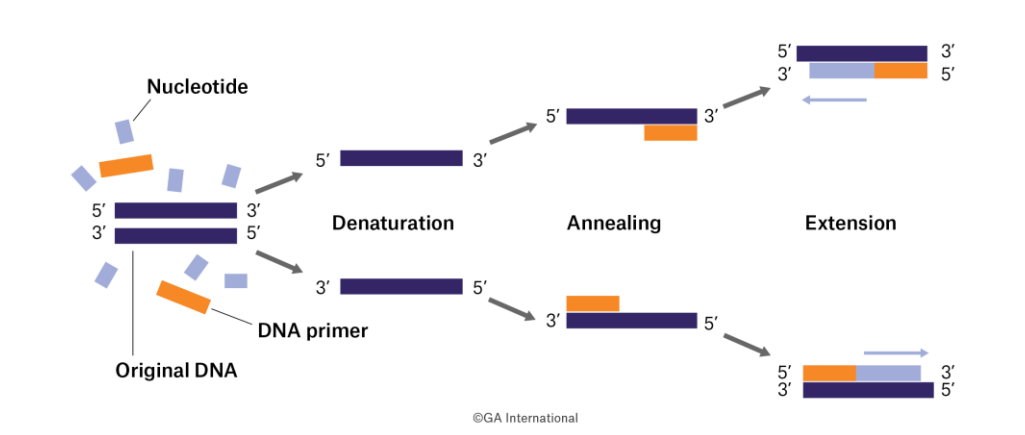


Chemical Engineering 4H03

Principal Component Regression (PCR)

Jake Nease
McMaster University



Portions of this work are copyright of ConnectMV

<https://blog.labtag.com/a-brief-history-of-pcr-and-its-derivatives/>

Think Back to Our Data Sets...

- What kinds of **output** variables did we have?
 - Perhaps user scores or ratings for pastries?
 - Pizza prices?
 - Others?
- Outcomes can be **continuous** or **categorical**
 - Continuous
 - Density
 - Reactor yield
 - Strain required for fracture
 - Absorbance/reflectance of surface material
 - Final course grade (as %)
 - Categorical
 - Pass/fail (also for course grades but also lots of processes!)
 - Good/marginal/bad (or another Likert-style rating)

