# Across the Great Divide: from ML Theory to Practice

**Aditya Krishna Menon**

Google NYC
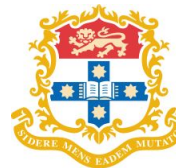
Google Research

# Introduction

Research Scientist at Google NYC

Working on machine learning algorithm design and analysis

Past lives:

- USyd
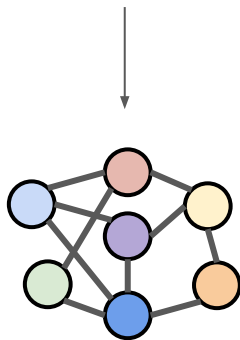
- UCSD

- NICTA/CSIRO Data61/ANU

# Supervised learning in theory

Training data



Model training



Model predictions



| 0.2 | 0.1 | 0.4 | 0.2 | 0.1 |
|-----|-----|-----|-----|-----|

# Supervised learning in theory

Training data



$$\{(x_n, y_n)\}_{n=1}^N$$

Model training



$$\min_{f \in \mathcal{F}} \frac{1}{N} \sum_{n \in N} \ell(y_n, f(x_n))$$

Model predictions



| 0.2 | 0.1 | 0.4 | 0.2 | 0.1 |

$$f(x^*)$$

# Supervised learning in practice

Training data

$$\{(x_n, y_n)\}_{n=1}^{N}$$

Model training

What if the model size is **too large**?

$$\min_{f \in \mathcal{F}} \frac{1}{N} \sum_{n \in N} \ell(y_n, f(x_n))$$

Model predictions

| 0.2 | 0.1 | 0.4 | 0.2 | 0.1 |

$$f(x^*)$$
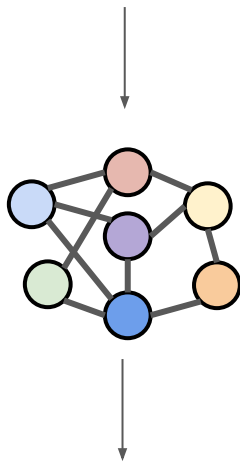
# Supervised learning in practice

Training data

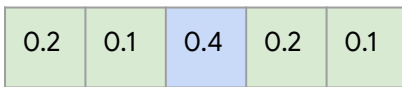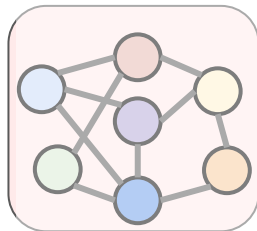$$\{(x_n, y_n)\}_{n=1}^N$$

Model training

$$\min_{f \in \mathcal{F}} \frac{1}{N} \sum_{n \in N} \boxed{\ell(y_n, f(x_n))}$$

What if this loss is **expensive** to compute?

Model predictions

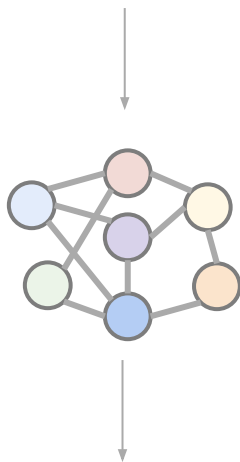| 0.2 | 0.1 | 0.4 | 0.2 | 0.1 |

$$f(x^*)$$

# Supervised learning in practice
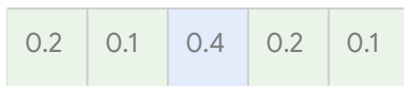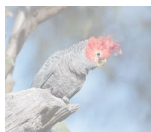
Training data

$$\{(x_n, y_n)\}_{n=1}^N$$

Model training

$$\min_{f \in \mathcal{F}} \frac{1}{N} \sum_{n \in N} \ell(y_n, f(x_n))$$

What if this operation is **stochastic**?

