

28] SVM maximizes $M = \text{width of margin}$
 $\beta_0, \dots, \beta_p, \epsilon_1, \dots, \epsilon_n$

subject to $\sum_{j=1}^p \beta_j^2 = 1$, $\epsilon_i \geq 0$, $\sum_{i=1}^n \epsilon_i \leq D$

omit β_0 from constraint, $\|\underline{\beta}\|_2 = 1$

and $z_i (\beta_0 + \underline{\beta}^T \underline{x}_i) \geq M (1 - \epsilon_i)$,

compare 20].

29] A slack variable $\epsilon_i = 0$ if \underline{x}_i is on the correct side of the margin. If $\epsilon_i > 0$, \underline{x}_i is on the wrong side of the hyperplane.

So $z_i (\beta_0 + \underline{\beta}^T \underline{x}_i) \geq m$ has $\epsilon_i = 0$ and is necessary for \underline{x}_i to be on the correct side of the margin. If $z_i (\beta_0 + \underline{\beta}^T \underline{x}_i) \geq m(1 - \epsilon_i)$ with $\epsilon_i > 0$ (but not if $\epsilon_i = 0$), then

\underline{x}_i is on the wrong side of the hyperplane (misclassified). See 17].

30) Let the kernel function be $k(\underline{x}_i, \underline{x}_j)$.

A linear kernel is $k(\underline{x}_i, \underline{x}_j) = \underline{x}_i^T \underline{x}_j$.

A polynomial kernel of degree d is

$$k(\underline{x}_i, \underline{x}_j) = [1 + \underline{x}_i^T \underline{x}_j]^d.$$

A radial kernel is $k(\underline{x}_i, \underline{x}_j) =$

$$\exp\left[-\gamma \sum_{k=1}^p (x_{ik} - x_{jk})^2\right] = \exp\left[-\gamma \|\underline{x}_i - \underline{x}_j\|_2^2\right].$$

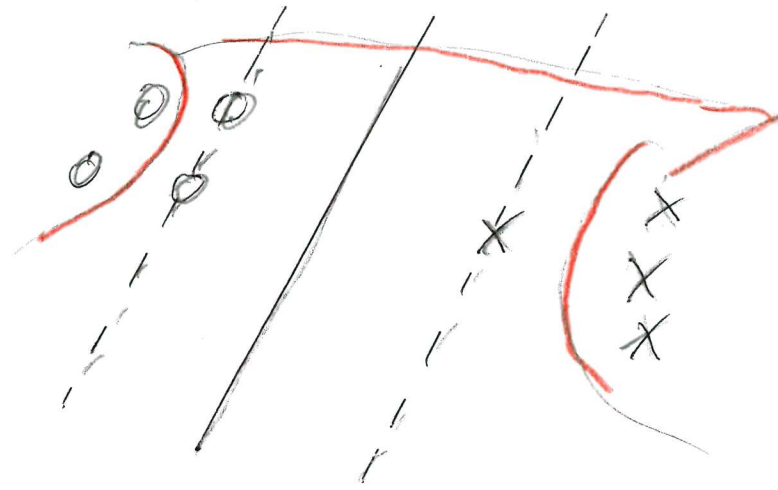
31) A SVM uses $SP = f(\underline{x}) =$

$$\beta_0 + \sum_{i=1}^n \alpha_i k(\underline{x}, \underline{x}_i) = \beta_0 + \sum_{i \in S} \alpha_i k(\underline{x}, \underline{x}_i)$$

$= SP = SP(\underline{x})$. where S is the index set of support vectors.

Note: the support vectors determine the (boundaries) hyperplane and margin: if they are moved, the hyperplane moves, too.

see sketch above 28).

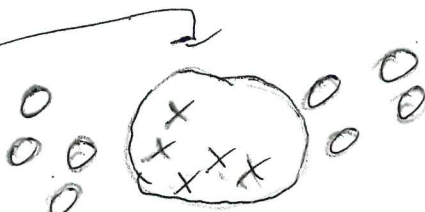


these points are not support vectors.

