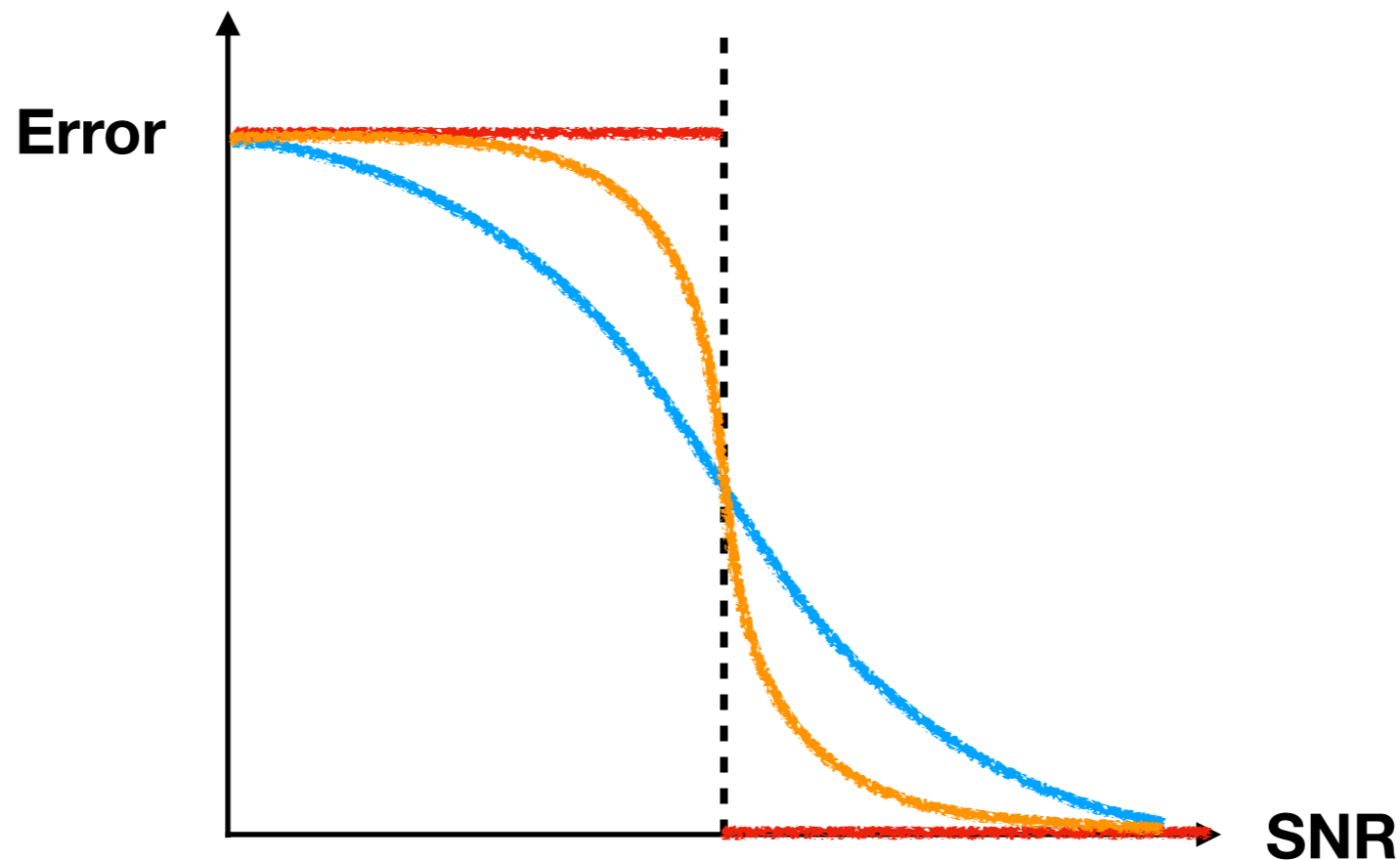


The « all-or-nothing phenomenon »

*A peculiar phase transition in sparse
high-dimensional inference & learning*

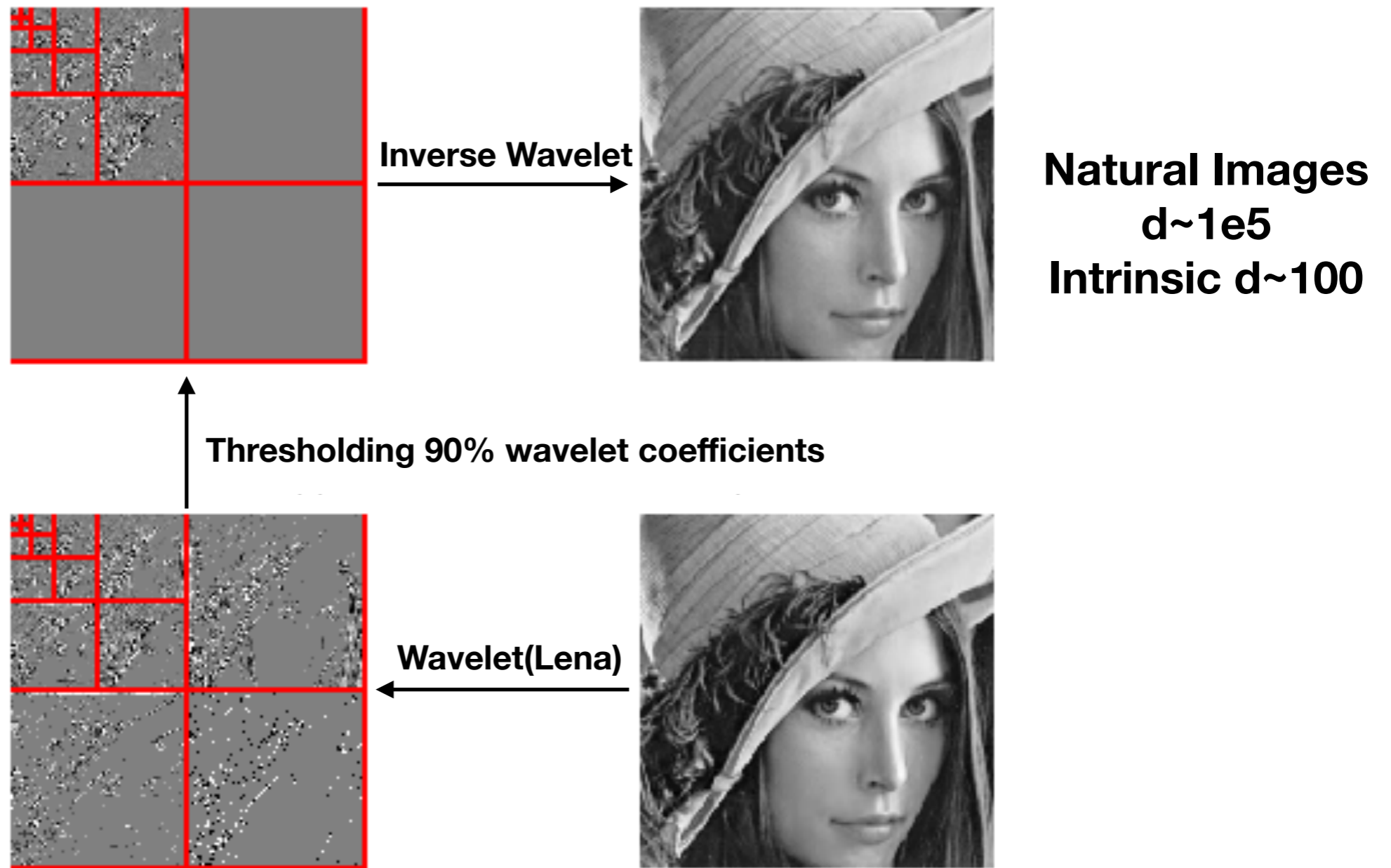
Jean Barbier

International Center for Theoretical Physics, Trieste



Joint with **Nicolas Macris** (EPFL), **Clément Luneau** (EPFL)
& **Cynthia Rush** (Columbia University)

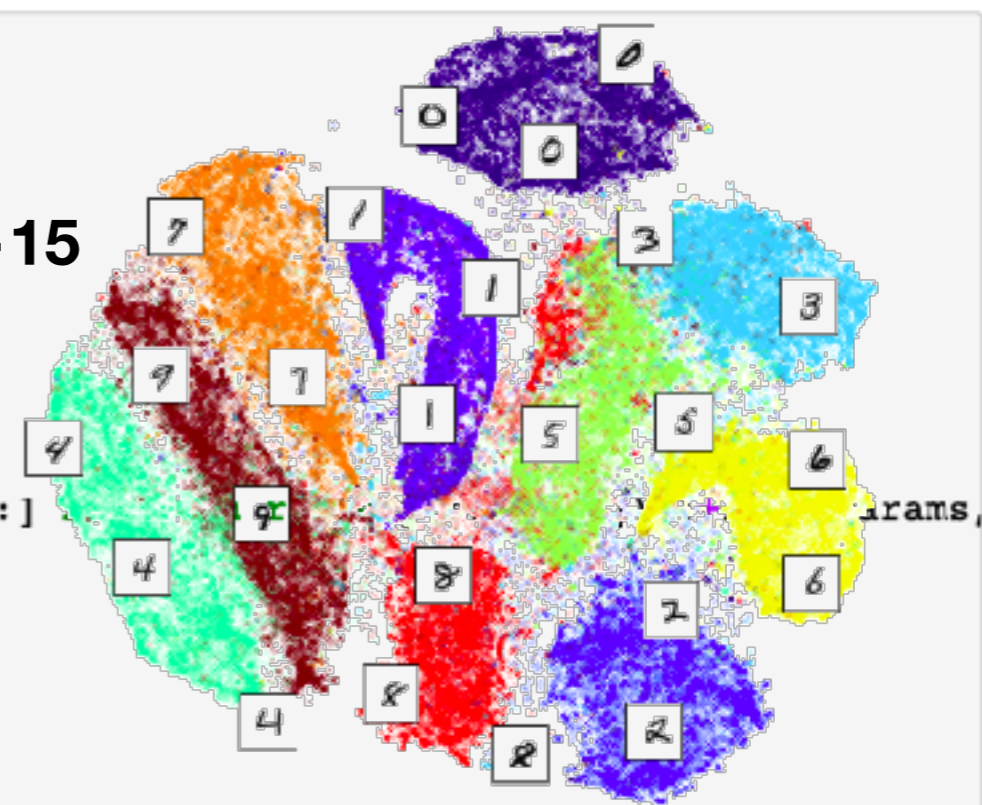
Low-dimensional manifolds underlying high-dimensional data



