**Problem 1** : more robust against noisy samples.

**Problem 2 :**

$$J(a) = \sum_{i=1}^{n} a_i - \frac{1}{2} \sum_{i=1} \sum_{j=1} a_i \cdot a_j \cdot y_i \cdot y_j \cdot \underline{\langle x_i, x_j \rangle}$$

↑ Similarity measure.

$\Downarrow$

① polynomial k

② RBF $K. = e^x$ ; $x = -\frac{\|x_i - x_j\|^2}{2 \cdot \sigma^2}$

$k(x_i, x_j)$

**Problem 3** : SV = defines the hyperplane that seperates +s from −s

significance = reduced complexity, not relying on non SVs.

**Problem 4 :**

$$\min_{\theta} \boxed{C} \left[ \sum y_i \text{cost}_1(\theta)^T x_i + (1-y_i) \text{cost}_0(\theta)^T x_i \right] + \frac{1}{2} \cdot \sum_{j=1}^{d} \theta_j^2$$

↳ classification error           max margin

if $C = \infty \Rightarrow$ only care about min classification error

$C = \boxed{-\infty} 0 \Rightarrow \frac{1}{2} \sum_{j=1}^{d} \theta_j^2 \Rightarrow$ only care about max the margin.

↓
✗

**Problem 5 :**



hyperplane.

dimensionality of HP = d of samples − 1

HP = decision boundary.

**Problem 6 :** kernel tricks replaces $\langle x_i, x_j \rangle$ in the Lagragian dual with $k(x_i, x_j)$ : raises lower-dimensional samples to higher dimensions without explictly doing the mapping.

**Problem 7 :** $a_j = 0.1 \Rightarrow x_j$ is a support vector.

Problem 8: Primal $= \frac{1}{2}\sum_{j=1}^{d}\theta_j^2$, st. $y_i\underbrace{(\theta^T x_i + b)}_{\text{Constraints}} \geq 1$

why? do we want the dual $\Rightarrow$ the dual internalized the constraints into the objective function.

Problem 9: kernel tricks.

Problem 10: Pros: reduced comp times/complexity

achieved linear seperability

Cons: assumption that raising to higher dim. is meaningful/ preserving data integrity.

Problem 11 what is $H(X,Y)$

$$H(X,Y) = -\sum_{x}\sum_{y} p(x,y)\log_2 p(x,y)$$

$$= -\left(\sum_{x} P(x,0)\log_2 P(x,0) + P(x,1)\log_2 P(x,1)\right)$$

$$= -P(0,0)\log_2 P(0,0) - P(0,1)\log_2 P(0,1)$$
$$\quad - P(1,0)\log_2 P(1,0) - P(1,1)\log_2 P(1,1)$$

cloudiness

Problem 12. $H(Y|X=1) = -\sum_{y\in Y} p(y|x=1)\log_2 p(y|x=1)$

$P(y|x) = \frac{P(x,y)}{P(x)}$ $\quad = -\sum_{y\in Y}\frac{P(x=1,y)}{P(x)}\log_2\frac{P(x=1,y)}{P(x)}$

$\hookrightarrow = -\frac{P(x=1,y=0)}{P(x)}\log_2\frac{P(x=1,y=0)}{P(x)}$ $\qquad P(x) =$

$\quad - \frac{P(x=1,y=1)}{P(x)}\log_2\frac{P(x=1,y=1)}{P(x)}$

Problem 13    $H(Y|X) = \sum\limits_{x \in X} p(x) H(Y|x)$

$H(Y|X=0)$

previously, we have $H(Y|X=1)$; use the same approach to get

$H(Y|X) = P(X=0) H(Y|X=0) + P(X=1) H(Y|X=1)$

$= \dfrac{75}{100} \cdot H(Y|X=0) + \dfrac{25}{100} \cdot H(Y|X=1)$

Problem 14.    $IG(Y|X) = H(Y) - H(Y|X)$

Problem 15.    $IG(Y|X)$, when x is useless $\Rightarrow IG(\cancel{\frac{Y}{X}}|X) = 0$

Problem 16.    $IG(Y|X)$, when x is equivalent to Y) $\Rightarrow IG(Y|X) = H(Y)$

problem 17.    Overfitting. $\Rightarrow$ potentially, train acc = 100%, with one leaf node correspond to one each sample.

Problem 18 - Decision Tree slides 8 & 12.

Problem 19    Decision tree is a greedy heuristics. At each split (down the tree) it picks the feature split that guarantees max $IG$.

problem 20.    pro: very interpretable.
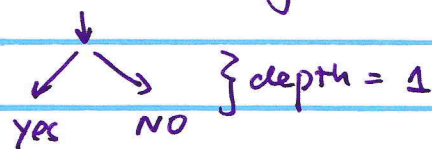         cons: 1. prone to overfit
               2. greedy heuristic $\Rightarrow$ tries to find only local optimum
               3. time & computation complexity
                  (compared to other non NN approaches)

**Problem #1.** a weak leaner = a model with > 50% accuracy.

example: decision stumps



$\}$ depth = 1

yes     NO

**Problem #2.** misclassified samples get bigger weights

correctly classified "      "    smaller weights

**Problem #3.** ① Simplicity

② Interpretability

**Problem #4.** $W_{t,i}$.     $t$: iteration / trial #

$i$ = the index of a sample, $x_i$.

$W_{t,i} \neq$ (not usually) $W_{t+1,i}$

$\sim$ cross entropy loss

**Problem #5** $J_{reg_t}(\theta) = -\sum_{i=1}^{n} W_{i,t} y_i \log h_\theta(x_i) + (1 - y_i) \log(1 - h_\theta(x_i))$

$+ \lambda \| \theta_{[1:d]} \|_2^2$ $\rightarrow$ regularization

( want to minimize $\theta$, to prevent

regularized cost for the $t$th iteration.    any index from getting

too big)