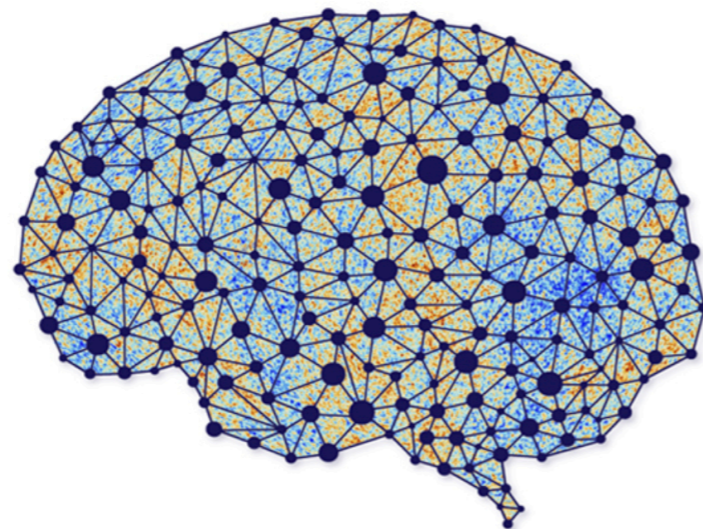


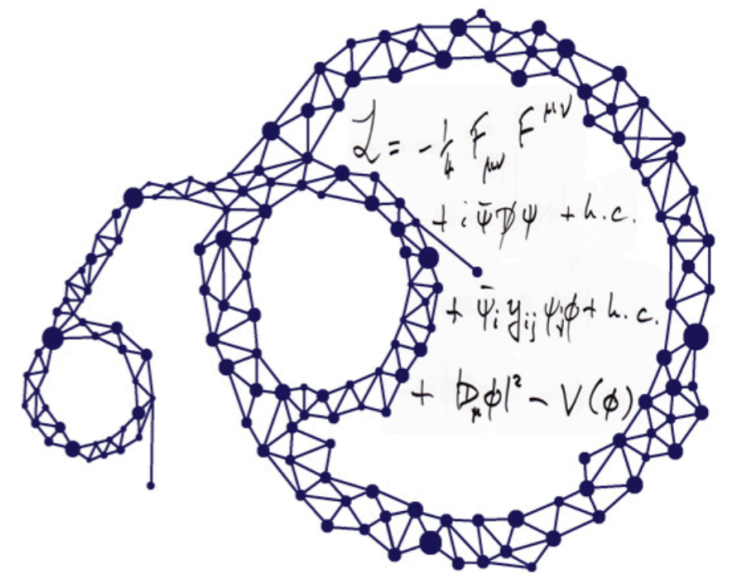
# Physics 361 - Machine Learning in Physics

## Lecture 4 – Basics of Machine Learning

Jan. 29<sup>th</sup> 2025



AI  
∩  
Universe



Moritz Münchmeyer

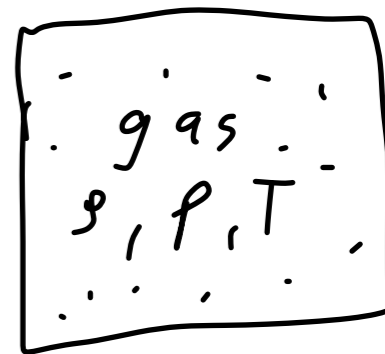
# Unit 1: Background

## 1.4 Information theory background (cont.)

# Entropy

- In statistical physics (thermodynamics) the **entropy** is given by

$$S = (K_B) \log \Omega$$



↑ here we set it to 1

$\Omega$ : number of microstates (equally likely)

- We can re-write this as

$$S = -\log p \quad \text{where} \quad p = \frac{1}{\Omega}$$

is the probability of each m. state

- The general definition of entropy

**Shannon entropy**

$$S = - \sum_i p_i \log p_i$$

- $\log$  here is base  $e$   $\rightarrow$  entropy is in "nats"  
(base 2  $\rightarrow$  " in bits)

# Properties of the entropy

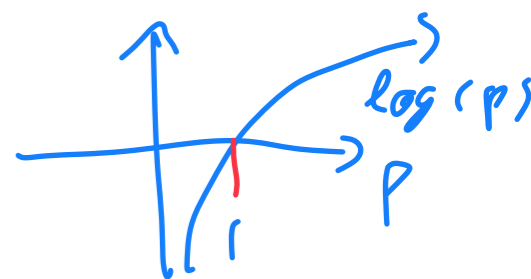
- discrete PDF

$$S = - \sum_i p_i \log(p_i)$$

$-\log(p_i)$  is the „self-information“ of event  $i$

$$= -\mathbb{E}_{x \sim P(x)} [\log(P(x))]$$

A unlikely event has a large selfinfo.



self information  $I(x) = -\log(P(x))$

A certain event has self inform. 0

maximum possible entropy is

- continuous PDF

$$S = - \int dx p(x) \log(p(x))$$

$$= -\mathbb{E}_{x \sim p(x)} [\log(P(x))]$$

$$S = \log(N)$$

More uniform distributions have a higher entropy.  
(spread out)


# Kullback - Leibler divergence

- The relative entropy = KL divergence provides a measure of the similarity of two probab. distr.

