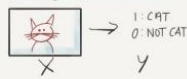


NEURAL NETWORKS: 1. DEEP LEARNING 2. CONVOLUTIONAL 3. RECURRENT 4. TRANSFER LEARNING 5. HYBRID 6. AUTOENCODERS 7. GAN 8. VAE 9. RL 10. ...

## BINARY CLASSIFICATION

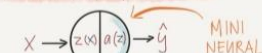


## LOGISTIC REGRESSION AS A NEURAL NET

### FINDING THE MINIMUM WITH GRADIENT DESCENT

- FIND THE DIRECTION (USING DERIVATIVES)
- WALK (UPDATE  $w$  &  $b$ ) AT A LEARNING RATE
- REPEAT (UNTIL YOU REACH BOTTOM (CONVERGE))

### PUTTING IT ALL TOGETHER



$$Z(x) = w \cdot x + b$$
$$\hat{y} = a(z) = \sigma(\text{SIGMOID}(z))$$

- FORWARD PROPAGATION
- CALCULATE  $\hat{y}$
- BACKWARD PROPAGATION
- GRADIENT DESCENT
- UPDATE  $w$  &  $b$
- REPEAT UNTIL IT CONVERGES



FIND THE MINIMUM

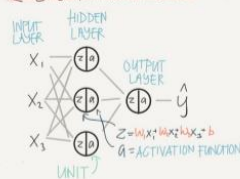
THE TASK IS TO LEARN  $w$  &  $b$  BUT HOW?

1. OPTIMIZE HOW GOOD THE GUESS IS BY MINIMIZING THE DIFF. BETWEEN GUESS ( $\hat{y}$ ) AND TRUTH ( $y$ )

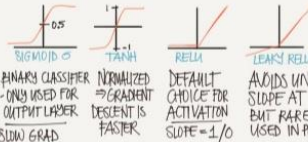
$$\text{LOSS} = \mathcal{L}(\hat{y}, y)$$
$$\text{COST} = J(w, b) = \frac{1}{n} \sum_{i=1}^n \mathcal{L}(\hat{y}^{(i)}, y^{(i)})$$

COST = LOSS FOR THE ENTIRE DATASET

## 2 LAYER NEURAL NET



### ACTIVATION FUNCTIONS



**SIGMOID**  $\sigma$  **TANH** **RELU** **LEAKY RELU**

**BINARY CLASSIFIER** - ONLY USED FOR OUTPUT LAYER

**NORMALIZED**  $\Rightarrow$  GRADIENT DESCENT IS FASTER

**DEFAULT CHOICE FOR ACTIVATION** SLOPE = 1/0

**AVOIDS UNDEF SLOPE AT 0** BUT RARELY USED IN PRACTICE

## SHALLOW NEURAL NETS

### WHY ACTIVATION FUNCTIONS?

EX. WITH NO ACTIVATION -  $a = z$

$$a^{(1)} = z^{(1)} = w^{(1)}x + b^{(1)} \quad \text{LAYER 1}$$
$$a^{(2)} = z^{(2)} = w^{(2)}a^{(1)} + b^{(2)} \quad \text{LAYER 2}$$

PLUG IN  $a^{(1)}$

$$a^{(2)} = w^{(2)}(w^{(1)}x + b^{(1)}) + b^{(2)}$$
$$= w^{(2)}w^{(1)}x + w^{(2)}b^{(1)} + b^{(2)}$$
$$w^{(2)}w^{(1)}x + b^{(2)}$$

**LINEAR FUNCTION**

### INITIALIZING $w$ & $b$

WHAT IF  $\rightarrow$  INIT TO 0

THIS WILL MAKE ALL THE UNITS TO BE THE SAME AND LEARN EXACTLY THE SAME FEATURES

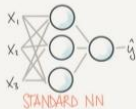
**SOLUTION:** RANDOM INIT BUT ALSO WANT THEM SMALL SO RAND \* 0.01

**HYPERPARAM**

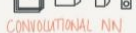
## INTRO TO DEEP LEARNING

### SUPERVISED LEARNING

INPUT: X	OUTPUT: y	NN TYPE
HOME FEATURES	PRICE	STANDARD NN
AD+USER INFO	WILL CLICK ON AD (y/n)	STANDARD NN
IMAGE	OBJECT (1...1000)	CONV. NN (CNN)
AUDIO	TEXT TRANSCRIPT CHINESE	RECURRENT NN (RNN)
ENGLISH	POS OF OTHER CARS	CUSTOM/HYBRID
IMAGE/RADAR		



STANDARD NN



CONVOLUTIONAL NN



RECURRENT NN

### NETWORK ARCHITECTURES

NNs CAN DEAL WITH BOTH STRUCTURED & UNSTRUCTURED DATA



STRUCTURED UNSTRUCTURED

### WHY NOW?



ONE OF THE BIG BREAKTHROUGHS HAS BEEN MOVING FROM SIGMOID TO RELU FOR FASTER GRADIENT DESCENT

**CHANGED** **RELU**

**FASTER (COMPUTATION IS IMPORTANT TO SPEED UP THE ITERATIVE PROCESS)**

**IDEA** **CODE** **EXPERIM**

# Two Forest Jump

## Комбинированный отбор признаков с использованием двухлесового метода

Афанасьев Сергей  
КБ «Ренессанс Кредит»

8 сентября 2020 г.  
Москва

# Для чего нужен отбор признаков?

1

Избежать проклятья размерности

2

Снизить переобучение модели

3

Упростить модель для интерпретации

4

Сократить время обучения модели

# Фильтры

(Filter method)



- Матрица корреляций
- Таргет-корреляция
- PCA
- ...

# Обертки

(Wrapper method)



- Stepwise Regression
- Random Forest
- RFE
- ...

