# Implicit regularisation, Benign interpolation, and Over-parametrisation

**Main references:**

- M. Belkin, D. Hsu, S. Ma, S. Mandal. "Reconciling modern machine learning and the bias-variance trade-off. (2018)
- P. Bartlett, P. Long, G. Lugosi, A. Tsigler. "Benign overfitting in Linear Regression". (2019)
- T. Hastie, A. Montanari, S. Rosset, R. Tibshirani. "Surprises in high-dimensional ridgeless least-square interpolation". (2019)
- L. Chizat, F. Bach. "Implicit bias of gradient descent for wide 2-layer neural networks trained with the logistic loss" (2020)
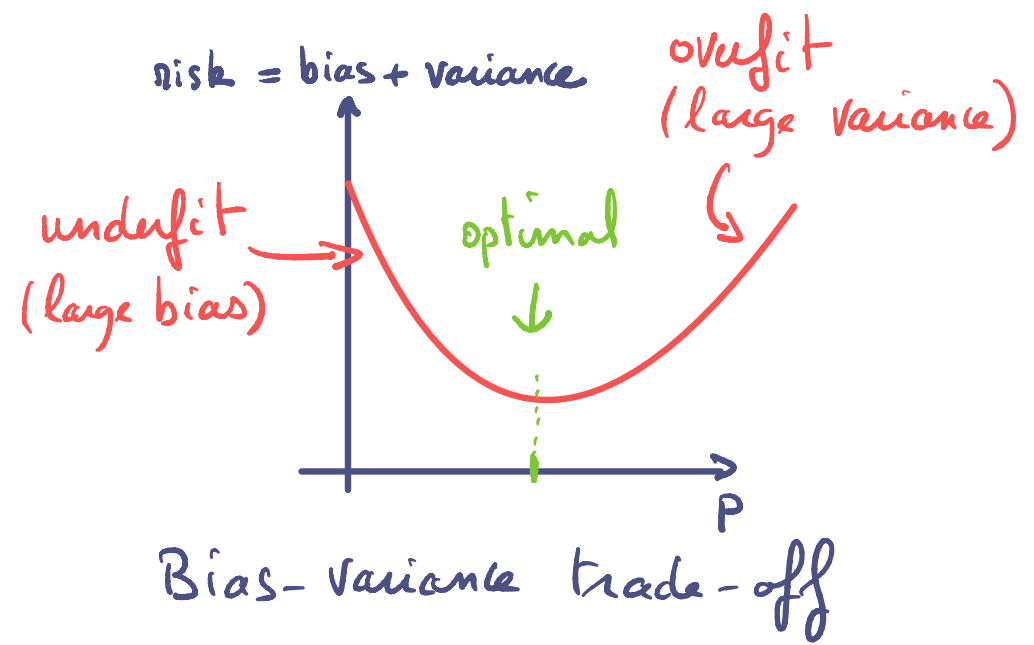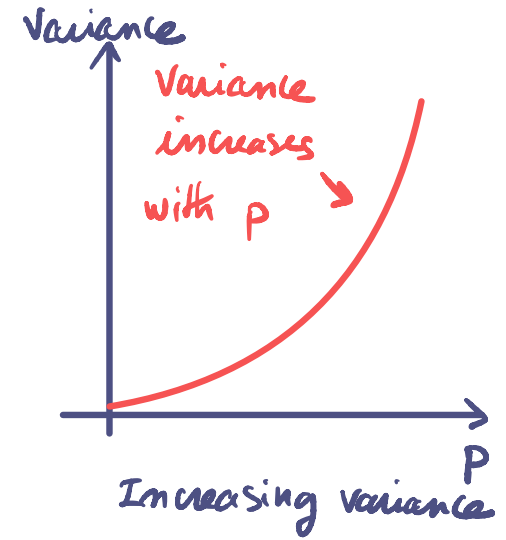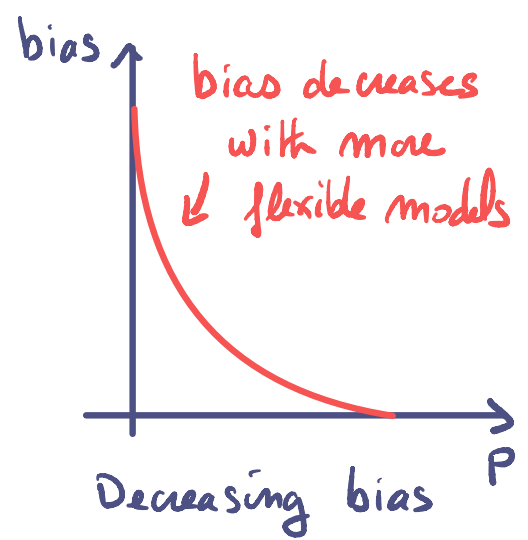
**Lecture notes:**

Lecture notes are available on the web page provided in the description of the video.

---

**Double descent phenomenon**

**Textbook theory:** risk = bias + variance

Let $p$ = number of parameters in the model.



bias decreases with more flexible models

Decreasing bias $P$

Variance increases with $p$

Increasing variance $P$

risk = bias + variance

underfit (large bias)

optimal

overfit (large variance)

$P$

Bias-variance trade-off

# Empirical observations (neural networks)

risk



the variance decreases after the interpolation point

↑ interpolation

P

Double descent phenomenon

# Four parts

① Implicit regularisation (a.k.a. implicit bias) of gradient descent

② Benign interpolation: A) intuitions

③ Benign interpolation: B) mathematical analysis

④ Benign overparametrisation

# 1. Implicit regularization of Gradient Descent

## a) Implicit regularisation for L.S.

- **Linear model**: $y_i = \langle x_i, \underset{\in \mathbb{R}^p}{\underline{\beta^*}} \rangle + \varepsilon_i$, $i = 1, \dots, m$

  with $\varepsilon_i \overset{iid}{\sim} N(0, \sigma^2)$.

- **Notation**: in vectorial notation

$$\underbrace{\begin{bmatrix} y_1 \\ \vdots \\ y_m \end{bmatrix}}_{=: Y} = \underbrace{\begin{bmatrix} x_1^T \\ \vdots \\ x_m^T \end{bmatrix}}_{=: X} \beta^* + \underbrace{\begin{bmatrix} \varepsilon_1 \\ \vdots \\ \varepsilon_m \end{bmatrix}}_{=: \varepsilon}$$

- **Least-Square**: MLE estimation

  amounts to minimise the squares

$$\hat{\beta}^{MLE} \in \underset{\beta \in \mathbb{R}^p}{\arg\min} \; \|Y - X\beta\|^2 \quad \text{(LS)}$$

- **High-dimensional regime**:

  Assume that $x_i \in \mathbb{R}^p$ with $p > m$.

  then, $\dim(\ker(X)) = p - \text{rank}(X) \geq p - m > 0$

  $\Rightarrow$ no unique solution to (LS)

- **Interpolation regime**: $p \geq m$

  when in addition $\text{rank}(X) = m$, then

  any solution $\hat{\beta}$ of (LS) fulfills

  $Y = X\hat{\beta}$ i.e $y_i = x_i^T \hat{\beta}$ for $i = 1, \dots, m$.

