

The replica method and its applications in biomedical modelling and data analysis

Statistical Physics Approaches to Systems Biology, Havana, Feb 2019

ACC Coolen

King's College London and Saddle Point Science



replica method

A clever trick that enables the analytical calculation of averages that are normally impossible to do, except numerically.

is particularly useful for

Complex heterogeneous systems composed of *many* interacting variables, and with *many* parameters on which we have only statistical information.
(too large for numerical averages to be computationally feasible)

gives us

Analytical predictions for the behaviour of *macroscopic* quantities in *typical* realisations of the systems under study.

note on biomedical applications

The 'large systems' could describe actual *biochemical processes* (folding proteins, proteome, transcriptome, immune or neural networks, etc), or *analysis algorithms* running on large biomedical data sets

1

The replica method

- Exponential families and generating functions
- The replica trick
- The replica trick and algorithms
- Alternative forms of the replica identity

2

Application: information storage in neural networks

- Attractor neural networks
- The replica calculation
- Replica symmetry
- Replica symmetric solution

3

Application: overfitting transition in linear separators

- Linear separability of data – version space
- The replica calculation
- Gardner's replica symmetric theory

1

The replica method

- Exponential families and generating functions
- The replica trick
- The replica trick and algorithms
- Alternative forms of the replica identity

2

Application: information storage in neural networks

- Attractor neural networks
- The replica calculation
- Replica symmetry
- Replica symmetric solution

3

Application: overfitting transition in linear separators

- Linear separability of data – version space
- The replica calculation
- Gardner's replica symmetric theory

Exponential distributions

Often we study stochastic processes for $\mathbf{x} \in X \subseteq \mathbb{R}^N$,
that evolve to a stationary state, with prob distribution $p(\mathbf{x})$
many are of the following form:

- stationary state is *minimally informative*,
subject to a number of constraints

$$\sum_{\mathbf{x} \in X} p(\mathbf{x}) \omega_1(\mathbf{x}) = \Omega_1 \quad \dots \dots \quad \sum_{\mathbf{x} \in X} p(\mathbf{x}) \omega_L(\mathbf{x}) = \Omega_L$$

This is enough to calculate $p(\mathbf{x})$:

- information content of \mathbf{x} : Shannon entropy
hence

$$\text{maximize } S = - \sum_{\mathbf{x} \in X} p(\mathbf{x}) \log p(\mathbf{x})$$

$$\text{subject to : } \begin{cases} p(\mathbf{x}) \geq 0 \quad \forall \mathbf{x}, \quad \sum_{\mathbf{x} \in X} p(\mathbf{x}) = 1 \\ \sum_{\mathbf{x} \in X} p(\mathbf{x}) \omega_\ell(\mathbf{x}) = \Omega_\ell \quad \text{for all } \ell = 1 \dots L \end{cases}$$

- solution using Lagrange's method:

$$\frac{\partial}{\partial p(\mathbf{x})} \left\{ \lambda_0 \sum_{\mathbf{x}' \in X} p(\mathbf{x}') + \sum_{\ell=1}^L \lambda_\ell \sum_{\mathbf{x}' \in X} p(\mathbf{x}') \omega_\ell(\mathbf{x}') - \sum_{\mathbf{x}' \in X} p(\mathbf{x}') \log p(\mathbf{x}') \right\} = 0$$

$$\lambda_0 + \sum_{\ell=1}^L \lambda_\ell \omega_\ell(\mathbf{x}) - 1 - \log p(\mathbf{x}) = 0 \quad \Rightarrow \quad p(\mathbf{x}) = e^{\lambda_0 - 1 + \sum_{\ell=1}^L \lambda_\ell \omega_\ell(\mathbf{x})}$$

$(p(\mathbf{x}) \geq 0 \text{ automatically satisfied})$

- ‘exponential distribution’:

$$p(\mathbf{x}) = \frac{e^{\sum_{\ell=1}^L \lambda_\ell \omega_\ell(\mathbf{x})}}{Z(\boldsymbol{\lambda})}, \quad Z(\boldsymbol{\lambda}) = \sum_{\mathbf{x} \in X} e^{\sum_{\ell=1}^L \lambda_\ell \omega_\ell(\mathbf{x})}$$

$$\boldsymbol{\lambda} = (\lambda_1, \dots, \lambda_L) : \quad \text{solved from} \quad \sum_{\mathbf{x} \in X} p(\mathbf{x}) \omega_\ell(\mathbf{x}) = \Omega_\ell \quad (\ell = 1 \dots L)$$

example:

physical systems in thermal equilibrium

$L = 1$, $\omega(\mathbf{x}) = E(\mathbf{x})$ (energy), $\lambda = -1/k_B T$

$$p(\mathbf{x}) = \frac{e^{-E(\mathbf{x})/k_B T}}{Z(T)}, \quad Z(T) = \sum_{\mathbf{x} \in X} e^{-E(\mathbf{x})/k_B T}$$

Generating functions

$$p(\mathbf{x}) = \frac{e^{\sum_{\ell=1}^L \lambda_\ell \omega_\ell(\mathbf{x})}}{Z(\boldsymbol{\lambda})}, \quad Z(\boldsymbol{\lambda}) = \sum_{\mathbf{x} \in X} e^{\sum_{\ell=1}^L \lambda_\ell \omega_\ell(\mathbf{x})}, \quad \langle f \rangle = \sum_{\mathbf{x} \in X} p(\mathbf{x}) f(\mathbf{x})$$

Idea behind generating functions:
reduce nr of state averages to be calculated ...

- define

$$F(\boldsymbol{\lambda}) = \log Z(\boldsymbol{\lambda}) \quad \frac{\partial F(\boldsymbol{\lambda})}{\partial \lambda_k} = \frac{\sum_{\mathbf{x} \in X} \omega_k(\mathbf{x}) e^{\sum_{\ell=1}^L \lambda_\ell \omega_\ell(\mathbf{x})}}{\sum_{\mathbf{x} \in X} e^{\sum_{\ell=1}^L \lambda_\ell \omega_\ell(\mathbf{x})}} = \langle \omega_k(\mathbf{x}) \rangle$$

- how to calculate
arbitrary state average $\langle \psi \rangle$?

$$F(\boldsymbol{\lambda}, \mu) = \log \left[\sum_{\mathbf{x} \in X} e^{\mu \psi(\mathbf{x}) + \sum_{\ell=1}^L \lambda_\ell \omega_\ell(\mathbf{x})} \right]$$

$$\langle \psi \rangle = \lim_{\mu \rightarrow 0} \frac{\partial F(\boldsymbol{\lambda}, \mu)}{\partial \mu}, \quad \langle \omega_\ell \rangle = \lim_{\mu \rightarrow 0} \frac{\partial F(\boldsymbol{\lambda}, \mu)}{\partial \lambda_\ell}$$

1

The replica method

- Exponential families and generating functions
- **The replica trick**
- The replica trick and algorithms
- Alternative forms of the replica identity

2

Application: information storage in neural networks

- Attractor neural networks
- The replica calculation
- Replica symmetry
- Replica symmetric solution

3

Application: overfitting transition in linear separators

- Linear separability of data – version space
- The replica calculation
- Gardner's replica symmetric theory

The replica trick

first appearance: Marc Kac 1968

first application in physics: Sherrington & Kirkpatrick 1975

first application in biology: Amit, Gutfreund & Sompolinsky 1985

- Consider processes with many fixed (pseudo-)random parameters ξ , distributed according to $\mathcal{P}(\xi)$

$$p(\mathbf{x}|\xi) = \frac{e^{\sum_{\ell=1}^L \lambda_{\ell} \omega_{\ell}(\mathbf{x}, \xi)}}{Z(\lambda, \xi)}, \quad Z(\lambda, \xi) = \sum_{\mathbf{x} \in X} e^{\sum_{\ell=1}^L \lambda_{\ell} \omega_{\ell}(\mathbf{x}, \xi)}$$

- calculating state averages $\langle f \rangle_{\xi}$ for each realisation of ξ is usually impossible
- we are mostly interested in *typical* values of state averages
- for $N \rightarrow \infty$ macroscopic averages will not depend on ξ , only on $\mathcal{P}(\xi)$,
'self-averaging': $\lim_{N \rightarrow \infty} \langle f \rangle_{\xi}$ indep of ξ

so focus on

$$\overline{\langle f \rangle_{\xi}} = \sum_{\xi} \mathcal{P}(\xi) \langle f \rangle_{\xi} = \sum_{\xi} \mathcal{P}(\xi) \left\{ \sum_{\mathbf{x} \in X} p(\mathbf{x}|\xi) f(\mathbf{x}, \xi) \right\}$$

- new generating function:

$$\bar{F}(\lambda, \mu) = \sum_{\xi} \mathcal{P}(\xi) \log Z(\lambda, \mu, \xi), \quad Z(\lambda, \mu, \xi) = \sum_{x \in X} e^{\mu \psi(x, \xi) + \sum_{\ell} \lambda_{\ell} \omega_{\ell}(x, \xi)}$$

$$\begin{aligned} \lim_{\mu \rightarrow 0} \frac{\partial}{\partial \mu} \bar{F}(\lambda, \mu) &= \lim_{\mu \rightarrow 0} \sum_{\xi} \mathcal{P}(\xi) \left\{ \frac{\sum_{x \in X} \psi(x, \xi) e^{\mu \psi(x, \xi) + \sum_{\ell} \lambda_{\ell} \omega_{\ell}(x, \xi)}}{\sum_{x \in X} e^{\mu \psi(x, \xi) + \sum_{\ell} \lambda_{\ell} \omega_{\ell}(x, \xi)}} \right\} \\ &= \sum_{\xi} \mathcal{P}(\xi) \left\{ \frac{\sum_{x \in X} \psi(x, \xi) e^{\sum_{\ell} \lambda_{\ell} \omega_{\ell}(x, \xi)}}{\sum_{x \in X} e^{\sum_{\ell} \lambda_{\ell} \omega_{\ell}(x, \xi)}} \right\} = \overline{\langle \psi \rangle}_{\xi} \end{aligned}$$

- main obstacle in calculating \bar{F} :

the logarithm ...

$$\text{replica identity : } \overline{\log Z} = \lim_{n \rightarrow 0} \frac{1}{n} \log \overline{Z^n}$$

proof:

$$\begin{aligned} \lim_{n \rightarrow 0} \frac{1}{n} \log \overline{Z^n} &= \lim_{n \rightarrow 0} \frac{1}{n} \log \overline{[e^{n \log Z}]} = \lim_{n \rightarrow 0} \frac{1}{n} \log \overline{[1 + n \log Z + \mathcal{O}(n^2)]} \\ &= \lim_{n \rightarrow 0} \frac{1}{n} \log [1 + n \overline{\log Z} + \mathcal{O}(n^2)] = \overline{\log Z} \end{aligned}$$

- apply $\overline{\log Z} = \lim_{n \rightarrow 0} \frac{1}{n} \log \overline{Z^n}$
(simplest case $L = 1$)

$$\begin{aligned}
 \overline{F}(\lambda) &= \sum_{\xi} \mathcal{P}(\xi) \log \left[\sum_{x \in X} e^{\lambda \omega(x, \xi)} \right] = \lim_{n \rightarrow 0} \frac{1}{n} \log \sum_{\xi} \mathcal{P}(\xi) \left[\sum_{x \in X} e^{\lambda \omega(x, \xi)} \right]^n \\
 &= \lim_{n \rightarrow 0} \frac{1}{n} \log \sum_{\xi} \mathcal{P}(\xi) \left[\sum_{x^1 \in X} \dots \sum_{x^n \in X} e^{\lambda \sum_{\alpha=1}^n \omega(x^\alpha, \xi)} \right] \\
 &= \lim_{n \rightarrow 0} \frac{1}{n} \log \left[\sum_{x^1 \in X} \dots \sum_{x^n \in X} \sum_{\xi} \mathcal{P}(\xi) e^{\lambda \sum_{\alpha=1}^n \omega(x^\alpha, \xi)} \right]
 \end{aligned}$$

- notes:

- impossible ξ -average converted into simpler one ...
- calculation involves n ‘replicas’ x^α of original system
- but $n \rightarrow 0$ at the end ... ?
- penultimate step true only for *integer* n ,
so limit requires *analytical continuation* ...

since then: alternative (more tedious) routes,
these confirmed correctness of the replica method!



