

Data Preprocessing

Data Preprocessing

Why?

- Reduce scale effects
- Reduce noise
- Reduce redundancy

Data Preprocessing

Reducing scale effects

Normalization – map every feature to the $[0,1]$ interval

For every feature i in data set:

- $\text{max_val} = \max(X_{\text{train}}[:,i])$
- $\text{min_val} = \min(X_{\text{train}}[:,i])$
- $X_{\text{train_n}}[:,i] = (X_{\text{train}}[:,i] - \text{min_val}) / (\text{max_val} - \text{min_val})$
- $X_{\text{test_n}}[:,i] = (X_{\text{test}}[:,i] - \text{min_val}) / (\text{max_val} - \text{min_val})$

Advantages:

- Simplicity
- Same interval for all feature values

Disadvantages:

- Sensitive to outliers

Data Preprocessing

Reducing scale effects

Standardization (a.k.a. z-normalization) – scale every feature to have zero mean and unit standard deviation

For every feature i in data set:

- $m = \text{mean}(X_{\text{train}}[:,i])$
- $s = \text{s}(X_{\text{train}}[:,i])$
- $X_{\text{train}_n}[:,i] = (X_{\text{train}}[:,i] - m) / s$
- $X_{\text{test}_n}[:,i] = (X_{\text{test}}[:,i] - m) / s$

Advantages:

- Simplicity
- Tends to produce good model performance

Disadvantages:

- Sensitive to outliers (not as much as normalization)

Data Preprocessing

Reducing noise and redundancy

Principal Component Analysis

Idea: Project data from n -dimensions (attributes) to m -dimensions (with $n > m$) while preserving as much information as possible

This is achieved by a simple combination of translation and rotation

Why: Features tend to be correlated – reduce redundancy

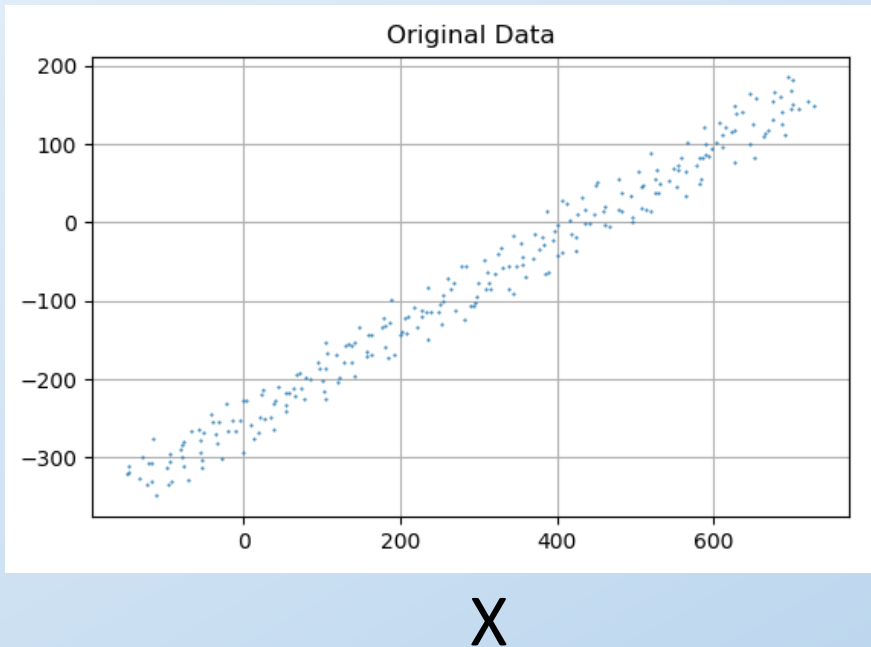
Why: Least important features in projected space approximate noise

