

Chemical Engineering 4H03

Projection of Latent Structures (NIPALS)

Jake Nease
McMaster University

Portions of this work are copyright of ConnectMV

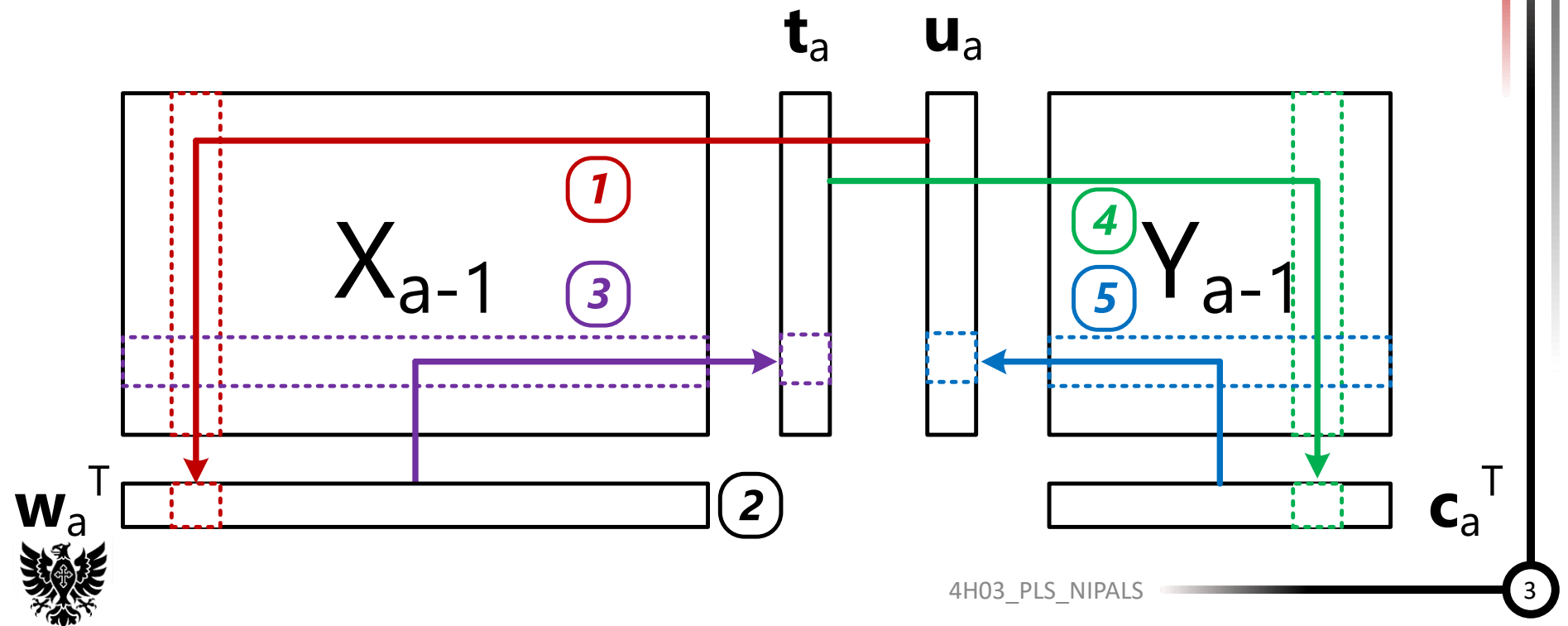
Where are We?

- We have derived PLS from “first principles”
 - Well, not really, but we know the **objective**
- Unfortunately, without optimization we have no way of reliably computing \mathbf{c}_a and \mathbf{w}_a to $\max \psi = \mathcal{C}(\mathbf{t}_a, \mathbf{u}_a)$
 - I’ll gently remind you here that $\max \psi$ is equivalent to finding the “pareto-optimal” of our three competing objectives
- This leads us to adapting our NIPALS algorithm for PLS



NIPALS: The Idea

- Imagine creating PCA models for the X and Y spaces separately using NIPALS
 - What would that look like?
 - What is the main thing missing?
- PLS NIPALS looks very similar to PCA:



NIPALS Algorithm for PLS

Naturally, **I**'d **P**refer **A**nother **L**east-**S**quares

PLS NIPALS Algorithm

- NIPALS begins with X_0 as before AND Y_0
 - Recall X_a and Y_a are the data sets explained by a components

FOR $a = 1, 2, \dots, A$:

1. Select an (arbitrary) column as \mathbf{u}_a
2. In a loop until convergence:
 - I. Regress columns from X_{a-1} onto \mathbf{u}_a to get weights \mathbf{w}_a
 - II. Normalize \mathbf{w}_a to unit length
 - III. Regress rows from X_{a-1} onto \mathbf{w}_a to get \mathbf{t}_a
 - IV. Regress columns from Y_{a-1} onto \mathbf{t}_a to get loadings \mathbf{c}_a
 - V. Regress rows from Y_{a-1} onto \mathbf{c}_a to get \mathbf{u}_a
3. Deflate component from X_{a-1} and Y_{a-1} to get X_a and Y_a

END



NIPALS for PLS Step By Step

- **STEP 1** Select an arbitrary column for \mathbf{u}_a
 - Any individual column in Y_0
 - A column of normally distributed random numbers
 - Basically anything except the zero column (WHY??)



NIPALS for PLS Step By Step

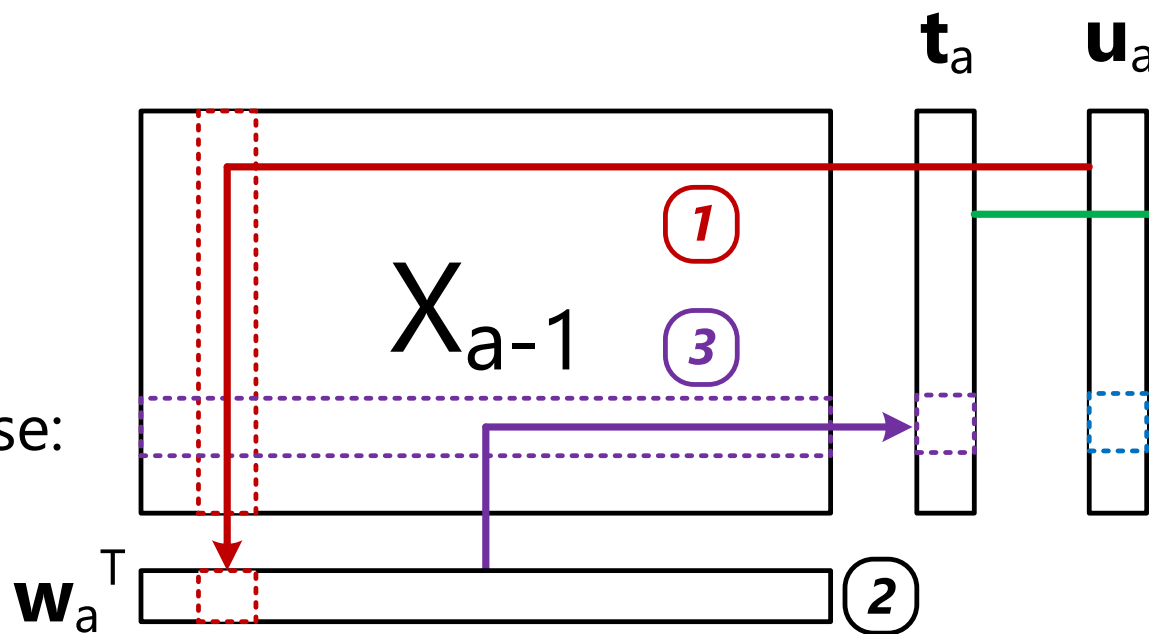
- **STEP 2.1** Regress each column of X_{a-1} (\mathbf{x}_k) onto \mathbf{u}_a
 - Perform a LS regression of \mathbf{x}_k onto \mathbf{u}_a (*regress y onto x*)
 - Store the regression coefficient as $w_{k,a}$