# Вложения

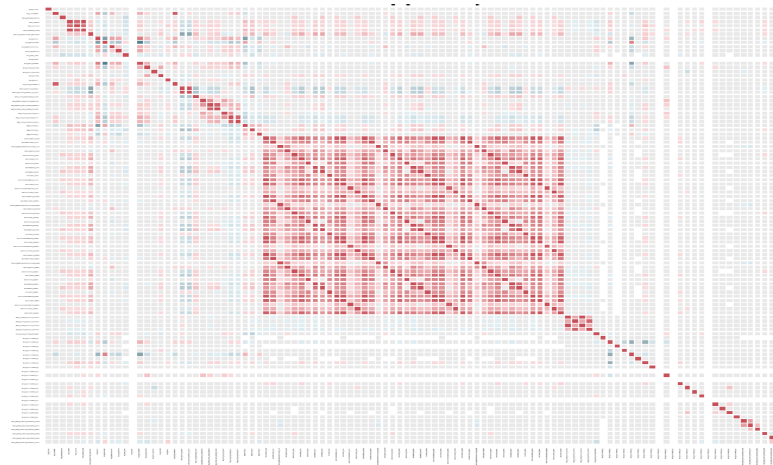
(Embedded method)



- LASSO (L1)
- Ridge regression (L2)
- ElasticNet
- ...

# Матрица корреляций

$$R_x = \begin{pmatrix} 1 & r_{x_1x_2} & \dots & r_{x_1x_n} \\ r_{x_2x_1} & 1 & \dots & r_{x_2x_n} \\ \dots & \dots & \dots & \dots \\ r_{x_nx_1} & r_{x_nx_2} & \dots & 1 \end{pmatrix}$$



# Principal Component Analysis (PCA)

## Метод главных компонент (PCA)

**Постановка задачи PCA (principal component analysis):**

$f_1(x), \dots, f_n(x)$  — исходные числовые признаки

$g_1(x), \dots, g_m(x)$  — новые числовые признаки,  $m \leq n$ ;

Требование: старые признаки должны линейно восстанавливаться по новым:

$$\hat{f}_j(x) = \sum_{s=1}^m g_s(x) u_{js}, \quad j = 1, \dots, n, \quad \forall x \in X,$$

как можно точнее на обучающей выборке  $x_1, \dots, x_l$ :

$$\sum_{i=1}^l \sum_{j=1}^n (\hat{f}_j(x_i) - f_j(x_i))^2 \rightarrow \min_{\{g_s(x_i), u_{js}\}}$$

## Матричные обозначения

Матрицы "объекты-признаки", старая и новая

$$F_{l \times n} = \begin{pmatrix} f_1(x_1) & \dots & f_n(x_1) \\ \dots & \dots & \dots \\ f_1(x_l) & \dots & f_n(x_l) \end{pmatrix}; \quad G_{l \times m} = \begin{pmatrix} g_1(x_1) & \dots & g_m(x_1) \\ \dots & \dots & \dots \\ g_1(x_l) & \dots & g_m(x_l) \end{pmatrix}.$$

Матрица линейного преобразования новых признаков в старые:

$$U_{m \times n} = \begin{pmatrix} u_{11} & \dots & u_{1n} \\ \dots & \dots & \dots \\ u_{n1} & \dots & u_{nm} \end{pmatrix}; \quad \hat{F} = GU^T \approx F.$$

Найти: и новые признаки  $G$ , и преобразование  $U$

$$\sum_{i=1}^l \sum_{j=1}^n (\hat{f}_j(x_i) - f_j(x_i))^2 = \|GU^T - F\|^2 \rightarrow \min_{G, U}$$

## Основная теорема метода PCA

Если  $m \leq \text{rk } F$ , то минимум  $\|GU^T - F\|^2$  достигается, когда

столбцы  $U$  — это собственные векторы матрицы  $F^T F$ , соответствующие  $m$  максимальным собственным значениям  $\lambda_1, \dots, \lambda_m$ , а матрица  $G - F U$ .

При этом:

- матрица  $U$  ортонормированна:  $U^T U = I_m$ ;
- матрица  $G$  ортогональна:  $G^T G = \Lambda = \text{diag}(\lambda_1, \dots, \lambda_m)$ ;
- $U\Lambda = F^T F U$ ;  $GA = FF^T G$ ;
- $\|GU^T - F\|^2 = \|F\|^2 - \text{tr } \Lambda = \sum_{j=m+1}^n \lambda_j$ .

## Связь с сингулярным разложением

Если взять  $m = n$ , то:

- $\|GU^T - F\|^2 = 0$ ;
- представление  $\hat{F} = GU^T = F$  точное и совпадает с сингулярным разложением при  $G = V\sqrt{\Lambda}$ :  
 $F = GU^T = V\sqrt{\Lambda}U^T$ ;  $U^T U = I_m$ ;  $V^T V = I_m$ .
- линейное преобразование  $U$  работает в обе стороны:  
 $F = GU^T$ ;  $G = FU$ .

Поскольку новые признаки некоррелированы ( $G^T G = \Lambda$ ), преобразование  $U$  называется декоррелирующим (или преобразованием Карунена-Лоэва).

## Эффективная размерность выборки

Упорядочим с.з.  $F^T F$  по убыванию:  $\lambda_1 \geq \dots \geq \lambda_n \geq 0$ .

Эффективная размерность выборки — это наименьшее целое  $m$ , при котором:

$$E_m = \frac{\|GU^T - F\|^2}{\|F\|^2} = \frac{\lambda_{m+1} + \dots + \lambda_n}{\lambda_1 + \dots + \lambda_n} \leq \varepsilon$$

## Решение задачи НК для МЛР в новых признаках

Задача наименьших квадратов для МЛР:  $\|F\alpha - y\|^2 \rightarrow \min_{\alpha}$ .

Заменим  $F_{l \times n}$  на ее приближение  $G_{l \times m} \bullet U_{m \times n}^T$ , предполагая  $m \leq n$ :

$$\|GU^T \alpha - y\|^2 = \|G\beta - y\|^2 \rightarrow \min_{\beta}$$

Связь нового и старого вектора коэффициентов:

$$\beta = U^T \alpha; \quad \alpha = U\beta.$$

Решение задачи наименьших квадратов относительно  $\beta$  (единственное отличие —  $m$  слагаемых вместо  $n$ ):

$$\beta^* = D^{-1}V^T y; \quad \alpha^* = UD^{-1}V^T y = \sum_{j=1}^m \frac{1}{\sqrt{\lambda_j}} u_j (v_j^T y);$$