MNIST

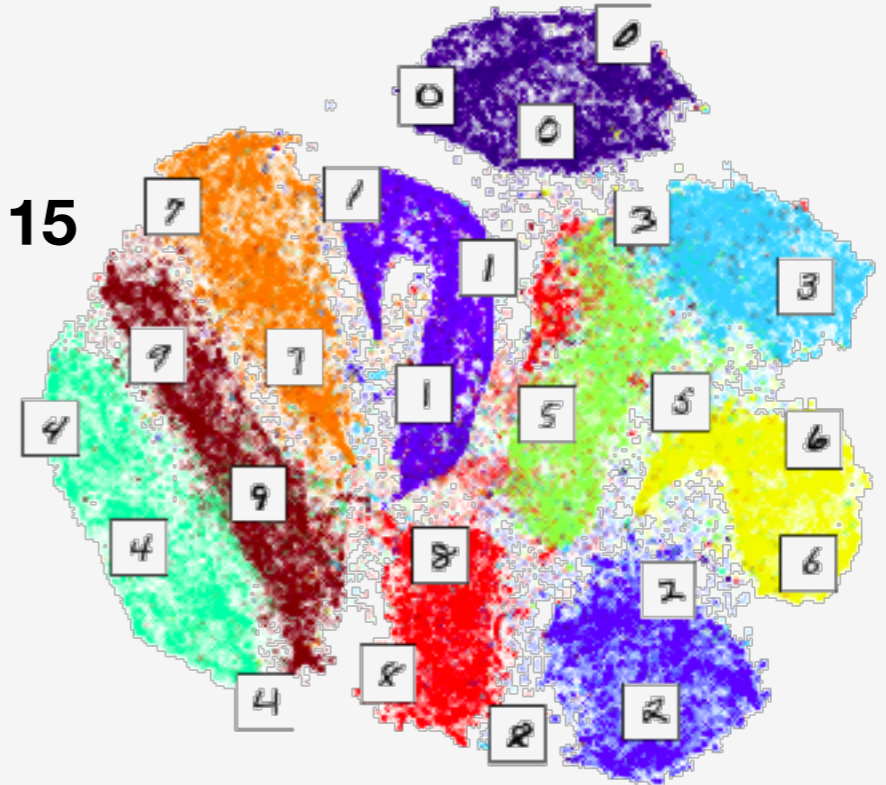
d=784

Intrinsic d~10–15



```
1 # logistic function
2 def logistic(x):
3     return 1/(1+np.exp(-x))
4
5 # classification error
6 def class_error(X,y,params):
7     return 100*np.mean(np.abs((logistic(np.dot(X,params))>0.5)-y))
8
9 # gradient of logistic regression with L2 penalization
10 def grad_log(X,y,params,lamb):
11     return np.sum(np.array([(logistic(np.dot(X[i,:],params))-y[i])*X[i,:]]
12
13 # divide the data in mini batches
14 def mini_batch(X,y,size_batch,number_batch):
15     perm=np.random.permutation(X.shape[0])
16     randX=X[perm,:]
17     randy=y[perm]
18
19     X_mini=np.array([randX[i*size_batch:min(i*size_batch+size_batch,X.shape[0]),:] for i in range(number_batch)])
20     y_mini=np.array([randy[i*size_batch:min(i*size_batch+size_batch,X.shape[0])] for i in range(number_batch)])
21     return X_mini,y_mini
22
23 # SGD with momentum
24 def SGD(X,y,grad,init,epochs=20,learn_rate=10**-4,momentum=0.9,L2_pen=0.01,size_batch=10,ground_truth=0):
25     params=np.array(init)
26     v=0
27     number_batch=int(np.ceil(X.shape[0]/size_batch))
28
29     for j in range(epochs):
30         X_mini,y_mini=mini_batch(X,y,size_batch,number_batch)
31
32         for k in range(number_batch):
33             v=momentum*v+learn_rate*np.array(grad(X_mini[k],y_mini[k],params,L2_pen))
34             params=params-v
35
36         train_er=class_error(X,y,params)
37         print(['epoch',j,'train error',train_er])
38
39     return params
```

MNIST
d=784
Intrinsic d~10–15



```
1 n_train=5000
2 n_test=8000
3 learn_rate=1e-4
4 size_batch=20
5 L2_pen=0.01
6 epochs=20
7
8 # select subset of two numbers
9 number_1='1'
10 number_2='9'
11
12 #####
13
14 N=784 # number of pixels/features
15
16 # select subset of two numbers
17 select=np.where((y_mnist==number_1)|(y_mnist==number_2))
18 X,y=np.array(X_mnist[select[0],:],np.array(y_mnist[select[0]]))
19 y[y==number_1]=0
20 y[y==number_2]=1
21 y=y.astype('float')
22
23 # pick training and test data sets
24 perm=np.random.permutation(X.shape[0])
25 X,y=X[perm,:],y[perm]
26 X_train,X_test,y_train,y_test=train_test_split(X,y,train_size=n_train,test_size=n_test)
27
28 # standardize data
29 scaler=StandardScaler()
30 X_train=scaler.fit_transform(X_train)
31 X_test=scaler.transform(X_test)
32
33 # SGD
34 params=SGD(X_train,y_train,grad_log,np.zeros(N),epochs=epochs,learn_rate=learn_rate,momentum=0.9,L2_pen=L2_pen,size
35
36 # check accuracy
37 print(['train: % misclassified', class_error(X_train,y_train,params)])
38 print(['test: % misclassified', class_error(X_test,y_test,params)])
```

```
['epoch', 17, 'train error', 0.24]
['epoch', 18, 'train error', 0.24]
['epoch', 19, 'train error', 0.24]
['train: % misclassified', 0.24]
['test: % misclassified', 0.525]
```



“Bet on sparsity principle”

Intrinsic / effective low-dimensionality is often a crucial ingredient for the interpretability of high-dimensional data

Statistical physics approach:

Idealised models of high-dimensional data with low intrinsic / effective dimension

- Average case / ensemble analysis -> **typical behaviour**
- Information-theoretic and algorithmic **phase transitions**
- **Exact asymptotic thresholds and formulas in « high-dimensional » regime**

free energy / mutual information, minimum mean-square and/or generalisation errors

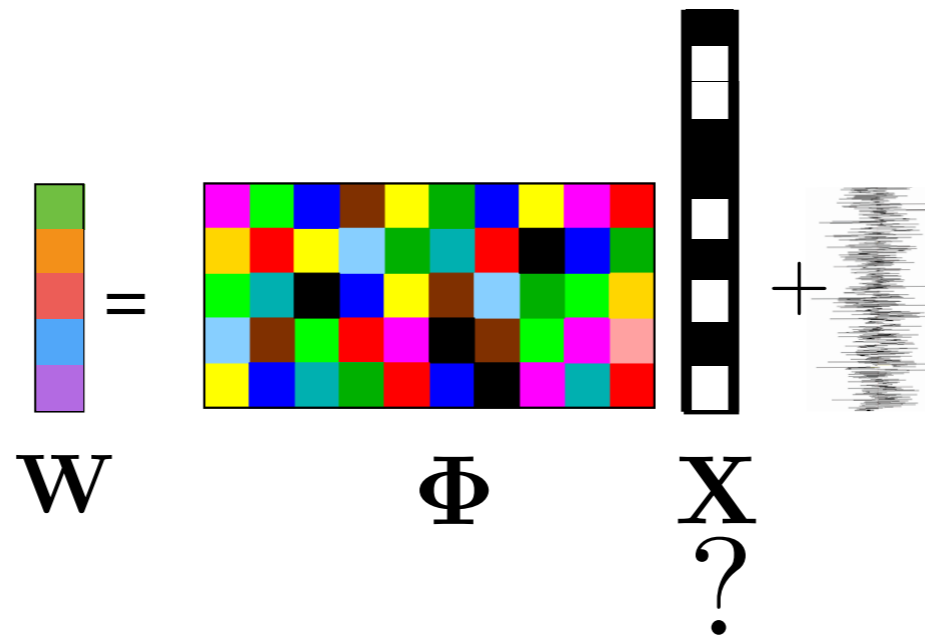
Paradigmatic models for high-dimensional inference & learning

