



UiO • Matematisk institutt

Det matematisk-naturvitenskapelige fakultet

STK-4051/9051 Computational Statistics Spring 2022 IRLS and ADMM

Instructor: Odd Kolbjørnsen, oddkol@math.uio.no



Last time

- Introduction – (things you will learn)
- Gradient based methods
 - Newton: $\mathbf{B} = \mathbf{J}(\boldsymbol{\theta}^{(t)})^{-1}$ $\boldsymbol{\theta}^{(t+1)} = \boldsymbol{\theta}^{(t)} + \mathbf{B}\mathbf{s}(\boldsymbol{\theta}^{(t)})$
 - Fisher scoring, $\mathbf{B} = \mathbf{I}(\boldsymbol{\theta}^{(t)})^{-1} = \mathbf{E} \left(\mathbf{J}(\boldsymbol{\theta}^{(t)}) \right)^{-1} = \text{Var} \left(\mathbf{s}(\boldsymbol{\theta}^{(t)}) \right)^{-1}$
 - Secant, \mathbf{B} : discrete approximation of $\mathbf{J}(\boldsymbol{\theta}^{(t)})^{-1}$
 - BFGS, (Quasi newton, optim in R) $\mathbf{B} = -\alpha \mathbf{M}$
 - Ascent, $\mathbf{B} = \alpha \mathbf{I}$, $\alpha > 0$, but small enough
 - Gauss – Newton , linearize around theta, update using linear regression
- Gauss Seidel: Iterate one coordinate at the time
- Other alternatives
 - Fixed point iterations (can also be gradient based) contraction
 - Nelder – Mead (optim in R)
- Know when to stop

GAUSS-NEWTON revisited

Regression: $\min_{\beta} \sum_{i=1}^n (y_i - \beta^T x_i)^2$

$$\mathbf{X} = \begin{bmatrix} \mathbf{x}_1 \\ \vdots \\ \mathbf{x}_n \end{bmatrix}$$

$$\min_{\beta} (\mathbf{y} - \mathbf{X}\beta)^T (\mathbf{y} - \mathbf{X}\beta)$$

$$\frac{\partial}{\partial \beta} : 2 \cdot \mathbf{X}^T (\mathbf{y} - \mathbf{X}\beta) = 0$$

$$\beta = (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \mathbf{y}$$

$$\min_{\theta} \sum_{i=1}^n (y_i - f(z_i, \theta))^2$$

NL-LS

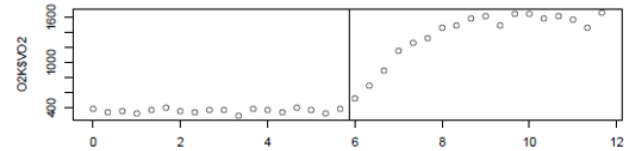
Approx :

$$\min_{\theta} \sum_{i=1}^n \underbrace{(y_i - f(z_i, \theta_k))}_{\text{residual}} - \underbrace{(\theta - \theta_k)^T \nabla f(z_i, \theta_k)}_{\text{linear approx}}^2$$

$$\mathbf{f}(\mathbf{z}; \boldsymbol{\theta}) = \begin{bmatrix} f(\mathbf{z}_1; \boldsymbol{\theta}) \\ \vdots \\ f(\mathbf{z}_n; \boldsymbol{\theta}) \end{bmatrix}$$

$$\mathbf{A}^{(k)} = \begin{bmatrix} \nabla_{\theta} f(\mathbf{z}_1, \theta_k) \\ \vdots \\ \nabla_{\theta} f(\mathbf{z}_n, \theta_k) \end{bmatrix}$$

$$\theta_{k+1} = \theta_k + \left(\mathbf{A}^{(k)T} \mathbf{A}^{(k)} \right)^{-1} \mathbf{A}^{(k)T} (\mathbf{y} - \mathbf{f}(\mathbf{z}; \boldsymbol{\theta}))$$



Solving the problem in R

- All parameters are positive, log-transformed first

```
data(O2K, package="nlstools")
tresh = 5.883
```

```
fit.vo2 = function(param, dat, lambda)
{
  tau1 = exp(param[1])
  beta0 = exp(param[2])
  beta1 = exp(param[3])
  fit = (dat$t <= lambda)*beta0 + (dat$t > lambda) *
    (beta0 + (beta1-beta0)*(1-exp(-(dat$t-lambda) / tau1)))
}
```

Mapping from
constrained to
unconstrained
optimization

```
minusloglik = function(param, dat, tresh=5.883)
{
  fit = fit.vo2(param, dat, tresh)
  sig2 = exp(param[4])
  rss = sum((dat$VO2-fit)^2)
  n = ncol(dat)
  loglik = -0.5*n*log(2*pi) - 0.5*n*log(sig2) - 0.5*rss / sig2
  -loglik
}
```

Function to be optimized (where first input is the parameters to be optimized)

→ res = **optim**(c(1, log(400), log(1600), 0), minusloglik, dat=O2K, tresh=tresh)

Initial guess of parameters
to be estimated

Control variables
(fixed during optimization)

Discussion

- Why is Newton's method particularly suited for optimization of log likelihoods?
- Why do you think the estimator obtained by **one** iteration of Newton's method is asymptotically just as good as the ML?
- Hint: What is the asymptotic result for ML?

