

Principal Components

$$\text{Let } \bar{C} = \bar{X}^T \bar{X} = \sum_{i=1}^n \vec{x}_i \vec{x}_i^T$$

\bar{C} has eigenvalue decomposition

$$\bar{C} = \bar{V} \bar{D} \bar{V}^T \quad \bar{D} = \begin{bmatrix} \lambda_1 & & 0 \\ & \ddots & \\ 0 & & \lambda_d \end{bmatrix}, \quad V V^T = I$$

$\lambda_i \geq 0$ because \bar{C} is positive semidefinite

$$(\text{if } \bar{C} \bar{z} = \lambda_i \bar{z}, \quad \bar{z}^T \bar{C} \bar{z} \geq 0 \Rightarrow \lambda_i \geq 0)$$

Recall that if $d \leq n$ and \bar{X} has full rank

$$\vec{w}_{OLS} = (\bar{X}^T \bar{X})^{-1} \bar{X}^T \vec{y}$$

$$\text{Now: } \vec{w}_{OLS} = \sum_{i=1}^d \frac{\alpha_i}{\lambda_i} \vec{v}_i$$

$$\text{Where } \alpha_i = \sum_{j=1}^n y_j \underbrace{\langle \vec{x}_j, \vec{v}_i \rangle}_{\text{value of } \vec{x}_j}$$

similarity
of
 \vec{x}_j and \vec{v}_i

Decomposes \vec{w}_{OLS} into d -feature directions which are a basis. Directions weighted by

λ_i : presence of signal in data

α_i : interpolated value in v_i direction.

NOTE: Small λ_i get amplified by $\frac{1}{\lambda_i}$!

SVD: $n \geq d$

$$\bar{X} = \bar{U} \bar{S} \bar{V}^T$$

$$\bar{U} \quad n \times n, \quad \bar{V} \quad d \times d$$

$$\bar{S} = \begin{bmatrix} \sigma_1 & & & \\ & \ddots & & \\ & & \sigma_d & \\ & & & 0 \end{bmatrix} \}_{n-d}$$

$$\bar{C} = \bar{X}^T \bar{X} = \bar{V} (\bar{S}^T \bar{S}) \bar{V}^T \quad \Rightarrow \quad \sigma_i^2 = \lambda_i$$

same \bar{V} !

Note: $\vec{w}_{OLS} = \bar{C}^{-1} \bar{X}^T \vec{y}$

$$\Rightarrow \vec{w}_{OLS} = \sum_{i=1}^d \frac{1}{\sigma_i} \langle \vec{u}_i, \vec{y} \rangle \vec{v}_i$$

Again, small σ_i get amplified.

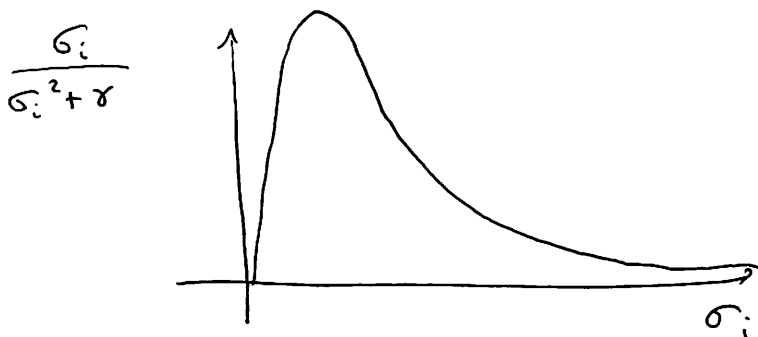
If $\sigma_i \approx 0$ should probably ignore it!

RIDGE REGRESSION: Pick $\gamma > 0$ and define

$$\vec{w}_{RIDGE} = \sum_{i=1}^d \frac{\sigma_i}{\sigma_i^2 + \gamma} \langle \vec{u}_i, \vec{y} \rangle \vec{v}_i$$

If σ_i is very small, $\frac{\sigma_i}{\sigma_i^2 + \gamma}$ is small

If σ_i is very large, $\frac{\sigma_i}{\sigma_i^2 + \gamma} \approx \frac{1}{\sigma_i}$



Underdetermined regression

If \bar{X} is not full rank, ridge still works.

$$\sigma_i = 0 \Rightarrow \frac{\sigma_i}{\sigma_i^2 + \gamma} = 0$$

Ignores directions in $\text{null}(\bar{X})$.

