## Chapter 9: Exploratory Data Analysis (EDA)

* dot product: $\langle u.v \rangle = u_1 v_1 + u_2 v_2 + \ldots +$
* norm : $\|v\| = \sqrt{\langle v, v \rangle}$
* unit vector: $\|v\| = 1$
* orthogonal : $\langle v, u \rangle = 0$

* projection $= \dfrac{\langle x, u_1 \rangle}{\langle u_1, u_1 \rangle} \cdot u_1 + \dfrac{\langle x, u_2 \rangle}{\langle u_2, u_2 \rangle} \cdot u_2 + \ldots$

   First check that, $u_1$ & $u_2$ are orthogonal

* PCA
* Transpose: $d_1 \times d_2$ and $d_2 \times d_1$
* Square matrix: $A' = A$
* $A.u = \lambda.u$, $u =$ eigenvector and $\lambda =$ eigenvalue

* Centroid $(\bar{x}) = \frac{1}{n}(x_1 + \ldots + x_n)$ } Centering
* Replace each $x_i$ by $x_i - \bar{x}$
* Normalization : $\dfrac{\text{values}}{\text{standard deviation}}$    $\boxed{sid = \sqrt{\text{Variance}}}$

   Variance $= \frac{1}{n}(x_1^2 + x_2^2 + \ldots)$

* subspace spanned by $u$
* When the variance along $u$ is large indicate that the projection along $u$ are interesting.
* Maximizing variance.
* let $u$ a unit vector
  - The projection of $x$ along vector $u$ is $\boxed{\langle x, u \rangle.u}$ (coefficient)
  - The variance is
  $$\frac{1}{N} \sum_{i=1}^{N} \langle x_i, u \rangle^2$$
  - N times the variance is
  $$\sum_{i=1}^{N} \langle x_i, u \rangle^2 = \sum_{i=1}^{N} u' x_i (x_i)' u = u' \left( \sum_{i=1}^{N} x_i (x_i)' \right) u$$

  so, $S = \sum_{i=1}^{N} x_i (x_i)'$  • the scatter matrix

Variance = eigevalue.

- The unit vector u that maximizes $u^t S u$ is same as
  - the eigenvector with largest eigenvalue.
  - the eigenvector is called the first principal component.

- Smmary of PCA
  - centering and normalization
  - finding a unit vector u, such that the variance of the projections along direction u is maximized
  - the variance is $\frac{1}{N} \cdot u^t S u$.
  - Computing the unit vector u which maximizes $u^t S u$
    . eigenvector of S with the largest eigen value

1. centering and normalization
2. compute the scatter matrix
3. compute the eigenvectors & eigenvalues of the scatter matrix.

$$S = \sum_{i=1}^{n} x_i (x_i)^t$$