Today

- Iteratively reweighted least square
- ADMM
- ADMM for LASSO

- Combinatorial optimization

Iteratively reweighted least square

- Assume you have data with unknown mean and known variable variance and want to estimate mean

$$\min_{\beta} \sum \frac{(y_i - x_i^T \beta_i)^2}{\sigma_i^2} = \min_{\beta} \sum w_i (y_i - x_i^T \beta_i)^2$$

$$\hat{\beta}_{WLS} = (X^T W X)^{-1} X^T W y$$

$$W = \text{diag}(\mathbf{w})$$

- Now: what if variance is unknown and a function of the parameter ?

$$\min_{\beta} \sum w_i(\beta, x_i) (y_i - x_i^T \beta_i)^2$$

- Then: update the weights with previous estimate of parameter.
- Useful trick, can be used for GLM generalized linear models

Weighted least square

$$\min_{\beta} \sum_{i=1}^n w_i (y_i - \beta^T x_i)^2$$

$$\min_{\beta} (\mathbf{y} - \mathbf{X}\beta)^T \mathbf{W} (\mathbf{y} - \mathbf{X}\beta)$$

$$\frac{\partial}{\partial \beta} : 2 \cdot \mathbf{X}^T \mathbf{W} (\mathbf{y} - \mathbf{X}\beta) = 0$$

$$\mathbf{X}^T \mathbf{W} \mathbf{X} \beta = \mathbf{X}^T \mathbf{W} \mathbf{y}$$

$$\beta = (\mathbf{X}^T \mathbf{W} \mathbf{X})^{-1} \mathbf{X}^T \mathbf{W} \mathbf{y}$$

$$\mathbf{W} = \begin{bmatrix} w_1 & 0 & 0 \\ 0 & \ddots & 0 \\ 0 & 0 & w_n \end{bmatrix}$$

Iteratively reweighted least square (linear case)

$$\min_{\boldsymbol{\beta}} \sum w_i(\boldsymbol{\beta}, \mathbf{x}_i) (y_i - \mathbf{x}_i^T \boldsymbol{\beta}_i)^2$$

1. $\mathbf{w}^{(0)} = \mathbf{1}$

2. For $k=1$ «until convergence»

1. $\boldsymbol{\beta}^{(k)} = \operatorname{argmin} \left\{ \sum w_i^{(k-1)} (y_i - \mathbf{x}_i^T \boldsymbol{\beta}_i)^2 \right\}$, *i.e.* $\boldsymbol{\beta}^{(k)} = \widehat{\boldsymbol{\beta}}_{WLS}(\mathbf{W}^{(k-1)})$

2. $w_i^{(k)} = w_i(\boldsymbol{\beta}^{(k)}, \mathbf{x}_i)$

Example: L1-regression robustness to outliers

- Quadratic loss:

$$\arg \min_{\beta} \sum_{i=1}^n (y_i - \beta^T x_i)^2$$

- Absolute loss:

$$\arg \min_{\beta} \sum_{i=1}^n |y_i - \beta^T x_i|$$

$$\sum_{i=1}^n |y_i - \beta^T x_i| = \sum_{i=1}^n w_i(\beta) (y_i - \beta^T x_i)^2 \quad \boxed{w_i(\beta) = \frac{1}{|y_i - \beta^T x_i|}}$$

If « β is known» we can do weighted least squares regression

* Start with $w_i^{(0)} = 1$. (= least squares regression) to get $\beta^{(0)}$

* In iteration k set $w_i^{(k)} = \frac{1}{|y_i - \beta^{(k-1)T} x_i|}$ or $\min \left\{ \frac{1}{|y_i - \beta^{(k-1)T} x_i|}, W_{\max} \right\}$

