



This page contains handwritten notes on various topics in statistics and machine learning:

- REGRESSION**:
  - Training set  $D = \{(x_1, y_1), \dots, (x_n, y_n)\}$ , and best predictor  $\hat{y} = h(x)$
  - Minimizing empirical MSE:  $(\hat{y} - y)^2 = (\sum_i (h(x_i) - y_i))^2 / n$
  - How do we select  $h$ ? VALIDATION
  - Validation set:  $\{(x_{n+1}, y_{n+1}), \dots, (x_m, y_m)\}$
  - MSE validation:  $(MSE) = \frac{1}{m-n} \sum_{i=n+1}^m (\hat{y}_i - y_i)^2$
  - Once we get  $h_{\text{best}}$ , retrain on full training set to find best  $h_{\text{best}}$ .
  - $X, Y$  can be independent if their graph is a rectangle (possible values of  $Y$  should not be dependent on  $X$ )
- Parametric model (generative)**:
  - $y = f(x) + \epsilon$ ,  $X = R$
  - $f(y) = f_1(y_1) \dots f_m(y_m)$
  - $P(y) = f_1(y_1) \dots f_m(y_m)$
  - $p(x|y) = \frac{1}{V^m} e^{-\frac{\|x-y\|^2}{2V^2}} \rightarrow N(\mu_y, V^2)$
  - $p(x|y) = \frac{1}{V^m} e^{-\frac{\|x-y\|^2}{2V^2}} \rightarrow N(\mu_x, V^2)$
  - So,  $p(x,y) = p(y)p(x|y) = (c - \|x-y\|^2)^{-\frac{m}{2}}$
  - Unknowns:  $(\mu_x, \mu_y, V)$  = 3 real parameters
  - Given  $m$  examples  $(x_1, y_1), \dots, (x_m, y_m) \sim p(x, y)$
  - If you know  $\mu_x, \mu_y, V$ ,  
 $h_{\text{MAP}}(x) = E[Y|x] = E[p(Y|x)] = \text{prior}(p(x)) + \text{posterior}(p(y|x))$
  - Posterior (MAP) for  $m=2$ :
    - Violates Gaussian assumption not made, optimality criterion changed.
    - Model: Given  $y_1, x_1, y_2, x_2$  we can  $E[y|y_1, x_1, y_2, x_2]$
    - Given  $y_1, x_1$  make a linear prediction with threshold, the projection  $w^T x$  to max class decision.
    - Given  $y_1, x_1$  projection w/  $x_2$  to max class?
    - Given  $y_1, x_1$  projection w/  $x_2$  is a scalar  $R^2$  or mean  $w^T x_1 y_1$  and variance  $w^T x_1 w$
    - Between-class variance (signal) of projection:  
 $\text{S}_{\text{between}}^2 = \text{Cov}(x_1, x_2) - w^T x_1 w^T x_2$
    - Within-class variance (noise) of projection:  
 $\text{S}_{\text{within}}^2 = w^T x_1 x_1 w^T x_1 + w^T x_2 x_2 w^T x_2$
    - Measure of separation between classes:  
 $\text{FNR}(w) = \Pr(Y=0 | w^T x < 0)$
    - Best  $w$  is the one that minimizes FNR
    - Wishart =  $\text{Beta}(m, m)$
  - MAIN RESULT:  $E[\hat{w}]$  is invertible,  $\text{Var}(\hat{w}) = \frac{1}{m} (X^T X)^{-1} (X^T Y)$
  - Final decision rule:  $\hat{w}^T x > 0 \Rightarrow C_1$  (if  $\hat{w}^T x > 0$ )
