

DATASET

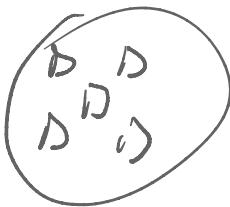
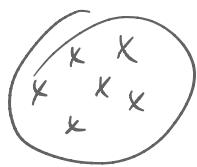
CLASSIFICATION

$X \xrightarrow{h} Y$ LABEL OF A CLASS
INPUT (DOG/CAT)

REGRESSION

$X \xrightarrow{h} Y \in \mathbb{R}$

CLUSTERING \rightarrow UNSUPERVISED LEARNING



BIDIMENSION REDUCTION

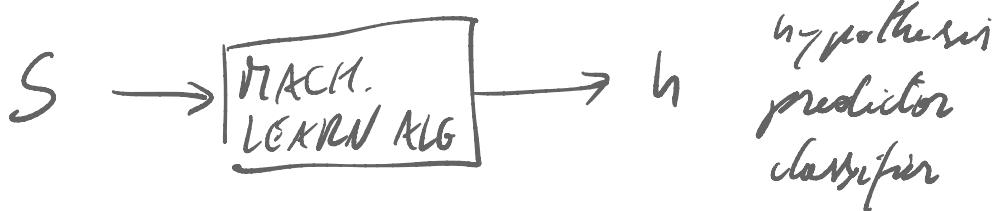
REINFORCEMENT LEARNING

CLASSIF / REGR PROBLEMS

$$\Rightarrow S = \{(x_i, t_i), i=1, 2, \dots, n\}$$

SUPERVISED LEARNING

$$x \xrightarrow{h} y$$



How good/bad h ?

$$h(x) \longleftrightarrow y$$

(x, y)
CLASS

$$\boxed{1I(h(x) \neq y)}$$

indicator
function

LOSS FUNCTION

$$l(h, (x, y)) = \begin{cases} 1I(h(x) \neq y) & \text{CLASSIF.} \\ (h(x) - y)^2 & \text{REGR.} \end{cases}$$

TOMORROW (x, y) FOLLOW DISTRIBUTION δ

$\mathbb{E} [l(h, (x, y))]$ \triangleq expected loss/risk
 ~~$(x, y) \sim D$~~ $L_\delta(h)$ $R(h)$
 $R_\delta(h)$

FIND A GOOD PREDICTOR \rightarrow MINIMIZE
THE EXPECTED LOSS

minimize $L_\delta(h) = \min_{h \in H} \mathbb{E} [l(h, (x, y))]$
 δ is UNKNOWN

H : class of hypotheses/predictors.
ML IS "DISTRIBUTION-FREE"

ML \neq STATISTICS

ML ASSUMPTION:

WE DO NOT KNOW δ

BUT S IS DRAWN ACCORDING
TO δ

EMPIRICAL LOSS

$$L_S(h) = \frac{1}{n} \sum_{i=1}^n \ell(h, (x_i, y_i)) \simeq \underset{(x,y) \sim D}{\mathbb{E}} [\ell(h, (x, y))]$$

MORE CORRECT
 THE LARGER IS n

Empirical Risk Minimization (ERM)

minimize $L_S(h) = \frac{1}{n} \sum_{i=1}^n \ell(h, (x_i, y_i))$

$$\underset{h \in H}{\text{minimize}}$$

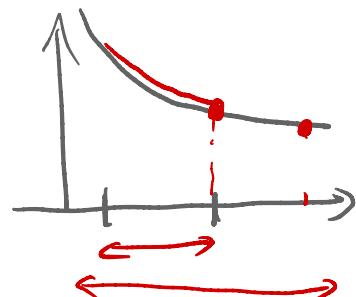
$$H_1 = \{a x + b, a, b \in \mathbb{R}\}$$

$$H_2 = \{a x^2 + b x + c, a, b, c \in \mathbb{R}\}$$

$$H_1 \subset H_2$$

$$h_1^* \in \underset{h \in H_1}{\text{argmin}} L_S(h) \quad h_2^* \in \underset{h \in H_2}{\text{argmin}} L_S(h)$$

