

MVE137

Probability and Statistical Learning Using Python

Linear methods for regression

Alexandre Graell i Amat

`alexandre.graell@chalmers.se`

<https://sites.google.com/site/agraellamat>

September 23 and 27, 2022



CHALMERS

Linear regression

Linear regression: Assumes linear model for $f(\mathbf{x})$,

$$\hat{y} = f(\mathbf{x}) = \mathbf{x}^\top \mathbf{w},$$

i.e., $\tilde{f}(\mathbf{x}) = \mathbb{E}_{y|\mathbf{x}}[y|x] \approx \tilde{\mathbf{x}}^\top \mathbf{w}^*$

- **Simple** and **interpretable**
- Can **outperform** non-linear methods when small number of training samples, very noisy data, sparse data
- Can be used to model **nonlinear relationships** using basis functions $\phi(\mathbf{x}) = (\phi_1(\mathbf{x}), \dots, \phi_M(\mathbf{x}))$
- Many nonlinear techniques direct generalizations of linear methods

Linear regression IBM

Linear regression and least squares

We assume:

- A probabilistic model

$$y = \tilde{f}(\mathbf{x}) + \varepsilon$$

with $\varepsilon \sim \mathcal{N}(0, \sigma^2)$

- A linear model

$$f(\mathbf{x}) = \mathbf{x}^\top \mathbf{w},$$

with $\mathbf{w} = (w_0, w_1, \dots, w_p)^\top$ and $\mathbf{x} = (1, x_1, \dots, x_p)^\top$

- A data set $\mathcal{D} = \{(\mathbf{x}_1, y_1), \dots, (\mathbf{x}_N, y_N)\} = (\mathbf{X}_{N \times (d+1)}, \mathbf{y}_{N \times 1})$
- Residual sum of squares criterion:

$$\text{RSS}(\mathbf{w}) = \sum_{i=1}^N (y_i - \mathbf{x}_i^\top \mathbf{w})^2$$

Linear regression and least squares

Optimal solution:

$$\mathbf{w}^* = (\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{X}^\top \mathbf{y}$$

obtained from differentiating the RSS and setting derivative to zero,

$$\mathbf{X}^\top (\mathbf{y} - \mathbf{X}\mathbf{w}) = 0$$

Predicted value corresponding to $\mathbf{x} = (1, x_1, \dots, x_d)$:

$$\hat{y} = \mathbf{x}^\top \mathbf{w}^*$$

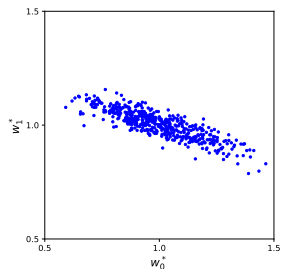
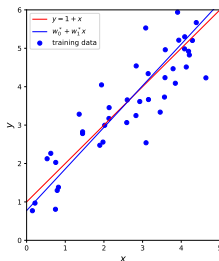
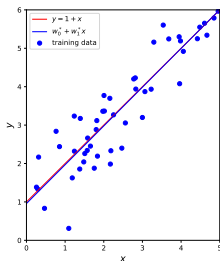
Properties of w^*

w^* depends on the training data $\rightarrow w^*$ is a **random variable**!

Example:

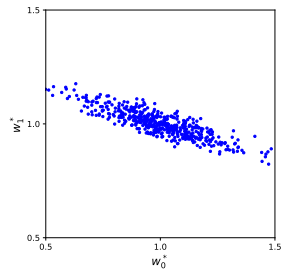
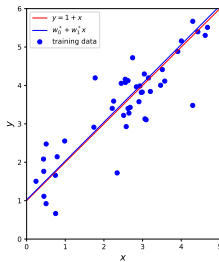
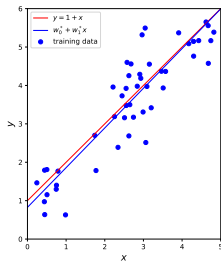
$$y = w_0 + w_1x + \varepsilon,$$

with $w = (1, 1)^T$, i.e., $y = 1 + x + \varepsilon$, and $\varepsilon = \mathcal{N}(0, 0.49)$