- Recall that a linear equation with zero intercept is $\hat{y} = a_1 x$

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

- Therefore in this case:

$$w_{k,a} = \frac{\mathbf{u}_a^T \mathbf{x}_k}{\mathbf{u}_a^T \mathbf{u}_a}$$



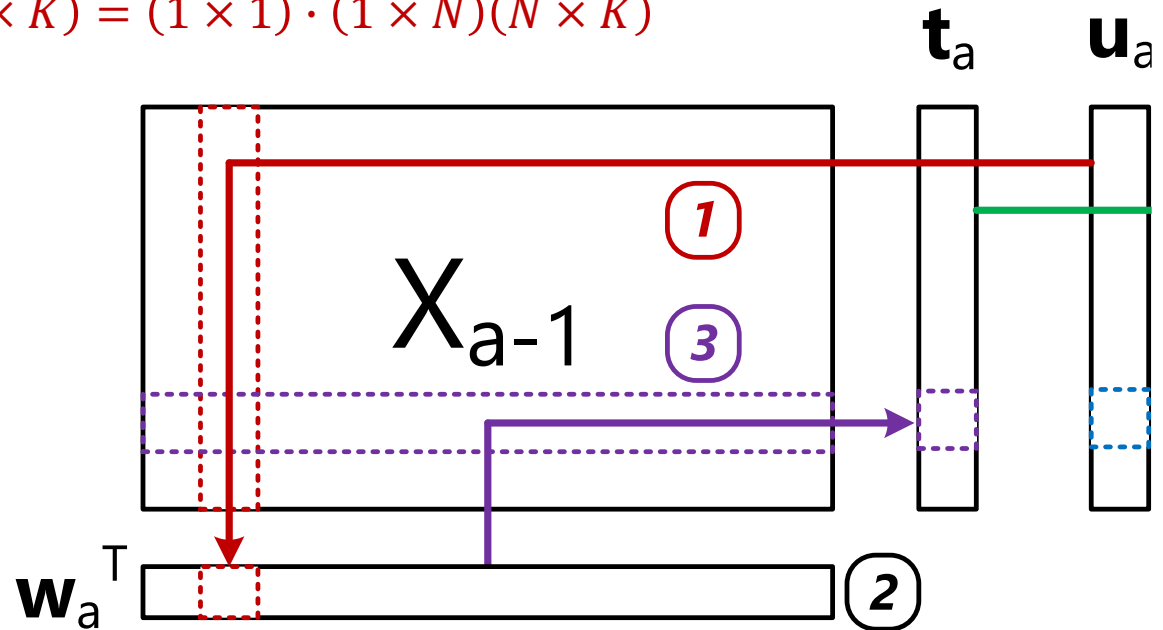
NIPALS for PLS Step By Step

- **STEP 2.1** Repeat regression for all columns in X_{a-1}
- Can compute all regressions in one go (HOW!?)

$$\mathbf{w}_a^T = \frac{1}{\mathbf{u}_a^T \mathbf{u}_a} \cdot \mathbf{u}_a^T X_{a-1}$$

$$(1 \times K) = (1 \times 1) \cdot (1 \times N)(N \times K)$$

- $\mathbf{u}_a \in \mathbb{R}^{N \times 1}$
- $X_{a-1} \in \mathbb{R}^{N \times K}$
- $\mathbf{w}_a \in \mathbb{R}^{K \times 1}$



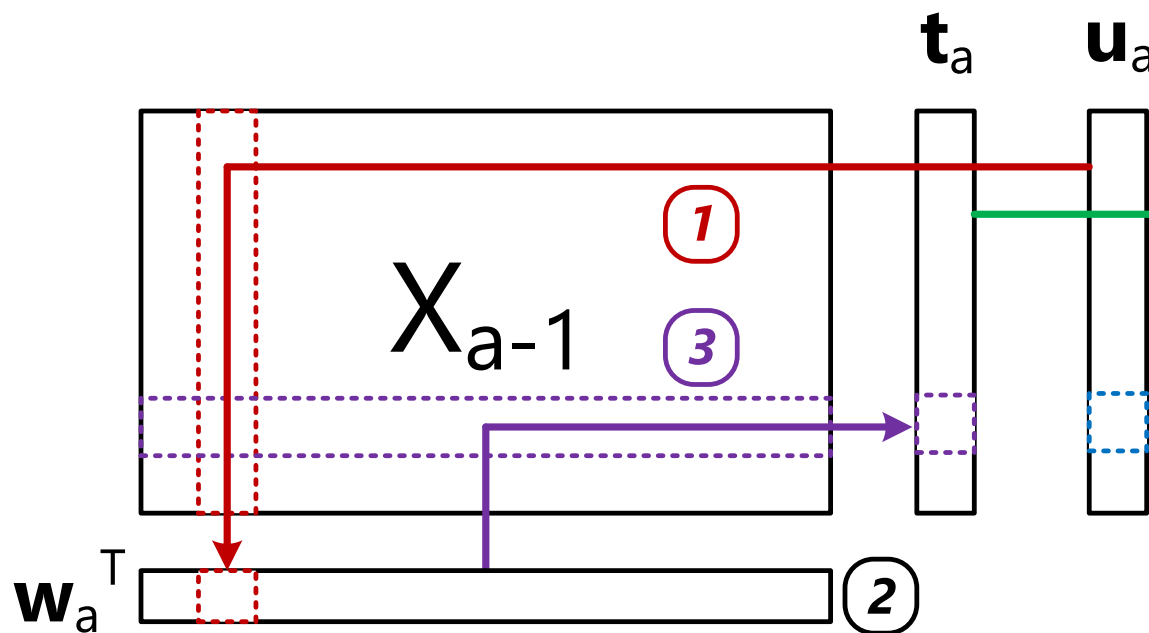
NIPALS for PLS Step By Step

- STEP 2.2 Normalize the weights

- w_a will not have unit length
- Rescale to magnitude of 1.0:

$$w_a = \frac{w_a}{\|w_a\|}$$

----- I'll gently point out the nomenclature difference here



NIPALS for PLS Step By Step

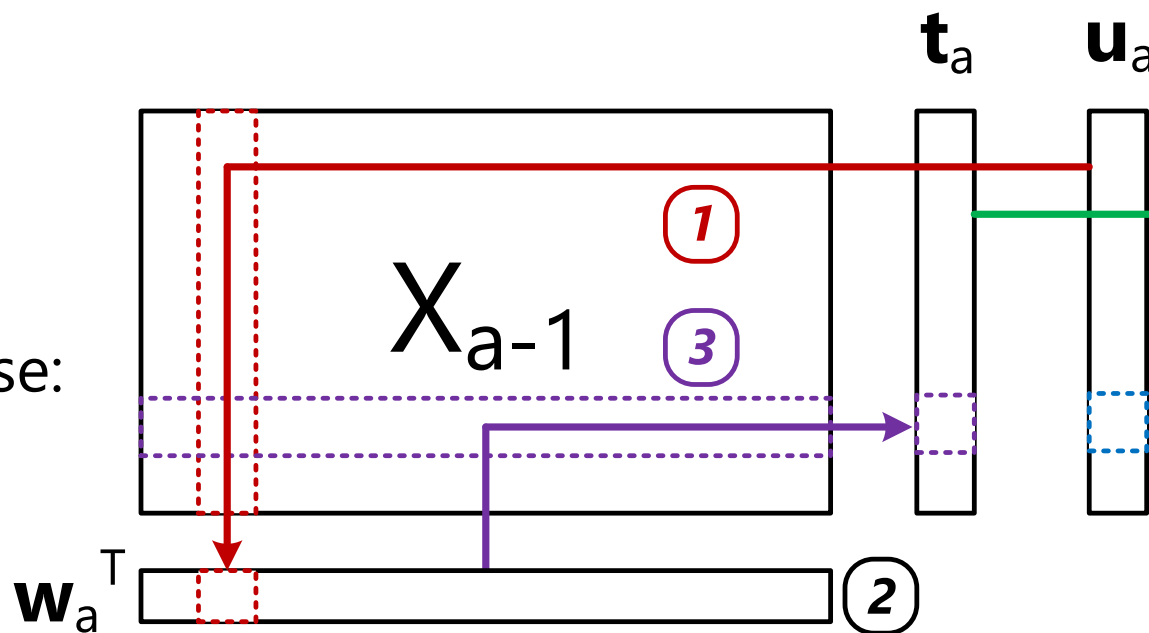
- **STEP 2.3** Regress each row of X_{a-1} (\mathbf{x}_n) onto \mathbf{w}_a^T
 - Perform a LS regression of \mathbf{x}_n onto \mathbf{w}_a^T (*regress y onto x*)
 - Store the regression coefficient as $t_{n,a}$

