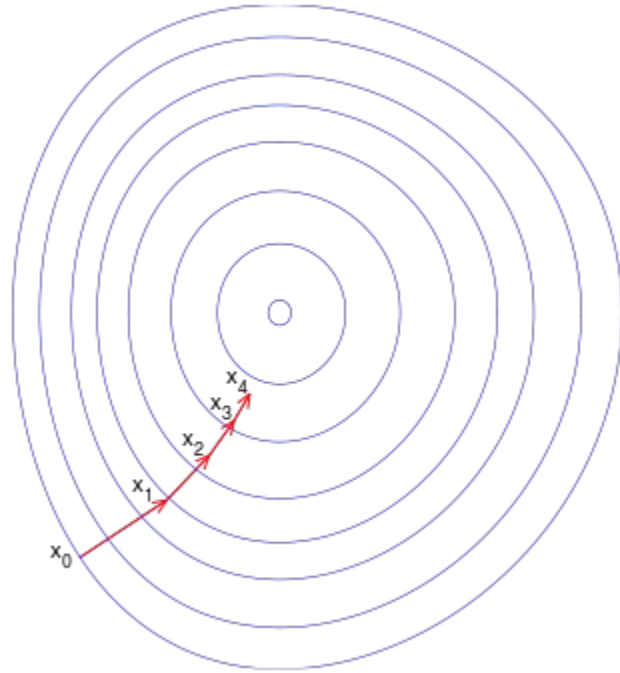# Lecture 7: Bias-variance decomposition
# & Feature importances

MIPT, 2019

# Outline

1. Gradient boosting recap

2. Bias-variance decomposition.

3. Feature importances estimation

We use gradient descent in *space of our models*

$$\hat{f}(x) = \sum_{i=0}^{t-1} \hat{f}_i(x),$$

$$r_{it} = -\left[\frac{\partial L(y_i, f(x_i))}{\partial f(x_i)}\right]_{f(x)=\hat{f}(x)}, \quad \text{for } i = 1, \dots, n,$$

$$\theta_t = \arg\min_{\theta} \sum_{i=1}^{n} (r_{it} - h(x_i, \theta))^2,$$

$$\rho_t = \arg\min_{\rho} \sum_{i=1}^{n} L(y_i, \hat{f}(x_i) + \rho \cdot h(x_i, \theta_t))$$

In linear regression case with MSE loss:

$$r_{it} = -\left[\frac{\partial L(y_i, f(x_i))}{\partial f(x_i)}\right]_{f(x)=\hat{f}(x)} = -2(\hat{y}_i - y_i) \propto \hat{y}_i - y_i$$

# Gradient boosting: recap

What we need:

- Data.
- Loss function and its gradient.
- Family of algorithms (with constraints on hyperparameters if necessary).
- Number of iterations M.
- Initial value (GBM by Friedman): constant.

The dataset $X = (x_i, y_i)_{i=1}^{\ell}$ with $y_i \in \mathbb{R}$ for regression problem.

# Bias-variance decomposition

The dataset $X = (x_i, y_i)_{i=1}^{\ell}$ with $y_i \in \mathbb{R}$ for regression problem.

Denote loss function $L(y, a) = \big(y - a(x)\big)^2$ .

# Bias-variance decomposition

The dataset $X = (x_i, y_i)_{i=1}^{\ell}$ with $y_i \in \mathbb{R}$ for regression problem.

Denote loss function $L(y, a) = (y - a(x))^2$ .

The corresponding risk estimation is

$$R(a) = \mathbb{E}_{x,y}\left[(y - a(x))^2\right] = \int_{\mathbb{X}} \int_{\mathbb{Y}} p(x, y)(y - a(x))^2 dx dy.$$

Let's show that $\quad a_*(x) = \mathbb{E}[y \mid x] = \displaystyle\int_{\mathbb{Y}} yp(y \mid x)dy$

# Bias-variance decomposition

Let's show that $\quad a_*(x) = \mathbb{E}[y \mid x] = \displaystyle\int_{\mathbb{Y}} y p(y \mid x) dy = \arg\min_a R(a).$

Let's show that $\quad a_*(x) = \mathbb{E}[y \mid x] = \int_{\mathbb{Y}} y p(y \mid x) dy = \arg\min_a R(a).$

$$L(y, a(x)) = (y - a(x))^2$$

# Bias-variance decomposition

Let's show that $\quad a_*(x) = \mathbb{E}[y \mid x] = \int_{\mathbb{Y}} y p(y \mid x) dy = \arg\min_a R(a).$

$$L(y, a(x)) = (y - a(x))^2 = (y - \mathbb{E}(y \mid x) + \mathbb{E}(y \mid x) - a(x))^2 =$$

Let's show that $a_*(x) = \mathbb{E}[y \mid x] = \int_{\mathbb{Y}} y p(y \mid x) dy = \arg\min_{a} R(a).$

$$L(y, a(x)) = (y - a(x))^2 = (y - \mathbb{E}(y \mid x) + \mathbb{E}(y \mid x) - a(x))^2 =$$
$$= (y - \mathbb{E}(y \mid x))^2 + 2(y - \mathbb{E}(y \mid x))(\mathbb{E}(y \mid x) - a(x)) + (\mathbb{E}(y \mid x) - a(x))^2.$$

Let's return to the risk estimation:

# Bias-variance decomposition

Let's show that $\quad a_*(x) = \mathbb{E}[y \mid x] = \int_{\mathbb{Y}} yp(y \mid x)dy = \arg\min_{a} R(a).$

$$L(y, a(x)) = (y - a(x))^2 = (y - \mathbb{E}(y \mid x) + \mathbb{E}(y \mid x) - a(x))^2 =$$
$$= (y - \mathbb{E}(y \mid x))^2 + 2(y - \mathbb{E}(y \mid x))(\mathbb{E}(y \mid x) - a(x)) + (\mathbb{E}(y \mid x) - a(x))^2.$$

Let's return to the risk estimation:

$$R(a) = \mathbb{E}_{x,y} L(y, a(x))$$

# Bias-variance decomposition

Let's show that
$$a_*(x) = \mathbb{E}[y \mid x] = \int_{\mathbb{Y}} y p(y \mid x) dy = \arg\min_a R(a).$$

$$L(y, a(x)) = (y - a(x))^2 = (y - \mathbb{E}(y \mid x) + \mathbb{E}(y \mid x) - a(x))^2 =$$
$$= (y - \mathbb{E}(y \mid x))^2 + 2(y - \mathbb{E}(y \mid x))(\mathbb{E}(y \mid x) - a(x)) + (\mathbb{E}(y \mid x) - a(x))^2.$$