IRLS.R

```
IRLS_L1 = function(x,y)
{
  betaHat = solve(t(x)%*%x) %*%(t(x)%*%y )
  pred0 = x%*%betaHat
  res0 = (y-pred)

  betaPrev=c(0,0)
  betaWHat=betaHat
  it=0

  while(sum(abs(betaPrev-betaWHat))>0.0001 & it<100)
  {
    maxw=10
    it=it+1
    betaPrev=betaWHat
    pred= xdata%*%betaPrev
    res=(y-pred)

    w = 1/abs(res)
    w[w>maxw]=maxw # adjustment to avoid super large numbers [ size relative to problem]
    w = diag(as.vector(w))

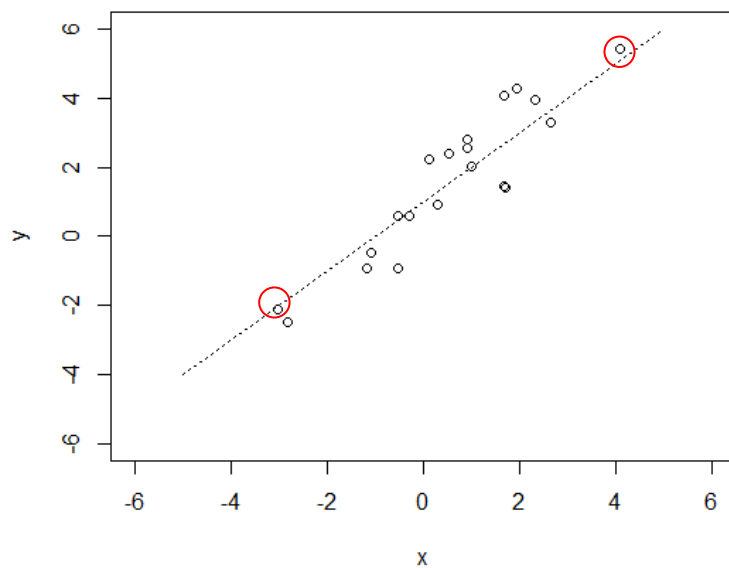
    betaWHat = solve(t(xdata)%*%w%*%xdata) %*%(t(xdata)%*%w%*%y )
    #show(betaWHat)
  }
  betaWHat
}
```

$$w_i^{(k)} = \min \left\{ \frac{1}{|y_i - \beta^{(k-1)T} x_i|}, W_{\max} \right\}$$

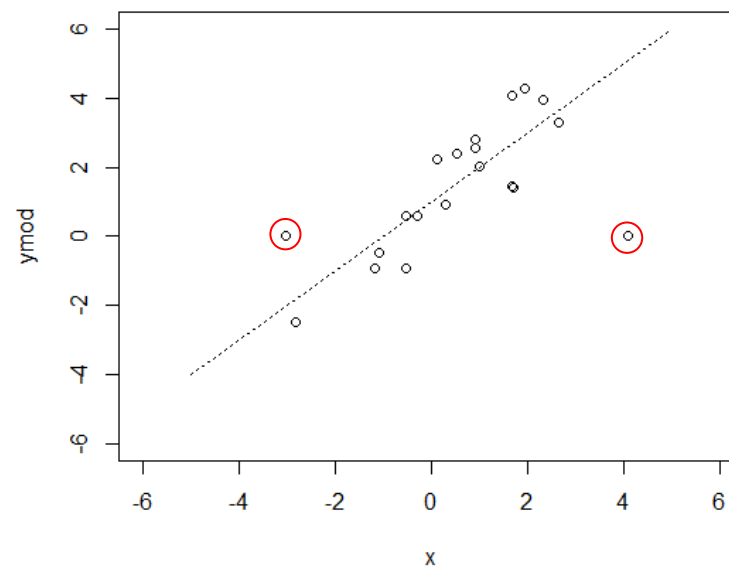
$$\beta = (X^T W X)^{-1} X^T W y$$

Example n=20

Case 1



Case 2



Example n=20

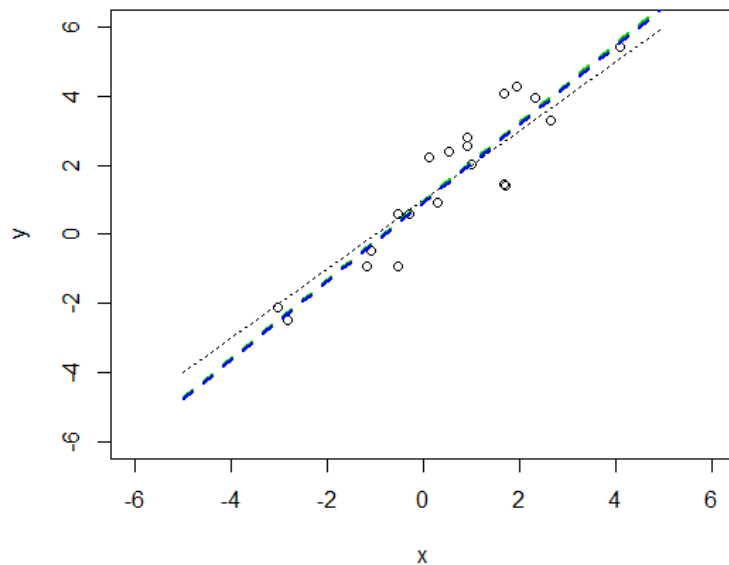
Quadratic loss - - - - -

Absolute loss - - - - -

```
> betaHat_L2  
[ ,1]  
0.9504518  
x 1.1409776
```

```
> betaHat_L1  
[ ,1]  
0.8972019  
x 1.1343827
```

