

# Overview of Machine Learning from Optimizational Point of View

Konstantin Vorontsov

CMCM 2021 • MIPT • October 27–29, 2021

## 1 Supervised Learning

- Regression and Classification
- Regularization
- Learning to Rank

## 2 Unsupervised Learning

- Density Estimation
- Clustering and Semi-Supervised Learning
- Representation Learning and Autoencoders

## 3 Multicriteria and Multimodel Learning

- Transfer Learning and Multi-task Learning
- Learning a model from another model
- Generative Adversarial Net

## General optimization problem for many Machine Learning tasks

**Given:** a training set of objects  $\{x_i\}_{i=1}^{\ell}$

**Find:** parameters  $w$  of the predictive model  $a(x, w)$

**Minimize** the empirical risk

$$\sum_{i=1}^{\ell} L_i(w) \rightarrow \min_w$$

where  $L_i(w)$  is a loss function of the model  $a(x, w)$  at the object  $x_i$

More generally, minimize the regularized empirical risk

$$\sum_{i=1}^{\ell} L_i(w) + \sum_{j=1}^r \tau_j R_j(w) \rightarrow \min_w$$

where  $R_j$  is regularization criterion,  $\tau_j$  is regularization coefficient

## Regression as optimization problem

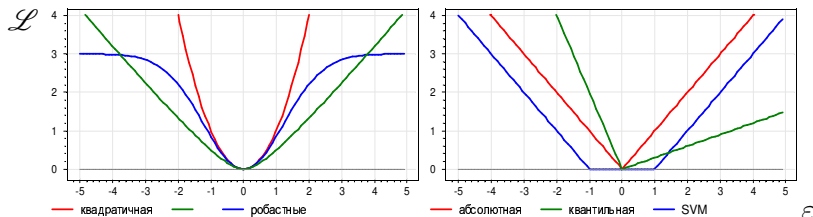
**Given:** a training set of objects  $(x_i, y_i)_{i=1}^{\ell}$ ,  $y_i \in \mathbb{R}$

**Find:** parameters  $w$  of the regression model  $a(x, w)$

**Minimize** the empirical risk

$$\sum_{i=1}^{\ell} \mathcal{L}(a(x_i, w) - y_i) \rightarrow \min_w$$

Unimodal loss function  $\mathcal{L}(\varepsilon)$  of the difference  $\varepsilon = a(x, w) - y$ :



## Classification as optimization problem

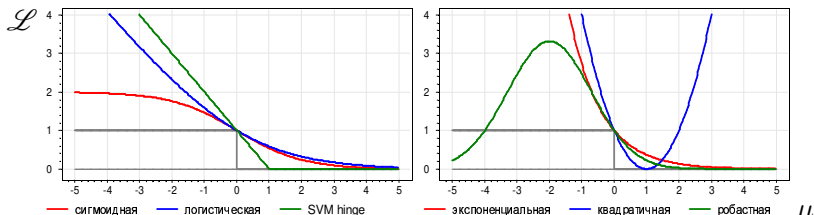
**Given:** a training set of objects  $(x_i, y_i)_{i=1}^{\ell}$ ,  $y_i \in \{-1, +1\}$

**Find:**  $w$  of the classification model  $a(x, w) = \text{sign } g(x, w)$

**Minimize** the empirical risk

$$\sum_{i=1}^{\ell} [g(x_i, w)y_i < 0] \leq \sum_{i=1}^{\ell} \mathcal{L}(g(x_i, w)y_i) \rightarrow \min_w$$

Decreasing loss function  $\mathcal{L}(\mu)$  of the margin  $\mu = g(x, w)y$ :



## Multi-class classification as optimization problem

**Given:** a training set of objects  $(x_i, y_i)_{i=1}^{\ell}$ ,  $y_i \in Y$ ,  $|Y| < \infty$

**Find:**  $w_y$  of the classification model  $a(x, w) = \arg \max_{y \in Y} g(x, w_y)$

The model of the class probability for a given object:

$$P(y|x, w) = \frac{\exp g(x, w_y)}{\sum_{z \in Y} \exp g(x, w_z)} = \text{SoftMax}_{y \in Y} g(x, w_y), \quad y \in Y$$

where  $\text{SoftMax}: \mathbb{R}^Y \rightarrow \mathbb{R}^Y$  is a smooth transformation of a vector into a normalized vector of a discrete distribution.