Let's return to the risk estimation:

$$R(a) = \mathbb{E}_{x,y} L(y, a(x)) =$$
$$= \mathbb{E}_{x,y}(y - \mathbb{E}(y \mid x))^2 + \mathbb{E}_{x,y}(\mathbb{E}(y \mid x) - a(x))^2 +$$
$$+ 2\mathbb{E}_{x,y}(y - \mathbb{E}(y \mid x))(\mathbb{E}(y \mid x) - a(x)).$$

$$R(a) = \mathbb{E}_{x,y}L(y, a(x)) =$$
$$= \mathbb{E}_{x,y}(y - \mathbb{E}(y \mid x))^2 + \mathbb{E}_{x,y}(\mathbb{E}(y \mid x) - a(x))^2 +$$
$$+ 2\mathbb{E}_{x,y}\left(y - \mathbb{E}(y \mid x)\right)\left(\mathbb{E}(y \mid x) - a(x)\right).$$

Focus on the last term:

$$R(a) = \mathbb{E}_{x,y} L(y, a(x)) =$$
$$= \mathbb{E}_{x,y}(y - \mathbb{E}(y \mid x))^2 + \mathbb{E}_{x,y}(\mathbb{E}(y \mid x) - a(x))^2 +$$

Focus on the last term:

$$+ 2\mathbb{E}_{x,y}\left(y - \mathbb{E}(y \mid x)\right)\left(\mathbb{E}(y \mid x) - a(x)\right).$$

$$\mathbb{E}_x \mathbb{E}_y \left[\left(y - \mathbb{E}(y \mid x)\right)\left(\mathbb{E}(y \mid x) - a(x)\right) \mid x\right] =$$

$$R(a) = \mathbb{E}_{x,y} L(y, a(x)) =$$
$$= \mathbb{E}_{x,y}(y - \mathbb{E}(y \mid x))^2 + \mathbb{E}_{x,y}(\mathbb{E}(y \mid x) - a(x))^2 +$$
$$+ 2\mathbb{E}_{x,y}\big(y - \mathbb{E}(y \mid x)\big)\big(\mathbb{E}(y \mid x) - a(x)\big).$$

Focus on the last term:

Does not depend on y

$$\mathbb{E}_x \mathbb{E}_y \left[ \big(y - \mathbb{E}(y \mid x)\big)\big(\mathbb{E}(y \mid x) - a(x)\big) \mid x \right] =$$

$$R(a) = \mathbb{E}_{x,y} L(y, a(x)) =$$
$$= \mathbb{E}_{x,y} (y - \mathbb{E}(y \mid x))^2 + \mathbb{E}_{x,y} (\mathbb{E}(y \mid x) - a(x))^2 +$$
$$+ 2\mathbb{E}_{x,y} (y - \mathbb{E}(y \mid x)) (\mathbb{E}(y \mid x) - a(x)).$$

Focus on the last term:

Does not depend on y

$$\mathbb{E}_x \mathbb{E}_y \left[ (y - \mathbb{E}(y \mid x)) \boxed{(\mathbb{E}(y \mid x) - a(x))} \mid x \right] =$$
$$= \mathbb{E}_x \left( (\mathbb{E}(y \mid x) - a(x)) \mathbb{E}_y \left[ (y - \mathbb{E}(y \mid x)) \mid x \right] \right) =$$

$$R(a) = \mathbb{E}_{x,y} L(y, a(x)) =$$
$$= \mathbb{E}_{x,y}(y - \mathbb{E}(y \mid x))^2 + \mathbb{E}_{x,y}(\mathbb{E}(y \mid x) - a(x))^2 +$$
$$+ 2\mathbb{E}_{x,y}\big(y - \mathbb{E}(y \mid x)\big)\big(\mathbb{E}(y \mid x) - a(x)\big).$$

Focus on the last term:

$$\mathbb{E}_x \mathbb{E}_y \left[ \big(y - \mathbb{E}(y \mid x)\big)\big(\mathbb{E}(y \mid x) - a(x)\big) \mid x \right] =$$
$$= \mathbb{E}_x \left( \big(\mathbb{E}(y \mid x) - a(x)\big) \mathbb{E}_y \left[ \big(y - \mathbb{E}(y \mid x)\big) \mid x \right] \right) =$$
$$= \mathbb{E}_x \left( \big(\mathbb{E}(y \mid x) - a(x)\big)\big(\mathbb{E}(y \mid x) - \mathbb{E}(y \mid x)\big) \right) =$$

$$R(a) = \mathbb{E}_{x,y} L(y, a(x)) =$$
$$= \mathbb{E}_{x,y}(y - \mathbb{E}(y \mid x))^2 + \mathbb{E}_{x,y}(\mathbb{E}(y \mid x) - a(x))^2 +$$
$$+ 2\mathbb{E}_{x,y}(y - \mathbb{E}(y \mid x))(\mathbb{E}(y \mid x) - a(x)).$$

Focus on the last term:

$$\mathbb{E}_x \mathbb{E}_y \left[ (y - \mathbb{E}(y \mid x))(\mathbb{E}(y \mid x) - a(x)) \mid x \right] =$$
$$= \mathbb{E}_x \left( (\mathbb{E}(y \mid x) - a(x)) \mathbb{E}_y \left[ (y - \mathbb{E}(y \mid x)) \mid x \right] \right) =$$
$$= \mathbb{E}_x \left( (\mathbb{E}(y \mid x) - a(x))(\mathbb{E}(y \mid x) - \mathbb{E}(y \mid x)) \right) =$$
$$= 0$$

$$R(a) = \mathbb{E}_{x,y} L(y, a(x)) =$$
$$= \mathbb{E}_{x,y}(y - \mathbb{E}(y \mid x))^2 + \mathbb{E}_{x,y}(\mathbb{E}(y \mid x) - a(x))^2+$$
$$+ 2\mathbb{E}_{x,y}(y - \mathbb{E}(y \mid x))(\mathbb{E}(y \mid x) - a(x)).$$

Focus on the last term:

0

$$\mathbb{E}_x \mathbb{E}_y \Big[ (y - \mathbb{E}(y \mid x))(\mathbb{E}(y \mid x) - a(x)) \mid x \Big] =$$
$$= \mathbb{E}_x \Big( (\mathbb{E}(y \mid x) - a(x)) \mathbb{E}_y \Big[ (y - \mathbb{E}(y \mid x)) \mid x \Big] \Big) =$$
$$= \mathbb{E}_x \Big( (\mathbb{E}(y \mid x) - a(x))(\mathbb{E}(y \mid x) - \mathbb{E}(y \mid x)) \Big) =$$
$$= 0$$

So the risk takes form:

Does not depend on a(x)

$$R(a) = \boxed{\mathbb{E}_{x,y}(y - \mathbb{E}(y\,|\,x))^2} + \mathbb{E}_{x,y}(\mathbb{E}(y\,|\,x) - a(x))^2.$$

The minimum is reached when $a(x) = \mathbb{E}(y\,|\,x).$

So the optimal regression model with square loss is

$$a_*(x) = \mathbb{E}(y\,|\,x) = \int_{\mathbb{Y}} y\, p(y\,|\,x)dy.$$

Denote $\mu : (\mathbb{X} \times \mathbb{Y})^\ell \to \mathcal{A}$, where $\mathcal{A}$ is some family of algorithms.

Denote $\mu : (\mathbb{X} \times \mathbb{Y})^\ell \to \mathcal{A}$ , where $\mathcal{A}$ is some family of algorithms.

So $L(\mu) = \mathbb{E}_X \left[ \mathbb{E}_{x,y} \left[ (y - \mu(X)(x))^2 \right] \right]$ , where X dataset.
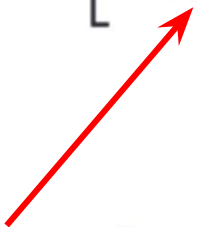
Denote $\mu : (\mathbb{X} \times \mathbb{Y})^{\ell} \to \mathcal{A}$ , where $\mathcal{A}$ is some family of algorithms.

So $L(\mu) = \mathbb{E}_X \left[ \mathbb{E}_{x,y} \left[ \left( y - \mu(X)(x) \right)^2 \right] \right]$ , where X dataset.

If X is fixed, then

$$\mathbb{E}_{x,y} \left[ \left( y - \mu(X) \right)^2 \right] = \mathbb{E}_{x,y} \left[ \left( y - \mathbb{E}[y \mid x] \right)^2 \right] + \mathbb{E}_{x,y} \left[ \left( \mathbb{E}[y \mid x] - \mu(X) \right)^2 \right].$$

Denote $\mu : (\mathbb{X} \times \mathbb{Y})^\ell \to \mathcal{A}$, where $\mathcal{A}$ is some family of algorithms.

So $L(\mu) = \mathbb{E}_X \left[ \mathbb{E}_{x,y} \left[ \left( y - \mu(X)(x) \right)^2 \right] \right]$, where X dataset.

If X is fixed, then

$$\mathbb{E}_{x,y} \left[ \left( y - \mu(X) \right)^2 \right] = \mathbb{E}_{x,y} \left[ \left( y - \mathbb{E}[y \mid x] \right)^2 \right] + \mathbb{E}_{x,y} \left[ \left( \mathbb{E}[y \mid x] - \mu(X) \right)^2 \right].$$

Let's combine the latter equations:

Denote $\mu : (\mathbb{X} \times \mathbb{Y})^{\ell} \to \mathcal{A}$ , where $\mathcal{A}$ is some family of algorithms.

So $L(\mu) = \mathbb{E}_X \left[ \mathbb{E}_{x,y} \left[ (y - \mu(X)(x))^2 \right] \right]$ , where X dataset.

If X is fixed, then

$$\mathbb{E}_{x,y} \left[ (y - \mu(X))^2 \right] = \mathbb{E}_{x,y} \left[ (y - \mathbb{E}[y \mid x])^2 \right] + \mathbb{E}_{x,y} \left[ (\mathbb{E}[y \mid x] - \mu(X))^2 \right].$$

Let's combine the latter equations:

$$L(\mu) = \mathbb{E}_X \left[ \underbrace{\mathbb{E}_{x,y} \left[ (y - \mathbb{E}[y \mid x])^2 \right]}_{\text{Does not depend on X}} + \mathbb{E}_{x,y} \left[ (\mathbb{E}[y \mid x] - \mu(X))^2 \right] \right]$$

$$L(\mu) = \mathbb{E}_X \left[ \underbrace{\mathbb{E}_{x,y} \left[ \left( y - \mathbb{E}[y \mid x] \right)^2 \right]}_{\text{Does not depend on X}} + \mathbb{E}_{x,y} \left[ \left( \mathbb{E}[y \mid x] - \mu(X) \right)^2 \right] \right] =$$

Does not depend on X

$$L(\mu) = \mathbb{E}_X \left[ \underbrace{\mathbb{E}_{x,y} \left[ \left( y - \mathbb{E}[y \mid x] \right)^2 \right]}_{\text{Does not depend on X}} + \mathbb{E}_{x,y} \left[ \left( \mathbb{E}[y \mid x] - \mu(X) \right)^2 \right] \right] =$$

$$= \mathbb{E}_{x,y} \left[ \left( y - \mathbb{E}[y \mid x] \right)^2 \right] + \mathbb{E}_{x,y} \left[ \mathbb{E}_X \left[ \left( \mathbb{E}[y \mid x] - \mu(X) \right)^2 \right] \right].$$

$$L(\mu) = \mathbb{E}_X\left[\mathbb{E}_{x,y}\left[\left(y - \mathbb{E}[y \mid x]\right)^2\right] + \mathbb{E}_{x,y}\left[\left(\mathbb{E}[y \mid x] - \mu(X)\right)^2\right]\right] =$$

$$= \mathbb{E}_{x,y}\left[\left(y - \mathbb{E}[y \mid x]\right)^2\right] + \boxed{\mathbb{E}_{x,y}\left[\mathbb{E}_X\left[\left(\mathbb{E}[y \mid x] - \mu(X)\right)^2\right]\right]}.$$