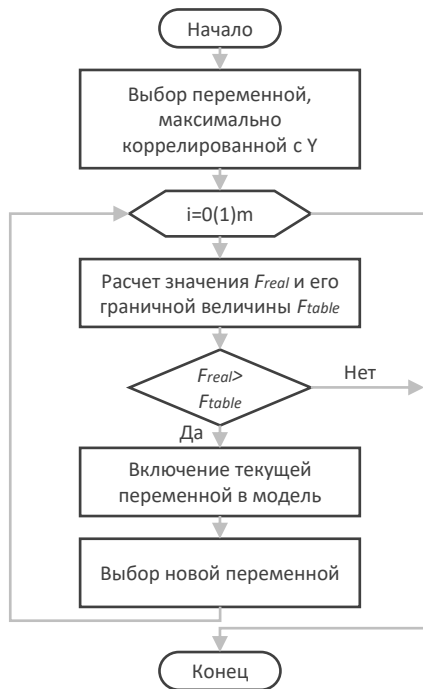
$$G\beta^* = VV^T y = \sum_{j=1}^m v_j (v_j^T y)$$

## Резюме:

- Метод главных компонент позволяет приближать матрицу ее низкоранговым разложением;
- Для этого достаточно взять из SVD-разложения первые  $m$  сингулярных чисел и векторов матрицы;
- Этот прием широко используется в анализе данных — в задачах регрессии, классификации, сжатия данных и др.

# Stepwise Regression

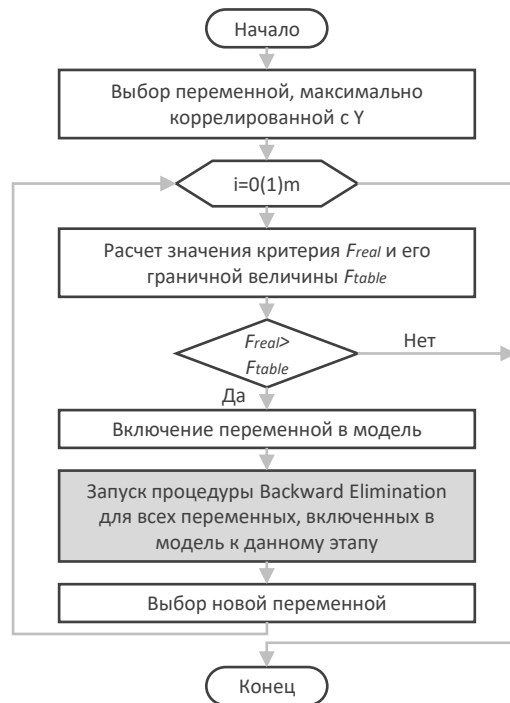
## Forward Selection



## Backward Elimination



## Bidirectional elimination

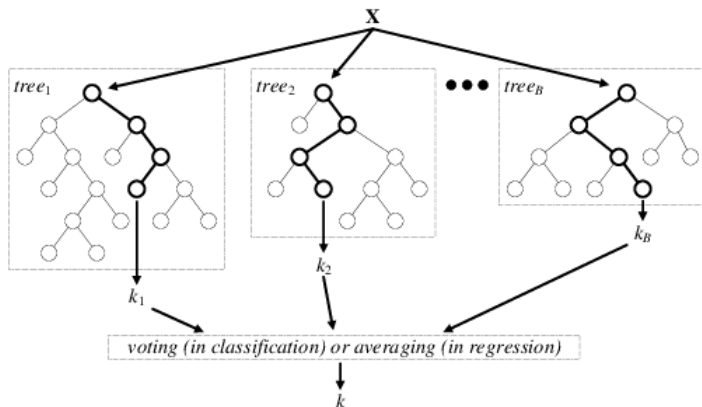


# Random Forest

## Подход 1. Важность на основе уменьшения неоднородности

1. Для каждого дерева случайного леса вычисляем сумму уменьшений неоднородности на всех ветвлениях, связанных с данным предиктором;
2. Итоговую сумму уменьшений неоднородности делим на общее количество деревьев;
3. Повторяем шаги 1-2 для всех переменных.

**Важность** = частота использования переменной в качестве предиктора ветвления.



## Подход 2. Важность на основе уменьшения качества прогнозирования при случайной перестановке (пермутации)

1. Обучить модель;
2. На тестовом/OOB множестве посчитать ошибку;
3. Зафиксировать переменную/группу переменных, случайно переставить значения на тестовом/OOB множестве;
4. По новой выборке посчитать ошибку;
5. Вычислить разность ошибок на исходном множестве и множестве с перестановкой.

Var 1	Var 2	Var 3	Target	Predict	Right
1	2	101	0	0	x
2	3	102	1	0	
3	5	103	1	1	x
4	7	104	0	0	x

accuracy 75%

Var 1	Var 2	Var 3	Target	Predict	Right
1	5	101	0	1	
2	7	102	1	0	
3	2	103	1	1	x
4	3	104	0	0	x

accuracy 50%

# Регуляризация L1 (LASSO) и L2 (Ridge)

## Регрессия LASSO

LASSO — Least Absolute Shrinkage and Selection Operator, два эквивалентных варианта постановки задачи:

$$Q(\alpha) = \|F\alpha - y\|^2 \rightarrow \min_{\alpha} \quad \text{при} \quad \sum_{j=1}^n |\alpha_j| \leq \chi;$$

$$Q(\alpha) = \|F\alpha - y\|^2 + \tau \sum_{j=1}^n |\alpha_j| \rightarrow \min_{\alpha};$$

После замены переменных

$$\begin{cases} \alpha_j = \alpha_j^+ - \alpha_j^-; \\ |\alpha_j| = \alpha_j^+ + \alpha_j^-; \end{cases} \quad \alpha_j^+ \geq 0; \quad \alpha_j^- \geq 0.$$

ограничения принимают канонический вид:

$$\sum_{j=1}^n \alpha_j^+ + \alpha_j^- \leq \chi; \quad \alpha_j^+ \geq 0; \quad \alpha_j^- \geq 0.$$

Чем меньше  $\chi$ , тем больше  $j$  таких, что  $\alpha_j^+ = \alpha_j^- = 0$ .

## Резюме:

LASSO обнуляет веса и приводит к отбору признаков в линейных моделях.

Источник: Coursera, К.В. Воронцов

## Гребневая регрессия (Ridge Regression)

Штраф за увеличение нормы вектора весов  $\|\alpha\|$ :

$$Q_{\tau}(\alpha) = \|F\alpha - y\|^2 + \frac{\tau}{2} \|\alpha\|^2,$$

где  $\tau$  — неотрицательный параметр регуляризации.