Changing training input data

Properties of w^*



Fixed training input data

Distribution of w^*

Assumptions:

- $y = x^\top w + \varepsilon$ (for training set: $y = Xw + \varepsilon$)
- Observations y_i **uncorrelated** and with constant variance σ^2
- X **fixed**

Expectation:

$$\begin{aligned}\mathbb{E}[w^*] &= \mathbb{E}[(X^\top X)^{-1} X^\top y] \\ &= (X^\top X)^{-1} X^\top \mathbb{E}[y]\end{aligned}$$

Fixed X : $\mathbb{E}[y] = Xw$

Thus,

$$\begin{aligned}\mathbb{E}[w^*] &= (X^\top X)^{-1} X^\top Xw \\ &= w\end{aligned}$$

Distribution of w^*

Variance:

$$\begin{aligned}w^* - \mathbb{E}[w^*] &= (X^T X)^{-1} X^T y - w \\&= (X^T X)^{-1} X^T y - (X^T X)^{-1} (X^T X) w \\&= (X^T X)^{-1} X^T (y - X w) \\&= (X^T X)^{-1} X^T \varepsilon\end{aligned}$$

Then:

$$\begin{aligned}\text{Var}[w^*] &= \mathbb{E} [(w^* - \mathbb{E}[w^*])(w^* - \mathbb{E}[w^*])^T] \\&= \mathbb{E} [(X^T X)^{-1} X^T \varepsilon \varepsilon^T X (X^T X)^{-1}] \\&= (X^T X)^{-1} X^T \text{Var}[\varepsilon] X (X^T X)^{-1} \\&= (X^T X)^{-1} X^T (\sigma^2 I_N) X (X^T X)^{-1} \\&= \sigma^2 (X^T X)^{-1} X^T X (X^T X)^{-1} \\&= (X^T X)^{-1} \sigma^2\end{aligned}$$

Distribution of w^*

And σ^2 ? Estimate it by the **sample variance**!

$$\hat{\sigma}^2 = \frac{1}{N} \sum_{i=1}^N (y_i - \hat{y}_i)^2$$

We typically use

$$\hat{\sigma}^2 = \frac{1}{N - d - 1} \sum_{i=1}^N (y_i - \hat{y}_i)^2$$

so that $\mathbb{E}[\hat{\sigma}^2] = \sigma^2$.

Distribution of w^*

Distribution:

From $y = Xw + \varepsilon$ and $w^* = (X^T X)^{-1} X^T y$,

$$\begin{aligned} w^* &= (X^T X)^{-1} X^T (Xw + \varepsilon) \\ &= w + (X^T X)^{-1} X^T \varepsilon \end{aligned}$$

w^* is a linear transformation of a multivariate Gaussian (ε) $\longrightarrow w^*$
multivariate Gaussian!

$$w^* \sim \mathcal{N}(w, (X^T X)^{-1} \sigma^2)$$

Interpretability of the model

Machine learning algorithms often a **black box**

Explainable AI: Understand decisions or predictions made by the AI

Linear regression allows for **interpretability**!

Interpretation of the estimated weights

$$\mathbf{w}^* \sim \mathcal{N}(\mathbf{w}, (\mathbf{X}^\top \mathbf{X})^{-1} \sigma^2)$$

Interpretation of (w_0^*, \dots, w_d^*) estimated by least squares?

- Which w 's are probably **zero**? \longrightarrow associated features are **irrelevant**

If true $w_j = 0$:

$$w_j^* \sim \mathcal{N}(0, \sigma^2 v_j)$$

- If $w_j^* > \sigma^2 v_j$: **highly improbable** $w_j = 0$

Test hypothesis $w_j = 0$ (z -score):

$$z_j = \frac{w_j^*}{\hat{\sigma} \sqrt{v_j}}$$

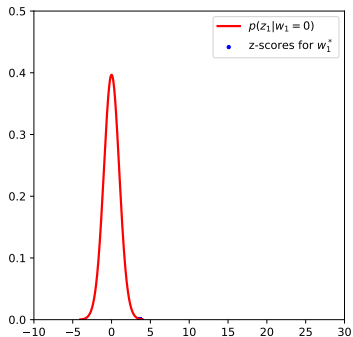
- $w_j = 0$: z_j has a **student's t -distribution** with $N - d - 1$ degrees of freedom

