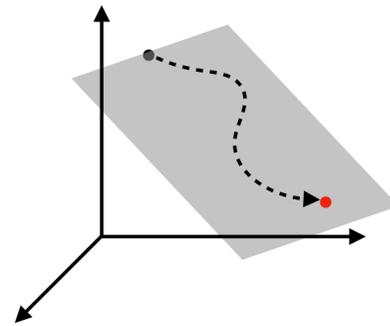


A geometric view of *intrinsic dimension* and *lottery subspaces*

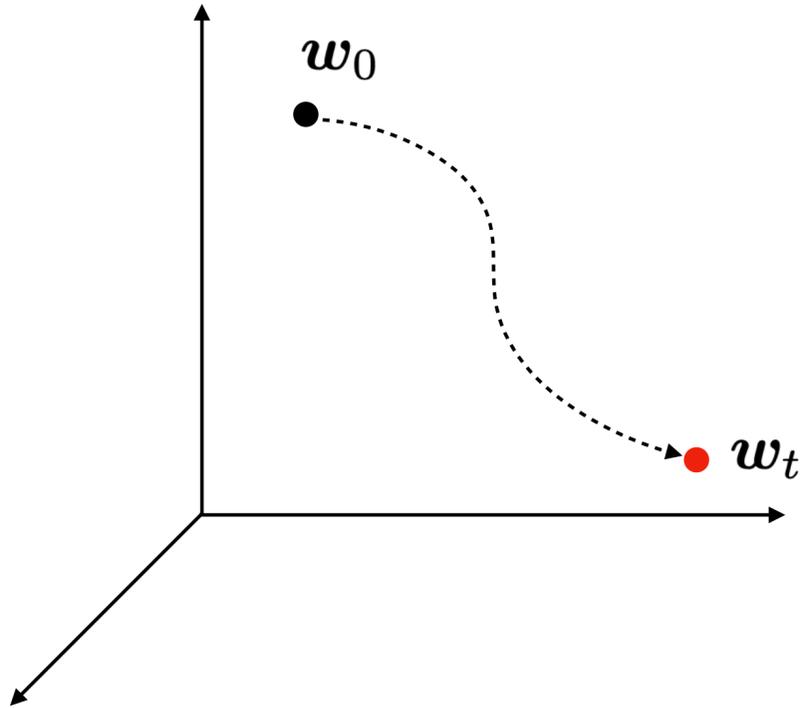


“Tony” Runzhe Yang

<https://runzhe-yang.science>

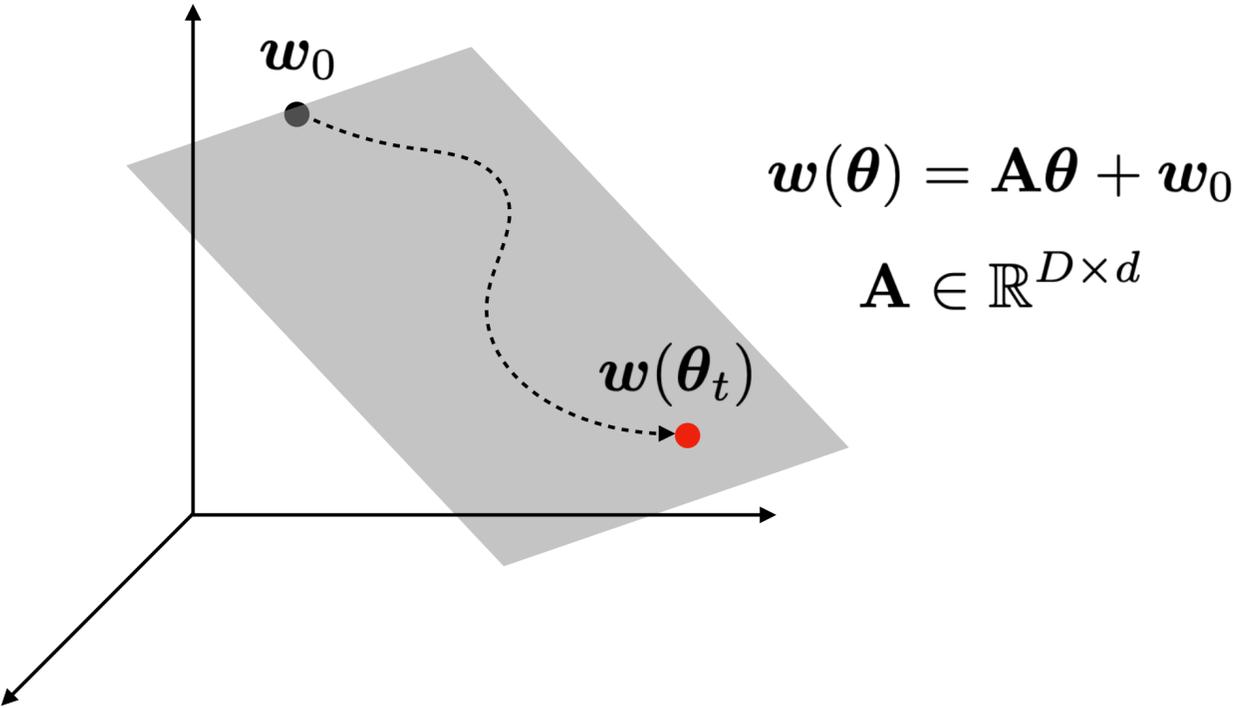
What happens if we limit the degrees of freedom of training a neural network?

Training a model with D parameters



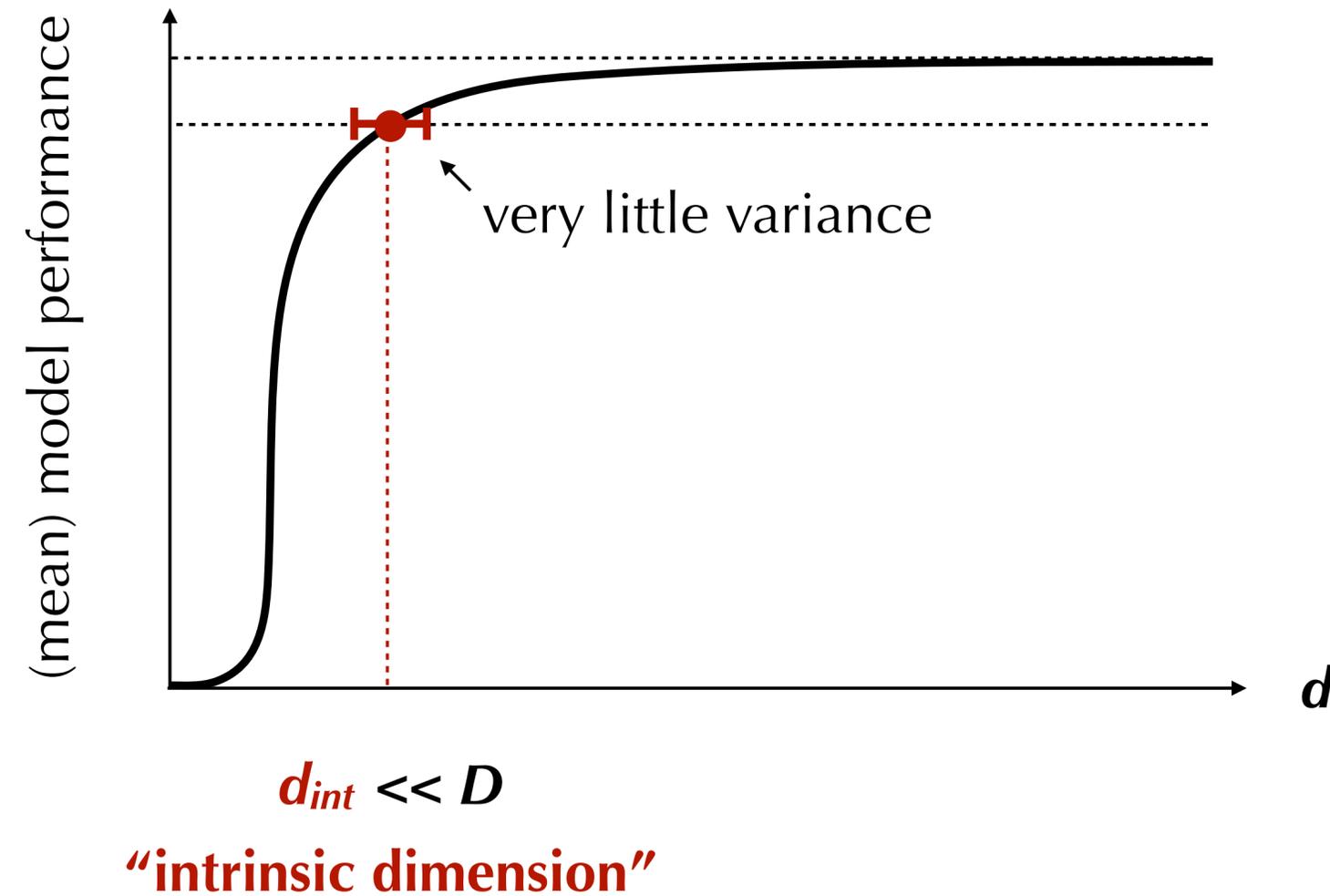
$$w_t \leftarrow w_{t-1} - \eta \nabla_w \mathcal{L}(w_{t-1})$$

