

Machine Learning

Regularization and Feature Selection

Fabio Vandin

November 5, 2020

Learning Model

- A : learning algorithm for a machine learning task
- S : m i.i.d. pairs $z_i = (x_i, y_i)$, $i = 1, \dots, m$, with $z_i \in Z = \mathcal{X} \times Y$, generated from distribution $\mathcal{D} \Rightarrow$ training set available to A to produce $A(S)$;
- \mathcal{H} : the hypothesis (or model) set for A
- loss function: $\ell(h, (x, y))$, $\ell : \mathcal{H} \times Z \rightarrow \mathbb{R}^+$
- $L_S(h)$: empirical risk or training error of hypothesis $h \in \mathcal{H}$

$$L_S(h) = \frac{1}{m} \sum_{i=1}^m \ell(h, z_i)$$

- $L_{\mathcal{D}}(h)$: true risk or generalization error of hypothesis $h \in \mathcal{H}$:

$$L_{\mathcal{D}}(h) = \mathbb{E}_{z \in \mathcal{D}}[\ell(h, z)]$$

Learning Paradigms

We would like A to produce $A(S)$ such that $L_{\mathcal{D}}(A(S))$ is *small*, or at least close to the smallest generalization error $L_{\mathcal{D}}(h^*)$ achievable by the “best” hypothesis h^* in \mathcal{H} :

$$h^* = \arg \min_{h \in \mathcal{H}} L_{\mathcal{D}}(h)$$

We have seen a *learning paradigms*: Empirical Risk Minimization

We will now see another learning paradigm...

Regularized Loss Minimization

Assume h is defined by a vector $w = (w_1, \dots, w_d)^T \in \mathbb{R}^d$ (e.g., linear models)

Regularization function $R : \mathbb{R}^d \rightarrow \mathbb{R}$

Regularized Loss Minimization (RLM): pick h obtained as

$$\arg \min_w (L_S(w) + R(w))$$

Intuition: $R(w)$ is a “measure of complexity” of hypothesis h defined by w

\Rightarrow regularization balances between low empirical risk and “less complex” hypotheses

We will see some of the most common regularization function

ℓ_1 Regularization

Regularization function: $R(\mathbf{w}) = \lambda \|\mathbf{w}\|_1$

- $\lambda \in \mathbb{R}, \lambda > 0$
- ℓ_1 norm: $\|\mathbf{w}\|_1 = \sum_{i=1}^d |w_i|$

Therefore the *learning rule* is: pick

$$A(S) = \arg \min_{\mathbf{w}} (L_S(\mathbf{w}) + \lambda \|\mathbf{w}\|_1)$$

Intuition:

- $\|\mathbf{w}\|_1$ measures the “complexity” of hypothesis defined by \mathbf{w}
- λ regulates the tradeoff between the empirical risk ($L_S(\mathbf{w})$) or overfitting and the complexity ($\|\mathbf{w}\|_1$) of the model we pick

LASSO

Linear regression with squared loss + ℓ_1 regularization \Rightarrow LASSO
(*least absolute shrinkage and selection operator*)

LASSO: pick

$$w = \arg \min_w \lambda \|w\|_1 + \sum_{i=1}^m (\langle w, x_i \rangle - y_i)^2$$

How?

Notes:

- no closed form solution!
- ℓ_1 norm is a convex function and squared loss is convex
 \Rightarrow problem can be solved efficiently! (true for every convex loss function)

Tikhonov regularization

Regularization function: $R(\mathbf{w}) = \lambda \|\mathbf{w}\|^2$

- $\lambda \in \mathbb{R}, \lambda > 0$
- ℓ_2 norm: $\|\mathbf{w}\|^2 = \sum_{i=1}^d w_i^2$

Therefore the *learning rule* is: pick

$$A(S) = \arg \min_{\mathbf{w}} (L_S(\mathbf{w}) + \lambda \|\mathbf{w}\|^2)$$

Intuition:

- $\|\mathbf{w}\|^2$ measures the “complexity” of hypothesis defined by \mathbf{w}
- λ regulates the tradeoff between the empirical risk ($L_S(\mathbf{w})$) or overfitting and the complexity ($\|\mathbf{w}\|^2$) of the model we pick

Ridge Regression

Linear regression with squared loss + Tikhonov regularization

\Rightarrow *ridge regression*

Linear regression with squared loss:

- **given:** training set $S = ((x_1, y_1), \dots, (x_m, y_m))$, with $x_i \in \mathbb{R}^d$ and $y_i \in \mathbb{R}$
- **want:** w which minimizes empirical risk:

$$w = \arg \min_w \frac{1}{m} \sum_{i=1}^m (\langle w, x_i \rangle - y_i)^2$$

equivalently, find w which minimizes the *residual sum of squares* $RSS(w)$

$$w = \arg \min_w RSS(w) = \arg \min_w \sum_{i=1}^m (\langle w, x_i \rangle - y_i)^2$$

Linear regression: pick

$$w = \arg \min_w RSS(w) = \arg \min_w \sum_{i=1}^m (\langle w, x_i \rangle - y_i)^2$$

Ridge regression: pick

$$w = \arg \min_w \left(\lambda ||w||^2 + \sum_{i=1}^m (\langle w, x_i \rangle - y_i)^2 \right)$$

RSS: Matrix Form

Let

$$\mathbf{X} = \begin{bmatrix} \cdots & \mathbf{x}_1 & \cdots \\ \cdots & \mathbf{x}_2 & \cdots \\ \cdots & \vdots & \cdots \\ \cdots & \mathbf{x}_m & \cdots \end{bmatrix}$$