Random linear estimation / compressive sensing

[Barbier Macris Dia Krzakala 16], [Reeves Pfister 16]

$$\mathbf{W} = \sqrt{\frac{\lambda}{n}} \Phi \mathbf{X} + \mathbf{Z}$$

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



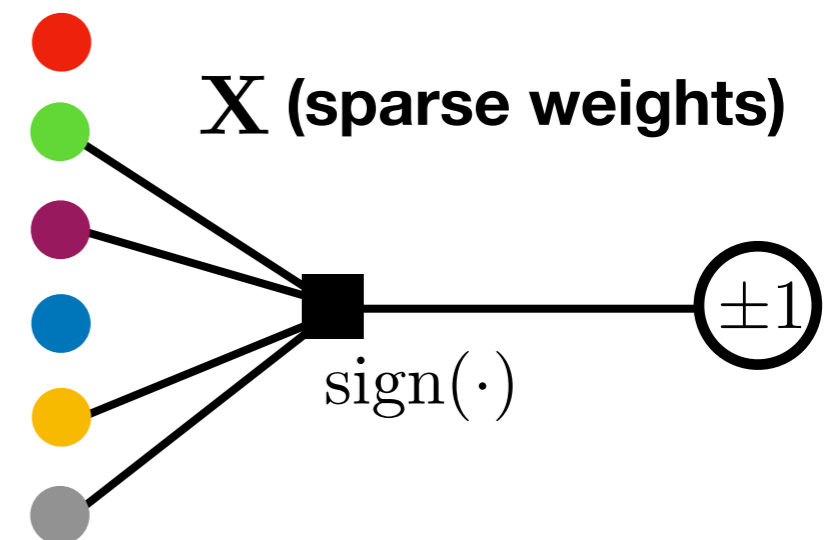
Generalized linear models

[Barbier Miolane Macris Krzakala Zdeborová 18]

$$\mathbf{W} = \varphi \left(\sqrt{\frac{1}{n}} \Phi \mathbf{X} \right)$$

Perceptron neural net $\mathbf{W} = \text{sign}(\Phi \mathbf{X})$

« Teacher-student » scenario
= realisable rule setting in ML



Signal processing

$$\mathbf{Y} = \mathbf{\Phi}\mathbf{X} + \mathbf{Z}\sqrt{\Delta}$$

Compressed sensing (signal processing)
CDMA (multi-user communication theory)
Superposition codes (coding theory)

$$\mathbf{Y} = |\mathbf{\Phi}\mathbf{X}|$$

Phase retrieval (signal processing)

$$\mathbf{Y} = \text{sign}(\mathbf{\Phi}\mathbf{X} + \mathbf{Z}\sqrt{\Delta})$$

1-bit compressed sensing (signal processing)
Perceptron (statistics/ML)

Learning

$$\mathbf{Y} = \max(0, \mathbf{\Phi}\mathbf{X})$$

Rectified Linear Unit (ReLU) (ML/neural nets)

$$P(Y_{\mu} = 1) = \frac{1}{1 + \exp(-\lambda\mathbf{\Phi}_{\mu} \cdot \mathbf{X})}$$

Sigmoid/Logistic regression (ML/neural nets)

Low-rank matrix factorisation

Sparse principal components analysis

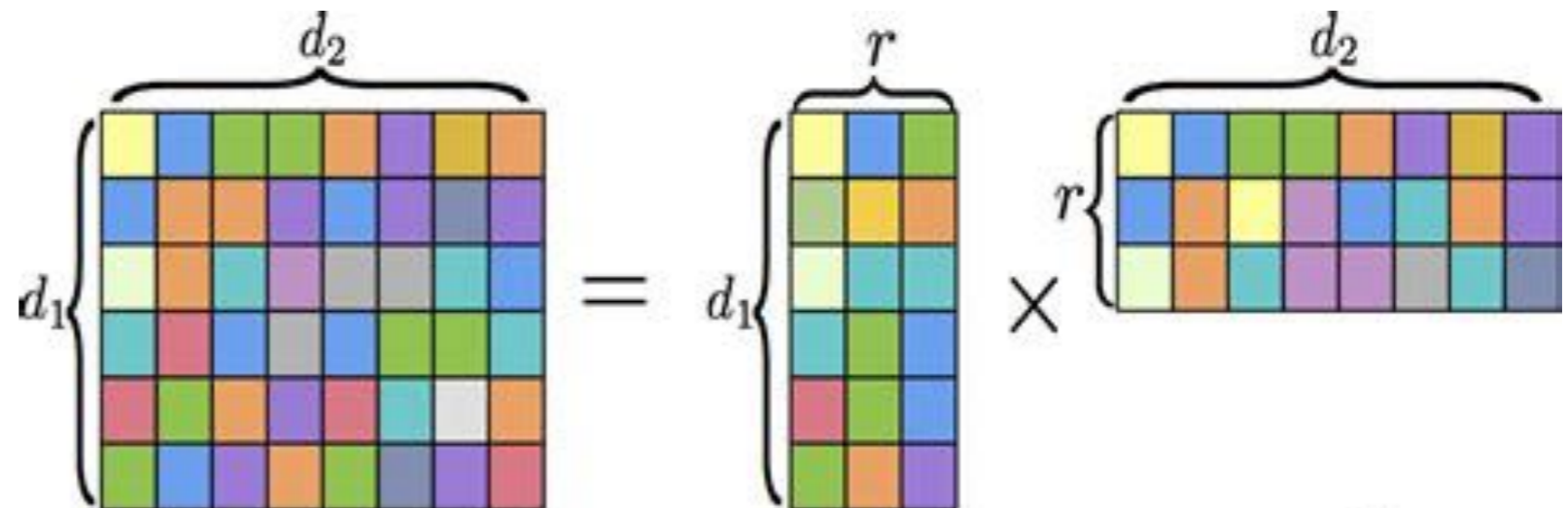
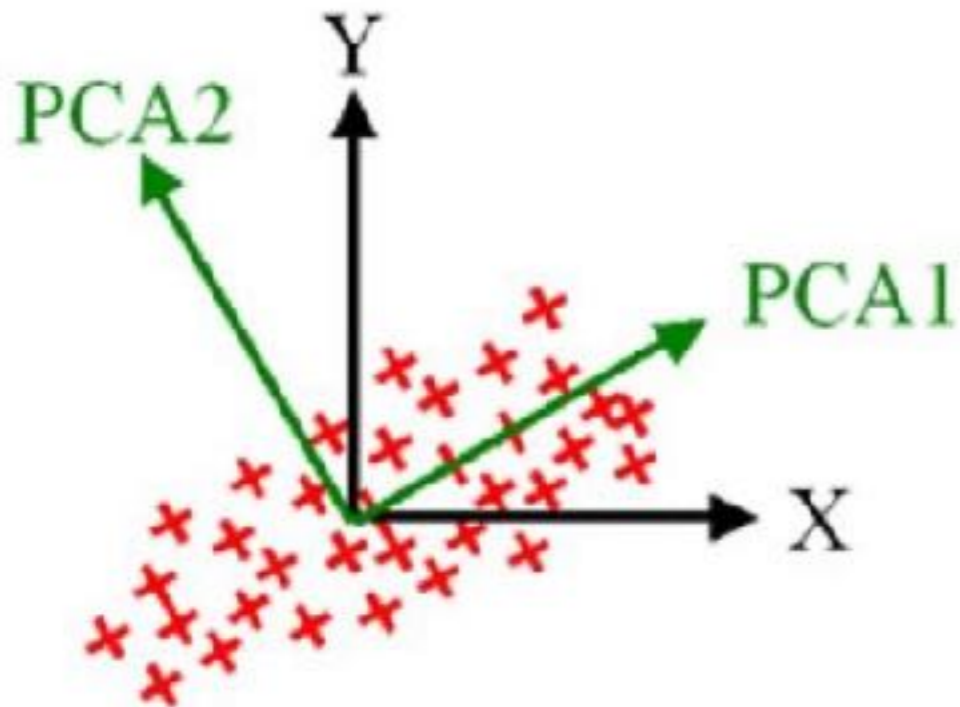
[Korada Macris 09]

[Deshpande Montanari 14]

[Barbier Macris Dia Lesieur Krzakala Zdeborová 16]

[Miolane Lelarge 17]

