Jong-Han Kim

## **Non-Quadratic Losses**

Jong-Han Kim

## EE787 Fundamentals of machine learning Kyung Hee University

## Penalty functions and error histograms

#### Loss and penalty functions

- ▶ empirical risk (or average loss) is  $\mathcal{L}(\theta) = \frac{1}{n} \sum_{i=1}^{n} \ell(\theta^{\top} x^{i}, y^{i})$
- $\blacktriangleright$  the loss function  $\ell(\hat{y},y)$  penalizes deviation between the predicted value  $\hat{y}$  and the observed value y
- ▶ common form for loss function:  $\ell(\hat{y}, y) = p(\hat{y} y)$
- ▶ p is the penalty function
- ▶ e.g., the square penalty  $p^{
  m sqr}(r) = r^2$
- $r = \hat{y} y$  is the *prediction error* or *residual*

## **Penalty functions**

- the penalty function tells us how much we object to different values of prediction error
- $\blacktriangleright$  usually p(0)=0 and  $p(r)\geq 0$  for all r
- ▶ if p is symmetric, i.e., p(-r) = p(r), we care only about the magnitude (absolute value) of prediction error
- ▶ if p is asymmetric, i.e.,  $p(-r) \neq p(r)$ , it bothers us more to over- or underestimate

#### Square versus absolute value penalty



 $\blacktriangleright$  for square penalty  $p^{
m sqr}(r)~=r^2$ 

▶ for small prediction errors, penalty is very small (small squared)

▶ for large prediction errors, penalty is very large (large squared)

- $\blacktriangleright$  for absolute penalty  $p^{\mathsf{abs}}(r) = |r|$ 
  - for small prediction errors, penalty is large (compared to square)
  - for large prediction errors, penalty is small (compared to square)

## Predictors and choice of penalty function

- choice of penalty function depends on how you feel about large, small, positive, or negative prediction errors
- ▶ different choices of penalty function yield different predictor parameters
- choice of penalty function shapes the histogram of prediction errors, i.e.,

 $r^1, \ldots, r^n$ 

(usually divided into bins and displayed as bar graph distribution)

## Histogram of residuals



▶ artificial data with n = 300 and d = 30, using 50/50 test/train split

- plots show histogram of residuals  $r^1, \ldots, r^n$
- $\blacktriangleright$  tilted loss results in distribution with most residuals  $r^i <$  0, i.e., predictor prefers  $\hat{y}^i < y^i$

## Robust fitting

#### Outliers

- ▶ in some applications, a few data points are 'way off', or just 'wrong'
- ▶ occurs due to transcription errors, error in decimal point position, etc.
- ▶ these points are called *outliers*
- even a few outliers in a data set can result in a poor predictor
- several standard methods are used to remove outliers, or reduce their impact
- ▶ one simple method:
  - create predictor from data set
  - flag data points with large prediction errors as outliers
  - remove them from the data set and repeat

## **Robust penalty functions**

- ▶ we say a penalty function is *robust* if it has low sensitivity to outliers
- robust penalty functions grow more slowly for large prediction error values than the square penalty
- and so 'allow' the predictor to have a few large prediction errors (presumably for the outliers)
- so they handle outliers more gracefully
- ▶ a *robust predictor* might fit, *e.g.*, 98% of the data very well

## Huber loss



▶ the *Huber* penalty function is

$$p^{\mathsf{hub}}(r) = egin{cases} r^2 & ext{if } |y| \leq lpha \ lpha(2|r|-lpha) & ext{if } |r| > lpha \end{cases}$$

```
\triangleright \alpha is a parameter
```

 $\blacktriangleright$  quadratic for small r, affine for large r

### Huber loss

 $\blacktriangleright$  linear growth for large r makes fit less sensitive to outliers

▶ ERM with Huber loss is called a *robust* prediction method



## Log Huber



 $\blacktriangleright$  quadratic for small y, logarithmic for large y

$$p^{\mathsf{dh}}(y) = egin{cases} y^2 & ext{if } |y| \leq lpha \ lpha^2(1-2\log(lpha)+\log(y^2)) & ext{if } |y| > lpha \end{cases}$$

 $\blacktriangleright$  diminishing incremental penalty at large y

## Log Huber



## ▶ even less sensitive to outliers than Huber

## **Error distribution**







# Quantile regression

#### Absolute penalty

- $\blacktriangleright$  absolute penalty  $p^{\mathsf{abs}}(r) = |r|$
- ▶ the best constant predictor  $(d = 1, x_1 = 1)$  minimizes  $\frac{1}{n} \sum_{i=1}^{n} |\theta_1 y^i|$
- solution is  $\hat{y} = heta_1 = \mathsf{median}\{y^1, \dots, y^n\}$
- ▶ (cf. best constant predictor with square loss, which is the average)
- ▶ rough idea:

$$rac{d}{d heta_1}\sum_{i=1}^n | heta_1-y^i| = ig( {\sf number of } y^i {\sf s} < heta_1 ig) - ig( {\sf number of } y^i {\sf s} > heta_1 ig)$$

 in general case, with no regularization on constant feature, median of errors is zero

#### Tilted absolute penalty

• tilted absolute penalty: for  $0 < \tau < 1$ ,

$$p^{\mathsf{tlt}}(z) = au(z)_+ + (1- au)(z)_- = (1/2)|z| + ( au - 1/2)z$$

- ightarrow au= 0.5: equal penalty for over- and under-estimating
- ▶  $\tau = 0.1$ : 9× more penalty for under-estimating
- ▶  $\tau = 0.9$ : 9× more penalty for over-estimating

## Tilted absolute penalty



- ▶ best constant predictor for  $\tau$  minimizes  $\frac{1}{n} \sum_{i=1}^{n} p^{\text{tlt}}(\theta_1 y^i)$
- $\blacktriangleright$  fraction au of training data satisfies  $heta_1 < y^i$
- τ-quantile of training residuals is zero
- ▶ solution is  $\hat{y} = \theta_1 = \text{the } (1 \tau)$ -quantile of  $\{y^1, \dots, y^n\}$
- plots show histogram of residuals for training data

## Quantile regression

- quantile regression uses penalty  $p^{tlt}$
- $\blacktriangleright$  in general case, with no regularization on constant feature,  $\tau\text{-quantile}$  of optimal errors is zero
- hence the name quantile regression

## Example: Quantile regression



▶ fit training data with loss 
$$l(\hat{y}, y) = p^{t|t}(\hat{y} - y)$$

## **Example: Quantile regression**



▶ three quite different predictors

## Example: Quantile regression



 $\blacktriangleright$   $\tau$ -quantile of training residuals is zero