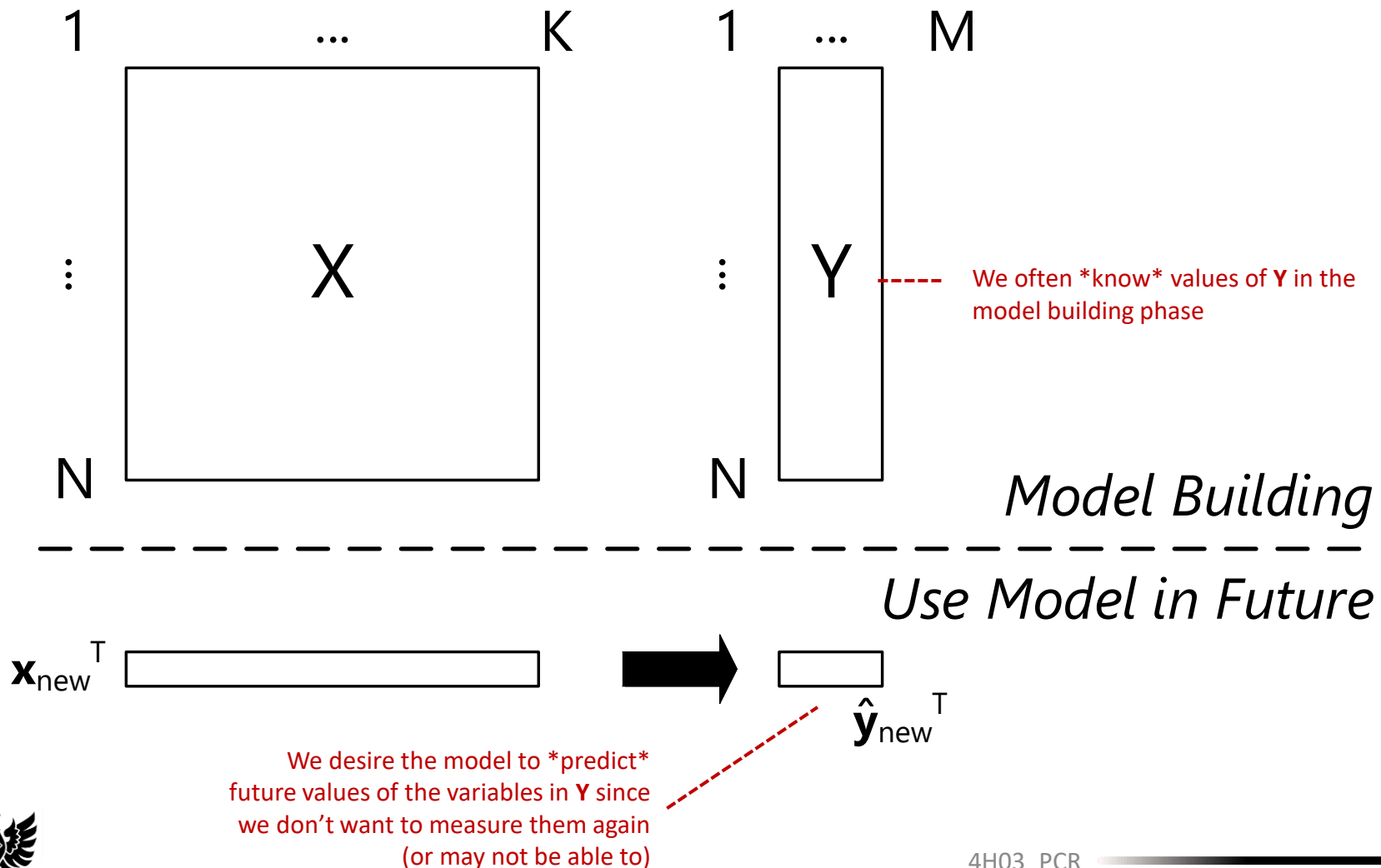


A Primary Goal of Modeling

- We always want to learn from and visualize data
 - Tools like PCA help us do this nicely
- However, we frequently want to **predict or monitor an outcome based on known quantities**
 - For this, we typically build a model or regression
- Lucky for us, there are TWO things PCA/LVMs can give us in this regard:
 - Dimension reduction → regression
 - Dimension reduction **and** regression simultaneously



A Primary Goal of Modeling



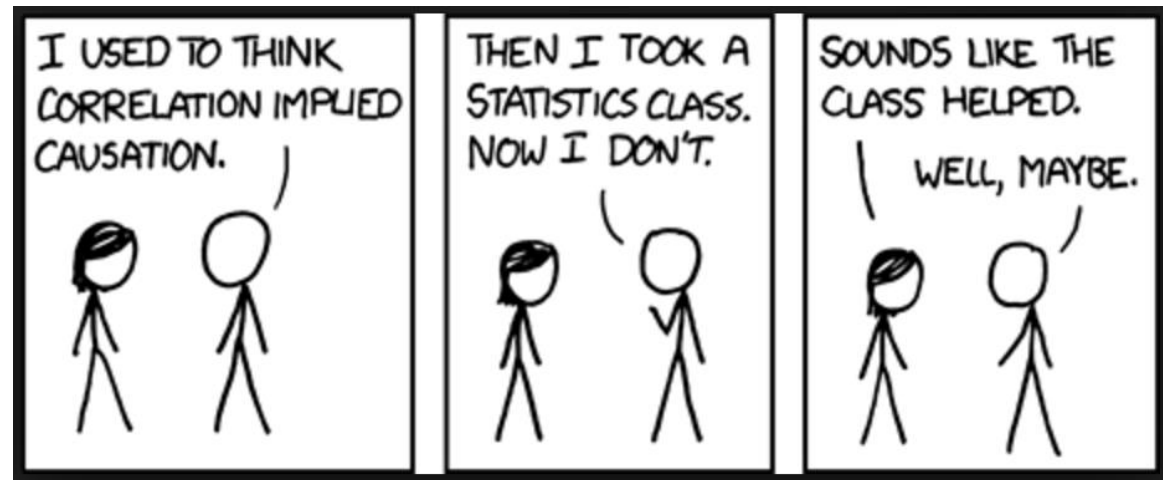
Objectives for this Class

- Now that we know how to compute a PCA model, can we use that information to predict an *outcome* **Y**?
- We sure can. There are two main ways:
 - **P**roincipal **C**omponent **R**egression (PCR) is the poor cousin to...
 - **P**rojection of **L**atent **S**tructures (PLS)
- **This class:** PCR
 1. Mathematical review of covariance
 2. Advantages of PCR to straight-up linear regression
 3. Methods of calculating (long way and shortcuts)



Review: Covariance

... But, like, an applied version of covariance



Let's Consider a Data Set

	Temperature (K)	Pressure (kPa)	Humidity (%)
	273	1600	42
	285	1670	48
	297	1730	45
	309	1830	49
	321	1880	41
	333	1920	46
	345	2000	48
	357	2100	48
	369	2170	45
	381	2200	49
MEAN	327	1910	46.1
VARIANCE	1188	38940	7.29



Review: Covariance

- The covariance of two variables x and y is:

$$\mathcal{C}\{x, y\} = \mathcal{E}\{(x - \bar{x})(y - \bar{y})\}$$

Recall \bar{z} is the mean
for a random variable

- $\mathcal{E}\{z\}$ is the so-called “expected outcome” of z and is \bar{z}
- The covariance of x with itself is its **variance**

$$\mathcal{C}\{x, x\} = \mathcal{V}(x) = \mathcal{E}\{(x - \bar{x})(x - \bar{x})\} \text{ ----- } \frac{\sum_i (x_i - \bar{x})^2}{n}$$

- (Co)variance of an centered vector is precisely equal to the (co)variance of an uncentered vector (why?)
 - Covariance attempts to describe tendencies of two variables



Example: Covariance Matrix $\mathcal{C}\{x, y\} = \mathcal{E}\{(x - \bar{x})(y - \bar{y})\}$

- A covariance matrix is a **real, symmetric** matrix
 - Shows **variances** on the diagonal
 - Shows **covariances** between observations on off-diagonals
- For our observations:

	Temperature	Pressure	Humidity
Temperature	1188	6780	35.4
Pressure	6780	38940	202
Humidity	35.4	202	7.29

```
>> C = cov(A,1)
```

```
C =
```

```
1.0e+04 *
```

```
0.1188    0.6780    0.0035  
0.6780    3.8940    0.0202  
0.0035    0.0202    0.0007
```

