

Statistical mechanics of learning: the physicist's view on learning theory

Summer School: Machine Learning in Quantum Physics and Chemistry
University of Warsaw
Aug 2021

Marylou Gabrié
(NYU Center for Data Science/Flatiron Institute CCMathematics)

Physics and Machine Learning

▶ Machine learning for physics

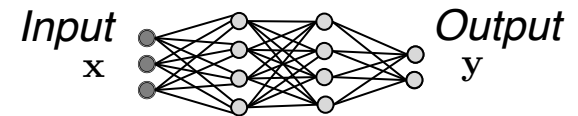
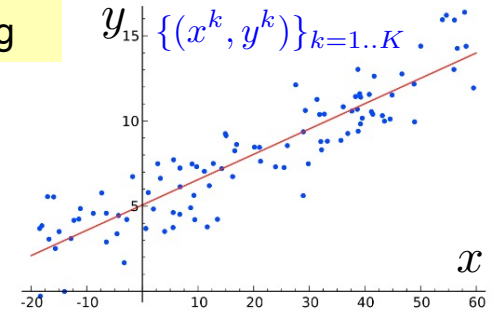
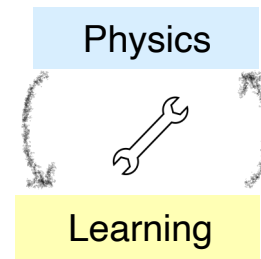
- Since the very beginning ...
 - Fitting models to experimental data
 - Data denoising

– e.g. linear regression $\hat{y} = ax + b$

- Recent deep learning “revolution”
 - Many sophisticated application (main topic of this school)

▶ Statistical physics for machine learning

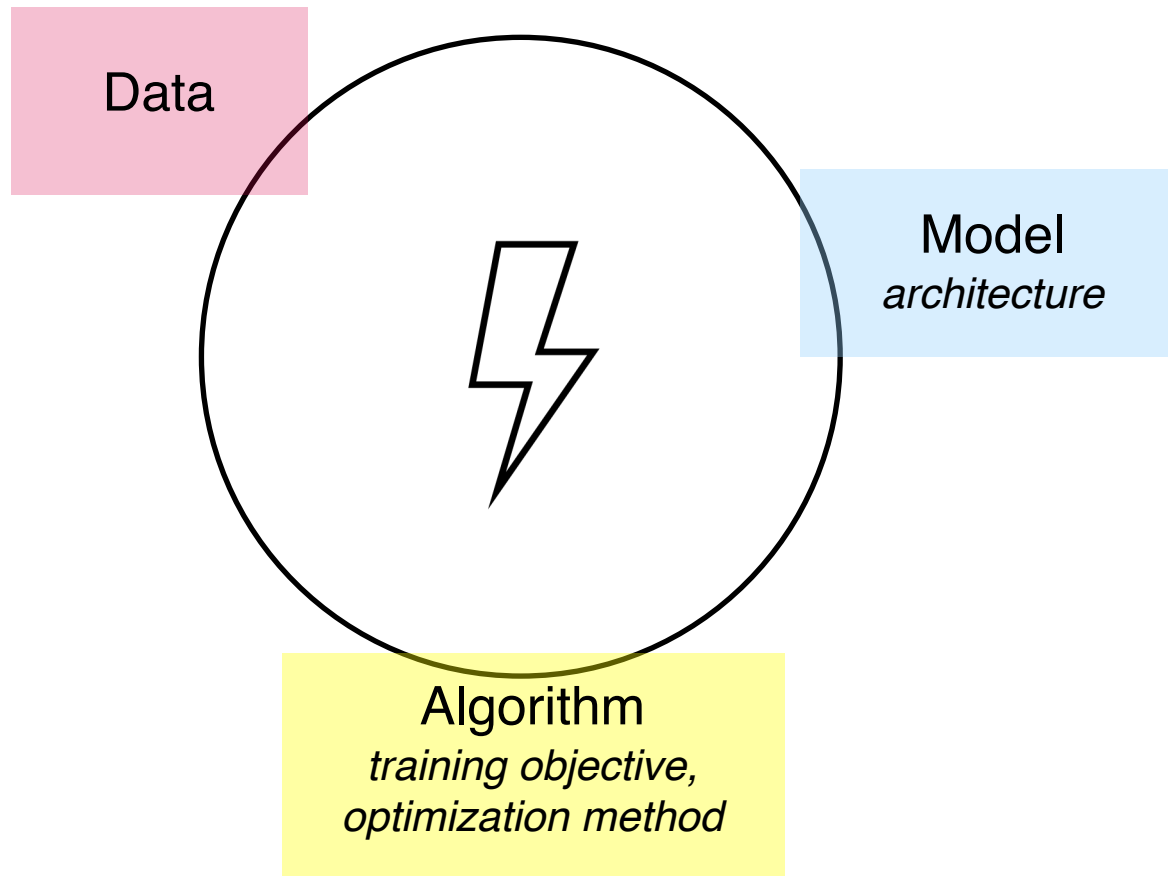
- In the 80s
 - Use the statistical physics toolbox
 - to study toy models of learning (mainly)
 - to design new learning algorithms
- Since mid 2010s: New wave of interest



thermodynamic limit approximate computations

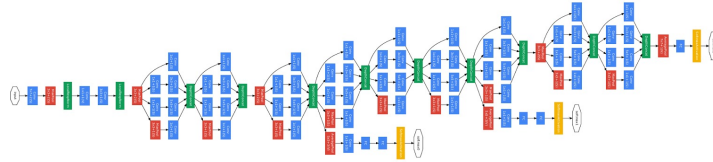


Understanding the successes of learning



Some puzzles of modern machine learning

Google le Net – inputs 112×122 pixels – ~ 7 millions parameters

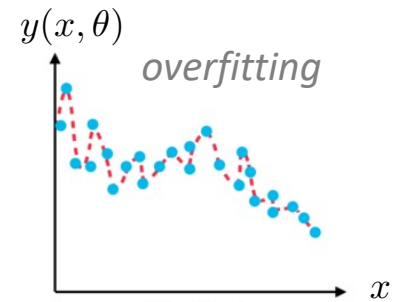
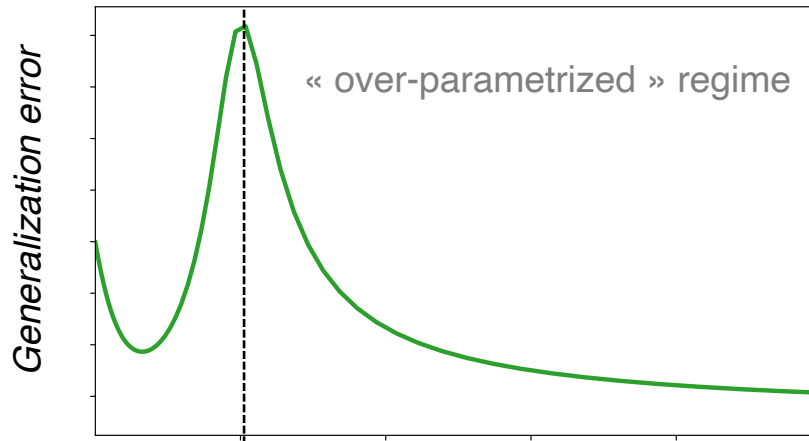
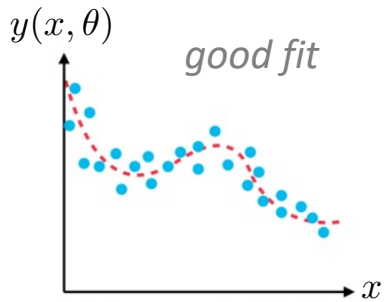
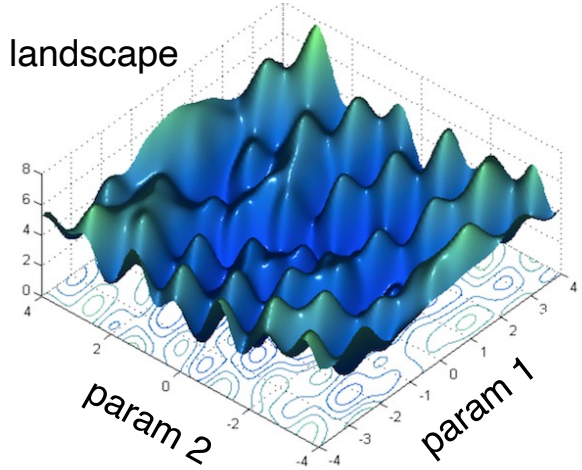


Szegedy et al, CVPR 2015

▷ Efficient optimization in non-convex settings

▷ Generalization in the « over-parametrized » regime

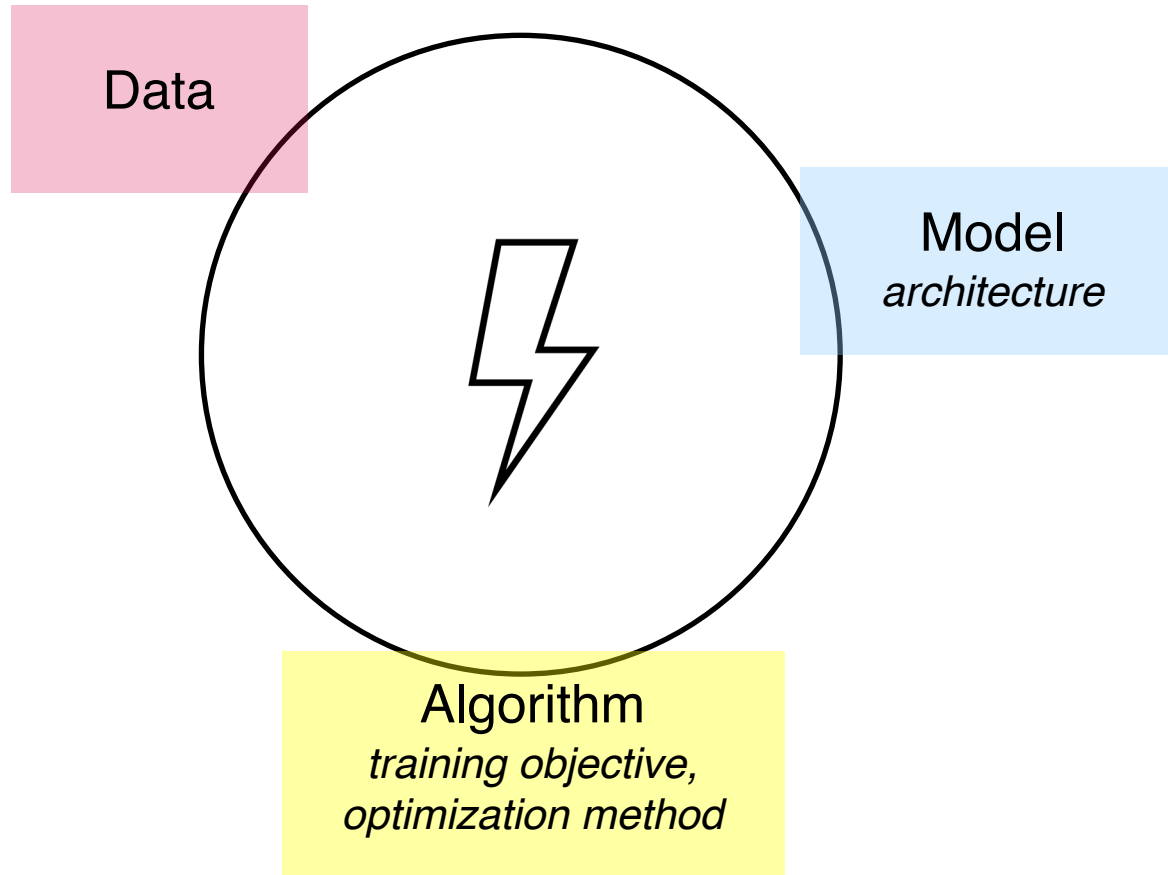
Loss landscape



« classical regime »
dilemma biais-variance

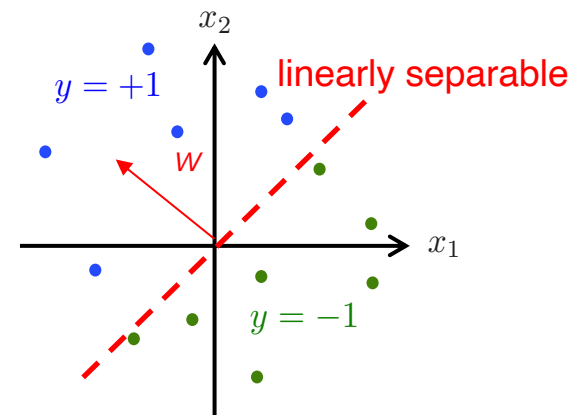
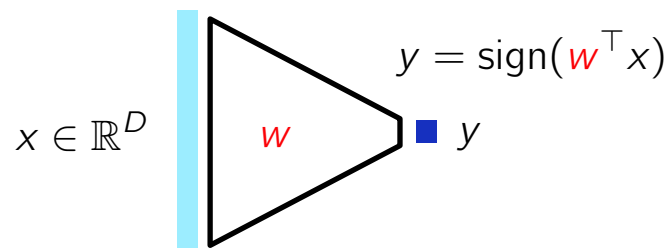
parameters / model flexibility

Understanding the successes of learning



Study simple 'solvable' models: the strategy of physicists (but not only!)

▷ Example: The perceptron – one neuron model



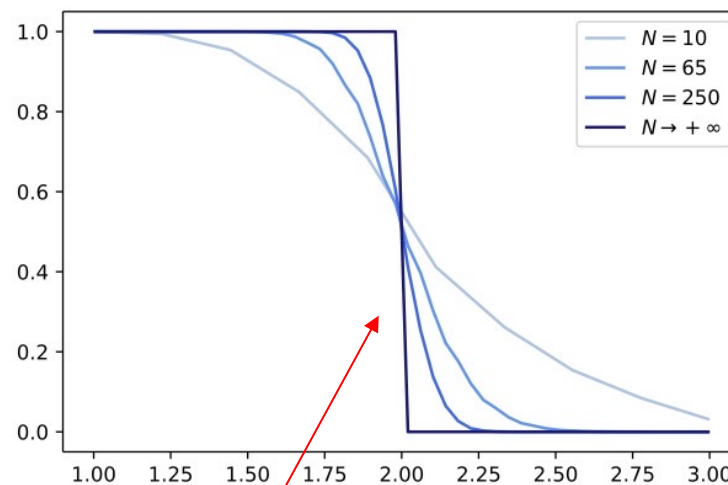
▷ First possible question: What is its capacity? How many data points can it fit?

