

Ансамблирование нейронных сетей

Бекян Артём, 201

Содержание

- 1) Введение
- 2) Что такое ансамбли
- 3) Зачем нужны ансамбли
- 4) Где используют ансамбли
- 5) Какие бывают ансамбли

Введение

Недостатки нейронных сетей

- Обучение нейронных сетей долгое
- Нет гарантии хорошего результата
- Высокая дисперсия

Что такое ансамбли

Что такое ансамбли: принцип работы

Очевидный вариант:

- 1) Обучим несколько моделей
- 2) Возьмем среднее от их предсказаний

Идея: объединение предсказаний нескольких моделей даст более хороший результат, чем предсказание одной модели

Что такое ансамбли: теорема Кондорсе

Формулировка:

Если каждый член жюри присяжных имеет независимое мнение, и если вероятность правильного решения члена жюри больше 0.5, то тогда вероятность правильного решения присяжных в целом возрастает с увеличением количества членов жюри, и стремится к единице.

Если же вероятность быть правым у каждого из членов жюри меньше 0.5, то вероятность принятия правильного решения присяжными в целом монотонно уменьшается и стремится к нулю с увеличением количества присяжных.

Что такое ансамбли: теорема Кондорсе

Доказательство:

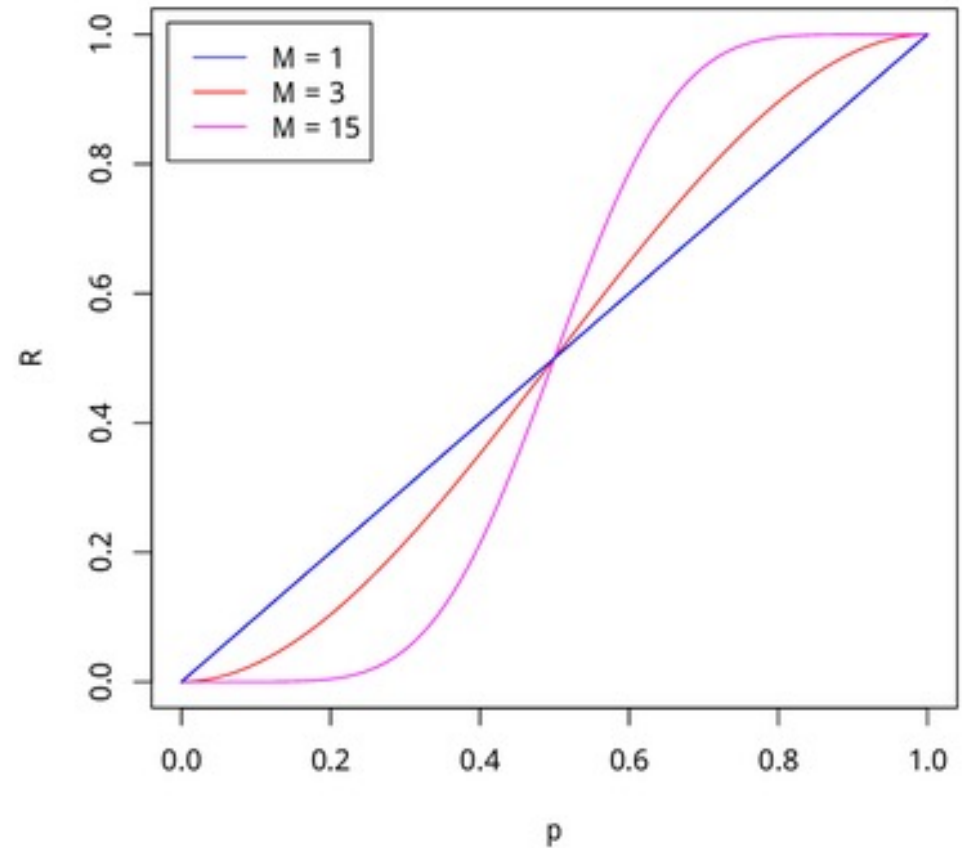
Пусть M – количество присяжных

p – вероятность правильного решения одного присяжного

R – вероятность правильного решения всех присяжных

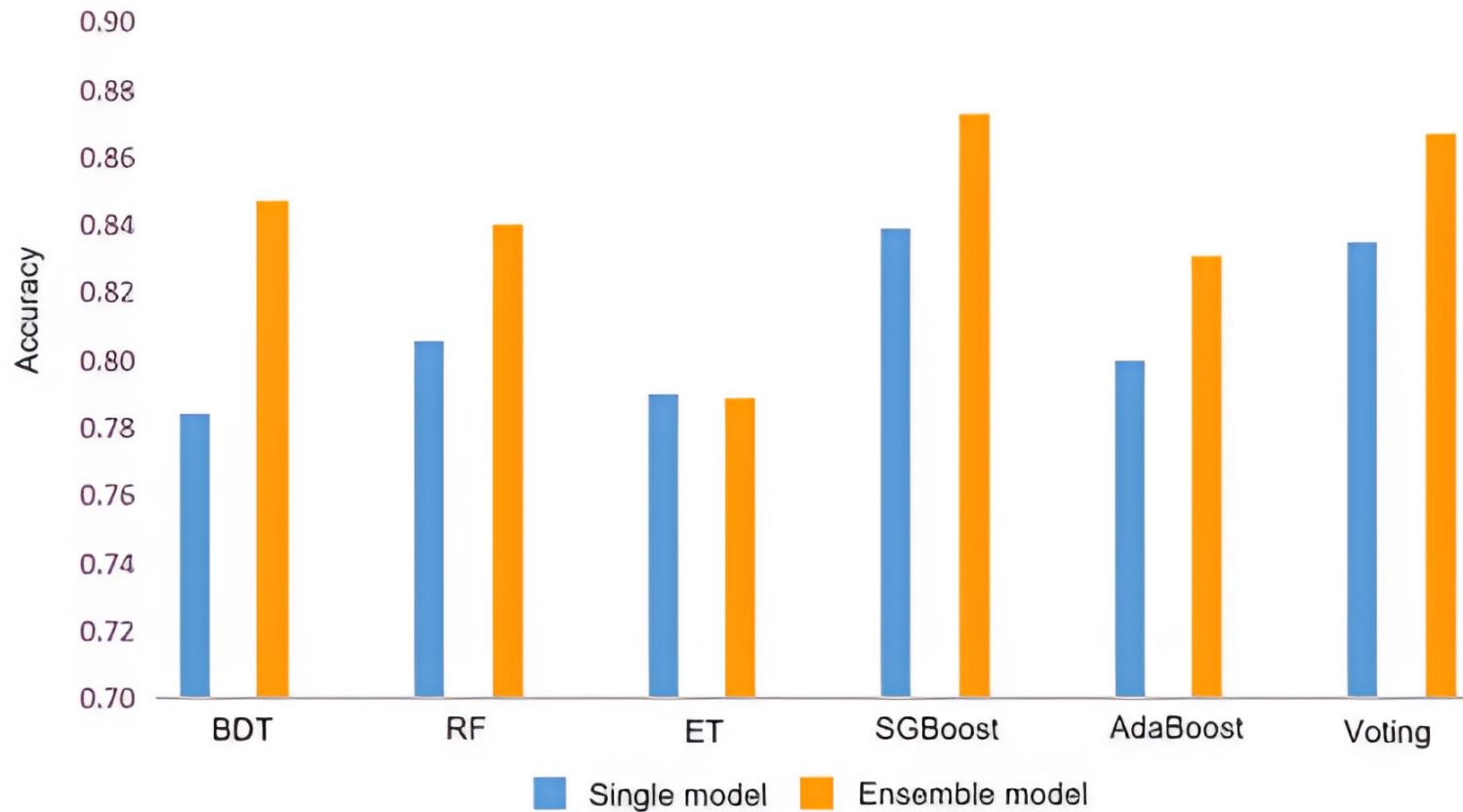
$m = \left\lfloor \frac{M}{2} \right\rfloor + 1$ – минимальное большинство членов жюри

Тогда $R = \sum_{i=m}^M C_M^i p^i (1-p)^{M-i}$



Зачем нужны ансамбли

Зачем нужны ансамбли: улучшение точности предсказаний



Accuracy comparison between single model and ensemble model. Where BDT, RF and ET represent Bagged Decision Tree, Random Forest and Extra Tree respectively

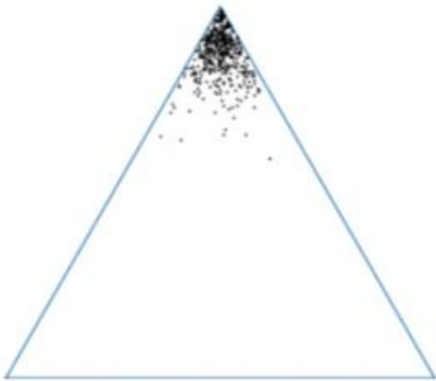
Зачем нужны ансамбли: подсчет неопределенности

Неопределенность – степень неуверенности модели в своем предсказании

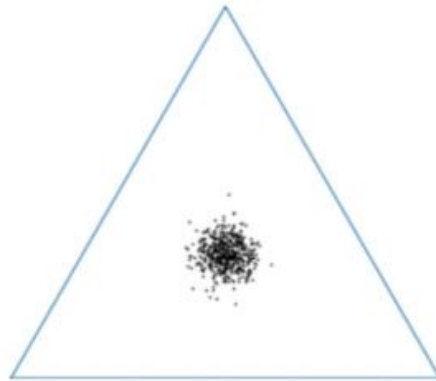
- **Неопределенность модели (эпистемическая)** – невозможность определить лучшую модель
- **Неопределенность данных (алеаторическая)** – присутствие шума в данных

