

Ансамблирование нейронных сетей

Панфилов Борис

Введение

Ансамбли в классическом ML

Deep ensembles

Snapshot ensembles

Fast Geometric ensembles

Dropout ensembles

Заключение

Введение

Ансамбли в классическом ML

Deep ensembles

Snapshot ensembles

Fast Geometric ensembles

Dropout ensembles

Заключение

Мотивация

- Улучшить качество модели
- Научиться определять случаи, когда модель не уверена в ответе

Возможные ограничения

- Время и стоимость обучения
- Память
- Инференс

Идея ансамблирования

1. Обучить K моделей
2. На этапе теста в качестве итогового ответа брать среднее арифметическое результатов этих моделей

Введение

Ансамбли в классическом ML

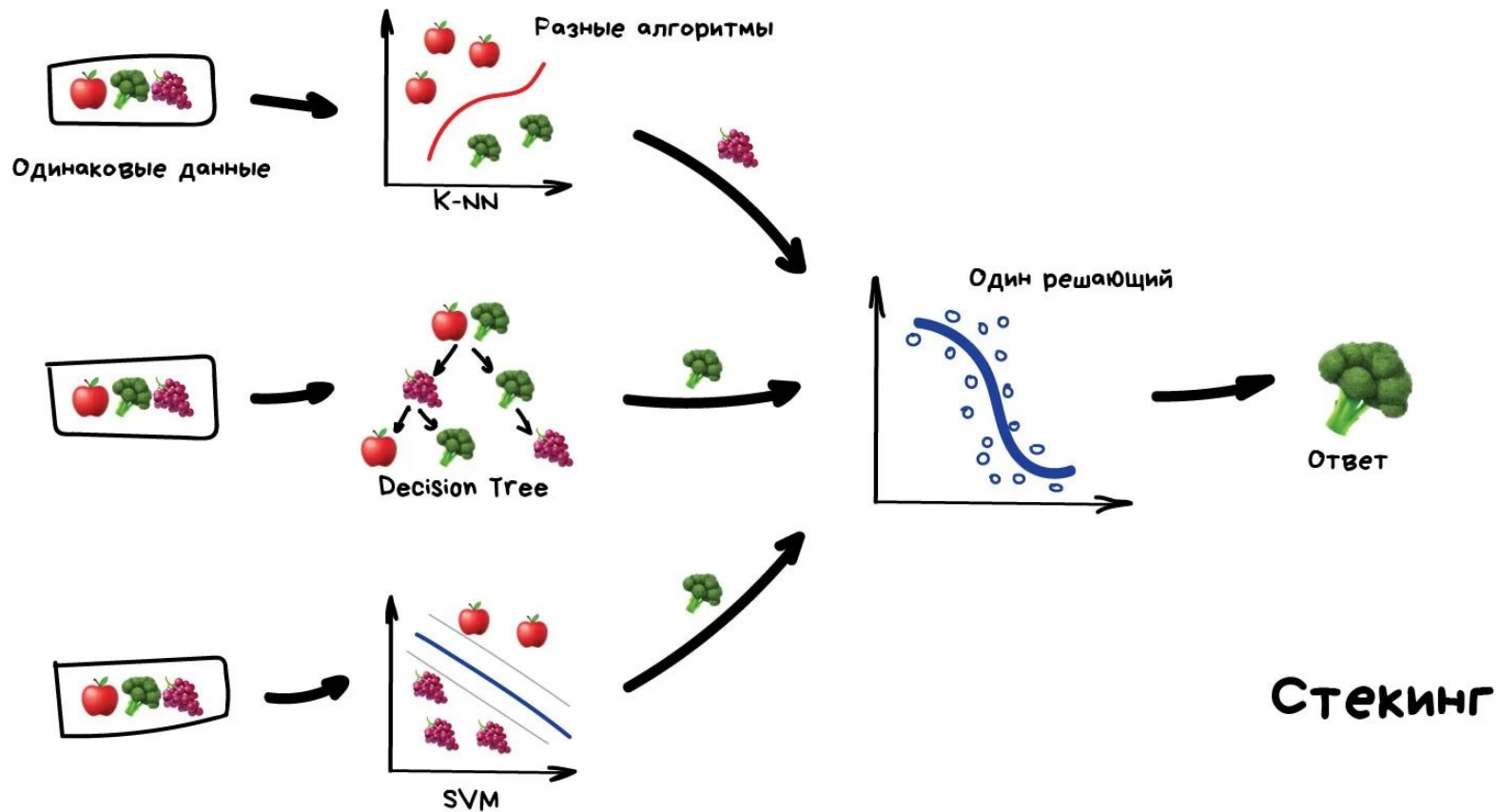
Deep ensembles

Snapshot ensembles

Fast Geometric ensembles

Dropout ensembles

Заключение



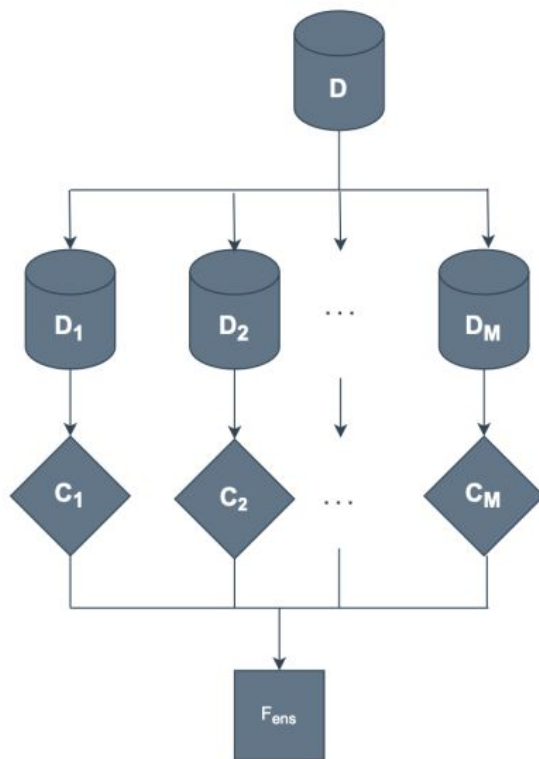


Figure 2: Bagging

$$b_1(x) := \arg \min_{b \in \mathcal{A}} \frac{1}{2} \sum_{i=1}^{\ell} (b(x_i) - y_i)^2$$

$$s_i^{(1)} = y_i - b_1(x_i)$$

$$b_2(x) := \arg \min_{b \in \mathcal{A}} \frac{1}{2} \sum_{i=1}^{\ell} (b(x_i) - s_i^{(1)})^2$$