- Training data: in general position
(« linearly independent ») $\{x^{(k)}\}_{k=1}^N$

- Probability random labelling
is linearly separable $\{y^{(k)}\}_{k=1}^N$
(Cover 1965)

- As a function of

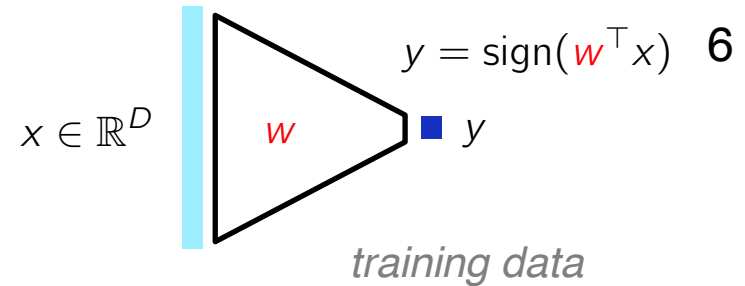
$$\alpha = \frac{\# \text{ training points}}{\# \text{ parameters}} = \frac{N}{D}$$



phase transition! $\alpha_c = 2$

$$\alpha = \frac{N}{D}$$

Perceptron's capacity reloaded: Statistical physics computation



▷ Volume of weight parameter space satisfying the labelling

$$\mathcal{D} = \{y^{(k)}, x^{(k)}\}_{k=1}^N$$

$$V_{D,N} = \int_{\Omega(w)} d\mathbf{w} \prod_{k=1}^N \delta(y(x^{(k)}; \mathbf{w}) - y^{(k)}) > 0$$

each new data point
= constraint/interaction for the weights

at least one solution!

$$V_{D,N} \propto \lim_{\beta \rightarrow +\infty} \int_{\Omega(w)} d\mathbf{w} e^{-\beta \sum_{k=1}^N (y(x^{(k)}, \mathbf{w}) - y^{(k)})^2}$$

$$V_{D,N} \propto \lim_{\beta \rightarrow +\infty} \int_{\Omega(w)} d\mathbf{w} e^{-\beta E(\mathbf{w}, \{x^{(k)}, y^{(k)}\})}$$

effective energy function

▷ Why is this hard? high-dimensional + depends on training set

Disordered systems physics and applications to learning

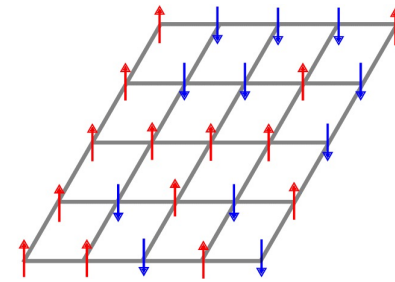
▷ Two types of random variables

- state of the system $s \in \mathbb{R}^d$
spins
- disordered interactions $J \in \mathbb{R}^N$
couplings

« learning parameters »

« data $\{x_i, y_i\}_{i=1}^N$

= constraint on parameters to learn »



▷ Boltzmann weight: $\rho(s|J) = e^{-E(s;J)} / \mathcal{Z}_J$

normalisation: $\mathcal{Z}_J = \int_{\mathbb{R}^d} ds e^{-E(s;J)}$

- ex: spin glass $E(s; J) = - \sum_{(i,j)} J_{ij} s_i s_j$ $J_{ij} \sim \mathcal{N}(J_{ij}; 0, 1/\sqrt{N})$ Sherrington - Kirkpatrick (1975)
disordered couplings

▷ Asymptotic computations

+ concentration on typical states

$$\lim_{d \rightarrow \infty} \mathbb{E}_J \left[\frac{1}{d} \log \mathcal{Z}_J \right] \approx \frac{1}{d} \log \mathcal{Z}_J \quad d \gg 1$$

disorder average

Replica computation, variational mean-field methods, high-temperature expansion ...

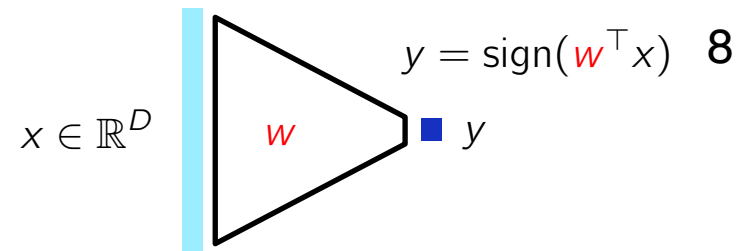
Physics: Edwards & Anderson (1975), Thouless Anderson Palmer (1977) Mézard, Parisi, Virasoro (1986), etc ...

Mathematics: Talagrand (2006), Panchenko (2014), Subag (2017) etc ...

▷ Application to learning since(80s-90s)

Amit, Gutfreund & Sompolinsky (1985), Gardner (1987), Derrida & Nadal (1987), Peterson & Anderson (1987), Krauth & Mézard (1987), Györgyi (1990). Opper, M Haussler, D. (1991) etc ...

Perceptron's capacity reloaded: Statistical physics computation



▷ Volume of weight parameter space satisfying the labelling

$$V_{D,N} \propto \lim_{\beta \rightarrow +\infty} \int_{\Omega(w)} d\mathbf{w} e^{-\beta E(\mathbf{w}, \{x^{(k)}, y^{(k)}\})}$$

disorder variables?

▷ Asymptotic computation using the replica method:

○ Fix scaling $\alpha = \frac{\# \text{ training points}}{\# \text{ parameters}} = \frac{N}{D}$

$$V(\alpha) = \lim_{N \rightarrow +\infty} \mathbb{E}_{x,y} [V_{D,M}]$$

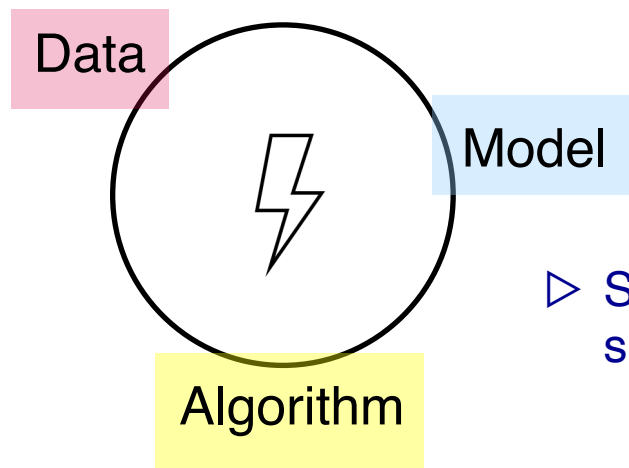
○ Fix distribution of data/disorder

$$V(\alpha_c) = 0$$

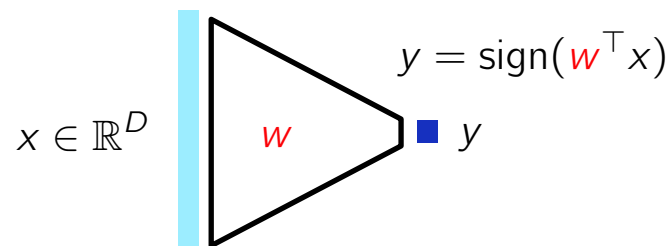
Gaussian inputs	$p(x_i^{(k)}) = \mathcal{N}(x_i; 0, 1)$	$\alpha_c = 2$
Binary inputs	$p(x_i^{(k)}) \pm 1$, binary weights $w_i \in \{-1, 1\}$	$\alpha_c \approx 0.83$

Introduction recap

- ▶ Take into account 3 elements to understand the successes of deep learning

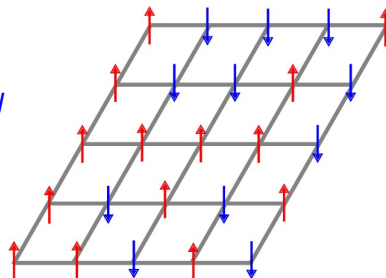


- ▶ Statistical mechanics of learning focuses on simple solvable models



- ▶ With the help of methods developed for disordered systems

$$\lim_{d \rightarrow \infty} \mathbb{E}_J \left[\frac{1}{d} \log \mathcal{Z}_J \right] \approx \frac{1}{d} \log \mathcal{Z}_J$$



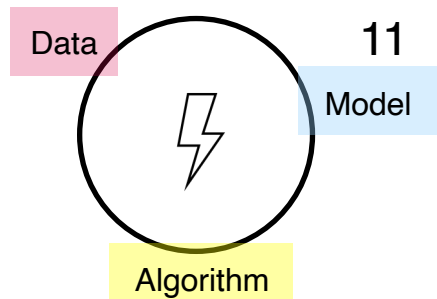
first example:
Compute the capacity of perceptron,
how many samples randomly sampled
can it classify?

Scope of today's lecture

Discuss classical and recent literature to give you an idea of what can be studied in machine learning with the statistical mechanics point of view.

- ▷ The teacher-student paradigm
- ▷ Models of data structure
- ▷ Dynamics of learning

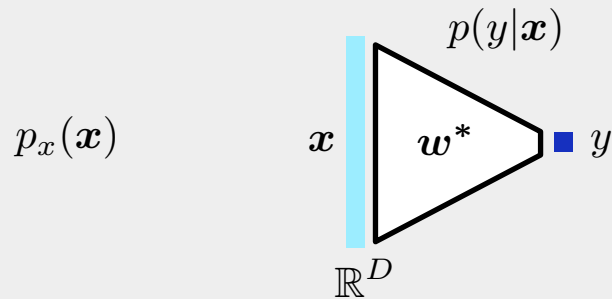
Teacher-Student paradigm: A toy model to study generalization



▷ Data-generating model: Teacher

Data

input distribution + input-output rule



ground truth parameter distribution $p_w(\mathbf{w}^*)$

training data

$$\mathcal{D} = \{y^{(k)}, \mathbf{x}^{(k)}\}_{k=1}^N$$

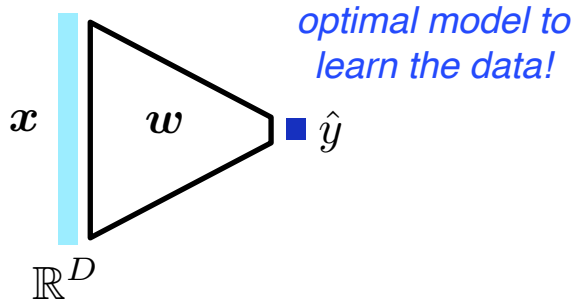
- Can we fit the data?
- Can we recover the teacher rule?

1. Choose a model
2. Choose an algorithm

▷ Learning model: Student

Model

ex: matching teacher



▷ Learning strategy?

Algorithm

- Loss minimization $\min_{\mathbf{w} \in \mathbb{R}^D} \sum_{k=1}^N \ell(y^{(k)}, \hat{y}(\mathbf{x}^{(k)}, \mathbf{w}))$

- Bayesian posterior

$$p(\mathbf{w} | \mathcal{D}) \propto \prod_{k=1}^N p(y^{(k)} | \mathbf{w}, \mathbf{x}^{(k)}) p(\mathbf{w})$$

Given the training data and a priori on the student, what information do we have about the student?

Bayes optimal generalization error (perceptron + Gaussian or binary input)

- ▷ When student model = teacher model, the setting is Bayes optimal

assuming square loss

- ▷ Bayes optimal generalization error: $\mathcal{E}_g = \mathbb{E}_{w|\mathcal{D}} \mathbb{E}_{x,y} [(y - \hat{y}(x, w))^2]$

« Mean error over the data if drawing the student parameters from the posterior distribution »

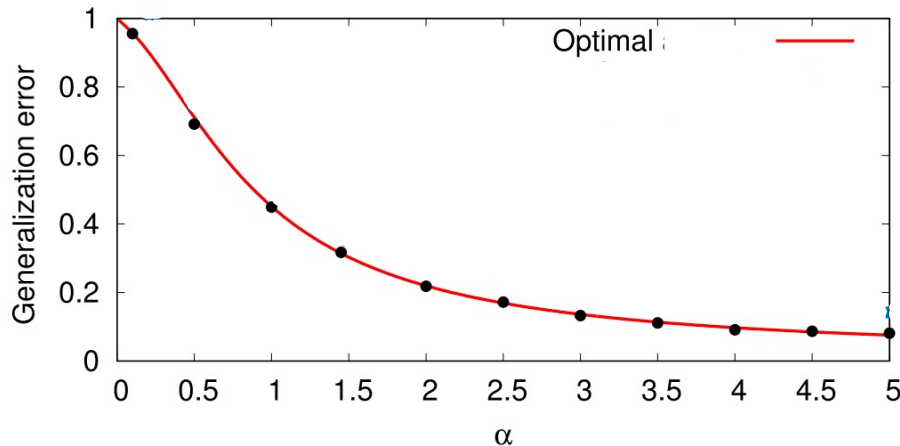
- ▷ Using the same tool as discussed before

(disorder average, thermodynamic limit, replica computation)

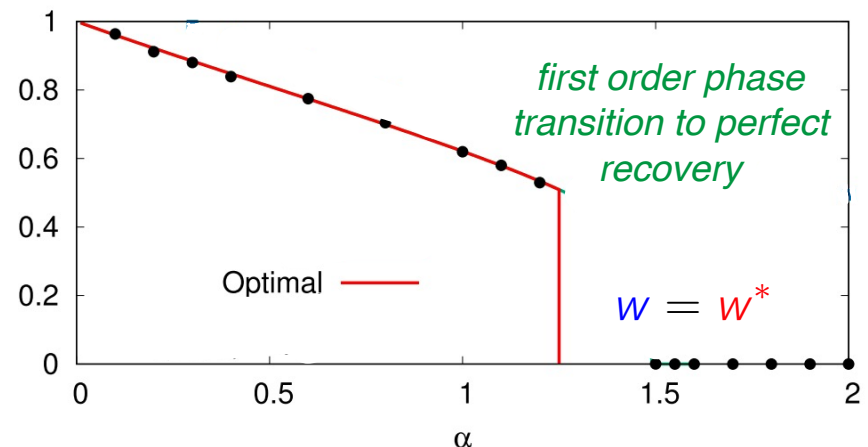
$$\lim_{N \rightarrow +\infty} \mathbb{E}_{\mathcal{D}} \left[\mathcal{E}_g(\mathcal{D}) \right] \xrightarrow{N \rightarrow \infty} \mathcal{E}_g(\alpha)$$

$$\alpha = \frac{\# \text{ training points}}{\# \text{ parameters}} = \frac{N}{D}$$

$w_i \in \mathbb{R}$ sparse



$w_i \pm 1$



Györgyi, G. (1990). *First-order transition to perfect generalization in a neural network with binary synapses*

Krauth, W., & Mézard, M. (1989). *Storage capacity of memory networks with binary couplings*

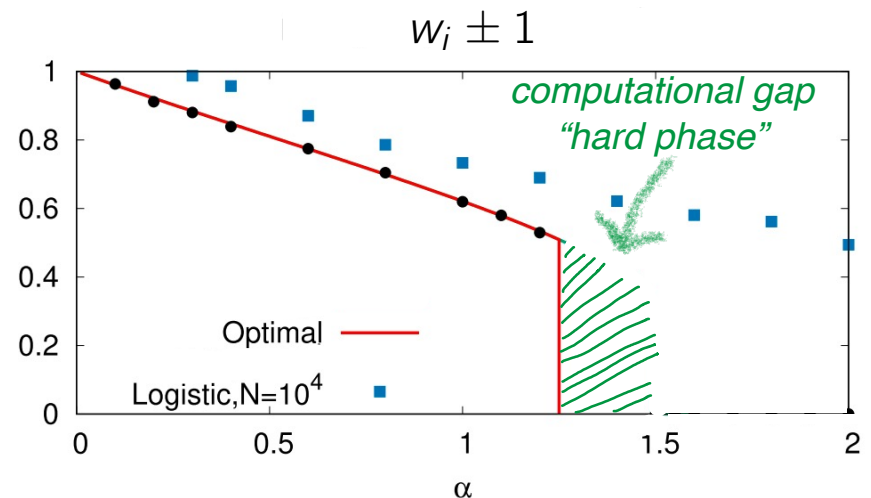
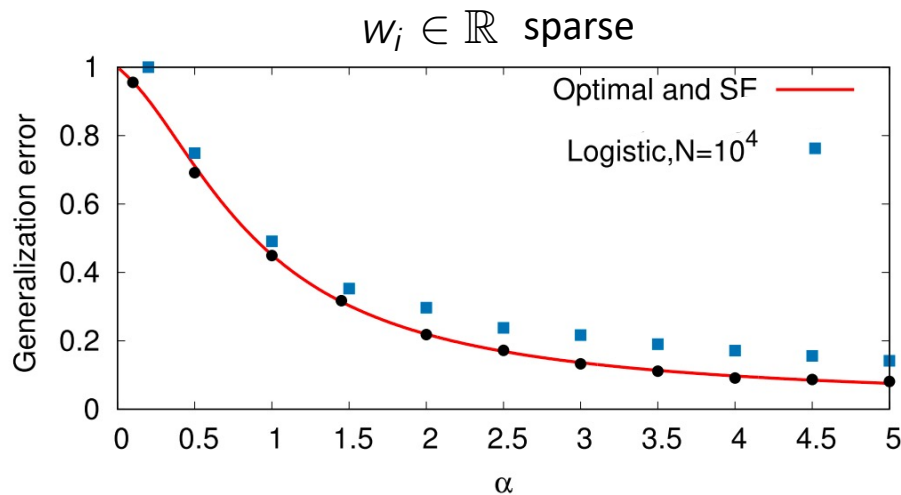
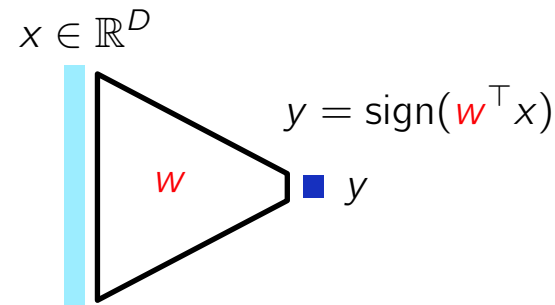
Barbier, J. et al. (2018). *Phase Transitions, Optimal Errors and Optimality of Message-Passing in Generalized Linear Models*

What about other (practical) algorithms?

