

Incorporating known risk factors into models

Jiaxuan Wang

<2017-02-23 Thursday>

1 objective

Incorporating known risk factors with unknown risk factors in predicting outcome. In the case of choosing between correlated variables, the model should favor known risk factors.

2 measuring success

Fixing the level of performance, the task of learning is to allocate weights to features so that desired structures are kept. A familiar example is lasso regression where the diamond shaped contour forces weights to be sparse so that they are interpretable in a sense. Another example is ridge regression where it corresponds to having a gaussian prior on weights. Learning is really a combination of data and domain expertise. Each model implicitly carries assumptions into the learning task and there's no universal consistent algorithm that works best under all distributions (no free lunch thm). However, there is one thing that for sure we want no matter what assumptions we make, that is the credibility of the model. In addition to having good performance, we want the model to be consistent with the current literature. More concretely, we want the model to place high weights on relevant and known features while keeping the unknown relevant features sparse. This whole process should be done in a data driven way so that the known risk factors are merely suggestions for the model to consider instead of forced constraints. In this spirit, we list the following properties of a credible model:

a) credibility should not come at the cost of performance b) irrelevant features whether known or unknown should have low weights c) within a group of dependent features, known risk factors should be given more weights and their weights should be dense d) within a group of dependent features of all unknown risk factors, the weights should be sparse

criteria a) are achieved by grid searching over validation set so that models in consideration has similar level of performance. b) is achieved by constraining on the size of parameters which all regularizations do.

For c) and d) we measure the KL divergence in each group of dependent features.

Here I give a typical examples of what I mean by measuring KL divergence in a group of dependent features

assume $r = [1, 1, 0, 0]^T$ and $\theta = [0.1, 0.2, -0.01, 0.02]^T$ (θ excluding b term), we first normalize each vector so that their $\|\cdot\|_1$ is 1.

$$r' = [0.5, 0.5, 0, 0]^T, \theta' = [0.32258065, 0.64516129, 0.03225806, 0.06451613]^T$$

To avoid 0 appearing in log of KL divergence calculation, a small smooth factor of $1e-6$ is added to any vector with 0, renormalizing giving

$$r'' = [4.99999000e-01, 4.99999000e-01, 9.99996000e-07, 9.99996000e-07]^T, \theta'' = [0.32258065, 0.64516129, 0.03225806, 0.06451613]^T$$

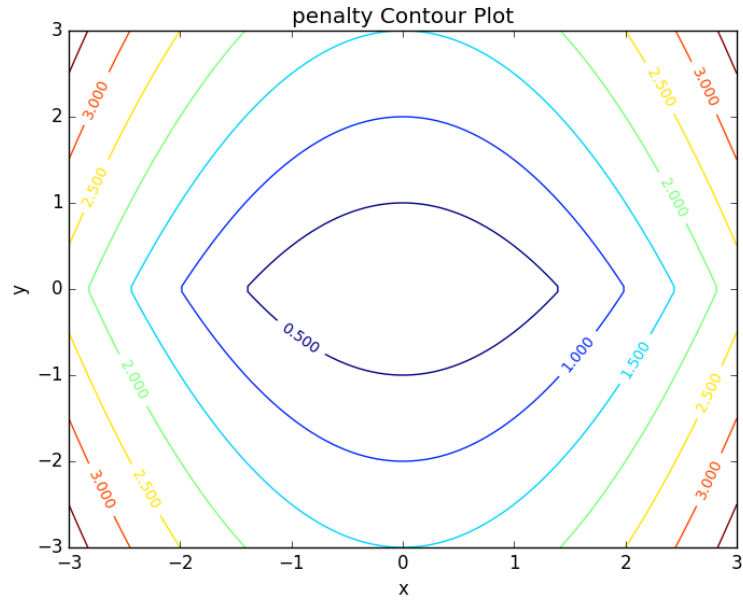
Then $KL(r''||\theta'')$ is the reported result in each dependent group, where $KL(x||y) = \sum_i p(x_i) \log \frac{p(x_i)}{p(y_i)}$

In the case where r is all 0 in relevant feature group, I give $\min_{v \in \text{one hot vectors}} KL(v||\theta'')$ as a loss as to encourage sparse feature.

3 approaches taken

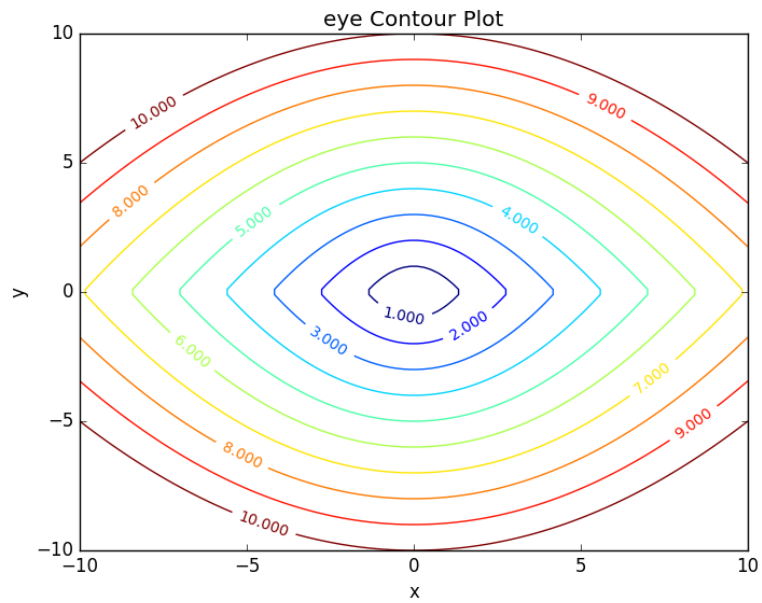
3.1 old approach

$$0.5 * \lambda_2 \|r * \theta\|_2^2 + \lambda_1 \|(1-r) * \theta\|_1 \text{ where } r \in \{0,1\}^d, \theta \in \mathbb{R}^d$$

