

# **MTTTS16 Learning from Multiple Sources**

**5 ECTS credits**

Autumn 2019, University of Tampere  
Lecturer: Jaakko Peltonen

**Lecture 1: Canonical correlation analysis**

## On this lecture:

- Canonical correlation analysis

# Part 1: Canonical Correlation Analysis

# Motivation for CCA

- Recap: **correlation** is one of the most basic statistical measures of dependency between two real-valued scalar variables.
- Covariance between two variables:  $\text{cov}(x,y) = E[(x-E[x])(y-E[y])]$   
**Depends on scales of the variables.**
- **Correlation coefficient** (Pearson's product-moment correlation coefficient) between two variables:

$$\text{corr}(x, y) = \frac{\text{cov}(x, y)}{\text{stdev}(x) \text{stdev}(y)} = \frac{E[(x - E[x])(y - E[y])]}{(E[(x - E[x])^2])^{1/2} (E[(y - E[y])^2])^{1/2}}$$

- For zero-mean variables: Terms can be estimated from samples  $\{(x^i, y^i)\}_{i=1, \dots, N}$

$$\text{corr}(x, y) = \frac{E[xy]}{(E[x^2]E[y^2])^{1/2}} \quad E_{ML}^{\hat{}}[xy] = \frac{1}{N} \sum_{i=1}^N x^i y^i$$

# Motivation for CCA

- If two variables have zero correlation (uncorrelated), they are **linearly independent**, otherwise they are **linearly dependent**.
- If two variables are linearly dependent, a linear predictor  $y = ax+b$  can predict one variable from the other with some accuracy.
- Example: People's height and weight are positively correlated; poverty rate and education level might be negatively correlated
- If two variables have zero correlation, a linear predictor cannot predict one variable from the other.
- Two statistically independent variables have zero correlation.

# Motivation for CCA

## Caveats:

- Two variables with zero correlation can be dependent in nonlinear ways. If two variables have nonzero correlation, they may also have nonlinear dependency that correlation is unable to measure. Nonlinear dependencies can be exploited by a nonlinear predictor.
- Correlation does not imply causation: if  $x$  is correlated with  $y$ , then  $x$  might cause  $y$ ,  $y$  might cause  $x$ , or  $x$  and  $y$  might be caused by other variables, some of which affect both  $x$  and  $y$ .

# Motivation for CCA

- **More complicated situation 1:** What about when we have several scalar variables  $x_1, x_2, \dots, x_K$ , and we are not sure how they are related to  $y$ ?
- It can happen that **individual** variables  $x_i$  are not strongly correlated with  $y$ , but there is a linear **combination** of variables that is strongly correlated to  $y$ .
- Example: total time of multi-hop bus journeys depends on the sum of journey lengths. The length of the first journey is not enough by itself to predict the total time.
- Idea: **which** linear combination of  $x_1, x_2, \dots, x_K$  is most strongly correlated with  $y$ ? Can we find the best linear combination  $a_1 x_1 + a_2 x_2 + \dots + a_K x_K$ , that is, the best weights  $a_1, \dots, a_K$ ?
- Yes: this turns out to be the same as **ordinary least-squares linear regression!**

# Canonical Correlation Analysis

- **More complicated situation 2:** What about when we have several scalar variables  $x_1, x_2, \dots, x_K$ , and also several scalar variables  $y_1, y_2, \dots, y_L$ , ? How can we now find what relationships exist between the  $x$  variables and the  $y$  variables?
- **Canonical correlation analysis (CCA)** is a method of correlating two multidimensional variables  $\mathbf{x} = [x_1, x_2, \dots, x_K]$  and  $\mathbf{y} = [y_1, y_2, \dots, y_L]$
- Proposed by H. Hotelling. (*Hotelling, H. Relations between two sets of variates. Biometrika, vol. 28, pages 312–377, 1936.*)
- Idea: not all of  $\mathbf{x}$  needs to be correlated with all of  $\mathbf{y}$  or vice versa, as long as there is something correlated between  $\mathbf{x}$  and  $\mathbf{y}$ .
- Problem definition: given two sets of variables, find **basis vectors** (linear transformations), one for each set of variables, so that the correlations between the linear projections of the variables onto the basis vectors are mutually maximized (the projected coordinates are maximally correlated).



