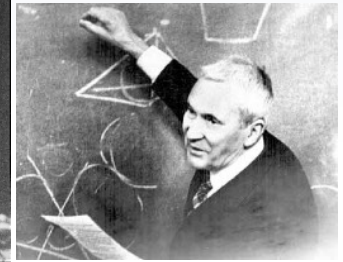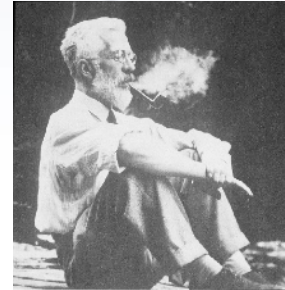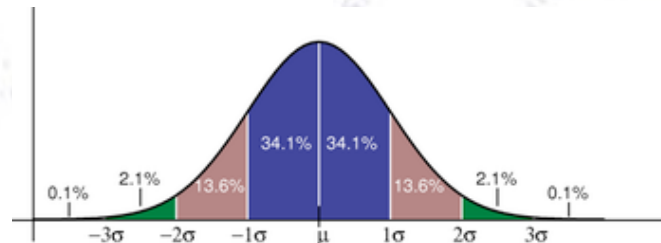# Applied ML

## Loss Functions

### Troels C. Petersen (NBI)

*"Statistics is merely a quantisation of common sense - Machine Learning is a sharpening of it!"*
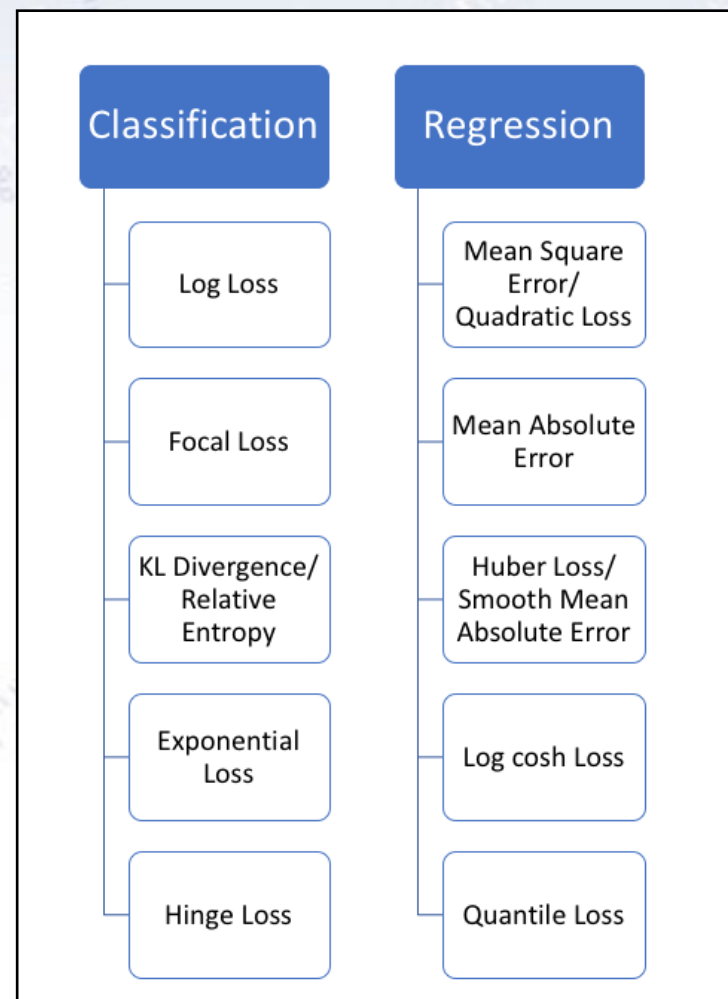
# What loss function to use?

The choice of loss function depends on the problem at hand, and in particular what you find important!

In classification:
- Do you care how wrong the wrong are?
- Do you want pure signal or high efficiency?
- Does it matter what type of errors you make?

In regression:
- Do you care about outliers?
- Do you care about size of outliers?
- Is core resolution vital?

