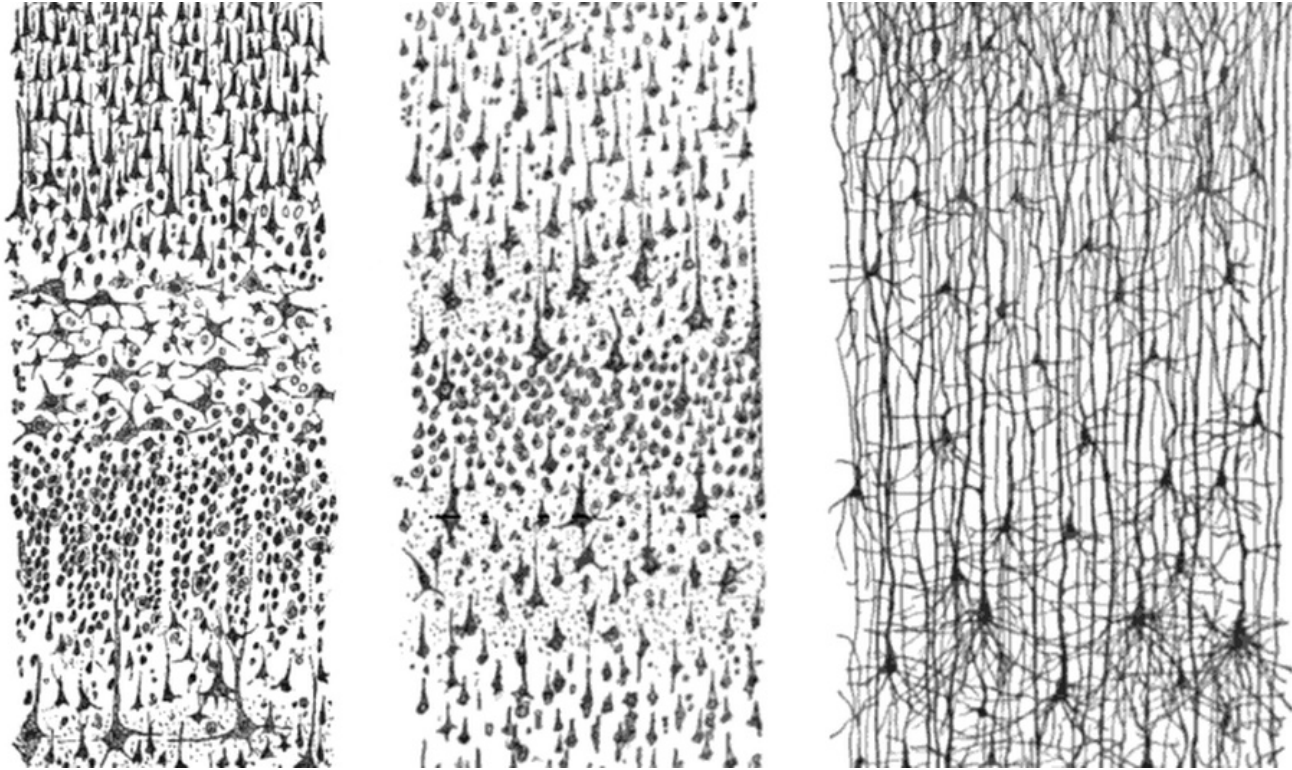
“ for any logical expression satisfying certain conditions, one can find a net behaving in the fashion it describes ”

A LOGICAL CALCULUS OF THE
IDEAS IMMANENT IN NERVOUS ACTIVITY

WARREN S. MCCULLOCH AND WALTER PITTS

1943

The messy reality of neuronal structure



Random vs stereotyped is not a dichotomy

Random vs stereotyped is not a dichotomy

- Some connections are stereotyped: CPGs, *Drosophila* Gal4 lines

Random vs stereotyped is not a dichotomy

- Some connections are stereotyped: CPGs, *Drosophila* Gal4 lines
- Some connections are nearly random: mushroom body, cortex?
 - Sophie Caron et al.

Random vs stereotyped is not a dichotomy

- Some connections are stereotyped: CPGs, *Drosophila* Gal4 lines
- Some connections are nearly random: mushroom body, cortex?
 - Sophie Caron et al.

Hypotheses:

Random vs stereotyped is not a dichotomy

- Some connections are stereotyped: CPGs, *Drosophila* Gal4 lines
- Some connections are nearly random: mushroom body, cortex?
 - Sophie Caron et al.

Hypotheses:

- Connections in many large networks are well-described by *random distributions with structure*

Random vs stereotyped is not a dichotomy

- Some connections are stereotyped: CPGs, *Drosophila* Gal4 lines
- Some connections are nearly random: mushroom body, cortex?
 - Sophie Caron et al.

Hypotheses:

- Connections in many large networks are well-described by *random distributions with structure*
- Plasticity and evolution modify these distributions

Overview

- Context: structure, randomness, & neuroscience
- **Kernel theory:**
 - **Mathematical framework**
 - Network sparsity \rightarrow additive functions
 - Tuning curves \rightarrow basis change
- Conclusions



Kernels: introduction & definition

X input vector l -dimensional $x \in X \subseteq \mathbb{R}^l$ compact

$k(x, x')$ kernel function computes "similarity"

requirements:

- $k(x, x') = k(x', x)$

- $K = (k(x_i, x_j))_{i,j=1}^n$

$n \times n$ kernel matrix

symmetric & positive definite

$$c^T K c > 0$$

- k continuous

$$k(x, x') = \langle \phi(x), \phi(x') \rangle$$





Examples of kernels



Examples of kernels

1) Linear

$$k(x, x') = \langle x, x' \rangle = \sum_{i=1}^l x_i x'_i$$



Examples of kernels

1) Linear

$$k(x, x') = \langle x, x' \rangle = \sum_{i=1}^l x_i x'_i$$

2) Polynomial

$$k(x, x') = (\underline{c} + \langle x, x' \rangle)^d$$



Examples of kernels

1) Linear

$$k(x, x') = \langle x, x' \rangle = \sum_{i=1}^l x_i x'_i$$

2) Polynomial

$$k(x, x') = (c + \langle x, x' \rangle)^d$$

3) Radial basis function

$$k(x, x') = \exp \left(-\frac{\|x - x'\|^2}{2\underline{\sigma^2}} \right) \quad \text{universal}$$

bandwidth

Eigendecomposition: Mercer's theorem

For a kernel which is symmetric, positive definite, continuous on a compact domain:

PD matrix

$$k(x, x') = \sum_{i=1}^{\infty} \lambda_i \psi_i(x) \psi_i(x')$$

eigenvalues
 ≥ 0

eigenvectors
orthonormal basis
for $L_2(\mathcal{X})$

Utility of kernel algorithms

Utility of kernel algorithms

- Kernels implicitly represent inputs in higher-dimensional features space called *reproducing kernel Hilbert space (RKHS)*

$$k(x, x') = \langle \underbrace{\phi(x)}, \phi(x') \rangle$$

Utility of kernel algorithms

- Kernels implicitly represent inputs in higher-dimensional features space called *reproducing kernel Hilbert space (RKHS)*

$$k(x, x') = \langle \phi(x), \phi(x') \rangle$$

- Algorithms can leverage this and just work with kernel matrix:
SVM, ridge regression, PCA, etc.

$K_{n \times n}$

$$\underline{f(x)} = y \underbrace{(K + \alpha I)^{-1}}_{\substack{\text{matrix} \\ \text{solve}}} \kappa(x), \quad \kappa(x) = (k(x_i, x))_{i=1}^n$$

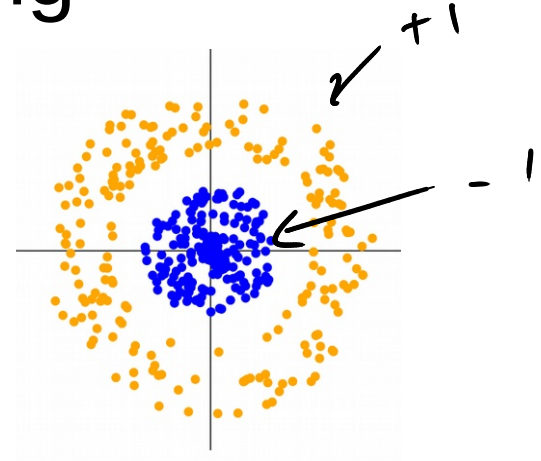
$$\text{Tr} \left(K (K + \alpha \underline{I})^{-1} \right)$$

$\swarrow \quad \searrow$
dim

Kernels and learning

Given: examples $(\mathbf{x}^i, y_i)_{i=1}^n$

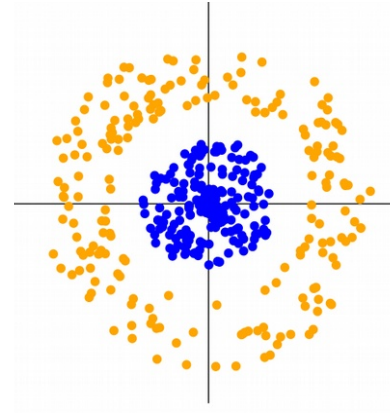
Find: f so that $f(\mathbf{x}) \approx y$