    - can be chosen as in LDA, but in practice it's treated as a hard classifier (it only outputs 0 or 1)
    - can't calculate probabilities (but can calculate  $\Pr(Y=1 | x)$ )
    - is not necessarily symmetric (unit norm length)
    - if  $A = ABC$ ,  $A^{-1}$  is symmetric (not true for circular bending)
    - Conditional probability:  $\text{PC}(A|B) = \text{PC}(B|A)$
    - Bayes Rule:  $\text{PC}(A|B) = \text{PC}(B|A) \text{PC}(A)$
    - Linear total probability:  $P(B) = \Pr(A \cap B) + \Pr(A^c \cap B)$
    - PCA =  $\text{PC}(B|A) \text{PC}(A)$
  - Independence of Random Variables:  $P(z|x) = P(z) Y$  and  $P(x|z) = P(x) Z$
  - Conditional Independence:  $\text{PC}(A|B, C) = \text{PC}(A|B) \text{PC}(C)$
  - Expectation:  $E[X|Y] = \sum x_i p(x_i|y)$
  - Variance:  $\text{Var}(X) = E[(X - E[X])^2] = \sigma^2$
  - $E[X^2] = E[X] + \text{Var}(X)$
  - $\text{Var}[X] = \text{Var}[X]^2$
  - + Gaussian distribution:
    - $N(\mu, \sigma^2)$ ,  $\mu \in R$ ,  $\sigma \geq 0$
    - $f(x|\mu, \sigma^2) = \frac{1}{\sqrt{2\pi}\sigma} e^{-\frac{(x-\mu)^2}{2\sigma^2}}$
    - $\Phi(x) = \int_{-\infty}^x \frac{1}{\sqrt{2\pi}\sigma} e^{-\frac{(t-\mu)^2}{2\sigma^2}} dt = 1 - \Pr(X > \mu)$
    - Probability that  $X > \mu$  is  $\Pr(X > \mu) = 1 - \Phi(\mu)$
    - $X$  is a Gaussian RV w/ mean  $\mu$ , variance  $\sigma^2$
    - $X = Y - \mu$  is standard normal and  $\Pr(Y > y) = \Pr(X > y - \mu) = \Phi(y - \mu)$
    - where  $\tau = y - \mu$
  - Conditional density:  $f_{Y|X}(y|x) = \frac{f_{X,Y}(x,y)}{f_X(x)}$
  - Wishart joint continuous & independent:
    - $f_{X,Y}(x,y) = f_X(x)f_Y(y)$  for all  $y: f_Y(y) > 0$
    - conditional expectation of  $X$  given  $Y=y$ :  
 $E[X|Y=y] = \int_{-\infty}^{\infty} x f_X(x|y) dx$
    - $f_{X,Y}(x,y) = E[X|Y=y]$  is a deterministic function of  $y$
    - Law of iterated expectation:  
 $E[E[X|Y]] = E[X]$  whenever  $E[X|Y]$  is finite.
  - $E[X|Z] = E[E[X|Y|Z]]$
  - Condition:  $R_{xy} = E[XY] = \int \int xy f_{X,Y}(x,y) dx dy$
  - $E[X^2] = E[E[X^2|Y]]$
  - $E[X^2] = E[E[X|Y]E[Y]]$
  - If  $X, Y$  indep.,  $E[XY] = E[X]E[Y]$
  - Correlation coefficient:  $\rho_{XY} = \frac{\text{Cov}(X,Y)}{\sqrt{\text{Var}(X)\text{Var}(Y)}}$
  - Chebyshev-Schwarz Inequality:  
 $|E[X|Y]| \leq E[|X| | Y]$
  - Triangle Inequality:  
 $|E[(X+Y)^2]| \leq \sqrt{E[X^2]} + \sqrt{E[Y^2]}$







