

# Implicit regularisation, Benign overfitting, and Over-parametrisation



## Lecture 5

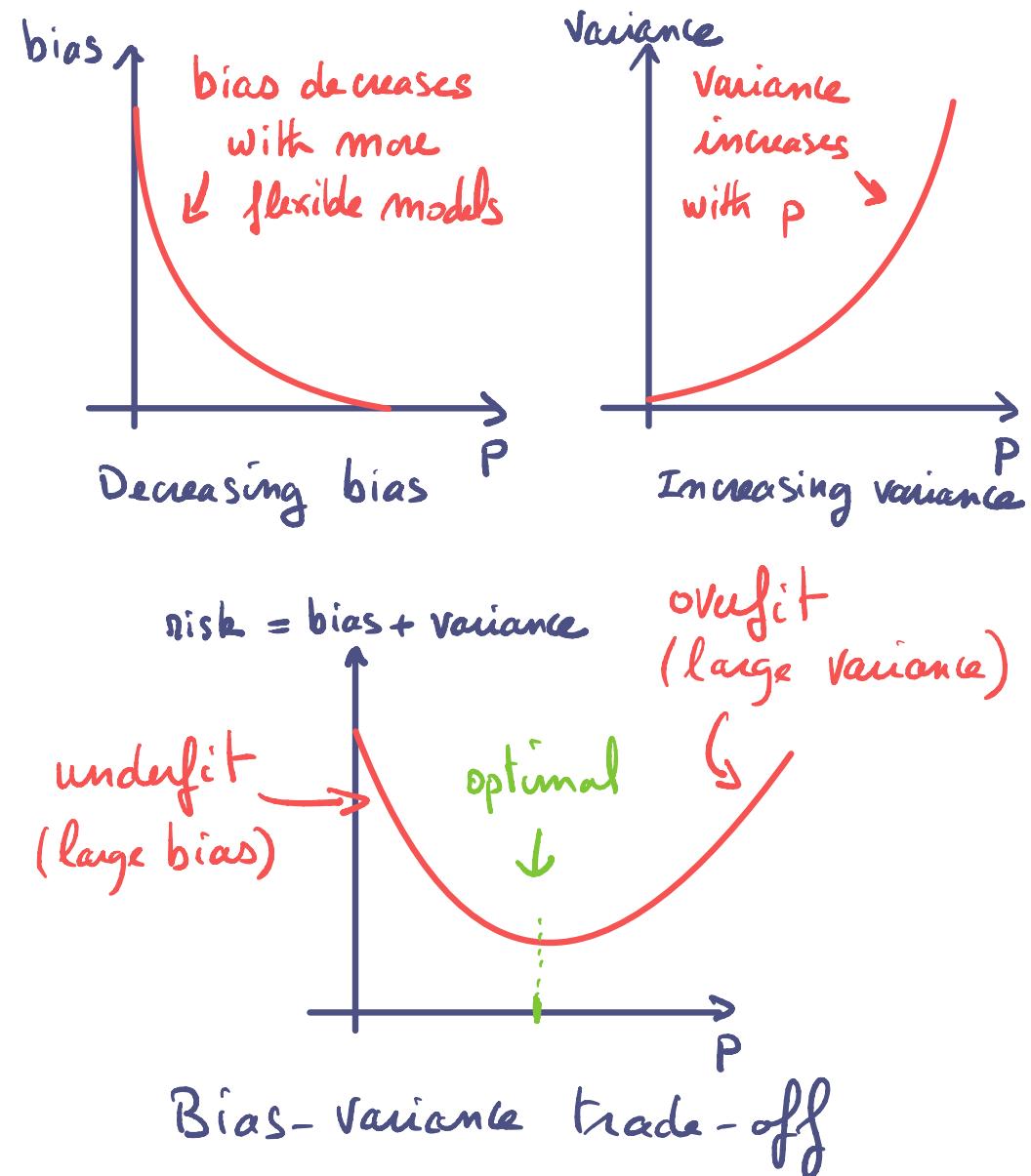
## Main references for today:

- M. Belkin, D. Hsu, S. Ma, S. Nandakumar. "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)