$$\mathbf{W} = \sqrt{\frac{\lambda}{n}} \mathbf{X} \mathbf{X}^T + \mathbf{Z}$$

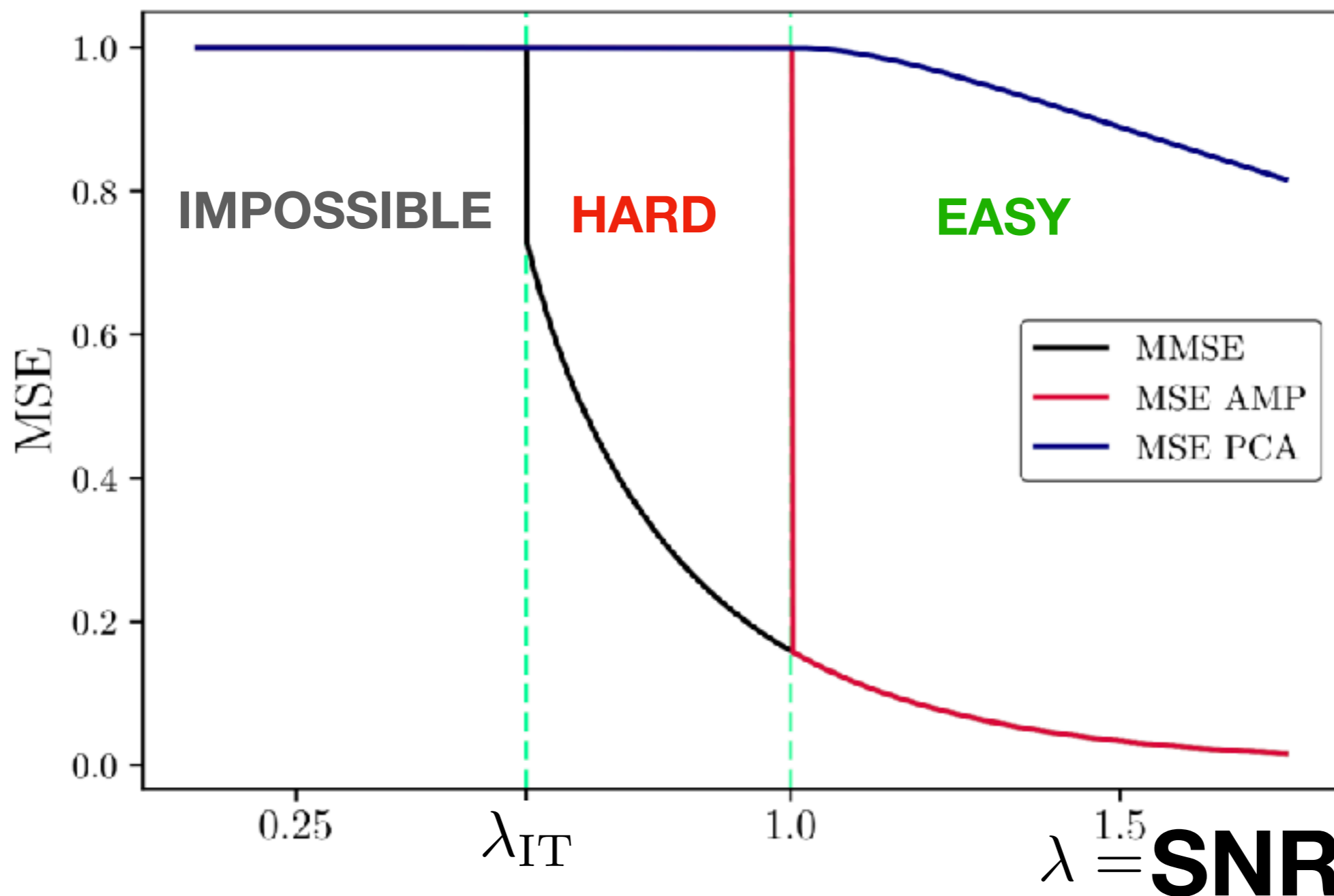


Generic scenario: discontinuous phase transition & computational gap

$$\mathbf{W} = \sqrt{\frac{\lambda}{n}} \mathbf{X}\mathbf{X}^\top + \mathbf{Z}$$

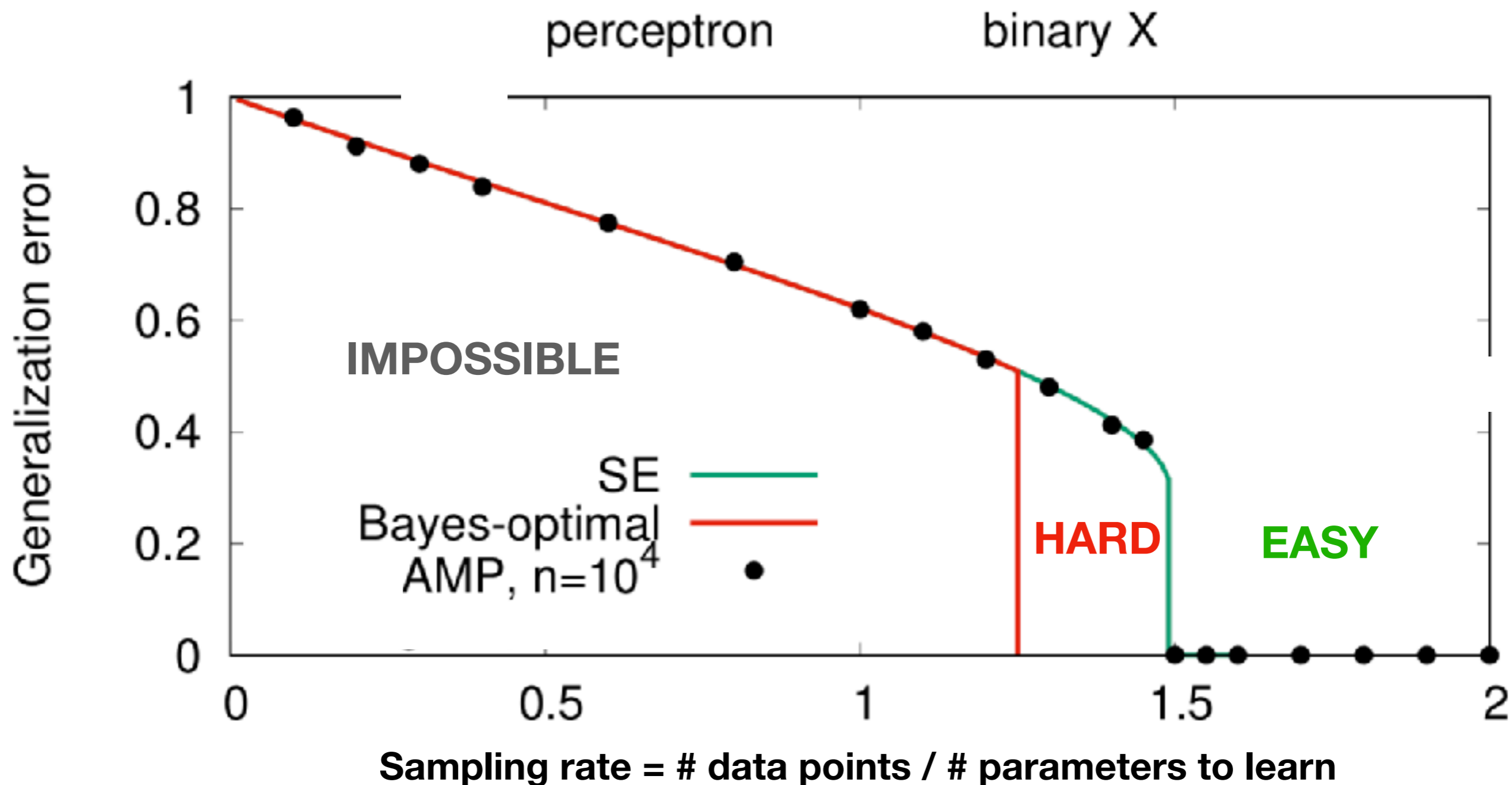
Binary X

$$\lim_{n \rightarrow \infty} \frac{1}{n^2} \mathbb{E} \|\mathbf{X}\mathbf{X}^\top - \mathbb{E}[\mathbf{X}\mathbf{X}^\top | \mathbf{W}]\|_F^2$$

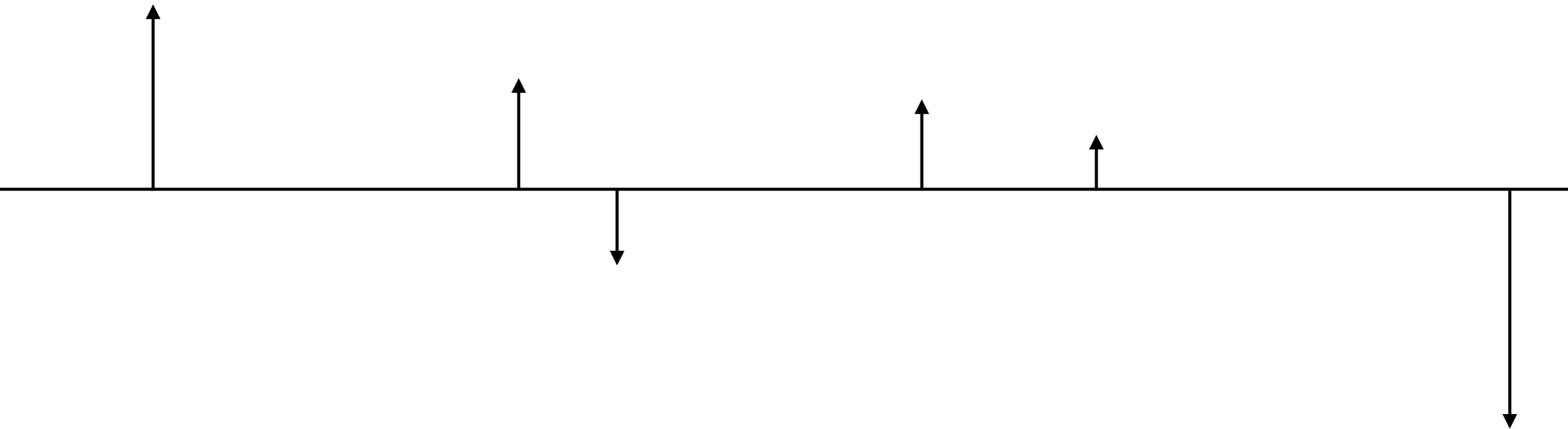


Generic scenario: discontinuous phase transition & computational gap

$$\mathbf{W} = \text{sign}(\Phi \mathbf{X}) \quad \lim_{n \rightarrow \infty} \mathbb{E} \left[\left(W_{\text{new}} - \mathbb{E}[W_{\text{new}} | \mathbf{W}, \Phi, \Phi_{\text{new}}] \right)^2 \right]$$



**Anything special in the
very high sparsity / low effective
dimension regime?**



Compressed sensing: low measurement-rate regime

$$W_\mu = \sqrt{\frac{\lambda}{n}} \Phi_\mu \cdot \mathbf{X} + Z_\mu \quad \mu = 1, \dots, m$$