**ICA**

Given  $Z = (\zeta_1 \dots \zeta_n)$ ,  $k \leq \text{rank}(Z)$

Eigenvector decomposition of Empirical Correlation Matrix:

$$Z^T Z = U \Lambda U^T \quad U = (\zeta_1 \dots \zeta_k), \quad \Lambda^{-1} = U^T$$

$$\Lambda = \begin{pmatrix} \lambda_1 & & \\ & \ddots & \\ & & \lambda_k \end{pmatrix}, \quad \lambda_1 \geq \lambda_2 \dots \geq \lambda_k > 0$$

Best k-dimensional Subspace:  $V_k = \text{span}(\zeta_1, \dots, \zeta_k)$

$$F = (\zeta_1 \dots \zeta_k) \quad F^T F = I_k$$

Encoder/Encoded Test Sample:  
 $\tilde{x} \in \mathbb{R}^d \rightarrow y = F^T \tilde{x} \in \mathbb{R}^k$

Reconstruct/Decode:  $y \rightarrow \hat{x} = Fy = F\Lambda^{-1}x \in \mathbb{R}^d$   
 $= \text{Proj}_{V_k}(x)$

Approximation Error:  $\|x - \hat{x}\|^2 = \|x\|^2 - \|\Lambda^{-1}x\|^2$   
 $= \|x\|^2 - \|U\Lambda^{-1}U^T x\|^2$   
 $= \|x\|^2 - \|U\Lambda^{-1}\|^2 \|U^T x\|^2$

Total Approx Error:  $\sum_j \|x_j - \hat{x}_j\|^2 = \sum_j \lambda_j$

Dual/kernel PCA

Given  $Z, k$

Gram/Kernel Matrix:  $Z^T Z = V \Lambda V^T, V^T = V^T$

Eigen Decomp:  $\Lambda = (\lambda_1 \dots \lambda_n), \lambda_i \geq \dots \geq \lambda_k > 0$

Best k-dimensional Subspace:  $V_k = \text{span}(\zeta_1 \dots \zeta_k)$

$F = Z [(\zeta_1 \dots \zeta_k)]^{-1/2}$  or  $(V \Lambda V^T)^{-1/2}$

$F^T F = I_k, V_k$  bests  $\text{span}(\zeta_1 \dots \zeta_k)$

Embed/Encode:  
 $x \in \mathbb{R}^d \rightarrow y = F^T x = \frac{1}{\sqrt{k}} \begin{pmatrix} -1 & \dots & -1 \\ \vdots & \ddots & \vdots \\ -1 & \dots & -1 \end{pmatrix} (Z^T x)$

Reconstruction: only if kernel feature mapping explicitly known  
 $\text{Proj}_{V_k}(x) = \hat{x} = Fy = F\Lambda^{-1}x$

Approx Error:  $\|x - \hat{x}\|^2 = \|x\|^2 - \|\Lambda^{-1}x\|^2$   
 $= \|x\|^2 - \|U\Lambda^{-1}U^T x\|^2$   
 $= \|x\|^2 - \|U\Lambda^{-1}\|^2 \|U^T x\|^2$

Total Approx:  $\sum_j \|x_j - \hat{x}_j\|^2 = \sum_j \lambda_j$

Encoding matrix Path Set:  $F^T Z = \Lambda^{-1/2} V_k^T$

**HW #9 Clustering**

1) Consider a 1-dimensional dataset of  $n=2m+1$  points ( $m \in \mathbb{N}$ ) spread on the real line with one point at  $0$ ,  $m$  points at  $a$ , and  $m$  points at  $-a$ , where  $0 < a \leq \frac{1}{2}$ . Now suppose we run k-means on the data points  $w_1, w_2, \dots, w_n$  and all centroid choices other than the random  $(0, a, -a)$  perform all distinct clusterings so that the k-means algorithm could converge to:  
 clustering 1:  $\{0, a, -a, \dots, a\}$ ,  $\{c_1, c_2, \dots, c_{m+1}\}$   
 clustering 2:  $\{0\}$ ,  $\{a, -a, \dots, a\} \cup \{w_m, w_{m+1}\}$

2) For each solution, provide expressions for centroids and WCSS in terms of  $m, n, a$ .

