

При постепенном увеличении строгости модели (движении  $PD'$  по убыванию) будут возрастать исследуемые показатели. В их возрастании можно выделить два этапа:

1. До пика распределения  $N$  – быстро возрастает  $TP$  с нарастающим темпом,  $FP$  также возрастает с нарастающим темпом но не так быстро как  $TP$ ;
2. После пика распределения  $N$  – рост  $TP$  начинает постепенно замедляться, рост  $FP$  ускоряется а  $FPR$  постепенно начинает догонять ушедший вперед  $TPR$ .

Закачивает этот процесс тем, что  $TP$  включает в себя все наблюдения с проявившимся признаком, а  $FP$  все наблюдения в которых признака не наблюдалось. Из того следует, что  $FPR = 1$ ,  $TPR = 1$ .

Теперь, отложим на оси абсцисс  $FPR$  а на оси ординат  $TPR$  и получим самый популярный способ оценки классифицирующей способности модели – ROC кривую.

Опираясь на рассуждения выше можно сказать, что чем более непохожие распределения по  $PD$  у наблюдений с проявившимся признаком и без, тем круче будет ROC кривая и, при том, она проходит через точки  $(0,0)$  и  $(1,1)$ . В подтверждение своих слов приводим рисунок Е.3 – на нем, слева, распределения по  $PD$  наблюдений с проявившимся признаком и без представлены в виде нормальных распределений с фиксированной дисперсией, математические ожидания этих распределений постепенно приближаются друг к другу. Справа для каждой из ситуации построена ROC кривая – видно, как по степени сближения распределений, ROC кривая выпрямляется.

В описанных условиях удобной метрикой, в одно число, классифицирующей способности может выступить AUC ROC кривой – площадь под ROC кривой. Далее, для краткости, описанный показатель, будем называть AUC.

Еще одной связанной метрикой качества модели выступает статистика Колмогорова-Смирнова (KS-статистика). Обычно эта статистика используется для проверки гипотезы о соответствии некоторого наблюдаемого ряда предполагаемому теоритическому распределению, но в нашем случае сверять мы будем распределение наблюдений с проявлением признака и без. В данном случае статистика определяется по формуле:

$$KS = \max_{PD'} |\hat{F}_P(PD') - \hat{F}_N(PD')|, \quad (2.3)$$

где  $\hat{F}_P(PD')$  – эмпирическая функция распределения наблюдений с проявлением признака по  $PD'$ ;

$\hat{F}_N(PD')$  – эмпирическая функция распределения наблюдений без проявления признака по  $PD'$ .

Графически эту статистику удобно представлять себе в виде рисунка 2.2.

Что из формулы, что из рисунка становится очевидно – KS статистика показывает на сколько максимум накопленные промежуточным итогом наблюдения одного класса отстают от наблюдений другого класса.