# Regression + Multivariate Methods

- <u>Regression</u>: relating a series of measurements to one or more possible factors.

<u>extensions</u> * regularized regression, * logistic regression

- <u>Principal Components Analysis</u>: finding a reduced set of variables to summarize a multivariate quantity
  ~ Synonyms: Factor analysis, Karhonen - Loeve decomposition, singular value decomposition

<u>extension</u> "rotations" - varimax, ...
  ICA (Independent components analysis)

- <u>Discriminant Analysis</u> (Fisher discriminant): finding a combination of variables that separate two sets of observations

<u>extension</u>   GIFA   Generalized induced factor analysis

Combination of the above:
  Procrustes analysis: find a linear transformation from one set of multivariate quantities to another

  Canonical correlation analysis: find a linear transformation that best relates (best correlates) one set of multivariate q-ities and another.

<u>Multidimensional Scaling</u>   Embed some points in a vector space to recover specified distances;

All of these can begin with a linear algebra setup, + some can be solved via matrix inversion, some as eigenvalue problems, + some only iteratively.

Basic setup for regression (extensible to PCA)

A "design matrix" $X = \{x_{mn}\}$, known

Observations $Y = \{y_m\}$, known (viewed as a column)

Find the best set of loadings $\{a_n\}$ for which

$$\sum_{n=1}^{N} x_{mn} a_n \approx y_m \qquad\qquad XA \approx Y$$

Convenient to write $y_m^{fit} = \sum_{n=1}^{N} x_{mn} a_n$,

"Best", by default, means that we want to minimize $R^2$

$$R = \sum_m |y_m^{fit} - y_m|^2 = \sum_m \left( \sum_n x_{mn} a_n - y_m \right)^2$$
$$= tr\left( (Y - XA)^T (Y - XA) \right)$$

$$[\text{ note } tr\, M^T M = \sum_{i,j} M_{ij}^T M_{ji} = \sum_{i,j} M_{ij}^2 \,]$$

Could use some other $R$ [ logistic regression ]
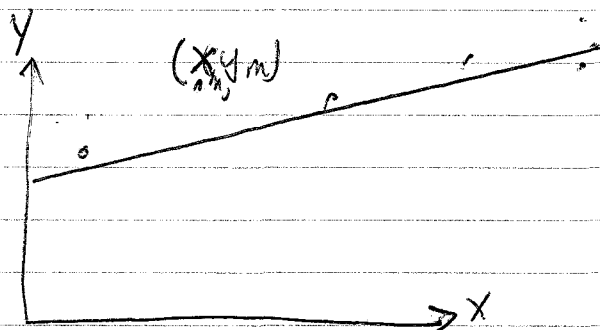
Could put priors on $A$ [ regularized regression ]

---

Applic. to curve-fitting

Find the best line through the data, on,



$$y_m^{fit} = p\, x_m^2 + q\, x_m + r \qquad - \text{ How to put this in above form?}$$

$$x_{m,1} = 1, \quad x_{m,2} = x_m, \quad x_{m,3} = x_m^2 \;\; ; \;\; r = a_1, \; q = a_2, \; p = a_3.$$

③

<u>Apply to fMRI signal analysis</u> ( one pixel at a time )