Модифицированное МНК-решение ( $\tau I_n$  — "гребень")

$$\alpha_{\tau}^* = (F^T F + \tau I_n)^{-1} F^T y.$$

Преимущество сингулярного разложения:

можно подбирать параметр  $\tau$ , вычислив SVD только один раз.

## Регуляризованный МНК через сингулярное разложение

Вектор регуляризованного МНК-решения  $\alpha_{\tau}^*$  и МНК-аппроксимация целевого вектора  $F\alpha_{\tau}^*$ :

$$\alpha_{\tau}^* = U(D^2 + \tau I_n)^{-1} D V^T y = \sum_{j=1}^n \frac{\sqrt{\lambda_j}}{\lambda_j + \tau} u_j (v_j^T y);$$

$$F\alpha_{\tau}^* = V D U^T \alpha_{\tau}^* = V \text{diag} \left( \frac{\sqrt{\lambda_j}}{\lambda_j + \tau} \right) V^T y = \sum_{j=1}^n \frac{\lambda_j}{\lambda_j + \tau} v_j (v_j^T y);$$

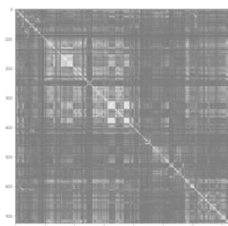
$$\|\alpha_{\tau}^*\|^2 = \|(D^2 + \tau I_n)^{-1} D V^T y\|^2 = \sum_{j=1}^n \frac{\lambda_j}{(\lambda_j + \tau)^2} (v_j^T y)^2$$

$F\alpha_{\tau}^* \neq F\alpha^*$ , но зато решение становится гораздо устойчивее.

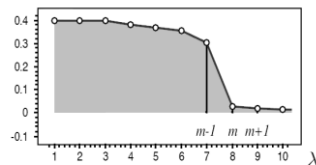


# КЛАССИКИ

## Matrix (Scoring)



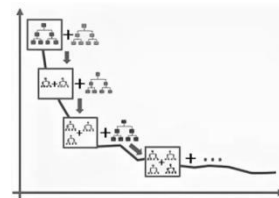
## PCA (ML)



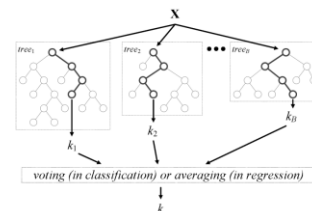
$$E_m = \frac{\|GU^T - F\|^2}{\|F\|^2} = \frac{\lambda_{m+1} + \dots + \lambda_n}{\lambda_1 + \dots + \lambda_n} \leq \varepsilon$$

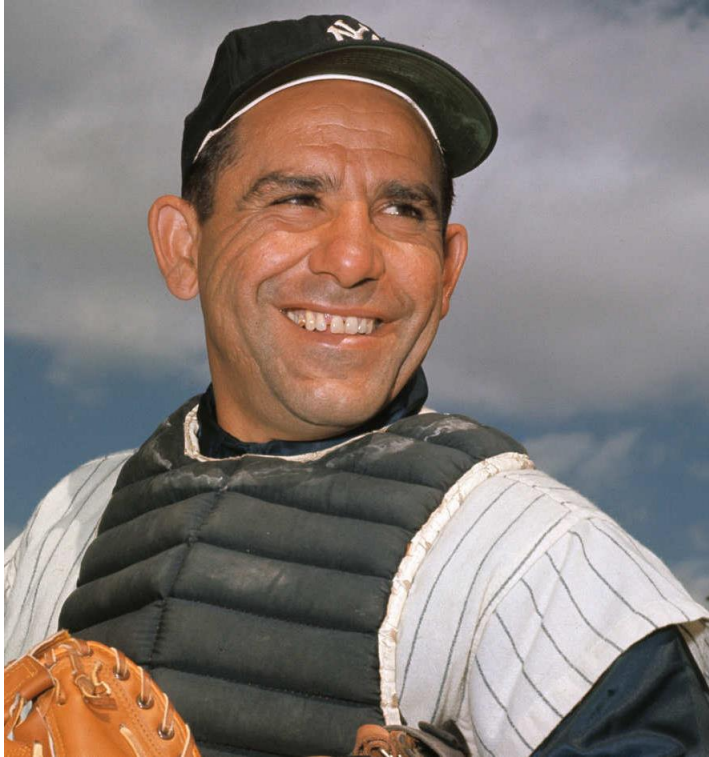
# НЕОФИТЫ

## Boosting (Kaggle)



## Forest (Science)

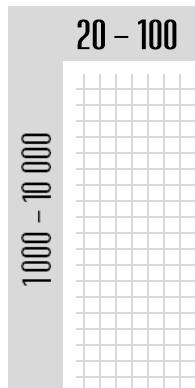
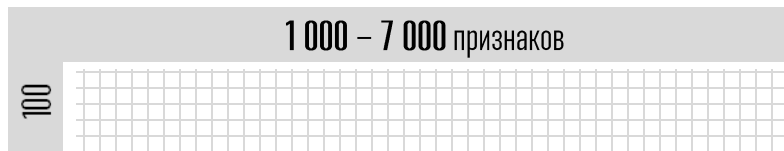




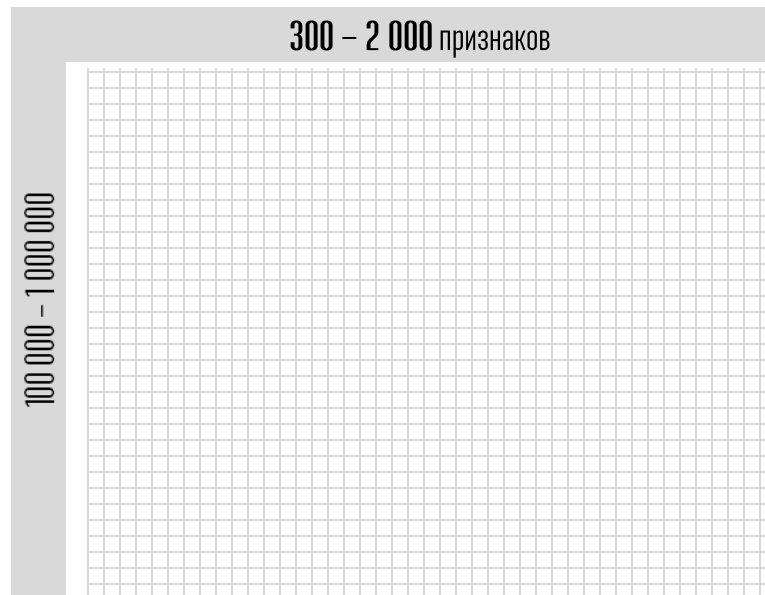
– В теории нет разницы  
между теорией и практикой.  
А на практике есть.

