

A statistical perspective on distillation

Aditya Krishna Menon



Ankit Singh Rawat



Seungyeon Kim



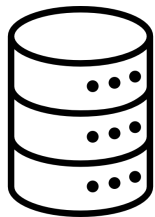
Sashank Reddi



Sanjiv Kumar



Supervised learning: standard



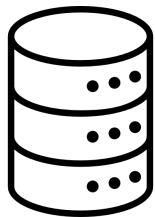
Labelled data

(x_1, y_1)

\vdots

(x_n, y_n)

Supervised learning: standard



Labelled data

(x_1, y_1)
 \vdots
 (x_n, y_n)



Learner

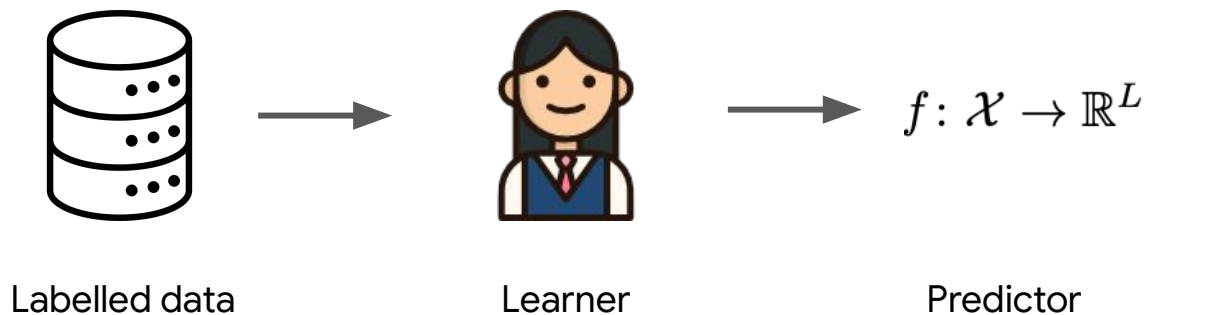
$$\min_f \frac{1}{N} \sum_{n \in [N]} e_{y_n}^\top \ell(f(x_n))$$

Loss vector

One-hot encoding of labels

$[\ell(1, f(x_n)), \dots, \ell(L, f(x_n))]$

Supervised learning: standard



(x_1, y_1)
 \vdots
 (x_n, y_n)

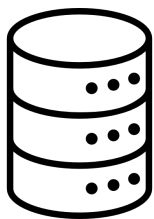
$$\min_f \frac{1}{N} \sum_{n \in [N]} e_{y_n}^\top \ell(f(x_n))$$

Loss vector

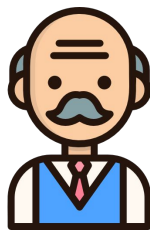
One-hot encoding of labels

$[\ell(1, f(x_n)), \dots, \ell(L, f(x_n))]$

Supervised learning: distillation



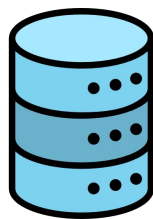
Labelled data

$$\begin{pmatrix} (x_1, y_1) \\ \vdots \\ (x_n, y_n) \end{pmatrix}$$


“Teacher”
model

$$\min_f \frac{1}{N} \sum_{n \in [N]} e_{y_n}^\top \ell(p^\dagger(x_n))$$

Distribution
over labels



Teacher-labelled data

$$\begin{pmatrix} (x_1, p^\dagger(x_1)) \\ \vdots \\ (x_n, p^\dagger(x_n)) \end{pmatrix}$$


“Student”
model

$$\min_f \frac{1}{N} \sum_{n \in [N]} p^\dagger(x_n)^\top \ell(f(x_n))$$

Replaces
one-hot labels

Successes of distillation

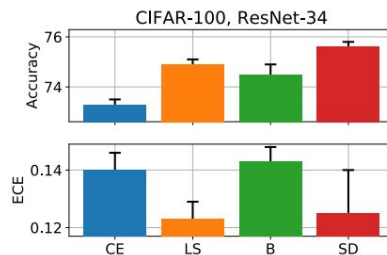
Empirically, distillation has demonstrated wide success:

System	Test Frame Accuracy
Baseline	58.9%
10xEnsemble	61.1%
Distilled Single model	60.8%

Hinton et al., 2015

Network	Teacher	BAN
DenseNet-112-33	18.25	16.95
DenseNet-90-60	17.69	16.69
DenseNet-80-80	17.16	16.36
DenseNet-80-120	16.87	16.00

Furlanello et al., 2018



Zhang et al., 2020

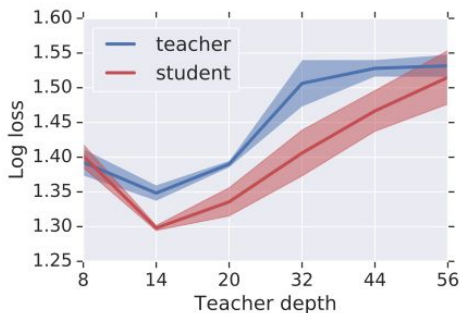
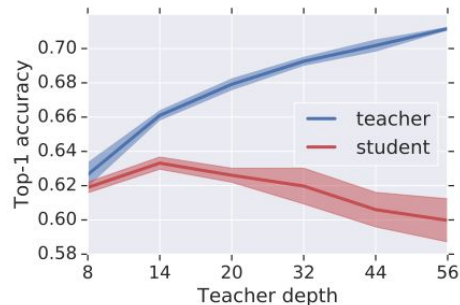
But *why* does distillation help?

Summary of our work

Q: **Why does distillation help?**

A: We provide a **statistical perspective**:

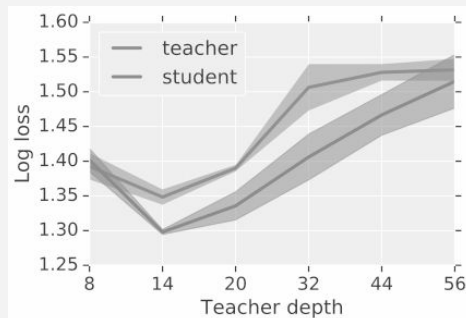
- Teacher approximates **Bayes probabilities**
- Exact Bayes probabilities → **reduce variance** of objective
- Approximate Bayes probabilities → **bias-variance tradeoff**



Distilling from a Bayes teacher



Distilling from an imperfect teacher



Applications

$$\widehat{\text{PD}}(f) \propto \sum_{i \in S, j \in S - \{i\}} p^t(x_i) \cdot (1 - p^t(x_j)) \cdot \mathbb{I}[f(x_i) < f(x_j)].$$

$$z(x_n) \doteq \log \left(\sum_{y' \in [L]} \alpha_{y'}(x_n) \cdot e^{f_{y'}(x_n)} \right)$$

$$\alpha_{y'}(x_n) \doteq 1 - \mathbb{I}[y' \neq y_n] \cdot p^t(y' | x_n).$$

Statistical learning setup

Suppose our training samples $(x, y) \sim \mathbb{P}$

Underlying “Bayes” distribution

Student goal: minimise the population **risk**, i.e., expected loss:

$$R(f) = \mathbb{E}_x \mathbb{E}_{y|x} [\ell(y, f(x))]$$

Statistical learning setup

Suppose our training samples $(x, y) \sim \mathbb{P}$

Underlying “Bayes” distribution

Student goal: minimise the population **risk**, i.e., expected loss:

$$\begin{aligned} R(f) &= \mathbb{E}_x \mathbb{E}_{y|x} [\ell(y, f(x))] \\ &= \mathbb{E}_x [p^*(x)^\top \ell(f(x))] \end{aligned}$$

 $(P^*(y = 1 | x), \dots, P^*(y = L | x))$

Inherently smooths loss by Bayes-probabilities!

“Bayes teacher” distillation



“Bayes-distilled” training loss:

$$\hat{R}_*(f) = \frac{1}{N} \sum_{n \in [N]} p^*(x_n)^\top \ell(f(x_n))$$

→ Predictions from a “Bayes teacher”

Like standard empirical loss, $\mathbb{E}[\hat{R}_*(f)] = R(f)$

But has an important advantage...

Why does “Bayes distillation” help?