$P(x)$  and  $Q(x)$ :

$$D_{KL}(P||Q) = \mathbb{E}_{x \sim P} \left[ \log \frac{P(x)}{Q(x)} \right] = \mathbb{E}_{x \sim P} [\log P(x) - \log Q(x)]$$

difference of self information  
of sample  $x$



- If  $P$  and  $Q$  are the same then  $D_{KL} = 0$ .

e.g. in generative ML

$P(x)$ : true (unknown) PDF of the training images

$Q(x)$ : PDF of generated images

- $D_{KL}$  is easy to evaluate by sampling (if  $P$  and  $Q$  are known).
- $D_{KL} \geq 0$

- Interpretation: If our data is described by  $P$  and we have a theory/model  $Q$  that is meant to describe the data, at what rate will collection of data inform me that theory  $Q$  differs from  $P$ .  
 $=$  KL divergence

- Minimizing the KL divergence is the same as maximizing the Likelihood of the data given the model  $Q$ . We define the **cross-entropy** as the part of the KL divergence that depends on the parameters of  $Q$ .

$$\text{Cross entropy } H(P, Q) = -E_{x \sim P(x)} \log[Q(x)]$$

- The KL divergence is not a true distance metric.  
E.g. it is not symmetric in  $P, Q$ .  
Other alternatives: - Wasserstein distance  
- F-divergence

- Another interesting concept: mutual information

$$I(X, Y) = D_{KL}[P(X, Y) \parallel P(X)P(Y)]$$

vanishes if  $X, Y$  are independent.

All of these information theory concepts are frequently used in machine learning.

# Unit 2: Basics of Machine Learning

Sources: e.g. [deeplearningbook.org](https://deeplearningbook.org)

# Overview

Most machine Learning algorithms have the following elements:

- dataset
  - { training data
  - { test data
  - { validation data
- cost function / Loss function / training objective
- model / architecture
- optimization procedure

We first discuss these at the example of curve fitting.

Then we train a simple neural network.

# Unit 2: Machine Learning Basics

## 2.1 Machine Learning concepts using the example of Linear and Polynomial Regression

# Linear regression

- The simplest machine learning algorithm.
- Predict 1 number from  $N$  features:

$$\hat{y} = \vec{w}^T \vec{x} + b$$

This can be re-written as

$$\hat{y} = \vec{w}^T \vec{x}$$

where we

include  $b$  by adding  
an element 1 to  $x$

$$x = \begin{pmatrix} \vdots \\ 1 \end{pmatrix}$$



- notation: We define the design matrix as

$$X = \begin{matrix} & \text{features} \\ \text{example} & \begin{pmatrix} | & | & | & | \end{pmatrix} \end{matrix}$$

$$X^{\text{train}} : \text{training data}$$

$$\vec{y} : \text{training labels}$$

vector over examples

- training set  $\{X^{\text{train}}, \vec{y}^{\text{train}}\}$
- test set  $\{X^{\text{test}}, \vec{y}^{\text{test}}\}$

- The typical cost function for such regression problems is the Mean Squared Error (MSE).

$$MSE = \frac{1}{m} \sum_i (\hat{y}_i - y_i)^2$$

for linear regression:  $\hat{y}_i = \vec{w}^T x_i$

prediction      truth

- Notation:  $L_p$  norm of a vector is

$$\|\vec{x}\|_p = \left( \sum_i |x_i|^p \right)^{1/p}$$

Thus we can rewrite

$$MSE = \frac{1}{m} \|\vec{\hat{y}} - \vec{y}\|_2^2$$

- Goal is to minimize the MSE.  
train to minimize  $MSE_{train}$  wrt.  $\vec{w}$ .

→ hope  $MSE_{test}$  will also be small.

