

Any reproduction or distribution of this document, in whole or in part, is prohibited unless permission is granted by the authors

EE-559

Deep Learning

slido



Have you completed the practice exercises for the first two lectures?

① Start presenting to display the poll results on this slide.

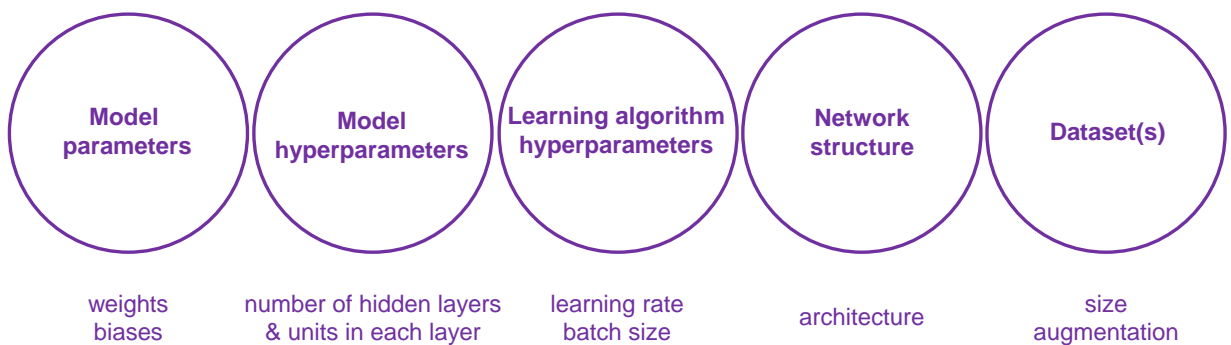
What's on today?

- **The design space**: the (many) choices for a model
- **Training the model**: how to fit the model to a dataset
- **Model performance**: how well a model predicts
- **Classification**: binary and multi-class problems
- **Real-world impact**: the cost of model fitting and deployment
- **Exercises**: the full training pipeline and the impact of data

The design space

$$y = f(x; \Theta)$$

The design space: choices for a model



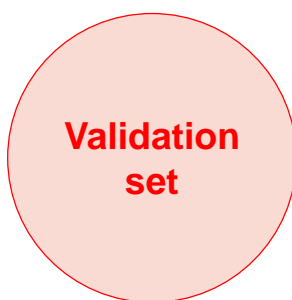
$$\{x_i, y_i\}_{i=1}^N$$

Datasets



to determine the
model parameters

minimize the loss



to determine the
hyperparameters

heuristics or search

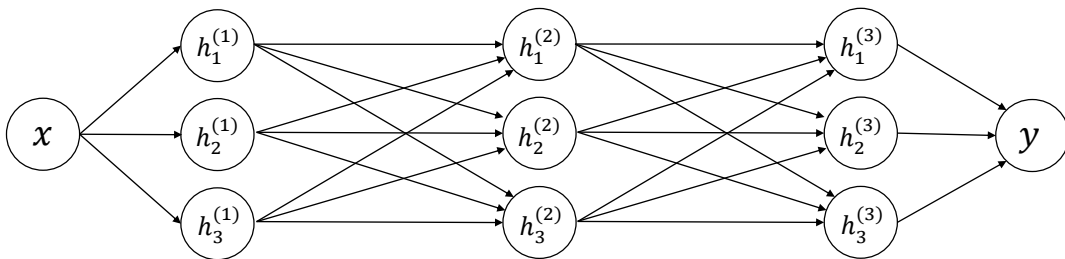


to measure the
ability **to generalize**

model performance

Training the model

Our 'deep' neural network



$$y = \theta_0^{(3)} + \boldsymbol{\theta}^{(3)} a \left[\theta_0^{(2)} + \boldsymbol{\theta}^{(2)} a \left[\theta_0^{(1)} + \boldsymbol{\theta}^{(1)} a \left[\theta_0^{(0)} + \boldsymbol{\theta}^{(0)} x \right] \right] \right]$$

how to determine the value of the model parameters?

