

Notes on L1 Regularized Logistic Regression

Max Leung

November 22, 2022

Disclaimer: This note can contain mistake, typo, error, etc.

1 Subgradient Optimality

If f is a convex function of x (not necessarily differentiable) and $\text{dom}(f) = R^n$, we have

$$x^* = \text{argmin}_x f(x) \iff 0 \in \partial f(x^*)$$

Proof: (\Leftarrow) $0 \in \partial f(x^*) := \{g \in R^n : f(y) \geq f(x^*) + g'(y - x) \text{ for } \forall y\} \Rightarrow f(y) \geq f(x^*) + 0'(y - x) \text{ for } \forall y \Rightarrow f(y) \geq f(x^*) \text{ for } \forall y \Rightarrow x^* \text{ is minimizer.}$

(\Rightarrow) $x^* = \text{argmin}_x f(x) \Rightarrow f(y) \geq f(x^*) \text{ for } \forall y \Rightarrow f(y) \geq f(x^*) + 0'(y - x) \text{ for } \forall y \Rightarrow 0 \in \partial f(x^*)$ Q.E.D.

2 Logistic Regression Without Constraint

2.1 Log Likelihood Function

$$\ln L(\beta; \mathbf{y}, \mathbf{X}) = \ln \prod_{i=1}^N \Pr(y_i | \mathbf{x}_i; \beta) = \sum_{i=1}^N \ln \Pr(y_i | \mathbf{x}_i; \beta) \quad \text{assume independence of } y_i | \mathbf{x}_i$$

$$\begin{aligned} \Pr(y_i = 1 | \mathbf{x}_i; \beta) &= \text{logistic}(\mathbf{x}_i' \beta) := \frac{e^{\mathbf{x}_i' \beta}}{1 + e^{\mathbf{x}_i' \beta}} \text{ and } \Pr(y_i = 0 | \mathbf{x}_i; \beta) = 1 - \Pr(y_i = 1 | \mathbf{x}_i; \beta) = 1 - \frac{e^{\mathbf{x}_i' \beta}}{1 + e^{\mathbf{x}_i' \beta}} = \frac{1}{1 + e^{\mathbf{x}_i' \beta}} = \frac{e^0}{1 + e^{\mathbf{x}_i' \beta}} \\ &= \sum_{i=1}^N \ln \frac{e^{\mathbf{x}_i' \beta y_i}}{1 + e^{\mathbf{x}_i' \beta}} \\ &= \sum_{i=1}^N \{\mathbf{x}_i' \beta y_i - \ln(1 + e^{\mathbf{x}_i' \beta})\} \end{aligned}$$

2.2 Gradient Vector

$$\begin{aligned} \frac{\partial \ln L(\beta; \mathbf{y}, \mathbf{X})}{\partial \beta} &= \frac{\partial \sum_{i=1}^N \{\mathbf{x}_i' \beta y_i - \ln(1 + e^{\mathbf{x}_i' \beta})\}}{\partial \beta} \\ &= \sum_{i=1}^N \left\{ y_i \frac{\partial \mathbf{x}_i' \beta}{\partial \beta} - \frac{\partial \ln(1 + e^{\mathbf{x}_i' \beta})}{\partial \beta} \right\} \\ &= \sum_{i=1}^N \left\{ y_i \mathbf{x}_i - \frac{1}{1 + e^{\mathbf{x}_i' \beta}} e^{\mathbf{x}_i' \beta} \mathbf{x}_i \right\} \\ &= \sum_{i=1}^N \{y_i \mathbf{x}_i - \text{logistic}(\mathbf{x}_i' \beta) \mathbf{x}_i\} \\ &= \sum_{i=1}^N \{y_i - \underbrace{\Pr(y_i = 1 | \mathbf{x}_i; \beta)}_{p_i}\} \mathbf{x}_i \quad \{y_i - \Pr(y_i = 1 | \mathbf{x}_i; \beta)\} \mathbf{x}_i \text{ is score function} \\ &= \sum_{i=1}^N \begin{pmatrix} (y_i - p_i)x_{1i} \\ \vdots \\ (y_i - p_i)x_{pi} \end{pmatrix} = \begin{pmatrix} \sum_{i=1}^N (y_i - p_i)x_{1i} \\ \vdots \\ \sum_{i=1}^N (y_i - p_i)x_{pi} \end{pmatrix} = \begin{pmatrix} \mathbf{x}'_1(\mathbf{y} - \mathbf{p}) \\ \vdots \\ \mathbf{x}'_p(\mathbf{y} - \mathbf{p}) \end{pmatrix} = \begin{pmatrix} \mathbf{x}'_1 \\ \vdots \\ \mathbf{x}'_p \end{pmatrix} (\mathbf{y} - \mathbf{p}) = \mathbf{X}'(\mathbf{y} - \mathbf{p}) \end{aligned}$$