**Йоги Берра**

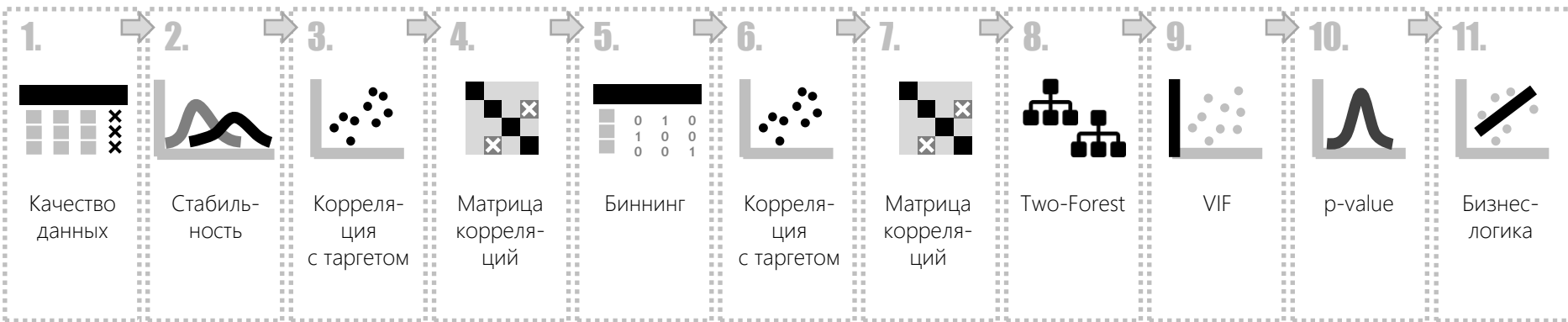
# В научных статьях



# В реальности



# Схема отбора переменных



# Шаг 1. Анализ качества данных

## Критерии качества данных

- Точность (accuracy)
- Полнота (completeness)
- Согласованность (consistency)
- Достоверность (credibility)
- Правильность (correctness)
- Доступность (availability)
- ...

## Этапы оценки качества данных

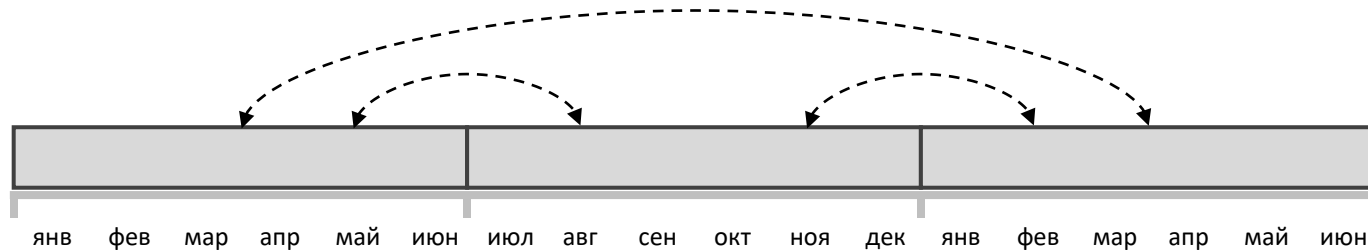
1. На стадии сбора данных
- 2. На стадии разработки модели**
3. На стадии эксплуатации модели

## Проблемы снижения качества данных

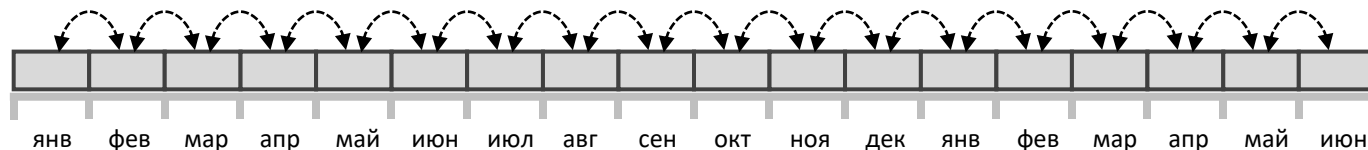
- Пропуски или неполнота;
- Орфографические ошибки;
- Аномалии;
- Фиктивные значения;
- Логические несоответствия;
- Закодированные значения;
- Несоответствие форматов;
- Дублирование;
- Избыточность информации;
- ...

# Шаг 2. Проверка стабильности

Большие  
периоды



Маленькие  
периоды



## Шаг 2. Проверка стабильности

	Стабильна	Слегка нестабильна	Нестабильна
KS	[0%; 10%)	[10%; 25%)	[25%; 100%]
S	[0%; 10%)	[10%; 25%)	[25%; 100%]
PSI	[0; 10)	[10; 25)	[25; ∞)
<b>Присвоенный вес</b>	<b>0</b>	<b>0,5</b>	<b>1</b>
Большие периоды: $\text{Big} = \max(KS_i) + \max(S_i) + \max(PSI_i)$	[0; 1]	[1; 2)	[2; 3)
Маленькие периоды: $\text{Small} = \text{Avg}(KS_j + S_j + PSI_j)$	[0; 1]	[1; 1,5)	[1,5; 3)
<b>Присвоенный вес</b>	<b>0</b>	<b>0,5</b>	<b>1</b>
Максимум (консервативный подход): $\text{Max}(\text{Big}; \text{Small})$	{0}	{0,5}	{1}
Усреднение (лояльный подход): $0,5 * \text{Big} + 0,5 * \text{Small}$	[0; 0,2)	[0,2; 0,6)	[0,6; 1]
<b>Присвоенный вес</b>	<b>0</b>	<b>0,5</b>	<b>1</b>

# Шаг 3. Корреляция с целевой переменной

## Бинарная целевая переменная

- Для непрерывных признаков рассчитываются **статистика Стьюдента** (если признак распределен нормально) или **тест Манна-Уитни**;
- Для категориальных и бинарных признаков рассчитывается **Хи-квадрат** критерий Пирсона.

