

LECTURE 17: LASSO AND COORDINATE DESCENT

STAT 598Z: INTRODUCTION TO COMPUTING FOR STATISTICS


Vinayak Rao

Department of Statistics, Purdue University

March 23, 2017

BIAS-VARIANCE AND REGULARIZATION

Problem: Given training data $(X, y) \equiv \{x_i, y_i\}$,
minimize $\mathcal{L}(w) = \frac{1}{2}(Y - Xw)^2$

$$y = x^T w$$


Problem: Given training data $(\mathbf{X}, \mathbf{y}) \equiv \{\mathbf{x}_i, y_i\}$,
minimize $\mathcal{L}(\mathbf{w}) = \frac{1}{2}(\mathbf{Y} - \mathbf{X}\mathbf{w})^2$

To reduce variance (i.e. sensitivity to small changes in training data) , add a penalty $\Omega(\mathbf{w})$:

$$\hat{\mathbf{w}} = \operatorname{argmin} \mathcal{L}(\mathbf{w}) + \lambda\Omega(\mathbf{w})$$

BIAS-VARIANCE AND REGULARIZATION

Problem: Given training data $(\mathbf{X}, \mathbf{y}) \equiv \{\mathbf{x}_i, y_i\}$,
minimize $\mathcal{L}(\mathbf{w}) = \frac{1}{2}(\mathbf{Y} - \mathbf{X}\mathbf{w})^2$

To reduce variance (i.e. sensitivity to small changes in training data), add a penalty $\Omega(\mathbf{w})$:

$$\hat{\mathbf{w}} = \operatorname{argmin} \mathcal{L}(\mathbf{w}) + \lambda\Omega(\mathbf{w})$$

Ridge regression/ L_2 regression:

- $\Omega(\mathbf{w}) = \|\mathbf{w}\|_2^2$
- $\hat{\mathbf{w}} = (\mathbf{X}^\top \mathbf{X} + \lambda \mathbf{I})^{-1} \mathbf{X}^\top \mathbf{y}$ (Shrinkage)

Problem: Given training data $(\mathbf{X}, \mathbf{y}) \equiv \{\mathbf{x}_i, y_i\}$,
minimize $\mathcal{L}(\mathbf{w}) = \frac{1}{2}(\mathbf{Y} - \mathbf{X}\mathbf{w})^2$

To reduce variance (i.e. sensitivity to small changes in training data) , add a penalty $\Omega(\mathbf{w})$:

$$\hat{\mathbf{w}} = \operatorname{argmin} \mathcal{L}(\mathbf{w}) + \lambda\Omega(\mathbf{w})$$

LASSO:

- $\Omega(\mathbf{w}) = \|\mathbf{w}\|_1$
- Shrinkage and selection
(\mathbf{w} is sparse with some components equal to 0)
- No simple closed-form solution

REGULARIZATION AS CONSTRAINED OPTIMIZATION

$\operatorname{argmin}(\mathbf{y} - \mathbf{X}^\top \mathbf{w})^2 + \lambda \|\mathbf{w}\|_2^2$ is equivalent to

$\operatorname{argmin}(\mathbf{y} - \mathbf{X}^\top \mathbf{w})^2 \quad \text{s.t.} \quad \|\mathbf{w}\|_2^2 \leq \gamma$

(Note: γ will depend on data)

REGULARIZATION AS CONSTRAINED OPTIMIZATION

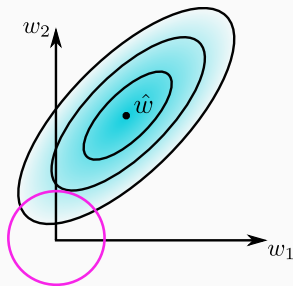
$\operatorname{argmin}(\mathbf{y} - \mathbf{X}^\top \mathbf{w})^2 + \lambda \|\mathbf{w}\|_2^2$ is equivalent to

$\operatorname{argmin}(\mathbf{y} - \mathbf{X}^\top \mathbf{w})^2 \quad \text{s.t.} \quad \|\mathbf{w}\|_2^2 \leq \gamma$

(Note: γ will depend on data)

