

An Intro to Machine Learning

Avrim Blum
Lecture #12

Plan for today

- Machine Learning intro: models and basic issues
- How much data do I need to see to be confident in generalizations I make from it?
- Connections of this to notion of Occam's razor
- A cool idea: "shatter coefficients", VC-dimension, and a very nice probabilistic argument.

Plan for Monday

- An interesting algorithm for online decision making. Problem of "combining expert advice"
- Algorithms for online decision making from very limited feedback. The "multi-armed bandit problem"

Machine learning can be used to...

- recognize speech,
- identify patterns in data,
- steer a car,
- play games,
- adapt programs to users,
- improve web search, ...

From a scientific perspective: can we develop models to understand learning as a computational problem, and what types of guarantees might we hope to achieve?

A typical setting

- Imagine you want a computer program to help filter which email messages are spam and which are important.
- Might represent each message by n features. (e.g., return address, keywords, spelling, etc.)
- Take sample S of data, labeled according to whether they were/weren't spam.
- Goal of algorithm is to use data seen so far produce good prediction rule (a "hypothesis") $h(x)$ for future data.

The concept learning setting

E.g., money pills Mr. bad spelling known-sender | spam?

The concept learning setting

E.g.,

	money	pills	Mr.	bad spelling	known-sender	spam?
a positive example	Y	N	Y	Y	N	Y
	N	N	N	Y	Y	N
	N	Y	N	N	N	Y
a negative example	Y	N	N	N	Y	N
	N	N	Y	N	Y	N
	Y	N	N	Y	N	Y

Given data, some reasonable rules might be:

- Predict SPAM if \neg known AND (money OR pills)
- Predict SPAM if money + pills - known > 0 .
- ...

Big questions

- (A) How might we automatically generate rules that do well on observed data?
[algorithm design]
- (B) What kind of confidence do we have that they will do well in the future?
[confidence bound / sample complexity]

for a given learning alg, how much data do we need...

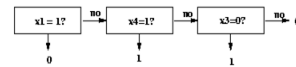
Natural formalization (PAC)

Email msg Spam or not?

- We are given sample $S = \{(x,y)\}$.
- View labels y as being produced by some target function f .
- Alg does optimization over S to produce some hypothesis (prediction rule) h .
- Assume S is a random sample from some probability distribution D . Goal is for h to do well on new examples also from D .

$$\text{I.e., } \Pr_D[h(x) \neq f(x)] < \epsilon.$$

Example of analysis: Decision Lists



Say we suspect there might be a good prediction rule of this form.

1. Design an efficient algorithm A that will find a consistent DL if one exists.
2. Show that if sample S is of reasonable size, $\Pr[\text{exists consistent DL } h \text{ with } \text{err}(h) > \epsilon] < \delta$.
3. This means that A is a good algorithm to use if f is, in fact, a DL.

(a bit of a toy example since would want to extend to "mostly consistent" DL)

How can we find a consistent DL?

x_1	x_2	x_3	x_4	x_5	label
1	0	0	1	1	+
0	1	1	0	0	-
1	1	1	0	0	+
0	0	0	1	0	-
1	1	0	1	1	+
1	0	0	0	1	-

- if ($x_1=0$) then -, else
- if ($x_2=1$) then +, else
- if ($x_4=1$) then +, else -

Decision List algorithm

- Start with empty list.
- Find if-then rule consistent with data.
(and satisfied by at least one example)
- Put rule at bottom of list so far, and cross off examples covered. Repeat until no examples remain.

If this fails, then:

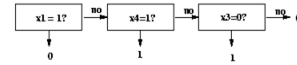
- No rule consistent with remaining data.
- So no DL consistent with remaining data.
- So, no DL consistent with original data.

OK, fine. Now why should we expect it to do well on future data?

Confidence/sample-complexity

- Consider some DL h with $\text{err}(h) > \epsilon$, that we're worried might fool us.
 - Chance that h survives $|S|$ examples is at most $(1-\epsilon)^{|S|}$.
 - Let $|H|$ = number of DLs over n Boolean features. $|H| < (4n+2)^n$. (really crude bound)
- So, $\Pr[\text{some DL } h \text{ with } \text{err}(h) > \epsilon \text{ is consistent}] < |H|(1-\epsilon)^{|S|}$.
- This is < 0.01 for $|S| > (1/\epsilon)[\ln(|H|) + \ln(100)]$ or about $(1/\epsilon)[n \ln n + \ln(100)]$

Example of analysis: Decision Lists



Say we suspect there might be a good prediction rule of this form.

