**CSAIL**

# Machine learning: lecture 11

Tommi S. Jaakkola
MIT CSAIL
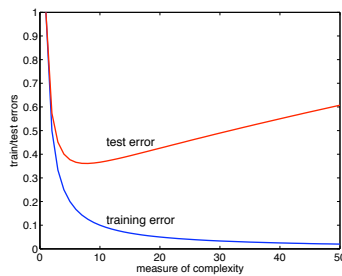*tommi@csail.mit.edu*

---

**CSAIL**

# Topics

- Complexity and generalization
  - finite set of classifiers
  - VC-dimension, learning

---

**CSAIL**

# Why care about "complexity"?



- We need a quantitative measure of complexity in order to be able to relate the training error (which we can observe) and the test error (that we'd like to optimize)

---

**CSAIL**

# Finite case

- We'll start by considering only a finite number of possible classifiers, $h_1(\mathbf{x}), \ldots, h_M(\mathbf{x})$ (e.g., randomly chosen linear classifiers)

- Key questions:

  1. Given $n$ training examples and $M$ possible classifiers how far can the training and test errors be?

  2. How many training examples do we need so that the errors are close?

  The answers will depend on $M$.

---

**CSAIL**

# Finite case: definitions

$$\hat{\mathcal{E}}_n(i) = \frac{1}{n}\sum_{t=1}^{n} \overbrace{\mathsf{Loss}(y_t, h_i(\mathbf{x}_t))}^{=0,1} = \text{empirical error of } h_i(\mathbf{x})$$

$$\mathcal{E}(i) = E_{(\mathbf{x},y)\sim P}\{\,\mathsf{Loss}(y, h_i(\mathbf{x}))\,\} = \text{expected error of } h_i(\mathbf{x})$$

---

**CSAIL**

# Finite case: definitions

$$\hat{\mathcal{E}}_n(i) = \frac{1}{n}\sum_{t=1}^{n} \overbrace{\mathsf{Loss}(y_t, h_i(\mathbf{x}_t))}^{=0,1} = \text{empirical error of } h_i(\mathbf{x})$$

$$\mathcal{E}(i) = E_{(\mathbf{x},y)\sim P}\{\,\mathsf{Loss}(y, h_i(\mathbf{x}))\,\} = \text{expected error of } h_i(\mathbf{x})$$

- Suppose we choose the classifier that minimizes the training error, $\hat{i}_n = \mathrm{argmin}_{i=1,\ldots,M}\,\hat{\mathcal{E}}_n(i)$, then

$$\text{Training error} = \hat{\mathcal{E}}_n(\hat{i}_n)$$
$$\text{Test error} = \mathcal{E}(\hat{i}_n)$$

## Finite case: errors

- The training and test errors,

$$\text{Training error} = \hat{\mathcal{E}}_n(\hat{i}_n)$$
$$\text{Test error} = \mathcal{E}(\hat{i}_n)$$

  are necessarily close if we can show that the errors are close for all the classifiers in our set:

$$|\hat{\mathcal{E}}_n(i) - \mathcal{E}(i)| \leq \epsilon, \ \text{ for all } i = 1, \ldots, M$$

- We can now express our key questions more formally in terms of $n$, $M$, and $\epsilon$

## Finite case: key questions revisited

- Key questions (rewritten):

  1. Given $n$ training examples and $M$ possible classifiers, what is the smallest $\epsilon$ such that

$$\max_{i=1,\ldots,M} |\hat{\mathcal{E}}_n(i) - \mathcal{E}(i)| \leq \epsilon$$

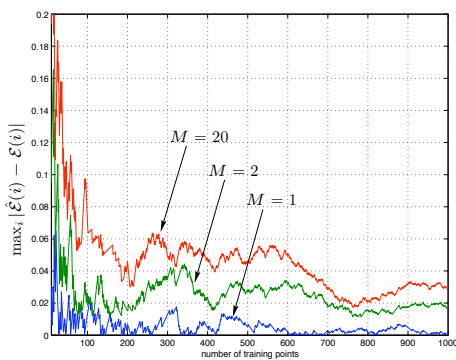  2. For a given $\epsilon$ how many training examples do we need so that

$$\max_{i=1,\ldots,M} |\hat{\mathcal{E}}_n(i) - \mathcal{E}(i)| \leq \epsilon$$

  Since training examples are sampled at random from some underlying distribution, we can only answer these questions probabilistically.