Datasets



to determine the
model parameters

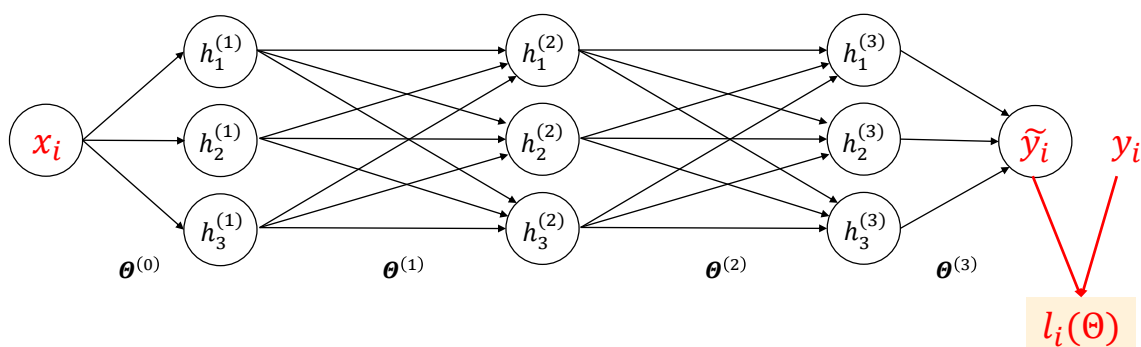


to determine the
hyperparameters



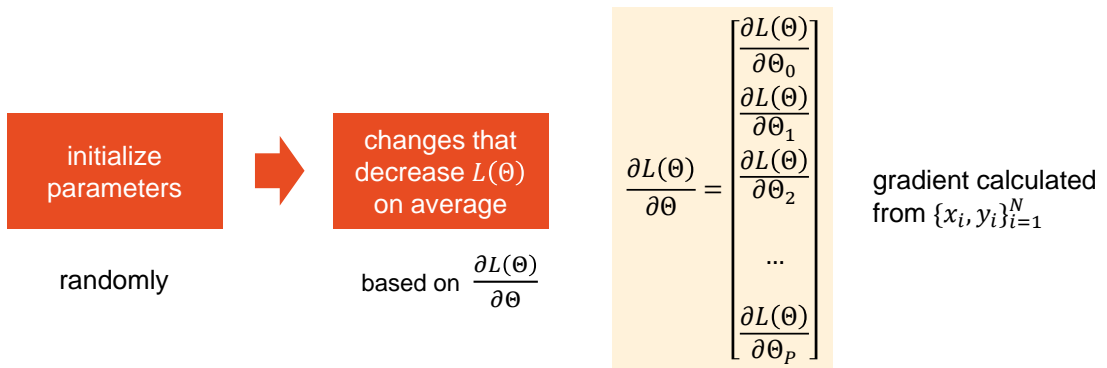
to measure the
ability to generalize

Computing the loss, the training signal



$$\{x_i, y_i\}_{i=1}^N \quad L(\theta) = \sum_{i=1}^N l_i(\theta) \quad \text{e.g. } L(\theta) = \sum_{i=1}^N (f(x_i; \theta) - y_i)^2$$

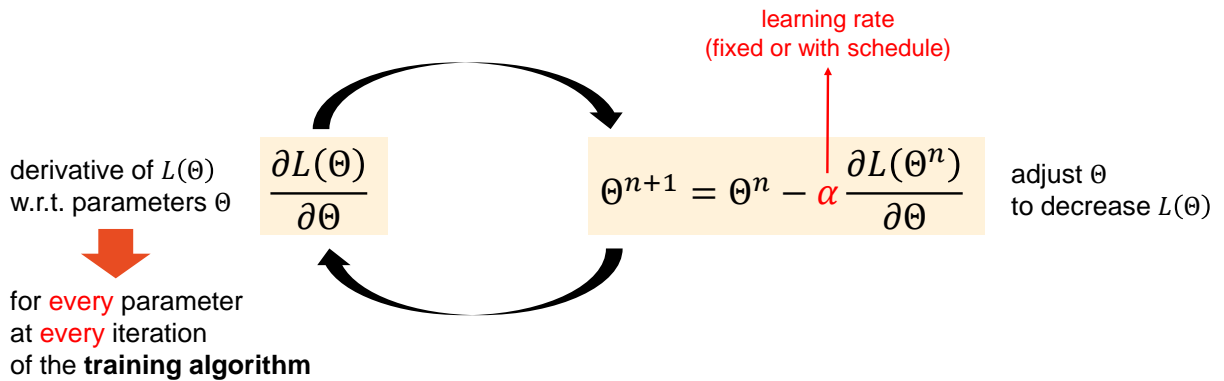
Determining the model parameters



Concepts: Initialization, vanishing gradient, gradient explosion, gradient of the loss with respect to the parameters (at current position)

1,000,000,000,000

Gradient descent



how to calculate the gradients efficiently?

Concepts: Global minimum, local minima, saddle points, convexity, concavity, trajectory of the gradient, exploration of the parameter space

Stochastic gradient descent (SGD)

gradient calculated from $\{x_i, y_i\}_{i=1}^N$

$$\theta^{n+1} = \theta^n - \alpha \sum_{i=1}^N \frac{\partial l_i(\theta^n)}{\partial \theta}$$

gradient estimate calculated from a **batch** (random subset of data) B^n with $|B^n| \in \{1, 2, \dots, N\}$ where $|\cdot|$ is the cardinality of a set

$$\theta^{n+1} = \theta^n - \alpha \sum_{i \in B^n} \frac{\partial l_i(\theta^n)}{\partial \theta}$$

set of batch indices at iteration n

Concepts:

Batch, epoch, noise (randomness) in the gradient, deterministic descent on a changing loss, normalisation