Зачем нужны ансамбли: подсчет неопределенности

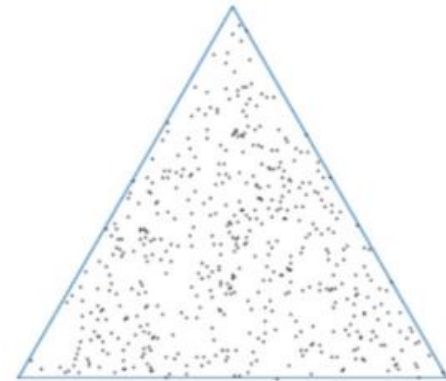
Ensemble $\{P(y|\mathbf{x}^*, \boldsymbol{\theta}^{(m)})\}_{m=1}^M$ can be visualized on a [simplex](#)



(a) Confident



(b) Data Uncertainty



(c) Knowledge Uncertainty

Где используют ансамбли

Где используют ансамбли

- Kaggle
- MNIST database of handwritten digits
- Защита от DoS-атак
- Обнаружение вредоносных программ
- ...
- И многое другое

Какие бывают ансамбли

Общие принципы

Какие бывают ансамбли: общий принципы.

Размер

- Универсальной верной формулы не существует
- Размер зависит от метода сборки ансамбля, использованных алгоритмов и т. д.
- Недавние исследования показывают, что может существовать оптимальный размер такой, что ансамбли с меньшим или большим размером имеют меньшую точность. Эти оптимальные значения близки к количеству признаков в выборке ([Bonab, Can, 2018](#))