Data Preprocessing

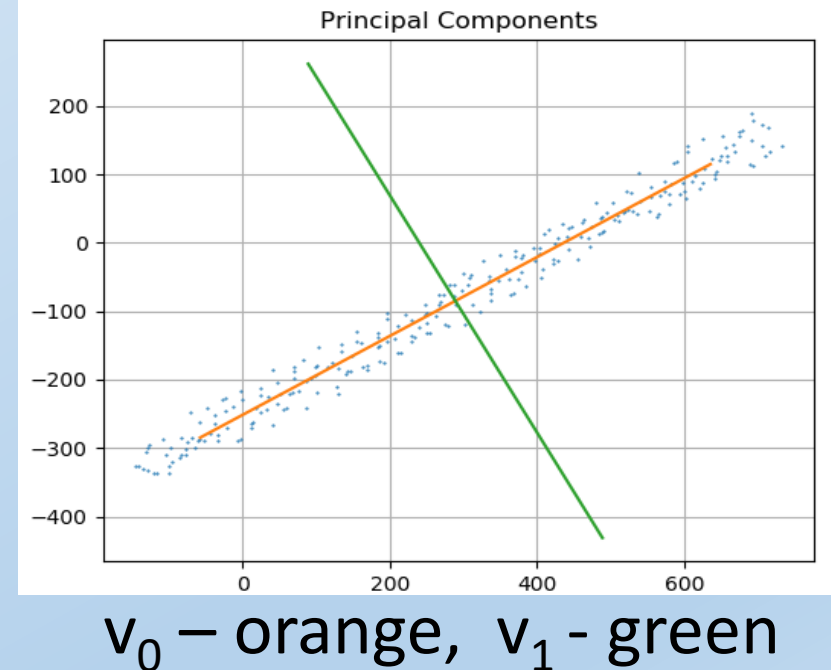
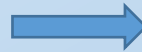
Reducing noise and redundancy

Principal Component Analysis

The principal components of a dataset X are the vectors v_0, \dots, v_n such that v_i is the vector that best fits X and is perpendicular to each of v_0, \dots, v_{i-1}



PCA

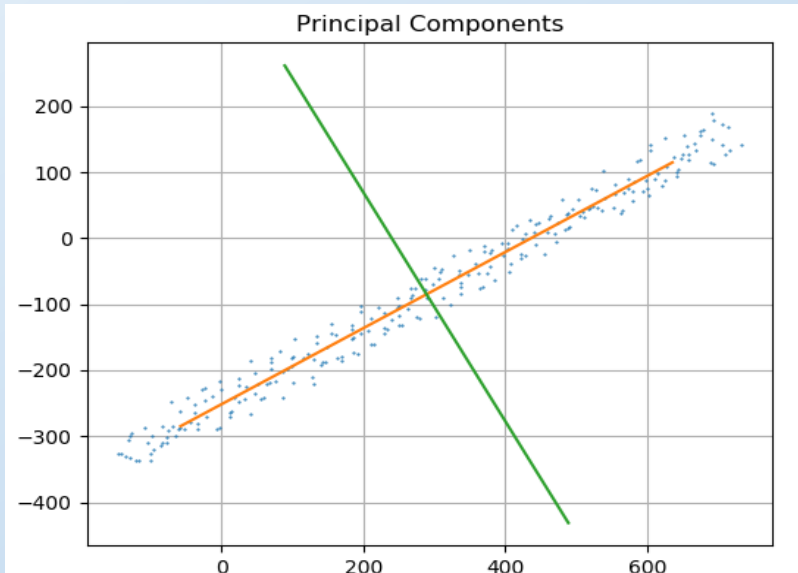


Data Preprocessing

Reducing noise and redundancy

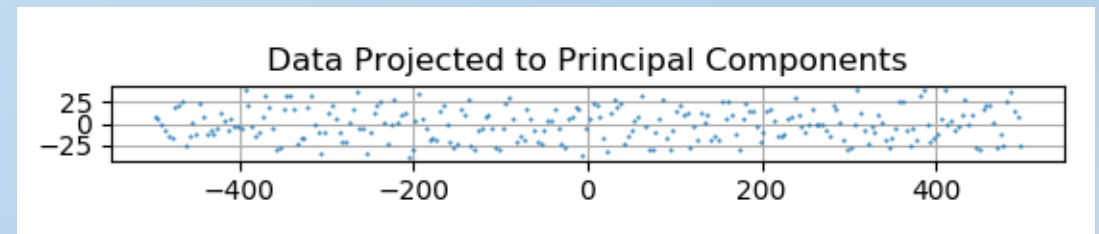
Principal Component Analysis

The principal components of a dataset X are the vectors v_0, \dots, v_n such that v_i is the vector that best fits X and is perpendicular to each of v_0, \dots, v_{i-1}