Model predictions

| 0.2 | 0.1 | 0.4 | 0.2 | 0.1 |

$$f(x^*)$$

# Agenda

Google Research

**01**

# Background

# Neural networks for classification

"Cockatoo"

# Neural networks for classification

Embedding layers

x

Cockatoo Parrot Pigeon Magpie

Embedding vector

Class weight matrix

| Cockatoo | Parrot | Pigeon | Magpie |
|---|---|---|---|
| 5.0 | 2.2 | -1.1 | 0.3 |

Output logits

# Neural networks for classification

Google Research

Training objective: minimise **softmax cross-entropy**

$$\log \sum \exp$$

| | Cockatoo | Parrot | Pigeon | Magpie | | |
|---|---|---|---|---|---|---|
| | 5.0 | 2.2 | -1.1 | 0.3 | − | 5.0 |

This approximately minimises the (negative) **prediction margin**:

$$\max$$

| | Parrot | Pigeon | Magpie | | |
|---|---|---|---|---|---|
| | 2.2 | -1.1 | 0.3 | − | 5.0 |

Highest score of
"wrong" label

# Neural networks for classification

Training objective: minimise **softmax cross-entropy**

$$\log \sum \exp \quad \boxed{\overset{\textit{Cockatoo}}{5.0} \; \overset{\textit{Parrot}}{2.2} \; \overset{\textit{Pigeon}}{-1.1} \; \overset{\textit{Magpie}}{0.3}} \quad - \quad \boxed{5.0}$$

This equivalently minimises the **KL divergence**:

$$\text{KL}(\mathbf{e}_y \,\|\, \mathbf{p}(x))$$

$$p_i(x) = \frac{\exp(s_i(x))}{\sum_{j \in [L]} \exp(s_j(x))}$$

"One-hot" label vector $\quad \boxed{1.0 \;\; 0.0 \;\; 0.0 \;\; \ldots}$

$\boxed{0.3 \;\; 0.5 \;\; 0.1 \;\; \ldots} \quad$ "Softmax" probability vector

**02**

# Distillation

# Supervised learning in theory

Training data



$$\{(x_n, y_n)\}_{n=1}^{N}$$

Model training



$$\min_{f \in \mathcal{F}} \frac{1}{N} \sum_{n \in N} \ell(y_n, f(x_n))$$

Model predictions



| 0.2 | 0.1 | 0.4 | 0.2 | 0.1 |

$$f(x^*)$$

# Supervised learning in practice

Training data



$$\{(x_n, y_n)\}_{n=1}^N$$

Model training

What if the model
size is **too large**?



$$\min_{f \in \mathcal{F}} \frac{1}{N} \sum_{n \in N} \ell(y_n, f(x_n))$$

Model predictions



| 0.2 | 0.1 | 0.4 | 0.2 | 0.1 |

$$f(x^*)$$

# Why        increase model size?

😃 Can **work better**!

    Particularly for complex tasks, e.g.,
    language modelling

Belkin et al., '19. Reconciling modern machine learning practice and the bias-variance trade-off.
Kaplan et al. '20. Scaling Laws for Neural Language Models.

# Why (not) increase model size?

😀 Can **work better**!

　　Particularly for complex tasks, e.g.,
　　language modelling

😞 More expensive to **train**

😣 More expensive to **predict**

Belkin et al., '19. Reconciling modern machine learning practice and the bias-variance trade-off.
Kaplan et al. '20. Scaling Laws for Neural Language Models.

# Idea: model compression

Ideally, **compress** our model while **preserving** performance



Many options: quantisation, architecture optimisation, distillation,....

# Distillation in a nutshell

Train a "student" model using **soft predictions** from "teacher" model



Predictions from "teacher"

| 0.3 | 0.5 | 0.1 | ... |

| 0.0 | 0.0 | 0.9 | ... |

| 0.2 | 0.4 | 0.4 | ... |

| 0.1 | 0.1 | 0.1 | ... |

| 0.5 | 0.1 | 0.2 | ... |

Trained on one-hot labels

Trained on **teacher predictions**

# Distillation loss function

Minimise                              softmax cross-entropy

$$\log \sum \exp \quad \boxed{\begin{array}{c} \text{Cockatoo} \\ 5.0 \end{array}} \; \boxed{\begin{array}{c} \text{Parrot} \\ 2.2 \end{array}} \; \boxed{\begin{array}{c} \text{Pigeon} \\ -1.1 \end{array}} \; \boxed{\begin{array}{c} \text{Magpie} \\ 0.3 \end{array}} \quad - \quad \boxed{5.0}$$