Kernels and learning

Given: examples $(\mathbf{x}^i, y_i)_{i=1}^n$

Find: f so that $f(\mathbf{x}) \approx y$



Key parameters

RKHS norm of the target function

$$f \in \mathcal{H} : \|f\|_{\mathcal{H}}$$

Dimensionality of the kernel matrix a.k.a. *Rademacher complexity*

significant eigenvalues of K \leftarrow depended on χ

Kernel theory of networks

Related review: “Randomness in neural networks” by Scardapane & Wang (2017)

Kernel theory of networks

Developed for big datasets, i.e. n huge

$$\vec{y} (K + \alpha I)^{-1}$$

- Rahimi & Recht (2008) **random features** ... sketching K
- older work in Gaussian processes by Neal (1996), Williams (1997)

Kernel theory of networks

Developed for big datasets, i.e. n huge

- Rahimi & Recht (2008) **random features** ... sketching K
- older work in Gaussian processes by Neal (1996), Williams (1997)

Exciting work tries to understand success of ANNs trained via gradient descent

- neural tangent kernel (NTK) – Jacot, Gabriel, Hongler, 2018; Arora et al, 2019
- convolutional kernel networks (CKN) – Mallat; Bruna; Harchaoui; Chizat et al
- interpolation & double descent – Belkin et al; Mei & Montanari

Kernel theory of networks

Developed for big datasets, i.e. n huge

- Rahimi & Recht (2008) **random features** ... sketching K
- older work in Gaussian processes by Neal (1996), Williams (1997)

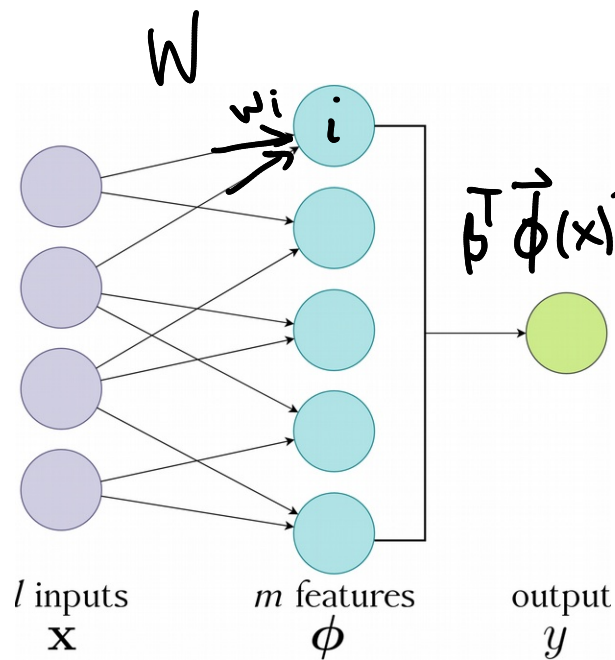
Exciting work tries to understand success of ANNs trained via gradient descent

- neural tangent kernel (NTK) – Jacot, Gabriel, Hongler, 2018; Arora et al, 2019
- convolutional kernel networks (CKN) – Mallat; Bruna; Harchaoui; Chizat et al
- interpolation & double descent – Belkin et al; Mei & Montanari

Barely applied in neuroscience, so far



From random network to kernel



class. ± 1
reg. \mathbb{R}

$$\phi_i(x) = h(w_i^T x)$$

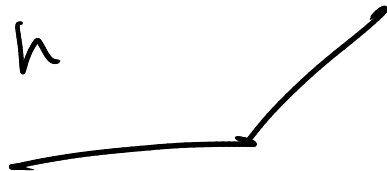
Assume $w_i \sim \mu$ iid random

Then assuming h Lipschitz, 2nd moments

$$\begin{aligned} \frac{1}{m} \sum_{i=1}^m \phi_i(x) \phi_i(x') \\ = \frac{1}{m} \vec{\phi}(x)^T \vec{\phi}(x') \end{aligned}$$

$$\begin{aligned} m \rightarrow \infty \rightarrow \mathbb{E}[\phi(x) \phi(x')] &= \int \overset{\text{scalars}}{h(\underbrace{w^T x}_{\text{scalar}}) h(\underbrace{w^T x'}_{\text{scalar}})} d\mu(w) \\ &\equiv k(x, x') \end{aligned}$$

Convergence rate to kernel



Harris, 2019 (informal version)

Claim The random map $\frac{1}{m}\phi(\mathbf{x})^\top\phi(\mathbf{x}')$ with κ -Lipschitz nonlinearity uniformly approximates $k_D^{\text{dist}}(\mathbf{x}, \mathbf{x}')$ to error ϵ using $m = \Omega(\frac{l\kappa^2}{\epsilon^2} \log \frac{C}{\epsilon})$ many features

Examples of random feature kernels

Examples of random feature kernels

Gaussian $w \sim N(0, \sigma^{-2}I)$ + Fourier $\exp(iw^\top x)$ FT of Gaussian
= Gaussian

$$k(x, x') = \exp\left(-\frac{\|x - x'\|^2}{2\sigma^2}\right)$$

RBK kernel
Rahimi & Recht, 2008

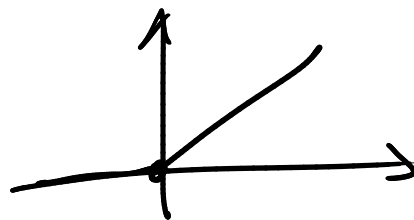
Examples of random feature kernels

Gaussian $w \sim N(0, \sigma^{-2}I)$ + Fourier $\exp(iw^\top x)$

$$k(x, x') = \exp\left(-\frac{\|x - x'\|^2}{2\sigma^2}\right)$$

RBf kernel
Rahimi & Recht, 2008

Gaussian + rectified polynomial



$$k(x, x') = \underbrace{\|x\|^p}_{\sim} \underbrace{\|x'\|^p}_{\sim} J_p \left(\arccos \frac{x^\top x'}{\|x\| \|x'\|} \right)$$

dot product kernel
(RBF on sphere)
Cho & Saul, 2009

Examples of random feature kernels

Gaussian $w \sim N(0, \sigma^{-2}I)$ + Fourier $\exp(iw^\top x)$

$$k(x, x') = \exp\left(-\frac{\|x - x'\|^2}{2\sigma^2}\right)$$

RBf kernel
Rahimi & Recht, 2008

Gaussian + rectified polynomial

$$k(x, x') = \|x\|^p \|x'\|^p J_p\left(\arccos \frac{x^\top x'}{\|x\| \|x'\|}\right)$$

dot product kernel
(RBF on sphere)
Cho & Saul, 2009

Almost always take uncorrelated, Gaussian weights
No network or input structure

The rest of the talk

Build in structure that occurs in neural systems

Show what the kernel theory says

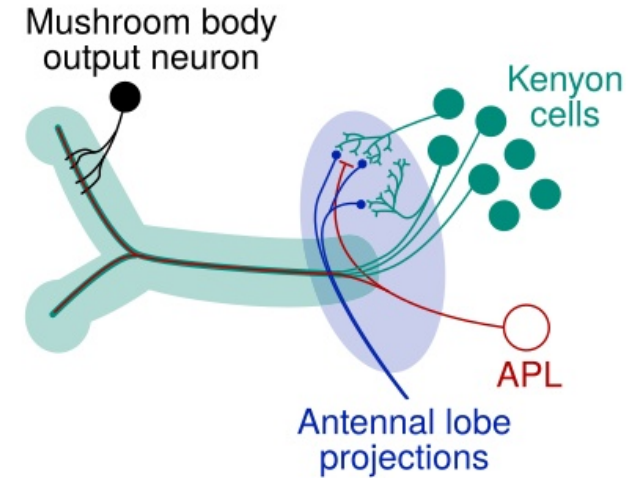
Overview

- Context: structure, randomness, & neuroscience
- **Kernel theory:**
 - Mathematical framework
 - **Network sparsity** → **additive functions**
 - Tuning curves → basis change
- Conclusions

Harris. “Additive function approximation in the brain.” NeurIPS Neuro/AI workshop, 2019

Litwin-Kumar, Harris, Axel, Sompolinsky, Abbott. “Optimal degrees of synaptic connectivity.” Neuron, 2017

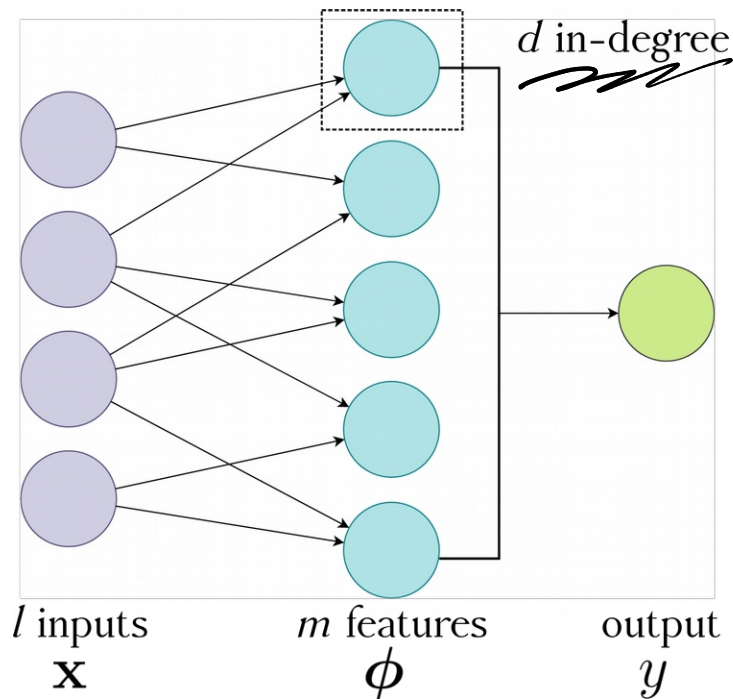
Olfactory network of *Drosophila*



Output neuron decides:
good smell or bad smell?