2.3 Hessian Matrix

$$\begin{aligned}
\frac{\partial^2 \ln L}{\partial \beta \partial \beta'} &= \frac{\partial \sum_{i=1}^N \{y_i \mathbf{x}_i - \frac{1}{1+e^{\mathbf{x}_i' \beta}} e^{\mathbf{x}_i' \beta} \mathbf{x}_i\}}{\partial \beta'} \\
&= \sum_{i=1}^N \frac{\partial \{y_i \mathbf{x}_i - \frac{1}{1+e^{\mathbf{x}_i' \beta}} e^{\mathbf{x}_i' \beta} \mathbf{x}_i\}}{\partial \beta'} \\
&= - \sum_{i=1}^N \mathbf{x}_i \frac{\partial (1 + e^{-\mathbf{x}_i' \beta})^{-1}}{\partial \beta'} \\
&= - \sum_{i=1}^N \mathbf{x}_i [-(1 + e^{-\mathbf{x}_i' \beta})^{-2} e^{-\mathbf{x}_i' \beta} (-1) \mathbf{x}_i'] \\
&= - \sum_{i=1}^N \mathbf{x}_i (1 + e^{-\mathbf{x}_i' \beta})^{-1} \frac{e^{-\mathbf{x}_i' \beta}}{1 + e^{-\mathbf{x}_i' \beta}} \mathbf{x}_i' \\
&= - \sum_{i=1}^N \mathbf{x}_i \text{logistic}(\mathbf{x}_i' \beta) \frac{1}{e^{\mathbf{x}_i' \beta} + 1} \mathbf{x}_i' \\
&= - \sum_{i=1}^N \mathbf{x}_i \text{logistic}(\mathbf{x}_i' \beta) \frac{1 + e^{\mathbf{x}_i' \beta} - e^{\mathbf{x}_i' \beta}}{1 + e^{\mathbf{x}_i' \beta}} \mathbf{x}_i' \\
&= - \sum_{i=1}^N \mathbf{x}_i \text{logistic}(\mathbf{x}_i' \beta) (1 - \text{logistic}(\mathbf{x}_i' \beta)) \mathbf{x}_i' \\
&= - \sum_{i=1}^N \mathbf{x}_i \mathbf{x}_i' \overbrace{\text{Pr}(y_i = 1 | \mathbf{x}_i; \beta)}^{p_i} (1 - \text{Pr}(y_i = 1 | \mathbf{x}_i; \beta)) \\
&= - \sum_{i=1}^N \begin{pmatrix} x_{1i} x_{1i} p_i (1 - p_i) & \cdots & x_{1i} x_{pi} p_i (1 - p_i) \\ \vdots & \ddots & \vdots \\ x_{pi} x_{1i} p_i (1 - p_i) & \cdots & x_{pi} x_{pi} p_i (1 - p_i) \end{pmatrix} \\
&= - \begin{pmatrix} \mathbf{x}_1' \mathbf{W} \mathbf{x}_1 & \cdots & \mathbf{x}_1' \mathbf{W} \mathbf{x}_p \\ \vdots & \ddots & \vdots \\ \mathbf{x}_p' \mathbf{W} \mathbf{x}_1 & \cdots & \mathbf{x}_p' \mathbf{W} \mathbf{x}_p \end{pmatrix} \quad \text{where } (\mathbf{W})_{ii} = p_i(1 - p_i) \text{ and } (\mathbf{W})_{ij} = 0 \text{ for } \forall i \neq j \\
&= - \begin{pmatrix} \mathbf{x}_1' \\ \vdots \\ \mathbf{x}_p' \end{pmatrix} \mathbf{W} (\mathbf{x}_1 \cdots \mathbf{x}_p) = -\mathbf{X}' \mathbf{W} \mathbf{X}
\end{aligned}$$

