Intuition : we need a "fat" margin



thin margin          fat margin.

The hyperplane defining the boundary
is given by $\underline{\Theta}^T \underline{x} = 0$

Change of notations / habits: no $x_0 = 1$ ⚠

$\Theta_0 \longrightarrow b$

$\underline{\Theta}^T \underline{x} = 0 \quad \longrightarrow \quad \underline{a}^T \underline{x} + b = 0$



$\underline{x} \in$ plane : $\underline{a}^T \underline{x} + b = 0$

what is the distance between $\underline{x}^{(i)}$ and the
plane : $d^{(i)} = \dfrac{\underline{a}^T}{\|\underline{a}\|} (\underline{x}^{(i)} - \underline{x})$ ← same for any $\underline{x}$

$$d^{(i)} = \frac{1}{\|\underline{a}\|} \left( \underline{a}^T \underline{x}^{(i)} + b - \underline{a}^T \underline{x} - b \right)$$

$$d^{(i)} = \frac{1}{\|\underline{a}\|} \left( \underline{a}^T \underline{x}^{(i)} + b \right)$$

Here, we can constrain $|\underline{a}^T \underline{x}^{(i)} + b| \geq 1$
with $\underline{a}^T \underline{x}^{(i)} + b = 1$ for the closest point(s)
(we can do that because $\underline{a}$ and $b$ are determined
up to a multiplicative constant).
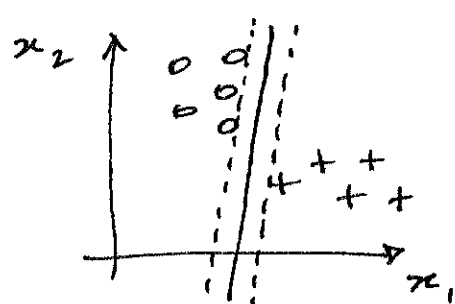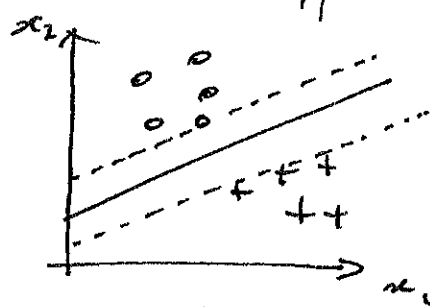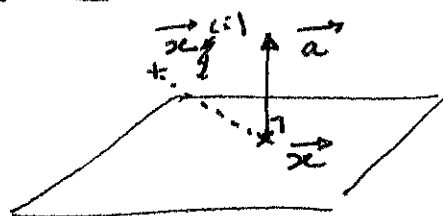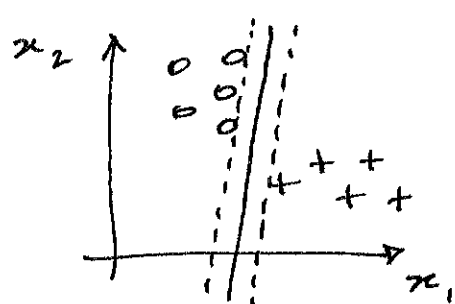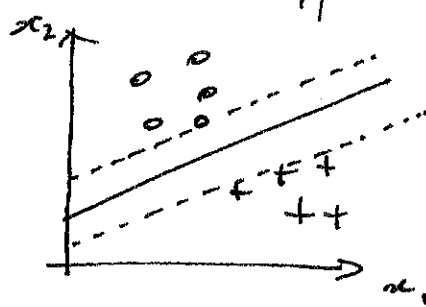
# Support vector machines

__Intuition__ : we need a "fat" margin



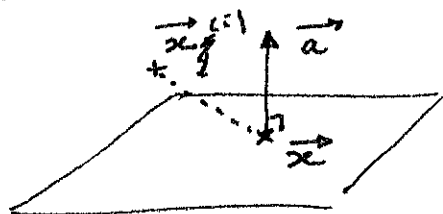thin margin                    fat margin.