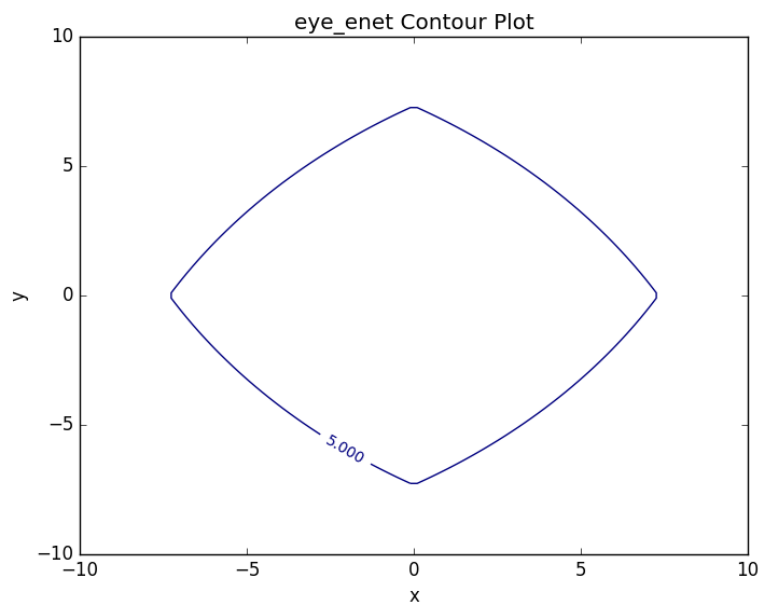


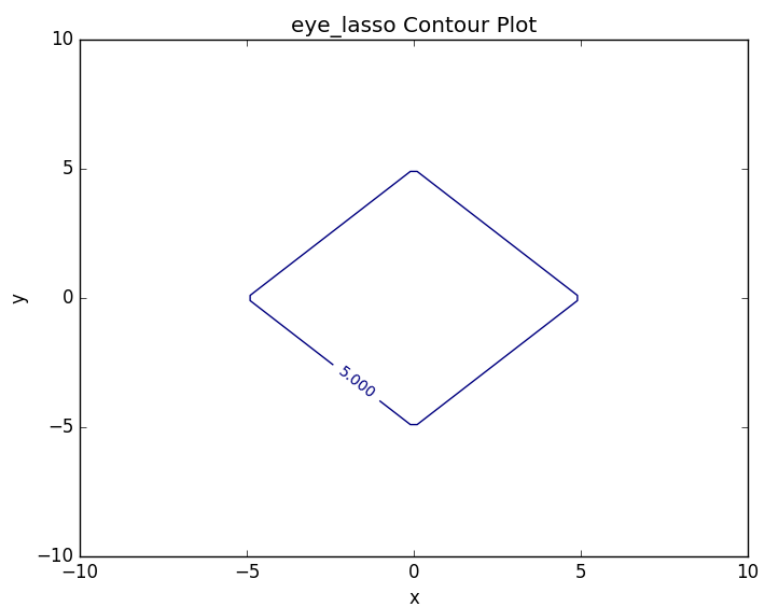
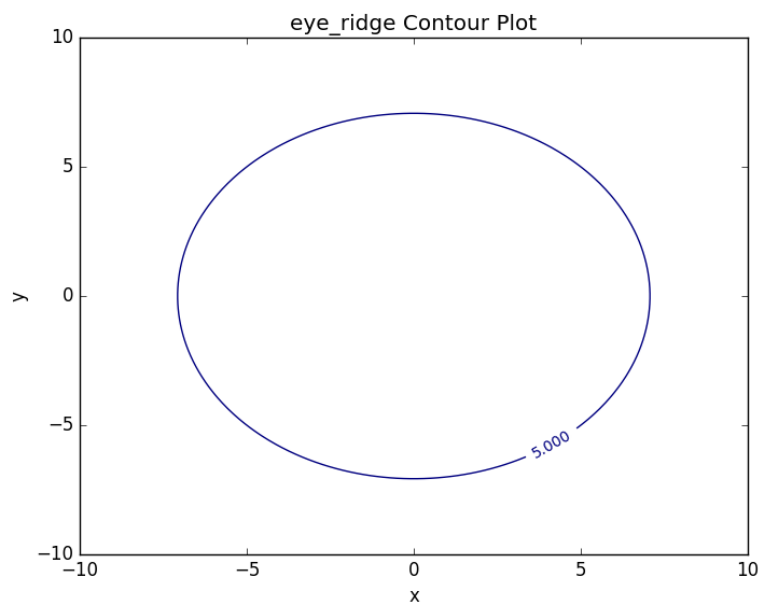
3.2 new approach

Fix a convex body that have property of the previous contour plot such that the angle at the end point is 45 degree. The following is the contour plot of its induced norm



In fact, by relaxing the constraint of r over binary to float, we can recover enet (setting $r=0.5 * 1$). Even without extending r , we can recover ridge ($r=1$) and lasso ($r=0$)





4 experiments

4.1 set up

4.1.1 1d data generation (genPartitionData)

Data $n = 5000$

n relevant groups (nrgroups) = 11

n irrelevant group (nirgroups) = 11

correlated variables pergroup (npergroup) = 10

$h_i = \text{Uniform}(-3, 1, n)$

$\theta_i = 1 \ \forall i$

$x_{i,j} \sim \text{Uniform}(1..2) h_i + N(0, 0.2)$ for $i \in [n]$ for $j \in [npergroup]$

$y = \frac{\sum_{i=1}^{nrgroups} h_i \theta_i}{\sum_{i=1}^{nrgroups} |\theta_i|} > -1$

r (known risk factors): for each correlated variable group, putting in one more known risk factor than the previous group

Loss function is the negative loss likelihood of the logistic regression model.

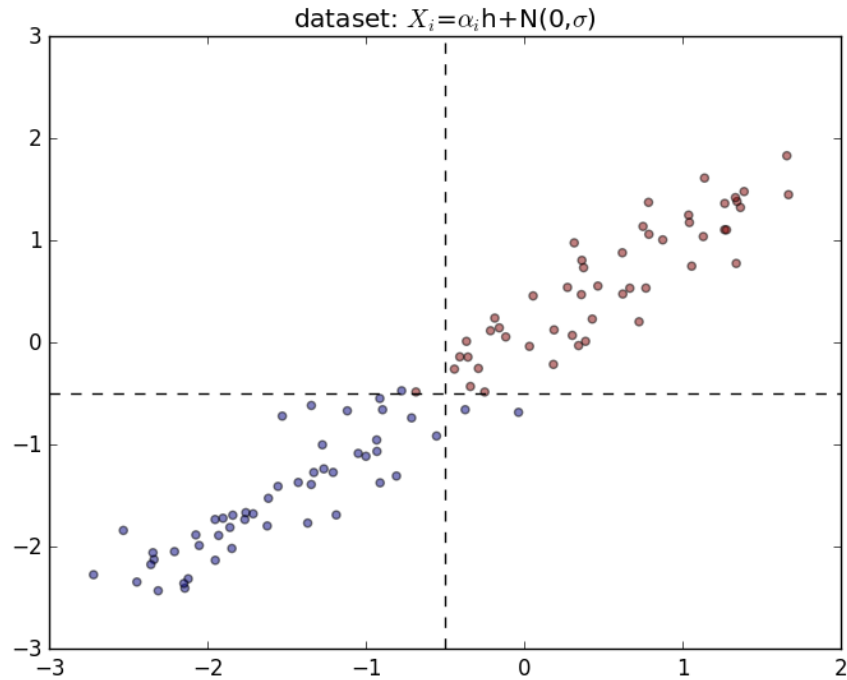
Optimizer: AdaDelta

Number of Epoch: 1000

Regularizers: elastic net, lasso, ridge, OWL, weighted lasso, weighted ridge, eye penalty

4.1.2 2d data generation

Data $n = 100$:



$h = \text{linspace}(-2.5, 1, n)$

$x_0 \sim \text{Uniform}(1..4) \cdot h + N(0, 0.2)$

$x_1 \sim \text{Uniform}(1..4) \cdot h + N(0, 0.2)$

$y = h > 0.5$

r (known risk factors) = [1, 0]

Loss function is the negative log likelihood of the logistic regression model.

Optimizer: AdaDelta

Number of Epoch: 1000

Regularizers: elastic net, lasso, ridge, OWL, weighted lasso, weighted ridge, penalty, eye penalty

4.1.3 eye penalty

$$q(\theta) := 2\beta \| (1-r) * \theta \|_1 + (1-\beta) \| r * \theta \|_2^2$$

$$\text{pena}(\theta) := \alpha q(\theta)$$

where $r \in \{0,1\}^d$, $\theta \in \mathbb{R}^d$, $\alpha \in \mathbb{R}_+$, $\beta \in (0,1)$ (β is also called $l1_{\text{ratio}}$ in this text)

For any constant c

$$\text{pena}(\theta) = c$$

is convex because pena is convex (addition of positively weighted norms)
 similarly, $q(\theta) = c$ is also convex

c can be chosen so that slope in the first quadrant between known risk factor x and unknown risk factor is -1

we define eye norm as a an atomic norm $\|\cdot\|_A$ as introduced in Venkat et al.

$$\|x\|_A := \inf\{t > 0 | x \in t \text{conv}(A)\}$$

Let $A = \{x | q(x) = \frac{\beta^2}{1-\beta}\}$, we get the eye penalty

Note that A is already a convex set, equivalently we write

$$\text{eye}(x) = \inf\{t > 0 | x \in t\{x | q(x) = \frac{\beta^2}{1-\beta}\}\}$$

1. derivation

The main intuition is to set c so that the slope in the first quadrant between known risk factor x and unknown risk factor is -1. Since we only care about this interaction between known and unknown risk factors and that $\{x | \text{pena}(x) = c\}$ is symmetric about origin, WLOG, we let y be the unknown feature and x be the known risk factor with constraint $y \geq 0, x \geq 0$.

$$\alpha[2\beta y + (1 - \beta)x^2] = c \tag{1}$$

$$\rightarrow 2\beta y + (1 - \beta)x^2 = \frac{c}{\alpha} \tag{2}$$

$$\rightarrow y = \frac{c}{2\alpha\beta} - \frac{(1 - \beta)x^2}{2\beta} \tag{3}$$

$$\rightarrow y = 0 \Rightarrow x = \sqrt{\frac{c}{\alpha(1 - \beta)}} \tag{4}$$

$$\rightarrow f'(x) = -\frac{(1 - \beta)}{\beta}x \tag{5}$$

$$\rightarrow f'\left(\sqrt{\frac{c}{\alpha(1 - \beta)}}\right) = -\frac{1 - \beta}{\beta} \sqrt{\frac{c}{\alpha(1 - \beta)}} = -1 \tag{6}$$

$$\rightarrow c = \frac{\alpha\beta^2}{1 - \beta} \tag{7}$$

$$\rightarrow 2\beta y + (1 - \beta)x^2 = \frac{\beta^2}{1 - \beta} \tag{8}$$

Thus, we just need $q(x) = \frac{\beta^2}{1-\beta}$

2. properties: a) A is symmetric about origin ($x \in A$ then $-x \in A$), so this is a norm

$$(a) \text{ eye}(t \theta) = |t| \text{ eye}(\theta)$$

$$(b) \text{ eye}(\theta + \beta) \leq \text{eye}(\theta) + \text{eye}(\beta)$$

$$(c) \text{ eye}(\theta) = 0 \text{ iff } \theta = 0$$

b) β doesn't affect the shape of contour, so no need to search over β
proof:

consider the contour $B_1 = \{x: \text{eye}_{\beta_1}(x) = t\}$ and $B_2 = \{x: \text{eye}_{\beta_2}(x) = t\}$

