## 3.4. Maximum Likelihood Estimation

$x_1, \ldots, x_n$ independent sample from pdf

$$f(x, \vartheta) \quad, \quad \vartheta \text{ a parameter.}$$

$$L(x; \vartheta) = \prod_{i=1}^{n} f(x_i, \vartheta) \qquad \underline{\text{likelihood function}}$$

$$\ell(x; \vartheta) = \log L(x; \vartheta) = \sum_{i=1}^{n} \log f(x_i, \vartheta)$$

$$\underline{\log\text{-likelihood fct.}}$$

Find $\vartheta$ which fits the data best, i.e.,

$$\hat{\vartheta} = \arg\max_{\vartheta} \ell(x; \vartheta)$$

$\hat{\vartheta}$ is called ML estimator. (MLE)

<u>Th. 3.6.</u> $X \sim N_p(\mu, \Sigma)$, $x_1, \ldots, x_n$ i.i.d. sample of $X$.
The MLEs of $\mu, \Sigma$ are:

$$\hat{\mu} = \frac{1}{n} \sum_{i=1}^{n} x_i = \bar{x} \quad, \quad \hat{\Sigma} = \frac{1}{n} \sum_{i=1}^{n} (x_i - \bar{x})(x_i - \bar{x})^T = S_n.$$

<u>Proof.</u>   Density of $N_p(\mu, \Sigma)$

$$f(x) = \frac{1}{(2\pi)^{p/2} |\Sigma|^{1/2}} \exp\left(-\frac{1}{2}(x - \mu)^T \Sigma^{-1}(x - \mu)\right), \quad x \in \mathbb{R}^p$$

$$\ell(x_1, \dots, x_n; \mu, \Sigma)$$

$$= \sum_{i=1}^{n} \left[ \log \frac{1}{(2\pi)^{p/2} |\Sigma|^{1/2}} - \frac{1}{2}(x_i - \mu)^T \Sigma^{-1}(x_i - \mu) \right]$$

$$= n \underbrace{\log \frac{1}{(2\pi)^{p/2}}}_{\text{Constant}} + \frac{n}{2} \log |\Sigma^{-1}| - \frac{1}{2} \sum_{i=1}^{n}(x_i - \mu)^T \Sigma^{-1}(x_i - \mu)$$

Leave the constant, set $\Lambda = \Sigma^{-1}$

$$\ell^*(\mu, \Sigma) = \frac{n}{2} \log \Lambda - \frac{1}{2} \sum_{i=1}^{n}(x_i - \mu)^T \Lambda (x_i - \mu)$$

$$= \frac{n}{2} \log |\Lambda| - \frac{1}{2} \sum_{i=1}^{n} tr\left( \Lambda (x_i - \mu)(x_i - \mu)^T \right)$$

$$= \frac{n}{2} \log |\Lambda| - \frac{1}{2} tr\left( \Lambda \sum_{i=1}^{n}(x_i - \mu)(x_i - \mu)^T \right)$$

Steiners rule:

$$\sum_{i=1}^{n}(x_i - \mu)(x_i - \mu)^T$$

$$= \underbrace{\sum_{i=1}^{n}(x_i - \bar{x})(x_i - \bar{x})^T}_{n S_n} + (\bar{x} - \mu)(\bar{x} - \mu)^T$$

$$\geq n S_n \quad (\text{equality if } \mu = \bar{x})$$

$$\leq \frac{n}{2} \log |\Lambda| - \frac{n}{2} tr(\Lambda S_n) = \ell^*(\mu^*, \Lambda)$$

$$\max \ \ell^*(\mu^*, \Lambda)$$

Need $\quad \dfrac{\partial}{\partial \Lambda} \log |\Lambda| = (\Lambda^{-1})^T$

$$\dfrac{\partial}{\partial \Lambda} tr(\Lambda A) = A^T$$

$$\dfrac{\partial}{\partial \Lambda} \ell^*(\mu^*, \Lambda) = \dfrac{n}{2} \Lambda^{-1} - \dfrac{n}{2} S_n \overset{!}{=} 0_{p \times p}$$

$$\Longleftrightarrow \quad \Sigma^* = S_n \quad \boxed{\exists}$$

## 4. Dimensionality Reduction

Represent data in a low dimensional space
    (high dim.)
in an "optimal" way. Dim. 1, 2, 3 allow for
visualization.

### 4.1. Principal Component Analysis (PCA)

Loose as little information as possible.

Given data $x_1, \ldots, x_n \in \mathbb{R}^p$.

a) Find a $k$-dim. subspace such that the
   projections of $x_1, \ldots, x_n$ thereon represent the
   data ~~it~~ on its best.

b) Preserve as much variance as possible.

   a) and b) are equivalent. $\to$ later

$x_1, \ldots, x_n$ independently sampled from some distribution.