Interpretation of the estimated weights

Example:

$$y = w_0 + w_1x + \varepsilon,$$

with $\mathbf{w} = (1, 1)^\top$, i.e., $y = 1 + x + \varepsilon$, and $\varepsilon = \mathcal{N}(0, 0.49)$



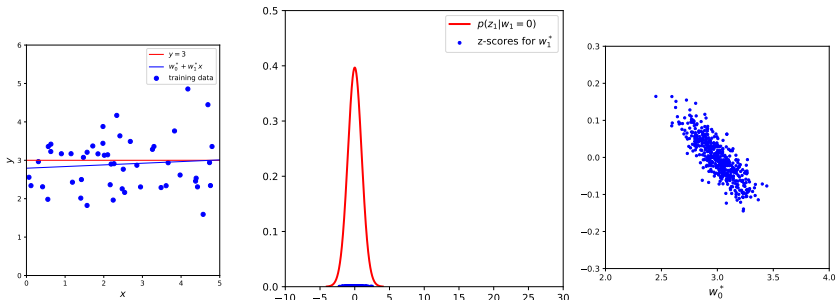
- $N = 50$ training examples
- blue circles: z -scores for w_1^* for 500 training data sets

Interpretation of the estimated weights

Example:

$$y = w_0 + w_1x + \varepsilon,$$

with $\mathbf{w} = (3, 0)^\top$, i.e., $y = 1 + x + \varepsilon$, and $\varepsilon = \mathcal{N}(0, 0.49)$



- $N = 50$ training examples
- $N = 500$ different training sets (right)

The Gauss-Markov theorem

The least squares estimate w^* has the smallest variance among all linear unbiased estimates.

We consider estimation of $\theta = a^T w$. Least squares estimate:

$$\theta^* = a^T w^* = a^T (X^T X)^{-1} X^T y$$

Assuming model $y = x^T w + \varepsilon$ is correct $\rightarrow \mathbb{E}[y|X] = Xw$ and

$$\begin{aligned}\mathbb{E}[a^T w^*] &= \mathbb{E} [a^T (X^T X)^{-1} X^T y] \\ &= a^T (X^T X)^{-1} X^T \mathbb{E}[y] \\ &= a^T (X^T X)^{-1} X^T X w \\ &= a^T w \\ &= \theta\end{aligned}$$

Gauss-Markov theorem. Any other linear estimator $\tilde{\theta} = c^T y$ that is unbiased for $a^T w$ has variance

$$\text{Var}[a^T w^*] \leq \text{Var}[c^T y]$$

The Gauss-Markov theorem: Implications

Mean-squared error of estimator $\tilde{\theta}$ of θ ,

$$\begin{aligned}\text{MSE}(\tilde{\theta}) &= \mathbb{E}[(\tilde{\theta} - \theta)^2] \\&= \mathbb{E}[(\tilde{\theta} - \mathbb{E}[\tilde{\theta}] + \mathbb{E}[\tilde{\theta}] - \theta)^2] \\&= (\tilde{\theta} - \mathbb{E}[\tilde{\theta}])^2 + 2(\tilde{\theta} - \mathbb{E}[\tilde{\theta}])(\mathbb{E}[\tilde{\theta}] - \theta) + (\mathbb{E}[\tilde{\theta}] - \theta)^2 \\&= (\tilde{\theta} - \mathbb{E}[\tilde{\theta}])^2 + (\mathbb{E}[\tilde{\theta}] - \theta)^2 \\&= \underbrace{\text{Var}[\tilde{\theta}]}_{\text{variance}} + \underbrace{(\mathbb{E}[\tilde{\theta}] - \theta)^2}_{\text{bias}}\end{aligned}$$

The Gauss-Markov theorem: Implications

Mean-squared error of estimator $\tilde{\theta}$ of θ ,

$$\text{MSE}(\tilde{\theta}) = \underbrace{\text{Var}[\tilde{\theta}]}_{\text{variance}} + \underbrace{(\mathbb{E}[\tilde{\theta}] - \theta)}_{\text{bias}}^2$$

Gauss-Markov theorem. For all linear estimators with zero bias, the least-squares estimator has the smallest MSE!

