l_1 regularization algorithms

Machine Learning: a probabilistic perspective: Kevin P Murphy (Chapter 13.3.2, 13.4) In the case of linear regression, the ℓ_1 objective becomes

$$f(\mathbf{w}) = \sum_{i=1}^{N} -\frac{1}{2\sigma^2} (y_i - (w_0 + \mathbf{w}^T \mathbf{x}_i))^2 + \lambda ||\mathbf{w}||_1$$

$$= \text{RSS}(\mathbf{w}) + \lambda' ||\mathbf{w}||_1$$
(13.36)

where $\lambda' = 2\lambda\sigma^2$. This method is known as **basis pursuit denoising** or **BPDN** (Chen et al. 1998). $\min_{\mathbf{w}} \text{RSS}(\mathbf{w}) + \lambda ||\mathbf{w}||_1$ (13.37)

We can rewrite this as a constrained but smooth objective (a quadratic function with linear constraints):

$$\min_{\mathbf{w}} \text{RSS}(\mathbf{w}) \quad \text{s.t.} \quad ||\mathbf{w}||_1 \le B \tag{13.38}$$

where *B* is an upper bound on the ℓ_1 -norm of the weights: a small (tight) bound *B* corresponds to a large penalty λ , and vice versa.² Equation 13.38 is known as **lasso**, which stands for "least absolute shrinkage and selection operator" (Tibshirani 1996).

> Similarly, we can write ridge regression
> $$\begin{split} & \min_{\mathbf{w}} RSS(\mathbf{w}) + \lambda ||\mathbf{w}||_2^2 \\ & \text{or as a bound constrained form:} \\ & \min_{\mathbf{w}} RSS(\mathbf{w}) \quad \text{ s.t. } ||\mathbf{w}||_2^2 \leq B \end{split}$$



In Figure 13.3, we plot the contours of the RSS objective function, as well as the contours of the ℓ_2 and ℓ_1 constraint surfaces. From the theory of constrained optimization, we know that the optimal solution occurs at the point where the lowest level set of the objective function intersects the constraint surface (assuming the constraint is active). It should be geometrically clear that as we relax the constraint B, we "grow" the ℓ_1 "ball" until it meets the objective; the corners of the ball are more likely to intersect the ellipse than one of the sides, especially in high dimensions, because the corners "stick out" more. The corners correspond to sparse solutions, which lie on the coordinate axes. By contrast, when we grow the ℓ_2 ball, it can intersect the objective at any point; there are no "corners", so there is no preference for sparsity.

 $\mathbf{X}^T \mathbf{X} = \mathbf{I}$. In this case, the RSS is given by

$$RSS(\mathbf{w}) = ||\mathbf{y} - \mathbf{X}\mathbf{w}||^2 = \mathbf{y}^T \mathbf{y} + \mathbf{w}^T \mathbf{X}^T \mathbf{X}\mathbf{w} - 2\mathbf{w}^T \mathbf{X}^T \mathbf{y}$$
(13.59)
$$= \operatorname{const} + \sum w_k^2 - 2\sum \sum w_k x_{ik} y_i$$
(13.60)

so we see this factorizes into a sum of terms, one per dimension. Hence we can write down the MAP and ML estimates analytically, as follows:

MLE The OLS solution is given by

$$\hat{w}_k^{OLS} = \mathbf{x}_{:k}^T \mathbf{y} \tag{13.61}$$

where $\mathbf{x}_{:k}$ is the *k*'th column of **X**. This follows trivially from Equation 13.60. We see that \hat{w}_k^{OLS} is just the orthogonal projection of feature *k* onto the response vector (see Section 7.3.2).

Ridge One can show that the ridge estimate is given by

$$\hat{w}_k^{ridge} = \frac{\hat{w}_k^{OLS}}{1+\lambda} \tag{13.62}$$

• Lasso From Equation 13.55, and using the fact that $a_k = 2$ and $\hat{w}_k^{OLS} = c_k/2$, we have

$$\hat{w}_{k}^{lasso} = \operatorname{sign}(\hat{w}_{k}^{OLS}) \left(|\hat{w}_{k}^{OLS}| - \frac{\lambda}{2} \right)_{+}$$
(13.63)

This corresponds to soft thresholding, shown in Figure 13.5(a).

 Subset selection If we pick the best K features using subset selection, the parameter estimate is as follows

$$\hat{w}_{k}^{SS} = \begin{cases} \hat{w}_{k}^{OLS} & \text{if } \operatorname{rank}(|w_{k}^{OLS}|) \leq K \\ 0 & \text{otherwise} \end{cases}$$
(13.64)

where rank refers to the location in the sorted list of weight magnitudes. This corresponds to hard thresholding, shown in Figure 13.5(b).

