

Deep Ensembles

Куйбида Всеволод БПМИ212

Содержание

1. Введение
2. Что такое ансамбли?
3. Зачем нужны ансамбли?
4. Как тупо делать ансамбли нейросетей? Deep Ensembles
5. Как умно делать ансамбли нейросетей?
6. Насколько большими должны быть ансамбли?
7. Статистика, графики
8. Заключение

Введение

Введение

Глубинные нейронные сети - **мощный современный инструмент, минусы -**

1. Проблемы с обучением
2. Высокая дисперсия оценок
3. Нет гарантий, что мы сойдемся, куда ходим

Решение проблем (частичное) - **ансамблинг**

Что такое ансамбли?

Что такое ансамбли?

Ансамблинг - метод машинного обучения, когда множество моделей объединяются для предсказания какого-то одного таргета

Очевидный пример:

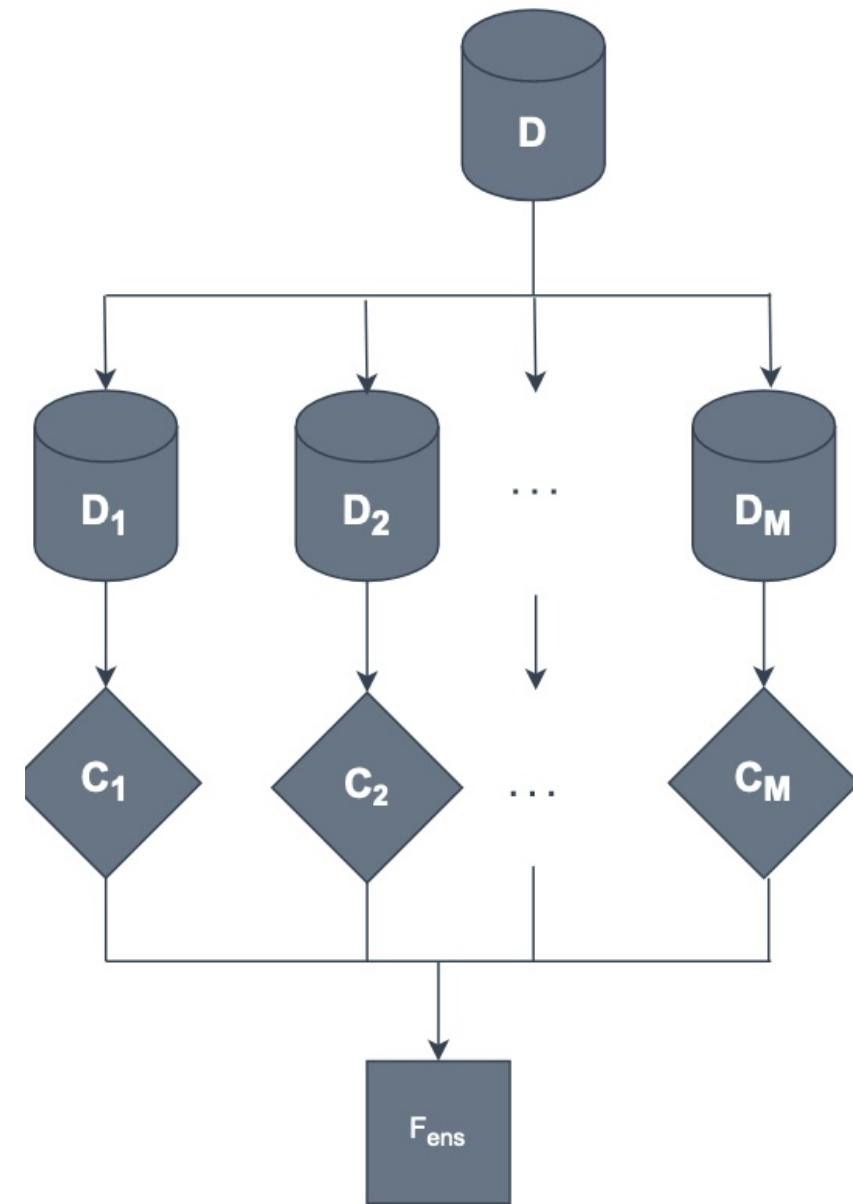
1. Обучим несколько моделей
2. В качестве ответа отдадим среднее арифметическое (или то, за что больше моделей "проголосуют")