Case 1

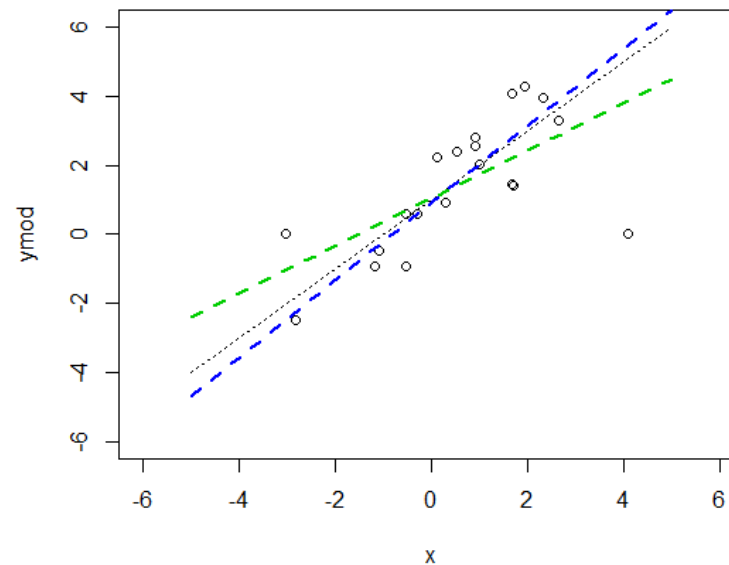


Results almost identical

```
> betaHat_L2m  
[ ,1]  
1.0240550  
x 0.6887221
```

```
> betaHat_L1m  
[ ,1]  
0.9009297  
x 1.1211309
```

Case 2



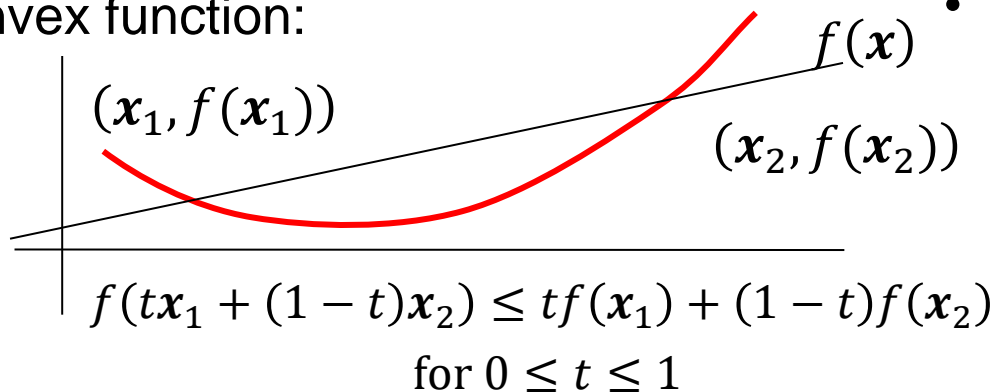
Quadratic loss is sensitive to outliers

Convex optimization

- minimize $_x$ $\{ f(\mathbf{x}) \}$
subject to $\mathbf{a}_i^T \mathbf{x} = b_i$
and $g_j(\mathbf{x}) \leq 0$
- $i = 1, \dots, n_a$
- $j = 1, \dots, n_c$
- $f(\mathbf{x}), g_j(\mathbf{x})$ convex

https://web.stanford.edu/~boyd/cvxbook/bv_cvxbook.pdf

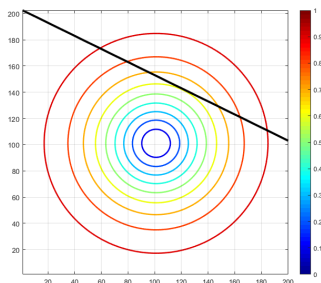
Convex function:



Interpolation always above function

- A field on its own,...
- Covers many important problems
 - Least squares
 - Linear programming
 - Convex quadratic minimization with linear or convex quadratic constraints
 - Conic optimization
 - Second order cone programming
 - ...
- Many algorithms open source
 - Bundle methods
 - Subgradient projection methods (Polyak),
 - Interior-point methods
 - Cutting-plane methods
 - Ellipsoid method
 - Subgradient method
 - Dual subgradients
 - ...