1. **DONE** Design an efficient algorithm A that will find a consistent DL if one exists.
2. **DONE** Show that if $|S|$ is of reasonable size, then $\Pr[\text{exists consistent DL } h \text{ with } \text{err}(h) > \epsilon] < \delta$.
3. So, if f is in fact a DL, then whp A 's hypothesis will be approximately correct. "PAC model"

Confidence/sample-complexity

- What's great is there was nothing special about DLs in our argument.
- All we said was: "if there are not *too* many rules to choose from, then it's unlikely one will have fooled us just by chance."
- And in particular, the number of examples needs to only be proportional to $\log(|H|)$.
(the "log" is important here)

Occam's razor

William of Occam (~1320 AD):

"entities should not be multiplied unnecessarily" (in Latin)

Which we interpret as: "in general, prefer simpler explanations".

Why? Is this a good policy? What if we have different notions of what's simpler?

Occam's razor (contd)

A computer-science-ish way of looking at it:

- Say "simple" = "short description".
- At most 2^s explanations can be $< s$ bits long.
- So, if the number of examples satisfies:

Think of as 10x #bits to write down h .

$$m > (1/\epsilon)[s \ln(2) + \ln(100)]$$

Then it's unlikely a bad simple explanation will fool you just by chance.

Occam's razor (contd)?

Nice interpretation:

- Even if we have different notions of what's simpler (e.g., different representation languages), we can both use Occam's razor.
- Of course, there's no guarantee there will be a short explanation for the data. That depends on your representation.

Extensions

We said: if $|S| \geq (1/\epsilon)[\ln(|H|) + \ln(1/\delta)]$, then with probability $\geq 1-\delta$, all $h \in H$ with $\text{err}_D(h) \geq \epsilon$ have $\text{err}_S(h) > 0$.

What if no perfect rule, and best we find is rule with error (say) 10% on training set? What can we say?

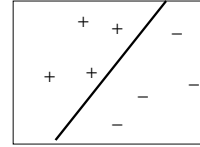
Thm: If $|S| \geq (1/(2\epsilon^2))[\ln(|H|) + \ln(2/\delta)]$, then with prob $\geq 1-\delta$, all $h \in H$ have $|\text{err}_D(h) - \text{err}_S(h)| < \epsilon$.

Proof: apply Hoeffding bounds.

- Chance of failure at most $2|H|e^{-2|S|\epsilon^2}$.
- Set to δ and solve.

One more extension

- What about something like the class H of linear separators? What is $|H|$?

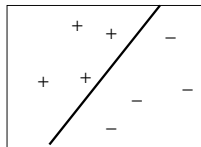


- There are infinitely many linear separators, but not that many really different ones.
- Union bound is too weak.

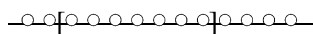
A cool idea: shatter coefficient

- Let $H[S]$ be the number of ways of splitting set S using functions in H .

- Let $H[m] = \max_{|S|=m} H[S]$.



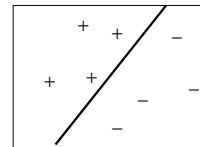
- E.g., linear separators in \mathbb{R}^d : $H[m] = O(m^d)$.
- E.g., intervals on a line: $H[m] = O(m^2)$.



A cool idea: shatter coefficient

- Let $H[S]$ be the number of ways of splitting set S using functions in H .

- Let $H[m] = \max_{|S|=m} H[S]$.



- E.g., linear separators in \mathbb{R}^d : $H[m] = O(m^d)$.

Thm: if $m = |S| \geq (2/\epsilon)[\lg(2H[2m]) + \lg(1/\delta)]$, then with probability $\geq 1-\delta$, all $h \in H$ with $\text{err}_D(h) \geq \epsilon$ have $\text{err}_S(h) > 0$.

A cool idea: shatter coefficient

Thm: if $m \geq (2/\epsilon)[\lg(2H[2m]) + \lg(1/\delta)]$, then with probability $\geq 1-\delta$, all $h \in H$ with $\text{err}_D(h) \geq \epsilon$ have $\text{err}_S(h) > 0$.

Note 1: For linear separators in \mathbb{R}^d , $H[2m] = O(m^d)$, so bound is $O(1/\epsilon)[d \lg(1/\epsilon) + \lg(1/\delta)]$

Note 2: $\text{VC-dimension}(H) = \max$ value m such that $H[m] = 2^m$

Sauer's lemma: $H[m] = O(m^{\text{VCdim}(H)})$.

