# Binäres / Dichotomes Merkmal Y, ein oder mehrere nominalskalierte Merkmale X

(nur zwei Kategorien)

Setting: $y \in \{0, 1\}$ binary target variable; feature vector X

often $y = 1$ as "positive case", $y = 0$ as "negative case"

Ziel: Vorhersage / Diagnose $\hat{y}^{(i)}$ based on score/probability $f(x^{(i)})$ with some threshold $c$ (scalar or vector) e.g. using logistic regression

Prognosegüte [Performance Measure] gesucht für diese binary classification task

## Confusion Matrix

| Predicted Class | | True Class | | |
|---|---|---|---|---|
| | | Y pos.(1) | Y neg.(0) | |
| | $\hat{y}$ pos. | True Positive (TP) | False Positive (FP) | # pos. pred. |
| | $\hat{y}$ neg. | False Negative (FN) | True Negative (TN) | # neg. pred. |
| | | # positive cases $(n_+)$ | # negative cases $(n_-)$ | n |

$P(y=1)$ Prevalence $= \dfrac{\text{\# positive cases}}{\text{Total population (n)}}$

### ROC-metrics

**Sensitivity / Recall**
True positive Rate (TPR): $P_{TPR} = \dfrac{TP}{TP+FN}$
↳ Proportion of TP among pos. cases $P(\hat{y}=1|y=1)$

**Specifity**
True Negative Rate (TNR): $P_{TNR} = \dfrac{TN}{TN+FP}$
↳ Proportion of TN among neg. cases $P(\hat{y}=0|y=0)$

**Precision** $P(y=1|\hat{y}=1)$
Positive Predictive Value: $PPV = \dfrac{TP}{TP+FP}$ If we predict $\hat{y}=1$, how likely is it a true 1? → Proportion of "real"(correctly class.) pos. among all pos. pred.
(jeweils $1 - \dots$ für False-Anteil)

Negative Predictive Value: $NPV = \dfrac{TN}{TN+FN}$ If we predict $\hat{y}=0$, how likely is it true 0? → Proportion of correctly classified neg. among all neg. pred.

$ACC = P(\hat{y}=y)$
Accuracy (ACC): $P_{ACC} = \dfrac{TP+TN}{n}$ Proportion of correct predictions $ACC = P(\hat{y}=1, y=1) + P(\hat{y}=0, y=0)$

$ACC = \underbrace{P(\hat{y}=1|y=1)}_{Sensitivity} \cdot \underbrace{P(y=1)}_{prevalence} + \underbrace{P(\hat{y}=0|y=0)}_{Specifity} \cdot \underbrace{P(y=0)}_{1-prevalence}$ Fixed weight of Sens & Spec

Bad metric for small prevalence "Accuracy Paradox" → use e.g. F1-score instead for highly class-imbalance

**Fall-out**
False Positive Rate: $FPR = \dfrac{FP}{TN+FP}$ ↓ # neg. cases $= 1 - TNR$

**Miss Rate**
False Negative Rate: $FNR = \dfrac{FN}{TP+FN}$ ↓ # pos. cases $= 1 - TPR$

Don't confuse with False Discovery Rate: $FDR = FP / (TP+FP) = 1-PPV$

### F1-score: Its very difficult to achieve high PPV and high TPR simultaneously

A classifier predicting more positives will be more sensitive (higher TPR), but will also tend to give more false positives (lower TNR, lower PPV)

A classifier predicting more negatives will be more precise (higher PPV), but will also produce more false negatives (lower TPR)

F1 score aims to maximize PPV and TPR $\quad P_{F1} = \dfrac{1}{\frac{1}{2}\left(\frac{1}{PPV} + \frac{1}{TPR}\right)} = 2\dfrac{PPV \cdot TPR}{PPV+TPR} = \dfrac{2TP}{2TP+FN+FP}$ harmonic mean of precision & recall  doesn't account for number of True Negatives!