## Finite case: errors

## Finite case: probabilistic statement

- We can relate $n$, $M$, and $\epsilon$ by requiring that with high probability, the empirical errors of all the classifiers in our set are $\epsilon$-close to their expected errors:

$$P\Big( \max_{i=1,\ldots,M} |\hat{\mathcal{E}}_n(i) - \mathcal{E}(i)| \leq \epsilon \Big) \geq 1 - \delta$$

  The probability is taken over the choice of the training set and $1 - \delta$ specifies our confidence in the probabilistic statement.

## Finite case: probabilistic statement

- We can relate $n$, $M$, and $\epsilon$ by requiring that with high probability, the empirical errors of all the classifiers in our set are $\epsilon$-close to their expected errors:

$$P\Big( \max_{i=1,\ldots,M} |\hat{\mathcal{E}}_n(i) - \mathcal{E}(i)| \leq \epsilon \Big) \geq 1 - \delta$$

  The probability is taken over the choice of the training set and $1 - \delta$ specifies our confidence in the probabilistic statement.

- Equivalently, we can bound the probability that the empirical error of some classifier in our set deviates more than $\epsilon$ from the expected error:

$$P\Big( \max_{i=1,\ldots,M} |\hat{\mathcal{E}}_n(i) - \mathcal{E}(i)| > \epsilon \Big) \leq \delta$$

## Finite case cont'd

- Let's fix $n$, $M$, and $\epsilon$ and try to find $\delta$ so that

$$P\Big( \max_{i=1,\ldots,M} |\hat{\mathcal{E}}_n(i) - \mathcal{E}(i)| > \epsilon \Big) \leq \delta$$

  still holds. The probability is take over the choice of the training set.

- Let's fix $n$, $M$, and $\epsilon$ and try to find $\delta$ so that

$$P\left(\max_{i=1,\ldots,M}|\hat{\mathcal{E}}_n(i) - \mathcal{E}(i)| > \epsilon\right) \leq \delta$$

still holds. The probability is take over the choice of the training set.

By using the fact that $P(A\,\text{or}\,B) \leq P(A) + P(B)$ we get

$$P\left(\max_i|\hat{\mathcal{E}}_n(i) - \mathcal{E}(i)| > \epsilon\right) \leq \sum_{i=1}^{M} P\left(|\hat{\mathcal{E}}_n(i) - \mathcal{E}(i)| > \epsilon\right)$$

---

- Let's fix $n$, $M$, and $\epsilon$ and try to find $\delta$ so that

$$P\left(\max_{i=1,\ldots,M}|\hat{\mathcal{E}}_n(i) - \mathcal{E}(i)| > \epsilon\right) \leq \delta$$

still holds. The probability is take over the choice of the training set.

By using the fact that $P(A\,\text{or}\,B) \leq P(A) + P(B)$ we get

$$P\left(\max_i|\hat{\mathcal{E}}_n(i) - \mathcal{E}(i)| > \epsilon\right) \leq \sum_{i=1}^{M} P\left(|\hat{\mathcal{E}}_n(i) - \mathcal{E}(i)| > \epsilon\right)$$

$$\leq \sum_{i=1}^{M} 2\exp(-2n\epsilon^2) \quad \text{(Chernoff)}$$

---

- Let's fix $n$, $M$, and $\epsilon$ and try to find $\delta$ so that

$$P\left(\max_{i=1,\ldots,M}|\hat{\mathcal{E}}_n(i) - \mathcal{E}(i)| > \epsilon\right) \leq \delta$$

still holds. The probability is take over the choice of the training set.

By using the fact that $P(A\,\text{or}\,B) \leq P(A) + P(B)$ we get

$$P\left(\max_i|\hat{\mathcal{E}}_n(i) - \mathcal{E}(i)| > \epsilon\right) \leq \sum_{i=1}^{M} P\left(|\hat{\mathcal{E}}_n(i) - \mathcal{E}(i)| > \epsilon\right)$$

$$\leq \sum_{i=1}^{M} 2\exp(-2n\epsilon^2) \quad \text{(Chernoff)}$$

$$= M \cdot 2\exp(-2n\epsilon^2) = \delta$$

---

- We are now able to relate $n$, $M$, $\epsilon$, and $\delta$:

$$M \cdot 2\exp(-2n\epsilon^2) = \delta, \quad \text{or} \quad \epsilon = \sqrt{\frac{\log(M) + \log(2/\delta)}{2n}}$$

- We can restate our result in terms of a bound on the expected error of any classifier in our set.

**Theorem:** With probability at least $1 - \delta$ over the choice of the training set, for all $i = 1, \ldots, M$

$$\mathcal{E}(i) \leq \hat{\mathcal{E}}_n(i) + \epsilon(n, M, \delta)$$

where $\epsilon = \epsilon(n, M, \delta)$ is a "complexity penalty".