Minimization under an equality constraint and the augmented Lagrangian multiplier

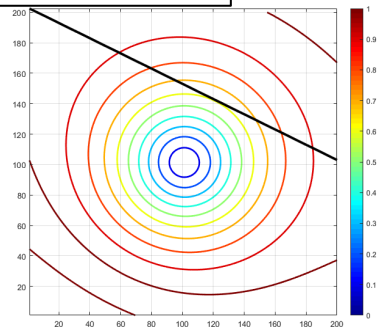


$\rho: 1 \times 1$
is a scalar

$$\begin{aligned} &\text{minimize}_x \quad \{ f(x) \} \\ &\text{subject to} \quad Ax = b \end{aligned}$$

$x: (p \times 1)$
 $b: (q \times 1)$
 $A: (q \times p)$

$$\begin{aligned} &\Leftrightarrow \\ &\text{minimize}_x \quad \left\{ f(x) + \frac{\rho}{2} \|Ax - b\|^2 \right\} \\ &\text{subject to} \quad Ax = b \end{aligned}$$

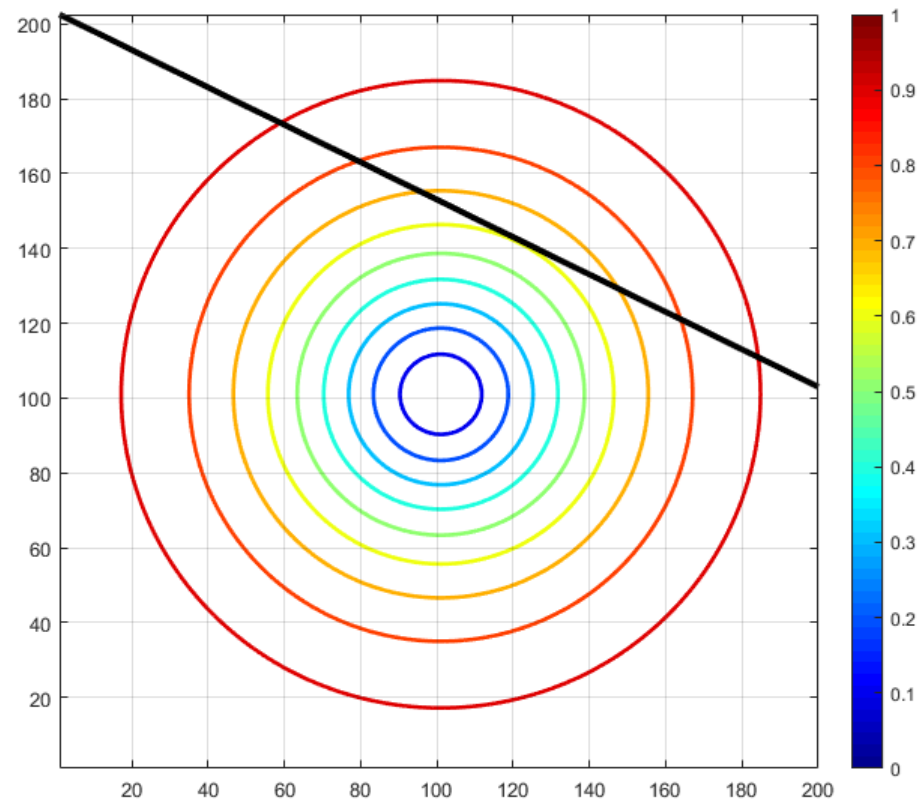


$$\begin{aligned} &\Leftrightarrow \\ &\text{minimize}_{x,\lambda} \quad \left\{ f(x) + \frac{\rho}{2} \|Ax - b\|^2 + \lambda^T (Ax - b) \right\} \end{aligned}$$

$\lambda : (q \times 1)$ is the
Lagrangian multiplier

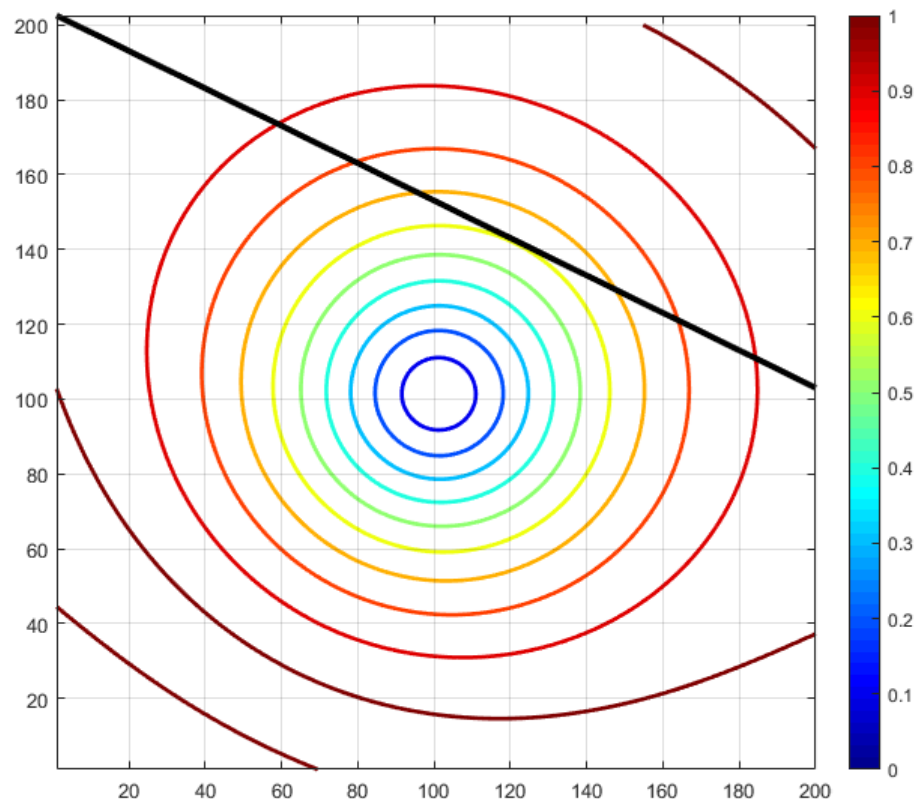
Geometric interpretation:
Minimum occurs
when the contour
line is tangent to
constraint