First problem: regularized optimization

Second problem: constrained optimization



REGULARIZATION AS CONSTRAINED OPTIMIZATION

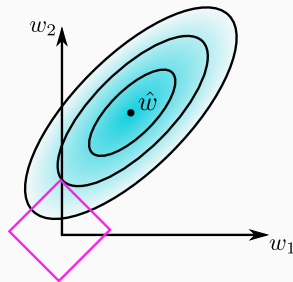
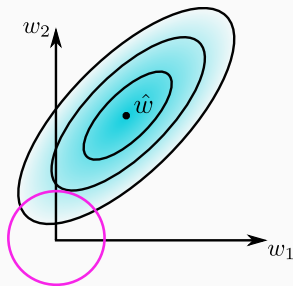
$\operatorname{argmin}(\mathbf{y} - \mathbf{X}^\top \mathbf{w})^2 + \lambda \|\mathbf{w}\|_2^2$ is equivalent to

$\operatorname{argmin}(\mathbf{y} - \mathbf{X}^\top \mathbf{w})^2 \quad \text{s.t.} \quad \|\mathbf{w}\|_2^2 \leq \gamma$

(Note: γ will depend on data)

First problem: regularized optimization

Second problem: constrained optimization



$$\operatorname{argmin}(\mathbf{y} - \mathbf{X}^\top \mathbf{w})^2 \quad \text{s.t. } \|\mathbf{w}\|_1 \leq \gamma$$

$\|\mathbf{w}\|_1 = \sum_{j=1}^p |w_j|$ is the ℓ_1 -norm.

$$\operatorname{argmin}(\mathbf{y} - \mathbf{X}^\top \mathbf{w})^2 \quad \text{s.t. } \|\mathbf{w}\|_1 \leq \gamma$$

$\|\mathbf{w}\|_1 = \sum_{j=1}^p |w_j|$ is the ℓ_1 -norm.

Lasso: least absolute shrinkage and selection operator.

$$\hat{\mathbf{w}} = \operatorname{argmin} \sum_{i=1}^n (y_i - \mathbf{x}_i^\top \mathbf{w})^2 + \lambda \|\mathbf{w}\|_1$$

$$\operatorname{argmin}(\mathbf{y} - \mathbf{X}^\top \mathbf{w})^2 \quad \text{s.t. } \|\mathbf{w}\|_1 \leq \gamma$$

$\|\mathbf{w}\|_1 = \sum_{j=1}^p |w_j|$ is the ℓ_1 -norm.

Lasso: least absolute shrinkage and selection operator.

$$\hat{\mathbf{w}} = \operatorname{argmin} \sum_{i=1}^n (y_i - \mathbf{x}_i^\top \mathbf{w})^2 + \lambda \|\mathbf{w}\|_1$$

- Penalizes small w_j more than ridge regression.
- Tolerates larger w_j more than ridge regression.

$$\operatorname{argmin}(\mathbf{y} - \mathbf{X}^\top \mathbf{w})^2 \quad \text{s.t. } \|\mathbf{w}\|_1 \leq \gamma$$

$\|\mathbf{w}\|_1 = \sum_{j=1}^p |w_j|$ is the ℓ_1 -norm.

Lasso: least absolute shrinkage and selection operator.

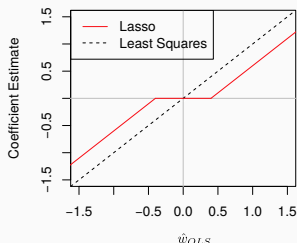
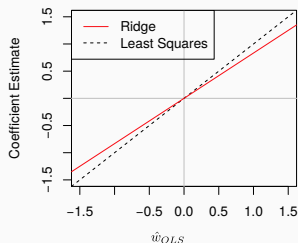
$$\hat{\mathbf{w}} = \operatorname{argmin} \sum_{i=1}^n (y_i - \mathbf{x}_i^\top \mathbf{w})^2 + \lambda \|\mathbf{w}\|_1$$

- Penalizes small w_j more than ridge regression.
- Tolerates larger w_j more than ridge regression.

Result:

- $\hat{\mathbf{w}}_{\text{LASSO}}$ has some components *exactly* equal to zero.
- Performs feature selection.