[Gamarnik Zadik 17]

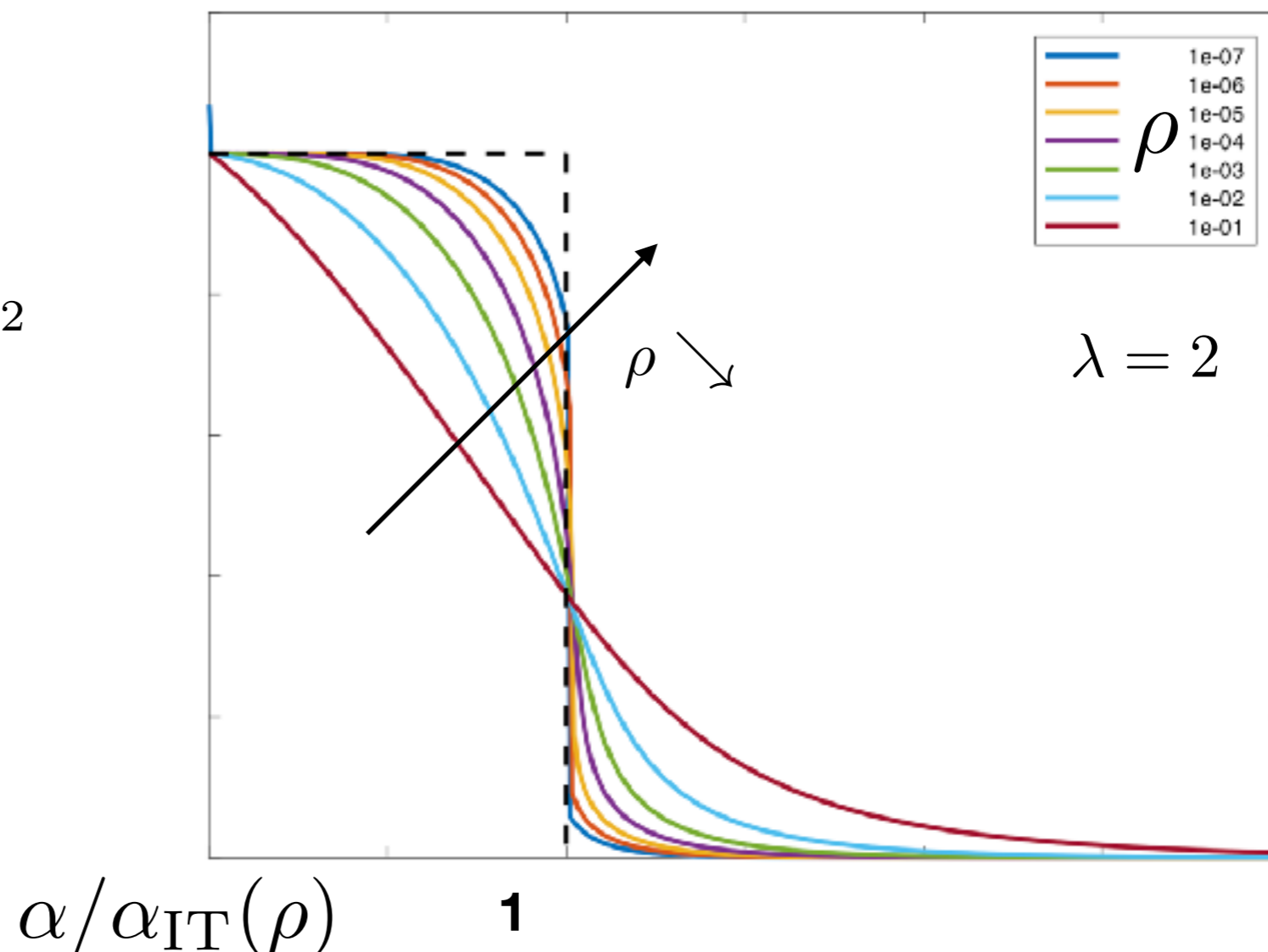
[Reeves Xu Zadik 19]

$$X_i \sim \text{Ber}(\rho) \quad i = 1, \dots, n$$

(informal) Theorem [Barbier Macris Dia Krzakala 16] [Reeves Pfister 16]:

$$\lim_{n \rightarrow \infty} \frac{1}{\rho n} \mathbb{E} \|\mathbf{X} - \mathbb{E}[\mathbf{X} | \mathbf{W}]\|^2 = \frac{1}{\rho} \text{mmse}(q^*, \lambda)$$

$$\lim_{n \rightarrow \infty} \frac{1}{\rho n} \mathbb{E} \|\mathbf{X} - \mathbb{E}[\mathbf{X} | \mathbf{W}]\|^2$$



$$W_\mu = \sqrt{\frac{\lambda_n}{\rho_n n}} \Phi_\mu \cdot \mathbf{X} + Z_\mu \quad \mu = 1, \dots, \alpha_n n$$

$$X_i \sim \text{Ber}(\rho_n) \quad i = 1, \dots, n \quad \lambda_n = \Omega(1)$$

(informal) Theorem [Reeves Xu Zadik 19]: « **All-or-nothing** » phenomenon

$$\rho_n = o(n^{-1/2}) \text{ and } \alpha < \alpha_{\text{IT},n} \text{ with } \alpha_{\text{IT},n} = \frac{2\rho_n |\ln \rho_n|}{\ln(1 + \lambda_n)}$$

→ **Weak recovery impossible:** $\liminf_{n \rightarrow \infty} \frac{1}{\rho_n n} \mathbb{E} \|\mathbf{X} - \hat{\mathbf{X}}_{\text{opt}}\|^2 = 1$

$$\rho_n = o(1) \text{ and } \alpha > \alpha_{\text{IT},n}$$

→ **Strong recovery possible:** $\limsup_{n \rightarrow \infty} \frac{1}{\rho_n n} \mathbb{E} \|\mathbf{X} - \hat{\mathbf{X}}_{\text{opt}}\|^2 = 0$

All-or-nothing in sparse principal components analysis

$$\mathbf{W} = \sqrt{\frac{\lambda_n}{n}} \mathbf{X} \otimes \mathbf{X} + \mathbf{Z}$$

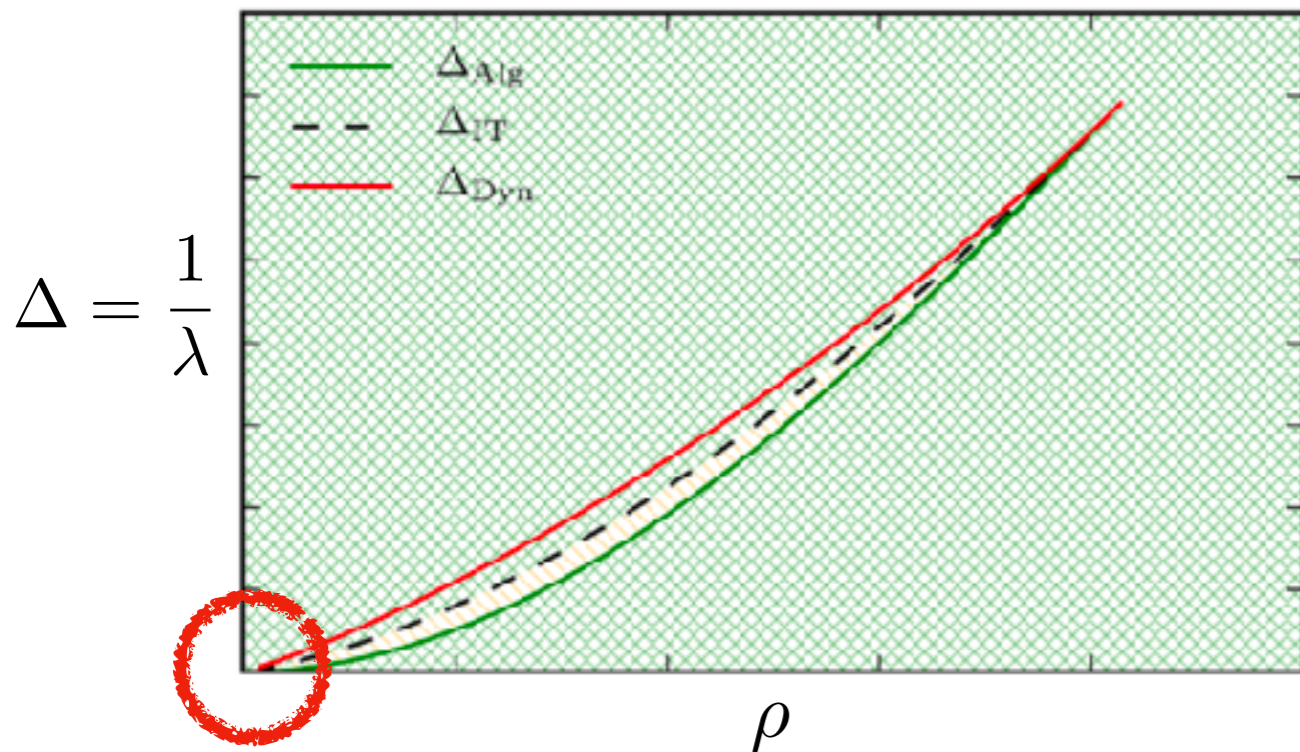
$$X_i \sim \text{Ber}(\rho_n) \quad i = 1, \dots, n$$