- For Linear regression we can solve for  $\vec{w}$  analytically:

$$\nabla_w \text{MSE}_{\text{train}} = 0$$

$$\Rightarrow \nabla_w \frac{1}{m} \|\hat{\mathbf{y}}^{(\text{train})} - \mathbf{y}^{(\text{train})}\|_2^2 = 0$$

$$\Rightarrow \frac{1}{m} \nabla_w \|\mathbf{X}^{(\text{train})} \mathbf{w} - \mathbf{y}^{(\text{train})}\|_2^2 = 0$$

$$\Rightarrow \nabla_w \left( \mathbf{X}^{(\text{train})} \mathbf{w} - \mathbf{y}^{(\text{train})} \right)^\top \left( \mathbf{X}^{(\text{train})} \mathbf{w} - \mathbf{y}^{(\text{train})} \right) = 0$$

$$\nabla_w \left( \mathbf{w}^\top \mathbf{X}^{(\text{train})\top} \mathbf{X}^{(\text{train})} \mathbf{w} - 2 \mathbf{w}^\top \mathbf{X}^{(\text{train})\top} \mathbf{y}^{(\text{train})} + \mathbf{y}^{(\text{train})\top} \mathbf{y}^{(\text{train})} \right) = 0$$

$$\Rightarrow 2 \mathbf{X}^{(\text{train})\top} \mathbf{X}^{(\text{train})} \mathbf{w} - 2 \mathbf{X}^{(\text{train})\top} \mathbf{y}^{(\text{train})} = 0$$

$$\Rightarrow \mathbf{w} = \left( \mathbf{X}^{(\text{train})\top} \mathbf{X}^{(\text{train})} \right)^{-1} \mathbf{X}^{(\text{train})\top} \mathbf{y}^{(\text{train})}$$

we need matrix multiplication  
and inversion.

# Background on matrix calculus

(required to understand the last slide, but not often used in this course)

## derivative of quadratic forms:

1. For a symmetric matrix  $A$ :

$$\nabla_w (w^T A w) = 2Aw.$$

2. For a general (not necessarily symmetric) matrix  $A$ :

$$\nabla_w (w^T A w) = (A + A^T)w.$$

3. Gradient with respect to  $w^T$ :

$$\nabla_{w^T} (w^T A w) = w^T (A + A^T).$$

### Properties of Transposition

- $(AB)^T = B^T A^T$
- $(A^T)^T = A$
- If  $A$  is square and invertible, then  $(A^{-1})^T = (A^T)^{-1}$ .

## derivatives of inner products

- If  $f(w) = b^T w$  for a vector  $b$ , then:

$$\nabla_w (b^T w) = b$$

- If  $f(w) = w^T b$ , then

$$\nabla_w (w^T b) = b$$

since bilinear forms are scalars we have:

$$w^T A b = (w^T A b)^T$$

much more about this:

# Linear regression (cont'd)

Very simple ML model:

$$\hat{y} = \vec{w}^T \vec{x}$$

$\vec{x}$ : input  
 $\vec{y}$ : output

training data:

$$\begin{aligned} X^{\text{train}} &= \begin{pmatrix} & \text{features} \end{pmatrix} \left\{ \begin{array}{l} \text{examples} \end{array} \right\} \\ \vec{y}^{\text{train}} &= \text{desired output} = \text{label} \end{aligned}$$

Loss function:

$$MSE = \frac{1}{n} \|\vec{\hat{y}} - \vec{y}\|_2^2$$

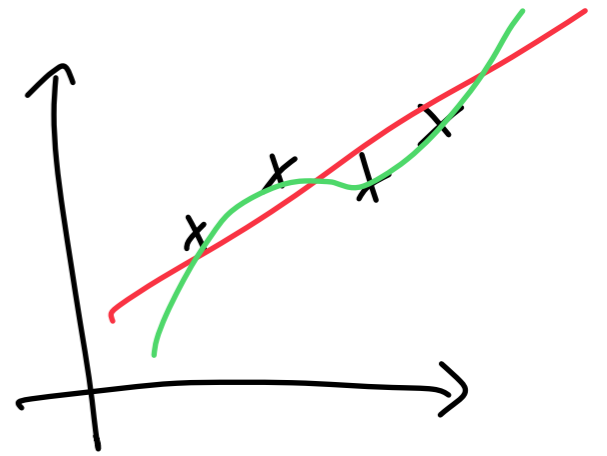
← estimate

← training "label"

Can be minimized w.r.t. to  $\vec{w}$ :

$$\vec{w} = \left( X^{\text{train}T} X^{\text{train}} \right)^{-1} X^{\text{train}T} \vec{y}^{\text{train}}$$

# Polynomial regression



- In polynomial regression we fit a higher order polynomial.