Biased estimates may give a **smaller MSE** \longrightarrow **bias-variance trade-off!**

In practice: Any model will be biased \longrightarrow strike the right bias-variance trade-off.

MSE and prediction error

Prediction error:

$$L(\hat{y}) = \sigma^2 + \underbrace{(\tilde{f}(\mathbf{x}) - \mathbb{E}_{\mathcal{D}}[f(\mathbf{x})])^2}_{\text{bias}} + \text{Var}_{\mathcal{D}} [f(\mathbf{x})]$$

With $f(\mathbf{x}) \equiv \tilde{\theta}$ and $\tilde{f}(\mathbf{x}) \equiv \theta$:

$$\begin{aligned} L(\hat{y}) &= \sigma^2 + \underbrace{(\theta - \mathbb{E}[\tilde{\theta}])^2}_{\text{bias}} + \text{Var}_{\mathcal{D}} [\tilde{\theta}] \\ &= \sigma^2 + \text{MSE}(\tilde{\theta}) \\ &= \sigma^2 + \text{MSE}(f(\mathbf{x})) \end{aligned}$$

- σ^2 : independent of model, irreducible
- $\text{MSE}(f(\mathbf{x}))$: error in the model

Minimizing the MSE minimizes the expected prediction error

Multiple outputs

Predict $K > 1$ output variables $\mathbf{y} = (y_1, \dots, y_K)^\top$

Linear model for each output:

$$\begin{aligned}y_j &= w_{j,0} + \sum_{\ell=1}^d x_\ell w_{j,\ell} + \varepsilon_i \\&= w_{j,0} + \mathbf{w}_j^\top \mathbf{x} + \varepsilon_i \\&= f_j(\mathbf{x}) + \varepsilon_i\end{aligned}$$

- $\mathcal{D} = \{(\mathbf{x}_1, \mathbf{y}_1), \dots, (\mathbf{x}_N, \mathbf{y}_N)\}$
- \mathbf{Y} : $N \times K$ matrix of outputs \mathbf{Y} , with \mathbf{y}_i as i -th row

Linear regression:

$$\mathbf{Y} = \mathbf{X}\mathbf{W} + \mathbf{E}$$

\mathbf{W} : $(d+1) \times K$ matrix of coefficients w

Multiple outputs

Single output:

$$\text{RSS}(\mathbf{w}) = \sum_{i=1}^N (y_i - \mathbf{x}_i^T \mathbf{w})^2$$

Multiple outputs:

$$\begin{aligned} \text{RSS}(\mathbf{W}) &= \sum_{i=1}^N \sum_{j=1}^K (y_{i,j} - f_j(\mathbf{x}_i))^2 \\ &= \sum_{i=1}^N \sum_{j=1}^K \left(y_{i,j} - \left(w_{j,0} + \sum_{\ell=1}^d x_{i,\ell} w_{j,\ell} \right) \right)^2 \\ &= \|\mathbf{Y} - \mathbf{XW}\|_{\text{F}}^2 \end{aligned}$$

with

$$\|\mathbf{A}\|_{\text{F}}^2 = \sum_i \sum_j a_{i,j}^2$$

Multiple outputs

Least squares solution:

$$\begin{aligned}\mathbf{W}^* &= \arg \min_{\mathbf{W}} \text{RSS}(\mathbf{W}) \\ &= (\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{X}^\top \mathbf{Y}\end{aligned}$$

and

$$\hat{\mathbf{Y}}^* = \mathbf{X} \mathbf{W}^* = \mathbf{X} (\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{X}^\top \mathbf{Y}$$

Coefficients for i -th output $y_{1,i}, \dots, y_{N,i}$ solution to least squares regression of $y_{1,i}, \dots, y_{N,i}$ on columns of $\mathbf{X} \rightarrow$ multiple outputs do not affect one another's estimates.