# Distillation loss function

Minimise **teacher-weighted** softmax cross-entropy

$$p^t(\ ) \quad \times \quad \log \sum \exp$$

| Cockatoo | Parrot | Pigeon | Magpie |
|:---:|:---:|:---:|:---:|
| 5.0 | 2.2 | -1.1 | 0.3 |

$-$ | 5.0 | $+$

$$p^t(\ ) \quad \times \quad \log \sum \exp$$

| 5.0 | 2.2 | -1.1 | 0.3 |
|:---:|:---:|:---:|:---:|

$-$ | 2.2 | $+$

$$p^t(\ ) \quad \times \quad \log \sum \exp$$

| 5.0 | 2.2 | -1.1 | 0.3 |
|:---:|:---:|:---:|:---:|

$-$ | -1.1 | $+$

$$p^t(\ ) \quad \times \quad \log \sum \exp$$

| 5.0 | 2.2 | -1.1 | 0.3 |
|:---:|:---:|:---:|:---:|

$-$ | 0.3

# Distillation loss function: formally

Suppose the teacher's predictions are $p^t$

Then, we may minimise:


Input data

$$\frac{1}{N}\sum_{n=1}^{N}\left[(1-\alpha)\cdot\mathrm{KL}(\mathbf{e}_{y_n}\,\|\,p(x_n)) + \alpha\cdot\mathrm{KL}(p^{\mathrm{t}}(x_n)\,\|\,p(x_n))\right]$$

Mixing weight

| 1.0 | 0.0 | 0.0 | ... |
|-----|-----|-----|-----|

"One-hot" label

| 0.3 | 0.5 | 0.1 | ... |
|-----|-----|-----|-----|

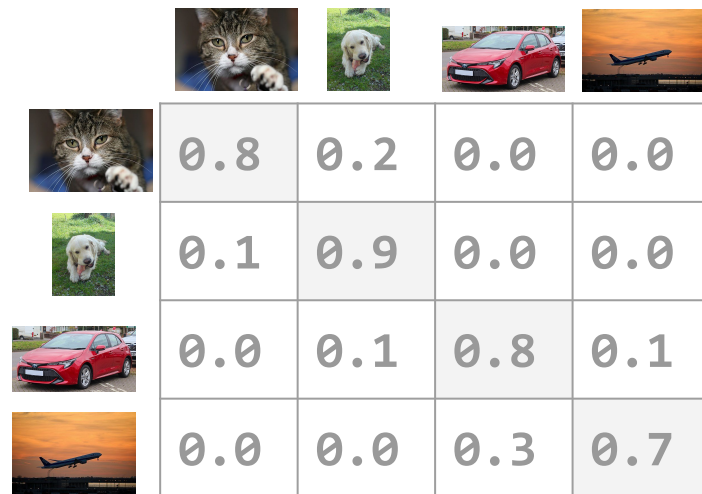"Soft" label

# Why does distillation help?

Transfers **class relationship** information
"Dark knowledge"
Learns which errors to penalise more

Per-sample **label smoothing**
Prevents over-confident predictions

Can be used on **unlabelled samples**
Form of semi-supervised learning!



|  | cat | dog | car | plane |
|---|---|---|---|---|
| cat | 0.8 | 0.2 | 0.0 | 0.0 |
| dog | 0.1 | 0.9 | 0.0 | 0.0 |
| car | 0.0 | 0.1 | 0.8 | 0.1 |
| plane | 0.0 | 0.0 | 0.3 | 0.7 |

# Beyond probability matching

Can match **more structure** in teacher model

e.g., match embeddings, pairwise similarities, …



(a) Teacher and Student Networks    (b) Hints Training

Romero et al., '15. FitNets: hints for thin deep nets.

# Do we need complex teachers?

**No.** You can "self-distill" **(!)**

Can give non-trivial gains



Step 0 → Step 1 → Step K → Ensemble

Why does this help?

Mostly an active area of research

One view: sample-dependent regularisation

Furlanello et al., '18. Born-again neural networks.

