

RIDGE REGULARIZATION: AN ESSENTIAL CONCEPT IN DATA SCIENCE

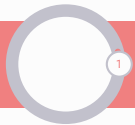
BASED ON THE ARTICLE OF TREVOR HASTIE 2020

04-2021

Bascou Florent
Lefort Tanguy

University of Montpellier





Ridge computational cost

Kernel trick

Data augmentation

Dropout regularization

Double descent

Rank selection for matrix

We denote $X \in \mathbb{R}^{n \times p}$, $y \in \mathbb{R}^n$ and $\beta \in \mathbb{R}^p$ such that $y \simeq X\beta$.

Ordinary least squares:

$$\hat{\beta} = \arg \min_{\beta} \|y - X\beta\|_2^2 \iff \hat{\beta} = (X^T X)^{-1} X^T y,$$

Problems that can happen:

- ▶ $X^T X$ may be ill conditioned ($\kappa = \frac{\text{largest singular value}}{\text{smallest singular value}} \gg 1$)
- ▶ $p > n$ leads to infinite number of solutions for OLS.

¹Tikhonov (1943); Hoerl and Kennard (1970)

We denote $X \in \mathbb{R}^{n \times p}$, $y \in \mathbb{R}^n$ and $\beta \in \mathbb{R}^p$ such that $y \simeq X\beta$.

Ordinary least squares:

$$\hat{\beta} = \arg \min_{\beta} \|y - X\beta\|_2^2 \iff \hat{\beta} = (X^T X)^{-1} X^T y,$$

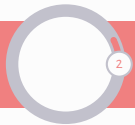
Problems that can happen:

- ▶ $X^T X$ may be ill conditioned ($\kappa = \frac{\text{largest singular value}}{\text{smallest singular value}} \gg 1$)
 - ▶ Shift spectrum by λ using $X^T X + \lambda \text{Id}$.
- ▶ $p > n$ leads to infinite number of solutions for OLS.
 - ▶ Add penalty to recover unicity.

¹Tikhonov (1943); Hoerl and Kennard (1970)

RIDGE¹ REGULARIZATION

WHY USE IT?



We denote $X \in \mathbb{R}^{n \times p}$, $y \in \mathbb{R}^n$ and $\beta \in \mathbb{R}^p$ such that $y \simeq X\beta$.

Ordinary least squares:

$$\hat{\beta} = \arg \min_{\beta} \|y - X\beta\|_2^2 \iff \hat{\beta} = (X^T X)^{-1} X^T y,$$

Problems that can happen:

- ▶ $X^T X$ may be ill conditioned ($\kappa = \frac{\text{largest singular value}}{\text{smallest singular value}} \gg 1$)
 - ▶ Shift spectrum by λ using $X^T X + \lambda \text{Id}$.
- ▶ $p > n$ leads to infinite number of solutions for OLS.
 - ▶ Add penalty to recover unicity.

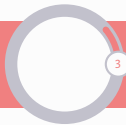
Ridge estimator

$$\hat{\beta}_{\text{ridge}} = \arg \min_{\beta} \|y - X\beta\|_2^2 + \lambda \|\beta\|_2^2 \iff \hat{\beta}_{\text{ridge}} = (X^T X + \lambda \text{Id})^{-1} X^T y.$$

¹Tikhonov (1943); Hoerl and Kennard (1970)

RIDGE REGULARIZATION

SOME GENERALITIES

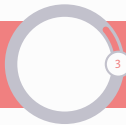


- ▶ Coefficients with smaller ℓ_2 norm,
- ▶ Handles the multicollinearity issue,
- ▶ Don't apply only to linear models, in general:

$$\arg \min_w f(w) + \lambda \|w\|_2^2 .$$

RIDGE REGULARIZATION

SOME GENERALITIES



- ▶ Coefficients with smaller ℓ_2 norm,
- ▶ Handles the multicollinearity issue,
- ▶ Don't apply only to linear models, in general:

$$\arg \min_w f(w) + \lambda \|w\|_2^2 .$$

Some cons (because there must be some)

- ▶ Alone won't give sparse solutions, . . .
- ▶ **BUT** can be combined with a LASSO² which gives the Elastic-Net³:

$$\arg \min_{\beta} \|y - X\beta\|_2^2 + \lambda_1 \|\beta\|_1 + \lambda_2 \|\beta\|_2^2 .$$

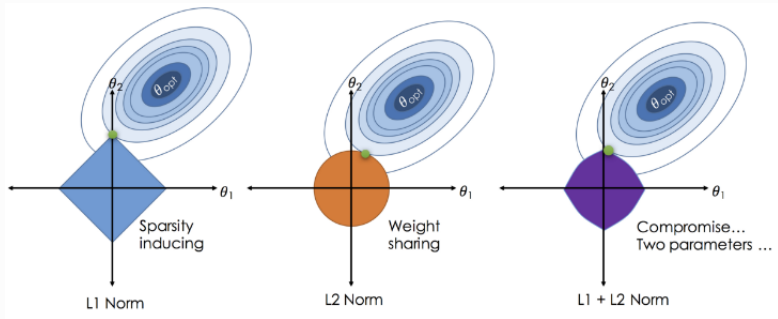
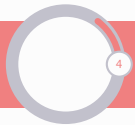
- ▶ Introduce an hyperparameter that needs tuning.