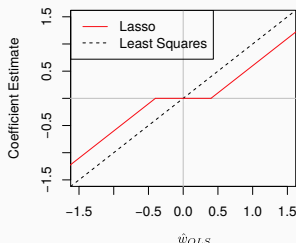
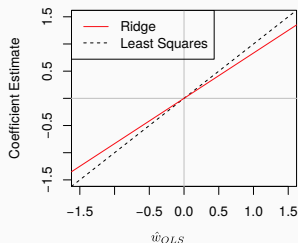
THE 1-D CASE



In the 1-d case, $(\mathbf{x}, \mathbf{y}) \equiv \{x_i, y_i\}$

Least-squares solution: $\hat{w}_{ols} = \frac{\mathbf{x}^\top \mathbf{y}}{\mathbf{x}^\top \mathbf{x}}$

THE 1-D CASE

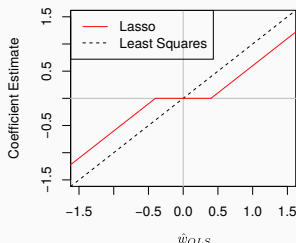
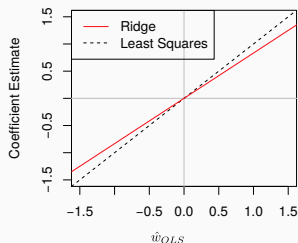


In the 1-d case, $(\mathbf{x}, \mathbf{y}) \equiv \{x_i, y_i\}$

Least-squares solution: $\hat{w}_{ols} = \frac{\mathbf{x}^\top \mathbf{y}}{\mathbf{x}^\top \mathbf{x}}$

Ridge regression solution: $\hat{w}_{ridge} = \frac{\mathbf{x}^\top \mathbf{x}}{\mathbf{x}^\top \mathbf{x} + \lambda}$

THE 1-D CASE



In the 1-d case, $(\mathbf{x}, \mathbf{y}) \equiv \{x_i, y_i\}$

Least-squares solution: $\hat{w}_{ols} = \frac{\mathbf{x}^\top \mathbf{y}}{\mathbf{x}^\top \mathbf{x}}$

Ridge regression solution: $\hat{w}_{ridge} = \frac{\mathbf{x}^\top \mathbf{x}}{\mathbf{x}^\top \mathbf{x} + \lambda}$

LASSO solution?

Use the `optim` function

Syntax:

```
optim(par, fn, gr = NULL, ...,  
      method = c('Nelder-Mead', 'BFGS', 'CG', 'L-BFGS-B', 'SANN',  
                  'Brent'),  
      lower = -Inf, upper = Inf,  
      control = list(), hessian = FALSE)
```

`fn`: function to be optimized

`gr`: gradient function (calculate numerically if `NULL`)

`par`: initial value of parameter to be optimized (should be first argument of `fn`)

$$\hat{w} = \operatorname{argmin} \mathcal{L}(w) = \operatorname{argmin} \frac{1}{2} \sum_{i=1}^n (y_i - wx_i)^2 + \lambda |w|$$

$$\hat{w} = \operatorname{argmin} \mathcal{L}(w) = \operatorname{argmin} \frac{1}{2} \sum_{i=1}^n (y_i - wx_i)^2 + \lambda |w|$$

At the minimum,

$$\frac{d\mathcal{L}}{dw} = 0$$

$$\hat{w} = \operatorname{argmin} \mathcal{L}(w) = \operatorname{argmin} \frac{1}{2} \sum_{i=1}^n (y_i - wx_i)^2 + \lambda |w|$$

At the minimum,

$$\begin{aligned} \frac{d\mathcal{L}}{dw} &= 0 \\ - \sum_{i=1}^n (y_i - wx_i)x_i + \lambda \frac{d|w|}{dw} &= 0 \end{aligned}$$

$$\hat{w} = \operatorname{argmin} \mathcal{L}(w) = \operatorname{argmin} \frac{1}{2} \sum_{i=1}^n (y_i - wx_i)^2 + \lambda |w|$$