Ансамбли с точки зрения статистики, примеры и виды

Известно, что ансамблирование **уменьшает** дисперсию/смещение в bias-variance decomposition или в их аналогах

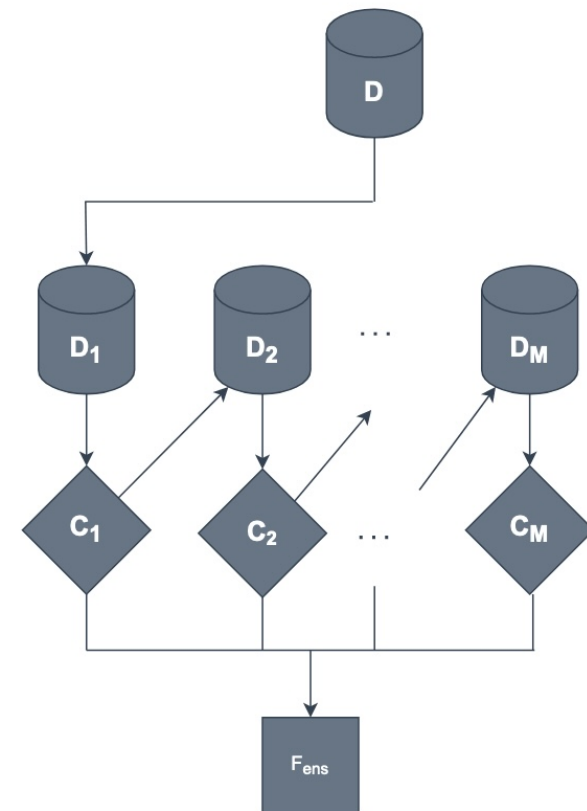
Традиционные виды ансамблей:

1. **Бэггинг** - уменьшает итоговую дисперсию, более робустный, чем модели внутри, используется в методе случайного леса



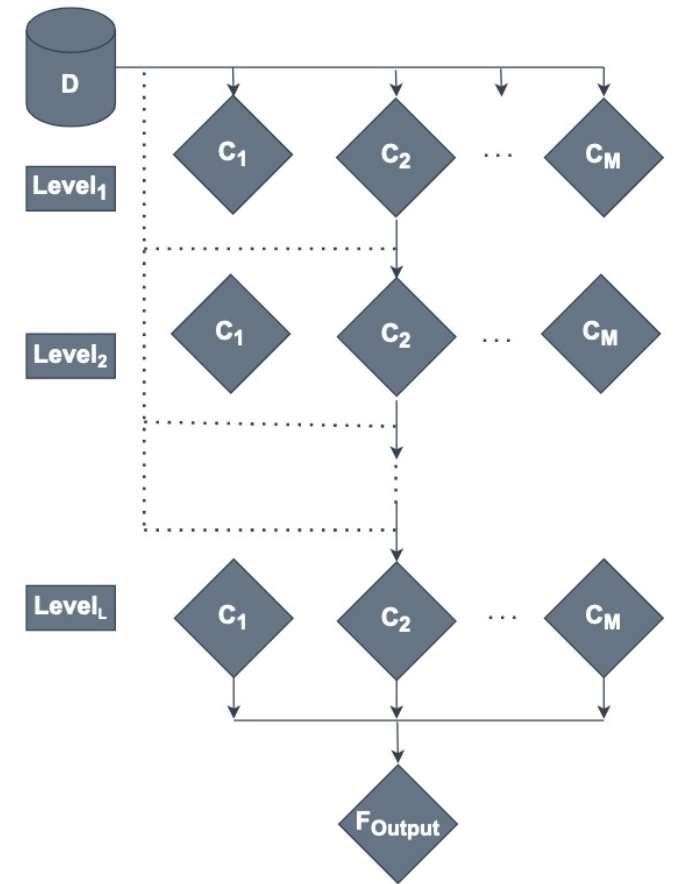
Традиционные виды ансамблей:

2. Бустинг - следующая модель улучшает предыдущую, используется в XGBoost, CatBoost, ...