minimize $\{ f(\mathbf{x}) \}$
subject to $A\mathbf{x} = \mathbf{b}$



$$\begin{aligned} & \text{minimize}_{\mathbf{x}, \lambda} \quad \left\{ f(\mathbf{x}) + \frac{\rho}{2} \|\mathbf{Ax} - \mathbf{b}\|^2 \right\} \\ & \text{subject to} \quad \mathbf{Ax} = \mathbf{b} \end{aligned}$$

Does not
change
along the
black line



Algorithm for solution: Method of multipliers

$$\begin{array}{ll} \text{minimize}_{\mathbf{x}, \lambda} & \left\{ f(\mathbf{x}) + \frac{\rho}{2} \|\mathbf{Ax} - \mathbf{b}\|^2 \right\} \\ \text{subject to} & \mathbf{Ax} = \mathbf{b} \end{array}$$

$$1. \quad \boldsymbol{\lambda}^{(0)} = \mathbf{0}$$

2. For $i=1$ «until convergence»

$$1. \quad \mathbf{x}^{(i)} = \operatorname{argmin} \left\{ f(\mathbf{x}) + \frac{\rho}{2} \|\mathbf{Ax} - \mathbf{b}\|^2 + \boldsymbol{\lambda}^{(i-1)T} (\mathbf{Ax} - \mathbf{b}) \right\}$$

$$2. \quad \boldsymbol{\lambda}^{(i)} = \boldsymbol{\lambda}^{(i-1)} + \rho (\mathbf{Ax}^{(i)} - \mathbf{b})$$

A version of **dual ascent**, details given in reference for the clever and mathematically inclined student (part of the article is syllabus for STK9051)

ADMM Alternating Direction Method of Multipliers

$$\begin{aligned} & \text{minimize} && \{ f(\mathbf{x}) + g(\mathbf{z}) \} \\ & \text{subject to} && \mathbf{Ax} + \mathbf{Bz} = \mathbf{c} \end{aligned}$$

$$\text{minimize} \left\{ f(\mathbf{x}) + g(\mathbf{z}) + \lambda^T (\mathbf{Ax} + \mathbf{Bz} - \mathbf{c}) + \frac{\rho}{2} \|\mathbf{Ax} + \mathbf{Bz} - \mathbf{c}\|^2 \right\}$$

$$1. \quad \lambda^{(0)} = \mathbf{0}, \mathbf{z}^{(0)} = \mathbf{0}$$

2. For $i=1$ «until convergence»

$$1. \quad \mathbf{x}^{(i)} = \operatorname{argmin} \left\{ f(\mathbf{x}) + \frac{\rho}{2} \|\mathbf{Ax} + \mathbf{Bz}^{(i-1)} - \mathbf{c}\|^2 + \lambda^{(i-1)T} (\mathbf{Ax} + \mathbf{Bz}^{(i-1)} - \mathbf{c}) \right\}$$

$$2. \quad \mathbf{z}^{(i)} = \operatorname{argmin} \left\{ g(\mathbf{z}) + \frac{\rho}{2} \|\mathbf{Ax}^{(i)} + \mathbf{Bz} - \mathbf{c}\|^2 + \lambda^{(i-1)T} (\mathbf{Ax} + \mathbf{Bz} - \mathbf{c}) \right\}$$

$$3. \quad \lambda^{(i)} = \lambda^{(i-1)} + \rho (\mathbf{Ax}^{(i)} + \mathbf{Bz}^{(i)} - \mathbf{c})$$

Alternate between solving \mathbf{x}
solving \mathbf{z} and update λ

We do not need to deal with
 $f(\mathbf{x})$ and $g(\mathbf{z})$ simultaneously!

