

Day 6, 09/25/2018

statistical regularization: ridge regression

$$y = X\delta + \epsilon, \text{Var}(\epsilon) = \sigma^2 I$$

Recall: Normal equations $\hat{\delta} = \arg \min_{\delta} (Y - X\delta)^T(Y - X\delta)$

Analytical solution: $\frac{d J(\delta)}{d \delta} = \frac{d}{d \delta} (Y - X\delta)^T(Y - X\delta) \stackrel{\text{set } \frac{d}{d \delta}}{=} 0$

$$X^T X \delta - X^T Y = 0 \Rightarrow \hat{\delta}_{LS} = (X^T X)^{-1} X^T Y$$

Properties of $\hat{\delta}$: $E\{\hat{\delta}_{LS}\} = E\{(X^T X)^{-1} X^T (X\delta + \epsilon)\} = \delta$

$$\text{Var}\{\hat{\delta}_{LS}\} = \text{Var}\{(X^T X)^{-1} X^T (X\delta + \epsilon)\} = \sigma^2 (X^T X)^{-1} X^T I X (X^T X)^{-1} = \sigma^2 (X^T X)^{-2}$$

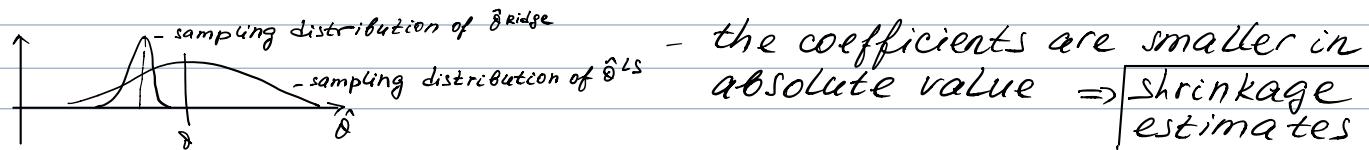
I.e.: $\hat{\delta}$ is a linear combination of y , is unbiased; $\text{Var} = \sigma^2 (X^T X)^{-2}$

- can show that has min variance among all linear unbiased estimators (Gauss-Markov theorem)

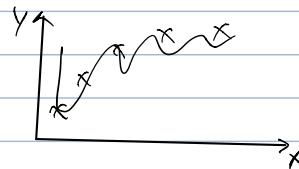
Can improve the estimation when $X^T X$ is (nearly) singular:
 define $\hat{\delta}_{Ridge} = (X^T X + \lambda I)^{-1} X^T y$
 - biased, but smaller variance
 - invert pre-defined constant $\lambda > 0$

$$E\{\hat{\delta}_{Ridge}\} = E\{(X^T X + \lambda I)^{-1} X^T (X\delta + \epsilon)\} = (X^T X + \lambda I)^{-1} X^T X \delta + 0 \neq \hat{\delta}_{LS}$$

$$\text{Var}\{\hat{\delta}_{Ridge}\} = \text{Var}\{(X^T X + \lambda I)^{-1} X^T (X\delta + \epsilon)\} = (X^T X + \lambda I)^{-1} X^T X (X^T X + \lambda I)^{-1} \underset{\text{symmetric}}{\leq} \hat{\delta}_{LS}$$



Intuition:
 polynomial regression



The same solution can be obtained using
 (1) penalized estimation
 (2) constrained optimization

(1) Penalized estimation:

$$\hat{\delta} = \arg \min_{\delta} (Y - X\delta)^T(Y - X\delta) + \lambda \delta^T \delta$$

New loss function $J_{\lambda}(\delta)$

- minimize SSE
 - penalize large δ_j

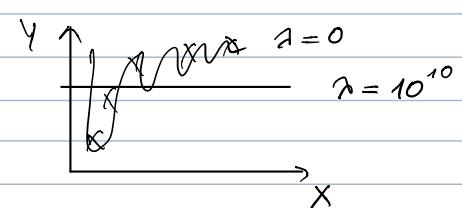
$$= \|\delta\|_2^2 = \sum_{j=1}^p \delta_j^2$$

Solution as function of λ :

$$\begin{aligned} \frac{d}{d \lambda} J_{\lambda}(\delta) &= \frac{d}{d \lambda} \left[Y^T Y - (X\delta)^T Y - Y^T X\delta + (X\delta)^T (X\delta) + \lambda \delta^T \delta \right] \\ &= -2X^T Y + 2X^T X\delta + 2\lambda \delta \stackrel{\text{set } \frac{d}{d \lambda} = 0}{=} 0 \end{aligned}$$

$$\hat{\delta}_{Ridge} = (X^T X + \lambda I)^{-1} X^T Y$$

λ is a hyperparameter



(2) Constrained optimization (see CB Appendix E for reference)