32] Using $K(\underline{x}, \underline{x}_i)$ leads to nonlinear decision boundaries if K is nonlinear. The kernel is a bivariate transformation.

There are $\binom{n}{2} = \frac{n(n-1)}{2}$ distinct pairs $(\underline{x}_i, \underline{x}_j)$ that are needed to estimate β_0 and the α_i .

An SVM with $\hat{f}(\underline{x}) = \hat{\beta}_0 + \sum_{i=1}^n \hat{\alpha}_i K(\underline{x}, \underline{x}_i)$ = ESP = ESP(\underline{x}) is a competitor for QDA while the SVM with $\hat{f}(\underline{x}) = \hat{\beta}_0 + \underline{\beta}^T \underline{x}$ is a competitor for LDA.

ex] boundaries  radial kernel

33] If \underline{x} is far from the \underline{x}_i , then

$\|x - x_i\|_2^2$ is large so

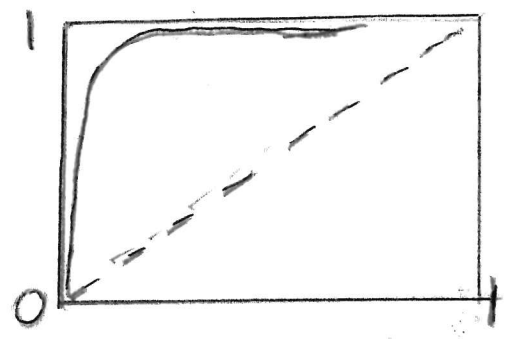
$k(x, x_i) = \exp(-\gamma \|x - x_i\|_2^2)$ is tiny, and x_i has almost no contribution in $f(x)$,
(Analogy KNN.)

(← from quality control)

34] A (receiver operating characteristic) ROC curve is used to evaluate binary classifiers. The overall performance is summarized by the area under the ROC curve (AUC). An ideal ROC curve is close to the top left corner of the plot, so the larger the AUC, the better the classifier.

		truth		
		-		
	-	true negative TN	false negative FN	total N^*
p predict		FP false positive	TP true positive	P^*
	total	N	P	

true positive rate



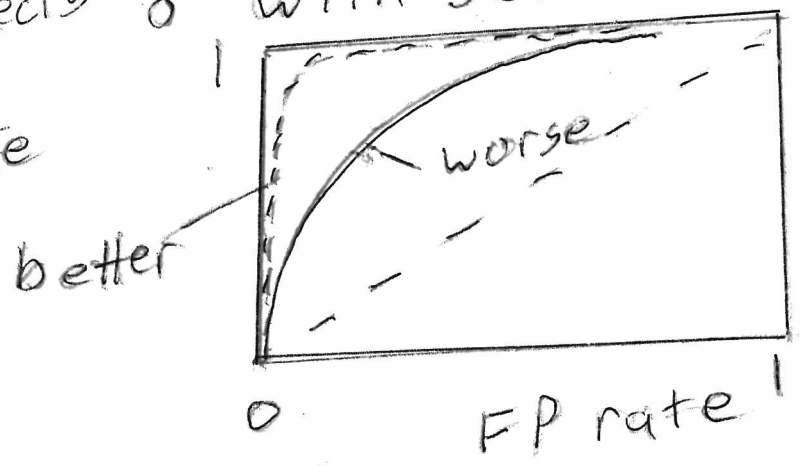
False positive rate

$$0 \leq AUC \leq 1$$

A classifier with $AUC = 0.5$ does no better than chance.

35] Plot varies γ for radial kernel and selects γ with best ROC curve

TP rate



better

FP rate

36] ROC from test data or validation data is better than ROC from training data.

37] The true positive rate is called the sensitivity and the false positive rate is $1 - \underline{\text{specificity}}$.