Optimal solution to LASSO

• Let the lasso problem be defined as $f(w) = RSS(w) + \lambda ||w||_{1}$ Where $RSS(w) = ||Xw - y||_{2}^{2}$ It can be shown that (Exercise 13.1) $\frac{\partial}{\partial w_{j}}RSS(w) = a_{j}w_{j} - c_{j}, where$ $a_{j} = 2\sum_{i=1}^{N} x_{ij}^{2}, = 2 ||x_{:,j}||^{2}$

$$c_j = 2\sum_{i=1}^N x_{ij} (y_i - \boldsymbol{w}_{-j}^T x_{i,-j}) = 2\mathbf{x}_{:,j}^T \boldsymbol{r}_j$$

where w_{-j} is w without the component j and similarly for $x_{i,-j}$, $r_j = y_i - w_{-j}^T x_{i,-j}$ Adding in the penalty term, we find that the subderivative¹ is given by

$$\frac{\partial}{\partial w_{j}} f(\mathbf{w}) = (a_{j}w_{j} - c_{j}) + \lambda \frac{\partial}{\partial w_{j}} ||\mathbf{w}||_{1}$$

$$= \begin{cases} \{a_{j}w_{j} - c_{j} - \lambda\} & \text{if } w_{j} < 0 \\ [-c_{j} - \lambda, -c_{j} + \lambda] & \text{if } w_{j} = 0 \\ a_{j}w_{j} - c_{j} + \lambda & \text{if } w_{j} > 0 \end{cases}$$

- In the second case, when $w_j = 0, -\lambda \le c_j \le \lambda$.
- Thus, depending on the value of c_j , the solution to $\frac{\partial}{\partial w_j} f(\mathbf{w}) = 0$ can occur at 3 different values of w_j
- Subderivatives are computed because of the discontinuity at the corners imposed by constraints causing certain w_i's going to 0.

1. Subderivative defined in Murphy – Sec. 13.3.2

Soft thresholding

$$\widehat{w}_{j}(c_{j}) = \begin{cases} \frac{c_{j} + \lambda}{a_{j}} & \text{if } c_{j} < -\lambda \\ 0 & \text{if } c_{j} \in [-\lambda, \lambda] \\ \frac{c_{j} - \lambda}{a_{j}} & \text{if } c_{j} > \lambda \end{cases}$$

This can be written as

$$\widehat{w}_{j} = soft\left(\frac{c_{j}}{a_{j}}, \frac{\lambda}{a_{j}}\right)$$
where $soft(\frac{c_{j}}{a_{j}}, \frac{\lambda}{a_{j}}) \equiv sign\left(\frac{c_{j}}{a_{j}}\right)\left(\left|\frac{c_{j}}{a_{j}}\right| - \frac{\lambda}{a_{j}}\right)_{+}$
 $c_{i} = \lambda$

 $\left(\left|\frac{c_j}{a_j}\right| - \frac{\lambda}{a_j}\right)_+ = \max\left(\left(\left|\frac{c_j}{a_j}\right| - \frac{\lambda}{a_j}\right), 0\right)$, the positive part. This is called

soft thresholding

Soft thresholding vs hard thresholding

This is illustrated in Figure 13.5(a), plotting $\hat{w}_j vs c_j$. The black line is the line $w_j = c_j/a_j$ corresponding to the least squares fit. The red line, which represents the regularized estimate shifts the red line down (or up) by λ , except when $-\lambda \leq c \leq \lambda$ in which case it sets $w_j = 0$.



• Figure 13.5 Left: soft thresholding. The flat region is the interval $[-\lambda, +\lambda]$. Right: hard thresholding.

Hard thresholding

- By contrast, in Figure 13.5(b), hard thresholding is illustrated. This sets values of w_j to 0 if $-\lambda \le c_j \le \lambda$, but it does not shrink the values of w_j outside of this interval
- The slope of the soft thresholding line does not coincide with the diagonal, which means that even large coefficients are shrunk towards zero; consequently lasso is a biased estimator.
- This is undesirable, since if the likelihood indicates (via c_j) that the coefficient w_j should be large, we do not want to shrink it.

l_1 regularization algorithms

- The various algorithms solving l_1 regularization problems are
 - Coordinate descent
 - LARS
 - Proximal Gradient

Coordinate descent algorith,

• Lasso objective has the form

$$f(w) = RSS(w) + \lambda ||w||_{1}$$
Where $RSS(w) = ||Xw - y||_{2}^{2}$

- The coordinate descent aims at obtaining w_j one at a time
- In particular, we can solve for the j'th coefficient with all the others held fixed.

Coordinate descent algorithm

Algorithm 13.1: Coordinate descent for lasso (aka shooting algorithm)

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

2 repeat

$$\begin{array}{c|c} \mathbf{s} & \text{for } j = 1, \dots, D \text{ do} \\ \mathbf{4} & a_j = 2 \sum_{i=1}^n x_{ij}^2; \\ \mathbf{5} & c_j = 2 \sum_{i=1}^n x_{ij} (y_i - \mathbf{w}^T \mathbf{x}_i + w_j x_{ij}); \\ \mathbf{6} & w_j = \text{soft}(\frac{c_j}{a_j}, \frac{\lambda}{a_j}); \end{array}$$

7 until converged;

Disadvantage : Since the variables are solved one at a time, the process is time consuming. Thus, methods such as LARS and proximal gradient are used

