

REFERENCES:

- BISHOP: PATTERN RECOGNITION AND MACHINE LEARNING
 - MY PHD THESIS
 - MITCHELL: MACHINE LEARNING
-

DEFINITION:

LEARNING = IMPROVING WITH EXPERIENCE
AT SOME TASK

EXPERIENCE: „DATA“

TASK: REGRESSION / CLASSIFICATION /
CLUSTERING

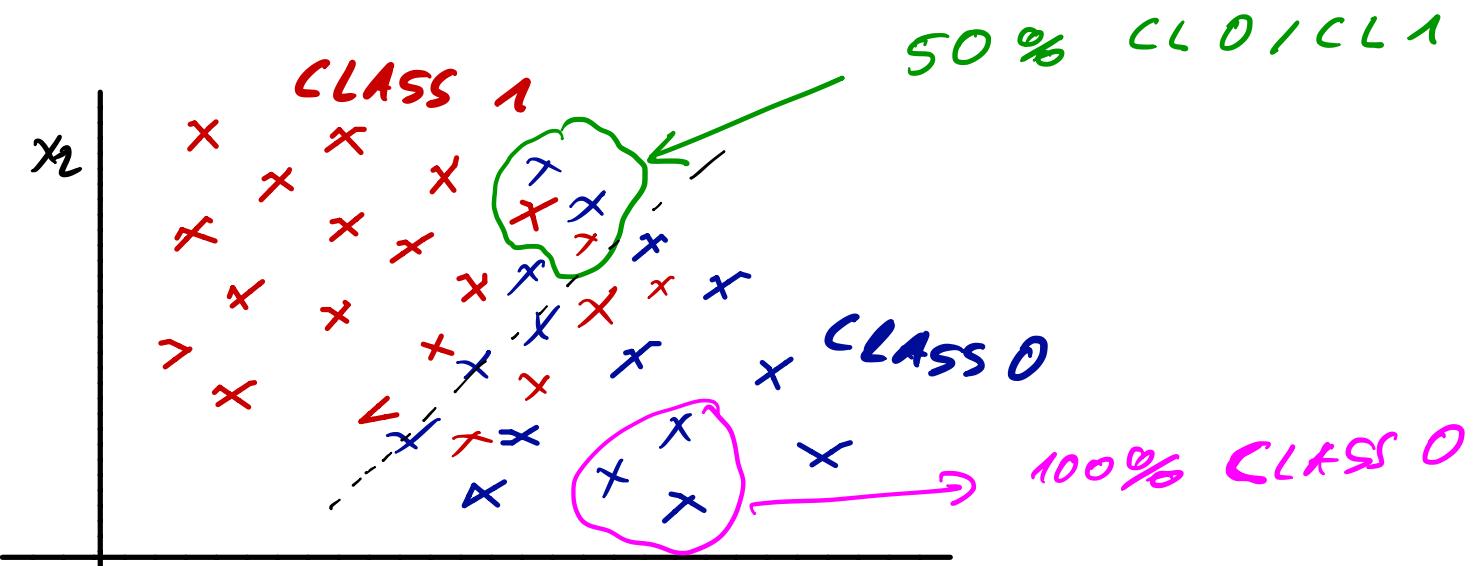
BASIC IDEA:

① FORMULATE LIKELIHOOD FOR
PROBLEM (CLASSIFICATION /
REGRESSION)

② FIND SUITABLE FUNCTIONS TO
PARAMETERIZE STATISTICS OF
LIKELIHOOD

③ FIT DATA

CLASSIFICATION



"COIN-FLIP EXPERIMENT WITH SUCCESS PROBABILITY THAT DEPENDS ON POSITION IN INPUT SPACE"

EXPERIMENT: COIN FLIP WITH STATES:
 $S = \{0, 1\}$

$$P(X=0) = \mu \quad P(X=1) = 1-\mu$$

$$P(X=s) = \mu^s (1-\mu)^{1-s}$$

LIKELIHOOD:

$$P(D|\mu) = \prod_{i=1}^N \mu^{x_i} (1-\mu)^{1-x_i}$$

LOG LIKELIHOOD:

$$\log(p(D|\mu)) = \sum_{i=1}^N x_i \log(\mu) + (1-x_i) \log(1-\mu)$$

"CROSS-ENTROPY COST FUNKT"