Focus on the second term:

$$L(\mu) = \mathbb{E}_X\left[\mathbb{E}_{x,y}\left[(y - \mathbb{E}[y \mid x])^2\right] + \mathbb{E}_{x,y}\left[(\mathbb{E}[y \mid x] - \mu(X))^2\right]\right] =$$

$$= \mathbb{E}_{x,y}\left[(y - \mathbb{E}[y \mid x])^2\right] + \boxed{\mathbb{E}_{x,y}\left[\mathbb{E}_X\left[(\mathbb{E}[y \mid x] - \mu(X))^2\right]\right]}.$$

Focus on the second term:

$$\mathbb{E}_{x,y}\left[\mathbb{E}_X\left[(\mathbb{E}[y \mid x] - \mu(X))^2\right]\right] =$$

$$L(\mu) = \mathbb{E}_X \left[ \mathbb{E}_{x,y} \left[ (y - \mathbb{E}[y \mid x])^2 \right] + \mathbb{E}_{x,y} \left[ (\mathbb{E}[y \mid x] - \mu(X))^2 \right] \right] =$$

$$= \mathbb{E}_{x,y} \left[ (y - \mathbb{E}[y \mid x])^2 \right] + \mathbb{E}_{x,y} \left[ \mathbb{E}_X \left[ (\mathbb{E}[y \mid x] - \mu(X))^2 \right] \right].$$

Focus on the second term:

$$\mathbb{E}_{x,y} \left[ \mathbb{E}_X \left[ (\mathbb{E}[y \mid x] - \mu(X))^2 \right] \right] =$$

$$= \mathbb{E}_{x,y} \left[ \mathbb{E}_X \left[ (\mathbb{E}[y \mid x] - \mathbb{E}_X[\mu(X)] + \mathbb{E}_X[\mu(X)] - \mu(X))^2 \right] \right]$$

$$\mathbb{E}_{x,y}\left[\mathbb{E}_X\left[\left(\mathbb{E}[y\mid x]-\mu(X)\right)^2\right]\right]=$$

$$=\mathbb{E}_{x,y}\left[\mathbb{E}_X\left[\left(\mathbb{E}[y\mid x]-\mathbb{E}_X\left[\mu(X)\right]+\mathbb{E}_X\left[\mu(X)\right]-\mu(X)\right)^2\right]\right]=$$

$$\mathbb{E}_{x,y}\left[\mathbb{E}_X\left[\left(\mathbb{E}[y \mid x] - \mu(X)\right)^2\right]\right] =$$

$$= \mathbb{E}_{x,y}\left[\mathbb{E}_X\left[\left(\mathbb{E}[y \mid x] - \mathbb{E}_X\left[\mu(X)\right] + \mathbb{E}_X\left[\mu(X)\right] - \mu(X)\right)^2\right]\right] =$$

$$= \mathbb{E}_{x,y}\left[\mathbb{E}_X\left[\underbrace{\left(\mathbb{E}[y \mid x] - \mathbb{E}_X\left[\mu(X)\right]\right)^2}\right]\right] + \mathbb{E}_{x,y}\left[\mathbb{E}_X\left[\left(\mathbb{E}_X\left[\mu(X)\right] - \mu(X)\right)^2\right]\right] +$$

$$+ 2\mathbb{E}_{x,y}\left[\mathbb{E}_X\left[\left(\mathbb{E}[y \mid x] - \mathbb{E}_X\left[\mu(X)\right]\right)\left(\mathbb{E}_X\left[\mu(X)\right] - \mu(X)\right)\right]\right].$$

$$\mathbb{E}_{x,y}\left[\mathbb{E}_X\left[\left(\mathbb{E}[y\mid x]-\mu(X)\right)^2\right]\right]=$$

$$=\mathbb{E}_{x,y}\left[\mathbb{E}_X\left[\left(\mathbb{E}[y\mid x]-\mathbb{E}_X\left[\mu(X)\right]+\mathbb{E}_X\left[\mu(X)\right]-\mu(X)\right)^2\right]\right]=$$

$$=\mathbb{E}_{x,y}\left[\mathbb{E}_X\left[\underbrace{\left(\mathbb{E}[y\mid x]-\mathbb{E}_X\left[\mu(X)\right]\right)^2}\right]\right]+\mathbb{E}_{x,y}\left[\mathbb{E}_X\left[\left(\mathbb{E}_X\left[\mu(X)\right]-\mu(X)\right)^2\right]\right]+$$

<span style="color:red">Does not depend on X</span>

$$+\,2\mathbb{E}_{x,y}\left[\mathbb{E}_X\left[\left(\mathbb{E}[y\mid x]-\mathbb{E}_X\left[\mu(X)\right]\right)\left(\mathbb{E}_X\left[\mu(X)\right]-\mu(X)\right)\right]\right].$$

$$\mathbb{E}_{x,y}\left[\mathbb{E}_X\left[\left(\mathbb{E}[y \mid x] - \mu(X)\right)^2\right]\right] =$$

$$= \mathbb{E}_{x,y}\left[\mathbb{E}_X\left[\left(\mathbb{E}[y \mid x] - \mathbb{E}_X\left[\mu(X)\right] + \mathbb{E}_X\left[\mu(X)\right] - \mu(X)\right)^2\right]\right] =$$

$$= \mathbb{E}_{x,y}\left[\mathbb{E}_X\left[\underbrace{\left(\mathbb{E}[y \mid x] - \mathbb{E}_X\left[\mu(X)\right]\right)^2}\right]\right] + \mathbb{E}_{x,y}\left[\mathbb{E}_X\left[\left(\mathbb{E}_X\left[\mu(X)\right] - \mu(X)\right)^2\right]\right] +$$

Does not depend on X

$$+ 2\mathbb{E}_{x,y}\left[\mathbb{E}_X\left[\left(\mathbb{E}[y \mid x] - \mathbb{E}_X\left[\mu(X)\right]\right)\left(\mathbb{E}_X\left[\mu(X)\right] - \mu(X)\right)\right]\right].$$