**03**

# Extreme classification

# Supervised learning in theory
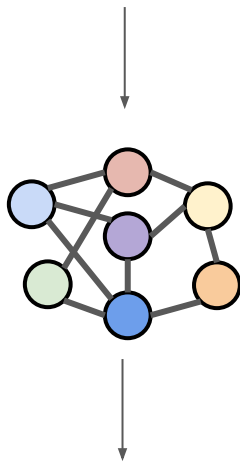
Training data



$$\{(x_n, y_n)\}_{n=1}^N$$

Model training



$$\min_{f \in \mathcal{F}} \frac{1}{N} \sum_{n \in N} \ell(y_n, f(x_n))$$

Model predictions



| 0.2 | 0.1 | 0.4 | 0.2 | 0.1 |

$$f(x^*)$$

# Supervised learning in practice

Training data



$$\{(x_n, y_n)\}_{n=1}^{N}$$

Model training



$$\min_{f \in \mathcal{F}} \frac{1}{N} \sum_{n \in N} \boxed{\ell(y_n, f(x_n))}$$

What if this loss is **expensive** to compute?
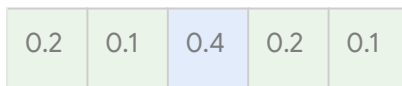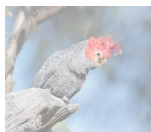
Model predictions



| 0.2 | 0.1 | 0.4 | 0.2 | 0.1 |
|-----|-----|-----|-----|-----|

$$f(x^*)$$

# Neural networks for　classification

# Neural networks for **extreme** classification

Embedding
vector

Class weight
matrix

# Neural networks for **extreme** classification

Google Research

Training objective: minimise **softmax cross-entropy**



$$\log \sum \exp$$

| | Cockatoo | Parrot | Pigeon | Magpie | Sparrow | Crow | Albatross | | Bird of Paradise |
|---|---|---|---|---|---|---|---|---|---|
| | 5.0 | 2.2 | -1.1 | 0.3 | 0.1 | -4.8 | -1.9 | ... | 0.5 |

\-   5.0

**Hard to compute** even for a single sample!

# Negative sampling

Select a subset of "**negative**" labels to contrast against "**positive**"

**"Positive" label**

**"Negative" labels**

# Negative sampling

Select a subset of "**negative**" labels to contrast against "**positive**"

**"Positive" label**

**"Negative" labels**



Ideally, we would like the sampling to:

- Be **easy** to compute

- Result in **informative** negatives

# Choosing the sampling distribution

**Solution #1**: within-batch negatives



"Positive" label

"Negative" labels

$(x_1, \quad)$

$(x_2, \quad)$

$(x_3, \quad)$

😃 Easy to compute

😞 Biased towards frequent labels

# Choosing the sampling distribution

**Solution #2**: uniform random negatives

"Positive" label

"Negative" labels

$(x_1, $  $)$

$(x_2, $  $)$

$(x_3, $  $)$

😃 Easy to compute

😃 Not biased towards any label

😞 May not be informative

# Choosing the sampling distribution

**Solution #3**: hard negative mining

"Positive" label

"Negative" labels

$(x_1,$  $)$

$(x_2,$  $)$

$(x_3,$  $)$

😃 Maximally informative

😞 Hard to compute

# Finding hard-negatives

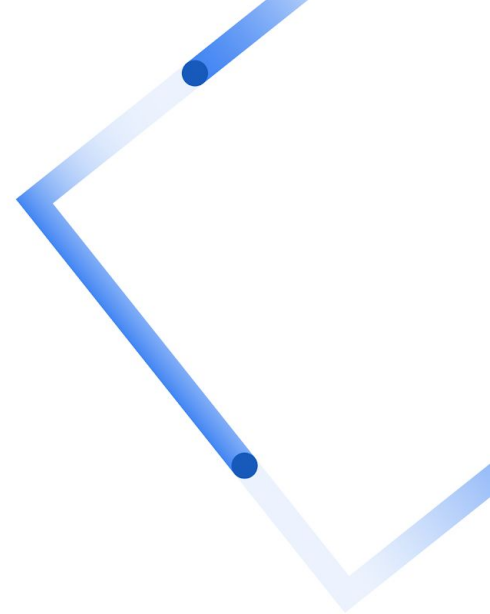Ideally, find labels that are **maximally confusing** for model

😞  this set changes as training progresses

😖  finding these exactly still requires sweeping over all labels!