The hyperplane defining the boundary
is given by $\underline{\theta}^T \underline{x} = 0$

Change of notations/habits: no $x_o = 1$ ⚠

$$\theta_o \longrightarrow b$$

$$\underline{\theta}^T \underline{x} = 0 \longrightarrow \underline{a}^T \underline{x} + b = 0$$



$$\vec{x} \in \text{plane} : \underline{a}^T \underline{x} + b = 0$$

what is the distance between $\underline{x}^{(i)}$ and the
plane : $d^{(i)} = \dfrac{\underline{a}^T}{\|\underline{a}\|} (\underline{x}^{(i)} - \underline{x})$ ⟵ same for any $\underline{x}$

$$d^{(i)} = \frac{1}{\|\underline{a}\|} \left( \underline{a}^T \underline{x}^{(i)} + b - \underline{a}^T \underline{x} - b \right)$$

$$d^{(i)} = \frac{1}{\|\underline{a}\|} \left( \underline{a}^T \underline{x}^{(i)} + b \right)$$

Here, we can constrain $\left| \underline{a}^T \underline{x}^{(i)} + b \right| \geq 1$
with $\underline{a}^T \underline{x}^{(i)} + b = 1$ for the closest point(s)
(we can do that because $\underline{a}$ and $b$ are determined
up to a multiplicative constant).

## Optimization problem

Find $\max \dfrac{1}{\|\underline{a}\|}$ subject to $\min\limits_{i=1..N}\left|\underline{a}^T\underline{x}^{(i)}+b\right|=1$

Or find $\min\|\underline{a}\|^2$ subject to ———

Be aware that $\underline{a}^T\underline{x}^{(i)}+b = \underline{\Theta}^T\underline{x}^{(i)} = h_\Theta(\underline{x}^{(i)})$

and we want $\begin{cases} h_\Theta(x^{(i)}) \geq 1 & \text{for } y^{(i)}=1 \\ h_\Theta(x^{(i)}) \leq -1 & \text{for } y^{(i)}=-1 \end{cases}$

$\left|\underline{a}^T\underline{x}^{(i)}+b\right| = y^{(i)}\left(\underline{a}^T\underline{x}^{(i)}+b\right)$

Find $\min\dfrac{\|\underline{a}\|^2}{2}$ subject to $y^{(i)}\left(\underline{a}^T\underline{x}^{(i)}+b\right)\geq 1$

for all $i$'s

## Quadratic programming

Lagrange formulation

$$\mathcal{L}(\underline{a}, b, \underline{\alpha}) = \frac{1}{2}\|\underline{a}\|^2 - \sum_{i=1}^{N}\underset{\geq 0}{\underline{\alpha_i}}\underbrace{\left(y^{(i)}\left(\underline{a}^t\underline{x}^{(i)}+b\right)-1\right)}_{\geq 0}$$

We want to find the $\min \mathcal{L}$ w.r.t. $\underline{a}$ and $b$

and maximize ~~w.r.t.~~ $\alpha_i \geq 0$ (KKT $\neq$ Lagrange multiplier)

$$\overrightarrow{\text{grad}}_{\underline{a}}\mathcal{L} = \underline{a} - \sum_{i=1}^{N}\alpha_i y^{(i)}\underline{x}^{(i)} = \underline{0} \quad\left.\begin{array}{l} \\ \end{array}\right\} \begin{array}{l}\min \\ \text{w.r.t} \\ \underline{a}, b\end{array}$$

$$\frac{\partial\mathcal{L}}{\partial b} = -\sum_{i=1}^{N}\alpha_i y^{(i)} = 0$$

Then $\underline{a} = \sum_{i=1}^{N}\alpha_i y^{(i)}\underline{x}^{(i)}$

$$\sum_{i=1}^{N}\alpha_i y^{(i)} = 0$$

Substituting in $\mathcal{L}$:

$$\mathcal{L} = \mathcal{L}(\alpha) = \frac{1}{2}\underbrace{\sum_{i=1}^{N}\sum_{j=1}^{N} y^{(i)} y^{(j)} \alpha_i \alpha_j \underline{x}^{(i)T}\underline{x}^{(j)}}_{\|\underline{a}\|^2}$$

$$- \sum_{i=1}^{N} \alpha_i \left(-1 + b\, y^{(i)} + y_i\, x^{(i)T} \sum_{j=1}^{N} \alpha_j y^{(j)} \underline{x}^{(j)}\right)$$

$$\mathcal{L} = \sum_{i=1}^{N} \alpha_i - \frac{1}{2}\sum_{i=1}^{N}\sum_{j=1}^{N} y^{(i)} y^{(j)} \alpha_i \alpha_j \underline{x}^{(i)T}\underline{x}^{(j)}$$

we want to find $\max(\mathcal{L})$ w.r.t $\alpha$.

