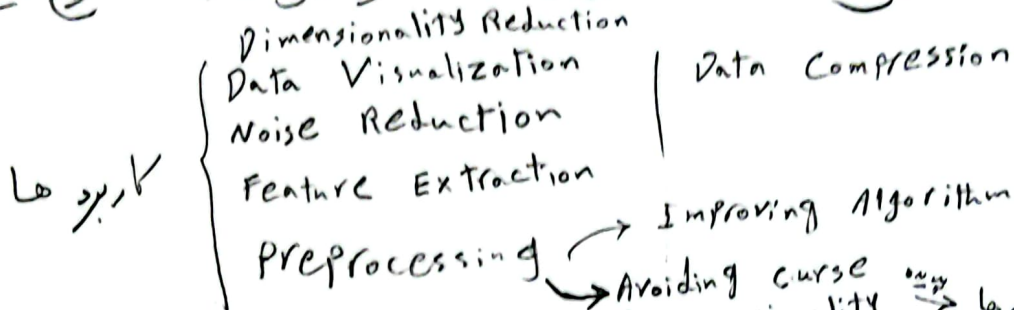


# >> Principal Component Analysis (PCA)

PCA چیست؟ (؟)

تکنیک آماری ← کاهش ابعاد + فشردگی داده + استخراج ویژگی‌ها



$$X = [x_1, x_2, \dots, x_d]$$

$$XW = Z$$

$$Z = [z_1, z_2, \dots, z_k]$$

1. Standardize d-dim dataset

$$\frac{x - \mu}{\sigma} = X_{std}$$

2. make covariance matrix

$$\text{Cov}(X_{std}) = \frac{1}{n-1} X_{std}^T X_{std}$$

3. Calculate EigenValues & EigenVectors

$$\text{Cov}(X_{std}) \cdot V = \lambda V$$

4. Sort eigenVals desc & related eigenVect

$$\text{EigenVal\_EigenVect\_Pairs} = \{(\lambda_1, v_1), (\lambda_2, v_2), \dots, (\lambda_m, v_m)\}$$

5. Form the projection matrix

$$W = [v_1, v_2, \dots, v_k]$$

6. Transform the data

$$X_{reduced} = X_{std} \cdot W$$

ویژگی‌های مهم ← Unsupervised  
(؟)

الگوریتم‌های clustering  
k-nearest neighbors  
در ابعاد بالا بهتر هستند.  
چون در ابعاد بالا داده‌ها پراکنده می‌شوند.

# مراحل

## >> Dimensionality Reduction in Machine Learning

- روش ها {
1. Principal Component Analysis (PCA)
  2. Linear Discriminant Analysis (LDA)
  3. t-Distributed Stochastic Neighbor Embedding (t-SNE)
  4. Autoencoders
  5. Factor Analysis
  6. Kernel PCA

### ▷ Linear Discriminant Analysis (LDA)

هدف: Project the data onto a lower-dimensional space with good class separability → to {

- avoid overfitting
- reduce computational costs

Supervised ← مسئله دسته بندی

# گام ها

0 - استاندارد سازی داده ها

1 - compute mean vectors

$$\mu_i = \frac{1}{N_i} \sum_{x \in X_i} x$$

برای هر کلاس

2 - compute within-class scatter matrix ( $S_w$ ):

$$S_w = \sum_{i=1}^K \sum_{x \in X_i} (x - \mu_i)(x - \mu_i)^T$$

3 - Between-class scatter matrix ( $S_b$ ):

$$S_b = \sum_{i=1}^K N_i (\mu_i - \mu)(\mu_i - \mu)^T$$

4 - Compute eigenval & eigenvectors  $\Rightarrow S_w^{-1} S_b v = \lambda v$

5 - Select the top  $K$  eigenvectors  $\rightarrow w = [v_1, v_2, \dots, v_K]$

6 - Transform the data  $\rightarrow Y = X W$

# >> t-distributed stochastic neighbor embedding (t-SNE)

کاربرد اصلی } dimensionality reduction  
data visualization

← از نوع غیر خطی کاهش ابعاد

# چگونه کار می کند ؟

## 1. Pairwise Similarities

- در فضای ابعاد بالا، شباهت ها و زوجهی نقطه داده های ورودی را اندازه گیری می کند
- استفاده از توزیع گاوسی برای محاسبه شباهت ها (احتمالات شرطی) بین نقاط، جایگزین نقاط شبیه دارای احتمالات بالا هستند

