# Optimization for inference
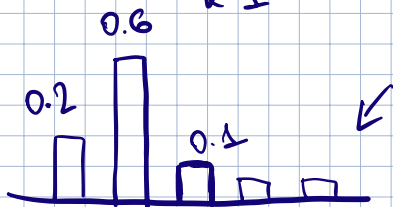
- Calibration & Uncertainty estimation
- Distillation
- Pruning
- Quantization

## Calibration
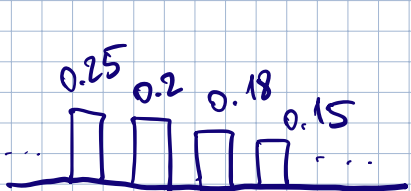
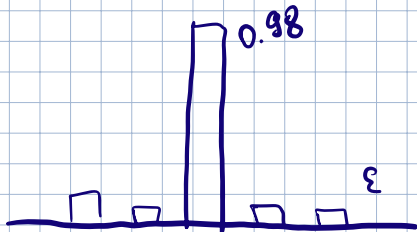$$P_i = \frac{\exp(x_i)}{\sum_{k=1}^{C} \exp(x_k)}$$
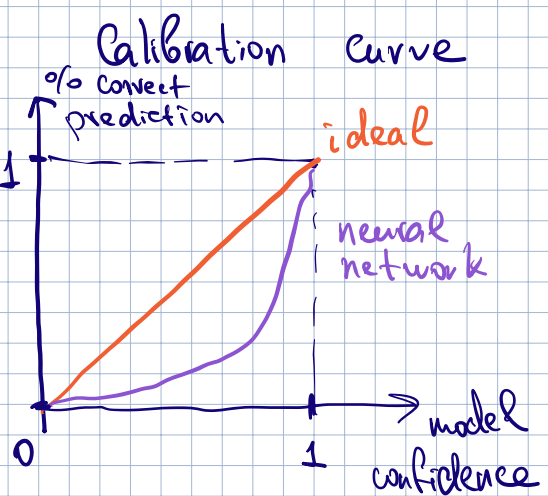
x - logits
p - probabilities

confident prediction

0.6
0.2
0.1

unconfident prediction

0.25  0.2  0.18
0.15

reality

0.98

ε

### Calibration curve

% correct prediction

ideal

neural network

1

0

1

model confidence

- Temperature softmax $(\tau > 0)$

$$p_i^\tau = \frac{\exp(x_i/\tau)}{\sum_{k=1}^{C} \exp(x_k/t)} \qquad i = \text{argmax } x$$

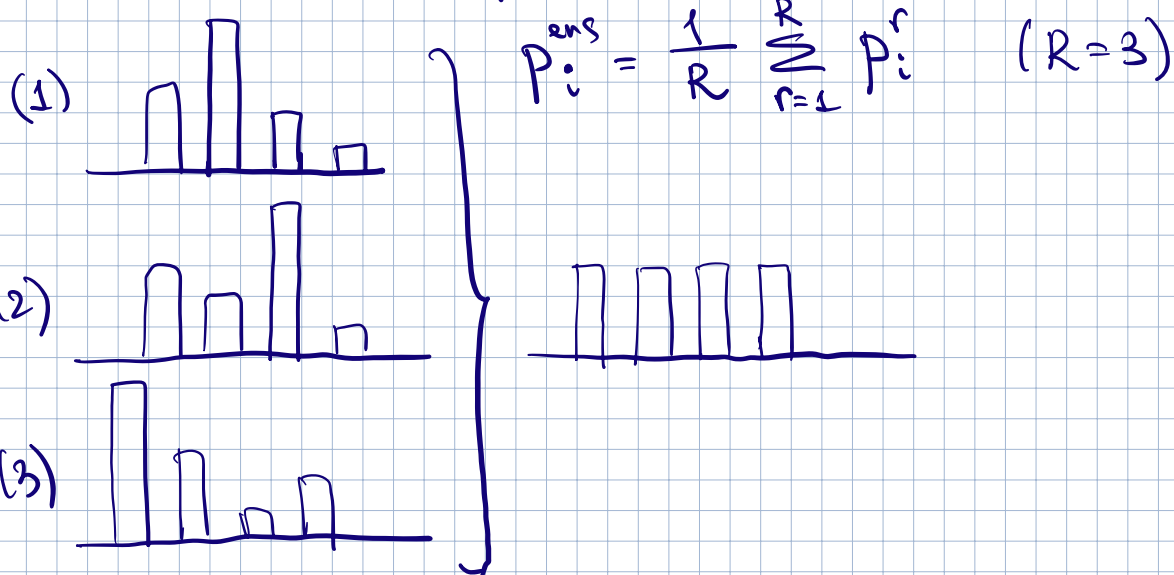1) $\tau \to 0 \qquad p^\tau \to \delta(\text{argmax } x) = (0, 0, .., 0, 1, 0, ..., 0)$

2) $\tau \to \infty \qquad p^\tau \to U(\{1, ..., C\})$

$\tau > 1$ usually in practice

$$t^* = \underset{\tau}{\text{argmin}} \left( -\sum_{n=1}^{N} \sum_{k=1}^{C} [y_n = k] \log p_{nk}^\tau \right)$$

$\underbrace{\hphantom{-\sum_{n=1}^{N} \sum_{k=1}^{C} [y_n = k] \log p_{nk}^\tau}}_{\text{on validation set}}$