Training a model with D parameters in a d -dimensional subspace



$$\theta_t \leftarrow \theta_{t-1} - \eta \nabla_{\theta} \mathcal{L}(w(\theta_{t-1}))$$

A small random subspace is enough to train a good model



The intrinsic dimension measures invariance of objective landscape

Fully connected networks, on MNIST

The intrinsic dimension is invariant to model size.

huge “solution set”

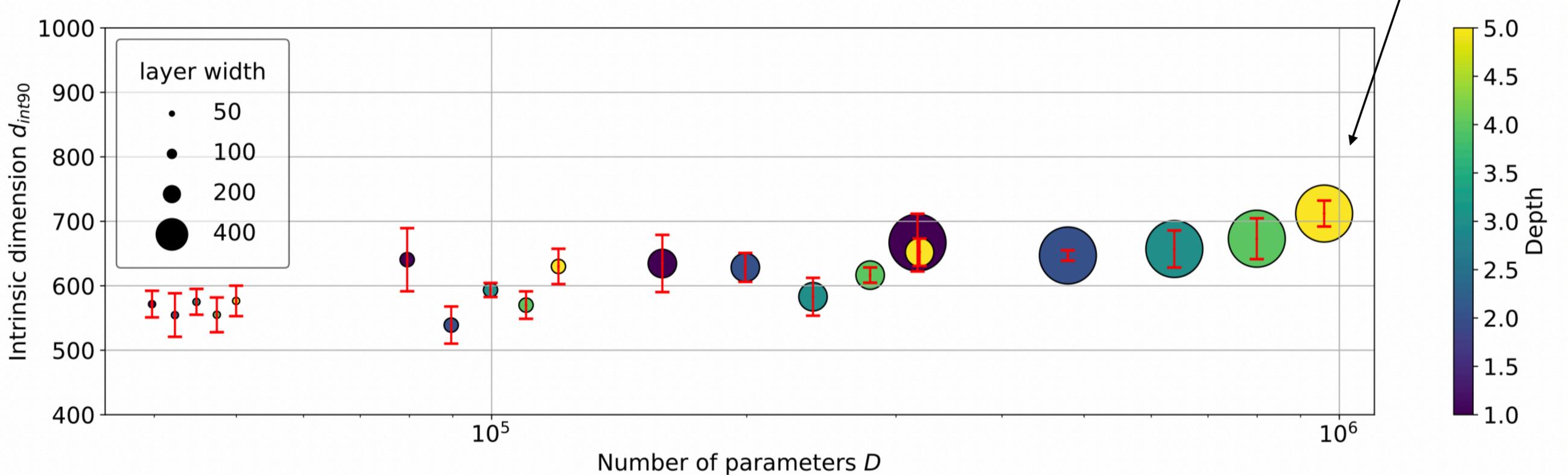
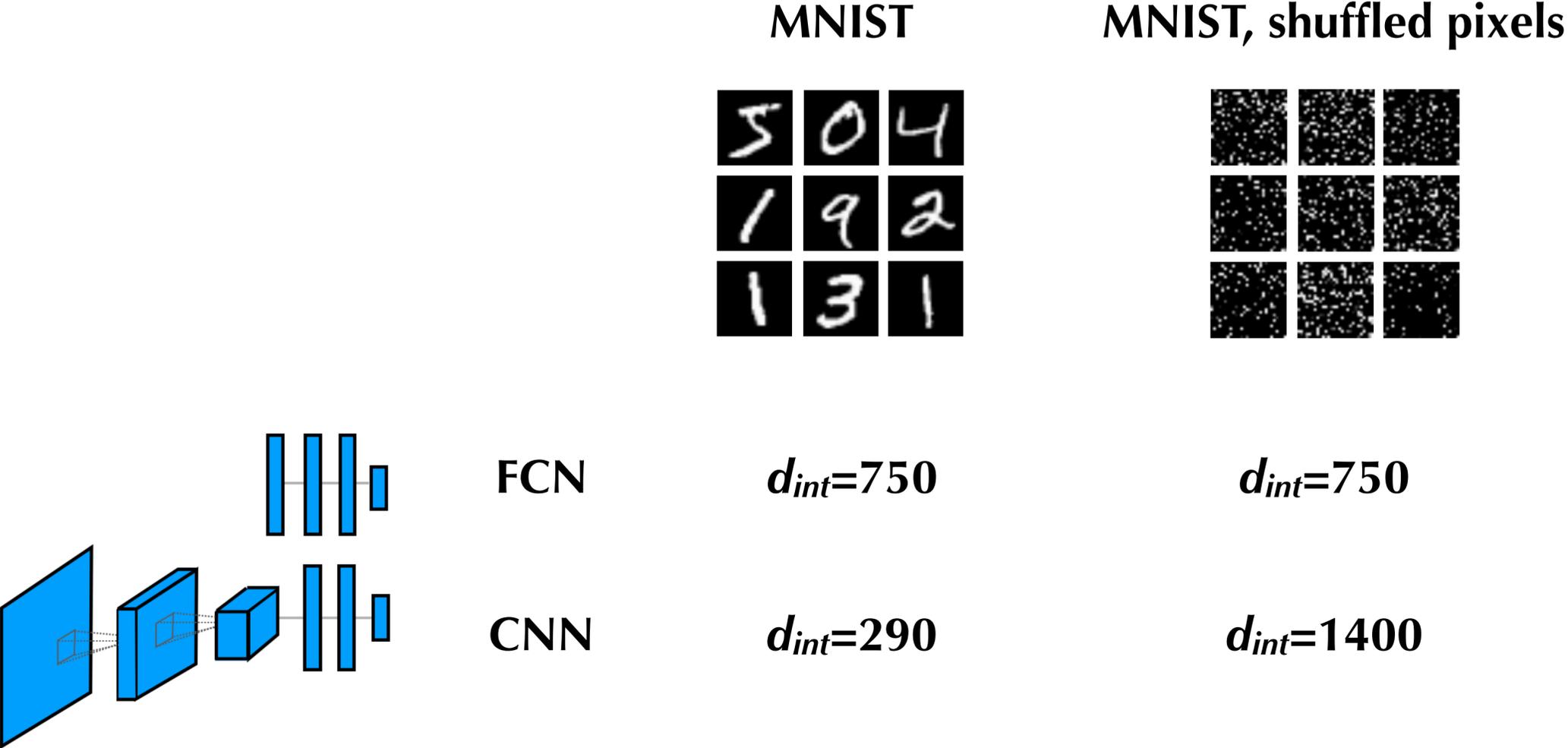