2.4 Iteratively Reweighted Least Squares (IRLS)

Note that Newton Method is $\mathbf{x}^+ = \mathbf{x} - (\nabla^2 f(\mathbf{x}))^{-1} \nabla f(\mathbf{x})$. If we apply it here:

$$\begin{aligned}
\beta^+ &= \beta - (\nabla^2 \ln L(\beta))^{-1} \nabla \ln L(\beta) \\
&= \beta - (-\mathbf{X}' \mathbf{W} \mathbf{X})^{-1} \mathbf{X}' (\mathbf{y} - \mathbf{p}) \quad \text{using above results} \\
&= \mathbf{I} \beta + (\mathbf{X}' \mathbf{W} \mathbf{X})^{-1} \mathbf{X}' \mathbf{I} (\mathbf{y} - \mathbf{p}) \\
&= (\mathbf{X}' \mathbf{W} \mathbf{X})^{-1} (\mathbf{X}' \mathbf{W} \mathbf{X}) \beta + (\mathbf{X}' \mathbf{W} \mathbf{X})^{-1} \mathbf{X}' \mathbf{W} \mathbf{W}^{-1} (\mathbf{y} - \mathbf{p}) \\
&= (\mathbf{X}' \mathbf{W} \mathbf{X})^{-1} \mathbf{X}' \mathbf{W} \underbrace{(\mathbf{X} \beta + \mathbf{W}^{-1} (\mathbf{y} - \mathbf{p}))}_z
\end{aligned}$$

Thus, β^+ is a Weighted Least Squares estimator i.e.,

$$\beta^+ = \underset{\beta}{\operatorname{argmin}} (\mathbf{z} - \mathbf{X} \beta)' \mathbf{W} (\mathbf{z} - \mathbf{X} \beta) = \underset{\beta}{\operatorname{argmin}} \sum_{i=1}^N w_i (z_i - \mathbf{x}_i' \beta)^2 \quad \text{where } w_i = p_i(1 - p_i)$$

We loop $\beta^+ = (\mathbf{X}' \mathbf{W} \mathbf{X})^{-1} \mathbf{X}' \mathbf{W} \mathbf{z}$ until β converge. Note that \mathbf{W} and \mathbf{z} is a function of β . Thus, minimizing Unpenalized Logistic Regression's negative log likelihood function is the same as repeatedly minimizing weighted least squares functions.

3 L1 Regularized Logistic Regression

3.1 Minimize Negative Log Likelihood Function with L1 Constraint

$$\begin{aligned} & \min_{\beta} \{-\ln L(\beta; \mathbf{y}, \mathbf{X}) + \lambda \|\beta\|_1\} \\ \iff & \min_{\beta} \left\{ -\sum_{i=1}^N (\mathbf{x}'_i \beta y_i - \ln(1 + e^{\mathbf{x}'_i \beta})) + \lambda \|\beta\|_1 \right\} \end{aligned}$$

3.2 Apply Subgradient Optimality

As $-\sum_{i=1}^N (\mathbf{x}'_i \beta y_i - \ln(1 + e^{\mathbf{x}'_i \beta})) + \lambda \|\beta\|_1$ is a convex function of β , Subgradient Optimality implies that

$$\beta^* = \operatorname{argmin}_{\beta} \left\{ -\sum_{i=1}^N (\mathbf{x}'_i \beta y_i - \ln(1 + e^{\mathbf{x}'_i \beta})) + \lambda \|\beta\|_1 \right\} \iff 0 \in \partial \left\{ -\sum_{i=1}^N (\mathbf{x}'_i \beta^* y_i - \ln(1 + e^{\mathbf{x}'_i \beta^*})) + \lambda \|\beta^*\|_1 \right\}$$

$$\begin{aligned} 0 & \in \partial \left\{ -\sum_{i=1}^N (\mathbf{x}'_i \beta^* y_i - \ln(1 + e^{\mathbf{x}'_i \beta^*})) + \lambda \|\beta^*\|_1 \right\} \\ & \in -\partial \sum_{i=1}^N (\mathbf{x}'_i \beta^* y_i - \ln(1 + e^{\mathbf{x}'_i \beta^*})) + \lambda \partial \|\beta^*\|_1 \\ & \in -\mathbf{X}'(\mathbf{y} - \mathbf{p}(\beta^*)) + \lambda \partial \|\beta^*\|_1 && \text{Apply Gradient Vector result above} \\ \mathbf{X}'(\mathbf{y} - \mathbf{p}(\beta^*)) & \in \lambda \partial \|\beta^*\|_1 \end{aligned}$$