²Tibshirani (1996)

³Zou and Hastie (2005)

DIFFERENT PENALTIES

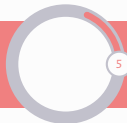
AND DIFFERENT SHAPES



Source: <https://towardsdatascience.com/>

RIDGE ESTIMATOR

HOW TO COMPUTE IT EFFICIENTLY

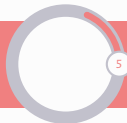


Solution is given by : $\hat{\beta}_\lambda = (X^T X + \lambda \text{Id})^{-1} X^T y$

Problem : λ is a tuning parameter \implies computing many $\hat{\beta}_\lambda$

RIDGE ESTIMATOR

HOW TO COMPUTE IT EFFICIENTLY



Solution is given by : $\hat{\beta}_\lambda = (X^T X + \lambda \text{Id})^{-1} X^T y$

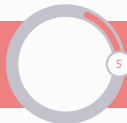
Problem : λ is a tuning parameter \implies computing many $\hat{\beta}_\lambda$

Use ONE SVD ($X = UDV^T$) to compute many estimations

$$\hat{\beta}_\lambda = (X^T X + \lambda \text{Id})^{-1} X^T y = V(D^T D + \lambda \text{Id})^{-1} D^T U^T y = \sum_{j=1}^{\text{rg}(X)} v_j \frac{d_j}{d_j^2 + \lambda} \langle u_j | y \rangle$$

RIDGE ESTIMATOR

HOW TO COMPUTE IT EFFICIENTLY



Solution is given by : $\hat{\beta}_\lambda = (X^T X + \lambda \text{Id})^{-1} X^T y$

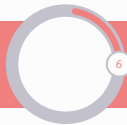
Problem : λ is a tuning parameter \implies computing many $\hat{\beta}_\lambda$

Use ONE SVD ($X = UDV^T$) to compute many estimations

$$\hat{\beta}_\lambda = (X^T X + \lambda \text{Id})^{-1} X^T y = V(D^T D + \lambda \text{Id})^{-1} D^T U^T y = \sum_{j=1}^{\text{rg}(X)} v_j \frac{d_j}{d_j^2 + \lambda} \langle u_j | y \rangle$$

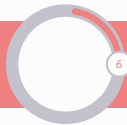
Use ONE SVD ($X = UDV^T$) to compute many predictions

$$\hat{y}_\lambda = X \hat{\beta}_\lambda = U D (D^T D + \lambda \text{Id})^{-1} D^T U^T y = \sum_{j=1}^{\text{rg}(X)} u_j \frac{d_j^2}{d_j^2 + \lambda} \langle u_j | y \rangle$$



Let us denote :

- ▶ $\hat{\beta}_\lambda^{(-i)}$ the estimated coefficients without using the pair (x_i, y_i) .
- ▶ $R^\lambda = X(X^\top X + \lambda \text{Id})^{-1} X^\top$ the Ridge operator matrix

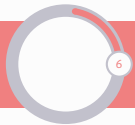


Let us denote :

- ▶ $\hat{\beta}_\lambda^{(-i)}$ the estimated coefficients without using the pair (x_i, y_i) .
- ▶ $R^\lambda = X(X^\top X + \lambda \text{Id})^{-1} X^\top$ the Ridge operator matrix

Easy computation for LOO-CV

$$\text{LOO}_\lambda = \sum_{i=1}^n (y_i - x_i^\top \hat{\beta}_\lambda^{(-i)})^2 = \sum_{i=1}^n \frac{(y_i - x_i^\top \hat{\beta}_\lambda)^2}{(1 - R_{ii}^\lambda)^2}$$



Let us denote :

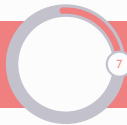
- ▶ $\hat{\beta}_\lambda^{(-i)}$ the estimated coefficients without using the pair (x_i, y_i) .
- ▶ $R^\lambda = X(X^\top X + \lambda \text{Id})^{-1} X^\top$ the Ridge operator matrix

Easy computation for LOO-CV

$$\text{LOO}_\lambda = \sum_{i=1}^n (y_i - \mathbf{x}_i^\top \hat{\beta}_\lambda^{(-i)})^2 = \sum_{i=1}^n \frac{(y_i - \mathbf{x}_i^\top \hat{\beta}_\lambda)^2}{(1 - R_{ii}^\lambda)^2},$$

$$R^\lambda = X(X^\top X + \lambda \text{Id})^{-1} X^\top = U(D^\top D + \lambda \text{Id})^{-1} D^\top U = US(\lambda)U.$$

with $S(\lambda)$ the diagonal shrinkage matrix with elements $\frac{d_j^2}{d_j^2 + \lambda}$.



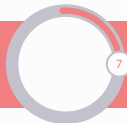
Suppose that the data arises from a linear model with *i.i.d* centered errors ε_i

$$y_i = \mathbf{x}_i^\top \beta + \varepsilon_i, \quad i = 1, \dots, n$$

Then, the Ridge estimate $\hat{\beta}_\lambda$ is a biased estimate of β .

RIDGE ESTIMATOR

BIAS AND VARIANCE



Suppose that the data arises from a linear model with *i.i.d* centered errors ε_i

$$y_i = \mathbf{x}_i^\top \beta + \varepsilon_i, \quad i = 1, \dots, n$$

Then, the Ridge estimate $\hat{\beta}_\lambda$ is a biased estimate of β .

Bias - Covariance matrix

If the x_i are assumed fixed, $n > p$ and X has full column rank, we get :

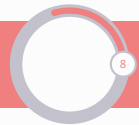
$$\text{Bias}(\hat{\beta}_\lambda) = \sum_{j=1}^p v_j \frac{\lambda}{d_j^2 + \lambda} \langle v_j | \beta \rangle,$$

$$\text{Var}(\hat{\beta}_\lambda) = \sigma^2 \sum_{j=1}^p \frac{d_j^2}{(d_j^2 + \lambda)^2} v_j v_j^\top \quad \text{with } \sigma^2 = \text{Var}(\varepsilon_i).$$

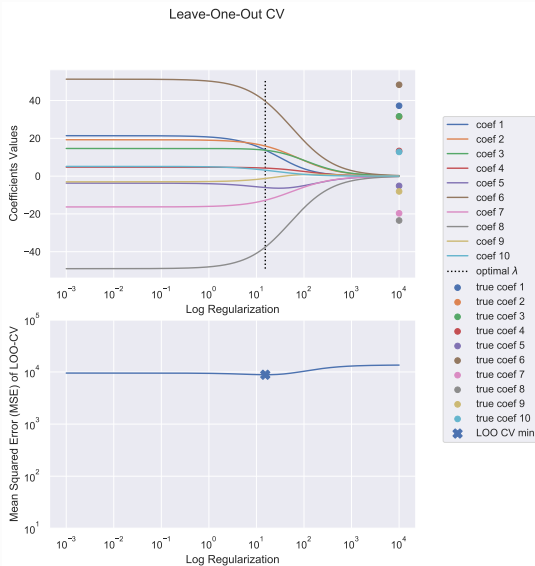
The smaller the j^{th} singular value is, the bigger the shrinkage associated is.