We want to show B_1 is similar to B_2

case1: $t = 0$, then $B_1 = B_2 = \{0\}$ by property a3

case2: $t \neq 0$

we can equivalently write B_1 and B_2 as: (by definition and a1 and q convex)

$$B_1 = t \{x: x \in \{x \mid q_{\beta_1}(x) = \frac{\beta_1^2}{1-\beta_1}\}\}$$

$$B_2 = t \{x: x \in \{x \mid q_{\beta_2}(x) = \frac{\beta_2^2}{1-\beta_2}\}\}$$

$$\text{let } B_{1'} = \{x: x \in \{x \mid q_{\beta_1}(x) = \frac{\beta_1^2}{1-\beta_1}\}\} \text{ and } B_{2'} = \{x: x \in t \{x \mid q_{\beta_2}(x) = \frac{\beta_2^2}{1-\beta_2}\}\}$$

$$\text{Claim: } B_{2'} = \frac{\beta_2(1-\beta_1)}{\beta_1(1-\beta_2)} B_{1'}$$

It should be clear that if this claim is true then B_1 is similar to B_2 and we are done

take $x \in B_{1'}$,

$$\text{then } q_{\beta_1}(x) = 2 \beta_1 \|(1-r) * x\|_1 + (1-\beta_1) \|r * x\|_2^2 = \frac{\beta_1^2}{1-\beta_1}$$

$$\text{let } x' = \frac{\beta_2(1-\beta_1)}{\beta_1(1-\beta_2)} x$$

$$q_{\beta_2}(x') = 2\beta_2\|(1-r)*x'\|_1 + (1-\beta_2)\|r*x'\|_2^2 \quad (9)$$

$$= \frac{2\beta_2^2(1-\beta_1)}{\beta_1(1-\beta_2)}\|(1-r)*x\|_1 + \frac{\beta_2^2(1-\beta_1)^2}{\beta_1^2(1-\beta_2)}\|r*x\|_2^2 \quad (10)$$

$$= \frac{\beta_2^2(1-\beta_1)}{\beta_1^2(1-\beta_2)}(2\beta_1\|(1-r)*x\|_1 + (1-\beta_1)\|r*x\|_2^2) \quad (11)$$

$$= \frac{\beta_2^2(1-\beta_1)}{\beta_1^2(1-\beta_2)} \frac{\beta_1^2}{1-\beta_1} \quad (12)$$

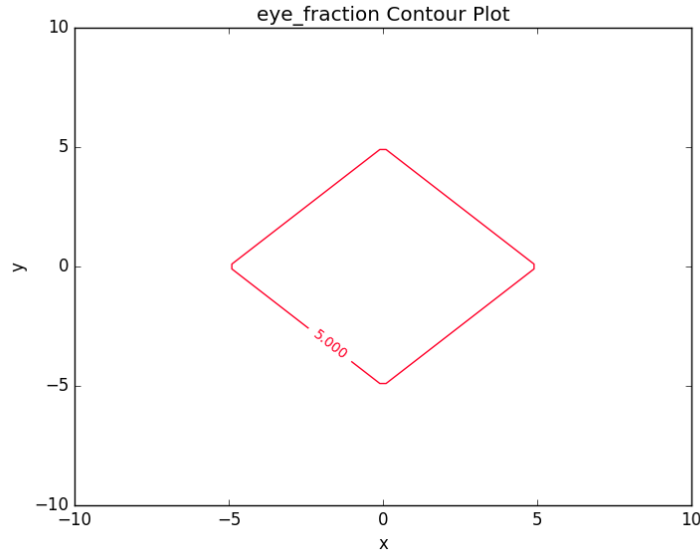
$$= \frac{\beta_2^2}{1-\beta_2} \quad (13)$$

so $x' \in B_2$. Thus $\frac{\beta_2(1-\beta_1)}{\beta_1(1-\beta_2)} B_1 \subset B_2$. The other direction is similarly proven.

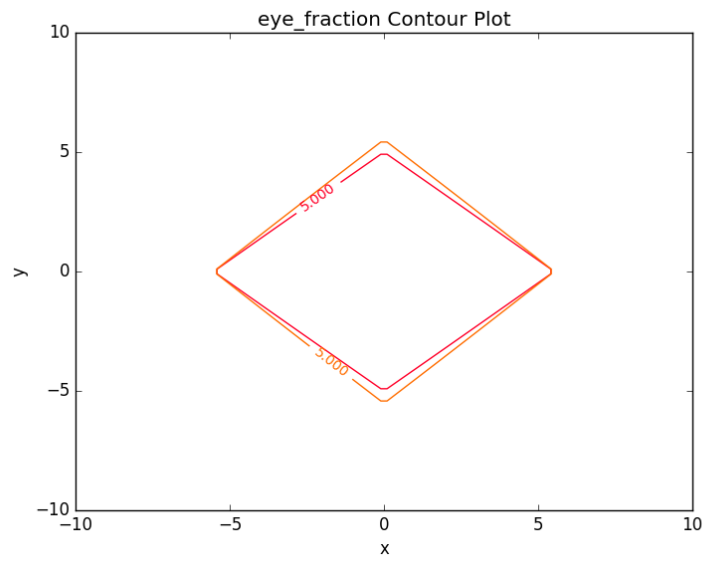
3. extending r to $[0,1]^d$ At times, it makes sense for risk factor to be fractionally weighted (eg. odds ratio in medical documents)