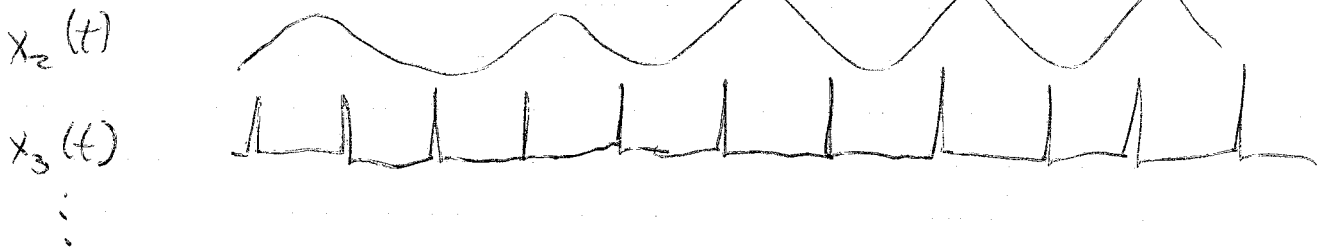Pixel signal is some $y(t)$, discretized as $y_m$

~~~

Each experimental variable is an $X_n(t)$, discretized as $X_{mn}$.

$X_1(t)$ ⌿

⋮

or nuisance variables (EKG, resp)

$X_2(t)$

$X_3(t)$

⋮

Want to write $\quad y(t) \approx \sum a_n X_n(t).$

Our basic regression problem is to minimize

$$R = tr\left[(Y - Y^{fit})^T (Y - Y^{fit})\right] = tr\left((Y - XA)^T (Y - XA)\right)$$

for



$$X = \begin{pmatrix} | & \xrightarrow{n} \\ | & \\ | & m \\ | & \\ \downarrow \end{pmatrix} \begin{pmatrix} | \\ | & n \\ \downarrow & \\ & D \end{pmatrix} \approx \begin{pmatrix} | \\ | & m \\ | & \\ \downarrow \end{pmatrix}$$

$$A \qquad\qquad Y$$

$\#(M) \geq \#(n)$. to hope to find unique $A$, ideally $\#(m) >> \#(n)$.

$Y^{fit}$ can be viewed as projection of $Y$ into the space spanned by the columns of $X$'s.

$$(Y^{fit} - Y) \perp Y^{fit}, \text{ i.e.,} \qquad tr\left((Y^{fit} - Y)^T Y^{fit}\right) = 0$$

$$tr\left(Y^{fit\,T} Y^{fit}\right) = tr\left(Y^T Y^{fit}\right)$$

$$R = tr\left((Y - Y^{fit})^T Y\right) = tr\left(Y^T Y\right) - tr\left(Y^{fit\,T} Y^{fit}\right)$$

So minimizing $R$ is the same as maximizing $tr\left(Y^{fit\,T} Y^{fit}\right)$, i.e., maximizing the length of the projection of $Y$.

We'll minimize $R$ by setting $\dfrac{\partial R}{\partial A}$ to 0.

⑤

$$R = \text{tr}(Y^T Y) - \text{tr}((XA)^T Y) - \text{tr}(Y^T XA) + \text{tr}((XA)^T XA)$$

What is $\dfrac{\partial}{\partial a_k}(\text{tr}\, QA)$ for $A$ a column?

$$\text{tr}(QA) = \sum_i q_{ik}\, a_{k1}, \text{ so } \frac{\partial}{\partial a_k}(\text{tr}\, QA) = q_{ik}$$

Think of the $\dfrac{\partial}{\partial a_k}(\text{tr}\, QA)$ is forming a column:

$$\begin{pmatrix} \frac{\partial}{\partial a_1}\,\text{tr}(QA) \\ \vdots \\ \frac{\partial}{\partial a_n}\,\text{tr}(QA) \end{pmatrix} = \begin{pmatrix} q_{11} \\ \vdots \\ q_{1n} \end{pmatrix}$$

$$\left[ \frac{\partial}{\partial a_k}\,\text{tr}(QA) \right] = Q^T \qquad (Q\ a\ \text{row}).$$

Note $\text{tr}(Y^T XA) = \text{tr}(A^T X^T Y) = \text{tr}((XA)^T Y)$

Also, we can use the product rule to find that $\dfrac{\partial}{\partial a_k}(\text{tr}\, A^T G A)$

$$\text{a column of } 2(A^T G)^T = 2\, G^T A$$

So $\dfrac{\partial R}{\partial a_k} = -2(Y^T X)^T + 2\, \underbrace{X^T X}_{n \times n} A \qquad (G = X^T X)$

$$\frac{\partial R}{\partial a_k} = 0 \implies X^T X A = (Y^T X)^T$$

$$A = (X^T X)^{-1} X^T Y$$

From this,

$$y^{fit} = XA = \left[ X(X^TX)^{-1}X^T \right] Y$$

the projection into the space spanned by the columns of $X$.