Momentum term

$$\Theta^{n+1} = \Theta^n - \alpha \mathbf{m}^{n+1} \quad \alpha \text{ learning rate}$$

$$\mathbf{m}^{n+1} = \beta \mathbf{m}^n + (1 - \beta) \sum_{i \in B^n} \frac{\partial l_i(\Theta^n)}{\partial \Theta} \quad \beta \in [0,1) \quad \text{controls the degree of smoothing of the gradient over time}$$

effective learning rate $\left\{ \begin{array}{ll} \text{increases} & \text{if gradient directions over time are aligned} \\ \text{decreases} & \text{if gradient directions over time change} \end{array} \right.$

Concepts:

Weighted combination of gradients, reduced oscillations at valleys, Nesterov accelerated momentum

Normalizing the gradient

$$\Theta^{n+1} = \Theta^n - \alpha \frac{\partial L(\Theta^n)}{\partial \Theta}$$



$$\Theta^{n+1} = \Theta^n - \alpha \frac{\mathbf{m}^{n+1}}{\sqrt{\mathbf{v}^{n+1}} + \epsilon}$$

$$\begin{aligned} \mathbf{m}^{n+1} &= \frac{\partial L(\Theta^n)}{\partial \Theta} \\ \mathbf{v}^{n+1} &= \left(\frac{\partial L(\Theta^n)}{\partial \Theta} \right)^2 \end{aligned}$$

Concept:

Point-wise operations

Adaptive moment estimation (Adam)

$$\Theta^{n+1} = \Theta^n - \alpha \frac{\tilde{m}^{n+1}}{\sqrt{\tilde{v}^{n+1} + \epsilon}}$$

$$\tilde{m}^{n+1} = \frac{m^{n+1}}{1 - (\beta)^{n+1}}$$

$$m^{n+1} = \beta m^n + (1 - \beta) \sum_{i \in B^n} \frac{\partial l_i(\Theta^n)}{\partial \Theta}$$

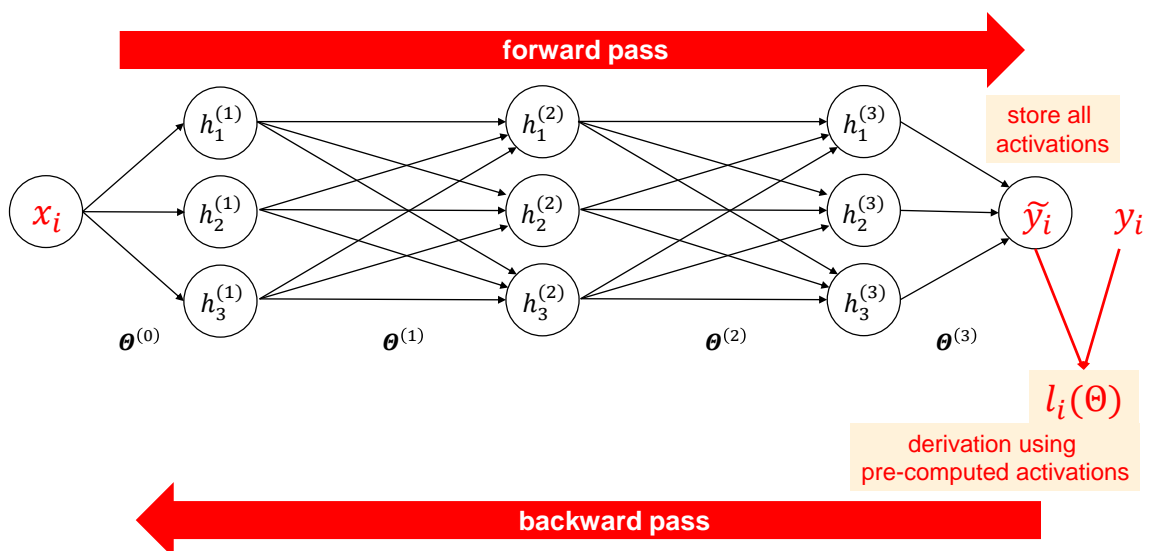
$$\tilde{v}^{n+1} = \frac{v^{n+1}}{1 - (\gamma)^{n+1}}$$

$$v^{n+1} = \gamma v^n + (1 - \gamma) \sum_{i \in B^n} \left(\frac{\partial l_i(\Theta^n)}{\partial \Theta} \right)^2$$

