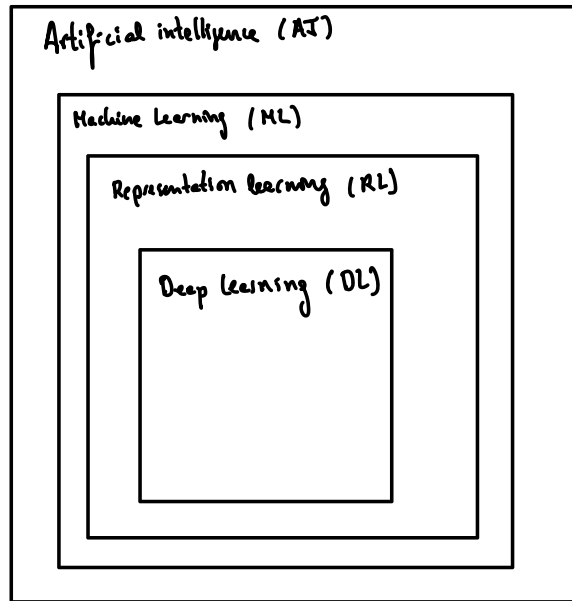


## ② Machine learning: Introduction

\* Machine learning / Artificial intelligence / Deep learning / Representation learning



AI: { Knowledge Base  
Example: Inferring after some statements/topics  
"Reasoning"

ML: { extracting patterns from raw data  
Example: logistic regression

RL: { Finding convenient ways of representing the data (features)  
Example: PCA / autoencoders

DL: { Representations expressed in terms of simpler representations  
(large number of unknowns)  
Example: ANN

## \* Supervised/Unsupervised learning

- \* Supervised:  $\mathcal{D} = \{ \underset{\text{features}}{\tilde{x}^i}, \underset{\text{outputs}}{\tilde{y}^i} \}$  \* Example  $\left\{ \begin{array}{l} y = \text{temperature} \\ x_1 = \text{day}; \\ x_2 = \text{place}; \end{array} \right.$  \* Cases  $\left\{ \begin{array}{l} \text{regression} \\ \text{linear} \\ \text{polynomial} \\ \text{classification} \end{array} \right.$
- \* Unsupervised:  $\mathcal{D} = \{ \tilde{x}^i \}$  \* Cases  $\left\{ \begin{array}{l} \text{Clustering} \\ \text{Dimension reduction} \end{array} \right.$

## \* Predictions - Representation matters

\*  $y_i \approx f_{\theta}(x_i);$

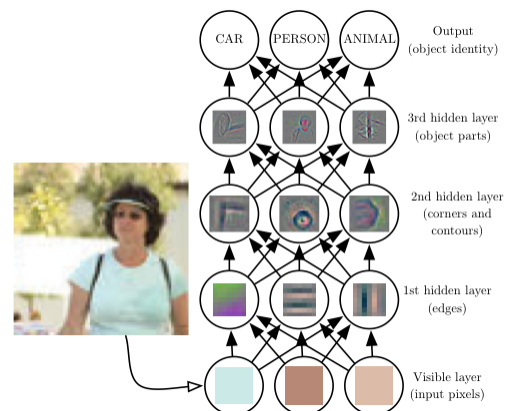
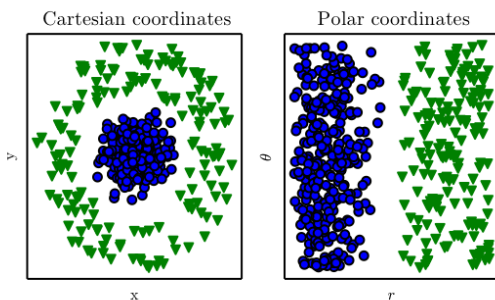
\* Linear predictor:  $f_{\theta}(x_i) = \theta_0 + \theta_1 x_1 + \theta_2 x_2$

\* Polynomial predictor:  $f_{\theta}(x_i) = \theta_0 + \theta_1 x_1 + \theta_2 x_2 + \theta_3 x_1^2 + \theta_4 x_1 x_2 + \theta_5$   
(Machine learning)

¿ why? exp?