At the minimum,

$$\begin{aligned} \frac{d\mathcal{L}}{dw} &= 0 \\ - \sum_{i=1}^n (y_i - wx_i)x_i + \lambda \frac{d|w|}{dw} &= 0 \\ - \sum_{i=1}^n y_i x_i + w \sum_{i=1}^n x_i^2 + \lambda \frac{d|w|}{dw} &= 0 \end{aligned}$$

$$\hat{w} = \operatorname{argmin} \mathcal{L}(w) = \operatorname{argmin} \frac{1}{2} \sum_{i=1}^n (y_i - wx_i)^2 + \lambda |w|$$

At the minimum,

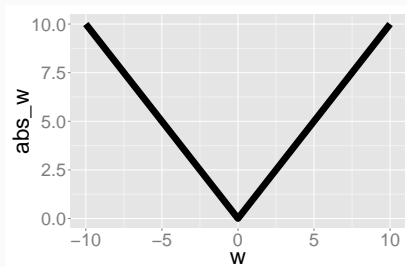
$$\begin{aligned} \frac{d\mathcal{L}}{dw} &= 0 \\ - \sum_{i=1}^n (y_i - wx_i)x_i + \lambda \frac{d|w|}{dw} &= 0 \\ - \sum_{i=1}^n y_i x_i + w \sum_{i=1}^n x_i^2 + \lambda \frac{d|w|}{dw} &= 0 \\ w &= \frac{\sum_{i=1}^n y_i x_i - \lambda \frac{d|w|}{dw}}{\sum_{i=1}^n x_i^2} \end{aligned}$$

SUBGRADIENTS

$$w = \frac{\sum_{i=1}^n y_i x_i - \lambda \frac{d|w|}{dw}}{\sum_{i=1}^n x_i^2} = \frac{\mathbf{y}^\top \mathbf{x} - \lambda \frac{d|w|}{dw}}{\mathbf{x}^\top \mathbf{x}} : \quad \text{What is } \frac{d|w|}{dw}?$$

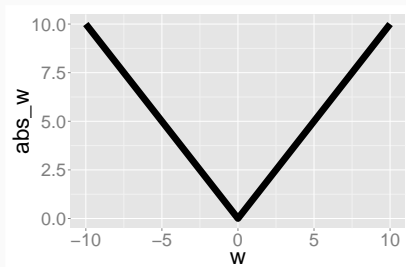
SUBGRADIENTS

$$w = \frac{\sum_{i=1}^n y_i x_i - \lambda \frac{d|w|}{dw}}{\sum_{i=1}^n x_i^2} = \frac{y^\top x - \lambda \frac{d|w|}{dw}}{x^\top x} : \quad \text{What is } \frac{d|w|}{dw}?$$



SUBGRADIENTS

$$w = \frac{\sum_{i=1}^n y_i x_i - \lambda \frac{d|w|}{dw}}{\sum_{i=1}^n x_i^2} = \frac{y^\top x - \lambda \frac{d|w|}{dw}}{x^\top x} : \quad \text{What is } \frac{d|w|}{dw}?$$



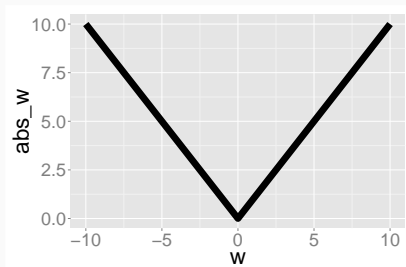
$$w > 0 \quad \Leftrightarrow \frac{d|w|}{dw} = 1$$

$$w < 0 \quad \Leftrightarrow \frac{d|w|}{dw} = -1$$

$$w = 0 \quad \Leftrightarrow \frac{d|w|}{dw} \in (-1, 1)$$

SUBGRADIENTS

$$w = \frac{\sum_{i=1}^n y_i x_i - \lambda \frac{d|w|}{dw}}{\sum_{i=1}^n x_i^2} = \frac{y^\top x - \lambda \frac{d|w|}{dw}}{x^\top x} : \quad \text{What is } \frac{d|w|}{dw}?$$



