



Machine learning

Introduction to **Learning Theory**

Mohammad-Reza A. Dehaqani

dehaqani@ut.ac.ir



Overfitting:

- Huge feature space with kernels, what about **overfitting**?
- Maximizing margin leads to **sparse set** of support vectors
- Some **interesting theory** says that SVMs search for **simple** hypothesis with **large** margin and it is often **robust** to **overfitting**
- We have explored many ways of **learning from data**:
- But...
 - **How good** is our classifier, really?
 - **How much data** do I need to make it “good enough”?

How likely is a bad hypothesis to get m data points right?



- Classification with m i.i.d **data points** and **finite** number
- A learner finds a **hypothesis h** that is **consistent** with training data
 - Gets zero error in training: $\text{error}_{\text{train}}(h) = 0$ (**$\text{error}_D(h)$**)
- The probability that h has **more** than ϵ error in test data ($\text{error}_{\text{true}}(h) \geq \epsilon$ (**$\text{error}_X(h) \geq \epsilon$**)
- **Even** if h makes **zero** errors in **training data**, may make errors in **test**
- If $\text{error}_{\text{true}}(h) \geq \epsilon$; the probability that h gets **one data point right**

$$\leq 1 - \epsilon$$

- Probability that h gets m data points **right**

$$\leq (1 - \epsilon)^m$$



How likely is a learner to pick a bad classifier?

- Usually there are many (say k) **bad** classifiers in **model class**

$$h_1, h_2, \dots, h_k \text{ s.t. } \text{error}_{\text{true}}(h_i) \geq \epsilon \quad i = 1, \dots, k$$

- Probability that learner **picks** a bad classifier
= Probability that some bad classifier **gets 0 training error**

- $\text{Prob}(h_1 \text{ gets 0 training error OR } h_2 \text{ gets 0 training error OR } \dots \text{ OR } h_k \text{ gets 0 training error})$

$$\leq \text{Prob}(h_1 \text{ gets 0 training error}) + \text{Prob}(h_2 \text{ gets 0 training error}) + \dots + \text{Prob}(h_k \text{ gets 0 training error})$$

$$\leq k (1-\epsilon)^m$$

- Probability that learner **picks a bad classifier**

$$\leq k (1-\epsilon)^m \leq \underset{\substack{\uparrow \\ \text{Size of model class}}}{|H|} (1-\epsilon)^m \leq |H| e^{-\epsilon m}$$



- **Theorem [Haussler'88]:** **Model class H** finite, dataset D with m i.i.d. samples, $0 < \epsilon < 1$: for any **learned classifier** h that gets **0 training error**:

$$P(\text{error}_{\text{true}}(h) \geq \epsilon) \leq |H|e^{-m\epsilon} \leq \delta$$

- Equivalently, with probability $\geq 1-\delta$

$$\text{error}_{\text{true}}(h) \leq \epsilon$$

- Important: **PAC bound** holds for all h with 0 training error, but **doesn't** guarantee that algorithm **finds best** h

Using a PAC bound



- Typically, 2 use cases:
 - Pick ϵ and δ , give you m
 - Pick m and δ , give you ϵ

$$|H|e^{-m\epsilon} \leq \delta$$

- Given ϵ and δ , yields **sample complexity**

Number of training data:

$$m \geq \frac{\ln |H| + \ln \frac{1}{\delta}}{\epsilon}$$

- Given m and d , yields **error bound**

Error:

$$\epsilon \geq \frac{\ln |H| + \ln \frac{1}{\delta}}{m}$$

Limitations of Haussler's bound



$$P(\text{error}_{\text{true}}(h) \geq \epsilon) \leq |H|e^{-m\epsilon}$$

- **Consistent** classifier:

- Only consider classifiers with 0 training error
 - $\text{error}_{\text{train}}(h) = 0$

- **Size** of hypothesis space

$$m \geq \frac{\ln |H| + \ln \frac{1}{\delta}}{\epsilon}$$

- what if $|H|$ **too big** or H is **continuous** (e.g. linear classifiers)?

What if our classifier does not have zero error on the training data?



