

Last time:

How to get primal variables from dual?

- argmax to conjugate of g

Example:

$$\begin{aligned}
 p(x) &= f(Ax) + \frac{\lambda}{2} \|x\|^2 \\
 -D(y) &= f^*(x) + g^*(A^T y) \\
 &= f^*(x) + \frac{1}{2\lambda} y^T A A^T y \Rightarrow \sup_x \left\{ y^T A - \frac{\lambda}{2} \|x\|^2 \right\} \\
 &\quad \downarrow \\
 &\quad x = \frac{1}{\lambda} A^T y
 \end{aligned}$$

Kernel trick:  $\hat{b} = \hat{A}x$   
 $= \frac{1}{\lambda} \hat{A} A^T y$   
 $= \frac{1}{\lambda} K(\hat{A}, A) y$

Admin:

- A2 / A3: pick up ready
- A4 marked version due Monday
- A5 due Wednesday
- Project Proposal due Nov 3rd.

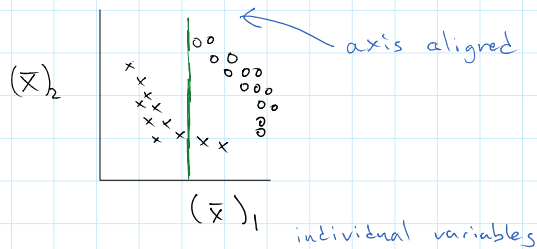
If have  $y^*$ ;  $x^* = A^T y^*$   
 assuming strong duality holds

Ensemble Methods

- "models that use other models."
- can give better performance than individual models
- e.g. decision trees, bootstrap, bagging, boosting, random forests  
 stacking, jackknives, etc...

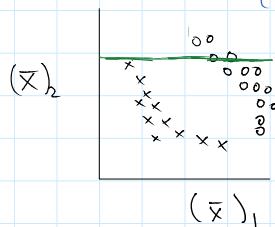
Decision Stumps

- find variable 'j' and threshold 'T',  
 such that classifier  $[(\bar{x}_i)_j \geq T]$   $O(np)$   
 maximizes some score  
 (e.g. Training accuracy) or other scores

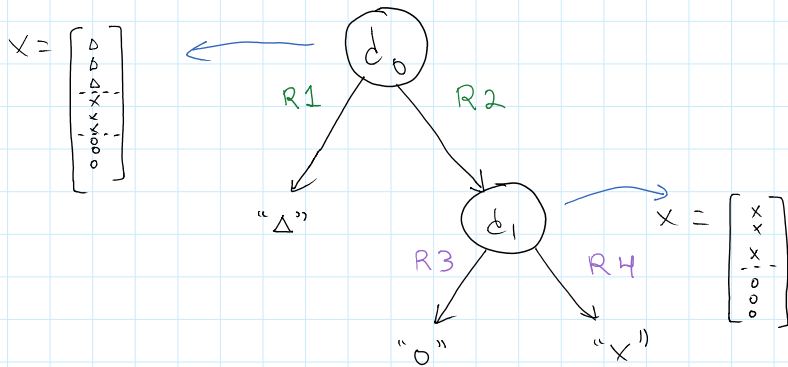
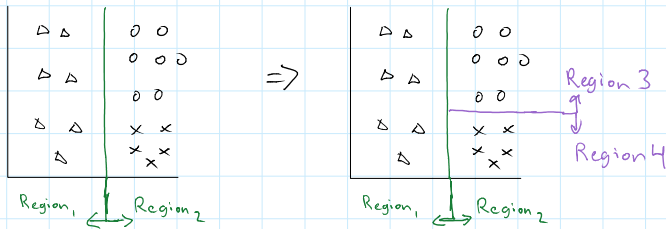


Randomized Decision Stump

- choose a random subset  $\{1, \dots, d\}$   
 only choose "j" from subset  
 (not recommended for accuracy, but increases speed)



# Decision Trees



## Issues:

1. Very interpretable (for shallow depth)
2. Finding optimal tree is NP-Hard
  - start w/ full data set (mostly done greedily)
  - learn decision stump
  - partition into smaller data sets
3. Choosing Score
  - last layer: counting
  - intermediate layers "info gain"  
"mutual information"
4. Choosing classifier
  - [randomized stump] linear classifier, non-linear classifier
  - cost increases  $\longrightarrow$
5. Pruning: prune nodes that don't decrease validation error
6. CART (Breiman) C4.5 (Quinlan)
7. tend to be worse than linear models in high dimensions