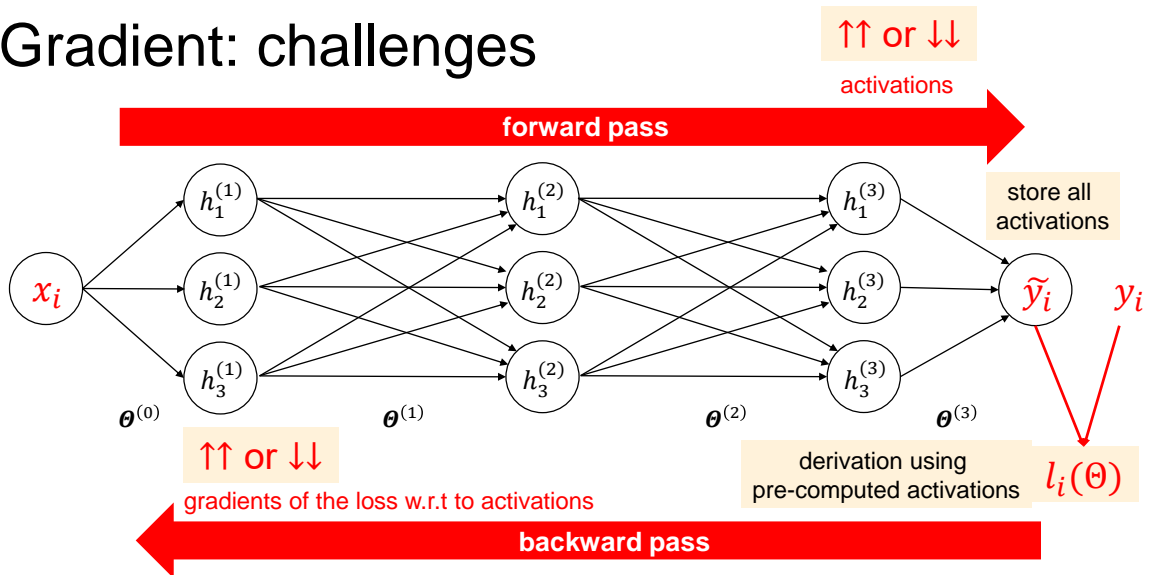
$$\gamma, \beta \in [0, 1)$$

SGD is a special case of Adam ($\gamma = \epsilon = 0$)

Backpropagation algorithm



Gradient: challenges



Concepts: Different starting point for the optimization for each different initialization, gradients of the activations may “vanish” or “explode” (worsens with backward pass), problem worsens with depth

Exploding or vanishing gradient

$$y = \theta_0^{(3)} + \theta^{(3)} a \left[\theta_0^{(2)} + \theta^{(2)} a \left[\theta_0^{(1)} + \theta^{(1)} a \left[\theta_0^{(0)} + \theta^{(0)} x \right] \right] \right]$$

$$\begin{aligned} f^{(n)} &= \theta_0^{(n)} + \theta^{(n)} h^{(n)} \\ &= \theta_0^{(n)} + \theta^{(n)} a[f^{(n-1)}] \end{aligned}$$

With ReLU, range of $h^{(n)}$ is half of range of $f^{(n-1)}$ (clipping below 0)

$$\theta_0^{(n)} = 0$$

$$\theta^{(n)} \sim N(\mu = 0, \sigma^2)$$

If $\sigma^2 \ll 1 \rightarrow$ much smaller magnitude than the input
If $\sigma^2 \gg 1 \rightarrow$ much larger magnitude than the input

Objective: to maintain the variance of $h^{(n)}$ the same across every layer

Concepts: magnitudes of the gradient decrease (increase) uncontrollably during the backward pass, updates to the model become vanishingly small (vanishing) or unstable (“exploding”)

Model performance

Under and overfitting



underperforming on the training set

Model is too simple
Model needs more training time and/or data
High bias



underperforming on the test set

Low-quality and/or small size training data
Model capacity used to model training data noise
High variance

Model performance: cause of errors

Choice of the model



unable to describe the true function



bias \uparrow

Choice of the training data



limited, noisy



variance \uparrow

Concepts: Sampling of the input space, choice of the hyperparameters, stochastic learning algorithm may not converge to the same solution each time

How to reduce bias and variance?

