
Калибровка моделей и uncertainty estimation

Смирнов Тимофей

Как выдавать предикт

$$b(x) = (b_1(x), \dots, b_l(x)) \in [0, 1]_{\|\cdot\|_1=1}^l, \\ a(x) = \arg \max_j (b_j).$$

Вероятность верности предикта

$$b_{a(x)}(x) \approx \mathbf{P}(y(x) = a(x)).$$

Как это посчитать

$$\frac{1}{|B|} \sum_{x \in B} b_k(x) \text{ vs } \frac{1}{|B|} \sum_{x \in B} I[y(x) = k].$$

Диаграмма калибровки

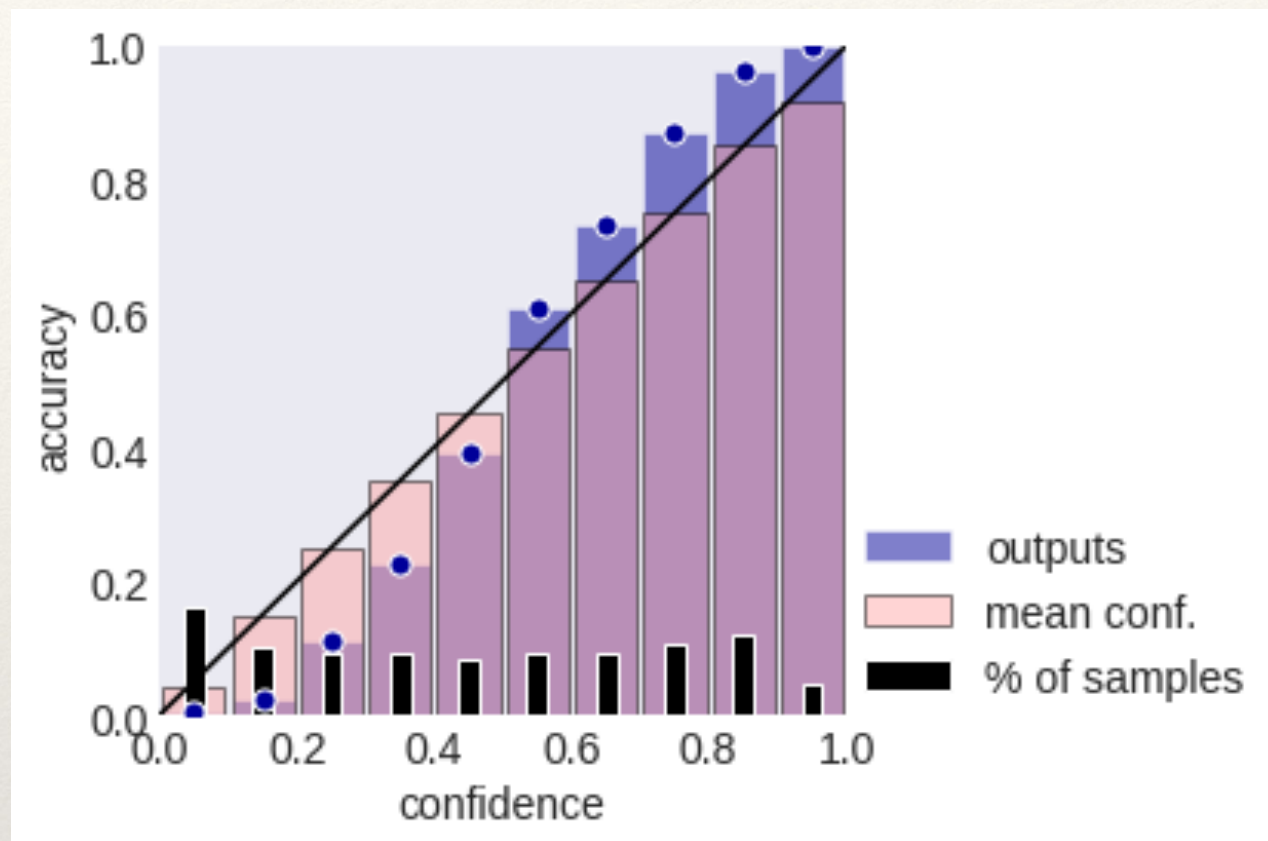
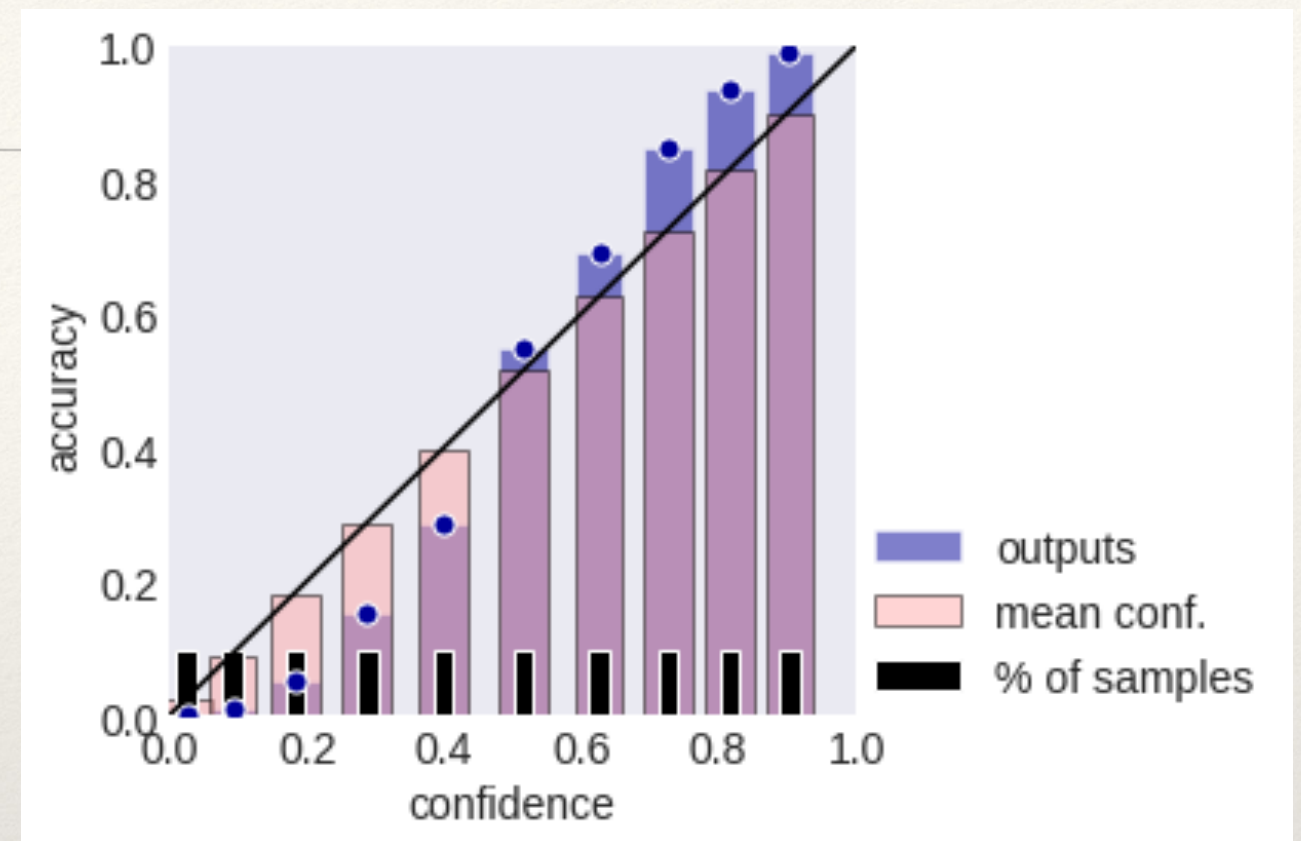