## Категориальная целевая переменная

- Для непрерывных признаков проводится тест **ANOVA**;
- Для категориальных и бинарных признаков рассчитывается **Хи-квадрат** критерий Пирсона.

## Непрерывная целевая переменная

- Для непрерывных признаков рассчитывается **корреляция Пирсона**;
- Для категориальных признаков проводится тест **ANOVA**;
- Для бинарных признаков рассчитываются **статистика Стьюдента** (если признак распределен нормально) или **тест Манна-Уитни**.



# Шаг 4. Матрица корреляций

$$R_x = \begin{pmatrix} 1 & r_{x_1x_2} & \dots & r_{x_1x_n} \\ r_{x_2x_1} & 1 & \dots & r_{x_2x_n} \\ \dots & \dots & \dots & \dots \\ r_{x_nx_1} & r_{x_nx_2} & \dots & 1 \end{pmatrix}$$

- Для непрерывных признаков рассчитывается **корреляция Пирсона**;
- Для категориальных и бинарных признаков рассчитывается **корреляция Спирмена**

# Шаг 5-7. Бинаризация и корреляции

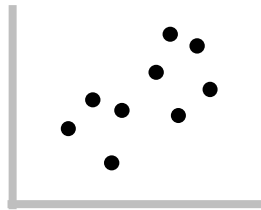
## Шаг 5

Бинаризация признаков

	0	1	0
	1	0	0
	0	0	1
	1	1	0

## Шаг 6

Корреляция с таргетом



## Шаг 7

Матрица корреляций

$$R_x = \begin{pmatrix} 1 & r_{x_1x_2} & \dots & r_{x_1x_n} \\ r_{x_2x_1} & 1 & \dots & r_{x_2x_n} \\ \dots & \dots & \dots & \dots \\ r_{x_nx_1} & r_{x_nx_2} & \dots & 1 \end{pmatrix}$$

# Шаг 8. Two-Forest

Важность – уменьшение качества прогнозирования при случайной перестановке

$$VI_j = P(Y \neq f(X_1, \dots, X_j^*, \dots, X_p)) - P(Y \neq f(X_1, \dots, X_j, \dots, X_p))$$

1. Исходная выборка делится на две подвыборки;
2. На каждой из двух подвыборок строится случайный лес;
3. Для построенных моделей для каждой  $j$ -ой переменной вычисляется важность при перестановке  $VI_j$  на множестве, которое не использовалось для построения модели;
4. Определяются множества:  
 $M1 = \{\text{все отрицательные важности}\},$   
 $M2 = \{\text{все нулевые важности}\},$   
 $M3 = \{\text{все отрицательные важности} * (-1)\};$
5. На  $M = M1 \cup M2 \cup M3$  строится плотность распределения  $F$ ;
6. Для каждой  $j$ -ой переменной рассчитывается  $p\text{-value} = 1 - F(VI_j)$



# Шаг 9. Мультиколлинеарность (VIF)

- 1 Для каждого признака  $X_i$  обучается линейная регрессия, где  $X_i$  является функций от всех остальных признаков:

$$X_i = \beta_0 + \sum_{j=1}^k \beta_j X_j, \quad i \neq j$$

- 2 Для каждого обученного признака рассчитывается коэффициент VIF:

$$VIF_i = \frac{1}{1 - R_i^2}$$

- 3 Признаки со значением  $VIF > 10$  относятся к мультиколлинеарным. Из всех мультиколлинеарных признаков удаляется признак с максимальным значением VIF;
- 4 Шаги 1-3 итерационно повторяются до тех пор, пока максимальное значение VIF по всем оставшимся признакам не станет меньше или равно 10.

# Шаг 10. Статистическая значимость

Проверку статистической значимости признаков можно делать с помощью различных тестов, среди которых тест отношения правдоподобия, тест Вальда и тест множителей Лагранжа (все три теста асимптотически эквивалентны).

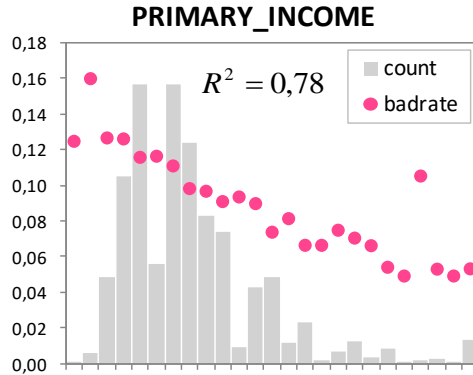
**Тест отношения правдоподобия:** для проверки нулевой гипотезы сравниваются функции правдоподобия полной модели и укороченной модели без тестируемого признака. Для этого рассчитывается статистика отношения правдоподобия:

$$LR = 2 \cdot (L_l - L_s) = 2 \cdot \ln \frac{L_l}{L_s}$$

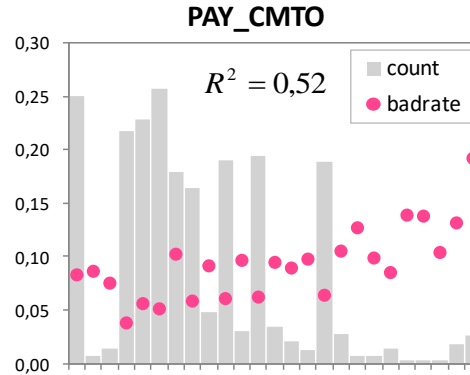
Где  $L_l$  – значение логарифмической функции правдоподобия полной модели,  
 $L_s$  – значение логарифмической функции правдоподобия укороченной модели.