  $\Rightarrow$ $\hat{\beta}$ interpolates the learning data

  $(x_i, y_i)_{i = 1, \dots, m}$.

  $\Rightarrow$ overfitting ?!?

---

- **Classical practice**: add a regularisation

  term in (LS) like $\lambda \|\beta\|^2$ in order

  to get a strictly convex objective function

- **Recent practice in N.N.**: apply GD

  on (LS) and use the solution selected

  by GD. What is this solution?

· **G.D. on least-squares:**

We apply G.D. on $\beta \to \|Y - X\beta\|^2$, started from $\hat{\beta}^0 = 0$ :

$$\hat{\beta}^{t+1} \overset{GD}{=} \hat{\beta}^t - 2\eta \, X^T(X\hat{\beta}^t - Y)$$

$$\overset{\hat{\beta}^0 = 0}{=} \sum_{k=0}^{t} (I - 2\eta X^T X)^k \cdot 2\eta X^T Y$$

$$\overset{t \to +\infty}{\underset{2\eta < \|X\|_{op}^{-2}}{\longrightarrow}} \left(I - (I - 2\eta X^T X)\right)^+ \cdot 2\eta X^T Y$$

$$= X^+ Y$$

$\uparrow$ Moore Penrose pseudo inverse (Appendix C)

So $\hat{\beta}^{LS} = X^+ Y$.

· Is there something special with this solution?

---

**Lemma**: In the interpolation regime, where $\text{rank}(X) = n$, we have
$$\hat{\beta}^{LS} = \underset{Y = X\beta}{\text{argmin}} \, \|\beta\|^2$$

$\leftarrow$ <span style="color:red">regularisation of G.D.</span>

**Proof**: · we first observe that
$$X X^+ Y = Y \qquad \text{so} \quad Y = X \hat{\beta}^{LS}$$
and any solution of $Y = X\beta$ can be decomposed as $\beta = \beta_0 + X^+ Y$, with $\beta_0 \in \ker(X)$

· Then, we notice that $\text{range}(X^+) = \text{range}(X^T)$ so we have
$$\mathbb{R}^p = \ker(X) \oplus \text{range}(X^T) = \ker(X) \oplus \text{range}(X^+)$$
and in particular $\|\beta\|^2 = \|\beta_0\|^2 + \|X^+ Y\|^2$.

So $X^+ Y = \underset{Y = X\beta}{\text{argmin}} \, \|\beta\|^2$  □

## b) Implicit regularization for logistic

- When the labels $y_i \in \{-1, +1\}$, we can predict $y$ by $\text{sign}(\langle \hat{\beta}, x \rangle)$, where $\hat{\beta}$ is learnt by applying GD on the empirical logistic risk

$$\mathcal{L}(\beta) := \frac{1}{m} \sum_{i=1}^{m} \ell(-y_i \langle \beta, x_i \rangle)$$

with $\ell(-\mathfrak{z}) = \log(1 + e^{-\mathfrak{z}})$

- **Interpolation regime:** when the learning data $(x_i | y_i)_{i=1,\dots,m}$ can be separated by an hyperplan.
- **GD in interpolation regime:**

Assume that at step $t$ of GD, $\hat{\beta}^t$ perfectly classifies the training data $y_i \langle \hat{\beta}^t, x_i \rangle > 0$, for $i = 1, \dots, m$.

- Since

$$\mathcal{L}(\hat{\beta}^t) = \frac{1}{m} \sum_{i=1}^{m} \ell\left(-\|\hat{\beta}^t\| \cdot y_i \left\langle \frac{\hat{\beta}^t}{\|\hat{\beta}^t\|}, x_i \right\rangle\right)$$

with $\mathfrak{z} \to \ell(-\mathfrak{z})$ decreasing, the loss $\mathcal{L}$ can be further decreased by sending $\|\hat{\beta}^t\| \longrightarrow +\infty$.

- Since $\ell(-\mathfrak{z}) \overset{\mathfrak{z}\to+\infty}{\sim} e^{-\mathfrak{z}}$, we get

$$\mathcal{L}(\hat{\beta}^t) \underset{\nearrow}{\sim} \frac{1}{m} \sum_{i=1}^{m} \exp\left(-\|\hat{\beta}^t\| \, y_i \left\langle \frac{\hat{\beta}^t}{\|\hat{\beta}^t\|}, x_i \right\rangle\right)$$

when $\|\hat{\beta}^t\| \gg 1$

$$\hookrightarrow \approx \frac{N_{min}}{m} \exp\left(-\|\hat{\beta}^t\| \cdot \min_{i=1\dots m} y_i \left\langle \frac{\hat{\beta}^t}{\|\hat{\beta}^t\|}, x_i \right\rangle\right)$$

which suggests that $u^t = \hat{\beta}^t / \|\hat{\beta}^t\|$ tends to solve the max-margin problem

$$(\pi\pi) \quad \max_{\|u\|=1} \quad \underbrace{\min_{i=1,\dots,m} \, y_i \langle u, x_i \rangle}_{\text{margin}},$$

which eventually occurs (Soudry et. al. 2017)

**Theorem** (informal)

The normalized solution $\hat{\beta}^t/\|\hat{\beta}^t\|$ of GD on the empirical logistic risk converges to the max-margin classifier **(MM)** when the data is linearly separable.

---

Take home message:

when interpolation is possible, GD selects some specific interpolating solutions

↗ minimal norm interpolating solutions in $L^2$ regression

↘ max-margin solution in logistic regression

⤳ implicit regularisation of GD.

These results can be generalized to more complex models like Neural Networks.

Example (Informal):

When training an homogeneous NN with GD and logistic loss:

if the normalized weights converge, then, at the limit, they solve the max-margin problem.

---

Next videos:

⤳ on the statistical benefit of this implicit regularisation for L.S.

2. Benign interpolation with high-dimensional input.

# A/ Intuitions

**Reminder :**

• We consider the linear model:
$$y_i = \langle x_i, \underbrace{\beta^*}_{\in \mathbb{R}^p} \rangle + \varepsilon_i, \quad i=1,\ldots,m$$

with $\varepsilon_i \overset{iid}{\sim} \mathcal{N}(0,\sigma^2)$. We set $X = \begin{bmatrix} x_1^T \\ \vdots \\ x_m^T \end{bmatrix}$.

• We compute $\hat{\beta}$ by minimizing the Least Square $\beta \to \|Y - X\beta\|^2$ by GD:
$$\rightsquigarrow \hat{\beta} = X^+ Y.$$

• **Prediction:** for $x \in \mathbb{R}^p$
$$\langle \hat{\beta}, x \rangle = \langle X^+ Y, x \rangle$$

$Y = X\beta^* + \varepsilon \rightarrow$
$$= \langle \underbrace{X^+ X}_{= P_{X^T} \text{ (projection on range}(X^T))} \beta^*, x \rangle + \langle X^+ \varepsilon, x \rangle$$

$$= \underbrace{\langle \beta^*, x \rangle}_{\text{target}} - \underbrace{\langle (I - P_{X^T})\beta^*, x \rangle}_{\text{bias}} + \underbrace{\langle \varepsilon, (X^T)^+ x \rangle}_{\mathcal{N}(0, x^T(X^TX)^+ x)}$$

cond. on $X$

**Remarks:** $\underset{\|}{\text{range}}(X^+)$

1) if $x \perp \text{range}(X^T)$: $\langle \hat{\beta}, x \rangle = 0$
   $\Rightarrow$ large bias, but no variance

2) $\text{rank}(X) = m \ll p$:
   $$\|P_{X^T} x\| \ll \|x\| \quad \text{if} \quad x \sim \mathcal{N}(0, \Sigma)$$
   w.h.p. with $\Sigma \simeq I_p$.