Balances them, but leans towards the smaller value
PPV or TPR = 0: F1 = 0  by convention (also for better 0 or NP) Always predict neg: F1 = 0

⇒ max. $P_{F1} \in [0,1]$ best case: PPV=TPR=1 (no FP,FN) ↝ F1 = 1  Always predict pos: F1 = 2·PPV/(PPV+1) = $2n_+/(n_+ + n)$ (TPR=1) big for big $n_+$, small for small $n_+$

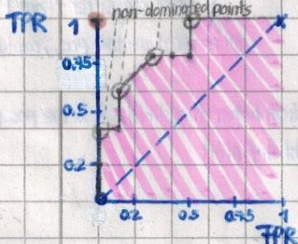Indicates Reliability of Predictions $\hat{y} = 1$ (Sensitivity/Recall, 1-Specifity) Interesting: FPR, TPR stay the same if we scale ratio $n_+/n_-$ at predict time

Insensitive to scaling of true class distr. ratio (test, but not true for training)

## ROC Curve  Comparing classifiers by plotting (TPR, FPR) values for all possible thresholds $c$, want max TPR & min. FPR

could also use two other trade-off metrics such as (TPR, PPV) where we want to max both



ROC curve invariant to monotonic transf. of scores $\quad h(x) = [\hat{\pi}(x) \geq c]$ at all pos. pred., $c=0$, always ein in (1,1)

1.) Rank test observations on decreasing score/prob. 2.) Start with $c = 1$ for probabilistic classif., so we start in (0,0) & no pos. predictions

3.) Iterate through all possible thresholds $c\downarrow$, discrete steps in between one or several obs with same score
If corresponding true $y = 1$: Move TPR up by $1/n_+$ as we have one TP point more (pos)
If corresponding true $y = 0$: Move FPR right by $1/n_-$ as we have one FP point more (neg)
We get diagonal step if, for certain score, we have more than one obs. and they have different y ⇒ we add TP, FP points at the same time

- The closer the ROC curve is to the left axis with ideal dream point (TPR=1, FPR=0) the better our classifier is  top left corner  ROC of "perfect" classifier, no FP! (until below all pos.)
- Diagonal represents worst possible classifier that produces random label predictions e.g. 20% pos, 80% neg  all thresholds $c \in (0,1)$ perfectly separate y-classes, AUC=1
  ↳ (non-discriminant) random prediction means TP=FP, FN=TN ⇒ TPR=FPR [AUC=0.5]  "Scorewise ranks all pos. obs $y^{(i)}$ > neg. obs. y"
- Curve below diagonal is impossible. We can then just flip predicted labels (0→1 and 1→0) which reflects curve at diagonal. Curve below diag. is such an awful classifier that we would never implement it. Flipping makes it so better  $F^{(0)}PR_{new} = 1 - F^{(1)}PR$

Probabilistic classif: Average predicted prob. should be ≈ prevalence (well calibrated model) ⇒ pred.prob. estimate real risk/chance  captures prevalence!

Finding best threshold? Small thresholds (ROC-points in top right corner) will very liberally predict pos. class ⇒ higher TPR, but also (potentially) higher FPR
for specific classifier  Big Thresholds (ROC-points in bottom left corner) will very conservatively predict pos. class ⇒ lower FPR, however lower TPR as well

If class imbalances & FP vs. FN costs are unknown, which is often the case, we need to decide threshold manually by visual cue
⇒ identify non-dominated points, assess their TPR, FPR combo, decide which combo is best for task & domain knowledge, and choose respective threshold

; treats Sensitivity & Specifity equally and ignores PPV, NPV!

concordant pairs $(y^{(i)}$pos, $f(x^{(i)}))$, $(y^{(j)}$neg, $f(x^{(j)}))$ if $f$  !concord. if same, P pairs mit Ties in y

**AUC** Area under the ROC-curve  Rank-based interpretation: Fraction of all pairs $(y^{(i)}$pos$, y^{(j)}$neg$)$ where $f(x^{(i)}) > f(x^{(j)})$ $\quad AUC = \dfrac{N_C + N_C/2}{N} = \dfrac{N_C}{N_{unique\, y}}$
technically 0 possible
AUC $\in [0.5, 1]$ Larger AUC = better classifier  $P(\text{score}(\text{random pos. } y^{(i)}) > \text{score}(\text{random neg. } y^{(j)}))$ "Probability of ranking random obs. correctly"

AUC integrates over all feasible thresholds bzw. corresponding points (TPR, FPR)   $AUC = \int_{0}^{1} ROC_f(t)\,dt$,   empirical $\hat{AUC} = \frac{N_c}{N_{unique\,y}}$

## Partial AUC
If certain value for TPR or FPR is no acceptable, then it might be useful to fix TPR or FPR to required value minimum and optimize the other metric based on that constraint. Focus only on relevant parts of ROC-curve for our use-case!