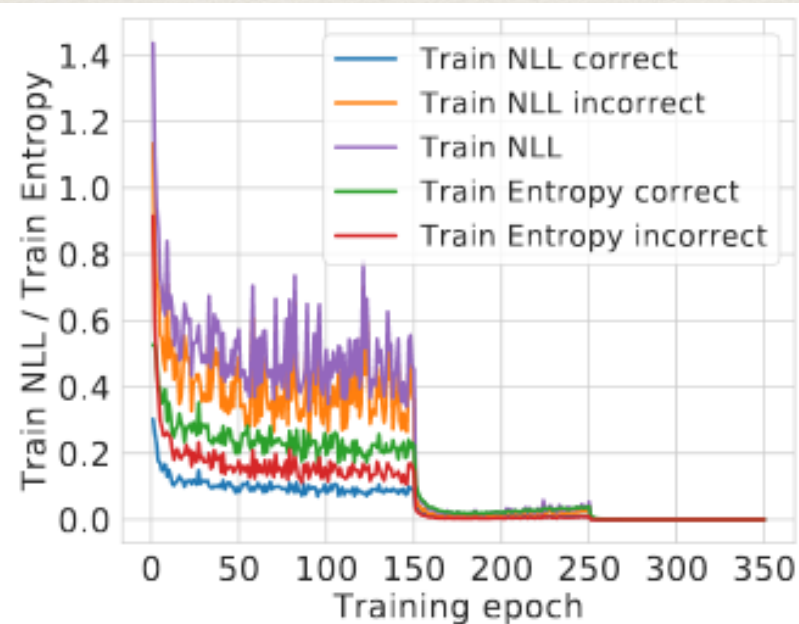
Диаграмма калибровки с
равномощными бинами