Какие бывают ансамбли

Базовые алгоритмы

Какие бывают
ансамбли:
базовые
алгоритмы.
Бэггинг

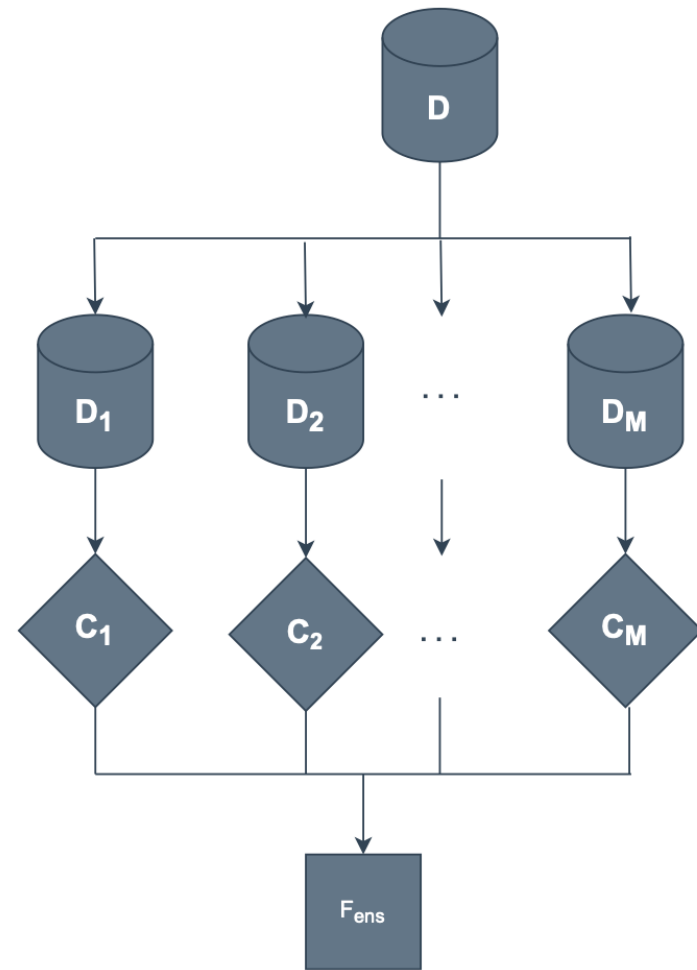


Figure 2: Bagging

Какие бывают
ансамбли:
базовые
алгоритмы.
Бустинг

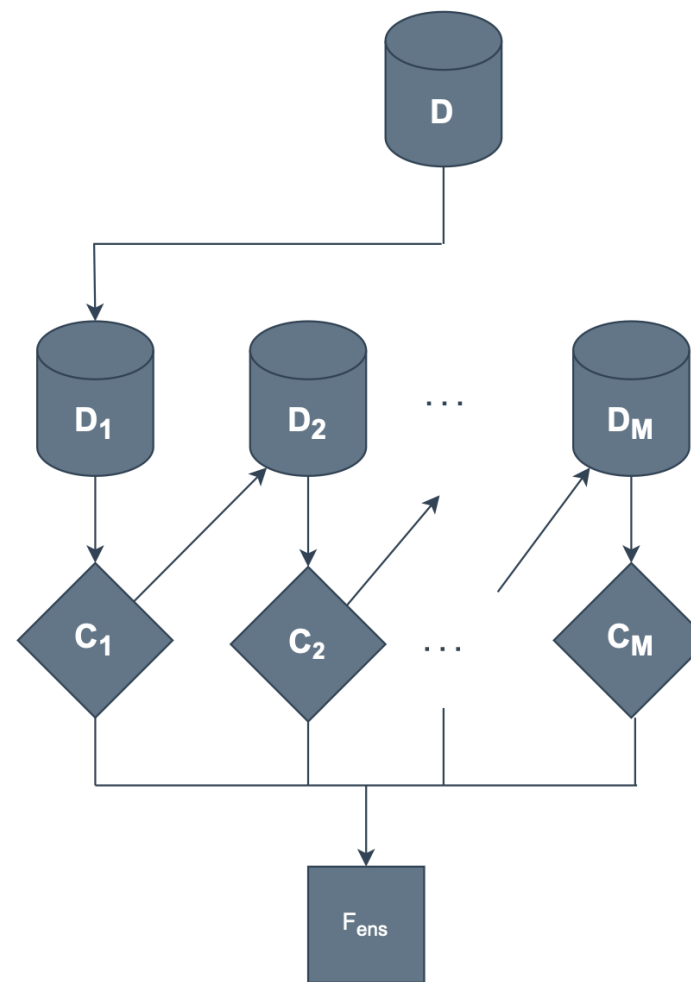


Figure 3: Boosting

Какие бывают
ансамбли:
базовые
алгоритмы.
Стэкинг

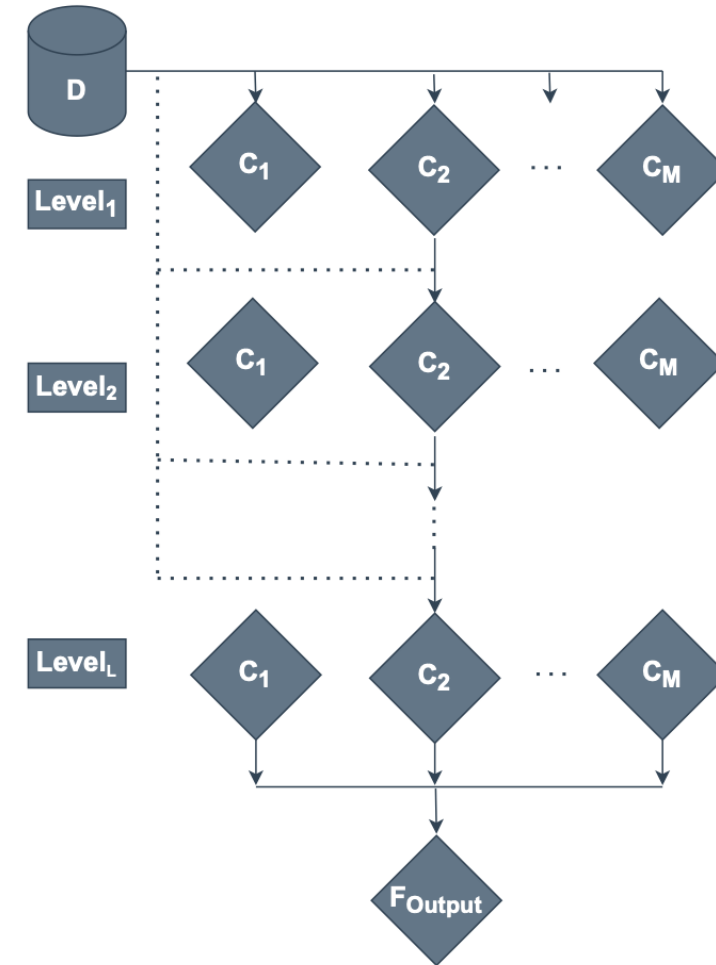


Figure 4: Stacking

Какие бывают ансамбли

Алгоритмы нейросетей

Какие бывают ансамбли: алгоритмы нейросетей. Деер ансамбли

Принцип работы:

- 1) Обучаем несколько независимых нейронных сетей на всем наборе данных
- 2) Для классификации считаем вероятность принадлежности к каждому из классов, для регрессии – мат. ожидание и квадрат дисперсии
- 3) Результат выбираем как среднее арифметическое предсказаний