😀  can **approximate**: find hardest labels **within a large batch** of uniformly sampled labels

$(x_1,$           $)$

Reddi et al., '18. Stochastic-negative mining for learning in large output spaces.
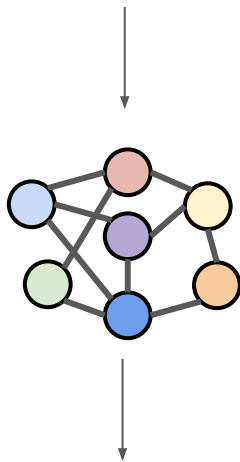
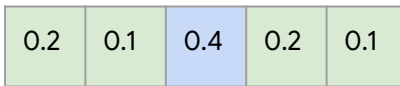**04**

# Model churn

# Supervised learning in theory

Training data

$$\{(x_n, y_n)\}_{n=1}^N$$

Model training

$$\min_{f \in \mathcal{F}} \frac{1}{N} \sum_{n \in N} \ell(y_n, f(x_n))$$

Model predictions

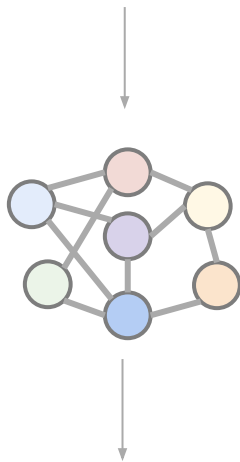| 0.2 | 0.1 | 0.4 | 0.2 | 0.1 |

$$f(x^*)$$

# Supervised learning in practice
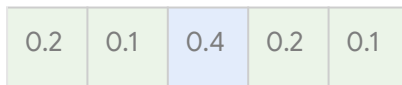
Training data

$$\{(x_n, y_n)\}_{n=1}^N$$

Model training

$$\min_{f \in \mathcal{F}} \frac{1}{N} \sum_{n \in N} \ell(y_n, f(x_n))$$

What if this operation is **stochastic?**

Model predictions

| 0.2 | 0.1 | 0.4 | 0.2 | 0.1 |

$$f(x^*)$$

# Churn in a nutshell

**Model prediction disagreement** under different training and/or inference conditions



Test example

Base model,
e.g., ResNet

Same model,
re-trained

**Power drill**
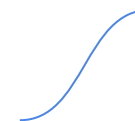(probability 0.48)

**Sewing machine**
(probability 0.68)

# Churn for classification

Suppose we have two classification models, $M_1$ and $M_2$
> e.g., two independently trained models on the same data

The corresponding churn is the probability of disagreement:

$$\text{Churn}(M_1, M_2) = \textbf{Pr}(M_1(x) \neq M_2(x))$$
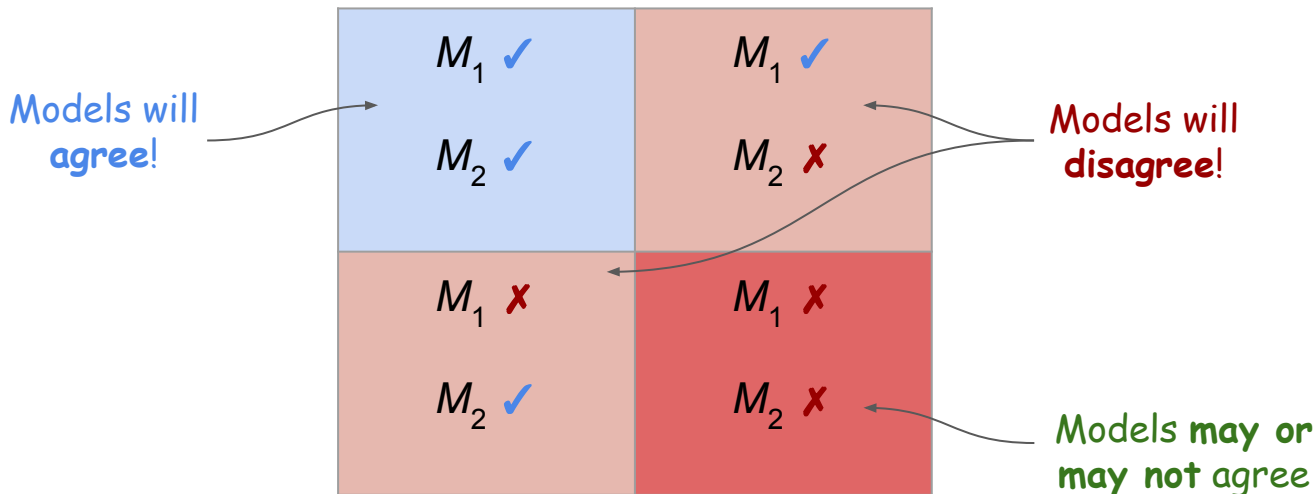
Fraction of times
they predict a
different label