A cool idea: shatter coefficient

Thm: if $m \geq (2/\epsilon)[\lg(2H[2m]) + \lg(1/\delta)]$, then with probability $\geq 1-\delta$, all $h \in H$ with $\text{err}_D(h) \geq \epsilon$ have $\text{err}_S(h) > 0$.

Proof of Thm:

- Consider drawing 2 sets S, S' of m examples each.
- Let A be the event: exists $h \in H$ with $\text{err}_D(h) \geq \epsilon$ and $\text{err}_{S'}(h) = 0$.
- Let B be the event: exists $h \in H$ with $\text{err}_S(h) \geq \epsilon/2$ and $\text{err}_{S'}(h) = 0$.
- Claim 1: $\Pr[A]/2 \leq \Pr[B]$ (because $\Pr[B|A] \geq \frac{1}{2}$)
- So, just need to show $\Pr[B]$ is low.

A cool idea: shatter coefficient

Thm: if $m \geq (2/\epsilon)[\lg(2H[2m]) + \lg(1/\delta)]$, then with probability $\geq 1-\delta$, all $h \in H$ with $\text{err}_D(h) \geq \epsilon$ have $\text{err}_S(h) > 0$.

Proof cont'd:

- Consider drawing 2 sets S, S' of m examples each.
- Let B be the event: exists $h \in H$ with $\text{err}_{S'}(h) \geq \epsilon/2$ and $\text{err}_S(h) = 0$. Suffices to show $\Pr[B]$ is low.
- Now, define T, T' as follows:
 - For $i=1$ to m , flip a fair coin:
 - If heads, put i th element of S into T and i th element of S' into T' .
 - If tails, do it other way around.

A cool idea: shatter coefficient

Thm: if $m \geq (2/\epsilon)[\lg(2H[2m]) + \lg(1/\delta)]$, then with probability $\geq 1-\delta$, all $h \in H$ with $\text{err}_D(h) \geq \epsilon$ have $\text{err}_{S'}(h) > 0$.

Proof cont'd:

- Consider drawing 2 sets S, S' of m examples each.
- Let C be the event: exists $h \in H$ with $\text{err}_{T'}(h) \geq \epsilon/2$ and $\text{err}_T(h) = 0$. Suffices to show $\Pr[C]$ is low.
- Now, define T, T' as follows:
 - For $i=1$ to m , flip a fair coin:
 - If heads, put i th element of S into T and i th element of S' into T' .
 - If tails, do it other way around.

A cool idea: shatter coefficient

Thm: if $m \geq (2/\epsilon)[\lg(2H[2m]) + \lg(1/\delta)]$, then with probability $\geq 1-\delta$, all $h \in H$ with $\text{err}_D(h) \geq \epsilon$ have $\text{err}_S(h) > 0$.

Proof cont'd:

- Will show that for all S, S' , $\Pr_{\text{swap}}[C]$ is low.
- Let C be the event: exists $h \in H$ with $\text{err}_{T'}(h) \geq \epsilon/2$ and $\text{err}_T(h) = 0$. Suffices to show $\Pr[C]$ is low.
- Now, define T, T' as follows:
 - For $i=1$ to m , flip a fair coin:
 - If heads, put i th element of S into T and i th element of S' into T' .
 - If tails, do it other way around.

A cool idea: shatter coefficient

Thm: if $m \geq (2/\epsilon)[\lg(2H[2m]) + \lg(1/\delta)]$, then with probability $\geq 1-\delta$, all $h \in H$ with $\text{err}_D(h) \geq \epsilon$ have $\text{err}_{S'}(h) > 0$.

Proof cont'd:

- Will show that for all S, S' , $\Pr_{\text{swap}}[C]$ is low.
- Let C be the event: exists $h \in H$ with $\text{err}_{T'}(h) \geq \epsilon/2$ and $\text{err}_T(h) = 0$. Suffices to show $\Pr[C]$ is low.
- Fix some splitting h of $S \cup S'$ (at most $H[2m]$)
- If for any i , h makes mistake on i th element of both S and S' , then $\Pr[C_h] = 0$. Also, if h makes fewer than $\epsilon m/2$ mistakes on $S \cup S'$, then $\Pr[C_h] = 0$.
- Else, $\Pr[C_h] \leq 2^{-\epsilon m/2}$. Set $H[2m] \times 2^{-\epsilon m/2} = \delta/2$. Done!