How to train in practice?

▷ Gradient descent Algorithm

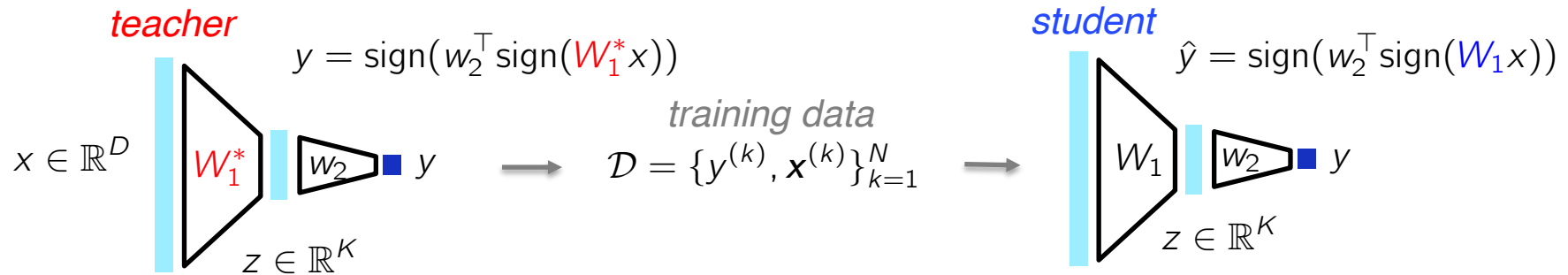
For instance binary classification problem
e.g. Logistic regression



▷ Message passing algorithms (among them belief propagation) Algorithm

- Related to some mean-field computation methods
- Performance predicted with the same tools (think disorder averaged) **state evolution (SE)**
- **Optimal among polynomial time algorithms for perceptron → computational gap**

A first 2-layer architecture: committee machines



▷ High-dimensional scaling $D \rightarrow +\infty$

$$\alpha = \frac{\# \text{ training points}}{\# \text{ parameters}} = \frac{N}{D} = O(1)$$

$$K = \# \text{ hidden layer units} = O(1)$$

▷ Bayes optimal generalization

$$\lim_{N \rightarrow +\infty} \mathbb{E}_{\mathcal{D}} \left[\mathcal{E}_g(\mathcal{D}) \right] \xrightarrow{N \rightarrow \infty} \mathcal{E}_g(\alpha)$$

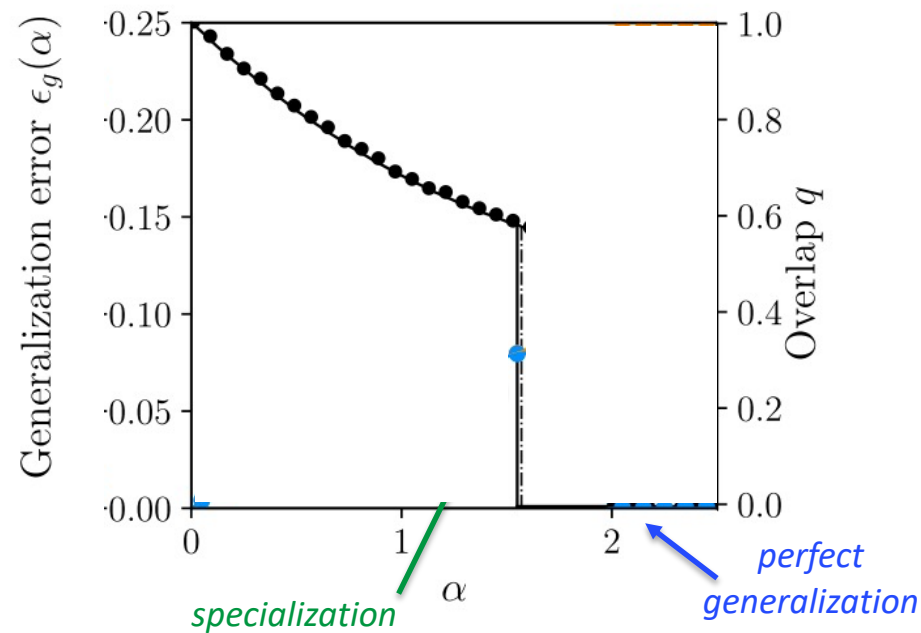
▷ Message passing algorithm

specialization transition!

$K = 2$, binary weights

$$q_{00} = \langle (W_1^*)_{0,\cdot}; (W_1)_{0,\cdot} \rangle$$

$$q_{01} = \langle (W_1^*)_{0,\cdot}; (W_1)_{1,\cdot} \rangle$$



Schwarze (1993). *Learning a Rule in a Multilayer Neural-Network.*

Schwarze & Hertz (1993). *Generalization in Fully Connected Committee Machines.*

Monasson et al (2004). *Learning and Generalization Theories of Large Committee-Machines*

Aubin et al (2018). *The committee machine: Computational to statistical gaps in learning a two-layers neural network*

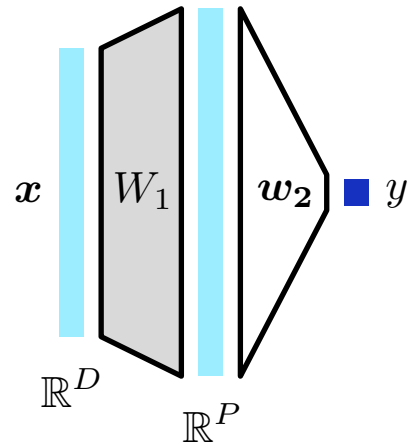
Another multi-layer example: Random features + not matched student and teacher

▷ Student: Alike a 2-layer neural network

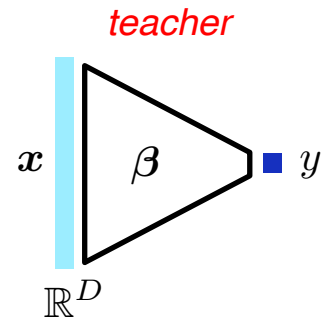
Model *student*

Rahimi, A. and Recht, B. *NeurIPS* 2018

- Learn a map $\hat{y} = \hat{f}(\mathbf{w}_2^\top \sigma(W_1 \mathbf{x}))$
- Weight parameters $\mathbf{w}_2 \in \mathbb{R}^P, W_1 \in \mathbb{R}^{P \times D}$
- First weight matrix fixed and random (e.g., i.i.d. Gaussian entries)



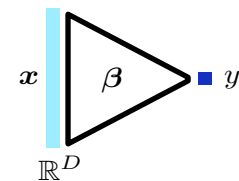
▷ Teacher: **Data** perceptron $y = f_0(\beta^\top \mathbf{x})$



▷ Training objective $\min_{\mathbf{w} \in \mathbb{R}^P} \sum_{\mu=1}^N \ell(y^\mu, \mathbf{x}^\mu, \mathbf{w}) + \frac{\lambda}{2} \|\mathbf{w}\|_2^2$

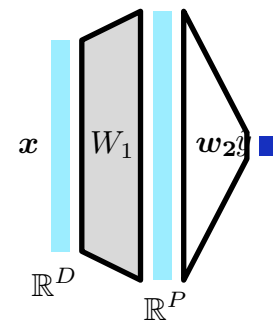
▷ Gibbs posterior formulation
$$\begin{aligned} \mu_\beta(\mathbf{w} \mid \{\mathbf{x}^\mu, y^\mu\}) &= \frac{1}{\mathcal{Z}_\beta} e^{-\beta [\sum_{\mu=1}^N \ell(y^\mu, \mathbf{x}^\mu \cdot \mathbf{w}) + \frac{\lambda}{2} \|\mathbf{w}\|_2^2]} \\ &= \frac{1}{\mathcal{Z}_\beta} \underbrace{\prod_{\mu=1}^N e^{-\beta \ell(y^\mu, \mathbf{x}^\mu \cdot \mathbf{w})}}_{\equiv P_y(\mathbf{y} \mid \mathbf{w}, \mathbf{x}^\mu)} \underbrace{\prod_{i=1}^P e^{-\frac{\beta \lambda}{2} w_i^2}}_{\equiv P_w(\mathbf{w})} \end{aligned}$$

An important intermediary result: Gaussian equivalence principle



$$y = f^0(\beta^\top x)$$

▷ Generalization error in the high-dimensional limit



$$\hat{y} = \hat{f}(w_2^\top \sigma(W_1 x))$$

Gaussian input	$x \sim \mathcal{N}(0, I_D)$	
input dimension	$D \rightarrow \infty$	$N/D = \alpha = O(1)$
# parameters	$P \rightarrow \infty$	$D/P = \gamma = O(1)$
# training samples	$N \rightarrow \infty$	

$$\lim_{N \rightarrow \infty} \mathcal{E}_g = \mathbb{E}_{\lambda, \nu} \left[\left(f^0(\nu) - \hat{f}(\lambda) \right)^2 \right]$$

$$\text{“}\nu = \beta^\top x\text{”}$$

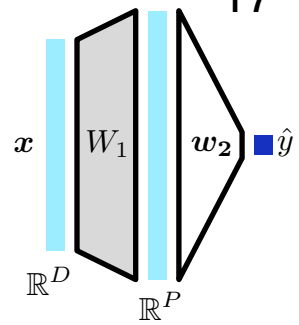
where the scalar fields are jointly Gaussian.

$$\text{“}\lambda = w_2^\top \sigma(W_1 x)\text{”}$$

▷ Compute covariance of Gaussian fields?

- random matrix theory tools or replica method from statistical physics!

Generalization in the random features model



$$\min_{\mathbf{w} \in \mathbb{R}^P} \sum_{\mu=1}^N \frac{1}{2} (y^\mu - \mathbf{w}^\top \sigma(W_1 \mathbf{x}^\mu))^2 + \frac{\lambda}{2} \|\mathbf{w}\|_2^2$$

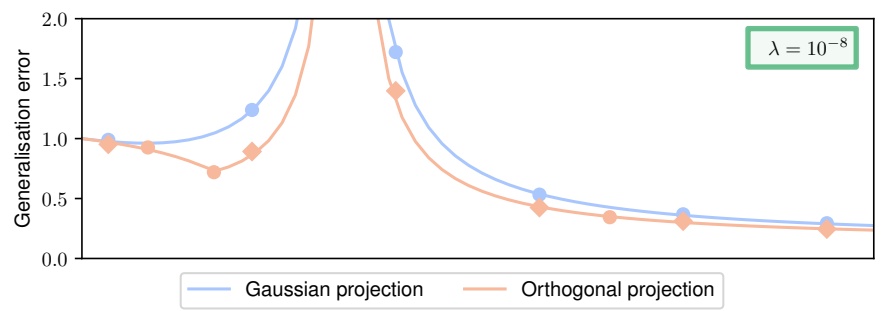
Regression teacher Data

$$y = \boldsymbol{\beta}^\top \mathbf{x} + \epsilon \quad \text{with } \epsilon \sim \mathcal{N}(0, \Delta)$$

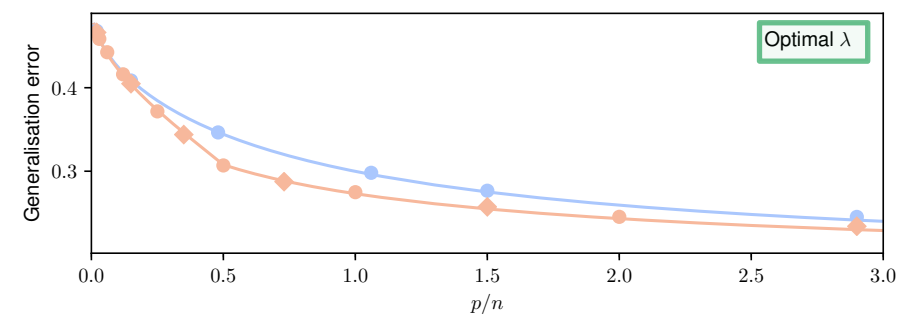
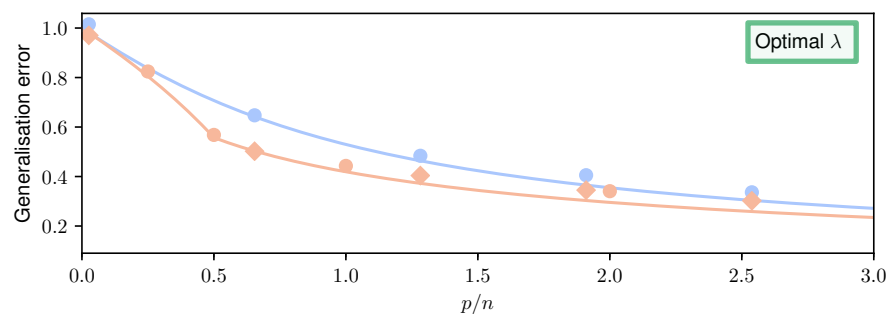
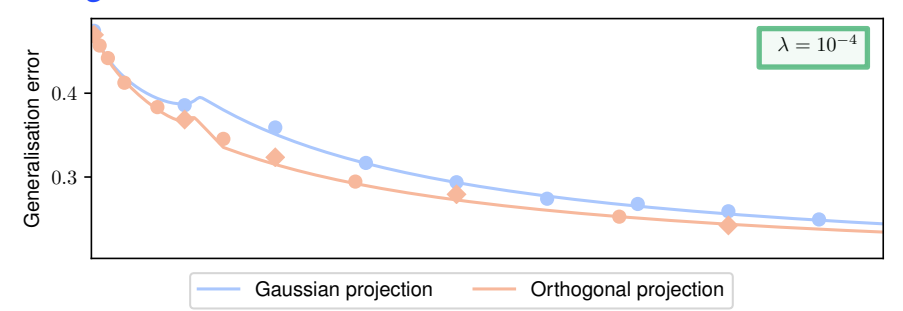
Classification teacher Data

$$y = \begin{cases} \text{sign}(\boldsymbol{\beta}^\top \mathbf{x}) & \text{with prob. } \Delta \\ -\text{sign}(\boldsymbol{\beta}^\top \mathbf{x}) & \text{with prob. } 1 - \Delta \end{cases}$$

double descent!