clustering 3:  $w_1 = 0, w_2 = \frac{a}{m+1} + a \frac{m+1}{m+1}$

$$\mu_1 = a$$

$$\text{WCSS}_1 = \frac{1}{2} (0-a)^2 + m(a-x_1)^2 + m(a+x_1)^2 = \frac{1}{2} + ma^2 + \frac{m}{2} \frac{a^2}{m+1} + ma^2$$

clustering 2:  $w_1 = 0, w_2 = \frac{a}{2}$

$$\text{WCSS}_2 = \frac{1}{2} (0-a)^2 + m(a-x_1)^2 + m(a+x_1)^2 = \frac{1}{2} + ma^2 + \frac{m}{2} \frac{a^2}{m+1} + ma^2$$

3) Compute the ratio of the largest and smallest WCSS as a function of  $m, n, a$  and comment on the implications.

Clustering 1)  $\text{WCSS}_1 \leq \text{WCSS}_2$

$$\text{WCSS}_1 = \frac{ma^2}{m+1} \leq \frac{ma^2}{2} \leq \frac{mb^2}{2} = \text{WCSS}_2$$

Ratio  $\text{WCSS}_1 / \text{WCSS}_2 = \frac{1}{m+1}$  This ratio can become arbitrarily large for suitably large  $m$ .

So, if a dataset, for which k-means can converge to a solution whose WCSS is arbitrarily higher than the best solution.

Take away: necessary to run k-means w/ multiple initializations of centroids. Exchange solutions resulting in the smallest WCSS.

4) Suppose two initial centroids are chosen independently uniformly at random over  $[0, a]$ .  
 Compute the probability of converging to each possible solution.  
 This is determined by whether the mid-point of the initial centroids lies to the left or right of  $a$ .

Cluster 2)  $P(\text{cluster } 2) = P\left(\frac{a_1 + a_2}{2} < a\right)$   
 where  $a_1, a_2 \sim \text{Uniform}[0, a]$

$$\frac{1}{2} \times 2a \times 2a \times \frac{1}{a^2} \frac{1}{a^2} = \frac{2a^2}{a^2} = \frac{1}{2}$$

So,  $P(\text{cluster 1}) = 1 - \frac{1}{2} = \frac{1}{2}$

$\mu_2(a) = \frac{a_1 + a_2}{2} < a$

$\mu_1(a) = \frac{a_1 + a_2}{2} > a$

$\text{Cov}(X, Y) = \frac{1}{n} \sum_{i=1}^n (X_i - \bar{X})(Y_i - \bar{Y})$

**HW #8 C (uniting)**

8.1 a) Let  $X = \text{diag}([-\pi/2, \pi/2])$  and  $Y = \sin(X)$ . Compute  $C_{WLMPLS}, b_{WLMPLS}$  and  $\text{argmin}_W E[(Y-WX-b)^2]$

For each value of  $W$ , we have  $b$  which minimizes  $E[(Y-WX-b)^2]$  is given by  $E[(Y-WX-b)^2] = \text{tr}(Y^T W^T W Y) - 2\text{tr}(Y^T W^T b) + \text{tr}(b^T b)$ . Also,  $M_W = 0$  since  $X$  is symmetrically distributed around 0.  $E[\sin(W)] = 0$  since antisymmetric function. So,  $b_{WLMPLS} = 0$

Best  $W$  obtained by orthogonality principle  $E[Y - \text{col}(W^T W X + b)]^2 = 0$ . This gives us  $\text{col}(W^T W X) = 0$ .  $E[X^T W^T W X] = \text{Var}(W) = \frac{\pi^2}{12}$

Column sum of  $W$  is  $\text{tr}(W) = \text{tr}(W^T W) = \frac{\pi^2}{12}$

$E[\sin(W)] = \frac{1}{2} \sum_{j=1}^{\pi/2} \sin(j) = \frac{1}{2} \sum_{j=1}^{\pi/2} \frac{1}{2} \sin(2j) = \frac{\pi^2}{12}$

b) Consider the following sub  $D = \{(x_1, y_1), (x_2, y_2), \dots, (x_n, y_n)\} = \{(-1, 1), (1, -1), (-1, -1)\}$

i) Compute  $C_{WLMPLS}, b_{WLMPLS} = \text{argmin}_W E[(Y - WX - b)^2]$

$\hat{W} = \sum_{j=1}^n \hat{x}_j y_j^T, b_{WLMPLS} = \hat{W} - \hat{W} M_W / \text{tr}(W)$

