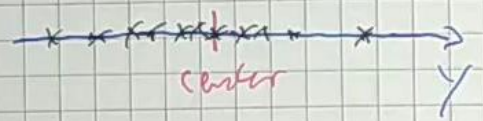# Robust loss functions

- put lower emphasize on outliers $\Rightarrow$ regression solution is influenced (much) less
- illustration: find the "center" of $N$ 1-dimensional data points

  - compute mean: $\bar{y} = \frac{1}{N} \sum_i y_i$



  center

  - compute median: $y_{med} = \{ y \mid \#\{y_i < y\} = \#\{y_i > y\}\}$

  ( is unique, when $N$ is odd: $y_{med}$ middle point in sorted order )

Q: which is better in the presence of outliers - mean or median?

construct toy problem:  inlier distribution : $y \sim N(0, \sigma^2)$

outlier distr. $y \sim N(0, \tau^2)$    $\tau^2 > \sigma^2$

superposition: "contaminated" distribution: $y \sim (1-\varepsilon) N(0, \sigma^2) + \varepsilon N(0, \tau^2)$

$(1-\varepsilon)$ :  inlier fraction    $(0 \le \varepsilon \le 1)$

$\frac{\tau}{\sigma} \ge 1$ :  scale ratio of the two distr.

- setting: • choose TS of size $N$ $\Rightarrow$ determine mean and median
  • repeat this infinitely often with different TS $\Rightarrow$ determine variance
  between the means of each TS and medians of each TS

$$\text{Var}_{TS}(\bar{y}) = \frac{\sigma^2}{N}\left(1 - \varepsilon + \varepsilon \frac{\tau^2}{\sigma^2}\right) \quad , \quad \text{Var}_{TS}(y_{med}) = \frac{\sigma^2}{N} \frac{\pi}{2\left(1 - \varepsilon + \varepsilon \frac{\tau}{\sigma}\right)^2}$$

$$\frac{\text{Var}(y_{med})}{\text{Var}(\bar{y})} = \begin{cases} > 1 & \Rightarrow \text{ mean is better} \\ < 1 & \Rightarrow \text{ median is better} \end{cases}$$

| $\dfrac{Var(Y_{med})}{Var(\bar{Y})}$ | $\varepsilon = 0$ | $\varepsilon = 5\%, \ \frac{\tau}{\sigma} = 4$ | $\varepsilon = 5\%, \ \frac{\tau}{\sigma} = 10$ | $\varepsilon = 10\%, \ \frac{\tau}{\sigma} = 3$ | $\varepsilon = 10\%, \ \frac{\tau}{\sigma} = 10$ |
|---|---|---|---|---|---|
| | 1.57 | 1 | 0.29 | 1 | 0.17 |
| | ↑ mean better | both are equal | median better | equal | median much better |

⟶ conclusion: - median more robust in the presence of outliers

         - mean more accurate for the inliers (when outliers are eliminated)

Q: Can we have both? Yes, <u>Huber loss</u>

desirable: results of a machine learning alg. do not change much if we picked

         a different training set $\hat{=}$ "variance of an estimator over all possible TS."

⟹ low variance $\hat{=}$ low probability to pick unfortunate TS (with high resulting error)

simplest estimator: select a representative for set of 1-D points ("center")

- mean : lowest variance without outliers     } toy problem: conta-
- median : lowest variance in presence of outliers     } minated Gaussian distribution
- Huber loss: is as good as mean w/o outliers, as median with outliers

| $\dfrac{Var(Y_{Huber})}{Var(\bar{Y})}$ | 1.05 | | 0.21 | | 0.13 |
|---|---|---|---|---|---|

Idea of Huber loss:

- penalize inliers like mean     } in terms of their distance from the
         outliers like median     } representative

- reformulate computation of mean and median as a optimisation problem
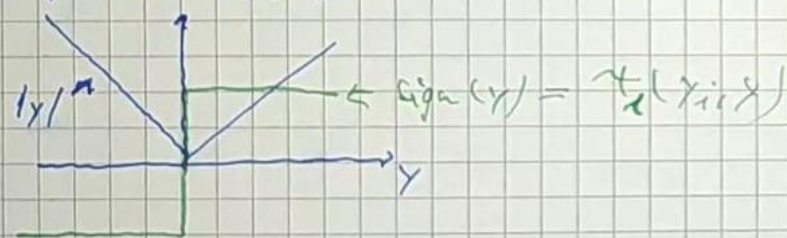
$$\hat{Y}_{center} = \arg\min_{Y} \sum_i loss(Y_i, Y)$$

