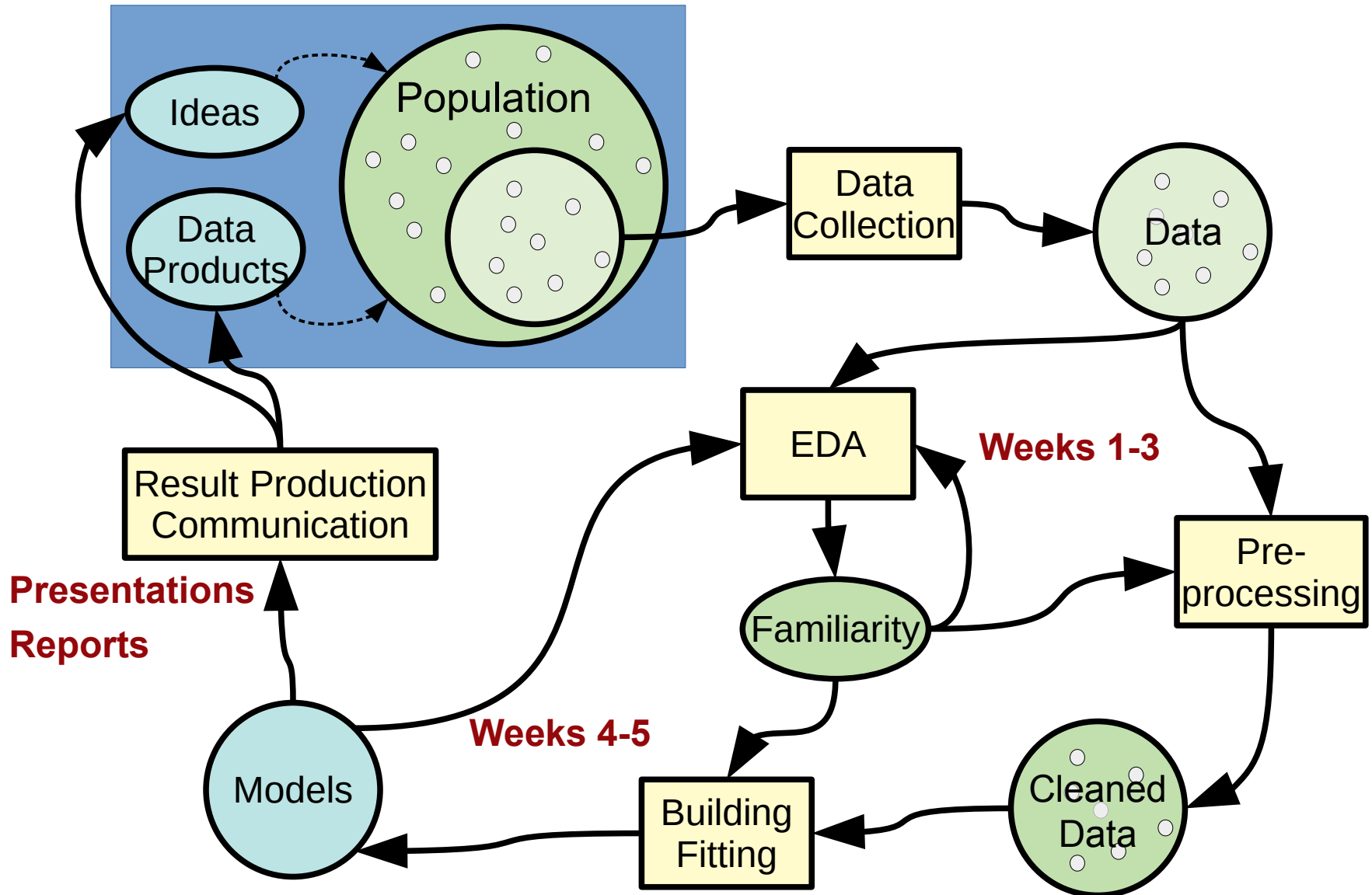


# Course Content



# Application of numerical data description

- Cheap  $\rightarrow$  apply plenty (robust and non-robust)
- Report in a meaningful way: take selection bias into account
- Understand structure of data

## Do this on training set

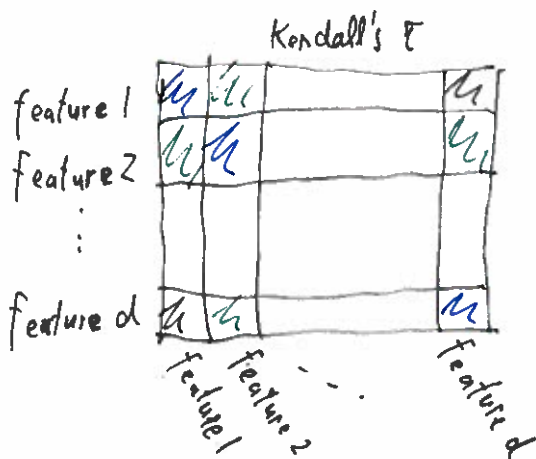
First step in data analysis:

Set aside held-out set

Then understand structure of remaining data.