## 2. Low-Dimensional Mapping

- در فضای ابعاد کم، الگوریتم مذکور شباهت ها را مدل می کند (با استفاده Student's t-distribution with one degree of freedom)
- باعث می شود که "crowding Problem" رخ ندهد (برخورد نقاط در ابعاد پایین تر)

## 3. Cost Function

- t-SNE سعی می کند تفاوت بین ۲ توزیع احتمال (بعد بالا و بعد کوچک) را با تابع هزینه ای بنام واگرایی Kullback-Leibler کم کند.

## 4. Optimization

- استفاده از گرادیان گاهشی برای پیدا کردن بهترین نمایش low-dim به طوری که ساختار high-dim تا جایی که امکان دارد حفظ شود

معایب {  
- محاسبات بالا (هر pair باید محاسبه شود)  
- حساسیت به hyperparameter  $\alpha$

- توجه بیشتر به Local Structure به جای Global Structure

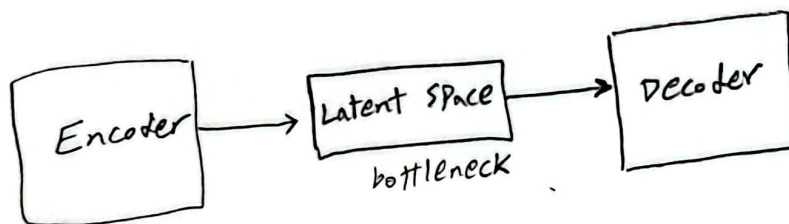


طبقه بکر - ← perplexity خوب است که بین 5 تا 50 باشد

## >> Auto Encoder

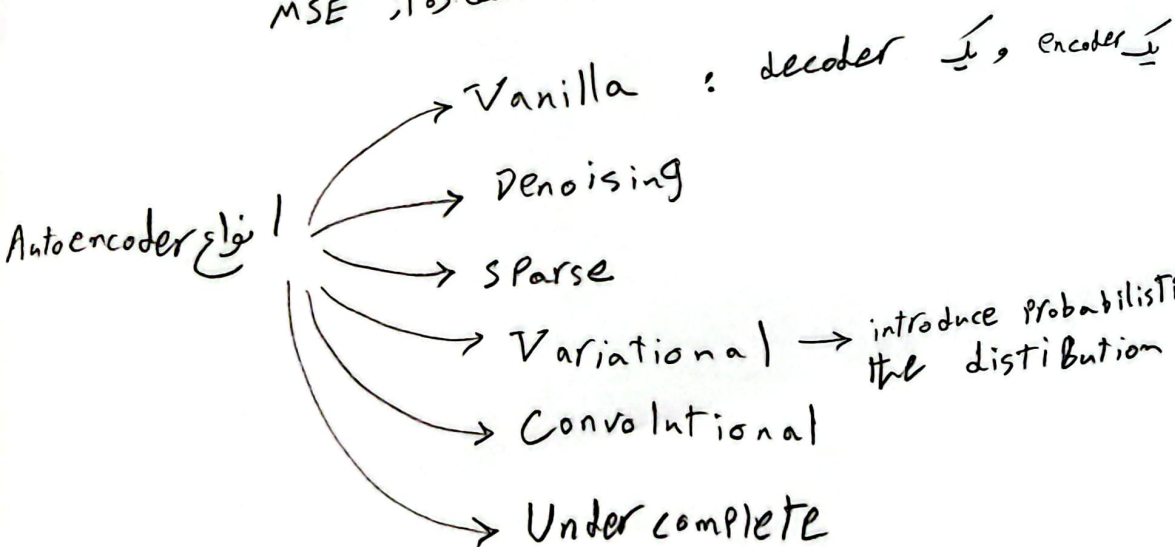
→ unsupervised learning

کاربردها { dimensionality reduction  
feature learning  
data compression



هدف: کمینه کردن اختلاف بین ورودی و خروجی (reconstruction error)

