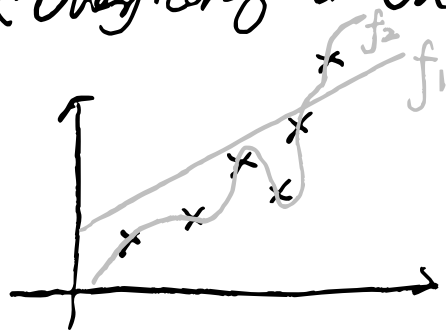


§ Introduction

$$\text{Statistical Learning} : Y = \underset{\uparrow}{f}(X) + \varepsilon$$

* Overfitting & Underfitting



f_1 or f_2 ?

* Training data & Testing data

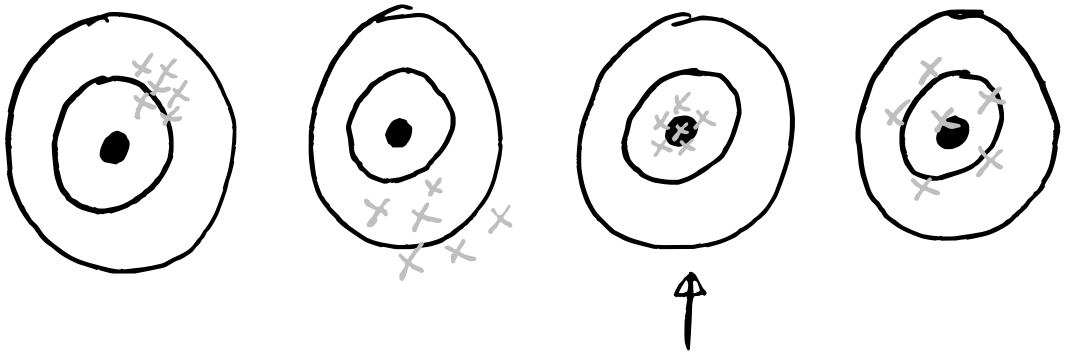
traditionally : 80% vs. 20% etc.
DL era : 99% vs. 1% etc.

* Evaluation of fit

$$MSE = \frac{1}{n} \sum_{i=1}^n (y_i - \hat{f}(x_i))^2 \quad \left\{ \begin{array}{l} \text{training} \\ \text{testing} \end{array} \right.$$

$$\mathbb{E} [(y_0 - f(x_0))^2] = \text{Var}(\hat{f}(x_0)) + [\text{bias}(\hat{f}(x_0))]^2 + \text{Var}(\varepsilon)$$

* Bias-Variance tradeoff



* Interpretation vs. Prediction

* Regression vs. Classification
(Y is continuous) (Y is discrete)

* Supervised vs. Unsupervised
(X, Y) (Only $X \rightarrow$ clustering)

§ Logistic Regression (LR)

$$y = \begin{cases} 1, & \text{if event occurs} \\ 0, & \text{otherwise} \end{cases}$$

$$\mathbb{P}(y=1) = p \rightarrow \mathbb{E}(y) = p$$

❖ How to model "p"?

$$f(p) = \beta_0 + \beta_1 x_1 + \dots + \beta_p x_p \in (-\infty, +\infty)$$

so, $f(p) = p$ is problematic

$$\rightarrow \text{logit}(p) = \log\left(\frac{p}{1-p}\right)$$

$$\rightarrow p = \frac{e^{\beta_0 + \beta_1 x_1 + \dots + \beta_p x_p}}{1 + e^{\beta_0 + \beta_1 x_1 + \dots + \beta_p x_p}}$$

$$\text{❖ Prediction: } \hat{y}_i = \begin{cases} 1, & \text{if } p_i \geq 0.5 \\ 0, & \text{otherwise} \end{cases}$$

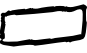

❖ Estimation: MLE + Newton-Raphson

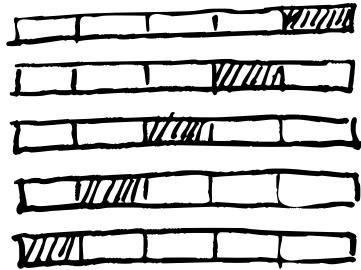
§ Cross validation, Evaluation measures

* Cross validation (CV)

{ waste less training data
| test on 100% of the data

e.g. 5-fold CV

{  training
|  testing



e.g. leave-one-out (LOO)