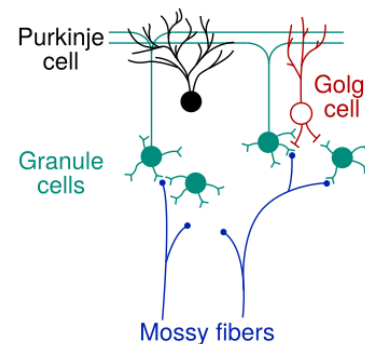
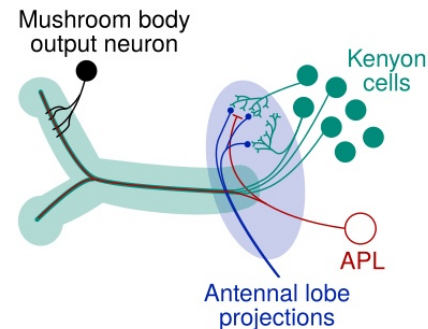
Common brain network structure: 2-layer sparse expansion

$$d < l$$



Mushroom body
 $d = 7$

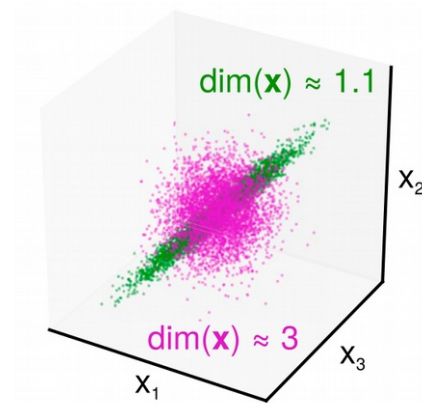
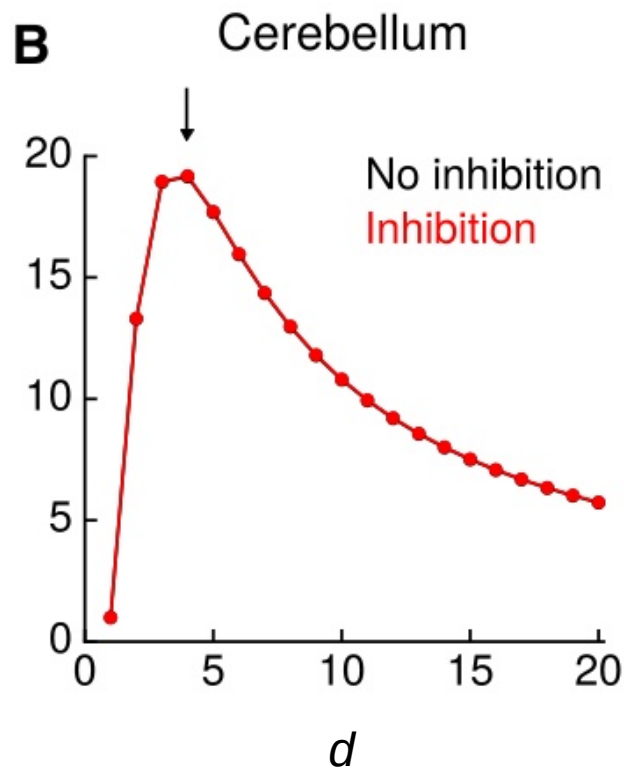
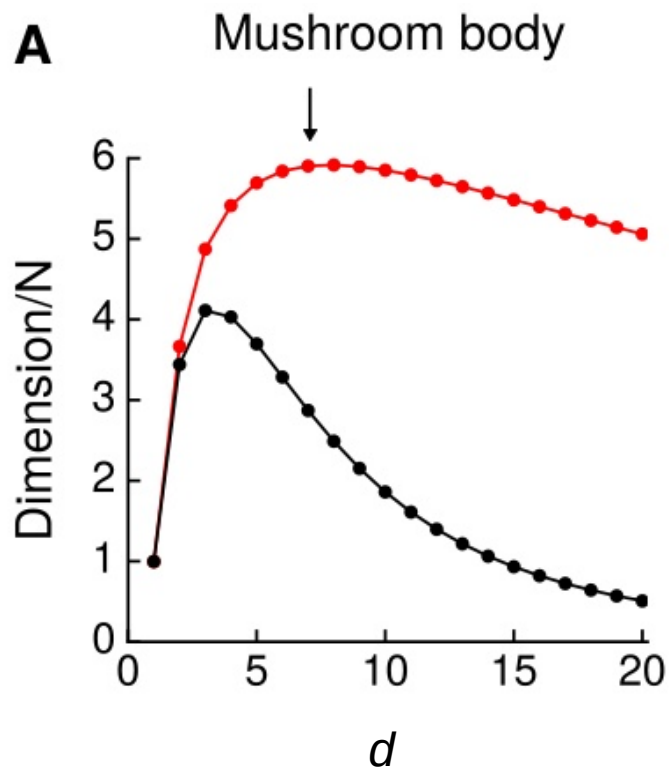
Cerebellum
 $d = 4$



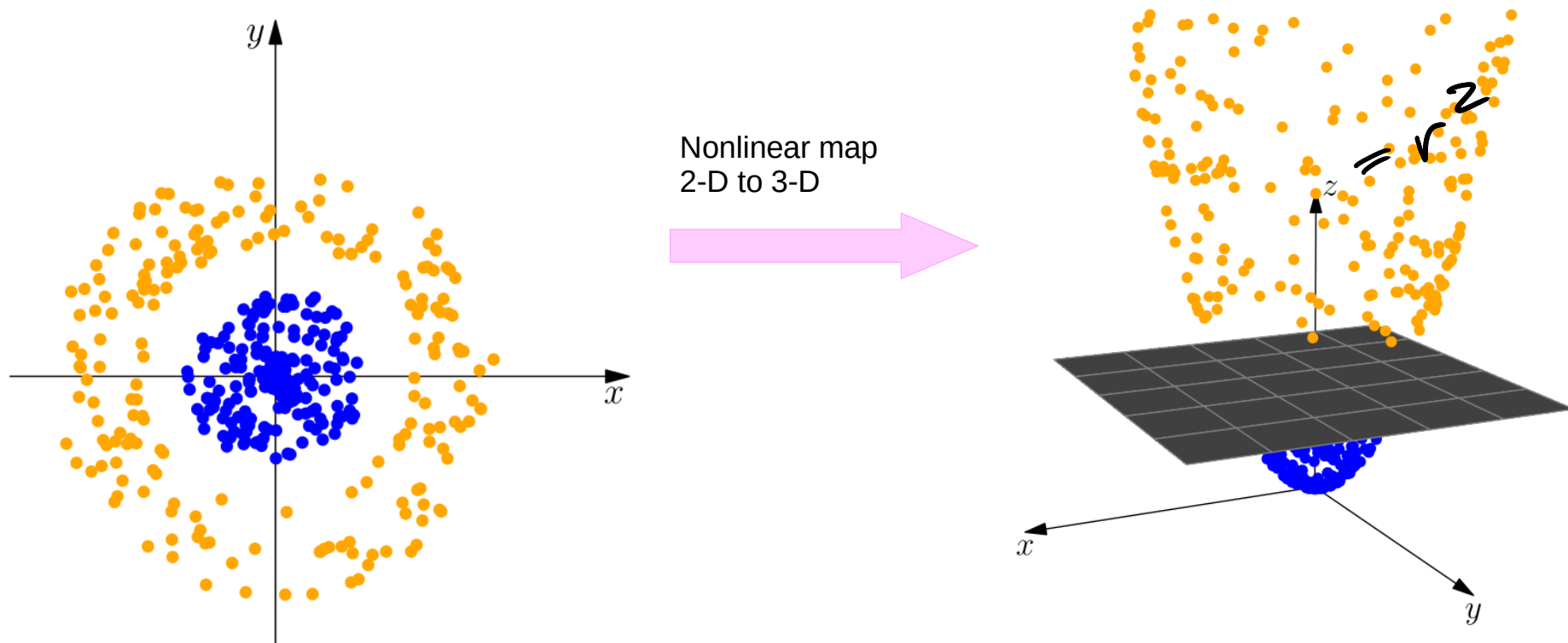
Litwin-Kumar, Harris, Axel, Sompolinsky, Abbott