# Churn versus accuracy

Churn can only occur when one or both models is wrong

The better the individual models, the lower the churn



Models will **agree!**

$M_1$ ✓

$M_2$ ✓

$M_1$ ✓

$M_2$ ✗

Models will **disagree!**

$M_1$ ✗

$M_2$ ✓

$M_1$ ✗

$M_2$ ✗

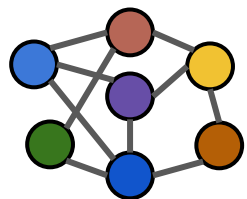Models **may or may not** agree

# Churn versus accuracy variation

Churn ≠ variation in accuracy

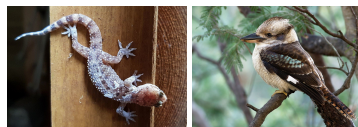Models may disagree here!



60% accuracy

60% accuracy

# Churn from model training

Churn exists even when training on the **same** data, due to several sources of randomness:



Training examples

Shuffle → Initialise params → Train

Shuffle → Initialise params → Train

Can produce different predictions!

# Churn from computing platform

Inherent non-determinism in GPU and TPU

Floating-point addition is not associative!



```
[1]  (0.1+0.2)+0.3

     0.6000000000000001

[2]  0.1+(0.2+0.3)

     0.6
```

# Do neural models exhibit churn?

Unfortunately, **yes**

Predictions from 5x independently trained ResNet models on ImageNet
  76.0% accuracy with 0.1% standard deviation
  Disagreement on **15%** of examples!



power drill - 0.48
sewing machine - 0.68
sewing machine - 0.28
sewing machine - 0.53
power drill - 0.87

wooden spoon - 0.24
spaghetti squash - 0.71
**French loaf** - 0.67
French loaf - 0.57
French loaf - 0.63

swing - 0.82
**lawn mower** - 0.56
tricycle - 0.49
balance beam - 0.75
lawn mower - 0.45

fountain pen - 0.46
can opener - 0.28
crossword - 0.62
hammer - 0.22
crossword - 0.5

How do we mitigate such prediction differences?

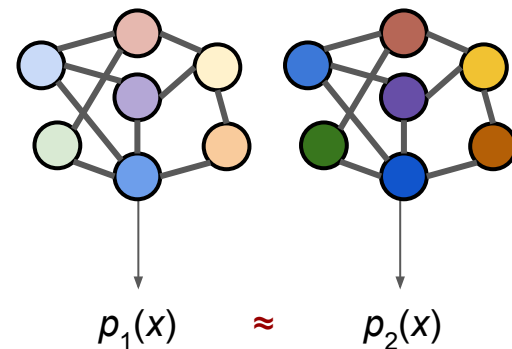Bhojanapalli et al., '20. On the reproducibility of neural network predictions

# Co-distillation

**Motivation:** churn is partly a result of randomness in training

**Idea:** explicitly try to smooth out this randomness!

**Approach:** train two independent models, and encourage their predictions to be similar to each other

> Can be seen as "co-distillation"
>
> **Bonus:** also improves performance!