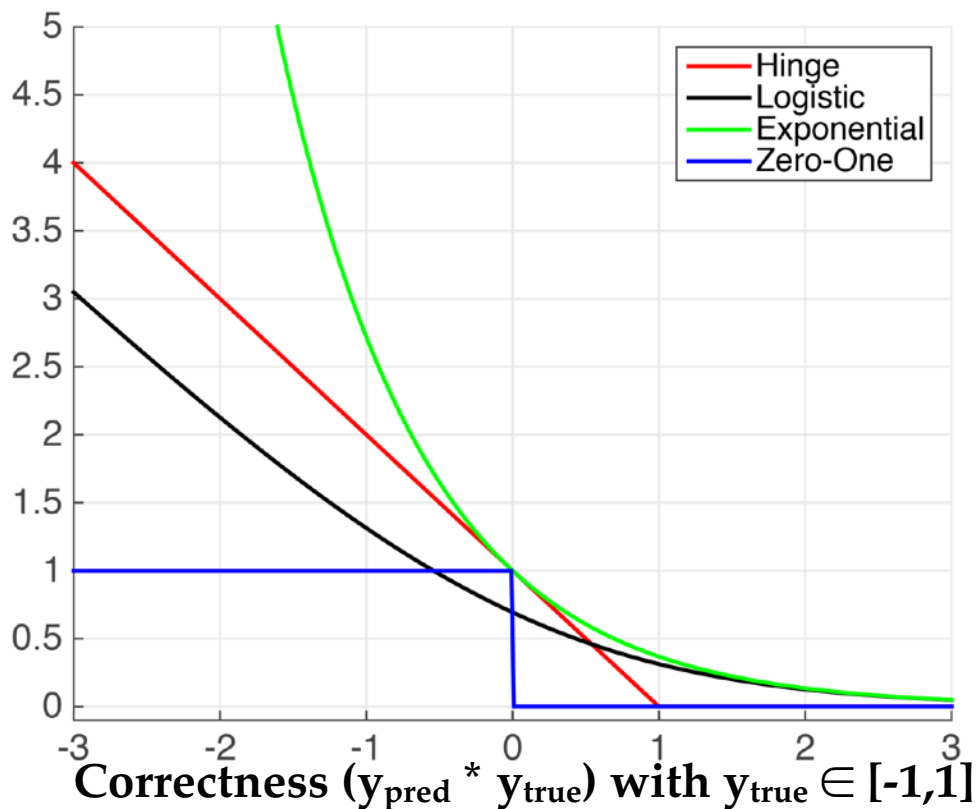| Classification | Regression |
| --- | --- |
| Log Loss | Mean Square Error/ Quadratic Loss |
| Focal Loss | Mean Absolute Error |
| KL Divergence/ Relative Entropy | Huber Loss/ Smooth Mean Absolute Error |
| Exponential Loss | Log cosh Loss |
| Hinge Loss | Quantile Loss |

# Classification

# What loss function to use?

The choice of loss function depends on the problem at hand, and in particular what you find important!