Just a bit further, we are almost there

$$\mathbb{E}_{x,y}\left[\mathbb{E}_X\left[\left(\mathbb{E}[y\mid x]-\mu(X)\right)^2\right]\right]=$$

$$=\mathbb{E}_{x,y}\left[\mathbb{E}_X\left[\left(\mathbb{E}[y\mid x]-\mathbb{E}_X\left[\mu(X)\right]+\mathbb{E}_X\left[\mu(X)\right]-\mu(X)\right)^2\right]\right]=$$

$$=\mathbb{E}_{x,y}\left[\mathbb{E}_X\left[\left(\mathbb{E}[y\mid x]-\mathbb{E}_X\left[\mu(X)\right]\right)^2\right]\right]+\mathbb{E}_{x,y}\left[\mathbb{E}_X\left[\left(\mathbb{E}_X\left[\mu(X)\right]-\mu(X)\right)^2\right]\right]+$$

$$+2\mathbb{E}_{x,y}\left[\mathbb{E}_X\left[\left(\mathbb{E}[y\mid x]-\mathbb{E}_X\left[\mu(X)\right]\right)\left(\mathbb{E}_X\left[\mu(X)\right]-\mu(X)\right)\right]\right].$$

Focus on this term

$$\mathbb{E}_X \left[ \left( \mathbb{E}[y \mid x] - \mathbb{E}_X \left[ \mu(X) \right] \right) \left( \mathbb{E}_X \left[ \mu(X) \right] - \mu(X) \right) \right] =$$

$$\mathbb{E}_X \left[ \left( \mathbb{E}[y \mid x] - \mathbb{E}_X \left[ \mu(X) \right] \right) \left( \mathbb{E}_X \left[ \mu(X) \right] - \mu(X) \right) \right] =$$

$$= \left( \mathbb{E}[y \mid x] - \mathbb{E}_X \left[ \mu(X) \right] \right) \mathbb{E}_X \left[ \mathbb{E}_X \left[ \mu(X) \right] - \mu(X) \right] =$$

$$\mathbb{E}_X \left[ \left( \mathbb{E}[y \mid x] - \mathbb{E}_X \left[ \mu(X) \right] \right) \left( \mathbb{E}_X \left[ \mu(X) \right] - \mu(X) \right) \right] =$$

$$= \left( \mathbb{E}[y \mid x] - \mathbb{E}_X \left[ \mu(X) \right] \right) \mathbb{E}_X \left[ \mathbb{E}_X \left[ \mu(X) \right] - \mu(X) \right] =$$

$$= \left( \mathbb{E}[y \mid x] - \mathbb{E}_X \left[ \mu(X) \right] \right) \left[ \mathbb{E}_X \left[ \mu(X) \right] - \mathbb{E}_X \left[ \mu(X) \right] \right] =$$

$$\mathbb{E}_X \Big[ \big( \mathbb{E}[y \mid x] - \mathbb{E}_X \big[ \mu(X) \big] \big) \big( \mathbb{E}_X \big[ \mu(X) \big] - \mu(X) \big) \Big] =$$

$$= \big( \mathbb{E}[y \mid x] - \mathbb{E}_X \big[ \mu(X) \big] \big) \mathbb{E}_X \Big[ \mathbb{E}_X \big[ \mu(X) \big] - \mu(X) \Big] =$$

$$= \big( \mathbb{E}[y \mid x] - \mathbb{E}_X \big[ \mu(X) \big] \big) \Big[ \mathbb{E}_X \big[ \mu(X) \big] - \mathbb{E}_X \big[ \mu(X) \big] \Big] =$$

$$= 0.$$

$$\mathbb{E}_{x,y}\left[\mathbb{E}_X\left[\left(\mathbb{E}[y\mid x]-\mu(X)\right)^2\right]\right] =$$

$$= \mathbb{E}_{x,y}\left[\mathbb{E}_X\left[\left(\mathbb{E}[y\mid x]-\mathbb{E}_X\left[\mu(X)\right]+\mathbb{E}_X\left[\mu(X)\right]-\mu(X)\right)^2\right]\right] =$$

$$= \mathbb{E}_{x,y}\left[\mathbb{E}_X\left[\left(\mathbb{E}[y\mid x]-\mathbb{E}_X\left[\mu(X)\right]\right)^2\right]\right] + \mathbb{E}_{x,y}\left[\mathbb{E}_X\left[\left(\mathbb{E}_X\left[\mu(X)\right]-\mu(X)\right)^2\right]\right] +$$

$$+ 2\mathbb{E}_{x,y}\left[\mathbb{E}_X\left[\left(\mathbb{E}[y\mid x]-\mathbb{E}_X\left[\mu(X)\right]\right)\left(\mathbb{E}_X\left[\mu(X)\right]-\mu(X)\right)\right]\right].$$

0

$$L(\mu) = \underbrace{\mathbb{E}_{x,y}\left[\left(y - \mathbb{E}[y \mid x]\right)^2\right]}_{\text{noise}} +$$

$$+ \underbrace{\mathbb{E}_x\left[\left(\mathbb{E}_X\left[\mu(X)\right] - \mathbb{E}[y \mid x]\right)^2\right]}_{\text{bias}} + \underbrace{\mathbb{E}_x\left[\mathbb{E}_X\left[\left(\mu(X) - \mathbb{E}_X\left[\mu(X)\right]\right)^2\right]\right]}_{\text{variance}}.$$
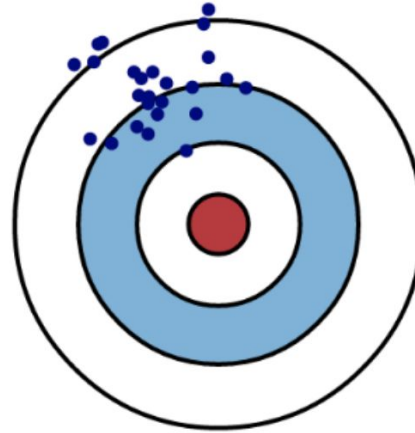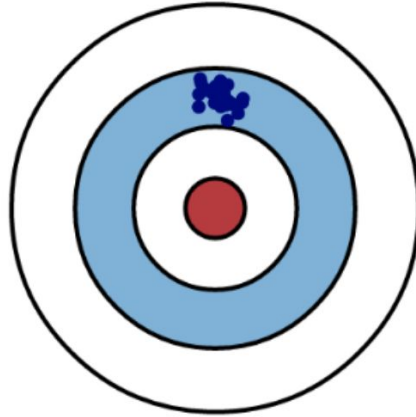
Low Variance    High Variance

Low Bias

High Bias

46

$$L(\mu) = \underbrace{\mathbb{E}_{x,y}\left[\left(y - \mathbb{E}[y \mid x]\right)^2\right]}_{\text{noise}} +$$

$$+ \underbrace{\mathbb{E}_x\left[\left(\mathbb{E}_X\left[\mu(X)\right] - \mathbb{E}[y \mid x]\right)^2\right]}_{\text{bias}} + \underbrace{\mathbb{E}_x\left[\mathbb{E}_X\left[\left(\mu(X) - \mathbb{E}_X\left[\mu(X)\right]\right)^2\right]\right]}_{\text{variance}}.$$

This exact form of bias-variance decomposition is correct for square loss in regression.