   So only a small part of $x$ is used for the prediction $\langle \hat{\beta}, x \rangle = \langle \hat{\beta}, P_{X^T} x \rangle$

---

Below we discuss two cases

a) isotropic case: $\Sigma = I_p$

b) spike model $\Sigma = I_k + \underbrace{\rho}_{\rho \ll 1} I_{p-k}$

**Average prediction error:**

- We have $\hat{\beta} = X^+ Y = \underbrace{X^+ X}_{P_{X^T}} \beta^* + X^+ \mathcal{E}$   with $X^+ X =$ orthogonal projection on $\text{range}(X^T)$

- Let us compute the average prediction error, conditionally on the design $X$

$$R_X = \mathbb{E}\left[ \langle x, \hat{\beta} - \beta^* \rangle^2 \mid X \right] \quad \text{where } x \sim \mathcal{N}(0, \Sigma) \text{ and } \mathcal{E} \sim \mathcal{N}(0, \sigma^2 I_n)$$

$$R_X = \mathbb{E}_{x,\mathcal{E}}\left[ \langle x, \hat{\beta} - \beta^* \rangle^2 \right] \overset{\hat{\beta} = P_{X^T}\beta^* + X^+\mathcal{E}}{=} \mathbb{E}_{x,\mathcal{E}}\left[ \left( \langle \underbrace{(I - P_{X^T})}_{P_{\ker(x)}} \beta^*, x \rangle - \langle X^+\mathcal{E}, x \rangle \right)^2 \right]$$

$P_{\ker(x)}$ (projection on $\ker(x)$)

$$= \mathbb{E}_x\left[ \beta^{*T} P_{\ker(X)} x x^T P_{\ker(X)} \beta^* \right] + \sigma^2 \mathbb{E}_x\left[ x^T \underbrace{X^+(X^+)^T}_{} x \right] + 0$$

cross-term

$$= \langle x x^T, X^+(X^+)^T \rangle_F$$

$$= \underbrace{\beta^{*T} P_{\ker(x)} \Sigma P_{\ker(x)} \beta^*}_{B_X} + \underbrace{\sigma^2 \langle \Sigma, (X^T X)^+ \rangle_F}_{V_X} \; .$$

bias term                          variance term

a) **isotropic input**: $\Sigma = I_p$

· **Linear model**: $y_i = \langle x_i, \beta^* \rangle + \varepsilon_i$, $i = 1, \dots, m$

with $x_i \overset{iid}{\sim} N(0, \Sigma)$ and $\varepsilon_i \overset{iid}{\sim} N(0, \sigma^2)$

· **singular values of random matrices**: with $p > m$

Let $X = \begin{bmatrix} x_1^T \\ \vdots \\ x_m^T \end{bmatrix} \in \mathbb{R}^{m \times p}$

$$= \sum_{k=1}^{m} \hat{\sigma}_k \hat{u}_k \hat{v}_k^T$$

we have:

$$\mathbb{E}\left(\|X\|_F^2\right) = \begin{cases} \sum_{i,j} \mathbb{E}\left(X_{ij}^2\right) = mp \\[2mm] \sum_{k=1}^{m} \mathbb{E}\left[\hat{\sigma}_k^2\right] \leq m\, \mathbb{E}\left[\hat{\sigma}_1^2\right] \end{cases}$$

$$\Rightarrow \mathbb{E}\left[\hat{\sigma}_1^2\right] \geq p.$$

This suggests that $\hat{\sigma}_k$ scales like $\sqrt{p}$ when $p \gg m$.

It eventually occurs:

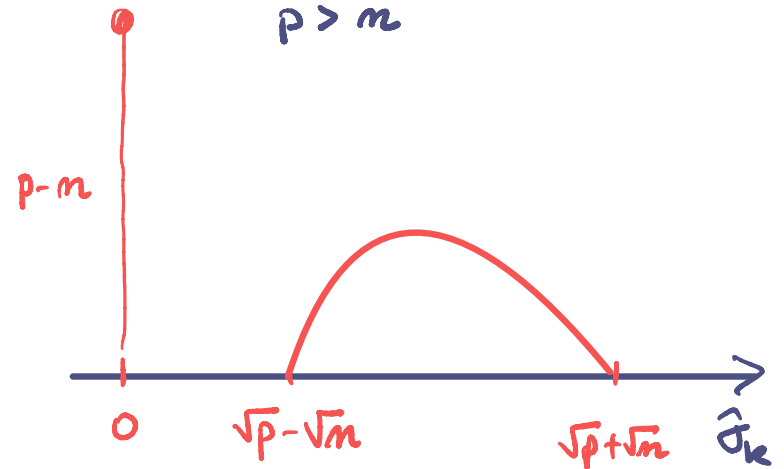## Lemma 8.3 (Davidson & Szarek)

For $p > m$, we have for $k = 1, \ldots, m$

$$\sqrt{p} - \sqrt{m} \leq \mathbb{E}[\hat{\sigma}_k] \leq \sqrt{p} + \sqrt{m}$$

Histogram of the $\hat{\sigma}_k$ for $p > m$



Furthermore, since $X \to \sigma_k(X)$ is 1-Lipschitz, by Gaussian concentration inequality, with high-probability

$$\hat{\sigma}_k \sim \sqrt{p} \qquad \text{for } p \gg m$$
$$\forall k = 1 \ldots m$$

$$\Rightarrow X^T X \overset{p \gg m}{\simeq} p \cdot P_{X^T}$$

$$\text{and } (X^T X)^+ \simeq \frac{1}{p} P_{X^T}$$

[ See Proposition 12.9 in the Lecture Notes ]

## Average prediction error :

$$\cdot \; R_x = \beta^{*T} P_{ker(x)} \Sigma \, P_{ker(x)} \beta^* + \sigma^2 \langle \Sigma, (x^T x)^+ \rangle_F$$

$$\overset{\Sigma = I_p}{=} \underbrace{\| P_{ker(x)} \beta^* \|^2}_{= B_x} + \underbrace{\sigma^2 \, Tr((x^T x)^+)}_{= V_x}$$

## Variance :

$$\cdot \; V_x = \sigma^2 \, Tr(x^T x)^+ = \sigma^2 \sum_{k=1}^{m} \frac{1}{\hat{\sigma}_k^2}$$

$$\hat{\sigma}_n^2 \approx p$$

w.h. $\mathbb{P}_x$ when $p \gg m$

$$\boxed{\sigma^2 \frac{m}{p}} \longrightarrow 0 \quad \text{when} \quad p \gg m \quad !!!$$

$\rightsquigarrow$ In the interpolation regime ($p \geqslant m$), high-dimensional inputs kill the variance of G.D.-Least square

[See next video for the proof of this result]