RIDGE ESTIMATOR

AN EXAMPLE OF LOO CV AND BIAS

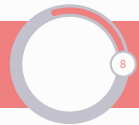


- ▶ Isotropic data:
 $X \sim \mathcal{N}(0, \text{Id})$,
- ▶ $n = 50, p = 10$
- ▶ **SNR = 1**
- ▶ $y = X\beta^* + \varepsilon$ with
 $\varepsilon \sim \mathcal{N}(0, \sigma^2 \text{Id})$

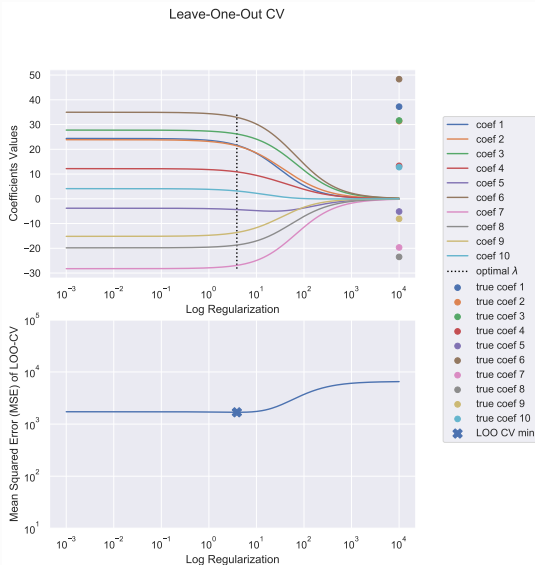


RIDGE ESTIMATOR

AN EXAMPLE OF LOO CV AND BIAS

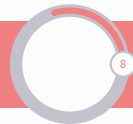


- ▶ Isotropic data:
 $X \sim \mathcal{N}(0, \text{Id})$,
- ▶ $n = 50, p = 10$
- ▶ **SNR = 2**
- ▶ $y = X\beta^* + \varepsilon$ with
 $\varepsilon \sim \mathcal{N}(0, \sigma^2 \text{Id})$

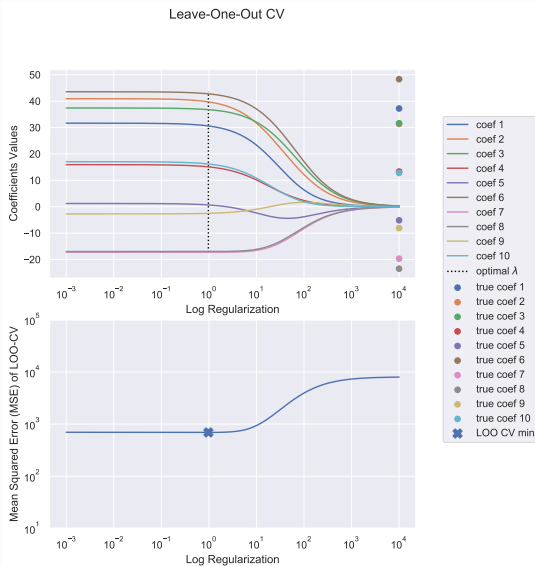


RIDGE ESTIMATOR

AN EXAMPLE OF LOO CV AND BIAS

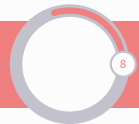


- ▶ Isotropic data:
 $X \sim \mathcal{N}(0, \text{Id})$,
- ▶ $n = 50, p = 10$
- ▶ **SNR = 3**
- ▶ $y = X\beta^* + \varepsilon$ with
 $\varepsilon \sim \mathcal{N}(0, \sigma^2 \text{Id})$

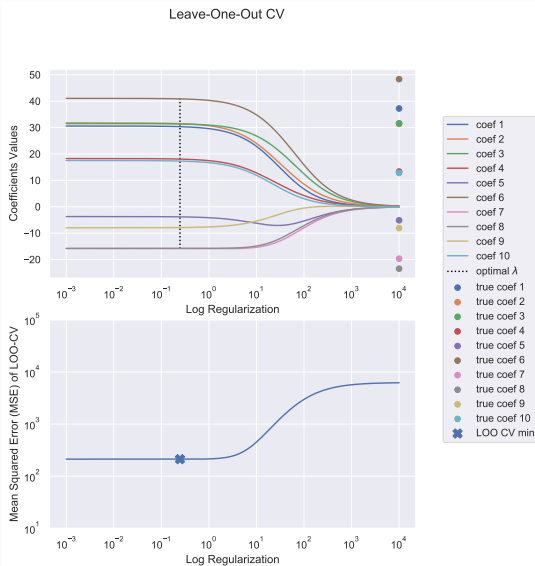


RIDGE ESTIMATOR

AN EXAMPLE OF LOO CV AND BIAS

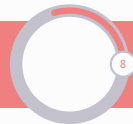


- ▶ Isotropic data:
 $X \sim \mathcal{N}(0, \text{Id})$,
- ▶ $n = 50, p = 10$
- ▶ **SNR = 5**
- ▶ $y = X\beta^* + \varepsilon$ with
 $\varepsilon \sim \mathcal{N}(0, \sigma^2 \text{Id})$

