

## ANNOUNCEMENTS

1. HW5 late due tomorrow, 5pm ; HW6 due next Mon, 5pm
  2. Grading policies updated on course website
- 

Turn your non-note-taking devices off now

Last week: ERM, bias-variance tradeoff

find  $\hat{h} \in \mathcal{H}$  s.t.  $\arg \min_{h \in \mathcal{H}} \mathcal{E}(h)$   $\xrightarrow{\text{training error}}$

$$\mathcal{E}(\hat{h}) \leq \mathcal{E}(h^*) + 2 \sqrt{\frac{1}{2n} \log \frac{2k}{\dots}}$$

gen err of  $\hat{h}$       gen err of  $h^*$        $\downarrow$  num of samples       $\xrightarrow{\text{complexity of } \mathcal{H}}$   $k = |\mathcal{H}|$

if  $\mathcal{H}$  becomes more complex — {linear, poly, quad},  
then  $\mathcal{E}(h^*) \downarrow$

TRADE OFF b/w  $\mathcal{E}(h^*)$ ,  $2 \sqrt{\dots}$  when  $\mathcal{H} \uparrow$   $\xrightarrow{\text{complexity of}}$

# Complexity - regularized ERM

$$\hat{h} = \underset{h \in \mathcal{H}}{\operatorname{argmin}} \hat{e}(h) + \lambda \Omega(h)$$

measures complexity of " $\underline{h}$ "

Bayesian  
Viewpoint  
of ERM!

linear rego.

$$\min \underbrace{\frac{1}{n} \sum (y - \hat{y})^2}_{\text{cost fn}} + \lambda^2 \|a\|^2$$

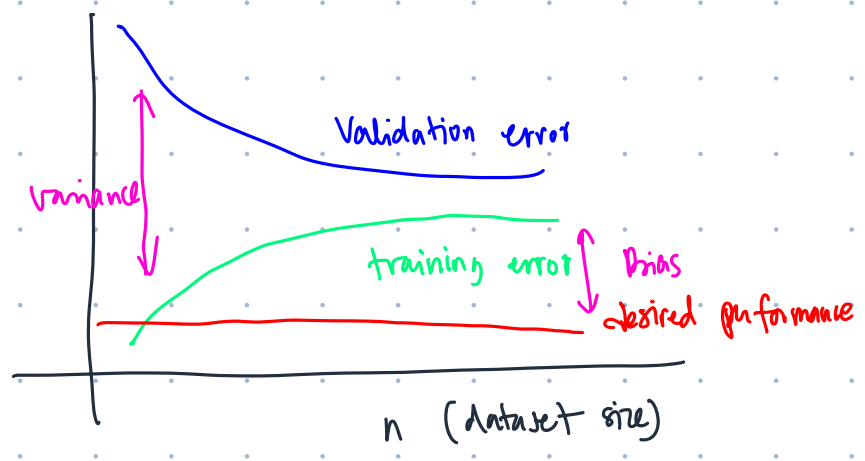


MAP estimate



# Overfitting vs. underfitting

overfitting  $\leftarrow$  



# Bad optimization vs. Bad objective

LP -  $l(\theta)$  objective

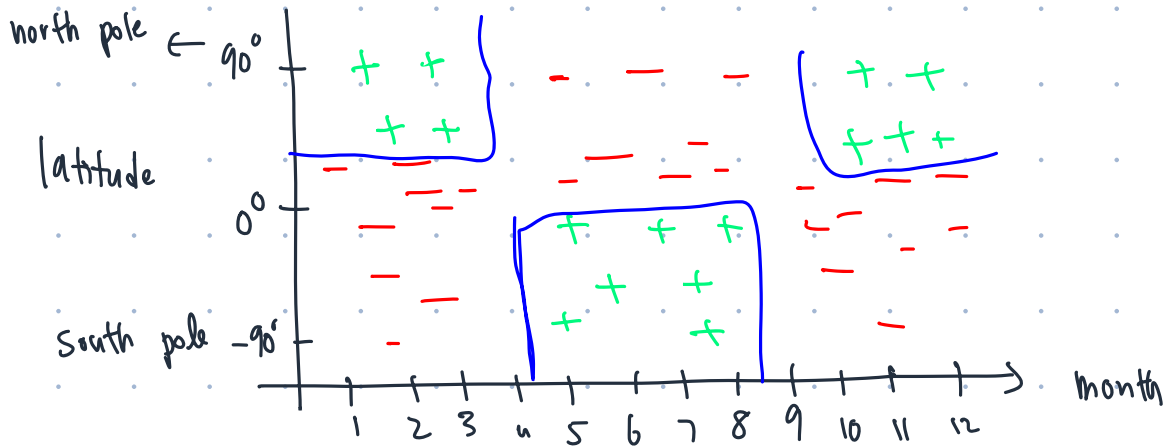
Validation set -  $acc(LP) < acc(SVM)$   
 $l(\theta_{LP}) < l(\theta_{SVM})$