# Шаг 11. Экспертный анализ

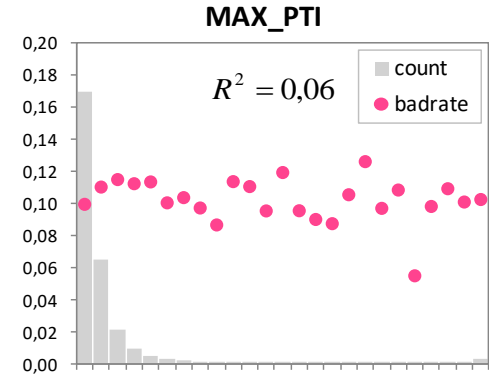
Сильная зависимость



Умеренная зависимость



Слабая зависимость



# Matrix

# Forward

# Two-Forest

- 1 Качество
- 2 Стабильность
- 3 Таргет-корреляция
- 4 Матрица корреляций
- 5 Биннинг
- 6 Таргет-корреляция
- 7 Матрица корреляций
- 8 **Gini**
- 9 VIF
- 10 P-value

- 1 Качество
- 2 Стабильность
- 3 Таргет-корреляция
- 4 Матрица корреляций
- 5 Биннинг
- 6 Таргет-корреляция
- 7 Матрица корреляций
- 8 **Forward**
- 9 VIF
- 10 P-value

- 1 Качество
- 2 Стабильность
- 3 Таргет-корреляция
- 4 Матрица корреляций
- 5 Биннинг
- 6 Таргет-корреляция
- 7 Матрица корреляций
- 8 **Two-Forest**
- 9 VIF
- 10 P-value

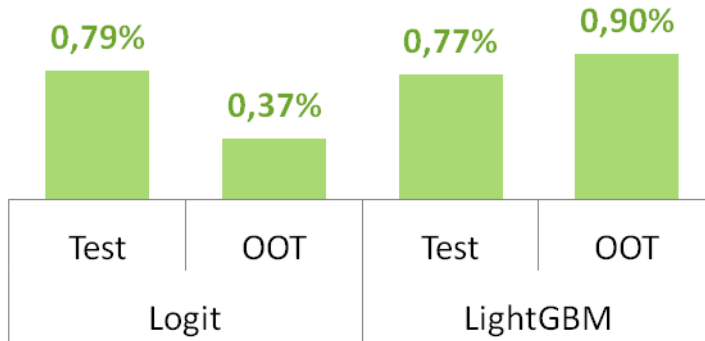
# Бизнес-задачи

	Кол-во наблюдений	Кол-во признаков
PTB (CRM)	303 220	1 222
Application PD (Scoring)	497 063	423
Behavioral PD (Scoring)	588 385	1 087
Allocation (Collection)	172 250	162

