## Unsupervised Learning

- recap:  supervised learning :- we possess a training set of features and desired response

$$TS = \{ (x_i, y_i) \}_{i=1}^{N} \Rightarrow \text{learn the mapping } Y = f(x)$$

- basic properties of the mapping ("structure") of is already given in TS, e.g. we know that we are looking for a regression plane (least squares regression) or a decision boundary (classification)

- unsupervised learning: - no pre-defined problem structure, TS without response

$$TS \{ x_i \}_{i=1}^{N}$$

- why?  • labeled data (with Y's) is very expensive, unlabelled data is cheap
  ⇒ what can we do without labels?

  • in the initial state of a research project, we may not even know what to look for ≘ "data mining"

- ultimately, unsupervised learning is equivalent to <u>artificial general intelligence</u> (AGI)
  ⇒ "AI scientist" : can define research goals, data collection procedures / experiments, models, optimization / learning alg. all on it's own

- for now:  search for certain types of useful "default" structures in the data
  ≙ pattern discovery, where types of patterns are implicitly defined by the model & learning alg.

## Unsupervised Learning

- recap:   supervised learning :- we possess a training set of features and desired response

$$TS = \{ (x_i, y_i) \}_{i=1}^{N} \quad \Rightarrow \text{ learn the mapping } Y = f(x)$$

- basic properties of the mapping ("structure") of is already
  given in TS, e.g. we know that we are looking for
  a regression plane (least-squares regression) or a
  decision boundary (classification)

- unsupervised learning: - no pre-defined problem structure, TS without response

$$TS \{ x_i \}_{i=1}^{N}$$

- why?   • labeled data (with Y's) is very expensive, unlabelled data is cheap
  ⇒ what can we do without labels?
  • in the initial state of a research project, we may not even know what to
  look for ≙ "data mining"

- ultimately, unsupervised learning is equivalent to <u>artificial general intelligence</u> (AGI)
  ⇒ "AI scientist" : can define research goals, data collection procedures/experiments,
  models, optimization/learning alg. all on it's own

- for now:   search for certain types of useful "default" structures in the data
  ≙ pattern discovery, where types of patterns are implicitly defined
  by the model & learning alg.

## Principal Component Analysis (PCA)

- find linear variable transform, such that the dimension is reduced without loosing information: $z = \varphi(x) = X \cdot V^T \quad \in \mathbb{R}^{D'} \quad x \in$

$$x_i \in \mathbb{R}^D \quad , \quad z_i \in \mathbb{R}^{D'} \quad D' < D \quad (\text{ideally } D' \ll D)$$

all the essential properties of $x$ are still in $z$, but the "noise" is gone

- what's the optimal projection matrix $V$ ?

- extreme case : $D' = 1 \qquad z_i$ is a single scalar feature

embed new feature into original feature space: define a 1-dimensional line in $\mathbb{R}^D$

$$\tilde{x}_i = \underset{\uparrow \text{ mean}}{\mu} + \underset{\uparrow \text{ direction}}{V} \cdot z_i$$



the $\tilde{x}_i$ should be approximations of $x_i$

$$\Rightarrow \hat{\mu}, \hat{v}, \{\hat{z}_i\} = \underset{\mu, v, \{z_i\}}{\arg\min} \; \sum_i (x_i - \tilde{x}_i)^2 \qquad s.t. \; v v^T = 1$$

$$= \underset{\mu, v, \{z_i\}, \lambda}{\arg\min} \; \sum_i (x_i - \mu - v z_i)^2 + \lambda (v v^T - 1)$$

$$\frac{\partial \text{Loss}}{\partial \mu} = -2 \sum_i (x_i - \mu - v z_i) \overset{!}{=} 0 \quad \Rightarrow \quad \mu = \underbrace{\frac{1}{N} \sum_i x_i}_{= \bar{x}} - v \cdot \underbrace{\frac{1}{N} \sum_i z_i}_{= \bar{z}}$$

