# Guarantees : generalization / robustness
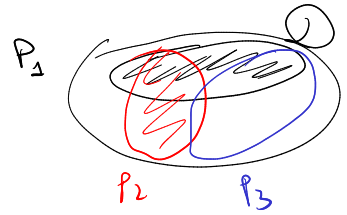
## I. Generalization

### I.A) Ensemble methods

- agglomerate different predictors

  ↓

  average
  all predictors

  or

  more advanced ways
  (use the validation set)

  ↳ simple case: $x \mapsto \sum_k f_k(x) / \#R$
  
  (predictors)

→ different techniques : → SVM
  → random forests
  → deep learning

→ same technique,
  different parameters
  → different kernels
  → deep: different architectures
    { hyper-parameters

→ same technique,
  same parameters → different parts
  of the dataset

$P_1$

$P_2$  $P_3$

[Distilling the knowledge in a Neural Network] "Teacher-student" / "Distillation"

- hard task
- "small" network : training → poor performances
  several " s : → " "

set of predictors

ensemble method (average) → good performance

issue : many networks
  ↳ lot of space
  ↳ " " computational resources

train a "small" network $g$
to imitate the ensemble method

task : $g \simeq \frac{1}{n} \sum_k f_k$

student →          ← teacher

↳ this training works → good performance

### Other ensemble techniques :

- Boosting : combining weak classifiers into a strong one
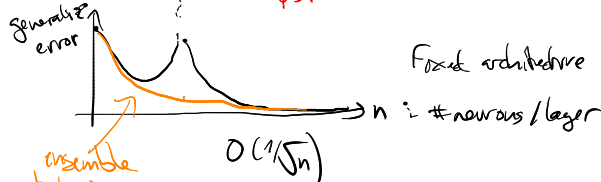
wrongly
classified
samples
(by $f_t$)

$f_t \longrightarrow f_{t+1}$ : different weights to samples
  - ↗ weight for wrongly-classified
  ↳ train a weak classifier with these weights
  ↳ agglomerate : $f_{t+1} = f_t + w c_t$
                                    $(x)$

### I-B  Generalization without regularization

no $l_2$ penalty on weights
no functional norm
→ loss : $\sum_{\text{samples } i} \| \hat{y}_i - y_i \|^2$ or cross-entropy

"double gradient descent"

generalize
error

ensemble
technique
over 20 networks
of same size

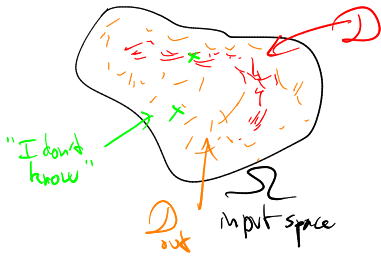$O(1/\sqrt{n})$

$n$ := #neurons / layer

Fixed architecture

→ Neural Tangent Kernel
  ↳ infinite-width layers

# I-C  Generalize only when relevant

[Towards N.N. that provably know when they don't know]



└→ use 2 distributions

Matthias Hein's team

→ $\mathcal{D}$ : data "in"     ex: CIFAR

→ $\mathcal{D}_{out}$ : to model samples outside of $\mathcal{D}$
     ex: ImageNet

"I don't know"

$\mathcal{D}_{out}$   input space

└→ kernels → "distance" to $\mathcal{D}_{in}$
              → " to $\mathcal{D}_{out}$

Training:
- samples in $\mathcal{D}_{in}$ : loss (data loss) → $\begin{pmatrix} 0 \\ 0 \\ 1 \\ 0 \end{pmatrix}$ target (label)

- " in $\mathcal{D}_{out}$ : loss: $\hat{y} \simeq \begin{pmatrix} 1/u \\ 1/u \\ 1/u \\ \vdots \end{pmatrix}$ uniform distribution over $u$ classes

---

# II  Learning from noisy data

## Denoising auto-encoder  [ICML 2008]

- dealing with noisy data $\simeq$ noise modeling

- originally:  $x \rightarrow \square\square\square \cdot \square \rightarrow \tilde{x}$   auto-encoder
             usually: narrow middle layer to concentrate information

- set-up:

  input:
  noisy version of $x$ → $x^{av}$ → $\square \square \square \square \square$ → goal: reconstruct $x$   i.e. get rid of noise

                        └→ learn | more robust features | to denoise

## Classification with noisy labels                    [D.L. is robust to massive label noise; 2017]

- true distribution $(x, y)$
- given noisy " $(x, \tilde{y})$ ⟩ → sometimes $y$
                              → sometimes random labels

