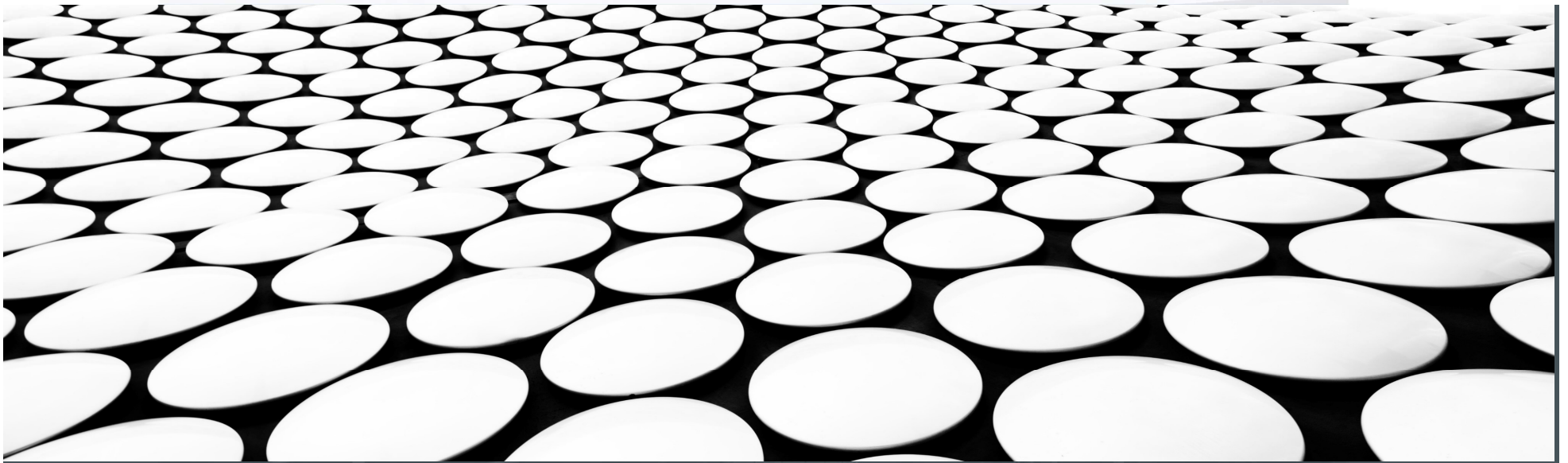

INTRODUCTION TO STATISTICAL LEARNING THEORY

ROBERTO BASILI, UNIVERSITÀ DI ROMA, TOR VERGATA

DEEP LEARNING, MARCH 2024





OUTLINE

- Statistical Learning Theory
 - PAC learnability
 - VC dimension
 - Learning Machines
 - Model Optimization and Concept Class complexity
 - Model Optimization via Cross-Validation
 - Towards perceptrons and SVMs