- Recall that a linear equation with zero intercept is $\hat{y} = a_1 x$

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

- Therefore in this case:

$$t_{n,a} = \frac{\mathbf{w}_a^T \mathbf{x}_n}{\mathbf{w}_a^T \mathbf{w}_a}$$



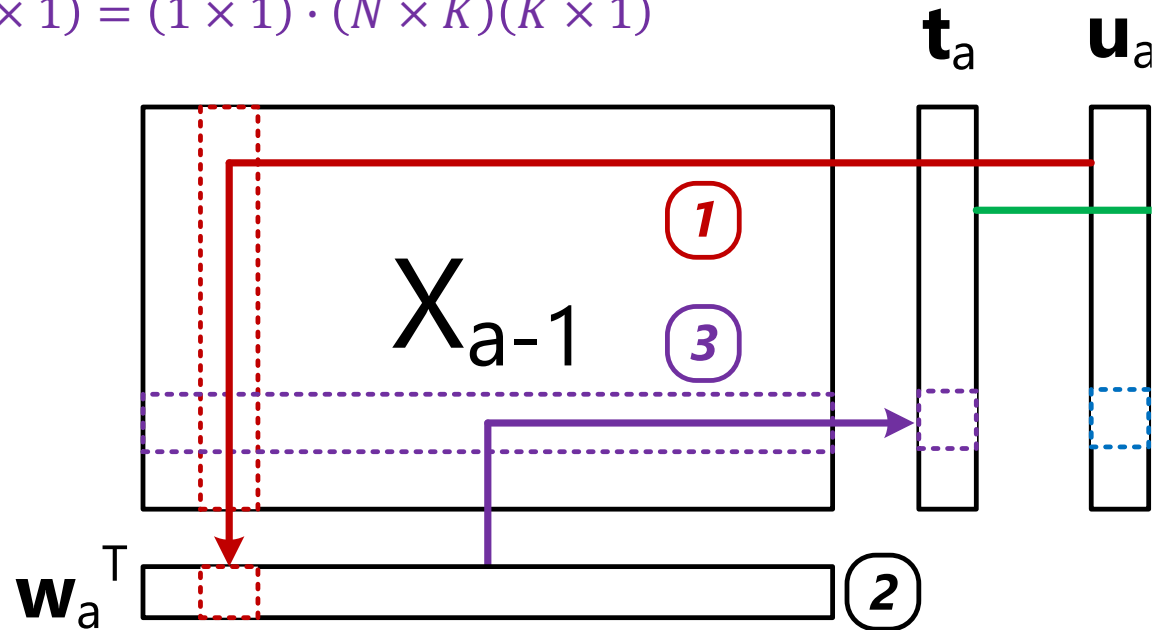
NIPALS for PLS Step By Step

- **STEP 2.3** Repeat regression for all rows in X_{a-1}
- Can compute all regressions in one go

$$\mathbf{t}_a = \frac{1}{\mathbf{w}_a^T \mathbf{w}_a} \cdot X_{a-1} \mathbf{w}_a$$

$$(N \times 1) = (1 \times 1) \cdot (N \times K)(K \times 1)$$

- $\mathbf{t}_a \in \mathbb{R}^{N \times 1}$
- $X_{a-1} \in \mathbb{R}^{N \times K}$
- $\mathbf{w}_a \in \mathbb{R}^{K \times 1}$



NIPALS for PLS Step By Step

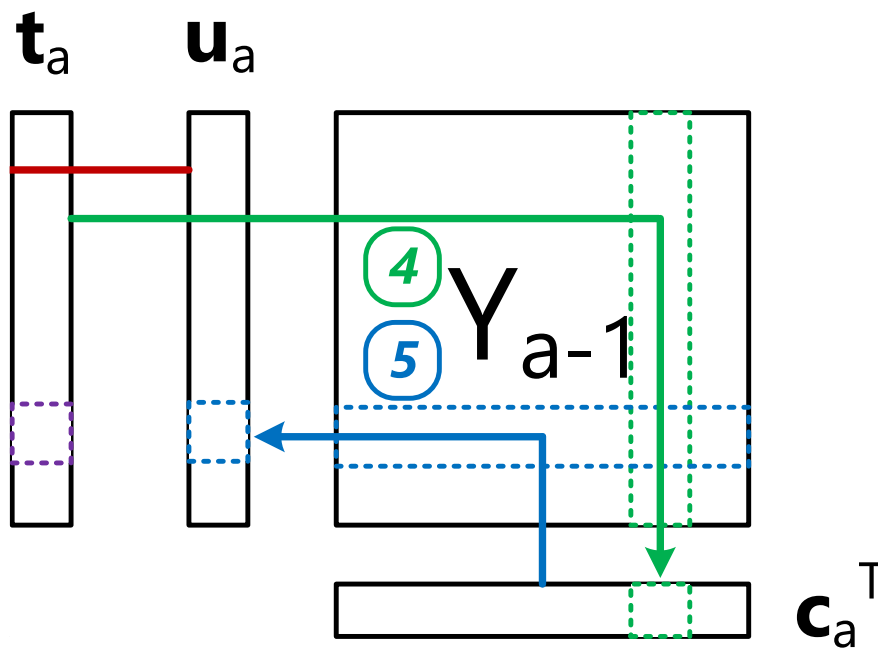
- **STEP 2.4** Regress each column of Y_{a-1} (\mathbf{y}_m) onto \mathbf{t}_a
 - Perform a LS regression of \mathbf{y}_m onto \mathbf{t}_a^T (*regress y onto x*)
 - Store the regression coefficient as $c_{m,a}$

- Recall that a linear equation with zero intercept is $\hat{y} = a_1 x$

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

- Therefore in this case:

$$c_{m,a} = \frac{\mathbf{t}_a^T \mathbf{y}_m}{\mathbf{t}_a^T \mathbf{t}_a}$$



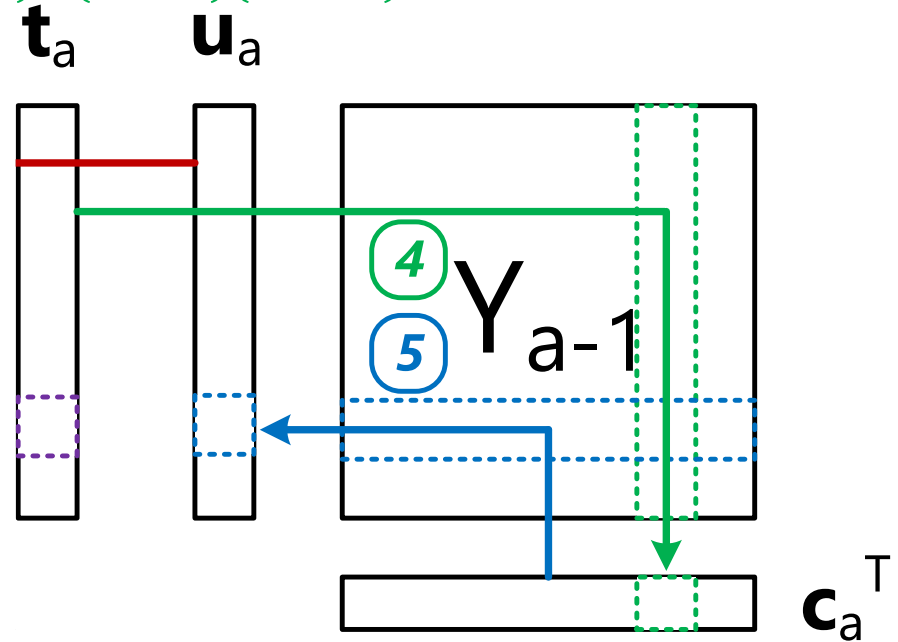
NIPALS for PLS Step By Step

- **STEP 2.4** Repeat regression for all columns in Y_{a-1}
- Can compute all regressions in one go

$$\mathbf{c}_a^T = \frac{1}{\mathbf{t}_a^T \mathbf{t}_a} \cdot \mathbf{t}_a^T Y_{a-1}$$

$$(1 \times M) = (1 \times 1) \cdot (1 \times N)(N \times M)$$

- $\mathbf{t}_a \in \mathbb{R}^{N \times 1}$
- $Y_{a-1} \in \mathbb{R}^{N \times M}$
- $\mathbf{c}_a \in \mathbb{R}^{M \times 1}$