However, it is much more general. See extra materials for more exotic cases.

# Bagging = Bootstrap aggregating

Denote dataset $\tilde{X}$ bootstrapped from $X$.

Denote $\mu$: $\tilde{\mu}(X) = \mu(\tilde{X})$. Let $b_n(x)$ be basic algorithm.

Denote the ensemble:

$$a_N(x) = \frac{1}{N} \sum_{n=1}^{N} b_n(x) = \frac{1}{N} \sum_{n=1}^{N} \tilde{\mu}(X)(x).$$

$$L(\mu) = \underbrace{\mathbb{E}_{x,y}\left[\left(y - \mathbb{E}[y \mid x]\right)^2\right]}_{\text{noise}} +$$

$$+ \underbrace{\mathbb{E}_x\left[\left(\mathbb{E}_X\left[\mu(X)\right] - \mathbb{E}[y \mid x]\right)^2\right]}_{\text{bias}} + \underbrace{\mathbb{E}_x\left[\mathbb{E}_X\left[\left(\mu(X) - \mathbb{E}_X\left[\mu(X)\right]\right)^2\right]\right]}_{\text{variance}}.$$

The bias term takes the following form:

$$\mathbb{E}_{x,y}\left[\left(\mathbb{E}_X\left[\frac{1}{N}\sum_{n=1}^{N}\tilde{\mu}(X)(x)\right] - \mathbb{E}[y \mid x]\right)^2\right] =$$

The bias term takes the following form:

$$\mathbb{E}_{x,y}\left[\left(\mathbb{E}_X\left[\frac{1}{N}\sum_{n=1}^{N}\tilde{\mu}(X)(x)\right]-\mathbb{E}[y\,|\,x]\right)^2\right]=$$

$$=\mathbb{E}_{x,y}\left[\left(\frac{1}{N}\sum_{n=1}^{N}\mathbb{E}_X[\tilde{\mu}(X)(x)]-\mathbb{E}[y\,|\,x]\right)^2\right]=$$

The bias term takes the following form:

$$\mathbb{E}_{x,y}\left[\left(\mathbb{E}_X\left[\frac{1}{N}\sum_{n=1}^{N}\tilde{\mu}(X)(x)\right] - \mathbb{E}[y\,|\,x]\right)^2\right] =$$

$$= \mathbb{E}_{x,y}\left[\left(\frac{1}{N}\sum_{n=1}^{N}\mathbb{E}_X[\tilde{\mu}(X)(x)] - \mathbb{E}[y\,|\,x]\right)^2\right] =$$

$$= \mathbb{E}_{x,y}\left[\left(\mathbb{E}_X[\tilde{\mu}(X)(x)] - \mathbb{E}[y\,|\,x]\right)^2\right].$$

The bias term takes the following form:

$$\mathbb{E}_{x,y}\left[\left(\mathbb{E}_X\left[\frac{1}{N}\sum_{n=1}^{N}\tilde{\mu}(X)(x)\right] - \mathbb{E}[y\,|\,x]\right)^2\right] =$$

$$= \mathbb{E}_{x,y}\left[\left(\frac{1}{N}\sum_{n=1}^{N}\mathbb{E}_X[\tilde{\mu}(X)(x)] - \mathbb{E}[y\,|\,x]\right)^2\right] =$$

$$= \mathbb{E}_{x,y}\left[\left(\mathbb{E}_X[\tilde{\mu}(X)(x)] - \mathbb{E}[y\,|\,x]\right)^2\right].$$

One algorithm bias

The variance:
$$\mathbb{E}_{x,y}\left[\mathbb{E}_X\left[\left(\frac{1}{N}\sum_{n=1}^{N}\tilde{\mu}(X)(x)-\mathbb{E}_X\left[\frac{1}{N}\sum_{n=1}^{N}\tilde{\mu}(X)(x)\right]\right)^2\right]\right].$$

$$\left(\frac{1}{N}\sum_{n=1}^{N}\tilde{\mu}(X)(x)-\mathbb{E}_X\left[\frac{1}{N}\sum_{n=1}^{N}\tilde{\mu}(X)(x)\right]\right)^2 =$$

$$= \frac{1}{N^2}\left(\sum_{n=1}^{N}\left[\tilde{\mu}(X)(x)-\mathbb{E}_X\left[\tilde{\mu}(X)(x)\right]\right]\right)^2 =$$

$$= \frac{1}{N^2}\sum_{n=1}^{N}\left(\tilde{\mu}(X)(x)-\mathbb{E}_X\left[\tilde{\mu}(X)(x)\right]\right)^2 +$$