Fig S7 from Li et al. Measuring the intrinsic dimension of objective landscapes. ICLR'18

Intrinsic dimension approximates minimal description length (MDL)

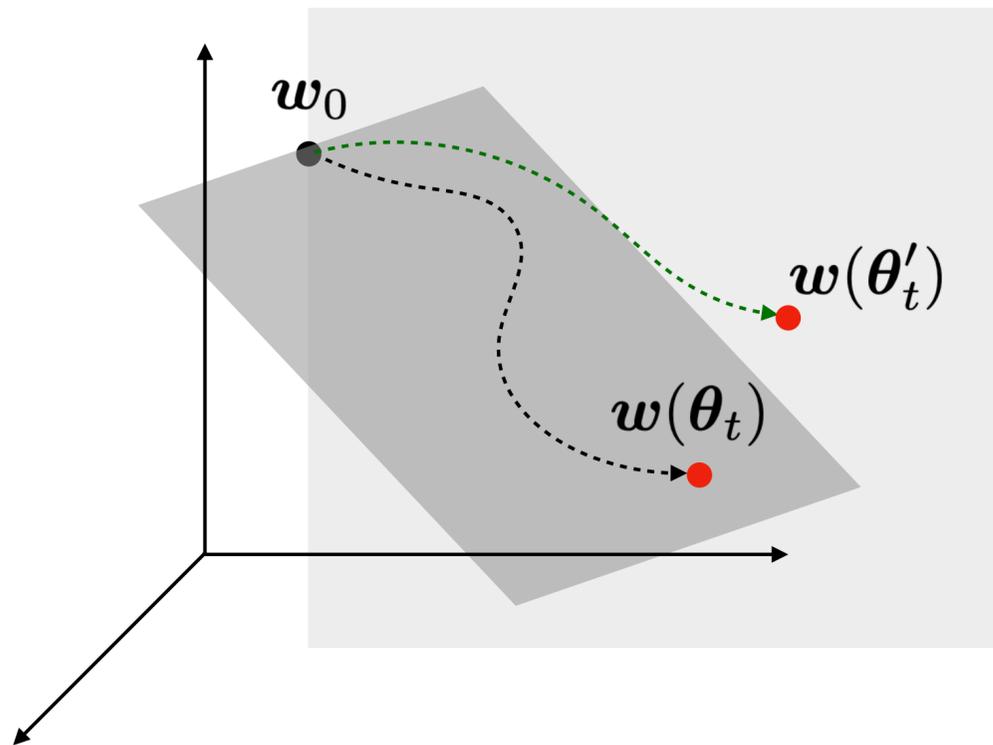


Modified from Blog: <https://eng.uber.com/intrinsic-dimension/>

Why do we always get a “lucky” subspaces when $d \geq d_{int}$?

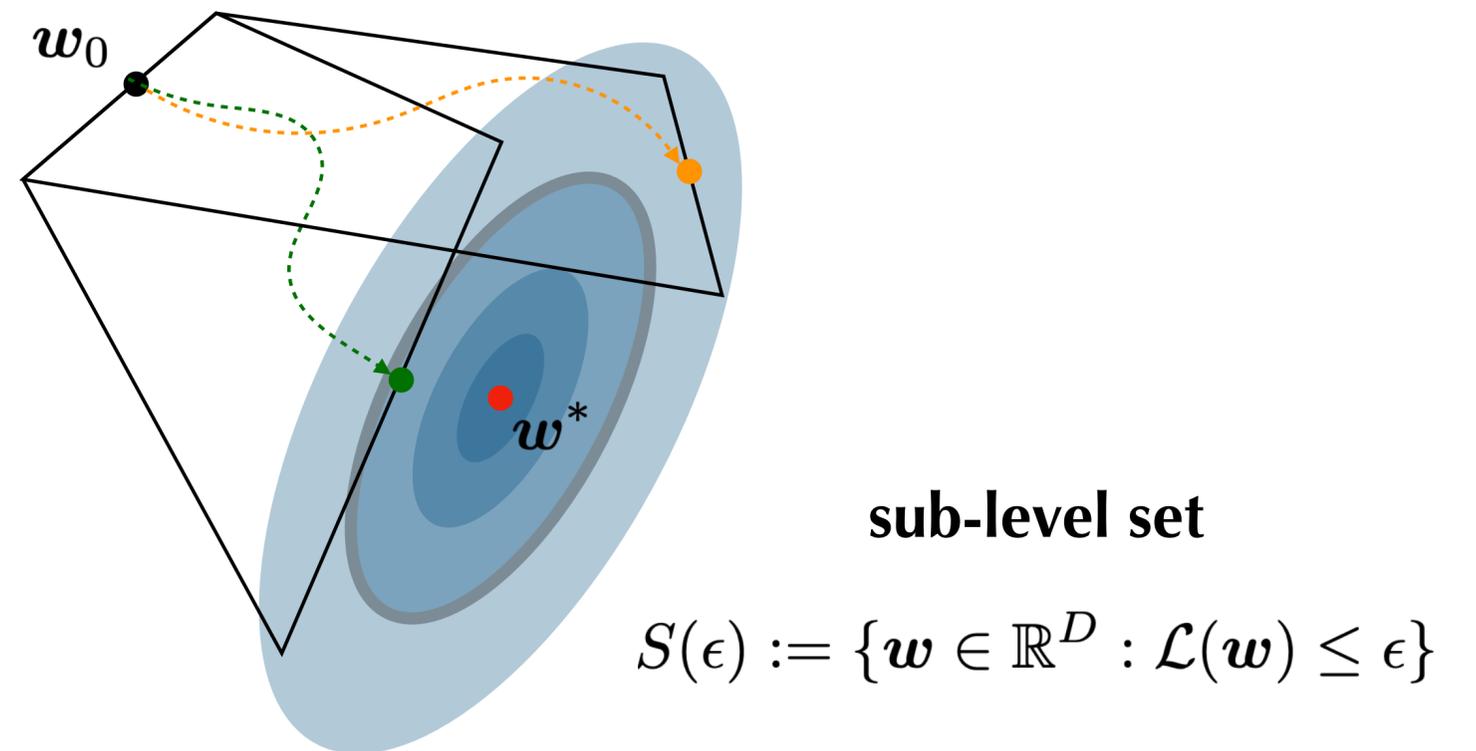
(and why always unlucky subspaces when $d < d_{int}$)