with $\sum_{i=1}^{N} \alpha_i y^{(i)} = 0$ and $\alpha_i \geqslant 0$

This is a quadratic optimization problem
for which there is plenty of methods...
One of them is a penalty method (iterative):
at each step, we maximize:

$$\mathcal{L} - \sigma_k \sum_i \min(\alpha_i, 0) - \sigma_k \left(\sum_i \alpha_i y^{(i)}\right)^2$$

with $\sigma_k$ increasing at each step (10-fold, for instance)

## Soft margin

If data is not linearly separable, we
use a "hinge loss" function

$$\max\left(0, 1 - y^{(i)}\left(\underline{a}^T \underline{x}^{(i)} + b\right)\right)$$

this function is $0$ if the output $y^{(i)} = \pm 1$
is correctly predicted, otherwise it is prop. to the
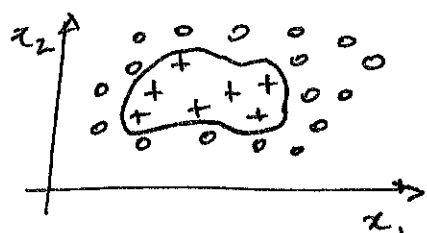distance to the margin.

We want to minimize:

$$\lambda \| \underline{a} \|^2 + \frac{1}{n} \sum_{i=1}^{n} \max \left( 0, 1 - y^{(i)} \left( \underline{a}^T \underline{x}^{(i)} + b \right) \right)$$

$\lambda$ is the trade-off between the margin size and ensuring that $\underline{x}^{(i)}$'s lie on the correct side.

In the limit $\lambda \longrightarrow 0$, it converges towards a hard-margin SVM.
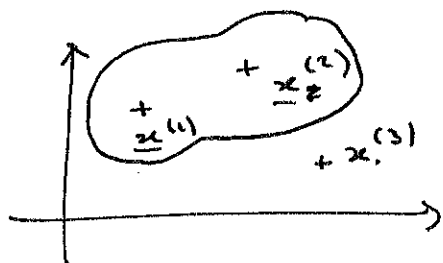
## SVM with kernels



non-linear decision boundary

$$\theta_0 + \theta_1 x_1 + \theta_2 x_2 + \theta_3 x_1^2 + \theta_4 x_1 x_2 + \theta_5 x_2^2 + \cdots \geq 0$$

$$f_1 = x_1, \quad f_2 = x_2, \quad f_3 = x_1^2 \quad : \text{ features}$$

Is there a better choice of features?
| an alternative

Gaussian kernel : $\quad f_1 = \exp \left( - \frac{\| \underline{x} - \underline{x}^{(1)} \|^2}{2 \sigma^2} \right)$



Intuition if

$$\theta_0 = -0.5 \qquad \theta_1 = \theta_2 = 1$$
$$\theta_3 = 0$$

Q. how many landmarks?
how to choose them?

$\longrightarrow$ all $\underline{x}^{(i)}$'s

$$\underline{x}^{(i)} \longrightarrow \begin{bmatrix} f_1 = \text{sim}(\underline{x}^{(i)}, \underline{x}^{(1)}) \\ \vdots \\ f_i = 1 \\ \vdots \\ f_N \end{bmatrix} \underbrace{\phantom{xxx}}_{\underline{f}^{(i)}}$$

$$\underline{x}^{(i)} \in \mathbb{R}^{n\,(+1)}$$
$$\underline{f}^{(i)} \in \mathbb{R}^{N\,(+1)}$$

# Logistic vs. SVMs

$n$ : number of features

$N$ : ———————— training examples

- $n \geqslant N$ ——→ Logistic or SVM without kernel

- $n$ small
    - $N$ intermediate ($\in [10 \ 10000]$) ——→ Gaussian kernel
    - $N$ large ($> 50000$) ——→ Logistic or SVM without kernel

      create more features and ten
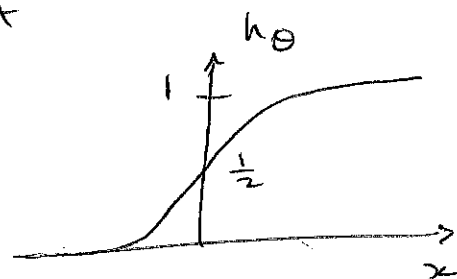
——→ neural network is an alternative in all cases.

