# naive bayes theorem

$X = (\text{Refund} = \text{No}, \text{Divorced}, \text{Income} = 120\text{K})$

| Tid | Refund | Marital Status | Taxable Income | Evade |
|-----|--------|----------------|----------------|-------|
| 1 | Yes | Single | 125K | No |
| 2 | No | Married | 100K | No |
| 3 | No | Single | 70K | No |
| 4 | Yes | Married | 120K | No |
| 5 | No | Divorced | 95K | Yes |
| 6 | No | Married | 60K | No |
| 7 | Yes | Divorced | 220K | No |
| 8 | No | Single | 85K | Yes |
| 9 | No | Married | 75K | No |
| 10 | No | Single | 90K | Yes |

**Using Bayes Theorem:**

$P(\text{Yes} \mid X) = \dfrac{P(X \mid \text{Yes}) P(\text{Yes})}{P(X)} \rightarrow 3/10$

$P(\text{No} \mid X) = \dfrac{P(X \mid \text{No}) P(\text{No})}{P(X)} \rightarrow 7/10$

since they are same, no need to compute

$\rightarrow P(\text{Refound}=\text{No}\mid\text{Yes}) \times P(\text{Divorced}\mid\text{Yes}) \times P(\text{income}=120k\mid\text{Yes})$

continuous attributes

$\hookrightarrow$ discretization: partition into bins

$\hookrightarrow$ probability density estimation: assume normal distri.
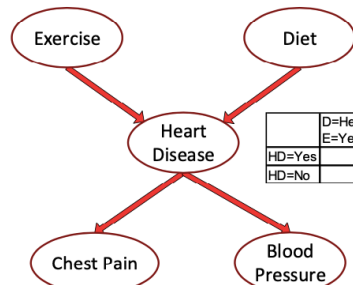calculate mean and std dev $\rightarrow$ calculate probability

✳ if one probability is zero all becomes zero
$\hookrightarrow$ laplace estimate $= \dfrac{n_c+1}{n+v}$   $v=$ total number of att val $x_i$ can take ( martial status $=3$ )

bayesian belief network=

| | Exercise=Yes | 0.7 |
|---|---|---|
| | Exercise=No | 0.3 |

| | Diet=Healthy | 0.25 |
|---|---|---|
| | Diet=Unhealthy | 0.75 |

Exercise   Diet

Heart Disease

|  | D=Healthy E=Yes | D=Healthy E=No | D=Unhealthy E=Yes | D=Unhealthy E=No |
|---|---|---|---|---|
| HD=Yes | 0.25 | 0.45 | 0.55 | 0.75 |
| HD=No | 0.75 | 0.55 | 0.45 | 0.25 |

Chest Pain   Blood Pressure

|  | HD=Yes | HD=No |
|---|---|---|
| CP=Yes | 0.8 | 0.01 |
| CP=No | 0.2 | 0.99 |

|  | HD=Yes | HD=No |
|---|---|---|
| BP=High | 0.85 | 0.2 |
| BP=Low | 0.15 | 0.8 |

$P(\text{HD}\mid E=\text{No}, D=\text{Yes}, CP=\text{Yes}, BP=\text{High})$

$\text{Yes} = P(\text{HD}=\text{Yes}\mid E=\text{No}, D=\text{Yes}) \times P(\text{Yes}\mid CP=\text{Yes}) \times P(\text{Yes}\mid \text{High})$

$\qquad 0.45 \qquad \times \qquad 0.8 \qquad \times \qquad 0.85$

$= 0.374$

$N_o = 0.0009$    $\text{Yes} > \text{No} \rightarrow$ classification $\rightarrow$ Yes

# artificial neural networks

update rule for learning weights $= w_j^{(k+1)} = w_j^{(k)} + \lambda (y_i - \hat{y}_i^{(k)}) x_{ij}$ (gradient descent in multilayer)
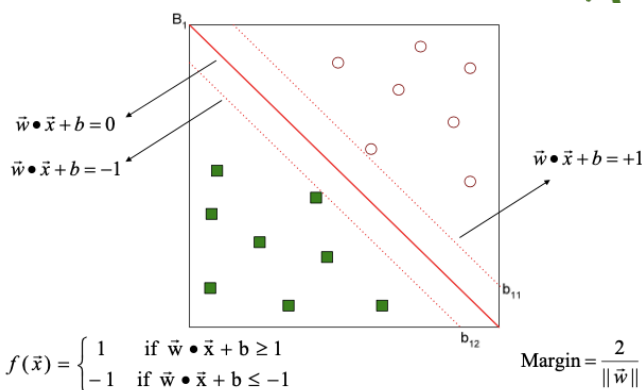
• if data is not linearly separable, perception not enough $\rightarrow$ multilayer