**Maximize** the log-likelihood of the data (log-loss):

$$-\sum_{i=1}^{\ell} \ln P(y_i|x_i, w) \rightarrow \min_w$$

## Regularizers that penalize the complexity of a linear model

*Regularizer* is an additive complexity penalty to the main criterion:

$$\sum_{i=1}^{\ell} \mathcal{L}(\langle x_i, w \rangle, y_i) + \tau \text{penalty}(w) \rightarrow \min_w$$

where  $\mathcal{L}(a, y)$  is a loss function,  $\tau$  is regularization coefficient

*L<sub>2</sub>-regularization* (ridge regression, SVM):

$$\text{penalty}(w) = \|w\|_2^2 = \sum_{j=1}^n w_j^2.$$

*L<sub>1</sub>-regularization* (LASSO, ElasticNet for feature selection):

$$\text{penalty}(w) = \|w\|_1 = \sum_{j=1}^n |w_j|.$$

*L<sub>0</sub>-regularization* (Akaike/Bayes Information Criteria AIC/BIC):

$$\text{penalty}(w) = \|w\|_0 = \sum_{j=1}^n [w_j \neq 0].$$

## Non-smooth regularizers for feature selection

A general form of a regularizer with selectivity parameter  $\mu$ :

$$\sum_{i=1}^{\ell} \mathcal{L}(\langle x_i, w \rangle, y_i) + \tau \sum_{j=1}^n R_{\mu}(w_j) \rightarrow \min_w .$$

Regularizer with grouping effect for multi-collinear features:

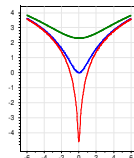
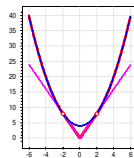
**Elastic Net:**  $R_{\mu}(w) = \mu|w| + w^2$

**Support Features Machine (SFM):**

$$R_{\mu}(w) = \begin{cases} 2\mu|w|, & |w| \leq \mu; \\ \mu^2 + w^2, & |w| \geq \mu; \end{cases}$$

**Relevance Features Machine (RFM):**

$$R_{\mu}(w) = \ln(\mu w^2 + 1)$$





## Learning to Rank

**Given:** a training set of objects  $\{x_i\}_{i=1}^{\ell}$

$i \prec j$  — partial order relation on object pairs  $(x_i, x_j)$

**Find:** parameters  $w$  of the ranking model  $a(x, w)$

$$i \prec j \Rightarrow a(x_i, w) < a(x_j, w)$$

**Minimize** number of misordered pairs  $(x_i, x_j)$  or approximated pairwise empirical risk:

$$\sum_{i \prec j} [a(x_j, w) < a(x_i, w)] \leq \sum_{i \prec j} \mathcal{L} \left( \underbrace{a(x_j, w) - a(x_i, w)}_{\mu_{ij}(w)} \right) \rightarrow \min_w$$

where  $\mathcal{L}(\mu)$  is a decreasing loss function of pairwise margin  $\mu_{ij}(w)$

## Density Estimation

**Given:** a training set of objects  $\{x_i\}_{i=1}^{\ell}$

**Find:** parameters  $\theta$  of the density model  $p(x|\theta)$

**Minimize** Likelihood Estimation (MLE)

$$\sum_{i=1}^{\ell} \ln p(x_i|\theta) \rightarrow \max_{\theta}$$

or Maximum A Posteriori (MAP) estimation:

$$\sum_{i=1}^{\ell} \ln p(x_i|\theta) + \ln p(\theta|\gamma) \rightarrow \max_{\theta}$$

