## Bias - variance разложение

В статистике: $\hat{\theta}$ — оценка $\theta$

$$\text{MSE}_{\hat{\theta}}(\theta) := E_\theta(\hat{\theta} - \theta)^2 = \underbrace{(E_\theta \hat{\theta} - \theta)^2}_{\substack{\text{bias}^2 \\ \text{смещение}}} + \underbrace{D_\theta \hat{\theta}}_{\substack{\text{variance} \\ \text{разброс}}}$$

Обозначения: 1) $\hat{y}$ — случайная ф-ция $\mathcal{X} \to \mathcal{Y}$

Случайность может быть обусловлена:

- погрешностью при изм. отклика (случайные $y_i$)

- случайные признаки

- случайность модели $\left(\begin{array}{l}\text{например, выбор случайного}\\ \text{набора пр-ков в вершине}\end{array}\right)$

2) $x$ — неслучайный объект

$y$ — случайный отклик, т.е. $y = f(x, \varepsilon)$,

где $f$ — детерм ф-ция

$\varepsilon$ — случайный шум

$\hat{y}$, $\varepsilon$ независимы

$$MSE_{\hat{y}}(x) := E\left(\hat{y}(x) - y\right)^2 = E\left(\hat{y}(x) - f(x, \varepsilon)\right)^2 =$$

$$= E\left(\hat{y}(x) \pm Ef(x, \varepsilon) - f(x, \varepsilon)\right)^2 =$$

$$= E\left(\hat{y}(x) - Ef(x, \varepsilon)\right)^2 + E\left(f(x, \varepsilon) - Ef(x, \varepsilon)\right)^2 - 2\underbrace{E\overbrace{\left(\hat{y}(x) - Ef(x, \varepsilon)\right)}^{\text{зависит от } \hat{y}}\overbrace{\left(f(x, \varepsilon) - Ef(x, \varepsilon)\right)}^{\text{зависит от } \varepsilon}}_{} =$$

под «b-v разл. для оценок»  под «дисперсия $f(x, \varepsilon)$»  $E\left(f(x, \varepsilon) - Ef(x, \varepsilon)\right) = 0$

$$= \underbrace{\left(E\,\hat{y}(x) - Ef(x, \varepsilon)\right)^2}_{\text{bias}^2} + \underbrace{D\,\hat{y}(x)}_{\text{variance}} + \underbrace{D\,f(x, \varepsilon)}_{\text{noise}}$$

$$E\left(\hat{y}(x) - y\right)^2 = \left(E\hat{y}(x) - Ef(x, \varepsilon)\right)^2 + D\,\hat{y}(x) + D\,f(x, \varepsilon)$$

# Следствие 1

Пусть шум аддитивен и несмещ. $\quad f(x, \varepsilon) = h(x) + \varepsilon$

$$E\varepsilon = 0 \qquad D\varepsilon = \sigma^2$$

Тогда $\quad E\left(\hat{y}(x) - y\right)^2 = \left(E\hat{y}(x) - h(x)\right)^2 + D\hat{y}(x) + \sigma^2$

# Следствие 2

Пусть $X$ — новый объект, но теперь он случаен

Тогда $\quad E\left(\hat{y}(X) - y\right)^2 = E\left[\underbrace{E\left(\left(\hat{y}(X) - y\right)^2 \mid X\right)}_{\substack{X- \text{фикс.} \Rightarrow \\ \text{уже посчитано}}}\right] =$

$$D(x \mid y) = E(x^2 \mid y) - E(x \mid y)^2$$

УД

УМО

$$E(x \mid y) = f(y)$$

$$E(x \mid y = y) = f(\overset{\uparrow}{y})$$

$$EE(x \mid y) = Ex$$

$$E(\hat{y}(x) - y)^2 = \left(E(\hat{y}(x) \mid x) - E(f(x, \varepsilon) \mid x)\right)^2 + ED(\hat{y}(x) \mid x) + ED(f(x, \varepsilon) \mid x)$$