Традиционные виды ансамблей:

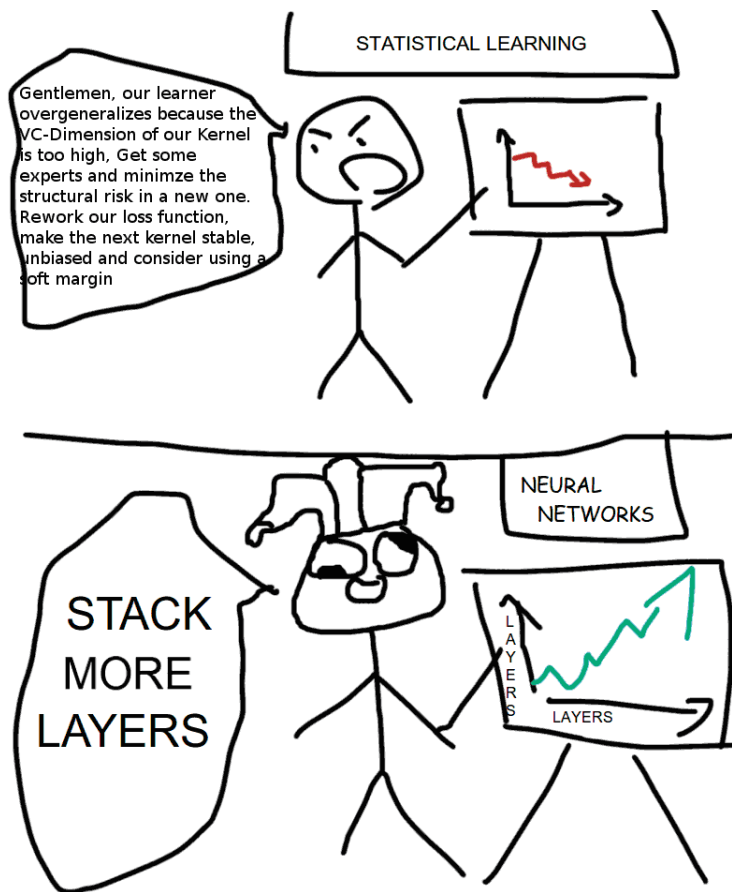
3. **Стэкинг** - ансамблирование разнообразных и так (чаще всего) сильных моделей, любимчик чемпионов Kaggle



Зачем нужны ансамбли?

Зачем нужны ансамбли?

Можем наstackать бесконечно много слоёв и радоваться жизни



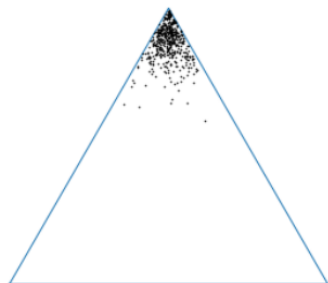
Зачем нужны ансамбли?

Введём понятие неопределенности:

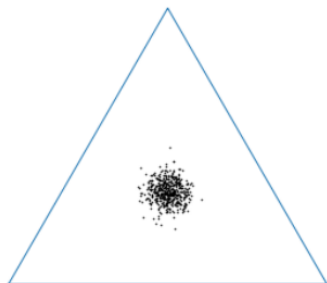
1. **Неопределенность данных** (алеаторическая) - присутствие шума в данных
2. **Неопределенность модели** (эпистемическая) - невозможность модели
восстановить истинный ответ

Визуализации неопределённостей

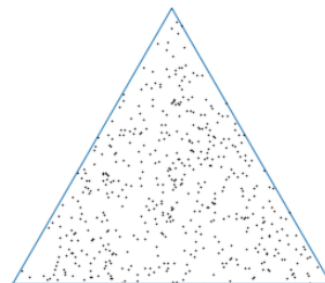
Ensemble $\{P(y|\mathbf{x}^*, \boldsymbol{\theta}^{(m)})\}_{m=1}^M$ can be visualized on a [simplex](#)