- if some samples are mislabeled (in the training set) → not much change of accuracy
- a significant part of — — — — — — — → ?
    90%
    99%      → still possible to get reasonable results
             provided that data is available in large quantities

             └→ what matters is the number of correct labels

                → noise 90% → ×10 data
                    99% → ×100 data

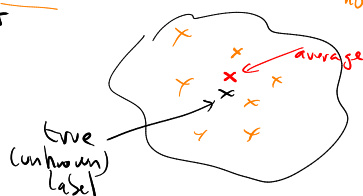## Regression with noisy labels

- dataset $(x, y + \varepsilon)$
                 ⟨ noise iid centered
                   variance

- simplistic case:
    - always the same input $x$

noise $\varepsilon$



true (unknown) label

average

noisy target
$\tilde{y}_i = y(x) + \varepsilon_i$

$L^2$ loss

$\sum_i \| \hat{y}(x) - \tilde{y}_i \|^2$
     └── always same predictions → average of $\tilde{y}_i$
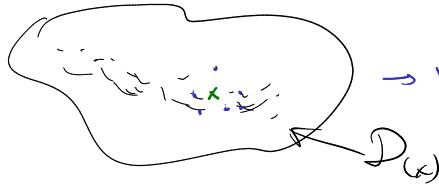
[Noise2Noise]    denoising     <⟸    classical statistics
                 effect:                N samples ⟹ $|\mu - y| \sim \frac{1}{\sqrt{N}} \times \sigma$
                 $1/\sqrt{N}$ factor
                                                    ↑     ↑
                                                 $\frac{1}{\#i}\sum_i \hat{y}_i$  true
                                                               label

, real case : x varies



→ N.N. will perform an "average" over similar points

$\sim$   $\frac{d}{\uparrow \#\text{similar examples}}$

↳ define "similarity"
   according to the network:
   based on the ability of the N.N. to distinguish samples

↳ similarity between samples x and x':   $\frac{\nabla_\theta \hat{y}_\theta(x)}{\|\;\;\|} \cdot \frac{\nabla_\theta \hat{y}_\theta(x')}{\|\;\;\|}$   $\in \mathbb{R}$ if $\hat{y} \in \mathbb{R}$

                                                                $\in [-1, 1]$

↳ how does label noise in training set
   affect prediction at a given test point?

      ↳ can compute a denoising factor

          noise σ   →  prediction noise : σ/factor
          training                              ↑ depends on the test sample

[Understanding blackbox predictions via Influence functions, ICML2017]
   ↳ based on classical statistics
         → noise on x
         → noise on the sampling distribution : probability to pick x
      ↳ same kind of formulas
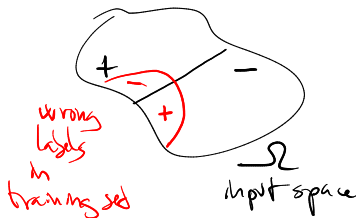         → ∇ Loss instead of ∇ŷ
         → inverse of Hessian



increase
sampling
rate in
specific area

Classification :

   [Learning with noisy labels]
      ↳ Ass: know the amount of noise
         ex: 20% of labels are inverted    (binary classifier)



wrong
labels
in
training set

input space

→ build a loss : unbiased
   ↳ on average, correct loss!
              as if    with true labels
   ↕
   samples
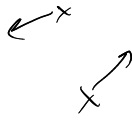
→ variance?
   Rademacher bound !
   ↳ notion of complexity

# III Formal proofs

- to prove (formally) that the N.N. will never mistake
- express a property → check it

$$\beta : \forall \text{ input samples } x \in \chi, \ |\hat{y}(x) - y(x)| \leq \varepsilon$$

completely specified

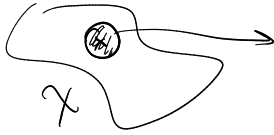- ACAS-Xu :

plane

given : locations
-speeds    (no perception)

prove : no collision

small network
3 inputs ? → Few outputs
1000 neurons
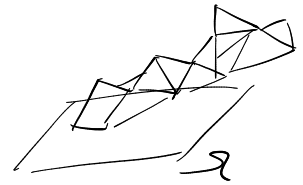
- local properties : locally robust to adversarial attacks

→ no adversarial
ex in that ball

$\chi$

- keywords : - abstraction
            - ReLuplex    → N.N. with ReLU
                          ↳ piecewise-affine $f$

- naive application
of traditional provers → complexity $\sim 2^{\# neurons}$