number of nodes in input/output layer $=$ • one for binary attribute/class

• $k$ or $\log_2 k$ nodes for $k$-class attribute

each combination for one of $k$

• if network is too large $\rightarrow$ overfitting

# support vector machines

$\vec{w} \bullet \vec{x} + b = 0$
$\vec{w} \bullet \vec{x} + b = -1$
$\vec{w} \bullet \vec{x} + b = +1$

$f(\vec{x}) = \begin{cases} 1 & \text{if } \vec{w} \bullet \vec{x} + b \geq 1 \\ -1 & \text{if } \vec{w} \bullet \vec{x} + b \leq -1 \end{cases}$   $\text{Margin} = \dfrac{2}{\|\vec{w}\|}$

• if not linearly separable
$\hookrightarrow$ use penalties in loss function
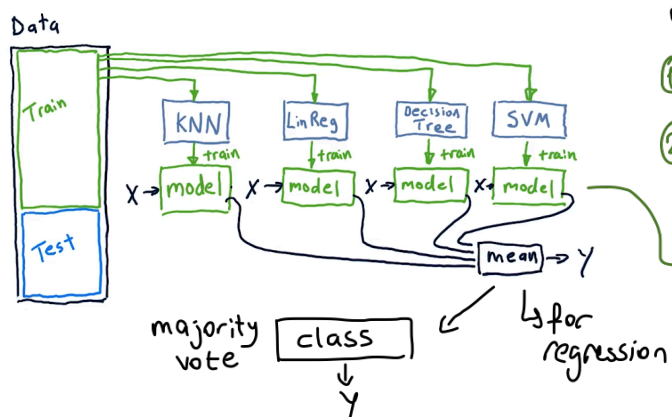$\hookrightarrow$ kernels $\rightarrow$ transformation to higher dimension space

# ensemble methods

works better than a single classifier if=
① all classifiers are independent
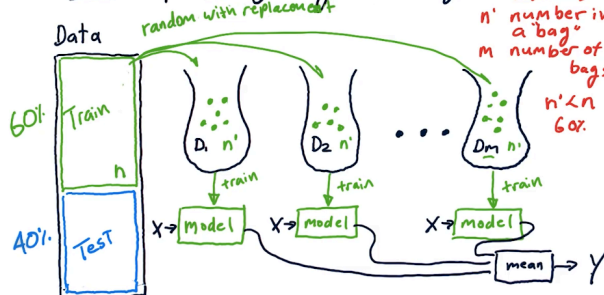② all classifiers perform better than random guessing
$\hookrightarrow$ ex: error rate $< 0.5$ for binary classification
• these models can be same algorithms with different parameters

Data
Train
KNN   LinReg   Decision Tree   SVM
↓train   ↓train   ↓train   ↓train
X→model   X→model   X→model   X→model
Test
mean → Y
$\hookrightarrow$ for regression
majority vote   class
↓
Y

# bagging=

Bootstrap aggregating - bagging

$n$ number of instances
$n'$ number in a bag
$m$ number of bags
$n' < n$
$6a$.

random with replacement

Data
60% Train
$D_1 n'$   $D_2 n'$   $\cdots$   $D_m n'$
↓train   ↓train   ↓train
40% Test
X→model   X→model   X→model
mean → Y

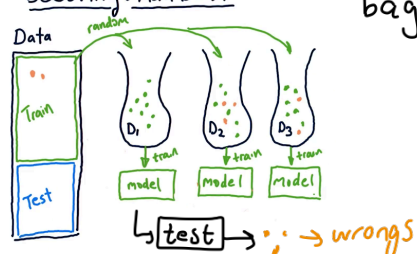$\rightarrow$ building ensemble classes using the same algorithm but training each learner on a different set of the data

# boosting=

Boosting: Ada Boost
Data
random
Train
$D_1$   $D_2$   $D_3$
↓train   ↓train   ↓train
model   model   model
Test
$\hookrightarrow$ test $\rightarrow$ → wrongs

bagging + train data for next model is weighted
$\hookrightarrow$ poorly classified in previous model has greater change to be chosen for next model to be trained