# Unsupervised Learning

what is it?
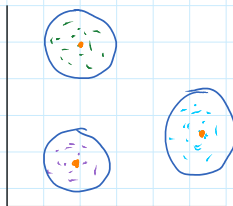- no labels
- discover patterns in data /structures
- clustering

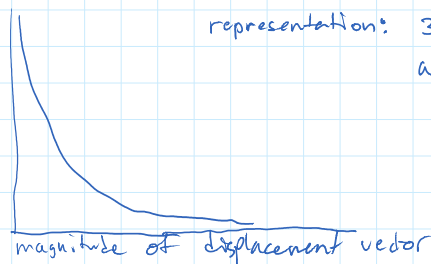Just a bunch of vectors, What to do with them?

What are our Goals?
- understand data
- sanity check
- remove uninformative part of data  (dimensionality reduction) (compression)
- learning features    (and  feature  selection)
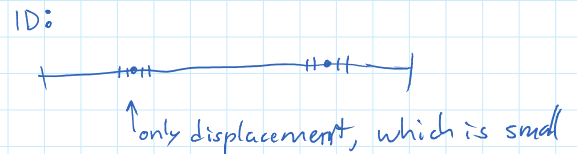- visualization

clustering

K-means: adds points to clusters

representation: 3 centers,
                assignment    (lossy compression)

1D:

magnitude of displacement vector

only displacement, which is small

- new points are likely in the cluster

# Principle Component Analysis (PCA)

- dimensionality reduction method

    approximate data set in some dimension $m$, $m < d$

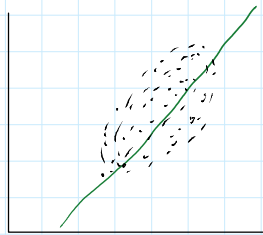PCA is a linear model

$$v_n = W x_n$$

$v \in \mathbb{R}^M$ &larr; new dimensionality

$W \in \mathbb{R}^{M \times D}$ &larr; original dimensionality

$x \in \mathbb{R}^D$

$X$ : whole data set $\mathbb{R}^{D \times N}$ &larr; amount of data

&larr; one dimensional subspace  (project)

approximate representation

"maximize variance" of projected values

         - equivalent to small L2 error / fitting

$M = 1$  (for now)

$u \in \mathbb{R}^D$   &larr; learning $u$

$$u^T u = 1$$

variance of projected data:

$$= \frac{1}{N} \sum_{n=1}^{N} \left( u^T x_n \right)^2 \quad \text{(sum of squares)}$$

$$= u^T S u, \quad \text{where } S \text{ is covariance matrix } XX^T$$

$$\max_u \quad u^T S u + \lambda \left( 1 - u^T u \right) \quad \text{lagrange multiplier: ie, } 1 = u^T u$$

$$S u - \lambda u = 0$$
$$S u = \lambda u \quad \text{(eigenvector of covariance matrix of the data)}$$

         - but which eigenvector?

variance is now:
$$u^T S u = u^T \lambda u = \lambda u^T u = \lambda$$
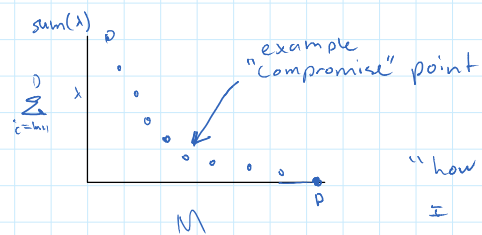
         - pick biggest eigenvalue

$m > 1$?  repeat with remaining variance

Singular Value Decomposition  (SVD)

$$X = U \Sigma V^\top$$

eigenvalues in $U$, pick biggest

another way of doing eigenvalue decomposition

How to pick  m?



$M$

"how much  of the variance am
I  throwing  away?"

Probabilistic  PCA          — generative, higher dimension

$$X_n = W_{z_n} + \varepsilon \quad \leftarrow \text{noise}$$

$\leftarrow$ latent space

$$z_n \in \mathbb{R}^m, \quad x_n \in \mathbb{R}^D \qquad W \in \mathbb{R}^{D \times M}$$

$$z_n \sim N(0, \mathbb{I}) \quad \text{iid}$$

$$\varepsilon \sim N(0, \sigma^2 \mathbb{I}) \quad \text{iid} \qquad (\text{isotropic noise})$$

$$p(x \mid z) \sim N(Wz, \sigma^2 \mathbb{I})$$

$$p(x) \sim N(\mu, C)$$

$$\mu = E[x] = E[Wz + \varepsilon] = 0$$

$$C = \text{cov}[x] = E\left[ (Wz + \varepsilon)(Wz + \varepsilon)^\top \right]$$

$$= E[Wzz^\top W^\top] + E[\varepsilon \varepsilon^\top]$$

$$= WW^\top + \sigma^2 \mathbb{I}$$

MLE is PCA

$W$ only appears as $WW^\top$

$\overset{\text{orthogonal}}{\swarrow}$

$$\tilde{W} = WR$$

$$\tilde{W}\tilde{W} = WRR^\top W^\top = WW^\top$$

could be considered undesireable, any orthogonal
transformation  on  $W$  gives the same model

Cocktail Party Problem
     microphones listening to sum of 5 voices
     decompose signal to 5 sources

## Blind Source Seperation

treat as iid data

$$Z = \begin{bmatrix} \text{person 1} \\ \text{person 2} \\ \vdots \end{bmatrix}$$

infer distances, dissambuguate people

observed
↓
$$x = Wz$$

$t =$ # time samples
$n_s =$ # of speakers
$n_m =$ # of microphones

# unknowns: $n_s n_m + t n_s$
# equations: $t n_m$

if $t n_m >$ # unknowns
we can hopefully solve it

PCA: disaster!  z means 5 people
          PCA could be multiplied by some orthogonal matrix
          the problem has a "correct answer"

## Independant Component Analysis  (ICA)

- linear model (like PCA)
- NOT Gaussian prior on z
- $p(z) = \prod_{m=1}^{M} p(z_n)$

## Factor Analysis

- like Probabilistic PCA, but: instead of variance $\sigma^2 I$ variance,
                            each element can have its own.