E.g. 2nd order pol.  $\hat{y} = b + w_1 x + w_2 x^2$

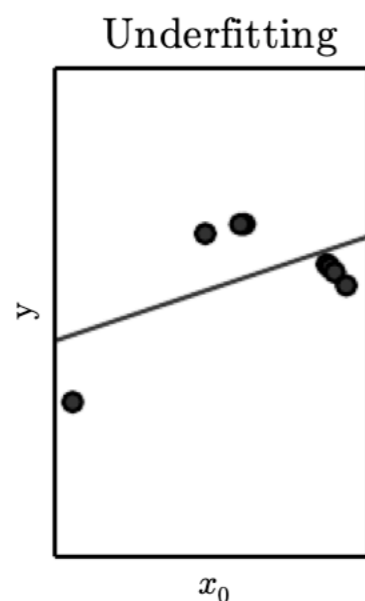
General:  $\hat{y} = b + \sum_{i=1}^n w_i x^i$

- This model can be solved for  $\vec{w}$  with the same equations as linear regression because it is still linear in the parameters  $w_i$ .

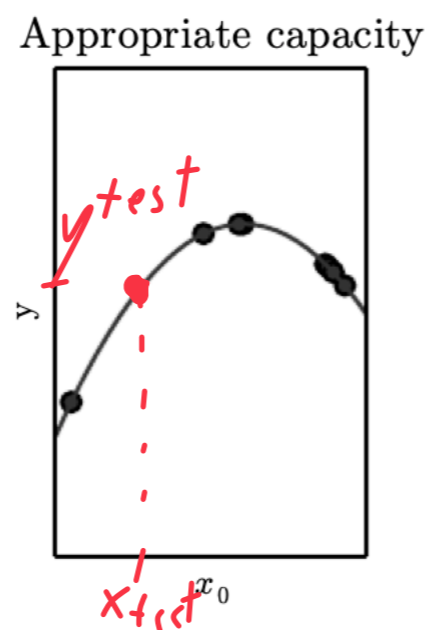
# Capacity, Overfitting, Underfitting

- How well does a model trained/fitted to the training data work on novel data,  
 $\Rightarrow$  Generalization error
- We usually assume that training data and test data are drawn from the same distribution. "i.i.d." data : independent identically distributed
- Generalization is closely related to overfitting and underfitting.
  - underfit : training error is large
  - overfit : training error is small but test error is large.