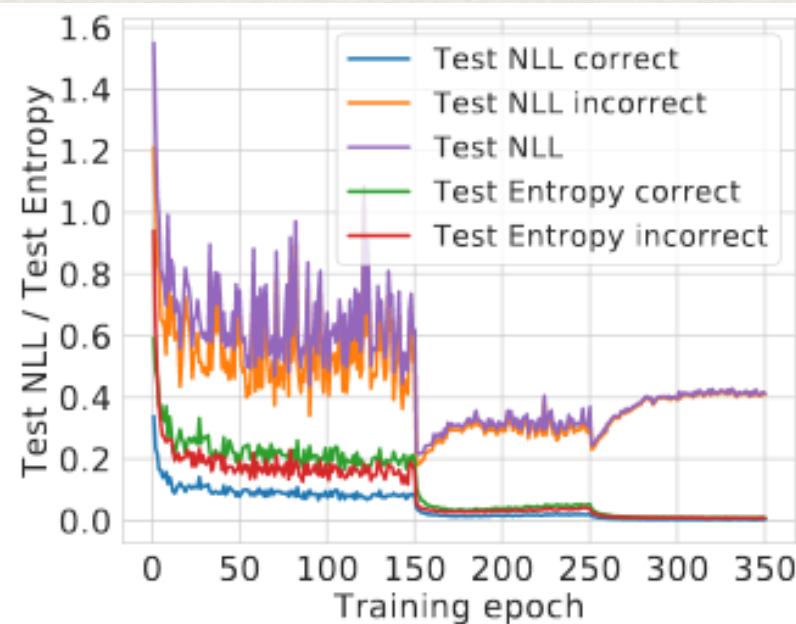
Expected Calibration Error

$$\begin{aligned} & \frac{1}{|\{B\}|} \sum_{B \in \{B\}} \frac{|B|}{m} \left| \frac{1}{|B|} \sum_{x \in B} b_k(x) - \frac{1}{|B|} \sum_{x \in B} I[y(x) = k] \right| = \\ & = \frac{1}{|\{B\}| m} \sum_{B \in \{B\}} \left| \sum_{x \in B} b_k(x) - \sum_{x \in B} I[y(x) = k] \right| \end{aligned}$$

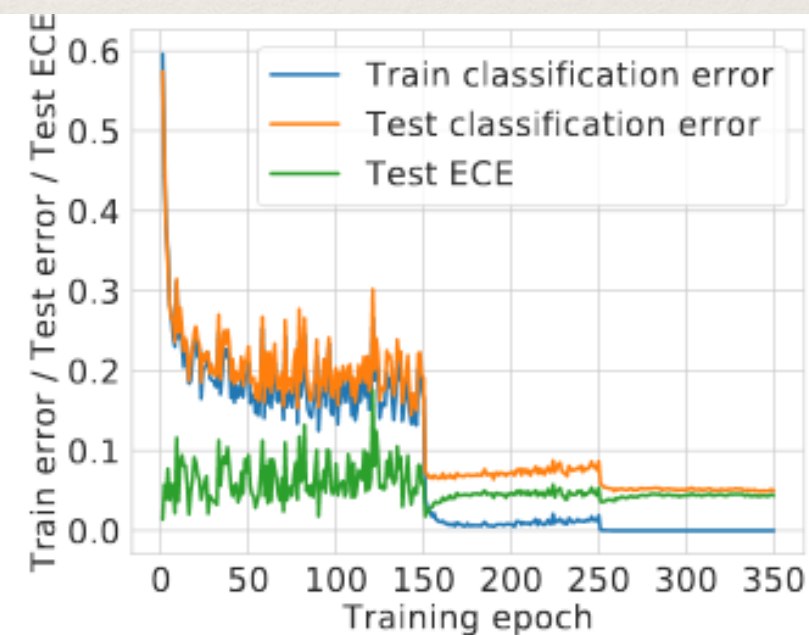
$$\text{NLL} = -\frac{1}{m} \sum_{j=1}^m y_j \log b_j \quad \text{MSE} = -\frac{1}{m} \sum_{j=1}^m (y_j - b_j)^2$$



(a)