Simple extension: residual regression + related.

Why do we use $R = \sum_i \left( y^{fit} - y \right)^2$ ?

Ⓐ It leads to a linear problem we can solve

Ⓑ $e^{-R/2\sigma^2}$ can be interpreted as the probability of the observations $Y$, given that $y^{fit}$ should have been observed (i.e., $XA$ is the model), or each observation $Y$ independently deviates from $y^{fit}$, assuming the measured error drawn from a Gaussian of stddev $\sigma$.

Thus, minimizing $R$ maximizes the a posteriori probability of the $A$'s.

But say that we knew that the $A$'s came from a distribution with covariance $C_{A}$, i.e.,

$$p(A) \sim e^{-A^T(C_A)^{-1}A/2}$$

[ For example, $C_A = \sigma_A^2 I$ — $A$'s not "too big" ]

And the noise might not be independent:

$$\rho(y^{fit}-y) \sim e^{-(y^{fit}-y)C_y^{-1}(y^{fit}-y)/2}$$

Now we need to maximize

$$e^{-A^T(C_A)^{-1}A/2} \; e^{-(y^{fit}-y)C_y^{-1}(y^{fit}-y)/2}$$

i.e, minimize

$$(y^{fit}-y)^T C_y^{-1}(y^{fit}-y) + A^T C_A^{-1} A$$

$$= (XA-y)^T C_y^{-1}(XA-y) + A^T C_A^{-1} A$$

$\dfrac{\partial}{\partial a_n} = 0$ leads to

$$-(y^T C_y^{-1} X)^T + X^T C_y^{-1} XA + C_A^{-1}A = 0$$

$$A = \left(X^T C_y^{-1} X + C_A^{-1}\right)^{-1}\left(X^T C_y^{-1} y\right)$$

$C_A$ large $\Rightarrow$ its effect goes away

$C_A$ small $\Rightarrow$ A forced to be small.

$C_y$ decorrelates the errors.

Note that if we have a series of regression problems
with the same X's but diff, Y's, we consider
them in parallel.

$$X = \left(\Bigg|\xrightarrow{\;n\;}\Bigg\downarrow_m\;\right)\left(\Bigg|\xrightarrow{\;r\;}\Bigg\downarrow_n\;\right) = \left(\Bigg|\xrightarrow{\;r\;}\Bigg\downarrow_m\;\right)$$

$$\qquad\qquad\qquad\qquad A\qquad\qquad\qquad Y$$

Each column of Y is a separate regression. $R = \text{sum } R's$ for each
                                                            column
Still have same solution, column by column, since the columns don't
interact:

$$A = (X^T X)^{-1} X^T Y.$$

$$\boxed{\text{Regression} \longrightarrow PCA}$$

Say we want to deduce a good set of X's, i.e.

write $\quad y^{fit}_{mr} = \sum x_{mn} a_{nr}$ , with $n$ small.

Can view columns of Y as time series, each col. is a pixel
                                          or each col. is an electrode

OR exchange rows & columns.
   Each col is a "snapshot"

The solution will have to be ambiguous.

If $Y^{fit} = XA$, and $Q$ is any $n \times n$ invertible matrix,

$$Y^{fit} = (XQ)(Q^{-1}A) = X'A', \quad \text{for}$$

$$X' = XQ, \quad A' = Q^{-1}A.$$

We could partially resolve this ambiguity by requiring $X$ to be orthogonal (i.e, apply Gram-Schmidt if it wasn't).