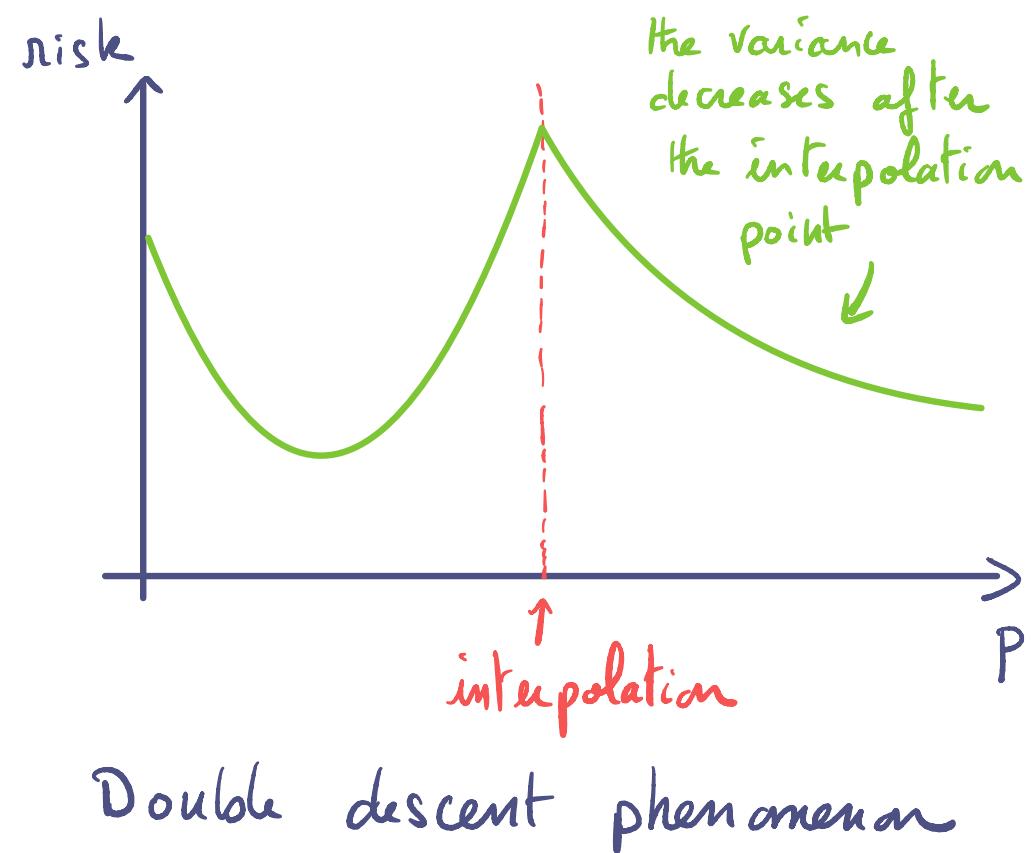
## Double descent phenomenon

Text book theory: risk = bias + variance

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



## Empirical observations (neural networks)



Today:

- ① Implicit regularisation (a.k.a. implicit bias) of gradient descent
- ② The magic of high-dimensional input: benign overfitting in high-dimensional (linear) regression
- ③ Benign overparametrisation

# ① Implicit regularisation of G.D.

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

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

- High-dimension:  $x_i \in \mathbb{R}^p$  with  $p > m$

- G.D. on least-squares:

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

$$\hat{\beta}^{t+1} \stackrel{\text{GD}}{=} \hat{\beta}^t - 2\gamma X^T(X\hat{\beta}^t - Y)$$

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

$$\begin{aligned} t \rightarrow +\infty \\ \rightarrow 2\gamma \underbrace{< 1}_{\text{cond. on } X_{\text{op}}} \underbrace{X^T X}_{\text{invertible}}^{-1} (I - (I - 2\gamma X^T X))^+ \cdot 2\gamma X^T Y \end{aligned}$$

$$= X^+ Y$$

↑ Moore-Penrose pseudo inverse

$$\text{So } \hat{\beta} = X^+ Y = \underset{\substack{\uparrow \\ Y = X\beta}}{\underset{\substack{\uparrow \\ \text{if } \text{rank}(X) = m \\ (\text{interpolation})}}{\underset{\substack{\uparrow \\ \text{regularisation} \\ \text{of G.D.}}}{\arg\min \| \beta \|^2}}$$