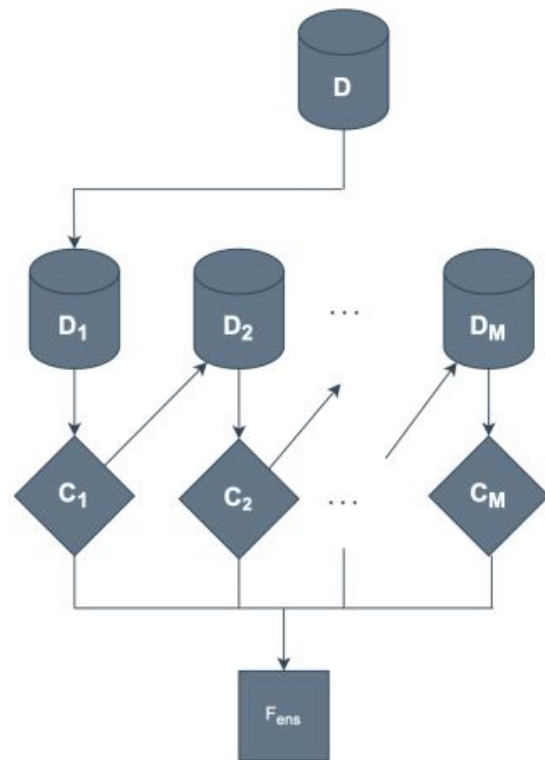


Figure 3: Boosting

Введение

Ансамбли в классическом ML

Deep ensembles

Snapshot ensembles

Fast Geometric ensembles

Dropout ensembles

Заключение

Идея

1. Обучаем много независимых сеток
2. Ансамблируем путем усреднения их предсказаний

Плюсы:

- Очень хорошее качество модели
- Очень хорошая оценка неопределенности

Минусы:

- Требуется в K раз больше времени и ресурсов
- Требуется в K раз больше памяти

Введение

Ансамбли в классическом ML

Deep ensembles

Snapshot ensembles

Fast Geometric ensembles

Dropout ensembles

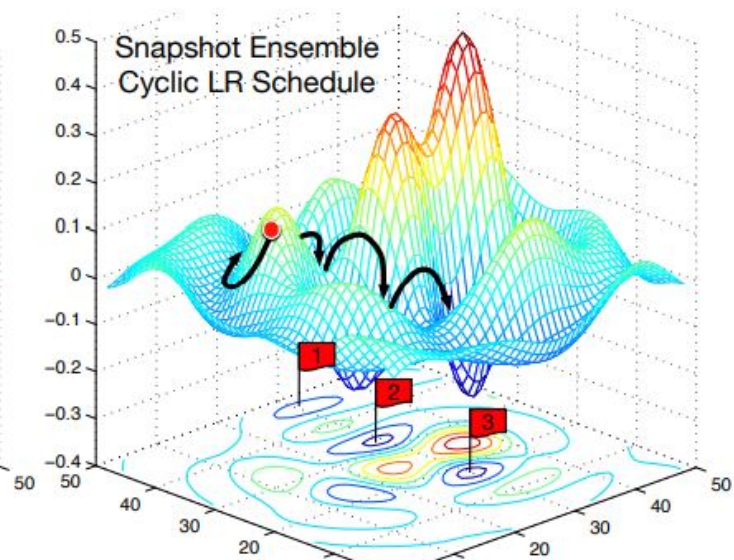
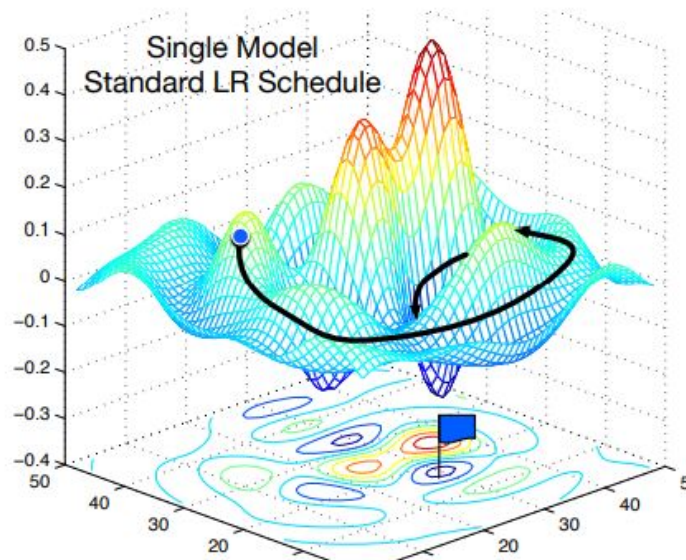
Заключение

Мотивация

- Хотим ансамбль, но без дополнительных затрат на ресурсы и время обучения

Идея

- Делаем несколько снэпшотов во время обучения, при этом обучаем, используя циклическое расписание lr
- Активно используем, что во время такого обучения модель может попасть в разные локальные минимумы