add **hidden units** and/or **hidden layers**

$D \uparrow \Rightarrow$ **bias** \downarrow

increase quantity and/or quality of **training data**

$N \uparrow \Rightarrow$ **variance** \downarrow

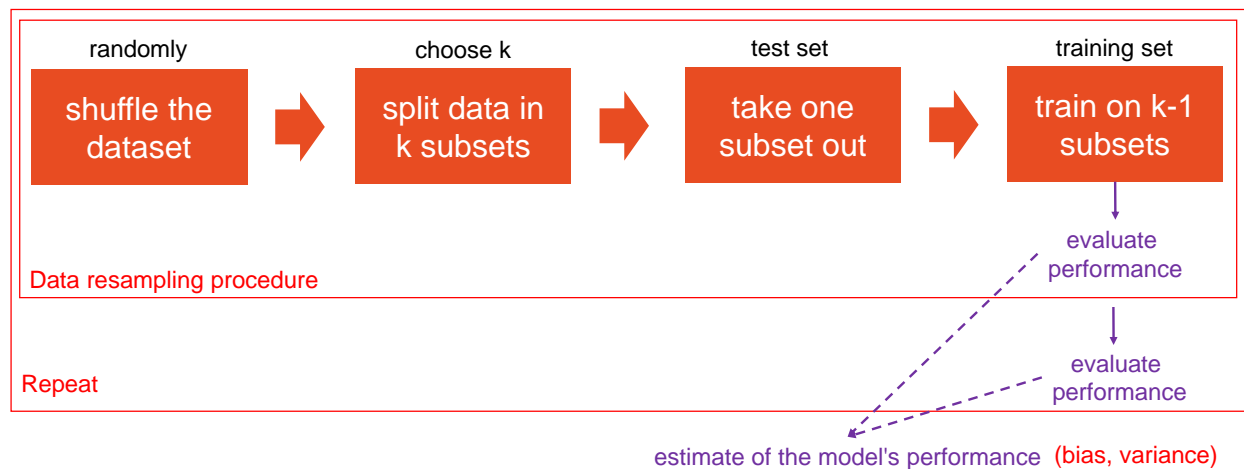
Given a fixed N

if $D \uparrow \Rightarrow$ **bias** \downarrow but **variance** $\uparrow \Rightarrow$ trade-off

Concepts:

Bias-variance trade-off, a model capacity increase does not (necessarily) imply a test error reduction

k-fold cross validation



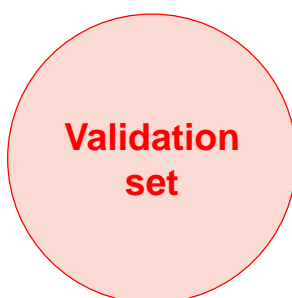
Concepts:

Aggregation of performance results, leave-p-out cross validation, nested cross-validation

Datasets



to determine the
model parameters

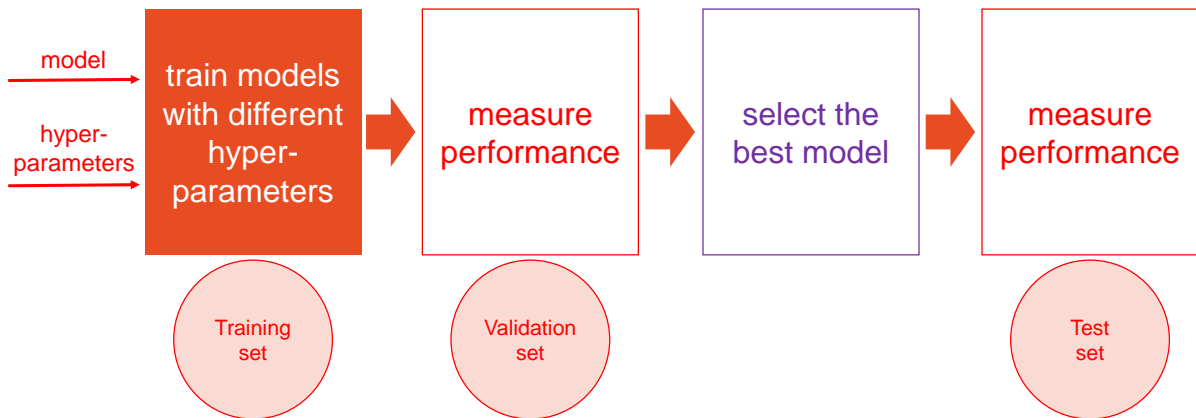


to determine the
hyperparameters



to measure the
ability to generalize

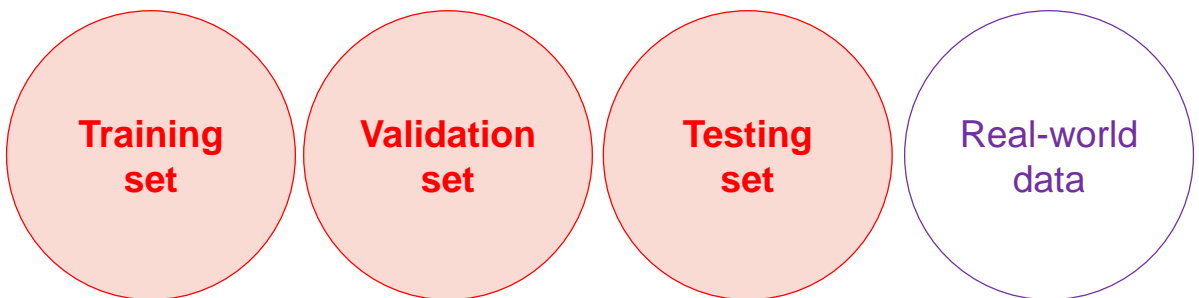
Hyperparameters



Concepts:

Hyperparameter search, empirical choice, sampling the hyperparameter space

Data

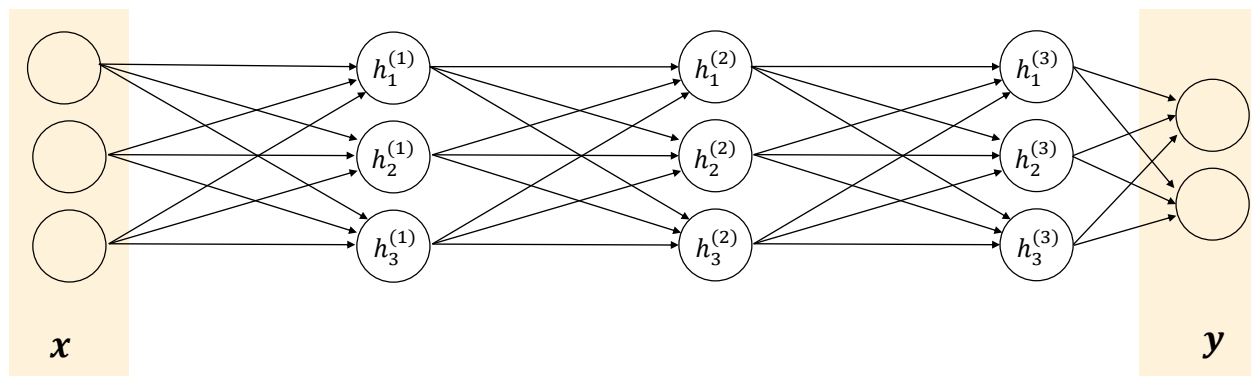


Classification

Recap: negative log-likelihood criterion

$$\theta^* = \arg \min_{\theta} \left[- \sum_{i=1}^N \log[P(y_i | f(x_i; \theta))] \right]$$

Recap: network diagram



Binary classification

$y_i \in \{0,1\}$ 2 classes

Probability distribution (*Bernoulli*) over the output space:

$$P(y|\lambda) = \begin{cases} 1 - \lambda & y_i = 0 \\ \lambda & y_i = 1 \end{cases} \quad \lambda \in [0,1]$$

$$P(y|\lambda) = \lambda^y (1 - \lambda)^{1-y}$$

Concepts:
Categories, labels

Binary cross-entropy loss

$$\text{sig}(z) = \frac{1}{1 + e^{-z}} \quad \text{logistic sigmoid } (\mathbb{R} \rightarrow [0,1])$$

$$P(y|x) = \text{sig}(f(x; \Theta))^y (1 - \text{sig}(f(x; \Theta)))^{1-y}$$

$$L(\Theta) = \sum_{i=1}^N -y_i \log[\text{sig}(f(x_i; \Theta))] - (1 - y_i) \log[1 - \text{sig}(f(x_i; \Theta))]$$

inference $\hat{y} = 1$ if $\lambda > 0.5$ (and 0 otherwise)

Concepts:

Setting $f(x; \Theta)$ to predict λ of the Bernoulli distribution, negative log-likelihood of the training set

Multi-class classification

$y_i \in \{1, 2, \dots, \mathbf{c}, \dots, C\}$ C classes

Probability distribution (*categorical*) over the output space:

$P(y = \mathbf{c}) = \lambda_{\mathbf{c}}$ $\lambda_1, \lambda_2, \dots, \lambda_{\mathbf{c}}, \dots, \lambda_C$ probability of each category

$$\lambda_{\mathbf{c}} \in [0,1]$$

$$\sum_{\mathbf{c}=1}^C \lambda_{\mathbf{c}} = 1$$



$$\text{softmax}_{\mathbf{c}}(\mathbf{z}) = \frac{e^{z_{\mathbf{c}}}}{\sum_{\mathbf{c}'=1}^C e^{z_{\mathbf{c}'}}} \quad \text{softmax function}$$

$$P(y = \mathbf{c}|x) = \text{softmax}_{\mathbf{c}}(f(x; \Theta))$$

Concepts:

Exponential function ensures positivity, sum in the denominator ensures “sum to 1”

Multi-class cross-entropy loss

$$L(\Theta) = - \sum_{i=1}^N \log[\text{softmax}_{y_i}(f(x_i; \Theta))]$$

$$= - \sum_{i=1}^N f_{y_i}(x_i; \Theta) - \log \left[\sum_{c'=1}^C e^{f_{c'}(x_i; \Theta)} \right]$$

output c' of the neural network

inference $\hat{y} = \arg \max_c P(y = c | f(x; \Theta^*))$

Concepts:

Categorical distribution of C possible classes, point estimate, most probable category