$\hat{A} = \hat{W} M_W = 0, \text{ so } b_{WLMPLS} = 0$

$\hat{X} = \frac{1}{3} \sum_{j=1}^3 x_j y_j^T = \frac{1}{3}$

$\hat{X}^T \hat{X} = \frac{1}{3} \sum_{j=1}^3 x_j^T x_j = \frac{2}{3}$

$\hat{X}^T Y = \frac{1}{3} \sum_{j=1}^3 x_j^T y_j = 1/40$

So,  $b_{WLMPLS} = 2 \times 1/40 = 1/20$

ii) Compute  $C_{WLMPLS}, b_{WLMPLS} = \text{argmin}_W E[(Y - WX - b)^2]$

$\hat{W} = \{(x_1, y_1), (x_2, y_2), (x_3, y_3)\}, W = (w_{ij}, w_{ij})^T$

For  $j=1, 2, 3, \hat{W}_j = \begin{pmatrix} -1 & 1 \\ 1 & -1 \end{pmatrix}, \begin{pmatrix} 1 & -1 \\ -1 & -1 \end{pmatrix}, \begin{pmatrix} -1 & -1 \\ -1 & -1 \end{pmatrix}$

$\hat{M}_W = \frac{1}{3} (\hat{W}_1 + \hat{W}_2 + \hat{W}_3) = \begin{pmatrix} 0 & 0 \\ 0 & 0 \end{pmatrix}$

For  $j=1, 2, 3, \hat{W}_j^T \hat{W}_j = \begin{pmatrix} 1 & -1 \\ -1 & -1 \end{pmatrix}, \begin{pmatrix} -1 & 1 \\ 1 & -1 \end{pmatrix}, \begin{pmatrix} -1 & -1 \\ -1 & -1 \end{pmatrix}$

$\hat{W} = \frac{1}{3} \sum_{j=1}^3 \hat{W}_j \hat{W}_j^T = \frac{1}{3} \left\{ \begin{pmatrix} 1 & -1 \\ -1 & -1 \end{pmatrix} (-1) + \begin{pmatrix} -1 & 1 \\ 1 & -1 \end{pmatrix} (-1) + \begin{pmatrix} -1 & -1 \\ -1 & -1 \end{pmatrix} (-1) \right\}$

$= \begin{pmatrix} 1/3 & 0 \\ 0 & 1/30 \end{pmatrix}$

$\hat{X}^T \hat{Y} = \frac{1}{3} \sum_{j=1}^3 \hat{x}_j^T y_j = \frac{1}{3} \left\{ \begin{pmatrix} -1 & 1 \\ 1 & -1 \end{pmatrix} (-1) + \begin{pmatrix} -1 & 1 \\ 1 & -1 \end{pmatrix} (-1) + \begin{pmatrix} -1 & -1 \\ -1 & -1 \end{pmatrix} (-1) \right\}$

$= \begin{pmatrix} 1/3 & 0 \\ 0 & 1/30 \end{pmatrix}$

So,  $\hat{W} = \hat{X}^T \hat{Y} = \begin{pmatrix} 1/3 & 0 \\ 0 & 1/30 \end{pmatrix} = \begin{pmatrix} 1 \\ 0 \end{pmatrix} = \begin{pmatrix} 1 \\ 0 \end{pmatrix}$

$\hat{b}_{WLMPLS} = \hat{W} M_W = \begin{pmatrix} 1/3 & 0 \\ 0 & 1/30 \end{pmatrix} \begin{pmatrix} 0 & 0 \\ 0 & 0 \end{pmatrix} = \begin{pmatrix} 0 \\ 0 \end{pmatrix}$

