

Exercises from ‘Regression and SVD’

1. Show that for V_r defined in the SVD of ϕ , $V_r V_r^T$ is a projection matrix and that $\text{tr}(V_r V_r^T) = r$.

$$(V_r V_r^T)^2 = V_r V_r^T$$

$$\text{We know } \phi^T \phi = U_r \Lambda_r U_r^T \Rightarrow (U_r^T)^{-1} \Lambda_r^{-1} U_r^{-1} = (\phi^T \phi)^{-1}$$

$$\text{We also know } \phi = V_r \Lambda_r^{1/2} U_r^T$$

$$\Rightarrow V_r = \phi (U_r^T)^{-1} \Lambda_r^{-1/2}$$

$$\Rightarrow V_r^T = \Lambda_r^{-1/2} (U_r)^{-1} \phi^T$$

$$\text{Because } (\Lambda_r^{-1/2})^T = \Lambda_r^{-1/2} \text{ and } ((U_r^T)^{-1})^T = (U_r)^{-1}$$

$$(V_r V_r^T)^2 = V_r V_r^T V_r V_r^T$$

$$\begin{aligned} V_r V_r^T V_r V_r^T &= V_r (\Lambda_r^{-1/2} (U_r)^{-1} \phi^T) (\phi (U_r^T)^{-1} \Lambda_r^{-1/2}) V_r^T \\ &= V_r (\Lambda_r)^{-1/2} (U_r)^{-1} (\phi^T \phi) (U_r^T)^{-1} (\Lambda_r)^{-1/2} V_r^T \\ &= V_r (\Lambda_r)^{-1/2} (U_r)^{-1} U_r \Lambda_r U_r^T (U_r^T)^{-1} (\Lambda_r)^{-1/2} V_r^T \\ &= V_r (\Lambda_r)^{-1/2} \Lambda_r (\Lambda_r)^{-1/2} V_r^T \\ &= V_r V_r^T \end{aligned}$$

$$\text{Show the trace} = V_r V_r^T = r$$

$$\text{tr}(V_r V_r^T) = \text{tr}(V_r^T V_r)$$

$$\text{tr}(V_r^T V_r) = \text{tr}(\Lambda_r^{-1/2} (U_r)^{-1} \phi^T \phi (U_r^T)^{-1} \Lambda_r^{-1/2})$$

$$= \text{tr}(\Lambda_r^{-1/2} (U_r)^{-1} U_r \Lambda_r U_r^T (U_r^T)^{-1} \Lambda_r^{-1/2})$$

$$= \text{tr}(\Lambda_r^{-1/2} \Lambda_r \Lambda_r^{-1/2})$$

$$= \text{tr}(I_{(r)}) = r$$

2. Weighted least squares: Consider the following small variation of the regression problem. The observations are not homoscedastic, that is $\text{var}(y_i)$ changes with x_i . Mathematically, the linear regression model now writes as $y = \phi^T w + \epsilon$ where D is a known precision matrix (typically diagonal in practice). Show that the normal equations in this case become:

$$\begin{aligned} N(t|x, \phi w, (qD)^{-1}) &= 1/(2\pi)^{D/2} (qD)^{1/2} e^{(-1/2(t-\phi w)^T qD(t-\phi w))} \\ \log(N(t|x, \phi w, (qD)^{-1})) &= \log(1/(2\pi)^{D/2}) + (1/2)\log(qD) - (1/2)(t-\phi w)^T qD(t-\phi w) \\ &= \log(1/(2\pi)^{D/2}) + (1/2)\log(qD) - (1/2)(q)(t^T D t - t^T D \phi w - w^T \phi^T D t + w^T \phi^T D \phi w) \end{aligned}$$

Differentiating with respect to w and equating to 0. $(-1/2)$ and (q) get cancelled off.

$$0 = -t^T D \phi - t^T D^T \phi + w^T (\phi^T D \phi + (\phi^T D \phi)^T)$$

$$2t^T D^T \phi = w^T (2\phi^T D \phi) \text{ because } D^T = D \text{ and } (\phi^T D \phi)^T = (\phi^T D^T (\phi^T)^T) = \phi^T D \phi$$