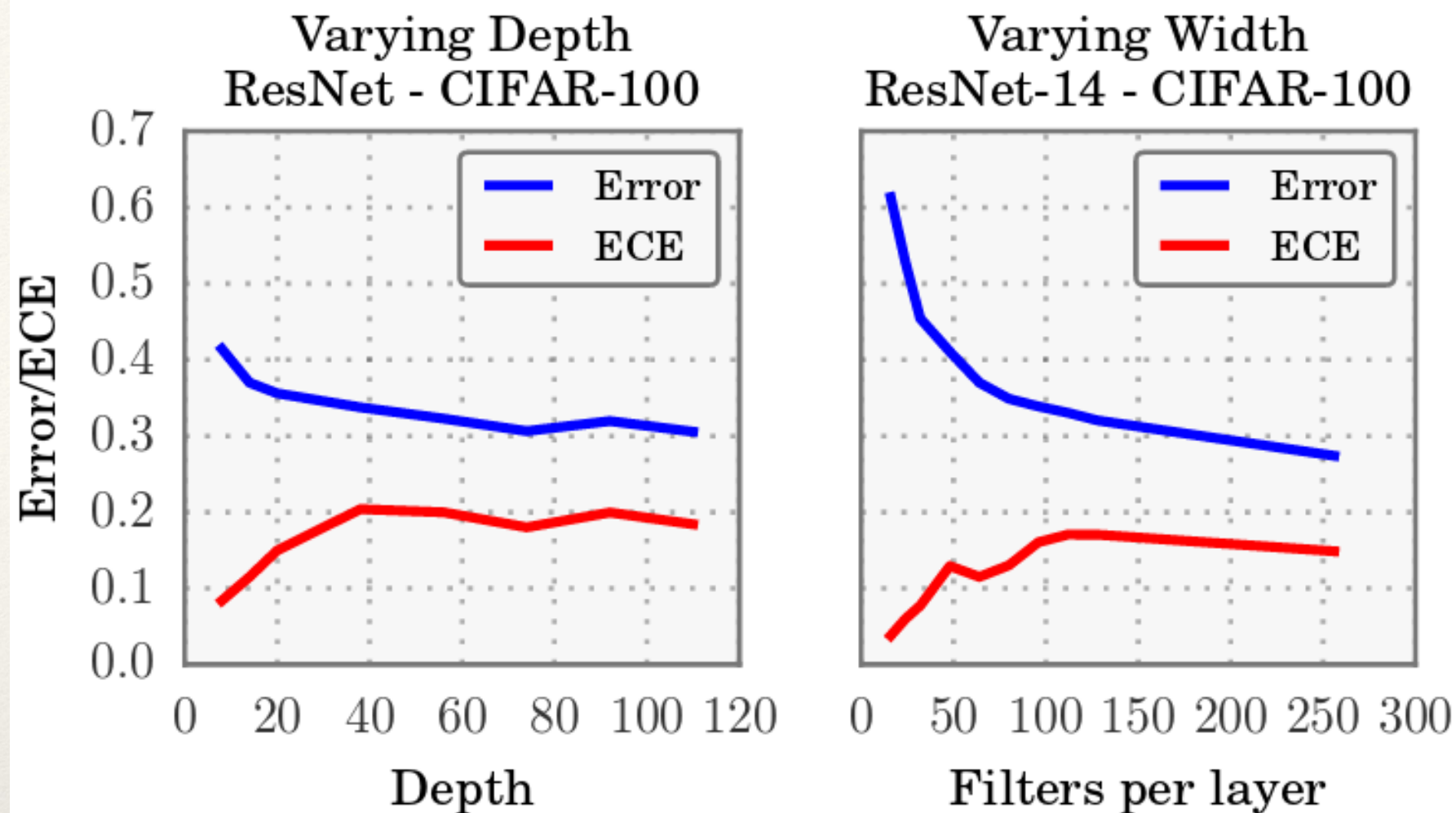


(b)

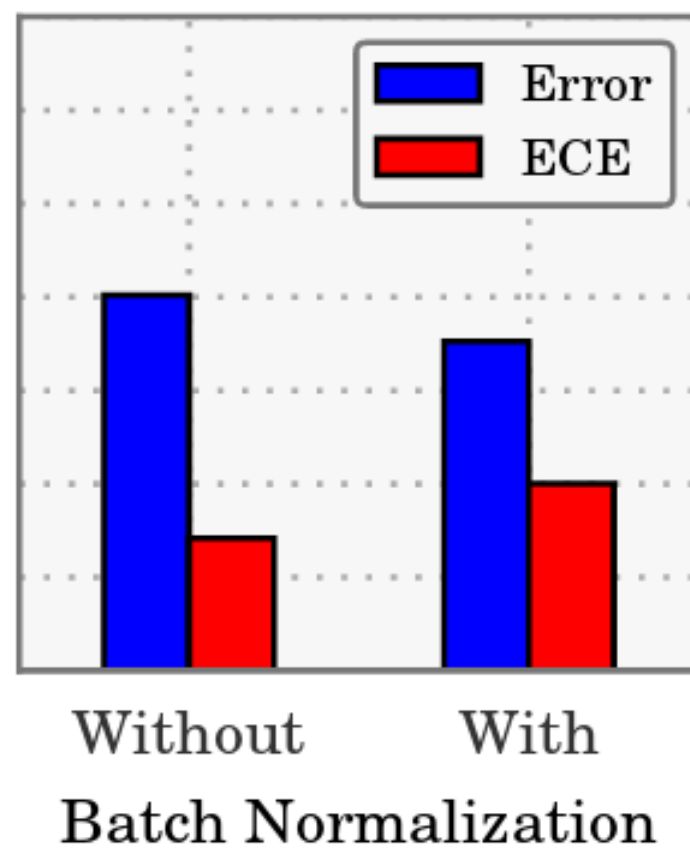


(c)