We want to pick $\mathbf{v} \in \partial \|\beta^*\|_1$ such that $\mathbf{X}'(\mathbf{y} - \mathbf{p}(\beta^*)) = \lambda \mathbf{v}$

$$\begin{aligned} \mathbf{X}'(\mathbf{y} - \mathbf{p}(\beta^*)) &= \lambda \mathbf{v} \\ \begin{pmatrix} \mathbf{x}'_1 \\ \vdots \\ \mathbf{x}'_p \end{pmatrix} (\mathbf{y} - \mathbf{p}(\beta^*)) &= \lambda \mathbf{v} \\ \begin{pmatrix} \mathbf{x}'_1(\mathbf{y} - \mathbf{p}(\beta^*)) \\ \vdots \\ \mathbf{x}'_p(\mathbf{y} - \mathbf{p}(\beta^*)) \end{pmatrix} &= \lambda \begin{pmatrix} v_1 \\ \vdots \\ v_p \end{pmatrix} \end{aligned}$$

So, we have $\mathbf{x}'_j(\mathbf{y} - \mathbf{p}(\beta^*)) = \lambda v_j$ for $\forall j \in \{1, \dots, p\}$

As $\mathbf{v} \in \partial \|\beta^*\|_1$, we have

$$v_j \in \partial |\beta_j^*| = \begin{cases} \{\operatorname{sign}(\beta_j^*)\} & \text{if } \beta_j^* \neq 0 \\ [-1, 1] & \text{if } \beta_j^* = 0 \end{cases}$$

if $\beta_j^* \neq 0$

$$v_j \in \partial |\beta_j^*| = \{\operatorname{sign}(\beta_j^*)\}$$

that is

$$v_j = \operatorname{sign}(\beta_j^*) \quad \quad \quad \{\operatorname{sign}(\beta_j^*)\} \text{ is singleton}$$

Thus, we have

$$\mathbf{x}'_j(\mathbf{y} - \mathbf{p}(\beta^*)) = \lambda \operatorname{sign}(\beta_j^*)$$

This is Equation 4.32 on page 126 of the famous book The Elements of Statistical Learning (2nd ed.)

4 Algorithm for Solving L1 Regularized Logistic Regression

4.1 Proximal-Newton Iterative Approach in R package glmnet (Friedman et al., 2010)

As mentioned in Section 2.4, Unpenalized Logistic Regression Problem can be solved by repeatedly solving Weighted Least Squares Problems with working response $\mathbf{z} = \mathbf{X}\boldsymbol{\beta} + \mathbf{W}^{-1}(\mathbf{y} - \mathbf{p})$ and weight matrix \mathbf{W} where $(\mathbf{W})_{ii} = p_i(1 - p_i)$ and $(\mathbf{W})_{ij} = 0$ for $\forall i \neq j$. Similarly, solving

$$\min_{\boldsymbol{\beta}} \left\{ -\sum_{i=1}^N (\mathbf{x}'_i \boldsymbol{\beta} y_i - \ln(1 + e^{\mathbf{x}'_i \boldsymbol{\beta}})) + \lambda \|\boldsymbol{\beta}\|_1 \right\} \quad \text{L1 Regularized Logistic Regression Problem}$$

is equivalent as repeatedly solving

$$\min_{\boldsymbol{\beta}} \left\{ \sum_{i=1}^N w_i (z_i - \mathbf{x}'_i \boldsymbol{\beta})^2 + \lambda \|\boldsymbol{\beta}\|_1 \right\} \quad \text{Weighted Lasso Regression Problem}$$

with $\tilde{\boldsymbol{\beta}}$ from last step, where $z_i = \mathbf{x}'_i \tilde{\boldsymbol{\beta}} + \frac{y_i - p_i}{w_i}$ and $w_i = p_i(1 - p_i)$ and $p_i = \widetilde{Pr(y_i = 1 | \mathbf{x}_i; \boldsymbol{\beta})} = \text{logistic}(\mathbf{x}'_i \tilde{\boldsymbol{\beta}})$. Here, negative log likelihood function is locally and quadratically approximated by weighted least squares function.