# Canonical Correlation Analysis

- For  $\mathbf{x}$ , find a projection  $w_{x,1}x_1 + w_{x,2}x_2 + \dots + w_{x,K}x_K$  where  $\mathbf{w}_x = [w_{x,1}, w_{x,2}, \dots, w_{x,K}]$  is the projection basis.
- For  $\mathbf{y}$ , find a projection  $w_{y,1}y_1 + w_{y,2}y_2 + \dots + w_{y,L}y_L$  where  $\mathbf{w}_y = [w_{y,1}, w_{y,2}, \dots, w_{y,L}]$  is the projection basis.
- Find the projection bases by maximizing the correlation between the projections: maximize

$$\text{corr}(\mathbf{w}_x^T \mathbf{x}, \mathbf{w}_y^T \mathbf{y}) = \frac{E[\mathbf{w}_x^T \mathbf{x} \mathbf{w}_y^T \mathbf{y}]}{(E[(\mathbf{w}_x^T \mathbf{x})^2] E[(\mathbf{w}_y^T \mathbf{y})^2])^{1/2}}$$

with respect to  $\mathbf{w}_x$  and  $\mathbf{w}_y$ .

This definition assumes  $\mathbf{x}$  and  $\mathbf{y}$  are zero-mean, otherwise subtract the means as in the original correlation definition.

- For a finite data set: maximize the sample correlation

$$\hat{\text{corr}}(\mathbf{w}_x^T \mathbf{x}, \mathbf{w}_y^T \mathbf{y}) = \frac{\hat{E}_{ML}[\mathbf{w}_x^T \mathbf{x} \mathbf{w}_y^T \mathbf{y}]}{(\hat{E}_{ML}[(\mathbf{w}_x^T \mathbf{x})^2] \hat{E}_{ML}[(\mathbf{w}_y^T \mathbf{y})^2])^{1/2}}$$

Same definition as before

$$\hat{E}_{ML}[x y] = \frac{1}{N} \sum_{i=1}^N x^i y^i$$

# Canonical Correlation Analysis

- How can we find the maximal correlation with respect to  $\mathbf{w}_x$  and  $\mathbf{w}_y$ ?
- It turns out this can be rewritten as a **generalized matrix eigenvalue problem**. **Let's show how that is done.**
- First, rewrite the correlation to show it depends on covariance matrices:

$$\begin{aligned}
 \hat{c}orr(\mathbf{w}_x^T \mathbf{x}, \mathbf{w}_y^T \mathbf{y}) &= \frac{\hat{E}_{ML}[\mathbf{w}_x^T \mathbf{x} \mathbf{w}_y^T \mathbf{y}]}{(\hat{E}_{ML}[(\mathbf{w}_x^T \mathbf{x})^2] \hat{E}_{ML}[(\mathbf{w}_y^T \mathbf{y})^2])^{1/2}} \\
 &= \frac{\hat{E}_{ML}[\mathbf{w}_x^T \mathbf{x} \mathbf{y}^T \mathbf{w}_y]}{(\hat{E}_{ML}[\mathbf{w}_x^T \mathbf{x} \mathbf{x}^T \mathbf{w}_x] \hat{E}_{ML}[\mathbf{w}_y^T \mathbf{y} \mathbf{y}^T \mathbf{w}_y])^{1/2}} \\
 &= \frac{\mathbf{w}_x^T \hat{E}_{ML}[\mathbf{x} \mathbf{y}^T] \mathbf{w}_y}{(\mathbf{w}_x^T \hat{E}_{ML}[\mathbf{x} \mathbf{x}^T] \mathbf{w}_x \mathbf{w}_y^T \hat{E}_{ML}[\mathbf{y} \mathbf{y}^T] \mathbf{w}_y)^{1/2}} \\
 &= \frac{\mathbf{w}_x^T \hat{C}_{x,y} \mathbf{w}_y}{(\mathbf{w}_x^T \hat{C}_x \mathbf{w}_x \mathbf{w}_y^T \hat{C}_y \mathbf{w}_y)^{1/2}}
 \end{aligned}$$