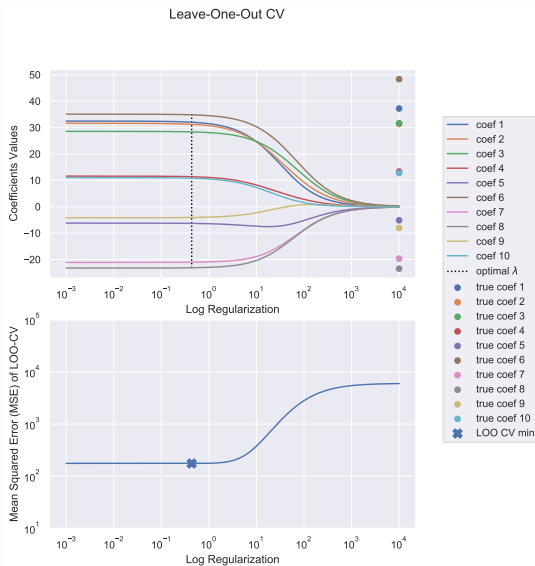


RIDGE ESTIMATOR

AN EXAMPLE OF LOO CV AND BIAS

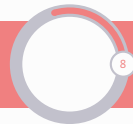


- ▶ Isotropic data:
 $X \sim \mathcal{N}(0, \text{Id})$,
- ▶ $n = 50, p = 10$
- ▶ **SNR = 7**
- ▶ $y = X\beta^* + \varepsilon$ with
 $\varepsilon \sim \mathcal{N}(0, \sigma^2 \text{Id})$

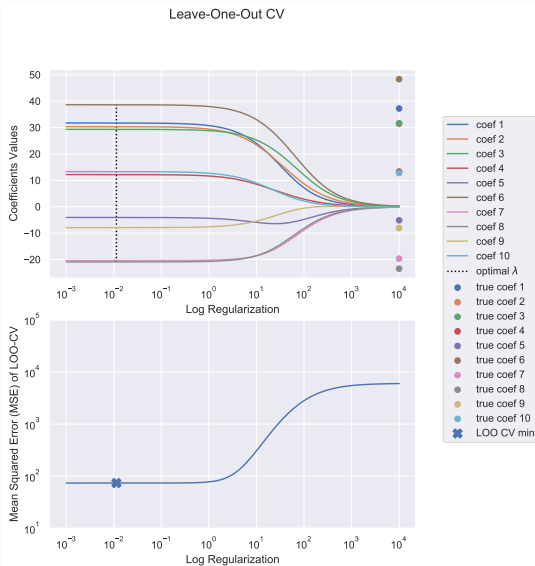


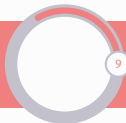
RIDGE ESTIMATOR

AN EXAMPLE OF LOO CV AND BIAS



- ▶ Isotropic data:
 $X \sim \mathcal{N}(0, \text{Id})$,
- ▶ $n = 50, p = 10$
- ▶ **SNR = 10**
- ▶ $y = X\beta^* + \varepsilon$ with
 $\varepsilon \sim \mathcal{N}(0, \sigma^2 \text{Id})$





Reduce the complexity

If $p > n$, $(X^T X + \lambda \text{Id})^{-1} \in \mathbb{R}^{p \times p}$ is costly. But we can actually only solve a $n \times n$ system thanks to the relation (*proof with SVD*):

$$X^T (X X^T + \lambda \text{Id})^{-1} y = (X^T X + \lambda \text{Id})^{-1} X^T y .$$

Reduce the complexity

If $p > n$, $(X^T X + \lambda \text{Id})^{-1} \in \mathbb{R}^{p \times p}$ is costly. But we can actually only solve a $n \times n$ system thanks to the relation (*proof with SVD*):

$$X^T (X X^T + \lambda \text{Id})^{-1} y = (X^T X + \lambda \text{Id})^{-1} X^T y .$$

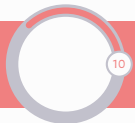
So we can write $\hat{\beta} = X^T u$ with $u \in \mathbb{R}^n$. Denote $K = X X^T$ the **Gram** matrix,

$$\begin{aligned} \hat{y} &= X \hat{\beta} = X X^T u \\ &= K (K + \lambda \text{Id})^{-1} y \end{aligned}$$

- ▶ New ridge problem smaller to solve.
- ▶ **Opens the door for a lot more!**

KERNEL RIDGE REGRESSION

NON LINEAR RELATION



Suppose $y_i = \varphi(X_i)$, then $\hat{\varphi} = \arg \min_{\varphi \in \mathcal{H}} \|y - \varphi(X)\|_2^2 + \lambda \|\varphi\|_{\mathcal{H}}^2$.

Representer theorem⁴

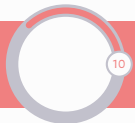
Take \mathcal{H} RKHS of kernel K of map Φ s.t. $K(x, y) = \langle \Phi(x), \Phi(y) \rangle = (\Phi(x))(y)$.
For $f \in \mathcal{H}$, $f(x) = \langle f, \Phi(x) \rangle$. Then,

$$\varphi = \sum_{i=1}^n \alpha_i K(x_i, \bullet) \ .$$

⁴Schölkopf et al. (2001)

KERNEL RIDGE REGRESSION

NON LINEAR RELATION



Suppose $y_i = \varphi(X_i)$, then $\hat{\varphi} = \arg \min_{\varphi \in \mathcal{H}} \|y - \varphi(X)\|_2^2 + \lambda \|\varphi\|_{\mathcal{H}}^2$.