varying r_1 and r_2

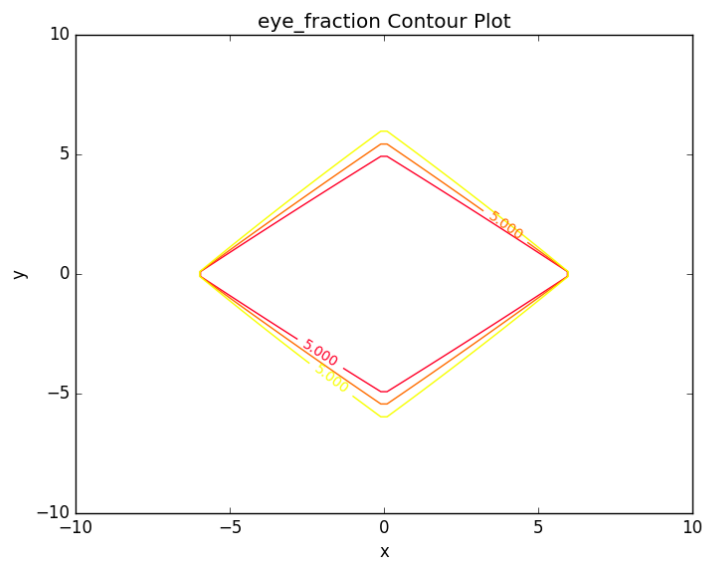
$r_1 = 0.0$



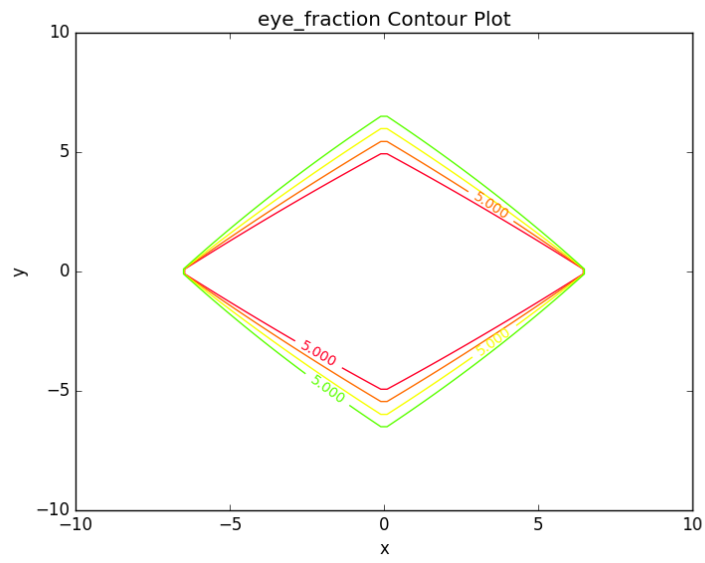
$r_1 = 0.1$



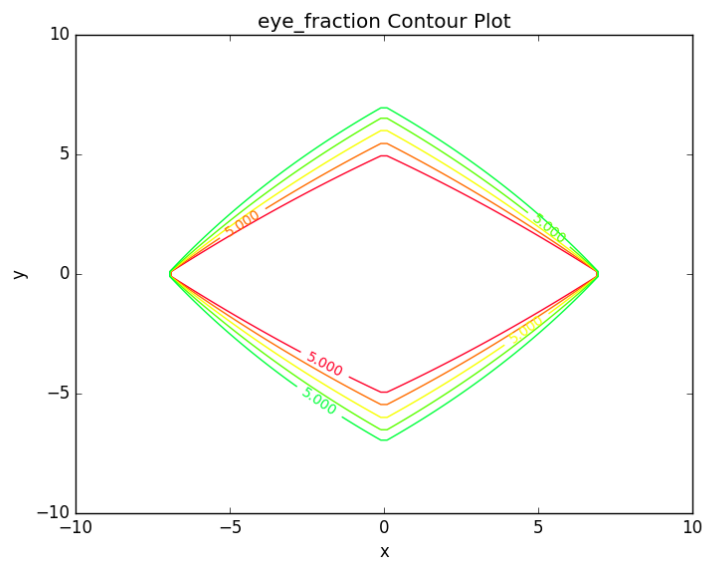
$$r_1 = 0.2$$



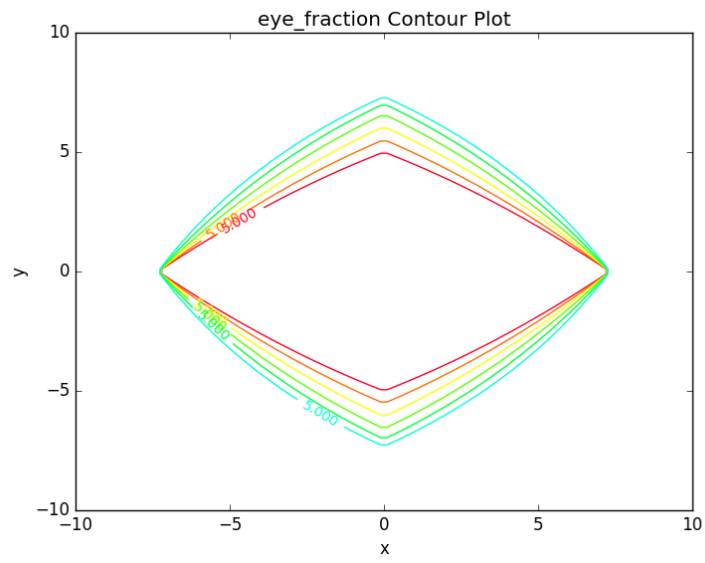
$$r_1 = 0.3$$



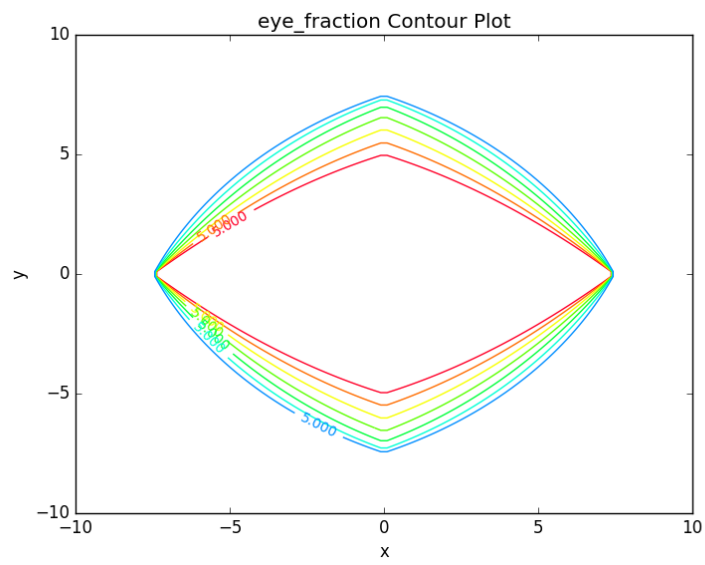
$$r_1 = 0.4$$



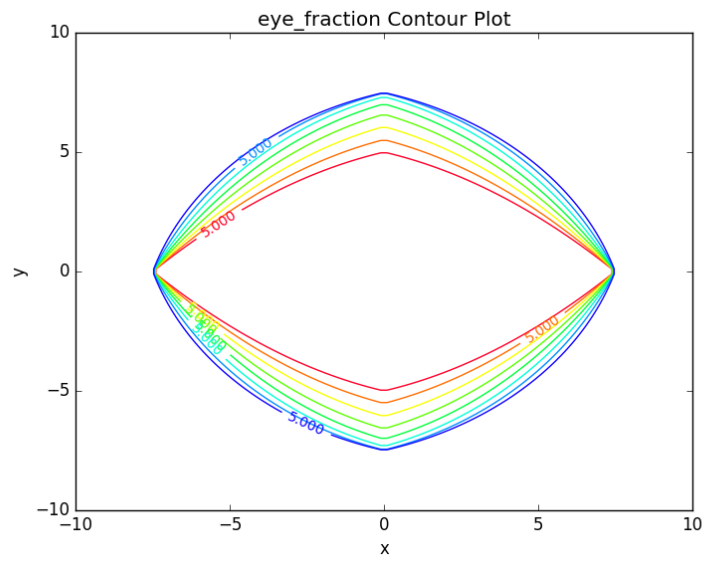
$$r_1 = 0.5$$