**Discussion :** we observe that for a given $x \in \mathbb{R}^p$, with $\|x\|^2 = p$ :

$$\to \operatorname{Var}_\varepsilon \langle \hat{\beta}, x \rangle = \sigma^2 x^T (X^T X)^+ x = \sigma^2 x^T \left( \sum_{k=1}^m \frac{1}{\hat{\sigma}_k^2} \hat{v}_k \hat{v}_k^T \right) x = \frac{\sigma^2}{p} \underbrace{x^T P_{X^T} x}_{= \|P_{X^T} x\|^2}$$

$$\to \text{ if } x \in \operatorname{range}(X^T): \qquad \operatorname{Var}_\varepsilon \langle \hat{\beta}, x \rangle \cong \sigma^2 \underbrace{\frac{\|x\|^2}{p}}_{\substack{p \gg m \\ \text{for } \|x\|^2 = p}} \approx \sigma^2 \leftarrow \quad \textcolor{red}{\text{strong overfit}}$$

$$\textcolor{red}{\text{on range}(X^T)}$$

$\to$ but $\|P_{X^T} x\|^2 \ll \|x\|^2$ w.h. $\mathbb{P}_x$ : due to the high-dimension of the input space, w.h. $\mathbb{P}_x$, a new $x$ is almost not correlated with the learning points $x_1, \dots, x_m$ (columns of $X^T$), so

$$\operatorname{Var}_\varepsilon \langle \hat{\beta}, x \rangle \cong \sigma^2 \underbrace{\frac{\|P_{X^T} x\|^2}{p}}_{} \quad \text{small}$$

$\to$ <u>in words</u> : we overfit on a small space = span $\{x_1, \dots, x_m\}$ spanned by the learning inputs, but we underfit everywhere else

$$\Rightarrow \text{ large bias and small variance for } p \gg m \, !$$

# Bias:

$$\mathbb{R}^p = \ker(x) \oplus \text{range}(x^T)$$

- $B_x = \| P_{\ker(x)} \beta^* \|^2 = \| \beta^* \|^2 - \| P_{X^T} \beta^* \|^2 \underset{\text{w.h. } \mathbb{P}_x}{\widetilde{\approx}} \| \beta^* \|^2 (1 - \frac{m}{p})$

the larger $p$, the larger the bias $B_x$ !!

$$\text{Range}(X^T) \overset{\text{unif.}}{\sim} \dim(m) \subset \mathbb{R}^p$$

$$\text{so} \quad \| P_{X^T} \beta \|^2 \approx \frac{m}{p} \| \beta \|^2$$

# Prediction risk:

- $R_x = B_x + V_x \overset{p \gg m}{\approx} \| \beta^* \|^2 + \frac{m}{p} (\sigma^2 - \| \beta^* \|^2)$

$\uparrow$

large

( mainly induced by the bias)

small since $m \ll p$

b) Anisotropic input (intuitions)

Prototypical example: $\Sigma = I_k + \rho\, I_{>k} = \begin{bmatrix} 1 & & & \\ & \ddots & 1 & \\ & & \rho & \\ & & & \ddots \\ & & & & \rho \end{bmatrix}$ with $\cdot\ k \ll m \ll p$
$\cdot\ \rho \ll 1$

where $\underbrace{\rho}_{\ll 1}$

Then $\frac{1}{m} X^T X = \frac{1}{m} \sum_{i=1}^{m} x_i x_i^T \overset{k \ll m}{\underset{\simeq}{\downarrow}} \underbrace{I_k}_{\text{low dimensional}} + \text{"something of rank } m-k\text{"}$

low dimensional
$\Rightarrow$ good estimation

$\text{Trace}(\Sigma) = k + \rho(p-k)$

$\simeq I_k + \rho\, \frac{p-k}{m-k}\, \overset{\downarrow}{\underbrace{\hat{I}_{m-k}}_{\simeq \text{ projection of rank } m-k}}$

$\cdot$ So $X^T X \overset{k \ll m \ll p}{\underset{\simeq}{\downarrow}} m\, I_k + \rho p\, \hat{I}_{m-k}$ and $(X^T X)^+ \overset{k \ll m \ll p}{\underset{\simeq}{\downarrow}} \frac{1}{m} I_k + \frac{1}{\rho p} \hat{I}_{m-k}$ — Hence:

$V_X = \sigma^2 \left\langle (X^T X)^+, \Sigma \right\rangle_F \cong \frac{\sigma^2}{m} \underbrace{\| I_k \|_F^2}_{= k} + \frac{\sigma^2}{p} \underbrace{\left\langle I_{>k}, \hat{I}_{m-k} \right\rangle_F}_{= \text{Tr}(\hat{I}_{m-k}) = m-k \simeq m}$

$\simeq \sigma^2 \left( \underbrace{\frac{k}{m}}_{\substack{\text{small due} \\ \text{to low rank } k}} + \underbrace{\frac{m}{p}}_{\substack{\text{small due} \\ \text{to high-dimensional input}}} \right)$

# Next video:

→ we consider a general $\Sigma$

→ we prove (sharp) upper-bound on $V_x$

( recovering today's "results" )

Benign interpolation

with high-dimensional input.

# B/ Mathematical analysis

Reminder:

- Learning data: $(x_i, y_i)_{i=1\dots n}$ with $x_i \sim N(0, \Sigma)$ and $y_i = \langle x_i, \beta^* \rangle + \underbrace{\varepsilon_i}_{\overset{iid}{\sim} N(0, \sigma^2)}$

- GD - least square: $\hat{\beta} = X^+ Y$

- Average prediction risk:

$$R_X := \mathbb{E}\left[ \langle \hat{\beta} - \beta^*, x \rangle^2 \mid X \right]$$

$$= \underbrace{\beta^{*T} P_{ker(X)} \Sigma P_{ker(X)} \beta^*}_{B_X} + \underbrace{\sigma^2 \langle \Sigma, (X^T X)^+ \rangle_F}_{V_X}$$

$$\underset{\text{bias term}}{B_X} \qquad\qquad \underset{\text{variance term}}{V_X}$$

$\rightsquigarrow$ to understand $V_X$, we need to understand $(X^T X)^+$.

a) spectrum of $XX^T$ : $p > m$

- $spec^*(XX^T) = spec^*(X^TX)$

- $rank(\underbrace{XX^T}_{m \times m}) = m$ for $p > m$

$\rightsquigarrow$ easier to analyse $XX^T$

- Setting : we assume (with no loss of generality) that $\Sigma = \begin{pmatrix} \lambda_1 & 0 \\ & \ddots \\ 0 & \lambda_p \end{pmatrix}$ with $\lambda_1 \geqslant \lambda_2 \geqslant \cdots \geqslant \lambda_p > 0$.