$\hat{b}_{WLMPLS} = \hat{W} - \hat{W} M_W / \text{tr}(W) = 0$

So,  $\hat{Y}_{PLS} = \hat{X}^T \hat{Y} = \begin{pmatrix} 1/3 & 0 \\ 0 & 1/30 \end{pmatrix} = \begin{pmatrix} 1/3 & 0 \\ 0 & 1/30 \end{pmatrix} = \begin{pmatrix} 1 \\ 0 \end{pmatrix} = \begin{pmatrix} 1 \\ 0 \end{pmatrix}$

c) Let  $X = y + z$  where  $y \perp z, z \sim N(0, \sigma^2), \sigma > 0$ , and  $\text{Var}(y) = \frac{1}{2} \sigma^2$ . Denote  $\text{hmap}(x)$ , MAP estimate of  $y$  given  $x$

$\text{hmap}(x) = \frac{1}{2} e^{-\frac{|x|}{\sigma^2}} \frac{1}{\sqrt{2\pi}} e^{-\frac{(x-y)^2}{2\sigma^2}}$

$\text{hmap}(x) = \text{argmin}_y \text{hmap}(x, y)$

$= \text{argmin}_y \left[ (x-y)^2 + 2\sigma^2 z^2(y) \right]$

$\text{hmap}(y)$  is a continuous function that is differentiable everywhere except at  $y=0$ .

$\frac{d}{dy} \text{hmap}(y) = \begin{cases} 2(x-y) - 1 + 2\sigma^2 z^2(y) & y \neq 0 \\ \text{not diff} & y=0 \end{cases}$

$\frac{d^2}{dy^2} \text{hmap}(y) = \begin{cases} 2(2x-2y) & y \neq 0 \\ 2(2y-x+1\sigma^2 z^2) & y=0 \end{cases}$

Case 1:  $|x| \leq 1/\sigma^2$

$\frac{d}{dy} \text{hmap}(y) = \begin{cases} < 0 & y < 0, \text{ so } y \text{ strictly decreasing when } y < 0 \\ > 0 & 0 < y, \text{ increasing when } y > 0 \\ \text{not diff} & y=0 \end{cases}$

Case 2:  $|x| > 1/\sigma^2$

$\frac{d}{dy} \text{hmap}(y) = \begin{cases} < 0 & y < 0 \\ \text{not diff} & y=0 \\ < 0 & 0 < y < x-1/\sigma^2 \\ \geq 0 & x-1/\sigma^2 \leq y \end{cases}$

So,  $y$  decreasing for  $y < x-1/\sigma^2$ , increasing for  $x-1/\sigma^2 < y$

4) minimized at  $y = x-1/\sigma^2$  if  $x > 1/\sigma^2$

hmap(x) =  $x-1/\sigma^2$  if  $x > 1/\sigma^2$

Case 3:  $x < -1/\sigma^2$

$\frac{d}{dy} \text{hmap}(y) = \begin{cases} < 0 & y < 0 \\ \geq 0 & 0 < y < x+1/\sigma^2 \\ \text{not diff} & y=0 \\ > 0 & x+1/\sigma^2 < y \end{cases}$

So,  $y$  decreasing for  $y < x+1/\sigma^2$  and increasing for  $x+1/\sigma^2 < y$ . Therefore minimized at  $y = x+1/\sigma^2$

hmap(x) =  $x+1/\sigma^2$  if  $x < -1/\sigma^2$

So,

$\text{hmap}(x) = \begin{cases} 0 & |x| \leq 1/\sigma^2 \\ x-1/\sigma^2 & 1/\sigma^2 < x < x+1/\sigma^2 \\ x+1/\sigma^2 & x > x+1/\sigma^2 \end{cases}$

$\text{hmap}(x) = \begin{cases} 0 & |x| \leq 1/\sigma^2 \\ x-1/\sigma^2 & 1/\sigma^2 < x < x+1/\sigma^2 \\ x+1/\sigma^2 & x > x+1/\sigma^2 \end{cases}$