# Logistic regression
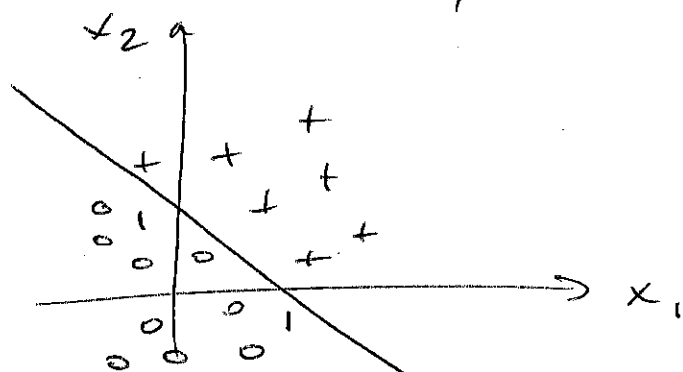
Useful for classification



hypothesis : $h_\theta = \dfrac{1}{1 + e^{-\theta^T \underline{x}}}$



(logistic function)

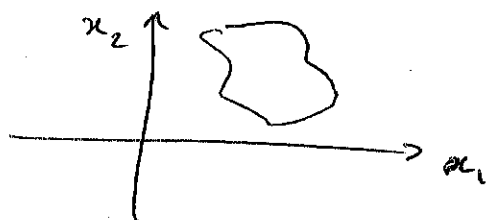Interpretation : $h_\theta = p(y = 1 \mid \underline{x} ; \theta)$

Decision boundary



$x_1 + x_2 \geqslant 1$
$x_1 + x_2 \leq 1$

$\left.\phantom{\begin{matrix}a\\b\end{matrix}}\right)$ $\underline{\theta}^T \underline{x} = -1 + x_1 + x_2$
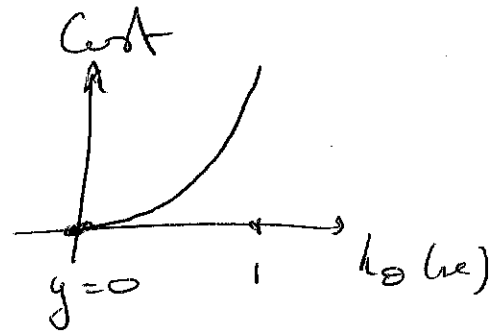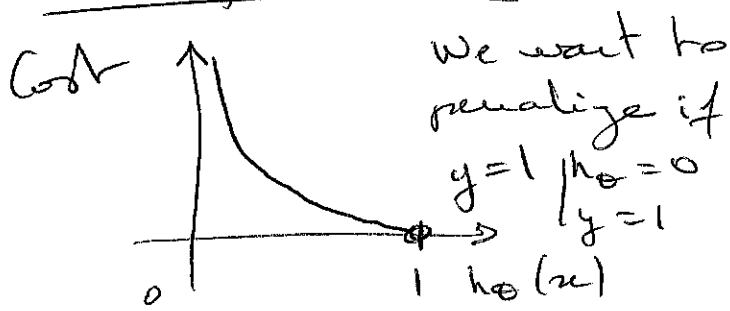
Works with larger order polynomials if
$x_2 = x_1^2$ for instance.



$\underline{x} = \begin{pmatrix} x_1 \\ x_2 \\ x_1^2 \\ x_2^2 \\ \vdots \end{pmatrix}$

# Cost function

Cost

We want to penalize if
$y=1$ $h_\theta = 0$
$y=1$

Cost

$y=0$ $h_\theta(x)$

$$
\text{Cost}: \begin{cases} -\log(h_\theta) & \text{if } y=1 \\ -\log(1-h_\theta) & \text{if } y=0 \end{cases}
$$

Generalization:

$$
J(\underline{\Theta}) = -\frac{1}{N} \sum_{i=1}^{N} \left[ y^{(i)} \log\left(h_\theta(\underline{x}^{(i)})\right) + (1-y^{(i)}) \log\left(1-h_\theta(\underline{x}^{(i)})\right) \right]
$$

(Analogy with Shannon entropy).