But the above still is an ambiguity, since $X' = XR$ is orthogonal for $R$ any unitary matrix.

$$\left[ X \text{ has orthonormal cols} \Leftrightarrow X^T X = I_n \right.$$

$$(X')^T X' = (XR)^T XR = R^T X^T X R = R^T R = I.$$

So we really should think of this as a search for the subspace spanned by the columns of $X$.

Same argument can be made for $A$ – its rows can always be made orthogonal.

Leads to a more symmetric statement of the problem:

$$Y^{fit} = B^T \Lambda A \quad \times \quad \text{i.e,} \quad Y^{fit}_{mr} = \sum_{i=1}^{r} \lambda_n b_{nm} a_{nr}$$

$$B = b_{nm}, \quad \Lambda = \begin{pmatrix} \lambda_1 & & \\ & \ddots & \\ & & \lambda_r \end{pmatrix}, \quad A = \begin{pmatrix} \xrightarrow{n} \\ r \end{pmatrix}$$

$$BB^T = I$$

$$AA^T = I. \quad \text{Soln is unique up to An order of}$$

$$\text{the } \lambda\text{'s.}$$

How to do it?

Minimize $\text{tr}\left[(Y-XA)^T(Y-XA)\right]$

$\left(\downarrow_m \overrightarrow{n}\right)\left(\downarrow_n \overrightarrow{r}\right) = \left(\downarrow_m \overrightarrow{r}\right)$

$\quad\quad X \quad\quad A \quad\quad\quad\quad Y$

over $X$ and $A$. With $X$ known, $A = (X^TX)^{-1}X^TY$

We can keep $X$ orthonormal (in columns) so $X^TX = I_{n \times n}$.

so $A = X^TY$.

$(Y - XA)^T(Y-XA) = Y^TY - Y^TXA - A^TX^TY + A^TX^TXA$

$\quad = Y^TY - Y^TXX^TY - Y^TXX^TY + (Y^TX)(X^TX)(X^TY)$

$\quad = Y^TY - Y^TXX^TY$

$\text{tr}\left((Y-XA)^T(Y-XA)\right) = \text{tr}\left(Y^TY - Y^TXX^TY\right)$

minimizing this now maximizing $\text{tr}(Y^TXX^TY) = \text{tr}(YY^TXX^T)$

$\quad\quad\quad\quad\quad\quad\quad = \text{tr}(X^TYY^TX)$

$YY^T$ is $m \times m$, and symmetric. Let's write it via its eigenvectors + eigenvalues.

$YY^T = \sum_{h=1}^{r} \lambda_h \phi_h \phi_h^T$ , $\phi$'s orthonormal, with

$\quad\quad\quad\quad (YY^T)\phi_h = \lambda_h \phi_h$.

Orthonormality means $\phi_k^T \phi_\ell = \delta_{k\ell}$.

Let's order the $\lambda$'s so $\lambda_1 \geq \lambda_2 \geq \lambda_3 \cdots$

Now, let's consider the $n=1$ -case.

$$X = \sum_h z_h \phi_h.$$

$$X^T Y Y^T = \sum_h z_h \phi_h^T \sum_k \lambda_k \phi_k \phi_k^T$$

$$= \sum_{h,k} z_h \lambda_k \phi_h^T \phi_k \phi_k^T$$

$$= \sum_k z_k \lambda_k \phi_k^T.$$

$$\boxed{\phi_h^T \phi_k = \delta_{hk}} \ !$$

So $X^T Y Y^T X = \left( \sum_k z_k \lambda_k \phi_k^T \right) \sum_h z_h \phi_h$

$$= \sum_{k,h} z_k \lambda_k z_h \phi_k^T \phi_h$$

$$= \sum_k z_k^2 \lambda_k$$

How do we maximize $\sum_k z_k^2 \lambda_k$ subj to $\sum z_k^2 = 1$?

Take $z_1 = 1$   (largest $\lambda \doteq \lambda_1$)