X in original feature space and principal components

Projection

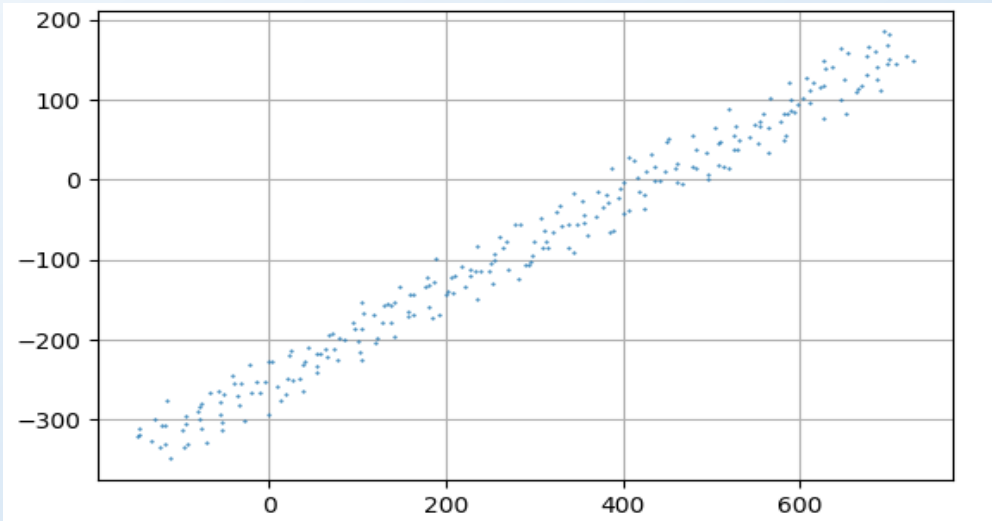



X in new feature space

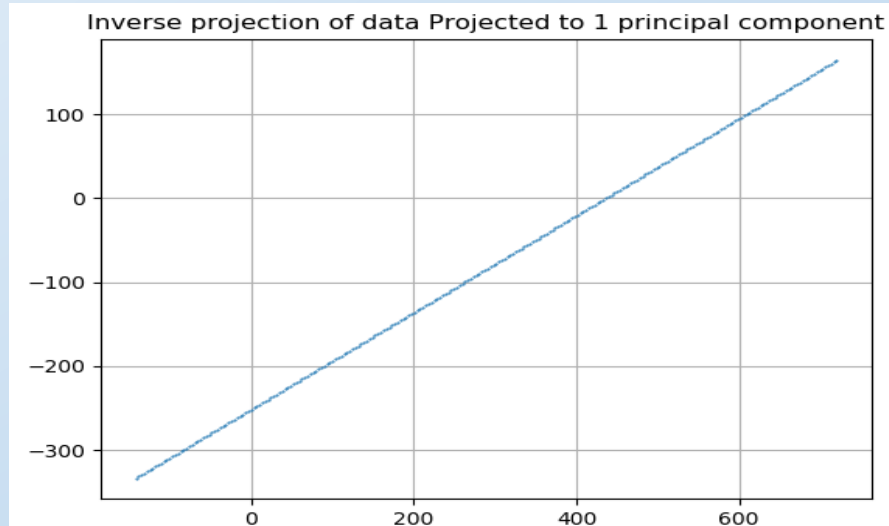
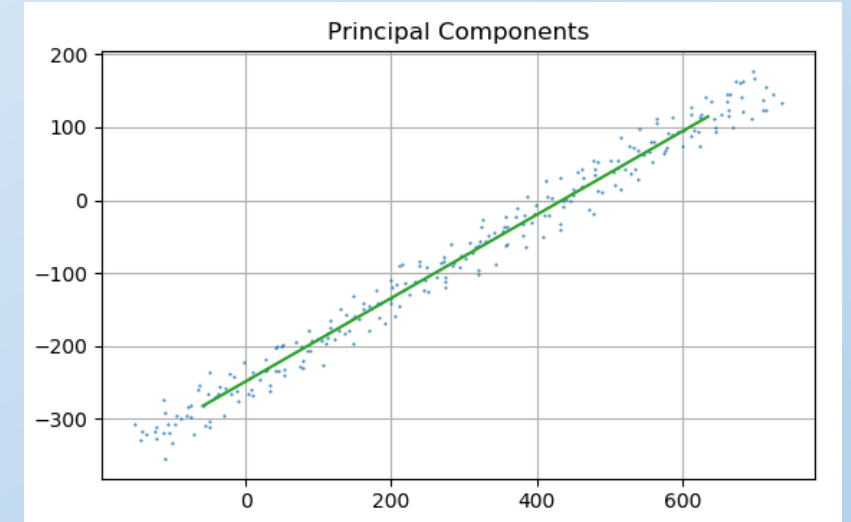
Features with smallest variance can be discarded

Data Preprocessing

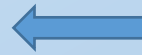
Principal Component Analysis for dimensionality reduction




`pca.fit(X,n_components=1)`



$X_f =$
`pca.inverse_transform(P)`



`P = pca.transform(X)`

