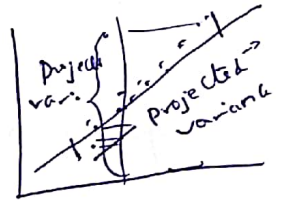# Principal Component Analysis.

Objective: Find a linear (orthogonal) Projection of the data such that the Projected variance is maximized.

$$D = \{x_1, x_2 \ldots, x_n\}$$

$$\underline{x_i} \in R^d.$$

Let $X = \begin{bmatrix} x_1, x_2 \ldots x_n \end{bmatrix}_{d \times n}$ matrix constructed from the data - Data matrix.

Let us start by projecting the data onto a 1D manifold [line].

Let $u_1 \in R^d$ be the line on to which we are seeking the projection.

(101)

Let $z$ represent the projected points.

we have

$$Z_{\bullet} = u_1^T X$$
$$\substack{1 \times n} \quad \substack{1 \times d \quad d \times n}$$

Every component of $z$ corresponds to the projection

of $x$ on $u_1$.

objective of PCA

$$u_1^* \overset{\bullet}{Z} = \underset{\partial u_1}{\arg\max} \ \text{var}(u_1^T x). \quad -(1)$$

$$\text{var}(u_1^T X) = u_1^T S u_1 \qquad\qquad \underset{d \times d}{S} = (X - \bar{X})(X - \bar{x})^T$$

$$\bar{X} = \frac{1}{n} \sum_{i=1}^{n} x_i$$

observe that $u_1^T S u_1$ is a scalar & $S$ is

PSD & symmetric.

Thus (1) is ill-defined if it is unconstrained

(102)

However since we are only interested in the direction of projection we can constrain the problem by fixing the norm of $u_i$ to any constant.

Let $u_i^T u_i = 1$.

Thus, $u_i^* = \underset{u_i}{\text{argmax}} \left( u_i^T S u_i \right)$

$$\text{s.t } u_i^T u_i = 1.$$

$$L = u_i^T S u_i - d \left( u_i^T u_i - 1 \right)$$

$$\frac{\partial L}{\partial u_i} = 0 \quad \Rightarrow \quad S u_i = d u_i$$

$$\therefore u_i^T S u_i = u_i^T d u_i$$

$$= d.$$

Thus, $u_i^* = \underset{u_i}{\text{argmax}} \left( d \right)$.

$(1 \to 2)$

But we have $d$ as the Eval of S.

Thus, $u_1$ will be the direction corresponding of the Evec of S.
to the maximum Eval of S.

Extending this, suppose the Evals of S are
ordered in accordance to their value.

$$d_1 > d_2 > d_3 > \cdots d_d.$$
$$u_1 \quad u_2 \quad \cdots \quad u_d.$$

Define
$$U = [\psi_1 \ v_2 \cdots \psi_d]_{d \times d}$$

$$\underset{d \times n}{Z} = \underset{d \times d}{U^T} \underset{d \times n}{X} \qquad - \text{ represents the projection}$$

of X on to a new-co-ordinate system
where the variance is maximized.

Since S is symmetric, U will be a orthonormal
matrix — columns of U are called principal
vectors.

$(104)$.

Now if we know that the original data 'effectively lies' in a $p$-dimensional subspace of $d$ [usually $p << d$]

Then one can consider

$$\hat{Z}_{P \times n} = \hat{U}^T X_{P \times d \ d \times n}$$

where $\hat{U} = [\mu_1, \mu_2, \dots \mu_P]_{d \times P}$

Now $\hat{Z}$ will be a new set of data points lying in a $p$-dimensional subspace.

Thus, PCA can be used for dimensionality reduction.

$Z = U^T X$   since $U$ is orthogonal,

$UZ = U U^T X$
$\underline{= X}$   Thus, one can recover back the original data from projections.