oths $= 0$.    So $X = \phi_1$.

$n \geq 2$

$$X = \left( \sum_h z_{h,1} \phi_h \ \middle| \ \sum_h z_{h,2} \phi_h \ \middle| \ \cdots \ \sum_h z_{h,n} \phi_h \right)$$

$$X^T Y Y^T = \begin{pmatrix} \sum_k z_{k,1} \lambda_k \phi_k^T \\ \sum_k z_{k,2} \lambda_k \phi_k^T \\ \vdots \\ \sum_k z_{k,n} \lambda_k \phi_k^T \end{pmatrix}$$

$$tr(X^T Y Y^T X) = \sum_k z_{k,1}^2 \lambda_k + \sum_k z_{k,2}^2 \lambda_k + \cdots + \sum_k z_{k,n}^2 \lambda_k$$

Coef of $\quad \lambda_1$ is $z_{1,1}^2 + \cdots + z_{1,n}^2$ , max possible 1.

$\lambda_2$ is $z_{2,1}^2 + \cdots + z_{2,n}^2$

So we can make the first $n$ $\lambda$'s have a cf of 1, by choosing

$$X = \left( \phi_1 \ \middle| \ \phi_2 \ \middle| \ \cdots \ \phi_n \right), \text{the } \overset{\text{first } n}{\text{eigenvectors of } YY^T}$$

$YY^T$ is $m \times m$.

$$\boxed{YY^T X = X \begin{pmatrix} \lambda_1 & & \\ & \ddots & \\ & & \lambda_n \end{pmatrix}}$$

$A = X^T Y$

Rows of $A$ are Left eigenvectors of $Y^T Y$. $\qquad$ because cols of $X$ are eiv's of $YY^T$

$$A Y^T Y = (X^T Y) Y^T Y = X^T (Y Y^T) Y = \begin{pmatrix} \lambda_1 & & \\ & \lambda_2 & \\ & & \ddots \\ & & & \lambda_n \end{pmatrix} X^T Y = \begin{pmatrix} \lambda_1 & & \\ & \ddots & \\ & & \lambda_n \end{pmatrix} A$$

(13)

Symmetric form of sdd(m):

$$Y^{fit} = B^T \Lambda A \quad \text{where}$$

A: $n \times r$
B: $n \times m$
$\Lambda$: $n \times n$, diagonal

$n$ rows of $A$ are left eigenvecs of $Y^T Y$ $(r \times r)$, $AA^T = I$.

$n$ cols of $B^T$ are right eigenvecs of $Y Y^T$ $(m \times m)$, $BB^T = I_n$

$(\Rightarrow n$ rows of $B$ are left eigenvecs of $Y Y^T$

$n$ cols of $A^T$ are right eigenvecs of $Y^T Y)$

$$\Lambda = \begin{pmatrix} \sqrt{\lambda_1} & & \\ & \ddots & \\ & & \sqrt{\lambda_n} \end{pmatrix}, \text{ since } (B^T \Lambda A)(B^T \Lambda A)^T = Y Y^T$$

$$= B^T \Lambda A A^T \Lambda B$$

$$= B^T \Lambda^2 B$$

so eigvals $\Lambda^2$ must be eigs $Y Y^T$.

If $m$ r are very different, $Y Y^T$, $Y^T Y$ are very different in size, + time to diagonalize differs. Always diagonalize the smaller one!

The Elegant Approach (avoiding coordinates) to PCA.

We'll "need" "Lagrange Multipliers" - a general approach for solving constrained minimization problems. Here, it will turn quadratic minimization with quadratic constraints into eigenvalue problems.
[ Also, by role in stat. physics + into theory ]

<u>Recipe</u>  Say you want to minimize $F(x_1, x_2, \cdots, x_k)$
(or maximize)

subject to constraints $C_1(x_1, \cdots, x_k) = 0$
$$\vdots$$
$$C_L(x_1, \cdots, x_k) = 0.$$