■

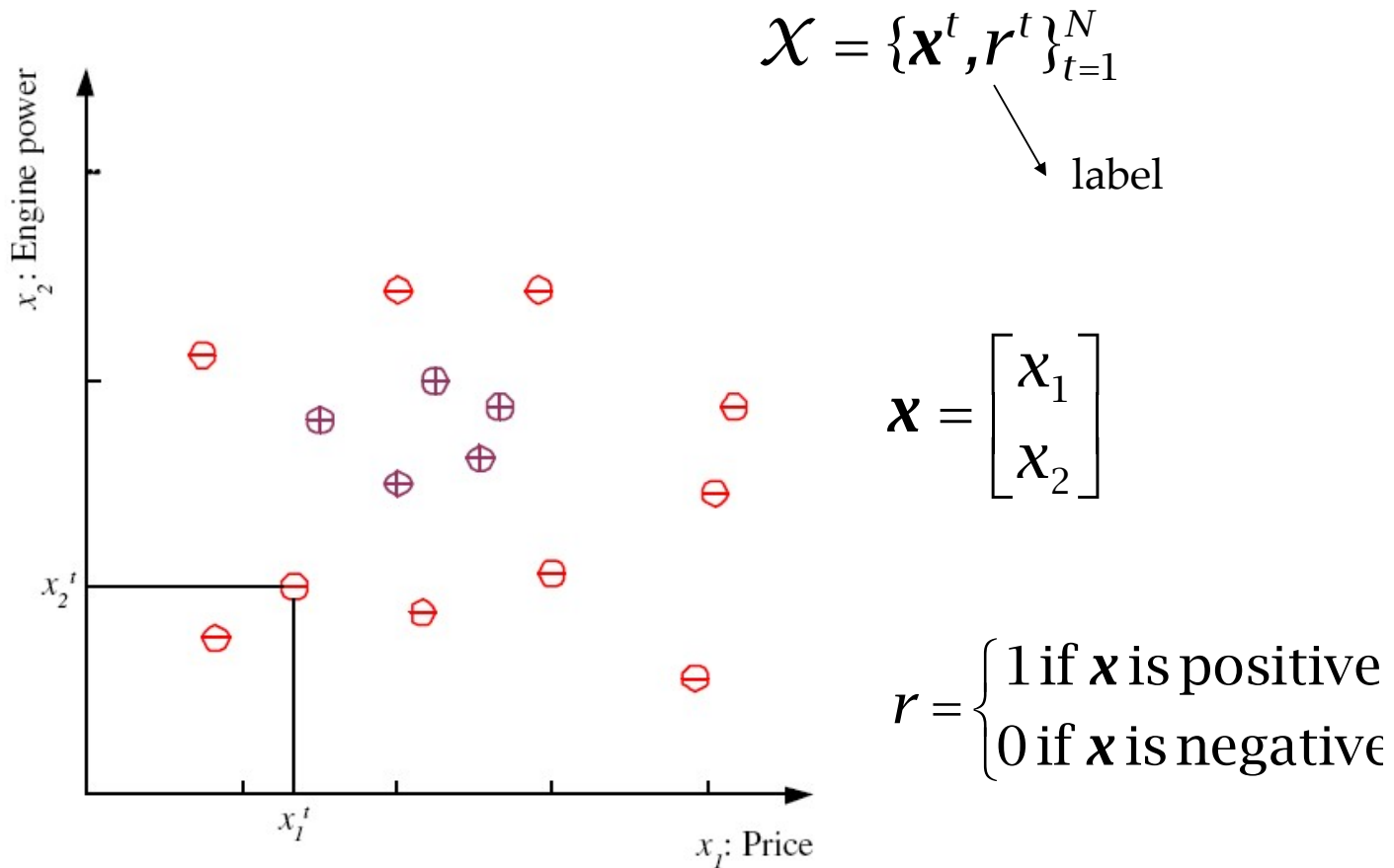


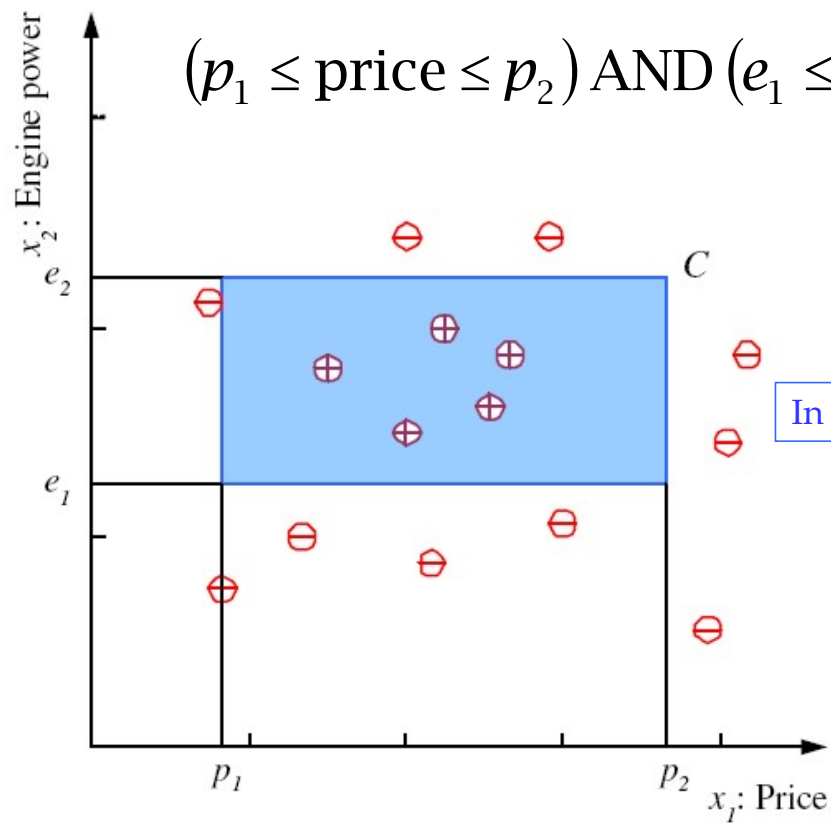
FROM STATISTICAL LEARNING THEORY TO SVMs