$$+ \frac{1}{N^2}\sum_{n_1\neq n_2}\left(\tilde{\mu}(X)(x)-\mathbb{E}_X\left[\tilde{\mu}(X)(x)\right]\right)\left(\tilde{\mu}(X)(x)-\mathbb{E}_X\left[\tilde{\mu}(X)(x)\right]\right)$$
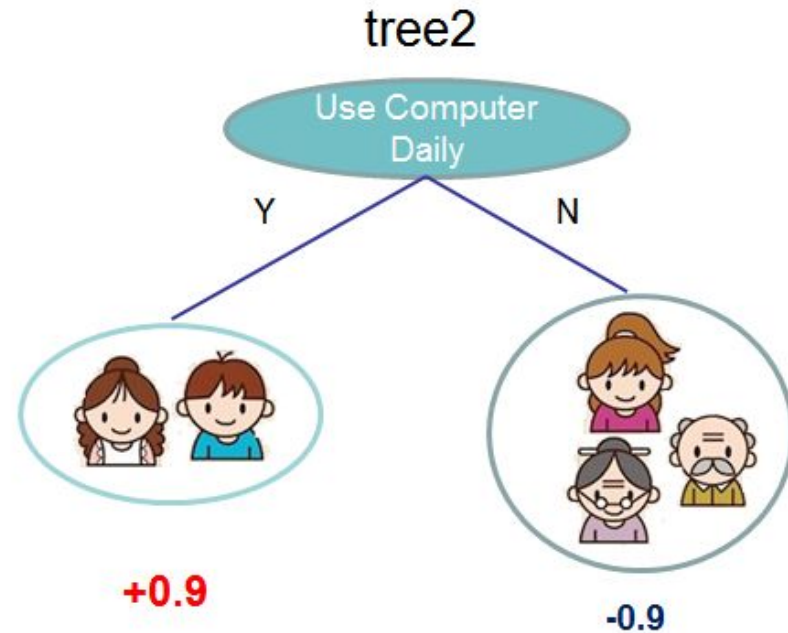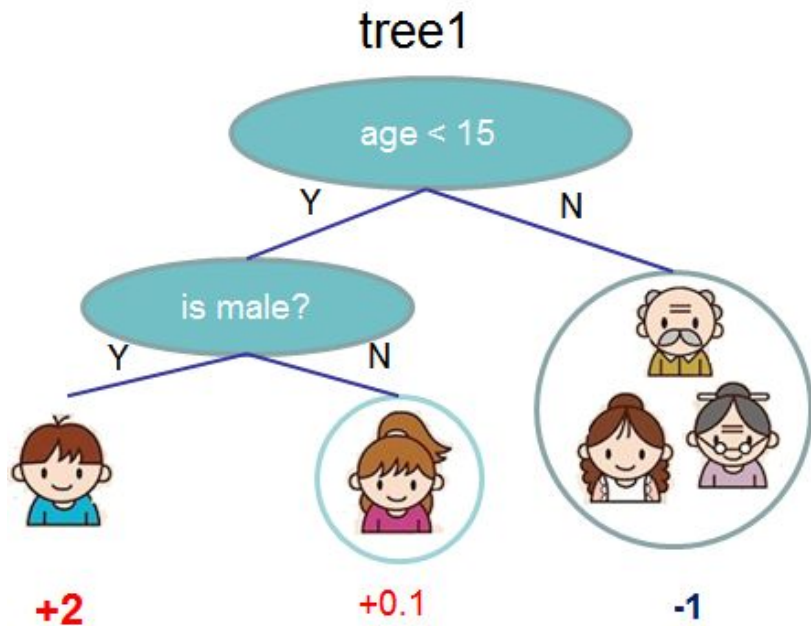
The variance:

$$
\mathbb{E}_{x,y}\Big[\mathbb{E}_X\Big[\frac{1}{N^2}\sum_{n=1}^{N}\big(\tilde{\mu}(X)(x)-\mathbb{E}_X[\tilde{\mu}(X)(x)]\big)^2+
$$

$$
+\frac{1}{N^2}\sum_{n_1\neq n_2}\big(\tilde{\mu}(X)(x)-\mathbb{E}_X[\tilde{\mu}(X)(x)]\big)\big(\tilde{\mu}(X)(x)-\mathbb{E}_X[\tilde{\mu}(X)(x)]\big)\Big]\Big]=
$$

$$
=\frac{1}{N^2}\mathbb{E}_{x,y}\Big[\mathbb{E}_X\Big[\sum_{n=1}^{N}\big(\tilde{\mu}(X)(x)-\mathbb{E}_X[\tilde{\mu}(X)(x)]\big)^2\Big]\Big]+
$$

$$
+\frac{1}{N^2}\mathbb{E}_{x,y}\Big[\mathbb{E}_X\Big[\sum_{n_1\neq n_2}\big(\tilde{\mu}(X)(x)-\mathbb{E}_X[\tilde{\mu}(X)(x)]\big)\times
$$

$$
\times\big(\tilde{\mu}(X)(x)-\mathbb{E}_X[\tilde{\mu}(X)(x)]\big)\Big]\Big]=
$$

One algorithm variance * 1/N

$$
=\boxed{\frac{1}{N}\mathbb{E}_{x,y}\Big[\mathbb{E}_X\Big[\big(\tilde{\mu}(X)(x)-\mathbb{E}_X[\tilde{\mu}(X)(x)]\big)^2\Big]\Big]}+
$$

$$
+\frac{N(N-1)}{N^2}\mathbb{E}_{x,y}\Big[\mathbb{E}_X\Big[\big(\tilde{\mu}(X)(x)-\mathbb{E}_X[\tilde{\mu}(X)(x)]\big)\times
$$

$$
\times\big(\tilde{\mu}(X)(x)-\mathbb{E}_X[\tilde{\mu}(X)(x)]\big)\Big]\Big]
$$

The variance:

$$\mathbb{E}_{x,y}\Big[\mathbb{E}_X\Big[\frac{1}{N^2}\sum_{n=1}^{N}\Big(\tilde{\mu}(X)(x)-\mathbb{E}_X\big[\tilde{\mu}(X)(x)\big]\Big)^2+$$

$$+\frac{1}{N^2}\sum_{n_1\neq n_2}\Big(\tilde{\mu}(X)(x)-\mathbb{E}_X\big[\tilde{\mu}(X)(x)\big]\Big)\Big(\tilde{\mu}(X)(x)-\mathbb{E}_X\big[\tilde{\mu}(X)(x)\big]\Big)\Big]\Big]=$$

$$=\frac{1}{N^2}\mathbb{E}_{x,y}\Big[\mathbb{E}_X\Big[\sum_{n=1}^{N}\Big(\tilde{\mu}(X)(x)-\mathbb{E}_X\big[\tilde{\mu}(X)(x)\big]\Big)^2\Big]\Big]+$$

$$+\frac{1}{N^2}\mathbb{E}_{x,y}\Big[\mathbb{E}_X\Big[\sum_{n_1\neq n_2}\Big(\tilde{\mu}(X)(x)-\mathbb{E}_X\big[\tilde{\mu}(X)(x)\big]\Big)\times$$

$$\times\Big(\tilde{\mu}(X)(x)-\mathbb{E}_X\big[\tilde{\mu}(X)(x)\big]\Big)\Big]\Big]=$$ One algorithm variance * 1/N

$$=\boxed{\frac{1}{N}\mathbb{E}_{x,y}\Big[\mathbb{E}_X\Big[\Big(\tilde{\mu}(X)(x)-\mathbb{E}_X\big[\tilde{\mu}(X)(x)\big]\Big)^2\Big]\Big]}+$$

$$+\frac{N(N-1)}{N^2}\mathbb{E}_{x,y}\Big[\mathbb{E}_X\Big[\Big(\tilde{\mu}(X)(x)-\mathbb{E}_X\big[\tilde{\mu}(X)(x)\big]\Big)\times$$

$$\times\Big(\tilde{\mu}(X)(x)-\mathbb{E}_X\big[\tilde{\mu}(X)(x)\big]\Big)\Big]\Big]$$