Instead of trying to write $x_1 = G_1(x_{L+1}, \cdots, x_k)$     $\begin{bmatrix} \text{The "obvious" approach -} \\ \text{use constraints to} \\ \text{eliminate variables} \end{bmatrix}$
$$\vdots$$
$$x_L = G_L(x_{L+1}, \cdots, x_k)$$
$\underbrace{\qquad\qquad\qquad}_{\text{L variables eliminated via constraints}}$ $\underbrace{\qquad}_{\text{K-L free variables}}$

+ minimizing $\underbrace{F(G_1(x_{L+1}, \cdots, x_k), \cdots, G_L(x_{L+1}, \cdots, x_k), x_{L+1}, \cdots, x_k)}$

L.M. says:  minimize $\mathcal{F}(x_1, \cdots, x_k) + \sum_{\ell=1}^{L} \lambda_\ell \, C_\ell(x_1, \cdots, x_k)$

and find the $\lambda$'s which satisfy the constraints.


<u>Example</u>  Maximize $\sum x_i a_i$ subject to $\sum b_i x_i^2 = 1$  $\begin{bmatrix} \text{one} \\ \text{constraint} \end{bmatrix}$

$\mathcal{F} = \sum x_i a_i + \lambda(\sum b_i x_i^2 - 1)$     $\dfrac{\partial F}{\partial x_i} = a_i + 2 x_i \lambda b_i.$

So $\dfrac{\partial F}{\partial x_i} = 0 \Rightarrow x_i = -a_i / 2\lambda b_i.$     $\lambda = \pm \frac{1}{2}\sqrt{\sum \dfrac{a_i^2}{b_i}}$

Now find $\lambda$.  $\sum b_i x_i^2 = 1 \Rightarrow \dfrac{1}{(2\lambda)^2} \sum \dfrac{a_i^2}{b_i} = 1 \Rightarrow$

Obs: 1. sometimes you in turn need to find $\lambda$.  $x_i/x_j = \dfrac{a_i}{b_i} \cdot \dfrac{b_j}{a_j}$  $\begin{bmatrix} \text{sometimes only} \\ \lambda \text{ is hard} \end{bmatrix}$
2. Problem stays symmetric.

Why does it work? Try simpler: 2 variables, one constraint.

Extremize $F(x,y)$ subj. to $G(x,y) = 0$. Say $G(x,y) = 0 \Rightarrow x = H(y)$.

"Standford" approach: set $\frac{dF}{dy} = 0$.

$$\frac{dF}{dy} = \frac{\partial}{\partial y}\left(F(H(y),y)\right) = \frac{\partial F}{\partial x}\frac{\partial H}{\partial y} + \frac{\partial F}{\partial y}.$$

$$G(x,y) = 0 \Rightarrow G(H(y),y) = 0 \Rightarrow \frac{\partial G}{\partial y} = 0 \Rightarrow \frac{\partial G}{\partial x}\frac{\partial H}{\partial y} + \frac{\partial G}{\partial y} = 0.$$

$$\Rightarrow \frac{\partial H}{\partial y} = -\frac{\partial G}{\partial y}\Big/\frac{\partial G}{\partial x}.$$

So we need to solve $\frac{dF}{dy} = 0$ $\boxed{\text{i.e., } -\frac{\partial F}{\partial x}\frac{\partial G}{\partial y} + \frac{\partial F}{\partial y}\frac{\partial G}{\partial x} = 0.}$

L.M. method: Solve $\frac{\partial}{\partial x}\left(F(x,y) + \lambda G(x,y)\right) = 0$

$$\frac{\partial}{\partial y}\left(F(x,y) + \lambda G(x,y)\right) = 0$$

$$\begin{cases} \dfrac{\partial F}{\partial x} + \lambda \dfrac{\partial G}{\partial x} = 0 \\[2em] \dfrac{\partial F}{\partial y} + \lambda \dfrac{\partial G}{\partial y} = 0 \end{cases} \Rightarrow \lambda = -\frac{\partial F}{\partial x}\Big/\frac{\partial G}{\partial y}$$