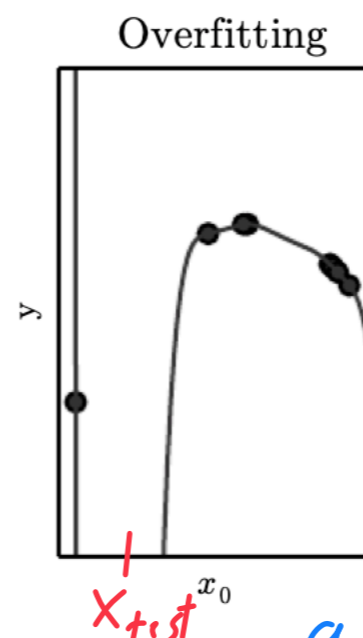
# Example of polynomial regression



linear  
model



quadratic  
fit



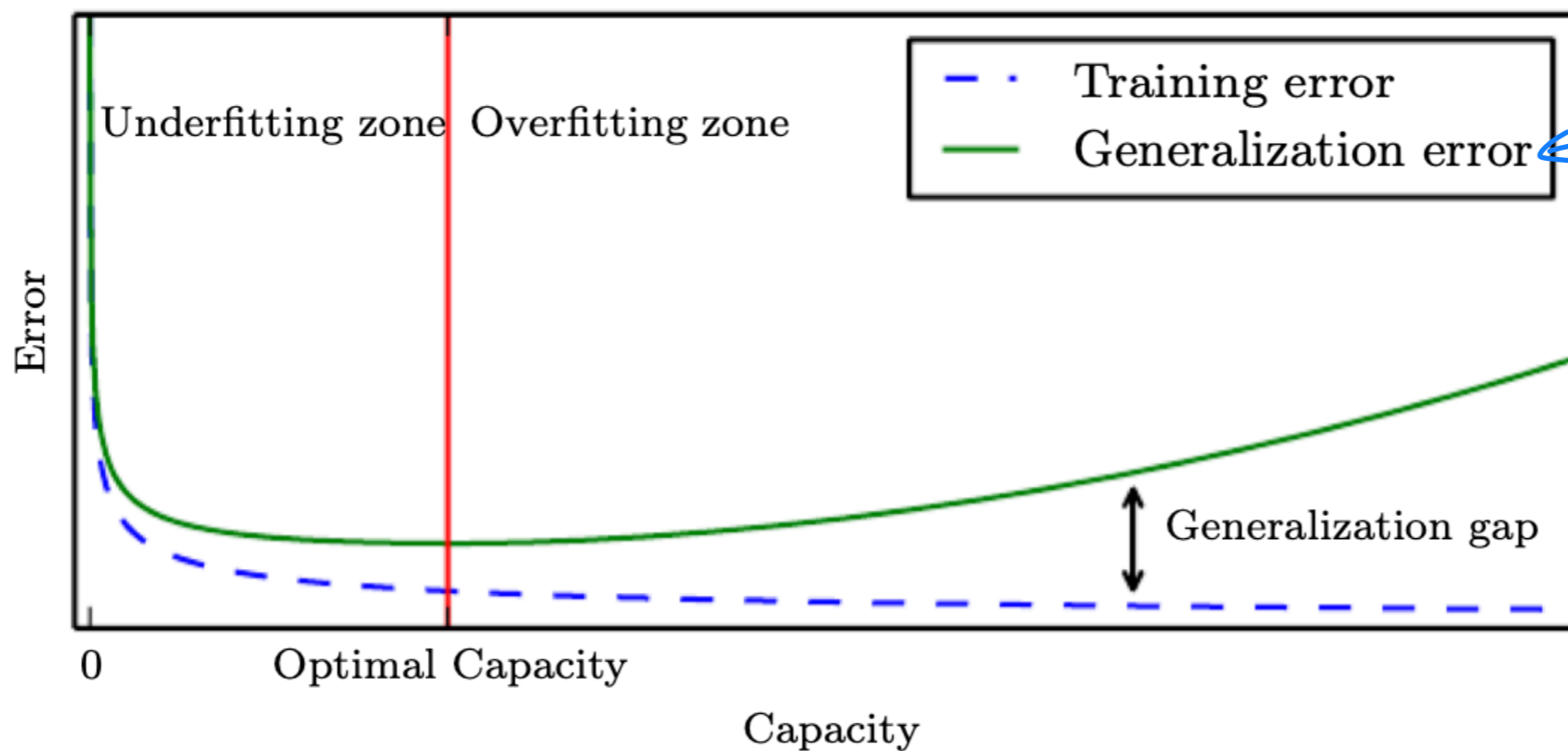
[deeplearningbook.org](http://deeplearningbook.org)

9th order  
polynomial

$y_{test}$

- The space of functions that the model can represent is called the **capacity**.
  - Occam's razor: choose the simplest "that fits"
  - Lower capacity  $\rightarrow$  tends to generalize better
  - Higher capacity  $\rightarrow$  reduces training error,
- $\Rightarrow$  Need to optimize the capacity or "regularize."

Typical behavior of error vs capacity  
(after training out the model)



error on the  
test  
set

[deeplearningbook.org](http://deeplearningbook.org)

We could e.g. train models with various orders of polynomials (or neural network parameter numbers).

Usually this is not the main approach.  
Instead we use "regularization" to avoid overfitting.

# Regularization

- Idea: Keep the model capacity fix, but encourage simpler solutions.

- Common method: **weight decay**

New Loss  
to minimize

$$J(\vec{w}) = \text{MSE}_{\text{train}} + \lambda \underbrace{\vec{w}^T \vec{w}}_{\text{L-2 norm}}$$

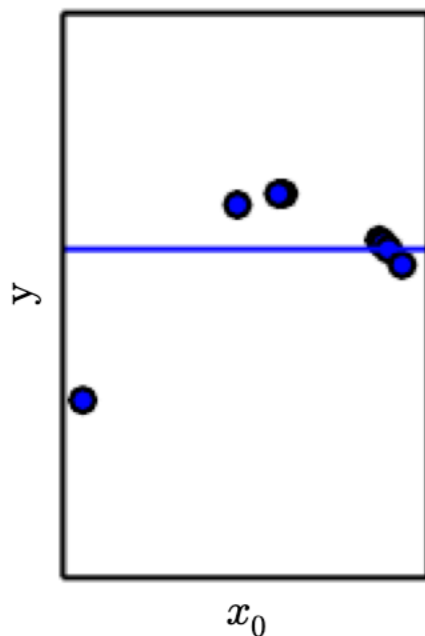
parameter  
to choose  
"lambda"