Sparsity under constraints = max dimensionality

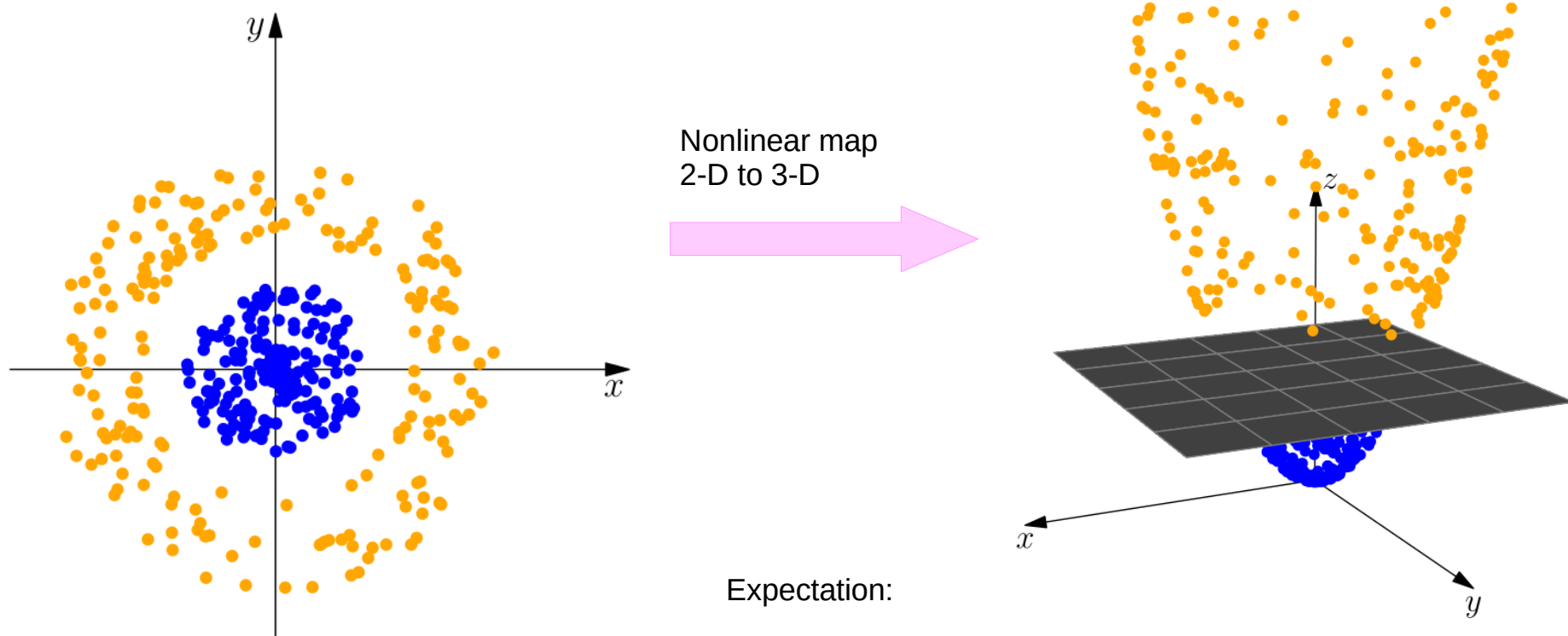
Arrows = avg degree observed in brains



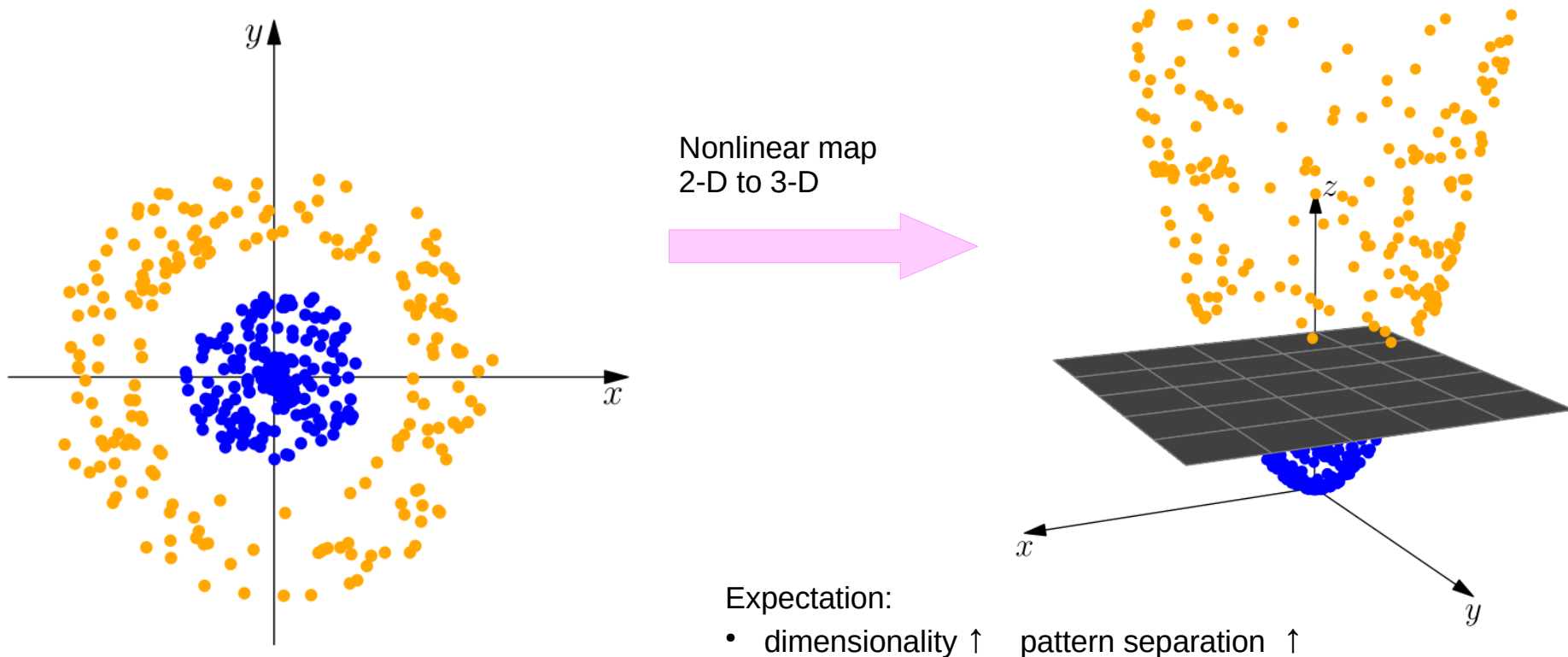
Dimensionality can help **or hurt** learning



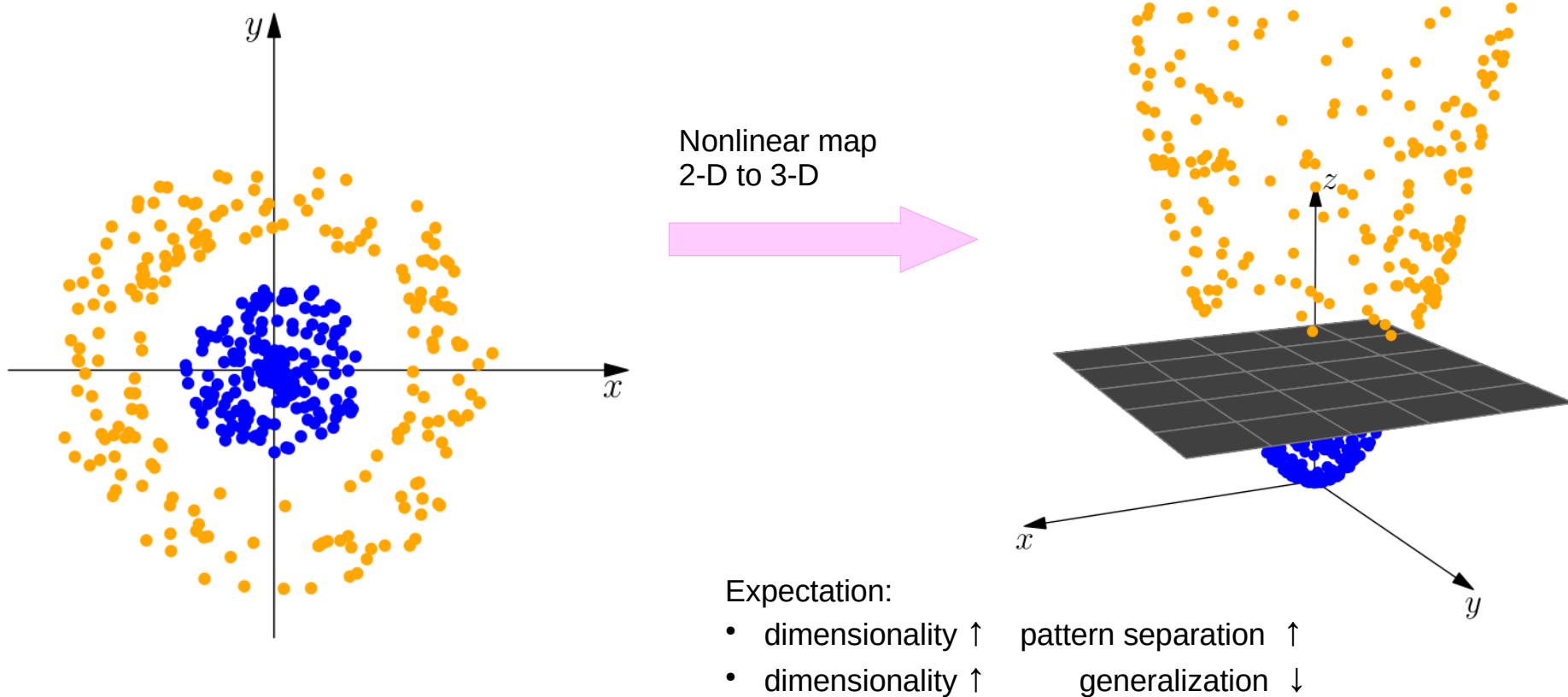
Dimensionality can help **or hurt** learning