- Easy to do in MATLAB: use `cov(A)` on matrix A
 - **NOTE:** the "1" means population covariance (normalized n) rather than sample covariance (normalized to $n - 1$)



Review: Correlation

- Covariance depends on **units** of x and y
- Correlation $r\{x, y\}$ removes this effect by scaling covariance by the root-product of each variance

$$r\{x, y\} = \frac{\mathcal{E}\{(x - \bar{x})(y - \bar{y})\}}{\sqrt{\mathcal{V}(x)\mathcal{V}(y)}} = \frac{\mathcal{C}\{x, y\}}{\sqrt{\mathcal{V}(x)\mathcal{V}(y)}}$$

- A couple of notes:
 - $-1 \leq r\{x, y\} \leq 1$
 - Result is dimensionless



Example: Correlation Matrix $r\{x, y\} = \frac{C\{x, y\}}{\sqrt{V(x)V(y)}}$

- A correlation matrix is a **real, symmetric** matrix
 - Shows **unity** on the diagonal (why?)
 - Shows **correlation** between observations on off-diagonals
- For our observations:

	Temperature	Pressure	Humidity
Temperature	1	0.997	0.38
Pressure	0.997	1	0.379
Humidity	0.38	0.379	1

```
>> R = corrcoef(A)
```

```
R =
```

```
1.0000    0.9968    0.3804
0.9968    1.0000    0.3791
0.3804    0.3791    1.0000
```

- Easy to do in MATLAB: use `corrcoef(A)` on matrix A



Reminder: Linear Regression

- You may recall the LSOE that can be solved for the regression coefficients in our regression $\hat{y} = a_0 + a_1x$:

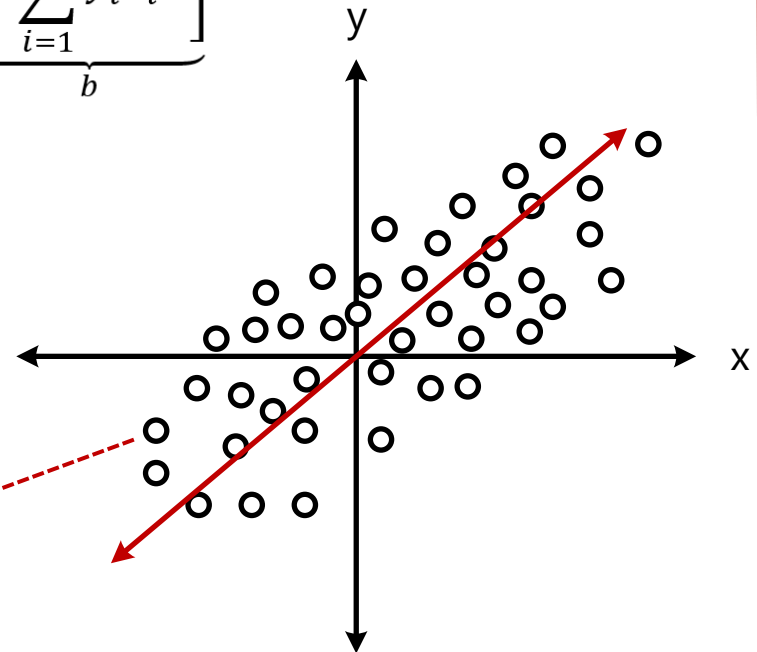
$$\underbrace{\begin{bmatrix} \sum_{i=1}^N 1 & \sum_{i=1}^N x_i \\ \sum_{i=1}^N x_i & \sum_{i=1}^N x_i^2 \end{bmatrix}}_A \underbrace{\begin{bmatrix} a_0 \\ a_1 \end{bmatrix}}_x = \underbrace{\begin{bmatrix} \sum_{i=1}^N y_i \\ \sum_{i=1}^N y_i x_i \end{bmatrix}}_b$$

- RECALL**

- We will center our data
- Thus $a_0 = 0$

Distance from point to model (line) is ϵ_i , the error between observation y_i and prediction \hat{y}_i

THIS SHOULD BE REVIEW AT THIS POINT. STOP ME IF IT IS STILL UNCLEAR

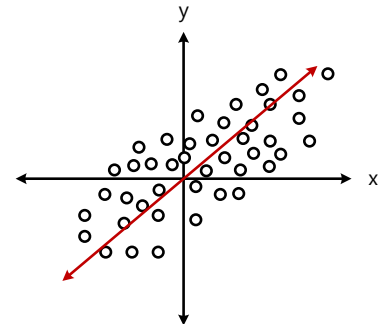


Reminder: LR in Vector Form

- Our “system of equations” in that case simply boils down to one equation and one unknown:

$$\underbrace{\begin{bmatrix} \sum_{i=1}^N x_i^2 & \sum_{i=1}^N x_i y_i \\ \sum_{i=1}^N x_i y_i & \sum_{i=1}^N y_i^2 \end{bmatrix}}_A \underbrace{\begin{bmatrix} a_1 \\ x \end{bmatrix}}_x = \underbrace{\begin{bmatrix} \sum_{i=1}^N y_i x_i \\ \sum_{i=1}^N y_i^2 \end{bmatrix}}_b$$