Какие бывают ансамбли: алгоритмы нейросетей. Деер ансамбли

Преимущества:

- Лучшее соотношение размер ансамбля / качество на данный момент
- Одно из лучших значений точности

Недостатки:

- Дорогое обучение

Какие бывают ансамбли: алгоритмы нейросетей. Snapshot ансамбли

Принцип работы:

- 1) Обучаем одну нейронную сеть
- 2) В момент обучения с помощью SGD приходим в локальный минимум функции метрики качества на тренировочной выборке
- 3) Сохраняем модель и вручную ставим в этой точке большое значение функции метрики качества
- 4) Повторяем пункты 2 – 3 пока не получим нужное количество моделей
- 5) Результат выбираем через среднее арифметическое

Какие бывают ансамбли: алгоритмы нейросетей. Snapshot ансамбли

Преимущества:

- Сильное снижение затрат на обучения ансамбля
- Значительное улучшение точности предсказаний
- Одно из лучших соотношений размер ансамбля / качество

Недостатки:

- Зачастую проигрывает deep ансамблям в точности

Какие бывают ансамбли: алгоритмы нейросетей. Dropout ансамбли

Принцип работы:

- 1) Каждому нейрону присваивается p – вероятность того, что этот нейрон не выкинется
- 2) Проходим n раз, каждый раз каждый нейрон оставляем с вероятностью p или выкидываем с вероятностью $1 - p$
- 3) На тестовой выборке не выкидываем нейроны, а просто умножаем выход на p
- 4) Результатом является выход нейросети

Какие бывают ансамбли: алгоритмы нейросетей. Dropout ансамбли

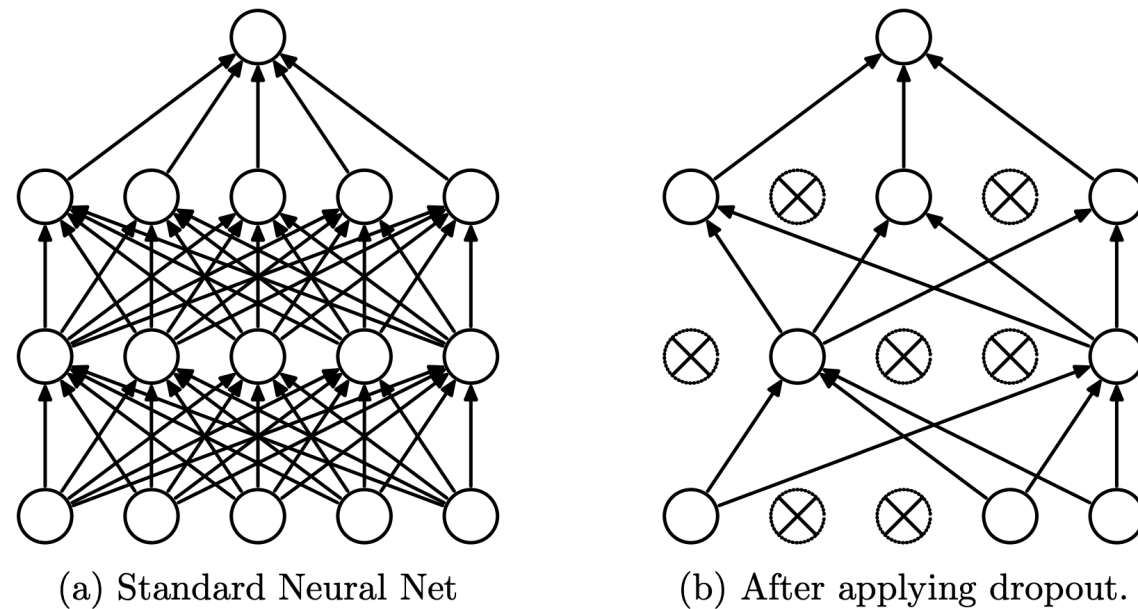


Figure 1: Dropout Neural Net Model. **Left:** A standard neural net with 2 hidden layers. **Right:** An example of a thinned net produced by applying dropout to the network on the left. Crossed units have been dropped.