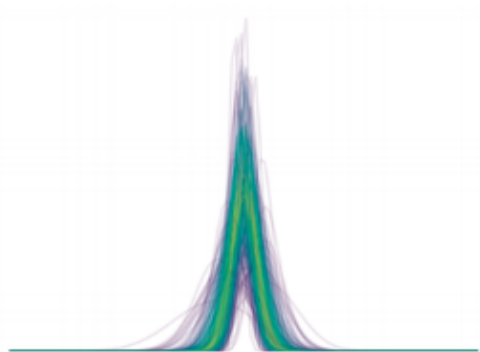
(a) Confident



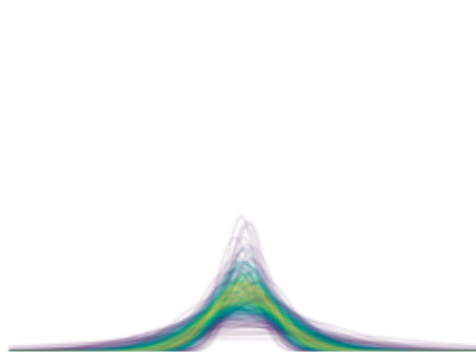
(b) Data Uncertainty



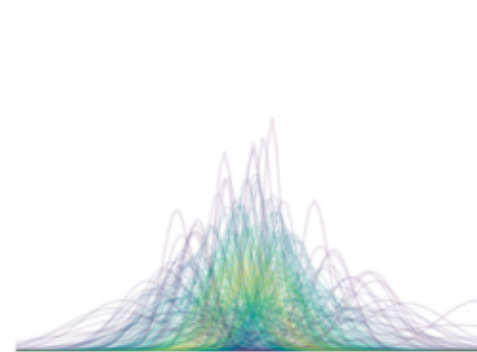
(c) Knowledge Uncertainty



(a) Low uncertainty



(b) Data uncertainty



(c) Knowledge Uncertainty

Зачем нужны ансамбли?

1. Уменьшаем эпистемическую неопределенность
2. Можем **параллельно** тренировать относительно неглубокие модели

Как тупо делать ансамбли нейросетей? Deep Ensembles

Как тупо делать ансамбли нейросетей? Deep Ensembles

Что делаем?

1. Тренируем параллельно M нейросетей на всем наборе данных
2. Для классификации считаем вероятность принадлежности к каждому из классов, для регрессии – мат. ожидание и дисперсию
3. Итоговый результат – среднее арифметическое ответов нейросетей

Как тупо делать ансамбли нейросетей? Deep Ensembles

Плюсы:

1. очень хорошее значение точности
2. лучшее отношение размера ансамбля к качеству в данный момент

Минусы - очень дорого учить

Deep Ensembles как практический инструмент

Применяют:

1. в соревновательном Data Science
2. в медицине
3. в антифроде
4. в компьютерном зрении
5. в работе с речью
6. ...

Deep Ensembles как теоретический инструмент: Simple and Scalable Predictive Uncertainty Estimation using Deep Ensembles

Рассмотрим статью Lakshminarayanan et al. 2016 ([ссылка](#)):

Проблема: хотим измерять неуверенность модели, SOTA подходы предполагают использование байесовских нейронных сетей со своими сложностями

Решение: брать ансамбли нейронных сетей, которые дают сравнимые результаты с меньшими издержками

Итоги: подход даёт сравнимые результаты, при этом обучение Deep Ensembles проще, требуется меньше гиперпараметров, и сам процесс обучения в разы проще реализовывать

Как умно делать ансамбли нейросетей?

Как умно делать ансамбли нейросетей?

Рассмотрим два способа делать ансамбли **умнее**, чем Deep Ensembles:

1. **Snapshot Ensembles**
2. **Dropout Ensembles**

Snapshot ensembles

Что делаем?

1. Обучаем одну нейронную сеть циклами обучения
2. В конце цикла обучения, когда мы “приходим” в локальный минимум, мы делаем “снэпшот” модели - сохраняем в другую модель веса
3. Повторяем шаги 1-2 до тех пор, пока не получим M моделей
4. Наш ответ на тесте - усреднённый ответ M моделей

Snapshot ensembles

Плюсы:

1. сильное уменьшение затрат на обучение ансамбля,
2. улучшение точности предсказания,
3. одно из лучших отношений качества к размеру модели

Минусы:

1. часто проигрывают Deep Ensembles в точности

Snapshot ensembles - визуализации

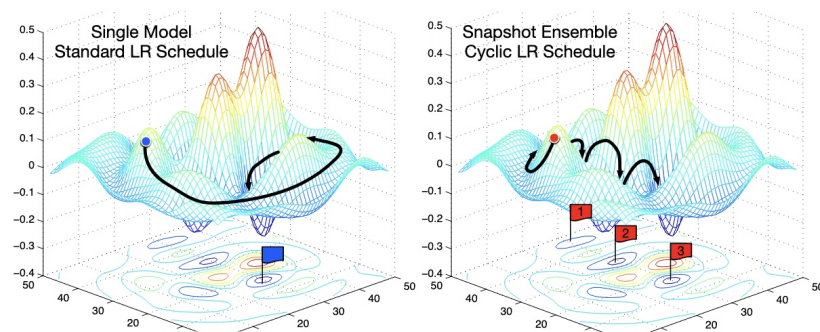


Figure 1: **Left:** Illustration of SGD optimization with a typical learning rate schedule. The model converges to a minimum at the end of training. **Right:** Illustration of Snapshot Ensembling. The model undergoes several learning rate annealing cycles, converging to and escaping from multiple local minima. We take a snapshot at each minimum for test-time ensembling.

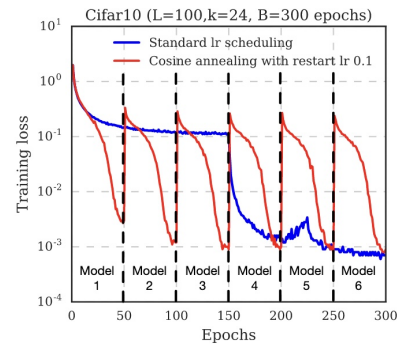


Figure 2: Training loss of 100-layer DenseNet on CIFAR10 using standard learning rate (blue) and $M = 6$ cosine annealing cycles (red). The intermediate models, denoted by the dotted lines, form an ensemble at the end of training.

Dropout ensembles

Что делаем (в самом простом случае - нейросети)?

1. Фиксируем p - вероятность того, что каждый нейрон будет "выключен" при обучении
2. При обучении на каждом шаге с вероятностью p "выключаем" нейрон (то есть не тренируем), при этом умножаем выход модели на $1 / (1 - p)$ для того, чтобы матожидание результата не поменялось
3. В "тестовом" режиме не выкидываем нейроны

Метод можно рассматривать как **простой и очень общий способ** регуляризации + **легко обобщается**

Dropout ensembles

Плюсы:

1. Метод предоставляет хороший, простой и очень общий способ регуляризации
2. Скорость предсказаний хороша
3. Легко и хорошо обобщается

Минусы:

1. требуется больше времени на обучение
2. увеличение размера ансамбля дает меньший прирост качества, нежели в других случаях

Насколько большими должны быть ансамбли?

Насколько большими должны быть ансамбли?

1. Универсальной верной формулы не существует
2. Размер зависит от метода ансамблирования, используемых алгоритмов и т. д.
3. В статье Bonab et al. 2018 ([ссылка](#)) рассказывается об оптимальном кол-ве моделей в ансамбле - авторы предлагают брать число, схожее с кол-вом классов в мультиклассовой классификации