## Intogain

$$I(X, Y) = H(X) - H(X|Y)$$

↓ entropy before split      ↑ entropy after split

$$H(X) = - \sum_{x \in X} p(x) \log_2(p(x))$$

"entropy" - measure of randomness

$$H(X|Y) = \sum_{y \in Y} p(y) H(X|Y=y)$$

↑ probability of assigning region 'R'      ← zero if split perfectly classifies (otherwise higher)

- Comparison to log sum exp:

$$f(x) = \log \left( \sum_i \exp(x_i) \right)$$

$$f^*(y) = \sum_i y_i \log_e(y_i) \quad \text{s.t.} \quad y_i \geq 0 \\ \sum_i y_i = 1$$

similar to entropy,  
but with conditions

## Stacking (Model Averaging)

- train  $m$  different classifiers

naïve bayes, knn, logistic, neural nets ...

↓                      ↓                      ↓                      ↓  
stacking: new classifier that combines the outputs

$$\hat{y}_i = f \left( \sum_{j=1}^m w_j h_j(\bar{x}_i) \right)$$

Special Cases

1. Linear  $h_j(\bar{x}_i) = (\bar{x}_i)_j$

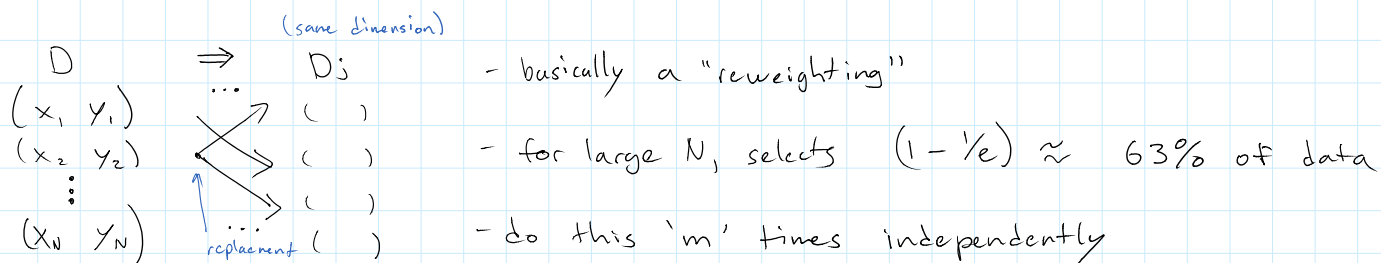
2. Neural Nets  $h_j(\bar{x}_i) = \sigma(\bar{w}_{0j}^T x_i)$

3. Generalized Additive  $h_j(\bar{x}_i) = g_j((\bar{x}_i)_j)$

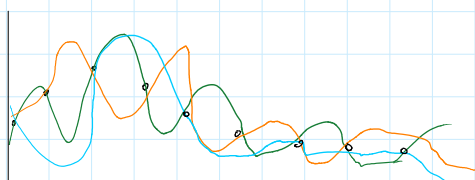
\* Winner of Netflix prize (\$1 million)

## Bootstrap / Bagging

- Input: a high variance classifier (over-fitting)
- Output: a lower variance classifier



Bagging: train high variance classifier on each sample, average results



(averaging will provide a good fit)

## Random Forests

- Bagging
  - Decision Trees
  - Info gain
  - Random Decision stumps
- ↳ why? - speed  
- reduce correlation between classifier

Kinect: RF to predict body part or background at each pixel.

# Boosting

Input: a 'weak' learner,  
binary classifier w/ accuracy  $> 50\%$  ex: decision stump,  
decision trees.

Output: a 'strong' learner

Ada Boost (Freund & Schapire)

Set datapoint weights  $z_i = \frac{1}{N}$  (weight: examples worth more)  
✓ number of weak classifiers to learn

for  $j = 1 : m$

- train 'weak' classifier with weight  $\bar{z}$

- choose optimal  $w_j$  (exponential loss)

- re-weight  $z_i = z_i \exp(w_j \mathbb{I}[y_i \neq h_j(\bar{x}_i)])$

when to stop? cross validation

increase weight of  
points you got "wrong"

$$\hat{y} = f\left(\sum_j w_j h_j(\bar{x}_i)\right)$$

Issues:

- can overfit

- doesn't work if classifiers lack diversity