Sample estimate of the covariance matrix

$$\hat{C}_{x,y} = \frac{1}{N} \sum_{i=1}^N \mathbf{x}^i \mathbf{y}^{iT}$$

$$\hat{C}_x = \frac{1}{N} \sum_{i=1}^N \mathbf{x}^i \mathbf{x}^{iT}$$

$$\hat{C} = \begin{bmatrix} \hat{C}_x & \hat{C}_{x,y} \\ \hat{C}_{y,x} & \hat{C}_y \end{bmatrix}$$

# Canonical Correlation Analysis

## Proof continues...

- Notice that the scales of  $\mathbf{w}_x$  and  $\mathbf{w}_y$  do not affect the correlation, so we can set any requirement that can be satisfied by changing the scales alone. For example, we can require  $\mathbf{w}_x^T \hat{C}_x \mathbf{w}_x = 1$  ,  $\mathbf{w}_y^T \hat{C}_y \mathbf{w}_y = 1$
- Then the optimization becomes a constrained optimization problem:

$$\max_{\mathbf{w}_x, \mathbf{w}_y} [\mathbf{w}_x^T \hat{C}_{x,y} \mathbf{w}_y] \quad \text{such that} \quad \mathbf{w}_x^T \hat{C}_x \mathbf{w}_x = 1, \quad \mathbf{w}_y^T \hat{C}_y \mathbf{w}_y = 1$$

- Constrained optimization problems can be solved by the method of **Lagrange multipliers**, that is we maximize the Lagrangian

$$\max_{\mathbf{w}_x, \mathbf{w}_y} L(\mathbf{w}_x, \mathbf{w}_y, \lambda_x, \lambda_y)$$

$$L(\mathbf{w}_x, \mathbf{w}_y, \lambda_x, \lambda_y) = \left( \mathbf{w}_x^T \hat{C}_{x,y} \mathbf{w}_y - \frac{\lambda_x}{2} (\mathbf{w}_x^T \hat{C}_x \mathbf{w}_x - 1) - \frac{\lambda_y}{2} (\mathbf{w}_y^T \hat{C}_y \mathbf{w}_y - 1) \right)$$

# Canonical Correlation Analysis

Proof continues...

- At the optimum, the derivative of the Lagrangian must be zero with respect to  $\mathbf{w}_x$  and  $\mathbf{w}_y$ . This gives

$$\frac{dL}{d\mathbf{w}_x} = \hat{C}_{x,y} \mathbf{w}_y - \lambda_x \hat{C}_x \mathbf{w}_x = 0$$

$$\mathbf{w}_x^T \hat{C}_x \mathbf{w}_x = 1 \quad \mathbf{w}_y^T \hat{C}_y \mathbf{w}_y = 1$$

$$\frac{dL}{d\mathbf{w}_y} = \hat{C}_{y,x} \mathbf{w}_x - \lambda_y \hat{C}_y \mathbf{w}_y = 0$$

subtract the  $\mathbf{w}_x$  times the second equation from  $\mathbf{w}_y$  times the first:

$$\begin{aligned} \mathbf{w}_x^T \hat{C}_{x,y} \mathbf{w}_y - \lambda_x \mathbf{w}_x^T \hat{C}_x \mathbf{w}_x - \mathbf{w}_y^T \hat{C}_{y,x} \mathbf{w}_x + \lambda_y \mathbf{w}_y^T \hat{C}_y \mathbf{w}_y \\ = -\lambda_x \mathbf{w}_x^T \hat{C}_x \mathbf{w}_x + \lambda_y \mathbf{w}_y^T \hat{C}_y \mathbf{w}_y \\ = -\lambda_x + \lambda_y \\ = 0 \end{aligned}$$

because we  
required

$$\mathbf{w}_x^T \hat{C}_x \mathbf{w}_x = 1$$

$$\mathbf{w}_y^T \hat{C}_y \mathbf{w}_y = 1$$

# Canonical Correlation Analysis

Proof continues...