$\Rightarrow \mu$ is not unique, can choose $\bar{z}$ freely $\Rightarrow \bar{z} = 0$ (new feature is centered)
if the original are also centered, $\bar{x} = 0 \Rightarrow \mu = 0$

$\Rightarrow$ assume data is centered, $\bar{x} = 0$, set $\mu = 0$ $\Rightarrow$ $\bar{z} = 0$ is also centered

simplified opt. problem

$$\hat{v}, \{\hat{z}_i\} = \arg\min_{v, \{z_i\}, \lambda} \sum_i (X_i - v z_i)^2 + \lambda (v v^T - 1)$$

$$\frac{d \text{loss}}{d z_i} = -2 \sum_i (x_i - v z_i) \cdot v^T \stackrel{!}{=} 0 \qquad x_i v^T = \underbrace{v v^T}_{=1} z_i \qquad z_i = \frac{x_i v^T}{v v^T}$$

$$\boxed{z_i = x_i v^T}$$

simplify again:

$$\hat{v}, \hat{\lambda} = \arg\min_{v, \lambda} \sum_i (X_i - v \cdot x_i v^T)^2 + \lambda (v v^T - 1)$$

expand the square:

$$(x_i - (x_i v^T) \cdot v)^2 = (x_i - (x_i v^T) v)(x_i - (x_i v^T) v)^T$$

$$= x_i x_i^T - (x_i v^T)(v x_i^T) - (x_i v^T)(x_i v^T) + (x_i v^T) \underbrace{v v^T}_{=1} (x_i v^T)$$

$$= x_i x_i^T - (x_i v^T)(v x_i^T)$$

$$\underbrace{= x_i x_i^T}$$

independent of the solution $v$ $\Rightarrow$ drop

$$\hat{v}, \hat{\lambda} = \arg\max_{v, \lambda} \sum_i v x_i^T x_i v^T + \lambda (v v^T - 1)$$

$$\boxed{\hat{v}, \hat{\lambda} = \arg\max_{v, \lambda} v S v^T + \lambda (v v^T - 1)}$$

$$S = \sum_i x_i^T x_i = X^T X$$

scatter matrix

$\Rightarrow$ assume data is centered, $\bar{x} = 0$, set $\mu = 0$ $\Rightarrow$ $\bar{z} = 0$ is also centered

simplified opt. problem

$$\hat{v}, \{\hat{z}_i\} = \arg\min_{v, \{z_i\}, \lambda} \sum_i (x_i - v z_i)^2 + \lambda(v v^T - 1)$$

$$\frac{d\text{loss}}{dz_i} = -2 \sum_k (x_i - v z_i) \cdot v^T \doteq 0 \qquad x_i v^T = \underbrace{v v^T}_{=1} z_i \qquad z_i = \frac{x_i v^T}{v v^T}$$

$$\boxed{z_i = x_i v^T}$$

simplify again:

$$\hat{v}, \hat{\lambda} = \arg\min_{v, \lambda} \sum_i (x_i - v \cdot x_i v^T)^2 + \lambda(v v^T - 1)$$

expand the square:

$$(x_i - (x_i v^T) \cdot v)^2 = (x_i - (x_i v^T) v)(x_i - (x_i v^T) v)^T$$

$$= x_i x_i^T - (x_i v^T)(v x_i^T) - (x_i v^T)(x_i v^T) + (x_i v^T) \underbrace{v v^T}_{=1} (x_i v^T)$$

$$= x_i x_i^T - (x_i v^T)(v x_i^T)$$

$$\underbrace{\underbrace{x_i x_i^T}}_{\text{independent of the solution } v \Rightarrow \text{drop}}$$

$$\hat{v}, \hat{\lambda} = \arg\max_{v, \lambda} \sum_i v x_i^T x_i v^T - \lambda(v v^T - 1)$$

$$\boxed{\hat{v}, \hat{\lambda} = \arg\max_{v, \lambda} v S v^T + \lambda \frac{(v v^T - 1)}{(1 - v v^T)}}$$

$$S = \sum_i x_i^T x_i = X^T X$$
scatter matrix

$$\frac{dLoss}{dv} = 2Sv^T - 2\lambda v^T \overset{!}{=} 0 \qquad \boxed{Sv^T = \lambda v^T}$$