1

The replica method

- Exponential families and generating functions
- The replica trick
- The replica trick and algorithms**
- Alternative forms of the replica identity

2

Application: information storage in neural networks

- Attractor neural networks
- The replica calculation
- Replica symmetry
- Replica symmetric solution

3

Application: overfitting transition in linear separators

- Linear separability of data – version space
- The replica calculation
- Gardner's replica symmetric theory

The replica trick and algorithms

Suppose we have data D , with prob distr $\mathcal{P}(D)$

and an algorithm which minimises an error function $E(D, \theta)$

(maximum likelihood, Cox & Bayesian regression, SVM, perceptron, ...)

- algorithm outcome:

$$\theta^*(D) = \arg \min_{\theta} E(D, \theta), \quad E_{\min}(D) = \min_{\theta} E(D, \theta)$$

typical performance:

$$\theta^* = \sum_D \mathcal{P}(D) \theta^*(D) = \overline{\theta^*(D)} \quad E_{\min} = \sum_D \mathcal{P}(D) E_{\min}(D) = \overline{E_{\min}(D)}$$

- steepest descent identity & replica trick:

$$E_{\min}(D) = \min_{\theta} E(D, \theta) = - \lim_{\beta \rightarrow \infty} \frac{1}{\beta} \log \int d\theta e^{-\beta E(D, \theta)}$$

$$E_{\min} = \overline{E_{\min}(D)} = - \lim_{\beta \rightarrow \infty} \frac{1}{\beta} \overline{\log \int d\theta e^{-\beta E(D, \theta)}}$$

$$= - \lim_{\beta \rightarrow \infty} \lim_{n \rightarrow 0} \frac{1}{\beta n} \log \left[\overline{\int d\theta e^{-\beta E(D, \theta)}} \right]^n$$

$$= - \lim_{\beta \rightarrow \infty} \lim_{n \rightarrow 0} \frac{1}{\beta n} \log \int d\theta^1 \dots \theta^n \overline{e^{-\beta \sum_{\alpha=1}^n E(D, \theta^\alpha)}}$$

1

The replica method

- Exponential families and generating functions
- The replica trick
- The replica trick and algorithms
- Alternative forms of the replica identity

2

Application: information storage in neural networks

- Attractor neural networks
- The replica calculation
- Replica symmetry
- Replica symmetric solution

3

Application: overfitting transition in linear separators

- Linear separability of data – version space
- The replica calculation
- Gardner's replica symmetric theory

Alternative forms of the replica identity

suppose we need averages, but for
a $p(\mathbf{x}|\xi)$ that is not of an exponential form?

or we need to average quantities that we
don't want in the exponent of $Z(\lambda\xi)$?

$$p(\mathbf{x}|\xi) = \frac{W(\mathbf{x}, \xi)}{\sum_{\mathbf{x}' \in X} W(\mathbf{x}', \xi)}, \quad \overline{\langle f \rangle}_{\xi} = \overline{\sum_{\mathbf{x} \in X} p(\mathbf{x}|\xi) f(\mathbf{x}, \xi)}$$

- main obstacle here:
the fraction ...

$$\begin{aligned}\overline{\langle f \rangle}_{\xi} &= \overline{\left[\frac{\sum_{\mathbf{x} \in X} W(\mathbf{x}, \xi) f(\mathbf{x}, \xi)}{\sum_{\mathbf{x} \in X} W(\mathbf{x}, \xi)} \right]} = \overline{\left[\sum_{\mathbf{x} \in X} W(\mathbf{x}, \xi) f(\mathbf{x}, \xi) \right]} \overline{\left[\sum_{\mathbf{x} \in X} W(\mathbf{x}, \xi) \right]^{-1}} \\ &= \lim_{n \rightarrow 0} \overline{\left[\sum_{\mathbf{x} \in X} W(\mathbf{x}, \xi) f(\mathbf{x}, \xi) \right]} \overline{\left[\sum_{\mathbf{x} \in X} W(\mathbf{x}, \xi) \right]^{n-1}} \\ &= \lim_{n \rightarrow 0} \sum_{\mathbf{x}^1 \in X} \dots \sum_{\mathbf{x}^n \in X} \overline{f(\mathbf{x}^1, \xi) W(\mathbf{x}^1, \xi) \dots W(\mathbf{x}^n, \xi)}\end{aligned}$$

(again: used integer n , but $n \rightarrow 0 \dots$)

- equivalence between two forms
of replica identity, if

$$W(\mathbf{x}, \xi) = e^{\sum_{\ell} \lambda_{\ell} \phi_{\ell}(\mathbf{x}, \xi)}$$

proof:

$$\begin{aligned}\overline{\langle f \rangle_{\xi}} &= \lim_{n \rightarrow 0} \overline{\sum_{\mathbf{x}^1 \in X} \dots \sum_{\mathbf{x}^n \in X} f(\mathbf{x}^1, \xi) W(\mathbf{x}^1, \xi) \dots W(\mathbf{x}^n, \xi)} \\ &= \lim_{n \rightarrow 0} \overline{\sum_{\mathbf{x}^1 \in X} \dots \sum_{\mathbf{x}^n \in X} f(\mathbf{x}^1, \xi) e^{\sum_{\alpha=1}^n \sum_{\ell} \lambda_{\ell} \phi_{\ell}(\mathbf{x}^{\alpha}, \xi)}} \\ &= \lim_{n \rightarrow 0} \frac{1}{n} \overline{\sum_{\mathbf{x}^1 \in X} \dots \sum_{\mathbf{x}^n \in X} \left[\sum_{\alpha=1}^n f(\mathbf{x}^{\alpha}, \xi) \right] e^{\sum_{\alpha=1}^n \sum_{\ell} \lambda_{\ell} \phi_{\ell}(\mathbf{x}^{\alpha}, \xi)}} \\ &= \lim_{n \rightarrow 0} \frac{1}{n} \lim_{\mu \rightarrow 0} \frac{\partial}{\partial \mu} \overline{\sum_{\mathbf{x}^1 \in X} \dots \sum_{\mathbf{x}^n \in X} e^{\sum_{\alpha=1}^n \sum_{\ell} \lambda_{\ell} \phi_{\ell}(\mathbf{x}^{\alpha}, \xi) + \mu \sum_{\alpha=1}^n f(\mathbf{x}^{\alpha}, \xi)}} \\ &= \lim_{\mu \rightarrow 0} \frac{\partial}{\partial \mu} \lim_{n \rightarrow 0} \frac{1}{n} \overline{\sum_{\mathbf{x}^1 \in X} \dots \sum_{\mathbf{x}^n \in X} e^{\sum_{\alpha=1}^n [\sum_{\ell} \lambda_{\ell} \phi_{\ell}(\mathbf{x}^{\alpha}, \xi) + \mu f(\mathbf{x}^{\alpha}, \xi)]}} \\ &= \lim_{\mu \rightarrow 0} \frac{\partial}{\partial \mu} \lim_{n \rightarrow 0} \frac{1}{n} \overline{Z^n(\lambda, \mu, \xi)}, \quad Z(\lambda, \mu, \xi) = \sum_{\mathbf{x} \in X} e^{\sum_{\ell} \lambda_{\ell} \phi_{\ell}(\mathbf{x}, \xi) + \mu f(\mathbf{x}, \xi)}\end{aligned}$$

1

The replica method

- Exponential families and generating functions
- The replica trick
- The replica trick and algorithms
- Alternative forms of the replica identity

2

Application: information storage in neural networks

• Attractor neural networks

- The replica calculation
- Replica symmetry
- Replica symmetric solution

3

Application: overfitting transition in linear separators

- Linear separability of data – version space
- The replica calculation
- Gardner's replica symmetric theory

Attractor neural networks

$N \sim 10^{12-14}$ brain cells (neurons),
each connected with $\sim 10^{3-5}$ others

- **neurons**

two states:

$$\begin{aligned}\sigma_i = 1 & \quad (i \text{ fires electric pulses}) \\ \sigma_i = -1 & \quad (i \text{ is at rest})\end{aligned}$$

- **dynamics of firing states**

$$\sigma_i(t+1) = \operatorname{sgn} \left[\underbrace{\sum_{j=1}^N J_{ij} \sigma_j(t)}_{\text{activation signal}} + \underbrace{\theta_i + z_i(t)}_{\text{threshold, noise}} \right]$$

$\theta_i \in \mathbb{R}$: *firing threshold of neuron i*
 $J_{ij} \in \mathbb{R}$: *synaptic connection $j \rightarrow i$*

learning = adaptation of $\{J_{ij}, \theta_i\}$



*non-local 'distributed' storage of
'program' and 'data'*

attractor neural networks

models for associative memory in the brain

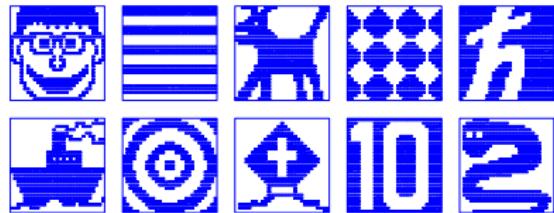
- the neural code

represent ‘patterns’ as

$$\text{micro-states } \xi = (\xi_1, \dots, \xi_N)$$

○: $\sigma_i = -1$, ●: $\sigma_i = 1$

e.g. $N=400$,
10 patterns:



- information storage

modify synapses $\{J_{ij}\}$ such that ξ is
stable state (attractor) of the neuronal dynamics

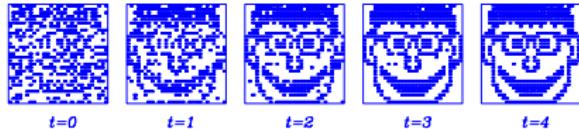
- information recall

initial state $\sigma(t=0)$:

evolution to nearest attractor

if $\sigma(0)$ close (i.e. similar) to ξ :

$$\sigma(t=\infty) = \xi$$



- learning rule: recipe for storing patterns via modification of $\{J_{ij}\}$

Hebb (1949): $\Delta J_{ij} \propto \xi_i \xi_j$

choose $J_{ij} = J_0 \xi_i \xi_j$, $\theta_i = 0$,

update randomly drawn i at each step:

$$\begin{aligned}\sigma_i(t+1) &= sgn \left[\sum_{j=1}^N J_{ij} \sigma_j(t) + z_i(t) \right] = sgn \left[J_0 \xi_i \left(\overbrace{\sum_{j=1}^N \xi_j \sigma_j(t)}^{pattern\ overlap} \right) + z_i(t) \right] \\ &= \xi_i sgn \left[J_0 \sum_{j=1}^N \xi_j \sigma_j(t) + \xi_i z_i(t) \right]\end{aligned}$$