- We define :

$\rightarrow k^* = \min \{k : \lambda_{k+1} < \frac{\rho}{m} \sum_{j>k} \lambda_j \}$
  with $\rho = \frac{1}{142}$

$\rightarrow X = [X_\leqslant \; X_>]$ and $\Sigma = \begin{pmatrix} \Sigma_\leqslant & 0 \\ 0 & \Sigma_> \end{pmatrix}$

$\underset{k^*}{\overset{\longrightarrow}{\quad}} \quad \underset{p - k^*}{\overset{\longrightarrow}{\quad}}$

So that $XX^T = X_\leqslant X_\leqslant^T + X_> X_>^T$

Lemma (spec) ───────

① if $k^* \leqslant \frac{m}{600}$, then with proba $\geqslant 1 - 2e^{-m/84}$

$$\left| \Sigma_\leqslant^{-1/2} (X_\leqslant^T X_\leqslant) \Sigma_\leqslant^{-1/2} - m I_{k^*} \right|_{op} \leqslant \frac{3m}{4}$$

② with proba $\geqslant 1 - 2e^{-m/84}$

$$\left| spec^*(X_> X_>^T) - \sum_{j>k^*} \lambda_j \right| \leqslant \frac{3}{4} \sum_{j>k^*} \lambda_j$$

Comments : the above lemma shows that:

$\rightarrow$ the low dim. part $\frac{1}{m} X_\leqslant^T X_\leqslant$ remains "close" to $\Sigma_\leqslant$

$\rightarrow$ the low signal part $X_> X_>^T$ has a flat spectrum around $\sum_{j>k^*} \lambda_j$

To prove this lemma, we need a concentration bound on quadratic form of Gaussian vectors.

Theorem B.8. (Hanson-Wright)

Let $\varepsilon \sim \mathcal{N}(0, I_d)$ and $\alpha \in \mathbb{R}^d$.

There exists $\mathfrak{z} \sim \text{Exp}(1)$ such that

$$\sum_j \alpha_j \varepsilon_j^2 - \alpha^T \mathbb{1} \leq \sqrt{8\|\alpha\|^2 \mathfrak{z}} \vee (8|\alpha|_\infty \mathfrak{z})$$

$$\leq \frac{|\alpha|_1}{4} + 8|\alpha|_\infty \mathfrak{z}$$

Proof of H.W. (by Chernov method).

Since $-\log(1-x) \leq x + x^2$ for $|x| \leq 1/2$,

we have for any $|\lambda| \leq 1/4$

$$\mathbb{E}\left[\exp(\lambda(\varepsilon_j^2 - 1))\right] = \frac{e^{-\lambda}}{(1-2\lambda)^{1/2}} \overset{(*)}{\leq} e^{2\lambda^2}.$$

So for any $|s| \leq 1/4|\alpha|_\infty$, we have

Chernov

$$\mathbb{P}\left[\sum_j \alpha_j(\varepsilon_j^2-1) > t\right] \leq e^{-st} \prod_j \mathbb{E}\left[e^{s\alpha_j(\varepsilon_j^2-1)}\right]$$

$$\underset{(*)}{\leq} \exp(2\alpha_j^2 s^2)$$

$$\leq \exp(-st + 2\|\alpha\|^2 s^2)$$

Setting $s = \frac{t}{4\|\alpha\|^2} \wedge \frac{1}{4|\alpha|_\infty}$ we get

$$\mathbb{P}\left[\sum_j \alpha_j(\varepsilon_j^2-1) > t\right] \leq \exp\left(-\frac{1}{8}\left(\frac{t^2}{\|\alpha\|^2} \wedge \frac{t}{|\alpha|_\infty}\right)\right)$$

i.e $\sum_j \alpha_j \varepsilon_j^2 - \alpha^T \mathbb{1} \leq \sqrt{8\|\alpha\|^2 \mathfrak{z}} \vee (8|\alpha|_\infty \mathfrak{z})$.

Furthermore

$2ab \leq a^2 + b^2$

$$\sqrt{8\|\alpha\|^2 \mathfrak{z}} \leq 2\sqrt{\frac{|\alpha|_1}{4} \cdot |\alpha|_\infty 8\mathfrak{z}} \leq \frac{|\alpha|_1}{4} + 8|\alpha|_\infty \mathfrak{z}$$ ☐

To control the operator norm

$$|A|_{op} = \max_{\|u\|=1} \|Au\|$$

for A symmetric $= \max_{\|u\|=1} |\langle Au, u \rangle|$

We use the following discretization lemma

**Lemma 12.11-12:** For any $\delta > 0$, there exists $\mathcal{N}_\delta \subset \{u \in \mathbb{R}^d : \|u\| = 1\}$ with cardinality

$$|\mathcal{N}_\delta| \leq \left(1 + \frac{2}{\delta}\right)^d$$

such that, for any $A \in \mathbb{R}^{d \times d}$ symmetric

$$|A|_{op} \leq \frac{1}{1-2\delta} \max_{u \in \mathcal{N}_\delta} |\langle Au, u \rangle|$$

**Proof:** i) _construction of $\mathcal{N}_\delta$_ : Let us set $\mathcal{S}_{d-1} = \{u \in \mathbb{R}^d : \|u\| = 1\}$, and pick iteratively

$$u_1 \in \mathcal{S}_{d-1} \; ; \; u_2 \in \mathcal{S}_{d-1} \setminus B(u_1, \delta) \; ; \; \dots$$

$$u_k \in \mathcal{S}_{d-1} \setminus \bigcup_{j \leq k-1} B(u_j, \delta) \; ; \; \text{until impossible.}$$

By construction, we have

(a) $\|u\| = 1$ and $\|u - u'\| > \delta \quad \forall u \neq u' \in \mathcal{N}_\delta$

(b) $\forall v \in \mathcal{S}_{d-1}, \exists u \in \mathcal{N}_\delta \text{ s.t. } \|u - v\| \leq \delta$

• From (a), we get that:

$$\bigsqcup_{u \in \mathcal{N}_\delta} B(u, \delta/2) \subset B(0, 1 + \delta/2)$$

so comparing the volumes

$$|\mathcal{N}_\delta| \cdot \left(\frac{\delta}{2}\right)^d V_d(1) \leq \left(1 + \frac{\delta}{2}\right)^d V_d(1)$$

$$\Rightarrow |\mathcal{N}_\delta| \leq \left(1 + \frac{2}{\delta}\right)^d$$

ii) _approximation_ : from (b) we have

$$|A|_{op} = |\langle Au^*, u^* \rangle| \qquad \text{with } u^* \in \mathcal{S}_{d-1}$$

with $u \in \mathcal{N}_\delta$
$\|u - u^*\| \leq \delta$

$$= |\langle Au, u \rangle + \langle A(u^* - u), u \rangle + \langle Au^*, u^* - u \rangle|$$

$$\leq |\langle Au, u \rangle| + |A|_{op} \delta + |A|_{op} \delta. \qquad \square$$

---

We are ready to prove Lemma (Spec)

**Proof of Lemma (Spec):**

① Let us set $Z := X_{\leq} \Sigma_{\leq}^{-1/2} \in \mathbb{R}^{m \times k^\circ}$.