- A learner with zero training errors may make mistakes in test set
- What about a learner with $\text{error}_{\text{train}}(h) \neq 0$ in training set? ($\text{error}_{\text{train}}(\text{error}_D(h))$ relates $\text{error}_{\text{true}}(\text{error}_X(h))$)
- The error of a classifier is a **Bernoulli** random variable
- like estimating the **parameter of a coin!**

$$\text{error}_X(h) = \text{error}_{\text{true}}(h) := P(h(X) \neq Y) \equiv P(H=1) =: \theta$$

$$\text{error}_D(h) = \text{error}_{\text{train}}(h) := \frac{1}{m} \sum_i \mathbf{1}_{h(X_i) \neq Y_i} \equiv \frac{1}{m} \sum_i x_i$$



Hoeffding's bound

- Consider m i.i.d. **Bernoulli random** variable $\mathbf{x}_1, \dots, \mathbf{x}_m$, ($x_i \in \{0,1\}$;
flip a coin with **parameter** θ) . For $0 < \epsilon < 1$:

$$P \left(\left| \theta - \frac{1}{m} \sum_i x_i \right| \geq \epsilon \right) \leq 2e^{-2m\epsilon^2}$$

- For **a** single classifier h_i

$$P (|\text{error}_{true}(h_i) - \text{error}_{train}(h_i)| \geq \epsilon) \leq 2e^{-2m\epsilon^2}$$

- For **any** learned classifier $h \in H$ (we are comparing $|H|$ classifiers)

Union bound

$$P (|\text{error}_{true}(h) - \text{error}_{train}(h)| \geq \epsilon) \leq 2|H|e^{-2m\epsilon^2} \leq \delta$$

$$\textbf{Probability of mistake} \leq 2|H|e^{-2m\epsilon^2}$$

Recall tail bounds in probability

Hoeffding's inequality (1963)



- The bounds when there are the range of the variables, but not the variances.

$$\left. \begin{array}{l} X_1, \dots, X_n \text{ independent} \\ X_i \in [a_i, b_i] \\ \varepsilon > 0 \end{array} \right\} \Rightarrow$$
$$\Rightarrow \left\{ \begin{array}{l} \mathbb{P}(|\frac{1}{n} \sum_{i=1}^n (X_i - \mathbb{E}X_i)| > \varepsilon) \leq 2 \exp \left(\frac{-2n\varepsilon^2}{\frac{1}{n} \sum_{i=1}^n (b_i - a_i)^2} \right) \\ \text{two-sided} \\ \\ \mathbb{P}(\frac{1}{n} \sum_{i=1}^n (X_i - \mathbb{E}X_i) > \varepsilon) \leq \exp \left(\frac{-2n\varepsilon^2}{\frac{1}{n} \sum_{i=1}^n (b_i - a_i)^2} \right) \\ \text{one-sided} \end{array} \right.$$

PAC bound and Bias-Variance tradeoff



test set mistake

$$P(|\text{error}_{\text{true}}(h) - \text{error}_{\text{train}}(h)| \geq \epsilon) \leq 2|H|e^{-2m\epsilon^2} \leq \delta$$

- After moving some terms around, at least with **probability 1-δ**:

$$\text{error}_{\text{true}}(h) \leq \text{error}_{\text{train}}(h) + \sqrt{\frac{\ln |H| + \ln \frac{2}{\delta}}{2m}}$$

test set mistake
Bias (training mistake)
variance

- Fixed m**

Model Class		
Complex ($ H $ is Big)	Small	Large
Simple ($ H $ is small)	Large	Small

The size of the model class



$$2|H|e^{-2m\epsilon^2} \leq \delta$$

- Sample complexity

$$m \geq \frac{1}{2\epsilon^2} \left(\ln |H| + \ln \frac{2}{\delta} \right)$$

- How large is the model class?