- Proven convergence under general assumptions
- Ideal for splitting problem into smaller «solvable» problems
- Syllabus ADMM
 - No proofs
 - Operational use
 - Questions about related stuff

Reference

Foundations and Trends® in
Machine Learning
Vol. 3, No. 1 (2010) 1–122
© 2011 S. Boyd, N. Parikh, E. Chu, B. Peleato
and J. Eckstein
DOI: 10.1561/22000000016



Distributed Optimization and Statistical Learning via the Alternating Direction Method of Multipliers

Stephen Boyd¹, Neal Parikh², Eric Chu³
Borja Peleato⁴ and Jonathan Eckstein⁵

Available online

https://web.stanford.edu/~boyd/papers/admm_distr_stats.html

Solving lasso using ADMM

- Lasso: minimize $\left\{ \frac{1}{2} \|X\boldsymbol{\beta} - \mathbf{y}\|^2 + \gamma \|\boldsymbol{\beta}\|_1 \right\}$
 \Updownarrow
 minimize $\left\{ \frac{1}{2} \|X\boldsymbol{\beta} - \mathbf{y}\|^2 + \gamma \|\mathbf{z}\|_1 \right\}$
 subject to $\boldsymbol{\beta} - \mathbf{z} = \mathbf{0}$

The factor $\frac{1}{2}$
is used to simplify
expressions

In ADMM

$$1. \quad \boldsymbol{\beta}^{(i)} = \operatorname{argmin} \left\{ \frac{1}{2} \|X\boldsymbol{\beta} - \mathbf{y}\|^2 + \frac{\rho}{2} \|\boldsymbol{\beta} - \mathbf{z}^{(i-1)}\|^2 + \boldsymbol{\lambda}^{(i-1)T} (\boldsymbol{\beta} - \mathbf{z}^{(i-1)}) \right\}$$

First problem has only quadratic and linear terms, thus it is solved by linear equations

$$2. \quad \mathbf{z}^{(i)} = \operatorname{argmin} \left\{ \gamma \|\mathbf{z}\|_1 + \frac{\rho}{2} \|\boldsymbol{\beta}^{(i)} - \mathbf{z}\|^2 + \boldsymbol{\lambda}^{(i-1)T} (\boldsymbol{\beta}^{(i)} - \mathbf{z}) \right\}$$

Second problem is non linear but can be solved for each component separately

First problem

$$\operatorname{argmin}_{\boldsymbol{\beta}} \left\{ \frac{1}{2} \|\mathbf{X}\boldsymbol{\beta} - \mathbf{y}\|^2 + \frac{\rho}{2} \|\boldsymbol{\beta} - \mathbf{z}^{(i-1)}\|^2 + \boldsymbol{\lambda}^{(i-1)T} (\boldsymbol{\beta} - \mathbf{z}^{(i-1)}) \right\}$$

$$(\mathbf{X}\boldsymbol{\beta} - \mathbf{y})^T (\mathbf{X}\boldsymbol{\beta} - \mathbf{y}) + \frac{\rho}{2} (\boldsymbol{\beta} - \mathbf{z}^{(i-1)})^T (\boldsymbol{\beta} - \mathbf{z}^{(i-1)}) + \boldsymbol{\lambda}^{(i-1)T} (\boldsymbol{\beta} - \mathbf{z}^{(i-1)})$$

$$\frac{\partial}{\partial \boldsymbol{\beta}}: \quad \mathbf{X}^T (\mathbf{X}\boldsymbol{\beta} - \mathbf{y}) + \rho (\boldsymbol{\beta} - \mathbf{z}^{(i-1)}) + \boldsymbol{\lambda}^{(i-1)} = 0$$

$$(\mathbf{X}^T \mathbf{X} + \rho \mathbf{I}) \boldsymbol{\beta} = \mathbf{X}^T \mathbf{y} + \rho \left(\mathbf{z}^{(i-1)} - \frac{\boldsymbol{\lambda}^{(i-1)}}{\rho} \right)$$

$$\boldsymbol{\beta} = (\mathbf{X}^T \mathbf{X} + \rho \mathbf{I})^{-1} \left(\mathbf{X}^T \mathbf{y} + \rho \left(\mathbf{z}^{(i-1)} - \frac{\boldsymbol{\lambda}^{(i-1)}}{\rho} \right) \right)$$