The well-known solution for Weighted Lasso Regression is (I may write a note on this if I have time):

$$\begin{aligned} \beta_j^* &= S_{\lambda / \mathbf{x}'_j \mathbf{W} \mathbf{x}_j} \left(\frac{\mathbf{x}'_j \mathbf{W} (\mathbf{z} - \overbrace{\mathbf{X}_{-j} \boldsymbol{\beta}_{-j}^*}^{\mathbf{r}^{(j)}})}{\mathbf{x}'_j \mathbf{W} \mathbf{x}_j} \right) & \mathbf{x}_j \text{ is column } j \text{ of } \mathbf{X} \\ &= S_{\lambda / \langle \mathbf{x}_j, \mathbf{x}_j \rangle_w} \left(\frac{\langle \mathbf{x}_j, \mathbf{r}^{(j)} \rangle_w}{\langle \mathbf{x}_j, \mathbf{x}_j \rangle_w} \right) & \langle \cdot, \cdot \rangle_w \text{ is an inner product with weight} \\ &= \text{sign} \left(\frac{\langle \mathbf{x}_j, \mathbf{r}^{(j)} \rangle_w}{\langle \mathbf{x}_j, \mathbf{x}_j \rangle_w} \right) \left(\left| \frac{\langle \mathbf{x}_j, \mathbf{r}^{(j)} \rangle_w}{\langle \mathbf{x}_j, \mathbf{x}_j \rangle_w} \right| - \lambda / \langle \mathbf{x}_j, \mathbf{x}_j \rangle_w \right)_+ \end{aligned}$$

where $S_{\lambda / \langle \mathbf{x}_j, \mathbf{x}_j \rangle_w}(\cdot)$ is soft-threshold operator with parameter $\lambda / \langle \mathbf{x}_j, \mathbf{x}_j \rangle_w$

As $(z)_+ = 0$ if $z \leq 0$, It is obvious that

$$\beta_j^* = \begin{cases} 0 & \text{if } \left| \frac{\langle \mathbf{x}_j, \mathbf{r}^{(j)} \rangle_w}{\langle \mathbf{x}_j, \mathbf{x}_j \rangle_w} \right| \leq \lambda / \langle \mathbf{x}_j, \mathbf{x}_j \rangle_w \\ S_{\lambda / \langle \mathbf{x}_j, \mathbf{x}_j \rangle_w} \left(\frac{\langle \mathbf{x}_j, \mathbf{r}^{(j)} \rangle_w}{\langle \mathbf{x}_j, \mathbf{x}_j \rangle_w} \right) & \text{if } \left| \frac{\langle \mathbf{x}_j, \mathbf{r}^{(j)} \rangle_w}{\langle \mathbf{x}_j, \mathbf{x}_j \rangle_w} \right| > \lambda / \langle \mathbf{x}_j, \mathbf{x}_j \rangle_w \end{cases}$$

If variables are standardized such that $\langle \mathbf{x}_j, \mathbf{x}_j \rangle_w = 1$, we have a special case

$$\beta_j^* = \begin{cases} 0 & \text{if } |\langle \mathbf{x}_j, \mathbf{r}^{(j)} \rangle_w| \leq \lambda \\ S_{\lambda}(\langle \mathbf{x}_j, \mathbf{r}^{(j)} \rangle_w) & \text{if } |\langle \mathbf{x}_j, \mathbf{r}^{(j)} \rangle_w| > \lambda \end{cases}$$

This is Equation 16.15 on page 315 of the book “Computer Age Statistical Inference”.

Given λ and β_l^* for $l \in \{1, \dots, p\} \setminus \{j\}$, β_j^* is computed by firstly computing the weighted inner product $\langle \mathbf{x}_j, \mathbf{r}^{(j)} \rangle_w$ and then put it in soft-threshold operator $S_{\lambda}(\cdot)$. The process repeats for $\forall j \in \{1, \dots, p\}$ again and again until convergence. It is called Coordinate Descent, which is fast as for each step it essentially only requires calculating an inner product and performing soft-threshold operation.