Training a model with D parameters
in a d -dimensional subspace



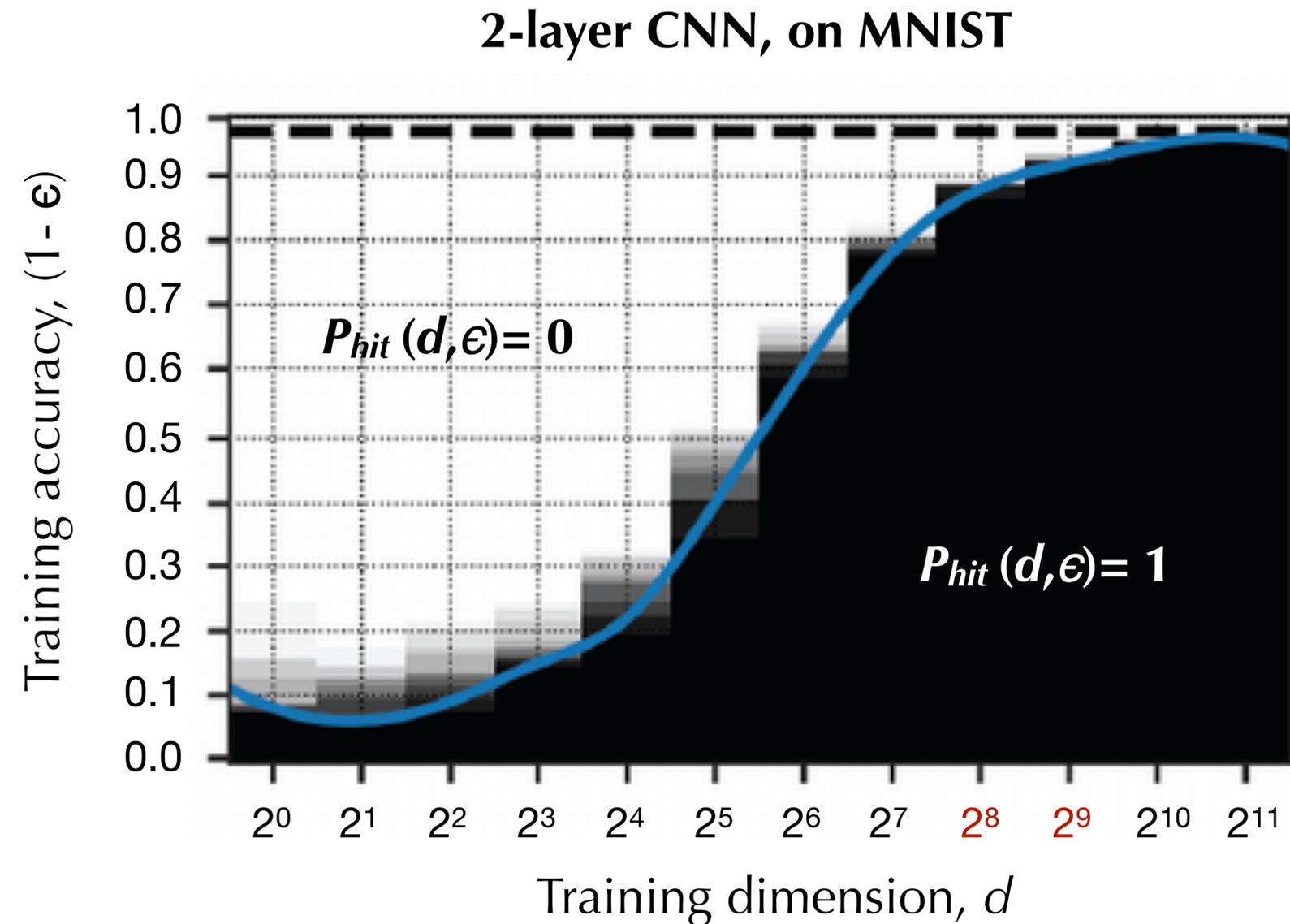
$$\theta_t \leftarrow \theta_{t-1} - \eta \nabla_{\theta} \mathcal{L}(w(\theta_{t-1}))$$

Success probability in hitting
a sub-level set



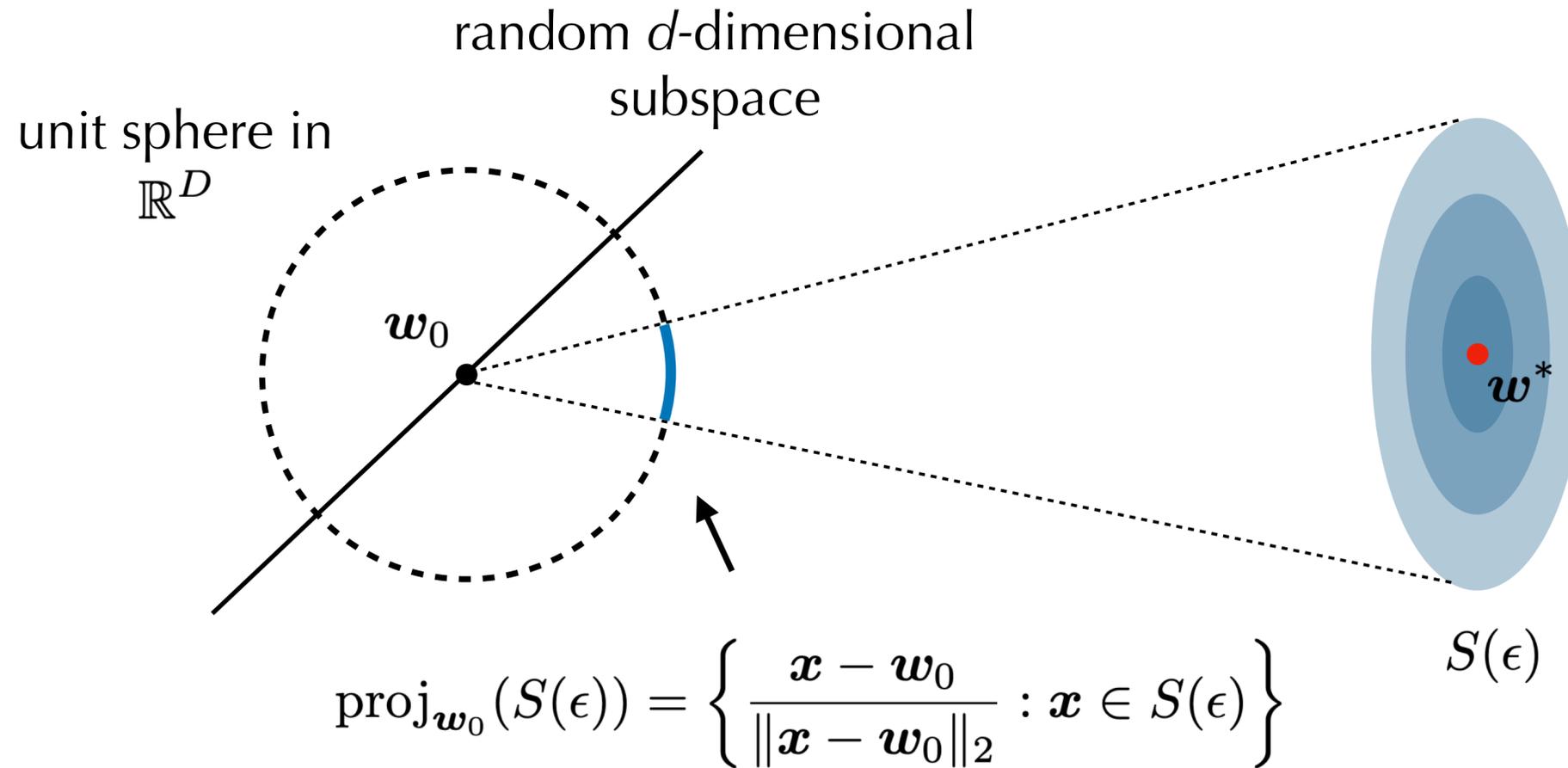
$$P_{hit}(d, \epsilon) = \mathbb{P}[S(\epsilon) \cap \{w_0 + \text{span}(\mathbf{A})\} \neq \emptyset]$$

