Claim: $H$ nonuniform learnable with $A$, $m^{NUL}$ $\Rightarrow$ $H = \bigcup_{n\geq 1} H_n$,
each $H_n$ is agnostic PAC learnable

Recall Corollary 6.4 (NFL)
$X$ domain, $VCdim(X) \geq 2n$ ($\Rightarrow \exists C \leq X, |C|=2n$, $H$ shatters)
$\Rightarrow \forall$ learning algorithms $A$, $\exists D$: $\exists h\in H$: $L_D(h)=0$
with prob $\geq \frac{1}{2}$ over $S\sim D^n$: $L_D(A(S)) \geq 1/8$
[Holds also for $A$ which learn one on $H'\geq H$]

Proof: $H_n = \{h\in H: m^{NUL}_H (0.1, 0.1, h) \leq n\}$
$\underset{\Rightarrow}{}$ $D$ realizable with $H_n$     $S\sim D$
$\forall h\in H_n$: $L_D(A(S)) \leq L_D(h) + 0.1$
$\quad" \quad \leq \underbrace{\min_{h\in H_n} L_D(h)}_{=0} +0.1 =0.1$    $L_D(A(S)) \leq 0.1$

Assume $VCdim(H_n) \geq 2n$
So apply Cor.6.4 $\forall$ learning alg $A'$, $\exists D$ with realiz. wrt $H_n$
s.t. with prob $\geq 1/2$ over $S\sim D^n$ $\underline{L_D(A'(S)) \geq \frac{1}{8}}$
$\quad\downarrow$ It should hold for $A$, as well
$\quad\Rightarrow VCdim(H_n) \not\geq \, < 2n$
$\underset{\Rightarrow}{\text{Finit Thm}}$ $H_n$ is agnostic PAC-learnable

Recall  SRM: $\underset{h\in H}{argmin}$ $L_S(h) + \varepsilon_{n(h)} (m, w(n(h)) \delta)$
$\qquad n(h):= \min\{n\in N: h\in H_n\}$
$\qquad w(n)$ weight fct: $\sum_{n\geq 1} w(n)\geq 1$
$\qquad \varepsilon_n(m,\delta) = \min \{\varepsilon\in(0,1): m^{UC}_{H_n} (\varepsilon,\delta)\leq m\}$

MDL and Occam's Razor (as a special case of SRM)
$H$ countable $\Rightarrow \underset{n\in N}{\overset{finite}{\bigcup}} \{h_n\}=H$
$m^{UC}_{\{h_n\}} \leq \frac{\log(\frac{2}{\delta})|H_n|)}{2\varepsilon^2} = \frac{\log(2/\delta)}{2\varepsilon^2}$
$\varepsilon_n(m,\delta)= \sqrt{\frac{\log(2/\delta)}{2m}}$ , $n(h)=n$

SRM becomes  $\underset{h\in H}{argmin}$ $L_S(h)+ \sqrt{\frac{-\log(w(n)+\log\frac{2}{\delta})}{2m}}$
$\qquad = \quad " \quad + \sqrt{\frac{-\log(w(h))+\log\frac{2}{\delta}}{2m}}$

Want to use a $w$ based on a description
language:
$\qquad \Sigma$: finite alphabet (e.g. $\{0,1\}$)
$\qquad \Sigma^*$: all finite strings over $\Sigma$
$\qquad d: H \to \Sigma^*$   description language
Focus on prefix-free languages:
$\qquad \forall h\neq h'\in H: d(h)$ is not a prefix of $d(h')$
Kraft inequality: $S \leq \{0,1\}^*$ prefix free
$$\sum_{\sigma\in S} \frac{1}{2^{|\sigma|}} \leq 1$$
Proof: Draw $0,1$ uniformly at random (and stop if it equals a string in $S$)
$\qquad \forall \sigma\in S: P(\sigma)= \frac{1}{2^{|\sigma|}}$
$$\sum_{\sigma\in S} \frac{1}{2^{|\sigma|}} = P(S) \leq 1$$

We can use $d(h)$ as a weight $w(h)=\frac{1}{2^{|h|}}$
MDL: $\underset{h\in H}{argmin}$ $\boxed{L_S(h)}+ \sqrt{\frac{\boxed{|h|}+\ln(2/\delta)}{2^m}}$
$\qquad\qquad$ tradeoff between emp. risk at "complexity"
$\qquad\qquad$ of describing $h$
Occam's Razor:
$\qquad$ "A short explanation tends to be more
$\qquad\qquad$ valid than a long one"

Consistency
Algorithm $A$ is consistent wrt. $H$ and $P$ (=set of distributions)
if : $\exists m^{con}_H (0,1)^2 \times H\times P \to N$ s.t.
$\qquad \forall \varepsilon,S, \forall h\in H, \forall D\in P: L_D(A(S)) \leq L_D(h)+\varepsilon$

Example: MEMORIZE is consistent
Given a test instance $x$, MEMORIZE returns the
majority of the labels of instances in the sample,
which are equal to $x$. ( just predict the majority of all labels if there is no $x$)

$\Rightarrow$ This notion is too weak to capture "learning"

Comparison

| | Bounds on true error by the emp. risk | How many exampl. are needed to be as best to any hypo. in $H$ | Encode prior knowledge |
|---|---|---|---|
| (agn) PAC | $\checkmark$ | $\checkmark$ (in advance) | $\checkmark$ (specify $H$) |
| Nonuniform | $\checkmark$ (Thm 7.4) | $\checkmark$ depends on the best $h\in H$ | $\checkmark$ (weights) |
| Consistent | $\times$ | $\times$ | $\times$ |

Runtime of Algorithms

Input size?
$\qquad$ e.g. sample size is a bad idea.
$\qquad \to$ depend on $\varepsilon$ and $\delta$
Computational Complexity for Learning Algorithms
$\cdot$ $A$ solves a learning task $(Z, H, l)$
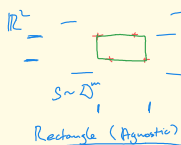$\quad$ in time $O(f(\varepsilon,\delta))$, if:
$\qquad \cdot A$ terminates in $O(f(\varepsilon,\delta))$ time.
$\qquad \cdot$ output of $A$ should applicable to
$\qquad\quad$ (new) instances in $O(f(\varepsilon,\delta))$ time.
$\qquad \cdot$ agn. PAC learn $(Z,H,l)$
$A$ solves a sequence of learning problems
$\quad (Z_n, H_n, l_n)_{n=1}^\infty$ in $O(f_n(\varepsilon,\delta))$ if:
$\quad$ for each fixed $n$, $A$ solves $(Z_n, H_n, l_n)$ in time
$\qquad\qquad O(f_n(\varepsilon,\delta))$   (efficient if $f_n=O(\text{poly } c(\frac{1}{\varepsilon},\frac{1}{\delta},n)))$
Example: Rectangles (realizable) in $\mathbb{R}^n$

$\mathbb{R}^2$



$S\sim D^m$

For each dim:
find min and max in $O(m)$ time
Total time $O(nm)$ $=O(n\cdot m\, c(\varepsilon,\delta))$

Rectangle (Agnostic)



It's NP-hard to compute the ERM-rectangle
One can learn it in $O(m^{O(n)})$.