Goal: find the stationary point of a function subject to a constraint "Budget"

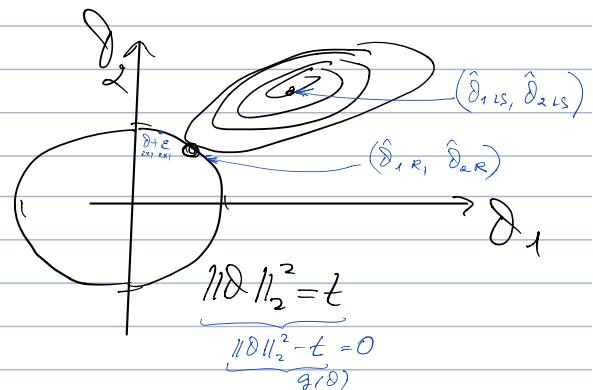
$\min_{\delta} J(\delta)$, such that $\|\delta\|_2^2 \leq t$, equivalently, $\|\delta\|_2^2 - t = 0$

At any point on the constraint surface $g(\delta)$,
 $\nabla g(\delta) \perp g(\delta)$ (i.e., gradient is \perp to surface)

From Taylor approximation around δ :

$$g(\delta + \varepsilon) \approx g(\delta) + (\delta + \varepsilon - \delta)^T \nabla g(\delta) = 0$$

↓ another point on the surface ↓ on the surface ↓ on the surface → for small ε
 $\nabla g(\delta)$ is orthogonal to the surface



Next, look for point $\hat{\delta}$ on the surface, to minimize $J(\delta)$)

At $\hat{\delta}$, $\nabla J(\delta)$ is also orthogonal to the surface (i.e., cannot decrease $J(\delta)$ by moving along the surface)

$\Rightarrow \nabla J(\delta)$ and $\nabla g(\delta)$ are (anti)parallel vectors \Rightarrow there exists λ s.t.

$$\nabla J(\delta) + \lambda \nabla g(\delta) = 0$$

→ Lagrange multiplier, $\lambda \neq 0$ (any sign)

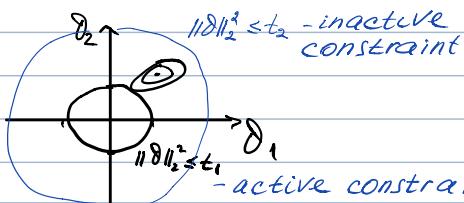
Lagrangian function: $L(\delta, \lambda) = J(\delta) + \lambda g(\delta)$

Stationary point: $\frac{dL(\delta, \lambda)}{d\delta} = \nabla J(\delta) + \lambda \nabla g(\delta) \stackrel{\text{set to 0}}{=} 0$ and $\frac{dL(\delta, \lambda)}{d\lambda} = g(\delta) \stackrel{\text{set to 0}}{=} 0$

$\frac{dL(\delta, \lambda)}{d\lambda} = g(\delta) = 0$ checks the constraint

Extension to inexact (i.e., inequality) constraints:

$\min_{\delta} J(\delta)$, such that $\|\delta\|_2^2 \leq t$, or $\|\delta\|_2^2 - t \leq 0$



Inactive constraint:

- stationary point solves $\nabla J(\delta) = 0$ (i.e., $\lambda = 0$)

Active constraint:

- $L(\delta, \lambda) = J(\delta) + \lambda g(\delta)$

sign is important
gradient is oriented away from $g(\delta) \leq 0$ (anti-parallel vectors)

$$\nabla J(\delta) = -\lambda \nabla g(\delta) \text{ for } \lambda > 0$$

\Rightarrow Solution of $\min_{\delta} J(\delta)$ s.t. $\|\delta\|_2^2 - t \leq 0$ is obtained by

minimizing $L(\delta, \lambda)$ subject to constraints

constraint: $\|\delta\|_2^2 - t \leq 0$

direction: $\lambda \geq 0$

active / $\lambda(\|\delta\|_2^2 - t) = 0$

inactive / $= 0$
if inactive at the optimum if active

} Karush-Kuhn-Tucker (KKT) conditions

i.e., optimize first; then check that the solution satisfies the KKT conditions

Connection between (1) and (2):

$$(1) \min_{\delta} J(\delta) + \lambda \|\delta\|_2^2$$

$$(2) \min_{\delta} J(\delta), \text{ s.t. } \|\delta\|_2^2 - t \leq 0 \quad \rightarrow \quad L(\delta, \lambda) = J(\delta) + \lambda (\|\delta\|_2^2 - t)$$

rename λ to α for the purpose of comparing two approaches
different from (1)