~~At least ridge still works
Result: $\bar{X}^T \bar{X} = 0$
 $\bar{X}^T \bar{X} \bar{X}^T = 0$~~

$$\vec{w}_{\text{PINV}} = \sum_{i=1}^d r_i \langle \vec{u}_i, \vec{y} \rangle \vec{v}_i$$

$$r_i = \begin{cases} \frac{1}{\sigma_i} & \sigma_i > 0 \\ 0 & \sigma_i = 0 \end{cases}$$

Limit of ridge as $\gamma \rightarrow 0$.

same as $\vec{w}_{\text{PINV}} = \bar{X}^T \bar{y}$

pseudoinverse of \bar{X} times \bar{y} .

Ignores $\text{null}(\bar{X})$, ~~works only in range~~.

Ridge regression

Optimization View

Least squares: minimize $\|\bar{X}\vec{w} - \vec{y}\|^2$

if $n < d$ or \bar{X} rank deficient, there are infinitely many solutions.

if ~~\vec{w}_*~~ \vec{w}_* is a minimizer, so is $\vec{w}_* + \vec{u}$ for $\vec{u} \in \text{null}(\bar{X})$.

Consider: minimize $\|\bar{X}\vec{w} - \vec{y}\|^2 + \underbrace{\lambda \|\vec{w}\|^2}_{\text{penalty term on } \vec{w}}$

Write $\vec{w} = \vec{w}_x + \vec{w}_\perp$ where

$$\vec{w}_x \in \text{range}(\bar{X}^T) = \text{span}(\{\vec{X}_i\})$$

$$\vec{w}_\perp \in \text{null}(\bar{X})$$

Then $\|\bar{X}\vec{w} - \vec{y}\| + \lambda \|\vec{w}\|^2$

$$= \|\bar{X}\vec{w}_x - \vec{y}\|^2 + \lambda \|\vec{w}_x\|^2 + \lambda \|\vec{w}_\perp\|^2$$

minimizer sets ~~$\vec{w}_\perp = 0$~~ $\vec{w}_\perp = 0$

Ridge solution:

$$\| \bar{X}(\vec{w}_* + \Delta \vec{w}) - \vec{y} \|^2 + \lambda \|\vec{w}_* + \Delta \vec{w}\|^2 \geq \text{XXXXXXXXXX}$$

$$\| \bar{X} \vec{w}_* - \vec{y} \|^2 + \lambda \|\vec{w}_*\|^2 \quad \forall \Delta \vec{w}$$

$$\Rightarrow 2 \langle \bar{X} \vec{w}_* - \vec{y}, \bar{X} \Delta \vec{w} \rangle + \|\bar{X} \Delta \vec{w}\|^2 + 2\lambda \langle \vec{w}_*, \Delta \vec{w} \rangle + \lambda \|\Delta \vec{w}\|^2 \geq 0 \quad \forall \Delta \vec{w}$$

$$\Rightarrow \bar{X}^T (\bar{X} \vec{w}_* - \vec{y}) + \lambda \vec{w}_* = 0$$

$$\Rightarrow (\bar{X}^T \bar{X} + \lambda \bar{I}) \vec{w}_* = \bar{X}^T \vec{y}$$

$$\Rightarrow \vec{w}_* = (\bar{X}^T \bar{X} + \lambda \bar{I})^{-1} \bar{X}^T \vec{y}$$

\vec{w}_{RIDGE} is unique.

Alternative Solutions

We observed $\vec{w}_{\text{RIDGE}} = \bar{X}^T \vec{\beta}$ for some $\vec{\beta}$

We can optimize for $\vec{\beta}$:

$$\underset{\vec{\beta}}{\text{minimize}} \quad \|\bar{X} \bar{X}^T \vec{\beta} - \vec{y}\|^2 + \lambda \|\bar{X}^T \vec{\beta}\|^2$$

Define $\bar{G} = \bar{X} \bar{X}^T$ $n \times n$ $G_{ij} = \langle x_i, x_j \rangle$

Then $\|\bar{G} \vec{\beta} - \vec{y}\|^2 + \lambda \vec{\beta}^T \bar{G} \vec{\beta}$ is

minimized if

$$\bar{G} (\bar{G} \vec{\beta} - \vec{y} + \lambda \vec{\beta}) = 0$$

$$\vec{\beta} = (\bar{G} + \lambda \bar{I}_n)^{-1} \vec{y}$$

~~coefficient of~~

$$\vec{w}_{\text{RIDGE}}^T \vec{x} = \sum_{i=1}^n \beta_i \langle \vec{x}_i, \vec{x} \rangle$$