where  $\gamma$  is a hyperparameter of a prior distribution

## Mixture Density Estimation

**Given:** a training set of objects  $\{x_i\}_{i=1}^{\ell}$

**Find:** parameters  $w_j, \theta_j$  of the mixture  $p(x|\theta, w) = \sum_{j=1}^K w_j p(x|\theta_j)$

**Minimize** Likelihood Estimation (MLE)

$$\sum_{i=1}^{\ell} \ln p(x_i|\theta, w) \rightarrow \max_{\theta, w}$$

or Maximum A Posteriori (MAP) estimation:

$$\sum_{i=1}^{\ell} \ln p(x_i|\theta, w) + \ln p(\theta, w|\gamma) \rightarrow \max_{\theta, w}$$

where  $\gamma$  is a hyperparameter of a prior distribution

# Clustering

**Given:** a training set of objects  $\{x_i \in \mathbb{R}^n : i = 1, \dots, \ell\}$

**Find:**

- centers of clusters  $\mu_j \in \mathbb{R}^n, j = 1, \dots, K$
- what cluster  $a_i \in \{1, \dots, K\}$  each object  $x_i$  pertains to

**Minimize** the average intra-cluster distances:

$$\sum_{i=1}^{\ell} \|x_i - \mu_{a_i}\|^2 \rightarrow \min_{\{a_i\}, \{\mu_j\}}$$

in the case of the Euclidean metric

$$\|x_i - \mu_j\|^2 = \sum_{d=1}^n (x_{id} - \mu_{jd})^2$$

## Semi-Supervised Learning

**Given:** labeled  $(x_i, y_i)_{i=1}^k$  and unlabeled  $(x_i)_{i=k+1}^\ell$  data

**Find:** classification  $(a_i)_{i=k+1}^\ell$  of unlabeled objects

**Minimize** the combined clustering/classification criterion:

- with no classification model (Transductive Learning):

$$\sum_{i=1}^{\ell} \|x_i - \mu_{a_i}\|^2 + \lambda \sum_{i=1}^k [a_i \neq y_i] \rightarrow \min_{\{a_i\}, \{\mu_j\}}$$

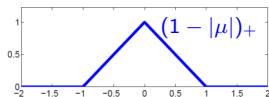
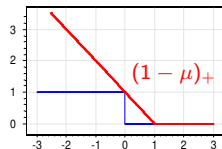
- with classification model,  $a_i = a(x_i, w)$ :

$$\sum_{i=1}^{\ell} \|x_i - \mu_{a_i}\|^2 + \lambda \sum_{i=1}^k \mathcal{L}(a(x_i, w), y_i) \rightarrow \min_{\{a_i\}, \{\mu_j\}, w}$$

## Transductive Learning of a margin-based classifier

$\mu_i(w) = g(x_i, w)y_i$  is margin of the  $x_i$  object

- loss function  $\mathcal{L}(\mu) = (1 - \mu)_+$  penalizes labeled objects for margin decreasing
- loss function  $\mathcal{L}(\mu) = (1 - |\mu|)_+$  penalizes unlabeled objects for falling into the gap between classes



**Minimize** the combined clustering/classification criterion:

$$\sum_{i=1}^k (1 - \mu_i(w))_+ + \gamma \sum_{i=k+1}^{\ell} (1 - |\mu_i(w)|)_+ \rightarrow \min_w$$

## Low-rank matrix factorization

- Generation of better feature vector representation of objects
- Recovering missing values in a matrix

**Given**  $Z = \|z_{ij}\|_{n \times m}$  matrix,  $(i, j) \in \Omega \subseteq \{1..n\} \times \{1..m\}$

**Find:** matrixes  $X = \|x_{it}\|_{n \times k}$  и  $Y = \|y_{tj}\|_{k \times m}$

**Minimize**

$$\|Z - XY\| = \sum_{(i,j) \in \Omega} \mathcal{L}\left(z_{ij} - \sum_t x_{it} y_{tj}\right) \rightarrow \min_{X, Y}$$

Why the classic SVD is abandoned in practice:

- non-square loss function  $\mathcal{L}$
- non-negative matrix factorization:  $x_{it} \geq 0, y_{tj} \geq 0$
- sparse data:  $|\Omega| \ll nm$
- orthogonality is unnecessary or not interpretable