Basic algorithms covariance

56

# Feature importance estimation

Image source: https://habr.com/ru/post/428213/
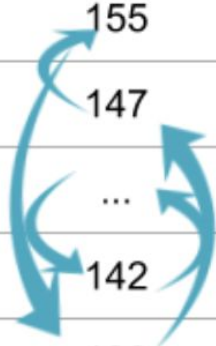
# Feature importance estimation

1. Permutation importance
2. Partial Dependence Plots (PDP)
3. Tree specific:
    a. Gain
    b. Frequency (Split Count)
    c. Cover (weighted Split Count)
4. Shap

# Permutation importance

| Height at age 20 (cm) | Height at age 10 (cm) | ... | Socks owned at age 10 |
|---|---|---|---|
| 182 | 155 | ... | 20 |
| 175 | 147 | ... | 10 |
| ... | ... | ... | ... |
| 156 | 142 | ... | 8 |
| 153 | 130 | ... | 24 |

Image source: https://www.kaggle.com/dansbecker/permutation-importance

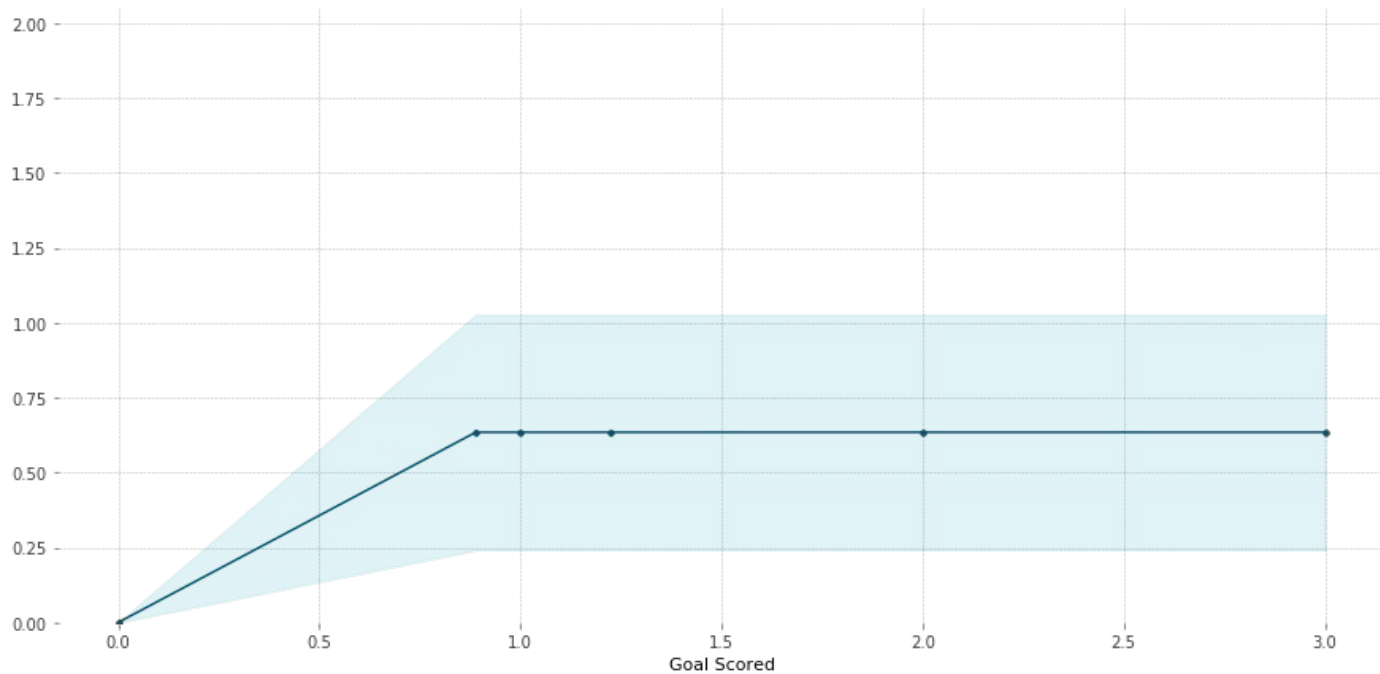| Height at age 20 (cm) | Height at age 10 (cm) | ... | Socks owned at age 10 |
|---|---|---|---|
| 182 | 155 | ... | 20 |
| 175 | 147 | ... | 10 |
| ... | ... | ... | ... |
| 156 | 142 | ... | 8 |
| 153 | 130 | ... | 24 |

Train model

Observe changes caused by feature random permutations

Image source: https://www.kaggle.com/dansbecker/permutation-importance

# Partial Dependence Plots



PDP for feature "Goal Scored"
Number of unique grid points: 6

Image source: https://www.kaggle.com/dansbecker/partial-plots

# Importance estimation problems



Image source: "Consistent Individualized Feature Attribution for Tree Ensembles" paper

Consider *i-th* feature. Shap value will be

$$\phi_i(p) = \sum_{S \subseteq N/\{i\}} \frac{|S|!(n - |S| - 1)!}{n!} (p(S \cup \{i\}) - p(S))$$

where $p(S \cup \{i\})$ is model prediction on feature subset S with *i-th* feature added.

Consider *i-th* feature. Shap value will be

$$\phi_i(p) = \sum_{S \subseteq N/\{i\}} \frac{|S|!(n - |S| - 1)!}{n!} (p(S \cup \{i\}) - p(S))$$

where $p(S \cup \{i\})$ is model prediction on feature subset S with *i-th* feature added.

*SHAP values are the only consistent and locally accurate*

*individualized feature attributions*

See [Consistent Individualized Feature Attribution for Tree Ensembles](#) paper for more info

1. Remember the bias-varience decomposition
2. Consider using SHAP values to estimate feature importances.

Great demo: http://arogozhnikov.github.io/2016/06/24/gradient_boosting_explained.html

Based on ml-mipt, HSE 2018 by Evgeny Sokolov, mlcourse_open by ODS, kaggle posts and blog posts referred on slides