Proximal gradient method

• Consider a convex objective function of the form

 $f(\boldsymbol{w}) = L(\boldsymbol{w}) + R(\boldsymbol{w})$

where L(w), representing loss, is convex and differentiable, whereas R(w), representing regularization, is convex and maybe non-differentiable

- Solution is to split the functions individually to solve the loss function and the regularization function (using proximal gradient method).
- Named "proximal" because the non-smooth function is solved using its 'proximity operator'

Proximal operator

- For example, consider L(w) = RSS(w), let $RSS(w) = ||y Xw||_2^2$
- For the the design matrix X = I, the objective function becomes $f(w) = R(w) + \frac{1}{2} ||w - y||_2^2$
- The minimizer of the above function is given by the **proximal** operator $(prox_R(w))$ for the convex set R.

$$prox_{R}(\mathbf{y}) = argmin_{Z}\left(R(z) + \frac{1}{2}\left|\left|z - \mathbf{y}\right|\right|_{2}^{2}\right)$$

 Intuitively, a point that minimizers R and also close (proximal) to wy is returned

Examples of Proximal operators

- If $R(w) = \lambda ||w||_1$, the proximal operator is given by component wise soft thresholding $prox_R(w) = soft(w, \lambda)$
- If $R(w) = \lambda ||w||_0$, the proximal operator is given by component wise hard thresholding $prox_R(w) = hard(w, \sqrt{2\lambda})$
- If $R(w) = I_C(w)$, the proximal operator is given by the projection on to the set C: $prox_R(w) = argmin_{z \in C} ||z - w||_2^2 = proj_C(w)$

Examples of Projection (proximal) operators

- For some convex sets, it is easy to compute the **projection operator**.
- For example, to project onto the rectangular set defined by the box constraints $C = \{w: l_j \le w_j \le u_j\}$, we can use

$$proj_{C}(\boldsymbol{w})_{j} = \begin{cases} l_{j} & w_{j} \leq l_{j} \\ w_{j} & l_{j} \leq w_{j} \leq u_{j} \\ u_{j} & w_{j} \geq u_{j} \end{cases}$$

Examples of Projection (proximal) operators

• To project onto the Euclidean ball

$$C = \{ \boldsymbol{w} : \left| |\boldsymbol{w}| \right|_{2} \le 1 \} \text{ we can use}$$
$$proj_{C}(\boldsymbol{w}) = \begin{cases} \frac{\boldsymbol{w}}{\left| |\boldsymbol{w}| \right|_{2}} & \left| |\boldsymbol{w}| \right|_{2} > 1 \\ \boldsymbol{w} & \left| |\boldsymbol{w}| \right|_{2} \le 1 \end{cases}$$

• To project onto the 1-norm ball

 $C = \{ \boldsymbol{w} : ||\boldsymbol{w}||_1 \le 1 \} \text{ we can use}$ $proj_C(\boldsymbol{w}) = \operatorname{soft}(\boldsymbol{w}, \lambda), \text{ where}$

$$\lambda = 0,$$
 if $||w||_1 \le 1$ and
 λ is the soln of: $\sum_{j=1}^{D} \max(|w_j| - \lambda, 0) = 1$, otherwise

Proximal gradient method

- Proximal operator used in gradient descent routine
- Consider the function f(w) = L(w) + R(w)
- Basic idea: minimize a simple quadratic approximation to the loss function, centered on w^k for k^{th} iteration.
- Consider the Taylor series expansion of L(z)

$$L(z) = L(\mathbf{w}^k) + g_k^T(z - \mathbf{w}^k) + \frac{1}{2t_k} \left\| |z - \mathbf{w}^k| \right\|_2^2$$

where $g_k = \nabla L(\mathbf{w}^k)$ and a simple approximation of Hessian of the loss $\nabla^2 L(\mathbf{w}^k) \approx \frac{1}{t_k} I$

Thus,

$$\mathbf{w}^{k+1} = argmin_{z}(R(z) + L(\mathbf{w}^{k}) + g_{k}^{T}(z - \mathbf{w}^{k}) + \frac{1}{2t_{k}} \left\| z - \mathbf{w}^{k} \right\|_{2}^{2}$$

Proximal gradient method

 Dropping the terms independent of z, multiplying by t_k and rewriting the above equation in terms of proximal operator for kth iteration:

$$w^{k+1} = argmin_{z} \left[t_{k}R(z) + \frac{1}{2} ||z - u_{k}||_{2}^{2} \right] = prox_{t_{k}R}(u_{k})$$
$$u_{k} = w^{k} - t_{k}g_{k}$$
$$g_{k} = \nabla L(w^{k})$$
$$(Verify: g_{k}^{T}(z - w^{k}) + \frac{1}{2t_{k}} ||z - w^{k}||_{2}^{2} = \frac{1}{2t_{k}} ||z - u_{k}||_{2}^{2})$$

- If R(w) = 0, this is equivalent to gradient descent.
- If R(w) = I_C(w), this method is equivalent to projected gradient descent
- If $R(w) = \lambda ||w||_1$, the method is known as **iterative soft thresholding**

References

 Machine Learning – A probabilistic perspective, Kevin P Murphy