## Autoencoders: unsupervised learning

**Given** a training set of objects  $\{x_i\}_{i=1}^{\ell}$

**Find:**

encoder  $f: X \rightarrow Z$  that produces code vector  $z = f(x, \alpha)$

decoder  $g: Z \rightarrow X$  that reconstructs vector  $\hat{x} = g(z, \beta)$  from  $z$

**Minimize**

the reconstruction error under square loss  $\mathcal{L}(\hat{x}, x) = \|\hat{x} - x\|^2$ :

$$\sum_{i=1}^{\ell} \mathcal{L}(g(f(x_i, \alpha), \beta), x_i) \rightarrow \min_{\alpha, \beta}$$

**Examples of autoencoders:**

$$f(x, A) = \underset{m \times n}{A} x, \quad g(z, B) = \underset{n \times m}{B} z \quad - \text{linear}$$

$$f(x, A) = \sigma(Ax), \quad g(z, B) = \sigma(Bz) \quad - \text{neural}$$



## Autoencoders for supervised learning

Given labeled  $(x_i, y_i)_{i=1}^k$  and unlabeled  $(x_i)_{i=k+1}^{\ell}$  data

Find:

$$z_i = f(x_i, \alpha) \text{ — encoder}$$

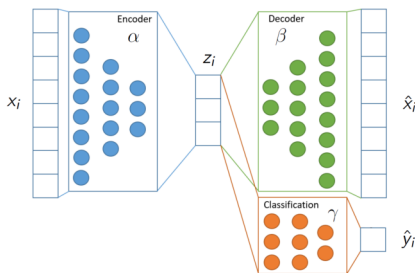
$$\hat{x}_i = g(z_i, \beta) \text{ — decoder}$$

$$\hat{y}_i = \hat{y}(z_i, \gamma) \text{ — predictor}$$

Loss function:

$$\mathcal{L}(\hat{x}_i, x_i) \text{ — for reconstruction}$$

$$\tilde{\mathcal{L}}(\hat{y}_i, y_i) \text{ — for prediction}$$



Minimize the combined reconstruction/prediction criterion:

$$\sum_{i=1}^{\ell} \mathcal{L}(g(f(x_i, \alpha), \beta), x_i) + \lambda \sum_{i=1}^k \tilde{\mathcal{L}}(\hat{y}(f(x_i, \alpha), \gamma), y_i) \rightarrow \min_{\alpha, \beta, \gamma}$$

## Graph Factorization

**Given** a set  $(i, j) \in E$  of edges of the graph  $\langle V, E \rangle$ ,  
similarities  $S_{ij}$  between vertices of the edge  $(i, j)$

For example,  $S_{ij} = [(i, j) \in E]$  is binary adjacency matrix

**Find:** vector representation (embedding) of vertices such that  
adjacent vertices would have similar vectors

**Minimize** the reconstruction error of graph edges:

- in the case of undirected graph and symmetric  $S$  matrix

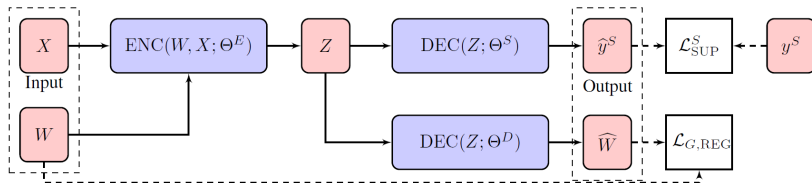
$$\sum_{(i,j) \in E} (\langle z_i, z_j \rangle - S_{ij})^2 \rightarrow \min_Z, \quad Z \in \mathbb{R}^{V \times d}$$

- in the case of directed graph and asymmetric  $S$  matrix

$$\sum_{(i,j) \in E} (\langle \varphi_i, \theta_j \rangle - S_{ij})^2 \rightarrow \min_{\Phi, \Theta}, \quad \Phi, \Theta \in \mathbb{R}^{V \times d}$$

# GraphEDM: a big family of graph autoencoders

*Graph Encoder Decoder Model* generalizes more than 30 models:



$W \in \mathbb{R}^{V \times V}$  is input data about edges

$X \in \mathbb{R}^{V \times n}$  is input feature data about vertices

$Z \in \mathbb{R}^{V \times d}$  is vector representation (embedding) of vertices

$\text{DEC}(Z; \Theta^D)$  is decoder reconstructing the edge data

$\text{DEC}(Z; \Theta^S)$  is decoder solving an applied supervised task

$y^S$  is (semi-)supervised data about vertices or edges

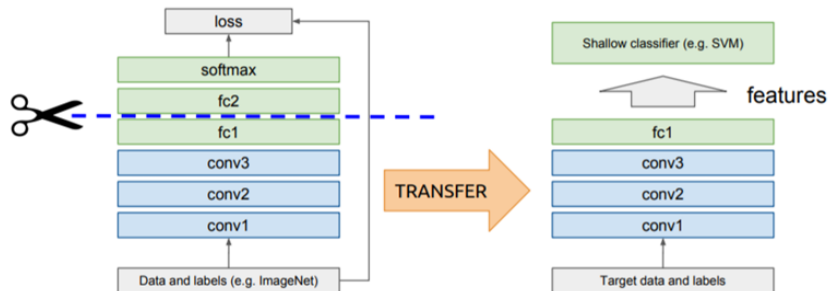
$\mathcal{L}$  is loss function

*I. Chami et al.* Machine learning on graphs: a model and comprehensive taxonomy. 2020.

## Pre-training of neural networks

*Convolutional Neural Network (CNN)* for image classification:

- $z = f(x, \alpha)$  is convolutional layers for image vectorization
- $y = g(z, \beta)$  is feedforward layers for vector classification



Jason Yosinski, Jeff Clune, Yoshua Bengio, Hod Lipson. How transferable are features in deep neural networks? 2014.

## Transfer learning

$f(x, \alpha)$  is the universal part of the model (object vectorization)

$g(x, \beta)$  is a specific part of the model targeted for an applied task

Base task on a dataset  $\{x_i\}_{i=1}^{\ell}$  with loss  $\mathcal{L}_i$ :

$$\sum_{i=1}^{\ell} \mathcal{L}_i(f(x_i, \alpha), g(x_i, \beta)) \rightarrow \min_{\alpha, \beta}$$

Target task on another dataset  $\{x'_i\}_{i=1}^m$ , with another  $\mathcal{L}'_i, g'$ :

$$\sum_{i=1}^m \mathcal{L}'_i(f(x'_i, \alpha), g'(x'_i, \beta')) \rightarrow \min_{\beta'}$$

if  $m \ll \ell$  then pre-training vectorizer  $f(x_i, \alpha)$  could be better than

$$\sum_{i=1}^m \mathcal{L}'_i(f(x'_i, \alpha), g'(x'_i, \beta')) \rightarrow \min_{\alpha, \beta'}$$

## Multi-task learning

$f(x, \alpha)$  is the universal part of the model (object vectorization)

$g_t(x, \beta)$  is a specific part of the model targeted for the task  $t \in T$

Joint training of the model  $f$  from datasets  $X_t$  of tasks  $t \in T$ :

$$\sum_{t \in T} \sum_{i \in X_t} \mathcal{L}_{ti}(f(x_{ti}, \alpha), g_t(x_{ti}, \beta_t)) \rightarrow \min_{\alpha, \{\beta_t\}}$$

The property of *learnability*: we learn the task  $\langle X_t, \mathcal{L}_t, g_t \rangle$  better by augmenting data size  $|X_t|$

*Learning to learn*: we learn each of the tasks  $\langle X_t, \mathcal{L}_t, g_t \rangle$  better by augmenting the number of tasks  $|T|$

*Few-shot learning*: to solve the problem  $t$ , a small number of examples may be enough, sometimes even one

