

# Model Selection

## Содержание

- Модель

$$\boldsymbol{\eta} = \mathbf{f}(\boldsymbol{\xi}) + \boldsymbol{\epsilon}, \quad \text{cov } \boldsymbol{\epsilon} = \sigma^2 \mathbf{I}, \quad \epsilon_i \perp \epsilon_j$$

- Функция потерь

$$L(f(\boldsymbol{\xi}), \eta) = |f(\boldsymbol{\xi}) - \eta|^d, \quad d \in \{1, 2\}.$$

- Пусть модель  $\hat{f}$  обучена на тренировочном множестве  $\mathcal{T}$ ; тогда хотелось бы знать её generalization (test, prediction) error

$$\text{Err}_{\mathcal{T}} = \mathbb{E} [L(\hat{f}(\boldsymbol{\xi}), \eta) \mid \mathcal{T}]$$

т.е. ожидание функции потерь на генеральной совокупности или

$$\text{Err} = \mathbb{E} [\text{Err}_{\mathcal{T}}] = \mathbb{E} [L(\hat{f}(\boldsymbol{\xi}), \eta)]$$

т.е. ожидание еще и по всем тренировочным выборкам. Тогда из всех возможных моделей могли бы выбрать  $\hat{f}$ , минимизирующую  $\text{Err}_{\mathcal{T}}$  (или  $\text{Err}$ ?)

- Справедливо для квадратичной  $L$ ,

$$\begin{aligned} \text{Err} &= \mathbb{E} [L(\hat{f}(\boldsymbol{\xi}), \eta)] = \mathbb{E} [(\hat{f}(\boldsymbol{\xi}) - \eta)^2] = \mathbb{E} \left[ \left( (\hat{f}(\boldsymbol{\xi}) - f(\boldsymbol{\xi})) + \epsilon \right)^2 \right] \\ &= \mathbb{E} \left[ (\hat{f}(\boldsymbol{\xi}) - f(\boldsymbol{\xi}))^2 + 2\epsilon (\hat{f}(\boldsymbol{\xi}) - f(\boldsymbol{\xi})) + \epsilon^2 \right] \\ &= \mathbb{E} \left[ (\hat{f}(\boldsymbol{\xi}) - f(\boldsymbol{\xi}))^2 \right] + \sigma^2. \end{aligned}$$

Видно, что  $\mathbb{E} \left[ (\hat{f}(\boldsymbol{\xi}) - f(\boldsymbol{\xi}))^2 \right]$  потенциально можно свести к 0, тогда как фиксированную  $\sigma^2$  — нет. Значит, нулевой ошибки для модели  $\hat{f}$  добиться физически невозможно (её минимальной величиной и будет  $\sigma^2$ ).

- Для  $\mathbf{x}_0 \notin \mathcal{T}$ , и квадратичной функции потерь справедливо<sup>1</sup> разложение

$$\text{Err}(\mathbf{x}_0) = \mathbb{E} [L(\hat{f}(\mathbf{x}_0), \eta) \mid \boldsymbol{\xi} = \mathbf{x}_0] = \sigma^2 + \underbrace{(\mathbb{E} \hat{f}(\mathbf{x}_0) - f(\mathbf{x}_0))^2}_{\text{bias}^2} + \underbrace{\mathbb{E} (\hat{f}(\mathbf{x}_0) - \mathbb{E} \hat{f}(\mathbf{x}_0))^2}_{D\hat{f}(\mathbf{x}_0)=\text{variance}}$$

**bias** показывает, насколько ожидание модели отличается от истинного значения,

**variance** показывает, насколько изменится оценка  $\hat{f}$ , если изменится тренировочная выборка  $\mathcal{T}$ .

---

<sup>1</sup><http://robjhyndman.com/hyndsight/files/2015/08/2-biasvardecomp.pdf>

С ростом гибкости модели уменьшается смещение, но растет разброс (модель начинает аппроксимировать шум  $\epsilon$ ) и наоборот. Стоит задача по тренировочной выборке подобрать такую гибкость модели, что  $\text{Err}_{\mathcal{T}}$  минимальна. Эмпирический способ это сделать — в противовес аналитическим методам вроде информационных критериев — для модели с заданной гибкости оценить  $\text{Err}_{\mathcal{T}}$ .

- Пусть  $\mathcal{S}$  — test set; величина

$$\widehat{\text{Err}}_{\mathcal{T}}(\mathcal{S}) = \frac{1}{|\mathcal{S}|} \sum_{(\mathbf{x}_i, y_i) \in \mathcal{S}} L(\hat{f}(\mathbf{x}_i), y_i) = \text{mean}_{\mathcal{S}} L$$

будет иметь играть роль оценки  $\text{Err}_{\mathcal{T}}$  (model assessment).

- Для нахождения оптимальной гибкости могли бы непосредственно посчитать training error

$$\overline{\text{err}}_{\mathcal{T}} = \frac{1}{N} \sum_{i=1}^N L(\hat{f}(\mathbf{x}_i), y_i) = \text{mean}_{\mathcal{T}} L.$$

Однако  $\overline{\text{err}}_{\mathcal{T}}$  будет убывать с увеличением гибкости модели, в то время, как  $\text{Err}_{\mathcal{T}}$  с какого-то момента начнет расти.

- Поэтому тренировочную выборку следует разделить на собственно тренировочную (по которой строится модель) и валидационную  $\mathcal{V}$ , по которой оценивать  $\text{Err}_{\mathcal{T}}$ . Изменением гиперпараметров модели (регулирующих гибкость), минимизировать

$$\widehat{\text{Err}}_{\mathcal{T}}(\mathcal{V}) = \text{mean}_{\mathcal{V}} L.$$

Так осуществляется задача model selection.

- Если выборка мала и мало  $\mathcal{V}$ , следует воспользоваться  $K$ -fold cross-validation: случайным образом разбить  $\mathcal{T}$  на  $\mathcal{T}_1, \dots, \mathcal{T}_K$ , затем обучать модель  $\hat{f}^{(-k)}$  на  $\mathcal{T} \setminus \mathcal{T}_k$ , а валидировать на  $\mathcal{T}_k$ ; результаты затем усреднить:

$$\frac{1}{K} \sum_{k=1}^K \left( \frac{1}{|\mathcal{T}_k|} \sum_{(\mathbf{x}_i, y_i) \in \mathcal{T}_k} \left| \hat{f}^{(-k)}(\mathbf{x}_i) - y_i \right|^d \right) = \text{mean}_{\{\mathcal{T}_1, \dots, \mathcal{T}_K\}} \left( \text{mean}_{\mathcal{T}_k} L_k \right).$$

- При  $K = N$  тренировочные выборки почти не отличаются друг от дружки, поэтому, по сравнению с другими вариантами, разброс будет большим как среднее скоррелированных величин, а смещение малым.
- При  $K \rightarrow 1$ , наоборот, разброс падает, смещение растёт.