Online learning

- What if we don't want to make assumption that data is coming from some fixed distribution? Or any assumptions on data?
- Can no longer talk about past performance predicting future results.
- Can we do better than just guessing at all??

Idea: recursively select the best hypothesis as best performing hypothesis as well as best performing hypothesis.



Using "expert" advice

- Say we want to predict the stock market.
- We solicit n "experts" for their advice. (Will the market go up or down?)
- We then want to use their advice somehow to make our prediction. E.g.,

Expt 1	Expt 2	Expt 3	neighbor's dog	truth
down	up	up	up	up
down	up	up	down	down
...

Basic question: Is there a strategy that allows us to do nearly as well as best of these in hindsight?

["expert" = someone with an opinion. Not necessarily someone who knows anything.]

Simpler question

- We have n "experts".
- One of these is perfect (never makes a mistake). We just don't know which one.
- Can we find a strategy that makes no more than lg(n) mistakes?

Answer: sure. Just take majority vote over all experts that have been correct so far.

➤ Each mistake cuts # available by factor of 2.

➤ Note: this means ok for n to be very large.

What if no expert is perfect?

Intuition: Making a mistake doesn't completely disqualify an expert. So, instead of crossing off, just lower its weight.

Weighted Majority Alg:

- Start with all experts having weight 1.
- Predict based on weighted majority vote.
- Penalize mistakes by cutting weight in half.

	prediction		correct
weights	1	1	1
predictions	Y	Y	N
weights	1	1	.5
predictions	Y	N	Y
weights	1	.5	.5

Analysis: do nearly as well as best expert in hindsight

- M = # mistakes we've made so far.
- m = # mistakes best expert has made so far.
- W = total weight (starts at n).
- After each mistake, W drops by at least 25%. So, after M mistakes, W is at most $n(3/4)^M$.
- Weight of best expert is $(1/2)^m$. So,

$$\begin{aligned} (1/2)^m &\leq n(3/4)^M \\ (4/3)^M &\leq n2^m \\ M &\leq 2.4(m + \lg n) \end{aligned}$$

So, if m is small, then M is pretty small too.

Randomized Weighted Majority

$2.4(m + \lg n)$ not so good if the best expert makes a mistake 20% of the time. Can we do better? Yes.

- Instead of taking majority vote, use weights as probabilities. (e.g., if 70% on up, 30% on down, then pick 70:30) Idea: smooth out the worst case.
- Also, generalize $\frac{1}{2}$ to $1 - \epsilon$.

Solves to: $M \leq \frac{-m \ln(1 - \epsilon) + \ln(n)}{\epsilon} \approx (1 + \epsilon/2)m + \frac{1}{\epsilon} \ln(n)$

M = expected #mistakes	$M \leq 1.39m + 2 \ln n \leftarrow \epsilon = 1/2$
	$M \leq 1.15m + 4 \ln n \leftarrow \epsilon = 1/4$
	$M \leq 1.07m + 8 \ln n \leftarrow \epsilon = 1/8$

Analysis

- Say at time t we have fraction F_t of weight on experts that made mistake.
- So, we have probability F_t of making a mistake, and we remove an ϵF_t fraction of the total weight.
 - $W_{final} = n(1 - \epsilon F_1)(1 - \epsilon F_2) \dots$
 - $\ln(W_{final}) = \ln(n) + \sum_t [\ln(1 - \epsilon F_t)] \leq \ln(n) - \epsilon \sum_t F_t$
(using $\ln(1-x) \leq -x$)
 $= \ln(n) - \epsilon M$. ($\sum F_t = E[\# \text{ mistakes}]$)
- If best expert makes m mistakes, then $\ln(W_{final}) > \ln((1-\epsilon)^m)$.
- Now solve: $\ln(n) - \epsilon M > m \ln(1-\epsilon)$.

$$M \leq \frac{-m \ln(1 - \epsilon) + \ln(n)}{\epsilon} \approx (1 + \epsilon/2)m + \frac{1}{\epsilon} \log(n)$$

Additive regret

- So, have $M \leq OPT + \epsilon OPT + 1/\epsilon \log(n)$.
- Say we know we will play for T time steps. Then can set $\epsilon = (\log(n) / T)^{1/2}$. Get $M \leq OPT + 2(T * \log(n))^{1/2}$.
- If we don't know T in advance, can guess and double.
- These are called "additive regret" bounds.

$$M \leq \frac{-m \ln(1 - \epsilon) + \ln(n)}{\epsilon} \approx (1 + \epsilon/2)m + \frac{1}{\epsilon} \log(n)$$

Extensions

- What if experts are actions? (rows in a matrix game, choice of deterministic alg to run,...)
- At each time t , each has a loss (cost) in $\{0,1\}$.
- Can still run the algorithm
 - Rather than viewing as "pick a prediction with prob proportional to its weight",
 - View as "pick an expert with probability proportional to its weight"
- Same analysis applies.