$$w > 0 \quad \Leftrightarrow \frac{d|w|}{dw} = 1$$

$$w > 0 \quad \Leftrightarrow w = \frac{y^\top x - \lambda}{x^\top x}$$

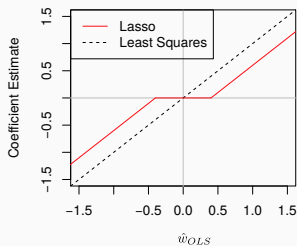
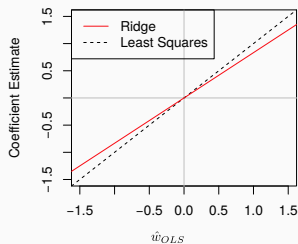
$$w < 0 \quad \Leftrightarrow \frac{d|w|}{dw} = -1$$

$$w < 0 \quad \Leftrightarrow w = \frac{y^\top x + \lambda}{x^\top x}$$

$$w = 0 \quad \Leftrightarrow \frac{d|w|}{dw} \in (-1, 1)$$

$$w = 0 \quad \Leftrightarrow w = \text{otherwise}$$

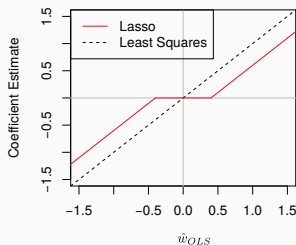
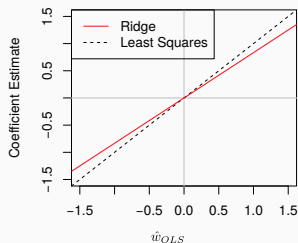
THE 1-D CASE



LASSO

First calculate: $\hat{w}_{ols} = \frac{\mathbf{y}^\top \mathbf{x}}{\mathbf{x}^\top \mathbf{x}}$

THE 1-D CASE



LASSO

First calculate: $\hat{w}_{ols} = \frac{\mathbf{y}^\top \mathbf{x}}{\mathbf{x}^\top \mathbf{x}}$

Soft threshold: $\hat{w}_{LASSO} = \text{sign}(\hat{w}_{ols})(|\hat{w}_{ols}| - \frac{\lambda}{\mathbf{x}^\top \mathbf{x}})_+$

$(x)_+ = x$ if $x > 0$, else 0, and

$\text{sign}(x) = +1$ if $x > 0$ else -1

LASSO IN HIGHER (P) DIMENSIONS

Find \mathbf{w} by coordinate descent

$$\mathcal{L}(\mathbf{w}) = \sum_{i=1}^n (y_i - \mathbf{w}^\top \mathbf{x}_i)^2 + \lambda \|\mathbf{w}\|_1 \quad (1)$$

(3)

LASSO IN HIGHER (P) DIMENSIONS

Find \mathbf{w} by coordinate descent

$$\mathcal{L}(\mathbf{w}) = \sum_{i=1}^n (y_i - \mathbf{w}^\top \mathbf{x}_i)^2 + \lambda \|\mathbf{w}\|_1 \quad (1)$$

$$= \sum_{i=1}^n (y_i - \sum_{j=1}^p w_j x_{ij})^2 + \lambda \sum_{j=1}^p |w_j| \quad (2)$$

$$(3)$$

LASSO IN HIGHER (P) DIMENSIONS

Find \mathbf{w} by coordinate descent

$$\mathcal{L}(\mathbf{w}) = \sum_{i=1}^n (y_i - \mathbf{w}^\top \mathbf{x}_i)^2 + \lambda \|\mathbf{w}\|_1 \quad (1)$$

$$= \sum_{i=1}^n (y_i - \sum_{j=1}^p w_j x_{ij})^2 + \lambda \sum_{j=1}^p |w_j| \quad (2)$$

$$= \sum_{i=1}^n (r_{id} - w_d x_{id})^2 + \lambda |w_d| + C \quad (3)$$

LASSO IN HIGHER (P) DIMENSIONS

Find \mathbf{w} by coordinate descent

$$\mathcal{L}(\mathbf{w}) = \sum_{i=1}^n (y_i - \mathbf{w}^\top \mathbf{x}_i)^2 + \lambda \|\mathbf{w}\|_1 \quad (1)$$