Bayes-distilled loss **lowers variance** over empirical loss:

$$\mathbb{V}_{S \sim \mathbb{P}^N} [\hat{R}_*(f)] \leq \mathbb{V}_{S \sim \mathbb{P}^N} [\hat{R}(f)]$$

→ Variance over draws of training set

Lower variance → better generalisation bound:

$$R(\mathbf{f}) \leq \hat{R}_*(\mathbf{f}; S) + \mathcal{O}\left(\sqrt{\mathbb{V}_N^*(\mathbf{f})/N} \cdot \sqrt{\log(\mathcal{M}_N^*/\delta)} + \log(\mathcal{M}_N^*/\delta)/N\right),$$

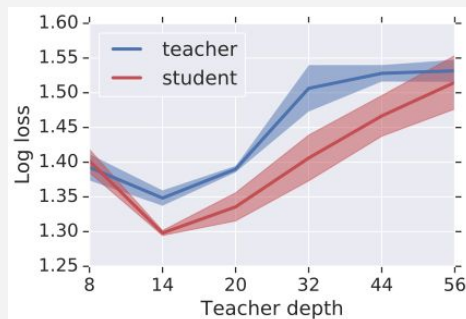
“Bayes distillation” can improve generalisation!

See paper
for more!

Distilling from a Bayes teacher



Distilling from an imperfect teacher



Applications

$$\widehat{\text{PD}}(f) \propto \sum_{i \in S, j \in S - \{i\}} p^t(x_i) \cdot (1 - p^t(x_j)) \cdot \mathbb{I}[f(x_i) < f(x_j)].$$

$$z(x_n) \doteq \log \left(\sum_{y' \in [L]} \alpha_{y'}(x_n) \cdot e^{f_{y'}(x_n)} \right)$$
$$\alpha_{y'}(x_n) \doteq 1 - \mathbb{I}[y' \neq y_n] \cdot p^t(y' | x_n).$$

Distilling from imperfect teacher

“Bayes teacher” helps; but what about other teachers?

Better approximation of p^* \rightarrow better generalisation:

$$\begin{aligned} \mathbb{E}[(\tilde{R}(\mathbf{f}; S) - R(\mathbf{f}))^2] &\leq \frac{1}{N} \cdot \mathbb{V} [\mathbf{p}^t(x)^\top \ell(\mathbf{f}(x))] \\ &+ \mathcal{O}(\mathbb{E} \|\mathbf{p}^t(x) - \mathbf{p}^*(x)\|_2^2) \end{aligned}$$

Distilling from imperfect teacher

“Bayes teacher” helps; but what about other teachers?

Better approximation of p^* \rightarrow better generalisation:

$$\mathbb{E}[(\tilde{R}(\mathbf{f}; S) - R(\mathbf{f}))^2] \leq \frac{1}{N} \cdot \mathbb{V} [\mathbf{p}^t(x)^\top \ell(\mathbf{f}(x))] \\ + \mathcal{O}\left(\|\mathbb{E}[\mathbf{p}^t(x)] - \mathbf{p}^*(x)\|_2^2 + \mathbb{V}[\mathbf{p}^t(x)]\right).$$

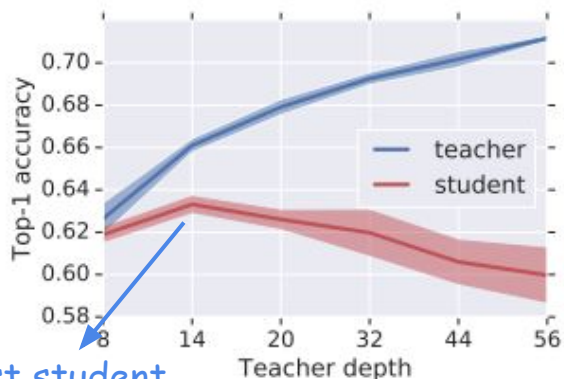
 **Bias-variance tradeoff**
for modelling p^*