\* Representation learning example:  $f_{\theta}(x_i) = \theta_i \underbrace{f_{\theta}^i(x)}_{\text{learn also the representation}}$

\* Deep-learning example:  $f_{\theta}(x_i) = \theta_i f_{\theta}^i(f_{\theta}(f_{\theta}(\dots(x))). \dots)$



\* First example: Linear regression  
in  $l^2$  norm

$$\hat{y} = X\theta; \quad X = \begin{bmatrix} 1 & x_1 & x_2 \end{bmatrix} \left\{ \begin{array}{l} \text{all data} \\ \text{samples} \end{array} \right\}; \quad y = \begin{bmatrix} y_1 \\ y_2 \\ y_3 \end{bmatrix} \text{ output}$$

$$\theta = \{\theta_1, \theta_2, \theta_3\}$$

$$\min_{\theta} J(\theta) = \text{MSE}(\theta) = \|y - X\theta\|_2^2 \quad A = \begin{bmatrix} \end{bmatrix}$$

- ↳ Normal equations:  $\nabla J(\theta) = 0 \Rightarrow X^T X \theta = X^T y$
- ↳ Gradient method:  $\theta_{k+1} = \theta_k - \alpha X^T [X\theta_k - y]$
- ↳ Geometric interpretation:  $\min_{\theta, u} \|y - u\|_2^2$   
 $\left( \begin{array}{l} \text{Projection in the linear} \\ \text{space generated by} \\ \text{columns of } A \end{array} \right) \quad u = X\theta$

\* Second example: Linear regression  
in  $l^1$  norm

$$\min_{\theta} J(\theta) = \|y - X\theta\|_1 = \sum_{i=1}^n |y_i - x_i \theta|$$

- ↳ Convex problem / not differentiable
- ↳ Adding auxiliary variables (double of variables)

$$\begin{cases} \min_{\theta, t} & 1^T t \\ & -t \leq y - X\theta \leq t \end{cases}$$

↳ Canonical form

$$\begin{cases} \min_{v} & c^T v \\ & Av \leq b \end{cases}$$

$$\text{with } v = [\theta, t]^T, \quad A = \begin{bmatrix} X & -I \\ -X & -I \end{bmatrix}$$

$$c^T = [0, 1^T]; \quad b = \begin{bmatrix} y \\ y \end{bmatrix}$$

\* Third example: Linear regression  
(Chebyshev approx) in  $l^\infty$  norm

$$\min_{\theta} \phi(\theta) = \|y - X\theta\|_\infty = \max_i |y_i - x_i \theta|$$

↳ convex problem / not differentiable

↳ Adding auxiliary variables (adding 1 variable)

$$\begin{cases} \min_{\theta, t} & t \\ \text{s.t.} & -t \leq y - X\theta \leq t \end{cases}$$

↳ Canonical form

$$\begin{cases} \min_{\mathbf{v}} & \mathbf{c}^T \mathbf{v} \\ \text{s.t.} & A\mathbf{v} \leq \mathbf{b} \end{cases}$$

with  $\mathbf{v} = [0, t]^T$ ,  $A = \begin{bmatrix} X & -1 \\ -X & -1 \end{bmatrix}$   
 $\mathbf{c}^T = [0, 1]$ ;  $\mathbf{b} = \begin{bmatrix} y \\ -y \end{bmatrix}$

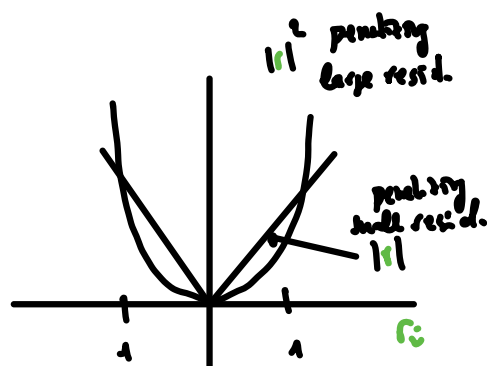
\* Different norms: interpretation in terms of penalty function