## Simplification

$$
h_\theta(\underline{x}) = \frac{1}{1+e^{-\underline{\Theta}^T \underline{x}}}
$$

$$
\log h_\theta = -\log\left(1+e^{-\underline{\Theta}^T \underline{x}}\right)
$$

$$
\frac{\partial}{\partial \Theta_j} \log h_\theta = -\frac{1}{1+e^{-\underline{\Theta}^T x}}\left(-x_j e^{-\underline{\Theta}^T x}\right) = \frac{x_j}{1+e^{\Theta^T x}}
$$

$$
\log(1-h_\theta) = +\log \frac{e^{-\Theta^T x}}{1+e^{-\Theta^T x}} = -\log\left(\frac{1}{e^{-\Theta^T x}}+1\right)
$$

$$
\frac{\partial}{\partial \Theta_j} \log(1-h_\theta) = -\frac{1}{1+e^{\Theta^T x}}\left(x_j e^{+\Theta^T x}\right)
$$

$y^{(i)} 2b.$

$$\frac{\partial}{\partial \theta_i} \left[ y \log h_\theta + (1-y) \log (1-h_\theta) \right]$$

$$= x_j \left[ \frac{1}{1+e^{\theta^T x}} \left( y \mathrel{\overline{\neq}} (1-y) e^{\theta^T x} \right) \right)$$
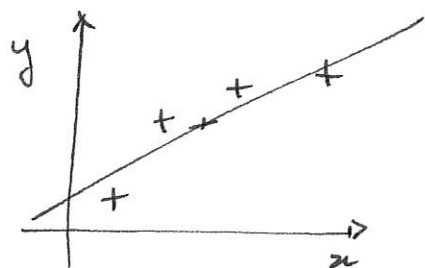
$$= x_j \left[ \frac{-1}{1+e^{-\theta^T x}} + y \right]$$

$$\Longrightarrow \theta_j := \theta_j - \alpha \frac{1}{N} \sum_{i=1}^{N} \underbrace{\left( h_\theta(x^{(i)}) - y^{(i)} \right)}_{\frac{1}{1+e^{-\theta^T x}}} x_j^{(i)}$$
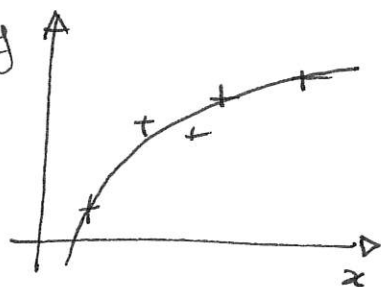
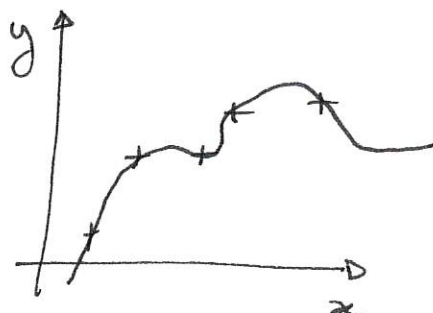Same formula as linear regression!

# Regularization (Pb of overfitting)

$\Theta_o + \Theta_1 x$

"underfitting"

$\Theta_o + \Theta_1 x + \Theta_2 x^2$

"feels 'just'"

$\Theta_o + \Theta_1 x + \Theta_2 x^2 + \Theta_3 x^3 + \Theta_4 x^4$

"overfitting"

Same if very large number of $x_j$'s
Same with logistic regression...