- Therefore the Lagrange multipliers are equal, and

$$\frac{dL}{d\mathbf{w}_y} = \hat{C}_{y,x} \mathbf{w}_x - \lambda \hat{C}_y \mathbf{w}_y = 0 \rightarrow \mathbf{w}_y = (1/\lambda) \hat{C}_y^{-1} \hat{C}_{y,x} \mathbf{w}_x$$

- Inserting that into the other derivative equation gives

$$\hat{C}_{x,y} \hat{C}_y^{-1} \hat{C}_{y,x} \mathbf{w}_x = \lambda^2 \hat{C}_x \mathbf{w}_x$$

- This is a generalized eigenvalue equation which we can solve to get  $\mathbf{w}_x$ , and the previous equation then gives  $\mathbf{w}_y$  from  $\mathbf{w}_x$ .
- If the covariance of  $\mathbf{x}$  is invertible, this can also be written as a standard eigenproblem with some additional steps

# Canonical Correlation Analysis

## Generalizations of CCA:

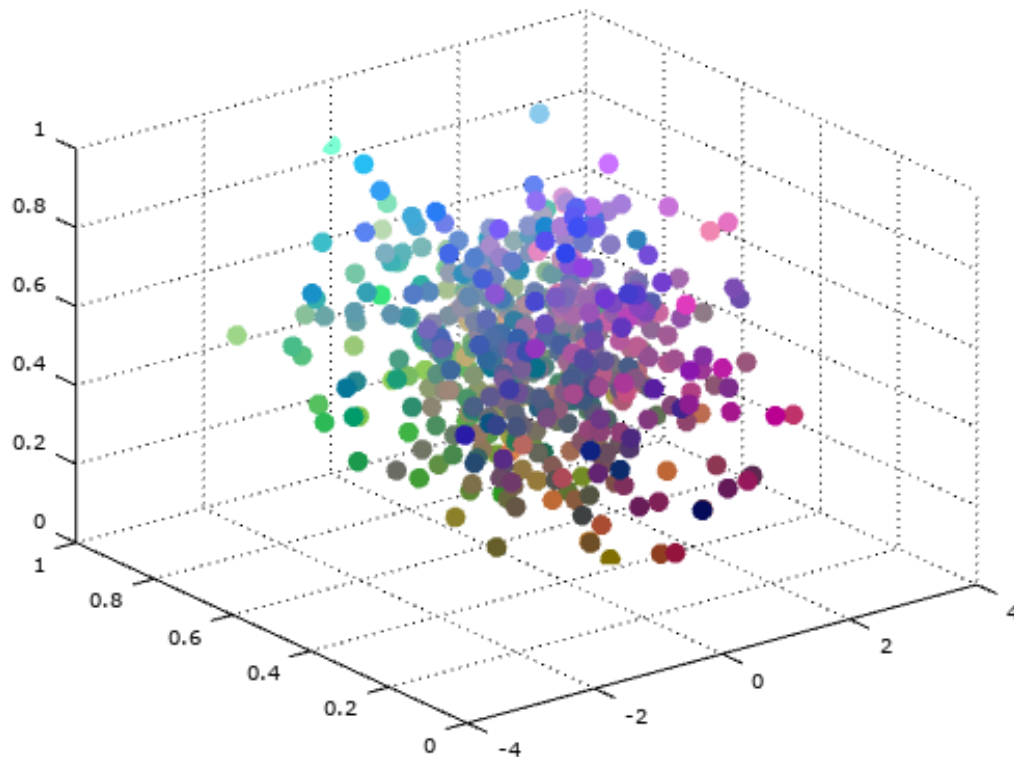
- Several possible generalizations of CCA to more than two sets of variables have been proposed. One is (Hardoon, Szedmak, Taylor):

$$\min_{\mathbf{W}^{(1)}, \dots, \mathbf{W}^{(M)}} \sum_{m, n=1, m \neq n}^M \left\| \mathbf{W}^{(m)T} \mathbf{x}^{(m)} - \mathbf{W}^{(n)T} \mathbf{x}^{(n)} \right\|_F$$
$$\mathbf{W}^{(m)T} \hat{\mathbf{C}}_{mm} \mathbf{W}^{(m)} = \mathbf{I} \quad \mathbf{W}^{(m)T} \hat{\mathbf{C}}_{mn} \mathbf{W}^{(n)} = 0$$