Dimensionality can help **or hurt** learning



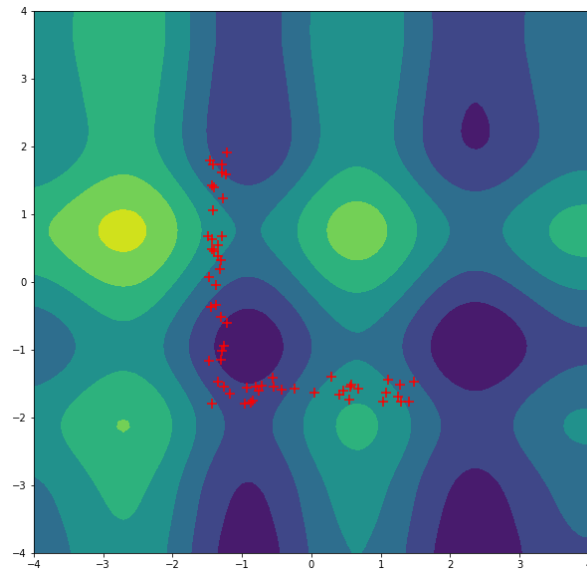
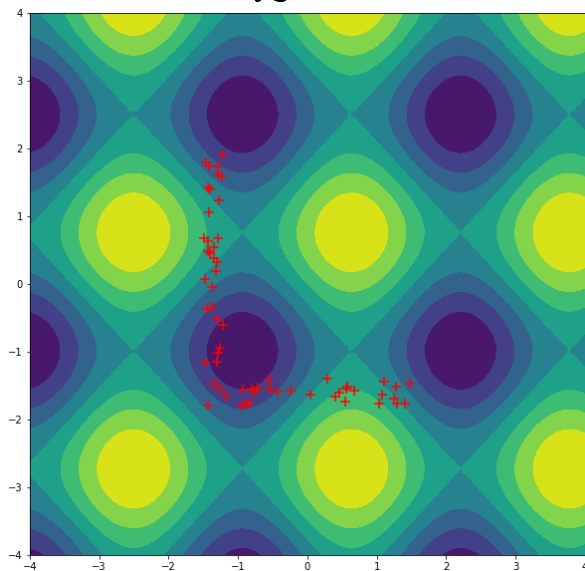
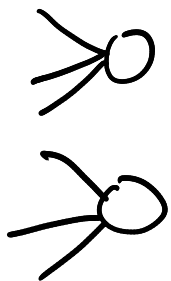
Dimensionality can help **or hurt** learning



Kernels: *sparse networks* = **additive** functions

Additive functions are constrained, hence low-dimensional (Stone, 1985 & '86)

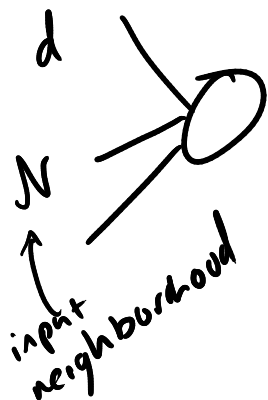
$$d = 3: \quad \underset{\substack{= \\ \ell\text{-dim}}}{f(\mathbf{x})} = f_1(x_1, x_3, x_4) + f_2(x_1, x_4, x_{11}) + \dots$$





Sparse network kernels

$$k(x, x') = \mathbb{E} [\phi(x) \phi(x')] = \mathbb{E} [\mathbb{E} [\phi(x_N) \phi(x'_N) | \mathcal{N}]]$$



$$d=2 \quad \binom{4}{2}$$

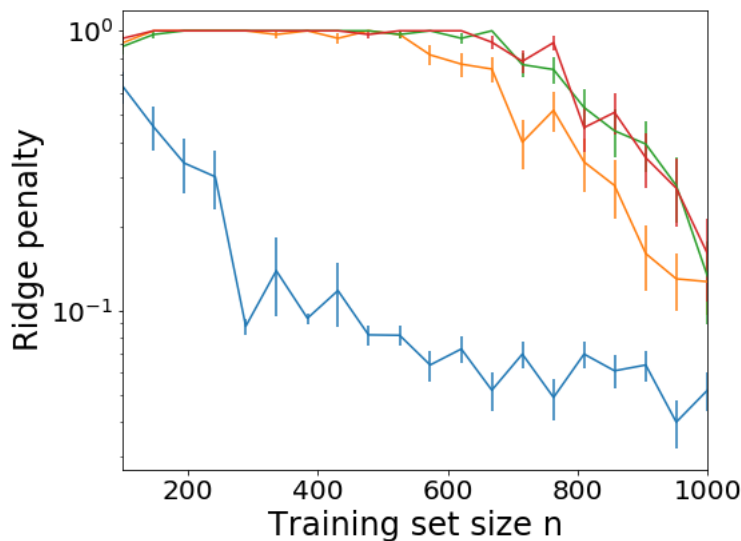
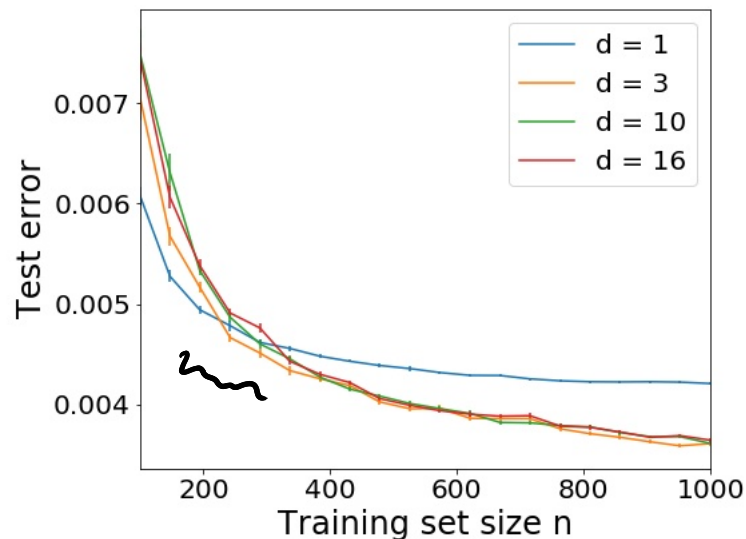
$$k_d^{\text{reg}}(x, x') = \frac{1}{\binom{l}{d}} \sum_{\substack{\mathcal{N}: |\mathcal{N}|=d \\ \mathcal{N} \in \mathcal{C}^l}} k_d(x_N, x'_N)$$

additive model
order d

$$d \sim \mathcal{D}$$

$$\mathcal{H} = \mathcal{H}_1 \oplus \mathcal{H}_2 \oplus \dots$$

Simulations confirm sparsity advantage



Target function:
Random linear + degree 3 polynomial

Overview

- Context: structure, randomness, & neuroscience
- Kernel theory:
 - Mathematical framework
 - Network sparsity → additive functions
 - **Tuning curves** → **basis change**
- Conclusions

“Random features for structured input”

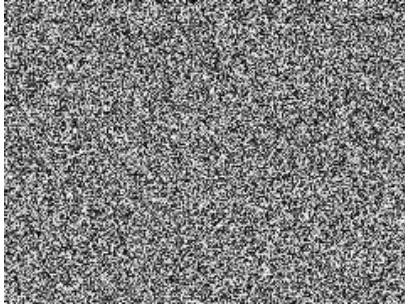


Biraj Pandey
NSF Grad Fellow
Applied Math

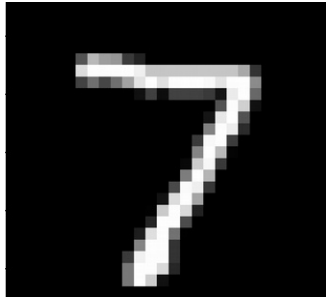


Bing Brunton
Biology

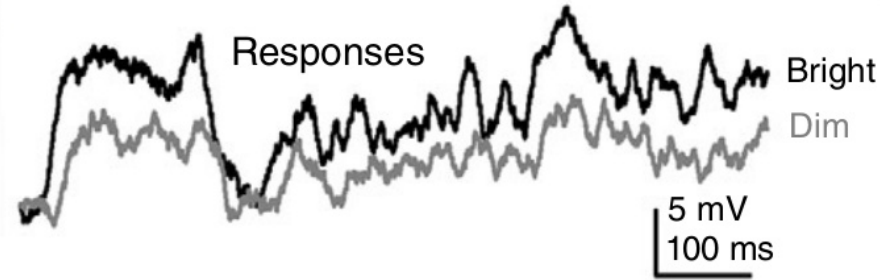
Defining “structured input”