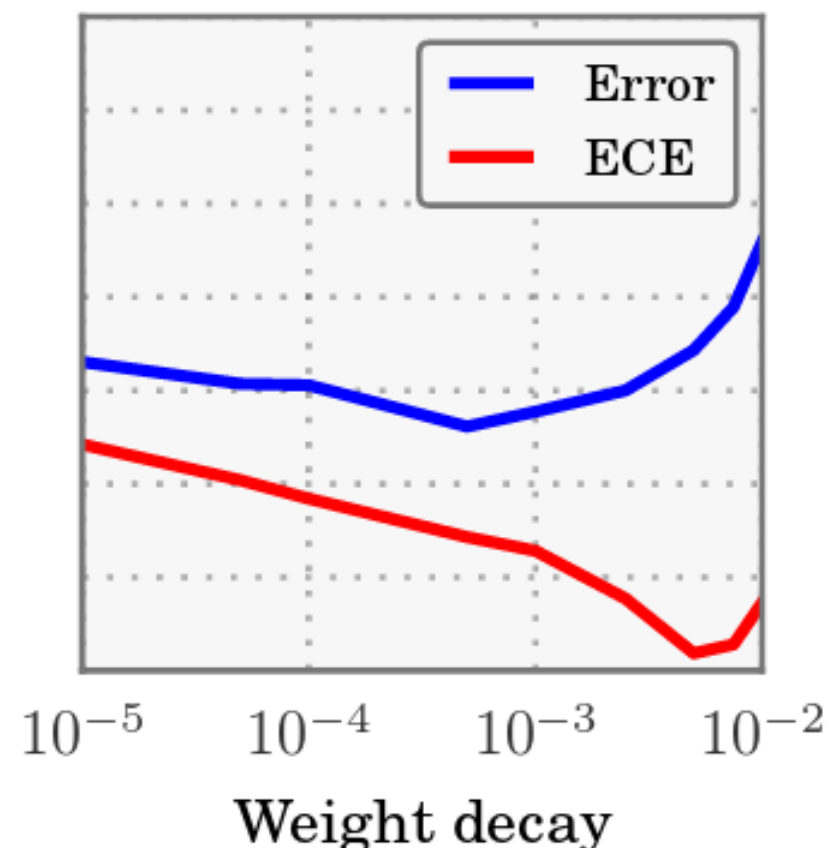
Figure 2: How metrics related to model calibration change whilst training a ResNet-50 network on CIFAR-10.