$|H|$ is large \Rightarrow need many training examples

What about continuous hypothesis spaces



$$\text{error}_{\text{true}}(h) \leq \text{error}_{\text{train}}(h) + \sqrt{\frac{\ln |H| + \ln \frac{2}{\delta}}{2m}}$$

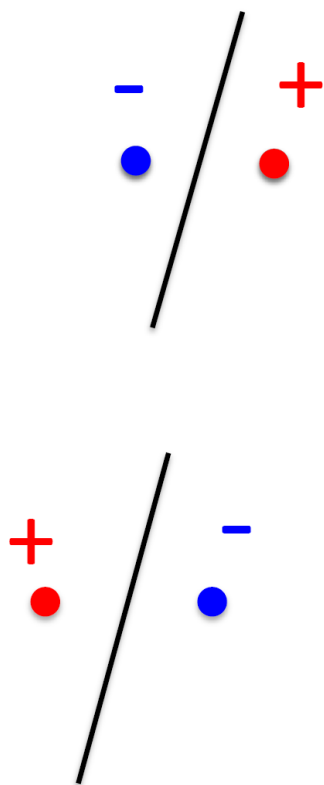
- Continuous model class (e.g. linear classifiers):

$$|H| = \infty$$

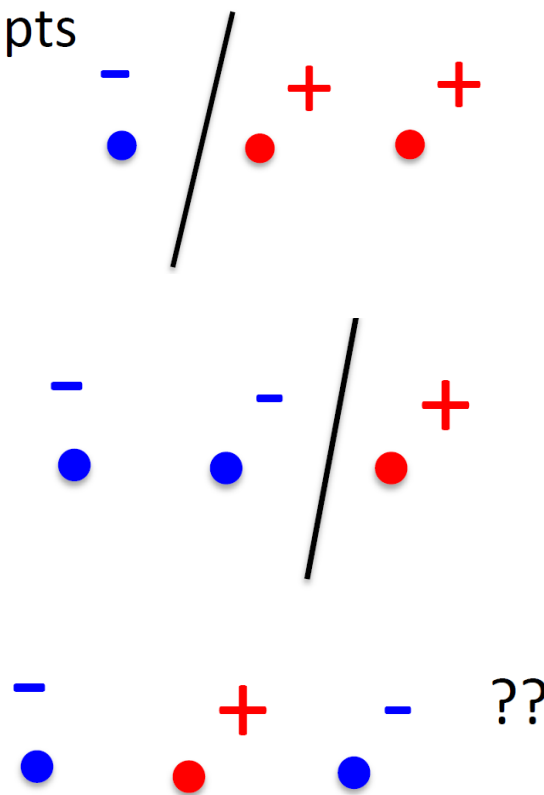
- Complexity of **model class** can depends on **maximum number of points** that can be **classified** exactly (and not **necessarily its size**)

How many points can a linear boundary classify exactly? (1-D)

2 pts



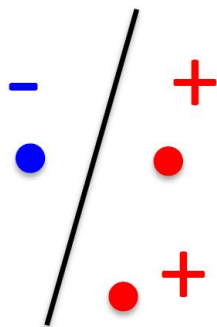
3 pts



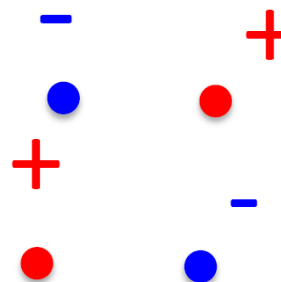
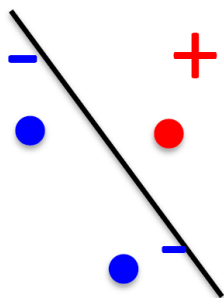
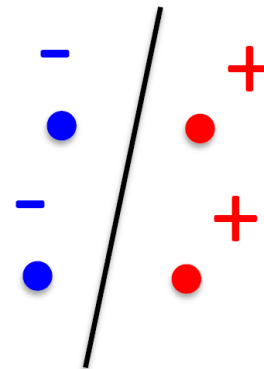
- There exists **placement** s.t. all labelings can be classified
- Complexity of **model class** is 2

How many points can a linear boundary classify exactly? (2-D)

3 pts



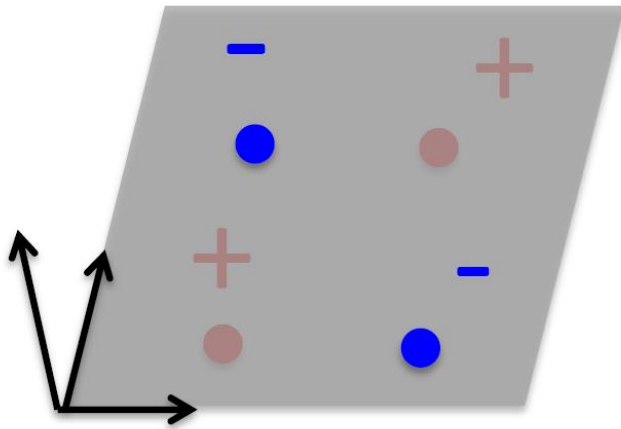
4 pts



??