A sharp phase transition in ϵ by d plane



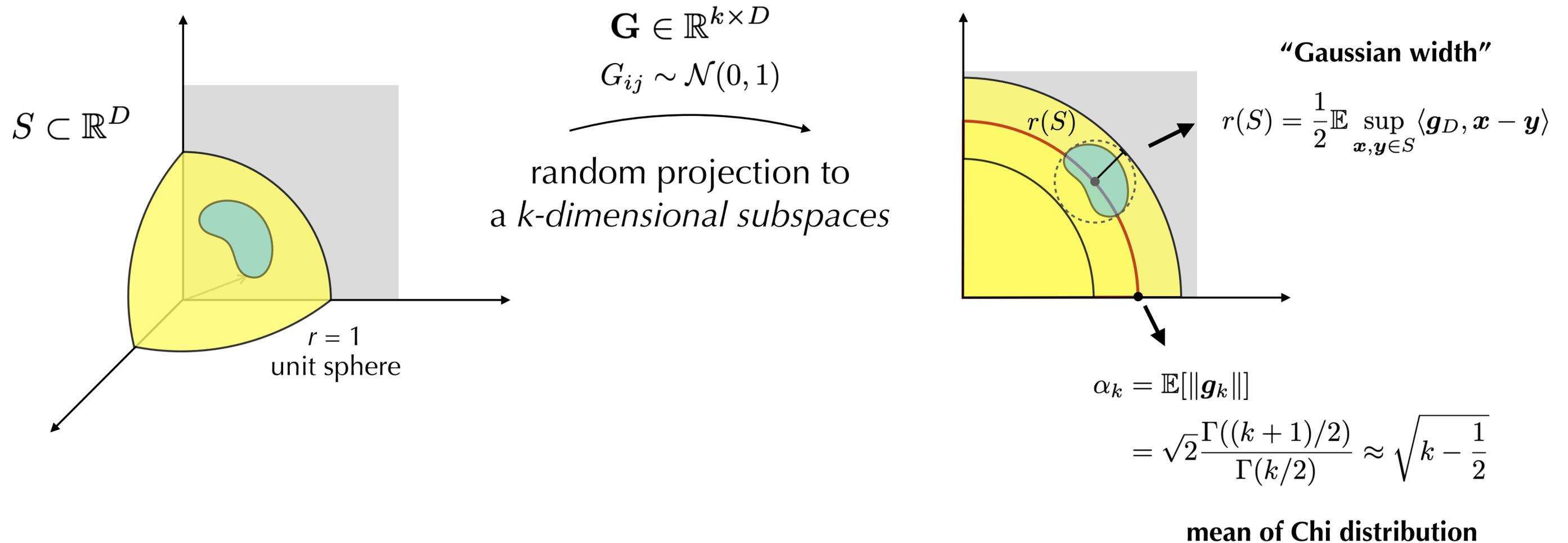
Modified from Fig.2 in Larsen et al. How many degrees of freedom do we need to train deep networks? A loss landscape perspective, ICLR'22

Equivalence to hitting the sub-level projection on a unit sphere



$$\Rightarrow P_{hit}(d, \epsilon) = \mathbb{P} [\text{proj}_{w_0}(S(\epsilon)) \cap \text{span}(\mathbf{A}) \neq \emptyset]$$

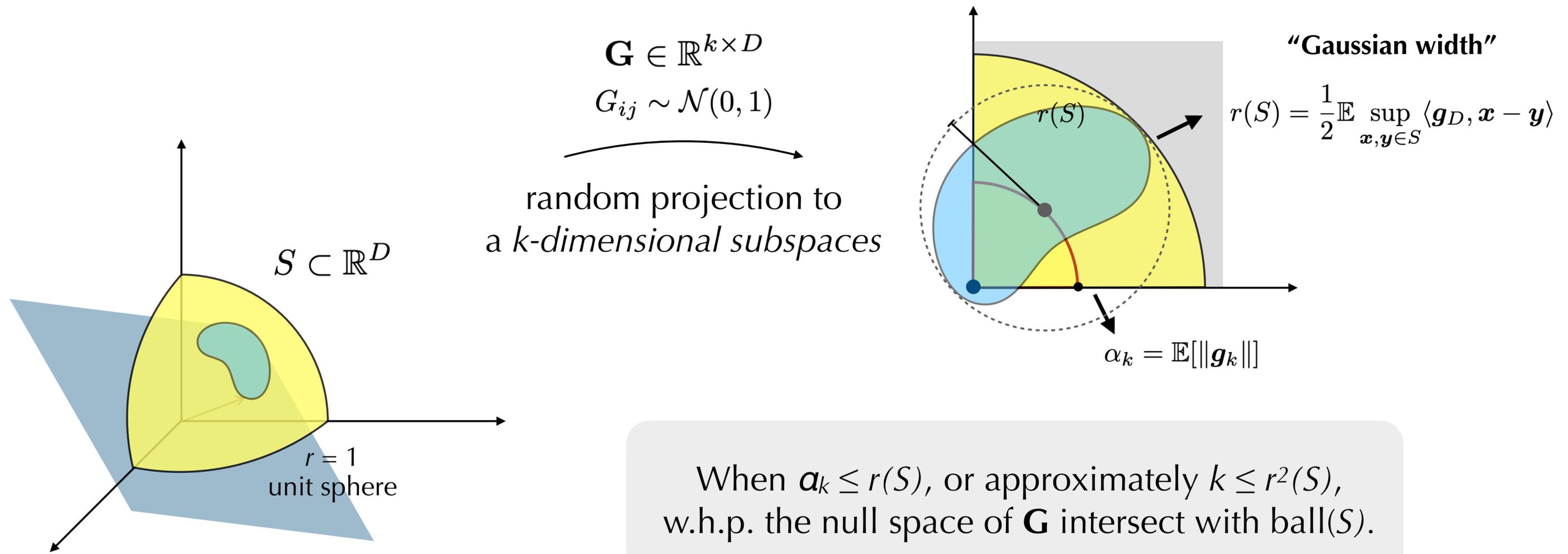
Picturing random projections in a high dimension space



$$\mathbb{E} \left[\min_{\mathbf{x} \in S} \|\mathbf{G}\mathbf{x}\|_2 \right] \geq \alpha_k - r(S)$$

$$\mathbb{E} \left[\max_{\mathbf{x} \in S} \|\mathbf{G}\mathbf{x}\|_2 \right] \leq \alpha_k + r(S)$$

High probability of intersection with the null space (when $\alpha_k \leq r(S)$)



When $\alpha_k \leq r(S)$, or approximately $k \leq r^2(S)$, w.h.p. the null space of \mathbf{G} intersect with ball(S).

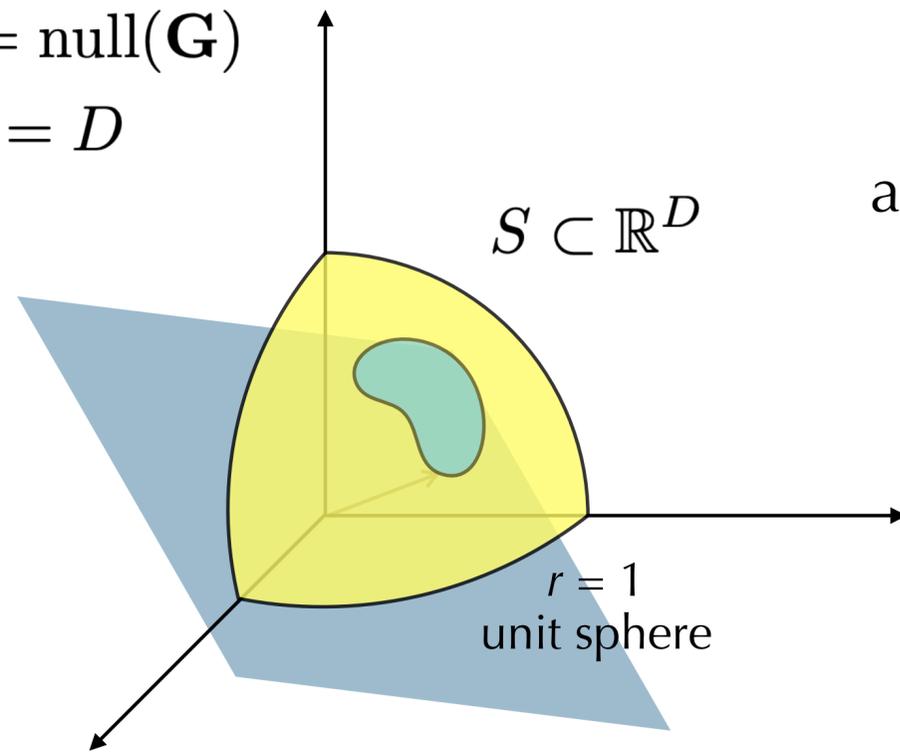
$\text{span}(\mathbf{A}) = \text{null}(\mathbf{G})$
 $d + k = D$