L-2 norm

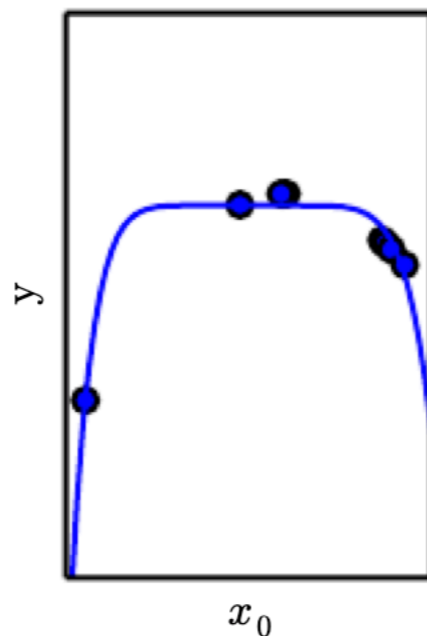
L-2 regularization

This discourages  
large weights

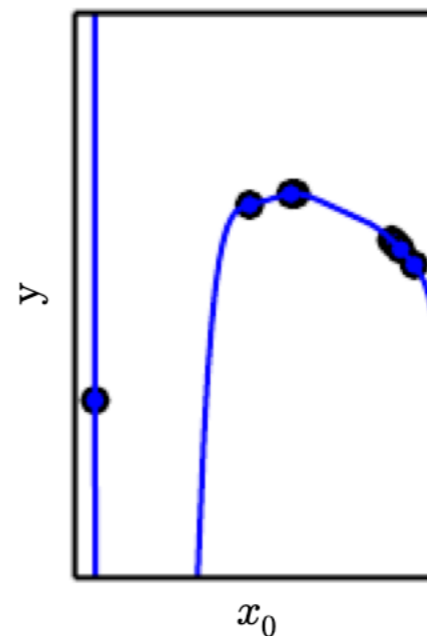
Underfitting  
(Excessive  $\lambda$ )



Appropriate weight decay  
(Medium  $\lambda$ )



Overfitting  
( $\lambda \rightarrow 0$ )



L-2 regularization is one of several popular methods of regularization

Regularization : any method meant to prevent overfitting without changing the model capacity.

Two other popular methods which we discuss later:

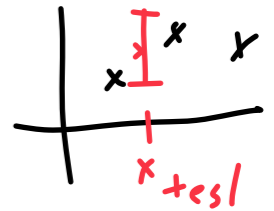
- drop-out
- early stopping.

# Hyperparameters

- Hyperparameters are parameters that are *NOT* part of the fitting/optimization procedure.
- Examples:
  - model architecture
    - e.g. - maximum order of polynomials
    - depth of a neural network
  - Learning parameters
    - e.g. Learning rate
  - Regularization parameters e.g.  $\lambda$
- To choose the best hyperparameters one often splits off a *validation data set* from the training data ( $\sim 20\%$  thereof).

# Relation of MSE and maximum Likelihood

- In Linear regression we were Learning  $y$  from  $x$ .  
Instead now we want to Learn  $P(y|x)$ .  
 $\Rightarrow$  we get an error bar.



- If we assume that the error is Gaussian we want to Learn:

$$P(y|x) = \mathcal{N}(y | \hat{y}(\vec{x}, \vec{w}), \sigma^2)$$

mean of prediction

width of prediction error.

- The loss is now the Likelihood of the training data:

$$\mathcal{L}(\theta) = \sum_{i=1}^N \log P(y^{(i)} | x^{(i)}; \theta)$$

model parameters,  
here  $\vec{\theta} = \vec{w}$  is linear regression.

$$= -n \log \sigma - \frac{n}{2} \log(2\pi) -$$

$$\sum_{i=1}^N \frac{|\hat{y}_{(\theta)}^{(i)} - y^{(i)}|^2}{\sigma^2}$$

MSE Loss

- We want to maximize this with respect to  $\theta$ .

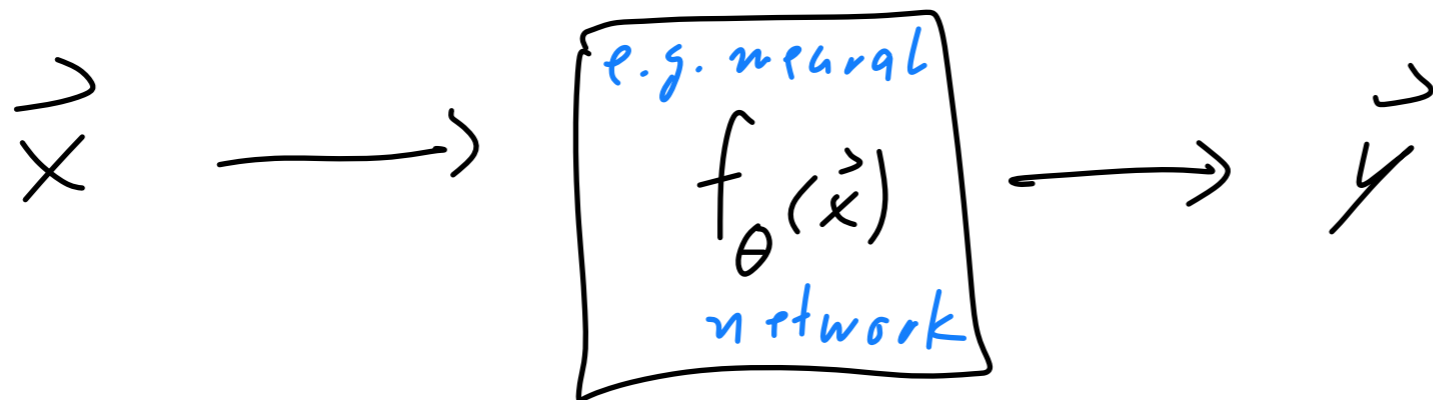
Maximize  $\mathcal{L}$  is the same as minimizing MSE.