Representer theorem⁴

Take \mathcal{H} RKHS of kernel K of map Φ s.t. $K(x, y) = \langle \Phi(x), \Phi(y) \rangle = (\Phi(x))(y)$.
For $f \in \mathcal{H}$, $f(x) = \langle f, \Phi(x) \rangle$. Then,

$$\varphi = \sum_{i=1}^n \alpha_i K(x_i, \bullet) .$$

Denote $\mathbf{K} = (\mathbf{K}_{ij})_{ij} = (K(x_i, x_j))_{ij}$. The problem is now

$$\hat{\alpha} = \arg \min_{\alpha} \|y - \mathbf{K}\alpha\|_2^2 + \lambda \alpha^\top \mathbf{K}\alpha .$$

- ▶ first order conditions: $\nabla = (\mathbf{K}\mathbf{K} + \lambda\mathbf{K})\hat{\alpha} - \mathbf{K}y = 0$,
- ▶ solution:

$$\hat{\alpha} = (\mathbf{K} + \lambda\text{Id})^{-1}y .$$

⁴Schölkopf et al. (2001)

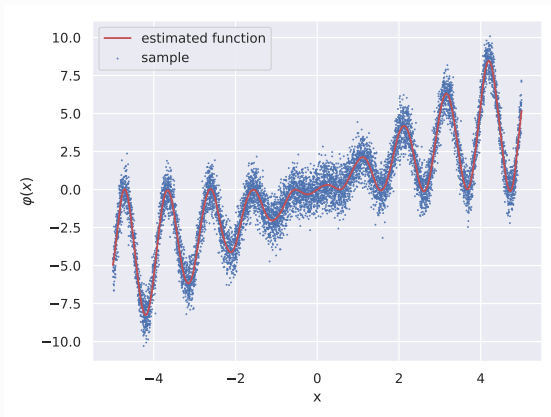
KERNEL RIDGE REGRESSION

EXAMPLE WITH KeOps PACKAGE



Gaussian kernel to estimate $\varphi(t) = t + t \cos(6t)$ on noised data ($\sigma = 0.8$),

$$K(x, y) = \exp \left\{ -\gamma \|x - y\|_2^2 \right\}, \gamma = \frac{1}{2 \cdot 0.2^2} .$$



Example of a small dog from CIFAR10 dataset



Example of a small dog from CIFAR10 dataset



⁵Chollet and Allaire (2018)

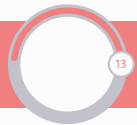
Example of a small dog from CIFAR10 dataset



⁵Chollet and Allaire (2018)

DATA AUGMENTATION

AND WITH CLOUD POINTS?

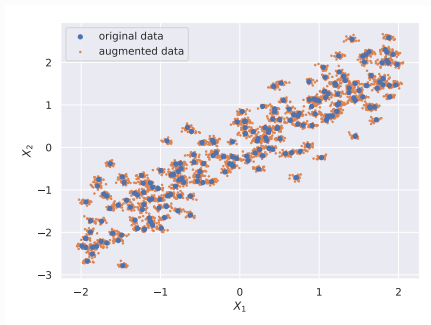


- ▶ **Goal:** Use the data we have to create new points,
- ▶ **But** can't flip it / make a small rotation, . . .

One way to do so: create perturbed points from the data we have:

$$\mathbf{x}_{ij} = \mathbf{x}_i + \varepsilon_{ij}, \varepsilon_{ij} \sim \mathcal{N}\left(0, \frac{\lambda}{n} \text{Id}\right), i \in [n], j \in [m],$$

with associated response $y_{ij} = y_i$.

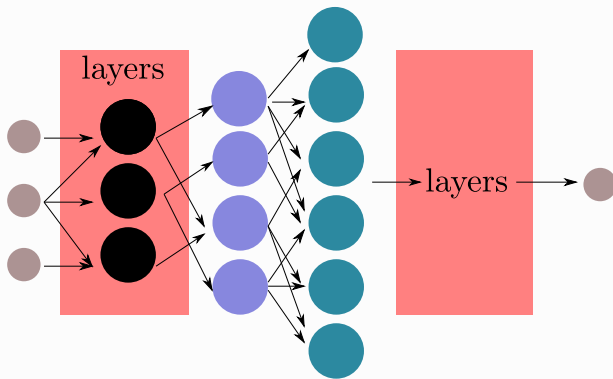


Added points compensate each other

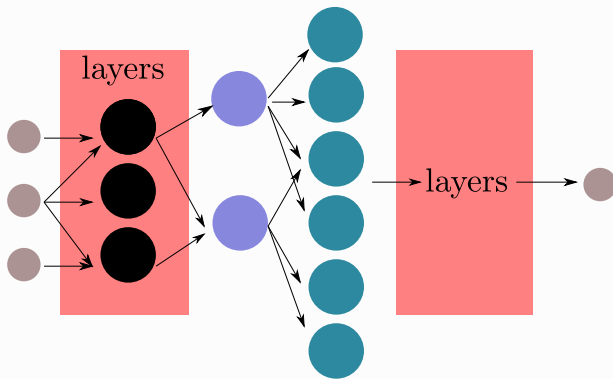
$$\sum_i \frac{1}{m} \sum_j \mathbf{x}_{ij} \mathbf{x}_{ij}^T \simeq (\mathbf{X}\mathbf{X}^T + \lambda \text{Id})$$

OLS with $\mathbf{X}^{\text{augmented}} \simeq$ Ridge with \mathbf{X}

- ▶ Randomly set units of a layer to 0 with probability ϕ to avoid overfitting,
- ▶ Inflate surviving ones by $1/(1 - \phi)$ factor as compensation.



- ▶ Randomly set units of a layer to 0 with probability ϕ to avoid overfitting,
- ▶ Inflate surviving ones by $1/(1 - \phi)$ factor as compensation.