MINIMUM OF LOG LIKELIHOOD:

$$\mu_M = \frac{1}{N} \sum_i x_i$$

INTRODUCE DEPENDENCE ON INPUTS

$$P(X=0) = \mu \longrightarrow P(X=0 | \vec{x}) = \frac{1}{1 + \exp(-f(\vec{x}))}$$

CONDITIONAL

ON POSITION
IN INPUT SPACE

$$P(X=0 | \vec{x}) = \frac{1}{1 + \exp(-f(\vec{x}))}$$

\Rightarrow INTERPRET AS PROBABILITY SINCE
BETWEEN $[0, 1]$. NOTE: OTHER FNCTS.
POSSIBLE

$f(\vec{x})$ IS SOME FLEXIBLE FNKT:

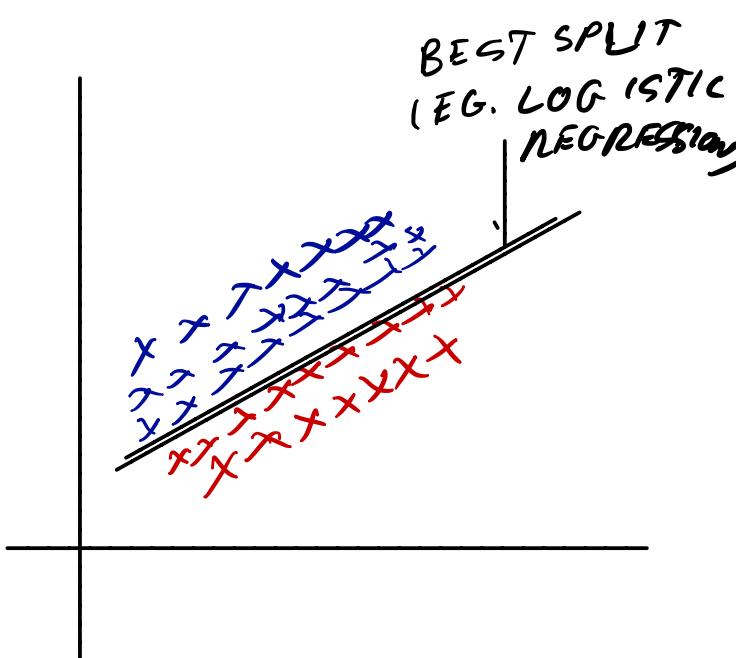
\times LINEAR BASIS FNKT: $f(\vec{x}) = \sum_i w_i \phi_i(\vec{x})$

\times NEURAL NETWORK: $f(\vec{x}) = \sum_{j=0}^M w_{kj}^{(2)} h\left(\sum_{i=0}^D w_{ji}^{(1)} x_i\right)$

h : SIGMOID / TANH.

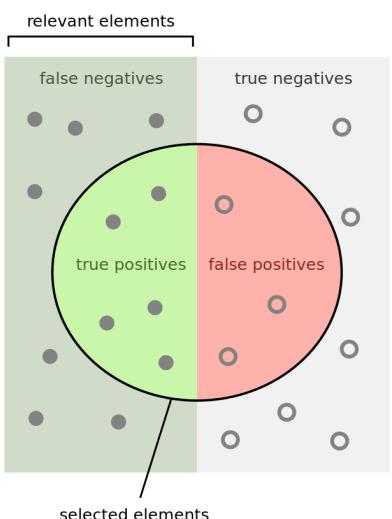
NOTE: DIFFERENT ML MODELS IMPOSE
DIFFERENT $f(\vec{x})$

EXAMPLE :



DIFFERENT DATA \Rightarrow DIFFERENT MODEL
PERFORMS BETTER

HOW TO EVALUATE CLASSIFIERS?