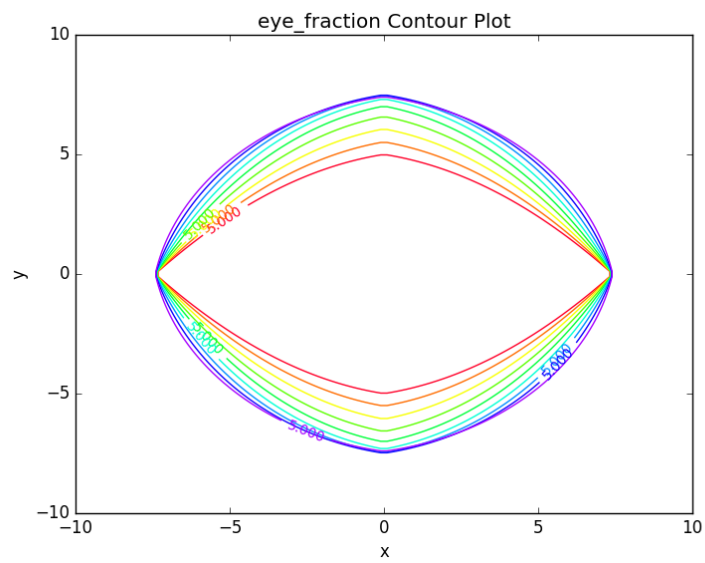
$$r_1 = 0.6$$



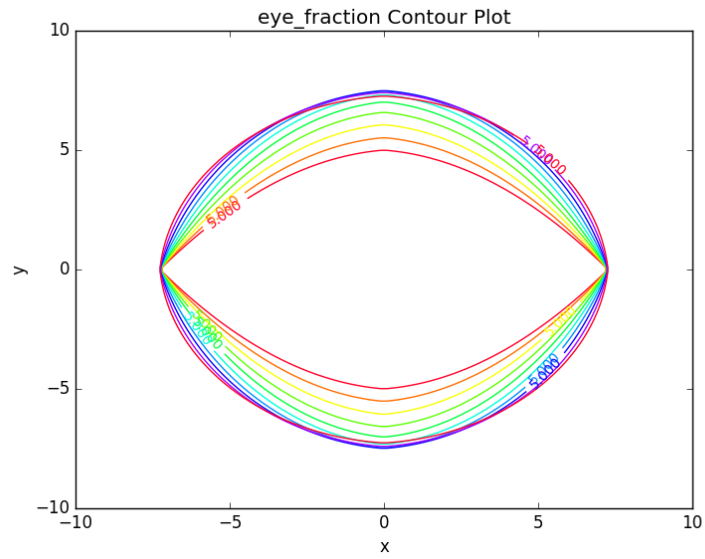
$$r_1 = 0.7$$



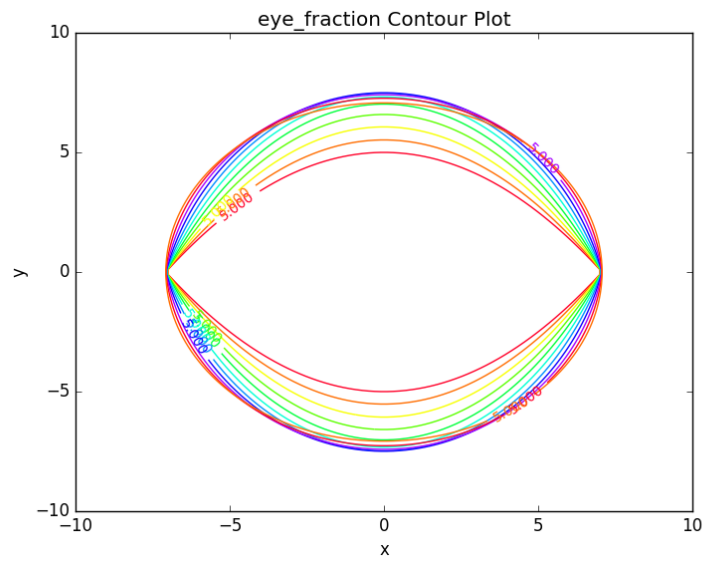
$$r_1 = 0.8$$



$$r_1 = 0.9$$

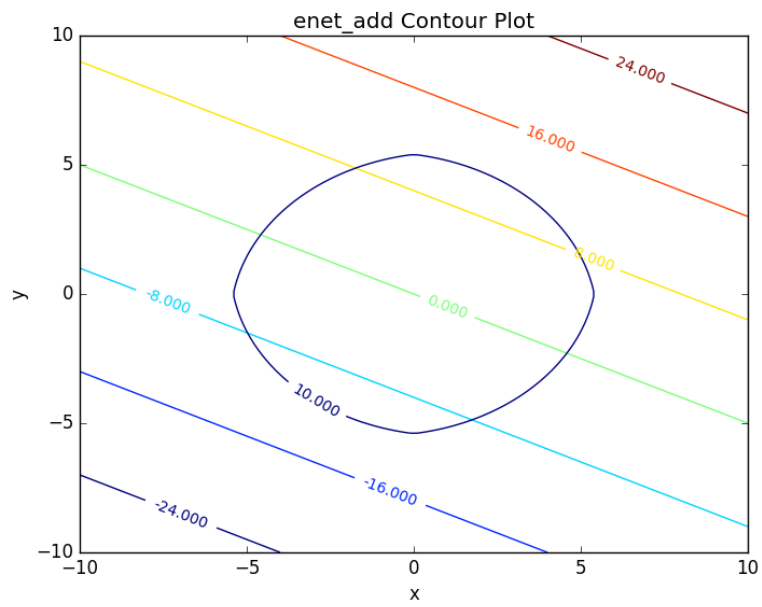


$$r_1 = 1.0$$



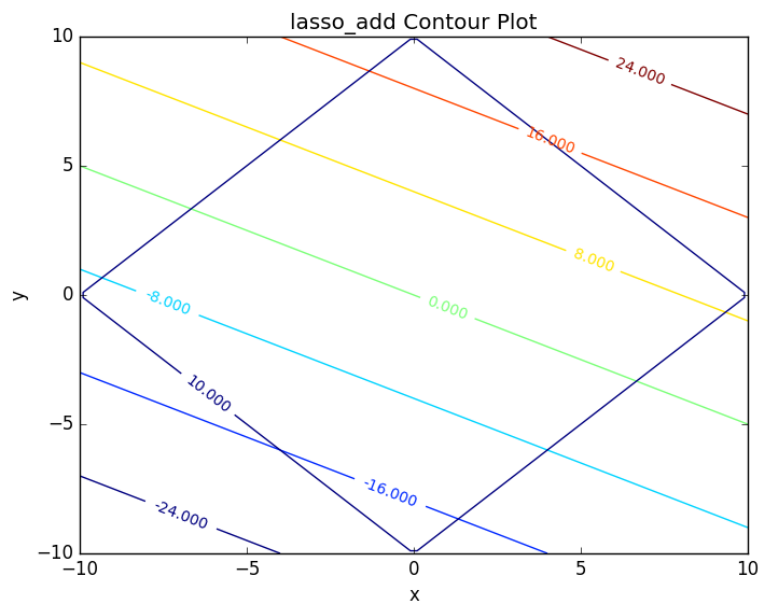
4.1.4 elastic net

$\alpha * (c * ||\theta||_1 + 0.5 * (1 - c) * ||\theta||_2^2)$ where c is a scaler