- CCA can be shown to be a generative model of data (next on this lectures).
- CCA can be generalized to analyze nonlinear correlations through a kernel mapping (later in the lectures)
- CCA can be incorporated with sparsity and other prior desires for the projection bases

# Canonical Correlation Analysis, example

- In this example  $\mathbf{x}$  and  $\mathbf{y}$  are both 3-dimensional.
- $\mathbf{x} = [x_1 \ x_2 \ x_3]$  is randomly distributed according to a three-dimensional isotropic (spherical) Gaussian
- $\mathbf{y} = [y_1 \ y_2 \ y_3] = [\text{red green blue}]$ , where “red” is proportional to  $x_1 + x_2 + x_3$  and green and blue are randomly normally distributed.



In this picture  $\mathbf{x}$  is shown as 3D spatial coordinates, and  $\mathbf{y}$  is shown as red-green-blue (RGB) color components

Matlab/octave code:

```
x=randn(500,3);
y=randn(500,3);
y(:,1)=sum(x,2);
for k=1:3,
y(:,k)=(y(:,k)-min(y(:,k)))/(max(y(:,k))-
min(y(:,k))); end;
```

```
scatter3(x(:,1),x(:,2),x(:,3), ...
10*ones(size(x,1),1),y,'filled');
```

# Canonical Correlation Analysis, example

- Let's compute and solve the eigenvalue equation.

```
C = cov([x y]);
Cx=C(1:3,1:3);
Cy=C(4:6,4:6);
Cxy=C(1:3,4:6);
A=Cxy*inv(Cy)*Cxy';
B=Cx;
[Wx,D]=eig(A,B); % A*Wx=B*Wx*D
Wy=inv(Cy)*Cxy'*Wx*inv(D);
```

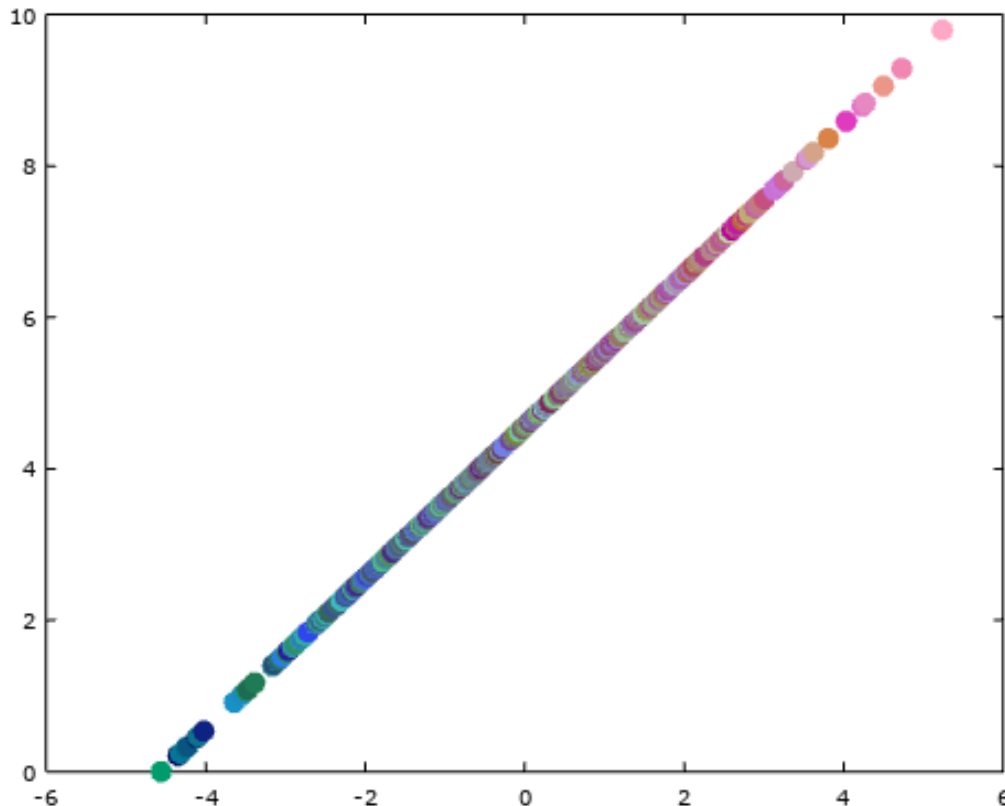
```
Wx =
1.00000  0.82944 -0.70922
1.00000  0.30068  1.00000
1.00000 -1.00000 -0.40489
```

```
D =
1.0000e+00    0    0
0  1.6523e-04    0
0    0  5.8796e-03
```