Prediction: for  $x \in \mathbb{R}^p$

$$\langle \hat{\beta}, x \rangle = \langle X^+ Y, x \rangle$$

$$\begin{aligned} Y = X\beta^* + \varepsilon \rightarrow & \langle X^+ X \beta^*, x \rangle + \langle X^+ \varepsilon, x \rangle \\ & = P_{X^T} \quad (\text{projection on } \text{range}(X^T)) \end{aligned}$$

$$= \underbrace{\langle \beta^*, x \rangle}_{\text{target}} + \underbrace{\langle (I - P_{X^T}) \beta^*, x \rangle}_{\text{bias}} + \underbrace{\langle \varepsilon, (X^T)^+ x \rangle}_{N(0, \sigma^2 (X^T X)^+)} \quad \text{cond. on } x$$

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

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

- $\text{rank}(X) \leq m < p$ :

$\|P_{X^T} x\| \ll \|x\|$  if  $x \sim N(0, \Sigma)$   
with  $\Sigma \succeq I_p$ .

## ② Benign overfitting with high-dimensional input

a) (Quasi-) isotropic input:  $\Sigma = I_p$  ( $\text{or } \Sigma \asymp I_p$ )

Variance analysis:  $X \stackrel{\text{SVD}}{=} \sum_k \hat{\sigma}_k \hat{u}_k \hat{v}_k^T$

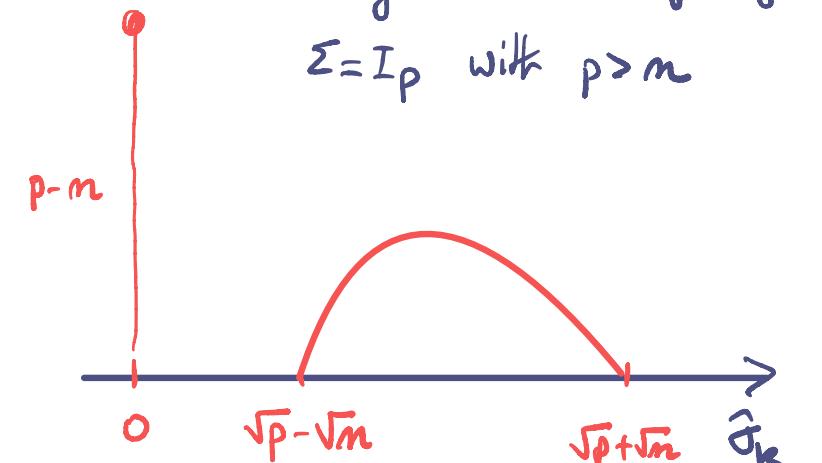
$$\text{Var}_{\Sigma} \langle \hat{\beta}, x \rangle = \sigma^2 x^T (X^T X)^+ x$$

$$= \sigma^2 x^T \left( \sum_{k=1}^m \frac{1}{\hat{\sigma}_k^2} \hat{v}_k \hat{v}_k^T \right) x$$

$$\begin{aligned} \hat{\sigma}_k^2 &\approx p \\ \text{w.h. } P_X &\sim \frac{\sigma^2}{p} x^T P_{X^T} x \\ \text{when } p \gg m &= \|P_{X^T} x\|^2 \quad \text{if } x \sim N(0, I_p) \\ &\approx m \\ &\text{w.h. } P_{\mathcal{C}} \end{aligned}$$

$$\approx \boxed{\sigma^2 \frac{m}{p}} \quad \text{with high-probability}$$

$\rightarrow 0$  when  $p \gg m$  !