$$L_S(h_1^*) \geq L_S(h_2^*)$$





STATISTICAL LEARNING

$$|L_S(h) - L_\delta(h)| > C \sqrt{\frac{d + \ln \frac{1}{\delta}}{n}}$$

with probability δ

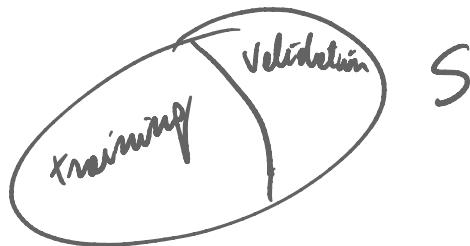
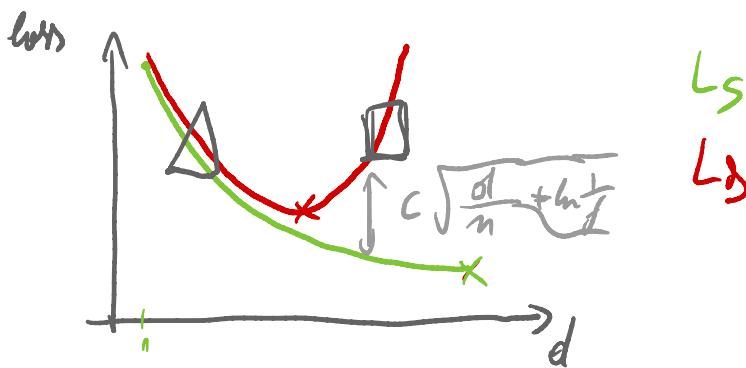
d : VC-dimension
measure of the complexity
of the class H

Vapnik

Chervonenkis

for the class of polynomials in one variable with degree m :

$$d = m + 1$$



$$S = S_T \cup S_V$$

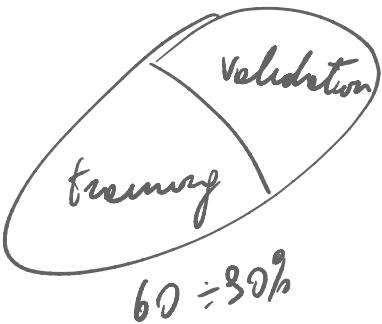
$$h_1^* \in \underset{h \in H_1}{\operatorname{arg\,min}} L_{S_T}(h)$$

$$h_2^* \in \underset{h \in H_2}{\operatorname{arg\,min}} L_{S_T}(h)$$

$$L_{S_T}(h_1^*) \geq L_{S_T}(h_2^*)$$

$$L_{S_V}(h_1^*) \leq L_{S_V}(h_2^*)$$

- pick h_1^*
 h_1^* is OVERFITTING
- pick h_2^*
 h_2^* is UNDERFITTING



K -fold cross validation



$$\hat{h}_1^{*(1)} \in \underset{h \in H_1}{\operatorname{argmin}} L_{S_1}(h)$$

$$\hat{h}_1^{*(1)} \in \underset{h \in H_1}{\operatorname{argmin}} L_{S_2 \cup \dots \cup S_K}(h)$$

$$\hat{h}_1^{*(2)} \in \underset{h \in H_1}{\operatorname{argmin}} L_{S_1 \cup S_3 \cup \dots \cup S_K}(h)$$

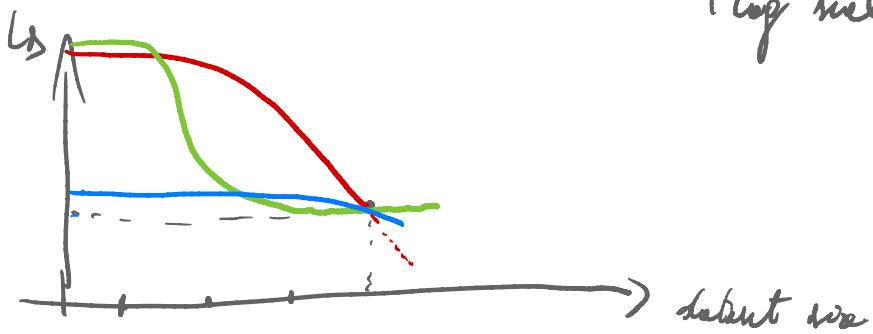
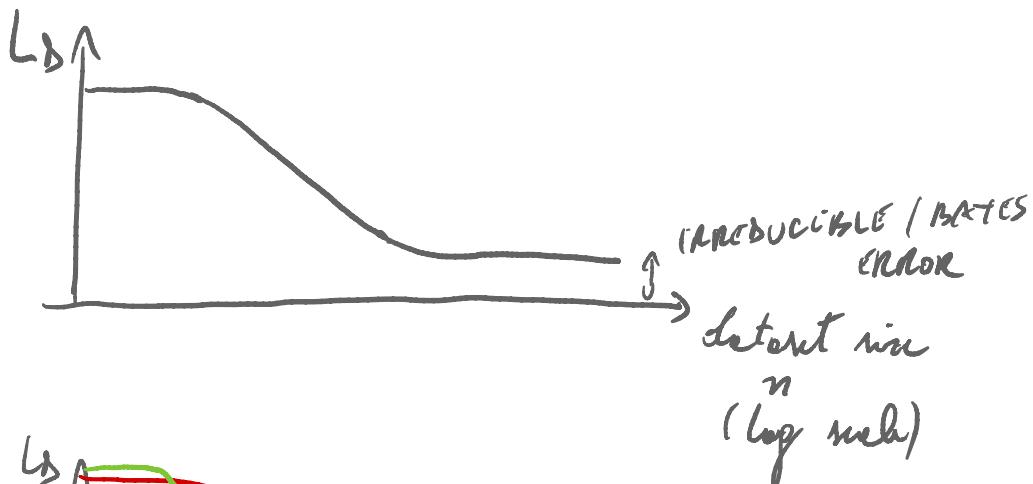
$$\hat{h}_1^{*(u)} \in \underset{h \in H_1}{\operatorname{argmin}} L_{S_1 \cup S_2 \cup \dots \cup S_{u-1}}(h)$$