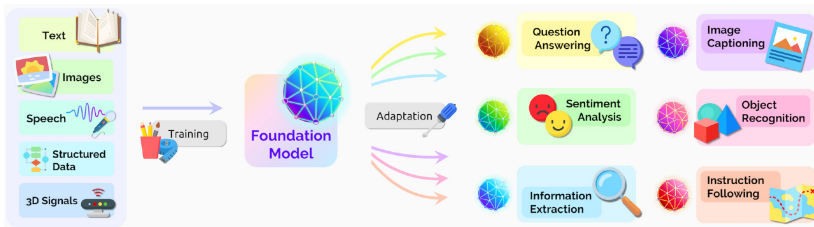
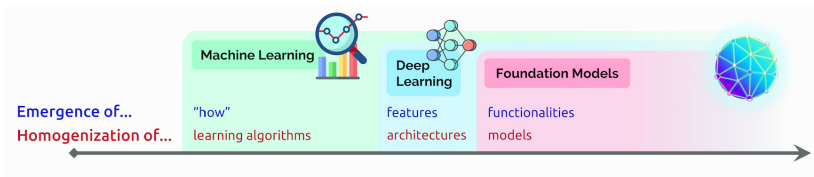
---

M. Crawshaw. Multi-task learning with deep neural networks: a survey. 2020

Y. Wang et al. Generalizing from a few examples: a survey on few-shot learning. 2020

## Foundation Models

Multi-task learnable data vectorization is a recent trend in AI/ML



*R. Bommasani et al. (Center for Research on Foundation Models, Stanford University)*  
On the opportunities and risks of foundation models // CoRR, 20 August 2021.

## Distillation and surrogate modeling

Learning a resource intensive heavy model  $a(x, w)$

$$\sum_{i=1}^{\ell} \mathcal{L}(a(x_i, w), y_i) \rightarrow \min_w$$

Learning a light model  $b(x, w')$ , possibly on other dataset

$$\sum_{i=1}^k \mathcal{L}(b(x'_i, w'), a(x'_i, w)) \rightarrow \min_{w'}$$

### Examples:

- approximation of a heavy model, which is calculated on a supercomputer for months (climate, aerodynamics, etc.), by a light surrogate model
- approximation of a heavy neural network which learns for weeks on big data, by a light neural network with fewer neurons and connections



## Learning Using Privileged Information (LUPI)

$x_i^*$  — information about  $x_i$  available only for training

Student model and **teacher model** are learned separately:

$$\sum_{i=1}^{\ell} \mathcal{L}(a(x_i, w), y_i) \rightarrow \min_w \quad \sum_{i=1}^{\ell} \mathcal{L}(a(x_i^*, w^*), y_i) \rightarrow \min_{w^*}$$

Student model learns from responses of the **teacher model**:

$$\sum_{i=1}^{\ell} \mathcal{L}(a(x_i, w), y_i) + \mu \mathcal{L}(a(x_i, w), a(x_i^*, w^*)) \rightarrow \min_w$$

Student model and **teacher model** are learned together:

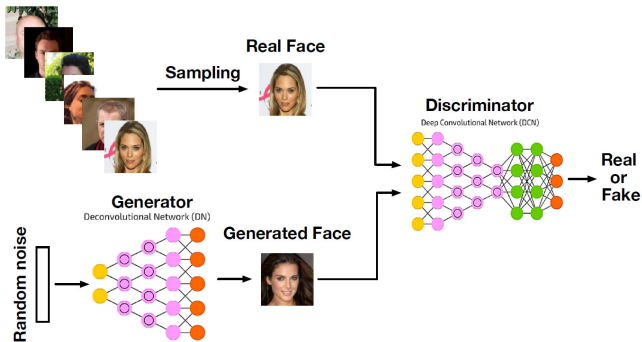
$$\sum_{i=1}^{\ell} \mathcal{L}(a(x_i, w), y_i) + \lambda \mathcal{L}(a(x_i^*, w^*), y_i) + \mu \mathcal{L}(a(x_i, w), a(x_i^*, w^*)) \rightarrow \min_{w, w^*}$$

---

*D.Lopez-Paz, L.Bottou, B.Scholkopf, V.Vapnik.* Unifying distillation and privileged information. 2016.