* **Intuition** $\longrightarrow$ we want to keep $\Theta_j$'s small

* **Solution :** new cost function

$$J(\underline{\Theta}) = \frac{1}{2N}\left[ \sum_{i=1}^{N} (h_\Theta(x^{(i)}) - y^{(i)})^2 + d \underbrace{\sum_{j=1}^{n} \Theta_j^2}_{} \right]$$

new term
⚠ starts at $j = 1$
does not include $\Theta_o$

$d$ : regularisation parameter

* **Application to linear regression**

$$\frac{\partial J}{\partial \Theta_o} = \frac{1}{N} \sum_{i=1}^{N} (h_\Theta(x^{(i)}) - y^{(i)}) \underbrace{x_o^{(i)}}_{=1} \quad \left\} \begin{array}{l} \text{idem} \\ \text{sans} \\ \text{regularis}^\circ \end{array}\right.$$

$$\frac{\partial J}{\partial \Theta_j} = \frac{1}{N} \sum_{i=1}^{N} (h_\Theta(x^{(i)}) - y^{(i)}) x_j^{(i)} + \frac{d}{N} \Theta_j$$

# Gradient descent

$$\Theta_j := \Theta_j - \alpha \left\{ \frac{1}{N} \sum_{i=1}^{N} \left( h_\theta(x^{(i)}) - y^{(i)} \right) x_j^{(i)} + \frac{\lambda}{N} \Theta_j \right\}$$

$$\Theta_j := \Theta_j \underbrace{\left( 1 - \frac{\alpha \lambda}{N} \right)}_{\in [0, 1]} - \alpha \frac{1}{N} \sum_{i=1}^{N} \left( h_\theta(x^{(i)}) - y^{(i)} \right) x_j^{(i)}$$

typically 0.99

# Normal equation

$$\Theta = \left( \underline{\underline{X}}^T \underline{\underline{X}}^{-1} + \lambda \begin{bmatrix} 0, & 0 \\ 0 & \searrow, 1 \end{bmatrix} \right)^{-1} \underline{\underline{X}}^T \underline{y}$$

Can make the first matrix invertible when the number of examples is less than the number of features.

# Logistic equation

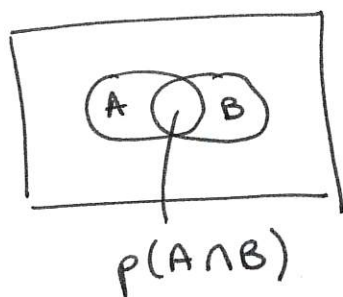Same gradient descent as linear regression.

# Naive Bayes

Classifier based on conditionnal probability

$$p(C_k | \underline{x})$$

Probability of being of class $C_k$ among $K$ possible classes, given $\underline{x}$

Bayes theorem:

$$p(C_k | \underline{x}) = \frac{p(C_k) \overset{prior}{} p(\underline{x} | C_k)}{p(\underline{x})} \quad \substack{\text{support that} \\ C_k \text{ provides} \\ \text{for } \underline{x}}$$

"Demonstration":



$$p(A \cap B) = p(A) p(B|A) = p(B) p(A|B)$$

$$\Rightarrow p(A|B) = p(A) \frac{p(B|A)}{p(B)}$$

$$p(A \cap B)$$

$$p(C_k , \underline{x}) = p(C_k | \underline{x}) p(\underline{x}) = p(C_k) p(\underline{x} | C_k)$$

$$p(C_k , \underline{x}) = p(x_1, \ldots, x_n, C_k)$$

$$= p(x_1 | x_2, \ldots, x_n, C_k) p(x_2, \ldots x_n, C_k)$$

$$= p(x_1 | x_2, \ldots x_n, C_k) \cdots p(x_n | C_k) p(C_k)$$

("naive" independent assumption

$$= p(x_1 | C_k) \cdots p(x_n | C_k) p(C_k)$$

$$p(C_k | \underline{x}) = \frac{p(C_k, \underline{x})}{p(\underline{x})}$$

$$= \frac{p(C_k)}{p(\underline{x})} \prod_{i=1}^{n} p(x_i | C_k)$$

This can be used to construct a classifier

$$y = \underset{k \in [1..K]}{\text{argmax}} \; p(C_k) \prod_{i=1}^{n} p(x_i | C_k)$$

(we can drop $p(\underline{x})$ which is a constant)

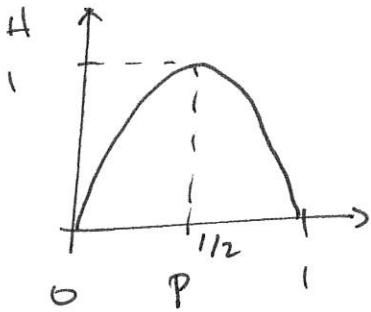$p(C_k)$ and $p(x_i | C_k)$ can be learned from training dataset.