$(105)$ ?

we can also reconstruct the data after we reduce the dimensions.

$$\hat{z} = \hat{U}^T x$$

Now $U\hat{z}$ (call it $\hat{x}$) will be the reconstruction of $x$ based on the first $P$ - principal components.

Lets consider the second formulation.

Error in the reconstruction is

$$e = \|x - \hat{x}\|_2 .$$

$$e = \frac{1}{n} \sum_{i=1}^{n} \| x_i - \hat{x}_i \|_2^2 \qquad - (1)$$

$$\hat{x}_i = \sum_{j=1}^{P} \alpha_{ij} u_j + \sum_{j=P+1}^{D} \beta_i u_j$$

$$\sum_{2}^{n} \qquad \ast \quad (\text{info thin})$$

substituting for $\hat{x}_i$ in (1),

$$\alpha_{ij} = x_i^T u_j \quad , \quad j = 1, \ldots P$$

$$\beta_j = \bar{x}_i^T u_j \qquad \bar{x} = \frac{1}{n} \sum_{j=1}^{n} x_j,$$

$$j = P+1, \ldots D.$$

NOW, $\quad \cancel{x_i - \frac{}{n}} \quad$ we have

$$x_i = \sum_{j=1}^{D} (x_i^T u_j) u_j$$

$$\therefore \quad x_i - \hat{x}_i = \sum_{j=P+1}^{D} \left[ (x_i - \bar{x})^T u_j \right] u_j$$

NOW, $\quad e = \frac{1}{n} \sum_{j=P+1}^{D} u_j^T S u_j \qquad \therefore \text{ minimizing } e$

$$(107).$$

$$e = \sum_{j=p+1}^{D} u_j^T S u_j$$

Similar to the previous case, $e$ will be minimized when $u_j$s are the ~~ta~~ Evecs of $S$ corresponding to last $D-P$ Evals..

# Linear models.

$h(x)$ is of the form

$$h(x) = W^T \phi(x)$$

where $\phi$ is a fixed function of $x$.

$\phi$ can be polynomial, logistic sigmoid etc.

$h(x)$ is called linear because its a linea. function in the parameter space $W$.

~~In all~~ We shall consider some of the families of these linear models.

Before we go to linear models, let us look at a general important result.

With our usual notations.

Let $R$ denote the risk associated with a classifier $h(x)$.

let us consider the squared error loss.

$$L(y, h(x)) = (h(x) - y)^2$$

$$R(h) = \int \int L(y, h(x)) P(x, y) \, dx \, dy$$

$$= \int \int (h(x) - y)^2 P(x, y) \, dx \, dy$$

Goal of ML : find $h$ such that

$$h^* = \underset{h}{\arg\min} \; R(h)$$

$$h^* = \arg\min \int \int (h(x) - y)^2 P(x, y) \, dx \, dy .$$

(110).

$$\frac{\partial R}{\partial h} = 2 \int \left( h(x) - y \right) P(x,y) \, dy$$

$$h(x) = \frac{\int y \, P(x,y) \, dy}{\int \, P(x,y) \, dy} = \frac{\int y \, P(x,y) \, dy}{P(x)}$$

$$= \int y \, P(y|x) \, dy = E_y \left[ y | x \right]$$

Thus, the optimal classifier is the Conditional expectation of the labels given the data for squared error loss.

when $y \in \{0,1\}$, $E_y \left[ y|x \right] = P\{y=1|x\} = q_1(x)$

Giving us back the Baye's classifier.

(iii).

Now, lets come backs to the linear models.

$$y \quad h(x) = w^T \phi(x).$$

Consider that $\phi$ true $y$ is, a deterministic function of $x$ $h(x)$ with an uncertainity $\epsilon$. being approximated by

$$\therefore \quad y = h(x) + \epsilon$$

$$= w^T \phi(x) + \epsilon$$

Let us assume that $\epsilon \sim N(0, \sigma^2)$.

Now since we know that the optimal estimator is the conditional expectation,

we need $\quad E[y/x] = \int y f(y|x) \, dy$

(112).