Learning rate

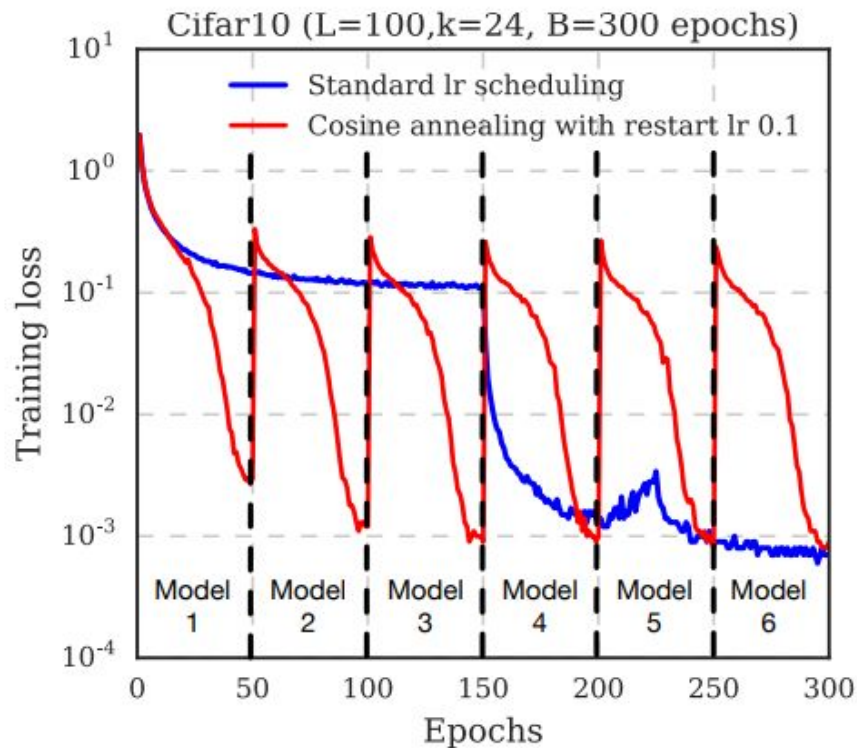
$$\alpha(t) = \frac{\alpha_0}{2} \left(\cos \left(\frac{\pi \text{mod}(t - 1, \lceil T/M \rceil)}{\lceil T/M \rceil} \right) + 1 \right)$$

t - номер итерации

T - всего итераций

M - количество циклов

α_0 - изначальный lr, гиперпараметр



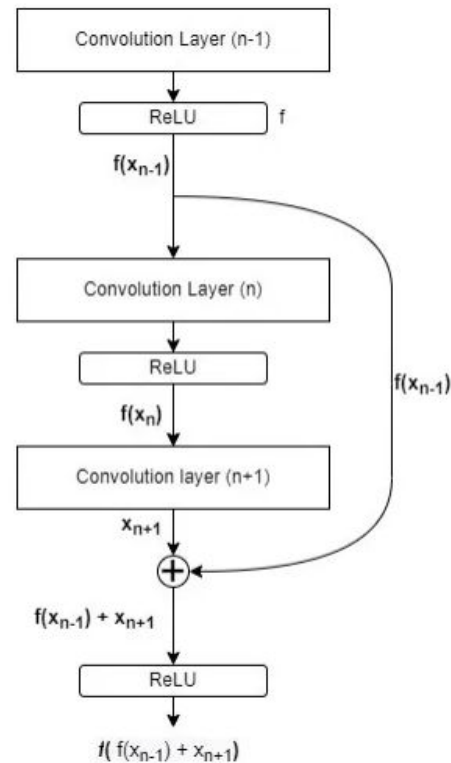
Ликбез про датасеты и модели

1. Датасеты

- 1.1. CIFAR-10
- 1.2. CIFAR-100
- 1.3. SVHN
- 1.4. ImageNet
- 1.5. Tiny ImageNet

2. Модели

- 2.1. VGG
- 2.2. ResNet
- 2.3. Wide ResNet
- 2.4. DenseNet



	Method	C10	C100	SVHN	Tiny ImageNet
ResNet-110	Single model	5.52	28.02	1.96	46.50
	NoCycle Snapshot Ensemble	5.49	26.97	1.78	43.69
	SingleCycle Ensembles	6.66	24.54	1.74	42.60
	Snapshot Ensemble ($\alpha_0 = 0.1$)	5.73	25.55	1.63	40.54
	Snapshot Ensemble ($\alpha_0 = 0.2$)	5.32	24.19	1.66	39.40
Wide-ResNet-32	Single model	5.43	23.55	1.90	39.63
	Dropout	4.68	22.82	1.81	36.58
	NoCycle Snapshot Ensemble	5.18	22.81	1.81	38.64
	SingleCycle Ensembles	5.95	21.38	1.65	35.53
	Snapshot Ensemble ($\alpha_0 = 0.1$)	4.41	21.26	1.64	35.45
	Snapshot Ensemble ($\alpha_0 = 0.2$)	4.73	21.56	1.51	32.90
DenseNet-40	Single model	5.24*	24.42*	1.77	39.09
	Dropout	6.08	25.79	1.79*	39.68
	NoCycle Snapshot Ensemble	5.20	24.63	1.80	38.51
	SingleCycle Ensembles	5.43	22.51	1.87	38.00
	Snapshot Ensemble ($\alpha_0 = 0.1$)	4.99	23.34	1.64	37.25
	Snapshot Ensemble ($\alpha_0 = 0.2$)	4.84	21.93	1.73	36.61
DenseNet-100	Single model	3.74*	19.25*	-	-
	Dropout	3.65	18.77	-	-
	NoCycle Snapshot Ensemble	3.80	19.30	-	-
	SingleCycle Ensembles	4.52	18.38	-	-
	Snapshot Ensemble ($\alpha_0 = 0.1$)	3.57	18.12	-	-
	Snapshot Ensemble ($\alpha_0 = 0.2$)	3.44	17.41	-	-

Table 1: Error rates (%) on CIFAR-10 and CIFAR-100 datasets. All methods in the same group are trained for the same number of iterations. Results of our method are colored in **blue**, and the best result for each network/dataset pair are **bolded**. * indicates numbers which we take directly from Huang et al. (2016a).