← معیار استاندارد از MSE



کاربردها

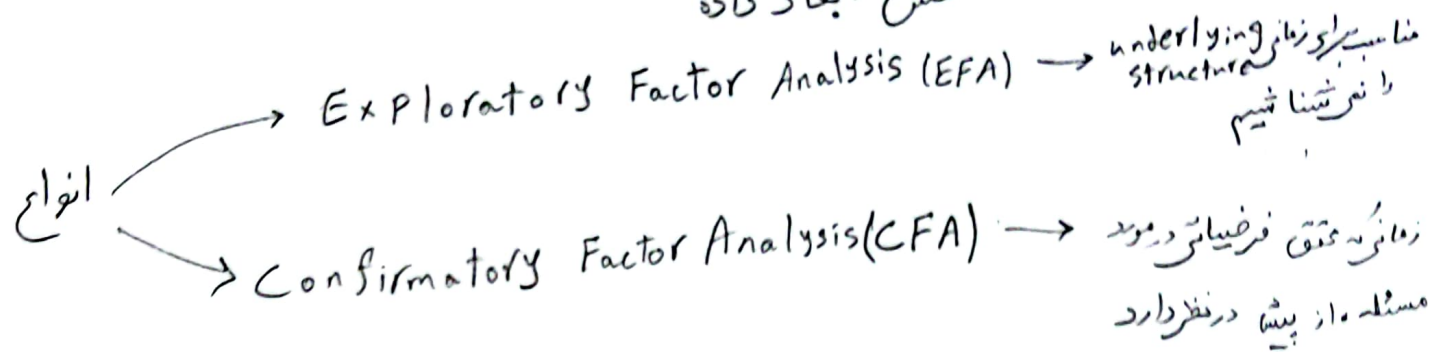
- کاهش ابعاد ← مانند PCA ولی دارای قابلیت گرفتن روابط غیرخطی
- حذف نویز
- تشخیص ناخواسته
- ایجاد data sample های جدید

# >> Factor Analysis

- روش آماری برای یافتن روابط بین متغیرها

اهداف ← شناسایی ساختار درونی داده‌های مشاهده شده

کاهش ابعاد داده



Latent variables → Factors

# گام‌ها

1- آماده سازی داده‌ها

- استناد به داده‌ها

2- انتخاب تعداد فاکتورها

رویکرد ها ← معیار گایز (eigenval > 1)  
(استفاده از) ← scree plot

← آنالیز موازی

3- استخراج فاکتورها

{ Principal Component Analysis (PCA)

{ Common Factor Analysis → Principal Axis Factoring (PAF)  
→ Maximum Likelihood

4- چرخش

امکان برای آنکه output بیشتر تفسیر پذیر باشد ← با رسیدن به ساختار فاکتور ساده‌تر و با معیار  
تفسیر می‌شود فاکتورها، unrotated → orthogonal Rotation (Varimax)  
Oblique Rotation (Promax) → Allows factor to be correlated

5 - تفسیر شناسایی و نام گذاری فاکتورها بر اساس متغیرهای مربوط به آنها

6 - اعتبار، سببی

بررسی پایداری و اعتبار ساختن فاکتور ← استفاده از  
split-sample validation  
cross-validation

یگر از پیش فرضهای مهم ← فشرده بودن رابطه بین متغیرها و فاکتورها

kernel PCA

توسعه یافته PCA برای non-linear mapping of data

Project's high-dimensional data into lower-dimensional space  
PCA شامل موارد زیر است:

1. Centering the data
2. Computing the covariance matrix
3. Eigen decomposition (of cov matrix)
4. Projection

در K-PCA داریم:

map the data into a higher-dimensional feature space

1. Choosing a kernel  $(K(x, y))$  معرّفی →
  - Linear kernel:  $K(x, y) = x^T y$
  - Polynomial kernel:  $K(x, y) = (x^T y + c)^d$
  - Gaussian (RBF) kernel  
 $K(x, y) = \exp\left(-\frac{\|x - y\|^2}{2\sigma^2}\right)$   
 $\Rightarrow \exp(-\gamma \|x - y\|^2)$
2. Computing the kernel matrix (K)  
 $K_{ij} = K(x_i, y_j)$

3. Centering the kernel Matrix

$$K_c = K - \mathbf{1}_n K - K \mathbf{1}_n + \mathbf{1}_n K \mathbf{1}_n$$

$\mathbf{1}_n \rightarrow n \times n$  matrix with all elements equal to  $\frac{1}{n}$

4. Eigen decomposition (on  $K_c$ )

5. Projection (using eigenvect from step 4)