\Rightarrow

$d \geq D - \frac{r^2(\text{proj}_{\mathbf{w}_0}(S(\epsilon)))}{\text{"local angular dimension"}}$
 $P_{hit}(d, \epsilon) \geq 1 - \exp\{-m_{d, \epsilon}\}$

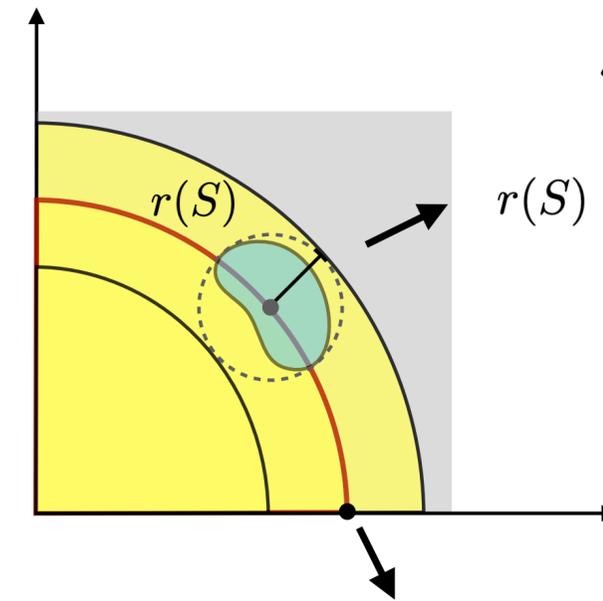
Gordon's escape from a mesh theorem (when $\alpha_k > r(S)$)

$\text{span}(\mathbf{A}) = \text{null}(\mathbf{G})$
 $d + k = D$



$\mathbf{G} \in \mathbb{R}^{k \times D}$
 $G_{ij} \sim \mathcal{N}(0, 1)$

random projection to a k -dimensional subspaces



"Gaussian width"

$r(S) = \frac{1}{2} \mathbb{E} \sup_{\mathbf{x}, \mathbf{y} \in S} \langle \mathbf{g}_D, \mathbf{x} - \mathbf{y} \rangle$

$\alpha_k = \mathbb{E}[\|\mathbf{g}_k\|]$

$\mathbb{P}[S \cap \{\text{span}(\mathbf{A})\} = \emptyset] = \mathbb{P} \left[\min_{\substack{\mathbf{x} \in S \\ \mathbf{y} \in \text{null}(\mathbf{G})}} \|\mathbf{x} - \mathbf{y}\|_2 > 0 \right]$

$\geq \lim_{\epsilon \rightarrow 0} \mathbb{P} \left[\min_{\substack{\mathbf{x} \in S \\ \mathbf{y} \in \text{null}(\mathbf{G})}} \|\mathbf{G}(\mathbf{x} - \mathbf{y})\|_2 \geq \epsilon(1 + \delta) \mathbb{E}\|\mathbf{G}\|_2 \right]$

$-\exp\{-\delta^2 \alpha_k^2 / 2\}$

$\geq 1 - \frac{7}{2} \exp \left\{ -\frac{1}{18} (\alpha_k - r(S))^2 \right\}$

\Rightarrow

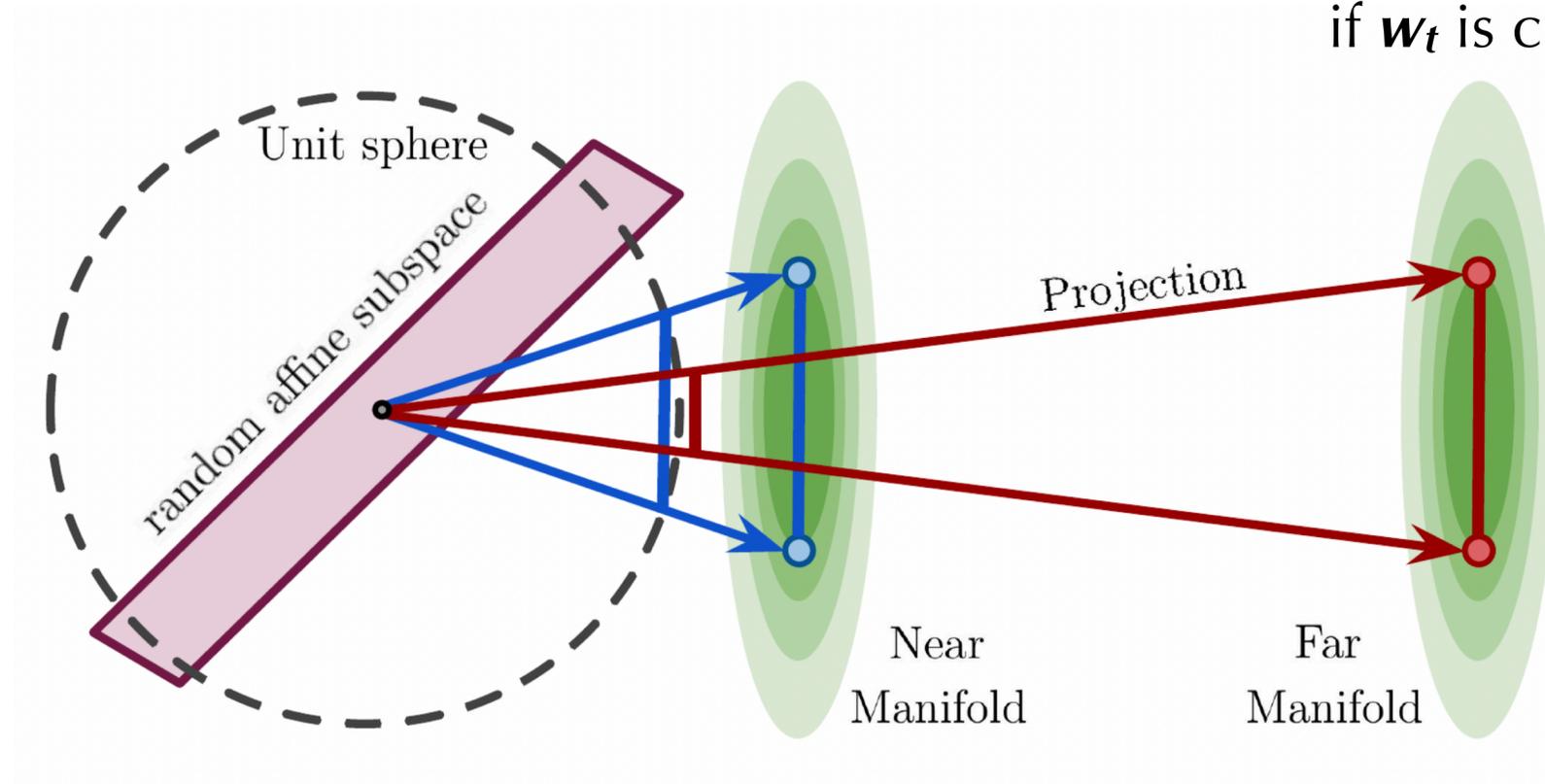
$d < D - \frac{r^2(\text{proj}_{\mathbf{w}_0}(S(\epsilon)))}{\text{"local angular dimension"}}$