$$\rho_n \rightarrow 0_+$$

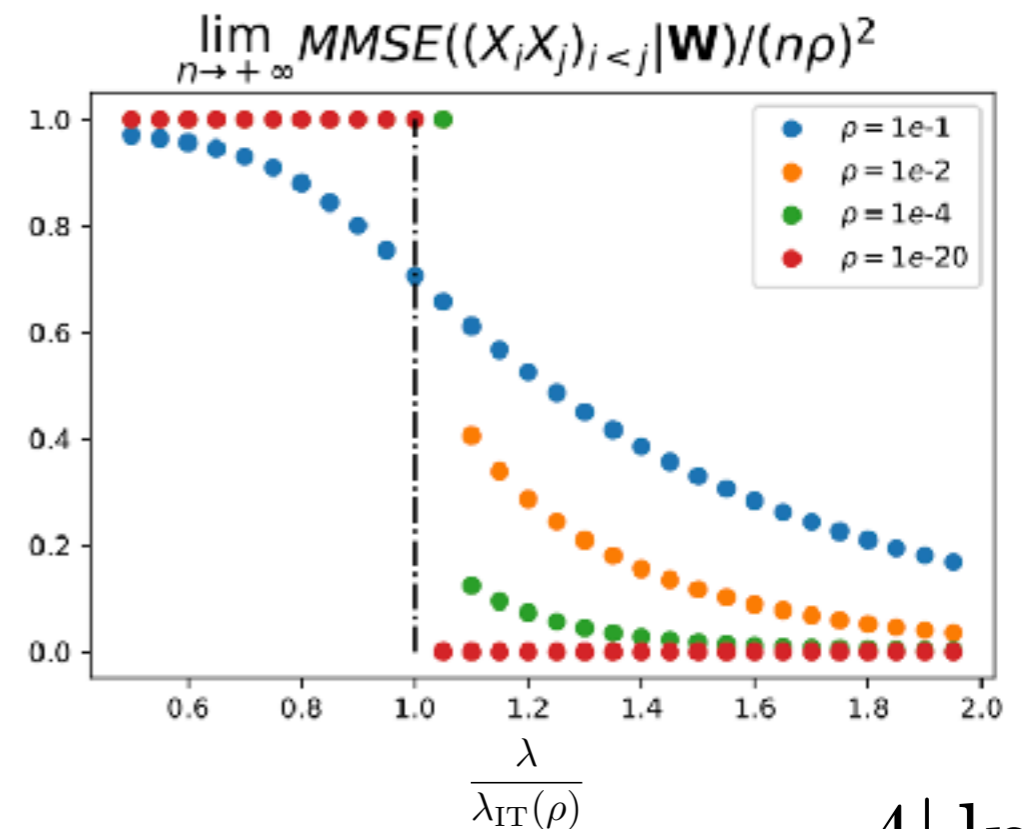
$$\lambda_n \rightarrow +\infty$$

Theorem [Korada Macris 09], [Deshpande Montanari 14],
[Barbier Macris Dia Lesieur Krzakala Zdeborová 16], [Miolane Lelarge 17]:

$$\lim_{n \rightarrow \infty} \frac{1}{(\rho n)^2} \mathbb{E} \|\mathbf{X} \otimes \mathbf{X} - \mathbb{E}[\mathbf{X} \otimes \mathbf{X} | \mathbf{W}]\|_F^2 = \text{mmse}(q^*)$$



[Lesieur Krzakala Zdeborová 15]



$$\lambda_{\text{IT}}(\rho) \xrightarrow{\rho \rightarrow 0} \frac{4 |\ln \rho|}{\rho}$$

Finite size bound for the mutual information

Theorem [Barbier Macris Rush 20]: based on the adaptive interpolation method

Let $P_{0,n} = \rho_n p_0 + (1 - \rho_n)\delta_0$. Let $\lambda_n = \Theta(|\ln \rho_n|/\rho_n)$ which is the appropriate scaling to observe a phase transition, and $\rho_n = \Theta(n^{-\beta})$ with $\beta \in [0, 1/6)$. Then

$$\frac{1}{\rho_n |\ln \rho_n|} \left| \frac{1}{n} I(\mathbf{X}; \mathbf{W}) - \inf_{q \in [0,1]} i_n^{\text{pot}}(q; \lambda_n, \rho_n) \right| \leq C \frac{(\ln n)^{1/3}}{n^{(1-6\beta)/7}}$$

with $i_n^{(\text{pot})}(q; \lambda_n, \rho_n) = \frac{\lambda_n \rho_n^2}{4} (q - 1)^2 + I_n(X; \sqrt{\lambda_n \rho_n q} X + Z)$ $X \sim P_{0,n}$, $Z \sim \mathcal{N}(0, 1)$.

Corollary [Barbier Macris Rush 20]:

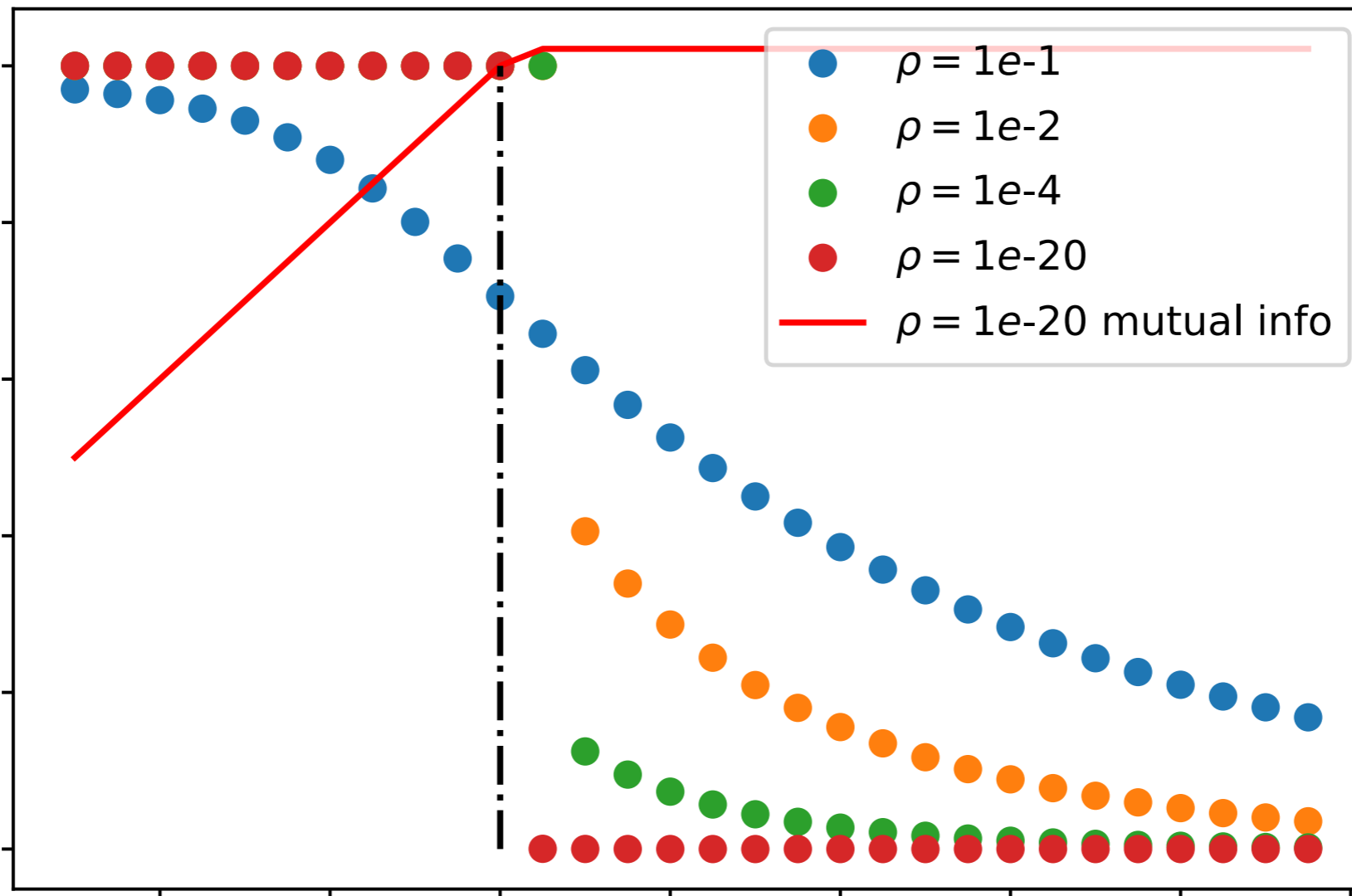
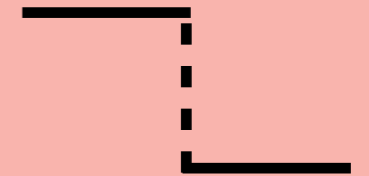
Let $\frac{1}{2} m_n(\lambda, \rho_n) \equiv \rho_n^{-2} \frac{d}{d\lambda} \inf_{q \in [0, \rho_n]} i_n^{\text{pot}}(q, \lambda, \rho_n)$. Let $\epsilon > 0$ and $\beta \in [0, 1/13)$. There exists $C' > 0$ independent of n such that

$$m_n(\lambda_n + \epsilon, \rho_n) - \frac{C'}{\epsilon} \frac{(\ln n)^{4/3}}{n^{(1-13\beta)/7}} \leq \frac{\text{MMSE}((X_i X_j)_{i < j} | \mathbf{W})}{(n \rho_n)^2} \leq m_n(\lambda_n - \epsilon, \rho_n) + \frac{C'}{\epsilon} \frac{(\ln n)^{4/3}}{n^{(1-13\beta)/7}}.$$