$$\Rightarrow \frac{\partial F}{\partial x} + \left(-\frac{\partial F}{\partial y}\Big/\frac{\partial G}{\partial y}\right)\frac{\partial G}{\partial x} = 0$$

(in the box).

Everything goes through with multiple variables or constraints, $\frac{\partial G}{\partial x} \to$ matrix of partials.

Applying LM's to the PCA problem:

Maximize $tr(YY^TXX^T)$ subject to $X^TX = I_{n\times n}$

View $X^TX = I_{n\times n}$ as a symmetric matrix of constraints, $\vec{x}_i \cdot \vec{x}_j = \delta_{ij}$

Each constraint is paired with a $\lambda_{ij}$, so $\Lambda$ = matrix of $\lambda_{ij}$'s is symmetric

LM formulation = to maximize $tr(YY^TXX^T) - tr(\Lambda X^TX) = \mathcal{I}$

View $\frac{\partial}{\partial x_{uv}} \mathcal{I} = 0$ as a matrix of equations.

what's $\frac{\partial}{\partial x_{uv}} (tr\, M X^TX)$?

$$\frac{\partial}{\partial x_{uv}} tr(MX^TX) = \frac{\partial}{\partial x_{uv}} \left( \sum_{i,j,k} m_{ij} (X^TX)_{ji} \right) = \frac{\partial}{\partial x_{uv}} \sum_{i,j,k} m_{ij} x_{kj} x_{ki}$$

$$= \sum_{i,j,k} m_{ij} \left( \frac{\partial}{\partial x_{uv}} x_{kj} \right) x_{ki} + \sum_{i,j,k} m_{ij} x_{kj} \frac{\partial}{\partial x_{uv}} (x_{ki})$$

$$= \sum_{\substack{i,j,k \\ k=u, j=v}} m_{ij} x_{ki} + \sum_{\substack{i,j,k \\ k=u \\ i=v}} m_{ij} x_{kj} = \sum_i m_{iv} x_{ui} + \sum_j m_{vj} x_{uj}$$

$$= (XM + XM^T)_{uv}. \quad \text{Take } M = \Lambda \,(=M)$$

Similarly, $\frac{\partial}{\partial x_{uv}} (tr\, M XX^T) = (MX + M^TX)_{uv}$.  Take $M = YY^T \,(=M^T)$

So the LM formulation,  $YY^TX = X\Lambda$, ~~subject~~ ~~$X^TX = I$~~

simultaneous with $X^TX = I$. This solves for $\Lambda = diag(eigs\, of\, YY^T)$
and $X$ = eigenvectors of $YY^T$.

("Guess" that $\Lambda$ is diagonal).

Another example of quadratic quantity to extremize, ā quadratic constraint.
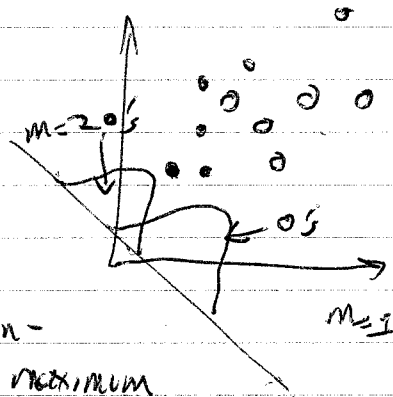
Fisher Discriminant & Canonical Variates.

Setup: multivariate quantities $Y =$

$$m \left\Vert \begin{pmatrix} \vert \\ \vert \\ \end{pmatrix} \begin{pmatrix} \vert \\ \vert \\ \end{pmatrix} \cdots \begin{pmatrix} \\ \\ \end{pmatrix} \right. \quad \text{i.e., observations } \vec{y}_1, \cdots, \vec{y}_r$$

Say we know a priori that some of the $\vec{y}$'s are in category $\underline{1}$

etc. for $\underline{c}$ categories.

We want to find linear combinations of the coordinates $m$ that do the best job of segregating the $\vec{y}$'s.

E.g., $c = 2$ (Fisher case)