- if $loss(Y_i, Y) = \frac{1}{2}(Y_i - Y)^2 \Rightarrow$ mean:

$$\frac{\partial}{\partial y} \sum_i loss(Y_i, Y) = \frac{\partial}{\partial y} \sum_i \frac{1}{2}(Y_i - Y)^2$$

$$= -\sum_i (Y_i - Y) \doteq 0$$

$$\hat{Y} = \frac{1}{N} \sum_i Y_i$$

$y^2$

$Y = Y_2(Y_i, Y)$

- if $loss(Y_i, Y) = |Y_i - Y| \Rightarrow$ median:

$$\frac{\partial}{\partial y} \sum_i |Y_i - Y| = -\sum_i sign(Y_i - Y) \doteq 0$$

$|y|$

$\leftarrow sign(y) = \frac{\partial}{\partial y}(Y_i, Y)$

$$= \sum_{i: Y_i - Y > 0} 1 \quad + \sum_{i: Y_i - Y < 0} (-1) \doteq 0$$

$\Rightarrow$ positive and negative sums must cancel

$$\#\{Y_i < \hat{Y}\} = \#\{Y_i > \hat{Y}\}$$

introduce: influence function of an instance: $\psi(Y_i, Y) = \frac{\partial}{\partial Y} loss(Y_i, Y)$

$\uparrow$ potential

$\hat{=}$ force, how an instance pulls or pushes     force
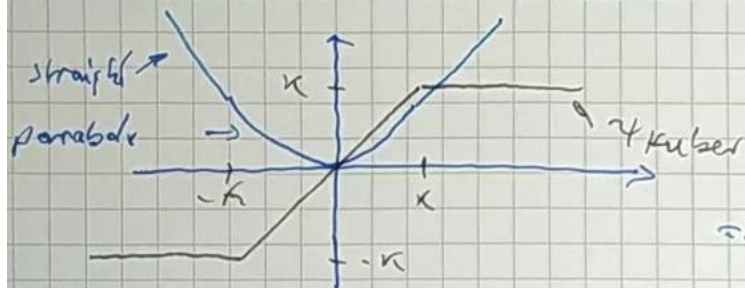
the representative

squared loss:   instances with large distance from the representative exert proportionally

larger forces   $\Rightarrow$ good for inliers

absolute loss:   all instances have the same force (except for sign)

$\Rightarrow$ good for outliers

combine these behaviors:



straight → 
parabola →

$$\Psi_{Huber}(Y_i, Y) = \begin{cases} Y_i - Y & \text{if } |Y_i - Y| \leq \kappa \leftarrow \text{inlier threshold} \\ \kappa \, \text{sign}(Y_i - Y) & \text{if } |Y_i - Y| \geq \kappa \leftarrow \text{outliers} \end{cases}$$
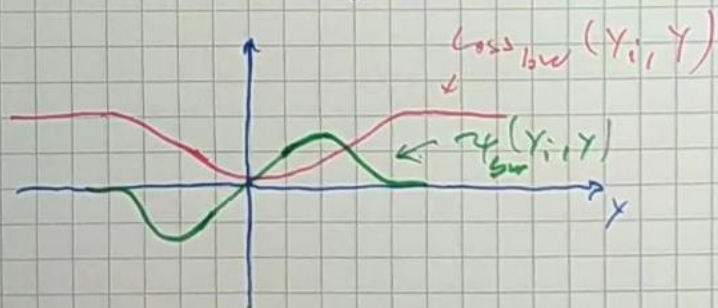
$\Rightarrow$ find loss by integration $\quad loss_{Huber}(Y_i - Y) =$

$$loss_{Huber}(Y_i, Y) = \begin{cases} \frac{1}{2}(Y_i - Y)^2 & \text{if } |Y_i - Y| \leq \kappa \quad \text{squared loss for inliers} \\ \kappa |Y_i - Y| - \frac{\kappa^2}{2} & \text{if } |Y_i - Y| \geq \kappa \end{cases}$$

$\kappa$: hyperparameter, if the inlier distribution is Gaussian,

optimal choice $\quad \kappa = 1.37 \, \sigma \quad$ std. dev.

interpretation: the force of outliers doesn't increase with distance $\Rightarrow$ more stable estimates

Q: can we eliminate the outlier influence entirely $\leftarrow$ far outliers have <u>zero</u> force