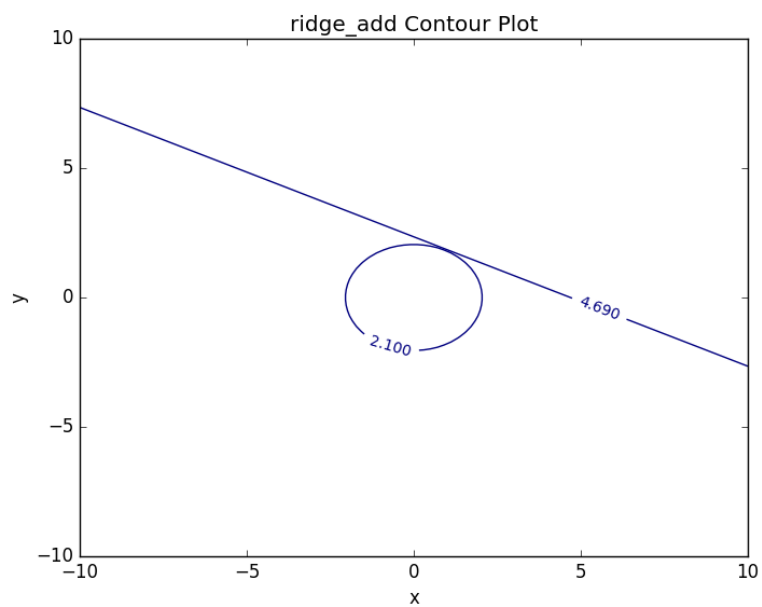
4.1.5 lasso

$$\alpha^* ||\theta||_1$$



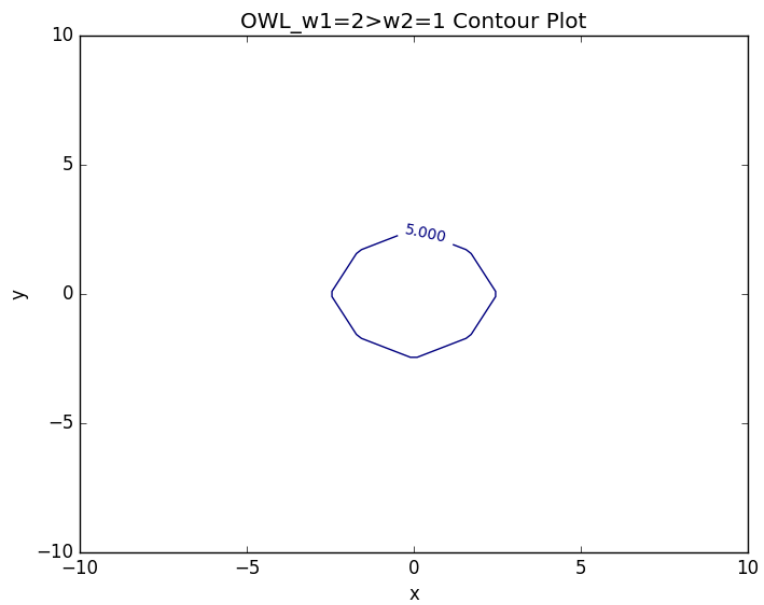
4.1.6 ridge

$$0.5 * \alpha * ||\theta||_2^2$$

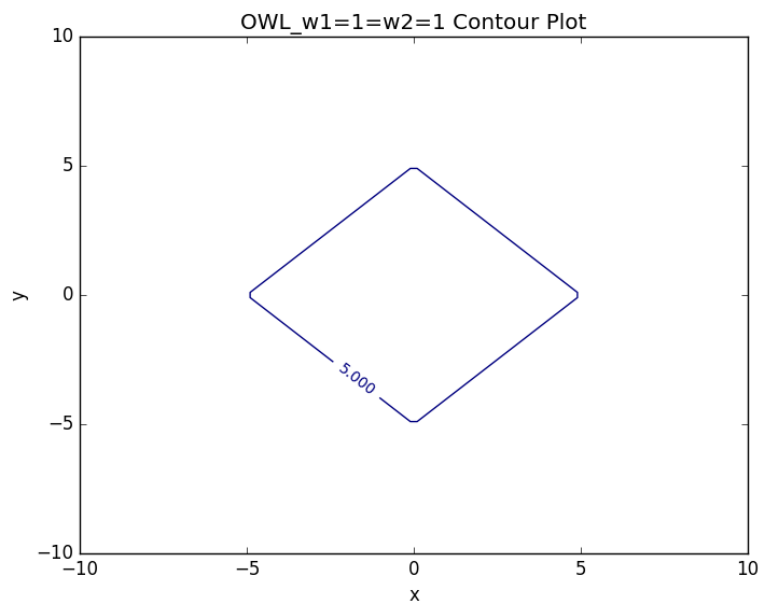


4.1.7 OWL

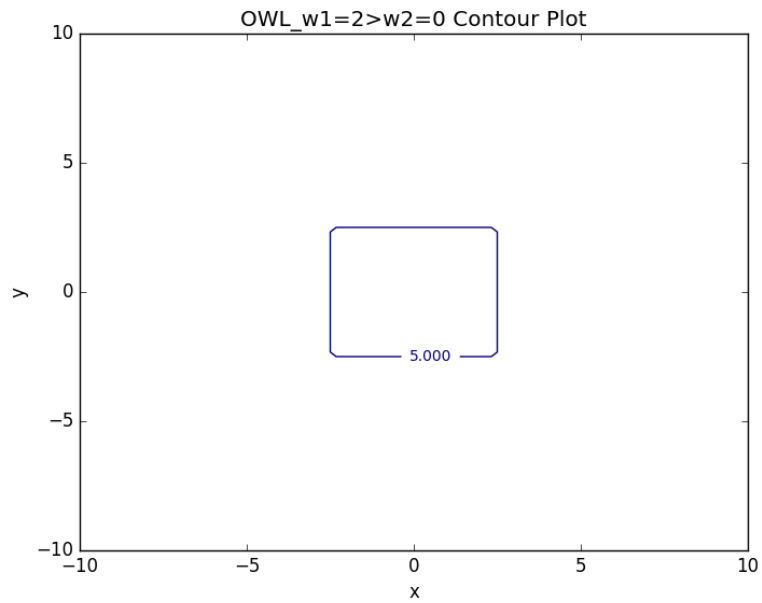
$$\alpha * \sum_{i=1}^n w_i |x|_{[i]} \text{ where } w \in K_{m+} \text{ (monotone nonnegative cone)}$$



degenerated case: back to lasso



degenerated case: back to l_{inf}



some properties:

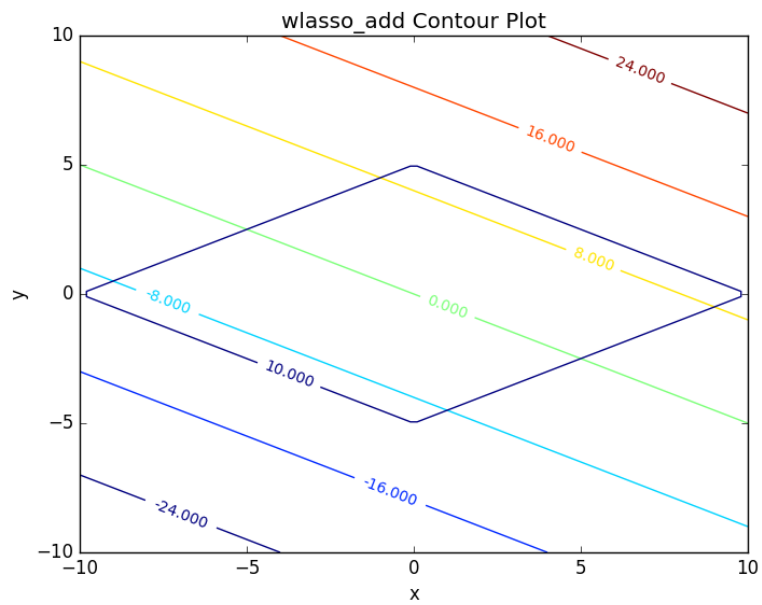
generalization of OSCAR norm

symmetry with respect to signed permutations

in the regular case, the minimal atomic set for this norm is known (the corners are easily calculated)

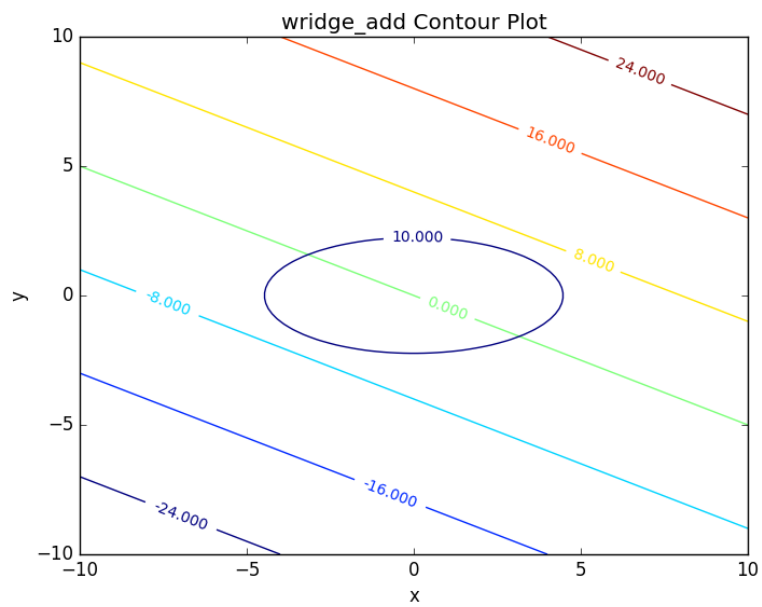
4.1.8 weighted lasso

$\alpha * ||w * \theta||_1$ where $w \in \mathbb{R}_{+}^d$



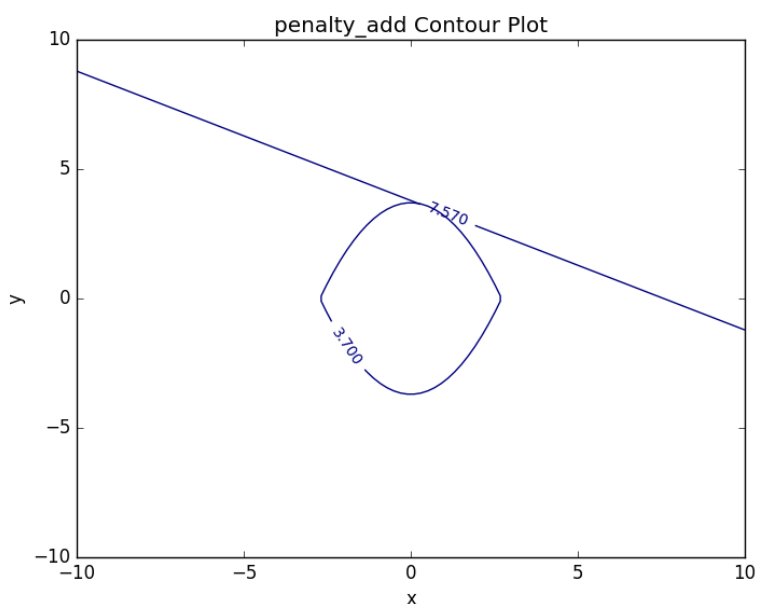
4.1.9 weighted ridge

$0.5 * \alpha * ||w * \theta||_2^2$ where $w \in \mathbb{R}_+^d$



4.1.10 old penalty

$\alpha * (0.5 * (1-c) * ||r * \theta||_2^2 + c * ||(1-r) * \theta||_1)$ where $r \in \{0,1\}^d$, $\theta \in \mathbb{R}^d$, $\alpha \in \mathbb{R}$, $c \in \mathbb{R}$



4.2 running procedure

4.2.1 first run (regularized b)

b regularized

- fix hyperparameters to predefined value

- repeat the following 100 times:

- generate data ($x_2 = 2x_1$), run the selected regularizers, record θ

4.2.2 second run (unregularized b, validation)

b unregularized

- generate two datasets ($x_2 = 2x_1$), one for training, one for validation

- parameter search over the different hyperparams of the regularizers

- for each regularizer, use the hyperparameters that achieves the minimal

loss

- repeat the following 100 times:

- generate data, run the selected regularizers, record θ

4.2.3 third run (data normalized, eye penalty)

b unregularized

generate two datasets ($x_2 = 2 \times 1$), one for training, one for validation

normalize the data to 2 mean and 2 variance (validation data is normalized using mean and variance for the training data)