Статистика, графики

Статистика, графики - Deep Ensemble Equivalent

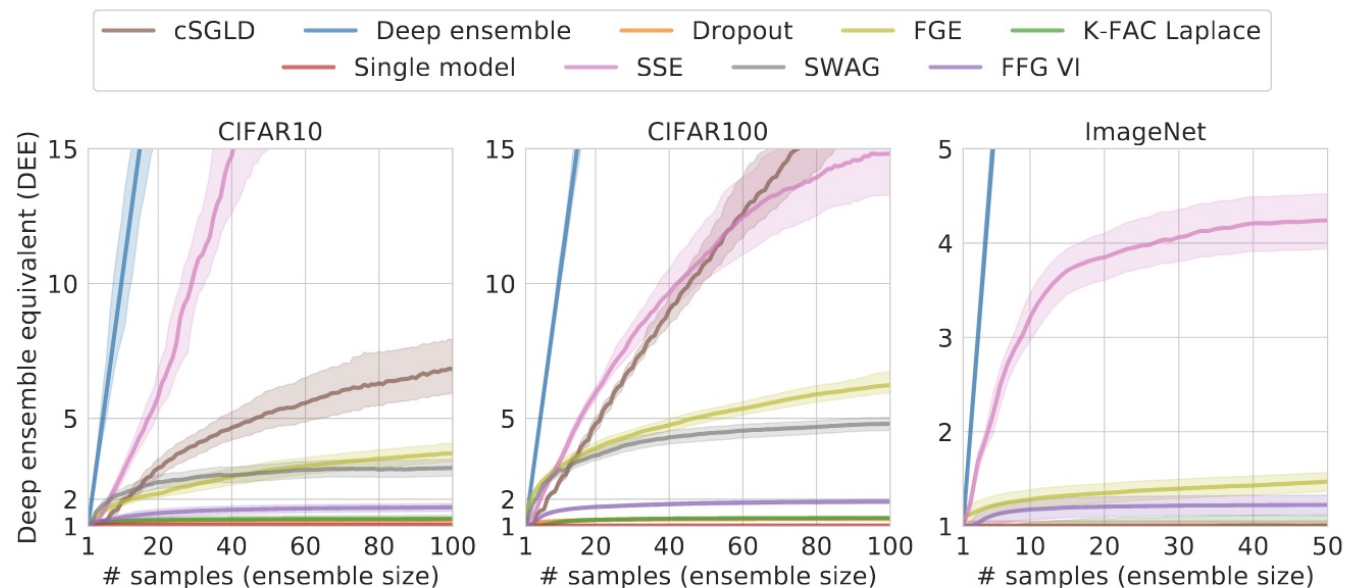
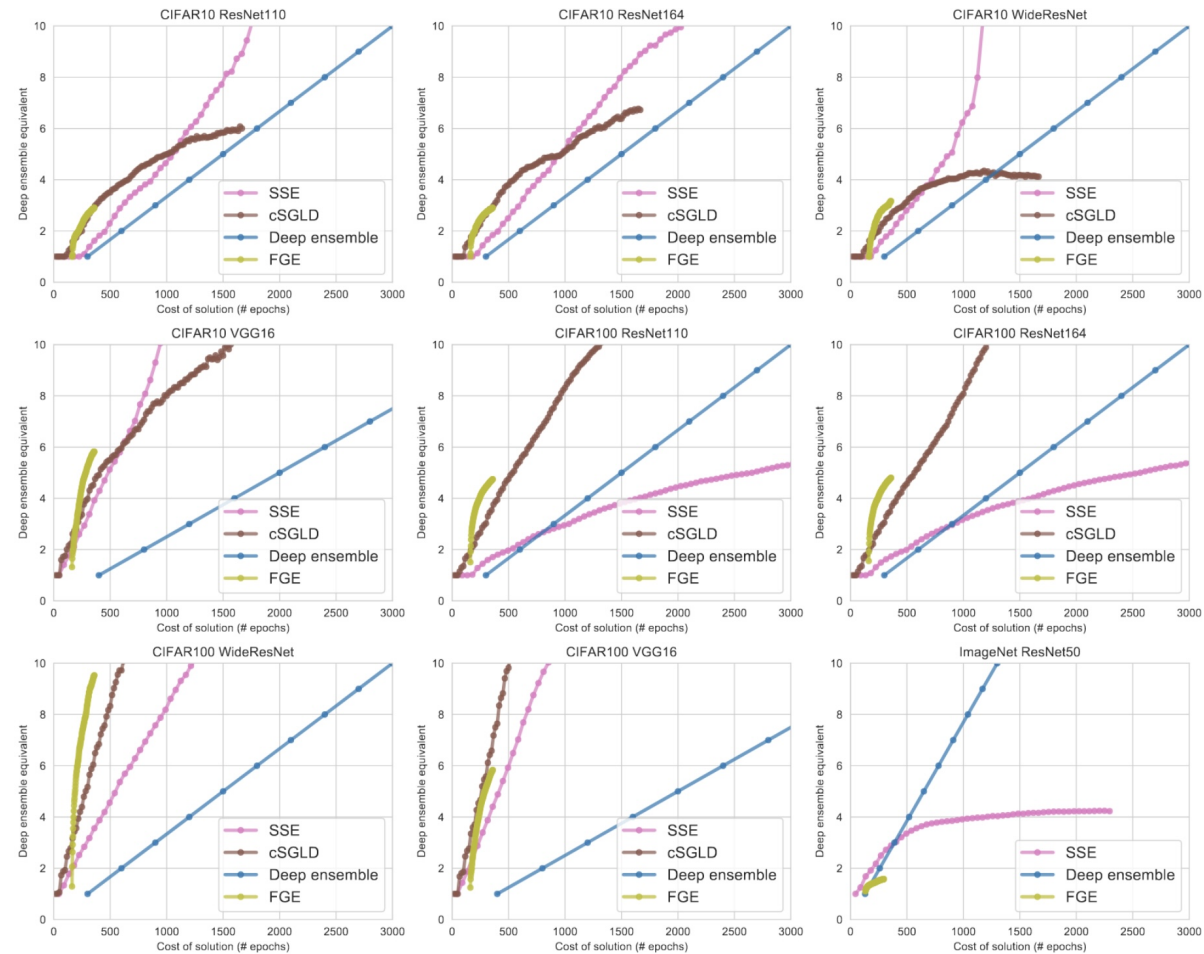