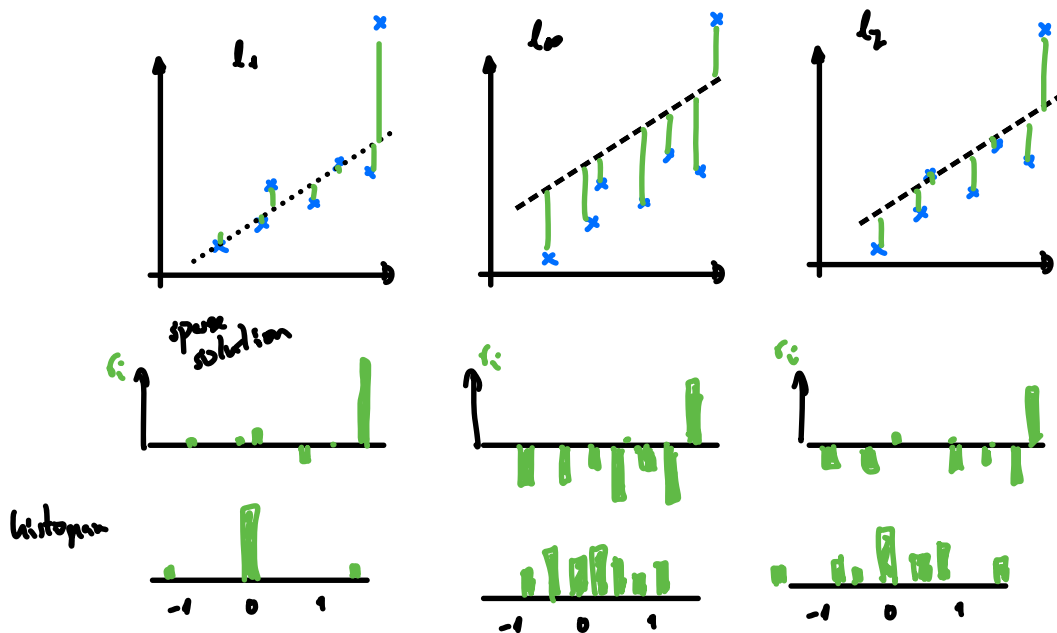
Defining  $r_i = y_i - x_i \theta$

$$L_2: [\sum \phi(r_i)]^2 \leq \sum r_i^2$$

$$L_1: \sum \phi(r_i) \leq \sum |r_i|$$

$$L_\infty: \max_i |r_i|$$





Other penalties:

$$\phi(r) = \begin{cases} r^T & \\ \phi(r) = \begin{cases} -a^2 \log(1 - (r/a)^2) & |r| < a \\ \infty & |r| \geq a \end{cases} & \\ \phi(r) = \begin{cases} 0 & |r| < a \\ |r| - a & |r| \geq a \end{cases} & \end{cases}$$

\* Hypothesis space / capacity / overfitting

\* Why not polynomial regression? Hypothesis space:

$$\hat{y}_i = \theta_0 + \theta_1 x_{i1} + \theta_2 x_{i2} + \theta_3 x_{i1}^2 + \theta_4 x_{i2}^2 + \theta_5 x_{i1} x_{i2} + \dots$$

\* : Which dimension of polynomial?

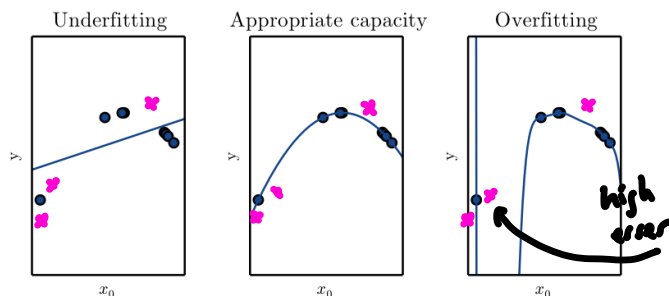
Is always useful to increase the dimension/capacity?

\* Divide data in train (80%) and test (20%)

{ Train : to minimize MSE and obtain  $\theta$  ( $MSE^{train}$ )  
Test : validate ( $MSE^{test}$ )

\* Optimization vs ML { Optimization goal: find best minimizers  
ML goal: predict

\* Capacity

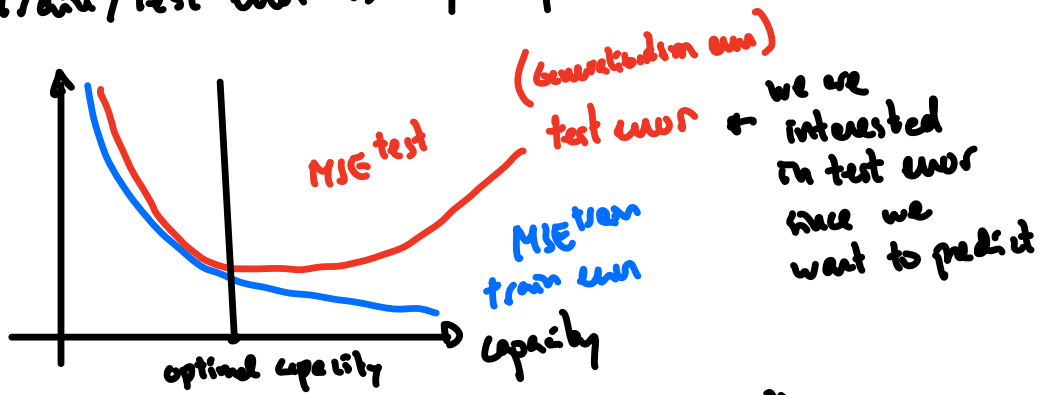


{ Too low capacity : { High error in train  
underfitting { Too few variables  $\theta$  in optim.  
Test error also high