Loss functions for classification



**Correctness ($y_{pred}$ * $y_{true}$) with $y_{true} \in$ [-1,1]**

Hinge
Logistic
Exponential
Zero-One

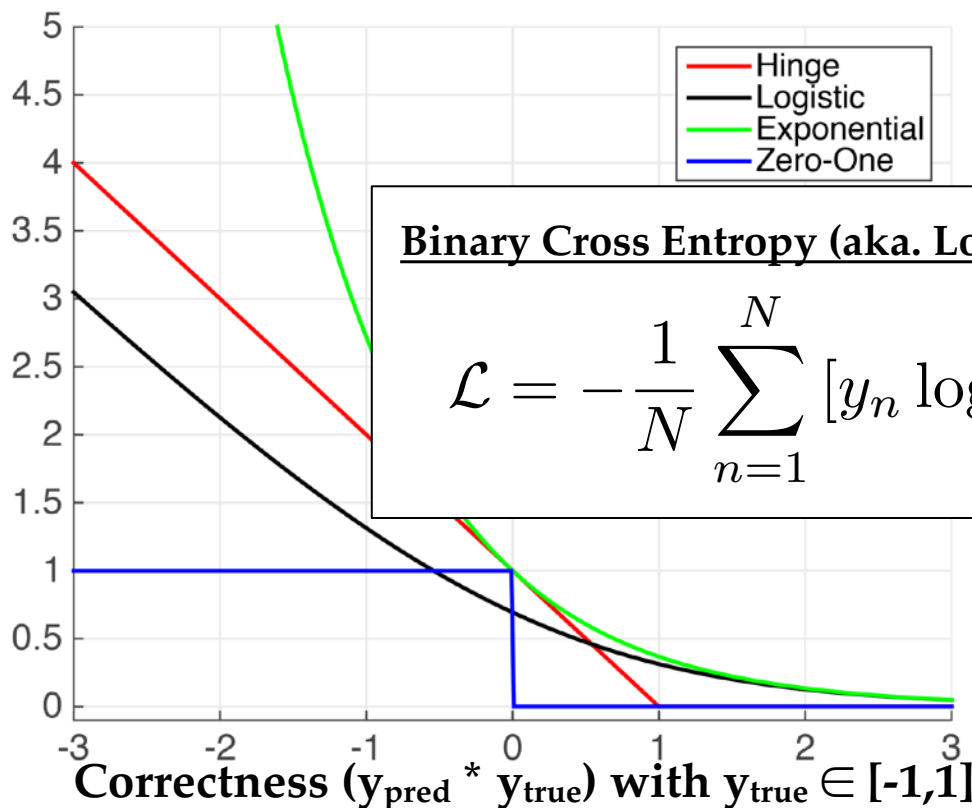| Classification | Regression |
|---|---|
| Log Loss | Mean Square Error/ Quadratic Loss |
| Focal Loss | Mean Absolute Error |
| KL Divergence/ Relative Entropy | Huber Loss/ Smooth Mean Absolute Error |
| Exponential Loss | Log cosh Loss |
| Hinge Loss | Quantile Loss |

# What loss function to use?

The choice of loss function depends on the problem at hand, and in particular what you find important!

Loss functions for classification



**Binary Cross Entropy (aka. LogLoss or Logistic Loss) with $y_{true} \in [-1,1]$:**

$$\mathcal{L} = -\frac{1}{N} \sum_{n=1}^{N} [y_n \log \hat{y}_n + (1 - y_n) \log(1 - \hat{y}_n)]$$

**Correctness ($y_{pred} * y_{true}$) with $y_{true} \in [-1,1]$**

# Unbalanced data

If the data is unbalanced, that is if one outcome/target is much more abundant than the alternative, case has to be taken.

Example: You consider data with 19600 (98%) healthy and 400 (2%) ill patients. An algorithm always predicting "healthy" would get an accuracy score of 98%!
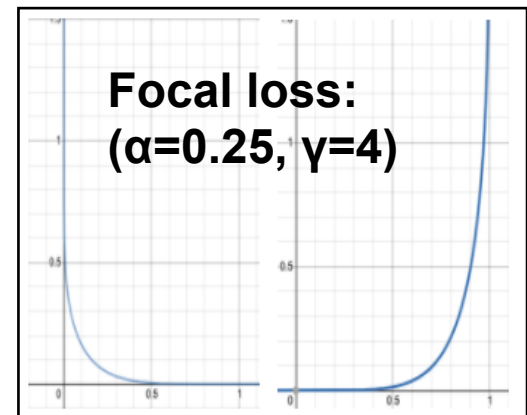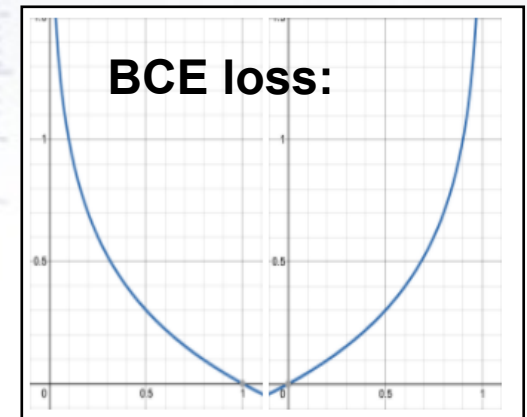
In this case, using Area Under Curve (AUC) or F1 for loss is better. An alternative is "focal loss", which focuses on the lesser represented cases:

**BCE loss:**



Binary Cross Entropy loss:

$$\mathcal{L} = -\frac{1}{N} \sum_{n=1}^{N} [y_n \log \hat{y}_n + (1 - y_n) \log(1 - \hat{y}_n)]$$

Focal loss:

$$\mathcal{L} = -\frac{1}{N} \sum_{n=1}^{N} [(1 - \alpha) y_n^\gamma \log \hat{y}_n + (1 - y_n)^\gamma \log \alpha (1 - \hat{y}_n)]$$

**Focal loss:**
**(α=0.25, γ=4)**

# Multi-Classification

If the problem contains multiple classes (i.e. not just two), then an extension of the loss function is needed.

This is simply done by including a term for each class of events as follows:

Binary Cross Entropy loss:

$$\mathcal{L}(\hat{y}, y) = -\frac{1}{N} \sum_{n=1}^{N} \left[ y \log \hat{y} + (1 - y) \log(1 - \hat{y}) \right]$$

Multi-Classification loss:

$$\mathcal{L}(\hat{y}, y) = -\frac{1}{N} \sum_{n=1}^{N} \sum_{k}^{K} y^{(k)} \log \hat{y}^{(k)}$$

# Regression

# What loss function to use?

The choice of loss function depends on the problem at hand, and in particular what you find important!

Loss functions for regression



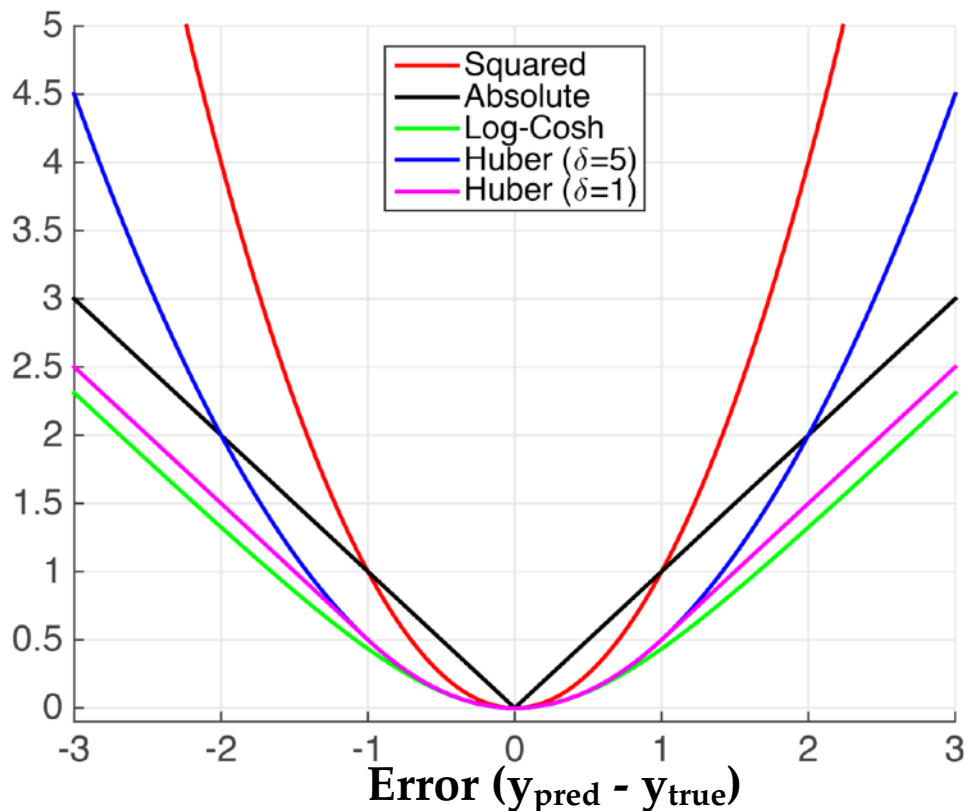Error ($y_{pred} - y_{true}$)



Discussion of regression loss functions

# What loss function to use?

The choice of loss function depends on the problem at hand, and in particular what you find important!

