Landau symbol ("big O"): $O(\cdot)$ for order of magnitude of complexity

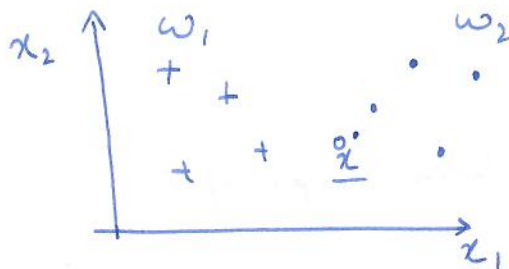$$O(N^x) \rightarrow \text{exponential computational complexity}$$

$$\lim_{N \to \infty} \frac{O(N^x)}{N^x} = A \neq 0, \infty$$

i.e. $O(N^2)$ means $AN^2$ is the biggest contributor to the complexity

$$AN^2 + O(N), \quad A \neq 0 = \text{constant}$$

## 4.2.2  k NN (k - Nearest Neighbours)

Idea: Follow the majority of k nearest training samples.



$k \in \mathbb{N}$

$k = 1$: look for the closest neighbour.
In this case $W_2$

$k = 3$: We have 2 from $W_1$ and 1 from $W_2$.
So $\hat{y}(\underline{x}) = W_1$

$k = 1$ leads to overfitting and should not be used

## 4.3 Bayes Plug-in

We use Bayes decision theory here Ch 2.

Bayesian decision theory:

$$P(\underline{x} \mid w_j), \ P(w_j) \longrightarrow MBR/MAP/ML/NP \text{ decision rule}$$

In practice
$$P(\underline{x} \mid w_j), P(w_j) \to \text{unknown}$$

Idea: $P(\underline{x} \mid w_j)$ known except for some unknown parameters $\underline{\nu}_j$

i.e. $P(\underline{x} \mid w_j ; \underline{\nu}_j)$

$$\downarrow \atop \%$$

eg. $\underline{x} \mid w_j \sim N(\underline{\mu}_j, \ \underline{\underline{\varsigma}}_j)$

$$\underline{\nu}_j = \{ \underline{\mu}_j, \underline{\underline{\varsigma}}_j \}$$

Training:

Estimate $\underline{\nu}_j, P(w_j)$ from $N_j$ training samples of class $w_j$

$$\Rightarrow \ \hat{\underline{\nu}}_j, \ \hat{P}(w_j)$$

Classification

use $P(\underline{x} \mid w_j, \hat{\underline{\nu}}_j)$ and $\hat{P}(w_j)$ in Bayesian decision

How to ~~use~~ estimate $\underline{\nu}_j, P(w_j)$?

## 4.3.1 ML parameter estimation

A crash course : see SASP.

- given sample $\underline{x}$ of pdf $p(\underline{x}; \underline{v})$

  are known        not known

- ML estimation:

$$\hat{\underline{v}}_{ML}(\underline{x}) = \arg\max_{\underline{v}} p(\underline{x}; \underline{v}) \longrightarrow \text{likelihood}$$

$$= \arg\max_{\underline{v}} \ln p(\underline{x}; \underline{v}) \longrightarrow \text{log-likelihood}$$

- Necessary condition:

$$\nabla \ln p(\underline{x}; \underline{v})\Big|_{\underline{v} = \hat{\underline{v}}_{ML}} \overset{!}{=} \underline{0}$$

$$\nabla = \begin{bmatrix} \frac{\partial}{\partial v_1} \\ \vdots \\ \frac{\partial}{\partial v_M} \end{bmatrix} \longrightarrow \text{gradient vector}$$

$\Rightarrow$ a linear / non-linear eq. system

$\Rightarrow$ analytical / numerical sol.

E4.4 ML estimation for Gaussian distributions.

$\underline{x}_n \in \mathbb{R}^d$ $(1 \leq n \leq N)$, iid $N(\underline{\mu}, \underline{C})$

$$p(\underline{x}_n; \underline{v}) = \frac{1}{(2\pi)^{d/2} \sqrt{|\underline{C}|}} \exp\left(-\frac{1}{2}(\underline{x}_n - \underline{\mu})^T \underline{C}^{-1}(\underline{x}_n - \underline{\mu})\right)$$

$$p(\underline{x}_1, \dots \underline{x}_N; \underline{v}) \overset{iid}{=} \prod_{i=1}^{N} p(\underline{x}_i; \underline{v})$$

log-likelihood:

$$\mathcal{L}(\underline{\gamma}) = \ln p(\underline{x}_1, \cdots \underline{x}_N; \underline{\gamma})$$

$$\cong \text{constant} - \frac{N}{2} \ln |\underline{\underline{C}}| - \frac{1}{2} \sum_{n=1}^{N} (\underline{x}_n - \underline{\mu})^T \underline{\underline{C}}^{-1} (\underline{x}_n - \underline{\mu})$$

$$\underline{\gamma} = \{\underline{\mu}, \underline{\underline{C}}\} \text{ contains } d \text{ elements of } \underline{\mu} \; \square$$

$$\text{and} \qquad \frac{d(d+1)}{2} \text{ „ } \qquad \text{ „ } \quad \underline{\underline{C}} = \underline{\underline{C}}^T \; \triangledown$$

$$\underline{\nabla} \mathcal{L}(\underline{\gamma}) \overset{!}{=} \underline{0}$$