- Ensemble (Deep Ensemble)

$$p_i^{ens} = \frac{1}{R} \sum_{r=1}^{R} p_i^r \qquad (R = 3)$$

(1)

(2)

(3)

# (Knowledge) Distillation

$P_{nk}^{T}$ — teacher probability $\Big\}$ with temperature
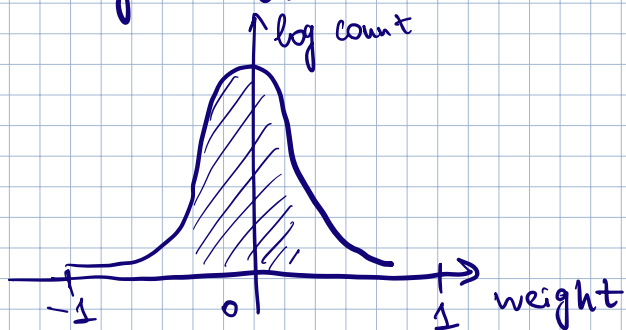
$P_{nk}^{S}$ — student probability

CE loss : $\mathcal{L}_{CE} = -\sum_{n=1}^{N} \sum_{k=1}^{C} [y_n = k] \log P_{nk}^{S}$

Distill. loss : $\mathcal{L}_{D} = -\sum_{n=1}^{N} \sum_{k=1}^{C} P_{nk}^{T} \log P_{nk}^{S}$

$$\mathcal{L} = \alpha \cdot \mathcal{L}_{D} + (1-\alpha) \mathcal{L}_{CE} \quad , \quad \alpha \in (0,1)$$
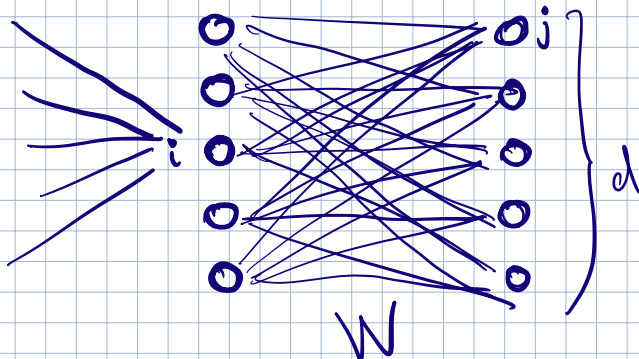
## Pruning

Weights distribution



drop $\alpha\%$ weights
with lowest magnitude
- unstructured pruning

- structured pruning



$$P_i^{\alpha} = \sqrt[\alpha]{\sum_{j=1}^{d} |w_{ij}|^{\alpha}}$$

Several pruning stages:

1. Drop $\alpha$% weight
2. Fine-tune pruned network

SVD of linear layers

$$W \in \mathbb{R}^{d_1 \times d_2}$$

$$W = U \Sigma V \quad - \text{SVD of weight matrix}$$
$$U \in \mathbb{R}^{d_1 \times s} \quad \Sigma \in \mathbb{R}^{s \times s}$$
$$V \in \mathbb{R}^{s \times d_2}$$

Quantization

float 32 $\rightarrow$ int 8

$$W \in \mathbb{R}^{d_1 \times d_2} \quad - \text{weight matrix}$$

$w$ — one scalar weight (float 32)
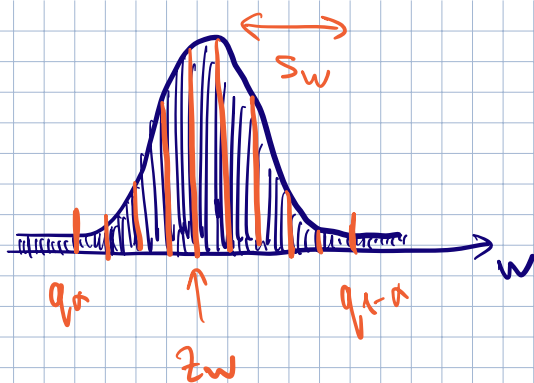
$$w = \boxed{S_w} \, ( w^q - \boxed{Z_w} )$$

$S_w$ — float 32 — scale parameter
$Z_w$ — int 8 — zero parameter
$w^q$ — int 8 — quantized weight

$$w^q = Z_w + \frac{w}{S_w}$$

$$w^q = \text{clip}\left( \text{round}\left(\frac{w}{S_w}\right) + Z_w, \, -128, \, 127 \right)$$

$$q_\alpha \Rightarrow -128$$
$$q_{1-\alpha} \Rightarrow 127$$

$$X \in \mathbb{R}^{B \times d_1} \qquad Y \in \mathbb{R}^{B \times d_2} \qquad W \in \mathbb{R}^{d_1 \times d_2}$$

$$Y = XW \qquad Y_{ij} = \sum_{k=1}^{d_1} X_{ik} \cdot W_{kj}$$

$$\boxed{S_x, Z_x \quad ; \quad S_y, Z_y \quad ; \quad S_w, Z_w}$$

$$Y_{ij} = \sum_{k=1}^{d_1} S_x (X_{ik}^q - Z_x) \cdot S_w (W_{kj}^q - Z_w) =$$

$$= S_x S_w \left[ \sum_{k=1}^{d_1} X_{ik}^q \cdot W_{ik}^q - Z_x \sum_{k=1}^{d_1} W_{kj}^q - Z_w \sum_{k=1}^{d_1} X_{ik}^q + d_1 Z_x Z_w \right]$$

$$\overset{\shortparallel}{S_y (Y_{ij}^q - Z_y)}$$

matrix mult in int8