Conditions for optimum:

$$(1) \nabla [J(\delta) + \lambda \|\delta\|_2^2] = 0$$

$$(2) \nabla [J(\delta) + \lambda (\|\delta\|_2^2 - t)] = 0 \quad \text{and} \quad \lambda (\|\delta\|_2^2 - t) = 0$$

Let's say that solution of (1) is $\hat{\delta}(t)$

Then (2) has the same solution wrt δ if $t = \|\hat{\delta}(t)\|_2^2$

(i.e., λ and t are connected but not the same)

$\nabla [J(\delta) + \lambda \|\delta\|_2^2]$ and $\nabla [J(\delta) + \lambda (\|\delta\|_2^2 - t)]$ are the same wrt $\delta \Rightarrow$ same optimum $\hat{\delta}(t)$

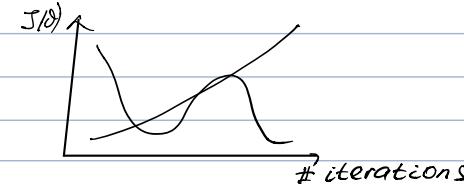
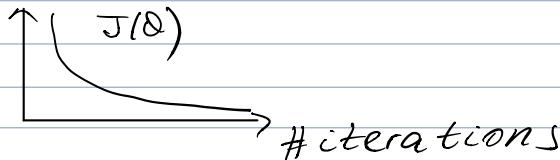
- if $t = \|\hat{\delta}(t)\|_2^2$, then the KKT condition verifies

Parameter estimate by gradient descent: new loss $J(\delta) + \lambda \|\delta\|_2^2$

replace the old loss function

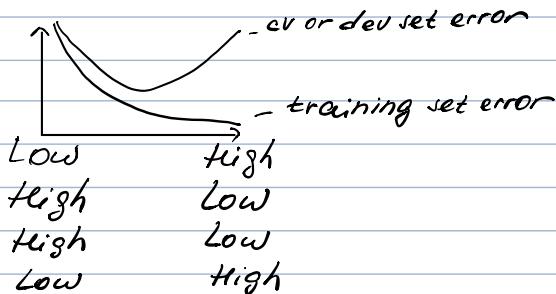
Practical considerations:

- scaling features (parameter values are comparable)
- centering features (no intercept)
- trace plots of $J(\delta)$



Q: how to select λ ? Try a range of λ ; use dev set or cross-validation to find the optimal point

model complexity (#params)



$N \uparrow \rightarrow$ fix high variance

$P \downarrow \rightarrow$ -"

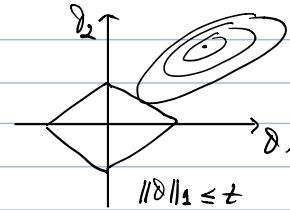
$P \uparrow \rightarrow$ fix high bias

$\lambda \downarrow \rightarrow$ fix high bias

$\lambda \uparrow \rightarrow$ fix high variance

Modifications to penalized least squares:

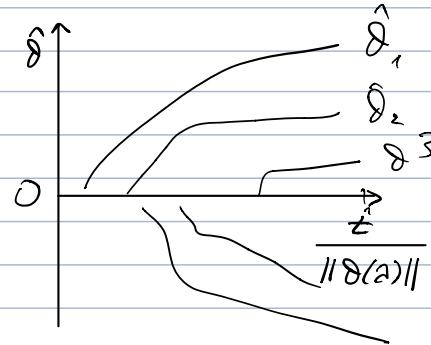
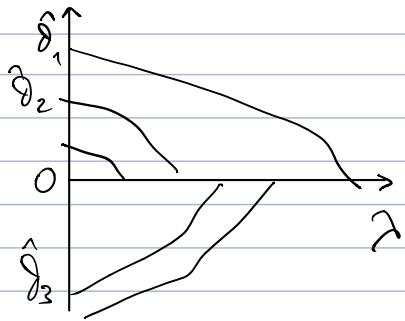
Lasso: $\hat{\delta} = \underset{\delta}{\operatorname{argmin}} J(\delta) + \lambda \sum_{j=1}^p |\delta_j|$



Elastic net: $\hat{\delta} = \underset{\delta}{\operatorname{argmin}} J(\delta) + \lambda \sum_{j=1}^p [\alpha \delta_j^2 + (1-\alpha)|\delta_j|]$

- Solutions are non-linear in Y, no closed-form expression
- Efficient algos exist
- Penalized optimization, while setting small $\hat{\delta}_j = 0$

Standardize predictors to have mean 0 and sd 1



Option for methods/ implementation project :

implement Lasso

(HT ALGO 3.2 and 3.2a)