Figure 3: The deep ensemble equivalent score (DEE) for different numbers of samples on CIFAR-10, CIFAR-100, and ImageNet datasets averaged across different deep convolutional architectures. A deep ensemble equivalent score (DEE) of a model is equal to the minimum size of a deep ensemble (an ensemble 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.

Статистика, графики - Performance in DEE



Статистика, графики - сравнительные результаты моделей

Model	Method	Error (%)				Negative calibrated log-likelihood			
		1	5	10	100	1	5	10	100
VGG16 CIFAR-10	Dropout	5.86±0.09	5.81±0.08	5.82±0.06	5.79±0.07	0.232±0.005	0.225±0.004	0.224±0.004	0.223±0.003
	SWA-Gaussian	7.03±0.50	5.66±0.08	5.49±0.12	5.25±0.13	0.230±0.014	0.182±0.003	0.171±0.002	0.160±0.002
	Cyclic SGLD	7.37±0.16	6.56±0.09	5.71±0.06	4.84±0.04	0.234±0.004	0.196±0.004	0.176±0.003	0.147±0.003
	Fast Geometric Ens.	6.52±0.16	5.95±0.16	5.69±0.16	5.10±0.13	0.213±0.005	0.187±0.003	0.178±0.003	0.155±0.004
	Deep Ensembles	5.95±0.14	4.79±0.11	4.57±0.07	4.39±NA	0.226±0.001	0.158±0.002	0.148±0.001	0.134±NA
	Single model	5.83±0.11	5.83±0.11	5.83±0.11	5.83±0.11	0.223±0.002	0.223±0.002	0.223±0.002	0.223±0.002
	Variational Inf. (FFG)	6.57±0.09	5.63±0.13	5.50±0.10	5.46±0.03	0.239±0.002	0.192±0.002	0.184±0.002	0.175±0.001
	KFAC-Laplace	6.00±0.13	5.82±0.12	5.82±0.19	5.80±0.19	0.210±0.005	0.203±0.007	0.201±0.007	0.200±0.008
	Snapshot Ensembles	7.76±0.22	5.52±0.13	5.00±0.10	4.54±0.05	0.247±0.005	0.176±0.001	0.160±0.001	0.137±0.001
ResNet110 CIFAR-10	SWA-Gaussian	5.77±0.45	4.56±0.17	4.46±0.12	4.34±0.13	0.178±0.009	0.143±0.004	0.139±0.003	0.131±0.003
	Cyclic SGLD	6.18±0.20	5.32±0.15	4.55±0.13	3.83±0.02	0.185±0.006	0.156±0.005	0.138±0.002	0.115±0.001
	Fast Geometric Ens.	5.52±0.09	4.83±0.08	4.73±0.10	4.28±0.05	0.163±0.002	0.141±0.003	0.137±0.003	0.126±0.002
	Deep Ensembles	4.66±0.11	3.77±0.11	3.63±0.07	3.53±NA	0.148±0.004	0.117±0.002	0.112±0.002	0.106±NA
	Single model	4.69±0.11	4.69±0.11	4.69±0.11	4.69±0.11	0.150±0.002	0.150±0.002	0.150±0.003	0.150±0.002
	Variational Inf. (FFG)	5.57±0.26	4.91±0.15	4.72±0.13	4.60±0.03	0.178±0.003	0.149±0.001	0.144±0.001	0.140±0.000
	KFAC-Laplace	5.81±0.39	5.14±0.15	4.90±0.14	4.78±0.08	0.187±0.014	0.160±0.007	0.153±0.005	0.147±0.003