$P_{hit}(d, \epsilon) \leq \exp\{-m_{1-d, 1-\epsilon}\}$

Intrinsic dimension as a function of initial point and sub-level set

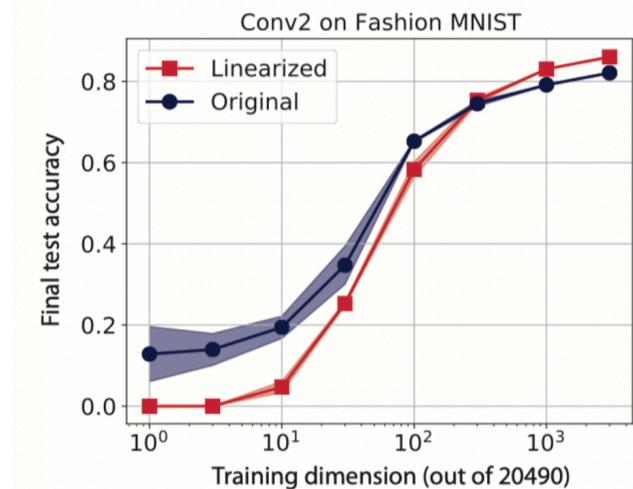
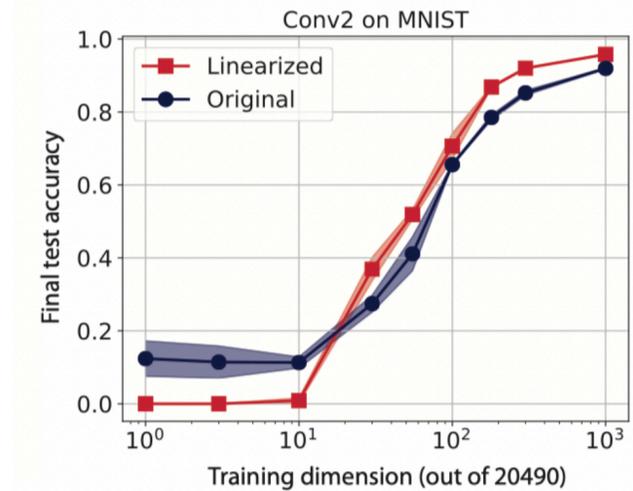
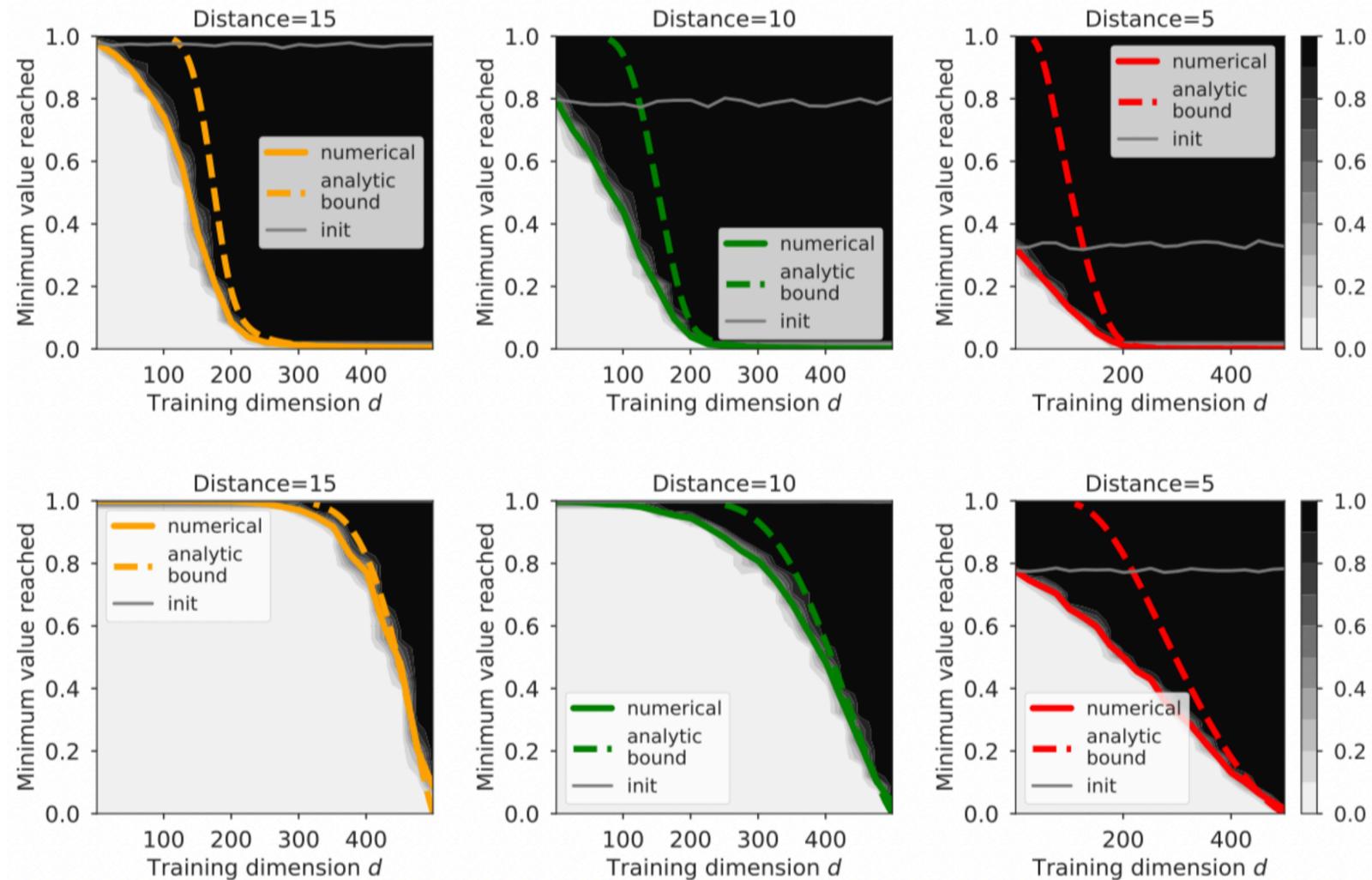
$$d_{int} = D - r^2(\text{proj}_{\mathbf{w}_0}(S(\epsilon)))$$

$\nearrow \epsilon \uparrow, d_{int} \downarrow$
 \searrow random subspace training after \mathbf{w}_t
if \mathbf{w}_t is closer to $S(\epsilon)$, $d_{int} \downarrow$



Modified from Fig.4 in Larsen et al. How many degrees of freedom do we need to train deep networks? A loss landscape perspective, ICLR'22

Deriving an analytical upper bound on intrinsic dimension



Optimizing $\mathcal{L}(w) = w^T H w$, w/ different init. An analytical upper bound on d_{int} .

Linearized networks using NTK.

Modified from Fig.5 and Fig.9 in Larsen et al. How many degrees of freedom do we need to train deep networks? A loss landscape perspective, ICLR'22

An analogy to lottery tickets (Frankle & Carbin, 2019)

	Training not used	Training used for init. only	Training used for pruning only	Training used for init. and pruning
Axis-aligned subspaces	Random weight pruning	Random weight pruning at step t	Lottery tickets, rewound to init.	Lottery tickets, rewound to step t
General subspaces	Random affine subspaces	Burn-in affine at step t	Lottery subspaces	Lottery subspaces at step t