LEARNING A CLASS FROM EXAMPLES

- Class C of a “family car”
 - **Prediction** Is car x a “family car”?
 - **Knowledge extraction** What do people expect from a family car?
- Output:
 - Positive (+) and negative (-) examples
- Input representation:
 - x_1 : price, x_2 : engine power

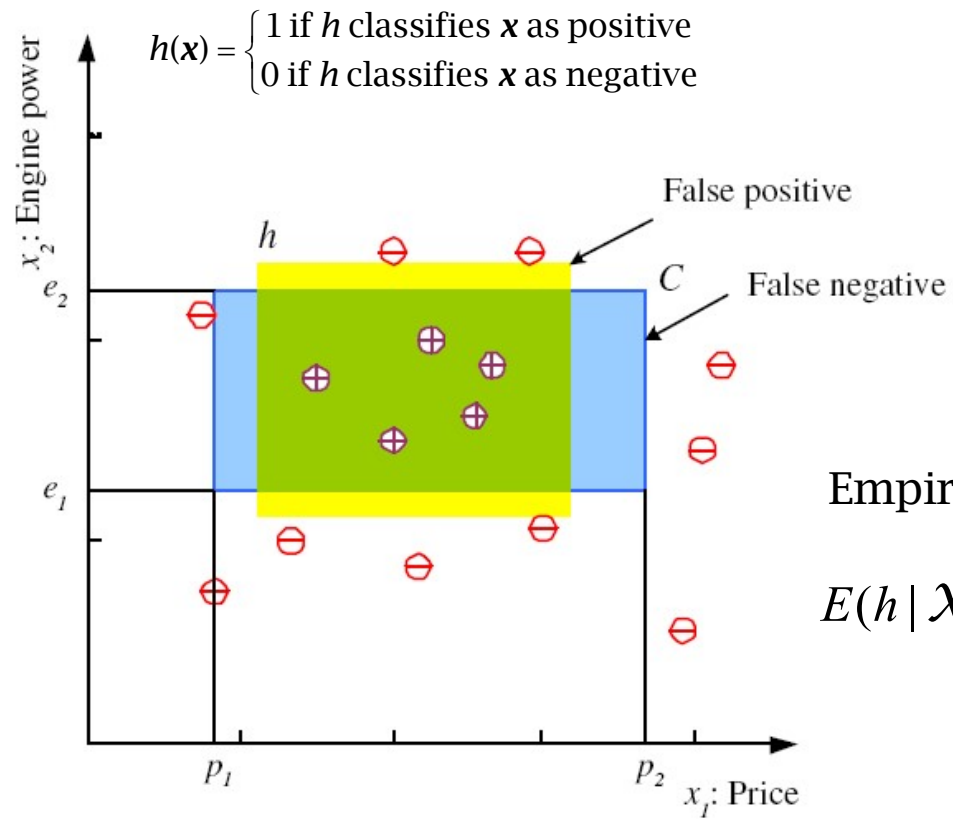
TRAINING SET \mathcal{X}



CLASS C 

In general we do not know $C(x)$.