	Snapshot Ensembles	8.41±0.27	4.85±0.11	4.16±0.16	3.52±0.10	0.252±0.006	0.153±0.002	0.132±0.002	0.107±0.001
ResNet164 CIFAR-10	SWA-Gaussian	5.41±0.71	4.21±0.19	4.21±0.23	4.02±0.14	0.171±0.028	0.130±0.004	0.128±0.004	0.121±0.002
	Cyclic SGLD	5.80±0.21	4.97±0.12	4.30±0.08	3.66±0.06	0.178±0.004	0.149±0.004	0.131±0.003	0.110±0.001
	Fast Geometric Ens.	5.22±0.07	4.49±0.06	4.36±0.07	4.09±0.12	0.157±0.003	0.134±0.002	0.130±0.001	0.119±0.002
	Deep Ensembles	4.53±0.11	3.51±0.09	3.50±0.06	3.34±NA	0.147±0.002	0.113±0.001	0.107±0.001	0.100±NA
	Single model	4.52±0.11	4.52±0.11	4.52±0.11	4.52±0.11	0.144±0.002	0.144±0.003	0.144±0.002	0.144±0.003
	Variational Inf. (FFG)	5.62±0.14	4.78±0.05	4.66±0.05	4.55±0.08	0.183±0.004	0.151±0.001	0.146±0.001	0.141±0.001
	KFAC-Laplace	5.23±0.29	4.77±0.23	4.65±0.17	4.60±0.09	0.168±0.008	0.151±0.007	0.146±0.005	0.142±0.004
	Snapshot Ensembles	8.06±0.10	4.50±0.04	3.89±0.09	3.50±0.05	0.241±0.004	0.144±0.003	0.124±0.002	0.104±0.001
WideResNet CIFAR-10	Dropout	3.88±0.12	3.70±0.18	3.63±0.19	3.64±0.17	0.130±0.002	0.120±0.002	0.119±0.001	0.117±0.002
	SWA-Gaussian	4.98±1.17	3.53±0.09	3.34±0.14	3.28±0.10	0.157±0.036	0.111±0.004	0.105±0.003	0.101±0.002
	Cyclic SGLD	4.78±0.16	4.09±0.11	3.63±0.13	3.19±0.04	0.155±0.003	0.128±0.002	0.114±0.001	0.099±0.002
	Fast Geometric Ens.	4.86±0.17	3.95±0.07	3.77±0.10	3.34±0.06	0.148±0.003	0.120±0.002	0.113±0.002	0.102±0.001
	Deep Ensembles	3.65±0.02	3.11±0.10	3.01±0.06	2.83±NA	0.123±0.002	0.097±0.001	0.095±0.001	0.090±NA
	Single model	3.70±0.15	3.70±0.15	3.70±0.15	3.70±0.15	0.124±0.005	0.124±0.005	0.125±0.005	0.124±0.005
	Variational Inf. (FFG)	5.61±0.04	4.15±0.15	3.94±0.10	3.64±0.07	0.189±0.002	0.134±0.002	0.127±0.002	0.117±0.001
	KFAC-Laplace	4.03±0.19	3.90±0.15	3.88±0.22	3.83±0.16	0.134±0.004	0.124±0.004	0.122±0.005	0.120±0.003
	Snapshot Ensembles	5.56±0.15	3.68±0.09	3.33±0.10	2.89±0.07	0.179±0.005	0.119±0.001	0.105±0.001	0.090±0.001

Заключение

Заключение

- Ансамблирование нейросетей – **хороший инструмент**, который **активно развивается**
- За свои плюсы часто приходится платить **дороговизной** обучения
- Имеет как **практические** и **инженерные** приложения, так и более **теоретические** (uncertainty estimation)