The 1<sup>st</sup> column of  $W_x$  projects  $x$  to  $x_1+x_2+x_3$  (corresponds to largest eigenvalue in  $D$ )

```
Wy =
9.7948e+00  2.2388e+01 -1.5861e+00
-3.0531e-16 -6.2843e+02 -4.2152e+01
5.2736e-16 -2.2199e+02  8.7614e+01
```

1<sup>st</sup> column of  $W_y$  projects  $y$  to  $y_1$



Resulting projections match perfectly (correlation 1)

```
xp=x*Wx(:,1);
yp=y*Wy(:,1);
```

```
figure;
scatter(xp,yp,10*ones(size(x,1),1),y,'filled');
```



# Canonical Correlation Analysis, example

## Connection to least-squares regression:

- If  $y$  is one-dimensional (scalar):

let  $\mathbf{X}_{all} = [\mathbf{x}^1 \ \mathbf{x}^2 \ \dots \ \mathbf{x}^N]$  and  $\mathbf{y}_{all} = [y^1 \ y^2 \ \dots \ y^N]^T$ , then we have

$$\hat{C}_{x,y} = \frac{1}{N} \sum_{i=1}^N \mathbf{x}^i y^i = \frac{1}{N} \mathbf{X}_{all} \mathbf{y}_{all}$$

$$\hat{C}_x = \frac{1}{N} \sum_{i=1}^N \mathbf{x}^i \mathbf{x}^{iT} = \frac{1}{N} \mathbf{X}_{all} \mathbf{X}_{all}^T \quad \hat{C}_y = \frac{1}{N} \sum_{i=1}^N y^i y^i = \frac{1}{N} \mathbf{y}_{all}^T \mathbf{y}_{all}$$

- The ordinary (least-squared) linear regression solution is

$\mathbf{w}_x = (\mathbf{X}_{all} \mathbf{X}_{all}^T)^{-1} \mathbf{X}_{all} \mathbf{y}_{all}$ , we show it satisfies the eigenvalue equation.

- We must show that the solution satisfies for some scalar  $\lambda^2$   $\hat{C}_{x,y} \hat{C}_y^{-1} \hat{C}_{y,x} \mathbf{w}_x = \lambda^2 \hat{C}_x \mathbf{w}_x$

# Canonical Correlation Analysis, example

Denote the singular value decomposition  $\mathbf{X}_{all} = \mathbf{U}\mathbf{D}\mathbf{V}$

- Right-hand side of the eigenvalue equation:

$$\lambda^2 \hat{C}_x \mathbf{w}_x = \lambda^2 \left( \frac{1}{N} \mathbf{X}_{all} \mathbf{X}_{all}^T \right) \left( \mathbf{X}_{all} \mathbf{X}_{all}^T \right)^{-1} \mathbf{X}_{all} \mathbf{y}_{all} = \frac{1}{N} \lambda^2 \mathbf{X}_{all} \mathbf{y}_{all}$$