Second problem

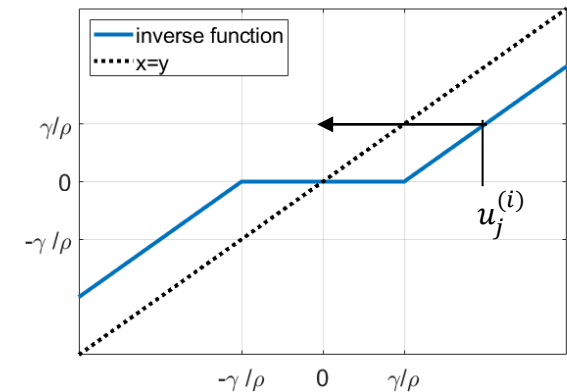
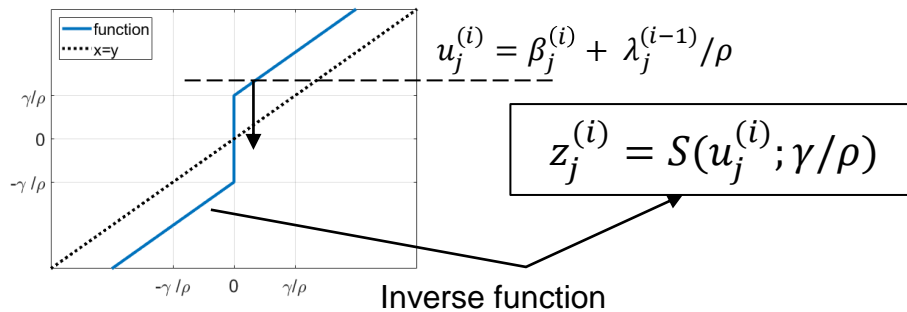
$$\operatorname{argmin}_{\mathbf{z}} \left\{ \gamma \|\mathbf{z}\|_1 + \frac{\rho}{2} \|\boldsymbol{\beta}^{(i)} - \mathbf{z}\|^2 + \boldsymbol{\lambda}^{(i-1)T} (\boldsymbol{\beta}^{(i)} - \mathbf{z}) \right\}$$

$$\sum_{j=1}^p \left(\gamma |z_j| + \frac{\rho}{2} (\beta_j^{(i)} - z_j)^2 + \lambda_j^{(i-1)} (\beta_j^{(i)} - z_j) \right)$$

$$\frac{\partial}{\partial z_j} : \quad \gamma \operatorname{sign}(z_j) - \rho (\beta_j^{(i)} - z_j) - \lambda_j^{(i-1)} = 0$$

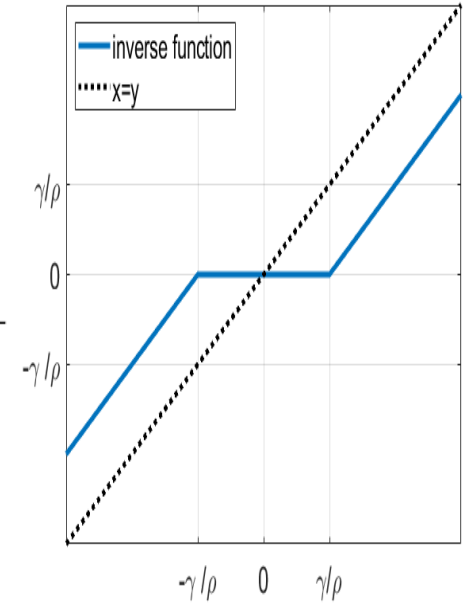
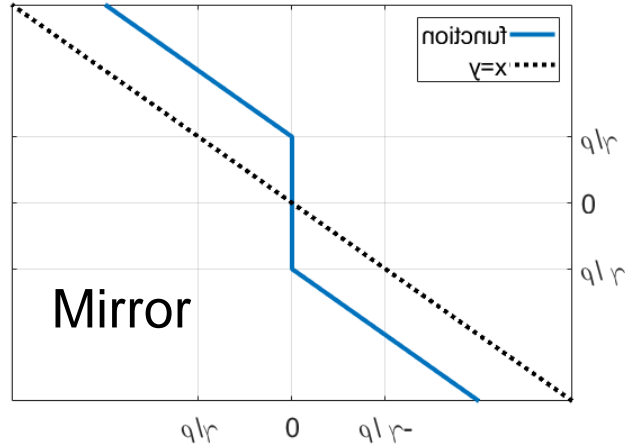
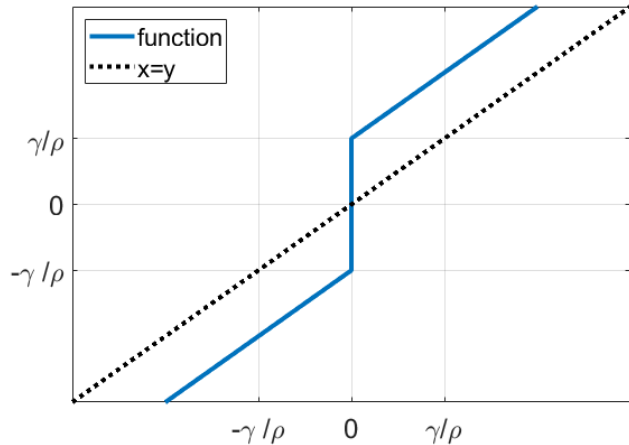
$$\left(\frac{\gamma}{\rho} \cdot \operatorname{sign}(z_j) + z_j \right) = \beta_j^{(i)} + \lambda_j^{(i-1)} / \rho$$

Solve for each component independently



$$S(u; \gamma/\rho) = \operatorname{sign}(u) \cdot \max(|u| - \gamma/\rho, 0)$$

Inverting a function in powerpoint



$$\frac{\gamma}{\rho} \cdot \text{sign}(z) + z$$

$$\text{sign}(u) \cdot \max(|u| - \gamma/\rho, 0)$$

Baseball example R