$M(t) = \sum_{j=1}^N \xi_j \sigma_j(t)$ sufficiently large: $\sigma_i(t+1) = \xi_i$

now $M(t+1) \geq M(t) \dots$

will continue until $\sigma = \xi$

- proper analysis:

noise: $P(z) = \frac{\beta}{2} [1 - \tanh^2(\beta z)]$,

symmetric synapses: $J_{ij} = J_{ji}$, $J_{ii} = 0$

sequential updates of σ_i

$$p(\sigma) = \frac{e^{-\beta H(\sigma)}}{Z(\beta)}, \quad H(\sigma) = -\frac{1}{2} \sum_{i \neq j} \sigma_i J_{ij} \sigma_j - \sum_i \theta_i \sigma_i$$

a more realistic model,
solvable via the replica method

- **storage of a pattern** $\xi = (\xi_1, \dots, \xi_N) \in \{-1, 1\}^N$
on background of zero-average Gaussian synapses

$$J_{ij} = \frac{J_0}{N} \xi_i \xi_j + \frac{J}{\sqrt{N}} \textcolor{blue}{z}_{ij}, \quad \bar{z}_{ij} = 0, \quad \bar{z}_{ij}^2 = 1, \quad J, J_0 \geq 0, \quad \theta_i = 0$$

to be averaged over: background synapses $\{z_{ij}\}$

pattern overlap: $m(\sigma) = \frac{1}{N} \sum_k \sigma_k \xi_k$

$$\begin{aligned} H(\sigma) &= -\frac{1}{2} \sum_{i \neq j} \sigma_i \sigma_j \left\{ \frac{J_0}{N} \xi_i \xi_j + \frac{J}{\sqrt{N}} \textcolor{blue}{z}_{ij} \right\} \\ &= -\frac{J_0}{2N} \sum_{ij} \sigma_i \sigma_j \xi_i \xi_j + \frac{J_0}{2N} \sum_i 1 - \frac{J}{2\sqrt{N}} \sum_{i \neq j} \sigma_i \sigma_j \textcolor{blue}{z}_{ij} \\ &= -\frac{1}{2} N J_0 m^2(\sigma) + \frac{1}{2} J_0 - \frac{J}{\sqrt{N}} \sum_{i < j} \sigma_i \sigma_j \textcolor{blue}{z}_{ij} \end{aligned}$$

- **generating function**

$$\overline{F} = \overline{\log Z(\beta)} = \lim_{n \rightarrow 0} \frac{1}{n} \log \overline{Z^n(\beta)} = \lim_{n \rightarrow 0} \frac{1}{n} \log \left[\sum_{\sigma^1 \dots \sigma^n} \overline{e^{-\beta \sum_{\alpha=1}^n H(\sigma^\alpha)}} \right]$$

1

The replica method

- Exponential families and generating functions
- The replica trick
- The replica trick and algorithms
- Alternative forms of the replica identity

2

Application: information storage in neural networks

- Attractor neural networks
- **The replica calculation**
- Replica symmetry
- Replica symmetric solution

3

Application: overfitting transition in linear separators

- Linear separability of data – version space
- The replica calculation
- Gardner's replica symmetric theory

The replica calculation

short-hands: $m(\sigma) = \frac{1}{N} \sum_i \xi_i \sigma_i$, $DZ = (2\pi)^{-1/2} e^{-z^2/2} dz$

Gaussian integral: $\int DZ e^{xz} = e^{\frac{1}{2}x^2}$

• average over random synapses

$$\begin{aligned}\overline{Z^n(\beta)} &= \sum_{\sigma^1 \dots \sigma^n} \overline{e^{-\beta \sum_{\alpha=1}^n H(\sigma^\alpha)}} \\ &= \sum_{\sigma^1 \dots \sigma^n} \overline{e^{-\beta \sum_{\alpha=1}^n \left[\frac{1}{2} J_0 - \frac{1}{2} N J_0 m^2(\sigma^\alpha) - \frac{J}{\sqrt{N}} \sum_{i < j} \sigma_i^\alpha \sigma_j^\alpha z_{ij} \right]}} \\ &= e^{-\frac{1}{2} n \beta J_0} \sum_{\sigma^1 \dots \sigma^n} e^{\frac{1}{2} N \beta J_0 \sum_{\alpha=1}^n m^2(\sigma^\alpha)} \overline{e^{\frac{\beta J}{\sqrt{N}} \sum_{\alpha=1}^n \sum_{i < j} \sigma_i^\alpha \sigma_j^\alpha z_{ij}}} \\ &= e^{-\frac{1}{2} n \beta J_0} \sum_{\sigma^1 \dots \sigma^n} e^{\frac{1}{2} N \beta J_0 \sum_{\alpha=1}^n m^2(\sigma^\alpha)} \prod_{i < j} \int Dz e^{\frac{\beta J}{\sqrt{N}} \sum_{\alpha=1}^n \sigma_i^\alpha \sigma_j^\alpha z} \\ &= e^{-\frac{1}{2} n \beta J_0} \sum_{\sigma^1 \dots \sigma^n} e^{\frac{1}{2} N \beta J_0 \sum_{\alpha=1}^n m^2(\sigma^\alpha)} \prod_{i < j} e^{\frac{\beta^2 J^2}{2N} \left[\sum_{\alpha=1}^n \sigma_i^\alpha \sigma_j^\alpha \right]^2} \\ &= e^{-\frac{1}{2} n \beta J_0} \sum_{\sigma^1 \dots \sigma^n} e^{N \left[\frac{1}{2} \beta J_0 \sum_{\alpha=1}^n m^2(\sigma^\alpha) + \frac{1}{2} (\beta J)^2 \sum_{\alpha, \gamma=1}^n \left(N^{-2} \sum_{i < j} \sigma_i^\alpha \sigma_j^\alpha \sigma_i^\gamma \sigma_j^\gamma \right) \right]}\end{aligned}$$

- **complete square** in sums over neurons

$$\begin{aligned} \sum_{i < j} \sigma_i^\alpha \sigma_j^\alpha \sigma_i^\gamma \sigma_j^\gamma &= \frac{1}{2} \sum_{i \neq j} \sigma_i^\alpha \sigma_j^\alpha \sigma_i^\gamma \sigma_j^\gamma = \frac{1}{2} \sum_{ij} \sigma_i^\alpha \sigma_j^\alpha \sigma_i^\gamma \sigma_j^\gamma - \frac{1}{2} \sum_i 1 \\ &= \frac{1}{2} \left(\sum_i \sigma_i^\alpha \sigma_i^\gamma \right)^2 - \frac{1}{2} N \end{aligned}$$

hence

$$\begin{aligned} \overline{Z^n(\beta)} &= e^{-\frac{1}{2}n\beta J_0} \sum_{\sigma^1 \dots \sigma^n} e^{N \left[\frac{1}{2}\beta J_0 \sum_{\alpha=1}^n m^2(\sigma^\alpha) + \frac{1}{4}(\beta J)^2 \sum_{\alpha, \gamma=1}^n \left(\left(\frac{1}{N} \sum_i \sigma_i^\alpha \sigma_i^\gamma \right)^2 - \frac{1}{N} \right) \right]} \\ &= e^{-\frac{1}{2}n\beta J_0 - \frac{1}{4}n(\beta J)^2} \sum_{\sigma^1 \dots \sigma^n} e^{N \left[\frac{1}{2}\beta J_0 \sum_{\alpha=1}^n m^2(\sigma^\alpha) + \frac{1}{4}(\beta J)^2 \sum_{\alpha, \gamma=1}^n \left(\frac{1}{N} \sum_i \sigma_i^\alpha \sigma_i^\gamma \right)^2 \right]} \end{aligned}$$

- insert:

$$1 = \prod_{\alpha=1}^n \int dm_\alpha \delta \left(m_\alpha - \frac{1}{N} \sum_i \xi_i \sigma_i^\alpha \right), \quad 1 = \prod_{\alpha, \gamma=1}^n \int dq_{\alpha\gamma} \delta \left(q_{\alpha\gamma} - \frac{1}{N} \sum_i \sigma_i^\alpha \sigma_i^\gamma \right)$$

$\mathbf{m} \in \mathbb{R}^n$, $\mathbf{q} \in \mathbb{R}^{n^2}$:

$$\begin{aligned} \overline{Z^n(\beta)} &= e^{-\frac{1}{2}n\beta J_0 - \frac{1}{4}n(\beta J)^2} \int d\mathbf{m} d\mathbf{q} e^{N \left[\frac{1}{2}\beta J_0 \sum_{\alpha=1}^n m_\alpha^2 + \frac{1}{4}(\beta J)^2 \sum_{\alpha, \gamma=1}^n q_{\alpha\gamma}^2 \right]} \\ &\times \sum_{\sigma^1 \dots \sigma^n} \left[\prod_{\alpha=1}^n \delta \left(m_\alpha - \frac{1}{N} \sum_i \xi_i \sigma_i^\alpha \right) \right] \left[\prod_{\alpha, \gamma=1}^n \delta \left(q_{\alpha\gamma} - \frac{1}{N} \sum_i \sigma_i^\alpha \sigma_i^\gamma \right) \right] \end{aligned}$$

remember: $\delta(x) = (2\pi)^{-1} \int d\hat{x} e^{ix\hat{x}}$

- the sum over neuron state variables

$$\begin{aligned}
 & \sum_{\sigma^1 \dots \sigma^n} \left[\prod_{\alpha=1}^n \delta\left(m_\alpha - \frac{1}{N} \sum_i \xi_i \sigma_i^\alpha\right) \right] \left[\prod_{\alpha, \gamma=1}^n \delta\left(q_{\alpha\gamma} - \frac{1}{N} \sum_i \sigma_i^\alpha \sigma_i^\gamma\right) \right] \\
 &= \sum_{\sigma^1 \dots \sigma^n} \int \frac{d\hat{\mathbf{m}} d\hat{\mathbf{q}}}{(2\pi)^{n^2+n}} e^{i \sum_{\alpha=1}^n \hat{m}_\alpha \left[m_\alpha - \frac{1}{N} \sum_i \xi_i \sigma_i^\alpha\right] + i \sum_{\alpha, \gamma=1}^n \hat{q}_{\alpha\gamma} \left[q_{\alpha\gamma} - \frac{1}{N} \sum_i \sigma_i^\alpha \sigma_i^\gamma\right]} \\
 &= \int \frac{d\hat{\mathbf{m}} d\hat{\mathbf{q}}}{(2\pi)^{n(n+1)}} e^{i \sum_\alpha \hat{m}_\alpha m_\alpha + i \sum_{\alpha\gamma} \hat{q}_{\alpha\gamma} q_{\alpha\gamma}} \sum_{\sigma^1 \dots \sigma^n} \prod_{i=1}^N e^{-\frac{i}{N} \left[\sum_\alpha \hat{m}_\alpha \xi_i \sigma_i^\alpha + \sum_{\alpha\gamma} \hat{q}_{\alpha\gamma} \sigma_i^\alpha \sigma_i^\gamma\right]} \\
 &= \int \frac{d\hat{\mathbf{m}} d\hat{\mathbf{q}}}{(2\pi)^{n(n+1)}} e^{i \sum_\alpha \hat{m}_\alpha m_\alpha + i \sum_{\alpha\gamma} \hat{q}_{\alpha\gamma} q_{\alpha\gamma}} \prod_i \sum_{\sigma_1 \dots \sigma_n} e^{-\frac{i}{N} \left[\sum_\alpha \hat{m}_\alpha \xi_i \sigma_\alpha + \sum_{\alpha\gamma} \hat{q}_{\alpha\gamma} \sigma_\alpha \sigma_\gamma\right]}
 \end{aligned}$$

transform: $\hat{\mathbf{m}} \rightarrow N\hat{\mathbf{m}}$, $\hat{\mathbf{q}} \rightarrow N\hat{\mathbf{q}}$, $\sigma_\alpha \rightarrow \xi_i \sigma_\alpha$:

$$\begin{aligned}
 \sum_{\sigma^1 \dots \sigma^n} [\dots] [\dots] &= \int \frac{d\hat{\mathbf{m}} d\hat{\mathbf{q}}}{(2\pi/N)^{n(n+1)}} e^{iN[\hat{\mathbf{m}} \cdot \mathbf{m} + \text{Tr}(\hat{\mathbf{q}} \mathbf{q})]} \left[\sum_{\sigma \in \{-1, 1\}^n} e^{-i\hat{\mathbf{m}} \cdot \boldsymbol{\sigma} - i\boldsymbol{\sigma} \cdot \hat{\mathbf{q}} \boldsymbol{\sigma}} \right]^N \\
 &= \int \frac{d\hat{\mathbf{m}} d\hat{\mathbf{q}}}{(2\pi/N)^{n(n+1)}} e^{iN\hat{\mathbf{m}} \cdot \mathbf{m} + iN\text{Tr}(\hat{\mathbf{q}} \mathbf{q}) + N \log \sum_{\boldsymbol{\sigma}} \exp(-i\hat{\mathbf{m}} \cdot \boldsymbol{\sigma} - i\boldsymbol{\sigma} \cdot \hat{\mathbf{q}} \boldsymbol{\sigma})}
 \end{aligned}$$