- Left-hand side of the eigenvalue equation:  $C_{x,y} \hat{C}_y^{-1} C_{y,x} \mathbf{w}_x =$

$$\left( \frac{1}{N} \mathbf{X}_{all} \mathbf{y}_{all} \right) \left( \frac{1}{N} \mathbf{y}_{all}^T \mathbf{y}_{all} \right)^{-1} \left( \frac{1}{N} \mathbf{X}_{all} \mathbf{y}_{all} \right)^T \left( \mathbf{X}_{all} \mathbf{X}_{all}^T \right)^{-1} \mathbf{X}_{all} \mathbf{y}_{all} =$$

$$\left( \left( N \mathbf{y}_{all}^T \mathbf{y}_{all} \right)^{-1} \right) \mathbf{X}_{all} \mathbf{y}_{all} \mathbf{y}_{all}^T \mathbf{X}_{all}^T \left( \mathbf{X}_{all} \mathbf{X}_{all}^T \right)^{-1} \mathbf{X}_{all} \mathbf{y}_{all} =$$

$$\left( \left( N \mathbf{y}_{all}^T \mathbf{y}_{all} \right)^{-1} \right) \mathbf{U}\mathbf{D}\mathbf{V} \mathbf{y}_{all} \mathbf{y}_{all}^T \mathbf{V}^T \mathbf{D}\mathbf{U}^T \left( \mathbf{U}\mathbf{D}\mathbf{V} \mathbf{V}^T \mathbf{D}\mathbf{U}^T \right)^{-1} \mathbf{U}\mathbf{D}\mathbf{V} \mathbf{y}_{all} =$$

$$\left( \left( N \mathbf{y}_{all}^T \mathbf{y}_{all} \right)^{-1} \right) \mathbf{U}\mathbf{D}\mathbf{V} \mathbf{y}_{all} \mathbf{y}_{all}^T \mathbf{V}^T \mathbf{D}\mathbf{U}^T \left( \mathbf{U} \mathbf{D}^2 \mathbf{U}^T \right)^{-1} \mathbf{U}\mathbf{D}\mathbf{V} \mathbf{y}_{all} =$$

$$\left( \left( N \mathbf{y}_{all}^T \mathbf{y}_{all} \right)^{-1} \right) \mathbf{U}\mathbf{D}\mathbf{V} \mathbf{y}_{all} \mathbf{y}_{all}^T \mathbf{V}^T \mathbf{D}\mathbf{U}^T \mathbf{U} \mathbf{D}^{-2} \mathbf{U}^T \mathbf{U}\mathbf{D}\mathbf{V} \mathbf{y}_{all} =$$

$$\left( \left( N \mathbf{y}_{all}^T \mathbf{y}_{all} \right)^{-1} \right) \mathbf{U}\mathbf{D}\mathbf{V} \mathbf{y}_{all} \mathbf{y}_{all}^T \mathbf{V}^T \mathbf{D} \mathbf{D}^{-2} \mathbf{D}\mathbf{V} \mathbf{y}_{all} =$$

$$\left( \left( N \mathbf{y}_{all}^T \mathbf{y}_{all} \right)^{-1} \right) \mathbf{U}\mathbf{D}\mathbf{V} \mathbf{y}_{all} \mathbf{y}_{all}^T \mathbf{V}^T \mathbf{V} \mathbf{y}_{all} = \left( \left( N \mathbf{y}_{all}^T \mathbf{y}_{all} \right)^{-1} \right) \mathbf{U}\mathbf{D}\mathbf{V} \mathbf{y}_{all} \mathbf{y}_{all}^T \mathbf{y}_{all} =$$

$$\left( N^{-1} \right) \mathbf{U}\mathbf{D}\mathbf{V} \mathbf{y}_{all} = \left( N^{-1} \right) \mathbf{X}_{all} \mathbf{y}_{all} \quad \text{Equality satisfied with lambda=1}$$