singular values of  $x$  for  
 $\Sigma = I_p$  with  $p > m$

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

Average prediction error:  $\hat{\beta} = P_{X^\perp} \beta^* + X^+ \varepsilon$   $P_{\text{ker}(X)}$  (projection on  $\text{ker}(X)$ )

$$\begin{aligned}
 R_x &= \mathbb{E}_{x, \varepsilon} [\langle x, \hat{\beta} - \beta^* \rangle^2] \stackrel{\downarrow}{=} \mathbb{E}_{x, \varepsilon} \left[ \left( \langle (I - P_{X^\perp}) \beta^*, x \rangle + \langle X^+ \varepsilon, x \rangle \right)^2 \right] \\
 &= \mathbb{E}_x \left[ \beta^{*T} P_{\text{ker}(X)} x x^T P_{\text{ker}(X)} \beta^* \right] + \sigma^2 \mathbb{E}_x \left[ \underbrace{x^T x^+ (X^+)^T x}_{\text{cross-term}} \right] + O \\
 &= \underbrace{\beta^{*T} P_{\text{ker}(X)} \sum P_{\text{ker}(X)} \beta^*}_{B_x} + \underbrace{\sigma^2 \langle \sum, (X^T X)^+ \rangle_F}_{V_x} = \langle x x^T, x^+ (X^+)^T \rangle_F
 \end{aligned}$$

$\cdot \sum \equiv I_p$ :

$$\cdot B_x = \|P_{\text{ker}(X)} \beta^*\|^2 = \|\beta^*\|^2 - \|P_{X^\perp} \beta^*\|^2 \underset{\text{w.h. } P_X}{\approx} \|\beta^*\|^2 \left(1 - \frac{m}{p}\right)$$

The larger  $p$ ,  
the larger the  
bias  $B_x$

$$\text{Range}(X^T) \xrightarrow{\text{unit. dim}(n)} \mathbb{R}^p$$

$$\begin{aligned}
 \cdot V_x &= \sigma^2 \text{Tr}(X^T X)^+ = \sigma^2 \sum_{k=1}^m \frac{1}{\lambda_k^2} \quad \text{so} \quad \|P_{X^\perp} \beta^*\|^2 \simeq \frac{m}{p} \|\beta^*\|^2 \\
 &\simeq \sigma^2 \frac{m}{p}
 \end{aligned}$$

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

$\Sigma \asymp I_p$ : More complex formula, based on Stieltjes transform of Dauchenko-Pastur distribution, but exhibits the same behavior.

Recap at this stage: ( $\Sigma \asymp I_p$ )

$$\rightarrow \text{if } x \in \text{range}(X^T): \quad \text{var}_{\Sigma} \langle \hat{\beta}, x \rangle = \sigma^2 \frac{\|x\|^2}{p} \approx \sigma^2 \leftarrow \begin{matrix} \text{strong overfit} \\ \text{on range}(X^T) \end{matrix}$$

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

$$\text{var}_{\Sigma} \langle \hat{\beta}, x \rangle \approx \sigma^2 \frac{\|P_{X^T} x\|^2}{p} \quad \text{small}$$

$\rightarrow$  in words: we overfit on a small space =  $\text{Span}\{x_1, \dots, x_n\}$  spanned by the learning inputs, but we underfit everywhere else

$\Rightarrow$  large bias and small variance for  $p \gg n$ !

b) Anisotropic case:

Prototypical example:  $\Sigma = I_k + \rho I_{p-k} = \begin{bmatrix} I_k & 0 \\ 0 & \rho I_{p-k} \end{bmatrix}$  with  $k \ll m \ll p$ ,  $\rho \ll 1$

Then  $\frac{1}{m} X^T X = \frac{1}{m} \sum_{i=1}^m x_i x_i^T \approx I_k + \rho \frac{I_{p-k}}{m-k} \hat{I}_{m-k}$   
low dimensional  $\Rightarrow$  good estimation  $\approx$  projection of rank  $m-k$   
 $\text{Trace} = k + \rho(p-k)$

$k \ll m \ll p$   
So  $(X^T X)^+ \approx \frac{1}{m} I_k + \frac{1}{\rho p} \hat{I}_{m-k}$  and

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

$$\approx \sigma^2 \left( \frac{k}{m} + \frac{m}{\rho p} \right)$$

small since

small rank  $k$

small due

to high-dimensional input

General case :  $\Sigma = \sum_{j=1}^p d_j \underbrace{v_j v_j^T}_{\downarrow}$

Theorem (informal)

. Define  $k^* = \min\{k : m \lambda_{k+1} \leq \frac{c}{\sqrt{d_j}} \sum_{j>k} d_j\}$

. Then

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

low dim.