Information-theoretic all-or-nothing transition in sparse PCA

Setting $\lambda_n \equiv \gamma \lambda_{n,IT}(\rho_n)$ where $\lambda_{n,IT}(\rho_n) \equiv 4 \frac{|\ln \rho_n|}{\rho_n}$

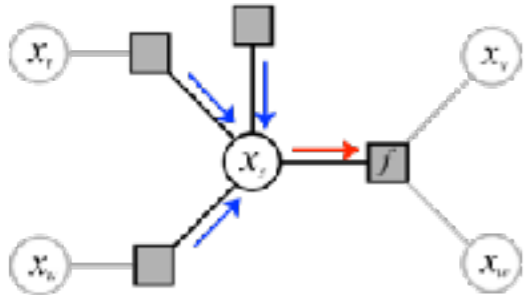
$$\frac{1}{(n\rho_n)^2} \mathbb{E} \|\mathbf{X} \otimes \mathbf{X} - \mathbb{E}[\mathbf{X} \otimes \mathbf{X} | \mathbf{W}]\|_F^2 \xrightarrow{n \rightarrow \infty} \mathbb{I}(\gamma \leq 1)$$



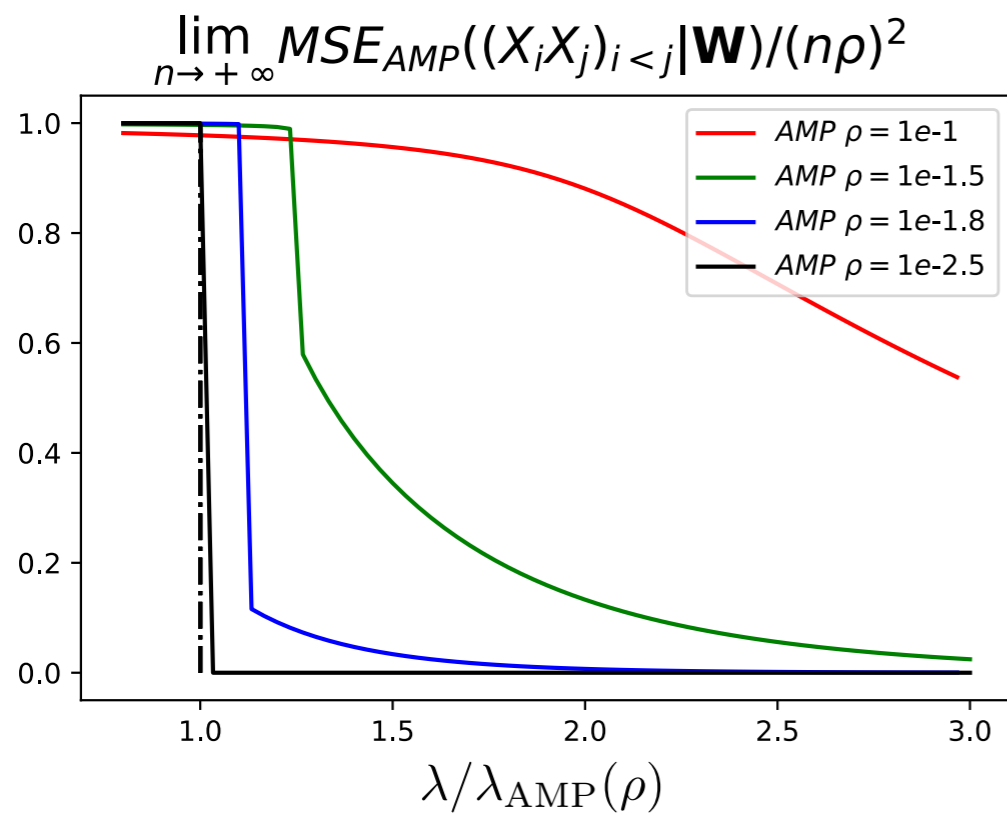
Algorithmic all-or-nothing transition in sparse PCA

Belief Propagation (BP) / Sum-product algo for DENSE factor graphs:
Approximate Message-Passing algo (AMP), originally introduced for CS

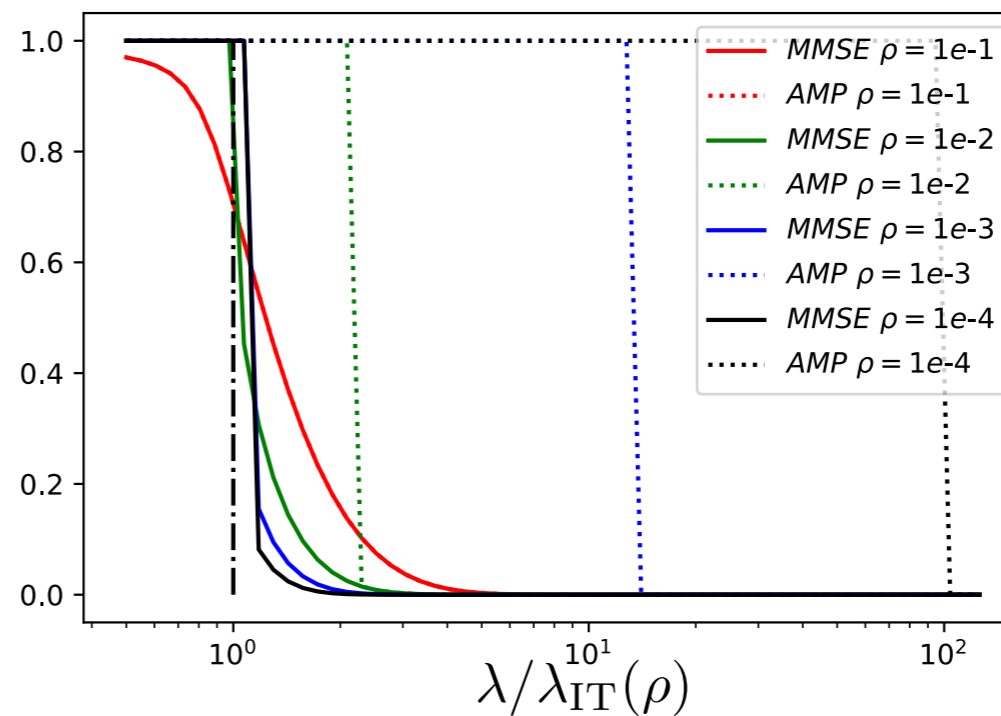
AMP:
$$\mathbf{x}^{t+1} = \frac{1}{\sqrt{n}} \mathbf{W} f_t(\mathbf{x}^t) - \mathbf{b}_t f_{t-1}(\mathbf{x}^{t-1}), \quad \mathbf{b}_t = \frac{1}{n} \sum_{i=1}^n f_t'(x_i^t), \quad \hat{\mathbf{X}}^t = f_t(\mathbf{x}^t)$$



- Exact characterisation in a proper asymptotic limit thanks to **STATE EVOLUTION** [Bayati, Montanari 12']
- **Conjectured to be optimal among all polynomial time algo**



MMSE vs MSE of AMP

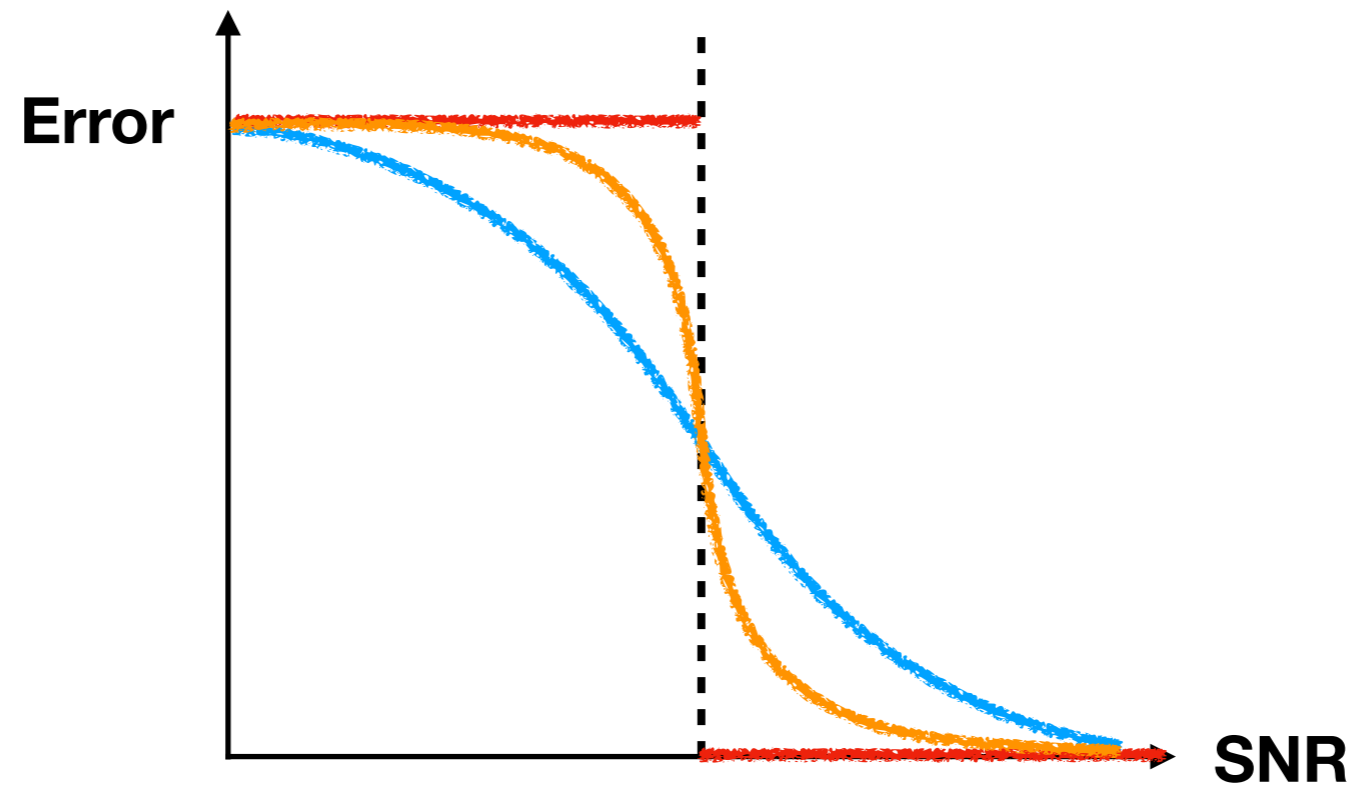


$$\lambda_{AMP}(\rho) = 1/(e\rho)^2$$

$$\gg \lambda_{IT}(\rho) = 4 \frac{|\ln \rho|}{\rho}$$

(Informal) Theorem [Barbier Macris Rush 20]:
AMP finite sample analysis for low sparsity