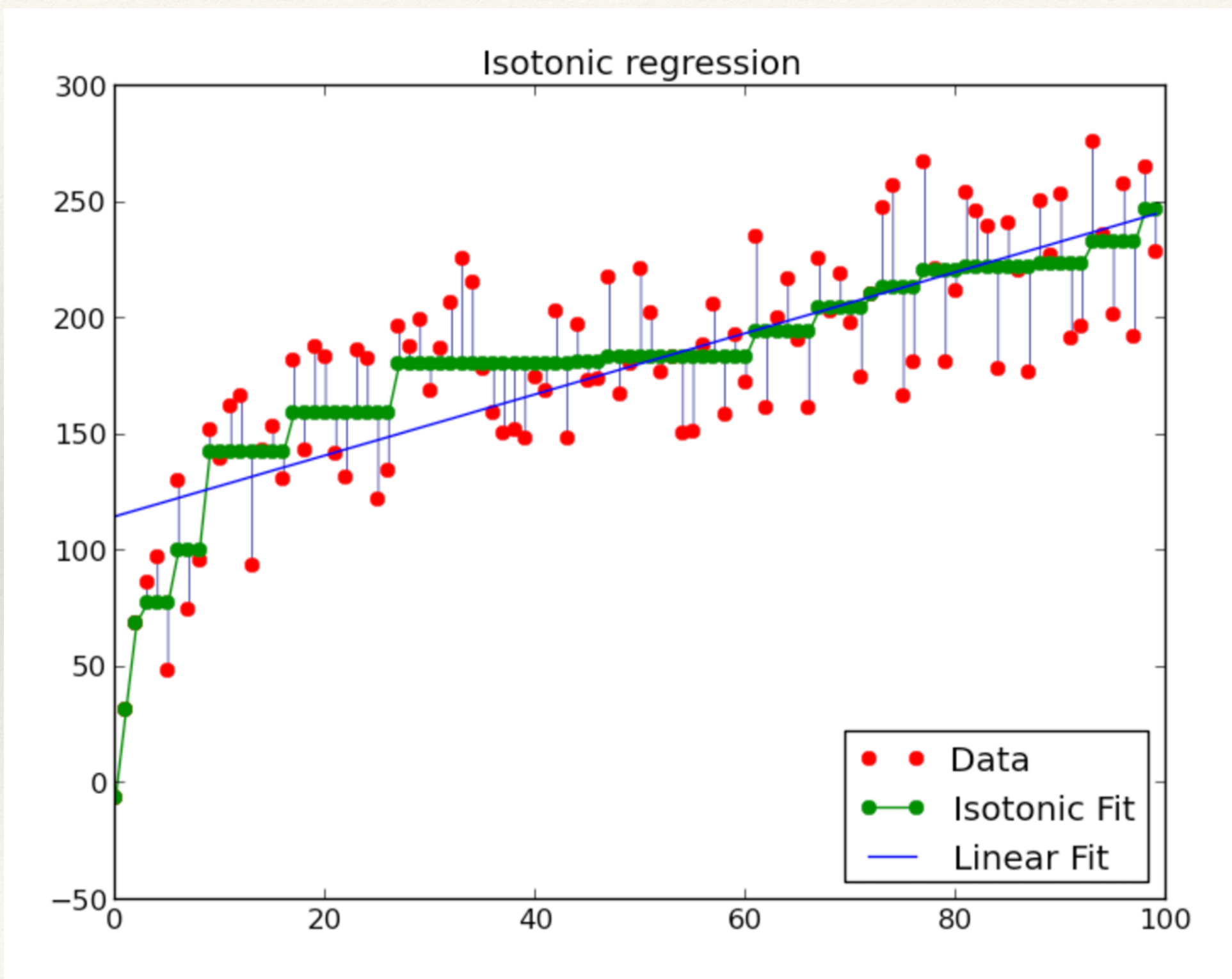
Using Normalization
ConvNet - CIFAR-100



Varying Weight Decay
ResNet-110 - CIFAR-100

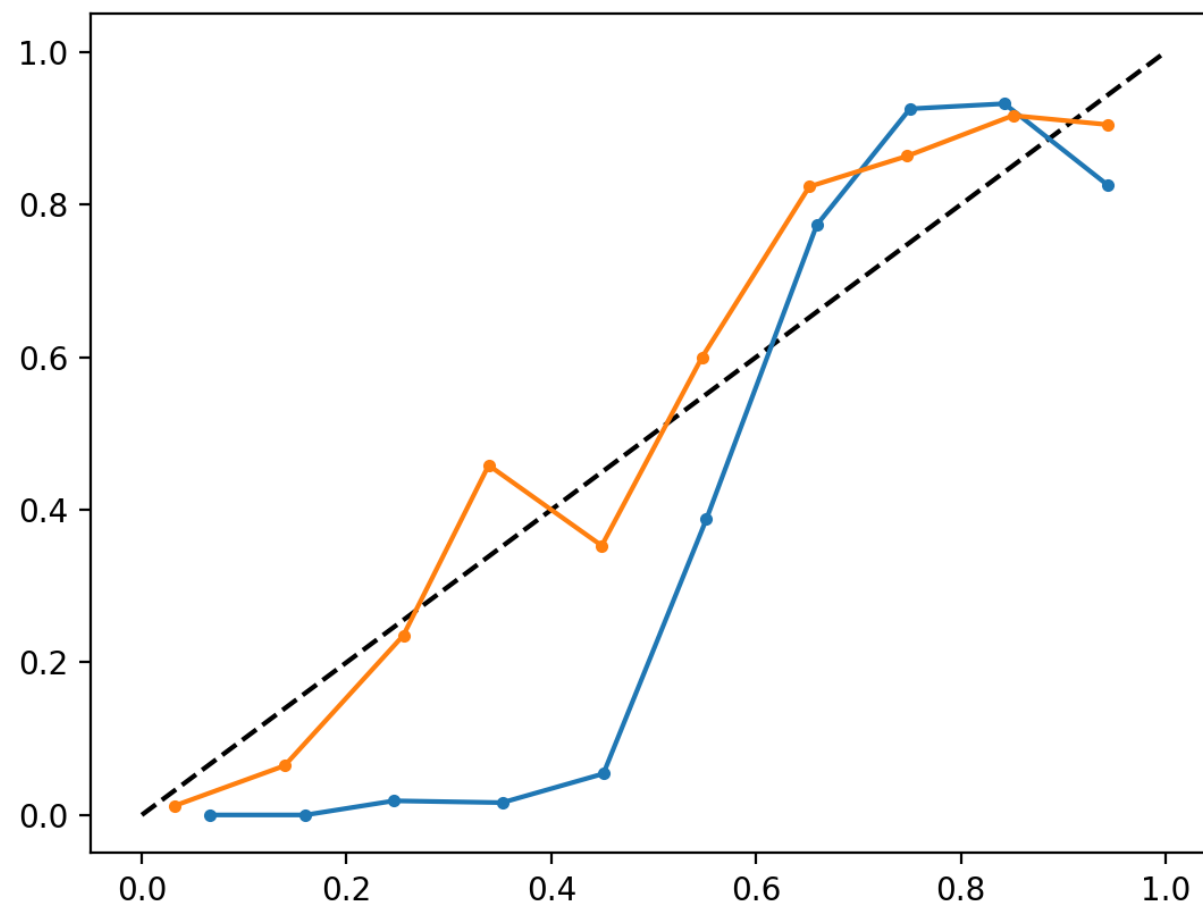


Изотоническая регрессия



Калибровка Платта