(n-1) training & 1 testing

* Evaluation measures

Confusion matrix &
$$F = \frac{2 \times TP}{2 \times TP + FP + FN}$$

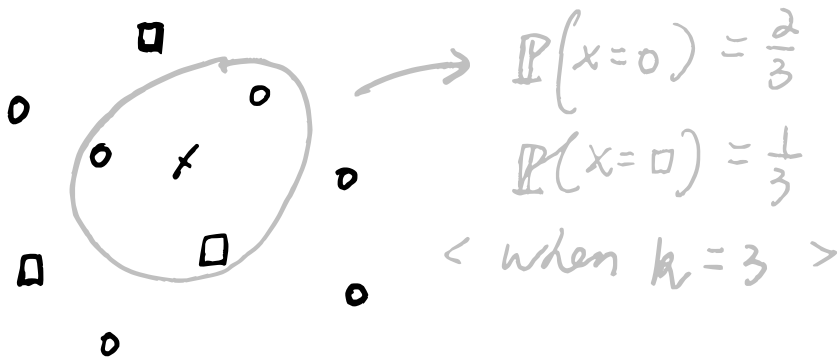
	Predicted True	Predicted False
Actual True	TP	FN
Actual False	FP	TN

§ k-Nearest Neighbor

$$\cdot x: P(Y=j' | X=x_0) = \frac{1}{K} \sum_{i \in N_0} \mathbb{1}\{y_i = j'\}$$

where $\begin{cases} x_0: \text{a new observation} \end{cases}$

$\begin{cases} N_0: \text{the set of the nearest } K \text{ observations} \end{cases}$



· x: How to break the classification tie?

- 1) at random
- 2) increase k until tie is broken
- 3) use 1NN as tie breaker

· x: Distances: Euclidean & Cosine

§ Naive Bayes

* Bayes' rule:
$$P(A|B) = \frac{P(B|A) \times P(A)}{P(B)}$$

*
$$P(y|x_1, \dots, x_n) = \frac{P(x_1|y)P(x_2|y) \dots P(x_n|y)P(y)}{P(x_1)P(x_2) \dots P(x_n)}$$

$$\rightarrow P(y|x_1, \dots, x_n) \propto P(y) \prod_{i=1}^n P(x_i|y)$$

* Laplace Smoothing

eliminate probabilities of 0 or 1

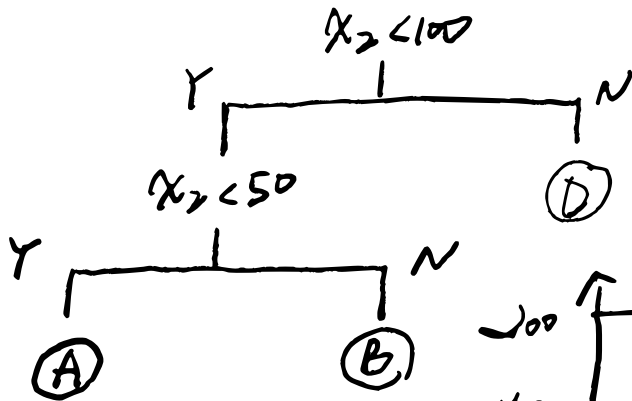
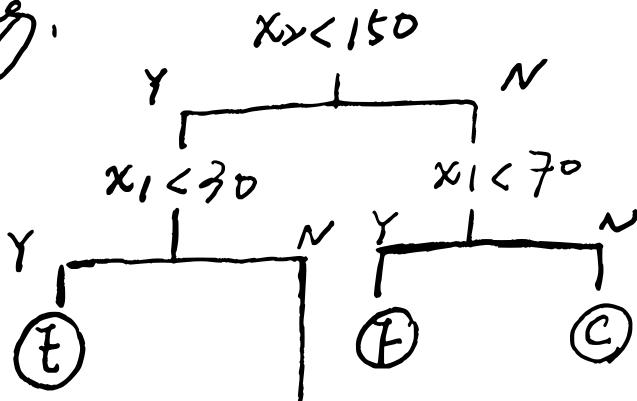
* Example: Stolen or not?

		Y	N	
1.° color	Red	<u>3/5</u>	<u>2/5</u>	Red SUV Domestic
	Yellow	2/5	3/5	
2.° type	Sports	4/5	2/5	$P(Y) = \frac{2}{5} \times \frac{1}{5} \times \frac{2}{5} \times \frac{1}{2}$
	SUV	<u>1/5</u>	<u>3/5</u>	
3.° origin	Domestic	<u>2/5</u>	<u>3/5</u>	$P(N) = \frac{2}{5} \times \frac{2}{5} \times \frac{2}{5} \times \frac{1}{2}$ $\Rightarrow \text{No!}$
	Imported	3/5	2/5	