A: Yes, <u>biweight function</u>



$loss_{bw}(Y_i, Y)$

$\leftarrow \Psi_{bw}(Y_i, Y)$

$$\Psi_{bw}(Y_i, Y) = \begin{cases} (Y_i - Y)\left(1 - \left(\frac{Y_i - Y}{\kappa}\right)^2\right)^2 & \text{if } |Y_i - Y| \leq \kappa \\ 0 & \text{if } |Y_i - Y| > \kappa \end{cases}$$

$$loss_{bw}(Y_i, Y) = \begin{cases} 1 - \left(1 - \left(\frac{Y_i - Y}{\kappa}\right)^2\right)^3 & \text{if } |Y_i - Y| \leq \kappa \\ 1 & \text{if } |Y_i - Y| > \kappa \end{cases}$$
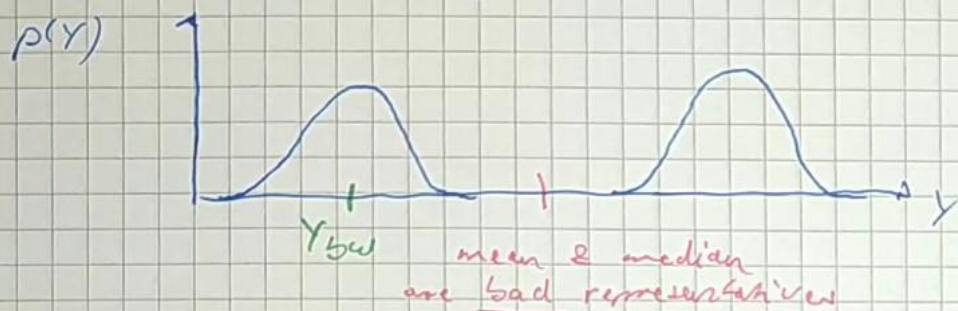
• advantage: only inliers and moderate outliers have influence

• disadvantage: optimization is non-convex, has many local optima

example for behavior:   suppose   the   true   distribution   is   bi-modal

$p(Y)$



$Y_{bw}$   mean & median   are **bad** representatives

bi-weight loss usually picks either cluster = two representative possibilities ($\equiv$ local optima)

$\Rightarrow$ gives a good solution for one camp, but entirely misses the other camp

$\Rightarrow$ in this case, no loss searching for a single representative works well, need a set of **diverse** representatives or a full estimate of the prob. density

unsolved problem: how many representatives are needed?

## analytical outlier detection / removal

- if the true model is linear:   $Y_i = X_i \beta + \varepsilon_i$ , we can derive analytical formula we learned in context of the bias-variance trade-off. the variance of $\beta$

$$\hat{\beta} \sim N(\beta^* , (X^TX)^{-1}\sigma^2) \qquad \varepsilon_i \sim N(0,\sigma^2)$$

$\rightarrow$ exponent of the Gaussian defines a distance

$$d(\hat{\beta}, \beta^*) = \frac{1}{D}(\hat{\beta} - \beta^*)^T \frac{X^TX}{\sigma^2} (\hat{\beta} - \beta^*) \qquad r_i = Y_i - X_i\hat{\beta}$$

↑ normalization for feature space dimension

approximate $\sigma^2$ by the mean squared error :   $\sigma^2 \approx \frac{1}{N-D}\sum_i r_i^2 = MSE$

- approximate distance by leave-one-out cross-validation :

$\hat{\beta}$ (full TS) as estimator for $\beta^*$   $\Rightarrow$ Cook's distance $d_i = \frac{(\hat{\beta}_{-i} - \hat{\beta})^T X^T X (\hat{\beta}_{-i} - \hat{\beta})}{D \cdot MSE}$

$\hat{\beta}_{-i}$ (solution of TS without instance i) for comparison

$d_i$ large $\Rightarrow$ instance $i$ has a very big influence on $\hat{\beta} \Rightarrow$ outlier

$d_i$ small ( from an $F$-distribution) $\Rightarrow$ normal behavior $\Rightarrow$ inliers

- we can compute $d_i$ analytically, after training rgly with full TS

define "hat matrix" $H = X \cdot \underbrace{(X^T X)^{-1} X^T}_{\text{pseudo inverse}}$

it puts the hat on the $Y_i$s : $\hat{Y_i} = X_i \hat{\beta} = X_i (X^T X)^{-1} X^T Y_i = H \cdot Y_i$

interpretation: $Y_i$ : noisy response $\quad \hat{Y_i}$ : corrected response on the regression line

$$d_i = \frac{r_i^2}{D \cdot MSE} \cdot \frac{H_{ii}}{(1 - H_{ii})^2} \qquad H_{ii} : \text{diagonal elems of } H \hat{=} \text{ "leverage" of instance } i$$

- outlier detection : place threshold (hyperparameter) on $d_i$ (ex.: $d_{max} = 1$
  $$d_{max} = \frac{4}{a} )$$

- for general non-linear problems, outlier detection is a largely unsolved problem

- many heuristic solutions, that work sometimes

- learn the inlier prob. density and define outliers as improbable under this

  density $\Rightarrow$ hard in high-dimensions $\qquad \Rightarrow$ Advanced Machine Learning