Применим разложение для аддитивного несмещ. шума.

$$\hat{y}(x) = x^T \hat{\theta} \qquad h(x) = x^T \theta \qquad y = X\theta + \varepsilon$$

$$E\varepsilon = 0$$
$$D\varepsilon = \sigma^2 I_n$$

$$\hat{\theta} = \left( X^T X + \lambda I_d \right)^{-1} X^T y$$

$$E\hat{y}(x) = x^T E\hat{\theta} = x^T \left( X^T X + \lambda I_d \right)^{-1} X^T X \theta$$

$$D\hat{y}(x) = x^T D\hat{\theta}\, x = x^T \left( X^T X + \lambda I_d \right)^{-1} X^T X \left( X^T X + \lambda I_d \right)^{-1} x\, \sigma^2$$

$$bias^2 = \left( E\hat{y}(x) - h(x) \right)^2 = \ldots$$
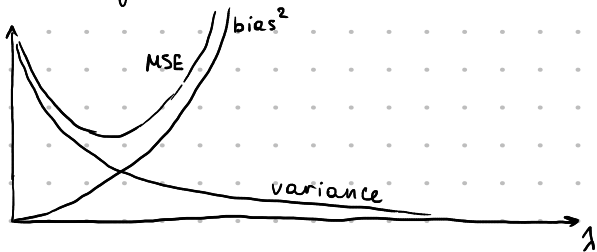
$$variance = D\hat{y}(x) = \ldots$$

# МНК $(\lambda = 0)$

$bias^2 = 0$

$variance = x^T (X^T X)^{-1} x \, \sigma^2$

**Вывод:** $X^T X$ близка к вырожденной $\Rightarrow$ в МНК разброс большой

$$Ridge - var < MHK - var$$

$$\hat{y}(x) = \frac{1}{T} \sum_{t=1}^{T} \hat{y}_t(x)$$

$$corr(x,y) = \frac{cov(x,y)}{\sqrt{Dx \, Dy}}$$

$\hat{y}_t$ ~~одинаково распределены~~   $cov(\hat{y}_a, \hat{y}_b) = cov(\hat{y}_c, y_d)$

$$bias^2 = \left( E\hat{y}(x) - h(x) \right)^2 = \left( E\hat{y}_1(x) - h(x) \right)^2$$

т.е. смещение композиции = смещ. одной модели

variance $\quad D\hat{y}(x) = \dfrac{1}{T^2} \sum_{t_1=1}^{T} \sum_{t_2=1}^{T} cov\left( \hat{y}_{t_1}(x), \hat{y}_{t_2}(x) \right) =$

$$= \frac{1}{T} D\hat{y}_1(x) + \frac{T-1}{T} cov\left( \hat{y}_1(x), \hat{y}_2(x) \right) =$$

$$D\hat{y}_1(x) \quad corr\left( \hat{y}_1(x), \hat{y}_2(x) \right)$$

$$= D\,\hat{y}_1(x)\left[\frac{1}{T} + \frac{T-1}{T}\,corr\left(\hat{y}_1(x),\hat{y}_2(x)\right)\right]$$

**Вывод:** Разброс тем меньше, чем меньше корреляция моделей.

# Bagging

Композиция вида $\quad \hat{y}(x) = \dfrac{1}{T} \sum\limits_{t=1}^{T} \hat{y}_t(x)$

где $\hat{y}_t$ построена по случ. подвыборке обучающей выборки с возвращениями.

Замечание: $\hat{y}_t$ могут быть из разных семейств (лин/дерево/...)

Частный случай — Random Forest

Bagging:

$\left.\begin{array}{ll} \text{bias} & \text{малое} \\ \text{varianse} & \text{одной модели большой} \\ \text{corr} & \text{моделей низкая} \end{array}\right\} \Rightarrow \begin{array}{l} \text{bias маленький} \\ \text{variance маленький} \end{array}$

# Как получить модели с малой корреляцией?

- разные типы моделей, разные гиперпараметры
- обучать на разных признаках
- разная предобработка данных