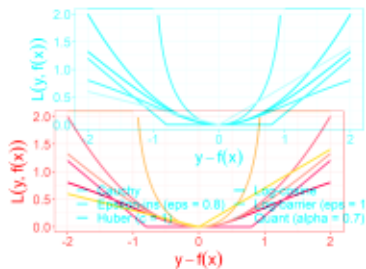


Introduction to Machine Learning

Advanced Risk Minimization

Advanced Regression Losses



— Cauchy
— Epsilon-ins (eps = 0.8)
— Huber (c = 1)
— Log-cosine
— Log-barrier (eps = 1)
— Quant (alpha = 0.7)

Learning goals

- Know the Huber loss

Learning goals

- Know the log-cosh loss

- Know the Cauchy loss

- Know the log-barrier loss

- Know the ϵ -insensitive loss

- Know the quantile loss

- Know the ϵ -insensitive loss

- Know the quantile loss

ADVANCED LOSS FUNCTIONS

Special loss functions can be used to estimate non-standard posterior components, to measure errors **customarily or which are designed to** have special properties like robustness.



Examples:

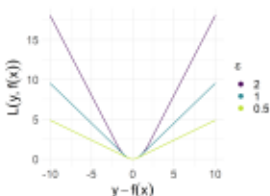
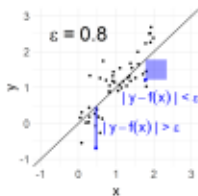
- Quantile loss: Overestimating a clinical parameter might not be as bad as underestimating it.
- Log-barrier loss: Extremely under- or overestimating demand in production would put company profit at risk.
- ϵ -insensitive loss: A certain amount of deviation in production does no harm, larger deviations do.

HUBER LOSS

$$L(y, f) = \begin{cases} \frac{1}{2}(y-f)^2 & \text{if } |y-f| \leq \epsilon \\ \epsilon|y-f| - \frac{1}{2}\epsilon^2 & \text{otherwise} \end{cases}, \quad \epsilon > 0$$



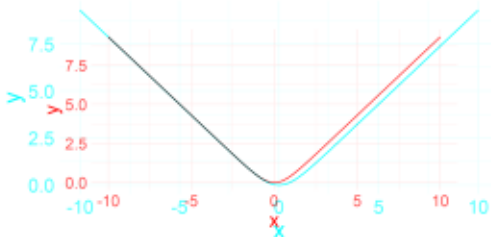
- Piece-wise combination of $L1/L2$ to have robustness/smoothness
- Analytic properties: convex, differentiable (once)



- Risk minimizer and optimal constant do not have a closed-form solution. To fit a model numerical optimization is necessary.
- Solution behaves like **trimmed mean**: a (conditional) mean of two (conditional) quantiles.

$$L(y, f) = \log(\cosh(|y - f|)) \quad \text{where } \cosh(x) := \frac{e^x + e^{-x}}{2}$$

- Logarithm of the hyperbolic cosine of the residual.
- Approximately equal to $0.5(|y - f|)^2$ for small f and to $|y - f| - \log 2$ for large f ; meaning it works mostly like $L2$ loss but is less outlier-sensitive.
- Has all the advantages of Huber loss and is, moreover, twice differentiable everywhere.



What is the idea behind the log-cosh loss?

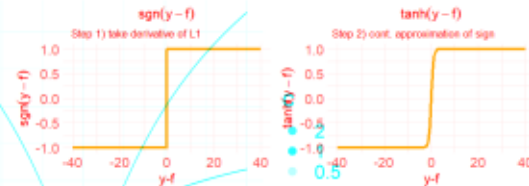
Essentially, we

$$L(y, f) = \frac{c^2}{2} \log \left(1 + \left(\frac{|y-f|}{c} \right)^2 \right), \quad c \in \mathbb{R}$$

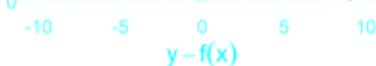
- take derivative of L1 loss w.r.t. $y - f$, which is the $\text{sgn}(y - f)$ function

- eliminate discontinuity at 0 by approximating $\text{sgn}(y - f)$ using the cont. differentiable $\tanh(y - f)$

- finally integrate the smoothed sign function "up again" to obtain smoothed L1 loss $\log(\cosh(y - f)) = \log(\cosh(|y - f|))$



The log-cosh approach to obtain a differentiable approximation of the L1 loss can also be extended to differentiable quantile/pinball losses.



The cosh(θ, σ) distribution:

The (normalized) reciprocal cosh(x) defines a pdf by its positivity on \mathbb{R} and since $\int_{-\infty}^{\infty} \frac{1}{\pi \cosh(x)} dx = 1$. if $|y - f| > \epsilon$

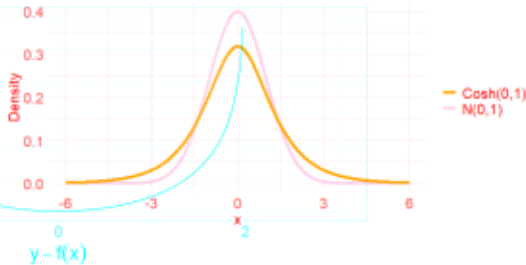
We can define a location-scale family of distributions (using θ and σ) resembling Gaussians with **heavier tails**.