we have

$\cdot \; Z_{ij} \overset{iid}{\sim} \mathcal{N}(0,1)$

$\cdot \; Z^T Z = \Sigma_{\leq}^{-1/2} X_{\leq}^T X_{\leq} \Sigma_{\leq}^{-1/2}$

$$\left|\Sigma_{\leq}^{-1/2} X_{\leq}^T X_{\leq} \Sigma_{\leq}^{-1/2} - m I_{k^*}\right|_{op} = \left|Z^T Z - m I_k\right|_{op}$$

$$\leq 2 \max_{u \in \mathcal{N}_{1/4} \subset S_{k-1}} \left|\langle (Z^T Z - m I_{k^*}) u, u\rangle\right|$$

$$= 2 \max_{u \in \mathcal{N}_{1/4}} \left|\|Zu\|^2 - m\right|$$

For any $u \in S_{k-1}$ we have

$$[Zu]_j = \sum_k z_{jk} u_k \overset{iid}{\sim} \mathcal{N}(0,1)$$

So $Zu \sim \mathcal{N}(0, I_m)$ and by Hanson-Wright

$\exists \; \zeta_u, \zeta_u' \sim Exp(1)$ such that

$$\left|Z^T Z - m I_{k^*}\right|_{op} \leq 2 \max_{u \in \mathcal{N}_{1/4}} \left(\frac{m}{4} + 8(\zeta_u \vee \zeta_u')\right)$$

$\alpha_1 = \cdots = \alpha_m = 1$

Furthermore,

union bound

$$\mathbb{P}\left[8 \max_{u \in \mathcal{N}_{1/4}} (\zeta_u \vee \zeta_u') \geq \frac{m}{8}\right] \overset{\downarrow}{\leq} 2|\mathcal{N}_{1/4}| \; \mathbb{P}\left[\zeta \geq \frac{m}{8 \cdot 8}\right]$$

$$\leq 9^{k^*}$$

$$\leq e^{-m/64}$$

for $k^* \leq \frac{m}{600} \longrightarrow \leq 2 \exp\left(-\frac{m}{84}\right)$

which proves ①

② We set now $Z := X_> \Sigma_>^{-1/2} \in \mathbb{R}^{m \times (p-k^*)}$

With the same reasoning as before, we have

$$\left|X_> X_>^T - T_n(\Sigma_>) I_m\right|_{op} = \left|Z \Sigma_> Z^T - T_n(\Sigma_>) I_m\right|_{op}$$

$$\leq 2 \max_{u \in \mathcal{N}_{1/4}} \left(\frac{T_n(\Sigma_>)}{4} + 8 \lambda_{k^*+1} (\zeta_u \vee \zeta_u')\right)$$

$$\subset S_{m-1}$$

Furthermore,

$$\mathbb{P}\left[8 \lambda_{k^*+1} \max_{u \in \mathcal{N}_{1/4}} \zeta_u \vee \zeta_u' \geq \frac{T_n(\Sigma_>)}{8}\right]$$

def. of $k^*$

$$\leq \mathbb{P}\left[\max_{u \in \mathcal{N}_{1/4}} \zeta_u \vee \zeta_u' \geq \frac{m}{64 \rho}\right]$$

union bound

$$\leq 2 \cdot 9^m \cdot \exp\left(-\frac{m}{64 \rho}\right) \leq 2 \exp\left(-\frac{m}{84}\right)$$

$$\frac{1}{\rho} = 142$$

This proves ②.

Remark: we have proved that

$$X^T X \asymp m \Sigma_{\leq} + \lambda^* \hat{I}_{m-k}$$

random projector

where $\lambda^* = \sum_{j > k^*} \lambda_j$

## b) Bounding the variance $V_X$:

We are ready to prove the main result.

**Theorem** (Bartlett et. al. 2019) —————

- $k^* = \min\left\{ k : m \lambda_{k+1} \leq \beta \sum_{j > k} \lambda_j \right\}$

  with $\beta = \frac{1}{142}$

- If $k^* \leq m/600$, then

$$V_X \leq 24 \, \sigma^2 \left( \underbrace{\frac{k^*}{m}}_{\text{low. rank}} + \underbrace{\frac{m \sum_{j > k^*} \lambda_j^2}{\left( \sum_{j > k^*} \lambda_j \right)^2}}_{\color{red}{\text{H.D. regularization}}} \right)$$

with $\mathbb{P}_X \geq 1 - 5 \exp(-m/84)$.

**Proof:**

We have $V_X = \sigma^2 \langle (X^T X)^+, \Sigma \rangle$

Next lemma disentangles the low rank part and the H.D. part.

**Lemma:** —————

$$\langle (X^T X)^+, \Sigma \rangle \leq \langle (X_{\leq}^T X_{\leq})^+, \Sigma_{\leq} \rangle + \langle (X_{>}^T X_{>})^+, \Sigma_{>} \rangle$$

**Proof:** We have $(X^T X)^+ = X^T (X X^T)^{2+} X$

and $X X^T = X_{\leq} X_{\leq}^T + X_{>} X_{>}^T$. So

$$\langle (X^T X)^+, \Sigma \rangle = \langle (X_{\leq} X_{\leq}^T + X_{>} X_{>}^T)^{2+}, \underbrace{X \Sigma X^T}_{\color{green}{= X_{\leq} \Sigma_{\leq} X_{\leq}^T + X_{>} \Sigma_{>} X_{>}^T}} \rangle$$

we have:

$X X^T \succeq X_{\leq} X_{\leq}^T$ and $X_{>} X_{>}^T$ so

$$\begin{cases} (X_{\leq} X_{\leq}^T)^{2+} \succeq (X X^T)^{2+} & \text{on range } (X_{\leq}) \\ (X_{>} X_{>}^T)^{2+} \succeq (X X^T)^{2+} & \text{on range } (X_{>}) \end{cases}$$

Hence

$$\langle (X^T X)^+, \Sigma \rangle \leq \langle (X_{\leq} X_{\leq}^T)^{2+}, X_{\leq} \Sigma_{\leq} X_{\leq}^T \rangle + \langle (X_{>} X_{>}^T)^{2+}, X_{>} \Sigma_{>} X_{>}^T \rangle$$

$$= \langle (X_{\leq}^T X_{\leq})^+, \Sigma_{\leq} \rangle + \langle (X_{>}^T X_{>})^+, \Sigma_{>} \rangle \qquad \square$$

**Low rank term:** We have

$$\langle (X_{\leq}^T X_{\leq})^+, \Sigma_{\leq} \rangle_F \stackrel{k^* \leq m}{=} \langle (\Sigma_{\leq}^{-1/2} X_{\leq}^T X_{\leq} \Sigma_{\leq}^{-1/2})^{-1}, I_{k^*} \rangle_F$$

(full rank)   full rank ↓

• From Lemma (Spec) ① we have

$$\text{Spec}\left( \Sigma_{\leq}^{-1/2} X_{\leq}^T X_{\leq} \Sigma_{\leq}^{-1/2} \right) \geq \frac{m}{4}$$

with proba $\geq 1 - 2\exp(-m/84)$

• On this event $\text{Spec}\left( (\Sigma_{\leq}^{-1/2} X_{\leq}^T X_{\leq} \Sigma_{\leq}^{-1/2})^{-1} \right) \leq \frac{4}{m}$

So $\langle (X_{\leq}^T X_{\leq})^+, \Sigma_{\leq} \rangle_F \leq k^* \cdot \frac{4}{m}$

with proba $\geq 1 - 2\exp(-\frac{m}{84})$.