} OPT was the problem!

$acc(LP) < acc(SVM)$   
 $l(\theta_{LP}) \geq l(\theta_{SVM})$

} OPT worked, but objective didn't

Today : DECISION TREES — sledding on the slope!



No "line", but kernels could help!

but, region-based segregation!!!  
seems like a good idea

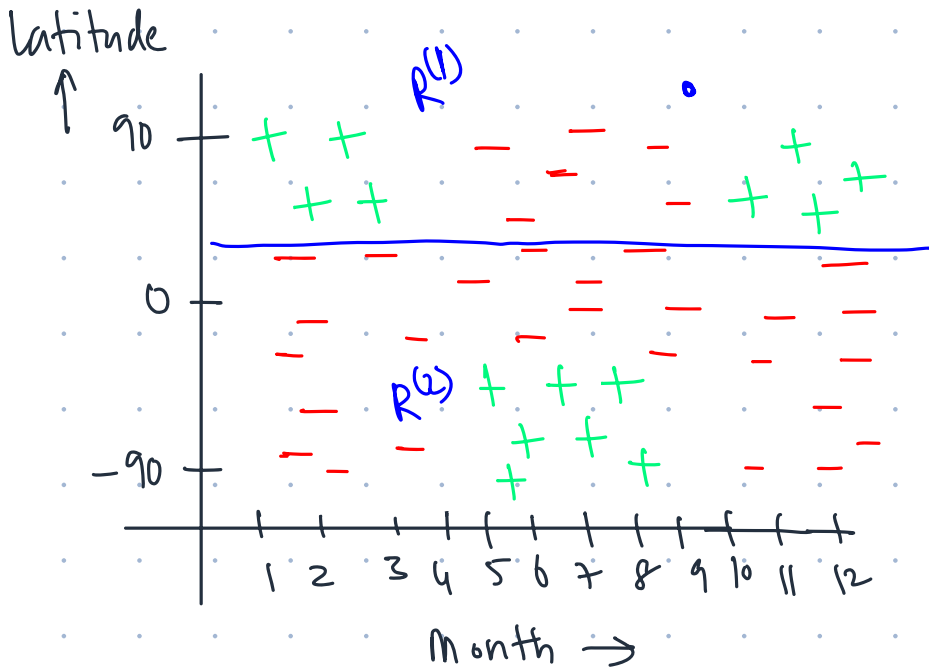
$$X = R^{(1)} \cup R^{(2)} \dots \cup R^{(r)}$$

splitting  $X$  into maximally-compact regions is NP-hard!