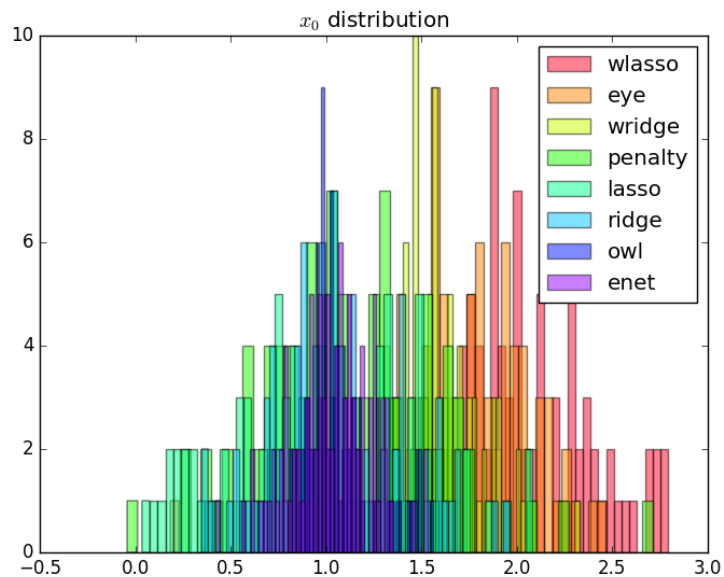
parameter search over the different hyperparams of the regularizers

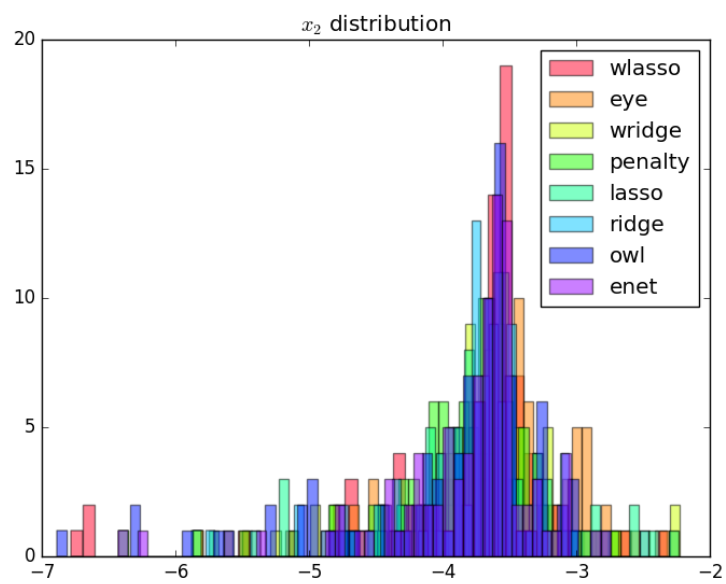
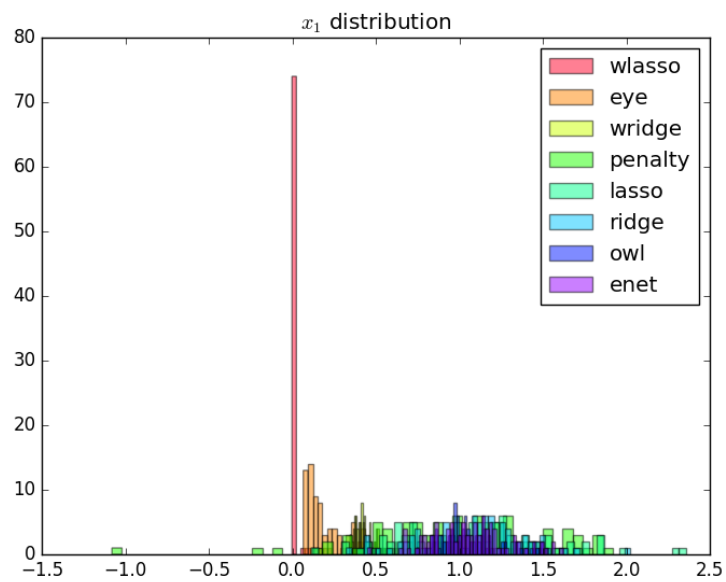
for each regularizer, use the hyperparameters that achieves the minimal loss

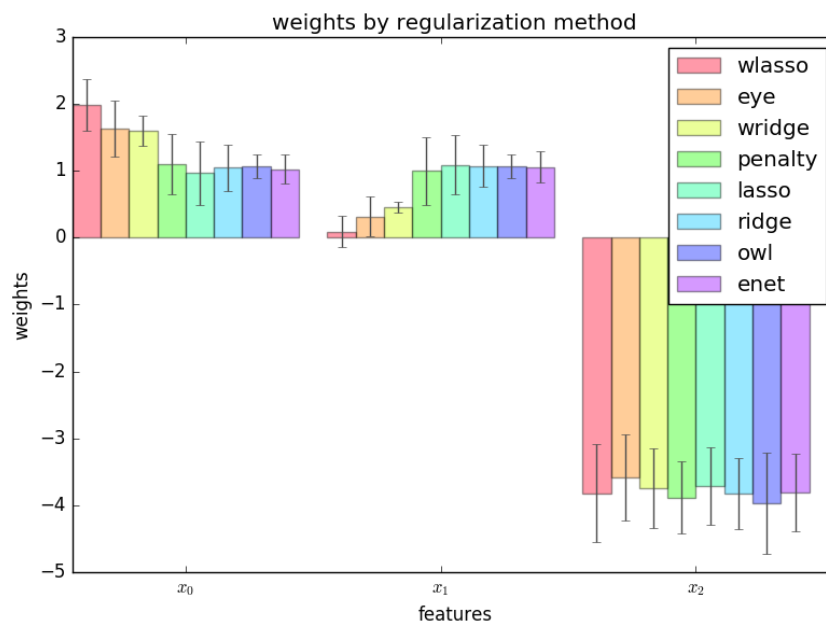
repeat the following 100 times:

generate data, normalize data, run the selected regularizers, record θ

The choosing criteria is still loss b/c AUROC is always going to be 1 in the deterministic case:







4.2.4 Fourth run (noise added)

b unregularized

generate two datasets, one for training, one for validation

normalize the data to 2 mean and 2 variance (validation data is normalized using mean and variance for the training data)

parameter search over the different hyperparams of the regularizers

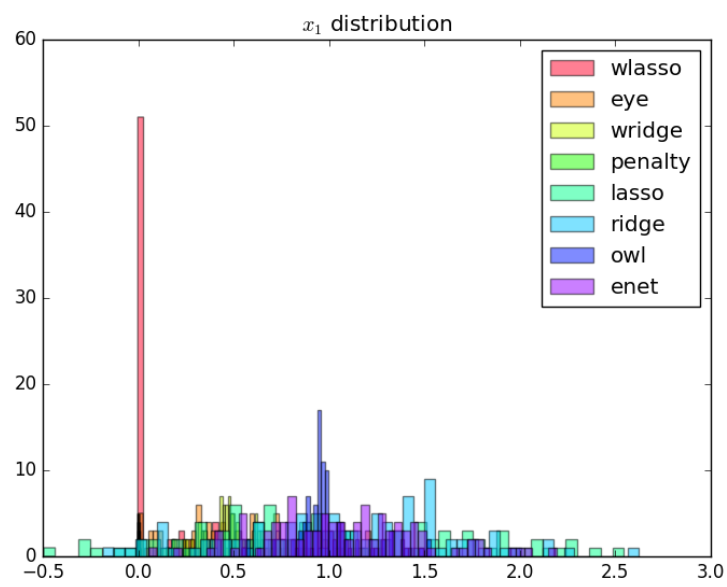
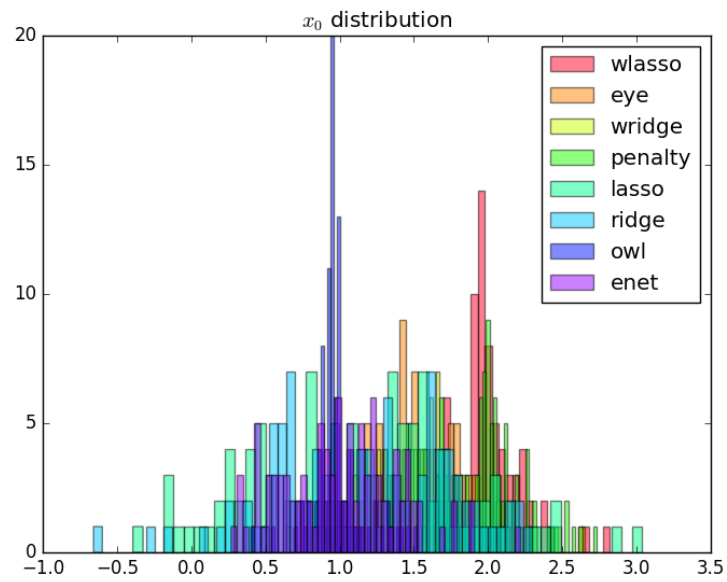
for each regularizer, use the hyperparameters that achieves the minimal