Subset selection

Least squares estimates:

- Low bias
- High variance: may penalize prediction accuracy

Accuracy can be improved by setting some w_j to zero.

Subset selection: Determine a subset of $k < d$ features which is the most informative → Increased interpretability, improved accuracy.

Best-subset selection

For each $k \in \{1, \dots, d\}$ find subset of size k that gives **smallest RSS**:

$$\text{RSS}_{\text{BSS}}(j_1, \dots, j_k) = \min_{w_0, w_{j_1}, \dots, w_{j_k}} \sum_{i=1}^N \left(y_i - w_0 - \sum_{\ell=1}^k w_{j_\ell} x_{i, j_\ell} \right)^2$$

with $j_i \in \{1, \dots, d\}$

Observations:

- $\binom{d}{k}$ different subsets
- $d \leq 40$: computationally feasible algorithms
- $d > 40$: infeasible to solve exactly, heuristic algorithms
- How to choose k ? **bias-variance trade-off**

Forward- and backward-stepwise selection

Idea: consider a **greedy** approach

Forward-stepwise selection:

Start with **intercept**, then augment model sequentially by adding **predictor** that **improves fit most**:

$$j_\ell = \arg \min_{j \in \mathcal{I}} \min_{w_0, w_{j_1}, \dots, w_{j_{\ell-1}}, \mathbf{w}_j} \sum_{i=1}^N \left(y_i - w_0 - \sum_{m=1}^{\ell-1} w_{j_m} x_{i,j_m} - \mathbf{w}_j x_{i,j} \right)^2$$

with $\mathcal{I} = \{1, \dots, d\} \setminus \{j_1, \dots, j_{\ell-1}\}$.

Backward-stepwise selection:

Start with **complete model**, and sequentially delete **predictor** that contributes **least to fit** → At each step, we remove variable with **smallest z -score**

Subset selection

Principle: Set some w_j 's to 0, let others be estimated using least squares.

- Interpretable
- Probably lower prediction

Drawbacks:

- **Discrete procedure:** w_j retained or discarded \rightarrow high variance
- **High complexity** (and heuristic solutions suboptimal)

Idea (Regularization): Allow all estimates to be positive, but constrain them to not become too big.

Shrinkage methods: Ridge regression

Idea: Shrink regression coefficients w_j by imposing a penalty on their size.

$$\begin{aligned} \mathbf{w}_{\text{ridge}}^* &= \arg \min_{\mathbf{w}} \underbrace{\sum_{i=1}^N \left(y_i - w_0 - \sum_{j=1}^d x_{i,j} w_j \right)^2}_{\text{RSS}} + \lambda \underbrace{\sum_{j=1}^d w_j^2}_{\text{reg. function}} \\ &= \arg \min_{\mathbf{w}} \|\mathbf{y} - \mathbf{X}\mathbf{w}\|^2 + \lambda \|\tilde{\mathbf{w}}\|^2 \end{aligned}$$

Observations:

- Larger λ : more shrinkage
- $\lambda = \infty$: $\tilde{\mathbf{w}}_{\text{ridge}}^* = 0$
- $\lambda = 0$: conventional linear regression

Shrinkage methods: Ridge regression

Ridge regression in an equivalent form:

$$\begin{aligned} \mathbf{w}_{\text{ridge}}^* = \arg \min_{\mathbf{w}} \sum_{i=1}^N \left(y_i - w_0 - \sum_{j=1}^d x_{i,j} w_j \right)^2 \\ \text{subject to } \|\tilde{\mathbf{w}}\|^2 \leq t \end{aligned}$$

Shrinkage methods: Ridge regression

Ridge regression in another equivalent form:

$$\mathbf{w}_c^* = \arg \min_{\mathbf{w}_c} \sum_{i=1}^N \left(y_i - w_{c,0} - \sum_{j=1}^d \tilde{x}_{i,j} w_{c,j} \right)^2 + \lambda \sum_{j=1}^d w_{c,j}^2,$$

with

$$\tilde{x}_{i,j} = x_{i,j} - \bar{x}_j = x_{i,j} - \frac{1}{N} \sum_{\ell=1}^N x_{\ell,j}$$