**H.D. term:** Again from Lemma (Spec) ②

with proba $\geq 1 - 2\exp(-\frac{m}{84})$, we have

$$\text{Sp}(X_> X_>^T) \subset \left[ \frac{1}{4} \text{Tr}(\Sigma_>), 2\text{Tr}(\Sigma_>) \right].$$

Hence, on this event:

$$\langle (X_>^T X_>)^+, \Sigma_> \rangle_F = \langle (X_> X_>^T)^{2+}, X_> \Sigma_> X_>^T \rangle_F$$

$$\leq \left( \frac{4}{\text{Tr}(\Sigma_>)} \right)^2 \langle I_m, X_> \Sigma_> X_>^T \rangle_F$$

It remains to evaluate the trace

$$\langle I_m, X_> \Sigma_> X_>^T \rangle_F = \langle I_m, Z \Sigma_>^2 Z \rangle_F$$

$$= \sum_{i=1}^{m} (Z \Sigma_>^2 Z^T)_{ii} = \sum_{i=1}^{m} \sum_{j=1}^{p-k^*} Z_{ij}^2 \, \lambda_{k^*+j}^2$$

H.W.

$$\leq \left( 1 + \frac{1}{4} \right) m \sum_{j > k^*} \lambda_j^2 + 8 \lambda_{k^*+1}^2 \Bigg\}$$

So with probability $\geq 1 - e^{-m/32}$, we have

$$\langle I_m, X_> \Sigma_> X_>^T \rangle_F \leq \frac{5m}{4} \sum_{j > k^*} \lambda_j^2 + \frac{m}{4} \underline{\lambda_{k^*+1}^2}$$

$$\leq \sum_{j > k^*} \lambda_j^2$$

$$\leq \frac{6}{4} m \sum_{j > k^*} \lambda_j^2$$

Hence, with proba $\geq 1 - 3\exp(-\frac{m}{84})$, we have

$$\langle (X_>^T X_>)^+, \Sigma_> \rangle_F \leq 24 \, m \, \frac{\sum_{j > k^*} \lambda_j^2}{(\sum_{j > k^*} \lambda_j)^2}$$

□

# Remarks:

→ a matching lower bound shows that

$$V_x \ \simeq \ \sigma^2 \left( \frac{k^*}{m} + \frac{m \sum_{j > k^*} \lambda_j^2}{\left( \sum_{j > k^*} d_j \right)^2} \right)$$

with $k^* = \min \{ k : m \lambda_{k+1} \leq \rho \sum_{j > k} d_j \}$, where $\rho = 1/142$

→ for $\Sigma = I_p$ and $p \gg m$, we have $k^* = 0$ and

$$V_x \ \simeq \ \sigma^2 \cdot \frac{mp}{p^2} = \sigma \frac{m}{p} \quad \leadsto \text{We recover our guess}$$

→ for $\Sigma = I_k + \alpha I_{\geq k}$ with $\alpha \ll 1$ and $k \ll m \ll p$,
we have $k^* = k$ and

$$V_x \ \simeq \ \sigma^2 \left( \frac{k}{m} + \frac{m \alpha^2 (p - k)}{\alpha^2 (p - k)^2} \right) \overset{k \ll p}{\simeq} \sigma^2 \left( \frac{k}{m} + \frac{m}{p} \right)$$

$\leadsto$ We recover our guess.

What about the bias $B_X$ ?

It has been shown that

$$B_X \simeq \|\beta^*_\leq\|^2_{\Sigma_\leq^{-1}} \times \left(\frac{\sum\limits_{j>k^*} \lambda_j}{m}\right)^2 + \|\beta^*_>\|^2_{\Sigma_>}$$

$\underbrace{\qquad\qquad\qquad\qquad\qquad\qquad}$ small bias for the low dimensional part

$\underbrace{\qquad\qquad}$ high bias for the high-dimensional part.

- Remark :

  The bias for the low dimensional part $\approx$ bias of ridge regression for $\beta^*_\leq$ with

  $$\lambda := \frac{1}{m} \sum_{j>k^*} \lambda_j$$

- Take home message : in the interpolation regime, for the GD least square estimator $\hat{\beta}^{LS} = X^+ Y$

  $\rightarrow$ the low dimensional part $\beta^*_\leq$ can be well estimated.

  $\rightarrow$ the high-dimensional part $\beta^*_>$ is almost estimated by 0 $\nearrow$ large bias $\searrow$ small variance

# 3. Benign Overparametrisation

**Recap**: so far we have seen that

→ GD induces an implicit regularization in the interpolation regime.

→ For linear regression, the interpolating estimator selected by G.D. on the least-square problem, is not suffering from an overly high variance in very high dimension, thanks to the implicit regularization of GD. But it may suffer from a large bias.

→ What does it tell us on "overparametrisation"

(having much more parameters than the sample size)

In linear regression $y_i = \langle x_i, \beta \rangle + \varepsilon_i$

$p$ is both → input dimension ERP
→ number of parameters

⤳ we must look at a more complex model to disentangle over-parametrisation from high input dimension.

___

Is over-parametrisation new in statistics?

Certainly not! Think to classical non-parametric regression

$$y_i = f^*(x_i) + \varepsilon_i, \quad i = 1, ..., n$$

where $f^*$ is some (say) Sobolev function

what is new then? In non-parametric statistics, we use some explicit regularisation

For example, cubic splines are solution to the regularized Least-Square problem

$$\hat{f}^{spline} \in \underset{\int (f'')^2 < +\infty}{argmin} \left\{ \sum_{i=1}^{m} (Y_i - f(x_i))^2 + \lambda \int (f'')^2 \right\}$$

<span style="color:red">regularization by mean curvature</span>

by G.D.

We will illustrate this on a simple example.

_____

What is new in Neural Network practice is that some people use very complex over-parametrized models without any explicit regularization.

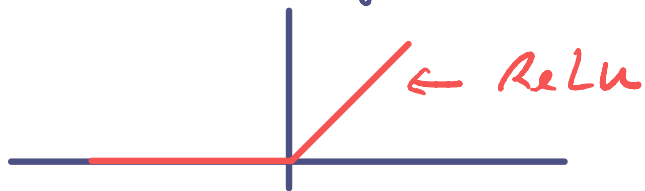• Why <u>does-it work</u>? Thanks to implicit regularization

# Over-parametrized 1-hidden layer Neural Network

$f_{\beta,\omega} : \mathbb{R}^d \longrightarrow \mathbb{R}$, defined by

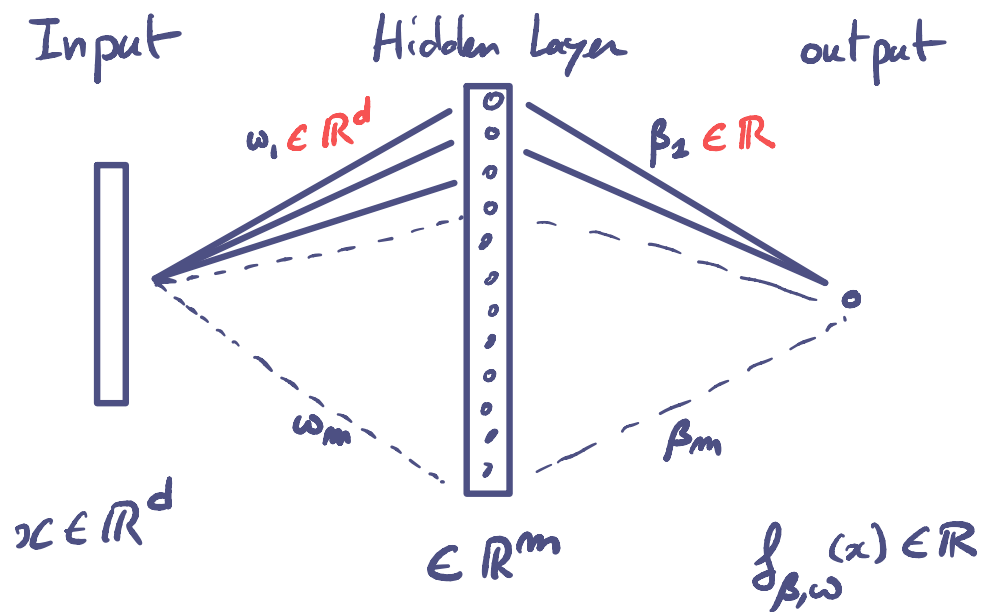$$f_{\beta,\omega}(x) := \frac{1}{m} \sum_{j=1}^{m} \beta_j \varphi(\langle \omega_j, x \rangle)$$