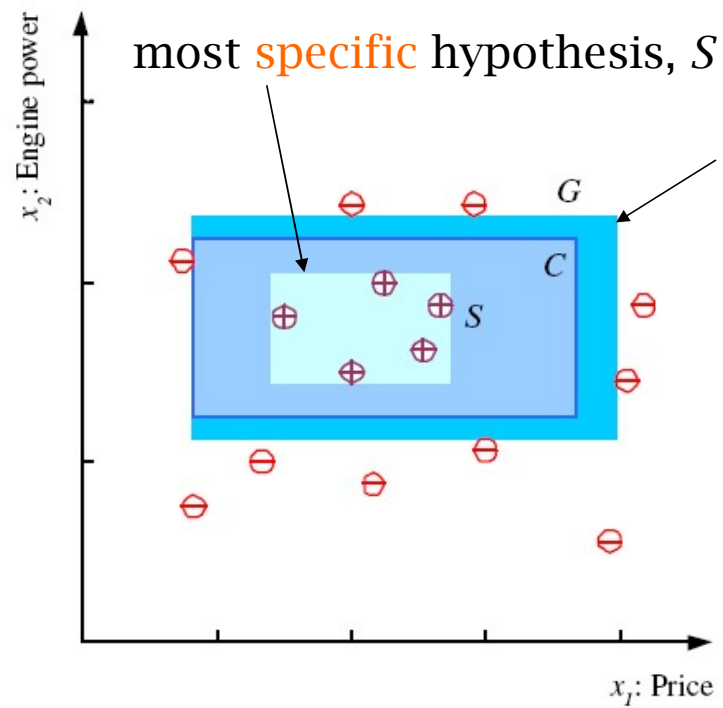
HYPOTHESIS CLASS \mathcal{H}



Empirical error:

$$E(h | \mathcal{X}) = \sum_{t=1}^N 1(h(\mathbf{x}^t) \neq r^t)$$



S, G, AND THE VERSION SPACE



$h \in \mathcal{H}$, between S and G is **consistent**
and make up the **version space**

(Mitchell, 1997)

PROBABLY APPROXIMATELY CORRECT (PAC) LEARNING

- How many training examples are needed so that the tightest rectangle S which will constitute our hypothesis, will **probably** be **approximately correct**?
 - We want to be **confident** (*above a level*) that 
 - ... the **error probability is bounded** by some value 

- A concept class C is called **PAC-learnable** if there exists a PAC-learning algorithm such that, for any $\epsilon > 0$ and $\delta > 0$, there exists a fixed sample size such that, for any concept $c \in C$ and for any probability distribution on X , the learning algorithm produces a probably-approximately-correct hypothesis h
- a (PAC) **probably-approximately-correct hypothesis** h is one that has error at most ϵ with probability at least $1-\delta$.

PROBABLY APPROXIMATELY CORRECT (PAC) LEARNING

- In PAC learning, given a class C and examples drawn from some unknown but fixed distribution $p(x)$, we want to find the number of examples N , such that with probability at least $1-\delta$, h has error at most ε ? (Blumer et al., 1989)

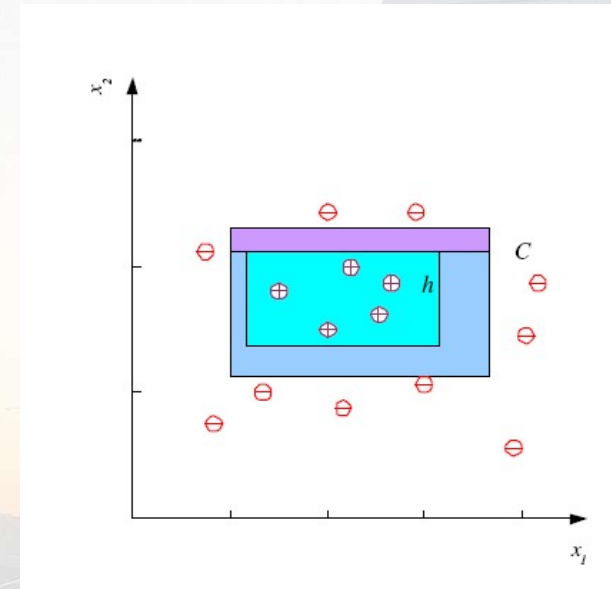
- $$P(C\Delta h \leq \varepsilon) \geq 1-\delta$$

- where $C\Delta h$ is (area of the) “the region of difference between C and h ”, and $\delta>0, \varepsilon>0$.