NIPALS for PLS Step By Step

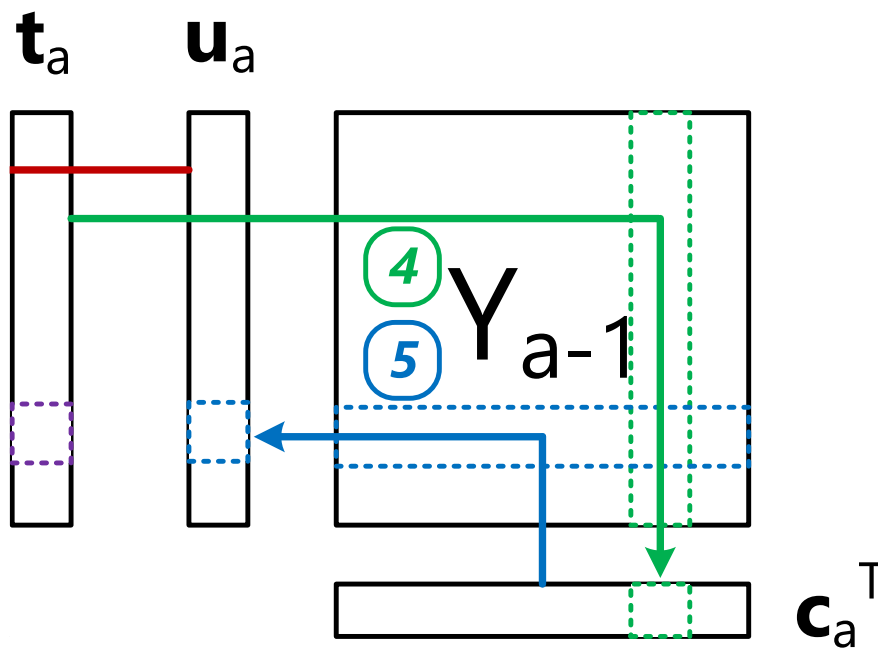
- **STEP 2.5** Regress each row of Y_{a-1} (\mathbf{y}_n) onto \mathbf{c}_a
 - Perform a LS regression of \mathbf{y}_n onto \mathbf{c}_a (*regress y onto x*)
 - Store the regression coefficient as $u_{n,a}$

- Recall that a linear equation with zero intercept is $\hat{y} = a_1 x$

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

- Therefore in this case:

$$u_{n,a} = \frac{\mathbf{c}_a^T \mathbf{y}_n}{\mathbf{c}_a^T \mathbf{c}_a}$$



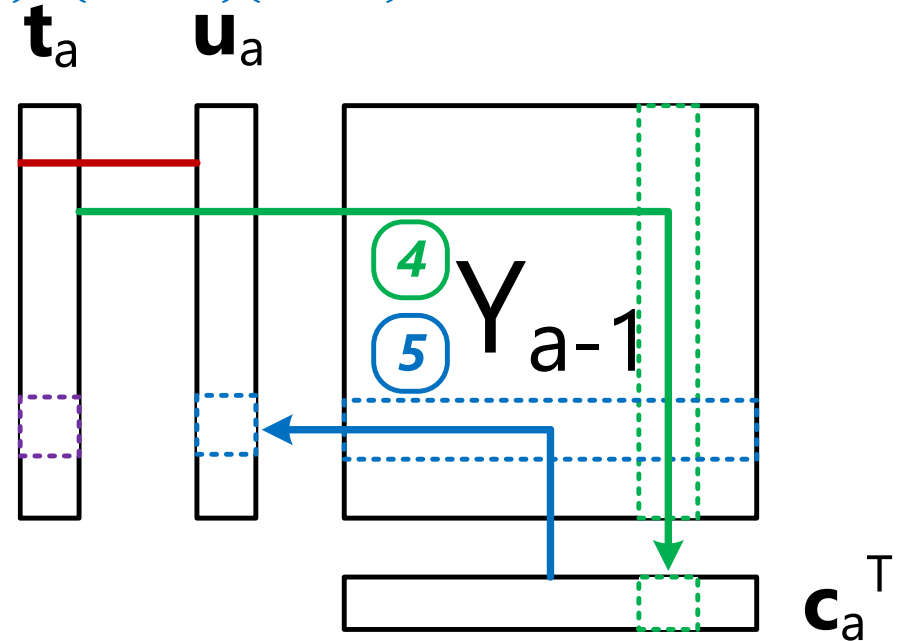
NIPALS for PLS Step By Step

- **STEP 2.5** Repeat regression for all rows in Y_{a-1}
- Can compute all regressions in one go

$$\mathbf{u}_a = \frac{1}{\mathbf{c}_a^T \mathbf{c}_a} \cdot Y_{a-1} \mathbf{c}_a$$

$$(N \times 1) = (1 \times 1) \cdot (N \times M)(M \times 1)$$

- $\mathbf{u}_a \in \mathbb{R}^{N \times 1}$
- $Y_{a-1} \in \mathbb{R}^{N \times M}$
- $\mathbf{c}_a \in \mathbb{R}^{M \times 1}$



NIPALS Step By Step

- CHECK FOR CONVERGENCE

- Compare \mathbf{u}_a to \mathbf{u}_a from previous iteration
- Stop if $\Delta \mathbf{u}_a \leq \sqrt{\epsilon} \approx 10^{-8}$
- Change could mean $\|\mathbf{u}_a^i - \mathbf{u}_a^{i-1}\| \leq \sqrt{\epsilon}$ [absolute]
- Change could mean $\frac{\|\mathbf{u}_a^i - \mathbf{u}_a^{i-1}\|}{\|\mathbf{u}_a^{i-1}\|} \leq \sqrt{\epsilon}$ [relative]
- ALSO probably want an iterations limit in case (300 is good)

- AT CONVERGENCE

- $\mathbf{t}_a, \mathbf{w}_a, \mathbf{u}_a$ and \mathbf{c}_a are the a^{th} component
- Store as columns in T, W, U , and C , respectively



OBSERVATION

- Something **sneaky** has happened here
 - Since we never regressed t_a onto X_{a-1} , we actually never computed the **loadings** of X_{a-1} as related to p_a
- WORKSHOP: How then do we find the **residuals** E_a of X_{a-1} ? It is NOT fair to use w_a (our weights) [WHY?]



NIPALS Step By Step

- **Step 3** Deflation
- 3.1. Compute loadings $\mathbf{p}_a = \frac{1}{\mathbf{t}_a^T \mathbf{t}_a} \mathbf{X}_{a-1}^T \mathbf{t}_a$
 - Regressing columns of \mathbf{X}_{a-1} on converged scores \mathbf{t}_a
- 3.2. Remove predicted variance in \mathbf{X}_{a-1} and \mathbf{Y}_{a-1}
 - Deflate \mathbf{X}_{a-1} :
 - $\mathbf{E}_a = \mathbf{X}_{a-1} - \hat{\mathbf{X}}_{a-1}$ $\mathbf{E}_a = \mathbf{X}_{a-1} - \mathbf{t}_a \mathbf{p}_a^T$
 - Set new deflated \mathbf{X} $\mathbf{X}_a = \mathbf{E}_a$
 - Deflate \mathbf{Y}_{a-1} :
 - $\mathbf{F}_a = \mathbf{Y}_{a-1} - \hat{\mathbf{Y}}_{a-1}$ $\mathbf{F}_a = \mathbf{Y}_{a-1} - \mathbf{t}_a \mathbf{c}_a^T$
 - Set new deflated \mathbf{Y} $\mathbf{Y}_a = \mathbf{F}_a$