H.D.  
regularisation

Example:  $\Sigma = I_{k^*} + \rho I_{m-k^*}$

$$V_X \asymp \sigma^2 \left( \frac{k^*}{m} + \frac{m}{p-k^*} \right)$$

$$\simeq \sigma^2 \left( \frac{k^*}{m} + \frac{m}{p} \right) \quad \text{if } k^* \ll m$$

Sketch of proof:

$$z_{ij} \stackrel{\text{iid}}{\sim} N(0, 1)$$

$$\text{Set } z_j := \frac{1}{\sqrt{d_j}} \times v_j = \frac{1}{\sqrt{d_j}} Z \sum^{1/2} v_j$$

$$= Z v_j \stackrel{\text{iid}}{\sim} N(0, I_m)$$

$$\underbrace{XX^T}_{m \times m} = \sum_{j=1}^p \underbrace{v_j v_j^T}_{\text{part related to } V_x} X^T = \sum_{j=1}^p d_j z_j z_j^T$$

$$= S_{\leq k^*} + S_{>k^*}$$

part related to  $V_x$

part related to  $V_x^\perp$

Key Lemma:

$$\text{Sp}(S_{>k^*}) \asymp \tilde{\lambda}^* := \sum_{j>k^*} d_j$$

→ delocalisation of small energy spectrum

Proof: For  $\|\omega\|=1$

$$\omega^T S_{>k^*} \omega = \sum_{j>k^*} d_j (\omega^T z_j)^2$$

$\stackrel{\text{iid}}{\sim} N(0, 1)$

So

$$\max_{\|\boldsymbol{\sigma}\|=1} \left| \boldsymbol{\sigma}^T S_{>k^*} \boldsymbol{\sigma} - \sum_{j>k^*} \lambda_j \right| = \max_{\|\boldsymbol{\sigma}\|=1} \left| \sum_{j>k^*} \lambda_j ((\boldsymbol{\sigma}^T \mathbf{z}_j)^2 - 1) \right|$$

$\stackrel{iid N(0,1)}{\sim}$

Hanson-Wright  
+  $\varepsilon$ -net discretisation  $\rightarrow$

$$\begin{aligned} &\leq m \|S_{>k^*}\|_{op} + \sqrt{m \sum_{j>k^*} \lambda_j^2} \\ &\leq \left(1 + \frac{1}{\alpha}\right) m \underbrace{\|S_{>k^*}\|_{op}}_{= \lambda_{k^*+1}} + \alpha \underbrace{\sum_{j>k^*} \lambda_j}_{= \lambda^{**}} \\ &\leq c \left(1 + \frac{1}{\alpha}\right) \lambda^* \end{aligned}$$

from  $k^*$

so

$$(1 - c') \lambda^* \leq \text{sp}(S_{>k^*}) \leq (1 + c') \lambda^*$$

where  $c' \asymp c \left(1 + \frac{1}{\alpha}\right) + \alpha$  can be small if  $c \ll 1$  and  $\alpha = \sqrt{c}$ .

□

For  $S_{\leq k^*}$ : similar arguments show that

$$\text{Sp}\left(\frac{1}{m} S_{\leq k^*}\right) \approx \{\lambda_1, \dots, \lambda_{k^*}\}$$

← low rank estimation part.