white noise



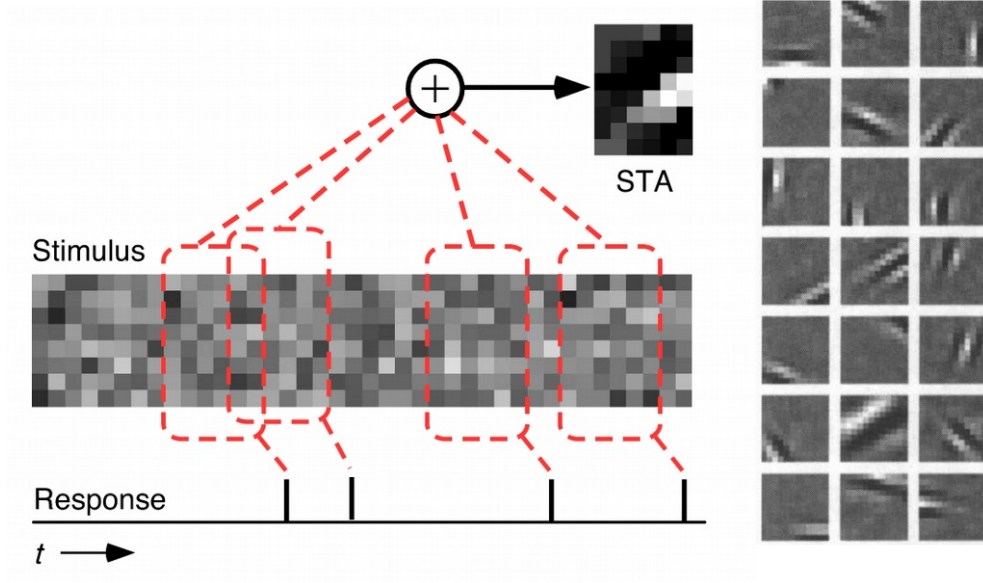
MNIST



locust photoreceptor, natural stimulus

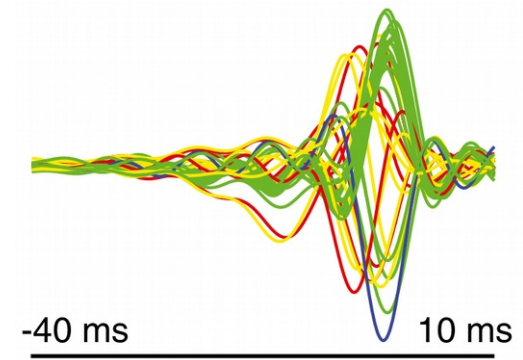
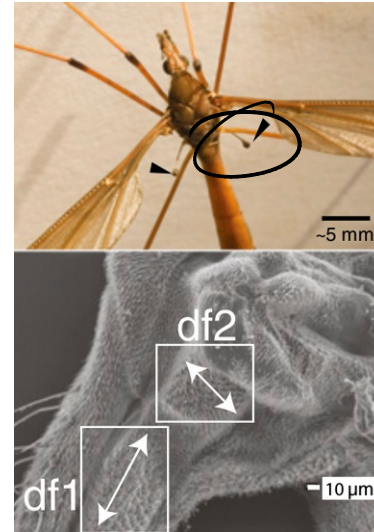
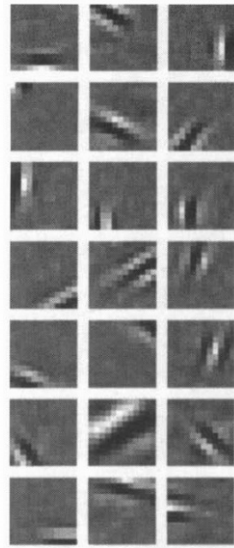
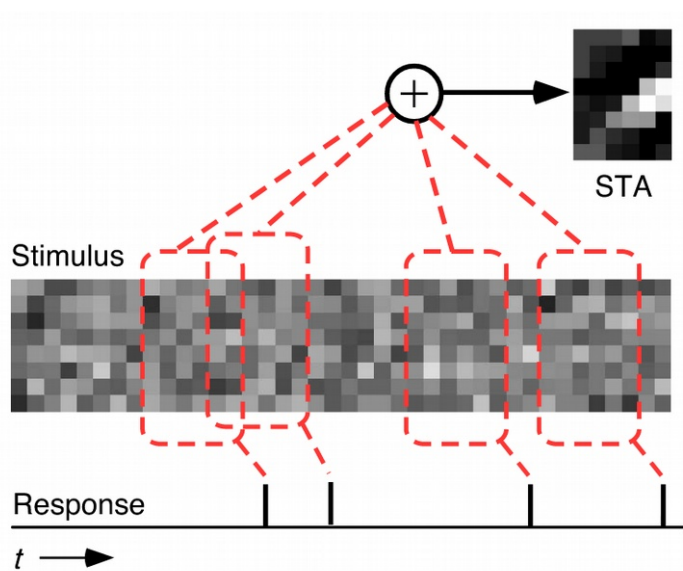
Tuning curves occurring in nature

Stimulus that best drives a neuron, i.e. its **receptive field**



Tuning curves occurring in nature

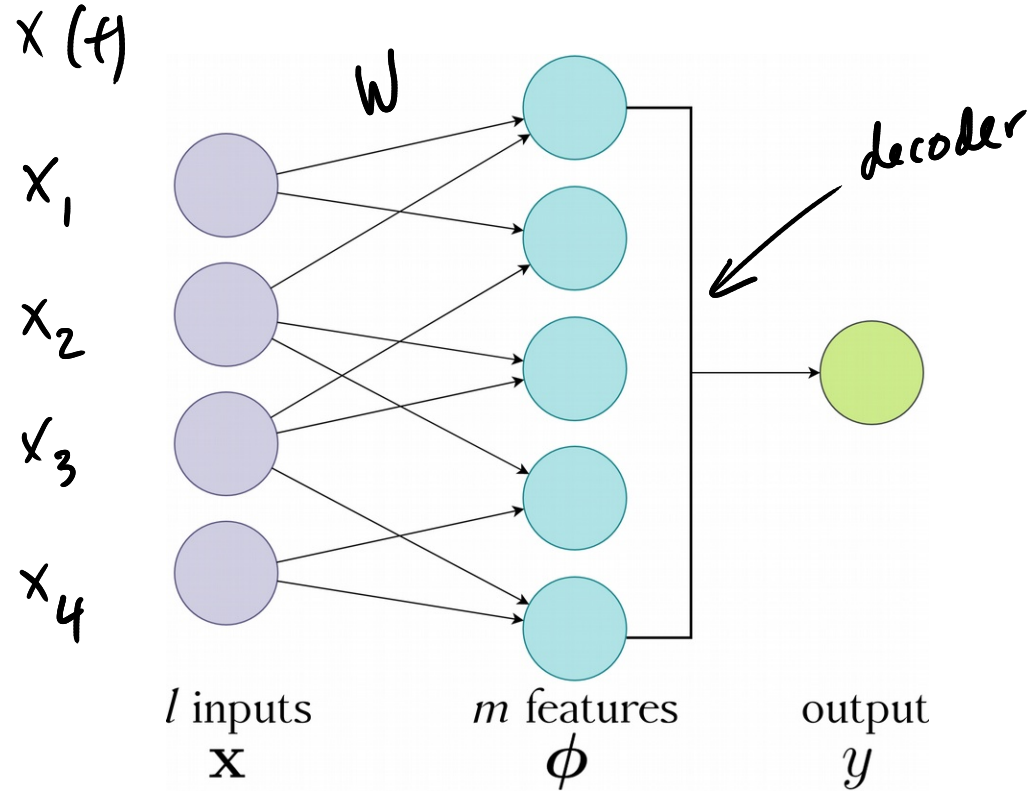
Stimulus that best drives a neuron, i.e. its **receptive field**



Linear-nonlinear model neurons



weights onto inputs
= linear filters

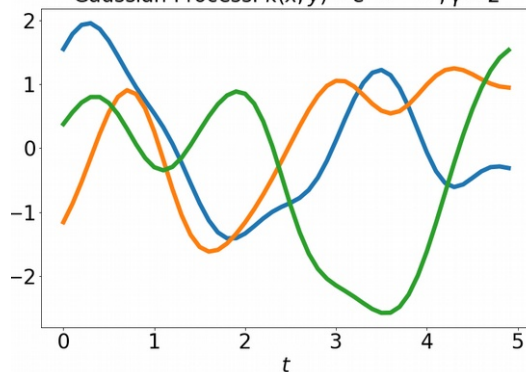




Theory of random tuning curves

Model **tuning curves** as random function drawn from a Gaussian process (GP)