$$a_1 = \frac{\sum_{i=1}^N y_i x_i}{\sum_{i=1}^N x_i^2}$$



- Fun fact, these sums of the elements in each vector can pretty easily be written as dot-products

$$\sum_{i=1}^N y_i x_i = y^T x$$

$$\sum_{i=1}^N x_i^2 = x^T x$$

- We can thus collapse our equation for a_1 into:

$$a_1 = \frac{x^T y}{x^T x}$$

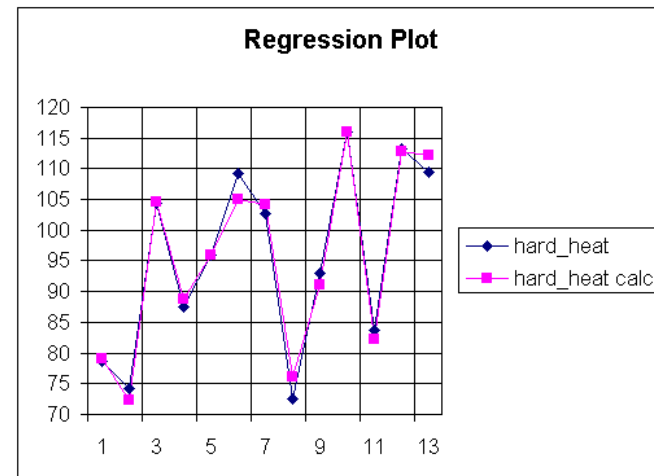
At this point, I should not need to convince you that this expression is identical to the one for a_1 above



Multi-Linear Regression

It's like what we did before but simplified

BOY there sure are some brutal examples of this stuff 😊



Multi-Linear Regression (MLR)

- The general form of a multi-linear regression (MLR) model for some observation i is a simplified version of our multidimensional basis function regression:
 - Each column of \mathbf{X} has been centered (maybe scaled)

$$y_i = a_1 x_1 + a_2 x_2 + \cdots + a_K x_K + \epsilon_i$$

$$y_i = [x_1 \quad x_2 \quad \cdots \quad x_K] \begin{bmatrix} a_1 \\ a_2 \\ \vdots \\ a_K \end{bmatrix} + \epsilon_i$$

These are all the same

The model for y_i represents a hyperplane in K dimensions (but you may have already known that)

$$y_i = \mathbf{x}^T \mathbf{a} + \epsilon_i$$



Multi-Linear Regression (MLR)

- For multiple rows:

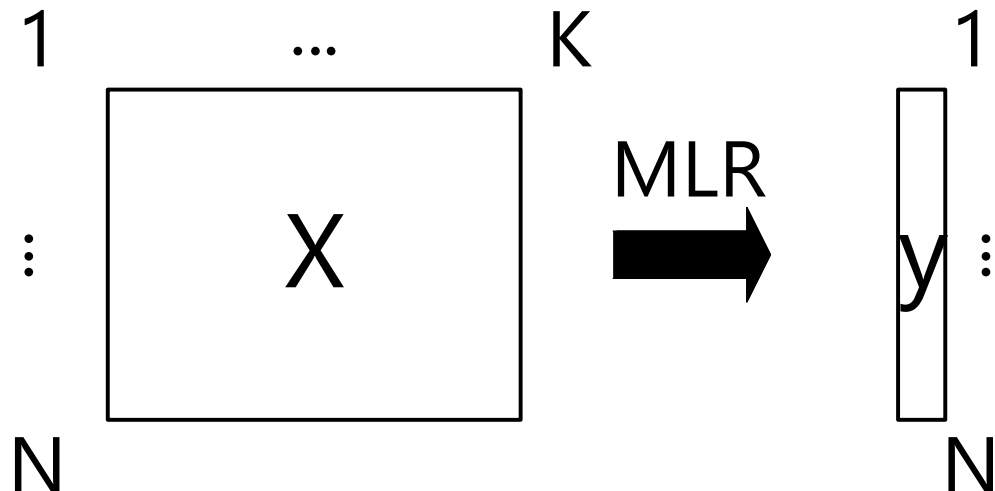
IMPORTANT

X and a DO NOT need to have the same length (actually, they SHOULD NOT have the same length)... WHY?

$$\begin{bmatrix} y_1 \\ y_2 \\ \vdots \\ y_N \end{bmatrix} = \begin{bmatrix} x_{1,1} & x_{1,2} & \cdots & x_{1,K} \\ x_{2,1} & x_{2,2} & \cdots & x_{2,K} \\ \vdots & \vdots & \ddots & \vdots \\ x_{N,1} & x_{N,2} & \cdots & x_{N,K} \end{bmatrix} \begin{bmatrix} a_1 \\ a_2 \\ \vdots \\ a_K \end{bmatrix} + \begin{bmatrix} \epsilon_1 \\ \epsilon_2 \\ \vdots \\ \epsilon_N \end{bmatrix}$$

$$y = Xa + \epsilon$$

- $y \in \mathbb{R}^{N \times 1}$
- $X \in \mathbb{R}^{N \times K}$
- $a \in \mathbb{R}^{K \times 1}$
- $\epsilon \in \mathbb{R}^{N \times 1}$



MLR Calculations

- Solving for the coefficients is simple for us now:
- We can do it on the board!