PAC LEARNING

How many training examples m should we have, such that with probability at least $1 - \delta$, h has error at most ϵ ? (Blumer et al., 1989)

- Let prob. of a + ex. in each strip be at most $\epsilon/4$
- Pr that a random ex. misses a strip: $1 - \epsilon/4$
- Pr that m random instances miss a strip:
 $(1 - \epsilon/4)^m$
- Pr that m random instances miss 4 strips:
 $4(1 - \epsilon/4)^m$
- We want $1 - 4(1 - \epsilon/4)^m \geq 1 - \delta$ or $4(1 - \epsilon/4)^m \leq \delta$
- Using $1 - x \leq e^{-x}$ an even stronger condition is:
 $[(1 - \epsilon/4)^m \leq \exp(-\epsilon m/4) \text{ so } (1 - \epsilon/4)^m \leq \exp(-\epsilon m/4) = \exp(-\epsilon m/4)]$
 $4e^{-\epsilon m/4} \leq \delta$ OR
- Divide by 4, take $\ln \dots$ and show that $m \geq (4/\epsilon) \ln(4/\delta)$



PROBABLY APPROXIMATELY CORRECT (PAC) LEARNING

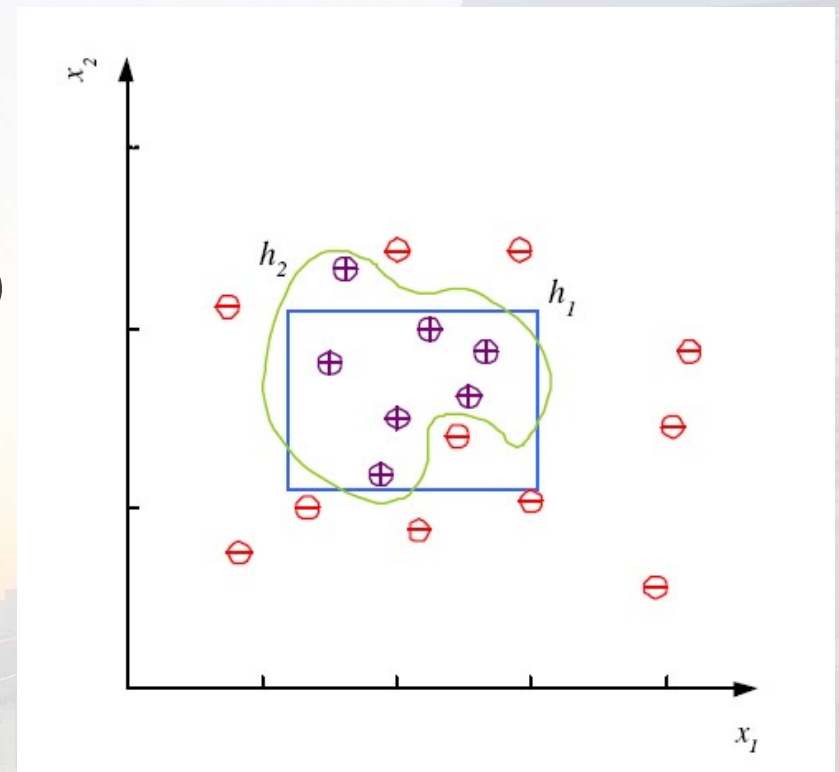
How many training examples m should we have, such that with probability at least $1 - \delta$, our hypothesis h has error at most ϵ ? (Blumer et al., 1989)