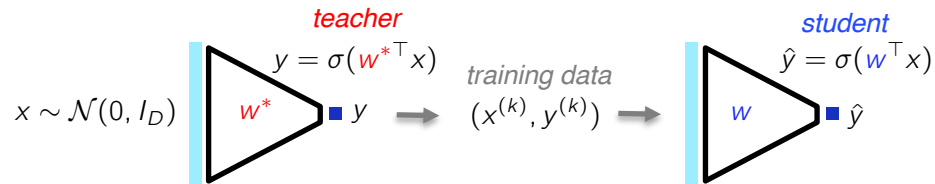
slight double descent



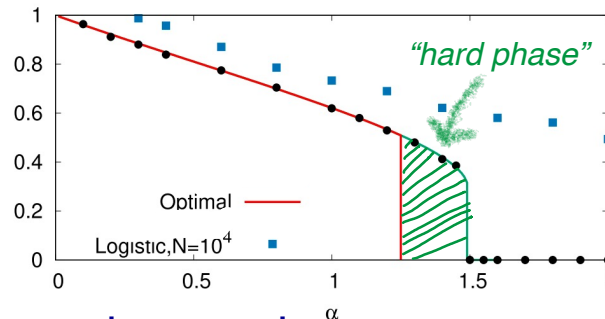
double descent cancelled by appropriate regularization

The teacher-student paradigm: Recap

- ▷ The teacher-student scenario consist in creating simple synthetic problem of learning



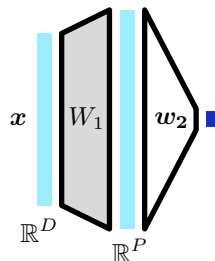
- ▷ Matched teacher and student: **Bayes optimal setting.** Information theoretic limit is not always achievable in reasonable time: **computational gaps.**



examples:
perceptron +
committee machine (2-layer)

- ▷ In mis-matched case we can observe the double descent phenomenon in overparametrized model.

example:
perceptron teacher w.
random feature student



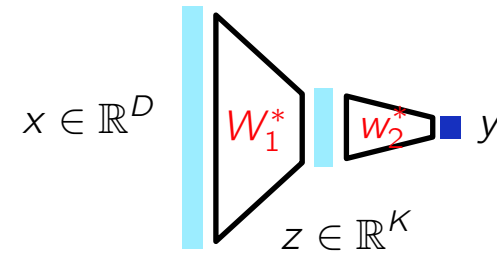
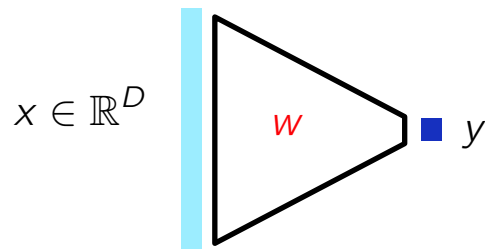
QUESTIONS?

Scope of today's lecture

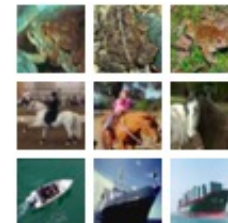
Discuss classical and recent literature to give you an idea of what can be studied in machine learning with the statistical mechanics point of view.

- ▷ The teacher-student paradigm
- ▷ Models of data structure
- ▷ Dynamics of learning

- ▷ Structure in between the input and outputs
(e.g. Neural network teacher models)



- ▷ Structure in the input



solvability /
interpretability



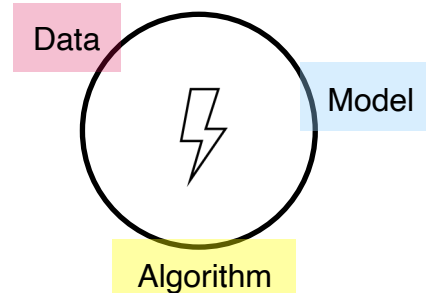
white noise

?

CIFAR images

degree of structure /
realism

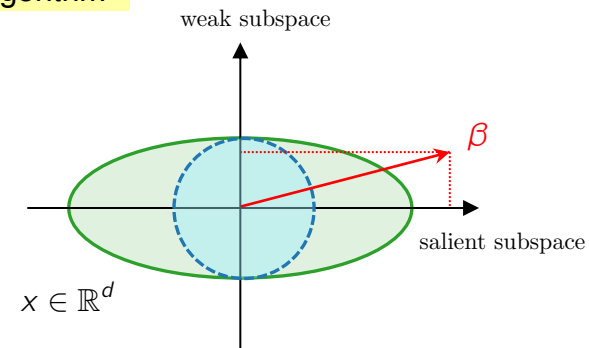
Simple model for data structure: Salient and weak features model



▷ Structure in the inputs Data $x \sim \mathcal{N}(0, \Sigma_x)$

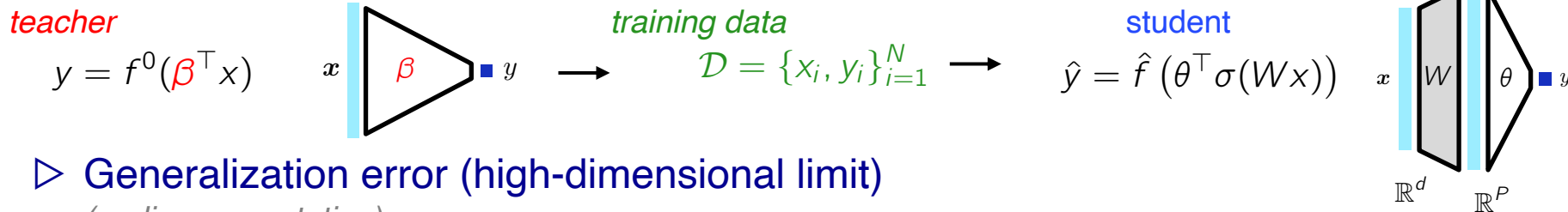
2 cases: - isotropiques $\Sigma_x = I_d$

- anisotropique structurée $\Sigma_x = \begin{bmatrix} \sigma_{x,1} I_{\phi_1 d} & 0 \\ 0 & \sigma_{x,2} I_{\phi_2 d} \end{bmatrix}$



▷ Structure input-output (teacher) and learning model (student)

(as we just discussed)

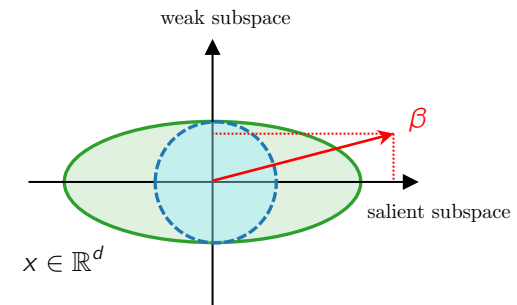
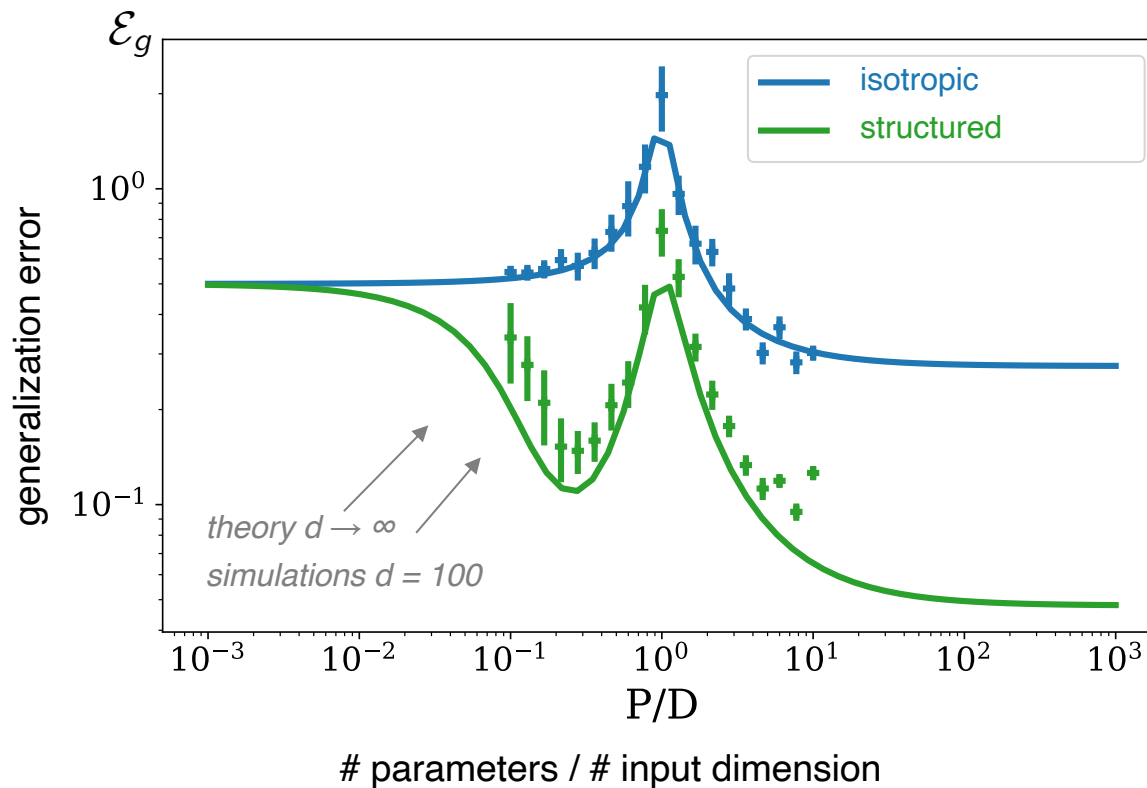


▷ Generalization error (high-dimensional limit)

(replica computation)

$$\mathcal{E}_g = \mathbb{E}_x \left[(y(x) - \hat{y}(x; \theta))^2 \right] \approx \lim_{d \rightarrow \infty} \mathbb{E}_\beta \mathbb{E}_D [\mathcal{E}_g] = \mathcal{E}_{g,\ell}(\alpha, \gamma, \sigma_{x,1,2}, \sigma_{\beta,1,2})$$

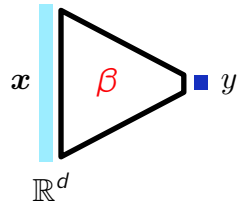
$$(d \rightarrow \infty, P \rightarrow \infty, N \rightarrow \infty, \gamma = P/d = O(1), \alpha = N/d = O(1))$$



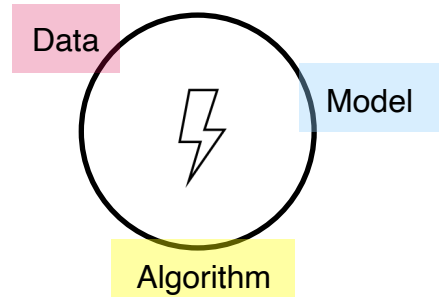
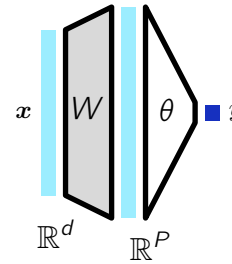
- ▷ The structure in the data is detected during training before the overfitting peak
- ▷ Double descent is exacerbated by the data structure when aligned with the task

Interplay between data structure and loss function in classification 1/2

teacher
 $y = \text{sign}(\beta^\top x)$

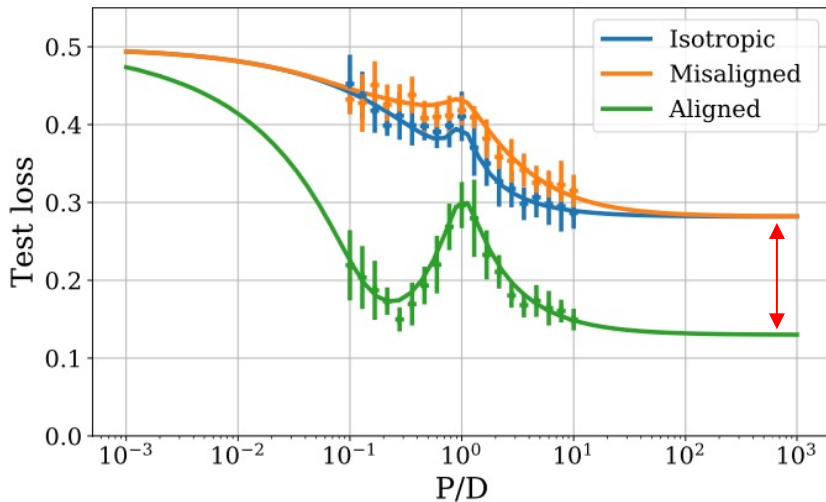


student
 $\hat{y} = \hat{f}(\theta^\top \sigma(Wx))$

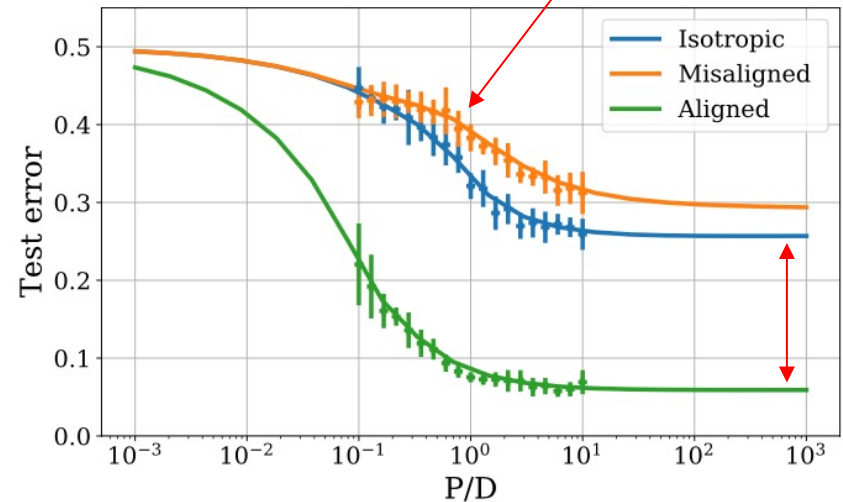


▷ Training objective $\min_{\theta \in \mathbb{R}^p} \sum_{\mu=1}^N \ell(y^\mu, x^\mu, \theta) + \frac{\lambda}{2} \|\theta\|_2^2$

square loss

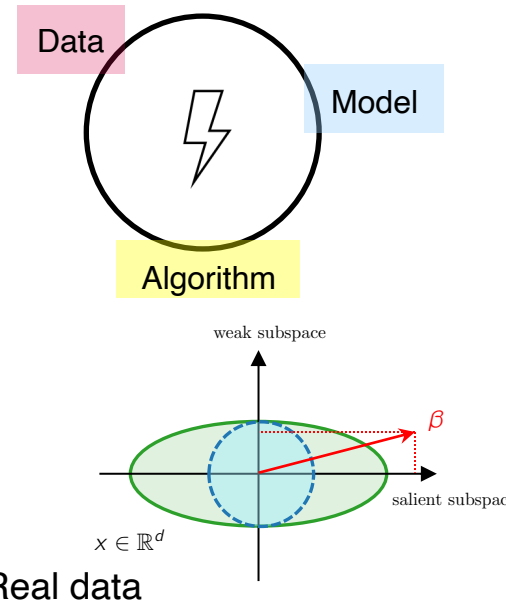
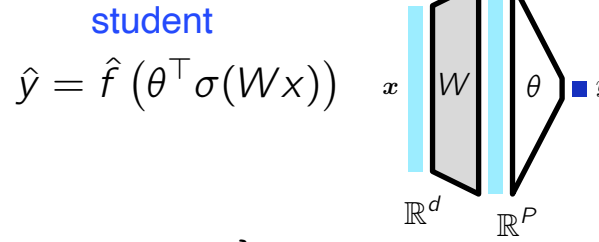
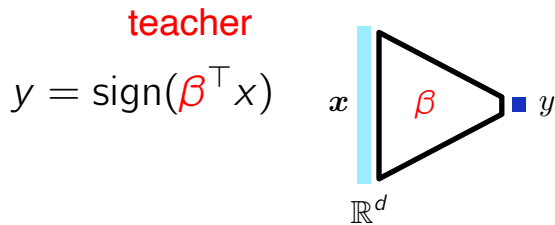


logistic loss



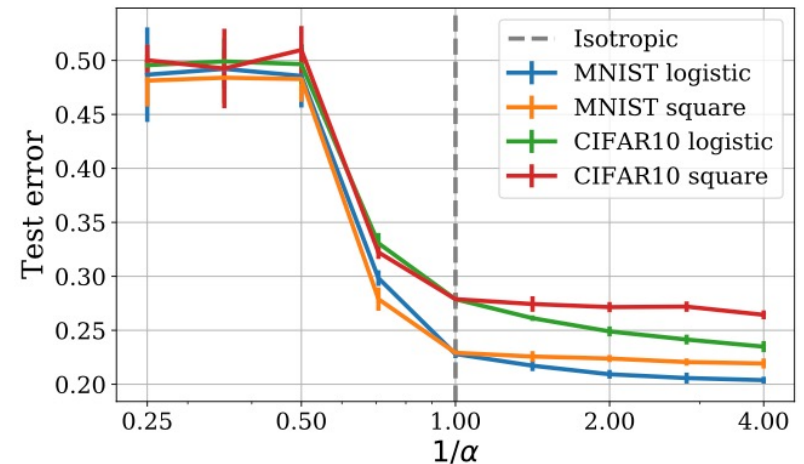
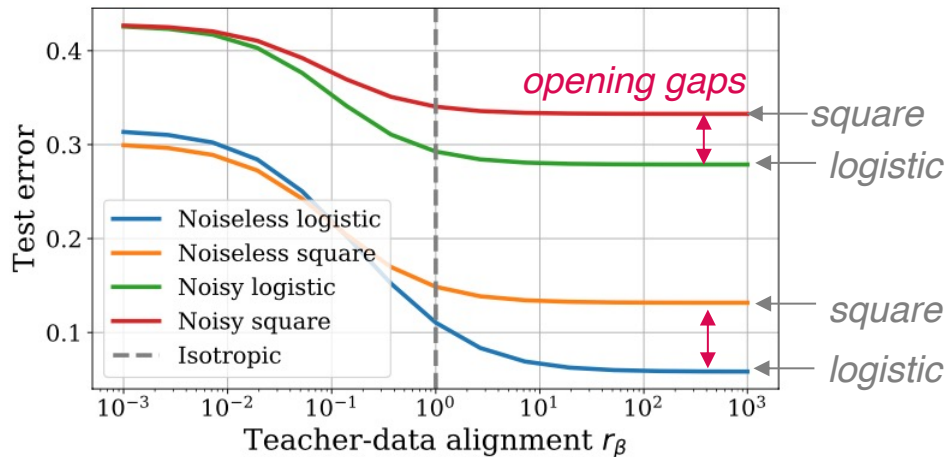
logistic loss takes most advantage of structure