# Unit 2: Machine Learning Basics

## 2.2 Neural Network Basics (Supervised Learning)

# Supervised machine Learning

We want to learn a complicated non-linear function to map input  $\vec{x}$  to output  $\vec{y}$ .



$\theta$ : model parameters = weights

We want to find the model parameters by minimizing some loss function on our training data.

$$\text{E.g. } L_{(\theta)}^{\text{MSE}} = \frac{1}{N} \sum_{i=1}^N (y_i - f_{\theta}(\vec{x}_i))^2$$

# A very simple neural network

- We need some way to specify the function  $f_{\theta}(\vec{x})$ .  
It turns out that a very simple "architecture" works very well in practice:

Linear transformations combined  
with element-wise non-linearities

- A Linear transformation / Linear Layer\* is  
given by  $\vec{y} = f(\vec{x}) = W\vec{x} + \vec{b}$   

$\vec{x} : N \text{ dim}$   
 $\vec{y} : M \text{ dim}$

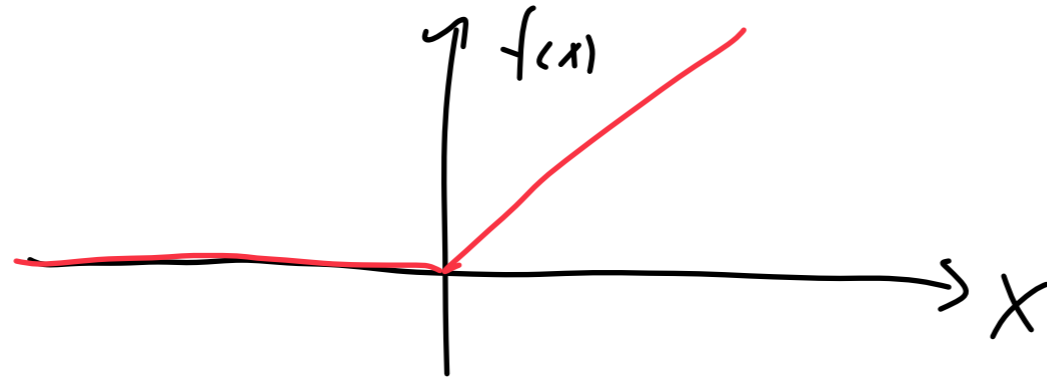
$\uparrow$   
matrix  
of weights

$\uparrow$   
bias

- The most common non-linearity / activation function  
is the ReLU "rectified Linear unit".

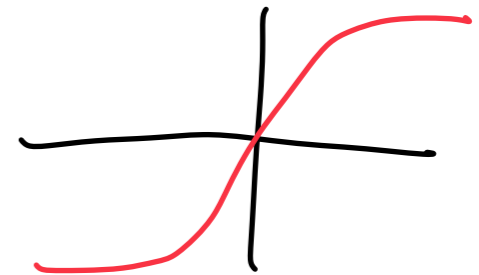
\* fully connected Linear Layer

- ReLU is a 1-dimensional function

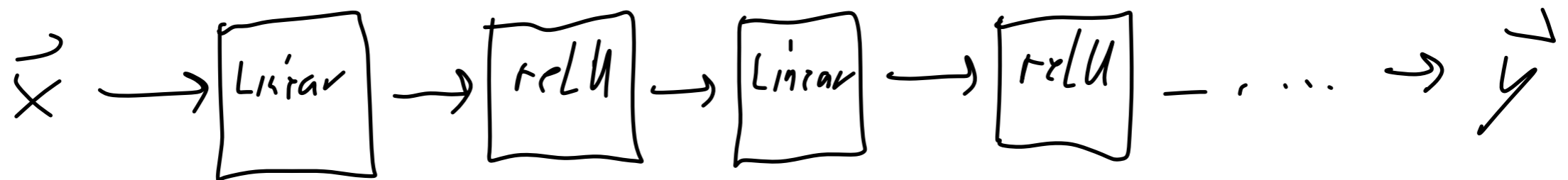


Piecewise linear function.