- combine everything ...

$$\overline{Z^n(\beta)} = e^{-\frac{1}{2}n\beta J_0 - \frac{1}{4}n(\beta J)^2 - n(n+1)\log(2\pi/N)} \int d\mathbf{m}d\mathbf{q}d\hat{\mathbf{m}}d\hat{\mathbf{q}} e^{N\Psi(\mathbf{m}, \mathbf{q}, \hat{\mathbf{m}}, \hat{\mathbf{q}})}$$

$$\Psi(\dots) = \frac{1}{2}\beta J_0 \mathbf{m}^2 + \frac{1}{4}(\beta J)^2 \text{Tr}(\mathbf{q}^2) + i\hat{\mathbf{m}} \cdot \mathbf{m} + i\text{Tr}(\hat{\mathbf{q}}\mathbf{q}) + \log \sum_{\sigma} e^{-i\hat{\mathbf{m}} \cdot \sigma - i\sigma \cdot \hat{\mathbf{q}}\sigma}$$

Hence

$$\begin{aligned}\overline{F} &= \lim_{n \rightarrow 0} \frac{1}{n} \log \overline{Z^n(\beta)} \\ &= -\frac{1}{2}\beta J_0 - \frac{1}{4}(\beta J)^2 - \log\left(\frac{2\pi}{N}\right) + \lim_{n \rightarrow 0} \frac{1}{n} \log \int d\mathbf{m}d\mathbf{q}d\hat{\mathbf{m}}d\hat{\mathbf{q}} e^{N\Psi(\mathbf{m}, \mathbf{q}, \hat{\mathbf{m}}, \hat{\mathbf{q}})}\end{aligned}$$

- Since $\overline{F} = \mathcal{O}(N)$,
large N behaviour follows from

$$\bar{f} = \lim_{N \rightarrow \infty} \overline{F}/N = \lim_{N \rightarrow \infty} \lim_{n \rightarrow 0} \frac{1}{nN} \log \int d\mathbf{m}d\mathbf{q}d\hat{\mathbf{m}}d\hat{\mathbf{q}} e^{N\Psi(\mathbf{m}, \mathbf{q}, \hat{\mathbf{m}}, \hat{\mathbf{q}})}$$

- assume limits commute,
steepest descent integration:

$$\bar{f} = \lim_{n \rightarrow 0} \frac{1}{n} \text{extr}_{\mathbf{m}, \mathbf{q}, \hat{\mathbf{m}}, \hat{\mathbf{q}}} \Psi(\mathbf{m}, \mathbf{q}, \hat{\mathbf{m}}, \hat{\mathbf{q}})$$

$$\begin{aligned}\Psi(\dots) = & \frac{1}{2}\beta J_0 \sum_{\alpha} m_{\alpha}^2 + \frac{1}{4}(\beta J)^2 \sum_{\alpha\gamma} q_{\alpha\gamma}^2 + i \sum_{\alpha} \hat{m}_{\alpha} m_{\alpha} + i \sum_{\alpha\gamma} \hat{q}_{\alpha\gamma} q_{\alpha\gamma} \\ & + \log \sum_{\sigma} e^{-i \sum_{\lambda} \hat{m}_{\lambda} \sigma_{\lambda} - i \sum_{\lambda\zeta} \sigma_{\lambda} \hat{q}_{\lambda\zeta} \sigma_{\zeta}}\end{aligned}$$

● saddle-point eqns

$$\frac{\partial \Psi}{\partial m_{\alpha}} = 0, \quad \frac{\partial \Psi}{\partial q_{\alpha\gamma}} = 0 : \quad \beta J_0 m_{\alpha} + i \hat{m}_{\alpha} = 0, \quad \frac{1}{2}(\beta J)^2 q_{\alpha\gamma} + i \hat{q}_{\alpha\gamma} = 0$$

$$\frac{\partial \Psi}{\partial \hat{m}_{\alpha}} = 0 : \quad i m_{\alpha} - i \frac{\sum_{\sigma} \sigma_{\alpha} e^{-i \sum_{\lambda} \hat{m}_{\lambda} \sigma_{\lambda} - i \sum_{\lambda\zeta} \sigma_{\lambda} \hat{q}_{\lambda\zeta} \sigma_{\zeta}}}{\sum_{\sigma} e^{-i \sum_{\lambda} \hat{m}_{\lambda} \sigma_{\lambda} - i \sum_{\lambda\zeta} \sigma_{\lambda} \hat{q}_{\lambda\zeta} \sigma_{\zeta}}} = 0$$

$$\frac{\partial \Psi}{\partial \hat{q}_{\alpha\gamma}} = 0 : \quad i q_{\alpha\gamma} - i \frac{\sum_{\sigma} \sigma_{\alpha} \sigma_{\gamma} e^{-i \sum_{\lambda} \hat{m}_{\lambda} \sigma_{\lambda} - i \sum_{\lambda\zeta} \sigma_{\lambda} \hat{q}_{\lambda\zeta} \sigma_{\zeta}}}{\sum_{\sigma} e^{-i \sum_{\lambda} \hat{m}_{\lambda} \sigma_{\lambda} - i \sum_{\lambda\zeta} \sigma_{\lambda} \hat{q}_{\lambda\zeta} \sigma_{\zeta}}} = 0$$

● eliminate (\hat{m}, \hat{q})

$$m_{\alpha} = \frac{\sum_{\sigma} \sigma_{\alpha} e^{\beta J_0 \sum_{\lambda} m_{\lambda} \sigma_{\lambda} + \frac{1}{2}(\beta J)^2 \sum_{\lambda \neq \zeta} \sigma_{\lambda} q_{\lambda\zeta} \sigma_{\zeta}}}{\sum_{\sigma} e^{\beta J_0 \sum_{\lambda} m_{\lambda} \sigma_{\lambda} + \frac{1}{2}(\beta J)^2 \sum_{\lambda \neq \zeta} \sigma_{\lambda} q_{\lambda\zeta} \sigma_{\zeta}}}$$

$$q_{\alpha\gamma} = \frac{\sum_{\sigma} \sigma_{\alpha} \sigma_{\gamma} e^{\beta J_0 \sum_{\lambda} m_{\lambda} \sigma_{\lambda} + \frac{1}{2}(\beta J)^2 \sum_{\lambda \neq \zeta} \sigma_{\lambda} q_{\lambda\zeta} \sigma_{\zeta}}}{\sum_{\sigma} e^{\beta J_0 \sum_{\lambda} m_{\lambda} \sigma_{\lambda} + \frac{1}{2}(\beta J)^2 \sum_{\lambda \neq \zeta} \sigma_{\lambda} q_{\lambda\zeta} \sigma_{\zeta}}}$$

trivial soln: $\mathbf{m} = \mathbf{q} = \mathbf{0}$,
any others?

1

The replica method

- Exponential families and generating functions
- The replica trick
- The replica trick and algorithms
- Alternative forms of the replica identity

2

Application: information storage in neural networks

- Attractor neural networks
- The replica calculation
- Replica symmetry**
- Replica symmetric solution

3

Application: overfitting transition in linear separators

- Linear separability of data – version space
- The replica calculation
- Gardner's replica symmetric theory

Replica symmetry

- $\beta = 0$ (infinite noise level):

$$m_\alpha = \frac{\sum_{\sigma} \sigma_\alpha e^0}{\sum_{\sigma} e^0} = 0, \quad q_{\alpha\gamma} = \frac{\sum_{\sigma} \sigma_\alpha \sigma_\gamma e^0}{\sum_{\sigma} e^0} = 0$$

$\mathbf{m} = \mathbf{q} = \mathbf{0}$ if $\beta = 0$

- bifurcations from trivial soln:

$$\begin{aligned} m_\alpha &= \frac{\sum_{\sigma} \sigma_\alpha \left[1 + \beta J_0 \sum_{\lambda} m_{\lambda} \sigma_{\lambda} + \frac{1}{2} (\beta J)^2 \sum_{\lambda \neq \zeta} \sigma_{\lambda} q_{\lambda\zeta} \sigma_{\zeta} \right]}{\sum_{\sigma} \left[1 + \beta J_0 \sum_{\lambda} m_{\lambda} \sigma_{\lambda} + \frac{1}{2} (\beta J)^2 \sum_{\lambda \neq \zeta} \sigma_{\lambda} q_{\lambda\zeta} \sigma_{\zeta} \right]} + \mathcal{O}(\mathbf{m}, \mathbf{q})^2 \\ &= \frac{2^n \beta J_0 m_\alpha}{2^n} + \dots = \beta J_0 m_\alpha + \dots \end{aligned}$$

$\mathbf{m} \neq \mathbf{0}$ if $\beta J_0 > 1$

$$\begin{aligned} q_{\alpha\gamma} &= \frac{\sum_{\sigma} \sigma_\alpha \sigma_\gamma \left[1 + \beta J_0 \sum_{\lambda} m_{\lambda} \sigma_{\lambda} + \frac{1}{2} (\beta J)^2 \sum_{\lambda \neq \zeta} \sigma_{\lambda} q_{\lambda\zeta} \sigma_{\zeta} \right]}{\sum_{\sigma} \left[1 + \beta J_0 \sum_{\lambda} m_{\lambda} \sigma_{\lambda} + \frac{1}{2} (\beta J)^2 \sum_{\lambda \neq \zeta} \sigma_{\lambda} q_{\lambda\zeta} \sigma_{\zeta} \right]} + \mathcal{O}(\mathbf{m}, \mathbf{q})^2 \\ &= \frac{2^n (\beta J)^2 q_{\alpha\gamma} + \dots}{2^n} + \dots = (\beta J)^2 q_{\alpha\gamma} + \dots \end{aligned}$$

$\mathbf{q} \neq \mathbf{0}$ if $\beta J > 1$

how to find form of nontrivial solns $\{m_\alpha, q_{\alpha\gamma}\}$?

need their physical interpretation!