The remaining question is how to determine λ . In this algorithm, we pre-determine 100 λ s:

$$\lambda_{max} := \lambda_1 > \lambda_2 > \dots > \lambda_{100} = \epsilon \lambda_{max} > 0$$

where $0 < \epsilon < 1$

We choose λ_{max} such that $\beta_j^* = 0$ for $\forall j \in \{1, \dots, p\}$. This can be achieved by $|\langle \mathbf{x}_j, \mathbf{r}^{(j)} \rangle_w| \leq \lambda_{max}$ for $\forall j \in \{1, \dots, p\}$ i.e., λ_{max} is the upper bound of the set $\{|\langle \mathbf{x}_j, \mathbf{r}^{(j)} \rangle_w|\}_{j=1}^p$. One obvious candidate for λ_{max} is the least upper bound $\sup_j |\langle \mathbf{x}_j, \mathbf{r}^{(j)} \rangle_w|$. If all betas are zero, $\mathbf{r}^{(j)} = \mathbf{W}^{-1}(\mathbf{y} - \bar{y}\mathbf{1})$ and $\langle \mathbf{x}_j, \mathbf{r}^{(j)} \rangle_w = \mathbf{x}'_j \mathbf{W} \mathbf{r}^{(j)} = \mathbf{x}'_j \mathbf{W} \mathbf{W}^{-1}(\mathbf{y} - \bar{y}\mathbf{1}) =$

$\mathbf{x}'_j(\mathbf{y} - \bar{y}\mathbf{1}) = \langle \mathbf{x}_j, \mathbf{y} - \bar{y}\mathbf{1} \rangle$. Thus, $\lambda_{max} = \sup_j | \langle \mathbf{x}_j, \mathbf{y} - \bar{y}\mathbf{1} \rangle |$. The set of non-zero betas (active set) becomes larger and larger when λ decreases from λ_1 to λ_{100} .

To summarize, the Pathwise Coordinate Descent algorithm is:

For $k = 1$ to 100:

For $m = 1 \cdots$ until $\boldsymbol{\beta}$ converge

Use $\boldsymbol{\beta}_{m-1}$ or initial value if $m = 1$

update $p_i = Pr(y_i = 1 | \mathbf{x}_i; \boldsymbol{\beta}_{m-1}) = \text{logistic}(\mathbf{x}'_i \boldsymbol{\beta}_{m-1})$

update $w_i = p_i(1 - p_i)$

update $z_i = \mathbf{x}'_i \boldsymbol{\beta}_{m-1} + \frac{y_i - p_i}{w_i}$

Find $\boldsymbol{\beta}_m$ by Coordinate Descent:

Use $\boldsymbol{\beta}_{m-1}$ as initial value, repeat below process for $\forall j \in \{1, \cdots, p\}$ again and again until convergence.

$$\beta_{j,m} = \begin{cases} 0 & \text{if } | \langle \mathbf{x}_j, \mathbf{r}^{(j)} \rangle_w | \leq \lambda_k \\ S_{\lambda_k}(\langle \mathbf{x}_j, \mathbf{r}^{(j)} \rangle_w) & \text{if } | \langle \mathbf{x}_j, \mathbf{r}^{(j)} \rangle_w | > \lambda_k \end{cases}$$

where $\mathbf{r}^{(j)} = \mathbf{z} - \mathbf{X}_{-j} \boldsymbol{\beta}_{-j,m}$

5 References

- Efron, B., & Hastie, T. (2016). *Computer Age Statistical Inference: Algorithms, Evidence, and Data Science* (Institute of Mathematical Statistics Monographs). Cambridge: Cambridge University Press. doi:10.1017/CBO9781316576533
- Friedman, J. H., Hastie, T., & Tibshirani, R. (2010). Regularization Paths for Generalized Linear Models via Coordinate Descent. *Journal of Statistical Software*, 33(1), 1–22. <https://doi.org/10.18637/jss.v033.i01>
- Hastie, T., Tibshirani, R., & Friedman, J. H. (2009). *The Elements of Statistical Learning: Data mining, Inference, and Prediction*. 2nd ed. New York: Springer.
- Hastie, T., Tibshirani, R., & Wainwright, M. (2015). *Statistical Learning with Sparsity: The Lasso and Generalizations* (1st ed.). Chapman and Hall/CRC. <https://doi.org/10.1201/b18401>
- Tibshirani, Ryan (2019). *Convex Optimization Lecture Notes*. <https://www.stat.cmu.edu/~ryantibs/convexopt/>