Extensions

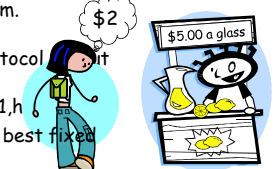
- What if losses (costs) in $[0,1]$?
 - Here is a simple way to extend the results.
 - Given cost vector c , view c_i as bias of coin. Flip to create boolean vector c' , s.t. $E[c'_i] = c_i$. Feed c' to alg A .

world
→
 c
→
\$
→
 c'
→
A
 - For any sequence of vectors c' , we have:
 - $E_A[\text{cost}(A)] \leq \min_i \text{cost}'(i) + [\text{regret term}]$
 - So, $E_{\$}[E_A[\text{cost}(A)]] \leq E_{\$}[\min_i \text{cost}'(i)] + [\text{regret term}]$
 - LHS is $E_A[\text{cost}(A)]$.
 - RHS $\leq \min_i E_{\$}[\text{cost}'(i)] + [\text{r.t.}] = \min_i [\text{cost}(i)] + [\text{r.t.}]$
- In other words, costs between 0 and 1 just make the problem easier...

What can we use this for?

- Can use to combine multiple algorithms to do nearly as well as best in hindsight.
 - E.g., do nearly as well as best strategy in hindsight in repeated play of matrix game.
- Extension: "sleeping experts". E.g., one for each possible keyword. Try to do nearly as well as best "coalition".
- More extensions: "bandit problem", movement costs.

Online pricing

- Say you are selling lemonade (or a cool new software tool, or bottles of water at the world expo).
 - Protocol #1: for $t=1,2,\dots,T$
 - Seller sets price p^t
 - Buyer arrives with valuation v^t
 - If $v^t \geq p^t$, buyer purchases and pays p^t , else doesn't.
 - v^t revealed to algorithm.
 - repeat
 - Protocol #2: same as protocol #1 but without v^t revealed.
 - Assume all valuations in $[1, h]$
 - Goal: do nearly as well as best fixed price in hindsight.
- 

Online pricing

- Say you are selling lemonade (or a cool new software tool, or bottles of water at the world expo).
- Protocol #1: for $t=1,2,\dots,T$
 - Seller sets price p^t
 - Buyer arrives with valuation v^t
 - If $v^t \geq p^t$, buyer purchases and pays p^t , else doesn't.
 - v^t revealed to algorithm.
- Bad algorithm: "best price in past"
 - What if sequence of buyers = $1, h, 1, \dots, 1, h, 1, \dots, 1, h, \dots$
 - Alg makes T/h , OPT makes T .

Factor of h worse!

Online pricing

- Say you are selling lemonade (or a cool new software tool, or bottles of water at the world expo).
 - Protocol #1: for $t=1,2,\dots,T$
 - Seller sets price p^t
 - Buyer arrives with valuation v^t
 - If $v^t \geq p^t$, buyer purchases and pays p^t , else doesn't.
 - v^t revealed to algorithm.
 - Good algorithm: Randomized Weighted Majority!
 - Define one expert for each price $p = (1+\epsilon)^i \in [1, h]$.
 - Best price of this form gives profit $\geq \text{OPT}/(1+\epsilon)$.
 - Run RWM algorithm. Get expected gain at least:

$$\begin{aligned} & (\text{best expert})/(1+\epsilon) - O(\epsilon^{-1} h \log n) \\ & = \text{OPT}/(1+\epsilon)^2 - O(\epsilon^{-1} h \log(\epsilon^{-1} \log h)) \end{aligned}$$
- [extra factor of h coming from range of gains]

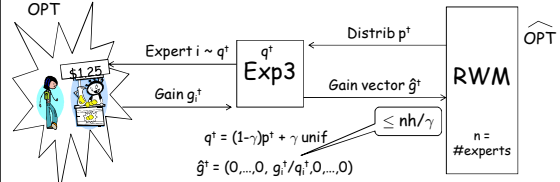
Online pricing

- Say you are selling lemonade (or a cool new software tool, or bottles of water at the world expo).
- What about Protocol #2? [just see accept/reject decision]
 - Now we can't run RWM directly since we don't know how to penalize the experts!
 - Called the "adversarial multiarmed bandit problem"
 - How can we solve that?