use alternative form(s) of replica identity:

$$\overline{\langle f(\sigma) \rangle} = \lim_{n \rightarrow 0} \frac{1}{n} \sum_{\gamma=1}^n \sum_{\sigma^1} \dots \sum_{\sigma^n} \overline{f(\sigma^\gamma)} e^{-\beta \sum_{\alpha=1}^n H(\sigma^\alpha)}$$

$$\overline{\langle\langle f(\sigma, \sigma') \rangle\rangle} = \lim_{n \rightarrow 0} \frac{1}{n(n-1)} \sum_{\alpha \neq \gamma=1}^n \sum_{\sigma^1} \dots \sum_{\sigma^n} \overline{f(\sigma^\alpha, \sigma^\gamma)} e^{-\beta \sum_{\alpha=1}^n H(\sigma^\alpha)}$$

apply to

$$P(m|\sigma) = \delta \left[m - \frac{1}{N} \sum_{i=1}^N \xi_i \sigma_i \right], \quad P(q|\sigma, \sigma') = \delta \left[q - \frac{1}{N} \sum_{i=1}^N \sigma_i \sigma'_i \right]$$

repeat steps of previous calculation,
gives expressions in terms of
saddle-point soln $\{m_\alpha, q_{\alpha\gamma}\}$:

$$\lim_{N \rightarrow \infty} \overline{\langle P(m|\sigma) \rangle} = \lim_{n \rightarrow 0} \frac{1}{n} \sum_{\alpha=1}^n \delta[m - m_\alpha]$$

$$\lim_{N \rightarrow \infty} \overline{\langle\langle P(q|\sigma, \sigma') \rangle\rangle} = \lim_{n \rightarrow 0} \frac{1}{n(n-1)} \sum_{\alpha \neq \gamma=1}^n \delta[q - q_{\alpha\gamma}]$$

ergodic mean-field systems

fluctuations in quantities like $\frac{1}{N} \sum_{i=1}^N \xi_i \sigma_i$
or $\frac{1}{N} \sum_{i=1}^N \sigma_i \sigma'_i$ scale as $\mathcal{O}(N^{-1/2})$

hence

$$\lim_{N \rightarrow \infty} \langle P(m|\sigma) \rangle = \lim_{N \rightarrow \infty} \left\langle \delta \left[m - \frac{1}{N} \sum_{i=1}^N \xi_i \sigma_i \right] \right\rangle = \delta \left[m - \frac{1}{N} \sum_{i=1}^N \xi_i \langle \sigma_i \rangle \right]$$
$$\lim_{N \rightarrow \infty} \langle\langle P(q|\sigma, \sigma') \rangle\rangle = \lim_{N \rightarrow \infty} \left\langle\left\langle \delta \left[q - \frac{1}{N} \sum_{i=1}^N \sigma_i \sigma'_i \right] \right\rangle\right\rangle = \delta \left[q - \frac{1}{N} \sum_{i=1}^N \langle \sigma_i \rangle^2 \right]$$

hence

$$\forall \alpha : \quad m_\alpha = m = \lim_{N \rightarrow \infty} \frac{1}{N} \sum_{i=1}^N \xi_i \overline{\langle \sigma_i \rangle}$$

$$\forall \alpha \neq \gamma : \quad q_{\alpha\gamma} = q = \lim_{N \rightarrow \infty} \frac{1}{N} \sum_{i=1}^N \overline{\langle \sigma_i \rangle^2}$$

replica-symmetric solution

1

The replica method

- Exponential families and generating functions
- The replica trick
- The replica trick and algorithms
- Alternative forms of the replica identity

2

Application: information storage in neural networks

- Attractor neural networks
- The replica calculation
- Replica symmetry
- **Replica symmetric solution**

3

Application: overfitting transition in linear separators

- Linear separability of data – version space
- The replica calculation
- Gardner's replica symmetric theory

Replica symmetric solution

$m_\alpha = m$, $q_{\alpha \neq \beta} = q$, now find m and q ...

- **RS saddle-point eqns**

insert RS form and use $\exp(\frac{1}{2}x^2) = \int Dz e^{xz}$

$$\begin{aligned}m &= \frac{\sum_{\sigma} \sigma_\alpha e^{\beta J_0 m \sum_\lambda \sigma_\lambda + \frac{1}{2} (\beta J)^2 q \sum_{\lambda \neq \zeta} \sigma_\lambda \sigma_\zeta}}{\sum_{\sigma} e^{\beta J_0 m \sum_\lambda \sigma_\lambda + \frac{1}{2} (\beta J)^2 q \sum_{\lambda \neq \zeta} \sigma_\lambda \sigma_\zeta}} = \frac{\sum_{\sigma} \sigma_\alpha e^{\beta J_0 m \sum_\lambda \sigma_\lambda + \frac{1}{2} (\beta J)^2 q [\sum_\lambda \sigma_\lambda]^2}}{\sum_{\sigma} e^{\beta J_0 m \sum_\lambda \sigma_\lambda + \frac{1}{2} (\beta J)^2 q [\sum_\lambda \sigma_\lambda]^2}} \\&= \frac{\int Dz \sum_{\sigma} \sigma_\alpha \prod_{\lambda=1}^n e^{\beta(J_0 m + Jz\sqrt{q})\sigma_\lambda}}{\int Dz \sum_{\sigma} \prod_{\lambda=1}^n e^{\beta(J_0 m + Jz\sqrt{q})\sigma_\lambda}} \\&= \frac{\int Dz \sinh[\beta(J_0 m + Jz\sqrt{q})] \cosh^{n-1}[\beta(J_0 m + Jz\sqrt{q})]}{\int Dz \cosh^n[\beta(J_0 m + Jz\sqrt{q})]}\end{aligned}$$

similarly

$$q = \frac{\int Dz \sinh^2[\beta(J_0 m + Jz\sqrt{q})] \cosh^{n-2}[\beta(J_0 m + Jz\sqrt{q})]}{\int Dz \cosh^n[\beta(J_0 m + Jz\sqrt{q})]}$$

- **the limit $n \rightarrow 0$**

$$m = \int Dz \tanh[\beta(J_0 m + Jz\sqrt{q})], \quad q = \int Dz \tanh^2[\beta(J_0 m + Jz\sqrt{q})]$$

RS equations for $m = \lim_{N \rightarrow \infty} \frac{1}{N} \sum_{i=1}^N \xi_i \overline{\langle \sigma_i \rangle}$
and $q = \lim_{N \rightarrow \infty} \frac{1}{N} \sum_{i=1}^N \overline{\langle \sigma_i \rangle^2}$

$$m = \int Dz \tanh[\beta(J_0 m + Jz\sqrt{q})], \quad q = \int Dz \tanh^2[\beta(J_0 m + Jz\sqrt{q})]$$

- bifurcations away from $(m, q) = (0, 0)$:

$$\begin{aligned} m &= \int Dz [\beta J_0 m + \beta J z \sqrt{q} + \mathcal{O}(m, \sqrt{q})^3] = \beta J_0 m + \dots \\ q &= \int Dz [\beta J_0 m + \beta J z \sqrt{q} + \mathcal{O}(m, \sqrt{q})^3]^2 = \int Dz (\beta J)^2 z^2 q + \dots \\ &= (\beta J)^2 q + \dots \end{aligned}$$

hence:

first continuous bifurcations away from $\mathbf{q} = \mathbf{m} = \mathbf{0}$,
as identified earlier, are the RS solutions

$$m = \lim_{N \rightarrow \infty} \frac{1}{N} \sum_{i=1}^N \xi_i \overline{\langle \sigma_i \rangle}, \quad q = \lim_{N \rightarrow \infty} \frac{1}{N} \sum_{i=1}^N \overline{\langle \sigma_i \rangle^2}$$

$$m = \int Dz \tanh[\beta(J_0 m + Jz\sqrt{q})], \quad q = \int Dz \tanh^2[\beta(J_0 m + Jz\sqrt{q})]$$

phase diagram

P: $m = q = 0$

random neuronal firing

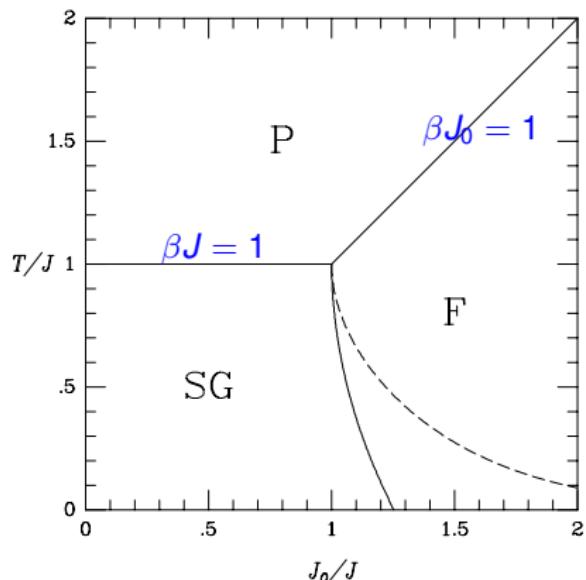
SG: $m = 0, q > 0$

stable firing patterns, but
not related to stored pattern

F: $m, q > 0$

recall of stored information

$T = 1/\beta$ (noise strength)



1

The replica method

- Exponential families and generating functions
- The replica trick
- The replica trick and algorithms
- Alternative forms of the replica identity

2

Application: information storage in neural networks

- Attractor neural networks
- The replica calculation
- Replica symmetry
- Replica symmetric solution

3

Application: overfitting transition in linear separators

- Linear separability of data – version space
- The replica calculation
- Gardner's replica symmetric theory

Linear separability of data and version space

Dimension mismatch and overfitting

two clinical outcomes (A,B),

4 patients, 60 expression levels ...

A : (1001010010100101010010001010111001001001001001000011111)
A : (01000100001010100101010100101000111100101001001010101000)
B : (0010100011101011011001001001110011100101001010101000101010)
B : (101011001010110010100100111100100101100111010111010001010010)

prognostic signature!

A : (1001010010100101010010001010111001001001001001000011111)
A : (010001000010101001010100101000111100101001001010101000)
B : (0010100011101011011001001001110011100101001010101000101010)
B : (101011001010110010100100111100100101100111010111010001010010)

shuffle outcome labels ...