Optimization can be done in two steps

1. Fit $w_{c,0}$ separately,

$$w_{c,0}^* = \bar{y} = \frac{1}{N} \sum_{i=1}^N y_i$$

2. Optimize all other $w_{c,j}$ by ridge regression without intercept

Shrinkage methods: Ridge regression

Assuming **centered inputs** $\tilde{x}_{i,j} = x_{i,j} - \bar{x}_j$ and $\tilde{y}_i = y_i - \bar{y}$, $w_{c,j}$, $j = 1, \dots, d$, obtained solving

$$\begin{aligned} \mathbf{w}_c^* &= \arg \min_{\mathbf{w}_c} \sum_{i=1}^N \left(y_i - \sum_{j=1}^d x_{i,j} w_{c,j} \right)^2 + \lambda \sum_{j=1}^d w_{c,j}^2 \\ &= \arg \min_{\mathbf{w}_c} (\mathbf{y} - \mathbf{X} \mathbf{w}_c)^\top (\mathbf{y} - \mathbf{X} \mathbf{w}_c) + \lambda \|\mathbf{w}_c\|^2 \end{aligned}$$

with $\mathbf{y} = (y_1, \dots, y_N)^\top$, $\mathbf{w}_c = (w_{c,1}, \dots, w_{c,d})^\top$, and

$$\mathbf{X} = \begin{pmatrix} x_{1,1} & x_{1,2} & \dots & x_{1,d} \\ x_{2,1} & x_{2,2} & \dots & x_{2,d} \\ \vdots & \vdots & \dots & \vdots \\ x_{N,1} & x_{N,2} & \dots & x_{N,d} \end{pmatrix}$$

$$\mathbf{w}_{\text{ridge}}^* = (\mathbf{X}^\top \mathbf{X} + \lambda \mathbf{I})^{-1} \mathbf{X}^\top \mathbf{y}$$

Shrinkage methods: Ridge regression

We assume **centered** input matrix \mathbf{X} (**no intercept**)

Singular value decomposition of \mathbf{X} :

$$\mathbf{X} = \mathbf{U}\mathbf{D}\mathbf{V}^\top$$

\mathbf{U} : $N \times d$ **semi-orthogonal** matrix; columns span column space of \mathbf{X}

\mathbf{V} : $d \times d$ **orthogonal** matrix

\mathbf{D} : $d \times d$ **diagonal** matrix with diagonal entries $d_1 \geq d_2 \geq \dots \geq d_d$ **singular values** of \mathbf{X}

Can write **least squares** fitted vector as:

$$\begin{aligned}\mathbf{X}\mathbf{w}^{\text{ls}} &= \mathbf{X}(\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{X}^\top \mathbf{y} \\ &= \mathbf{U}\mathbf{U}^\top \mathbf{y} \\ &= \sum_{i=1}^d \mathbf{u}_i (\mathbf{u}_i^\top \mathbf{y})\end{aligned}$$

Shrinkage methods: Ridge regression

Least squares solution:

$$\begin{aligned}\hat{\mathbf{y}} &= \mathbf{X}\mathbf{w}^{\text{ls}} = \mathbf{U}\mathbf{U}^{\top}\mathbf{y} \\ &= \sum_{i=1}^d \mathbf{u}_i(\mathbf{u}_i^{\top}\mathbf{y})\end{aligned}$$

Observation:

- $\mathbf{U}^{\top}\mathbf{y}$: coordinates of \mathbf{y} with respect to $\mathbf{U} \rightarrow$ least squares solution
 $\mathbf{U}\mathbf{U}^{\top}\mathbf{y}$ is closest approximation to \mathbf{y} in subspace spanned by columns of \mathbf{U}

Shrinkage methods: Ridge regression

Least squares solution:

$$\hat{\mathbf{y}} = \mathbf{X}\mathbf{w}^{\text{ls}} = \mathbf{U}\mathbf{U}^{\text{T}}\mathbf{y} = \sum_{i=1}^d \mathbf{u}_i (\mathbf{u}_i^{\text{T}}\mathbf{y})$$