Cross-entropy loss

Kullback-Leibler (KL) divergence: distance between $q(z)$ and $p(z)$

$$D[q||p] = \int_{-\infty}^{+\infty} q(z) \log[q(z)] dz - \int_{-\infty}^{+\infty} q(z) \log[p(z)] dz$$

empirical data distribution observed at $\{y_i\}_{i=1}^N$

$$q(y) = \frac{1}{N} \sum_{i=1}^N \delta[y - y_i]$$

Dirac delta function

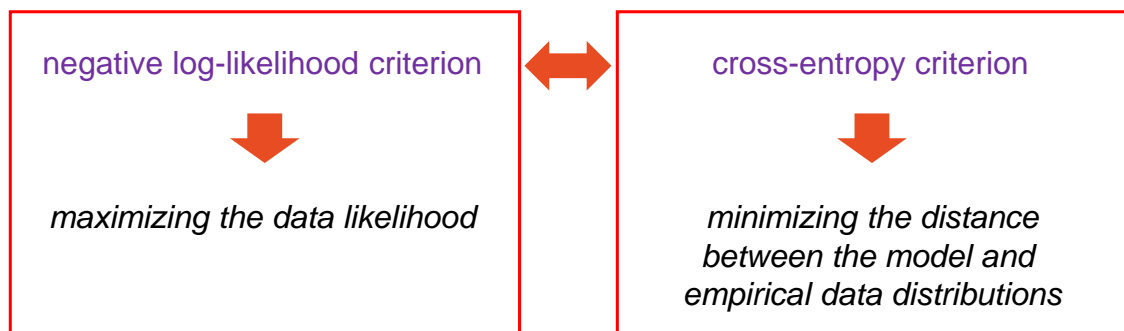
Concept:

Weighted sum of point masses

Cross-entropy loss

$$\begin{aligned}
 \Theta^* &= \arg \min_{\Theta} \left[\int_{-\infty}^{+\infty} q(y) \log[q(y)] dy - \int_{-\infty}^{+\infty} q(y) \log[P(y|\Theta)] dy \right] \\
 &= \arg \min_{\Theta} \left[- \int_{-\infty}^{+\infty} q(y) \log[P(y|\Theta)] dy \right] \quad \rightarrow \text{cross-entropy: amount of uncertainty in a distribution after taking into account what we already know from the other} \\
 &= \arg \min_{\Theta} \left[- \int_{-\infty}^{+\infty} \left(\frac{1}{N} \sum_{i=1}^N \delta[y - y_i] \right) \log[P(y|\Theta)] dy \right] \\
 &= \arg \min_{\Theta} \left[- \frac{1}{N} \sum_{i=1}^N \log[P(y_i|\Theta)] \right] = \arg \min_{\Theta} \left[- \sum_{i=1}^N \log[P(y_i|f(x_i; \Theta))] \right] \\
 &\quad \text{negative log-likelihood criterion}
 \end{aligned}$$

Negative log-likelihood and cross-entropy

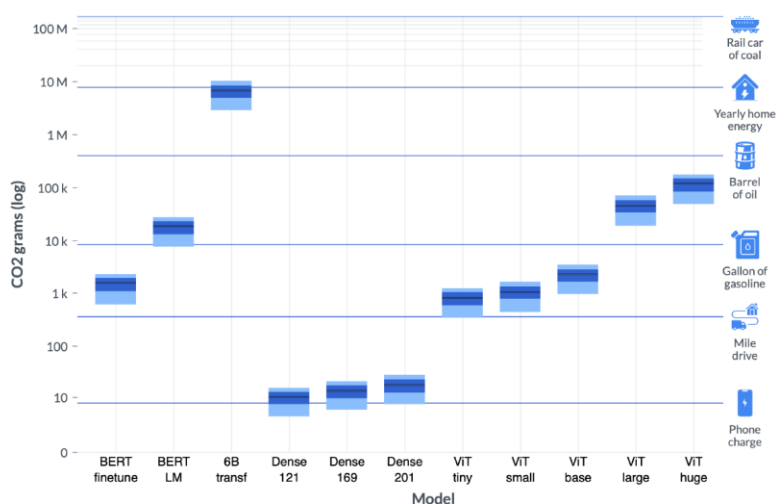


Concept:

Equivalence between negative log-likelihood and cross-entropy criteria

Real-world impact

The cost of model fitting



CO2 Relative Size Comparison

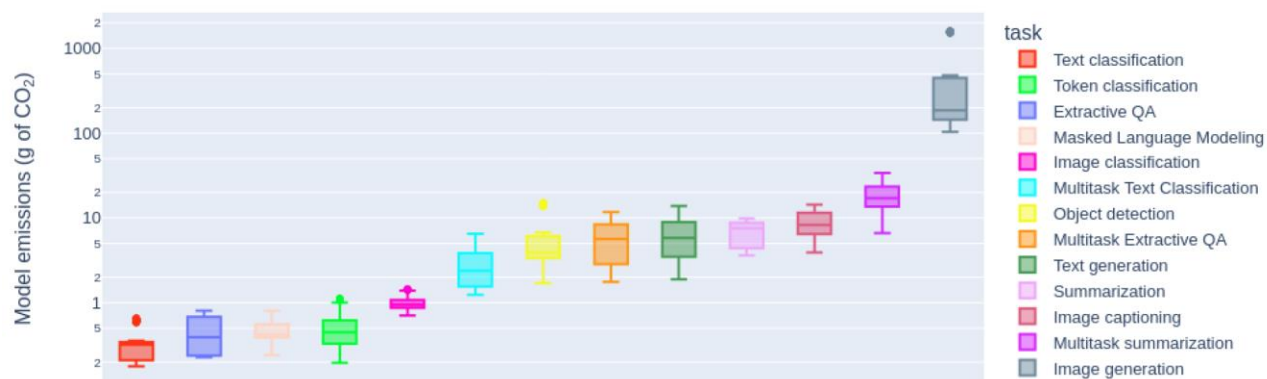
Measuring the carbon intensity of AI in cloud instances