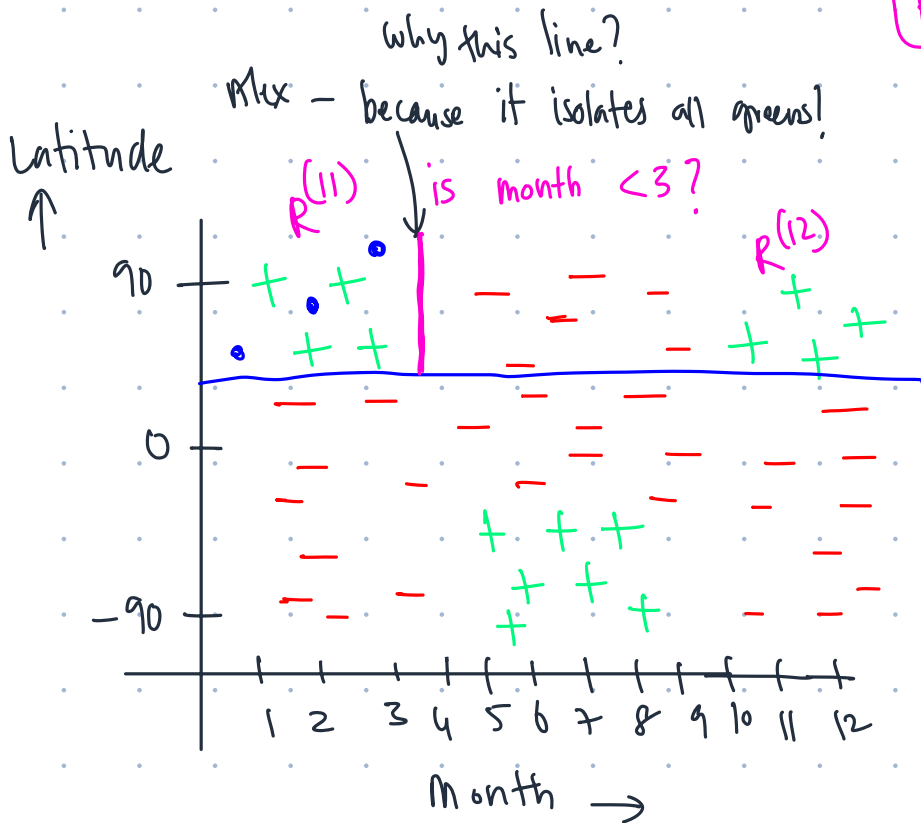
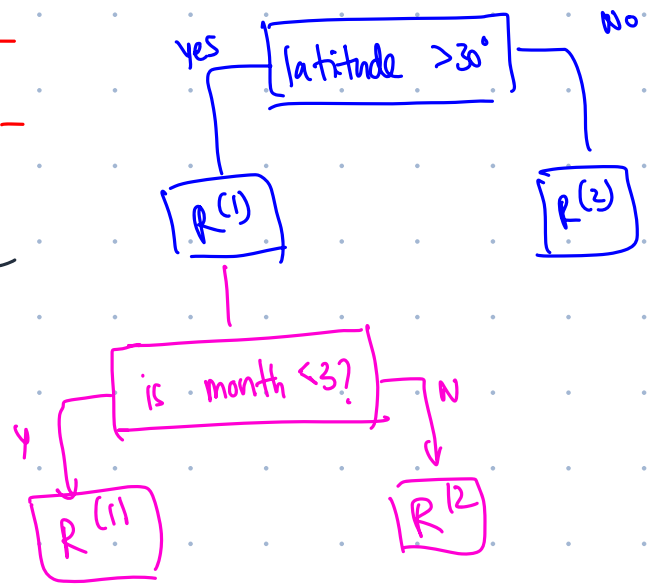
GREEDY WAY!

Idea: Top-down, greedy, recursive partitioning

We're gonna play 20 Q with this dataset



is latitude  $> 30^\circ$ ?



Bonus: FREE INTERPRETABILITY



## Two main "points of interest"

(1) what does it mean to ask the "right" questions?

(2) when do I stop asking questions?

when we have a "pure" region, further splitting is useless!

if not possible to get a pure node, just take majority vote!

Asking the "right" questions  
↳ = greedy best

Is latitude  $> 30^\circ$ ?  
vs.  
Is month  $> 6$ ?  
vs.  
Is latitude  $> 60^\circ$ ?

↳ what is a question?