$$\min_a \phi = \epsilon^T \epsilon$$

$$\min_a \phi = (\mathbf{y} - X\mathbf{a})^T (\mathbf{y} - X\mathbf{a})$$

$$\min_a \phi = \mathbf{y}^T \mathbf{y} - 2\mathbf{y}^T X\mathbf{a} - \mathbf{a}X^T X\mathbf{a}$$

$$\frac{\partial \phi}{\partial \mathbf{a}} = \mathbf{0} \quad \Rightarrow \quad \mathbf{a} = (X^T X)^{-1} X^T \mathbf{y}$$

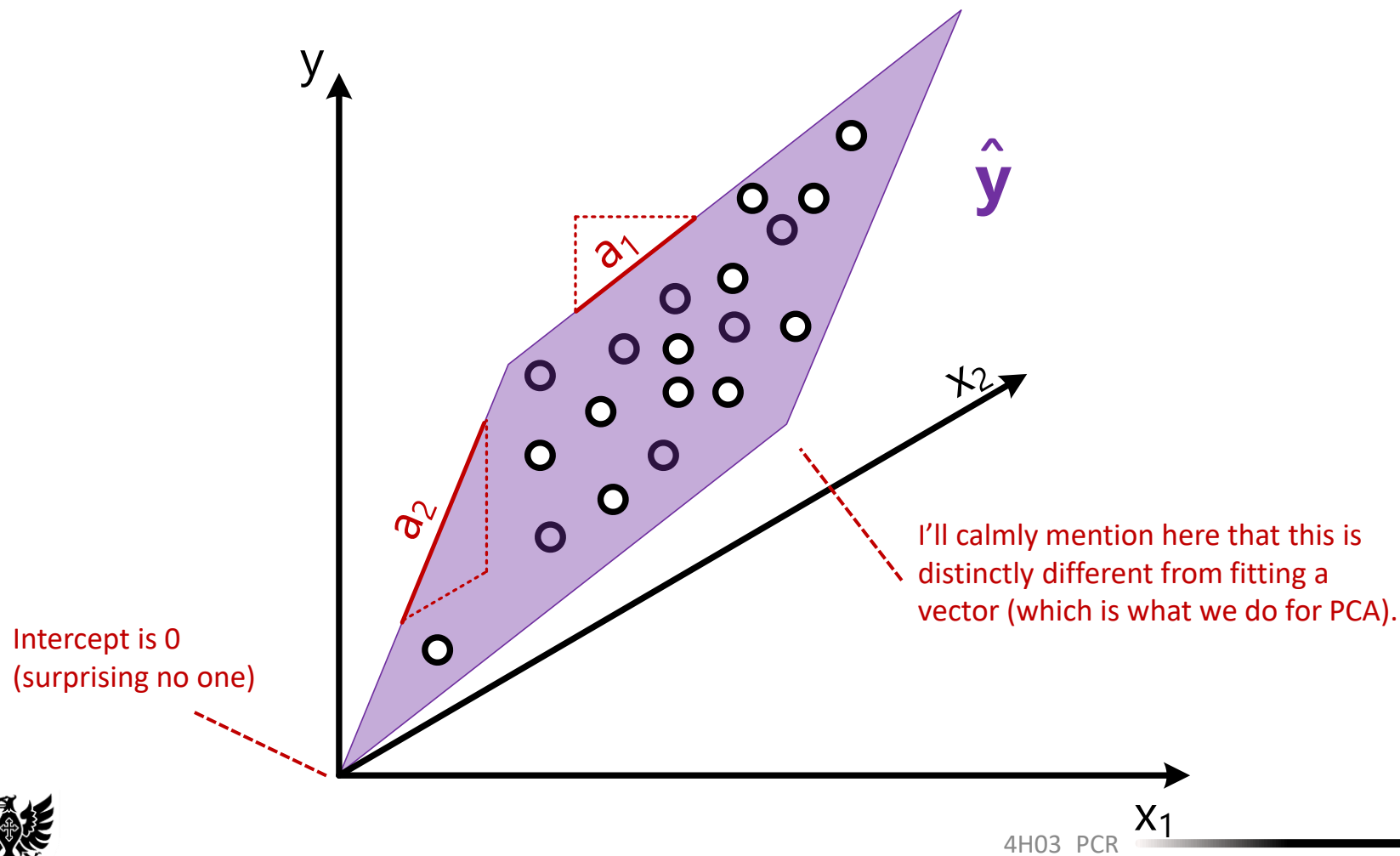
IMO, when you derive stuff like this from the optimization perspective, things make way more sense than when you are given random equations.

This ought to look VERY familiar with how we computed the loadings \mathbf{p} from the NIPALS algorithm :3



Interpretation of Coefficients

- Unsurprisingly, the coefficients represent the **projected slope** of variable k on the x_k, y plane



MLR: What Could Go Wrong?

- Missing values
 - What if one of our x values is missing?
 - Do we replace with 0? How about the mean (trick question :3)
 - This could be VERY dangerous
 - In reality, there is nothing we can do
 - **We can't use the model**
- $\hat{y}_i = a_1x_1 + a_2x_2 + \dots + a_Kx_K$



MLR: What Could Go Wrong?

- Correlated variables in \mathbf{X}
 - Leads to ALL KINDS of problems...