$$b_{\text{new}}(x) = \text{sigmoid}(\alpha \cdot r(x) + \beta)$$



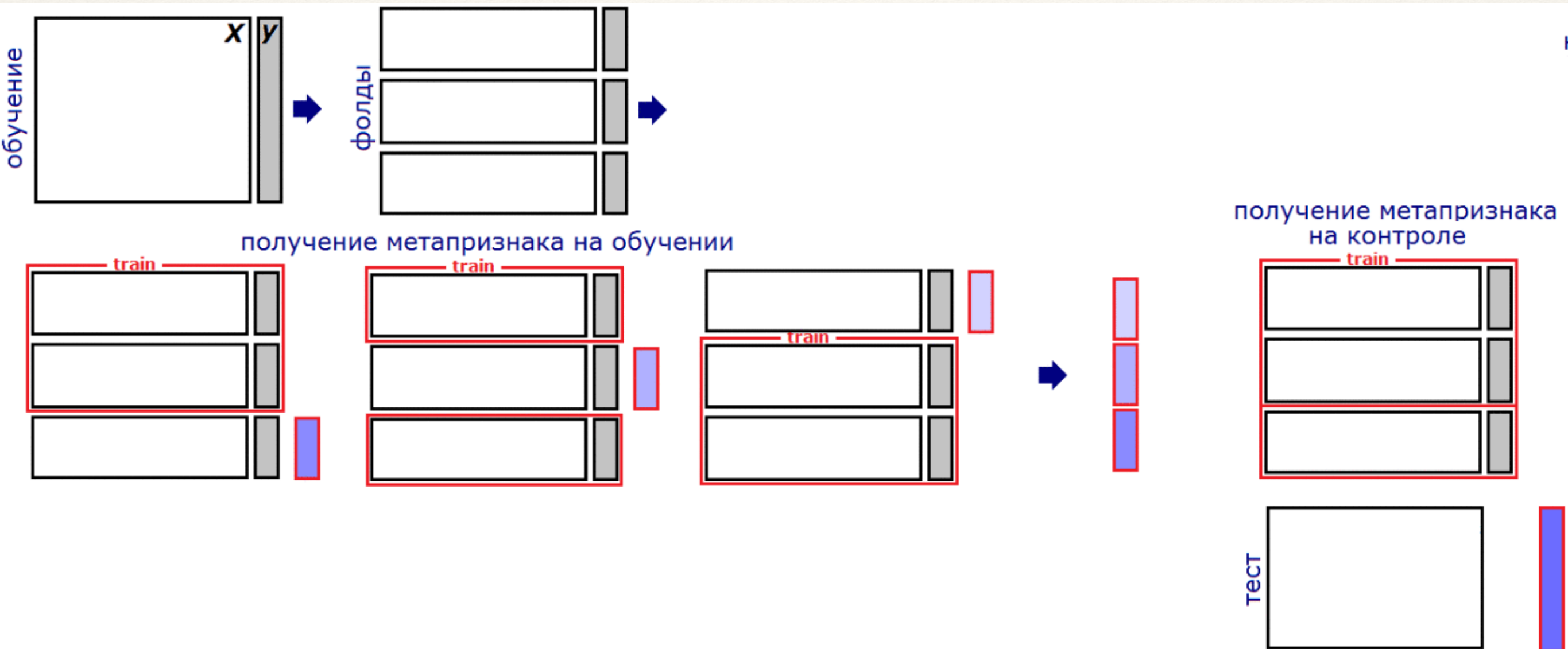
Температурное шкалирование

$$a(x) = \text{softmax}(b_1 / T, \dots, b_l / T)$$

Когда нужна калибровка ?