False positive rate = $\frac{FP}{N} \approx$ type I error, $1 - \text{specificity}$

True positive rate = $\frac{TP}{N} \approx$ $1 - \text{type II error}$, power, sensitivity, recall

Positive predicted value $\frac{TP}{P^*} = \text{precision} = 1 - \text{false discovery proportion}$

negative predicted value = $\frac{TN}{N}$.

denoted by $f_i(x)$. Then

$$\hat{y}(x) = d \text{ where } \hat{f}_d(x) = \max(\hat{f}_1(x), \dots, \hat{f}_G(x)).$$

ESPs

41] Rules 39] and 40] can be applied to any binary classifier, eg logistic regression.

Classification trees and related methods

1] Classification trees are worse than several ^{sums} alternatives, but are building blocks for a good HD classifier.

2] A classification tree is a "flexible" method for classification that produces a graph called a tree. Each branch has a label like $x_i > 7.56$
 $\underbrace{\hspace{10em}}_{x_i \text{ quantitative}}$

or $x_j = a, c$
 $\underbrace{\hspace{10em}}$

x_j a factor taking on levels a, b, c, d, e, f, say.
(categorical)

$x_1 = \text{term} = \text{Sentence length in years}$ HD 90

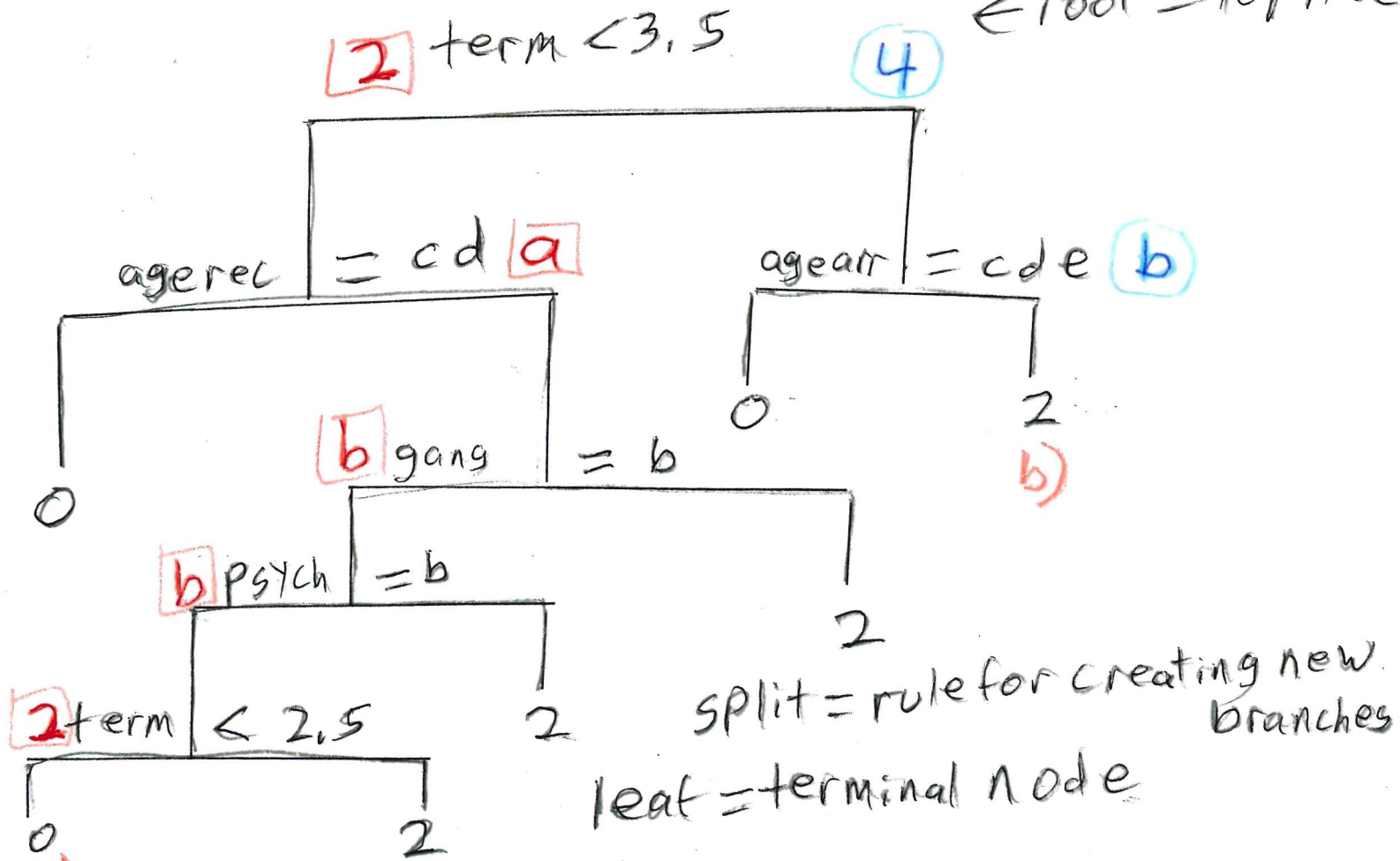
$x_2 = \text{agerec}$ (enter prison at)