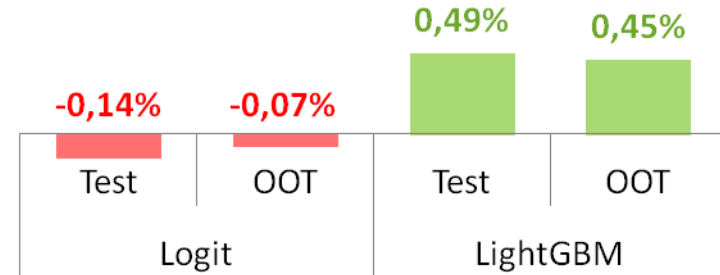


# Результаты (разница точности, Gini)

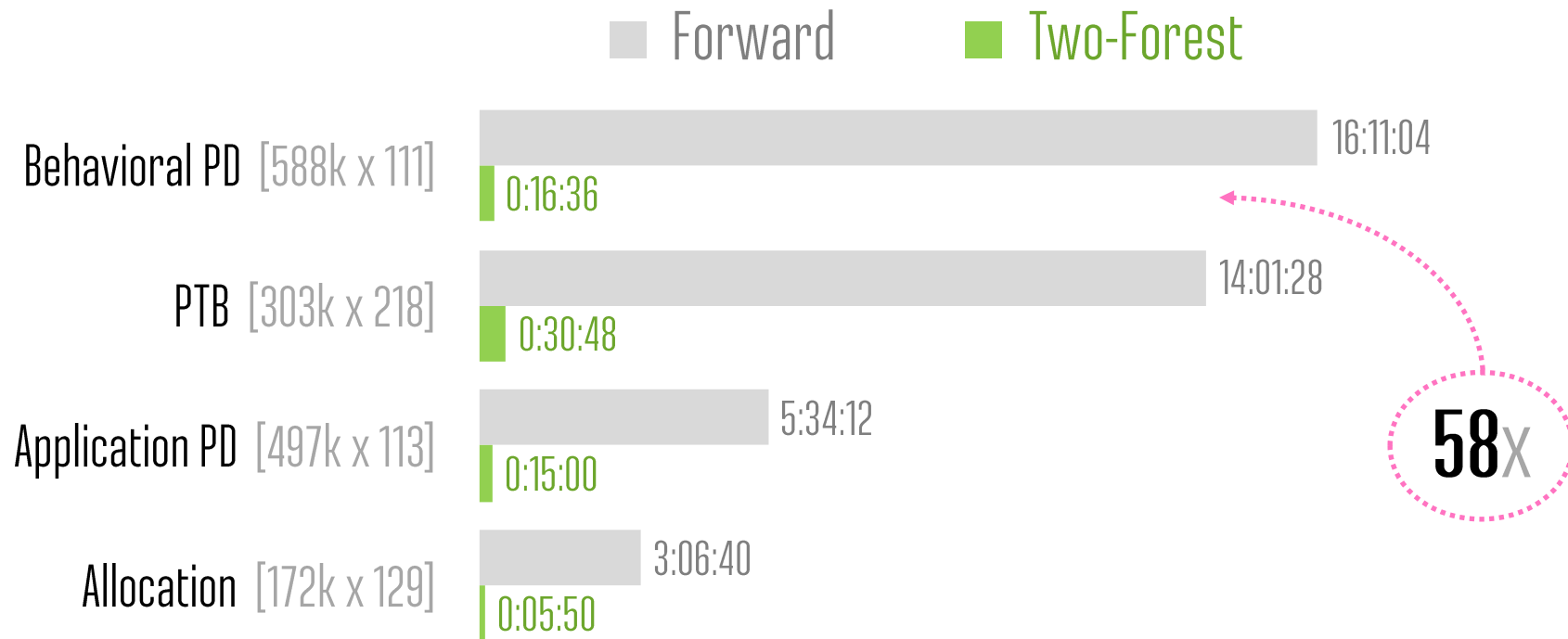
$\Delta(\text{TwoForest-Matrix})$



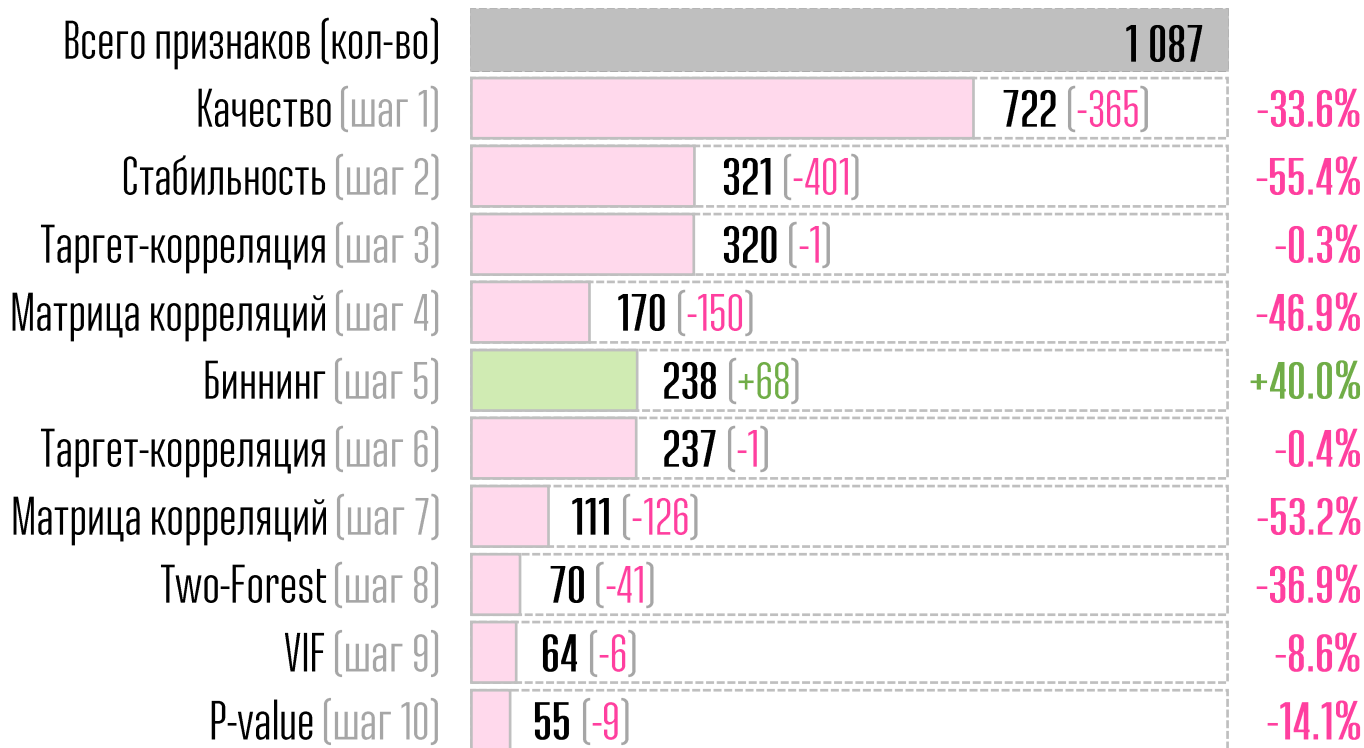
$\Delta(\text{TwoForest-Forward})$



# Результаты (скорость работы методов)



# Воронка отбора признаков (Behavioral PD)



# Выводы

1

Комбинированный отбор хорошо работает с большим количеством **признаков**: фильтры отсеивают мусорные признаки, обертки – учитывают многомерные зависимости

2

Проверка качества и стабильности переменных **позволяет использовать один большой универсальный набор признаков** для разных типов задач.

3

**Two-Forest** в сравнении с **Forward** **лучше по качеству для нелинейных моделей** и сопоставим по качеству для линейных моделей. По скорости **Two-Forest** в **десятки раз быстрее** регрессионных методов.

4

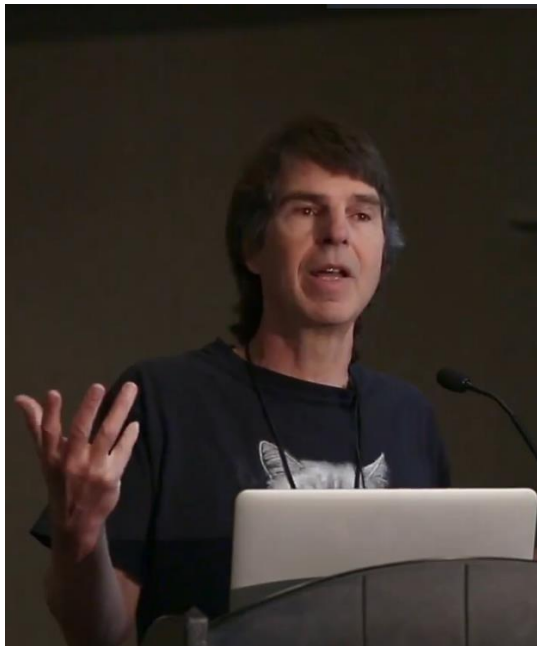
**Включение** в комбинированную схему отбора **нескольких оберток** позволяет контролировать корректность работы методов на предыдущих шагах



– С точки зрения академической науки, когда я начинал заниматься машинным обучением 25 лет назад, тот научный коллектив, в который я пришел еще студентом, в общем-то жил с полной уверенностью (и она была основана на примерно 30-летнем опыте предыдущих исследований), что **задачу можно решать любым методом.**

**Константин Воронцов, 2018**

# Теорема о бесплатных завтраках



В среднем по всем возможным порождающим определениям у любого алгоритма классификации частота ошибок классификации ранее не наблюдавшихся примеров одинакова. Самый изощренный алгоритм, который мы только можем придумать, в среднем (по всем возможным задачам) дает такое же качество, как простейшее предсказание: все точки принадлежат одному классу.

**David H. Wolpert, 1996**

Спасибо за  
внимание!

Афанасьев Сергей

Исполнительный директор

Начальник управления статистического анализа

КБ «Ренессанс Кредит»