Interplay between data structure and loss function in classification 2/2



▷ **Training objective**

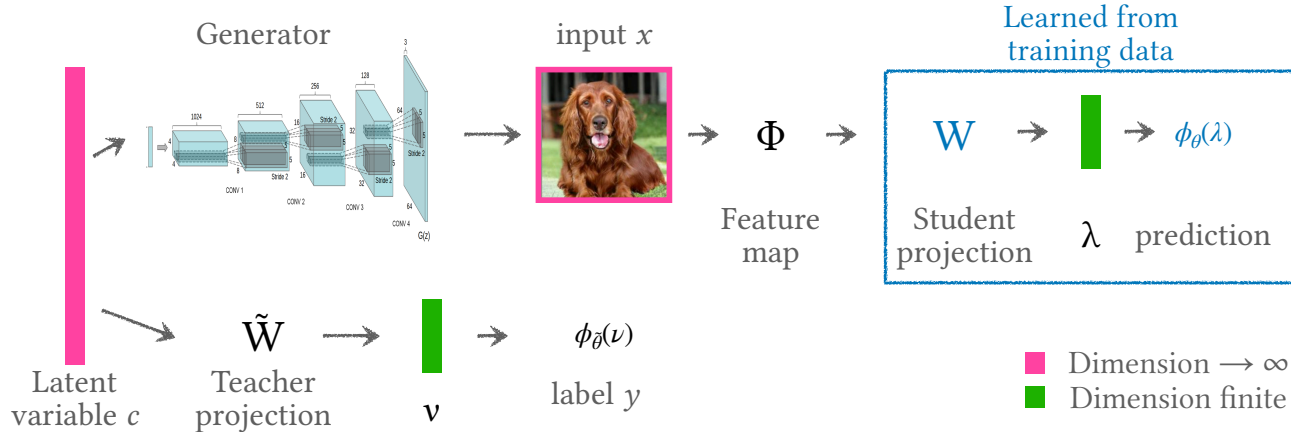
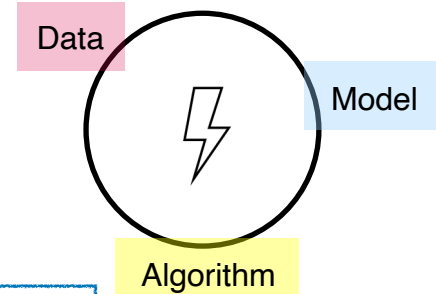
$$\min_{\theta \in \mathbb{R}^P} \sum_{\mu=1}^N \ell(y^\mu, x^\mu, \theta) + \frac{\lambda}{2} \|\theta\|_2^2$$



▷ Simple solvable model helps intuition further confirmed by real data experiments

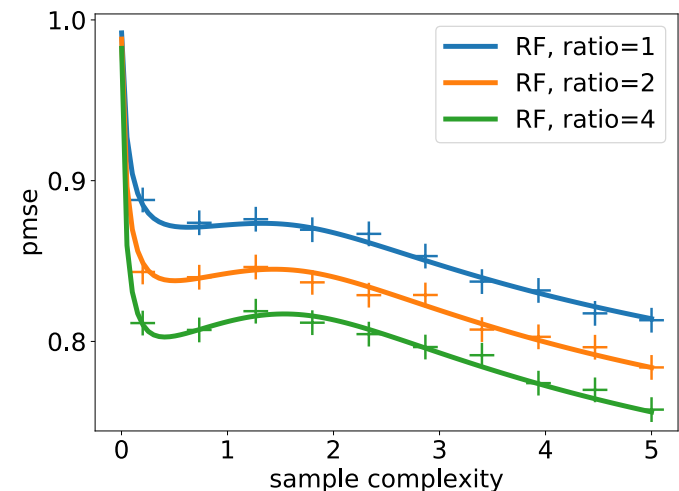
Towards more and more realistic data in the analysis: trained generative models

▷ Using pretrained generators in synthetic data models



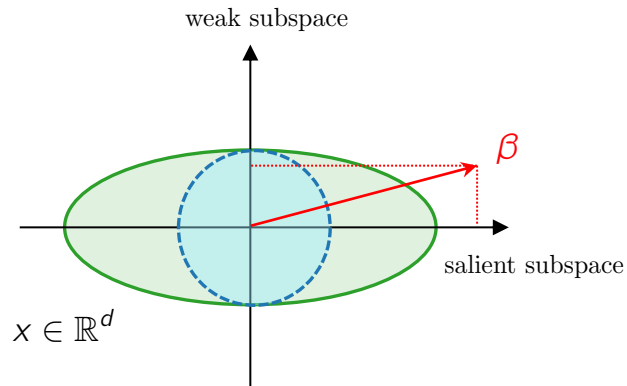
Gaussian equivalence principle here crucial!

▷ Example:
Predicting random features learning on CIFAR

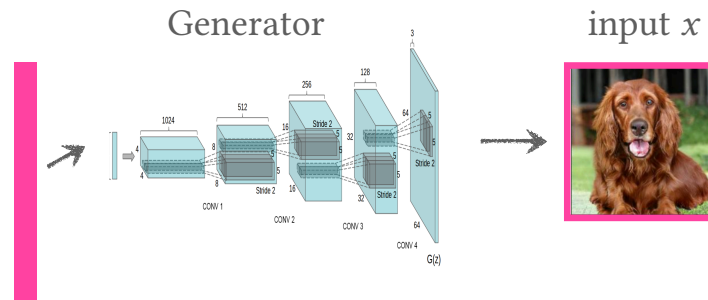


Models of data structures: Recap

- ▶ With simple models of data structures we can start to understand its impact on generalization

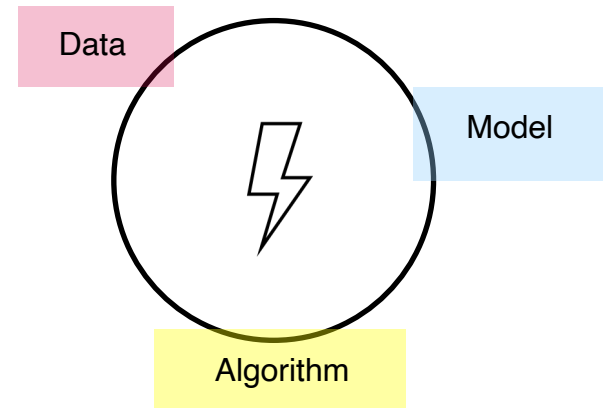


- ▶ Synthetic (trained) model can mimic real data and remain amenable to analysis



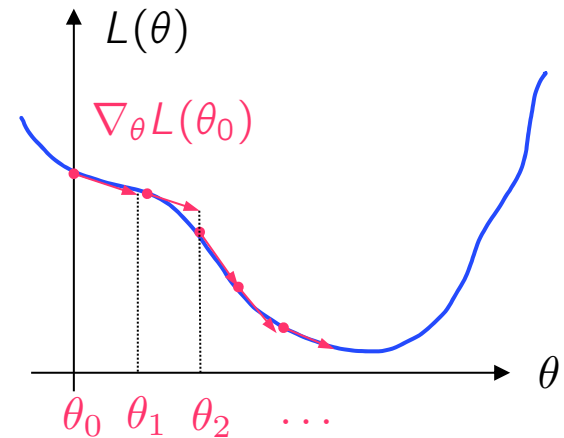
QUESTIONS?

More analysis of algorithms?



▷ (Stochastic) Gradient descent??

- Loss
$$L(\theta) = \frac{1}{N} \sum_{k=1}^N \ell[y^{(k)}, \hat{y}(x^{(k)}, \theta)]$$
- Gradient update
$$\theta^{t+1} \leftarrow \theta^t - \eta \nabla_{\theta} L(\theta^t)$$



Scope of today's lecture

Discuss classical and recent literature to give you an idea of what can be studied in machine learning with the statistical mechanics point of view.

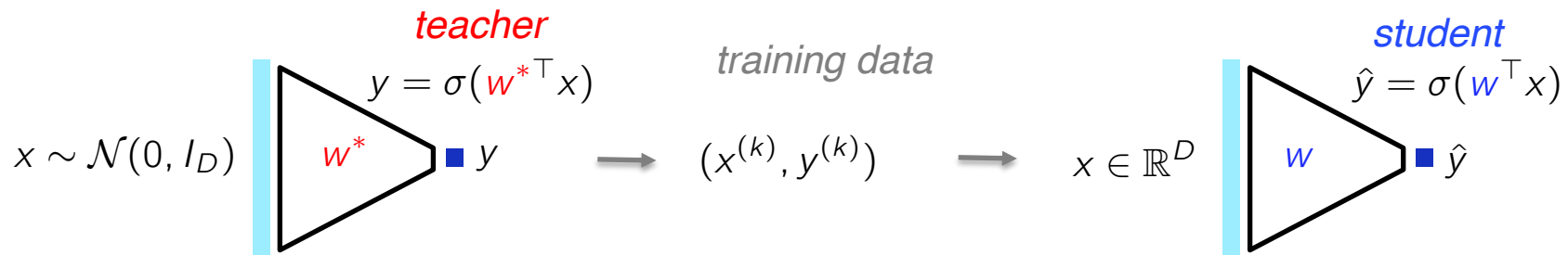
- ▷ The teacher-student paradigm
- ▷ Models of data structure
- ▷ Dynamics of learning

As soon as in the 90s: Analysis of online learning

▷ Online gradient descent

- Stream of data $(x^{(1)}, y^{(1)}), (x^{(2)}, y^{(2)}), \dots, (x^{(k)}, y^{(k)}), \dots$ *each sample seen once, size of dataset = # iterations*
- Parameter update with each sample $w^{k+1} = w^k - \eta \nabla_w L(w^k; (x^{(k)}, y^{(k)}))$

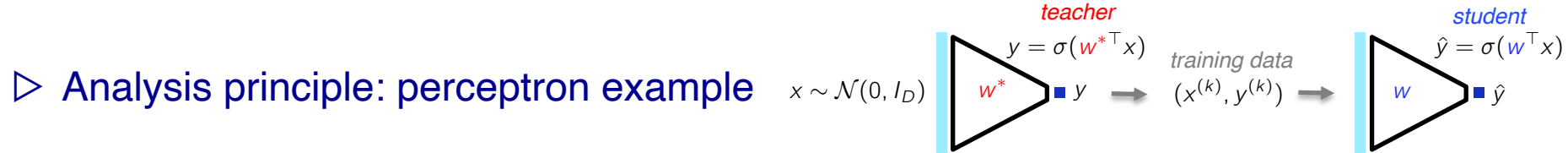
▷ Analysis principle: perceptron example



- Square loss for one sample $L(x; w, w^*) = \frac{1}{2} (y(x) - \hat{y}(x))^2 = \frac{1}{2} (\sigma(w^{*\top} x) - \sigma(w^\top x))^2$
- Parameter update $w^{t+1} - w^t = \eta (\sigma(w^{*\top} x) - \sigma(w^\top x)) \nabla \sigma(w^\top x)$
- Continuous time limit + high dimensional limit + disorder variable average

made possible by having uncorrelated samples at each update !

As soon as in the 90s: Analysis of online learning



- Follow training in terms of overlaps
 - $Q^{(k)} = w^{(k)\top} w^{(k)}$ « self overlap »
 - $R^{(k)} = w^{*\top} w^{(k)}$ « teacher-student overlap »
- Continuous time limit + high dimensional limit + disorder variable average

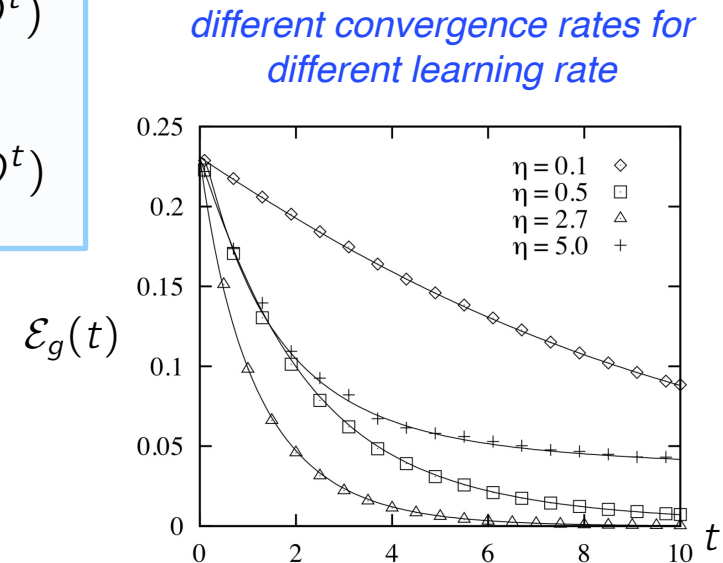
$$w^{t+1} - w^t = \eta dt (\sigma(w^{*\top} x) - \sigma(w^\top x)) \nabla \sigma(w^\top x)$$

→ closed set of equations on the overlaps

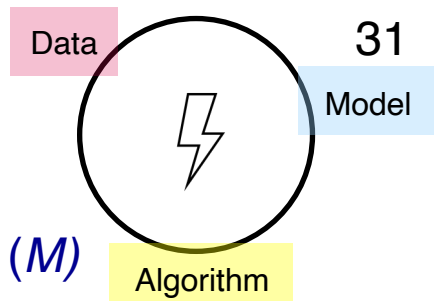
$$\begin{cases} \frac{dR}{dt} = g_R(R^t, Q^t) \\ \frac{dQ}{dt} = g_Q(R^t, Q^t) \end{cases}$$

- Generalization error:

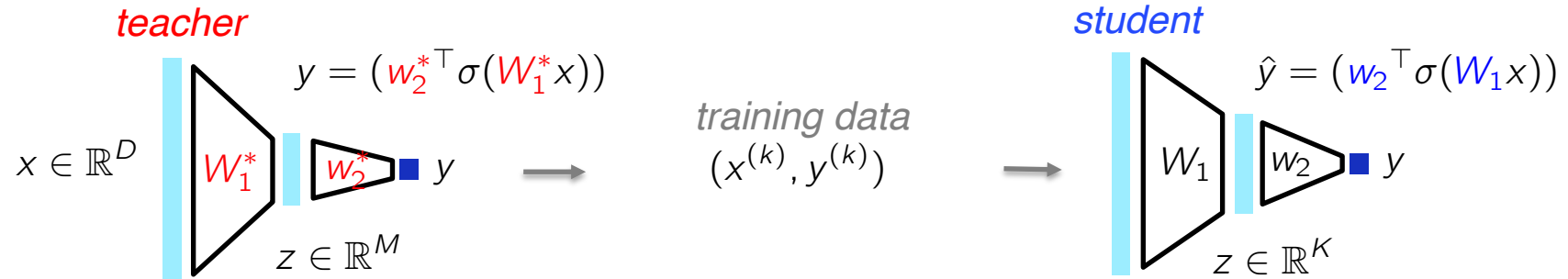
$$\begin{aligned} \mathcal{E}_g(w^{(k)}, w^*) &= \mathbb{E}_x \left[\frac{1}{2} \left(y^*(x) - \hat{y}(x, w^{(k)}) \right)^2 \right] \\ &= \mathcal{E}_g(R^{(k)}, Q^{(k)}) \xrightarrow{N \rightarrow +\infty, dt \rightarrow 0} \mathcal{E}_g(t) \end{aligned}$$



Committee machines online dynamics



- ▷ Two layer model with finite number of hidden units + possibly more/less hidden units in student (K) than teacher (M)



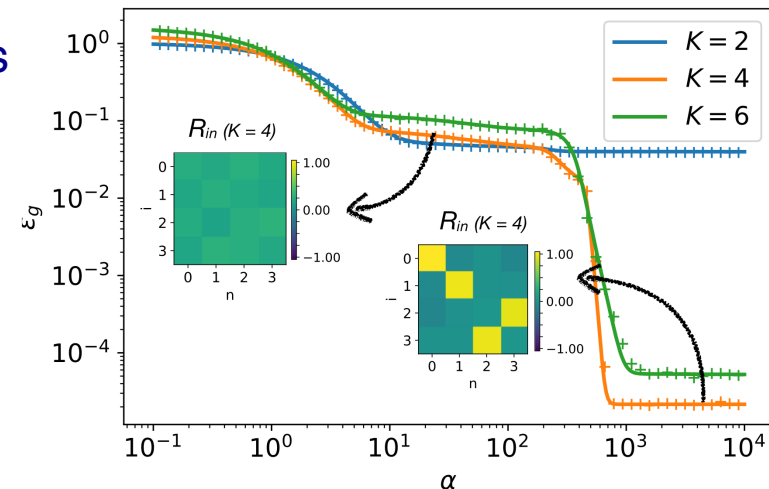
- ▷ Define overlaps again $R^{(k)} = W_1^* W_2^{(k)\top} \in \mathbb{R}^{K \times M}$
 $Q^{(k)} = W^{(k)} W^{(k)\top} \in \mathbb{R}^{K \times K}$

generalization $\mathcal{E}_g(t)$

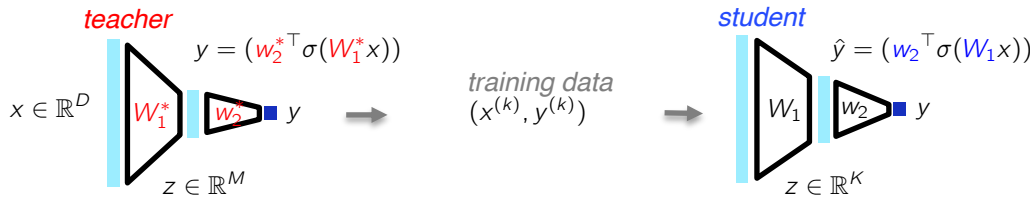
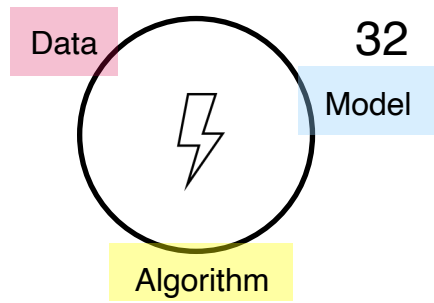
- ▷ First only learning first layer + matched models

$$W_2 = W_1$$

*specialization transition:
 student neurons match teacher neurones
 → escape plateau in dynamics*

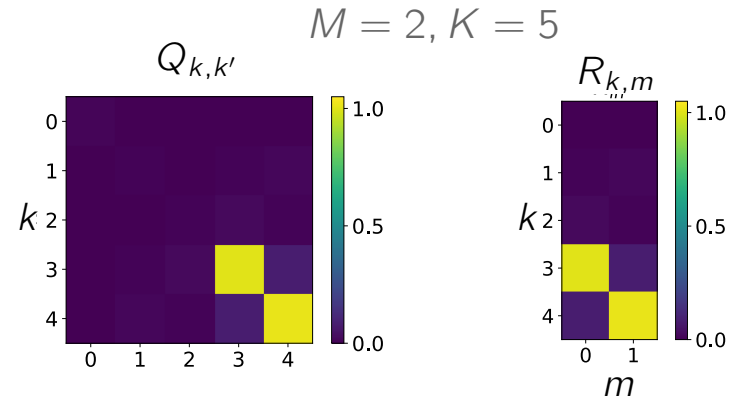
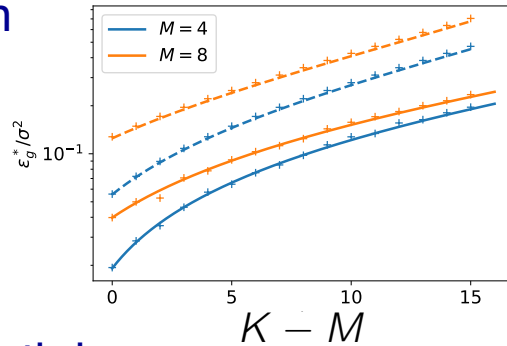


Committee machines online dynamics: overparametrization and feature learning



Over parametrization

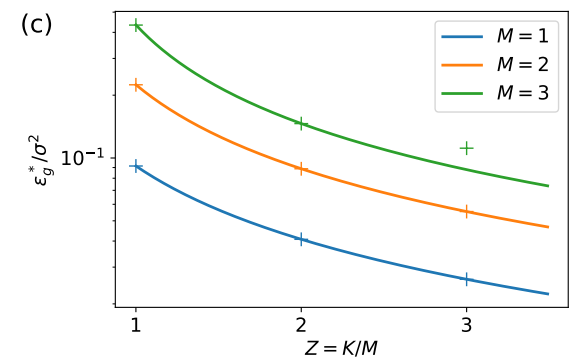
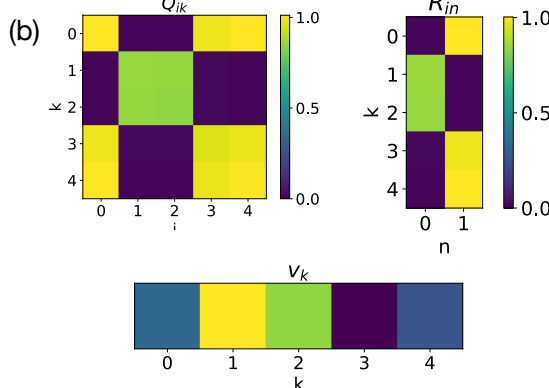
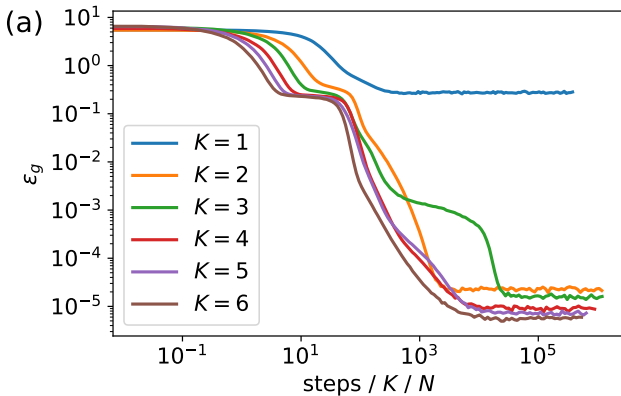
not helping if the second layer is not learned!



Feature learning = both layers

in the analysis

$$w_2^* = (1, 1, \dots, 1) \in \mathbb{R}^K$$



overparametrization beneficial by denoising!

Can we go beyond online gradient descent?

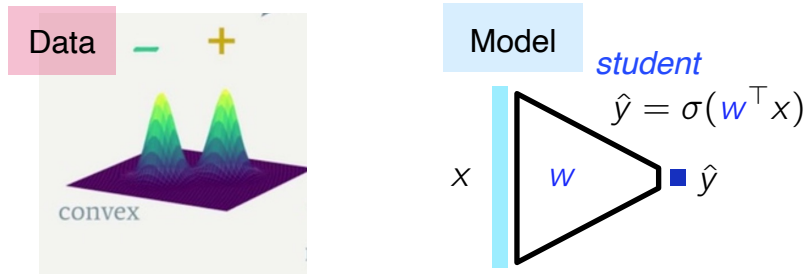
Analysis of multi-pass stochastic gradient descent

Involved but possible in simple models

- ▷ Use of dynamical mean-field theory (DMFT)

Mézar, Parisi, Virasoro (1987), Sompolinsky, Crisanti, Sommers (1988) etc...

- ▷ Example: Perceptron classification of Gaussian mixtures



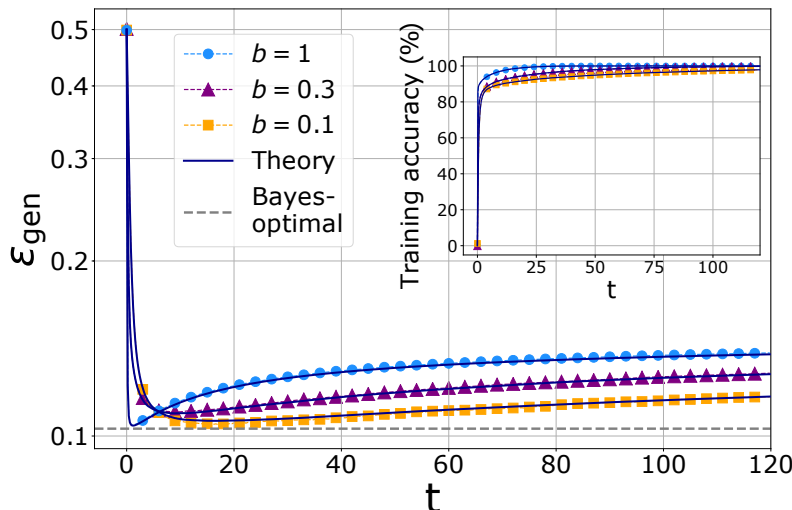
Algorithm

$$\min_w \sum_{i=1}^N L(y, \sigma(w^T x)) + \lambda \|w\|^2$$

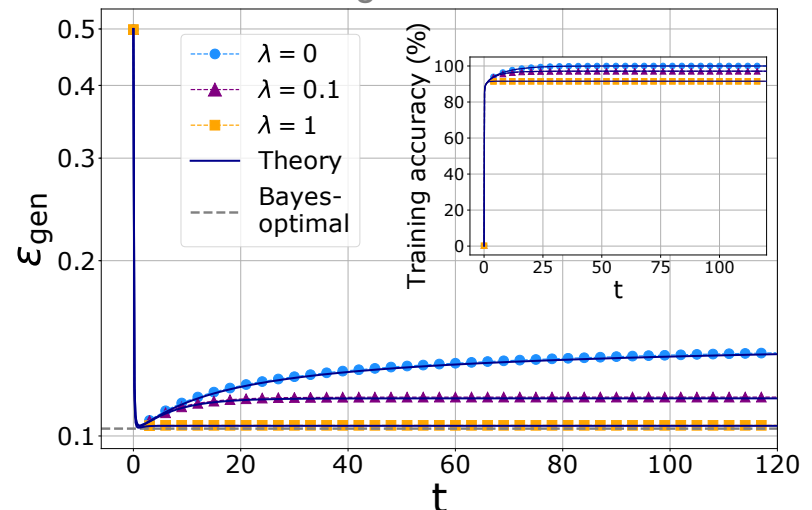
At each time sample in the mini-batch with probability b :

$$w^{t+1} - w^t = \eta dt \sum_{k=1}^N s_k \nabla L(y, \sigma(w^T x^{(k)}))$$

Effect of the batchsize



Effect of regularization

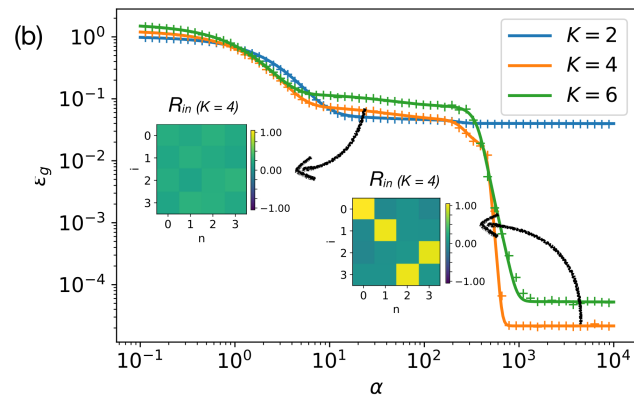
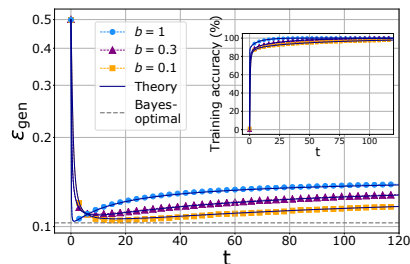
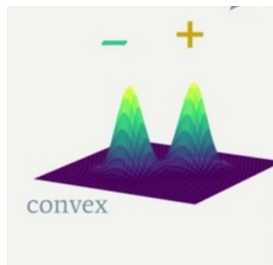


- ▷ The dynamics of online gradient descent can be studied in simple models by a set of closed ODEs on the overlaps

$$\begin{cases} \frac{dR}{dt} = g_R(R^t, Q^t) \\ \frac{dQ}{dt} = g_Q(R^t, Q^t) \end{cases} \quad \mathcal{E}_g(R^{(k)}, Q^{(k)}) \xrightarrow{N \rightarrow +\infty, dt \rightarrow 0} \mathcal{E}_g(t)$$

- ▷ In the committee machine training, we can describe the onset of the specialization transition

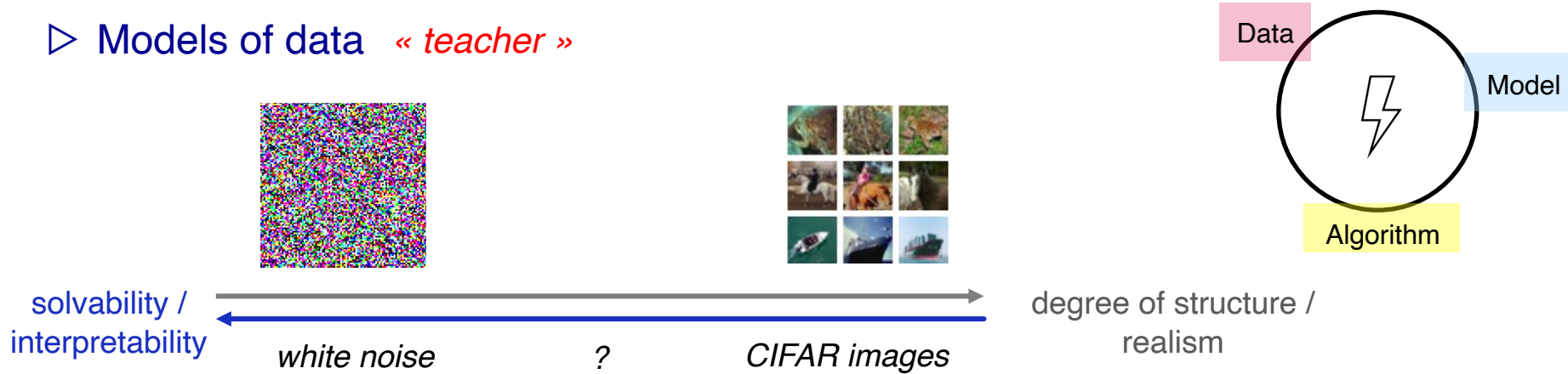
- ▷ Dynamical mean-field theory can help study batch gradient descent and understand the impact of batch-size



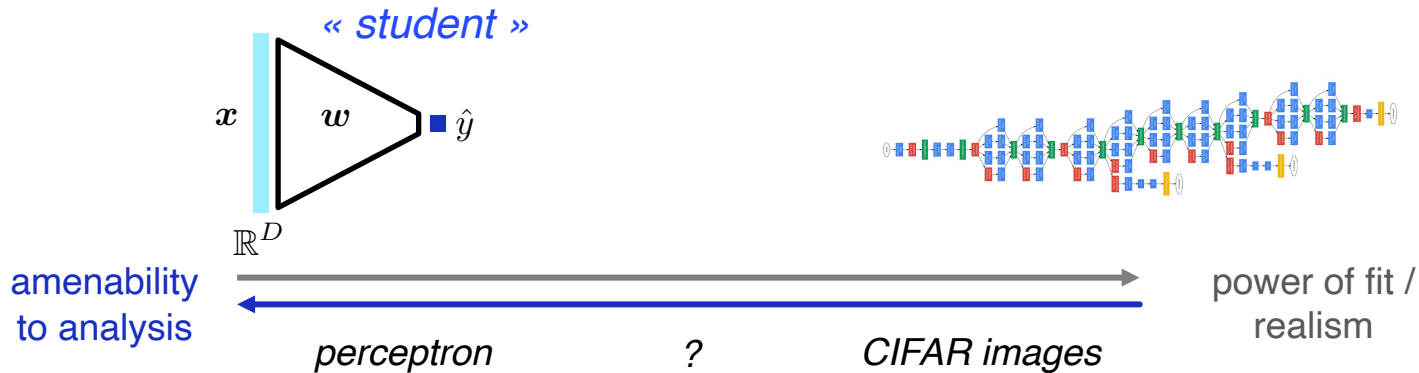
QUESTIONS?