4-21

4.29.1

$$\hat{\underline{\mu}} = \frac{1}{N} \sum_{n=1}^{N} \underline{x}_n \qquad \text{for } \underline{\mu} = E(\underline{x})$$

$$\hat{\underline{\underline{C}}} = \frac{1}{N} \sum_{n=1}^{N} (\underline{x}_n - \hat{\underline{\mu}})(\underline{x} - \hat{\underline{\mu}})^T \quad \text{for } \underline{\underline{C}} = E\left[(\underline{x}-\underline{\mu})(\underline{x}-\underline{\mu})^T\right]$$

slide 4.29

4.29.1

E4.5 cont : Solution

Method of Lagrange multiplier :

$$\underset{P_j, \lambda}{\max} \quad \tilde{\mathcal{L}}(P_1 \cdots P_c, \lambda) = \mathcal{L}(P_1 \cdots P_c) + \lambda (P_1 + \cdots + P_c - 1)$$

$$\left. \frac{\partial \tilde{\mathcal{L}}}{\partial P_j} \right|_{P = \hat{P}_j} = \frac{N_j}{P_j} + \lambda \overset{!}{=} 0 \quad \Rightarrow \quad \hat{P}_j = -\frac{1}{\lambda} N_j \sim N_j$$

$$\sum_{j=1}^{C} \hat{P}_j = 1 \quad -\frac{1}{\lambda} \sum_{j=1}^{C} N_j = \frac{-N}{\lambda} = 1 \quad \Rightarrow \quad \hat{P}_j = \frac{N_j}{N}$$

## 4.3.2    Gaussian classifier

Idea: • Bayes plug-in method
      • Assumption: Gaussian likelihood

$$\underline{x} \mid \omega_j \sim N(\underline{\mu}_j, \underline{\underline{C}}_j)$$

ML estimate of $\underline{\mu}_j$, $\underline{\underline{C}}_j$, $P(\omega_j)$ according to E 4.4, E4.5 from training examples.

slide 4-30

### 4-30.1

Comparision to nearest mean (4.2.1)
• nearest mean:

$$\min_j \quad D_{mahe} = (\underline{x} - \underline{\mu}_j)^T \underline{\underline{C}}_j^{-1} (\underline{x} - \underline{\mu}_j)$$

or $\quad \max_j \quad -D_{mahe} = -(\underline{x} - \underline{\mu}_j)^T \underline{\underline{C}}_j^{-1} (\underline{x} - \underline{\mu}_j)$

• Gaussian classifier $\neq$ MAP decision rule:

$$P(\underline{x} \mid \omega_j ; \tilde{\underline{\nu}}_j) \sim N(\hat{\underline{\mu}}_j, \hat{\underline{\underline{C}}}_j)$$

$$\max_j \quad \ln P(\underline{x} \mid \omega_j ; \underline{\nu}_j) \cdot \hat{P}(\omega_j)$$

$$= -\frac{1}{2}(\underline{x} - \underline{\mu}_j)^T \underline{\underline{C}}_j^{-1} (\underline{x} - \underline{\mu}_j) - \frac{1}{2}\ln |\hat{\underline{\underline{C}}}_j| + \ln \hat{P}_j$$
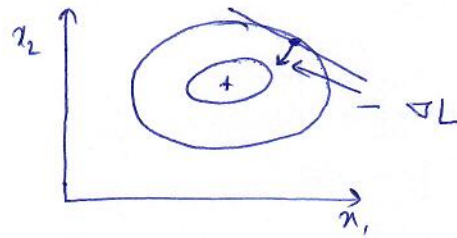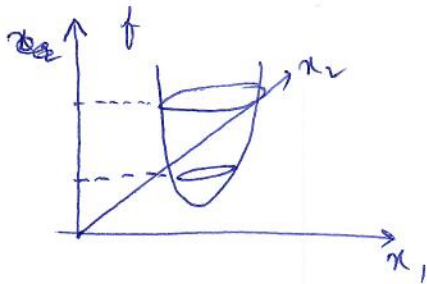
# ch5 : Advanced Optimization Techniques

## 5.1 Challenges in Optimization

Slide 5.1

### S15.4.1

2D visualization $f(\underline{x}) = f(x_1, x_2)$



Contour lines of $f(\underline{x})$
$$= \{ \underline{x} \mid f(\underline{x}) = \text{constant} \}$$

Slide 5-5