Loss functions for regression



**Error (y$_{pred}$ - y$_{true}$)**

### Squared Loss:
- Most popular regression loss function
- Estimates <u>Mean</u> Label
- ADVANTAGE: Differentiable everywhere
- DISADVANTAGE: Sensitive to outliers

### Absolute Loss:
- Also a very popular loss function
- Estimates <u>Median</u> Label
- ADVANTAGE: Less sensitive to noise
- DISADVANTAGE: Not differentiable at 0

### Huber Loss:
- ADVANTAGE: "Best of Both Worlds" of <u>Squared</u> and <u>Absolute</u> Loss.
- DISADVANTAGE: Only once-differentiable

### LogCosh Loss:
- ADVANTAGE: "Best of Both Worlds" of <u>Squared</u> and <u>Absolute</u> Loss.
- ADVANTAGE: Similar to Huber Loss, but twice differentiable everywhere.
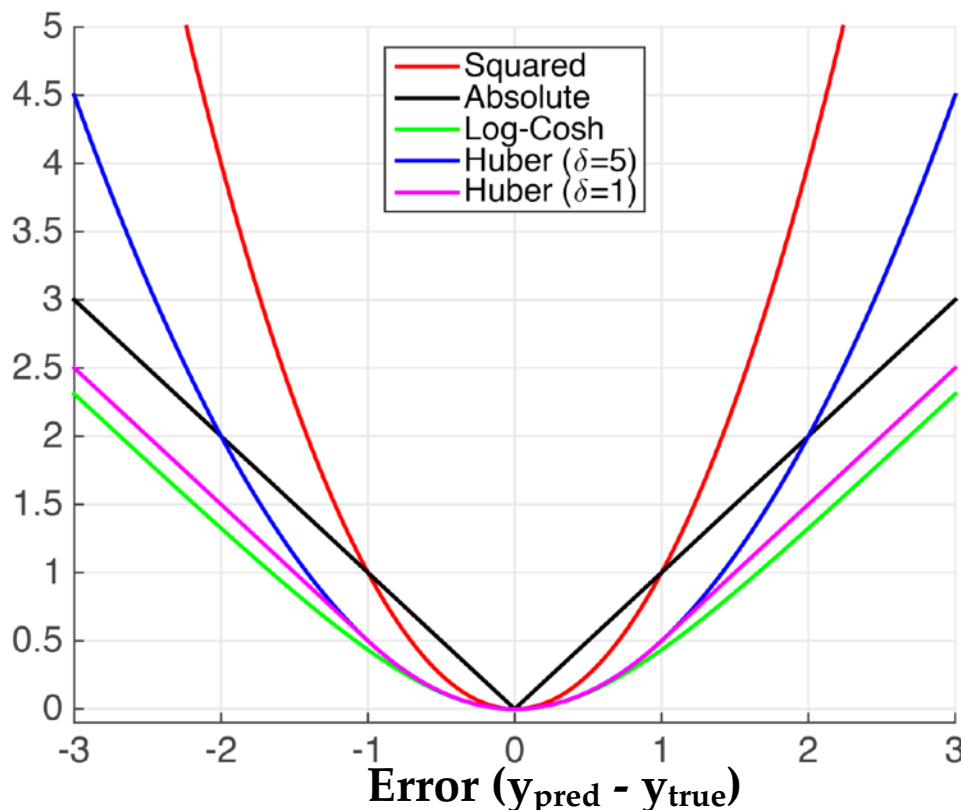
<u>Discussion of regression loss functions</u>