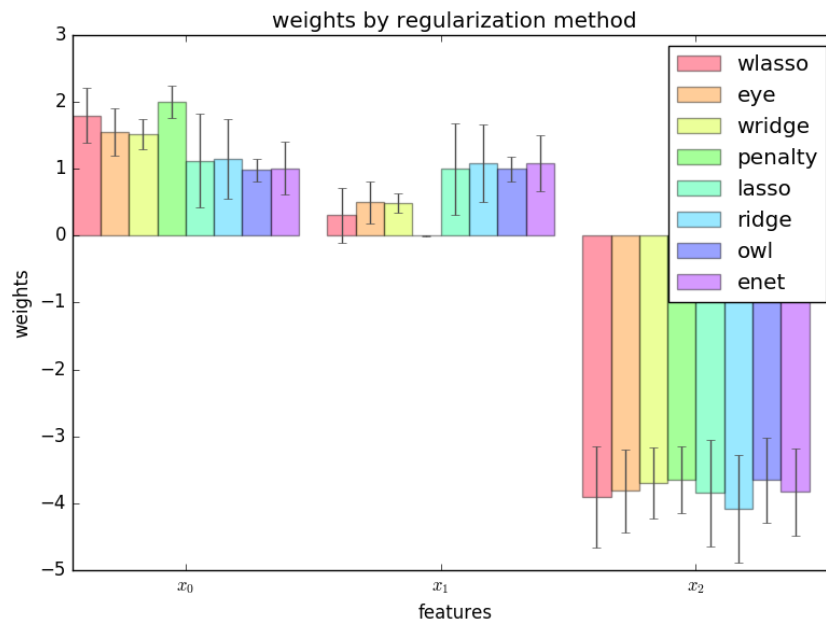
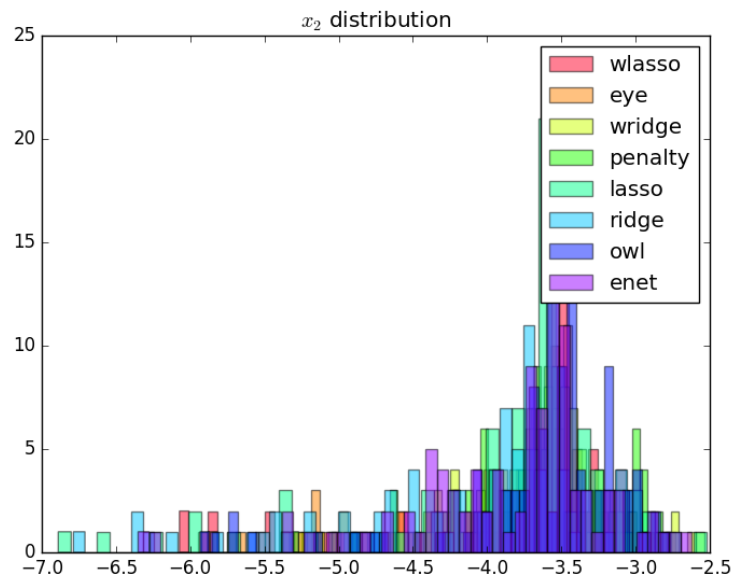
loss

repeat the following 100 times:

generate data ($x_i = \text{Uniform}(0..4) h + N(0,0.2)$), normalize data, run the selected regularizers, record θ

The choosing criteria is loss



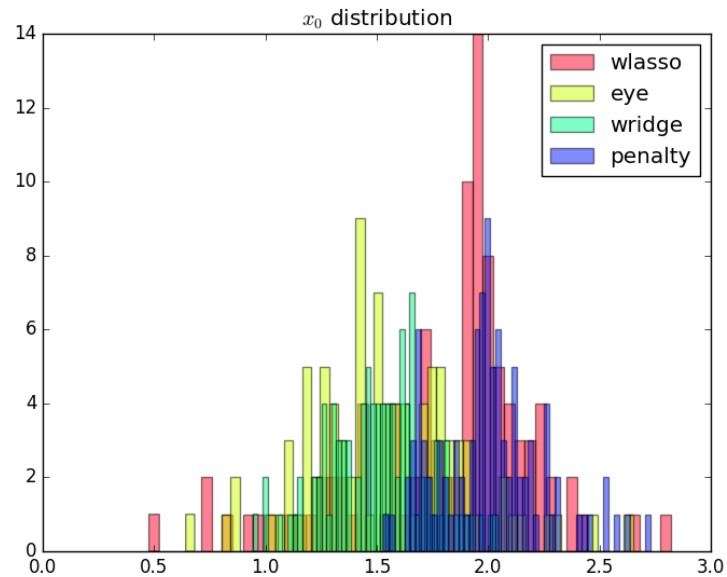


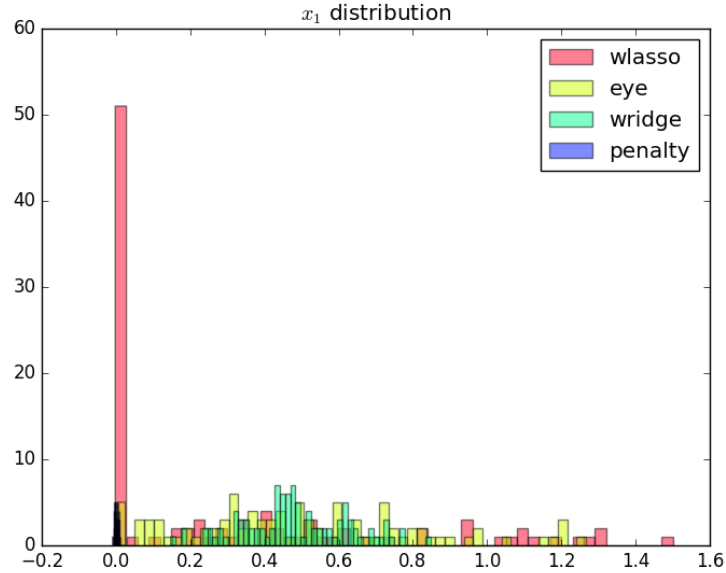
hyper parameter used:

- enet(0.01, 0.2)
- eye(array([1., 0.]), 0.01, 0.4)

- `lasso(0.0001)`
- `OWL([2, 1], 0.01)`
- `penalty(array([1., 0.]), 0.1, 1.0)`
- `ridge(0.001)`
- `weightedLasso(array([1., 2.]), 0.01)`
- `weightedRidge(array([1., 2.]), 0.01)`

The sparsity in penalty can be explained as I placed no constraint on known risk factor (l1 ratio is 1), so it only regularizes x_1 not x_0





4.2.5 fifth run (nd data)

b unregularized

generate two datasets, one for training, one for validation

normalize the data to 2 mean and 2 variance (validation data is normalized using mean and variance for the training data)

parameter search over the different hyperparams of the regularizers (each of the final candidate has loss around 0.083)

for each regularizer, use the hyperparameters that achieves the minimal loss

repeat the following 10-20 times:

generate data (detailed in nd data generation section), normalize data,

run the selected regularizers, record θ

The choosing criteria is loss

KL divergence metric filtering for relevant features:

eye: 2.5722261048

wlasso: 5.18104309657

wridge: 6.8364694347

lasso: 18.9613782735

ridge: 12.7547711529

owl: 13.5265637342

enet: 17.7231341012
KL divergence metric including irrelevant features:
eye: 13.1307145901
wlasso: 7.55507729218
wridge: 11.5881850514
lasso: 31.1710069808
ridge: 16.9635832109
owl: 17.5479982613
enet: 30.2439873411
kl_{metricvisual} (generated using `genresult.py:genndlosscsv`)

5 issues encountered

the validation sweep tend to pick very small regularization which essentially unconstrained (maybe increase the noise could help?)

6 next

1. develop a more general metric (related to 4)
2. look for all unknown but different weights (balanced data)
3. extend risk to fractional
4. generate data in a more diverse manner: eg. $x_3 = h_1 + h_2$