A : (100101001010010101010010001010111001001001001000011111)
B : (0100010000101010010101010100101000111100101001001010101000)
A : (0010100011101011011001001001110011100101001010101000101010)
B : (101011001010110010100100111100100101100111010111010001010010)

*overfitting, no reproducibility ...
how about overfitting in regression?*

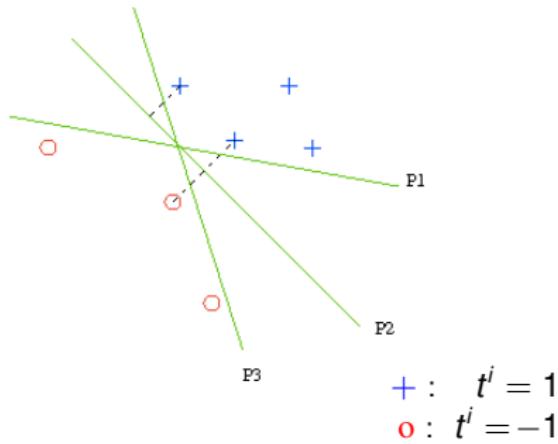
Suppose we have data D on N patients,
pairs of covariate vectors + clinical outcome labels

$$D = \{(\mathbf{x}^1, t^1), \dots, (\mathbf{x}^N, t^N)\}, \quad \mathbf{x}^i \in \{-1, 1\}^p, \quad t^i \in \{-1, 1\}, \quad p, N \gg 1$$

e.g. \mathbf{x}^i = gene expressions of i (on/off)
 t^i = treatment response of i (yes/no)

- assumed model:

$$\begin{aligned} t(\mathbf{x}) &= \begin{cases} 1 & \text{if } \sum_{\mu=1}^p \theta_\mu x_\mu > 0 \\ -1 & \text{if } \sum_{\mu=1}^p \theta_\mu x_\mu < 0 \end{cases} \\ &= \operatorname{sgn}\left[\sum_{\mu=1}^p \theta_\mu x_\mu\right] \end{aligned}$$



- regression/classification task:
find parameters $\boldsymbol{\theta} = (\theta_1, \dots, \theta_p)$
such that

$$\text{for all } i = 1 \dots N : \quad t^i = \operatorname{sgn}\left[\sum_{\mu=1}^p \theta_\mu x_\mu^i\right]$$

- data D explained perfectly by θ if

$$\text{for all } i = 1 \dots N : \quad t^i = \text{sgn}[\theta \cdot \mathbf{x}^i], \quad \text{i.e. } t^i(\theta \cdot \mathbf{x}^i) > 0$$

separating plane in input space : $\theta \cdot \mathbf{x} = 0$

distance Δ_i between \mathbf{x}^i and separating plane : $d_i = t^i(\theta \cdot \mathbf{x}^i)/|\theta|$

$|\theta|$ irrelevant, so choose $|\theta|^2 = p$

• version space

all θ that solve above eqns

with distances κ or larger

volume of version space:

$$V(\kappa) = \int d\theta \delta(\theta^2 - p) \prod_{i=1}^N \theta \left[\frac{t^i(\theta \cdot \mathbf{x}^i)}{\sqrt{p}} > \kappa \right]$$

- high dimensional data: p large, $\alpha = N/p$

$V(\kappa)$ scales exponentially with p , so

$$F = \frac{1}{p} \log V(\kappa) = \frac{1}{p} \log \int d\theta \delta(\theta^2 - p) \prod_{i=1}^N \theta \left[\frac{t^i(\theta \cdot \mathbf{x}^i)}{\sqrt{p}} > \kappa \right]$$

$F = -\infty$: no solutions θ exist, data D not linearly separable

F finite: solutions θ exist, data D linearly separable

overfitting: find parameters θ that ‘explain’ random patterns
 what if we choose **random data D** ?

$$D = \{(\mathbf{x}^1, t^1), \dots, (\mathbf{x}^N, t^N)\}, \quad \mathbf{x}^i \in \{-1, 1\}^p, \quad t^i \in \{-1, 1\}, \quad \text{fully random}$$

typical classification

performance:

$$\begin{aligned}\bar{F} &= \frac{1}{p} \log \int d\theta \delta(p - \theta^2) \prod_{i=1}^N \theta \left[\frac{t^i(\theta \cdot \mathbf{x}^i)}{\sqrt{p}} > \kappa \right] \\ &= \frac{1}{p} \log \int \frac{dz}{2\pi} e^{izp} \int d\theta e^{-iz\theta^2} \prod_{i=1}^N \theta \left[\frac{t^i(\theta \cdot \mathbf{x}^i)}{\sqrt{p}} > \kappa \right]\end{aligned}$$

transport data vars to harmless place,
 using δ -functions, by inserting

$$1 = \int dy_i \delta \left[y_i - \frac{t^i(\theta \cdot \mathbf{x}^i)}{\sqrt{p}} \right] = \int \frac{dy_i d\hat{y}_i}{2\pi} e^{i\hat{y}_i y_i - i\hat{y}_i t^i(\theta \cdot \mathbf{x}^i) / \sqrt{p}}$$

gives

$$\bar{F} = \frac{1}{p} \log \int \frac{dz dy d\hat{y} d\theta}{(2\pi)^{N+1}} e^{izp + i\hat{y} \cdot \mathbf{y} - iz\theta^2} \left(\prod_{i=1}^N \theta(y_i - \kappa) e^{-i\hat{y}_i t^i(\theta \cdot \mathbf{x}^i) / \sqrt{p}} \right)$$

1

The replica method

- Exponential families and generating functions
- The replica trick
- The replica trick and algorithms
- Alternative forms of the replica identity

2

Application: information storage in neural networks

- Attractor neural networks
- The replica calculation
- Replica symmetry
- Replica symmetric solution

3

Application: overfitting transition in linear separators

- Linear separability of data – version space
- **The replica calculation**
- Gardner's replica symmetric theory

The replica calculation

large p , large N ,
 $N = \alpha p$:

$$\overline{F} = \lim_{p \rightarrow \infty} \frac{1}{p} \log \overline{\int \frac{dz d\mathbf{y} d\hat{\mathbf{y}} d\theta}{(2\pi)^{N+1}} e^{izp + i\hat{\mathbf{y}} \cdot \mathbf{y} - iz\theta^2} \left(\prod_{i=1}^N \theta(y_i - \kappa) e^{-i\hat{y}_i \textcolor{blue}{t}^i (\theta \cdot \mathbf{x}^i) / \sqrt{p}} \right)}$$

• replica identity

$$\overline{\log Z} = \lim_{n \rightarrow 0} n^{-1} \log \overline{Z^n}$$

$$\begin{aligned} \overline{F} &= \lim_{p \rightarrow \infty} \lim_{n \rightarrow 0} \frac{1}{pn} \log \overline{\left[\int \frac{dz d\mathbf{y} d\hat{\mathbf{y}} d\theta}{(2\pi)^{N+1}} e^{izp + i\hat{\mathbf{y}} \cdot \mathbf{y} - iz\theta^2} \left(\prod_{i=1}^N \theta(y_i - \kappa) e^{-i\hat{y}_i \textcolor{blue}{t}^i (\theta \cdot \mathbf{x}^i) / \sqrt{p}} \right) \right]^n} \\ &= \lim_{p \rightarrow \infty} \lim_{n \rightarrow 0} \frac{1}{pn} \log \int \prod_{\alpha=1}^n \left[\frac{dz^\alpha d\mathbf{y}^\alpha d\hat{\mathbf{y}}^\alpha d\theta^\alpha}{(2\pi)^{N+1}} e^{ipz^\alpha + i\hat{\mathbf{y}}^\alpha \cdot \mathbf{y}^\alpha - iz^\alpha \theta^\alpha} \right] \prod_{i=1}^N \theta[y_i^\alpha - \kappa] \\ &\quad \times e^{-i \sum_{i=1}^N \sum_{\alpha=1}^n \hat{y}_i^\alpha \textcolor{blue}{t}^i (\theta^\alpha \cdot \mathbf{x}^i) / \sqrt{p}} \end{aligned}$$

● average over data D :

$$\begin{aligned}
 \Xi &= \overline{e^{-i \sum_{i=1}^N \sum_{\alpha=1}^n \hat{y}_i^\alpha \textcolor{blue}{t}^\mu (\boldsymbol{\theta}^\alpha \cdot \mathbf{x}_i^\mu) / \sqrt{p}}} = \overline{e^{-i \sum_{\mu=1}^p \sum_{i=1}^N \textcolor{blue}{t}^\mu \hat{y}_i^\mu \sum_{\alpha=1}^n \hat{y}_i^\alpha \theta_\mu^\alpha / \sqrt{p}}} \\
 &= \prod_{\mu=1}^p \prod_{i=1}^N \overline{e^{-i \textcolor{blue}{t}^\mu \hat{y}_i^\mu \sum_{\alpha=1}^n \hat{y}_i^\alpha \theta_\mu^\alpha / \sqrt{p}}} = \prod_{\mu=1}^p \prod_{i=1}^N \cos \left[\frac{1}{\sqrt{p}} \sum_{\alpha=1}^n \hat{y}_i^\alpha \theta_\mu^\alpha \right] \\
 &= \prod_{\mu=1}^p \prod_{i=1}^N \left\{ 1 - \frac{1}{2p} \left(\sum_{\alpha=1}^n \hat{y}_i^\alpha \theta_\mu^\alpha \right)^2 + \mathcal{O}\left(\frac{1}{p^2}\right) \right\} = e^{-\frac{1}{2p} \sum_{\mu=1}^p \sum_{i=1}^N \sum_{\alpha,\beta=1}^n \hat{y}_i^\alpha \hat{y}_i^\beta \theta_\mu^\alpha \theta_\mu^\beta} + \mathcal{O}(p^0)
 \end{aligned}$$

● giving

$$\begin{aligned}
 \mathcal{F} &= \lim_{p \rightarrow \infty} \lim_{n \rightarrow 0} \frac{1}{pn} \log \int \prod_{\alpha=1}^n \left[\frac{dz^\alpha d\mathbf{y} d\hat{\mathbf{y}}^\alpha d\boldsymbol{\theta}^\alpha}{(2\pi)^{N+1}} e^{ipz^\alpha + i\hat{\mathbf{y}}^\alpha \cdot \mathbf{y}^\alpha - iz^\alpha (\boldsymbol{\theta}^\alpha)^2} \prod_{i=1}^N \theta(y_i^\alpha - \kappa) \right] \\
 &\quad \times e^{-\frac{1}{2p} \sum_{\mu=1}^p \sum_{i=1}^N \sum_{\alpha,\beta=1}^n \hat{y}_i^\alpha \hat{y}_i^\beta \theta_\mu^\alpha \theta_\mu^\beta} + \mathcal{O}(p^0) \\
 &= -\alpha \log(2\pi) + \lim_{p \rightarrow \infty} \lim_{n \rightarrow 0} \frac{1}{pn} \log \int \prod_{\alpha=1}^n \left(dz^\alpha d\boldsymbol{\theta}^\alpha e^{ipz^\alpha - iz^\alpha (\boldsymbol{\theta}^\alpha)^2} \right) \\
 &\quad \times \prod_{i=1}^N \int \prod_{\alpha=1}^n \left[dy_i^\alpha d\hat{y}_i^\alpha e^{i \sum_{\alpha} \hat{y}_i^\alpha y_i^\alpha} \theta[y_i^\alpha - \kappa] \right] e^{-\frac{1}{2} \sum_{\alpha,\beta} \hat{y}_i^\alpha \hat{y}_i^\beta [\frac{1}{p} \sum_{\mu=1}^p \theta_\mu^\alpha \theta_\mu^\beta]}
 \end{aligned}$$

- so, with $\mathbf{y} = (y_1, \dots, y_n)$, $\hat{\mathbf{y}} = (\hat{y}_1, \dots, \hat{y}_n)$,
 $\mathbf{z} = (z_1, \dots, z_n)$:

$$\begin{aligned}\bar{F} &= -\alpha \log(2\pi) + \lim_{p \rightarrow \infty} \lim_{n \rightarrow 0} \frac{1}{pn} \log \int d\mathbf{z} \left(\prod_{\alpha=1}^n d\theta^\alpha e^{ipz^\alpha - iz^\alpha(\boldsymbol{\theta}^\alpha)^2} \right) \\ &\quad \times \left\{ \int d\mathbf{y} d\hat{\mathbf{y}} e^{i\hat{\mathbf{y}} \cdot \mathbf{y}} \prod_{\alpha=1}^n \theta^\alpha [y^\alpha - \kappa] e^{-\frac{1}{2} \sum_{\alpha, \beta} \hat{y}^\alpha \hat{y}^\beta [\frac{1}{p} \sum_{\mu=1}^p \theta_\mu^\alpha \theta_\mu^\beta]} \right\}^N\end{aligned}$$