Gaussian Process: $k(x, y) = e^{-\gamma|x - y|^2}$; $\gamma = 2$



$$w(t) \quad \mathbb{E}[w(t)] = 0$$

$$C(t, t') = \mathbb{E}[w(t) w(t')]$$

Mercer's theorem: $w(t) = \sum_{i=1}^{\infty} \sqrt{\lambda_i} G_i \psi_i(t) \leftarrow \text{eigenfunctions}$

$$G_i \sim \mathcal{N}(0, 1)$$

Discrete

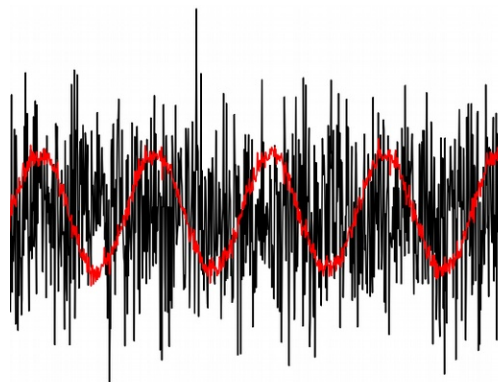
$$w = U D g$$

↑
orthogonal basis

↙ diagonal weights

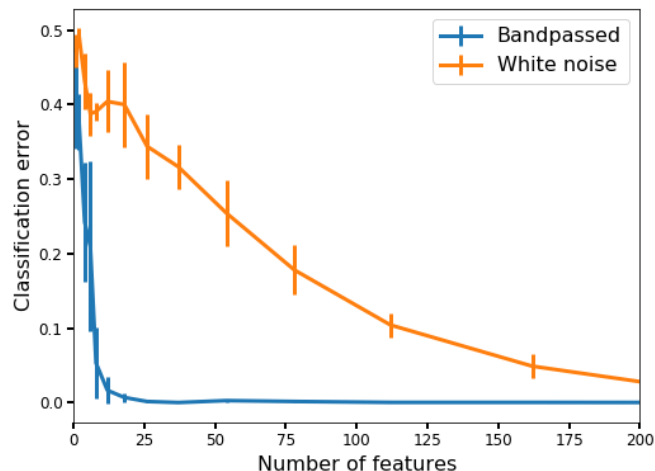
$$w^T x = g^T \underbrace{D U^T x}_{\tilde{x}} = g^T \tilde{x} \quad \text{new basis}$$

Example: frequency detection in timeseries

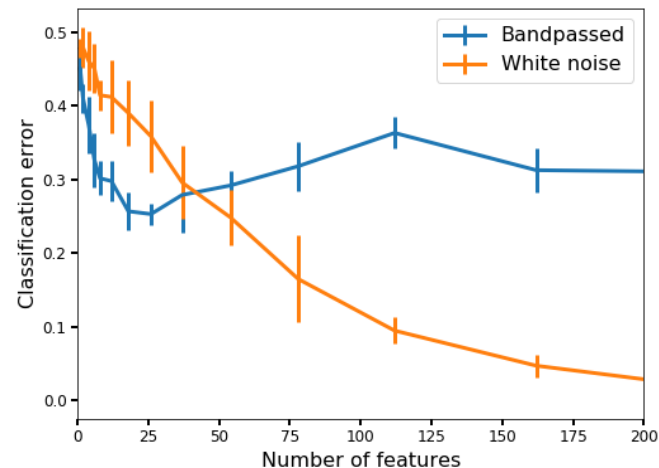


50 Hz signal
noise

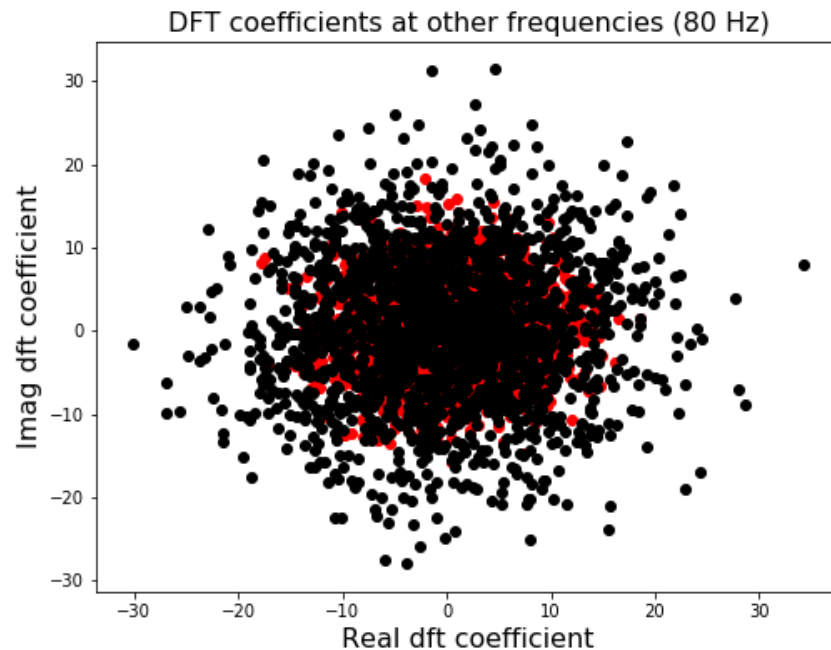
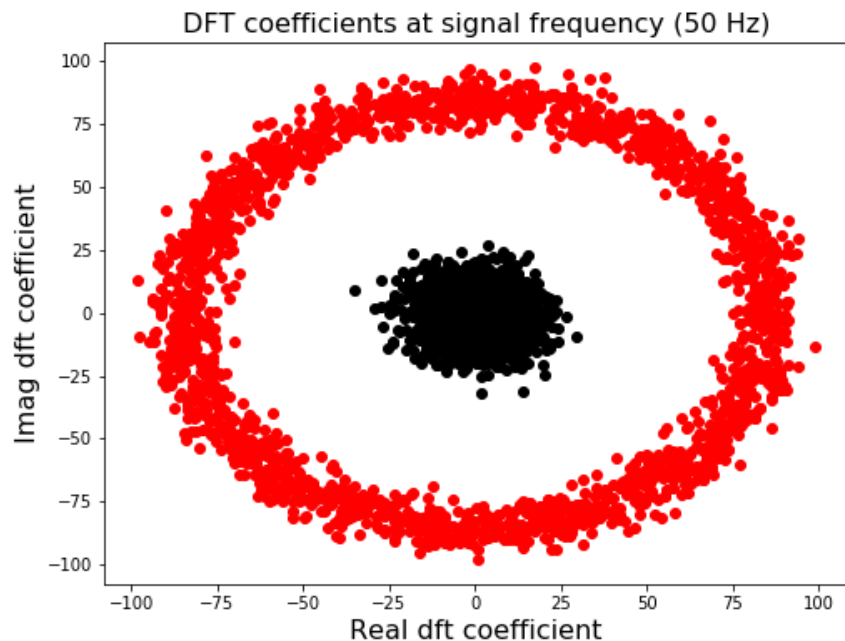
20 – 80 Hz
bandpass



70 – 120 Hz
bandpass

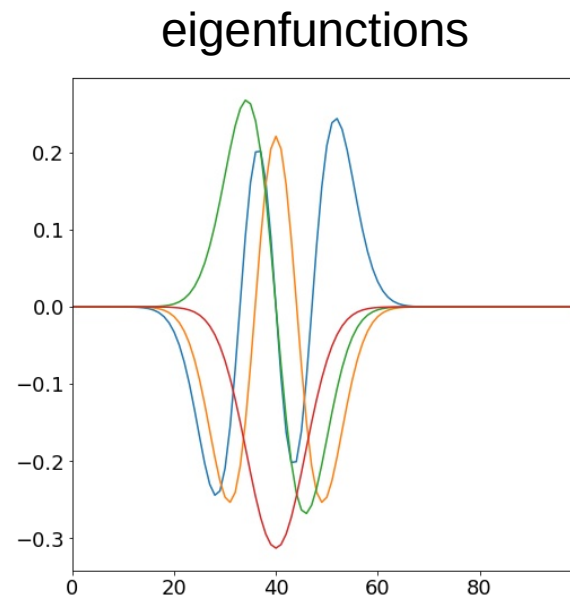
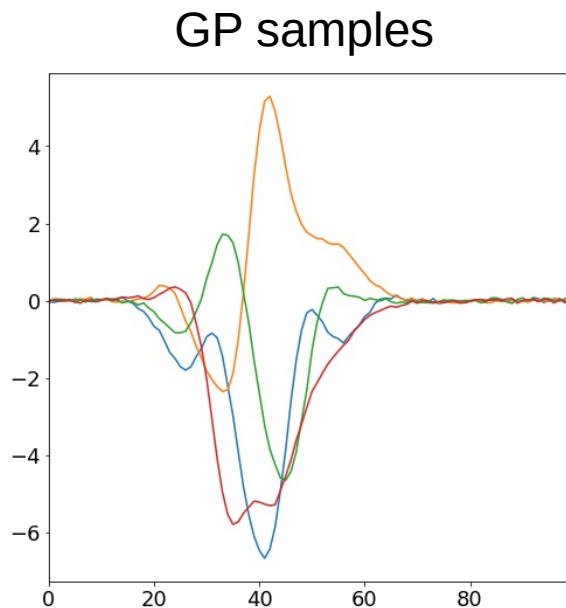
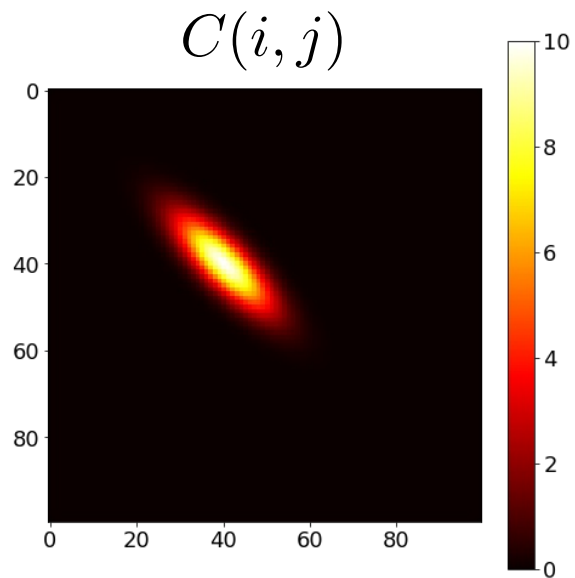