Implications: bias-variance bound

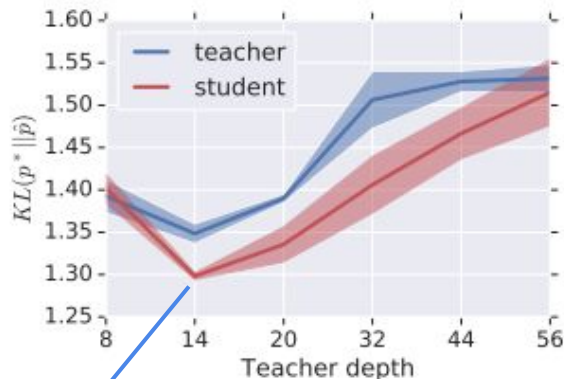
Bound is **not** on teacher **accuracy**

- Teacher can be accurate but poorly calibrated
- cf. finding that accurate teachers may distill poorer [Muller et al., 2019]

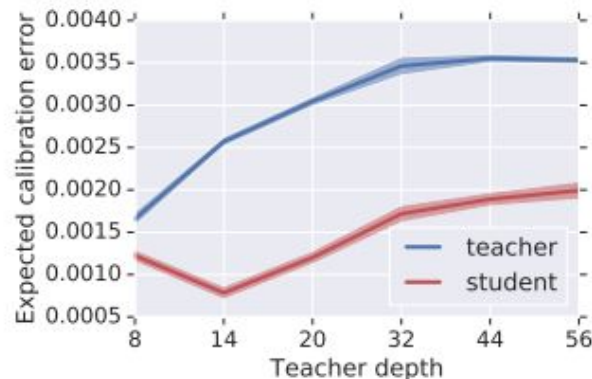
$$\mathbb{E}[(\tilde{R}(\mathbf{f}; S) - R(\mathbf{f}))^2] \leq \frac{1}{N} \cdot \mathbb{V}[\mathbf{p}^t(x)^\top \ell(\mathbf{f}(x))] + \mathcal{O}\left(\|\mathbb{E}[\mathbf{p}^t(x)] - \mathbf{p}^*(x)\|_2^2 + \mathbb{V}[\mathbf{p}^t(x)]\right).$$



(a) Top-1 accuracy.



(b) Log-loss.



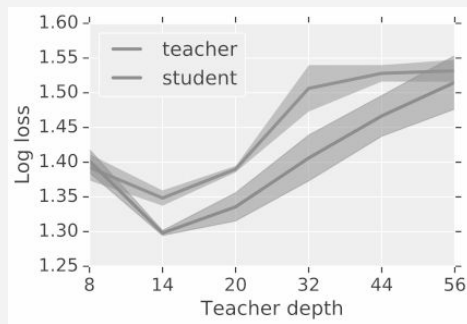
(c) Expected calibration error.

See paper for more!

Distilling from a Bayes teacher



Distilling from an imperfect teacher



Applications

$$\widehat{\text{PD}}(f) \propto \sum_{i \in S, j \in S - \{i\}} p^t(x_i) \cdot (1 - p^t(x_j)) \cdot \mathbb{I}[f(x_i) < f(x_j)].$$

$$z(x_n) \doteq \log \left(\sum_{y' \in [L]} \alpha_{y'}(x_n) \cdot e^{f_{y'}(x_n)} \right)$$
$$\alpha_{y'}(x_n) \doteq 1 - \mathbb{I}[y' \neq y_n] \cdot p^t(y' | x_n).$$

Applications of statistical view

Use teacher's estimates in place of Bayes probabilities $p^*(x)$

Example: bipartite ranking, where goal is to minimise

$$\text{PD}(f) = \mathbb{P}_{x|y=+1} \mathbb{P}_{x|y=-1} (f(x) < f(x'))$$

Distilled bipartite risk:

$$\widetilde{\text{PD}}(f) \propto \sum_{i \in S, j \in S - \{i\}} p^t(x_i) \cdot (1 - p^t(x_j)) \cdot \llbracket f(x_i) < f(x_j) \rrbracket.$$

Additional
weighting on
"negatives"

See paper
for more!

Summary of our work

Q: **Why does distillation help?**

A: We provide a **statistical perspective**:

- Teacher approximates **Bayes probabilities**
- Exact Bayes probabilities → **reduce variance** of objective
- Approximate Bayes probabilities → **bias-variance tradeoff**

See paper
for more!