---

# Measures of complexity

- Typically the set of classifiers is not a finite nor a countable set (e.g., the set of linear classifiers)

- There are still many ways of trying to capture the "effective" number of classifiers in such a set:

  - degrees of freedom (number of parameters)
  - Vapnik-Chervonenkis (VC) dimension
  - description length
    etc.

---

# VC-dimension: preliminaries

- **A set of classifiers F:** For example, this could be the set of all possible linear classifiers, where $h \in F$ means that

$$h(\mathbf{x}) = \text{sign}\left(w_0 + \mathbf{w}_1^T \mathbf{x}\right)$$

for some values of the parameters $w_0, \mathbf{w}_1$.

## VC-dimension: preliminaries

- **Complexity:** how many different ways can we label $n$ training points $\{\mathbf{x}_1, \ldots, \mathbf{x}_n\}$ with classifiers $h \in F$?

  In other words, how many distinct binary vectors

  $$[h(\mathbf{x}_1)\ h(\mathbf{x}_2)\ \ldots\ h(\mathbf{x}_n)]$$

  do we get by trying out each $h \in F$ in turn?

  $$\begin{array}{cccccc} [ & -1 & 1 & \ldots & 1 & ]\ \ h_1 \\ [ & 1 & -1 & \ldots & 1 & ]\ \ h_2 \\ & & & & & \ldots \end{array}$$
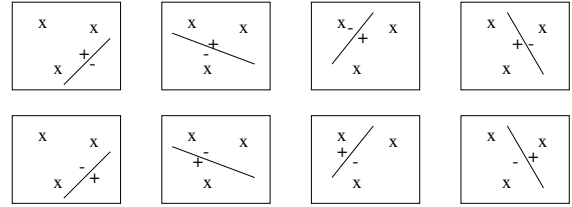
---

## VC-dimension: shattering

- A set of classifiers $F$ *shatters* $n$ points $\{\mathbf{x}_1, \ldots, \mathbf{x}_n\}$ if

  $$[h(\mathbf{x}_1)\ h(\mathbf{x}_2)\ \ldots\ h(\mathbf{x}_n)], \quad h \in F$$

  generates all $2^n$ distinct labelings.

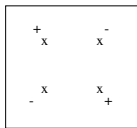- Example: linear decision boundaries shatter (any) 3 points in 2D



  but not any 4 points...

---

## VC-dimension: shattering cont'd

- We cannot shatter any set of 4 points in 2D with linear classifiers. For example, we cannot generate the following XOR-labeling:
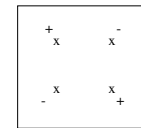


- More generally: the set of all $d$-dimensional linear classifiers can shatter exactly $d + 1$ points

---

## VC-dimension: shattering cont'd

- We cannot shatter any set of 4 points in 2D with linear classifiers. For example, we cannot generate the following XOR-labeling:



- More generally: the set of all $d$-dimensional linear classifiers can shatter exactly $d + 1$ points

- **Definition:** The VC-dimension $d_{VC}$ of a set of classifiers $F$ is the number of points $F$ can shatter

---

## Learning and VC-dimension

- We learn something only after we no longer can shatter the training points (have more than $d_{VC}$ training examples)

  **Rationale:** suppose we have $n$ training examples and labels $(\mathbf{x}_1, y_1), \ldots, (\mathbf{x}_n, y_n)$ and $n < d_{VC}$. Does the training set constrain our prediction for $\mathbf{x}_{n+1}$?

  Because we expect to be able to shatter $n+1$ points ($\leq d_{VC}$) it follows that we can find $h_1, h_2 \in F$, both consistent with training labels, but
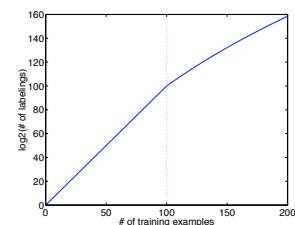
  $$h_1(\mathbf{x}_{n+1}) = 1, \quad h_2(\mathbf{x}_{n+1}) = -1$$

  We therefore cannot determine which label to predict for $\mathbf{x}_{n+1}$.

---

## Learning and VC-dimension

- We learn something only after we no longer can shatter the training points (have more than $d_{VC}$ training examples)



  $$n \leq d_{VC}: \quad \text{\# of labelings} = 2^n$$

  $$n > d_{VC}: \quad \text{\# of labelings} \leq \left(\frac{en}{d_{VC}}\right)^{d_{VC}}$$