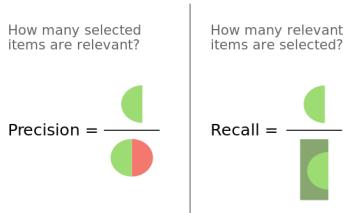
PROBLEM: IMBALANCED
DATASET : 1% LABEL 1
99% " 2
 \Rightarrow CLASSIFIER " THAT
ALWAYS PREDICTS 2
HAS 99% ACCURACY

\Rightarrow F1 SCORE:

$$F1 = \frac{2}{\frac{1}{RECALL} + \frac{1}{PRECISION}}$$

$$= 2 \frac{PREC \cdot RECALL}{PREC + RECALL}$$

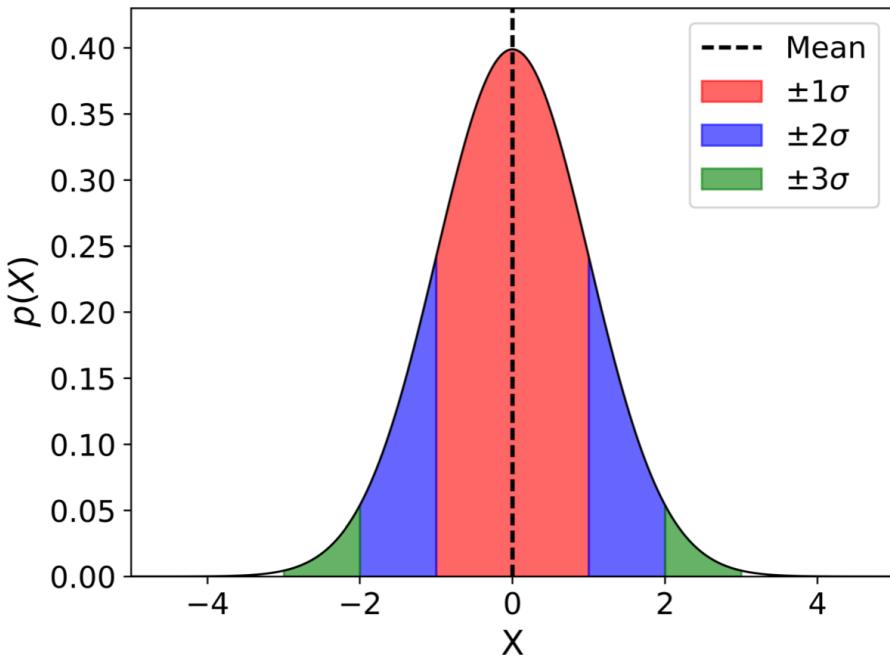
FROM WIKIPEDIA



REGRESSION

GAUSSIAN DISTRIBUTION: 2 PARAMETERS:

μ : MEAN
 σ : STANDARD DEVIATION



LIKELIHOOD:

$$P(D|\mu, \sigma) = \prod_{i=1}^N N(x_i | \mu, \sigma)$$

$$\log(P(D|\mu, \sigma)) = -\frac{1}{2\sigma^2} \sum_{i=1}^N (x_i - \mu)^2 - \frac{N}{2} \log(\sigma^2) - \frac{N}{2} \log(2\pi)$$

$$\mu_{ML} = \frac{1}{N} \sum_{i=1}^N x_i$$

$$\sigma_{ML} = \sqrt{\frac{1}{N} \sum_{i=1}^N (x_i - \mu_{ML})^2}$$

MEAN SQUARED ERROR

INTRODUCE DEPENDENCE ON INPUT

$$P(y) \longrightarrow P(y|x)$$

$\mu \rightarrow f(x) \Rightarrow$ MINIMIZE THE MSE
USE "SAME" FUNKTS. AS FOR
CLASSIFICATION.

x HOW COMPLEX SHOULD THE FUNKT. BE?

x WHAT IF $P(y|x)$ IS NOT GAUSSIAN?