- Recall $\mathbf{a} = (X^T X)^{-1} X^T \mathbf{y}$

- Variance of coefficient vector $\mathcal{V}(\mathbf{a}) = (X^T X)^{-1} S_\epsilon^2$

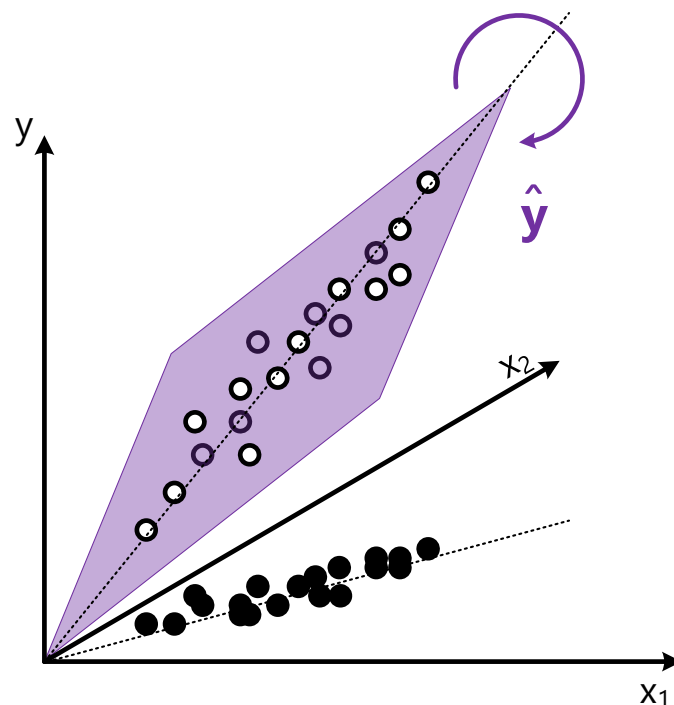
- S_ϵ^2 is the sample variance of ϵ with $N - K$ DOF: $S_\epsilon^2 = \frac{\epsilon^T \epsilon}{N - K}$

This is one of those things you can just take my word on ☺

- Leads to **lack of confidence** in coefficients of \mathbf{a}

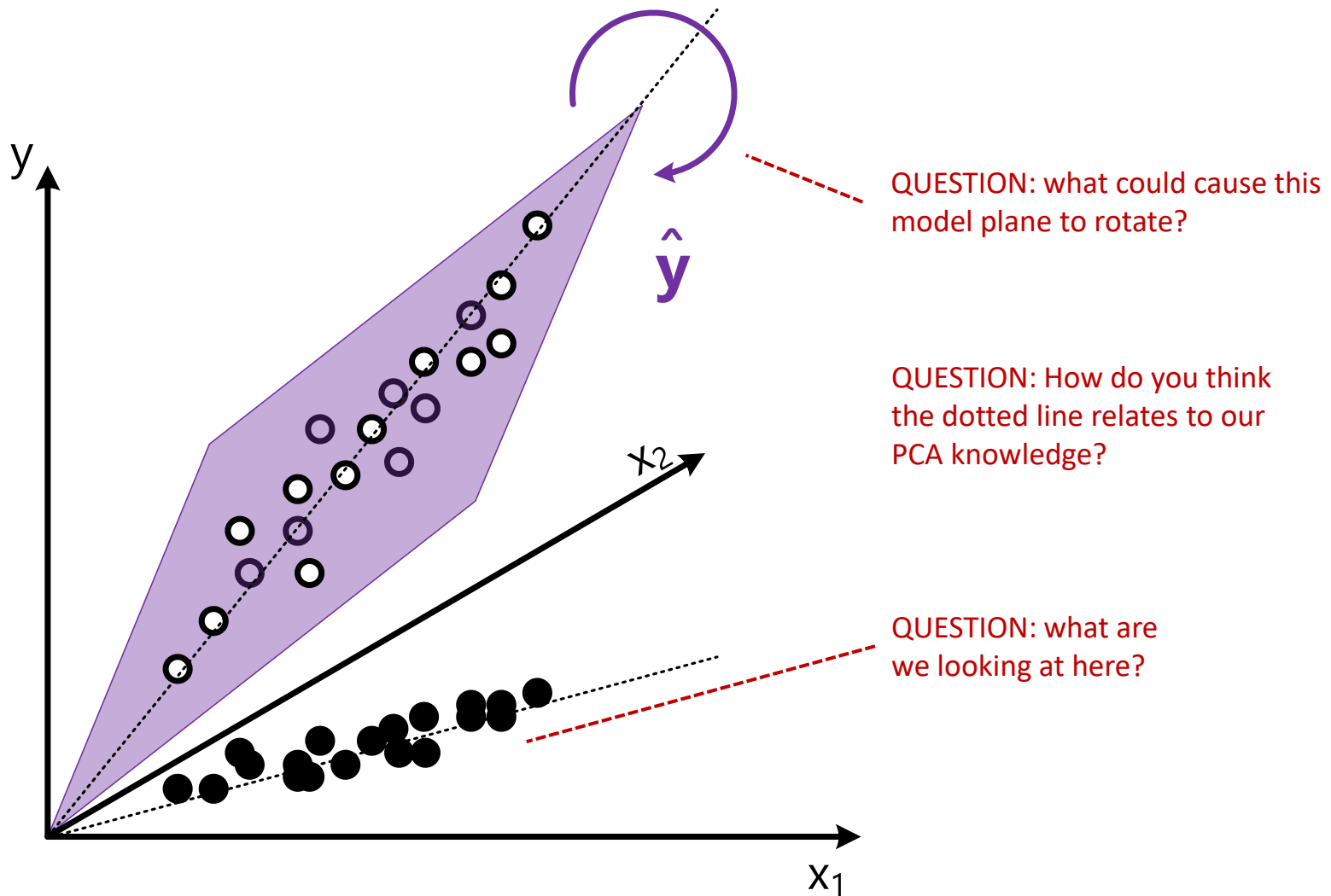
- Geometrically, it means the model plane could “spin” around and basically give the same results (see right)