Alternatives: tanh, sigmoid



- A multilayer perceptron (MLP) is a "deep neural network" made out of a stack of the 2 building blocks,



The combination of ReLU and Linear Layer  
can be written as

$$f^i(x) = \max(0, W_{ij}^i x^j + b^i)$$

ReLU

component notation

next Layer uses  $f(x)$  as  $x$  input.

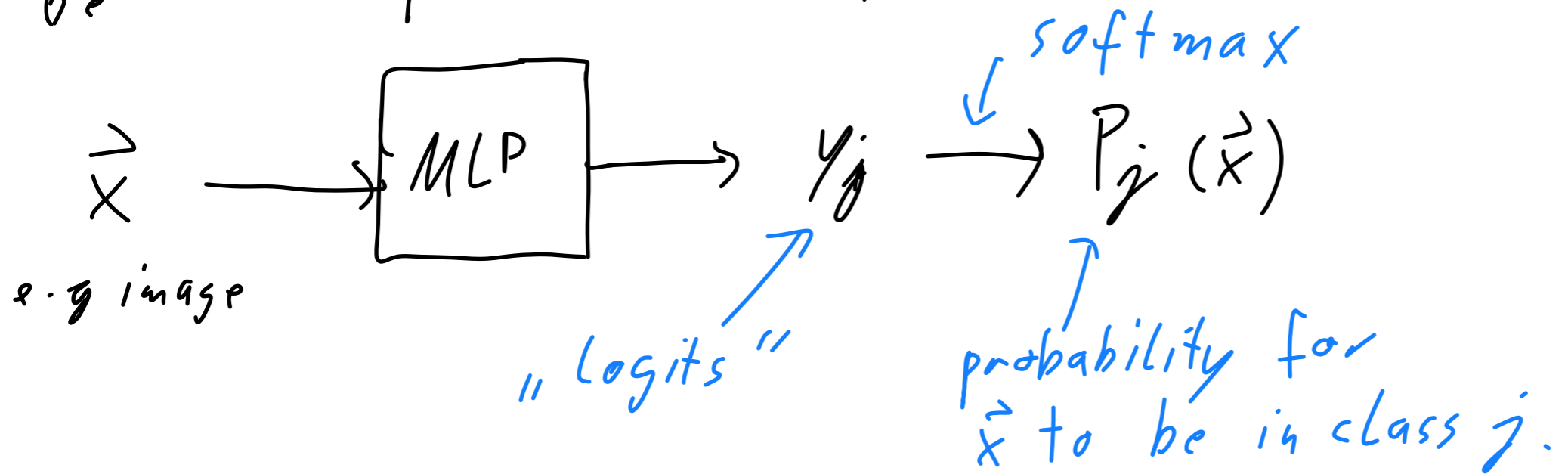
If we stack these layers:

$$\max(0, W_{n,j}^i \max(0, W_{n-1,k}^j \max(0, \dots) + b_{n-1}^j) + b_n^i)$$

Now we have a function that we can  
"fit" to the training data.

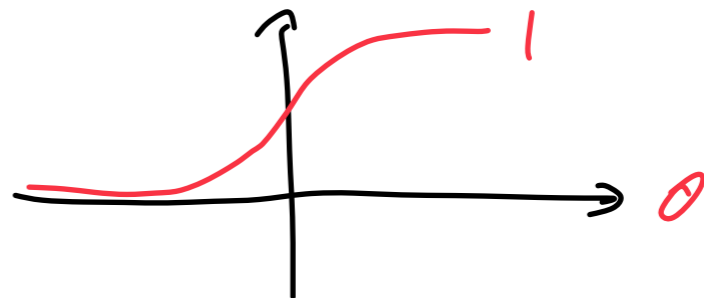
# Loss for classification

For classification we want the NN output to be the probabilities of the various classes.



A probability distribution need to have only positive probabilities and they must sum to 1. To achieve this we need the softmax function

$$P_j(y_j) = \frac{e^{y_j}}{\sum_i e^{y_i}}$$



The typical loss function for classification is the negative Log Likelihood = cross-entropy:

$$L(\theta) = - \sum_{i=1}^{N_{\text{examples}}} \sum_{j=1}^{N_{\text{classes}}} y_j^i \log P_j(\vec{x}_i)$$

$\uparrow$   
parameters of NN:  
 $W, b$

$y_i = \begin{pmatrix} 0 \\ \vdots \\ 1 \\ \vdots \\ 0 \end{pmatrix}$  one-hot vector which is 1 for the right class.

predicted by NN

# Course logistics

- **Reading for this lecture:**
  - **For example:** [Deeplearningbook.org](https://deeplearningbook.org) chapter 5.
- **Problem set:** First problem due Sunday night.