$\text{hmap}(x) = \text{sign}(x) \cdot \max(0, |x|-1/\sigma^2)$

Distance b/w a point & a line:  
 $\text{distance}(ax + by + c = 0, (x_0, y_0)) = \frac{|ax_0 + by_0 + c|}{\sqrt{a^2 + b^2}}$

**e.g. Find eigenvalues & eigenvectors of a  $2 \times 2$  matrix**

$A = \begin{pmatrix} 0 & 1 \\ -2 & -1 \end{pmatrix}$

$|A - \lambda I| = \begin{vmatrix} 0-\lambda & 1 \\ -2 & -1-\lambda \end{vmatrix} = \begin{vmatrix} 0 & 1 \\ -2 & -1 \end{vmatrix} = 0$

$\begin{vmatrix} 0 & 1 \\ -2 & -1 \end{vmatrix} = 1^2 + 2 = 0 \rightarrow \lambda_1 = -1, \lambda_2 = -2$

$A \cdot v_1 = \lambda_1 \cdot v_1$

$(A - \lambda_1 I) \cdot v_1 = 0$

$\begin{bmatrix} -1 & 1 \\ -2 & -1 \end{bmatrix} \begin{bmatrix} v_{11} \\ v_{12} \end{bmatrix} = 0 \rightarrow \begin{bmatrix} v_{11} \\ v_{12} \end{bmatrix} = \begin{bmatrix} 1 \\ 2 \end{bmatrix}$

$(A - \lambda_2 I) \cdot v_2 = 0$

$\begin{bmatrix} 1 & 1 \\ -2 & -1 \end{bmatrix} \begin{bmatrix} v_{21} \\ v_{22} \end{bmatrix} = 0 \rightarrow \begin{bmatrix} v_{21} \\ v_{22} \end{bmatrix} = \begin{bmatrix} 1 \\ -1 \end{bmatrix}$

(Same procedure for second eigenvalue)

proportion of variance for 1st:  $\frac{1}{1+2} = \frac{1}{3}$

**RCA**

1) Subtract mean from each data dimension (get  $\bar{x}$ ) → center data to get  $\tilde{x}$

2) Calculate covariance matrix

3) Calculate eigenvals/evecs

4) Eigenvect w/ highest eigenval is the PC of the data set

5) Final data =  $(eig_1 \ eig_2 \dots)^T (x_1 \ x_2 \dots)$

6) Back to orig:  $x_i = (eig_1 \ eig_2 \dots)(\text{Final data})$

def: # of features  
 $n$ : # of samples  
 $X, X_i \rightarrow \tilde{Y} \rightarrow$  new transformed matrix  
 $P$ : transformation mat  
 $V = (eig_1 \ eig_2)$

$PX = Y \rightarrow \text{goal}$

$S_k = \frac{1}{n-1} \sum_{i=1}^n \tilde{x}_i \tilde{x}_i^T = \text{covmat}$

choose  $P = VT$

Final Data: data rows in columns, dims on rows

USE UNIT EIGENVECS

Orig data =  $(eig_1 \ eig_2 \dots) \times (\text{Final data}) + \text{orig means}$

avg. measure for grade 1:  $\frac{35}{6}$   
 grade 2:  $\frac{36}{6}$   
 $\text{fit} = \text{project data onto it, minimize distances}$

x	10	11	8	3	21
y	6	4	5	3	2.5

Analy:

$n_{ij} = \# \text{ of points in cluster } i \text{ belonging to class } j$

$n_j = \sum_{i=1}^k n_{ij} = \text{total # of points in cluster } i$

purity of a cluster =  $p_i = \frac{\max n_{ij}}{n_i}$

purity of clustering:  $P = \frac{1}{k} \sum_{i=1}^k p_i$

$= \frac{1}{5} \max_{j=1}^2 n_{ij}$