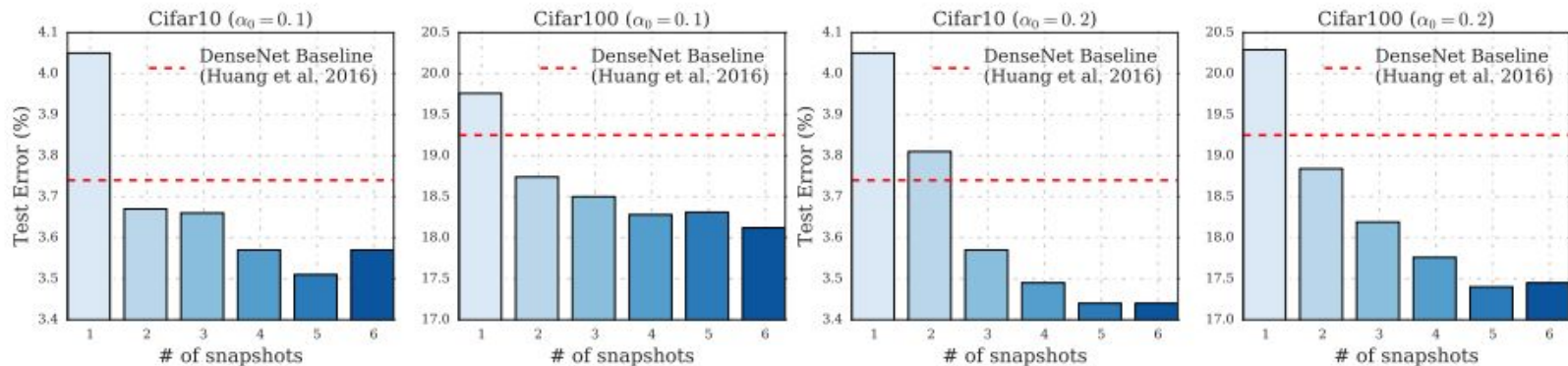
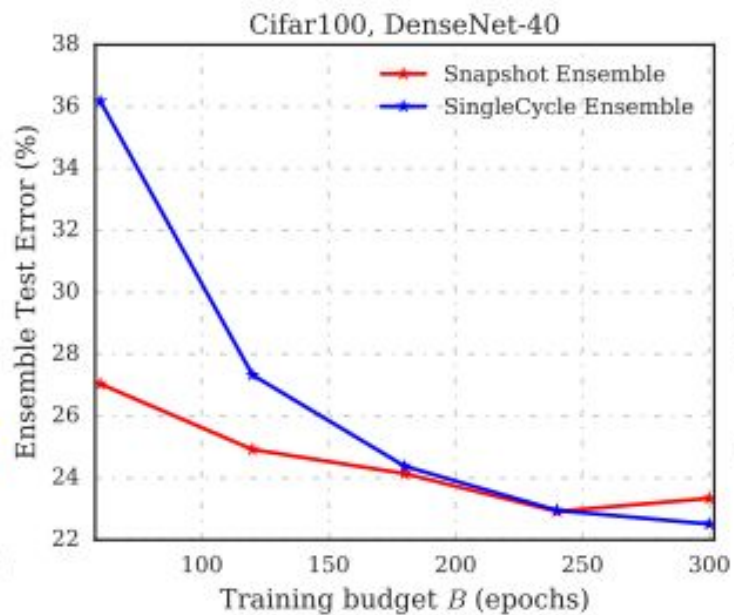
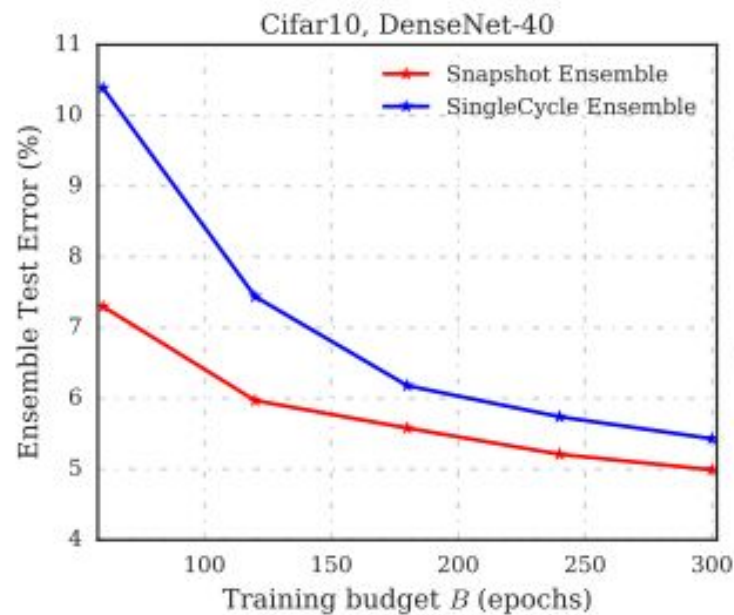


Figure 3: DenseNet-100 Snapshot Ensemble performance on CIFAR-10 and CIFAR-100 with restart learning rate $\alpha_0 = 0.1$ (left two) and $\alpha_0 = 0.2$ (right two). Each ensemble is trained with $M = 6$ annealing cycles (50 epochs per each).