- insert

$$1 = \int d\mathbf{q}_{\alpha\beta} \delta \left[\mathbf{q}_{\alpha\beta} - \frac{1}{p} \sum_{\mu=1}^p \theta_\mu^\alpha \theta_\mu^\beta \right] = \int \frac{d\mathbf{q}_{\alpha\beta} d\hat{\mathbf{q}}_{\alpha\beta}}{2\pi/p} e^{ip\hat{\mathbf{q}}_{\alpha\beta} \left[\mathbf{q}_{\alpha\beta} - \frac{1}{p} \sum_{\mu=1}^p \theta_\mu^\alpha \theta_\mu^\beta \right]}$$

to get

$$\begin{aligned}\bar{F} &= -\alpha \log(2\pi) + \lim_{p \rightarrow \infty} \lim_{n \rightarrow 0} \frac{1}{pn} \log \int d\mathbf{z} d\mathbf{q} d\hat{\mathbf{q}} e^{ip \sum_{\alpha\beta=1}^n \hat{q}_{\alpha\beta} q_{\alpha\beta} + ip \sum_{\alpha=1}^n z_\alpha} \\ &\quad \times \left\{ \int d\mathbf{y} d\hat{\mathbf{y}} e^{i\hat{\mathbf{y}} \cdot \mathbf{y}} \prod_{\alpha=1}^n \theta^\alpha [y^\alpha - \kappa] e^{-\frac{1}{2} \hat{\mathbf{y}} \cdot \mathbf{q} \hat{\mathbf{y}}} \right\}^N \int \prod_{\alpha=1}^n \left(d\theta^\alpha e^{-iz^\alpha(\boldsymbol{\theta}^\alpha)^2} \right) e^{-i \sum_{\mu=1}^p \sum_{\alpha\beta} \hat{q}_{\alpha\beta} \theta_\mu^\alpha \theta_\mu^\beta}\end{aligned}$$

- so, with $\theta = (\theta_1, \dots, \theta_n)$:

(remember: $N = \alpha p$)

$$\begin{aligned}\bar{F} &= -\alpha \log(2\pi) + \lim_{p \rightarrow \infty} \lim_{n \rightarrow 0} \frac{1}{pn} \log \int d\mathbf{z} d\mathbf{q} d\hat{\mathbf{q}} e^{ip \sum_{\alpha \beta=1}^n \hat{q}_{\alpha \beta} q_{\alpha \beta} + ip \sum_{\alpha=1}^n z_{\alpha}} \\ &\quad \times \left\{ \int d\mathbf{y} d\hat{\mathbf{y}} e^{i\hat{\mathbf{y}} \cdot \mathbf{y}} \prod_{\alpha=1}^n \theta[y^{\alpha} - \kappa] e^{-\frac{1}{2} \hat{\mathbf{y}} \cdot \mathbf{q} \hat{\mathbf{y}}} \right\}^{\alpha p} \left\{ \int d\theta e^{-i \sum_{\alpha=1}^n z^{\alpha} \theta_{\alpha}^2 - i \theta \cdot \mathbf{q} \theta} \right\}^p\end{aligned}$$

$$= \lim_{p \rightarrow \infty} \lim_{n \rightarrow 0} \frac{1}{pn} \log \int d\mathbf{z} d\mathbf{q} d\hat{\mathbf{q}} e^{p\Psi(\mathbf{z}, \mathbf{q}, \hat{\mathbf{q}})}$$

$$\begin{aligned}\Psi(\dots) &= i \sum_{\alpha \beta=1}^n \hat{q}_{\alpha \beta} q_{\alpha \beta} + i \sum_{\alpha=1}^n z_{\alpha} + \log \int d\theta e^{-i \sum_{\alpha=1}^n z^{\alpha} \theta_{\alpha}^2 - i \theta \cdot \mathbf{q} \theta} \\ &\quad + \alpha \log \int d\mathbf{y} d\hat{\mathbf{y}} e^{i\hat{\mathbf{y}} \cdot \mathbf{y}} \prod_{\alpha=1}^n \theta[y^{\alpha} - \kappa] e^{-\frac{1}{2} \hat{\mathbf{y}} \cdot \mathbf{q} \hat{\mathbf{y}}} - \alpha n \log(2\pi)\end{aligned}$$

- assume limits $n \rightarrow 0$ and $p \rightarrow \infty$ commute,
steepest descent integration

$$\bar{F} = \lim_{n \rightarrow 0} \frac{1}{n} \text{extr}_{\mathbf{z}, \mathbf{q}, \hat{\mathbf{q}}} \Psi(\mathbf{z}, \mathbf{q}, \hat{\mathbf{q}})$$

$$\begin{aligned}
\Psi(\mathbf{z}, \mathbf{q}, \hat{\mathbf{q}}) &= i \sum_{\alpha \beta=1}^n \hat{q}_{\alpha \beta} q_{\alpha \beta} + i \sum_{\alpha=1}^n z_{\alpha} + \log \int d\boldsymbol{\theta} e^{-i \sum_{\alpha=1}^n z^{\alpha} \theta_{\alpha}^2 - i \boldsymbol{\theta} \cdot \mathbf{q} \boldsymbol{\theta}} \\
&\quad + \alpha \log \int d\mathbf{y} d\hat{\mathbf{y}} e^{i \hat{\mathbf{y}} \cdot \mathbf{y}} \prod_{\alpha=1}^n \theta[y^{\alpha} - \kappa] e^{-\frac{1}{2} \hat{\mathbf{y}} \cdot \mathbf{q} \hat{\mathbf{y}}} - \alpha n \log(2\pi)
\end{aligned}$$

- transform $\hat{q}_{\alpha \beta} = -\frac{1}{2} i k_{\alpha \beta} - z_{\alpha} \delta_{\alpha \beta}$,
and integrate over $\hat{\mathbf{y}}$:

$$\begin{aligned}
\Psi(\mathbf{z}, \mathbf{q}, \mathbf{k}) &= \frac{1}{2} \sum_{\alpha \beta=1}^n k_{\alpha \beta} q_{\alpha \beta} + i \sum_{\alpha=1}^n z_{\alpha} (1 - q_{\alpha \alpha}) + \log \int d\boldsymbol{\theta} e^{-\frac{1}{2} \boldsymbol{\theta} \cdot \mathbf{k} \boldsymbol{\theta}} \\
&\quad + \alpha \log \int d\mathbf{y} \prod_{\alpha=1}^n \theta[y^{\alpha} - \kappa] \int d\hat{\mathbf{y}} e^{i \hat{\mathbf{y}} \cdot \mathbf{y} - \frac{1}{2} \hat{\mathbf{y}} \cdot \mathbf{q} \hat{\mathbf{y}}} - \alpha n \log(2\pi) \\
&= \frac{1}{2} \sum_{\alpha \beta=1}^n k_{\alpha \beta} q_{\alpha \beta} + i \sum_{\alpha=1}^n z_{\alpha} (1 - q_{\alpha \alpha}) + \log \frac{(2\pi)^{n/2}}{\sqrt{\text{Det} \mathbf{k}}} \\
&\quad + \alpha \log \int d\mathbf{y} \prod_{\alpha=1}^n \theta[y^{\alpha} - \kappa] \frac{(2\pi)^{n/2}}{\sqrt{\text{Det} \mathbf{q}}} e^{-\frac{1}{2} \mathbf{y} \cdot \mathbf{q}^{-1} \mathbf{y}} - \alpha n \log(2\pi)
\end{aligned}$$

- re-organise:

$$\begin{aligned}\Psi(\mathbf{z}, \mathbf{q}, \mathbf{k}) = & \frac{1}{2} \sum_{\alpha\beta=1}^n k_{\alpha\beta} q_{\alpha\beta} + i \sum_{\alpha=1}^n z_\alpha (1 - q_{\alpha\alpha}) - \frac{1}{2} \log \text{Det } \mathbf{k} - \frac{1}{2} \alpha \log \text{Det } \mathbf{q} \\ & + \alpha \log \int d\mathbf{y} \prod_{\alpha=1}^n \theta[y^\alpha - \kappa] e^{-\frac{1}{2} \mathbf{y} \cdot \mathbf{q}^{-1} \mathbf{y}} + \frac{1}{2} n(1-\alpha) \log(2\pi)\end{aligned}$$

- extremise with respect to \mathbf{z} :

$$\partial\Psi/\partial z_\alpha = 0 : q_{\alpha\alpha} = 0 \text{ for all } \alpha$$

$$\begin{aligned}\Psi(\mathbf{q}, \mathbf{k}) = & \frac{1}{2} n(1-\alpha) \log(2\pi) + \frac{1}{2} \sum_{\alpha\beta=1}^n k_{\alpha\beta} q_{\alpha\beta} - \frac{1}{2} \log \text{Det } \mathbf{k} - \frac{1}{2} \alpha \log \text{Det } \mathbf{q} \\ & + \alpha \log \int d\mathbf{y} \prod_{\alpha=1}^n \theta[y^\alpha - \kappa] e^{-\frac{1}{2} \mathbf{y} \cdot \mathbf{q}^{-1} \mathbf{y}}\end{aligned}$$

next: ergodicity assumption,
replica-symmetric form for \mathbf{q} and \mathbf{k} ...

1

The replica method

- Exponential families and generating functions
- The replica trick
- The replica trick and algorithms
- Alternative forms of the replica identity

2

Application: information storage in neural networks

- Attractor neural networks
- The replica calculation
- Replica symmetry
- Replica symmetric solution

3

Application: overfitting transition in linear separators

- Linear separability of data – version space
- The replica calculation
- Gardner's replica symmetric theory

Gardner's replica symmetric theory

$$\begin{aligned}\Psi(\mathbf{q}, \mathbf{k}) = & \frac{1}{2}n(1-\alpha)\log(2\pi) + \frac{1}{2} \sum_{\alpha,\beta=1}^n k_{\alpha\beta}q_{\alpha\beta} - \frac{1}{2}\log \text{Det } \mathbf{k} - \frac{1}{2}\alpha \log \text{Det } \mathbf{q} \\ & + \alpha \log \int d\mathbf{y} \prod_{\alpha=1}^n \theta[y^\alpha - \kappa] e^{-\frac{1}{2}\mathbf{y} \cdot \mathbf{q}^{-1} \mathbf{y}}\end{aligned}$$

RS saddle-points

$$q_{\alpha\beta} = \delta_{\alpha\beta} + (1-\delta_{\alpha\beta})q, \quad k_{\alpha\beta} = K\delta_{\alpha\beta} + (1-\delta_{\alpha\beta})k$$

- eigenvalues:

$$\mathbf{x} = (1, \dots, 1) : \quad (\mathbf{kx})_\alpha = \sum_{\beta=1}^n [k + (K-k)\delta_{\alpha\beta}]x_\beta = nk + K - k$$

eigenvalue : $\lambda = nk + K - k$

$$\sum_{\alpha=1}^n x_\alpha = 0 : \quad (\mathbf{kx})_\alpha = \sum_{\beta=1}^n [k + (K-k)\delta_{\alpha\beta}]x_\beta = (K-k)x_\alpha$$

eigenvalue : $\lambda = K - k$ ($n-1$ fold)

hence

$$\text{Det } \mathbf{k} = (nk + K - k)(K - k)^{n-1}, \quad \text{Det } \mathbf{q} = (nq + 1 - q)(1 - q)^{n-1}$$