Denoting $(l_{ij})_{ij}$ the dropout mask on $X = (x_{ij})_{ij}$ then cost function is:

$$L(\beta) = \frac{1}{2} \sum_{i=1}^n \left(y_i - \sum_{j=1}^p x_{ij} l_{ij} \beta_j \right)^2 .$$

Note that $\mathbb{E}[l_{ij}] = 0\phi + \frac{1}{1-\phi}(1-\phi) = 1$, thus:

$$\mathbb{E} \left[\frac{\partial L(\beta)}{\partial \beta} \right] = \mathbf{X}\mathbf{X}^T \beta - \mathbf{X}^T \mathbf{y} + \frac{\phi}{1-\phi} \text{diag}(\|\mathbf{x}_j\|^2)_{j=1}^p \beta ,$$

⁷Wager et al. (2013)

Denoting $(I_{ij})_{ij}$ the dropout mask on $X = (x_{ij})_{ij}$ then cost function is:

$$L(\beta) = \frac{1}{2} \sum_{i=1}^n \left(y_i - \sum_{j=1}^p x_{ij} I_{ij} \beta_j \right)^2 .$$

Note that $\mathbb{E}[I_{ij}] = 0\phi + \frac{1}{1-\phi}(1-\phi) = 1$, thus:

$$\mathbb{E} \left[\frac{\partial L(\beta)}{\partial \beta} \right] = \mathbf{X}\mathbf{X}^T \beta - \mathbf{X}^T \mathbf{y} + \frac{\phi}{1-\phi} \text{diag}(\|x_j\|^2)_{j=1}^p \beta ,$$

Link with ridge

Solving first order conditions:

$$\hat{\beta}^{\text{dropout}} = \left(\mathbf{X}^T \mathbf{X} + \frac{\phi}{1-\phi} \text{diag}(\|x_j\|^2) \right)^{-1} \mathbf{X}^T \mathbf{y} .$$

Normalizing out data leads **in average** to the ridge estimator for $\lambda = \frac{\phi}{1-\phi}$.

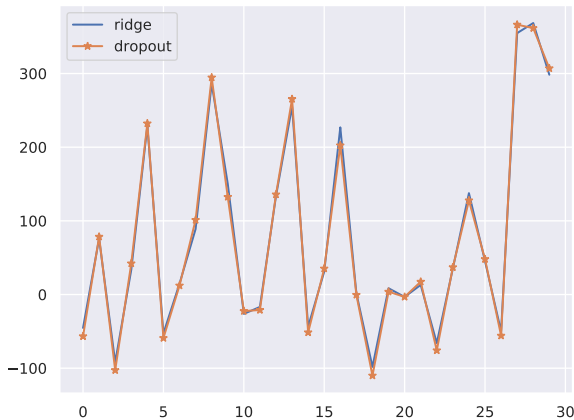
⁷Wager et al. (2013)

No hassle to implement a dropout layer in a Neural Network with PyTorch with this inflation rate:

```
1 class MyModel(nn.Module):
2     def __init__(self, phi):
3         super(MyModel, self).__init__()
4         # all your great layers
5         self.drop = nn.Dropout(phi)
6
7     def forward(self, x):
8         # forward x until the desired layer
9         out = self.drop(current_x)
10        # forward out until the output layer
11        return out
```

DROPOUT EXPERIMENT

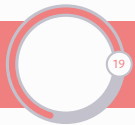
With $\phi = 0.5$ (generally $0.3 \leq \phi \leq 0.5$) and 5 repetitions:



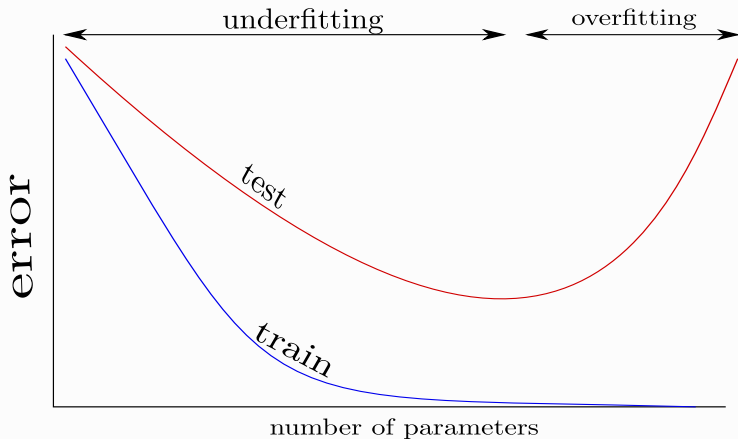
Coefficients of β for the ridge and dropout method with $n = 80$, $p = 30$ and ridge penalty $\lambda = \phi/(1 - \phi) = 1$ are close.

DOUBLE DESCENT WITH (NAKKIRAN ET AL., 2020)

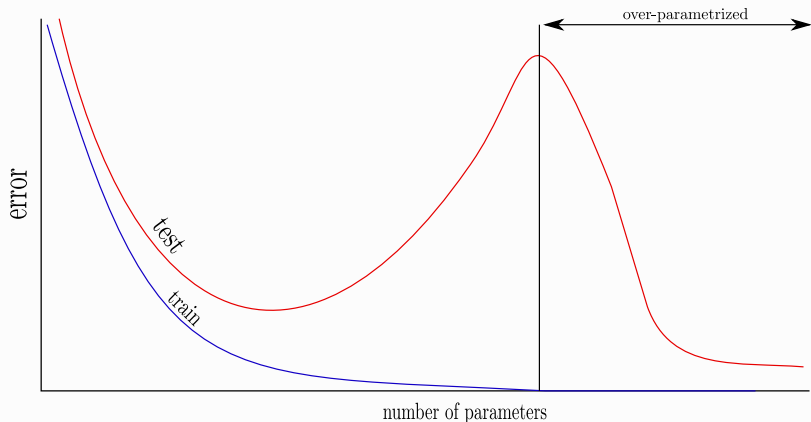
OVER PARAMETRIZATION



Phenomenon observed in Deep Learning, Random Forests ...



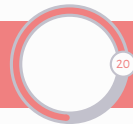
Phenomenon observed in Deep Learning, Random Forests ...



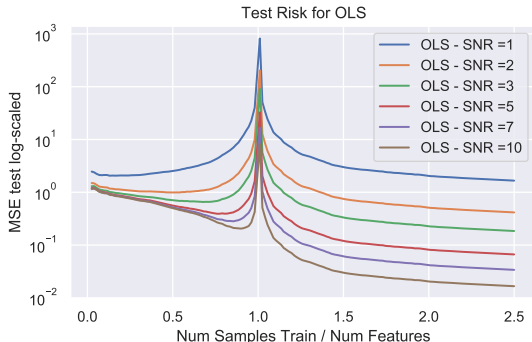
- ▶ interpolate the data $\implies l_2$ norm of $\hat{\beta}$ is high,
- ▶ smaller l_2 norm for $\hat{\beta}$ generalizes better and we keep zero training error.