\mathbf{X} : *design matrix*

$$\mathbf{y} = \begin{bmatrix} y_1 \\ y_2 \\ \vdots \\ y_m \end{bmatrix}$$

\Rightarrow we have that RSS is

$$\sum_{i=1}^m (\langle \mathbf{w}, \mathbf{x}_i \rangle - y_i)^2 = (\mathbf{y} - \mathbf{X}\mathbf{w})^T (\mathbf{y} - \mathbf{X}\mathbf{w})$$

Ridge Regression: Matrix Form

Linear regression: pick

$$\arg \min_w (\mathbf{y} - \mathbf{X}\mathbf{w})^T (\mathbf{y} - \mathbf{X}\mathbf{w})$$

Ridge regression: pick

$$\arg \min_w \left(\lambda \|\mathbf{w}\|^2 + (\mathbf{y} - \mathbf{X}\mathbf{w})^T (\mathbf{y} - \mathbf{X}\mathbf{w}) \right)$$

Want to find \mathbf{w} which minimizes

$$f(\mathbf{w}) = \lambda \|\mathbf{w}\|^2 + (\mathbf{y} - \mathbf{X}\mathbf{w})^T (\mathbf{y} - \mathbf{X}\mathbf{w}).$$

How?

Compute gradient $\frac{\partial f(\mathbf{w})}{\partial \mathbf{w}}$ of objective function w.r.t \mathbf{w} and compare it to 0.

$$\frac{\partial f(\mathbf{w})}{\partial \mathbf{w}} = 2\lambda \mathbf{w} - 2\mathbf{X}^T (\mathbf{y} - \mathbf{X}\mathbf{w})$$

Then we need to find \mathbf{w} such that

$$2\lambda \mathbf{w} - 2\mathbf{X}^T (\mathbf{y} - \mathbf{X}\mathbf{w}) = 0$$

$$2\lambda \mathbf{w} - 2\mathbf{X}^T(\mathbf{y} - \mathbf{X}\mathbf{w}) = 0$$

is equivalent to

$$(\lambda \mathbf{I} + \mathbf{X}^T \mathbf{X}) \mathbf{w} = \mathbf{X}^T \mathbf{y}$$

Note:

- $\mathbf{X}^T \mathbf{X}$ is positive semidefinite
- $\lambda \mathbf{I}$ is positive definite

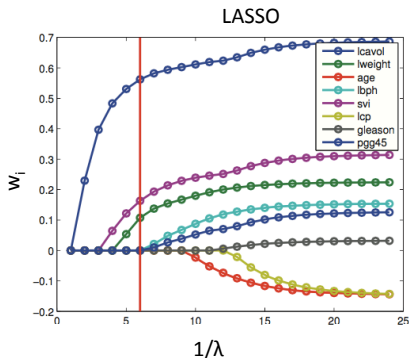
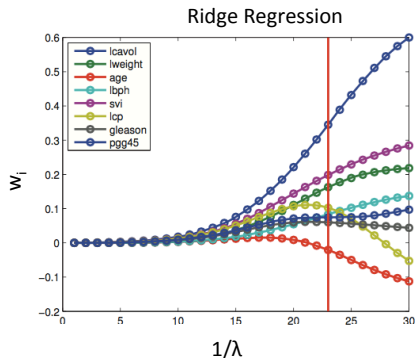
$\Rightarrow \lambda \mathbf{I} + \mathbf{X}^T \mathbf{X}$ is positive definite

$\Rightarrow \lambda \mathbf{I} + \mathbf{X}^T \mathbf{X}$ is invertible

Ridge regression solution:

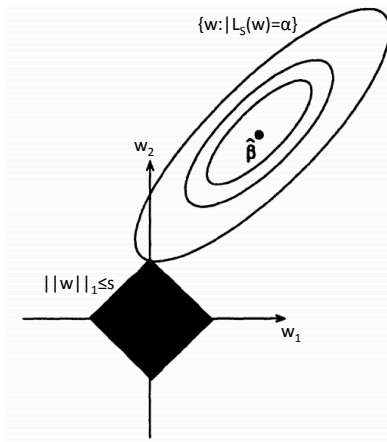
$$\mathbf{w} = (\lambda \mathbf{I} + \mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \mathbf{y}$$

Ridge Regression vs LASSO: Sparsity of Solutions

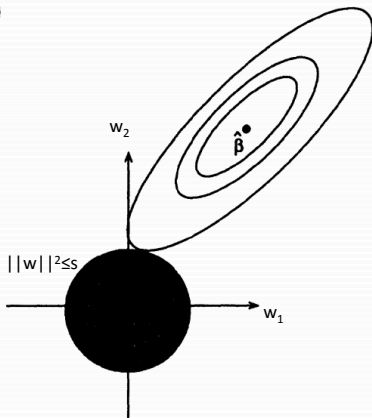


Ridge Regression vs LASSO

LASSO



RIDGE REGRESSION



ℓ_1 regularization performs a sort of **feature selection**

Exercise 5

Consider the ridge regression problem

$\arg \min_w \lambda \|w\|^2 + \sum_{i=1}^m (\langle w, x_i \rangle - y_i)^2$. Let: h_S be the hypothesis obtained by ridge regression on with training set S ; h^* be the hypothesis of minimum generalization error among all linear models.

- (A) Draw, in the plot below, a *typical* behaviour of (i) *the training error* and (ii) *the test/generalization error* of h_S as a function of λ .
- (B) Draw, in the plot below, a *typical* behaviour of (i) $L_D(h_S) - L_D(h^*)$ and (ii) $L_D(h_S) - L_S(h_S)$ as a function of λ .