Sample mean: $\bar{x} = \frac{1}{n} \sum_{i=1}^{n} x_i$

Sample covariance matrix: $S_n = \frac{1}{n-1} \sum_{i=1}^{n} (x_i - \bar{x})(x_i - \bar{x})^T$

$$\begin{pmatrix} \bar{x}: & \text{unbiased estimator of } E(X) \\ S_n: & \text{unbiased estimator of } \text{Cov}(X) \end{pmatrix}$$
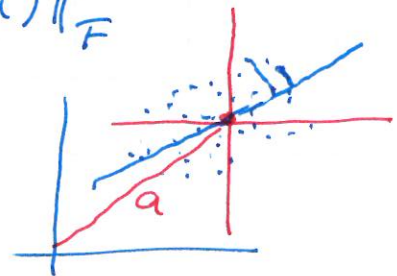
( $\underline{\text{Ex.}}$ MNIST data, $n = 500$, $p = 28 \cdot 28 = 784$ )

### 4.1.1. Find the best projection.

Consider the opt. problem

$$\min_{\substack{a \in \mathbb{R}^p, \\ Q \text{ orth. proj. on a } k\text{-dim. subspace}}} \sum_{i=1}^{n} \| x_i - a - Q(x_i - a) \|_F^2$$

$Q$ orth. proj. on a $k$-dim. subspace



$$\min_{a, Q} \sum_{i=1}^{n} \| x_i - a - Q(x_i - a) \|^2$$

$$= \min_{a, Q} \sum_{i=1}^{n} \| (I - Q)(x_i - a) \|^2$$

$$= \min_{a, R} \sum_{i=1}^{n} \| R(x_i - a) \|^2, \quad R = I - Q \text{ (orth. proj. as well)}$$

$$= \min_{a,R} \sum_{i=1}^{u} (x_i - a)^T R^T \dot{R} (x_i - a)$$

$$= \min_{a,R} \sum_{i=1}^{u} tr\left( (x_i - a)^T R (x_i - a) \right)$$

$$= \min_{a,R} \sum_{i=1}^{u} tr\left( R (x_i - a)(x_i - a)^T \right)$$

$$= \min_{a,R} tr\left( R \sum_{i=1}^{u} (x_i - a)(x_i - a)^T \right)$$

$$\geq \min_{R} tr\left( R \sum_{i=1}^{u} (x_i - \bar{x})(x_i - \bar{x})^T \right) \quad \left( \begin{array}{l} \text{see MLE for } N(\mu, \Sigma) \\ \text{equality if } a = \bar{x} \end{array} \right)$$

$$= \min_{R} tr\left( R (u-1) S_u \right)$$

$$= \min_{Q} (u-1)\, tr\, S_u (I - Q)$$

It remains to solve

$$\max_{Q} tr(S_u Q) \quad , \quad Q \text{ orth. proj., } Q = \sum_{i=1}^{k} q_i q_i^T, \; q_i \text{ orth.}$$

$$Q = \tilde{Q} \tilde{Q}^T, \; \tilde{Q} = (q_1, \dots, q_k)$$

$$= \max_{\tilde{Q}^T \tilde{Q} = I_k} tr\left( \tilde{Q}^T S_n \tilde{Q} \right) = \sum_{i=1}^{k} \lambda_i (S_u) \quad (\text{Ky Fan, Th. 2.4})$$

where $\lambda_1(S_u) \geq \cdots \geq \lambda_k(S_u) \geq \cdots \geq \lambda_p(S_u)$ are the eigenvalues of $S_u$ in decreasing order.

The max is attained if $q_1, \dots, q_k$ are the orthogonal eigenvector corresponding to $\lambda_1(S_u), \dots, \lambda_k(S_u)$.

## 4.1.2 Preserve most variance

Seek a hyperplane so that the proj. data has most variance.

$$\max_Q \; \sout{\sum_{i=1}^{u} \| Qx_i - \frac{1}{n} \sum_{\ell=1}^{u} Qx_\ell \|^2}$$

$$\max_Q \; \sum_{i=1}^{u} \| Qx_i - \frac{1}{n} \sum_{\ell=1}^{u} Qx_\ell \|^2 \;,\quad Q = \tilde{Q}\hat{Q}^T, \quad \hat{Q}^T\tilde{Q}^{\#} = I_k$$

orth. proj.

$$= \max_Q \; \sum_{i=1}^{u} \| Qx_i - Q\bar{x} \|^2$$

$$= \max_Q \; \sum_{i=1}^{u} \| Q(x_i - \bar{x}) \|^2$$

$$= \max_Q \; \sum_{i=1}^{u} \operatorname{tr}(x_i - \bar{x})^T Q (x_i - \bar{x})$$

$$= \max_Q \; \operatorname{tr} Q \sum_{i=1}^{u} (x_i - \bar{x})(x_i - \bar{x})^T$$

$$= \max_Q \; (u-1)\, \operatorname{tr} Q \, S_n$$

with the same solution as above.

4.1.3 How to carry out PCA

Given $x_1, \ldots, x_n \in \mathbb{R}^p$, fix $k \ll p$

Compute $S_n = \frac{1}{n-1} \sum_{i=1}^{n} (x_i - \bar{x})(x_i - \bar{x})^T$

$$S_n = V \Lambda V^T, \quad \Lambda = \text{diag}(\lambda_1, \ldots, \lambda_p)$$

$\lambda_1 \geq \cdots \geq \lambda_p, \quad V = (v_1, \ldots, v_p) \in \mathcal{O}(p)$ spectral decomposition

$v_1, \ldots, v_k$ are called the $k$ <u>principal eigenvectors</u> to the <u>principal eigenvalues</u> $\lambda_1, \ldots, \lambda_k$.

Projected points $\hat{x}_i = \begin{pmatrix} v_1^T \\ \vdots \\ v_k^T \end{pmatrix} x_i, \quad i = 1, \ldots, n.$