- Algebraically, it means that our columns are **dependent** and thus $X^T X$ risks singularity (woof!)



MLR: What Could Go Wrong?

- Let's discuss this figure a little more



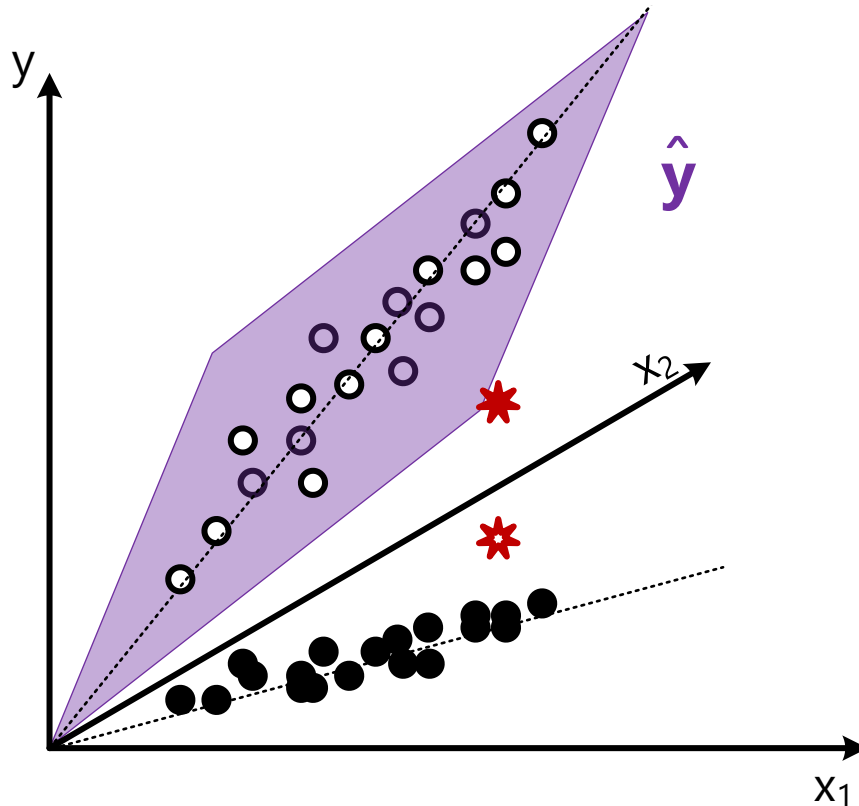
MLR: What Could Go Wrong?

- Noise in **X**
 - Recall the LS model is $\hat{y} = a_1x + \epsilon$
 - This assumes the error of the model is **contained in y**
- We therefore model our error in the **y**-space
 - We can compare our **y** error to S_ϵ^2
 - Akin to finding the R^2
 - See if something is wrong, and fix it
- CRITICALLY, LS assumes that **X** is exact (no noise)
 - No model space for the **X** error \Rightarrow loss of model fidelity



MLR: What Could Go Wrong?

- Nonsense observation for \mathbf{X}
 - Since we assume no error in \mathbf{X} , can result in:
 1. A “bad” outlier in \mathbf{y} that is not the model’s fault
 2. A prediction for \mathbf{y} that is completely unrelated to \mathbf{X}
 - In both cases, it is the measurement that is bad



MLR: What Could Go Wrong?

- MLR requires $N > K$
 - Some data sets have more columns than rows!
- Multiple **Y** variables
 - Need a separate regression model for each one
- Strategies to improve these issues?
 - Use basis functions if process is nonlinear
 - Visualization
 - Elimination of correlated columns in **X** (how, though?)
- For many data sets, problems solved using...



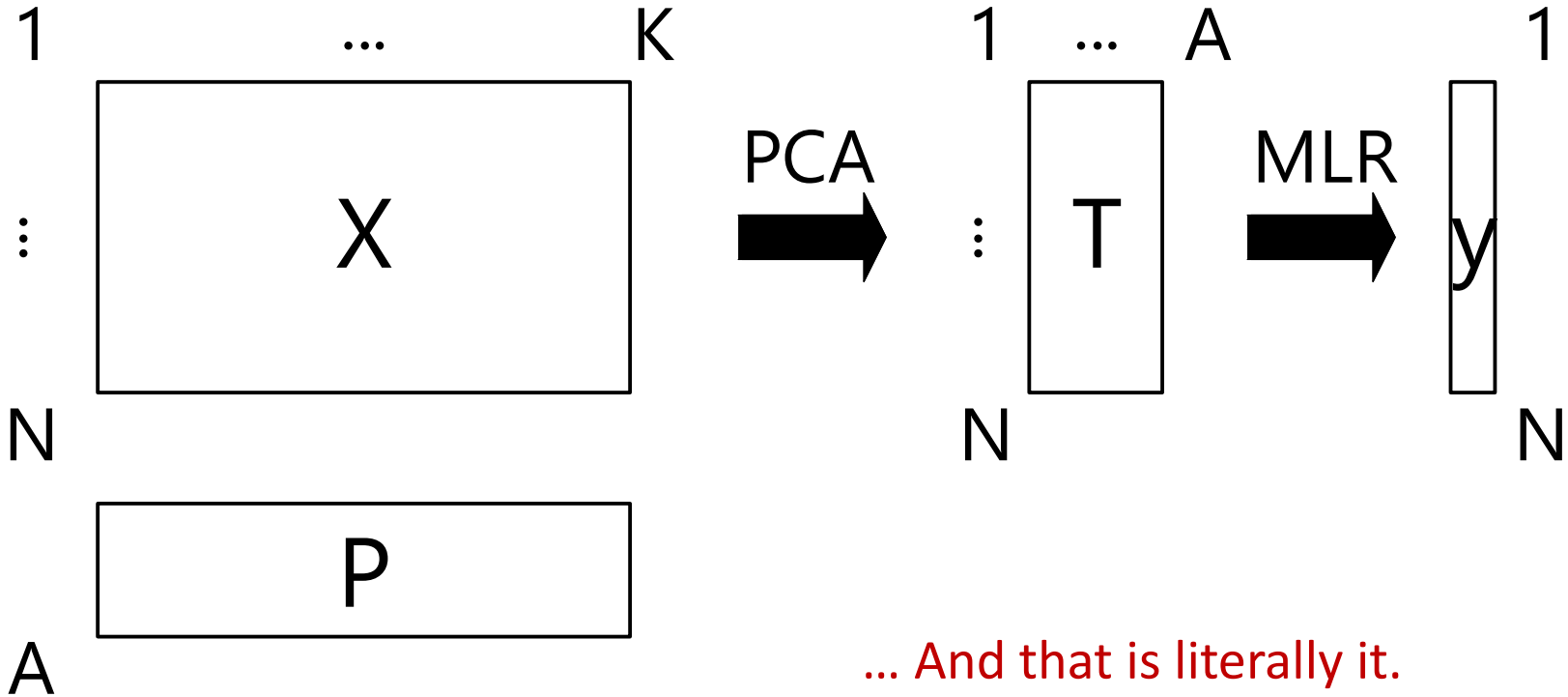
Principal Component Regression

And now for 4 slides on the actual topic :3



PCR in a Nutshell

- The idea (you saw this coming)
 1. Perform PCA on your data set \mathbf{X} to get loadings \mathbf{P}
 2. Regress \mathbf{Y} onto the SCORES \mathbf{T} to get regression coefficients



PCR in a Nutshell

- Calculations:

$$T = XP$$

$$\hat{y} = T\mathbf{a}$$

Don't forget to **center** (x AND y) and **scale** (really only needed for x)

- Coefficients in \mathbf{a} are computed as:

$$\mathbf{a} = (T^T T)^{-1} T^T \mathbf{y}$$

- Future predictions:

- Use x_{new} to get t_{new}
- Use t_{new} and \mathbf{a} to get \hat{y}_{new}

Don't forget to **center** and **scale** x_{new} here as well! Also, your \hat{y} will be centered so add \bar{y} to get the “real” y



PCR ADVANTAGES

- \mathbf{T} is orthogonal, thus $(\mathbf{T}^T\mathbf{T})^{-1}$ is easy to compute
 - Related point: correlated columns in \mathbf{X} will have one \mathbf{T}
- PCA step (using NIPALS) handles any missing data
- \mathbf{T} has considerably less error than \mathbf{X}
 - **Although** we do lose some variance in our PCA model...
 - **ALTHOUGH** we like to believe it is pure noise...
- PCA scores contain $A < K$ columns, thus the $N > K$ requirement for regression is **likely** to be met
 - QUESTION: What happens to the PCR output is $K = A$?
- CRITICAL: Each new point can be tested against:
 - SPE (is the \mathbf{X} data out of whack?)
 - T^2 (is this an extreme point on the PCA model?)



PCR DISADVANTAGES

- PCA components calculated without knowledge of \mathbf{Y}
 - PCA and regression steps are sequential, not simultaneous
 - The data in \mathbf{X} are not necessarily predictive of \mathbf{y}
- To ensure all predictors of \mathbf{Y} are contained in \mathbf{X} ...
 - We often fit extra components to PCA model
 - Violate our cross-validation heuristics
 - PRESS and Q^2 improvement rules, for example
- Would it not be better to fit ONE model that tried to relate the PCA components to \mathbf{Y} immediately?
 - It sure would



Final Remarks

- A quick review of what we have done:
 - Review of covariance and correlation
 - Review of linear regression with no intercept
 - Covered methods of computing coefficients to MLR
 - Caveats and shortcomings of MLR
 - PCR \rightarrow a better method of dealing with high correlation in \mathbf{X}
- Next up: **Projection of Latent Structures**
 - Idea: fit loadings for \mathbf{X} and \mathbf{Y} together
 - Maximize **covariance** between \mathbf{X} and \mathbf{Y} columns