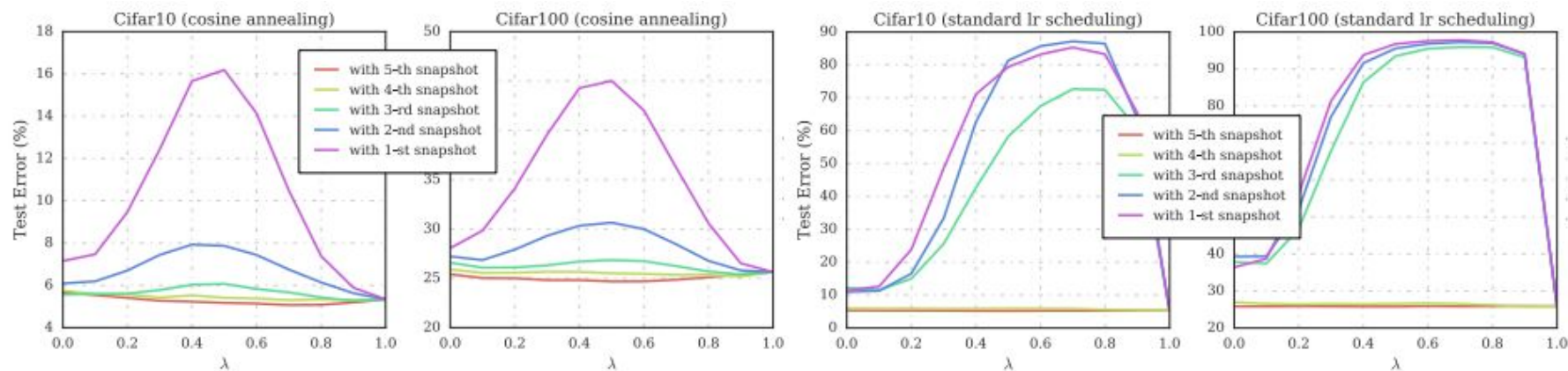


Figure 5: Interpolations in parameter space between the final model (sixth snapshot) and all intermediate snapshots. $\lambda = 0$ represents an intermediate snapshot model, while $\lambda = 1$ represents the final model. **Left:** A Snapshot Ensemble, with cosine annealing cycles ($\alpha_0 = 0.2$ every $B/M = 50$ epochs). **Right:** A NoCycle Snapshot Ensemble, (two learning rate drops, snapshots every 50 epochs).

Введение

Ансамбли в классическом ML

Deep ensembles

Snapshot ensembles

Fast Geometric ensembles

Dropout ensembles

Заключение

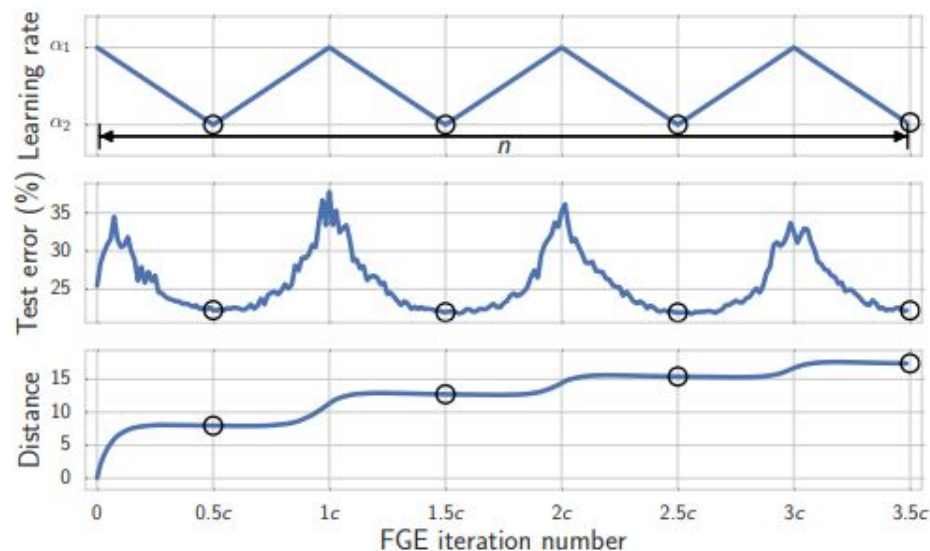
Мотивация

- Между любыми двумя локальными минимумами есть кривая, такая что вдоль нее значение функции потерь слабо изменяется

Идея

1. Сходимся в какой-то локальный минимум
2. Копируем модель и для копии проделываем еще несколько эпох с триангулярным l_r , сохраняя чекпоинты в середине каждой эпохи

Triangular learning rate



$$\alpha(i) = \begin{cases} (1 - 2t(i))\alpha_1 + 2t(i)\alpha_2 & 0 < t(i) \leq \frac{1}{2} \\ (2 - 2t(i))\alpha_2 + (2t(i) - 1)\alpha_1 & \frac{1}{2} < t(i) \leq 1 \end{cases},$$

где $t(i) = \frac{1}{c} * (\text{mod}(i - 1, c) + 1)$,

$\alpha_1 > \alpha_2$ - learning rates,

c - количество итераций в одном цикле.