Yes it actually is \mathbf{t}_a , that is not a typo



Interpreting the Weights

- Like in PCA, NIPALS computes scores based on deflated matrices
- However, UNLIKE in PCA, those scores are (likely) not the scores that explain the greatest variance
 - Recall we are trading-off the explanations of X , Y and $r\{X, Y\}$
- Therefore the weights \mathbf{w}_a we calculate can actually only be applied to the **deflated matrix** X_{a-1} used during the NIPALS algorithm
 - We could do this in PCA as well, but it would not change anything



Interpreting the Weights

- So, ANY scores t_a must be calculated using X_{a-1}
 - $t_1 = X_0 \mathbf{w}_1$ [Nothing spiced here]
 - $t_2 = X_1 \mathbf{w}_2$ [Note we have to use X_1 here]
- The weights \mathbf{w}_2 really only have meaning for the DELFATED matrix X_1 . This is annoying for two reasons:
 1. It is hard to interpret because \mathbf{w}_2 and \mathbf{w}_1 can be directly compared like \mathbf{p}_2 and \mathbf{p}_1 could, for example
 2. To calculate ALL scores, we are NOT allowed to use $T = XW$ since X means something “different” for each column of W
- We would likely prefer a variant of W that has the **same interpretive meaning** as P from PCA



Interpreting the Weights: W^*

- Consider some way to get t_a using X_0 only. We call it w^*
 - $t_1 = X_0 w_1^*$ [Same as before]
 - $t_2 = X_0 w_2^*$ [Note we have to use X_0 here]
 - $t_A = X_0 w_A^*$ [For all components]
- We can actually obtain this collection of w_a^* values by exploiting the loadings we calculated:

$$W^* = W(P^T W)^{-1}$$

$$(N \times A) = (N \times A) \cdot ([A \times K][K \times A])$$

- So we get $T = XW^*$ in one go (note X is just X_0 for us)
 - This is a cleaner representation of the relationships in X
 - Note that $w_1^* = w_1$ BUT $w_a^* \neq w_a \forall a > 1$



DETOUR: Why do we use $\hat{X} = t p^T$?

- You might be a little freaked out that we want to calculate \hat{X} based on the scores in t only
 - Instead, you might realize that we regressed the columns of X_{a-1} onto w_a to get u_a
 - It therefore might make sense to estimate $\hat{X}_{a-1} = u_a w_a^T$
- This leads to a **CRITICAL** point about PLS: Our objective is to train the PLS model so we can *predict* Y in the future... Which means... **We don't have Y**
 - So, we **can't** predict the Y -space scores to get $U = YC$
 - If we don't have U , we **can't** get $\hat{X} = UW^T$
 - INSTEAD, we would like to have some way to estimate \hat{X} using T only, since with a new data point we can get $T_{new} = X_{new}W^*$
 - This is **important** because we need \hat{X} to estimate SPE
 - Result: We actually don't end up using U ever again...

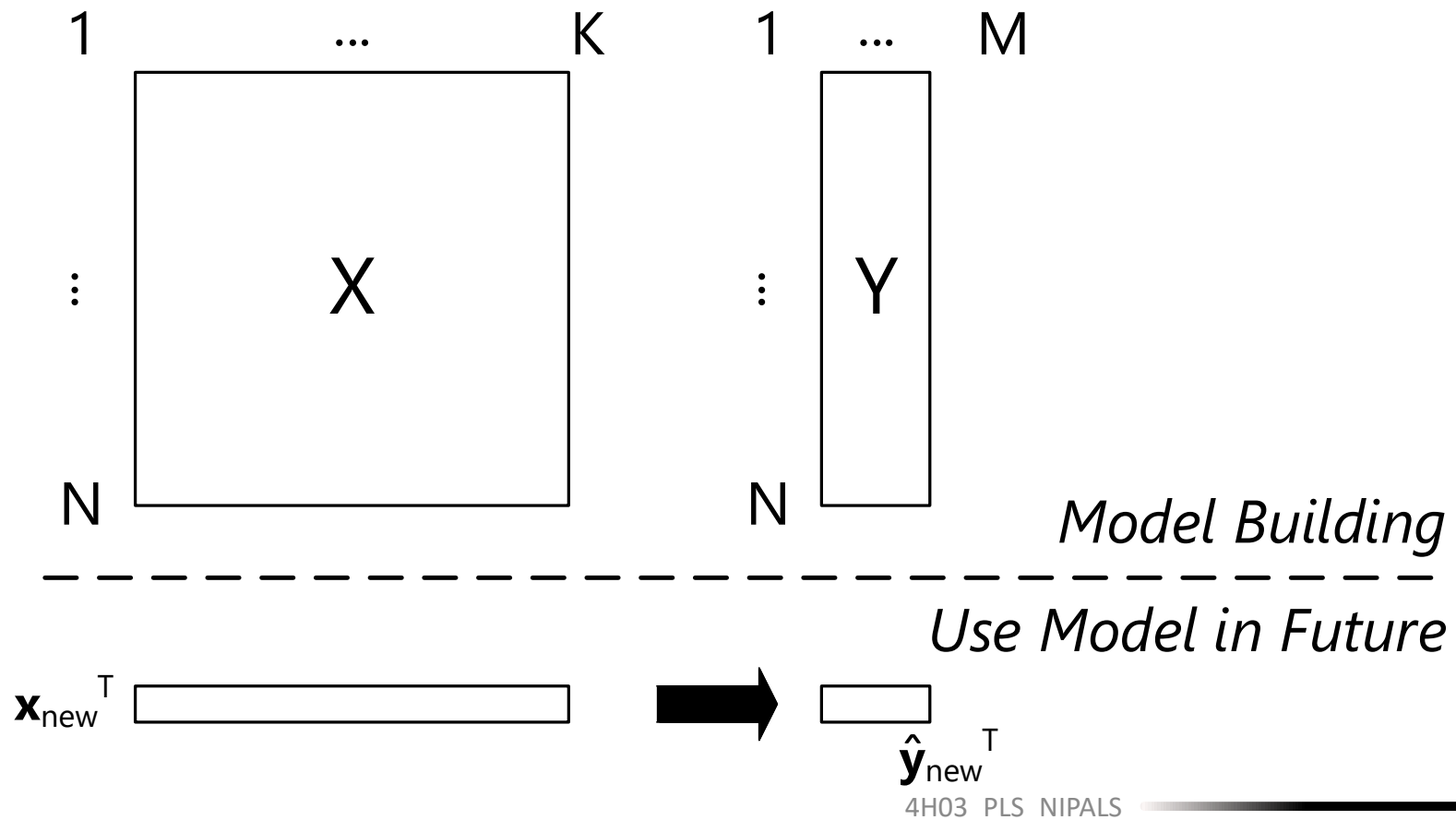


Using PLS on New Data

Y predict Y if you are not going to predict Y?

Using PLS on New Data

- Our objective here is to train the model using KNOWN values of Y so we can predict those outcomes later
 - To do this, we sacrifice explanation of X more than PCA



Using PLS on New Data

FIRST: Get Scores

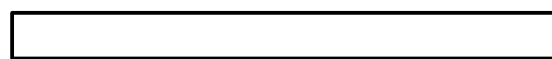
$$t_{1,new} = \mathbf{x}_{new}^T \mathbf{w}_1$$

$$\mathbf{x}_{new}^T = \mathbf{x}_{new}^T - t_{1,new} \mathbf{p}_1^T$$

$$t_{2,new} = \mathbf{x}_{new}^T \mathbf{w}_2$$

$$\mathbf{x}_{new}^T = \mathbf{x}_{new}^T - t_{2,new} \mathbf{p}_2^T$$

...