Multi-armed bandit problem

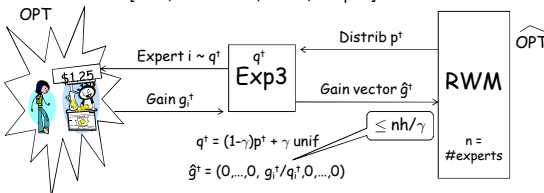
Exponential Weights for Exploration and Exploitation (exp³)
[Auer,Cesa-Bianchi,Freund,Schapire]



1. RWM believes gain is: $p^t \cdot \hat{g}^t = p_i^t (g_i^t / q_i^t) \equiv g_{RWM}^t$
2. $\sum_t g_{RWM}^t \geq \widehat{OPT} / (1+\epsilon) - O(\epsilon^{-1} nh / \gamma \log n)$
3. Actual gain is: $g_i^t = g_{RWM}^t (q_i^t / p_i^t) \geq g_{RWM}^t (1-\gamma)$
4. $E[\widehat{OPT}] \geq OPT$. Because $E[\hat{g}_j^t] = (1 - q_j^t)0 + q_j^t(g_j^t/q_j^t) = g_j^t$, so $E[\max_j \sum_t \hat{g}_j^t] \geq \max_j [E[\sum_t \hat{g}_j^t]] = OPT$.

Multi-armed bandit problem

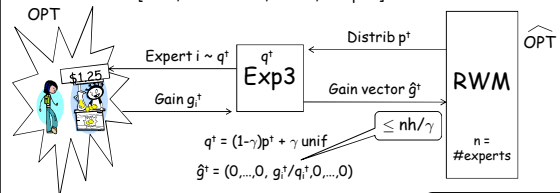
Exponential Weights for Exploration and Exploitation (exp³)
[Auer,Cesa-Bianchi,Freund,Schapire]



2. $\sum_t g_{RWM}^t \geq \widehat{OPT} / (1+\epsilon) - O(\epsilon^{-1} nh / \gamma \log n)$
3. Actual gain is: $g_i^t = g_{RWM}^t (q_i^t / p_i^t) \geq g_{RWM}^t (1-\gamma)$
4. $E[\widehat{OPT}] \geq OPT$. Because $E[\hat{g}_j^t] = (1 - q_j^t)0 + q_j^t(g_j^t/q_j^t) = g_j^t$, so $E[\max_j \sum_t \hat{g}_j^t] \geq \max_j [E[\sum_t \hat{g}_j^t]] = OPT$.

Multi-armed bandit problem

Exponential Weights for Exploration and Exploitation (exp³)
[Auer,Cesa-Bianchi,Freund,Schapire]



Conclusion ($\gamma = \epsilon$):
 $E[\text{Exp3}] \geq \widehat{OPT} / (1+\epsilon)^3 - O(\epsilon^{-2} h n \log(n))$

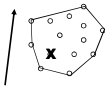
Can even reduce ϵ^2 to ϵ^1 with more care in analysis.

Almost as good as protocol 1!

Extensions (of expert or bandit problem)

[KV] setting:

- Implicit set S of feasible points in R^m . (E.g., $m = \#edges$, $S = \{\text{indicator vectors } 011010010 \text{ for possible paths}\}$)
- Assume have oracle for offline problem: given vector c , find $x \in S$ to minimize $c \cdot x$. (E.g., shortest path algorithm)
- Use to solve online problem: on day i , must pick $x_i \in S$ before c_i is given.
- $(c_1 \cdot x_1 + \dots + c_T \cdot x_T) / T \rightarrow \min_{x \in S} (c_1 + \dots + c_T) / T$.



[Z] setting:

- Assume S is convex.
- Allow $c(x)$ to be a convex function over S .
- Assume given any y not in S , can algorithmically find nearest $x \in S$.

Other models in learning

Lots of other models considered as well for different kinds of problems.

- "Active learning": have large unlabeled sample and alg may choose among these.
 - E.g., web pages, image databases.
- "Membership query learning": Algorithm can construct its own examples.
 - E.g., features represent variable-settings in some experiment, label represents outcome.
- "Semi-supervised learning": use of labeled+unlabeled data in passive setting.