↓ integrate over region after cut-off

- Example: FPR > 0.2 is not acceptable
  ⇒ cut out region FPR > 0.2 (vertical ☐ region remains), focus on FPR ≤ 0.2

- Or TPR < 0.8 is so awful that we would never use a classifier in that region.

Rescale Partial AUC metric so that not only does it normalize to [0,1] but also make random labeling correspond to $AUC_{partial} = 0.5$

$$f_{AUC,corrected} = \frac{1}{2}\left(1 + \frac{AUC - AUC_{min}}{AUC_{max} - AUC_{min}}\right)$$
with $AUC_{min} \triangleq$ Area under diagonal in constrained region
and $AUC_{max} \triangleq$ Area of constrained region

- We could even put constraints on both TPR, FPR and calculate the remaining small area (2way partial AUC)

## Multiclass AUC
ROC-curves and therefore AUC are only defined for binary classification tasks

Define conditional $AUC(k\,|\,L)$ where $k$ is positive and $L$ is negative class
↳ For computation subset rows with $y = k$ or $y = L$ from dataset

Careful: Unlike binary case where $AUC(1\,|\,0) = AUC(0\,|\,1)$ it is $AUC(k\,|\,L) \neq AUC(L\,|\,k)$!
E.g with probability $1 - \pi_1 = \pi_0$ but with three or more classes that is generally not true

$$\Rightarrow \text{using 1-vs-1 scheme}: AUC_{MC} = \frac{1}{g(g-1)} \sum_{k \neq \ell} AUC(k\,|\,\ell) \in [0,1]$$

Sidenote] There are other schemes like 1-vs-rest where we only have to average over $g$ AUC values, $AUC(1\,|\,Rest) \ldots \ldots AUC(g\,|\,Rest)$
This is computationally easier, but create imbalanced classes by design even if original classes were rel. balanced
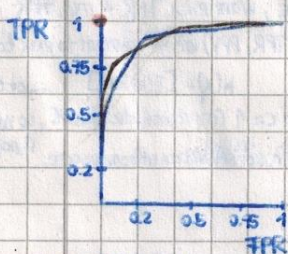
Precision Recall kind of ignoring TN(R) and FPR

## Precision-Recall-Curves
Comparing classifiers by plotting (PPV, TPR) instead for all possible thresholds $c$, want max. both equally weighted
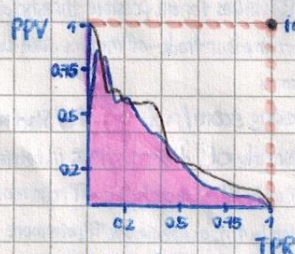Might be better than ROC for highly class-imbalanced data $n_- \gg n_+$
often then we are less interested in high TNR ⟷ low FPR, but more in high PPV which reacts more sensitive to abs. $\Delta FP$

For fixed given classifier:
ROC curves for very differently balanced classes could look very similar. PR curves change drastically from balanced to imbalanced classes



- ideal dream point (PPV = 1, TPR = 1) in top right corner

$$PPV = \frac{TP}{TP + FP} \qquad TPR = \frac{TP}{TP + FN}$$

ideal classifier from PR perspective
all pos. pred are true pos., capture all TP

PR Plot started links oben bei (1, 0) $\triangleq$ extrem hoher threshold s.d. keine pos. pred. erfolgen
PPV dann per Konvention = 1, weil 0/0 nicht definiert

For high TPR: PPV way below 0.5 so even worse than a coinflip guess
ROC plot could make us think that we have two good classifiers here (and similar). PR plot shows them to be quite different and actually pretty bad!

Good: Curve dominates in ROC ⟷ Curve dominates in PR

- In practice, if we only compare a few models on a single task, its not a bad idea to plot both curves and observe
- (Partial) PR AUC can be used for tuning