- There exists **placement** s.t. all labelings can be classified
- Complexity of **model class** is 3

How many points can a linear boundary classify exactly? (d-D)



- **Number** of training points that can be classified exactly is **VC dimension**

- How many parameters in linear Classifier in d-Dimensions?

$$w_0 + \sum_{i=1}^d w_i x_i$$

- d+1 parameters need d+1 **constraints**
- d+1 pts

PAC bound using VC dimension



- Measures relevant size of hypothesis space **using VC dimension**
- Bound for **infinite** dimension hypothesis spaces:
- In statistical learning, the **Vapnik-Chervonenkis (VC) dimension** is a popular measure of the **capacity** of a classifier.
- The **VC dimension** can predict a **probabilistic** upper bound on the **generalization error** of a classifier.

w.p. $\geq 1-\delta$

$$\text{error}_{\text{true}}(h) \leq \underbrace{\text{error}_{\text{train}}(h)}_{\text{Bias (training mistake)}} + \underbrace{\sqrt{\frac{VC(H) \left(\ln \frac{2m}{VC(H)} + 1 \right) + \ln \frac{4}{\delta}}{m}}}_{\text{variance}}$$

linear classifiers

2D	Small $VC(H)$
10,000 D	Big $VC(H)$

large
small

small
large

New game: Picking right $VC(H)$

VC (Vapnik-Chervonenkis) dimension



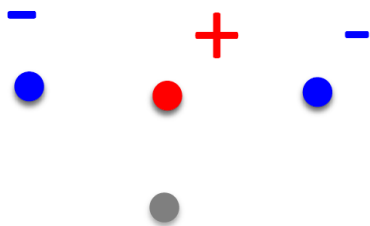
- Definition: VC dimension of a **hypothesis space H** is the **maximum number of points** such that **there exists a hypothesis** in H that is **consistent** with (can **correctly** classify) **any labeling** of the points.
 - You pick set of points
 - Adversary assigns labels
 - You find a hypothesis in H consistent with the labels
- If $VC(H) = k$, then for all $k+1$ points, there exists a labeling that **cannot be shattered** (can't find a hypothesis in H consistent with it)
- Definition: a set of instances S is **shattered** by hypothesis space H if and only if for **every dichotomy** (+/- labeling) of S there exists **some hypothesis** in H **consistent** with this dichotomy (labeling)
- $VC(H)$, of hypothesis space H defined over instances space X is the size of the **largest** finite subset of X **shattered** by H .
- If the **arbitrarily large** finite sets of X can be shattered by H , then $VC(H) \equiv \infty$



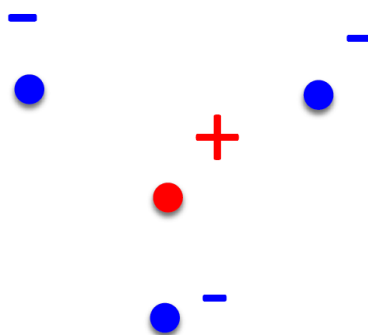
VC dim. example – What can't we shatter

- What's the VC dim. of decision in 2D?
- If $VC(H) = 3$, then **for all placements of** 4 pts, there exists a labeling that can't be shattered
- **Linear classifiers:**
 - $VC(H) = d+1$, for d features plus constant term (3 if $d=2$)