with

- $\varphi$ = activation function e.g. ReLu


$\leftarrow$ ReLu

- $m$ = number of hidden neurons

- $\theta = (\beta_j, \omega_j)_{j=1\ldots m} \in \mathbb{R}^{m(d+1)}$

- over-parametrization: $m \longrightarrow \infty$
(with fixed input dimension)

Input     Hidden Layer     output



$\omega_1 \in \mathbb{R}^d$     $\beta_1 \in \mathbb{R}$

$\omega_m$     $\beta_m$

$x \in \mathbb{R}^d$     $\in \mathbb{R}^m$     $f_{\beta,\omega}(x) \in \mathbb{R}$

learning: let $(X_i, Y_i)_{i=1,\ldots,m} \in (\mathbb{R}^d \times \{-1, 1\})^n$ and $l$ be a convex loss function. $(\hat{\beta}, \hat{\omega})$ are learnt from GD on the empirical loss

$$\mathcal{L}(\beta, \omega) = \frac{1}{n} \sum_{i=1}^{n} l(-Y_i f_{\beta,\omega}(x_i))$$

⚠️ Even if $l$ is convex, the function $\mathcal{L}$ is <u>not</u> convex

- <u>over-parametrisation</u>: $m \to +\infty$

the limit $m \to \infty$, corresponds to the parametrisation

$$f_\mu(x) := \int_{\beta, \omega} \underbrace{\beta}_{=\theta} \varphi(\langle \omega, x \rangle) \, d\mu(\beta, \omega)$$

with $\mu \in \mathcal{P}(\mathbb{R}^{d+1})$

- Setting $\phi(\theta, x) = \beta \varphi(\langle \omega, x \rangle)$

we have

$$f_\mu(x) = \langle \phi(\cdot, x), \mu \rangle \quad \text{linear!}$$

where $\langle g, \mu \rangle := \int g(\theta) \, d\mu(\theta)$

$$\ddot\smile \quad \mathcal{L}(\mu) = \frac{1}{m} \sum_{i=1}^{m} \ell(-y_i f_\mu(x_i))$$

is convex !

$\leadsto$ here, over-parametrisation <u>helps</u> for the optimisation landscape

$\leadsto$ is there a statistical price for it?

<u>Example</u>: $\varphi = \text{Relu}$, $\ell = $ logistic loss

- <u>reparametrisation</u>:

$$\phi(\lambda\theta, x) = \lambda\beta \, \varphi(\langle \lambda\omega, x \rangle)$$
$$= \lambda^2 \beta \varphi(\langle \omega, x \rangle)) = \lambda^2 \phi(\theta, x)$$

so with $\theta = \lambda u$, $\lambda > 0$ and $u \in S_d$

$$f_\mu(x) = \int_{\substack{\lambda > 0 \\ u \in S_d}} \phi(\lambda u, x) \, d\mu(\lambda u)$$

$$= \int_{u \in S_d} \phi(u, x) \underbrace{\int_{\lambda > 0} \lambda^2 \, d\mu(\lambda u)}_{=: d\pi_\mu(u)}$$

$$= \langle \phi(\cdot, x), \pi_\mu \rangle \qquad \in \mathcal{M}_+(S_d)$$

**Overfitting:** We can represent any Lipschitz function $f$ by a $f_\mu$. So if no point $x_i$ has two different labels $y_i$ in the learning data, then $\exists \mu$ such that

$$y_i = \text{sign}\left(f_\mu(x_i)\right) \quad \substack{i=1\ldots m \\ \textcolor{red}{(perfect\ fit)}}$$

**Max-margin:** assume that at some stage $t$ of G.D., $f_{\hat\mu^t}$ perfectly fit the data i.e. for $i = 1, \ldots, m$

$$y_i \, f_{\hat\mu^t}(x_i) = \langle y_i \phi(x_i, \cdot), \underbrace{\Pi_{\hat\mu^t}}_{\textcolor{red}{=: \hat\Pi^t}} \rangle > 0.$$

Then, since $\ell(-z) = \log(1 + e^{-z})$ decreases we can always decrease $\mathcal{L}(\hat\mu^t)$ by simply increasing the mass of $\hat\Pi^t$

$\leadsto$ consequence $|\hat\Pi^t| \to +\infty$

$\leadsto$ since $\ell(-z) \approx e^{-z}$ as $z \to +\infty$

$$\mathcal{L}(\hat\mu^t) \approx \frac{1}{m} \sum_i \exp\left(-|\hat\Pi^t| \, \langle y_i \phi(x_i, \cdot), \frac{\hat\Pi^t}{|\hat\Pi^t|} \rangle\right)$$

$$\approx \frac{N_{min}}{m} \exp\left(-\underbrace{\textcolor{green}{|\hat\Pi^t|}}_{\textcolor{green}{\to +\infty}} \underbrace{\min_i \langle y_i \phi(x_i, \cdot), \frac{\hat\Pi^t}{|\hat\Pi^t|} \rangle}_{\textcolor{red}{margin\ of\ f_{\hat\Pi^t/|\hat\Pi^t|}}}\right)$$

**Guess:** for $t$ large

$$\frac{\hat\Pi^t}{|\hat\Pi^t|} \approx \hat\nu \in \underset{\nu \in \mathcal{P}(S_d)}{\text{argmax}} \ \min_i \ y_i \underbrace{f_\nu(x_i)}_{\textcolor{red}{\langle \phi(x_i, \cdot), \nu \rangle}}$$

**Theorem** (informal)

Under some mild conditions, the above guess holds true

It is then possible to prove that the classifier $\hat{h}(x) := \text{sign}\left(f_{\hat{V}}(x)\right)$ has some nice statistical properties. For example, it is able to adapt to low-dimensional structures ...
(see Chizat and Bach 2020)

⤳ in this case

implicit regularisation of G.D.

+ overparametrisation

⇓

1) Nice optimisation landscape

2) Nice statistical behavior

°⌣°

**Take home messages:**

→ in the interpolation regime, G.D. has a regularizing effect by selecting some specific interpolating solutions

→ when the input space is very large, overfitting only occurs on domains which are rarely sampled. So overfitting does not harm prediction risk

→ over-parametrisation can be harmless, and even beneficial.

⚠ All these phenomena have to be better understood