Ridge regression solution:

$$\begin{aligned}\hat{\mathbf{y}}_{\text{ridge}} &= \mathbf{X}\mathbf{w}_{\text{ridge}}^* \\ &= \mathbf{U}\mathbf{D}\mathbf{V}^{\text{T}}\mathbf{V}(\mathbf{D}^2 + \lambda\mathbf{I})^{-1}\mathbf{D}\mathbf{U}^{\text{T}}\mathbf{y} \\ &= \mathbf{U}\mathbf{D}(\mathbf{D}^2 + \lambda\mathbf{I})^{-1}\mathbf{D}\mathbf{U}^{\text{T}}\mathbf{y}\end{aligned}$$

Observations:

- $\mathbf{D}(\mathbf{D}^2 + \lambda\mathbf{I})^{-1}\mathbf{D}$: diagonal matrix with elements $\frac{d_i^2}{d_i^2 + \lambda}$
- $\mathbf{U}^{\text{T}}\mathbf{y}$: coordinates of vector \mathbf{y} in the basis spanned by \mathbf{U}

$$\hat{\mathbf{y}}_{\text{ridge}} = \sum_{i=1}^d \mathbf{u}_i \frac{d_i^2}{d_i^2 + \lambda} \mathbf{u}_i^{\text{T}}\mathbf{y}$$

Principal components and meaning of the d_i^2 's

- A collection of points in a \mathbb{R}^d (d -dimensional vectors)
- Want to **summarize** them by projecting them onto a **q -dimensional subspace**

Most informative summary: directions $\mathbf{b} \in \mathbb{R}^d$ with **large spread** of the data

Principal components: q orthogonal directions in space along which projections of the original vectors have **largest variance**.

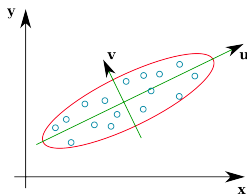
Principal components: q orthogonal directions that **minimize average** (mean-squared) **distance** between original vectors and their projections.

Principal components and meaning of the d_i^2 's

Principal components: Sequence of q (unit) vectors, with i -th vector the direction of a line that **best fits data** while being **orthogonal** to first $i - 1$ vectors

Best-fitting line \equiv one that **minimizes average squared distance** from points to line \equiv **maximizes variance**

1. 1st PC: direction in space along which projections have largest variance
2. 2nd PC: direction maximizing variance among all directions orthogonal to 1st PC
3. k -th PC: variance-maximizing direction orthogonal to previous $k - 1$ components



Principal components and meaning of the d_i^2 's

How do we find the principal components?

- Direction with **largest variation**: Eigenvector with **largest eigenvalue**
- **Second** direction with **largest variation**: Eigenvector with **second largest eigenvalue**
- ...

Back to our **regression problem** (with centered \mathbf{X}):

Sample covariance matrix:

$$\Sigma = \frac{1}{N} \mathbf{X}^\top \mathbf{X} = \frac{1}{N} \mathbf{V} \mathbf{D}^2 \mathbf{V}^\top$$

The **eigenvectors** \mathbf{v}_i (columns of \mathbf{V}) are the **principal components** of \mathbf{X} !

Diagonal entries of \mathbf{D}^2 , $d_1^2 \geq d_2^2 \geq \dots \geq d_d^2$, are the **eigenvalues** of $\mathbf{X} \mathbf{X}^\top$

Principal components and meaning of the d_i^2 's

Project \mathbf{x}_i 's onto first principal component (\mathbf{v}_1):

$$z_i^{(1)} = \mathbf{v}_1^\top \mathbf{x}_i$$

Variance of $z_i^{(1)}$:

$$\begin{aligned} \frac{1}{N} \sum_{i=1}^N \left(z_i^{(1)} \right)^2 &= \frac{1}{N} \sum_{i=1}^N \mathbf{v}_1^\top \mathbf{x}_i \mathbf{x}_i^\top \mathbf{v}_1 \\ &= \frac{1}{N} \mathbf{v}_1^\top \mathbf{X}^\top \mathbf{X} \mathbf{v}_1 \\ &= \frac{1}{N} \mathbf{v}_1^\top \mathbf{V} \mathbf{D}^2 \mathbf{V}^\top \mathbf{v}_1 \\ &= \frac{d_1^2}{N} \end{aligned}$$

We have: variance $z_i^{(1)} > \text{variance } z_i^{(2)}, \dots \text{variance } z_i^{(j)} > \text{variance } z_i^{(j+1)}$

Small d_j 's correspond to directions having small variance. Ridge regression shrinks these directions the most!