→
smaller intrinsic dimension required

“Lottery subspace”: $w(\theta) = \mathbf{U}_d \theta + w_t$

↙
top d principal component of a training trajectory

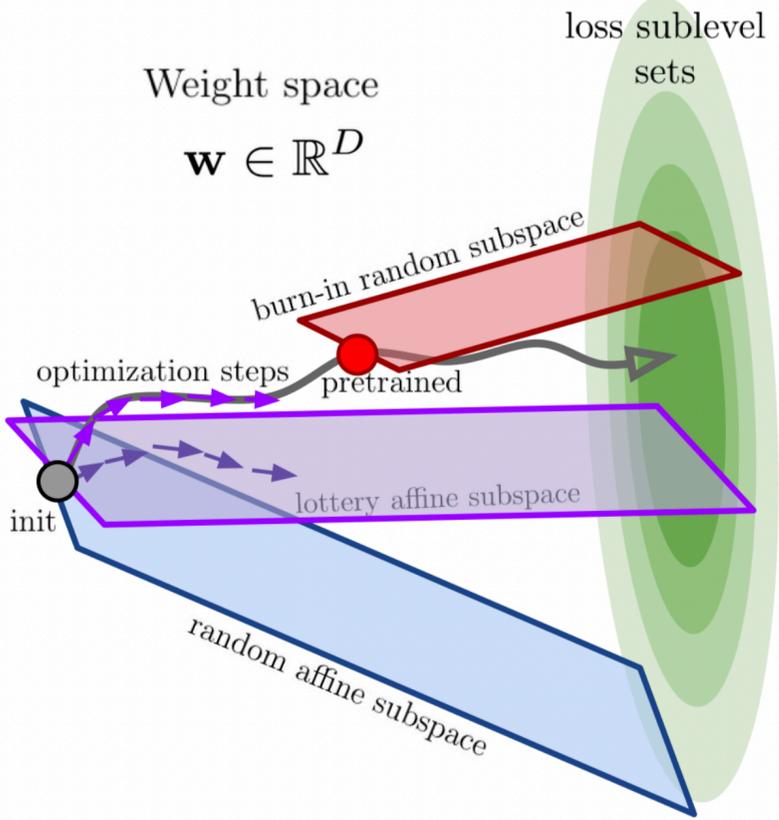


Fig.1 and Table 1 from Larsen et al. How many degrees of freedom do we need to train deep networks? A loss landscape perspective, ICLR'22

High “compression ratio” (D/d) compared with pruning

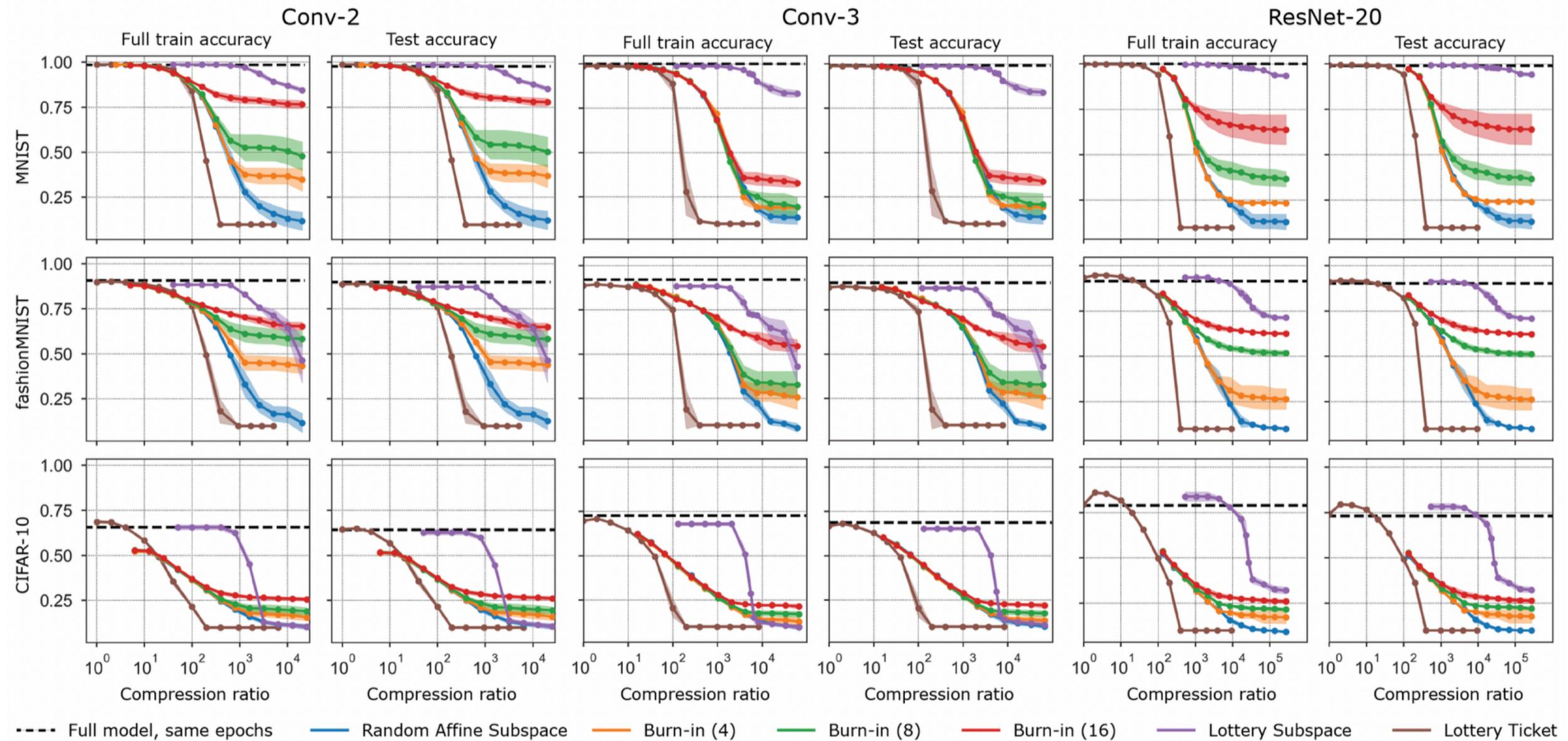
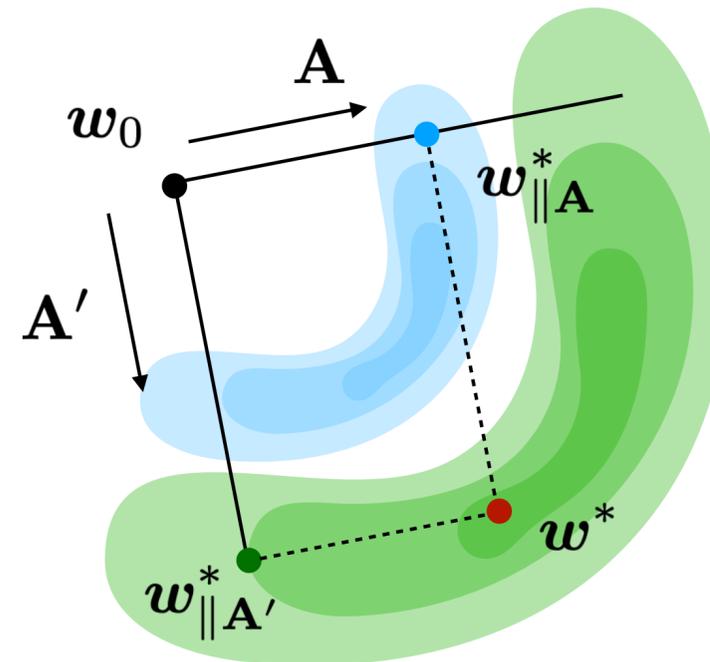


Fig.6 from Larsen et al. How many degrees of freedom do we need to train deep networks? A loss landscape perspective, ICLR'22

Discussion

- When $d \geq d_{int}$ and $d \ll D$, w.h.p. two random affine subspace $\text{span}\{A\}$ and $\text{span}\{A'\}$ hit the solution
 - A and A' are nearly orthogonal: $\Pr[\text{span}\{A\} \perp \text{span}\{A'\}] \approx 0$
 - The training trajectory in A has zero projection in A' , vice versa
 - Can we make use of the redundant degrees of freedom?

$$w(\{\theta^{(i)}\}_{i=1}^K) = \sum_{i=1}^K A^{(i)} \theta^{(i)} + w_0$$



- ~4 x subnets: Havasi et al. MIMO. ICLR'21
~40 x inputs: Murahari et al. DataMux. 2022
- Larger D/d_{int} for (transformers, text)?
Warm-up allows more subnets?
Disentangled readout from all layers / via (soft) weight masking?