Simple solvable models – what to focus on?

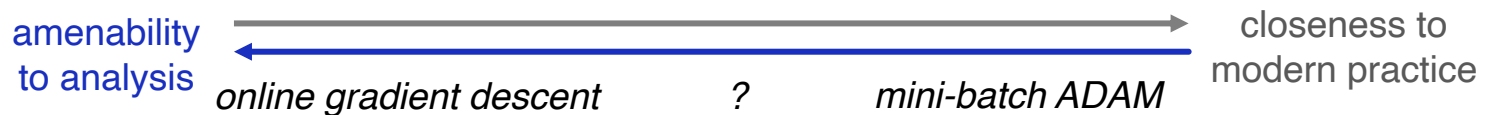
▷ Models of data « *teacher* »



▷ Architectures



▷ Algorithms



Broadening use of “mean-field approximations” in deep learning literature

▷ “take advantage of concentration effects thanks to randomness in the large size limit”

* Analysis of statistical inference performance

Reviews: - Zdeborová & Krzakala (2016) *Statistical physics of inference: Thresholds and algorithms*.
- **Gabrielé** (2020) *Mean field inference methods for neural networks*.

* Signal propagation in deep neural networks

- Trainability of very deep network at init. e.g., Schoenholz et al.(2017). *Deep Information Propagation*.
- Separation of structured data
e.g., Cohen, et al (2020). *Separability and geometry of object manifolds in deep neural networks*.

* Role of over-parametrization in trainability with Gradient Descent methods

- Convergence of SGD for 2-layers neural networks
Chizat & Bach (2018), Mei, Montanari & Nguyen (2018), Rotskoff & Vanden-Eijnden (2018)
- Neural Tangent Kernels, Equivalence to Gaussian processes, “Lazy training”
Jacot et al (2018), Lee et al (2019), review: Bahri et al (2020) *Statistical Mechanics of Deep Learning*
- Online learning e.g., Goldt, et al (2019). Dynamics of stochastic gradient descent for two-layer neural networks in the teacher-student setup

* Landscape, algorithms and interactions

Dauphin et al (2014). Identifying and attacking the saddle point problem in high-dimensional non-convex optimization
Sarao Mannelli & Zdeborová (2020). *Thresholds of descending algorithms in inference problems*.

▷ Review papers (among others):

- **classical** Watkin, T. L. H., Rau, A., & Biehl, M. (1993). The statistical mechanics of learning a rule
- **physics & modern M.L.** Carleo, G., et. al (2019). *Machine learning and the physical sciences*.
Gabrié, M. (2019). *Mean-field inference methods for neural networks*.
Bahri, Y., et al. (2020). *Statistical Mechanics of Deep Learning*
- **relevant statistical physics methods**
Zdeborová, L., & Krzakala, F. (2016). *Statistical physics of inference: Thresholds and algorithms*
Castellani, T., Cavagna, A., Fisica, D., & Moro, P. A. (2005). *Spin-Glass Theory for Pedestrians*

▷ Books:

- **classical stat. mech. of learning** Engel, A., & Van den Broeck, C. (2001). *Statistical Mechanics of Learning*.
Oppen, M., & Saad, D. (2001). *Advanced mean field methods: Theory and practice*.
- **spin glass** Mézard, M., Parisi, G., & Virasoro, M. (1986). *Spin Glass Theory and Beyond*
- **message passing** Mézard, M., & Montanari, A. (2009). *Information, Physics, and Computation*.

THANKS !

1. Mézard, M., Parisi, G., & Virasoro, M. (1986). *Spin Glass Theory and Beyond* (Vol. 9). WORLD SCIENTIFIC. <https://doi.org/10.1142/0271>
2. Gardner, E. (1987). Maximum Storage Capacity in Neural Networks. *Europhysics Letters (EPL)*, 4(4), 481–485. <https://doi.org/10.1209/0295-5075/4/4/016>
3. Gardner, E., & Derrida, B. (1989). Three unfinished works on the optimal storage capacity of networks. *Journal of Physics A: Mathematical and General*, 22(12), 1983–1994. <https://doi.org/10.1088/0305-4470/22/12/004>
4. Krauth, W., & Mézard, M. (1989). Storage capacity of memory networks with binary couplings. *Journal de Physique*, 50(20), 3057–3066. <https://doi.org/10.1051/jphys:0198900500200305700>
5. Györgyi, G. (1990). First-order transition to perfect generalization in a neural network with binary synapses. *Physical Review A*, 41(12), 7097–7100. <https://doi.org/10.1103/PhysRevA.41.7097>
6. Seung, H. S., Sompolinsky, H., & Tishby, N. (1992). Statistical mechanics of learning from examples. *Physical Review A*, 45(8), 6056–6091. <https://doi.org/10.1103/PhysRevA.45.6056>
7. Watkin, T. L. H., Rau, A., & Biehl, M. (1993). The statistical mechanics of learning a rule. *Reviews of Modern Physics*, 65(2), 499–556. <https://doi.org/10.1103/RevModPhys.65.499>
8. Biehl, M., & Riegler, P. (1994). On-Line Learning with a Perceptron. *Europhysics Letters (EPL)*, 28(7), 525–530. <https://doi.org/10.1209/0295-5075/28/7/012>
9. Saad, D., & Solla, S. A. (1995). On-line learning in soft committee machines. *Physical Review E*, 52(4), 4225–4243.
10. Saad, D., & Solla, S. A. (1995). Exact solution for on-line learning in multilayer neural networks. *Physical Review Letters*, 74(21), 4337–4340. <https://doi.org/10.1103/PhysRevLett.74.4337>
11. Biehl, M., & Schwarze, H. (1995). Learning by on-line gradient descent. *J. Phys. A. Math. Gen.*, 28(3), 643–656. <https://doi.org/10.1088/0305-4470/28/3/018>
12. Saad, D. (1999). *On-Line Learning in Neural Networks* (D. Saad (ed.)). Cambridge University Press. <https://doi.org/10.1017/CBO9780511569920>
13. Opper, M., & Saad, D. (2001). *Advanced mean field methods: Theory and practice*. MIT press.
14. Engel, A., & Van den Broeck, C. (2001). *Statistical Mechanics of Learning*. Cambridge University Press. <https://doi.org/10.1017/CBO9781139164542>

15. Castellani, T., Cavagna, A., Fisica, D., & Moro, P. A. (2005). Spin-Glass Theory for Pedestrians. *Journal of Statistical Mechanics: Theory and Experiment*, 5, 215–266. <https://doi.org/10.1088/1742-5468/2005/05/P05012>
16. Talagrand, M. (2006). The Parisi formula. *Annals of Mathematics*, 163(1), 221–263. <https://doi.org/10.4007/annals.2006.163.221>
17. Mézard, M., & Montanari, A. (2009). *Information, Physics, and Computation*. Oxford University Press.
18. Zdeborová, L., & Krzakala, F. (2016). Statistical physics of inference: Thresholds and algorithms. *Advances in Physics*, 65(5), 453–552. <https://doi.org/10.1080/00018732.2016.1211393>
19. Barbier, J., Dia, M., Macris, N., & Krzakala, F. (2017). The mutual information in random linear estimation. *54th Annual Allerton Conference on Communication, Control, and Computing, Allerton 2016*, 0(1), 625–632. <https://doi.org/10.1109/ALLERTON.2016.7852290>
20. Barbier, J., Krzakala, F., Macris, N., Miolane, L., & Zdeborová, L. (2018). Phase Transitions, Optimal Errors and Optimality of Message-Passing in Generalized Linear Models. *Proceedings of the 31st Conference On Learning Theory, PMLR 75*, 728–731. <http://arxiv.org/abs/1708.03395>
21. Gabrié, M., Manoel, A., Luneau, C., Barbier, Jean, Macris, N., Krzakala, F., & Zdeborová, L. (2018). Entropy and mutual information in models of deep neural networks. *Advances in Neural Information Processing Systems 31* (Issue Nips, pp. 1826--1836). Curran Associates, Inc. <http://papers.nips.cc/paper/7453-entropy-and-mutual-information-in-models-of-deep-neural-networks.pdf>
22. Barbier, J., Macris, N., Maillard, A., & Krzakala, F. (2018). The Mutual Information in Random Linear Estimation Beyond i.i.d. Matrices. *IEEE International Symposium on Information Theory - Proceedings, 2018-June(3)*, 1390–1394. <https://doi.org/10.1109/ISIT.2018.8437522>
23. Aubin, B., Maillard, A., Barbier, J., Krzakala, F., Macris, N., & Zdeborová, L. (2018). The committee machine: Computational to statistical gaps in learning a two-layers neural network. *Neural Information Processing Systems 2018, NeurIPS*, 1–44. <http://arxiv.org/abs/1806.05451>
24. Reeves, G. (2018). Additivity of information in multilayer networks via additive Gaussian noise transforms. *55th Annual Allerton Conference on Communication, Control, and Computing, Allerton 2017, 2018-Janua*, 1064–1070. <https://doi.org/10.1109/ALLERTON.2017.8262855>
25. Gabrié, M. (2019). Mean-field inference methods for neural networks. *ArXiv Preprint, 1911.00890*. <http://arxiv.org/abs/1911.00890>
26. Carleo, G., Cirac, I., Cranmer, K., Daudet, L., Schuld, M., Tishby, N., Vogt-Maranto, L., & Zdeborová, L. (2019). Machine learning and the physical sciences. *Reviews of Modern Physics*, 91(4). <https://doi.org/10.1103/revmodphys.91.045002>
27. Bahri, Y., Kadmon, J., Pennington, J., Schoenholz, S., Sohl-Dickstein, J., & Ganguli, S. (2020). Statistical Mechanics of Deep Learning. *Annual Review of Condensed Matter Physics*, 11(1), 501–528. <https://doi.org/10.1146/annurev-conmatphys-031119-050745>

▷ Replica trick

$$\mu_\beta(\mathbf{w} \mid \{\mathbf{x}^\mu, y^\mu\}) = \frac{1}{\mathcal{Z}_\beta} e^{-\beta[\sum_{\mu=1}^N \ell(y^\mu, \mathbf{x}^\mu \cdot \mathbf{w}) + \frac{\alpha}{2} \|\mathbf{w}\|_2^2]} \quad \mathbb{E}_{\mathcal{D}} \log \mathcal{Z}_\beta = \lim_{r \rightarrow 0^+} \frac{\partial}{\partial r} \mathbb{E}_{\mathcal{D}} \mathcal{Z}_\beta^r$$

▷ Rewritten in terms of order parameters

$$\mathbb{E}_{\{\mathbf{x}^\mu, y^\mu\}} \mathcal{Z}_\beta^r = \int \frac{d\rho d\hat{\rho}}{2\pi} \int \prod_{a=1}^r \frac{dm_s^a d\hat{m}_s^a}{2\pi} \int \prod_{1 \leq a \leq b \leq r} \frac{dq_s^{ab} d\hat{q}_s^{ab}}{2\pi} \frac{dq_w^{ab} d\hat{q}_w^{ab}}{2\pi} e^{\mathbf{P}\Phi^{(r)}}$$

diverging dimension ↙ replica potential ↘

▷ Saddle point equations

$$\begin{cases} \hat{r}_{s,i} = -2\sigma_{x,i} \frac{\alpha}{\gamma} \partial_{r_{s,i}} \Psi_y(R, Q, M) & r_{s,i} = -\frac{2}{\gamma} \partial_{\hat{r}_{s,i}} \Psi_w(\hat{r}_{s,i}, \hat{q}_{s,i}, \hat{m}_{s,i}, \hat{r}_w, \hat{q}_w) \\ \hat{q}_{s,i} = -2\sigma_{x,i} \frac{\alpha}{\gamma} \partial_{q_{s,i}} \Psi_y(R, Q, M) & q_{s,i} = -\frac{2}{\gamma} \partial_{\hat{q}_{s,i}} \Psi_w(\hat{r}_{s,i}, \hat{q}_{s,i}, \hat{m}_{s,i}, \hat{r}_w, \hat{q}_w) \\ \hat{m}_{s,i} = \sigma_{x,i} \frac{\alpha}{\gamma} \partial_{m_{s,i}} \Psi_y(R, Q, M) & m_{s,i} = \frac{1}{\gamma} \partial_{\hat{m}_{s,i}} \Psi_w(\hat{r}_{s,i}, \hat{q}_{s,i}, \hat{m}_{s,i}, \hat{r}_w, \hat{q}_w) \\ \hat{r}_w = -2\alpha \partial_{r_w} \Psi_y(R, Q, M) & r_w = -2\partial_{\hat{r}_w} \Psi_w(\hat{r}_{s,i}, \hat{q}_{s,i}, \hat{m}_{s,i}, \hat{r}_w, \hat{q}_w) \\ \hat{q}_w = -2\alpha \partial_{q_w} \Psi_y(R, Q, M) & q_w = -2\partial_{\hat{q}_w} \Psi_w(\hat{r}_{s,i}, \hat{q}_{s,i}, \hat{m}_{s,i}, \hat{r}_w, \hat{q}_w) \end{cases}$$