$$\text{Validation error} = \frac{1}{K} \sum_{i=1}^K L_{S_i}(\hat{h}_1^{*(i)})$$

- learning step / rate
- batch size
- regularization coefficient

HYPER PARAMETERS

LEARNING PROBLEM \equiv SOLVING MANY
OPTIMIZATION
PROBLEMS



minimize $L_S(h)$
 $h \in H$

h is parameterized by a vector of
parameters $w \in \mathbb{R}^d$

$$h(x) = \sum_i w_i x_i$$

$$L_S(h) = \sum_{i=1}^n \ell(h, (x_i, y_i)) =$$

ex
quadratic
loss

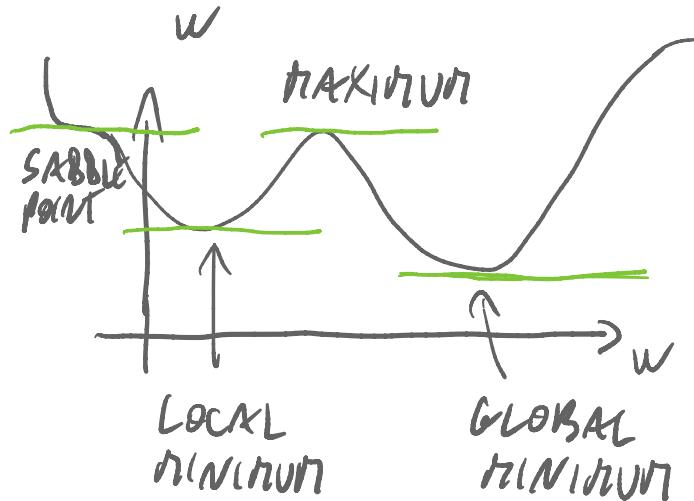
$$= \sum_{i=1}^n (h(w, x_i) - y_i)^2$$

$$\ell(h(x_i, y_i)) = f(w, i)$$

$$L_S(h) \longleftrightarrow F(w)$$

$$\begin{array}{l} \text{minimize}_{w \in W} F(w) \rightarrow \text{minimize}_{w \in \mathbb{R}^d} F(w) \end{array}$$

minimum $F(w)$



$F'(w)$

$$\nabla F(w) = \begin{bmatrix} \frac{\partial F}{\partial w_1} \\ \vdots \\ \frac{\partial F}{\partial w_d} \end{bmatrix} \in \mathbb{R}^d$$

$$F(w) = w_1^2 + w_2^2$$

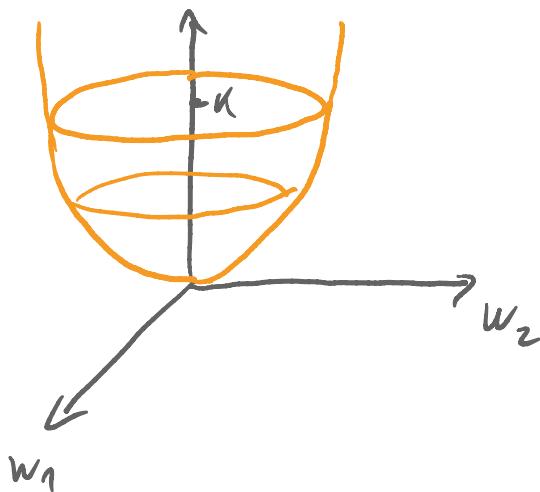
$$\frac{\partial F}{\partial w_1} = 2w_1 + 0 = 2w_1$$

$$\frac{\partial F}{\partial w_2} = 2w_2$$

$$\nabla F(w) = \begin{bmatrix} 2w_1 \\ 2w_2 \end{bmatrix}$$

AT A CRITICAL POINT $\Leftrightarrow \nabla F = 0$

$$\begin{bmatrix} 2w_1 \\ 2w_2 \end{bmatrix} = \begin{bmatrix} 0 \\ 0 \end{bmatrix} \rightarrow w_1 = w_2 = 0$$



$$F(w_1, w_2) = K$$

$$\underline{w_1^2 + w_2^2 = K}$$