{ Appropriate capacity : { low train error  
low test error

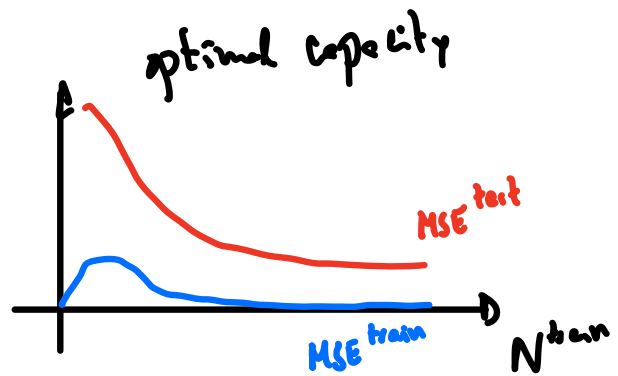
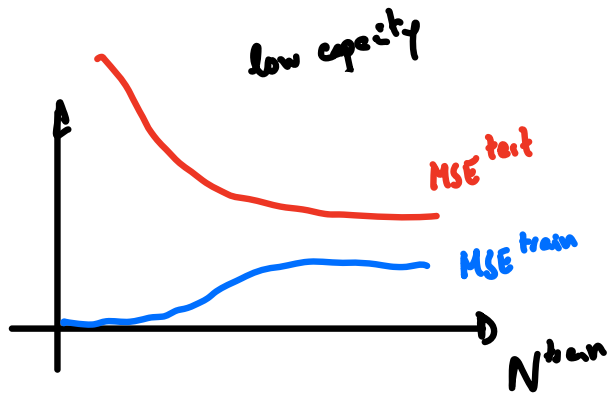
{ High capacity : { Very low train error  
overfitting { High test error

# \* Train/test error vs capacity



$$\| MSE^{test} - MSE^{train} \|^2 \leq f(c, N^{train})$$

$\uparrow$  capacity       $\uparrow$  number of data



\* Underfitting and overfitting from the optimization point of view

$$\min \|y - X\theta\|^2$$

$$\rightarrow (X^T X)\theta = X^T y$$

$$\left\{ \begin{array}{l} \text{Underfitting: } X = \begin{bmatrix} \vdots \\ \vdots \end{bmatrix}; \end{array} \right.$$

$$X^T X = \begin{bmatrix} \vdots \\ \vdots \end{bmatrix}; \quad \theta \text{ is too small in comparison with data}$$

$$\left\{ \begin{array}{l} \text{Overfitting: } X = \begin{bmatrix} \vdots \\ \vdots \end{bmatrix}; \quad X^T X = \begin{bmatrix} \vdots \\ \vdots \end{bmatrix} \\ \text{not full rank} \end{array} \right.$$

$\theta$  is too large in comparison with data

## \* Regularization

\*  $L^2$  regularization (Tikhonov regularization / Ridge regression)

$$\min \|y - X\theta\|_2^2 + \lambda \underbrace{\|\theta\|_2^2}_{\theta^T \theta}$$

$$\underbrace{(X^T X + \lambda I)}_{\text{full rank}} \theta = X^T y$$

- $\lambda > 0$
- Prior knowledge of  $\theta$
- Penalize large values of  $\theta$
- Also useful when  $X$  is square but bad conditioned i.e. elasticity?

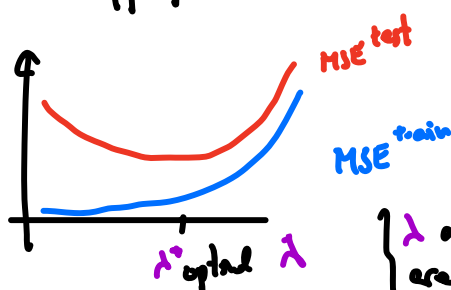
\* Other possible regularization

$$\left\{ \begin{array}{l} L^1 \text{ regularization } |\theta| \\ \|\nabla \theta\|_2^2 \rightarrow \text{smoothing (Laplacian)} \\ |\nabla \theta| \rightarrow \text{Total variation} \end{array} \right.$$

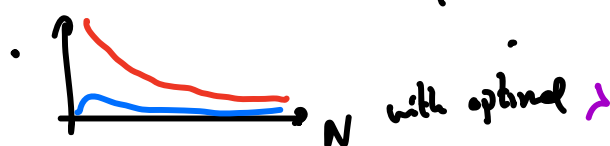


## \* Effect of regularization

- Possible to use large capacity with appropriate  $\lambda$  value



$\lambda$  and dimensions are called hyperparameters



with optimal  $\lambda$

## \* Cross-validation

- \* Useful when data is limited
- \* Divide data  $D = \cup D_i$  with  $D_i \cap D_j = \emptyset$  (around 5/6 times)
- \* Compute mean and standard deviation
  - $\mu = \frac{1}{L} \sum_i MSE_i^{test}$
  - $\sigma^2 = \frac{1}{L} \sum_i (MSE_i^{test} - \mu)^2$
  - we plot  $\mu \pm 1.96 \sigma$