$$t^T D^T \phi = w^T \phi^T D \phi$$

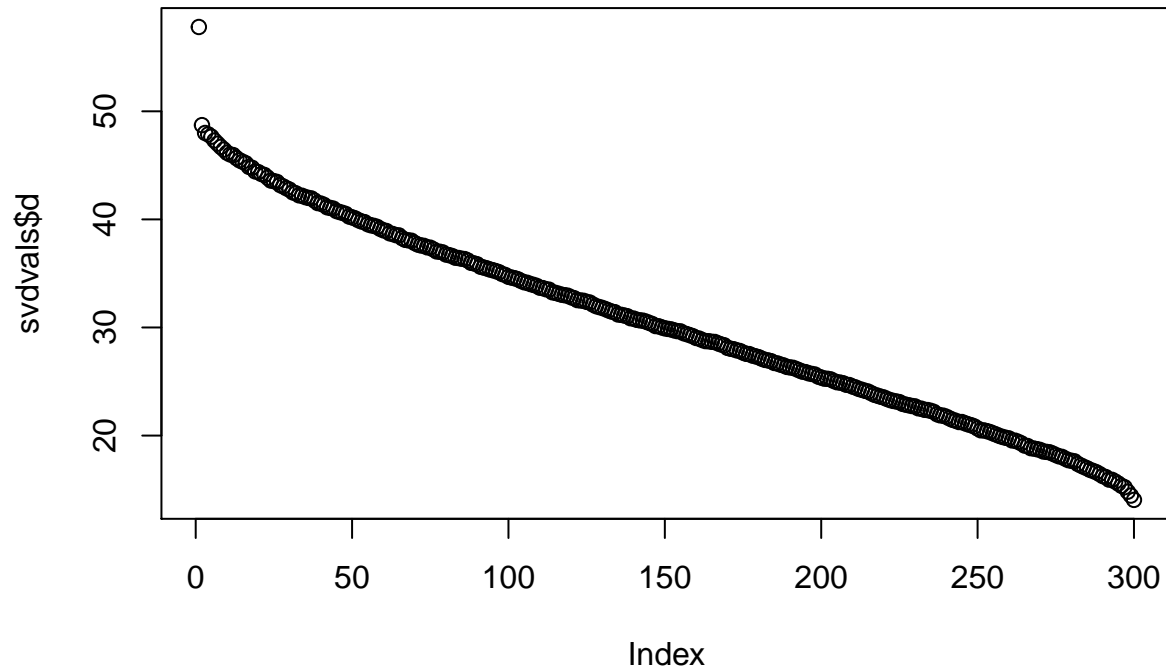
Taking Transpose on both sides:

$$\phi^T D(t^T)^T = \phi^T D^T (\phi^T)^T (w^T)^T$$

$$\phi^T D t = \phi^T D \phi w$$

3. Extract the first 300 rows and 1001 columns from the `synthetic_regression.txt` dataset; hence you have 300 replications and 1000 input variables.

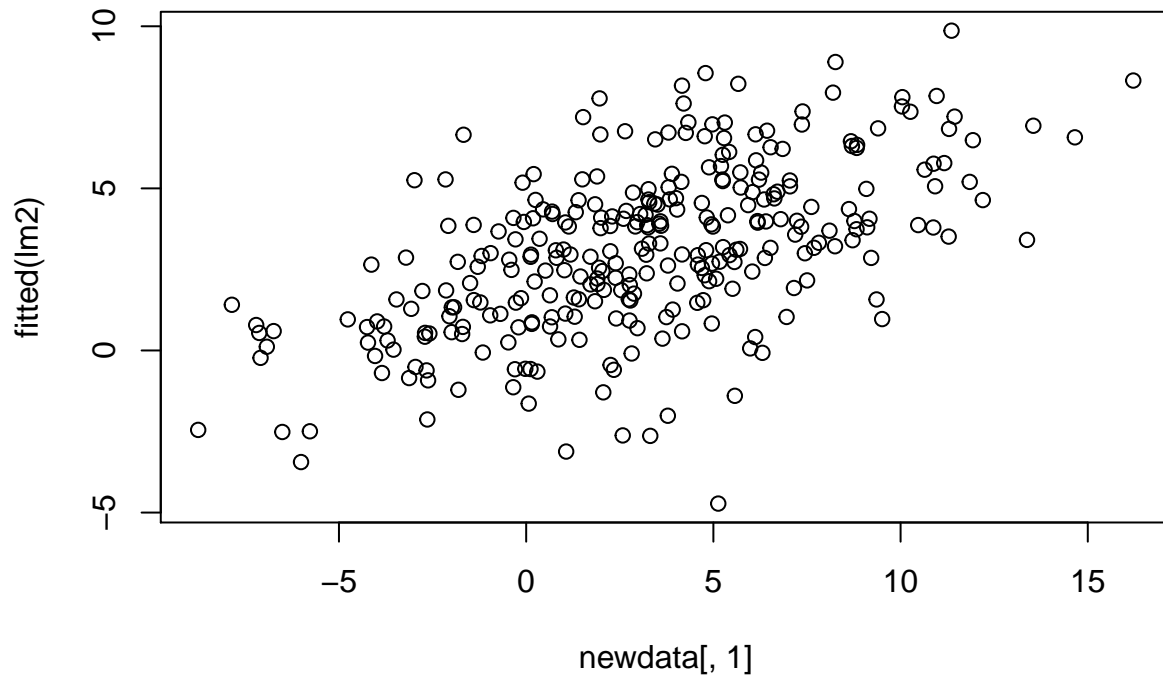
3.1 Plot the non-zero singular values of the input matrix.



3.2 Determine r , the rank of the input matrix.

300

3.3 Do a principal component regression based on the largest 30 eigenvalues of the input matrix and by including an intercept term (hence 31 features in total). Plot fitted vs observed values according to this model. Compare the fit of this model to the one you obtained in a previous exercise using the first (in order of appearance) 30 input variables and the intercept. Is it better or worse? Explain in at most 3 lines your findings and explanation.



The fit of the PCA model is slightly worse than the fit of the linear model.

The reason is the linear model minimizes the error with respect to the observed data whereas the PCA minimizes the error with respect to the model itself. So the linear model is a better fit for the observed data, where it is considered in the fitting, whereas the PCA is a better fit of the included features themselves.