\mathbf{x}_{new}^T

OH! And don't forget to
center and scale \mathbf{x}_{new}^T ...
MAN would you look silly :P



\mathbf{W}^*



\mathbf{t}_{new}^T

You get a row, but by
convention we'll call it
 t^T since any single
vector is a column for us

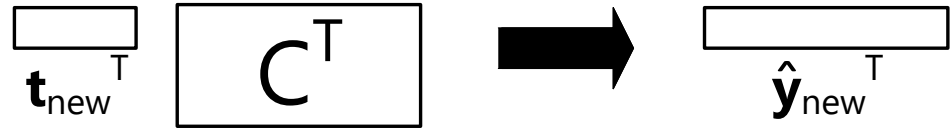
OR USE \mathbf{W}^* to get whole vector \mathbf{t}_{new} without deflation:

$$\mathbf{t}_{new}^T = \mathbf{x}_{new}^T \mathbf{W}^*$$



Using PLS on New Data

NEXT: Get Y



$$\hat{\mathbf{y}}_{new} = \mathbf{t}_{new}^T \mathbf{C}^T$$

ALSO remember to UNCENTER
and UNSCALE your \mathbf{y}_{new} lest
you look silly ☺

- Note that you **CAN** skip the middle step and go:

$$\hat{\mathbf{y}}_{new} = \mathbf{x}_{new}^T \mathbf{W}^* \mathbf{C}^T$$

- You should NEVER implement PLS this way
 - We WANT the scores to calculate our model stats (T^2 and SPE)



Cross-Validation

- One big advantage of predicting \hat{y} is that we can use it to determine the optimal number of components!
 - You may recall we did this in PCA using Q^2 and PRESS
 - The procedure is the same, but we will do it for Y instead!

- Recall for training data:
$$R^2 = 1 - \frac{\mathcal{V}(F_A)}{\mathcal{V}(Y)}$$
- DEFINE for testing data:
$$Q^2 = 1 - \frac{\mathcal{V}(F_A \text{ predicted})}{\mathcal{V}(Y)}$$

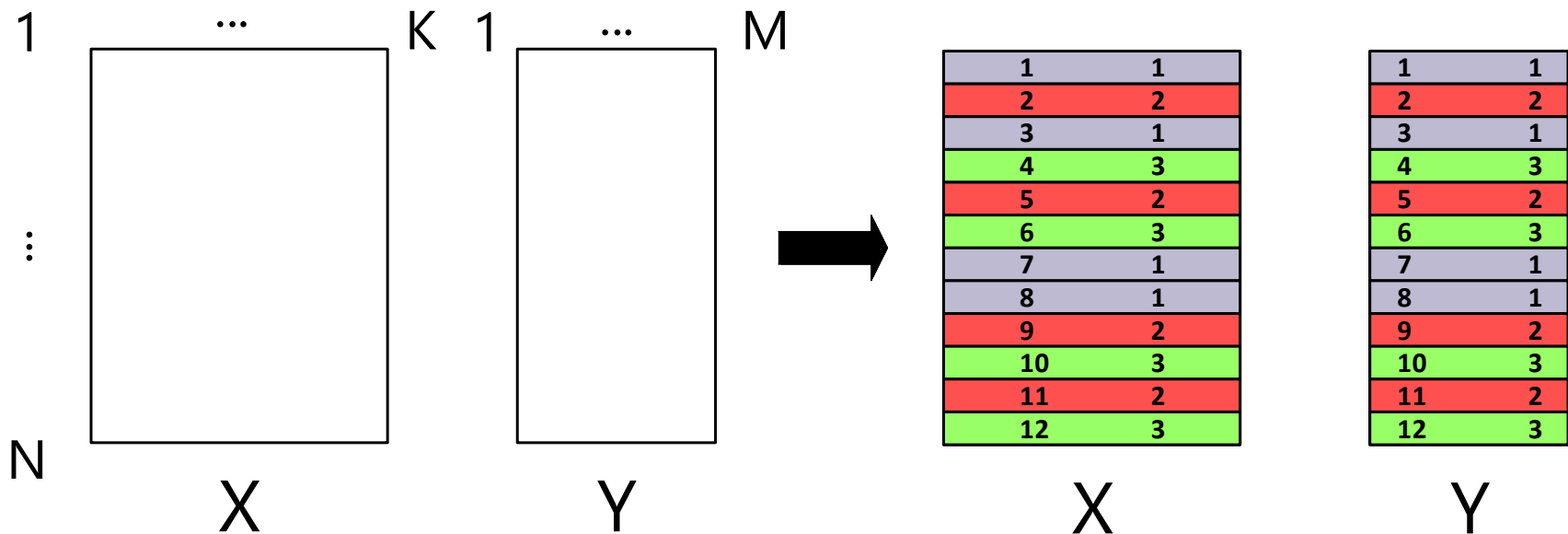
- $\mathcal{V}(E_A \text{ predicted})$ is known as the **PRESS**

Use the Y residuals and variance here



Cross-Validation

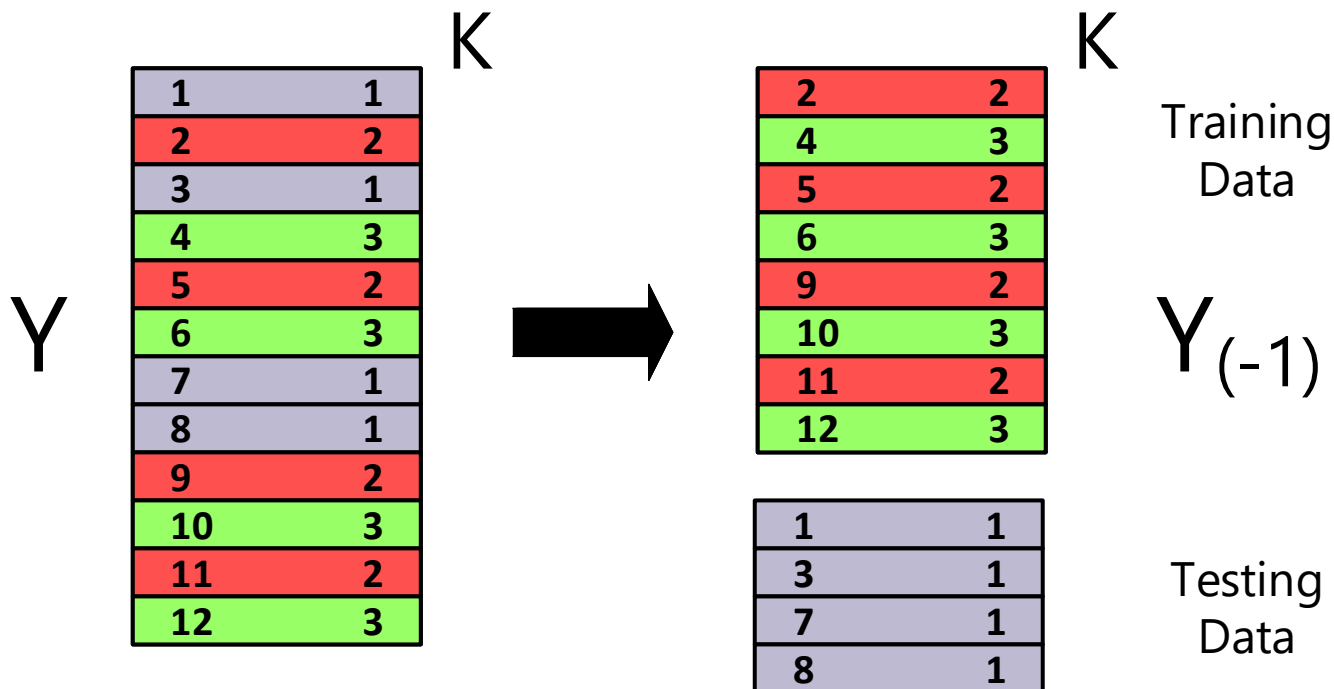
- Step 1: Split rows of X AND Y into \mathbf{G} groups (3 here)
 - Typically \mathbf{G} is 7 in software packages
 - Can be random or ordered (random in this example)
 - Note that this can depend on time relevance!