Какие бывают ансамбли: алгоритмы нейросетей. Dropout ансамбли

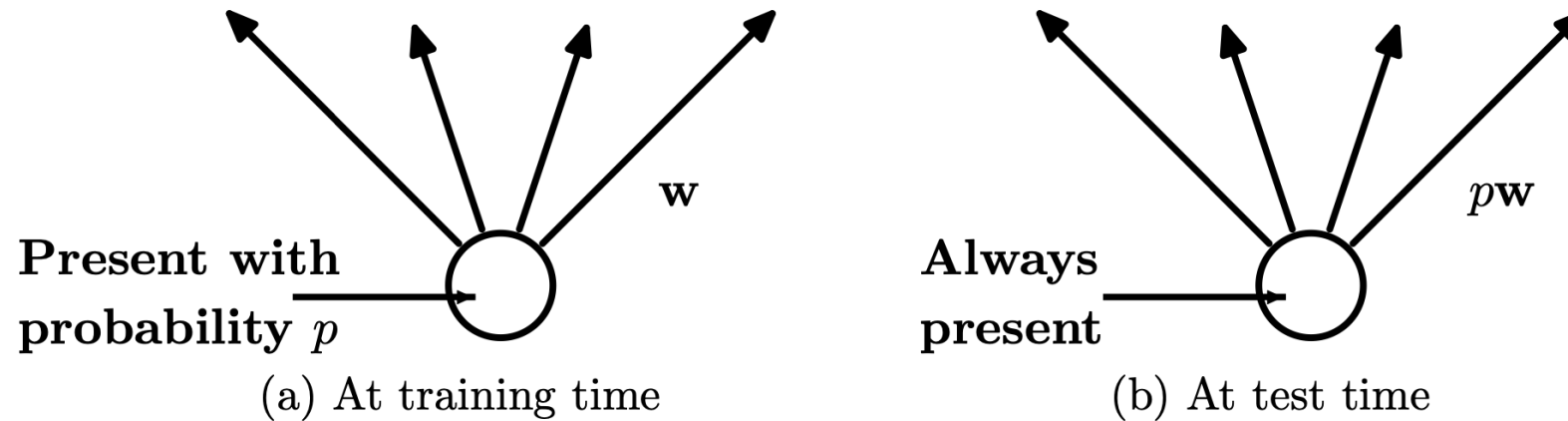


Figure 2: **Left:** A unit at training time that is present with probability p and is connected to units in the next layer with weights w . **Right:** At test time, the unit is always present and the weights are multiplied by p . The output at test time is same as the expected output at training time.

Какие бывают ансамбли: алгоритмы нейросетей. Dropout ансамбли

Преимущества:

- Скорость предсказаний
- По факту dropout – метод регуляризации, который может применяться в любых ансамблях

Недостатки:

- Увеличение размера ансамбля дает крайне малый прирост точности в сравнении с другими ансамблями