# What loss function to use?

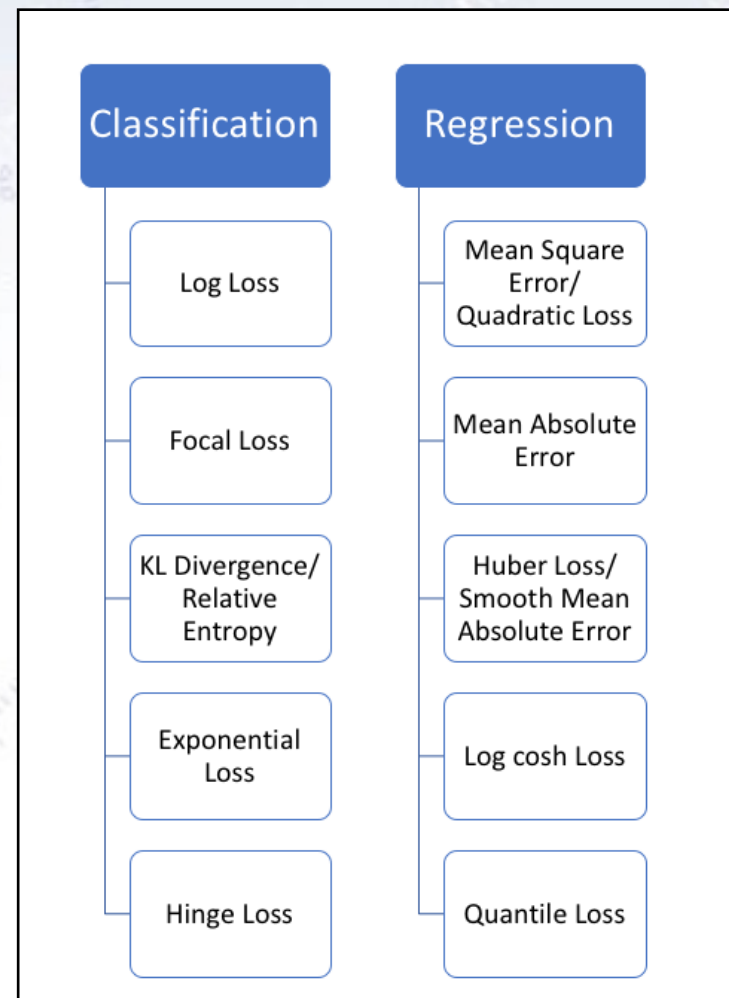The choice of loss function depends on the problem at hand, and in particular what you find important!

In classification:
- Do you care how wrong the wrong are?
- Do you want pure signal or high efficiency?
- Does it matter what type of errors you make?

In regression:
- Do you care about outliers?
- Do you care about size of outliers?
- Is core resolution vital?

Ultimately, the loss function should be tailored to match the wishes of the user. This is however not always that simple, as this might be hard to even know!

| Classification | Regression |
| --- | --- |
| Log Loss | Mean Square Error/ Quadratic Loss |
| Focal Loss | Mean Absolute Error |
| KL Divergence/ Relative Entropy | Huber Loss/ Smooth Mean Absolute Error |
| Exponential Loss | Log cosh Loss |
| Hinge Loss | Quantile Loss |

# XGBoost algorithm

In order to "punish" complexity, the cost-function has a regularised term also:

$$\mathcal{L}(\phi) = \sum_i l(\hat{y}_i, y_i) + \sum_k \Omega(f_k)$$

$$\text{where } \Omega(f) = \gamma T + \frac{1}{2}\lambda\|w\|^2$$

**Table 1: Comparison of major tree boosting systems.**

| System | exact greedy | approximate global | approximate local | out-of-core | sparsity aware | parallel |
|---|---|---|---|---|---|---|
| **XGBoost** | yes | yes | yes | yes | yes | yes |
| pGBRT | no | no | yes | no | no | yes |
| Spark MLLib | no | yes | no | no | partially | yes |
| H2O | no | yes | no | no | partially | yes |
| scikit-learn | yes | no | no | no | no | no |
| R GBM | yes | no | no | no | partially | no |

# XGBoost algorithm

In order to "punish" complexity, the cost-function has a regularised term also:

$$\mathcal{L}(\phi) = \sum_i l(\hat{y}_i, y_i) + \sum_k \Omega(f_k)$$

$$\text{where } \Omega(f) = \gamma T + \frac{1}{2}\lambda \|w\|^2$$

**Table 1: Comparison of major tree boosting systems.**

| System | exact greedy | approximate global | approximate local | out-of-core | sparsity aware | parallel |
|---|---|---|---|---|---|---|
| **XGBoo** | | | | | | |
| pGBRT | | | | | | |
| Spark M | | | | | | |
| H2O | | | | | | |
| scikit-learn | yes | no | no | no | no | no |
| R GBM | yes | no | no | no | partially | no |

Generally, all constraints or priors should be included into the model through additions to the loss function.