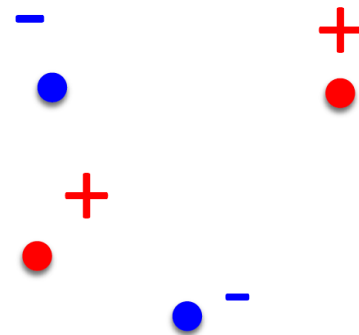
3 collinear



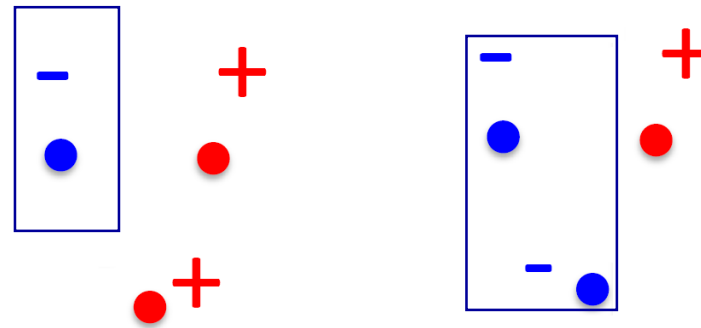
1 in convex hull
of other 3



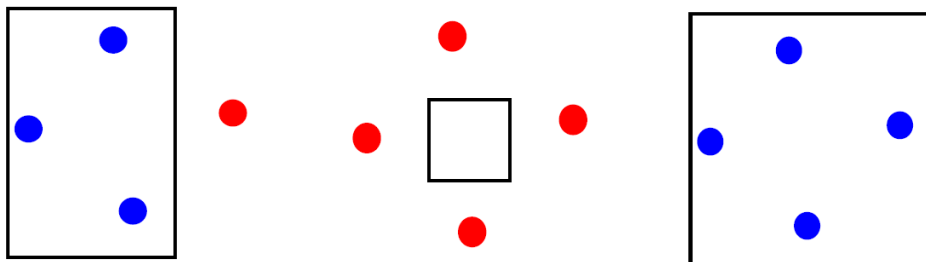
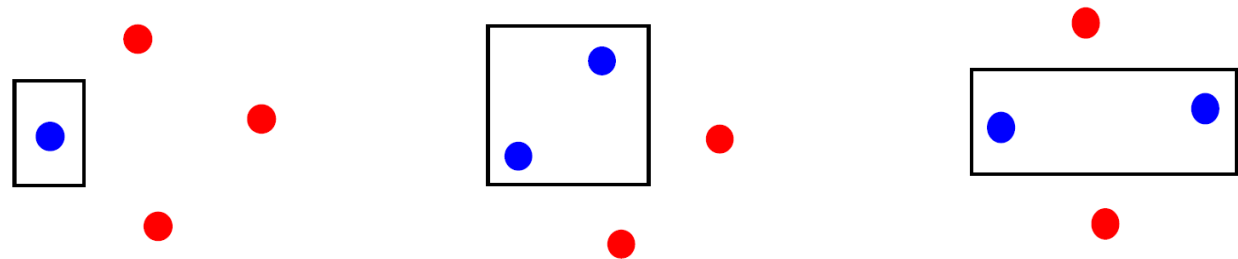
quadrilateral



VC dim. of axis parallel rectangles in 2D

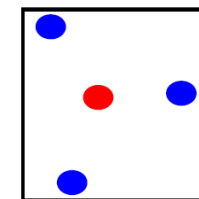


$$VC(H) \geq 3$$



$$VC(H) \geq 4$$

Some placement of 4 pts can't be **shattered**

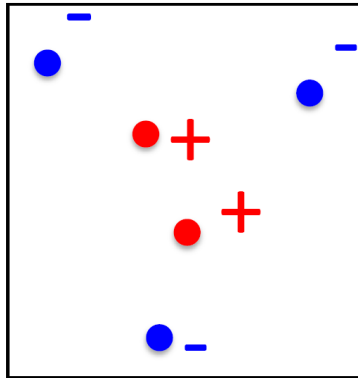


VC dim. of axis parallel rectangles in 2D is 4

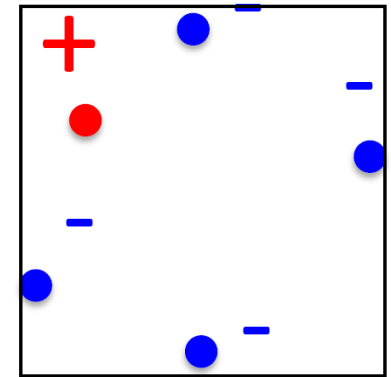


- If $VC(H) = 4$, then for **all placements** of 5 pts, **there exists** a **labeling** that can't be shattered