$\Rightarrow$ $v^T$ must be an eigenvector of the scatter matrix, which one?

$$\hat{\lambda} = \underset{\lambda}{\arg\max} \quad v \underbrace{S v^T}_{\lambda v^T} + \lambda(1 - v v^T)$$

$$= \underset{\lambda}{\arg\max} \quad \lambda v v^T + \lambda - \lambda v v^T = \underset{\lambda}{\arg\max} \quad \lambda$$

$\Rightarrow$ $\boxed{\lambda \text{ must be the biggest eigenvalue of } S}$ $\boxed{v \text{ the corresponding eigenvector}}$

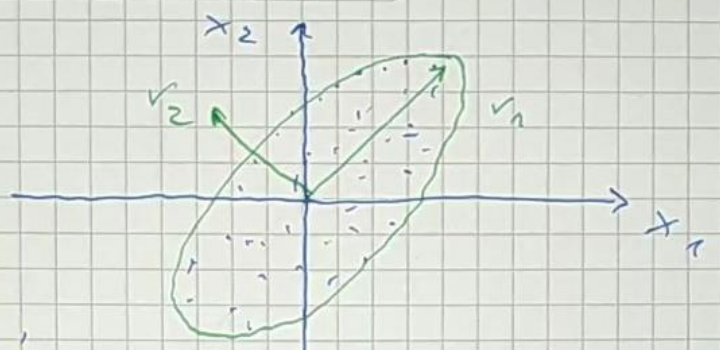can repeat the same procedure with residuals $x_i' = x_i - \hat{x_i} = x_i - v \cdot z_{i1}$

$\Rightarrow$ get the second-best feature $z_{i2}$, and so on until we have $D'$ features

## PCA algorithm

① center $X$ and compute scatter matrix $S = X^T X$

② compute eigen decomposition $S = V \Lambda V^T$ and sort by decreasing eigenvalue: $\Lambda = \begin{pmatrix} \lambda_1 & & 0 \\ & \ddots & \\ 0 & & \lambda_D \end{pmatrix}$ such that $\lambda_1 \geq \lambda_2 \geq \dots \geq \lambda_D \geq 0$

③ choose new feature dimension $D' < D$

④ compute new features: $z_{ij} = x_i v_j^T$ $\qquad$ $\underset{\text{eigenvectors}}{} $ for $j = 1 \dots D'$

<u>intuitive interpretation</u> :

• eigen-decomposition of $S$
approximates the data by an ellipse
( see QDA )



• eigenvectors are the ellipse axes ,

$\Rightarrow$ the most important new feature $z_{i,1}$ is along the longest axis of ellipse
( because it's largest eigenvalue )

$\Rightarrow$ PCA selects the coordinate system , where the data vary most

$\hat{=}$ the variance of the new features $Var_i(z_{ij})$ is maximized

$\hat{=}$ "most informative" means that the values of the features vary a lot

$(\Rightarrow$ choosing good units for the original features $X_i$ is important , otherwise
you can make the data vary arbitrarily $\Rightarrow$ PCA meaningless )

e.g. scale original features such that $Var_i(X_{ij}) = 1$

• important property of ① PCA: the new features $z_j$ are pair-wise <u>uncorrelated</u>

<u>Variants of PCA</u>

• keep the linear variable transform, but change the optimality criterion and/or constraints

   - Independent Component Analysis (ICA) : $z_j$ are pair-wise <u>independent</u> $\Big\}$ later

   - Non-negative matrix factorization (NMF): $z_j, v \geq 0$

• define non-linear PCA via the kernel trick

   - kernel - PCA