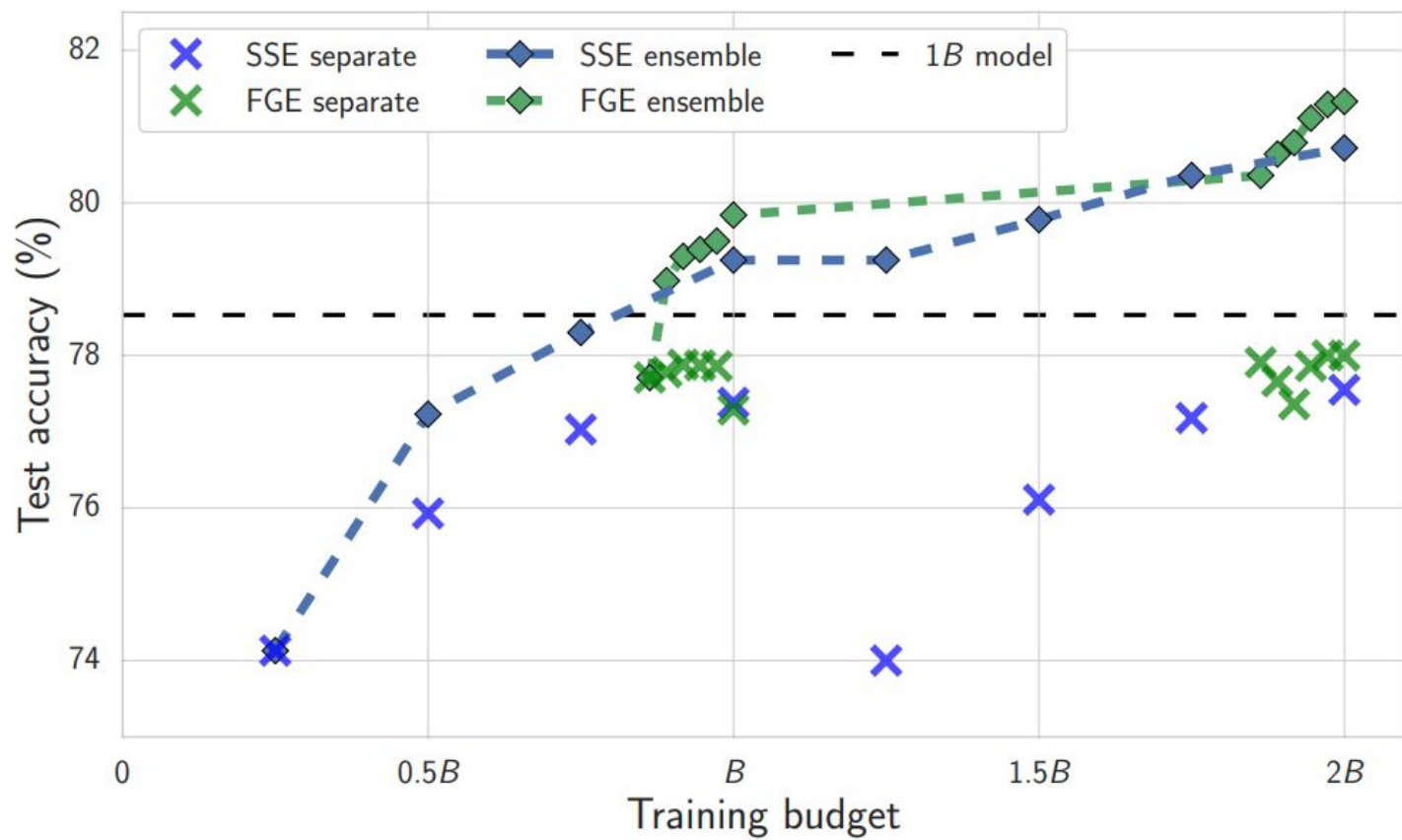


Table 1: Error rates (%) on CIFAR-100 and CIFAR-10 datasets for different ensembling techniques and training budgets. The best results for each dataset, architecture, and budget are **bolded**.

DNN (Budget)	method	CIFAR-100			CIFAR-10		
		1B	2B	3B	1B	2B	3B
VGG-16 (200)	Ind	27.4 ± 0.1	25.28	24.45	6.75 ± 0.16	5.89	5.9
	SSE	26.4 ± 0.1	25.16	24.69	6.57 ± 0.12	6.19	5.95
	FGE	25.7 ± 0.1	24.11	23.54	6.48 ± 0.09	5.82	5.66
ResNet-164 (150)	Ind	21.5 ± 0.4	19.04	18.59	4.72 ± 0.1	4.1	3.77
	SSE	20.9 ± 0.2	19.28	18.91	4.66 ± 0.02	4.37	4.3
	FGE	20.2 ± 0.1	18.67	18.21	4.54 ± 0.05	4.21	3.98
WRN-28-10 (200)	Ind	19.2 ± 0.2	17.48	17.01	3.82 ± 0.1	3.4	3.31
	SSE	17.9 ± 0.2	17.3	16.97	3.73 ± 0.04	3.54	3.55
	FGE	17.7 ± 0.2	16.95	16.88	3.65 ± 0.1	3.38	3.52

Введение

Ансамбли в классическом ML

Deep ensembles

Snapshot ensembles

Fast Geometric ensembles

Dropout ensembles

Заключение

Мотивация

- Хотим ансамбль, но без дополнительных затрат на хранение нескольких нейросетей

Идея

1. Обучаем модель, используя дропаут
2. На этапе тестирования несколько раз запускаем модель с включенным дропаутом и усредняем результаты