$$p_1(x) \quad \approx \quad p_2(x)$$

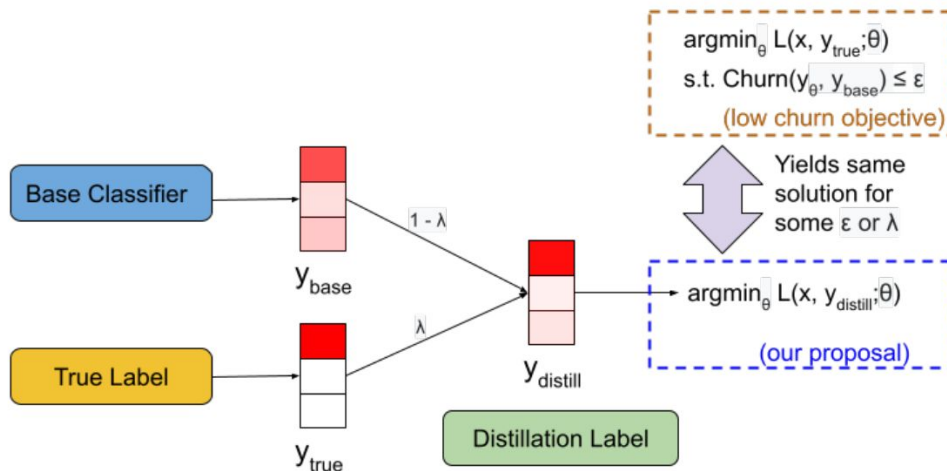Anil et al., '18. Large scale distributed neural network training through online distillation

# Distillation for churn

Churn can also occur more generally between model versions
    e.g., models trained on different weeks, with different architectures, …
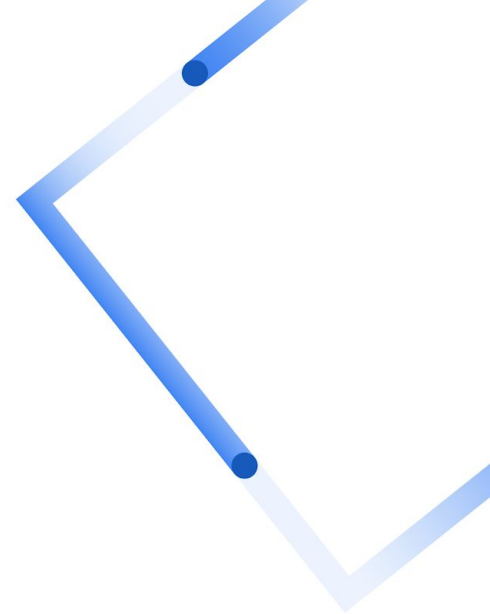
**Idea**: constrain predictions to be similar to original model

**Implementation**: distillation!



Jiang et al., '20. Churn Reduction via Distillation

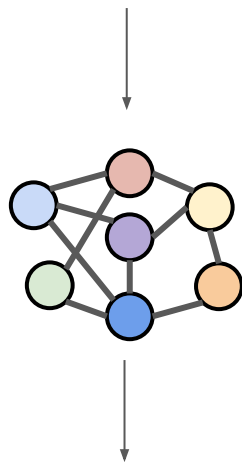**05**

# Summary

Google Research

# Supervised learning in practice!

Training data

$$\{(x_n, y_n)\}_{n=1}^{N}$$

Model training

$$\min_{f \in \mathcal{F}} \frac{1}{N} \sum_{n \in N} \ell(y_n, f(x_n))$$

Churn reduction     Negative mining     Distilled label

Model predictions

| 0.2 | 0.1 | 0.4 | 0.2 | 0.1 |
|-----|-----|-----|-----|-----|

$$f(x^*)$$

# Thank you!

**Aditya Krishna Menon**

Research Scientist

Google NYC

Google Research

# Acknowledgements

# Entropy regularisers

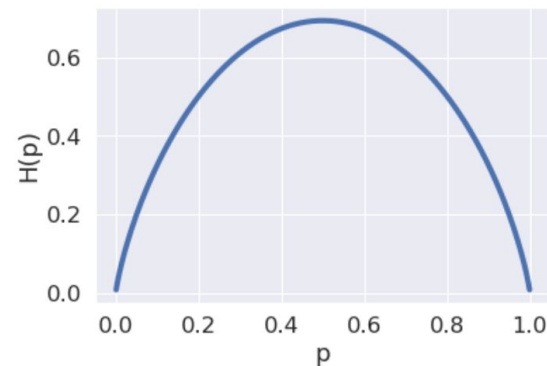**Motivation:** churn occurs when samples' labels flip



**Idea:** move examples away from the classifier boundary!

**Approach:** reduce prediction entropy: for logits $p$, penalise

$$H(p) = -\Sigma_i \, p_i \log p_i$$

discourage highly uncertain predictions

# Churn from data changes

Refreshes of the data can change the learned model



Training examples (January)

Train

Training examples (February)

Train

Can produce different predictions!