MODEL COMPLEXITY // BIAS - VARIANCE

TRADE OFF

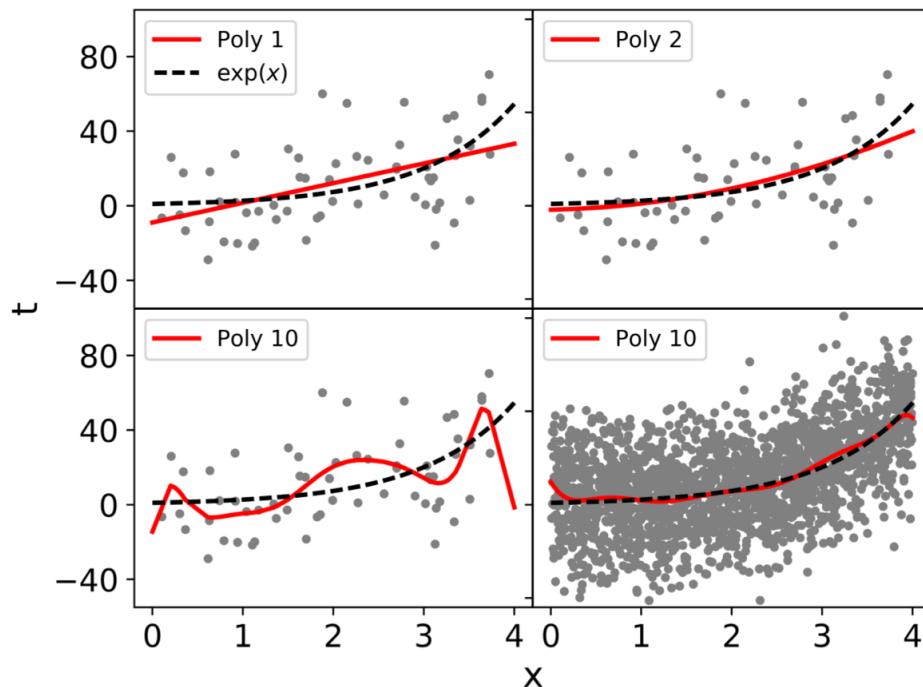


Figure 2.4: Illustration of the impact of model complexity on the accuracy of Machine Learning models. The black dashed line shows the exponential function from which $N = 60$ grey data points have been generated by adding zero mean Gaussian noise with $\sigma = 20$. The red curves show the polynomial models of order M that have been fitted to this data set. The bottom right panel shows the performance of a polynomial model with $M = 10$ fitted on a larger sample with $N = 2000$.

$$E[L] = \iint [y(x) - t]^2 p(\vec{x}, t) d\vec{x} dt$$

\vec{x} : INPUT

t : KNOWN TRAINING DATA RESPONSE

$y(x)$: FNKT. TO BE FITTED

$$\frac{\partial E[L]}{\partial y(\vec{x})} = 2 \iint [y(x) - t] p(\vec{x}, t) d\vec{x} dt = 0$$

$$\Rightarrow \underline{y(x) = \frac{\int t p(f|x) dt}{\text{CONDITIONAL MEAN}}}$$

IRRESPECTIVE
OF FORM OF
 $p(f|x)$!

FNKT. THAT
MINIMIZES MSE

\Rightarrow ESTIMATES
CONDITIONAL
MEAN

ASSUME WE HAVE SUCH AN ESTIMATE

$$\hat{g}(x)$$

$$E[L] = \iint [\hat{g}(x) - t]^2 p(x, t) dx dt$$

EXPAND AROUND OPTIMAL $y^*(x)$

$$E[L] = \int (\hat{g}(\vec{x}) - g(\vec{x}))^2 p(\vec{x}) d\vec{x} \quad \textcircled{A}$$

$$+ \iint (y^*(\vec{x}) - t)^2 p(\vec{x}, t) dx dt \quad \textcircled{B}$$

(A): HOW CLOSE IS $\hat{g}(x)$ TO OPTIMAL?