- invert \mathbf{q} , try $(\mathbf{q}^{-1})_{\alpha\beta} = r + (R - r)\delta_{\alpha\beta}$,
demand:

$$\begin{aligned}\delta_{\alpha\beta} &= (\mathbf{q}\mathbf{q}^{-1})_{\alpha\beta} = \sum_{\gamma} (q + (1-q)\delta_{\alpha\gamma})(r + (R-r)\delta_{\gamma\beta}) \\ &= nqr + q(R-r) + r(1-q) + (R-r)(1-q)\delta_{\alpha\beta}\end{aligned}$$

so $nqr + q(R-r) + r(1-q) = 0, \quad (R-r)(1-q) = 1$

$$R = r + \frac{1}{1-q}, \quad r = -\frac{q}{(1-q)(1-q+nq)}$$

- hence, using $\exp[\frac{1}{2}x^2] = \int Dz e^{xz}$

$$\begin{aligned}\log \int d\mathbf{y} \prod_{\alpha=1}^n \theta[y^\alpha - \kappa] e^{-\frac{1}{2}\mathbf{y} \cdot \mathbf{q}^{-1} \mathbf{y}} &= \log \int d\mathbf{y} \prod_{\alpha=1}^n \theta[y^\alpha - \kappa] e^{-\frac{1}{2} \sum_{\alpha\beta} y_\alpha [r + (R-r)\delta_{\alpha\beta}] y_\beta} \\ &= \log \int d\mathbf{y} \prod_{\alpha=1}^n \theta[y^\alpha - \kappa] e^{-\frac{1}{2} r [\sum_{\alpha\beta} y_\alpha]^2 - \frac{1}{2} (R-r) \sum_{\alpha} y_\alpha^2} \\ &= \log \int Dz \int d\mathbf{y} \prod_{\alpha=1}^n \theta[y^\alpha - \kappa] e^{z \sqrt{-r} \sum_{\alpha\beta} y_\alpha - \frac{1}{2} (R-r) \sum_{\alpha} y_\alpha^2} \\ &= \log \int Dz \left[\int_{\kappa}^{\infty} dy e^{z \sqrt{-r} y - \frac{1}{2} (R-r) y^2} \right]^n\end{aligned}$$

put everything together ...

$$\begin{aligned}
 \frac{1}{n} \Psi(\mathbf{q}, \mathbf{k}) &= \frac{1}{2}(1-\alpha) \log(2\pi) + \frac{1}{2}K + \frac{1}{2}(n-1)qk - \frac{1}{2n} \log[(nk+K-k)(K-k)^{n-1}] \\
 &\quad - \frac{\alpha}{2n} \log[(nq+1-q)(1-q)^{n-1}] + \frac{\alpha}{n} \log \int Dz \left[\int_{\kappa}^{\infty} dy e^{zy\sqrt{-r}y - \frac{1}{2}(R-r)y^2} \right]^n \\
 &= \frac{1}{2}(1-\alpha) \log(2\pi) + \frac{1}{2}(K-qk) - \frac{1}{2n} \log(1+\frac{nk}{K-k}) - \frac{1}{2} \log(K-k) \\
 &\quad - \frac{\alpha}{2n} \log(1+\frac{nq}{1-q}) - \frac{\alpha}{2} \log(1-q) + \mathcal{O}(n) \\
 &\quad + \frac{\alpha}{n} \log \int Dz \left[1 + n \log \int_{\kappa}^{\infty} dy e^{zy\sqrt{q}/(1-q) - \frac{1}{2(1-q)}y^2} + \mathcal{O}(n^2) \right]
 \end{aligned}$$

take limit $n \rightarrow 0$:

$$\begin{aligned}
 2\bar{F} &= (1-\alpha) \log(2\pi) + \text{extr}_{K,k,q} \left\{ K-qk - \frac{k}{K-k} - \log(K-k) \right. \\
 &\quad \left. - \frac{\alpha q}{1-q} - \alpha \log(1-q) + 2\alpha \int Dz \log \int_{\kappa}^{\infty} dy e^{zy\sqrt{q}/(1-q) - \frac{1}{2(1-q)}y^2} \right\}
 \end{aligned}$$

$$2\bar{F} = (1-\alpha) \log(2\pi) + \text{extr}_{K,k,q} \left\{ K - qk - \frac{k}{K-k} - \log(K-k) \right. \\ \left. - \frac{\alpha q}{1-q} - \alpha \log(1-q) + 2\alpha \int Dz \log \int_{\kappa}^{\infty} dy e^{zy\sqrt{q}/(1-q) - \frac{1}{2(1-q)}y^2} \right\}$$

• **extremise over K and k**

$$\begin{cases} \frac{\partial}{\partial K} = 0 : 1 + \frac{k}{(K-k)^2} - \frac{1}{K-k} = 0 \\ \frac{\partial}{\partial k} = 0 : -q - \frac{1}{K-k} - \frac{k}{(K-k)^2} + \frac{1}{K-k} = 0 \end{cases} \Rightarrow K = \frac{1-2q}{(1-q)^2}, \quad k = -\frac{q}{(1-q)^2}$$

result:

$$2\bar{F} = (1-\alpha) \log(2\pi) + \text{extr}_q \left\{ \frac{1}{1-q} - \frac{\alpha q}{1-q} + (1-\alpha) \log(1-q) \right. \\ \left. + 2\alpha \int Dz \log \int_{\kappa}^{\infty} dy e^{zy\sqrt{q}/(1-q) - \frac{1}{2(1-q)}y^2} \right\}$$

• write y -integral in terms of error function $\text{Erf}(x) = \frac{2}{\sqrt{\pi}} \int_0^x dt e^{-t^2}$:

$$\int_{\kappa}^{\infty} dy e^{zy\sqrt{q}/(1-q) - \frac{1}{2(1-q)}y^2} = e^{\frac{qz^2}{2(1-q)}} \int_{\kappa}^{\infty} dy e^{-\frac{|y-z\sqrt{q}|^2}{2(1-q)}} \\ = \sqrt{2(1-q)} e^{\frac{qz^2}{2(1-q)}} \frac{\sqrt{\pi}}{2} \left\{ 1 - \text{Erf} \left[\frac{K - z\sqrt{q}}{\sqrt{2(1-q)}} \right] \right\}$$

- insert previous integral:

$$2\bar{F} = \log \pi + (1-2\alpha) \log 2$$

$$+ \text{extr}_q \left\{ \frac{1}{1-q} + \log(1-q) + 2\alpha \int Dz \log \left[1 - \text{Erf} \left(\frac{\kappa - z\sqrt{q}}{\sqrt{2(1-q)}} \right) \right] \right\}$$

- extremisation with respect to q

short-hand $u(z, q) = (\kappa - z\sqrt{q}) / \sqrt{2(1-q)}$,

use $\text{Erf}'(x) = \frac{2}{\sqrt{\pi}} \exp[-x^2]$

$$\frac{d}{dq} = 0 : \quad \frac{1}{(1-q)^2} - \frac{1}{1-q} - 2\alpha \int Dz \left(\frac{\partial u}{\partial q} \right) \frac{\text{Erf}' u(z, q)}{1 - \text{Erf} u(z, q)} = 0$$

$$\frac{q}{(1-q)^2} = \frac{4\alpha}{\sqrt{\pi}} \int Dz \left(\frac{\partial u}{\partial q} \right) \frac{e^{-u^2(z,q)}}{1 - \text{Erf} u(z, q)}$$

work out:

$$\frac{\partial u}{\partial q} = \frac{1}{\sqrt{2}} \frac{\partial}{\partial q} \frac{\kappa - z\sqrt{q}}{(1-q)^{1/2}} = \dots = \frac{\kappa\sqrt{q} - z}{2\sqrt{2q}(1-q)^{3/2}}$$

insert into eqn for q :

$$q\sqrt{q} = \alpha \sqrt{\frac{2}{\pi}} \sqrt{1-q} \int Dz \frac{e^{-u^2(z,q)}(\kappa\sqrt{q}-z)}{1 - \text{Erf} u(z, q)}$$

$$2\bar{F} = \log \pi + (1-2\alpha) \log 2 + \frac{1}{1-q} + \log(1-q) + 2\alpha \int Dz \log [1 - \text{Erf } u(z, q)]$$

$$q\sqrt{q} = \alpha \sqrt{\frac{2}{\pi}} \sqrt{1-q} \int Dz \frac{e^{-u^2(z,q)} (\kappa\sqrt{q}-z)}{1-\text{Erf } u(z, q)}, \quad u(z, q) = \frac{\kappa-z\sqrt{q}}{\sqrt{2(1-q)}}$$

remember:

\bar{F} = finite: random data linearly separable with margin κ

$\bar{F} = -\infty$: random data not linearly separable with margin κ

- $\alpha = 0$ (so $1 \ll N \ll p$):

$q = 0, \quad 2\bar{F} = \log \pi + \log 2 + 1 \quad \text{random data linearly separable (overfitting)}$

- $\alpha > 0$ (so $1 \ll N \sim p$):

transition point: value of α where $q \rightarrow 1$

$$1 = \alpha_c(\kappa) \sqrt{\frac{2}{\pi}} \int Dz \lim_{q \rightarrow 1} \sqrt{1-q} \frac{e^{-[\frac{\kappa-z}{\sqrt{2(1-q)}}]^2} (\kappa-z)}{1 - \text{Erf} \left[\frac{\kappa-z}{\sqrt{2(1-q)}} \right]}$$

$$\alpha_c(\kappa) = \left[\frac{1}{\sqrt{\pi}} \int Dz (\kappa+z) \lim_{\gamma \rightarrow \infty} \frac{1}{\gamma} \frac{e^{-\gamma^2(\kappa+z)^2}}{1 - \text{Erf}[\gamma(\kappa+z)]} \right]^{-1}$$

- remaining limit:

$$\lim_{\gamma \rightarrow \infty} \frac{1}{\gamma} \frac{e^{-\gamma^2 Q^2}}{1 - \text{Erf}[\gamma Q]} = Q\sqrt{\pi} \theta(Q)$$

proof:

$$Q < 0 : \quad \text{Erf}[\gamma Q] \rightarrow -1 \quad \text{so} \quad \lim_{\gamma \rightarrow \infty} \frac{1}{\gamma} \frac{e^{-\gamma^2 Q^2}}{1 - \text{Erf}[\gamma Q]} = 0$$

$$Q > 0 : \quad \text{Erf}[\gamma Q] = 1 - \frac{1}{\gamma Q \sqrt{\pi}} e^{-\gamma^2 Q^2} \left(1 + \mathcal{O}\left(\frac{1}{\gamma^2 Q^2}\right) \right)$$

$$\lim_{\gamma \rightarrow \infty} \frac{1}{\gamma} \frac{e^{-\gamma^2 Q^2}}{1 - \text{Erf}[\gamma Q]} = \lim_{\gamma \rightarrow \infty} \frac{1}{\gamma} \frac{e^{-\gamma^2 Q^2}}{\frac{1}{\gamma Q \sqrt{\pi}} e^{-\gamma^2 Q^2} \left(1 + \mathcal{O}\left(\frac{1}{\gamma^2 Q^2}\right) \right)} = Q\sqrt{\pi}$$

final result

$$\alpha_c(\kappa) = \left[\int_{-\kappa}^{\infty} Dz (\kappa+z)^2 \right]^{-1} \quad \alpha_c(0) = \left[\int_0^{\infty} Dz z^2 \right]^{-1} = \left[\frac{1}{2} \right]^{-1} = 2$$

p covariates,
 N patients,
binary outcomes,
 p and N large

random data
(i.e. pure binary noise)
is *perfectly separable* if
 $N/p < \alpha_c(\kappa)$

algorithms (SVM etc)
will find pars $\theta_1 \dots \theta_p$
such that $t_i = \text{sgn}[\sum_{\mu=1}^p \theta_\mu x_\mu^i]$
for all $i = 1 \dots N$