# Generative Adversarial Net (GAN)

*Generator*  $G(z)$  learns to generate realistic objects  $x$  from noise  $z$   
*Discriminator*  $D(x)$  learns to distinguish is object real or fake



Antonia Creswell et al. Generative Adversarial Networks: an overview. 2017.

Zhengwei Wang et al. Generative Adversarial Networks: a survey and taxonomy. 2019.

Chris Nicholson. A Beginner's Guide to Generative Adversarial Networks.

<https://pathmind.com/wiki/generative-adversarial-network-gan>. 2019.

## Generative Adversarial Net (GAN)

**Given** a training set of objects  $\{x_i\}_{i=1}^{\ell}$

**Find** two probabilistic models:

- model  $G(z, \alpha)$  generates  $x \sim p(x|z, \alpha)$  from noise  $z$
- model  $D(x, \beta) = p(1|x, \beta)$  recognizes if object  $x$  is real

**Minimax** in the antagonistic game of generator vs. discriminator:

- *discriminator*  $D(x, \beta)$  learns to maximize log-likelihood in order to better distinguish real object  $x$  from the fake one
- *generator*  $G(z, \alpha)$  learns to minimize log-likelihood in order to generate realistic objects  $x$

$$\sum_{i=1}^{\ell} \ln D(x_i, \beta) + \ln(1 - D(G(z_i, \alpha), \beta)) \rightarrow \max_{\beta} \min_{\alpha}$$

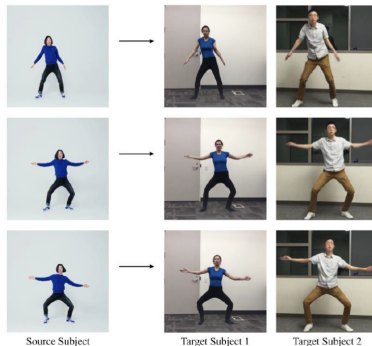
## Examples of how GAN generates realistic image and video



(d) input image

(e) output 3d face

(f) textured 3d face



*Chuan Li, Michael Wand.* Precomputed Real-Time Texture Synthesis with Markovian Generative Adversarial Networks. 2016.

*Xiaoxing Zeng, Xiaojiang Peng, Yu Qiao.* DF2Net: A Dense Fine Finer Network for Detailed 3D Face Reconstruction. ICCV-2019.

*Caroline Chan, Shiry Ginosar, Tinghui Zhou, Alexei A. Efros.* Everybody Dance Now. ICCV-2019.

- 1 Предварительная обработка (data preparation)
  - извлечение признаков (feature extraction)
  - отбор признаков (feature selection)
  - восстановление пропусков (missing values)
  - фильтрация выбросов (outlier detection)
- 2 Обучение с учителем (supervised learning)
  - классификация (classification)
  - регрессия (regression)
  - ранжирование (learning to rank)
  - прогнозирование (forecasting)
- 3 Обучение без учителя (unsupervised learning)
  - кластеризация (clustering)
  - восстановление плотности (density estimation)
  - поиск ассоциативных правил (association rule learning)
  - одноклассовая классификация (anomaly detection)
- 4 Частичное обучение (semi-supervised learning)
  - трансдуктивное обучение (transductive learning)
  - обучение с положительными примерами (PU-learning)

# Optimization criteria induce a typology of ML tasks

- 5 Обучение представлений (representation learning)
  - обучение признаков (feature learning)
  - матричные разложения (matrix factorization)
  - обучение многообразий (manifold learning)
- 6 Глубокое обучение (deep learning)
- 7 Обучение близости/связей (similarity/relational learning)
- 8 Перенос обучения (transfer learning)
- 9 Многозадачное обучение (multitask learning)
- 10 Привилегированное обучение (privileged learning, distilling)
- 11 Состязательное обучение (adversarial learning)
- 12 Обучение структуры модели (structure learning)
- 13 Динамическое обучение (online/incremental learning)
- 14 Активное обучение (active learning)
- 15 Обучение с подкреплением (reinforcement learning)
- 16 Мета-обучение (meta-learning, AutoML)