More formally, find $x_1, \cdots, x_m$ s.t.

$x^T \vec{y}$ 's have the minimum within-group variance & the maximum between-group variance. Equal-sized groups (for simplicity).

[in assume $\Sigma \vec{y} = 0$.]

$\underline{\text{Say}}$ $\vec{y}_1, \cdots, \vec{y}_{r_1}$ in category $\underline{1}$, with mean $\vec{\mu}_1 = \begin{pmatrix} \mu_{11} \\ \vdots \\ \mu_{1m} \end{pmatrix}$

$\vec{y}_{r_1+1}, \cdots, \vec{y}_{r_1+r_2}$ " " $\underline{2}$ " " $\vec{\mu}_2$

$\vec{y}_{r_1+r_2+\cdots+r_{c-1}+1}, \cdots, \vec{y}_{r_1+r_2+\cdots+r_c}$ in cat $\underline{c}$, with mean $\vec{\mu}_c = \begin{pmatrix} \mu_{c1} \\ \vdots \\ \mu_{cm} \end{pmatrix}$

Maximize $\sum_{jk} [x_j (\vec{\mu}_k - \vec{\mu})_j]^2$ subject to $\sum_{j,h} [x_j (\vec{y}_h - \vec{\mu}_{c(h)})_j]^2$

$[\vec{\mu} = \text{global mean, may not be 0 if groups are unequal}]$
$[c(h) = \text{category of } h]$

Previous strategy turns this into

$$S_g X = S_w X \Lambda \qquad [\text{"generalized eigenvalue problem"}]$$

where $S_g$ = covariance matrix of group means

$S_w$ = covariance matrix within group

$$S_g = \sum_k (\bar{\mu}_k - \bar{\mu})(\bar{\mu}_k - \bar{\mu})^T \;,\quad S_w = \sum_h (\bar{y}_h - \bar{\mu}_{c(h)})(\bar{y}_h - \bar{\mu}_{c(h)})^T$$

$\left(\text{Note that if } S_w = I, \text{ then } X = \sum_k (\bar{\mu}_k - \bar{\mu}) q_k \text{ will solve.}\right)$

Two "flavors" of interest: Ⓐ $c$ categories, top $c-1$ eigenvectors $X$ yield "best" ~~projection~~ linear map of data into a $c-1$-dimensional plane [in which categories separate best by between-group-variance]



This is "canonical variate", $c = 2$ is the "Fisher Discriminant".

$C = 2$: Further simplification, since $\bar{\mu}_1 \sim -\bar{\mu}_2$, so $S_g = 4(\mu^T \mu)$, recognize above as $S_w^{-1} S_g X = X \Lambda$ $\begin{bmatrix} \text{Many observations,} \\ \text{few covariates } m \end{bmatrix}$

Ⓑ Two categories, but consider more than just the leading eigenvector. [GIFA = generalized indicator factor analysis, Yahoo etal.]



This yields all of the linear mappings that discriminate the two categories. Each explains successively less of the variance (as $\lambda$ decreases). Choose some cutoff $\lambda_0$, select only the $\lambda$'s $> \alpha$, and construct

$$\sum X(\lambda) f(\gamma - \lambda)$$

$\Rightarrow$ the composite "discriminating image"
$\sim$ "Generalized" - replace $S_g$ by $S_g + \alpha S_w$

$\begin{bmatrix} \text{This is used in the imaging} \\ \text{context, where } S_w \text{ is} \\ \text{singular, so you can't} \\ \text{calculate } S_w^{-1}. \\ m \gg r \end{bmatrix}$