$$= \sum_{i=1}^n (y_i - \sum_{j=1}^p w_j x_{ij})^2 + \lambda \sum_{j=1}^p |w_j| \quad (2)$$

$$= \sum_{i=1}^n (r_{id} - w_d x_{id})^2 + \lambda |w_d| + C \quad (3)$$

LASSO IN HIGHER (P) DIMENSIONS

Find \mathbf{w} by coordinate descent

$$\mathcal{L}(\mathbf{w}) = \sum_{i=1}^n (y_i - \mathbf{w}^\top \mathbf{x}_i)^2 + \lambda \|\mathbf{w}\|_1 \quad (1)$$

$$= \sum_{i=1}^n (y_i - \sum_{j=1}^p w_j x_{ij})^2 + \lambda \sum_{j=1}^p |w_j| \quad (2)$$

$$= \sum_{i=1}^n (r_{id} - w_d x_{id})^2 + \lambda |w_d| + C \quad (3)$$

Here r_{id} is the residual of obs. i :

$$r_{id} = y_i - \sum_{j \neq d} w_j x_{ij}$$

LASSO IN HIGHER (P) DIMENSIONS

Find \mathbf{w} by coordinate descent

$$\mathcal{L}(\mathbf{w}) = \sum_{i=1}^n (y_i - \mathbf{w}^\top \mathbf{x}_i)^2 + \lambda \|\mathbf{w}\|_1 \quad (1)$$

$$= \sum_{i=1}^n (y_i - \sum_{j=1}^p w_j x_{ij})^2 + \lambda \sum_{j=1}^p |w_j| \quad (2)$$

$$= \sum_{i=1}^n (r_{id} - w_d x_{id})^2 + \lambda |w_d| + C \quad (3)$$

Here r_{id} is the **residual** of obs. i :

$$r_{id} = y_i - \sum_{j \neq d} w_j x_{ij}$$

Eq(3) is just 1d LASSO! Can solve for w_d by soft-thresholding.

Repeat

Initialize \mathbf{w} to some arbitrary value

For dimension d , calculate the residual $\mathbf{r}_d = (r_{1d}, \dots, r_{nd})$,

$r_{id} = y_i - \sum_{j \neq d} w_j x_{ij}$ for each observation i

Set $\hat{w}_{ols} = \frac{(\mathbf{x}_d)^\top \mathbf{r}_d}{(\mathbf{x}_d)^\top \mathbf{x}_d}$ where \mathbf{x}_d is the d th column of \mathbf{X} and we have:

$$\hat{w}_d = \text{sign}(\hat{w}_{ols}) \left(|\hat{w}_{ols}| - \frac{\lambda}{(\mathbf{x}_d)^\top \mathbf{x}_d} \right)_+$$

Initialize \mathbf{w} to some arbitrary value

For dimension d , calculate the residual $\mathbf{r}_d = (r_{1d}, \dots, r_{nd})$,

$r_{id} = y_i - \sum_{j \neq d} w_j x_{ij}$ for each observation i

Set $\hat{w}_{ols} = \frac{(\mathbf{x}_d)^\top \mathbf{r}_d}{(\mathbf{x}_d)^\top \mathbf{x}_d}$ where \mathbf{x}_d is the d th column of \mathbf{X} and we have:

$$\hat{w}_d = \text{sign}(\hat{w}_{ols}) \left(|\hat{w}_{ols}| - \frac{\lambda}{(\mathbf{x}_d)^\top \mathbf{x}_d} \right)_+$$

Repeat across dimensions d till convergence.

Initialize \mathbf{w} to some arbitrary value

For dimension d , calculate the residual $\mathbf{r}_d = (r_{1d}, \dots, r_{nd})$,

$r_{id} = y_i - \sum_{j \neq d} w_j x_{ij}$ for each observation i

Set $\hat{w}_{ols} = \frac{(\mathbf{x}_d)^\top \mathbf{r}_d}{(\mathbf{x}_d)^\top \mathbf{x}_d}$ where \mathbf{x}_d is the d th column of \mathbf{X} and we have:

$$\hat{w}_d = \text{sign}(\hat{w}_{ols}) \left(|\hat{w}_{ols}| - \frac{\lambda}{(\mathbf{x}_d)^\top \mathbf{x}_d} \right)_+$$

Repeat across dimensions d till convergence.