$$m \geq (4/\epsilon) \ln(4/\delta)$$

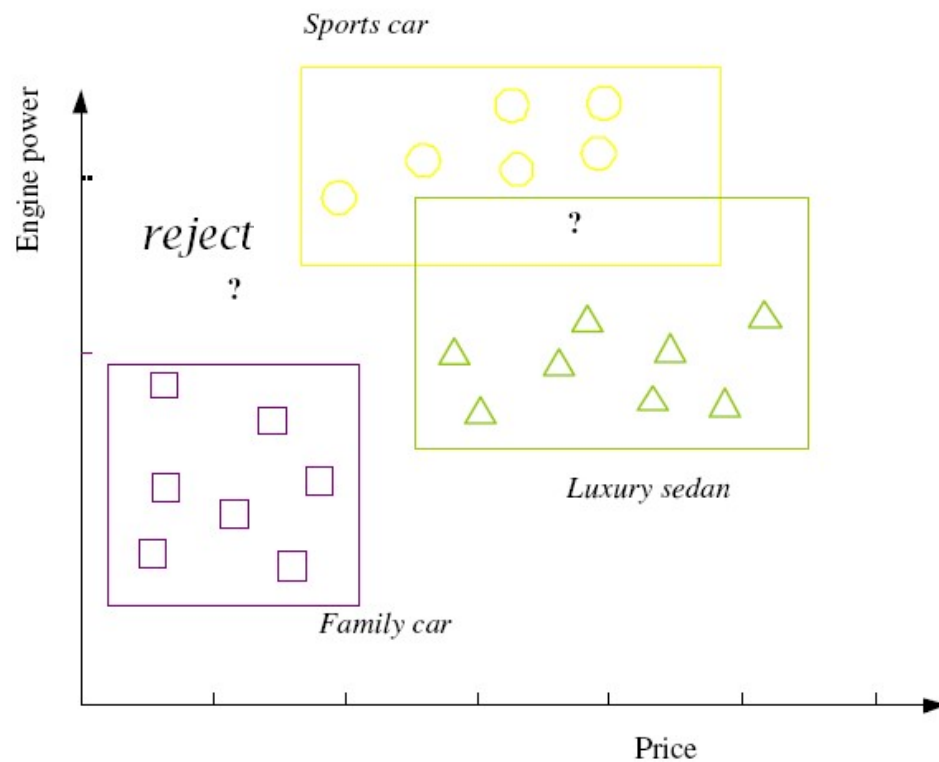
- m increases slowly with $1/\epsilon$ and $1/\delta$
- Say $\epsilon=1\%$ with confidence 95%, pick $m \geq 1752$
- Say $\epsilon=10\%$ with confidence 95%, pick $m \geq 175$

MODEL COMPLEXITY VS. NOISE

- Use the simpler one because
- Simpler to use (lower computational complexity)
- Easier to train (lower space complexity)
- Easier to explain (more interpretable)
- Generalizes better (lower variance – Occam's razor)



MULTIPLE CLASSES, $C_i, i=1, \dots, K$



$$\mathcal{X} = \{\mathbf{x}^t, r^t\}_{t=1}^N$$

$$r_i^t = \begin{cases} 1 & \text{if } \mathbf{x}^t \in C_i \\ 0 & \text{if } \mathbf{x}^t \in C_j, j \neq i \end{cases}$$

Train hypotheses
 $h_i(\mathbf{x}), i=1, \dots, K:$

$$h_i(\mathbf{x}^t) = \begin{cases} 1 & \text{if } \mathbf{x}^t \in C_i \\ 0 & \text{if } \mathbf{x}^t \in C_j, j \neq i \end{cases}$$

REGRESSION

$$X = \{x^t, r^t\}_{t=1}^N$$
$$r^t \in \mathfrak{R}$$
$$r^t = f(x^t) + \varepsilon$$

$$\underline{E(h'|X)} = \frac{1}{N} \sum_{t=1}^N [r^t - h'(x^t)]^2$$

$$\underline{E(h|X)} = E(w_1, w_0|X) = \frac{1}{N} \sum_{t=1}^N [r^t - (w_1 x^t + w_0)]^2$$

