

Machine Learning

Lecture Slides for

INTRODUCTION TO

Machine Learning

ETHEM ALPAYDIN © The MIT Press, 2010

In preparation of these slides, I have benefited from slides prepared by:

Analysis),

THEM ALPATERN

alpaydin@boun.edu.tr D. Bouchaffra and V. Murino (Pattern Classification and Scene Acaleria)

R. Gutierrez-Osuna (Texas A&M)

A. Moore (CMU)

CHAPTER 2:

Supervised Learning

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 X







Hypothesis class \mathcal{H}



S, G, and the Version Space





• Choose *h* with largest margin



VC Dimension and PAC Learning

Tools to analyze expected (test/generalization) error of hypothesis classes (i.e. classifiers)

- VC (Vapnik Chervonenkis) Dimension:
 - Given a hypothesis class (rectangles, circles, lines, neural networks) how many instances can it classify without any error.
 - Does not consider the input distribution
- PAC (Probably Approximately Correct) Learning
 - How many instances are needed to achieve a certain error with a certain probability?

VC Dimension

• N points can be labeled in 2^N ways as +/-

N=3 for this table

	R1	R2	R3	R4	R5	R6	R7	R8
x1	0	0	0	1	0	1	1	1
x2	0	0	1	0	1	0	1	1
x3	0	1	0	0	1	1	0	1

 \mathcal{H} shatters N if there exists $h \in \mathcal{H}$ consistent for any of these: $VC(\mathcal{H}) = N$



An axis-aligned rectangle shatters 4 points only ! Choose points in such a way that they can be shattered. e.g. Do not put all of them at the same spot.

Lecture Notes for E Alpaydın 2010 Introduction to Machine Learning 2e © The MIT Press (V1.0)

 X_{I}

VC Dimension Examples

We use g to denote hypothesis (like h) Question: Can the following g (line passing through origin) shatter the following points?

 $g(x,\theta) = sign(x, \theta) = sign(x_1\theta_1 + x_2\theta_2)$

Answer: No problem. There are four training sets to consider • +



Given machine g, the VC-dimension v is
The maximum number of points that can
be arranged so that g shatter them.
Example: For 2-d inputs, what' s VC-dim of g(x,θ,b) = sign(θ.x+b)?
Well, can g shatter these three points?

• • Yes, of course.

0

All -ve or all +ve is trivial

One +ve can be picked off by a line

One -ve can be picked off too.

Given machine **g**, the VC-dimension v is

The maximum number of points that can be arranged so that *g* shatter them.

Example: For 2-d inputs, what' s VC-dim of $g(x,\theta,b) = sign(\theta,x+b)$? Well, can we find four points that *g* can shatter?

Given machine **g**, the VC-dimension v is

The maximum number of points that can be arranged so that *g* shatter them.

Example: For 2-d inputs, what' s VC-dim of $g(x,\theta,b) = sign(\theta,x+b)$? Well, can g shatter these four points?



Can always draw six lines between pairs of four points.

Given machine **g**, the VC-dimension v is

The maximum number of points that can be arranged so that *g* shatter them.

Example: For 2-d inputs, what' s VC-dim of $g(x,\theta,b) = sign(\theta,x+b)$? Well, can g shatter these four points?



Can always draw six lines between pairs of four points.

Two of those lines will cross.

Given machine **g**, the VC-dimension v is

The maximum number of points that can be arranged so that *g* shatter them.

Example: For 2-d inputs, what' s VC-dim of $g(x, \theta, b) = sign(\theta.x+b)$? Well, can we find four points that *g* can shatter?



Can always draw six lines between pairs of four points.

Two of those lines will cross.

If we put points linked by the crossing lines in the same class they can't be linearly separated

So a line can shatter 3 points but not 4

So VC-dim of Line Machine is 3

If inputs are m dimensional, v=m+1(see A. Moore notes).

Why VC Dimension

 $R^{emp}(\theta) = \text{TRAINERR}(\theta) = \frac{\text{Fra}}{\text{Se}}$

Fraction Training Set misclassified

 $R(\theta) = \text{TESTERR}(\theta) =$

Probability of Misclassification

•Given some machine **g**, let *v* be its VC dimension.

v is a measure of g's power (*v* does not depend on the choice of training set)
Vapnik showed that with probability 1-η

$$\text{TESTERR}(\theta) \leq \text{TRAINERR}(\theta) + \sqrt{\frac{\nu(\log(2N/\nu) + 1) - \log(\eta/4)}{N}}$$

This bound is usually too big, so it may not be useful.



Probably Approximately Correct (PAC) Learning

How many training examples N should we have, such that with probability at least 1 – δ, h has error at most ε ?
 (Blumer et al., 1989)

- Each strip is at most ε/4
- Pr that we miss a strip $1 \epsilon/4$
- Pr that N instances miss a strip $(1 \varepsilon/4)^N$
- Pr that N instances miss 4 strips $4(1 \varepsilon/4)^N$
- $4(1 \epsilon/4)^N \le \delta$ and $(1 x) \le \exp(-x)$
- $4\exp(-\epsilon N/4) \le \delta$ and $N \ge (4/\epsilon)\log(4/\delta)$



Noise and Model Complexity

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)



Lecture Notes for E Alpaydın 2010 Introduction to Machine Learning 2e © The MIT Press (V1.0)

Multiple Classes, C_i i=1,...,K



21

Regression



Lecture Notes for E Alpaydin 2010 Introduction to Machine Learning 2e © The MIT Press (V1.0)

Model Selection & Generalization

- Learning is an ill-posed problem, i.e. data is not sufficient to find a unique solution
- The need for inductive bias, assumptions about ${\mathcal H}$
- Generalization: How well a model performs on new data
- Overfitting: \mathcal{H} more complex than C or f
- Underfitting: \mathcal{H} less complex than C or f

Model Selection & Generalization

- Learning is an ill-posed problem, i.e. data is not sufficient to find a unique solution
- The need for inductive bias, assumptions about ${\mathcal H}$
- Generalization: How well a model performs on new data
- Overfitting: \mathcal{H} more complex than C or f
- Underfitting: \mathcal{H} less complex than C or f

Data not sufficient to find a unique solution

Set of assumptions we make to make learning possible.

Triple Trade-Off

- There is a trade-off between three factors (Dietterich, 2003):
 - 1. Complexity of \mathcal{H} , c (\mathcal{H}),
 - 2. Training set size, N,
 - 3. Generalization error, E, on new data
- $\Box \quad \text{As } N \uparrow, E \downarrow$
- □ As $c(\mathcal{H})\uparrow$, first $E\downarrow$ and then $E\uparrow$

Cross-Validation

- To estimate generalization error, we need data unseen during training. We split the data as
 - Training set (50%)
 - Validation set (25%)
 - Test (publication) set (25%)

Resampling (e.g. bootstrapping) when there is few data

Cross Validation



Dimensions of a Supervised Learner

1. Model: $g(\mathbf{x} | \theta)$

2. Loss function:
$$E(\theta \mid X) = \sum_{t} L(r^{t}, g(\mathbf{x}^{t} \mid \theta))$$

3. Optimization procedure:

$$\theta^* = \arg\min_{\theta} E(\theta \,|\, \mathcal{X})$$