Cross-Validation

- Step 2.1: Fit PLS model
 - Use $\mathbf{X}_{(-1)}$ and $\mathbf{Y}_{(-1)}$ for fitting
 - Use $\mathbf{X}_{(1)}$ and $\mathbf{Y}_{(1)}$ for testing
- Compute $F_{(1)} = Y_{(1)} - \hat{Y}_{(1)}$

Important to remember that the main point of PLS is to predict \mathbf{Y} , so we are using the \mathbf{Y} residuals here



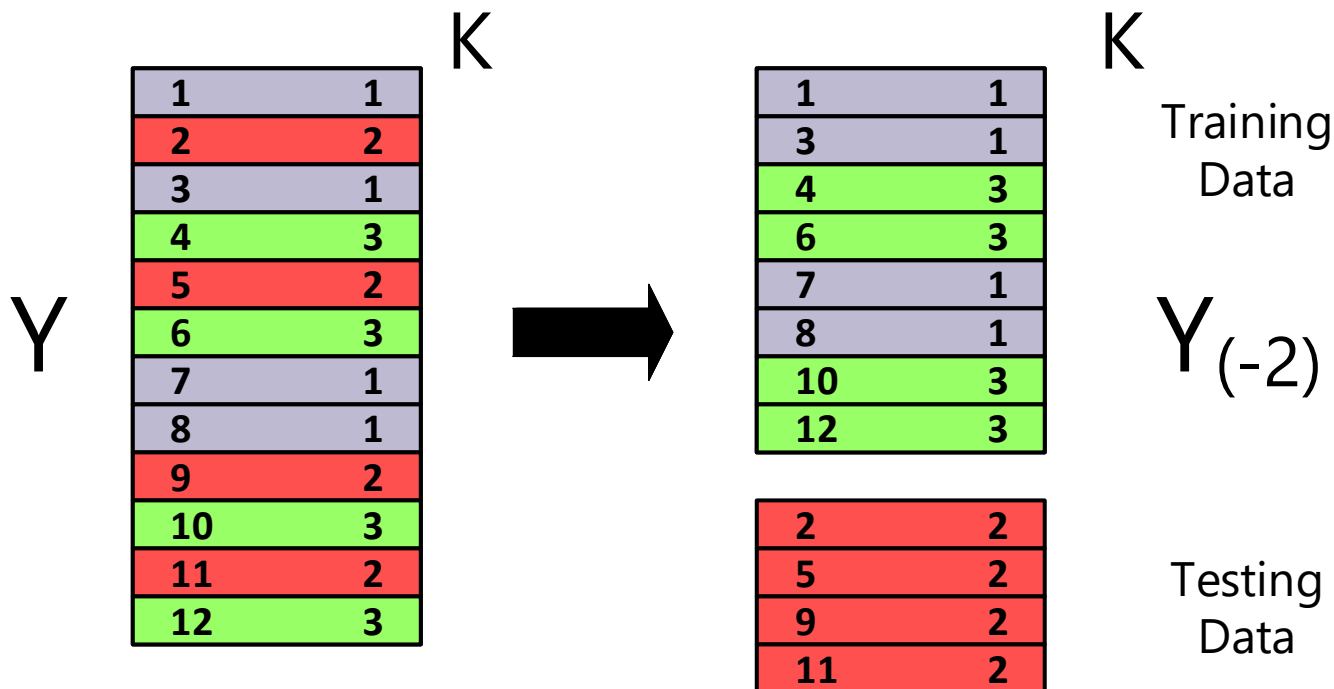
We also need \mathbf{X} here (not shown)



Cross-Validation

- Step 2.2: Fit PLS model
 - Use $\mathbf{X}_{(-2)}$ and $\mathbf{Y}_{(-2)}$ for fitting
 - Use $\mathbf{X}_{(2)}$ and $\mathbf{Y}_{(2)}$ for testing
- Compute $F_{(2)} = Y_{(2)} - \hat{Y}_{(2)}$

Important to remember that the main point of PLS is to predict \mathbf{Y} , so we are using the \mathbf{Y} residuals here



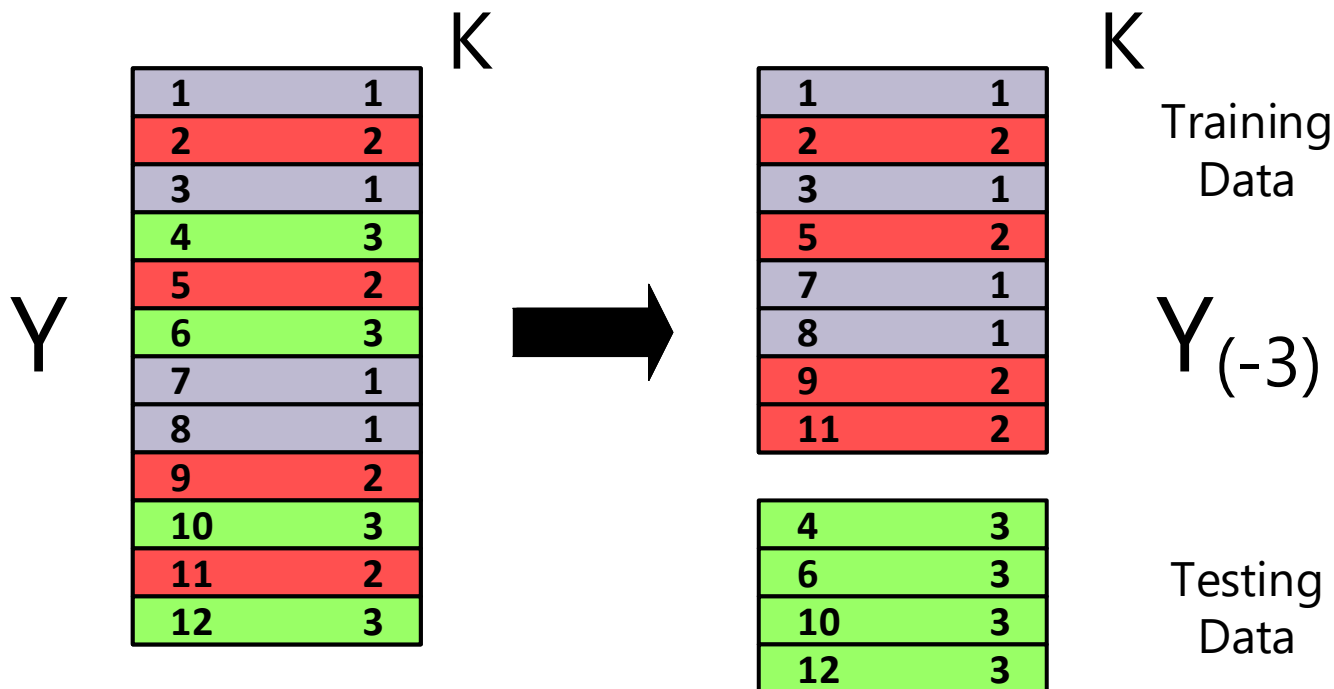
We also need \mathbf{X} here (not shown)



Cross-Validation

- Step 2.3: Fit PLS model
 - Use $\mathbf{X}_{(-3)}$ and $\mathbf{Y}_{(-3)}$ for fitting
 - Use $\mathbf{X}_{(3)}$ and $\mathbf{Y}_{(3)}$ for testing
- Compute $F_{(3)} = Y_{(3)} - \hat{Y}_{(3)}$

Important to remember that the main point of PLS is to predict \mathbf{Y} , so we are using the \mathbf{Y} residuals here



We also need \mathbf{X} here (not shown)



Cross-Validation

- Step 2.G: Fit PLS model
 - Use $\mathbf{X}_{(-G)}$ and $\mathbf{Y}_{(-G)}$ for fitting
 - Use $\mathbf{X}_{(G)}$ and $\mathbf{Y}_{(G)}$ for testing
- Compute $F_{(G)} = Y_{(G)} - \hat{Y}_{(G)}$

Juuuust pretend I have a nice “Gth”
group separated out here...