Введение

Ансамбли в классическом ML

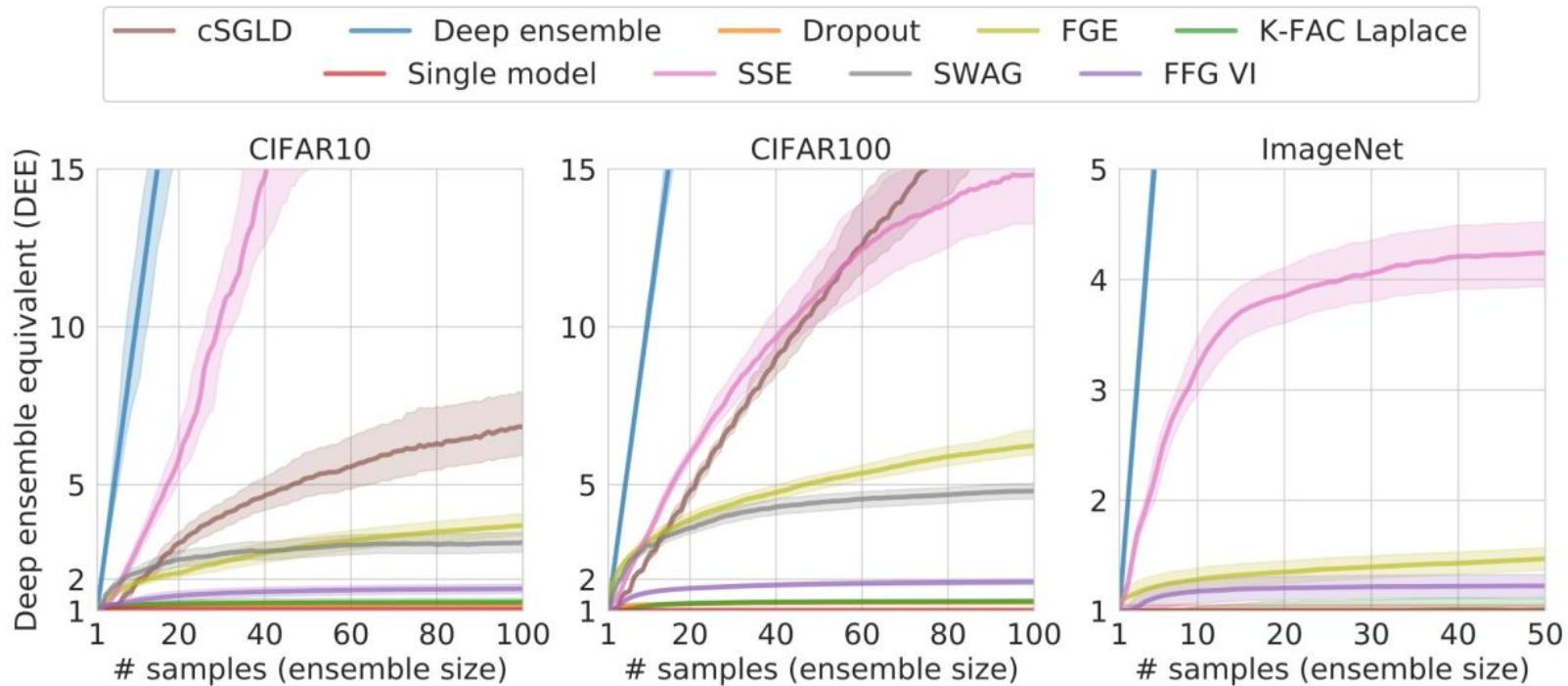
Deep ensembles

Snapshot ensembles

Fast Geometric ensembles

Dropout ensembles

Заключение



$$\text{DEE}_m(k) = \min \left\{ l \in \mathbb{R}, l \geq 1 \mid \text{CLL}_{DE}^{\text{mean}}(l) \geq \text{CLL}_m^{\text{mean}}(k) \right\},$$

Model	Method	Error (%)				Negative calibrated log-likelihood			
		1	5	10	100	1	5	10	100
VGG16 CIFAR-10	Dropout	5.86 \pm 0.09	5.81 \pm 0.08	5.82 \pm 0.06	5.79 \pm 0.07	0.232 \pm 0.005	0.225 \pm 0.004	0.224 \pm 0.004	0.223 \pm 0.003
	SWA-Gaussian	7.03 \pm 0.50	5.66 \pm 0.08	5.49 \pm 0.12	5.25 \pm 0.13	0.230 \pm 0.014	0.182 \pm 0.003	0.171 \pm 0.002	0.160 \pm 0.002
	Cyclic SGLD	7.37 \pm 0.16	6.56 \pm 0.09	5.71 \pm 0.06	4.84 \pm 0.04	0.234 \pm 0.004	0.196 \pm 0.004	0.176 \pm 0.003	0.147 \pm 0.003
	Fast Geometric Ens.	6.52 \pm 0.16	5.95 \pm 0.16	5.69 \pm 0.16	5.10 \pm 0.13	0.213 \pm 0.005	0.187 \pm 0.003	0.178 \pm 0.003	0.155 \pm 0.004
	Deep Ensembles	5.95 \pm 0.14	4.79 \pm 0.11	4.57 \pm 0.07	4.39 \pm NA	0.226 \pm 0.001	0.158 \pm 0.002	0.148 \pm 0.001	0.134 \pm NA
	Single model	5.83 \pm 0.11	5.83 \pm 0.11	5.83 \pm 0.11	5.83 \pm 0.11	0.223 \pm 0.002	0.223 \pm 0.002	0.223 \pm 0.002	0.223 \pm 0.002
	Variational Inf. (FFG)	6.57 \pm 0.09	5.63 \pm 0.13	5.50 \pm 0.10	5.46 \pm 0.03	0.239 \pm 0.002	0.192 \pm 0.002	0.184 \pm 0.002	0.175 \pm 0.001
	KFAC-Laplace	6.00 \pm 0.13	5.82 \pm 0.12	5.82 \pm 0.19	5.80 \pm 0.19	0.210 \pm 0.005	0.203 \pm 0.007	0.201 \pm 0.007	0.200 \pm 0.008
WideResNet CIFAR-10	Snapshot Ensembles	7.76 \pm 0.22	5.52 \pm 0.13	5.00 \pm 0.10	4.54 \pm 0.05	0.247 \pm 0.005	0.176 \pm 0.001	0.160 \pm 0.001	0.137 \pm 0.001
	Dropout	3.88 \pm 0.12	3.70 \pm 0.18	3.63 \pm 0.19	3.64 \pm 0.17	0.130 \pm 0.002	0.120 \pm 0.002	0.119 \pm 0.001	0.117 \pm 0.002
	SWA-Gaussian	4.98 \pm 1.17	3.53 \pm 0.09	3.34 \pm 0.14	3.28 \pm 0.10	0.157 \pm 0.036	0.111 \pm 0.004	0.105 \pm 0.003	0.101 \pm 0.002
	Cyclic SGLD	4.78 \pm 0.16	4.09 \pm 0.11	3.63 \pm 0.13	3.19 \pm 0.04	0.155 \pm 0.003	0.128 \pm 0.002	0.114 \pm 0.001	0.099 \pm 0.002
	Fast Geometric Ens.	4.86 \pm 0.17	3.95 \pm 0.07	3.77 \pm 0.10	3.34 \pm 0.06	0.148 \pm 0.003	0.120 \pm 0.002	0.113 \pm 0.002	0.102 \pm 0.001
	Deep Ensembles	3.65 \pm 0.02	3.11 \pm 0.10	3.01 \pm 0.06	2.83 \pm NA	0.123 \pm 0.002	0.097 \pm 0.001	0.095 \pm 0.001	0.090 \pm NA
	Single model	3.70 \pm 0.15	3.70 \pm 0.15	3.70 \pm 0.15	3.70 \pm 0.15	0.124 \pm 0.005	0.124 \pm 0.005	0.125 \pm 0.005	0.124 \pm 0.005
	Variational Inf. (FFG)	5.61 \pm 0.04	4.15 \pm 0.15	3.94 \pm 0.10	3.64 \pm 0.07	0.189 \pm 0.002	0.134 \pm 0.002	0.127 \pm 0.002	0.117 \pm 0.001
	KFAC-Laplace	4.03 \pm 0.19	3.90 \pm 0.15	3.88 \pm 0.22	3.83 \pm 0.16	0.134 \pm 0.004	0.124 \pm 0.004	0.122 \pm 0.005	0.120 \pm 0.003
	Snapshot Ensembles	5.56 \pm 0.15	3.68 \pm 0.09	3.33 \pm 0.10	2.89 \pm 0.07	0.179 \pm 0.005	0.119 \pm 0.001	0.105 \pm 0.001	0.090 \pm 0.001

Table 3: Classification error and negative calibrated log-likelihood for different models and numbers of samples on CIFAR-10/100.