- Behaves like L2 loss for small residuals.
- We use this if we don't want residuals larger than ϵ at all.
- No guarantee that the risk minimization problem has a solution.

It is easy to check that ERM using the log-cosh loss is equivalent to MLE of the cosh($\theta, 1$) distribution.

- Plot shows log-barrier loss for $\epsilon = 2$:

- $p(x|\theta, \sigma) = \frac{1}{\pi \sigma \cosh(\frac{x-\theta}{\sigma})}$
- $\mathbb{E}_{X \sim p}[X] = \theta$
- $\text{Var}_{X \sim p}[X] = \frac{1}{4}(\pi^2 \sigma^2)$
- $\hat{\theta}^{MLE} = \arg \max_{\theta} \prod_{j=1}^n \frac{1}{\pi \cosh(x_j - \theta)} \equiv$
- $\hat{\theta} = \arg \min_{\theta} \sum_{j=1}^n \log(\cosh(x_j - \theta))$

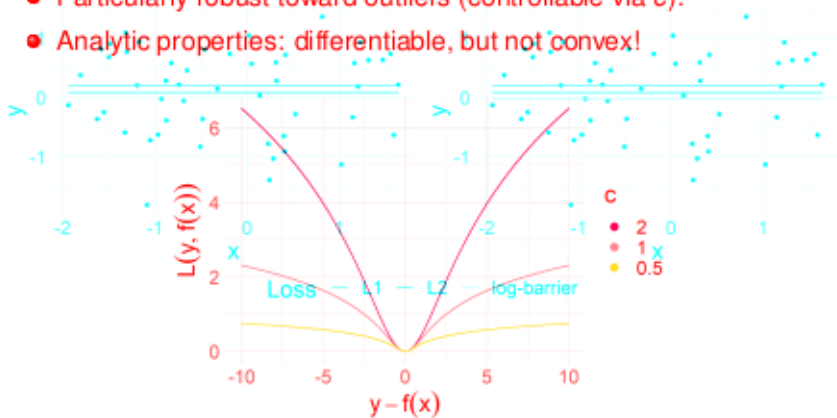


CAUCHY LOSS

- Note that the optimization problem has no (finite) solution if there is no way to fit a constant where all residuals are smaller than ϵ .

$$L(y, f) = \frac{c^2}{2} \log \left(1 + \left(\frac{|y - f|}{c} \right)^2 \right), \quad c \in \mathbb{R}$$

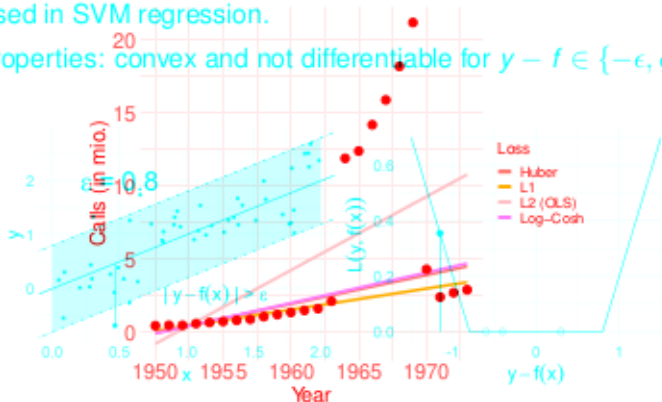
- Not feasible for $\epsilon = 1$. Feasible for $\epsilon = 2$.
- Particularly robust toward outliers (controllable via c).
- Analytic properties: differentiable, but not convex!



TELEPHONE DATA

We now illustrate the effect of using robust loss functions. The telephone data set contains the number of calls (in 10mio units) made in Belgium between 1950 and 1973 ($n = 24$). Outliers are due to a change in measurement without re-calibration for 6 years.

- Modification of L1 loss, errors below ϵ accepted without penalty.
- Used in SVM regression.
- Properties: convex and not differentiable for $y - f \in \{-\epsilon, \epsilon\}$.



LOG-BARRIER LOSS

$$L(y, f) = \begin{cases} (1 - \alpha)(f - y) & \text{if } y \leq f \\ \alpha(y - f) & \text{if } y \geq f \end{cases} \quad \text{if } |y - f| \leq \epsilon$$
$$L(y, f) = \begin{cases} -\frac{1}{\alpha} \log\left(1 - \left(\frac{y-f}{\epsilon}\right)^2\right) & \text{if } |y - f| \leq \epsilon \\ \infty & \text{if } |y - f| > \epsilon \end{cases}$$

- Extension of $L1$ loss (equal to $L1$ for $\alpha = 0.5$).
- Behaves like $L2$ loss for small residuals
- Weights either positive or negative residuals more strongly.
- We use this if we don't want residuals larger than ϵ at all
- $\alpha < 0.5$ ($\alpha > 0.5$) penalty to over-estimation (under-estimation)
- No guarantee that the risk minimization problem has a solution
- Risk minimizer is (conditional) α -quantile (median for $\alpha = 0.5$).
- Plot shows log-barrier loss for $\epsilon = 2$:

