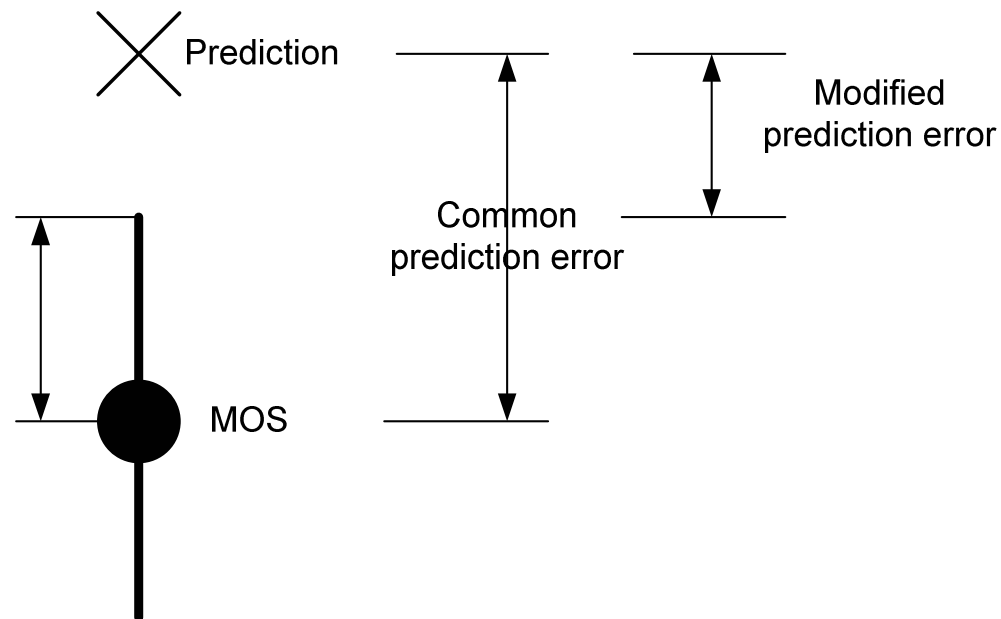


# Epsilon-insensitive r.m.s.e



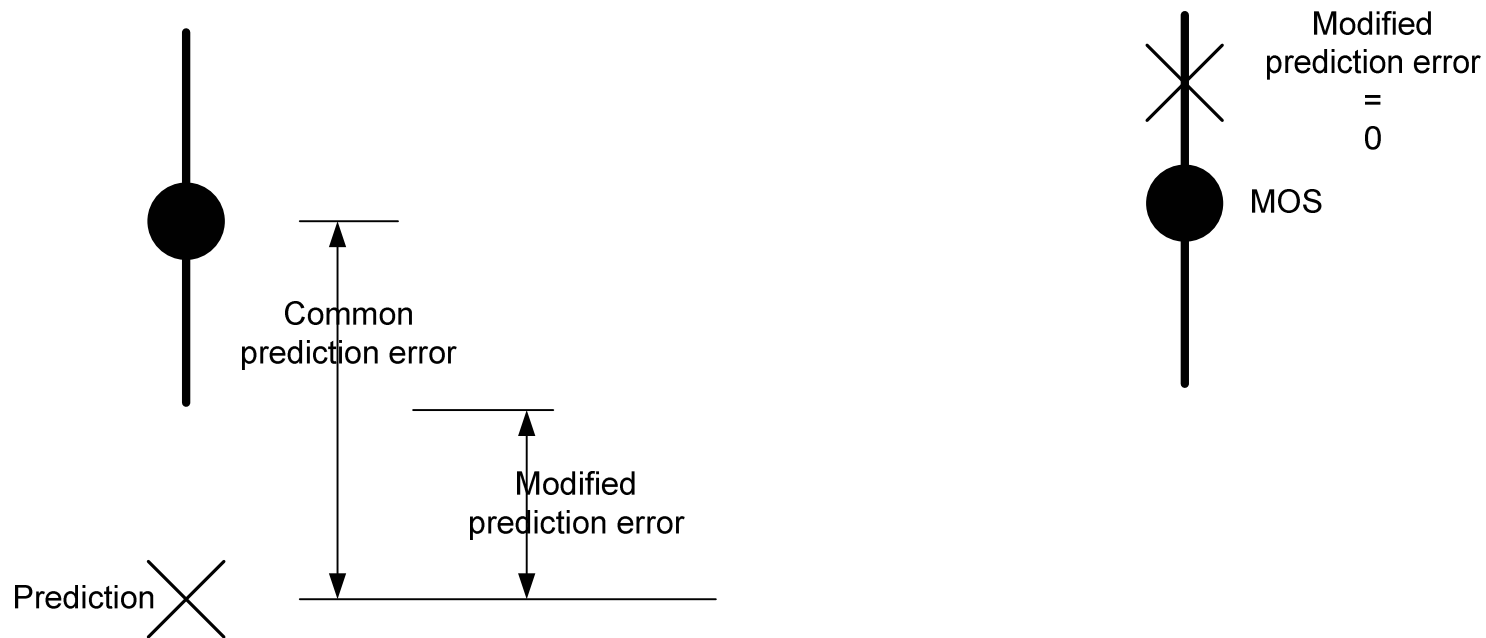
**Basic idea: Taking into account the ci95 confidence interval for calculating prediction errors and resulting rmse**



# Epsilon-insensitive r.m.s.e



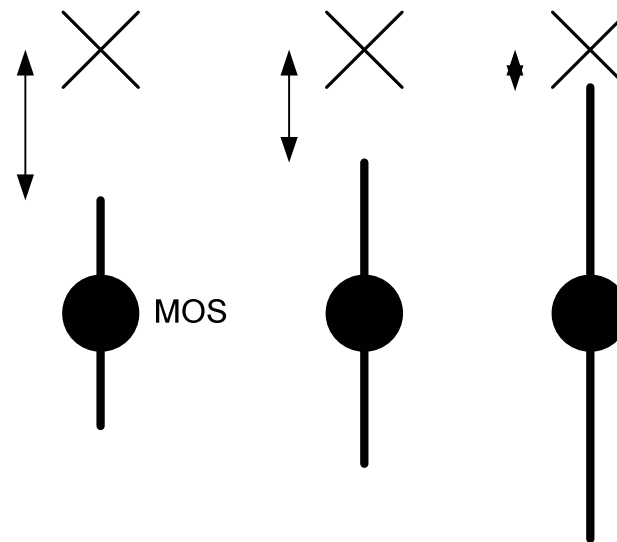
**Basic idea: Taking into account the ci95 confidence interval for calculating prediction errors and resulting rmse**



# Epsilon-insensitive r.m.s.e



**The Prediction Error becomes reduced in case the MOS is more inconfident (i.e. due to less votes or wide distribution of votes)**



# How to calculate rmse\*?



$$Perror(i) = \max(0, |MOSLQS(i) - MOSLQO(i)| - ci_{95}(i))$$

$$rmse^* = \sqrt{\left( \frac{1}{N-d} \sum_N Perror(i)^2 \right)}$$

Remark: There are some worth special rules in case the  $ci_{95} \rightarrow 0$  at the scale boundaries.

# How to compare models?



**For each data set and model a Distance  $d$  between the best performing model and the others are calculated. It considers statistical significance too.**

$$d_{k,v} = \max(0, rmse_{k,v}^* - rmse_{k,b}^* \times F(0.05, N_k, N_k))$$

# How to compare models?



**The Distances for each model are averaged across the datasets. If desired the data sets can be weighed.**

**Finally, models can be selected those are the best and statistically equivalent.**

$$p_v = \sum_{k=1}^M w_k \times d_{k,v}$$

$$t_v = \max\left(0, \frac{p_v}{(p_{\min} + c)} - F(0.05, K, K)\right)$$