Q² Calculations and Interpretation

- Step 3.1: Calculate **PRESS**
 - $\text{PRESS} = \text{ssq}(F_{(1)}) + \text{ssq}(F_{(2)}) + \dots + \text{ssq}(F_{(G)})$
- Step 3.2: Calculate Q²
 - $Q^2 = 1 - \frac{\mathcal{V}(\text{predeicted } F_A)}{\mathcal{V}(Y)} \Rightarrow Q^2 = 1 - \frac{\text{PRESS}}{\mathcal{V}(Y)}$
- Interpretations of Q²
 - Interpreted the same way as R² (higher the better, 1 is best)
 - You should find that $Q^2 \leq R^2$ (unless you are very lucky)
 - If $Q^2 \approx R^2$ for an ath component, the component is **useful**
 - If Q² is very small, likely fitting noise
 - Q² for an ath components CAN be negative (why?)
 - Q²_k can be calculated for specific variable k



EXAMPLE DATA SET

Polymer Reaction Process

Example: Polymers Reactor

- Consider the data set available on A2L called `LDPE.csv`
- This data set contains 14 columns in **X** that describe two tubular polymerization reactors in series:
 - T_{in} : inlet temperature to zone 1 of the reactor [K]
 - T_{max1} : maximum temperature along zone 1 [K]
 - T_{out1} : outlet temperature from zone 1 [K]
 - T_{max2} : maximum temperature along zone 2 [K]
 - T_{out2} : outlet temperature from zone 2 [K]
 - T_{cin1} : temperature of inlet coolant to zone [K]
 - T_{cin2} : temperature of inlet coolant to zone 2 [K]
 - $z1$: percentage along zone 1 where T_{max1} occurs [%]
 - $z2$: percentage along zone 2 where T_{max2} occurs [%]
 - F_{i1} : flow rate of initiators to zone 1 [g/s]
 - F_{i2} : flow rate of initiators to zone 2 [g/s]
 - F_{s1} : flow rate of solvent to zone 1 [% of ethylene]
 - F_{s2} : flow rate of solvent to zone 2 [% of ethylene]
 - $Press$: pressure in the reactor [atm]



Example: Polymers Reactor

- Consider the data set available on A2L called `LDPE.csv`
- This data set also contains 5 columns in **Y** that describe several quality metrics of the product out of the reactors:
 - Conv: quality variable: cumulative conversion
 - Mn: quality variable: number average molecular weight
 - Mw: quality variable: weight average molecular weight
 - LCB: quality variable: long chain branching per 1000 C atoms
 - SCB: quality variable: short chain branching per 1000 C atoms
- We would like to build a PLS model on this data using TWO components



Analyzing Data with PLS Plots

Y predict Y if you are not going to predict Y?

PLS Plots (Explored in Demo)

- Score plots
 - t and u show relationships between rows
- Weight plots
 - w shows **deflated** relationship between columns in X
 - w^* shows **un-deflated** relationship between columns in X (preferred)
- Loadings plots
 - c shows relationship between columns in Y
- Weight AND loadings plots
 - Super-imposed w^* and c plots show relationship between X and Y
- Modeling statistics plots
 - SPE (for X and Y spaces) to know if outlier off model plane
 - T^2 to know if outlier on model plane
 - R^2 plots for overall, by variable in X or Y , etc.
- Coefficient plot (described next)



Coefficient Plots

- Recall I mentioned a “quick” way to get \hat{y}_{new} :

$$\hat{y}_{new} = \mathbf{x}_{new}^T \mathbf{W}^* \mathbf{C}^T$$

- This actually kind of looks like a linear regression, no?

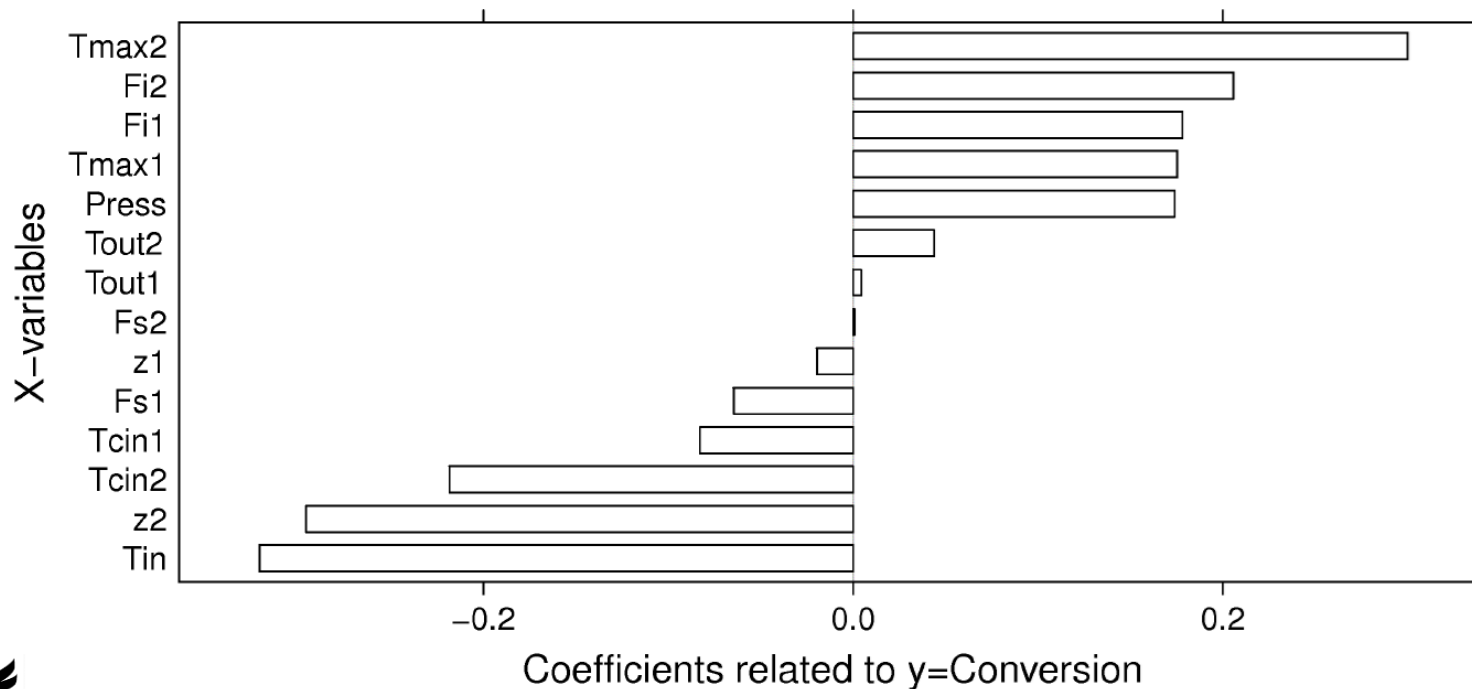
$$\hat{y}_{new} = \mathbf{x}_{new}^T \boldsymbol{\beta}$$

- $\boldsymbol{\beta} \in \mathbb{R}^{(K \times N)}$ is a collection of “importance” of each column in X as it relates to each column in Y
 - For one column in $\boldsymbol{\beta}$, we would just end up with the MLR equation for regressing y onto X !
 - However, it is critical to note that this is **DIFFERENT** than MLR



Example Coefficient Plot

- Consider our polymer reaction process
 - Below plot relates X to conversion (y) with 5 components
 - $\hat{y} = \beta_1 x_1 + \beta_2 x_2 + \dots + \beta_K x_K \quad \Leftrightarrow \quad \hat{y} = \mathbf{x}^T \boldsymbol{\beta}$
 - IMPORTANT**: this is good for visualization but you should **NEVER** implement PLS this way



Key Points

- PLS is an incredible modeling tool relating X to Y
 - Simultaneously explains variance and relationships
- NIPALS determines orthogonal weights and loadings that meet the objective of maximizing $\mathcal{C}(\mathbf{t}, \mathbf{u})$
- We don't want to rely on scores in the Y space in future
 - So, we calculate loadings \mathbf{p} from converged scores \mathbf{t} (fit using Y) so we can predict \mathbf{t} on new data
 - We also deflate X and Y based off of the scores in \mathbf{t}
- The weights \mathbf{w} are therefore annoying to analyze
 - Rely on deflated matrix X_a and are hard to interpret
 - Modified weights \mathbf{w}^* represent "projections" of weights \mathbf{w} onto X -space loadings \mathbf{p}
- Cross-validation is possible (just like PCA) for PLS
 - Want to add enough components for Q^2 to continue improving
- Many plots can be used to visualize PLS results
 - Visualization of components, data, scores, regression "coefficients" and more!



And Finally...

- We are done with Latent Variable Methods
 - More practice to come in A4!
- We will now move on to machine learning
 - Well, really just ANNs
- Focus for ANN section
 - Derivations
 - Key fundamental understanding of process/results
 - Perspective
 - Applications using built-in software