Pairwise measures:

- Visualise matrix



- Visualise histogram



- Minimise approximation error

$$\begin{aligned} & \text{minimise } E \left[ \left\| \underline{x} - \sum_{i=1}^k \underline{w}_i \underline{w}_i^T \underline{x} \right\|^2 \right] \\ & \underline{w}_1, \dots, \underline{w}_k \\ & \text{s.t. } \|\underline{w}_i\| = 1, i=1, \dots, k \text{ and } \underline{w}_i^T \underline{w}_j = 0 \end{aligned}$$

- Approximate data matrix  $X$

$$\begin{aligned} & \text{minimise } \sum_{i,j} (X - M)_{ij}^2 \\ & M \\ & \text{s.t. } \text{rank}(M) = k \end{aligned}$$

$$M = \sum_{i=1}^k s_i \underline{u}_i \underline{v}_i^T$$

$$\underline{w}_i = \underline{u}_i, Z_i = s_i \underline{v}_i^T$$

- Approximate sample covariance

$$\begin{aligned} & \text{minimise } \left\| \frac{1}{n} X X^T - M \right\|_F \\ & M \\ & \text{s.t. } \text{rank}(M) = k \text{ and } M^T = M \end{aligned}$$

$$\hat{\Sigma} = \sum_{i=1}^n \underline{u}_i \frac{s_i^2}{n} \underline{u}_i^T$$

- Approximation Gram matrix

$$\begin{aligned} & \text{minimise } \|X^T X - M\|_F \\ & M \\ & \text{s.t. } \text{rank}(M) = k \text{ and } M^T = M \end{aligned}$$

$$\hat{G} = \sum_{i=1}^n \underline{v}_i s_i^2 \underline{v}_i^T$$

$$\underline{Z} = \sqrt{\tilde{\Lambda}_k} V_k^T$$

## - Kernel PCA:

Define  $\tilde{G}$  with  $(\tilde{G})_{ij} = \phi(\underline{x}_i)^T \phi(\underline{x}_j)$   
 $= k(\underline{x}_i, \underline{x}_j)$

## - Classical MDS:

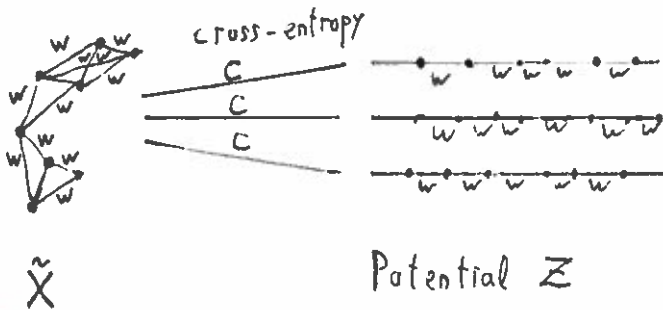
$$\begin{aligned} Z = \underset{M}{\text{minimise}} \quad & \|(-\frac{1}{2} C_n \Delta C_n) - M^T M\|_F \\ \text{s.t.} \quad & \text{rank}(M^T M) = k \end{aligned}$$

## - Isomap:

Classical MDS, but on geodesic distances



## - UMAP:



$$\underset{Z}{\text{minimise}} \quad C(W^{(\tilde{X})}, W^{(Z)})$$

# Estimating generalisation performance with hyperparameter selection

1. Split  $D$  into  $D^{\text{hyp}}$  and  $D^{\text{test}}$  (e.g. 80:20)

Lock away  $D^{\text{test}}$  until step 4

2. Optimise hyperparameters  $\underline{\lambda}$  on  $D^{\text{hyp}}$ :

held-out

Split  $D^{\text{hyp}}$  into  
 $D^{\text{train}}$  and  $D^{\text{val}}$

$$\hat{\underline{\lambda}} = \underset{\underline{\lambda}}{\operatorname{argmin}} \text{PL}(\underline{\lambda})$$

where

$$\text{PL}(\underline{\lambda}) = \hat{J}(\hat{h}_{\underline{\lambda}}, D^{\text{val}})$$

$$\text{and } \hat{h}_{\underline{\lambda}} = \mathcal{A}_{\underline{\lambda}}(D^{\text{train}})$$

cross validation

Split  $D^{\text{hyp}}$  into  
 $K$  folds and  
compute CV for  
each  $\underline{\lambda}$

$$\hat{\underline{\lambda}} = \underset{\underline{\lambda}}{\operatorname{argmin}} \text{EPL}(\underline{\lambda})$$

where

$$\text{EPL}(\underline{\lambda}) = \text{CV}_{\underline{\lambda}} = \hat{J}(\mathcal{A}_{\underline{\lambda}})$$

3. Reestimate  $\underline{\theta}$  using  $\hat{\underline{\lambda}}$ :  $\hat{h}^{\text{hyp}} = \mathcal{A}_{\hat{\underline{\lambda}}}(D^{\text{hyp}})$

4. Estimate prediction loss:  $\hat{J}(\hat{h}^{\text{hyp}}; D^{\text{test}})$

5. Reestimate  $\hat{h}$  on all data:  $\hat{h} = \mathcal{A}_{\hat{\underline{\lambda}}}(D)$

Do not estimate prediction loss on all data!

"Two-times held-out"

"CV and held-out"

2020 L15 (4)

# Loss functions

## - Regression:

square loss, absolute loss, Huber loss

## - Classification

- loss matrix, zero-one loss

→ true positive rate, true negative rate,  
false positive rate, false negative rate

- ROC: true positive rate vs. false positive rate

- differentiable loss functions  
based on margin  $y h(\underline{x})$