1. Для правильного понимания насколько можно доверять результатам модели(например человеческая обработка граничных случаев, краудсорсинг)
2. Для хорошего стакинга моделей(тк модели если модели по разному откалиброваны это может плохо сказаться на стакинге)
3. Когда на основе предиктов считается uplift или какая-то денежная выгода

Стакинг



Uplift-моделирование

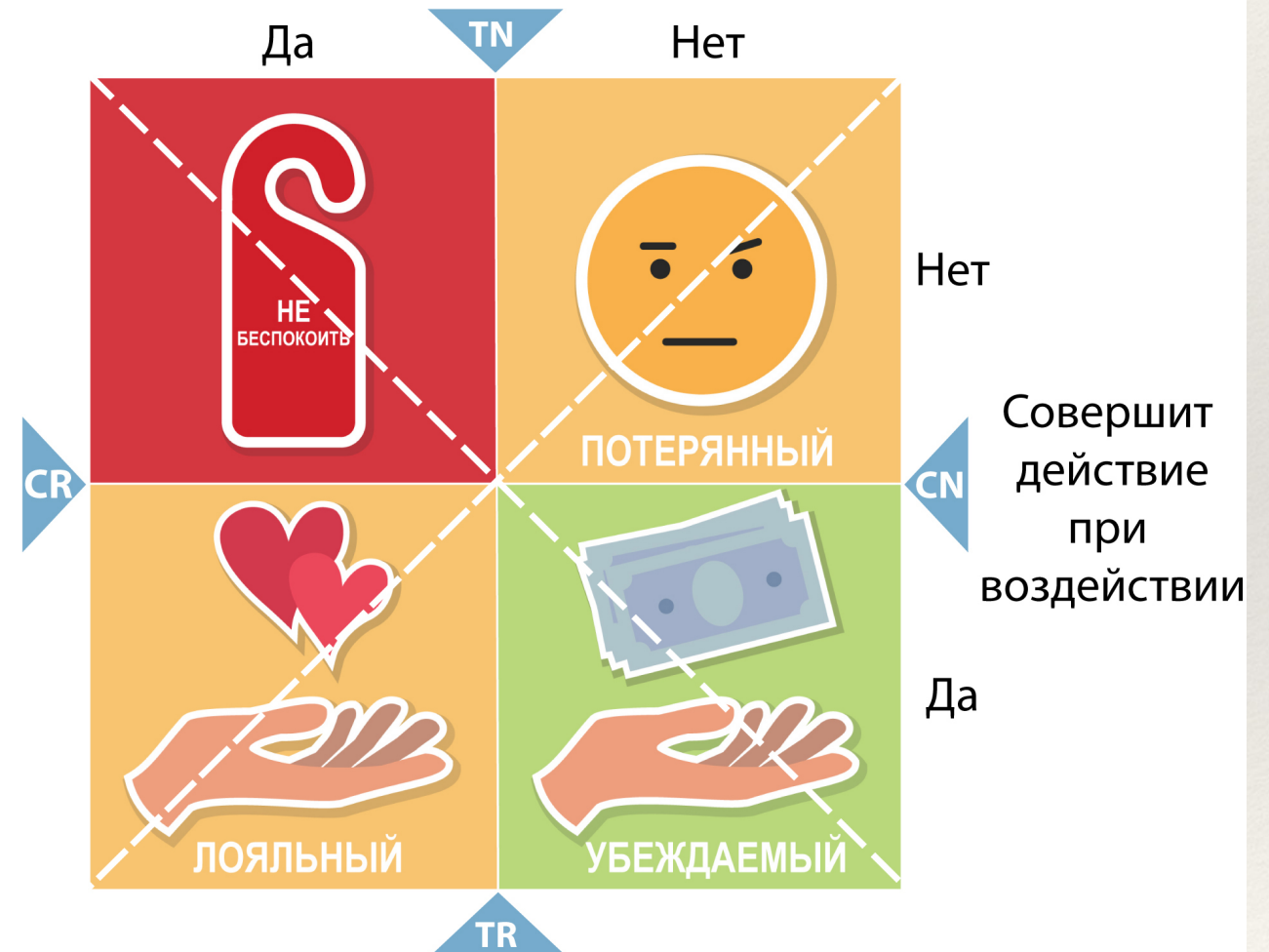
Что мы знаем:

Совершил действие

		Да	Нет
Было воздействие	Нет	<div>CR</div> $y=1$ $w=0$	<div>CN</div> $y=0$ $w=0$
	Да	<div>TR</div> $y=1$ $w=1$	<div>TN</div> $y=0$ $w=1$

Что мы хотим знать:

Совершит действие
без взаимодействия



Uplift-моделирование

Look-alike модель

$P(\text{целевого действия})$
на основе схожести

Response модель

$P(\text{целевого действия})$
при коммуникации

Uplift модель

$P(\text{целевого действия})$
при коммуникации

—

$P(\text{целевого действия})$
без коммуникации

Uplift * communication_cost vs marginality

Выводы

1. Калибровка нужна когда нам нужно интерпретировать результаты модели или использовать ее вероятности в бизнес метриках
2. Для нейросетей / SVM лучше всего использовать калибровку Платта / Температурное шкалирование