DOUBLE DESCENT

EXAMPLE WITH SAMPLES (NAKKIRAN ET AL., 2020) - 25 REPETITIONS



- ▶ Isotropic data: $X \sim \mathcal{N}(0, \text{Id})$,
- ▶ $n = 1000, p = 200,$
 $\|\beta^*\|_2 = 1$
- ▶ $y = X\beta^* + \varepsilon$ with
 $\varepsilon \sim \mathcal{N}(0, \sigma^2 \text{Id}),$



Ordinary Least Square produces a double descent phenomenon when training set has same size as the number of features.

DOUBLE DESCENT

EXAMPLE WITH FEATURES (NAKKIRAN ET AL., 2020) - 25 REPETITIONS

▶ To generate data :

▶ Isotropic data:

$$X \sim \mathcal{N}(0, \text{Id}),$$

▶ $n = p = 100, \|\beta^*\|_2 = 1$

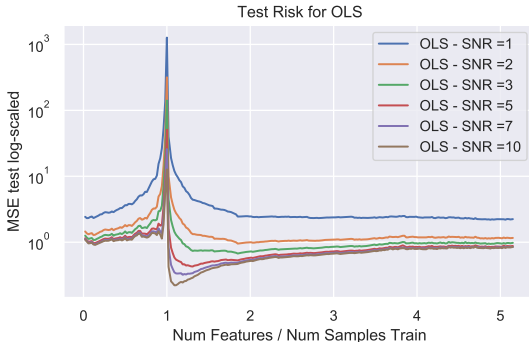
▶ $y = X\beta^* + \varepsilon$ with

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

▶ To learn model :

▶ if $p \leq n$, use X from data ;

▶ if $p > n$, merge columns to X from random distribution.



Ordinary Least Square produces a double descent phenomenon when the number of features is the same as the number of samples.

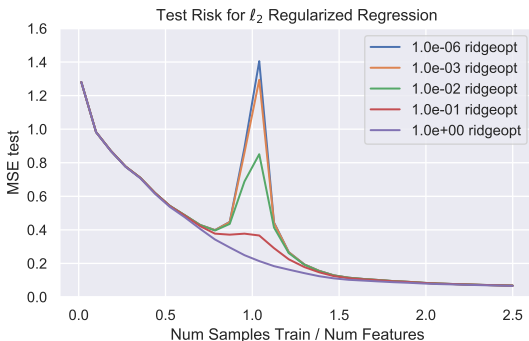
DOUBLE DESCENT

WITH SAMPLES (NAKKIRAN ET AL., 2020) - 25 REPETITIONS

22

In some cases ridge regularization can help get a monotonous error curve.

- ▶ Isotropic data: $X \sim \mathcal{N}(0, \text{Id})$,
- ▶ $n = 1000, p = 200,$
 $\|\beta^*\|_2 = 1$
- ▶ $y = X\beta^* + \varepsilon$ with
 $\varepsilon \sim \mathcal{N}(0, \sigma^2 \text{Id})$,
- ▶ $\lambda_{opt} = \sigma^2 p$

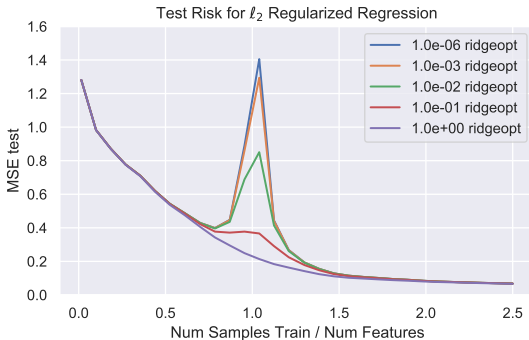


DOUBLE DESCENT

WITH SAMPLES (NAKKIRAN ET AL., 2020) - 25 REPETITIONS

In some cases ridge regularization can help get a monotonous error curve.

- ▶ Isotropic data: $X \sim \mathcal{N}(0, \text{Id})$,
- ▶ $n = 1000, p = 200,$
 $\|\beta^*\|_2 = 1$
- ▶ $y = X\beta^* + \varepsilon$ with
 $\varepsilon \sim \mathcal{N}(0, \sigma^2 \text{Id})$,
- ▶ $\lambda_{\text{opt}} = \sigma^2 p$



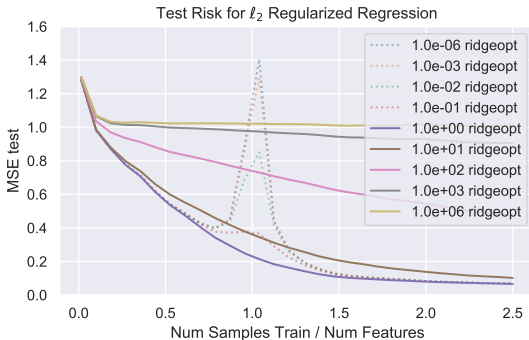
We can achieve monotonic test error decrease with ridge regularization varying p or n for the linear models with isotropic covariates.

DOUBLE DESCENT

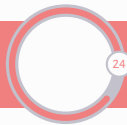
WITH SAMPLES (NAKKIRAN ET AL., 2020) - 25 REPETITIONS

In some cases ridge regularization can help get a monotonous error curve.

- ▶ Isotropic data: $X \sim \mathcal{N}(0, \text{Id})$,
- ▶ $n = 1000, p = 200,$
 $\|\beta^*\|_2 = 1$
- ▶ $y = X\beta^* + \varepsilon$ with
 $\varepsilon \sim \mathcal{N}(0, \sigma^2 \text{Id})$,
- ▶ $\lambda_{opt} = \sigma^2 p$



This property remains true if $\lambda \geq \lambda_{opt}$ varying n for the linear models with isotropic covariates.