Does this work?

DOES CO-ORDINATE DESCENT WORK?

For convex differentiable functions: yes

Convex function f : local optimum is a global minimum.

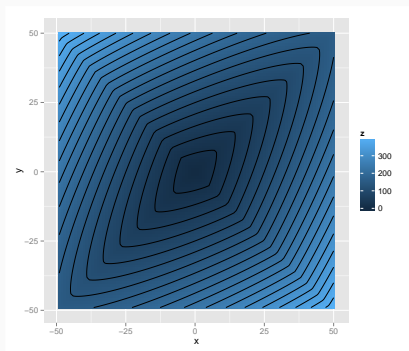
Local optimum for a differentiable function:

$$\nabla f(\mathbf{w}) = \left[\frac{\partial f}{\partial w_1}, \dots, \frac{\partial f}{\partial w_p} \right] = 0$$

At a stationary point of coordinate descent, the RHS is true.

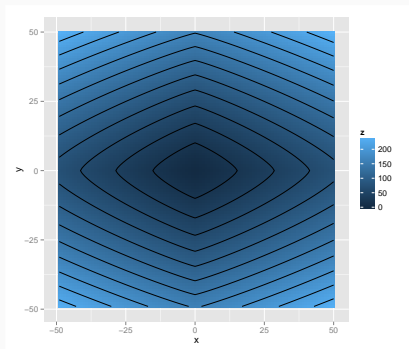
DOES CO-ORDINATE DESCENT WORK?

For convex non-differentiable functions: in general, no!



DOES CO-ORDINATE DESCENT WORK?

For functions of the form: $f(\mathbf{w}) = g(\mathbf{w}) + \sum_{i=1}^p h_i(w_i)$, where f is convex and differentiable, h_i 's are convex but not differentiable: yes



Competitive with state-of-the-art optimization for LASSO

Competitive with state-of-the-art optimization for LASSO

Since objective function is convex, any initialization works
(though some are better)

Competitive with state-of-the-art optimization for LASSO

Since objective function is convex, any initialization works (though some are better)

Obtains the solution $\hat{\mathbf{w}}$ for any λ

Competitive with state-of-the-art optimization for LASSO

Since objective function is convex, any initialization works (though some are better)

Obtains the solution $\hat{\mathbf{w}}$ for any λ

Can repeat for different λ 's (though some ways are better).

We want $\hat{\mathbf{w}}$'s for a set of λ 's

We want $\hat{\mathbf{w}}$'s for a set of λ 's

Pick a smallest and largest λ (latter corresponding to $\hat{\mathbf{w}} = 0$)

Divide into equidistant grid points (typ. on logscale)

We want $\hat{\mathbf{w}}$'s for a set of λ 's

Pick a smallest and largest λ (latter corresponding to $\hat{\mathbf{w}} = 0$)

Divide into equidistant grid points (typ. on logscale)

Start with the largest λ (solution = 0).

We want $\hat{\mathbf{w}}$'s for a set of λ 's

Pick a smallest and largest λ (latter corresponding to $\hat{\mathbf{w}} = 0$)

Divide into equidistant grid points (typ. on logscale)

Start with the largest λ (solution = 0).

Move to the next, using previous solution as initialization.

We want $\hat{\mathbf{w}}$'s for a set of λ 's

Pick a smallest and largest λ (latter corresponding to $\hat{\mathbf{w}} = 0$)

Divide into equidistant grid points (typ. on logscale)

Start with the largest λ (solution = 0).

Move to the next, using previous solution as initialization.

Converges after a few sweeps

Repeat

We want $\hat{\mathbf{w}}$'s for a set of λ 's

Pick a smallest and largest λ (latter corresponding to $\hat{\mathbf{w}} = 0$)

Divide into equidistant grid points (typ. on logscale)

Start with the largest λ (solution = 0).

Move to the next, using previous solution as initialization.

Converges after a few sweeps

Repeat

This kind of a guided search is often faster, even if we just want one λ .