Статистика

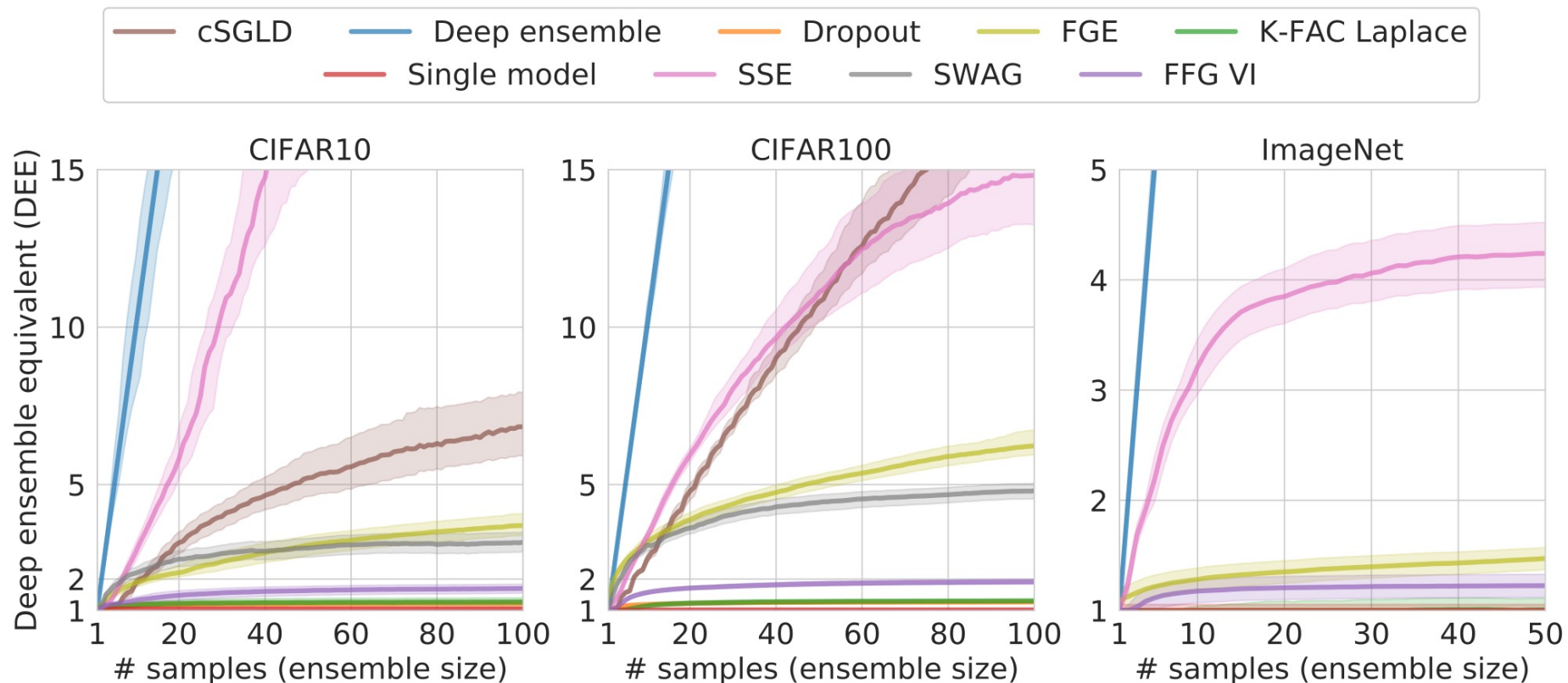
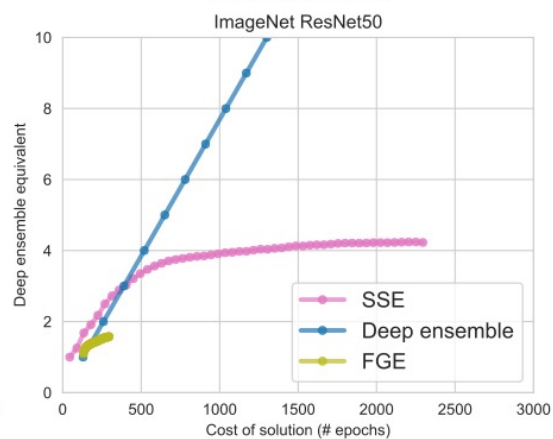
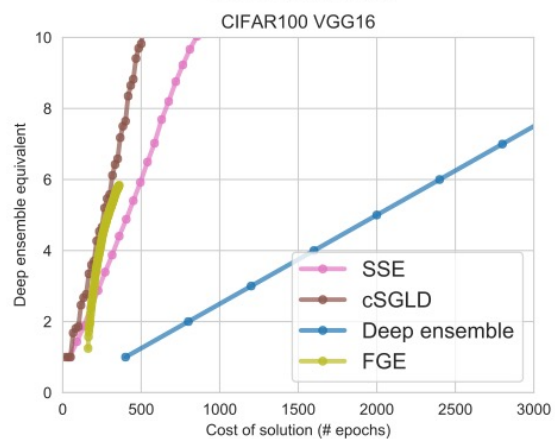
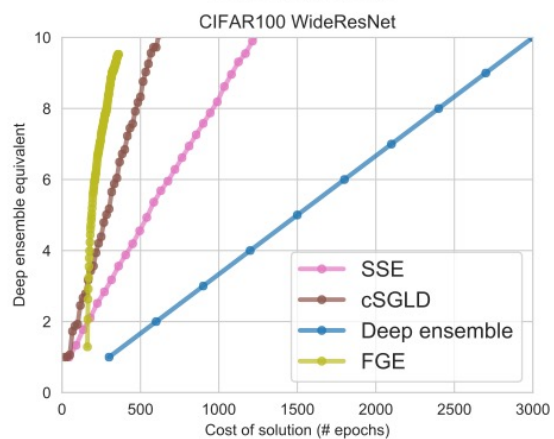
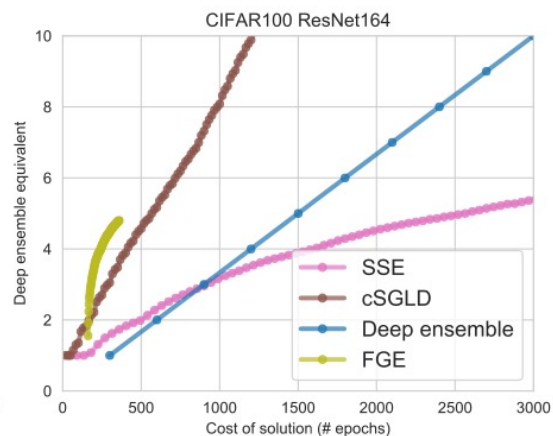
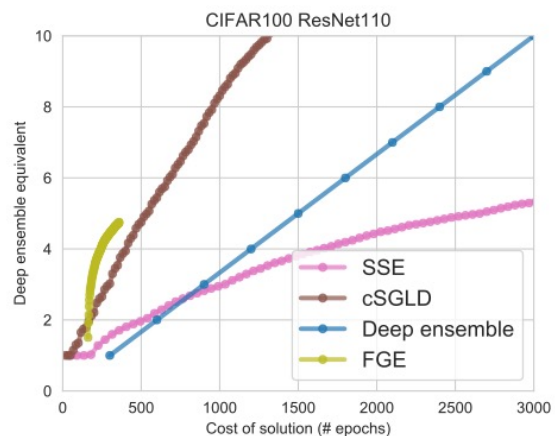
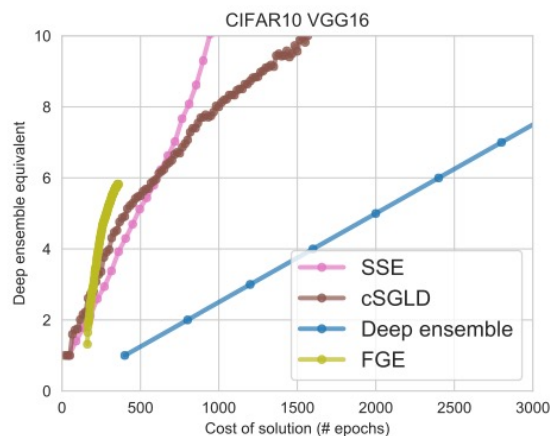
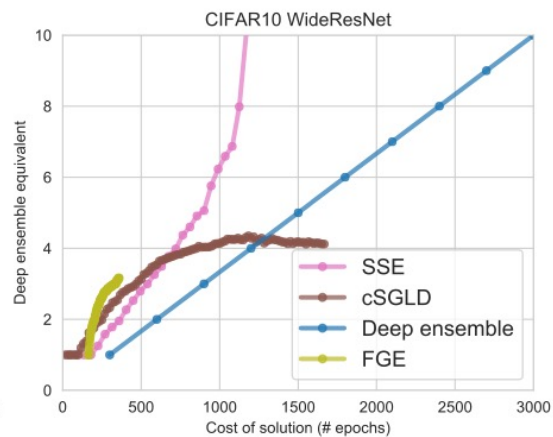
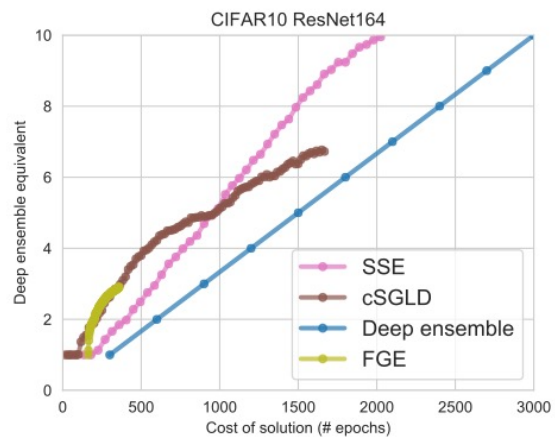
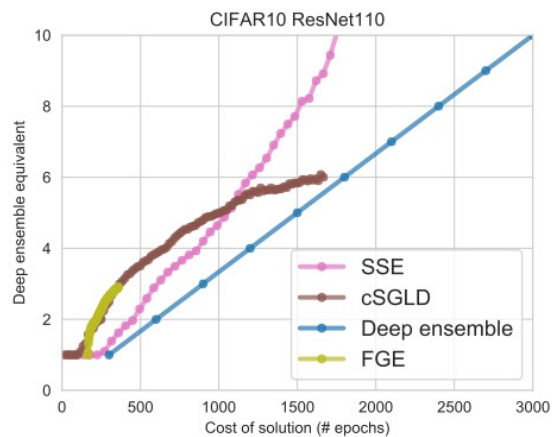


Figure 3: The deep ensemble equivalent score (DEE) for different numbers of samples on CIFAR-10, CIFAR-100, and ImageNet datasets averaged across different deep convolutional architectures. A deep ensemble equivalent score (DEE) of a model is equal to the minimum size of a deep ensemble (an ensemble of independently train networks) that achieves the same performance as the model under consideration. The score is measured in the number of models (higher is better). The area between average lower and upper bounds of DEE is shaded. **The plot demonstrates that all of the ensembling techniques are far less efficient than deep ensembles during inference and fail to produce the same level of performance as deep ensembles.** The comparison that is normalized on training time is presented in Appendix A.



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
VGG16 CIFAR-100	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
	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
WideResNet CIFAR-100	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
	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
WideResNet CIFAR-100	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.

Заключение

- Ансамбли зачастую долго обучаются
- Точность предсказаний ансамблей заметно выше, чем одиночных нейросетей
- С помощью ансамблей можно оценить неопределенность
- Тема ансамблирования нейронных сетей активно развивается