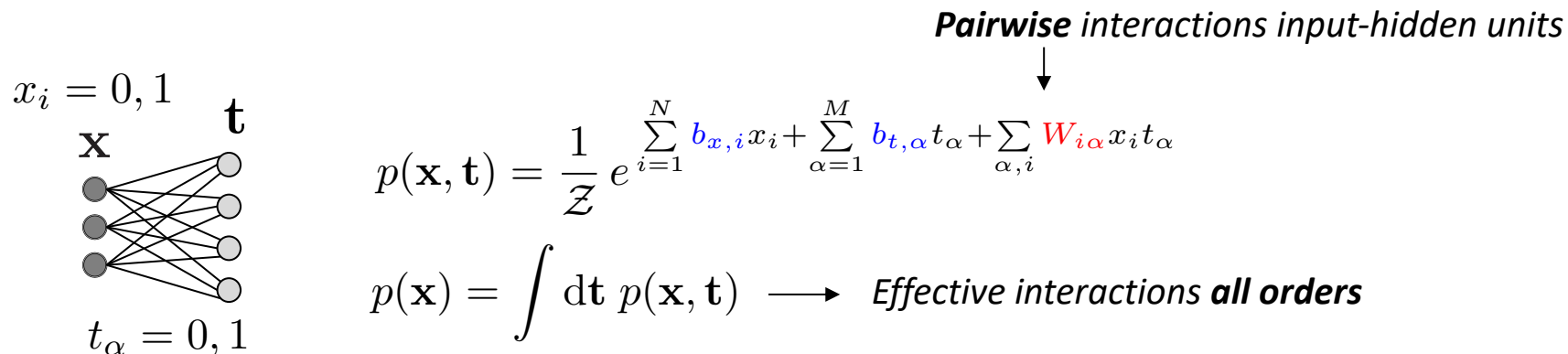
▷ Order parameters are precisely what we need

$$m_{s,i} = \lim_{D \rightarrow \infty} \frac{1}{D} \mathbb{E}_{\mu_\beta} [(\mathbf{w}_2 W_1)_i^\top \beta_i] \quad q_{s,i} = \lim_{D \rightarrow \infty} \frac{1}{D} \mathbb{E}_{\mu_\beta} [(\mathbf{w}_2 W_1)_i^\top (\mathbf{w}_2 W_1)_i], \quad q_w = \lim_{P \rightarrow \infty} \frac{1}{P} \mathbb{E}_{\mu_\beta} [\mathbf{w}^\top \mathbf{w}]$$

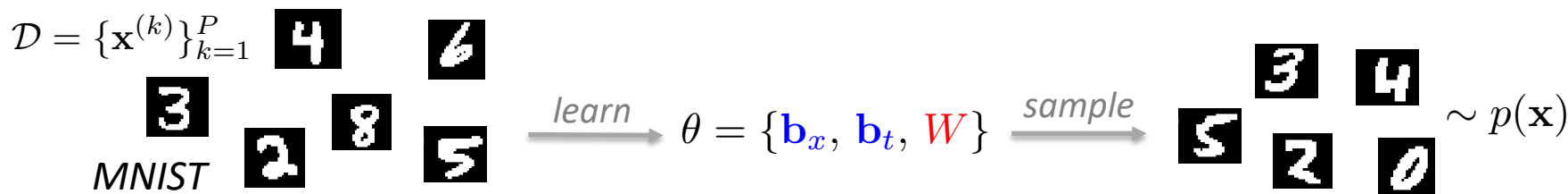
On a different note:

A statistical mechanics inspired learning algorithm

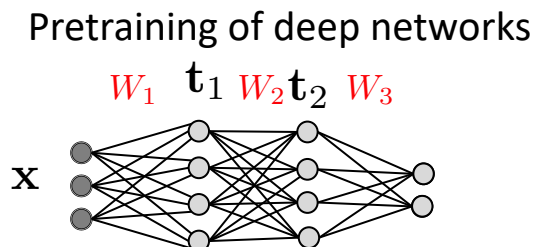
▷ Definition Restricted Boltzmann Machine (RBM):



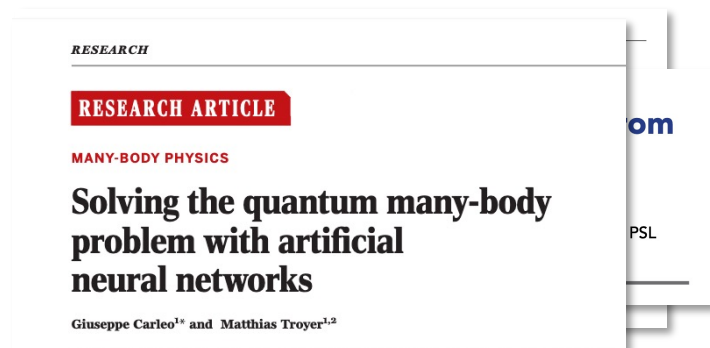
▷ Unsupervised learning:



▷ Applications:



- Biophysics models
- Quantum physics



Learning in Restricted Boltzmann Machines (RBMs)

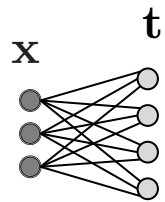
▷ Maximum likelihood learning

$$\ell(W, \mathbf{b}_x, \mathbf{b}_t) = \prod_{k=1}^P p(\mathbf{x}^{(k)})$$

Probability of training data according to RBM

+ gradient ascent !

▷ But intractable exact inference



$$\mathcal{Z} = \sum_{\mathbf{x}, \mathbf{t}} e^{\sum_{i=1}^N b_{x,i} x_i + \sum_{\alpha=1}^M b_{t,\alpha} t_\alpha + \sum_{\alpha,i} W_{i\alpha} x_i t_\alpha}$$

→ 2^{M+N} terms !

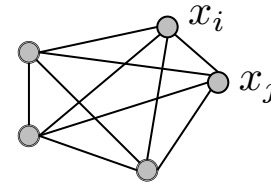
▷ Approximations

- ▷ Approximate Monte Carlo: Contrastive Divergence (Hinton 2002) *state-of-the-art*
- ▷ Mean-field strategies ?
 - Without hidden units: (Kappen, 1999; Cocco & Monasson, 2012; Ricci-Tersenghi 2012 etc ..)
 - With hidden units: (Tieleman et al. 2009; Salakhutdinov & Hinton, 2009) *only naive mean-field*

Can we use other mean-field methods to ease inference and learning ?

From Naive mean-field (NMF) to Thouless-Anderson-Palmer (TAP) approximation

$$E(\mathbf{x}) = - \sum_{i=1}^N b_i x_i - \sum_{i,j} W_{ij} x_i x_j \quad x_i = 0, 1$$



- Free energy functional at **fixed magnetizations** ($\langle x_i \rangle_{\mathbf{m}} = m_i$) + **without interactions**

≈ "β → 0" = easy to compute

$$-\beta G(\mathbf{m}) = \underbrace{H_{\text{NMF}}(\mathbf{m})}_{\text{entropy}} + \beta \underbrace{\langle \sum_{i=1}^N b_i x_i \rangle_{\mathbf{m}}}_{\text{NMF}} + \beta \sum_{i,j} \underbrace{W_{ij} m_i m_j}_{\text{NMF}} + \underbrace{\frac{\beta^2}{2} \sum_{i,j} W_{ij}^2 (m_i - m_i^2)(m_j - m_j^2)}_{\text{TAP}}$$

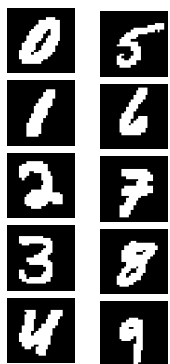
- Select magnetizations of minimum free energy

$$\frac{\partial G}{\partial m_i} = 0 \Rightarrow m_i = \sigma \left[\beta b_i + \sum_j \beta W_{ij} m_j \right] - \beta^2 W_{ij}^2 \left(m_i - \frac{1}{2} \right) (m_j - m_j^2)$$

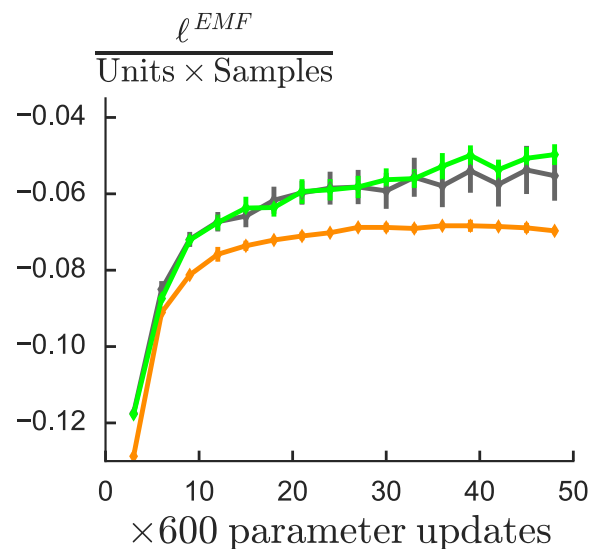
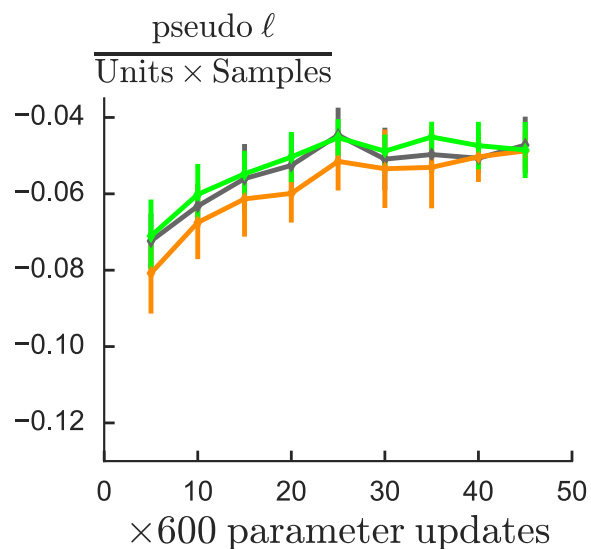
$$\sigma(x) = \frac{1}{1 + e^{-x}}$$

- ▷ High-temperature = small weights
- ▷ $\ln \mathcal{Z} \approx -G(\mathbf{m})$
- ▷ Compare TAP/MCMC (CD) approx of likelihood and gradients

Binarized
MNIST



pseudo log-likelihood log-likelihood TAP log-likelihood



+ TAP
+ NMF
+ CD

*comparable
running times*

Further checks:

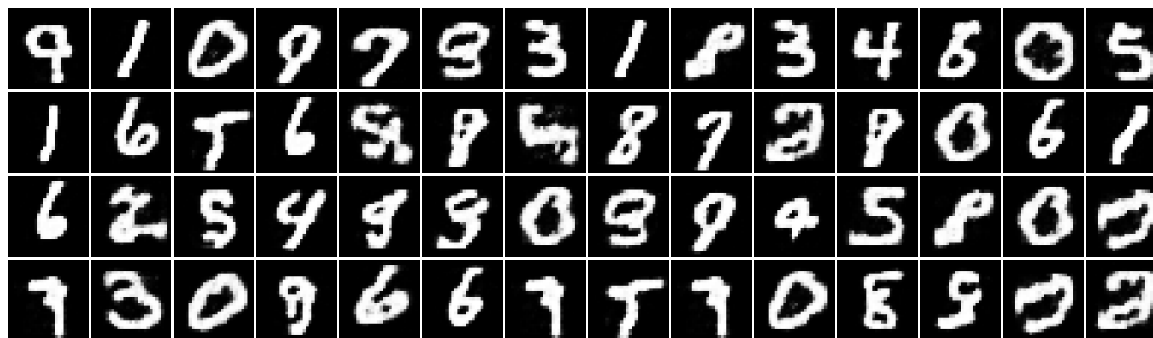
- ▷ No significant improvements with more orders / AdaTAP
- ▷ Extends to real-valued variables

$$m_i^x \leftarrow \sigma \left[\underbrace{\beta b_i + \sum_{\alpha} \beta W_{i\alpha} m_{\alpha}^t}_{\text{NMF}} - \underbrace{\beta^2 W_{ij}^2 \left(m_i^x - \frac{1}{2} \right) (m_{\alpha}^t - m_{\alpha}^{t^2})}_{\text{TAP}} \right]$$

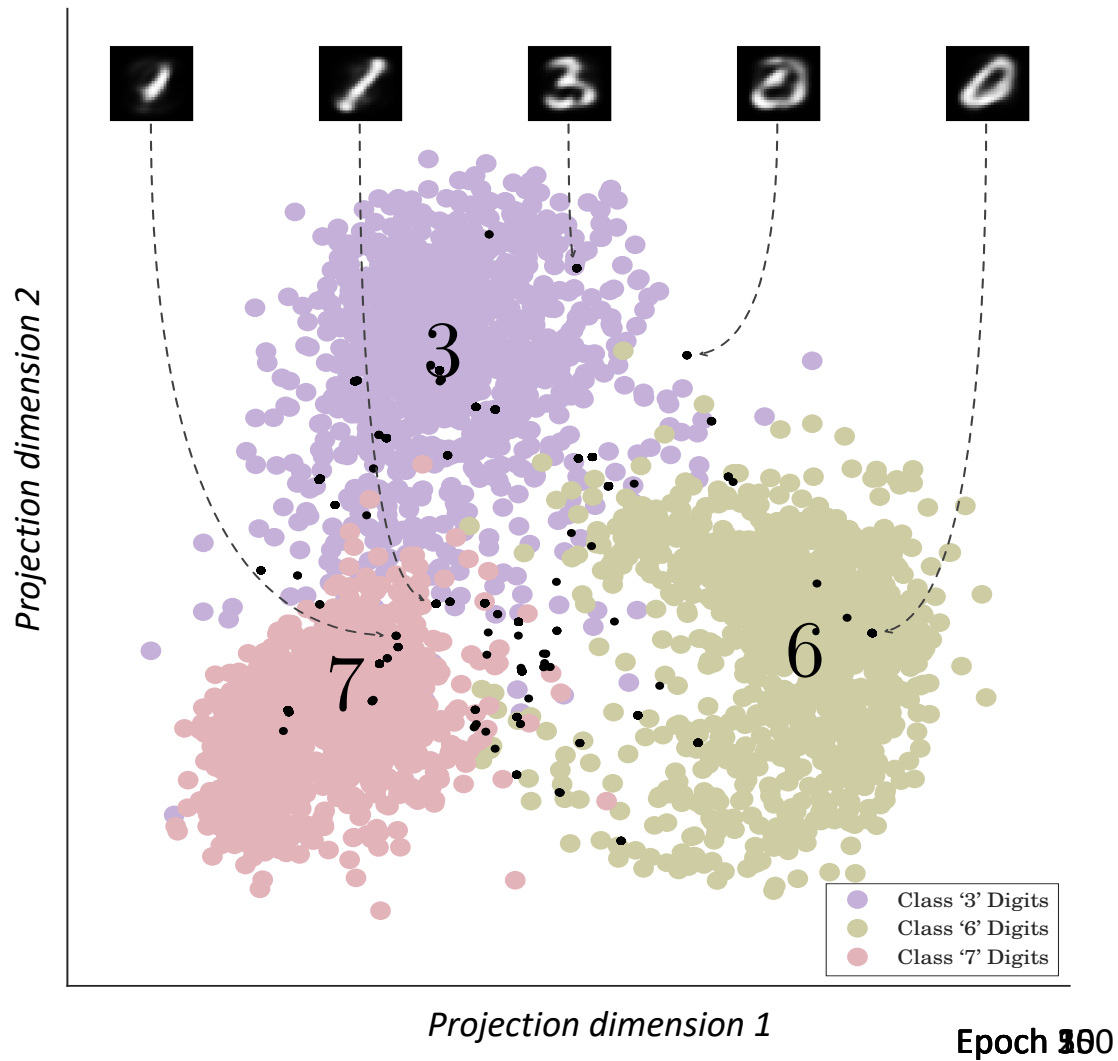
TAP fixed points with trained weights and biases for different initializations



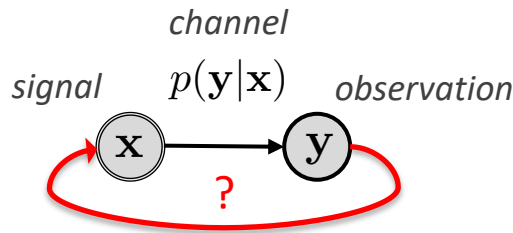
Now for a naive-MF machine



Evolution of magnetizations in configuration space (2d projection)



▷ Reconstruction problems



▷ E.g. compressed sensing (CS)

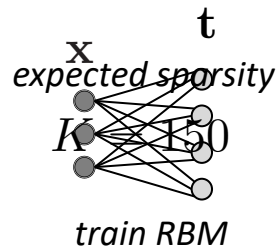
$$\mathbf{y} = F \mathbf{x} + \xi$$

\mathbb{R}^N sparse signal
 $K \ll N$

few observations
 $M \ll N$

▷ Exploit prior information

typical signals



AMP algorithm
 \approx TAP for CS
+
RBM TAP equations

**RBM learned representations
improve drastically reconstruction**