Shrinkage methods: The lasso

Idea: Shrink regression coefficients w_j by imposing a penalty on their size.

$$\begin{aligned} \mathbf{w}_{\text{lasso}}^* &= \arg \min_{\mathbf{w}} \underbrace{\sum_{i=1}^N \left(y_i - w_0 - \sum_{j=1}^d x_{i,j} w_j \right)^2}_{\text{RSS}} + \underbrace{\lambda \sum_{j=1}^d |w_j|}_{\text{reg. function}} \\ &= \arg \min_{\mathbf{w}} \|\mathbf{y} - \mathbf{X}\mathbf{w}\|^2 + \lambda \|\tilde{\mathbf{w}}\|_1 \end{aligned}$$

Equivalently,

$$\begin{aligned} \mathbf{w}_{\text{lasso}}^* &= \arg \min_{\mathbf{w}} \sum_{i=1}^N \left(y_i - w_0 - \sum_{j=1}^d x_{i,j} w_j \right)^2 \\ &\text{subject to } \|\tilde{\mathbf{w}}\|_1 \leq t \end{aligned}$$

Shrinkage methods: The lasso

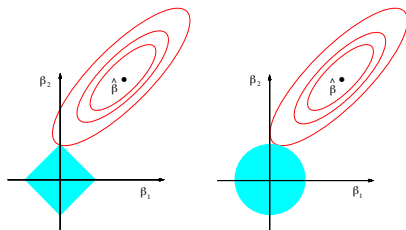
$$\begin{aligned} \mathbf{w}_{\text{lasso}}^* &= \arg \min_{\mathbf{w}} \underbrace{\sum_{i=1}^N \left(y_i - w_0 - \sum_{j=1}^d x_{i,j} w_j \right)^2}_{\text{RSS}} + \underbrace{\lambda \sum_{j=1}^d |w_j|}_{\text{reg. function}} \\ &= \arg \min_{\mathbf{w}} \|\mathbf{y} - \mathbf{X}\mathbf{w}\|^2 + \lambda \|\tilde{\mathbf{w}}\|_1 \end{aligned}$$

Observations:

- Solution not a simple linear function of the y_i 's
- No closed-form solution
- Computing the Lasso solution a quadratic programming problem
- We can standardize predictors: intercept can be fitted separately, then other parameters via lasso without intercept

Property: For λ sufficiently large, some coefficients w_j driven to zero \rightarrow can be interpreted as **continuous subset selection**

Ridge regression and the lasso: Graphical representation



- Elliptical contours: pairs (w_1, w_2) that yield a given RSS
- Ridge (right): $\|w\|^2 \leq t \rightarrow w_1^2 + w_2^2 \leq t$
- Lasso (left): $\|w\|_1 \leq t \rightarrow |w_1| + |w_2| \leq t$

Solution lies at intersection between RSS contours and constrained function.

- $\lambda = 0$: least-squares solution
- $\lambda = \infty$: $(0, 0)$
- Lasso: Intersection may occur in a vertex \rightarrow many zero coefficients

Generalization: Regularization of least-squares

$$w_{\text{lasso}}^* = \arg \min_w \underbrace{\sum_{i=1}^N \left(y_i - w_0 - \sum_{j=1}^d x_{i,j} w_j \right)^2}_{\text{RSS}} + \underbrace{\lambda \sum_{j=1}^d |w_j|^q}_{\text{reg. function}}$$

- $q = 1$: Lasso
- $q = 2$: Ridge regression
- $q = 0$: Subset selection

Regularization: allows complex models to be trained on data sets of limited size without severe overfitting.