Suppose $X \in \mathbb{R}^{m \times n}$. Then we get the rank selection for matrix problem :

$$\min_M \|X - M\|_F^2 \quad \text{s. t.} \quad \text{rank}(M) \leq q.$$

Suppose $X \in \mathbb{R}^{m \times n}$. Then we get the rank selection for matrix problem :

$$\min_M \|X - M\|_F^2 \quad \text{s. t.} \quad \text{rank}(M) \leq q.$$

- ▶ Solution: set all the singular values of D to zero, except the q largest.
- ▶ Problem: the rank constraint makes this problem non-convex

Suppose $X \in \mathbb{R}^{m \times n}$. Then we get the rank selection for matrix problem :

$$\min_M \|X - M\|_F^2 \quad \text{s. t.} \quad \text{rank}(M) \leq q.$$

- ▶ Solution: set all the singular values of D to zero, except the q largest.
- ▶ Problem: the rank constraint makes this problem non-convex

Convex relaxation ⁸- LASSO version of rank selection matrix

$$\tilde{M} = \arg \min_M \|X - M\|_F^2 + \lambda \|M\|_*$$

with $\|M\|_*$ denoting the nuclear norm - sum of singular values.

⁸Fazel (2002)

Suppose $X \in \mathbb{R}^{m \times n}$. Then we get the rank selection for matrix problem :

$$\min_M \|X - M\|_F^2 \quad \text{s. t.} \quad \text{rank}(M) \leq q.$$

- ▶ Solution: set all the singular values of D to zero, except the q largest.
- ▶ Problem: the rank constraint makes this problem non-convex

Convex relaxation ⁸- LASSO version of rank selection matrix

$$\tilde{M} = \arg \min_M \|X - M\|_F^2 + \lambda \|M\|_*$$

with $\|M\|_*$ denoting the nuclear norm - sum of singular values.

- ▶ Solution: soft-tresholding the singular values: $\max(d_i - \lambda, 0)$

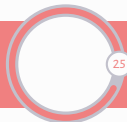
⁸Fazel (2002)

Let $X \in \mathbb{R}^{m \times n}$, $A \in \mathbb{R}^{m \times q}$ and $B \in \mathbb{R}^{n \times q}$. We get the double ridge problem :

$$\tilde{A}, \tilde{B} = \arg \min_{A, B} \|X - AB^T\|_F^2 + \lambda \|A\|_F^2 + \lambda \|B\|_F^2$$

⁹Srebro et al. (2005)

DOUBLE RIDGE PROBLEM⁹



Let $X \in \mathbb{R}^{m \times n}$, $A \in \mathbb{R}^{m \times q}$ and $B \in \mathbb{R}^{n \times q}$. We get the double ridge problem :

$$\tilde{A}, \tilde{B} = \arg \min_{A, B} \|X - AB^T\|_F^2 + \lambda \|A\|_F^2 + \lambda \|B\|_F^2$$

- ▶ This is a ℓ_2 bi-convex problem.
- ▶ The solution is the same as the ℓ_1 convex relation problem: $\tilde{A}\tilde{B}^T = \tilde{M}$

Interest of this property - when SVD fail

If X is massive and sparse, useful to compute a low-rank matrix approximation by alternating ridge regression.

Given A , we get B by:

$$B^T = \left(A^T A + \frac{\lambda}{2} \text{Id}_q \right)^{-1} A^T X$$

⁹Srebro et al. (2005)

The Ridge regularization is linked to several modern techniques, sometimes hidden behind them.

- ▶ LASSO methods are powerful, but **SO ARE RIDGE'S**,
- ▶ it is easy to implement,
- ▶ has computational speed-ups from theoretical results.

https://github.com/tanglef/ml_mtp

The Ridge regularization is linked to several modern techniques, sometimes hidden behind them.

- ▶ LASSO methods are powerful, but **SO ARE RIDGE'S**,
- ▶ it is easy to implement,
- ▶ has computational speed-ups from theoretical results.

One or the other?

The best of both worlds can be used with Elastic-Net regularization.

https://github.com/tanglef/ml_mtp

- François Chollet and Joseph J Allaire. Deep learning with r, ch. 5.4, 2018.
- M. Fazel. *Matrix rank minimization with applications*. PhD thesis, Stanford University, 2002.
- A. E. Hoerl and R. W. Kennard. Ridge regression: Biased estimation for nonorthogonal problems. *Technometrics*, 12(1):55–67, 1970.
- Preetum Nakkiran, Prayaag Venkat, Sham Kakade, and Tengyu Ma. Optimal regularization can mitigate double descent, 2020.
- Bernhard Schölkopf, Ralf Herbrich, and Alex J Smola. A generalized representer theorem. In *International conference on computational learning theory*, pages 416–426. Springer, 2001.
- Nathan Srebro, Jason Rennie, and Tommi S Jaakkola. Maximum-margin matrix factorization. In *Advances in neural information processing systems*, pages 1329–1336, 2005.
- Nitish Srivastava, Geoffrey Hinton, Alex Krizhevsky, Ilya Sutskever, and Ruslan Salakhutdinov. Dropout: a simple way to prevent neural networks from overfitting. *The journal of machine learning research*, 15(1):1929–1958, 2014.
- R. Tibshirani. Regression shrinkage and selection via the lasso. *J. R. Stat. Soc. Ser. B Stat. Methodol.*, 58(1):267–288, 1996.
- A. N. Tikhonov. On the stability of inverse problems. *Dokl. Akad. Nauk SSSR*, 39:176–179, 1943.
- Stefan Wager, Sida Wang, and Percy Liang. Dropout training as adaptive regularization. *arXiv preprint arXiv:1307.1493*, 2013.
- H. Zou and T. J. Hastie. Regularization and variable selection via the elastic net. *J. R. Stat. Soc. Ser. B Stat. Methodol.*, 67(2):301–320, 2005.