measure **operational carbon emissions** data per energy unit
 - location-based
 - time-specific marginal emissions

natural language processing
 computer vision applications

[arXiv:2206.05229](https://arxiv.org/abs/2206.05229)

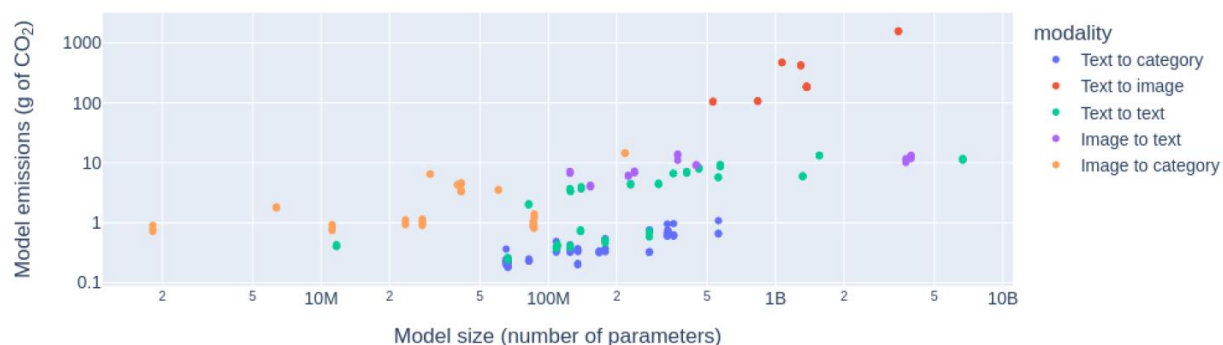
The cost of deployment



for 1000 queries

[arXiv:2311.16863](https://arxiv.org/abs/2311.16863)

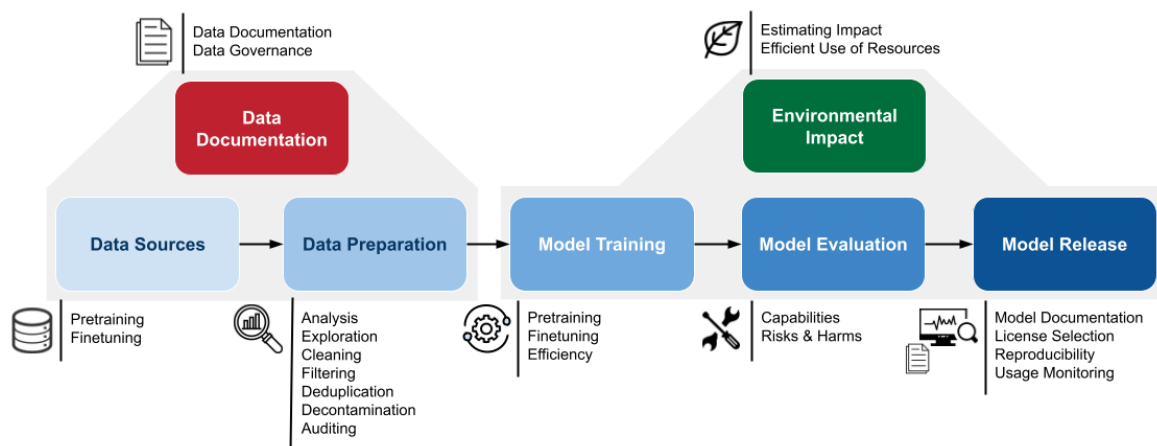
Carbon emitted for 1000 inferences



Note - logarithmic scales in both axes

[arXiv:2311.16863](https://arxiv.org/abs/2311.16863)

Model development



<https://fmcheatsheet.org>

Exercises

Today's practice exercises

You will experiment with a **full training pipeline**

- implementing the model architecture
- training and evaluation
- logging the metrics
- using different optimizers

Today's marked exercises

The impact of data on model performance

- data augmentation on a balanced dataset
- weighted sampler on an imbalanced dataset
- weighted cross-entropy loss to mitigate the imbalance

What did we learn today?

- The design space
- Training the model
- Model performance
- Classification
- Real-world impact
- Exercises

EE-559

Deep Learning

andrea.cavallaro@epfl.ch