# Decision trees

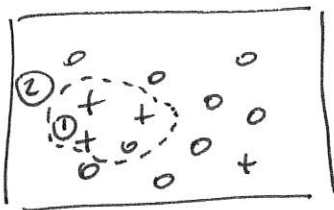* Crash course on information theory (Shannon)

Shannon entropy (in bits per symbol):

$$H = - \sum_i p_i \log_2 (p_i)$$

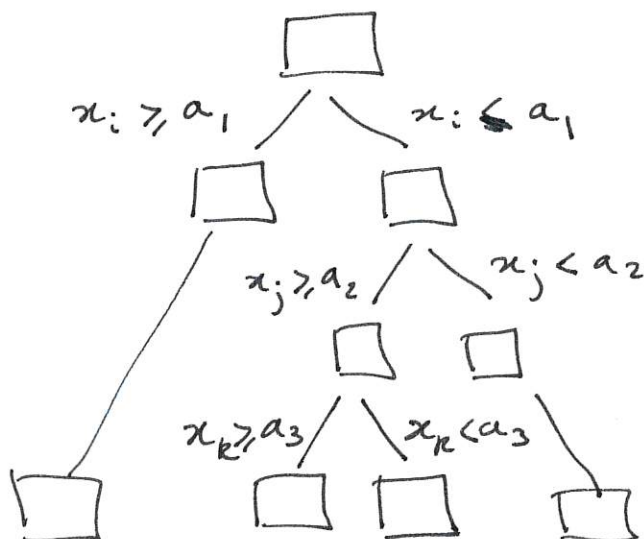binary entropy:

$$H(p) = -p \log_2 p - (1-p) \log_2 (1-p)$$

$H$ : entropy of the total set $N$

$$H_{split} = \frac{N_1}{N} H_1 + \frac{N_2}{N} H_2$$

If $H_1 = H_2 = H \implies H_{split} = H$ no gain
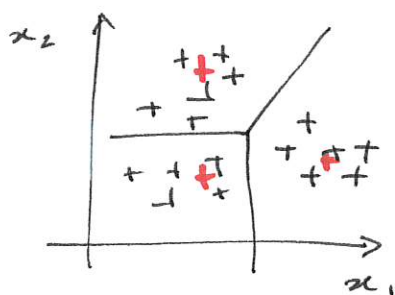
$H_1 = H_2 = 0$ (when only one label in each subset). $\implies H_{split} = H$

The idea is to maximize the entropy loss. (or information gain) at each split.

$x_i \geq a_1$ / \ $x_i \leq a_1$

$x_j \geq a_2$ / \ $x_j < a_2$

$x_k \geq a_3$ / \ $x_k < a_3$

# k-means

We want to partition the space into $k$ cells (Voronoi cells defined by their centers)

Pb computationally difficult (NP-hard)
k-means is an efficient heuristic algorithm.

The objective of k-means clustering is to find

$$\text{argmin} \sum_{j=1}^{K} \sum_{x \in S_j} \| \underline{x} - \underline{\mu}_j \|^2$$

where $\underline{\mu}_i$ is the mean of points in $S_i$

## Standard algorithm

- Randomly choose $K$ centroids $\underline{\mu}_j \in \mathbb{R}^n$
- For the $N$ $\underline{x}^{(i)}$ find the cluster $c^{(i)} \in \{1 \ldots K\}$ to which it belongs by computing $\| x^{(i)} - \underline{\mu}_j \|^2$
- Calculate the average of points assigned to each cluster
- update cluster centroid positions

repeat

Possibility to measure efficiency with objective function

$$J(\underline{\mu}_1 \ldots \underline{\mu}_K) = \frac{1}{N} \sum_{i=1}^{N} \| \underline{x}^{(i)} - \underline{\mu}_{c^{(i)}} \|^2$$

index of cluster

→ Random initialization and then select result with smallest J

→ Choosing number of cluster

$J$



elbow method

1  2  3  4  5      K