Model	Method	Error (%)				Negative calibrated log-likelihood			
		1	5	10	100	1	5	10	100
VGG16 CIFAR-100	Dropout	26.10 \pm 0.20	25.68 \pm 0.18	25.66 \pm 0.14	25.60 \pm 0.17	1.176 \pm 0.008	1.111 \pm 0.008	1.098 \pm 0.009	1.084 \pm 0.009
	SWA-Gaussian	27.74 \pm 1.87	24.53 \pm 0.09	23.64 \pm 0.28	22.97 \pm 0.20	1.109 \pm 0.073	0.931 \pm 0.007	0.879 \pm 0.007	0.826 \pm 0.005
	Cyclic SGLD	29.75 \pm 0.17	26.79 \pm 0.19	24.14 \pm 0.11	21.15 \pm 0.11	1.114 \pm 0.003	0.976 \pm 0.004	0.881 \pm 0.006	0.749 \pm 0.004
	Fast Geometric Ens.	27.07 \pm 0.24	25.35 \pm 0.29	24.68 \pm 0.40	22.78 \pm 0.22	1.057 \pm 0.010	0.965 \pm 0.003	0.930 \pm 0.003	0.827 \pm 0.004
	Deep Ensembles	25.72 \pm 0.17	21.60 \pm 0.13	20.79 \pm 0.16	19.88 \pm NA	1.092 \pm 0.004	0.840 \pm 0.005	0.794 \pm 0.002	0.723 \pm NA
	Single model	25.44 \pm 0.29	25.44 \pm 0.29	25.44 \pm 0.29	25.44 \pm 0.29	1.087 \pm 0.006	1.087 \pm 0.006	1.087 \pm 0.006	1.087 \pm 0.006
	Variational Inf. (FFG)	27.24 \pm 0.09	25.24 \pm 0.11	24.85 \pm 0.05	24.56 \pm 0.07	1.154 \pm 0.004	1.001 \pm 0.002	0.973 \pm 0.002	0.939 \pm 0.001
	KFAC-Laplace	27.11 \pm 0.59	25.98 \pm 0.21	25.84 \pm 0.38	25.70 \pm 0.38	1.174 \pm 0.037	1.089 \pm 0.007	1.069 \pm 0.005	1.050 \pm 0.008
WideResNet CIFAR-100	Snapshot Ensembles	31.19 \pm 0.33	23.87 \pm 0.18	22.31 \pm 0.31	21.03 \pm 0.10	1.170 \pm 0.012	0.899 \pm 0.004	0.834 \pm 0.005	0.751 \pm 0.003
	Dropout	20.19 \pm 0.11	19.41 \pm 0.17	19.36 \pm 0.12	19.22 \pm 0.15	0.823 \pm 0.008	0.768 \pm 0.005	0.760 \pm 0.006	0.751 \pm 0.005
	SWA-Gaussian	20.45 \pm 0.73	17.57 \pm 0.17	17.21 \pm 0.22	17.08 \pm 0.19	0.794 \pm 0.025	0.653 \pm 0.004	0.634 \pm 0.005	0.614 \pm 0.005
	Cyclic SGLD	21.42 \pm 0.32	19.42 \pm 0.28	17.88 \pm 0.16	16.29 \pm 0.10	0.813 \pm 0.010	0.713 \pm 0.009	0.654 \pm 0.005	0.583 \pm 0.004
	Fast Geometric Ens.	21.48 \pm 0.31	18.54 \pm 0.16	18.00 \pm 0.19	17.12 \pm 0.16	0.770 \pm 0.007	0.652 \pm 0.006	0.630 \pm 0.006	0.596 \pm 0.003
	Deep Ensembles	19.38 \pm 0.20	16.55 \pm 0.08	16.17 \pm 0.15	15.77 \pm NA	0.797 \pm 0.007	0.623 \pm 0.003	0.595 \pm 0.003	0.571 \pm NA
	Single model	19.31 \pm 0.24	19.31 \pm 0.24	19.31 \pm 0.24	19.31 \pm 0.24	0.797 \pm 0.010	0.797 \pm 0.010	0.797 \pm 0.010	0.797 \pm 0.010
	Variational Inf. (FFG)	24.38 \pm 0.27	20.17 \pm 0.15	19.28 \pm 0.09	18.74 \pm 0.08	1.004 \pm 0.011	0.767 \pm 0.004	0.727 \pm 0.003	0.685 \pm 0.002
	KFAC-Laplace	20.02 \pm 0.18	19.76 \pm 0.15	19.53 \pm 0.19	19.43 \pm 0.21	0.834 \pm 0.009	0.803 \pm 0.006	0.795 \pm 0.007	0.789 \pm 0.006
	Snapshot Ensembles	23.01 \pm 0.26	18.20 \pm 0.13	17.12 \pm 0.31	16.07 \pm 0.07	0.859 \pm 0.009	0.678 \pm 0.006	0.633 \pm 0.008	0.582 \pm 0.004

Table 3: Classification error and negative calibrated log-likelihood for different models and numbers of samples on CIFAR-10/100.

Выводы

- Ансамбли улучшают качество
- Ансамбли дают возможность оценить неопределенность ответа модели
- Ансамблирование активно развивается прямо сейчас