Consequence:  $\text{sp}(XX^T) \simeq \{m\lambda_1, \dots, m\lambda_{k^*}, \underbrace{\lambda^{*}, \dots, \lambda^{*}}_{m-k^* \text{ times}}\}$  and

$$\text{sp}(X^T X) \simeq \{m\lambda_1, \dots, m\lambda_{k^*}, \underbrace{\lambda^{*}, \dots, \lambda^{*}}_{m-k^*}, \underbrace{0, \dots, 0}_{p-m}\}$$

so  $\frac{1}{m} X^T X \simeq \underbrace{\sum_{\leq k^*}}_{\text{low rank part}} + \underbrace{\frac{\lambda^*}{m} \sum_{j=k^*+1}^m \hat{U}_j \hat{U}_j^T}_{\text{delocalised spectrum}}$  (compare to prototypical example)

$$\begin{aligned} V_X &:= \sigma^2 \langle \sum, (X^T X)^+ \rangle \simeq \frac{\sigma^2}{m} \langle \sum_{\leq k^*}, \sum_{\leq k^*}^+ \rangle + \frac{\sigma^2}{\lambda^*} \sum_{j=k^*+1}^m \hat{U}_j^T \sum_{>k^*} \hat{U}_j \\ &\simeq \frac{\sigma^2}{m} k^* + \frac{\sigma^2}{\lambda^*} \sum_{l>k^*} \lambda_l \sum_{j=k^*+1}^m \underbrace{\hat{U}_j^T U_l U_l^T \hat{U}_j}_{= (\hat{U}_j^T U_l)^2 \simeq \frac{\lambda_l}{\lambda^*}} \end{aligned}$$

$$\simeq \frac{k^*}{m} \sigma^2 + \frac{m-k^*}{(\lambda^*)^2} \sum_{l \geq k^*+1} \lambda_l^2$$

since  $\sum_l (\hat{U}_j^T U_l)^2 = 1$  and energy  
in direction  $U_l$  proportional  
to  $\lambda_l$

□

### ③ Benign over-parametrisation

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

$p$  is both  
 ↗ input dimension  $\in \mathbb{R}^p$   
 ↗ number of parameters

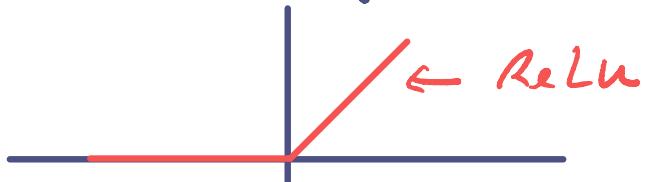
→ we must look at a more complex model to disentangle over-parametrisation from input dimension.

#### 1-Hidden Layer Neural Network

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

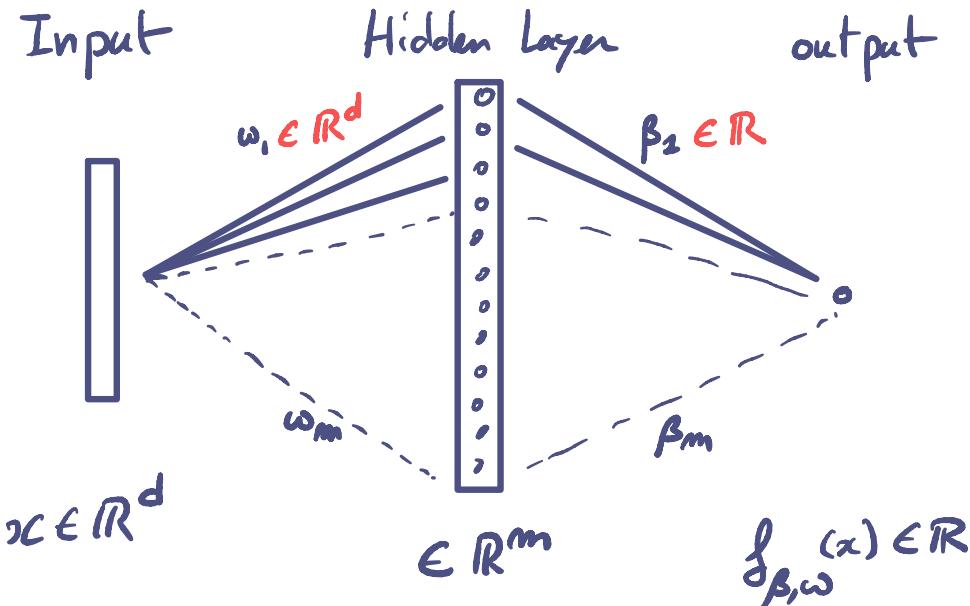
with

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



- $m$  = number of hidden neurons

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



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

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



Even if  $l$  is convex, the function  $\mathcal{L}$  is not convex

• over-parametrisation:  $m \rightarrow +\infty$   
 the limit  $m \rightarrow \infty$ , corresponds to the  
 parametrisation

$$f_\mu(x) := \int_{\beta, \omega} \beta \varphi(\langle \omega, x \rangle) d\mu(\beta, \omega) \stackrel{=}{\Theta}$$

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

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

we have

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

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

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

is convex!