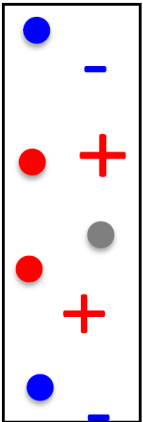
2 in convex hull of other 3



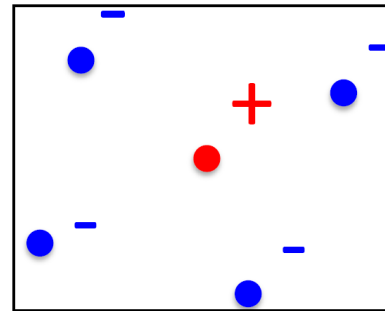
pentagon



4 collinear



1 in convex hull of other 4



Examples of VC dimension



- **Linear classifiers:**
 - $VC(H) = d+1$, for d features plus constant term
- **Axis parallel rectangles:**
 - $VC(H) = 2d$ (4 if $d=2$)
- **Nearest Neighbor:**
 - $VC(H) = \infty$
- $VC(H) \leq \log_2 |H|$ (So VC bound is tighter)
 - Given $|H|$ hypothesis can hope to shatter max $m = \log_2 |H|$ points
 - 2^m labelings $\Rightarrow |H| \geq 2^m$

PAC bound for SVMs



- SVM uses **a linear classifier**
- For d features, $\mathbf{VC(H)} = d+1$:

$$\text{error}_{\mathcal{X}}(h) \leq \text{error}_D(h) + \sqrt{\frac{(d+1) \left(\ln \frac{2m}{d+1} + 1 \right) + \ln \frac{4}{\delta}}{m}}$$

- **Polynomials** kernel:

- Number of features grows really fast \Rightarrow Bad bound

- Number of terms = $\binom{p+n-1}{p} = \frac{(p+n-1)!}{p!(n-1)!}$

n -input features

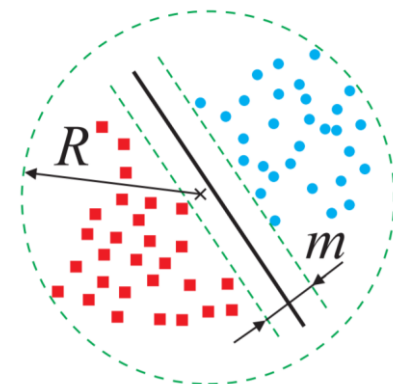
p -degree of polynomial

- **Doesn't take margin into account**
- **Gaussian kernels** can classify any set of points exactly ($\mathbf{VC(H)} = \infty$);
- Suggests Gaussian kernels and **deep nets (due to large number of parameters)** are really BAD!! But contradicts practice!

Margin-based VC dimension



- H : Class of linear classifiers: $\mathbf{w} \cdot \Phi(\mathbf{x})$ ($b=0$)
- Canonical form: $\min_j |\mathbf{w} \cdot \Phi(\mathbf{x}_j)| = 1$
- $VC(H) = R^2 \mathbf{w} \cdot \mathbf{w}$
 - Doesn't depend on number of features!
 - $R^2 = \max_j \Phi(\mathbf{x}_j) \cdot \Phi(\mathbf{x}_j)$ –magnitude of data (R be the radius of the smallest inclosing ball (sphere) of the data)
 - R^2 is **bounded** even for Gaussian kernels \rightarrow bounded VC dimension
- **Large** margin \Rightarrow **low** $\mathbf{w} \cdot \mathbf{w}$ \Rightarrow low VC dimension \Rightarrow low variance error
- SVMs minimize $\mathbf{w} \cdot \mathbf{w}$
- We require bound over **infinite number of possible** VC dimensions...



Structural risk minimization theorem



$$\text{error}_{\mathcal{X}}(h) \leq \text{Bias (empirical risk)} \quad \text{error}_D^\gamma(h) + \text{Variance (VC confidence)} \quad C \sqrt{\frac{\frac{R^2}{\gamma^2} \ln m + \ln \frac{1}{\delta}}{m}}$$

- γ is **margin**
- For a **family** of hyperplanes (with respective growing VC dimensions) with margin $\gamma > 0$
- SVMs **maximize** margin γ + **hinge loss**
- Optimize **tradeoff** training error (bias) versus margin γ (variance)