Fourier analysis of this test case



need a kernel which is good for
data in Fourier basis

Example: wavelet basis via non-stationary GPs



Overview

- Context: structure, randomness, & neuroscience
- Kernel theory:
 - Mathematical framework
 - Network sparsity \rightarrow additive functions
 - Tuning curves \rightarrow basis change
- **Conclusions**

Conclusions

Conclusions

- Kernel representation of random networks is powerful
 - statistical learning theory: **Which functions are easy to learn?**
 - neuro can learn from ML (dimensionality) and vice-versa (structure)

Conclusions

- Kernel representation of random networks is powerful
 - statistical learning theory: **Which functions are easy to learn?**
 - neuro can learn from ML (dimensionality) and vice-versa (structure)
- Sparse network leads to additive kernels
 - classical way to model “big data”

Conclusions

- Kernel representation of random networks is powerful
 - statistical learning theory: **Which functions are easy to learn?**
 - neuro can learn from ML (dimensionality) and vice-versa (structure)
- Sparse network leads to additive kernels
 - classical way to model “big data”
- Random tuning curves could explain variability seen
 - despite randomness, may represent inputs in Fourier/wavelet bases

Conclusions

- Kernel representation of random networks is powerful
 - statistical learning theory: **Which functions are easy to learn?**
 - neuro can learn from ML (dimensionality) and vice-versa (structure)
- Sparse network leads to additive kernels
 - classical way to model “big data”
- Random tuning curves could explain variability seen
 - despite randomness, may represent inputs in Fourier/wavelet bases
- Many future directions
 - feedback, temporal dynamics, unsupervised settings

Thank you!

- Funding:  **Washington Research**
F O U N D A T I O N
- Collaborators:
 - Biraj Pandey, Bing Brunton
 - Marjorie Xie, Ashok Litwin-Kumar, Larry Abbott, Richard Axel, Haim Sompolsky
- Thanks to Raj Rao, Kamesh Krishnamurthy, Yian Ma, & Francis Bach for discussions

Meaning of dimensionality in statistical learning

- Eigenvalue decay of kernel matrix K that depends on
 - kernel function
 - distribution of data x

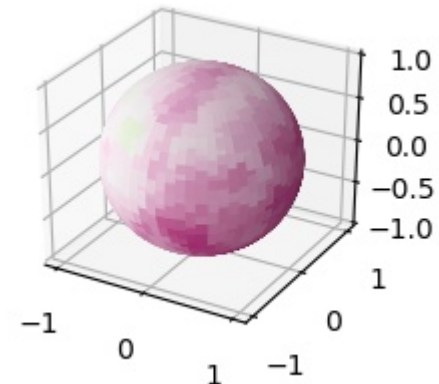
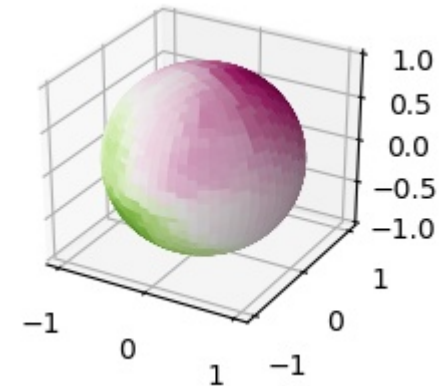
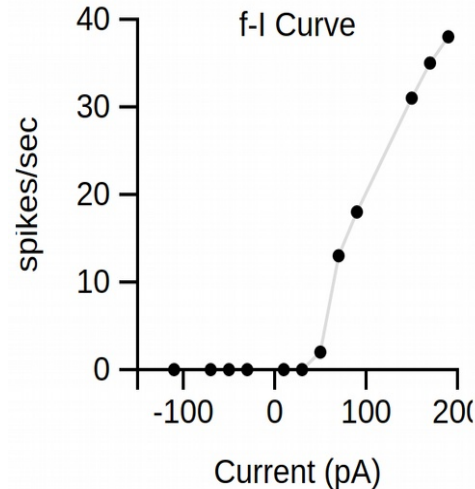
$$\mathbb{E}(y - \hat{y})^2 = (\text{bias})^2 + \underbrace{\text{variance}}_{\leq \frac{\sigma_y^2}{n} \mathbf{dim}}$$

* but not the
“participation ratio”

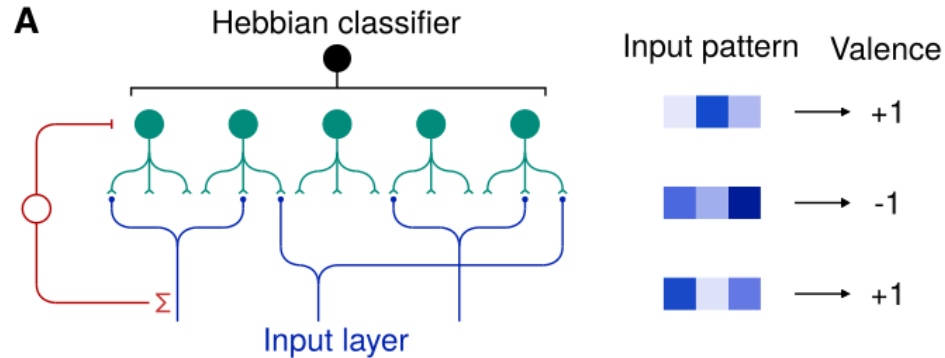
Kernels highlight importance of “preprocessing”

- Antennal lobe glomeruli provide
 - pooling of ORN inputs
 - divisive normalization
- RBFs on the unit sphere

$$\phi_i(\mathbf{x}) = h(\mathbf{w}_i^T \mathbf{x})$$



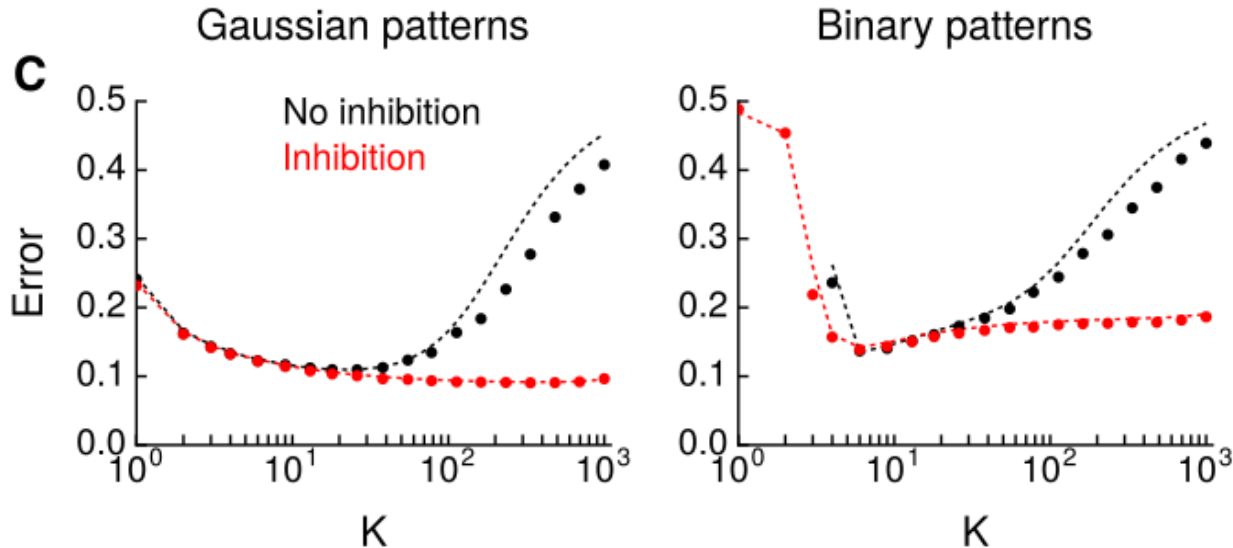
Sparsity can improve classification



Classifier must **remember & denoise**

Uncorrelated random

- Input patterns
- Binary valence
- Binary noise



$$\text{SNR} \approx \frac{\dim(\mathbf{m}) \cdot (1 - \Delta)^2}{P}$$

Noise amplification

Learning input-mixed weights most useful only in dense networks