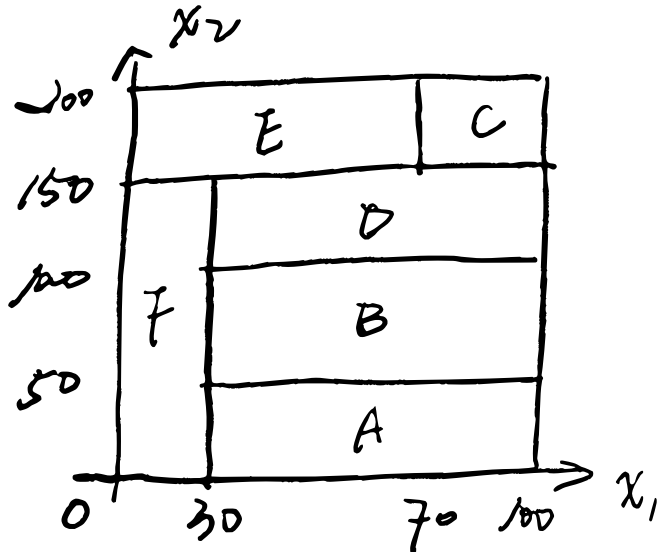
§ Trees & Random forests

* Classification And Regression Trees (CART)

eg.



"Rectangles"



* full tree vs. pruning

Criteria $\left\{ \begin{array}{l} \text{classification error} \\ \text{Gini index} \\ \text{cross-entropy} \end{array} \right.$

* Bootstrap: sample with replacement

* Bagging (Bootstrap Aggregation)

$$\rightarrow \hat{f}_{\text{bag}}(x) = \frac{1}{B} \sum_{b=1}^B \hat{f}_b(x)$$

\rightarrow reduce variance

$$\rightarrow \text{classifiers} = \hat{f}_{\text{bag}}(x) = \frac{1}{B} \sum_{b=1}^B \hat{f}_b(x)$$

* Random Forests "majority vote"

for $b = 1$ to B :

draw a bootstrap sample

grow a tree $\left\{ \begin{array}{l} \text{select } m \text{ out of } p \text{ variables} \\ \text{choose the best split} \end{array} \right.$

output ensemble

* Out-of-bag (OOB)

Each bootstrap sample on average only uses $2/3$ of the observations

- the remaining $1/3$ observations can be used to estimate the test error
- train / test split is not needed

§ Boosting

* Forward stagewise additive modelling

$$f_0(x) = 0$$

for $m = 1$ to M :

$$\min \sum_{i=1}^n \tilde{L}[f_{m-1}(x) + \beta_m b(x; \gamma_m)]$$

$$f_m(x) = f_{m-1}(x) + \beta_m b(x; \gamma_m)$$

* Exponential loss \rightarrow Adaboost

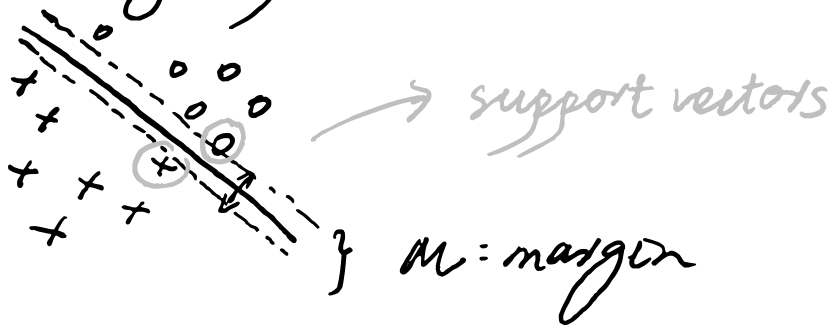
Binomial likelihood \rightarrow Logit-boost

* Regularization

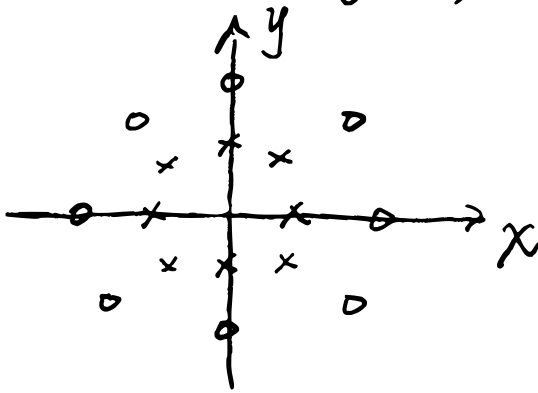
\rightarrow fit less aggressively to avoid overfitting
 { add penalty to loss function
 | shrinking

§ Support Vector Machines (SVM)

* Linearly Separable



* Not Linearly Separable



$$\begin{cases} x: x^2 + y^2 = 1 \\ 0: x^2 + y^2 = 2 \end{cases}$$

$$\rightarrow z = \sqrt{x^2 + y^2} \quad \begin{cases} x: z = 1 \\ 0: z = \sqrt{2} \end{cases}$$

"Linearly separable"