The state evolution recursion allowing to track the performance of AMP, and therefore the curves above, are rigorously valid as $\lambda_n = \Theta(\rho_n^{-2})$ and $\rho_n = \Theta(1/(\ln n)^\alpha)$ for any fixed $\alpha \geq 0$ and n sufficiently large.



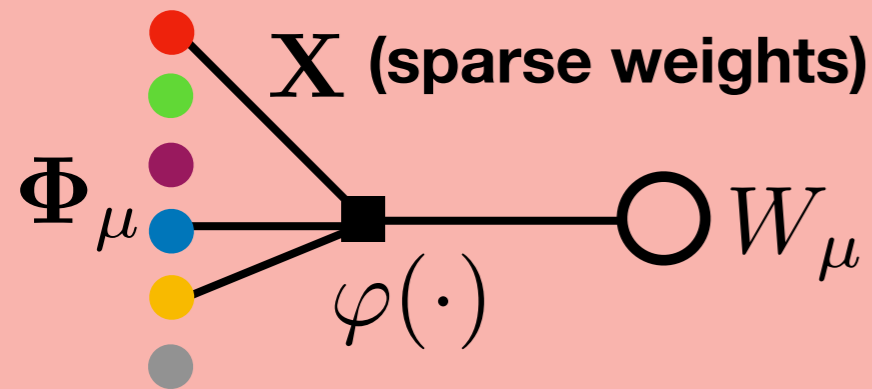
**All-or-nothing transition:
a « universal » phenomenon?**

Information-theoretic all-or-nothing transition in GLMs

Theorem [Luneau Macris Barbier 20]:

Suppose that the following hypotheses hold:

1. The entries of \mathbf{X} are bounded.
2. φ is bounded, and its first and second partial derivatives are bounded and continuous.
3. $(\Phi_{\mu i}) \sim_{\text{iid}} \mathcal{N}(0, 1)$.



Let $\rho_n = \Theta(n^{-\lambda})$ with $\lambda \in [0, 1/9)$ and $\alpha_n = \gamma \rho_n |\ln \rho_n|$ with $\gamma > 0$. Then for all n :

$$\left| \frac{I(\mathbf{X}; \mathbf{W} | \Phi)}{\alpha_n n} - \inf_{q \geq 0} \sup_{r \geq 0} i_{\text{RS}}(q, r; \alpha_n, \rho_n) \right| \leq \frac{\sqrt{C} |\ln n|^{1/6}}{n^{\frac{1}{12} - \frac{3\lambda}{4}}}$$

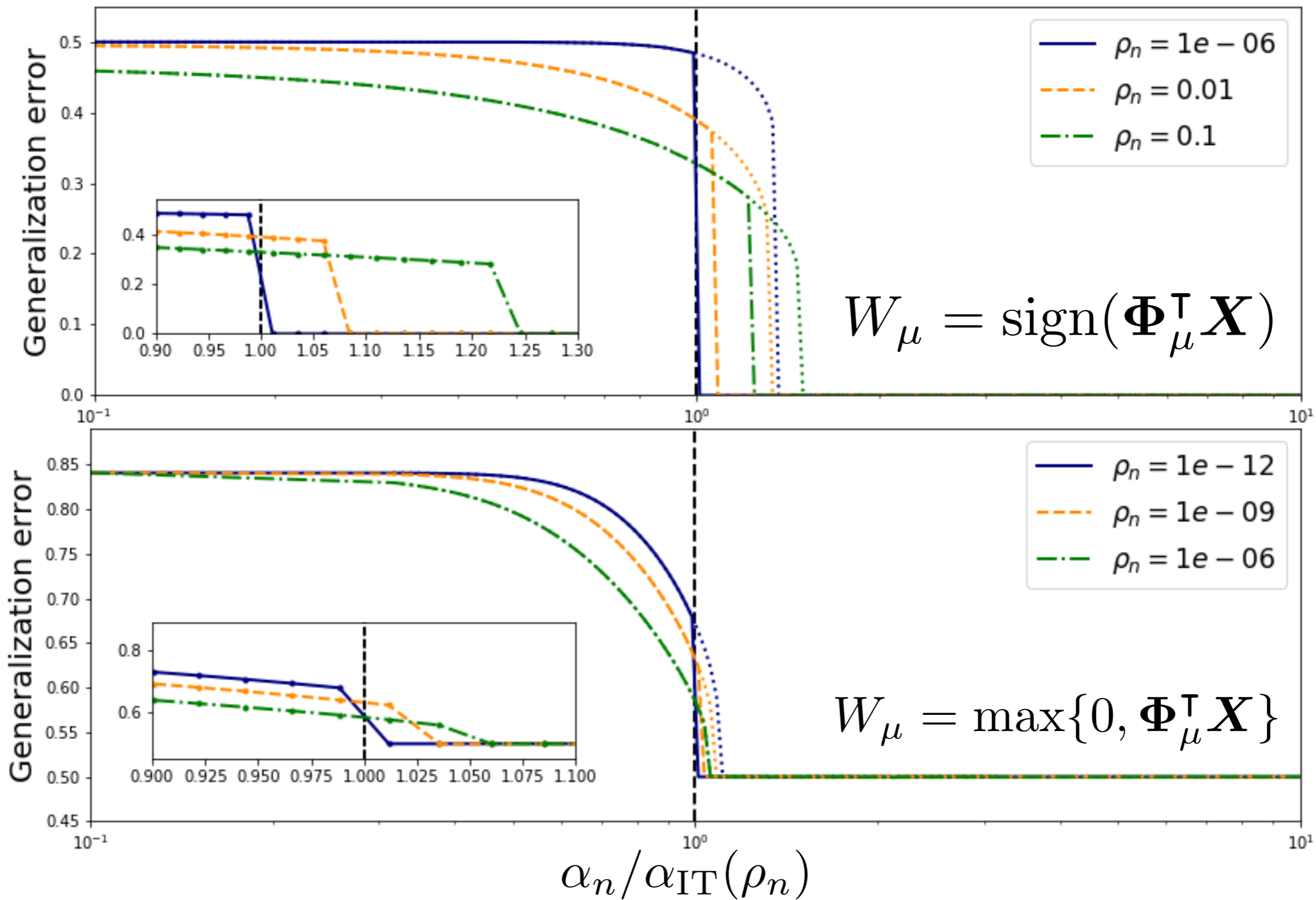
for some constant $C(\varphi) > 0$.

+ analytic exact expression for generalisation: $\lim_{n \rightarrow \infty} \mathbb{E} \left[\left(W_{\text{new}} - \mathbb{E}[W_{\text{new}} | \mathbf{W}, \Phi, \Phi_{\text{new}}] \right)^2 \right]$

Theorem [Luneau Macris Barbier 20]: In the Bernoulli signal case,

$$\lim_{\rho_n \rightarrow 0^+} I(\rho_n, \alpha_n) = \min \left\{ I_{P_{\text{out}}}(0, 1), \frac{1}{\gamma} \right\}.$$

Information-theoretic all-or-nothing transition in GLMs



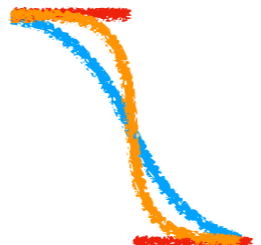
$$\alpha_{IT}(\rho_n) = \gamma_\varphi (H_b(\rho_n) + \rho_n H(P_X)) \xrightarrow{\rho_n \rightarrow 0} \gamma_\varphi \rho_n |\ln \rho_n|$$

- **Statistical physics predictions rigorously valid in « corners » of phase diagrams** (with « control parameters » depending on the system size).
- « All-or-nothing » seems quite generic: **« Universal » phenomenon in reconstruction of very sparse information (i.e., with « effective low dimension ») from much higher-dimensional data?**

-> [Niles-Weed Zadik 20]: information-theoretic all-or-nothing transition appears for a large class of additive Gaussian models with any sub-linear sparsity:

$$W = \sqrt{\lambda}X + Z$$

- **« All-or-nothing » for AMP algo (at least in sparse PCA). What about other algorithms (Markov chain monte-carlo, gradient-based, etc)?**
- **What about more complex models (multi layers etc)?**

Thank  You