(B): WHAT IS THE "INTRINSIC" ERROR
IN THE DATA

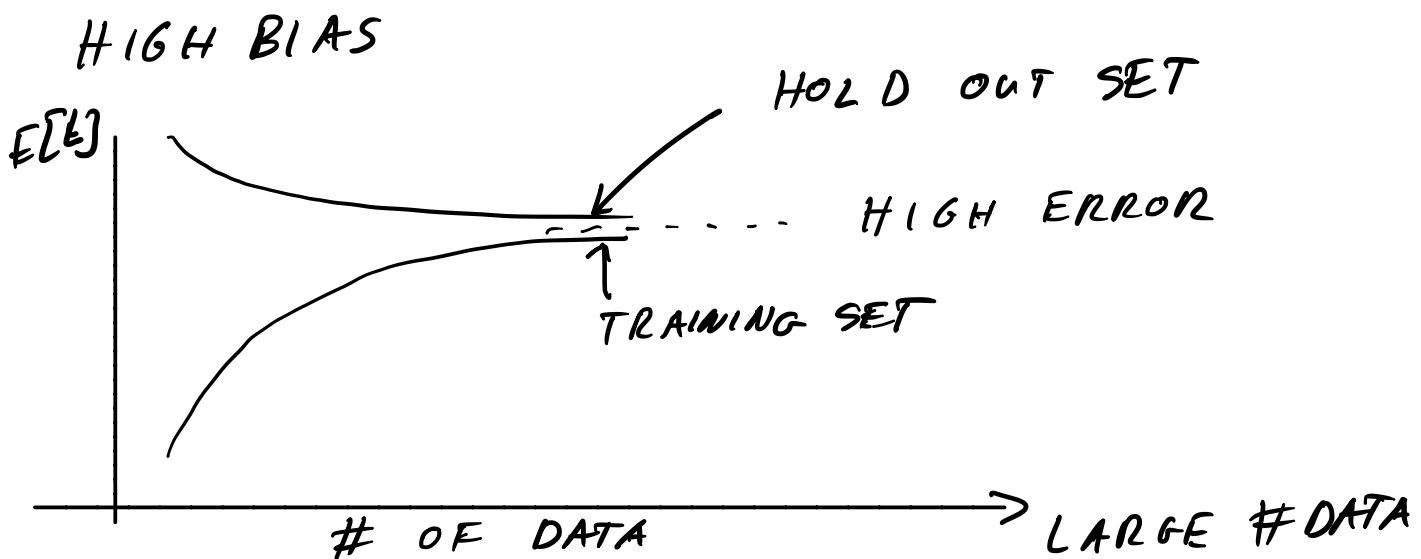
DECOMPOSE \textcircled{A} FURTHER:

$$\bar{E}[L] = (\text{BIAS})^2 + \underbrace{\text{VARIANCE} + \text{NOISE}}_{\textcircled{B}}$$

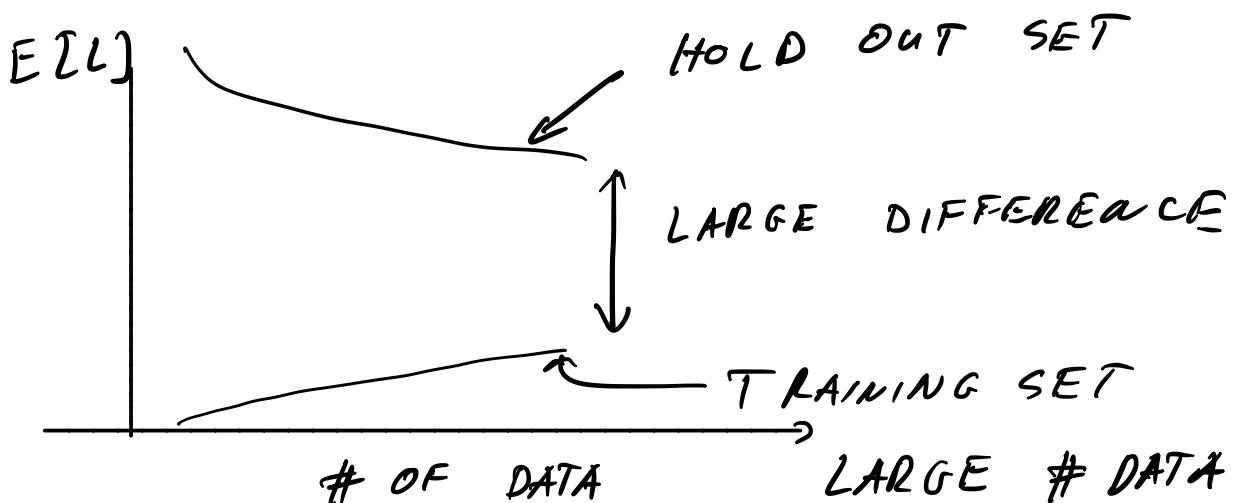
$$(\text{BIAS})^2 = \int (E_D[\hat{g}(\vec{x}; D)] - g(\vec{x}))^2 p(\vec{x}) d\vec{x}$$

$$\text{VARIANCE} = \int E_D[(\hat{g}(x; D) - E_D[\hat{g}(x; D)])^2] p(x) dx$$

2 SCENARIOS:



=) MAKE MODEL MORE COMPLEX
(FEATURE ENGINEERING, MORE LAYERS,
HIGHER ORDER POLYNOMIALS....)



=) ADD MORE DATA, MAKE MODEL LESS COMPLEX;
"REGULARIZATION", LESS WEIGHTS, ...