→ here, overparametrisation helps for  
 the optimisation landscape  
 → is there a statistical price for it?

Example:  $\varphi = \text{ReLU}$ ,  $l = \text{logistic loss}$

• reparametrisation:

$$\begin{aligned} \phi(\lambda \theta, x) &= \lambda \beta \varphi(\langle \lambda \omega, x \rangle) \\ &= \lambda^2 \beta \varphi(\langle \omega, x \rangle) = \lambda^2 \phi(\theta, x) \end{aligned}$$

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

$$\begin{aligned} f_\mu(x) &= \int_{\substack{\lambda > 0 \\ u \in S_d}} \phi(\lambda u, x) d\mu(\lambda u) \\ &= \int_{u \in S_d} \phi(u, x) \underbrace{\int_{\lambda > 0} \lambda^2 d\mu(\lambda u)}_{=: d\pi_{\mu}(u)} \\ &= \langle \phi(\cdot, x), \pi_\mu \rangle \quad \in \mathcal{M}_+(S_d) \end{aligned}$$

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

$$y_i = \text{sign}(f_\mu(x_i)) \quad \stackrel{i=1\dots n}{\text{(perfect fit)}}$$

Max-Margin: assume that at some stage  $t$  of G.O.,  $f_{\hat{\mu}^t}$  perfectly fit the data i.e. for  $i=1\dots n$

$$y_i f_{\hat{\mu}^t}(x_i) = \langle y_i \phi(x_i, \cdot), \frac{\pi_{\hat{\mu}^t}}{\|\pi_{\hat{\mu}^t}\|} \rangle > 0.$$

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

→ consequence  $\|\hat{\pi}^t\| \rightarrow +\infty$   
 → since  $l(-z) \approx -e^{-z}$  as  $z \rightarrow +\infty$

$$\begin{aligned} L(\hat{\mu}^t) &\approx \sum_i \exp\left(-\|\hat{\pi}^t\| \langle y_i \phi(x_i, \cdot), \frac{\pi_{\hat{\mu}^t}}{\|\pi_{\hat{\mu}^t}\|} \rangle\right) \\ &\approx \exp\left(-\underbrace{\|\hat{\pi}^t\|}_{\rightarrow +\infty} \min_i \underbrace{\langle y_i \phi(x_i, \cdot), \frac{\pi_{\hat{\mu}^t}}{\|\pi_{\hat{\mu}^t}\|} \rangle}_{\text{margin of } f_{\hat{\mu}^t}/\|\hat{\pi}^t\|}\right) \end{aligned}$$

Guess: for  $t$  large

$$\frac{\hat{\pi}^t}{\|\hat{\pi}^t\|} \approx \hat{v} \in \underset{v \in \mathcal{P}(S_d)}{\operatorname{argmax}} \min_i y_i \langle \phi(x_i, \cdot), v \rangle$$

Theorem (informal)

Under some mild conditions, the above guess holds true

It is then possible to prove that the classifier  $\hat{h}(x) := \text{sign}(\hat{f}_{\hat{\gamma}}(x))$  has some nice statistical properties - For example, it is able to adapt to low-dimensional structures ...

→ 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

## Messages from this course:

- in general, learning in high-dimensional spaces is hopeless
- but feasible, when there exist some low-dimensional structures in the data
- it is possible to adapt to low-dimensional structures, but it comes with computational challenges (convexification or iterative optimisation)
- bias from convexification can be an issue

→ G.D. induces some regularisation beneficial in highly dimensional input spaces and/or highly overparametrised models.

---