a: 16-20, b: 21-26, c: 27-35, d: 30-35, e: ≥ 35

$x_3 = \text{agearr}$ (1st arrest) a: 0-17, b: 18-21, c: 22-29, d: 30-35, e: ≥ 36

$x_4 = \text{gang}$ a = gang activity, b = no gang activity
 $x_5 = \text{psych}$ a = mental illness, b = no mental illness

← root = top node



a) classify if $x_1 = \text{term} = 2$, $x_2 = \text{agerec} = a$,
 $x_4 = \text{gang} = b$, $x_5 = \text{psych} = b$

From the red values, $\hat{y} = 0$.

90.5

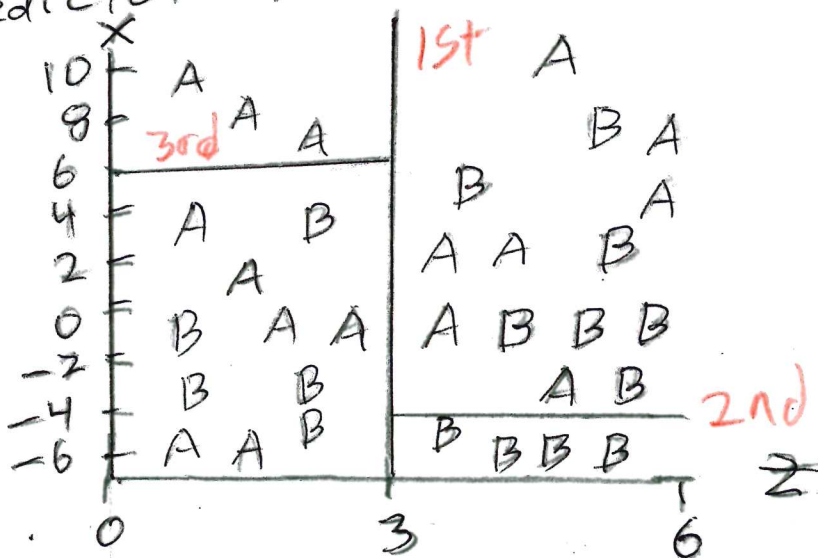
b) Classify if $x_1 = \text{term} = 4$, $x_3 = \text{age arr} = 6$

From the blue values, $\hat{y} = 2$.

4) Regression trees are similar, but the leaves give $\hat{y} \in \mathbb{R}$ rather than a classification label.

5) Trees that use recursive partitioning for classification and regression trees use the CART algorithm.

ex) $Y = A \text{ or } B$ (classification) with predictors X and Z



breaking up $X \mid Z$ is better than breaking up X if Y depends on X only through $X \mid Z$