↳ feature (eg. latitude, month)  
↳ threshold (eg.  $30^\circ$ , March)

What do we want from parent  $\rightarrow$  children?

some,  $J(R)$  = impure a region is

$$\max J(R^{(0)}) - \frac{J(R^{(1)}) + J(R^{(2)})}{2} \rightarrow \text{maximize the decrease in impurity!}$$

more generally, we weighted avg of children:

$$\max J(R^{(0)}) - \left[ \frac{|R^{(1)}| J(R^{(1)}) + |R^{(2)}| J(R^{(2)})}{|R^{(1)}| + |R^{(2)}|} \right]$$

↳ #samples in  $R^{(1)}$

# Idea-1: Misclassification error

Measuring impurity  
 $J(R)$

What if impurity of a node = misclassification error of that region

$$R - \underbrace{600+, 400-} \rightarrow \text{misclassification rate} = \frac{400}{1000} = 0.4$$

$R =$  not pure  
prediction = +ve

$$R^{(1)} - 400+, 0- \rightarrow \text{misclassification rate} = 0$$

$$900+, 100- \quad J=0.1$$

$(f, t)$

$$C1: 700+, 100- \quad C2: 200+, 0-$$

$$J = \frac{100}{800}$$

$$J = 0$$

$$900+, 100- \quad J=0.1$$

$(f', t')$

$$C1: 500+, 100- \quad C2: 400+, 0-$$
$$J = \frac{100}{600} \quad J=0$$

$$\text{OBJ} - 0.1 - \left( \frac{800 \left( \frac{100}{800} \right) + 200(0)}{1000} \right)$$

$$= 0.1 - \frac{100}{1000} = 0$$

$$0.1 - \frac{100}{1000} = 0$$

- ① Neither splits are better
- ② Neither splits gave any improvement over parent.

## Idea-2: Measure of randomness

Entropy: measure of randomness

$$J = - \sum_{y \in \{+1, -1\}} p_y \log_2 p_y$$

$p_y$  = proportion of samples in  $R$  with class label =  $y$ .

600+, 400-

$$p_+ = \frac{600}{1000}, \quad p_- = \frac{400}{1000}$$

$$J = - \left[ 0.6 \log_2 0.6 + 0.4 \log_2 0.4 \right]$$

# Entropy vs. misclassification error

## Binary classification

$$J_{\text{ent}} = -p_+ \log p_+ - (1-p_+) \log(1-p_+)$$

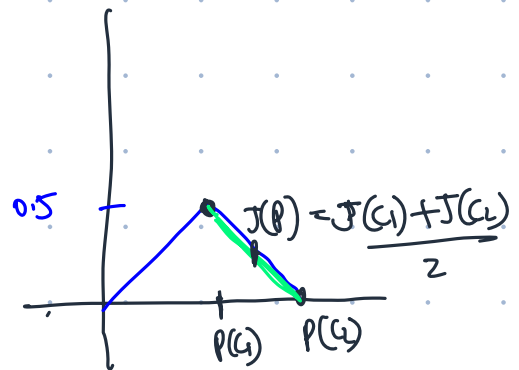
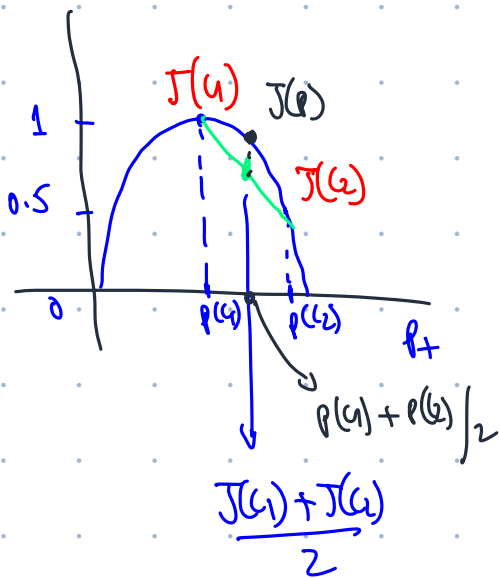
$$J_{\text{mis}} = 1 - \max(p_+, 1-p_+)$$

}  $p_- = 1-p_+$  in two-class

Assume: whatever splits we have, they have equal cardinalities

$$|R^{(c_1)}| = |R^{(c_2)}|$$

$$p_+^{(\text{parent})} = \frac{p_+^{(c_1)} + p_+^{(c_2)}}{2}$$



not strictly concave,  
children need not have  
lower error

strictly concave,  
children will always have lower  
entropy than parent

Aside: Categorical attributes?

Splitting on NetID?